

CSC 311 Fall 2022 Midterm B

Friday, October 28, 2022

Q1 Short answers [10 pts]

Q1a. (True/False) Decision Trees are a non-parametric method. [1pt]

Q1b. (True/False) When training a regression model, we should continue to optimize our model until the averaged training loss is exactly zero. [1pt]

Q1c. (True/False) Early stopping is used to prevent underfitting. [1pt]

Q1d. (True/False) A two-layer neural network with 2^D hidden units and a hard threshold activation function can approximate any boolean function over D input variables. [1pt]

Q1e. You have trained a logistic regression model on a very large dataset up to an accuracy of 77%.(a) Your colleague claims that by using a neural network you can outperform the Bayes Optimal classifier – are they correct? (Yes/No with short justification) (b) You now train a neural network to do binary classification via the backpropagation algorithm on the dataset - do you expect your accuracy to increase/stay the same or decrease? Provide a short justification. [2pt]

Q1f. Consider the following dataset representing the XNOR boolean function.

x_1	x_2	y
0	0	1
0	1	0
1	0	0
1	1	1

Is the data linearly separable? Briefly justify your choice on whether the data is (or is not) linearly separable. [2pts]

Q1g. Why can a neural network with the hard-threshold activation function correctly classify the XOR function as we saw in lecture? Your answer should refer to the hypothesis space of functions represented by each model. Highlight one technical challenge that can arise when using back-propagation to learn neural networks with a hard threshold activation function. [2pts]

Q2 [10pts] Gradient Descent

Consider a data set with three examples. For example i , $x^{(i)}$ is the value of the input feature and $t^{(i)}$ is the value of the target.

Example	$x^{(i)}$	$t^{(i)}$
1	1	0
2	2	4
3	3	-1

Table 1: Data Set for Linear Regression

We will fit a ridge regression model to the data above resulting in on weight w and one bias term b . Recall that in vector notation (with the bias term incorporated into the design matrix via concatenation into a vector of all 1s), this corresponds to training with:

$$\mathcal{J}(\mathbf{w}) = \frac{1}{2} \|X^T \mathbf{w} - \mathbf{t}\|^2 + \lambda \frac{1}{2} \|\mathbf{w}\|^2$$

Derive the gradient descent update rule using the parameter values below.

$$w = 2, b = -1, \alpha = 0.2, \lambda = 0.5$$

Assume that the cost function is the total loss over all the training examples. **We recommend breaking up the calculations into easy to understand steps rather than only writing down the final answer. We will award partial marks for correct intermediate steps.** You should leave your answer in the form of two equations:

$$b \leftarrow Ab + C \tag{1}$$

$$w \leftarrow Cw + D \tag{2}$$

where A, B, C, D are real valued numbers.

Q2 Answer

Q3 Feature maps and separability [5pts]

Q3a. [1pts] In lecture, we discussed polynomial regression in one variable. It is also possible to fit polynomials of two variables. A degree-2 polynomial in two features x_1 and x_2 is a sum of monomials which are at most quadratic in x_1 and x_2 . One such example is: $x_1^2 - 3x_1x_2 + 2x_2 + 12$. Design a feature map $\phi(u, v)$ that allows you to use a linear regression solver to fit *any* degree-2 polynomials. No justification is required.

Q3b. [4pts] Figure 1 presents a visual depiction of a binary dataset.

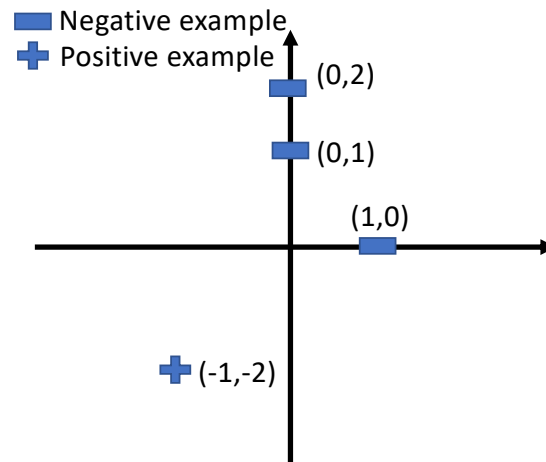


Figure 1: Binary Dataset

Recall the binary linear classification model with a decision rule:

$$z = w^T x$$
$$y = \begin{cases} 1, & \text{if } z \geq 0 \\ 0, & \text{if } z < 0 \end{cases}$$

Assume that the bias term is zero ($b = w_0 = 0$).

