

CS 6375 Machine Learning Midterm Examination

University of Texas at Dallas

10/20/2021

Name: _____

NetID: _____

Question	Topic	Points
1	Short Answers	20
2	True/False Questions	8
3	Naive Bayes	10
4	Perceptrons	10
5	SVMs	24
6	Logistic Regression	18
7	Decision Trees	10
Total		100

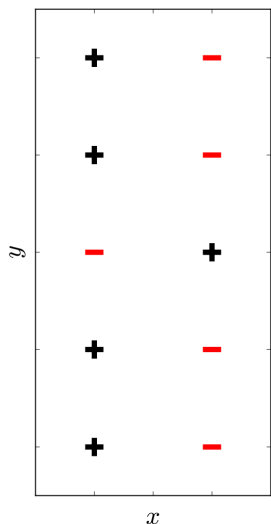
Instructions:

1. This examination contains 14 pages, including this page.
2. You have **two (2) hours** to complete the examination.
3. PLEASE SHOW ALL THE STEPS CLEARLY OF HOW YOU CAME UP WITH THE FINAL ANSWER. JUST THE FINAL ANSWER WITH NO WORK WILL BE ZERO POINTS!!
4. Either you can use this paper or separate set of sheets to fill in your answers. Write clearly so we can understand your handwriting. Please scan the answers once you are done (you can also take a photograph once you are done) and upload the scanned copy AS A COMBINED PDF to eLearning. You can also use a tablet device (like an Ipad) to write if you prefer.
5. Please do not search online for answers to the questions. If the answers are similar to something available online, you will get zero points on this examination.
6. The mid-term examination has to be done individually by everyone. If someone copies, the entire group of students involved will get a zero.
7. Work efficiently. Some questions are easier, some more difficult. Be sure to give yourself time to answer all of the easy ones, and avoid getting bogged down in the more difficult one.
8. All the Best!!

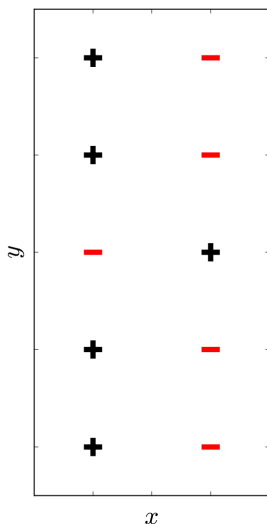
Question 1: Short Answers

[20 pts] Please provide short and clear answers for the questions below.

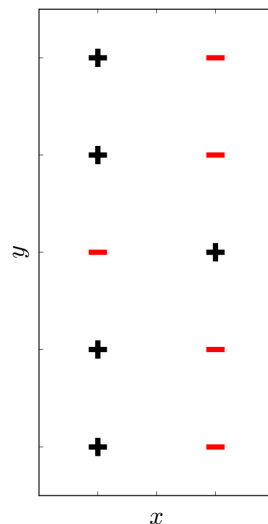
- (a) Given a 2 dimensional dataset below, draw the decision boundaries obtained by the following classifiers (4 points: a and b are 1 point each and c is for 2 points)



(a) Logistic regression ($\lambda = 0$)



(b) 1-NN



(c) 3-NN

- (b) (6 points: each sub-part is 2 points each) A random variable follows an exponential distribution with parameter $\lambda : \lambda > 0$, and has the following density:

$$p(t) = \lambda e^{-\lambda^2 t^2}, t \in [0, \infty] \quad (1)$$

This distribution models waiting times between events. Given a iid data: $T = (t_1, \dots, t_n)$, where each t_i is modeled as drawn from the exponential distribution with parameter λ . Then:

- Compute the log-likelihood $p(T|\lambda)$
- Solve for λ_{MLE}
- Suppose we have a prior distribution: $p(\lambda) \propto e^{-\mu\lambda}$. Obtain λ_{MAP} . Compare λ_{MLE} and λ_{MAP} as $n \rightarrow \infty$.

- (c) (6 points) Decision Trees: Using the dataset below, we want to build a decision tree which classifies Y as T or F given the binary variables A, B, C.

(Part 1) Draw the tree that would be learned by the greedy algorithm with zero training error. You do not need to show any computation.

