

USC CSCI 467
Intro to Machine Learning

Midterm Exam
October 10, 2023, 2:00-3:20pm

Fall 2023
Instructor: Robin Jia

Name: _____

USC e-mail: _____@usc.edu

Answer the questions in the spaces provided. **If you write solutions on the back of the pages, indicate this on the front of the pages so we know to look there, but please try to avoid this if possible.** You may use the backs of pages for scratch work. This exam has 5 questions, for a total of 100 points.

Question 1: Weighted Linear Regression (28 points)

Consider a modification of the standard linear regression setup where each datapoint is associated with an importance weight. Formally, we have a training dataset consisting of n examples $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$, where each datapoint is associated with an importance weight $r_i > 0$. The weighted residual sum of squares objective is defined as

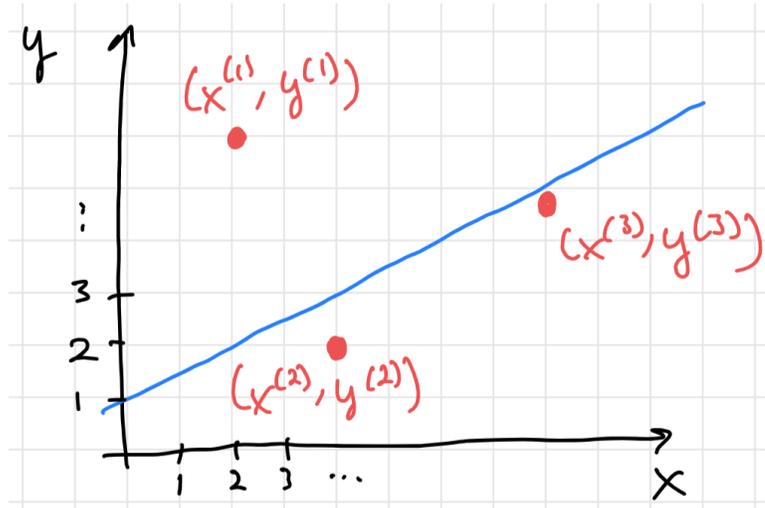
$$\text{WRSS}(w) = \sum_{i=1}^n r_i \left(w^\top x^{(i)} - y^{(i)} \right)^2.$$

(a) (5 points) Derive the expression for $\nabla \text{WRSS}(w)$.

(b) (6 points) Let \mathbf{X} be the $n \times d$ matrix whose i -th row is $x^{(i)\top}$, \mathbf{y} be the n -dimensional column vector whose i -th entry is $y^{(i)}$, and \mathbf{R} be the diagonal matrix where $\mathbf{R}_{ii} = r_i$ for all i and 0 for all other entries. Show that the WRSS objective can be written as follows in matrix form:

$$\text{WRSS}(w) = (\mathbf{X}w - \mathbf{y})^\top \mathbf{R}(\mathbf{X}w - \mathbf{y}).$$

- (c) The diagram below shows a training dataset with three examples $(x^{(1)}, y^{(1)})$, $(x^{(2)}, y^{(2)})$, and $(x^{(3)}, y^{(3)})$.



Here the inputs x are all one-dimensional. Soumya chooses some importance weights r_1 , r_2 , and r_3 and then trains a model to minimize the WRSS objective. The model's predictions are illustrated by the blue line.

- i. (2 points) What are the values of w and b that Soumya learned?

- ii. (4 points) Out of the values for r_1 , r_2 , and r_3 that Soumya chose, which one was the smallest? Explain your reasoning.

(d) (3 points) Explain in 1-2 sentences how allowing $r_i < 0$ for some i could lead to undesirable behavior.

(e) Let w^* be any value of w that achieves the lowest possible value of the standard linear regression loss discussed in class. Let \tilde{w} be any value of w that achieves the lowest possible WRSS loss when $r_i = 10$ for all i .

i. (5 points) Assume that the matrix $X^\top X$ is invertible. Are w^* and \tilde{w} guaranteed to be the same? Explain your reasoning.

