

HW 5

CS 216, Everything Data, Spring 2020

DUE: Monday Feb. 24 by 4:40 pm (class time)

Joe Cusano (jgc28) and Pierce Forte (ph7)

In this assignment, you will use Numpy to build three different basic classifiers for prediction. You will include all of your answers for this assignment within this notebook. You will then convert your notebook to a .pdf and a .py file to submit to gradescope (submission instructions are included at the bottom).

Please take note of the [course collaboration policy \(https://sites.duke.edu/compsci216s2020/policies/\)](https://sites.duke.edu/compsci216s2020/policies/). You may work alone or with a single partner. If you work with a partner, you may not split up the assignment; you should work together in-person or complete parts independently and come together to discuss your solutions. In either case, you are individually responsible for your work, and should understand everything in your submission.

Part 1: Getting Started with Numpy

Numpy is the standard library for scientific computing with Python, and Numpy arrays are the standard data structure for working with prediction tasks in Python. If you are unfamiliar with Numpy, we recommend that you start this assignment by reviewing the brief tutorial on Numpy at <http://cs231n.github.io/python-numpy-tutorial/#numpy> (<http://cs231n.github.io/python-numpy-tutorial/#numpy>) (just the Numpy section). The Numpy documentation itself also contains an expanded tutorial <https://numpy.org/doc/1.17/user/quickstart.html> (<https://numpy.org/doc/1.17/user/quickstart.html>). We mention a few particularly important notes here.

- You can easily turn any Python list into a Numpy array (for example, `ar = numpy.array([0,1,2])` creates a Numpy array named `ar` containing 1, 2, and 3).
- Indexing and slicing 1-d Numpy arrays is similar to indexing and slicing Python lists. For example, `ar[1:]` returns `[1,2]`.
- Indexing and slicing 2-d Numpy arrays is similar to using the `.iloc` method on a Pandas dataframe. For example, if you have a 2-d Numpy array `Mat` and you want the entry for row 0, column 5, you would just index `Mat[0,5]`. If you want column 5, you just write `Mat[:,5]`.
- Operations on Numpy arrays are element-wise *by default*. For example, `ar+5` would return `[5, 6, 7]`. `ar+ar` would return `[0, 2, 4]`.
- Building on the element-wise theme, writing a boolean condition will return a True/False array. For example, `ar==1` would return `[False, True, False]`.
- Similar to Pandas filtering, you can use a boolean mask of this sort to filter a Numpy array. For example, `ar[ar >= 1]` would return `[1, 2]`.
- Numpy comes with full support for nearly all mathematical computing needs. You can find the reference documentation at <https://docs.scipy.org/doc/numpy/reference> (<https://docs.scipy.org/doc/numpy/reference>). Of particular interest are the many mathematical functions implemented in Numpy <https://docs.scipy.org/doc/numpy/reference/routines.math> (<https://docs.scipy.org/doc/numpy/reference/routines.math>) and the many random sampling functions implemented in Numpy <https://docs.scipy.org/doc/numpy/reference/random/index> (<https://docs.scipy.org/doc/numpy/reference/random/index>).

Once you have familiarized yourself with the basics of Numpy, it's time to start building classifiers for prediction. Run the below codeblocks to import libraries and define a function to calculate the uniform error of a prediction. Recall that the uniform error is just the percentage of predictions that are not the same as the true class.

```
In [1]: import pandas as pd
import numpy as np
```

```
In [2]: def uniform_error(prediction, target):
n = len(prediction)
if(n != len(target)):
    print('Error: prediction and target should have same length')
    return(1)
else:
    return((n-np.sum(prediction==target))/n)
```

Part 2: Naive Bayes Classifier

In this part of the assignment, you will implement a simple Naive Bayes classifier for predicting political party membership of members of the house of representatives based on votes.

