

This homework is due **Friday, October 20 at 10pm.**

1 Getting Started

You may typeset your homework in latex or submit neatly handwritten and scanned solutions. Please make sure to start each question on a new page, as grading (with Gradescope) is much easier that way! Deliverables:

1. Submit a PDF of your writeup to assignment on Gradescope, “HW[n] Write-Up”
2. Submit all code needed to reproduce your results, “HW[n] Code”.
3. Submit your test set evaluation results, “HW[n] Test Set”.

After you've submitted your homework, be sure to watch out for the self-grade form.

- (a) Before you start your homework, write down your team. Who else did you work with on this homework? List names and email addresses. In case of course events, just describe the group. How did you work on this homework? Any comments about the homework?

None. I worked alone

Comments : Need the lecture notes for homework

- (b) Please copy the following statement and sign next to it:

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.

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.

Hanl

Problem # 2

$$a.) P(x|w_1) = \frac{1}{\sqrt{2\pi} 6} e^{-\frac{(x-\mu_1)^2}{26^2}} \quad P(x|w_2) = \frac{1}{\sqrt{2\pi} 6} e^{-\frac{(x-\mu_2)^2}{26^2}}$$

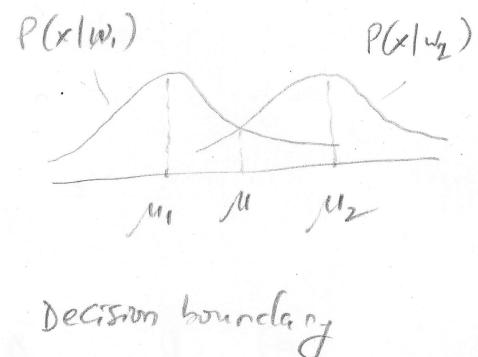
Decision boundary :

$$\frac{1}{\sqrt{2\pi} 6} e^{-\frac{(x-\mu_1)^2}{26^2}} \stackrel{w_1}{=} \frac{1}{\sqrt{2\pi} 6} e^{-\frac{(x-\mu_2)^2}{26^2}} \stackrel{w_2}{=}$$

$$\Rightarrow (x - \mu_1)^2 = (x - \mu_2)^2$$

$$\Rightarrow x - \mu_1 = \mu_2 - x$$

$$\Rightarrow x = \underbrace{\frac{\mu_1 + \mu_2}{2}}_{\bar{\mu}}$$



Decision rule : $x < \frac{\mu_1 + \mu_2}{2} \rightarrow x \in \text{category 1}$

$x > \frac{\mu_1 + \mu_2}{2} \rightarrow x \in \text{category 2}$.

$x = \underbrace{\frac{\mu_1 + \mu_2}{2}}_{\bar{\mu}} \rightarrow x \text{ can be of category 1 or 2}$

$$b.) P_e = P((\text{misclassified as } w_1) | w_2) P(w_2) + P((\text{misclassified as } w_2) | w_1) P(w_1)$$

$$= \frac{1}{2} \int_{\mu}^{+\infty} \frac{1}{\sqrt{2\pi} 6} e^{-\frac{(x-\mu_1)^2}{26^2}} dx + \frac{1}{2} \int_{-\infty}^{\mu} \frac{1}{\sqrt{2\pi} 6} e^{-\frac{(x-\mu_2)^2}{26^2}} dx$$

