

Lecture 2: Math Review

Learning goals:

- Review some basic math concepts the course is built upon
 - Topics:
 - Vector spaces, norms, inner products
 - Analysis + functions
 - Vector calculus
 - Linear algebra
-

- In our course, optimization variables will be elements of finite-dimensional vector spaces

- canonical example: \mathbb{R}^n

- n -dimensional Euclidean vector space

- $x \in \mathbb{R}^n \iff x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

- well-defined operations for adding vectors and multiply them by scalars (usually real numbers)

- another important example: $\mathbb{R}^{n \times m}$

- vector space of $n \times m$ real-valued matrices
- can be identified with a Euclidean vector space of dimension nm
- but get distinct and interesting spaces using matrix norms (more soon)

- we will not consider infinite-dimensional vector spaces such as $C = \{f: \mathbb{R} \rightarrow \mathbb{R} \mid f \text{ is continuous}\}$

- cannot represent any element as a linear combination of a finite set of ~~bases~~ vectors
- more mathematically intricate since certain properties of norms, inner products, etc. not the same as in finite-dim spaces
- can pose interesting optimization problems in these spaces, but they are mostly theoretically rather than computationally interesting

- We can use norms to endow our vector space with notions of magnitude and distance

- there are many different norms on vector spaces

Definition A function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a norm if

① $f(x) \geq 0 \quad \forall x \in \mathbb{R}^n$ (non negativity)

② $f(x) = 0 \iff x = 0$ (definiteness)

③ $f(tx) = |t| f(x) \quad \forall x \in \mathbb{R}^n, \forall t \in \mathbb{R}$ (homogeneity)

④ $f(x+y) \leq f(x) + f(y) \quad \forall x, y \in \mathbb{R}^n$ (triangle inequality)

Notation: $f(x) = \|x\|$ or $\|x\|_{\text{symbol}}$
↑
to indicate a specific norm

Can use norm to define distance between vectors:

$$\text{dist}(x, y) = \|x - y\|$$

Examples:

• 2-norm or Euclidean norm: $\|x\|_2 = \sqrt{\sum_i x_i^2}$

• 1-norm: $\|x\|_1 = \sum_i |x_i|$

• ∞ -norm: $\|x\|_\infty = \max_i |x_i| = \max\{|x_1|, \dots, |x_n|\}$

• p-norm: $\|x\|_p = \left(\sum_i |x_i|^p\right)^{\frac{1}{p}}, \quad p \geq 1$

• Quadratic norm: $\|x\|_Q = \sqrt{x^T Q x} = \left\| \overset{\substack{\uparrow \\ \text{matrix} \\ \text{square root}}}{Q^{\frac{1}{2}}} x \right\|_2, \quad Q = Q^T \succ 0$
($Q \in S_{++}^n$)

Matrix norms on $\mathbb{R}^{m \times n}$: $X = \begin{bmatrix} x_{11} & \dots & x_{1n} \\ \vdots & & \vdots \\ x_{m1} & \dots & x_{mn} \end{bmatrix} \in \mathbb{R}^{m \times n}$

- Frobenius norm: $\|X\|_F = \sqrt{\text{tr}(X^T X)} = \sqrt{\sum_{i=1}^m \sum_{j=1}^n x_{ij}^2}$
 - i.e. the Euclidean norm of the vector of all entries
- Sum-absolute value norm: $\|X\|_{\text{sav}} = \sum_{i=1}^m \sum_{j=1}^n |x_{ij}|$
- Max-absolute value norm: $\|X\|_{\text{max}} = \max_{i,j} |x_{ij}|$
- Induced matrix norms (aka operator norms):

Definition: An operator norm of $X \in \mathbb{R}^{m \times n}$, induced by the vector norms $\|\cdot\|_a$ and $\|\cdot\|_b$ is

$$\|X\|_{a,b} = \max_u \{ \|Xu\|_a \mid \|u\|_b \leq 1 \}$$

- When the same vector norm is used on both spaces we write $\|X\|_c$

• Examples:

• $\|X\|_2 = \sigma_{\text{max}}(X) = \sqrt{\lambda_{\text{max}}(X^T X)}$

\swarrow max singular value \nwarrow max eigenvalue

• $\|X\|_1 = \max_j \sum_{i=1}^m |x_{ij}| = \text{max column sum}$

• $\|X\|_{\infty} = \max_i \sum_{j=1}^n |x_{ij}| = \text{max row sum}$

• Nuclear norm: $\|X\|_* = \text{trace}(\sqrt{X^T X}) = \sum_{i=1}^{\min(m,n)} \sigma_i(X)$
 = sum of singular values

- The set of vectors with norm less than or equal to 1

$$B = \{ x \in \mathbb{R}^n \mid \|x\| \leq 1 \}$$

is called the unit ball of the norm $\|\cdot\|$.

- Dual norms: Let $\|\cdot\|$ be a norm on \mathbb{R}^n .

