

## Appendix C. Properties of Matrices

In this appendix, we gather together some useful properties and identities involving matrices and determinants. This is not intended to be an introductory tutorial, and it is assumed that the reader is already familiar with basic linear algebra. For some results, we indicate how to prove them, whereas in more complex cases we leave the interested reader to refer to standard textbooks on the subject. In all cases, we assume that inverses exist and that matrix dimensions are such that the formulae are correctly defined. A comprehensive discussion of linear algebra can be found in Golub and Van Loan (1996), and an extensive collection of matrix properties is given by Lütkepohl (1996). Matrix derivatives are discussed in Magnus and Neudecker (1999).

---

### Basic Matrix Identities

A matrix  $\mathbf{A}$  has elements  $A_{ij}$  where  $i$  indexes the rows, and  $j$  indexes the columns. We use  $\mathbf{I}_N$  to denote the  $N \times N$  identity matrix (also called the unit matrix), and where there is no ambiguity over dimensionality we simply use  $\mathbf{I}$ . The transpose matrix  $\mathbf{A}^T$  has elements  $(\mathbf{A}^T)_{ij} = A_{ji}$ . From the definition of transpose, we have

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T \quad (\text{C.1})$$

which can be verified by writing out the indices. The inverse of  $\mathbf{A}$ , denoted  $\mathbf{A}^{-1}$ , satisfies

$$\mathbf{AA}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}. \quad (\text{C.2})$$

Because  $\mathbf{ABB}^{-1}\mathbf{A}^{-1} = \mathbf{I}$ , we have

$$(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}. \quad (\text{C.3})$$

Also we have

$$(\mathbf{A}^T)^{-1} = (\mathbf{A}^{-1})^T \quad (\text{C.4})$$

which is easily proven by taking the transpose of (C.2) and applying (C.1).

A useful identity involving matrix inverses is the following

$$(\mathbf{P}^{-1} + \mathbf{B}^T \mathbf{R}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{R}^{-1} = \mathbf{P} \mathbf{B}^T (\mathbf{B} \mathbf{P} \mathbf{B}^T + \mathbf{R})^{-1}. \quad (\text{C.5})$$

which is easily verified by right multiplying both sides by  $(\mathbf{B} \mathbf{P} \mathbf{B}^T + \mathbf{R})$ . Suppose that  $\mathbf{P}$  has dimensionality  $N \times N$  while  $\mathbf{R}$  has dimensionality  $M \times M$ , so that  $\mathbf{B}$  is  $M \times N$ . Then if  $M \ll N$ , it will be much cheaper to evaluate the right-hand side of (C.5) than the left-hand side. A special case that sometimes arises is

$$(\mathbf{I} + \mathbf{A} \mathbf{B})^{-1} \mathbf{A} = \mathbf{A} (\mathbf{I} + \mathbf{B} \mathbf{A})^{-1}. \quad (\text{C.6})$$

Another useful identity involving inverses is the following:

$$(\mathbf{A} + \mathbf{B} \mathbf{D}^{-1} \mathbf{C})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{B} (\mathbf{D} + \mathbf{C} \mathbf{A}^{-1} \mathbf{B})^{-1} \mathbf{C} \mathbf{A}^{-1} \quad (\text{C.7})$$

which is known as the *Woodbury identity* and which can be verified by multiplying both sides by  $(\mathbf{A} + \mathbf{B} \mathbf{D}^{-1} \mathbf{C})$ . This is useful, for instance, when  $\mathbf{A}$  is large and diagonal, and hence easy to invert, while  $\mathbf{B}$  has many rows but few columns (and conversely for  $\mathbf{C}$ ) so that the right-hand side is much cheaper to evaluate than the left-hand side.

A set of vectors  $\{\mathbf{a}_1, \dots, \mathbf{a}_N\}$  is said to be *linearly independent* if the relation  $\sum_n \alpha_n \mathbf{a}_n = 0$  holds only if all  $\alpha_n = 0$ . This implies that none of the vectors can be expressed as a linear combination of the remainder. The rank of a matrix is the maximum number of linearly independent rows (or equivalently the maximum number of linearly independent columns).

---

## Traces and Determinants

Trace and determinant apply to square matrices. The trace  $\text{Tr}(\mathbf{A})$  of a matrix  $\mathbf{A}$  is defined as the sum of the elements on the leading diagonal. By writing out the indices, we see that

$$\text{Tr}(\mathbf{A} \mathbf{B}) = \text{Tr}(\mathbf{B} \mathbf{A}). \quad (\text{C.8})$$