A

B

$$\text{for A : let } z = \frac{x - \mu_1}{6} \Rightarrow dz = \frac{1}{6} dx$$

$$x = \mu \Rightarrow z = \frac{\mu - \mu_1}{6} = \frac{\mu_2 - \mu_1}{26} = a$$

$$\Rightarrow A = \frac{1}{2} \int_a^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz$$

for B: let $z = \frac{x-\mu_2}{\sigma} \Rightarrow dz = \frac{1}{\sigma} dx$

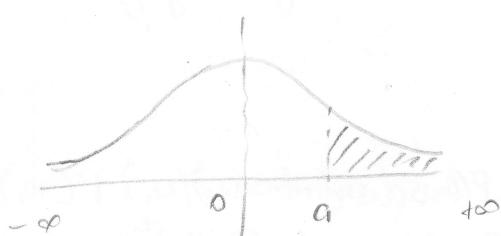
$$x = \mu_1 \Rightarrow z = \frac{\mu_1 - \mu_2}{\sigma} = \frac{\mu_1 - \mu_2}{\sigma} = -a$$

$$\Rightarrow B = \frac{1}{2} \int_{-\infty}^{-a} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz$$

$$= \frac{1}{2} \int_a^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz \quad (\text{since } e^{-\frac{z^2}{2}} \text{ is even function})$$

$$\Rightarrow P_e = A + B = \int_a^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz$$

c.) The integrand of P_e is $\frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$ of the standard normal distribution form



P_e is the area under the curve (shaded area)

if $a \rightarrow \infty \Rightarrow P_e \rightarrow 0$

Problem # 3

- a) Gaussian prior w/ smaller variance
- b) TLS allows errors in X & y . OLS only allows errors in y .
- c) Convex functions:

$$f_1(x) = \max\{-x, 0.1x\}$$

$$f_2(x) = x + \frac{x^2}{10}$$

$$f_3(x) = \frac{e^x + e^{-x}}{2} - 1$$

d)

$$\begin{bmatrix} 1 & -0.1 \\ -0.5 & 1 \end{bmatrix}$$

e.) Training error Bias

f.) Variance

g.) Validation error

Problem #4

a.) $\det(\Sigma) \geq 0$

$$4 - a^2 \geq 0$$

$$-2 \leq a \leq 2$$

b)

2

$$\begin{cases} 1 & \text{for } (x, y, z) \\ 1 & \text{for } (v_x, v_y, v_z) \end{cases}$$

Problem #5

$$P(X=x) = \frac{\lambda^x e^{-\lambda}}{x!}$$

(a) $P(x_1, x_2, \dots, x_n | \lambda) = \prod_{i=1}^n P(x_i | \lambda)$ (x_i 's are i.i.d.)

$$\begin{aligned} & \Rightarrow \ln P(x_1, x_2, \dots, x_n | \lambda) = \sum_{i=1}^n \ln P(x_i | \lambda) \\ &= \sum_{i=1}^n \ln \frac{\lambda^{x_i} e^{-\lambda}}{x_i!} = \sum_{i=1}^n (\ln \lambda^{x_i} + \ln e^{-\lambda} - \ln x_i!) \\ &= \sum_{i=1}^n (x_i \ln \lambda - \lambda - \ln x_i!) \\ &= \ln \lambda \sum_{i=1}^n x_i - n\lambda - \sum_{i=1}^n \ln x_i! \end{aligned}$$

(b) Consider $f(\lambda) = -\ln P = n\lambda - \ln \lambda \sum_{i=1}^n x_i - \sum_{i=1}^n \ln x_i!$

$$f'(\lambda) = n - \frac{1}{\lambda} \sum_{i=1}^n x_i$$

$$f''(\lambda) = \frac{1}{\lambda^2} \sum_{i=1}^n x_i > 0 \quad \forall \lambda \neq 0$$

$\Rightarrow f(\lambda)$ is convex

$\ln P$ max iff $f(\lambda)$ min

$$\Rightarrow f'(\lambda) = 0 \Rightarrow \lambda = \frac{1}{n} \sum_{i=1}^n x_i$$

(c)

$$P(\lambda | x_1, \dots, x_n) = \frac{P(x_1, \dots, x_n | \lambda) P(\lambda)}{P(x_1, \dots, x_n)}$$

$$= \frac{\prod_{i=1}^n P(x_i | \lambda) P(\lambda)}{\prod_{i=1}^n P(x_i)} \leftarrow \text{constant}$$

$$\Rightarrow \ln P(\lambda | x_1, \dots, x_n) = \ln \prod_{i=1}^n P(x_i | \lambda) + \ln P(\lambda) - \ln \prod_{i=1}^n P(x_i)$$

$$= \ln \lambda \sum_{i=1}^n x_i - n\lambda - \sum_{i=1}^n x_i! + \underbrace{\ln \lambda e^{-\lambda}}_{\ln \lambda - \lambda}$$

