# CS 2601 Linear and Convex Optimization
## 2. Math review

### Bo Jiang

John Hopcroft Center for Computer Science
Shanghai Jiao Tong University

### Fall 2022

# Outline

- First-order conditions for unconstrained local min


- Second-order conditions for unconstrained local min

## Review: Derivative

$x$ is an interior point of $X \subset \mathbb{R}^n$ if there exists $\epsilon > 0$ s.t. $B(x, \epsilon) \subset X$.

The interior of $X$, denoted by $\operatorname{int} X$, is the set of interior points of $X$.

A function $f : X \subset \mathbb{R}^n \to \mathbb{R}^m$ is differentiable at $x_0 \in \operatorname{int} X$, if there exists a matrix[1] $A \in \mathbb{R}^{m \times n}$ s.t.

$$\lim_{\Delta x \to 0} \frac{f(x_0 + \Delta x) - f(x_0) - A \Delta x}{\|\Delta x\|} = 0$$

i.e.

$$\Delta f := f(x_0 + \Delta x) - f(x_0) = A \Delta x + o(\|\Delta x\|)$$

The affine function $f(x_0) + A(x - x_0)$ is the first-order approximation of $f$ at $x_0$,

$$f(x) = f(x_0) + A(x - x_0) + o(\|x - x_0\|)$$

---

[1] More precisely, a linear transformation represented by matrix $A$

2

## Review: Derivative

The matrix $A$ is called the derivative of $f$ at $x_0$, and we write

$$f'(x_0) = Df(x_0) = A$$

The derivative is given by the Jacobian matrix of $f = (f_1, \ldots, f_m)^T$

$$f'(x_0) = \begin{bmatrix} \frac{\partial f_1(x_0)}{\partial x_1} & \frac{\partial f_1(x_0)}{\partial x_2} & \cdots & \frac{\partial f_1(x_0)}{\partial x_n} \\ \frac{\partial f_2(x_0)}{\partial x_1} & \frac{\partial f_2(x_0)}{\partial x_2} & \cdots & \frac{\partial f_2(x_0)}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m(x_0)}{\partial x_1} & \frac{\partial f_m(x_0)}{\partial x_2} & \cdots & \frac{\partial f_m(x_0)}{\partial x_n} \end{bmatrix}$$

i.e.

$$[f'(x_0)]_{ij} = \frac{\partial f_i(x_0)}{\partial x_j}, \quad i = 1, \ldots, m; j = 1, \ldots, n$$

Note

$$f_i(x_0 + \Delta x) = f_i(x_0) + \sum_{j=1}^{n} \frac{\partial f_i(x_0)}{\partial x_j} \Delta x_j + o(\|\Delta x\|), \quad i = 1, 2, \ldots, m$$

# Review: Derivative

Example. An affine function $f(x) = Ax + b$ from $\mathbb{R}^n$ to $\mathbb{R}^m$ has derivative $f'(x) = A$ at all $x$. In particular, when $m = 1, f(x) = a^T x + b$ has derivative $f'(x) = a^T$, which is a $1 \times n$ matrix, i.e. a row vector.

Proof. In component form,

$$f_i(x) = \sum_{k=1}^{n} A_{ik} x_k + b_i = A_{i1} x_1 + A_{i2} x_2 + \cdots + A_{in} x_n + b_i$$

so

$$\frac{\partial f_i(x_0)}{\partial x_j} = A_{ij} \implies f'(x_0) = A$$

Alternative proof.

$$f(x_0 + \Delta x) - f(x_0) = A \Delta x \implies f'(x_0) = A$$

4

## Review: Derivative

Example. For symmetric $\boldsymbol{A}$, $f(\boldsymbol{x}) = \boldsymbol{x}^T \boldsymbol{A} \boldsymbol{x} = \sum_{i=1}^{n} \sum_{j=1}^{n} A_{ij} x_i x_j$ has derivative

$$f'(\boldsymbol{x}) = 2\boldsymbol{x}^T \boldsymbol{A}$$

Proof.

$$\frac{\partial f}{\partial x_k} = \sum_{i=1}^{n} \sum_{j=1}^{n} A_{ij} \left( x_j \frac{\partial x_i}{\partial x_k} + x_i \frac{\partial x_j}{\partial x_k} \right) = \sum_{j=1}^{n} A_{kj} x_j + \sum_{i=1}^{n} A_{ik} x_i = 2 \sum_{i=1}^{n} x_i A_{ik}$$

Alternatively,

$$f(\boldsymbol{x}_0 + \Delta \boldsymbol{x}) - f(\boldsymbol{x}_0) = \boldsymbol{x}_0^T (\boldsymbol{A} + \boldsymbol{A}^T) \Delta \boldsymbol{x} + \underbrace{\Delta \boldsymbol{x}^T \boldsymbol{A} \Delta \boldsymbol{x}}_{=o(\|\Delta \boldsymbol{x}\|)}$$

Note. For general $\boldsymbol{A}$, $f'(\boldsymbol{x}) = \boldsymbol{x}^T (\boldsymbol{A} + \boldsymbol{A}^T)$. This can also be obtained by noting $\boldsymbol{x}^T \boldsymbol{A} \boldsymbol{x} = \boldsymbol{x}^T \tilde{\boldsymbol{A}} \boldsymbol{x}$ and $f'(\boldsymbol{x}) = 2\boldsymbol{x}^T \tilde{\boldsymbol{A}}$, where $\tilde{\boldsymbol{A}} = \frac{1}{2}(\boldsymbol{A} + \boldsymbol{A}^T)$.

