

Appendix A

Linear algebra and numerical techniques

A.1	Matrix decompositions	736
A.1.1	Singular value decomposition	736
A.1.2	Eigenvalue decomposition	737
A.1.3	QR factorization	740
A.1.4	Cholesky factorization	741
A.2	Linear least squares	742
A.2.1	Total least squares	744
A.3	Non-linear least squares	746
A.4	Direct sparse matrix techniques	747
A.4.1	Variable reordering	748
A.5	Iterative techniques	748
A.5.1	Conjugate gradient	749
A.5.2	Preconditioning	751
A.5.3	Multigrid	753

In this appendix, we introduce some elements of linear algebra and numerical techniques that are used elsewhere in the book. We start with some basic decompositions in matrix algebra, including the singular value decomposition (SVD), eigenvalue decompositions, and other matrix decompositions (factorizations). Next, we look at the problem of linear least squares, which can be solved using either the QR decomposition or normal equations. This is followed by non-linear least squares, which arise when the measurement equations are not linear in the unknowns or when robust error functions are used. Such problems require iteration to find a solution. Next, we look at direct solution (factorization) techniques for sparse problems, where the ordering of the variables can have a large influence on the computation and memory requirements. Finally, we discuss iterative techniques for solving large linear (or linearized) least squares problems. Good general references for much of this material include the work by Björck (1996), Golub and Van Loan (1996), Trefethen and Bau (1997), Meyer (2000), Nocedal and Wright (2006), and Björck and Dahlquist (2010).

A note on vector and matrix indexing. To be consistent with the rest of the book and with the general usage in the computer science and computer vision communities, I adopt a 0-based indexing scheme for vector and matrix element indexing. Please note that most mathematical textbooks and papers use 1-based indexing, so you need to be aware of the differences when you read this book.

Software implementations. Highly optimized and tested libraries corresponding to the algorithms described in this appendix are readily available and are listed in Appendix C.2.

A.1 Matrix decompositions

In order to better understand the structure of matrices and more stably perform operations such as inversion and system solving, a number of decompositions (or factorizations) can be used. In this section, we review singular value decomposition (SVD), eigenvalue decomposition, QR factorization, and Cholesky factorization.

A.1.1 Singular value decomposition

One of the most useful decompositions in matrix algebra is the *singular value decomposition* (SVD), which states that any real-valued $M \times N$ matrix \mathbf{A} can be written as

$$\mathbf{A}_{M \times N} = \mathbf{U}_{M \times P} \mathbf{\Sigma}_{P \times P} \mathbf{V}_{P \times N}^T \quad (\text{A.1})$$

$$= \left[\begin{array}{c|c|c} \mathbf{u}_0 & \cdots & \mathbf{u}_{p-1} \end{array} \right] \left[\begin{array}{ccc} \sigma_0 & & \\ & \ddots & \\ & & \sigma_{p-1} \end{array} \right] \left[\begin{array}{c} \mathbf{v}_0^T \\ \cdots \\ \mathbf{v}_{p-1}^T \end{array} \right],$$

where $P = \min(M, N)$. The matrices \mathbf{U} and \mathbf{V} are orthonormal, i.e., $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ and $\mathbf{V}^T \mathbf{V} = \mathbf{I}$, and so are their column vectors,

$$\mathbf{u}_i \cdot \mathbf{u}_j = \mathbf{v}_i \cdot \mathbf{v}_j = \delta_{ij}. \quad (\text{A.2})$$

The singular values are all non-negative and can be ordered in decreasing order

$$\sigma_0 \geq \sigma_1 \geq \cdots \geq \sigma_{p-1} \geq 0. \quad (\text{A.3})$$

A geometric intuition for the SVD of a matrix \mathbf{A} can be obtained by re-writing $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ in (A.2) as

$$\mathbf{A} \mathbf{V} = \mathbf{U} \mathbf{\Sigma} \quad \text{or} \quad \mathbf{A} \mathbf{v}_j = \sigma_j \mathbf{u}_j. \quad (\text{A.4})$$

This formula says that the matrix \mathbf{A} takes any basis vector \mathbf{v}_j and maps it to a direction \mathbf{u}_j with length σ_j , as shown in Figure A.1

If only the first r singular values are positive, the matrix \mathbf{A} is of *rank* r and the index p in the SVD decomposition (A.2) can be replaced by r . (In other words, we can drop the last $p - r$ columns of \mathbf{U} and \mathbf{V} .)

An important property of the singular value decomposition of a matrix (also true for the eigenvalue decomposition of a real symmetric non-negative definite matrix) is that if we truncate the expansion

$$\mathbf{A} = \sum_{j=0}^t \sigma_j \mathbf{u}_j \mathbf{v}_j^T, \quad (\text{A.5})$$

we obtain the best possible least squares approximation to the original matrix \mathbf{A} . This is used both in eigenface-based face recognition systems (Section 14.2.1) and in the separable approximation of convolution kernels (3.21).

A.1.2 Eigenvalue decomposition

If the matrix \mathbf{C} is symmetric ($m = n$),¹ it can be written as an eigenvalue decomposition,

$$\begin{aligned} \mathbf{C} &= \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T = \left[\begin{array}{c|c|c} \mathbf{u}_0 & \cdots & \mathbf{u}_{n-1} \end{array} \right] \left[\begin{array}{ccc} \lambda_0 & & \\ & \ddots & \\ & & \lambda_{n-1} \end{array} \right] \left[\begin{array}{c} \mathbf{u}_0^T \\ \cdots \\ \mathbf{u}_{n-1}^T \end{array} \right] \\ &= \sum_{i=0}^{n-1} \lambda_i \mathbf{u}_i \mathbf{u}_i^T. \end{aligned} \quad (\text{A.6})$$

¹ In this appendix, we denote symmetric matrices using \mathbf{C} and general rectangular matrices using \mathbf{A} .