By applying this formula multiple times to the product of three matrices, we see that

$$\text{Tr}(\mathbf{A} \mathbf{B} \mathbf{C}) = \text{Tr}(\mathbf{C} \mathbf{A} \mathbf{B}) = \text{Tr}(\mathbf{B} \mathbf{C} \mathbf{A}) \quad (\text{C.9})$$

which is known as the *cyclic* property of the trace operator and which clearly extends to the product of any number of matrices. The determinant  $|\mathbf{A}|$  of an  $N \times N$  matrix  $\mathbf{A}$  is defined by

$$|\mathbf{A}| = \sum (\pm 1) A_{1i_1} A_{2i_2} \cdots A_{Ni_N} \quad (\text{C.10})$$

in which the sum is taken over all products consisting of precisely one element from each row and one element from each column, with a coefficient  $+1$  or  $-1$  according

to whether the permutation  $i_1 i_2 \dots i_N$  is even or odd, respectively. Note that  $|\mathbf{I}| = 1$ . Thus, for a  $2 \times 2$  matrix, the determinant takes the form

$$|\mathbf{A}| = \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} = a_{11}a_{22} - a_{12}a_{21}. \quad (\text{C.11})$$

The determinant of a product of two matrices is given by

$$|\mathbf{AB}| = |\mathbf{A}||\mathbf{B}| \quad (\text{C.12})$$

as can be shown from (C.10). Also, the determinant of an inverse matrix is given by

$$|\mathbf{A}^{-1}| = \frac{1}{|\mathbf{A}|} \quad (\text{C.13})$$

which can be shown by taking the determinant of (C.2) and applying (C.12).

If  $\mathbf{A}$  and  $\mathbf{B}$  are matrices of size  $N \times M$ , then

$$|\mathbf{I}_N + \mathbf{AB}^T| = |\mathbf{I}_M + \mathbf{A}^T\mathbf{B}|. \quad (\text{C.14})$$

A useful special case is

$$|\mathbf{I}_N + \mathbf{ab}^T| = 1 + \mathbf{a}^T\mathbf{b} \quad (\text{C.15})$$

where  $\mathbf{a}$  and  $\mathbf{b}$  are  $N$ -dimensional column vectors.

## Matrix Derivatives

Sometimes we need to consider derivatives of vectors and matrices with respect to scalars. The derivative of a vector  $\mathbf{a}$  with respect to a scalar  $x$  is itself a vector whose components are given by

$$\left( \frac{\partial \mathbf{a}}{\partial x} \right)_i = \frac{\partial a_i}{\partial x} \quad (\text{C.16})$$

with an analogous definition for the derivative of a matrix. Derivatives with respect to vectors and matrices can also be defined, for instance

$$\left( \frac{\partial x}{\partial \mathbf{a}} \right)_i = \frac{\partial x}{\partial a_i} \quad (\text{C.17})$$

and similarly

$$\left( \frac{\partial \mathbf{a}}{\partial \mathbf{b}} \right)_{ij} = \frac{\partial a_i}{\partial b_j}. \quad (\text{C.18})$$

The following is easily proven by writing out the components

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{a}) = \frac{\partial}{\partial \mathbf{x}} (\mathbf{a}^T \mathbf{x}) = \mathbf{a}. \quad (\text{C.19})$$

Similarly

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{AB}) = \frac{\partial \mathbf{A}}{\partial \mathbf{x}} \mathbf{B} + \mathbf{A} \frac{\partial \mathbf{B}}{\partial \mathbf{x}}. \quad (\text{C.20})$$

The derivative of the inverse of a matrix can be expressed as

$$\frac{\partial}{\partial x} (\mathbf{A}^{-1}) = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1} \quad (\text{C.21})$$

as can be shown by differentiating the equation  $\mathbf{A}^{-1} \mathbf{A} = \mathbf{I}$  using (C.20) and then right multiplying by  $\mathbf{A}^{-1}$ . Also

$$\frac{\partial}{\partial x} \ln |\mathbf{A}| = \text{Tr} \left( \mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \right) \quad (\text{C.22})$$

which we shall prove later. If we choose  $x$  to be one of the elements of  $\mathbf{A}$ , we have

$$\frac{\partial}{\partial A_{ij}} \text{Tr}(\mathbf{AB}) = B_{ji} \quad (\text{C.23})$$

as can be seen by writing out the matrices using index notation. We can write this result more compactly in the form

$$\frac{\partial}{\partial \mathbf{A}} \text{Tr}(\mathbf{AB}) = \mathbf{B}^T. \quad (\text{C.24})$$

With this notation, we have the following properties

$$\frac{\partial}{\partial \mathbf{A}} \text{Tr}(\mathbf{A}^T \mathbf{B}) = \mathbf{B} \quad (\text{C.25})$$

$$\frac{\partial}{\partial \mathbf{A}} \text{Tr}(\mathbf{A}) = \mathbf{I} \quad (\text{C.26})$$

$$\frac{\partial}{\partial \mathbf{A}} \text{Tr}(\mathbf{ABA}^T) = \mathbf{A}(\mathbf{B} + \mathbf{B}^T) \quad (\text{C.27})$$

which can again be proven by writing out the matrix indices. We also have

$$\frac{\partial}{\partial \mathbf{A}} \ln |\mathbf{A}| = (\mathbf{A}^{-1})^T \quad (\text{C.28})$$

which follows from (C.22) and (C.26).

