

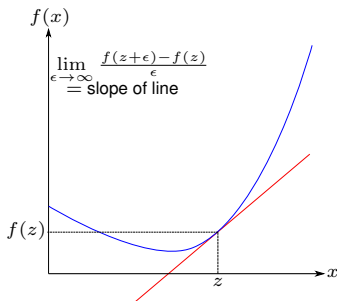
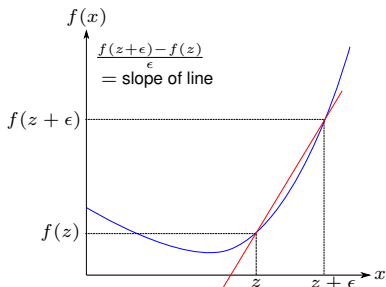
3. Derivatives

- scalar derivatives
- gradient and hessian
- multi-variable differentiation rules

Derivative definition

the *derivative* of $f : \mathbb{R} \rightarrow \mathbb{R}$ at the number z is defined as

$$f'(z) = \frac{df}{dz} = \lim_{\epsilon \rightarrow 0} \frac{f(z + \epsilon) - f(z)}{\epsilon}$$



- when $f'(x)$ is positive, the function $f(x)$ increases as x does
- $f'(x)$ is negative, $f(x)$ decreases as x increases

Common derivatives

$f(x)$	$f'(x)$
c	0
x^ℓ	$\ell x^{\ell-1}$
e^x	e^x
$\log(x), x > 0$	$\frac{1}{x}$
$\log_c(x), x > 0, c > 0$	$\frac{1}{x \ln(c)}$
$\sin(x)$	$\cos(x)$
$\cos(x)$	$-\sin(x)$

(we use $\log(\cdot) = \ln(\cdot)$ to denote the natural logarithm)

Derivative rules

- **Linearity:** for $f(x) = \alpha g(x) + \beta h(x)$:

$$f'(x) = \alpha g'(x) + \beta h'(x)$$

- **Product rule:** for $f(x) = g(x)h(x)$:

$$f'(x) = g'(x)h(x) + g(x)h'(x)$$

- **Quotient rule:** for $f(x) = \frac{g(x)}{h(x)}$:

$$f'(x) = \frac{g'(x)h(x) - g(x)h'(x)}{h(x)^2}$$

- **Chain rule:** for $f(x) = g(h(x))$ where $g : \mathbb{R} \rightarrow \mathbb{R}$ and $h : \mathbb{R} \rightarrow \mathbb{R}$:

$$f'(x) = h'(x)g'(h(x))$$

Second derivative

the *second derivative* of f at a point z is the derivative of the first derivative:

$$f''(z) = \frac{d^2 f}{dz^2} = \lim_{\epsilon \rightarrow 0} \frac{f'(z + \epsilon) - f'(z)}{\epsilon},$$

the second derivative conveys information about the curvature of the function

- when $f''(x) > 0$, then $f'(x)$ is increasing, which suggests the slope of the tangent line to f increases as x does yielding a concave-upwards shape
- if $f''(x)$ is negative, the function exhibits a concave-downwards curvature

Outline

- scalar derivatives
- **gradient and hessian**
- multi-variable differentiation rules

Gradient

the *gradient* of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ (at \mathbf{z}) is

$$\nabla f(\mathbf{z}) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(\mathbf{z}) \\ \frac{\partial f}{\partial x_2}(\mathbf{z}) \\ \vdots \\ \frac{\partial f}{\partial x_n}(\mathbf{z}) \end{bmatrix}$$

- entries $\frac{\partial f}{\partial x_i}(\mathbf{z})$ are *partial derivative* of f at point \mathbf{z} , with respect to x_i :

$$\frac{\partial f}{\partial x_i}(\mathbf{z}) = \lim_{\epsilon \rightarrow 0} \frac{f(z_1, \dots, z_{i-1}, z_i + \epsilon, z_{i+1}, \dots, z_n) - f(\mathbf{z})}{\epsilon}$$