$\ln P(\lambda | x_1, \dots, x_n)$ max iff $g(\lambda) = -\ln P(\lambda | x_1, \dots, x_n)$ min

$$g(\lambda) = (n+\alpha)\lambda - \ln \lambda \sum_{i=1}^n x_i + \sum_{i=1}^n x_i! - \ln \alpha$$

$$g'(\lambda) = (n+\alpha) - \frac{1}{\lambda} \sum_{i=1}^n x_i = 0$$

$$\Rightarrow \lambda = \frac{1}{n+\alpha} \sum_{i=1}^n x_i$$

$\lambda_{MAP} < \lambda_{MLE}$ since α appears in the denominator

$$n \rightarrow \infty : \lambda = \frac{\sum_{i=1}^n x_i / n}{1 + \frac{\alpha}{n}} = \mu \quad \text{where } \mu = \frac{1}{n} \sum_{i=1}^n x_i$$

$\lambda \rightarrow$ the mean of samples

thus $\lambda_{MAP} \rightarrow \lambda_{MLE}$ when $n \rightarrow \infty$

Problem # 6:

$$(a) \text{ OLS w/ } X' \text{ & } y': \frac{1}{2} \|X'w' - y'\|_2^2 \leftarrow \min$$

$$\text{the close form solution: } w' = (X'^T X)^{-1} X'^T y'$$

Transform OLS w/ $X' \text{ & } y'$ into Ridge Regression w/ $X \text{ & } y$:

$$\begin{aligned} \frac{1}{2} \|X'w' - y'\|_2^2 &= \frac{1}{2} \|(X_e + C)w' - y'\|_2^2 & X = \begin{bmatrix} X \\ Ce^T \\ Ce_d^T \end{bmatrix} = \begin{bmatrix} X \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ Ce^T \\ Ce_d^T \end{bmatrix} \\ &= \underbrace{\frac{1}{2} \|X_e w' - y'\|_2^2}_{\frac{1}{2} \|Xw' - y\|_2^2} + \underbrace{\frac{1}{2} \times 2 \cdot (X_e w' - y')^T C w'}_{= [(Xw' - y)^T 0 \dots 0]} \begin{bmatrix} 0 \\ \vdots \\ Ce^T w' \\ Ce_d^T w' \end{bmatrix} + \underbrace{\frac{1}{2} \|Cw'\|_2^2}_{\frac{c^2}{2} \|w'\|_2^2} & X_e \quad C \\ &= \frac{1}{2} \|Xw' - y\|_2^2 \end{aligned}$$

since the last d rows
are 0

$$= \frac{1}{2} \|Xw' - y\|_2^2 + \frac{c^2}{2} \|w'\|_2^2 \leftarrow \text{ridge regression}$$

Now transform w' into w :

$$\begin{aligned} w' &= ([X^T Ce_1 \dots Ce_d] \begin{bmatrix} X \\ Ce_1^T \\ \vdots \\ Ce_d^T \end{bmatrix})^{-1} [X^T Ce_1 \dots Ce_d^T] \begin{bmatrix} y \\ 0 \\ \vdots \\ 0 \end{bmatrix} \\ &\quad \underbrace{(X_e^T + C^T)(X_e + C)}_{X^T X} \quad \underbrace{y}_{X^T y} \\ &= X_e^T X_e + X_e^T C + C^T X_e + C^T C \end{aligned}$$

we know that:

$$X_e = \begin{bmatrix} X \\ 0^T \\ 0^T \\ \vdots \\ 0^T \end{bmatrix} \quad C = \begin{bmatrix} 0 \\ Ce_1^T \\ \vdots \\ Ce_d^T \end{bmatrix} \quad (X_e + C = X')$$

$$\Rightarrow X_e^T X_e = X^T X \quad X_e^T C = C^T X_e = 0 \quad C^T C = c^2 I_d$$

$$\Rightarrow \hat{w}' = (X^T X + c^2 I_d)^{-1} X^T y$$

$$\text{compared to } \hat{w} = (X^T X + \lambda I_d)^{-1} X^T y$$

$$\Rightarrow c^2 = \lambda \text{ or } c = \sqrt{\lambda}$$

thus, we can convert Ridge regression into OLS by extending X and y as in the question because the ridge regression problem is equivalent to the extended OLS and the ridge regression solution is equivalent to the OLS solution with $c = \sqrt{\lambda}$ as we prove above.

(b) Ridge regression \rightarrow OLS

$$\min_{w'} f(w') = \min_{w'} \frac{1}{2} \|X' w' - y'\|_2^2$$

The step size of gradient descent for the above OLS will be

$$\alpha' = \frac{2}{\lambda_{\min}(X'^T X') + \lambda_{\max}(X'^T X')}$$

$$\begin{aligned} \text{we know from (a) that: } X'^T X' &= X^T X + c^2 I \\ &= X^T X + \lambda I \end{aligned}$$

If λ_i is the eigen value of $X^T X$, then $\lambda_i + \lambda$ is the eigen value of $X^T X + \lambda I = X'^T X'$ since

$$X^T X \vec{v} = \lambda_i \vec{v} + \vec{v}$$

$$\text{and } \lambda I \vec{v} = \lambda \vec{v} + \vec{v}$$

$$\Rightarrow (X^T X + \lambda I) \vec{v} = (\lambda_i + \lambda) \vec{v} + \vec{v}$$

$$\text{Thus } \lambda_{\max}(X^T X') = \lambda_{\max}(X^T X) + \lambda \\ = m + \lambda$$

$$\lambda_{\min}(X^T X') = \lambda_{\min}(X^T X) + \lambda \\ = m + \lambda$$

\Rightarrow The step size for the extended OLS (ridge regression)

$$\beta' = \frac{1}{\lambda_{\max}(X^T X') + \lambda_{\min}(X^T X')} = \frac{1}{m+m+2\lambda}$$

Problem #7

$$Y = M^* + N$$

$$N_{ij} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1) \Rightarrow \Sigma_n = I$$

$$\Rightarrow P(M) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp\left(-\frac{1}{2}(M - 0)^T \Sigma^{-1} (M - 0)\right)$$

$$\Rightarrow P(Y|M) = \frac{1}{\sqrt{(2\pi)^d}} \exp\left(-\frac{1}{2}(Y - M)^T (Y - M)\right)$$

$$P(Y|M) = \prod_{i,j}^d P(y_{ij} | m_i) \quad (\text{since } N_{ij} \text{ i.i.d. so } y_{ij} \text{ i.i.d.})$$

$$\Rightarrow \ln P(Y|M) = \sum_{i,j}^d \ln \frac{1}{\sqrt{(2\pi)^d}} \exp\left(-\frac{1}{2}(y_{ij} - m_{ij})^2\right)$$

$$= \sum_{i,j}^d \left(-\frac{1}{2}(y_{ij} - m_{ij})^2 - \ln \sqrt{(2\pi)^d} \right) \quad \nwarrow \text{const}$$

$$= -\frac{1}{2} \|Y - M\|_F^2 + \text{const.}$$

$$\Rightarrow \max \ln P(Y|M) \text{ is equivalent to } \min \|Y - M\|_F^2$$

i.e.

$$\hat{M} = \underset{M \in \mathbb{R}^{d \times d}}{\operatorname{argmin}} \|Y - M\|_F^2$$