Each row corresponds to a member of congress. The first column has a 1 if that member is a republican, and a 0 if that member is a democrat. The next 16 columns contain information about sixteen different votes that were taken in that year; there is a 1 if the member voted positively on that issue, and a 0 if that member voted negatively on that issue. We want to predict the first column based on the next sixteen.

```
In [73]: df_v = pd.read_csv('votes.csv')
df_v.head()
```

Out[73]:

	republican	v_1	v_2	v_3	v_4	v_5	v_6	v_7	v_8	v_9	v_10
0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0
2	1	1	1	1	1	1	1	1	1	1	1
3	0	0	0	0	0	0	0	0	0	0	0
4	1	1	1	1	1	1	1	1	1	1	1

The first thing we do when learning a predictive model from data is split the dataset into training and test data. For each, we separate the target variable we want to predict (`y` and `y_test` below), and the variables we want to use to predict (`x` and `x_test` below). We will use our training data to learn our predictive model, and then use the test data to verify its accuracy (it is, of course, no great accomplishment to do well at prediction for the data points the model has already seen in training; you can always just memorize the training data).

Below, we randomly split 1/3 of the data into the test set (`x_test` and `y_test`) and the remaining 2/3 of the data into the training set (`x` and `y`). We do this randomly to ensure that the training and test sets look very similar. You will also note that we convert all of the different datasets into Numpy arrays; `.values` in Pandas converts a Pandas dataframe into a Numpy array. Note that `x` and `x_test` are 2-d arrays, whereas `y` and `y_test` are 1-d arrays. In both cases, `y[i]` corresponds to the row `x[i, :]`. We print the first five rows of `x` and the first five values of `y`.

```
In [74]: # Do not change this code, including the seed
np.random.seed(137486213)
test_indices = np.random.binomial(1, 0.33, df_v.shape[0])
df_v['test'] = test_indices

df_v_test = df_v[df_v['test']==1]
df_v_train = df_v[df_v['test']==0]

x = (df_v_train[['v_1', 'v_2', 'v_3', 'v_4', 'v_5', 'v_6', 'v_7', 'v_8',
'v_9', 'v_10']]).values
x_test = (df_v_test[['v_1', 'v_2', 'v_3', 'v_4', 'v_5', 'v_6', 'v_7', 'v_8',
'v_9', 'v_10']]).values

y = df_v_train['republican'].values
y_test = df_v_test['republican'].values

print(x[0:5,:])
print(y[0:5])

[[0 0 0 0 0 0 0 0 0 0]
 [0 0 0 0 0 0 0 0 0 0]
 [1 1 1 1 1 1 1 1 1 1]
 [0 0 0 0 0 0 0 0 0 0]
 [0 1 0 1 1 0 0 0 0 0]
 [0 0 1 0 1]
```

For a baseline, let's see what kind of error we get if we just guess randomly. Below, we generate a "prediction" for the test data by ignoring the test data altogether (we don't use `x_test` at all below) and just flipping a coin to predict 0 or 1 for each data point. Note that the prediction is also just a 1-d array where the entry at index `i` corresponds to a prediction for `y_test[i]`. As expected, this yields fairly high error (0.5 in expectation). In what follows, you will train a Naive Bayes classifier to improve on this performance.

```
In [5]: random_guess = np.random.randint(2, size=len(y_test))
print(uniform_error(random_guess, y_test))

0.44680851063829785
```

Problem A

Recall how Naive Bayes classifiers work: for a data point `x_test[i,:]` we want to predict the class `j` that maximizes the likelihood: the conditional probability of seeing `x[i,:]` given `j`. We make the assumption that the different features are independent given the class, that is, that we can break this probability up into the product over columns `c` of the probability of `x[i,c]` given `j`, all times the probability that `y=j`.

First, compute all of the conditional probabilities on the training data. That is, for every column `c`, compute the conditional probability for a random representative `i` that `x[i,c]=1` given `y=1`, and the same conditional probability given `y=0` (note that the conditional probability that `x[i,c]=0` given `y=1` is just one minus the probability that `x[i,c]=1` given `y=1`, and the same for `y=0`). You will also need to compute the probability that `y=1` and the probability that `y=0`. Print the resulting probabilities.

```

In [93]: # Write your code for Problem A here
zeros = 0
total = 0
for i in range(0, len(y)):
    total += 1
    if y[i] == 0:
        zeros += 1

prob_y_zero = zeros/total
prob_y_one = 1 - zeros/total
print('y=0: ' + str(100*prob_y_zero) + '%')
print('y=1: ' + str(100*prob_y_one) + '%')

cond_prob_zero_y_zero = {}
cond_prob_zero_y_one = {}
cond_prob_one_y_zero = {}
cond_prob_one_y_one = {}
for c in range(0, len(x[0])):
    cond_zero_y_zero = 0
    cond_zero_y_one = 0
    cond_one_y_zero = 0
    cond_one_y_one = 0
    total_y_zero = 0
    total_y_one = 0
    for i in range(0, len(y)):
        if y[i] == 0:
            total_y_zero += 1
            if x[i,c] == 0:
                cond_zero_y_zero += 1
            else:
                cond_one_y_zero += 1
        elif y[i] == 1:
            total_y_one += 1
            if x[i,c] == 1:
                cond_one_y_one += 1
            else:
                cond_zero_y_one += 1
    cond_prob_zero_y_zero[c] = (cond_zero_y_zero/total_y_zero)# * prob_y
    cond_prob_zero_y_one[c] = (cond_zero_y_one/total_y_one)# * prob_y_on
    cond_prob_one_y_zero[c] = (cond_one_y_zero/total_y_zero)# * prob_y_z
    cond_prob_one_y_one[c] = (cond_one_y_one/total_y_one)# * prob_y_one

print(cond_prob_zero_y_zero)
print(cond_prob_zero_y_one)
print(cond_prob_one_y_zero)
print(cond_prob_one_y_one)

