

# Optimization For ML

Machine learning 2021

Mansoor Rezghi

Regression  $\Rightarrow$

$$\min \left\{ \| \Phi w - y \|_2^2 + \lambda \| w \|_1 \right\}$$

$$\min f(x)$$

$$x \in \mathbb{R}^n$$

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$V \text{ vector space } \mathbb{R}^n$$

Introduction to Matrix Computer

$$x \in V : \left\{ \begin{array}{l} m \times n \\ w, v \in V \end{array} \right\}$$

$$v \in V \quad \alpha v \in V$$

$$x, y \in V \quad \forall \alpha, \beta \in \mathbb{R}$$

$$\alpha x + \beta y \in V$$

$$\|x\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{1/p} \rightarrow \| \cdot \|_p$$

$$p=2 \Rightarrow \|x\|_2 = \sqrt{\sum x_i^2}$$

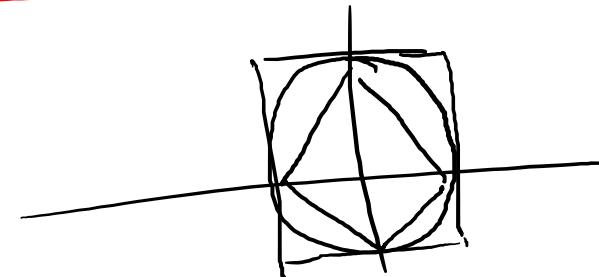
$$p=1 \Rightarrow \|x\|_1 = \sum |x_i|$$

$$p=\infty \Rightarrow \|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$$

$$\|x-y\|$$

$$f_{(m=1, n)}$$

$\|\alpha x\| > 0 \quad \|\alpha x\| = |\alpha| \|x\| \quad \|\alpha x + (1-\alpha)y\| \leq \|\alpha x\| + \|(1-\alpha)y\| \quad \|\alpha x + (1-\alpha)y\| \leq |\alpha| \|x\| + (1-\alpha) \|y\|$

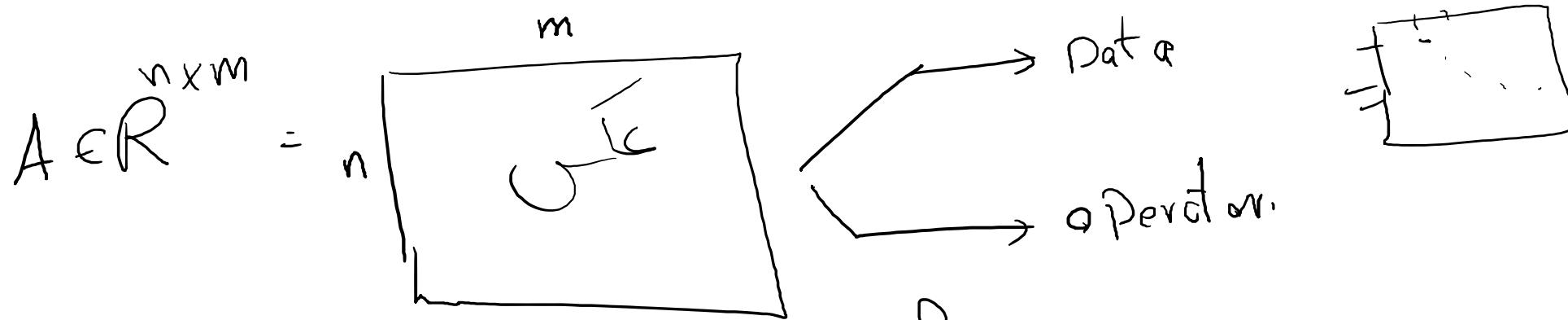


$$\|x\|_p = 1$$

$$\|x\|_p = \sqrt[p]{\left( \sum |x_i|^p \right)} = \max_{1 \leq i \leq n} |x_i|$$

$$\| \alpha x + (1-\alpha)y \| \leq \|\alpha x\| + \|(1-\alpha)y\| \leq |\alpha| \|x\| + (1-\alpha) \|y\|$$

$$f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$$



$$A: \mathbb{R}^m \rightarrow \mathbb{R}^m$$

$$y = Ax = f(x)$$

$$\Delta f(\alpha x + \beta y) = \alpha f(x) - \beta f(y)$$

$$\|A\|_F^2 = \sum_{i,j}^{n,m} a_{ij}^2 = \sum_{j=1}^m \|a_j\|_r^2 = \sum_{j=1}^m \|s_j\|_r^2$$

$$A = \begin{bmatrix} 0 & 0_1 & 0_2 & \dots & a_m \end{bmatrix}$$

$$- \begin{bmatrix} b_1 & s_1 \\ \vdots & \vdots \\ b_n & s_n \end{bmatrix}$$

$$\|A\|_1 = \sum_{i,j}^{n,m} |a_{ij}| \rightarrow \underline{ML}$$

~~Def~~ A As a lin op.

$A: \mathbb{R}^m \rightarrow \mathbb{R}^n$

$$\|A\|_q = \max_{x \neq 0} \frac{\|Ax\|_q}{\|x\|_q}$$

(G.W)  $\tilde{r}_j$

$$= \min_{x \neq 0} \|A(\frac{x}{\|x\|_q})\|_q = \min_{\|x\|_q=1} \|Ax\|_q$$

---

$\{x_1, \dots, x_n\} \Rightarrow$  linear independence  $\checkmark_i$

Linear independent vectors  $\sum \alpha_i x_i = 0 \Rightarrow \alpha_i = 0$

$$\sum \alpha_i x_i = 0 \quad \alpha_i \neq 0 \Rightarrow \alpha_i x_i = - \sum_{i \neq i} \alpha_i x_i \Rightarrow x_i = \cancel{\alpha_i} \sum \left(-\frac{\alpha_i}{\alpha_i}\right) x_i$$

V: Rank(V) = {مدى خطي متعارض, 0}

(R<sup>2</sup>)  $\begin{pmatrix} x \\ y \end{pmatrix} = x \begin{pmatrix} 1 \\ 0 \end{pmatrix} + y \begin{pmatrix} 0 \\ 1 \end{pmatrix}$

dim(R) = 2  $\begin{pmatrix} 0 \\ n \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix}$

$$\alpha \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \beta \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$\alpha = .$   
 $\beta = .$

$$\begin{pmatrix} 0 \\ n \end{pmatrix} = a \begin{pmatrix} 1 \\ 0 \end{pmatrix} + b \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

$\underbrace{\{x_1, \dots, x_n\}}_{\text{Base of } V} \rightarrow \dim(V) = n$

$x_1, \dots, x_n$  indepnt.

$x \in V$   $x = \sum \alpha_i x_i$

$V = \text{Span}\{x_1, \dots, x_n\}$

---

$\begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix}$	$R^2$
$\begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix}$	
$\begin{pmatrix} 1 \\ -1 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix}$	

$$A = [a_1 \dots a_n] \in \mathbb{R}^{m \times n}$$

$$a_i \in \mathbb{R}^m$$

$$= \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_m \end{bmatrix}$$

$A$  مکعبی متعال  $\Leftrightarrow$   $A$  میز متعال

$A$  مکعبی متعال  $\Leftrightarrow$  این سمعی

$$\text{گوئی} = \text{معنی}$$

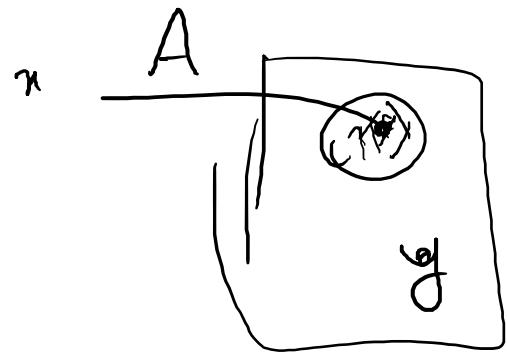
~~Rank~~ - Rank ( $A$ ) = { # of independent columns or rows }

$$\text{Rank } (A) \leq \min\{m, n\}$$

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

$$A: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$R(A) = \left\{ y \in \mathbb{R}^m \mid \exists \alpha \text{ s.t. } A\alpha = y \right\}$$



$$R(A) \subset \mathbb{R}^m$$

$$w \in V$$

$$y = A\alpha = \sum \alpha_i a_i$$

$$R(A) = \text{Span}\{a_1, \dots, a_n\}$$

$$\mathbb{R}^3 \rightarrow \mathbb{R}^3$$

$$\dim(R(A)) = \underline{\text{Rank}}(A) \leq \min\{m, n\}$$

$$A: \mathbb{R}^{m \times n} : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$N(A) = \left\{ \begin{matrix} x \in \mathbb{R}^n \\ \text{such that } Ax = 0 \end{matrix} \right\} \Rightarrow Ax = 0 \Rightarrow \sum x_i a_i = 0$$

columns of  $A$  independent  $\Rightarrow N(A) = \{0\}$

$$\begin{matrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{matrix}$$

$$\mathbb{R}^m = R(A) \oplus N(A^\top)$$

$$\mathbb{R}^n = R(A^\top) \oplus N(A)$$

$$\frac{R(A) \perp N(A^\top)}{\forall x \in R(A) \quad z \in N(A^\top)}$$

$$y^\top z = 0$$

$$y^\top z = (Ax)^\top z = x^\top A^\top z = 0$$

Defini:  $V = V_1 \oplus V_2 \Rightarrow V_1 \cap V_2 = \{0\} \Rightarrow x \in V \Rightarrow \exists! x_1 \in V_1, x_2 \in V_2$

$$V = V_1 + V_2 \quad \forall x \in V \quad \exists x_1 \in V_1, x_2 \in V_2 \quad x = x_1 + x_2$$

$$Ax = y$$

$$y \in R(A)$$

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

$$\text{Rank}(A) = n$$

$$\begin{array}{l} Ax_1 = y \\ Ax_2 = y \end{array} \quad \begin{array}{c} \downarrow \\ A(x_1 - x_2) = 0 \end{array} \quad \rightarrow x_1 = x_2$$

nonsingular.

