

Introduction to Machine Learning, Fall 2023

Homework 1

(Due Thursday, Oct. 26 at 11:59pm (CST))

October 11, 2023

1. [10 points] [Math review] Suppose $\{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n\}$ form a random sample from a multivariate distribution:

- (a) Prove that the covariance of \mathbf{X}_i is a semi positive definite matrix. [3 points]
- (b) Assuming $\mathbf{X}_i \sim \mathcal{N}(\mu, \Sigma)$ which is a multivariate normal distribution, and samples \mathbf{X}_i , derive the the log-likelihood $l(\mu, \Sigma)$ and MLE of μ . [4 points]
- (c) Suppose $\hat{\theta}$ is an unbiased estimator of θ and $\text{Var}(\hat{\theta}) > 0$. Prove that $(\hat{\theta})^2$ is not an unbiased estimator of θ^2 . [3 points]

(a). to prove the covariance of \mathbf{X}_i is a semi positive definite matrix, we have to show that $\mathbf{a}^T \Sigma \mathbf{a} \geq 0$ where \mathbf{a} can be any nonzero real vector

let Σ_i be the variance of \mathbf{X}_i , so that

$$\Sigma_i = E[(\mathbf{X}_i - \mu)(\mathbf{X}_i - \mu)^T], \text{ where } \mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n) \text{ and } \mu = (E[\mathbf{X}_1], E[\mathbf{X}_2], \dots, E[\mathbf{X}_n])$$

$$\begin{aligned} \text{then } \mathbf{a}^T \Sigma \mathbf{a} &= \mathbf{a}^T E[(\mathbf{X}_i - \mu)(\mathbf{X}_i - \mu)^T] \mathbf{a} \\ &= E[\mathbf{a}^T (\mathbf{X}_i - \mu)(\mathbf{X}_i - \mu)^T \mathbf{a}] \end{aligned}$$

$$\text{let } \mathbf{Y}_i = \mathbf{X}_i - \mu$$

$$\begin{aligned} \text{then } \mathbf{a}^T \Sigma \mathbf{a} &= E[\mathbf{a}^T \mathbf{Y}_i \mathbf{Y}_i^T \mathbf{a}] \\ &= \mathbf{a}^T E[\mathbf{Y}_i \mathbf{Y}_i^T] \mathbf{a} \end{aligned}$$

$$\because E[\mathbf{Y}_i \mathbf{Y}_i^T] \geq 0, \therefore \mathbf{a}^T \Sigma \mathbf{a} \geq 0.$$

$\therefore \Sigma_i$ is a semi positive definite matrix.

b. ① 对数似然函数 $l(\mu, \Sigma) = \sum [\log f(\mathbf{X}_i; \mu, \Sigma)]$ 其中 $f(\mathbf{X}_i; \mu, \Sigma)$ 表示多元正态分布的概率密度函数。

$$f(\mathbf{X}_i; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} e^{-\frac{1}{2} (\mathbf{X}_i - \mu)^T \Sigma^{-1} (\mathbf{X}_i - \mu)}$$

$$f(\mathbf{X}_i; \mu, \Sigma) = (2\pi)^{-\frac{d}{2}} |\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2} (\mathbf{X}_i - \mu)^T \Sigma^{-1} (\mathbf{X}_i - \mu)}$$

$$\log f(\mathbf{X}_i; \mu, \Sigma) = -\frac{1}{2} d \log(2\pi) - \frac{1}{2} \log|\Sigma| - \frac{1}{2} (\mathbf{X}_i - \mu)^T \Sigma^{-1} (\mathbf{X}_i - \mu)$$

$$\therefore l(\mu, \Sigma) = \sum_{i=1}^n \left(-\frac{1}{2} d \log(2\pi) - \frac{1}{2} \log|\Sigma| - \frac{1}{2} (\mathbf{X}_i - \mu)^T \Sigma^{-1} (\mathbf{X}_i - \mu) \right).$$

② To find MLE of μ . $\frac{\partial \ell(\mu, \Sigma)}{\partial \mu} = 0$.

$$\frac{\partial}{\partial \mu} \sum_{i=1}^n \left(-\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right) = 0.$$

$$\frac{\partial}{\partial \mu} \sum_{i=1}^n \left[-\frac{1}{2} (x_i^T \Sigma^{-1} x_i - x_i^T \Sigma^{-1} \mu - \mu^T \Sigma^{-1} x_i + \mu^T \Sigma^{-1} \mu) \right] = 0.$$

$$\therefore \frac{\partial (A^T \mu)}{\partial \mu} = A.$$

$$\frac{\partial (x^T A x)}{\partial x} = (A^T + A)x.$$

$$\therefore \text{Set } = \sum_{i=1}^n \left(-\frac{1}{2} (-x_i \Sigma^{-1} - x_i \Sigma^{-1} + 2 \Sigma^{-1} \mu) \right).$$

$$= \sum_{i=1}^n (x_i \Sigma^{-1} - \Sigma^{-1} \mu).$$

$$\because (\Sigma^{-1})^T = (\Sigma^{-1})$$

$$= \sum_{i=1}^n \Sigma^{-1} (x_i - \mu) = 0.$$

$$\therefore \sum_{i=1}^n (x_i) = \sum_{i=1}^n \mu = n \mu.$$

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

(c). suppose $\hat{\theta}$ is an unbiased estimator of θ , then $E(\hat{\theta}) = \theta$.

