DISEÑO OPTIMO DE SISTEMAS DE INGENIERIA

ANTECEDENTES MATEMATICOS Y NUMERICOS DE LAS
TECNICAS DE OPTIMIZACION

ANEXO

DR. JORGE ANGELES ALVAREZ
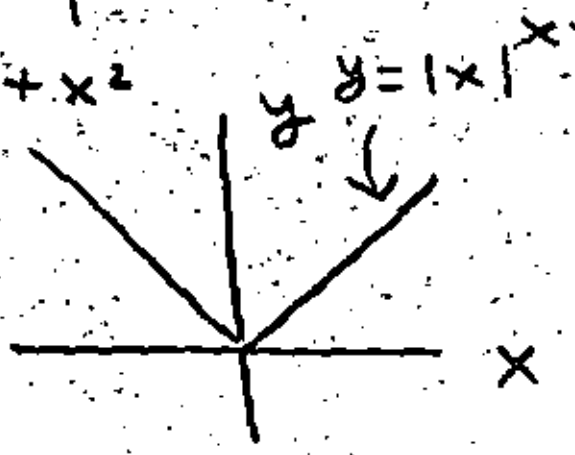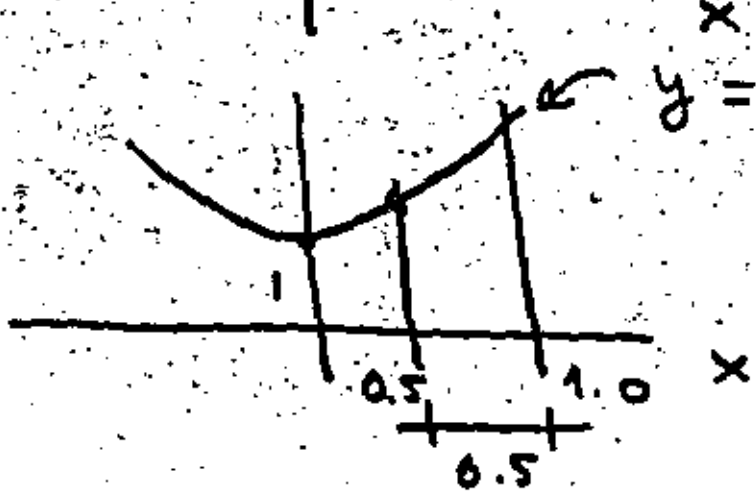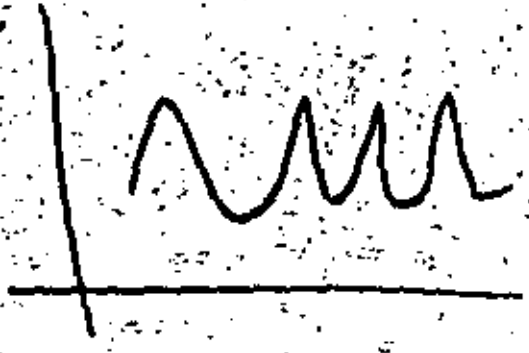
MARZO 1983

$y'(x) = 0$
$y''(x) < 0$

$y = f(x)$

$y'(x) = 0$
$y''(x) > 0$

$y'(x) = 0$

$y''(x) = 0$

$y = 1 + x^2$

$y = |x|$

0.5    1.0

0.5

$$y = y(x_1, x_2, \ldots, x_n) = y(\underset{\sim}{x})$$

$$\nabla y = \left[ \frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, \ldots, \frac{\partial y}{\partial x_n} \right]^T$$

$$\nabla y = \underset{\sim}{0}$$

$$\nabla\nabla y = \begin{bmatrix} \dfrac{\partial^2 y}{\partial x_1{}^2} & \dfrac{\partial^2 y}{\partial x_1 \partial x_2} & \cdots & \dfrac{\partial^2 y}{\partial x_1 \partial x_n} \\[2mm] \dfrac{\partial^2 y}{\partial x_1 \partial x_2} & \dfrac{\partial^2 y}{\partial x_2{}^2} & & \dfrac{\partial^2 y}{\partial x_2 \partial x_n} \\[2mm] \vdots & & & \\[2mm] \dfrac{\partial^2 y}{\partial x_1 \partial x_n} & \dfrac{\partial^2 y}{\partial x_2 \partial x_n} & & \dfrac{\partial^2 y}{\partial x_n{}^2} \end{bmatrix}$$

## Punto Silla



## Espacio Vectorial

Campo $F$ : conjunto de números

Si $x_1$ y $x_2 \in F$, $x_1 + x_2 \in F$ $(x_1 + x_2 = x_2 + x_1)$

$\exists \, 0 \, \ni \, x + 0 = x$, $\forall x$

$x_1 x_2 \in F$ $(x_1 x_2 = x_2 x_1)$

$\exists \, 1 \, \ni \, x \cdot 1 = 1 \cdot x = x$, $\forall x$

$$x_1(x_2 x_3) = (x_1 x_2)x_3$$

$$x_1 + (x_2 + x_3) = (x_1 + x_2) + x_3$$

$$x_1(x_2 + x_3) = x_1 x_2 + x_1 x_3$$

$$\forall x \; \exists \; -x \; \ni \; x + (-x) = 0$$

$$\forall x \; \exists \; x^{-1} \; \ni \; x x^{-1} = x^{-1}x = 1 \; (x \neq 0)$$

----//----

i) Aditivas

Si $\underset{\sim}{v_1}, \underset{\sim}{v_2}$ y $\underset{\sim}{v_3} \in V$ (esp. vect.)

$\underset{\sim}{v_1} + \underset{\sim}{v_2} = \underset{\sim}{v_2} + \underset{\sim}{v_1} \in V$ (cerradura)

$\exists \; \underset{\sim}{0} \; \ni \; \forall \underset{\sim}{v} \in V, \; \underset{\sim}{v} + \underset{\sim}{0} = \underset{\sim}{v}$

ii) Multiplicativas

Si $\alpha$ y $\beta \in F$

$\alpha \underset{\sim}{v} = \underset{\sim}{v} \alpha \in V$

iii) Distributiva

$$(\alpha + \beta)\underset{\sim}{v} = \alpha \underset{\sim}{v} + \beta \underset{\sim}{v}$$

$$(\alpha \beta)\underset{\sim}{v} = \alpha(\beta \underset{\sim}{v})$$

$$\alpha(\underset{\sim}{v_1} + \underset{\sim}{v_2}) = \alpha \underset{\sim}{v_1} + \alpha \underset{\sim}{v_2}$$

$$\underset{\sim}{v} = [v_1, v_2, \ldots, v_n]^T, \; v_i \in \mathbb{R}, \; \forall i$$

$$A = \{\underset{\sim}{v_1}, \underset{\sim}{v_2}, \ldots, \underset{\sim}{v_m}\}, \; \underset{\sim}{v_i} \in V$$

$$B = \{\alpha_1, \alpha_2, \ldots, \alpha_m\}, \; \alpha_i \in \mathbb{R}$$

$$\underset{\sim}{l} = \alpha_1 \underset{\sim}{v_1} + \alpha_2 \underset{\sim}{v_2} + \ldots + \alpha_m \underset{\sim}{v_m}$$

$A$ es $l.i.$ si y sólo si

$$\underset{\sim}{l} = \underset{\sim}{0} \iff \alpha_1 = \alpha_2 = \ldots = \alpha_m = 0$$

De otra forma, $A$ es $l.d.$

## Base de un espacio

$B = \{\underset{\sim}{v_1}, \ldots, \underset{\sim}{v_n}\} \in V$ es una base para $V$ si

i) $B$ es $l.i.$

ii) Cualquier $\underset{\sim}{v} \in V$ se puede expresar como

$$\underset{\sim}{v} = c_1 \underset{\sim}{v_1} + c_2 \underset{\sim}{v_2} + \ldots + c_n \underset{\sim}{v_n}$$

$$\underset{\sim}{v}_1 = \underset{\sim}{i} + 2j \; , \quad \underset{\sim}{v}_2 = -\underset{\sim}{i} + \underset{\sim}{j}$$

¿ $\{\underset{\sim}{v}_1, \underset{\sim}{v}_2\}$ es una base para el

espacio $E^2$ $(X, Y)$ ?



$$\underset{\sim}{v} = x\underset{\sim}{i} + y\underset{\sim}{j}$$

$$\underset{\sim}{\ell} = \alpha_1 \underset{\sim}{v}_1 + \alpha_2 \underset{\sim}{v}_2 = \alpha_1(\underset{\sim}{i} + 2j) + \alpha_2(-\underset{\sim}{i} + \underset{\sim}{j})$$

$$= (\alpha_1 - \alpha_2)\underset{\sim}{i} + (2\alpha_1 + \alpha_2)\underset{\sim}{j} = \underset{\sim}{0}$$

$$\begin{cases} \alpha_1 - \alpha_2 = 0 & (1) \\ 2\alpha_1 + \alpha_2 = 0 & (2) \end{cases}$$

$(1) \Rightarrow \alpha_2 = \alpha_1 \quad (3)$

$(3)$ en $(2) \Rightarrow 3\alpha_1 = 0 \Rightarrow \alpha_1 = 0 \Rightarrow \alpha_2 = 0$

$\Rightarrow \{\underset{\sim}{v}_1, \underset{\sim}{v}_2\}$ es $\ell.i.$

¿ Existen números reales $c_1$ y $c_2$

$\ni$

$$\underset{\sim}{v} = c_1 \underset{\sim}{v}_1 + c_2 \underset{\sim}{v}_2 \, ?$$

$$x\underset{\sim}{i} + y\underset{\sim}{j} = c_1(\underset{\sim}{i} + 2j) + c_2(-\underset{\sim}{i} + \underset{\sim}{j}) =$$

$$= (c_1 - c_2)\underset{\sim}{i} + (2c_2 + c_1)\underset{\sim}{j}$$

$$c_1 - c_2 = x \qquad (4)$$

$$2c_2 + c_1 = y \qquad (5)$$

$(4) \Rightarrow c_2 = c_1 - x \qquad (6)$

$(6)$ en $(5) \Rightarrow 2(c_1 - x) + c_1 = y$

$\Rightarrow 3c_1 - x = y \Rightarrow c_1 = \dfrac{x+y}{3}$

$\Rightarrow c_2 = \dfrac{x+y}{3} - x = \dfrac{-2x+y}{3}$

$\Rightarrow \{\underset{\sim}{U}_1 \text{ y } \underset{\sim}{U}_2\}$ $\underline{sí}$ generan a $E^2$

$\Rightarrow \{\underset{\sim}{U}_1 \text{ y } \underset{\sim}{U}_2\}$ son $\underline{\text{una base para } E^2}$

Si $\{\underset{\sim}{U}_1, \underset{\sim}{U}_2, \ldots, \underset{\sim}{U}_n\}$ es una base

para $V$, dim $V = n$

En el ejemplo anterior, si

$x = 5, y = 10, c_1 = 5, c_2 = 0$

$\underset{\sim}{U} = 5\underset{\sim}{i} + 10\underset{\sim}{j}$ tiene la $\underline{re\text{-}}$

$\underline{presentación}$ $[5, 10]^T$ en la

base $\{\underset{\sim}{i}, \underset{\sim}{j}\}$; pero tiene la

representación $[5, 0]^T$ en
la base $\{\underline{v}_1, \underline{v}_2\}$

---

↓ Un espacio vectorial $V$
una región $\Omega$ de $V$ $(\Omega \subseteq V)$ ↓



$\Omega$ (región
factible)

$f :$ un funcional definido
sobre $\Omega$

Problema: Hallar el punto
$\underline{x}^* \in \Omega$ donde $f$ alcance un
valor óptimo (máximo o mínimo)

$f :$ función objetivo

$V$ es de dimensión $\begin{cases} \text{finita (algebrai} \\ \text{infinita} \\ \text{(cálculo de} \\ \text{variaciones)} \end{cases}$

$f(z)$



Volumen mínimo

Area dada



$$V = \int_{x_1}^{x_2} y(x) \, ds \; ; \quad ds = \sqrt{dx^2 + dy^2}$$

$$V = \sum_{1}^{n-1} y_i \, \Delta_i \; , \quad \Delta_i = \sqrt{\Delta x_i^2 + \Delta y_i^2}$$

$$\Delta x_i = x_{i+1} - x_i \, , \; \Delta y_i = y_{i+1} - y_i$$

# Transformaciones lineales

$U^m, V^n$ : espacios vectoriales sobre $F$.

$T : U^m \longrightarrow V^n$

$\underline{u} \in U^m, \quad \underline{v} \in V^n, \quad \underline{v} = T(\underline{u})$

$\underline{v}$ está definido para cada $\underline{u}$

$T$ es lineal si

$T(\alpha_1 \underline{u}_1 + \alpha_2 \underline{u}_2) = \alpha_1 T(\underline{u}_1) + \alpha_2 T(\underline{u}_2)$

$\alpha_1, \alpha_2 \in F$

Si $T$ es lineal: $\underline{v} = \underline{\underline{T}}\,\underline{u}$

$\underline{\underline{T}}$ es <u>suproyectiva</u> si para cualquier $\underline{v} \in V^n$, existe(n) $\underline{u}(s)$
∂
$\underline{v} = \underline{\underline{T}}\,\underline{u}$

<u>Ejemplo:</u>

$P(x,y,z) \qquad \underline{v} = T(\underline{u})$

$U^m = V^n = E^3$

$P'(x,y,0)$

$T$ es **uno a uno** si y sólo si

$$u_1 \neq u_2 \Rightarrow T u_1 \neq T u_2$$

Si $T$ es suproyectiva y uno a uno, se dice que es invertible:

$$v = T u \Rightarrow u = T^{-1} v$$

$$B_u = \{u_1, u_2, \ldots, u_m\} \in U^m$$

$$B_v = \{v_1, v_2, \ldots, v_n\} \in V^n$$

$$T: U^m \to V^n$$

$$T u_i \in V^n$$

$$\Rightarrow T u_i = \alpha_{1i} v_1 + \alpha_{2i} v_2 + \cdots + \alpha_{ni} v_n, \quad i = 1, 2, \ldots, m$$

$$\alpha_{1i}, \alpha_{2i}, \ldots, \alpha_{ni} \in F$$

$$[\lambda] \equiv \begin{bmatrix} \alpha_{11} & \alpha_{12} & \cdots & \alpha_{1m} \\ \alpha_{21} & \alpha_{22} & & \alpha_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_{n1} & \alpha_{n2} & & \alpha_{nm} \end{bmatrix}$$

$$U^m = V^n = E^3$$

$$B_u = B_v = \{\underline{i}, \underline{j}, \underline{k}\}$$

$P(x,y,z)$

$P(x,y,-z)$

$$R\underline{i} = -\underline{i}, \quad R\underline{j} = \underline{j}, \quad R\underline{k} = -\underline{k}$$

$$[R] = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & -1 \end{bmatrix}$$

Dominio de una T.L.

$$T: U^m \to V^n \; ; \quad Tu = v$$

Codominio (range): el conjunto de vectores $\underline{v}$ para los cuales existe por lo menos una $\underline{u} \ni T\underline{u} = \underline{v}$

Espacio nulo de $T$: el conjunto de vectores $\underline{u}$ para los cuales

$$T\underline{u} = \underline{0}$$

Dominio de $T = U^m$

dim dom = dim cod + dim E.N.

$$\underset{\sim}{v} = \underset{\sim}{T}\, \underset{\sim}{r}$$



Si $\underset{\sim}{T}: U^m \longrightarrow V^n \ni U^m \cong V^n$,

$\underset{\sim}{T}$ es un <u>isomorfismo</u>

Dom = Cod

$$\underset{\sim}{T}\,\underset{\sim}{u} = \underset{\sim}{v} \in U^m$$

En general $\underset{\sim}{u}$ y $\underset{\sim}{v}$ son l.i.;

pero si sucede que sean l.d.,

esto es, si

$$\underset{\sim}{T}\,\underset{\sim}{u} = \lambda \underset{\sim}{u}, \quad \lambda \in F \quad (\underset{\sim}{u} \neq \underset{\sim}{0})$$

$\underline{u}$ es un vector característico (13) $\underline{T}$

$\lambda$ es un valor característico de $\underline{T}$

asociado a $\underline{u}$.

$$\underline{T}\,\underline{u} = \lambda\,\underline{u} \Rightarrow (\underline{T} - \lambda\underline{I})\,\underline{u} = \underline{0}$$

$$\Rightarrow \underbrace{det(\underline{T} - \lambda\,\underline{I}) = 0}$$

$$P_n(\lambda) = a_0 + a_1\lambda + a_2\lambda^2 + \cdots + a_n\lambda^n$$

Polinomio característico de $\underline{T}$.

Tiene $n$ raíces $\lambda_i$, $i = 1, \ldots, n$,

que son los valores característicos de $\underline{T}$.

$\underline{H}$ es hermitiana si

$$\underline{H} = \underline{H}^+ \quad (transp.\ conjugada)$$

Ejemplo:

$$\underline{H} = \begin{bmatrix} 2 & 1+i \\ 1-i & 3 \end{bmatrix} \quad \text{es hermitiana}$$

Teorema: Los valores car. de una
matriz hermitiana son reales y
sus vectores característicos son ortogo-
nales.

$$\det(\underline{H} - \lambda \underline{I}) = \begin{vmatrix} 2-\lambda & 1+i \\ 1-i & 3-\lambda \end{vmatrix} = (2-\lambda)(3-\lambda)$$

$$\underbrace{-(1-i)(1+i)}_{1-i^2 = 2} =$$

$$= \lambda^2 - 5\lambda + 6 - 2 = \lambda^2 - 5\lambda + 4 =$$

$$= (\lambda - 4)(\lambda - 1) \Rightarrow \lambda_1 = 1, \lambda_2 = 4$$

$\lambda_1 = 1:$

$$\begin{bmatrix} 1 & i+1 \\ 1-i & 2 \end{bmatrix} \overset{\underset{\underset{\sim}{e_1}}{}}{\begin{bmatrix} e_{11} \\ e_{21} \end{bmatrix}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$e_{11} + (1+i)e_{21} = 0 \Rightarrow e_{11} = -(1+i)e_{21}$$

$$\|\underset{\sim}{e_1}\|^2 = 1 \Rightarrow \underset{\sim}{e_1} \underset{\sim}{e_1^*} = 1 = \underset{\sim}{e_1^+}\underset{\sim}{e_1}$$

$$[\tilde{e}_{11}, \tilde{e}_{21}]\begin{bmatrix} e_{11} \\ e_{21} \end{bmatrix} = |e_{11}|^2 + |e_{21}|^2$$

$$|e_{11}|^2 = [-(1+i)e_{21}][-(1-i)\tilde{e}_{21}] =$$

$$= (2)|e_{21}|^2$$

$$\Rightarrow 2|e_{21}|^2 + |e_{21}|^2 = 3|e_{21}|^2 = 1$$

$$\Rightarrow |e_{21}|^2 = \frac{1}{3} \Rightarrow e_{21} = \frac{1}{\sqrt{3}} = \frac{\sqrt{3}}{3}$$

$$\Rightarrow e_{11} = -(1+i)\frac{\sqrt{3}}{3} \Rightarrow \underset{\sim}{e_1} = \frac{\sqrt{3}}{3}\begin{bmatrix} -(1+i) \\ 1 \end{bmatrix}$$

$$(\underline{H} - \lambda_i \underline{I}) \underline{e}_i = \underline{0}$$

$$\begin{bmatrix} -2 & 1+i \\ 1-i & -1 \end{bmatrix} \begin{bmatrix} e_{12} \\ e_{22} \end{bmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$-2e_{12} + (1+i) e_{22} = 0 \Rightarrow e_{12} = \frac{1+i}{2} e_{22}$$

$$[\tilde{e}_{12}, \tilde{e}_{22}] \begin{bmatrix} e_{12} \\ e_{22} \end{bmatrix} = 1$$

$$\Rightarrow |e_{12}|^2 + |e_{22}|^2 = 1 \quad ; \quad |e_{12}|^2 = \left| \frac{1+i}{2} \right|^2 |e_{22}|^2$$

$$\underbrace{\frac{1}{4} 2 = \frac{1}{2}}$$

$$\Rightarrow \frac{1}{2} |e_{22}|^2 + |e_{22}|^2 = 1$$

$$\Rightarrow \frac{3}{2} |e_{22}|^2 = 1 \Rightarrow |e_{22}|^2 = \frac{2}{3} \Rightarrow e_{22} = \sqrt{\frac{2}{3}}$$

$$\Rightarrow \underline{e}_2 = \sqrt{\frac{2}{3}} \begin{bmatrix} (1+i)/2 \\ 1 \end{bmatrix}$$

$$\underline{e}_1^* \underline{e}_2 = \frac{\sqrt{3}}{3} [-1+i, 1] \sqrt{\frac{2}{3}} \begin{bmatrix} (1+i)/2 \\ 1 \end{bmatrix} =$$

$$= \frac{\sqrt{6}}{3\sqrt{3}} \left[ -\frac{1}{2}(-2) + 1 \right] = 0$$

$$\#$$

$$\underline{L} : U^m \to V^n = U^m$$

$$B = \{ \underline{\beta}_1, \underline{\beta}_2, \ldots, \underline{\beta}_n \} \in U$$

$$\underline{v} = \underline{\underline{L}} \, \underline{u}$$

$$[\underline{\underline{L}}]_B = \begin{bmatrix} \ell_{11} & \ell_{12} & \cdots & \ell_{1n} \\ \ell_{21} & \ell_{22} & & \ell_{2n} \\ \vdots & & \vdots & \vdots \\ \ell_{n1} & \ell_{n2} & & \ell_{nn} \end{bmatrix}$$

$$[\underline{v}]_B = [v_1, v_2, \ldots, v_n]$$

$$\underline{v} = v_1 \underline{\beta}_1 + v_2 \underline{\beta}_2 + \cdots + v_n \underline{\beta}_n$$

Dada otra base

$$C = \{\underline{\gamma}_1, \underline{\gamma}_2, \ldots, \underline{\gamma}_n\}$$

¿$[\underline{\underline{L}}]_C = ?$, ¿$[\underline{v}]_C = ?$

$$[\underline{v}]_C = [v_1', v_2', \ldots, v_n']$$

$$\underline{v} = v_1' \underline{\gamma}_1 + v_2' \underline{\gamma}_2 + \cdots + v_n' \underline{\gamma}_n$$

$$\underline{\gamma}_1 = \alpha_{11} \underline{\beta}_1 + \alpha_{21} \underline{\beta}_2 + \cdots + \alpha_{n1} \underline{\beta}_n$$

$$\underline{\gamma}_2 = \alpha_{12} \underline{\beta}_1 + \alpha_{22} \underline{\beta}_2 + \cdots + \alpha_{n2} \underline{\beta}_n$$

$$\underline{\gamma}_n = \alpha_{1n} \underline{\beta}_1 + \alpha_{2n} \underline{\beta}_2 + \cdots + \alpha_{nn} \underline{\beta}_n$$

$$[\underline{A}]_B = \begin{bmatrix} \alpha_{11} & \alpha_{12} & \cdots & \alpha_{1n} \\ \alpha_{21} & \alpha_{22} & & \alpha_{2n} \\ \vdots & & \cdot & \vdots \\ \alpha_{n1} & \alpha_{n2} & & \alpha_{nn} \end{bmatrix}$$

$$\underline{v} = v'_1 \underline{\gamma}_1 + v'_2 \underline{\gamma}_2 + \cdots + v'_n \underline{\gamma}_n =$$

$$= v'_1 \sum_i^n \alpha_{i1} \underline{\beta}_i + v'_2 \sum_i^n \alpha_{i2} \underline{\beta}_i + \cdots + v'_n \sum_i^n \alpha_{in} \underline{\beta}_i$$

$$= \underbrace{\sum_{j=1}^n v'_j \sum_{i=1}^n \alpha_{ij} \underline{\beta}_i}_{\sum_{i=1}^n \sum_{j=1}^n \alpha_{ij} v'_j \underline{\beta}_i} = \sum_{i=1}^n v_i \underline{\beta}_i$$

$$v_i = \sum_{j=1}^n \alpha_{ij} v'_j$$

$$[\underline{v}]_C = [\underline{A}^{-1}]_B [\underline{v}]_B$$

$$[\underline{v}]_B = [\underline{A}]_B [\underline{v}]_C$$
$$\Uparrow$$

$$\begin{bmatrix} v_1 \\ v_2 \\ \cdot \\ \cdot \\ v_n \end{bmatrix} = \begin{bmatrix} \alpha_{11} & \alpha_{12} & \cdots & \alpha_{1n} \\ \alpha_{21} & \alpha_{22} & \cdots & \alpha_{2n} \\ \vdots & & & \\ \alpha_{n1} & \alpha_{n2} & \cdots & \alpha_{nn} \end{bmatrix} \begin{bmatrix} v'_1 \\ v'_2 \\ \cdot \\ \cdot \\ v'_n \end{bmatrix}$$

$$[\underset{\sim}{v}]_B = [\underset{\sim}{A}]_B [\underset{\sim}{v}]_C$$

$$[\underset{\sim}{v}]_C = [\underset{\sim}{A}^{-1}]_B [\underset{\sim}{v}]_B \checkmark$$

$$\underset{\sim}{v} = \underset{\sim}{L} \underset{\sim}{u}$$

$$\rightarrow [\underset{\sim}{v}]_B = [\underset{\sim}{L}]_B [\underset{\sim}{u}]_B$$

$$[\underset{\sim}{v}]_C = [\underset{\sim}{L}]_C [\underset{\sim}{u}]_C \qquad \underset{\nwarrow}{[\underset{\sim}{u}]_C}$$

$$[\underset{\sim}{A}^{-1}]_B [\underset{\sim}{v}]_B = [\underset{\sim}{L}]_C [\underset{\sim}{A}^{-1}]_B [\underset{\sim}{u}]_B$$

$$[\underset{\sim}{v}]_B = [\underset{\sim}{A}]_B \underbrace{[\underset{\sim}{L}]_C [\underset{\sim}{A}^{-1}]_B}_{[\underset{\sim}{L}]_B} [\underset{\sim}{u}]_B$$

$$[\underset{\sim}{L}]_B = [\underset{\sim}{A}]_B [\underset{\sim}{L}]_C [\underset{\sim}{A}^{-1}]_B$$

$$[\underset{\sim}{L}]_C = [\underset{\sim}{A}^{-1}]_B [\underset{\sim}{L}]_B [\underset{\sim}{A}]_B : \text{Transforma}$$

ción similar

$$\text{//}$$

$$\underset{n\times n}{\underset{\sim}{A}} \underset{\sim}{e}_i = \lambda_i \underset{\sim}{e}_i , \quad i = 1, 2, \ldots, n$$

$\underset{\sim}{A}$ es hermitiana

$\Rightarrow$ Tiene n vectores característicos, mutuamente

ortogonales.

$$e_i^T \, e_j = \delta_{ij} \begin{cases} = 1, & i = j \\ = 0, & i \neq j \end{cases}$$

$$Q = [e_1 \vdots e_2 \vdots e_3 \vdots \ldots \vdots e_n] =$$

$$= \begin{bmatrix} e_{11} & e_{12} & \cdots & e_{1n} \\ e_{21} & e_{22} & & e_{2n} \\ \vdots & \vdots & & \vdots \\ e_{n1} & e_{n2} & & e_{nn} \end{bmatrix}$$

$$A e_1 = \lambda_1 e_1, \quad A e_2 = \lambda_2 e_2, \ldots, \quad A e_n = \lambda_n e_n$$

$$\underbrace{A[e_1 \vdots e_2 \vdots \ldots \vdots e_n]}_{Q} = [\lambda_1 e_1 \vdots \lambda_2 e_2 \vdots \ldots \vdots \lambda_n e_n]$$

$$= Q \Lambda$$

$$\Lambda = diag(\lambda_1, \lambda_2, \ldots, \lambda_n)$$

$$\Rightarrow A Q = Q \Lambda \Rightarrow \Lambda = \underset{Q^{-1} = Q^T}{Q^{-1}} A Q$$

$$\Rightarrow \Lambda = Q^T A Q$$

Si $\underline{A}$ y $\underline{B}$ son $n \times n$ y están
relacionadas por

$$\underline{B} = \underline{P}^{-1} \underline{A} \underline{P} \quad (T.S.)$$

$$\Rightarrow \quad \underline{B}^k = \underline{P}^{-1} \underline{A}^k \underline{P}$$

$$\Rightarrow \quad \underline{Q}^T \underline{A}^k \underline{Q} = \underline{\Lambda}^k$$

$$\det(\underline{A} - \lambda \underline{I}) = \begin{vmatrix} 1-\lambda & 2 \\ 2 & 4-\lambda \end{vmatrix} = (1-\lambda)(4-\lambda) - 4$$

$$= \lambda^2 - 5\lambda + 4 - 4 = 0 \Rightarrow$$

$$\lambda^2 - 5\lambda = 0 \Rightarrow \lambda_1 = 0, \quad \lambda_2 = 5$$

$$\underline{e}_1 = \begin{bmatrix} e_{11} \\ e_{21} \end{bmatrix}, \quad \underline{e}_2 = \begin{bmatrix} e_{12} \\ e_{22} \end{bmatrix}$$

$$\underline{A}\,\underline{e}_1 = \lambda_1 \underline{e}_1 \Rightarrow (\underline{A} - \lambda_1 \underline{I}) \underline{e}_1 = \underline{0}$$

$$\begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} \begin{bmatrix} e_{11} \\ e_{21} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow \begin{array}{l} e_{11} + 2e_{21} = 0 \\ \Rightarrow e_{11} = -2e_{21} \end{array}$$

$$\|\underline{e}_1\|^2 = 1 \Rightarrow e_{11}^2 + e_{21}^2 = 1 \Rightarrow 4e_{21}^2 + e_{21}^2 = 1$$

$$\Rightarrow e_{21} = \frac{\sqrt{5}}{5}, \quad e_{11} = -\frac{2}{5}\sqrt{5}$$

$$\begin{bmatrix} -4 & 2 \\ 2 & -1 \end{bmatrix} \begin{bmatrix} e_{12} \\ e_{22} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

(21)

$$-4e_{12} + 2e_{22} = 0 \Rightarrow e_{22} = 2e_{12}$$

$$e_{12}^2 + e_{22}^2 = 1 \Rightarrow 4e_{12}^2 + e_{12}^2 = 1$$

$$\Rightarrow e_{12} = \frac{\sqrt{5}}{5} \Rightarrow e_{22} = \frac{2}{5}\sqrt{5}$$

$$\underline{Q} = [\underline{e}_1 \vdots \underline{e}_2] = \frac{\sqrt{5}}{5}\begin{bmatrix} -2 & \vdots & 1 \\ 1 & \vdots & 2 \end{bmatrix}$$

$$\underline{Q}\,\underline{Q}^T = \frac{1}{5}\begin{bmatrix} -2 & 1 \\ 1 & 2 \end{bmatrix}\begin{bmatrix} -2 & 1 \\ 1 & 2 \end{bmatrix} = \frac{1}{5}\begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix} = \underline{I}$$

$$\underline{\Lambda} = \begin{bmatrix} 0 & 0 \\ 0 & 5 \end{bmatrix} \Rightarrow \underline{\Lambda}^{100} = \begin{bmatrix} 0 & 0 \\ 0 & 5^{100} \end{bmatrix}$$

$$\underline{A}^{100} = \underline{Q}\,\underline{\Lambda}^{100}\underline{Q}^T = \frac{1}{5}\begin{bmatrix} -2 & 1 \\ 1 & 2 \end{bmatrix}\begin{bmatrix} 0 & 0 \\ 0 & 5^{100} \end{bmatrix}\begin{bmatrix} -2 & 1 \\ 1 & 2 \end{bmatrix}$$

Una matriz cuadrada de $n \times n$, $\underline{A}$
es positiva definida (semidefinida) si

$$f = \underline{x}^T \underline{A} \underline{x} > 0 \; (\geq 0) \; , \; \forall \; \underline{x} \neq \underline{0}$$

Si $\underline{A} = \underline{A}^* \Rightarrow f$ es real

$$\underline{A} = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \bigcirc \\ \bigcirc & & \ddots & \\ & & & \lambda_n \end{bmatrix}$$

$$\underline{x}^T \underline{A} \underline{x} = \sum_1^n \lambda_i x_i^2 > 0 \; (\geq 0) \iff \lambda_i > 0 \; (\geq), \; i=1,\ldots,n$$

Una matriz es p.d. ( p.s.d.) si todos
sus valores característicos son positivos
(no negativos)

$$f = f(\underline{x}) \in \mathbb{R}$$

Serie de Taylor alrededor de $\underline{x}_0$

$$f(\underline{x}) = f(\underline{x}_0) + f'(\underline{x}_0)^T(\underline{x} - \underline{x}_0) +$$
$$+ \frac{1}{2}(\underline{x} - \underline{x}_0)^T f''(\underline{x}_0)(\underline{x} - \underline{x}) + \mathbb{R}(\|\underline{x} - \underline{x}_0\|^3)$$

Si $\|\underline{x} - \underline{x}_0\| \to 0$, el término dominante es el lineal. Un punto estacionario de $f$ es aquel en el cual se anula el término lineal $\forall \ \underline{x} - \underline{x}_0$

$$f'(\underline{x}_0) = \nabla f \Big|_{\underline{x}_0}$$

Si es así, se tiene

$$f(\underline{x}) = f(\underline{x}_0) + \frac{1}{2}(\underline{x} - \underline{x}_0)^T f''(\underline{x}_0)(\underline{x} - \underline{x}_0)$$
$$+ \eta \to 0$$

$\underline{x}_0$ en el cual $f'$ se anula es un máximo si $\Delta f = f(\underline{x}) - f(\underline{x}_0) \leq 0$

$$f(\underline{x}) - f(\underline{x}_0) \approx \frac{1}{2}(\underline{x} - \underline{x}_0)^T \underbrace{f''(\underline{x}_0)}_{\text{Hessiana}}(\underline{x} - \underline{x}_0)$$

Para esto, $f''(\underline{x}_0) \leq 0$

Si $f'(\underline{x}_0)$ se anula si $f''(\underline{x}_0) > 0$, $\underline{x}_0$ es un mínimo

Ejemplo: $f = \sum_1^n x_i^4$

$$\nabla f = [4x_1^3, 4x_2^3, \ldots, 4x_n^3]^T$$

$$\nabla\nabla f = \begin{bmatrix} 12x_1^2 & 0 & \cdots & 0 \\ 0 & 12x_2^2 & \cdots & 0 \\ \vdots & & & \\ 0 & 0 & \cdots & 12x_n^2 \end{bmatrix}$$

$\nabla f = \underset{\sim}{0} \implies$ SNC: $\underline{f}(\underset{\sim}{x}) = \underset{\sim}{0}$

Ejemplo: $f(\underset{\sim}{x}) = e^{x_1} \operatorname{sen} x_2 \cosh x_3$

$$\nabla f = [e^{x_1}\operatorname{sen}x_2\cosh x_3, \; e^{x_1}\cos x_2\cosh x_3, \; e^{x_1}\operatorname{sen}x_2\operatorname{senh}x_3]^T$$

$\nabla f = 0 \implies$

$e^{x_1} \operatorname{sen} x_2 \cosh x_3 = 0$

$e^{x_1} \cos x_2 \cosh x_3 = 0$

$e^{x_1} \operatorname{sen} x_2 \operatorname{senh} x_3 = 0$

DISEÑO OPTIMO DE SISTEMAS DE INGENIERIA

ANEXO

Dr. Jorge Angeles Alvarez

MARZO, 1982

# AN OPTIMIZATION MODEL FOR THE PLANNING OF ELECTRICAL DISTRIBUTION NETWORKS

J.M.Cobián                                    M. Ruiz C.
Instituto de Investigaciones Eléctricas
Cuernavaca, Mor.

## INTRODUCTION

The problem of structuring a distribution network in space and time is formulated in this paper as a mixed-integer programming optimization model. The approach also allows for the solution of the problem of the optimal location of distribution substations and their capacity expansion, as well as the related configuration of the corresponding primary-feeder network. This same model can be used also to solve the problem of the optimal design in space and time of secondary networks with their corresponding distribution transformers. The optimization consists of minimizing the present worth of the costs of capital investment and energy losses in the network for the planning horizon to be considered.

Among the advantages of the approach taken here to formulate this problem, as compared to other approaches in the literature, are, that first of all, it considers basically the actual network and not an idealization of it. It also includes the dynamic aspects of the capacity expansion of the network in time, and is able to handle large-scale networks with a relatively small number of integer (binary) variables and automatically solves the important problem of obtaining the optimal distribution of power flows for a given configuration of a network. The application of this model can also help electrical utilities to evaluate the impact of changes in their cost components and in the existing normalization of their equipment in the network.(1,4,6)

The solution algorithm to the problem formulated above is a specialized branch and bound search that takes advantage of the particular structure of this problem to find the optimal solution more efficiently. The corresponding linear programming sub-problems are also of a particular type known as "transshipment" or "distribution" problems which can be solved very rapidly, making the overall algorithm very effective.

## FORMULATION OF THE PROBLEM

### NOMENCLATURE

| | |
|---|---|
| T | number of stages of the planning horizon. |
| NP | number of proposed expansion projects. |
| $D_t$ | peak demand of the entire system at time t. |
| N | {q/q node index of the corresponding circulatory network). |
| RN | {(i,j,m,p) / (i,j,m,p) directed arc from node i to j in the circulatory network corresponding to a segment m of the linearized power loss function, with an index p indicating whether or not the arc is been modified by the acceptance of project p). |
| PN | {(i,j,m,p) / (i,j,m,p) directed arc from node i to j in the circulatory network corresponding to a segment m of the linearized power loss function belonging to project p). |
| $OUT_q$ | {(i,j,m,p) / (i,j,m,p) $\in$ RN and i=q). |
| $rOUT_q$ | {(i,j,m,p) / (i,j,m,p) $\in$ PN and i=q). |

$D_q$    $\{(i,j,m,p) \, / \, (i,j,m,p) \, \epsilon RN \text{ and } j=q\}$.

$PIN_q$    $\{(i,j,m,p) \, / \, (i,j,m,p) \, \epsilon PN \text{ and } j=q\}$.

$x_{ijmpt}$ continuous decision variable representing the power flow of arc $(i,j,m,p) \, \epsilon RN$ at time stage t.

$w_{ijmpt}$ continuous decision variable representing the power flow of arc $(i,j,m,p) \, \epsilon FN$ at time stage t.

$y_{pt}$ binary decision variable indicating whether project p is in operation at time stage t.

$A_{ijmpt}$ unit cost at present value associated to arc $(i,j,m,p) \, \epsilon RN$.

$C_{pt}$ capital cost at present value of the required investment to start project p at time t.

$K_{ijmt}$ capacity of the directed arc $(i,j,m,p) \, \epsilon RN$ ( If j=n+1 the capacity is equal to the peak demand at load node i at stage t. If j=0 the capacity is equal to the total demand $D_t$ of the network at time t).

$KK_{ijmp}$ capacity of the arc $(i,j,m,p) \, \epsilon PN$.

NE    number of restrictions of logical and/or technical dependence between the different projects proposed.

$f_\ell(\cdot)$ functional relationship between the binary variables in the logical and/or technical dependence restrictions.

$B_{ijmpt}$ unit cost at present value associated to arc $(i,j,m,p) \, \epsilon PN$.

## MATHEMATICAL MODEL

It is assumed that: $\ell = 1,\ldots,NE$ ; $p=1,\ldots,NP$ ; $t=1,\ldots T$.

$$\text{Minimize } z = \sum_{t=1}^{T} \left\{ \sum_{(i,j,m,p) \, \epsilon RN} A_{ijmpt} \, x_{ijmpt} + \right.$$

$$\left. \sum_{(i,j,m,p) \, \epsilon PN} B_{ijmpt} \, w_{ijmpt} + \sum_{p=1}^{NP} C_{pt} (y_{pt} - y_{pt-1}) \right\} \quad (1)$$

subject to

$$\sum_{(i,j,m,p) \epsilon OUT_q} x_{ijmpt} + \sum_{(i,j,m,p) \, \epsilon POUT_q} w_{ijmpt} -$$

$$\sum_{(i,j,m,p) \, \epsilon IN_q} x_{ijmpt} - \sum_{(i,j,m,p) \, \epsilon PIN_q} w_{ijmpt} = 0 \, ; q \epsilon N \quad (2)$$

$$D_t \leq x_{n+1,0,0,0,t} \leq K_{n+1,0,0,0,t} \quad (3)$$

$$0 \leq x_{ijm0t} \leq K_{ijmt} \, ; \, (i,j,m,0) \, \epsilon FN \quad (4)$$

$$0 \leq x_{ijmpt} \leq K_{ijmt} (1-y_{pt}) \, ; \, (i,j,m,p) \, \epsilon RN \quad (5)$$

$$0 \leq w_{ijmpt} \leq KK_{ijmp} \, y_{pt} \quad ; \, (i,j,m,p) \, \epsilon PN \quad (6)$$

$$y_{p0} = 0 \quad (7)$$

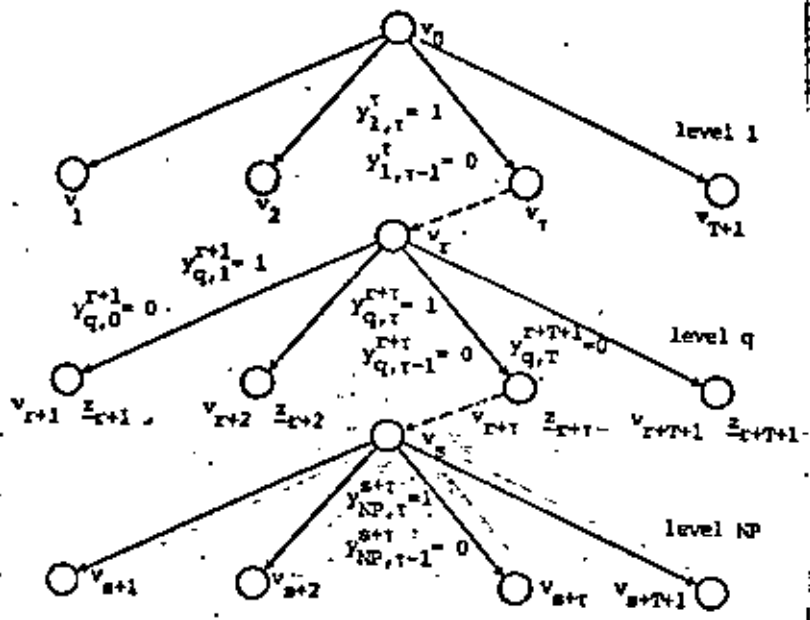$$y_{pt} - y_{pt-1} \leq 0 \quad (8)$$

$$y_{pT+1} = 1 \quad (9)$$

$$f_\ell (y_{1,1},\ldots,y_{NP,T}) = 0 \quad (10)$$

$$y_{pt} = 0 \text{ or } 1 \quad (11)$$

## SOLUTION ALGORITHM

A "branch and bound" search is conducted where to each level of the search tree is associated a project q. Any vertex $v_{r+\tau}$ at level q in the tree fixes all binary variables corresponding to projects with indices less than or equal to q. Each of the T+1 successor vertices of vertex $v_r$ considers that project q becomes operational at stage $\tau$ for $\tau=1,\ldots,$ T+1 (stage T+1 is an artificial stage to indicate that project q never becomes operational) (3).

At each vertex a lower bound is obtained by solving a relaxation of the original restricted problem which consists of setting to zero the capital costs of the free projects and



fixing their corresponding binary variables to one or zero, or both simultaneously, to permit a maximum relaxation of the power flow restrictions. The structure of the resulting mathematical problem is the one known as the transshipment or distribution problem (2,5). The objective function values of the relaxation for different values of $\tau$ are as follows:

$$z_{r+\tau} = \sum_{t=1}^{\tau-1} \min \zeta_{rt}^a + \sum_{t=\tau}^{T} \min \zeta_{rt}^b + \sum_{t=1}^{T} \sum_{p=1}^{q} C_{pt} \cdot (y_{pt}^{r+\tau} - y_{pt-1}^{r+\tau}) \qquad (12)$$

for $t=1,\ldots,T+1$, considering that $\sum_{t=1}^{0} \min \zeta_{rt}^a = 0$ and $\sum_{t=T+1}^{T} \min \zeta_{rt}^b = 0$

where $\min \zeta_{rt}^a$ and $\min \zeta_{rt}^b$ are defined as follows:

Minimize $\zeta_{rt}^a = \sum_{(i,j,m,p) \in RN} A_{ijmpt} x_{ijmpt} +$

$$\sum_{(i,j,m,p) \in PN} B_{ijmpt} w_{ijmpt} \qquad (13)$$

subject to $\sum_{(i,j,m,p) \in OUT_q} x_{ijmpt} + \sum_{(i,j,m,p) \in POUT_q} w_{ijmpt} -$

$$\sum_{(i,j,m,p) \in IN_q} x_{ijmpt} - \sum_{(i,j,m,p) \in PIN_q} w_{ijmpt} = 0 \quad ; q \in N \qquad (14)$$

$$0 \le x_{n+1,0,0,0,t} \le K_{n+1,0,0,0,t} \qquad (15)$$

5. (Upper bound updating). If $q_s \neq NP$ then set $q=q_s+1$, $r=s$, and go to 2. Otherwise, set $\bar{z} = \underline{z}_s$ ; store $y_{pt}^s$ $\forall$ p,t as the best current assignment.

6. (Fathoming by bounds). Any live vertices $v_u$ for which $\underline{z}_u \geq \bar{z}$ are fathomed.

7. ( Backtracking). If there are no live vertices go to 8. Otherwise, choose vertex $v_r$ where r is such that:

$$\underline{z}_r = \min \{ \underline{z}_u ; u/v_u \text{ is a live vertex} \} ; \text{ set } q=q_r+1; \text{ go to 2.}$$

8. (Termination). If $z = \infty$ the problem has no feasible solution. If $\bar{z} < \infty$ the best current assignment is the optimal assignment.

## COMPUTATIONAL EXPERIENCE

An application problem of the above formulation and solution procedure was considered and solved in a VAX 11/780 computer in an approximate total time of two minutes. The problem included the consideration of approximately 50 load nodes and 100 existing arcs, with 5 proposed projects in two-time stages which incorporated a total of 20 different arcs. The corresponding mathematical problem had approximately 900 continuous variables and 10 binary variables.

## REFERENCES

1. Adams, R.N. and M.A. Laughton (1974),"Optimal Planning of Power Networks Using Mixed-Integer Programming – Part 1 – Static and Time – Phased Network Synthesis", Proc. IEE, Vol. 121, No. 2, 139-147.

2. Bradley G.H., G.G. Brown and G.W. Graves (1974), "Design and Implementation of Large Scale Primal Transshipment Algorithms", Management Science, Vol. 20, 1-34.

3. Garfinkel, R.S. and G.L. Nemhauser (1972), Integer Programming, Wiley-Interscience.

4. Masud, E. (1974), "An Interactive Procedure for Sizing and Timing Distribution Substations Using Optimization Techniques", IEEE Trans.PAS, Vol. PAS-93, No. 5, 1281-1286.

5. Simonnard, M. (1966), Linear Programming, Prentice-Hall.

6. Thompson, G.L. and D.L. Wall (1980), "A Branch and Bound Model for Choosing Optimal Substation Locations", Presented at the IEEE-PES Winter Meeting, New York, NY, February 3-8, 1980.

$$0 \leq x_{1jm0t} \leq K_{1jmt} \; ; \; (i,j,m,0) \; \epsilon RN. \tag{16}$$

$$0 \leq x_{1jmpt} \leq K_{1jmt} (1-y_{pt}^{r}) \; ; \; (i,j,m,p) \; \epsilon RN \; ; \; p<q. \tag{17}$$

$$0 \leq x_{1jmqt} \leq K_{1jmt} ; (i,j,m,q) \; \epsilon RN. \tag{18}$$

$$0 \leq x_{1jmpt} \leq K_{1jmt} ; (i,j,m,p) \; \epsilon RN \; ; \; p>q. \tag{19}$$

$$0 \leq w_{1jmpt} \leq KK_{1jmp} \, y_{pt}^{r} ; (i,j,m,p) \; \epsilon PN \; ; \; p<q. \tag{20}$$

$$w_{1jmqt} = 0 \qquad ; (i,j,m,q) \; \epsilon PN. \tag{21}$$

$$0 \leq w_{1jmpt} \leq KK_{1jmp} ; (i,j,m,p) \; \epsilon PN \; ; \; p>q. \tag{22}$$

The problem $\zeta_{rt}^{b}$ is equivalent to the $\zeta_{rt}^{a}$ with the exceptions that set (18) changes to: $x_{1jmqt} = 0$ ; $(i,j,m,q)$ $\epsilon RN.$ $\qquad (23)$

and the set (21) changes to: $0 \leq w_{ijmqt} \leq KK_{1jmq}$; $(i,j,m,q)$ $\epsilon PN.$ $\qquad (24)$

At the last level of the search the relaxation is non-existing, so if it is possible to reach this level without violating restrictions $f_{\ell}(\cdot)$ and if $\underline{z}_{s+\tau}$ exists (see figure), a feasible solution is obtained to the original problem, which provides an upper bound for the search procedure. The solution algorithm is thus as follows:

1. (Initialization). Set $\bar{z} = \infty$ , $q=1$ and $r=0$.

2. (Feasibility phase). For $\tau=1,\ldots,T+1$ test the feasibility of $f_{\ell}(\cdot)=0$ ($\ell=1,\ldots,NE$), assuming that: $y_{qt}^{r+\tau}=0$ for $0 \leq t \leq \tau-1$, $y_{qt}^{r+\tau}=1$ for $\tau \leq t \leq T+1$, and $y_{pt} = y_{pt}^{r+\tau}$ $\forall t$ and $p \leq q$. If for some $\tau$ there is no feasibility, then $\bar{z}_{r+\tau} = \infty$.

3. (Calculation of lower bounds). (a) For $t=1,\ldots,T$ calculate $\min \zeta_{rt}^{a}$ and $\min \zeta_{rt}^{b}$. If $\min \zeta_{rt}^{a}$ does not exist , then let $\min \zeta_{rt}^{a}=\infty$. If $\min \zeta_{rt}^{b}$ does not exist , then let $\min \zeta_{rt}^{b} = \infty$ .

(b) For $\tau=1,\ldots,T+1$ calculate $\underline{z}_{r+\tau}$, except when $\underline{z}_{r+\tau}=\infty$. If $\underline{z}_{r+\tau} < \bar{z}$ then the live vertex $v_{r+\tau}$ is generated; set: $y_{qt}^{r+\tau} = 0$ for $0 \leq t \leq \tau-1$, $y_{qt}^{r+\tau}=1$ for $\tau \leq t \leq T+1$, $y_{pt}^{r+\tau} = y_{pt}^{r}$ $\forall t$ and $p<q$ and $q_{r+\tau} = q$ ; store : $q_{r+\tau}, \underline{z}_{r+\tau}$ and $y_{pt}^{r+\tau}$ $\forall t$ and $p \leq q$. Otherwise, continue.

4. (Branching) If no live vertex was generated in step 3, go to 7. Otherwise, choose vertex $v_{s}$, where s is such that $z_{s} = \min_{s_{\tau/\tau=1,\ldots,T+1}} \{\underline{z}_{r+\tau}\}$

DISEÑO OPTIMO DE SISTEMAS DE INGENIERIA

METODOS DE OPTIMACION CON RESTRICCIONES

Dra Susana Gómez Gómez

Marzo, 1982

# GRADIENT METHODS IN MATHEMATICAL PROGRAMMING
## PART 1 - REVIEW OF PREVIOUS TECHNIQUES

by

A. MIELE, H.Y. HUANG, and J.W. CANTRELL

RICE UNIVERSITY

1969

---

Gradient Methods in Mathematical Programming

Part 1 - Review of Previous Techniques [1]

by

A. MIELE[2], H.Y. HUANG[3], AND J.W. CANTRELL[4]

Abstract. This report is the first of a series on gradient methods in mathematical programming. It considers the problem of minimizing a function $f(x)$, where $f$ is a scalar function and $x$ is an n-vector whose components are unconstrained. For this problem, three previous methods are reviewed, namely, the ordinary gradient method, the conjugate-gradient method, and the variable-metric method. A new intuitive derivation of the last two algorithms is presented.

[2] Professor of Astronautics and Director of the Aero-Astronautics Group, Department of Mechanical and Aerospace Engineering and Materials Science, Rice University, Houston, Texas.

[3] Assistant Professor of Astronautics, Department of Mechanical and Aerospace Engineering and Materials Science, Rice University, Houston, Texas.

[4] Graduate Student in Aero-Astronautics, Department of Mechanical and Aerospace Engineering and Materials Science, Rice University, Houston, Texas.

### 1. Definitions

The following definitions are used throughout the paper:

(a) The symbol x denotes the position vector

$$x = \begin{bmatrix} x^1 \\ x^2 \\ \vdots \\ x^n \end{bmatrix} \qquad (1)$$

whose scalar components are $x^1, x^2, \ldots, x^n$.

(b) The symbol $f$ denotes a scalar function of the vector x, that is

$$f = f(x) \qquad (2)$$

(c) The symbol g denotes the column vector

$$g(x) = \begin{bmatrix} \partial f/\partial x^1 \\ \partial f/\partial x^2 \\ \vdots \\ \partial f/\partial x^n \end{bmatrix} \qquad (3)$$

whose components are the first partial derivatives of $f$ with respect to the scalar variables $x^1, x^2, \ldots, x^n$. This is the gradient of the function $f$.

(d) The symbol H denotes the square matrix

$$H(x) = \begin{bmatrix} \partial^2 f/\partial x^1 \partial x^1 & \partial^2 f/\partial x^1 \partial x^2 & \cdots\cdots & \partial^2 f/\partial x^1 \partial x^n \\ \partial^2 f/\partial x^2 \partial x^1 & \partial^2 f/\partial x^2 \partial x^2 & \cdots\cdots & \partial^2 f/\partial x^2 \partial x^n \\ \cdots\cdots\cdots\cdots\cdots\cdots\cdots\cdots\cdots\cdots \\ \partial^2 f/\partial x^n \partial x^1 & \partial^2 f/\partial x^n \partial x^2 & \cdots\cdots & \partial^2 f/\partial x^n \partial x^n \end{bmatrix} \qquad (4)$$

whose components are the second partial derivatives of the function $f$ with respect to the scalar variables $x^1, x^2, \ldots, x^n$.

(e) The symbol x denotes the nominal point. The symbol $\tilde{x}$ denotes the point following x. The symbol $\hat{x}$ denotes the point preceding x.

(f) The symbol $\delta(\ldots)$ denotes the displacement leading from a point to the next point. Therefore, the following relations hold:

$$\tilde{x} = x + \delta x$$
$$x = \hat{x} + \delta\hat{x} \qquad (5)$$

(g) The superscript T denotes the transpose of a matrix.

## 2. Introduction

A basic problem of mathematical programming is that of finding the minimum of a function

$$f = f(x) \tag{6}$$

where $f$ is a scalar function and $x$ is an $n$-vector. If the $n$ components of the vector $x$ are unconstrained, the extremum of (6) occurs when the following necessary condition is satisfied:

$$g(x) = 0 \tag{7}$$

where $g$ is the gradient of the function $f$ with respect to the vector $x$. For a minimum, the matrix of the second derivatives (4) must be positive definite at the point $x$ defined by (7).

If the function (6) is quadratic, the gradient $g(x)$ is linear with respect to $x$. Hence, Eq. (7) can be solved analytically. On the other hand, if (6) is nonquadratic, the gradient $g(x)$ is nonlinear. This being the case, approximate methods must be employed to solve Eq. (7). One possible method consists of quasilinearizing (7) about a nominal point. Another method, the descent method, consists of constructing corrections $\Delta x$ leading from a nominal point $x$ to a varied point $\bar{x}$ such that

$$f(\bar{x}) < f(x) \tag{8}$$

Thus, by an iterative procedure (that is, through successive decreases in the value of the function $f$), it is hoped that the minimum of $f$ is approached to any desired degree of accuracy.

This report is the first of a series on gradient methods in mathematical programming. It reviews three of the existing techniques, namely, the ordinary gradient method, the conjugate-gradient method (Refs. 1-3), and the variable-metric method (Refs. 4-5). For the last two methods, a new intuitive derivation is presented.

3. **Ordinary Gradient Method.**

To first-order terms, the values of the function (6) at the varied point and the nominal point are related by

$$I(\tilde{x}) \sim I(x) + \delta I(x) \tag{9}$$

where the first variation $\delta I(x)$ is given by

$$\delta I(x) = g^T(x)\delta x \tag{10}$$

with

$$\delta x = \tilde{x} - x \tag{11}$$

Also to first-order terms, the greatest decrease in the value of the function is achieved if the first variation (10) is minimized. Here, we limit our analysis to those variations $\delta x$ which satisfy the constraint

$$K = \delta x^T \delta x \tag{12}$$

where K is a prescribed quantity.

3.1. _Derivation of the Algorithm._ Standard methods of the theory of maxima and minima show that the fundamental function of this problem is the scalar function

$$F = g^T(x)\delta x + (1/2\alpha)\delta x^T \delta x \tag{13}$$

where $1/2\alpha$ is a constant Lagrange multiplier. The optimum system of equations must be such that

$$G(\delta x) = 0 \tag{14}$$

where G is the gradient of the function F with respect to the scalar variables $\delta x^1, \delta x^2, \ldots, \delta x^n$. In the light of (13), the explicit form of (14) is the following:

$$\delta x = -\alpha g(x) \tag{15}$$

and shows that the optimum correction $\delta x$ has the gradient direction. This is why the method is called the _ordinary gradient method._ Upon substituting (15) into (12), we see that

$$K = \alpha^2 g^T(x)g(x) \tag{16}$$

Therefore, a one-to-one correspondence exists between the value of the constant K and the value of $\alpha$. This being the case, one can bypass prescribing K and reason directly on $\alpha$, as in the considerations which follow.

3.2. _Descent Property._ Upon combining Eqs. (10) and (15), we see that the first variation becomes

$$\delta I(x) = -\alpha g^T(x)g(x) \tag{17}$$

and is negative for $\alpha > 0$. Therefore, if $\alpha$ is sufficiently small, the function I decreases. This guaranteed decrease of I at every step is the most important property of the gradient method.

## 4. Modifications of the Ordinary Gradient Method

The ordinary gradient method is conceptually simple and stable in that the function $f(x)$ is reduced at every iteration; however, it has the drawback of slow convergence. For this reason, methods have been developed to reduce the number of iterations required for convergence. In this connection, let the displacement vector $\Delta x$ be written in the form

$$\Delta x = -\alpha p \qquad (23)$$

where $p$ is the search direction. The following are particular forms of the vector $p$:

$$p = g(x) + q \qquad (24)$$

and

$$p = A g(x) \qquad (25)$$

where $q$ is an $n$-vector and $A$ is an $n \times n$ symmetric matrix. In Section 5, the conjugate-gradient algorithm is derived by reasoning on (24); in Section 6, the variable-metric algorithm is derived by reasoning on (25).

## 5. Conjugate-Gradient Method

In this section, we consider the algorithm

$$\tilde{x} = x + \Delta x, \quad \Delta x = -\alpha p, \quad p = g(x) + q \qquad (26)$$

where $q$ is an $n$-vector to be specified. The first variation of the function (6) is given by Eq. (10) which, in the light of (26), becomes

$$\delta f(x) = -\alpha[g^T(x)g(x) + g^T(x)q] \qquad (27)$$

We note that $g^T(x)g(x) > 0$. Therefore, for $\alpha > 0$, the descent property of this algorithm is ensured if one chooses $q$ so that

$$g^T(x)q = 0 \qquad (28)$$

If Eqs. (26-1) and (26-2) are combined, the position vector at the end of any iteration becomes

$$\tilde{x} = x - \alpha p \qquad (29)$$

For a given point $x$ and a given vector $p$, Eq. (29) defines a one-parameter family of points $\tilde{x}$ for which the function $f$ takes the form

$$f(\tilde{x}) = f(x - \alpha p) = F(\alpha) \qquad (30)$$

The greatest decrease in the function $F(\alpha)$ occurs if the parameter $\alpha$ satisfies the following necessary condition:

$$F_\alpha = 0 \qquad (31)$$

On account ... ), the following relation holds:

$$F_\alpha = -g^T(\bar{x})p \qquad (32)$$

Therefore, Eq. (11) becomes

$$g^T(\bar{x})p = 0 \qquad (33)$$

and shows that the gradient $g(\bar{x})$ is orthogonal to the search direction p.