Q3b Answer

1. Write out the inequalities that represent the constraints that the weight space needs to satisfy. [3pts]
2. If the problem is feasible, provide one assignment of the weights. [1pts]

Q4 [15pts] Decision Trees

Consider the data set in Table 2. There are 11 examples. There are 2 binary discrete features: colour and length. “Colour” has two values: dark and light. “Length” has two values: long and short. Each example has a binary label: True or False.

Example	Colour	Length	Label
1	Dark	Long	True
2	Dark	Long	True
3	Dark	Long	False
4	Dark	Short	True
5	Dark	Short	True
6	Dark	Short	True
7	Light	Short	False
8	Light	Short	False
9	Light	Short	False
10	Light	Short	False
11	Light	Short	True

Table 2: Data Set for Decision Tree

Q4a. [9pts] Complete the following decision tree for the data set.

- For each node, write down the number of true and false examples.
- Then, for each leaf (rectangle-shaped) node, write down the label/decision.

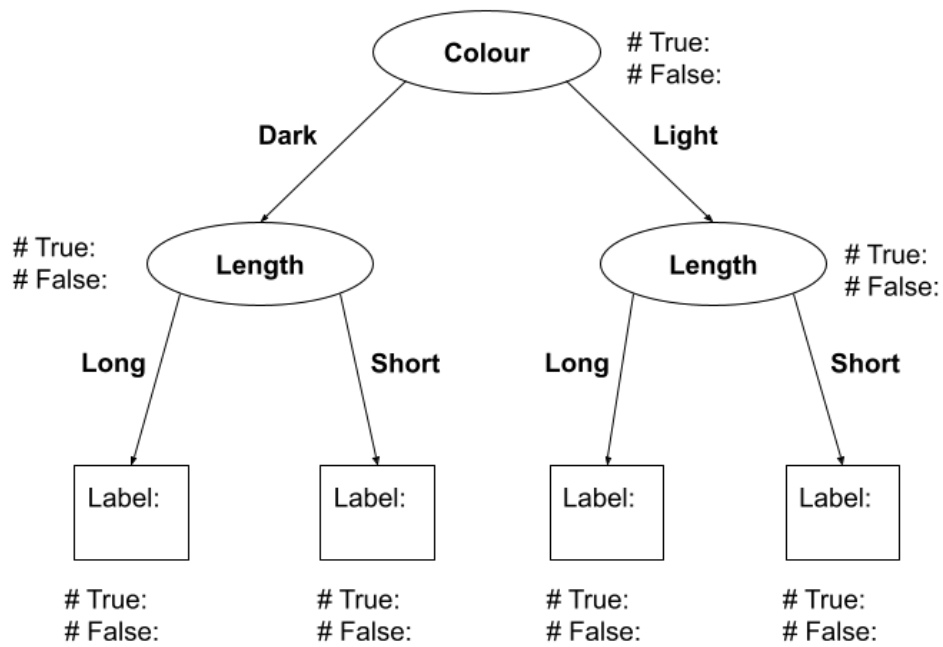


Figure 2: Complete this decision tree

Q4b. [6pts] Write down the formulas for calculating the expected information gain of testing Length at the root. H denotes entropy. You do not need to calculate the result and can leave numbers as fractions. i.e. for each $\frac{\boxed{}}{\boxed{}}$ you need to enter one number in the numerator and another in the denominator.

(2 pts) The entropy before testing Length

$$= H\left(\frac{\boxed{}}{\boxed{}}, \frac{\boxed{}}{\boxed{}}\right)$$

$$= -\frac{\boxed{}}{\boxed{}} \log_2 \frac{\boxed{}}{\boxed{}} - \frac{\boxed{}}{\boxed{}} \log_2 \frac{\boxed{}}{\boxed{}}$$

(4 pts) The expected conditional entropy after splitting on Length at the root

$$\begin{aligned}
 &= \frac{\boxed{}}{\boxed{}} H\left(\frac{\boxed{}}{\boxed{}}, \frac{\boxed{}}{\boxed{}}\right) \\
 &+ \frac{\boxed{}}{\boxed{}} H\left(\frac{\boxed{}}{\boxed{}}, \frac{\boxed{}}{\boxed{}}\right) \\
 &= \frac{\boxed{}}{\boxed{}} \left(-\frac{\boxed{}}{\boxed{}} \log_2 \frac{\boxed{}}{\boxed{}} - \frac{\boxed{}}{\boxed{}} \log_2 \frac{\boxed{}}{\boxed{}} \right) \\
 &+ \frac{\boxed{}}{\boxed{}} \left(-\frac{\boxed{}}{\boxed{}} \log_2 \frac{\boxed{}}{\boxed{}} - \frac{\boxed{}}{\boxed{}} \log_2 \frac{\boxed{}}{\boxed{}} \right)
 \end{aligned}$$

This page intentionally left blank