

466 / 566 : Sept 17 - Oct. 3

2019

$X$  observations / inputs , e.g.  $X \subseteq \mathbb{R}^d$

$y$  targets , e.g.  $y = \{0, 1\}$  or  $\mathcal{Y} = \mathbb{R}^d$

$f: X \rightarrow \mathcal{Y}$  cost :  $\mathcal{Y} \times \mathcal{Y} \rightarrow [0, \infty)$

cost( $\hat{y}, y$ ) ,  $\hat{y} = f(x)$   $X$  random variable.

instance  $x \in X$  ,  $y \in \mathcal{Y}$  ( $x, y$ )  $\mathcal{Y}$  random variable

$$C = \text{cost}(f(X), Y)$$

Goal : minimize  $E[C]$   
 $f \in \mathcal{F}$

Classification: (binary)  
or  
multi-class

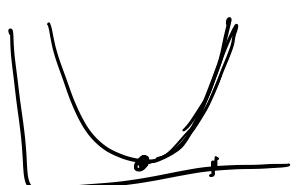
$$\underline{\text{cost}(\hat{y}, y)} = \begin{cases} 0 & \text{if } y = \hat{y} \\ 1 & y \neq \hat{y} \end{cases}$$

e.g. medical setting

- $\hat{y}$   $\neg$  have disease, no test
- $y$  have disease, test

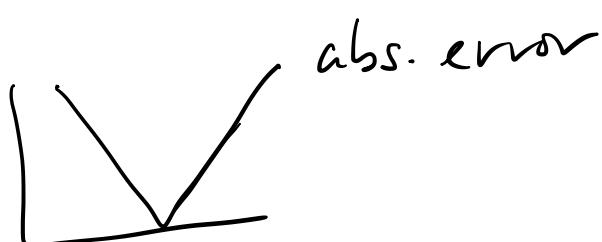
$\hat{y}$	$y$	$\neg$ have disease	have disease
0	0	0	Glansuit = 1000
1	1	$C_{\text{lab}} = 1$	$C_{\text{lab}} = 1$

Regression:  $y \in \mathbb{R}$



$$\text{cost}(\hat{y}, y) = (\hat{y} - y)^2 \quad (\text{squared error})$$

$$\text{cost}(\hat{y}, y) = |\hat{y} - y|$$



Classification, 0-1 cost  $X \subseteq \mathbb{R}^d$ ,  $y = \{1, 2, \dots, k\}$

$$\begin{aligned} E[C] &= \int_X \sum_{y \in Y} \underbrace{\text{cost}(f(x), y)}_{p(y|x) p(x)} \underbrace{p(y|x)}_{p(y|x) p(x)} dx \\ &= \int_X p(x) \left[ \sum_{y \in Y} \underbrace{\text{cost}(f(x), y)}_{E[C|X=x]} p(y|x) \right] dx \end{aligned}$$

$$\begin{aligned} f^*(x) &= \underset{y \in Y}{\operatorname{argmin}} E[C|X=x] & p(y|x) = \begin{cases} 0.9 & \\ 0.1 & \end{cases} \\ &= \underset{y \in Y}{\operatorname{argmin}} \sum_{y \in Y} \frac{\text{cost}(y, y)}{p(y|x)} & \downarrow \\ &= \underset{y \in Y}{\operatorname{argmax}} p(y|x) \end{aligned}$$

Regression: Conclusion:  $\hat{y} = E[Y|x]$   
 $\Rightarrow f^*(x) = E[Y|x]$

$$E[C] = \int_X \int_Y (f(x) - y)^2 p(y|x) dy \cdot p(x) dx$$

$\underbrace{(f(x) - y)^2}_{g(f(x))}$

$$g(\hat{y}) = \int_Y (\hat{y} - y)^2 p(y|x) dy$$

$$0 = \frac{dg(\hat{y})}{d\hat{y}} = \int_Y \frac{d(\hat{y} - y)^2}{d\hat{y}} p(y|x) dy = 2 \int_Y (\hat{y} - y) p(y|x) dy$$

$$= 2\hat{y} \underbrace{\int_Y p(y|x) dy}_{=1} - 2 \underbrace{\int_Y y p(y|x) dy}_{E[Y|x]}$$