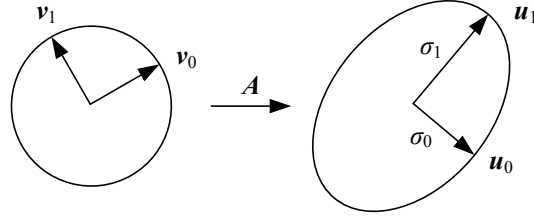


Figure A.1 The action of a matrix A can be visualized by thinking of the domain as being spanned by a set of orthonormal vectors v_j , each of which is transformed to a new orthogonal vector u_j with a length σ_j . When A is interpreted as a covariance matrix and its eigenvalue decomposition is performed, each of the u_j axes denote a principal direction (component) and each σ_j denotes one standard deviation along that direction.

(The eigenvector matrix U is sometimes written as Φ and the eigenvectors u as ϕ .) In this case, the eigenvalues

$$\lambda_0 \geq \lambda_1 \geq \dots \geq \lambda_{n-1} \quad (\text{A.7})$$

can be both positive and negative.²

A special case of the symmetric matrix C occurs when it is constructed as the sum of a number of outer products

$$C = \sum_i a_i a_i^T = A A^T, \quad (\text{A.8})$$

which often occurs when solving least squares problems (Appendix A.2), where the matrix A consists of all the a_i column vectors stacked side-by-side. In this case, we are guaranteed that all of the eigenvalues λ_i are non-negative. The associated matrix C is *positive semi-definite*

$$x^T C x \geq 0, \quad \forall x. \quad (\text{A.9})$$

If the matrix C is of full rank, the eigenvalues are all positive and the matrix is called *symmetric positive definite* (SPD).

Symmetric positive semi-definite matrices also arise in the statistical analysis of data, since they represent the *covariance* of a set of $\{x_i\}$ points around their mean \bar{x} ,

$$C = \frac{1}{n} \sum_i (x_i - \bar{x})(x_i - \bar{x})^T. \quad (\text{A.10})$$

In this case, performing the eigenvalue decomposition is known as *principal component analysis* (PCA), since it models the principal directions (and magnitudes) of variation of the point

² Eigenvalue decompositions can be computed for non-symmetric matrices but the eigenvalues and eigenvectors can have complex entries in that case.

distribution around their mean, as shown in Section 5.1.1 (5.13–5.15), Section 14.2.1 (14.9), and Appendix B.1.1 (B.10). Figure A.1 shows how the principal components of the covariance matrix C denote the principal axes \mathbf{u}_j of the uncertainty ellipsoid corresponding to this point distribution and how the $\sigma_j = \sqrt{\lambda_j}$ denote the standard deviations along each axis.

The eigenvalues and eigenvectors of C and the singular values and singular vectors of A are closely related. Given

$$A = U\Sigma V^T, \quad (\text{A.11})$$

we get

$$C = AA^T = U\Sigma V^T V\Sigma U^T = U\Lambda U^T. \quad (\text{A.12})$$

From this, we see that $\lambda_i = \sigma_i^2$ and that the left singular vectors of A are the eigenvectors of C .

This relationship gives us an efficient method for computing the eigenvalue decomposition of large matrices that are rank deficient, such as the scatter matrices observed in computing eigenfaces (Section 14.2.1). Observe that the covariance matrix C in (14.9) is exactly the same as C in (A.8). Note also that the individual difference-from-mean images $\mathbf{a}_i = \mathbf{x}_i - \bar{\mathbf{x}}$ are long vectors of length P (the number of pixels in the image), while the total number of exemplars N (the number of faces in the training database) is much smaller. Instead of forming $C = AA^T$, which is $P \times P$, we form the matrix

$$\hat{C} = A^T A, \quad (\text{A.13})$$

which is $N \times N$. (This involves taking the dot product between every pair of difference images \mathbf{a}_i and \mathbf{a}_j .) The eigenvalues of \hat{C} are the squared singular values of A , namely Σ^2 , and are hence also the eigenvalues of C . The eigenvectors of \hat{C} are the right singular vectors V of A , from which the desired eigenfaces U , which are the left singular vectors of A , can be computed as

$$U = AV\Sigma^{-1}. \quad (\text{A.14})$$

This final step is essentially computing the eigenfaces as linear combinations of the difference images (Turk and Pentland 1991a). If you have access to a high-quality linear algebra package such as LAPACK, routines for efficiently computing a small number of the left singular vectors and singular values of rectangular matrices such as A are usually provided (Appendix C.2). However, if storing all of the images in memory is prohibitive, the construction of \hat{C} in (A.13) can be used instead.

How can eigenvalue and singular value decompositions actually be computed? Notice that an eigenvector is defined by the equation

$$\lambda_i \mathbf{u}_i = C\mathbf{u}_i \quad \text{or} \quad (\lambda_i I - C)\mathbf{u}_i = 0. \quad (\text{A.15})$$

(This can be derived from (A.6) by post-multiplying both sides by \mathbf{u}_i .) Since the latter equation is *homogeneous*, i.e., it has a zero right-hand-side, it can only have a non-zero (non-trivial) solution for \mathbf{u}_i if the system is rank deficient, i.e.,

$$|(\lambda \mathbf{I} - \mathbf{C})| = 0. \quad (\text{A.16})$$

Evaluating this determinant yields a *characteristic* polynomial equation in λ , which can be solved for small problems, e.g., 2×2 or 3×3 matrices, in closed form.

For larger matrices, iterative algorithms that first reduce the matrix \mathbf{C} to a real symmetric tridiagonal form using orthogonal transforms and then perform QR iterations are normally used (Golub and Van Loan 1996; Trefethen and Bau 1997; Björck and Dahlquist 2010). Since these techniques are rather involved, it is best to use a linear algebra package such as LAPACK (Anderson, Bai, Bischof *et al.* 1999)—see Appendix C.2.

Factorization with missing data requires different kinds of iterative algorithms, which often involve either hallucinating the missing terms or minimizing some weighted reconstruction metric, which is intrinsically much more challenging than regular factorization. This area has been widely studied in computer vision (Shum, Ikeuchi, and Reddy 1995; De la Torre and Black 2003; Huynh, Hartley, and Heyden 2003; Buchanan and Fitzgibbon 2005; Gross, Matthews, and Baker 2006; Torresani, Hertzmann, and Bregler 2008) and is sometimes called *generalized PCA*. However, this term is also sometimes used to denote algebraic subspace clustering techniques, which is the subject of a forthcoming monograph by Vidal, Ma, and Sastry (2010).