- ii. (3 points) Now assume that $X^T X$ is not invertible. Does this change your answer? Explain your reasoning.

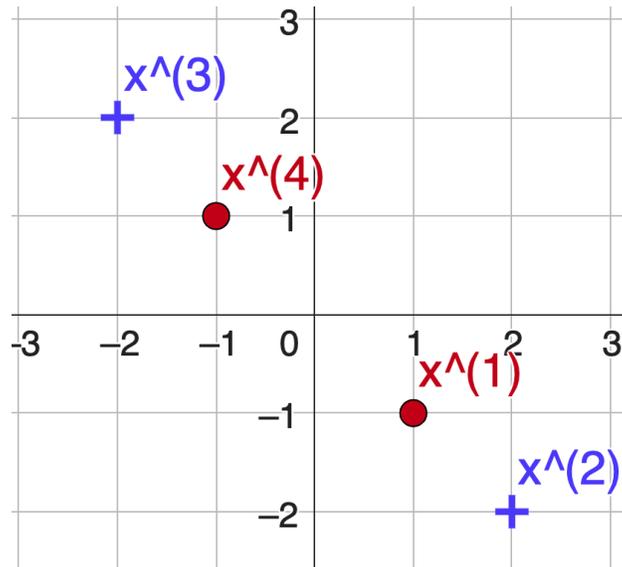
Question 2: Logistic Regression and Kernels (23 points)

In class, we discussed performing classification tasks using kernelized logistic regression. One popular choice of kernel function is the polynomial kernel, defined as

$$k(x, z) = (x^\top z + p)^q$$

where $p, q \in \mathbb{R}$ are hyperparameters.

In this problem, we will explore using this kernel for different values of p and q on the toy dataset shown below. Let $x \in \mathbb{R}^2$, where each point has 2 features x_1 and x_2 . In the plot below, the pluses are positive points and the circles are negative points. The plot shows x_1 on the horizontal axis and x_2 on the vertical axis.



- (a) (4 points) Assume that $p = 1$. What is the smallest possible positive integer value of q such that the model can learn a decision boundary that perfectly separates the positive and negative datapoints? Explain your reasoning.

(b) (7 points) Suppose that $q = 1$. **Prove** that for any value of p , the decision boundary learned by kernel logistic regression will be a straight line.

(c) (4 points) Is there a value of $q > 1000$ such that it is impossible to perfectly classify the four datapoints in the figure using kernelized linear regression? Explain your reasoning.

(d) (4 points) In general when using kernel methods, provide one potential drawback of using a very large value of q (e.g., $q = 1000$). Explain your reasoning.

(e) (4 points) Finally, let's see how other classification methods would perform on this same dataset. Suppose we ran 1-nearest neighbor using 4-fold cross-validation (i.e., each example is in its own fold). What would be the validation accuracy? Show your work. Use the Euclidean distance to compute nearest neighbors.

(c) (6 points) Lorena now wants to train her RNN to generate text. In particular:

- She has a training document with words $u_1, u_2, \dots, u_T, u_{T+1}$. u_1 is the special “beginning of sequence” token [BEGIN], and u_{T+1} is the special “end of sequence” token [END]. For example, if her document was “Fight on,” then we would have $T = 3$ and $u_1 = [\text{BEGIN}]$, $u_2 = \text{Fight}$, $u_3 = \text{on}$, and $u_4 = [\text{END}]$.
- During training, she feeds this sequence of words to the RNN one at a time (i.e., she uses teacher forcing), except for the final [END] token. This generates a sequence of hidden states h_1, \dots, h_T .
- The hidden states are fed directly into a softmax regression-type layer to predict the next word. This layer has a separate parameter vector $w^{(u)}$ for each word u in the vocabulary V .
- Lorena computes the loss on this document by summing up the losses for generating each “next word” from u_2 to u_{T+1} .

Write the mathematical expression for the training loss Lorena computes on this training document. It may help to remember that in softmax regression, given input x and C weight vectors $w^{(1)}, \dots, w^{(C)}$, the probability that $y = j$ is given by

$$P(y = j \mid x) = \frac{\exp(w^{(j)\top} x)}{\sum_{k=1}^C \exp(w^{(k)\top} x)}.$$

(d) (4 points) Bill doesn't like RNNs, so he comes up with an idea to use an MLP instead. Since he anticipates the maximum length of a document to be 200 words, he plans to build a model that can take in the concatenation of 200 word vectors as input. (If a document has less than 200 words, he will concatenate enough zero vectors to make it the same length as 200 word vectors.) He will use 50-dimensional word vectors, and an MLP with 500 hidden units. How many parameters will he need in the first layer of his MLP? Note: Include both the weight matrix and bias vector.

(e) (4 points) Using the terms "bias" and/or "variance," explain why having a large number of parameters in your model can make it harder to have high test accuracy.

Question 4: Short Response (12 points)

Answer the following questions and **explain your reasoning fully**. You may also draw explanatory diagrams when appropriate.

(a) (4 points) Is the function $f(x) = \max(x, 0)$ from $\mathbb{R} \rightarrow \mathbb{R}$ a convex function?

(b) (4 points) When training a neural network using stochastic gradient descent, is the loss on the training dataset guaranteed to decrease at every iteration?

(c) (4 points) Generative classifiers make predictions by estimating $P(y)$ and $P(x | y)$. Is the Naive Bayes assumption used when modeling $P(y)$, $P(x | y)$, or both?

Question 5: Multiple Choice (13 points)

In the following questions, circle the correct answer(s). There is no need to explain your answer.

- (a) (2 points) Suppose we are training a neural network with SGD using a batch size of 50, on a training dataset with 50,000 examples. How many total gradient updates would be performed after 5 epochs?
- A. 1,000
 - B. 5,000
 - C. 50,000
 - D. 250,000
- (b) (2 points) Which of the following situations is called overfitting?
- A. Low training error, low test error
 - B. High training error, high test error
 - C. High training error, low test error
 - D. Low training error, high test error
- (c) (3 points) Consider a training set $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$ and a probabilistic model $\Pr(y | x; w)$ which specifies the probability of seeing outcome y given x and parameterized by w . Which of the following is the value of w chosen by Maximum Likelihood Estimation (MLE)? Choose all that apply.
- A. $\arg \max_w \sum_{i=1}^n P(y^{(i)} | x^{(i)}; w)$
 - B. $\arg \max_w \prod_{i=1}^n P(y^{(i)} | x^{(i)}; w)$
 - C. $\arg \max_w \sum_{i=1}^n \log P(y^{(i)} | x^{(i)}; w)$
 - D. $\arg \max_w \prod_{i=1}^n \log P(y^{(i)} | x^{(i)}; w)$
- (d) (3 points) Which of the following would **not** be a valid loss function for linear regression? Choose all that apply.
- A. $\mathcal{L}(w) = \frac{1}{n} \sum_{i=1}^n (w^\top x^{(i)} - y^{(i)})^3$
 - B. $\mathcal{L}(w) = \frac{1}{n} \sum_{i=1}^n (w^\top x^{(i)} - y^{(i)})^2$
 - C. $\mathcal{L}(w) = \frac{1}{n} \sum_{i=1}^n |w^\top x^{(i)} - y^{(i)}|$
 - D. $\mathcal{L}(w) = \frac{1}{n} \sum_{i=1}^n (w^\top x^{(i)} - y^{(i)})$
- (e) (3 points) Which of the following strategies may help address overfitting? Choose all that apply.
- A. Acquire more training data.
 - B. Train the model with a small subset of the training dataset.
 - C. Apply regularization on the parameters.
 - D. Remove regularization on the parameters.