

# CPSC 340 Assignment 1 (due 2018-01-17 at 9:00pm)

## Instructions

Rubric: {mechanics:5}

**IMPORTANT! Before proceeding, please carefully read the general homework instructions at [https://github.ugrad.cs.ubc.ca/CPSC340-2017W-T2/home/blob/master/homework\\_instructions.md](https://github.ugrad.cs.ubc.ca/CPSC340-2017W-T2/home/blob/master/homework_instructions.md).** Please pay special attention to the section on visualizations as there were no visualizations required in a0.

**A note on the provided code:** in the *code* directory we provide you with a file called *main.py*. This file, when run with different arguments, runs the code for different parts of the assignment. For example,

```
python main.py -q 1.1
```

runs the code for Question 1.1. At present, this should do nothing, because the code for Question 1.1 still needs to be written (by you). But we do provide some of the bits and pieces to save you time, so that you can focus on the machine learning aspects. The file *utils.py* contains some helper functions. You're not required to read/understand the code in there (although you're welcome to!) and will not need to modify it.

## 1 Data Exploration

Your repository contains the file *fluTrends.csv*, which contains estimates of the influenza-like illness percentage over 52 weeks on 2005-06 by Google Flu Trends. Your *main.py* loads this data for you and stores it in a pandas DataFrame *X*, where each row corresponds to a week and each column corresponds to a different region. If desired, you can convert from a DataFrame to a raw numpy array with *X.values()*.

### 1.1 Summary Statistics

Rubric: {reasoning:2}

**Report the following statistics:** The minimum, maximum, mean, median, and mode of all values across the dataset.

1. The 5%, 25%, 50%, 75%, and 95% quantiles of all values across the dataset.
2. The names of the regions with the highest and lowest means, and the highest and lowest variances.

In light of the above, **is the mode a reliable estimate of the most “common” value? Describe another way we could give a meaningful “mode” measurement for this (continuous) data.** Note: the function *utils.mode()* will compute the mode value of an array for you.

