PRINCETON UNIVERSITY

ORFE

Instructor:
Amir Ali Ahmadi
Spring 2017

TA: Georgina Hall

- This lecture:

   Mathematical Background

      o Inner products and norms

      o Positive semidefinite matrices

      o Basics of differential calculus

   We also establish our notation.

---

### Inner products and norms

A function $\langle \cdot , \cdot \rangle : \mathbb{R}^n \times \mathbb{R}^n \longrightarrow \mathbb{R}$ is an __inner product__ if

① $\langle x, x \rangle \geq 0$, and $\langle x, x \rangle = 0 \Leftrightarrow x = 0$      (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 $\langle rx, y \rangle = r \langle x, y \rangle, \forall r \in \mathbb{R}$      (homogeneity)

o Additivity in the second argument follows:

$$\langle x, y+z \rangle \overset{②}{=} \langle y+z, x \rangle \overset{③}{=} \langle y, x \rangle + \langle z, x \rangle \overset{②}{=} \langle x, y \rangle + \langle x, z \rangle$$

o Homogeneity in the second argument follows:

$$\langle x, ry \rangle \overset{②}{=} \langle ry, x \rangle \overset{④}{=} r \langle y, x \rangle \overset{②}{=} r \langle x, y \rangle$$

__Examples__:

○ The standard inner product in $\mathbb{R}^n$:

$$\langle x, y \rangle = x^T y = \sum_{i=1}^{n} x_i y_i. \qquad (x, y \in \mathbb{R}^n)$$

---

○ The standard inner product between matrices:

$$\langle X, Y \rangle = Tr(X^T Y) = \sum_{i=1}^{m} \sum_{j=1}^{n} X_{ij} Y_{ij} \qquad (X, Y \in \mathbb{R}^{m \times n}).$$

<u>Notation</u>. $\mathbb{R}^{m \times n}$: The space of real $m \times n$ matrices.

$Tr(Z) =$ The trace of a (square) matrix $Z$; i.e, $\sum_{i} Z_{ii}$.

<u>Note</u>. The matrix inner product is the same as our original inner product

applied to two vectors of length $mn$ obtained by stacking the columns of our matrices.

○ An example of a less standard inner product in $\mathbb{R}^2$:

$$\langle x, y \rangle = 5 x_1 y_1 + 8 x_2 y_2 - 6 x_1 y_2 - 6 x_2 y_1$$

$$\begin{pmatrix} y_1 \\ y_2 \end{pmatrix}^T \begin{pmatrix} 5 & -6 \\ -6 & 8 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

Symmetry  ✓

homogeneity ✓

additivity ✓

Positivity: $\langle x, x \rangle = 5 x_1^2 + 8 x_2^2 - 12 x_1 x_2 = (x_1 - 2 x_2)^2 + (2 x_1 - 2 x_2)^2$

$$\langle x, x \rangle = 0 \quad \Rightarrow \quad \begin{cases} x_1 = 2 x_2 \\ x_1 = x_2 \end{cases} \Rightarrow x_1 = x_2 = 0.$$

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

Given an inner product $\langle \cdot, \cdot \rangle$, define the length of a vector $x$ to be:

$$\|x\| = \sqrt{\langle x, x \rangle}$$

Theorem (Cauchy-Schwarz inequality). $\left| \langle x, y \rangle \right| \leq \|x\| \, \|y\|$.

Proof. Suppose first $\|x\| = \|y\| = 1$.

$$\|y - x\|^2 \geq 0 \Rightarrow \langle y - x, y - x \rangle \geq 0 \Rightarrow \langle y, y \rangle + \langle x, x \rangle - 2 \langle x, y \rangle \geq 0 \Rightarrow 2 \geq 2 \langle x, y \rangle$$

$$\Rightarrow \langle x, y \rangle \leq 1.$$

Now consider general $x, y \in \mathbb{R}^n$, $x, y \neq 0$ (otherwise the inequality is trivial). We know

$$\left\langle \frac{x}{\|x\|}, \frac{y}{\|y\|} \right\rangle \leq 1 \Rightarrow \langle x, y \rangle \leq \|x\| \, \|y\|. \qquad ①$$

