

HOMework 3

DECISION TREES, K-NN, PERCEPTRON, REGRESSION *

10-301 / 10-601 INTRODUCTION TO MACHINE LEARNING (FALL 2022)

<http://mlcourse.org>

OUT: Sep. 21, 2022

DUE: Sep. 28, 2022

TAs: Brandon, Calvin, Monica, Shubham, Chu

Summary It's time to practice what you've learned! In this assignment, you will answer questions on topics we've covered in class so far, including Decision Trees, K-Nearest Neighbors, Perceptron, and Linear Regression. This assignment consists of a written component split into four sections, one for each topic. These questions are designed to test your understanding of the theoretical and mathematical concepts related to each topic. For each topic, you will also apply your understanding of the concept to the related ideas such as overfitting, error rates, and model selection. This homework is designed to help you apply what you've learned and solve a few concrete problems.

START HERE: Instructions

- **Collaboration Policy:** Please read the collaboration policy here: <http://www.cs.cmu.edu/~mgormley/courses/10601/syllabus.html>
- **Late Submission Policy:** For this homework, you will only have 2 late days instead of the usual 3. This allows us to provide feedback before the exam. See the late submission policy here: <http://www.cs.cmu.edu/~mgormley/courses/10601/syllabus.html>
- **Submitting your work:** You will use Gradescope to submit answers to all questions and code. Please follow instructions at the end of this PDF to correctly submit all your code to Gradescope.
 - **Written:** For written problems such as short answer, multiple choice, derivations, proofs, or plots, please use the provided template. Submissions can be handwritten onto the template, but should be labeled and clearly legible. If your writing is not legible, you will not be awarded marks. If your scanned submission misaligns the template, there will be a 5% penalty. Alternatively, submissions can be written in LaTeX. Each derivation/proof should be completed in the boxes provided. If you do not follow the template, your assignment may not be graded correctly by our AI assisted grader.

*Compiled on Sunday 25th September, 2022 at 18:19

Instructions for Specific Problem Types

For “Select One” questions, please fill in the appropriate bubble completely:

Select One: Who taught this course?

- ☒ Matt Gormley
- ☐ Marie Curie
- ☐ Noam Chomsky

If you need to change your answer, you may cross out the previous answer and bubble in the new answer:

Select One: Who taught this course?

- ☒ Henry Chai
- ☐ Marie Curie
- ☒ Noam Chomsky

For “Select all that apply” questions, please fill in all appropriate squares completely:

Select all that apply: Which are scientists?

- ☒ Stephen Hawking
- ☒ Albert Einstein
- ☒ Isaac Newton
- ☐ I don't know

Again, if you need to change your answer, you may cross out the previous answer(s) and bubble in the new answer(s):

Select all that apply: Which are scientists?

- ☒ Stephen Hawking
- ☒ Albert Einstein
- ☒ Isaac Newton
- ☒ I don't know

For questions where you must fill in a blank, please make sure your final answer is fully included in the given space. You may cross out answers or parts of answers, but the final answer must still be within the given space.

Fill in the blank: What is the course number?

10-601

10-~~6~~301

1 Latex Bonus Point (1 points)

1. (1 point) **Select one:** Did you use Latex for the entire written portion of this homework?

☐ Yes

☐ No

2 Decision Tree (Revisited) (15 points)

1. Consider the following 4×4 checkerboard pattern. Suppose our goal is to perfectly classify the range shown such that all black regions are labeled as $+1$ and all white regions are labeled as -1 . Let the horizontal axis denote feature x_1 and vertical axis denote feature x_2 .

NOTE: If a point is on the border or corner of a region, it has the same label as the region that is above it and/or to its right. For example, $(1, 0)$ has label $+1$, $(1, 1)$ has label -1 , and $(1, 1.5)$ has label -1 .

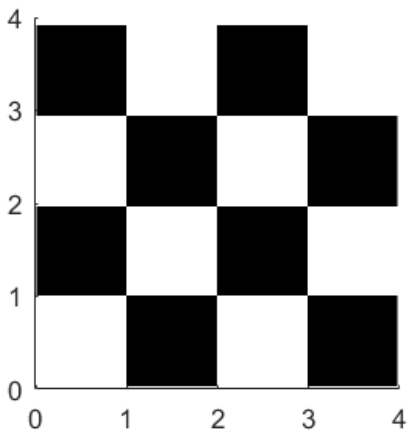


Figure 1: Checkerboard Pattern

- (a) (2 points) What is the minimum depth of a binary decision tree that perfectly classifies the colored regions in Figure 1, using features that only inspect either x_1 or x_2 but not both x_1 and x_2 ? Feel free to get creative in your use of the features x_1 and x_2 for the splits!

