## Q1. [Multiple choices question]

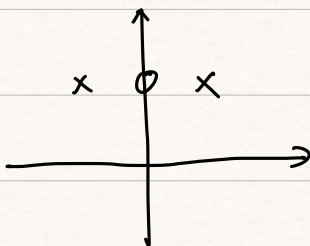(a)    False

(b)    False

<u>Reason:</u>     $\nabla J(w) = 0 \iff \underline{\Phi^T \Phi w = \Phi^T y}$

$\longrightarrow$ when $\Phi^T \Phi$ is not invertible, then we may

have <u>infinitely many</u> $\hat{w}$

(c)    False

Reason:



(d)    False, should use <u>validation set</u> to tune the hyper-params

and use <u>test set</u> to evaluate the model

(e)    True

---

## Q2. [linear basis regression]

$$\Phi := \begin{bmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_N)^T \end{bmatrix} \in \mathbb{R}^{N \times M} \qquad w := \begin{pmatrix} w_0 \\ \vdots \\ w_{M-1} \end{pmatrix} \in \mathbb{R}^M$$

(a)    $\hat{w} = \underset{w \in \mathbb{R}^M}{\text{argmin}} \ R_{emp}(w)$

$$= \underset{w \in \mathbb{R}^M}{\text{argmin}} \ \frac{1}{2N} (\Phi w - y)^T (\Phi w - y)$$

$$= \underset{w \in \mathbb{R}^M}{\text{argmin}} \ \underline{w^T \Phi^T \Phi w - 2 w^T \Phi^T y} := J(w)$$

$\Rightarrow \quad \hat{w} = \underset{w \in \mathbb{R}^M}{\text{argmin}} \ J(w)$, which is an <u>unconstrained optimization</u> problem

① since $\nabla_w J(w) = 2 \Phi^T \Phi w - 2 \Phi^T y$

$\nabla_w^2 J(w) = 2 \Phi^T \Phi \succeq 0$   (PSD)

then $J(w)$ is convex (function) w.r.t $\underline{w}$

② for convex & differentiable function $J(w)$,

$$\hat{w} = \underset{w \in \mathbb{R}^M}{\arg\min} \ J(w) \iff \nabla_w J(\hat{w}) = 0$$

$$\iff 2\Phi^T\Phi\hat{w} - 2\Phi^T y = 0$$

$$\iff \hat{w} = (\Phi^T\Phi)^{-1}\Phi^T y$$

(given that $\underline{\Phi^T\Phi}$ is invertible)

(b) [Weighted Least Square]

$$\hat{w}_{WLS} = \underset{w \in \mathbb{R}^M}{\arg\min} \ R_{emp}(w, a)$$

$$= \underset{w \in \mathbb{R}^M}{\arg\min} \ \sum a_i (w^T\phi(x_i) - y_i)^2 \qquad D(\Phi w - y) \cdot D$$

$$= \underset{w \in \mathbb{R}^M}{\arg\min} \ [D(\Phi w - y)]^T [D(\Phi w - y)] \qquad D = \text{diag}\{\sqrt{a_1}, \ldots, \sqrt{a_N}\}$$

$$= \underset{w \in \mathbb{R}^M}{\arg\min} \ (\Phi w - y)^T W (\Phi w - y) \qquad \underline{W = \text{diag}\{a_1, \ldots, a_N\}}$$

$$= \underset{w \in \mathbb{R}^M}{\arg\min} \ \underline{w^T\Phi^T W\Phi w - 2w^T\Phi^T W y} := J(w)$$

① since $\nabla_w J(w) = 2\Phi^T W \Phi w - 2\Phi^T W y$

$$\nabla_w^2 J(w) = 2\Phi^T W \Phi \succcurlyeq 0 \qquad (\text{since } a_i > 0 \text{ for } i = 1, 2, \ldots, N)$$

then $J(w)$ is convex (function) w.r.t $\underline{w}$

② for convex & differentiable function $J(w)$,

$$\hat{w}_{WLS} = \underset{w \in \mathbb{R}^M}{\arg\min} \ J(w) \iff \nabla_w J(\hat{w}_{WLS}) = 0$$

$$\iff 2\Phi^T W \Phi \hat{w}_{WLS} - 2\Phi^T W y = 0$$

$$\iff \hat{w}_{WLS} = (\Phi^T W \Phi)^{-1} \Phi^T W y$$

(c) [Application Scenario for WLS]

① When we have prior that some data points in dataset $\mathcal{D}$ are outliers, we can assign low weights (small $a_i$) to them

② Linear regression can be viewed as the Gaussian Model that

$$y \sim \text{Gaussian}(w^T \phi(x), \underline{\Sigma}) \quad \text{where } \underline{\Sigma} = 6^2 I. \quad \boxed{+ \text{ MLE}}$$

$\hookrightarrow$ isotropic Gaussian

Weighted Least Square can be viewed as the Gaussian Model that

$$y \sim \text{Gaussian}(w^T \phi(x), \underline{\Sigma}) \quad \text{where } \underline{\Sigma} = \begin{pmatrix} 6_1^2 & & \\ & \ddots & \\ & & 6_N^2 \end{pmatrix}$$

$$\left( a_i = \frac{1}{6_i^2} \right) \longleftarrow \qquad \boxed{+ \text{MLE}}$$

$\longrightarrow$ from this perspective, we can take <u>deviance of each data point</u> into consideration:

$\hookrightarrow$ <span style="color:blue">our <u>confidence</u> of each data</span>

o for those data points we have more confidence (due to the observation error), we can assign <u>small $6_i^2$ (large $a_i$)</u> to them;

o for those data points are more likely to be the noise, we can assign <u>large $6_i^2$ (small $a_i$)</u> to them.

---

Q3. [nearest-neighbour] $\longrightarrow$ <u>$k$NN with $k=1$</u>

a) $i(x) = \underset{j \in [N]}{\text{argmin}} \; \|x - x_j\|_2^2 = \underset{j \in [N]}{\text{argmin}} \; \langle x - x_j, x - x_j \rangle$

$\qquad = \underset{j \in [N]}{\text{argmin}} \; \langle x_j, x_j \rangle - 2\langle x, x_j \rangle$

$\qquad = \underset{j \in [N]}{\text{argmin}} \; \langle x_j - 2x, x_j \rangle$

b) <u>feature map</u> $\phi: x \in \mathbb{R}^d \longmapsto \phi(x) \in \mathbb{R}^M$, denote $\underline{K(x,y) = \langle \phi(x), \phi(y) \rangle}$

$f(x) = y_{i(\phi(x))}$ such that $i(\phi(x)) = \underset{j \in [N]}{\text{argmin}} \; \|\phi(x) - \phi(x_j)\|_2^2 = \underset{j \in [N]}{\text{argmin}} \; K(x_j - 2x, x_j)$

## Q4. [LR]

**(a)** we have $\underline{\mathbb{E}[\varepsilon_i] = 0}$   $\underline{\mathbb{E}[\varepsilon_i^2] = \text{Var}[\varepsilon_i] + \mathbb{E}[\varepsilon_i]^2 = \sigma^2}$

$$\min_{w_0, w_1} \mathbb{E}\left[ \frac{1}{2N} \sum_{i=1}^{N} [w_0 + w_1 x_i - y_i + w_1 \varepsilon_i]^2 \right]$$

$$\Longleftrightarrow \min_{w_0, w_1} \mathbb{E}\left[ \frac{1}{2N} \sum_{i=1}^{N} \left[ (w_0 + w_1 x_i - y_i)^2 + w_1^2 \varepsilon_i^2 + w_1 \varepsilon_i \cdot C_i \right] \right]$$

where $\underline{C_i = 2(w_0 + w_1 x_i - y_i)}$

$$\Longleftrightarrow \min_{w_0, w_1} \frac{1}{2N} \sum_{i=1}^{N} \left[ (w_0 + w_1 x_i - y_i)^2 + w_1^2 \underline{\mathbb{E}[\varepsilon_i^2]} + \underline{\mathbb{E}[\varepsilon_i]} w_1 C_i \right]$$
$$\qquad\qquad\qquad\qquad\qquad\qquad\quad \sigma^2 \qquad\qquad 0$$

$$\Longleftrightarrow \min_{w_0, w_1} \frac{1}{2N} \sum_{i=1}^{N} (w_0 + w_1 x_i - y_i)^2 + \frac{w_1^2}{2} \sigma^2$$

$$\Longleftrightarrow \min_{w_0, w_1} \underbrace{\frac{1}{2} \| Xw - y \|_2^2}_{\text{OLS}} + \underbrace{\frac{N \sigma^2}{2} w_1^2}_{\text{regularization on } w_1}$$

$w = \begin{pmatrix} w_0 \\ w_1 \end{pmatrix}$

$X = \begin{pmatrix} 1 & x_1 \\ & \vdots \\ 1 & x_n \end{pmatrix}$

$$\longrightarrow \min_{w} \frac{1}{2}\left( w^T X^T X w - 2 w^T X^T y \right) + \frac{1}{2} w^T K w \qquad K = \begin{pmatrix} 0 & \\ & N\sigma^2 \end{pmatrix}$$

$$\Longleftrightarrow \min_{w} \underline{\frac{1}{2} w^T (X^T X + K) w - w^T X^T y} := J(w)$$

$$\nabla J(w) = (X^T X + K) w - X^T y$$

$$\nabla^2 J(w) = X^T X + K \geq 0$$

$$\Longrightarrow \hat{w} \in \underset{w}{\arg\min} \, J(w) \Longleftrightarrow \nabla J(\hat{w}) = 0$$

$$\Longleftrightarrow \underline{\hat{w} = (X^T X + K)^{-1} X^T y}$$
$$\uparrow$$
$$\underline{\text{if } (X^T X + K) \text{ is invertible}}$$

**(b)** If $\mathbb{E}[\varepsilon_i] = b \neq 0$, then $\mathbb{E}[\varepsilon_i^2] = \sigma^2 + b^2$

$$\min_{w_0, w_1} \frac{1}{2N} \sum_{i=1}^{N} \left[ (w_0 + w_1 x_i - y_i)^2 + w_1^2 \mathbb{E}[\varepsilon_i^2] + \mathbb{E}[\varepsilon_i] w_1 C_i \right]$$

where $C_i = 2(w_0 + w_1 x_i - y_i)$

$\Leftrightarrow \min_w \dfrac{1}{2N}(Xw-y)^T(Xw-y) + \dfrac{1}{2}w_1^2(\delta^2+b^2) + \dfrac{1}{N}\cdot bw_1 \sum_{i=1}^{N}(w_0 + w_1 x_i - y_i)$

$\dfrac{1}{2N}(w^T X^T X w - 2w^T X^T y) + \dfrac{1}{2N} w^T K w + \dfrac{2b}{2N}\cdot \left[ w^T T_1 w - w^T T_2^T y \right]$

$K = \begin{pmatrix} 0 & \\ & N(\delta^2+b^2) \end{pmatrix}$

$w^T \begin{pmatrix} 0 \\ 1 \end{pmatrix} \mathbb{1}^T (Xw - y)$

$= w^T \begin{pmatrix} 0 \cdots 0 \\ 1 \cdots 1 \end{pmatrix} Xw - w^T \begin{pmatrix} 0 \cdots 0 \\ 1 \cdots 1 \end{pmatrix} y$

$\boxed{\mathbb{R}^{2 \times N}}$

$\Leftrightarrow \min_w \dfrac{1}{2} w^T X^T X w + \dfrac{1}{2} w^T K w + \dfrac{2b}{2} w^T T_1 w$

$\qquad - w^T X^T y - \dfrac{2b}{2} w^T T_2^T y$

$\qquad := J(w)$

$T_1 = \begin{pmatrix} 0 & 0 \\ N & \Sigma x_i \end{pmatrix} \in \mathbb{R}^{2 \times 2}$

$T_2 = \begin{pmatrix} 0 & 1 \\ \vdots & \vdots \\ 0 & 1 \end{pmatrix} \in \mathbb{R}^{n \times 2}$

$\left[ \begin{array}{l} \nabla J(w) = (X^T X + K + 2b T_1) w - (X + b T_2)^T y \\ \nabla^2 J(w) = X^T X + K + 2b T_1 \end{array} \right.$

Suppose $\nabla^2 J(w) \geq 0 \Rightarrow J(w)$ is convex on $w$

$\Rightarrow \hat{w} = \arg\min J(w) \Leftrightarrow \nabla J(\hat{w}) = 0$

$\qquad\qquad \Leftrightarrow \hat{w} = (X^T X + K + 2b T_1)^{-1} (X + b T_2)^T y$

$\boxed{\text{Conclusion}}$ : <u>Minimization Problem changed.</u> But we can still achieve

the closed-form solution <u>under some assumption</u>.

## Q5. [RBF kernel]

$$k(x,y) = \exp\left(-\dfrac{1}{2s^2}\|x-y\|_2^2\right) \qquad s > 0$$

Prove that $k(\cdot,\cdot)$ is a valid kernel function.

$\boxed{\text{Pf}}$ : $\|x-y\|_2^2 = \|x\|_2^2 - 2x^T y + \|y\|_2^2$

$\Rightarrow k(x,y) = \exp\left(-\dfrac{1}{2s^2}\|x\|_2^2\right) \underline{\exp\left(\dfrac{1}{s^2}x^T y\right)} \exp\left(-\dfrac{1}{2s^2}\|y\|_2^2\right)$

From Taylor Expansion,

$$\Rightarrow \exp\left(\frac{1}{s^2} x^T y\right) = \sum_{n=0}^{\infty} \frac{1}{n!} \left(\frac{1}{s^2} x^T y\right)^n$$

Then, we can show that, RBF kernel $K(\cdot, \cdot)$ is valid.

Reason can be shown as follows :

① $K(x, y) = x^T y \longrightarrow$ linear kernel, which is valid kernel

② use scaling property $\Rightarrow K(x, y) = \frac{1}{s^2} x^T y$ is valid kernel

③ use product property $\Rightarrow K_n(x, y) = \left(\frac{1}{s^2} x^T y\right)^n$ is valid kernels

④ use scaling property $\Rightarrow K_n(x, y) = \frac{1}{n!} \left(\frac{1}{s^2} x^T y\right)^n$ is valid kernels

⑤ use addition property $\Rightarrow K_n(x, y) = \sum_{k=0}^{n} \frac{1}{n!} \left(\frac{1}{s^2} x^T y\right)^n$ is valid kernels

⑥ use limit property $\Rightarrow K(x, y) = \lim_{n \to \infty} K_n(x, y)$

$$= \lim_{n \to \infty} \sum_{k=0}^{\infty} \frac{1}{n!} \left(\frac{1}{s^2} x^T y\right)^n$$

$$= \exp\left(\frac{1}{s^2} x^T y\right) \text{ is valid kernel}$$

⑦ use normalization property

$$\Rightarrow K(x, y) = \exp\left(-\frac{1}{2s^2} \|x\|_2^2\right) \exp\left(\frac{1}{s^2} x^T y\right) \exp\left(-\frac{1}{2s^2} \|y\|_2^2\right)$$

$$= \exp\left[-\frac{1}{2s^2} \left(\|x\|_2^2 - 2x^T y + \|y\|_2^2\right)\right]$$

$$= \exp\left(-\frac{1}{2s^2} \|x - y\|_2^2\right) \text{ is valid kernel}$$

That is, RBF kernel is a valid kernel !

#