Next, we apply Eq. (33) to the previous iteration and obtain

$$g^T(x)\bar{p} = 0 \qquad (34)$$

By comparing (28) and (34), we conclude that one possible choice of the vector q is the following:

$$q = \gamma\bar{p} \qquad (35)$$

where $\gamma$ is a constant. As a consequence, the algorithm (28) can be rewritten as

$$\bar{x} = x + \delta x, \quad \delta x = -\alpha p, \quad p = g(x) + \gamma\bar{p} \qquad (36)$$

The next step is to determine the constant $\gamma$. If Eqs. (36) are combined, the position vector at the end of any iteration becomes

$$\bar{x} = x - \alpha g(x) - \alpha\gamma\bar{p} \qquad (37)$$

For a given point x and a given vector $\bar{p}$, Eq. (37) defines a two-parameter family of points $\bar{x}$ for which the function f takes the form

$$f(\bar{x}) = f(x - \alpha g(x) - \alpha\gamma\bar{p}) = F(\alpha, \gamma) \qquad (38)$$

The greatest decrease in the function $F(\alpha, \gamma)$ occurs if the parameters ... satisfy the following necessary conditions:

$$F_\alpha = 0, \quad F_\gamma = 0 \qquad (39)$$

On account of (36-3) and (38), the following relations hold:

$$F_\alpha = -g^T(\bar{x})p, \quad F_\gamma = -\alpha g^T(\bar{x})\bar{p} \qquad (40)$$

Therefore, Eqs. (39) become

$$g^T(\bar{x})p = 0, \quad g^T(\bar{x})\bar{p} = 0 \qquad (41)$$

and show that the gradient $g(\bar{x})$ is orthogonal to the search directions p and $\bar{p}$. A mathematical consequence of Eqs. (36-3) and (41) is that

$$g^T(\bar{x})g(x) = 0 \qquad (42)$$

showing that the gradients $g(\bar{x})$ and $g(x)$ are orthogonal.

5.1. Quadratic Function. Now, consider the particular case of a quadratic function, that is, a function of the form

$$f(x) = a + b^T x + \frac{1}{2}x^T H x \qquad (43)$$

where a is a constant scalar, b is a constant n-vector, and H is a constant, symmetric n x n matrix. For this function, the gradient is a linear function of x, that is,

$$g(x) = b + Hx \qquad (44)$$

$$g(\bar{x}) = b + H\bar{x} \tag{45}$$

relations (44) and (45) imply that

$$g(\bar{x}) = g(x) + H\Delta x = g(x) - \alpha H p \tag{46}$$

Next, we introduce Eqs. (46) into (41) and, after laborious manipulations, obtain the solutions (Ref. 1)

$$\alpha = \frac{g^T(x)g(x)}{p^T H p} \quad , \quad \gamma = \frac{g^T(\bar{x})g(\bar{x})}{g^T(x)g(x)} \tag{47}$$

where p is given by (36-3).

For a quadratic function, Hestenes and Stiefel (Ref. 1) proved that, if the first step of the descent process is a gradient step, the following relations hold:

$$g^T(x)g(x_*) = 0 \quad , \quad g^T(x)p_* = 0 \quad , \quad p^T H p_* = 0 \tag{48}$$

where $x_*$ denotes any state preceding x. Equation (48-1) states that the gradient at each iteration is orthogonal to the gradient at every previous iteration. Equation (48-2) states that the gradient at each iteration is orthogonal to the search direction at every previous iteration. Finally, Eq. (48-3) states that the search direction at each iteration and the search direction at every previous iteration are conjugate with respect to the constant matrix H; this is why the algorithm is called the conjugate-gradient method. The algorithm (36) with $\alpha$ and $\gamma$ defined by (17) reduces the gradient to zero in no more than n steps; therefore, the minimum of f(x) is reached in no more than n steps.

3.2. Nonquadratic Function. For a nonquadratic function, solving Eqs. (41) for $\alpha$ and $\gamma$ requires a two-dimensional search (Ref. 6). The difficulty of this process can be avoided if one optimizes $\alpha$ exactly and uses an approximate value for $\gamma$, namely, that given by Eq. (47-2). This leads to the algorithm (Ref. 3)

$$\bar{x} = x + \delta x, \quad \delta x = -\alpha p \quad , \quad p = g(x) + \frac{g^T(x)g(x)}{g^T(\hat{x})g(\hat{x})} p \tag{49}$$

in which $\alpha$ is optimized by searching for the minimum of f along the direction defined by (49). Theoretically, therefore, the optimization of $\alpha$ requires that the relation (44-1) be satisfied.

For any iteration except the first, the complete algorithm can be stated as follows: (a) for a given nominal point x, the gradient g(x) is known; since the gradient g(x̂) and the search direction p̂ are known from the previous iteration, the search direction p can be determined with Eq. (49-3); (b) the optimum stepsize $\alpha$ must be determined by minimizing the function f along the search direction p, as in Section 7; (c) the correction $\delta x$ to the position vector x is determined using Eq. (49-2); and (d) the new position vector x̄ is computed through Eq. (49-1). Next, the position vector x̄ becomes the nominal point for the subsequent iteration, and the procedure is repeated until a predetermined stopping condition is satisfied (see Section 8). To start the algorithm, one bypasses (49-1) and sets p = g(x), equivalent to stating that the first step is a gradient step.

In closing, the following comments are pertinent: (a) in the conjugate-gradient method it is important that the stepsize $\alpha$ be determined accurately, while this is not the case with the ordinary gradient method; (b) theoretical considerations and numerical experience show the desirability of restarting the process every n or n + 1 iterations, that is, resetting p = g(x) every n or n + 1 iterations (Ref. 3).

## 6. Variable Metric Algorithm

In this section, we consider the algorithm

$$\hat{x} = x + \delta x, \qquad \delta x = -\alpha p, \qquad p = Ag(x) \tag{50}$$

where A is a symmetric $n \times n$ matrix to be specified. The first variation of the function (6) is given by Eq. (10) which, in the light of (50), becomes

$$\delta f(x) = -\alpha g^T(x) A g(x) \tag{51}$$

We note that, if the matrix A is positive definite, $g^T(x) A g(x) > 0$. Therefore, for $\alpha > 0$, the descent property of this algorithm is ensured.

Now, consider the points $\hat{x}$ and $x$. At point $\hat{x}$, the gradient $g(\hat{x})$ and the matrix $\hat{A}$ are known; at point $x$, the gradient $g(x)$ is known. Therefore, the differences

$$\delta \hat{x} = x - \hat{x}, \qquad \delta \hat{g} = g(x) - g(\hat{x}) \tag{52}$$

are available. We wish to determine the matrix A so that the relation

$$A \delta \hat{g} = \delta \hat{x} \tag{53}$$

is satisfied. After defining the matrix difference

$$\delta \hat{A} = A - \hat{A} \tag{54}$$

we combine (53)-(54) to obtain

$$\delta \hat{A} \delta \hat{g} = \delta \hat{x} - \hat{A} \delta \hat{g} \tag{55}$$

In which the only unknown is $\delta \hat{A}$. Equation (55) admits the solution

$$\delta \hat{A} = \frac{\delta \hat{x} y^T}{y^T \delta \hat{g}} - \frac{\hat{A} \delta \hat{g} z^T}{z^T \delta \hat{g}} \tag{56}$$

where y and z denote arbitrary n-vectors. Therefore, the matrix A must be updated according to the relation

$$A = \hat{A} + \frac{\delta \hat{x} y^T}{y^T \delta \hat{g}} - \frac{\hat{A} \delta \hat{g} z^T}{z^T \delta \hat{g}} \tag{57}$$

In particular, if one chooses

$$y = \delta \hat{x}, \qquad z = \hat{A} \delta \hat{g}$$

Eq. (57) becomes

$$A = \hat{A} + \frac{\delta \hat{x} \delta \hat{x}^T}{\delta \hat{x}^T \delta \hat{g}} - \frac{\hat{A} \delta \hat{g} \delta \hat{g}^T \hat{A}}{\delta \hat{g}^T \hat{A} \delta \hat{g}} \tag{58}$$

Note that the second and third matrices on the right-hand side of (58) are symmetric; therefore, if $\hat{A}$ is symmetric, A is also symmetric.

**6.1. Quadratic Function.** Now, consider the particular case of a function having the form (43). For this quadratic function, the following properties can be shown to hold (Refs. 4-5 and 7-8):

(a) If the initial matrix A is chosen to be the inverse of the second derivative matrix H, that is, if

$$A = H^{-1} \tag{59}$$

at the initial point, the variable-metric algorithm exhibits one-step convergence. This is due to the fact that the variable-metric algorithm becomes identical with quasilinearization.

(b) If the initial matrix A is chosen to be positive definite, any subsequent matrix A is also positive definite. With this understanding, the following relations hold:

$$g^T(x)p_* = 0, \qquad p^T H p_* = 0 \tag{61}$$

Equation (61-1) states that the gradient at each iteration is orthogonal to the search direction at every previous iteration. Equation (61-2) states that the search direction at each iteration and the search direction at every previous iteration are conjugate with respect to the constant matrix H. As the algorithm progresses, the matrix A tends to the inverse of the second derivative matrix H, and relation (60) becomes satisfied exactly when convergence is achieved. The algorithm (50), with A updated according to (59), reduces the gradient to zero in no more than n steps; therefore, the minimum of f(x) is reached in no more than n steps.

(c) As a particular case of (b), the initial matrix can be chosen to be

$$A = I \tag{62}$$

where I is the n x n identity matrix. Under these conditions, the variable-metric algorithm becomes identical with the conjugate-gradient algorithm of Section 5.

6.2. Nonquadratic Function. For a nonquadratic function, the variable-metric algorithm is represented by

$$\tilde{x} = x + \delta x, \quad \delta x = -\alpha p, \quad p = Ag(x) \tag{63}$$

with

$$A = \hat{A} + \frac{\delta \hat{x} \delta \hat{x}^T}{\delta \hat{x}^T \delta \hat{g}} - \frac{\hat{A} \delta \hat{g} \delta \hat{g}^T \hat{A}}{\delta \hat{g}^T \hat{A} \delta \hat{g}} \tag{64}$$

The stepsize α is to be optimized by searching for the minimum of f along the direction defined by (63).

For any iteration except the first, the complete algorithm can be stated as follows: (a) for a given nominal point x, the gradient g(x) is known; since g(x̂), Δx̂, Â are known from the previous iteration, the matrix A can be computed with (64) and the search direction p with (63-3); (b) the optimum stepsize α must be determined by minimizing the function f along the search direction p, as in Section 7; (c) the correction δx to the position vector x is determined using Eq. (63-2); and (d) the new position vector x̃ is computed through Eq. (63-1). Next, the new position vector x̃ becomes the nominal point for the subsequent iteration and the procedure is repeated until a predetermined stopping condition is satisfied (see Section 8). To start the algorithm, one bypasses (64) and sets A equal to any symmetric, positive-definite matrix (for instance, the identity matrix).

In closing, the following comments are pertinent: (a) in the variable-metric method, it is important that the stepsize α be determined accurately; (b) restarting the algorithm every n or n + 1 iterations is not necessary; however, restarting is indispensable whenever the positive-definiteness of the matrix A is violated, for example, if the stepsize α becomes negative.

## 7.   Search Technique

In each of the previous methods, the stepsize α must be optimized. In this section,

we present techniques to solve the equation

$$F_\alpha(\eta) = 0 \qquad (65)$$

that is, to find the minimum of the function $F(\alpha)$ given by Eq. (19) for the ordinary gradient

algorithm or Eq. (30) for the conjugate-gradient and variable-metric algorithms. Since

the techniques in question involve the consideration of the first derivative $F_\alpha$ and perhaps

the second derivative $F_{\alpha\alpha}$, we summarize these derivatives below.

For all of the previous methods, we have

$$F_\alpha(\eta) = -g^T(\tilde{x})p , \quad F_{\alpha\alpha}(\eta) = p^T H(\tilde{x})p \qquad (66)$$

where

$$\tilde{x} = x - \alpha p \qquad (67)$$

The search direction is given by $p = g(x)$ for the ordinary gradient method, Eq. (49-3)

for the conjugate-gradient method, and Eq. (63-3) for the variable-metric method. Of

course, Eq. (66-2) requires that the second-derivative matrix $H(x)$ be explicitly available.

If this is not the case, one can use the difference scheme

$$F_{\alpha\alpha}(\eta) = (1/2\eta)[F_\alpha(\eta + \theta) - F_\alpha(\eta - \theta)] \qquad (68)$$

$$= (1/2\eta)[g(\tilde{x} + \theta p) - g(\tilde{x} - \theta p)]^T p$$

In practice, one may choose

$$\theta = \epsilon_1 / |p| \qquad (69)$$

where $\epsilon_1$ is a small number.

7.1.   Cubic Interpolation.   Let the values of the function $F(\alpha)$ and derivative

$F_\alpha(\alpha)$ be computed for two different values of α, namely, $\alpha_1$ and $\alpha_2$, with $\alpha_2 > \alpha_1 > 0$.

If $\alpha_1$ and $\alpha_2$ are such that

$$F_\alpha(\alpha_1) < 0 , \quad F_\alpha(\alpha_2) > 0$$

then the minimum of the function $F(\alpha)$ occurs for some value α in the range

$$\alpha_1 < \alpha < \alpha_2 \qquad (71)$$

In this range, we represent the function $F(\alpha)$ with the cubic

$$F(\alpha) = A + B(\alpha - \alpha_1) + C(\alpha - \alpha_1)^2 + D(\alpha - \alpha_1)^3 \qquad (72)$$

whose first and second derivatives are given by

$$F_\alpha(\alpha) = B + 2C(\alpha - \alpha_1) + 3D(\alpha - \alpha_1)^2 , \quad F_{\alpha\alpha}(\alpha) = 2C + 6D(\alpha - \alpha_1) \qquad (73)$$

The scalar coefficients A, B, C, D are determined by requiring (72) to match the ordinate

and the slope of the curve $F(\alpha)$ at $\alpha_1$ and $\alpha_2$. Therefore, one has to solve the linear

equations

$$F(\alpha_1) = A , \quad F(\alpha_2) = A + B(\alpha_2 - \alpha_1) + C(\alpha_2 - \alpha_1)^2 + D(\alpha_2 - \alpha_1)^3 ,$$
$$F_\alpha(\alpha_1) = B , \quad F_\alpha(\alpha_2) = B + 2C(\alpha_2 - \alpha_1) + 3D(\alpha_2 - \alpha_1)^2 \qquad (74)$$

Once the coefficients of the cubic (72) are known, the optimum value of α is determined by

the condition (65). Therefore, in the light of (73), one arrives at the solution

$$\alpha = \alpha_1 + (1/3D)[-C + \sqrt{(C^2 - 3BD)}] \qquad (75)$$

At this point... recomputes the function $F(\alpha)$ and the derivative $F_\alpha(\alpha)$. Then, the process is iterated until a predetermined stopping condition is satisfied. For instance, one may require that

$$|F_\alpha(\alpha)| < \epsilon_2 \tag{76}$$

or that

$$|F_\alpha(\alpha)| \leq \epsilon_3 |F_\alpha(0)| \tag{77}$$

where $\epsilon_2$ and $\epsilon_3$ are prescribed small numbers.

7.2. Quasilinearization. An alternate technique for computing the optimum stepsize, that of quasilinearization with built-in safeguards to ensure that the function decreases at every step of the iterative search, is now presented. Let

$$\delta\alpha = \alpha - \alpha_0 \tag{78}$$

denote the correction to $\alpha$ starting from an arbitrary nominal value $\alpha_0$. If quasilinearization is applied to Eq. (65), one obtains the linear algebraic equation

$$F_{\alpha\alpha}(\alpha_0)\delta\alpha + F_\alpha(\alpha_0) = 0 \tag{79}$$

Next, we imbed Eq. (79) in the more general equation

$$F_{\alpha\alpha}(\alpha_0)\delta\alpha + u s F_\alpha(\alpha_0) = 0 \tag{80}$$

where u denotes a scaling factor and s a direction factor such that

$$0 < u \leq 1 \;,\; s = \pm 1 \tag{81}$$

Equation (80) admits the solution

$$\delta\alpha = -u s F_\alpha(\alpha_0)/F_{\alpha\alpha}(\alpha_0)$$

The direction factor s is determined in such a way that the first variation

$$\delta F(\alpha_0) = F_\alpha(\alpha_0)\delta\alpha \tag{83}$$

is negative. From (82)-(83), we obtain

$$\delta F(\alpha_0) = -u s F_\alpha^2(\alpha_0)/F_{\alpha\alpha}(\alpha_0) \tag{84}$$

Therefore, $\delta F(\alpha_0)$ is negative if the direction factor s is chosen as follows:

$$s = \text{sign } F_{\alpha\alpha}(\alpha_0)$$

Because of this choice, the correction (82) becomes

$$\delta\alpha = -u F_\alpha(\alpha_0)/|F_{\alpha\alpha}(\alpha_0)|$$

To perform the search, a nominal value must be given to $\alpha_0$. Then, one sets $u = 1$, computes $\delta\alpha$ from Eq. (86) and $\alpha$ from Eq. (78). If $F(\alpha) < F(\alpha_0)$, the scaling factor $u = 1$ is acceptable. If $F(\alpha) > F(\alpha_0)$, the previous value of u must be replaced by some smaller value in the range $0 < u \leq 1$ until the condition $F(\alpha) < F(\alpha_0)$ is met; this can be obtained through bisection, that is, by successively dividing the value of u by 2. At this point, the search step is completed. The value obtained for $\alpha$ becomes the nominal value $\alpha_0$ for the next search step, and the procedure is repeated until a desired degree of accuracy is obtained, that is, until Eqs. (76) or (77) is satisfied. In the absence of better information, the first step of the search procedure can be made with $\alpha_0 = 0$.

## 8. Termination of the Algorithm

One way to terminate the gradient algorithm is to impose a condition on the modulus of the gradient, for example,

$$g^T(x)g(x) \le \epsilon_4 \tag{87}$$

where $\epsilon_4$ is a prescribed small number. If the function $f(x)$ is rather flat in the neighborhood of the minimum, then Eq. (87) may not yield precise coordinates. In this case, the following additional condition is suggested:

$$\Delta x^T \Delta x \le \epsilon_5 \tag{88}$$

where $\epsilon_5$ is a prescribed small number.

## References

1. HESTENES, M.R., and STIEFEL, E., Methods of Conjugate Gradients for Solving Linear Systems, Journal of Research of the National Bureau of Standards, Vol. 49, No. 6, 1952.

2. BECKMAN, F.S., The Solution of Linear Equations by the Conjugate Gradient Method, Mathematical Methods for Digital Computers, Vol. 1, Edited by A. Ralston and H.S. Wilf, John Wiley and Sons, New York, 1960.

3. FLETCHER, R., and REEVES, C.M., Function Minimization by Conjugate Gradients, Computer Journal, Vol. 7, No. 2, 1964.

4. DAVIDON, W.C., Variable-Metric Method for Minimization, Argonne National Laboratory, Report No. ANL-5990, 1959.

5. FLETCHER, R., and POWELL, M.J.D., A Rapidly Convergent Descent Method for Minimization, Computer Journal, Vol. 6, No. 2, 1963.

6. MIELE, A., and CANTRELL, J.W., Gradient Methods in Mathematical Programming, Part 2, Memory Gradient Method, Rice University, Aero-Astronautics Report No. 56, 1969.

7. MYERS, G.E., Properties of the Conjugate Gradient and Davidon Methods, Journal of Optimization Theory and Applications, Vol. 2, No. 4, 1968.

8. PEARSON, J.D., On Variable Metric Methods of Minimization, Research Analysis Corporation, Technical Paper No. RAC-TP-303, 1968.

# UPDATING RULES FOR THE PENALTY CONSTANT USED IN THE PENALTY FUNCTION METHOD FOR MATHEMATICAL PROGRAMMING PROBLEMS

A. MIELE, G. M. COGGINS, and A. V. LEVY

Abstract. The problem of minimizing a function $f(x)$ subject to a constraint $\varphi(x)=0$ is considered, where $f$ is a scalar, $x$ an $n$-vector, and $\varphi$ a $q$-vector, with $q < n$. The penalty function method is investigated; that is, a sequence of unconstrained minimization problems is solved, each of which involves the penalty function $U(x, k) = f(x) + kP(x)$. Here, the penalty constant $k$ is a positive, scalar quantity, and $P(x) = \varphi^T(x)\varphi(x)$ is the norm squared of the constraint error.

Crucial to the penalty function method is the prediction of the rate $a = k_2/k_1$ at which the penalty constant must be increased when shifting from one cycle of the algorithm to the next. Here, $k_1$ denotes the penalty constant of the present cycle, and $k_2$ denotes the penalty constant of the next cycle.

In this paper, two variable-rate updating rules are developed: (i) the updating rule $a = \beta(\overline{P_1/P_*})$ and (ii) the updating rule $a = \beta(P_1/P_*)$, where $P_1$ is the constraint error at the end of the present cycle, $P_*$ is the constraint error allowed for convergence, and $\beta > 1$. Updating rule (i) tends to produce at the end of the next cycle a constraint error $P_2$ below the geometric mean of the constraint error at the end of the present cycle and that allowed for convergence; updating rule (ii) tends to produce at the end of the next cycle a constraint error $P_2$ below that allowed for convergence.

In order to evaluate these updating rules, six numerical examples are investigated. The first example deals with a quadratic function subject to linear constraints; the remaining examples deal with nonquadratic functions subject to nonlinear constraints. Each example is solved with the conjugate gradient algorithm for five starting values of the

penalty constant $k_0$, ranging between $10^{-1}$ and $10^3$. For each example, the variable-rate updating rules (i) and (ii), with $N = 10$, are compared with the constant-rate updating rules $x = 5$ and $x = 10$. From the numerical experiments, it is concluded that the variable-rate updating rules are superior to the constant-rate updating rules, in that they generally lead to convergence in a smaller number of iterations.

## 1. Introduction

Over the past several years, considerable work has been done on the problem of minimizing a function $f(x)$ subject to a constraint $\varphi(x) = 0$ using numerical methods. Here, $f$ is a scalar, $x$ an $n$-vector, and $\varphi$ a $q$-vector, with $q < n$.

The approaches employed are generally based on one of two basic ideas. One is to develop algorithms such that the constraints are satisfied, at least to first order, at the end of each iteration (Refs. 2-4). Another is to develop algorithms involving cycles in which the vector $x$ is viewed as unconstrained and a new function related to $f(x)$ and $\varphi(x)$ is minimized. The penalty function method (Refs. 5-8) is an approach of the latter type.

Crucial to the penalty function method is the prediction of the rate at which the penalty constant must be increased when shifting from one cycle of the algorithm to the next. This key question is considered in this paper, whose objective is the following: to improve the convergence characteristics of the penalty function method by automatically adjusting the penalty constant used in the method.

## 2. Statement of the Problem

We consider the problem of minimizing the function

$$f = f(x) \tag{1}$$

subject to the constraint

$$\varphi(x) = 0, \tag{2}$$

where $f$ is a scalar, $x$ an $n$-vector, and $\varphi$ a $q$-vector, with $q < n$. Here, all vectors are column vectors. It is assumed that the first and second partial derivatives of the functions $f(x)$ and $\varphi(x)$ exist and are continuous and that the constrained minimum exists.

2.1. *First-Order Conditions*. From theory of maxima and minima, it is known that the above problem is equivalent to that of minimizing the

augmented function

$$F(x, \lambda) = f(x) + \lambda^T \varphi(x) \tag{3}$$

subject to the constraint (2). Here, the $q$-vector $\lambda$ is the Lagrange multiplier and the superscript $T$ denotes the transpose of a matrix. If

$$F_x(x, \lambda) = f_x(x) + \varphi_x(x)\lambda \tag{4}$$

denotes the gradient of the augmented function, the optimum solution for $x$ and $\lambda$ must satisfy the relations

$$\varphi(x) = 0, \quad F_x(x, \lambda) = 0, \tag{5}$$

which are a system of $n + q$ equations in the $n + q$ components of $x$, $\lambda$. In Eqs. (4)-(5), the gradients $f_x$ and $F_x$ denote $n$-vectors and the matrix $\varphi_x$ is $n \times q$.

2.2. *Approximate Solutions*. Since the system (5) is generally nonlinear, approximate methods must be employed. In this connection, we introduce here the scalar performance indexes

$$P(x) = \varphi^T(x)\varphi(x), \quad Q(x, \lambda) = F_x^T(x, \lambda) F_x(x, \lambda), \tag{6}$$

which measure the errors in the constraint and the optimum condition, respectively. Then, we observe that $P = 0$ and $Q = 0$ for the optimum solution, while $P > 0$ and/or $Q > 0$ for any approximation to the solution. When approximate methods are used, they must ultimately lead to values of $x$ and $\lambda$ such that

$$P(x) \leq P_0, \quad Q(x, \lambda) \leq Q_0. \tag{7}$$

Alternatively, (7) can be replaced by

$$R(x, \lambda) \leq R_0, \tag{8}$$

where

$$R(x, \lambda) = P(x) + Q(x, \lambda)$$

denotes the cumulative error in the constraint and the optimum condition. In Eqs. (7)-(8), $P_0$, $Q_0$, $R_0$ are small, preselected numbers. Note that, if one chooses $P_0 = Q_0 = R_0$, satisfaction of Ineq. (8) implies satisfaction of Ineqs. (7).

## 3. Penalty Function Method

The penalty function method is based on the construction of a sequence of special functions having, in the limit, an unconstrained minimum point coincident with the solution of the original constrained minimization problem. Specifically, the penalty function is defined by

$$U(x, k) = f(x) + k P(x) = f(x) + k \varphi^T(x) \varphi(x) \tag{10}$$

and is obtained by adding to the function $f(x)$ a term quadratic in the constraint $\varphi(x)$, where $k > 0$ is the penalty constant.

The problem of minimizing the function (1) subject to the constraint (2) is replaced by a sequence of unconstrained minimization problems. In each element of the sequence or cycle, one minimizes the penalty function (10) with respect to $x$ for given $k$. Therefore, theoretically speaking, the following necessary condition must be satisfied at the end of a cycle:

$$U_x(x, k) = f_x(x) + k P_x(x) = f_x(x) + 2k \varphi_x(x) \varphi(x) = 0. \tag{11}$$

If one defines the Lagrange multiplier to be

$$\lambda = 2k \varphi(x), \tag{12}$$

Eq. (11) can be rewritten as

$$F_x(x, \lambda) = f_x(x) + \varphi_x(x) \lambda = 0, \tag{13}$$

meaning that the combination of $x$ and $\lambda$ obtained at the end of a cycle satisfies exactly the optimum condition (5-2). However, if the penalty constant $k$ is arbitrary, the vector $x$ which satisfies Eq. (11) generally violates the constraint condition (5-1).

In order to obtain constraint satisfaction, increasingly larger values of the penalty constant must be employed in successive cycles of the penalty function method. In this connection, let $k_1$ denote the penalty constant of the present cycle and $k_2$ denote the penalty constant of the next cycle, with $k_2 > k_1$. Because of the jump in $k$, the penalty function increases by the amount

$$U(x, k_2) - U(x, k_1) = (k_2 - k_1) P(x), \tag{14}$$

and the norm squared of the gradient of the penalty function takes on the value [1]

$$U_x^T(x, k_2) U_x(x, k_1) = (k_2 - k_1)^2 P_x^T(x) P_x(x) = (n - 1)^2 f_x^T(x) f_x(x), \tag{15}$$

where

$$P(x) = \varphi^T(x) \varphi(x), \qquad P_x(x) = 2\varphi_x(x) \varphi(x). \tag{16}$$

The positiveness of the right-hand side of Eq. (14) is the key to the mechanism on which the penalty function method is based.

After a sufficient number of cycles, the constraint error can be made as small as desired providing the penalty constant has become sufficiently large. Theoretically speaking, the condition $\varphi(x) = 0$ is desired at convergence; consequently, the multiplier $\lambda$ defined by Eq. (12) can be identical with the multiplier satisfying Eqs. (5), which is generally nonzero, only if $k \to \infty$. In a practical digital computer, this means that large values of $k$ are needed at convergence.

3.1. *Numerical Implementation.* From the above considerations, the following outline of the penalty function method emerges.

(a) The original constrained minimization problem is replaced by a sequence of unconstrained minimization problems.

(b) In each element of the sequence or cycle, the penalty function

$$U(x, k) = f(x) + k \varphi^T(x) \varphi(x) \tag{17}$$

is minimized with respect to $x$ for given $k$. The minimum of $U(x, k)$ is achieved when the following stopping condition is satisfied:

$$U_x^T(x, k) U_x(x, k) \leq Q_s, \tag{18}$$

where $Q_s$ is a small, preselected number.

(c) Upon termination of a cycle, one checks the following inequality:

$$\varphi^T(x) \varphi(x) + U_x^T(x, k) U_x(x, k) \leq R_s, \tag{19}$$

where $R_s$ is a small, preselected number. If Ineq. (19) is satisfied, the algorithm is terminated. If Ineq. (19) violated, the algorithm is continued by choosing the solution point of the present cycle as the starting point of the next cycle.

---

[1] Note that $U_x^T(x, k_1) U_x(x, k_1) = 0.$

(d) For the next cycle, a higher value of the penalty constant is selected, specifically,

$$k_2 = \pi k_1,\tag{20}$$

where $\pi > 1$ is the penalty constant ratio. After updating the penalty constant, one returns to (b) and continues iteratively.

REMARK 3.1. At convergence of a cycle, the stopping condition (18) can be written as

$$Q(x, \lambda) \leq Q_*,\tag{21}$$

where $\lambda$ is given by Eq. (12). Analogously, at convergence of the algorithm, the stopping condition (19) becomes

$$R(x, \lambda) = P(x) + Q(x, \lambda) \leq R_*,\tag{22}$$

where $\lambda$ is given by Eq. (12).

## 4. Updating Rules

Crucial to the penalty function method is the prediction of the rate

$$\pi = k_2/k_1\tag{23}$$

at which the penalty constant must be increased when shifting from one cycle of the algorithm to the next. Here, $k_1$ denotes the penalty constant of the present cycle, and $k_2$ denotes the penalty constant of the next cycle.

4.1. *Standard Technique.* The standard method (see, for instance, Ref. 7) consists of employing a constant value of $\pi$ throughout the algorithm, for example,

$$\pi = 5 \quad \text{or} \quad \pi = 10.\tag{24}$$

4.2. *Proposed Technique.* An alternate method consists of employing a variable value of $\pi$ throughout the algorithm. For instance, one might select $\pi$ so as to achieve a predetermined constraint error at the end of the next cycle of the algorithm.

Let the following definition be introduced:

$$Z = \lambda^T \lambda,\tag{25}$$

with $\lambda$ given by Eq. (12). Let Eqs. (12) and (25) be applied twice, once at the end of the present cycle (subscript 1) and once at the end of the next cycle (subscript 2). One obtains the relations

$$Z_1 = 4 k_1^2 P_1, \quad Z_2 = 4 k_2^2 P_2,\tag{26}$$

which imply that

$$\pi = \sqrt{(Z_2 P_1 / Z_1 P_2)}\tag{27}$$

Now, we introduce the basic assumption [*]

$$Z_1 = Z_2\tag{28}$$

that is, we neglect the change in the norm squared of the multiplier between the end of the present cycle and the end of the next cycle. With this understanding, Eq. (27) simplifies to

$$\pi = \sqrt{(P_1/P_2)}.$$

This relation enables one to predict the penalty constant needed for the next cycle in order to achieve a predetermined constraint error $P_2$.

As an example, assume that the expected constraint error at the end of the next cycle $P_2$ is below the geometric mean of the constraint error at the end of the present cycle $P_1$ and that allowed for convergence $P_*$; that is, assume that

$$P_2 = \sqrt{(P_1 P_*)}/\beta,\tag{30}$$

where $\beta > 1$. Then, Eq. (29) simplifies to

$$\pi = \beta \sqrt{(P_1/P_*)}.\tag{31}$$

As another example, assume that the expected constraint error at the end of the next cycle $P_2$ is below that allowed for convergence $P_*$; that is, assume that

$$P_2 = P_*/\beta,\tag{32}$$

where $\beta > 1$. Then, Eq. (29) simplifies to

$$\pi = \beta \sqrt{(P_1/P_*)}.\tag{33}$$

[*] Hypothesis (28) is key to this paper and is supported by the data exhibited in Tables 1-4.

## 5. Experimental Conditions

In order to evaluate the theory, six numerical examples were solved using a Burroughs B-5500 computer, double-precision arithmetic, and FORTRAN-IV program. Within each cycle of the penalty function method, the conjugate-gradient algorithm was employed (Ref. 1).

*Starting Point.* In all the examples, the nominal point chosen to start the penalty function method was defined by

$$x_1 = x_2 = \dots = x_n = 2, \tag{34}$$

where $x_1, x_2, \dots, x_n$ denote the components of the vector $x$.

*Starting Penalty Constant.* For the first cycle of the penalty function method, five values of the penalty constant were employed, namely,

$$k_1 = 10^{-1}, \quad k_2 = 10^{-1}, \quad k_3 = 10^0, \quad k_4 = 10^1, \quad k_5 = 10^2. \tag{35}$$

*Updating Rules.* For subsequent cycles, the penalty constant was determined in accordance with one of the following penalty constant ratios:

and

$$n = 5, \quad n = 10, \tag{36}$$

with

$$n = \beta \sqrt{(P_1/P_0)}, \quad n = \beta \sqrt{(P_1/P_0)}, \tag{37}$$

$$\beta = 10, \quad P_0 = 10^{-12}. \tag{38}$$

The penalty constant ratios (36) are representative of the standard method (see, for instance, Ref. 7) and the penalty constant ratios (37) are representative of the new method proposed here.

*Definition of Convergence.* Convergence of a cycle was defined as

$$Q(x, k) = U_x^T(x, k) U_x(x, k) \leq 10^{-11}, \tag{39}$$

and convergence of the complete algorithm was defined as

$$R(x, k) = P(x) + Q(x, k) = \varphi^T(x) \varphi(x) + U_x^T(x, k) U_x(x, k) \leq 10^{-11}, \tag{40}$$

where the multiplier $\lambda$ is given by Eq. (12).

*Search Technique.* Within each iteration of the conjugate gradient algorithm, let $x$ denote the nominal point, $\tilde{x}$ the varied point, and $p$ the search direction. Let $\tilde{U}(x)$ denote the function representing the behavior of the penalty function along the search direction, that is,

$$\tilde{U} = U(\tilde{x}, k) = U(x + \Delta x, k) = U(x - \alpha p, k) = \tilde{U}(\alpha). \tag{41}$$

This function was employed in order to determine the optimum stepsize at each iteration. Specifically, a precise one-dimensional search on the function $\tilde{U}(\alpha)$ was conducted such that, in any given step, the inequality (*)

$$\tilde{U}(\alpha) < \tilde{U}(\alpha_0) \tag{42}$$

was satisfied, where $\alpha_0$ is the nominal stepsize and $\alpha$ is the varied stepsize. The search was terminated when the stopping condition

$$\tilde{U}_\alpha^2(\alpha) \leq \tilde{U}_\alpha^2(0) \times 10^{-8} \tag{43}$$

was satisfied.

## 6. Numerical Examples

In this section, six numerical examples are described. The first example pertains to a quadratic function subject to linear constraints. The remaining examples pertain to nonquadratic functions subject to nonlinear constraints.

EXAMPLE 6.1. Consider the problem of minimizing the function

$$f = (x_1 - x_2)^2 + (x_2 + x_3 - 2)^2 + (x_4 - 1)^2 + (x_5 - 1)^2 \tag{44}$$

subject to the constraints

$$x_1 + 3x_2 = 0, \quad x_3 + x_4 - 2x_5 = 0, \quad x_2 - x_5 = 0. \tag{45}$$

___

(*) Due to the analytical nature of the examples considered here, analytical numerical quasilinearization was employed in order to ensure satisfaction of Ineq. (43). However, in a more realistic situation, one would use quadratic interpolation or cubic interpolation. For the details of the conjugate-gradient algorithm and the search technique employed, the reader should consult Ref. 1.

This function admits the relative minimum $f = 4.0030$ at the point defined by

$$x_1 = -0.7674, \quad x_2 = 0.2558, \quad x_3 = 0.0279, \quad x_4 = -0.1162, \quad x_5 = 0.2558 \quad (46)$$

and

$$\lambda_1 = 2.0165, \quad \lambda_2 = 3.2325, \quad \lambda_3 = -5.0534. \quad (47)$$

EXAMPLE 6.2. Consider the problem of minimizing the function

$$f = (x_1 - 1)^2 + (x_1 - x_2)^2 + (x_2 - x_3)^4 \quad (48)$$

subject to the constraint

$$x_1(1 + x_2^2) + x_3^4 - 4 - 3\sqrt{2} = 0. \quad (49)$$

This function admits the relative minimum $f = 0.3256 \times 10^{-1}$ at the point defined by

$$x_1 = 1.1048, \quad x_2 = 1.1966, \quad x_3 = 1.5352 \quad (50)$$

and

$$\lambda_1 = -0.1072 \times 10^{-1}. \quad (51)$$

EXAMPLE 6.3. Consider the problem of minimizing the function

$$f = (x_1 - 1)^2 + (x_1 - x_2)^2 + (x_2 - x_3)^2 + (x_3 - x_4)^4 + (x_4 - x_5)^4 \quad (52)$$

subject to the constraints

$$x_1 + x_2^2 + x_3^3 - 2 - 3\sqrt{2} = 0, \quad x_2 - x_3^2 + x_4 + 2 - 2\sqrt{2} = 0,$$
$$x_1 x_5 - 2 = 0. \quad (53)$$

The function admits the relative minimum $f = 0.7877 \times 10^{-1}$ at the point defined by

$$x_1 = 1.1911, \quad x_2 = 1.3626, \quad x_3 = 1.4728, \quad x_4 = 1.6350, \quad x_5 = 1.6790 \quad (54)$$

and

$$\lambda_1 = -0.3482 \times 10^{-1}, \quad \lambda_2 = -0.1072 \times 10^{-1}, \quad \lambda_3 = -0.2579 \times 10^{-1}. \quad (55)$$

EXAMPLE 6.4. Consider the problem of minimizing the function

$$f = 0.01(x_1 - 1)^2 + (x_2 - x_1^2)^2 \quad (56)$$

subject to the inequality constraint

$$x_1 \le -1. \quad (57)$$

Introduce the auxiliary variable $x_2$ defined by

$$x_1 + x_2^2 + 1 = 0.$$

Then, the previous problem can be recast as that of minimizing the function (56) subject to the equality constraint (58). The function (56) admits the relative minimum $f = 0.04$ at the point defined by

$$x_1 = -1, \quad x_2 = 1, \quad x_3 = 0$$

and

$$\lambda_1 = 0.04. \quad$$

EXAMPLE 6.5. Consider the problem of minimizing the function

$$f = -x_1 \quad (61)$$

subject to the inequality constraints

$$x_1 \ge x_2^2, \quad x_1 \le x_2^2 \quad (62)$$

Introduce the auxiliary variables $x_3$ and $x_4$ defined by

$$x_1 - x_2^2 - x_3^2 = 0, \quad x_2^2 - x_1 - x_4^2 = 0. \quad (63)$$

Then, the previous problem can be recast as that of minimizing the function (61) subject to the equality constraints (63). The function (61) admits a relative minimum $f = -1$ at the point defined by

$$x_1 = 1, \quad x_2 = 1, \quad x_3 = 0, \quad x_4 = 0 \quad (64)$$

and

$$\lambda_1 = -1, \quad \lambda_2 = -1. \quad (65)$$

EXAMPLE 6.6. Consider the problem of minimizing the function

$$f = \log x_1 - x_2 \quad (66)$$

subject to the equality constraint

$$x_1^2 + x_2^2 - 4 = 0 \quad (67)$$

and the inequality constraint

$$x_3 \geq 1. \qquad (68)$$

Introduce the auxiliary variable $x_4$ defined by

$$x_3 = 1 + x_4^2. \qquad (69)$$

Then, the previous problem can be recast as that of minimizing the function

$$f = \log(1 + x_1^2) - x_3 \qquad (70)$$

subject to the equality constraint

$$(1 + x_1^2)^2 + x_2^2 - 4 = 0. \qquad (71)$$

1 that $x_3$ has been eliminated from the problem and can be computed a posteriori with (69). The function (70) admits the relative minimum $f = -\sqrt{3}$ at the point defined by

$$x_1 = 0, \quad x_2 = \sqrt{3}, \quad x_3 = 1 \qquad (72)$$

and

$$\lambda_1 = 1/2\sqrt{3}. \qquad (73)$$

## 7. Results and Conclusions

The examples described in Section 6 were solved with the penalty function method in conjunction with the conjugate gradient algorithm according to the experimental conditions outlined in Section 5. The numerical results are presented in Tables 1-16. Tables 1-4 supply a justification of the basic assumption (25), and Tables 5-16 supply comparative data for the standard updating rules and the new updating rules.

*Basic Hypothesis.* Tables 1-4 refer to a particular example, namely, Example 6.5. They show the behavior of the penalty function method for each of the penalty constant ratios (36) and (37) assuming that the starting penalty constant is $k_0 = 1$. In the tables, the quantities $P$, $Q$, $R$, $Z$ are shown versus the iteration number $N$ and the cycle number $N_c$. To shorten the tables, these quantities are given only at the initial point and the final point of each cycle.

As the tables indicate, the stopping condition (33) is met at convergence of a cycle and the stopping condition (40) is met at convergence of

the algorithm. Of particular interest is the behavior of the quantity $Z$, defined by $Z = P^T P = k^2 P$. If the end conditions of successive cycles are compared, it is clear that $Z$ approaches an asymptotic value as the cycle number $N_c$ increases and the penalty constant $k$ becomes sufficiently large. This asymptotic value is independent of the particular penalty constant ratio employed and is the same in each of Tables 1-4, namely $Z = 3$. This is in agreement with the theoretical values of the multipliers at convergence, which are supplied by Eqs. (68). This result confirms the basic assumption (25) and justifies the reasoning leading to the general penalty constant ratio (20).

*Comparative Data.* Tables 5-16 refer to Examples 6.1 through 6.6. They show the number of iterations at convergence $N_i$ and the number of cycles at convergence $N_c$ versus the starting penalty constant $k_0$ for each of the penalty constant ratios (36)-(37). Since each iteration requires the same amount of computational work, it is clear that the number of iterations is representative of the computer time required for convergence.

For all of the examples and each starting penalty constant $k_0$, the variable-rate updating rules are superior to the constant rate updating rules in that they require a smaller number of iterations for convergence. This smaller number of iterations is achieved by employing a smaller number of cycles and by increasing the penalty constant at a higher rate than that given by Eqs. (36).

TABLE 1. Example 6.5, $k_a = 1$, $x = 5$.

| $N_a$ | $N$ | $k$ | $P$ | $Q$ | $R$ | $z$ |
|---|---|---|---|---|---|---|
| 1 | 0 | 0.10 E + 01 | 0.10 E + 03 | 0.57 E + 05 | 0.57 E + 05 | 0.43 E + 03 |
| 1 | 17 | 0.10 E + 01 | 0.95 E − 01 | 0.12 E − 16 | 0.56 E − 01 | 0.30 E + 00 |
| 2 | 17 | 0.50 E + 01 | 0.93 E − 01 | 0.16 E + 03 | 0.16 E + 03 | 0.95 E + 01 |
| 2 | 27 | 0.50 E + 01 | 0.92 E − 02 | 0.12 E − 12 | 0.52 E − 02 | 0.93 E + 00 |
| 3 | 27 | 0.25 E + 02 | 0.92 E − 02 | 0.16 E + 03 | 0.16 E + 03 | 0.25 E + 03 |
| 3 | 41 | 0.25 E + 02 | 0.63 E − 03 | 0.16 E − 14 | 0.62 E − 03 | 0.14 E + 01 |
| 4 | 41 | 0.13 E + 03 | 0.62 E − 03 | 0.16 E + 03 | 0.16 E + 03 | 0.59 E + 03 |
| 4 | 47 | 0.13 E + 03 | 0.30 E − 04 | 0.60 E − 15 | 0.30 E − 04 | 0.19 E + 01 |
| 5 | 47 | 0.63 E + 03 | 0.30 E − 04 | 0.16 E + 03 | 0.16 E + 03 | 0.17 E + 02 |
| 5 | 53 | 0.63 E + 03 | 0.13 E − 05 | 0.13 E − 17 | 0.13 E − 05 | 0.20 E + 01 |
| 6 | 53 | 0.31 E + 04 | 0.13 E − 05 | 0.16 E + 03 | 0.16 E + 03 | 0.49 E + 02 |
| 6 | 59 | 0.31 E + 04 | 0.51 E − 07 | 0.64 E − 17 | 0.51 E − 07 | 0.20 E + 01 |
| 7 | 59 | 0.16 E + 05 | 0.51 E − 07 | 0.16 E + 03 | 0.16 E + 03 | 0.50 E + 02 |
| 7 | 65 | 0.16 E + 05 | 0.20 E − 08 | 0.14 E − 16 | 0.20 E − 08 | 0.20 E + 01 |
| 8 | 65 | 0.78 E + 05 | 0.20 E − 08 | 0.16 E + 02 | 0.16 E + 03 | 0.50 E + 02 |
| 8 | 70 | 0.78 E + 05 | 0.82 E − 10 | 0.25 E − 13 | 0.82 E − 10 | 0.20 E + 01 |
| 9 | 70 | 0.39 E + 06 | 0.82 E − 10 | 0.16 E + 03 | 0.16 E + 03 | 0.50 E + 02 |
| 9 | 75 | 0.39 E + 06 | 0.33 E − 11 | 0.66 E − 17 | 0.33 E − 11 | 0.20 E + 01 |
| 10 | 75 | 0.20 E + 07 | 0.33 E − 11 | 0.16 E + 03 | 0.16 E + 03 | 0.50 E + 02 |
| 10 | 79 | 0.20 E + 07 | 0.11 E − 12 | 0.11 E − 13 | 0.15 E − 12 | 0.20 E + 01 |

TABLE 2. Example 6.5, $k_a = 1$, $x = 10$.

| $N_a$ | $N$ | $k$ | $P$ | $Q$ | $R$ | $z$ |
|---|---|---|---|---|---|---|
| 1 | 0 | 0.10 E + 01 | 0.10 E + 03 | 0.57 E + 05 | 0.57 E + 05 | 0.43 E + 03 |
| 1 | 17 | 0.10 E + 01 | 0.95 E − 01 | 0.12 E − 16 | 0.56 E − 01 | 0.30 E + 00 |
| 2 | 17 | 0.10 E + 02 | 0.95 E − 01 | 0.61 E + 03 | 0.61 E + 03 | 0.34 E + 02 |
| 2 | 30 | 0.10 E + 02 | 0.90 E − 02 | 0.13 E − 13 | 0.90 E − 02 | 0.17 E + 01 |
| 3 | 30 | 0.10 E + 03 | 0.90 E − 02 | 0.61 E + 03 | 0.61 E + 03 | 0.12 E + 03 |
| 3 | 40 | 0.10 E + 03 | 0.18 E − 04 | 0.73 E − 14 | 0.18 E − 01 | 0.19 E + 01 |
| 4 | 40 | 0.10 E + 04 | 0.18 E − 04 | 0.61 E + 02 | 0.61 E + 03 | 0.19 E + 03 |
| 4 | 47 | 0.10 E + 04 | 0.60 E − 06 | 0.50 E − 14 | 0.50 E − 06 | 0.20 E + 01 |
| 5 | 47 | 0.10 E + 05 | 0.60 E − 06 | 0.61 E + 02 | 0.61 E + 03 | 0.20 E + 03 |
| 5 | 53 | 0.10 E + 05 | 0.60 E − 08 | 0.13 E − 17 | 0.50 E − 08 | 0.20 E + 01 |
| 6 | 53 | 0.10 E + 06 | 0.60 E − 08 | 0.61 E + 02 | 0.61 E + 03 | 0.20 E + 03 |
| 6 | 59 | −0.10 E + 06 | 0.50 E − 10 | 0.27 E − 20 | 0.50 E − 10 | 0.20 E + 01 |
| 7 | 59 | 0.10 E + 07 | 0.60 E − 10 | 0.61 E + 02 | 0.61 E + 03 | 0.20 E + 03 |
| 7 | 64 | 0.10 E + 07 | 0.60 E − 12 | 0.19 E − 16 | 0.56 E − 12 | 0.20 E + 01 |

# Comparison of Multiplier and Quasilinearization Methods

**Alejandro V. Levy***

CIMAS, Universidad Nacional Autonoma de Mexico, Mexico City, Mexico and Department of Mechanical Engineering, Rice University, Houston, Texas

**Antonio Montalvo**

Department of Mathematics, Universidad Iberoamericana, Mexico City, Mexico

Miele et al. (1972) developed the method of multipliers for minimizing a function $f(x)$ subject to the constraint $c(x) = 0$, where $f$ is a scalar, $x$ is an $n$ vector, and $c$ is a $q$ vector, with $q < n$. In this paper, a comparison of the standard quasilinearization method and the method of multipliers is presented. The comparison is made by considering the relative stability and the relative speed of convergence of the two methods. Numerical results are presented for nine examples, each of which is solved for 100 different starting points. The examples show that, in general, the method of multipliers is a more economic and robust algorithm than the standard quasilinearization method. In particular, this advantage increases in proportion to the following characteristics of the problem being solved: (a) the ratio $q/n$ and the size of the problem and (b) the nonlinearity of the problem and the number of relative minima and maxima that may exist.

## 1. Introduction

In recent years, considerable attention has been given to the quasilinearization methods for minimizing a function $f(x)$ subject to a constraint $c(x) = 0$. Here, $f$ is a scalar, $x$ is an $n$-vector, and $c$ is a $q$-vector, with $q < n$.

In the *standard quasilinearization algorithm* (SQL), the augmented function $F(x,\lambda) = f(x) + \lambda^T c(x)$ is utilized, where $\lambda$, a $q$-vector, is an unknown Lagrange multiplier. Starting from nominal values of $x$ and $\lambda$, the method obtains corrections $\Delta x$ and $\Delta \lambda$ by solving a linear system of equations. In this way, updated values $\tilde{x} = x + \Delta x$ and $\tilde{\lambda} = \lambda + \Delta \lambda$ are obtained. Then, the process is repeated iteratively until convergence occurs. For the particular case of a quadratic function subject to a linear constraint, the method converges to the solution in one iteration. However, for a nonlinear function subject to nonlinear constraints, the method is rather unstable. Depending on the nominal values of $x$ and $\lambda$, it can lead to a relative minimum, a relative maximum, or an inflection point; occasionally, the method produces displacements which are so purposeless that overflow can occur in a given computer.

In the *modified quasilinearization algorithm* (MQL), Miele et al. (1973), the same basic approach is taken as in SQL, with one important modification: the stepsize $\alpha$, $0 < \alpha \leq 1$, is introduced in order to control the magnitude of the displacements. The result is an improved quasilinearization algorithm that (i) retains quadratic convergence for a linear-quadratic problem and (ii) improves the stability of SQL for nonlinear functions subject to nonlinear constraints.

Both the above mentioned methods require the solution of a linear system of equations of order $n + q$. Thus, if Gaussian elimination is employed, the computational effort is of the order of $(n + q)^3$ multiplications.

For large systems, one way to reduce the computational effort per iteration is to resort to penalty function methods, either in the standard format or the multiplier format. Hestenes (1969) suggested this approach using the augmented penalty function $W(x,\lambda,k) = F(x,\lambda) + k c^T(x) c(x)$, where the scalar $k$ is the penalty constant.

With reference to the method of multipliers, an important modification was presented by Miele et al. (1972). In this modification, called the *modified method of multi-*

*pliers* (MMM), the augmented penalty function is employed in connection with several minimization algorithms, namely, the ordinary gradient algorithm, the conjugate gradient algorithm, and the modified quasilinearization algorithm. Improved updating rules for the Lagrange multiplier and the penalty constant were developed in connection with each minimization algorithm.

When the modified method of multipliers is employed in conjunction with the modified quasilinearization algorithm, a linear system of equations of order $n$ must be solved at every iteration. Therefore, if Gaussian elimination is employed, the computational effort is of the order of $n^3$ multiplications.

Besides reducing the computational effort per iteration, MMM is more stable than SQL. This is because MMM has a descent property on the augmented penalty function $W(x,\lambda,k)$. While SQL might converge to a relative minimum, a relative maximum, or an inflection point, MMM generally leads to a relative minimum.

The standard quasilinearization algorithm was used by Luus and Jaakola (1973) with some success to solve several problems, using random numbers as starting values for $x$ and $\lambda$. The purpose of this report is to show that even greater success and computer time savings can be obtained if the modified method of multipliers is utilized instead of the standard quasilinearization algorithm.

In section 2, the statement of the problem is given together with the first-order optimality conditions. In sections 3 and 4, the standard quasilinearization algorithm and the modified method of multipliers are reviewed briefly. In section 5, a comparison of the computational effort per iteration is given. In sections 6 and 7, the experimental conditions and nine numerical examples are presented. In section 8, numerical results are discussed, and the conclusions are stated.

## 2. Statement of the Problem

We consider the problem of minimizing the function

$$f = f(x) \tag{1}$$

subject to the constraint

$$c(x) = 0 \tag{2}$$

TABLE 3. Example 6.5, $k_0 = 1$, $\pi = 10 \sqrt[4]{(P_1/P_0)}$

| $N_r$ | $N$ | $k$ | $P$ | $Q$ | $R$ | $S$ |
|---|---|---|---|---|---|---|
| 1 | 0 | 0.10 E + 01 | 0.10 E + 03 | 0.57 E + 06 | 0.87 E + 05 | 0.43 E + 03 |
| 1 | 17 | 0.10 E + 01 | 0.05 E − 01 | 0.13 E − 16 | 0.95 E − 01 | 0.39 E + 04 |
| 2 | 17 | 0.35 E + 04 | 0.25 E − 01 | 0.31 E + 05 | 0.81 E + 05 | 0.13 E + 04 |
| 3 | 33 | 0.65 E + 04 | 0.16 E − 07 | 0.13 E − 13 | 0.19 E − 07 | 0.30 E + 01 |
| 3 | 33 | 0.63 E + 06 | 0.16 E − 07 | 0.13 E + 05 | 0.13 E + 05 | 0.75 E + 05 |
| 3 | 36 | 0.63 E + 06 | 0.15 E − 11 | 0.33 E − 15 | 0.13 E − 11 | 0.20 E + 01 |
| 4 | 36 | 0.67 E + 07 | 0.13 E − 13 | 0.53 E + 07 | 0.05 E + 03 | 0.23 E + 05 |
| 4 | 42 | 0.87 E + 07 | 0.11 E − 19 | 0.34 E − 13 | 0.35 E − 13 | 0.20 E + 01 |

TABLE 4. Example 6.5, $k_0 = 1$, $\pi = 10 \sqrt[4]{(P_1/P_0)}$

| $N_r$ | $N$ | $k$ | $P$ | $Q$ | $R$ | $S$ |
|---|---|---|---|---|---|---|
| 1 | 0 | 0.10 E + 01 | 0.10 E + 03 | 0.57 E + 05 | 0.67 E + 05 | 0.42 E + 03 |
| 1 | 17 | 0.10 E + 01 | 0.06 E − 01 | 0.13 E − 13 | 0.95 E − 01 | 0.05 E + 00 |
| 2 | 17 | 0.33 E + 07 | 0.35 E − 01 | 0.06 E + 13 | 0.96 E + 13 | 0.80 E + 13 |
| 2 | 34 | 0.31 E + 07 | 0.53 E − 19 | 0.60 E − 14 | 0.63 E − 19 | 0.29 E + 03 |

TABLE 5. Example 6.1, number of iterations $N_i$.

| $k_0$ | $\pi = 5$ | $\pi = 10$ | $\pi = 10 \sqrt[4]{(P_1/P_0)}$ | $\pi = 10 \sqrt[4]{(P_1/P_0)}$ |
|---|---|---|---|---|
| $10^{-2}$ | 30 | 63 | 34 | 22 |
| $10^{-1}$ | 60 | 47 | 21 | 17 |
| $10^0$ | 55 | 43 | 21 | 19 |
| $10^1$ | 46 | 37 | 33 | 29 |
| $10^2$ | 40 | 33 | 30 | 28 |

TABLE 6. Example 6.1, number of cycles $N_c$.

| $k_0$ | $\pi = 5$ | $\pi = 10$ | $\pi = 10 \sqrt[4]{(P_1/P_0)}$ | $\pi = 10 \sqrt[4]{(P_1/P_0)}$ |
|---|---|---|---|---|
| $10^{-2}$ | 14 | 10 | 5 | 5 |
| $10^{-1}$ | 13 | 8 | 4 | 5 |
| $10^0$ | 11 | 6 | 4 | 5 |
| $10^1$ | 10 | 5 | 4 | 5 |
| $10^2$ | 8 | 5 | 4 | 5 |

TABLE 7. Example 6.2, number of iterations $N_e$.

| $b_i$ | $n=5$ | $n=10$ | $n=10\,f(\overline{F_i/F_a})$ | $n=10\,f(\overline{F_i/F_a})$ |
|---|---|---|---|---|
| $10^{-2}$ | 68 | 62 | 58 | 41 |
| $10^{-1}$ | 84 | 60 | 60 | 43 |
| $10^{0}$ | 59 | 60 | 58 | 55 |
| $10^{1}$ | 68 | 87 | 44 | 41 |
| $10^{2}$ | 50 | 80 | 77 | 76 |

TABLE 8. Example 6.2, number of cycles $N_c$.

| $b_i$ | $n=5$ | $n=10$ | $n=10\,f(\overline{F_i/F_a})$ | $n=10\,f(\overline{F_i/F_a})$ |
|---|---|---|---|---|
| $10^{-2}$ | 10 | 7 | 4 | 3 |
| $10^{-1}$ | 8 | 6 | 3 | 2 |
| $10^{0}$ | 7 | 5 | 3 | 2 |
| $10^{1}$ | 6 | 4 | 4 | 3 |
| $10^{2}$ | 4 | 3 | 8 | 2 |

TABLE 9. Example 6.3, number of iterations $N_e$.

| $b_i$ | $n=5$ | $n=10$ | $n=10\,f(\overline{F_i/F_a})$ | $n=10\,f(\overline{F_i/F_a})$ |
|---|---|---|---|---|
| $10^{-2}$ | (a) | (a) | 131 | 176 |
| $10^{-1}$ | 154 | 186 | 170 | 86 |
| $10^{0}$ | 186 | 162 | 90 | 100 |
| $10^{1}$ | 135 | 138 | 83 | 79 |
| $10^{2}$ | 186 | 131 | 107 | 102 |

TABLE 10. Example 6.3, number of cycles $N_c$.

| $b_i$ | $n=5$ | $n=10$ | $n=10\,f(\overline{F_i/F_a})$ | $n=10\,f(\overline{F_i/F_a})$ |
|---|---|---|---|---|
| $10^{-2}$ | (a) | (a) | 4 | 3 |
| $10^{-1}$ | 9 | 7 | 4 | 3 |
| $10^{0}$ | 8 | 8 | 3 | 2 |
| $10^{1}$ | 6 | 6 | 8 | 3 |
| $10^{2}$ | 5 | 4 | 8 | 3 |

(a) Number of iterations exceeded 200.

TABLE 11. Example 6.4, number of iterations $N_e$.

| $k_0$ | $a = 5$ | $a = 10$ | $a = 10\sqrt{(F_2/F_1)}$ | $a = 10\sqrt{(F_2/F_0)}$ |
|---|---|---|---|---|
| $10^{-2}$ | 65 | 65 | 91 | 90 |
| $10^{-1}$ | 68 | 67 | 83 | 81 |
| $10^0$ | 65 | 65 | 39 | 65 |
| $10^1$ | 81 | 75 | 70 | 68 |
| $10^2$ | 175 | 170 | 167 | 161 |

TABLE 12. Example 6.4, number of cycles $N_{r_0}$.

| $k_0$ | $a = 5$ | $a = 10$ | $a = 10\sqrt{(F_2/F_1)}$ | $a = 10\sqrt{(F_2/F_0)}$ |
|---|---|---|---|---|
| $10^{-2}$ | 11 | 8 | 4 | 3 |
| $10^{-1}$ | 9 | 7 | 8 | 3 |
| $10^0$ | 8 | 6 | 8 | 3 |
| $10^1$ | 6 | 5 | 8 | 2 |
| $10^2$ | 6 | 4 | 4 | 3 |

TABLE 13. Example 6.5, number of iterations $N_e$.

| $k_0$ | $a = 5$ | $a = 10$ | $a = 10\sqrt{(F_2/F_1)}$ | $a = 10\sqrt{(F_2/F_0)}$ |
|---|---|---|---|---|
| $10^{-2}$ | 116 | 93 | 61 | 44 |
| $10^{-1}$ | 99 | 78 | 47 | 37 |
| $10^0$ | 78 | 94 | 42 | 34 |
| $10^1$ | 67 | 81 | 35 | 39 |
| $10^2$ | 58 | 48 | 88 | 31 |

TABLE 14. Example 6.5, number of cycles $N_{r_0}$.

| $k_0$ | $a = 5$ | $a = 10$ | $a = 10\sqrt{(F_2/F_1)}$ | $a = 10\sqrt{(F_2/F_0)}$ |
|---|---|---|---|---|
| $10^{-2}$ | 18 | 9 | 4 | 3 |
| $10^{-1}$ | 11 | 8 | 4 | 2 |
| $10^0$ | 19 | 7 | 6 | 1 |
| $10^1$ | 8 | 8 | 8 | 3 |
| $10^2$ | 7 | 6 | 8 | 4 |

TABLE 15. Example 8.6, number of iterations $N_e$.

| $k_e$ | $n = 5$ | $n = 10$ | $n = 10\sqrt{(F_2/F_1)}$ | $n = 10\sqrt{(F_2/F_1)}$ |
|---|---|---|---|---|
| $10^{-2}$ | 27 | 21 | 15 | 10 |
| $10^{-1}$ | 23 | 16 | 13 | 9 |
| $10^0$ | 21 | 19 | 13 | 11 |
| $10^1$ | 20 | 16 | 13 | 11 |
| $10^2$ | 26 | 21 | 20 | 19 |

TABLE 16. Example 8.6, number of cycles $N_c$.

| $k_e$ | $n = 5$ | $n = 10$ | $n = 10\sqrt{(F_2/F_1)}$ | $n = 10\sqrt{(F_2/F_1)}$ |
|---|---|---|---|---|
| $10^{-2}$ | 11 | 8 | 4 | 2 |
| $10^{-1}$ | 10 | 8 | 4 | 2 |
| $10^0$ | 8 | 7 | 3 | 3 |
| $10^1$ | 7 | 6 | 3 | 3 |
| $10^2$ | 6 | 6 | 3 | 2 |

## REFERENCES

[1] MIELE, A., COGGINS, O. M., and LEVY, A. V., «Updating Rules for the Penalty Constant Used in the Penalty Function Method for Mathematical Programming Problems», Rice University, Aero-Astronautics Report No. 90, 1972.

