

Name:

Student ID:

## CS 189: Introduction to Machine Learning

### Homework 2

Due: February 18, 2016 at 11:59pm

### Instructions

- Homework 2 is completely a written assignment; no coding involved.
- We prefer that you typeset your answers using the  $\text{\LaTeX}$  template on bCourses. If there is not enough space for your answer, you may continue your answer on the next page. Make sure to start each question on a new page.
- Neatly handwritten and scanned solutions will also be accepted. Make sure your answers are readable!
- Submit a PDF with your answers to the Homework 2 assignment on Gradescope. You should be able to see CS 189/289A on Gradescope when you log in with your bCourses email address. Please make a Piazza post if you have any problems accessing Gradescope.
- While submitting to Gradescope, you will have to select the pages containing your answer for each question.
- The assignment covers concepts in probability, linear algebra, matrix calculus, and decision theory.
- **Start early. This is a long assignment. Some of the material may not have been covered in lecture; you are responsible for finding resources to understand it.**

**Problem 1: Expected Value.**

A target is made of 3 concentric circles of radii  $1/\sqrt{3}$ , 1 and  $\sqrt{3}$  feet. Shots within the inner circle are given 4 points, shots within the next ring are given 3 points, and shots within the third ring are given 2 points. Shots outside the target are given 0 points.

Let  $X$  be the distance of the hit from the center (in feet), and let the probability density function of  $X$  be

$$f(x) = \begin{cases} \frac{2}{\pi(1+x^2)} & x > 0 \\ 0 & \text{otherwise} \end{cases}$$

What is the expected value of the score of a single shot?

**Solution:**

**Problem 2: MLE.**

Assume that the random variable  $X$  has the exponential distribution

$$f(x; \theta) = \theta e^{-\theta x} \quad x \geq 0, \theta > 0$$

where  $\theta$  is the parameter of the distribution. Use the method of maximum likelihood to estimate  $\theta$  if 5 observations of  $X$  are  $x_1 = 0.9$ ,  $x_2 = 1.7$ ,  $x_3 = 0.4$ ,  $x_4 = 0.3$ , and  $x_5 = 2.6$ , generated i.i.d. (i.e., independent and identically distributed).

**Solution:**

**Definition.** Let  $A \in \mathbb{R}^{n \times n}$  be a symmetric matrix. We say that  $A$  is **positive definite** if  $\forall x \in \mathbb{R}^n \mid x \neq \vec{0}, x^\top Ax > 0$ . Similarly, we say that  $A$  is **positive semidefinite** if  $\forall x \in \mathbb{R}^n, x^\top Ax \geq 0$ .

**Problem 3: Positive Definiteness.**

Let  $x = [x_1 \ \cdots \ x_n]^\top \in \mathbb{R}^n$ , and let  $A \in \mathbb{R}^{n \times n}$  be the square matrix

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix}$$

- (a) Give an explicit formula for  $x^\top Ax$ . Write your answer as a sum involving the elements of  $A$  and  $x$ .
- (b) Show that if  $A$  is positive definite, then the entries on the diagonal of  $A$  are positive (that is,  $a_{ii} > 0$  for all  $1 \leq i \leq n$ ).

**Solution:**

**Problem 4: Short Proofs.**

$A$  is symmetric in all parts.

- (a) Let  $A$  be a positive semidefinite matrix. Show that  $A + \gamma I$  is positive definite for any  $\gamma > 0$ .
- (b) Let  $A$  be a positive definite matrix. Prove that all eigenvalues of  $A$  are greater than zero.
- (c) Let  $A$  be a positive definite matrix. Prove that  $A$  is invertible. (Hint: Use the previous part.)
- (d) Let  $A$  be a positive definite matrix. Prove that there exist  $n$  linearly independent vectors  $x_1, x_2, \dots, x_n$  such that  $A_{ij} = x_i^\top x_j$ . (Hint: Use the spectral theorem and what you proved in (b) to find a matrix  $B$  such that  $A = B^\top B$ .)

**Solution:**

**Problem 5: Derivatives and Norm Inequalities.**

Derive the expression for following questions. Do not write the answers directly.

- (a) Let  $\mathbf{x}, \mathbf{a} \in \mathbb{R}^n$ . Derive  $\frac{\partial(\mathbf{x}^T \mathbf{a})}{\partial \mathbf{x}}$ .
- (b) Let  $\mathbf{A} \in \mathbb{R}^{n \times n}, \mathbf{x} \in \mathbb{R}^n$ . Derive  $\frac{\partial(\mathbf{x}^T \mathbf{A} \mathbf{x})}{\partial \mathbf{x}}$ .
- (c) Let  $\mathbf{A}, \mathbf{X} \in \mathbb{R}^{n \times n}$ . Derive  $\frac{\partial \text{Trace}(\mathbf{X} \mathbf{A})}{\partial \mathbf{X}}$ .
- (d) Let  $\mathbf{x} \in \mathbb{R}^n$ . Prove that  $\|\mathbf{x}\|_2 \leq \|\mathbf{x}\|_1 \leq \sqrt{n} \|\mathbf{x}\|_2$ . (Note that  $\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$  and  $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$ .) (Hint: The Cauchy-Schwarz inequality may come in handy.)

**Solution:**

**Problem 6: Weighted Linear Regression.**

Let  $\mathbf{X}$  be a  $n \times d$  data matrix,  $\mathbf{Y}$  be the corresponding  $n \times 1$  target/label matrix and  $\mathbf{\Lambda}$  be the diagonal  $n \times n$  matrix containing a weight for each example. More explicitly, we have

$$\mathbf{X} = \begin{bmatrix} (\mathbf{x}^{(1)})^T \\ (\mathbf{x}^{(2)})^T \\ \vdots \\ (\mathbf{x}^{(n)})^T \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} \mathbf{y}^{(1)} \\ \mathbf{y}^{(2)} \\ \vdots \\ \mathbf{y}^{(n)} \end{bmatrix} \quad \mathbf{\Lambda} = \text{diag}(\lambda^{(1)}, \lambda^{(2)}, \dots, \lambda^{(n)})$$

where  $\mathbf{x}^{(i)} \in \mathbb{R}^d$ ,  $\mathbf{y}^{(i)} \in \mathbb{R}$ , and  $\lambda^{(i)} > 0 \quad \forall i \in \{1 \dots n\}$ .  $\mathbf{X}$ ,  $\mathbf{Y}$  and  $\mathbf{\Lambda}$  are fixed and known.

In this question, we will try to fit a weighted linear regression model for this data. We want to find the value of weight vector  $\mathbf{w}$  which best satisfies the following equation  $\mathbf{y}^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)} + \epsilon^{(i)}$ , where  $\epsilon$  is noise. This is achieved by minimizing the weighted noise for all the examples. Thus, our risk (cost) function is defined as follows:

$$\begin{aligned} R[\mathbf{w}] &= \sum_{i=1}^n \lambda^{(i)} (\epsilon^{(i)})^2 \\ &= \sum_{i=1}^n \lambda^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} - \mathbf{y}^{(i)})^2 \end{aligned}$$

- Write this risk function  $R[\mathbf{w}]$  in matrix notation (i.e., in terms of  $\mathbf{X}$ ,  $\mathbf{Y}$ ,  $\mathbf{\Lambda}$  and  $\mathbf{w}$ ).
- Find the weight vector  $\mathbf{w}$  that minimizes the risk function obtained in the previous part. You can assume that  $\mathbf{X}^T \mathbf{\Lambda} \mathbf{X}$  is full rank. (Hint: You may use the expression you derived in Question 5(b).)
- The  $L_2$  regularized risk function, for  $\gamma > 0$ , is

$$R[\mathbf{w}] = \sum_{i=1}^n \lambda^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} - \mathbf{y}^{(i)})^2 + \gamma \|\mathbf{w}\|_2^2$$

Rewrite this new risk function in matrix notation as in (a) and solve for  $\mathbf{w}$  as in (b).

- How does  $\gamma$  affect the regression model? How does this fit in with what you already know about  $L_2$  regularization? (Hint: Observe the different expressions for  $\mathbf{w}$  obtained in (b) and (c).)

**Solution:**

**Problem 7: Classification.**

Suppose we have a classification problem with classes labeled  $1, \dots, c$  and an additional doubt category labeled as  $c + 1$ . Let the loss function be the following:

$$\ell(f(x) = i, y = j) = \begin{cases} 0 & \text{if } i = j \quad i, j \in \{1, \dots, c\} \\ \lambda_r & \text{if } i = c + 1 \\ \lambda_s & \text{otherwise} \end{cases}$$

where  $\lambda_r$  is the loss incurred for choosing doubt and  $\lambda_s$  is the loss incurred for making a misclassification. Note that  $\lambda_r \geq 0$  and  $\lambda_s \geq 0$ .

Hint: The risk of classifying a new datapoint as class  $i \in \{1, 2, \dots, c + 1\}$  is

$$R(\alpha_i|x) = \sum_{j=1}^c \ell(f(x) = i, y = j)P(\omega_j|x)$$

- (a) Show that the minimum risk is obtained if we follow this policy: (1) choose class  $i$  if  $P(\omega_i|x) \geq P(\omega_j|x)$  for all  $j$  and  $P(\omega_i|x) \geq 1 - \lambda_r/\lambda_s$ , and (2) choose doubt otherwise.
- (b) What happens if  $\lambda_r = 0$ ? What happens if  $\lambda_r > \lambda_s$ ? Is this consistent with your intuition?

**Solution:**



**Problem 8: Gaussians.**

Let  $P(x | \omega_i) \sim \mathcal{N}(\mu_i, \sigma^2)$  for a two-category, one-dimensional classification problem with  $P(\omega_1) = P(\omega_2) = 1/2$ . Here, the classes are  $\omega_1$  and  $\omega_2$ . For this problem, we have  $\mu_2 \geq \mu_1$ .

- (a) Find the optimal Bayes decision boundary (i.e., find  $x$  such that  $P(\omega_1 | x) = P(\omega_2 | x)$ ). What is the corresponding decision rule?
- (b) Show that the Bayes error associated with this decision rule is

$$P_e = \frac{1}{\sqrt{2\pi}} \int_a^\infty e^{-z^2/2} dz$$

where  $a = \frac{\mu_2 - \mu_1}{2\sigma}$ . The Bayes error is the probability of misclassification:

$$P_e = P(\text{misclassified as } \omega_1 | \omega_2)P(\omega_2) + P(\text{misclassified as } \omega_2 | \omega_1)P(\omega_1).$$

**Solution:**