invertible.

$$\min_{\mathbf{x}} \|Ax - b\|_2^2$$

$$\min_{\mathbf{x}} \|Ax - b_R\|_2^2$$

$$\begin{array}{c} Ax = b_R \\ \cancel{A^T A} \end{array}$$

$$\det(A) \neq 0$$

$$b_R \in R(A)$$

$$b = b_R + b_N$$

$$b_N \in N(A^T)$$

$$\min_x \|Ax - b\|_2^2$$

$$b \in R(A)$$

$$\begin{array}{c} x \\ | \\ A \in \mathbb{R}^{m \times n} \end{array}$$

$$b \in \mathbb{R}^m$$

$$\min_x \|Ax - b_R - b_N\|_2^2 = \|A^T z - b_N\|_2^2 = (z^T - b_N^T)(z - b_N) = \frac{z^T z - 2z^T b_N + b_N^T b_N}{4}$$

$$\|x\|_F^2 = \sum x_i^2 = x^T x$$

$$\|Ax - b_R\|_2^2$$

$$\min \|Ax - b\|_2^2 \iff Ax = b_R \iff Ax - b_R = 0$$

$$A^T(Ax - b_R) = 0 \quad \begin{matrix} \nearrow \\ R(A) \perp N(A^T) \end{matrix}$$

$$\begin{aligned} A^T A x &= A^T b_R + A^T b_N \\ &= A^T(b_R + b_N) \\ &= A^T b \end{aligned}$$

$$A^T A x = A^T b \Rightarrow x = (A^T A)^{-1} A^T b$$

$$\min \|Ax - b\|^2 \iff x = (A^T A)^{-1} A^T b$$

$(A^T A)^{-1} A^T b$

Eigenvalue:

$$\forall \lambda \quad Ax = \lambda x$$

$A: \mathbb{R}^{n \times n}$

# $\exists \lambda$  if  $\exists x \neq 0$  s.t.  $Ax = \lambda x$

eigenvector

$Ax = \lambda x$

eigenvalue.

$$(A - \lambda I)x = 0 \Rightarrow x \in N(A - \lambda I) \Rightarrow \det(A - \lambda I) = 0$$

wishto  $\therefore$  if  $\lambda^r$  is a root  
 $\lambda^r (\lambda - \lambda_1)(\lambda - \lambda_2) \dots$

$$Bx =$$

$$\text{Tr}_{nc}(A) = \sum a_{ii}$$

$$\min \text{Tr}_{nc}(X^T A X) = \sum_{i=n-k+1}^n \lambda_i$$

A SPD

$$\leftarrow \text{s.t. } X^T X = I$$

$$A: \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$$

$$A \in \mathbb{R}^{n \times n} \quad X \in \mathbb{R}^{n \times k}$$

the rank of  $X$  is  $k$ .

$$\max_{\mathbf{x} \in \mathbb{R}^n} \text{Tr}(X^T A X) = \sum_{i=1}^k \lambda_i$$

$\lambda_1 > \lambda_2 > \dots > \lambda_n$

$$X^T X = I$$

$$X = [x_1, x_2, \dots, x_n]$$

$\lambda_i$  ← Corresponding eigenvectors  $x_i$ ;  $i=1-k$

---

$$\alpha \in \mathbb{R} \text{ s.t. } \alpha x = \alpha x^2 \geq 0$$

$\forall x \neq 0$

$A \in \mathbb{R}^{n \times n}$  semi-positive definite if  $\forall x \neq 0 \quad X^T A X \geq 0$

Positive definite  $\Leftrightarrow X^T A X > 0$

$$A \text{ SPD} \Leftrightarrow \lambda_i > 0$$

$X^T A X \geq 0$

$$A \text{ SPD} \rightarrow A x = \lambda x$$

$$\underline{X^T A X = \lambda}$$

Text:

Numerical Optimization Nocedal, springer 2006, Chap 2,3

Convex optimization, Boyd, [Cambridge University Press](#)

Deep Learning, chap4, chap 8

KM(section 13.3.2)

~~min<sub>x</sub> f(x)~~

min f(x)

$\arg \min_x f(x) \Leftrightarrow \arg \min_x -f(x)$

$\min_x f(x)$

min f(x)

$$\left[ \begin{array}{l} \sum \| \phi w - y_i \|_r \\ \sum \| \phi w - y_i \|_1 \end{array} \right] + \lambda \| w \|_2$$

- Unconstraint Optimization

logistic Regression, Perceptron, ...

- Constraint Optimization

SVM Classifier

$$\min f_0(x)$$

s.t.

$$f_i(x) \leq 0$$

$$h_i(x) = 0$$

$y \cdot (1)$ 

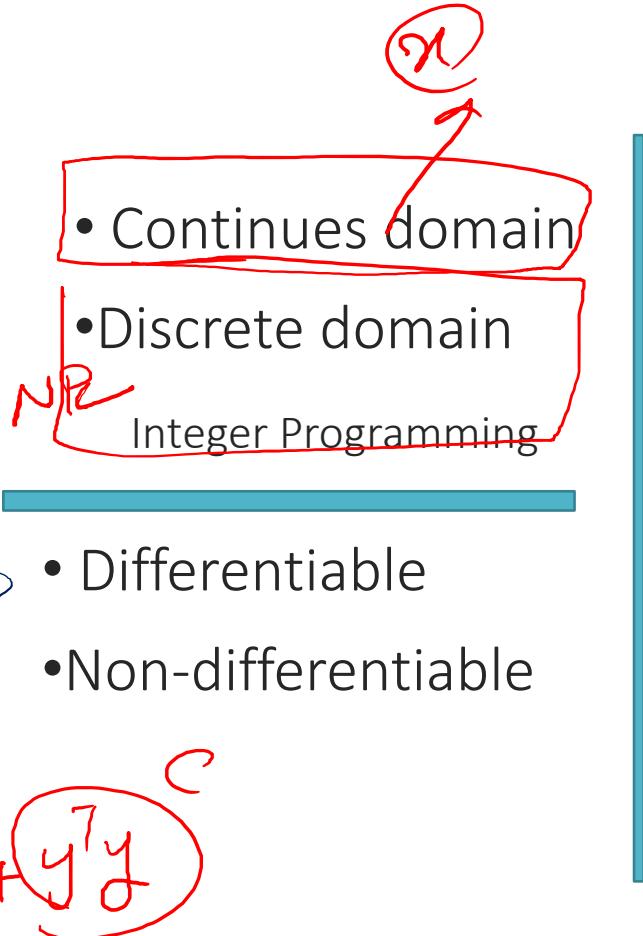
$$\begin{bmatrix} w_1 & \dots & w_d \end{bmatrix} \begin{bmatrix} \$1 \\ \$2 \\ \vdots \\ \$n \end{bmatrix} = \sum_{i=1}^n w_i \cdot \$i$$

$x \in R$   
 ~~$x \in \mathbb{N}, \mathbb{C}$~~

# Objective function and constraints

- Linear  
Linear Programming (LP)
- Non-linear
- Special case: Quadratic

$$\begin{aligned}
 & \min \| \Phi w - y \|^2_2 \\
 &= (\Phi w - y)^T (\Phi w - y) \\
 &= w^T \underbrace{\Phi^T \Phi}_A w + 2 y^T \underbrace{\Phi w}_b + \underbrace{y^T y}_c \\
 &= w^T A w + b^T w + c
 \end{aligned}$$



- Convex
- Non-Convex

$$f(w) = w^T A w + b^T w + c$$

$$= \underbrace{w^T w}_0 + \underbrace{w^T b}_0 + \underbrace{c}_0$$

# Un-constraint Constraint and optimization

- Unconstraint Optimization

$$\text{minimize} \quad f_0(x)$$

- Constraint optimization

$$\begin{aligned}\text{minimize} \quad & f_0(x) \\ \text{subject to} \quad & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & h_i(x) = 0, \quad i = 1, \dots, p\end{aligned}$$

$$f(\alpha x_1 + (1-\alpha)x_2) \leq \alpha f(x_1) + (1-\alpha)f(x_2)$$

$\theta$

# Convex Optimization

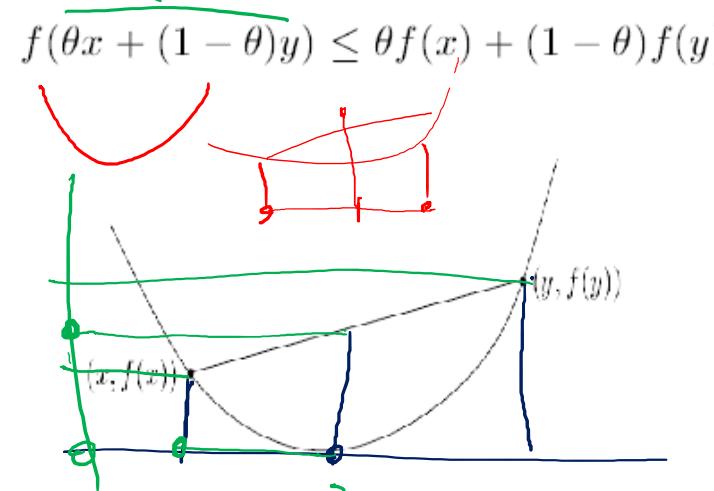
minimize  $f_0(x)$   
 subject to  $f_i(x) \leq 0, i = 1, \dots, m$   
 $h_i(x) = 0, i = 1, \dots, p$

