

Introduction to Mathematics for Political Science: Smooth Functions

3 September 2019

Linear functions are rarely sufficient to describe political phenomena of interest. Smooth functions are often better suited for this purpose. Smoothness is particularly important feature in optimization problems. You've encountered smooth functions and optimization in a 1-dimensional setting. Today, we generalize these intuitions to multi-dimensional spaces, preparing ourselves to approximate and optimize functions of many variables.

Smoothness in One Dimension (Univariate Calculus)

Consider a *policy production function*, $p : y \rightarrow \mathbb{R}_+$, which maps a legislator's staff size ($y \in \mathbb{R}_+$) to a number of bills produced during a session of congress.¹ Assume this function is continuous – small changes in staff size produce small changes in policy output. You learned what it means for this function to be *smooth* in your self-study of calculus. A function $f : \mathbb{R} \rightarrow \mathbb{R}$ was said to be *differentiable* at x if

$$\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

exists.² A function is said to be smooth if it is everywhere differentiable. Geometrically, a function is differentiable at x if we can draw a *line approximating* f that is tangent to $f(x)$ at x . The derivative at x is the slope of this line at x and is given by

$$\alpha = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

Rearranging gives

$$\begin{aligned} \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} - \alpha &= 0 \\ \lim_{h \rightarrow 0} \frac{f(x+h) - f(x) - \alpha h}{h} &= 0 \\ \lim_{h \rightarrow 0} \frac{f(x+h) - f(x) - g(x)}{h} &= 0 \end{aligned}$$

where $g(x) = \alpha h$ is a linear function approximating f in the neighborhood of x .³ The difference $f(x+h) - f(x) - g(x)$ is the *approximation error* that $g(x)$ seeks to minimize in the neighborhood of x

$$\epsilon(x, h) = f(x+h) - f(x) - g(x)$$

¹ We might parameterize this function as follows

$$p(y) = y^\beta$$

for $\beta > 0$

² Give an example of a continuous but not smooth function. How do these concepts differ from one another?

³ Note that the function $f(x) + g(x)$ is not linear but *affine*.

For the derivative to exist, we must have

$$\lim_{h \rightarrow 0} \epsilon(x, h) = 0 \quad (1)$$

Policy productivity is obviously more complicated than we postulated when writing $p(y)$. Consider the slightly more complex $p : \{x \times y\} \rightarrow \mathbb{R}_+$, where x is the legislator's policy expertise $x \in \mathbb{R}_+$ and y is again the number of staff members.⁴ What does it mean for this function to be smooth? What about a generic $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$? We want to generalize our one-dimensional notion of smoothness to these arbitrary spaces. We say a function is *smooth* if it *can be well approximated locally by a linear function*.

Smooth Functions and Differentiability

Formally, let $f : X \rightarrow Y$ be a function between normed linear spaces. We want to find a linear function $g(x)$ that well approximates f in the neighborhood of a particular element $x_0 \in X$ as we move away from x_0 in the direction of some x .

$$f(x_0 + x) \approx f(x_0) + g(x)$$

The approximation error in this setting is given by

$$\epsilon(x) = f(x_0 + x) - (f(x_0) + g(x))$$

We can scale the approximation error by its distance from x_0 , giving us a familiar quotient

$$\eta(x) = \frac{\epsilon(x)}{\|x\|} = \frac{f(x_0 + x) - (f(x_0) + g(x))}{\|x\|}$$

We want $\eta(x)$ to get small as x gets small. This presents a natural definition for differentiability.

Definition: A function $f : X \rightarrow Y$ is *differentiable* at $x_0 \in X$ if there exists a linear function $g : X \rightarrow Y$ such that for all $x \in X$,

$$f(x_0 + x) = f(x_0) + g(x) + \eta(x)\|x\|$$

and $x \rightarrow 0_X \implies \eta(x) \rightarrow 0_Y$.

Note: We'll call this function $g(x)$ the derivative of f at x_0 and denote it with $Df[x_0]$ or $f'[x_0]$.

