

H03-1

a) Name: Jihun Wan Email Address: jhun0324@berkeley.edu

Description of Team: Best Group Ever

How did I work?

Comments:

b) I certify that all solutions are entirely in my words and that I have not looked at another student's solutions. I have credited all external sources in this write up.

Qingyang Zhao

Homework

$$a) \int_{-\infty}^y P_{Y|X}(t|x) dt = P(Y \leq y | X=x) \quad X, Y \text{ denote RV. } x, y: \text{the value in Alphabet Set}$$

$$= P(Xw+b + Z \leq y | X=x)$$

$$= P(Z \leq y - xw - b) \quad = \int_{-\infty}^{y-xw-b} p(z) dz \quad z \sim N(0,1)$$

$$P_{Y|X}(t|x) = \frac{\int_{-\infty}^y P_{Y|X}(t|x) dt}{dy} = \frac{\int_{-\infty}^y p(z) dz}{dy} = \frac{1}{\sqrt{2\pi}} e^{-\frac{(y-xw-b)^2}{2}} \quad Y | X \sim N(xw+b, 1)$$

$$b) \underset{w,b}{\text{maximize}} \sum_{i=1}^n P_{Y|X_i}(y_i|x_i) \Leftrightarrow \underset{w,b}{\text{maximize}} \log \prod_{i=1}^n P_{Y|X_i}(y_i|x_i)$$

$$\text{Set } f(w,b) = \sum_{i=1}^n \log P_{Y|X_i}(y_i|x_i)$$

$$f(w,b) = \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}} e^{-\frac{(y_i-wx_i-b)^2}{2}} = \left( \sum_{i=1}^n -\frac{(y_i-wx_i-b)^2}{2} \right) - n \log \sqrt{2\pi}$$

Thus,  $\underset{w,b}{\text{maximize}} f(w,b)$  is equivalent to  $\underset{w,b}{\text{minimize}} \sum_{i=1}^n (y_i - wx_i - b)^2$ .

$$\underset{w,b}{\text{minimize}} \left\| \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} b \\ w \end{bmatrix} \right\|_2^2 \quad \text{let } \vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \vec{x} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}, \vec{w} = \begin{bmatrix} b \\ w \end{bmatrix}$$

$$\underset{w,b}{\text{minimize}} \| \vec{y} - \vec{x} \vec{w} \|_2 \quad \vec{w} = (\vec{x}^\top \vec{x})^{-1} \vec{x}^\top \vec{y} = \begin{bmatrix} b \\ w \end{bmatrix}$$

$$c) \quad Z \sim U[-0.5, 0.5] \quad P_2(z) = \begin{cases} 1, & z \in [-0.5, 0.5] \\ 0, & \text{otherwise.} \end{cases}$$

$$\int_{-\infty}^y P_{Y|X}(t|x) dt = P(Y \leq y | X=x) = \mathbb{P}(WZ + Z \leq y | X=x)$$

$$= \int_{-\infty}^y P_Z(z) dz = \mathbb{P}(Z \leq y - wz) = \frac{\mathbb{P}(P_{Y|X}(y|x))}{\mathbb{P}_Z(y)} = \frac{\int_{-\infty}^y P_{Y|X}(t|x) dt}{\mathbb{P}_Z(y)} = \frac{\int_{-\infty}^y P_{Y|X}(t|x) dt}{\mathbb{P}_Z(y-wx)}$$

$$P_{Y|X}(y|x) = P_Z(y-wx) = \begin{cases} 1, & y \in [-0.5+wX, 0.5+wX] \\ 0, & \text{otherwise} \end{cases}$$

$$Y|X \sim \cup [-0.5+wX, 0.5+wX]$$

(d) Criterion: find a "w", such that the number of points lies in  $[-0.5+wX_i, 0.5+wX_i]$  is maximized.

