

线性优化

Linear and Convex Optimization

与
凸
优
化

丘 - 润 Y. Qin

5200309/10/55

Introduction to Optimization

9/13

Def. $\min_{x \in X} f(x)$ i.e. minimize $f(x)$

subject to $x \in X$

where $f(x): \mathbb{R}^n \rightarrow \mathbb{R}$
is objective function
(目标函数)

We say " x is feasible" if $x \in X$.

x is infeasible if $x \notin X$.

x is the optimization/decision variable

X is the feasible/constraint set

We say the problem is unconstrained if $X = \mathbb{R}^n$.

and constrained if $X \neq \mathbb{R}^n$.

X is often specified by constrained functions, i.e. $\min f(x)$

s.t. $g_i(x) = 0 \quad i=1, 2, \dots, m$
 $h_i(x) \leq 0 \quad i=1, 2, \dots, k$

here. $X = \{x \in \mathbb{R}^n \mid g_i(x) = 0, i=1, \dots, m\}$
 $h_i(x) \leq 0, i=1, \dots, k\}$

Note. Under strict definition, only " $y = kx$ " is referred to as linear function
while " $y = kx + b (b \neq 0)$ " is named affine function

In this notebook, we use "linear function" to mean both of them.

e.g. Data Fitting: $\min_{k, b \geq 0} \sum \epsilon_i^2$ where ϵ_i is the "measurement error" (误差).

Def. A linear model predicts a response:

$$\hat{y} = f(x) = w^T x + b = \sum_{i=1}^n w_i x_i + b \sim \text{Target: } y$$

Prediction

$$\text{Target: } \min_{w \in \mathbb{R}^n, b \in \mathbb{R}} \sum_{i=1}^m (f(x_i) - y_i)^2 = \min_{w \in \mathbb{R}^n, b \in \mathbb{R}} \sum_{i=1}^m (x_i^T w + b - y_i)^2$$

$$=: \min_{w \in \mathbb{R}^n, b \in \mathbb{R}} \|Xw + b - y\|^2$$

where $X = (x_1 \ \dots \ x_m)^T = \begin{pmatrix} x_1^T \\ \vdots \\ x_m^T \end{pmatrix}, \ 1 = (1, \dots, 1)^T$

$$\|\alpha\| = \sqrt{\alpha^T \alpha} = \sqrt{\sum_{i=1}^n \alpha_i^2}$$

Norm $\overline{\mathbb{R}}^n$

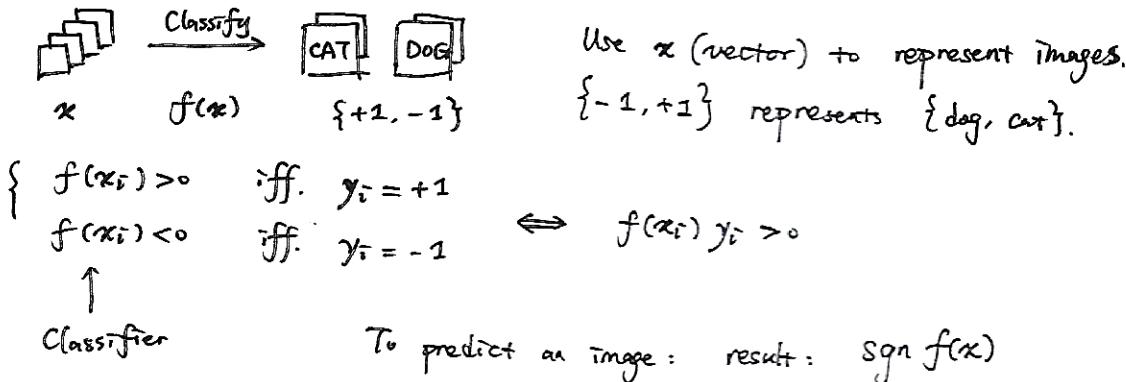
In Practice, we use $\tilde{x}_i = \begin{pmatrix} x_i \\ 1 \end{pmatrix}, \ \tilde{X} = \begin{pmatrix} x_1^T & 1 \\ \vdots & \vdots \\ x_m^T & 1 \end{pmatrix}, \ \tilde{w} = \begin{pmatrix} w \\ b \end{pmatrix}$

The target is then $\min_{\tilde{w} \in \mathbb{R}^{n+1}} \|\tilde{X}\tilde{w} - y\|^2$

Def. Linear Program 线性规划

$$\begin{array}{ll} \min f(x) & \text{s.t.} \\ & g_i(x) = 0 \quad i=1, 2, \dots, m \\ & h_i(x) \geq 0 \quad i=1, 2, \dots, k \\ & x_j \geq 0 \quad i, j \in \{ \} \times \{ \} \end{array}$$

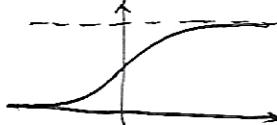
Def. Binary Classification 分类器



How to find "f"?

1) Logistic Regression

$$\sigma(z) := \frac{1}{1+e^{-z}}$$



$\sigma(z)$ can be used to mean any "S-curve" functions in fact.

$$p(y|x) = \sigma(yf(x)) = \sigma(y(w^T x + b))$$

$$\text{Target: } \max_{w,b} \mathcal{L}(w,b) = \max_{w,b} \prod_{i=1}^m p(y_i|x_i)$$

$$\Leftrightarrow \min_{w,b} \text{NLL}(w,b) = -\log \mathcal{L}(w,b) = \sum_{i=1}^m \log (1+e^{-y_i(w^T x_i + b)})$$

Negative Log Likelihood

product: unstable

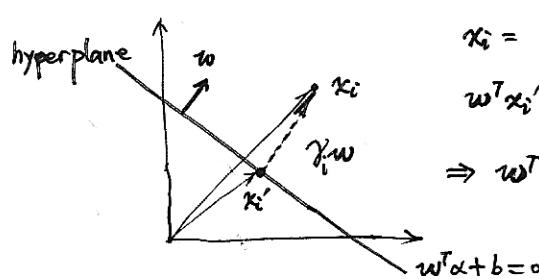
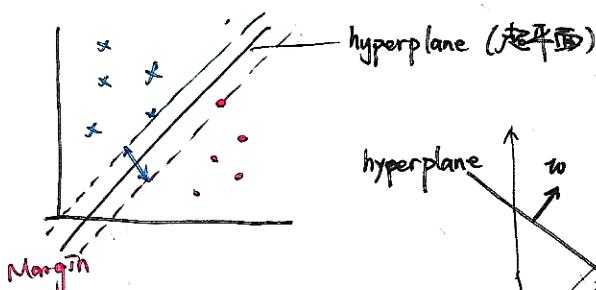
also difficult when you apply gradient. (e.g.: $(uv)' = u'v + v'u$)

more terms

a convex func. in fact!

log: turn the product into a sum. simplify the question.

2) Support Vector Machine



Thus, $\gamma_i = \frac{w^T x_i + b}{\|w\|}$ ← the "distance" of x_i to the hyperplane $w^T x + b = 0$

the "margin" $\gamma_0 = \min_{1 \leq i \leq m} \|y_i w\| \rightarrow$ the target: $\max_{w, b} \gamma_0 = \max_{w, b} \min_{1 \leq i \leq m} \|y_i w\|$

Reform

$$w^T x + b = 0 \quad \text{Let } \tilde{w} = \alpha w, \tilde{b} = \alpha b.$$

\tilde{w} and \tilde{b} determines a hyperplane.

$$x \in P \Leftrightarrow w^T x + b = 0 \Leftrightarrow \tilde{w}^T x + \tilde{b} = 0 \quad (\text{hyperplane})$$

Choosing α property, we can assume

$$\min_{1 \leq i \leq m} y_i (\tilde{w}^T x_i + \tilde{b}) = 1.$$

$$\| \tilde{w}^T x_i + \tilde{b} \| \quad (\text{since } y_i = \text{sgn}(\tilde{w}^T x_i + \tilde{b}))$$

$$\Rightarrow \max_{w, b} \gamma_0 = \max_{w, b} \frac{1}{\|w\|} \Leftrightarrow \min_{w, b} \frac{1}{2} \|w\|^2 \quad (\text{since } \min_{w, b} \|w\| \Leftrightarrow \min_{w, b} \frac{1}{2} \|w\|^2)$$

$$\text{Thus, Target: } \min_{w, b} \frac{1}{2} \|w\|^2 \quad \text{s.t. } y_i (\tilde{w}^T x_i + \tilde{b}) \geq 1 \quad \square$$

Soft Margin SVM

Introduce slack variables $\xi_0 = (\xi_1, \dots, \xi_n)^T$

$$\min_{w, b} \frac{1}{2} \|w\|^2 + c \sum_i \xi_i \quad (c > 0: \text{a hyper parameter})$$

$$\text{s.t. } y_i (\tilde{w}^T x_i + \tilde{b}) \geq 1 - \xi_i, \quad i=1, 2, \dots, n$$

$$\xi \geq 0 \quad (\text{i.e. } \xi_1 \geq 0, \dots, \xi_n \geq 0)$$

Def. Global Optima: $\hat{x} \in X$ is a global ~~minimum~~ minimum of f over X if

$$f(\hat{x}^*) \leq f(x). \quad \forall x \in X$$

\hat{x}^* is an optimal solution of $\min_{x \in X} f(x)$.

Extremum: Maximum and minimum are called extremum.

(might not exist)

e.g. $\inf f(x)$ ↗ unbounded below ($-\infty$)
not achievable

Def. (Linear Algebra Reprise)

Euclidean Inner Product: $\langle x, y \rangle = x^T y = \sum_{i=1}^n x_i y_i$

Euclidean Norm (2-norm, ℓ_2 -norm): $\|x\|_2 = \sqrt{x^T x} = \sqrt{\sum_{i=1}^n x_i^2}$.

A norm on \mathbb{R}^n is a function $\|\cdot\|: \mathbb{R}^n \rightarrow \mathbb{R}$ satisfying
范数

- 1) $\|x\| \geq 0 \quad \forall x \in \mathbb{R}^n$
- 2) $\|x\| = 0 \text{ iff. } x = 0$
- 3) $\|ax\| = |a|\|x\|, \quad \forall a \in \mathbb{R}, x \in \mathbb{R}^n$ (positive homogeneity)
- 4) $\|x+y\| \leq \|x\| + \|y\|, \quad \forall x, y \in \mathbb{R}^n$ (Triangle Inequality)

p -norm: $\|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}, \quad p \geq 1 \rightarrow$ * Minkowski Inequality
Used to prove Property 4).

1-norm: $\|x\|_1 = \sum_i |x_i|$
abs

2-norm: $\|x\|_2 = \sqrt{x^T x}$ (By default, $\|x\| = \|x\|_2$)
Euclidean Norm

∞ -norm: $\|x\|_\infty = \max_i |x_i| = \lim_{p \rightarrow \infty} \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}$

Def. Open ball $B(x_0, r) = \{x: \|x - x_0\| < r\}$

Closed ball $\bar{B}(x_0, r) = \{x: \|x - x_0\| \leq r\}$

A set is open if $\forall x \in S$ exists $\varepsilon > 0$ s.t. $B(x, \varepsilon) \subset S$
is closed if its complement S^c is open.

Convergence: a sequence $\{x_n\}$ converges to x : $x_n \rightarrow x$
收敛

$$\lim_{n \rightarrow \infty} x_n = x \quad (\text{if } \lim_{n \rightarrow \infty} \|x_n - x\| = 0)$$

Thm. if S is closed, \forall sequence $\{x_n\} \subset S$, $x_n \rightarrow x \in S$.

Def. Compactness

A set is bounded if exists $M < \infty$. s.t. $\|x\| < M, \forall x \in S$

A set is compact if it is closed and bounded.

Def. Continuity

A function $f: X \subset \mathbb{R}^n \rightarrow \mathbb{R}$ is continuous at $x \in X$ if

$\forall \varepsilon > 0, \exists \delta > 0$ s.t.

$$y \in X \cap B(x, \delta) \Rightarrow |f(y) - f(x)| < \varepsilon$$

i.e. $\forall \{x_n\} \subset X, x_n \rightarrow x \Rightarrow f(x_n) \rightarrow f(x)$

We say "f is continuous on X" if it is continuous at every $x \in X$.

Thm. If $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is continuous on \mathbb{R}^n , then

- $\forall c \in \mathbb{R}$,
- 1) $\{x : f(x) < c\}$ is open
 - 2) $\{x : f(x) \leq c\}$ is closed.

Thm. (Extreme Value Thm.)

If f is continuous on a compact set X, then

exists $x_1, x_2 \in X$ (not necessarily unique) s.t. $f(x_1) \leq f(x) \leq f(x_2)$ $(\forall x \in X)$

sufficient conditions but not necessary

e.g. $f(x) = x^2$ on $(-1, 2]$. ~~closed~~. global minimum exists. ($x=0$)

Corollary. If f is continuous on \mathbb{R}^n and $f(x) \rightarrow \infty$ as $\|x\| \rightarrow \infty$, then global minimum exists, i.e. exists x^* s.t. $f(x^*) \leq f(x)$. $\forall x$.

Proof. Since $\|x\| \rightarrow \infty$, $f(x) \rightarrow \infty$. exists $M > 0$ s.t. when $\|x\| > M$, $f(x) > f(0)$

The closed ball $[-M, M] =: \bar{B}(0, M) = \{x | \|x\| \leq M\}$ is compact.

By Extreme Value Thm., exists $x^* \in \bar{B}(0, M)$. s.t. $f(x^*) \leq f(x)$. $\forall x \in \bar{B}(0, M)$

On the other side, $f(x^*) \leq f(0) < f(x)$. $\forall x \in \mathbb{R}^n \setminus \bar{B}(0, M)$.

Thus, $\forall x$. $f(x^*) \leq f(x)$, i.e. global min exists.

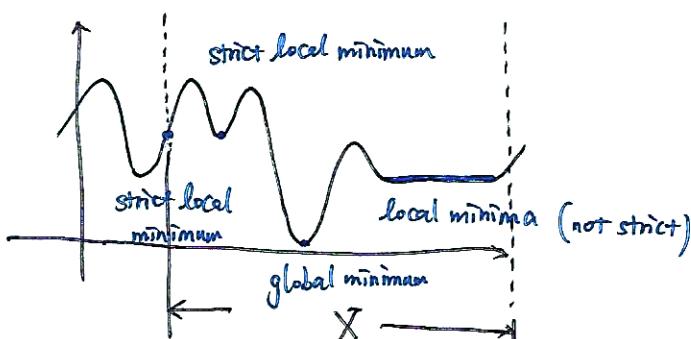
Def.

A function is called coercive if $f(x) \rightarrow \infty$ as $\|x\| \rightarrow \infty$

Def. $x^* \in X$ is a local minimum of f on X if exists $\varepsilon > 0$ s.t.

$$f(x^*) \leq f(x) \quad \forall x \in X \cap \bar{B}(x^*, \varepsilon) \text{ open.}$$

x^* is a strict local minimum if strict inequality holds for $x \neq x^*$.



(none in $X \cap B(x^*, \varepsilon) \setminus \{x^*\}$ s.t. $f(x) = f(x^*)$)

global min is always a local min.
not vice versa.

