



UNIVERSITY *of* LIMERICK
OLLSCOIL LUIMNIGH

College of Informatics and Electronics

END OF SEMESTER ASSESSMENT PAPER

MODULE CODE: MS4327

SEMESTER: Spring 2003

MODULE TITLE: Optimisation

DURATION OF EXAMINATION: 2 1/2 hours

LECTURER: Dr. J. Kinsella

PERCENTAGE OF TOTAL MARKS: 70%

EXTERNAL EXAMINER: Prof. J.D. Gibbon

INSTRUCTIONS TO CANDIDATES: **Answer four questions correctly for full marks; 70%.**

1 The Wolfe conditions for the step length α in a line search require that

$$f(x_k + \alpha p_k) \leq f(x_k) + c_1 \alpha p_k^T g(x_k), \quad (1a)$$

$$p_k^T g(x_k + \alpha p_k) \geq c_2 p_k^T g(x_k) \quad (1b)$$

where $g(x) \equiv \nabla f(x)$ and $0 < c_1 < c_2 < 1$. The **strong** Wolfe conditions replace (1b) by

$$|p_k^T g(x_k + \alpha p_k)| \leq c_2 |p_k^T g(x_k)|. \quad (2)$$

- (a) • **Either** show that: provided that the search direction p_k is a descent direction and f is bounded below along p_k then there exists an interval of α -values such that both Wolfe conditions hold. 15

- **or** prove: 15

Theorem 1 (Zoutendijk) Consider any iteration of the form $x_{k+1} = x_k + \alpha_k p_k$, where p_k is a descent direction and α_k satisfies the Wolfe conditions Eqs. 1a and 1b. Suppose that f is bounded below in \mathbb{R}^n and that f is C^1 in an open set \mathcal{N} containing the level set $\mathcal{L} \equiv \{x : f(x) \leq f(x_0)\}$, where x_0 is the starting point. Also assume that $g(x)$, the gradient of f , is **Lipschitz continuous** on \mathcal{N} , i.e. there exists a constant L such that

$$\|g(x) - g(\bar{x})\| \leq L \|x - \bar{x}\|, \quad \text{for all } x, \bar{x} \in \mathcal{N}. \quad (3)$$

Then

$$\sum_{k \geq 0} \cos^2 \theta_k \|g(x_k)\|^2 < \infty, \quad (4)$$

where θ_k is the angle between p_k and the steepest descent direction $-g(x_k)$.

- (b) Suppose that the search directions p_k are generated using a Newton-like method: $p_k = -B_k^{-1}g(x_k)$ where B_k is symmetric and positive definite. Show that if $\|B_k\| \|B_k^{-1}\| \leq M$ for all k then $\cos \theta_k \leq 1/M$ where θ_k is as defined above. 5

- (c) Use Zoutendijk's Theorem above to show that in this case $\lim_{k \rightarrow \infty} \|g(x_k)\| = 0$. 5

- 2 (a) Show that if the first Wolfe condition (1a) fails for some $\alpha_0 > 0$, then a quadratic approximation $\phi_Q(\alpha) = A_Q\alpha^2 + \phi'(0)\alpha + \phi(0)$ can be formed that interpolates $\phi(0)$, $\phi'(0)$ and $\phi(\alpha_0)$ and has a unique minimum at

$$\alpha_1 = \frac{\phi'(0)\alpha_0^2}{2[\phi(\alpha_0) - \phi(0) - \alpha_0\phi'(0)]} \quad (5)$$

where $\phi(\alpha) \equiv f(x_k + \alpha p_k)$. 5

- (b) Use the fact that the first Wolfe condition (1a) is not satisfied at α_0 to prove that: 3

$$\alpha_1 \leq \frac{1}{2(1 - c_1)}\alpha_0.$$

- (c) Show that $c_1 < \frac{1}{2}$ implies that $\alpha_1 < \alpha_0$. 2

- (d) Suppose that the first Wolfe condition (1a) is not satisfied at α_1 . Show that a cubic approximation $\phi_C(\alpha) = A_C\alpha^3 + B_C\alpha^2 + \phi'(0)\alpha + \phi(0)$ that interpolates $\phi(0)$, $\phi'(0)$, $\phi(\alpha_0)$ and $\phi(\alpha_1)$ can be constructed and that the constants A_C and B_C satisfy: 5

$$\begin{bmatrix} A_C \\ B_C \end{bmatrix} = \frac{1}{\alpha_0^2\alpha_1^2(\alpha_1 - \alpha_0)} \begin{bmatrix} \alpha_0^2 & -\alpha_1^2 \\ -\alpha_0^3 & \alpha_1^3 \end{bmatrix} \begin{bmatrix} D_1 \\ D_0 \end{bmatrix}, \quad (6)$$