Finally, since ① holds $\forall x, y$, replace $y$ with $-y$

$$\Rightarrow \langle x, -y \rangle \leq \|x\| \, \|-y\| \Rightarrow \langle x, y \rangle \geq - \|x\| \, \|y\| \qquad ②$$

$$① + ② \Rightarrow \left| \langle x, y \rangle \right| \leq \|x\| \, \|y\|. \qquad \square$$

**Norms:**

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

① $f(x) \geq 0 \quad \forall x, \ f(x) = 0 \Leftrightarrow x = 0$     (positivity)

② $f(\alpha x) = |\alpha| \, f(x), \ \forall \alpha \in \mathbb{R}$     (homogeneity)

③ $f(x + y) \leq f(x) + f(y)$     (Triangle inequality)

<u>Examples</u>:

$$\|x\|_2 = \sqrt{\sum x_i^2}, \qquad \|x\|_1 = \sum |x_i|, \qquad \|x\|_\infty = \max_i |x_i|$$

$$\|x\|_p = \left( \sum |x_i|^p \right)^{1/p}, \qquad p \geq 1.$$

<u>Lemma</u>. Let $\langle x, y \rangle$ be any inner product, then $f(x) = \sqrt{\langle x, x \rangle}$ is a norm.

<u>Proof</u>. Positivity follows from definition. Homogeneity:

$$f(\alpha x) = \sqrt{\langle \alpha x, \alpha x \rangle} = |\alpha| \sqrt{\langle x, x \rangle} = |\alpha| f(x).$$

Triangle inequality:

Suppose not $\Rightarrow x, y$ s.t. $\sqrt{\langle x+y, x+y \rangle} > \sqrt{\langle x, x \rangle} + \sqrt{\langle y, y \rangle}$

$$\Rightarrow \quad \langle x+y, x+y \rangle > \langle x, x \rangle + \langle y, y \rangle + 2\sqrt{\langle x, x \rangle \langle y, y \rangle}$$

$$\Rightarrow \quad 2 \langle x, y \rangle > 2\sqrt{\langle x, x \rangle}\sqrt{\langle y, y \rangle}$$

Contradicting the Cauchy Schwarz inequality. $\square$

<u>Note</u>: Not every norm comes from an inner product.

<u>Matrix norms</u>: One can also define norms on matrices.

$$\|X\|_F = Tr^{1/2}(X^T X) = \left( \sum_{i=1}^m \sum_{j=1}^n X_{ij}^2 \right)^{1/2}. \qquad \text{(the Frobenius norm)}$$

$$\|X\|_{sav} = \sum \sum |X_{ij}| \qquad \text{(sum-absolute-value)}$$

$$\|X\|_{mav} = \max_{i,j} |X_{ij}|$$

## Operator norms.

Let $\|\cdot\|_a$, $\|\cdot\|_b$ be norms on $\mathbb{R}^m$ and $\mathbb{R}^n$. We can define the induced matrix norm on $A \in \mathbb{R}^{m\times n}$ as

$$\|A\|_{a,b} = \max \quad \|Ax\|_a$$
$$s.t. \quad \|x\|_b \leq 1$$

This is indeed a norm. Proof of triangle inequality:

$$\|A+B\|_{a,b} = \max_{s.t. \|x\|_b \leq 1} \|Ax+Bx\|_a \leq \max_{s.t. \|x\|_b \leq 1} \|Ax\|_a + \|Bx\|_a$$

$$\leq \max_{\|x\|_b \leq 1} \|Ax\|_a + \max_{\|x\|_b \leq 1} \|Bx\|_b = \|A\|_{a,b} + \|B\|_{a,b} \cdot \square$$

○ Notation: $\|A\|_a := \|A\|_{a,a}$ ; i.e., the same vector norm is used in both spaces.

○ Three common induced norms:

$$\|A\|_2 = \sqrt{\lambda_{max}(A^TA)} \qquad\qquad \left(\|A\|_2 := \|A\|_{2,2}\right)$$

$$\|A\|_1 = \max_j \sum_i |A_{ij}| \qquad\qquad \text{(maximum column sum)}$$

$$\|A\|_\infty = \max_i \sum_j |A_{ij}| \qquad\qquad \text{(maximum row sum)}$$

Not every matrix norm is an induced norm: $\|A\|_F$ isn't (why?)

$$\|I\|_F = \sqrt{n}, \qquad \text{but identity always has induced norm one (why?)}$$

Induced norms are <u>submultiplicative</u>: $\|AB\| \leq \|A\| \|B\|$

Let's first show that $\forall A \in \mathbb{R}^{m \times n}$, $\forall x \in \mathbb{R}^n$ we have $\|Ax\| \leq \|A\| \|x\|$.

