

# Introduction to Machine Learning, Fall 2023

## Homework 1

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

October 13, 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  $X_i$ , derive the log-likelihood  $l(\mu, \Sigma)$  and MLE of  $\mu$  [4 points]
- (c) Suppose  $\hat{\theta}$  is an unbiased estimator of  $\theta$  and  $\text{Var}(\hat{\theta}) > 0$ . Prove that  $(\hat{\theta})^2$  is not an unbiased estimator of  $\theta^2$ . [3 points]

(a) As  $\mathbf{X}_i$ 's covariance matrix can be expressed as  $E[(\mathbf{X}_i - E(\mathbf{X}_i))(\mathbf{X}_i - E(\mathbf{X}_i))^T]$

For arbitrary column vector  $a$ , we want to calculate

$$a \cdot E[(\mathbf{X}_i - E(\mathbf{X}_i))(\mathbf{X}_i - E(\mathbf{X}_i))^T] a^T \\ = E[a(\mathbf{X}_i - E(\mathbf{X}_i))(\mathbf{X}_i - E(\mathbf{X}_i))^T a^T]$$

Because  $(\mathbf{X}_i - E(\mathbf{X}_i))^T a^T = (a(\mathbf{X}_i - E(\mathbf{X}_i)))^T$ ,

it can be rewritten as  $E[a(\mathbf{X}_i - E(\mathbf{X}_i)) \cdot (a(\mathbf{X}_i - E(\mathbf{X}_i)))^T]$ ,

which means that all elements in  $[ ]$  are bigger or equal to 0.

Thus the expected value of it is bigger or equal to 0.

Thus the covariance of  $\mathbf{X}_i$  is a semi positive definite matrix

$$(b) l(\mu, \Sigma) = -\frac{n}{2} \log(2\pi) - n \log \sqrt{\Sigma} - \frac{1}{2} \sum_{i=1}^n (\mathbf{X}_i - \mu)^2$$

$$\frac{\partial}{\partial \mu} l(\mu, \Sigma) = \frac{1}{2\sigma^2} \cdot 2 \cdot \sum_{i=1}^n (\mathbf{X}_i - \mu) = 0$$

$$\mu = \frac{\sum_{i=1}^n \mathbf{X}_i}{n}$$

(c) As  $\hat{\theta}$  is unbiased,

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n \hat{\theta}_i$$

$$E(\hat{\theta}^2) = \frac{1}{n^2} \left( \sum_{i=1}^n \hat{\theta}_i \right)^2 \neq \frac{1}{n^2} \sum_{i=1}^n \hat{\theta}_i^2 = E(\hat{\theta}^2).$$

Thus  $\hat{\theta}^2$  is biased.

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  $\epsilon$  is an independent variable, called a noise term, which is drawn from a Gaussian distribution with mean 0, and variance  $\sigma^2$ . 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) Finding  $a$  and  $b$  to make  $\log L(a, b) = \sum_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(y_i - ax_i - b)^2\right)$  reach the maximum

$$(b) \log L(a, b) = -\frac{n}{2} \log(2\pi) - n \cdot \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - (ax_i + b))^2$$

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

$$\frac{\partial}{\partial b} \log L(a, b) = -\frac{1}{2\sigma^2} \sum_{i=1}^n -2(y_i - ax_i - b) = \frac{1}{\sigma^2} \sum_{i=1}^n (-ax_i + y_i - b) \quad \text{②}$$

$$\begin{cases} \text{①} \\ \text{②} \end{cases} \Rightarrow \begin{cases} \sum_{i=1}^n x_i(-ax_i + y_i - b) = 0 \dots \text{③} \\ \sum_{i=1}^n (-ax_i + y_i - b) = 0 \dots \text{④} \end{cases} \Rightarrow \begin{cases} a = \frac{\sum_{i=1}^n x_i(\sum_{j=1}^n y_j) - n \sum_{i=1}^n x_i y_i}{(\sum_{i=1}^n x_i)^2 - n \sum_{i=1}^n x_i^2} \\ b = -a \cdot \bar{x} + \bar{y} \end{cases}$$

$$\text{③} \Rightarrow a \sum_{i=1}^n x_i^2 + \sum_{i=1}^n x_i y_i - b \sum_{i=1}^n x_i = 0$$

$$\text{④} \Rightarrow a \sum_{i=1}^n x_i + \sum_{i=1}^n y_i - n \cdot b = 0.$$

$$\Leftrightarrow -a \cdot \frac{1}{n} \sum_{i=1}^n x_i + \frac{1}{n} \sum_{i=1}^n y_i = b.$$

$$= -\frac{(\sum_{i=1}^n x_i)(\sum_{j=1}^n y_j) - n \sum_{i=1}^n x_i y_i \cdot (\sum_{i=1}^n x_i)}{(\sum_{i=1}^n x_i)^2 - n \sum_{i=1}^n x_i^2} + \frac{\sum_{i=1}^n y_i}{n}$$

(c) From ④ in (b), we know that  $-a \cdot \frac{1}{n} \sum_{i=1}^n x_i + \frac{1}{n} \sum_{i=1}^n y_i = b \Leftrightarrow -a \cdot \bar{x} + \bar{y} = b \Leftrightarrow \bar{y} = a \bar{x} + b$

which means that point  $(\bar{x}, \bar{y})$  is always on  $f(x)$

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  $\beta$  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  $\beta_*$  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

$$\begin{aligned}
 \text{(a)} & \text{ We suppose the regression function to be a linear function.} \\
 \text{(b)} & \text{ Let } f(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2 \\
 \frac{df(\beta)}{d\beta} &= 0 \Leftrightarrow -2\mathbf{X}^T(\mathbf{y} - \mathbf{X}\beta) + 2\lambda \mathbf{I}\beta = 0 \\
 &\Leftrightarrow (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})\beta = \mathbf{X}^T \mathbf{y}
 \end{aligned}$$

As  $\mathbf{X}^T \mathbf{X}$  is symmetric, let  $\mathbf{X}^T \mathbf{X} = \mathbf{Q} \mathbf{S} \mathbf{Q}^T$   
 $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} = \mathbf{Q} (\mathbf{S} + \lambda \mathbf{I}) \mathbf{Q}^T$

Because all elements on the diagonal line of  $\mathbf{S} + \lambda \mathbf{I} > 0$ ,  
 $\mathbf{S} + \lambda \mathbf{I}$  is invertible.

Thus  $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$  is invertible

$$\text{So } \beta = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

(c) Not linear separable.

If linear separable.

Suppose  $Ax + By + C = 0$ .

$(-2, 1)$ , and  $(1, 3)$  should be at the same side of the line.  
if  $(3, -6)$  at another side, the point  $(-3, 5)$  can't at the  
same side as  $(3, -6)$