# Errata and Corrections for Classical Numerical Analysis

November 7, 2023

# Chapter 1

# Linear Operators and Matrices

1. Page 5. The transpose operator for vectors acts as follows: (i) it converts column vectors into row vectors, $(\cdot)^{\mathsf{T}} : \mathbb{C}^{k \times 1} \to \mathbb{C}^{1 \times k}$, and (ii) it converts row vectors into column vectors, $(\cdot)^{\mathsf{T}} : \mathbb{C}^{1 \times k} \to \mathbb{C}^{k \times 1}$. Only one direction was specified in the text.

2. Page 8, proof of Theorem 1.21. Style consistency. Change the "$\exists y \in \mathbb{C}^m$" to "there is some $y \in \mathbb{C}^m$."

3. Page 9. We should explicitly define a matrix norm and provide more properties. For example, we should point out that the objects defined are all matrix norms, and the induced norm of the identity is 1. Otherwise, students are missing some key concepts and are confused.

4. Page 10, Proposition 1.31. In the hypotheses of the proposition, the extra $[a_{i,j}]$ in "$A = [a_{i,j}] = [a_{i,j}]$" is superfluous.

5. Page 11, Proposition 1.35. This result should be expanded to include both right and left multiplication by unitary matrices, and it should include the analogous results for the Frobenius norm.

6. Page 13, Proposition 1.43. The trace of a square matrix, denoted, $\text{tr}(A)$, is defined as the sum of the diagonal elements. Specifically, for $A = [a_{i,j}] \in \mathbb{C}^{n \times n}$, $\text{tr}(A) = \sum_{i=1}^{n} a_{i,i}$. The symbol $\det(A)$ stands for the determinant of $A$. See the references.

7. Page 14, Proposition 1.47. This should be Theorem 1.47.

8. Page 16, Problem 1.5. $C_A$ should just be $\text{col}(A)$, the previously introduced notation for the column space of $A$.

9. Page 16, Problem 1.19. Change the problem to the following: Suppose that $(\mathbb{V}, \|\cdot\|_{\mathbb{V}})$ and $(\mathbb{W}, \|\cdot\|_{\mathbb{W}})$ are finite-dimensional complex normed vector spaces. Suppose that $\|\cdot\|_{\mathfrak{L}(\mathbb{V}, \mathbb{W})} : \mathfrak{L}(\mathbb{V}, \mathbb{W}) \to \mathbb{R}$ is the induced norm. Then, $\|\cdot\|_{\mathfrak{L}(\mathbb{V}, \mathbb{W})}$ is a bona fide norm on the vector space $\mathfrak{L}(\mathbb{V}, \mathbb{W})$ and

$$\|A\|_{\mathfrak{L}(\mathbb{V}, \mathbb{W})} = \sup \{ \|Ax\|_{\mathbb{W}} \mid x \in \mathbb{V}, \ \|x\|_{\mathbb{V}} = 1 \}$$
$$= \sup \{ \|Ax\|_{\mathbb{W}} \mid x \in \mathbb{V}, \ \|x\|_{\mathbb{V}} \leq 1 \} .$$

Furthermore, for the identity operator $I \in \mathfrak{L}(\mathbb{V})$, we have $\|I\|_{\mathfrak{L}(\mathbb{V})} = 1$.

10. Page 17, Problem 1.24. This should come after Problem 1.29. Furthermore, we need the general result $\rho(A) \leq \|A\|$, for any induced norm, for any square matrix. This, fact, however, does not appear until Chapter 4, specifically, Theorem 4.3.

11. Page 17, Problem 1.26. The symbol $\text{tr}\,A$ should be $\text{tr}(A)$ for notational consistency. This problem needs $\|UA\| = \|A\|$ and $\|AV\| = \|A\|$ for the 2 and Frobenius norms, where $U$ and $V$ are unitary. Unfortunately, the results are only partially alluded to. See Proposition 1.35.

12. Page 17, Problem 1.29. The hint should refer to Problem 1.39 and Proposition/Theorem 1.47, specifically.

13. Page 17, Problem 1.30. Add the following to the end of the problem: "where

$$S_{\mathbb{C}^n}^{n-1} = \{x \in \mathbb{C}^n \mid \|x\|_{\mathbb{C}^n} = 1\}\,."$$

14. Page 18, Problem 1.32. The assumptions in parts (c) and should be corrected by adding "for all $i \leq i \leq n$."

# Chapter 2

# The Singular Value Decomposition

1. Page 22, Theorem 2.3. Add a remark. The proof of existence of the SVD can be replaced by a more elementary one. See Problem 2.9, page 30.

# Chapter 3

# Systems of Linear Equations

1. Page 35, Theorem 3.5, part 2 of theorem hypotheses. Em dash or en dash? Check the style guide.

2. Page 35, Theorem 3.5, part 6. $T_k^{-1}$ should be $T^{-1}$. The hypotheses should read as follows: If $[T]_{i,i} > 0$, for all $i = 1, \ldots, n$, then $\left[T^{-1}\right]_{i,i} = \frac{1}{[T]_{i,i}} > 0$, for all $i = 1, \ldots, n$.

3. Page 52, Proof of Theorem 3.24. In the proof of the second part, instead of
$$\frac{\|Ax\|_\infty}{\|x\|_\infty} \geq \delta, \quad \forall x \in \mathbb{C}^n,$$
the line should read
$$\frac{\|Ax\|_\infty}{\|x\|_\infty} \geq \delta, \quad \forall x \in \mathbb{C}_\star^n.$$
The subscript $\star$ was missing.

4. Page 66, Problem 3.15. Add hint: "Use the Gershgorin Circle Theorem."

5. Page 68, Listing 3.2. The variable `denominator` is spelled four different ways. The code has been corrected in the Github repo.

# Chapter 4

# Norms and Matrix Conditioning

1. Page 74, Proof of Theorem 4.3. The line
$$\frac{\|Ax\|_{\mathbb{C}^n}}{\|x\|_{\mathbb{C}^n}} \le \frac{C_2}{C_1}\frac{\|Ax\|_2}{\|x\|_2} \le C\,\|A\|_2 \le C\sqrt{\rho(A^HA)}$$
   should be replace by
$$\frac{\|Ax\|_{\mathbb{C}^n}}{\|x\|_{\mathbb{C}^n}} \le \frac{C_2}{C_1}\frac{\|Ax\|_2}{\|x\|_2} \le C\,\|A\|_2 = C\sqrt{\rho(A^HA)}.$$
   In other words, the last inequality should be an equality.