## Review: Gradient

For a real-valued function $f : \mathbb{R}^n \to \mathbb{R}$, the gradient of $f$ at $\boldsymbol{x}$, denoted by $\nabla f(\boldsymbol{x})$, is the transpose of $f'(\boldsymbol{x})$,

$$\nabla f(\boldsymbol{x}) = [f'(\boldsymbol{x})]^T, \quad [\nabla f(\boldsymbol{x})]_i = \frac{\partial f(\boldsymbol{x})}{\partial x_i}, \quad i = 1, \ldots, n$$

$\nabla f(\boldsymbol{x})$ is a column vector and satisfies

$$f'(\boldsymbol{x})\Delta\boldsymbol{x} = \langle \nabla f(\boldsymbol{x}), \Delta\boldsymbol{x} \rangle = \nabla f(\boldsymbol{x})^T \Delta\boldsymbol{x}$$

The first-order approximation of $f$ at $\boldsymbol{x}_0$ is

$$f(\boldsymbol{x}_0) + \nabla f(\boldsymbol{x}_0)^T(\boldsymbol{x} - \boldsymbol{x}_0)$$

Example. For symmetric $A$, the gradient of $f(\boldsymbol{x}) = \boldsymbol{x}^T A \boldsymbol{x} + \boldsymbol{b}^T \boldsymbol{x} + c$ is

$$\nabla f(\boldsymbol{x}) = 2A\boldsymbol{x} + \boldsymbol{b}$$

# Review: Gradient

$\nabla f(\pmb{x})$ is the direction of fastest rate of increase of $f$ at $\pmb{x}$,

$$f(\pmb{x} + \pmb{d}) - f(\pmb{x}) \approx \nabla f(\pmb{x})^T \pmb{d} \le \|\nabla f(\pmb{x})\| \cdot \|\pmb{d}\|$$

where equality holds in the last step iff $\pmb{d} = \alpha \nabla f(\pmb{x})$ for some $\alpha \ge 0$.

## Review: Chain rule

If $f : X \subset \mathbb{R}^n \to \mathbb{R}^m$ is differentiable at $x_0 \in X$, $g : Y \subset \mathbb{R}^m \to \mathbb{R}^p$ is differentiable at $y_0 = f(x_0)$, then the composition of $f$ and $g$ defined by $h(x) = g(f(x))$ is differentiable at $x_0$, and

$$h'(x_0) = g'(y_0)f'(x_0) = g'(f(x_0))f'(x_0)$$

Note. The order is important since $g'(y_0) \in \mathbb{R}^{p \times m}$ and $f'(x_0) \in \mathbb{R}^{m \times n}$ are matrices. In general $f'(x_0)g'(y_0)$ is undefined.

$$
\begin{array}{ccccc}
\mathbb{R}^n & \xrightarrow{\;f\;} & \mathbb{R}^m & \xrightarrow{\;g\;} & \mathbb{R}^p \\[4pt]
x_0 & \mapsto & y_0 = f(x_0) & \mapsto & z_0 = h(x_0) = g(y_0) \\[4pt]
\Delta x & \overset{f'}{\mapsto} & \Delta y \approx f'(x_0)\Delta x & \overset{g'}{\mapsto} & \Delta z \approx g'(y_0)\Delta y \approx g'(y_0)f'(x_0)\Delta x
\end{array}
$$

In component form,

$$[h'(x_0)]_{ij} = \frac{\partial h_i(x_0)}{\partial x_j} = \sum_{k=1}^{m} \frac{\partial g_i(y_0)}{\partial y_k} \cdot \frac{\partial f_k(x_0)}{\partial x_j} = \sum_{k=1}^{m} [g'(y_0)]_{ik} [f'(x_0)]_{kj}$$

8

# Review: Chain rule

Example. $h(x) = f(Ax + b)$ has derivative $h'(x_0) = f'(Ax_0 + b)A$. If $f$ is real-valued,

$$\nabla h(x_0) = A^T[f'(Ax_0 + b)]^T = A^T \nabla f(Ax_0 + b)$$

Example. Given $f : \mathbb{R}^n \to \mathbb{R}$ and $x, d \in \mathbb{R}^n$, define

$$g(t) = f(x + td)$$

Then

$$g'(t) = f'(x + td)d = \nabla f(x + td)^T d = d^T \nabla f(x + td)$$

Note. $g$ is the restriction of $f$ to the straight line through $x$ with direction $d$. We can often get useful information about $f$ by looking at $g$, which is usually easier to deal with.

# First-order necessary condition

Consider unconstrained optimization problem, i.e. $X = \mathbb{R}^n$.