Some feature examples: $\text{ceil}(x_1)$, $\text{floor}(x_2)\%4$, $2 < x_1 < 4$, $x_2 < 1$ or $x_2 > 3$

Your Answer

- (b) (2 points) What is the minimum depth of a binary decision tree that perfectly classifies the colored regions in Figure 1, using *only* features of the form $x_1 < c$ or of the form $x_2 < c$ for different constants c ?

Your Answer

- (c) (2 points) What is the minimum depth of a binary decision tree that perfectly classifies the colored regions in Figure 1, using ANY features that involve x_1 , x_2 , or both?

Your Answer

2. Consider the graph below analyzing the size of tree vs. accuracy for a decision tree which has been pruned back to the vertical (red) line. Assume all of our labeled data for this task was randomly divided into a training dataset D_{train} , a validation data set D_{val} , and a test dataset D_{test} . The full tree was trained only on D_{train} , then reduced-error pruning was applied using D_{val} .

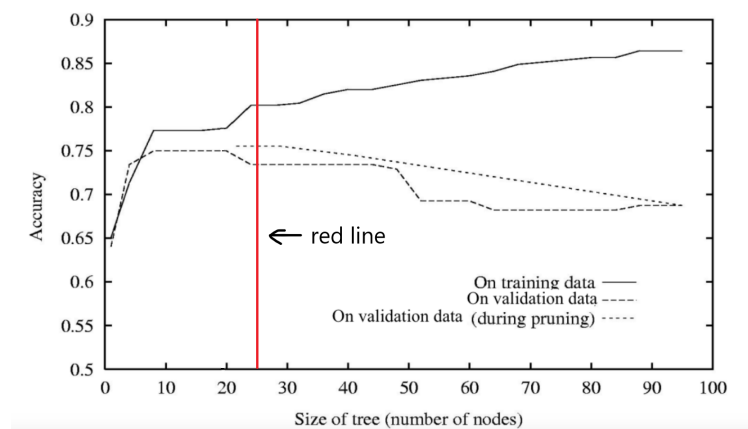


Figure 2: Pruned decision tree

- (a) (1 point) **Select one:** Refer to Figure 2. Note that D_{test} was not used for training or pruning. When the size of the pruned tree is at **25 nodes**, what is its accuracy on D_{test} likely to be?
- ☐ Slightly higher than the pruned tree's accuracy on D_{val}
 - ☐ The same as the pruned tree's accuracy on D_{val}
 - ☐ Slightly lower than the pruned tree's accuracy on D_{val}

(b) (1 point) **Select one:** Which of the following gives us the best approximation of the true error?

- ☐ Error corresponding to training data D_{train}
- ☐ Error corresponding to validation data D_{val}
- ☐ Error corresponding to test data D_{test}

3. (2 points) **Select all that apply:** Which of the following are valid ways to avoid overfitting?

- ☐ Decrease the training set size.
- ☐ Set a threshold for a minimum number of examples required to split at an internal node.
- ☐ Prune the tree so that cross-validation error is minimal.
- ☐ Increase the tree depth.
- ☐ None of the above.

4. Consider a binary classification problem using 1-nearest neighbors with the Euclidean distance metric. We have N 1-dimensional training points $x^{(1)}, x^{(2)}, \dots, x^{(N)}$ and corresponding labels $y^{(1)}, y^{(2)}, \dots, y^{(N)}$, with $x^{(i)} \in \mathbb{R}$ and $y^{(i)} \in \{0, 1\}$.

Assume the points $x^{(1)}, x^{(2)}, \dots, x^{(N)}$ are in ascending order by value. If there are ties during the 1-NN algorithm, we break ties by choosing the label corresponding to the $x^{(i)}$ with lower value.

(a) (2 points) **Select one:** Is it possible to build a decision tree that behaves exactly the same as the 1-nearest neighbor classifier? Assume that the decision at each node takes the form of “ $x \leq t$ ” or “ $x > t$ ”, where $t \in \mathbb{R}$.