where

$$D_0 = \phi(\alpha_0) - \phi(0) - \alpha_0\phi'(0) \quad \text{and} \quad (7)$$

$$D_1 = \phi(\alpha_1) - \phi(0) - \alpha_1\phi'(0). \quad (8)$$

- (e) Show that the minimiser of $\phi_C(\alpha)$ is: 2

$$\alpha_2 = \frac{-B_C + \sqrt{B_C^2 - 3A_C\phi'(0)}}{3A_C}. \quad (9)$$

- (f) Finally, show that $\alpha_2 \in [0, \alpha_1]$. 8

3 The trust region method has at its core the problem:

$$\min_{p \in \mathbb{R}^n} m(p) \equiv f_0 + g^T p + \frac{1}{2} p^T B p, \quad \text{such that } \|p\| \leq \Delta, \quad (10)$$

where f_0 is a fixed scalar, g a fixed vector in \mathbb{R}^n , B a fixed $n \times n$ matrix and Δ a fixed positive scalar. The “dogleg” method finds an approximate solution to (10) by replacing the (unknown) curved trajectory for $p^*(\Delta)$ with a path consisting of two line segments. The first line segment runs from the starting point to the unconstrained minimiser along the steepest descent direction defined by

$$p^U = -\frac{g^T g}{g^T B g} g \quad (11)$$

while the second line segment runs from p^U to $p^B \equiv -B^{-1}g$. We can define the trajectory as a path $\tilde{p}(\tau)$ parameterised by τ as follows:

$$\tilde{p}(\tau) = \begin{cases} \tau p^U, & 0 \leq \tau \leq 1, \\ p^U + (\tau - 1)(p^B - p^U), & 1 \leq \tau \leq 2. \end{cases} \quad (12)$$

(a) Show that if B is positive definite then

(i) $\|\tilde{p}(\tau)\|$ is an increasing function of τ and 5

(ii) $m(\tilde{p}(\tau))$ is a decreasing function of τ . 5

(b) Give a detailed argument which demonstrates that the path $\tilde{p}(\tau)$ intersects the trust region boundary $\|p\| = \Delta$ at exactly one point if $\|\tilde{p}(2)\| \equiv \|p^B\| \geq \Delta$ and nowhere otherwise. 2

(c) What value should be selected for p if $\|p^B\| \leq \Delta$? 3

(d) In the case where the vector p is chosen to be on the boundary, explain carefully how the appropriate value of the parameter τ is chosen. 5

(e) Finally, if $p_U = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$, $p_B = \begin{bmatrix} -1.5 \\ -0.5 \end{bmatrix}$ and $\Delta = 1.5$, find the appropriate value of τ . 5

- 4 “Nearly exact” trust region methods seek to solve the trust region problem (as defined in Eq. 10 above) as accurately as possible. Given that the following Theorem holds:

Theorem 2 *The vector p^* is a global solution of the problem (10) if and only if there is a scalar $\lambda \geq 0$ such that the following conditions are satisfied:*

$$(B + \lambda I)p^* = -g \quad (13a)$$

$$\lambda(\Delta - \|p^*\|) = 0 \quad (13b)$$

$$(B + \lambda I) \text{ is positive semi-definite.} \quad (13c)$$

and defining $p(\lambda) = -(B + \lambda I)^{-1}g$, we need to show that $\|p(\lambda)\| = \Delta$ may be solved for λ . Proceed as follows:

- (a) Use the fact that a symmetric matrix B can be written $B = Q\Lambda Q^T$ to show that: 5

$$p(\lambda) = -Q(\Lambda + \lambda I)^{-1}Q^Tg = -\sum_{j=1}^n \frac{q_j^T g}{\lambda + \lambda_j} q_j, \quad (14)$$

where $\lambda_1 < \lambda_2 < \dots < \lambda_n$ are the eigenvalues of B , $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ and the orthonormal matrix $Q = [q_1 q_2 \dots q_n]$ where the column vectors q_i are the eigenvectors of B .