Theorem. If $x^*$ is a local minimum of $f$ and $f$ is differentiable at $x^*$, then its gradient at $x^*$ vanishes, i.e.

$$\nabla f(x^*) = \left( \frac{\partial f(x^*)}{\partial x_1}, \ldots, \frac{\partial f(x^*)}{\partial x_n} \right)^T = \mathbf{0}.$$
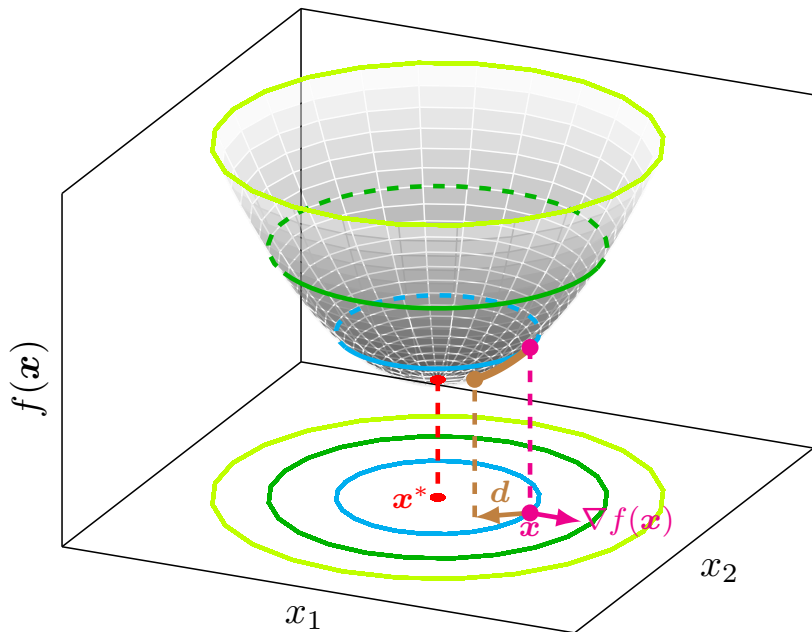
Proof. Let $d \in \mathbb{R}^n$. Define $g(t) = f(x^* + td)$.

- Since $x^*$ is a local minimum, $g(t) \geq g(0)$
- For $t > 0$,

$$\frac{g(t) - g(0)}{t} \geq 0 \implies g'(0) = \lim_{t \downarrow 0} \frac{g(t) - g(0)}{t} \geq 0$$
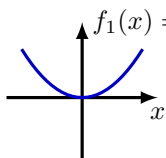
- By chain rule, $g'(0) = \sum_{i=1}^{n} d_i \frac{\partial f(x^*)}{\partial x_i} = d^T \nabla f(x^*) \geq 0$
- Setting $d = -\nabla f(x^*) \implies \|\nabla f(x^*)\|^2 \leq 0 \implies \nabla f(x^*) = \mathbf{0}$
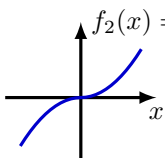
10

# First-order Necessary Condition (cont'd)

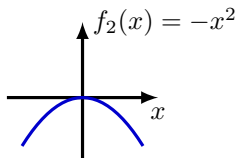A point $x^*$ with $\nabla f(x^*) = \mathbf{0}$ is called a stationary point of $f$.



| $f_1(x) = x^2$ | $f_2(x) = x|x|$ | $f_2(x) = -x^2$ |
|:---:|:---:|:---:|
| $x^* = 0$ | $x^* = 0$ | $x^* = 0$ |
| minimum | inflection point | maximum |



| $f(\boldsymbol{x}) = x_1^2 + x_2^2$ | $f(\boldsymbol{x}) = x_1^2 - x_2^2$ | $f(\boldsymbol{x}) = -x_1^2 - x_2^2$ | $f(\boldsymbol{x}) = -x_1|x_1| + x_2^2$ |
|:---:|:---:|:---:|:---:|
| minimum | saddle point | maximum | |

Note. Will see stationarity is sufficient for convex optimization.
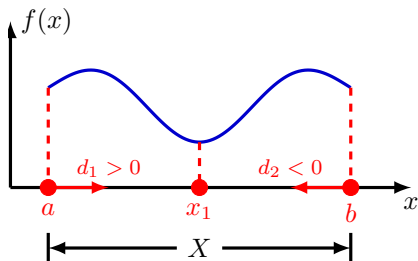
# First-order Necessary Condition (cont'd)

For constrained optimization problem, i.e. $X \neq \mathbb{R}^n$,

- if $\boldsymbol{x}^*$ is in the interior of $X$, i.e. $B(\boldsymbol{x}^*, \epsilon) \subset X$ for some $\epsilon > 0$, then the proof still works, so $\nabla f(\boldsymbol{x}^*) = \boldsymbol{0}$
- otherwise, the proof shows $\boldsymbol{d}^T \nabla f(\boldsymbol{x}^*) \geq 0$ for any feasible direction $\boldsymbol{d}$ at $\boldsymbol{x}^*$
  - $\boldsymbol{d}$ is a feasible direction at $x \in X$ if $\boldsymbol{x} + \alpha\boldsymbol{d} \in X$ for all sufficiently small $\alpha > 0$
- will revisit later