- $f_i, i=0, \dots, m$  convex function

- $h_i, i=1, \dots, p$  Affine function

$$h_i(y) = \mathbf{A}_i^T y + b_i$$

A function  $f : \mathbf{R}^n \rightarrow \mathbf{R}$  is *convex* if  $\text{dom } f$  is a **convex set** and if for all  $x, y \in \text{dom } f$ , and  $\theta$  with  $0 \leq \theta \leq 1$



$$\theta x_1 + (1 - \theta)x_2$$

$\theta$

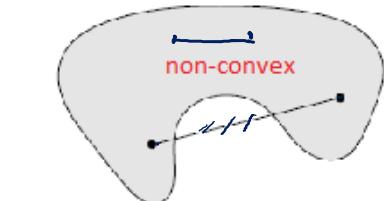
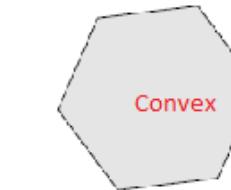
$$x_1 \quad x_2$$

$\theta x_1 + (1 - \theta)x_2 \in C$

$0 < \theta < 1$

A set  $C$  is *convex* if the line segment between any two points in  $C$  lies in  $C$ , i.e., if for any  $x_1, x_2 \in C$  and any  $\theta$  with  $0 \leq \theta \leq 1$

$$\theta x_1 + (1 - \theta)x_2 \in C.$$



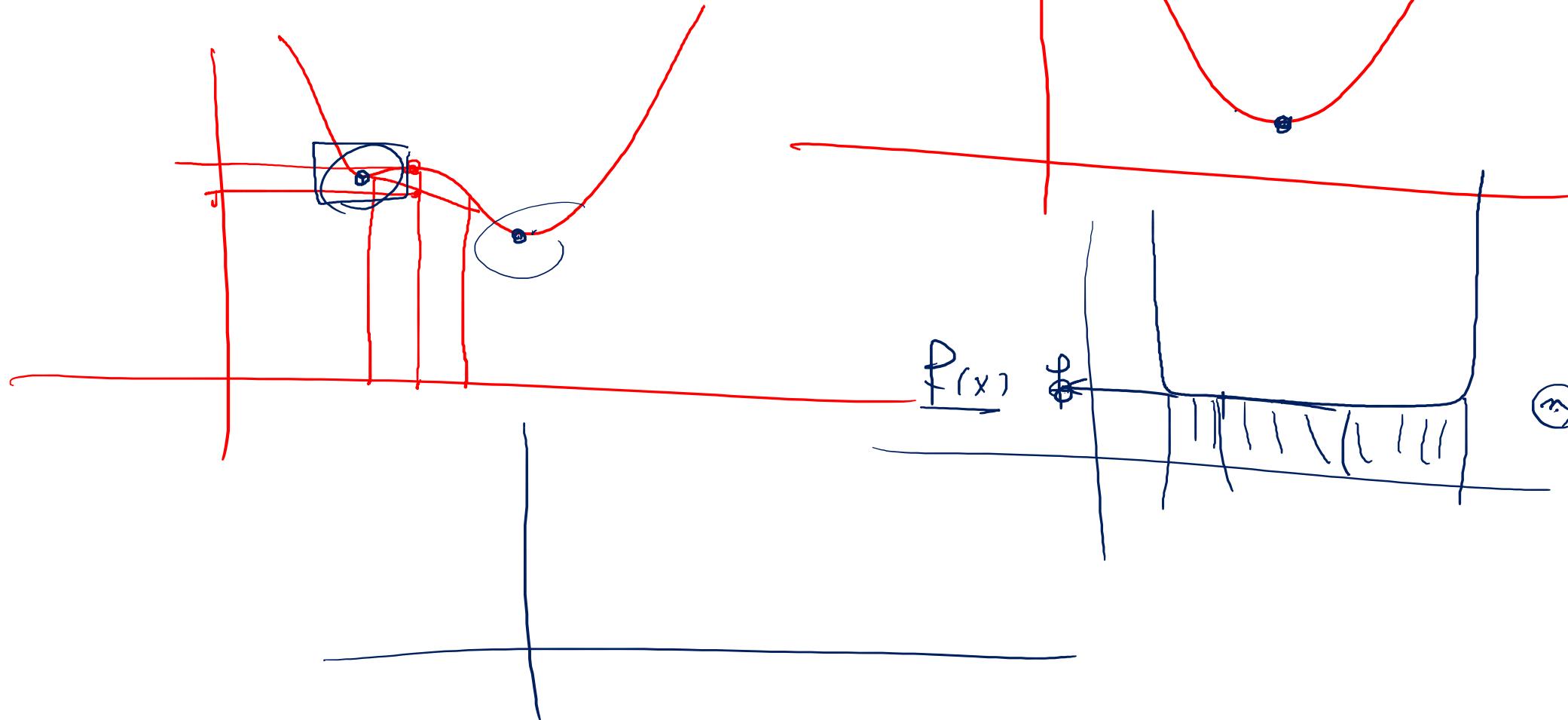
f:

$\partial_n f$

conv. set

contra.

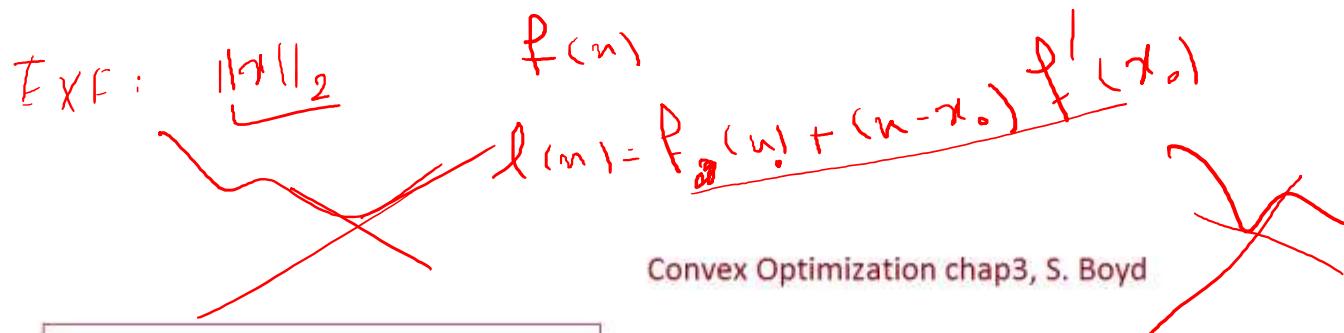
$$f(\theta x_1 + (1-\theta)x_2) \leq \theta f(x_1) + (1-\theta)f(x_2)$$



$f(x)$

$$f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$$

$f(u) = \|u\|_2 \rightarrow \text{norm. Conv}$



## First-order conditions

Suppose  $f$  is differentiable (i.e., its gradient  $\nabla f$  exists at each point in  $\text{dom } f$ , which is open). Then  $f$  is convex if and only if  $\text{dom } f$  is convex and