- ☐ Yes
- ☐ No

If your answer is yes, please explain how you will construct the decision tree. If your answer is no, explain why it's not possible.

Your answer:

- (b) (3 points) **Select all that apply:** Let's add a dimension! Now the training points are 2-dimensional where $\mathbf{x}^{(i)} = (x_1^{(i)}, x_2^{(i)}) \in \mathbb{R}^2$. For which of the following training sets is it possible to build a decision tree that behaves exactly the same as a 1-nearest neighbor classifier? Assume that the decision at each node takes the form of " $x_j \leq t$ " or " $x_j > t$ ", where $t \in \mathbb{R}$ and $j \in \{1, 2\}$.

- ☐ Points: $\{(-5, 5), (5, 5)\}$, Labels: $\{0, 1\}$
- ☐ Points: $\{(0, -9), (0, 10)\}$, Labels: $\{0, 1\}$
- ☐ Points: $\{(3, 4), (-3, -4)\}$, Labels: $\{0, 1\}$
- ☐ Points: $\{(0, 0), (0, 1), (1, 0)\}$, Labels: $\{0, 0, 1\}$
- ☐ Points: $\{(0, 0), (0, 1), (1, 0)\}$, Labels: $\{0, 1, 1\}$
- ☐ Points: $\{(1, 2), (3, 4), (5, 6)\}$, Labels: $\{0, 1, 1\}$
- ☐ None of the above

3 k -Nearest Neighbors (20 points)

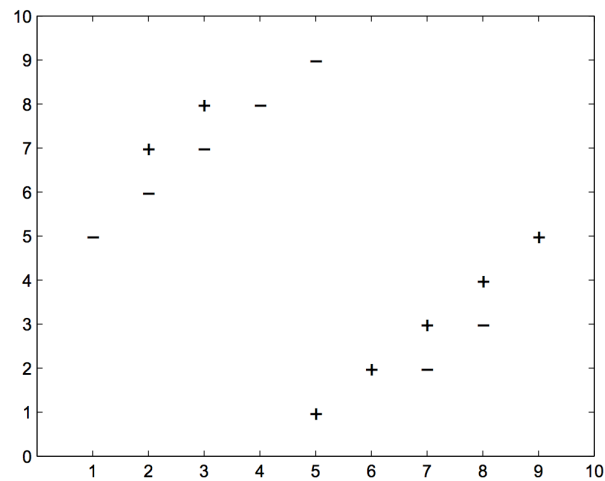


Figure 3: k-NN Dataset

1. Consider a k -nearest neighbors (k -NN) binary classifier which assigns the class of a test point to be the class of the majority of the k -nearest neighbors, according to the Euclidean distance metric. Assume that ties are broken by selecting one of the labels uniformly at random.
 - (a) (2 points) Using Figure 3 shown above to train the classifier and choosing $k = 6$, what is the classification error on the training set? **Round to 4 decimal places.**

Your answer:

- (b) (2 points) **Select all that apply:** Let's say that we have a new test point (not present in our training data) $\mathbf{x}^{\text{new}} = [3, 9]^T$ that we would like to apply our k -NN classifier to, as seen in figure 4.