A.1.3 QR factorization

A widely used technique for stably solving poorly conditioned least squares problems (Björck 1996) and as the basis of more complex algorithms, such as computing the SVD and eigenvalue decompositions, is the QR factorization,

$$\mathbf{A} = \mathbf{Q}\mathbf{R}, \quad (\text{A.17})$$

where \mathbf{Q} is an *orthonormal* (or *unitary*) matrix $\mathbf{Q}\mathbf{Q}^T = \mathbf{I}$ and \mathbf{R} is upper triangular.³ In computer vision, QR can be used to convert a camera matrix into a rotation matrix and an upper-triangular calibration matrix (6.35) and also in various self-calibration algorithms (Section 7.2.2). The most common algorithms for computing QR decompositions, modified Gram–Schmidt, Householder transformations, and Givens rotations, are described by Golub and Van Loan (1996), Trefethen and Bau (1997), and Björck and Dahlquist (2010) and are

³ The term “R” comes from the German name for the lower–upper (LU) decomposition, which is LR for “links” and “rechts” (left and right of the diagonal).

```

procedure Cholesky( $C, R$ ):

     $R = C$ 

    for  $i = 0 \dots n - 1$ 

        for  $j = i + 1 \dots n - 1$ 

             $R_{j,j:n-1} = R_{j,j:n-1} - r_{ij}r_{ii}^{-1}R_{i,j:n-1}$ 

         $R_{i,i:n-1} = r_{ii}^{-1/2}R_{i,i:n-1}$ 

```

Algorithm A.1 Cholesky decomposition of the matrix C into its upper triangular form R .

also found in LAPACK. Unlike the SVD and eigenvalue decompositions, QR factorization does not require iteration and can be computed exactly in $O(MN^2 + N^3)$ operations, where M is the number of rows and N is the number of columns (for a tall matrix).

A.1.4 Cholesky factorization

Cholesky factorization can be applied to any symmetric positive definite matrix C to convert it into a product of symmetric lower and upper triangular matrices,

$$C = LL^T = R^T R, \quad (\text{A.18})$$

where L is a lower-triangular matrix and R is an upper-triangular matrix. Unlike Gaussian elimination, which may require pivoting (row and column reordering) or may become unstable (sensitive to roundoff errors or reordering), Cholesky factorization remains stable for positive definite matrices, such as those that arise from normal equations in least squares problems (Appendix A.2). Because of the form of (A.18), the matrices L and R are sometimes called *matrix square roots*.⁴

The algorithm to compute an upper triangular Cholesky decomposition of C is a straightforward symmetric generalization of Gaussian elimination and is based on the decomposition (Björck 1996; Golub and Van Loan 1996)

$$C = \begin{bmatrix} \gamma & \mathbf{c}^T \\ \mathbf{c} & C_{11} \end{bmatrix} \quad (\text{A.19})$$

$$= \begin{bmatrix} \gamma^{1/2} & \mathbf{0}^T \\ \mathbf{c}\gamma^{-1/2} & \mathbf{I} \end{bmatrix} \begin{bmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & C_{11} - \mathbf{c}\gamma^{-1}\mathbf{c}^T \end{bmatrix} \begin{bmatrix} \gamma^{1/2} & \gamma^{-1/2}\mathbf{c}^T \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \quad (\text{A.20})$$

⁴ In fact, there exists a whole family of matrix square roots. Any matrix of the form LQ or QR , where Q is a unitary matrix, is a square root of C .

$$= \mathbf{R}_0^T \mathbf{C}_1 \mathbf{R}_0, \quad (\text{A.21})$$

which, through recursion, can be turned into

$$\mathbf{C} = \mathbf{R}_0^T \dots \mathbf{R}_{n-1}^T \mathbf{R}_{n-1} \dots \mathbf{R}_0 = \mathbf{R}^T \mathbf{R}. \quad (\text{A.22})$$

Algorithm A.1 provides a more procedural definition, which can store the upper-triangular matrix \mathbf{R} in the same space as \mathbf{C} , if desired. The total operation count for Cholesky factorization is $O(N^3)$ for a dense matrix but can be significantly lower for sparse matrices with low fill-in (Appendix A.4).

Note that Cholesky decomposition can also be applied to block-structured matrices, where the term γ in (A.19) is now a square block sub-matrix and \mathbf{c} is a rectangular matrix (Golub and Van Loan 1996). The computation of square roots can be avoided by leaving the γ on the diagonal of the middle factor in (A.20), which results in the $\mathbf{C} = \mathbf{L}\mathbf{D}\mathbf{L}^T$ factorization, where \mathbf{D} is a diagonal matrix. However, since square roots are relatively fast on modern computers, this is not worth the bother and Cholesky factorization is usually preferred.

A.2 Linear least squares

Least squares fitting problems are pervasive in computer vision. For example, the alignment of images based on matching feature points involves the minimization of a squared distance objective function (6.2),

$$E_{\text{LS}} = \sum_i \|\mathbf{r}_i\|^2 = \sum_i \|\mathbf{f}(\mathbf{x}_i; \mathbf{p}) - \mathbf{x}'_i\|^2, \quad (\text{A.23})$$

where

$$\mathbf{r}_i = \mathbf{f}(\mathbf{x}_i; \mathbf{p}) - \mathbf{x}'_i = \hat{\mathbf{x}}'_i - \tilde{\mathbf{x}}'_i \quad (\text{A.24})$$

is the *residual* between the measured location $\hat{\mathbf{x}}'_i$ and its corresponding current *predicted* location $\tilde{\mathbf{x}}'_i = \mathbf{f}(\mathbf{x}_i; \mathbf{p})$. More complex versions of least squares problems, such as large-scale structure from motion (Section 7.4), may involve the minimization of functions of thousands of variables. Even problems such as image filtering (Section 3.4.3) and regularization (Section 3.7.1) may involve the minimization of sums of squared errors.

