

Math. review

Real vectors and matrices (§2.1)

- \mathbb{R} : set of real numbers
- \mathbb{R}^n : set of real column vectors

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \quad x_i \in \mathbb{R}$$

- $\mathbb{R}^{m \times n}$: set of $m \times n$ real matrices

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$

- Treat $\mathbb{R}^{n \times 1}$ and \mathbb{R}^n as equivalent
- \mathbf{A}^T : transpose of \mathbf{A}

$$[1, 2]^T = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

Functions

- Function $f : X \rightarrow Y$
- f takes values in X and gives values in Y
 - f is Y -valued
 - $f(x)$ is the value of f at x , where $x \in X$
- Example: $f : \mathbb{R}^3 \rightarrow \mathbb{R}$

$$f(\mathbf{x}) = \frac{x_1^3 + 3 \log(x_2 x_3)}{x_2}$$

Linear independence (§2.1)

- A set of vectors $\{\mathbf{a}_1, \dots, \mathbf{a}_k\}$ is said to be *linearly independent* if the equality

$$\alpha_1 \mathbf{a}_1 + \alpha_2 \mathbf{a}_2 + \cdots + \alpha_k \mathbf{a}_k = \mathbf{0}$$

implies that all the scalar coefficients $\alpha_i, i = 1, \dots, k$, are equal to zero.

- A set of vectors is linearly dependent if and only if one of the vectors from the set is a linear combination of the remaining vectors.

Rank of a matrix (§2.2)

- The maximal number of linearly independent columns of \mathbf{A} is called the *rank* of the matrix \mathbf{A} , denoted $\text{rank } \mathbf{A}$.
- A square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is *nonsingular* or *invertible* if $\text{rank } \mathbf{A} = n$ (full rank).
- A matrix is nonsingular if and only if its determinant is nonzero.

Inner product and norm (§2.4)

- Given $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$.
- Define the *inner product* of \mathbf{x} and \mathbf{y} :

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = \sum_{i=1}^n x_i y_i$$

- What are the properties of inner product?
- Define the *norm* of \mathbf{x} :

$$\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sqrt{\mathbf{x}^T \mathbf{x}} = \sqrt{\sum_{i=1}^n x_i^2}.$$

- What are the properties of norm?
- \mathbf{x} and \mathbf{y} are *orthogonal* if $\langle \mathbf{x}, \mathbf{y} \rangle = 0$.

Eigenvalues and eigenvectors (§3.2)

- Let \mathbf{A} be an $n \times n$ square matrix.
- A scalar λ (possibly complex) and a nonzero vector \mathbf{v} satisfying the equation $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$ are said to be, respectively, an *eigenvalue* and *eigenvector* of \mathbf{A} .
- λ is an eigenvalue of \mathbf{A} if and only if $\lambda\mathbf{I} - \mathbf{A}$ is singular (i.e., $\det[\lambda\mathbf{I} - \mathbf{A}] = 0$).
- $\det[\lambda\mathbf{I} - \mathbf{A}]$ is called the *characteristic polynomial* of \mathbf{A} .
- What are the zeros of the characteristic polynomial of \mathbf{A} ?

Symmetric matrices (§3.4)

- Q is symmetric if $Q = Q^T$.
- A symmetric matrix Q is said to be *positive definite* if $x^T Q x > 0$ for all nonzero vectors x .
- It is *positive semidefinite* if $x^T Q x \geq 0$ for all x .
- Similarly define *negative definite* and *negative semidefinite*.
- How is definiteness related to eigenvalues?