```

```

y=0: 48.29931972789115%
y=1: 51.70068027210885%
{0: 0.6830985915492958, 1: 0.7323943661971831, 2: 0.528169014084507, 3:
0.8873239436619719, 4: 0.7253521126760564, 5: 0.5774647887323944, 6: 0.
5352112676056338, 7: 0.5563380281690141, 8: 0.6549295774647887, 9: 0.69
01408450704225}
{0: 0.3684210526315789, 1: 0.2631578947368421, 2: 0.40789473684210525,
3: 0.013157894736842105, 4: 0.06578947368421052, 5: 0.3421052631578947
5, 6: 0.42105263157894735, 7: 0.42105263157894735, 8: 0.230263157894736
84, 9: 0.4144736842105263}
{0: 0.31690140845070425, 1: 0.2676056338028169, 2: 0.47183098591549294,
3: 0.11267605633802817, 4: 0.2746478873239437, 5: 0.4225352112676056,
6: 0.4647887323943662, 7: 0.44366197183098594, 8: 0.34507042253521125,
9: 0.30985915492957744}
{0: 0.631578947368421, 1: 0.7368421052631579, 2: 0.5921052631578947, 3:
0.9868421052631579, 4: 0.9342105263157895, 5: 0.6578947368421053, 6: 0.
5789473684210527, 7: 0.5789473684210527, 8: 0.7697368421052632, 9: 0.58
55263157894737}

```

Problem B

Now that you have computed the marginals, use them to predict, for every test data point (that is, for each row in `x_test`), whether the given member is a republican (1) or a democrat (0). More concretely, you should compute an array of the same length as `y_test` where for each entry, your prediction is a 1 or 0. Do so by selecting, for each data point, the class (1 or 0) that maximizes the above likelihood under the independence assumption of the Naive Bayes classifier. You do not need to use Laplace smoothing for this example. Once you have your prediction, compute and print the uniform error of your prediction by comparing to `y_test`. For full credit, your classifier should have uniform error less than 0.15 on this particular data.

```

In [94]: # Write your code for Problem B here
predictions = []
for i in range(0, len(x_test)):
    prob_zero = 0
    prob_one = 0
    for c in range(0, len(x_test[0])):
        if x_test[i,c] == 0:
            prob_zero += cond_prob_zero_y_zero[c]
            prob_one += cond_prob_zero_y_one[c]
        else:
            prob_zero += cond_prob_one_y_zero[c]
            prob_one += cond_prob_one_y_one[c]
    prob_one = prob_one * prob_y_zero
    prob_zero = prob_zero * prob_y_zero
    if prob_zero > prob_one:
        predictions.append(0)
    else:
        predictions.append(1)

print(uniform_error(predictions, y_test))

```

```
0.1347517730496454
```

Part 3: Decision Tree Classification

In this part, we will tackle the same prediction task as in Part 2 with the same data. However, instead of a Naive Bayes classifier, we will use a very simple decision tree that simply chooses a single feature on which to split the data.

Problem C

Recall that decision trees split the training data into different sets based on their values for a given feature. Concretely, we will choose a single column c , and split the data into two sets, one for the data points with $x[i, c]=1$ and another for the data points with $x[i, c]=0$. Then, to make a prediction for a given test point, we simply check which of the sets it would go into, and predict the most common class ($y=1$ or $y=0$) among the training data split into that set.