$$f(y) \geq f(x) + \nabla f(x)^T(y - x)$$

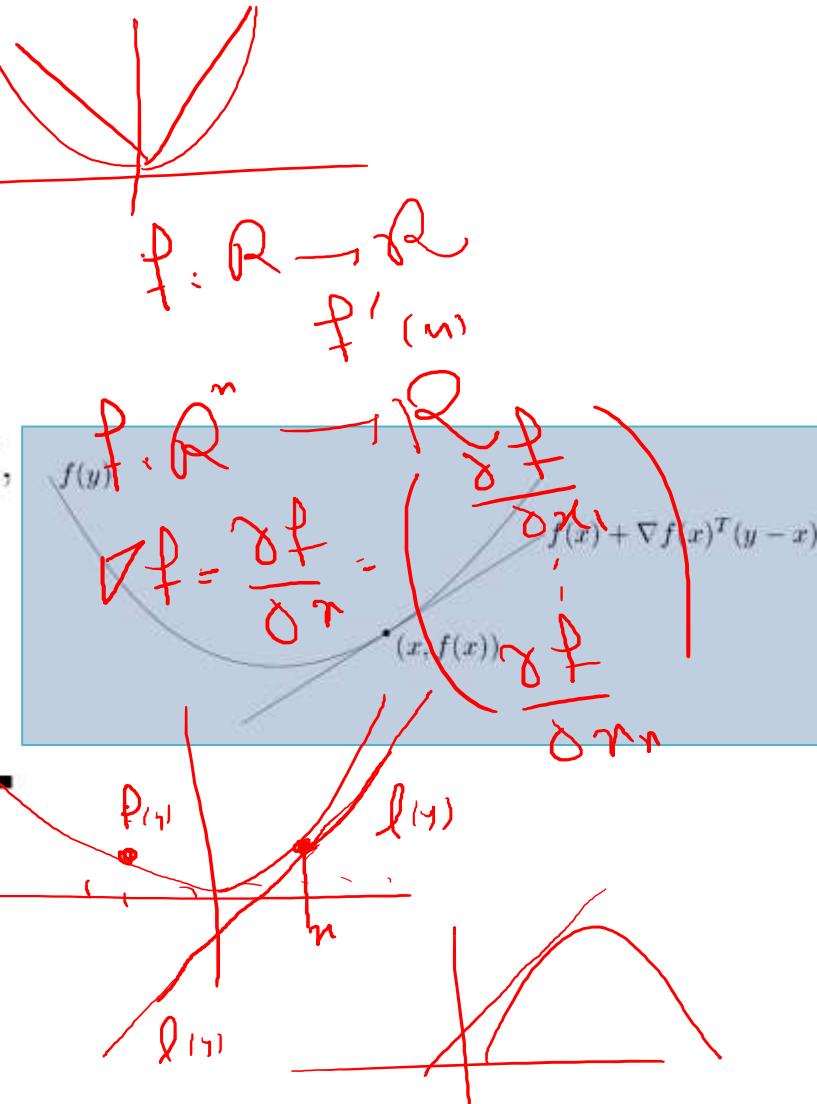
holds for all  $x, y \in \text{dom } f$

$$L_x(y) = f(x) + (y - x)^T \nabla f(x)$$

## Second-order conditions

We now assume that  $f$  is twice differentiable, that is, its *Hessian* or second derivative  $\nabla^2 f$  exists at each point in  $\text{dom } f$ , which is open. Then  $f$  is convex if and only if  $\text{dom } f$  is convex and its Hessian is positive semidefinite: for all  $x \in \text{dom } f$ ,

$$\nabla^2 f(x) \succeq 0.$$



$\alpha \neq x \in R$

$$f_{\text{int}} = h \frac{f_{x_0+h} - f_{x_0}}{h}$$

$$f: R^n \rightarrow R$$

$$x \in R^n$$

$$\nabla f = \frac{\partial f}{\partial x} =$$

$$\begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}$$

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

$$f(x) = \|x\|^r = \sum_i x_i^r$$

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

$$\begin{pmatrix} rx_1 \\ rx_2 \\ \vdots \\ rx_n \end{pmatrix} = r \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

$$\nabla^2 f = \frac{\partial^2 f}{\partial x^2} = r \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix} = r I$$

$$\nabla^2 f Z = Z^T Z - r I \geq 0$$

$$f(w) = \|\Phi w - y\|_2 \quad f(\theta w_1 + (1-\theta)w_2) = \|\Phi(\theta w_1 + (1-\theta)w_2) - [\theta y + (1-\theta)y]\|_2$$

## Linear regression

$$0 < \theta < 1$$

$$= \|\theta \left[ \frac{\Phi w_1 - y}{z_1} \right] + (1-\theta) \left[ \frac{\Phi w_2 - y}{z_2} \right]\|_2$$

$$\|\theta z_1 + (1-\theta)z_2\|_2$$

$$\leq \theta \|z_1\|_1 + (1-\theta) \|z_2\|_1$$

$$\leq \theta f(w_1) + (1-\theta) f(w_2)$$

$$\min_w \|\Phi w - y\|_2^2$$

$$\boxed{\min_w \|\Phi w - y\|_1}$$

### Nonnegative weighted sums

A nonnegative weighted sum of convex functions,  $f = w_1 f_1 + \dots + w_m f_m$  is convex. Boyd chap3

$$\min_w \|\Phi w - y\|_2^2 + \lambda \|w\|_2^2$$

$$\min_w \|\Phi w - y\|_2^2 + \lambda \|w\|_1$$

$$F(w) = \|\Phi w - y\|_2^2 = \underbrace{w^\top \Phi^\top \Phi w}_{h(w) = w^\top Z} - 2y^\top \Phi w + y^\top y$$

$$\frac{\partial h}{\partial w} = Z \mid \\ = 2^* w$$

$$A = \Phi^\top \Phi \\ Z^\top A Z = \frac{1}{2} \Phi^\top \Phi^2 = \|Z\|^2 y$$

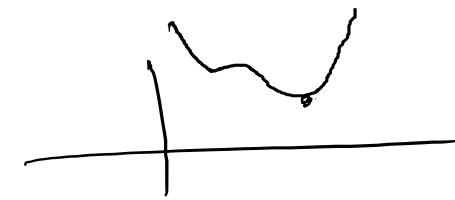
$$f(w) = w^\top Aw \\ \frac{\partial g}{\partial w} = 2Aw$$

$$\frac{\partial g}{\partial w} = 2A \\ = 2\Phi^\top \Phi$$

# Un-Constraint Optimization

$$\min_{\mathbf{w}} f(\mathbf{w})$$

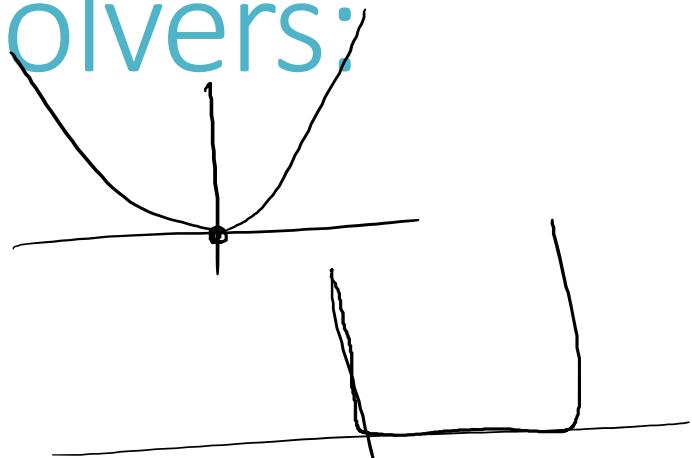
$$\min \|\mathbf{w} - \mathbf{y}\|_2^2$$



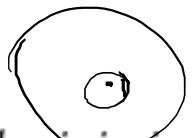
# Unconstraint optimization solvers:

Stim.

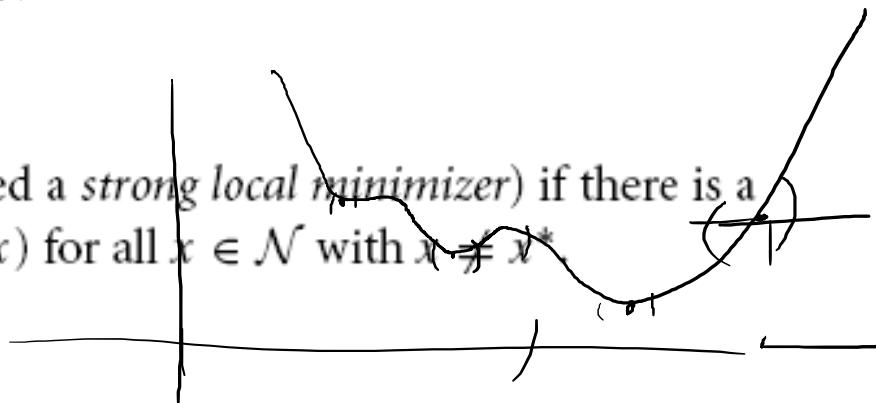
A point  $x^*$  is a *global minimizer* if  $f(x^*) \leq f(x)$  for all  $x$



A point  $x^*$  is a *local minimizer* if there is a neighborhood  $\mathcal{N}$  of  $x^*$  such that  $f(x^*) \leq f(x)$  for all  $x \in \mathcal{N}$ .



A point  $x^*$  is a *strict local minimizer* (also called a *strong local minimizer*) if there is a neighborhood  $\mathcal{N}$  of  $x^*$  such that  $f(x^*) < f(x)$  for all  $x \in \mathcal{N}$  with  $x \neq x^*$ .



$$\frac{\mathbf{y}^T \Phi \mathbf{w}}{\|\mathbf{z}\|_2} = \mathbf{z}^T \mathbf{w} \quad \mathbf{z} = \Phi^T \mathbf{y}$$

# Recognizing a Local Minimum

**Theorem** (Second-Order Necessary Conditions).

If  $x^*$  is a local minimizer of  $f$  then  $\nabla f(x^*) = 0$  and  $\nabla^2 f(x^*)$  is positive semidefinite.

**Theorem** (Second-Order Sufficient Conditions).

Suppose that  $\nabla f(x^*) = 0$  and  $\nabla^2 f(x^*)$  is positive definite.

Then  $x^*$  is a strict local minimizer of  $f$ .

$$\nabla^2 f(x^*) \succ 0$$

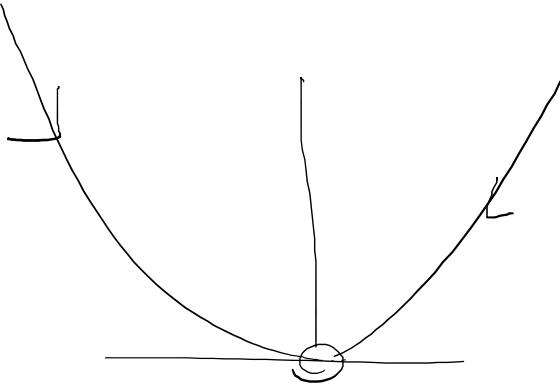
~~$$f(\mathbf{w}) = \|\Phi \mathbf{w} - \mathbf{y}\|_2^2 = \mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2 \mathbf{y}^T \Phi \mathbf{w} + \mathbf{y}^T \mathbf{y}$$~~

When  $f$  is convex, any local minimizer  $x^*$  is a global minimizer of  $f$ .