(← Connex. will reprise)

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

The interior of X , (denoted by $\text{int } X$) is the set of all interior points of X .

Def. A function $f: X \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ is differentiable at $x_0 \in \text{int } X$ if exists a matrix $A \in \mathbb{R}^{m \times n}$ s.t.

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

A: derivative We write $f'(x_0) = Df(x_0) = A$

First-order Approximation: $f(x) = f(x_0) + A(x - x_0) + o(\|x - x_0\|)$

$$\begin{array}{ccc} \uparrow & \uparrow & \uparrow \\ \mathbb{R}^m & \mathbb{R}^{m \times n} & \mathbb{R}^n \end{array}$$

The derivative is given by Jacobian matrix of $f = (f_1, \dots, f_m)$

$$[f'(x_0)]_{i,j} = \frac{\partial f_i(x_0)}{\partial x_j} \quad i \in \{1, 2, \dots, m\} \quad j \in \{1, 2, \dots, n\}$$

Prof. $f_i(x) - f_i(x_0) \approx \sum_k \frac{\partial f_i}{\partial x_k} \Delta x_k$.

$$f(x) - f(x_0) = \begin{pmatrix} f_1(x) - f_1(x_0) \\ \vdots \\ f_n(x) - f_n(x_0) \end{pmatrix} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \cdots & \frac{\partial f_n}{\partial x_n} \end{pmatrix} \begin{pmatrix} \Delta x_1 \\ \vdots \\ \Delta x_n \end{pmatrix}$$

e.g. $\vec{f}(\vec{x}) = A\vec{x} + \vec{b}$. $\vec{f}'(\vec{x}) = A$ (Can be proved).

For symmetric A . $\vec{f}(\vec{x}) = \vec{x}^T A \vec{x} = \sum_i \sum_j A_{ij} x_i x_j$. $\vec{f}'(\vec{x}) = 2\vec{x}^T A$

Prof. $f(\vec{x}) = \sum_i \sum_j A_{ij} x_i x_j$.

$$\frac{\partial f(\vec{x})}{\partial x_k} = 2A_{kk} x_k + \sum_{i \neq k} (A_{ki} + A_{ik}) x_i = 2(A_{kk} + \sum_{i \neq k} A_{ik} x_i)$$

Symmetric! \uparrow \uparrow $= 2 \sum_i A_{ki} x_i$

~~\vec{x}~~

$$\text{Thus, } f'(\vec{x}) = (x_1 \ \dots \ x_n) \begin{pmatrix} A_{11} & \cdots & A_{1n} \\ A_{21} & \ddots & \vdots \\ \vdots & \ddots & \vdots \\ A_{n1} & \cdots & A_{nn} \end{pmatrix} = \vec{x}^T A$$

Note. For general A (not necessarily symmetric),

$$\vec{x}^T A \vec{x} = \vec{x}^T \tilde{A} \vec{x} \text{ where } \tilde{A} = \frac{1}{2}(A + A^T)$$

Thus, $f'(\vec{x}) = 2\vec{x}^T \tilde{A} = \vec{x}^T (A + A^T)$.

* To speed up note-taking, we discard " \vec{x} ", " x " in the following chapters. But we'll note that it's a vec

Prof. Since $x^T A x$ is a number, we have $(x^T A x)^T = x^T A x$.
i.e. $x^T A^T x = x^T A x$

Thus, $x^T A x = \frac{1}{2} (x^T A x + x^T A^T x) = x^T \tilde{A} x$
 \uparrow symmetric!

Def. For a function $f: \mathbb{R}^n \rightarrow \mathbb{R}$, the gradient of f at x (denoted by $\nabla f(x)$) is the transpose of derivative.

i.e. $\nabla f(x) = [f'(x)]^T$. $[\nabla f(x)]_i = \frac{\partial f(x)}{\partial x_i}$

$\nabla f(x)$: a column vector s.t. $f'(x) \Delta x = \langle \nabla f(x), \Delta x \rangle = [\nabla f(x)]^T \Delta x$.
inner product

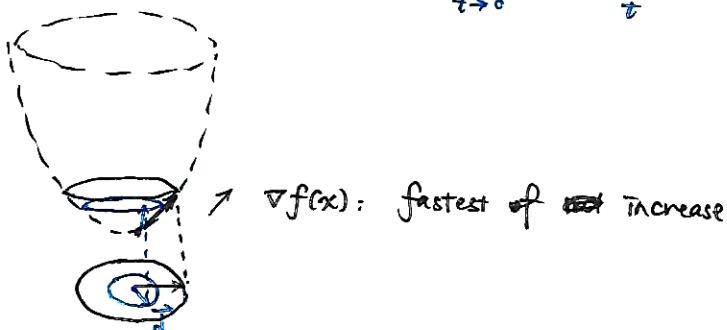
First-order Approximation: $f(x) = f(x_0) + \nabla f(x)_0^T \Delta x + o(\|\Delta x\|)$

$\nabla f(x)$ is the direction of fastest rate of increase of f at x .

$f(x+\vec{d}) = f(x) \doteq \nabla f(x)^T \vec{d} \leq \|\nabla f(x)\| \cdot \|\vec{d}\|$ (Let $\|\vec{d}\|=1$)

where equality holds in the last step iff. $\vec{d} = \alpha \nabla f(x)$ for $\alpha > 0$. (denoted by

Directive \Rightarrow Gradient: $\lim_{t \rightarrow 0} \frac{f(x+t\vec{d}) - f(x)}{t} = \vec{d}^T \nabla f(x)$ $\vec{d} \propto \nabla f(x)$



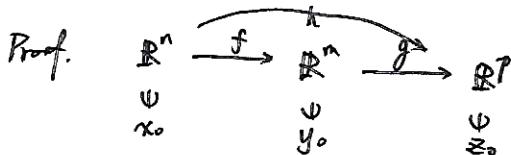
Thm. CHAIN RULE

If $f: X \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ is differentiable at $x_0 \in X$

$g: Y \subset \mathbb{R}^m \rightarrow \mathbb{R}^p$ is differentiable at $y_0 = f(x_0)$.

then the composition of f and g , denoted by $h(x) = g(f(x))$ is differentiable at x_0 ,

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



$\Delta x \rightarrow \Delta y = f'(x_0) \Delta x \rightarrow \Delta z = g'(y_0) \Delta y = g'(f(x_0)) f'(x_0) \Delta x$

$\Delta z = h'(x_0) \Delta x$

Thus, $h'(x_0) = g'(f(x_0)) f'(x_0)$

e.g. $g(x) = f(Ax+b)$

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

Thm. If x^* is a local minimum, then the gradient at x^* must vanish.

i.e. $\nabla f(x^*) = \left(\frac{\partial f(x^*)}{\partial x_1}, \dots, \frac{\partial f(x^*)}{\partial x_n} \right)^T = 0$.

Proof. Let $\vec{d} \in \mathbb{R}^n$. Define $g(t) = f(x^* + t\vec{d})$

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

- For $t > 0$, $\frac{g(t) - g(0)}{t} \geq 0 \Rightarrow g'(0) = \lim_{t \downarrow 0} \frac{g(t) - g(0)}{t} \geq 0$

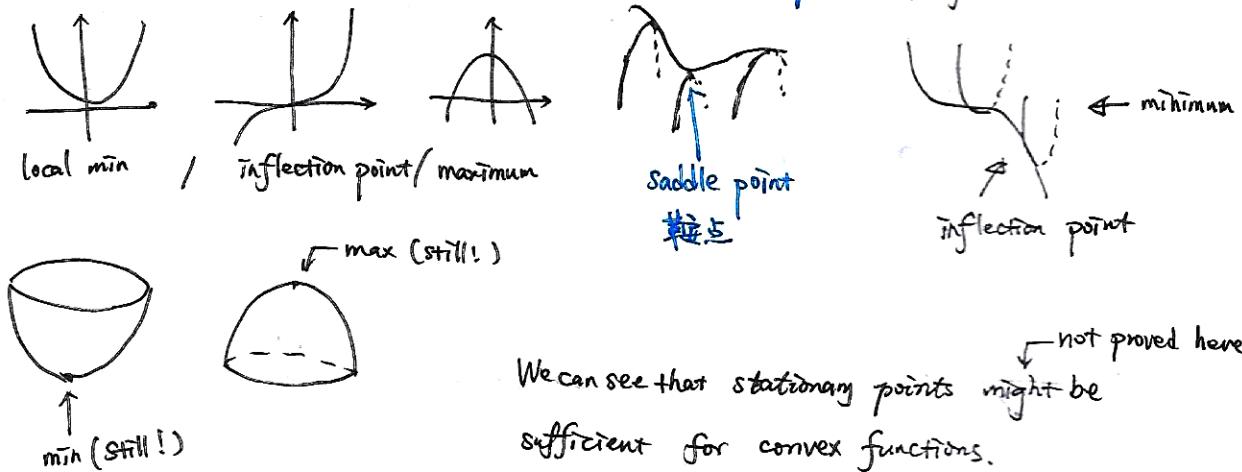
By CHAIN RULE, $g'(0) = \sum_{i=1}^n d_i \frac{\partial f(x^*)}{\partial x_i} = \vec{d}^T \nabla f(x^*) \geq 0$

Setting $d = -\nabla f(x^*) \Rightarrow -\|\nabla f(x^*)\| \geq 0$ Thus, $\|\nabla f(x^*)\| = 0$.

Therefore, $\nabla f(x^*) = 0$ \square

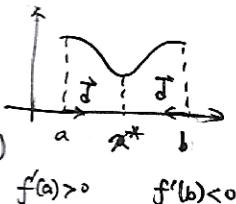
First-order Necessary Condition.

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



For all $x \in X$ except for x^* (min/max value) $\rightarrow \vec{d}^T f'(x) \geq 0$ holds

(for all feasible directions \vec{d})



Second-order Conditions

Def. The second-order partial derivatives of $f: X \subset \mathbb{R}^n \rightarrow \mathbb{R}$ at $x_0 \in \text{int } X$

are

$$\frac{\partial^2 f(x_0)}{\partial x_i \partial x_j}$$

Thm. If $\frac{\partial^2 f(x)}{\partial x_i \partial x_j}$ and $\frac{\partial^2 f(x)}{\partial x_j \partial x_i}$ exist in a neighborhood of x_0 and are continuous at x_0 ,

we have $\frac{\partial^2 f(x_0)}{\partial x_i \partial x_j} = \frac{\partial^2 f(x_0)}{\partial x_j \partial x_i}$

Hessian Matrix:

$$[\nabla^2 f(x_0)]_{ij} = \frac{\partial^2 f(x_0)}{\partial x_i \partial x_j}$$

$$\begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{pmatrix}$$

$$\text{e.g. } f(x) = x^T A x \rightarrow \nabla f'(x) = 2x^T A \rightarrow \nabla^2 f(x) = 2A$$

$$\text{Proof. } \nabla(x^T A). \quad \frac{\partial f}{\partial x_i} = 2 \sum_k A_{ki} x_k \quad \frac{\partial f}{\partial x_i x_j} = 2 \sum_k \frac{\partial x_i}{\partial x_j} A_{ki} = 2 A_{jk}$$

$$\text{CHAIN RULE. } g(x) = f(Ax+b) \rightarrow \nabla^2 g(x) = A^T \nabla^2 f(Ax+b) A$$

Thm. Second-order Taylor Expansion

$$\text{Taylor Expansion: } g(a+t) = g(a) + g'(a)t + \frac{1}{2}g''(a)t^2 + \frac{1}{3!}g'''(a)t^3 + \dots$$

$$\underline{\text{2nd-order T.E.}} \quad \underline{\frac{g(a+t)}{=}}$$

$$\frac{1}{2} \cdot t \cdot g''(a) \cdot t$$

$$\begin{aligned} \text{2nd-order T.E. } f(x+d) &= f(x) + \nabla f(x)^T d + \frac{1}{2} d^T \nabla^2 f(x) d + \underline{o(\|d\|^2)} \\ &= f(x) + \sum_i \frac{\partial f(x)}{\partial x_i} d_i + \frac{1}{2} \sum_{i,j} \frac{\partial^2 f(x)}{\partial x_i \partial x_j} d_i d_j + o(\|d\|^2) \end{aligned}$$

Denotation: $f(x+d) := f(x+\vec{d})$ where $\vec{d} = \|d\| \cdot \vec{t}$ at $x=0$.

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

$$\text{Lagrange Remains: } g(a+t) = g(a) + g'(a)t + \frac{1}{2}g''(a+\xi)t^2, \quad \xi \in (0, t).$$

(PSD)

Def. A matrix A is positive semidefinite, denoted by $A \succeq 0$. if.
半正定

1) A is symmetric, i.e. $A = A^T$

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

positive definite: if $\stackrel{1)}{x^T A x > 0} \quad (\forall x \in \mathbb{R}^n) \quad \text{and } x \neq 0$ denoted by $A > 0$.

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

indefinite: $\exists x_1, x_2 \in \mathbb{R}^n$ s.t. $x_1^T A x_1 > 0 > x_2^T A x_2$.

[Trick (Reprise)] A : not symmetric $\leftarrow \tilde{A} = \frac{1}{2}(A+A^T)$, symmetric.

[Trick] $A = BB^T$ must be positive semidefinite. (where B is a vector)

$$1) \quad A^T = (BB^T)^T = (B^T)^T B^T = BB^T = A.$$

$$2) \quad x^T A x = \underbrace{x^T B B^T x}_{\text{a vector}} = y^T y = \langle y, y \rangle = \|y\|^2 \geq 0.$$

let it be y

When is A definite? $\Rightarrow 2) \quad \|y\|^2 > 0 \quad (\forall x \in \mathbb{R}^n \text{ and } x \neq 0)$

$$y = B^T x = (b_1 \dots b_n) \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \quad \begin{cases} x \neq 0 \Rightarrow y \neq 0 \\ y = 0 \Rightarrow x = 0 \end{cases} \quad \text{Thus, } B \text{ must be linearly independent.}$$

Def. Vector x is called eigenvector of a matrix A with associated eigenvalue λ if

$$Ax = \lambda x$$

(No change in direction, no rotation, only change the "scale")

We can find all λ s by solving $\det(\lambda I - A) = 0$ where $I = \begin{pmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{pmatrix}$

Thm. Let A be a symmetric matrix.

$$1) A \geq 0 \text{ iff all } \lambda \geq 0 \quad 2) A > 0 \text{ iff all } \lambda > 0,$$

Def. Principal Submatrix 主子矩阵

a) p.s. of matrix A consists of k rows and k columns with exactly same indices

$$I = \{i_1 \leq i_2 \leq \dots \leq i_k\}, \text{ i.e. } A_I = \begin{pmatrix} a_{i_1 i_1} & a_{i_1 i_2} & \dots & a_{i_1 i_k} \\ \vdots & \ddots & & \vdots \\ a_{i_k i_1} & a_{i_k i_2} & \dots & a_{i_k i_k} \end{pmatrix}$$

$$\text{e.g. } \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix}. \text{ a } 2 \times 2 \text{ p.s. : } A_{\{1,3\}} = \begin{pmatrix} 1 & 3 \\ 7 & 9 \end{pmatrix}.$$

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

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

Thm. (Sylvester) A symmetric matrix

$$1) A \geq 0 \text{ iff } D_k(A) \geq 0 \text{ for all } k=1, 2, \dots, n$$

$$2) A \geq 0 \text{ iff } \det A_I \geq 0 \text{ for all } I \subset \{1, 2, \dots, n\}$$

* We'll use \succ to denote $>$ and \succeq to denote \geq in this notebook (since it writes faster)
longer than > longer than \geq

Def. 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 $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ $Q = (v_1, \dots, v_n)$ is an orthogonal matrix

i.e. $Q^T Q = Q Q^T = I = \text{diag}(1, \dots, 1)$, and $A v_i = \lambda_i v_i$.