2. Page 78, Corollary 4.9. Punctuation/grammatical error. The statement of the corollary should read as follows: "Let $M \in \mathbb{C}^{n \times n}$, and assume that, for some induced matrix norm $\|\cdot\| : \mathbb{C}^{n \times n} \to \mathbb{R}$,
$$\|M\| < 1.$$
   Then, $M$ is convergent to zero."

3. Page 79, Corollary 4.12. The assumption about the norm is unnecessary. That is, one should delete/neglect the following assumption: "Suppose that $\|\cdot\| : \mathbb{C}^{n \times n} \to \mathbb{R}$ is an induced matrix norm."

4. Page 80, Remark 4.15. Here we refer to the term "ill-conditioned," but it is not defined. We say a matrix is ill-conditioned if it has a large condition number, that is, significantly larger than 1.

5. Page 86, Theorem 4.21. The assumption $\|A^{-1}\delta A\| < 1$ should be replaced by $\|A^{-1}\| \cdot \|\delta A\| < 1$. The proof should be modified accordingly. In particular, to conclude that the perturbed coefficient matrix $A + \delta A$ is invertible, use the inequality
$$\|M\| = \|-A^{-1}\delta A\| \le \|A^{-1}\| \cdot \|\delta A\| < 1.$$
   Later in the proof, one can be assured that
$$\frac{1}{1 - \|M\|} \le \frac{1}{1 - \|A^{-1}\|\,\|\delta A\|} = \frac{1}{1 - \kappa(A)\frac{\|\delta A\|}{\|A\|}},$$
   an inequality that may fail if only $\|A^{-1}\delta A\| < 1$ is assumed.

6. Page 87, Theorem 4.22. Same correction as in the previous theorem.

# Chapter 5

# Linear Least Squares Problem

1. Page 90, proof of Lemma 5.4. The line

$$\left(A^H A\boldsymbol{x}, \boldsymbol{x}\right)_2 = (A\boldsymbol{x}, A\boldsymbol{x})_2 = \|A\boldsymbol{x}\|_2 \geq 0.$$

   should read

$$\left(A^H A\boldsymbol{x}, \boldsymbol{x}\right)_2 = (A\boldsymbol{x}, A\boldsymbol{x})_2 = \|A\boldsymbol{x}\|_2^2 \geq 0.$$

   The exponent 2 is missing on the norm.

2. Page 90, proof of Lemma 5.4. The sentence fragment "to reach a contradiction, that $\mathrm{rank}(A) < n$" should be "to reach a contradiction, that $\mathrm{rank}(\mathsf{A}) < n$". In other words the matrix should be written as $\mathsf{A}$ not $A$. Recall that the symbol $\mathsf{A}$ usually represents a matrix, while the symbol $A$ represents a linear operator.

3. Page 91, proof of Theorem 5.5. The line

$$= \Phi(\boldsymbol{x}) - \boldsymbol{r}^\mathsf{T} A\boldsymbol{y} - \boldsymbol{y} A^\mathsf{T} \boldsymbol{r} + \boldsymbol{y}^\mathsf{T} A^\mathsf{T} A\boldsymbol{y}$$

   should read

$$= \Phi(\boldsymbol{x}) - \boldsymbol{r}^\mathsf{T} A\boldsymbol{y} - \boldsymbol{y}^\mathsf{T} A^\mathsf{T} \boldsymbol{r} + \boldsymbol{y}^\mathsf{T} A^\mathsf{T} A\boldsymbol{y}.$$

   In other words, there is a missing $^\mathsf{T}$ on $\boldsymbol{y}$ in the term $-\boldsymbol{y} A^\mathsf{T} \boldsymbol{r}$.

4. Page 95, Remark 5.19. The first statement is incorrect. If $S_1$ and $S_2$ are generic orthogonal subspaces of $\mathbb{C}^n$, it is not necessarily true that $S_1 \oplus S_2 = \mathbb{C}^n$. In other words, $S_1$ and $S_2$ can be orthogonal and not complementary.

5. Page 96, Theorem 5.23. This should be listed as Corollary 5.23. Its proof follows directly from the Theorem 5.21.

6. Page 100, Proof of Theorem 5.26. In the "(2 $\implies$ 1)" part of the proof, specifically, last line to establish $\Phi(\boldsymbol{x}_o + \boldsymbol{w}) \geq \Phi(\boldsymbol{x}_o)$, the "$\geq \Phi(\boldsymbol{x}_o)$" should be on a separate line to follow style guidelines.

7. Page 111, Proof of Lemma 5.47. The line

$$\|\boldsymbol{x}\|_2 = \|\mathsf{H}_{\boldsymbol{w}} \boldsymbol{x}\|_2 = |k| \, \|\boldsymbol{x}\|_2 \, \|\boldsymbol{e}_j\|$$

11

should read
$$\|x\|_2 = \|H_w x\|_2 = |k| \, \|x\|_2 \, \|e_j\|_2 \, .$$

In other words, the term $\|e_j\|$ is missing the subscript 2.

8. Pages 112–114, Lemma 5.50 and Definition 5.51. The symbol $\hat{e}_1$ should be changed to $e_1$. Check notational consistency throughout. Do we use $\hat{e}_j$ or $e_j$ for canonical basis elements.

9. Page 113, Definition 5.51. Equation (5.19)

$$a_1^{(0)} \neq \pm e^{i\alpha_1} \left\| a_1^{(0)} \right\|_{s^2(\mathbb{C}^m)} \widehat{e}_1$$

is incorrect. The subscript of the norm should be $\ell^2(\mathbb{C}^m)$ instead of $s^2(\mathbb{C}^m)$. In other words, Equation (5.19) should read

$$a_1^{(0)} \neq \pm e^{i\alpha_1} \left\| a_1^{(0)} \right\|_{\ell^2(\mathbb{C}^m)} \widehat{e}_1.$$

10. Page 119, Listing 5.1. Lines 15 and 16 are incorrect. The lines should be changed from

```
15    m = size(A)(1);
16    n = size(A)(2);
```

to

```
15    m = size(A,1);
16    n = size(A,2);
```

A similar change has been made in Listing A.1. The code has been updated in the GitHub repo.