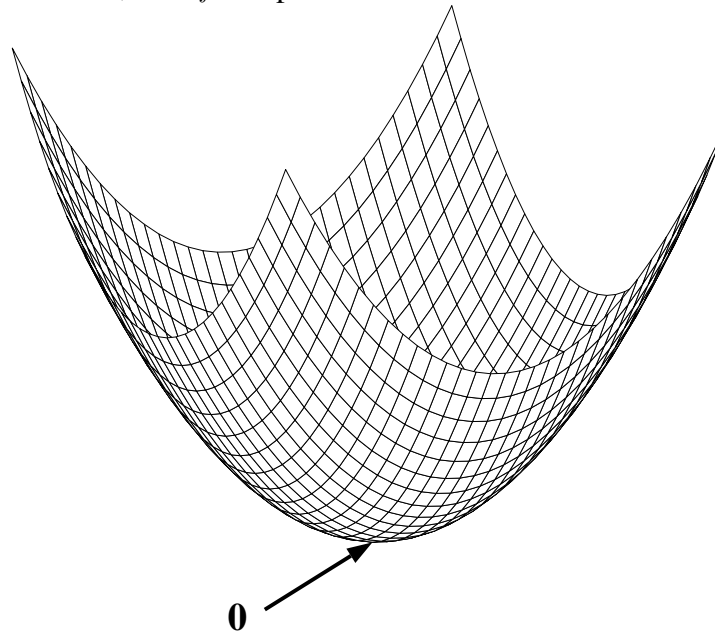
Quadratic functions (§3.4)

- $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a quadratic function if

$$f(x) = x^T Q x + b^T x + c,$$

where Q is symmetric.

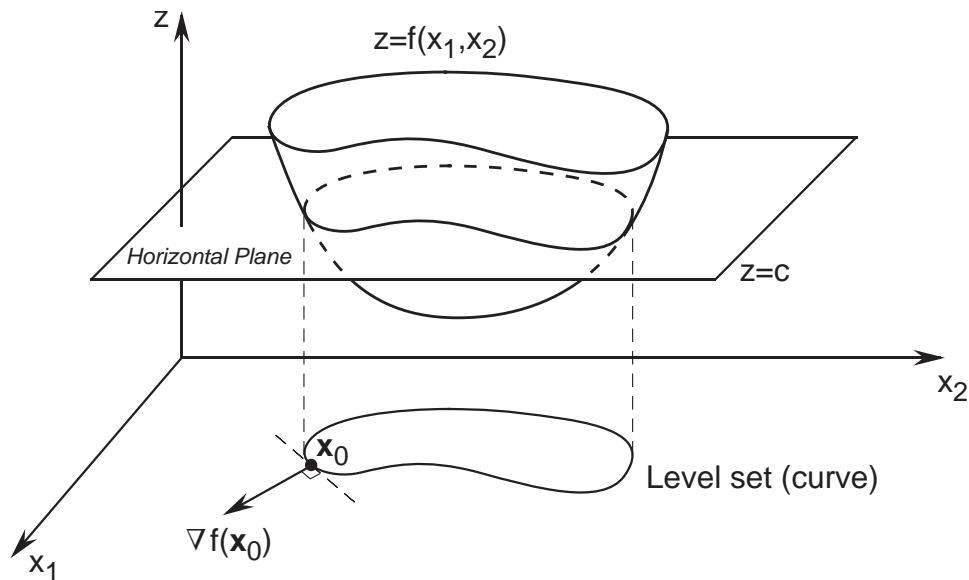
- If Q is positive definite, then f is a parabolic “bowl.”



- Quadratics are useful in the study of optimization.
- Often, objective functions are “close to” quadratic near the solution.
- It is easier to analyze the behavior of algorithms when applied to quadratics.
- Analysis of algorithms for quadratics gives insight into their behavior in general.

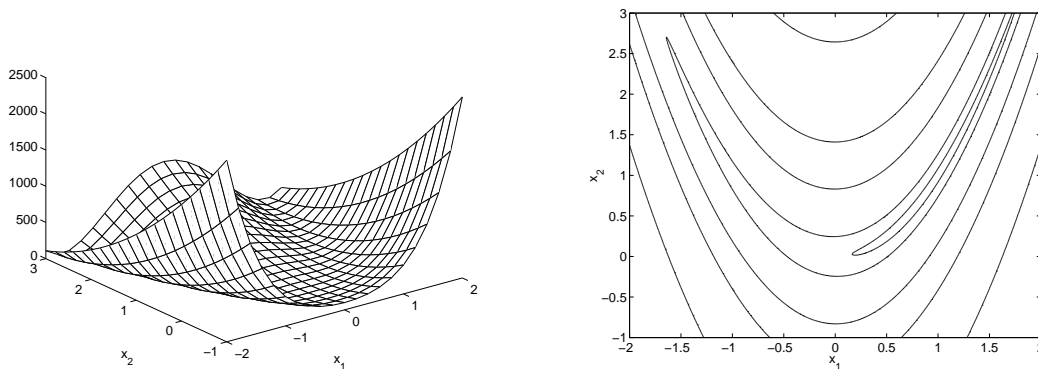
Level sets (§5.5)

- The *level set* of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at level c is the set of points $S = \{\mathbf{x} : f(\mathbf{x}) = c\}$.
- The level set of f is a subset of \mathbb{R}^n .



