

# CPSC 340 Assignment 2 (due Friday September 27 at 11:55pm)

The assignment instructions are the same as Assignment 1, except you have the option to work in a group of 2. If you work in a group, please only hand in *one* assignment. It is recommended that you work in groups as the assignment is quite long, but please only submit one assignment for the group and make sure that everyone's name/ID is on the front page.

## 1 Training and Testing

### 1.1 Training Error

Running `example_train.jl` fits decision trees of different depths using two different implementations: the “decisionStump” function from Assignment 1, and using a variant using a more sophisticated splitting criterion called the information gain. Describe what you observe. Can you explain the results?

Answer: The result shows the accuracy based decision tree have an lower bound on the training error of 0.11, but the information gain based decision tree can successfully classify all the sample in the training set (training error = 0).

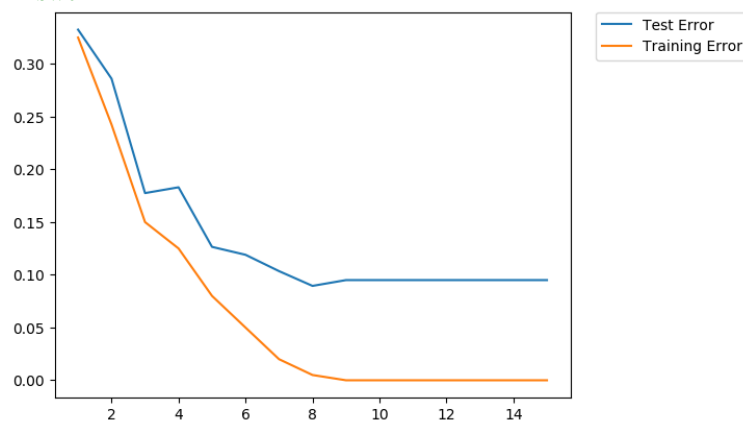
However, if the depth of the tree is shallow (depth  $\leq 4$  here), the training error of the information gain based tree have a higher training error.

### 1.2 Training and Testing Error Curves

Notice that the `citiesSmall.mat` file also contains test data, “Xtest” and “ytest”. Running `example_trainTest` trains a depth-2 decision tree and evaluates its performance on the test data. With a depth-2 decision tree, the training and test error are fairly close, so the model hasn't overfit much.

Make a plot that contains the training error and testing error as you vary the depth from 1 through 15. How do each of these errors change with the decision tree depth?

Answer:



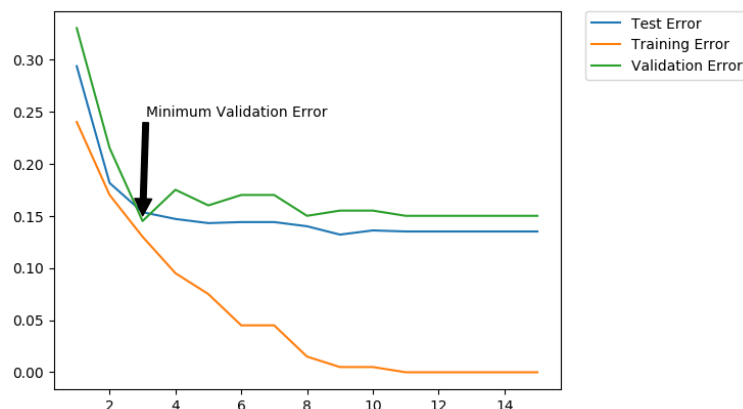
As the depth of the tree increased, the Test Error decreased first, and reached minimum at depth 8. The training error keeps decreasing after that but the test error increased a bit and settled at 0.095

Note: use the provided infogain-based decision tree code from the previous subsection.

