

# CPSC 340 Assignment 2 (due Friday September 30 at 11:55pm)

The assignment instructions are the same as Assignment 1, except you have the option to work in a group of 2. It is recommended that you work in groups as the assignment is quite long, but please only submit one assignment for the group.

Name(s) and Student ID(s):

## 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 training error with accuracy-based tree starts at a lower value, but eventually reaches a minimum error of 0.11. Whereas, the infogain-based decision tree starts at a high infogain value and eventually reaches 0 additional info gain.

An explanation for training error might be that with depth, the training error decreases as expected as the model becomes more detailed and specific. However, it reaches a peak as there are outliers that the decision tree is unable to fit to.

An explanation for infogain-based decision tree might be that in the first depth, there could be much information for it to learn as it is all new. However, as the depth increases, there is less new information for it to learn, and eventually hits 0 when it has been fitted and there is no more information to learn.

## 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?

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

Answer: The training error of the model decreases as the depth decreases to reach nearly 0. Whereas, the test error decreases but has a point where it stops decreasing at a certain depth.

### 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?

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

Answer:

## 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 12 training examples and 3 features:

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

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:

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

#### 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: 7/12

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

Answer: 5/12

#### 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: 1/7

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

Answer: 1/7

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

Answer: 4/7

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

Answer:  $4/5$

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

Answer:  $3/5$

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

Answer:  $1/5$

### 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:  $p(\text{spam} = 1 \mid <\text{your name}> = 1, \text{pharmaceutical} = 0, \text{PayPal} = 1) = p(<\text{your name}> = 1 \mid \text{spam}) * p(\text{pharmaceutical} = 0 \mid \text{spam}) * p(\text{PayPal} = 1 \mid \text{spam})$   
 $= 1/7 * 1/7 * 4/7 = 4/343$

$p(\text{Not spam} = 1 \mid <\text{your name}> = 1, \text{pharmaceutical} = 0, \text{PayPal} = 1) = 4/5 * 3/5 * 1/5 = 12/125$

The most likely label for the test example is "Not Spam" as the probability of being "Spam" is much smaller than the probability of the message being "Not 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. *Xtest* and *ytest*: 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 50 of *X*?

Answer: league

2. Which words are present in training example 500?

Answer: car, engine, evidence, problem, system

3. Which newsgroup name does training example 500 come from?

Answer: rec.\*

## 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 other 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: New test error = 0.599

## 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:  $O(tkd)$ .

This involves calculating all possible values of the class labels given the test example, which takes  $O(kd)$  time. Since look up and multiplication takes  $O(d)$  and it has to be repeated  $k$  times as there are  $k$  class labels. This will be the time to classify one test example. This will be repeated  $t$  times as there are  $t$  examples.



### 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:

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

Answer: When  $k = 1$ , training error = 0.000, test error = 0.065

When  $k = 2$ , training error = 0.032, test error = 0.092

When  $k = 3$ , training error = 0.028, test error = 0.066

3. Hand in the plot generated by *plot2Dclassifier* on the *citiesSmall.mat* dataset for  $k = 1$  on the training data.

Answer:

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?

Answer: Vancouver is located approximately in 49N and 123W, thus it would be in the blue state.

5. Why is the training error 0 for  $k = 1$ ?

Answer: As the nearest neighbour to any point in the same dataset is itself, by KNN's algorithm, its predicted result would also be itself. Therefore, the predicted result on the training dataset will always correct and cause the training error to be 0

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

Answer: A way to choose  $k$  would be to use

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.

## 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?

Answer: The training function terminates when there is a base-split, which occurs at the root of the tree. There is no longer any probabilities for the tree to split into, as they are as detailed as they can get.

2. Why does the random tree model, using infoGain and a depth of *Inf*, have a training error greater 0?

Answer: This is as the random tree only considers  $\sqrt{d}$  randomly-chosen features. Therefore, even at a depth of infinity, it will not fully train as the tree will overlook some features, resulting in a training error greater than 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).`

Answer:

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.

Answer: Training Error = 0.000, Test Error = 0.178

Yes, the answers are as expected. The errors are significantly lower than the error from a single random

---

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

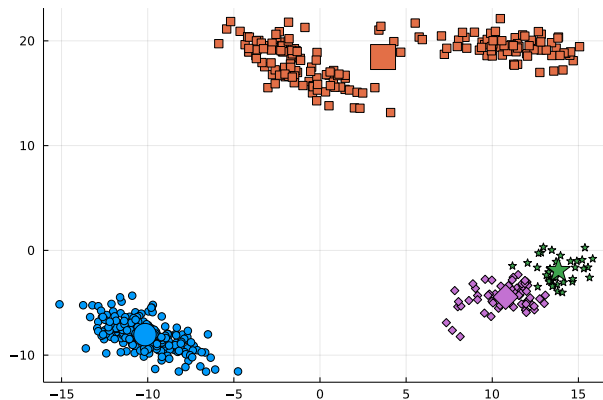
tree. The reasons are stated below.

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?