Example. $X = [a, b]$

- $f'(x_1) = 0$
- $d_1 f'(a) \geq 0 \implies f'(a) \geq 0$
- $d_2 f'(b) \geq 0 \implies f'(b) \leq 0$

# Outline

## Review: Second derivative

The second-order partial derivatives of $f : X \subset \mathbb{R}^n \to \mathbb{R}$ at $\boldsymbol{x}_0 \in \operatorname{int} X$ are

$$\frac{\partial^2 f(\boldsymbol{x}_0)}{\partial x_i \partial x_j}, \quad i, j = 1, 2, \ldots, n$$

The Hessian (matrix) of $f$ at $\boldsymbol{x}_0$, denoted by $\nabla^2 f(\boldsymbol{x}_0)$, is given by

$$[\nabla^2 f(\boldsymbol{x}_0)]_{ij} = \frac{\partial^2 f(\boldsymbol{x}_0)}{\partial x_i \partial x_j}, \quad i, j = 1, 2, \ldots, n$$

Note. Do not confuse with Jacobian matrix of vector-valued function.

If $\dfrac{\partial^2 f(\boldsymbol{x})}{\partial x_i \partial x_j}$ and $\dfrac{\partial^2 f(\boldsymbol{x})}{\partial x_j \partial x_i}$ exist in a neighborhood of $\boldsymbol{x}_0$ and are continuous at $\boldsymbol{x}_0$, then

$$\frac{\partial^2 f(\boldsymbol{x}_0)}{\partial x_i \partial x_j} = \frac{\partial^2 f(\boldsymbol{x}_0)}{\partial x_j \partial x_i}$$

so $\nabla^2 f(\boldsymbol{x}_0)$ is symmetric.

Will assume twice continuous differentiability when considering $\nabla^2 f$. 15

# Review: Second derivative

Example. For an affine function $f(\boldsymbol{x}) = \boldsymbol{b}^T \boldsymbol{x} + \boldsymbol{c}$

$$\nabla f^2(\boldsymbol{x}) = \boldsymbol{O}$$

Example. For a quadratic function $f(\boldsymbol{x}) = \boldsymbol{x}^T \boldsymbol{A} \boldsymbol{x}$ with a symmetric $\boldsymbol{A}$,

$$\nabla^2 f(\boldsymbol{x}) = 2\boldsymbol{A}$$

Proof. Recall $f'(\boldsymbol{x}) = 2\boldsymbol{x}^T \boldsymbol{A}$, i.e.

$$\frac{\partial f(\boldsymbol{x})}{\partial x_j} = 2 \sum_{k=1}^{n} x_k A_{kj}$$

so

$$\frac{\partial^2 f(\boldsymbol{x})}{\partial x_i \partial x_j} = 2 \sum_{k=1}^{n} \frac{\partial x_k}{\partial x_i} A_{kj} = 2A_{ij}$$

16

## Review: Chain rule for second derivative

The composition with affine function $g(\boldsymbol{x}) = f(\boldsymbol{A}\boldsymbol{x} + \boldsymbol{b})$ has Hessian

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

Proof. Let $\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x} + \boldsymbol{b}$, i.e. $y_k = \sum_i A_{ki} x_i$. Recall $\nabla g(\boldsymbol{x}) = \boldsymbol{A}^T \nabla f(\boldsymbol{y})$, i.e.

$$\frac{\partial g(\boldsymbol{x})}{\partial x_j} = \sum_k \frac{\partial f(\boldsymbol{y})}{\partial y_k}\frac{\partial y_k}{\partial x_j} = \sum_k \frac{\partial f(\boldsymbol{y})}{\partial y_k} A_{kj}$$

$$\frac{\partial^2 g(\boldsymbol{x})}{\partial x_i \partial x_j} = \sum_k \frac{\partial}{\partial x_i}\frac{\partial f(\boldsymbol{y})}{\partial y_k} A_{kj} = \sum_k \sum_\ell \frac{\partial^2 f(\boldsymbol{y})}{\partial y_\ell \partial y_k} A_{\ell i} A_{kj} = [\boldsymbol{A}^T \nabla^2 f(\boldsymbol{y})\boldsymbol{A}]_{ij}$$

Special case. For $g(t) = f(\boldsymbol{x} + t\boldsymbol{d})$,

$$g''(t) = \boldsymbol{d}^T \nabla^2 f(\boldsymbol{x} + t\boldsymbol{d})\boldsymbol{d}$$

Proof. Set $\boldsymbol{A} \leftarrow \boldsymbol{d}$, $\boldsymbol{x} \leftarrow t$, $\boldsymbol{b} \leftarrow \boldsymbol{x}$ in the general formula above.