11. Page 120, Listing 5.2. Lines 17 and 18 are incorrect. The lines should be changed from

```
17    m = size(A)(1);
18    n = size(A)(2);
```

to

```
17    m = size(A,1);
18    n = size(A,2);
```

A similar change has been made in Listing 5.1. The code has been updated in the GitHub repo.

# Chapter 6

# Linear Iterative Methods

1. Page 126. Notation. In this chapter we use the notation $[\boldsymbol{x}_k]_i = x_{i,k}$. But, it seems that we use the notation $[\boldsymbol{x}_k]_i = x_{k,i}$ in other places. Check the consistency.

2. Page 129, proof of Theorem 6.12. Half way down the page, the line

$$\left|\sum_{j+1}^{n} a_{i,j} x_j\right| \le \sum_{j+1}^{n} |a_{i,j}| \, \|\boldsymbol{x}\|_\infty$$

   should be replaced by

$$\left|\sum_{j=i+1}^{n} a_{i,j} x_j\right| \le \sum_{j=i+1}^{n} |a_{i,j}| \, \|\boldsymbol{x}\|_\infty.$$

   In other words, the lower summation indices are incorrect.

3. Page 129, proof of Theorem 6.12. Last line of the proof. The line

$$\|T_{\mathsf{GS}}\|_\infty \le \gamma < 1$$

   should read

$$\|\mathsf{T}_{\mathsf{GS}}\|_\infty \le \gamma < 1.$$

   In other words, the typeface of the $\mathsf{T}_{\mathsf{GS}}$ is incorrect.

4. Page 134, proof of Theorem 6.15. Halfway down the page, the line

$$\|\boldsymbol{e}_{k+1}\|_2 = \left\|\mathsf{T}_{\mathsf{R}}^k \boldsymbol{e}_0\right\|_2 \le \rho^k \|\boldsymbol{e}_0\|_2.$$

   should read

$$\|\boldsymbol{e}_k\|_2 = \left\|\mathsf{T}_{\mathsf{R}}^k \boldsymbol{e}_0\right\|_2 \le \rho^k \|\boldsymbol{e}_0\|_2.$$

   In other words, $\boldsymbol{e}_{k+1}$ should be $\boldsymbol{e}_k$.

5. Page 137, proof of Theorem 6.18. Top of the page, the proof that $(\mathsf{B}_\omega - \mathsf{A})\,\boldsymbol{y} = \lambda \mathsf{A}\boldsymbol{w}$ can be significantly simplified.

6. Page 138, proof of Theorem 6.19. In the first line, the sentence "Another method of proof is demonstrated in the next section." should read instead "Another method of proof is demonstrated in the Section 6.8."

7. Page 138, proof of Theorem 6.19. For the forward direction, a few more steps are required for the proof. The fact that $\boldsymbol{w}^{\mathsf{H}}\mathsf{A}\boldsymbol{w} > 0$ for all eigenvectors $\boldsymbol{w}$ of $\mathsf{T}$ is not, on its own, enough to show that $\mathsf{A}$ is HPD. See Problem 6.6.

8. Page 138, proof of Theorem 6.19. Last line. The last line

   "every eigenvector $\boldsymbol{w} \in \mathbb{C}_\star^n$. This proves that $\mathsf{A}$ must be HPD."

   should be replaced by

   "every eigenvector $\boldsymbol{w} \in \mathbb{C}_\star^n$ of $\mathsf{T}$. However, this is not enough to prove that $\mathsf{A}$ is HPD. For this direction see Problem 6.6 and the proof of Theorem 6.26 for inspiration."

9. Page 141, proof of Theorem 6.25. After the words "we obtain", the calculation should read

   $$
   \begin{aligned}
   0 &= (\mathsf{B}\boldsymbol{q}_{k+1}, \boldsymbol{q}_{k+1})_2 + (\mathsf{A}\boldsymbol{e}_k, \boldsymbol{q}_{k+1})_2 \\
   &= \left( \left( \mathsf{B} - \frac{1}{2}\mathsf{A} \right) \boldsymbol{q}_{k+1}, \boldsymbol{q}_{k+1} \right)_2 + \frac{1}{2}(\mathsf{A}\boldsymbol{e}_{k+1}, \boldsymbol{e}_{k+1})_2 - \frac{1}{2}(\mathsf{A}\boldsymbol{e}_k, \boldsymbol{e}_k)_2 + i\Im\left( (\mathsf{A}\boldsymbol{e}_k, \boldsymbol{e}_{k+1})_2 \right) \\
   &= (\mathsf{Q}\boldsymbol{q}_{k+1}, \boldsymbol{q}_{k+1})_2 + \frac{1}{2}\|\boldsymbol{e}_{k+1}\|_{\mathsf{A}}^2 - \frac{1}{2}\|\boldsymbol{e}_k\|_{\mathsf{A}}^2 + i\Im\left( (\mathsf{A}\boldsymbol{e}_k, \boldsymbol{e}_{k+1})_2 \right).
   \end{aligned}
   $$

   In other words, the term $i\Im\left( (\mathsf{A}\boldsymbol{e}_k, \boldsymbol{e}_{k+1})_2 \right)$ is missing in the text. After this point, the proof is correct.

10. Page 147, proof of Theorem 6.30 and preceding discussion. We tacitly assume that $\alpha_{k+1}$ is real.

11. Page 148, proof of Theorem 6.31. We tacitly assume that $\alpha_{k+1}$ is real.

12. Page 148, proof of Theorem 6.31. The calculations and identities

    $$
    \begin{aligned}
    (\mathsf{C}\boldsymbol{v}_k, \boldsymbol{v}_k)_2 &= (\mathsf{S}^{-1/2}\mathsf{A}\mathsf{S}^{-1/2}\boldsymbol{w}_k, \mathsf{S}^{-1/2}\boldsymbol{w}_k)_2 = (\mathsf{A}\boldsymbol{w}_k, \boldsymbol{w}_k)_2, \\
    \|\mathsf{C}\boldsymbol{v}_k\|_2^2 &= (\mathsf{S}^{-1/2}\mathsf{A}\boldsymbol{w}_k, \mathsf{S}^{-1/2}\mathsf{A}\boldsymbol{w}_k)_2 = \|\mathsf{A}\boldsymbol{w}_k\|_{\mathsf{S}^{-1}}^2, \\
    \|\boldsymbol{v}_{k+1}\|_2 &= (\mathsf{S}\boldsymbol{w}_{k+1}, \boldsymbol{w}_{k+1})_2 = \|\mathsf{A}\boldsymbol{e}_{k+1}\|_{\mathsf{S}^{-1}}^2,
    \end{aligned}
    $$