To decide on the best feature to use for splitting, we calculate the information gain of a split as the entropy of the original distribution (over classes $y=1$ and $y=0$) minus the weighted average of the entropy of the two distributions represented by the split. To begin, compute the information gain for every feature (that is, every column of x) in the training data. Which feature (that is, which column) has the greatest information gain?

```

In [95]: def entropy(zeroCount, oneCount, array):
    length = len(array)
    ret = (-(zeroCount/length) * np.log2(zeroCount/length)) + (-(oneCount/length) * np.log2(oneCount/length))
    return ret

original = 0
oneCount = 0
zeroCount = 0
total = len(y)

for i in range(len(y)):
    num = y[i]
    if num == 1:
        oneCount += 1
    else:
        zeroCount += 1
original = entropy(zeroCount, oneCount, x)
best = 0
bestDex = 0
for i in range(9):
    groupOne = []
    groupZero = []
    for j in range(len(x)):
        num = x[j,i]
        if num == 1:
            groupOne.append(j)
        if num == 0:
            groupZero.append(j)
    groupOneOneCount = 0
    groupOneZeroCount = 0
    groupZeroOneCount = 0
    groupZeroZeroCount = 0
    for z in groupOne:
        num = y[z]
        if num == 1:
            groupOneOneCount += 1
        else:
            groupOneZeroCount += 1
    for z in groupZero:
        num = y[z]
        if num == 1:
            groupZeroOneCount += 1
        else:
            groupZeroZeroCount += 1
    groupOneEntropy = entropy(groupOneZeroCount, groupOneOneCount, groupOne)
    groupZeroEntropy = entropy(groupZeroZeroCount, groupZeroOneCount, groupZero)
    infoGain = original - ((groupOneEntropy * (len(groupOne) / len(x))) + (groupZeroEntropy * (len(groupZero) / len(x))))
    if infoGain > best:
        best = infoGain
        bestDex = i
print(best)

```



```
print(bestDex)
```

```
0.6943661833808711
```

```
3
```

Problem D

Now that you have identified the feature with the greatest information gain, split the training data into two sets based on that feature. For each of the two sets, find the most common class ($y=1$ or $y=0$) among that set. Use that information make predictions for the test data as in Part 1. Also as in Part 1, once you have your prediction, compute and print the uniform error of your prediction by comparing to `y_test` . For full credit, your classifier should have uniform error less than 0.15 on this particular data.

```
In [96]: predictY = []
         for i in range(len(x_test)):
             vote = x_test[i, 3]
             if(vote == 1):
                 predictY.append(1)
             if(vote == 0):
                 predictY.append(0)
         print(uniform_error(predictY, y_test))
```

```
0.0851063829787234
```

Part 4: k-Nearest Neighbor Classification

In this part, we will look at a new dataset, `iris.csv` , and explore a different technique for classification, the k-nearest neighbor classifier. This dataset contains measurements of 150 flowers of three different types. The measurements (the first four columns) are all physical measurements in centimeters, and the flower types are designated by 0, 1, and 2 in the rightmost column. Our goal is to predict the flower type from measurements of the physical dimensions of the flower. Note that unlike in Parts 1 and 2, our features are no longer categorical, but instead are numerical. That would require some adapting of our previous methods, but k-nearest neighbor actually works very well with many numerical features all of which are on the same scale.

```
In [98]: df_iris = pd.read_csv('iris.csv')
         df_iris.head()
```

```
Out[98]:
```

	sepal_length	sepal_width	petal_length	petal_width	flower_type
0	5.1	3.5	1.4	0.2	0
1	4.9	3.0	1.4	0.2	0
2	4.7	3.2	1.3	0.2	0
3	4.6	3.1	1.5	0.2	0
4	5.0	3.6	1.4	0.2	0