# Review: Second-order Taylor expansion

The second-order Taylor expansion for $g : \mathbb{R} \to \mathbb{R}$ takes the form

$$g(a + t) = g(a) + g'(a)t + \frac{1}{2}g''(a)t^2 + o(|t|^2) \tag{T1}$$

The second-order Taylor expansion for $f : \mathbb{R}^n \to \mathbb{R}$ takes the form

$$f(\boldsymbol{x} + \boldsymbol{d}) = f(\boldsymbol{x}) + \nabla f(\boldsymbol{x})^T \boldsymbol{d} + \frac{1}{2}\boldsymbol{d}^T \nabla^2 f(\boldsymbol{x})\boldsymbol{d} + o(\|\boldsymbol{d}\|^2) \tag{T2}$$

i.e.

$$f(\boldsymbol{x} + \boldsymbol{d}) = f(\boldsymbol{x}) + \sum_{i=1}^{n} \frac{\partial f(\boldsymbol{x})}{\partial x_i} d_i + \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \frac{\partial^2 f(\boldsymbol{x})}{\partial x_i \partial x_j} d_i d_j + o(\|\boldsymbol{d}\|^2)$$

Note. (T2) can be obtained by applying (T1) to $g(t) = f(\boldsymbol{x} + t\hat{\boldsymbol{d}})$ at $a = 0$ and $t = \|\boldsymbol{d}\|$, where $\hat{\boldsymbol{d}}$ is the unit vector in the direction $\boldsymbol{d}$, i.e. $\boldsymbol{d} = \|\boldsymbol{d}\|\hat{\boldsymbol{d}}$,

$$g(\|\boldsymbol{d}\|) = g(0) + g'(0)\|\boldsymbol{d}\| + \frac{1}{2}g''(0)\|\boldsymbol{d}\|^2 + o(\|\boldsymbol{d}\|^2)$$

By the chain rule, $g'(0) = \nabla f(\boldsymbol{x})^T \hat{\boldsymbol{d}}, \quad g''(0) = \hat{\boldsymbol{d}}^T \nabla^2 f(\boldsymbol{x})\hat{\boldsymbol{d}}$

# Review: Second-order Taylor expansion

For a quadratic function $f(\boldsymbol{x}) = \boldsymbol{x}^T A \boldsymbol{x} + \boldsymbol{b}^T \boldsymbol{x} + c$, the second-order Taylor expansion is exact with no $o(\|\boldsymbol{d}\|^2)$ term, i.e.

$$f(\boldsymbol{x} + \boldsymbol{d}) = f(\boldsymbol{x}) + \nabla f(\boldsymbol{x})^T \boldsymbol{d} + \frac{1}{2} \boldsymbol{d}^T \nabla^2 f(\boldsymbol{x}) \boldsymbol{d}$$

Note. This can be used to find the expressions for $\nabla f$ and $\nabla^2 f$.

Assume $A$ is symmetric; otherwise, replace $A$ by $\tilde{A} = \frac{1}{2}(A + A^T)$.

$$\begin{aligned}
f(\boldsymbol{x} + \boldsymbol{d}) &= (\boldsymbol{x} + \boldsymbol{d})^T A (\boldsymbol{x} + \boldsymbol{d}) + \boldsymbol{b}^T (\boldsymbol{x} + \boldsymbol{d}) + c \\
&= \boldsymbol{x}^T A \boldsymbol{x} + \boldsymbol{d}^T A \boldsymbol{x} + \boldsymbol{x}^T A \boldsymbol{d} + \boldsymbol{d}^T A \boldsymbol{d} + \boldsymbol{b}^T \boldsymbol{x} + \boldsymbol{b}^T \boldsymbol{d} + c \\
&= f(\boldsymbol{x}) + (2A\boldsymbol{x} + \boldsymbol{b})^T \boldsymbol{d} + \frac{1}{2} \boldsymbol{d}^T (2A) \boldsymbol{d}
\end{aligned}$$

Comparison with the Taylor expansion shows that

$$\nabla f(\boldsymbol{x}) = 2A\boldsymbol{x} + \boldsymbol{b}, \quad \nabla^2 f(\boldsymbol{x}) = 2A.$$

19

## Review: Definite matrices

A matrix $A \in \mathbb{R}^{n \times n}$ is positive semidefinite, denoted by $A \succeq O$, if
1. it is symmetric, i.e. $A = A^T$
2. $x^T A x \geq 0$, $\forall x \in \mathbb{R}^n$

It is positive definite, denoted by $A \succ O$, if condition 2 is replaced by
2'. $x^T A x > 0$, $\forall x \in \mathbb{R}^n$ and $x \neq \mathbf{0}$.

Note. For a quadratic form $x^T A x$, can always assume $A$ is symmetric, since

$$x^T A x = x^T A^T x = x^T \left( \frac{A + A^T}{2} \right) x$$

$A$ is negative (semi)definite if $-A$ is positive (semi)definite.

$A$ is indefinite if it is neither positive semidefinite nor negative semidefinite, i.e. there exists $x_1, x_2 \in \mathbb{R}^n$ s.t.

$$x_1^T A x_1 > 0 > x_2^T A x_2$$

20

# Review: Test for positive definiteness

A vector $x$ is an eigenvector of a matrix $A$ with associated eigenvalue $\lambda$ if

$$Ax = \lambda x$$

We can find all eigenvalues by solving $\det(\lambda I - A) = 0$.

