

Basic Math for 16-720

August 23, 2002

1 Linear Algebra

1.1 Vectors and Matrices

First, a reminder of a few basic notations, definitions, and terminology:

- Unless indicated otherwise, vectors are always column vectors ($n \times 1$ matrices). The identity matrix is denoted by \mathbf{I} .
- The dot product of two column vectors is defined as $\mathbf{u} \cdot \mathbf{v} = \mathbf{u}^T \mathbf{v}$
- The norm of a vector is defined as $|\mathbf{u}|^2 = \mathbf{u}^T \mathbf{u}$
- If θ is the angle between \mathbf{u} and \mathbf{v} , then $\cos \theta = \frac{\mathbf{u}^T \mathbf{v}}{|\mathbf{u}| |\mathbf{v}|}$
- A matrix \mathbf{A} is symmetric if $\mathbf{A} = \mathbf{A}^T$
- The null space of a matrix \mathbf{A} is the set of vectors \mathbf{x} such that $\mathbf{A}\mathbf{x} = \mathbf{0}$
- The rank of \mathbf{A} is the dimension of its column(row) space
- A $n \times n$ matrix \mathbf{A} is non-singular (\mathbf{A}^{-1} exists) if its rank is n , i.e., $\det(\mathbf{A}) \neq 0$
- A $p \times p$ minor of \mathbf{A} is an $p \times p$ sub-matrix extracted from the rows and columns of \mathbf{A} . If \mathbf{A} is of rank $r < n$ then all of its $p \times p$ minors have zero determinant for $p > r$.
- The range of \mathbf{A} is the set of vector \mathbf{y} such that there exists a vector \mathbf{x} with $\mathbf{A}\mathbf{x} = \mathbf{y}$
- As an often-used notation, we denote by $[\mathbf{A}|\mathbf{B}]$ the $p \times (m + n)$ matrix obtain by concatenating the m columns of the $p \times m$ matrix \mathbf{A} and the n columns of the $p \times n$ matrix \mathbf{B} (with the obvious corresponding notation for concatenation of the rows.)

1.2 Matrix Decompositions

1.2.1 Eigenvectors and Eigenvalues

\mathbf{x} is an *eigenvector* of a $n \times n$ matrix \mathbf{A} if there exists a scalar λ such that: $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$; λ is the *eigenvalue* associated with the eigenvector \mathbf{x} ¹.

If \mathbf{A} is a symmetric matrix, then it has n real eigenvalues that are solution of the n th-degree polynomial equation $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$.

If \mathbf{A} is of rank r , then $n - r$ eigenvalues are equal to 0.

Note that, if \mathbf{x} is an eigenvector of \mathbf{A} , then so is $\alpha\mathbf{x}$ for any scalar α . Usually, when we talk about *THE* eigenvectors of the \mathbf{A} , we refer to the vector such that $\|\mathbf{x}\|^2 = 1$.

¹In MATLAB: `help eig`

Eigenvalues and eigenvectors are used for reducing a general symmetric matrix \mathbf{A} to a simpler form. Namely, if λ_i are the eigenvalues and \mathbf{v}_i are the corresponding eigenvectors, then \mathbf{A} can be decomposed (diagonalized) in the form:

$$\mathbf{A} = \mathbf{R}\mathbf{D}\mathbf{R}^T$$

where \mathbf{D} is a diagonal matrix whose i th diagonal elements is λ_i and \mathbf{R} is an orthonormal matrix ($\mathbf{R}^T\mathbf{R} = \mathbf{R}\mathbf{R}^T = \mathbf{I}$) whose i th column is \mathbf{v}_i .

Eigenvalues and eigenvectors have a simple (and very useful) geometric interpretation. Suppose that all the eigenvalues of \mathbf{A} are positive. In that case, the locus of points \mathbf{x} such that $\mathbf{x}^T\mathbf{A}\mathbf{x} = 1$ is an ellipsoid. The eigenvectors are the directions of the axes of the ellipsoid; the eigenvalues quantify the elongations in the directions of the axes. Specifically, if λ_i is the eigenvalue corresponding to the i th axis of the ellipsoid, the radius of the ellipsoid along this axis is $\frac{1}{\sqrt{\lambda_i}}$. Another interpretation of the decomposition is that \mathbf{R} is the rotation matrix that aligns the coordinate axis with the axes of the ellipsoid.

Two additional properties of the eigenvalues are often useful:

- The sum of the eigenvalues of \mathbf{A} is equal to the trace of \mathbf{A} , i.e., the sum of its diagonal elements
- The product of the eigenvalues is equal to the determinant of \mathbf{A}

The concept of eigenvalue can be generalized further:

- Given two symmetric $n \times n$ matrices \mathbf{A} and \mathbf{B} , \mathbf{x} is a *generalized* eigenvector associated with *generalized* eigenvalue λ if:

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{B}\mathbf{x}$$

The computation of the generalized eigenvalues is more difficult, of course, but is critical in some problems².

1.2.2 Rayleigh Quotient

It turns out that many problems in computer vision can be reduced to the mathematical problem statement:

- Given a symmetric matrix \mathbf{A} , find a vector \mathbf{x} such that
- $\mathbf{x}^T\mathbf{A}\mathbf{x}$ is maximum AND
- $\|\mathbf{x}\|^2 = 1$

This problem is a constrained maximization problem (constrained because of $\|\mathbf{x}\|^2 = 1$). This problem is strictly equivalent to the problem:

- Find \mathbf{x} such that $\frac{\mathbf{x}^T\mathbf{A}\mathbf{x}}{\mathbf{x}^T\mathbf{x}}$ is maximum.

(recall that $\|\mathbf{x}\|^2 = \mathbf{x}^T\mathbf{x}$.) The ratio above is sometimes called a *Rayleigh Quotient*.

The solution to this problem is given by the following theorem:

- $\frac{\mathbf{x}^T\mathbf{A}\mathbf{x}}{\mathbf{x}^T\mathbf{x}}$ reaches its absolute maximum when \mathbf{x} is an eigenvector of \mathbf{A} corresponding to the *largest* eigenvalue λ_{max} .

²In MATLAB, the same function `eig` solves the generalized problem.