Figure A.2a shows an example of a simple least squares line fitting problem, where the quantities being estimated are the line equation parameters (m, b) . When the sampled vertical values y_i are assumed to be noisy versions of points on the line $y = mx + b$, the optimal estimates for (m, b) can be found by minimizing the squared vertical residuals

$$E_{\text{VLS}} = \sum_i |y_i - (mx_i + b)|^2. \quad (\text{A.25})$$

Note that the function being fitted need not itself be linear to use linear least squares. All that is required is that the function be linear in the unknown parameters. For example, polynomial fitting can be written as

$$E_{\text{PLS}} = \sum_i |y_i - (\sum_{j=0}^p a_j x_i^j)|^2, \quad (\text{A.26})$$

while sinusoid fitting with unknown amplitude A and phase ϕ (but known frequency f) can be written as

$$E_{\text{SLS}} = \sum_i |y_i - A \sin(2\pi f x_i + \phi)|^2 = \sum_i |y_i - (B \sin 2\pi f x_i + C \cos 2\pi f x_i)|^2, \quad (\text{A.27})$$

which is linear in (B, C) .

In general, it is more common to denote the unknown parameters using \mathbf{x} and to write the general form of linear least squares as⁵

$$E_{\text{LLS}} = \sum_i |a_i \mathbf{x} - b_i|^2 = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2. \quad (\text{A.28})$$

Expanding the above equation gives us

$$E_{\text{LLS}} = \mathbf{x}^T (\mathbf{A}^T \mathbf{A}) \mathbf{x} - 2\mathbf{x}^T (\mathbf{A}^T \mathbf{b}) + \|\mathbf{b}\|^2, \quad (\text{A.29})$$

whose minimum value for \mathbf{x} can be found by solving the associated *normal equations* (Björck 1996; Golub and Van Loan 1996)

$$(\mathbf{A}^T \mathbf{A}) \mathbf{x} = \mathbf{A}^T \mathbf{b}. \quad (\text{A.30})$$

The preferred way to solve the normal equations is to use Cholesky factorization. Let

$$\mathbf{C} = \mathbf{A}^T \mathbf{A} = \mathbf{R}^T \mathbf{R}, \quad (\text{A.31})$$

where \mathbf{R} is the upper-triangular Cholesky factor of the Hessian \mathbf{C} , and

$$\mathbf{d} = \mathbf{A}^T \mathbf{b}. \quad (\text{A.32})$$

After factorization, the solution for \mathbf{x} can be obtained as

$$\mathbf{R}^T \mathbf{z} = \mathbf{d}, \quad \mathbf{R}\mathbf{x} = \mathbf{z}, \quad (\text{A.33})$$

which involves the solution of two triangular systems, i.e., forward and backward substitution (Björck 1996).

⁵ Be extra careful in interpreting the variable names here. In the 2D line-fitting example, x is used to denote the horizontal axis, but in the general least squares problem, $\mathbf{x} = (m, b)$ denotes the unknown parameter vector.

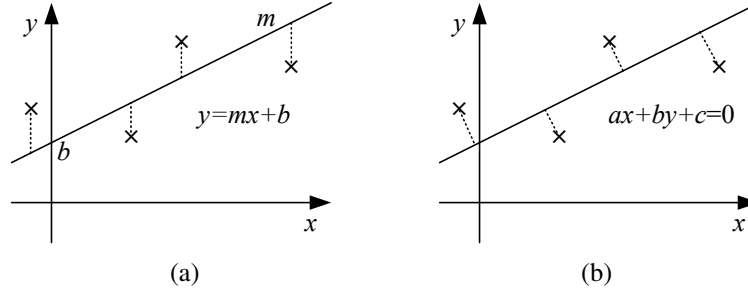


Figure A.2 Least squares regression. (a) The line $y = mx + b$ is fit to the four noisy data points, $\{(x_i, y_i)\}$, denoted by \times by minimizing the squared vertical residuals between the data points and the line, $\sum_i \|y_i - (mx_i + b)\|^2$. (b) When the measurements $\{(x_i, y_i)\}$ are assumed to have noise in all directions, the sum of orthogonal squared distances to the line $\sum_i \|ax_i + by_i + c\|^2$ is minimized using total least squares.

In cases where the least squares problem is numerically poorly conditioned (which should generally be avoided by adding sufficient regularization or prior knowledge about the parameters, (Appendix A.3)), it is possible to use QR factorization or SVD directly on the matrix \mathbf{A} (Björck 1996; Golub and Van Loan 1996; Trefethen and Bau 1997; Nocedal and Wright 2006; Björck and Dahlquist 2010), e.g.,

$$\mathbf{A}\mathbf{x} = \mathbf{Q}\mathbf{R}\mathbf{x} = \mathbf{b} \quad \longrightarrow \quad \mathbf{R}\mathbf{x} = \mathbf{Q}^T \mathbf{b}. \quad (\text{A.34})$$

Note that the upper triangular matrices \mathbf{R} produced by the Cholesky factorization of $\mathbf{C} = \mathbf{A}^T \mathbf{A}$ and the QR factorization of \mathbf{A} are the same, but that solving (A.34) is generally more stable (less sensitive to roundoff error) but slower (by a constant factor).

A.2.1 Total least squares

In some problems, e.g., when performing geometric line fitting in 2D images or 3D plane fitting to point cloud data, instead of having measurement error along one particular axis, the measured points have uncertainty in all directions, which is known as the *errors-in-variables* model (Van Huffel and Lemmerling 2002; Matei and Meer 2006). In this case, it makes more sense to minimize a set of homogeneous squared errors of the form

$$E_{\text{TLS}} = \sum_i (\mathbf{a}_i \mathbf{x})^2 = \|\mathbf{A}\mathbf{x}\|^2, \quad (\text{A.35})$$

which is known as *total least squares* (TLS) (Van Huffel and Vandewalle 1991; Björck 1996; Golub and Van Loan 1996; Van Huffel and Lemmerling 2002).