- the gradient $\nabla f(\mathbf{z})$ is a vector that is pointing to the direction where f increases the fastest at \mathbf{z}

Example 3.1

- a) partial derivatives of $f(\mathbf{x}) = \|\mathbf{x}\|^2 = x_1^2 + \cdots + x_n^2$ are $\frac{\partial f}{\partial x_i}(\mathbf{x}) = 2x_i$;
hence

$$\nabla f(\mathbf{x}) = (2x_1, \dots, 2x_n) = 2\mathbf{x}$$

- b) gradient of the function $f(\mathbf{x}) = 5x_1 + 8x_2 + x_1x_2 - x_1^2 - 2x_2^2$ is

$$\nabla f(\mathbf{x}) = (5 + x_2 - 2x_1, 8 + x_1 - 4x_2)$$

- c) gradient of $f(x) = x_1^2 + e^{-x_1} + \sin(x_2)$ is

$$\nabla f(\mathbf{x}) = \begin{bmatrix} 2x_1 - e^{-x_1} \\ \cos(x_2) \end{bmatrix}$$

Hessian

the *Hessian* of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at \mathbf{z} is defined as

$$\nabla^2 f(\mathbf{z}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(\mathbf{z}) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(\mathbf{z}) & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(\mathbf{z}) \\ \frac{\partial^2 f}{\partial x_2 \partial x_1}(\mathbf{z}) & \frac{\partial^2 f}{\partial x_2^2}(\mathbf{z}) & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n}(\mathbf{z}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1}(\mathbf{z}) & \frac{\partial^2 f}{\partial x_n \partial x_2}(\mathbf{z}) & \cdots & \frac{\partial^2 f}{\partial x_n^2}(\mathbf{z}) \end{bmatrix}$$

- we say that f is *twice differentiable* if $\nabla^2 f(\mathbf{x})$ exists for all $\mathbf{x} \in \mathbb{R}^n$
- the Hessian is a symmetric matrix $\nabla^2 f(\mathbf{z}) = \nabla^2 f(\mathbf{z})^T$ since

$$\frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{z}) = \frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{z}), \quad \text{for all } i, j = 1, \dots, n$$

Example 3.2

a) for $f(\mathbf{x}) = 5x_1 + 8x_2 + x_1x_2 - x_1^2 - 2x_2^2$:

$$\nabla f(\mathbf{x}) = \begin{bmatrix} 5 + x_2 - 2x_1 \\ 8 + x_1 - 4x_2 \end{bmatrix}, \quad \nabla^2 f(\mathbf{x}) = \begin{bmatrix} -2 & 1 \\ 1 & -4 \end{bmatrix}$$

b) for

$$f(\mathbf{x}) = e^{x_1+x_2-1} + e^{x_1-x_2-1} + e^{-x_1-1}$$

the gradient is

$$\nabla f(\mathbf{x}) = \begin{bmatrix} e^{x_1+x_2-1} + e^{x_1-x_2-1} - e^{-x_1-1} \\ e^{x_1+x_2-1} - e^{x_1-x_2-1} \end{bmatrix}$$

and the Hessian is

$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} e^{x_1+x_2-1} + e^{x_1-x_2-1} + e^{-x_1-1} & e^{x_1+x_2-1} - e^{x_1-x_2-1} \\ e^{x_1+x_2-1} - e^{x_1-x_2-1} & e^{x_1+x_2-1} + e^{x_1-x_2-1} \end{bmatrix}$$