Note that the same result holds if the ratio is minimized (this time the *smallest* eigenvalue is used, of course.)

This result can be generalized in a straightforward manner:

- Given two matrices \mathbf{A} and \mathbf{B} , $\frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{B} \mathbf{x}}$ reaches its absolute maximum when \mathbf{x} is a *generalized* eigenvector of \mathbf{A} corresponding to the largest *generalized* eigenvalue λ_{max} , that is: $\mathbf{A} \mathbf{x} = \lambda_{max} \mathbf{B} \mathbf{x}$

Note that this generalized problem is equivalent to maximizing $\mathbf{x}^T \mathbf{A} \mathbf{x}$ under the constraint $\mathbf{x}^T \mathbf{B} \mathbf{x} = 1$.

As a simple example of application of the Rayleigh Quotient theorem, consider the following (contrived) problem: Suppose that we have a set of vectors in two dimensions, \mathbf{u}_j with $j = 1 \dots m$. We wish to find the vector \mathbf{x} of magnitude 1 such that \mathbf{x} is as close as possible to the \mathbf{u}_j 's. The distance between \mathbf{x} and \mathbf{u}_j could be measured by the dot product $\mathbf{x}^T \mathbf{u}_j$. This dot product measures the angle between \mathbf{x} and \mathbf{u}_j and is maximum (equal to the magnitude of \mathbf{u}_j) if $\mathbf{x} = \frac{\mathbf{u}_j}{|\mathbf{u}_j|}$. Therefore, we want to find a unit vector \mathbf{x} that maximizes:

$$\sum_j (\mathbf{x}^T \mathbf{u}_j)^2 = \sum_j \mathbf{x}^T \mathbf{u}_j \mathbf{u}_j^T \mathbf{x} = \mathbf{x}^T \sum_j \mathbf{u}_j \mathbf{u}_j^T \mathbf{x}$$

This is exactly the Rayleigh Quotient problem, so the solution is the eigenvector \mathbf{x}_{max} corresponding to the largest eigenvalue λ_{max} of the matrix:

$$\mathbf{A} = \sum_j \mathbf{u}_j \mathbf{u}_j^T$$

Incidentally \mathbf{A} is the *scatter matrix* of the cloud of points formed by the vectors \mathbf{u}_j and it defines the best fit ellipsoid to that cloud of points. In particular, \mathbf{x}_{max} corresponds to the direction of largest variation of the cloud of points.

1.2.3 Singular Value Decomposition

The decomposition in eigenvalues (diagonalization) applies only to square matrices. A more general decomposition is the Singular Value Decomposition (SVD)³. A $m \times n$ matrix \mathbf{A} with $m \geq n$ can always be decomposed in the form:

$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$$

where:

- $\mathbf{\Sigma}$ is an $n \times n$ *diagonal* matrix. The elements $\sigma_1, \dots, \sigma_n$ of the diagonal are called the singular values of \mathbf{A} and $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$
- \mathbf{U} is a $m \times n$ matrix. The columns of \mathbf{U} are of unit norm and orthogonal to each other, that is: $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ (sometimes called *column-orthonormal*)
- \mathbf{V} is a $n \times n$ orthonormal matrix, that is: $\mathbf{V}^T \mathbf{V} = \mathbf{V} \mathbf{V}^T = \mathbf{I}$

The SVD is a powerful tool for analyzing matrices. The key properties of the SVD are:

- SVD and eigenvalues: If \mathbf{A} happens to be a square matrix $m = n$, then there is a simple relation between the eigenvalues α_i of $\mathbf{A}^T \mathbf{A}$ and the singular values σ_i : $\alpha_i = \sigma_i^2$. In particular, if \mathbf{A} is symmetric, then: $\lambda_i^2 = \sigma_i^2$ where λ_i is the eigenvalue of \mathbf{A} .
- Rank: Suppose that p singular values are non-0, i.e., $\sigma_1 \geq \dots \geq \sigma_p \geq 0$ and $\sigma_{p+1} = \dots = \sigma_{n-p} = 0$, then the rank of \mathbf{A} is p .

³In MATLAB: `help svd`

- Null Space: If \mathbf{A} is of rank p , then the $n - p$ lowest singular values are equal to 0 and the last $n - p$ columns of \mathbf{V} form an orthonormal basis of the null space of \mathbf{A} . This is extremely important in practice because it allows us to represent explicitly the set of vectors for which $\mathbf{A}\mathbf{x} = 0$, the null space.
- Range: This is the dual property to the null space property: the first p columns of \mathbf{U} form an orthonormal basis for the range of \mathbf{A} , that is, the set of vectors \mathbf{y} such that there exists a vector \mathbf{x} with $\mathbf{y} = \mathbf{A}\mathbf{x}$.

A key application of SVD is to solve the following problem:

- Given a $m \times n$ matrix \mathbf{A} with $m \geq n$.
- Given a number $p \leq n$
- Find the matrix \mathbf{A}_p of rank p that is closest to \mathbf{A} .

To solve this problem, we need to define what we mean by “closest”. We can define the distance between any two matrices \mathbf{M} and \mathbf{N} by: $\sum_{i,j} (m_{ij} - n_{ij})^2$, that is the sum of squared differences between the components of \mathbf{M} and \mathbf{N} . This defines a norm on the space of matrices, which is usually called the *Frobenius* norm. The Frobenius norm simply captures the intuitive notion that matrices are “close to each other” if they all their elements are close to each other.

Having defined the Frobenius norm, we can now state the key result (a result that we will use repeatedly in vision applications, by the way):

- The matrix \mathbf{A}_p of rank p that is the closest (in the Frobenius sense) to a given matrix \mathbf{A} with SVD decomposition $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ is given by:

$$\mathbf{A}_p = \mathbf{U}_p \mathbf{\Sigma}_p \mathbf{V}_p^T, \text{ where}$$

- $\mathbf{\Sigma}_p$ is a $p \times p$ diagonal matrix whose elements are the first p singular values of \mathbf{A} .
- \mathbf{U}_p is a $m \times p$ matrix formed by taking the first p columns of \mathbf{U} .
- \mathbf{V}_p is a $n \times p$ matrix formed by taking the first p columns of \mathbf{V} .