$$1. f(x) = E[Y|X] \quad \text{or} \quad 2. f(x) \neq E[Y|X]$$

$$E[C] = \int_X p(x) \int_y (E[Y|X=x] - y)^2 p(y|x) dy dx$$

$= \text{Var}[Y|X=x]$

$$2. f(x) \neq E[Y|X]$$

$$(f(x) - y)^2 = (\underbrace{f(x) - E[Y|X]}_{\cdot} + \underbrace{E[Y|X] - y}_{\cdot})^2$$

$$= (f(x) - E[Y|X])^2 + 2(f(x) - E[Y|X])(E[Y|X] - y) + (E[Y|X] - y)^2$$

call this  $g(x,y)$

$$\begin{aligned}
 E[g(x, Y)] &= E\left[\frac{(f(x) - E[Y|x])(E[Y|x] - Y)}{x}\right] \\
 &= (f(x) - E[Y|x]) E\left[\frac{E[Y|x] - Y}{x}\right] \\
 &= 0
 \end{aligned}$$

$E[Y|x] - E[Y|x]$   
 $= 0$

$$\begin{aligned}
 E[C] &= E[(f(x) - Y)^2] \\
 &= E[(f(x) - E[Y|x])^2] + E[(E[Y|x] - Y)^2]
 \end{aligned}$$

$E[(f(x) - E[Y|x])^2]$   $p(y|x)$   
 $E[(E[Y|x] - Y)^2]$  variance of  $Y$

reducible error irreducible error

$f \neq f^*$

$$\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d \quad y \in \mathbb{R}$$

$$Y = \sum_{j=1}^d w_j^* X_j + \varepsilon$$

(X, Y)

$$\varepsilon \sim N(\mu=0, \sigma^2)$$

$$P(y|x) = N(\mu=x^T w^*, \sigma^2)$$

$$x^T w + \begin{array}{c} \text{a bell-shaped curve} \\ \mu=0 \end{array} \rightarrow$$

Fixed  $\sigma^2 > 0$

$$P(y|x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - x^T w)^2}{2\sigma^2}\right)$$

$$\sum_{i=1}^n \log P(y_i|x_i, w)$$

$$= \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{(y - x^T w)^2}{2\sigma^2}$$

$$\arg \min_{w \in \mathbb{R}^d} - \sum_{i=1}^n \ln p(y_i | x_i)$$

$$= \arg \min_{w \in \mathbb{R}^d} \sum_{i=1}^n \underbrace{(y_i - x_i^T w)^2}_{\textcircled{2\sigma^2}} = \arg \min_{w \in \mathbb{R}^d} \sum_{i=1}^n (y_i - x_i^T w)^2$$

$$\sum_{j \in \{1, \dots, d\}}$$

$$\frac{\partial}{\partial w_j} (y - x^T w)^2 = -2(y_i - x^T w) \frac{\partial x^T w}{\partial w_j}$$

$$\frac{\partial x^T w}{\partial w_j} = \frac{\partial}{\partial w_j} \sum_{i=1}^d w_i x_i = \frac{\partial}{\partial w_j} (w_1 x_1 + \dots + \underbrace{w_j x_j + \dots + w_d x_d}_{\partial w_j})$$

$$= \frac{\partial w_j x_j}{\partial w_j} = x_j$$

$$\frac{\partial}{\partial w_j} \sum_{i=1}^n (y_i - x_i^T w)^2 = -2 \sum_{i=1}^n (y_i - x_i^T w) \underbrace{x_{i,j}}_{\forall j \in \{1, \dots, d\}} = 0$$

$$-\sum_{i=1}^n (y_i - x_i \cdot \bar{w}) x_{ij} = -\sum_{i=1}^n y_i x_{ij} + \underbrace{\sum x_{ij} x_i \cdot \bar{w}}$$

$$\frac{\partial L}{\partial w_1} = \left[ \sum_{i=1}^n y_i x_{ii} \right] \xrightarrow{\frac{\partial}{\partial w_1}} \dots \\ \vdots \\ \left[ \sum_{i=1}^n y_i x_{id} \right] \xrightarrow{\frac{\partial}{\partial w_d}}$$

$$\underbrace{\sum_{i=1}^n y_i x_i}_b$$

$$\left[ \begin{array}{c} \sum_{i=1}^n x_{ii} x_i \cdot \bar{w} \\ \vdots \\ \sum_{i=1}^n x_{id} x_i \cdot \bar{w} \end{array} \right] \begin{array}{c} j=1 \\ \vdots \\ j=d \end{array} = \textcircled{0}$$