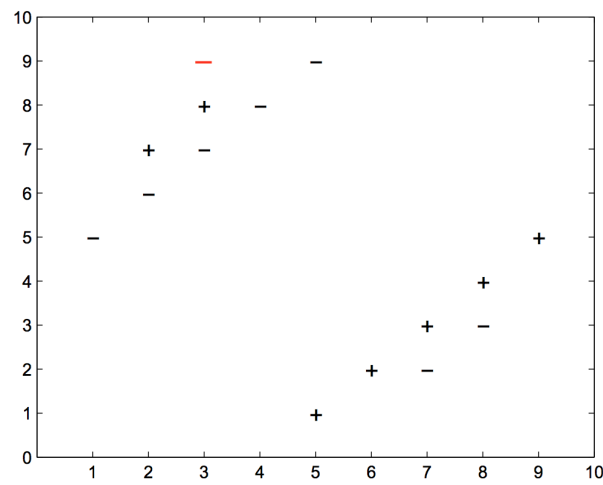


Figure 4: k-NN Dataset with Test Point

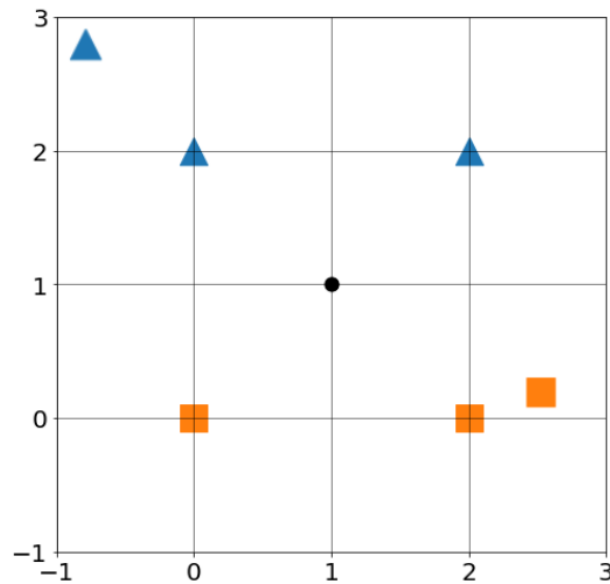
For which values of k is this test point correctly classified by the k -NN algorithm?

- ☐ $k = 1$
 - ☐ $k = 5$
 - ☐ $k = 9$
 - ☐ $k = 12$
 - ☐ None of the above
2. **Select one:** Assume we have a large labeled dataset that is randomly divided into a training set and a test set, and we would like to classify points in the test set using a k -NN classifier.
- (a) (1 point) In order to minimize the classification error on this test set, we should always choose the value of k which minimizes the training set error.
- ☐ True
 - ☐ False

- (b) (2 points) **Select one:** Instead of choosing the hyperparameters by merely minimizing the training set error, we instead consider splitting the training-all data set into a training and a validation data set, and choose the hyperparameters that lead to lower validation error. Is choosing hyperparameters based on validation error better than choosing hyper-parameters based on training error?
- ☐ Yes, lowering validation error instead of training error is better because lowering training error will not help generalize our model and may lead to overfitting.
 - ☐ Yes, lowering validation error is better for the model because cross-validation guarantees a better test error.
 - ☐ No, lowering training error instead of validation error is better because lowering validation error will not help generalize our model and may lead to overfitting.
 - ☐ No, lowering training error is better for the model because we have to learn the training set as well as possible to guarantee the best possible test error.
- (c) (2 points) **Select one:** Your friend Sally suggests that instead of splitting the original training set into separate training and validation sets, we should instead use the test set as the validation data for choosing hyperparameters. Is this a good idea? Justify your opinion with no more than 3 sentences.
- ☐ Yes
 - ☐ No

Your answer:

3. (2 points) **Select all that apply:** Consider a binary k -NN classifier where $k = 4$ and the two labels are “triangle” and “square”. Consider classifying a new point $\mathbf{x} = (1, 1)$, where two of the \mathbf{x} ’s nearest neighbors are labeled “triangle” and two are labeled “square” as shown below.



Which of the following methods can be used to break ties or avoid ties on this dataset?

- ☐ Assign \mathbf{x} the label of its nearest neighbor
 - ☐ Flip a coin to randomly assign a label to \mathbf{x} (from the labels of its 4 closest points)
 - ☐ Use $k = 3$ instead
 - ☐ Use $k = 5$ instead
 - ☐ None of the above.
4. (3 points) **Select all that apply:** Which of the following is/are correct statement(s) about k -NN models?
- ☐ A larger k tends to give a smoother decision boundary.
 - ☐ To reduce the impact of noise or outliers in our data, we should increase the value k .
 - ☐ If we make k too large, we could end up overfitting the data.
 - ☐ We can use cross-validation to help us select the value of k .
 - ☐ We should never select the k that minimizes the error on the validation dataset.
 - ☐ None of the above.

5. Consider the following data concerning the relationship between academic performance and salary after graduation. High school GPA and university GPA are two numerical features and salary is the numerical target. Note that salary is measured in thousands of dollars per year.