**Answer:**

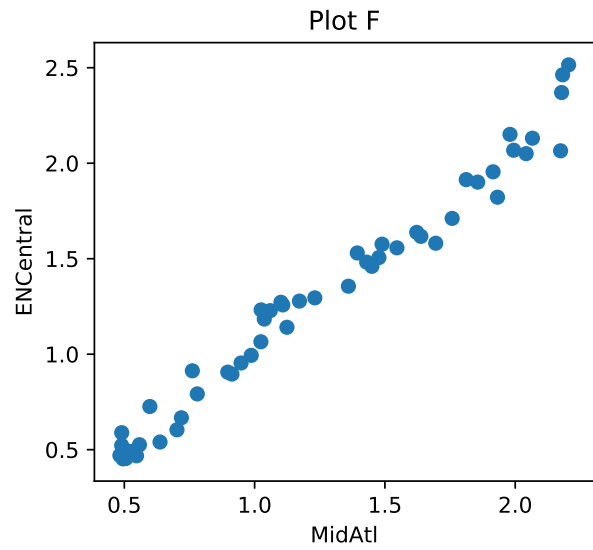
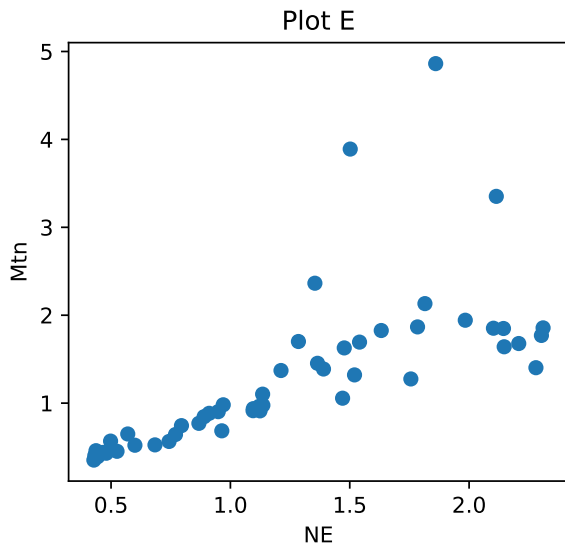
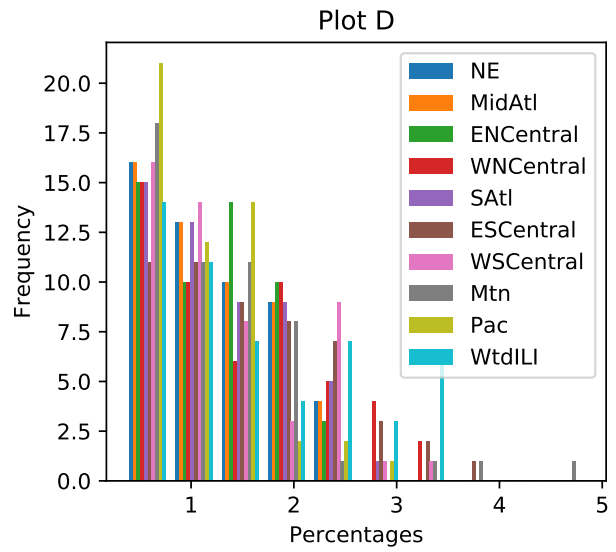
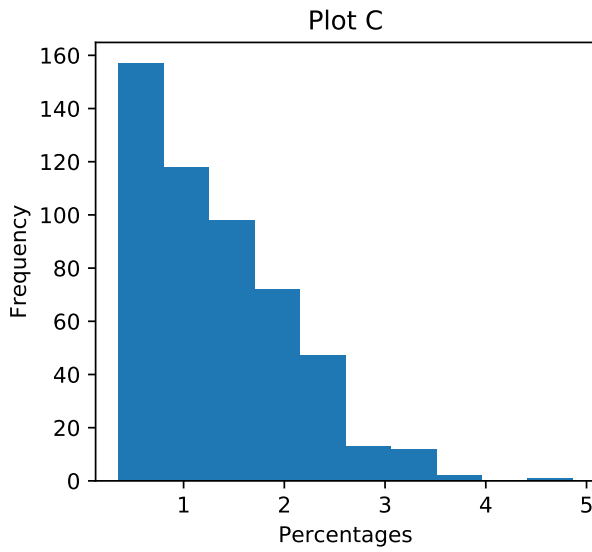
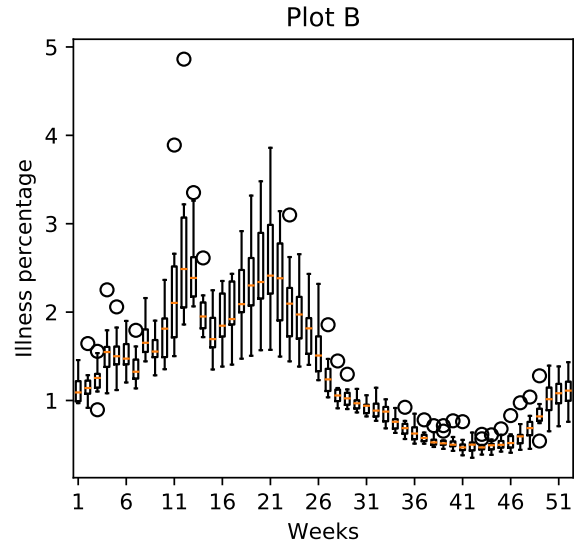
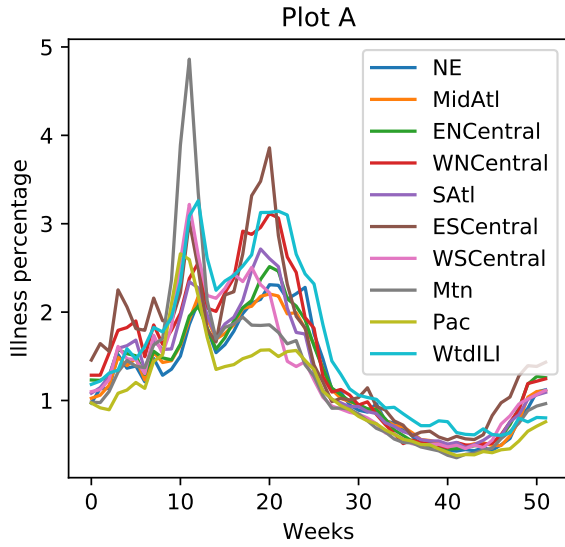
1. The minimum of all values across the dataset is 0.3520, the maximum is 4.8620, the mean is 1.3246, the mode is 0.7700, the median is 1.1590.
2. The 5%, 25%, 50%, 75%, and 95% quantiles of all values across the dataset are respectively 0.4650, 0.7180, 1.1590, 1.8133, 2.6240.

3. The region with the highest mean (1.5670) is WtdILI and the one with the lowest mean (1.0632) is Pac, and the one with the highest variance (0.7834) is Mtn and the one with the lowest variance (0.3158) is Pac.

## 1.2 Data Visualization

Rubric: {reasoning:3}

Consider the figure below.



**Figure 1:** visualize

The figure contains the following plots, in a shuffled order:

1. A single histogram showing the distribution of *each* column in  $X$ .
2. A histogram showing the distribution of each the values in the matrix  $X$ .
3. A boxplot grouping data by weeks, showing the distribution across regions for each week.
4. A plot showing the illness percentages over time.
5. A scatterplot between the two regions with highest correlation.
6. A scatterplot between the two regions with lowest correlation.

Match the plots (labeled A-F) with the descriptions above (labeled 1-6), with an extremely brief (a few words is fine) explanation for each decision.

Answer:

1. D. D is a histogram and each column in  $X$  (like NE) can be found in it.
2. C. C is a histogram and we can find the distribution of values.
3. B. B is the only boxplot with the x-axis set to be weeks.
4. A. A shows the relation between illness and weeks from the axis.
5. F. F is a scatterplot and it has higher values of correlation (dots around the line  $y=x$ ).
6. E. E is a scatterplot and it has lower correlation compared to F..

## 2 Decision Trees

If you run `python main.py -q 2`, it will load a dataset containing longitude and latitude data for 400 cities in the US, along with a class label indicating whether they were a “red” state or a “blue” state in the 2012 election.<sup>1</sup> Specifically, the first column of the variable  $X$  contains the longitude and the second variable contains the latitude, while the variable  $y$  is set to 1 for blue states and 2 for red states. After it loads the data, it plots the data and then fits two simple classifiers: a classifier that always predicts the most common label (1 in this case) and a decision stump that discretizes the features (by rounding to the nearest integer) and then finds the best equality-based rule (i.e., check if a feature is equal to some value). It reports the training error with these two classifiers, then plots the decision areas made by the decision stump.