(因正交) in fact (v_1, \dots, v_n) is a basis of $\mathbb{R}^{n \times 1}$.

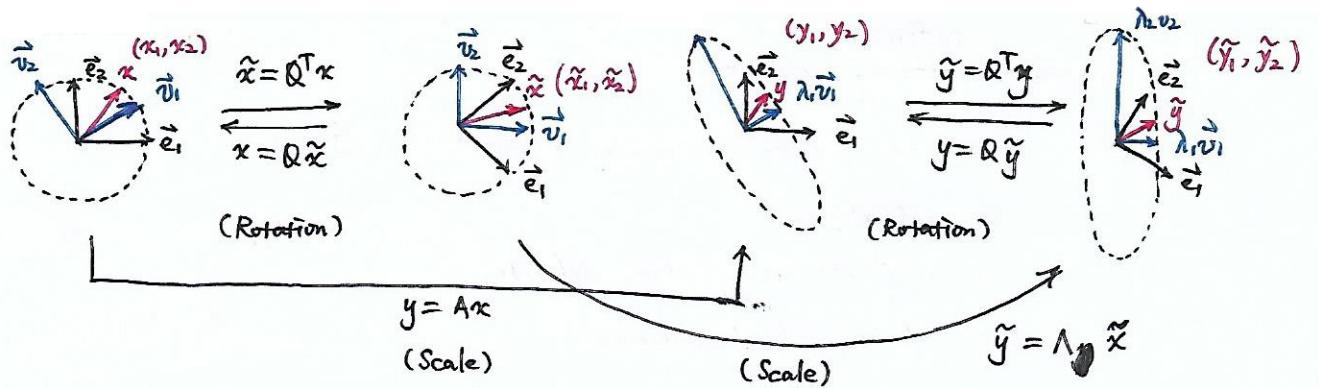
$$\text{Prof. } A Q = A (v_1, \dots, v_n) = (\lambda_1 v_1, \dots, \lambda_n v_n) = (v_1, \dots, v_n) \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} = Q \Lambda.$$

~~$A = Q \Lambda Q^T$~~ Meanwhile, $Q^T Q = Q Q^T = I$.

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

$$= (v_1, \dots, v_n) \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} \begin{pmatrix} v_1^T \\ \vdots \\ v_n^T \end{pmatrix} = \sum_i \lambda_i v_i v_i^T$$

□



$$y = Ax = Q \Lambda Q^T x \Rightarrow Q^T y = \Lambda Q^T x$$

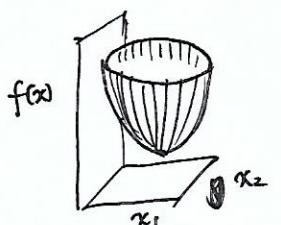
$\underbrace{\qquad}_{\downarrow} \quad \underbrace{\qquad}_{\downarrow}$

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

$$\tilde{x} = Q^T x \Leftrightarrow x = Q \tilde{x} \quad (\text{since } QQ^T = I)$$

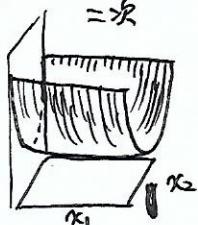
$$\tilde{x} = \tilde{x}_1 v_1 + \tilde{x}_2 v_2 + \dots + \tilde{x}_n v_n$$

Def. Geometry of Quadratic Forms



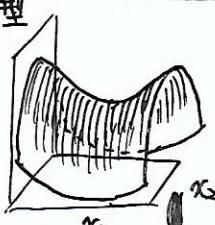
$$A = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

positive definite



$$A = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$$

positive semidefinite



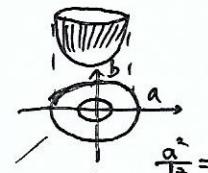
$$A = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$$

Indefinite

$$f(x) = \boxed{x^T A x}$$

$$A = \begin{pmatrix} \frac{1}{2} & 1 \\ 1 & 1 \end{pmatrix}$$

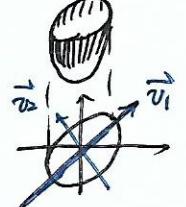
negative definite



$$\frac{\partial^2}{\partial x^2} = \frac{1}{2}$$

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

negative
semidefinite



(Rotation)

Def. Bounds on Quadratic Forms

A: symmetric matrix

$$\lambda_{\min} \|x\|_2^2 \leq x^T A x \leq \lambda_{\max} \|x\|_2^2$$

Proof. 1) $A = \text{diag}(\lambda_1, \dots, \lambda_n)$. Obvious. $\lambda_{\min} \|x\|_2^2 \leq \sum_i \lambda_i x_i^2 \leq \lambda_{\max} \|x\|_2^2$

2) $A = Q \Lambda Q^T$ (since A is symmetric) \leftarrow eigen decomposition

$$\lambda_{\min} \|\tilde{x}\|_2^2 \leq x^T A x = x^T Q \Lambda Q^T x = \tilde{x}^T \Lambda \tilde{x} \quad (\text{Let } \tilde{x} = Q^T x) \leq \lambda_{\max} \|\tilde{x}\|_2^2$$

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

QED. \square

- Second-order necessary conditions

Theorem. $f: \mathbb{R}^n \rightarrow \mathbb{R}$, twice continuously differentiable x^* : local minimum of f .

then its Hessian matrix $\nabla^2 f(x^*)$ is positive semidefinite. i.e. $d^T \nabla^2 f(x^*) d \geq 0$
 $(\forall d \in \mathbb{R}^n)$

Proof: Taylor Expansion. $f(x^* + t d) = f(x^*) + \frac{t^2}{2} d^T \nabla^2 f(x) d + o(t^2 \|d\|^2) \geq f(x^*)$

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

$\Rightarrow f''(\text{local minimum}) \geq 0$.

$$\xrightarrow[t \rightarrow 0]{} \frac{d^T \nabla^2 f(x^*) d}{t^2} \geq 0$$

- Second-order sufficient condition

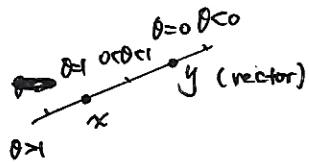
Thm. Suppose f is twice continuously differentiable. If

- 1) $\nabla f(x^*) = 0$
 - 2) $\nabla^2 f(x^*)$ is positive definite.
- then x^* is a local minimum.

Convex Functions and Convex Sets

- Convex Sets

Def. Lines, line segments and Rays



$$z = y + \theta(x-y) = \theta x + (1-\theta)y, \quad \theta \in \mathbb{R}$$

ray: endpoint y , $y \rightarrow \infty$ direction. $z = \theta x + (1-\theta)y, \theta \geq 0$

segment:

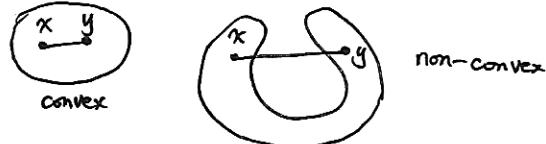
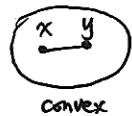
$$z = \theta x + (1-\theta)y, \quad \theta \in [0,1]$$

Sometimes use $\bar{\theta} := 1 - \theta$.

Def. Convex Sets

A set $C \subset \mathbb{R}^n$ is convex if the line segment between any two points $x, y \in C$ lies entirely in C , i.e.

$$\forall x \in C, y \in C, \theta \in [0,1] \Rightarrow \theta x + \bar{\theta}y \in C$$



Here, $\theta x + \bar{\theta}y$ is called a convex combination.

In a more symmetric form: $\theta_1 x + \theta_2 y$ where $\theta_1 \geq 0, \theta_2 \geq 0, \theta_1 + \theta_2 = 1$.

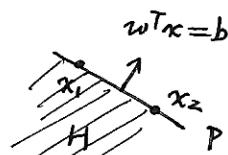
e.g. \emptyset, \mathbb{R}^n , singletons (points), lines, line segments and rays.

Hyperplane: $P = \{x \in \mathbb{R}^n : w^T x = b\}$

Proof. $x_1, x_2 \in P, \theta \in [0,1] \rightarrow \theta x_1 + \bar{\theta} x_2 \in P$.

i.e. to prove $w^T(\theta x_1 + \bar{\theta} x_2) \in P$

$$\begin{aligned} &= \theta \cdot (w^T x_1) + \bar{\theta} \cdot (w^T x_2) \\ &= \theta b + \bar{\theta} b = b. \end{aligned}$$



Q.E.D.

□

Halfspace. $H = \{x \in \mathbb{R}^n : w^T x \leq b\}$.

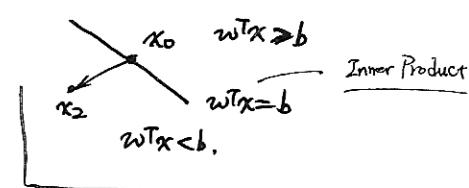
$\{f(x) \leq 0\}$ — sublevel set

Proposition. The intersection of two arbitrary convex sets is convex.

$\Rightarrow \{C_i\}_{i \in I}$ is a set of convex sets.

$C = \bigcap_{i \in I} C_i$ is also convex.

Easy to prove.



e.g. Affine spaces: $\{x \in \mathbb{R}^n : Ax = b\}$ is convex. \rightarrow Multi-hyperplanes' intersection.

[e.g.] polyhedron

$$\mathbb{R}^{m \times n}$$

$$(a_i^T x = b_i)$$

Polyhedra: $\{x \in \mathbb{R}^n : Ax \leq b\}$ is convex. \rightarrow Multi-halfspaces' intersection.

Def. The inequality " \leq ". " \geq " of vectors is interpreted component-wise.

Trick. $Ax = b \Leftrightarrow \begin{cases} Ax \geq b \\ Ax \leq b \end{cases} \Leftrightarrow \begin{pmatrix} A \\ -A \end{pmatrix} x \leq \begin{pmatrix} b \\ -b \end{pmatrix}$

(An affine space is in fact a polyhedron.)



1-norm balls

$$B_1^{(2)} = \{x \in \mathbb{R}^2 : |x_1| + |x_2| \leq a\}$$

↓
four faces — four constraints

$$B_1^{(3)} = \{x \in \mathbb{R}^3 : |x_1| + |x_2| + |x_3| \leq a\}$$

↓
eight faces

Theorem. A closed ball $\bar{B}(x_0, r) = \{x \in \mathbb{R}^n : \|x - x_0\| \leq r\}$ is convex.

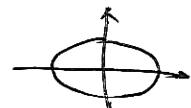
"Norm Ball"

Proof: $\forall x_1, x_2 \in \bar{B}(x_0, r)$ and $\forall \theta \in [0, 1]$

$$\begin{aligned} \|(\theta x_1 + \bar{\theta} x_2) - x_0\| &= \|\theta(x_1 - x_0) + \bar{\theta}(x_2 - x_0)\| \leq \theta \|x_1 - x_0\| + \bar{\theta} \|x_2 - x_0\| \leq \theta r + \bar{\theta} r \\ &= r. \end{aligned}$$

True for any norm $\|\cdot\|$, also true for open ball

e.g. An ellipsoid $\left\{x \in \mathbb{R}^2 : \frac{x_1^2}{\lambda_1^2} + \frac{x_2^2}{\lambda_2^2} \leq 1\right\}$ is convex.



Proof. Transformation $\begin{cases} x_1 = \lambda_1 u_1^* \\ x_2 = \lambda_2 u_2^* \end{cases}$. $\underline{\mathcal{E} \rightarrow \text{ball}}$ proved.

Cod.

Formal Proof. Transformation: $u_1 = \frac{x_1}{\lambda_1}, u_2 = \frac{x_2}{\lambda_2}, x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = \Lambda u$

~~Thus~~, $\bar{B} = \{u : \|u\|_2 \leq 1\}$ is convex.

Now we prove $\mathcal{E} = \{\Lambda u : u \in \bar{B}\}$ is convex.

$$\forall x_1, x_2 \in \mathcal{E}, \theta \in [0, 1]. \exists u_1, u_2 \in \bar{B} \text{ s.t. } x_1 = \lambda_1 u_1, x_2 = \lambda_2 u_2 = \lambda u_2$$

$$\text{Thus, } \theta x_1 + \bar{\theta} x_2 = \theta \lambda_1 u_1 + \bar{\theta} \lambda_2 u_2 = \lambda (\theta u_1 + \bar{\theta} u_2) \in \underline{\mathcal{E}}$$

$\frac{\mathcal{E}}{\bar{B}}$

e.g.'



$$\mathcal{E}' = \{x_0 + \Lambda u : \|u\|_2 \leq 1\}, \Lambda \in \mathbb{R}^{n \times n}, \Lambda \succ 0.$$

$$\Lambda = Q \Lambda Q^T, \text{ Let } \tilde{u} = Q^T u. \Rightarrow \mathcal{E}' = \{x_0 + Q \Lambda \tilde{u} : \|\tilde{u}\|_2 \leq 1\}$$

(shifted and rotated ball)

$$\downarrow \theta x_1 + \bar{\theta} x_2 = x_0 + A(\theta u_1 + \bar{\theta} u_2)$$

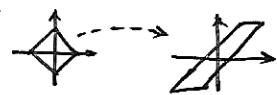
$$\text{Rewrite: } \mathcal{E}' = \{x : (x - x_0)^T P^{-1} (x - x_0) \leq 1\} \text{ where } P = \Lambda^2$$

$\frac{\mathcal{E}'}{\bar{B}}$

$$\text{Prof. } x = x_0 + \Lambda u \Rightarrow u = A^{-1}(x - x_0). \frac{u^T u \leq 1}{(\text{ball})} \Rightarrow (x - x_0)^T A^{-1} A^{-1} (x - x_0) \leq \frac{\downarrow}{P = \Lambda^2}$$

Thm. The image of a convex set under an affine transformation is also convex.

Proof. Let $C \subset \mathbb{R}^n$. $f(x) = Ax + b$ an affine transformation: $\mathbb{R}^n \rightarrow \mathbb{R}^m$.



$$\forall y_1, y_2 \in f(C) := \{f(x) : x \in C\}, \theta \in [0, 1]. \quad \text{By definition.}$$

$$\text{Since } y \text{ is affine, } \theta y_1 + \bar{\theta} y_2 = \theta(Ax_1 + b) + \bar{\theta}(Ax_2 + b) = A(\theta x_1 + \bar{\theta} x_2) + b.$$

Since C is convex, $\theta x_1 + \bar{\theta} x_2 \in C$. Thus, $\theta y_1 + \bar{\theta} y_2 \in f(C)$ \longrightarrow QED. \square

Proposition. The inverse image of a convex set under an affine transformation is convex.
逆映射

f^{-1}

Proof.