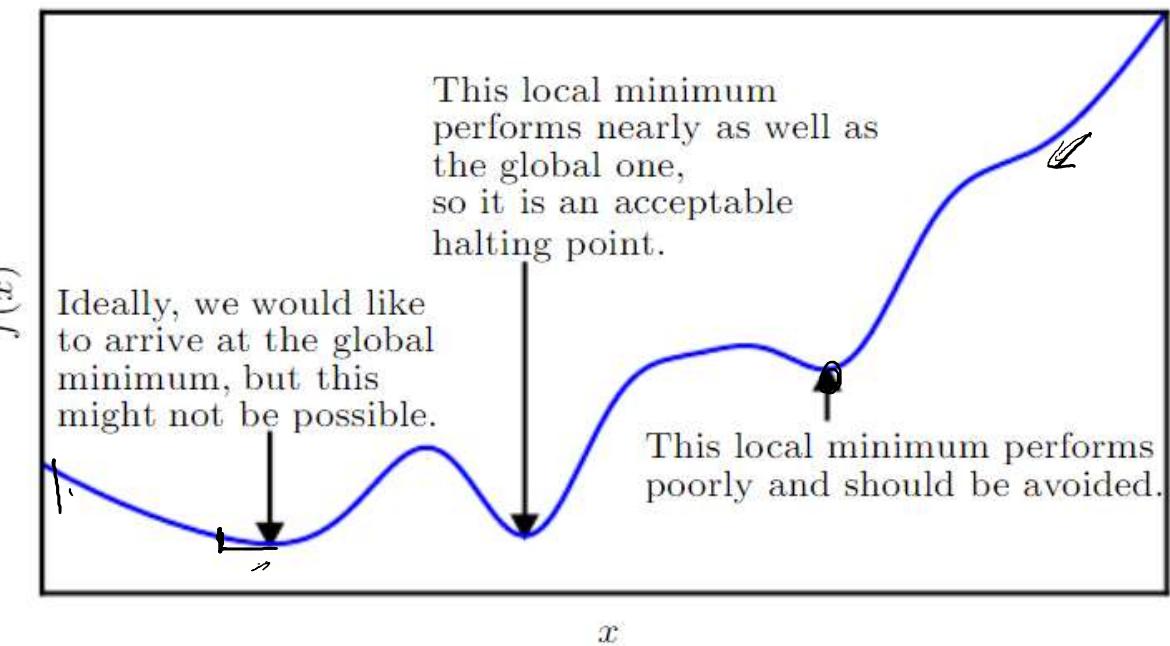
$$\nabla f = 2 \Phi^T \Phi \mathbf{w} - 2 \Phi^T \mathbf{y} = 0 \Rightarrow \Phi^T \Phi \mathbf{w} = \Phi^T \mathbf{y} \Rightarrow \mathbf{w} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

$$\nabla^2 f = 2 \Phi^T \Phi$$

$$\Phi = [\Phi_1 \dots \Phi_n] \quad \Phi^T \Phi = 2 \Phi^T \Phi$$



Convex

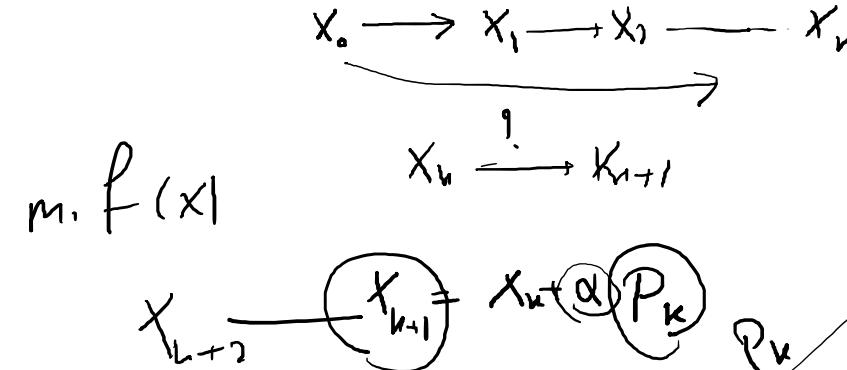


f ↗ ①

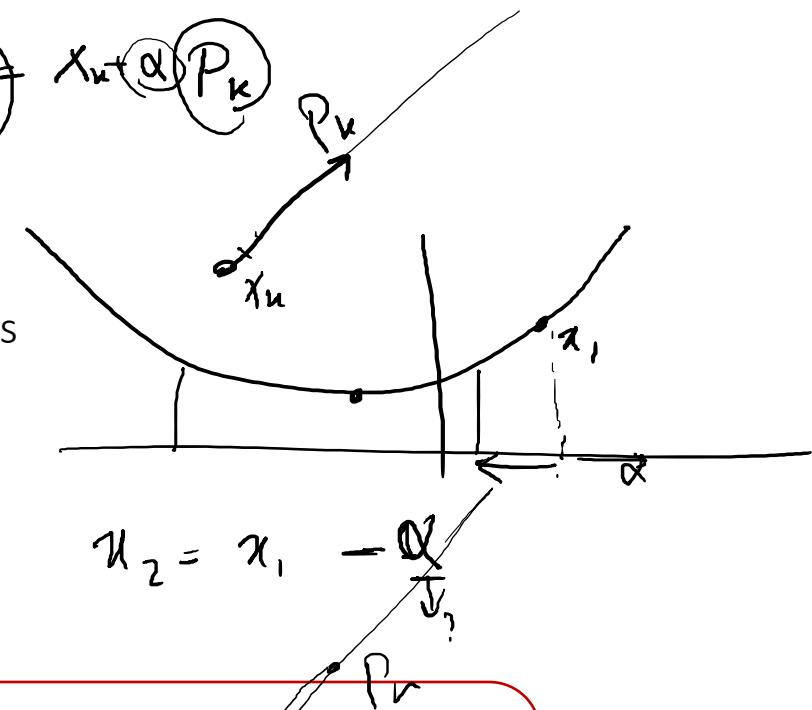
# Algorithms

- Line Search
- Trust Region
- Quasi Newton
- Conjugate gradient
- ...

## 1. Line search



- First order methods
- Second order-methods



✓ choose a direction  $p_k$  and search along this direction

$$\min_{\alpha > 0} f(x_k + \alpha p_k).$$

Line Search → descent method → Steepest descent  
Gradient descent

## SEARCH DIRECTIONS FOR LINE SEARCH METHODS

**Theorem** (Taylor's Theorem).

Suppose that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is continuously differentiable and that  $p \in \mathbb{R}^n$ . Then we have that

(\*)

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

$$f(x + p) = f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x + tp)p,$$

$$f(x_k + \alpha p) = f(x_k) + \alpha p^T \nabla f(x_k) + \frac{1}{2} \alpha^2 p^T \nabla^2 f(x_k + tp)p,$$

$$f(x_{k+1}) = f(x_k) + \alpha p^T \nabla f(x_k) + \frac{1}{2} \alpha^2 p^T \nabla^2 f(x_k + tp)p$$

$$f_{k+1} = f_k + \alpha p^T \nabla f_k \quad \text{if } p^T \nabla f_k < 0$$

$$f_{k+1} < f_k$$

$$\min f(x + \alpha p)$$

$$\nabla f_k^T P < 0$$

$$x^T y = \|x\| \|y\| \cos \theta$$

## Descent Methods

Descent direction  $p$ :  $\nabla f_k^T p < 0$

$$-\|x\| \|y\| \leq x^T y \leq \|x\| \|y\|$$

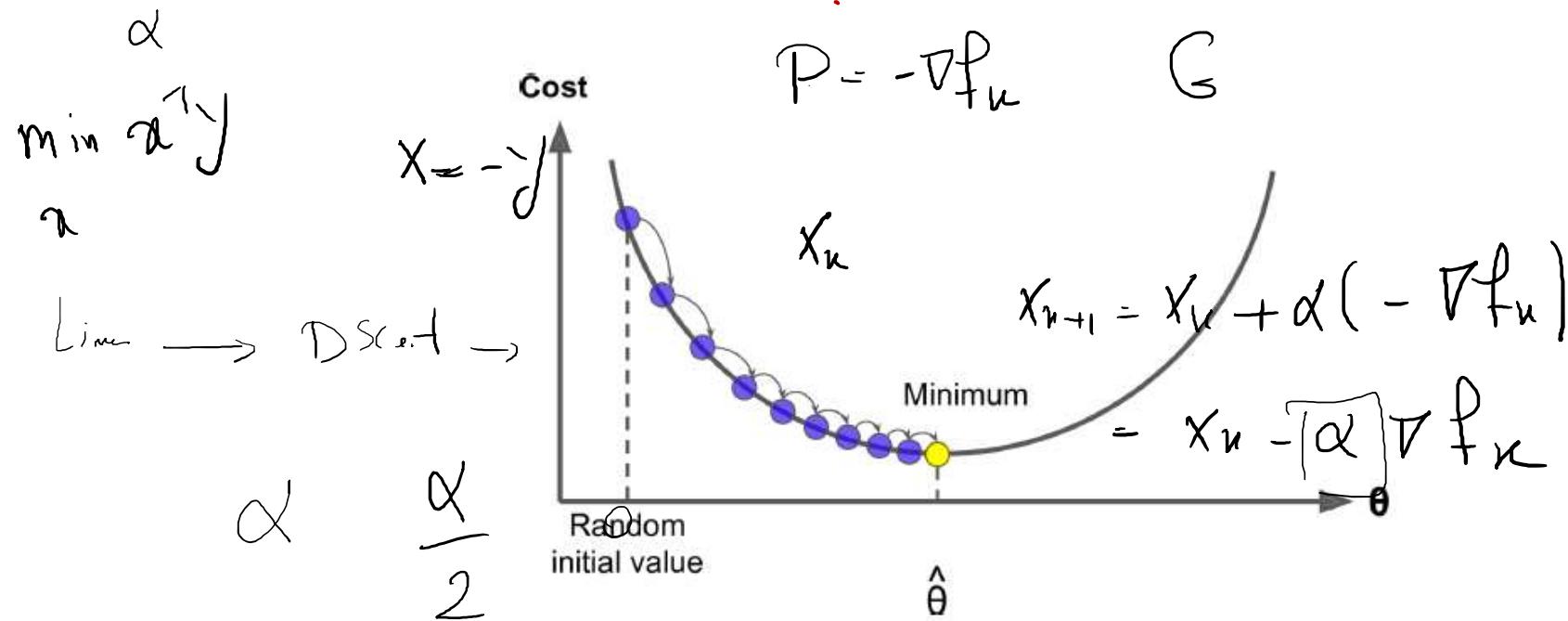
Descent direction is guaranteed to produce a decrease in  $f$ , provided that the step length is sufficiently small.