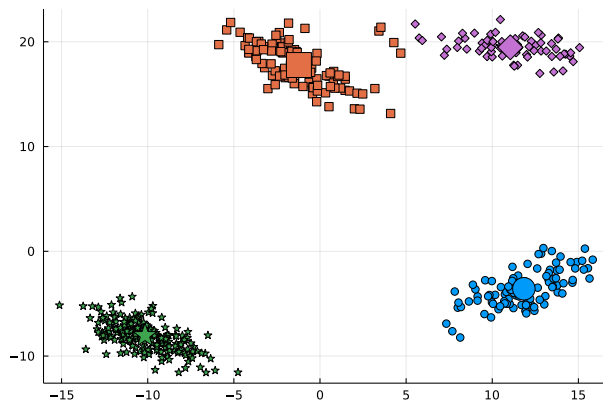
Answer: This is as random forest uses an ensemble voting method that is voted by multiple random trees. Thus, since each tree considers  $\sqrt{d}$  features, when all trees are considered individually, all features are as a by-product considered. Furthermore, since each tree generally has a  $<0.5$  training error, by probability and the binomial theorem concept, the training error should decrease for the random forest. Thus, with both considerations in mind, the random forest would have a lower training error than the random trees individually.

## 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:



### 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.

Answer:

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?

Answer: The *kMeansError* decreases after every iteration, but it reaches a point where it stops decreasing and the function stops and returns.

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  $k$ .

Answer: *kMeansError* will decrease as  $k$  increases. If there are more discrete points, the distance from any point to the nearest point will decrease, thus making it not a good measure to choose  $k$  with.

2. Explain why even evaluating the *kMeansError* function on test data still wouldn't be a suitable approach to choosing  $k$ .

Answer:

3. Hand in a plot of the minimum error found across 50 random initializations, as you vary  $k$  from 1 to 10.

Answer:

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.

Answer:

### 5.3 $k$ -Medians

The data in `clusterData2.mat` is the exact same as the above data, except it has 4 outliers that are very far away from the data.

1. Using the `clustering2Dplot` function, output the clustering obtained by running k-means 50 times (with  $k = 4$ ) on `clusterData2.mat` and taking the one with the lowest error. Are you satisfied with the result?
2. Hand in the elbow plot for this data. What values of  $k$  might be chosen by the elbow method for this dataset?
3. Instead of the squared distances between the examples  $x_i$  and their cluster centers, consider measuring the distance to the cluster centers in the L1-norm,

$$f(w_1, w_2, \dots, w_k, y_1, y_2, \dots, y_n) = \sum_{i=1}^n \|x_i - w_{y_i}\|_1 = \sum_{i=1}^n \sum_{j=1}^d |x_{ij} - w_{y_i j}|.$$

Hand in the elbow plot for k-means when we measure the error of the final model with the L1-norm. What value of  $k$  would be chosen by the elbow method?

4. The k-means algorithm tries to minimize the squared error and not the L1-norm error, so in the last question there is a mis-match between the what the learning algorithm tries to minimize and how we measure the final error. We can try to directly minimize the L1-norm with the *k-medians* algorithm, which assigns examples to the nearest  $w_c$  in the L1-norm and to updates the  $w_c$  by setting them to the “median” of the points assigned to the cluster (we define the  $d$ -dimensional median as the concatenation of the median of the points along each dimension). Implement the *k-medians* algorithm, and hand in your code and the the plot obtained by minimizing the L1-norm error across 50 random initializations of *k-medians* with  $k = 4$ .



## 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: Some feature transformation that can be applied would be feature aggregation, discretization, or feature selection. This will reduce the number of features such that there would be less “coupons” to collect, making it faster to learn.

2. What is one reason we would want to look at scatterplots of the data before doing supervised learning?

Answer: One reason could be to understand the relationship of the data or to identify outliers in the dataset for preprocessing.

3. When we fit decision stumps, why do we only consider  $>$  (for example) and not consider  $<$  or  $\geq$ ?

Answer: All of the signs results in the same model. However it is ideal to use the same condition such that it is consistent.

4. What is a reason that the data may not be IID in the email spam filtering example from lecture?

Answer: This is as some words may in general come together, such as “I” and “am”, which usually comes together in sentences.

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

Answer: The validation set comes from a part of the training dataset for hyper-parameter optimization, while the test set is a separate dataset that is set aside to obey the golden rule.

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

Answer: In the example for a decision tree, maximizing the depth will increase the training error, but overfit the model. In the general case, using training error to select hyper-parameter would cause the model to be overfitted and not being able to generalize to new data.

7. If you can fit one model in 10 sec., how long (in days) does it take to find the best among a set of 16 hyperparameter values using leave-one-out cross-validation on a dataset containing  $n = 4320$  examples?

Answer:  $10 \cdot 16 \cdot 4320 / 86400 = 8$  days

8. Naïve Bayes makes the assumption that all features are conditionally independent given the class label. Why is this assumption necessary and what would happen without it?

Answer: This assumption is necessary for the purpose of not having to require a huge amount of data. As the probability of the exact scenario happening is twice as unlikely in a small dataset, Naive Bayes gives a good estimate of its probability, thus reducing the data necessary.

9. Why is KNN considered a non-parametric method and what are two undesirable consequences of KNNs non-parametric design?

Answer: KNN is considered as non-parametric as it is not predicted using variables and its storage depends on ‘n’, which is its training examples. One undesirable consequences of KNN non-parametric design is that the number of parameters grows with ‘n’, thus it is not able to scale with infinite data as well. Another undesirable consequence is that the model gets more complicated as data increases, making it again not as scalable.

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

Answer: By increasing the number of training examples, in general, data would increase the training accuracy, but decrease the approximation accuracy. However, as parametric models are simple, at a certain point, more data may not be necessarily useful to train any further.

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?

Answer: A simple way would be to change the pitch or the bass of the sound slightly.

12. Both supervised learning and clustering models take in an input  $x_i$  and produce a label  $y_i$ . What is the key difference?

Answer: Clustering method does not require target feature in order to train the model.

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?

Answer: No, they are not guaranteed to be convex.