[2] MIELE, A., HUANG, H. Y., and HEIDEMAN, J. C., «Sequential Gradient-Restoration Algorithm for the Minimization of Constrained Functions, Ordinary and Conjugate Gradient Versions», Journal of Optimization Theory and Applications, Vol. 4, No. 4, 1969.

[3] MIELE, A., LEVY, A. V., and CRAGG, E. E., «Modifications and Extensions of the Conjugate Gradient-Restoration Algorithm for Mathematical Programming Problems», Journal of Optimization Theory and Applications, Vol. 7, No. 6, 1971.

[4] MIELE, A., HEIDEMAN, J. C., and LEVY, A. V., «Combined Conjugate Gradient-Restoration Algorithm for Mathematical Programming Problems», Ricerche di Automatica, Vol. 2, No. 2, 1971.

[5] KELLEY, H. J., «Methods of Gradients», Optimization Techniques, Edited by G. Leitmann, Academic Press, New York, 1961.

[6] BRYSON, A. E., Jr., and HO, Y. C., «Applied Optimal Controls», Blaisdell Publishing Company, Waltham, Massachusetts, 1969.

[7] FIACCO, A. V., and McCORMICK, G. P., «Nonlinear Programming: Sequential Unconstrained Minimization Techniques», John Wiley and Sons, New York, 1968.

[8] DAVIES, D., and SWANN, W. H., «Review of Constrained Optimization», Optimization, Edited by R. Fletcher, Academic Press, New York, 1969.

**G. A. Gabriele**
Graduate Assistant.

**K. M. Ragsdell**
Associate Professor.
Mem. ASME

School of Mechanical Engineering,
Purdue University,
West Lafayette, Ind.

# The Generalized Reduced Gradient Method: A Reliable Tool for Optimal Design

*This paper is a presentation of a method, called the Generalized Reduced Gradient Method, which has not received wide attention in the engineering design literature. Included is a theoretical development of the method, a description of the basic algorithm, and additional recommendations to produce an efficient code. A Fortran code employing this theory was written and tested on the Eason and Fenton [1][1] test problems, illustrating the method to be efficient and reliable.*

## Introduction

Engineering design is a multiphase process requiring constant decision making by the designer. Basing decisions on his knowledge and experience, he must arrive at a combination of design variables that best satisfies his objectives. As engineering design has matured so have the guidelines and methods that the designer has at his disposal to aid him in his choices. Although his experience may provide some answers, many of his choices are based on analytical methods. Designers have generally been receptive to new methods, when they simplify and/or otherwise enhance the overall synthesis process. The availability of high-speed digital computing power has encouraged the application of modern optimization methods (mathematical programming methods) to engineering design. This marriage of effort is generating developments in both fields, and has produced an important area of research known as optimal design.

Drawing on his experience the engineer is able to define variables, a design object is central a set of constraints that must be met in order that the design be a workable solution. By developing corresponding equations the design problem can be formulated into a standard form acceptable to mathematical programming techniques. This standard form is defined here

Minimize

$$ f(x); \quad x = [x_1, x_2, x_3, \ldots, x_N]^T, \quad x \in R^N \tag{1} $$

subject to

$$ \phi_k(x) \geq 0 \quad k = 1, 2, 3, \ldots, K \tag{2} $$

$$ \psi_l(x) = 0 \quad l = 1, 2, 3, \ldots, L \tag{3} $$

where

$x$ = a column vector of design variables

$N$ = total number of design variables

$f(x)$ = design criteria or objective function

$\phi_k(x)$ = $K$ inequality constraint functions; these functions define regions in the design space

$\psi_l(x)$ = $L$ equality constraint functions; these functions vastly reduce the number of candidate designs, because they require specific combinations of the design variables

When equations (1)–(3) are linear functions in the design variables the problem falls in the general class of linear programming (LP) problems. Much of the research in the past 20 years in the field of mathematical programming has been with linear programming. The simplex algorithm has become highly developed and most linear programs can be solved using this technique.

A more general occurrence in engineering design arises when expressions (1)–(3) are nonlinear. This is known as the nonlinear programming (NLP) problem. No general method has been developed to solve nonlinear problems in the sense that the simplex algorithm exists to solve the linear problem. Although many strategies have been suggested, comparative studies [1, 2] have shown that no method has been successfully applied to all problems.

One procedure for solving the NLP that has seen widespread usage is the penalty function approach. Consider the following function

$$ P(x) = f(x) + W(x, R, \phi, \psi) \tag{4} $$

which we call the generalized weighted penalty function. $W(x, R, \phi, \psi)$ is the penalty term, and is constructed in such a manner that feasible points (points satisfying (2) and (3)) are given favored treatment. The basic concept of the penalty function method is that successive minimizations of the unconstrained function $P(x)$ will converge to the constrained minimum of the original NLP. A survey of penalty function techniques is given by Fiacco and McCormick [3].

1

[1] Numbers in brackets designate References at end of paper.

An advantage of penalty function methods is their ease of implementation. Any unconstrained searching technique can be used to perform the successive minimizations of $P(x)$. However, there are several difficulties with penalty function methods that do not make them a generally reliable technique. First, if the successive minimizations of $P(x)$ are not trivial, then the time to solution may become impractical. Penalty function methods are also very sensitive to parameter adjustment and constraint scaling, and may require repetitive adjustment to obtain convergence to an optimum. Finally, penalty function methods generally do not possess the ability to move along equality constraints. Rather they tend to get caught on a point that satisfies the constraints, but is not the constrained optimum.

Because the state of the art for linear programming is so well developed, another approach to solving the constrained NLP is linearization of the NLP to tailor it to the LP methods. Basically this can be accomplished by replacing the nonlinear functions of (1)-(3) by first-order Taylor series approximations and solving the resulting LP. The Griffith and Stewart method [4] is based on this approach.

The advantage of this approach lies in being able to use existing linear programming theory. The method performs well on large-scale problems when the objective function and constraints are nearly linear. Difficulties arise when these functions are very nonlinear. Hence in order that the linear approximations hold, only small changes in the design variables are allowed, which may cause progress to the optimum to be slow. Indeed in the presence of highly nonlinear functions, the LP solution may give rise to a poor search direction, therefore hindering the optimization process.

Other methods of handling the constrained NLP are described in the literature [5, 6]; however, many of these methods do not solve the complete NLP described in (1)-(3). The two approaches described here do handle the complete problem and have seen wide acceptance in optimal engineering design today. In this paper another approach known as the Reduced Gradient Method will be described. The Reduced Gradient Method avoids many of the problems associated with penalty function and LP-like methods, producing one of the most powerful methods currently known for handling the constrained nonlinear programming problem.

## Reduced Gradient—Theory

The Reduced Gradient Method was originally given by Wolfe for a nonlinear objective function with linear constraints [7, 8]. A generalization of Wolfe's method to accommodate nonlinearities in both the objective function and constraints was first accomplished by Abadie [9]. Concurrently to both Wolfe and Abadie, Wilde and Beightler developed their differential algorithm based on the constrained derivative [10]. The constrained derivative and the reduced gradient employ much the same theoretical basis, but for purposes of this discussion the method shall be known as the Reduced Gradient Method.

The general constrained nonlinear programming problem of (1)-(3) can be restated in the following form:

Minimize

$$f(x); \qquad x = |x_1, x_2, \ldots, x_N|^T, \qquad x \in R^N \qquad (5)$$

subject to

$$\phi_m(x) = 0 \qquad m = 1, 2, \ldots, M \qquad (6)$$

$$\boxed{A \leq x \leq B} \qquad (7)$$

The $N \times 1$ vectors $A$ and $B$ represent upper and lower bounds on the design vector $x$. The inequality constraints have been included as equality constraints by using the following transformation:

$$\phi_k(x) = \phi_k(x) - S_k = 0$$

$$0 \leq S_k \leq \infty \qquad k = 1, 2, \ldots, K \qquad (8)$$

The variables $S_k$ are nonnegative slack variables which will be included in the original set of design variables. Hence the parameter $N$ represents the total number of design variables plus the number of slack variables used in the transformation of (8). The parameter $M$ represents the total number of constraints.

$$M = L + K \qquad (9)$$

It should be stressed that the constraints of (6) are nontrivial constraints; that is, functional constraints. Variable bounds are defined in (7) and will require separate handling.

Consider the following strategy. Divide the design vector of equation (5) into two classes which shall be known as the decision and state variables.

$$x = |z, y|^T \qquad (10)$$

$$z = |z_1, z_2, z_3, \ldots, z_Q|^T; \qquad \text{decision variables} \qquad (11)$$

$$y = |y_1, y_2, y_3, \ldots, y_M|^T; \qquad \text{state variables} \qquad (12)$$

$$Q = N - M \qquad (13)$$

The decision variables are completely independent, and the state variables are slaves to the decision variables used to satisfy the constraints $\phi_m(x)$.

We can see a motivation for this division of the $x$ vector into decision and state variables based on the elimination technique for handling equality constraints. Given the equality constrained NLP.

Minimize

$$f(x); \qquad x = |x_1, x_2, x_3, \ldots, x_N|^T \qquad (14)$$

subject to

$$\phi_i(x) = 0 \qquad i = 1, 2, 3, \ldots, L \qquad (15)$$

Reformulate each constraint in (15) to explicitly define a different design variable.

$$x_i = h_i(x) \qquad i = 1, 2, 3, \ldots, L \qquad (16)$$

Substitution into the objective function (14) would yield

$$f(h_1, h_2, \ldots, h_L, x_{L+1}, \ldots, x_N) \qquad (17)$$

This new objective function becomes an unconstrained function in $N - L$ design variables. Minimization of this new objective function would yield the constrained optimal values of the $N - L$ design variables describing (17) with the remaining design variables calculated by back substitution into (16).

In the foregoing technique the decision variables become the $N - L$ design variables describing (17) with the state variables calculated by (16). We can see here a basic approach for constrained optimization

## Nomenclature

$f(x), f$ = objective function

$x = x_1, x_2, \ldots, x_N$; set of design variables

$R^N$ = $N$-dimensional space of reals

$\phi_k(x), \phi$ = inequality constraint functions

$\phi_i(x), \phi$ = equality constraint functions

$N$ = number of design variables

$K$ = number of inequality constraints

$L$ = number of equality constraints

$P(x)$ = penalty function

$W(x, R, \phi, \psi)$ = ... penalty term

$S_k$ = slack variables

$M$ = total number of constraints

$Q$ = number of decision variables

$z$ = $Q$-dimensional column vector of decision variables

$y$ = $M$-dimensional column vector of state variables

$L(x, u, \theta)$ = Lagrangian function

$H, H_{ij}$ = transformation matrix, $M \times Q$

$\ldots$ = reduced gradient, $Q \times 1$

$P(z)$ = search direction of decision variable, $Q \times 1$

$P(y)$ = search direction providing first approximation of state variables, $M \times 1$

$\alpha$ = step length parameter

$\theta^*$ = first approximation of state variables

$y^*$ = adjusted values of state variables

$N_f$ = number of objective function evaluations $\ldots x$

$N_c$ = number of constraint set evaluations

using the decision and state dichotomy. We can perturb the decision variables in the newly formed unconstrained space to seek the minimum of the objective function while adjusting the state variables to maintain feasibility. However, use of the technique just described presupposes that the state variables can be defined explicitly in terms of the constraint functions. In engineering design problems this is generally the exception and not the rule. Some other method must be devised to utilize the decision and state variable dichotomy when nonlinearity is a factor.

The following notation will be useful in the discussion to follow:

$$g(z) = \left[\frac{\partial f(z)}{\partial z_1}, \frac{\partial f(z)}{\partial z_2}, \dots, \frac{\partial f(z)}{\partial z_Q}\right]^T$$

$$g(y) = \left[\frac{\partial f(z)}{\partial y_1}, \frac{\partial f(z)}{\partial y_2}, \dots, \frac{\partial f(z)}{\partial y_M}\right]^T$$

$$\frac{\partial \psi}{\partial z} = \begin{bmatrix} \frac{\partial \psi_1}{\partial z_1} & \frac{\partial \psi_1}{\partial z_2} & \cdots & \frac{\partial \psi_1}{\partial z_Q} \\ & & & \\ & & & \\ \frac{\partial \psi_M}{\partial z_1} & \frac{\partial \psi_M}{\partial z_2} & \cdots & \frac{\partial \psi_M}{\partial z_Q} \end{bmatrix} \quad (M \times Q)$$

$$\frac{\partial \psi}{\partial y} = \begin{bmatrix} \frac{\partial \psi_1}{\partial y_1} & \frac{\partial \psi_1}{\partial y_2} & \cdots & \frac{\partial \psi_1}{\partial y_M} \\ & & & \\ & & & \\ \frac{\partial \psi_M}{\partial y_1} & \frac{\partial \psi_M}{\partial y_2} & \cdots & \frac{\partial \psi_M}{\partial y_M} \end{bmatrix} \quad (M \times M)$$

Let us examine the first variation of $f(z)$ and $\psi(z)$,

$$df = g(z)^T dz + g(y)^T dy \tag{18}$$

$$d\psi = \frac{\partial \psi}{\partial z} dz + \frac{\partial \psi}{\partial y} dy = 0 \tag{19}$$

where

$dz = Q \times 1$ vector of differential displacements of $z$

$dy = M \times 1$ vector of differential displacements of $y$

Solving (19) for $dy$ yields,

$$dy = -\frac{\partial \psi}{\partial y}^{-1} \frac{\partial \psi}{\partial z} dz \tag{20}$$

Substituting (20) into (18) and rearranging will yield the following linear approximation to the *reduced gradient*:

$$g_r(z)^T = g(z)^T - g(y)^T \frac{\partial \psi}{\partial y}^{-1} \frac{\partial \psi}{\partial z} \tag{21}$$

The reduced gradient defines the rate of change of the objective function with respect to the decision variables with the state variables adjusted to maintain feasibility. Expression (20) gives the changes necessary in the states for a given change in the decisions for latent constraints. Geometrically the reduced gradient can be described as a projection of the original $N$-dimensional gradient onto the $(N - M)$-dimensional feasible region described by the decision variables.

A necessary condition for the existence of a minimum of an unconstrained nonlinear function is that the elements of the gradient vanish. Similarly, a minimum of the constrained nonlinear function occurs when the appropriate elements of the reduced gradient vanish. This conclusion can be verified by a comparison with the Kuhn-Tucker [11] conditions for the existence of a constrained relative minimum.

By first transforming the variable bounds into inequality constraints,

$$\phi_i(z) = z_i - a_i \geq 0 \left.\begin{array}{l} \\ \end{array}\right| \quad i = 1, 2, 3, \dots, N \tag{22}$$
$$\phi_{N+i}(z) = b_i - z_i \geq 0 \left.\begin{array}{l} \\ \end{array}\right|$$

we can form the following Lagrangian function,

$$L(z, a, \omega) = f(z) + \sum_{m=1}^{M} \omega_m \psi_m(z) - \sum_{j=1}^{2N} u_j \phi_j(z) \tag{24}$$

The following Kuhn-Tucker necessary conditions hold for a point to be a relative minimum $z^*$.

$$g(z^*)^T + \sum_{m=1}^{M} \omega_m^* \frac{\partial \psi_m}{\partial z} - \sum_{j=1}^{J} u_j^* \frac{\partial \phi_j}{\partial z} = 0 \tag{25}$$

and

$$\phi_m(z^*) = 0 \qquad m = 1, 2, \dots, M \tag{26}$$

$$\phi_j(z^*) \geq 0 \qquad j = 1, 2, \dots, J = 2N \tag{27}$$

$$u_j^* \phi_j(z^*) = 0 \qquad j = 1, 2, \dots, J = 2N \tag{28}$$

$$u_j^* \geq 0 \qquad j = 1, 2, \dots, J = 2N \tag{29}$$

$$\omega_m^* \neq 0 \qquad m = 1, 2, \dots, M \tag{30}$$

Introducing decision and state variables into (25) and decomposing we can obtain the following form:

$$g(z^*)^T + \omega^{*T} \frac{\partial \psi}{\partial z^*} - u^{*T} \frac{\partial \phi}{\partial z^*} = 0 \tag{31}$$

$$g(y^*)^T + \omega^{*T} \frac{\partial \psi}{\partial y^*} - u^{*T} \frac{\partial \phi}{\partial y^*} = 0 \tag{32}$$

For reasons to be examined in the next section, a state variable is not allowed to be equal or sufficiently close to either of its bounds. From expression (28) the elements of $u^*$ corresponding to the state variable bounds must be zero. Also those elements of $\partial \phi/\partial y$ corresponding to the decision variable bounds will be zero eliminating the last term of (32). Solving (32) for $\omega^*$ and substituting into (31) will produce the following expression:

$$g(z^*)^T - g(y^*)^T \frac{\partial \psi}{\partial y^*}^{-1} \frac{\partial \psi}{\partial z^*} - u^{*T} \frac{\partial \phi}{\partial z^*} = 0 \tag{33}$$

Rearranging (33) we obtain

$$u^{*T} \frac{\partial \phi}{\partial z^*} = g(z^*)^T - g(y^*)^T \frac{\partial \psi}{\partial y}^{-1} \frac{\partial \psi}{\partial z} \tag{34}$$

We recognize the right-hand side of (34) to be $g_r(z)$. By examining the possible values of the left-hand side of (34), a candidate point $z$ will be $z^*$ if

$$g_r(z)_i \begin{cases} > 0 \text{ if } z_i = a_i \\ < 0 \text{ if } z_i = b_i \\ = 0 \text{ if } a_i \leq z_i \leq b_i \end{cases}$$
$$i = 1, 2, 3, \dots, Q \tag{35}$$

## Reduced Gradient — Algorithm

A flow chart of the Reduced Gradient Method to be described here is shown in Fig. 1. The basic steps involved are similar to those described by Abadie. The major differences described here occur in the calculation of the reduced gradient and the line search. Also it should be noted here that the present implementation of the method uses numerical techniques to obtain the required derivatives. In the results reported by Abadie, analytical derivatives were employed. However, obtaining analytical derivatives in many engineering applications can prove to be a tedious task and is often prone to error. In the following sections, implementation and computational considerations for each step in Fig. 1 will be discussed.
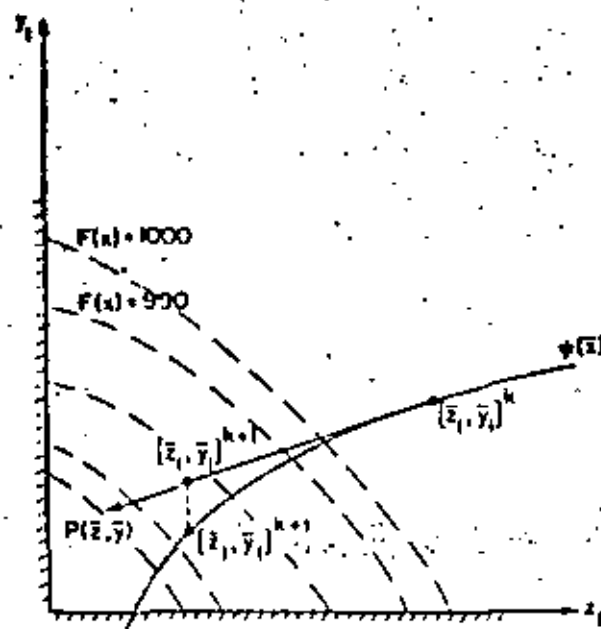
1 Specification of State and Decision Variables. In most

**Fig. 2** Adjustment of state variable to obtain a feasible point during the linear search
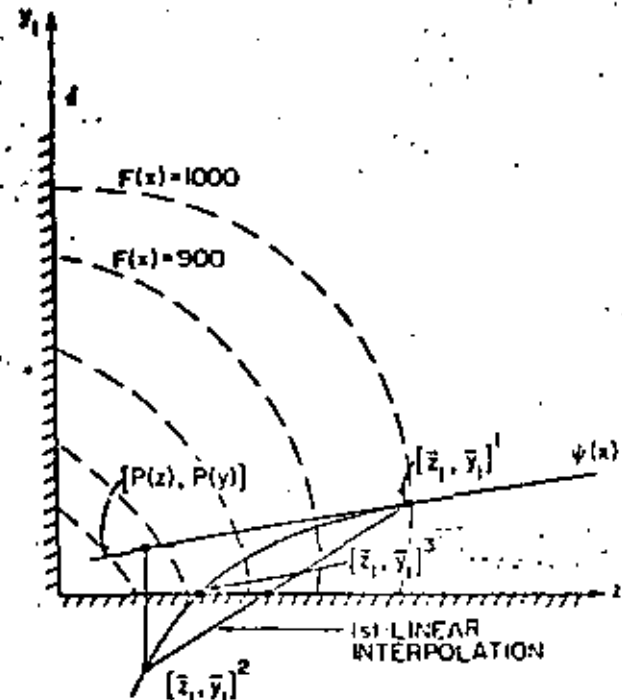


**Fig. 3** Adjustment of state variable to locate a feasible point

and

$$y_m^{k+1} = y_m^k + \alpha P^k r_m \qquad m = 1, 2, 3, \ldots, M \qquad (14)$$

where $\alpha$ = step length parameter.

Because of nonlinearities arising in the constraint functions, the point $[z, y]^{k+1}$ is likely to be infeasible. Holding the decision variables $z^{k+1}$ constant, the state variables $y^{k+1}$ are adjusted to obtain a feasible point, $[z, y]^{k+1}$. This situation is shown in Fig. 2 for the case $M = 1$, $U = 1$. This step is equivalent to the solution of $M$ nonlinear equations $\psi(z) = 0$ in $M$ unknowns $y_m$. A number of numerical techniques exist to perform this task. Newton's method [15] has proven to be an efficient technique as well as convenient since the necessary partial derivatives have already been calculated.

At the completion of the adjustment procedure a new point $[z, y]^{k+1}$ has been determined and the following possible results must be considered:

(a) If all elements of $y^{k+1}$ are within their specified bounds then $f(z)$ is evaluated at $[z, y]^{k+1}$ and the procedure for determining the minimum continues in the normal manner.

(b) If any element of $y^{k+1}$ is not within its specified bounds, then $[z, y]^{k+1}$ is infeasible. Success a linear interpolation is performed between the last feasible point $[z, y]^k$ and the point $[z, y]^{k+1}$ to determine the step length at which the nearest bound becomes inactive. Hence this step should conclude with a single state variable equal to one of its bounds and all other state variables within their operating bounds. Fig. 3 shows $[z, y]^2$ being out of bounds ($y_1 < 0$). Using successive false position, the bound point $[z, y]^3$ can be located. Supplementary tests are then performed to determine whether the local minimum lies at the bound or at some point before it. If the minimum lies at the bound, the line search is terminated. If it lies below the bound, then the minimum has been located and no supplementary tests can be used to locate the minimum.

(c) If the procedure fails to converge in a reasonable amount of time, the step length is reduced (to $\alpha/2$) and a new trial point generated.

**5. Respecify State and Decision Variables.** We see in the previous step that it was possible to end a line search with a state variable equal to one of its specified bounds (second result listed previously). To start another iteration the guidelines initially set forth for state variable selection must be satisfied. These state variables

that are bounded must be exchanged with decision variables that are unbounded and will not cause $\partial \psi / \partial y$ to become singular.

Assume that $y_s$ is a bound state variable. Abadie [16] suggests the following expression be *maximized* to determine the appropriate decision variable to exchange with $y_s$.

$$\text{Min}\{|D_{si}|(b_i - z_i), |D_{si}|(z_i - a_i)\} \qquad i = 1, 2, \ldots, U \qquad (14)$$

where

$$D_{si} = \text{element of } \frac{\partial \psi^{-1}}{\partial y} \frac{\partial \psi}{\partial z} \qquad (15)$$

**Further Considerations and Recommendations**

In the previous section the basic steps of a reduced gradient algorithm are described. In this section we would like to recommend some improvements and point out further modifications that can be made to the basic algorithm.

As with most nonlinear programming techniques, a large amount of effort is expended during the line search. Because of the reduced gradient nature of the state and decision variable dichotomy, much time is spent during the line search performing iterations of Newton's method to adjust the state variables. This time can be reduced through a modification of Newton's method given here:

$$y^{i+1} = y^i - \frac{\partial \psi^{-1}}{\partial y} \psi(z^k, y^i) \qquad (16)$$

where

$\frac{\partial \psi^{-1}}{\partial y}$ = the initial inverse used at the start to calculate the reduced gradient.

$i$ = $i$th iteration of Newton's method

Although this modification does not preserve the convergence rate of the mathematical method, it does offer a considerable saving in time by avoiding successive reformulations of the Jacobian inverse.

The most stringent requirement of Newton's method is the need for a good initial approximation to the solution. Without it there is no guarantee of convergence. Abadie has suggested linearization of

the constraints at the start to provide this initial approximation. This was described in the previous section. The following procedure offers a slight improvement:

1  From the initial starting point of the line search $[\hat{z}, \hat{y}]^0$, a step $\alpha_1$ is taken along the search directions with the approximation of the new state variables taken as before

$$y^1 = y^0 + \alpha_1 P^1(\hat{y})$$

2  Locate the feasible point $\hat{y}^1$ by Newton's method.

3  All other state variable approximations are generated by the following pattern step whenever Newton's method converges for the previous point

$$\hat{y}^k = y^k + \frac{m_k}{m_{k-1}}(\hat{y}^{k-1} - y^k) \tag{47}$$

where

$\hat{y}^{k-1}$ = the feasible point obtained at the previous step
$m_{k-1}$ = the step length used for $[\hat{z}, \hat{y}]^{k-1}$.

This improvement in providing a starting point for Newton's



Fig. 4  Starting points for Newton's method by linearization and by direct



$\hat{z}_{1,2} = 0 \; ; \; y_1 = 0$

Fig. 5  State variable adjustment strategies

method decreases the number of Newton iterations required during the bounding phase, but is more readily exhibited during movement to a bound or a local minimum. Both approaches are pictured in Fig. 4.

The procedure outlined in (b) in the previous section on line searching differs from an approach that has been suggested by Abadie. In Fig. 5, we move from the point $[\hat{z}, \hat{y}]^0$ to the point $[\hat{z}, \hat{y}]^1$ as $y_1$ is adjusted. Since $y_1$ is nearly violating its lower bound, Abadie suggests the following scheme. Set $y_1$ to its lower bound and reclassify it as a decision variable, choosing an appropriate decision variable to take its place. Now adjust the new state variable to maintain feasibility. This would result in the point $[\hat{z}, \hat{y}]^2$ shown in Fig. 5 with $y_1$ exchanged with $z_1$. Following the scheme given previously would result in the point $[\hat{z}, \hat{y}]^4$.

Both approaches seem to have advantages and it is unclear which approach is best. The main advantage of Abadie's line search is that it may more quickly locate those points where constraints intersect. However, if the modified Newton's method described previously is used, then a reformulation of the Jacobian inverse would be required every time a state variable was respecified. The approach described here leads to a less complicated algorithm and avoids frequent reformulation of the Jacobian inverse.

### Results of Numerical Experiments With OPT

OPT is the name of our Reduced Gradient Fortran code, which implements the theory discussed in the previous sections. We report here testing of the program, which involved applications of the code to each of the Eason and Fenton test problems. Eason and Fenton gave results for ten problems, three of which were originally given by Colville. Four of these ten problems are essentially unconstrained, in that the constraints are neither active at the solution nor encountered during the search. If constraints had been encountered during the search, the numerical results would have been included. Since they were not, the results indicate the performance of the unconstrained searching scheme only, which was not the subject of this study. Accordingly, these problems were not considered further. OPT has no difficulty in identifying interior optima, as is demonstrated in our solution of Eason and Fenton problem number 6. It is certainly worthy of note that the basic concept of the reduced gradient, the state decision dichotomy, can be employed with any unconstrained searching algorithm. But this will be the subject of later reports.

Table 1 gives the results obtained with OPT on Eason and Fenton problems number 1–3, 6, 7, and 9 [1]. In each case OPT terminated at a point which matched or exceeded the convergence criteria specified by Eason and Fenton, using standard program parameters. Listed is the number of objective and constraint function combinations, and standard time needed to achieve solution.

OPT will terminate whenever the $L$-norm of the projected reduced gradient is less than $10^{-4}$, or when the $L$-norm of the change in the design variables is less than $10^{-5}$. The gradient condition controlled termination in each case except problem 3. The line searches in OPT are conducted to an accuracy of $10^{-5}$, and the tight constraints are satisfied to within $10^{-5}$ at every step of the search. Note that OPT found solutions for all of the problems. Table 2 gives the number of

| Table 1 | | | |
|---|---|---|---|
| Problem Number | $N_f$ | $N_g$ | Standard Time |
| 1 | 34 | 121 | 0.0182 |
| 2 | 24 | 85 | 0.0005 |
| 3 | 16 | 110 | 0.0010 |
| 6 | 47 | 141 | 0.00012 |
| 7 | 27 | 107 | 0.0010 |
| 9 | 46 | 141 | |

# Table 2 Reduced gradient stages required for solution

| Problem Number | Stages |
|---|---|
| 1 | 6 |
| 2 | 2 |
| 3 | 3 |
| 5 | 6 |
| 7 | 3 |
| 9 | 10 |

# Table 3 Relative ranking of optimization codes

(table illegible)

reduced gradient stages required to find each solution. Table 3 contains the rating schemes proposed by Eason and Fenton. $N_s$ is the number of problems solved. The Eason and Fenton results were recomputed for the six problems considered here, and codes with an $N_s$ value greater than or equal to three are included in Table 3. $I_a$ and $I_b$ are indicators of code efficiency defined by Eason and Fenton. $T_s$ is the sum of the normalized times for each code on all problems. Our CDC 6500 system is rated at 51 s using Colville's standard timing program. Problems 6 and 7 have infeasible starting points. OPT handles this difficulty with the artificial variable approach suggested by Abadie. The artificial variable is used only to find an initial feasible starting point and plays no part in the remainder of the search.

The excellent performance demonstrated by OPT is the result of the ease with which the generalized reduced gradient algorithm identifies and follows active constraints. Additionally, the performance here has been enhanced by careful attention to detail in the program implementation. Our experience indicates that even the most outstanding algorithm can perform poorly if the programming stage is improperly attended to.

Additional numerical results are reported in the thesis by Gabriele [17].

## Conclusions

The following conclusions can be drawn from this work:

1 The state and decision variable dichotomy helps to reduce the constrained nonlinear programming problem to a simpler, more readily solved problem. By using the dichotomy the problem becomes, in simple terms, a search in the reduced unconstrained feasible region of the decision variables while maintaining feasibility through state variable adjustment.

2 The reduced gradient is a reliable and efficient tool for the solution of a wide variety of nonlinear programming problems, since OPT receives the top rating in each of the categories listed.

## Acknowledgments

## References

1 Eason, E. D., and Fenton, R. G., "A Comparison of Numerical Optimization Methods for Engineering Design," JOURNAL OF ENGINEERING FOR INDUSTRY, TRANS. ASME, Series B, Vol. 96, No. 1, Feb. 1974, pp. 196-200.

2 Colville, A. R., "A Comparative Study of Nonlinear Programming Codes," Proceedings of the Princeton Symposium on Mathematical Programming, Kuhn, H. W., ed., Princeton, N. J., 1970, pp. 487-501.

3 Fiacco, A. V., and McCormick, G. P., Nonlinear Programming, Wiley, New York, 1968.

4 Griffith, R. E., and Stewart, R. A., "A Nonlinear Programming Technique for the Optimization of Continuous Processing Systems," Management of Science, Vol. 7, 1961, pp. 379-392.

5 Fox, R. L., Optimization Methods for Engineering Design, Addison-Wesley, 1971.

6 Himmelblau, D. M., Applied Nonlinear Programming, McGraw-Hill, New York, 1972.

7 Wolfe, P., "Methods for Linear Constraints," Nonlinear Programming, Abadie, J., ed., North Holland, Amsterdam, 1967, pp. 99-131.

8 Wolfe, P., "Methods of Nonlinear Programming," Recent Advances in Mathematical Programming, Graves, R. L., and Wolfe, P., eds., McGraw-Hill, New York, 1963, pp. 76-77.

9 Abadie, J., and Carpenter, J., "Generalization of the Wolfe Reduced Gradient Method to the Case of Nonlinear Constraints," Optimization, Fletcher, R., ed., Academic Press, 1969, p. 37.

10 Wilde, D., and Beightler, C., Foundations of Optimization, Prentice-Hall, Englewood Cliffs, N. J., 1967, Chapters 2, 3.

11 Kuhn, W. W., and Tucker, A. W., "Nonlinear Programming," Proceedings, 2nd Berkeley Symposium on Mathematical Statistics and Probability, University of California Press, Berkeley, 1951, pp. 481-492.

12 Fletcher, R., and Reeves, C. M., "Function Minimization by Conjugate Gradients," Computer Journal, Vol. 6, No. 2, 1963, pp. 163-168.

13 Davidon, W. C., "Variable Metric Method for Minimization," Argonne National Laboratory, ANL-5990 Rev., University of Chicago, 1959.

14 Ragsdell, K. M., "A Survey of Some Useful Optimization Methods," Proceedings of the Design Engineering Technical Conference, New York, 1974, pp. 129-135.

15 Ragsdell, K. M., Root, D. D., and Gabriele, G. A., "Newton's Method: A Classical Technique of Contemporary Value," Proceedings of the Design Engineering Technical Conference, New York, 1974, pp. 137-115.

16 Abadie, J., "Application of the GRG Algorithm to Optimal Control Problems," Integer and Nonlinear Programming, Abadie, J., ed., North Holland, Amsterdam, 1970, p. 191.

17 Gabriele, G. A., "Application of the Reduced Gradient Method to Optimal Engineering Design," Master's thesis, School of Mechanical Engineering, Purdue University, Dec. 1975.

# Sequential Gradient-Restoration Algorithm
## for the Minimization of Constrained Functions—
## Ordinary and Conjugate Gradient Versions[1]

A. MIELE,[2] H. Y. HUANG,[3] AND J. C. HEIDEMAN[4]

**Abstract.** The problem of minimizing a function $f(x)$ subject to the constraint $\varphi(x) = 0$ is considered. Here, $f$ is a scalar, $x$ an $n$-vector, and $\varphi$ a $q$-vector. A *sequential algorithm* is presented, composed of the alternate succession of gradient phases and restoration phases.

In the *gradient phase*, a nominal point $x$ satisfying the constraint is assumed; a displacement $\Delta x$ leading from point $x$ to a varied point $y$ is determined such that the value of the function is reduced. The determination of the displacement $\Delta x$ incorporates information at only point $x$ for the *ordinary gradient* version of the method (Part 1) and information at both points $x$ and $\bar{x}$ for the *conjugate gradient* version of the method (Part 2).

In the *restoration phase*, a nominal point $y$ not satisfying the constraint is assumed; a displacement $\Delta y$ leading from point $y$ to a varied point $\tilde{x}$ is determined such that the constraint is restored to a prescribed degree of accuracy. The restoration is done by requiring the least-square change of the coordinates.

If the stepsize $\alpha$ of the gradient phase is of $O(\varepsilon)$, then $\Delta x \approx O(\varepsilon)$ and $\Delta y \approx O(\varepsilon^2)$. For $\varepsilon$ sufficiently small, the restoration phase preserves the descent property of the gradient phase: the function $f$ decreases between any two successive restoration phases.

The ordinary gradient version of the algorithm exhibits asymptotic convergence, but not quadratic convergence. On the other hand, the conjugate gradient version exhibits quadratic convergence in the neighborhood of the minimum point. In particular, for a quadratic function subject to a linear constraint, the minimum point is obtained in no more than $n - q$ iterations.

## 1. Introduction

In recent years, considerable attention has been given to the iterative algorithms for minimizing a function $f(x)$ subject to the constraint $\varphi(x) = 0$, where $f$ is a scalar, $x$ an $n$-vector, and $\varphi$ a $q$-vector. With reference to the gradient method, one possible approach is that of the penalty functions. The advantage of this approach is that the constrained minimal problem is replaced by a mathematically simpler, unconstrained minimal problem. The disadvantages are these: no clear-cut method exists for choosing the penalty constants; the algorithm must be repeated several times for increasing values of the penalty constants; the values of the function between iterations are not comparable, since the constraints are not satisfied; and, even when the algorithm is terminated, the constraints are generally not satisfied.

In this paper, we present a *sequential algorithm* constructed in such a way that the values of the function $f(x)$ between iterations are comparable. The algorithm is composed of the alternate succession of gradient phases and restoration phases (Fig. 1). In the *gradient phase*, a nominal point $x$ satisfying the constraint is assumed; then, a displacement $\Delta x$ leading from point $x$ to a varied point $y$ is determined by minimizing the first variation $\delta f(x)$ subject to the linearized constraint $\delta \varphi(x) = 0$ and a quadratic constraint on $\Delta x$. Due to the fact that the constraint is accounted for only to first order, the varied point $y$ may be such that $\varphi(y) \neq 0$. This being the case, a restoration phase is needed prior to starting the next gradient phase. In this *restoration phase*, a nominal point $y$ not satisfying the constraint is assumed; then, a displacement $\Delta y$ leading from point $y$ to a varied point $\tilde{x}$ is determined by minimizing the distance $\Delta y^T \Delta y$ subject to the linearized constraint $k\varphi(x) = \delta \varphi(y) = 0$, where $0 < k \leq 1$. If a single restoration cycle fails to produce the required degree of accuracy in the constraint, several cycles must be employed, as shown in Fig. 2. After the restoration phase is finished, the iteration is completed, and the next iteration is started using $\tilde{x}$ as the nominal point $x$.

Fig. 1



Fig. 2

Two versions of the method are presented: (a) the *sequential ordinary gradient-restoration algorithm*, in which the gradient phase uses information at point $x$ only, and (b) the *sequential conjugate gradient-restoration algorithm*, in which the gradient phase uses information at both points $x$ and $\hat{x}$, where $\hat{x}$ is the point preceding $x$. Method (a) is given in Part 1; Method (b) is given in Part 2; then, in Part 3, several numerical examples are presented.

## 2. Statement of the Problem

We consider the problem of minimizing the function

$$f = f(x) \tag{1}$$

subject to the constraint

$$\varphi(x) = 0 \tag{2}$$

In the above equation, $f$ is a scalar, $x$ an $n$-vector, and $\varphi$ a $q$-vector,[a] where $q < n$. It is assumed that the first and second partial derivatives of the function $f$ and the constraint $\varphi$ with respect to $x$ exist and are continuous; it is also assumed that the constrained minimum exists.

## PART 1: SEQUENTIAL ORDINARY GRADIENT-RESTORATION ALGORITHM

In this part, we present the ordinary gradient version of the sequential gradient-restoration algorithm. The gradient phase is treated in Section 3 and the restoration phase in Section 4. In Section 5, an order-of-magnitude analysis is presented. Finally, in Section 6, the sequential ordinary gradient-restoration algorithm is summarized.

## 3. Gradient Phase

Consider a displacement $\Delta x$ leading from a nominal point $x$ to a varied point $y$ such that

$$y = x + \Delta x \tag{3}$$

[a] All the vectors in this paper are column vectors.

Assume that the nominal point $x$ satisfies (2) exactly and that the varied point $y$ satisfies (2) to first order. The first-order change of the function (1) is given by

$$\delta f(x) = f_x^T(x)\, \Delta x \tag{4}$$

where $f_x(x)$ is the gradient of the scalar function $f$ with respect to the vector $x$ and the symbol $T$ denotes the transpose of a matrix. In turn, the first-order change of the constraint (2) is represented as

$$\delta \varphi(x) = \varphi_x^T(x)\, \Delta x = 0 \tag{5}$$

where $\varphi_x(x)$, an $n \times q$ matrix, denotes the gradient of the vectorial function $\varphi$ with respect to the vector $x$.

Next, consider the following quadratic constraint on the displacement $\Delta x$:

$$K = \Delta x^T \Delta x \tag{6}$$

where $K$ is a constant. With this understanding, we formulate the following problem: Find the variation $\Delta x$ which minimizes (4) subject to (5) (6).

3.1. Displacement $\Delta x$. Standard methods of the theory of maxima and minima show that the *fundamental function* of this problem is the scalar function

$$\Omega = f_x^T(x)\, \Delta x + \lambda^T \varphi_x^T(x)\, \Delta x + (1/2\alpha)\, \Delta x^T \Delta x \tag{7}$$

where $1/2\alpha$ is a scalar Lagrange multiplier and $\lambda$ a $q$-vector Lagrange multiplier. If one introduces the *augmented function*

$$F(x, \lambda) = f(x) + \lambda^T \varphi(x) \tag{8}$$

and observes that

$$F_x(x, \lambda) = f_x(x) + \varphi_x(x)\, \lambda \tag{9}$$

the fundamental function (7) becomes

$$\Omega = F_x^T(x, \lambda)\, \Delta x + (1/2\alpha)\, \Delta x^T \Delta x \tag{10}$$

In Eqs. (9)–(10), the symbol $F_x(x, \lambda)$ denotes the gradient of the augmented function $F$ with respect to the vector $x$. The optimal displacement $\Delta x$ satisfies the relation

$$\Omega_{\Delta x} = 0 \tag{11}$$

where $\Omega_{\Delta x}$ denotes the gradient of the fundamental function $\Omega$ with respect to the vector $\Delta x$. The explicit form of (11) is the following:

$$\Delta x = -\alpha F_x(x, \lambda) \qquad (12)$$

and shows that the displacement vector $\Delta x$ has the same direction as the gradient of the augmented function $F$. Note that (12) determines $\Delta x$ providing $\lambda$, $\alpha$ are specified.

3.2. Relation Between $K$ and $\alpha$. As Eq. (12) shows, the displacement $\Delta x$ is proportional to $\alpha$, the *stepsize* of the gradient phase. Upon substituting (12) into (6), we see that

$$K = \alpha^2 F_x{}^T(x, \lambda) F_x(x, \lambda) \qquad (13)$$

Therefore, a correspondence exists between the values of the constant $K$ and the values of the stepsize $\alpha$. This being the case, one can bypass prescribing $K$ and reason directly on $\alpha$, as in the considerations which follow.

3.3. Determination of $\lambda$. If Eqs. (5), (9), (12) are combined, we obtain the relation

$$\varphi_x{}^T(x) f_x(x) + \varphi_x{}^T(x) \varphi_x(x) \lambda = 0 \qquad (14)$$

Since the vector $f_x(x)$ and the matrix $\varphi_x(x)$ are known at the nominal point $x$, Eq. (14) supplies the multiplier $\lambda$; this is precisely the value which guarantees satisfaction of the constraint (2) to first order.

3.4. Descent Property of the Gradient Phase. The first variation of the augmented function is given by

$$\delta F(x, \lambda) = F_x{}^T(x, \lambda) \, \Delta x \qquad (15)$$

which, in the light of (12), can be written as

$$\delta F(x, \lambda) = -\alpha F_x{}^T(x, \lambda) F_x(x, \lambda) \qquad (16)$$

Equation (16) shows that, for $\alpha > 0$, the first variation of the augmented function is negative. Therefore, if $\alpha$ is sufficiently small, the augmented function $F$ decreases during the gradient phase.

Alternatively, the first variation of $F(x, \lambda)$ can be written as

$$\delta F(x, \lambda) = \delta f(x) + \lambda^T \delta \varphi(x) \qquad (17)$$

Because of (5), Eq. (17) reduces to the form

$$\delta f(x) = \delta F(x, \lambda) \qquad (18)$$

which states that the functions $f(x)$ and $F(x, \lambda)$ behave identically, to first order. Therefore, the descent property also holds for the function $f$.

3.5. Stepsize. The next step is to assign a value to the stepsize $\alpha$. If Eqs. (3) and (12) are combined, the position vector at the end of the gradient phase becomes

$$y = x - \alpha F_x(x, \lambda) \qquad (19)$$

Since $\lambda$ is known through Eq. (14), Eq. (19) defines a one-parameter family of points $y$ for which the augmented function $F$ takes the form

$$F(y, \lambda) = F(x - \alpha F_x(x, \lambda), \lambda) = \Psi(\alpha) \qquad (20)$$

The greatest decrease of the function $\Psi(\alpha)$ occurs if the parameter $\alpha$ satisfies the following necessary condition:

$$\Psi_\alpha(\alpha) = 0 \qquad (21)$$

After observing that

$$\Psi_\alpha(\alpha) = -F_x{}^T(y, \lambda) F_x(x, \lambda) \qquad (22)$$

we see that Eq. (21) can be written as

$$F_x{}^T(y, \lambda) F_x(x, \lambda) = 0 \qquad (23)$$

showing that the gradient of the augmented function at point $y$ is orthogonal to the gradient at point $x$.

To obtain satisfaction of (21) or (23), some one-dimensional search method must be employed. In particular, *cubic interpolation* (it employs first derivatives only) and *quasilinearization* (it employs both first and second derivatives) are powerful methods (Refs. 1 and 3). These methods are to be employed iteratively until Eq. (21) is satisfied to a desired degree of accuracy, that is, until

$$\Psi_\alpha{}^2(\alpha) \leq \theta_1 \qquad (24)$$

where $\theta_1$ is a small number. In practice, $\theta_1$ can be fixed or one may choose

$$\theta_1 = \epsilon_1 \Psi_\alpha{}^2(0) \qquad (25)$$

where $\epsilon_1$ is a small number.

## 4. Restoration Phase

At the end of the gradient phase, the point $y$ is known. If the constraint is linear, the relation $\varphi(y) = 0$ holds. On the other hand, if the constraint is nonlinear, $\varphi(y) \neq 0$, which means that some degree of dissatisfaction exists. Therefore, a restoration phase is needed prior to starting the next gradient phase. Specifically, one has to apply a small variation $\Delta y$ to $y$ to generate a new position vector

$$\tilde{x} = y + \Delta y \tag{26}$$

such that $\varphi(\tilde{x}) = 0$. While there are infinite ways to perform the restoration, the most logical is that developed in Ref. 4: the constraint is restored to a prescribed degree of accuracy with the least-square change of the position vector.

If quasilinearization is employed, Eq. (2) is approximated by

$$\varphi(y) + \varphi_x^T(y) \, \Delta y = 0 \tag{27}$$

In order to prevent the variation $\Delta y$ from becoming too large, we imbed Eq. (27) into the one-parameter family of equations

$$k\varphi(y) + \varphi_x^T(y) \, \Delta y = 0 \tag{28}$$

where

$$0 \leq k \leq 1 \tag{29}$$

denotes a scaling factor.

In the light of previous discussion, we seek the minimum of the function

$$J = \tfrac{1}{2} \Delta y^T \, \Delta y \tag{30}$$

subject to the linearized constraint (28). Standard methods of the theory of maxima and minima show that the *fundamental function* of this problem is given by

$$\omega = \tfrac{1}{2} \Delta y^T \, \Delta y + \sigma^T [k\varphi(y) + \varphi_x^T(y) \, \Delta y] \tag{31}$$

where $\sigma$, a $q$-vector, denotes an undetermined, constant Lagrange multiplier. The optimal change $\Delta y$ satisfies the relation

$$\omega_{\Delta y} = 0 \tag{32}$$

where $\omega_{\Delta y}$ denotes the gradient of the fundamental function $\omega$ with respect to the vector $\Delta y$. The explicit form of (32) is the following:

$$\Delta y = -\varphi_x(y) \, \sigma \tag{33}$$

The Lagrange multiplier $\sigma$ is obtained by combining (28) and (33) to eliminate $\Delta y$. This yields the relation

$$k\varphi(y) - \varphi_x^T(y) \, \varphi_x(y) \, \sigma = 0 \tag{34}$$

For any given $k$ in the range (29), Eq. (34) supplies the Lagrange multiplier vector $\sigma$. Once $\sigma$ is known, the correction $\Delta y$ is given by (33) and the corrected position vector $\tilde{x}$ is given by (26). Of course, the restoration phase must be performed iteratively until a desired degree of accuracy is obtained, that is, until the inequality

$$P(\tilde{x}) \leq \theta_2 \tag{35}$$

is satisfied, where $\theta_2$ is a small number and the performance index

$$P(x) = \varphi^T(x) \, \varphi(x) \tag{36}$$

measures the error in the constraint.

**4.1. Descent Property of the Performance Index.** The first variation of the performance index $P(y)$ is given by

$$\delta P(y) = 2\varphi^T(y) \, \varphi_x^T(y) \, \Delta y \tag{37}$$

and, because of Eq. (28), reduces to

$$\delta P(y) = -2k\varphi^T(y) \, \varphi(y) \tag{38}$$

which, in turn, can be written as

$$\delta P(y) = -2kP(y) \tag{39}$$

Since $P(y) > 0$, Eq. (39) shows that the first variation of the performance index is negative for $k > 0$. Therefore, if $k$ is sufficiently small, the decrease of the performance index is guaranteed. In practice, one can use $k = 1$. If this value of $k$ does not result in a decrease in $P$, then $k$ is successively bisected until $P$ is decreased.

## 5. Order-of-Magnitude Analysis

The position vector $\tilde{x}$ at the end of the restoration phase and the position vector $x$ at the beginning of the gradient phase are related by

$$\tilde{x} = x + \Delta x + \Delta y \tag{40}$$

where $\Delta x$ is the gradient displacement and $\Delta y$ is the restoration displacement. From Eq. (12), we see that, if the stepsize $\alpha$ is of $O(\epsilon)$, the gradient displacement has the order

$$\Delta x = O(\epsilon) \qquad (41)$$

We observe that, to second order,

$$\varphi(y) = \varphi(x) + \varphi_x{}^T(x)\,\Delta x + \tfrac{1}{2}\,\Delta x^T\,\varphi_{xx}(x)\,\Delta x \qquad (42)$$

where $\varphi_{xx}(x)$ denotes the array of the second partial derivatives of $\varphi$. Since the nominal point $x$ satisfies (2) exactly and the variation $\Delta x$ satisfies (2) to first order, the first two terms on the right-hand side of (42) vanish. Therefore, in the light of (41)-(42), we conclude that

$$\varphi(y) = O(\epsilon^2) \qquad (43)$$

Next, we turn our attention to Eqs. (33)-(34) and observe that, because of (43),

$$\Delta y = O(\epsilon^2) \qquad (44)$$

In conclusion, if the gradient displacement is of $O(\epsilon)$, the restoration displacement is of $O(\epsilon^2)$. This guarantees that, for sufficiently small $\epsilon$,

$$|\Delta y| \ll |\Delta x| \qquad (45)$$

**5.1. Descent Property of the Algorithm.** Finally, we consider the points $x$ and $\tilde{x}$, both satisfying the constraint (2). To first order, the difference of the values of the function $f$ at these points is given by

$$f(\tilde{x}) - f(x) = f_x{}^T(x)(\Delta x + \Delta y) \qquad (46)$$

On account of (45), the second term on the right-hand side of (46) can be neglected with respect to the first. Hence, Eq. (46) becomes

$$f(\tilde{x}) - f(x) = -\alpha F_x{}^T(x, \lambda) F_x(x, \lambda) \qquad (47)$$

Therefore, for $\epsilon$ sufficiently small, the restoration algorithm preserves the descent property of the gradient algorithm: the function $f$ decreases between any two successive restoration phases.

## 6. Summary of the Algorithm

The algorithm presented in Part 1 consists of the alternate succession of gradient phases and restoration phases. A summary of the algorithm is given below.

**6.1. Gradient Phase.** For this phase, the algorithm is represented by

$$\lambda = -[\varphi_x{}^T(x)\,\varphi_x(x)]^{-1}\,\varphi_x{}^T(x) f_x(x)$$
$$F_x(x, \lambda) = f_x(x) + \varphi_x(x)\lambda \qquad (48)$$
$$\Delta x = -\alpha F_x(x, \lambda)$$
$$y = x + \Delta x$$

The sequence of operations is as follows: (a) select a nominal point $x$ such that $\varphi(x) = 0$; (b) at this point, determine the vector $f_x(x)$ and the matrix $\varphi_x(x)$; (c) compute the multiplier $\lambda$ with (48-1) and the vector $F_x(x, \lambda)$ with (48-2); (d) determine the optimal stepsize $\alpha$ by a one-dimensional search in which either $\Psi(\alpha) = f(y)$ or $\Psi(\alpha) = F(y, \lambda)$ is minimized along the search direction $F_x(x, \lambda)$; the search is terminated when Ineq. (24) is satisfied; (e) compute the displacement $\Delta x$ with (48-3) and the varied point $y$ with (48-4).

**6.2. Restoration Phase.** For this phase, the algorithm is represented by

$$\sigma = \lambda[\varphi_x{}^T(y)\,\varphi_x(y)]^{-1}\,\varphi(y)$$
$$\Delta y = -\varphi_x(y)\sigma \qquad (49)$$
$$\tilde{x} = y + \Delta y$$

The sequence of operations is as follows: (a) at point $y$, compute the vector $\varphi(y)$ and the matrix $\varphi_x(y)$; (b) assuming $k = 1$, determine the multiplier $\sigma$ with (49-1), the displacement $\Delta y$ with (49-2), and the varied point $\tilde{x}$ with (49-3); (c) if $P(\tilde{x}) < P(y)$, the scaling factor $k = 1$ is acceptable; if $P(\tilde{x}) > P(y)$, the previous value of $k$ must be replaced by some smaller value in the range (29) until the condition $P(\tilde{x}) < P(y)$ is met; this can be achieved through successive bisections of $k$; (d) return to step (a) and repeat the restoration algorithm using $\tilde{x}$ as the starting point $y$ for the subsequent iteration; (e) terminate the restoration algorithm when the stopping condition (35) is satisfied; (f) once the restoration algorithm is completed, verify the inequality

$$f(\tilde{x}) < f(x) \qquad (50)$$

If Ineq. (50) is satisfied, start the next gradient phase. If Ineq. (50) is violated, return to the previous gradient phase and reduce the stepsize $\alpha$ until, after restoration, Ineq. (50) is satisfied.

### 6.3. Stopping Condition. The algorithm is terminated when

$$Q(\tilde{x}, \lambda) < \theta_1 \tag{51}$$

where $\theta_1$ is a small number and

$$Q(x, \lambda) = F_x{}^T(x, \lambda) F_x(x, \lambda) \tag{52}$$

measures the error in the optimal conditions.

## PART 2: SEQUENTIAL CONJUGATE GRADIENT-RESTORATION ALGORITHM

In this part, we present the conjugate gradient version of the sequential gradient-restoration algorithm. The gradient phase is discussed in Section 7 in general; the case of a quadratic function subject to a linear constraint is given in Section 8; then, the practical method to be used for a nonquadratic function subject to a nonlinear constraint is given in Section 9. The treatment of the restoration phase is omitted, since it is covered in Section 4. In Section 10, an order-of-magnitude analysis is presented. Finally, in Section 11, the sequential conjugate gradient-restoration algorithm is summarized.

## 7. Gradient Phase: General Discussion

Consider a displacement $\Delta x$ leading from a nominal point $x$ to a varied point $y$ given by Eq. (3). Assume that the nominal point $x$ satisfies (2) exactly and that the varied point $y$ satisfies (2) to first order. The first-order change of the function (1) is given by Eq. (4). In turn, the first-order change of the constraint (2) is represented by Eq. (5). Next, consider the following quadratic constraint on the displacement $\Delta x$:

$$K = (\Delta x - \beta \Delta \tilde{x})^T(\Delta x - \beta \Delta \tilde{x}) \tag{53}$$

where $K$ and $\beta$ are constants and $\Delta \tilde{x}$ is the displacement of the previous gradient phase, that is, the displacement leading from the nominal point $\tilde{x}$ to the varied point $y$ (see Figs. 1 and 2). With this understanding, we formulate the following problem: Find the displacement $\Delta x$ which minimizes (4) subject to (5) and (53).

### 7.1. Displacement $\Delta x$. Standard methods of the theory of maxima and minima show that the *fundamental function* of this problem is the scalar function

$$\Omega = f_x{}^T(x)\, \Delta x + \lambda^T \varphi_x{}^T(x)\, \Delta x + (1/2\alpha)(\Delta x - \beta \Delta \tilde{x})^T(\Delta x - \beta \Delta \tilde{x}) \tag{54}$$

where $1/2\alpha$ is a scalar Lagrange multiplier and $\lambda$ a $q$-vector Lagrange multiplier. If one introduces the *augmented function* (8) and its gradient (9), the fundamental function (54) becomes

$$\Omega = F_x{}^T(x, \lambda)\, \Delta x + (1/2\alpha)(\Delta x - \beta \Delta \tilde{x})^T(\Delta x - \beta \Delta \tilde{x}) \tag{55}$$

The optimal displacement $\Delta x$ satisfies the relation (11), whose explicit form is the following:

$$\Delta x = -\alpha F_x(x, \lambda) + \beta \Delta \tilde{x} \tag{56}$$

If one defines the search direction $p$ from

$$\Delta x = -\alpha p \tag{57}$$

and observes that, for the previous iteration,

$$\Delta \tilde{x} = -\tilde{\alpha} \tilde{p} \tag{58}$$

the following relation ensues from (56)–(58):

$$p = F_x(x, \lambda) + \gamma \tilde{p} \tag{59}$$

where

$$\gamma = \beta \tilde{\alpha}/\alpha \tag{60}$$

In conclusion, the displacement $\Delta x$ during the gradient phase is given by (57), with $p$ governed by Eq. (59), where $\tilde{p}$ is known from the previous iteration. Note that (57) and (59) determine $\Delta x$ providing $\lambda, \alpha, \gamma$ are specified.

### 7.2. Relation Between $K$ and $\alpha$. As Eq. (57) shows, the displacement $\Delta x$ is proportional to $\alpha$, the *stepsize* of the gradient phase. Upon substituting (56) into (53), we obtain Eq. (13). Therefore, a correspondence exists between the values of the constant $K$ and the values of the stepsize $\alpha$. This being the case, one can bypass prescribing $K$ and reason directly on $\alpha$, as in the considerations which follow.

**7.3. Determination of $\lambda$, $\alpha$, $\gamma$.** If Eqs. (5), (9), (57), (59) are combined, we obtain the relation

$$\varphi_x^T(x) f_x(x) + \varphi_x^T(x) \varphi_x(x) \lambda + \gamma \varphi_x^T(x) \dot{p} = 0 \tag{61}$$

which ensures satisfaction of the constraint (2) to first order. Equation (61) is a linear relation between $\lambda$ and $\gamma$ and admits the solution

$$\lambda = -[\varphi_x^T(x) \varphi_x(x)]^{-1}[\varphi_x^T(x) f_x(x) + \gamma \varphi_x^T(x) \dot{p}] \tag{62}$$

Therefore, the rate of change of $\lambda$ with respect to $\gamma$ is given by the vector equation

$$\lambda_\gamma = -[\varphi_x^T(x) \varphi_x(x)]^{-1} \varphi_x^T(x) \dot{p} \tag{63}$$

The next step is to assign values to $\alpha$ and $\gamma$. If Eqs. (3), (57), (59) are combined, the position vector at the end of the gradient phase becomes

$$y = x - \alpha F_x(x, \lambda) - \alpha \gamma \dot{p} \tag{64}$$

Since $\lambda$ depends on $\gamma$ through Eq. (62), Eq. (64) defines a two-parameter family of points $y$ for which the augmented function $F$ takes the form

$$F(y, \lambda) = F(x - \alpha F_x(x, \lambda) - \alpha \gamma \dot{p}, \lambda) = \Psi(\alpha, \gamma) \tag{65}$$

The greatest decrease of the function $\Psi(\alpha, \gamma)$ occurs if the parameters $\alpha$, $\gamma$ satisfy the following necessary conditions:

$$\Psi_\alpha(\alpha, \gamma) = 0, \qquad \Psi_\gamma(\alpha, \gamma) = 0 \tag{66}$$

After observing that

$$\Psi_\alpha(\alpha, \gamma) = -F_x^T(y, \lambda) \dot{p}$$
$$\Psi_\gamma(\alpha, \gamma) = -\alpha F_x^T(y, \lambda) \dot{p} + [\alpha F_x^T(y, \lambda) \varphi_x(x) - \varphi^T(y)] \lambda_\gamma \tag{67}$$

we see that Eqs. (66) can be written as

$$F_x^T(y, \lambda) \dot{p} = 0, \quad \alpha F_x^T(y, \lambda) \dot{p} + [\alpha F_x^T(y, \lambda) \varphi_x(x) - \varphi^T(y)] \lambda_\gamma = 0 \tag{68}$$

In the light of (59) and (62)–(64), Eqs. (68) constitute a system of two scalar equations in the unknowns $\alpha$ and $\gamma$. Once the stepsize $\alpha$ and the coefficient $\gamma$ are known from (68), the multiplier $\lambda$ follows from (62), the search direction $p$ from (59), the displacement $dx$ from (57), and the position vector $y$ from (3). Thus, the problem of determining $\lambda$, $\alpha$, $\gamma$ is solved in principle; computationally, however, further simplifications are needed to make the algorithm practical.