Similarly.

e.g. The set $S^+ = \{A \in \mathbb{R}^{n \times n} : A \succeq 0\}$.

\uparrow symmetric of course.

$$\rightarrow \theta A + \bar{\theta} B \in S^+, \theta \in [0, 1].$$

{ Property I. easy to prove <symmetry>

Property II. $x^T (\theta A + \bar{\theta} B)x \geq 0 \quad (\forall x)$

Def.

Def. A convex combination of $x_1, \dots, x_m \in \mathbb{R}^n$ is a point of the form

$$\sum_{i=1}^m \theta_i x_i = \theta_1 x_1 + \dots + \theta_m x_m \quad \text{where } \theta_i \geq 0 \quad (\forall i) \text{ and } \sum_{i=1}^m \theta_i = 1.$$

Thm. If C is convex, $x_1, \dots, x_m \in C$ then any convex combination $\sum_{i=1}^m \theta_i x_i \in C$.

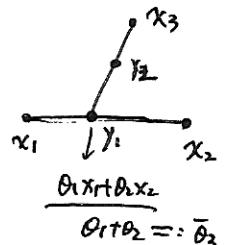
Proof. By Induction. when $n=2$. Obvious

when $n=k$ the proposition holds

$$n=k+1: \quad \sum_{i=1}^{k+1} \theta_i x_i = \overline{\theta_{k+1}} \left(\sum_{i=1}^k \frac{\theta_i}{\theta_{k+1}} x_i \right) + \theta_{k+1} x_{k+1}$$

$\underbrace{\quad}_{EC}$

(TBC)



QED. \square

Def. The convex hull of a set $S \subset \mathbb{R}^n$ (denoted $\text{conv } S$) is the smallest convex set containing S .

Thm. $\text{conv } S = \left\{ \sum_{i=1}^m \theta_i x_i : m \in \mathbb{N}; x_i \in S, \theta_i \geq 0, i=1, \dots, m; \sum_{i=1}^m \theta_i = 1 \right\}$.

Def. Affinely independent points

$(m+1)$ points $x_0, \dots, x_m \in \mathbb{R}^n$ are affinely independent if $x_1 - x_0, \dots, x_m - x_0$ are linearly independent.

e.g.



e.g. $\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ (affinely independent).

$$\begin{pmatrix} 1 & 1 & \dots & 1 \\ x_0 & x_1 & \dots & x_m \end{pmatrix} \begin{pmatrix} c_0 \\ c_1 \\ \vdots \\ c_m \end{pmatrix} = 0.$$

$\uparrow \text{rank } A = m+1$.

Proposition. x_0, \dots, x_m are affinely independent iff.

$$\begin{cases} \sum_{i=0}^m c_i x_i = 0 \\ \sum_{i=0}^m c_i = 0 \end{cases}$$

$$\Rightarrow c_i = 0 \quad (\forall i \in \{1, 2, \dots, m\})$$

Proof is easy (by definition.)

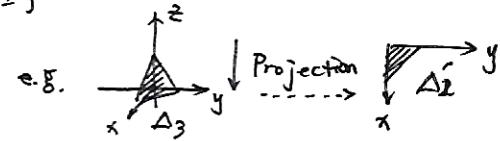
Def. An m -dimensional simplex $\Delta = \text{conv}\{x_0, \dots, x_m\}$, where x_0, \dots, x_m are affinely independent. (also called m -simplex)

\mathbb{R}^n only have m -simplices with $m \leq n$.

The probability n -simplex: $\Delta_n = \{\vec{\theta} \in \mathbb{R}^{n+1} : \vec{\theta} \geq \vec{0}, \vec{1}^T \vec{\theta} = 1\}$ i.e. determined by std vectors $\vec{e}_1, \vec{e}_2, \dots, \vec{e}_{n+1}$.

The unit n -simplex: $\Delta'_n = \{\vec{\theta} \in \mathbb{R}^{n+1} : \vec{\theta} \geq \vec{0}, \vec{1}^T \vec{\theta} \leq 1\}$

In fact: Δ'_k is the projection of Δ_{k+1} !



Relationship between Δ_m and general simplices:

$$\vec{\theta} = \sum_{i=1}^m \theta_i \vec{e}_i \mapsto \vec{x} = \sum_{i=1}^m \theta_i \vec{x}_i = X \vec{\theta} \quad \text{where } X = (x_0, \dots, x_m)_{n \times (m+1)}, \vec{\theta} = (\theta_0, \dots, \theta_m)^T \in \Delta_{m-1}$$

$$\vec{x} = \vec{x}_0 + \sum_{i=1}^m \theta_i (\vec{x}_i - \vec{x}_0) = \vec{x}_0 + B \vec{\theta}' \quad \text{where } \vec{\theta}' \in \Delta'_{m-1}$$

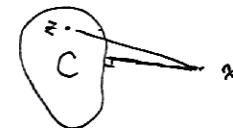
(since $\vec{1}^T \vec{\theta} = 1$)

$$B = (x_1 - x_0, \dots, x_m - x_0)$$

• Projection onto Convex Sets

Def. The Distance between a point x and C is (a set)(nonempty)

$$\text{dist}(x, C) := \inf_{z \in C} \|x - z\|$$



Thm. If $C \subset \mathbb{R}^n$ is nonempty, closed and convex, then for any x , exists a unique $\hat{x} \in C$ s.t. $\|x - \hat{x}\| = \text{dist}(x, C)$

\hat{x} is called the projection of x onto C . denoted by $P_C(x)$.

Specially, $P_C(x) = x$ iff. $x \in C$. (obvious by the definition)

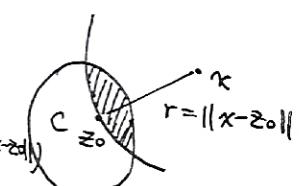
Prof. First show existence.

• Let $z_0 \in C$, $\text{dist}(x, C) \leq \|x - z_0\|$.

• Let $K = \{z \in C : \|x - z\| \leq \|x - z_0\|\} = C \cap \overline{B}(x, \|x - z_0\|)$.

• $\|x - z\|$ is continuous. (~~$\|x - z\| < \|x - y\| < \|x - z\|$~~)

$$(f(z) = \|x - z\|, |f(z) - f(y)| = |\|x - z\| - \|x - y\|| \leq \|z - y\|). \rightarrow 0 \rightarrow 0.$$



Also, K is compact $\Rightarrow \exists \hat{x} \in K$ s.t. $\text{dist}(x, C) = \|x - \hat{x}\|$.

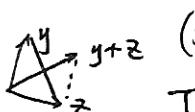
Now show uniqueness. (by contradiction)

If exists \hat{x}_1, \hat{x}_2 s.t. $\|x - \hat{x}_1\| = \|x - \hat{x}_2\| = \text{dist}(x, C) =: d$.

$$\hat{x}_C = \frac{\hat{x}_1 + \hat{x}_2}{2} \in C. \quad \|x - \hat{x}_C\| \leq \|\hat{x}_1 - \hat{x}_2\|$$

$$\text{Let } y = x - \hat{x}_1. \quad z = x - \hat{x}_2. \quad \text{We have } \|y + z\|^2 + \|y - z\|^2 = 2\|y\|^2 + 2\|z\|^2$$



 (Since $\|y \pm z\|^2 = \langle y \pm z, y \pm z \rangle = \|y\|^2 + \|z\|^2 \pm 2\langle y, z \rangle$)

Thus, we have, $\|2x - \hat{x}_1 - \hat{x}_2\|^2 + \|\hat{x}_1 - \hat{x}_2\|^2 = 2\|x - \hat{x}_1\|^2 + 2\|x - \hat{x}_2\|^2 = 4d^2$

$$\text{i.e. } 0 \leq \|\hat{x}_1 - \hat{x}_2\|^2 = 4d^2 - \cancel{\|2(x - x_C)\|^2} = 4d^2 - 4\|x - x_C\|^2 \leq 0$$

$$\cancel{\|2(x - x_C)\|^2} \Rightarrow \|\hat{x}_1 - \hat{x}_2\| = 0 \Rightarrow \hat{x}_1 = \hat{x}_2$$

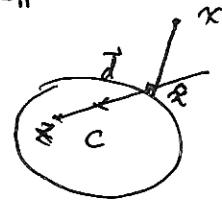
[Proposition]. Let $C \subset \mathbb{R}^n$ be nonempty, closed and convex. Given $\hat{x} \in C$. ~~Prove~~

$$\hat{x} = P_C(x). \text{ Iff. } \langle x - \hat{x}, z - \hat{x} \rangle \leq 0, \forall z \in C.$$

Proof. " \Rightarrow ". Assume $\hat{x} = P_C(x) = \arg\min_{z \in C} \|x - z\| = \arg\min_{z \in C} \|x - z\|^2$

$$\text{Let } f(z) = \|x - z\|^2, \forall z \in C, d = z - \hat{x}$$

$$g(t) = \|x - (\hat{x} + td)\|^2. \quad \Leftrightarrow \nabla f^T d \geq 0.$$



$$f(z) = (x - z)^T (x - z), \quad \langle f'(\hat{x}), d \rangle = \langle \hat{x} - x, z - \hat{x} \rangle \geq 0, \text{ i.e. } \langle x - \hat{x}, z - \hat{x} \rangle \leq 0.$$

" \Leftarrow ". Assume $\langle x - \hat{x}, z - \hat{x} \rangle \leq 0$ for any $z \in C$.

$$\text{Then } \|x - z\|^2 = \|x - \hat{x}\|^2 + \|z - \hat{x}\|^2 - 2\langle x - \hat{x}, z - \hat{x} \rangle \geq \|x - \hat{x}\|^2. \forall z \in C.$$

By definition, $\hat{x} = P_C(x)$.

□

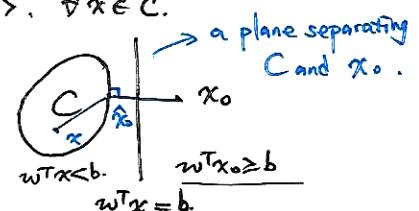
~~QED~~

[Corollary]. Let $C \subset \mathbb{R}^n$ be nonempty, closed and convex. For $x_0 \notin C$.

exists a $w \in \mathbb{R}^n \setminus \{0\}$, s.t. $\langle w, x \rangle < \langle w, x_0 \rangle, \forall x \in C$.

Proof. Let $w = x_0 - P_C(x_0)$. Let $\hat{x}_0 = P_C(x_0)$.

Since $x_0 \notin C$, $P_C(x_0) \in C$, $w \neq 0$.

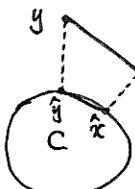


By proposition above, $\frac{\langle x_0 - \hat{x}_0, x - \hat{x}_0 \rangle}{w} \leq 0$.

$$\langle w, x \rangle \leq \langle w, \hat{x}_0 \rangle = \langle w, x_0 \rangle - \langle w, w \rangle < \langle w, x_0 \rangle \quad (w \neq 0, \langle w, w \rangle > 0).$$

□

[Corollary]. The projection operation is nonexpansive, i.e. $\|P_C(x) - P_C(y)\| \leq \|x - y\|$.



Proof. Let $\hat{y} = P_C(y)$, $\hat{x} = P_C(x)$.

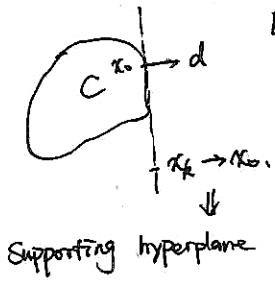
$$\begin{aligned} \|x - y\|^2 &= \|\hat{x} - \hat{y}\|^2 + \underbrace{\|x - y - (\hat{x} - \hat{y})\|^2}_{\geq 0} + 2\langle x - y - (\hat{x} - \hat{y}), \hat{x} - \hat{y} \rangle \\ &\geq \|\hat{x} - \hat{y}\|^2 + 2\langle x - y, \hat{x} - \hat{y} \rangle - 2\langle \hat{x} - \hat{y}, \hat{x} - \hat{y} \rangle \\ &\geq \|\hat{x} - \hat{y}\|^2 \end{aligned}$$

□

Def. Supporting Hyperplane.

Proof. Lemma. If C is convex, $\text{int } C$ and \bar{C} are convex.

Lemma. If $C \subset \mathbb{R}^n$ is convex and $x_0 \in \partial C$. exists a direction d s.t. $x_0 + td \notin \bar{C}$ for any $t > 0$.



[Proof] Replacing C with $C - x_0$. Assume $x_0 = 0$. Prove the Lemma by contradiction.

1) Let e_1, \dots, e_n be std basis vectors. $e_0 = -\sum_{i=1}^n e_i$

2) $\exists t > 0$ s.t. $d_i := te_i \in \bar{C}, (\forall i)$

3) $(d_0, d_1, d_2) \begin{pmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{pmatrix} = x \in \Delta(d_0, d_1, d_2)$

At least $0 \in \Delta(d_0, d_1, d_2)$

$$\begin{pmatrix} 1 & 1 & 1 \\ \underbrace{d_0+u_0}_{a_0} & \underbrace{d_1+u_1}_{a_1} & \underbrace{d_2+u_2}_{a_2} \end{pmatrix} \begin{pmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{pmatrix} = \begin{pmatrix} 1 \\ x \\ 1 \end{pmatrix}$$

$$\det \begin{pmatrix} 1 & 1 & 1 \\ a_0 & a_1 & a_2 \end{pmatrix} \neq 0$$

$$\Rightarrow \begin{pmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 \\ a_0 & a_1 & a_2 \end{pmatrix}^{-1} \begin{pmatrix} 1 \\ x \\ 1 \end{pmatrix} \geq 0. \text{ Then } x \in \Delta(a_0, a_1, a_2).$$

$\downarrow \theta_i \geq 0 \ (\forall i)$

Considering $(d_0, d_1, d_2)^{-1} \begin{pmatrix} 1 \\ x \\ 1 \end{pmatrix} > 0$. when $x \rightarrow 0$. $u_i \rightarrow 0$. $\begin{pmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{pmatrix} > 0$, i.e. $x \in \Delta(a_0, a_1, a_2)$.

Thus, $x \in \text{int } C$. Contradiction. (Since $x_0 \in \partial C$)

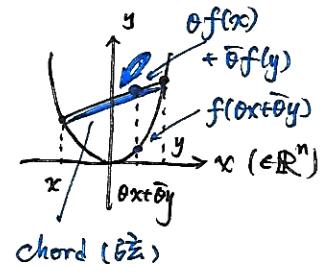
Convex Functions

Def. A function $f: S \subset \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if

1) its domain $\text{dom } f = S$ is a convex set.