## Eigenvector Equation

For a square matrix  $\mathbf{A}$  of size  $M \times M$ , the eigenvector equation is defined by

$$\mathbf{A} \mathbf{u}_i = \lambda_i \mathbf{u}_i \quad (\text{C.29})$$

for  $i = 1, \dots, M$ , where  $\mathbf{u}_i$  is an *eigenvector* and  $\lambda_i$  is the corresponding *eigenvalue*. This can be viewed as a set of  $M$  simultaneous homogeneous linear equations, and the condition for a solution is that

$$|\mathbf{A} - \lambda_i \mathbf{I}| = 0 \quad (\text{C.30})$$

which is known as the *characteristic equation*. Because this is a polynomial of order  $M$  in  $\lambda_i$ , it must have  $M$  solutions (though these need not all be distinct). The rank of  $\mathbf{A}$  is equal to the number of nonzero eigenvalues.

Of particular interest are symmetric matrices, which arise as covariance matrices, kernel matrices, and Hessians. Symmetric matrices have the property that  $A_{ij} = A_{ji}$ , or equivalently  $\mathbf{A}^T = \mathbf{A}$ . The inverse of a symmetric matrix is also symmetric, as can be seen by taking the transpose of  $\mathbf{A}^{-1} \mathbf{A} = \mathbf{I}$  and using  $\mathbf{A} \mathbf{A}^{-1} = \mathbf{I}$  together with the symmetry of  $\mathbf{I}$ .

In general, the eigenvalues of a matrix are complex numbers, but for symmetric matrices the eigenvalues  $\lambda_i$  are real. This can be seen by first left multiplying (C.29) by  $(\mathbf{u}_i^*)^T$ , where  $\star$  denotes the complex conjugate, to give

$$(\mathbf{u}_i^*)^T \mathbf{A} \mathbf{u}_i = \lambda_i (\mathbf{u}_i^*)^T \mathbf{u}_i. \quad (\text{C.31})$$

Next we take the complex conjugate of (C.29) and left multiply by  $\mathbf{u}_i^T$  to give

$$\mathbf{u}_i^T \mathbf{A} \mathbf{u}_i^* = \lambda_i^* \mathbf{u}_i^T \mathbf{u}_i^*. \quad (\text{C.32})$$

where we have used  $\mathbf{A}^* = \mathbf{A}$  because we consider only real matrices  $\mathbf{A}$ . Taking the transpose of the second of these equations, and using  $\mathbf{A}^T = \mathbf{A}$ , we see that the left-hand sides of the two equations are equal, and hence that  $\lambda_i^* = \lambda_i$  and so  $\lambda_i$  must be real.

The eigenvectors  $\mathbf{u}_i$  of a real symmetric matrix can be chosen to be orthonormal (i.e., orthogonal and of unit length) so that

$$\mathbf{u}_i^T \mathbf{u}_j = I_{ij} \quad (\text{C.33})$$

where  $I_{ij}$  are the elements of the identity matrix  $\mathbf{I}$ . To show this, we first left multiply (C.29) by  $\mathbf{u}_j^T$  to give

$$\mathbf{u}_j^T \mathbf{A} \mathbf{u}_i = \lambda_i \mathbf{u}_j^T \mathbf{u}_i \quad (\text{C.34})$$

and hence, by exchange of indices, we have

$$\mathbf{u}_i^T \mathbf{A} \mathbf{u}_j = \lambda_j \mathbf{u}_i^T \mathbf{u}_j. \quad (\text{C.35})$$

We now take the transpose of the second equation and make use of the symmetry property  $\mathbf{A}^T = \mathbf{A}$ , and then subtract the two equations to give

$$(\lambda_i - \lambda_j) \mathbf{u}_i^T \mathbf{u}_j = 0. \quad (\text{C.36})$$

Hence, for  $\lambda_i \neq \lambda_j$ , we have  $\mathbf{u}_i^T \mathbf{u}_j = 0$ , and hence  $\mathbf{u}_i$  and  $\mathbf{u}_j$  are orthogonal. If the two eigenvalues are equal, then any linear combination  $\alpha \mathbf{u}_i + \beta \mathbf{u}_j$  is also an eigenvector with the same eigenvalue, so we can select one linear combination arbitrarily,

and then choose the second to be orthogonal to the first (it can be shown that the degenerate eigenvectors are never linearly dependent). Hence the eigenvectors can be chosen to be orthogonal, and by normalizing can be set to unit length. Because there are  $M$  eigenvalues, the corresponding  $M$  orthogonal eigenvectors form a complete set and so any  $M$ -dimensional vector can be expressed as a linear combination of the eigenvectors.

