

# 4 Linear Algebra

## • Review

$$A = n \begin{bmatrix} & m \end{bmatrix}$$

$A$ :  $n \times m$  matrix

$$A = [a_1 \ a_2 \ \dots \ a_m] \quad a_i \in \mathbb{R}^n$$

def

$$\textcircled{1} \text{range}(A) = \text{column space of } A = \left\{ x = \sum_{i=1}^m \alpha_i a_i : \alpha_i \in \mathbb{R} \right\} \in \mathbb{R}^n$$

$$\textcircled{2} \text{rank}(A) = \dim(\text{range}(A))$$

$$\textcircled{3} \text{null}(A) = \{ x \in \mathbb{R}^m : Ax = 0 \}$$

$$\textcircled{4} W^\perp := \text{orthogonal complement of set } W \in \mathbb{R}^m \\ = \{ v \in \mathbb{R}^m : v^T w = 0 \quad \forall w \in W \}$$

$$\textcircled{5} W \text{ is subspace of } \mathbb{R}^n \text{ if}$$

$$\text{i) } 0 \in W$$

$$\text{ii) } w_1, w_2 \in W \Rightarrow w_1 + w_2 \in W$$

$$\text{iii) } w \in W \Rightarrow cw \in W \quad (\forall c \in \mathbb{R})$$

$$\text{Thms} \left( \begin{array}{l} \textcircled{1} \text{rank}(A) + \dim(\text{null}(A)) = n \\ \textcircled{2} \mathbb{R}^m = W \oplus W^\perp \quad (\text{if } W \text{ is subspace of } \mathbb{R}^m) \\ \textcircled{3} \begin{array}{l} (\text{row}(A))^\perp = \text{null}(A) \\ (\text{col}(A))^\perp = \text{null}(A^T) \end{array} \end{array} \right.$$

## • Norm

definition  $\|\cdot\|$  is a norm on  $\mathbb{R}^n$  if  $\|x\| \in [0, \infty)$   $\forall x \in \mathbb{R}^n$ ,

- ①  $\|x+y\| \leq \|x\| + \|y\| \quad \forall x, y \in \mathbb{R}^n$
- ②  $\|\alpha x\| = |\alpha| \|x\| \quad \forall \alpha \in \mathbb{R}, x \in \mathbb{R}^n$
- ③  $\|x\| = 0 \Leftrightarrow x = 0 \in \mathbb{R}^n$

## Famous Norms

- ①  $L_2$  norm (= Euclidean norm)

$$\|x\|_2 = (x_1^2 + \dots + x_n^2)^{\frac{1}{2}}$$

- ②  $L_p$  norm

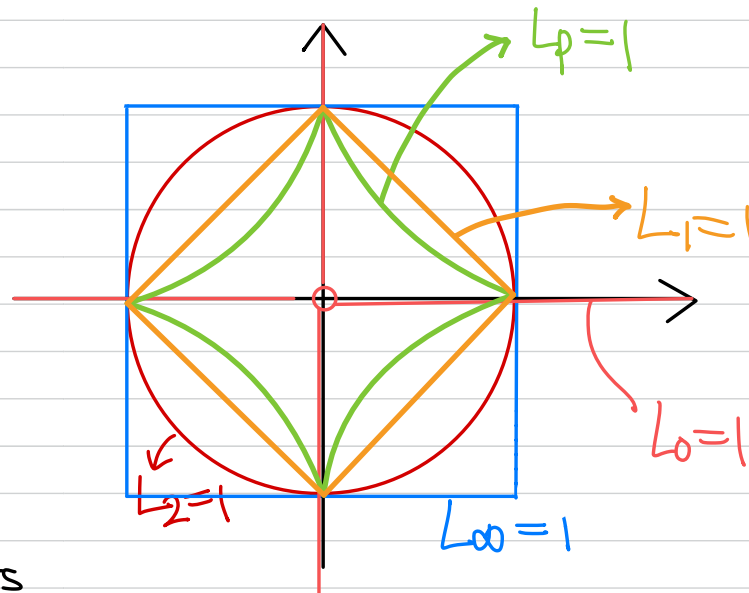
$$\|x\|_p = (x_1^p + \dots + x_n^p)^{\frac{1}{p}}$$