Student ID	High School GPA	University GPA	Salary
1	2.5	3.8	45
2	3.3	3.5	90
3	4.0	4.0	142
4	3.0	2.0	163
5	3.8	3.0	2600
6	3.3	2.8	67
7	3.9	3.8	unknown

- (a) (2 points) Among Students 1 to 6, who is the nearest neighbor to Student 7, using Euclidean distance? Answer the Student ID only.

Your answer:

- (b) (2 points) Now, our task is to predict the salary Student 7 earns after graduation. We apply k -NN to this regression problem: the prediction for the numerical target (salary in this example) is equal to the average of salaries for the top k nearest neighbors. If $k = 3$, what is our prediction for Student 7's salary? Be sure to use the same unit of measure (thousands of dollars per year) as the table above.

Round your answer to the nearest integer.

Your answer:

- (c) (2 points) **Select all that apply:** Suppose that the first 6 students shown above are only a subset of your full training data set, which consists of 10,000 students. We apply k -NN regression using Euclidean distance to this problem and we define training loss on this full data set to be the mean squared error (MSE) of salary. Now consider the possible consequences of modifying the data in various ways. Which of the following changes **could** have an effect on training loss on the full data set as measured by mean squared error (MSE) of salary?

- ☐ Rescaling only "High School GPA" to be a percentage of 4.0
- ☐ Rescaling only "University GPA" to be a percentage of 4.0
- ☐ Rescaling both "High School GPA" and "University GPA", so that each is a percentage of 4.0
- ☐ None of the above.

4 Perceptron (19 points)

1. (1 point) **Select one:** Consider running the online perceptron algorithm on some sequence of examples S (an example is a data point and its label). Let S' be the same set of examples as S , but presented in a different order.

True or False: The online perceptron algorithm is guaranteed to make the same number of mistakes on S as it does on S' .

- ☐ True
☐ False

2. (3 points) **Select all that apply:** Suppose we have a perceptron whose inputs are 2-dimensional vectors and each feature vector component is either -1 or 1, i.e., $x_i \in \{-1, 1\}$. The prediction function is $y = \text{sign}(w_1x_1 + w_2x_2 + b)$, and

$$\text{sign}(z) = \begin{cases} 1, & \text{if } z \geq 0 \\ -1, & \text{otherwise.} \end{cases}$$

Which of the following functions can be implemented with the above perceptron? That is, for which of the following functions does there exist a set of parameters w, b that correctly define the function.

- ☐ AND function, i.e., the function that evaluates to 1 if and only if all inputs are 1, and -1 otherwise.
☐ OR function, i.e., the function that evaluates to 1 if and only if at least one of the inputs are 1, and -1 otherwise.
☐ XOR function, i.e., the function that evaluates to 1 if and only if the inputs are not all the same. For example

$$\text{XOR}(1, -1) = 1, \text{ but } \text{XOR}(1, 1) = -1.$$

- ☐ None of the above.

3. (2 points) **Select one:** Suppose we have a dataset $\{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$, where $\mathbf{x}^{(i)} \in \mathbb{R}^M$, $y^{(i)} \in \{+1, -1\}$. We would like to apply the perceptron algorithm on this dataset. Assume there is no intercept term. How many parameter values is the perceptron algorithm learning?

- ☐ N
☐ $N \times M$
☐ M

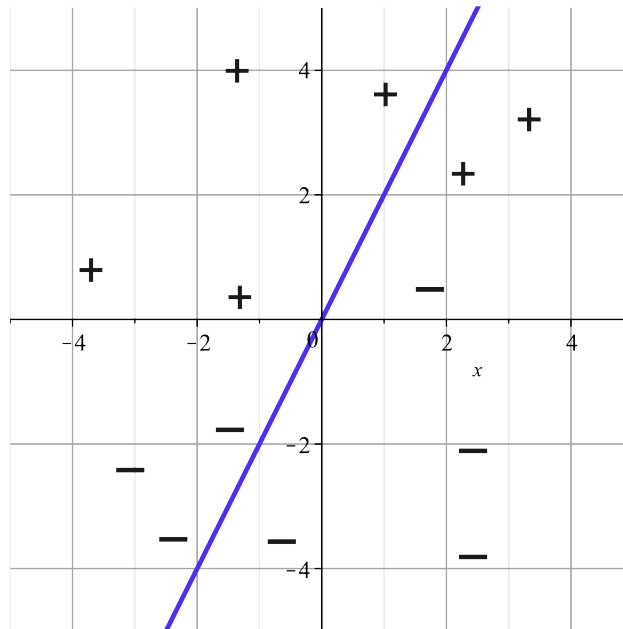
4. (2 points) **Select one:** The following table shows a data set and the number of times each point is misclassified during a run of the perceptron algorithm. What is the separating plane θ found by the algorithm, i.e. $\theta = [b, \theta_1, \theta_2, \theta_3]$? Assume that the initial weights are all zero.