We can take the eigenvectors  $\mathbf{u}_i$  to be the columns of an  $M \times M$  matrix  $\mathbf{U}$ , which from orthonormality satisfies

$$\mathbf{U}^T \mathbf{U} = \mathbf{I}. \quad (\text{C.37})$$

Such a matrix is said to be *orthogonal*. Interestingly, the rows of this matrix are also orthogonal, so that  $\mathbf{U}\mathbf{U}^T = \mathbf{I}$ . To show this, note that (C.37) implies  $\mathbf{U}^T \mathbf{U} \mathbf{U}^{-1} = \mathbf{U}^{-1} = \mathbf{U}^T$  and so  $\mathbf{U}\mathbf{U}^{-1} = \mathbf{U}\mathbf{U}^T = \mathbf{I}$ . Using (C.12), it also follows that  $|\mathbf{U}| = 1$ .

The eigenvector equation (C.29) can be expressed in terms of  $\mathbf{U}$  in the form

$$\mathbf{A}\mathbf{U} = \mathbf{U}\mathbf{\Lambda} \quad (\text{C.38})$$

where  $\mathbf{\Lambda}$  is an  $M \times M$  diagonal matrix whose diagonal elements are given by the eigenvalues  $\lambda_i$ .

If we consider a column vector  $\mathbf{x}$  that is transformed by an orthogonal matrix  $\mathbf{U}$  to give a new vector

$$\tilde{\mathbf{x}} = \mathbf{U}\mathbf{x} \quad (\text{C.39})$$

then the length of the vector is preserved because

$$\tilde{\mathbf{x}}^T \tilde{\mathbf{x}} = \mathbf{x}^T \mathbf{U}^T \mathbf{U} \mathbf{x} = \mathbf{x}^T \mathbf{x} \quad (\text{C.40})$$

and similarly the angle between any two such vectors is preserved because

$$\tilde{\mathbf{x}}^T \tilde{\mathbf{y}} = \mathbf{x}^T \mathbf{U}^T \mathbf{U} \mathbf{y} = \mathbf{x}^T \mathbf{y}. \quad (\text{C.41})$$

Thus, multiplication by  $\mathbf{U}$  can be interpreted as a rigid rotation of the coordinate system.

From (C.38), it follows that

$$\mathbf{U}^T \mathbf{A} \mathbf{U} = \mathbf{\Lambda} \quad (\text{C.42})$$

and because  $\mathbf{\Lambda}$  is a diagonal matrix, we say that the matrix  $\mathbf{A}$  is *diagonalized* by the matrix  $\mathbf{U}$ . If we left multiply by  $\mathbf{U}$  and right multiply by  $\mathbf{U}^T$ , we obtain

$$\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T \quad (\text{C.43})$$

Taking the inverse of this equation, and using (C.3) together with  $\mathbf{U}^{-1} = \mathbf{U}^T$ , we have

$$\mathbf{A}^{-1} = \mathbf{U} \mathbf{\Lambda}^{-1} \mathbf{U}^T. \quad (\text{C.44})$$

These last two equations can also be written in the form

$$\mathbf{A} = \sum_{i=1}^M \lambda_i \mathbf{u}_i \mathbf{u}_i^T \quad (\text{C.45})$$

$$\mathbf{A}^{-1} = \sum_{i=1}^M \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T. \quad (\text{C.46})$$

If we take the determinant of (C.43), and use (C.12), we obtain

$$|\mathbf{A}| = \prod_{i=1}^M \lambda_i. \quad (\text{C.47})$$

Similarly, taking the trace of (C.43), and using the cyclic property (C.8) of the trace operator together with  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ , we have

$$\text{Tr}(\mathbf{A}) = \sum_{i=1}^M \lambda_i. \quad (\text{C.48})$$

We leave it as an exercise for the reader to verify (C.22) by making use of the results (C.33), (C.45), (C.46), and (C.47).

A matrix  $\mathbf{A}$  is said to be *positive definite*, denoted by  $\mathbf{A} \succ 0$ , if  $\mathbf{w}^T \mathbf{A} \mathbf{w} > 0$  for all values of the vector  $\mathbf{w}$ . Equivalently, a positive definite matrix has  $\lambda_i > 0$  for all of its eigenvalues (as can be seen by setting  $\mathbf{w}$  to each of the eigenvectors in turn, and by noting that an arbitrary vector can be expanded as a linear combination of the eigenvectors). Note that positive definite is not the same as all the elements being positive. For example, the matrix

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \quad (\text{C.49})$$

has eigenvalues  $\lambda_1 \simeq 5.37$  and  $\lambda_2 \simeq -0.37$ . A matrix is said to be *positive semidefinite* if  $\mathbf{w}^T \mathbf{A} \mathbf{w} \geq 0$  holds for all values of  $\mathbf{w}$ , which is denoted  $\mathbf{A} \succeq 0$ , and is equivalent to  $\lambda_i \geq 0$ .

## Appendix D. Calculus of Variations

We can think of a function  $y(x)$  as being an operator that, for any input value  $x$ , returns an output value  $y$ . In the same way, we can define a *functional*  $F[y]$  to be an operator that takes a function  $y(x)$  and returns an output value  $F$ . An example of a functional is the length of a curve drawn in a two-dimensional plane in which the path of the curve is defined in terms of a function. In the context of machine learning, a widely used functional is the entropy  $H[x]$  for a continuous variable  $x$  because, for any choice of probability density function  $p(x)$ , it returns a scalar value representing the entropy of  $x$  under that density. Thus the entropy of  $p(x)$  could equally well have been written as  $H[p]$ .