we need to prove that $(\hat{\theta})^2$ is not an unbiased estimator of θ , then.

we need to prove that $E[(\hat{\theta})^2] \neq \theta^2$

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

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

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

$$\therefore \text{Var}(\hat{\theta}) > 0. \text{ \& } E(\hat{\theta}) = \theta$$

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

$$\therefore E[(\hat{\theta})^2] \neq \theta^2.$$

2. [10 points] Consider real-valued variables X and Y , in which Y is generated conditional on X according to

$$Y = aX + b + \epsilon, \text{ where } \epsilon \sim \mathcal{N}(0, \sigma^2).$$

Here ϵ is an independent variable, called a noise term, which is drawn from a Gaussian distribution with mean 0, and variance σ^2 . This is a single variable linear regression model, where a is the only weight parameter and b denotes the intercept. The conditional probability of Y has a distribution $p(Y|X, a, b) \sim \mathcal{N}(aX + b, \sigma^2)$, so it can be written as:

$$p(Y|X, a, b) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(Y - aX - b)^2\right).$$

- (a) Assume we have a training dataset of n i.i.d. pairs (x_i, y_i) , $i = 1, 2, \dots, n$, and the likelihood function is defined by $L(a, b) = \prod_{i=1}^n p(y_i|x_i, a, b)$. Please write the Maximum Likelihood Estimation (MLE) problem for estimating a and b . [3 points]
- (b) Estimate the optimal solution of a and b by solving the MLE problem in (a). [4 points]
- (c) Based on the result in (b), argue that the learned linear model $f(X) = aX + b$, always passes through the point (\bar{x}, \bar{y}) , where $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ and $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ denote the sample means. [3 points]

$$(a) \therefore p(y_i|x_i, a, b) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(y_i - ax_i - b)^2} \text{ and }.$$

$$L(a, b) = \prod_{i=1}^n p(y_i|x_i, a, b).$$

$$\therefore L(a, b) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2\sigma^2}(y_i - ax_i - b)^2\right]$$

$$\therefore L(a, b) = \frac{1}{(2\pi)^{\frac{n}{2}} \sigma^n} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - ax_i - b)^2\right]$$

$$\therefore \ell(a, b) = \log(L(a, b)) = \sum_{i=1}^n \left[-\frac{1}{2} \log 2\pi - \log \sigma - \frac{1}{2\sigma^2} (y_i - ax_i - b)^2 \right]$$

(b) To find the MLE problem for estimating a and b , we need to get the log-likelihood and $\frac{\partial l(a,b)}{\partial a} = 0$ and $\frac{\partial l(a,b)}{\partial b} = 0$.

$$l(a,b) = \log(L(a,b)) = \sum_{i=1}^n \left[-\frac{1}{2} \log 2\pi - \log \sigma - \frac{1}{2\sigma^2} (y_i - ax_i - b)^2 \right]$$

$$\frac{\partial l(a,b)}{\partial a} = \sum_{i=1}^n -\frac{1}{\sigma^2} \cdot 2 (y_i - ax_i - b) \cdot (-x_i)$$

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

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

$$\sum_{i=1}^n ax_i^2 = \sum_{i=1}^n y_i x_i - \sum_{i=1}^n bx_i \quad \textcircled{1}$$

$$\frac{\partial l(a,b)}{\partial b} = \sum_{i=1}^n -\frac{1}{\sigma^2} \cdot 2 (y_i - ax_i - b) \cdot (-1).$$

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

$$\therefore \sum_{i=1}^n b = \sum_{i=1}^n y_i - \sum_{i=1}^n ax_i.$$

$$nb = \sum_{i=1}^n (y_i - ax_i) \quad \textcircled{2}$$

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

Now from $\textcircled{1}$ $\textcircled{2}$ we have.

$$a \left(\sum_{i=1}^n x_i^2 \right) = \sum_{i=1}^n y_i x_i - \sum_{i=1}^n bx_i$$

$$a \left(\sum_{i=1}^n x_i^2 \right) = \sum_{i=1}^n y_i - nb.$$

$$a \sum_{i=1}^n x_i^2 = \sum_{i=1}^n y_i x_i - \left(\sum_{i=1}^n x_i \right) \frac{1}{n} \left(\sum_{i=1}^n y_i - \sum_{i=1}^n ax_i \right).$$

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

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

$$(c). \quad f(X) = \hat{a}X + \hat{b} \quad \text{A} \quad \hat{a} = \frac{\sum_{i=1}^n y_i x_i - \hat{b} \sum_{i=1}^n x_i}{\sum_{i=1}^n x_i^2} \quad \text{A}.$$

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

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

$\therefore f(x)$ always passes (\bar{x}, \bar{y}) .