The above error metric has a trivial minimum solution at $\mathbf{x} = 0$ and is, in fact, homogeneous in \mathbf{x} . For this reason, we augment this minimization problem with the requirement that $\|\mathbf{x}\|^2 = 1$, which results in the eigenvalue problem

$$\mathbf{x} = \arg \min_{\mathbf{x}} \mathbf{x}^T (\mathbf{A}^T \mathbf{A}) \mathbf{x} \quad \text{such that} \quad \|\mathbf{x}\|^2 = 1. \quad (\text{A.36})$$

The value of \mathbf{x} that minimizes this constrained problem is the eigenvector associated with the smallest eigenvalue of $\mathbf{A}^T \mathbf{A}$. This is the same as the last right singular vector of \mathbf{A} , since

$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}, \quad (\text{A.37})$$

$$\mathbf{A}^T \mathbf{A} = \mathbf{V} \mathbf{\Sigma}^2 \mathbf{V}, \quad (\text{A.38})$$

$$\mathbf{A}^T \mathbf{A} \mathbf{v}_k = \sigma_k^2, \quad (\text{A.39})$$

which is minimized by selecting the smallest σ_k value.

Figure A.2b shows a line fitting problem where, in this case, the measurement errors are assumed to be isotropic in (x, y) . The solution for the best line equation $ax + by + c = 0$ is found by minimizing

$$E_{\text{TLS-2D}} = \sum_i (ax_i + by_i + c)^2, \quad (\text{A.40})$$

i.e., finding the eigenvector associated with the smallest eigenvalue of⁶

$$\mathbf{C} = \mathbf{A}^T \mathbf{A} = \sum_i \begin{bmatrix} x_i \\ y_i \\ 1 \end{bmatrix} \begin{bmatrix} x_i & y_i & 1 \end{bmatrix}. \quad (\text{A.41})$$

Notice, however, that minimizing $\sum_i (\mathbf{a}_i \mathbf{x})^2$ in (A.35) is only statistically optimal (Appendix B.1.1) if all of the measured terms in the \mathbf{a}_i , e.g., the $(x_i, y_i, 1)$ measurements, have equal noise. This is definitely not the case in the line-fitting example of Figure A.2b (A.40), since the 1 values are noise-free. To mitigate this, we first subtract the mean x and y values from all the measured points

$$\hat{x}_i = x_i - \bar{x} \quad (\text{A.42})$$

$$\hat{y}_i = y_i - \bar{y} \quad (\text{A.43})$$

and then fit the 2D line equation $a(x - \bar{x}) + b(y - \bar{y}) = 0$ by minimizing

$$E_{\text{TLS-2Dm}} = \sum_i (a\hat{x}_i + b\hat{y}_i)^2. \quad (\text{A.44})$$

⁶ Again, be careful with the variable names here. The measurement equation is $\mathbf{a}_i = (x_i, y_i, 1)$ and the unknown parameters are $\mathbf{x} = (a, b, c)$.

The more general case where each individual measurement component can have different noise level, as is the case in estimating essential and fundamental matrices (Section 7.2), is called the *heteroscedastic errors-in-variable* (HEIV) model and is discussed by Matei and Meer (2006).

A.3 Non-linear least squares

In many vision problems, such as structure from motion, the least squares problem formulated in (A.23) involves functions $\mathbf{f}(\mathbf{x}_i; \mathbf{p})$ that are *not* linear in the unknown parameters \mathbf{p} . This problem is known as *non-linear least squares* or *non-linear regression* (Björck 1996; Madsen, Nielsen, and Tingleff 2004; Nocedal and Wright 2006). It is usually solved by iteratively re-linearizing (A.23) around the current estimate of \mathbf{p} using the gradient derivative (Jacobian) $\mathbf{J} = \partial \mathbf{f} / \partial \mathbf{p}$ and computing an incremental improvement $\Delta \mathbf{p}$.

As shown in Equations (6.13–6.17), this results in

$$E_{\text{NLS}}(\Delta \mathbf{p}) = \sum_i \|\mathbf{f}(\mathbf{x}_i; \mathbf{p} + \Delta \mathbf{p}) - \mathbf{x}'_i\|^2 \quad (\text{A.45})$$

$$\approx \sum_i \|\mathbf{J}(\mathbf{x}_i; \mathbf{p}) \Delta \mathbf{p} - \mathbf{r}_i\|^2, \quad (\text{A.46})$$

where the Jacobians $\mathbf{J}(\mathbf{x}_i; \mathbf{p})$ and residual vectors \mathbf{r}_i play the same role in forming the normal equations as \mathbf{a}_i and \mathbf{b}_i in (A.28).

Because the above approximation only holds near a local minimum or for small values of $\Delta \mathbf{p}$, the update $\mathbf{p} \leftarrow \mathbf{p} + \Delta \mathbf{p}$ may not always decrease the summed square residual error (A.45). One way to mitigate this problem is to take a smaller step,

$$\mathbf{p} \leftarrow \mathbf{p} + \alpha \Delta \mathbf{p}, \quad 0 < \alpha \leq 1. \quad (\text{A.47})$$

A simple way to determine a reasonable value of α is to start with 1 and successively halve the value, which is a simple form of *line search* (Al-Baali and Fletcher. 1986; Björck 1996; Nocedal and Wright 2006).

Another approach to ensuring a downhill step in error is to add a diagonal damping term to the approximate Hessian

$$\mathbf{C} = \sum_i \mathbf{J}^T(\mathbf{x}_i) \mathbf{J}(\mathbf{x}_i), \quad (\text{A.48})$$

i.e., to solve

$$[\mathbf{C} + \lambda \text{diag}(\mathbf{C})] \Delta \mathbf{p} = \mathbf{d}, \quad (\text{A.49})$$

where

$$\mathbf{d} = \sum_i \mathbf{J}^T(\mathbf{x}_i) \mathbf{r}_i, \quad (\text{A.50})$$

which is called a *damped Gauss–Newton* method. The damping parameter λ is increased if the squared residual is not decreasing as fast as expected, i.e., as predicted by (A.46), and is decreased if the expected decrease is obtained (Madsen, Nielsen, and Tingleff 2004). The combination of the Newton (first-order Taylor series) approximation (A.46) and the adaptive damping parameter λ is commonly known as the Levenberg–Marquardt algorithm (Levenberg 1944; Marquardt 1963) and is an example of more general *trust region methods*, which are discussed in more detail in (Björck 1996; Conn, Gould, and Toint 2000; Madsen, Nielsen, and Tingleff 2004; Nocedal and Wright 2006).