This looks complicated, but all it really says is that if you want to approximate a matrix by a matrix of rank p , all you need to do is to ignore the $n - p$ least important components of the matrix (corresponding to the $n - p$ smallest singular values), and that it is mathematically the optimal thing to do! See Figure 2 for a graphical illustration.

The reason why this is important is that, in many problems, we know the rank p that some matrix \mathbf{A} should have because of prior knowledge about the structure of the underlying problem. However, because of noise in the data, the real matrix that we observe \mathbf{A}' may have a higher rank. Using the theorem above, we can recover the “ideal” matrix of rank p that approximates the matrix observed from the actual data, \mathbf{A}' , in a simple way.

Let us take a look at a simple toy example. Suppose that we are working on a problem in which we want to exploit the fact that a matrix computed from the data is known to be of rank 2. For example, suppose that the ideal (unknown) value of the matrix is:

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 1 & 3 \\ 4 & 2 & 6 \end{bmatrix}$$

Because of noise in the data, e.g., the matrix may have been derived from pixel values in an image, the actual matrix that we observe from the data is perturbed and is instead:

$$\mathbf{A}' = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 1 & 3 \\ 4 & 2 & 6 + 10^{-4} \end{bmatrix}$$

\mathbf{A}' is now of full rank because of the noise perturbation, and therefore any algorithm that depends on the rank-2 assumption will fail. The SVD decomposition of \mathbf{A}' is:

$$\mathbf{A}' = \begin{bmatrix} -0.497 & -0.648 & -0.577 \\ -0.313 & 0.754 & -0.577 \\ -0.810 & 0.106 & 0.577 \end{bmatrix} \begin{bmatrix} 9.075 & 0 & 0 \\ 0 & 1.280 & 0 \\ 0 & 0 & 2.5 \times 10^{-5} \end{bmatrix} \begin{bmatrix} -0.497 & -0.648 & -0.577 \\ -0.313 & 0.754 & -0.577 \\ -0.810 & 0.106 & 0.577 \end{bmatrix}^T$$

The projection theorem above tells us that the closest rank-2 matrix to the actual data matrix \mathbf{A}' is:

$$\mathbf{A}'' = \begin{bmatrix} -0.497 & -0.648 \\ -0.313 & 0.754 \\ -0.810 & 0.106 \end{bmatrix} \begin{bmatrix} 9.075 & 0 \\ 0 & 1.280 \end{bmatrix} \begin{bmatrix} -0.497 & -0.648 \\ -0.313 & 0.754 \\ -0.810 & 0.106 \end{bmatrix}^T$$

What the theorem gives us is an easy way to find the matrix \mathbf{A}'' that best fits our model (rank 2) given the data, \mathbf{A}' .

1.3 2-D and 3-D Geometry

1.3.1 Rigid and Affine Transformations

Let's consider for a moment the cases of 2 and 3-dimensional spaces. The simplest geometric transformation in such spaces is a translation $\mathbf{y} = \mathbf{x} + \mathbf{t}$. The next type of transformation is a rotation $\mathbf{y} = \mathbf{R}\mathbf{x}$ where \mathbf{R} is an orthonormal matrix, $\mathbf{R}^T\mathbf{R} = \mathbf{R}\mathbf{R}^T = \mathbf{I}$.

In 3-D, \mathbf{R} is fully defined by three parameters, either three rotation angles about the three coordinate axes, or the direction of the rotation axis and the angle of rotation.

In 2-D, \mathbf{R} is fully parameterized by one rotation angle θ and has the general form:

$$\mathbf{R} = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

Rigid transformations are combinations of rotations of translations: $\mathbf{y} = \mathbf{R}\mathbf{x} + \mathbf{t}$. Rigid transformations preserve lengths and angles.

If we replace the rotation matrix \mathbf{R} by a general 2×2 (or 3×3) matrix \mathbf{A} , the resulting transformation $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{t}$ is called an *affine* transformation. Affine transformations do not preserve lengths or angles but they do preserve *ratios* of lengths and parallelism (parallel lines remain parallel after transformation.)

1.3.2 Homogeneous Coordinates

If $\mathbf{u} = [x \ y]^T$ is a 2-D vector, we could write the transformation $\mathbf{v} = \mathbf{A}\mathbf{u} + \mathbf{t}$ as:

$$\mathbf{v} = \mathbf{M} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}, \text{ where } \mathbf{M} = [\mathbf{A}|\mathbf{t}] \text{ is the } 2 \times 3 \text{ matrix obtained by concatenating the columns of } \mathbf{A} \text{ and } \mathbf{t}.$$

The interesting thing is that the transformation between \mathbf{v} and $[x \ y \ 1]^T$ is expressed by a single 2×3 matrix \mathbf{M} . In general, the vector $[x \ y]^T$ is said to be expressed in *homogeneous* coordinates as $[x \ y \ 1]^T$. More generally, any set of homogeneous coordinates $[a \ b \ c]^T$ such that $a/c = x$ and $b/c = y$ represent the same 2-D vector $[x \ y]^T$. Two vectors expressed in homogeneous coordinates, $\tilde{\mathbf{u}} = [a \ b \ c]^T$ and $\tilde{\mathbf{u}}' = [a' \ b' \ c']^T$, represent the same vector if $a/c = a'/c'$ and $b/c = b'/c'$. In that case we write $\tilde{\mathbf{u}} \equiv \tilde{\mathbf{u}}'$ (or $\tilde{\mathbf{u}} \propto \tilde{\mathbf{u}}'$, depending on the text!) to indicate that $\tilde{\mathbf{u}}$ and $\tilde{\mathbf{u}}'$ are equivalent sets of homogeneous coordinates.

Be careful: $\tilde{\mathbf{u}} \equiv \tilde{\mathbf{u}}'$ does not mean that the components of $\tilde{\mathbf{u}}$ and $\tilde{\mathbf{u}}'$ are equal, only that they are proportional.