2) $\forall x, y \in S, \forall \theta \in [0, 1]$ Jensen's Inequality holds,

$$\text{i.e. } f(\theta x + \bar{\theta} y) \leq \theta f(x) + \bar{\theta} f(y)$$



Note. Condition 1 guarantees $\theta x + \bar{\theta} y \in \text{dom } f$. Actually only need to check 2) for $x \neq y$ and $\theta \in (0, 1)$.

Def. f is strictly convex if 1)

2)' for any $x \neq y \in S$ and any $\theta \in (0, 1)$,

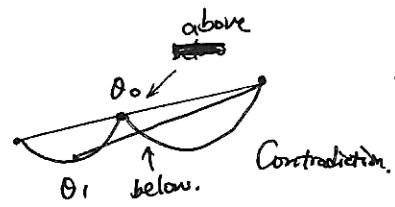
$$f(\theta x + \bar{\theta} y) < \theta f(x) + \bar{\theta} f(y).$$

[Proposition] Let f be convex.

If $f(\theta x + \bar{\theta} y) = \theta f(x) + \bar{\theta} f(y)$ for some $\theta_0 \in (0, 1)$,

then $\forall \theta \in [0, 1]$, $f(\theta x + \bar{\theta} y) = \theta f(x) + \bar{\theta} f(y)$

i.e. $g(\theta) = f(\theta x + \bar{\theta} y)$ is an affine function.



Def. f is (strictly) concave if $-f$ is (strictly) convex.

An affine function is both convex and concave, and is neither strictly convex nor strictly concave.

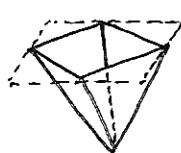
e.g. $f(x) = e^{ax}$ is convex, and strictly convex when $a \neq 0$.

$f(x) = \log_a x$ is $\begin{cases} \text{concave} & (a > 0) \\ \text{convex} & (a < 0) \end{cases}$

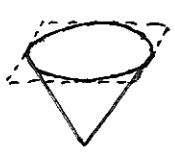
$f(x) = x^a$ is $\begin{cases} \text{convex on } \mathbb{R}^+ & \text{for } a \geq 1 \text{ or } a \leq 0 \\ \text{concave on } \mathbb{R}^+ & \text{for } 0 < a < 1 \end{cases}$



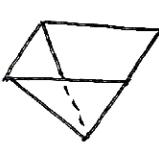
$f(x) = \|x\|$ is convex.



1-norm



2-norm



∞ -norm

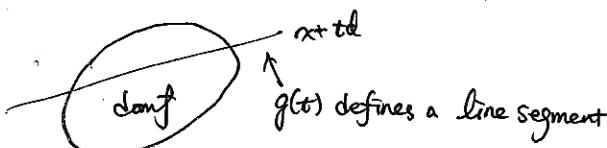
not strictly convex:

Let $x=0, y \neq 0$

$$\| \theta x + \bar{\theta} y \| = \| \bar{\theta} y \| = \bar{\theta} \| y \| = \theta \| x \| + \bar{\theta} \| y \| \quad \square$$

[Proposition] f is convex iff. $\forall x \in \text{dom } f, \text{ any direction } d, (\neq 0)$

$g(t) = f(x + t d)$ is convex on $\text{dom } g := \{t : x + t d \in \text{dom } f\}$



Proof. \Leftarrow : $x, y \rightarrow$ find line segment $\rightarrow g(t^*)$, convex.

\Rightarrow : Obvious

Prof. $\Rightarrow: \forall t_1, t_2 \in \text{dom}f, \theta \in [0, 1]. \bar{t} := \theta t_1 + \bar{\theta} t_2 \in \text{dom}f$

We have $x_1 = x + t_1 d \in \text{dom}f, x_2 = x + t_2 d \in \text{dom}f, \bar{x} = x + \bar{t} d \in \text{dom}f$

$$\text{iff. } (\because \bar{x} = x + (\theta t_1 + \bar{\theta} t_2) d = \theta(x + t_1 d) + \bar{\theta}(x + t_2 d)) \in \text{dom}f \quad \in \text{dom}f$$

$$g(\bar{t}) = f(x + \bar{t} d) = f(\theta x_1 + \bar{\theta} x_2) \leq \theta f(x_1) + \bar{\theta} f(x_2) = \theta g(t_1) + \bar{\theta} g(t_2). \quad \square$$

$\Leftarrow:$ Likewise.

Def. Extended-Value Extension

Given convex $f: S \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$, its extended-value extension

$$\tilde{f}: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}. \quad \tilde{f}(x) = \begin{cases} f(x), & x \in S \\ \infty, & x \notin S. \end{cases}$$

with extended arithmetic and ordering: $\begin{cases} a + \infty = \infty \text{ for } a > -\infty \\ a \cdot \infty = \infty \text{ for } a > 0. \\ 0 \cdot \infty = 0 \end{cases}$

- Effective domain: $\text{dom} \tilde{f} = \text{dom} f = S$

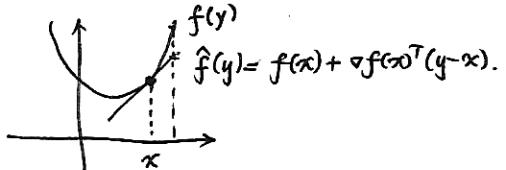
- f is convex iff. \tilde{f} is convex.

- * We can also extend a concave function using $-\infty$.

• First-order condition for convexity

Thm. A differentiable f with an open convex domain $\text{dom}f$ is convex

iff. $f(y) \geq f(x) + \nabla f(x)^T (y - x) \quad \forall x, y \in \text{dom}f$



* Actually Taylor approximation underestimates a convex function.

e.g. $f(x) = e^x$.

\rightarrow To prove $\forall x \forall y \quad e^y \geq e^x + e^x(y - x) \rightarrow e^{y-x} \geq 1 + y - x$. i.e. $e^z \geq 1 + z$

$$g(z) = e^z - z - 1 \geq g(0) = 0 \quad (\underset{-\infty}{z \rightarrow \infty}, \underset{+\infty}{g(z) \rightarrow +\infty} \text{ coercive, global min exists}), \quad g'(z) = e^z - 1 = 0 \Rightarrow z = 0$$

Prof. \Rightarrow Assume f is convex. $d \triangleq y - x. \bar{t} := 1 - t$.

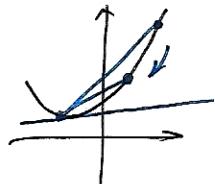
$$f(x + td) = f(x + t(y - x)) = f(ty + \bar{t}x) \leq t f(y) + \bar{t} f(x) \quad t \in (0, 1)$$

$$\text{i.e. } f(x + td) - f(x) = t f(y) - t f(x)$$

$$\Rightarrow \lim_{t \rightarrow 0} \frac{f(x + td) - f(x)}{t} \leq f(y) - f(x)$$

$$\nabla f(x)^T d. \quad \text{i.e. } \nabla f(x)^T (y - x) \leq f(y) - f(x).$$

$\frac{f(x + td) - f(x)}{t \|d\|}$ is actually the slope of the secant line through x and $x + td$.



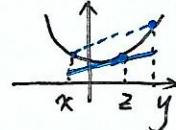
\Leftarrow . Assume the first-order condition holds.

$$z = \theta x + \bar{\theta} y \quad \text{then} \quad f(y) \geq f(z) + \nabla f(z)^T (y-z) \quad \textcircled{1}$$

$$f(x) \geq f(z) + \nabla f(z)^T (x-z) \quad \textcircled{2}$$

$$\begin{aligned} \bar{\theta} \cdot \textcircled{1} + \theta \cdot \textcircled{2} \text{ yields } \theta f(x) + \bar{\theta} f(y) &\geq f(z) + \nabla f(z)^T (\bar{\theta} y + \theta x) - \nabla f(z)^T z \\ &= f(z) = f(\theta x + \bar{\theta} y) \end{aligned}$$

i.e. f is convex.



□

Thm. $\dots \text{dom } f$ is strictly convex

$$\text{iff. } f(y) > f(x) + \nabla f(x)^T (y-x) \quad \forall x \neq y \in \text{dom } f.$$

Proof. ~~Likewise~~ Likewise.

$$\Rightarrow: d = y - x.$$

$$\begin{aligned} \text{If } 0 < t < s < 1, \quad f(x+td) &< \frac{t}{s} f(x+sd) + (1-\frac{t}{s}) f(x) \quad \text{strictly convex.} \\ \Rightarrow \frac{f(x+td) - f(x)}{t} &< \frac{f(x+sd) - f(x)}{s} < f(x+d) - f(x) \\ t \rightarrow 0, \text{ yield } \nabla f(x)^T d &\leq \frac{f(x+sd) - f(x)}{s} < f(y) - f(x). \end{aligned}$$

↑ Why we need "s"? To maintain " $<$ ".

□

[Corollary]. f : differentiable. $I \rightarrow \mathbb{R}$ defined on an open interval $I \subset \mathbb{R}$ is ~~convex~~ (strictly) convex iff. f' is (strictly) increasing on I .

Proof. $\Leftarrow: \forall x, y \in I, x < y$. By Mean Value Thm. can find $c \in (x, y)$ s.t. $f(y) - f(x) = f'(c)(y-x)$

$$\text{Since } \therefore f \text{ convex } (\text{By First-order condition}) \Rightarrow f'(x)(y-x) \geq f'(x)(y-x) \geq 0$$

$$\Rightarrow f(y) \geq f(x) + f'(x)(y-x), \quad f(x) \geq f(y) + f'(y)(x-y).$$

$$f'(x) \leq \frac{f(y) - f(x)}{y-x} \leq f'(y).$$

□

Optimality of stationary points

[Corollary.] $\nabla f(x^*) = 0$ is a global minimum if f is convex.

Proof. Obvious. $f(y) \geq f(x^*) + \underbrace{f'(x^*)(y-x)}_{=0} = f(x^*)$ unique global minimum \Rightarrow strict convex \Rightarrow global minimum

□

Concave \rightarrow global maximum.

• Second-order Condition for Convexity

Thm. A twice differentiable f with an open convex domain $\text{dom } f$ ~~continuous~~

continuously

is convex iff. ~~$\nabla^2 f(x) \succeq 0$~~ $\nabla^2 f(x) \succeq 0$ at every $x \in \text{dom } f$.

Proof. \Rightarrow : Assume f is convex,

$g(t) = f(x+td)$, is convex $\Rightarrow g'(t) \uparrow \Rightarrow g''(t) \geq 0$ for every $t \in \text{domain}$,

i.e. $d^T \nabla^2 f(x) d = g''(0) \geq 0$. for any $d \Rightarrow \nabla^2 f(x) \succeq 0$.

\Leftarrow Also use $g(t)$. Likewise. □

* Thm. f is strictly convex if $\nabla^2 f(x) > 0$. Proof is Likewise.

e.g. $f(x) = x^4$ is strictly convex, while $f''(x) = 0 @ x=0$.

* For Quadratic $f(x) = x^T Q x + b^T x + c$, it is necessary. Can be proved.

e.g. $f(x) = \log\left(\sum_{i=1}^n e^{x_i}\right)$. \rightarrow softmax. Smoothly approximates $\max_{1 \leq i \leq n} x_i$.

$$\begin{aligned} \log(e^{x_1} + e^{x_2}) &= \log(e^{x_1}(1 + e^{x_2 - x_1})) = x_1 + \log(1 + e^{-x_1}) \approx \begin{cases} 0 & x_1 \gg x_2 \\ \Delta x & x_1 \ll x_2 \end{cases} \\ &= \begin{cases} x_1 & x_1 \gg x_2 \\ x_2 & x_1 \ll x_2 \end{cases} \end{aligned}$$

Show that it is convex.

Prof. $g(t) := f(\vec{x} + t\vec{d})$ \uparrow univariate function $g(t) = \log\left(\sum_{i=1}^n e^{x_i + t d_i}\right) = \log s(t)$

$$g'(t) = \frac{s'(t)}{s(t)}, \quad g''(t) = \frac{s''(t)s(t) - s'(t)^2}{s(t)^2}, \quad g''(0) = d^T \nabla^2 f(x) d.$$

② Suffices to show $g''(0) \geq 0$. for every $x, d \in \mathbb{R}^n$.

$$\begin{aligned} g''(0) &= \frac{\sum_{i=1}^n e^{x_i} d_i^2}{\sum_{j=1}^n e^{x_j}} - \left(\frac{\sum_{i=1}^n e^{x_i} d_i}{\sum_{j=1}^n e^{x_j}} \right)^2 = \sum_{i=1}^n \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} d_i^2 - \left(\sum_{i=1}^n \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} d_i \right)^2 \\ &=: \sum_{i=1}^n p_i d_i^2 - \left(\sum_{i=1}^n p_i d_i \right)^2 \geq 0, \quad \sum_{i=1}^n p_i = 1. \end{aligned}$$

□

$$f(\alpha x + \bar{\alpha} y) \leq \alpha f(x) + \bar{\alpha} f(y)$$

$$f\left(\sum_{i=1}^n \theta_i x_i\right) \leq \sum_{i=1}^n \theta_i f(x_i), \quad \theta_i \geq 0, \sum_{i=1}^n \theta_i = 1.$$

Generalization of

Jensen's Inequality 下下及

$$h(x) := x^2. \quad \text{Then } \sum_{i=1}^n h(d_i) p_i - h\left(\sum_{i=1}^n \theta_i d_i\right) \geq 0. \quad \square$$

• Convexity-preserving Operations & Properties of Convex Functions

Thm. If x^* is a local minimum of a convex function f over S , then x^* is also a global minimum of f over S .

Proof. By Contradiction.

Suppose $\exists x \in S$ s.t. $f(x) < f(x^*)$. $\forall \theta \in (0, 1)$ $f(\theta x + \bar{\theta} x^*) = \theta f(x) + \bar{\theta} f(x^*) < f(x^*)$.

$$\|(\theta x + \bar{\theta} x^*) - x^*\| = \|\theta(x - x^*)\| = \theta \|x - x^*\| \rightarrow 0 \text{ when } \theta \rightarrow 0.$$

On the other side, since x^* is a local minimum. $\exists \delta > 0$. s.t. $\|x - x^*\| < \delta$. $f(x) \geq f(x^*)$.

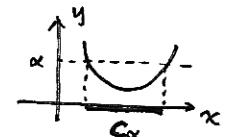
Contradiction.

[Proposition]. Strictly convex $f \rightarrow$ unique global minimum. (By contradiction.)

Def. The α -sublevel set of a function f is: $C_\alpha = \{x \in \text{dom } f : f(x) \leq \alpha\}$.

Thm. Sublevel sets of a convex function are convex.

Similarly, the superlevel set $\{x \in \text{dom } f : f(x) \geq \alpha\}$ of a concave f is convex.

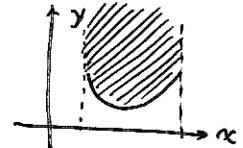


* Not-convex functions can have convex sublevel sets.