When the initial solution is far away from its quadratic region of convergence around a local minimum, *large residual methods*, e.g., *Newton-type methods*, which add a second-order term to the Taylor series expansion in (A.46), may converge faster. Quasi-Newton methods such as BFGS, which require only gradient evaluations, can also be useful if memory size is an issue. Such techniques are discussed in textbooks and papers on numerical optimization (Toint 1987; Björck 1996; Conn, Gould, and Toint 2000; Nocedal and Wright 2006).

A.4 Direct sparse matrix techniques

Many optimization problems in computer vision, such as bundle adjustment (Szeliski and Kang 1994; Triggs, McLauchlan, Hartley *et al.* 1999; Hartley and Zisserman 2004; Snavely, Seitz, and Szeliski 2008b; Agarwal, Snavely, Simon *et al.* 2009) have Jacobian and (approximate) Hessian matrices that are extremely sparse (Section 7.4.1). For example, Figure 7.9a shows the *bipartite* model typical of structure from motion problems, in which most points are only observed by a subset of the cameras, which results in the sparsity patterns for the Jacobian and Hessian shown in Figure 7.9b–c.

Whenever the Hessian matrix is sparse enough, it is more efficient to use sparse Cholesky factorization instead of regular Cholesky factorization. In such sparse direct techniques, the Hessian matrix \mathbf{C} and its associated Cholesky factor \mathbf{R} are stored in *compressed form*, in which the amount of storage is proportional to the number of (potentially) non-zero entries (Björck 1996; Davis 2006).⁷ Algorithms for computing the non-zero elements in \mathbf{C} and \mathbf{R} from the sparsity pattern of the Jacobian matrix \mathbf{J} are given by Björck (1996, Section 6.4), and algorithms for computing the numerical Cholesky and QR decompositions (once the sparsity pattern has been computed and storage allocated) are discussed by Björck (1996, Section 6.5).

⁷ For example, you can store a list of (i, j, c_{ij}) triples. One example of such a scheme is *compressed sparse row (CSR)* storage. An alternative storage method called *skyline*, which stores adjacent vertical spans of non-zero elements (Bathe 2007), is sometimes used in finite element analysis. Banded systems such as snakes (5.3) can store just the non-zero band elements (Björck 1996, Section 6.2) and can be solved in $O(nb^2)$, where n is the number of variables and b is the bandwidth.

A.4.1 Variable reordering

The key to efficiently solving sparse problems using direct (non-iterative) techniques is to determine an efficient *ordering* for the variables, which reduces the amount of *fill-in*, i.e., the number of non-zero entries in \mathbf{R} that were zero in the original \mathbf{C} matrix. We already saw in Section 7.4.1 how storing the more numerous 3D point parameters before the camera parameters and using the Schur complement (7.56) results in a more efficient algorithm. Similarly, sorting parameters by time in video-based reconstruction problems usually results in lower fill-in. Furthermore, any problem whose adjacency graph (the graph corresponding to the sparsity pattern) is a tree can be solved in linear time with an appropriate reordering of the variables (putting all the children before their parents). All of these are examples of good reordering techniques.

In the general case of unstructured data, there are many heuristics available to find good reorderings (Björck 1996; Davis 2006).⁸ For general adjacency (sparsity) graphs, *minimum degree orderings* generally produce good results. For planar graphs, which often arise on image or spline grids (Section 8.3), *nested dissection*, which recursively splits the graph into two equal halves along a *frontier* (or boundary) of small size, generally works well. Such *domain decomposition* (or *multi-frontal*) techniques also enable the use of parallel processing, since independent sub-graphs can be processed in parallel on separate processors (Davis 2008).

The overall set of steps used to perform the direct solution of sparse least squares problems are summarized in Algorithm A.2, which is a modified version of Algorithm 6.6.1 by Björck (1996, Section 6.6)). If a series of related least squares problems is being solved, as is the case in iterative non-linear least squares (Appendix A.3), steps 1–3 can be performed ahead of time and reused for each new invocation with different \mathbf{C} and \mathbf{d} values. When the problem is block-structured, as is the case in structure from motion where point (structure) variables have dense 3×3 sub-entries in \mathbf{C} and cameras have 6×6 (or larger) entries, the cost of performing the reordering computation is small compared to the actual numerical factorization, which can benefit from block-structured matrix operations (Golub and Van Loan 1996). It is also possible to apply sparse reordering and multifrontal techniques to QR factorization (Davis 2008), which may be preferable when the least squares problems are poorly conditioned.

A.5 Iterative techniques

When problems become large, the amount of memory required to store the Hessian matrix \mathbf{C} and its factor \mathbf{R} , and the amount of time it takes to compute the factorization, can become prohibitively large, especially when there are large amounts of fill-in. This is often

⁸Finding the optimal reordering with minimal fill-in is provably NP-hard.

procedure *SparseCholeskySolve*(C, d):

1. Determine symbolically the structure of C , i.e., the adjacency graph.
2. (Optional) Compute a reordering for the variables, taking into account any block structure inherent in the problem.
3. Determine the fill-in pattern for R and allocate the compressed storage for R as well as storage for the permuted right hand side \hat{d} .
4. Copy the elements of C and d into R and \hat{d} , permuting the values according to the computed ordering.
5. Perform the numerical factorization of R using Algorithm A.1.
6. Solve the factored system (A.33), i.e.,

$$R^T z = \hat{d}, \quad Rx = z.$$

7. Return the solution x , after undoing the permutation.

Algorithm A.2 Sparse least squares using a sparse Cholesky decomposition of the matrix C .

the case with image processing problems defined on pixel grids, since, even with the optimal reordering (nested dissection) the amount of fill can still be large.