### 2.1 Splitting rule

Rubric: {reasoning:1}

Is there a particular type of features for which it makes sense to use an equality-based splitting rule rather than the threshold-based splits we discussed in class?

Answer: Categorical features

### 2.2 Decision Stump Implementation

Rubric: {code:3}

The file `decision_stump.py` contains the class `DecisionStumpEquality` which finds the best decision stump using the equality rule and then makes predictions using that rule. Instead of discretizing the data and using a rule based on testing an equality for a single feature, we want to check whether a feature is above or below a threshold and split the data accordingly (this is the more sane approach, which we discussed in class). Create a `DecisionStump` class to do this, and report the updated error you obtain by using inequalities instead of discretizing and testing equality.

Hint: you may want to start by copy/pasting the contents `DecisionStumpEquality` and then make modifications from there.

Answer: In the `DecisionStumpEquality` function, it splits the dataset by “==” or “!=” value. To change it to the `inEquality` function, we have to change the “==” part to “>” and “<=” so as to meet the requirement

---

<sup>1</sup>The cities data was sampled from <http://simplemaps.com/static/demos/resources/us-cities/cities.csv>. The election information was collected from Wikipedia.

to check whether a feature is above or below a threshold. From the q2 part, the mode predictor error is 0.415 and the error of decision stump with equality rule is 0.380. When using inequality rule, the error is reduced to 0.253.

The figure below shows the situation with inequality rule.

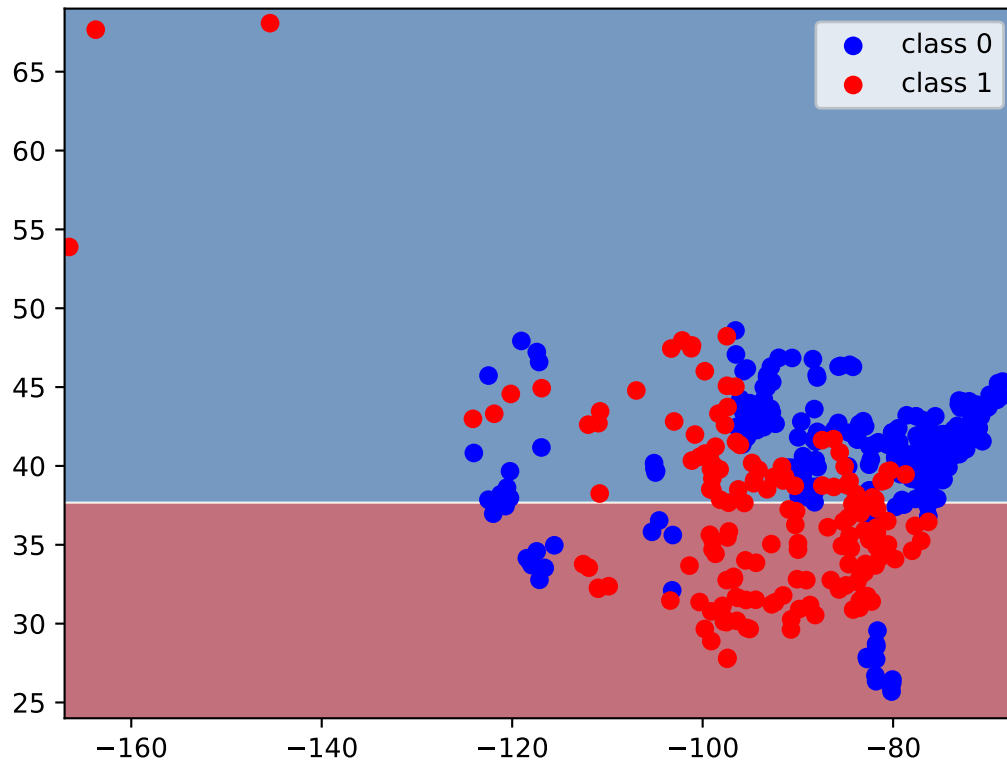


Figure 2: Decision boundary for inequality

## 2.3 Constructing Decision Trees

Rubric: {code:2}

Once your decision stump class is finished, running `python main.py -q 2.3` will fit a decision tree of depth 2 to the same dataset (which results in a lower training error). Look at how the decision tree is stored and how the (recursive) *predict* function works. Using the same splits as the fitted depth-2 decision tree, write an alternative version of the predict function for classifying one training example as a simple program using if/else statements (as in slide 6 of the Decision Trees lecture). Save your code in a new file called `simple_decision.py` (in the *code* directory) and make sure you link to this file from your README.

Note: this code should implement the specific, fixed decision tree which was learned after calling `fit` on this particular data set. It does not need to be a learnable model.

Answer: The URL of the python code is: [https://github.com/ugrad.cs.ubc.ca/CPSC340-2017W-T2/u7p1b\\_a1/blob/master/code/simple\\_decision.py](https://github.com/ugrad.cs.ubc.ca/CPSC340-2017W-T2/u7p1b_a1/blob/master/code/simple_decision.py).