With these definitions, a generic affine transformation $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{t}$ is expressed in homogeneous coordinates as a transformation of the form $\tilde{\mathbf{u}}' = \mathbf{M}\tilde{\mathbf{u}}$, where \mathbf{M} is the 3×3 matrix defined as:

$$\mathbf{M} = \begin{bmatrix} \mathbf{A} & \mathbf{t} \\ 0 & 0 & 1 \end{bmatrix}$$

For example, the rigid transformation of rotation angle θ and translation $[2 \ 3]^T$ is represented in homogeneous coordinates by:

$$\mathbf{M} = \begin{bmatrix} \cos \theta & -\sin \theta & 2 \\ \sin \theta & \cos \theta & 3 \\ 0 & 0 & 1 \end{bmatrix}$$

Homogeneous coordinates are defined in the same way in 3-D by adding one coordinate to the three-dimensional vectors. That is, two homogeneous vectors $\tilde{\mathbf{u}} = [a \ b \ c \ d]^T$ and $\tilde{\mathbf{u}}' = [a' \ b' \ c' \ d']^T$ represent the same 3-D vector if: $a/c = a'/d' = x$, $b/d = b'/d' = y$, and $c/d = c'/d' = z$. A generic affine transformation $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{t}$ is expressed in homogeneous coordinates as a transformation of the form $\tilde{\mathbf{u}}' = \mathbf{M}\tilde{\mathbf{u}}$, where \mathbf{M} is the 4×4 matrix defined as:

$$\mathbf{M} = \begin{bmatrix} \mathbf{A} & \mathbf{t} \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

1.3.3 Projective Transformations

Consider first the 2-D case. The last row of the transformation matrices defined above in homogeneous coordinates is $[0 \ 0 \ 1]^T$. We could replace those fixed values by arbitrary values to define a general transformation in homogeneous coordinates: $\tilde{\mathbf{u}}' = \mathbf{M}\tilde{\mathbf{u}}$, where \mathbf{M} is a general 3×3 matrix. Consider a vector $\mathbf{x} = [x \ y]^T$, a representation in homogeneous coordinates is, for example, $\mathbf{u} = [x \ y \ 1]^T$ and

$$\tilde{\mathbf{u}}' = \mathbf{M}\tilde{\mathbf{u}} = \begin{bmatrix} m_1^T \tilde{\mathbf{u}} \\ m_2^T \tilde{\mathbf{u}} \\ m_3^T \tilde{\mathbf{u}} \end{bmatrix}, \text{ where } m_i^T \text{ is the } i\text{th row of } \mathbf{M}.$$

The vector represented by the transformed set of homogeneous coordinates $\tilde{\mathbf{u}}'$ is obtained by dividing the first coordinates by the third one, yielding:

$$x' = \frac{m_1^T \tilde{\mathbf{u}}}{m_3^T \tilde{\mathbf{u}}} \text{ and } y' = \frac{m_2^T \tilde{\mathbf{u}}}{m_3^T \tilde{\mathbf{u}}}$$

Note that, if we use arbitrary numbers instead of $[0 \ 0 \ 1]^T$ in the last row of the matrix, the transformation becomes a non-linear transformation of the original coordinates, even though it is a *linear* transformation in homogeneous coordinates. Such a transformation, represented by a 3×3 matrix in homogeneous coordinates is called a *projective* transformation⁴. Projective transformations do not preserve lengths, angles, ratios, or parallelism. They do preserve colinearity (things that are on the same line remain on the same line after transformation) and incidence (curves that are tangent to each other remain tangent to each other.)

Projective transformations are defined in three dimensions in the same manner: A 4×4 matrix is applied in homogeneous coordinates. The vector coordinates corresponding to the homogeneous coordinates are recovered by dividing the first three coordinates by the fourth one. The special case in which the last row of the matrix is $[0 \ 0 \ 0 \ 1]^T$ corresponds to affine transformations just like in 2-D. It is important to note that a transformation expressed in homogeneous coordinates is defined only up to a scale factor because two homogeneous vectors are equivalent if they are equal up to a scale factor ($\mathbf{u} \equiv \tilde{\mathbf{u}}'$ means that $\tilde{\mathbf{u}} = \alpha \tilde{\mathbf{u}}'$ for some scalar α .) In other words, if \mathbf{M} is the matrix of a transformation in homogeneous coordinates, then $\alpha \mathbf{M}$ represents the same transformation for any $\alpha \neq 0$.

⁴If you want to sound sophisticated, the correct technical name is “collineation”, but projective transformation is good enough for now....

1.3.4 3-D to 2-D Transformations

Transformations that map a point in 3-D to a point in the 2-D plane are expressed in homogeneous coordinates by 3×4 matrices. A 3-D to 2-D transformation maps a 3-D point of coordinates $[x \ y \ z]^T$ to a point in the plane of coordinates $[x' \ y']^T$ by applying a 3×4 matrix \mathbf{M} to the homogeneous coordinates of the 3-D point, $[x \ y \ z \ 1]^T$, i.e., the coordinates x' and y' are computed by:

$$\begin{bmatrix} u \\ v \\ w \end{bmatrix} = \mathbf{M} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} \quad x' = \frac{u}{w} \quad y' = \frac{v}{w}$$

3-D to 2-D transformations come in three varieties. If the last row of \mathbf{M} is $[0 \ 0 \ 0 \ 1]^T$ and the first two rows are the first two rows of a 3-D rigid transformation matrix, the transformation is an *orthographic* projection. This type of transformation amounts to applying a rigid transformation in 3-D and keeping the first two coordinates of the transformed 3-D point as the result of the projection.

If the last row of \mathbf{M} is $[0 \ 0 \ 0 \ 1]^T$ and the first two rows are arbitrary, i.e., not necessarily from a rigid transformation, then \mathbf{M} is an *affine* projection. This type of transformation amounts to applying an affine transformation in 3-D and keeping the first two coordinates of the transformed 3-D point as the result of the projection.