- ③  $L_\infty$  norm

$$\|x\|_\infty = \max_i |x_i|$$

- ④  $L_0$  norm

$$\|x\|_0 = \# \text{ of nonzero elements}$$



ex  $L_1$  norm is less sensitive to outliers than  $L_2$  norm

## • Orthogonal

definition  $U \in \mathbb{R}^{n \times n}$  is orthogonal (=unitary) if transpose of it is inverse of it

$$\Leftrightarrow U^T U = I \Leftrightarrow U U^T = I$$

$\Leftrightarrow$  column vectors of  $U$  forms a orthonormal basis of  $\mathbb{R}^n$

## • Singular Vector Decomposition

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

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

$u_1, \dots, u_m$  : orthonormal basis of  $\mathbb{R}^m$

$v_1, \dots, v_n$  : orthonormal basis of  $\mathbb{R}^n$

$$\text{def } m \begin{bmatrix} n \\ A \end{bmatrix} = m \begin{bmatrix} m \\ U \end{bmatrix} m \begin{bmatrix} n \\ \Sigma \end{bmatrix} n \begin{bmatrix} n \\ V^T \end{bmatrix}$$

$$= [u_1 | u_2 \dots | u_m] \Sigma \begin{bmatrix} v_1^T \\ \vdots \\ v_n^T \end{bmatrix}$$

$$\begin{bmatrix} \sigma_1 & \dots & \sigma_r & \dots & 0 \end{bmatrix}$$

$$= u_1 \sigma_1 v_1^T + \dots + u_r \sigma_r v_r^T$$

$$= m \begin{bmatrix} u_1 & \dots & u_r \end{bmatrix} r \begin{bmatrix} \sigma_1 & \dots & \sigma_r \end{bmatrix} n \begin{bmatrix} v_1^T \\ \vdots \\ v_r^T \end{bmatrix}$$

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

## • Pseudo - Inverse

### definition

$$A^+ = V_C \Sigma_C^{-1} U_C \quad \Sigma_C^{-1} = \begin{bmatrix} 1/\sigma_1 & \dots & 1/\sigma_r \end{bmatrix}$$

$$\ast (U_C U_C^T \neq \bar{I}_d \quad V_C V_C^T \neq \bar{I}_d \\ U_C^T U_C = \bar{I}_d \quad V_C^T V_C = \bar{I}_d)$$

### property

minimum achiever of  $\|A\theta - B\|^2 : A = U \Sigma V^T$

then  $\theta = A^+ B \rightarrow$  minimum achiever



- **Expected Risk** (= True risk, population risk)

$$\text{def } R[f] = \mathbb{E}[l(f(x), y)]$$

→ minimize  $R[f]$

- **Bayes Risk** (= optimal risk)

$$\text{def } R^* = \inf_f R[f]$$

→ infimum over all possible functions

optimal achieving  $f^*$ : Bayes predictor

$$R^* = R[f^*]$$

- **Empirical Risk Minimization**

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N l(f_{\theta}(x^{(i)}), y^{(i)})$$

Actually, Training is Optimizing

→ minimizing  $L(\theta)$ : Gradient Descent

## • Some Math Knowledge

- def
- ① directional derivative =  $\nabla f(x) \cdot v = \langle \nabla f(x), v \rangle = \nabla f(x)^T v$
  - ② convexity  $\Leftrightarrow f''(x) \geq 0$  in  $\mathbb{R}^d \Leftrightarrow \nabla^2 f \geq 0$
  - ③  $f$  is  $L$ -Lipschitz if  $f(x_1) - f(x_2) \leq L \|x_1 - x_2\|$
  - ④  $f$  is  $L$ -smooth if  $\nabla f$  is  $L$ -Lipschitz.
  - ⑤  $f$  is  $\mu$ -strongly convex if  $f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|^2$