## 2.4 Decision Tree Training Error

Rubric: {reasoning:2}

Running `python main.py -q 2.4` fits decision trees of different depths using two different implementations: first, our own implementation using your `DecisionStump`, and second, the decision tree implementation from the popular Python ML library *scikit-learn*. The decision tree from sklearn uses a more sophisticated splitting criterion called the information gain, instead of just the classification accuracy. Run the code and look at the figure. Describe what you observe. Can you explain the results?

Note: we set the `random_state` because sklearn's `DecisionTreeClassifier` is non-deterministic. I'm guessing this is because it breaks ties randomly.

Note: the code also prints out the amount of time spent. You'll notice that sklearn's implementation is substantially faster. This is because our implementation is based on the  $O(n^2d)$  decision stump learning algorithm and sklearn's implementation presumably uses the faster  $O(nd \log n)$  decision stump learning algorithm that we discussed in lecture.

Answer: From the figure below,

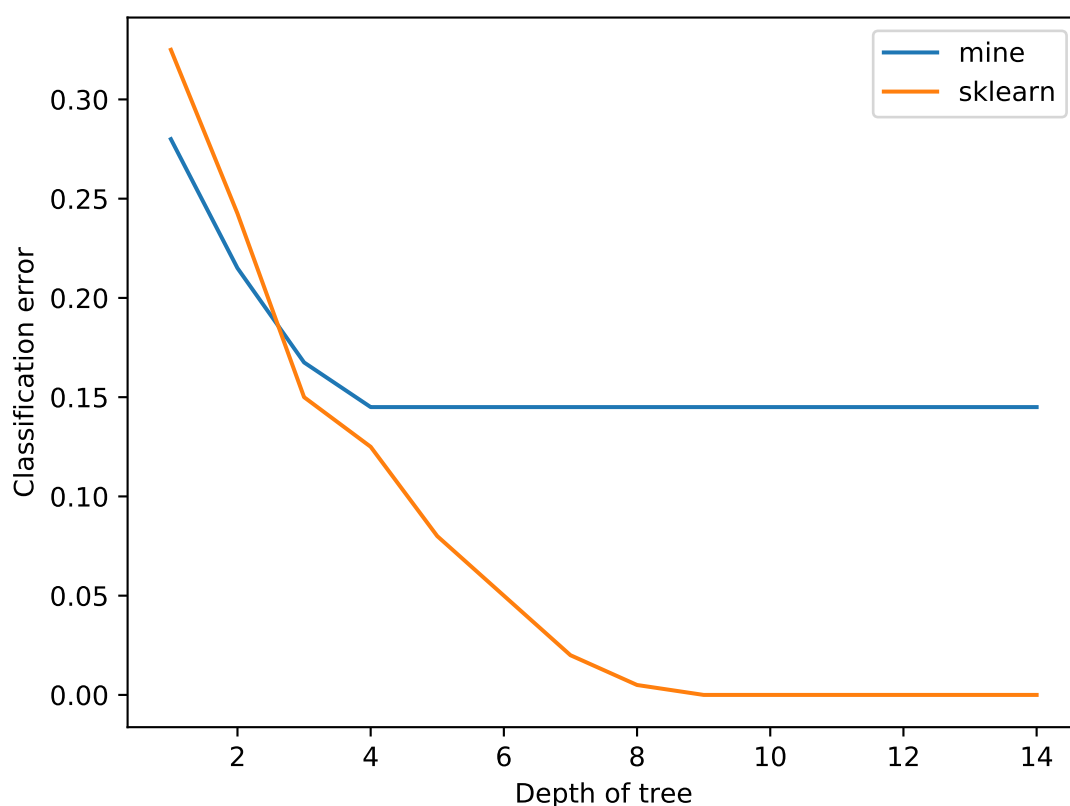


Figure 3: Tree errors

We can find that although when the depth is very small, the error of sklearn's `DecisionTreeClassifier` might be larger and the two errors show similar trend before error of my `DecisionStump` gets stable. After the error of my `DecisionStump` get stable in around 0.15, with depth increasing, the error of sklearn's `DecisionTreeClassifier` continues to decrease and results in zero after depth reaches 9. The reason should be that our decision stump does splitting works based on the mode of labels. Thus the split cannot help increase the accuracy when the modes of labels in two sides are the same.

## 2.5 Cost of Fitting Decision Trees

Rubric: {reasoning:3}

In class, we discussed how in general the decision stump minimizing the classification error can be found in  $O(nd \log n)$  time. Using the greedy recursive splitting procedure, what is the total cost of fitting a decision tree of depth  $m$  in terms of  $n$ ,  $d$ , and  $m$ ?

Hint: even though there could be  $(2^m - 1)$  decision stumps, keep in mind not every stump will need to go through every example. Note also that we stop growing the decision tree if a node has no examples, so we may not even need to do anything for many of the  $(2^m - 1)$  decision stumps.

Answer: As the depth is  $m$ , the total number of decision stumps is  $(2^m - 1)$ . This means that it seems the total cost should be  $O(2^m nd \log n)$ . However, not every stump in this decision tree will need to go through every example. In every stump, there is actually an if/else statement and each example should only choose one stump in each layer. As the depth is  $m$ , the total cost should be  $O(mnd \log n)$ .

## 3 Training and Testing

If you run `python main.py -q 3`, it will load *citiesSmall.pkl*. Note that this file contains not only training data, but also test data, `X_test` and `y_test`. After training a depth-2 decision tree it will evaluate the performance of the classifier on the test data.<sup>2</sup> With a depth-2 decision tree, the training and test error are fairly close, so the model hasn't overfit much.