A preferable approach to solving such linear systems is to use iterative techniques, which compute a series of estimates that converge to the final solution, e.g., by taking a series of downhill steps in an energy function such as (A.29).

A large number of iterative techniques have been developed over the years, including such well-known algorithms as successive overrelaxation and multi-grid. These are described in specialized textbooks on iterative solution techniques (Axelsson 1996; Saad 2003) as well as in more general books on numerical linear algebra and least squares techniques (Björck 1996; Golub and Van Loan 1996; Trefethen and Bau 1997; Nocedal and Wright 2006; Björck and Dahlquist 2010).

A.5.1 Conjugate gradient

The iterative solution technique that often performs best is conjugate gradient descent, which takes a series of downhill steps that are *conjugate* to each other with respect to the C matrix,

i.e., if the \mathbf{u} and \mathbf{v} descent directions satisfy $\mathbf{u}^T \mathbf{C} \mathbf{v} = 0$. In practice, conjugate gradient descent outperforms other kinds of gradient descent algorithm because its convergence rate is proportional to the square root of the *condition number* of \mathbf{C} instead of the condition number itself.⁹ Shewchuk (1994) provides a nice introduction to this topic, with clear intuitive explanations of the reasoning behind the conjugate gradient algorithm and its performance.

Algorithm A.3 describes the conjugate gradient algorithm and its related least squares counterpart, which can be used when the original set of least squares linear equations are available in the form of $\mathbf{A} \mathbf{x} = \mathbf{b}$ (A.28). While it is easy to convince yourself that the two forms are mathematically equivalent, the least squares form is preferable if rounding errors start to affect the results because of poor conditioning. It may also be preferable if, due to the sparsity structure of \mathbf{A} , multiplies with the original \mathbf{A} matrix are faster or more space efficient than multiplies with \mathbf{C} .

The conjugate gradient algorithm starts by computing the current residual $\mathbf{r}_0 = \mathbf{d} - \mathbf{C} \mathbf{x}_0$, which is the direction of steepest descent of the energy function (A.28). It sets the original descent direction $\mathbf{p}_0 = \mathbf{r}_0$. Next, it multiplies the descent direction by the quadratic form (Hessian) matrix \mathbf{C} and combines this with the residual to estimate the optimal step size α_k . The solution vector \mathbf{x}_k and the residual vector \mathbf{r}_k are then updated using this step size. (Notice how the least squares variant of the conjugate gradient algorithm splits the multiplication by the $\mathbf{C} = \mathbf{A}^T \mathbf{A}$ matrix across steps 4 and 8.) Finally, a new search direction is calculated by first computing a factor β as the ratio of current to previous residual magnitudes. The new search direction \mathbf{p}_{k+1} is then set to the residual plus β times the old search direction \mathbf{p}_k , which keeps the directions conjugate with respect to \mathbf{C} .

It turns out that conjugate gradient descent can also be directly applied to non-quadratic energy functions, e.g., those arising from non-linear least squares (Appendix A.3). Instead of explicitly forming a local quadratic approximation \mathbf{C} and then computing residuals \mathbf{r}_k , non-linear conjugate gradient descent computes the gradient of the energy function E (A.45) directly inside each iteration and uses it to set the search direction (Nocedal and Wright 2006). Since the quadratic approximation to the energy function may not exist or may be inaccurate, line search is often used to determine the step size α_k . Furthermore, to compensate for errors in finding the true function minimum, alternative formulas for β_{k+1} such as Polak–Ribière,

$$\beta_{k+1} = \frac{\nabla E(\mathbf{x}_{k+1})[\nabla E(\mathbf{x}_{k+1}) - \nabla E(\mathbf{x}_k)]}{\|\nabla E(\mathbf{x}_k)\|^2} \quad (\text{A.51})$$

are often used (Nocedal and Wright 2006).

⁹ The condition number $\kappa(\mathbf{C})$ is the ratio of the largest and smallest eigenvalues of \mathbf{C} . The actual convergence rate depends on the clustering of the eigenvalues, as discussed in the references cited in this section.

ConjugateGradient(C, d, x_0)

1. $r_0 = d - Cx_0$
2. $p_0 = r_0$
3. **for** $k = 0 \dots$
4. $w_k = Cp_k$
5. $\alpha_k = \|r_k\|^2 / (p_k \cdot w_k)$
6. $x_{k+1} = x_k + \alpha_k p_k$
7. $r_{k+1} = r_k - \alpha_k w_k$
- 8.
9. $\beta_{k+1} = \|r_{k+1}\|^2 / \|r_k\|^2$
10. $p_{k+1} = r_{k+1} + \beta_k p_k$

ConjugateGradientLS(A, b, x_0)

1. $q_0 = b - Ax_0, \quad r_0 = A^T q_0$
2. $p_0 = r_0$
3. **for** $k = 0 \dots$
4. $v_k = Ap_k$
5. $\alpha_k = \|r_k\|^2 / \|v_k\|^2$
6. $x_{k+1} = x_k + \alpha_k p_k$
7. $q_{k+1} = q_k - \alpha_k v_k$
8. $r_{k+1} = A^T q_{k+1}$
9. $\beta_{k+1} = \|r_{k+1}\|^2 / \|r_k\|^2$
10. $p_{k+1} = r_{k+1} + \beta_k p_k$

Algorithm A.3 Conjugate gradient and conjugate gradient least squares algorithms. The algorithm is described in more detail in the text, but in brief, they choose descent directions p_k that are conjugate to each other with respect to C by computing a factor β by which to discount the previous search direction p_{k-1} . They then find the optimal step size α and take a downhill step by an amount $\alpha_k p_k$.

A.5.2 Preconditioning

As we mentioned previously, the rate of convergence of the conjugate gradient algorithm is governed in large part by the condition number $\kappa(C)$. Its effectiveness can therefore be increased dramatically by reducing this number, e.g., by rescaling elements in x , which corresponds to rescaling rows and columns in C .