Linear and quadratic functions

Linear and affine functions: for $f(\mathbf{x}) = \mathbf{a}^T \mathbf{x} + b$:

$$\nabla f(\mathbf{x}) = \mathbf{a}$$

$$\nabla^2 f(\mathbf{x}) = 0$$

Quadratic functions: for $f(\mathbf{x}) = \mathbf{x}^T Q \mathbf{x} + \mathbf{r}^T \mathbf{x} + c$, where $Q = Q^T$ is symmetric:

$$\nabla f(\mathbf{x}) = 2Q\mathbf{x} + \mathbf{r}$$

$$\nabla^2 f(\mathbf{x}) = 2Q$$

Least-squares function

the *least-squares function* $f(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|^2$ can be expressed as

$$\begin{aligned}f(\mathbf{x}) &= \|\mathbf{Ax} - \mathbf{b}\|^2 \\&= (\mathbf{Ax} - \mathbf{b})^T(\mathbf{Ax} - \mathbf{b}) \\&= (\mathbf{x}^T\mathbf{A}^T - \mathbf{b}^T)(\mathbf{Ax} - \mathbf{b}) \\&= \mathbf{x}^T\mathbf{A}^T\mathbf{Ax} - \mathbf{b}^T\mathbf{Ax} - \mathbf{x}^T\mathbf{A}^T\mathbf{b} + \mathbf{b}^T\mathbf{b} \\&= \mathbf{x}^T\mathbf{A}^T\mathbf{Ax} - 2\mathbf{b}^T\mathbf{Ax} + \mathbf{b}^T\mathbf{b}\end{aligned}$$

this means that f is quadratic $f(\mathbf{x}) = \mathbf{x}^T\mathbf{Q}\mathbf{x} + \mathbf{r}^T\mathbf{x} + c$ with

$$\mathbf{Q} = \mathbf{A}^T\mathbf{A}, \quad \mathbf{r}^T = -2\mathbf{b}^T\mathbf{A}, \quad c = \mathbf{b}^T\mathbf{b}$$

hence,

$$\nabla f(\mathbf{x}) = 2\mathbf{A}^T\mathbf{Ax} - 2\mathbf{A}^T\mathbf{b}, \quad \nabla^2 f(\mathbf{x}) = 2\mathbf{A}^T\mathbf{A}$$

Outline

- scalar derivatives
- gradient and hessian
- **multi-variable differentiation rules**

Composition with affine function

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

- $f : \mathbb{R}^n \rightarrow \mathbb{R}, g : \mathbb{R}^m \rightarrow \mathbb{R}$
- A is an $m \times n$ matrix
- \mathbf{b} is an m vector

the gradient and Hessian are

$$\nabla f(\mathbf{x}) = A^T \nabla g(A\mathbf{x} + \mathbf{b})$$

and

$$\nabla^2 f(\mathbf{x}) = A^T \nabla^2 g(A\mathbf{x} + \mathbf{b}) A$$

Example 3.3

use the composition with affine function property to find the gradient and Hessian of

$$f(\mathbf{x}) = e^{x_1+x_2-1} + e^{x_1-x_2-1} + e^{-x_1-1}$$

we can express f as $f(\mathbf{x}) = g(A\mathbf{x} + \mathbf{b})$, where $g(\mathbf{y}) = e^{y_1} + e^{y_2} + e^{y_3}$, and

$$A = \begin{bmatrix} 1 & 1 \\ 1 & -1 \\ -1 & 0 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix}$$

the gradient and Hessian of g are

$$\nabla g(\mathbf{y}) = \begin{bmatrix} e^{y_1} \\ e^{y_2} \\ e^{y_3} \end{bmatrix}, \quad \nabla^2 g(\mathbf{y}) = \begin{bmatrix} e^{y_1} & 0 & 0 \\ 0 & e^{y_2} & 0 \\ 0 & 0 & e^{y_3} \end{bmatrix}$$

hence

$$\begin{aligned}\nabla f(\mathbf{x}) &= A^T \nabla g(A\mathbf{x} + \mathbf{b}) = \begin{bmatrix} 1 & 1 & -1 \\ 1 & -1 & 0 \end{bmatrix} \begin{bmatrix} e^{x_1+x_2-1} \\ e^{x_1-x_2-1} \\ e^{-x_1-1} \end{bmatrix} \\ &= \begin{bmatrix} e^{x_1+x_2-1} + e^{x_1-x_2-1} - e^{-x_1-1} \\ e^{x_1+x_2-1} - e^{x_1-x_2-1} \end{bmatrix}\end{aligned}$$