## 8. Gradient Phase: Quadratic Function, Linear Constraint

Now, consider the particular case of a quadratic function and a linear constraint given in the form

$$f(x) = a + b^T x + \tfrac{1}{2} x^T c x, \qquad \varphi(x) = d + s^T x \tag{69}$$

where $a$ is a scalar, $b$ is an $n$-vector, $c$ an $n \times n$ symmetric matrix, $d$ a $q$-vector, and $s$ an $n \times q$ matrix. Here, all the coefficients are constant. The gradients of the functions $f$ and $\varphi$ become

$$f_x(x) = b + cx, \quad \varphi_x(x) = s \tag{70}$$

Because of the linearity, the constraint is never violated during the gradient phase and, consequently,

$$dy = 0, \quad y = x + dx, \quad \varphi_x^T(x) \dot{p} = 0 \tag{71}$$

This being the case, Eq. (62) supplies the following expression for the multiplier:

$$\lambda = -[\varphi_x^T(x) \varphi_x(x)]^{-1} \varphi_x^T(x) f_x(x) \tag{72}$$

which is now independent of $\gamma$, that is,

$$\lambda_\gamma = 0 \tag{73}$$

The relations (68) optimizing $\alpha$ and $\gamma$ become

$$F_x^T(y, \lambda) \dot{p} = 0, \quad F_x^T(y, \lambda) \dot{p} = 0 \tag{74}$$

showing that the gradient of the augmented function at point $y$ is orthogonal to both the present and previous search directions. A mathematical consequence of (59) and (74) is that

$$F_x^T(y, \lambda) F_x(x, \lambda) = 0 \tag{75}$$

showing that the gradients at point $y$ and point $x$ are orthogonal. Furthermore, after laborious manipulations, Eqs. (74)·(75) lead to

$$F_x^T(x, \lambda) \dot{p} = 0, \quad F_x^T(x, \lambda) \dot{p} = 0, \quad F_x^T(x, \lambda) F_x(x, \lambda) = 0 \tag{76}$$

For a quadratic function subject to a linear constraint, the following relationship can be shown to hold:

$$F_x(y, \lambda) = F_x(x, \lambda) - \alpha c p \tag{77}$$

Invoking (69)–(77), we see that (74-1) yields the following solution for the optimal stepsize:

$$\alpha = F_x{}^T(x, \lambda)\, ^c{}_x(x, \lambda)/p^T c p \qquad (78)$$

Furthermore, (74-2) leads to

$$p^T c \dot{p} = 0 \qquad (79)$$

which states that the search directions $p$ and $\dot{p}$ are conjugate with respect to the matrix $c$. In turn, (79) yields the following explicit solution for $y$:

$$y = F_x{}^T(x, \lambda) F_x(x, \lambda)/F_x{}^T(\hat{x}, \hat{\lambda}) F_x(\hat{x}, \hat{\lambda}) \qquad (80)$$

In conclusion, for a given nominal point $x$, the multiplier $\lambda$ is supplied by (72), the coefficient $y$ by (80), the search direction $p$ by (59), the optimum stepsize $\alpha$ by (78), the displacement $\Delta x$ by (57), and the position vector $y$ by (3).

**8.1. Convergence Properties.** For a quadratic function subject to a linear constraint, the following relations can be shown to hold *providing the first step of the algorithm is a gradient step:*

$$F_x{}^T(x, \lambda) F_x(x_*, \lambda_*) = 0, \qquad F_x{}^T(x, \lambda) p_* = 0, \qquad p^T c p_* = 0 \qquad (81)$$

where $x_*$ denotes any state preceding $x$. Equations (81) can be derived from (76) and (79) through mathematical induction. The first of Eqs. (81) states that the gradient at each point is orthogonal to the gradient at every previous point. The second of Eqs. (81) states that the gradient at each point is orthogonal to the search direction at every previous point. Finally, the third of Eqs. (81) states that the search direction at each point and the search direction at every previous point are conjugate with respect to the constant matrix $c$; this is why the algorithm is called the *conjugate-gradient algorithm.*

Laborious manipulations, omitted for the sake of brevity, show that the algorithm defined by (3), (57), (59) with $\lambda$, $\alpha$, $y$ defined by (72), (78), (80) reduces the gradient $F_x(x, \lambda)$ to zero in no more than $n - q$ steps; therefore, the minimum of the function $f(x)$ subject to the constraint $q(x) = 0$ is reached in no more than $n - q$ steps. If the function is unconstrained, that is, if $q = 0$, then the minimum of $f(x)$ is reached in no more than $n$ steps (see Refs. 5 and 6).

## 9. Gradient Phase: Nonquadratic Function and/or Nonlinear Constraint

If the function $f(x)$ is nonquadratic and/or the constraint $q(x)$ is nonlinear, the relations (72), (78), (80) defining $\lambda$, $\alpha$, $y$ are not simultaneously valid. We observe that (72) and (80) involve first derivatives only, while (78) involves the second-derivative matrix $c$. This being the case, we choose to discard (78) and retain (72) and (80). Thus, we employ the algorithm

$$\lambda = -[\varphi_x{}^T(x)\,\varphi_x(x)]^{-1}\,\varphi_x{}^T(x) f_x(x)$$
$$F_x(x, \lambda) = f_x(x) + \varphi_x(x)\lambda$$
$$y = F_x{}^T(x, \lambda) F_x(x, \lambda)/F_x{}^T(\hat{x}, \hat{\lambda}) F_x(\hat{x}, \hat{\lambda}) \qquad (82)$$
$$p = F_x(x, \lambda) + y\dot{p}$$
$$\Delta x = -\alpha p$$
$$y = x + \Delta x$$

which, for the unconstrained case, reduces to the well-known Fletcher–Reeves algorithm (see Refs. 3 and 7). For the constrained case, the justification of the expressions for the multiplier $\lambda$ and the coefficient $y$ is presented in Section 10. Once the nominal point $x$ is given, the multiplier $\lambda$ is determined with (82-1), the gradient $F_x(x, \lambda)$ with (82-2), the coefficient $y$ with (82-3), and the search direction $p$ with (82-4); for a given stepsize $\alpha$, the displacement $\Delta x$ is computed with (82-5) and the varied point $y$ with (82-6). The determination of $\alpha$ is discussed in the following section.

**9.1. Stepsize.** If Eqs. (82-5) and (82-6) are combined, the position vector at the end of the gradient phase becomes

$$y = x - \alpha p \qquad (83)$$

For a given nominal point $x$, the vector $p$ is known through Eqs. (82-1)–(82-4); therefore, Eq. (83) defines a one-parameter family of points $y$ for which the augmented function $F(y, \lambda)$ takes the form

$$F(y, \lambda) = F(x - \alpha p, \lambda) = \Psi(\alpha) \qquad (84)$$

The greatest decrease in the function $\Psi(\alpha)$ occurs if the parameter $\alpha$ satisfies the following necessary condition:

$$\Psi_\alpha(\alpha) = 0 \qquad (85)$$

After observing that

$$\Psi_\alpha(\alpha) = -F_x^T(y, \lambda)\, p \tag{86}$$

we see that Eq. (85) can be written as

$$F_x^T(y, \lambda)\, p = 0 \tag{87}$$

showing that the gradient of the augmented function at point $y$ is orthogonal to the search direction $p$.

To obtain satisfaction of (85) or (87), some one-dimensional search method must be employed. In particular, *cubic interpolation* (it employs first-derivatives only) and *quasilinearization* (it employs both first and second derivatives) are powerful methods (Refs. 1–3). These methods are to be employed iteratively until Eq. (85) is satisfied to a desired degree of accuracy, that is, until Ineq. (24) is satisfied.

**9.2. Starting the Algorithm.** The algorithm (82) requires that the search direction $p$ be known from the previous iteration. Since this is not the case for the first iteration, some assumption concerning $\gamma p$ is needed in order to start the algorithm. To retain quadratic convergence, one must choose $p = 0$ or $\gamma = 0$. Consequently, for the first step, the search direction (82-4) becomes

$$p = F_x(x, \lambda) \tag{88}$$

meaning that the first step is a pure gradient step.

**9.3. Restarting the Algorithm.** For a quadratic function subject to a linear constraint, the present algorithm converges to the exact minimum in no more than $n - q$ steps. For the general case, this suggests the idea of restarting the algorithm every $\Delta N = n - q$ or $\Delta N = n - q + 1$ steps.

## 10. Order-of-Magnitude Analysis

The sequential conjugate gradient-restoration algorithm is represented by Eqs. (62) and (49). While Eqs. (49) are exact, Eqs. (82) are an approximation to the true optimal conditions. In this section, an estimate of the error involved in the computation of the Lagrange multiplier $\lambda$ is given; furthermore, a verification of the descent properties of the algorithm is presented.

**10.1. Lagrange Multiplier.** Here, we assume that, if $\alpha$ is a small quantity,

$$\alpha = O(\epsilon) \tag{89}$$

for every gradient phase. Concerning the multiplier $\lambda$, we note that the exact equation (62) has been replaced by the approximate equation (82-1); therefore, the $q$-vector

$$R_1 = \gamma \varphi_x^T(x)\, p \tag{90}$$

is representative of the order of magnitude of the error in $\lambda$. This is due to the fact that the terms $\varphi_x^T(x)\, f_x(x)$ and $\varphi_x^T(x)\, \varphi_x(x)$ in Eq. (62) can be regarded to be of $O(1)$.

If $\lambda$ is computed with the approximate equation (82-1), the first-order change of the constraint (2) is not exactly zero and, because of Eqs. (82), is given by

$$\Delta \varphi(x) = -\alpha \gamma \varphi_x^T(x)\, p \tag{91}$$

If one observes that $\varphi(x) = 0$ and that, to first order,

$$\Delta \varphi(x) = \varphi(y) - \varphi(x) = \varphi(y) \tag{92}$$

it follows that

$$\varphi(y) = -\alpha \gamma \varphi_x^T(x)\, p \tag{93}$$

If Eqs. (49) and (93) are combined, the displacement $\Delta y$ associated with the restoration phase can be written as

$$\Delta y = \alpha \gamma \varphi_x(y)[\varphi_x^T(y)\, \varphi_x(y)]^{-1} \varphi_x^T(x)\, p \tag{94}$$

To first order, the expansion of the gradient $\varphi_x(\tilde{x})$ is given by

$$\varphi_x(\tilde{x}) = \varphi_x(x) + \varphi_{xx}(x)(dx + dy) \tag{95}$$

where $\varphi_{xx}(x)$ denotes the array of the second partial derivatives of $\varphi$. If the transposes of both sides of (95) are postmultiplied by $p$ and if Eqs. (82-5) and (94) are accounted for, one deduces that

$$\varphi_x^T(\tilde{x})\, p = \varphi_x^T(x)\, p - \alpha p^T \varphi_{xx}(x)\, p + \alpha \gamma k p^T \varphi_{xx}(x)\, \varphi_x(y)[\varphi_x^T(y)]^{-1} \varphi_x(y)]^{-1} \varphi_x^T(x)\, p \tag{96}$$

Because of (82-1) and (82-4), the first term on the right-hand side of (96) can be written as

$$\varphi_x^T(x)\, p = \gamma \varphi_x^T(x)\, p \tag{97}$$

As a consequence, the third term on the right-hand side of (96) is negligible with respect to the first and Eq. (96) can be approximated by

$$\varphi_x^T(x)\, p = \gamma\varphi_x^T(x)\, \hat{p} - \alpha\rho^T\varphi_{xx}(x)\, p \qquad (98)$$

This is the key recurrence formula necessary to estimate the order of magnitude of the terms of the form $\varphi_x^T(x)\hat{p}$ or $\varphi_x^T(x)p$. Since $\alpha$ is of $O(\varepsilon)$, we see that, for the first step of the algorithm ($\gamma = 0$ or $\hat{p} = 0$),

$$\varphi_x^T(\hat{x})\, p = O(\varepsilon) \qquad (99)$$

Therefore, for any subsequent step, an analogous relation holds. Thus, at any point $x$, we conclude that

$$\varphi_x^T(x)\, \hat{p} = O(\varepsilon) \qquad (100)$$

Since $R_1$ given by (90) is proportional to the left-hand side of (100), we see that

$$R_1 = O(\varepsilon) \qquad (101)$$

Therefore, the Lagrange multiplier $\lambda$ computed with (82-1) is precise to $O(\varepsilon)$.

**10.2. Remark.** Because of (89), (91), (93), (100), we conclude that

$$\delta y(x) = O(\varepsilon^2), \qquad \nu(y) = O(\varepsilon^2) \qquad (102)$$

Furthermore, from (82-5), (89), (94), (100), we see that

$$\Delta x = O(\varepsilon), \qquad \Delta y = O(\varepsilon^2) \qquad (103)$$

**10.3. Descent Property of the Gradient Phase.** The first-order change of the function $F(x, \lambda)$ between points $x$ and $y$ is given by Eq. (15); in the light of (82), this can be written as

$$\delta F(x, \lambda) = -\alpha F_x^T(x, \lambda) F_x(x, \lambda) - R_1 \qquad (104)$$

where $R_1$ is a scalar given by

$$R_1 = \alpha\gamma F_x^T(v, \lambda)\, \hat{p} \qquad (105)$$

On account of (9), Eq. (105) can be rewritten as

$$R_1 = R_3 + R_4 \qquad (106)$$

where

$$R_3 = \alpha\gamma[f_x(x) + \varphi_x(\lambda)\hat{\lambda}]^T\hat{p}, \qquad R_4 = \alpha\gamma(\lambda - \hat{\lambda})^T\varphi_x^T(x)\, \hat{p} \qquad (107)$$

We observe that, to first order,

$$f_x(x) = f_x(\hat{y}) + f_{xx}(\hat{y})\,\Delta\hat{y}, \qquad \varphi_x(x) = \varphi_x(\hat{y}) + \varphi_{xx}(\hat{y})\,\Delta\hat{y} \qquad (108)$$

and that

$$\lambda - \hat{\lambda} = \Phi_x^T(\hat{x})(\Delta\hat{x} + \Delta\hat{y}) \qquad (109)$$

where $\Phi(x)$ denotes the right-hand side of (82-1). In the light of Eq. (87) applied to the previous iteration, Eqs. (107) become

$$R_3 = \alpha\gamma\,\Delta\hat{y}^T F_{xx}(\hat{y}, \lambda)\,\hat{p}, \qquad R_4 = \alpha\gamma(\Delta\hat{x} + \Delta\hat{y})^T \Phi_x(\hat{x})\,\varphi_x^T(x)\,\hat{p} \qquad (110)$$

If (89), (100), (103) are recalled, we see that

$$R_3 = O(\varepsilon^3), \qquad R_4 = O(\varepsilon^3) \qquad (111)$$

so that

$$R_1 = O(\varepsilon^3) \qquad (112)$$

Since the first term on the right-hand side of Eq. (104) is of $O(\varepsilon)$, $R_1$ can be neglected and Eq. (104) can be approximated by

$$\delta F(x, \lambda) = -\alpha F_x^T(x, \lambda) F_x(x, \lambda) \qquad (113)$$

This relationships shows that, if $\alpha > 0$, $\delta F(x, \lambda) < 0$. This is the descent property of the fundamental function during the gradient phase.

Because of definition (8), the relationship (17) can be established. Since $\delta F(x, \lambda)$ is of $O(\varepsilon)$ and $\delta y(x)$ is of $O(\varepsilon^2)$, Eq. (17) can be approximated by

$$\delta f(v) = \delta F(x, \lambda) \qquad (114)$$

which states that the functions $f(x)$ and $F(x, \lambda)$ behave identically, to first order. This and (113) establish the descent property of the function $f(x)$ during the gradient phase.

**10.4. Descent Property of the Algorithm.** Finally, we consider points $x$ and $\hat{x}$, both satisfying the constraint (2). To first order, the difference of the values of the function $f(x)$ at these points is given by

$$f(\hat{x}) - f(x) = f_x^T(x)(\Delta x + \Delta y) \qquad (115)$$

On account of (103), the second term on the right-hand side of (115) can be neglected with respect to the first. Hence, Eq. (115) becomes

$$f(\hat{x}) - f(x) = -\alpha F_x^T(x, \lambda) F_x(x, \lambda) \qquad (116)$$

Therefore, for $\epsilon$ sufficiently small, the restoration algorithm preserves the descent property of the gradient algorithm: the function $f$ decreases between any two successive restoration phases.

10.5. Coefficient $\gamma$. In the present algorithm, expression (82-3) has been used for the coefficient $\gamma$. The main justification for (82-3) is that it produces quadratic convergence in the terminal stage of the algorithm. The additional justification is that the terms (90) and (105) containing $\gamma$ are of $O(\epsilon)$ and $O(\epsilon^2)$, respectively. As a consequence, they do not seriously affect the computation of the multiplier $\lambda$ and the descent property of the algorithm.

## 11. Summary of the Algorithm

The algorithm presented in Part 2 consists of the alternate succession of gradient phases and restoration phases. A summary of the algorithm is given below.

11.1. Gradient Phase. For this phase, the algorithm is represented by

$$\lambda = -[\varphi_x^T(x)\,\varphi_x(x)]^{-1}\varphi_x^T(x)f_x(x)$$

$$F_x(x, \lambda) = f_x(x) + \varphi_x(x)\lambda$$

$$\gamma = F_x^T(x, \lambda)F_x(x, \lambda)/F_x^T(\hat{x}, \lambda)F_x(\hat{x}, \lambda) \tag{117}$$

$$p = F_x(x, \lambda) + \gamma\hat{p}$$

$$\Delta x = -\alpha p$$

$$y = x + \Delta x$$

The sequence of operations is as follows: (a) select a nominal point $x$ such that $\varphi(x) = 0$; (b) at this point, determine the vector $f_x(x)$ and the matrix $\varphi_x(x)$; (c) compute the multiplier $\lambda$ with (117-1), the vector $F_x(x, \lambda)$ with (117-2), the coefficient $\gamma$ with (117-3), and the search direction $p$ with (117-4), where $\hat{p}$ is known from the previous iteration; (d) determine the optimal stepsize $\alpha$ by a one-dimensional search in which either $\Psi(x) = f(y)$ or $\Psi(x) = F(y, \lambda)$ is minimized along the search direction $p$; the search is terminated when Ineq. (24) is satisfied; (e) compute the displacement $\Delta x$ with (117-5) and the varied point $y$ with (117-6).

11.2. Restoration Phase. For this phase, the algorithm is represented by

$$\sigma = k[\varphi_x^T(y)\,\varphi_x(y)]^{-1}\varphi(y)$$

$$\Delta y = -\varphi_x(y)\sigma \tag{118}$$

$$\hat{x} = y + \Delta y$$

The sequence of operations is as follows: (a) at point $y$, compute the vector $\varphi(y)$ and the matrix $\varphi_x(y)$; (b) assuming $k = 1$, determine the multiplier $\sigma$ with (118-1), the displacement $\Delta y$ with (118-2), and the varied point $\hat{x}$ with (118-3); (c) if $P(\hat{x}) < P(y)$, the scaling factor $k = 1$ is acceptable; if $P(\hat{x}) > P(y)$, the previous value of $k$ must be replaced by some smaller value in the range (29) until the condition $P(\hat{x}) < P(y)$ is met; this can be achieved through successive bisections of $k$; (d) return to step (a) and repeat the restoration algorithm using $\hat{x}$ as the starting point $y$ for the subsequent iteration; (e) terminate the restoration algorithm when the stopping condition (35) is satisfied; (f) once the restoration algorithm is completed, verify the inequality

$$f(\hat{x}) < f(x) \tag{119}$$

If Ineq. (119) is satisfied, start the next gradient phase. If Ineq. (119) is violated, return to the previous gradient phase and reduce the stepsize $\alpha$ until, after restoration, Ineq. (119) is satisfied.

11.3. Starting Condition. At the start of the algorithm, no information pertaining to the previous iteration is available; hence, we set $\hat{p} = 0$ or $\gamma = 0$. This means that the first step is a pure gradient step.

11.4. Restarting Condition. The algorithm must be restarted (a) when the optimal stepsize $\alpha$ cannot be employed due to violation of Ineq. (119) and (b) at the end of every $\Delta N = n - q$ or $\Delta N = n - q + 1$ iterations. The restarting is performed by setting $\hat{p} = 0$ or $\gamma = 0$.

11.5. Stopping Condition. The algorithm is terminated when Ineq. (51) is satisfied.

## PART 3: NUMERICAL EXAMPLES

In order to illustrate the theory, several numerical examples are developed using a Burroughs B-5500 computer and double-precision arithmetic. Concerning the *gradient phase*, the one-dimensional search is performed so as

to minimize either $\Psi(\alpha) = f(y)$ or $\Psi(\alpha) = F(y, \lambda)$ with respect to $\alpha$ along the search direction. Quasilinearization is used; the following stopping condition is employed:

$$\Psi_\alpha'(\alpha)/\Psi_\alpha'(0) \leqslant 10^{-4} \qquad (120)$$

corresponding to $\epsilon_1 = 10^{-4}$ in (25). Concerning the *restoration phase*, the restoration algorithm is employed iteratively until the error in the constraint satisfies the inequality

$$P(\tilde{x}) \leqslant 10^{-20} \qquad (121)$$

corresponding to $\theta_2 = 10^{-20}$ in (35). Finally, the *sequential gradient-restoration algorithm* is terminated when the error in the optimal conditions satisfies the inequality

$$Q(\tilde{x}, \lambda) \leqslant 10^{-14} \qquad (122)$$

corresponding to $\theta_4 = 10^{-14}$ in (51).

In the following sections, two groups of numerical examples are presented: (a) examples pertaining to quadratic function and linear constraint and (b) examples pertaining to nonquadratic function and/or nonlinear constraint. Both the ordinary gradient and the conjugate gradient versions of the algorithm are used. The following terminology is adopted: $N$ is the iteration number (each iteration includes a gradient phase and a restoration phase), $N_x$ is the number of restoration cycles per iteration, $N_*$ is the number of iterations for convergence, and $\Delta N$ is the number of iterations between successive restarting points. For simplicity, the symbols employed throughout Part 3 are scalar.

## 12. Examples: Quadratic Function, Linear Constraint

The first group of examples deals with functions of the form (69-1) subject to a constraint of the form (69-2). For these functions, the following properties hold: (a) no restoration is needed, that is, $N_x = 0$; (b) the search performed by quasilinearization yields the optimal value of $\alpha$ in one step; (c) the value of $\alpha$ computed by minimizing $\Psi = f$ is the same as that computed by minimizing $\Psi = F$; and (d) convergence to the minimum occurs in no more than $n - q$ iterations in the conjugate gradient version of the algorithm.

Example 12.1. We consider the problem of minimizing the function

$$f = (x + y)^2 + (y + z)^2 \qquad (123)$$

subject to the constraint

$$x + 2y + 3z - 1 = 0 \qquad (124)$$

Note that $n = 3$, $q = 1$, so that $n - q = 2$. This function admits the relative minimum $f = 0$ at the point defined by

$$x = \tfrac{1}{2}, \quad y = -\tfrac{1}{2}, \quad z = \tfrac{1}{2} \qquad (125)$$

The nominal point chosen for starting the algorithm is the point of coordinates

$$x = -4, \quad y = 1, \quad z = 1 \qquad (126)$$

consistent with (124). Convergence is achieved in $N_* = 19$ iterations with the ordinary gradient version of the algorithm and $N_* = 2$ iterations with the conjugate gradient version. For the latter, Table 1 shows the detailed results.

Table 1 $(n - q = 2)$

| $N$ | $N_x$ | $x$ | $y$ | $z$ | $f$ | $Q$ |
|---|---|---|---|---|---|---|
| 0 | — | −4.0000 | 1.0000 | 1.0000 | $0.13 \times 10^2$ | $0.33 \times 10^2$ |
| 1 | 0 | −1.6459 | 1.8759 | −0.3686 | $0.23 \times 10^1$ | $0.20 \times 10^1$ |
| 2 | 0 | 0.5000 | −0.5000 | 0.5000 | $0.93 \times 10^{-16}$ | $0.50 \times 10^{-16}$ |

Example 12.2. We consider the problem of minimizing the function

$$f = (x - 1)^2 + (y - z)^2 + (u - w)^2 \qquad (127)$$

subject to the constraints

$$x + y + z + u + w - 5 = 0, \qquad x - 2(u + w) + 3 = 0 \qquad (128)$$

Note that $n = 5$, $q = 2$, so that $n - q = 3$. This function admits the relative minimum $f = 0$ at the point defined by

$$x = 1, \quad y = 1, \quad z = 1, \quad u = 1, \quad w = 1 \qquad (129)$$

The nominal point chosen for starting the algorithm is the point of coordinates

$$x = 3, \quad y = 5, \quad z = -3, \quad u = 2, \quad w = -2 \qquad (130)$$

consistent with (128). Convergence is achieved in $N_* = 17$ iterations with the ordinary gradient version of the algorithm and $N_* = 3$ iterations with the conjugate gradient version. For the latter, Table 2 shows the detailed results.

Table 2 ($n - q = 3$)

| N | $N_*$ | $x$ | $y$ | $z$ | $u$ | $w$ | $f$ | $Q$ |
|---|---|---|---|---|---|---|---|---|
| 0 | — | 3.0000 | 5.0000 | −5.0000 | 2.0000 | −2.0000 | $0.84 \times 10^5$ | $0.62 \times 10^4$ |
| 1 | 0 | 1.6477 | 0.6651 | 0.6565 | 0.7909 | 1.0373 | $0.71 \times 10^4$ | $0.26 \times 10^4$ |
| 2 | 0 | 1.1904 | 0.8169 | 0.9417 | 1.1059 | 0.8549 | $0.10 \times 10^6$ | $0.63 \times 10^4$ |
| 3 | 0 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | $0.12 \times 10^{-11}$ | $0.69 \times 10^{-10}$ |

Example 12.3. We consider the problem of minimizing the function

$$f = (x - y)^2 + (y + z - 2)^2 + (u - 1)^2 + (w - 1)^2 \qquad (131)$$

subject to the constraints

$$x + 3y - 4 = 0, \qquad z + u - 2w = 0, \qquad y - u = 0 \qquad (132)$$

Note that $n = 5$, $q = 3$, so that $n - q = 2$. This function admits the relative minimum $f = 0$ at the point defined by

$$x = 1, \quad y = 1, \quad z = 1, \quad u = 1, \quad w = 1 \qquad (133)$$

The nominal point chosen for starting the algorithm is the point of coordinates

$$x = 1, \quad y = 1, \quad z = 2, \quad u = -1, \quad w = 1 \qquad (134)$$

consistent with (132). Convergence is achieved in $N_* = 13$ iterations with the ordinary gradient version of the algorithm and $N_* = 2$ iterations with the conjugate gradient version. For the latter, Table 3 shows the detailed results.

Table 3 ($n - q = 2$)

| N | $N_*$ | $x$ | $y$ | $z$ | $u$ | $w$ | $f$ | $Q$ |
|---|---|---|---|---|---|---|---|---|
| 0 | — | 2.5000 | 0.5000 | 2.6565 | −1.0000 | 0.5000 | $0.85 \times 10^5$ | $0.40 \times 10^4$ |
| 1 | 0 | 0.5055 | 2.0638 | 1.5795 | 0.5244 | 1.0635 | $0.73 \times 10^4$ | $0.37 \times 10^4$ |
| 2 | 0 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | $0.20 \times 10^{-11}$ | $0.30 \times 10^{-10}$ |

## 13. Examples: Nonquadratic Function and/or Nonlinear Constraint

The second group of examples deals with general functions and general constraints. Concerning the search, the value of $\alpha$ computed by minimizing $\Psi = f$ is different from that computed by minimizing $\Psi = F$. For the conjugate gradient version, the algorithm is restarted every $\Delta N$ iterations. Note that $\Delta N = 1$ corresponds to the ordinary gradient version.

Example 13.1. We consider the problem of minimizing the function

$$f = (x - y)^2 + (y - z)^4 \qquad (135)$$

subject to the constraint

$$x(1 + y^2) + z^4 - 3 = 0 \qquad (136)$$

Note that $n = 3$, $q = 1$, so that $n - q = 2$. This function admits the relative minimum $f = 0$ at the point defined by

$$x = 1, \quad y = 1, \quad z = 1 \qquad (137)$$

The nominal point chosen for starting the algorithm is the point of coordinates

$$x = -13/5, \quad y = 2, \quad z = 2 \qquad (138)$$

consistent with (136). Table 4 shows $N_*$ for several values of $\Delta N$ and both $\Psi = f$ and $\Psi = F$. The detailed results pertaining to $\Delta N = 2$ and $\Psi = F$ are presented in Table 5.

Table 4 ($n - q = 2$)

| $\Delta N$ | $N_*$ ($\Psi = f$) | $N_*$ ($\Psi = F$) |
|---|---|---|
| 1 | 3821 | 4569 |
| 2 | 8 | 13 |
| 3 | 12 | 13 |
| 4 | 7 | 16 |

Table 5 $(n - q = 2, \Delta N = 2, \Psi = F)$

| N | $N_a$ | x | y | z | f | Q |
|---|---|---|---|---|---|---|
| 0 | — | −2.0000 | 2.0000 | 2.0000 | $0.21 \times 10^1$ | $0.15 \times 10^1$ |
| 1 | 5 | −0.3769 | 0.0112 | 1.3556 | $0.33 \times 10^1$ | $0.82 \times 10^1$ |
| 2 | 3 | 1.3359 | 0.9512 | 0.4357 | $0.43 \times 10^0$ | $0.21 \times 10^1$ |
| 3 | 3 | 1.2026 | 1.1752 | 0.6076 | $0.10 \times 10^1$ | $0.86 \times 10^1$ |
| 4 | 6 | 0.7711 | 0.9447 | 1.1140 | $0.30 \times 10^{-1}$ | $0.22 \times 10^0$ |
| 5 | 2 | 0.8515 | 0.8659 | 1.1078 | $0.35 \times 10^1$ | $0.34 \times 10^1$ |
| 6 | 4 | 1.0766 | 1.0316 | 0.9214 | $0.72 \times 10^1$ | $0.23 \times 10^1$ |
| 7 | 1 | 1.0058 | 1.0641 | 0.9234 | $0.39 \times 10^{-2}$ | $0.19 \times 10^1$ |
| 8 | 3 | 1.0211 | 1.0225 | 0.9769 | $0.62 \times 10^1$ | $0.17 \times 10^1$ |
| 9 | 1 | 1.0216 | 1.0217 | 0.9770 | $0.39 \times 10^1$ | $0.17 \times 10^1$ |
| 10 | 2 | 1.0055 | 1.0056 | 0.9943 | $0.26 \times 10^{-1}$ | $0.85 \times 10^{-1}$ |
| 11 | 1 | 1.0055 | 1.0055 | 0.9943 | $0.15 \times 10^1$ | $0.43 \times 10^{-14}$ |
| 12 | 2 | 1.0008 | 1.0008 | 0.9991 | $0.12 \times 10^{-15}$ | $0.24 \times 10^{-15}$ |
| 13 | 0 | 1.0008 | 1.0008 | 0.9991 | $0.97 \times 10^{-16}$ | $0.65 \times 10^{-16}$ |

Example 13.2.   We consider the problem of minimizing the function

$$f = (x - y)^2 + (z - 1)^4 + (u - 1)^4 + (w - 1)^4 \qquad (139)$$

subject to the constraints

$$ux^2 + \sin(u - w) - 1 = 0, \qquad y + z^4 u^4 - 2 = 0 \qquad (140)$$

Note that $n = 5$, $q = 2$, so that $n - q = 3$. This function admits the relative minimum $f = 0$ at the point defined by

$$x = 1, \quad y = 1, \quad z = 1, \quad u = 1, \quad w = 1 \qquad (141)$$

The nominal point chosen for starting the algorithm is the point of coordinates

$$x = 1/\sqrt{2}, \quad y = 7/4, \quad z = 1, \quad u = 2, \quad w = 2 \qquad (142)$$

Table 6 $(n - q = 3)$

| $\Delta N$ | $N_*$ $(\Psi = f)$ | $N_*$ $(\Psi = F)$ |
|---|---|---|
| 1 | 30 | 112 |
| 2 | 10 | 10 |
| 3 | 8 | 11 |
| 4 | 10 | 12 |
| 5 | 12 | 12 |

Table 7 $(n - q = 3, \Delta N = 3, \Psi = F)$

| N | $N_a$ | x | y | z | u | w | f | Q |
|---|---|---|---|---|---|---|---|---|
| 0 | — | 0.7071 | 1.7500 | 0.5000 | 2.0000 | 2.0000 | $0.33 \times 10^1$ | $0.55 \times 10^1$ |
| 1 | 5 | 0.9072 | 1.3220 | 1.0192 | 0.7926 | 0.4377 | $0.20 \times 10^0$ | $0.11 \times 10^1$ |
| 2 | 3 | 1.0700 | 1.0988 | 1.1376 | 0.3334 | 0.5725 | $0.30 \times 10^{-1}$ | $0.38 \times 10^1$ |
| 3 | 3 | 1.0302 | 1.0431 | 1.0338 | 0.9032 | 0.1619 | $0.26 \times 10^0$ | $0.44 \times 10^1$ |
| 4 | 2 | 1.0387 | 1.0485 | 1.0105 | 0.9008 | 0.3728 | $0.18 \times 10^0$ | $0.38 \times 10^1$ |
| 5 | 1 | 1.0073 | 1.0074 | 1.0047 | 0.9861 | 0.9862 | $0.22 \times 10^0$ | $0.15 \times 10^1$ |
| 6 | 2 | 1.0063 | 1.0082 | 1.0038 | 0.9343 | 0.9389 | $0.18 \times 10^0$ | $0.67 \times 10^1$ |
| 7 | 1 | 1.0066 | 1.0069 | 1.0140 | 0.9443 | 0.9902 | $0.16 \times 10^0$ | $0.73 \times 10^1$ |
| 8 | 2 | 1.0030 | 1.0010 | 1.0000 | 0.9983 | 1.0028 | $0.31 \times 10^0$ | $0.19 \times 10^1$ |
| 9 | 1 | 1.0010 | 1.0010 | 1.0000 | 0.9943 | 1.0028 | $0.26 \times 10^0$ | $0.74 \times 10^1$ |
| 10 | 1 | 1.0030 | 1.0030 | 1.0000 | 0.9943 | 1.0028 | $0.22 \times 10^0$ | $0.13 \times 10^1$ |
| 11 | 1 | 1.0029 | 1.0029 | 1.0000 | 0.9983 | 1.0029 | $0.50 \times 10^{-11}$ | $0.61 \times 10^{-14}$ |

consistent with (140). Table 6 shows $N_*$ for several values of $\Delta N$ and both $\Psi = f$ and $\Psi = F$. The detailed results pertaining to $\Delta N = 3$ and $\Psi = F$ are presented in Table 7.

Example 13.3.   We consider the problem of minimizing the function

$$f = (x - y)^2 + (y - z)^2 + (z - u)^4 + (u - w)^4 \qquad (143)$$

subject to the constraints

$$x + y^2 + z^2 - 3 = 0, \qquad y - z^2 + u - 1 = 0, \qquad uw - 1 = 0 \qquad (144)$$

Note that $n = 5$, $q = 3$, so that $n - q = 2$. This function admits the relative minimum $f = 0$ at the point defined by

$$x = 1, \quad y = 1, \quad z = 1, \quad u = 1, \quad w = 1 \qquad (145)$$

The nominal point chosen for starting the algorithm is the point of coordinates

$$x = 2, \quad y = \sqrt{2}, \quad z = -1, \quad u = 2 - \sqrt{2}, \quad w = 1 \qquad (146)$$

consistent with (144). Table 8 shows $N_*$ for several values of $\Delta N$ and both $\Psi = f$ and $\Psi = F$. The detailed results pertaining to $\Delta N = 2$ and $\Psi = F$ are presented in Table 9.

Table 8 (n − q = 2)

| ΔN | $N_*$ (Ψ = f) | $N_*$ (Ψ = F) |
|---|---|---|
| 1 | 29 | 11 |
| 2 | 10 | 11 |
| 3 | 11 | 14 |
| 4 | 10 | 12 |

Table 9 (n − q = 2, ΔN = 2, Ψ = F)

| N | $N_r$ | x | y | z | u | w | f | Q |
|---|---|---|---|---|---|---|---|---|
| 0 | — | 2.0000 | 1.4142 | −1.0000 | 0.5357 | 0.5000 | $0.11 \times 10^0$ | $0.35 \times 10^0$ |
| 1 | 4 | 1.8040 | 1.0517 | −0.0536 | −0.0488 | 0.5279 | $0.20 \times 10^0$ | $0.48 \times 10^0$ |
| 2 | 4 | 1.0633 | 1.1390 | 0.3402 | −0.0233 | 0.6012 | $0.10 \times 10^0$ | $0.37 \times 10^0$ |
| 3 | 6 | 1.0921 | 0.7644 | 1.0979 | 1.4409 | 0.9156 | $0.30 \times 10^0$ | $0.27 \times 10^0$ |
| 4 | 4 | 0.9039 | 0.8443 | 1.1111 | 1.3877 | 1.1037 | $0.83 \times 10^{-1}$ | $0.30 \times 10^0$ |
| 5 | 3 | 0.8563 | 0.9538 | 1.0719 | 1.1542 | 1.1677 | $0.21 \times 10^{-1}$ | $0.22 \times 10^{-1}$ |
| 6 | 3 | 1.0133 | 1.0147 | 0.9758 | 0.9376 | 0.9606 | $0.21 \times 10^{-1}$ | $0.26 \times 10^{-1}$ |
| 7 | 2 | 1.0017 | 1.0099 | 0.9927 | 0.9755 | 0.9982 | $0.16 \times 10^{-1}$ | $0.16 \times 10^{-1}$ |
| 8 | 2 | 1.0017 | 0.9979 | 0.9994 | 0.9969 | 0.9982 | $0.12 \times 10^{-1}$ | $0.89 \times 10^{-1}$ |
| 9 | 1 | 1.0003 | 1.0004 | 0.9995 | 0.9986 | 0.9994 | $0.10 \times 10^{-4}$ | $0.95 \times 10^{-1}$ |
| 10 | 1 | 0.9993 | 1.0000 | 0.9999 | 0.7993 | 1.0000 | $0.70 \times 10^{-11}$ | $0.29 \times 10^{-14}$ |
| 11 | 0 | 0.9999 | 0.9999 | 1.0000 | 1.0000 | 1.0000 | $0.20 \times 10^{-14}$ | $0.21 \times 10^{-14}$ |

## 14. Discussion and Conclusions

In the previous sections, a sequential algorithm is developed for minimizing a function $f(x)$ subject to the constraint $\varphi(x) = 0$. The algorithm is composed of the alternate succession of gradient phases and restoration phases. For the gradient phase, two versions are presented, one in which information at only point $x$ is used (ordinary gradient version) and one in which information at both points $x$ and $\tilde{x}$ is used (conjugate gradient version). For the restoration phase, the criterion employed is that of the least-square change of the coordinates.

While the ordinary gradient version of the algorithm exhibits only asymptotic convergence, the conjugate gradient version exhibits quadratic convergence; this means that, for a quadratic function subject to a linear constraint, the minimum point is obtained in no more than $n − q$ iterations.

In order to illustrate the theory, several numerical examples are presented. The first three examples are concerned with a quadratic function subject to a linear constraint: the computer results confirm the quadratic convergence properties of the conjugate gradient version of the algorithm. The next three

examples are concerned with general functions subject to general constraints: the computer results show the rapidly convergent properties of the conjugate gradient version of the algorithm and its superiority with respect to the ordinary gradient version.

In the theory as well as in the examples, the function $\Psi(\alpha) = f(y)$ or the function $\Psi(\alpha) = F(y, \lambda)$ is minimized along the search direction. However, the occurrence of cases where $\Psi(\alpha)$ does not possess a relative minimum is conceivable. For these cases, an upper limit must be imposed on the stepsize $\alpha$ or the performance index $P(y)$.

In the numerical examples, the starting point was chosen so that $\varphi(x) = 0$. However, the present algorithms can be started even if $\varphi(x) \neq 0$. In this case, the first phase is a restoration phase rather than a gradient phase.

## References

1. MIELE, A., and HEIDEMAN, J. C., *Mathematical Programming for Constrained Minimal Problems, Part 1, Sequential Gradient-Restoration Algorithm*, Rice University, Aero-Astronautics Report No. 59, 1969.

2. MIELE, A., HUANG, H. Y., and HEIDEMAN, J. C., *Mathematical Programming for Constrained Minimal Problems, Part 2, Sequential Conjugate Gradient-Restoration Algorithm*, Rice University, Aero-Astronautics Report No. 61, 1969.

3. MIELE, A., HUANG, H. Y., and CANTRELL, J. W., *Gradient Methods in Mathematical Programming, Part 1, Review of Previous Techniques*, Rice University, Aero-Astronautics Report No. 55, 1969.

4. MIELE, A., HEIDEMAN, J. C., and DAMOULAKIS, J. N., *The Restoration of Constraints in Holonomic and Nonholonomic Problems*, Journal of Optimization Theory and Applications, Vol. 3, No. 5, 1969.

5. HESTENES, M. R., and STIEFEL, E., *Methods of Conjugate Gradients for Solving Linear Systems*, Journal of Research of the National Bureau of Standards, Vol. 49, No. 6, 1952.

6. BECKMAN, F. S., *The Solution of Linear Equations by the Conjugate Gradient Method*, Mathematical Methods for Digital Computers, Edited by A. Ralston and H. S. Wilf, John Wiley and Sons, New York, 1960.

7. FLETCHER, R., and REEVES, C. M., *Function Minimization by Conjugate Gradients*, Computer Journal, Vol. 7, No. 2, 1964.

## Additional Bibliography

ROSEN, J. B., *The Gradient Projection Method for Nonlinear Programming, Part 2, Nonlinear Constraints*, SIAM Journal on Applied Mathematics, Vol. 9, No. 4, 1961.
BRYSON, A. E., JR., and HO, Y. C., *Applied Optimal Control*, Chapter 1, Blaisdell Publishing Company, New York, 1969.

# On the Method of Multipliers
## for Mathematical Programming Problems[1]

A. MIELE,[2] P. E. MOSELEY,[3] A. V. LEVY,[4] AND G. M. COGGINS[5]

**Abstract.** In this paper, the numerical solution of the basic problem of mathematical programming is considered. This is the problem of minimizing a function $f(x)$ subject to a constraint $\varphi(x) = 0$. Here, $f$ is a scalar, $x$ is an $n$-vector, and $\varphi$ is a $q$-vector, with $q < n$.

The approach employed is based on the introduction of the augmented penalty function $W(x, \lambda, k) = f(x) + \lambda^T \varphi(x) + k \varphi^T(x) \varphi(x)$. Here, the $q$-vector $\lambda$ is an approximation to the Lagrange multiplier, and the scalar $k > 0$ is the penalty constant.

Previously, the augmented penalty function $W(x, \lambda, k)$ was used by Hestenes in his method of multipliers. In Hestenes' version, the method of multipliers involves cycles, in each of which the multiplier and the penalty constant are held constant. After the minimum of the augmented penalty function is achieved in any given cycle, the multiplier $\lambda$ is updated, while the penalty constant $k$ is held unchanged.

In this paper, two modifications of the method of multipliers are presented in order to improve its convergence characteristics. The improved convergence is achieved by (i) increasing the updating frequency so that the number of iterations in a cycle is shortened to $\Delta N = 1$ for the ordinary-gradient algorithm and the modified quasilinearization algorithm and $\Delta N = n$ for the conjugate-gradient algorithm, (ii) imbedding Hestenes' updating rule for the multiplier $\lambda$

into a one-parameter family and determining the scalar parameter $\beta$ so that the error in the optimum condition is minimized, and (iii) updating the penalty constant $k$ so as to cause some desirable effect in the ordinary-gradient algorithm, the conjugate-gradient algorithm, and the modified-quasilinearization algorithm. For the sake of identification, Hestenes' method of multipliers is called Method MM-1, the modification including (i) and (ii) is called Method MM-2, and the modification including (i), (ii), (iii) is called Method MM-3.

Evaluation of the theory is accomplished with seven numerical examples. The first example pertains to a quadratic function subject to linear constraints. The remaining examples pertain to non-quadratic functions subject to nonlinear constraints. Each example is solved with the ordinary-gradient algorithm, the conjugate-gradient algorithm, and the modified-quasilinearization algorithm, which are employed in conjunction with Methods MM-1, MM-2, and MM-3.

The numerical results show that (a) for given penalty constant $k$, Method MM-2 generally exhibits faster convergence than Method MM-1, (b) in both Methods MM-1 and MM-2, the number of iterations for convergence has a minimum with respect to $k$, and (c) the number of iterations for convergence of Method MM-3 is close to the minimum with respect to $k$ of the number of iterations for convergence of Method MM-2. In this light, Method MM-3 has very desirable characteristics.

# 1. Introduction

Over the past several years, considerable work has been done on the numerical solution of the constrained minimization problem. This is the problem of minimizing a function $f(x)$ subject to a constraint $\varphi(x) = 0$. Here, $f$ is a scalar, $x$ is an $n$-vector, and $\varphi$ is a $q$-vector, with $q < n$.

The methods employed are generally based on one of two basic ideas. One approach ensures constraint satisfaction, at least to first order, at the end of each iteration (see, for example, Refs. 1–3). The other approach depends on the construction of a sequence of special functions having, in the limit, an unconstrained minimum point coincident with the solution of the original constrained minimization problem. With regard to the latter approach, the standard penalty function method (see, for example, Refs. 4–6) and Hestenes' method of multipliers (Ref. 7) must be mentioned.

As the numerical experiments of Refs. 8–9 indicate, the method of multipliers generally exhibits faster convergence than the standard penalty function method. Crucial to the method of multipliers is the manner in which the multiplier $\lambda$ is estimated and the penalty constant $k$ is updated at the beginning of each cycle. These key questions are considered in this paper, whose objective is to enlarge the investigation of Ref. 10 and to develop techniques for improving the convergence characteristics of the method of multipliers. The resulting algorithm is called *modified method of multipliers*.

## 2. Statement of the Problem

We consider the problem of minimizing the function

$$f = f(x) \tag{1}$$

subject to the constraint

$$\varphi(x) = 0. \tag{2}$$

In the above equations, $f$ is a scalar, $x$ is an $n$-vector, and $\varphi$ is a $q$-vector, with $q < n$. Here, all vectors are column vectors. It is assumed that the first and second partial derivatives of the functions $f$ and $\varphi$ exist and are continuous and that the constrained minimum exists.

### 2.1. First-Order Conditions.
The previous problem can be recast as that of minimizing the augmented function

$$F(x, \lambda) = f(x) + \lambda^T \varphi(x) \tag{3}$$

subject to the constraint (2). The $q$-vector $\lambda$ is the Lagrange multiplier, and the superscript $T$ denotes the transpose of a matrix.

From theory of maxima and minima, it is known that the optimal solution must satisfy the relations

$$\varphi(x) = 0, \qquad F_x(x, \lambda) = f_x(x) + \varphi_x(x)\lambda = 0. \tag{4}$$

which are a system of $n + q$ equations in $x$ and $\lambda$. The subscript $x$ denotes the gradient of a function: in this case, $f_x$ and $F_x$ are $n$-vectors, and $\varphi_x$ is an $n \times q$ matrix, defined in such a way that its $i$th column is the gradient of the $i$th scalar component of $\varphi$ with respect to $x$.

### 2.2. Approximate Solutions.
In general, the system (4) is

nonlinear; consequently, approximate methods must be employed. Here, we introduce the scalar performance indexes

$$P(x) = \varphi^T(x)\, \varphi(x), \qquad Q(x, \lambda) = F_x^T(x, \lambda)\, F_x(x, \lambda), \qquad (5)$$

which measure the errors in the constraint and the optimum condition, respectively. At the solution point, $P \to 0$ and $Q = 0$, while $P > 0$ and/or $Q > 0$ for any approximation of the solution. When approximate methods are employed, they must ultimately lead to values of $x$, $\lambda$ such that

$$P(x) \leq \epsilon_1, \qquad Q(x, \lambda) \leq \epsilon_2. \qquad (6)$$

Alternatively, (6) can be replaced by

$$R(x, \lambda) \leq \epsilon_3, \qquad (7)$$

where

$$R(x, \lambda) = P(x) + Q(x, \lambda) \qquad (8)$$

denotes the cumulative error in the constraint and the optimum condition. Here, $\epsilon_1$, $\epsilon_2$, $\epsilon_3$ are small, preselected numbers. Note that satisfaction of Ineq. (7) implies satisfaction of Ineqs. (6) if one chooses $\epsilon_1 = \epsilon_2 = \epsilon_3$.

## 3. Review of Penalty Function Methods

The penalty function method is based on the construction of a sequence of special functions having, in the limit, an unconstrained minimum point coincident with the solution of the original constrained minimization problem. In this section, two versions of the penalty function method are reviewed: (i) the standard penalty function method and (ii) the method of multipliers. Method (i) is based on the standard penalty function (9) and Method (ii) is based on the augmented penalty function (22).

3.1. Standard Penalty Function Method. This method is based on the consideration of the standard penalty function

$$U(x, k) = f(x) + k\varphi^T(x)\varphi(x). \qquad (9)$$

This is obtained by adding to the function $f(x)$ a term quadratic in the constraint $\varphi(x)$, $k > 0$ being the penalty constant.

The problem of minimizing the function (1) subject to the con-

straint (2) is replaced by a sequence of unconstrained minimization problems. In each element of the sequence or cycle, one minimizes the function (9) with respect to $x$ for given $k$. Therefore, theoretically speaking, the following necessary condition must be satisfied at the end of each cycle:

$$U_x(x, k) = f_x(x) + 2k\varphi_x(x)\varphi(x) = 0. \qquad (10)$$

If the penalty constant $k$ is arbitrary, the vector $x$ which satisfies Eq. (10) is such that $\varphi(x) \neq 0$. However, if one defines the Lagrange multiplier to be

$$\lambda = 2k\varphi(x), \qquad (11)$$

Eq. (10) reduces to

$$F_x(x, \lambda) = f_x(x) + \varphi_x(x)\lambda = 0, \qquad (12)$$

meaning that the combination of $x$ and $\lambda$ thus obtained satisfies exactly the optimum condition.

In order to obtain constraint satisfaction, increasingly larger values of the penalty constant must be employed in successive cycles of the standard penalty function method. In this connection, let $k_1$ denote the penalty constant of the present cycle and $k_2$ denote the penalty constant of the next cycle, with $k_2 > k_1$. Because of the jump in $k$, the standard penalty function increases by the amount

$$U(x, k_2) - U(x, k_1) = (k_2 - k_1)\, P(x), \qquad (13)$$

and the norm of the gradient of the standard penalty function takes the value[a]

$$U_x^T(x, k_2)\, U_x(x, k_2) = (k_2 - k_1)^2\, P_x^T(x)\, P_x(x), \qquad (14)$$

where

$$P(x) = \varphi^T(x)\varphi(x), \qquad P_x(x) = 2\varphi_x(x)\varphi(x). \qquad (15)$$

The positiveness of the right-hand side of Eq. (13) is the key to the mechanism on which the standard penalty function method is based.

After a sufficient number of cycles, the constraint error can be made as small as desired providing the penalty constant has become sufficiently large. Theoretically speaking, the condition $\varphi(x) = 0$ is desired at convergence; consequently, the multiplier $\lambda$ defined by Eq. (11) can be identical with the multiplier satisfying Eqs. (4), which is

---

[a] Note that $U_x(x, k_1) = 0$.

generally nonzero, only if $k \to \infty$. In a practical digital computer, this means that very large values of $k$ are needed at convergence.

*Numerical Implementation.* From the above considerations, the following outline of the standard penalty function method emerges.

(a) The original constrained minimization problem is replaced by a sequence of unconstrained minimization problems.

(b) In each element of the sequence or cycle, the standard penalty function

$$U(x, k) = f(x) + k\varphi^T(x)\,\varphi(x) \tag{16}$$

is minimized with respect to $x$ for given $k$. The minimum of $U(x, k)$ is achieved when the following stopping condition is satisfied:

$$U_x^T(x, k)\,U_x(x, k) \leqslant \epsilon_1. \tag{17}$$

where $\epsilon_1$ is a small, preselected number.

(c) The solution point of any given cycle is chosen as the starting point of the next cycle of the standard penalty function method.

(d) For the next cycle, a higher value of the penalty constant is selected, one choice being

$$k_2 = \pi k_1. \tag{18}$$

where $\pi > 1$ is the penalty constant ratio.

(e) After updating the penalty constant, one returns to (b) and continues iteratively.

(f) The algorithm is terminated when the following stopping condition is satisfied:

$$\varphi^T(x)\varphi(x) + U_x^T(x, k)\,U_x(x, k) \leqslant \epsilon_2. \tag{19}$$

where $\epsilon_2$ is a small, preselected number.

*Remark.* At convergence of a cycle, the stopping condition (17) can be written as

$$Q(x, \lambda) \leqslant \epsilon_1. \tag{20}$$

where $\lambda$ is given by Eq. (11). Analogously, at convergence of the algorithm, the stopping condition (19) becomes

$$P(x, \lambda) = P(x) + Q(x, \lambda) \leqslant \epsilon_2. \tag{21}$$

where $\lambda$ is given by Eq. (11).

**3.2. Method of Multipliers.** This method is based on the consideration of the augmented penalty function

$$W(x, \lambda, k) = f(x) + \lambda^T\varphi(x) + k\varphi^2(x)\,\varphi(x). \tag{22}$$

This is obtained by adding to the penalty function $U(x, k)$ a term linear in the constraint $\varphi(x)$, the $q$-vector $\lambda$ being an approximation to the Lagrange multiplier. The use of this function was suggested by Hestenes (Ref. 7) in order to circumvent the numerical difficulties associated with the extremely large values of the penalty constant required by the standard penalty function method.

The problem of minimizing the function (1) subject to the constraint (2) is replaced by a sequence of unconstrained minimization problems. In each element of the sequence or cycle, one minimizes the function (22) with respect to $x$ for given $\lambda$ and $k$. Therefore, theoretically speaking, the following necessary condition must be satisfied at the end of each cycle:

$$W_x(x, \lambda, k) = f_x(x) + \varphi_x(x)\lambda + 2k\varphi_x(x)\,\varphi(x) = 0 \tag{23}$$

and is equivalent to

$$W_x(x, \lambda, k) = f_x(x) + \varphi_x(x)[\lambda + 2k\varphi(x)] = 0. \tag{24}$$

If the penalty constant $k$ and the multiplier $\lambda$ are arbitrary, the vector $x$ which satisfies Eq. (24) is such that $\varphi(x) \neq 0$. However, by means of a proper updating rule, a new Lagrange multiplier can be found such that the optimum condition is satisfied exactly. In this connection, let $\lambda_1$ denote the Lagrange multiplier of the present cycle and $\lambda_2$ denote the Lagrange multiplier of the next cycle. If $\lambda_2$ is chosen to be

$$\lambda_2 = \lambda_1 + 2k\varphi(x). \tag{25}$$

Eq. (24) reduces to

$$F_x(x, \lambda_2) = f_x(x) + \varphi_x(x)\lambda_2 = 0. \tag{26}$$

meaning that the combination of $x$ and $\lambda_2$ thus obtained satisfies exactly the optimum condition.

At the end of any given cycle, whenever the value of the Lagrange multiplier is changed from $\lambda_1$ to $\lambda_2$, the augmented penalty function increases by the amount

$$W(x, \lambda_2, k) - W(x, \lambda_1, k) = 2kP(x). \tag{27}$$

and the norm of the gradient of the augmented penalty function takes the value[7]

$$W_x^T(x, \lambda_d, k) \, W_x(x, \lambda_d, k) = k^2 P_x^T(x) \, P_x(x), \tag{28}$$

where $P(x)$ and $P_x(x)$ are given by Eqs. (15). The positiveness of the right-hand side of Eq. (27) is the key to the mechanism on which the method of multipliers is based.

The attention of the reader is called on the essential similarity between Eqs. (13) and (27). In the standard penalty function method, the drive toward constraint satisfaction is supplied by increasing the penalty constant from cycle to cycle. In the method of multipliers, the drive toward constraint satisfaction is supplied by changing the multiplier from cycle to cycle in accordance with Eq. (25).

While the standard penalty function method requires extremely large values of the penalty constant at convergence, this is not the case with the method of multipliers. In the latter method, convergence can be achieved even with moderate values of the penalty constant.

*Numerical Implementation.* From the above considerations, the following outline of the method of multipliers (Method MM-1) emerges.

(a) The original constrained minimization problem is replaced by a sequence of unconstrained minimization problems.

(b) In each element of the sequence or cycle, the augmented penalty function

$$W(x, \lambda, k) = f(x) + \lambda^T \varphi(x) + k \varphi^T(x) \varphi(x) \tag{29}$$

is minimized with respect to $x$ for given $\lambda$ and $k$. The minimum of $W(x, \lambda, k)$ is achieved when the following stopping condition is satisfied:

$$W_x^T(x, \lambda, k) \, W_x(x, \lambda, k) \leq \epsilon_1, \tag{30}$$

where $\epsilon_1$ is a small, preselected number.

(c) The solution point of any given cycle is chosen as the starting point of the next cycle of the method of multipliers.

(d) For the next cycle, the multiplier is updated according to the simple rule

$$\lambda_d = \lambda_1 + 2k\varphi(x). \tag{31}$$

(e) After updating the multiplier, one returns to (b) and continues iteratively.

(f) The algorithm is terminated when the following stopping condition is satisfied:

$$\varphi^T(x) \varphi(x) + W_x^T(x, \lambda, k) \, W_x(x, \lambda, k) \leq \epsilon_2, \tag{32}$$

where $\epsilon_2$ is a small, preselected number.

(g) To start the algorithm some assumption concerning the multiplier is necessary. The simplest assumption is

$$\lambda = 0 \tag{33}$$

and is equivalent to stating that the augmented penalty function (29) and the standard penalty function (16) are identical for the first cycle of the algorithm.

*Remark.* At convergence of a cycle, the stopping condition (30) can be written as

$$Q(x, \lambda_d) \leq \epsilon_1, \tag{34}$$

where $\lambda_d$ is given by Eq. (31). Analogously, at convergence of the algorithm, the stopping condition (32) becomes

$$R(x, \lambda_d) = P(x) + Q(x, \lambda_d) \leq \epsilon_2, \tag{35}$$

where $\lambda_2$ is given by Eq. (31).

## 4. Modifications of the Method of Multipliers

The method of multipliers described in Section 3 has one drawback: a sequence of unconstrained minimization problems must be solved, each possibly requiring a large number of iterations $\Delta N$. Consequently, the total number of iterations for convergence $N_* = \Sigma(\Delta N)$ may become excessive for practical applications.

In order to accelerate convergence, we explore here several modifications of the method of multipliers. These modifications are obtained by (i) shortening the length of a cycle, (ii) improving the estimate of the multiplier, and (iii) selecting the penalty constant in an appropriate fashion.

**4.1. Updating Frequency.** Let a cycle be defined as any sequence of iterations in which the multiplier $\lambda$ and the penalty constant

$k$ are held unchanged, while the vector $x$ is viewed as unconstrained. Let $\Delta N$ denote the number of iterations in a cycle, regardless of whether complete convergence or incomplete convergence is achieved.

To shorten the cycle, we assign *a priori* the value of $\Delta N$. Precisely, we choose $\Delta N$ as the *smallest* number of iterations compatible with the characteristics of the particular algorithm being considered. Therefore,

$$\Delta N = 1 \tag{36}$$

for the ordinary-gradient algorithm and the modified-quasilinearization algorithm and

$$\Delta N = n \tag{37}$$

for the conjugate-gradient algorithm. Therefore, the stopping condition (30) for a cycle is bypassed and is replaced by (36) for the ordinary-gradient algorithm and the modified-quasilinearization algorithm and by (37) for the conjugate-gradient algorithm.

4.2. Multiplier Estimate. Now, we assume that the cycle length $\Delta N$ is defined by Eq. (36) for the ordinary-gradient algorithm and the modified-quasilinearization algorithm and by Eq. (37) for the conjugate-gradient algorithm. Then, we inquire about ways in which the multiplier $\lambda$ can be estimated.

*First Estimate.* We assume that Hestenes' updating rule (31) is employed at the end of any cycle of $\Delta N$ iterations. We note that the updated error in the optimum condition can be larger or smaller than the error prior to updating. Since an increase in the error $Q(x, \lambda)$ is not desirable, the multiplier $\lambda$ might be chosen as follows:

$$\begin{aligned} \lambda_2 &= \lambda_1 && \text{if} && Q(x, \lambda_0) \geq Q(x, \lambda_1), \\ \lambda_2 &= \lambda_0 && \text{if} && Q(x, \lambda_0) < Q(x, \lambda_1), \end{aligned} \tag{38}$$

where $\lambda_0$ is given by

$$\lambda_0 = \lambda_1 + 2k\varphi(x). \tag{39}$$

*Second Estimate.* The previous estimate can be improved if Hestenes' updating rule (31) is renounced and is replaced with the more general updating rule

$$\lambda_2 = \lambda_1 + 2\beta\varphi(x). \tag{40}$$

where $\beta$ is a scalar parameter. This parameter must be determined so as to produce some optimum effect.