### 1.3 Validation Set

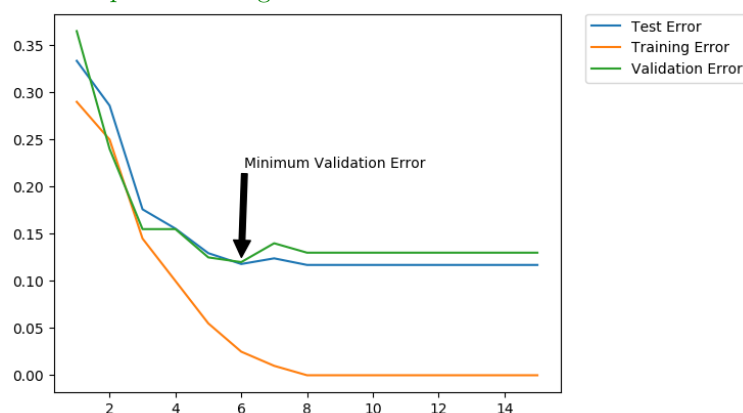
Suppose we're in the typical case where we don't have the labels for the test data. In this case, we might instead use a *validation* set. Split the training set into two equal-sized parts: use the first  $n/2$  examples as a training set and the second  $n/2$  examples as a validation set (we're assuming that the examples are already in a random order). What depth of decision tree would we pick if we minimized the validation set error? Does the answer change if you switch the training and validation set? How could we use more of our data to estimate the depth more reliably?

Answer:



The depth of decision tree is 3.

If we swap the training set and validation set:



The we would pick the depth as 6. The way to use the data is to use k-fold method to test on different validation sets.

Note: use the provided infogain-based decision tree code from the previous subsection.

## 2 Naive Bayes

In this section we'll implement naive Bayes, a very fast classification method that is often surprisingly accurate for text data with simple representations like bag of words.

## 2.1 Naive Bayes by Hand

Consider the dataset below, which has 10 training examples and 3 features:

$$X = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}, \quad y = \begin{bmatrix} \text{spam} \\ \text{spam} \\ \text{spam} \\ \text{spam} \\ \text{spam} \\ \text{spam} \\ \text{not spam} \\ \text{not spam} \\ \text{not spam} \\ \text{not spam} \end{bmatrix}.$$

The feature in the first column is <your name> (whether the e-mail contained your name), in the second column is “pharmaceutical” (whether the e-mail contained this word), and the third column is “PayPal” (whether the e-mail contained this word). Suppose you believe that a naive Bayes model would be appropriate for this dataset, and you want to classify the following test example:

$$\hat{x} = [1 \quad 1 \quad 0].$$

### 2.1.1 Prior probabilities

Compute the estimates of the class prior probabilities (you don’t need to show any work):

- $p(\text{spam})$

Answer: 0.6.

- $p(\text{not spam})$

Answer: 0.4.

### 2.1.2 Conditional probabilities

Compute the estimates of the 6 conditional probabilities required by naive Bayes for this example (you don’t need to show any work):

- $p(<\text{your name}> = 1 \mid \text{spam})$ .

Answer: 1/6

- $p(\text{pharmaceutical} = 1 \mid \text{spam})$ .

Answer: 5/6

- $p(\text{PayPal} = 0 \mid \text{spam})$ .

Answer:  $4/6 = 2/3$

- $p(<\text{your name}> = 1 \mid \text{not spam})$ .

Answer: 1

- $p(\text{pharmaceutical} = 1 \mid \text{not spam})$ .

Answer: 1/4

- $p(\text{PayPal} = 0 \mid \text{not spam})$ .

Answer: 1/4

### 2.1.3 Prediction

Under the naive Bayes model and your estimates of the above probabilities, what is the most likely label for the test example? (Show your work.)

Answer:

$$\begin{aligned} & P(\text{spam} \mid \langle \text{your name} \rangle = 1 \& \text{pharmaceutical} = 1 \& \text{PayPal} = 0) \\ & \approx \frac{p(\langle \text{your name} \rangle = 1 \mid \text{spam}) \times p(\text{pharmaceutical} = 1 \mid \text{spam}) \times p(\text{PayPal} = 0 \mid \text{spam}) \times P(\text{spam})}{C} \\ & = \frac{0.6 * 1/6 * 5/6 * 2/3}{C} \\ & = \frac{1/18}{C} \end{aligned}$$