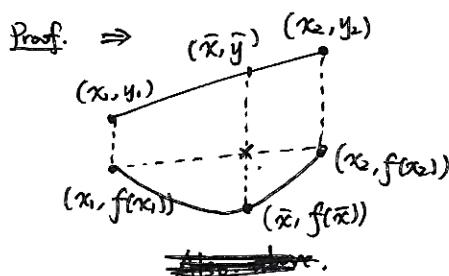
e.g. $f(x) = \sqrt{|x|}$. $f(x) = -e^x$

Def. The epigraph of $f: S \subset \mathbb{R}^n \rightarrow \mathbb{R}$ is: $\text{epi } f = \{(x, y) \in \mathbb{R}^{n+1} : x \in S, y \geq f(x)\}$.
epi = above

f and its extended-value extension \tilde{f} have the same epigraph.



Thm. $f: S \subset \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex function iff. $\text{epi } f$ is a convex set.



\Leftarrow : Consider $(x_1, f(x_1))$, $(x_2, f(x_2))$
 $\text{epi } f$ is convex: $\bar{x} \in \text{epi } f$.
Then $\theta f(x_1) + \bar{\theta} f(x_2) = \bar{y} \geq f(\bar{x})$.

(\bar{x}, \bar{y}) is also above $(\bar{x}, f(\bar{x}))$.

□

* The projection of a convex set C ($\{x : (x, y) \in C \text{ for some } y\}$) is also convex.
② Proof Likewise.

$$(x, y) \in \text{epi } f. \quad y \geq f(x) \geq f(x_0) + \nabla f(x_0)^T (x - x_0)$$

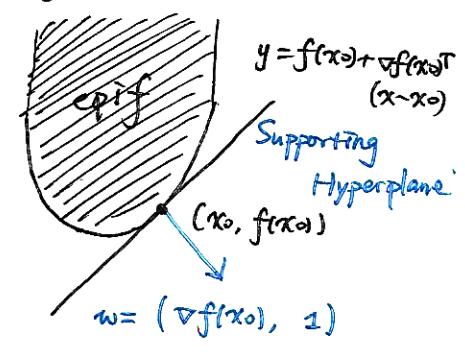
$$\text{i.e. } \nabla f(x_0)^T x - y \leq \nabla f(x_0)^T x_0 - f(x_0)$$

$$\text{i.e. } \begin{pmatrix} \nabla f(x_0) \\ -1 \end{pmatrix}^T \begin{pmatrix} x \\ y \end{pmatrix} \leq \begin{pmatrix} \nabla f(x_0) \\ -1 \end{pmatrix}^T \begin{pmatrix} x_0 \\ f(x_0) \end{pmatrix} \quad \blacksquare$$

$\forall (x, y) \in \text{epi } f$

$\therefore w = (\nabla f(x_0), 1)$

$\therefore w$. A Hyperplane



Thm. Jensen's Inequality

For convex function f , $\forall x_i \in \text{dom } f$. $\theta_i \geq 0$ s.t. $\sum_{i=1}^n \theta_i = 1$.

$$f\left(\sum_{i=1}^n \theta_i x_i\right) \leq \sum_{i=1}^n \theta_i f(x_i).$$

$$\text{e. g. } f(x) = x^2 \Rightarrow \frac{1}{n} \sum_{i=1}^n x_i^2 \leq \sqrt{\frac{1}{n} \sum_{i=1}^n x_i^2}.$$

e.g. $f(x) = \log x$: $\frac{1}{n} \sum_{i=1}^n x_i \geq \sqrt[n]{\prod_{i=1}^n x_i}$
 (concave)

\Rightarrow Hölder's Inequality

Let $p, q \in (1, +\infty)$ be conjugate exponents, i.e. $\frac{1}{p} + \frac{1}{q} = 1$. $\sum_{i=1}^n |x_i y_i| \leq \|x\|_p \|y\|_q$.

$$\underline{\text{Proof.}} \quad x^{\frac{1}{p}} y^{\frac{1}{q}} \leq \frac{1}{p}x + \frac{1}{q}y. \quad (x, y \geq 0) \Leftrightarrow \frac{1}{p} \log x + \frac{1}{q} \log y \leq \log \left(\frac{1}{p}x + \frac{1}{q}y \right).$$

Thus, we have

$$\sum_{i=1}^n \left(\frac{\sum_{j=1}^n |x_j|^p}{\sum_{j=1}^n |y_j|^p} \right)^{\frac{1}{p}} \left(\frac{\sum_{j=1}^n |y_j|^q}{\sum_{j=1}^n |x_j|^q} \right)^{\frac{1}{q}} \leq \frac{\sum_{i=1}^n |x_i|^p}{\sum_{j=1}^n |y_j|^p} + \frac{1}{q} \frac{\sum_{i=1}^n |y_i|^q}{\sum_{j=1}^n |x_j|^q}$$

i.e.

$$\frac{\sum_{i=1}^n |x_i y_i|}{\|x\|_p \|y\|_q} \leq 1 \Rightarrow \sum_{i=1}^n |x_i y_i| \leq \|x\|_p \|y\|_q.$$

\Rightarrow Minkowski's Inequality

$$\text{For } 1 < p < +\infty, \quad \|x+y\|_p \leq \|x\|_p + \|y\|_p.$$

$$\text{Proof: } \|(x+y)\|_p^p = \sum_{i=1}^n |x_i + y_i|^p \leq \sum_{i=1}^n (|x_i| + |y_i|)^{p-1} |x_i| + \sum_{i=1}^n (|x_i| + |y_i|)^{p-1} |y_i|$$

Convexity-preserving Operations:

$\sum_{i=1}^n c_i f_i(x)$, $\bullet f(Ax+b)$, $h(f_1(x), \dots, f_n(x))$, $\sup_{i \in I} f_i(x)$, $\inf_{y \in C} f(x,y)$ are convex
 $(c_i \geq 0)$ $\|g(x)\|$ $\|h(x)\|$ if f_i are convex.

$$\text{Ex. 8. } f(x_1, x_2) = (x_1 - 2x_2)^4 + 2e^{x_1 - 3x_2 + 5}.$$

$$\text{Let } g(y) = y_1^4 + 2e^{y_2} \quad . \quad y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} x_1 - 2x_2 \\ x_1 - 3x_2 + 5 \end{pmatrix} = \begin{pmatrix} 1 & -2 \\ 1 & -3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} 0 \\ 5 \end{pmatrix}.$$

$$f_1 - f_2 ? \quad 0 - f ; \quad x^4 - x^2 .$$

$$f_1/f_2? \quad \frac{f}{-1}; \frac{1}{1+\rho+x}$$

$$f(x) = h(g(x)), \quad h: \mathbb{R} \rightarrow \mathbb{R}, \quad g: \mathbb{R}^n \rightarrow \mathbb{R} \quad \text{Scalar Composition.}$$

- f is convex if h is convex and increasing, g is convex
- if h is convex and ↗ . g is concave
- f is concave if h is concave and ↗ . g is concave
- if h is concave and ↗ . g is convex

$$\left. \begin{array}{l} h^* = h(-x) \\ \text{or} \\ g^* = -g(+x) \end{array} \right\} \begin{array}{l} \text{Transformed} \\ \text{into Case 01.} \end{array}$$

Prof. (CASE 01). Consider $n=1$. $f''(x) = h''(g(x)) [g'(x)]^2 + h'(g(x)) g''(x) \geq 0$.

$\begin{matrix} \checkmark & \checkmark \\ 0 & 0 \\ (\text{convex}) & \end{matrix}$

$\begin{matrix} \checkmark & \checkmark \\ 0 & 0 \\ (\text{increasing}) & \end{matrix}$

Formal Proof. $\forall x, y \in \mathbb{R}^n, \theta \in [0, 1]$. Since g is convex. $g(\theta x + \bar{\theta} y) \leq \theta g(x) + \bar{\theta} g(y)$.

Since h is monotonic,
(increasing) $h(g(\theta x + \bar{\theta} y)) \leq h(\theta g(x) + \bar{\theta} g(y)) \leq \theta h(g(x)) + \bar{\theta} h(g(y))$

↓
 $\because h$ is convex

e.g. To prove $f(x) = \|x\|^2$ is convex.

QED. □

$$\rightarrow g(x) = \|x\|. \quad h(x) = \begin{cases} 0, & x < 0 \\ x^2, & x \geq 0 \end{cases} \quad (\text{Need to be monotonic})$$



Note. If ~~g~~ doing, don't are not \mathbb{R}/\mathbb{R}^n . use extended-value extensions \tilde{g}, \tilde{h} instead.

Vector Composition

Def. $h: \mathbb{R}^m \rightarrow \mathbb{R}$ is increasing if. $x \geq y$ (componentwise) $\Rightarrow h(x) \geq h(y)$.

decreasing

≤ .

Pointwise Maximum

If f_i is convex. $f := \max_{1 \leq i \leq m} f_i(x)$ is convex.

$$\text{dom } f = \bigcap_{i=1}^m \text{dom } f_i$$

$$f_i(x) < \infty \rightarrow \exists x \in \text{dom } f_i \quad (\forall i)$$

Pointwise Supreme: $\sup_{i \in I} f_i(x)$ I could be " \mathbb{R} "

→ We prove $\text{epi } f = \bigcap_i \text{epi } f_i$ Intersection Convex set.

$$\sup_{x \in \mathbb{R}^n} (f(x) + \lambda g(x)) = h(\lambda)$$

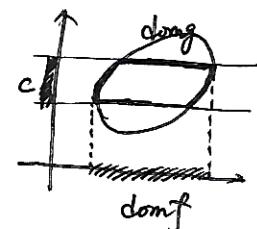
Prove. \uparrow index \uparrow variable \uparrow index \uparrow index

$\rightarrow f$ is convex function.

Partial Minimization

$f: \mathbb{R}^n \times \mathbb{R}^m \rightarrow (-\infty, \infty]$ is convex. $\emptyset = C \subset \mathbb{R}^m$ is convex.

then $f(x) = \inf_{y \in C} g(x, y)$ is convex, provided $f(x) > -\infty \quad (\forall x)$.



$$\text{dom } f := \{x \in \mathbb{R}^n : f(x) = \inf_{y \in C} g(x, y) < +\infty\} = \{x : \exists y \in C, \text{ s.t. } g(x, y) < \infty\}$$

$$= \{x : \exists y \in C \text{ s.t. } (x, y) \in \text{dom } g\} = \{x \in \mathbb{R}^n : (x, y) \in \text{dom } g \cap \mathbb{R}^n \times C\}$$

Proof. $\forall x_1, x_2 \in \text{dom } f, \theta \in [0, 1]$

By definition, $f(\theta x_1 + \bar{\theta} x_2) = \inf_{y \in C} g(\theta x_1 + \bar{\theta} x_2, y) \leq g(\theta x_1 + \bar{\theta} x_2, \theta y_1 + \bar{\theta} y_2)$.

↑ since $\theta y_1 + \bar{\theta} y_2 \in C$.

$\leq \theta g(x_1, y_1) + \bar{\theta} g(x_2, y_2) \quad \text{convex} \quad \theta(f(x_1) + \varepsilon) + \bar{\theta}(f(x_2) + \varepsilon)$

$= \theta f(x_1) + \bar{\theta} f(x_2) + \varepsilon$

↑ $g(x_i, y_i) < f(x_i) + \varepsilon$.
(can find ε)

Convex Optimizations

Std Form for Optimization Probs

$$\min_{\mathbb{R}^n} f(x) \quad \text{s.t.} \quad \begin{aligned} h_i(x) &= 0 \quad (i=1, 2, \dots, k) \\ g_i(x) &\leq 0 \quad (i=1, 2, \dots, m) \end{aligned}$$

$i=1, \dots, k$ $j=1, \dots, m$

Solution Set $X = \{x \in D : h_i(x)=0, g_j(x) \leq 0\}$

$$D = \text{dom } f \cap \bigcap_{i=1}^k \text{dom } h_i \cap \bigcap_{j=1}^m \text{dom } g_j$$

Def. Optimal Value (Generalized Definition)

The optimal value of Problem (P) is $f^* = \inf_{x \in X} f(x)$. allow f^* to take $\pm\infty$.

- * $f^* = \infty$ if (P) is infeasible, i.e. $X = \emptyset$. $\sup \emptyset = -\infty$. $\inf \emptyset = +\infty$.
- * $f^* = -\infty$: unbounded below. $\exists \{x_i\} \in X$ s.t. $f(x_i) \rightarrow -\infty$ as $i \rightarrow +\infty$
- * x^* is an optimal point or solves (P). if $x^* \in X$ and $f(x^*) = f^*$. Note that f^* is not always attainable.
- x_0 is an ϵ -suboptimal if $x_0 \in X$ and $f(x_0) \leq f^* + \epsilon$.
- locally optimal. omitted

Def. Convex Optimization Problem

(P) is a convex optimization problem if. (1) f, g_i are convex functions

(2) h_i are affine functions, i.e. $h_i(x) = a_i^T x - b_i$

- * $\max_x f(x)$ where f is concave is also a convex optimization problem.
 \downarrow
 $\min_x -f(x)$, $-f$ is convex.

We transform the prob. into std form.

$$D = \text{dom } f \cap \bigcap_{i=1}^m \text{dom } g_i \quad (\text{Since } \text{dom } h_i = \mathbb{R}^n)$$

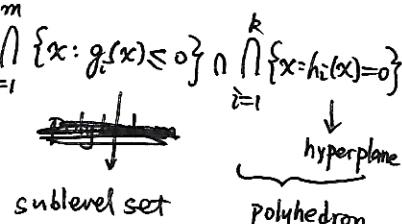
$$X = \{x \in D : g_i(x) \leq 0, 1 \leq i \leq m; h_i(x) = 0, 1 \leq i \leq k\} = D \cap \bigcap_{i=1}^m \{x : g_i(x) \leq 0\} \cap \bigcap_{i=1}^k \{x : h_i(x) = 0\}$$

Both D and X are convex.

$$X_{\text{opt}} := \{x^* \in X : f(x^*) \leq f(x), \forall x \in X\}.$$

Any local minimum is a global minimum.

If f is strictly convex. $|X_{\text{opt}}| \leq 1$.



First-order Optimality Condition

Thm. a convex prob. with differentiable f . $x^* \in X$ is optimal iff. $\nabla f(x^*)^T (x - x^*) \geq 0, \forall x \in X$.

Proof. \Leftarrow : $f(x) \geq f(x^*) + \nabla f(x^*)^T (x - x^*) \geq f(x^*), \forall x \in X$. $\nabla f(x^*)^T (x - x^*) \geq 0$.
 \Updownarrow
 f is convex

\Rightarrow : Assume x^* is optimal. $\alpha \in [0, 1]$. $x^* + \alpha(x - x^*) = \alpha x + \alpha x^* \in X$

$$g(\alpha) := f(x^* + \alpha(x - x^*)) \quad d := x - x^* \text{ is a direction.}$$

$$g(\alpha) - g(0) \geq 0 \rightarrow \lim_{\alpha \rightarrow 0} \frac{g(\alpha) - g(0)}{\alpha} \geq 0. \text{ i.e. } \nabla f(x^*)^T (x - x^*) \geq 0$$

Corollary. $x^* \in \text{int } X$. then $\nabla f(x^*) = 0$

[Prof.] By Thm. $\nabla f(x^*)^T d \geq 0$. $\Rightarrow \nabla f(x^*)^T (-d) \leq 0$

$$\frac{1}{2} \int_{-\infty}^{\infty} e^{-x^2} dx = \frac{1}{2} \sqrt{\pi}$$

$$\nabla f(x^*)^T d = 0 \Rightarrow \nabla f(x^*) = 0$$

(Can choose $d = \nabla f(x^*)$)

$\Leftrightarrow \nabla f(x^*) = 0$ iff. $x^* \in \partial X$.

$$\|\nabla f(x^*)\| < \delta.$$

$$[\text{Prof.}] \quad w = \nabla f(x^*)^\top. \quad w^\top (x - x^*) \leq 0 \Rightarrow w^\top x \leq w^\top x^*$$

Thus, $w^T x = w^T x^*$ is a supporting hyperplane. $\forall x \in \mathcal{X}$

Linear Program

Def. A linear program is an optimization problem of $\min_{\mathbf{x}} \mathbf{c}^T \mathbf{x}$ s.t. $\mathbf{Bx} \leq \mathbf{d}$
 $\mathbf{Ax} = \mathbf{b}$.