- thms
- ①  $f$ : twice diff  $L$ -smooth  $\Leftrightarrow \nabla^2 f \preceq LI$
  - ②  $f(y) \leq f(x) + \langle \nabla f(x), y-x \rangle + \frac{L}{2} \|y-x\|^2$
  - ③  $f$ : minimizer  $x_0$  then  
 $\frac{1}{2L} \|\nabla f(z)\|^2 \leq f(z) - f(x_0) \leq \frac{L}{2} \|z-x_0\|^2$
  - ④  $f$ :  $\mu$ -strongly-convex  $\Rightarrow \nabla^2 f \succeq \mu I \Leftrightarrow \nabla^2 f - \mu I \succeq 0$
  - ⑤  $f$ : convex  $\Rightarrow (\nabla f(x) - \nabla f(y))^T (x-y) \geq 0$
  - ⑥  $f$ :  $\mu$ -strongly convex  $\Rightarrow (\nabla f(x) - \nabla f(y))^T (x-y) \geq \mu \|x-y\|^2$

## • $C_0$ - Coercivity : name of inequality

$f$ : convex,  $L$ -smooth

$$\Rightarrow (\nabla f(x) - \nabla f(y))^T (x-y) \geq \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|^2$$

## • Polyak - Łojasiewicz (PL)

if  $f$  is  $\mu$ -strongly convex, then  $f$  satisfies PL condition

$$\|\nabla f(x)\|^2 \geq 2\mu (f(x) - f(x_{opt}))$$

# • Gradient Descent

$f: \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $L$ -smooth,  $x^{(0)}$ : pick randomly

$$x^{(k+1)} = x^{(k)} - \alpha \nabla f(x^{(k)}) \quad \alpha: \text{learning rate}$$

If  $f$  has global maximum,  $\exists \alpha > 0$  s.t. Gradient Descent converges

## • Gradient Descent + convex + $L$ -smooth

$$f(x^T) - f(x_{\text{opt}}) \leq \frac{1}{T} \left( \frac{1}{2\alpha} \|x^{(0)} - x_{\text{opt}}\|^2 \right) = O\left(\frac{1}{T}\right)$$

## • Gradient Descent + convex + $\mu$ -strongly convex + $L$ -smooth

(Lemma  $f: \mu$ -strongly convex &  $L$ -smooth  
 $\Rightarrow h(x) = f(x) - \frac{\mu}{2} \|x\|^2$  is convex &  $(L - \mu)$ -smooth

$$f(x^T) - f(x_{\text{opt}}) \leq \frac{L}{2} C^T \|x^{(0)} - x_{\text{opt}}\|^2 = O(C^T)$$

$C < 1$ : exponentially fast convergence

conclusion

① GD + convex		$O\left(\frac{1}{T}\right)$
② GD + strong convex		$O(e^{-\alpha T})$
③ GD + convex		$\exists \text{ algorithm } O\left(\frac{1}{T^2}\right)$



## • Operator

$$T: \mathbb{R}^d \rightarrow \mathbb{R}^d$$

- def
- ①  $T$  is **nonexpansive** if  $\|Tx - Ty\| \leq \|x - y\|$
  - ②  $T$  is **contractive** if  $\|Tx - Ty\| \leq L\|x - y\|$  ( $L < 1$ )
  - ③  $T$  is  **$\theta$ -averaged** if  $T = (1-\theta)I + \theta S$ ,  $I = \text{id}$ ,  $S$  nonexpansive
  - ④  $x$  is **fixed point** of  $T$  if  $Tx = x$