If the last row of \mathbf{M} is not $[0 \ 0 \ 0 \ 1]^T$, \mathbf{M} is a projective projection. Projective transformations from 3-D to 2-D are simply arbitrary 3×4 matrices applied to homogeneous coordinates. One important example that illustrates the use of projective projection is the representation of the perspective projection used in the pinhole camera model. Recall that, if $[x \ y \ z]^T$ is a point in space, the coordinates of the corresponding point $[x' \ y']^T$ after projection through a pinhole are given by:

$$x' = f \frac{x}{z} \quad \text{and} \quad y' = f \frac{y}{z}$$

In homogeneous coordinates, this becomes a linear transformation:

$$\begin{bmatrix} a \\ b \\ c \end{bmatrix} = \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}, \quad \text{with } x' = \frac{a}{c} \quad \text{and} \quad y' = \frac{b}{c}$$

This is the reason why the projective (homogeneous representation) is so important to us: It enables us to convert the non-linear, perspective projection problems to linear problems⁵

1.3.5 Cross-Product and Skew-Symmetric Matrices

In three dimensions, the cross product of two vectors $\mathbf{u} = [u_x \ u_y \ u_z]^T$ and $\mathbf{v} = [v_x \ v_y \ v_z]^T$ is defined by:

$$\mathbf{u} \times \mathbf{v} = \begin{bmatrix} u_y v_z - u_z v_y \\ u_z v_x - u_x v_z \\ u_x v_y - u_y v_x \end{bmatrix}$$

A few basic properties of the cross product are useful and worth remembering:

- $\mathbf{u} \times \mathbf{v} = -\mathbf{v} \times \mathbf{u}$ (antisymmetry)
- $\mathbf{w} \cdot (\mathbf{u} \times \mathbf{v})$ is the determinant of the matrix formed by the three vectors $\mathbf{u} \ \mathbf{v} \ \mathbf{w}$

⁵Of course, there is a much broader theory of projective representation behind this, which we will need later, but this is sufficient for now.

| Transformation | Vector Coordinates | Homogeneous Coordinates | Degrees of Freedom | Invariants |
|----------------|--|--|--------------------|--------------------------------|
| Translation | $\mathbf{y} = \mathbf{x} + \mathbf{t}$ | $\begin{bmatrix} \mathbf{I} & \mathbf{t} \\ 0 & 1 \end{bmatrix}$ | 2 | lengths, angles |
| Rotation | $\mathbf{y} = \mathbf{R}\mathbf{x}$ $\mathbf{R}^T \mathbf{R} = \mathbf{R}\mathbf{R}^T = \mathbf{I}$ | $\begin{bmatrix} \mathbf{R} & \mathbf{0} \\ 0 & 1 \end{bmatrix}$ | 1 | lengths, angles |
| Rigid | $\mathbf{y} = \mathbf{R}\mathbf{x} + \mathbf{t}$ | $\begin{bmatrix} \mathbf{R} & \mathbf{t} \\ 0 & 1 \end{bmatrix}$ | 3 | lengths, angles |
| Affine | $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{t}$ | $\begin{bmatrix} \mathbf{A} & \mathbf{t} \\ 0 & 1 \end{bmatrix}$ | 6 | ratios of lengths, parallelism |
| Projective | | 3×3 matrix \mathbf{M} | 8 | colinearity, incidence |

Table 1: Classification of 2-D transformations

| Transformation | Vector Coordinates | Homogeneous Coordinates | Degrees of Freedom | Invariants |
|----------------|--|--|--------------------|--------------------------------|
| Translation | $\mathbf{y} = \mathbf{x} + \mathbf{t}$ | $\begin{bmatrix} \mathbf{I} & \mathbf{t} \\ 0 & 0 & 1 \end{bmatrix}$ | 3 | lengths, angles |
| Rotation | $\mathbf{y} = \mathbf{R}\mathbf{x}$ $\mathbf{R}^T \mathbf{R} = \mathbf{R}\mathbf{R}^T = \mathbf{I}$ | $\begin{bmatrix} \mathbf{R} & \mathbf{0} \\ 0 & 0 & 1 \end{bmatrix}$ | 3 | lengths, angles |
| Rigid | $\mathbf{y} = \mathbf{R}\mathbf{x} + \mathbf{t}$ | $\begin{bmatrix} \mathbf{R} & \mathbf{t} \\ 0 & 0 & 1 \end{bmatrix}$ | 6 | lengths, angles |
| Affine | $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{t}$ | $\begin{bmatrix} \mathbf{A} & \mathbf{t} \\ 0 & 0 & 1 \end{bmatrix}$ | 12 | ratios of lengths, parallelism |
| Projective | | 4×4 matrix \mathbf{M} | 15 | colinearity, incidence |

Table 2: Classification of 3-D transformations

| Transformation | Homogeneous Coordinates | Degrees of Freedom |
|----------------|--|--------------------|
| Orthographic | $\begin{bmatrix} \mathbf{R} & \mathbf{t} \\ 0 & 0 & 0 & 1 \end{bmatrix}$ $\mathbf{R}\mathbf{R}^T = \mathbf{I}$ (row-orthonormal) | 5 |
| Affine | $\begin{bmatrix} \mathbf{A} & \mathbf{t} \\ 0 & 0 & 0 & 1 \end{bmatrix}$ | 8 |
| Projective | 3×4 matrix \mathbf{M} | 11 |

Table 3: Classification of 3-D to 2-D transformations

- $\mathbf{a} \times (\mathbf{b} \times \mathbf{c}) = (\mathbf{a} \cdot \mathbf{b})\mathbf{c} - (\mathbf{a} \cdot \mathbf{c})\mathbf{b}$
- $|\mathbf{u} \times \mathbf{v}| = \sin \theta |\mathbf{u}| |\mathbf{v}|$

Given a vector \mathbf{u} , The mapping that associates $\mathbf{u} \times \mathbf{x}$ to \mathbf{x} is obviously a linear mapping. Therefore, there must exist a matrix \mathbf{M} such that: $\mathbf{u} \times \mathbf{x} = \mathbf{M}\mathbf{x}$ for any \mathbf{x} . Such a matrix \mathbf{M} is called a *skew-symmetric matrix*. The skew-symmetric matrix associated with a vector \mathbf{u} is denoted typically by \mathbf{u}_\times or $[\mathbf{u}]_\times$ (sorry, no standardized notations!) The expression for the skew-symmetric matrix is:

$$\mathbf{u}_\times = \begin{bmatrix} 0 & -u_z & u_y \\ u_z & 0 & -u_x \\ -u_y & u_x & 0 \end{bmatrix}$$

So, for any \mathbf{x} , $\mathbf{u} \times \mathbf{x} = \mathbf{u}_\times \mathbf{x}$. The matrix is called skew-symmetric because $\mathbf{u}_\times^T = -\mathbf{u}_\times$ owing to the antisymmetry of the cross product.