Std form: $\min_x C^T x$ st. $Ax = b$
 $x \geq 0$

Inequality form: $\min_{\mathbf{x}} \mathbf{c}^T \mathbf{x}$ s.t. $\mathbf{A}\mathbf{x} \leq \mathbf{b}$

$$[\text{Conversion}] \quad Bx \leq d \rightarrow \underset{\emptyset}{\textcircled{X}} \left\{ \begin{array}{l} Bx + s = d \\ s \geq 0 \end{array} \right. \quad \textcircled{2} \quad x_i = x_i^+ - x_i^- \quad x_i^+, x_i^- \in \mathbb{R}_+$$

$$x^+ = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

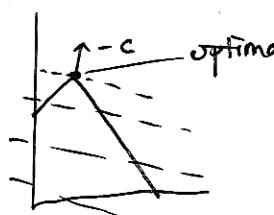
$$x^- = \begin{cases} 0 & \text{if } x \geq 0 \\ -x & \text{if } x < 0 \end{cases}$$

\Rightarrow std form.

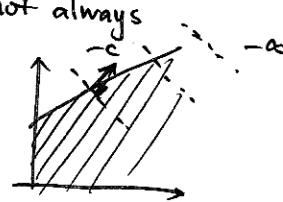
$$[\text{Conversion}] \quad \text{Method 1).} \quad Ax = b \Rightarrow \begin{pmatrix} A \\ -A \end{pmatrix} x \leq \begin{pmatrix} b \\ -b \end{pmatrix}$$

Method 2). Solve one $x_i = f(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_m)$

fewer variables



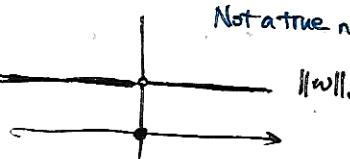
optimal solution exists?
optimal solution unique?
not always
(parallel) 



e.g. $y \in \mathbb{R}^n$. $X \in \mathbb{R}^{n \times p}$. $\text{rank } X = n < p$. Underdetermined Sys.: $Xw = y$. has infinitely many solutions

$$\min_w \|w\|_0 \quad \text{s.t. } Xw = y. \quad \text{where "norm" } \|w\|_0 = \sum_{j=1}^p \mathbb{I}\{w_j \neq 0\}$$

not convex prob.! Why?



Not a true norm

$$\left\{ \begin{array}{l} 1. \quad w_j \neq 0 \\ 0. \quad w_j = 0. \end{array} \right.$$

The ℓ_1 approximation, called basis pursuit is convex.

$$\min_w \|w\|_1 = \sum_{j=1}^p |w_j| \quad \text{s.t. } Xw=y. \quad (\text{not an LP})$$

Transform

$$\textcircled{1} \quad t_j = |w_j| : \quad \min_{t, w} \sum_{j=1}^p t_j \quad \text{s.t. } Xw=y. \quad \frac{t_j = |w_j|}{\text{not affine!}}$$

$$\textcircled{2} \quad \mathbf{1}^T t = \min_{t, w} \sum_{j=1}^p t_j \quad \text{s.t. } Xw=y, \quad t_j \geq |w_j|.$$

$$\text{Why equivalent? } \min_{t, w} f(t, w) = \min_t \min_w f(t, w) = \dots$$

$$\Rightarrow \min_{t, w} \mathbf{1}^T t \quad \text{s.t. } Xw=y, \quad -t \leq w \leq t.$$

e.g. Piecewise Linear Minimization

$$\min_x f(x) = \max_{1 \leq i \leq m} a_i^T x + b$$

$$\Rightarrow \min_{x, t} t \quad \text{s.t. } t \geq a_i^T x + b \quad (1 \leq i \leq m)$$

Quadratic Program

$$\text{Def. } \min_x \frac{1}{2} x^T Q x + c^T x \quad \text{s.t. } Bx \leq d, \quad \begin{matrix} \text{affine} \\ \downarrow \end{matrix} \quad Ax=b \quad (Q \succeq 0) \quad = \text{QP}$$

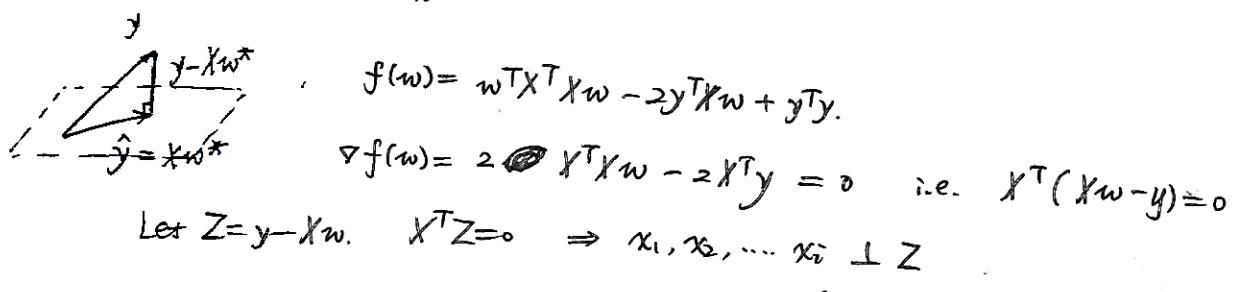
convex when Q is positively semidefinite

reduce to LP when $Q=0$.

$$\text{Def. } \min_x \frac{1}{2} x^T Q x + c^T x \quad \text{s.t. } \frac{1}{2} x^T Q_i x + c_i^T x + d_i \leq 0, \quad Ax=b \quad \text{Quadratically Constrained QP (QCQP)}$$

convex iff. Q, Q_i are positively semidefinite. reduce to QP when all $Q_i=0$.

e.g. Linear Least Squares $\min_w \|y - Xw\|_2^2$.



i.e. $y - Xw^*$ is perpendicular to the column space of X .

CASE 01. $\text{rank } X = p$ (full column rank). $X^T X > 0$. unique solution $w^* = (X^T X)^{-1} X^T y$.

CASE 02. $\text{rank } X = r < p$. $X = (X_1, X_2)$ $\hat{y} = Xw = \begin{pmatrix} X_1 & X_2 \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} = X_1 w_1 + X_2 w_2$

There is a solution with the last $(p-r)$ components being 0.

Thus, if w_1^* solves $\min_{w_1 \in \mathbb{R}^r} \|y - X_1 w_1\|$, $\begin{pmatrix} w_1^* \\ 0 \end{pmatrix}$ solves $\min_{w \in \mathbb{R}^p} \|y - Xw\|$. $w^* = (X_1^T X_1)^{-1} X_1^T y$.
<not unique> \hat{y} is unique

General Unconstrained QP

$$\min_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c \rightarrow \nabla f(\mathbf{x}) = \mathbf{Q} \mathbf{x} + \mathbf{b} = 0, \quad (\mathbf{Q} \succ 0)$$

① $\mathbf{Q} \succ 0$. $\mathbf{x}^* = \mathbf{Q}^{-1} \mathbf{b}$

$\mathbf{Q} \mathbf{x} = -\mathbf{b}$. $-\mathbf{b} \in \text{column space of } \mathbf{Q}$
 $\Rightarrow \mathbf{b} \in -$

② $\det \mathbf{Q} = 0$. $\mathbf{b} \in \text{column space of } \mathbf{Q}$. ~ infinitely many solutions

③ $\det \mathbf{Q} = 0$. $\mathbf{b} \notin \text{column space of } \mathbf{Q}$. ~ no solution $f^* = -\infty$.

e.g. $\mathbf{Q} = \text{diag}(\lambda_1, \lambda_2, 0) \rightarrow f(\mathbf{x}) = \frac{1}{2} \lambda_1 x_1^2 + \frac{1}{2} \lambda_2 x_2^2 + b_1 x_1 + b_2 x_2 + \underline{\underline{b_3 x_3}}$ unconstrained
 $b_3 \neq 0$
 $x_3 \rightarrow -\infty, f(\mathbf{x}) \rightarrow -\infty$.

\mathbf{Q} is not diagonal. \rightarrow use orthogonal matrix!

$$\mathbf{Q} = \underbrace{(\mathbf{u} \Lambda \mathbf{u}^T)}_{\mathbf{u}^T \mathbf{u} = E} \Rightarrow \mathbf{y} = \mathbf{u}^T \mathbf{x} \quad \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} = \frac{1}{2} \mathbf{x}^T \mathbf{u} \Lambda \mathbf{u}^T \mathbf{x} = \frac{1}{2} \mathbf{y}^T \Lambda \mathbf{y}$$

$$\text{Thus, } "f(\mathbf{x})" \rightarrow "g(\mathbf{y})" = \frac{1}{2} \mathbf{y}^T \Lambda \mathbf{y} + \underbrace{b^T \mathbf{u} \mathbf{y} + c}_{\mathbf{b}^T}$$

$$\det \mathbf{Q} = 0 \Leftrightarrow \det \Lambda = 0.$$

$\frac{(\det \mathbf{u})^2 \det \Lambda}{E}$

$$\mathbf{Q} \mathbf{x} + \mathbf{b} = 0 \Leftrightarrow \Lambda \mathbf{y} + \tilde{\mathbf{b}} = 0$$

$\frac{\mathbf{u}^T \mathbf{u} \Lambda \mathbf{u}^T \mathbf{u} \mathbf{y} + \mathbf{u}^T \mathbf{b}}{E}$

(Since $\mathbf{u}^T \mathbf{u} = E$)

$$\text{For ③. } \exists i_0. \lambda_{i_0} = 0, \tilde{b}_{i_0} \neq 0 \rightarrow f(\mathbf{x}) = \left(\sum \text{rest} \right) + \tilde{b}_{i_0} y_{i_0} \xrightarrow{y_{i_0} \rightarrow -\infty/\infty} -\infty$$

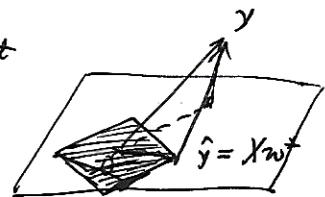
e.g. Lasso (Least Absolute Shrinkage and Selection Optimization)

$$\text{Given } \mathbf{y} \in \mathbb{R}^n, \mathbf{X} \in \mathbb{R}^{n \times p}, t > 0, \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \text{ s.t. } \|w\|_1 \leq t$$

Ridge Regression



$$\min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \text{ s.t. } \|w\|_2^2 \leq t$$



unique solution if $n \geq p$ and $\mathbf{X}^T \mathbf{X} \succ 0$.

CASE 1. Constrain $\rightarrow "(x_0)"$

\Rightarrow On boundary.

A/B/C/D/...

CASE 2. The global optimal of

objective function is

inside the constrained

space

\Rightarrow 约束无用.

CASE 3.

unconstrained optimum

$$\langle \nabla f(x^*), x - x^* \rangle \geq 0.$$

all:
acute angle! (锐角!)

$$\nabla f(x^*)$$

only possible direction for $\nabla f(x^*)$

Ridge: More difficult to satisfy

Lasso: easier to satisfy

Geometric Program

Def. A monomial : $f(x) = \gamma x_1^{a_1} x_2^{a_2} \dots x_n^{a_n}$ ($\gamma > 0$, $x \in \mathbb{R}_{++}^n$)

A posynomial $f(x) = \sum_{k=1}^p \gamma_k x_1^{a_{k1}} x_2^{a_{k2}} \dots x_n^{a_{kn}}$

Geometric Program : $\min_x f(x)$ s.t. $g_i(x) \leq 1$, $\gamma_j x_1^{G_1} x_2^{G_2} \dots x_n^{G_n} = 1$
 $\Leftrightarrow \min \log(\sum e^D)$ s.t. $\log(\sum e^D) \leq 1$. affine.

$$x_k \rightarrow e^{u_k}$$

Gradient Descent

Consider $\min_x f(x)$ where f is convex and differentiable on \mathbb{R}^n .
 $\nabla f(x^*) = 0$.

Descent Method

choose $x_0 \in \mathbb{R}^n$
repeat
descent direction $d_k \in \mathbb{R}^n$, step size $t_k > 0$
 $x_{k+1} \leftarrow x_k + t_k d_k$ s.t. $f(x_{k+1}) < f(x_k)$
until (stopping criterion is satisfied.)

2. descent direction.

- if $g(t) := f(x_k + t d_k) < f(x_k) = g(0)$ for all small enough t . we call d_k a ~.
- d_k is a descent direction $\Rightarrow g'(0) = d_k^T \nabla f(x_k) \leq 0$
 - $g'(0) = d_k^T \nabla f(x_k) < 0 \Rightarrow d_k$ is a descent direction.

Choose $d_k = -\nabla f(x_k)$. $d_k^T \nabla f(x_k) = -\|\nabla f(x_k)\|_2^2 < 0$. unless $\nabla f(x_k) = 0$.
(fastest rate of descent)
i.e. $x_{k+1} \leftarrow x_k - t_k \nabla f(x_k)$.

$x \leftarrow x_0 \in \mathbb{R}^n$
while $\|\nabla f(x)\| < \delta$. do
 $x \leftarrow x - t \nabla f(x)$
end while
return x

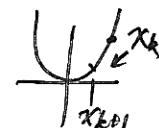
practical then $\nabla f(x) = 0$
also: $\|f(x_{\text{new}}) - f(x_{\text{old}})\| \leq \delta$
 $\frac{\|f(x_{\text{new}}) - f(x_{\text{old}})\|}{\|f(x_{\text{old}})\|} < \delta$, etc.

How to choose a decent "t" (step)?

e.g. $f(x) = \frac{1}{2} ax^2$. ($a > 0$) \Rightarrow ~~$x_{k+1} = x_k - t f'(x_k) = (1-ax) x_k$~~

$$(1-ax)^k x_k \rightarrow 0$$

$$|1-ax| < 2$$



Note: f satisfies $|f'(x) - f'(y)| = \alpha|x-y| \rightarrow$ Lipschitz continuity (f').
 $f''(x) = \alpha$ Lipschitz constant

Def. Lipschitz continuity: $\|f(x) - f(y)\| \leq L\|x-y\| \quad \forall x, y$ "L-Lipschitz"
 || implies usually 2-norm.