As before, we split the dataset into training and test datasets, and separate the target y from the data we want to use to predict the target x . Review the discussion in Part 2 for more details.

```
In [99]: # Do not change this code, including the seed
np.random.seed(137486213)
test_indices = np.random.binomial(1, 0.33, df_iris.shape[0])
df_iris['test'] = test_indices

df_train = df_iris[df_iris['test']==0]
df_test = df_iris[df_iris['test']==1]

x = df_train[['sepal_length', 'sepal_width', 'petal_length', 'petal_width']].values
x_test = df_test[['sepal_length', 'sepal_width', 'petal_length', 'petal_width']].values

y = df_train['flower_type'].values
y_test = df_test['flower_type'].values

print(x[0:5,:])
print(y[0:5])

[[4.9 3. 1.4 0.2]
 [4.6 3.1 1.5 0.2]
 [5. 3.6 1.4 0.2]
 [5. 3.4 1.5 0.2]
 [4.4 2.9 1.4 0.2]]
[0 0 0 0 0]
```

Note that the uniform error of a random guess on this dataset is higher than in Parts 2 and 3, because here we are trying to classify into three possible classes (y can be 0, 1, or 2) instead of just two classes.

```
In [100]: random_guess = np.random.randint(3, size=len(y_test))
print(uniform_error(random_guess, y_test))
```

0.6875

Problem E

Recall that the k-nearest neighbor algorithm works by searching for points from the training data that are similar to the point for which we want to make a prediction. We quantify this notion by employing a distance function. The simplest example is the Euclidean distance function (that is, the distance function from your high school geometry and physics) which, given two points u and v in d dimensions, is given by

$$\text{dist}(\mathbf{u}, \mathbf{v}) = \sqrt{\sum_{i=1}^d (u_i - v_i)^2}.$$

Write a function that computes the Euclidean distance between two points represented as Numpy arrays.

```
In [101]: def euc(array1, array2):
    total = 0
    for i in range(len(array1)):
        num1 = array1[i]
        num2 = array2[i]
        diff = num1 - num2
        total += np.power(diff, 2)
    return np.sqrt(total)
```

Problem F

Now, to make predictions for a given point p (for example, $x_test[0,:]$ would be one such point), we first find the k points in the training data set x (note that each row in x is a data point) with minimum distance to the point p . k is a parameter that can be set to different values, but for the purpose of this assignment, let's use $k=5$. We then predict the class ($y=2$, $y=1$, or $y=0$) that is most common among these k points in the training data set.

Use the k-nearest neighbor algorithm to make predictions of the flower type for all of the test data. More concretely, you should compute an array of the same length as y_test where for each entry, your prediction is a 2, 1, or 0. Once you have your prediction, compute and print the uniform error of your prediction by comparing to y_test . For full credit, your classifier should have uniform error less than 0.1 on this particular data.

```
In [102]: closest = []
    predictY = []
    for i in range(len(x_test)):
        array1 = x_test[i,:]
        for j in range(len(x)):
            array2 = x[j,:]
            if(len(closest) < 5):
                closest.append(j)
            else:
                maxi = closest[0]
                for v in closest:
                    tmp = x[v,:]
                    dist = euc(array1, tmp)
                    maxArray = x[maxi,:]
                    if dist > euc(array1,maxArray):
                        maxi = v
                if euc(array2, array1) < euc(array1, x[maxi,:]):
                    closest.remove(maxi)
                    closest.append(j)

    total = 0
    for dex in closest:
        num = y[dex]
        total += num
    average = total / 5
    predictY.append(np rint(average))
    print(uniform_error(predictY, y_test))
```

0.0

Submitting HW 5

1. Double check that you have written all of your answers along with your supporting work in this notebook. Make sure you save the complete notebook.
2. Double check that your entire notebook runs correctly and generates the expected output. To do so, you can simply select Kernel -> Restart and Run All.
3. You will download two versions of your notebook to submit, a .pdf and a .py. To create a PDF, we recommend that you select File --> Download as --> HTML (.html). Open the downloaded .html file; it should open in your web browser. Double check that it looks like your notebook, then print a .pdf using your web browser (you should be able to select to print to a pdf on most major web browsers and operating systems). Check your .pdf for readability: If some long cells are being cut off, go back to your notebook and split them into multiple smaller cells. To get the .py file from your notebook, simply select File -> Download as -> Python (.py) (note, we recognize that you may not have written any Python code for this assignment, but will continue the usual workflow for consistency).
4. Upload the .pdf to gradescope under hw 5 report and the .py to gradescope under hw 5 code. If you work with a partner, only submit one document for both of you, but be sure to add your partner using the [group feature on gradescope \(https://www.gradescope.com/help#help-center-item-student-group-members\)](https://www.gradescope.com/help#help-center-item-student-group-members).

In []: