

Midterm Exam

Midterm exam will be done individually. Use of partial or entire solutions obtained from others or online is strictly prohibited.

- There will be 7 pages in this exam (including this cover sheet)
- This is a **CLOSED-BOOK** exam. You can use your one-sided half-page cheatsheet. You **CANNOT** use materials brought by other students.
- If you need more room to work out your answer to a question, use the back of the page and clearly mark this on the front of the page.
- Work efficiently and independently.
- You have 150 minutes.
- Good luck!

Question	Topic	Max. score	Score
1	Bayesian Decision Theory	15	
2	Logistic Regression	15	
3	Perceptron Algorithm	15	
4	Nonparametric Method	15	
5	Maximum Likelihood Estimator	20	
6	Short Answer Questions	20	
Total		100	

1. **Bayesian Decision Theory** (15 points)

- (a) (5 pts) Assume we have c classes $\omega_1, \dots, \omega_c$, and feature vector \mathbf{x} , write down the Bayesian decision rule for classification *in terms of* prior probabilities of classes, i.e., $P(\omega_i)$, and class conditional densities of \mathbf{x} i.e., $p(\mathbf{x}|\omega_i)$.

- (b) (10 pt) Consider two-class classification, suppose $p(\mathbf{x}|\omega_1)$ is standard normal distribution and $p(\mathbf{x}|\omega_2)$ is uniform distribution over $[-\frac{1}{2}, \frac{1}{2}]$, i.e.,

$$\begin{aligned} p(\mathbf{x}|\omega_1) &= \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \\ p(\mathbf{x}|\omega_2) &= 1, \mathbf{x} \in [-\frac{1}{2}, \frac{1}{2}] \end{aligned}$$

Assuming zero-one loss and $P(\omega_1) = P(\omega_2)$, derive the corresponding decision rule.

-
2. **Logistic Regression** (15 points) In the lecture, we consider the logistic regression for binary classification on N observations $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$:

$$f(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$$

where $y_i \in \{0, 1\}$, $\sigma(\cdot)$ is sigmoid function, and denote $\mathbf{y} = (y_1, y_2, \dots, y_N)^T$. We showed to learn \mathbf{w} , maximizing the likelihood $p(\mathbf{y}|\mathbf{w})$ is equivalent to minimize the following *cross entropy* error function:

$$E(\mathbf{w}) = \sum_{i=1}^N -\{y_i \log f(\mathbf{x}_i) + (1 - y_i) \log(1 - f(\mathbf{x}_i))\}$$

Now consider the penalized version with the following error function:

$$E_p(\mathbf{w}) = \sum_{i=1}^N -\{y_i \log f(\mathbf{x}_i) + (1 - y_i) \log(1 - f(\mathbf{x}_i))\} + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- (a) (7 pts) Write down the gradient of $E_p(\mathbf{w})$ w.r.t \mathbf{w} , i.e., $\nabla_{\mathbf{w}} E_p(\mathbf{w})$.

- (b) (8 pts) In terms of probabilities, minimizing the above penalized error function $E_p(\mathbf{w})$ is equivalent to maximizing the corresponding posterior $p(\mathbf{w}|\mathbf{y})$. Show that it is indeed the case by writing down the likelihood $p(\mathbf{y}|\mathbf{w})$ and prior $p(\mathbf{w})$.

-
3. **Perceptron Algorithm** (15 points) Assume that you are given observations (x_1, x_2) with their labels in the following order:

Instance	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
label	+	-	+	-	+	-	+	+
Data (x_1, x_2)	(10,10)	(0,0)	(8,4)	(3,3)	(4,8)	(0.5, 0.5)	(4,3)	(2,5)

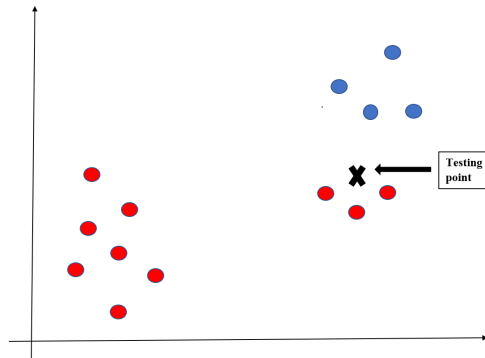
Show the action of the perceptron algorithm for the above sequence of observations. We start with $x_1 + x_2 = 0$. For points on the decision boundary, the model will predict the point to be *positive*. Note that the detailed steps of perceptron algorithm are needed.

4. **Nonparametric Methods** (15 points)

- (a) (8 pts)[1D Parzen window] Suppose we have 5 points $\{0, 1, 2, 4, 5\}$ that coming from 1D unknown distribution $p(x)$. Use Parzen window to estimate density value at $x = 3$. Assume we use *unit interval* as our window function $k(u)$, and the window width (interval length) $h = 5$:

$$k(u) = \begin{cases} 1, & |u| \leq 1/2 \\ 0, & \text{otherwise} \end{cases}$$

- (b) (7 pts)[K nearest neighbour] Suppose we have the following data points that comes from two classes (red vs blue), X is the testing point. Distance is measured using Euclidean distance. Consider K nearest neighbour approach for testing point classification.



- (2 pts) If $K = 1$, what is the class of X ?
- (2 pts) If $K = 7$, what is the class of X ?
- (3 pts) In general, use large value of K might be a good choice. Based on our data points, would you recommend use $K = 9$ nearest neighbours? Why or why not?

5. **Maximum Likelihood Estimator** (20 points) Suppose we have training samples $\{x_1, x_2, \dots, x_n\}$. Consider the following distributions: (please write down the derivation steps)

(a) (10 pts) $f(x; \theta) = 2\theta x e^{-\theta x^2}$, $x > 0$, find MLE for θ

(b) (10 pts) $Unif[\theta_1, \theta_2]$, $Unif$ represents uniform distribution, find the MLE for θ_1 and θ_2 .

6. Short Answer Questions (20 points)

- (a) (5 pts) Describe main differences between generative approach, discriminative approach and discriminant function approach for classification.

- (b) (5 pts) Describe main differences between maximum likelihood estimation and Bayesian estimation.

- (c) (5 pts) Briefly compare the Principal Component Analysis (PCA) and Fisher's Linear Discriminant (FLD).

- (d) (5 pts) Describe the main difference between batch gradient descent, mini-batch gradient descent and stochastic gradient descent. What are the pros and cons for batch gradient descent?