### 3.1 Training and Testing Error Curves

Rubric: {reasoning:2}

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: it's OK to reuse the provided code from part 2.4 as a starting point.

Answer: As the plot shown below, the training error decreases to 0 after the depth increases to 9. The testing error also shows decreasing trend. But when the depth increases to 9, it shows a slight increase and stays stable with the value of 0.0835.

---

<sup>2</sup>The code uses the "information gain" splitting criterion; see the Decision Trees bonus slides for more information.

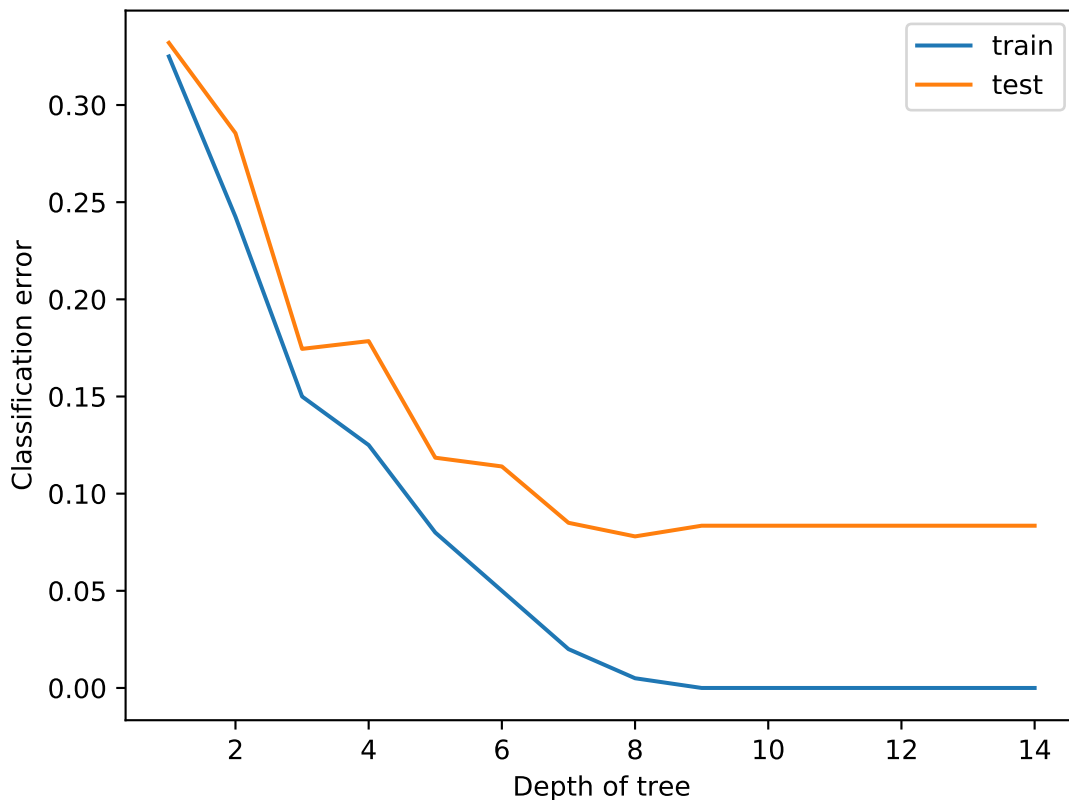


Figure 4: Train and test errors

### 3.2 Validation Set

Rubric: {reasoning:3}

Suppose that we didn't have an explicit test set available. 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 to minimize the validation set error? Does the answer change if you switch the training and validation set? How could use more of our data to estimate the depth more reliably?

Answer: The depth is 8 when the validation set error is minimized with the value 0.100. After switch the validation set and the training set, the depth becomes 9 when the minimum validation set error appears as 0.145. To make it more reliable, we could use cross validation method.

The training error and validation error before and after switch is drawn in one graph, the URL is [https://github.ugrad.cs.ubc.ca/CPSC340-2017W-T2/u7p1b\\_a1/blob/master/figs/q3\\_2\\_train\\_vali\\_errors.pdf](https://github.ugrad.cs.ubc.ca/CPSC340-2017W-T2/u7p1b_a1/blob/master/figs/q3_2_train_vali_errors.pdf).

## 4 K-Nearest Neighbours

In the *citiesSmall* dataset, nearby points tend to receive the same class label because they are part of the same U.S. state. This indicates that a  $k$ -nearest neighbours classifier might be a better choice than a decision



tree. The file *knn.py* has implemented the training function for a  $k$ -nearest neighbour classifier (which is to just memorize the data).

## 4.1 KNN Prediction

Rubric: {code:3, reasoning:4}

Fill in the *predict* function in *knn.py* so that the model file implements the  $k$ -nearest neighbour prediction rule. You should Euclidean distance, and may numpy's *sort* and/or *argsort* functions useful. You can also use *utils.euclidean\_dist\_squared*, which computes the squared Euclidean distances between all pairs of points in two matrices.