when  $X_i < 0$

$$-0.5+wX_i \leq Y_i \leq 0.5+wX_i, \quad y_i - 0.5 \leq wX_i \leq y_i + 0.5 \quad y_i - 0.5 \leq w \leq \frac{y_i - 0.5}{X_i} \quad | X_i < 0$$

$$\frac{y_i - 0.5}{X_i} \leq w \leq \frac{y_i + 0.5}{X_i} \quad \frac{y_i}{X_i} - \frac{0.5}{|X_i|} \leq w \leq \frac{y_i}{X_i} + \frac{0.5}{|X_i|}$$

$$y_i > 0 \quad \frac{y_i - 0.5}{X_i} \leq w \leq \frac{y_i + 0.5}{X_i} \quad \text{Thus we have } \max \frac{y_i}{X_i} - \frac{0.5}{|X_i|}, \min \frac{y_i}{X_i} + \frac{0.5}{|X_i|}$$

(e)

As n gets large, the area of w getting smaller

$$f_1: P_{W|X_i,W}(\omega | \{x_i = x_i, y_i = y_i\}) = \frac{P(\{x_i, y_i\}, \omega)}{\int_{\mathbb{R}^n} P(\{x_i, y_i\}, \omega) d\omega} = \frac{\{x_i, y_i\} \text{ denote } \{x_j, y_j, \dots, x_n, y_n\}}{\int_{\mathbb{R}^n} P(\{x_i, y_i\}, \omega) d\omega} \quad \begin{matrix} \text{noise are ind.} \\ \text{P} \text{ denote PDF of R.V.} \end{matrix}$$

$$P_{(X_i, Y_i)|W}(\{x_i, y_i\}, \omega) = \prod_{j=1}^n \left( P_{Y_j|X_j, W} \right) P_{X_i, Y_i|W} = \prod_{j=1}^n \left( P_{Y_j|X_j, W} \right) \prod_{j=1}^n P_{X_j} (P_W) \quad (1)$$

$$\int_{\mathbb{R}^n} P_{(X_i, Y_i), W}(\{x_i, y_i\}, \omega) d\omega = \prod_{j=1}^n \int_{\mathbb{R}} P_{X_j} \int_{-\infty}^{\infty} P_{Y_j|X_j, W} P_W d\omega \quad (2)$$

$$P_{W|X_i, Y_i} = \frac{(1)}{(2)} = \frac{\prod_{j=1}^n \int_{\mathbb{R}} P_{X_j} P_{Y_j|X_j, W} P_W d\omega}{\int_{\mathbb{R}^n} \prod_{j=1}^n \int_{\mathbb{R}} P_{X_j} P_{Y_j|X_j, W} P_W d\omega} \leftarrow \text{Some constant} \quad !$$

Thus, [next Page]

103-2

$$\Rightarrow f) \text{ Thus, } P_{W|\{x_i, y_i\}} = \frac{1}{n} e^{-\frac{(w - \mu_0)^2}{2\sigma_0^2}}$$

All we need to do is to calculate  $\{\mu_0, \sigma_0\}$   
which is inside  $\prod_{i=1}^n (P_{Y_i|X_i, W}) P_W$

$$\text{Some constant } A_0 \frac{1}{\prod_{i=1}^n \sigma_i^2} = A_0 e^{\frac{-\sum_{i=1}^n (y_i - x_i w)^2}{2\sigma_0^2}}$$

I don't care

Goal:

$$\frac{(w - \mu_0)^2}{2\sigma_0^2}$$

Given:

$$\begin{aligned} & \frac{\sum_{i=1}^n (y_i - x_i w)^2}{2} - \frac{w^2}{2\sigma^2} \\ & - \frac{\sum_{i=1}^n y_i^2 + \sum_{i=1}^n x_i^2 w^2 - 2 \sum_{i=1}^n x_i y_i w}{2} - \frac{w^2}{2\sigma^2} \\ & - \frac{(\sum_{i=1}^n x_i^2 + \frac{1}{\sigma^2}) w^2 + 2 \sum_{i=1}^n x_i y_i w}{2} + \text{const} \end{aligned}$$