Steepest descent(Gradient descent)

$$p = -\nabla f_k$$

$$x_{k+1} = x_k - \alpha \nabla f_k$$

Causes error



$$f(x + \alpha p) = f(x) + \alpha \nabla f(x)^T p + \underbrace{\frac{1}{2} \alpha^2 p^T \nabla^2 f(x) p}_{\text{quadratic term}}$$

## Newton direction

Second-order Taylor series approximation

$$f(x_k + p) \approx f_k + \underbrace{p^T \nabla f_k + \frac{1}{2} p^T \nabla^2 f_k p}_{\text{quadratic term}} \stackrel{\text{def}}{=} m_k(p).$$

$$\frac{\partial m_k(p)}{\partial p}$$

$$p_k^N = -(\nabla^2 f_k)^{-1} \nabla f_k.$$

$$P_k = -\nabla^2 f_k^{-1} \nabla f_k$$

The Newton direction is reliable when the difference between the true function  $f(x_k + p)$  and its quadratic model  $m_k(p)$  is not too large.

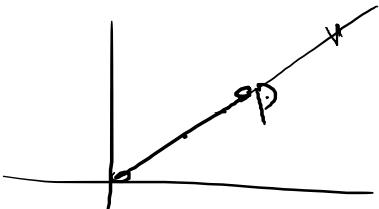
$$f(v_k + \alpha p_k^N)$$

$$\min_P f(x_k + P) \approx f(x_k) + \bar{P}^\top \nabla f_k + \frac{1}{2} \bar{P}^\top \nabla^2 f_k \bar{P}$$

$$\min_P b^\top f_k + \bar{P}^\top \nabla f_k + \frac{1}{2} \bar{P}^\top \nabla^2 f_k \bar{P}$$

$$\frac{\partial m_k(P)}{\partial P} = 0 \Rightarrow \nabla f_k + \nabla^2 f_k P = 0 \quad \cancel{P} = -(\nabla^2 f_k)^{-1} \nabla f_k$$

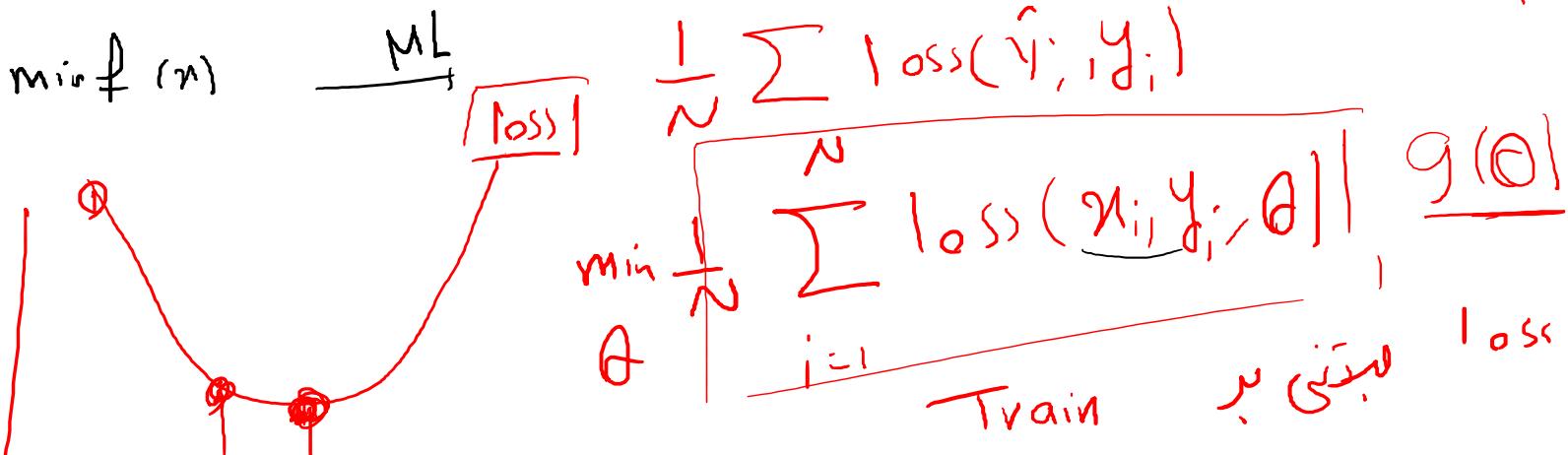
$$f(x_k + \alpha P_k)$$



$$x_{n+1} = x_n + \alpha P \quad \text{if} \quad P = -\nabla f_n \Rightarrow x_{n+1} = x_n - \alpha \nabla f_n$$

GD  
SD

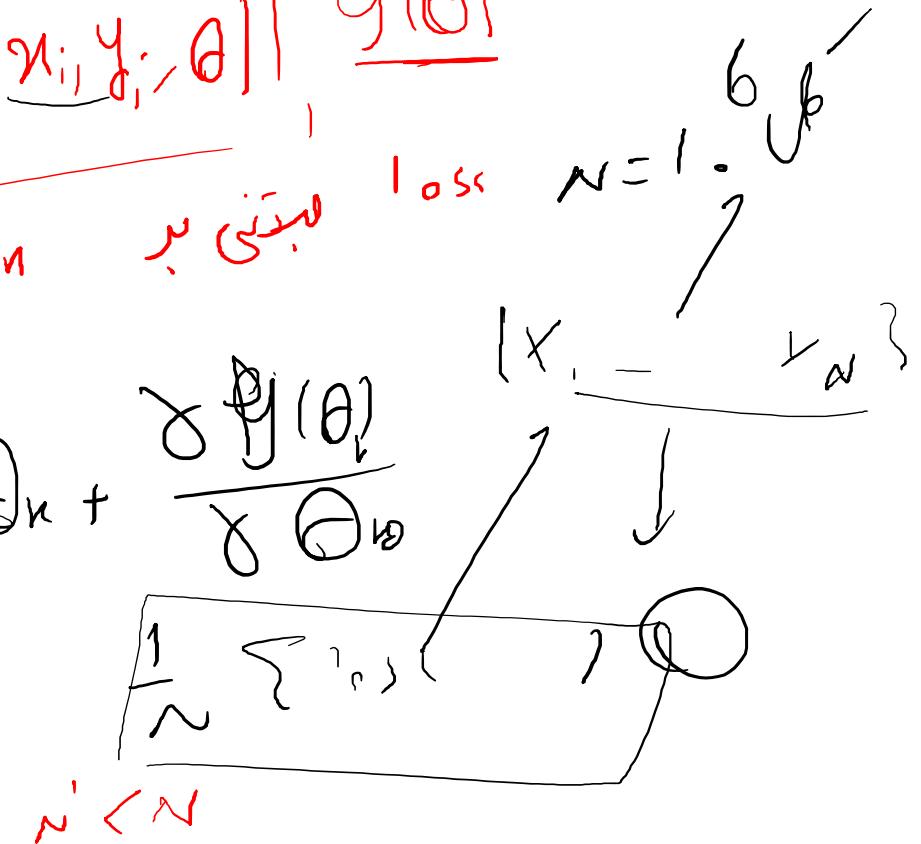
$$\hat{y}_i = f(x_i; \theta)$$



Early Stopping → Regularization

$$\frac{\partial g}{\partial \theta} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \text{Loss}(x_i, y_i, \theta)}{\partial \theta}$$

$$-\frac{1}{N} \sum_{i=1}^N \frac{\partial \text{Loss}(x_i, y_i, \theta)}{\partial \theta}$$