For given values of $x$ and $\lambda_1$, a change in $\beta$ causes a change in the updated multiplier $\lambda_2$. Consequently, the updated error in the optimum condition

$$Q(x, \lambda_2) = F_x^T(x, \lambda_2) F_x(x, \lambda_2) \tag{41}$$

changes. The optimum value of $\beta$ is that which gives $Q(x, \lambda_2)$ the smallest value for given $x$ and $\lambda_1$. After combining (40)–(41), we obtain the relations

$$\begin{aligned} Q(x, \lambda_1, \beta) &= [F_x(x, \lambda_1) + \beta P_x(x)]^T [F_x(x, \lambda_1) + \beta P_x(x)], \\ Q_\beta(x, \lambda_1, \beta) &= 2 P_x^T(x) [F_x(x, \lambda_1) + \beta P_x(x)], \\ Q_{\beta\beta}(x, \lambda_1, \beta) &= 2 P_x^T(x) P_x(x), \end{aligned} \tag{42}$$

the second of which vanishes for

$$\beta = - P_x^T(x) F_x(x, \lambda_1) / P_x^T(x) P_x(x). \tag{43}$$

This value of $\beta$ minimizes $Q(x, \lambda_1, \beta)$, since $Q_{\beta\beta}(x, \lambda_1, \beta) > 0$ provided $P_x(x)$ does not vanish.

The method of multipliers with cycle stopping condition (30) replaced by (36) or (37) and with multiplier updating rule (31) replaced by (40) and (43) is called modified method of multipliers or Method MM-2. For this method, the following comments are pertinent.

(i) Equations (40) and (43) imply that

$$P_x^T(x) F_x(x, \lambda_2) = 0. \tag{44}$$

Therefore, the optimum value of the parameter $\beta$ is such that the gradient of the constraint error and the gradient of the updated augmented function are orthogonal.

(ii) The relation between the present updating rule and Hestenes' updating rule can be obtained as follows. Let the gradient of the augmented penalty function prior to updating be rewritten as

$$W_x(x, \lambda_1, k) = F_x(x, \lambda_1) + k P_x(x). \tag{45}$$

Then, combining (43) and (45) yields the relation

$$\beta = k - P_x^T(x) W_x(x, \lambda_1, k) / P_x^T(x) P_x(x), \tag{46}$$

which shows that

$$\beta = k, \tag{47}$$

providing

$$W_i(x, \lambda_1, k) = 0 \quad \text{or} \quad P_x^T(x) W_x(x, \lambda_1, k) = 0. \tag{48}$$

Clearly, the present updating rule and Hestenes' updating rule are identical if applied at complete convergence, that is, at a point where Ineq. (30) is satisfied.

(iii) Equations (40) and (43) are to be employed at the beginning of each cycle of $JN$ iterations, including the first cycle. However, for the first cycle, $\lambda_1$ is not defined. This being the case, an assumption is necessary, and the simplest assumption is

$$\lambda_1 = 0. \tag{49}$$

(iv) Whenever (40) and (43) are employed, the algorithm should be started at a point where $\varphi(x) \leq 0$.

*Third Estimate.* In the previous section, the Lagrange multiplier was estimated within the frame of the one-parameter family (40). An even better estimate can be obtained by removing this limitation and minimizing the performance index $Q(x, \lambda)$ with respect to the $q$-vector $\lambda$ for given $x$. After observing that

$$Q(x, \lambda) = [f_x(x) + \varphi_x(x)\lambda]^T [f_x(x) + \varphi_x(x)\lambda],$$
$$Q_\lambda(x, \lambda) = 2\varphi_x^T(x)[f_x(x) + \varphi_x(x)\lambda], \tag{50}$$
$$Q_{\lambda\lambda}(x, \lambda) = 2\varphi_x^T(x)\varphi_x(x),$$

we see that the optimal multiplier is defined by the relation

$$\varphi_x^T(x)\varphi_x(x)\lambda + \varphi_x^T(x)f_x(x) = 0, \tag{51}$$

which is a system of $q$ scalar equations in the $q$ components of the multiplier $\lambda$. The solution of (51) minimizes $Q$, since the matrix (50-3) is positive definite. The following comments are pertinent.

(i) On premultiplication by $\varphi^T(x)$, Eq. (51) leads to

$$P_x^T(x)P_x(x, \lambda) = 0. \tag{52}$$

Therefore, the optimum value of the multiplier $\lambda$ is such that the gradient of the constraint error and the gradient of the updated augmented function are orthogonal.

(ii) Obviously, Eq. (51) supplies the best estimate of the Lagrange

multiplier $\lambda$ compatible with any given position vector $x$ (Ref. 3). However, its use requires the solution of a system of $q$ linear equations in $q$ unknowns. This being the case, the decision on whether to use (40) and (43) or (51) should be made on the basis of the particular algorithm employed.

(iii) For the ordinary-gradient algorithm and the conjugate-gradient algorithm, the displacement vector $\Delta x$ is computed without solving a system of linear equations. Hence, the estimation of $\lambda$ should be made with (40) and (43) rather than (51).

(iv) For the modified-quasilinearization algorithm, the displacement $\Delta x$ is computed by solving a system of $n$ linear equations in $n$ unknowns. Hence, it is optional to estimate $\lambda$ with (40) and (43) or (51). In the sequel, the estimate given by (40) and (43) is employed.

**4.3. Penalty Constant Estimate.** Now, the question arises as to whether the penalty constant can be selected in such a manner as to improve the convergence characteristics of Method MM-2. In this section, two techniques are presented for updating the penalty constant at the end of a cycle, one suitable for the ordinary-gradient algorithm and one suitable for the conjugate-gradient algorithm and the modified-quasilinearization algorithm. Method MM-2 with Eqs. (40) and (43) completed by a relation updating the penalty constant is called Method MM-3.

*Ordinary-Gradient Algorithm.* When this algorithm is employed, the multiplier updating rule (40) and (43) induces an interesting characteristic: a descent property in the constraint error $P(x)$. This characteristic is utilized in this section to establish an updating rule for the penalty constant.

While the ordinary-gradient algorithm is described in detail in Section 5, we note here that the displacement $\Delta x$ leading from the nominal point $x$ to the varied point $\tilde{x}$ produces a change in the constraint error $P(x)$. To first order, this change is given by

$$\delta P(x) = -\alpha k P_x^T(x) P_x(x), \tag{53}$$

where $\alpha$ is the stepsize, and $k$ is the penalty constant. Therefore, $\delta P(x) < 0$, since $\alpha > 0$ and $k > 0$. This result guarantees that

$$P(\tilde{x}) < P(x), \tag{54}$$

providing $\alpha$ is sufficiently small.

Among all the values which can be attributed to the penalty con-

stant, we select $k$ in such a way that, on the average, the constraints are satisfied to first order. Thereby, we determine $k$ from the relation

$$\delta P(x) = -2 \lambda P(x). \tag{55}$$

Comparing (53) and (55), we see that the appropriate value of the penalty constant should be

$$k = 2P(x)/P_x^T(x) \, P_x(x). \tag{56}$$

Since both the numerator and the denominator of the right-hand side of Eq. (56) contain equal powers of $q(x)$, the penalty constant $k$ varies slowly along the algorithm and is finite at convergence. For the particular case of a single scalar constraint ($q = 1$), Eq. (56) reduces to

$$k = 1/2 \varphi_x^T(x) \, \varphi_x(x). \tag{57}$$

where $\varphi_x(x)$ is now an $n$-vector.

*Conjugate-Gradient Algorithm and Modified-Quasilinearization Algorithm.* The penalty constant estimate developed for the ordinary-gradient algorithm is based on the descent property (53) and the descent requirement (55). It produces a slowly varying penalty constant (56), which is finite at convergence.

For the conjugate-gradient algorithm and the modified-quasilinearization algorithm, the penalty constant (56) is not desirable for the reasons indicated below. Consider a quadratic function $J(x)$ and a linear constraint $q(x)$. Regardless of the value of $k$, the augmented penalty function $W(x, \lambda, k)$ is quadratic in $x$. Theoretically speaking, the optimality condition (24) is satisfied exactly after $n$ iterations of the conjugate-gradient algorithm and after one iteration of the modified quasilinearization algorithm; hence, the theoretical optimum value of the penalty constant should be $k = \infty$, in that it guarantees simultaneous satisfaction of the constraint equation $q(x) = 0$ at the end of a cycle. In a practical digital computer, this result means that large values of the penalty constant should be employed if fast convergence is desired.

If the function $J(x)$ is nonquadratic and/or the constraint $q(x)$ is nonlinear, the above reasoning is approximately true near the solution of the constrained minimization problem. This leads to the concept of programming $k$ so as to achieve moderate values far away from the solution and large values near the solution, even though these large values are not as large as those needed with the standard penalty function method.

Possible choices of the penalty constant satisfying the above property are the following:

$$k_0 = |\lambda_2^T \varphi(x)|/P(x) \tag{58}$$

or

$$k_0 = \sqrt{[\varphi_x(x) \lambda_2]^T [\varphi_x(x) \lambda_2]/P_x^T(x) \, P_x(x)}. \tag{59}$$

If Eq. (58) is employed, the order of magnitude of the linear term and the quadratic term appearing in the augmented penalty function is the same. If Eq. (59) is employed, the order of magnitude of the gradient of the linear term and the gradient of the quadratic term appearing in the augmented penalty function is the same.[*]

If the convergence requirements on the constraint error $P(x)$ and the error in the optimum condition $Q(x, \lambda)$ are the same, it might be desirable to reduce $P(x)$ and $Q(x, \lambda)$ at approximately the same rate. With this idea in mind, we propose the following updating rule for the penalty constant:

$$k_2 = \min(k_0, k_1) \qquad \text{if} \quad P(x) \le Q(x, \lambda_2), \tag{60}$$
$$k_2 = \max(k_0, \pi k_1) \qquad \text{if} \quad P(x) > Q(x, \lambda_2).$$

where $\pi \ge 1$, $k_0$ is given by (58) or (59), $k_1$ is the penalty constant prior to updating, and $k_2$ is the penalty constant after updating. Obviously, (60-1) prevents an increase of the penalty constant if the constraint error is relatively small. Conversely, (60-2) prevents a decrease of the penalty constant if the constraint error is relatively large.

Equation (60) is to be employed at the beginning of each cycle of $\Delta N$ iterations, including the first cycle. However, for the first cycle, $k_1$ is not defined. This being the case, an assumption is necessary, and the simplest assumption is

$$k_1 = k_0. \tag{61}$$

### 4.4. Summary of Methods.

In this section, we summarize the combination of methods arising from the previous discussion, as follows.

(i) Method MM-1 is characterized by updating decision (30), multiplier estimate (31), and penalty constant unchanged throughout a particular algorithm.

(ii) Method MM-2 is characterized by updating decision (36) or

---

(37), multiplier estimate (40) and (43), and penalty constant unchanged throughout a particular algorithm.

(iii) Method MM-3 is characterized by updating decision (36) or (37), multiplier estimate (40) and (43), and penalty constant estimate (56) or (58) and (60).

*Remark.* The updating decision (36) is to be employed with the ordinary-gradient algorithm and the modified-quasilinearization algorithm, and the updating decision (37) is to be employed with the conjugate-gradient algorithm. Also, the penalty constant estimate (56) is to be employed with the ordinary-gradient algorithm, and the penalty constant estimate (58) and (60) is to be employed with the conjugate-gradient algorithm and the modified-quasilinearization algorithm.

## 5. Unconstrained Minimization Algorithms

In this section, the unconstrained minimization algorithms employed to compute the displacement vector $\Delta x$ in connection with Methods MM-2 and MM-3 are described. They are the ordinary-gradient algorithm, the conjugate-gradient algorithm, and the modified-quasilinearization algorithm. All of these algorithms make use of the augmented penalty function

$$W(x, \lambda, k) \quad F(x, \lambda) + kP(x), \tag{62}$$

where

$$F(x, \lambda) \quad f(x) + \lambda^T \varphi(x), \qquad P(x) = \varphi^T(x)\varphi(x), \tag{63}$$

and are employed with this understanding: in each cycle of $\Delta N$ iterations, the multiplier $\lambda$ and the penalty constant $k$ are held unchanged, and the vector $x$ is viewed as unconstrained.

**5.1. Ordinary-Gradient Algorithm.** Let $x$ denote the nominal point, $\tilde{x}$ the varied point, $\Delta x$ the displacement leading from the nominal point to the varied point, and $\alpha$ the stepsize. With this understanding, the ordinary-gradient algorithm is represented by

$$F_x(x, \lambda) \quad f_x(x) + \varphi_x(x)\lambda. \tag{64-1}$$

$$P_x(x) \quad 2\varphi_x(x)\varphi(x). \tag{64-2}$$

$$W_x(x, \lambda, k) \quad F_x(x, \lambda) + kP_x(x). \tag{64-3}$$

$$p \quad W_x(x, \lambda, k). \tag{64-4}$$

$$\Delta x = -\alpha p, \tag{64-5}$$

$$\tilde{x} = x + \Delta x. \tag{64-6}$$

For given nominal point $x$, multiplier $\lambda$, and penalty constant $k$, Eqs. (64) constitute a complete iteration leading to the varied point $\tilde{x}$, providing one specifies the stepsize $\alpha$.

*Descent Properties.* To first order, the changes in the function $W(x, \lambda, k)$ and $P(x)$ are given by

$$\delta W(x, \lambda, k) = W_x^T(x, \lambda, k)\Delta x, \qquad \delta P(x) = P_x^T(x)\Delta x, \tag{65}$$

which, in the light of (64), become

$$\delta W(x, \lambda, k) = -\alpha W_x^T(x, \lambda, k) W_x(x, \lambda, k).$$

$$\delta P(x) = -\alpha P_x^T(x)[F_x(x, \lambda) + kP_x(x)]. \tag{66}$$

Recalling Eq. (44), we see that the following orthogonality condition holds:

$$P_x^T(x) F_x(x, \lambda) = 0, \tag{67}$$

so that

$$\delta W(x, \lambda, k) = -\alpha W_x^T(x, \lambda, k) W_x(x, \lambda, k),$$

$$\delta P(x) = -\alpha k P_x^T(x) P_x(x). \tag{68}$$

Since the right-hand sides of (68) are negative, the following inequalities can be enforced for $\alpha$ sufficiently small:

$$W(\tilde{x}, \lambda, k) < W(x, \lambda, k), \qquad P(\tilde{x}) < P(x). \tag{69}$$

While enforcement of (69-1) is mandatory, enforcement of (69-2) is optional, but can be used in order to give greater stability to the ordinary-gradient algorithm.

**5.2. Conjugate-Gradient Algorithm.** Let $x$ denote the nominal point, $\hat{x}$ the previous point, $\tilde{x}$ the varied point, $\Delta x$ the displacement leading from the nominal point to the varied point, $p$ the present search direction, $\hat{p}$ the previous search direction, $\gamma$ the directional coefficient, and $\alpha$ the stepsize. With this understanding, the conjugate-gradient algorithm is represented by

$$F_x(x, \lambda) \quad f_x(x) + \varphi_x(x)\lambda, \tag{70-1}$$

$$P_x(x) \quad 2\varphi_x(x)\varphi(x). \tag{70-2}$$

$$W_s(x, \lambda, k) = F_s(x, \lambda) + kP_s(x), \tag{70-3}$$

$$\gamma = W_x^T(x, \lambda, k)\, W_x(x, \lambda, k)/W_{\hat{x}}^T(\hat{x}, \lambda, k)\, W_{\hat{x}}(\hat{x}, \lambda, k), \tag{70-4}$$

$$p = W_x(x, \lambda, k) + \gamma \hat{p}, \tag{70-5}$$

$$\Delta x = - \alpha p, \tag{70-6}$$

$$\hat{x} = x + \Delta x. \tag{70-7}$$

For given nominal point $x$, multiplier $\lambda$, directional coefficient $\gamma$, and penalty constant $k$, Eqs. (70) constitute a complete iteration leading to the varied point $\hat{x}$, providing one specifies the stepsize $\alpha$. For the first iteration of a cycle, Eq. (70-4) is bypassed and is replaced by $\gamma = 0$.

*Descent Properties.* To first order, the changes in the function $W(x, \lambda, k)$ and $P(x)$ are given by Eqs. (65) which, in the light of (70), become

$$\delta W(x, \lambda, k) = - \alpha W_x^T(x, \lambda, k)\, [W_x(x, \lambda, k) + \gamma \hat{p}], \tag{71}$$

$$\delta P(x) = - \alpha P_x^T(x)\, [F_x(x, \lambda) + kP_x(x) + \gamma \hat{p}].$$

For the first iteration of a cycle, relation (67) must be applied in conjunction with $\gamma = 0$, leading to

$$\delta W(x, \lambda, k) = - \alpha W_x^T(x, \lambda, k)\, W_x(x, \lambda, k), \tag{72}$$

$$\delta P(x) = - \alpha k P_x^T(x)\, P_x(x).$$

For subsequent iterations, relation (67) does not hold and $\gamma \neq 0$. However, since the previous stepsize is optimized, the following orthogonality relation can be invoked:

$$W_x^T(x, \lambda, k)\, \hat{p} = 0, \tag{73}$$

with the consequence that

$$\delta W(x, \lambda, k) = - \alpha W_x^T(x, \lambda, k)\, W_x(x, \lambda, k), \tag{74}$$

$$\delta P(x) = - \alpha P_x^T(x)\, [F_x(x, \lambda) + kP_x(x) + \gamma \hat{p}].$$

Inspection of (72) and (74) shows that the descent property on the augmented penalty function

$$W(\hat{x}, \lambda, k) < W(x, \lambda, k) \tag{75}$$

can be enforced for all iterations of a cycle regardless of the value of $k$.

On the other hand, the descent property on the constraint error

$$P(\hat{x}) < P(x) \tag{76}$$

can be enforced for the first iteration of a cycle regardless of the value of $k$ and for subsequent iterations only if $k$ is sufficiently large.

**5.3. Modified-Quasilinearization Algorithm.** Let $x$ denote the nominal point, $\hat{x}$ the varied point, $\Delta x$ the displacement leading from the nominal point to the varied point, $p$ the search direction, $\rho = \pm 1$ the direction factor, and $\alpha$ the stepsize. With this understanding, the modified-quasilinearization algorithm is represented by

$$F_x(x, \lambda) = f_x(x) + \varphi_x(x)\lambda, \tag{77-1}$$

$$P_x(x) = 2\varphi_x(x)\, \varphi(x), \tag{77-2}$$

$$W_x(x, \lambda, k) = F_x(x, \lambda) + kP_x(x), \tag{77-3}$$

$$F_{xx}(x, \lambda) = f_{xx}(x) + \varphi_{xx}(x)\lambda, \tag{77-4}$$

$$P_{xx}(x) = 2[\varphi_{xx}(x)\, \varphi(x) + \varphi_x(x)\, \varphi_x^T(x)], \tag{77-5}$$

$$W_{xx}(x, \lambda, k) = F_{xx}(x, \lambda) + kP_{xx}(x), \tag{77-6}$$

$$W_{xx}(x, \lambda, k)\, \Delta + W_x(x, \lambda, k) = 0, \tag{77-7}$$

$$p = \text{sign}[W_x^T(x, \lambda, k)\, \Delta], \tag{77-8}$$

$$p = \rho \Delta, \tag{77-9}$$

$$\Delta x = - \alpha p, \tag{77-10}$$

$$\hat{x} = x + \Delta x. \tag{77-11}$$

For given nominal point $x$, multiplier $\lambda$, and penalty constant $k$, Eqs. (77) constitute a complete iteration leading to the varied point $\hat{x}$, providing one specifies the stepsize $\alpha$.

*Descent Properties.* To first order, the changes in the functions $W(x, \lambda, k)$ and $P(x)$ are given by Eqs. (65) which, in the light of (77), become

$$\delta W(x, \lambda, k) = - \alpha\, \text{sign}[W_x^T(x, \lambda, k)\, \Delta]\, W_x^T(x, \lambda, k)\, \Delta, \tag{78}$$

$$\delta P(x) = - \alpha\, \text{sign}[F_x^T(x, \lambda)\, \Delta + kP_x^T(x)\, \Delta]\, P_x^T(x)\, \Delta.$$

Inspection of (78) shows that the descent property on the augmented penalty function

$$W(\hat{x}, \lambda, k) < W(x, \lambda, k) \tag{79}$$

can be enforced regardless of the value of $k$. On the other hand, the descent property on the constraint error

$$P(\tilde{x}) < P(x) \tag{80}$$

can be enforced for any $k$ if

$$P_x{}^T(x, \lambda)\, A/P_x{}^T(x)\, A > 0 \tag{81}$$

and only for $k$ sufficiently large if

$$P_x{}^T(x, \lambda)\, A/P_x{}^T(x)\, A < 0. \tag{82}$$

## 6. Stepsize Determination

For all of the previous algorithms, the position vector at the end of any step can be written as

$$\tilde{x} = x - \alpha p, \tag{83}$$

where $p$ denotes the search direction. This is a one-parameter family of varied points $\tilde{x}$, for which the augmented penalty function (62) takes the form

$$W(\tilde{x}, \lambda, k) = W(x - \alpha p, \lambda, k) = W(\alpha). \tag{84}$$

A precise search to be employed with the conjugate-gradient algorithm and an approximate search to be employed with the ordinary-gradient algorithm and the modified-quasilinearization algorithm are described below.

*Precise Search.* We now assume that a minimum of $W(\alpha)$ exists. Then, we employ some one-dimensional search scheme (for instance, quadratic interpolation, cubic interpolation, or quasilinearization) to determine the value of $\alpha$ for which

$$W_\alpha(\alpha) = 0. \tag{85}$$

Ideally, this procedure should be used iteratively until the modulus of the slope satisfies any of the following inequalities:

$$|W_\alpha(\alpha)| \le \epsilon_4 \quad \text{or} \quad |W_\alpha(\alpha)| \le \epsilon_5 |W_\alpha(0)|, \tag{86}$$

where $\epsilon_4$ and $\epsilon_5$ are small, preselected numbers. Of course, the value of $\alpha$ satisfying Ineq. (86) must be such that

$$W(\alpha) < W(0). \tag{87}$$

*Approximate Search.* Since the rigorous determination of $\alpha$ might require excessive computing time, one might replace solving Eq. (85) with a particular degree of precision and determine the stepsize in a noniterative fashion. For instance, one might employ a bisection procedure on $\alpha$, starting from a reference value $\alpha = \alpha_R$, until satisfaction of Ineq. (87) occurs. For the ordinary gradient algorithm, the reference stepsize $\alpha_R$ can be chosen to be the first optimum value of $\alpha$ supplied by the search procedure. For the modified-quasilinearization algorithm, the reference stepsize $\alpha_R$ can be chosen to be $\alpha_R = 1$.

*Remark.* Optionally, Ineq. (87) can be completed by the additional inequality

$$P(\alpha) < P(0), \tag{88}$$

which is designed to give greater stability to the algorithm. Inequality (88) can be enforced in the ordinary-gradient algorithm for any $k$. In the conjugate-gradient algorithm and the modified-quasilinearization algorithm, Ineq. (88) can be enforced only for $k$ sufficiently large.

## 7. Experimental Conditions

In order to evaluate the theory, seven numerical examples were explored. Each example was solved with the ordinary-gradient algorithm, the conjugate-gradient algorithm, and the modified-quasilinearization algorithm, which were employed in conjunction with Methods MM-1, MM-2, and MM-3. All of the algorithms were programmed in FORTRAN IV, and the numerical results were obtained using a Burroughs B-5500 computer and double-precision arithmetic.

*Starting Point of the Algorithm.* For all of the examples, the nominal point chosen to start the algorithm was defined by

$$x_1 = x_2 = \cdots = x_n = 2, \tag{89}$$

where $n$ denotes the dimension of the vector $x$.

*Convergence of the Algorithm.* Convergence of an algorithm was defined through the inequality

$$P(x) + Q(x, \lambda) \le 10^{-8} \tag{90}$$

for the ordinary-gradient algorithm and the inequality

$$P(x) + Q(x, \lambda) \le 10^{-11} \tag{91}$$

for the conjugate-gradient algorithm and the modified-quasilinearization algorithm.

*Nonconvergence of the Algorithm.* Conversely, nonconvergence of an algorithm was defined by means of the inequalities

(a) $\begin{cases} N > 1000 & \text{for the ordinary-gradient algorithm,} \\ N > 200 & \text{for the conjugate-gradient algorithm,} \\ N > 100 & \text{for the modified-quasilinearization algorithm,} \end{cases}$ (92)

or

(b)    $N_s > 20$ (93)

or

(c)    $M > 0.4 \times 10^{69}$. (94)

Here, $N$ is the iteration number, $N_s$ is the number of bisections of the stepsize $\alpha$ required to satisfy Ineq. (87) [and optionally Ineq. (88)]; and $M$ is the modulus of any of the quantities employed in the algorithm. Satisfaction of Ineq. (92) indicates divergence or extreme slowness of convergence; satisfaction of Ineq. (93) indicates extreme smallness of the displacement $\Delta x$; and satisfaction of Ineq. (94) indicates exponential overflow. Each of these situations is undesirable.

*Convergence of a Cycle.* When Method MM-1 was employed, convergence of a cycle was defined through the inequality

$$W_x^T(x, \lambda_1, k)\, W_x(x, \lambda_1, k) \leqslant 10^{-8} \qquad (95)$$

for the ordinary-gradient algorithm and the inequality

$$W_x^T(x, \lambda_1, k)\, W_x(x, \lambda_1, k) \leqslant 10^{-14} \qquad (96)$$

for the conjugate-gradient algorithm and the modified-quasilinearization algorithm.

When Methods MM-2 and MM-3 were employed, convergence of a cycle was defined by

$$\Delta N = 1 \qquad (97)$$

for the ordinary-gradient algorithm and the modified-quasilinearization algorithm and by either

$$\Delta N = n \qquad (98)$$

or

$$W_x^T(x, \lambda_1, k)\, W_x(x, \lambda_1, k) \leqslant 10^{-14}. \qquad (99)$$

whichever occurred first, for the conjugate-gradient algorithm.

*Search Technique.* For the ordinary-gradient algorithm, an approximate search was employed. This consisted of one-step, corrected quasilinearization, followed by a bisection process until the inequality

$$W(\alpha) < W(0) \qquad (100)$$

was satisfied.

For the conjugate-gradient algorithm, a precise search was employed. This consisted of multistep, corrected quasilinearization such that, in any given step, the inequality

$$W(\alpha) < W(\alpha_0) \qquad (101)$$

was satisfied, where $\alpha_0$ is the nominal stepsize and $\alpha$ is the varied stepsize. The search was started with $\alpha_0 = 0$ and was terminated when the following stopping condition was satisfied:

$$W_\alpha^T(\alpha) \leqslant W_\alpha^T(0) \times 10^{-4}. \qquad (102)$$

For the modified-quasilinearization algorithm, an approximate search was employed. This consisted of assigning the value

$$\alpha = 1 \qquad (103)$$

to the stepsize, followed by a bisection process until Ineq. (100) was satisfied.

*Penalty Constant Estimate.* For Method MM-3, the penalty constant $k$ was estimated with Eq. (56) for the ordinary-gradient algorithm and with Eqs. (58) and (60), with $n = 1$, for the conjugate-gradient algorithm and the modified-quasilinearization algorithm.

## 8. Numerical Examples

In this section seven numerical examples are described. The first example pertains to a quadratic function subject to linear constraints. The remaining examples pertain to nonquadratic functions subject to nonlinear constraints.

Example 8.1. Consider the problem of minimizing the function

$$f = (x_1 - x_2)^2 + (x_2 + x_3 - 2)^2 + (x_4 - 1)^2 + (x_5 - 1)^2 \qquad (104)$$

subject to the constraints

$$x_1 + 3x_2 = 0, \qquad x_3 + x_4 - 2x_5 = 0, \qquad x_2 - x_5 = 0. \qquad (105)$$

This function admits the relative minimum $f = 4.0930$ at the point defined by

$$x_1 = -0.7674, \quad x_2 = 0.2558, \quad x_3 = 0.6279,$$
$$x_4 = -0.1162, \quad x_5 = 0.2558 \tag{106}$$

and

$$\lambda_1 = 2.0165, \quad \lambda_2 = 2.2325, \quad \lambda_3 = -5.9534. \tag{107}$$

Example 8.2. Consider the problem of minimizing the function

$$f = (x_1 - 1)^2 + (x_1 - x_2)^2 + (x_2 - x_3)^2 \tag{108}$$

subject to the constraint

$$x_1(1 + x_2^2) + x_3^4 - 4 - 3\sqrt{2} = 0. \tag{109}$$

This function admits the relative minimum $f = 0.3256 \times 10^{-1}$ at the point defined by

$$x_1 = 1.1048, \quad x_2 = 1.1966, \quad x_3 = 1.5352 \tag{110}$$

and

$$\lambda_1 = -0.1072 \times 10^{-1}. \tag{111}$$

Example 8.3. Consider the problem of minimizing the function

$$f = (x_1 - 1)^2 + (x_1 - x_2)^2 + (x_3 - 1)^2 + (x_4 - 1)^4 + (x_5 - 1)^6 \tag{112}$$

subject to the constraints

$$x_1^2 x_4 + \sin(x_4 - x_5) - 2\sqrt{2} = 0, \quad x_2 + x_3^4 x_4^2 - 8 - \sqrt{2} = 0. \tag{113}$$

This function admits the relative minimum $f = 0.2415$ at the point defined by

$$x_1 = 1.1661, \quad x_2 = 1.1821, \quad x_3 = 1.3802,$$
$$x_4 = 1.5064, \quad x_5 = 0.6109 \tag{114}$$

and

$$\lambda_1 = -0.8553 \times 10^{-1}, \quad \lambda_2 = -0.3187 \times 10^{-1}. \tag{115}$$

Example 8.4. Consider the problem of minimizing the function

$$(x_1 - 1)^2 + (x_1 - x_2)^2 + (x_2 - x_3)^2 + (x_3 - x_4)^4 + (x_4 - x_5)^4 \tag{116}$$

subject to the constraints

$$x_1 + x_2^2 + x_3^3 - 2 - 3\sqrt{2} = 0,$$
$$x_2 - x_3^2 + x_4 + 2 - 2\sqrt{2} = 0, \quad x_1 x_5 - 2 = 0. \tag{117}$$

This function admits the relative minimum $f = 0.7877 \times 10^{-1}$ at the point defined by

$$x_1 = 1.1911, \quad x_2 = 1.3626, \quad x_3 = 1.4728,$$
$$x_4 = 1.6350, \quad x_5 = 1.6790 \tag{118}$$

and

$$\lambda_1 = -0.3882 \times 10^{-1}, \quad \lambda_2 = -0.1672 \times 10^{-1},$$
$$\lambda_3 = -0.2879 \times 10^{-3}. \tag{119}$$

Example 8.5. Consider the problem of minimizing the function

$$f = 0.01(x_1 - 1)^2 + (x_2 - x_1^2)^2 \tag{120}$$

subject to the inequality constraint

$$x_1 \leqslant -1. \tag{121}$$

Introduce the auxiliary variable $x_3$ defined by

$$x_1 + x_3^2 + 1 = 0. \tag{122}$$

Then, the previous problem can be recast as that of minimizing the function (120) subject to the equality constraint (122). The function (120) admits the relative minimum $f = 0.04$ at the point defined by

$$x_1 = -1, \quad x_2 = 1, \quad x_3 = 0 \tag{123}$$

and

$$\lambda_1 = 0.04. \tag{124}$$

Example 8.6. Consider the problem of minimizing the function

$$f = -x_1 \tag{125}$$

subject to the inequality constraints

$$x_2 \geqslant x_1^3, \quad x_2 \leqslant x_1^3. \tag{126}$$

Introduce the auxiliary variables $x_3$ and $x_4$ defined by

$$x_2 - x_1^3 - x_3^2 = 0, \quad x_1^3 - x_2 - x_4^2 = 0. \tag{127}$$

Then, the previous problem can be recast as that of minimizing the function (125) subject to the equality constraints (127). The function (125) admits the relative minimum $f = -1$ at the point defined by

$$x_1 = 1, \quad x_2 = 1, \quad x_3 = 0, \quad x_4 = 0 \tag{128}$$

and

$$\lambda_1 = -1, \quad \lambda_2 = -1. \tag{129}$$

**Example 8.7.** Consider the problem of minimizing the function

$$f = \log x_3 - x_2 \tag{130}$$

subject to the equality constraint

$$x_2^2 + x_3^2 - 4 = 0 \tag{131}$$

and the inequality constraint

$$x_3 \geq 1. \tag{132}$$

Introduce the auxiliary variable $x_1$ defined by

$$x_3 = 1 + x_1^2. \tag{133}$$

Then, the previous problem can be recast as that of minimizing the function

$$f = \log(1 + x_1^2) - x_2 \tag{134}$$

subject to the equality constraint

$$(1 + x_1^2)^2 + x_2^2 - 4 = 0. \tag{135}$$

Note that $x_3$ has been eliminated from the problem and can be computed a posteriori with (133). The function (134) admits the relative minimum $f = -\sqrt{3}$ at the point defined by

$$x_1 = 0, \quad x_2 = \sqrt{3}, \quad x_3 = 1 \tag{136}$$

and

$$\lambda_1 = 1/2\sqrt{3}. \tag{137}$$

## 9. Results and Conclusions

The examples described in Section 8 were solved with Hestenes' method of multipliers (Method MM-1) and the modified method of multipliers (Methods MM-2 and MM-3) according to the experimental

Table 1. Method MM-1, ordinary-gradient algorithm, number of iterations $N_*$.

| $h$ | Examples | | | | | | |
|---|---|---|---|---|---|---|---|
| | 8.1 | 8.2 | 8.3 | 8.4 | 8.5 | 8.6 | 8.7 |
| $10^{-2}$ | (a) | 24 | 565 | 580 | (a) | (a) | 44 |
| $10^{-1}$ | 664 | 73 | 847 | 193 | 516 | (a) | 17 |
| $10^0$ | 350 | 646 | (a) | 431 | 58 | 946 | 35 |
| $10^1$ | 863 | (a) | (a) | 377 | 761 | 762 | 155 |
| $10^2$ | (a) | (a) | (a) | (a) | 92 | 746 | 774 |

(a) Number of iterations exceeded 1000.

Table 2. Method MM-2, ordinary-gradient algorithm, number of iterations $N_*$.

| $h$ | Examples | | | | | | |
|---|---|---|---|---|---|---|---|
| | 8.1 | 8.2 | 8.3 | 8.4 | 8.5 | 8.6 | 8.7 |
| $10^{-2}$ | (a) | 14 | 135 | 102 | 138 | 455 | 54 |
| $10^{-1}$ | 394 | 21 | 433 | 94 | 39 | 318 | 11 |
| $10^0$ | 94 | 661 | (a) | 306 | 33 | 278 | 32 |
| $10^1$ | 266 | 998 | (a) | (a) | 161 | 129 | 49 |
| $10^2$ | (a) | (a) | (a) | (a) | 194 | (a) | 133 |

(a) Number of iterations exceeded 1000.

Table 3. Method MM-3, ordinary-gradient algorithm, number of iterations $N_*$.

| | Examples | | | | | | |
|---|---|---|---|---|---|---|---|
| $h$ | 8.1 | 8.2 | 8.3 | 8.4 | 8.5 | 8.6 | 8.7 |
| Eq. (56) | 88 | 13 | 193 | 101 | 20 | 275 | 17 |

Table 4.  Method MM-1, conjugate-gradient algorithm, number of iterations $N_*$.

| $h$ | Examples | | | | | | |
|---|---|---|---|---|---|---|---|
| | 8.1 | 8.2 | 8.3 | 8.4 | 8.5 | 8.6 | 8.7 |
| $10^{-4}$ | (a) | 25 | (a) | (a) | 93 | (a) | 58 |
| $10^{-2}$ | (a) | 28 | 142 | 138 | (a) | (a) | 16 |
| $10^0$ | (a) | 37 | 103 | 84 | 40 | 184 | 14 |
| $10^2$ | 49 | 54 | 145 | 81 | 71 | 32 | 11 |
| $10^4$ | 20 | 76 | 176 | 95 | 164 | 33 | 18 |

(a) Number of iterations exceeded 200.

Table 5.  Method MM-2, conjugate-gradient algorithm, number of iterations $N_*$.

| $h$ | Examples | | | | | | |
|---|---|---|---|---|---|---|---|
| | 8.1 | 8.2 | 8.3 | 8.4 | 8.5 | 8.6 | 8.7 |
| $10^{-4}$ | (a) | 46 | (a) | (a) | (a) | (a) | 55 |
| $10^{-2}$ | (a) | 16 | 76 | 97 | 63 | (a) | 13 |
| $10^0$ | (a) | 33 | 79 | 42 | 49 | 150 | 10 |
| $10^2$ | 47 | 48 | 110 | 50 | 61 | 36 | 12 |
| $10^4$ | 20 | 71 | 172 | 49 | (a) | 28 | 16 |

(a) Number of iterations exceeded 200.

Table 6.  Method MM-3, conjugate-gradient algorithm, number of iterations $N_*$.

| $h$ | Examples | | | | | | |
|---|---|---|---|---|---|---|---|
| | 8.1 | 8.2 | 8.3 | 8.4 | 8.5 | 8.6 | 8.7 |
| Eqs. (59), (60) | 43 | 15 | 75 | 55 | 32 | 50 | 14 |

Table 7.  Method MM-1, modified-quasilinearization algorithm, number of iterations $N_*$.

| $h$ | Examples | | | | | | |
|---|---|---|---|---|---|---|---|
| | 8.1 | 8.2 | 8.3 | 8.4 | 8.5 | 8.6 | 8.7 |
| $10^{-4}$ | (a) | 14 | (a) | (a) | 36 | (a) | 66 |
| $10^{-2}$ | (a) | 11 | 31 | 34 | 21 | (a) | 25 |
| $10^0$ | 46 | 14 | 37 | 17 | 14 | 86 | 20 |
| $10^2$ | 10 | 13 | 19 | 13 | 26 | 29 | 23 |
| $10^4$ | 4 | 13 | 25 | 22 | 46 | 20 | 35 |

(a) Number of iterations exceeded 100.

Table 8.  Method MM-2, modified-quasilinearization algorithm, number of iterations $N_*$.

| $h$ | Examples | | | | | | |
|---|---|---|---|---|---|---|---|
| | 8.1 | 8.2 | 8.3 | 8.4 | 8.5 | 8.6 | 8.7 |
| $10^{-4}$ | (a) | 9 | 90 | (a) | (c) | (a) | 46 |
| $10^{-2}$ | (a) | 8 | 20 | 26 | (c) | (a) | 36 |
| $10^0$ | 46 | 11 | (d) | 12 | 16 | 38 | 14 |
| $10^2$ | 10 | 10 | 17 | 11 | 21 | 19 | 20 |
| $10^4$ | 4 | 12 | 23 | 21 | 45 | 16 | 33 |

(a) Number of iterations exceeded 100.
(c) Exponential overflow.
(d) Algorithm converged to a different relative minimum.

Table 9.  Method MM-3, modified-quasilinearization algorithm, number of iterations $N_*$.

| $h$ | Examples | | | | | | |
|---|---|---|---|---|---|---|---|
| | 8.1 | 8.2 | 8.3 | 8.4 | 8.5 | 8.6 | 8.7 |
| Eqs. (59), (60) | 9 | 9 | 13 | 9 | 33 | 21 | 12 |

conditions outlined. These methods were employed in conjunction with the ordinary-gradient algorithm, the conjugate-gradient algorithm, and the modified-quasilinearization algorithm. For Methods MM-1 and MM-2, several values of the penalty constant, ranging between $10^{-2}$ and $10^{2}$, were considered.

The numerical results are presented in Tables 1-9, where the number of iterations required for convergence $N_*$ is shown. Tables 1-3 refer to the ordinary-gradient algorithm, Tables 4-6 refer to the conjugate-gradient algorithm, and Tables 7-9 refer to the modified-quasilinearization algorithm.

Comparison of Methods MM-1 and MM-2 shows that, for given $k$, Method MM-2 generally exhibits faster convergence than Method MM-1. For both Methods MM-1 and MM-2, the number of iterations for convergence has a minimum with respect to $k$.

Concerning Method MM-3, the number of iterations for convergence is close to the minimum with respect to $k$ of the number of iterations for convergence of Method MM-2. In this light, Method MM-3 has very desirable characteristics.

Remark 9.1. In Section 4.2, several ways to estimate the Lagrange multiplier $\lambda$ were presented. The estimate represented by Eqs. (40) and (43) was preferred to that represented by Eq. (51) in order to avoid solving a set of $q$ linear equations in $q$ unknowns. The above statement is significant if one employs the ordinary-gradient algorithm and the conjugate-gradient algorithm, since the computation of the displacement vector $\Delta x$ is made bypassing the solution of sets of linear equations. On the other hand, if one employs the modified-quasilinearization algorithm, a case can be made for using the estimate (51) for the multiplier, since the computation of the displacement vector $\Delta x$ requires the solution of a set of $n$ linear equations in $n$ unknowns.

For the sake of discussion, let the modified-quasilinearization algorithm be employed and let the following terminology be used: (a) Method MM-3 is the modified method of multipliers with Lagrange multiplier represented by (40) and (43) and penalty constant represented by (58) and (60); and (b) Method MM-4 is the modified method of multipliers with Lagrange multiplier represented by (51) and penalty constant represented by (58) and (60).

As a point of interest, the seven numerical examples of Section 8 were solved employing both Methods MM-3 and MM-4. The comparative results are given in Table 10, where the number of iterations at convergence $N_*$ is shown. As expected, for a single scalar constraint ($q = 1$), Methods MM-3 and MM-4 converge to the solution in the

Table 10. Modified-quasilinearization algorithm, number of iterations $N_*$.

| Method | Examples | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | 8.1 | 8.2 | 8.3 | 8.4 | 8.5 | 8.6 | 8.7 |
| MM-3 | 9 | 9 | 13 | 9 | 13 | 21 | 12 |
| MM-4 | 8 | 9 | 10 | 11 | 13 | 15 | 12 |

same number of iterations. This is precisely the case with Examples 8.2, 8.5, and 8.7. On the other hand, if several constraints are present, Method MM-4 generally converges to the solution in a smaller number of iterations than Method MM-3.

Remark 9.2. Both Methods MM-3 and MM-4 employed in conjunction with the modified-quasilinearization algorithm are interesting in that they generally lead to a minimum point (or at most a saddle point), while standard quasilinearization may also lead to a maximum point.[*] To verify this concept, the following example was considered: Extremize the function

$$f = \sin(\pi x_1/12) \cos(\pi x_2/16) \tag{138}$$

subject to the constraint

$$4x_1 - 3x_2 = 0. \tag{139}$$

This function has a relative minimum $f = -0.5$ at

$$x_1 = -3, \quad x_2 = -4, \quad \text{and} \quad \lambda_1 = -\pi/96. \tag{140}$$

and a relative maximum $f = 0.5$ at

$$x_1 = 3, \quad x_2 = 4, \quad \text{and} \quad \lambda_1 = -\pi/96. \tag{141}$$

The following nominal point was considered:

$$x_1 = 2, \quad x_2 = 2. \tag{142}$$

It was found that Methods MM-3 and MM-4 employed in conjunction with the modified-quasilinearization algorithm led to the relative

---

[*] Standard quasilinearization involves the solution of a system of $n + q$ linear equations in $n + q$ unknowns (see, for example, Ref. 3).

minimum (140) in $N_s \approx 5$ iterations. On the other hand, standard quasilinearization (Ref. 3) led to the relative maximum (141) in $N_s \approx 3$ iterations.

Remark 9.3. In Ref. 11, it was suggested that Hestenes' method of multipliers could be accelerated by employing Eq. (31) at the beginning of each iteration, rather than at the beginning of each cycle. With particular reference to the conjugate-gradient algorithm, this technique has the disadvantage of producing discontinuities in the augmented penalty function $W(x, \lambda, k)$ within a cycle of $\Delta N = n$ iterations. To verify this point, numerical experiments were carried out for the seven examples of Section 8 and one value of the penalty constant, namely, $k = 1$. It was found that, in the majority of the examples, the algorithm of Ref. 11 had not converged to the solution within the required accuracy (91) in the imposed limit of 200 iterations.

## References

1. MIELE, A., HUANG, H. Y., and HEIDEMAN, J. C., *Sequential Gradient-Restoration Algorithm for the Minimization of Constrained Functions, Ordinary and Conjugate Gradient Versions*, Journal of Optimization Theory and Applications, Vol. 4, No. 4, 1969.

2. MIELE, A., HEIDEMAN, J. C., and LEVY, A. V., *Combined Conjugate Gradient-Restoration Algorithm for Mathematical Programming Problems*, Ricerche di Automatica, Vol. 2, No. 2, 1971.

3. MIELE, A., and LEVY, A. V., *Modified Quasilinearization and Optimal Initial Choice of the Multipliers, Part 1, Mathematical Programming Problems*, Journal of Optimization Theory and Applications, Vol. 6, No. 5, 1970.

4. KELLEY, H. J., *Method of Gradients*, Optimization Techniques, Edited by G. Leitmann, Academic Press, New York, 1962.

5. BRYSON, A. E., JR., and HO, Y. C., *Applied Optimal Control*, Blaisdell Publishing Company, Waltham, Massachusetts, 1969.

6. FIACCO, A. V., and McCORMICK, G. P., *Nonlinear Programming: Sequential Unconstrained Minimization Techniques*, John Wiley and Sons, New York, 1968.

7. HESTENES, M. R., *Multiplier and Gradient Methods*, Journal of Optimization Theory and Applications, Vol. 4, No. 5, 1969.

8. MIELE, A., MOSELEY, P. E., and CRAGG, E. E., *Numerical Experiments on Hestenes' Method of Multipliers for Mathematical Programming Problems*, Rice University, Aero-Astronautics Report No. 85, 1971.

*Penalty Constant Used in the Penalty Function Method for Mathematical Programming Problems*, Rice University, Aero-Astronautics Report No. 90, 1972.

10. MIELE, A., MOSELEY, P. E., and CRAGG, E. E., *A Modification of the Method of Multipliers for Mathematical Programming Problems*, Techniques of Optimization, Edited by A. V. Balakrishnan, Academic Press, New York, 1972.

11. TRIPATHI, S. S., and NARENDRA, K. S., *Constrained Optimization Using Multiplier Methods*, Journal of Optimization Theory and Applications, Vol. 9, No. 1, 1972.

# COMPARISON OF SEVERAL GRADIENT ALGORITHMS FOR MATHEMATICAL PROGRAMMING PROBLEMS

by

Angelo MIELE, J. L. TIETZE and A. V. LEVY
Department of Mechanical and Aerospace
Engineering and Material Science
Rice University
Houston, Texas

# COMPARISON OF SEVERAL GRADIENT ALGORITHMS
## FOR MATHEMATICAL PROGRAMMING PROBLEMS

by

ANGELO MIELE, J.L. TIETZE and A.V. LEVY

*This work is dedicated by the senior
author to Professor Carlo Ferrari with deep
friendship and profound admiration for his
many achievements in aerospace engineering
and applied mathematics throughout the years.*

## Abstract

In this paper, the numerical solution of the basic problem of mathematical programming is considered. This is the problem of minimizing a function $f(x)$ subject to a constraint $\varphi(x) = 0$. Here, $f$ is a scalar, $x$ an n-vector, and $\varphi$ a q-vector, with $q < n$.

Six variations of the sequential gradient-restoration algorithm and the combined gradient-restoration algorithm are considered, and their relative efficiency (in terms of number of iterations for convergence) is evaluated. The variations being considered are as follows:

(i)   SGRA-CR, sequential gradient-restoration algorithm, complete restoration,

(ii)  SGRA-IR, sequential gradient-restoration algorithm, incomplete restoration,

(iii) SGRA-OR, sequential gradient-restoration algorithm, optional restoration,

(iv)  CGRA-NR, combined gradient-restoration algorithm, no restoration,

(v)   CGRA-AR, combined gradient-restoration algorithm, alternate restoration,

(vi)  CGRA-OR, combined gradient-restoration algorithm, optional restoration.

Evaluation of these algorithms is accomplished through eight numerical examples. The first two examples pertain to quadratic functions subject to linear constraints. The remaining examples pertain to nonquadratic functions subject to nonlinear constraints. The results indicate that (a) the inclusion of a restoration phase is necessary for

rapid convergence and (b) the algorithms with alternate restoration or optional restoration are the most efficient among those considered here.

## 1. Introduction

In previous papers (refs.1,2,3), two basic algorithms for the minimization of constrained functions were developed: the sequential gradient-restoration algorithm (SGRA) and the combined gradient-restoration algorithm (CGRA). The former is an iterative algorithm which consists of the alternate succession of gradient phases and restoration phases ; the latter is an iterative algorithm in which the gradient phase and the restoration phase are combined in a single phase.

In the gradient phase of SGRA, one generates a displacement $dx$ lowering the value of the function, while avoiding excessive constraint violation; in the restoration phase of SGRA, one generates a displacement $dx$ restoring the constraint to a predetermined accuracy, while avoiding excessive change in the value of the function. On the other hand, in the gradient-restoration phase of CGRA, one generates a displacement $dx$ lowering the value of the augmented function, while simultaneously reducing the constraint violation.

In this paper, six variations of the sequential gradient-restoration algorithm and the combined gradient-restoration algorithm are considered, and their relative efficiency (in terms of number of iterations for convergence) is evaluated through eight numerical examples. The variations being considered are indicated below:

(i) SGRA-CR, sequential gradient-restoration algorithm, complete restoration,

(ii) SGRA-IR, sequential gradient-restoration algorithm, incomplete restoration,

(iii) SGRA-OR, sequential gradient-restoration algorithm, optional restoration,

(iv) CGRA-NR, combined gradient-restoration algorithm, no restoration,

(v) CGRA-AR, combined gradient-restoration algorithm, alternate restoration,

(vi) CGRA-OR, combined gradient-restoration algorithm, optional restoration.

## 2. Statement of the Problem

We consider the problem of minimizing the function

$$(1) \qquad f = f(x) ,$$

subject to the constraint

$$(2) \qquad \varphi(x) = 0 ,$$

where $f$ is a scalar, $x$ an n-vector, and $\varphi$ a q-vector, with $q < n$. Here, all vectors are column vectors. It is assumed that the first and second partial derivatives of the functions $f(x)$ and $\varphi(x)$ exist and are continuous and that the constrained minimum exists.

2.1. First-Order Conditions. From theory of maxima and minima, it is known that the above problem is equivalent to that of minimizing the augmented function

$$(3) \qquad F(x,\lambda) = f(x) + \lambda^T \varphi(x) ,$$

subject to the constraint (2). Here, the q-vector $\lambda$ is the Lagrange multiplier and the superscript $T$ denotes the transpose of a matrix. If

$$(4) \qquad F_x(x,\lambda) = f_x(x) + \varphi_x(x)\lambda$$

denotes the gradient of the augmented function, the optimum solution for $x$ and $\lambda$ must satisfy the relations

$$(5) \qquad \varphi(x) = 0 , \qquad F_x(x,\lambda) = 0 ,$$

which are a system of $n + q$ equations in the $n + q$ components of $x$ and $\lambda$. In Eqs. (4)-(5), the gradients $f_x$ and $F_x$ denotes n-vectors and the matrix $\varphi_x$ is $n \times q$.

2.2. Approximate Solutions. Since the system (5) is generally nonlinear, approximate methods must be employed. In this connection, we introduce here the scalar performance indexes

$$(6) \qquad P(x) = \varphi^T(x)\varphi(x) , \qquad Q(x,\lambda) = F_x^T(x,\lambda)F_x(x,\lambda) ,$$

which measure the errors in the constraint and the optimum condition, respectively. Then, we observe that $P = 0$ and $Q = 0$ for the optimum solution, while $P > 0$ and/or $Q > 0$ for any approximation to the solution. When approximate methods are used, they must ultimately lead to

values of  x  and $\lambda$  such that

(7)                P(x) ≤ $r_1$ ,       Q(x, $\lambda$) ≤ $r_2$  .

Alternatively, (7) can be replaced by

(8)                      R(x, $\lambda$) ≤ $r_3$  ,

where

(9)                   R(x, $\lambda$) = P(x) + Q(x, $\lambda$)

denotes the cumulative error in the constraint and the optimum condition.
In (7)-(8), $r_1$ , $r_2$ , $r_3$  are small, preselected numbers.  Note that,
if one chooses $r_1 = r_2 = r_3$ , satisfaction of Ineq.(8) implies satisfac-
tion of Ineqs.(7).


3. . Description of the Algorithms

      In this section, the algorithms being investigated are described.

      SGRA-CR    Sequential gradient-restoration algorithm, com-
plete restoration.    This algorithm consists of the alternate succession
of gradient phases and restoration phases.
      The gradient phase is started providing

(10)                        P(x) ≤ $r_1$  .

It involves a single iteration, in which the augmented function  is re-
duced subject to an upper limit for the constraint error, that is,

(11)               F($\tilde{x}$, $\lambda$) < F(x, $\lambda$)  ,  P($\tilde{x}$) ≤ $r_4$  .

The symbol  x  denotes the nominal point,  $\tilde{x}$  the varied point, and $\lambda$
the Lagrange multiplier.
      The restoration phase is started providing

(12)                        P(x) > $r_1$  .

It involves several iterations, in each of which the constraint error is
reduced, that is,

(13)                      P($\tilde{x}$) < P(x)  .

The restoration phases is terminated whenever Ineq.(10) is satisfied.

      Remark.    The algorithm is started with a gradient phase if Ineq.
(10) is satisfied or a restoration phase if Ineq.(10) is violated.
Normally, a gradient phase is followed by a restoration phase.   Occa-
sionally, the gradient phase is followed by another gradient phase, that
is, the restoration phase is bypassed: this is precisely the case when-
ever Ineq.(10) is satisfied.

      SGRA-IR.   Sequential gradient-restoration algorithm, incom-
plete restoration.    This algorithm consists of the alternate succession
of gradient phases and restoration phases.
      The gradient phase is started regardless of whether Ineq.(10) is sat-
isfied.    It involves a single iteration, in which the augmented function
is reduced subject to an upper limit on the constraint error, that is,

(14)               F($\tilde{x}$, $\lambda$) < F(x, $\lambda$) ,  P($\tilde{x}$) ≤ P(x) + $r_4$  .

      The restoration phase is started only if Ineq.(12) is satisfied.
It involves a single iteration, in which the constraint error is reduced
in accordance with Ineq.(13).
      The starting condition and the bypassing condition for SGRA-IR are
identical with those of SGRA-CR (see Remark).

      SGRA-OR.   Sequential gradient-restoration algorithm, op-
tional restoration.    This algorithm consists of the alternate succession
of gradient phases and restoration phases.
      The gradient phase is started providing

(15)                      Z(x, $\lambda$) ≤ 1 ,

where the parameter  Z  is defined by

(16)                   Z = r P(x)/Q(x, $\lambda$) ,

with

(17)                      r = $r_2$/$r_1$  .

It involves a single iteration, in which the augmented function is  re-
duced in accordance with Ineqs.(14).
      The restoration phase is started providing

(18)                      Z(x, $\lambda$) > 1  .

It involves several iterations, in each of which the constraint error is reduced in accordance with Ineq. (13). The restoration phase is terminated whenever Ineq. (15) is satisfied.

The bypassing condition for SGRA-OR is identical with that of SGRA-CR (see Remark)

CGRA-NR. *Combined gradient-restoration algorithm, no restoration.* In this algorithm, the gradient phase and the restoration phase are combined together in a single phase. It involves a single iteration, in which the augmented function is reduced in accordance with Ineqs. (14).

CGRA-AR. *Combined gradient-restoration algorithm, alternate restoration.* This algorithm consists of the alternate succession of combined gradient-restoration phases and restoration phases.

The combined gradient-restoration phase is started regardless of whether Ineq. (10) is satisfied. It involves a single iteration, in which the augmented function is reduced in accordance with Ineqs. (14).

The restoration phase is started only if Ineq. (12) is satisfied. It involves a single iteration, in which the constraint error is reduced in accordance with Ineq. (13).

The starting condition and the bypassing condition for CGRA-AR are identical with those of SGRA-CR (see Remark).

CGRA-OR. *Combined gradient-restoration algorithm, optional restoration.* This algorithm consists of the alternate succession of combined gradient-restoration phases and restoration phases.

The combined gradient-restoration phase is started providing Ineq. (15) is satisfied. It involves a single iteration, in which the augmented function is reduced in accordance with Ineqs. (14).

The restoration phase is started providing Ineq. (18) is satisfied. It involves several iterations, in each of which the constraint error is reduced in accordance with Ineq. (13). The restoration phase is terminated whenever Ineq. (15) is satisfied.

The bypassing condition for CGRA-OR is identical with that of SGRA-CR (see Remark).

Remark. For the algorithms with optional restoration, the multiplier $\lambda$ appearing in (15)-(18) is computed as follows. For SGRA-OR, Eq. (19-1) must be solved with $C_1 = 1$ and $C_2 = 0$. For CGRA-OR, Eq. (19-1) must be solved with $C_1 = 1$ and $C_2 = 1$.

## 4. Generalized Algorithm

Let $x$ denote the nominal point, $\tilde{x}$ the varied point, $\Delta x$ the displacement leading from the nominal point to the varied point, and $\alpha$ the stepsize. With this understanding, the previous algorithms can be represented in the following generalized form:

$$(19\text{-}1) \qquad \varphi_x^T(x)\varphi_x(x)\lambda + C_1 \varphi_x^T(x) f_x(x) - C_2 \varphi(x) = 0 ,$$

$$(19\text{-}2) \qquad p = C_1 f_x(x) + \varphi_x(x)\lambda ,$$

$$(19\text{-}3) \qquad \Delta x = -\alpha p ,$$

$$(19\text{-}4) \qquad \tilde{x} = x + \Delta x .$$

For given nominal point $x$ and constants $C_1$ and $C_2$, Eqs. (19) constitute a complete iteration leading to the varied point $\tilde{x}$, providing one specifies the stepsize $\alpha$. The constants $C_1$ and $C_2$ depend on the particular algorithm and take the values given in Table 1. The detailed derivation of Eqs. (19) is presented in refs. 1,2,3 and hence, is not repeated here.

Table 1. Characteristic constants.

| Algorithm | Phase | $C_1$ | $C_2$ |
|---|---|---|---|
| SGRA | Gradient | 1 | 0 |
| | Restoration | 0 | 1 |
| CGRA | Gradient-restoration | 1 | 1 |
| | Restoration | 0 | 1 |

## 5. Stepsize Determination

For all of the previous algorithms, the position vector at the end of any step can be written as

$$(20) \qquad \tilde{x} = x - \alpha p ,$$

where $p$ denotes the search direction, which is given by (19-2). This is a one-parameter family of varied points $\tilde{x}$, for which the augmented function (3), the constraint error (6-1), and the error in the optimum condition (6-2) take the form

$$(21) \qquad F(\tilde{x}, \lambda) = F(x - \alpha p, \lambda) = \tilde{F}(\alpha) \; ,$$

$$(22) \qquad P(\tilde{x}) = P(x - \alpha p) = \tilde{P}(\alpha) \; ,$$

$$(23) \qquad Q(\tilde{x}, \lambda) = Q(x - \alpha p, \lambda) = \tilde{Q}(\alpha) \; .$$

For the gradient phase of a SGRA-algorithm or the combined gradient-restoration phase of a CGRA-algorithm, Ineqs. (11) and (14) can be written in the general form

$$(24) \qquad \tilde{F}(\alpha) < \tilde{F}(0) \; , \quad \tilde{P}(\alpha) < \tilde{P}(0) + \epsilon_4 \; .$$

Their satisfaction can be ensured by employing a bisection process, starting from a suitably chosen reference stepsize

$$(25) \qquad \alpha = \alpha_0 \; .$$

For the determination of the reference stepsize, see Section 6.

For the restoration phase of a SGRA-algorithm or a CGRA-algorithm, Ineq. (13) can be written as

$$(26) \qquad \tilde{P}(\alpha) < \tilde{P}(0) \; .$$

Its satisfaction can be ensured by employing a bisection process, starting from the reference stepsize

$$(27) \qquad \alpha = 1 \; .$$

This value reduces the constraint error $P(x)$ to zero, if the constraint function $\varphi(x)$ is linear in $x$.

## 6. Reference Stepsize

The search technique outlined in Section 5 for the gradient stepsize employs a bisection process, starting from the reference stepsize (25), until satisfaction of Ineqs.(24) occurs. A procedure useful to determine this reference stepsize is outlined here and is based on a quadratic

representation of the augmented function associated with the one-parameter family of solutions (20).

Let the function $\tilde{F}(\alpha)$ be represented in the quadratic form

$$(28) \qquad \tilde{F}(\alpha) = k_0 + k_1 \alpha + k_2 \alpha^2 \; ,$$

and let the coefficients of the quadratic be determined so as to match the values of the ordinate and the slope at $\alpha = 0$ and the value of the ordinate at $\alpha = 1$. This yields the relations

$$(29) \qquad \tilde{F}(0) = k_0 \; , \quad \tilde{F}_\alpha(0) = k_1 \; , \quad \tilde{F}(1) = k_0 + k_1 + k_2 \; ,$$

which imply that

$$(30) \qquad k_0 = \tilde{F}(0) \; , \quad k_1 = -\tilde{Q}(0) \; , \quad k_2 = \tilde{F}(1) - \tilde{F}(0) + \tilde{Q}(0) \; .$$

With the coefficients known, the following possibilities arise:

$$(31) \qquad (i) \; k_2 > 0 \quad or \quad (ii) \; k_2 < 0 \; .$$

In Case (i), the quadratic function (28) has a minimum for the following value of the gradient stepsize:

$$(32) \qquad \alpha = - k_1/2k_2 \; .$$

In Case (ii), the quadratic function (28) decreases monotonically with $\alpha$. This suggests the use of the following reference values for the gradient stepsize:

$$\alpha_0 = - k_1/2k_2 \quad if \quad k_2 > 0 \; ,$$

$$(33)$$

$$\alpha_0 = 1 \qquad if \quad k_2 < 0 \; .$$

## 7. Experimental Conditions

In order to evaluate the previous algorithms, eight numerical examples were considered. The first two examples pertain to quadratic functions subject to linear constraints. The remaining examples pertain to nonquadratic functions subject to nonlinear constraints. Each example was solved with the three versions of SGRA and the three versions of CGRA

outlined in Section 3.   All of the algorithms were programmed in FORTRAN IV, and the numerical results were obtained using a Burroughs B - 5500 computer and double-precision arithmetic.