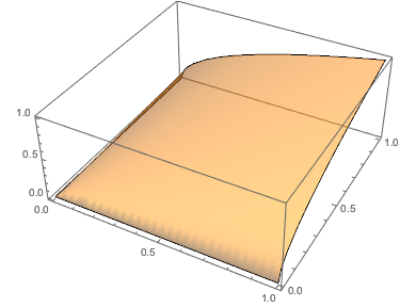
Definition: A function $f : X \rightarrow Y$ is differentiable if it is differentiable at all $x_0 \in X$. The derivative defines a function $Df : X \rightarrow Y$ that is a bounded linear function between normed linear spaces X and Y .

We call $Df[x_0]$ the derivative of f at x_0 and $Df[x_0](x)$ the derivative of f at x_0 in the direction of x .

⁴ We might parameterize this function with

$$p(x, y) = x^\alpha y^\beta$$

Such a function would look something like this



Note: If Df is a continuous function then we say f is *continuously differentiable* or C^1 . If Df itself is differentiable and its derivative is a continuous function then we say f is *twice continuously differentiable* or C^2 , and so on.

Proposition: If a function is differentiable at x_0 then it is continuous at x_0 .

Proof: Take a sequence $x^m \rightarrow x_0$. For a function f to be continuous at x_0 , we must have $x^m \rightarrow x_0 \implies f(x^m) \rightarrow f(x_0)$.⁵ Let $x^m = x_0 + x^n$ with $x^n \rightarrow 0_X \implies x^m \rightarrow x_0$. If f is differentiable, then, $x^n \rightarrow 0_X \implies \eta(x^n) \rightarrow 0_Y$ and

$$\begin{aligned} f(x_0 + x^n) - (f(x_0) + g(x^n)) &\rightarrow 0_Y \\ f(x^m) - (f(x_0) + g(x^n)) &\rightarrow 0_Y \\ f(x^m) &\rightarrow f(x_0) + g(x^n) \\ f(x^m) &\rightarrow f(x_0) \end{aligned}$$

Since $x^n \rightarrow 0_X$ and g is a linear function. This demonstrates $x^m \rightarrow x_0 \implies f(x^m) \rightarrow f(x_0)$ as desired.

Proposition (Chain Rule): Let X , Y , and Z , be normed linear spaces and $g : X \rightarrow Y$ and $f : Y \rightarrow Z$ be functions. If f and g are differentiable at x_0 then $f \circ g$ is also differentiable at x_0 with

$$D_{f \circ g}[x_0] = Df[g(x_0)]Dg[x_0]$$

Before proceeding to the proof, we state the following helper result without proof⁶

Proposition (Helper): If $f : X \rightarrow Y$ is a continuous linear function between normed vector spaces X and Y then there exists a C such that

$$\|f(x)\|_Y \leq C\|x\|_X$$

for all $x \in X$.

Proof: We need to show we can approximate $h = f \circ g$ as follows

$$f(g(x_0 + x)) = f(g(x_0)) + Df[g(x_0)]Dg[x_0](x) + \eta_Z(x)$$

with $x \rightarrow 0_X \implies \eta_Z(x) \rightarrow 0_Z$.

Since f and g are differentiable, we have $x \rightarrow 0_X \implies \epsilon_Y(x) \rightarrow 0_Y$ and $\epsilon_Z(y) \rightarrow 0_Z$ for $x \in X$ and $y \in Y$ with

$$\epsilon_Y(x) = g(x_0 + x) - g(x_0) - Dg[x_0](x)$$

and

$$\epsilon_Z(y) = f(g(x_0) + y) - f(g(x_0)) - Df[g(x_0)](y)$$

Then,

$$\begin{aligned} \eta_Z(x) &= f(g(x_0 + x)) - f(g(x_0)) - Df[g(x_0)]Dg[x_0](x) \\ &= f(g(x_0) + Dg[x_0](x) + \epsilon_Y(x)) - f(g(x_0)) - Df[g(x_0)]Dg[x_0](x) \end{aligned}$$

⁵ See Definition 3 in the notes on continuous functions.

⁶ See Ok, *Real Analysis with Economic Applications* p. 240.