x_1	x_2	x_3	y	Times Misclassified
2	1	5	1	10
5	3	3	1	5
1	6	2	1	8
7	2	1	-1	2
3	2	6	-1	3

- ☐ [18, 25, 14, 34]
☐ [18, 30, 63, 61]
☐ [16, 56, 18, 47]
☐ [18, 52, 19, 47]
5. (2 points) **Select all that apply:** Which of the following is/are correct statement(s) about the mistake bound of the perceptron algorithm?
- ☐ If the minimum distance from any data point to the separating hyperplane is increased, without any other change to the data points, the mistake bound will also increase.
☐ If the whole dataset is shifted away from origin, then the mistake bound will also increase.
☐ If the pair-wise distance between data points is increased, i.e. the data is scaled by some constant value, then the mistake bound will also increase.
☐ The mistake bound is linearly inverse-proportional to the minimum distance of any data point to the separating hyperplane of the data.
☐ None of the above.
6. (2 points) **Select one:** Suppose we have data whose examples are of the form $[x_1, x_2]$, where $x_1 - x_2 = 0$. We do not know the label for each element. Suppose the perceptron algorithm starts with $\theta = [3, 5]$; which of the following values will θ never take on in the process of running the perceptron algorithm on the data?
- ☐ $[-1, 1]$
☐ $[4, 6]$
☐ $[-3, 0]$
☐ $[-6, -4]$

7. (2 points) **Select all that apply:** Consider the linear decision boundary below and the dataset shown. Which of the following are valid weights θ and their corresponding error on this dataset? (Note: Assume the decision boundary is fixed and does not change while evaluating error.)

- ☐ $\theta = [-2, 1]$, error = 5/13
- ☐ $\theta = [2, -1]$, error = 8/13
- ☐ $\theta = [2, -1]$, error = 5/13
- ☐ $\theta = [-2, 1]$, error = 8/13
- ☐ None of the above.



8. The following problem will walk you through an application of the Perceptron Mistake Bound. The following table shows a linearly separable dataset, and your task will be to determine the mistake bound for the dataset.

NOTE: The proof of the perceptron mistake bound requires that the optimal linear separator passes through the origin. To make the linear separator pass through the origin, we fold the bias into the weights and prepend a 1 to each training example's input. The original data is on the left, and the result of this prepending is shown on the right. **Be sure to use the modified dataset on the right in your calculations.**

x_1	x_2	y
-2	2	1
-1	-3	-1
-2	-3	-1
0	1	1
2	-1	1

x_0	x_1	x_2	y
1	-2	2	1
1	-1	-3	-1
1	-2	-3	-1
1	0	1	1
1	2	-1	1

- (a) (2 points) Compute the radius R of the “circle” centered at the origin that bounds the data points.
Round to 4 decimal places after the decimal point.

Radius:

- (b) (2 points) Assume that the linear separator with the largest margin is given by

$$\theta^{*T} \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} = 0, \text{ where } \theta^* = \begin{bmatrix} 6 \\ 3 \\ 4 \end{bmatrix}$$

Now, compute the margin of the dataset.

Round to 4 decimal places after the decimal point.

Margin:

- (c) (1 point) Based on the above values, what is the theoretical perceptron mistake bound for this dataset, given this linear separator?

Round to 4 decimal places after the decimal point.

Mistake Bound:

5 Linear Regression (19 points)

1. Consider the following dataset:

x	9.0	2.0	6.0	1.0	8.0
y	1.0	0.0	3.0	0.0	1.0

Let \mathbf{x} be the vector of datapoints and \mathbf{y} be the label vector. Here, we are fitting the data using gradient descent, and our objective function is $\frac{1}{N} \sum_{i=1}^N (wx_i + b - y_i)^2$ where N is the number of data points, w is the weight, and b is the intercept.

