**Solutions**

**1 a)** Denote $r = Ax - b$; $\quad x = \arg\min \|r\|_2 = \arg\min r^T r = \arg\min(x^T A^T A x - 2(A^T b)^T x + b^T b)$

Differentiating with respect to x and setting the result to zero, we get $2A^T Ax - 2A^T b = 0$, or

$$A^T Ax = A^T b \qquad\qquad (2)$$

$A$ is of full column rank $\Rightarrow A^T A$ is nonsingular $\Rightarrow$ (2) has a unique solution.

$A^T A$ is symmetric and positive definite $\Rightarrow$ we can compute its Cholesky factorization

$$A^T A = LL^T,$$

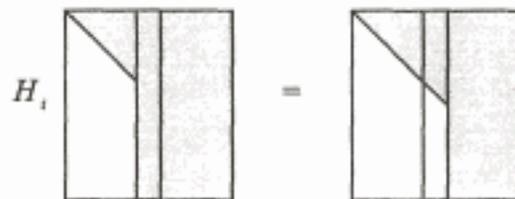where $L$ - lower triangular. This allows us to reduce (2) to solving the triangular systems

$$Ly = A^T b, \qquad L^T x = y.$$

**b)** First we want to decompose $A$ into the product $Q\begin{bmatrix} R \\ 0 \end{bmatrix}$, $Q : m \times m$ orthogonal, $R : n \times n$ upper-triangular; since $A$ is of full column rank, $R$ is nonsingular.

To reduce A to the upper-triangular form, we successively apply Householder transformations

$$H_n \ldots H_1 A = \begin{bmatrix} R \\ 0 \end{bmatrix}, \qquad \text{where } H_i = I - 2\frac{v_i v_i^T}{v_i^T v_i} \qquad \Rightarrow \quad \text{orthogonal and symmetric.}$$

$H_i$ reflects a vector against the hyperplane $v_i^\perp$. We choose $v_i$ so that $H_i$ zeros out the subdiagonal part of the ith column $\tilde{a}_i$ of the current state of $A$ - $(H_{i-1} \ldots H_1 A)$



If $\tilde{a}_i' = \{\tilde{a}_i$ with the upper i-1 entries set to 0\}, then $v_i = \tilde{a}_i' \pm \|\tilde{a}_i'\|_2 e_i$, where the sign is chosen so as to avoid cancellation.

Now, using this decomposition $Q^T A = \begin{bmatrix} R \\ 0 \end{bmatrix}$, $Q^T = H_n \ldots H_1$,

$$x = \arg\min \|Ax - b\|_2 = \arg\min \|Q^T Ax - Q^T b\|_2 = \arg\min \left\| \begin{bmatrix} R \\ 0 \end{bmatrix} x - Q^T b \right\|_2;$$

let $Q^T b = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$, $\quad b_1 : n \times 1$, $b_2 : (m-n) \times 1$;

$$x = \arg\min \left\| \begin{bmatrix} R \\ 0 \end{bmatrix} x - \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \right\|_2^2 = \arg\min(\|Rx - b_1\|_2^2 + \|b_2\|_2^2) = \arg\min \|Rx - b_1\|_2 = R^{-1} b_1$$

**c)**

| | Computations | Accuracy |
|---|---|---|
| Normal equations | $A^T A$ - $\approx mn^2$ flops;<br><br>Cholesky factorization - $\approx \dfrac{n^3}{3}$ flops<br><br>Triangular systems - $O(n^2)$<br><br>$\approx mn^2 + \dfrac{n^3}{3}$ flops | relative error in $x$ is proportional to $(cond(A))^2$ |
| Householder transformations | $\approx 2mn^2 + \dfrac{2}{3}n^3$ flops | relative error in $x$ is proportional to $cond(A) + \|r\|_2 (cond(A))^2$ |

For nearly square problems, $m \approx n$, the two methods require about the same amounts of work, but for $m \gg n$ normal equations method is about 2 times cheaper than Householder method. On the other hand, the Householder method is more accurate.

**2 a)** Via elementary symbolic calculations we obtain $y(t) = e^{\lambda t}$. As $t \to \infty$, $y(t) \to +0$.

**b)** $y_{k+1}(1 - \dfrac{\lambda h}{2}) = y_k(1 + \dfrac{\lambda h}{2}),$ $\qquad y_{k+1} = \dfrac{1 + \lambda h/2}{1 - \lambda h/2} y_k,$ $\qquad y_k = \left(\dfrac{1 + \lambda h/2}{1 - \lambda h/2}\right)^k y_0$

$y_k \to 0 \quad \Leftrightarrow \quad \left|\dfrac{1 + \lambda h/2}{1 - \lambda h/2}\right| < 1, \qquad (2 + \lambda h)^2 < (2 - \lambda h)^2, \qquad 4\lambda h < -4\lambda h, \qquad \lambda h < 0,$

and since we assume $h > 0$, this is true for all negative $\lambda$.