A common problem in conventional calculus is to find a value of  $x$  that maximizes (or minimizes) a function  $y(x)$ . Similarly, in the calculus of variations we seek a function  $y(x)$  that maximizes (or minimizes) a functional  $F[y]$ . That is, of all possible functions  $y(x)$ , we wish to find the particular function for which the functional  $F[y]$  is a maximum (or minimum). The calculus of variations can be used, for instance, to show that the shortest path between two points is a straight line or that the maximum entropy distribution is a Gaussian.

If we weren't familiar with the rules of ordinary calculus, we could evaluate a conventional derivative  $dy/dx$  by making a small change  $\epsilon$  to the variable  $x$  and then expanding in powers of  $\epsilon$ , so that

$$y(x + \epsilon) = y(x) + \frac{dy}{dx}\epsilon + O(\epsilon^2) \quad (\text{D.1})$$

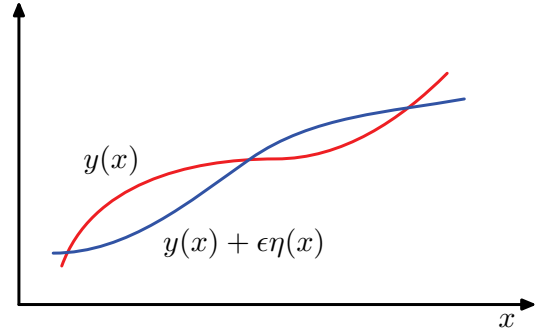
and finally taking the limit  $\epsilon \rightarrow 0$ . Similarly, for a function of several variables  $y(x_1, \dots, x_D)$ , the corresponding partial derivatives are defined by

$$y(x_1 + \epsilon_1, \dots, x_D + \epsilon_D) = y(x_1, \dots, x_D) + \sum_{i=1}^D \frac{\partial y}{\partial x_i} \epsilon_i + O(\epsilon^2). \quad (\text{D.2})$$

The analogous definition of a functional derivative arises when we consider how much a functional  $F[y]$  changes when we make a small change  $\epsilon\eta(x)$  to the function



**Figure D.1** A functional derivative can be defined by considering how the value of a functional  $F[y]$  changes when the function  $y(x)$  is changed to  $y(x) + \epsilon\eta(x)$  where  $\eta(x)$  is an arbitrary function of  $x$ .



$y(x)$ , where  $\eta(x)$  is an arbitrary function of  $x$ , as illustrated in Figure D.1. We denote the functional derivative of  $E[f]$  with respect to  $f(x)$  by  $\delta F/\delta f(x)$ , and define it by the following relation:

$$F[y(x) + \epsilon\eta(x)] = F[y(x)] + \epsilon \int \frac{\delta F}{\delta y(x)} \eta(x) dx + O(\epsilon^2). \quad (D.3)$$

This can be seen as a natural extension of (D.2) in which  $F[y]$  now depends on a continuous set of variables, namely the values of  $y$  at all points  $x$ . Requiring that the functional be stationary with respect to small variations in the function  $y(x)$  gives

$$\int \frac{\delta E}{\delta y(x)} \eta(x) dx = 0. \quad (D.4)$$

Because this must hold for an arbitrary choice of  $\eta(x)$ , it follows that the functional derivative must vanish. To see this, imagine choosing a perturbation  $\eta(x)$  that is zero everywhere except in the neighbourhood of a point  $\hat{x}$ , in which case the functional derivative must be zero at  $x = \hat{x}$ . However, because this must be true for every choice of  $\hat{x}$ , the functional derivative must vanish for all values of  $x$ .

Consider a functional that is defined by an integral over a function  $G(y, y', x)$  that depends on both  $y(x)$  and its derivative  $y'(x)$  as well as having a direct dependence on  $x$

$$F[y] = \int G(y(x), y'(x), x) dx \quad (D.5)$$

where the value of  $y(x)$  is assumed to be fixed at the boundary of the region of integration (which might be at infinity). If we now consider variations in the function  $y(x)$ , we obtain

$$F[y(x) + \epsilon\eta(x)] = F[y(x)] + \epsilon \int \left\{ \frac{\partial G}{\partial y} \eta(x) + \frac{\partial G}{\partial y'} \eta'(x) \right\} dx + O(\epsilon^2). \quad (D.6)$$

