

Derivative as linear approximation

Linear functions

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is said to be *linear* if the following two conditions are satisfied:

- i) (Homogeneity) For all $\lambda \in \mathbb{R}$ and for all $\mathbf{x} \in \mathbb{R}^n$, $f(\lambda\mathbf{x}) = \lambda f(\mathbf{x})$,
- ii) (Additivity) For all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$.

By taking $\lambda = 0$ in i), we see that $f(0) = 0$ for all linear functions. In the case of $n = 1$, this rules out functions whose graphs are straight lines that do not go through the origin. In this simplest setting, i) actually implies ii), and fixing $f(1)$ determines the entire function.

For $n > 1$, requirement ii) has bite. Observe that we can write $\mathbf{x} = \sum_{i=1}^n x_i \mathbf{e}^i$. By i), $f(x_i \mathbf{e}^i) = x_i f(\mathbf{e}^i)$ for all i, x_i . By ii),

$$f(\mathbf{x}) = f\left(\sum_{i=1}^n x_i \mathbf{e}^i\right) = \sum_{i=1}^n x_i f(\mathbf{e}^i).$$

Hence a linear function is completely determined by n values $f(\mathbf{e}^i)$. If $m = 1$, then $f(\mathbf{e}^i) \in \mathbb{R}$ for all i and letting $f(\mathbf{e}^i) = a_i$ we see that all real linear functions from \mathbb{R}^n are given by inner products $\mathbf{a} \cdot \mathbf{x} = \sum_{i=1}^n a_i x_i$.

If $m > 1$, then each $f(\mathbf{e}^i)$ is an m -dimensional vector. If we denote $\mathbf{a}^i = f(\mathbf{e}^i) \in \mathbb{R}^m$, we have as before $f(\mathbf{x}) = \sum_{i=1}^n x_i \mathbf{a}^i$. Writing $\mathbf{A} = [\mathbf{a}^1, \dots, \mathbf{a}^n]$ for the matrix consisting of columns \mathbf{a}^i . But this means that

$$f(\mathbf{x}) = \mathbf{A}\mathbf{x}.$$

Many of the properties of linear functions also extend to *affine* functions of the form

$$f(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{b},$$

for some $\mathbf{b} \in \mathbb{R}^m$. Actually this is not so bad because by shifting the origin to $(0, f(0))$, $\hat{f}(\mathbf{x}) := f(\mathbf{x}) - f(0) = \mathbf{A}\mathbf{x}$ is linear.

Why are linear functions so much simpler than non-linear? i) A change in \mathbf{x} has the same effect regardless of the starting point:

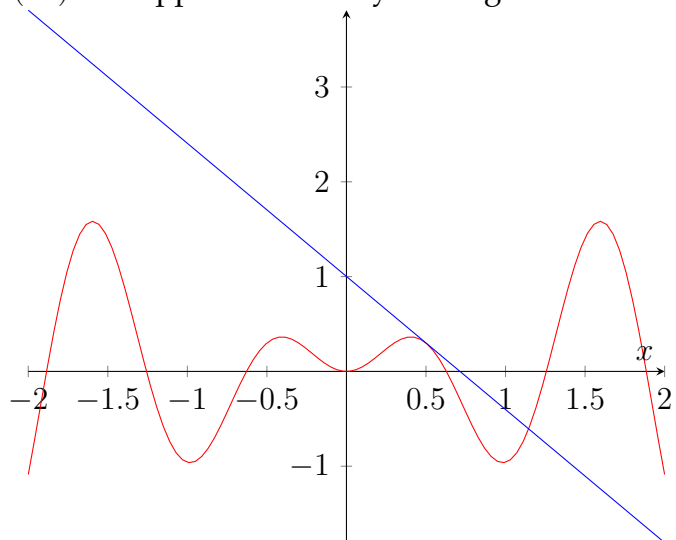
$$f(\mathbf{x}) - f(\hat{\mathbf{x}}) = \mathbf{A}(\mathbf{x} - \hat{\mathbf{x}}).$$