**c)** We have $\qquad y_{k+1} = y_k + \dfrac{f(t_k, y_k) + f(t_{k+1}, y_{k+1})}{2} h \qquad\qquad (1)$

For the true solution $y(t)$ we can write

$$y(t_{k+1}) = y(t_k) + \dfrac{y'(t_k) + y'(t_{k+1})}{2} h + \xi, \qquad\qquad (2)$$

where $\xi$ is some unknown term. From Taylor expansion we have two formulas:

$\dfrac{f(x+h) + f(x-h)}{2} = f(x) + \dfrac{\theta h^2}{2}, \qquad |\theta| \le \|f''\|_c ;$

$\dfrac{f(x+h) - f(x-h)}{2h} = f'(x) + \dfrac{\theta_1 h^2}{6}, \qquad |\theta_1| \le \|f^{(3)}\|_c .$

Regrouping (2) and applying these formulas, we get:

$\dfrac{\xi}{h} = \dfrac{y(t_{k+1}) - y(t_k)}{h} - \dfrac{y'(t_{k+1}) + y'(t_k)}{2} = y'(t_{k+1/2}) + \dfrac{\eta h^2}{6} - y'(t_{k+1/2}) - \dfrac{\eta_1 h^2}{2} = (\dfrac{\eta}{6} - \dfrac{\eta_1}{2})h^2,$ where

$|\eta| \le \|y^{(3)}\|_c, \quad |\eta_1| \le \|y''\|_c.$ Hence $|\xi| \le \dfrac{1}{2}(\|y''\|_c + \|y^{(3)}\|_c)h^3.$ $\qquad (3)$

Now denoting $\delta_k = y_k - y(t_k)$ and subtracting (1)–(2), we get

$\delta_{k+1} = \delta_k + \dfrac{f(t_k, y_k) - f(t_k, y(t_k))}{2} h + \dfrac{f(t_{k+1}, y_{k+1}) - f(t_{k+1}, y(t_{k+1}))}{2} h - \xi.$

Then using Lipschitz-continuity of $f$,

$|\delta_{k+1}| \le |\delta_k| + \dfrac{Lh}{2}|\delta_k| + \dfrac{Lh}{2}|\delta_{k+1}| + |\xi|, \qquad (1 - \dfrac{Lh}{2})|\delta_{k+1}| \le (1 + \dfrac{Lh}{2})|\delta_k| + |\xi|,$

$$|\delta_{k+1}| \le \left|\frac{1+Lh/2}{1-Lh/2}\right||\delta_k| + \frac{|\xi|}{1-Lh/2} \le \left|\frac{1+Lh/2}{1-Lh/2}\right||\delta_k| + 2|\xi| \qquad \text{for} \quad h \le \frac{1}{L}.$$

To simplify this, consider the general situation: $x_{k+1} \le ax_k + b$. Applying this inequality recursively,

we get $x_k \le a^k x_0 + b(a^{k-1} + a^{k-2} + \ldots + 1) = a^k x_0 + b\frac{a^k - 1}{a - 1}$. The first term in our case disappears,

because $y(0) = y_0 = 1$; therefore $|\delta_k| \le 2|\xi| \dfrac{\left|\dfrac{1+Lh/2}{1-Lh/2}\right|^k - 1}{\left|\dfrac{1+Lh/2}{1-Lh/2}\right| - 1}$.

Now since $k = t_k/h$, $\quad \left|\dfrac{1+Lh/2}{1-Lh/2}\right|^k \xrightarrow[h\to 0]{} \dfrac{e^{\frac{Lt_k}{2}}}{e^{-\frac{Lt_k}{2}}} = e^{Lt_k}$, and hence

$$|\delta_k| \le 2|\xi|\frac{e^{Lt_k} + \varepsilon(h) - 1}{Lh}(1 - Lh/2) \le \frac{2|\xi|}{h}\frac{e^{Lt_k} + \varepsilon(h) - 1}{L},$$

where $\dfrac{2|\xi|}{h} \le (\|y''\|_c + \|y^{(3)}\|_c)h^2$, and $\varepsilon(h) \to 0$ as $h \to 0$. Thus $|\delta_k| \le Ch^2$ for sufficiently small $h$.

**3 a)** Let $l(x) = f(a) + (x-a)\dfrac{f(b) - f(a)}{b - a}$, then

$$I(f) = \int_a^b l(x)dx = f(a)(b-a) + \frac{(x-a)^2}{2}\bigg|_a^b \frac{f(b)-f(a)}{b-a} = (b-a)(f(a) + \frac{f(b)-f(a)}{2}) = (b-a)\frac{f(a)+f(b)}{2}$$

---

**b)** Let $I_0 = \int_a^b f(x)dx$, $\quad \Delta = b - a$.

$$f(a+t) = f(a) + f'(a)t + f''(\xi)\frac{t^2}{2}, \quad I_0 = f(a)\Delta + f'(a)\frac{\Delta^2}{2} + A\frac{\Delta^3}{6}, \qquad |A| \le \|f''\|_c$$

$$f(b-t) = f(b) - f'(b)t + f''(\xi_1)\frac{t^2}{2}, \quad I_0 = f(b)\Delta - f'(b)\frac{\Delta^2}{2} + B\frac{\Delta^3}{6}, \qquad |B| \le \|f''\|_c .$$

Averaging the two expressions for $I_0$, we obtain

$$I_0 = \frac{f(a)+f(b)}{2}\Delta + \frac{A+B}{2}\frac{\Delta^3}{6} = I + C\frac{\Delta^3}{6}, \quad \text{where} \quad C = \frac{A+B}{2}. \text{ This implies } |I_0 - I| \le \|f''\|_c\frac{\Delta^3}{6}.$$

**c)** Let $I_0^{(i)}$, $I^{(i)}$ be the true integral and our estimate correspondingly for the $i$th subinterval. Then

$$I_0 = \sum_{i=1}^n I_0^{(i)}, \quad I = \sum_{i=1}^n I^{(i)}.$$

$$|I_0 - I| \le \sum_{i=1}^n |I_0^{(i)} - I^{(i)}| \underset{(b)}{\le} \sum_{i=1}^n |C_i|\frac{\Delta_i^3}{6} \le \sum_{i=1}^n \|f''\|_c\frac{1}{6}\left(\frac{\Delta}{n}\right)^3 = \|f''\|_c\frac{\Delta^3}{6n^2} = \frac{\|f''\|_c(b-a)^3}{6n^2}.$$