(b) Applying the Eckart-Young theorem directly

$$\| Y - M_k \|_F \leq \| Y - M \|_F$$

thus

$$\hat{M} = \underset{\text{rank}(M)=d-1}{\arg \min} \| Y - M \|_F = M_{d-1}$$

where $M_{d-1} = \sum_{i=1}^{d-1} \sigma_i \vec{u}_i \vec{v}_i^T$ σ_i is the $(d-1)$ largest eigenvalues of $Y^T Y$

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

i.e.

$$\Sigma_{d-1} = \begin{bmatrix} \sigma_1 & & \\ & \sigma_2 & \\ & & \ddots \\ & & & \sigma_{d-1} \\ & & & 0 \end{bmatrix}$$

that is we decompose Y into $U \Sigma V^T$. Then take $(d-1)$ largest singular value of Y and corresponding \vec{u}_i, \vec{v}_i^T to build $M_k = \sum_{i=1}^{d-1} \sigma_i \vec{u}_i \vec{v}_i^T$

$$c) Y = U \Sigma V^{-1}$$

$M_k = \sum_{i=1}^k \sigma_i \vec{u}_i \vec{v}_i^\top$ where σ_i is the k largest singular values of Y and \vec{v}_i, \vec{u}_i are the corresponding singular vectors.

Problem #8

(a)

$$C = \Sigma_{xx}^{-1/2} \Sigma_{x\alpha} \Sigma_{\alpha\alpha}^{-1/2}$$

$$\hat{x} = U^T \Sigma_{xx}^{-1/2} x \quad \hat{\alpha} = V^T \Sigma_{\alpha\alpha}^{-1/2} \alpha$$

$$E[\hat{x}\hat{\alpha}^T] = E[U^T \underbrace{\Sigma_{xx}^{-1/2} x}_{C = U\Lambda U^T} \underbrace{\alpha^T \Sigma_{\alpha\alpha}^{-1/2} V}_{}]$$

$$= E[U^T U \Lambda V^T V] = E[\Lambda] = 1$$

$$E[\hat{x}\hat{x}^T] = E[U^T \Sigma_{xx}^{-1/2} x x^T \Sigma_{xx}^{-1/2} U]$$

$$= U^T \underbrace{\Sigma_{xx}^{-1/2}}_I E[x x^T] \underbrace{\Sigma_{xx}^{-1/2} U}_I$$

$$= U^T \underbrace{\Sigma_{xx}^{-1}}_I U = U^T I U = I$$

$$E[\hat{\alpha}\hat{\alpha}^T] = E[V^T \Sigma_{\alpha\alpha}^{-1/2} \alpha \alpha^T \Sigma_{\alpha\alpha}^{-1/2} V]$$

$$= V^T \Sigma_{\alpha\alpha}^{-1/2} \underbrace{E[\alpha\alpha^T]}_I \Sigma_{\alpha\alpha}^{-1/2} V$$

$$= V^T \Sigma_{\alpha\alpha}^{-1/2} \Sigma_{\alpha\alpha}^{-1/2} V = V^T I V = I$$

$$\begin{aligned}
 b) E[(y - w^T \hat{x})^2] &= E[(y - w^T \hat{x})^T (y - w^T \hat{x})] \\
 &= E[y^T y - y^T w^T \hat{x} - \hat{x}^T w y + \hat{x}^T w w^T \hat{x}] \\
 &= \underbrace{E[y^T y]}_{E[y^2]} - \underbrace{E[y^T w^T \hat{x}]}_{w^T E[y^T \hat{x}]} - \underbrace{E[\hat{x}^T w y]}_{w E[\hat{x}^T y]} + \underbrace{E[\hat{x}^T w w^T \hat{x}]}_{\|w\|_2^2 E[\hat{x}^T \hat{x}]} \\
 &\quad \text{I.}
 \end{aligned}$$

$$E[y^T x] = E[\hat{x}^T y] = E[y \hat{x}]$$

since $y \in \mathbb{R}$

$$\Rightarrow E[(y - w^T \hat{x})^2] = E[y^2] - 2w^T E[y \hat{x}] + \|w\|_2^2.$$

$$\begin{aligned}
 c) \hat{w} &= \underset{\hat{w}}{\operatorname{argmin}} E[(y - w^T \hat{x})^2] + \|w\|_{\text{CCD}}^2 \\
 &= \underset{\hat{w}}{\operatorname{argmin}} E[y^2] - 2w^T E[y \hat{x}] + \|w\|_2^2 + \sum_i^d \frac{1 - \lambda_i}{\lambda_i} (w_i)^2 \\
 &\quad \underbrace{\|w\|_2^2 + \sum_i^d \frac{1}{\lambda_i} (w_i)^2 - \|w\|_2^2}_{f(\vec{w})}
 \end{aligned}$$

$$f'(w_i) = 0 - 2 E[y(x)_i] + 2 w_i \frac{1}{\lambda_i} = 0$$

$$\Rightarrow w_i = \lambda_i E[y(x)_i]$$

(d)

$$\mathbb{E}[\|\tilde{\omega} - \bar{\omega}\|_2^2] = \mathbb{E}[\|\tilde{\omega} - E(\tilde{\omega})\|_2^2]$$

$$= \sigma(\tilde{\omega}) = \sqrt{\gamma}$$

$$\tilde{\omega}_i = \lambda_i \frac{1}{n} \sum_{j=1}^n y^j (\tilde{x}^j)_i$$

$$\sigma(\tilde{\omega}) = \mathbb{E}[\tilde{\omega}^2] + \cancel{\mathbb{E}[\tilde{\omega}]^2}$$

$$= \sum_i^d \lambda_i^2 \frac{1}{n^2} \left(\sum_j^n y^j (\tilde{x}^j)_i \right)^2$$

$$= \sum_i^d \lambda_i^2 \frac{1}{n^2} n Y^2$$

$$= \sum_i^d \lambda_i^2 \frac{1}{n} Y^2 \quad Y^2 \leq 1$$

$$\leq \sum_i^d \frac{\lambda_i^2}{n}$$

$$\Rightarrow \mathbb{E}[\|\tilde{\omega} - \bar{\omega}\|_2^2] \leq \sum_i^d \frac{\lambda_i^2}{n}$$

Problem # 9

Multiple choice

a) Select function convex over \mathbb{R} (3 pts)

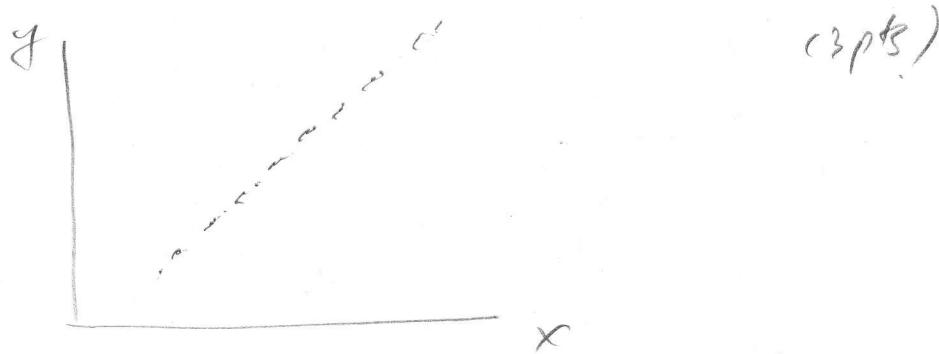
(i) $f(x) = x^3$

(iii) $f(x) = \cosh(x)$

(ii) $f(x) = \tanh(x)$

(iv) $f(x) = \sinh(x)$

b) Give the best covariance matrix of the diagram



c) True or false? (4 pts)

(i) PCA can be used in case of noise with large variance?

(ii) CCA is affected by noise variance?

Solution

a) $f(x) = \sinh(x)$

b) $\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$

c) (i) False (ii) False