ii) A function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is said to be *surjective* (or *onto*) if for all $\mathbf{b} \in \mathbb{R}^m$, there is an $\mathbf{x} \in \mathbb{R}^n$ such that $f(\mathbf{x}) = \mathbf{b}$. f is said to be *injective* or *one-to-one* if for all $\mathbf{x} \neq \mathbf{x}'$, $f(\mathbf{x}) \neq f(\mathbf{x}')$. f is said to be *bijective* if it is injective and surjective. Bijective functions f have an *inverse function* $f^{-1} : \mathbb{R}^m \rightarrow \mathbb{R}^n$ such that $f^{-1}(f(\mathbf{x})) = \mathbf{x}$ and $f(f^{-1}(\mathbf{y})) = \mathbf{y}$. In matrix algebra, we saw that if $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$, then f is bijective if and only if \mathbf{A} has full rank. Gaussian elimination (or the determinant) gives an easy way of determining when linear functions are bijective and computing the inverse function $f^{-1}(\mathbf{x}) = \mathbf{A}^{-1}\mathbf{x}$.

The derivative: Inducing a linear approximation

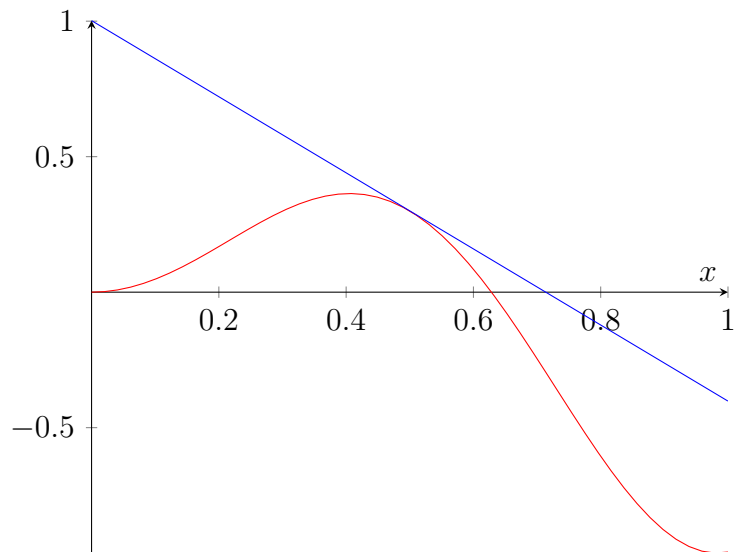
Motivation

Consider the graph of the following highly non-linear function $f(x) = x \sin(5x)$ and approximate it by its tangent at $x = 0.5$. Not a great success:



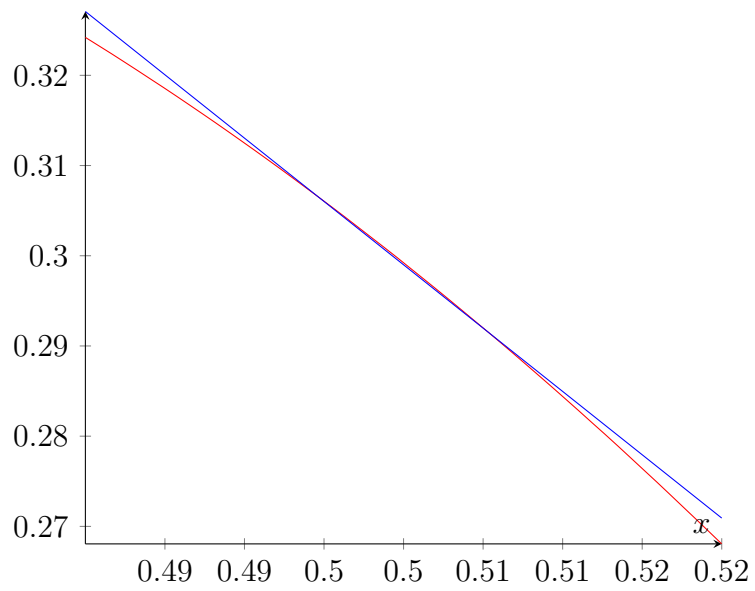
It is a bit less variable over the interval $[0, 1]$:

$$f(x) = x \sin(5x)$$



On the interval $[0.48, 0.52]$, it looks almost linear:

$$f(x) = x \sin(5x)$$



As we zoom closer to a fixed point (here $\hat{x} = 0.5$), the original highly irregular function starts resembling a linear function. The point of computing derivatives is exactly this: we want to get good linear approximations to non-linear functions near a given point $(\hat{x}, f(\hat{x}))$ on the graph of f .