Suppose not: $\|Ax\| > \|A\| \|x\|$. $\implies \left\| A \frac{x}{\|x\|} \right\| > \|A\|$, contradicting

the definition of $\|A\|$ as $\max \|Ay\|$.
$$\text{s.t. } \|y\| \leq 1$$

Now, $\|AB\| = \max_{\text{s.t. } \|x\| \leq 1} \|ABx\| \leq \max_{\text{s.t. } \|x\| \leq 1} \|A\| \|Bx\| = \|A\| \max_{\text{s.t } \|x\| \leq 1} \|Bx\| = \|A\| \|B\|$.

○ Not all norms are submultiplicative:

$$\text{e.g., } \|A\|_{mav} = \max_{i,j} |A_{ij}|.$$

Let $A = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \implies A^2 = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix}$. For any submultiplicative norm $\|.\|$ we must have

$$\|A^2\| \leq \|A\|^2. \qquad \text{But } \|A^2\|_{mav} = 2 > \|A\|^2_{mav} = 1.$$

○ But not every submultiplicative norm is an operator norm: e.g., $\|A\|_F$

is submulticative (why?)

## Dual norms.

Let $\|\cdot\|$ be any norm. Its __dual norm__ is defined as

$$\|x\|_* = \max_y \; x^T y .$$
$$\text{s.t. } \|y\| \leq 1$$

So you can think of this as the operator norm of $x^T$.

○ The dual norm is a norm:

$$\| x+z \|_* = \max_y \; x^T y + z^T y \; \leq \; \max_y \; x^T y \; + \; \max_y \; z^T y$$
$$\text{s.t. } \|y\| \leq 1 \qquad \text{s.t. } \|y\| \leq 1 \qquad \text{s.t. } \|y\| \leq 1$$

$$= \| x \|_* + \| z \|_* .$$

Other properties also easy to check.

○ Dual of common norms:
$$\|x\|_{1*} \overset{①}{=} \|x\|_\infty, \quad \|x\|_{2*} \overset{②}{=} \|x\|_2, \quad \|x\|_{\infty *} \overset{③}{=} \|x\|_1 .$$

__Proof of ③__
$$\|x\|_{\infty *} = \max_y \; x^T y$$
$$\|y\|_\infty \leq 1$$

$$y_{opt} = \text{sign}(x) \qquad \Rightarrow \text{ optimal value} = \|x\|_1 .$$

__Proof of ②:__
$$\|x\|_{2*} = \max_y \; x^T y$$
$$\|y\|_2 \leq 1$$

Cauchy-Schwarz $\Rightarrow$ $x^T y \leq \|x\| \|y\| \leq \|x\|.$

But $y = x$ acheives this bound.

__Proof of ①: Exercise.__

Positive semidefinite matrices

Given a matrix $A \in \mathbb{R}^{n \times n}$, we'll be looking often at the quadratic form $x^T A x$. Whenever you see $x^T A x$, w.l.o.g. you may assume $A$ is symmetric. Here's why:

$$A = \overbrace{\left( \frac{A + A^T}{2} \right)}^{B} + \overbrace{\left( \frac{A - A^T}{2} \right)}^{C} , \qquad \text{Note: } x^T C x = 0.$$

$\underbrace{\hspace{2cm}}$ Symmetric part of $A$ $\qquad$ $\underbrace{\hspace{2cm}}$ Anti-symmetric part of $A$

Notation. $S^{n \times n}$: the space of symmetric (real) $n \times n$ matrices.

Definition. A matrix $A \in S^{n \times n}$ is said to be

- Positive semidefinite (psd) if $x^T A x \geq 0 \;\; \forall x \in \mathbb{R}^n$. Notation: $A \succeq 0$.

- Positive definite (pd) if $x^T A x > 0 \;\; \forall x \in \mathbb{R}^n, x \neq 0$. " : $A \succ 0$.

- Negative semidefinite (nsd) if $-A$ is psd. " : $A \preceq 0$.

- Negative definite (nd) if $-A$ is pd. " : $A \prec 0$.

- Indefinite, if it's neither psd nor nsd.