and

$$\begin{aligned}\nabla^2 f(\mathbf{x}) &= A^T \nabla^2 g(A\mathbf{x} + \mathbf{b}) A \\ &= \begin{bmatrix} 1 & 1 & -1 \\ 1 & -1 & 0 \end{bmatrix} \begin{bmatrix} e^{x_1+x_2-1} & 0 & 0 \\ 0 & e^{x_1-x_2-1} & 0 \\ 0 & 0 & e^{-x_1-1} \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & -1 \\ -1 & 0 \end{bmatrix} \\ &= \begin{bmatrix} e^{x_1+x_2-1} + e^{x_1-x_2-1} + e^{-x_1-1} & e^{x_1+x_2-1} - e^{x_1-x_2-1} \\ e^{x_1+x_2-1} - e^{x_1-x_2-1} & e^{x_1+x_2-1} + e^{x_1-x_2-1} \end{bmatrix}\end{aligned}$$

Jacobian

let $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$:

$$f(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} f_1(x_1, \dots, x_n) \\ \vdots \\ f_m(x_1, \dots, x_n) \end{bmatrix}$$

where f_i is a scalar-valued function of \mathbf{x}

the **Jacobian** or **derivative matrix** of f at \mathbf{z} is the $m \times n$ matrix:

$$Df(\mathbf{z}) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(\mathbf{z}) & \frac{\partial f_1}{\partial x_2}(\mathbf{z}) & \cdots & \frac{\partial f_1}{\partial x_n}(\mathbf{z}) \\ \frac{\partial f_2}{\partial x_1}(\mathbf{z}) & \frac{\partial f_2}{\partial x_2}(\mathbf{z}) & \cdots & \frac{\partial f_2}{\partial x_n}(\mathbf{z}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1}(\mathbf{z}) & \frac{\partial f_m}{\partial x_2}(\mathbf{z}) & \cdots & \frac{\partial f_m}{\partial x_n}(\mathbf{z}) \end{bmatrix} = \begin{bmatrix} \nabla f_1(\mathbf{z})^T \\ \nabla f_2(\mathbf{z})^T \\ \vdots \\ \nabla f_m(\mathbf{z})^T \end{bmatrix}$$

- if $m = 1$, then $Df(\mathbf{z}) = \nabla f(\mathbf{z})^T$
- Jacobian of the gradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is its Hessian

Example 3.4

a) the Jacobian of

$$f(\mathbf{x}) = \begin{bmatrix} x_1 + x_2^2 \\ -x_1 + x_1x_2 \end{bmatrix}$$

is

$$Df(\mathbf{x}) = \begin{bmatrix} 1 & 2x_2 \\ -1 + x_2 & x_1 \end{bmatrix}$$

b) the derivative matrix or Jacobian of $f(\mathbf{x}) = A\mathbf{x}$ is

$$Df(\mathbf{x}) = A$$

Rules

Sum of two functions: if $f(\mathbf{x}) = f_1(\mathbf{x}) + f_2(\mathbf{x})$, then

$$\nabla f(\mathbf{x}) = \nabla f_1(\mathbf{x}) + \nabla f_2(\mathbf{x}), \quad \nabla^2 f(\mathbf{x}) = \nabla^2 f_1(\mathbf{x}) + \nabla^2 f_2(\mathbf{x})$$

Scalar multiplication: if $f(\mathbf{x}) = \alpha g(\mathbf{x})$, where α is a scalar, then

$$\nabla f(\mathbf{x}) = \alpha \nabla g(\mathbf{x}), \quad \nabla^2 f(\mathbf{x}) = \alpha \nabla^2 g(\mathbf{x})$$

Multivariable product rule: let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be

$$f(\mathbf{x}) = g(\mathbf{x})^T h(\mathbf{x}),$$

where $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $h : \mathbb{R}^n \rightarrow \mathbb{R}^m$, then

$$\nabla f(\mathbf{x}) = Df(\mathbf{x})^T = Dg(\mathbf{x})^T h(\mathbf{x}) + Dh(\mathbf{x})^T g(\mathbf{x})$$