Real valued functions of a single variable

From elementary calculus, the derivative of $f : \mathbb{R} \rightarrow \mathbb{R}$ at $\hat{x} \in \mathbb{R}$ is defined as:

$$Df(\hat{x}) = f'(\hat{x}) = \left. \frac{df(x)}{dx} \right|_{x=\hat{x}} = \lim_{h \rightarrow 0} \frac{f(\hat{x} + h) - f(\hat{x})}{h}.$$

You may recall that this numerical value (whenever the limit exists) can be interpreted as the slope of the tangent to f at \hat{x} .

This is of course fine, but an alternative way to think about the derivative at \hat{x} is as a linear approximation through origin $(\hat{x}, f(\hat{x}))$ that best approximates f near \hat{x} . If the limit exists, we can write:

$$f(\hat{x} + h) - f(\hat{x}) = Df(\hat{x})h + \varepsilon(h),$$

where $\lim_{h \rightarrow 0} \frac{\varepsilon(h)}{h} = 0$. Sometimes, the reminder $\varepsilon(h)$ is just written as higher order terms or h.o.t.

The reason for insisting on this seemingly trivial point is that this view of the derivative generalizes immediately to multivariate functions $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$. It also tells you immediately that for small h you get a good approximations for the change in the value $f(\hat{x} + h) - f(\hat{x}) = Df(\hat{x})h$. Notice also that for functions of a single variable, there is only a single direction h that we need to consider.

At this point, it may be a good idea to refresh the rules for computing derivatives (a handout on this is in the materials for week 2).

Real valued functions of a multiple variables

Consider next real valued functions defined in the $x - y$ plane, in other words functions $f : \mathbb{R}^2 \rightarrow \mathbb{R}$. Since we need a third dimension to graph the value corresponding to the value of the function at each point of the plane, we run into some difficulties in representing a function on a two-dimensional screen.

We have some options available:

1. Drawing a 3-d graph:

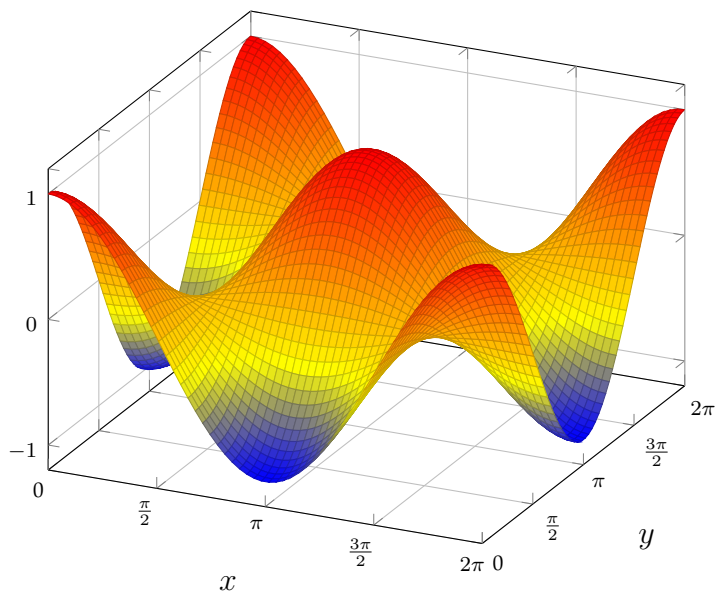


Figure 1: The graph of $f(x, y) = \cos(x)\cos(y)$.