1. Write the *predict* function.

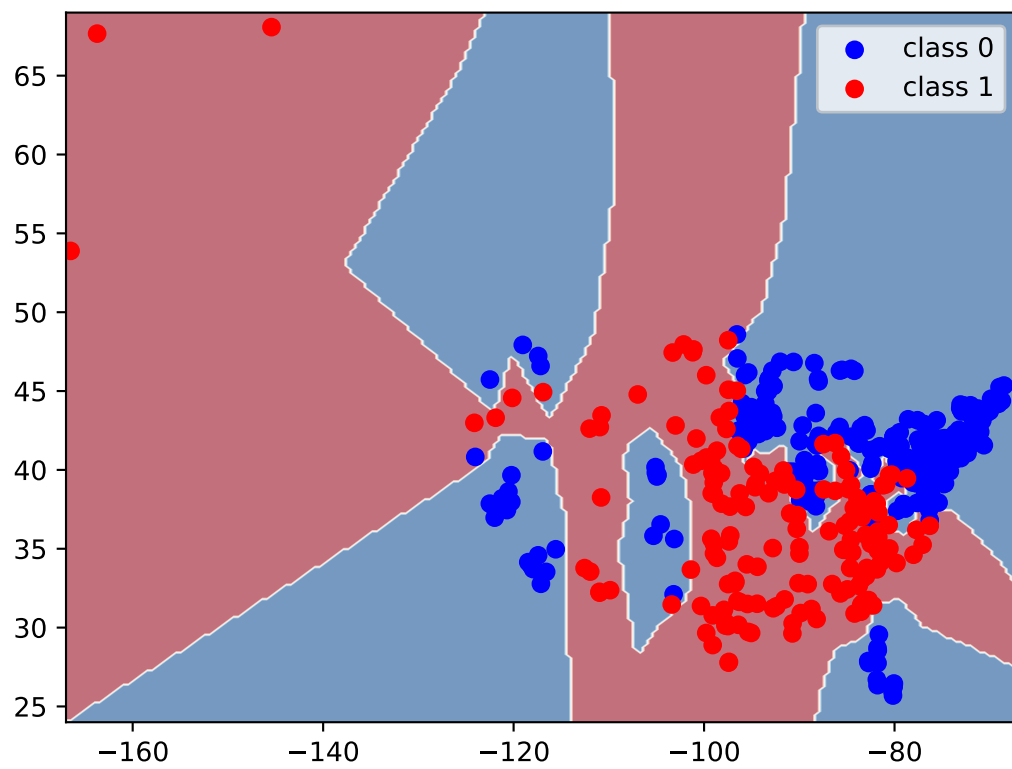
Answer: The URL of this function is: [https://github.ugrad.cs.ubc.ca/CPSC340-2017W-T2/u7p1b\\_a1/blob/master/code/knn.py](https://github.ugrad.cs.ubc.ca/CPSC340-2017W-T2/u7p1b_a1/blob/master/code/knn.py); The URL of the prediction result is: [https://github.ugrad.cs.ubc.ca/CPSC340-2017W-T2/u7p1b\\_a1/blob/master/figs/q4\\_1\\_decisionBoundary.pdf](https://github.ugrad.cs.ubc.ca/CPSC340-2017W-T2/u7p1b_a1/blob/master/figs/q4_1_decisionBoundary.pdf).

2. Report the training and test error obtained on the *citiesSmall* dataset for  $k = 1$ ,  $k = 3$ , and  $k = 10$ . How do these numbers compare to what you got with the decision tree?

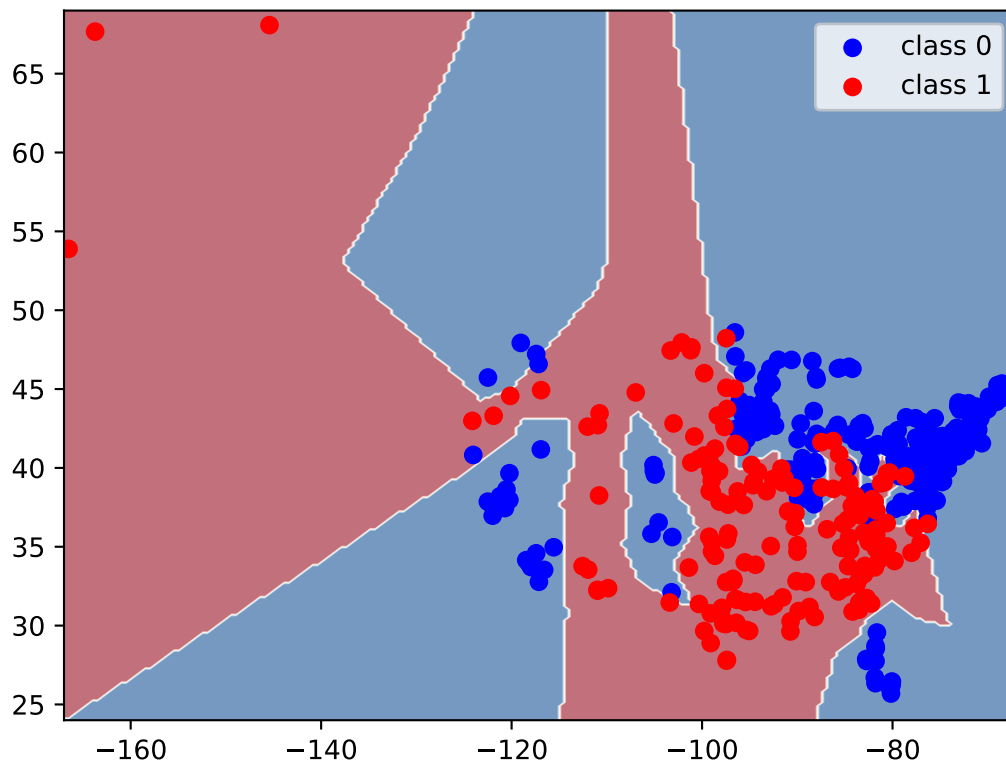
Answer: when  $k = 1$ , training error=0.0, testing error=0.085; when  $k = 3$ , training error=0.040, testing error=0.055; when  $k = 10$ , training error=0.125, testing error=0.185. Compared with the results from the decision tree, the training error may not be 0 and increase with  $k$  increasing. Also, with  $k$  increasing, the testing error also shows increase and the value could be larger than that from the decision tree with proper depth.