## • Fixed point iterator (Picard iteration)

$x^{(0)}$ : starting point,  $T$  is  $\theta$ -averaged,  $x_f$ : fixed point

$$\|x^{(k+1)} - x^{(k)}\| \leq \frac{1}{k+1} \cdot \frac{\theta}{1-\theta} \|x^{(0)} - x_f\|^2$$

pf by  $\| (1-\theta)x + \theta y \|^2 = (1-\theta)\|x\|^2 + \theta\|y\|^2 - \theta(1-\theta)\|x-y\|^2$

$$\Rightarrow \text{converging rate} : O\left(\frac{1}{k}\right)$$

ex> Gradient descent is fixed point iteration of  $\theta$ -averaged operator

$$T = I - \alpha \nabla f$$

$$S = I - \frac{2}{L} \nabla f \quad \theta = \frac{\alpha L}{2}$$

+ if  $f$  is  $\mu$ -strongly convex,

$$\|Tx - Ty\| \leq L\|x - y\|, L < 1$$

$\therefore T$  is contractive operator

# • Nesterov's accelerated gradient method (AGM)

$x_0$ : random initialization     $f$ : convex,  $L$ -smooth  
 $y_0 = x_0$

$$① \quad x_{k+1} = y_k - \frac{1}{L} \nabla f(y_k)$$

$$② \quad y_{k+1} = x_{k+1} + \frac{k-1}{k+2} (x_{k+1} - x_k)$$

## \* Equivalent formula

$$① \quad x_{k+1} = y_k - \frac{1}{L} \nabla f(y_k)$$

$$② \quad z_{k+1} = z_k - \frac{k+1}{2L} \nabla f(y_k)$$

$$③ \quad y_{k+1} = \left(1 - \frac{2}{k+2}\right) x_{k+1} + \frac{2}{k+2} z_{k+1}$$

$$\& \quad x_0 = y_0 = z_0$$

## \* Converging rate

global minimizer  $x^*$

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

## \* First order Method

any iterate algorithm that selects  $x_{k+1}$  in the set

$$\{x_0 + \text{span}\{\nabla f(x_0), \nabla f(x_1), \dots, \nabla f(x_k)\}\}$$

$$\forall k \in \frac{d-1}{2}, \exists f, \forall \text{ 1st order method, } f(x_k) - f(x^*) \geq \frac{L \|x_0 - x^*\|^2}{32(k+1)}$$

( $d$ : dimension)

- **Momentum**

$$V_0 = 0$$

$$V_{k+1} = \beta V_k + \alpha \nabla f(x_k)$$

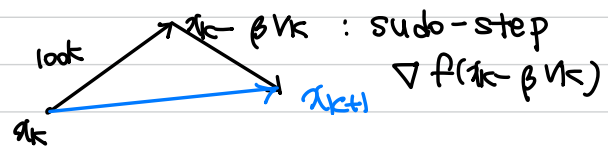
$$x_{k+1} = x_k - V_{k+1}$$

$$(\beta=0 : \text{GD})$$

- **Nesterov oscillation**

$$V_{k+1} = \beta V_k + \alpha \nabla f(x_k - \beta V_k)$$

$$x_{k+1} = x_k - V_{k+1}$$



$$gd : O\left(\frac{1}{k^2}\right)$$

$$gd + \text{momentum} : O\left(\frac{1}{k}\right)$$

# • Stochastic gradient descent

1st SGD: minimize  $f(x)$

$$\min f(\theta) = \frac{1}{N} \sum_{i=1}^N f_i(\theta)$$

① SGD:  $i_k$  is randomly picked from  $\{1, 2, \dots, N\}$

$$\theta_{k+1} = \theta_k - \alpha \cdot \nabla f_{i_k}(\theta)$$

sample gradients

② minibatch:  $I_k \subseteq \{1, 2, \dots, N\}$ ,  $|I_k| = B$

$$\theta_{k+1} = \theta_k - \alpha \cdot \frac{1}{B} \sum_{i \in I_k} \nabla f_i(\theta)$$

$g_k$ : stochastic gradient

$$\textcircled{3} \mathbb{E}[f(\bar{\theta}_K) - f(\theta^*)] \leq \frac{\sqrt{L^2 \sigma^2} \|\theta_0 - \theta^*\|^2}{\sqrt{K+1}}$$