The associated dual norm, denoted $\|z\|_*$, is defined as

$$\|z\|_* = \max \{ z^T x \mid \|x\| \leq 1 \}$$

- Examples:

- $\|x\|_{1*} = \|x\|_\infty$

- $\|x\|_{2*} = \|x\|_2$

- $\|x\|_{\infty*} = \|x\|_1$

- More generally, $\|x\|_{p*} = \|x\|_q$ where q satisfies

$$\frac{1}{p} + \frac{1}{q} = 1$$

- We can give further structure to a vector space by using an inner product, which gives notions of angle and orthogonality

Definition: A function $\langle \cdot, \cdot \rangle : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ is an inner product if

- ① $\langle x, x \rangle \geq 0$, $\langle x, x \rangle \Leftrightarrow x = 0$ (positivity)
- ② $\langle x, y \rangle = \langle y, x \rangle$ (symmetry)
- ③ $\langle x+y, z \rangle = \langle x, z \rangle + \langle y, z \rangle$ (additivity)
- ④ $\langle tx, y \rangle = t \langle x, y \rangle \quad \forall t \in \mathbb{R}$ (homogeneity)

• standard inner product on \mathbb{R}^n :

$$\langle x, y \rangle = x^T y = \sum x_i y_i$$

• standard inner product on $\mathbb{R}^{m \times n}$:

$$\langle X, Y \rangle = \text{trace}(X^T Y) = \sum_i \sum_j X_{ij} Y_{ij}$$

• the angle between nonzero vectors $x, y \in \mathbb{R}^n$ is

$$\angle(x, y) = \cos^{-1} \left(\frac{x^T y}{\|x\|_2 \|y\|_2} \right)$$

we say x and y are orthogonal if $\langle x, y \rangle = 0$

• Any inner product defines a norm given by $f(x) = \sqrt{\langle x, x \rangle}$

Theorem (Cauchy-Schwarz) For any $x, y \in \mathbb{R}^n$

$$|\langle x, y \rangle| \leq \|x\| \|y\|$$

• Analysis / Topology

- Open and closed sets

Definitions :

- an element x of a subset $C \subseteq \mathbb{R}^n$ is called an interior point of C if $\exists \varepsilon > 0$ for which

$$\{y \in \mathbb{R}^n \mid \|y - x\| \leq \varepsilon\} \subseteq C$$

- the set of all interior points of C is called the interior of C , denoted $\text{int } C$

- a set $C \subseteq \mathbb{R}^n$ is called open if $C = \text{int } C$

- a set $C \subseteq \mathbb{R}^n$ is called closed if its complement $\bar{C} = \mathbb{R}^n \setminus C = \{x \in \mathbb{R}^n \mid x \notin C\}$ is open

- the closure of a set C is defined as

$$\text{cl } C = \mathbb{R}^n \setminus \text{int}(\mathbb{R}^n \setminus C)$$

- a set C is closed iff it contains the limit point of every convergent sequence in it, and the closure is the set of all limit points of convergent sequences on C

- the boundary of a set C is defined as

$$\text{bd } C = \text{cl } C \setminus \text{int } C$$

- The supremum of $C \subseteq \mathbb{R}$ is its least upper bound
 - We take $\sup \emptyset = -\infty$ and $\sup C = \infty$ if C is unbounded above
- The infimum of $C \subseteq \mathbb{R}$ is its greatest lower bound
 - We take $\inf \emptyset = \infty$ and $\inf C = -\infty$ if C is unbounded below

Functions

- The notation $f: A \rightarrow B$ means that f is a function on a subset of A into the set B
- A function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is continuous at $x \in \mathbb{R}^n$ if $\forall \varepsilon > 0, \exists \delta > 0$ such that for $y \in \mathbb{R}^n$

$$\|y - x\| \leq \delta \Rightarrow \|f(y) - f(x)\|_2 \leq \varepsilon$$
- A function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is called closed if $\forall \alpha \in \mathbb{R}$ the sublevel set $\{x \in \mathbb{R}^n \mid f(x) \leq \alpha\}$ is closed
 Equivalently, f is closed iff the epigraph of f

$$\text{epi } f = \{(x, t) \in \mathbb{R}^{n+1} \mid x \in \mathbb{R}^n, f(x) \leq t\}$$
- The level set of a function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is the set

$$\{x \in \mathbb{R}^n \mid f(x) = \alpha\}$$

Derivatives

Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$

- The partial derivative with respect to x_i is

$$\frac{\partial f}{\partial x_i} = \lim_{t \rightarrow 0} \frac{f(x + te_i) - f(x)}{t}$$

- The gradient of f is the (column) vector of partial derivatives

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

Note: the derivative of f is the row vector

$$Df(x) = \nabla f(x)^T$$

- The Hessian of f , denoted $\nabla^2 f(x)$, is the $n \times n$ symmetric matrix of second derivatives

$$(\nabla^2 f)_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

- The Jacobian of $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$, where $f = \begin{bmatrix} f_1(x) \\ \vdots \\ f_m(x) \end{bmatrix}$ is the $m \times n$ matrix of (first) derivatives

$$J_f(x) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \nabla f_1(x)^T \\ \vdots \\ \nabla f_m(x)^T \end{bmatrix}$$