We now have to cast this in the form (D.3). To do so, we integrate the second term by parts and make use of the fact that  $\eta(x)$  must vanish at the boundary of the integral (because  $y(x)$  is fixed at the boundary). This gives

$$F[y(x) + \epsilon\eta(x)] = F[y(x)] + \epsilon \int \left\{ \frac{\partial G}{\partial y} - \frac{d}{dx} \left( \frac{\partial G}{\partial y'} \right) \right\} \eta(x) dx + O(\epsilon^2) \quad (D.7)$$

from which we can read off the functional derivative by comparison with (D.3). Requiring that the functional derivative vanishes then gives

$$\frac{\partial G}{\partial y} - \frac{d}{dx} \left( \frac{\partial G}{\partial y'} \right) = 0 \quad (\text{D.8})$$

which are known as the *Euler-Lagrange* equations. For example, if

$$G = y(x)^2 + (y'(x))^2 \quad (\text{D.9})$$

then the Euler-Lagrange equations take the form

$$y(x) - \frac{d^2 y}{dx^2} = 0. \quad (\text{D.10})$$

This second order differential equation can be solved for  $y(x)$  by making use of the boundary conditions on  $y(x)$ .

Often, we consider functionals defined by integrals whose integrands take the form  $G(y, x)$  and that do not depend on the derivatives of  $y(x)$ . In this case, stationarity simply requires that  $\partial G / \partial y(x) = 0$  for all values of  $x$ .

If we are optimizing a functional with respect to a probability distribution, then we need to maintain the normalization constraint on the probabilities. This is often most conveniently done using a Lagrange multiplier, which then allows an unconstrained optimization to be performed.

The extension of the above results to a multidimensional variable  $\mathbf{x}$  is straightforward. For a more comprehensive discussion of the calculus of variations, see Sagan (1969).

## Appendix E. Lagrange Multipliers

*Lagrange multipliers*, also sometimes called *undetermined multipliers*, are used to find the stationary points of a function of several variables subject to one or more constraints.

Consider the problem of finding the maximum of a function  $f(x_1, x_2)$  subject to a constraint relating  $x_1$  and  $x_2$ , which we write in the form

$$g(x_1, x_2) = 0. \quad (\text{E.1})$$

One approach would be to solve the constraint equation (E.1) and thus express  $x_2$  as a function of  $x_1$  in the form  $x_2 = h(x_1)$ . This can then be substituted into  $f(x_1, x_2)$  to give a function of  $x_1$  alone of the form  $f(x_1, h(x_1))$ . The maximum with respect to  $x_1$  could then be found by differentiation in the usual way, to give the stationary value  $x_1^*$ , with the corresponding value of  $x_2$  given by  $x_2^* = h(x_1^*)$ .

One problem with this approach is that it may be difficult to find an analytic solution of the constraint equation that allows  $x_2$  to be expressed as an explicit function of  $x_1$ . Also, this approach treats  $x_1$  and  $x_2$  differently and so spoils the natural symmetry between these variables.

A more elegant, and often simpler, approach is based on the introduction of a parameter  $\lambda$  called a **Lagrange multiplier**. We shall motivate this technique from a geometrical perspective. Consider a  $D$ -dimensional variable  $\mathbf{x}$  with components  $x_1, \dots, x_D$ . The constraint equation  $g(\mathbf{x}) = 0$  then represents a  $(D-1)$ -dimensional surface in  $\mathbf{x}$ -space as indicated in Figure E.1.

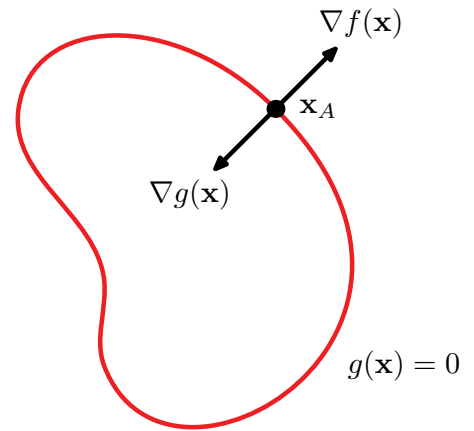
We first note that at **any point on the constraint surface the gradient  $\nabla g(\mathbf{x})$  of the constraint function will be orthogonal to the surface**. To see this, consider a point  $\mathbf{x}$  that lies on the constraint surface, and consider a nearby point  $\mathbf{x} + \epsilon$  that also lies on the surface. If we make a Taylor expansion around  $\mathbf{x}$ , we have

$$g(\mathbf{x} + \epsilon) \simeq g(\mathbf{x}) + \epsilon^T \nabla g(\mathbf{x}). \quad (\text{E.2})$$

Because both  $\mathbf{x}$  and  $\mathbf{x} + \epsilon$  lie on the constraint surface, we have  $g(\mathbf{x}) = g(\mathbf{x} + \epsilon)$  and hence  $\epsilon^T \nabla g(\mathbf{x}) \simeq 0$ . In the limit  $\|\epsilon\| \rightarrow 0$  we have  $\epsilon^T \nabla g(\mathbf{x}) = 0$ , and because  $\epsilon$  is

反证法：否则会在该曲面上移动到另一个点，其值大于0

**Figure E.1** A geometrical picture of the technique of Lagrange multipliers in which we seek to maximize a function  $f(\mathbf{x})$ , subject to the constraint  $g(\mathbf{x}) = 0$ . If  $\mathbf{x}$  is  $D$  dimensional, the constraint  $g(\mathbf{x}) = 0$  corresponds to a subspace of dimensionality  $D - 1$ , indicated by the red curve. The problem can be solved by optimizing the Lagrangian function  $L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$ .



then parallel to the constraint surface  $g(\mathbf{x}) = 0$ , we see that the vector  $\nabla g$  is normal to the surface.