Starting Point.   For all of the examples, the nominal point chosen to start an algorithm was defined by

$$(34) \qquad x_1 = x_2 = \ldots = x_n = 2,$$

where $n$ denotes the dimension of the vector $x$.

Search Technique.   The determination of the gradient stepsize and the restoration stepsize was performed in accordance with Sections 5 and 6.   For the gradient phase, the stepsize $\alpha$ was subject to the inequalities

$$(35) \qquad \tilde{P}(\alpha) < \tilde{P}(0), \ \tilde{P}(\alpha) \leq \tilde{P}(0) + 1 .$$

For the restoration phase, the stepsize was subject to the inequality

$$(36) \qquad \tilde{P}(\alpha) < \tilde{P}(0) .$$

Convergence.   Convergence of an algorithm was defined through the inequalities

$$(37) \qquad P(x) \leq 10^{-8}, \ Q(x,\lambda) \leq 10^{-4}.$$

Nonconvergence.   Conversely, nonconvergence of an algorithm was defined by means of the inequalities

$$(38\text{-}1) \qquad (a) \qquad N > 100 ,$$

or

$$(38\text{-}2) \qquad (b) \qquad N_s > 20 ,$$

or

$$(38\text{-}3) \qquad (c) \qquad M > 0.4 \times 10^{69} .$$

Here, $N$ is the iteration number, $N_s$ is the number of bisections of the stepsize $\alpha$ required to satisfy Ineq. (35) or (36), and $M$ is the modulus of any of the quantities employed in the algorithm. Satisfaction of Ineq. (38-1) indicates divergence or extreme slowness of convergence; satisfaction of Ineq. (38-2) indicates extreme smallness of the displacement $dx$ ; and satisfaction of Ineq. (38-3) indicates exponential overflow.   Each of these situations is undesirable.

## 8.   Numerical Examples

In this section, eight numerical examples are described. The first two examples pertain to quadratic functions subject to linear constraints. The remaining examples pertain to nonquadratic functions subject to nonlinear constraints.

Example 8.1.   Consider the problem of minimizing the function

$$(39) \qquad f = (x_1 + x_2)^2 + (x_2 + x_3 - 2)^2 + (x_4 - 1)^2 + (x_5 - 1)^2 ,$$

subject to the constraints

$$(40) \qquad x_1 + 3x_2 = 0 , \ x_3 + x_4 - 2x_5 = 0, \ x_2 - x_5 = 0 .$$

This function admits the relative minimum $f = 4.0930$ at the point defined by

$$(41) \quad x_1 = -0.7674, \ x_2 = 0.2558, \ x_3 = 0.6279, \ x_4 = -0.1162, \ x_5 = 0.2558$$

and

$$(42) \qquad \lambda_1 = 2.0465 , \ \lambda_2 = 2.2325 , \ \lambda_3 = -5.9534 .$$

Example 8.2.   Consider the problem of minimizing the function

$$(43) \qquad f = (4x_1 - x_2)^2 + (x_2 + x_3 - 2)^2 + (x_4 - 1)^2 + (x_5 - 1)^2 ,$$

subject to the constraints

$$(44) \qquad x_1 + 3x_2 = 0 , \ x_3 + x_4 - 2x_5 = 0, \ x_2 - x_5 = 0 .$$

This function admits the relative minimum $f = 5.3266$ at the point defined by

$$(45) \qquad x_1 = -0.9455 \times 10^{-1}, \ x_2 = 0.3151 \times 10^{-1}, \ x_3 = 0.5157 ,$$
$$x_4 = -0.4527, x_5 = 0.3151 \times 10^{-1}$$

and

$$(46) \qquad \lambda_1 = 3.2779, \lambda_2 = 2.9054 , \lambda_3 = -7.7478 .$$

Example  8.3.   Consider the problem of minimizing the function

(47)           $f = (x_1 - 1)^2 + (x_1 - x_2)^2 + (x_2 - x_3)^4$ ,

subject to the constraint

(48)            $x_1(1 + x_2^2) + x_3^4 - 4 - 3\sqrt{2} = 0$ .

This function admits the relative minimum  $f = 0.3256 \times 10^{-1}$   at  the
point defined by

(49)           $x_1 = 1.1048, \ x_2 = 1.1966, \ x_3 = 1.5352$

and

(50)                $\lambda_1 = - 0.1072 \times 10^{-1}$ .


Example  8.4.   Consider the problem of minimizing the function

(51)      $f = (x_1 - 1)^2 + (x_1 - x_2)^2 + (x_3 - 1)^2 + (x_4 - 1)^4 + (x_5 - 1)^4$ ,

subject to the constraints

(52)     $x_1^2 x_4 + \sin(x_4 - x_5) - 2\sqrt{2} = 0, \ x_2 + x_3^4 x_4^2 - 8 - \sqrt{2} = 0.$

This function admits the relative minimum  $f = 0.2415$ at the point  defined by

(53) $x_1 = 1.1661, \ x_2 = 1.1821, \ x_3 = 1.3802, \ x_4 = 1.5060, \ x_5 = 0.6109$

and

(54)            $\lambda_1 = -0.8553 \times 10^{-1}, \ \lambda_2 = -0.3187 \times 10^{-1}.$


Example  8.5.   Consider the problem of minimizing the function

(55)   $f = (x_1 - 1)^2 + (x_1 - x_2)^2 + (x_2 - x_3)^2 + (x_3 - x_4)^4 + (x_4 - x_5)^4$ ,

subject to the constraints

(56) $x_1 + x_2^2 + x_3^3 - 2 - 3\sqrt{2} = 0, \ x_2 - x_3^2 + x_4 + 2 - 2\sqrt{2} = 0, \ x_1 x_5 - 2 = 0$ .

This function admits the relative minimum  $f = 0.7877 \times 10^{-1}$  at  the

point defined by

(57)  $x_1 = 1.1911, \ x_2 = 1.3626, \ x_3 = 1.4728, \ x_4 = 1.6350, \ x_5 = 1.6790$

and

(58)   $\lambda_1 = -0.3882 \times 10^{-1}, \ \lambda_2 = -0.1672 \times 10^{-1}, \ \lambda_3 = -0.2879 \times 10^{-1}.$


Example  8.6.   Consider the problem of minimizing the function

(59)            $f = 0.01(x_1 - 1)^2 + (x_2 - x_1^2)^2$ ,

subject to the inequality constraint

(60)                  $x_1 \le -1$ .

Introduce the auxiliary variable  $x_3$  defined by

(61)               $x_1 + x_3^2 + 1 = 0$ .

Then, the previous problem can be recast as that of minimizing the func-
tion (59) subject to the equality constraint (61).  The function (59)
admits the relative minimum  $f = 0.04$  at the point defined by

(62)            $x_1 = -1, \ x_2 = 1, \ x_3 = 0$

and

(63)                 $\lambda_1 = 0.04$ .


Example  8.7.   Consider the problem of minimizing the function

(64)                 $f = -x_1$ ,

subject to the inequality constraints

(65)             $x_2 \ge x_1^3, \ x_2 \le x_1^2.$

Introduce the auxiliary variables  $x_3$  and  $x_4$  defined by

(66)          $x_2 - x_1^3 - x_3^2 = 0, \ x_1^2 - x_2 - x_4^2 = 0$ .

Then, the previous problem can be recast as that of  minimizing   the

functic subject to the equality constraints (66). The function (64) admits the relative minimum $f = -1$ at the point defined by

(67) $$x_1 = 1, \ x_2 = 1, \ x_3 = 0, \ x_4 = 0$$

and

(68) $$\lambda_1 = -1, \ \lambda_2 = -1.$$

Example 8.8. Consider the problem of minimizing the function

(69) $$f = \log x_3 - x_2,$$

subject to the equality constraint

(70) $$x_2^2 + x_3^2 - 4 = 0$$

and the inequality constraint

(71) $$x_3 \geq 1.$$

Introduce the auxiliary variable $x_1$ defined by

(72) $$x_3 = 1 + x_1^2.$$

Then, the previous problem can be recast as that of minimizing the function

(73) $$f = \log(1 + x_1^2) - x_2,$$

subject to the equality constraint

(74) $$(1 + x_1^2)^2 + x_2^2 - 4 = 0.$$

Note that $x_3$ has been eliminated from the problem and can be computed a posteriori with (72). The function (73) admits the relative minimum $f = -\sqrt{3}$ at the point defined by

(75) $$x_1 = 0, \ x_2 = \sqrt{3}, \ x_3 = 1$$

and

(76) $$\lambda_1 = 1/2\sqrt{3}.$$

## 9. Results and Conclusions

The examples described in Section 8 were solved with the three versions of SGRA and the three versions of CGRA described in Section 3. The numerical results are presented in Tables 2-3, where the number of iterations for convergence $N_e$ is shown. For the eight examples considered, Table 4 shows the cumulative number of iterations for convergences $\Sigma N_e$. From the tables, the following conclusions arise: (a) a restoration of some form is necessary for rapid convergence; and (b) while SGRA-CR is the most stable among the algorithms considered here, rapidity of convergence can be increased somewhat if one employs algorithms with alternate restoration or optional restoration.

Table 2. Number of iterations for convergence $N_e$.

| Example | SGRA-CR | SGRA-IR | SGRA-OR |
|---|---|---|---|
| 8.1 | 5 | 5 | 5 |
| 8.2 | 8 | 8 | 8 |
| 8.3 | 18 | 14 | 16 |
| 8.4 | 56 | 51 | 42 |
| 8.5 | 8 | 7 | 7 |
| 8.6 | 15 | 12 | 16 |
| 8.7 | 9 | 15 | 9 |
| 8.8 | 11 | 11 | 10 |

Table 3. Number of iterations for convergence $N_e$.

| Example | CGRA-NR | CGRA-AR | CGRA-OR |
|---|---|---|---|
| 8.1 | 17 | 5 | 5 |
| 8.2 | 65 | 8 | 8 |
| 8.3 | 22 | 16 | 16 |
| 8.4 | 36 | 54 | 43 |
| 8.5 | 7 | 7 | 7 |
| 8.6 | >100 | 19 | 13 |
| 8.7 | 13 | 7 | 9 |
| 8.8 | 15 | 8 | 10 |

Table 4.   Cumulative number of iterations for convergence  $\Sigma n_\phi$ .

| Algorithm | $\Sigma n_\phi$ |
|-----------|-----------------|
| SGRA-CR | 130 |
| SGRA-1R | 123 |
| SGRA-OR | 113 |
| CGRA-NR | > 275 |
| CGRA-AR | 124 |
| CGRA-OR | 111 |

## REFERENCES

1.  Miele, A., Huang, H.Y., and Heideman, J.C., *Sequential Gradient-Restoration Algorithm for the Minimization of Constrained Functions, Ordinary and Conjugate Gradient Versions*, Journal of Optimization Theory and Applications, Vol.4, No.4, 1969.

2.  Miele, A., and Heideman, J.C., *Mathematical Programming for Constrained Minimal Problems, Part 1, Sequential Gradient-Restoration Algorithm*, Rice University, Aero-Astronautics Report No.59, 1969.

3.  Miele, A., Heideman, J.C., and Levy, A.V., *Mathematical Programming for Constrained Minimal Problems, Part 2, Combined Gradient-Restoration Algorithm*, Rice University, Aero-Astronautics Report No.68, 1970.

# Comparison of Multiplier and Quasilinearization Methods

Alejandro V. Levy

Antonio Montalvo

In the above equations, $f$ is a scalar, $x$ an $n$-vector and $c$ a $q$-vector (all vectors are column vectors) where $q < n$. It is assumed that the first and second partial derivatives of the functions $f$ and $c$ with respect to $x$ exist and are continuous; it is also assumed that the constrained minimum exists.

2.1. *Exact First-Order Conditions.* From theory of maxima and minima, it is known that the previous problem can be recast as that of minimizing the augmented function

$$F(x,\lambda) = f(x) + \lambda^T c(x) \tag{3}$$

subject to the constraint (2). Here, $\lambda$ is a $q$-vector Lagrange multiplier and the superscript $T$ denotes the transpose of a matrix. If

$$F_x(x,\lambda) = f_x(x) + c_x(x)\lambda \tag{4}$$

denotes the gradient of the augmented function (In eq 4, the gradients $f_x$ and $F_x$ denote $n$-vectors and the matrix $c_x$ is $n \times q$), the optimum solution $x,\lambda$ must satisfy the simultaneous equations

$$c(x) = 0; \quad F_x(x,\lambda) = 0 \tag{5}$$

2.2. *Approximate Solutions.* In general, the system (5) is nonlinear; consequently, approximate methods must be employed. These are of two kinds: first-order methods (see, for instance, Miele et al. (1969)) and second-order methods. Here, we introduce the scalar quantities

$$P(x) = c^T(x)c(x); \quad Q(x,\lambda) = F_x^T(x,\lambda)F_x(x,\lambda) \tag{6}$$

which measure the errors in the constraint and the optimum condition, respectively. We observe that $P = 0$ and $Q = 0$ for the optimum solution, while $P > 0$ and/or $Q > 0$ for any approximation to the solution. When approximate methods are used, they must ultimately lead to values of $c,\lambda$ such that

$$P(x) \le \epsilon_1; \quad Q(x,\lambda) \le \epsilon_2 \tag{7}$$

Alternately, (7) can be replaced by

$$R(x,\lambda) \le \epsilon_3 \tag{8}$$

where

$$R(x,\lambda) = P(x) + Q(x,\lambda) \tag{9}$$

denotes the cumulative error in the constraint and the optimum condition. Here, $\epsilon_1,\epsilon_2,\epsilon_3$ are small, preselected numbers. Note that satisfaction of inequality 8 implies satisfaction of inequalities 7, if one chooses $\epsilon_1 = \epsilon_2 = \epsilon_3$.

3. Review of the SQL-Algorithm

Here, a review of SQL is given. Let $x,\lambda$ denote the nominal values, and let $\Delta x, \Delta\lambda$ denote the displacements leading from the nominal values to the varied values $\tilde{x}, \tilde{\lambda}$. With this understanding, SQL can be summarized as follows.

(i) For given $x$ and $\lambda$, compute $f(x)$, $c(x)$, $f_x(x)$, $c_x(x)$.

(ii) Compute

$$F_x(x,\lambda) = f_x(x) + c_x(x)\lambda \tag{10-1}$$

$$Q(x,\lambda) = F_x^T(x,\lambda)F_x(x,\lambda) \tag{10-2}$$

$$P(x) = c^T(x)c(x) \tag{10-3}$$

(iii) If $P(x) \le \epsilon_1$ and $Q(x,\lambda) \le \epsilon_2$, convergence is achieved and the algorithm is stopped; otherwise, go to step (iv).

(iv) Compute $f_{xx}(x)$, $c_{xx}(x)$ and obtain

$$F_{xx}(x,\lambda) = f_{xx}(x) + c_{xx}(x)\lambda \tag{11}$$

(v) Solve the linear system of equations (of order $n + q$)

$$\begin{bmatrix} F_{xx}(x,\lambda) & c_x(x) \\ c_x^T(x) & 0 \end{bmatrix}\begin{bmatrix} \Delta x \\ \Delta\lambda \end{bmatrix} + \begin{bmatrix} F_x(x,\lambda) \\ c(x) \end{bmatrix} = 0 \tag{12}$$

(vi) Obtain the varied point

$$\tilde{x} = x + \Delta x \tag{13-1}$$

$$\tilde{\lambda} = \lambda + \Delta\lambda \tag{13-2}$$

and go back to step (i).

4. Review of the MMM-Algorithm

Here, a review of MMM is given in conjunction with the modified quasilinearization algorithm. Let $x$ denote the nominal point, $\tilde{x}$ the varied point, $\Delta x$ the displacement leading from the nominal point to the varied point, $p$ the search direction, $\rho = \pm 1$ the direction factor, and $\alpha$ the step size. With this understanding, the modified method of the multipliers can be summarized as follows.

(i) Choose a nominal value for $x$ not satisfying the constraint. Compute $f(x)$ and $c(x)$.

(ii) Compute $f_x(x)$, $c_x(x)$ and obtain

$$P(x) = c^T(x)c(x) \tag{14-1}$$

$$P_x(x) = 2c_x(x)c(x) \tag{14-2}$$

$$F_x(x,\lambda_1) = f_x(x) + c_x(x)\lambda_1 \tag{14-3}$$

$$\beta = -P_x^T(x)F_x(x,\lambda_1)/P_x^T(x)P_x(x) \tag{14-4}$$

where $\lambda_1$ is the Lagrange multiplier of the previous iteration.

(iii) Update the Lagrange multiplier to the value $\lambda_2$ given by

$$\lambda_2 = \lambda_1 + 2\beta c(x) \tag{15}$$

(iv) Compute

$$F_x(x,\lambda_2) = f_x(x) + c_x(x)\lambda_2 \tag{16-1}$$

$$Q(x,\lambda_2) = F_x^T(x,\lambda_2)F_x(x,\lambda_2) \tag{16-2}$$

(v) Obtain a tentative penalty constant $k_0$

$$k_0 = |\lambda_2^T c(x)|/P(x) \tag{17-1}$$

and the updated penalty constant $k_2$ as follows

$$k_2 = \min(k_0,k_1) \text{ if } P(x) \le Q(x,\lambda_2) \tag{17-2}$$

$$k_2 = \max(k_0,k_1) \text{ if } P(x) > Q(x,\lambda_2) \tag{17-3}$$

where $k_1$ is the penalty constant of the previous iteration.

(vi) Compute the augmented penalty function and its gradient

$$W(x,\lambda_2,k_2) = f(x) + \lambda_2^T c(x) + k_2 P(x) \tag{18-1}$$

$$W_x(x,\lambda_2,k_2) = f_x(x) + c_x(x)\lambda_2 + k_2 P_x(x) \tag{18-2}$$

(vii) Compute the second derivative matrices $f_{xx}(x)$, $c_{xx}(x)$ and obtain

$$F_{xx}(x,\lambda_2) = f_{xx}(x) + c_{xx}(x)\lambda_2 \tag{19-1}$$

$$P_{xx}(x) = 2[c_{xx}(x)c(x) + c_x(x)c_x^T(x)] \tag{19-2}$$

$$W_{xx}(x,\lambda_2,k_2) = F_{xx}(x,\lambda_2) + k_2 P_{xx}(x) \tag{19-3}$$

(viii) Solve the linear system of equations (of order $n$) and obtain the search direction $p$ from the relation

$$W_{xx}(x,\lambda_2,k_2)p + W_x(x,\lambda_2,k_2) = 0 \tag{20}$$

(ix) Set the step size at the value $\alpha = 1$, set the proper direction factor $\rho$ and the search direction $p$ from the relations

$$\rho = \text{sign} \; [W_1^T(x, \lambda_2, k_2) A] \qquad (21\text{-}1)$$

$$p = \rho_1 \qquad (21\text{-}2)$$

(x) Compute the displacement $\Delta x$, the varied point $\tilde{x}$, and the new value of the augmented penalty function with

$$\Delta x = -\alpha p \qquad (22\text{-}1)$$

$$\tilde{x} = x + \Delta x \qquad (22\text{-}2)$$

$$W(\tilde{x}, \lambda_2, k_2) = f(\tilde{x}) + \lambda_2^T \varphi(\tilde{x}) + k_2 P(\tilde{x}) \qquad (22\text{-}3)$$

(xi) If

$$W(\tilde{x}, \lambda_2, k_2) < W(x, \lambda_2, k_2) \qquad (22\text{-}4)$$

accept the present value of the step size $\alpha$ and go back to step (ii); otherwise, go to step (xii).

(xii) If

$$W(\tilde{x}, \lambda_2, k_2) > W(x, \lambda_2, k_2) \qquad (22\text{-}5)$$

bisect the step size as many times as needed (The bisections are started from $\alpha = 1$.) and go back to step (x) until inequality (22-4) is satisfied.

Remark. For the first iteration, the previous values of $\lambda_1$ and $k_1$ are not known. This being the case, one sets

$$\lambda_1 = 0, k_1 = k_0 \qquad (23)$$

## 5. Operational Count

In this section, we give the number of arithmetic operations performed by each algorithm per iteration. This operational count gives a measure of the computer time needed per iteration. The computer time needed to compute the functions $f(x)$ and $\varphi(x)$ and the derivatives $f_x(x)$, $c_x(x)$, $f_{xx}(x)$, $\varphi_{xx}(x)$ is not considered in the comparison. In the absence of step size bisections, SQL and MMM require the same number of function evaluations per iteration. Therefore, the number of function evaluations is not important for the comparison if the basic criterion is the computer time difference between the two algorithms.

We note that the summation and subtraction operations require much less computer time than the multiplication and division operations. Therefore, the operational count is done only on the basis of the number of multiplications and divisions required to complete one iteration.

5.1. Operational Count for the SQL-Algorithm. Let $\mu_1$ denote the number of multiplications and divisions required by the SQL-algorithm per iteration in order to solve the linear system of equations of order $n + q$ (Isaacson and Keller (1966)). Let $\mu_2$ denote the number of remaining operations. These numbers are given by

$$\mu_1 = (n + q)^3/3 + (n + q)^2 - (n + q)/3 \qquad (24\text{-}1)$$

$$\mu_2 = n^2 q/2 + 3nq/2 + n + q \qquad (24\text{-}2)$$

As a consequence, their sum

$$\mu = \mu_1 + \mu_2 \qquad (25)$$

becomes

$$\mu = (n + q)^3/3 + (n + q)^2 - (n + q)/3 +$$
$$n^2/2 + 3nq/2 + n + q \qquad (26)$$

5.2. Operational Count for the MMM-Algorithm. Let $\mu_1$ denote the number of multiplications and divisions required by the MMM-algorithm per iteration in order to solve a linear system of $W_{xx}$ ... of order $n$ (Isaacson and Keller (1966)). Let $\mu_2$ denote the number of remaining operations ... given by

$$\mu_1 = n^3/3 + n^2 - n/3 \qquad (27)$$

---

Table 1. Critical Value of the Number of Constraints

| $n$ | $q_c$ | $n$ | $q_c$ | $n$ | $q_c$ |
|-----|-------|------|-------|------|-------|
| 2 | 2 | 20 | 5 | 200 | 14 |
| 3 | 3 | 30 | 6 | 300 | 17 |
| 4 | 3 | 40 | 7 | 400 | 20 |
| 5 | 3 | 50 | 7 | 500 | 22 |
| 6 | 3 | 60 | 8 | 600 | 24 |
| 7 | 4 | 70 | 9 | 700 | 26 |
| 8 | 4 | 80 | 9 | 800 | 28 |
| 9 | 4 | 90 | 10 | 900 | 30 |
| 10 | 4 | 100 | 10 | 1000 | 32 |

$$\mu_2 = 3n^2 q/2 + 9nq/2 + 3(n + q) + n^2 + 6n + 2q + 5 \qquad (28)$$

As a consequence, their sum (29) becomes

$$\mu = n^3/3 + 3n^2 q/2 + 2n^2 + 9nq/2 + 26n/3 + 5q + 5 \qquad (29)$$

5.3. Comparison of Operational Counts. Comparison between SQL and MMM shows that the former requires a larger number of operations to solve the linear system and a smaller number of operations to perform the remaining tasks. It is of interest to compute the difference $\Delta$ between the overall number of operations, that is

$$\Delta = (\mu)_{SQL} - (\mu)_{MMM} \qquad (30)$$

From eq 26, 29, and 30, we obtain

$$\Delta = q^3/3 + q^2(n + 1) - (n^2 + nq + 8n + 13q/3 + 5) \qquad (31)$$

Now, let $\Delta = 0$, so that eq 31 becomes

$$q^3/3 + q^2(n + 1) - (n^2 + nq + 8n + 13q/3 + 5) = 0 \qquad (32)$$

For given $n$, let $q_*$ denote the solution of eq 32 and let $q_c$ denote the closest integer (from below) to $q_*$, subject to the limitation $q_c \le n$. Inspection of eq 31 and 32 shows that, for any given $n$

$$\Delta < 0; 0 \le q \le q_c \qquad (33\text{-}1)$$

$$\Delta > 0; q_c + 1 \le q \le n \qquad (33\text{-}2)$$

The values of $q_c$ are given in Table 1 as a function of the number of variables $n$. As an example, consider the case $n = 10$, for which $q_c = 4$. If $q \le 4$, SQL requires a smaller number of operations than MMM. On the other hand, if $q \ge 5$, the opposite is true. From the table, it appears that MMM might become superior to SQL for large systems (large $n$), even when these large systems are moderately constrained.

## 6. Experimental Conditions

In order to illustrate the characteristics of the SQL and MMM algorithms, nine numerical examples were solved using an IBM 370/155 computer and double precision arithmetic. Both algorithms were programmed in Fortran.

Nominal Values. For both algorithms, the nominal values of ... following values of $\alpha$

For the SQL-algorithm, the nominal values of the Lagrange

multipliers were chosen as $\lambda_i = \gamma$, $i = 1,2,3,\ldots,q$, with the following values of $\gamma$

$$\gamma = \pm10, \pm8, \pm6, \pm4, \pm2 \qquad (34\text{-}2)$$

For the MMM-algorithm, one starts with $\lambda = 0$ (section 4). Therefore, 100 runs were made for each problem with the SQL algorithm, and 10 runs were made for each problem with the MMM algorithm.

Stopping Conditions. The convergence criterion for both algorithms was chosen to be

$$P(x) + Q(x,\lambda) \leq 10^{-6} \qquad (35)$$

Conversely, the nonconvergence criteria were chosen as follows:

$$\text{(a)} \quad N \geq 100 \qquad (36\text{-}1)$$

$$\text{(b)} \quad N_b \geq 20 \qquad (36\text{-}2)$$

$$\text{(c)} \quad M \geq 10^{75} \qquad (36\text{-}3)$$

Here, $N$ is the iteration number, $N_b$ is the number of bisections of the stepsize $\alpha$ required to satisfy inequality (22-4), and $M$ is the modulus of any of the quantities employed in the algorithm. Condition (a) implies slow convergence; condition (b) denotes extreme smallness of the displacements and, consequently, slow convergence; and, finally, condition (c) denotes an overflow in the number range of the computer.

Percentage of Success. Each problem was solved several times starting with the nominal values given by eq 34. Let $N_t$ denote the total number of runs made for a given problem, using a given method; let $N_s$ denote the number of runs for which the algorithm succeeded in converging to a relative minimum. Then, the percentage of success $p$ is given by

$$p = N_s/N_t \qquad (37)$$

It is noted that $0 \leq p \leq 1$ and that the higher the value of $p$, the more robust the algorithm is, since it means that it has a larger range of successful convergence.

Actual Computer Time per Run. Let $T_i$ denote the CPU time in seconds employed in the $i$th run to solve a given problem with a given algorithm. For a given problem and a given algorithm, the average CPU time per successful run is given by

$$T_{av} = (1/N_s) \sum_{i=1}^{N_s} T_i \qquad (38)$$

Effective Computer Time per Run. The main characteristics of the algorithms considered in this report, namely, robustness (high $p$) and speed (low $T_{av}$) can be combined in a single parameter by considering the effective CPU time per run, which, for a given problem and a given algorithm, is given by

$$T_{eff} = T_{av}/p \qquad (39)$$

Thus, an algorithm which is fast and has a high percentage of success has a smaller effective computer time than an algorithm which is slow and has a small percentage of success.

Relative Efficiency of Two Algorithms. In order to compare the SQL and MMM class of algorithms and arrive at a significant parameter which is machine independent, we introduce the relative efficiency index

$$E = (T_{eff})_{SQL}/(T_{eff})_{MMM} \qquad (40)$$

Therefore, $E < 1$ implies that SQL is more efficient than MMM, while the opposite is true for $E > 1$.

## 7. Numerical Examples

In this section nine numerical examples are described. Examples 7.1 through 7.7 were considered by Miele et al. (1972), while example 7.8 was considered by Luus and Jaakola (1973). For simplicity, scalar notation is used.

Example 7.1. Consider the problem of minimizing the function

$$f = (x_1 - x_2)^2 + (x_2 + x_3 - 2)^2 + (x_4 - 1)^2 + (x_5 - 1)^2 \qquad (41)$$

subject to the constraints

$$x_1 + 3x_2 = 0; \ x_3 + x_4 - 2x_5 = 0; \ x_2 - x_5 = 0 \qquad (42)$$

The function (41) admits the relative minimum $f = 0.4093 \, E + 01$ at the point defined by

$$x_1 = -0.7674; \ x_2 = 0.2558; \ x_3 = 0.6279;$$
$$x_4 = -0.1162; \ x_5 = 0.2558 \qquad (43)$$
$$\lambda_1 = 0.2046 \, E + 01; \ \lambda_2 = 0.2232 \, E + 01;$$
$$\lambda_3 = -0.3933 \, E + 01 \qquad (44)$$

Example 7.2. Consider the problem of minimizing the function

$$f = (x_1 - 1)^2 + (x_1 - x_2)^2 + (x_2 - x_3)^4 \qquad (45)$$

subject to the constraint

$$x_1(1 + x_2^2) + x_3^4 - 4 - 3\sqrt{2} = 0 \qquad (46)$$

(i) At the point defined by

$$x_1 = 0.1104 \, E + 01; \ x_2 = 0.1196 \, E + 01;$$
$$x_3 = 0.1535 \, E + 01 \qquad (47)$$
$$\lambda_1 = -0.1072 \, E - 01 \qquad (48)$$

the function (45) has the relative minimum $f = 0.3256 \, E - 01$.

(ii) At the point defined by

$$x_1 = 0.9861 \, E = 01; \ x_2 = -0.8954; \ x_3 = -0.1685 \, E + 01 \qquad (49)$$
$$\lambda_1 = -0.3029 \qquad (50)$$

the function (45) has the relative minimum $f = 0.2169 \, E + 01$.

Example 7.3. Consider the problem of minimizing the function

$$f = (x_1 - 1)^2 + (x_1 - x_2)^2 + (x_3 - 1)^2 + (x_4 - 1)^4 + (x_5 - 1)^6 \qquad (51)$$

subject to the constraints

$$x_1^2 x_4 + \sin(x_4 - x_5) - 2\sqrt{2} = 0;$$
$$x_2 + x_3^4 x_4^2 - 8 - \sqrt{2} = 0 \qquad (52)$$

(i) At the point defined by

$$x_1 = 0.1166 \, E + 01;$$
$$x_2 = 0.1182 \, E + 01;$$
$$x_3 = 0.1382 \, E + 01;$$
$$x_4 = 0.1506 \, E + 01; \ x_5 = 0.6102 \qquad (53)$$
$$\lambda_1 = -0.5353 \, E - 01; \ \lambda_2 = -0.3187 \, E - 01 \qquad (54)$$

the function (51) has the relative minimum $f = 0.2415$.

(ii) At the point defined by

$$x_1 = 0.1089\,E + 01;\, x_2 = 0.1177\,E + 01;$$
$$x_3 = -0.1281\,E + 01;\, x_4 = 0.1747\,E + 01;$$
$$x_5 = 0.8913 \quad (55)$$

and

$$\lambda_1 = -0.1489\,E - 03;\, \lambda_2 = -0.1774 \quad (56)$$

the function (51) has the relative minimum $f = 0.5523\,E + 01$.

(iii) At the point defined by

$$x_1 = -0.1025\,E + 01;\, x_2 = -0.1017\,E + 01;$$
$$x_3 = 0.1334\,E + 01;\, x_4 = 0.1760;\, x_5 = 0.4531 \quad (57)$$

and

$$\lambda_1 = -0.1126\,E + 01;\, \lambda_2 = -0.2301\,E - 01 \quad (58)$$

the function (51) has the relative minimum $f = 0.4602\,E + 01$.

(iv) At the point defined by

$$x_1 = -0.9263;\, x_2 = -0.9142;\, x_3 = 0.1302\,E + 01;$$
$$x_4 = 0.1803\,E + 01;\, x_5 = 0.4975 \quad (59)$$

and

$$\lambda_1 = -0.1102\,E + 01;\, \lambda_2 = -0.1452 \quad (60)$$

the function (51) has the relative minimum $f = 0.9908\,E + 01$.

Example 7.4. Consider the problem of minimizing the function

$$f = (x_1 - 1)^2 + (x_1 - x_2)^2 + (x_2 - x_3)^2 +$$
$$(x_3 - x_4)^4 + (x_4 - x_5)^4 \quad (61)$$

subject to the constraints

$$x_1 + x_2^2 + x_3^3 - 2 - 3\sqrt{2} = 0;$$
$$x_2 - x_3^2 + x_4 + 2 - 2\sqrt{2} = 0;\, x_1 x_5 - 2 = 0 \quad (62)$$

(i) At the point defined by

$$x_1 = 0.1191\,E + 01;\, x_2 = 0.1362\,E + 01;$$
$$x_3 = 0.1472\,E + 01;$$
$$x_4 = 0.1635\,E + 01;\, x_5 = 0.1679\,E + 01 \quad (63)$$

and

$$\lambda_1 = -0.5882\,E - 01;\, \lambda_2 = -0.1672\,E - 01;$$
$$\lambda_3 = -0.2873\,E - 03 \quad (64)$$

the function (61) has the relative minimum $f = 0.7877\,E - 01$.

(ii) At the point defined by

$$x_1 = 0.2717\,E + 01;\, x_2 = 0.2033\,E + 01;\, x_3 = -0.8179;$$
$$x_4 = -0.4859;\, x_5 = 0.7359 \quad (65)$$

and

$$\lambda_1 = -0.2825\,E + 01;\, \lambda_2 = 0.7106\,E + 01;$$
$$\lambda_3 = -0.2651\,E + 01 \quad (66)$$

the function (61) has the relative minimum $f = 0.1396\,E + 02$.

(iii) At the point defined by

$$x_1 = -0.7651;\, x_2 = 0.2636\,E + 01;\, x_3 = -0.4681;$$
$$x_4 = -0.5911\,E + 01;\, x_5 = -0.2614\,E + 01 \quad (67)$$

and

$$\lambda_1 = -0.2875\,E + 01;\, \lambda_2 = 0.2203\,E + 01;$$
$$\lambda_3 = -0.5085\,E + 01 \quad (68)$$

the function (61) has the relative minimum $f = 0.2715\,E + 02$.

(iv) At the point defined by

$$x_1 = -0.1246\,E + 01;\, x_2 = 0.2422\,E + 01;$$
$$x_3 = 0.1174\,E + 01;\, x_4 = -0.2132;$$
$$x_5 = -0.1604\,E + 03 \quad (69)$$

and

$$\lambda_1 = -0.2016\,E + 01;\, \lambda_2 = -0.6231\,E - 01;$$
$$\lambda_3 = -0.5632\,E + 01 \quad (70)$$

the function (61) has the relative minimum $f = 0.2752\,E + 02$.

(v) At the point defined by

$$x_1 = 0.9194;\, x_2 = -0.2296\,E + 01;\, x_3 = 0.5377;$$
$$x_4 = 0.3384\,E + 01;\, x_5 = 0.2106\,E + 01 \quad (71)$$

and

$$\lambda_1 = -0.2481\,E + 02;\, \lambda_2 = -0.1006\,E + 03;$$
$$\lambda_3 = 0.8790\,E + 01 \quad (72)$$

the function (61) has the relative minimum $f = 0.5652\,E + 02$.

(vi) At the point defined by

$$x_1 = -0.2702\,E + 01;\, x_2 = -0.2989\,E + 01;\, x_3 = 0.1719;$$
$$x_4 = 0.3847\,E + 01;\, x_5 = -0.7401 \quad (73)$$

and

$$\lambda_1 = -0.9898\,E + 02;\, \lambda_2 = -0.5850\,E + 03;$$
$$\lambda_3 = -0.1429\,E + 03 \quad (74)$$

the function (61) has the relative minimum $f = 0.6495\,E + 03$.

Example 7.5. Consider the problem of minimizing the function

$$f = 0.01(x_1 - 1)^2 + (x_2 - x_1^2)^2 \quad (75)$$

subject to constraint

$$x_1 + x_3^2 + 1 = 0 \quad (76)$$

At the point defined by

$$x_1 = -1;\, x_2 = 1;\, x_3 = 0 \quad (77)$$
$$\lambda_1 = 0.01 \quad (78)$$

this function (75) has the relative minimum $f = 0.04$.

Example 7.6. Consider the problem of minimizing the function

$$f = -x_1 \quad (79)$$

subject to constraints

$$x_2 - x_1^3 - x_3^2 = 0;\, x_1^2 - x_2 - x_4^2 = 0 \quad (80)$$

The function (79) admits the relative minimum $f = -1$ at the point defined by

$$x_1 = 1;\, x_2 = 1;\, x_3 = 0;\, x_4 = 0 \quad (81)$$
$$\lambda_1 = -1;\, \lambda_2 = -1 \quad (82)$$

Example 7.7. Consider the problem of minimizing the function

$$f = \log(1 + x_1^2) - x_2 \quad (83)$$

subject to the constraint

$$(1 + x_1^2)^2 + x_2^2 - 4 = 0 \quad (84)$$

(i) At the point defined by

Table II. Results for the Examples

| Example | $N_r$ | SQL-algorithm $N_s$ | $\sum_{i=1}^{N_s} f_i$ | MMM-algorithm $N_r$ | $N_s$ | $\sum_{i=1}^{N_s} f_i$ |
|---|---|---|---|---|---|---|
| 7.1 | 100 | 100 | 5.60 | 10 | 10 | 1.55 |
| 7.2 | 100 | 93 | 12.92 | 10 | 9 | 1.02 |
| 7.3 | 100 | 91 | 67.65 | 10 | 8 | 5.27 |
| 7.4 | 100 | 76 | 33.32 | 10 | 10 | 3.03 |
| 7.5 | 100 | 100 | 13.59 | 10 | 7 | 1.54 |
| 7.6 | 100 | 90 | 27.79 | 10 | 10 | 3.07 |
| 7.7 | 100 | 51 | 4.45 | 10 | 10 | 1.17 |
| 7.8 | 100 | 19 | 1.57 | 10 | 10 | 1.15 |
| 7.9 | 100 | 22 | 50.72 | 10 | 10 | 27.70 |

Table III. Results for the Examples

| Example | SQL-algorithm $p$ | $T_{av}$ | $T_{eff}$ | MMM-algorithm $p$ | $T_{av}$ | $T_{eff}$ | $E$ |
|---|---|---|---|---|---|---|---|
| 7.1 | 1.00 | 0.058 | 0.058 | 1.00 | 0.133 | 0.133 | 0.43 |
| 7.2 | 0.95 | 0.136 | 0.143 | 0.90 | 0.113 | 0.125 | 1.14 |
| 7.3 | 0.91 | 0.743 | 0.816 | 0.80 | 0.658 | 0.623 | 0.99 |
| 7.4 | 0.76 | 0.427 | 0.547 | 1.00 | 0.303 | 0.303 | 1.80 |
| 7.5 | 1.00 | 0.135 | 0.135 | 0.70 | 0.220 | 0.314 | 0.42 |
| 7.6 | 0.90 | 0.368 | 0.542 | 1.00 | 0.307 | 0.307 | 1.11 |
| 7.7 | 0.51 | 0.087 | 0.170 | 1.00 | 0.117 | 0.117 | 1.45 |
| 7.8 | 0.19 | 0.103 | 0.545 | 1.00 | 0.115 | 0.115 | 4.73 |
| 7.9 | 0.22 | 2.300 | 10.450 | 1.00 | 2.771 | 2.771 | 3.77 |

$$x_1 = 0; \quad x_2 = \sqrt{3} \tag{85}$$

$$\lambda_1 = \sqrt{3}/6 \tag{86}$$

the function (83) has the relative minimum $f = -\sqrt{3}$.

(ii) At the point defined by

$$x_1 = 0; \quad x_2 = \sqrt{3} \tag{87}$$

$$\lambda_1 = -\sqrt{3}/6 \tag{88}$$

the function (83) has the relative minimum $f = \sqrt{3}$.

Example 7.6. Consider the problem of minimizing the function

$$f = -(x_1^2 + x_2^2 + x_3^2) \tag{89}$$

subject to the constraints

$$x_1 + 2x_2 + 3x_3 - 1 = 0; \quad x_1^2 + x_2^2/2 + x_3^2/4 - 4 = 0 \tag{90}$$

(i) At the point defined by

$$x_1 = -0.1305; x_2 = -0.2429 \, E + 01; x_3 = 0.2014 \, E + 01 \tag{91}$$

$$\lambda_1 = 0.5209; \lambda_2 = 0.2403 \, E + 01 \tag{92}$$

the function (89) has the relative minimum $f = -0.9555 \, E + 01$.

(ii) At the point defined by

$$x_1 = 0.1540; x_2 = 0.2022 \, E + 01; x_3 = -0.1460 \, E + 01$$

$$\lambda_1 = -0.4121; \lambda_2 = 0.2311 \, E + 01 \tag{94}$$

the function (89) has the relative minimum $f = -0.6408 \, E + 01$.

Example 7.9. Consider the problem of minimizing the function

$$f = -\sum_{i=1}^{10} x_i^2 \tag{95}$$

subject to the constraints

$$\sum_{i=1}^{10} (x_i/c_{ij}) - 1 = 0; j = 1, 2, \ldots, 8 \tag{96}$$

$$\sum_{i=1}^{10} (x_i^2/a_i^2) - 4 = 0 \tag{97}$$

where the coefficients $c_{ij}$ and $a_i$ are given by

$$c_{ij} = 1; i \neq j$$
$$c_{ij} = 2; i = j \tag{98}$$

and

$$a_i^2 = 1 + 3(i-1)/9; i = 1, 2, \ldots, 10 \tag{99}$$

(i) At the point defined by

$$x_1 = -0.9158 \, E - 01; x_2 = -0.9158 \, E - 01;$$
$$x_3 = -0.9158 \, E - 01; x_4 = -0.9158 \, E - 01;$$
$$x_5 = -0.9158 \, E - 01; x_6 = -0.9158 \, E - 01;$$
$$x_7 = -0.9158 \, E - 01;$$
$$x_8 = -0.9158 \, E - 01; x_9 = -0.1814 \, E + 01;$$
$$x_{10} = 0.3501 \, E + 01; \tag{100}$$

and

$$\lambda_1 = -0.6332; \lambda_2 = -0.2779; \lambda_3 = -0.6489 \, E - 01;$$
$$\lambda_4 = 0.7732 \, E - 01; \lambda_5 = 0.1753; \lambda_6 = 0.2549;$$
$$\lambda_7 = 0.3141; \lambda_8 = 0.3615; \lambda_9 = 0.3579 \, E + 01 \tag{101}$$

the function (95) has the relative minimum $f = -0.1562 \, E + 01$.

(ii) At the point defined by

$$x_1 = 0.1064; x_2 = 0.1064; x_3 = 0.1064;$$
$$x_4 = 0.1064; x_5 = 0.1064; x_6 = 0.1064;$$
$$x_7 = 0.1064; x_8 = 0.1064; x_9 = 0.2843 \, E + 01;$$
$$x_{10} = -0.2642 \, E + 01 \tag{102}$$

and

$$\lambda_1 = 0.7254; \lambda_2 = 0.6156; \lambda_3 = 0.7462 \, E - 01;$$
$$\lambda_4 = -0.8208 \, E - 01; \lambda_5 = -0.2042; \lambda_6 = -0.2914;$$
$$\lambda_7 = -0.3592; \lambda_8 = -0.4134; \lambda_9 = 0.3519 \, E - 01 \tag{103}$$

the function (95) has the relative minimum $f = -0.1516 \, E + 01$.

## 8. Numerical Results and Conclusions

The examples described in section 7 were solved using both the standard quasilinearization algorithm (SQL) and the modified method of multipliers (MMM) in accordance with the experimental conditions outlined in section 6.

Table II shows, for each example and each algorithm, the total number of runs $N_r$, the number of successful runs $N_s$, and the sum of the minima over the successful runs. Table III shows the probability of success $p$, the average time $T_{av}$, the effective time $T_{eff}$, and the ratio of the efficiencies $E$ derived.

On the average for SQL we note in particular for the examples 7.7 and 7.8 that the percentage of successful runs is lower than 50% for MMM. On the other hand, for SQL, the

percentage of success became as low as 19% in example 7.8 and 22% in example 7.9.

On the upper side, MMM achieved a percentage of success of 100% in six examples, while SQL achieved a percentage of success of 100% in two examples (one of which is the obvious linear-quadratic example 7.1).

For the nine examples, the cumulative percentage of success was 90% for MMM and 74% for SQL. From all these data, the higher robustness of MMM is apparent.

(ii) Computer Time per Run. Inspection of the average computer time per run $T_{pr}$ shows that SQL is superior to MMM in five examples and inferior in four examples. However, when one looks at the effective computer time per run $T_{en}$, the situation is just the opposite: MMM is superior to SQL in six examples and inferior in three examples.

(iii) Relative Efficiency Index. In section 6, the relative efficiency index $E$ was introduced as a way of combining percentage of success with average computer time per run, while arriving at a parameter which is machine independent, to some degree. This efficiency index $E$ was defined as follows

$$E = \frac{(T_{en})_{SQL}}{(T_{en})_{MMM}} = \frac{(T_{pr})_{SQL}}{(T_{pr})_{MMM}} \frac{(p)_{MMM}}{(p)_{SQL}} \qquad (104)$$

This relative efficiency index is defined so that $E < 1$ indicates superiority of SQL with respect to MMM, while $E > 1$ indicates inferiority of SQL with respect to MMM. Inspection of Table III shows that $E < 1$ in three examples and $E > 1$ in six examples. Thus, from the examples investigated, we conclude that MMM compares favorably with SQL.

(iv) The results of Tables II and III must be taken with a grain of salt, since computer times are precise only to 20%. This being the case, values of $E$ in the range $0.8 \leq E$ < 1.2 are not significant. In particular, this applies to examples 7.2, 7.3, and 7.6.

If one excludes these examples, then conclusion (iii) must be modified as follows: $E < 0.8$ in two examples and $E > 1.2$ in four examples. Therefore, even accounting for possible imprecision in computer time measurements, MMM compares favorably with SQL.

(v) As shown by eq 31 and 32 and Table I, the relative advantage of MMM with respect to SQL should become more apparent for large systems (large $n$) involving many constraints (large $q$). Example 7.9, which includes $n = 10$ variables and $q = 9$ constraints supports this point of view (see Table III).

## Acknowledgment

## Literature Cited

Hestenes, M. R., J. Optim. Theory Appl., 4 (5), 303–320 (1969).

Isaacson, E., Keller, H. B., "Analysis of Numerical Methods," pp 34–37, Wiley, New York, N. Y., 1966.

Luus, R., Jaakola, H. I., Ind. Eng. Chem., Process Des. Dev., 12, 380–383 (1973).

Miele, A., Huang, H. Y., Heideman, J. L., J. Optim. Theory Appl., 4, (4), 213–343 (1969).

Miele, A., Levy, A. V., Iyer, R. R., Wall, K. H., "Modified Quasilinearization Method for Mathematical Programming Problems and Optimal Control Problems," "Control and Dynamic Systems, Advances in Theory and Applications," Vol. 9, pp 239–307, C. T. Leondes, Ed., Academic Press, New York, N. Y., 1973.

Miele, A., Mosseley, P. E., Levy, A. V., Coggins, G. M., J. Optim. Theory Appl., 10, (1), 1–30 (1972).

Does problem have simple serial structure with low dimensionality at each stage?

Is problem unconstrained?

Are constraints probabilistic?

Are nonlinearities only in objective function?

Are the constraints a linear problem?

No — Yes

Classical calculus
Chap. 2

Gradient and search procedures
Chap. 3

Newton
Ord. Grad
Conj. Grad
Quasi-Newton

Can constraints be replaced by failure costs?

Gradient projection
Sec. 5.3

Linearization must laws
Sec. 5.2

Is problem simple and of very low overall dimensionality?

Failure penalty approach
Sec. 9.5 & 9.6

Penalty techniques combined with gradient and search procedures
Sec. 5.4
Sec. 3, 2.5, 3.3

Are there integer requirements on the variables?

Does problem have simple serial structure with no feedback or recycle?

Stochastic dynamic programming
Sec. 9.7

Discrete max principle
Sec. 8.9

Is problem linear?

Can optimization be expressed as allocation problem?

Linear programming simplex method
Sec. 6.4, 4.10

Integer linear programming
Sec. 6.3 & 6.4

Stochastic elements in problem?

Dynamic programming
Sec. 8.1–8.7

Stochastic elements in constraints?

Ordinary optimization

Transportation Algorithm
Sec. 6.2

Is problem too large for available program?

Decomposition
Sec. 6.5

Chance-constrained programming
Sec. 9.4

Is objective function expressed as an integral?

Does problem have simple serial structure without feedback or recycle?

Start

Optimization of functional

Discrete problem

Are all functions continuous and differentiable?

Are u and x low-dimensional vectors and φ and G very simple functions?

Do linear differential equations result from application of continuous max principle?

Continuous max principle
Sec. 7.5–7.8

Finite difference approximation
Sec. 7.9

Variational problem

Does problem have differential constraints?

Can objective function and integral constraints be expressed as functions of u and t only?

Calculus of variations Euler-Lagrange equations
Sec. 7.1–7.4

Reference: Introduction to Optimization theory
B.S. Gottfried and J. Weisman

# Diseño óptimo mediante computadora y su aplicación a la ingeniería mecánica

*Enrique Chicurel Uziel*

## 8.1 El problema de diseño

El problema de diseño tradicional se refiere a la concepción y determinación de las características físicas de un sistema en forma tal que satisfaga básicamente los requisitos de funcionamiento cumpliendo con ciertas restricciones.

El sistema a diseñar puede ser, entre otros, una estructura, un aparato, una máquina, un sistema de control, un circuito, o bien una combinación de dos o más de los mismos.

Las restricciones se pueden referir a las limitaciones en las características mecánicas o eléctricas de los materiales, a los tamaños y capacidades de los componentes que se encuentran comercialmente, al espacio disponible, al medio ambiente, a consideraciones de seguridad y a muchas otras.

## 8.2 Necesidad de optimizar en diseño

Por lo general, no basta con satisfacer los requisitos de funcionamiento sino que, además, hay que satisfacerlos de acuerdo con algún criterio que depende del objeto que se persigue y la aplicación del sistema, como por ejemplo, el diseño debe ser lo más económico posible, o bien lo más eficiente, o lo más compacto, o lo más seguro. Es decir que, esencialmente, el problema de diseño es un problema de optimización, proceso que requiere una gran cantidad de operaciones aritméticas.

Cabe señalar que un diseño se mejora durante la evolución del producto, es decir, mediante el desarrollo de nuevos prototipos en donde se utilizan las experiencias logradas en la producción y aplicación de los anteriores. Este es un proceso lento y costoso pero inevitable. Sin embargo, las técnicas matemáticas de optimización y la computadora electrónica han hecho posible abreviar y abaratar dicho proceso.

Estas circunstancias, aunadas a la existencia de la fuerte competencia entre fabricantes de los países altamente industrializados ha traído como consecuencia un gran auge en el desarrollo de dichas técnicas.

Existen tres clases de optimización, a saber:

1. Optimización de magnitud
2. Optimización de forma
3. Optimización de configuración

Lo que comúnmente se entiende por optimización se refiere a la primera categoría. Sin embargo, todas son dignas de consideración.

## 8.3 Optimización de magnitud

Para concretar las ideas anteriores y lograr hacer optimizaciones prácticas a la mayor brevedad, echemos una ojeada a unos ejemplos de optimización de magnitud.

### Ejemplo 1

Se desea diseñar una caja de cartón cilíndrica lo más económica en material posible para contener 1000 cm³ de cierto producto. Los anaqueles que se acostumbran utilizar para almacenar dichas cajas tienen una altura de 6 cm entre repisas.

Determinar las proporciones óptimas.

### Solución

*Formulación inicial:*

$A = 2 (\pi rh + \pi r^2)$ — Criterio de optimización. (A, función a minimizar).

$V = \pi r^2 h$ — Requisito de funcionalidad.

$h \leq h_{lim}$ — Limitación.

Se especifican: V, $h_{lim}$

Variables: r, h

Variable restringida: h

Variable libre: r

Para obtener la formulación final, eliminamos la variable libre r

### Formulación final:

$A = 2(\sqrt{\pi V h^\dagger} + V h^{-1})$ Criterio de optimización.

$h \leq h_{lim}$          Limitación.

Se especifican: V y $h_{lim}$

Variable: h

Recurriendo al cálculo se obtiene:

$$\frac{dA}{dh} = \sqrt{\pi V}\, h_c^{-\dagger} - 2V h_c^{-2} = 0$$

$h_c = 1.085V^{\frac{1}{2}} = 10.85\, cm$

puesto que $h_{lim} = 6$ cm

$h_c > h_{lim}$

tanto $h_c$ no es la altura óptima y según podemos en la gráfica de la figura 1

$h_{opt} = h_{lim} = 6$ cm

y por lo tanto

$$r_{opt} = \sqrt{\frac{V}{\pi h_{opt}}} = \sqrt{\frac{1000}{\pi(6)}} = 7.25\ cm$$

$$A_{opt} = 2(\pi r_{opt} h_{opt} + \pi r_{opt}^2) =$$
$$2[\pi(7.25)(6) + \pi(7.25)^2] = 609\ cm^2$$

En la figura 1 apreciamos que cualquier punto en la región de diseño factible satisface el requisito de funcionalidad de la caja: contener 1000 $cm^3$. Sin embargo, sólo cuando h = 6 cm se satisface el criterio de optimización: máxima economía.

Un diseño en donde no se utilizara ningún criterio de optimización podría resultar pésimo; considérese, por ejemplo, la parte izquierda de la región de diseño factible

El método que se ilustra es debido al profesor R. C. Johnson, del Instituto Politécnico de Worcester.[17]

### Ejemplo 2

Una viga de acero simplemente apoyada de 100" de largo, de sección rectangular, con 3" de peralte y 1" de ancho, soporta una carga concentrada de 720 lb en el centro.

1. Calcular $\sigma$ máximo, $\sigma$ = esfuerzo normal a una sección.
2. Calcular $\tau_t$ máximo, $\tau_t$ = esfuerzo cortante transversal.



**Figura 1.** Optimización de caja de cartón cilíndrica.

3. Establecer una relación entre el esfuerzo cortante máximo $\tau$ en función de $\sigma$ y $\tau$, para cualquier punto.
4. Comparar los valores numéricos obtenidos en los incisos 1 y 2 y, a la luz de ello, simplificar la relación que se estableció en el inciso 3.
5. Considerar ahora que las dimensiones de la sección de la viga no han sido determinadas. Utilizando la relación del inciso 4 determinar las dimensiones de la sección, tal que:

>   El costo (peso) sea mínimo
>   $\tau$ no exceda a 8000 lb/in²
>   El peralte $h$ no exceda a 3"
>   La deflexión máxima no exceda 0.45"

Proceder de la siguiente manera:

   a) Obtener la formulación inicial completa.
   b) A partir de la inicial obtener la formulación final completa.
   c) Dibujar a escala la región de diseño factible.
   d) Calcular los valores óptimos de: el área seccional (índice del costo), las dimensiones, y el esfuerzo $\tau$.

6. Si se exige ahora que la deflexión no exceda a 0.2",

   a) Dibujar la nueva región de diseño factible.
   b) Recalcular los valores requeridos en el inciso 5d.

## olución

### Nomenclatura

$A$ = área seccional de la viga.
$a$ = área seccional de la viga arriba (o abajo) del punto para el cual se calcula $\tau$.
$b$ = ancho de la viga.
$c$ = distancia medida desde el eje neutro al punto más alejado del mismo.
$\Delta$ = deflexión máxima de la viga.
$E$ = módulo de Young.
$h$ = peralte de la viga.
$I$ = momento de inercia del área seccional de la viga.
$L$ = longitud de la viga.
$M$ = momento flexionante.
$P$ = carga transversal sobre la viga.
$Q$ = primer momento del área "a" respecto al eje neutro
$V$ = fuerza cortante.

$\bar{y}$ = distancia centroidal del área a, medida desde el eje neutro.
$(\ )_L$ = valor límite.

1. $\sigma_{max} = \dfrac{M_{max}\,c}{I} \qquad c = \dfrac{h}{2} \qquad I = \dfrac{bh^3}{12}$

$M_{max} = \dfrac{PL}{4} = \dfrac{(720)(70)}{4} = 12\,600\ \text{lb-in.}$

$\therefore\ \sigma_{max} = 6\,\dfrac{M_{max}}{bh^2} = \dfrac{6\,(16\,000)}{1\,h\,5} = 12\,200\ \text{lb in}^2$

2. $\tau_{max} = \dfrac{V_{max}\,Q_{max}}{I\,b}$

$V_{max} = \dfrac{P}{2}$

$Q_{max} = a\,\bar{y} = \left(\dfrac{bh}{2}\right)\dfrac{h}{4} = \dfrac{bh^2}{8}$

$\tau_{max} = \dfrac{3}{4}\dfrac{P}{bh} = \dfrac{3\,(72)}{4\,(1)\,3} = 250\ \text{in}^2$



**Figura 2.** Ejemplo 2. Área "a" y su distancia centroidal.

3. Para cualquier punto, mediante el círculo de Mohr se obtiene

$$\tau = \sqrt{\left(\dfrac{\sigma}{2}\right)^2 + \tau^2}$$



**Figura 3.** Ejemplo 2. Círculo de Mohr.

4. $\tau_1 <<< \sigma$ y además

cuando $y = c$ $\quad \sigma$ es máximo y $\tau_1 = 0$

cuando $y = 0$ $\quad \sigma = 0$ $\quad$ y $\tau_1$ es máximo

$\therefore \tau_1$ se puede considerar nulo y

$$\tau = \frac{\sigma}{2}$$

5. Consideraciones preliminares:

$$\tau = \frac{\sigma}{2} = \frac{3M}{bh^2}$$

$$\Delta = \frac{PL^3}{48EI} = \frac{ML^2}{Ebh^3}$$

### Formulación inicial

$A = bh$ Criterio de optimización
$\quad\quad$ (A, función a minimizar)

$$\left. \begin{aligned} \tau &= 3\frac{M}{bh^2} \\ \Delta &= \frac{ML^2}{Ebh^3} \end{aligned} \right\} \text{Requisitos de funcionalidad}$$

$$\left. \begin{aligned} \tau &\leq \tau_L \\ \Delta &\leq \Delta_L \\ h &\leq h_L \end{aligned} \right\} \text{Limitaciones}$$

Se especifican: E, M, L, $\tau_L$, $\Delta_L$, $h_L$.

Variables: b, h, $\tau$, $\Delta$

Variables libres: b

Nótese que *las variables libres son sencillamente las que no se indican como limitadas.*

Eliminando las variables libres se obtiene la:

### Formulación final

$A = \dfrac{3\,M}{\tau h}$ Criterio de optimización (A, función a minimizar)

$= \dfrac{L^2\tau}{3\,Eh}$ Requisito de funcionalidad.

$$\left. \begin{aligned} \tau &\leq \tau_L \\ \Delta &\leq \Delta_L \\ h &\leq h_L \end{aligned} \right\} \text{Limitaciones}$$

Se especifican: los parámetros E, M, L y los valores límites $\tau_L$, $\Delta_L$, $h_L$.

Variables independientes: $\tau$, h

Variable dependiente: $\Delta$

Observamos que, a diferencia del ejemplo 1, ahora tenemos dos variables independientes y que, por lo tanto, tendremos una región de diseño factible en dos dimensiones.

Determinemos pues las fronteras de dicha región.

Las fronteras referentes a las variables independientes quedan definidas inmediatamente por las limitaciones, de donde se obtienen sus ecuaciones:

$$\tau = \tau_L = 8\,000$$
$$h = h_L = 3$$

La frontera referente a la variable dependiente se obtiene sustituyendo la limitación correspondiente en la ecuación referente al requisito de funcionalidad

$$\tau = \frac{3\,E\Delta_L}{L^2}\,h = \frac{3(30 \times 10^6)(0.45)}{(100)^2}\,h$$
$$\tau = 4\,060\,h$$

La región de diseño factible se muestra en la figura 4.



Figura 4. Ejemplo 2. Región de diseño factible para $\Delta_L = 0.45''$.

Puesto que la función objetiva es

$$A = \frac{3 M}{r h} - \frac{54\,000}{r h}$$

A es mínimo cuando h y r son ambas máximas dentro de la región de diseño factible o sea h = 3", r = 8000 lh/in².

$$\therefore A_{\text{optimo}} = \frac{54\,000}{3(8\,000)} = 2.23 \text{ in}^2$$

y el valor correspondiente de b es:

$$b = \frac{A}{h} = \frac{2.23}{3} = 0.75"$$

6. Para $\Delta_r = 0.2"$

$$r = 1\,800\,h$$

y la nueva región de diseño factible se muestra en la figura 5.