- The first and second order Taylor expansions of a function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ at a point x_0 are

$$f(x) \approx f(x_0) + \nabla f(x_0)^T (x - x_0) + o(\|x - x_0\|)$$

$$f(x) \approx f(x_0) + \nabla f(x_0)^T (x - x_0) + \frac{1}{2} (x - x_0)^T \nabla^2 f(x_0) (x - x_0) + o(\|x - x_0\|^2)$$

"little o" notation: $f = o(g(x))$
 if $\lim_{x \rightarrow 0} \frac{|f(x)|}{|g(x)|} = 0$

Differentiation Rules

- Product Rule. Let $f, g: \mathbb{R}^n \rightarrow \mathbb{R}^m$, $h(x) = f^T(x)g(x)$

$$J_h(x) = f^T(x) J_g(x) + g^T(x) J_f(x), \quad \nabla h(x) = J_h(x)^T$$

- Chain Rule. Let $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $g: \mathbb{R}^m \rightarrow \mathbb{R}^p$.
 Define the composition $h: \mathbb{R}^n \rightarrow \mathbb{R}^p$ by $h(x) = g(f(x))$

Then

$$J_h(x) = J_g(f(x)) J_f(x)$$

Common functions

• Linear functions: $f(x) = c^T x$, $c \in \mathbb{R}^n$, $c \neq 0$

• Affine functions: $f(x) = c^T x + b$, $c \in \mathbb{R}^n$, $b \in \mathbb{R}$

$$\nabla f(x) = c, \quad \nabla^2 f(x) = 0$$

• Quadratic functions: $f(x) = x^T Q x + c^T x + b$, $Q = Q^T$
 $Q \in \mathbb{R}^{n \times n}$

$$\nabla f(x) = 2Qx + c$$

$$\nabla^2 f(x) = 2Q$$

• Important composition example:

$$A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m$$

Let $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ ~~be affine~~ be affine: $f(x) = Ax + b$

Let $g: \mathbb{R}^m \rightarrow \mathbb{R}^p$ and define $h(x) = g(f(x)) = g(Ax + b)$

Then applying the chain rule gives

$$J_h(x) = J_g(Ax + b) A$$

When $m = 1$, we get the gradient formula

$$\nabla g(x) = A^T \nabla f(Ax + b)$$

Linear Algebra

Let $A \in \mathbb{R}^{m \times n}$

- The range of A is $R(A) = \{Ax \mid x \in \mathbb{R}^n\} \subseteq \mathbb{R}^m$
- The nullspace (or kernel) of A is $N(A) = \{x \in \mathbb{R}^n \mid Ax = 0\} \subseteq \mathbb{R}^n$
- These are subspaces of \mathbb{R}^m and \mathbb{R}^n , resp.

- If V is a ~~sub~~ subspace of \mathbb{R}^n , its orthogonal complement is

$$V^\perp = \{x \in \mathbb{R}^n \mid z^T x = 0 \quad \forall z \in V\}$$

- A fundamental result of linear algebra:

$$N(A) = R(A^T)^\perp, \quad R(A) = N(A^T)^\perp$$

Symmetric Eigenvalue Decomposition

Let $A \in S^n$ (i.e. $A = A^T \in \mathbb{R}^{n \times n}$).

Then A can be factored/decomposed as

$$A = Q \Lambda Q^T$$

where Q is orthogonal (i.e. $Q^T Q = I$)

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$$

The eigenvalues of $A \in S^n$ are real and can be ordered as

$$\lambda_{\max}(A) = \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n = \lambda_{\min}(A)$$

• We have

$$\det A = \prod_{i=1}^n \lambda_i, \quad \operatorname{tr}(A) = \sum_{i=1}^n \lambda_i$$

$$\|A\|_2 = \max_i |\lambda_i|, \quad \|A\|_F = \sqrt{\sum_{i=1}^n \lambda_i^2} \\ = \max(\lambda_1, -\lambda_n)$$

• A matrix $A \in S^n$ is called positive semidefinite ($A \succeq 0$) if

$$x^T A x \geq 0 \quad \forall x \in \mathbb{R}^n$$

• " " " " positive definite ($A \succ 0$) if

$$x^T A x > 0 \quad \forall x \in \mathbb{R}^n, x \neq 0 \quad (A \succ 0)$$

• " " " " negative semidefinite if $-A \succeq 0$

• " " " " negative definite if $-A \succ 0$
($A \prec 0$)

• There is a partial ordering of symmetric matrices

• we write $A \succeq B$ when $A - B \succeq 0$

• in general, for $A, B \in S^n$ it's possible that $A \not\succeq B$ and $B \not\succeq A$

Singular Value Decomposition

Let $A \in \mathbb{R}^{m \times n}$ with $\operatorname{rank} A = r$. Then A can be factored as

$$A = U \Sigma V^T$$

where $U \in \mathbb{R}^{m \times r}$ and $U^T U = I$, $V \in \mathbb{R}^{n \times r}$ and $V^T V = I$, and

$\Sigma = \operatorname{diag}(\sigma_1, \dots, \sigma_r)$ with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ } singular values