

1) prove X_i 's covariance is a semi positive definite matrix
 covariance $\Sigma = E[(X - \bar{X})(X - \bar{X})^T] \leftarrow (\bar{X} = \bar{X}_i \text{ is the average of } X_i)$

then for \forall real vector u , we can have

$$\begin{aligned} u^T \Sigma u &= u^T E[(X - \bar{X})(X - \bar{X})^T] u \\ &= E[u^T (X - \bar{X})(X - \bar{X})^T u] \end{aligned}$$

$$\text{we can see } u^T (X - \bar{X}) = G_i = ((X - \bar{X})^T u)^T$$

\therefore we can have

$$u^T \Sigma u = E[G_i G_i^T] = \underbrace{E[G_i^2]}_{\geq 0} = E[\|u^T (X - \bar{X})\|_2^2]$$

so finally we get it.

2) Because $X_i \sim N(\mu, \Sigma)$ which is a mul normal distribution

\therefore we can directly get it PDF

$$f(x) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

then we use log-likelihood to get the $\ell(\mu, \Sigma)$

$$\begin{aligned} \Rightarrow \ell(\mu, \Sigma) &= \sum_{i=1}^N \log f(x_i) = \sum_{i=1}^N \log \left[\frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)\right) \right] \\ &= \sum_{i=1}^N \log \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} - \frac{1}{2} \sum_{i=1}^N (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \end{aligned}$$

then we can use it to get the MLE of μ .

$$\Rightarrow \frac{\partial \ell(\mu, \Sigma)}{\partial \mu} = \frac{\partial \left[\sum_{i=1}^N \log \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} - \frac{1}{2} \sum_{i=1}^N (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right]}{\partial \mu} = -\frac{1}{2} \frac{\partial \sum_{i=1}^N (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)}{\partial \mu}$$

not corresponding

$$\begin{aligned} &= -\frac{1}{2} \cdot (-2) \cdot \sum_{i=1}^N \Sigma^{-1} (x_i - \mu) \\ &= \sum_{i=1}^N \Sigma^{-1} (x_i - \mu) \end{aligned}$$

Let it equal to 0

$$\therefore \sum_{i=1}^N \Sigma^{-1} (x_i - \mu) = 0 \quad \therefore \hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\therefore \sum_{i=1}^N (x_i - \mu) = 0$$

No.

Date. / /

3) $\therefore \hat{\theta}$ is an unbiased estimator of θ &

$$\therefore E[\hat{\theta}] = \theta$$

$$\therefore \text{Var}(X) = E[X^2] - E[X]^2$$

$$\therefore \text{Var}(\hat{\theta}) = E[\hat{\theta}^2] - E[\hat{\theta}]^2$$

$$= E[\hat{\theta}^2] - \theta^2 > 0$$

$$\therefore E[\hat{\theta}^2] > \theta^2$$

$\therefore \hat{\theta}^2$ is not an unbiased estimator of θ^2

王 2021/5/33 13T

multiply
No.

calculate

MLE problem

$$\frac{\partial L(a,b)}{\partial a} \text{ and } \frac{\partial L(a,b)}{\partial b}$$

are too complicated to

2) The likelihood function is $L(a,b) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}b} \exp(-\frac{1}{2b^2}(y_i - ax_i - b)^2)$ ✓

then we can use log-likelihood to get the simplified expression so that we can estimate a, b more easily $\Rightarrow \therefore$ (add)

$$\Rightarrow \arg\max_{a,b} \ell(a,b) \triangleq \arg\max_{a,b} L(a,b)$$

$$\begin{aligned} \text{then } \ell(a,b) &= \log \prod_{i=1}^n \frac{1}{\sqrt{2\pi}b} \exp(-\frac{1}{2b^2}(y_i - ax_i - b)^2) \\ &= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}b} \exp(-\frac{1}{2b^2}(y_i - ax_i - b)^2) \\ &= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}b} - \frac{1}{2b^2} \sum_{i=1}^n (y_i - ax_i - b)^2 \end{aligned}$$