Note: Showing your work in these questions is optional, but it is recommended to help us understand where any misconceptions may occur. We may give partial credit for correct work if your answer is incorrect.

- (a) (2 points) If we initialize the weight as 3.0 and intercept as 0.0, what is the gradient of the loss function with respect to the weight w , calculated over all the data points, in the first step of the gradient descent update?

Round to 4 decimal places after the decimal point.

Gradient:

Work

- (b) (2 points) What is the gradient of the loss function with respect to the intercept b , calculated over all the data points, in the first step of the gradient descent update?

Gradient:

Work

- (c) (2 points) Let the learning rate be 0.01. Perform one step of gradient descent on the data. Fill in the following blanks with the value of the weight and the value of the intercept after this step. **Round to 4 decimal places after the decimal point.**

Weight:

Intercept:

2. Consider a dataset $\mathcal{D}_1 = \{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$. Assume the linear regression model that minimizes the mean-squared error on \mathcal{D}_1 is $y = w_1x + b_1$.

(a) (2 points) **Select one:** Now, suppose we have the dataset $\mathcal{D}_2 = \{(x^{(1)} + \alpha, y^{(1)} + \beta), \dots, (x^{(N)} + \alpha, y^{(N)} + \beta)\}$ where $\alpha > 0, \beta > 0$ and $w_1\alpha \neq \beta$. Assume the linear regression model that minimizes the mean-squared error on \mathcal{D}_2 is $y = w_2x + b_2$. Select the correct statement about w_1, w_2, b_1, b_2 below. Note that the statement should hold no matter what values α, β take on within the specified constraints.

- ☐ $w_1 = w_2, b_1 = b_2$
- ☐ $w_1 \neq w_2, b_1 = b_2$
- ☐ $w_1 = w_2, b_1 \neq b_2$
- ☐ $w_1 \neq w_2, b_1 \neq b_2$

- (b) (2 points) We decide to ask a friend to analyze \mathcal{D}_1 ; however, he makes a mistake by duplicating a subset of the rows in \mathcal{D}_1 . Explain why the linear regression parameters that minimize mean-squared error on the duplicated data *may* differ from the parameters learned on \mathcal{D}_1 , i.e. w_1 and b_1 .

Your answer:

3. We wish to learn a linear regression model on the dataset $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$ where $\mathbf{x} \in \mathbb{R}^k$. For this question, define the *log-cosh loss* as

$$\ell(\hat{y}, y) = \log(\cosh(\hat{y} - y))$$

In particular, for a given point $\mathbf{x}^{(i)}$, the log-cosh loss of a model with parameters $\boldsymbol{\theta}$ is

$$J^{(i)}(\boldsymbol{\theta}) = \log\left(\cosh\left(\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}\right)\right)$$

We are interested in minimizing loss over our training data, so we minimize the average log-cosh loss over all points in \mathcal{D} .

- (a) (2 points) What is the objective $J(\boldsymbol{\theta})$ in this setting?

Your answer:

- (b) (3 points) What is the partial derivative of $J^{(i)}(\boldsymbol{\theta})$ with respect to the j^{th} parameter, θ_j ? It may be helpful to know that $\frac{d}{dx} \cosh(x) = \sinh(x)$.

Your answer:

- (c) (2 points) What is the gradient of $J^{(i)}(\boldsymbol{\theta})$ with respect to the entire parameter vector $\boldsymbol{\theta}$?

Your answer:

- (d) (2 points) To find the optimal parameter vector $\boldsymbol{\theta}^*$ that minimizes $J(\boldsymbol{\theta})$, we again decide to use gradient descent. Write pseudocode that performs gradient descent for one iteration. Set learning rate $\alpha = 0.1$ and initialize $\boldsymbol{\theta}$ to be the zero vector. You may use `gradient[i]` as a variable that contains your answer to part (c) in your pseudocode. Limit your answer to 10 lines.

Your Answer

6 Collaboration Questions

After you have completed all other components of this assignment, report your answers to these questions regarding the collaboration policy. Details of the policy can be found [here](#).

1. Did you receive any help whatsoever from anyone in solving this assignment? If so, include full details.
2. Did you give any help whatsoever to anyone in solving this assignment? If so, include full details.
3. Did you find or come across code that implements any part of this assignment? If so, include full details.

Your Answer