$$\begin{cases} \mu_0 = -\frac{\sum_{i=1}^n x_i y_i}{2} / \left( \sum_{i=1}^n x_i^2 + \frac{1}{\sigma^2} \right) \\ \sigma^2 = \frac{1}{\sum_{i=1}^n x_i^2 + \frac{1}{\sigma^2}} \end{cases}$$

$$E[P_{W|\{x_i, y_i\}}] = -\frac{\sum_{i=1}^n x_i y_i}{2}$$

g) maximize  $\prod_{i=1}^n P_{Y_i|X_i, W}(y_i | x_i)$   $\Leftrightarrow$  maximize  $\log \prod_{i=1}^n P_{Y_i|X_i, W}(y_i | x_i)$

$$\max_w \sum_{i=1}^n \log P_{Y_i|X_i, W}(y_i | x_i) \Leftrightarrow \min_w \sum_{i=1}^n (y_i - w^\top x_i)^2$$

$$\Leftrightarrow \min_w \sum_{i=1}^n \|y_i - w^\top x_i\|^2$$

$$\hat{w} = X(X^\top X)^{-1} Y$$

h)  $P_{\text{ML}}(x_i, y_i)$  is a multivariable Gaussian

It is just the expansion of f)

$$w \sim N(\mu, \Sigma)$$

Proof: consider  $\prod_{i=1}^n P_{\text{ML}}(x_i, y_i)$

$$\begin{aligned} & -\sum_{i=1}^n (y_i - w^T x_i)^2 - \frac{1}{2} w^T \sum_{i=1}^n w \\ & = C_0 \cdot \frac{\sum_{i=1}^n (w^T x_i)^2 - 2 \sum_{i=1}^n w^T x_i y_i}{2} - \frac{1}{2} w^T \sum_{i=1}^n w + C_1 \\ & \Leftarrow -\frac{1}{2} w^T \left( \sum_{i=1}^n x_i x_i^T + \sum_{i=1}^n w \right) w + \sum_{i=1}^n w^T x_i y_i + C_1 \end{aligned}$$

$$\bar{\mu} = \left( \sum_{i=1}^n x_i x_i^T + \sum_{i=1}^n w \right)^{-1} \sum_{i=1}^n x_i y_i \quad \Sigma = \left( \sum_{i=1}^n x_i x_i^T + \sum_{i=1}^n w \right)^{-1}$$

i) As  $\sigma^2$  becomes large, choice of  $w$  becomes large, because more  $w$  can be chosen with the same probability.

As training sample becomes large, the estimation becomes better.

But  $n \gg 0$  and the ~~will have a~~ will have a better estimate.

1102 - 2

$$a) \left[ E\left[ \frac{X_1 + X_2 + \dots + X_n}{n} - \mu \right] \right] = E[X] - \mu = 0$$

$$E\left[ \frac{X_1 + X_2 + \dots + X_n}{n+1} - \mu \right] = \frac{n}{n+1} E[X] - \mu = -\frac{1}{n+1} \mu$$

$$E\left[ \frac{X_1 + X_2 + \dots + X_n}{n+n_0} - \mu \right] = \frac{n}{n+n_0} E[X] - \mu = -\frac{n_0}{n+n_0} \mu$$

$$E[\sigma - \mu] = -\mu$$

$$b) \quad \text{var}\left[ \frac{X_1 + X_2 + \dots + X_n}{n} \right] = \frac{1}{n^2} \cdot n \text{Var}[x] = \frac{1}{n} \text{Var}[x] = \frac{1}{n} \sigma^2$$

$$\text{Var}\left[ \frac{X_1 + X_2 + \dots + X_n}{n+1} \right] = \frac{n}{(n+1)^2} \text{Var}[x] = \frac{n}{(n+1)^2} \sigma^2$$

$$\text{Var}\left[ \frac{X_1 + X_2 + \dots + X_n}{n+n_0} \right] = \frac{n}{(n+n_0)^2} \text{Var}[x] = \frac{n}{(n+n_0)^2} \sigma^2$$

$$\text{Var}[0] = 0$$

c) Expected total Error

$$\text{Error}_1 = \text{Var}[\hat{x}] + E[(\hat{x} - \mu)]^2 = \frac{1}{n} \sigma^2 \quad \text{Error}_2 = \frac{n}{(n+n_0)^2} \sigma^2 + \frac{n_0^2}{(n+n_0)^2} \mu^2$$

$$\text{Error}_3 = \frac{n}{(n+1)^2} \sigma^2 + \left(\frac{1}{n+1}\right)^2 \mu^2 \quad \text{Error}_4 = \mu^2$$