Thus, we can get the correct a, b to maximize $\ell(a,b)$ through $\frac{\partial \ell(a,b)}{\partial a}$ or $\frac{\partial \ell(a,b)}{\partial b}$, and this way's effect is equal to maximize $L(a,b)$ ✓

$$\begin{aligned} 2) \text{ (i) } \frac{\partial \ell(a,b)}{\partial a} &= -\frac{1}{2b^2} \frac{\partial \sum_{i=1}^n (y_i - ax_i - b)^2}{\partial a} = -\frac{1}{2b^2} (-2) \sum_{i=1}^n x_i (y_i - ax_i - b) \\ &= \frac{1}{b^2} \sum_{i=1}^n x_i (y_i - ax_i - b) \end{aligned}$$

Let $\frac{\partial \ell(a,b)}{\partial a}$ be equal to 0

$$\frac{1}{b^2} \sum_{i=1}^n x_i (y_i - ax_i - b) = 0$$

$$\sum_{i=1}^n x_i (y_i - b) = \sum_{i=1}^n ax_i^2$$

$$\frac{\sum_{i=1}^n x_i (y_i - b)}{\sum_{i=1}^n x_i^2} = \hat{a}$$

$$\begin{aligned} \text{ii) } \frac{\partial \ell(a,b)}{\partial b} &= -\frac{1}{2b^2} \frac{\partial \sum_{i=1}^n (y_i - ax_i - b)^2}{\partial b} = -\frac{1}{2b^2} \cdot (-2) \sum_{i=1}^n (y_i - ax_i - b) \\ &= \frac{1}{b^2} \sum_{i=1}^n (y_i - ax_i - b) \end{aligned}$$

Let $\frac{\partial \ell(a,b)}{\partial b}$ be equal to 0

$$\frac{1}{b^2} \sum_{i=1}^n (y_i - ax_i - b) = 0 \quad \Rightarrow \quad \sum_{i=1}^n (y_i - ax_i) = \hat{b}$$

$$\begin{aligned} \Rightarrow \hat{a} &= \frac{\sum_{i=1}^n x_i (y_i - \hat{b})}{\sum_{i=1}^n x_i^2} \quad \hat{b} = \frac{\sum_{i=1}^n (y_i - ax_i)}{n} \\ &\quad \text{Not simplify} \end{aligned}$$

$$\begin{aligned} \Rightarrow \hat{b} &= \frac{1}{n} \sum_{i=1}^n y_i - \hat{a} \frac{1}{n} \sum_{i=1}^n x_i \\ \hat{b} &= \bar{y} - \hat{a} \bar{x} = \bar{y} - \hat{a} \frac{\sum_{i=1}^n x_i}{n} \\ \hat{a} &= \frac{\sum_{i=1}^n x_i (y_i - \bar{y} + \hat{a} \bar{x})}{\sum_{i=1}^n x_i^2} \\ \hat{a} &= \frac{\sum_{i=1}^n (y_i - \bar{y}) x_i}{\sum_{i=1}^n (x_i - \bar{x}) x_i} \quad \text{simplify} \end{aligned}$$

Wengu 文谷

No.

Date. / /

$$3) f(X) = aX + b$$

let we substitute \hat{a}, \hat{b} in this expression

$$f(X) = \hat{a}X + \hat{b}$$

and then substitute \bar{x} in it.

$$f(\bar{x}) = \hat{a}\bar{x} + \hat{b}$$

$$= \frac{\hat{a}}{n} \sum_{i=1}^n x_i + \hat{b}$$

$$= \frac{\hat{a}}{n} \sum_{i=1}^n x_i + \frac{1}{n} \sum_{i=1}^n (y_i - \hat{a}x_i)$$

$$= \frac{\hat{a}}{n} \sum_{i=1}^n x_i + \frac{1}{n} \sum_{i=1}^n y_i - \frac{\hat{a}}{n} \sum_{i=1}^n x_i$$

$$= \frac{1}{n} \sum_{i=1}^n y_i = \bar{y}$$

$$\therefore \boxed{f(\bar{x}) = \bar{y}}$$

Therefore, based on the result in (b): \hat{a}, \hat{b} the ^{learned} linear model $f(X) = \hat{a}X + \hat{b}$ always passes through the point (\bar{x}, \bar{y}) .