    are incorrect. The correct calculations and identities are

    $$
    \begin{aligned}
    (\mathsf{C}\boldsymbol{v}_k, \boldsymbol{v}_k)_2 &= (\mathsf{S}^{-1/2}\mathsf{A}\mathsf{S}^{-1/2}\mathsf{S}^{1/2}\boldsymbol{w}_k, \mathsf{S}^{1/2}\boldsymbol{w}_k)_2 = (\mathsf{A}\boldsymbol{w}_k, \boldsymbol{w}_k)_2, \\
    \|\mathsf{C}\boldsymbol{v}_k\|_2^2 &= (\mathsf{S}^{-1/2}\mathsf{A}\boldsymbol{w}_k, \mathsf{S}^{-1/2}\mathsf{A}\boldsymbol{w}_k)_2 = \|\mathsf{A}\boldsymbol{w}_k\|_{\mathsf{S}^{-1}}^2, \\
    \|\boldsymbol{v}_{k+1}\|_2^2 &= (\mathsf{S}\boldsymbol{w}_{k+1}, \boldsymbol{w}_{k+1})_2 = \|\mathsf{A}\boldsymbol{e}_{k+1}\|_{\mathsf{S}^{-1}}^2.
    \end{aligned}
    $$

# Chapter 7

# Variational and Krylov Subspace Methods

1. Page 159, Definition 7.3. The sentence fragment "We say that B is self-adjoint positive definite with respect to the inner product ..." should read "We say that B is **self-adjoint positive definite** with respect to the inner product ..." In other words, "self-adjoint positive definite" should appear in bold letters.

2. Page 164, Proof of Theorem 7.16. In the last line of the proof the sentence "Combining (7.6) and(7.7), we get the desired result." should read "Combining (7.6) and (7.7), we get the desired result." In other words, a space should be placed between "and" and "(7.7)".

3. Page 171, Theorem 7.31. The theorem statement should be modified to read as follows:

   **Theorem 7.31** (convergence). Let $A \in \mathbb{C}^{n \times n}$ be HPD, $f \in \mathbb{C}_\star^n$, and $x = A^{-1}f$. Suppose that $\{x_k\}_{k=1}^\infty$ is the sequence generated by the zero-start CG algorithm. Then, there is an $m_\star \in \{1, \ldots, n\}$ for which
   $$x_k \neq x, \quad 1 \leq k \leq m_\star - 1, \quad x_k = x, \quad k \geq m_\star,$$
   and $\dim \mathcal{K}_k(A, f) = k$, for $k = 1, \ldots, m_\star$.

4. Page 172, proof of Theorem 7.31. The line "... and notice that, since $f \neq 0$, ..." should read "... and notice that, since $f \neq \mathbf{0}$, ...".

5. Page 172, proof of Theorem 7.31. The line "Assume now that, for all $m = 1, \ldots, k$ with $k < n - 1$ we have $\dim \mathcal{K}_k = k$ and $x_k \neq x$." should be replaced by "Assume now that, for all $m = 1, \ldots, k$, with $k < n - 1$, we have $\dim \mathcal{K}_m = m$ and $x_m \neq x$."

6. Page 175, statement of Theorem 7.36. The line "...for all $j = 1, \ldots, n$, with the orthogonality relations..." should read "...for all $j = 1, \ldots, m$, with the orthogonality relations...". In other words, the $n$ should be an $m$.

7. Page 176, proof of Theorem 7.36. In the expansion of $\phi^2(z)$, the term $2w^H A e_j$ should be replaced by $2\Re\left(w^H A e_j\right)$, and the term $2w^H r_j$ should be replaced by $2\Re\left(w^H r_j\right)$, since the computation is done over the complex field.

# Chapter 8

# Eigenvalue Problems

1. Page 200, statement of Theorem 8.2. The expression "$\sigma(A) \subset \bigcup_{i=1}^{n} D_i$" should be changed to "$\sigma(A) \subseteq \bigcup_{i=1}^{n} D_i$". It is possible to have set equality when the Gerschgorin radii are all zeros.

2. Page 209, statement of Theorem 8.19. The line

$$\lambda_r = \operatorname*{argmin}_{j=1}^{n} |\lambda_j - \mu|, \quad \lambda_s = \operatorname*{argmin}_{\substack{j=1 \\ j \neq r}}^{n} |\lambda_j - \mu|$$

should be

$$r = \operatorname*{argmin}_{j=1}^{n} |\lambda_j - \mu|, \quad s = \operatorname*{argmin}_{\substack{j=1 \\ j \neq r}}^{n} |\lambda_j - \mu|.$$

3. Page 218, proof of Theorem 8.26. At the top of the page, the line

$$\tilde{Q}_k, \to I \qquad \tilde{R}_k \to I, \qquad k \to \infty$$

should read

$$\tilde{Q}_k \to I \qquad \tilde{R}_k \to I, \qquad k \to \infty.$$

There is an errant comma.

# Chapter 15

# Solution of Nonlinear Equations

1. Page 432, proof of Theorem 15.26. There is a format error 2/3 of the way down the page. The multiline equation should begin a new line with the second equals sign, and there is a missing comma. The separate line should read
$$= f^{(m)}(\zeta_k) \frac{(x_k - \xi)^{m-1}}{(m-1)!},$$

2. Page 434, proof of Theorem 4.27. Format error, similar to that above. In the multiline equation 1/2 down the page, the last equals sign should begin a new line. The separate line should read
$$= \frac{1}{2}|x_k - \xi|.$$

3. Page 435, Theorem 15.26. Add to the assumptions of the theorem that $m \geq 2$. Otherwise, the rate of convergence is not exactly linear.