**Figura 5.** Ejemplo 2. Región de diseño factible para $\Delta_r = 0.2"$.

Una vez más el óptimo se obtiene de la relación:

$$A_{\text{opt}} = \frac{54\,000}{r_{\text{MAX}}\,h_{\text{MAX}}}$$

como se aprecia en la figura 5

$$h_{\text{MAX}} = 3"$$

y $r_{\text{MAX}}$ se obtiene de la intersección de las fronteras:

$$r = 1\,800\,h$$
$$h = h_L = 3"$$

es decir

$$r_{\text{MAX}} = 1\,800\,h_L = 5\,400 \text{ lb/in}^2$$

$$\therefore A_{\text{opt}} = \frac{54\,000}{(5\,400)\,(3)} = 3.33"$$

$$b_{\text{opt}} = \frac{A}{h} = \frac{3.33}{3} = 1.11"$$

Resumiendo, los valores óptimos son:

| $\Delta_r$, in | A, in² | h, in | b, in | $\Delta$, in | r, lb/in² |
|---|---|---|---|---|---|
| 0.45" | 2.23 | 3 | 0.75 | 0.296 | 8 000 |
| 0.20" | 3.33 | 3 | 1.11 | 0.2 | 5 400 |

### 8.4 Consideraciones algebraicas

Si en los anteriores ejemplos se consideraron exclusivamente los requisitos de funcionalidad, se tiene un sistema con un número mayor de variables que de ecuaciones.

A continuación se resumen estas cantidades en relación a los dos primeros ejemplos.

| Ejemplo | Formulación | N° de variables | N° de ecuaciones requisito de funcionalidad |
|---|---|---|---|
| 1 | Inicial | 2 | 1 |
| 1 | Final | 1 | 0 |
| 2 | Inicial | 4 | 2 |
| 2 | Final | 3 | 1 |

Lo anterior quiere decir que existe un número infinito de soluciones de dicho sistema.

Si se considera ahora el conjunto de las ecuaciones referentes a los requisitos de funcionalidad y las desigualdades referentes a las limitaciones, se tiene que se han reducido los valores que pueden asumir las variables pero el número de soluciones sigue siendo infinito.

En diseño tradicional, no se define ningún criterio para seleccionar una sola de las soluciones posibles.

En cambio, en diseño óptimo sí se fija con toda precisión dicho criterio. El proceso de optimización consiste en la búsqueda de dicha solución.

### Ejemplo 3

Se requiere diseñar un resorte helicoidal para un convertidor de par. Ya ha sido seleccionada la cons-

*óptimo mediante computadora*

...resorte, así como la fuerza máxima de

...a montado en una flecha cuyo diámetro
...: lo que el diámetro interior $D_i$ queda
...las por consideraciones de espacio existe
... no debe exceder el diámetro exterior

...rminó el material, por lo cual ya queda-
...de el esfuerzo permisible $\tau_p$ y el módulo
... Sin embargo, el alambre viene única-
...diferentes diámetros $d_i$.
...acer mínima la longitud cerrada del re-
...do todas las espiras están en contacto)
... a la carga máxima Q.
...nes que se emplean en el diseño de re-
...ales son:

$$\frac{G d^4}{8 N D^3_m}$$

$$\frac{Q D_n W}{\pi d^3}$$

$$\frac{4 D_m - d}{4 D_m - d_1} - 0.615 \frac{d}{D_m}$$

...retro del alambre.
...retro medio de la espira.
...es de espiras.
...or de concentración de esfuerzo de Wahl.

...la formulación inicial.

...de la formulación inicial obtener la final.

...dos los siguientes valores:

...88 lb/in

...31½ lb

...2 in

...4 in

...40.000 lb/in²

...12 × 10⁶ lb/in²

$$\frac{1}{32}, \frac{1}{16}, \frac{3}{32}, \dots, 1 \text{ in}$$

*...ción inicial*

Criterio de optimización
($L_c$ función a minimizar)

...e

---

$$k = \frac{G d^4}{8 D'_m N}$$

$$\tau = \frac{8 Q D_m W}{\pi d^3}$$

$$W = \frac{4 D_m - d}{4(D_m - d)} + 0.615 \frac{d}{D_m}$$

$$D^*_m = \frac{D_i + D_e}{2}$$

} Requisitos de funcionalidad

$$D_e = D_i + 2d$$
$$D_i \geqq D_f$$
$$\tau \leq \tau_p$$
$$D_e \leq D_{e\,max}$$
$$d = d_1, d_2, \dots d_M$$

} limitaciones

Se especifican: $k$, $Q$, $G$, $D_f$, $\tau_p$, $D_{e\,max}$

Variables: $N$, $d$, $\tau$, $D_m$, $W$, $D_i$, $D_e$

Variables libres: $N$, $W$, $D_m$.

Eliminando las variables libres se obtiene:

*Formulación final*

$$L_c = \frac{G d^4}{8 k (D_i + d)^3}$$

Criterio de optimización.
$L_c$, función de minimizar.

$$\tau = \frac{4 Q (D_i + d)}{\pi d^3} \left[ \frac{4 D_i + 3d}{2 D_i} + 1.23 \frac{d}{D_i + d} \right]$$

} Requisitos de funcionalidad

$$D_e = D_i + 2d$$

$$D_i \geq D_f$$
$$\tau \leq \tau_p$$
$$D_e \leq D_{e\,max}$$
$$d = d_1, d_2, d_3, \dots d_M$$

} Limitaciones

Se especifican los parámetros: $k$, $Q$, $G$, así como los
valores límite: $D_f$, $\tau_p$, $D_{e\,max}$, $d_1, d_2, \dots d_M$

Variables independientes: $D_i$, $d$

Variables dependientes: $\tau$, $D_e$

Substituyendo valores numéricos obtenemos:

*Formulación final*

$$L_c = 3.84 \times 10^4 \frac{d^4}{(D_i + d)^3}$$

Criterio de optimización
($L_c$ función a optimizar)

$$\tau = 531 \frac{(D_1 + d)}{d^3}\left[\frac{4D_1 + 3d}{2D_1} + \frac{1.23d}{D_1 + d}\right]$$

$$D_e = D_1 + 2d$$

$$D_1 \geq 2$$

$$\tau \leq 4 \times 10^4$$

$$D_e \leq 4$$

$$d = \frac{1}{32}, \frac{1}{16}, \frac{3}{32}, \cdots 1$$

} Requisitos de funcionalidad

} Limitaciones

Variables independientes: $D_1$, d

Variables dependientes: $\tau$, $D_e$.

Determinamos primero la región de diseño factible en el sistema de coordenadas $D_1$, d. Se obtiene una frontera de la limitación en $D_1$ y es la recta

$$D_1 = 2$$

Nótese que por referirse a una variable independiente esta frontera es constante.

Las otras fronteras se obtienen de substituir las limitaciones de las variables dependientes en las ecuaciones referentes a los requisitos de funcionalidad.

$$4 = D_1 + 2d$$

$$4 \times 10^4 = 531 \frac{(D_1 + d)}{d^3}\left(\frac{4D_1 + 3d}{2D_1} + \frac{1.23d}{D_1 + d}\right)$$

Como estas fronteras se refieren a las variables dependientes, no son constantes, y, en general, son curvas. Las fronteras se muestran en la figura 6. Analizando las limitaciones se concluye que la región de diseño factible es la que se muestra encerrada por las líneas sólidas.

Puesto que la función a optimizar ya no es tan sencilla, no se puede determinar el punto óptimo mediante una mera inspección de la ecuación referente al criterio de optimización.

Si no se tiene idea de optimización, lo primero que se ocurre es obtener los contornos de nivel de $L_e$. La ecuación se obtiene de la relación correspondiente al criterio de optimización

$$D_1 = \left(3.84 \times 10^4 \frac{d^3}{L_e}\right)^{\frac{1}{4}} - d$$

En la figura 6 se muestran cuatro contornos correspondientes a los valores de $L_e = 40, 70, 150, 400$ plg.

Si imaginamos a la región de diseño factible llena de contornos, llegamos a la conclusión que el óptimo



**Figura 6.** Ejemplo 3. Región de diseño factible y contornos de nivel de $L_e$.

se encuentra en el extremo derecho inferior de la región de diseño factible, es decir, $d = 0.48$ plg.

Pero consideremos ahora a d como realmente es: una variable discreta. Esto significa que la región de diseño factible consta realmente de los segmentos de recta que se muestran en la figura 7.

Como los valores de $L_e$ suben más lentamente a lo largo de la frontera inferior que a lo largo de la frontera derecha y como

$$\frac{15''}{32} < d < \frac{1''}{2}$$

el óptimo debe quedar en el extremo derecho del segmento inferior, por lo que las proporciones óptimas son:

$$d = 15/32''$$

$$D_1 = 2.7''$$

$$D_e = 3.6''$$

$$L_e = 28''$$

calculadas a partir de la formulación final

**Figura 7.** Ejemplo 3. Región de diseño factible tomando en cuenta los valores discretos de d.

## 8.5   Métodos de exploración local

La solución anterior, aunque muy ilustrativa, es sumamente ineficiente pues requiere de una gran cantidad de cómputos: por un lado, para determinar las fronteras (de hecho, la frontera inferior impuesta por la limitación en τ se calculó por tanteos); y, por otro lado, para determinar suficientes curvas de nivel para obtener una idea del comportamiento de la función a optimizar. El método sería muy difícil y laborioso si tuviéramos que lidiar con tres variables independientes y prácticamente imposible si fueran más.

Todas estas dificultades resultan del hecho que se trató de examinar la topografía de la región completa.

Hay métodos mucho más eficientes. Para entender la idea fundamental recurramos a una analogía.

Baldomero y Agripina juegan. Baldomero se halla con los ojos vendados en un punto en el interior de un predio grande, cercado, en una región montañosa y árida. Agripina le pide que, sin quitarse la venda, encuentre el punto más bajo del predio.

Baldomero se sienta en el suelo y con sus manos palpa el terreno en derredor suyo y escoge la dirección a lo largo de la cual parece descender más el terreno. Camina un trecho como de dos metros y se vuelve a sentar para palpar su derredor una vez más, repitiendo el proceso varias veces hasta encontrar un punto tal que en todas direcciones sube el terreno o bien se topa con cerca. Eufórico, le grita a Agripina que ya encontró

el lugar señalado. Agripina le responde que no es cierto, pero lo conduce a otro punto inicial para darle una oportunidad más. Baldomero procede una segunda vez de idéntica manera hasta llegar a otro punto desde donde pregunta a Agripina si ése es el más bajo del predio. Agripina le responde afirmativamente y premia a Baldomero con un chocolate.

Evidentemente, en su primer intento, Baldomero llegó a un mínimo relativo. Está claro que, dadas suficientes oportunidades, Baldomero siempre podrá encontrar el mínimo absoluto, y esto, sin haber visto nunca el terreno.

Baldomero empleó un método que podríamos llamar de exploración local.

Una técnica muy útil de exploración local en optimización es el método de Box.[3]

Hemos utilizado una terminología propia de diseño pero que difiere de la utilizada por los especialistas en optimización. Así es que, primero que nada, presentamos un pequeño glosario de términos equivalentes.

| *Diseño* | *Optimización* |
|---|---|
| Función a optimizar | = Función objetivo |
| Requisitos de funcionalidad | = Restricciones de igualdad |
| Limitaciones | = Restricciones de desigualdad |
| Variables independientes | = Variables de decisión |
| Limitaciones a las variables independientes | = Restricciones explícitas |
| Limitaciones a las variables dependientes | = Restricciones implícitas |

Examinemos ahora el método de Box para dos variables de decisión. Supongamos concretamente que la optimización se refiere a minimización.

Se escogen cuatro puntos que satisfagan todas las restricciones, es decir, que se encuentren dentro de la región factible. Se aísla el punto más alto (el de valor más alto de la función objetivo). Se encuentra el centroide de los tres puntos restantes y se "refleja" el punto aislado a través de dicho centroide.

El punto reflejado substituye ahora al punto inicialmente aislado. De los cuatro puntos que quedan se vuelve a aislar el más alto y se repite todo el proceso anterior. De esta manera el conjunto de puntos o "Simplex" se va desplazando hacia abajo, substituyendo un punto del conjunto con cada desplazamiento. Si un punto reflejado viola una restricción explícita, se regresa a la frontera; es decir, que, a la variable que excedió su limitación, se le reasigna su valor limitativo. Si un punto reflejado viola una restricción implícita, se regresa medio camino.

# DISEÑO

## OPTIMO

## DE

## FILTROS

## DIGITALES

Horacio Martínez C.

A.bril, 82

## Procesamiento digital de señales [1, 2]

● Representación de señales por medio de secuencias de números y el efectuar algunas transformaciones sobre ellas para obtener algún resultado deseado.

Ejs: interpolación
integración
diferenciación
separar distintas bandas de frecuencia
quitar ruido
estimar parámetros
estimación de espectro
identificación de sistemas

Aplicaciones en: Comunicaciones
señales de voz
señales de audio
bioingeniería
simulación
señales sísmicas
radar
sonar

señal continua: definida para todos los valores de la variable independiente.



señal discreta: definida para valores discretos de la variable independiente. Generalmente el espaciamiento es uniforme.

Señal digital : tiempo y amplitud discretas

Señal analógica : tiempo y amplitud continuos

Teorema del muestreo: $\dfrac{1}{\Delta t} > 2 f_{max}$

Sistema discreto



$$y(n) = \emptyset \, [x(n)]$$

Propiedades : Causalidad
           Linealidad
           Invariancia en el tiempo
           Estabilidad.

La salida de un sistema lineal e invariante
en el tiempo es

$$y(n) = x(n) * h(n)$$

$$= \sum_{m=0}^{n} x(n) \, h(n-m)$$

$h(\cdot) \triangleq$ respuesta a pulso

Relación entrada salida de los sistemas que
vamos a estudiar

$$y(n) = -\sum_{j=1}^{n} b_j \, y(n-i) + \sum_{i=1}^{N} a_i \, x(n-i) \Rightarrow \underline{\text{filtro digital}}$$

Respuesta en frecuencia

$$y(n) = \sum_{m=0}^{n} h(m) x(n-m) \quad \text{si la entrada es}$$
$$x(n) = e^{j\omega n}$$

$$y(n) = \sum_{m=0}^{n} h(m) \, e^{j\omega(m-n)} = e^{j\omega n} \sum_{m=0}^{n} h(m) e^{-j\omega m}$$

$$H(e^{j\omega}) \triangleq \sum_{m=0}^{n} h(m) e^{-j\omega m}$$

entrada senoidal $\longrightarrow$ salida sonoidal misma frecuencia con magnitud $|H(e^{j\omega})|$ y ángulo $\underline{/H(e^{j\omega})}$

$$\text{Retraso de grupo} \triangleq -\frac{d\, \underline{/H(e^{j\omega})}}{d\omega}$$

Transformada $z$

$$X(z) = \sum_{n=0}^{\infty} x(n) z^{-n} \qquad x(n) = \frac{1}{2\pi j} \oint X(z) z^{n-1} dz$$

Propiedades:

$$z[x(n)] = X(z) \quad \Rightarrow \quad z[x(n-k)] = z^{-k} X(z)$$

$$z[y(n)] = z[h(m) * x(m)] \quad \Rightarrow \quad Y(z) = H(z) X(z)$$

$H(z)$ función de transferencia

$$y(n) = -\sum_{i=1}^{N} b_i\, y(n-i) + \sum_{i=1}^{m} a_i\, x(n-i)$$

$$= \left(-\sum_{i=1}^{N} b_i\, z^{-i}\right) Y(z) + \left(\sum_{i=1}^{m} a_i\, z^{-i}\right) X(z)$$

$$\frac{Y(3)}{X(3)} = \frac{\sum_{i=0}^{M} a_i 3^{-i}}{\sum_{i=0}^{N} b_i 3^{-i}}$$

$\longleftarrow$ ceros

$\longleftarrow$ polos

Filtros **R**espuesta a **I**mpulso de duración **F**inita

(RIF)

$$H(3) = \sum_{n=0}^{L} h(n) 3^{-n}$$

Filtros **R**espuesta a **I**mpulso de duración **I**nfinita

(RII)

$$H(3) = \frac{\sum_{i=0}^{M} a_i 3^{-i}}{\sum_{i=0}^{N} b_i 3^{-i}}$$

# Problema de aproximación [3,4]

$f(x)$    función aproximada

$F(A,x)$    función aproximante

$A$    vector de parametros

$P$    espacio de los parametros $A$    $A \subset P$

$\rho$    función de distancia

cumple con las sigs. propiedades

$$\rho[g(x)] \geqslant 0 \quad y \quad \rho[g(x)] = 0 \quad s \cdot y \cdot s \quad g(x) = 0$$

$$\rho[c \, g(x)] = |c| \, \rho[g(x)] \quad \text{para cualquier } c \text{ real}$$

$$\rho[g(x) + h(x)] \leq \rho[g(x)] + \rho[h(x)]$$

$X$    intervalo de aproximación

≡ Sea $f(x)$ una función real continua en $X$, y $F(A,x)$ una función real continua en $X$ que depende de $n$ parámetros $A$. Dada la función $\rho$ determinar los parámetros $A^* \in P$ tales que

$$\rho[F(A^*,x), f(x)] \leq \rho[F(A,x), f(x)]$$

$$\forall \; A \in P$$

# Normas $L_p$

$$L_p(f) = \left[ \int_0^1 |f(x)|^p \, dx \right]^{1/p} \qquad p \geq 1$$



Comportamiento
de
$X^p$

$p=1$

$p=2$

$p \to \infty$

Diseño de filtras en el dominio de la frecuencia



$\dfrac{\pi}{2}$

$\omega$

$$L_\infty = \min_{A \in P} \ \max_{x \in X} \ \rho \left[ F(A,x), f(x) \right]$$

# Filtros RIF [5]

$$H(z) = \sum_{n=0}^{L-1} h(n) z^{-n}$$

L impar

$$h(\ell) = h(L-1-\ell)$$



0 1 2 3 4 5 6

$$H(e^{j\omega}) = \sum_{n=0}^{\frac{L-1}{2}} a(n) \cos \omega n$$

## Teorema de caracterización:

Si $P(\omega)$ es una función de la forma $\sum_{n=0}^{r} a(n) \cos(\omega n)$
una condición necesaria y suficiente para que $P(\omega)$
sea la mejor aproximación (en la norma $\infty$) en el in-
tervalo $X \in [0, \pi]$ ($X$ es compacto) a una función
$D(\omega)$ continua en $X$ es que el error tenga por
lo menos $(r+2)$ valores extremos en $X$, es decir,
debe Haber por lo menos $r+2$ puntos $\omega_i$
que cumplan con las sigs. condiciones

$$0 \leq \omega_1 < \omega_2 \ldots\ldots < \omega_{r+1} < \pi$$

$$\omega_i \in X$$

y $\quad E(\omega_i) = [D(\omega_i) - P(\omega_i)] = \max_{\omega \in X} |E(\omega)|$

$$E(\omega_i) \cdot E(\omega_{i+1}) < 1 \quad i = 0, 1, \ldots r+1$$

Ejemplo

$L = 9$

6 extremos



Algoritmo de Remez (2°)



valores iniciales de frecuencias extremas $\omega_i$

se resuelve el problema
$E(\omega_i) = \pm \delta$
$\omega_i = 0, 1, \dots, r+1$

se encuentran extremos locales de la función de error $\hat{\omega}_i$

$\hat{\omega}_i \overset{?}{=} \omega_i$

$\omega_i = \hat{\omega}_i$

NO

Si

Aprox. Optima

$$E(\omega_i) = D(\omega_2) - P(\omega_i) = \pm \rho$$

$$
\begin{bmatrix}
1 & \cos \omega_0 & \cos 2\omega_0 & \cdots & \cos r\omega_0 & 1 \\
1 & \cos \omega_1 & \cos 2\omega_1 & \cdots & \cos r\omega_1 & -1 \\
\vdots & & & & & \\
\vdots & & & & & \\
\end{bmatrix}
\begin{bmatrix}
a_0 \\
a_1 \\
\vdots \\
a_r \\
\rho
\end{bmatrix}
=
\begin{bmatrix}
D(\omega_0) \\
D(\omega_1) \\
\vdots \\
D(\omega_{r+1})
\end{bmatrix}
$$

# Filtros  R I I

$$H(z) = \frac{\displaystyle\sum_{k=0}^{M} \tilde{a}_k z^{-k}}{1 + \displaystyle\sum_{k=1}^{N} \tilde{b}_k z^{-k}} = K \frac{\displaystyle\prod_{i=1}^{M}(z - z_i)}{\displaystyle\prod_{i}(z - P_i)}$$

Magnitud al cuadrado.

$$H(z)H(z^{-1}) = \frac{c(0) + \displaystyle\sum_{k=1}^{M} c(k)(z^k + z^{-k})}{d(0) + \displaystyle\sum_{k=1}^{N} d(k)(z^k + z^{-k})}$$

$$c(i) = \sum_{k=0}^{M-i} a(k)\, a(k+i)$$

$$d(i) = \sum_{k=0}^{N-i} b(k)\, b(k+i)$$

$$|H(e^{j2\pi f})|^2 = \frac{c(0) + 2\sum_{k=1}^{M} c(k)\cos(2\pi f k)}{d(0) + 2\sum_{k=1}^{N} d(k)\cos(2\pi f k)}$$

$$= \frac{\hat{c}(0) + \sum_{k=1}^{M} \hat{c}(k)\cos^k(2\pi f)}{1 + \sum_{k=1}^{N} \hat{a}(k)\cos^k(2\pi f)}$$

Minimización de normas $L_p$ [6]

la norma $L_p$ $\left\{ \int_{X} |F(A,\omega) - f(\omega)|^p \, d\omega \right\}^{1/p}$

Se aproxima por

$$L_{2p} = \sum_{k=1}^{K} |(F(A, \omega_k) - f(\omega_k)|^{2p}$$

Algoritmo de Polya

$$\lim_{p \to \infty} A_p^* = A^*$$

donde $F(A^*, X)$ es la solución óptima en la
norma de Chebychou.

Forma de la función de transferencia

$$H(z) = b_0 \prod_{i=1}^{N} \frac{z^2 + a_{1i} z + a_{2i}}{z^2 + b_{1i} z + b_{2i}}$$

evaluando en el círculo unitario

$$\alpha(A,\omega) = H(e^{j\omega}) = b_0 \prod_{i=1}^{N} \frac{(1+a_{2i})\cos\omega + a_{1i} + j(1-a_{2i})\sin\omega}{(1+b_{2i})\cos\omega + b_{1i} + j(1-b_{2i})\sin\omega}$$

$$\tau(A,\omega) = \sum_{i=1}^{N} \left[ Re\left( \frac{2\cos\omega + b_{1i} + j\, 2\sin\omega}{(1+b_{2i})\cos\omega + b_{1i} + j(1-b_{2i})\sin\omega} \right) \right.$$

$$\left. - Re\left( \frac{2\cos\omega + a_{1i} + j\, 2\sin\omega}{(1+a_{2i})\cos\omega + a_{1i} + j(1-a_{2i})\sin\omega} \right) \right]$$

más conveniente usar coordenadas polares

Se puede aproximar simultaneamente magnitud
y fase

p. ej. $\quad L_{2p,2q}^{\alpha,\tau}(A,\tau_0) = \delta \sum_{k=1}^{K} (\alpha(A,\omega_k) - \alpha_d(\omega_k))^{2p}$

$$+ (1-\delta) \sum_{r=1}^{S} (\tau(A,\omega_r) - \tau_d(\omega_r) - \tau_0)^{2q} +$$

Uso de programación lineal [7]

Se trabaja con la magnitud al cuadrado

$$\left[c_0 + \sum_{i=1}^{n} 2c_i \cos(\omega i)\right] \Big/ \left[d_0 + \sum_{i=1}^{n} 2d_i \cos(\omega i)\right]$$

$$= \hat{N}(\omega) \,/\, \hat{D}(\omega)$$

$F(\omega)$ función deseada

$$-\epsilon(\omega) \le \frac{\hat{N}(\omega)}{\hat{D}(\omega)} - F(\omega) \le \epsilon(\omega)$$

$\epsilon(\omega)$ función de tolerancia

$$\hat{N}(\omega) - \hat{D}(\omega) F(\omega) \le \epsilon(\omega)\hat{D}(\omega)$$
$$-\hat{N}(\omega) + \hat{D}(\omega) F(\omega) \le \epsilon(\omega)\hat{D}(\omega)$$

$$\Downarrow$$

$$\hat{N}(\omega) - \hat{D}(\omega)[F(\omega) + \epsilon(\omega)] \le 0$$

$$-\hat{N}(\omega) + \hat{D}(\omega)[F(\omega) - \epsilon(\omega)] \le 0$$

$$-\hat{N}(\omega) \le 0$$

$$-\hat{D}(\omega) \le 0$$

Se resta la variable auxiliar $v$ de las
ecs. anteriores y se minimiza

$$Z = v$$

Solución si $\nu = 0$

Si $\nu > 0$. hay que modificar $E(\omega)$ o $F(\omega)$

Algoritmo de Remez para el caso racional

$$\sum_{i=0}^{m} c_i \cos \omega_i / \sum_{i=0}^{n} b_i \cos \omega_i - f(\omega) = \pm \delta$$

$$\sum_{i=0}^{m} c_i \cos \omega_i - f(\omega) \left[ b_0 + \sum_{i=1}^{n} b_i \cos \omega_i \right] =$$

$$(-1)^i \delta \left[ b_0 + \sum_{i=1}^{n} b_i \cos \omega_i \right]$$

Teorema de caracterización [4]

Sea $R = P/Q$ donde $P$ es un polinomio de grado $m$, $Q$ es un polinomio de grado $n$, $P/Q$ es irreducible y $Q > 0$ en $X$. $R$ es una mejor aproximación a $f \in C[X]$ si la función de error tiene al menos $m+n+2$ extremos en $X$

Solución iterando entre diseño de numerador y denominador [8]

## Referencias:

[1]   L. Rabiner + B. Gold, Theory and applications of digital signal processing, Prentice-Hall, 75

[2]   A. Oppenheim + R. Schafer, Digital Signal Processing, Prentice-Hall, 75

[3]   J. R. Rice, The approximation of fonction, Addison Wesley, 1964

[4]   E. W. Cheney, Introduction to approximation theory, Mc Graw-Hill, 1966

[5]   L. R. Rabiner, J. H. McClellan + T. W. Parks, FIR digital filter design techniques using weighted Chebyshev approximation, Proc IEEE, Abril 1975.

[6]   A. G. Deczky, Synthesis of recursive Digital Filters Using the minimum p. Error Criterion, IEEE Trans on Audio Electroacoust, vol AU-20, Oct. 72

[7]   L. R. Rabiner, N. Y. Graham + H. O. Helms, Linear programming design of IIR Digital Filters with arbitrary magnitude function, IEEE Trans. on ASSP, Abril 74

[8]   H. G. Martinez + T. W. Parks, "Design of recursive filters with optimum magnitude and attenuation poles on the unit circle", IEEE Trans on ASSP, Abril 78

DISEÑO OPTIMO DE SISTEMAS DE INGENIERIA

CONDICIONES DE OPTIMALIDAD

Dr. Antonio Montalvo

MARZO, 1982

## I. EL PROBLEMA GENERAL DE PROGRAMACION NO LINEAL

En el sentido más amplio, el problema general no lineal es el de encontrar un extremo (máximo o mínimo) de una función objetivo sujeta a restricciones de igualdad $y/_o$ no lineales. Sin embargo, en las si guientes secciones de estas notas han quedado excluidos los dos siguientes problemas:

*desigualdad, las cuales pueden ser lineales y/o*

    a) Las variables estan restringidas a valores enteros (programación no lineal entera)

    b) Las restricciones incluyen al parámetro tiempo en la forma de una ecuación diferencial (control óptimo).

En lo siguiente se supondrá que la función objetivo $f(x)$ es continua, $h_1(x), \ldots, h_m(x)$ denotan las restricciones de igualdad y $g_{m+1}(x), \ldots, g_p(x)$ las restricciones de desigualdad, donde: $x = \left(x_1, \ldots, x_n\right)^T$ es un vector columna de componentes $x_1, \ldots, x_n$ en un espacio euclidiano n-dimensional. (Las variables $x_1, x_2, \ldots, x_n$ pueden ser parámetros de diseño, ajuste de controles, lecturas de

instrumentos, etc., mientras que la función objetivo podría representar
el costo, peso, ganancias, etc.; finalmente, las restricciones pueden re
presentar requerimientos técnicos, condiciones de operación, etc., del
proceso).

El problema de programación no lineal se puede establecer
formalmente como :

$$\text{Minimizar}: \quad f(x), \quad x \in E^n \tag{1.1}$$

sujeta a m restricciones de igualdad, lineales $^y/_o$ no lineales,

$$h_j(x) = 0 \qquad\qquad j = 1, \ldots, m \tag{1.2}$$

y (p - m) restricciones de desigualdad, lineales $^y/_o$ no lineales,

$$g_j(x) \geq 0 \qquad\qquad j = m+1, \ldots, p \tag{1.3}$$

O en forma alterna como :

$$\text{Minimizar}: \quad \{f(x) | x \in R\} \tag{1.4}$$

donde R es el dominio de x para el cual (1.2) y (1.3), se satisfacen,
es decir :

$$R = \{x | h_j(x) = 0, \; g_j(x) \geq 0, \; \text{para toda } j\} \tag{1.5}$$

Un ejemplo sencillo de programación no lineal es
el que se muestra en la figura 1, y está dado por :

Minimice: $f(x) = x_1^2 + x_2^2 + 2 \, \widehat{x_2}$    $x_2$

sujeto a: $h_1(x) = x_1^2 + x_2^2 - 1 = 0$

$g_2 = x_1 + 2\widehat{x_2} - \frac{1}{2} \geq 0$
  $x_2$

$g_3 = x_1 \geq 0$

$g_4 = x_2 \geq 0$



FIGURA 1. Representación geométrica de un problema de programación no lineal

## II. NOTACION Y TERMINOLOGIA

El vector columna $x^* = \left( x_1^*, \ldots, x_n^* \right)^T$ que satisface
(1.1) - (1.3) se denomina punto óptimo, y el valor de $f(x^*)$ que le
corresponde se denomina valor óptimo de la función objetivo. La pareja,
$x^*$, $f(x^*)$ constituye la solución óptima. Para algunos problemas, pue
den existir varias categorías de soluciones óptimas si la función objetivo
no es unimodal (exhibe mas de un punto extremo) como se ilustra en la
figura 2. La solución global óptima representa el valor más pequeño de
$f(x)$, mientras que una solución óptima local (o relativa) representa
el valor más pequeño de $f(x)$ en una cierta vecindad del vector $x$:
es decir,

óptimo global $x^*$ satisface $\qquad f(x^*) < f(x) \, \forall x \in R$

óptimo local $x^*$ satisface $\qquad f(x^*) < f(x) \, | \, ||x-x^*|| \leq \delta(x^*)$

### 2.1. Concavidad y Convexidad

Los conceptos de concavidad y convexidad ayudan a determinar
bajo que condiciones una solución óptima local es tambien solución óptima
global.

Una función $\phi(x)$ se dice que es convexa en el dominio $R$,
si para cualesquiera dos vectores $x_1$ y $x_2$ $\in R$,

$$\phi\left(\theta x_1 + (1-\theta)x_2\right) \le \theta\phi(x_1) + \phi(x_2)(1-\theta) \tag{1.6}$$

donde $\theta$ es un escalar $0 \le \theta \le 1$. Además, $\phi(x)$ es _estrictamente_ _convexa_ si, para $x_1 \neq x_2$, el signo $\le$ en (1.6) se puede rem plazar por el signo de desigualdad $(<)$. Si en (1.6) la desigualdad contraria es la válida, se dice que la función $\phi(x)$ es cóncava $(\ge)$ o _estrictamente_ cóncava $(>)$. Note que si $\phi(x)$ es cóncava (cónvexa), $-\phi(x)$ es cónvexa (cóncava). (Las funciones lineales son, simultánea mente, convexas y cóncavas)..

Una función convexa diferenciable posee las siguientes propie dades.

a) $\phi(x_2) - \phi(x_1) \ge \nabla^T\phi(x_1)(x_2-x_1)$     para toda $x_1$ y $x_2$.

b) La matriz de segundas derivadas parciales de $\phi(x)$ con respecto a $x$ (matriz Hessiana) es positiva definida (o positiva semidefinida) para toda $x$ si $\phi(x)$ es estric tamente convexa (o convexa).

c) Sobre el dominio de $R$; $\phi(x)$ posee un sólo mínimo.

$$(1.7)$$

Un conjunto de puntos (o región) se define como _conjunto_ convexo en un espacio $n$-dimensional si, para toda pareja de puntos $x_1$ y $x_2$ en el conjunto, la línea recta que los une pertenece completamente al conjunto. Es decir, $R$ es convexo si para toda $x_1$ y $x_2$ $\varepsilon$ $R$

$$x = 0x_1 + (1-0)x_2 \in R$$

De los conceptos de convexidad emerge un resultado importante en programación matemática: para el problema de programación no lineal conocido como el problema de programación convexa

$$\text{Minimizar}: \quad f(x)$$

$$\text{sujeta a}: \quad g_j(x) \geq 0 \qquad j = 1, \ldots, p$$

$$x \geq 0$$

en el cual $(a)$, $f(x)$ es una función convexa y $(b)$ cada restricción de desigualdad es una función cóncava (las restricciones forman un conjunto convexo), se puede demostrar el siguiente resultado: el mínimo local también es mínimo global (Usando argumentos opuestos, el resultado opuesto, máximo, también es cierto).

## 2.2. Factibilidad

Cualquier vector $x$ que satisface tanto las restricciones de desigualdad como las de igualdad se llama punto factible. El conjunto de todos los puntos que satisfacen las restricciones constituyen el dominio factible de $f(x)$, y se denotará por $R$; cualquier punto no en $R$ se llama punto no factible.

Un óptimo restringido es uno para el cual el óptimo local cae en la frontera de la región factible. Si las restricciones son únicamente

de igualdad, un punto x factible debe caer en la intersección de todas

las hipersuperficies que satisfacen $h_j(x) = 0$

Con respecto a las restricciones de desigualdad, un punto x

se puede clasificar como punto interior (factible), punto frontera

(factible) o punto exterior (no factible). Los puntos interiores son

aquellos para los cuales $\bar{g}_j(x) > 0$ ; para un punto frontera, $g_j$

$g_j(x) = 0$ para al menos una restricción; y un punto exterior,

$g_j(x) < 0$ para al menos una restricción. Las restricciones se llaman

activas (o de atadura) si $g_j(x) = 0$.

Una región R de vectores admisibles puede ser convexa o

no convexa, según se describió con anterioridad, pero además puede ser

simplemente conexa ó no-simplemente conexa (ver Figura 2).



(a) simplemente conexa
(no convexa)

(b) no simplemente conexa
(obviamente no convexa)

FIGURA 2. Ejemplos de tipos de región

## 2.3. El Gradiente

El conjunto de puntos para los cuales una función $f(x)$ exhibe un valor constante, se llaman contornos de $f(x)$. Si la función $f(x)$ es continua y diferenciable, el gradiente de la función existe y está definido como el vector columna formado por las primeras derivadas parciales de $f(x)$ con respecto a $x$ es decir:

$$\nabla f(x) = \begin{bmatrix} \dfrac{\partial f(x)}{\partial x_1} \\ \cdot \\ \cdot \\ \cdot \\ \dfrac{\partial f(x)}{\partial x_n} \end{bmatrix} \tag{1.8}$$

Se puede demostrar que en el espacio métrico euclidiano, el gradiente de una función escalar apunta en la dirección de máximo incremento en el valor de la función, máximo ascenso, y que es, además, ortogonal a las líneas de contorno. El negativo del gradiente apunta en la dirección del máximo descenso de $f(x)$. Finalmente, cualquier vector V, ortogonal a $\nabla f(x)$, tal como la superficie tangente a $f(x)$, está definido por

$$V^T \nabla f(x) = 0 \qquad \text{(Figura 3)}$$

FIGURA 3. El gradiente, y la dirección de máximo descenso.

## 2.4. Aproximación de funciones

Algunos de los procedimientos de programación matemática que se discutirán más tarde requieren de aproximaciones lineales o cuadráticas para $f(x)$, $g(x)$ y $h(x)$.

Una aproximación lineal, o de primer orden, para una función $f(x)$, se puede hacer truncando la serie de Taylor, alrededor de un punto $x_0$, como

$$f(x) \approx f(x_0) + \nabla^T f(x_0)(x-x_0) \tag{1.9}$$

Para obtener una aproximación cuadrática, en la misma serie de Taylor se puéden despreciar los términos mayores e iguales a tercer orden, obteniéndose

$$f(x) = f(\bar{x}_0) + \nabla^T f(x_0)(x-x_0) + \frac{1}{2}(x-x_0)^T \nabla^2 f(x_0)(x-x_0) \qquad (1.10)$$

donde $\nabla^2 f(x)$ es la <u>matriz Hessiana</u> _de_ $f(x)$, $H(x)$, es decir

$$H(x_0) = \{h_{ij}(x_0)\} \quad i,j = 1, \ldots, n \qquad (1.11a)$$

donde $$\dot{h}_{ij} = \frac{\partial^2 f(x)}{\partial x_i \partial x_j} \qquad (1.11b)$$

Se observa facilmente que $H(x)$ es una matriz simétrica.

## 2.5. Condiciones Necesarias y Suficientes para que Una Solución sea Solución Optima

Para algunas clases especiales del problema general de programa
ción no lineal $\left[\text{Ecs }(1.1) - (1.3)\right]$ ha sido posible establecer criterios
de optimalidad. Sin embargo, para funciones completamente generales, no ha
sido posible establecer criterios de optimalidad precisos. En consecuencia,
unicamente se describirán algunos casos especiales, los cuales, sin embargo,
son bastante comunes y de importancia práctica. Las condiciones que determi
nan si un vector $x$ resuelve o no el problema de programación no lineal serán
presentadas en una serie de teoremas los cuales no serán demostrados (las
demostraciones están fuera de los objetivos de este curso).

### 2.5.1. Programación no-lineal sin restricciones

El problema es el siguiente:

$$\text{Minimizar:} \quad f(x), \, x \in E^n \qquad (1.12)$$

Las condiciones necesarias para que $x^*$ sea un mínimo local del problema (1.12) son:

1. $f(x)$ diferenciable en $x^*$.

2. $\nabla f(x^*) = 0$, es decir, existe un punto estacionario de $f(x)$ en $x^*$

Las condiciones suficientes para que $x^*$ sea un mínimo local del problema (1.12) son:

3. $\nabla^2 f(x) > 0$; es decir, la matriz Hessiana es positiva definida.

(Las condiciones para la existencia de un máximo son iguales, excepto que el hessiano deberá ser negativo definido).

### 2.5.2. Programación no-lineal con restricciones de Igualdad y Desigualdad

En este caso el problema es el siguiente:

$$\text{Minimizar:} \quad f(x) \qquad\qquad x \in E^n \qquad \Big)$$
$$\Big)$$
$$\text{Sujeta a:} \quad h_j(x) = 0 \qquad\quad j = 1, \ldots, m \qquad \Big) \; (1.13)$$
$$\Big)$$
$$g_j(x) \geqslant 0 \qquad\quad j = m+1, \ldots, p \quad \Big)$$

Las condiciones necesarias para que $x^*$ sea un mínimo local se establecen en dos teoremas, el primero de los cuales (teorema 2) puede ser llamado condiciones de primer orden (debido a que las funciones que intervienen se consideran una vez diferenciables). El segundo teorema (teorema 3) se le denomina condiciones de segundo orden (se considera que las funciones son dos veces diferenciables).

Para establecer las condiciones necesarias, empezaremos con el siguiente concepto: si $x^*$ es un mínimo local de $f(x^*)$, ésta no puede decrecer a lo largo de ningun arco "suave" dirigido desde $x^*$ hacia la región factible. Sea el vector $V$ tangente al arco que empieza en $x^*$. Usando los conceptos de Fiacco y Mc Cormick, se asignan tres diferentes categorías o clases al vector $V$, donde cada conjunto $V_i$ incluye el conjunto de $V$ tales que: (ver tabla I).

Todas las posibles perturbaciones de $x^*$ caen en la unión de $V_1$ y $V_2$ y si $V \in V_2$, $f(x)$ decrece, mientras que si $V \in V_1$, $f(x)$ se incrementa o es constante. En esencia, las condiciones necesarias de primer orden imponen el requerimiento de que el conjunto $V_2$ esté vacío.

| CLASE | DESIGUALDAD | IGUALDAD | FUNCION OBJETIVO |
|---|---|---|---|
| $V_1$ | $\left\{\begin{array}{l} v^T \nabla g_j (x^*) \geq 0 \\ \text{para las restricciones} \\ \text{activas} \end{array}\right\}$ y | $\left\{\begin{array}{l} v^T \cdot \nabla h_j (x^*) = 0 \\ j = 1,2, \ldots, m \end{array}\right\}$ y | $\left\{ v^T \nabla f (x^*) \geq 0 \right\}$ |
| $V_2$ | $\left\{\begin{array}{l} v^T \nabla g_j (x_i) \geq 0 \\ \text{para las restricciones} \\ \text{activas} \end{array}\right\}$ y | $\left\{\begin{array}{l} v^T \nabla h_j (x^*) = 0 \\ j = 1,2, \ldots, m \end{array}\right\}$ y | $\left\{ v^T \nabla f (x^*) < 0 \right\}$ |
| $V_3$ | $\left\{\begin{array}{l} v^T \nabla g_j (x^*) < 0 \\ \text{para, al menos, una} \\ \text{restricción activa} \end{array}\right\}$ o | $\left\{\begin{array}{l} v^T \nabla h_j (x^*) \neq 0 \\ \text{para, al menos, una} \\ \text{restricción} \end{array}\right\}$ | |

TABLA I. Clasificación de los conjuntos $V_i$

Si $V_2$ está vacío, se puede demostrar la existencia de los multiplicadores

de Lagrange, resultando el siguiente teorema.

TEOREMA I.

Si (a) $x^*$ satisface el problema (1.13), (b) las funciones

$f(x)$, $g_j(x)$ son una vez diferenciables, y (c) en $x^*$ $V_2$ está

vacío, entonces existen los vectores $u^*$ y $w^*$ (multiplicadores de

Lagrange) tales que, $(u^*, w^*, x^*)$ satisfacen

(1)   $h_j(x^*) = 0$                     $j = 1, \ldots, m$

(2)   $g_j(x^*) \geq 0$                  $j = m + 1, \ldots, p$

(3)   $u_j^* g_j(x) = 0$                 $j = m + 1, \ldots, p$

(4)   $u_j^* \geq 0$                     $j = m + 1, \ldots, p$

(5)   $\nabla L(x^*, u^*, w^*) = 0$

donde la función

$$L(x, u, w) = f(x) + \sum_{j=1}^{m} w_j h_j(x) - \sum_{j=m+1}^{p} u_j g_j(x)$$

puede ser considerada como una función Lagrangiana generalizada, asociada

al problema (1.13).

Con el propósito de establecer bajo que circunstancias el conjunto

$V_2$ está vacío, se necesita calificar, a primer orden, a las restricciones.

Sea $x^*$ un punto factible del problema (1.13) y supóngase que $h_1(x), \ldots,$ $h_m(x)$, $g_{m+1}(x), \ldots, g_p(x)$ son funciones una vez diferenciables. La calificación a primer orden de las restricciones es una condición que se impone unicamente sobre las restricciones (sin importar la función objetivo) y que consiste en que para cada punto frontera formado por el conjunto de restricciones de igualdad y las activas de desigualdad, debe de existir una curva suave que termine en el punto frontera y que pertenezca completamente al conjunto de las restricciones. Si $x^*$ es un mínimo local de $f(x)$, ésta no puede decrecer a lo largo de tal curva, dirigida desde $x^*$ hacia la región factible. Una <u>condición suficiente para la calificación a primer orden de las restricciones</u> que debe cumplirse es que todos los gradientes de las res tricciones de desigualdad activas y los gradientes de las restricciones de igualdad evaluados en $x^*$, sean linealmente independientes. Esto último se establece en el siguiente teorema:

<u>TEOREMA 2.</u>

Si las funciones $h_1(x), \ldots, h_m(x)$, $g_{m+1}(x), \ldots, g_p(x)$ son una vez diferenciables en $x^*$, y si la calificación a primer orden de las restricciones es válida en $x^*$, entonces la condición necesaria para que $x^*$ sea un mínimo local del problema (1.7), es que existan multiplicadores de Lagrange $u^*$ y $v^*$ tales que $(u^*, v^*, x^*)$ satisfagan las ecuaciones (1) - (5) del Teorema 1.

Para tomar en cuenta la curvatura de las funciones en el problema (1.13), Mc Cormick estableció las condiciones necesarias de segundo orden para que $x^*$ sea un mínimo local. Supongase que las funciones $f(x)$, $h_1(x), \ldots, h_m(x)$, $g_{m+1}(x), \ldots, g_p(x)$ son dos veces diferenciables en $x^*$, un punto que satisface al problema (1.13). Sea V cualquier vector no cero tal que

$$V^T \nabla g_j(x) = 0 \qquad \text{para las restricciones de desigualdad activas.}$$

$$V^T \nabla h_j(x) = 0. \qquad \text{para las restricciones de igualdad}$$

Entonces, si V es la tangente a una curva $\psi(\theta)$, $\theta \geq 0$, dos veces diferenciable, a lo largo de la cual $g_j(\psi(\theta)) = 0$ para todas las restricciones de desigualdad activas, y $h_j(\psi(\theta)) = 0$ para todas las restricciones de igualdad, la calificación a segundo orden de las restricciones en $x^*$ es válida. Una condición suficiente para la calificación a segundo orden de las restricciones que debe cumplirse es que los gradientes de las restricciones de desigualdad activas en $x^*$ y los gradientes de las restricciones de igualdad en $x^*$ sean linealmente independientes.

Las condiciones necesarias de segundo orden se pueden establecer como sigue.

TEOREMA 3.

(a) Si las funciones $f(x)$, $h_1(x)$, ..., $h_m(x)$, $g_{m+1}(x)$, ..., $g_p(x)$ son dos veces diferenciables en $x^*$, y (b) si la calificación a primer orden de las restricciones es válida en $x^*$, entonces las condiciones necesarias para que $x^*$ sea un mínimo local del problema (1.13) son que exis_tan $u^*$ y $v^*$ tales que (c) ecuaciones (1) - (5) del Teorema 1 se sa_tisfagan, y (d) para cada vector no cero $V$, para el cual $V^T \nabla g_j(x^*) = 0$, para las restricciones de desigualdad activas, y $V^T \nabla h_j(x^*) = 0$, para las restricciones de igualdad, se cumple lo siguiente:

$$(6) \quad V^T \nabla^2 L(x^*, u^*, v^*) V \geq 0$$

Las condiciones suficientes para que $x^*$ sea un mínimo local aisla_do del problema (1.13) son las mismas (a), (b) y (c) del Teorema 3, excepto la parte (d) ( ecuación (6) ), la cual debe ser substituida por:

(d') Para cada vector $V$ no cero para el cual $V^T \nabla g_j(x^*) = 0$ para las restricciones de desigualdad activas, $V^T \nabla g_j(x^*) \gtreqless 0$ para las restricciones de desigualdad no activas y $V^T \nabla h_j(x^*) = 0$ para las restricciones de igualdad; lo siguiente es verdadero:

$$(6') \quad V^T \nabla^2 L(x^*, u^*, v^*) V > 0$$

Ejemplo 1.    Condiciones necesarias y suficientes con

restricciones  de  desigualdad.

Minimizar :      $f(x) = x_1^2 + x_2$

Sujeta  a:       $g_1(x) = (-x_1^2 + x_2^2) + 9 \geq 0$

                 $g_2(x) = -x_1 - x_2 + 1 \geq 0$.



FIGURA 4.- Región admisible y curvas de contorno del Ejemplo 1.

Se observa que  $g_1(x)$  es una restricción activa

mientras que  $g_2(x)$. no lo es.

Ya que sólo una de las restricciones está activa, no se necesita

comprobar la calificación a primer - y segundo orden de las restricciones

(Note que $f'(x)$, $g_1(x)$ y $g_2(x)$ son dos veces diferenciables).

De acuerdo a los Teoremas (1) y (2) se necesita demostrar que exi

ten $u^*$ y $x^*$ tales que

(2) $\quad g_j(x^*) \geq 0$

$$- x_1^{*2} - x_2^{*2} + 4 \geq 0$$

$$- x_1^* + (- x_1^*) + 1 \geq 0$$

(3) $\quad u_j^* \, g_j(x^*) = 0$

$$u_1^* \, (- x_1^{*2} - x_2^{*2} - 4) = 0$$

$$u_2^* \, (- x_1^* - x_2^* + 1) = 0$$

(4) $\quad u_j^* \geq 0$

$$u_1^* \geq 0$$

$$u_2^* \geq 0$$

(5) $\quad \nabla L(x^*, u^*) = 0 \qquad (L = f(x) - u_1 g_1(x) - u_2 g_2(x))$

$$\begin{pmatrix} 2 x_1^* \\ 1 \end{pmatrix} - u_1^* \begin{pmatrix} - 2x_1^* \\ - 2x_2^* \end{pmatrix} - u_2^* \begin{pmatrix} - 1 \\ - 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

Resolviendo las ecuaciones anteriores se puede verificar que :

$$x^* = [\, 0, \, -3\, ]^T$$

$$y \qquad u^* = (\tfrac{1}{6}, \, 0\, )^T$$

La condición de segundo orden que debe ser satisfecha es :

$$v^T \, \nabla g_j \, (x^*) = 0 \qquad\qquad g_j (x^*) : \text{ restricción activa}$$

es decir: 
$$(v_1, v_2) \begin{pmatrix} -2x_1^* \\ -2x_2^* \end{pmatrix} = v_1 (0) + v_2 (6) = 0$$

de donde $V_1$ puede tomar cualquier valor, y $V_2 = 0$, substituyendo

en ( 6), se tiene

$$(6) \qquad v^T \, \nabla^2 L \, (x^*, \, u^*) \; v \geq 0$$

$$\text{donde } \nabla^2 L = \begin{bmatrix} 2 (1 + u_1) & 0 \\ 0 & 2u_2 \end{bmatrix}$$

Ejemplo 2. Condiciones necesarias y suficientes con restricciones de

Igualdad y de Desigualdad

$$\text{Minimice:} \qquad f(x) = x_1^2 + x_2$$

$$\text{Sujeta a:} \qquad h_1(x) = x_1^2 + x_2^2 - 9 = 0 \qquad j$$

$$g_2(x) = -(x_1 + x_2^2) + 1 \geq 0$$

$$g_3(x) = -(x_1 + x_2) + 1 \geq 0$$

FIGURA 5. Región admisible y curvas de nivel del Ejemplo 2.

De acuerdo al Teorema 3, $h_1(x)$, $g_2(x)$ y $g_3(x)$ deben ser dos veces diferenciables. Además, los gradientes de las restricciones activas deberán ser linealmente independientes para que satisfagan la califi cación a primer- y segundo orden. Suponga que A, localizado en $x^* = (-2.37, -1.84)^T$ en la intersección de $h_1(x^*) = g_2(x^*) = 0$ sea un candidato a mínimo local. Si se forma una combinación lineal entre los gradientes de $h_1(x)$ y $g_2(x)$, en $x^*$, se tiene que, para que sean

linealmente independientes

$$C_1 \nabla h_1(x^*) + C_2 \nabla g_2(x^*) = 0$$

$$C_1 = C_2 = 0$$

o

$$C_1 \begin{pmatrix} 2x_1^* \\ 2x_2^* \end{pmatrix} + C_2 \begin{pmatrix} -1 \\ -2x_2^* \end{pmatrix} = 0$$

y como

$$\det \begin{pmatrix} 2x_1^* & - & 1 \\ 2x_2^* & - & 2x_2^* \end{pmatrix} \neq 0 \implies C_1 = C_2 = 0$$

es decir, los gradientes de las restricciones activas son linealmente independientes.

De acuerdo a los Teoremas (1) y (2) se debe cumplir que.

(1)  $h_1^-(x^*) = 0$,    $x_1^{*2} + x_2^{*2} = 0$.

(2)  $g_j(x^*) = 0$    $-(x_1^* + x_2^{*2}) + 9 \geq 0$

$-(x_1^* + x_2^*) + 1 \geq 0$

(3)  $u_j^* g_j(x^*) = 0$    $u_2^* \left[ -(x_1^{*2} + x_2^{*2}) + 9 \right] = 0$

$u_3^* \left[ -(x_1^* + x_2^*) + 1 \right] = 0$

(4)  $u_j^* \geq 0$    $u_2^* \geq 0$

$u_3^* \geq 0$

(5) $\nabla L (\overset{*}{x}, \overset{*}{u}, \overset{*}{w}) = 0$    $L = f(\overset{*}{x}) + \overset{*}{w_1} h_1(\overset{*}{x}) - \overset{*}{u_2} g_2(\overset{*}{x}) - \overset{*}{u_3} g_3 (x$

$$\begin{pmatrix} 2 x_1^* \\ \\ 1 \end{pmatrix} + w_1^* \begin{pmatrix} 2 x_1^* \\ \\ 2 x_2^* \end{pmatrix} - u_2^* \begin{pmatrix} -1 \\ \\ -2 x_2^* \end{pmatrix} - u_3^* \begin{pmatrix} -1 \\ \\ -1 \end{pmatrix} = \begin{pmatrix} 0 \\ \\ 0 \end{pmatrix}$$

de donde se obtiene que

$$w_1^* = -0,779$$

$$u_1^* = 1,05$$

$$u_3^* = 0$$

De acuerdo al Teorema 3, se debe encontrar V, tal que

$v^T \nabla g_j(x^*) = 0$  para las restricciones de desigualdad activas, y

además  $v^T \nabla h_1(x^*) = 0$,  es decir

$$(v_1, v_2) \begin{pmatrix} -1 \\ \\ -2 x_2^* \end{pmatrix} = 0 \qquad v_1 + 2 x_2^* v_2 = 0$$

$$\begin{bmatrix} v_1, v_2 \end{bmatrix} \begin{bmatrix} 2 x_1^* \\ \\ 2 x_2^* \end{bmatrix} = 0 \qquad 2 x_1^* v_1 + 2 x_2^* v_2 = 0$$

de donde se obtiene que $V = (v_1, v_2)^T = (0, 0)^T$ y entonces

existe una única solución al problema en la intersección de $g_2(x)$ y $h_1(x)$.

Note que en este problema en particular

$$v^T \nabla^2 L (x^*, u^*, w^*) \ V \geq 0$$

para toda $V \neq 0$, es decir $\nabla^2 L(x^*, u^*, w^*)$ es una matriz positiva

definida. ( condición (6') ).

· Los dos ejemplos analizados en este capitulo tienen el propósito

de ilustrar las condiciones de optimalidad de primero y segundo orden.

Ahora bien, los problemas que se pueden tratar analiticamente son demasiado

sencillos, sin que esto quiera decir que para problemas muy grandes $^y/_0$

complicados lo anterior no sea válido.

# B I B L I O G R A F I A

1. A.V. Fiacco and G.P. Mc Cormick. "Non linear Programming",

    John Wiley and Sons, Inc., New York, 1968.


2. W.W. Kuhn and A. W. Tucker, "Non linear Programming. Proc. 2nd.

    Berdeley Symp. on Mathematical Statistics and Probability,

    Univ. of Calif. Press, Berkeley, 1951, pp. 481-492.


3. G.P. Mc. Cormick. SIAM J. Appl. Math., 15 : 641 (1967)


4. O.L. Mangasarian, "Non linear Programming". Mc. Graw-Hill,

    New York, 1969.

DISEÑO OPTIMO DE SISTEMAS DE INGENIERIA

METODOS DE OPTIMACION CON RESTRICCIONES

Dra. Susana Gómez Gómez

MARZO, 1982

$$\boxed{\text{M E T O D O L O G I A}}$$

## I. BUSQUEDAS UNIDIRECCIONALES

Casi la totalidad de las técnicas de minimización que se describi-
rán en otras secciones más adelante, requieren de técnicas de minimización
unidimensionales las cuales tienen como propósito el localizar el mínimo
local de una función de una variable. La implementación de los métodos
mencionados requieren del conocimiento previo de un cierto intervalo $\Delta^{(o)}$ e'
cual contiene al mínimo de la función objetivo $f(x)$, y además se supone
que en el intervalo prescrito la función es unimodal. En la Tabla II
se mencionan algunos de los métodos más conocidos para lograr la minimiza-
ción deseada; todos los cuales tienen como propósito reducir el tamaño del
intervalo $\Delta^{(o)}$ hasta un tamaño $\Delta^{(n)}$. Para comparar la rapidez re-
lativa de los métodos, Wilde define una eficiencia para n evaluaciones de
la función como :

$$\text{Eficiencia} = E = \frac{\Delta^{(n)}}{\Delta^{(o)}}$$

En la tabla II se comparan los valores de E para varios métodos.

$E$

| Métodos no-secuenciales | $E$ | Secuenciales | $E$ |
|---|---|---|---|
| Búsqueda uniforme | $\dfrac{2}{n+1}$ | Búsqueda secuencial (Dicotomus) | $\dfrac{1}{2^{n/2}}$ |
| Búsqueda uniforme (Dicotomus) | $\dfrac{1}{\left(\dfrac{n}{2}\right)+1}$ | Búsqueda Fibonacci | $\dfrac{1}{F\,n}$ |
| | | Sección Dorada | $(0.618)^{1-n}$ |

(*): Número de Fibonacci para $n$ evaluaciones

TABLA II. Eficiencia de Técnicas de Búsqueda Unidimensional

En la Tabla III se compara el número de evaluaciones de la función que se requiere para reducir un intervalo inicial de $5 \times 10^{-1}$ a uno menor.

| | NO SECUENCIAL | | SECUENCIAL | | |
|---|---|---|---|---|---|
| $\Delta^{(n)}$ | Uniforme | Dicotomus | Dicotomus | Fibonacci | Sección Dorada |
| $5 * 10^{-3}$ | 199 | 198 | 14 | 11 | 11 |
| $5 * 10^{-5}$ | 19,999 | 19,998 | 28 | 21 | 21 |

TABLA III. Número de evaluación de la función para reducir $\Delta^{(o)} = 5 \times 10^{-1}$

A continuación se describe el método de la sección dorada. Este método esta basado en la división de una linea en dos segmentos tales que la relación del tamaño original de la linea al segmento mayor es la misma

que la relación del segmento mayor al menor, es decir :

$$F_1 + F_2 = 1$$



$$\frac{1}{F_2} = \frac{F_2}{F_1} \quad ; \quad F_2^2 = F_1$$

de donde

$$F_1 = \frac{3 - \sqrt{5}}{2} \simeq 0.38$$

$$F_2 = \frac{\sqrt{5} - 1}{2} \simeq 0,62$$

Para iniciar la búsqueda del mínimo de $f(x)$, se necesita especificar (o averiguar) en que dirección ésta se llevará a cabo. (Se supondrá que se conoce).

Como primer paso se debe encontrar un intervalo $\Delta$ donde se encuentra el mínimo de $f(x)$ usando, por ejemplo, una serie de pasos cada vez más grande sobre la variable independiente. Suponga que esto se ha hecho y que los últimos tres puntos obtenidos en $x$ son los siguientes : $x_3^{(o)}$, el último, $x_2^{(o)}$ y $x_1^{(o)}$, donde $f(x_3^{(o)}) \geq f(x_2^{(o)})$ y sea $\Delta^{(k)} = x_3^{(k)} - x_1^{(k)}$ (Fig. 6). (Lo anterior) Una vez hecho la k-ésima etapa modifica el intervalo de la siguiente forma.

$$y_1^{(k)} = x_1^{(k)} + F_1 \Delta^{(k)}$$

$$y_2^{(k)} = x_1^{(k)} + F_2 \Delta^{(k)} = x_3^{(k)} - F_1 \Delta^{(k)}$$

si $f(y_1^{(k)}) < f(y_2^{(k)})$ : $\Delta^{(k+1)} = (y_2^{(k)} - x_1^{(k)})$, y $x_1^{(k+1)} = x_1^{(k)}$, $x_3^{(k+1)} = y_2^{(k)}$

si $f(y_1^{(k)}) > f(y_2^{(k)})$ : $\Delta^{(k+1)} = (x_3^{(k)} - y_1^{(k)})$, y $x_1^{(k+1)} = y_1^{(k)}$, $x_3^{(k+1)} = x_3^{(k)}$

si $f(y_1^{(k)}) = f(y_2^{(k)})$ : $\Delta^{(k+1)} = (y_2^{(k)} - x_1^{(k)}) = (x_3^{(k)} - y_1^{(k)})$, y

$$x_1^{(k+1)} = x_1^{(k)}, \quad x_3^{(k+1)} = y_2^{(k)}, \quad \text{o bien}$$

$$x_1^{(k+1)} = y_1^{(k)}, \quad x_3^{(k+1)} = x_3^{(k)}$$



$$y_1^{(o)} = x_1^{(o)} + 0.38 \Delta^{(o)}$$

$$y_2^{(o)} = x_1^{(o)} + 0.62 \Delta^{(o)}$$

FIGURA 6. Búsqueda mediante sección Dorada.