Now let $\mathbf{y} = Dg[\mathbf{x}_0](\mathbf{x}) + \epsilon_Y(\mathbf{x})$ and apply the equation for $\epsilon_Z(\mathbf{y})$ above to see

$$\begin{aligned}\eta_Z(\mathbf{x}) &= f(g(\mathbf{x}_0)) + Df[g(\mathbf{x}_0)](Dg[\mathbf{x}_0](\mathbf{x}) + \epsilon_Y(\mathbf{x})) + \epsilon_Z(Dg[\mathbf{x}_0](\mathbf{x}) + \epsilon_Y(\mathbf{x})) - f(g(\mathbf{x}_0)) - Df[g(\mathbf{x}_0)](\mathbf{x})Dg[\mathbf{x}_0](\mathbf{x}) \\ &= Df[g(\mathbf{x}_0)](\epsilon_Y(\mathbf{x})) + \epsilon_Z(Dg[\mathbf{x}_0](\mathbf{x}) + \epsilon_Y(\mathbf{x}))\end{aligned}$$

We now just need to show $\mathbf{x} \rightarrow 0$ implies

$$\frac{Df[g(\mathbf{x}_0)](\epsilon_Y(\mathbf{x})) + \epsilon_Z(Dg[\mathbf{x}_0](\mathbf{x}) + \epsilon_Y(\mathbf{x}))}{\|\mathbf{x}\|} \rightarrow 0$$

Breaking this into pieces, we need first

$$\frac{Df[g(\mathbf{x}_0)](\epsilon_Y(\mathbf{x}))}{\|\mathbf{x}\|} \rightarrow \mathbf{0}_Y$$

as $\mathbf{x} \rightarrow \mathbf{0}_X$. Note that $Df[g(\mathbf{x}_0)]$ is a continuous linear function, so our helper Theorem applies:

$$\frac{Df[g(\mathbf{x}_0)](\epsilon_Y(\mathbf{x}))}{\|\mathbf{x}\|} \leq C \left\| \frac{\epsilon_Y(\mathbf{x})}{\|\mathbf{x}\|} \right\|$$

$\mathbf{x} \rightarrow \mathbf{0}_X \implies \left\| \frac{\epsilon_Y(\mathbf{x})}{\|\mathbf{x}\|} \right\| \rightarrow \mathbf{0}_Y$ by the differentiability of g so we have $\frac{Df[g(\mathbf{x}_0)](\epsilon_Y(\mathbf{x}))}{\|\mathbf{x}\|} \rightarrow \mathbf{0}_Y$ as well.⁷

Now we show $\mathbf{x} \rightarrow \mathbf{0}_X$ implies

$$\frac{\epsilon_Z(Dg[\mathbf{x}_0](\mathbf{x}) + \epsilon_Y(\mathbf{x}))}{\|\mathbf{x}\|} \rightarrow \mathbf{0}_Z$$

Let $\mathbf{y} = Dg[\mathbf{x}_0](\mathbf{x}) + \epsilon_Y(\mathbf{x})$. Since $Dg[\mathbf{x}_0](\mathbf{x})$, we can find a C such that

$$C\|\mathbf{y}\| \leq \|\mathbf{x}\|$$

which implies

$$\frac{\epsilon_Z(Dg[\mathbf{x}_0](\mathbf{x}) + \epsilon_Y(\mathbf{x}))}{\|\mathbf{x}\|} \leq \frac{1}{C} \left\| \frac{\epsilon_Z(\mathbf{y})}{\|\mathbf{y}\|} \right\|$$

where $\left\| \frac{\epsilon_Z(\mathbf{y})}{\|\mathbf{y}\|} \right\| \rightarrow \mathbf{0}_Z$ as $\mathbf{y} \rightarrow \mathbf{0}_Y$. This completes the proof. ■

Note the analogue here with linear regression. Estimation of a regression coefficient entails searching for a linear function that minimizes the distance itself between data (\mathbf{x}) . Formally, outcomes y are modeled as a function of the linear approximation, an error term, and an intercept

$$\begin{aligned}y &= f(\mathbf{x}_0) + g(\mathbf{x}) + \epsilon \\ &= \alpha + \mathbf{x}^T \boldsymbol{\beta} + \epsilon\end{aligned}$$

Derivatives help us construct affine approximations of smooth functions. The approximate change in a function f in the neighborhood of \mathbf{x}_0 is

$$df = f(\mathbf{x}_0 + d\mathbf{x}) - f(\mathbf{x}_0) \approx Df[\mathbf{x}_0](d\mathbf{x}) \approx f(\mathbf{x}_0) + Df[\mathbf{x}_0](d\mathbf{x})$$

⁷ Recall the exercise that proves this fact from the lecture Normed Linear Spaces.