3. Hand in the plot generated by *utils.plotClassifier* on the *citiesSmall* dataset for  $k = 1$ , using both your implementation of KNN and the *KNeighborsClassifier* from *scikit-learn*.

Answer: The figure below shows the result of *utils.plotClassifier*:



**Figure 5:** Decision boundary with KNN



**Figure 6:** Decision boundary with K-N-Classifer

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

Answer: The training error is 0 when  $k=1$  because every training example is the 1NN of itself. Thus the prediction process is just copy the labels from the training set.

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

Answer: The training error increases as  $k$  increases, which means that the training error cannot be used to choose  $k$ . The method should be to split the data into training and validation set or use cross-validation method.

## 4.2 Condensed Nearest Neighbours

Rubric: {code:3, reasoning:5}

The dataset *citiesBig1* 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). If the subset is small, then prediction will be faster.

The most common variation of this algorithm works as follows:

```

initialize subset with first training example;
for each subsequent 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.py*.

Your tasks:

1. The point of this algorithm is to be faster than KNN. Try running the condensed NN 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 – how long did it take before you panicked and went for CTRL-C... or did it actually finish?

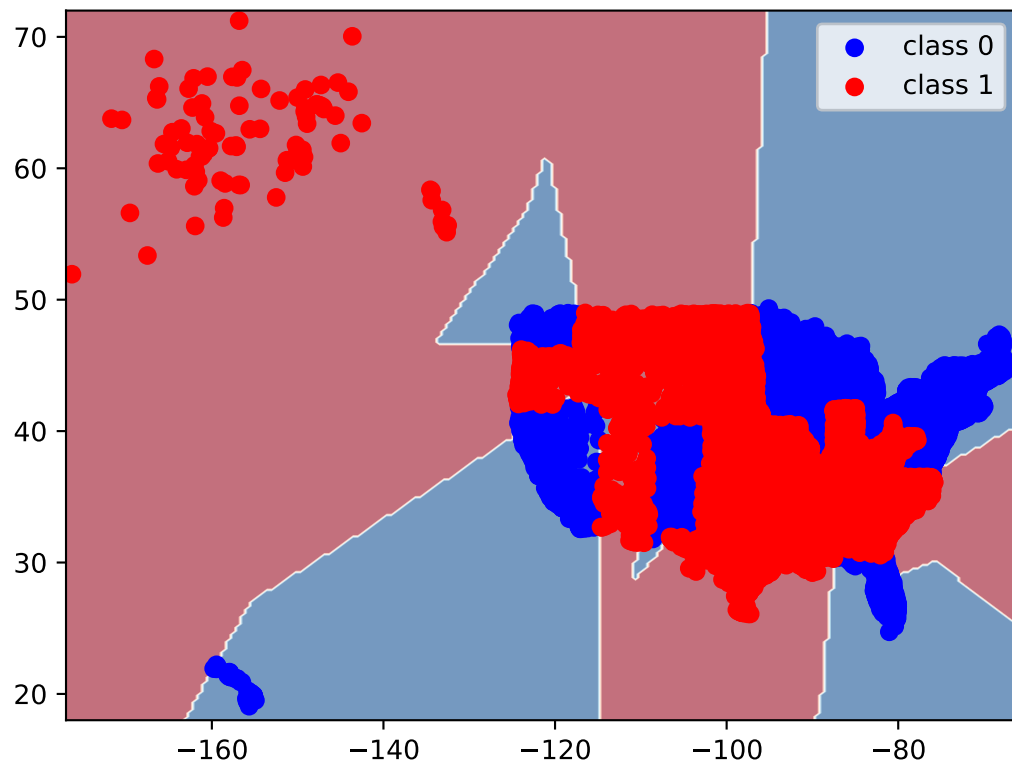
Answer: When split the data into 9/10 of training set and 1/10 of validation set, the prediction for training set takes 10.1798s. This is very fast compared to using KNN, which is endless without using Ctrl-C.

2. Report the training and testing errors for condensed NN, as well as the number of variables in the subset, on the *citiesBig1* dataset with  $k = 1$ .

Answer: The training error is 0.0087; The testing error is 0.0209; The number of variables in the subset is 311.

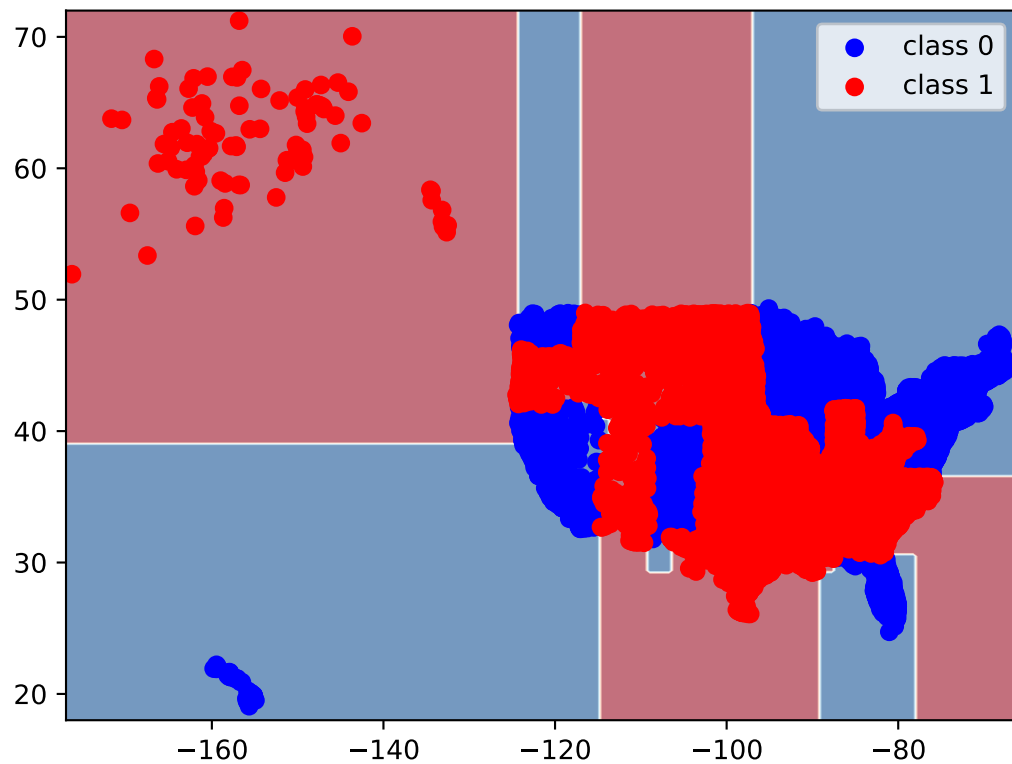
3. Hand in the plot generated by *utils.plotClassifier* on the *citiesBig1* dataset for  $k = 1$ .

Answer: The figure below shows the plot:



**Figure 7:** Decision boundary for CNN

4. Why is the training error with  $k = 1$  now greater than 0?  
 Answer: In this case, we do not store examples that were correctly classified during the training phase. When we add more examples to the subset, these examples can become incorrectly classified. We can go through the dataset for more times to fix these errors.
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  $s$ ?  
 Answer: Originally we calculate  $n$  distances to  $t$  test points. But now, we calculate  $s$  distances to  $t$  test points. Since to calculate one distance takes  $O(d)$ , the cost is reduced from  $O(ndt)$  to  $O(sdt)$ .
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?  