$$\begin{aligned} & P(\text{not spam} \mid \langle \text{your name} \rangle = 1 \& \text{pharmaceutical} = 1 \& \text{PayPal} = 0) \\ & \approx \frac{p(\langle \text{your name} \rangle = 1 \mid \text{not spam}) \times p(\text{pharmaceutical} = 1 \mid \text{not spam}) \times p(\text{PayPal} = 0 \mid \text{not spam}) \times P(\text{not spam})}{C} \\ & = \frac{0.4 * 1 * 1/4 * 1/4}{C} \\ & = \frac{1/40}{C} \end{aligned}$$

$\frac{P(\text{spam} \mid \langle \text{your name} \rangle = 1 \& \text{pharmaceutical} = 1 \& \text{PayPal} = 0)}{P(\text{not spam} \mid \langle \text{your name} \rangle = 1 \& \text{pharmaceutical} = 1 \& \text{PayPal} = 0)} = \frac{20}{9} > 1$   
Therefore, the email is more likely to be a spam.

## 2.2 Bag of Words

If you run the script *example\_bagOfWods.jl*, it will load the following dataset:

1. *X*: A sparse binary matrix. Each row corresponds to a newsgroup post, and each column corresponds to whether a particular word was used in the post. A value of 1 means that the word occurred in the post.
2. *wordlist*: The set of words that correspond to each column.
3. *y*: A vector with values 1 through 4, with the value corresponding to the newsgroup that the post came from.
4. *groupnames*: The names of the four newsgroups.
5. *Xvalid* and *yvalid*: the word lists and newsgroup labels for additional newsgroup posts.

Answer the following:

1. Which word is present in the newsgroup post if there is a 1 in column 30 of *X*?

Answer: format

2. Which words are present in training example 200?

Answer: food, world

3. Which newsgroup name does training example 200 come from?

Answer: talk.\*

## 2.3 Naive Bayes Implementation

If you run the function `example_decisionTree_newsgroups.jl` it will load the newsgroups dataset and report the test error for decision trees of different sizes (it may take a while for the deeper trees, as this is a sub-optimal implementation). On the other hand, if you run the function `example_naiveBayes.jl` it will fit the basic naive Bayes model and report the test error.

While the `predict` function of the naive Bayes classifier is already implemented, the calculation of the variable `p_xy` is incorrect (right now, it just sets all values to  $1/2$ ). [Modify this function so that `p\_xy` correctly computes the conditional probabilities of these values based on the frequencies in the data set.](#) Hand in your code and report the test error that you obtain.

Answer:

```
# We will store  $p(x(i,j) = 1 \mid y(i) = c)$  in p_xy(1,j,c)
# We will store  $p(x(i,j) = 0 \mid y(i) = c)$  in p_xy(2,j,c)
p_xy = zeros(2,d,k)
for i in 1:n
    for j in 1:d
        if X[i,j] == 1
            p_xy[1,j,y[i]] += 1
        else
            p_xy[2,j,y[i]] += 1
        end
    end
end

for j in 1:d
    for i in 1:k
        p_xy[1,j,i] = p_xy[1,j,i]/counts[i]
        p_xy[2,j,i] = p_xy[2,j,i]/counts[i]
    end
end
```

And the test error is 0.188

## 2.4 Runtime of Naive Bayes for Discrete Data

Assume you have the following setup:

- The training set has  $n$  objects each with  $d$  features.
- The test set has  $t$  objects with  $d$  features.
- Each feature can have up to  $c$  discrete values (you can assume  $c \leq n$ ).
- There are  $k$  class labels (you can assume  $k \leq n$ )

You can implement the training phase of a naive Bayes classifier in this setup in  $O(nd)$ , since you only need to do a constant amount of work for each  $X[i,j]$  value. (You do not have to actually implement it in this way for the previous question, but you should think about how this could be done). [What is the cost of classifying  \$t\$  test examples with the model?](#)

Answer: For each test sample, we need to calculate the factor of each  $y$ . To do that we need to multiply the probability of all the features in that sample, which is total of  $d$  features. That is we need  $O(kd)$  for one sample. Therefore, the total cost will be  $O(kdt)$ .