Warning: If you have never seen this skew-symmetric construct before, you should make sure you understand it. Most of the derivations in multi-camera geometry will rely on this notation.

A final note in this discussion of the cross product has to do with its use in expressing parallelism between vectors. Specifically, let \mathbf{a} and \mathbf{b} be two parallel vectors. The fact that they are parallel can be expressed by saying that they are colinear, i.e., there exists a scalar λ such that: $\mathbf{a} = \lambda\mathbf{b}$. This is fine, except that we have introduced an extra variable λ which really is a dummy variable which adds absolutely no information if all that we are interested in is the parallelism of the vectors. This seemingly trivial observation will become critical in some of the problems we'll have to solve. Specifically, in solving any problem that involves parallelism, we would be forced to introduce spurious extra degrees of freedom, the λ 's, which would greatly complicate the algorithms.

A better way to describe parallelism is to observe that \mathbf{a} and \mathbf{b} are parallel iff $\mathbf{a} \times \mathbf{b} = \mathbf{0}$, or even better: $\mathbf{a}_\times \mathbf{b} = \mathbf{0}$, that way we have converted the parallelism condition to a simple set of linear equations without any extra parameters. In the future, we will automatically convert problems like "those rays are parallel to this direction" to a set of linear equations using this trick. In particular, given two homogeneous vectors in the plane $\tilde{\mathbf{u}} = [a \ b \ c]^T$ and $\tilde{\mathbf{u}}' = [a' \ b' \ c']^T$, we have seen that $\tilde{\mathbf{u}} \equiv \tilde{\mathbf{u}}'$ means that the two vectors are proportional to each other, in other words: $\tilde{\mathbf{u}} \times \tilde{\mathbf{u}}' = \mathbf{0}$. A little trick that will help a lot.

1.3.6 Lines and Planes

Here are some very simple (but sometimes forgotten) facts about lines and planes:

- A line in the plane (resp. plane in space) has equation $ax + by + c = 0$ (resp. $ax + by + cz + d = 0$)

- In homogeneous coordinates, the equation becomes $\mathbf{l}^T \mathbf{v}$ (resp. $\mathbf{p}^T \mathbf{v}$) with $\mathbf{v} = [x \ y \ 1]^T$ and $\mathbf{l} = [a \ b \ c]^T$ (resp. $\mathbf{v} = [x \ y \ z \ 1]^T$ and $\mathbf{p} = [a \ b \ c \ d]^T$.)
- The normal direction to the line (resp. plane) is the vector $\frac{1}{\sqrt{a^2+b^2}}[a \ b]^T$ (resp. $\frac{1}{\sqrt{a^2+b^2+c^2}}[a \ b \ c]^T$)
- The distance of a generic point \mathbf{v} to the line \mathbf{l} (resp. plane \mathbf{p}) is $\frac{1}{\sqrt{a^2+b^2}}|\mathbf{l}^T \mathbf{v}|$ (resp. $\frac{1}{\sqrt{a^2+b^2+c^2}}|\mathbf{p}^T \mathbf{v}|$.)
- The intersection of two lines \mathbf{l}_1 and \mathbf{l}_2 is the point of homogeneous coordinates $\mathbf{l}_1 \times \mathbf{l}_2$ ⁶

2 Optimization

2.1 Linear Least-Squares

Given a m -dimensional vector \mathbf{b} and a $m \times n$ matrix \mathbf{A} , the problem is to find the n -dimensional vector \mathbf{x} such that the residual $r = \|\mathbf{Ax} - \mathbf{b}\|^2$ is minimum⁷. This is the standard least-squares problem which we need to solve whenever we have m n -dimensional input data vectors \mathbf{a}_i , output values b_i and we want to fit a linear model of the form $\mathbf{a}_i^T \mathbf{x} = b_i$ by minimizing $\sum_1^m |\mathbf{a}_i^T \mathbf{x} - b_i|^2$. The matrix \mathbf{A} is obtained by stacking the m row vectors \mathbf{a}_i^T and \mathbf{b} is obtained by concatenating the b_i 's into a column vector.

The solution to the least-squares problem is given by:

$$\mathbf{x} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$$

The matrix $\mathbf{A}^+ = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ is the *pseudo-inverse* of \mathbf{A} . This solution assumes that the problem is over-constrained, i.e., there are more independent equations than unknowns ($\text{rank}(\mathbf{A}) = n$ and $m > n$). If the problem is underconstrained, the pseudo-inverse cannot be computed.

A solution that works with both over- and under-constrained problems can be derived using the SVD. This solution can also be more efficient and is more robust numerically.

Let $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ be the SVD of \mathbf{A} . Define the matrix $\mathbf{A}^+ = \mathbf{V}\mathbf{\Sigma}^+\mathbf{U}^T$, where $\mathbf{\Sigma}^+$ is the matrix formed by taking the inverse of the non-zero elements of $\mathbf{\Sigma}$ (and leaving the 0 elements unchanged.) With this definition of \mathbf{A}^+ , the solution to the least-squares problem is:

$$\mathbf{x}_o = \mathbf{A}^+ \mathbf{b}.$$

If the system is overconstrained, this is the same solution as with the previous definition of the pseudo-inverse.

If the system is under-constrained, and therefore has multiple solutions, this gives the solution \mathbf{x}_o of minimum norm. Note that, any other solution will be of the form $\mathbf{x}_o + \mathbf{V}'\mathbf{y}$, where \mathbf{V}' is the matrix formed from the columns corresponding to the null singular values of \mathbf{A} (i.e., the columns of \mathbf{V}' span the null space of \mathbf{A}) and \mathbf{y} is an arbitrary vector.

2.2 Unconstrained Non-Linear Problems

The problem now is to find the n -vector \mathbf{x} that is an extremum (minimum or maximum) of an objective function $z = F(\mathbf{x})$. We will encounter many such problems and we can generally use a standard minimization ⁸ as a black box. However, it is still interesting to understand the general characteristics of minimization approaches in order to assess the feasibility of a particular solution. For example, the dimensionality of some problems is such that, without some manipulation of the data or some approximation to the problem, optimization of the objective function will not be feasible. To help understand

⁶A similar relation exists for the intersection of planes but not with the cross product, which works only for 3-D vectors. One needs to use a different representation of lines in 3D, the Plucker coordinates, see Forsyth and Ponce, p. 165.