4. Page 439, proof of Theorem 15.33. There is an error in the proof. **The Following lines**

   Thus,
   $$x_{k+1} - \xi = x_k - \xi - \frac{f'(\gamma_k)(x_k - \xi)}{f'(\eta_k)} = (x_k - \xi)\left[1 - \frac{f'(\gamma_k)}{f'(\eta_k)}\right] \leq \frac{2}{5}(x_k - \xi).$$

   If $|x_0 - \xi| \leq \delta$ and $|x_1 - \xi| \leq \delta$ then, by induction, we see that for $k \geq 2$,
   $$|x_k - \xi| \leq \left(\frac{2}{5}\right)^{k-1} \delta.$$

   **should be replaced by**

   Thus,
   $$x_{k+1} - \xi = x_k - \xi - \frac{f'(\gamma_k)(x_k - \xi)}{f'(\eta_k)} = (x_k - \xi)\left[1 - \frac{f'(\gamma_k)}{f'(\eta_k)}\right].$$

   Since
   $$-\frac{2}{3} \leq 1 - \frac{f'(\gamma_k)}{f'(\eta_k)} \leq \frac{2}{5},$$

it follows that

$$|x_{k+1} - \xi| \leq \frac{2}{3}|x_k - \xi|.$$

If $|x_0 - \xi| \leq \delta$ and $|x_1 - \xi| \leq \delta$ then, by induction, we see that for $k \geq 2$,

$$|x_k - \xi| \leq \left(\frac{2}{3}\right)^{k-1} \delta.$$

5. Pages $440 - 441$. All references to the $i^{\text{th}}$ component of the vector $\boldsymbol{f}$ should be $f_i$, not $\boldsymbol{f}_i$. Components of a vector function should be unbolded.

6. Page 443, proof of Theorem 15.37. Format error. In the multiline equation/inequality $1/3$ of the way down the page, the subsequent equals signs and less than equals signs should each begin a new line.

7. Page 448, problem 15.28. The problem steps are incorrect. They should read as follows:

   a) Let $\boldsymbol{e}_k = \boldsymbol{\xi} - \boldsymbol{x}_k$ be the error. Establish an iteration error equation of the form

   $$\begin{bmatrix} \frac{\partial f}{\partial x_1}(\boldsymbol{\xi}) & 0 \\ \frac{\partial g}{\partial x_1}(\boldsymbol{\xi}) & \frac{\partial g}{\partial x_2}(\boldsymbol{\xi}) \end{bmatrix} \boldsymbol{e}_{k+1} = \begin{bmatrix} \frac{\partial f}{\partial x_1}(\boldsymbol{\xi})\boldsymbol{e}_{1,k+1} \\ \frac{\partial g}{\partial x_1}(\boldsymbol{\xi})\boldsymbol{e}_{1,k+1} + \frac{\partial g}{\partial x_2}(\boldsymbol{\xi})\boldsymbol{e}_{2,k+1} \end{bmatrix} = \boldsymbol{r}_{k+1}.$$

   Give a precise expression for the remainder term, $\boldsymbol{r}_{k+1}$.

   b) Give sufficient conditions for the convergence of the scheme.

# Chapter 16

# Convex Optimization

1. Page 452, Example 16.2. The definition of the inner product

$$(p, q)_{L^2(-1,1)} = \int_0^1 p(x)q(x)\,dx, \quad \forall\, p, q \in \mathbb{P}_n$$

   is incorrect. The definition should read

$$(p, q)_{L^2(-1,1)} = \int_{-1}^1 p(x)q(x)\,dx, \quad \forall\, p, q \in \mathbb{P}_n$$

   In other words, the lower limit of the integral should be $-1$, not $0$.

2. Page 462, proof of Theorem 16.27. About halfway through the proof, the line

$$\alpha < y_n < \alpha + \frac{1}{n}$$

   should instead read

$$\alpha \le y_n < \alpha + \frac{1}{n}.$$

3. Page 467, after the proof of Proposition 16.46. The line

   "If $E$ is strongly convex and the Lipschitz smooth, then ..."

   should read

   "If $E$ is strongly convex and Lipschitz smooth, then ..."

   There is a superfluous "the".

# Chapter 17

# Initial Value Problems for Ordinary Differential Equations

1. Page 510, Definition 17.4. The fragment "for all $t \in I$ and for all $\boldsymbol{v}_1, \boldsymbol{v}_2 \in \Omega$" should be replaced with "for all $t \in I$ and for all $\boldsymbol{v}_1, \boldsymbol{v}_2 \in \overline{\Omega}$." In other words, the $\Omega$ should be $\overline{\Omega}$.

2. Page 511, Theorem 17.5. The value of $\delta_1$ is incorrect. It should read

$$\delta_1 = \max\left\{\delta_0, \frac{1}{2L}, \frac{\beta}{M}\right\}.$$

3. Page 512, proof of Theorem 17.5. The interval $I_1$ is incorrect. It should read

$$I_1 = \left[t_0 - \frac{1}{2L}, t_0 + \frac{1}{2L}\right] \cap J.$$

4. Page 517, Proposition 17.11. Add to the initial hypotheses the following: "Suppose that $S = [0, T] \times \overline{\Omega}$, where $\Omega \subseteq \mathbb{R}^d$ is open and convex."

5. Page 517, Proposition 17.11. To the end of the proposition statement add the following: "Moreover, if $S = [0, T] \times \mathbb{R}^d$, then $\boldsymbol{f}$ is globally $\boldsymbol{u}$-Lipschitz.

6. Page 517, Remark 17.12. The sentence "The assumption that $\boldsymbol{f} \in F^1(S)$ is not often verified in practice." should be replaced by "The assumption that $\boldsymbol{f} \in F^1(S)$, when $S = [0, T] \times \mathbb{R}^d$, is not often verified in practice. In fact, it often fails to be true. If $\boldsymbol{f} \in F^1([0, T] \times \mathbb{R}^d)$, then $\boldsymbol{f}$ would be globally $\boldsymbol{u}$-Lipschitz. In many important, real-world problems the slope function is only locally Lipschitz."

7. Page 518, Theorem 17.13. The assumptions about the solution must be changed. Replace the sentence

   "Then the unique classical solution on $I$ to (17.1), which we denote $\boldsymbol{u} \in C^1\left(I; \mathbb{R}^d\right)$, actually belongs to $C^{m+1}\left(I; \mathbb{R}^d\right)$."

   with the following sentences:

   "Assume that $\boldsymbol{u} \in C^1(I, \Omega)$ is a classical solution to (17.1). Then $\boldsymbol{u} \in C^{m+1}(I; \Omega)$."