: SGD rate

$$\therefore O\left(\frac{1}{\sqrt{K}}\right)$$

when  $\alpha = \frac{\|\theta_0 - \theta^*\|^2}{\sqrt{L^2 \sigma^2} \sqrt{K+1}}$ ,  $\bar{\theta}_K = \frac{1}{K+1} \sum_{k=0}^K \theta_k$ ,

$L(\theta)$ :  $L$ -Lipschitz, continuous, convex

$\theta_0$ : starting point  $K$ : total iteration count

- **Regularization** purpose: prevent from overfitting

$$r(\theta) = \sum_{i=1}^d \theta_i^2$$

$$\sum_{i=1}^N \ell(f_\theta(x^i), y^i) + \lambda \|\theta\|^2$$

↪ **Ridge Regression** : ① linear + ②  $L_2$  penalty ( $\lambda \|\theta\|^2$ )

ridge regression

goal: minimizing  $\|X\theta - Y\|^2 + \lambda \|\theta\|^2$

if  $X$ : full rank  $\Rightarrow$  minimum achiever is  $\theta^* = (X^T X + \lambda I)^{-1} X^T Y$

- $f$ :  $G$ -Lipschitz, continuous, convex

$x^*$ : minimizer of  $f(x) + \frac{\mu}{2\|x\|^2}$

$$\mathbb{E}_k[g_k] = \nabla f(x_k), \text{Var}_k(g_k) \leq \sigma^2$$

$$\alpha_k = \frac{1}{\mu(k+1)}, \bar{x}_k = \frac{1}{k+1} \sum_{i=0}^k x_i$$

$$\rightarrow f(\bar{x}_k) - f(x^*) \leq \frac{2(G^2 + \sigma^2)}{\mu} \cdot \frac{1 + \log(k+1)}{k+1}$$

$\therefore$  convergence rate is  $O\left(\frac{\log k}{k}\right)$  which is smaller than SGD  $O\left(\frac{1}{\sqrt{k}}\right)$

- comparison

	GD	SGD
convex	$O(\frac{1}{k})$	$O(\frac{1}{\sqrt{k}})$
strongly convex	$O(c^k)$ ( $c < 1$ )	$O(\frac{\log k}{k}) \xrightarrow{\text{can become}} O(\frac{1}{k})$

# • RMSProp

$$S_t = r S_{t-1} + (1-r) (\nabla L(x_t))^2 \quad \text{element wise!}$$

: vector

$$\begin{pmatrix} 4 \\ 2 \\ 1 \end{pmatrix}^2 = \begin{pmatrix} 16 \\ 4 \\ 1 \end{pmatrix}$$

$$x_{t+1} = x_t - \frac{\alpha}{\sqrt{S_t + \epsilon}} \nabla L(x_t)$$

elementwise

$$\begin{pmatrix} 4 \\ 2 \\ 1 \end{pmatrix} / \begin{pmatrix} 16 \\ 4 \\ 1 \end{pmatrix} = \begin{pmatrix} 1/4 \\ 1 \\ 1 \end{pmatrix}$$

# • Adam (Adaptive moment Estimation)

$$V_t = \beta_1 V_{t-1} + (1-\beta_1) \nabla L(x_t)$$

$$S_t = \beta_2 S_{t-1} + (1-\beta_2) (\nabla L(x_t))^2$$

$$\hat{V}_t = \frac{V_t}{1-\beta_1} \quad \hat{S}_t = \frac{S_t}{1-\beta_2}$$

$$x_{t+1} = x_t - \frac{\alpha}{\sqrt{S_t + \epsilon}} \hat{V}_t$$

thm Adam converges to  $\forall$  convex fcn  
 $\Rightarrow$  it is not proved

- SAM (Sharpness Aware Minimization)

finding flat minima

$$\min_x \left\{ \max_{\varepsilon: \|\varepsilon\| \leq \rho} L(x + \varepsilon) \right\}$$