Otra clase de métodos para búsquedas unidireccionales, localizan un punto $x$, cercano a $x^*$ (el mínimo), mediante interpolación y extrapolación. En estos métodos se usan interpolaciones cuadráticas y cúbicas para aproximar el valor de la función. A continuación se describen un par de algoritmos que al usarlos en forma conjunta generan un algoritmo bastante poderoso para la localización de mínimos locales en una dirección: estos dos algoritmos son los siguientes-.

a) Davis-Swann-- Compey (DSC) para definir el intervalo $-\Delta-$ donde se encuentra el mínimo, y

b) Powell, para definir la localización "exacta" del mínimo.

En el método DSC, se toman pasos de tamaño cada vez mayor hasta que se sobrepasa la localización del mínimo. A partir de ese momento se usa interpolación cuadrática (Figura 7). Los pasos que se siguen en este algoritmo son los siguientes

1. Evaluar $f(x)$ en el punto inicial $x^{(o)}$.
   Si $f(x^{(o)} + \Delta x) \leq f(x^{(o)})$ se continua con el paso 2.
   Si $f(x^{(o)} + \Delta x) > f(x^{(o)})$ se define $\Delta x = -\Delta x$ (se cambia la dirección de búsqueda) y se continua con el paso 2

2. $x^{(k+1)} = x^{(k)} + \Delta x$

3. Calcular $f(x^{(k+1)})$

4. Si $f(x^{(k+1)}) \leq f(x^{(k)})$, se duplica el tamaño de paso $\Delta x$ y se repite el procedimiento desde el paso 2. Si $f(x^{(k+1)}) > f(x^{(k)})$ sea $x^{(m)} = x^{(k+1)}$, $x^{(m-1)} = x^{(k)}$, etc., redúzcase $\Delta x$ a la mitad y regrese a los pasos 2 y 3 una vez más.

5. De los cuatro puntos igualmente espaciados $x$, en el conjunto $\{x^{(m+1)}, x^{(m)}, x^{(m-1)}, x^{(m-2)}\}$, se eliminan o $x^{(m)}$ o $x^{(m-2)}$, el que este más lejos de la $x$ con el valor de la función más baja. Los tres puntos restantes se denotan como $x^{(a)}$, $x^{(b)}$ y $x^{(c)}$, donde $x^{(b)}$ es el punto central y $x^{(a)} = x^{(b)} - \Delta x$ y $x^{(c)} = x^{(b)} + \Delta x$.

6. Use interpolación cuadrática para estimar $x^*$,

$$x^* = \bar{x}^* = x^{(b)} + \frac{\Delta x \left[f(x^{(a)}) - f(x^{(c)})\right]}{2\left[f(x^{(a)}) - 2f(x^{(b)}) + f(x^{(c)})\right]}$$

Los pasos anteriores completan la primera etapa del método DSC. Para continuar, se reinicia en $\bar{x}^*$ o $x^{(c)}$, si $f(x^{(c)}) < f(\bar{x})$, se reduce $\Delta x$ y se empieza desde el paso 1 nuevamente.

FIGURA 7. Método DSC para minimización unidimensional.

En el método de Powell, se usa una aproximación cuadrática usando los tres primeros puntos obtenidos en la dirección de búsqueda dada. La x correspondiente al mínimo de la función cuadrática se usa para efectuar una nueva aproximación y se continua de esta forma hasta localizar el mínimo de f (x) (Figura 8).



FIGURA 8. Método de Powell para minimización unidimensional.

Los pasos que deben seguirse para implementar el método de Powell son:

1. dados $x^{(1)}$ y $\Delta$, calcule $x^{(2)} = x^{(1)} + \Delta$

2. calcule $f(x^{(1)})$ y $f(x^{(2)})$

3. Si $f(x^{(1)}) > f(x^{(2)})$, $x^{(3)} = x^{(1)} + 2\Delta$

   Si $f(x^{(1)}) \leq f(x^{(2)})$, $x^{(3)} = x^{(1)} - \Delta$

4. Calcule $f(x^{(3)})$

5. Estime el valor de x en el mínimo de $f(x)$, $\bar{x}^*$, como

$$x^* \simeq \bar{x}^* = \frac{1}{2} \frac{\left[(x^{(2)})^2 - (x^{(3)})^2\right] f(x^{(1)}) + \left[(x^{(3)})^2 - (x^{(1)})^2\right] f(x^{(2)}) + \left[(x^{(1)})^2 - (x^{(2)})^2\right] f(x^{(3)})}{(x^{(2)} - x^{(3)}) f(x^{(1)}) + (x^{(3)} - x^{(1)}) f(x^{(2)}) + (x^{(1)} - x^{(2)}) f(x^{(3)})}$$

6. Si $\bar{x}^*$ o alguno de entre $\{x^{(1)}, x^{(2)}, x^{(3)}\}$ que corresponda al valor más pequeño de $f(x)$, difiere en menos que la exactitud prescrita para x, o la exactitud en el valor de $f(x)$, se termina la búsqueda. En caso contrario, evalúe $f(\bar{x}^*)$ y se elimina del conjunto $\{x^{(1)}, x^{(2)}, x^{(3)}\}$ del que tenga el valor más alto de la función $f(x)$, a menos que se pierda el intervalo de x donde se encuentra el mínimo, en cuyo caso se debe eliminar una x tal que el intervalo adecuado no se pierda. Se repita el procedimiento desde el paso 5.

**2**

1. MINIMIZACION SIN RESTRICCIONES USANDO DERIVADAS

El problema general de minimización sin restricciones se puede plantear como :

$$\text{Minimizar} : \quad f(x), \quad x \in E^n \tag{1}$$

donde $f(x)$ es la función objetivo. Según se vió en el capitulo anterior, se pretende encontrar un punto $x^*$ tal que $\nabla f(x^*) = 0$. En la presente sección se atacará el problema definido en (1) mediante e. uso de métodos que hagan *empleen* ~~uso de las~~ primeras y segundas derivadas parciales $\left(\nabla f(x) \ y \ \nabla^2 f(x)\right)$ de la función objetivo.

## 2.1. Método del Máximo Descenso (Gradiente).

Según se recordará del capitulo anterior, el gradiente de la función objetivo $f(x)$, en cualquier punto x, es un vector dirigido en la dirección del máximo incremento local en $f(x)$. Resulta obvio que se *puede* escoger como dirección de búsqueda, para minimizar $f(x)$, la dirección opuesta al gradiente, $-\nabla f(x)$, esto es, en la dirección de máximo descen so. Si se escoge como dirección de búsqueda la ya señalada, resultará lo siguiente.

a) sea $d = -\nabla f(x_k)$ la dirección de búsqueda en el k-ésimo paso del algoritmo.

b) Si $\Delta x_k = + \alpha_k S_k = - \alpha_k \nabla f(x_k)$, es el desplazamiento

del punto $x_k$ al $x_{k+1}$, es decir

$$x_{k+1} = x_k + \Delta x_k \qquad (\alpha_k \text{ un escalar positivo})$$

c) Entonces, la aproximación a primer orden de $f(x)$ quedará

como

$$f(x_k + \Delta x_k) = f(x_k) + \nabla^T f(x_k) \Delta x_k$$

o bien

$$f(x_k + \Delta x_k) - f(x_k) = - \alpha_k \nabla^T f(x_k) \nabla f(x_k) < 0$$

es decir, se puede garantizar que para $\alpha_k$ suficientemente

pequeño, el valor de la función en $x_{k+1}$ decrecerá con res-

pecto al valor previo en $x_k$, siempre y cuando $\nabla f(x_k) \neq 0$.

(En los puntos (a) — (c), $\alpha_k$ se le conoce como tamaño de

paso).

Existen varias alternativas para escoger el tamaño de paso $\alpha_k$

de las cuales se pueden mencionar las siguientes

a) Fijar el valor de $\alpha_k$ de antemano e igual para todas las

iteraciones del método. La desventaja de proceder de esta

forma es que no se puede garantizar que de una iteración a

la siguiente, el valor de la función decrezca, ya que esta

propiedad sólo es válida cuando $\alpha_k \longrightarrow 0$.

Debido a lo anterior, el método puede presentar las siguientes

desventajas. En primer lugar, que si $\alpha_k$ no es lo "suficiente

mente pequeña" el método ascilará. Por otro lado, si $\alpha_k$ se

escoge "demasiado pequeña" la convergencia puede volverse

demasiado ~~pequeña~~ *lenta*

b) Una segunda alternativa es escoger $\alpha_k$ en cada iteración de

forma tal que el valor de la función objetivo se reduzca para

algún cierto valor de $\alpha_k$. Para lograr lo anterior se puede

proceder de la siguiente forma en cada etapa

o) se escoge $\alpha_{ref} > 0$ y $0 < \pi < 1$ un factor multiplicativo

i) se fija $\beta_0 = \alpha_{ref}$, $i = 0$

ii) $y^{(i)} = x_k + \beta_0 \Delta x_k$

iii) calcular $f(y^i)$

iv) si $f(y^i) < f(x_k)$ continuar con el paso vii)

v) $i \longleftarrow i + 1$

vi) $\beta_{i+1} = \pi \beta_i$ ; repetir el procedimiento desde (ii)

vii) $x_{k+1} = y^{(i)}$ ; $\alpha_k = \beta_i$

$f(x_{k+1}) = f(y^i)$

El procedimiento anterior garantiza que, si se permite un número

ilimitado de iteraciones $\beta_{i+1} = \pi \beta_i$, en cada etapa del método de

máximo descenso, se obtendrá un descenso en la función objetivo. Tiene

el defecto a que puede consumir demasiado tiempo en la búsqueda de un ta

maño de paso adecuado $\alpha_k$.

c) Una tercera alternativa es la siguiente. Supongase que de

conocida

~~alguna forma~~ la dirección de búsqueda $S_k$, es decir, el nuevo iterando

$x_{k+1}$ caerá sobre la ecuación de un parámetro

$$x_{k+1} = x_k + \alpha S_k$$

siendo $\alpha$ el parámetro. La diferencia entre este procedimiento y los dos

anteriores es que, en el presente se pide que el tamaño de paso $\alpha$ sea tal

que $f(x_{k+1}) = f(x_k + \alpha S_k)$ adquiera su mínimo valor, formalmente

$$\frac{d f(x_k + \alpha S_k)}{d\alpha} = 0$$

Por ejemplo, si $f(x)$ es una función cuadrática

$$f(x) = a + b^T x + \tfrac{1}{2} x^T Q x \quad (Q : \text{positiva definida} \\ \text{simétrica})$$

entonces

$$S_k = -\nabla f(x_k) = -b - Q x_k$$

$$x_{k+1} = \chi_k - \alpha_k \ (b + Q x_k)$$

$$\frac{df[x_k - \alpha_k(b+Qx_k)]}{d\alpha} \quad \xcancel{f(x_k - \alpha_k(b+Qx_k))} = 0 = \nabla^T f(x_k) \ s_k + s_k^T Q(\alpha_k s_k)$$

de donde

$$\alpha_k = -\frac{\nabla^T f(x_k) s_k}{s_k^T Q s_k}$$

Si se supone que $f(x)$ no es una función cuadrática de $x$, entonces se pueden emplear cualquiera de las técnicas de búsquedas unidimensionales descritas en la sección anterior.

Una característica interesante de este procedimiento de minimización es que el gradiente en el nuevo punto, $\nabla f(x_{k+1})$ es ortogonal a la dirección de búsqueda empleada para localizar $x_{k+1}$, es decir:

$\nabla^T f(x_{k+1}) \ s_k = 0$, lo cual se demuestra como sigue. Supóngase la misma función cuadrática ya empleada entonces

$$\nabla f(x_k) = b + Q x_k$$

y de la expresión para $\dfrac{df(\alpha)}{d\alpha} = 0$ se obtiene

$$(b + Q x_k)^T \ s_k + s_k^T Q \ \alpha_k \ s_k = 0$$

y como $\quad \alpha x_{k+1} - x_k = \alpha_k s_k$, entonces

$$s_k^T (b + Q x_k) + s_k^T Q (x_{k+1} - x_k) = 0.$$

o bien

$$s_k^T (b + H x_{k+1}) = s_k^T \quad \nabla f (x_{k+1}) \equiv 0.$$

Lo anterior se ilustra en la Figura 9.



FIGURA 9. Ortogonalidad de las direcciones de Búsqueda.

La implementación práctica del algoritmo de máximo descenso, con cualquiera de las tres alternativas discutidas para determinar $\alpha_k$, quedaría de la siguiente forma.

1. $k = 0$

2. Estimar $x_0$ (punto de arranque).

3. Calcular $f(x_k)$, $\nabla f(x_k)$

4. Si $\nabla^T f(x_k)\ \nabla f(x_k) \leq$ tolerancia, se ha encontrado un punto estacionario de la función $f(x)$

5. $S_k = -\nabla f(x_k)$.

6. Cálculo de $a_k$ (ver texto)

   como resultado se calcula $x_{k+1}$ y $f(x_{k+1})$

7. Calcular $\nabla f(x_{k+1})$.

8. $k \longrightarrow k+1$ ; se repite el procedimiento desde el paso 4.

Para finalizar la discusión sobre el método del máximo descenso es conveniente hacer notar que bajo algunas condiciones, por cierto no muy frecuentes, el algoritmo puede ser atraído por un punto silla ya que también en esta clase de puntos se satisface que $\nabla^T f(x)\ \nabla f(x) = 0$, no exis tiendo forma de detectar a priori que tal cosa sucederá. Ahora bien, para clasificar el tipo de punto donde se detuvo el algoritmo, es necesario anali zar la matriz de segundas derivadas, de acuerdo a :

i) H, positiva definida — mínimo

ii) H, negativa definida — máximo

iii) H, semi-positiva o semi-negativa definida — punto silla.

## 2.2. Método de Newton

El método de Newton que a continuación se presenta esta basado en una aproximación a segundo orden de la función objetivo $f(x)$, es decir, usa información de segundo orden (matriz hessiana), de aqui que se le clasifique como método de segundo orden.

Considérese la aproximación a segundo orden de $f(x)$

$$f(x + \Delta x) = f(x) + \nabla^T f(x) \Delta x + \tfrac{1}{2} \Delta x^T \nabla^2 f(x) \Delta x$$

Si se supone que la aproximación anterior es buena, como de hecho lo es en la vecindad de un mínimo y de un máximo, entonces, al derivar parcialmente $f(x + \Delta x)$ con respecto a cada uno de los elementos de $\Delta x$ se obtiene

$$\frac{\partial f(x + \Delta k)}{\partial \Delta x} = \nabla f(x) + \nabla^2 f(x) \Delta x = 0$$

de donde

$$\Delta x = - H^{-1}(x) \nabla f(x)$$

donde $H(x) = \nabla^2 f(x)$ es la matriz hessiana $\left( hij = \frac{\partial^2 f(x)}{\partial x_i \partial x_j} \right)$.

y $\Delta X$ será la dirección del desplazamiento desde un punto $x_k$ a $x_{k+1}$, es decir

$$x_{k+1} = x_k - H^{-1}(x_k) \nabla f(x_k)$$

El empleo de la última fórmula para generar los iterandos del méto-do de Newton puede provocar los siguientes problemas si es que el método se está empleando para minimizar una función $f(x)$. El problema es el siguiente. Note que tanto un máximo como un mínimo, así como los puntos silla, satisfacen $\nabla f = 0$ ~~∂f(x + Δx)/∂x = 0~~, por lo que no se puede garantizar que el método converga a un mínimo, si no se modifica adecuadamente Esta modificación consiste en lo siguiente :

sea $\Delta x_k = - H^{-1}(x_k) \; \nabla f(x_k)$ la dirección de avance del punto $x_k$ al $x_{k+1}$ y considérase la aproximación a primer orden de $f(x)$ como sigue :

$$f(x_k + \Delta x_k) = f(x_k) + \alpha_k(\nabla^T f(x_k) \; \Delta x_k)\rho_k$$

o bien

$$\Delta f_k = f(x_k + \Delta x_k) - f(x_k) = \alpha_k ( \nabla^T f(x_k) \; \Delta x_k) \rho_k$$

donde $\alpha_k$ es el tamaño de paso, descrito en el inciso anterior, y $\rho_k$ es el sentido de la dirección de búsqueda, el cual se determina como sigue : ya que se desea que $\Delta f_k < 0$ ( $f(x_{k+1}) < f(x_k)$ ) y como $\alpha_k > 0$ , entonces

$$\rho_k = \begin{cases} 1 & \text{si } \nabla^T f(x_k)H^{-1}(x_k)\nabla f(x_k) > 0 \\ -1 & \text{si } \nabla^T f(x_k)H^{-1}(x_k)\nabla f(x_k) < 0 \end{cases}$$

Y con esto se garantiza que si $\alpha_k(>0)$ es suficientemente

pequeña entonces $f(x_{k+1}) < f(x_k)$ y el método de Newton convergerá

a un mínimo, o en el peor de los casos a un punto silla, pero nunca a un

máximo (se omite la discusión sobre el tamaño de paso $\alpha_k$ , por ser idén

tica a la presentada con anterioridad).

Tomando en cuenta los puntos anteriores, la implementación del

método de Newton queda como sigue

1. $k = 0$

2. Estimar $x_0$ (punto de arranque)

3. Calcular $f(x)$,

4. Calcular $\nabla f(x_k)$;

5. Si $\nabla^T f(x_k)$ $\nabla f(x_k) \leq$ tolerancia, se ha encontrado

un punto estacionario de $f(x)$ (mínimo o punto silla).

6. Calcular $H(x_k)$ y $H^{-1}(x_k)$.

7.
$$\rho_k = \begin{cases} 1 & \text{si.} \quad \nabla^T f(x_k) \ H^{-1}(x_k) \ \nabla f(x_k) > 0 \\ -1 & \text{si} \quad \nabla^T f(x_k) \ H^{-1}(x_k) \ \nabla f(x_k) < 0 \end{cases}$$

8. $\Delta x_k = - \rho_k H^{-1}(x_k)\nabla f(x_k)$

9. Calcular $\alpha_k$ (inciso 2.1)

como resultado se obtiene $x_{k+1}$, $f(x_{k+1})$

10. $k \longleftarrow k + 1$. Se repite el procedimiento desde (4).

Es fácil ver que si la función objetivo es cuadrática, p. ej.

$$f(x) = a + b^T x + \tfrac{1}{2} x^T Q x$$

el método de Newton converge en una sola iteración, ya que

$$\nabla f(x_k) = b + Q x_k$$

$$y \qquad \nabla^2 f(x_k) = Q$$

entonces

$$x_{k+1} = x_k - \left[\nabla^2 f(x)\right]^{-1} \nabla f(x_k) = x_k - Q^{-1}(b + Qx_k)$$

$$\boxed{Q^{-1} b = x^*}$$

convergencia cuadrática

## 2.3. Conjugancia y Direcciones Conjugadas

Como se verá más adelante, una función objetivo cuadrática de n - variables que exhibe un mínimo, puede ser minimizada en n pasos (o menos) si estos pasos se toman en lo que se ha llamado direcciones conjugadas. En el inciso siguiente se limitará la discusión a funciones cuadráticas del tipo

$$f(x) = a + b^T x + \tfrac{1}{2} x^T H x \qquad\qquad (o)$$

donde H es una matriz positiva definida.

### 2.3.1. Conjugancia

Supóngase que la minimización de $f(x)$ empieza en $x_o$ en la dirección $\bar{S}_o$, escogida arbitrariamente (o por algún algoritmo); se supondrá $(\bar{S}_o)^T \bar{S}_o = 1.0$. El siguiente punto generado por el algoritmo será

$$x_1 = x_o + \lambda_o \bar{S}_o \qquad\qquad (1)$$

donde el tamaño de paso $\lambda_o$ se determina minimizando $f\left(x_o + \lambda_o \bar{S}_o\right)$ con respecto a $\lambda$, es decir

$$\frac{df(x_o + \lambda_o \bar{S}_o)}{d\lambda} = 0 = \nabla^T f(x_o)\bar{S}_o + (\bar{S}_o)^T \nabla^2 f(x_o)\left(\bar{S}_o \lambda\right)$$

de donde

$$\lambda_o = -\frac{\nabla^T f(x_o)\bar{S}_o}{\bar{S}_o^T \nabla^2 f(x_o)\bar{S}_o} \qquad\qquad (2)$$

Una vez que se encuentra el siguiente iterando, $x_1$, se deberá seleccionar una nueva dirección de búsqueda para la minimización de $f(x)$. La nueva dirección $\bar{S}_1$ se dice que es <u>conjugada</u> con $\bar{S}_o$ si

$(\bar{S}_1)^T \nabla^2 f(x_o)\bar{S}_o = 0$. (En general, un conjunto de n direcciones independientes de búsqueda $s_o$, $s_1$, ....$s_{n-1}$ son conjugadas con respecto a una matriz Q, positiva definida, si

$$S_i^T QS_j = 0 \qquad 0 \le i \ne j \le n-1 \qquad (3)$$

Q, podría ser, por ejemplo, la matriz hessiana de la función objetivo H. Note además que si $Q = I$, la matriz unitaria, conjugancia y ortogonalidad son sinónimos).

Debido a que los vectores $S_i$, son linealmente independientes además de conjugados, cualquier vector $V \in E^n$, se puede representar en términos de aquellos, como:

$$V = \sum_{j=0}^{n-1} \nu_j S_j$$

$$\nu_j = \frac{S_j^T H(x) V}{S_j^T H(x) S_j}$$

Otra relación importante que se utilizará más adelante es la siguiente. Considérese la matriz $P$ definida por

$$P = \sum_{j=0}^{n-1} a_j S_j S_j^T$$

Resulta obvio que si las $a_j$ se escogan de manera tal que

$$P H S_k = S_k \quad \text{mayúscula}$$

entonces $P = H^{-1}$. Ahora bien.

$$P \, H \, S_k = \left( \sum_{j=0}^{n-1} \alpha_j S_j S_j^T \right) H S_k = \sum_{j=0}^{n-1} \alpha_j S_j (S_j^T H S_k)$$

$$= \alpha_k S_k \; (S_k^T \, H \, S_k)$$

de donde si

$$\alpha_k = (S_k^T \, H \, S_k)^{-1}, \text{ entonces } P = H^{-1},$$

es decir

$$H^{-1} = \sum_{j=0}^{n-1} \frac{S_j \, S_j^T}{S_j^T \, H \, S_j}$$

Si las direcciones de búsqueda empleadas en la minimización de $f(x)$ se escogen conjugadas (esto se demostrará más adelante) a continuación se demuestra que: cualquier función cuadrática de $n$ variables que exhibe un mínimo, puede ser minimizada en $n$ pasos, si se emplean direcciones conjugadas, una dirección diferente en cada paso. Además, el orden en que se usan las direcciones de búsqueda es irrelevante para alcanzar el mínimo.

Demonstración   Sea $f \cdot (x) = a + b^T x + \frac{1}{2} x^T H x$, $\nabla f(x) = b + Hx$

$\nabla f(x) = H$, y en el mínimo de $x$, $\nabla f(x^*) = 0$, es decir $x^* = -H^{-1} b$

Para la $n$-ésima etapa se tiene; usando (1) y (2), que

$$x_n = x_0 + \sum_{k=0}^{n-1} \lambda_k \bar{S}_k$$

y como en cada etapa se usó el valor óptimo de $\lambda_k$, dado por la ecuación (2),

$$x_n = X_o - \sum_{k=0}^{n-1} \frac{(\bar{S}_k)^T \nabla f(x_k)}{(\bar{S}_k)^T H \bar{S}_k} \bar{S}_k \qquad (5a)$$

Por otro lado

$$(\bar{S}_k)^T \nabla f(x_k) = (\bar{S}_k)^T (Hx_k + b)$$
$$= (\bar{S}_k)^T \left[ H(X_o + \sum_{i=1}^{n-1} \lambda_i \bar{S}_i) + b \right]$$
$$= (\bar{S}_k)^T (Hx_o + b) \Bigg\} \text{ por conjugancia de las } \bar{S}_i$$

Entonces

$$x_n = x_o - \sum_{k=0}^{n-1} \frac{(\bar{S}_k)^T (H x_o + b) \bar{S}_k}{(\bar{S}_k)^T H \bar{S}_k} \qquad (5b)$$

Usando la relación (4 a) se obtiene que

$$X_o = \sum_{k=0}^{n-1} \frac{(\bar{S}_k)^T H x_o}{(\bar{S}_k)^T H \bar{S}_k} \bar{S}_k$$

y por lo tanto

$$x_n = \sum_{k=0}^{n-1} \frac{(\bar{S}_k)^T b \bar{S}_k}{(\bar{S}_k)^T H \bar{S}_k} = - \sum_{k=0}^{n-1} \frac{(\bar{S}_k)^T H (H^{-1} b) \bar{S}_k}{(\bar{S}_k)^T H \bar{S}_k}$$

y usando (4 a) nuevamente, se obtiene

$$x_n = -H^{-1}b$$

l.q.e.d.

Un método para el cual se garantiza que alcanza el mínimo de una función objetivo cuadrática en un número específico de pasos, se dice que tiene la propiedad de terminación cuadrática. (El método de gradiente conjugado necesita de n pasos, mientras que el de Newton uno sólo).

### 2.3.2. Método de Gradiente conjugado

El método de Fletcher - Reeves de gradiente conjugado, que a continuación se describe, genera una secuencia de direcciones de búsqueda que son combinación lineal de $- \nabla f(x_k)$, la dirección de máximo descenso en el último punto, y de las k direcciones de búsqueda anteriores, $s_0$, $s_1$, ..., $s_{k-1}$, usando factores de peso $\alpha_k$ tales que $S_k$ sea conjugada a las direcciones anteriores.

Para ilustrar el método, sea $S_0 = - \nabla f(x_0)$, y

$$x_1 = x_0 + \lambda_0^* S_0 \; ; \; sea$$

$$S_1 = - \nabla f(x_1) + \alpha_1 S_0 \qquad (6)$$

donde $\alpha_1$ se escoge de forma tal que $S_0$ y $S_1$ sean conjugadas con respecto a H, es decir

$$S_0^T H S_1 = 0 \qquad (7)$$

Para eliminar $S_o$, considérese la aproximación a primer orden del gradiente, es decir

$$\nabla f(x_1) - \nabla f(x_o) = \nabla^2 f(x_o)(x_1 - x_o)$$

$$= \lambda_o^* H S_o \tag{8}$$

$$S_o = \frac{1}{\lambda_o^*} H^{-1}\left[\nabla f(x_1)\nabla f(x_o)\right]$$

y como $H$ es simétrica, entonces

$$S_o^T = \frac{\left[\nabla f(x_1) - \nabla f(x_o)\right]^T H^{-1}}{\lambda_o^*} \tag{9}$$

Substituyendo (6) y (9) en (7) se obtiene

$$\left[\nabla f(x_1) - \nabla f(x_o)\right]^T \left(-\nabla f(x_1) + \alpha_1 S_o\right) = 0 \tag{10}$$

ya que, según se vió con anterioridad $\nabla^T f(x_o)\,\nabla f(x_1) = \nabla^T f(x_1) S_o = 0$, entonces

$$\alpha_1 = -\frac{\nabla^T f(x_1)\nabla f(x_1)}{\nabla^T f(x_o)S_o}$$

o bien,

$$\alpha_1 = -\frac{\nabla^T f(x_1)\nabla f(x_1)}{\nabla^T f(x_o)\nabla f(x_o)} \tag{11}$$

La dirección de búsqueda $S_2$ se forma como una combinación lineal de $-\nabla f(x_2)$, $S_1$ y $S_o$, y se fuerza a que sea conjugada a $S_1$ y $S_o$

con lo que se obtiene la siguiente expresión para los factores de peso $\alpha_k$

$$\alpha_k = \frac{\nabla^T f(x_k) \, \nabla f(x_k)}{\nabla^T f(x_{k-1}) \, \nabla f(x_{k-1})} \qquad (12)$$

La implementación del algoritmo de gradiente conjugado de Fletcher - Reeves incluye los siguientes pasos :

1. Estimar $x_o$ (punto de arranque)

2. $S_o = . - \nabla f(x_o)$

3. En la $k$ - ésima etapa del algoritmo, se determina el mínimo unidireccional de $f(x)$, a lo largo de la dirección de búsqueda $S_k$. Con esto se localiza $x_{k+1}$.

4. La nueva dirección de búsqueda se determina como

$$S_{k+1} = - \nabla f(x_{k+1}) + \frac{\nabla^T f(x_{k+1}) \, \nabla f(x_{k+1})}{\nabla^T f(x_k) \, \nabla f(x_k)} \, S_k$$

Después de $(n+1)$ iteraciones $(k = n)$, se empieza un nuevo ciclo del algoritmo, es decir, $x_{n+1}$ se convierte en $x_o$.

5. La búsqueda se da por terminada cuando en alguna iteración sucede que

$$S_k^T \, S_k \leq \text{tolerancia}$$

Note que al igual que en el método de gradiente ordinario, no se necesita la inversión de matriz alguna, lo cual es una ventaja.

## 2.4. Métodos de Métrica Variable

Los métodos de Métrica Variable o Cuasi-Newton son métodos que aproximan el hessiano, o su inversa, usando unicamente información acerca del gradiente. La mayoria de estos métodos usan direcciones conjugadas, con forme avanzan, lo cual hacen siguiendo el esquema general

$$x_{k+1} = x_k + \lambda_k \dot{s}_k = x_k - \alpha_k \eta(x_k) \nabla f(x_k) \tag{1}$$

donde $\eta(x_k)$ representa una aproximación a $H^{-1}(x_k)$. (En el método de Gradiente ordinario $\eta(x_k) = I$, mientras que el método de Newton toma $\eta(x_k) = H^{-1}(x_k)$, con la desventaja de que hay que invertir el hessiano).

En una serie de métodos de Cuasi-Newton, $H^{-1}(x_{k+1})$ se aproxima de la información disponible en la $k$-ésima etapa, como

$$H^{-1}(x_{k+1}) \approx \omega \eta_{k-1} = \omega(\eta_k + \Delta\eta_k) \tag{2}$$

donde $\eta$ es una aproximación a $H^{-1}$, $\Delta\eta_k$ es una matriz que se especifica de acuerdo al método, y $\omega$ una constante de escalamiento que frecuentemente se fija en $1$. La selección de $\Delta\eta_k$ determina, esencialmente, el método

de método variable. Para garantizar convergencia, $\omega n_{k+1}$ debe ser posi tiva definida y debe satisfacer la siguiente relación cuando reemplaza a $H^{-1}$

$$x_{k+1} - x_k = H^{-1}(x_k)\left[\nabla f(x_{k+1}) - \nabla f(x_k)\right] \qquad (3\ a)$$

que es una aproximación a primer orden del gradiente.

En la $(k+1)$ - ésima etapa de cualquier método se conocen

$$x_k, \quad \nabla f(x_k), \quad x_{k+1}, \quad \nabla f(x_{k+1}) \quad y \quad n_k \quad , \quad y \text{ se}$$

desea calcular $n_{k+1}$

De la relación obtenida con (2) y (3) como

$$n_{k+1}\Delta g_k = \frac{1}{\omega}\Delta x_k \qquad (3\ b)$$

donde $\Delta g_k = \nabla f(x_{k+1}) - \nabla f(x_k)$. Sea $\Delta n_k = n_{k+1} - n_k$, por lo que la ecuación

$$\Delta n_k \Delta g_k = \frac{1}{\omega}\Delta x_k - n_k \Delta g_k \qquad (3\ c)$$

debe ser resuelta para $\Delta n_k$ , y esta se obtiene como sigue. Si el lado derecho de (3 c) se multiplica y divide por $y^T \Delta g_k$ , el primer término, y $z^T\Delta g_k$ el segundo término se obtiene que:

$$\left[\Delta n_k - \left(\frac{1}{\omega}\frac{\Delta x_k y^T}{y^T \Delta g_k} - \frac{n_k \Delta g_k z^T}{z^T \Delta g_k}\right)\right]\Delta g_k = 0$$

o bien,

$$\Delta\eta_k = \frac{1}{w}\left[\frac{\Delta x_k y^T}{y^T \Delta g_k} - \frac{\eta_k \Delta g_k\, z^T}{z^T \Delta g_k}\right] \qquad (4)$$

donde los vectores columna y, Z, son arbitrarios, al igual que w. Si

por ejemplo se escogen

$$w = 1$$

$$y = Z = \Delta x_k - \eta_k \Delta g_k$$

se genera el algoritmo de Broyden, mientras que si se escogen

$$y = \Delta x_k$$

$$Z = \eta_k \Delta g_k$$

entonces la matriz $\eta_{k+1}$ se actualiza de acuerdo al método de Davidon –

Fletcher – Powell. Ya que los vectores y, Z son arbitrarios se pueden

efectuar varias selecciones, las que se discuten más adelante. Si los

pasos $\Delta x_k$ se determinan mediante minimizaciones unidireccionales de

f (x) en la dirección $S_k$, todos los métodos que calculan una $\eta_{k+1}$

simétrica que satisfagan (3 b), generan direcciones que son mutuamente

ortogonales (para funciones cuadráticas)..

## 2.4.1. $\Delta\eta_k$ de Rango 1

Broyden demostró que si $\Delta\eta_k$ es simétrica con rango 1,

la relación $\eta_{k+1}\Delta g_k = \Delta x_k$ se satisface, la única posibilidad

de escoger $\Delta\eta_k$ es

$$\Delta \eta_k = \frac{[\Delta x_k - \eta_k \Delta g_k][\Delta x_k - \eta_k \Delta g_k]^T}{[\Delta x_k - \eta_k \Delta g_k]^T \Delta g_k} \qquad (5)$$

El algoritmo funcionaría de la siguiente forma. Se escogen $x_o$ y $\eta_o > 0$ y se usan una forma secuencial (1), (5) y (2) hasta que por ejemplo , $\nabla^T f(x_k) \nabla f(x_k) \leq \varepsilon$ . Por otro lado, si se usan minimizaciones unidireccionales, el método genera direcciones conjugadas y bajo algunas condiciones más o menos restrictivas, se puede demostrar que el algoritmo converge a la solución. Una característica atractiva de este método es que $\alpha_k$ en (1) no necesariamente tiene que ser un parámetro que minimize $f(x)$ a lo largo de $S_k$. El mismo Broyden demuestra que $\alpha$ puede tomar cualquier valor con la única condición de que no provoque que $\eta$ se haga singular (denominador en (5)).

Si la función objetivo no es cuadrática, algunos de los aspectos poco satisfactorios al usar (5), son los siguientes :

1. $\eta$ puede dejar de ser positiva definida, en cuyo caso es necesario recurrir a alguna otra estrategia que lo garantice.

2. La corrección $\Delta \eta_k$ puede no quedar acotada (generalmente por errores de redondeo, incluso para funciones cuadráticas)

3. Si por coincidencia $\Delta x_k = - \alpha_k \eta(x_k) \nabla f(x_k)$ queda en la dirección del paso anterior, $\eta(x_{k-1})$ se vuelve singular.

E consecuencia, en el algoritmo de Broyden, si sucede que

$$n_k \Delta g_k = \Delta x_k$$

o

$$(n_k \Delta g_k - \Delta x_k)^T \Delta g_k^* = 0$$

se fuerza a que

$$n_{k-1} = n_k \qquad (\Delta n_k = 0)$$

### 2.4.2. Método de Davidon – Fletcher – Powell

En este método la matriz $\Delta n$ se escoge que tenga rango 2. La $n$ inicial normalmente se toma como $n = I$. (se puede usar cualquier otra matriz simétrica positiva definida), con lo que el método arranca con la dirección del máximo descenso. Conforme el método avanza, va existiendo un cambio del método de gradiente a Newton con lo que se obtiene una gran ventaja al usar las mejores características de ambos métodos.

Como se mencionó con anterioridad la relación para $\Delta n_k$ en el método de Davidon – Fletcher – Powell, $y_k = \Delta x_k$ y $Z_k = n_k \Delta g_k$ con lo que al substituir en (4) se obtiene

$$n_{k-1} = n_k + A_k + B_k$$

$$= n_k + \frac{\Delta x_k (\Delta x_k)^T}{(\Delta x_k)^T \Delta x_k} - \frac{n_k \Delta g_k (\Delta g_k)^T n_k^T}{(\Delta g_k)^T n_k \Delta g_k} \tag{5}$$

en donde las matrices $A_k$ y $B_k$ son simétricas y si, además, $n_k$ es

también simétrica, entonces $\eta_{k+1}$ también lo será. La relación anterior
(Ec. (5)) produce resultados satisfactorios en la práctica siempre y cuando

1. El error al evaluar $\nabla f\left(x_k\right)$ no sea grande

2. $\eta_k$ no se haga mal-condicionada

El papel de la matriz $A_k$ en la ecuación (5) es garanti
zar que $\eta \longrightarrow H^{-1}$, mientras que la matriz $B_k$ garantiza que
$\eta_{k+1}$ sea positiva definida en todos los pasos, y en el límite se cancela
con $\eta_0$. Esto se puede ver como sigue

$$\eta_1 = I + A_o - B_o$$

$$\eta_2 = \eta_1 + A_1 - B_1 = I + (A_o + A_1) - (B_o - B_1)$$

$$\eta_{k+1} = I + \sum_{i=0}^{k} A_k - \sum_{i=0}^{k} B_k$$

Para una función cuadrática la suma de las matrices $A_i$ debe ser igual a
$H^{-1}$ cuando $k = n - 1$, y la suma de las matrices $B_i$ deberá cancelar $\eta_o$
(I en este caso), se puede decir que el método de Davidon – Fletcher – Powell
refleja, en cierta forma, toda la información ganada en iteraciones anterio
res, a través de $\eta$.

Debe señalarse que el método que se está ~~descubriendo~~ describiendo usa
direcciones conjugadas si la función objetivo es cuadrática. Para que la
última dirección, $S_{n-1}$, sea conjugada a todas las anteriores, se debe

cumplir que :

$$X_{n-1}^T \; H \, S_{n-1} = 0$$

si se substituye que $S_{n-1} = - \eta_{n-1} \; \nabla f (x_{n-1})$, entonces

$$X_{n-1}^T \; H \, \eta_{n-1} \; \nabla f (x_{n-1}) = 0 \qquad (6)$$

donde $X_{n-1} = \left[ \Delta x_0 , \Delta x_1 , \ldots , \Delta x_{n-1} \right]^T$. Si $H \, \eta_{n-1} = I$

( $\eta_{n-1} = H^{-1}$ ), entonces $\nabla f (x_{n-1})$ es conjugada a todas las direcciones

de búsqueda anteriores dadas por $\Delta x_0 , \Delta x_1 , \ldots , \Delta x_{n-1}$. Sabiendo que

todas las direcciones de búsqueda son conjugadas, se puede demostrar que

$$\sum_{i=0}^{n-1} A_i = H^{-1}, \; \text{como sigue.} \quad \text{Ya que} \quad \Delta g_k = H \, \Delta x_k , \; \text{entonces el numerado}$$

y denominador de cada $A_i$ es

$$( \Delta x_k ) ( \Delta x_k )^T = ( \alpha_k \, S_k ) ( \alpha_k \, S_k )^T = \alpha_k^2 \; S_k \, S_k^T$$

$$( x_k )^T \; g_k = ( \alpha_k \, S_k )^T ( H \alpha_k \, S_k ) = \alpha^2 \; S_k^T \; H \; S_k$$

de donde

$$\sum_{i=0}^{n-1} A_i = \sum_{i=0}^{n-1} \frac{S_i S_i^T}{S_i^T H S_i} = H^{-1} \qquad (7)$$

que es la fórmula (4 b). Séc. 2,3, obtenida con anterioridad.

Para terminar la presentación de este método, se hacen dos
cementarios sobre la implementación práctica del mismo

1. En algunos problemas, los métodos de métri y variable
   fallan en alcanzar el mínimo de la función objetivo si
   el grado de precisión en la búsqueda unidimencional no
   es suficientemente $fina$. Se recomienda que la pre
   cisión en la búsqueda unidireccional sea al menos equi
   valente que en que se requiere para detener al algoritmo
   completo.

2. La búsqueda por el mínimo se debe detener, si al evaluar
   los vectores $-n_k \, \delta f(x_k)$ y $-a_k \, n_k \, \nabla f(x_k)$
   ocurre cualquiera de los dos siguientes puntos.

   a) <u>Cada</u> componente en ambos vectores es menor que una
      tolerancia dada.

   b) Las longitudes predichas al mínimo, de cualquiera
      de los dos vectores es inferior a una cierta tole
      rancia.

### 2.4.3. Algoritmos de Pearson

Pearson propuso una serie de algoritmos para calcular $\eta$, usando direcciones que fueran conjugadas. Los algoritmos de Pearson se pueden obtener empleando diferentes vectores y, Z en la ecuación (4) del inciso 2.4., según se muestra a continuación

1. <u>Pearson No. 2</u> Sea $y = Z = \Delta x_k$ y $w = 1$. Entonces

$$\eta_{k+1} = \eta_k + \frac{(\Delta x_k - \eta_k \, \Delta g_k) \cdot (\Delta x_k)^T}{(\Delta x_k)^T \, \Delta g_k}$$

$$\eta_0 = R_0$$

donde $R_0$ es cualquier simétrica positiva definida. Este algoritmo generalmente conduce a matrices mal condicionadas.

2. <u>Pearson No. 3</u> Sea $y = Z = \eta_k \Delta g_k$, con $w = 1$. Entonces

$$\eta_{k+1} = \eta_k + (\Delta x_k - \eta_k \, \Delta g_k) \frac{(\eta_k \Delta g_k)^T}{(\Delta g_k)^T \eta_k \, \Delta g_k}$$

$$\eta_0 = R_0$$

Este algoritmo se comporta bastante parecido al de Davidon-Fletcher-Powell, excepto que el tamaño del paso es, en general, inferior al de este último.

### 3. Newton – Raphson Proyectado.

Pearson propuso este otro algoritmo al cual se puede obtener haciendo que $\chi \longrightarrow \infty$ y $Z = n_k \, \Delta g_k$, con lo que se obtiene.

$$n_{k+1} = n_k - \frac{(n_k \, \Delta g_k)(n_k \, \Delta g_k)^T}{(\Delta g_k)^T n_k \, \Delta g_k}$$

$$n_0 = R_0$$

Este método incluye la siguiente regla de reinicio cada n etapas, donde n es el número de variables independientes, sobre la matriz $n_k$.

$$R_{k-1} = R_k + \frac{(\Delta x_k - R_k \Delta g_k)(n_k \Delta g_k)^T}{(\Delta g_k)^T n_k \Delta g_k}$$

es decir, cada n etapas se toma $n_k = R_k$.

# BIBLIOGRAFIA

1. C.G. Broyden, Math. Computation, $\underline{21}$ : 368 (1967)

2. M.J. Box, D. Davis, and W.H. Swann. "Non-linear Optimization Techniques", Chemical Industries Monograph 5, Oliver and Boyd, Edinburgh, 1970

3. G.E. Coggins, Univariate Search Methods, Imperial Chemical Industries, Ltd., Central Instr. Lab. Res., Note 64/11, 1964

4. W.C. Davidon, USAEC Doc. ANL-5990 (Rev.), Nov., 1959

5. W.C. Davidon, Computes J., 10 : 406 (1968) ; Chap. 2 in R. Fletcher (ed.), "Optimization" Academic Press Inc. N.Y., 1969

6. R. Fletcher and M.J.D. Powell, Computer J., $\underline{6}$ : 163 (1963)

7. R. Fletcher and C.M. Reeves, Computer J., $\underline{7}$ : 149 (1964)

8. M.R. Hestenes, The Conjugate Gradient Method for Solving Linear Systems, in Proceedings of the Symposium on Applied Mathematics, Vol. VI, Mc. Graw-Hill Book Company, New York, 1956, pp. 83-102.

9. M.R. Hestenes and E.L. Steifel, J. Res. Nat. Bur. Std., B 49 : 409 (1952)

10. J.D. Pearson, Computer J., $\underline{13}$ : 171 (1969)

11. M.J.D. Powell, Computer J., $\underline{7}$ : 155 (1964)

12. M.J.D. Powell; Rank One Methods for Unconstrained Minimization, AERE Rept., TP 372, 1969

13. D.J. Wilde, "Optimum Seeking Methods", Prentice Hall, Inc., Englewood Cliffs, N.J., 1964.

_o_o_o_o_o_o_o_o_o_o_o_o_o_o_c_

## OPTIMIZACION DE UN TREN DE CAMBIADORES DE CALOR

Se desea calentar un fluido desde una temperatura $T_{in}$ hasta una temperatura de salida $T_{out}$, mediante intercambio de calor con tres corrientes liquidas calientes, en un sistema de tres cambiadores de calor que operan en contra corriente, según se muestra en la figura



Para el proceso anterior se especifican los siguientes parámetros :

$WC_A$ : producto del flujo en masa por la capacidad calorifica (se supone idéntica en todas las corrientes)

$t_i$ : temperatura de entrada de las corrientes calientes $(i = 1, 2, 3)$

$U_i$ : coeficiente total de transferencia de calor en cada uno de los tres cambiadores $( i = 1, 2, 3)$

Si se supone que la inversión total requerida para el sistema de cambiadores es proporcional al area total de los mismos, el problema se puede plantear como el de seleccionar las areas $A_i$ $(i = 1,2,3)$ de manera

tal que :

$$A_T = A_1 + A_2 + A_3$$

sea mínima.

1°) Formular el problema como una optimización en estado estacionario. Identifique la función objetivo, variables de decisión y restricciones. Los siguientes parámetros se consideran fijos

| | | |
|---|---|---|
| Tin = To = 100 | $t_1$ = 300 | $U_1$ = 120 |
| Tout = $T_3$ = 500 | $t_2$ = 400 | $U_2$ = 80 |
| WC = 100,000 | $t_3$ = 600 | $U_3$ = 40 |

2°) Suponga que en el problema anterior las corrientes de salida caliente de los cambiadores 1 y 2 deben mezclarse y la temperatura resultante no deberá exceder los 230° . Plantee este nuevo problema tomando en cuenta la restricción planteada.

## Modelo del cambiador de calor

Considerése el cambiador de calor a contra corriente que se muestra a continuación

Las ecuaciones que describen su funcionamiento son :

$Q_n$ : rapidez de transferencia de calor en el n-ésimo cambiador

$$= WC \ (T_n - T_{n-1})$$

$$= WC \ (t_n - t_n^*)$$

$$= U_n A_n \ (t_n - T_n)$$

$$= U_n A_n \ (t_n^* - T_{n-1})$$

siendo las temperaturas de salida las siguientes :

$$T_n = \frac{T_{n-1} + \alpha \ t}{1 + \alpha \ h}$$

donde

$$\alpha_n = \frac{U_n \ A_n}{WC}$$

y $\quad t_n^* = t_n - (T_n - T_{n-1})$

## Función objetivo

Según se mencionó con anterioridad, se trata de minimizar el área total del sistema

$$A_T = A_1 + A_2 + A_3$$

## Variables de decisión y restricciones

1- Sin restricciones en el mezclado (Parte 1)

a) $T_1$ y $T_2$ como variables de decisión

$$Q_1 = W C (T_1 - T_0)$$

$$A_1 = Q_1 / U_1 (t_1 - T_1)$$

$$Q_2 = W C (T_2 - T_1)$$

$$A_2 = Q_2 / U_2 (t_2 - T_2)$$

$$Q_3 = W C (T_3 - T_2)$$

$$A_3 = Q_3 / U_3 (t_3 - T_3)$$

Restricciones sobre las variables independientes

$$T_0 \leq T_1 \leq t_1$$

$$T_1 \leq T_2 \leq t_2$$

Restricciones sobre las variables dependientes

$$Q_i \geq 0 \qquad i = 1, 2, 3$$

o  equivalentemente

$$A_i \geq 0 \qquad i = 1, 2, 3$$

b) Areas $A_1$ y $A_2$ como variables independientes

$$T_1 = \frac{T_0 + a_1 t_1}{1 + a_1} \qquad ; \qquad a_1 = U_1 A_1 / W C$$

$$T_2 = \frac{T_1 + a_2 t_2}{1 + a_2} \qquad ; \qquad a_2 = U_2 A_2 / W C$$

$Q_1$, $Q_2$ y $Q_3$ como en el caso anterior, al igual que $A_3$.

Restricciones en las variables independientes

$$0 \leq A_1 \leq A_1^* \qquad \text{(dato)}$$

$$0 \leq A_2 \leq A_2^* \qquad \text{(dato)}$$

restricciones en las variables dependientes :  no hay.

c) tomar técnicas $Q_1$ y $Q_2$ como variables de depende...

$$T_1 = Q_1 / W C + T_o \qquad ; \qquad A_1 = Q_1 / U_1 \ (t_1 - T_1)$$

$$T_2 = Q_2 / W C + T_1 \qquad ; \qquad A_2 = Q_2 / U_2 \ (t_2 - T_2)$$

$$Q_3 = W C \ (T_3 - T_2) \qquad ; \qquad A_3 = Q_3 / U_3 \ (t_3 - T_3)$$

Restricciones sobre las variables independientes

$$0 \le Q_1 \le W C \ (t_1 - T_o)$$

$$0 \le Q_2 \le W C \ (t_2 - T_o)$$

Restricciones sobre las variables dependientes

$$0 \le Q_1 + Q_2 \le W C \ (t_2 - T_o)$$

2- **Con restricciones en la corriente de mezcla (Parte   ).**

Las ecuaciones y restricciones analizadas en la parte 1 siguen

siendo válidas aunque en este caso existe la necesidad de añadir otra res

tricción, según se analiza a continuación



Variable dependiente :  $t_m^* = \dfrac{t_1^* + t_2^*}{2}$

restricciones sobre $t_m$ :

Si $A_1 = \infty$ 　　y $A_2 = \infty$ $\implies$ $t_1^* = T_0$, $T_1 = t_1$

$$t_2^* = T_1, \quad T_2 = t_2$$

entonces $\dfrac{t_1 + T_o}{2} \leq t_m^*$

y por lo tanto $\dfrac{t_1 + T_o}{2} \leq t_m^* \leq 230$ (restriccioón del problema)

⌐ (restricción por área infinita en $A_1$ y $A_2$)

Las soluciones óptimas resultan ser :

**Primera parte** 　　Sin restricción en $t_m^*$

$T_1 = 186.2$, $T_2 = 292.7$, $Q_1 = 8.62 * 10^6$, $Q_2 = 10.64 * 10^6$, $Q_3 = 20.73 \cdot 10^6$

$t_1 = 631.8$ 　$A_2 = 1239.1$, 　$A_3 = 5183.8$, 　$A_T = 7054.8$

$A_1 = 631.8$

**Segunda parte** 　　Con restricción en $t_m^*$

$T_1 = 210.0$, $T_2 = 340.1$, $Q_1 = 11.00 * 10^8$, $Q_2 = 13.01 * 10^6$, $Q_3 = 16.01 \cdot 10^6$

$t_m^* = 229.9$, $A_1 = 1017.9$, $A_2 = 2715.7$, $A_3 = 4009.3$, $A_T = 7743.0$

## AMPLIACION DE UNA PLANTA QUIMICA

El presente estudio tiene el propósito de diseñar un nuevo reactor catalítico y fijar nuevas condiciones de operación para aumentar la capacidad de una planta que produce ES.

La planta fue originalmente diseñada para producir 91 T/D de ES y se pretende aumentar la capacidad para producir 1*0 T/D. Dentro de las varias alternativas analizadas se pretende efectuar la optimización sobre el diagrama que se muestra a continuación:

C1, C2, C3 : calentadores

E1, E2, E3 : cambiadores

R : reactor catalítico

## Bases de diseño

1. Composición de la carga fresca ( % )

   B —— 0.43

   T —— 0.86

   E B —— 98.46  ⟵  (materia prima)

   P E B —— 0.22

2. Composición de la carga total al reactor (sin vapor)

   B —— 0.166

   T —— 1.866

   E B —— 95.217

   E S —— 2.668  ⟵  (producto final)

   P E B —— 0.082

## Descripción del Flujo

El EB fresco se une con la corriente de recirculación proveniente de otra sección de la planta. La corriente resultante se bombea al sistema de precalentamiento de carga, constituido por los cambiadores E1 y E2 ; antes de entrar a estos cambiadores la carga combinada se mezcla con aproximadamente el 9% del vapor total usado en la reacción. La mezcla, parcialmente vaporizada, se alimenta al cambiador E2 a una temperatura de 316°F, saliendo del mismo a 692°F, completamente vaporizada. Del cambiador E2 pasa al E1 de donde sale a 1092°F. El calentamiento final de la mezcla vapor-

...hidrocarburos que suministra en el calentador C1, procuraba... ción de vapor adicional proveniente del cambiador E3 y del calentador C2, de donde salió a una temperatura de 1150° F.

El 21% restante del vapor requerido por el proceso se precalienta en el cambiador E3. Este vapor entra a 366° F y sale del cambiador E3 a 748° F. El vapor asi calentado, parte se alimenta al calentador C3, de donde sale a 1300° F y el resto se alimenta al C2 para más tarde mezclarse con la corriente de vapor-hidrocarburos antes de entrar al calentador C1. Por otro lado, el vapor que sale del calentador C3 se divide en dos corrientes, a saber: una parte sirve para dar la temperatura final y la relación (vapor/hidrocarburos) a la entrada del reactor R, mientras que la otra se usa para elevar la temperatura de los gases de reacción entre los hechos catalíticos del reactor.

La mezcla vapor-hidrocarburos que sale del reactor R a 1119° F se aprovecha para precalentar los hidrocarburos y el vapor que entran al proceso.

(Nota: La mayoría de los datos dados con anterioridad son resultado de la optimización que se describe más adelante).

Datos

Carga fresca : cantidad máxima disponible ,

composición y temperatura

Recirculación : Cantidad / composición y temperatura

Calentador C3 : temperatura de salida

Vapor : temperatura de entrada

ES en producto a separación. : cantidad ( 140 T/D)

Se requiere calcular lo siguiente:

1.. Calcular el volúmen de catalizador en el lecho L1

2.. Calcular el volúmen de catalizador en el lecho L2

3.. La cantidad (%) del flujo de salida del reactor R que pasa por los cambiadores E1 y E3

4.. La cantidad total de vapor

5.. La cantidad de vapor que se inyecta al cambiador E2

6.. La cantidad de vapor que pasa por el calentador C2 para que la mezcla a la salida del calentador C1 sea de 1150°F

7.. La relación (vapor/hidrocarburos) a la entrada del lecho L1 del reactor R y a la entrada del lecho L2 del mismo reactor.


El objetivo es el siguiente

1. Minimizar la cantidad total de catalizador en el reactor R (lechos L1 y L2

2. Que la producción sea la ~~mínima~~ de 140 T/D.

3. Que el consumo de carga fresca sea la menor posible (C)

4. Que el producto a separación esté lo más frio posible (T).

Con lo anterior se puede plantear la siguiente función objetivo :

$$F = K_1 (L_1 + L_2) + K_2 (prod - 100)^2 + K_3 (G) + K_4 (T)$$

donde $K_1$, $K_2$, $K_3$ y $K_4$ son constantes que se usan para "condicionar adecuadamente" la función objetivo.

## Restricciones

Debido a condiciones de operación de los equipos, se deben tomar en cuenta las siguientes restricciones :

$$1000° F \leq \text{Temp. entrada al reactor} \leq 1280° F$$

$$\text{Temp. salida del primer lecho (L1)} \leq \text{Temp. entrada al segundo lecho (L2)} \leq 1250° F$$

$$1.0 \leq \text{Relación (Vap./hidrocarburos) entrada al primer lecho} \leq 3.0$$

$$\text{Relación (Vap./hidrocarburos) entrada al } 1^{er} \text{ lecho} \leq \text{Relación (Vap./hidrocarburos) entrada al } 2^a \text{ lecho} \leq 3.0$$

## Modelo de los equipos

Los equipos que es necesario simular y/o dimensionar son los siguientes :

1. Calentadores    (C1, C2, C3)

2. Cambiadores de calor

   a) Gas – Gas    (E1 y E3)

   b) Gas – Liquido con evaporación  (E2)

3. Reactor químico catalítico ( R )

. De los equipos mencionados sólo se describirá el modelo usado para el reactor catalítico, por ser este el equipo más importante de la planta. Para el resto de los equipos los modelos matemáticos son bastante convencionales (Ver., p.ej., D.D. Kern ("Process Heat Transfer", Mc Graw-Hill Book Co., New York).

## Cinética del sistema reaccionante

Se considera que en los lechos catalíticos  L1  y L2  se llevan a cabo las siguientes reacciones.. ( A : vapor de agua).

1)      E B  $\rightleftharpoons$  E S   +   H           (catalítica)

2)      E B  $\longrightarrow$  B   +   E           (catalítica)

3)   H  +  E B  $\longrightarrow$  T   +   M        (catalítica)

4)   M  +  A  $\longrightarrow$ C O  +  H          (catalítica)

5)   C O +  A  $\longrightarrow$ C O$_2$  +  H      (catalítica)

6)      E  $\longrightarrow$  A C  +  H       (descomposición en fase vapor)

Las expresiones para la velocidad y la presión de cada una de las reacciones anteriores, son las siguientes

$$V_1 = \left[ P_{EB} - \frac{P_{ES} \; P_H}{K_p} \right] \exp \left( - \frac{5715}{T} - 6.16 \right)$$

$$V_2 = \exp \left( - \frac{25600}{T} + 12.8 \right) \; P_{EB}$$

$$V_3 = \exp \left( - \frac{11000}{T} - 1.8 \right) \; P_{EB} \; P_H$$

$$V_4 = \exp \left( - \frac{7900}{T} - 3.36 \right) \; P_M$$

$$V_5 = {}^P \exp \left( - \frac{8850}{T} + 3.8 \right) \; P_{co} \; P_A$$

$$V_6 = 2.5 * 10^4 \; \exp \left( - \frac{38000}{T} \right) \; P_E / T$$

donde:

$$K_p = T^{0.549} \exp \left( - \frac{14516}{T} + 11.41 \right)$$

P : Presión en atmósferas

T : Temperatura $(^\circ K)$

$P_i$ : Presión parcial a cada componente

El calor liberado por cada una de las reacciones se tomó igual a

$$\Delta H_1 = 28,843 + 1.09 \; T$$

$$\Delta H_2 = 25,992 - 1.09 \; T$$

$$\Delta H_3 = 12,702 - 3.15 \; T$$

$$\overline{a}\overline{b} \cdot \overline{b} \cdot \overline{c} = 111,05 \cdot H_a$$

$$T_{(n-1)} = 86,30 \cdot {}_{i}t_a$$

con sus correspondientes velocidades asociadas ( $v = 0$ )

$$\Delta H_4 = 33,015 + 3.96 \ T$$

$$\Delta H_5 = 10,802 + 2.5 \ T$$

$$\Delta H_6 = 38,278 + 11.45 \ T$$

Con la anterior información, es posible establecer las relaciones

que describen la variación, con la longitud, de cada uno de los componentes

así como de la presión y temperatura, las cuales tendrían la forma general

$$\frac{d \ n_i}{d \ z} = f_i \ (n_1, \ n_2, \ \ldots, \ n_{NC}, \ T, \ P) \qquad i = 1,2, \ \ldots, N. \ de$$
componentes (NC)

$$\frac{d \ T}{d \ z} = g \ (n_1, \ n_2, \ \ldots, \ n_{NC}, \ T, \ P)$$

$$\frac{d \ P}{d \ z} = h \ (n_1, \ n_2, \ \ldots, \ n_{NC}, \ T, \ P)$$

con sus condiciones iniciales asociadas ( $z = 0$ )

Hay que tomar en cuenta que el volúmen de catalizador en cada uno

de los lechos es una incognita, por lo que el límite superior de integración

en cada lecho ( $z = z_1$ y $z = z_2$ ) son variables, además de que entre

ambos lechos existe una sección de mezclado con vapor.

## Resultados

Los resultados que produjeron un valor mínimo de la función objetivo

fueron los siguientes :

| | |
|---|---|
| Temperatura de entrada a la primera cama | 1243° F |
| Temperatura de entrada a la segunda cama | 1145° F |
| Longitud de la primera cama | 211.4 cm |
| Longitud de la segunda cama | 157.1 cm |
| (Vapor/hidrocarburos) primera cama | 2,738 lb/lb |
| Gasto de hidrocarburos en la primera cama | 249.31 lbmol/h |
| (Vapor/hidrocarburos) segunda cama | 2,923 lb/lb |
| ES a purificación | 150.4 T/D |
| Conversión total | 65.06 % |
| Selectividad total | 81.84 % |
| Fracción de la salida del reactor a E1 y E3 | 64 % |
| Temperatura del producto a separación | 562° F |
| Fracción del vapor total al cambiador E 2 | 8.3 % |
| Fracción del vapor total al calentador C 3 | 0 % |

## Comentarios Finales

El problema anteriormente descrito resulta ser, desde el punto de vista de minimización de funciones, uno de los más complejos ya que una sola evaluación de la función objetivo requiere de los siguientes cálculos

a) Solución de sistemas de ecuaciones diferenciales ordinarias no lineales

b) Cálculos iterativos en los cambiadores de calor ya que estos existen de antemano y los flujos y temperaturas de operación deben ajustarse al diseño mecánico que tienen.

c) Existe dentro del proceso un paso de recuperación de energía (cambiadores de calor), lo cual también genera que se hagan cálculos iterativos dentro de cada evaluación de la función objetiva.

En opinión del autor de este ejemplo, el modificar las condiciones de operación de una planta para ajustarlas a nuevos requerimientos es, con mucho, bastante más complejo que el diseño de una nueva planta ya que en este último caso se tienen muchos grados de libertad.