# Appendix B. 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 $\mathrm{H}[x]$ 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 $\mathrm{H}[p]$.

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 were not familiar with the rules of ordinary calculus, we could evaluate a conventional derivative $\mathrm{d}y/\mathrm{d}x$ by making a small change $\epsilon$ to the variable $x$ and then expanding in powers of $\epsilon$, so that
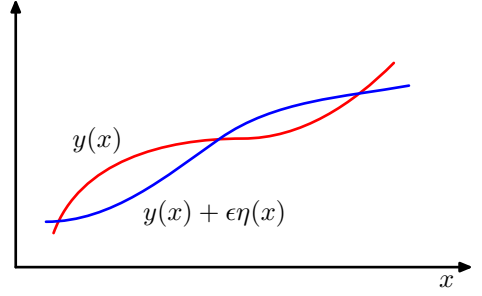
$$y(x + \epsilon) = y(x) + \frac{\mathrm{d}y}{\mathrm{d}x}\epsilon + \mathcal{O}(\epsilon^2) \tag{B.1}$$

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

$$y(x_1 + \epsilon_1, \ldots, x_D + \epsilon_D) = y(x_1, \ldots, x_D) + \sum_{i=1}^{D} \frac{\partial y}{\partial x_i}\epsilon_i + \mathcal{O}(\epsilon^2). \tag{B.2}$$

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 B.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 B.1. We denote the functional derivative of $F[y]$ with respect to $y(x)$ by $\delta F/\delta y(x)$ and define it by the following relation:

$$F[y(x) + \epsilon\eta(x)] = F[y(x)] + \epsilon \int \frac{\delta F}{\delta y(x)}\eta(x)\,\mathrm{d}x + \mathcal{O}(\epsilon^2). \qquad \text{(B.3)}$$

This can be seen as a natural extension of (B.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

$$\int \frac{\delta F}{\delta y(x)}\eta(x)\,\mathrm{d}x = 0. \qquad \text{(B.4)}$$

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 $\widehat{x}$, in which case the functional derivative must be zero at $x = \widehat{x}$. However, because this must be true for every choice of $\widehat{x}$, the functional derivative must vanish for all values of $x$.

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

$$F[y] = \int G(y(x), y'(x), x)\ \mathrm{d}x \qquad \text{(B.5)}$$

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

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

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

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

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

$$\frac{\partial G}{\partial y} - \frac{\mathrm{d}}{\mathrm{d}x}\left(\frac{\partial G}{\partial y'}\right) = 0, \tag{B.8}$$

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

$$G = y(x)^2 + \left(y'(x)\right)^2 \tag{B.9}$$

then the Euler–Lagrange equations take the form

$$y(x) - \frac{\mathrm{d}^2 y}{\mathrm{d}x^2} = 0. \tag{B.10}$$

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 $x$.

*Appendix C*

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 multi-dimensional variable $\mathbf{x}$ is straightforward. For a more comprehensive discussion of the calculus of variations, see Sagan (1969).