Partial Derivatives, the Gradient Vector, and the Hessian Matrix

Considering the set of *functionals* allows us to consider how the value of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ changes as a function of a single variable, holding all other variables constant. Let

$$h_i(t) = f(x_1^0, \dots, t, \dots, x_n^0)$$

denote the value of the i th partial of f at x^0 .

Definition: The i th partial derivative of a differentiable function f at x^0 is

$$\lim_{t \rightarrow 0} \frac{h(x_i^0 + t) - h(x_i^0)}{t}$$

We denote these partial derivatives with

$$\frac{\partial f[x^0]}{\partial x_i}$$

Geometrically, we're asking for the slope of the function when sliced at the i th cross section of f through x^0 .⁸ Why should we care about partial derivatives?

Example: Consider the policy production function again, $p(y, x)$. Suppose the candidates policy expertise is fixed. What is the marginal effect of adding staff members on the policy output? This is precisely what the partial derivative gives us.

Iteratively taking partial derivatives of a function f would tell us how the function is changing in every direction. We could collect these partial derivatives into a vector, called the *gradient*, which summarizes all of these changes.⁹

Definition: The *gradient* of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at x^0 is

$$\nabla f(x^0) = \left(\frac{\partial f[x^0]}{\partial x_1}, \dots, \frac{\partial f[x^0]}{\partial x_n} \right)$$

and

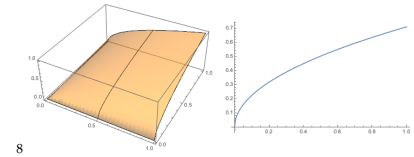
$$Df[x^0](x) = \sum_{i=1}^n \frac{\partial f[x^0]}{\partial x_i} x_i^0$$

We see that for functionals on \mathbb{R}^n , we can decompose the derivative into a sum of partial derivatives.¹⁰ The gradient is a vector, and it has interesting geometric properties. In particular, it points in the direction of greatest change.¹¹ To get a better understanding of this geometry, it's useful to define contours.

Definition: The contour of a functional f through $c = f(x^0)$ is¹²

$$f^{-1}(c) = \{x \in X | f(x) = c\}$$

Think about an "election function" $p(x, y)$ that returns a probability of winning an election as a function of a policy choice x and

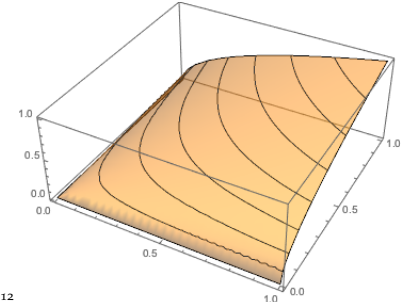


8

⁹ The ∇ you see in the notes can be written backslash "nabla" in L^AT_EX

¹⁰ Note that the gradient can be expressed as the inner product $\nabla f(x^0)^T x$

¹¹ If you were a campaign strategist advising a candidate facing the election function $p(y, x)$, knowing the gradient would be useful for improving your chances of winning.



12

campaign advertisement spending y . The concept of the contour is simple in this setting. $p^{-1}(c)$ gives us the set of $\{x \times y\}$ that give the candidate a c chance of victory. There are likely to be many potential $x \times y$ combinations of campaign spending and policy that could produce this result. Contours are sometimes useful for visualizing high-dimensional functions.

Remark: The gradient is orthogonal to the contour¹³

$$\nabla f(\mathbf{x}^0)^T f^{-1}(c) = 0$$

If a function is twice differentiable, we can compute its *second-order partial derivatives* by differentiating the partial derivatives with respect to another variable. We write these

$$\frac{\partial^2 f[\mathbf{x}^0]}{\partial x_i \partial x_j}$$

When $i \neq j$, we call this the *cross partial* of f at \mathbf{x}^0 . If $i = j$ we write

$$\frac{\partial^2 f[\mathbf{x}^0]}{\partial x_i^2}$$

Example: The cross partial of the election function would be written

$$\frac{\partial p(y, x)}{\partial y \partial x}$$

and describes how the effectiveness of campaign spending changes with respect to the candidate's policy position (or vice versa).

