

# Homework 4

Ayman Tawaalai

February 2024

## 1 Introduction

The purpose of the given assignment is to study and implement Naive Bayes through practical means. The assignment states that the student must one hot encode and ngram the patient sequence data from homework 1. Once that is completed, it will serve as the input features of x. Next, K means clustering must be done and it will serve as the truth labels, y. Once the x and y values are attained, multinomial and gaussian naive bayes will be implemented separately.

## 2 Naive Bayes

Naive Bayes is a supervised machine learning classification model that uses probability. It works by taking in probabilities for predicted outcomes and uses those probabilities to determine if the given inputs are what is the expected output. An example of this can be seen through spam email classification. In this scenario, emails are scanned and each word in the email will weigh the probability of it being real or spam. Whichever has the higher probability will be sorted into the respective category. There are two types of Naive Bayes. The email example most similarly fits with multinomial naive bayes because of the use of discrete variables. Gaussian Naive Bayes focus on continuous variables. It works by getting the probability dense function for each class and then uses it for classification purposes.

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$

Figure 1: Naive Bayes Formula

```
In [103]: datahot = datahot.astype(int)
print(datahot.values)

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

Figure 2: One Hot Conversion

```
In [137]: n = 30
def generate_ngrams(data, n):
    ngrams_list = [tuple(data[i:i+n]) for i in range(len(data)-n+1)]
    return ngrams_list

ngrams = generate_ngrams(patientdata, n)
ngramscdf = pd.DataFrame(ngrams, columns=[f'Ngram_{i+1}' for i in range(n)])
print(ngramscdf)

   Ngram_1 Ngram_2 Ngram_3 Ngram_4 Ngram_5 Ngram_6 Ngram_7 Ngram_8 \
0    58.0   33.0   34.0   62.0   33.0   58.0   33.0   33.0
1    33.0   34.0   62.0   33.0   48.0   58.0   33.0   34.0
2    34.0   62.0   33.0   48.0   58.0   33.0   34.0   33.0
3    62.0   33.0   48.0   58.0   33.0   34.0   33.0   62.0
4    33.0   48.0   58.0   33.0   34.0   33.0   62.0   33.0
...
29296   34.0   33.0   34.0   48.0   58.0   33.0   34.0   33.0
29297   33.0   34.0   33.0   34.0   48.0   58.0   33.0   34.0
29298   34.0   33.0   34.0   48.0   58.0   33.0   34.0   33.0
29299   33.0   34.0   48.0   58.0   33.0   34.0   33.0   34.0
29300   34.0   48.0   58.0   33.0   34.0   33.0   34.0   33.0

   Ngram_9 Ngram_10 ... Ngram_21 Ngram_22 Ngram_23 Ngram_24 \
0    34.0   33.0 ...   33.0   62.0   48.0   33.0
1    33.0   62.0 ...   62.0   48.0   33.0   58.0
2    62.0   33.0 ...   48.0   33.0   58.0   33.0
3    33.0   58.0 ...   33.0   58.0   33.0   34.0
4    58.0   33.0 ...   58.0   33.0   34.0   33.0
...
29296   34.0   33.0 ...   34.0   48.0   58.0   33.0
29297   33.0   34.0 ...   48.0   58.0   33.0   34.0
29298   34.0   33.0 ...   58.0   33.0   34.0   33.0
29299   33.0   34.0 ...   33.0   34.0   33.0   34.0
29300   34.0   62.0 ...   34.0   33.0   34.0   48.0

   Ngram_25 Ngram_26 Ngram_27 Ngram_28 Ngram_29 Ngram_30
0    58.0   33.0   34.0   33.0   62.0   33.0
1    33.0   34.0   33.0   62.0   33.0   48.0
2    34.0   33.0   62.0   33.0   48.0   33.0
3    33.0   62.0   33.0   48.0   33.0   58.0
```

Figure 3: NGram conversion of original patient data

### 3 Implementation

Since the working data is housed in the third column, we will load the files using pandas. Once all the strings are parsed to floats, dummies is used for one hot conversion. The next step will be to conduct ngram as shown in figure three. Since the x values have now been generated it is now time to get the y values. The one hot was sent into the k means function and the user would have to input how many clusters they would like and how many iterations the model should run for. After K means is completed the y values have been generated. With both x and y values, it is time to begin Naive Bayes.

The first implementation will be for multinomial. The main work horse is the fit function. The function will take in both x and y values, split into training sets and testing sets, and uses Laplace smoothing to calculate the class probabilities.