$$\left[ \sum_{i=1}^n x_i x_i^\top \bar{w} \right] = \textcircled{0}$$

$$b - Aw = \textcircled{0} \Rightarrow Aw = b$$

If  $A$  invertible

$w = A^{-1}b.$

$$A = \sum_{i=1}^n x_i x_i^\top \quad \begin{array}{l} n = \# \text{samples} \\ d = \# \text{features} \end{array} \quad X = \begin{bmatrix} \vec{x}_1^\top \\ \vdots \\ \vec{x}_n^\top \end{bmatrix} = \begin{bmatrix} x_{11} & \cdots & x_{1d} \\ x_{21} & \cdots & x_{2d} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nd} \end{bmatrix}$$

$$b = \sum_{i=1}^n x_i y_i = X^\top y$$

$$A = X^\top X = \sum_{i=1}^n x_i x_i^\top$$

$$X = U \Sigma V^\top \quad \begin{array}{c} n \times d \\ n \times n \\ n \times d \\ d \times d \end{array}$$

$$U^\top U = I$$

$$U U^\top = I$$

$$d < n$$

$$\Sigma = \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_d \\ \hline & 0 & \end{bmatrix} \quad \sigma_d = 0$$

zero matrix

$$X^T X = (\mathcal{U} \Sigma V^T)^T (\mathcal{U} \Sigma V^T) \quad (AB)^T = B^T A^T$$

$$= \sqrt{\Sigma^T \underbrace{\mathcal{U}^T \mathcal{U}}_{=I} \Sigma} V^T \quad \mathcal{U}^T \mathcal{U} = I$$

$$= \underbrace{\sqrt{\Sigma^T \Sigma}}_{d \times d, d \times n} \underbrace{\sqrt{V^T V}}_{n \times d, d \times d}$$

$$\Sigma^T \Sigma = \begin{bmatrix} \sigma_1^2 & & & 0 \\ & \ddots & & \\ 0 & & \ddots & 0 \\ & & & \sigma_d^2 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & & & 0 \\ & \ddots & & \\ 0 & & \ddots & \sigma_d^2 \\ & & & 0 \end{bmatrix}$$

$$A = X^T X$$

$$= \sqrt{\Lambda} V^T$$

$$\Lambda = \Sigma^T \Sigma$$

$$\sigma_d^2 = 0$$

$$Aw = \sqrt{(\Sigma^T \Sigma)} \underbrace{V^T w}_{\tilde{w}} = \sqrt{\begin{bmatrix} \sigma_1^2 & & & 0 \\ & \ddots & & \\ 0 & & \ddots & \sigma_d^2 \end{bmatrix}} \tilde{w}$$

$$= \sqrt{\begin{bmatrix} \sigma_1^2 \tilde{w}_1 \\ \vdots \\ \sigma_d^2 \tilde{w}_d \end{bmatrix}} = \sqrt{\begin{bmatrix} \sigma_1^2 \tilde{w}_1 \\ \vdots \\ \sigma_{d-1}^2 \tilde{w}_{d-1} \\ 0 \end{bmatrix}}$$

$$Aw = b$$

$$\min_w \|(Aw - b)\|_2^2$$

$$\underline{Aw = b}$$

$$\min_w \|Xw - y\|_2^2$$

$$Xw \neq y$$
$$\hat{y} \neq y$$

$$\sum_d i^{-2} \sum_{d \times n}^T u_{n \times n}$$

$$\sum^T u^T = \sum_d u_d^T$$

$$\sum_d [0] u^T = \sum_d u_d^T$$

$$\begin{bmatrix} u_d^T \\ \vdots \\ u_{d+1:n}^T \end{bmatrix}$$

$$\min_w \sum_{i=1}^n (x_i^T w - y_i)^2 = \min_w \|Xw - y\|_2^2 \leftarrow \|\hat{y} - y\|_2^2$$