Example: $\begin{bmatrix} 5 & 1 \\ 1 & -2 \end{bmatrix}$ is indefinite: Consider $x = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ and $x = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$.

Notation comment: $A \succeq 0$ means $A$ is psd; $A \geq 0$ means $A_{ij} \geq 0 \;\; \forall ij$.

Eigenvalue characterization

---

**Thm.** Eigenvalues of a real symmetric matrix are real.

**Proof.** Let $Ax = \lambda x$, where $\lambda \in \mathbb{C}$ and $x \in \mathbb{C}^n$.

$\Rightarrow x^* A x = \lambda x^* x$ ①, where $x^*$ is the conjugate transpose.

Let's now take the conjugate of both sides, remembering that $A \in S^{n \times n}$:

$$x^* A^T x = \bar{\lambda} x^* x \Rightarrow x^* A x = \bar{\lambda} x^* x \quad ② \quad (\bar{\lambda} \text{ is the conjugate of } \lambda)$$

$① + ② \Rightarrow (\lambda - \bar{\lambda}) \underset{\substack{\text{evec} \\ \text{non-zero}}}{x^* x} = 0 \Rightarrow \lambda = \bar{\lambda} \Rightarrow \lambda \text{ is real.} \quad \square$

**Thm.** $A$ is psd $\iff$ all eigenvalues of $A$ are nonnegative.

$A$ is pd $\iff$ " " " " " positive.

**Proof.** We only prove the "psd case". The pd claim is similar.

($\Rightarrow$) Suppose some eigenvalue $\lambda$ is negative.

$$A x = \lambda x \Rightarrow x^T A x = \underset{<0}{\lambda} \underset{>0}{x^T x} < 0 \Rightarrow A \text{ not psd.}$$

($\Leftarrow$) For any symmetric matrix we can pick a set of eigenvectors $v_1, \ldots, v_n$ that form an orthogonal basis for $\mathbb{R}^n$.

Pick any $x \in \mathbb{R}^n$.

$$x^T A x = \left( \alpha_1 v_1 + \cdots + \alpha_n v_n \right)^T A \left( \alpha_1 v_1 + \cdots + \alpha_n v_n \right)$$

$$\underset{\substack{v_i^T v_j = 0 \\ i \neq j}}{=} \sum_{i=1}^{n} \alpha_i^2 \, v_i^T A v_i = \sum \underbrace{\alpha_i^2}_{\geq 0} \underbrace{\lambda_i}_{\geq 0} \underbrace{v_i^T v_i}_{> 0} \geq 0. \qquad \square$$

## Sylvester's characterization.

__Thm.__   $A \succeq 0 \iff$ All $2^n - 1$ principal minors are nonnegative.

   $A \succ 0 \iff$ All $n$ leading principal minors are positive.

Minors are determinants of subblocks of $A$. Principal minors are the ones where the block comes from the same row & column index set. Leading principal minors are the ones with index set $1, \ldots, K$, for $K = 1, \ldots, n$.

__2x2:__

$$Q = \begin{bmatrix} a & b \\ b & c \end{bmatrix}, \quad \left[ Q \succ 0 \iff \begin{matrix} a > 0 \\ \underset{det Q}{ac - b^2 > 0} \end{matrix} \right], \quad \left[ Q \succeq 0 \iff \begin{matrix} a \geq 0, \; c \geq 0 \\ ac - b^2 \geq 0 \end{matrix} \right]$$

__3x3:__

$$Q = \begin{bmatrix} a & b & c \\ b & d & e \\ c & e & f \end{bmatrix}, \quad \left[ Q \succ 0 \iff \begin{matrix} a > 0 \\ ad - b^2 > 0 \\ \det Q > 0 \end{matrix} \right], \quad \left[ Q \succeq 0 \iff \begin{matrix} a \geq 0, \; d \geq 0, \; f \geq 0 \\ ad - b^2 \geq 0, \; af - c^2 \geq 0, \; df - e^2 \geq 0 \\ \det Q \geq 0 \end{matrix} \right]$$

__Proof of the theorem.__ We only proved ($\Rightarrow$). Principal submatrices of psd matrices should be psd (why?). The determinant of psd matrices is nonnegative

(why?).

## Differential Calculus

- <u>Continuity</u>. A function $f: \mathbb{R}^n \to \mathbb{R}^m$ is continuous at $a \in \mathbb{R}^n$ if

$$\forall \epsilon > 0, \exists \delta > 0 \quad \text{s.t.} \quad \|x - a\| \leq \delta \implies \|f(x) - f(a)\| \leq \epsilon.$$

  where the choice of the norm is yours.