## 3 K-Nearest Neighbours

In *citiesSmall* dataset, nearby points tend to receive the same class label because they are part of the same state. This indicates that a  $k$ -nearest neighbours classifier might be a better choice than a decision tree (while naive Bayes would probably work poorly on this dataset). The file *knn.jl* has implemented the training function for a  $k$ -nearest neighbour classifier (which is to just memorize the data) but the predict function always just predicts 1.

### 3.1 KNN Prediction

Fill in the *predict* function in *knn.jl* so that the model file implements the  $k$ -nearest neighbour prediction rule. You should use Euclidean distance.

Hint: although it is not necessary, you may find it useful to pre-compute all the distances (using the *distancesSquared* function in *misc.jl*) and to use the *sortperm* command.

1. Hand in the predict function.

Answer:

```
function knn_predict(Xhat,X,y,k)
    (n,d) = size(X)
    (t,d) = size(Xhat)
    k = min(n,k) # To save you some debuggin
    dist = zeros(t,n)
    dist = transpose(distancesSquared(X , Xhat))
    minimumDistance = zeros(t,k)
    for i in 1:t
        minimumDistance[i,:] = y[sortperm(dist[i,:])[1:k]]
    end
    minY = zeros(t)
    for i in 1:t
        minY[i] = mode(minimumDistance[i,:]);
    end
    return minY
end
```

2. Report the training and test error obtained on the *citiesSmall.mat* dataset for  $k = 1$ ,  $k = 3$ , and  $k = 10$ . (You can use *example\_knn.jl* to get started.)

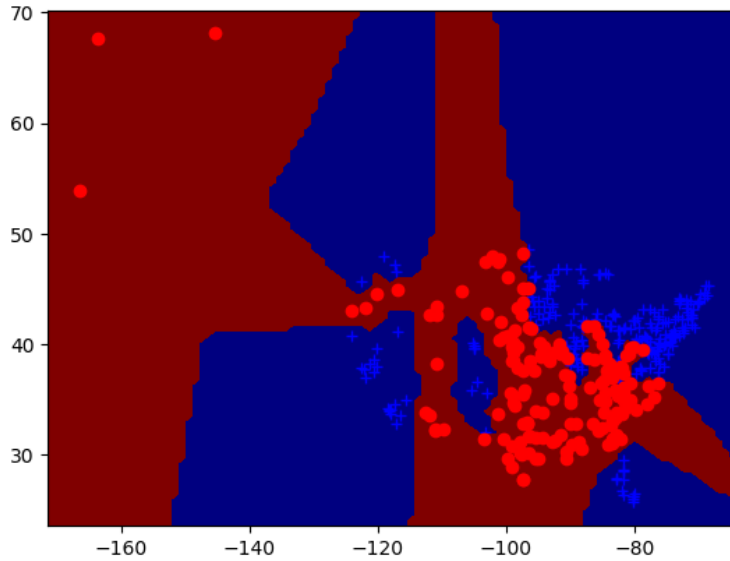
Answer:  $k = 1$ : Training Error : 0, Test Error: 0.065

$k = 3$ : Training Error : 0.028, Test Error: 0.066

$k = 10$ : Training Error : 0.072, Test Error: 0.097

3. Hand in the plot generated by classifier2Dplot on the *citiesSmall.mat* dataset for  $k = 1$  on the training data.
4. Why is the training error 0 for  $k = 1$ ?

Answer:



5. If you didn't have an explicit test set, how would you choose  $k$ ?

Answer: Use part of the training set as validation set.

Hint: when writing a function, it is typically a good practice to write one step of the code at a time and check if the code matches the output you expect. You can then proceed to the next step and at each step test is if the function behaves as you expect. You can also use a set of inputs where you know what the output should be in order to help you find any bugs. These are standard programming practices: it is not the job of the TAs or the instructor to find the location of a bug in a long program you've written without verifying the individual parts on their own.