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} = \begin{bmatrix} x_{11} & \cdots & x_{1d} \\ x_{21} & \cdots & x_{2d} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nd} \end{bmatrix} \quad \begin{array}{l} x_i \in \mathbb{R}^d \\ \text{sgn}(x_{ij}) \in \mathbb{R} \\ x_{ij} \end{array} \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$\|z\|_2 = \sqrt{\sum_{i=1}^n z_i^2} \quad \|z\|_2^2 = \sum_{i=1}^n z_i^2$$

$$Xw = \hat{y} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$X = U \Sigma V^T \quad \text{where} \quad \begin{array}{c} n \times n \\ n \times d \\ d \times d \end{array}$$

$$\frac{U^T U}{V^T V} = I$$

System:  $\underbrace{X^T X}_{A \text{ } d \times d} \underbrace{w}_{b} = \underbrace{X^T y}_{\text{diagonal matrix}}$

$$\Sigma = \begin{bmatrix} \Sigma_d \\ 0 \end{bmatrix}$$

$$\Sigma_d = \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_d \end{bmatrix}$$

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d > 0$$

If  $A$  is invertible,  $A^{-1}$  is inverse,  $w = A^{-1}b$

---

(i.e.  $A^{-1}A = I$ ,  $Aw = b \Rightarrow \underbrace{A^{-1}Aw}_{I} = A^{-1}b \Rightarrow w = A^{-1}b$ )

$$X^T X = (U \Sigma V^T)^T (U \Sigma V^T) \stackrel{\text{dxd, dxn, nxn, dxd}}{=} V \underbrace{\Sigma^T \Sigma}_{\substack{\text{dxd} \\ \text{dxd}}} V^T = V \sum_d \sigma_d^2 V^T$$

$$= V \underbrace{\Sigma^T \Sigma}_{\substack{= I}} V^T = \underbrace{V \sum_d \sigma_d^{-2} V^T}_{\text{check.}} \quad (A^{-1}A = ? I)$$

$$(X^T X)^{-1} X^T = \underbrace{V \sum_d \sigma_d^{-2} V^T}_{\substack{= I}} V \Sigma^T U^T \quad X^T = V \sum_d U^T$$

$$= \underbrace{V \sum_d \sigma_d^{-2} \Sigma^T U^T}_{\substack{\text{pseudo-inv}}} \quad \begin{bmatrix} \sigma_1^{-2} & & & \\ & \ddots & & \\ & & \sigma_d^{-2} & \\ & & & \end{bmatrix} \begin{bmatrix} 0_{1 \times 1} & & & \\ & \ddots & & \\ & & 0_d & \\ & & & 0 \end{bmatrix}$$

$$= V \sum_d \sigma_d^{-1} U_d^T \quad \leftarrow X^+ \quad U_d = \text{first } \frac{\sum_d \sigma_d^{-2}}{d \text{ columns}}$$

$$w = V \sum_d \sigma_d^{-1} U_d^T y \quad \leftarrow \quad \begin{bmatrix} \Sigma^T \\ X^T X = I \end{bmatrix} \quad \text{check}$$

$$w = \underline{X^T y} \quad \text{if } X \text{ full rank} \quad \sigma_1 \geq \dots \geq \sigma_d > 0$$

or  $X$  not full rank, then  $w = X^T y$   
is still solution

$$w = X^T y = \left( \sum_{j=1}^r \frac{v_j u_j^\top}{\sigma_j} \right) y$$

$$r = \text{rank}(X) \leq d.$$

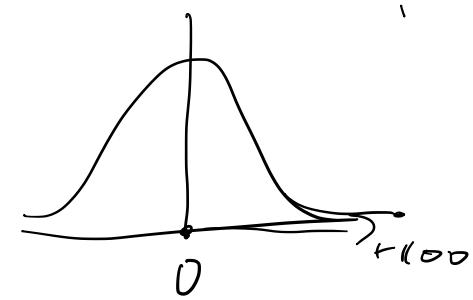
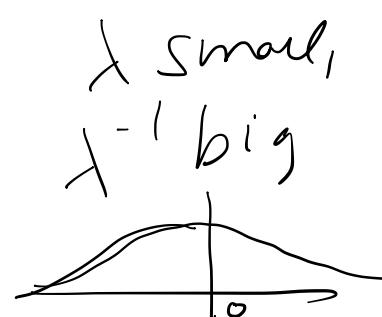
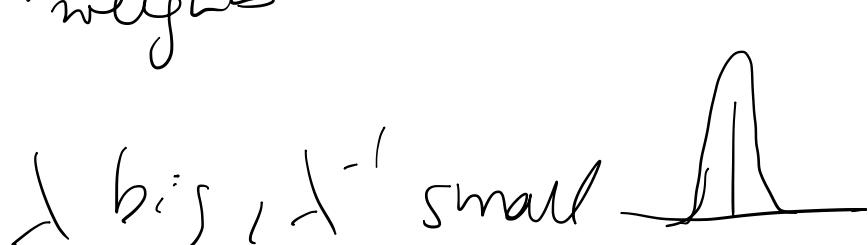
$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r \geq \sigma_{r+1} = \sigma_{r+2} = \dots = \sigma_d = 0$$