Uniform continuity $\forall x. \exists \delta(\varepsilon). \text{ s.t. } \|f(x) - f(y)\| < \varepsilon. \forall y \text{ s.t. } \|x-y\| < \delta.$
一致連續

e.g. $f(x) = \alpha^T x$ is $\|\alpha\|$ -Lipschitz. $\|\alpha^T x - \alpha^T y\| = \|\alpha(x-y)\| \leq \|\alpha\| \|x-y\|$
 $f(x) = \frac{1}{2}\|x\|^2$ is $\|\alpha\|$ -Lipschitz.

Cauchy-Schwarz Inequality

~~Q. $f(x) = Qx$ with $Q \geq 0$. is $\lambda_{\max}(Q)$ -Lipschitz.~~

Proof: $f(x) - f(y) = Q(x-y) = Qd. \quad (d := x-y)$

$$\|f(x) - f(y)\|^2 = (Qd)^T (Qd) = d^T Q^T Q d \leq \lambda_{\max}(Q^T Q) \|d\|^2$$

$$\underbrace{Q \geq 0}_{\text{symmetric.}} \rightarrow Q^T = Q. \quad \overbrace{Qx = \lambda x}^{\text{symmetric.}} \quad = (\lambda_{\max}(Q) \|d\|)^2$$

Lipschitz constant

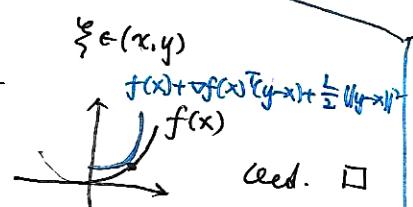
For general cases, $f(x) = Qx$ is $\sqrt{\lambda_{\max}(Q^T Q)}$ -Lipschitz.

Def. L-smoothness: ∇f is L-Lipschitz. ← upper bound the rate of change of ∇f .

e.g. for convex (twice cont. diff.) $f: \mathbb{R}^n \rightarrow \mathbb{R}$. L-smooth iff. $\nabla^2 f(x) \preceq L E$.

f is L-smooth $\Rightarrow f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|y-x\|^2.$ $L \text{ diag}(1, \dots, 1)$

Proof: $f(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2}(y-x)^T \nabla^2 f(\xi)(y-x)$
 $\leq f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} \lambda_{\max}(\nabla^2 f(\xi)) \|y-x\|^2$
 $\leq f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|y-x\|^2.$



Another proof. Prove D case.

$$g(t) := f(x + t(y-x)). \quad g(1) = g(0) + g'(0) + \frac{1}{2} \Delta g.$$

$$|g'(t) - g'(s)| = |(\nabla f(x + t(y-x)) - \nabla f(x + s(y-x)))^T (y-x)|$$

$$\leq L \|y-x\| \|y-x\| |t-s| \rightarrow \Delta g = L \|y-x\|^2.$$

Back to

Gradient Descent: L-smooth f.

$$\{x_k\}: \quad f(x_{k+1}) \leq f(x_k) - t \left(1 - \frac{L^2}{2}\right) \|\nabla f(x_k)\|^2$$

$$(x_{k+1} = x_k - t \nabla f(x_k)) \Rightarrow f(x_k) - f(x_{k+1}) \geq t \left(1 - \frac{L^2}{2}\right) \|\nabla f(x_k)\|^2$$

As long as $0 < t \leq \frac{1}{L}$ $f(x_k) - f(x_{k+1}) \geq \frac{t}{2} \|\nabla f(x_k)\|^2 \geq 0$

always descent! \rightarrow which is good

f = convex & L -smooth. x^* = a minimum of f . for step size $t \in (0, \frac{1}{L}]$

$\{x_k\}$ produced by gradient descent satisfies

$$f(x_k) - f(x^*) \leq \frac{\|x_0 - x^*\|^2}{2tk}.$$

Note.

- $f(x_k) \downarrow f^*$ as $k \rightarrow +\infty$.
- $\lim_{k \rightarrow +\infty} x_k =: \hat{x}_k$. \hat{x}_k is a optimal solution.
- Problem: $f(x_k) - f(x^*) \leq \frac{C}{k}$. difference $< \varepsilon \Rightarrow k > \frac{1}{\varepsilon} \cdot \frac{C}{L}$. $k = \frac{10^P}{\varepsilon}$!!
- Faster convergence?

→ Extend to: $t \in (0, \frac{2}{L})$. $f(x_k) - f(x_{k+1}) \geq t \|\nabla f(x_k)\|^2 \geq 0$.

$$f(x_k) - f(x^*) \leq \frac{\|x_0 - x^*\|^2}{4tk}$$

- $\|x_0 - x^*\|^2$ smaller. might help.

Proof.

$$\begin{aligned} \|x_{k+1} - x^*\|^2 &= \|(x_k - x^*) - t \nabla f(x_k)\|^2 = \|x_k - x^*\|^2 + t^2 \|\nabla f(x_k)\|^2 - 2t \nabla f(x_k)^T (x_k - x^*) \\ f(x_k) - f(x_{k+1}) &\geq \frac{t}{2} \|\nabla f(x_k)\|^2. \end{aligned}$$

First-order condition of Convexity: $f(x^*) - f(x_k) \geq \nabla f(x_k)^T (x^* - x_k)$

$$\begin{aligned} \text{Thus, } \|x_{k+1} - x^*\|^2 &\leq \|x_k - x^*\|^2 + 2t [f(x_k) - f(x_{k+1})] + 2t [f(x^*) - f(x_k)] \\ &= \|x_k - x^*\|^2 + 2t [f(x^*) - f(x_{k+1})] \end{aligned}$$

$$\sum_{k=0}^{N-1} f(x_{k+1}) - f(x^*) \leq \sum_{k=0}^{N-1} \frac{1}{2t} (\|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2)$$

$$\text{i.e. } (N-t)(f(x_N) - f(x^*)) \leq \frac{1}{2t} (\|x_0 - x^*\|^2 - \|x_N - x^*\|^2) \leq \frac{1}{2t} \|x_0 - x^*\|^2$$

$$\text{i.e. } f(x_N) - f(x^*) \leq \frac{1}{2Nt} \|x_0 - x^*\|^2.$$

QED. □

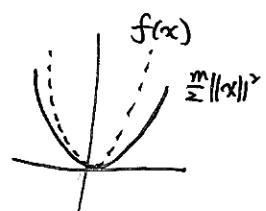
Faster Convergence? — Strong Convexity

Def. Strong Convexity.

$f(x)$ is strongly convex with m (m -strongly convex) if

$$\tilde{f}(x) = f(x) - \frac{m}{2} \|x\|^2 \text{ is convex.}$$

e.g. x^* is not strongly convex (for any m)



→ Def. First-order Condition for Strong Convexity

f is m -strongly convex iff. $f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{m}{2} \|x-y\|^2 \quad \forall x, y$

strong convexity \Rightarrow strict convexity \Rightarrow convexity



* m -strong convexity and L -smoothness together imply $\frac{m}{2} \|x-y\|^2 \leq f(y) - f(x) - \nabla f(x)^T(y-x) \leq \frac{L}{2} \|x-y\|^2$

First-order condition for strong convexity.

Proof. $\Rightarrow: \tilde{f}(x) = f(x) - \frac{m}{2} \|x\|^2$ convex. $\Rightarrow f(x) = \tilde{f}(x) + \frac{m}{2} \|x\|^2$

$$\tilde{f}(y) \geq \tilde{f}(x) + \langle \nabla \tilde{f}(x), y-x \rangle$$

$$\|y\|^2 = \|y-x\|^2 + \|x\|^2 + 2\langle x, y-x \rangle$$

$$\begin{aligned} \text{Thus, } f(y) &= \tilde{f}(y) + \frac{m}{2} \|y\|^2 \geq \tilde{f}(x) + \frac{m}{2} \|x\|^2 + \frac{m}{2} \|y-x\|^2 \\ &\quad + \langle \nabla \tilde{f}(x) + mx, y-x \rangle \\ &\geq f(x) + \langle \nabla f(x), y-x \rangle + \frac{m}{2} \|y-x\|^2. \end{aligned}$$

$\Leftarrow:$ Obvious. \square

Second-order condition for Strong Convexity

Thm. twice continuously differentiable f . f is m -strongly convex iff. $\nabla^2 f(x) \succcurlyeq mE$. $\forall x$.

$$\text{iff. } \lambda_{\min}(\nabla^2 f) \geq m, \forall x$$

Proof. $\tilde{f}(x) = f(x) - \frac{m}{2} \|x\|^2$ is convex iff $\nabla^2 \tilde{f} = \nabla^2 f - mE \succeq 0$. \square

$$\rightarrow \lambda_{\min} \leq \lambda \leq \lambda_{\max}$$

$$\begin{array}{c} \text{Strong} \\ \uparrow \\ \text{restrict} \end{array} \begin{array}{c} \text{eigen-} \\ \text{value} \end{array} \begin{array}{c} \uparrow \\ \text{L-smoothness} \\ \text{restrict} \end{array}$$

Convergence Restricted by Smoothness and Strong Convexity

Thm. f : m -strongly convex, L -smooth. x^* : minimum of f

step size $t \in (0, \frac{1}{L})$. $\{x_k\}$ produced by G. D. satisfies:

- $f(x_k) - f(x^*) \leq \frac{L(1-mt)^k}{2} \|x_0 - x^*\|^2$
- $\|x_k - x^*\|^2 \leq (1-mt)^k \|x_0 - x^*\|^2$.

Note: $f(x_k) \rightarrow f(x^*)$ exponentially fast.

$f(x_k) - f(x^*) \leq \varepsilon \rightarrow \cancel{\text{something}} O(\log \frac{1}{\varepsilon})$. faster convergence.

$$\nabla f(x^*) = 0 \Rightarrow \frac{m}{2} \|x_k - x^*\|^2 \leq f(x_k) - f(x^*) \leq \frac{L}{2} \|x_k - x^*\|^2.$$

(First condition)
-order

Proof. $\|x_{k+1} - x^*\|^2 = \|x_k - t \nabla f(x_k) - x^*\|^2 = \|x_k - x^*\|^2 + t^2 \|\nabla f(x_k)\|^2 + 2t \nabla f(x_k)^T (x_k - x^*)$

By L -smoothness: $t^2 \|\nabla f(x_k)\|^2 \leq 2t (f(x_k) - f(x_{k+1}))$.

By m -strong convexity: $\nabla f(x_k)^T (x^* - x_k) \leq f(x^*) - f(x_k) - \frac{m}{2} \|x_k - x^*\|^2$

$$\Rightarrow \|x_{k+1} - x^*\|^2 \leq (1-mt) \|x_k - x^*\|^2 + 2t [f(x^*) - f(x_{k+1})]$$

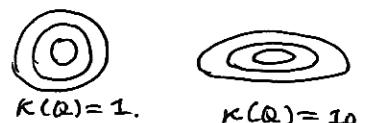
$$f(x^*) < f(x_{k+1}) \Rightarrow \|x_{k+1} - x^*\|^2 \leq (1-mt) \|x_k - x^*\|^2$$

$$\text{Thus, } \|x_k - x^*\|^2 \leq (1-mt)^k \|x_0 - x^*\|^2.$$

Note. 2-smoothness $\Rightarrow f(x_k) - f(x^*) \leq \frac{\|x_k - x^*\|^2}{2t}$. $t \uparrow$ faster.
 Converge $\rightarrow 0 < t \leq \frac{1}{L}$. $\rightarrow t_{\max} = \frac{1}{L}$. fastest descending rate.

Def. Condition number.

$$Q \in \mathbb{R}^{n \times n}, Q \succ 0. \quad \kappa(Q) = \frac{\lambda_{\max}(Q)}{\lambda_{\min}(Q)}.$$



$$\kappa(Q) = 1.$$

well-conditioned

ill-conditioned

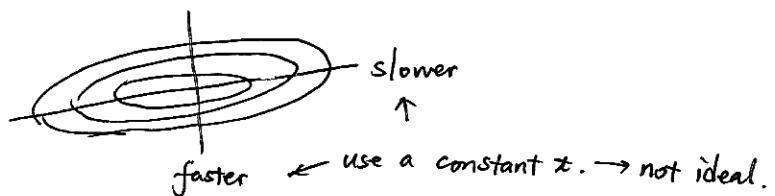
$m \leq \lambda \leq L$, $f(x_k)$ dominates the slower one in $(1-mt)^{2k}$ and $(1-Lt)^{2k}$.

Convergence Analysis.

$$\frac{L}{2} (1 - \frac{m}{L})^k \|x_0 - x^*\|^2 \leq \varepsilon \Rightarrow (1 - \frac{m}{L})^k = \frac{\varepsilon}{c}. \quad 1 - \frac{m}{L} \approx e^{-\frac{m}{L}} \text{ (if } \frac{m}{L} \text{ is small)}$$

$$\text{Thus, } k \approx \frac{L}{m} \log \frac{c}{\varepsilon} \geq \kappa(Q) \log \frac{c}{\varepsilon}.$$

$$\kappa(Q) \leq \frac{L}{m}$$

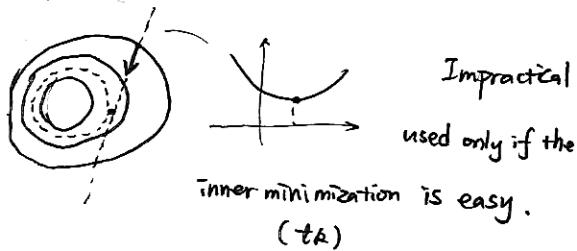


Step Size.

$$x_{k+1} \leftarrow x_k - t_k \nabla f(x_k)$$

- constant step size: t
- exact line search: $t_k = \arg \min_s f(x_k - s \nabla f(x_k))$
- backtracking line search (Armijo's rule): t_k satisfies $f(x_k) - f(x_k - t_k \nabla f(x_k)) \geq \alpha t_k \|\nabla f(x_k)\|^2$.

Exact Line Search



Note. Can prove that

$$0 = h'(t_k) = \nabla (f(x_k) + t_k \nabla f(x_k))^T \nabla f(x_k)$$

$$= \nabla f(x_{k+1})^T \nabla f(x_k)$$

相邻两次的下降率是垂直的。

Thm. exact line search. $f(x_k) - f(x^*) \leq (1 - \frac{m}{L})^k [f(x_0) - f(x^*)]$.

Proof. L-smoothness: $f(x_k - t \nabla f(x_k)) - f(x_k) \leq -t (1 - \frac{L}{2}) \|\nabla f(x_k)\|^2$

$$\downarrow \min_t \quad \downarrow \min_t$$

$$f(x_{k+1}) - f(x_k) \leq -\frac{1}{2} \|\nabla f(x_k)\|^2$$

m -strong convexity: $f(x) \geq f(x_k) + \nabla f(x_k)^T (x - x_k) + \frac{m}{2} \|x - x_k\|^2$.

$$\min_x \quad \min_x$$

$$f(x^*) \geq f(x_k) - \frac{1}{2m} \|\nabla f(x_k)\|^2.$$

Thus, $f(x_{k+1}) - f(x^*) \leq (1 - \frac{m}{L}) [f(x_k) - f(x^*)]$. \square