A	B	C	Y
F	F	F	F
T	F	T	T
T	T	F	T
T	T	T	F

(Part 2) Is this tree optimal (i.e. does it get zero training error with minimal depth)? Explain in less

than two sentences. If it is not optimal, draw the optimal tree as well.

(d) (4 points) Compute the dual of the following problem. Minimize $x^2 + 1$ such that $(x - 1)(x - 3) \leq 2$.

Question 2: True or False with Explanations

[8 pts] Each question below is for 2 points each. Please do not just write true or false. You also need to provide explanations. Just true or false will yield no points.

(a) Given a function f , the set of sub-gradients is always non-empty.

(b) k -nearest neighbor models can be used for regression.

(c) No classifier can do better than a naive Bayes classifier if the distribution of the data is known.

(d) Maximizing the likelihood of logistic regression yields multiple local optimums

Question 3: Naive Bayes

[10 pts] Given the following training data points, provide the output of naive bayes classifier (trained with MLE) for test data points $z1$ and $z2$. Give the same answer if we have consider Laplace smoothing (i.e., trained with MAP) with $\alpha = 1$:

$x1 = (0, 0, 0, 1, 0, 0, 1), y1 = 1$
 $x2 = (0, 0, 1, 1, 0, 0, 0), y2 = 1$
 $x3 = (1, 1, 0, 0, 0, 1, 0), y3 = -1$
 $x4 = (1, 0, 0, 0, 1, 1, 0), y4 = -1$
 $x5 = (0, 1, 1, 1, 1, 1, 1), y5 = 1$
 $x6 = (0, 0, 0, 0, 0, 0, 0), y6 = 1$
 $x7 = (0, 0, 1, 1, 1, 1, 1), y7 = -1$
 $z1 = (1, 0, 0, 0, 0, 1, 0)$
 $z2 = (0, 1, 1, 0, 0, 1, 1)$

Question 4: Perceptrons

[10 pts] Demonstrate how the perceptron without bias (i.e. we set the parameter $b = 0$ and keep it fixed) updates its parameters given the following training sequence:

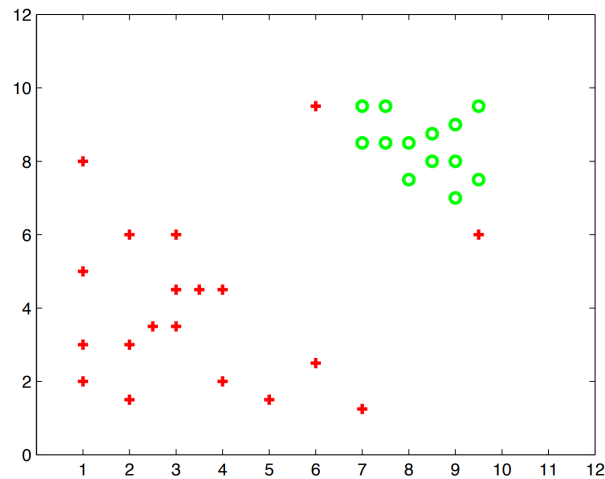
$$\begin{aligned}x1 &= (0, 0, 0, 1, 0, 0, 1), \quad y1 = 1 \\x2 &= (1, 1, 0, 0, 0, 1, 0), \quad y2 = -1 \\x3 &= (0, 0, 1, 1, 0, 0, 0), \quad y3 = 1 \\x4 &= (1, 0, 0, 0, 1, 1, 0), \quad y4 = -1 \\x5 &= (1, 0, 0, 0, 0, 1, 0), \quad y5 = -1\end{aligned}$$

Start with the weights all zeros. Is the data linearly separable? Add one more point to the dataset such that the resulting dataset will not be linearly separable any more.

Question 5: SVMs

[24 pts] This question will cover support vector machines.

- (a) (15 points) Given the following dataset (shown in the figure below), assume we are training the SVM with a quadratic kernel. The slack penalty C will determine the location of the separating hyper-plane. Each question below is for 3 points each.



- Where would the decision boundary be if $C = 10^{10}$?
- For $C = 10^{-10}$, where would the decision boundary be?
- If we know that we cannot fully trust the obtained data points, which of the scenarios would we

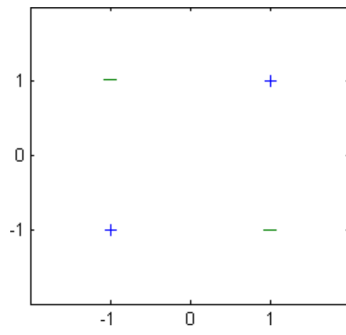
prefer to train the model? $C = 10^4$ or $C = 10^{-4}$ or neither?

- Given $C = 10^5$, draw an additional data point which *will not* change the decision boundary. Justify.

- Given $C = 10^5$, draw an additional data point which *will* change the decision boundary. Justify.

(b) (9 points) Consider the data set below. Under which of the following feature vectors is the data linearly separable? For full credit, you must justify your answer by either providing a linear separator or explaining why such a separator does not exist. Each part is 1.8 points each

- $\phi(x_1, x_2) = [x_1 + x_2 + 1, x_1 - x_2 - 1]$
- $\phi(x_1, x_2) = [x_1^2, x_2^2, x_1 x_2]$
- $\phi(x_1, x_2) = [\exp(x_1) + 1, \exp(x_2) - 1]$
- $\phi(x_1, x_2) = [x_1 \sin x_2, x_1]$
- $\phi(x_1, x_2) = [x_1 x_2, x_1]$



Question 6: Logistic Regression

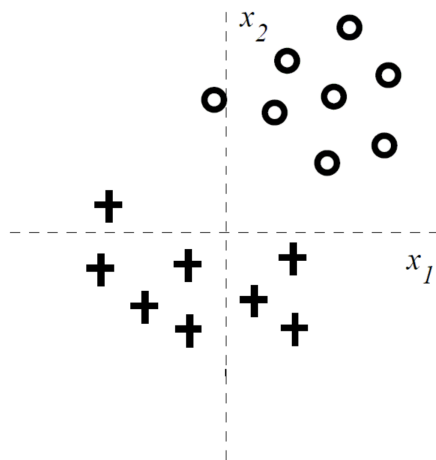
[18 pts] Given a two dimensional dataset shown in the figure below, we attempt to solve a simple binary classification using a linear logistic regression. The model is:

$$p(y = 1|x, w_0, w_1, w_2) = \frac{1}{1 + \exp(-w_0 - w_1x_1 - w_2x_2)} \quad (2)$$

Consider training a regularized linear logistic regression model where we try to maximize:

$$\sum_{i=1}^n \log P(y_i|x_i, w_0, w_1, w_2) - C_0w_0^2 - C_1w_1^2 - C_2w_2^2 \quad (3)$$

Note we have different regularization parameters for each coordinate. Draw the approximate decision bound-



aries in the following cases (each sub-question is 3 points). Also explain what will be the training error in each case.

(a) When $C_0, C_1, C_2 = 10^{-10}$

(b) When $C_0 = 10^{10}$ and $C_1, C_2 = 10^{-10}$

(c) When $C_1 = 10^{10}$ and $C_0, C_2 = 10^{-10}$

(d) When $C_2 = 10^{10}$ and $C_0, C_1 = 10^{-10}$

(e) Consider the case when $C_1 = C_2 = 10^{10}$ and $C_0 = 10^{-10}$. What is the value of w_0 that we expect to obtain if the dataset is balanced?

(f) In the above case, assume we have more negatives compared to positives in this dataset (i.e. make it imbalanced). What is the value/range of w_0 we expect? (you can give a range of values of w_0 if you prefer).

Question 7: Decision Trees

[10 pts] Consider the dataset shown below. We will use this dataset to learn a decision tree which predicts if people pass machine learning (Yes or No), based on their previous GPA (High, Medium, or Low) and whether or not they studied. For this problem, assume that $\log_2 3 \approx 1.6$ (it is ok if you leave the answers in

GPA	Studied	Passed
L	F	F
L	T	T
M	F	F
M	T	T
H	F	T
H	T	T

terms of the logs as well).

- (a) (6 points) Compute $H(\text{Passed})$, $H(\text{Passed} \mid \text{GPA})$ and $H(\text{Passed} \mid \text{Studied})$

(b) (4 points) Draw the full decision tree that would be learned for this dataset.