### 3.2 Condensed Nearest Neighbours

The file *citiesBig1.mat* contains a version of this dataset with more than 30 times as many cities. KNN can obtain a lower test error if it's trained on this dataset, but the prediction time will be very slow. A common strategy for applying KNN to huge datasets is called *condensed nearest neighbours*, and the main idea is to only store a *subset* of the training examples (and to only compare to these examples when making predictions). A simple variation of this algorithm would be:

```
initialize subset with first training example;
for each training example do
    if the example is incorrectly classified by the KNN classifier using the current subset then
        add the current example to the subset;
    else
        do not add the current example to the subset (do nothing);
    end
end
```

#### Algorithm 1: Condensed Nearest Neighbours

You are provided with an implementation of this *condensed nearest neighbours* algorithm in *knn.jl*.

1. The point of this algorithm is to be faster than KNN. Try running the condensed nearest neighbours (called "cknn" in the code) on the *citiesBig1* dataset and report how long it takes to make a prediction. What about if you try to use KNN for this dataset?

2. Report the training and testing errors for condensed NN, as well as the number of training examples in the subset, on the *citiesBig1* dataset with  $k = 1$ .
3. Why is the training error with  $k = 1$  now greater than 0?
4. If we entered the coordinates of Vancouver into the predict function, would it be predicted to be in a blue state or a red state?
5. If you have  $s$  examples in the subset, what is the cost of running the predict function on  $t$  test examples in terms of  $n$ ,  $d$ ,  $t$ , and  $s$ ?
6. Try out your function on the dataset *citiesBig2*. Why are the test error *and* training error so high (even for  $k = 1$ ) for this method on this dataset?

## 4 Random Forests

### 4.1 Implementation

The file *vowels.jld* contains a supervised learning dataset where we are trying to predict which of the 11 “steady-state” English vowels that a speaker is trying to pronounce.

You are provided with a `randomTree` function in *randomTree.jl* (based on information gain). The random tree model differs from the decision tree model in two ways: it takes a bootstrap sample of the data before fitting and when fitting individual stumps it only considers  $\lfloor \sqrt{d} \rfloor$  randomly-chosen features.<sup>1</sup> In other words, `randomTree` is the model we discussed in class that is combined to make up a random forest.

If you run *example\_randomTree.jl*, it will fit both models to the dataset, and you will notice that it overfits badly.

1. If you set the *depth* parameter to *Inf*, why do the training functions terminate?
2. Why does the random tree model, using infoGain and a depth of *Inf*, have a training error greater 0?
3. Create a function `randomForest` that takes in hyperparameters `depth` and `nTrees` (number of trees), and fits `nTrees` random trees each with maximum depth `depth`. For prediction, have all trees predict and then take the mode. Hand in your function. Hint: you can define an array for holding 10 *GenericModel* types using:  
`subModels = Array{GenericModel}(undef,10).`
4. Using 50 trees, and a depth of  $\infty$ , report the training and testing error. Compare this to what we got with a single `DecisionTree` and with a single `RandomTree`. Are the results what you expected? Discuss.
5. Why does a random forest typically have a training error of 0, even though random trees typically have a training error greater than 0?

## 5 K-Means Clustering

If you run the function *example\_Kmeans*, it will load a dataset with two features and a very obvious clustering structure. It will then apply the  $k$ -means algorithm with a random initialization. The result of applying the algorithm will thus depend on the randomization, but a typical run might look like this: (Note that the colours are arbitrary due to the label switching problem.) But the ‘correct’ clustering (that was used to make the data) is something more like this:

---

<sup>1</sup>The notation  $\lfloor x \rfloor$  means the “floor” of  $x$ , or “ $x$  rounded down”.



## 5.1 Selecting among k-means Initializations

If you run the demo several times, it will find different clusterings. To select among clusterings for a *fixed* value of  $k$ , one strategy is to minimize the sum of squared distances between examples  $x_i$  and their means  $w_{y_i}$ ,

$$f(w_1, w_2, \dots, w_k, y_1, y_2, \dots, y_n) = \sum_{i=1}^n \|x_i - w_{y_i}\|_2^2 = \sum_{i=1}^n \sum_{j=1}^d (x_{ij} - w_{y_{ij}})^2.$$

where  $y_i$  is the index of the closest mean to  $x_i$ . This is a natural criterion because the steps of k-means alternately optimize this objective function in terms of the  $w_c$  and the  $y_i$  values.