In other words, we should assume that $\boldsymbol{u} \in C^1(I, \Omega)$ not $\boldsymbol{u} \in C^1(I, \mathbb{R}^d)$ to conform to the general definition of $S$.

# Chapter 18

# Single-Step Methods

1. Page 527, before Definition 18.2, add the following remark:

   **Remark 18.1.** We will assume throughout this and the next few chapters that the slope function satisfies
   $$\boldsymbol{f} \in F^1(S), \quad \text{where} \quad S = [0, T] \times \mathbb{R}^d,$$
   which implies that $\boldsymbol{f}$ is globally $\boldsymbol{u}$-Lipschitz continuous. This simplification guarantees the existence of a classical solution. More importantly, it allows us to apply the Lipschitz estimate, with a single constant $L$, with either the solution values or with the approximate solution values, without worrying about whether those values are in some bounded open set $\Omega$.

# Chapter 19

# Runge-Kutta Methods

1. Page 536. The first Taylor expansion,

$$u(t+s) = u(t) + su'(t) + \frac{s^2}{2}u''(t) + \frac{s^3}{6}u''(t) + \mathcal{O}(|s|^4)$$

   is incorrect. It should read as follows:

$$u(t+s) = u(t) + su'(t) + \frac{s^2}{2}u''(t) + \frac{s^3}{6}u'''(t) + \mathcal{O}(|s|^4)$$

   In other words, the term $\frac{s^3}{6}u''(t)$ should be $\frac{s^3}{6}u'''(t)$.

2. Page 338, proof of Theorem 19.1. The last estimate on the page,

$$(1+\tau L)^m < e^{m\tau L} \leq e^{K\tau L} = e^{TL},$$

   is incorrect. It should read as follows:

$$\left(1 + \tau L + \frac{\tau^2 L^2}{2}\right)^m < e^{m\tau L} \leq e^{K\tau L} = e^{TL}.$$

3. Page 541, Remark 19.8. The definition of the local truncation error,

$$\tau \mathcal{E}[u](t,s) = u(t) - u(t-s) - s\sum_{i=1}^{r} b_i f(t - \tau + c_i\tau, \boldsymbol{\xi}_{e,i}),$$

   is incorrect. It should read as follows:

$$\mathcal{E}[u](t,s) = \frac{u(t) - u(t-s)}{s} - \sum_{i=1}^{r} b_i f(t - s + c_i s, \boldsymbol{\xi}_{e,i}).$$

4. Page 545, proof of Theorem 19.13, second sentence. The statement $\boldsymbol{\rho} \in [\mathbb{P}_r]^d$ is incorrect. It should read $\boldsymbol{\rho} \in [\mathbb{P}_{r-1}]^d$. In other words, the $r$ should be $r-1$.

# Chapter 20

# Linear Multi-step Methods

1. Page 569, Definition 20.16. The statement

   > We say that solutions to (20.8) are **stable** iff given any starting values $\{\zeta_k\}_{k=0}^{q-1} \subset \mathbb{R}$, the sequence $\{\zeta_k\}_{k=0}^{\infty} \subset \mathbb{R}$ is bounded by a constant $C > 0$ that only depends upon the starting values.

   is incorrect. It should read

   > We say that solutions to (20.8) are **stable** iff given any starting values $\{\zeta_k\}_{k=0}^{q-1} \subset \mathbb{C}$, the sequence $\{\zeta_k\}_{k=0}^{\infty} \subset \mathbb{C}$ is bounded by a constant $C > 0$ that only depends upon the starting values.

   In other words, $\mathbb{R}$ should be replaced by $\mathbb{C}$.

2. Page 570, Example 20.9. The equation

$$w^K = \frac{2^K - 1}{K} \to \infty.$$

   should be replaced by

$$\zeta^K = \frac{2^K - 1}{K} \to \infty.$$

# Chapter 21

# Stiff Systems of Ordinary Differential Equations and Linear Stability

1. Bottom of page 589, proof of Theorem 21.17. The fragment "for some coefficients $\alpha_{j,i} \in \mathbb{C}$, $q$ a number, ..." is incorrect. It should read "for some coefficients $\alpha_{j,i} \in \mathbb{C}$, $q$ in number, ...". In other words, the word "a" should be replaced by the word "in".

# Chapter 24

# Finite Difference Methods for Elliptic Problems

1. Page 676, Proposition 24.20. The BDF operator notation introduced in

$$(\text{BDF } w)_i = \frac{1}{2h} \left(3w_i - 4w_{i-1} + w_{i-2}\right),$$

   is incorrect and should be changed to

$$(\text{BDF}_2 \, w)_i = \frac{1}{2h} \left(3w_i - 4w_{i-1} + w_{i-2}\right),$$

2. Page 683, proof of Theorem 24.30, last estimate of the proof. There are two errors. The estimate

$$0 = -h^2 \Delta v_1 = 2v_1 - v_0 - v_2 > 0$$

   should be replaced by

$$0 \geq -h^2 \Delta_h v_1 = 2v_1 - v_0 - v_2 > 0.$$

   In particular, the "=" should be "≥" and the "Δ" should be "$\Delta_h$".

3. Page 684, proof of Theorem 24.31, next to last estimate of the proof. The estimate

$$\pm w \leq \pm \psi \leq \max\left\{\psi_\pm(0), \psi_\pm(1)\right\} \leq \max_{j \in \{0,1\}} |g_j| + \frac{\|f_h\|_{L_h^\infty}}{8}$$

   should be replaced by

$$\pm w \leq \psi_\pm \leq \max\left\{\psi_\pm(0), \psi_\pm(1)\right\} \leq \max_{j \in \{0,1\}} |g_j| + \frac{\|f_h\|_{L_h^\infty}}{8}.$$

   In other words, "$\pm \psi$" should be "$\psi_\pm$".