Multivariable chain rule: let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be the composition

$$f(\mathbf{x}) = g(h(\mathbf{x})) = g(h_1(\mathbf{x}), \dots, h_p(\mathbf{x}))$$

where $g : \mathbb{R}^p \rightarrow \mathbb{R}$ and $h : \mathbb{R}^n \rightarrow \mathbb{R}^p$ are differentiable functions

- using the chain rule, the partial derivatives of f are

$$\frac{\partial f}{\partial x_j}(\mathbf{x}) = \frac{\partial h_1}{\partial x_j}(\mathbf{x}) \frac{\partial g}{\partial y_1}(h(\mathbf{x})) + \dots + \frac{\partial h_p}{\partial x_j}(\mathbf{x}) \frac{\partial g}{\partial y_p}(h(\mathbf{x}))$$

for $j = 1, \dots, n$

- the gradient can be compactly represented as the vector-matrix product:

$$\nabla f(\mathbf{x}) = Df(\mathbf{x})^T = Dh(\mathbf{x})^T \nabla g(h(\mathbf{x}))$$

Example 3.5

a) use the chain-rule to find the gradient of

$$f(\mathbf{x}) = (\sin(x_1) + x_2^2)^2 + (\sin(x_1) + x_2^2)(x_1 + x_2)^2$$

- we can write f as $f(\mathbf{x}) = g(h(\mathbf{x}))$ where

$$g(\mathbf{y}) = y_1^2 + y_1 y_2^2, \quad h(\mathbf{x}) = \begin{bmatrix} \sin(x_1) + x_2^2 \\ x_1 + x_2 \end{bmatrix}$$

- the gradient of g is $\nabla g(\mathbf{y}) = \begin{bmatrix} 2y_1 + y_2^2 \\ 2y_1 y_2 \end{bmatrix}$ and the derivative of h is

$$Dh(\mathbf{x}) = \begin{bmatrix} \cos(x_1) & 2x_2 \\ 1 & 1 \end{bmatrix}$$

- hence,

$$\begin{aligned} \nabla f(\mathbf{x}) &= Dh(\mathbf{x})^T \nabla g(h(\mathbf{x})) \\ &= \begin{bmatrix} \cos(x_1) & 1 \\ 2x_2 & 1 \end{bmatrix}^T \begin{bmatrix} 2\sin(x_1) + 2x_2^2 + (x_1 + x_2)^2 \\ 2(\sin(x_1) + x_2^2)(x_1 + x_2) \end{bmatrix} \end{aligned}$$

b) *nonlinear least-squares function:*

$$f(\mathbf{x}) = \|\mathbf{h}(\mathbf{x})\|^2 = \sum_{j=1}^p h_j(\mathbf{x})^2$$

we have $f(\mathbf{x}) = g(\mathbf{h}(\mathbf{x}))$ where $g(\mathbf{y}) = \|\mathbf{y}\|^2$

using $\nabla g(\mathbf{y}) = 2\mathbf{y}$ and the chain rule, we get

$$\nabla f(\mathbf{x}) = D\mathbf{h}(\mathbf{x})^T \nabla g(\mathbf{h}(\mathbf{x})) = 2D\mathbf{h}(\mathbf{x})^T \mathbf{h}(\mathbf{x})$$

References and further readings

- L. Vandenberghe. *EE133A Lecture Notes*, UCLA.
(<http://www.seas.ucla.edu/~vandenbe/ee133a.html>)
- Edwin KP Chong and Stanislaw H Zak. *An Introduction to Optimization*, John Wiley & Sons, 2013, chapters 2,3.
- Stephen Boyd and Lieven Vandenberghe. *Introduction to Applied Linear Algebra: Vectors, Matrices, and Least Squares*, Cambridge University Press, 2018, chapters 3,5,8.
- Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*, Cambridge University Press, 2004, appendices A.1, A.5, C.5.