- d) Simply see the Total Error.  
 Error<sub>1</sub> is when  $n_0 = 0$ . Error<sub>2</sub> is  $n_0 = 1$   
 Error<sub>3</sub> is  $n_0 = n$ .

c) minimize.  $\left(\frac{n\alpha}{1+n\alpha}\right)^2 \mu^2 + \frac{n}{(1+n\alpha)^2} \sigma^2$

||

$$f(\alpha) = \left(\frac{\alpha}{1+\alpha}\right)^2 \mu^2 + \frac{1}{n} \frac{1}{(1+\alpha)^2} \sigma^2$$

$$\frac{d(f(\alpha))}{d\alpha} = \frac{2\mu^2\alpha - \frac{2}{n}\sigma^2}{(1+\alpha)^3}$$

$$\alpha \in (-\infty, -1) \quad \frac{df}{d\alpha} > 0$$

$$(-1, -\frac{\sigma^2}{n\mu^2}) \quad < 0$$

$$(-\frac{\sigma^2}{n\mu^2}, +\infty) \quad > 0$$

Thus when  $\alpha = \frac{\sigma^2}{n\mu^2}$   $f(\alpha)$  reaches minimum.

f)  $\alpha$  becomes rather big. If we know  $\mu$  should be small,  $\sigma$ , large we should put more sample  $X=0$  to estimate a better mean, with a low total error.

g)  $X' = X - \mu_0 \quad E[X'] = E[X] - \mu_0 \rightarrow 0$

Or  
 $X' = \frac{X_1 - \mu_0 + X_2 - \mu_0 + \dots + X_n - \mu_0}{n-n_0}$  Actually this is the estimator of  $X - \mu_0$

h)  $\alpha$  can be viewed as a regularization term,

Because when we have prior knowledge about  $\mu$  to be close to 0, we will set  $\alpha$  really big, to drag the estimate back to 0.

Similarly, in Ridge Regression,  $\lambda$  is also used to bound the parameter, if you believe it really small, a large  $\lambda$  should be set to constrain the norm of parameter.

$$\hat{A}\hat{x} - y^* \quad \text{---} \quad \hat{y} \rightarrow \hat{A}\hat{x}$$

$\text{HOP}_k - \psi$

a)

$$\|\hat{A}\hat{x} - y^*\|_2^2 = \|\hat{A}\hat{x} - \hat{Y}\bar{w}\|_2^2 = \|\hat{Y} - \hat{Y}\bar{w}\|_2^2 = \|\vec{z} - \vec{w}\|_2^2$$

$\rightarrow$  column space of  $A$

$$\text{Proj}_A(Y) = \hat{Y} - \vec{z} \quad \vec{z} \perp C(A)$$

$$\text{Proj}_A(y^* + \bar{w}) = \text{Proj}_A(\bar{w}) + y^* = \hat{Y} - \vec{z}$$

$$\text{Proj}_A(\bar{w}) = \bar{w} - \vec{z} \quad \|\text{Proj}_A(\bar{w})\|_2 = \|\bar{w} - \vec{z}\|_2$$