4. Page 695, proof of Theorem 24.50. The estimate

$$v(kj, \ell h) = \max_{(ih, jh) \in \Omega_h} v(ih, jh) > \max_{(ih, jh) \in \partial \Omega_h} v(ih, jh).$$

is incorrect. It should read

$$v(kh, \ell h) = \max_{(ih,jh) \in \Omega_h} v(ih, jh) > \max_{(ih,jh) \in \partial\Omega_h} v(ih, jh).$$

In other words, "$v(kj, \ell h)$" should be "$v(kh, \ell h)$".

5. Page 695, proof of Theorem 24.50. In the definition of $\tilde{\Omega}_h$ near the end of the proof, the equation

$$\tilde{\Omega}_h = \Omega_h \setminus \{(h, h), (h, Nh), (Nh, h), N, N)\}$$

should be replaced by
$$\tilde{\Omega}_h = \Omega_h \setminus \{(h, h), (h, Nh), (Nh, h), Nh, Nh)\}.$$

In other words, "$(N, N)$" should be "$(Nh, Nh)$".

6. Page 695 proof of Theorem 24.51. The statement "Define the grid function $\phi_{i,j} = \Phi(ih, jh)$" should be "Define the grid function $\Phi_{i,j} = \Phi(ih, jh)$." In other words, $\phi_{i,j}$ should be $\Phi_{i,j}$.

7. Pages 695 and 696, proof of Theorem 24.51. Consistency error. Should use lower case $\psi$ as in the 1D proof from section 24.3.

8. Page 695, proof of Theorem 24.51. The last equation on the page

$$-\Delta_h \Psi = \pm f - \|f\|_{L_h^\infty} \leq 0$$

should read
$$-\Delta_h \Psi_\pm = \pm f - \|f\|_{L_h^\infty} \leq 0.$$

In other words, the subscript $\pm$ is missing.

9. Page 696, proof of Theorem 24.51. The estimate

$$\pm w_{i,j} \leq \Psi_{i,j} \leq \max_{\partial\Omega_h} \Psi \leq \max_{\partial\Omega_h} w + \frac{\|f\|_{L_h^\infty}}{8}$$

should read
$$\pm w_{i,j} \leq \Psi_{\pm i,j} \leq \max_{\partial\Omega_h} \Psi_\pm \leq \max_{\partial\Omega_h} w + \frac{\|f\|_{L_h^\infty}}{8}.$$

In other words, two $\pm$ subscripts are missing.

10. Page 699, Problem 22. The dimension of space should be 1, in particular, $\Omega = (0, 1)$.

# Chapter 25

# Finite Element Methods for Elliptic Problems

1. Page 702, top of the page. The expression $v_n = \sum_{i=1}^{n} v_i \phi_i$ may cause some confusion, and it will be corrected in the next edition. In particular, the object $v_n$ on the left is in the Hilbert space $\mathcal{H}$, whereas the object $v_n$ on the right is a real number, the $n^{\text{th}}$ coordinate in the basis expansion.

2. Page 704, proof of Theorem 25.6. The last equality

$$\|u - u_h\|_{\mathcal{H}} = \min_{v \in \mathcal{H}_n} \|u - v\|_{\mathcal{H}}$$

   is incorrect. It should read

$$\|u - u_h\|_{\mathcal{H}} \leq \min_{v \in \mathcal{H}_n} \|u - v\|_{\mathcal{H}}.$$

   In other words, the equality should be replaced by "$\leq$". Also, the infimum is achieved, but not necessarily at $u_h$.

3. Page 707, Theorem 25.15. In the second estimate, the norm $\|v' - \Pi_h v'\|_{L^2(I_j)}$ should be replaced by $\left\|(v - \Pi_h v)'\right\|_{L^2(I_j)}$

4. Page 707, Theorem 25.15. Use the simpler proof from the prelim packet. Students have struggled to understand the logic of the current proof.

5. Page 709. The spacing is inadequate in Equation (25.5). The line should read

$$\mathcal{A}(v, z_g) = \int_0^1 g(x) v(x) \, \mathrm{d}x, \quad \forall \, v \in H_0^1(0, 1).$$

# Chapter 28

# Finite Difference Methods for Parabolic Problems

1. Page 779, proof of Corollary 28.10. The equation
$$(I + \tau \Delta_h) w^{k+1} = \tau f_h^{k+1} + w^k = \tilde{f}_h^{k+1}$$
has a sign error. It should read
$$(I - \tau \Delta_h) w^{k+1} = \tau f_h^{k+1} + w^k = \tilde{f}_h^{k+1}.$$

2. Page 779, proof of Corollary 28.10. The estimate
$$\pm w^{k+1} \leq \psi_\pm^{k+1} \leq \max_{\partial_p \mathcal{C}_h^\tau} \psi \leq \|u_{0,h}\|_{L_h^\infty} + C \|f\|_{L_\tau^\infty(L_h^\infty)}$$
is incorrect. It should read
$$\pm w \leq \psi_\pm \leq \max_{\partial_p \mathcal{C}_h^\tau} \psi \leq \|u_{0,h}\|_{L_h^\infty} + C \|f\|_{L_\tau^\infty(L_h^\infty)}.$$
In other words, the superscripts $k+1$ are unnecessary.

3. Page 783, Remark 28.15. The first line of the estimate
$$(1 + \mu) |w_i^{k+1}| = |1 - \mu| |w_i^k| + \frac{\mu}{2} |w_{i-1}^{k+1}| + \frac{\mu}{2} |w_{i+1}^{k+1}| + \frac{\mu}{2} |w_{i-1}^k| + \frac{\mu}{2} |w_{i+1}^k|$$
is incorrect. It should read
$$(1 + \mu) |w_i^{k+1}| \leq |1 - \mu| |w_i^k| + \frac{\mu}{2} |w_{i-1}^{k+1}| + \frac{\mu}{2} |w_{i+1}^{k+1}| + \frac{\mu}{2} |w_{i-1}^k| + \frac{\mu}{2} |w_{i+1}^k|.$$
In other words, the equal sign should be $\leq$ after application of the triangle inequality.

4. Page 787, proof of Theorem 28.18, bottom third of the page. The estimate beginning with
$$\left\| e^{k+1} \right\|_2 \leq \left\| A e^n \right\|_2 + \tau \left\| A_2 \mathcal{E}_h^\tau [u]^{k+1} \right\|_2$$
is incorrect. It should read
$$\left\| e^{k+1} \right\|_2 \leq \left\| A e^k \right\|_2 + \tau \left\| A_2 \mathcal{E}_h^\tau [u]^{k+1} \right\|_2.$$
In other words, the superscript $n$ should be changed to $k$.