If  $\sigma_r = 10^{-6}$   
 $\Rightarrow$  might cause big  $w$

$$\Sigma_d^+ = \begin{bmatrix} \sigma_1^{-1} & & & \\ & \ddots & & \sigma_r^{-1} \\ & & 0 & 0 \\ & & 0 & 0 \end{bmatrix}$$

---

prior on weights :  $w_j \sim N(\mu=0, \lambda^{-1}) \quad \forall j \in \{1, \dots, d\}$



$$\min_w - \sum_{i=1}^n \ln p(y_i | x_i, w) - \ln p(w)$$

*drop constants*

$$-\ln p(w_j) = \boxed{-\ln \frac{1}{\sqrt{2\pi\lambda^{-1}}} + \frac{\lambda w_j^2}{2}}$$

$$p(w_j) = \frac{1}{\sqrt{2\pi\lambda^{-1}}} \exp\left(-\frac{(w_j - 0)^2}{2\lambda}\right)$$

*prior*  
 $\lambda_j$

$$p(\vec{w}) = \prod_{j=1}^d p(w_j) \Rightarrow \ln p(w) = \sum_{j=1}^d \ln p(w_j)$$

$$\text{MAP: } \arg \min_w \sum_{i=1}^n (x_i^\top w - y_i)^2 + \frac{\lambda}{2} \underbrace{\sum_{j=1}^d w_j^2}_{\|w\|_2^2}$$

*Gaussian prior*

$$\equiv \arg \max_w P(\mathcal{D}|w) p(w)$$

*Gaussian p(y|x)*

Procedure: ① take log ② drop constants ③ min negative instead of max

MLE  $\equiv$  min negative log likelihood, MAP: min - log posterior

$$w_{MLE} = \underline{(X^T X)^{-1} X^T y}$$

or  $w_{MLE} = X^T y$

$$w_{MAP} = (X^T X + \lambda I)^{-1} X^T y$$

$$\lambda I = \begin{bmatrix} \lambda & & & \\ & \lambda & & 0 \\ & & \ddots & \\ 0 & & & \lambda \end{bmatrix}$$

$$\underline{X^T X + \lambda I} = V \sum_d \sigma_d^2 V^T + \lambda \underbrace{I}_{V V^T} = V \sum_d \sigma_d^2 V^T + V(\lambda I) V^T$$

$$= \underline{V \left( \sum_d \sigma_d^2 + \lambda I \right) V^T}$$

$$w_{MLE} = \sum_{j=1}^r v_j u_j^T y$$

$\boxed{\sigma_j}$  could be small

$$w_{MAP} = \sum_{j=1}^d \sigma_j (v_j u_j^T) y$$

$\boxed{\sigma_j^2 + \lambda}$  ?

$$\frac{(X^T X + \lambda I)^{-1} X^T}{\left( \sum_d \sigma_d^2 + \lambda I \right)^{-1} \sum_d}$$

Exercise : 16  
 $\lambda > 0$ , is  $(X^T X + \lambda I)$   
 invertible

$$X \rightarrow \text{scale} \quad 2X \quad A^2 = AA$$

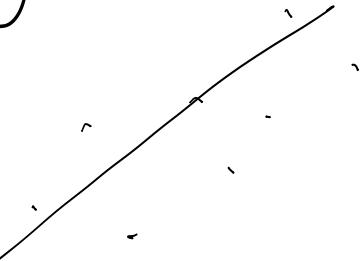
$n \times d \quad n \times d$

~~$\Sigma\Sigma$~~

$$X_w \approx y \quad 2X_w \approx y \quad \Sigma^\top \Sigma \neq \cancel{\Sigma^2}$$