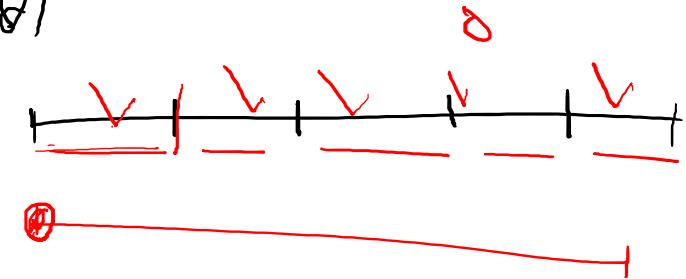


GD



$$\text{Loss}(f(x_i, \theta), y_i) = \text{Loss}(x_i, y_i, \theta)$$

(1)  $\theta_0 \theta_1 \dots$



$$\text{Loss}(\theta) = \frac{1}{N} \sum_{i=1}^N \text{Loss}(x_i, y_i, \theta) \rightarrow \frac{\partial \text{Loss}(\theta)}{\partial \theta} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \text{Loss}(x_i, y_i, \theta)}{\partial \theta}$$

• **Batch:** Uses the whole training set to compute the gradients at every step, which makes it very slow when the training set is large

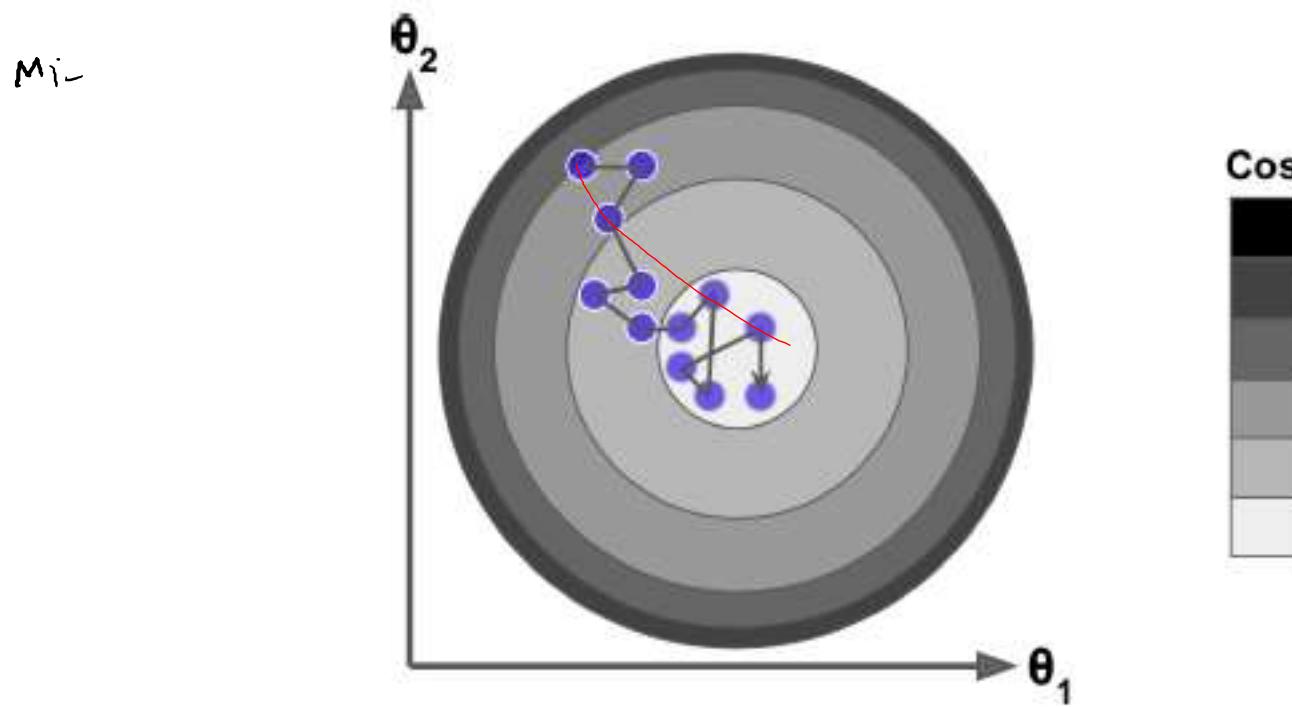
• **Stochastic GD(SGD):** Picks a random instance in the training set at every step and computes the gradients based only on that single instance.

- SGD is much less regular than Batch Gradient Descent
- Instead of gently decreasing until it reaches the minimum, the cost function will bounce up and down, decreasing only on average.

- **Mini- Batch Gradient Descent**

Mini-batch GD computes the gradients on small random sets of instances called mini-batches

# Stochastic Gradient Descent

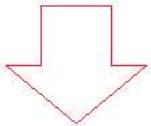


final parameter values are good, but not optimal

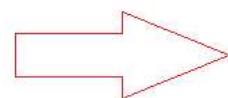


# SGD

$$J(\boldsymbol{\theta}) = \mathbb{E}_{y \sim \hat{p}} L(\mathbf{x}, y, \boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m L(\mathbf{x}^{(i)}, y^{(i)}, \boldsymbol{\theta})$$



$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} L(\mathbf{x}^{(i)}, y^{(i)}, \boldsymbol{\theta}).$$



$$\mathbf{g} = \frac{1}{m'} \nabla_{\boldsymbol{\theta}} \sum_{i=1}^{m'} L(\mathbf{x}^{(i)}, y^{(i)}, \boldsymbol{\theta})$$

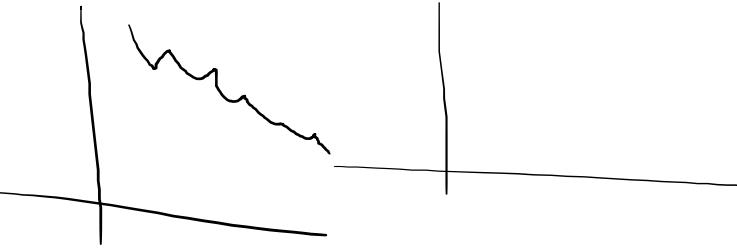
**Estimation of gradient**

---

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \epsilon \mathbf{g},$$

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + \lambda \text{SL}(\mathbf{x}, \mathbf{w})$$

Early Stopping



# Stochastic Gradient

---

**Algorithm 8.1** Stochastic gradient descent (SGD) update at training iteration  $k$

---

**Require:** Learning rate  $\epsilon_k$ .

**Require:** Initial parameter  $\theta$

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

    Compute gradient estimate:  $\hat{\mathbf{g}} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)} ; \theta), \mathbf{y}^{(i)})$

    Apply update:  $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$

**end while**

---

**Condition for Convergence**

$$\sum_{k=1}^{\infty} \epsilon_k = \infty \quad \text{and} \quad \sum_{k=1}^{\infty} \epsilon_k^2 < \infty \quad \epsilon_k = \frac{1}{k}$$

**In Practice (One example):**  $\epsilon_k = (1 - \alpha)\epsilon_0 + \alpha\epsilon_{\tau} \quad \alpha = \frac{k}{\tau}$

After iteration  $\tau$ , it is common to leave  $\epsilon$  constant

# Unconstraint Examples in Learning

## Least squares

$$E = \sum_{n=1}^N (\mathbf{w}^T \phi_n - t_n)^2$$

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N \phi_n \phi_n^T = \Phi^T \Phi \quad \nabla E(\mathbf{w}) = \sum_{n=1}^N (\mathbf{w}^T \phi_n - t_n) \phi_n = \Phi^T \Phi \mathbf{w} - \Phi^T \mathbf{t}$$



$$\begin{aligned}\mathbf{w}^{(\text{new})} &= \mathbf{w}^{(\text{old})} - (\Phi^T \Phi)^{-1} \{ \Phi^T \Phi \mathbf{w}^{(\text{old})} - \Phi^T \mathbf{t} \} \\ &= (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}\end{aligned}$$

$$\left\{ \begin{array}{l} f(\omega) = \omega^T \Phi^T \Phi \omega - 2y^T \Phi \omega + y^T y + \lambda \omega^T \omega \\ \frac{\partial f}{\partial \omega} = 2 \Phi^T \Phi \omega - 2 \Phi^T y + 2\lambda \omega = 0 \end{array} \right. \quad \begin{aligned} (\Phi^T \Phi + \lambda I) \omega &= \Phi^T y \\ \omega &= (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y \end{aligned}$$

$$\frac{\partial f}{\partial \omega^i} = 2 \Phi^T \Phi \cancel{\omega^i} \\ = 2(\cancel{\Phi^T \Phi} + \lambda \delta)$$

$$\min_{\omega} \frac{f(\omega)}{\|\Phi \omega - y\|_2^2 + \lambda \|\omega\|_2^2}$$

GD

$$\omega_{k+1} = \omega_k - \alpha \left( 2 \Phi^T \Phi \omega_k - 2 \Phi^T y + 2 \lambda \omega_k \right)$$

Newton

$$\omega_{k+1} = \omega_k + P$$

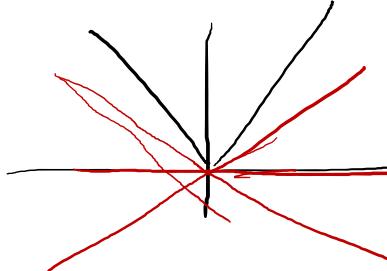
$$P = -\nabla^2 f^{-1} \Rightarrow P = -(\Phi^T \Phi + \lambda I)^{-1} \left[ \Phi^T \Phi \omega_k - \Phi^T y \right]$$

$$= -(\Phi^T \Phi + \lambda I)^{-1} \left[ (\Phi^T \Phi + \lambda I) \omega_k - \Phi^T y \right] = -\cancel{\omega_k} \cancel{+} \omega_{k+1} = \omega_k - (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$$