5. Page 791, proof of Corollary 28.23, first line.  The equation

$$\delta_\tau e^{k+1} + b\mathring{\delta}_h e^k - \Delta_h e^k = \mathcal{E}_h^\tau[u]^k, \quad k = 1, \ldots, K,$$

in incorrect.  It should be

$$\overline{\delta}_\tau e^{k+1} + b\mathring{\delta}_h e^k - \Delta_h e^k = \mathcal{E}_h^\tau[u]^k, \quad k = 1, \ldots, K,$$

In other words, the term $\delta_\tau e^{k+1}$ should be either $\overline{\delta}_\tau e^{k+1}$ or $\delta_\tau e^k$.

# Chapter 29

# Finite difference methods for hyperbolic problems

1. Page 811, Example 29.1. The equation

$$\delta_\tau w^{k+1} + \delta_h w^k = 0.$$

   is incorrect. It should be

$$\overline{\delta}_\tau w^{k+1} + \delta_h w^k = 0.$$

   In other words, $\delta_\tau w^{k+1}$ should be $\overline{\delta}_\tau w^{k+1}$.

2. Page 812, Definition 29.1. Equation (29.1),

$$\begin{cases} \delta_\tau w^{k+1} + c\overline{\delta}_h w^k = 0, & k = 0, \ldots, K-1, \\ w^0 = u_{0,h}; \end{cases}$$

   is incorrect. It should read

$$\begin{cases} \overline{\delta}_\tau w^{k+1} + c\overline{\delta}_h w^k = 0, & k = 0, \ldots, K-1, \\ w^0 = u_{0,h}; \end{cases}$$

3. Page 812, Definition 29.1. Equation (29.2),

$$\begin{cases} \delta_\tau w^{k+1} + c\mathring{\delta}_h w^k = 0, & k = 0, \ldots, K-1, \\ w^0 = u_{0,h}; \end{cases}$$

   is incorrect. It should read

$$\begin{cases} \overline{\delta}_\tau w^{k+1} + c\mathring{\delta}_h w^k = 0, & k = 0, \ldots, K-1, \\ w^0 = u_{0,h}; \end{cases}$$

4. Page 813, Definition 29.1. Equation (29.4),

$$\begin{cases} \delta_\tau w^{k+1} + c\mathring{\delta}_h w^k - \dfrac{c^2\tau}{2}\Delta_h w^k = 0, & k = 0, \ldots, K-1, \\ w^0 = u_{0,h}; \end{cases}$$

is incorrect. It should read

$$\begin{cases} \overline{\delta}_\tau w^{k+1} + c\mathring{\delta}_h w^k - \dfrac{c^2\tau}{2}\Delta_h w^k = 0, & k = 0, \dots, K-1, \\ w^0 = u_{0,h}; \end{cases}$$

5. Page 813, Definition 29.1. Equation (29.5),

$$\begin{cases} \delta_\tau w^{k+1} + c\,\mathrm{BDF}_h\,w^k - \dfrac{c^2\tau}{2}\Delta_h w^k = 0, & k = 0, \dots, K-1, \\ w^0 = u_{0,h}, \end{cases}$$

is incorrect. It should read

$$\begin{cases} \overline{\delta}_\tau w^{k+1} + c\,\mathrm{BDF}_2\,w^k - \dfrac{c^2\tau}{2}\Delta_h w^k = 0, & k = 0, \dots, K-1, \\ w^0 = u_{0,h}. \end{cases}$$

In other words, $\delta_\tau w^{k+1}$ should be $\overline{\delta}_\tau w^{k+1}$ and $\mathrm{BDF}_h\,w^k$ should be $\mathrm{BDF}_2\,w^k$.

6. Page 813, Definition 29.1. Equation (29.6),

$$\begin{cases} \delta_\tau w^{k+1} + \dfrac{c}{2}\mathring{\delta}_h\left(w^{k+1} + w^k\right) = 0, & k = 0, \dots, K-1, \\ w^0 = u_{0,h}. \end{cases}$$

is incorrect. It should read

$$\begin{cases} \overline{\delta}_\tau w^{k+1} + \dfrac{c}{2}\mathring{\delta}_h\left(w^{k+1} + w^k\right) = 0, & k = 0, \dots, K-1, \\ w^0 = u_{0,h}. \end{cases}$$

7. Page 815, Remark 29.7. The equation

$$\delta_\tau w^{k+1} + c\mathring{\delta}_h w^k - \dfrac{ch}{2}\Delta_h w^k = 0.$$

is incorrect. It should read

$$\overline{\delta}_\tau w^{k+1} + c\mathring{\delta}_h w^k - \dfrac{ch}{2}\Delta_h w^k = 0.$$

In other words, $\delta_\tau w^{k+1}$ should be $\overline{\delta}_\tau w^{k+1}$.

8. Page 820, proof of Theorem 29.14, first line. The equation

$$\delta_\tau w^{k+1} + c\mathring{\delta}_h w^k - \dfrac{ch}{2}\Delta_h w^k = 0$$

is incorrect. It should read

$$\overline{\delta}_\tau w^{k+1} + c\mathring{\delta}_h w^k - \dfrac{ch}{2}\Delta_h w^k = 0.$$

In other words, $\delta_\tau w^{k+1}$ should be $\overline{\delta}_\tau w^{k+1}$.

# Appendix A

# Linear Algebra Review

1. Page 852, Listing A.1. Lines 16 and 17 are incorrect. The lines should be changed from

   ```
   16    m = size(W)(1);
   17    n = size(W)(2);
   ```

   to

   ```
   16    m = size(W,1);
   17    n = size(W,2);
   ```

   A similar change correction applies to Listing 5.1. The code has been updated in the GitHub repo.

# Appendix B

# Basic Analysis Review

1. Page 858, Definition B.16. $x \in \mathbb{C}$ should be $x \in \mathbb{C}^d$. The proper dimension $d$ is missing.

2. Page 865, Definition B.42. The opening sentence "Let $[a, b] \subset \mathbb{R}$ be a compact integral." Should be "Let $[a, b] \subset \mathbb{R}$ be a compact interval.". In other words "integral" should be changed to "interval."