1. Write a new function called *kMeansError* that takes in a dataset  $X$ , a set of cluster assignments  $y$ , and a set of cluster means  $W$ , and computes this objective function. Hand in your code.
2. Instead of printing the number of labels that change on each iteration, what trend do you observe if you print the value of *kMeansError* after each iteration of the k-means algorithm?
3. Using the *clustering2Dplot* file, output the clustering obtained by running k-means 50 times (with  $k = 4$ ) and taking the one with the lowest error. Note that the k-means training function will run much faster if you set `doPlot = false` or just remove this argument.

## 5.2 Selecting $k$ in k-means

We now turn to the task of choosing the number of clusters  $k$ .

1. Explain why the *kMeansError* function should not be used to choose  $k$ .
2. Explain why even evaluating the *kMeansError* function on test data still wouldn't be a suitable approach to choosing  $k$ .
3. Hand in a plot of the minimum error found across 50 random initializations, as you vary  $k$  from 1 to 10.
4. The *elbow method* for choosing  $k$  consists of looking at the above plot and visually trying to choose the  $k$  that makes the sharpest “elbow” (the biggest change in slope). What values of  $k$  might be reasonable according to this method? Note: there is not a single correct answer here; it is somewhat open to interpretation and there is a range of reasonable answers.

## 6 Very-Short Answer Questions

Write a short one or two sentence answer to each of the questions below. Make sure your answer is clear and concise.

1. What is a feature transformation that you might do to address a “coupon collecting” problem in your data?  
Answer: Discretization
2. What is one reason we would want to look at scatterplots of the data before doing supervised learning?  
Answer: Because summary statistics may trick us.
3. When we fit decision stumps, why do we only consider  $>$  (for example) and not consider  $<$  or  $\geq$ ?  
Answer: Because they will give the same result.
4. What is a reason that the data may not be IID in the email spam filtering example from lecture?

Answer: There might be the similar email from the same sender that will result in similarity and duplication which is not IID.

5. What is the difference between a validation set and a test set?

Answer: We will never look into the test set during our training. But we might use the validation set to determine our hyper-parameter. We are appropriate test error using validation error and we are using test error to approximate the error in real life.

6. Why can't we (typically) use the training error to select a hyper-parameter?

Answer: Because this will introduce a bias that we are picking the best test error. Which means we might overfit the test data.

7. What is an advantage and a disadvantage of using a large  $k$  value in  $k$ -fold cross-validation.

Answer: It will be more accurate with a higher  $k$  score. But it will also be more expensive

8. Why is naive Bayes called "naive"?

Answer: Because the naive bayes assumes each feature are independent, which means we are assuming  $P(f_1, f_2, f_3 \dots | y_i) = P(f_1 | y_i) * P(f_2 | y_i) * P(f_3 | y_i) \dots$

9. What is the effect of  $k$  in KNN on the two parts (training error and approximation error) of the fundamental trade-off. Hint: think about the extreme values.

Answer: If  $K$  is small, the training error will be small, but the approximation error might be big. And with a higher  $K$ , the training error will be small but the approximation error will be big.

Answer: Add similar translations to the training data as additional training data.

10. For any parametric model, how does increasing number of training examples  $n$  affect the two parts of the fundamental trade-off.

Answer: For parametric model, increasing the training examples will decrease the training error and also the approximation error.

11. Suppose we want to classify whether segments of raw audio represent words or not. What is an easy way to make our classifier invariant to small translations of the raw audio?
12. Both supervised learning and clustering models take in an input  $x_i$  and produce a label  $y_i$ . What is the key difference?
13. In  $k$ -means clustering the clusters are guaranteed to be convex regions. Are the areas that are given the same label by KNN also convex?