Notice that there are many cross partials we can take as n gets large.¹⁴ The gradient vector can be differentiated with respect to each of its elements. We store these derivatives in the *Hessian* matrix.

Definition: The *Hessian* of a function f at \mathbf{x}^0 is a matrix of second-order derivatives, with

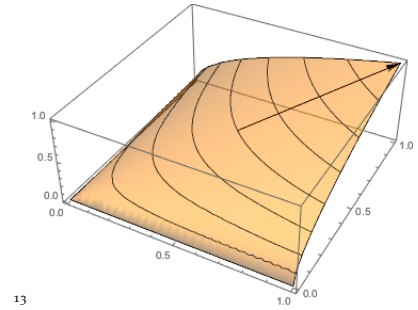
$$H_f(\mathbf{x}^0) = \begin{bmatrix} \frac{\partial^2 f[\mathbf{x}^0]}{\partial x_1^2} & \dots & \frac{\partial^2 f[\mathbf{x}^0]}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f[\mathbf{x}^0]}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f[\mathbf{x}^0]}{\partial x_n^2} \end{bmatrix}$$

We can describe the concavity or convexity of a high-dimensional function using its Hessian.

Proposition: A twice-differentiable function f is strictly locally convex at \mathbf{x} iff $H_f(\mathbf{x})$ is positive definite. The function is strictly locally concave at \mathbf{x} iff $H_f(\mathbf{x})$ is negative definite.

The Jacobian

We often represent political systems as systems of equations. In a multi-candidate election, for example, we might have a family of



¹⁴ n^2 to be precise.

policy production $p_i(\mathbf{y}, \mathbf{x})$ that map *every* legislator's expertise and staff size into policy output.¹⁵ The *Jacobian* is a matrix that stores a linear approximation of this system.

More formally, consider a system of m functionals, with each $f_m : \mathbb{R}^n \rightarrow \mathbb{R}$.¹⁶

Definition: The *Jacobian* of a system of functionals f is a matrix of partial derivatives, where each element ij is the i th partial of the j th function evaluated at \mathbf{x}^0

$$J_f(\mathbf{x}^0) = \begin{bmatrix} \frac{\partial f_1[\mathbf{x}^0]}{\partial x_1} & \dots & \frac{\partial f_1[\mathbf{x}^0]}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m[\mathbf{x}^0]}{\partial x_1} & \dots & \frac{\partial f_m[\mathbf{x}^0]}{\partial x_n} \end{bmatrix}$$

Notice that the Jacobian can also be written as a vector of gradients stacked on top of one another

$$J_f(\mathbf{x}^0) = \begin{pmatrix} \nabla f_1(\mathbf{x}^0) \\ \vdots \\ \nabla f_m(\mathbf{x}^0) \end{pmatrix}$$

More generally the derivative of any function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ can be expressed by the Jacobian where each element is the derivative of the m th output with respect to the n th input.¹⁷

Example: Calculate the Jacobian of the function

$$f(x, y) = \begin{pmatrix} x^2 + y \\ 2xy \end{pmatrix}$$

The Jacobian is

$$J_f(x, y) = \begin{pmatrix} 2x & 1 \\ 2y & 2x \end{pmatrix}$$

References

1. Carter, Michael. *Foundations of Mathematical Economics*. Chapter 4.
2. Ok, Efe. *Real Analysis with Economic Applications*
3. Lipshitz, Robert. "Linear Maps, the Total Derivative, and the Chain Rule"

¹⁵ Notice that the arguments of the function p are now vectors.

¹⁶ Note that the system $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, because each of the m functionals returns a single real number.

¹⁷ Remember any linear function mapping \mathbb{R}^n to \mathbb{R}^m can be encoded as an $m \times n$ matrix $A_{m \times n}$ with $f(\mathbf{x}) = A\mathbf{x}$