Theorem. Let $A$ be a symmetric matrix.

- $A \succ O$ iff all its eigenvalues $\lambda > 0$.
- $A \succeq O$ iff all its eigenvalues $\lambda \geq 0$.

Exmaple. $A = \begin{pmatrix} 1 & 2 \\ 2 & 5 \end{pmatrix}$ is positive definite.

$$\det(\lambda I - A) = (\lambda - 1)(\lambda - 5) - 4 = 0 \implies \lambda = 3 \pm 2\sqrt{2} > 0$$

Exmaple. $A = \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix}$ is positive semidefinite.

$$\det(\lambda I - A) = (\lambda - 1)(\lambda - 4) - 4 = 0 \implies \lambda_1 = 0, \lambda_2 = 5$$

# Review: Test for positive definiteness

Given matrix $A = (a_{ij}) \in \mathbb{R}^{n \times n}$, a $k \times k$ principal submatrix of $A$ consists of $k$ rows and $k$ columns with the same indices $I = \{i_1 < i_2 < \cdots < i_k\}$,

$$A_I = \begin{pmatrix} a_{i_1 i_1} & \cdots & a_{i_1 i_k} \\ \vdots & \ddots & \vdots \\ a_{i_k i_1} & \cdots & a_{i_k i_k} \end{pmatrix}$$

A principal minor of order $k$ of $A$ is $\det A_I$ for some $I$ with $|I| = k$.

If $I = \{1, 2, \ldots, k\}$, $D_k(A) \triangleq \det A_I$ is called the leading principal minor of order $k$.

Theorem (Sylvester). Let $A$ be a symmetric matrix.

- $A \succ O$ iff $D_k(A) > 0$ for $k = 1, 2, \ldots, n$.
- $A \succeq O$ iff $\det A_I \geq 0$ for all $I \subset \{1, 2, \ldots, n\}$

Note. For positive semidefiniteness, we need to check all principal minors, not just the leading principal minors.

# Review: Test for positive definiteness

Exmaple. $A = \begin{pmatrix} 1 & 2 \\ 2 & 5 \end{pmatrix}$ is positive definite.

$$D_1(\boldsymbol{A}) = \det(1) = 1 > 0, \quad D_2(\boldsymbol{A}) = \det \boldsymbol{A} = 1 > 0$$

Example. $A = \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix}$ is positive semidefinite.

$$D_1(A) = \det(1) = 1, \ \det \boldsymbol{A}_{\{2\}} = \det(4) = 4, \ D_2(A) = \det \boldsymbol{A} = 0$$

Note. It is not enough to check $D_k(A) \geq 0$ for all $k$!

Example. $A = \begin{pmatrix} 0 & 0 \\ 0 & -2 \end{pmatrix}$ is negative semidefinite,

$$D_1(\boldsymbol{A}) = \det(0) = 0, \quad D_2(\boldsymbol{A}) = \det \boldsymbol{A} = 0,$$

but

$$\det \boldsymbol{A}_{\{2\}} = \det(-2) = -2 < 0$$

23

## Review: Test for positive definiteness

Exmaple. $A = \begin{pmatrix} 1 & 2 \\ 2 & 5 \end{pmatrix}$ is positive definite.

- Use definition,

$$\boldsymbol{x}^T A \boldsymbol{x} = x_1^2 + 4x_1x_2 + 5x_2^2 = (x_1 + 2x_2)^2 + x_2^2 \geq 0, \quad \forall \boldsymbol{x} \in \mathbb{R}^2$$

with equality $\iff \begin{cases} x_1 + 2x_2 = 0 \\ x_2 = 0 \end{cases} \iff \boldsymbol{x} = 0$

- Find eigenvalues by solving $\det(\lambda \boldsymbol{I} - \boldsymbol{A}) = 0$

$$\det \begin{pmatrix} \lambda - 1 & -2 \\ -2 & \lambda - 5 \end{pmatrix} = (\lambda - 1)(\lambda - 5) - 4 = 0 \implies \lambda = 3 \pm 2\sqrt{2} > 0$$

- Check leading principal minors

$$D_1(\boldsymbol{A}) = \det(1) = 1 > 0, \quad D_2(\boldsymbol{A}) = \det \boldsymbol{A} = 1 > 0$$

# Review: Test for positive definiteness

Exmaple. $A = \begin{pmatrix} 1 & 2 & 1 \\ 2 & 5 & 8 \\ 1 & 8 & 1 \end{pmatrix}$ is not positive definite.