- (b) Derive a formula for $\|p(\lambda)\|^2$. 3
- (c) Sketch the relevant features of the graph of $\|p(\lambda)\|^2$ for $\lambda \geq -\lambda_1$. 2
- (d) Carefully explain — referring to your graph if you wish but using mathematical arguments — why (provided that $q_1^T g \neq 0$) there must be a single root $\lambda = -\lambda^*$ in the interval $(-\lambda_1, \infty)$. 3
- (e) We could use Newton’s method for root-finding (solve $F(x) = 0$ using $x_{n+1} = x_n - F(x_n)/F'(x_n)$) to solve the equation $\|p(\lambda)\| = \Delta$. Explain why this is not satisfactory and why the equation $1/\|p(\lambda)\| = 1/\Delta$ is a better choice. 2

- (f) Finally show that using Newton's method for root-finding to solve $1/\|p(\lambda)\| = 1/\Delta$ is equivalent to using the procedure:

10

Algorithm 1 (Exact Trust Region)

begin

Given $\lambda_0 > 0, \Delta > 0, \varepsilon > 0$

while $l < l_{max} \wedge \text{abs}(\|p_i(\lambda)\| - \Delta) > \varepsilon$ **do**

Factor $B + \lambda^{(i)}I = R^T R$

Solve $R^T R p_i = -g, R^T q_i = p_i$

$\lambda^{(i+1)} := \lambda^{(i)} + \left(\frac{\|p_i\|}{\|q_i\|}\right)^2 \left(\frac{\|p_i(\lambda)\| - \Delta}{\Delta}\right)$

$i := i + 1$

end

end

Hints:

$$\begin{aligned} \frac{d}{d\lambda} \left(\frac{1}{\|p(\lambda)\|} \right) &= \frac{d}{d\lambda} (\|p(\lambda)\|^2)^{-1/2} \\ &= -\frac{1}{2} (\|p(\lambda)\|^2)^{-3/2} \frac{d}{d\lambda} \|p(\lambda)\|^2 \\ \frac{d}{d\lambda} \|p(\lambda)\|^2 &= -2 \sum_{j=1}^n \frac{(q_j^T g)^2}{(\lambda + \lambda_j)^3}. \end{aligned}$$

Finally,

$$\|q\| = \|R^T p\|^2 = p^T (B + \lambda I)^{-1} p = \sum_{j=1}^n \frac{(q_j^T g)^2}{(\lambda + \lambda_j)^3}.$$

5 The Fletcher-Reeves (FR) version of the non-linear conjugate gradient algorithm is:

Algorithm 2

begin

Given x_0 .

set $r_0 \leftarrow \nabla f_0, p_0 \leftarrow -r_0, k \leftarrow 0$;

while $r_k \neq 0$ **do**

$\alpha_k \leftarrow$ Result of line search along p_k ;

$x_{k+1} \leftarrow x_k + \alpha_k p_k$;

$r_{k+1} \leftarrow \nabla f_{k+1}$

$\beta_{k+1}^{FR} \leftarrow \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k} \equiv \frac{\|r_{k+1}\|^2}{\|r_k\|^2}$;

$p_{k+1} \leftarrow -r_{k+1} + \beta_{k+1}^{FR} p_k$;

$k \leftarrow k + 1$;

end (while)

end

- **Either** prove:

25

Lemma 1 Suppose that Alg. 2 is implemented with a step length α_k that satisfies the strong Wolfe conditions with $0 < c_2 < \frac{1}{2}$. Then the method generates descent directions p_k that satisfy the following inequalities:

$$-\frac{1}{1-c_2} \leq \frac{p_k^T g_k}{\|g_k\|^2} \leq \frac{2c_2-1}{1-c_2}, \quad \text{for all } k = 0, 1, \dots \quad (15)$$

- **or** — using the Dai-Yuan (DY) version of the algorithm which replaces β^{FR} in Alg. 2 by

$$\beta_{k+1}^{DY} = \frac{\|g_{k+1}\|^2}{p_k^T y_k}, \quad (16)$$

where $y_k \equiv g_{k+1} - g_k$ — show that:

25

Theorem 3 (DY-convergence) Under the same assumptions as those of Zoutendijk's Theorem 1, (in particular that the weak Wolfe conditions hold and that f is bounded below) — with the exception that we need **not** assume that the p_k are descent directions, we have that the algorithm either stops at a stationary point ($\|g_k\| = 0$) or $\liminf \|g_k\| = 0$.

- 6 The BFGS method updates an estimate H_k of the inverse Hessian in the form:

$$H_{k+1} = (I - \gamma_k s_k y_k^T) H_k (I - \gamma_k y_k s_k^T) + \gamma_k s_k s_k^T. \quad (17)$$

25

Derive this update rule by solving the problem:

$$\min_H \|H - H_k\| \quad (17a)$$

$$\text{subject to } H = H^T, H y_k = s_k \quad (17b)$$

A choice of norm that allows easy solution of 17a is the “weighted Frobenius norm”:

$$\|A\|_W \equiv \|W^{\frac{1}{2}} A W^{\frac{1}{2}}\|_F, \quad (18)$$

where $\|C\|_F^2 \equiv \sum_{i=1}^n \sum_{j=1}^n C_{ij}^2$. The weight W can be chosen as **any** positive definite matrix satisfying $W s_k = y_k$.