$$\|A \cdot (ATA)^{-1} A^T \bar{w}\|_2 = \|A\hat{x} - y^*\|_2.$$

$$b) \quad \|\hat{A}\hat{x} - y^*\|_2^2 = \|A(ATA)^{-1}A^T w\|_2$$

$$= \|\bar{U}\Sigma V^T (V\bar{\Sigma}^2 V^T)^{-1} V\Sigma U^T w\|_2$$

$$= \|\bar{U}\sum_{i=1}^n \bar{U}^T V\bar{\Sigma}^2 V^T V\Sigma U^T w\|_2$$

$$= \|\bar{U}\sum_{i=1}^n \bar{U}^T \underbrace{\bar{U}^T V\bar{\Sigma}^2 V^T}_{\text{nxd}} \underbrace{V\Sigma U^T w\|_2}_{\text{dxd}} = \|\bar{U}^T w\|_2$$

$$c) \quad \frac{1}{n} E[\|\hat{A}\hat{x} - y^*\|_2^2] = \frac{1}{n} E[\|\bar{U}^T w\|_2^2]$$

$$= \frac{1}{n} E\left[\left\|\begin{bmatrix} u_1^T \\ u_2^T \\ \vdots \\ u_d^T \end{bmatrix} w\right\|_2^2\right] = \frac{1}{n} E\left[(u_1^T w)^2 + (u_2^T w)^2 + \dots + (u_d^T w)^2\right]$$

$$\|u_i\|_2^2 = 1$$

$$= \frac{\sigma^2}{n} E\left[\left(\frac{u_1^T w}{\sigma}\right)^2 + \left(\frac{u_2^T w}{\sigma}\right)^2 + \dots + \left(\frac{u_d^T w}{\sigma}\right)^2\right] = \frac{u^T w}{n} \sim N(0, 1)$$

$$= \frac{\sigma^2 d}{n}$$

$$d) \frac{\sigma^2 d+1}{n} \leq \zeta \quad n \geq \frac{d+1}{\zeta} \sigma^2$$

when Model Complexity increase,  $d \uparrow$ ,  $n \uparrow$

Thus  $n$  is proportional to model complexity

e) As  $D$  becomes large average error is bigger.

As  $n$  becomes large average error becomes smaller.

103-5

a) LinAlg Err: Singular Matrix

Because  $X^T X$  is  $2700 \times 2700$ , samples are only 91.  
 $X^T X$  is singular.

b)  $\lambda = 0.1 \quad 1.2567 \times 10^{-15}$   
1.0  $1.2567 \times 10^{-13}$   
10  $1.2566 \times 10^{-11}$   
100  $1.2556 \times 10^{-9}$   
1000  $1.2460 \times 10^{-7}$

c)  $\lambda = 0.1 \quad 3.2557 \times 10^{-7}$   
1.0  $2.9105 \times 10^{-5}$   
10  $1.5904 \times 10^{-3}$   
100  $3.4773 \times 10^{-2}$   
1000  $2.5440 \times 10^{-1}$

d) (Test data)  $\lambda = 0.1 \quad 0.86808$   
1.0  $0.86210$   
10  $0.82751$   
100  $0.72465$   
1000  $0.72501$

$N$  goes large. Model Complexity goes down, bias goes up ↑,  
variance goes down ↓

e) without  $N$   $2778523071065691$

$N=100 \quad 197781$

## H03 - 6 Own Question

Via training error & true error in bias / variance decomposition

training:  $E_{\text{train}}[f(x_i; D) - y_i]^2$

$$Y = f(x) + N$$

$$= E \left[ \left( h(x_i; D) - E[Y] \right)^2 \right] + V[h(x_i; D)] + V[N]$$

$\downarrow$   
Bias Variance Irreducible Error

True Error:

$$E \left[ \left( h(x_i; D) - \hat{y}_i \right)^2 \right]$$

$A\hat{x}$

$$= E \left[ \left( h(x_i; D) - y^* \right)^2 \right] + V[h(x_i; D) - y^*] + V[y^*]$$

$\downarrow$   
Variance Bias

Thus Training Error and True Error in difference of Irreducible error  
(over-fitting)

As Model become Complex, models are trained to fit Irreducible error

That's why training error becomes small but true error increase (variance)

Error  $\rightarrow$  True Error

$\nearrow$  Training error

$\rightarrow$  Model Complexity