2. Drawing 2-d slices of the 3-d graph:

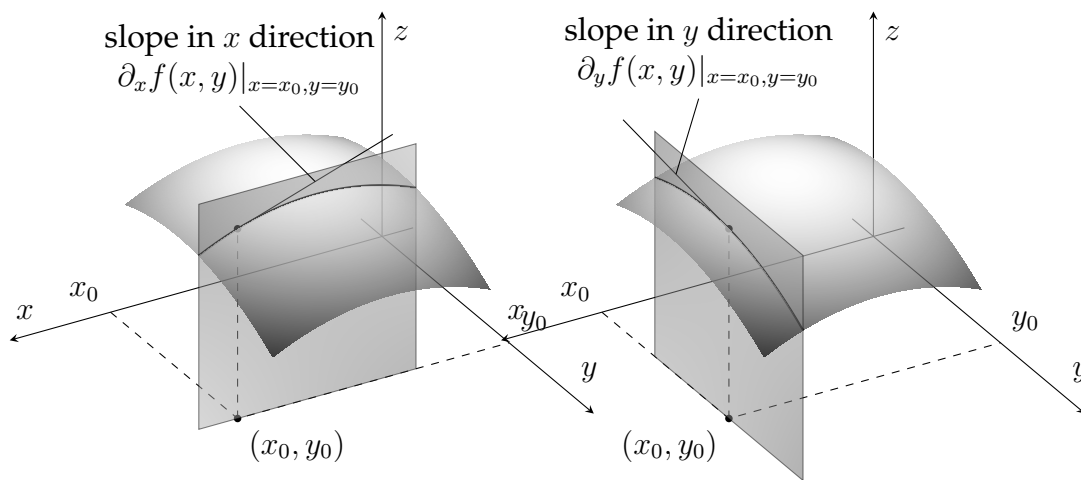


Figure 2: Cross sections of f in the x and y direction at (x_0, y_0) .

3. Drawing level curves of the function in its domain:

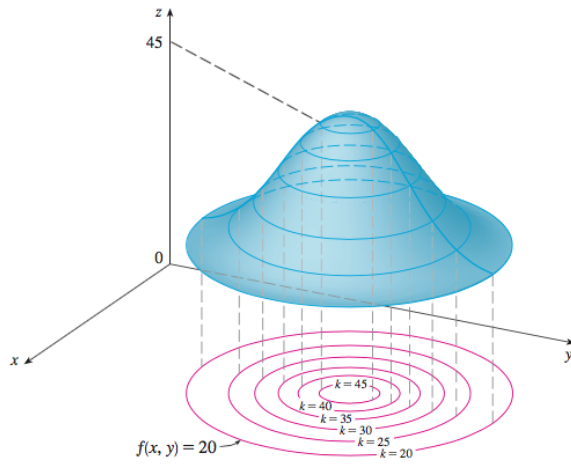


Figure 3: Some level curves of f .

4. combining all three ways in one picture:

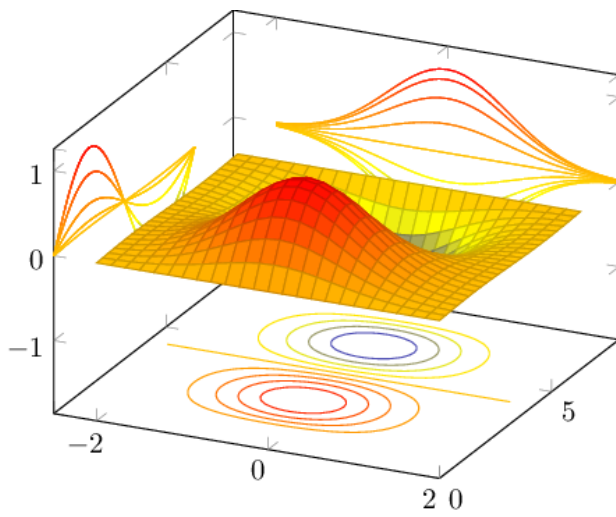


Figure 4: The graph of f together with some of its cross sections and level curves.