Next we seek a point  $\mathbf{x}^*$  on the constraint surface such that  $f(\mathbf{x})$  is maximized. Such a point must have the property that the vector  $\nabla f(\mathbf{x})$  is also orthogonal to the constraint surface, as illustrated in Figure E.1, because otherwise we could increase the value of  $f(\mathbf{x})$  by moving a short distance along the constraint surface. Thus  $\nabla f$  and  $\nabla g$  are parallel (or anti-parallel) vectors, and so there must exist a parameter  $\lambda$  such that

$$\nabla f + \lambda \nabla g = 0 \quad (\text{E.3})$$

where  $\lambda \neq 0$  is known as a *Lagrange multiplier*. Note that  $\lambda$  can have either sign.

At this point, it is convenient to introduce the *Lagrangian* function defined by

$$\text{加減都可以, 因为}\lambda\text{符号无约束} \quad L(\mathbf{x}, \lambda) \equiv f(\mathbf{x}) + \lambda g(\mathbf{x}). \quad (\text{E.4})$$

The constrained stationarity condition (E.3) is obtained by setting  $\nabla_{\mathbf{x}} L = 0$ . Furthermore, the condition  $\partial L / \partial \lambda = 0$  leads to the constraint equation  $g(\mathbf{x}) = 0$ .

Thus to find the maximum of a function  $f(\mathbf{x})$  subject to the constraint  $g(\mathbf{x}) = 0$ , we define the Lagrangian function given by (E.4) and we then find the stationary point of  $L(\mathbf{x}, \lambda)$  with respect to both  $\mathbf{x}$  and  $\lambda$ . For a  $D$ -dimensional vector  $\mathbf{x}$ , this gives  $D + 1$  equations that determine both the stationary point  $\mathbf{x}^*$  and the value of  $\lambda$ . If we are only interested in  $\mathbf{x}^*$ , then we can eliminate  $\lambda$  from the stationarity equations without needing to find its value (hence the term ‘undetermined multiplier’).

As a simple example, suppose we wish to find the stationary point of the function  $f(x_1, x_2) = 1 - x_1^2 - x_2^2$  subject to the constraint  $g(x_1, x_2) = x_1 + x_2 - 1 = 0$ , as illustrated in Figure E.2. The corresponding Lagrangian function is given by

$$L(\mathbf{x}, \lambda) = 1 - x_1^2 - x_2^2 + \lambda(x_1 + x_2 - 1). \quad (\text{E.5})$$

The conditions for this Lagrangian to be stationary with respect to  $x_1$ ,  $x_2$ , and  $\lambda$  give the following coupled equations:

$$-2x_1 + \lambda = 0 \quad (\text{E.6})$$

$$-2x_2 + \lambda = 0 \quad (\text{E.7})$$

$$x_1 + x_2 - 1 = 0. \quad (\text{E.8})$$

补上这个式子才和原问题等价

本文从梯度分析推导出  
Lagrange multiplier

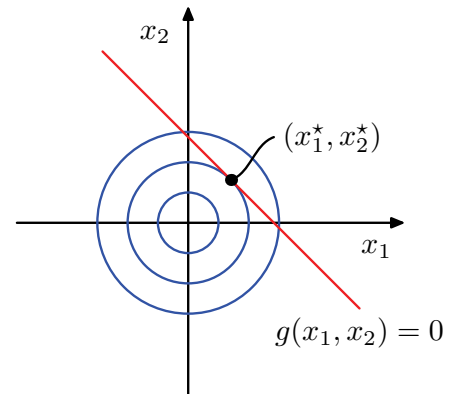
1. 在约束曲面 $g(\mathbf{x})=0$ 上, 任意一点, 其梯度为0

2. 梯度为0是 $f$ 取极值的必要条件

3. 在 $f$ 与 $g(\mathbf{x})=0$ 的交汇点上,  $f$ 与 $g$ 的梯度都为0, 说明二者梯度正交—注意和KKT一样  
Lagrange乘子仅仅是必要条件

对 $\lambda$ 求导, 梯度取0, 可得到  
 $g(\mathbf{x})=0$

**Figure E.2** A simple example of the use of Lagrange multipliers in which the aim is to maximize  $f(x_1, x_2) = 1 - x_1^2 - x_2^2$  subject to the constraint  $g(x_1, x_2) = 0$  where  $g(x_1, x_2) = x_1 + x_2 - 1$ . The circles show contours of the function  $f(x_1, x_2)$ , and the diagonal line shows the constraint surface  $g(x_1, x_2) = 0$ .