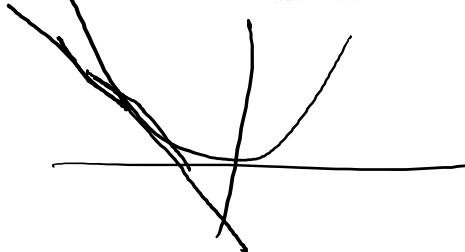
$$\|\mathbf{w}\|_1 = \sum |\mathbf{w}_i|$$

# SGD for Lasso

$$\min_{\mathbf{w}} \|\Phi \mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1$$



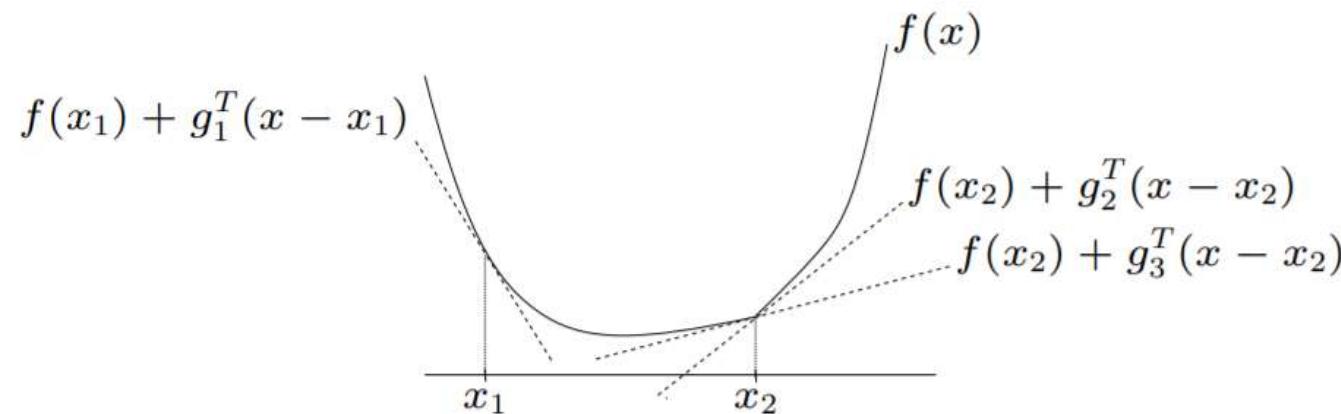
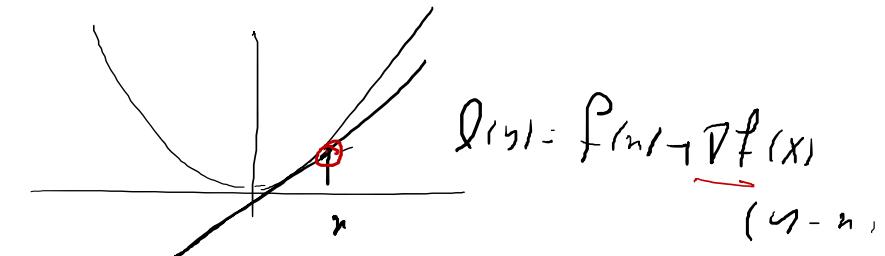
Boyd-  
Murphy-13.3.2-13.3.5



$g$  is a **subgradient** of  $f$  (not necessarily convex) at  $x$  if

$$f(y) \geq f(x) + g^T(y - x) \quad \text{for all } y$$

$\hat{f}(y) = f(x) + g^T(y - x)$

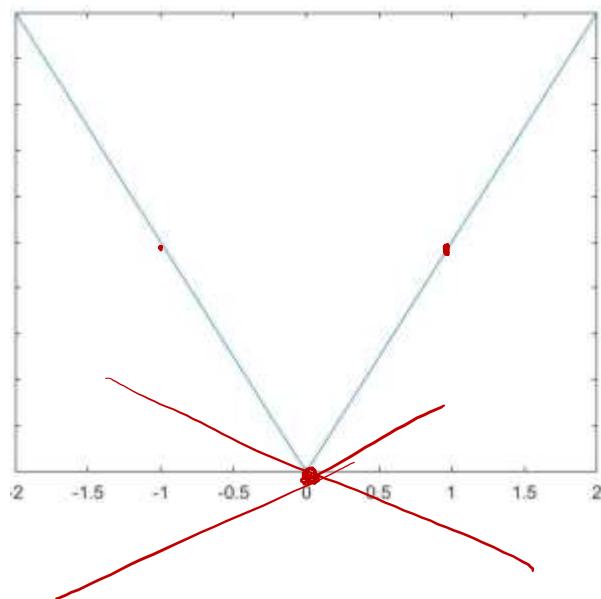


if  $f$  is convex and differentiable,  $\nabla f(x)$  is a subgradient of  $f$  at  $x$

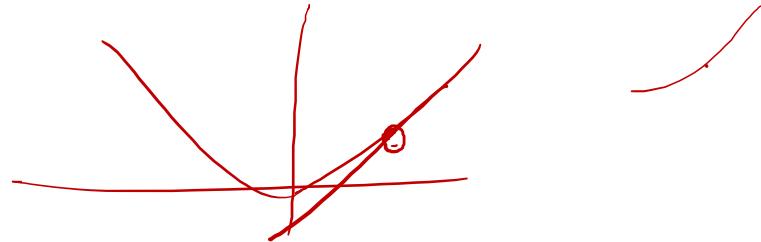
# Subgradient: Example

$$f(\theta) = |\theta| \quad \longrightarrow \quad \partial f(\theta) = \begin{cases} \{-1\} & \text{if } \theta < 0 \\ [-1, 1] & \text{if } \theta = 0 \\ \{+1\} & \text{if } \theta > 0 \end{cases}$$

1      subgradient



$f(\theta) = |\theta|$

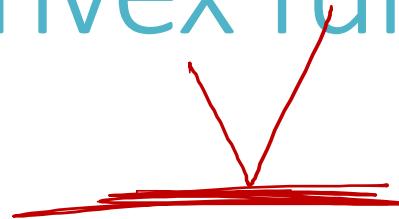


if  $f$  is convex,

- $\partial f(x)$  is nonempty, for  $x \in \text{relint } \text{dom } f$
- $\partial f(x) = \{\nabla f(x)\}$ , if  $f$  is differentiable at  $x$
- if  $\partial f(x) = \{g\}$ , then  $f$  is differentiable at  $x$  and  $g = \nabla f(x)$

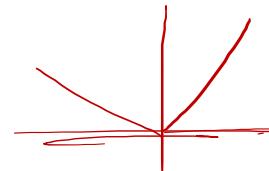
# Optimality condition for convex function

recall for  $f$  convex, differentiable,



$$f(x^*) = \inf_x f(x) \iff 0 = \nabla f(x^*)$$

generalization to nondifferentiable convex  $f$ :



$$f(x^*) = \inf_x f(x) \iff 0 \in \partial f(x^*)$$

Boyd

---

Suppose  $f : \mathbf{R}^n \rightarrow \mathbf{R}$  is convex. To minimize  $f$ , the subgradient method uses the iteration

$$x^{(k+1)} = x^{(k)} - \alpha_k g^{(k)}.$$

Here  $x^{(k)}$  is the  $k$ th iterate,  $g^{(k)}$  is *any* subgradient of  $f$  at  $x^{(k)}$ , and  $\alpha_k > 0$  is the  $k$ th step size.

## Coordinate descent for lasso

$$\sum_{i=1}^N \frac{-\text{RSS}(\mathbf{w})}{(y_i - (\mathbf{w}_0 + \mathbf{w}^T \mathbf{x}_i))^2 + \lambda \|\mathbf{w}\|_1}$$

$\min_{\mathbf{w}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{w}\|_1$

$\downarrow$

$$\frac{\partial}{\partial w_j} \text{RSS}(\mathbf{w}) = a_j w_j - c_j$$

$$a_j = 2 \sum_{i=1}^n x_{ij}^2$$

$$c_j = 2 \sum_{i=1}^n x_{ij} (y_i - \mathbf{w}_{-j}^T \mathbf{x}_{i,-j})$$

$$\begin{aligned} \partial_{w_j} f(\mathbf{w}) &= (a_j w_j - c_j) + \lambda \partial_{w_j} \|\mathbf{w}\|_1 \\ &= \begin{cases} \{a_j w_j - c_j - \lambda\} & \text{if } w_j < 0 \\ [-c_j - \lambda, -c_j + \lambda] & \text{if } w_j = 0 \\ \{a_j w_j - c_j + \lambda\} & \text{if } w_j > 0 \end{cases} \end{aligned}$$



$$\hat{w}_j(c_j) = \begin{cases} (c_j + \lambda)/a_j & \text{if } c_j < -\lambda \\ 0 & \text{if } c_j \in [-\lambda, \lambda] \\ (c_j - \lambda)/a_j & \text{if } c_j > \lambda \end{cases}$$

$$\hat{w}_j = \text{soft}\left(\frac{c_j}{a_j}; \frac{\lambda}{a_j}\right)$$

$$\text{soft}(a; \delta) \triangleq \text{sign}(a) (|a| - \delta)_+$$

$x_+ = \max(x, 0)$  is the positive part of  $x$ .

- LARS
- Nesterov
- Proximal Gradient method
- ADMM
- FISTA