⁷MATLAB: `help regress`

⁸`help optim` in MATLAB will list pretty much all the non-linear optimization techniques known to mankind.

key issues and limitations, this section, along with the next two sections, provides a very quick overview of the basic approaches commonly used to solve non-linear optimization problems.

An optimization algorithm starts from an initial choice of the argument, \mathbf{x}_o and iteratively generates successive values $\mathbf{x}_1, \dots, \mathbf{x}_n$ corresponding to values of the objective function $F(\mathbf{x}_1), \dots, F(\mathbf{x}_n)$ that converge to an extremum. It should be clear from the outset that the best we can hope for is to find a *local* extremum of the objective function. It is not possible in general to guarantee that the algorithm converges to a globally optimal value. It is also generally not possible to guarantee that the local extremum is a “deep” minimum or maximum. Moreover, many techniques may fail to converge if the starting point \mathbf{x}_o is too far from a local extremum. It is therefore critical to understand how good (or poor) the choice of \mathbf{x}_o is.

Different algorithms may have dramatically different convergence rates, as measured by the rate of change of the objective function as a function of the number of iterations. Typically, the fastest algorithms have quadratic convergence rate. Unfortunately, it is often the case that the computations involved at each step are prohibitively expensive. Therefore, assessing the feasibility of solving a problem through non-linear optimization involves understanding how good is the choice of the starting point, the convergence properties of the optimization algorithms, and the amount of computation to be performed at each iteration - and whether the computation can be carried out at all in practice.

To compute \mathbf{x}_{k+1} from \mathbf{x}_k , we need to first consider the Taylor expansion of F around \mathbf{x}_k :

$$F(\mathbf{x}_{k+1}) = F(\mathbf{x}_k) + \nabla F^T(\mathbf{x}_{k+1} - \mathbf{x}_k) + \frac{1}{2}(\mathbf{x}_{k+1} - \mathbf{x}_k)^T \mathbf{H}(\mathbf{x}_{k+1} - \mathbf{x}_k)$$

∇F is the gradient of F at \mathbf{x}_k , and \mathbf{H} is the Hessian, i.e., matrix of second derivatives, of F at \mathbf{x}_k : $\mathbf{H} = \nabla^2 F(\mathbf{x}_k)$. Optimization algorithms attempt to find the update vector $(\mathbf{x}_{k+1} - \mathbf{x}_k)$ that would yield the largest change in F based on this Taylor expansion.

We assume from now on that we want to find a minimum of F (a few sign reversals are needed if we want a maximum of F .) The first approach is based on the observation that the fastest variation of F is in the direction opposite to that of the gradient (*steepest descent*), i.e., $-\nabla F$. In other words, the iteration is defined by: $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \nabla F$, where α is chosen to maximize the change in F . Substituting $\alpha \nabla F$ for $(\mathbf{x}_{k+1} - \mathbf{x}_k)$ in the Taylor expansion yields:

$$F(\mathbf{x}_{k+1}) = F(\mathbf{x}_k) - \alpha |\nabla F|^2 + \frac{1}{2} \alpha^2 \nabla F^T \mathbf{H} \nabla F$$

From this expression it is easy to see that the α that minimizes the right-hand side is $\alpha = \frac{|\nabla F|^2}{\nabla F^T \mathbf{H} \nabla F}$, and the update rule at each iteration is:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{|\nabla F|^2}{\nabla F^T \mathbf{H} \nabla F} \nabla F$$

This approach is the *gradient descent* algorithm. Gradient descent has the advantage that it does not require any expensive matrix inversion. As shown here, it requires the computation of the second derivatives of F . However, it should be noted that other techniques can be used to choose α that do not require the second derivatives (with potentially slower convergence.) Gradient descent has two drawbacks: It has slow (linear) convergence and it is not guaranteed to converge if the starting point is far from the minimum.

An alternative to gradient descent is to not try to force the iterations to follow the gradient and to, instead, find the difference $(\mathbf{x}_{k+1} - \mathbf{x}_k)$ that minimizes the right-hand side of the Taylor expansion. Setting the derivative of the right-hand side with respect to $(\mathbf{x}_{k+1} - \mathbf{x}_k)$ to zero yields directly:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{H}^{-1} \nabla F$$

This approach is the *Newton* method. The advantage of the Newton method is that it has fast (quadratic) convergence. It still requires the computation of the second derivatives. More importantly, a major drawback is that it requires the inversion of the Hessian matrix which may cause serious computational and numerical problems. From a computational standpoint, if d is the dimension of \mathbf{x} , then \mathbf{H} is a

$d \times d$ matrix and its inversion is $O(d^3)$. In real problems (in vision), d could be on the order of several hundreds or thousands, thus requiring a large amount of computation at each iteration, irrespective of the convergence rate.

Another problem is that \mathbf{H} may be close to singular, in which case the inversion is numerically unstable especially if d is large. To see why this is a problem in practice, it is useful to understand the geometric interpretation of the Newton method. The right-hand side of the Taylor expansion of F is an approximation of the surface $z = F(\mathbf{x})$ by a quadric surface. The Newton method simply finds the lowest on that approximating surface. Assuming that we have translated the axes so that the minimum is at $\mathbf{x} = \mathbf{0}$, the Taylor expansion around $\mathbf{0}$ becomes $F(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{H} \mathbf{x}$, which is a quadratic approximation of the objective function.

The eigenvalues of \mathbf{H} give us an indication of the local shape of $z = F(\mathbf{x})$ (Figure 3) Specifically, large eigenvalues indicate directions along which the function varies quickly. Along these directions, the minimum is sharp and can be localized easily. Furthermore, the inversion of \mathbf{H} is numerically stable. On the other hand, small eigenvalues correspond to directions along which the function varies slowly, i.e., the minimum of the surface is very shallow. Along these directions, it is more difficult to accurately localize the minimum and the inversion of \mathbf{H} becomes numerically unstable.