Check leading principal minors

$$D_1(A) = \det(1) = 1 > 0, \quad D_2(A) = \det \begin{pmatrix} 1 & 2 \\ 2 & 5 \end{pmatrix} = 1 > 0$$

$$D_3(A) = \det A = 1 \times \begin{vmatrix} 5 & 8 \\ 8 & 1 \end{vmatrix} - 2 \times \begin{vmatrix} 2 & 8 \\ 1 & 1 \end{vmatrix} + 1 \times \begin{vmatrix} 2 & 5 \\ 1 & 8 \end{vmatrix} = -36 < 0$$

Can also check eigenvalues, e.g. using `numpy.linalg.eig`,

$$\lambda_1 = 11.69585173, \quad \lambda_2 = 0.58307572, \quad \lambda_3 = -5.27892745$$

## Review: Eigendecomposition

A symmetric matrix $A \in \mathbb{R}^{n \times n}$ has the following eigendecomposition

$$A = Q\Lambda Q^T = \sum_{i=1}^{n} \lambda_i v_i v_i^T$$

where $\Lambda = \text{diag}\{\lambda_1, \ldots, \lambda_n\}$, $Q = (v_1, \ldots, v_n)$ is an orthogonal matrix, i.e. $Q^T Q = Q Q^T = I$, and $A v_i = \lambda_i v_i$.

Example. $A = \frac{1}{4}\begin{pmatrix} 3 & -1 \\ -1 & 3 \end{pmatrix}$ has eigenvalues $\lambda_1 = \frac{1}{2}$ and $\lambda_2 = 1$, with corresponding eigenvectors $v_1 = \frac{1}{\sqrt{2}}(1,1)^T$ and $v_2 = \frac{1}{\sqrt{2}}(-1,1)^T$. The eigendecomposition is

$$A = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{-1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} = \frac{1}{2}\begin{pmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix}\begin{pmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix}^T + \begin{pmatrix} \frac{-1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix}\begin{pmatrix} \frac{-1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix}^T$$

## Review: Eigendecomposition

Consider the linear transformation $x \mapsto y = Ax$.

Recall $v_1, \ldots, v_n$ form an orthonormal basis of $\mathbb{R}^n$, so

$$x = Q\tilde{x} = \sum_{i=1}^{n} \tilde{x}_i v_i, \quad y = Q\tilde{y} = \sum_{i=1}^{n} \tilde{y}_i v_i$$

where

$$\tilde{x} = Q^T x, \quad \tilde{y} = Q^T y,$$

Thus

$$y = Ax \iff Q^T y = Q^T A Q \tilde{x} \iff \tilde{y} = \Lambda \tilde{x}$$

In components,

$$\tilde{x}_i = v_i^T x, \quad \tilde{y}_i = v_i^T y$$

so

$$y = Ax = \sum_{i=1}^{n} \lambda_i v_i v_i^T x = \sum_{i=1}^{n} (\lambda_i \tilde{x}_i) v_i \iff \tilde{y}_i = \lambda_i \tilde{x}_i$$

# Review: Eigendecomposition



$$y = Ax$$

corercive

$$\tilde{x} = Q^T x \qquad x = Q\tilde{x}$$

$$\tilde{y} = Q^T y \qquad y = Q\tilde{y}$$

$$\tilde{y} = \Lambda \tilde{x}$$

# Review: Geometry of quadratic forms

Quadratic form $f(\boldsymbol{x}) = \boldsymbol{x}^T \boldsymbol{A} \boldsymbol{x}$ in $\mathbb{R}^2$



$\boldsymbol{A} = \mathsf{diag}\{1, 1\}$
positive definite

$\boldsymbol{A} = \mathsf{diag}\{0, 1\}$
positive semidefinite

$\boldsymbol{A} = \mathsf{diag}\{1, -1\}$
indefinite

$\boldsymbol{A} = \mathsf{diag}\{-1, -1\}$
negative definite

$\boldsymbol{A} = \mathsf{diag}\{-1, 0\}$
negative semidefinite

# Review: Geometry of quadratic forms

Quadratic form $f(\boldsymbol{x}) = \boldsymbol{x}^T \boldsymbol{A} \boldsymbol{x}$ in $\mathbb{R}^2$



$\boldsymbol{A} = \mathsf{diag}\{1, 1\}$

$\boldsymbol{A} = \mathsf{diag}\{\frac{1}{2}, 1\}$

$\boldsymbol{A} = \frac{1}{4}\begin{pmatrix} 3 & -1 \\ -1 & 3 \end{pmatrix}$

## Review: Bounds on quadratic forms

Proposition. For a symmetric matrix $A \in \mathbb{R}^{n \times n}$,

$$\lambda_{\min} \|x\|_2^2 \leq x^T A x \leq \lambda_{\max} \|x\|_2^2, \quad \forall x \in \mathbb{R}^n$$

where $\lambda_{\max}$ and $\lambda_{\min}$ are the largest and the smallest eigenvalues of $A$, respectively.

Proof. Recall that $A$ can be orthogonally diagonalized, i.e. $A = Q \Lambda Q^T$, where $\Lambda = \mathrm{diag}\{\lambda_1, \ldots, \lambda_n\}$ and $Q^T Q = I$. Let $x = Q\tilde{x}$.

$$x^T A x = \tilde{x}^T (Q^T A Q)\tilde{x} = \tilde{x}^T \Lambda \tilde{x} = \sum_{i=1}^n \lambda_i \tilde{x}_i^2 \leq \sum_{i=1}^n \lambda_{\max} \tilde{x}_i^2 = \lambda_{\max} \|\tilde{x}\|_2^2$$

Then use the fact that orthogonal transformations preserve 2-norm, i.e.

$$\|x\|_2^2 = x^T x = (Q\tilde{x})^T (Q\tilde{x}) = \tilde{x}^T (Q^T Q)\tilde{x} = \tilde{x}^T \tilde{x} = \|\tilde{x}\|_2^2.$$

Similarly for $x^T A x \geq \lambda_{\min} \|x\|_2^2$.

