

CPSC 340 Midterm (Fall 2015)

Name:

Student Number:

Please enter your information above, turn off cellphones, space yourselves out throughout the room, and wait until the official start of the exam to begin. You can raise your hand to ask a question, and please look up occasionally in case there are clarifications written on the projector (the time will also be written on the projector). You are welcome to (quietly) leave early if you finish early, and we will ask everyone to stop at 3:55.

The midterm consists of 5 questions, and they will all be equally weighted in the marking scheme. Note that some question have multiple parts, written as **(a)**, **(b)**, and in some cases **(c)**. All parts are equally weighted. Clearly mark where you are answering each part, and make sure to check that you answered all parts before handing in your midterm.

Good luck!

2 KNN and Decision Stumps

Consider the dataset below, which has 10 training examples and 2 features:

$$X = \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 1 \\ 1 & 0 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}.$$

Suppose that you want to classify the following test example:

$$\hat{x} = [1 \quad 1].$$

- (a) What class label would we assign to the test example if we used a k -nearest neighbours classifier, with $k = 3$ and the Euclidean distance measure?
- (b) Suppose we want to fit a decision stump to this dataset. What is the decision rule (based on a single variable) that minimizes the classification error? (Show your reasoning for the two possible splits.)
- (c) Under the decision rule you estimated from part (b), what is the most likely label for the test example?

3 Less-Naive Bayes

In class, we talked about naive Bayes models. Given a set of features $\{x_1, x_2, x_3, \dots, x_d\}$, naive Bayes approximates the conditional probability of a label y using

$$\begin{aligned} p(y|x_1, x_2, x_3, \dots, x_d) &\propto p(y)p(x_1, x_2, x_3, \dots, x_d|y) \\ &= p(y)p(x_1|y)p(x_2|x_1, y)p(x_3|x_1, x_2, y) \dots p(x_d|x_1, x_2, x_3, \dots, x_{d-1}, y) \\ &\approx p(y)p(x_1|y)p(x_2|y)p(x_3|y) \dots p(x_d|y). \end{aligned}$$

Naive Bayes makes a strong independence assumption (“ \approx ”), and there are various ways to relax it. For example, consider a ‘less-naive’ Bayes model that depends on a parameter k and assumes that x_j given y is independent of all variables *except* the up to k largest values j where $j < i$. For example, if $k = 3$ then x_6 is conditionally independent of all other variables given y (as in the usual naive Bayes model) *and* given x_3, x_4 , and x_5 : $p(x_6|x_1, x_2, x_3, \dots, x_d, y) = p(x_6|x_3, x_4, x_5, y)$.

Naive Bayes corresponds to $k = 0$, and we make weaker assumptions as k grows. As another example, if $k = 2$ then we use

$$\begin{aligned} p(y|x_1, x_2, x_3, \dots, x_d) &\propto p(y)p(x_1, x_2, x_3, \dots, x_d|y) \\ &= p(y)p(x_1|y)p(x_2|x_1, y)p(x_3|x_1, x_2, y) \dots p(x_d|x_1, x_2, x_3, \dots, x_{d-1}, y) \\ &\approx p(y)p(x_1|y)p(x_2|x_1, y)p(x_3|x_2, x_1, y) \dots p(x_d|x_{d-1}, x_{d-2}, y). \end{aligned}$$

Instead of simply estimating conditionals like $p(x_5 = 1|y = 1)$, this will now involve estimating conditionals like $p(x_5 = 1|x_4 = 1, x_3 = 0, y = 1)$.

(a) What are the two parts of the fundamental trade-off in machine learning?

(b) For the less-naive Bayes model, how would the choice of k affect the two parts of the fundamental trade-off?

4 Runtime of K-Means

One of the outputs of the k-means algorithm is a set of cluster means μ_c . To find the cluster of a new data point \hat{x} , we can perform the following:

- For each cluster c , compute the Euclidean distance between \hat{x} and the cluster mean μ_c .
- Assign \hat{x} to the cluster c with the minimum distance.

Following our usual convention, we'll use:

1. d as the length of an \hat{x} and μ_c .
2. k as the number of clusters.
3. t as the number of test examples.

If we use the above two steps to find the cluster of t test examples, what is the total cost in terms of d , k , and t ?

5 Regularized Linear Regression in 1D

Consider the problem of performing linear regression in 1-dimension where:

- We use the squared error as our loss function.
- We use an ℓ_2 -regularizer with weight λ .

This gives us the following objective:

$$\operatorname{argmin}_w \frac{1}{2} \sum_{i=1}^n [(y_i - wx_i)^2] + \frac{\lambda}{2} w^2.$$

- (a) Compute the derivative of this objective function with respect to w .
- (b) By equating the derivative of this objective (which is convex and quadratic) with 0, compute the solution of this problem in terms of the x_i , y_i , and λ (show your work).