In general, preconditioning is usually thought of as a change of basis from the vector x to a new vector

$$\hat{x} = Sx. \quad (\text{A.52})$$

The corresponding linear system being solved then becomes

$$AS^{-1}\hat{x} = S^{-1}b \quad \text{or} \quad \hat{A}\hat{x} = \hat{b}, \quad (\text{A.53})$$

with a corresponding least squares energy (A.29) of the form

$$E_{\text{PLS}} = \hat{x}^T (\mathbf{S}^{-T} \mathbf{C} \mathbf{S}^{-1}) \hat{x} - 2\hat{x}^T (\mathbf{S}^{-T} \mathbf{d}) + \|\hat{\mathbf{b}}\|^2. \quad (\text{A.54})$$

The actual preconditioned matrix $\hat{\mathbf{C}} = \mathbf{S}^{-T} \mathbf{C} \mathbf{S}^{-1}$ is usually not explicitly computed. Instead, Algorithm A.3 is extended to insert \mathbf{S}^{-T} and \mathbf{S}^T operations at the appropriate places (Björck 1996; Golub and Van Loan 1996; Trefethen and Bau 1997; Saad 2003; Nocedal and Wright 2006).

A good preconditioner \mathbf{S} is easy and cheap to compute, but is also a decent approximation to a square root of \mathbf{C} , so that $\kappa(\mathbf{S}^{-T} \mathbf{C} \mathbf{S}^{-1})$ is closer to 1. The simplest such choice is the square root of the diagonal matrix $\mathbf{S} = \mathbf{D}^{1/2}$, with $\mathbf{D} = \text{diag}(\mathbf{C})$. This has the advantage that any scalar change in variables (e.g., using radians instead of degrees for angular measurements) has no effect on the range of convergence of the iterative technique. For problems that are naturally block-structured, e.g., for structure from motion, where 3D point positions or 6D camera poses are being estimated, a block diagonal preconditioner is often a good choice.

A wide variety of more sophisticated preconditioners have been developed over the years (Björck 1996; Golub and Van Loan 1996; Trefethen and Bau 1997; Saad 2003; Nocedal and Wright 2006), many of which can be directly applied to problems in computer vision (Byröd and Åström 2009; Jeong, Nistér, Steedly *et al.* 2010; Agarwal, Snavely, Seitz *et al.* 2010). Some of these are based on an *incomplete Cholesky* factorization of \mathbf{C} , i.e., one in which the amount of fill-in in \mathbf{R} is strictly limited, e.g., to just the original non-zero elements in \mathbf{C} .¹⁰ Other preconditioners are based on a sparsified, e.g., tree-based or clustered, approximation to \mathbf{C} (Koutis 2007; Koutis and Miller 2008; Grady 2008; Koutis, Miller, and Tolliver 2009), since these are known to have efficient inversion properties.

For grid-based image-processing applications, *parallel* or *hierarchical* preconditioners often perform extremely well (Yserentant 1986; Szeliski 1990b; Pentland 1994; Saad 2003; Szeliski 2006b). These approaches use a change of basis transformation \mathbf{S} that resembles the pyramidal or wavelet representations discussed in Section 3.5, and are hence amenable to parallel and GPU-based implementations. Coarser elements in the new representation quickly converge to the low-frequency components in the solution, while finer-level elements encode the higher-frequency components. Some of the relationships between hierarchical preconditioners, incomplete Cholesky factorization, and multigrid techniques are explored by Saad (2003) and Szeliski (2006b).

¹⁰ If a complete Cholesky factorization $\mathbf{C} = \mathbf{R}^T \mathbf{R}$ is used, we get $\hat{\mathbf{C}} = \mathbf{R}^{-T} \mathbf{C} \mathbf{R}^{-1} = \mathbf{I}$ and all iterative algorithms converge in a single step, thereby obviating the need to use them, but the complete factorization is often too expensive. Note that incomplete factorization can also benefit from reordering.

A.5.3 Multigrid

One other class of iterative techniques widely used in computer vision is *multigrid* techniques (Briggs, Henson, and McCormick 2000; Trottenberg, Oosterlee, and Schuller 2000), which have been applied to problems such as surface interpolation (Terzopoulos 1986a), optical flow (Terzopoulos 1986a; Bruhn, Weickert, Kohlberger *et al.* 2006), high dynamic range tone mapping (Fattal, Lischinski, and Werman 2002), colorization (Levin, Lischinski, and Weiss 2004), natural image matting (Levin, Lischinski, and Weiss 2008), and segmentation (Grady 2008).

The main idea behind multigrid is to form coarser (lower-resolution) versions of the problems and use them to compute the low-frequency components of the solution. However, unlike simple coarse-to-fine techniques, which use the coarse solutions to initialize the fine solution, multigrid techniques only *correct* the low-frequency component of the current solution and use multiple rounds of coarsening and refinement (in what are often called “V” and “W” patterns of motion across the pyramid) to obtain rapid convergence.

On certain simple homogeneous problems (such as solving Poisson equations), multigrid techniques can achieve optimal performance, i.e., computation times linear in the number of variables. However, for more inhomogeneous problems or problems on irregular grids, variants on these techniques, such as *algebraic multigrid* (AMG) approaches, which look at the structure of C to derive coarse level problems, may be preferable. Saad (2003) has a nice discussion of the relationship between multigrid and parallel preconditioners and on the relative merits of using multigrid or conjugate gradient approaches.

Appendix B

Bayesian modeling and inference

B.1	Estimation theory	757
B.1.1	Likelihood for multivariate Gaussian noise	757
B.2	Maximum likelihood estimation and least squares	759
B.3	Robust statistics	760
B.4	Prior models and Bayesian inference	762
B.5	Markov random fields	763
B.5.1	Gradient descent and simulated annealing	765
B.5.2	Dynamic programming	766
B.5.3	Belief propagation	768
B.5.4	Graph cuts	770
B.5.5	Linear programming	773
B.6	Uncertainty estimation (error analysis)	775