Overall goal: Linear approximation

For functions $f : \mathbb{R}^n \rightarrow \mathbb{R}$ variables, a good approximation does well in all n directions in the domain of f . Ideally we would end up with a formula that allows us to approximate f at $\hat{x} \in \mathbb{R}^n$ with a linear function. Recall that a linear function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ take the form of an inner product (multiplication by a row vector or a $(1 \times n)$ matrix. So we would like to have for some $\mathbf{a} \in \mathbb{R}^n$:

$$f(\hat{x} + \Delta\mathbf{x}) - f(\hat{x}) = \mathbf{a} \cdot \Delta\mathbf{x} + \text{h.o.t.} ,$$

whenever $\Delta\mathbf{x}$ is small.

If we can find such a \mathbf{a} , then computing the approximation in all possible directions is simply a matter of taking inner products of vectors and we are back to linear algebra as long as we stay close to the original point \hat{x} .

How to get there: Partial Derivatives

Recall that we denote unit vectors in the standard coordinate system of \mathbb{R}^n by e^i . Moving away from \hat{x} in the direction of the first coordinate axis is then denoted by moving to $\hat{x} + \Delta\mathbf{x}$ with $\Delta\mathbf{x} = he^1$ for some $h \in \mathbb{R}$. But this means that all the other components in \mathbf{x} remain fixed.

When we consider changes in the direction of a coordinate axis, we are really analyzing a function of a single real variable since all the other components in \hat{x} do not change. But this means that we can compute an approximation for f in the direction of e^i by computing the derivative:

$$\lim_{h \rightarrow 0} \frac{f(\hat{x} + he^i) - f(\hat{x})}{h} ,$$

exactly as we did before since we now have a function of a single variable x_i the x_j for $j \neq i$ are 'fixables'.

We call this limit the partial derivative of f at \hat{x} and denote it by:

$$D_{x_i} f(\hat{x}) := \frac{\partial f(\hat{x})}{\partial x_i} := \lim_{h \rightarrow 0} \frac{f(\hat{x} + he^i) - f(\hat{x})}{h} .$$

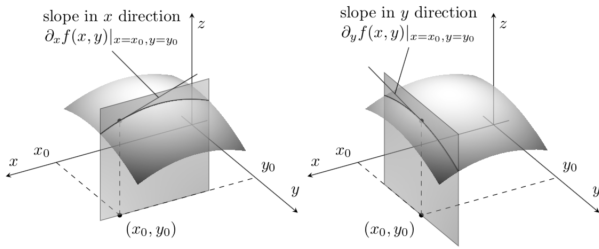


Figure 5: Partial derivatives of f at (x_0, y_0) .

But this is all we need if a linear approximation exists! A linear approximation in the direction $\Delta \mathbf{x} = \mathbf{e}^i$ must coincide with $\frac{\partial f(\hat{\mathbf{x}})}{\partial x_i}$. But each direction $\Delta \mathbf{x}$ can be written as:

$$\Delta \mathbf{x} = \sum_{i=1}^n x_i \mathbf{e}^i,$$

so by linearity we get for all $\Delta \mathbf{x} = h(x_1, \dots, x_n)$ that

$$f(\hat{\mathbf{x}} + \Delta \mathbf{x}) - f(\hat{\mathbf{x}}) = D_{\mathbf{x}} f(\hat{\mathbf{x}}) \cdot \Delta \mathbf{x} + \text{h.o.t.},$$

where $D_{\mathbf{x}} f(\hat{\mathbf{x}})$ is the row vector of partial derivatives

$$D_{\mathbf{x}} f(\hat{\mathbf{x}}) = \left(\frac{\partial f(\hat{\mathbf{x}})}{\partial x_1}, \dots, \frac{\partial f(\hat{\mathbf{x}})}{\partial x_n} \right).$$

Figure 6 shows a plane approximating a two-dimensional surface in three-dimensional space