$w = \frac{1}{2} w_{\text{best}}$

$\Sigma^3?$



$w_{\text{best}}$

$\lambda = 1$

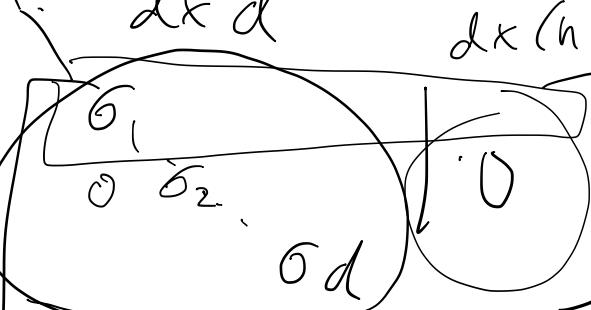
$$2X \frac{1}{2} w_{\text{best}} = Xw_{\text{best}}$$

$\lambda$  should be  $\sqrt{2}$

$$\Sigma^\top \Sigma =$$

$\overset{\text{Ed.}}{\circlearrowleft} \quad \overset{\text{Ed.}}{\circlearrowright}$

$d \times d \quad d \times (n-d)$



$\begin{bmatrix} [\sigma_1^2] & \cdots & [d] \\ \vdots & \ddots & \vdots \\ [0] & \cdots & [0] \end{bmatrix} \quad \begin{bmatrix} [I_d] & \cdots & [0] \\ \vdots & \ddots & \vdots \\ [0] & \cdots & [0] \end{bmatrix}_{d \times (n-d)}$

$$= \begin{bmatrix} \sigma_1^2 & & & \\ 0 & \ddots & & \\ 0 & & \ddots & \\ 0 & 0 & \cdots & \sigma_d^2 \end{bmatrix} = \Sigma_d^2$$

$w_{MLE}(\mathcal{D})$

random variable

$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ ,  $X = [x_1^\top \dots x_n^\top]^\top$ ,  $y = [y_1 \dots y_n]^\top$

random variable  
because it's a sample  $\sim N(0, \sigma^2)$

$$w_{MLE}(\mathcal{D}) = (X^\top X)^{-1} X^\top y$$

unbiased?

$$E[w_{MLE}(\mathcal{D})] = E[(X^\top X)^{-1} X^\top y]$$

$w$  = true underlying weight

$$= E[(X^\top X)^{-1} X^\top (Xw + \varepsilon)]$$

$\varepsilon$  noise in targets

$$= E[\underbrace{(X^\top X)^{-1} X^\top X w}_\text{= w} + \underbrace{E[(X^\top X)^{-1} X^\top \varepsilon]}_\text{= 0}]$$

$$= E[w] + \underbrace{E[(X^\top X)^{-1} X^\top]}_{X^+} \underbrace{E[\varepsilon]}_\text{= 0}$$

independent RVs A, B

$$E[AB] = E[A]E[B]$$

$$= w$$

$$w_{MLE}(\mathcal{D}) = w + X^+ \varepsilon$$

$$\begin{aligned}
\text{Cov}(w_{MLE}(\mathcal{D})) &= E \left[ (w_{MLE}(\mathcal{D}) - w) (w_{MLE}(\mathcal{D}) - w)^T \right] \\
&= E \left[ w_{MLE}(\mathcal{D}) w_{MLE}(\mathcal{D})^T \right] - w w^T \\
&= E \left[ (w + X^+ \varepsilon) (w + X^+ \varepsilon)^T \right] - w w^T \\
&= E(w w^T) + \underbrace{E[X^+ \varepsilon w^T]}_{\substack{E[X^+] \\ E[\varepsilon] \\ = 0}} + \underbrace{E[w (X^+ \varepsilon)^T]}_{\substack{w \\ E[\varepsilon^T] \\ = 0}} + \underbrace{E[X^+ \varepsilon (X^+ \varepsilon)^T]}_{-w w^T} \\
&= \cancel{w w^T} + \underbrace{E[X^+ \varepsilon \varepsilon^T X^{+T}]}_{\substack{E[X^+] \\ E[\varepsilon] \\ = 0}} - \cancel{w w^T} \\
&= E[X^+ \varepsilon \varepsilon^T X^{+T}] \\
&= E[E[X^+ \varepsilon \varepsilon^T X^{+T} | X]] \\
&= E[X^+ E[\varepsilon \varepsilon^T | X] X^{+T}] \\
&= \sigma^2 E[X^+ X^{+T}] = \sigma^2 E \left[ \sum_{j=1}^n v_j v_j^T \right] = \sigma^2 I
\end{aligned}$$