 Answer: In the *citiesBig2*, the dataset including the training set and the test set is sorted (by state), so the data is not even close to IID.  
 The reason that the test error is high: there are some states that are not present in the training set.  
 The reason that the training error is high: the CNN method depends on the order of the examples thus all the cities in a state might be correctly classified as we visit all at once but then they can become incorrectly classified since we are trying to get other states right and don't visit them again.
7. Try running a decision tree on *citiesBig1* using scikit-learn's `DecisionTreeClassifier` with default hyperparameters. Does it work? Is it fast? Any other thoughts? Overall, do you prefer decision trees of (condensed) KNN for this data set? Include the usual plot of the decision surface.  
 Answer: The plot using decision tree is:



**Figure 8:** Decision boundary for decision tree

The method using the decision tree is pretty faster than using the CNN. And the testing error is smaller (0.013 and 0.025). Thus I prefer the decision tree.

## 5 Very-Short Answer Questions

Rubric: {reasoning:8}

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

Answer: There might be some odd data in the dataset. Or the data might be simple and show obvious features.

2. What is a reason that the examples in a training and test set might not be IID?

Answer: The examples are sorted.

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

Answer: A validation set is a part of the training set and is to use to approximate test error after training. A test set is a set for testing which cannot have any influence on training. A test set may contain new data.

4. What is the main advantage of non-parametric models?

Answer: The more the data is, the more complicated the model will be.

5. A standard pre-processing step is “standardization” of the features: for each column of  $X$  we subtract its mean and divide by its variance. Would this pre-processing change the accuracy of a decision tree classifier? Would it change the accuracy of a KNN classifier?  
Answer: For decision tree, it would not as it does not change the relative ordering of the features. For KNN classifier, it would as it changes the importance of each dimension.
6. Does increasing  $k$  in KNN affect the training or prediction asymptotic runtimes?  
Answer: No. The runtime is not related to  $k$  in this case for both fitting and predicting process.
7. How does increase the parameter  $k$  in  $k$ -nearest neighbours affect the two parts (training error and approximation error) of the fundamental trade-off (hint: think of the extreme values).  
Answer: Increasing  $k$  will increase the training error but decrease the approximate error.
8. For any parametric model, how does increasing number of training examples  $n$  affect the two parts of the fundamental trade-off.  
Answer: Increasing  $n$  will lead to a relatively high training error and decrease the approximate error.