```
import numpy as np
labels = []
def kmeans(x,k, iters):
    centroids = x[np.random.choice(len(x), k, replace=False)]
    for z in range(iters):
        labels = np.argmax(np.linalg.norm(x[:, np.newaxis] - centroids, axis=2), axis=1)
        new_centroids = np.array([x[labels == i].mean(axis=0) for i in range(k)])
        if np.all(centroids == new_centroids):
            break
    centroids = new_centroids
    return labels
y = kmeans(datahot.values, 3, 1000)
print(y)
```

Figure 4: KMeans implementation

```

classprob = counts / amp(y)
featureprob = []

for c in classes:
    labelidx = (y == c)
    classcount = np.sum(labelidx)
    curfeatureprob = (labelidx == 1) / (classcount + x.shape[1])
    featureprob.append(curfeatureprob)

return classprob, featureprob, classes

def predict_naivebayes(x, classprobfinal, featureprobfinal):
    predictions = []
    for xs in x:
        classprobs = np.log(classprobfinal)
        for i, ix in enumerate(classes):
            curfeatureprob = np.log(featureprobfinal[i])
            curprediction = np.sum(classprobs * curfeatureprob)
            classprob[i] += curprediction
        prediction = classes[np.argmax(classprobs)]
        predictions.append(prediction)
    return np.array(predictions)

xtrain, xtest, ytrain, ytest = train_test_split(dataset.values, y, test_size=0.2, random_state=42)
classprobfinal, featureprobfinal, classes = fit(xtrain, ytrain)
print(classprobfinal, featureprobfinal)
pred = predict_naivebayes(xtest, classprobfinal, featureprobfinal)
print(pred)

[1.] [array([[4.25740119e-05, 4.25740119e-05, 4.25740119e-05, ...,
4.25740119e-05, 4.25740119e-05, 4.25740119e-05],
[4.25740119e-05, 4.25740119e-05, 4.25740119e-05, ...,
4.25740119e-05, 4.25740119e-05, 4.25740119e-05],
[4.25740119e-05, 4.25740119e-05, 4.25740119e-05, ...,
4.25740119e-05, 4.25740119e-05, 4.25740119e-05],
[4.25740119e-05, 4.25740119e-05, 4.25740119e-05, ...,
4.25740119e-05, 4.25740119e-05, 4.25740119e-05],
[4.25740119e-05, 4.25740119e-05, 4.25740119e-05, ...,
4.25740119e-05, 4.25740119e-05, 4.25740119e-05],
[4.25740119e-05, 4.25740119e-05, 4.25740119e-05, ...,
4.25740119e-05, 4.25740119e-05, 4.25740119e-05]])]

```

Figure 5: Multinomial Naive Bayes implementation

```

In [10]: M def pdf(x, mean, std):
    varlen = len(x)
    exponent = np.exp(-(x - mean)**2) / (2 * (std**2 * np.pi))
    P = (1 / (varlen**2 * np.pi)) * (std**2 * np.pi) * exponent
    return np.where(P == 0, epsilon, P)

def gaussianfit(x, y):
    classes, counts = np.unique(y, return_counts=True)
    classprob = counts / len(y)
    featuremeans = []
    featuresstds = []

    for c in classes:
        labelidx = (y == c)
        data = x[labelidx]
        featuremeans.append(np.mean(data, axis=0))
        featuresstds.append(np.std(data, axis=0))
    return classes, classprob, featuremeans, featuresstds

def predict_gaussian(x, classes, classprobfinal, featuremeansfinal, featuresstdsfinal):
    predictions = []
    for xs in x:
        classprobs = np.log(classprobfinal)
        for i, ix in enumerate(classes):
            if np.isnan(classprobs[i]) == 0 or np.isnan(featuremeansfinal[i]) == 0 or np.isnan(featuresstdsfinal[i]) == 0:
                # print('Found a 0 at', ix)
                curfeature = np.nan * logpdf(xs, featuremeansfinal[i], featuresstdsfinal[i])
                classprobs[i] += curfeature
            predictionclass = classes[np.argmax(classprobs)]
            predictions.append(predictionclass)
    return np.array(predictions)

# Example usage
# Loading 'dataset', 'xtrain', 'ytrain' are defined
classes, classprobfinal, featuremeansfinal, featuresstdsfinal = gaussianfit(xtrain, ytrain)
print(classprobfinal, "0-----", featuremeansfinal, "x10", featuresstdsfinal)
pred = predict_gaussian(xtest, classes, classprobfinal, featuremeansfinal, featuresstdsfinal)
print(pred)

```

Figure 6: Gaussia Naive Bayes implementation

It places all the feature probabilities into an array and will return that array along with the class probability and classes itself. The predict function will take the final probabilities, along with the test data, and then predict the test labels for each test. Accuracy is then computed using SK learn.

The second implementation will be for gaussian. The same x and y values would be used, but a new fit function would have to be made. The gaussianfit function will return the feature means and standard deviation but is different from the fit function where it will return the feature probabilities. The predict function will take in the finals returned by the fit function and use the pdf function for calculation. The pdf function includes an epsilon to avoid divide by zero errors which were an issue. It will use that calculation to predict the class labels and append that to the predictions array, which is what is returned.

## 4 WEKA

To load the data into weka, some changes were made to the pd loading code. The new changes will now have the third column, along with the key value of