3. [10 points] [Regression and Classification]

- (a) When we talk about linear regression, what does 'linear' regard to? [2 points]
- (b) Assume that there are n given training examples $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, where each input data point x_i has m real valued features. When $m > n$, the linear regression model is equivalent to solving an under-determined system of linear equations $\mathbf{y} = \mathbf{X}\beta$. One popular way to estimate β is to consider the so-called ridge regression:

$$\underset{\beta}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2$$

for some $\lambda > 0$. This is also known as Tikhonov regularization.

Show that the optimal solution β_* to the above optimization problem is given by

$$\beta_* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

Hint: You need to prove that given $\lambda > 0$, $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$ is invertible. [5 points]

- (c) Is the given data set linear separable? If yes, construct a linear hypothesis function to separate the given data set. If no, explain the reason. [3 points]

Data	(1,3)	(4,4)	(3,-6)	(-2,1)	(-3,5)	(-6,-4)
Label	+1	-1	-1	+1	-1	-1

(a). linear refers to the relationship between the independent variables (predictors) and the dependent variable (response). It means that the model assumes a linear relationship between the independent variables and the expected value of the dependent variable.

(b). $L(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2$

$$\begin{aligned} \frac{dL(\beta)}{d\beta} &= \frac{d[(\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda \beta^T \beta]}{d\beta} \\ &= \frac{d[\mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\beta - (\mathbf{X}\beta)^T \mathbf{y} + \beta^T \mathbf{X}^T \mathbf{X} \beta + \lambda \beta^T \beta]}{d\beta} \\ &= -\mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \beta + 2\lambda \beta. \\ &= 2\mathbf{X}^T (\mathbf{X}\beta - \mathbf{y}) + 2\lambda \beta = 0. \\ \therefore \quad \mathbf{X}^T \mathbf{X} \beta - \mathbf{X}^T \mathbf{y} + \lambda \beta &= 0. \\ (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})\beta &= \mathbf{X}^T \mathbf{y} \end{aligned}$$

if $\lambda > 0$, then $X^T X + \lambda I$ is invertible. —

we can prove it by contradiction: we suppose that it is ~~invertible~~ ^{not invertible},
there $\exists v$. v is a nonzero vector, s.t. $(X^T X + \lambda I)v = 0$.

$$\therefore X^T X v + \lambda v = 0.$$

$$v^T X^T X v + v^T \lambda v = 0.$$

$$(Xv)^T Xv + \lambda v^T v = 0.$$

$$\therefore \|Xv\|_2^2 + \lambda \|v\|_2^2 = 0$$

$\therefore v$ is a nonzero vector and $\lambda > 0$

$$\therefore \lambda \|v\|_2^2 > 0.$$

$\therefore \|Xv\|_2^2 + \lambda \|v\|_2^2 > 0$, our suppose is contradict.

$$\therefore \forall v. (X^T X + \lambda I)v \neq 0.$$

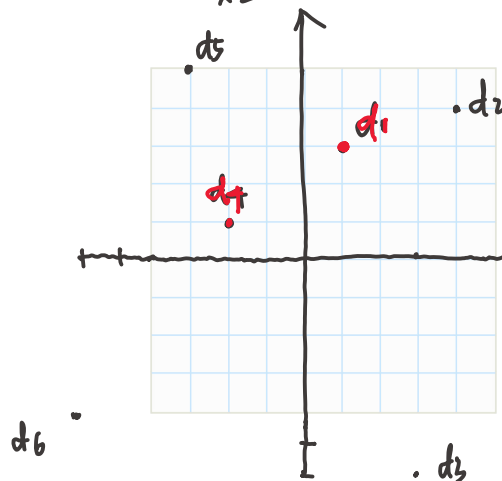
$X^T X + \lambda I$ is invertible. (by def).

$$\therefore \beta = (X^T X + \lambda I)^{-1} X^T y.$$

$$\therefore \beta^* = (X^T X + \lambda I)^{-1} X^T y$$

(C) - linear separable means there is a hyperplane that can separate the data of different classes.

so we can suppose a linear function $h(x) = \theta_1 x_1 + \theta_2 x_2 + \theta_0$.



there isn't a linear which can separate d_1, d_4 and d_2, d_3, d_5, d_6 , so the x_1 data set is not linear separable.