- Example (Rosenbrock's function):

$$f(\mathbf{x}) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2, \quad \mathbf{x} = [x_1, x_2]^T.$$



Derivatives (§5.1–5.2)

- Given $f : \mathbb{R} \rightarrow \mathbb{R}$
- The *derivative* of f is a function $f' : \mathbb{R} \rightarrow \mathbb{R}$ given by

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

if the limit exists.

- Also written $\frac{df}{dx}$
- If the derivative exists, we say that f is *differentiable*.
- If f' is continuous, we say that f is *continuously differentiable*.
- Given $f : \mathbb{R}^n \rightarrow \mathbb{R}$
- The *gradient* of f is a function $\nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ given by

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

- At each \mathbf{x} , $\nabla f(\mathbf{x})$ is a vector in \mathbb{R}^n .
- Given $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $\mathbf{f} = [f_1, \dots, f_m]^T$.
- The derivative of \mathbf{f} is a function $D\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^{m \times n}$ given by

$$D\mathbf{f}(\mathbf{x}) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(\mathbf{x}) & \dots & \frac{\partial f_1}{\partial x_n}(\mathbf{x}) \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1}(\mathbf{x}) & \dots & \frac{\partial f_m}{\partial x_n}(\mathbf{x}) \end{bmatrix}.$$

- Sometimes called *Jacobian*.
- At each \mathbf{x} , $D\mathbf{f}(\mathbf{x})$ is an $m \times n$ matrix.
- If $D\mathbf{f}$ is continuous, we say that \mathbf{f} is *continuously differentiable*.
- We write $\mathbf{f} \in \mathcal{C}^1$.
- Note that for $f : \mathbb{R}^n \rightarrow \mathbb{R}$, we have

$$\nabla f(\mathbf{x}) = Df(\mathbf{x})^T.$$

- If the derivative of ∇f exists, we say that f is *twice differentiable*.

Write the second derivative as $D^2 f$ (or \mathbf{F}), and call it the *Hessian* of f .

$$\mathbf{F} = D^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_2 \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_n \partial x_1} \\ \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_2^2} & \dots & \frac{\partial^2 f}{\partial x_n \partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_n} & \frac{\partial^2 f}{\partial x_2 \partial x_n} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}.$$

- If \mathbf{F} is continuous, we write $f \in \mathcal{C}^2$.

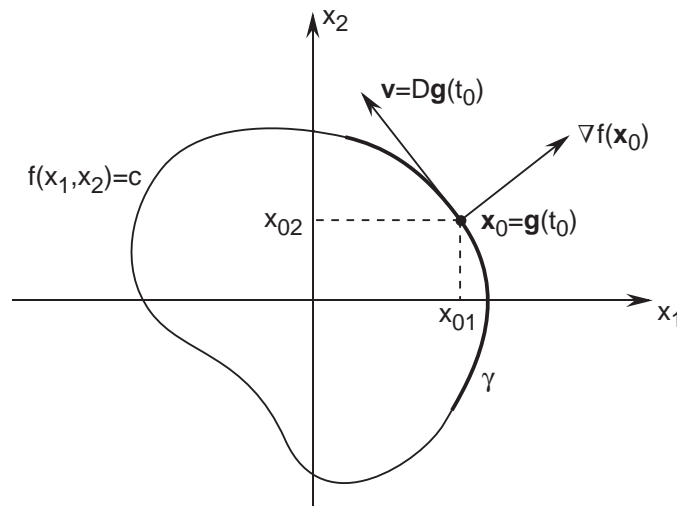
Chain rule (§5.4)

- Given $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $g : \mathbb{R} \rightarrow \mathbb{R}^n$, both differentiable.
- Suppose $f \in \mathcal{C}^1$.
- Define the composite function $F : \mathbb{R} \rightarrow \mathbb{R}$ by $F(t) = f(g(t))$.
- Then, F is differentiable, and

$$\begin{aligned}
 F'(t) &= Df(g(t)) \cdot Dg(t) \\
 &= \nabla f(g(t))^T g'(t) \\
 &= g'(t)^T \nabla f(g(t)).
 \end{aligned}$$

Gradients and level sets (§5.5)

- Given $f : \mathbb{R}^n \rightarrow \mathbb{R}$.
- Fact: $\nabla f(x_0)$ is orthogonal to the level set at x_0



Proof of fact:

- Imagine a particle traveling along the level set.
- Let $g(t)$ be the position of the particle at time t , with $g(0) = x_0$.
- Note that $f(g(t)) = \text{constant}$ for all t .
- Velocity vector $g'(t)$ is tangent to the level set.