Recall  
 $E[A] = E[E[A|B]]$   
 $E[\varepsilon \varepsilon^T | X] = E[\varepsilon \varepsilon^T]$   
 $= E[\varepsilon \varepsilon^T] - E[\varepsilon] E[\varepsilon]^T$   
 $= \text{cov}(\varepsilon) = \sigma^2 I$

$$w_{MAP}(\vartheta) = (X^T X + \lambda I)^{-1} X^T Y$$

$$\begin{aligned} E[w_{MAP}(\vartheta)] &= E[(X^T X + \lambda I)^{-1} X^T (Xw + \varepsilon)] \\ &= \underbrace{E[(X^T X + \lambda I)^{-1} X^T X] w}_{\neq I} + 0 \\ &\neq w \end{aligned}$$

$\lambda = 0$ , unbiased  
MLE

$n$  large, close to  
unbiased

$$X^T X = \sum_{i=1}^n x_i x_i^T$$

$$X^T X + \lambda I = \underbrace{\sum_{i=1}^n x_i x_i^T}_{\perp \sum_{i=1}^n x_i y_i} + \lambda I$$

$$\frac{1}{n} X^T X = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

$$w_{MLE}(\vartheta) = (\frac{1}{n} X^T X)^{-1} (\frac{1}{n} X^T y)$$

$$w_{MAP}(\vartheta) = (\frac{1}{n} X^T X + \frac{\lambda}{n} I)^{-1} (\frac{1}{n} X^T y)$$

$$\text{Bias}(\hat{w}_{MAP}(\vartheta))^2 = \|\mathbb{E}[\hat{w}_{MAP}(\vartheta)] - w\|_2^2$$

$$\begin{aligned}\text{Bias}(\hat{w}_{MLE}(\vartheta))^2 &= \|\underbrace{\mathbb{E}[\hat{w}_{MLE}(\vartheta)]}_w - w\|_2^2 \\ &= \|w - w\|_2^2 = 0\end{aligned}$$

Assume  $Y = f(X) + \varepsilon$   $\varepsilon \sim N(\mu=0, \sigma^2)$

$$\hat{f}(x) = x^T \hat{w}$$

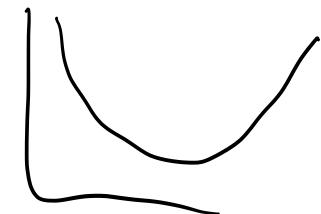
$w$  true parameters for  $f$

MSE = mean squared error

$$MSE(\hat{w}) = E[\|\hat{w} - w\|_2^2]$$

$\hat{w}$  is an RV

$$= \underbrace{\|E[\hat{w}] - w\|_2^2}_{\text{Bias}(\hat{w})^2} + \text{Var}(\hat{w})$$



more generally: mse for one sample  $(x, y)$

$$E[(\hat{f}(x) - f(x))^2] = \underbrace{(E[\hat{f}(x)] - f(x))^2}_{\text{Bias}} + \text{Var}(\hat{f}(x))$$

$$MSE(\hat{f}) = E[(\hat{f}(X) - f(X))^2]$$

reducible

irreducible

$$E[(Y - \hat{f}(X))^2] = E[(\hat{f}(X) - f(X))^2] + \sigma^2$$

Objective function  $c: \mathbb{R}^d \rightarrow \mathbb{R}$   $\min_w c(w)$

Starting simple:  $d=1$

$$c(w) = \sum_{n=0}^{\infty} \frac{c^{(n)}(w_0)}{n!} (w-w_0)^n$$

$$\nabla c(w) = 0$$

Second-order

locally around  $w_0$

$$\frac{d}{dw} c(w) = 0$$

Start at  $w_0$

$$c(w) \approx \hat{c}_0(w) = c(w_0) + \underline{(w-w_0)} \underline{c'(w_0)} + \frac{1}{2} (w-w_0)^2 c''(w_0)$$

$$\frac{d}{dw} \hat{c}_0(w) = 0 + c'(w_0) + \underline{(w-w_0)} \underline{c''(w_0)} = 0$$