In the extreme, an eigenvalue of 0 corresponds to a direction in which the function does not vary at all, i.e., the surface is flat in that direction. In that case, of course, \mathbf{H} becomes singular and the iteration rule above can no longer be used⁹. This discussion illustrates why the performance of the Newton method may vary widely depending on the shape of the objective function. Since the computation of the Hessian and its inversion may be problematic, other approaches (*quasi-Newton*) use various approximations of the Hessian. We review one class of techniques that is particularly important for our applications in the next section.

2.3 Non-Linear Least Squares

A particularly important case is the case in which F has the special form:

$$F(\mathbf{x}) = \sum_{j=1}^m r_j^2(\mathbf{x})$$

The functions r_j are usually termed *residuals*. This special form of the objective function is very common. For example, suppose that we observe the projections $\mathbf{p}_j = [x'_j \ y'_j]^T$ in m images of an unknown point $\mathbf{P} = [x \ y \ z]^T$. To recover the \mathbf{P} from its image projections, we need to find \mathbf{P} that minimizes the distance between \mathbf{p}_j and $g_j(\mathbf{P})$, where g_j is the function that maps a point in the scene to its projection in image j (a non-linear function that is a general form of the perspective projection: $\mathbf{p}_j = f_j[\frac{x}{z} \ \frac{y}{z}]^T$.) In other words, we need to minimize the objective function:

$$F(\mathbf{P}) = \sum_{j=1}^m |\mathbf{p}_j - g_j(\mathbf{P})|^2$$

This formulation will be central to solving several problems in structure recovery from multiple images. In these problems, there may be hundreds of variables, i.e., hundreds of points \mathbf{P} , and direct computation of the Hessian may become problematic. Instead, we can use an approximation of the Hessian as shown below.

Let \mathbf{J} denote the matrix whose rows are the gradients of the residuals, i.e., \mathbf{J} is obtained by stacking the row vectors $\nabla r_j(\mathbf{x})^T$, and let $\mathbf{r}(\mathbf{x})$ denote the column vector $[r_1(\mathbf{x}), \dots, r_m(\mathbf{x})]^T$. We can write the gradient of the objective function as:

$$\nabla F(\mathbf{x}) = \sum_{j=1}^m \nabla r_j(\mathbf{x}) r_j(\mathbf{x}) = \sum_{j=1}^m \mathbf{J}^T r_j(\mathbf{x}) = \mathbf{J}^T \mathbf{r}(\mathbf{x})$$

Differentiating this expression one more time, we get the Hessian of F :

⁹A third case occurs when the eigenvalues are negative, in which case the point is a saddle point at which the gradient vanishes but is not an extremum of the function.

$$\mathbf{H} = \nabla^2 F(\mathbf{x}) = \mathbf{J}^T \mathbf{J} + \sum_{j=1}^m \nabla^2 r_j(\mathbf{x}) r_j(\mathbf{x})$$

The key observation is that only the second term in this expression of $\nabla^2 F(\mathbf{x})$ requires the knowledge of the second derivatives. Therefore, a natural approach is to drop this hard to compute second term and to approximate the Hessian by the first term only:

$$\mathbf{H} \approx \mathbf{J}^T \mathbf{J}$$

In particular, replacing this approximation of the Hessian in the update rule used in the Newton method, $\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{H}^{-1} \nabla F$, we obtain:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - (\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \mathbf{r}$$

This approach is the *Gauss-Newton* method, which has many advantages. Gauss-Newton does not require the computation of any second derivatives and, in many problems the first term used in the approximation of \mathbf{H} is indeed much larger than the second term. It appears that we still need a matrix inversion (and a matrix multiplication to evaluate $\mathbf{J}^T \mathbf{J}$.) In fact, $(\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T$ is the pseudo-inverse of \mathbf{J} , as defined in Section 2.1, and can be computed directly from \mathbf{J} using more efficient SVD techniques as described in 2.1 or other linear least squares techniques not listed here. The bottom line is that one can avoid completely any explicit computation of $\mathbf{J}^T \mathbf{J}$ or of its inverse!

Gauss-Newton provides a computationally efficient approach in the case where the objective function is a sum of squared residuals. The fact still remains that Gauss-Newton is based on an approximation of the Hessian. Presumably, if this approximation could be refined, we ought to be able to achieve faster convergence. A general approach to refining the Hessian approximation is the *Levenberg-Marquart* algorithm in which the approximation $\mathbf{H} \approx \mathbf{J}^T \mathbf{J}$ is replaced by:

$$\mathbf{H} \approx \mathbf{J}^T \mathbf{J} + \lambda \mathbf{I}$$

In this approximation, λ is a scalar that is adjusted at each iteration based on the rate of decrease of F at prior iterations. The heuristics used for controlling λ are rather involved, but the general idea of the algorithm is that the order of magnitude of the discrepancy between the true Hessian and its approximation is evaluated at each iteration and the Hessian approximation is corrected by adding a diagonal matrix of that magnitude. Levenberg-Marquart combines the good convergence properties of the original Newton method and it avoids completely the computation of the Hessian while maintaining a better estimate than Gauss-Newton.

2.4 Constrained Problems

In their simplest forms, constrained problems involve finding an extremum of an objective function $z = F(\mathbf{x})$ under a constraint $G(\mathbf{x}) = 0$. The problem is solved by optimizing a combined function, the *Lagrangian*:

$$L(\mathbf{x}) = F(\mathbf{x}) - \lambda G(\mathbf{x})$$

λ is an additional variable, the *Lagrange multiplier*. The gradient of L vanishes at the solution, yielding the general relation' between the gradients of the objective function and of the constraint:

$$\nabla F(\mathbf{x}) - \lambda \nabla G(\mathbf{x}) = 0$$

For example, if the objective function is $F(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$, where \mathbf{A} is a $n \times n$ matrix, and the constraint is $\mathbf{x}^T \mathbf{x} = 1$, the relation above between the gradients becomes: $\mathbf{A} \mathbf{x} - \lambda \mathbf{x} = 0$. Therefore, the solution to this constrained problem is an eigenvector of \mathbf{A} since $\mathbf{A} \mathbf{x} = \lambda \mathbf{x}$. This is basically the Rayleigh Quotient theorem from 1.2.2.