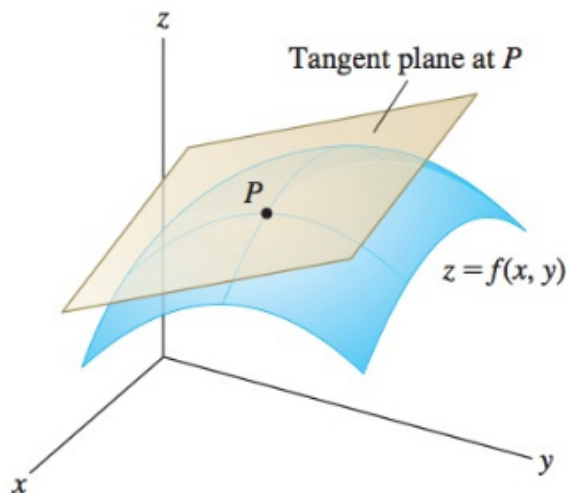


Figure 6: Linear approximation to f at point P .

Extra material: Existence of linear approximation

It can be shown that such a linear approximation exists if all the partial derivatives evaluated at \mathbf{x} are continuous functions of \mathbf{x} , i.e. the point at which they are evaluated.

The mere existence of partial derivatives is not enough to guarantee even continuity of the function. For an example, you can consider the function around $(x, y) = (0, 0)$:

$$f(x, y) = \begin{cases} \frac{xy}{x^2+y^2} & \text{if } (x, y) \neq (0, 0), \\ 0 & \text{if } (x, y) = (0, 0). \end{cases}$$

The requirement of continuous partial derivatives obviously fails in this case. We will not prove the intuitive result that with continuous partial derivatives, you get the existence of the linear approximation. In this case, we say that f is differentiable at $\hat{\mathbf{x}}$. The proof can be found in any book on advanced calculus.

Vector valued functions of multiple real variables

A function $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a vector of functions $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$:

$$\mathbf{f}(\mathbf{x}) = \begin{pmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{pmatrix}.$$

Since each of the component functions is a real valued function of n variables, we have from the previous subsection that a linear approximation of f_i at $\hat{\mathbf{x}}$ is given by the derivative $D_x f_i(\hat{\mathbf{x}})$. If all the partial derivatives of all component functions exist and are continuous at $\hat{\mathbf{x}}$, then the derivative of \mathbf{f} at $\hat{\mathbf{x}}$ is the $m \times n$ matrix:

$$D_x \mathbf{f}(\mathbf{x}) = \begin{pmatrix} \frac{\partial f_1(\hat{\mathbf{x}})}{\partial x_1} & \cdots & \frac{\partial f_1(\hat{\mathbf{x}})}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m(\hat{\mathbf{x}})}{\partial x_1} & \cdots & \frac{\partial f_m(\hat{\mathbf{x}})}{\partial x_n} \end{pmatrix}.$$

Since partial derivatives can be viewed as standard derivatives in a fixed direction, the rules for computing derivatives remain valid for multivariate functions. In particular, we have the chain rule for $h(\mathbf{x}) := f(g(\mathbf{x}))$, where $f : \mathbb{R}^k \rightarrow \mathbb{R}^m$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}^k$. Let $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_k) = (g_1(\hat{\mathbf{x}}), \dots, g_k(\hat{\mathbf{x}}))$:

$$D_x h(\mathbf{x}) = D_y f(\hat{\mathbf{y}}) D_x g(\hat{\mathbf{x}}).$$

Writing this matrix multiplication explicitly gives the i_j^{th} element of $D_x h(\mathbf{x})$ as:

$$\frac{\partial h_i(\hat{\mathbf{x}})}{\partial x_j} = \sum_{k=1}^k \frac{\partial f_i(\hat{\mathbf{y}})}{\partial y_k} \frac{\partial g_k(\hat{\mathbf{x}})}{\partial x_j}.$$

To recap: the derivative of \mathbf{f} at $\hat{\mathbf{x}}$ is the matrix of its partial derivatives evaluated at that point. This matrix (of fixed numbers) generates the linear function

$$\mathbf{f}(\hat{\mathbf{x}} + \Delta \mathbf{x}) - \mathbf{f}(\hat{\mathbf{x}}) = D_x \mathbf{f}(\hat{\mathbf{x}}) \Delta \mathbf{x} + h.o.t.$$