$$\Rightarrow w_1 = w_0 - \frac{c'(w_0)}{c''(w_0)}$$

$w_1 \neq \min_w c(w)$   
 $c''(w_1) > 0$

$$w_{t+1} = w_t - \frac{c'(w_t)}{c''(w_t)}$$

eventually stationary  
 $w_t \rightarrow$  point of  $c$

First-order:  $\hat{c}_0(w) = c(w_0) + (w - w_0)c'(w_0) + \frac{1}{2\eta} (w - w_0)^2$

$$\Rightarrow w_{t+1} = w_t - \underbrace{\eta}_{\text{stepsize.}} c'(w_t)$$

$O((w - w_0)^2)$   
approximates

Stepsize for 1<sup>st</sup>:  $\eta$ .  
 Stepsize for 2<sup>nd</sup>:  $\frac{1}{c''(w_t)}$

$$c(w)$$



$c''(w)$  big.

stepsize small.

$$c(w)$$



$c''(w)$  small.

stepsize big.

$$-\frac{c''(w_0)}{2} (w - w_0)^2$$

Recall: 2<sup>nd</sup>.

$$w_{t+1} = w_t - \frac{c'(w_t)}{c''(w_t)}$$

$$\nabla c(w) = \begin{bmatrix} \frac{\partial c}{\partial w_1}(w) \\ \frac{\partial c}{\partial w_2}(w) \\ \vdots \\ \frac{\partial c}{\partial w_d}(w) \end{bmatrix} \quad H_{c(w)} = \begin{bmatrix} \frac{\partial^2 c}{\partial w_1^2}(w) & \dots & \frac{\partial^2 c}{\partial w_1 \partial w_d}(w) \\ \frac{\partial^2 c}{\partial w_2 \partial w_1}(w) & \ddots & \vdots \\ \vdots & \ddots & \frac{\partial^2 c}{\partial w_d^2}(w) \end{bmatrix}$$

$$H_{c(w)}[i,j] = \frac{\partial^2 c}{\partial w_i \partial w_j}(w) \quad \nabla c(w) = 0$$

$$c(w) \approx \hat{c}_0(w) = c(w_0) + \underbrace{\nabla c^\top(w_0)(w - w_0)}_{+} + \frac{1}{2} (w - w_0)^\top \underbrace{H_{c(w_0)}}_{D} (w - w_0)$$

$$\sum_{j=1}^d \frac{\partial c}{\partial w_j}(w_0)(w_j - w_0[j])$$

$$\frac{\|w - w_0\|_{H_{c(w_0)}}^2}{2}$$

$$c(w) \approx \hat{c}_0(w) = c(w_0) + \nabla c(w_0)^T (w - w_0) + \frac{1}{2} \eta^{-1} (w - w_0)^T (w - w_0)$$

First order

Now  $H_{c(w_0)} \approx \frac{1}{\eta} I, n > 0$

$$\Rightarrow 0 = \nabla \hat{c}_0(w) = \nabla c(w_0) + \eta^{-1} (w - w_0)$$

$$\Rightarrow w = w_0 - \eta \nabla c(w_0)$$



How pick  $\eta$ ? Why do first order?

Cost of one iteration is 1st order (in terms of n and d)

$$c(w) = \sum_{i=1}^n c_i(w) \quad \text{e.g. } c_i(w) = \frac{1}{2} (x_i^T w - y_i)^2$$

$$\nabla c_i(w) = (x_i^T w - y_i) x_i$$

Compute  $\nabla c(w_0) = \sum_{i=1}^n \underline{\nabla c_i(w_0)}$   $O(d)$

$\nearrow$  costs  $O(nd)$

Cost of 2<sup>nd</sup>-order:

$$w_t = w_{t-1} - \underbrace{H_c^{-1}(w_{t-1})}_{\text{costs } O(nd)} \nabla_c(w_{t-1})$$

dxd

$$H_c(w) = \sum_{i=1}^n H_{c_i}(w)$$

Exercise: Show it costs  $O(nd^2)$

Computing inverse: about  $O(d^3)$

2nd-order  
Cost:  $O(nd^2 + d^3)$

First order:  
 $O(nd)$

Superlinear, takes 2 iterations

Linear, take 15 iterations

Big d, bad for 2<sup>nd</sup> order

Small d worth it to do 2<sup>nd</sup> order

Huge n  $\rightarrow$  neither, we should <sup>do</sup> stochastic gradient descent