# Second-order necessary condition

Theorem. If $f : \mathbb{R}^n \to \mathbb{R}$ is twice continuously differentiable and $\boldsymbol{x}^*$ is a local minimum of $f$, then its Hessian matrix $\nabla^2 f(\boldsymbol{x}^*)$ is positive semidefinite, i.e.

$$\boldsymbol{d}^T \nabla^2 f(\boldsymbol{x}^*) \boldsymbol{d} \geq 0, \quad \forall \boldsymbol{d} \in \mathbb{R}^n$$

Proof. Fix $\boldsymbol{d} \in \mathbb{R}^n$. By the first-order necessary condition, $\nabla f(\boldsymbol{x}^*) = \boldsymbol{0}$. By the second-order Taylor expansion, for any $t > 0$,

$$f(\boldsymbol{x}^* + t\boldsymbol{d}) = f(\boldsymbol{x}^*) + \frac{t^2}{2} \boldsymbol{d}^T \nabla^2 f(\boldsymbol{x}) \boldsymbol{d} + o(t^2 \|\boldsymbol{d}\|^2) \geq f(\boldsymbol{x}^*)$$

So

$$\frac{1}{2} \boldsymbol{d}^T \nabla^2 f(\boldsymbol{x}) \boldsymbol{d} + o(\|\boldsymbol{d}\|^2) \geq 0$$

Taking the limit $t \to 0$ yields $\boldsymbol{d}^T \nabla f(\boldsymbol{x}^*) \boldsymbol{d}^T \geq 0$.

Note. Can apply the same argument to $g(t) = f(\boldsymbol{x}^* + t\boldsymbol{d})$ with local minimum $t^* = 0$ and use chain rule to obtain $g''(0) = \boldsymbol{d}^T \nabla^2 f(\boldsymbol{x}^*) \boldsymbol{d} \geq 0$.

# Second-order sufficient condition

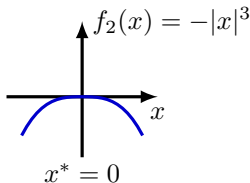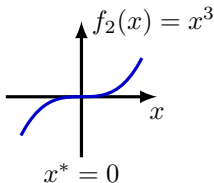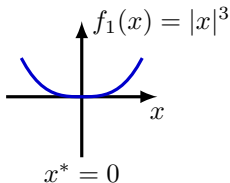Theorem. Suppose $f$ is twice continuously differentiable. If

1. $\nabla f(\boldsymbol{x}^*) = 0$
2. $\nabla^2 f(\boldsymbol{x}^*)$ is positive definite, i.e.

$$\boldsymbol{d}^T \nabla^2 f(\boldsymbol{x}^*)\boldsymbol{d} > 0, \quad \forall \boldsymbol{d} \neq \boldsymbol{0}$$
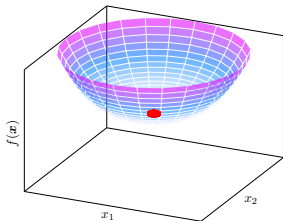
then $\boldsymbol{x}^*$ is a local minimum.

Proof. Use second-order Tayler expansion.

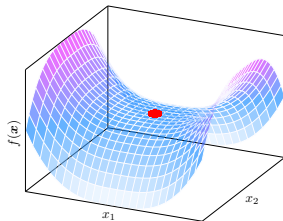Note. In condition 2, positive definiteness cannot be replaced by positive semidefiniteness.



$f_1(x) = |x|^3$  $\quad$  $f_2(x) = x^3$  $\quad$  $f_2(x) = -|x|^3$

$x^* = 0$  $\qquad$  $x^* = 0$  $\qquad$  $x^* = 0$
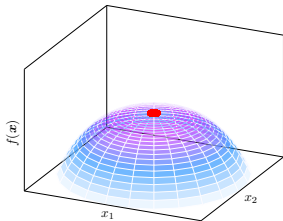
# Second-order sufficient condition (cont'd)

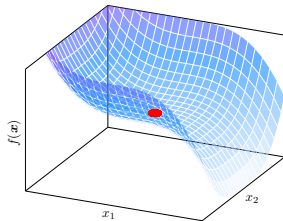$\nabla f(\mathbf{0}) = \mathbf{0}$ and $\nabla^2 f(\mathbf{0}) = \boldsymbol{O}$ for all examples below.



$f(\boldsymbol{x}) = |x_1|^3 + |x_2|^3$



$f(\boldsymbol{x}) = |x_1|^3 - |x_2|^3$



$f(\boldsymbol{x}) = -|x_1|^3 - |x_2|^3$



$f(\boldsymbol{x}) = -x_1^3 + |x_2|^3$