- <u>Jacobians, gradients, and Hessians</u>

  o Let $f: \mathbb{R}^n \to \mathbb{R}$. $\quad \dfrac{\partial f}{\partial x_i} = \lim_{\epsilon \to 0} \dfrac{f(x + \epsilon e_i) - f(x)}{\epsilon}$, where $e_i$ is the $i^{th}$ standard basis vector.

  o For $f: \mathbb{R}^n \to \mathbb{R}^m$, the <u>Jacobian</u> matrix $J_f$ is the $m \times n$ matrix of first partial derivatives:

$$J_f(x) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{pmatrix}.$$

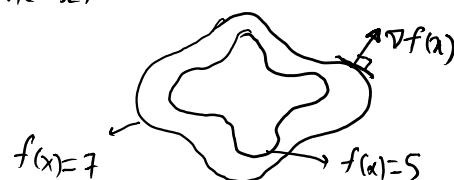  o For $f: \mathbb{R}^n \to \mathbb{R}$, the <u>gradient</u> vector $\nabla f$ is defined as

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

  o For $f: \mathbb{R}^n \to \mathbb{R}$, the <u>Hessian</u> $\nabla^2 f$ is the symmetric matrix of partial derivatives:

$$\left[ \nabla^2 f(x) \right]_{ij} = \frac{\partial f}{\partial x_i \partial x_j}$$

  o For a function $f: \mathbb{R}^n \to \mathbb{R}$, the <u>$\alpha$-level set</u> is the set

$$S_\alpha = \{ x \mid f(x) = \alpha \}.$$

$\nabla f(x)$

$f(x) = 7$

$f(x) = 5$

Basic functions we encounter frequently:

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

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

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

- Quadratic: $f(x) = x^T Q x + c^T x + b$

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

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

Differentiation Rules :

- Product rule: $f, g : \mathbb{R}^n \to \mathbb{R}^m$, $\quad h(x) = f^T(x) g(x)$

then, $J_h(x) = f^T(x) J_g(x) + g^T(x) J_f(x)$, $\quad \nabla_h(x) = J_h^T(x)$.

- Chain rule: $f: \mathbb{R} \to \mathbb{R}^m$, $g: \mathbb{R}^n \to \mathbb{R}$, $h(t) = g(f(t))$.

$$h'(t) = \nabla g^T(f(t)) \begin{bmatrix} f_1'(t) \\ \vdots \\ f_n'(t) \end{bmatrix}$$

- Important special case.

Fix $x, y \in \mathbb{R}^n$. Consider $g: \mathbb{R}^n \to \mathbb{R}$ and let.

$$h(t) = g(x + ty).$$

Then,

$$h'(t) = y^T \nabla g(x + ty).$$

<u>Taylor expansion</u>:

Let $f \in C^m$ (m times continuously differentiable)

Taylor expansion around $a$ (in one variable):

$$f(b) = f(a) + \frac{h}{1!} f'(a) + \frac{h^2}{2!} f''(a) + \cdots + \frac{h^m}{m!} f^{(m)}(a) + o(h^m).$$

where $h := b - a$, and the "little o" is defined as follows:

$$f(x) = o(g(x)) \text{ if } \lim_{x \to 0} \frac{|f(x)|}{|g(x)|} = 0 \quad (\text{"f goes to zero faster than g".})$$

In many dimensions, the following pop up often:

$$f: \mathbb{R}^n \to \mathbb{R}$$

<u>$1^{st}$ order</u>: $\quad f(x) = f(x_0) + \nabla f(x_0)^T (x - x_0) + o(\|x - x_0\|)$

<u>$2^{nd}$ order</u>: $\quad f(x) = f(x_0) + \nabla f^T(x_0)(x - x_0) + \frac{1}{2}(x - x_0)^T \nabla^2 f(x_0)(x - x_0) + o(\|x - x_0\|^2).$

Notes

For more background material see Appendeix A of [BV04].

References

- [BV04] S. Boyd and L. Vandenberghe. Convex Optimization. Cambridge Press, 2004.

- [CZ13] E.K.P. Chong and S.H. Zak. An Introduction to Optimization. Fourth Edition. Wiley, 2013.