3. D "linear" regards that the input features and the outputs have a linear relation.

2) From ^{what} we have learned in class. we define

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \operatorname{PRSS}(\lambda, \beta) = \underset{\beta}{\operatorname{argmin}} \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2$$

$$\therefore \operatorname{PRSS}(\lambda, \beta) = \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2$$

$$= (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$$

$$= y^T y - \beta^T X^T y - y^T X \beta + \beta^T X^T X \beta + \lambda \beta^T \beta$$

to minimise $\operatorname{PRSS}(\lambda, \beta)$ we take the derivative of it

$$\begin{aligned} \frac{\partial \operatorname{PRSS}(\lambda, \beta)}{\partial \beta} &= -X^T y - X^T y + 2X^T X \beta + 2\lambda I_p \beta \\ &= -2X^T y + 2(X^T X + \lambda I_p) \beta \end{aligned}$$

$$\text{Let it be equal to } \frac{\partial \operatorname{PRSS}(\lambda, \beta)}{\partial \beta} = 0$$

~~In this expression, λ is the parameter which we have determined before this step~~

$$\Rightarrow \hat{\beta}^{\text{ridge}} = (X^T X + \lambda I_p)^{-1} X^T y = \beta_*$$

Thus, we simply get this expression. but there is also a problem we need to prove. that is \downarrow should be invertible, given $\lambda > 0$

$$\Rightarrow X^T X + \lambda I_p = (C U D V^T)^T (C U D V^T) + \lambda I_p$$

$$\text{SVD decomposition. } = V D U^T U D V^T + \lambda I_p$$

$$\left\{ \begin{array}{l} X = U D V^T \\ V^T V = U^T U = I_p \end{array} \right. \quad = V D^2 V^T + \lambda V I_p V^T$$

$$= V (D^2 + \lambda I_p) V^T$$

\downarrow
this is semi positive definite

$\lambda > 0$ and I_p is a Identity matrix

$\therefore D^2 + \lambda I_p$ is positive definite

$\therefore (X^T X + \lambda I_p)$ is invertible

\therefore Finally, we get it \checkmark

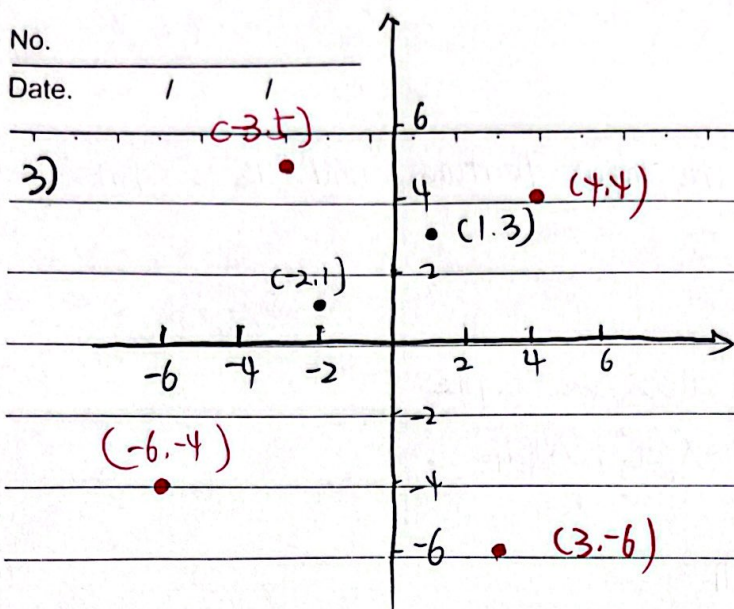
No.

Date.

label

● means -1

● means +1



Therefore, from the graph above, we can know the given data set is not linear separable. We can't use one ^{straight} line to separate these two labels from each other.

Instead, we can use a circle to separate them. which is Multivariate nonlinear ✓