- Consider $F(t) = f(\mathbf{g}(t))$. We have $F'(0) = 0$. By the chain rule,

$$F'(0) = (\mathbf{g}'(0))^T \nabla f(\mathbf{g}(0)).$$

- Hence, $\nabla f(\mathbf{x}_0)$ and $\mathbf{g}'(0)$ are orthogonal.

Taylor's formula (§5.5)

- Suppose $f : \mathbb{R} \rightarrow \mathbb{R}$ is in \mathcal{C}^1 .

- Then,

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + o(x - x_0).$$

- $o(h)$ is a term such that $o(h)/h \rightarrow 0$ as $h \rightarrow 0$.
- Around x_0 , f can be approximated by a linear function, and the approximation gets better the closer we are to x_0 .

- Suppose $f : \mathbb{R} \rightarrow \mathbb{R}$ is in \mathcal{C}^2 .

- Then,

$$\begin{aligned} f(x) &= f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2 \\ &\quad + o((x - x_0)^2). \end{aligned}$$

- Around x_0 , f can be approximated by a quadratic function.

- Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$.

- If $f \in \mathcal{C}^1$, then

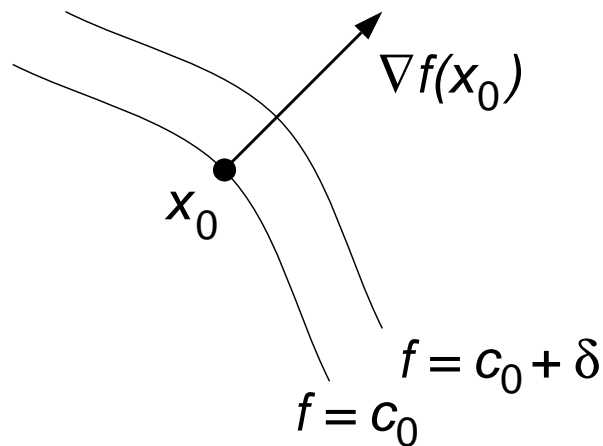
$$f(\mathbf{x}) = f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0) + o(\|\mathbf{x} - \mathbf{x}_0\|).$$

- If $f \in \mathcal{C}^2$, then

$$\begin{aligned} f(\mathbf{x}) &= f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0) \\ &\quad + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \mathbf{F}(\mathbf{x}_0) (\mathbf{x} - \mathbf{x}_0) \\ &\quad + o(\|\mathbf{x} - \mathbf{x}_0\|^2). \end{aligned}$$

In what direction does a gradient point?

- Given $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{x}_0 \in \mathbb{R}^n$.
- We already know that $\nabla f(\mathbf{x}_0)$ is orthogonal to the level set at \mathbf{x}_0 .
- Suppose $\nabla f(\mathbf{x}_0) \neq \mathbf{0}$.
- Fact: ∇f points in the direction of increasing f .



Proof of fact:

- Consider $\mathbf{x}_\alpha = \mathbf{x}_0 + \alpha \nabla f(\mathbf{x}_0)$, $\alpha > 0$.
- By Taylor's formula,

$$\begin{aligned} f(\mathbf{x}_\alpha) &= f(\mathbf{x}_0) + (\mathbf{x}_\alpha - \mathbf{x}_0)^T \nabla f(\mathbf{x}_0) + o(\|\mathbf{x}_\alpha - \mathbf{x}_0\|) \\ &= f(\mathbf{x}_0) + \alpha \|\nabla f(\mathbf{x}_0)\|^2 + o(\alpha). \end{aligned}$$

- Therefore, for sufficiently small α , $f(\mathbf{x}_\alpha) > f(\mathbf{x}_0)$.