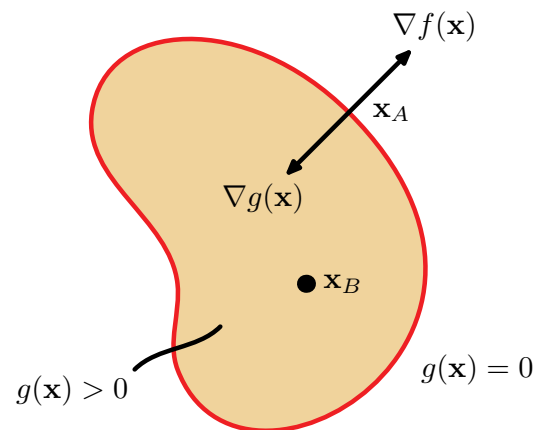
Solution of these equations then gives the stationary point as  $(x_1^*, x_2^*) = (\frac{1}{2}, \frac{1}{2})$ , and the corresponding value for the Lagrange multiplier is  $\lambda = 1$ .

So far, we have considered the problem of maximizing a function subject to an *equality constraint* of the form  $g(\mathbf{x}) = 0$ . We now consider the problem of maximizing  $f(\mathbf{x})$  subject to an *inequality constraint* of the form  $g(\mathbf{x}) \geq 0$ , as illustrated in Figure E.3.

There are now two kinds of solution possible, according to whether the constrained stationary point lies in the region where  $g(\mathbf{x}) > 0$ , in which case the constraint is *inactive*, or whether it lies on the boundary  $g(\mathbf{x}) = 0$ , in which case the constraint is said to be *active*. In the former case, the function  $g(\mathbf{x})$  plays no role and so the stationary condition is simply  $\nabla f(\mathbf{x}) = 0$ . This again corresponds to a stationary point of the Lagrange function (E.4) but this time with  $\lambda = 0$ . The latter case, where the solution lies on the boundary, is analogous to the equality constraint discussed previously and corresponds to a stationary point of the Lagrange function (E.4) with  $\lambda \neq 0$ . Now, however, the sign of the Lagrange multiplier is crucial, because the function  $f(\mathbf{x})$  will only be at a maximum if its gradient is oriented away from the region  $g(\mathbf{x}) > 0$ , as illustrated in Figure E.3. We therefore have  $\nabla f(\mathbf{x}) = -\lambda \nabla g(\mathbf{x})$  for some value of  $\lambda > 0$ .

For either of these two cases, the product  $\lambda g(\mathbf{x}) = 0$ . Thus the solution to the

**Figure E.3** Illustration of the problem of maximizing  $f(\mathbf{x})$  subject to the inequality constraint  $g(\mathbf{x}) \geq 0$ .



problem of maximizing  $f(\mathbf{x})$  subject to  $g(\mathbf{x}) \geq 0$  is obtained by optimizing the Lagrange function (E.4) with respect to  $\mathbf{x}$  and  $\lambda$  subject to the conditions

$$g(\mathbf{x}) \geq 0 \quad (\text{E.9})$$

$$\lambda \geq 0 \quad (\text{E.10})$$

$$\lambda g(\mathbf{x}) = 0 \quad (\text{E.11})$$

These are known as the *Karush-Kuhn-Tucker* (KKT) conditions (Karush, 1939; Kuhn and Tucker, 1951).

Note that if we wish to minimize (rather than maximize) the function  $f(\mathbf{x})$  subject to an inequality constraint  $g(\mathbf{x}) \geq 0$ , then we minimize the Lagrangian function  $L(\mathbf{x}, \lambda) = f(\mathbf{x}) - \lambda g(\mathbf{x})$  with respect to  $\mathbf{x}$ , again subject to  $\lambda \geq 0$ .

Finally, it is straightforward to extend the technique of Lagrange multipliers to the case of multiple equality and inequality constraints. Suppose we wish to maximize  $f(\mathbf{x})$  subject to  $g_j(\mathbf{x}) = 0$  for  $j = 1, \dots, J$ , and  $h_k(\mathbf{x}) \geq 0$  for  $k = 1, \dots, K$ . We then introduce Lagrange multipliers  $\{\lambda_j\}$  and  $\{\mu_k\}$ , and then optimize the Lagrangian function given by

$$L(\mathbf{x}, \{\lambda_j\}, \{\mu_k\}) = f(\mathbf{x}) + \sum_{j=1}^J \lambda_j g_j(\mathbf{x}) + \sum_{k=1}^K \mu_k h_k(\mathbf{x}) \quad (\text{E.12})$$

subject to  $\mu_k \geq 0$  and  $\mu_k h_k(\mathbf{x}) = 0$  for  $k = 1, \dots, K$ . Extensions to constrained functional derivatives are similarly straightforward. For a more detailed discussion of the technique of Lagrange multipliers, see Nocedal and Wright (1999).