

Prepared by: Abu Naweem Khan & Sayed Muhammad Saifuddin

Part I: Naïve Bayes for Classification

Objective:

The main goal of this part is to implement the Naïve Bayes algorithm to create a basic classifier.

Background:

In this part of the project, we have implemented a Naïve Bayes algorithm that uses probabilities to perform classification. The probabilities are estimated based on training data which is given as a sample dataset named “fishing.data” for which the value of the classification is known. This dataset consists of the 14 data points (examples). Our target is to find the class of a new instance using the algorithm which has implemented based on the training data. To do this, our algorithm needs to follow the below operations:

Learn (applied to the Training Data)

1. Estimate the probability of each class:
 $P(c_j) = \#c_j / \# \text{training examples}$
2. Estimate the probability of each attribute value a_i , given a class of type j
 $P(a_i | c_j) = \#a_i / \#c_j$

Classify (applied to the new Instance)

1. Return classification C_{NB} for new instance

$$C_{NB} = \max_{c_j \in C} \left(P(c_j) \prod_i P(a_i | c_j) \right)$$

where c_j denotes the class, a_i denotes the value of i^{th} attribute

We have tried to create a generalized function for training with any type of textual dataset like .csv/.tsv and with/without header in file. It will also count first row as name of the attributes if the header is present in the dataset, otherwise it will count first row as instance.

```
def train(self, class_at_column=1, header=True, filepath='', delimiter=','):
    """
    self.class_col = class_at_column - 1

    with open(filepath, mode='r', encoding='utf8') as fin:
        """
        Take the first line and based on header(True/False)
        Create Prior, Likelihood and Variable name lists
        """

        first_line = fin.readline()
        parts = first_line.strip().split(delimiter)

        if header:
            i = 1
            for col in parts:
                if col in self.attribute_names:
                    raise IOError("Column Names not Unique.")
                self.attribute_names.append(col.strip())
                i += 1
        else:
            self.total_instances += 1

            for i in range(1, len(parts) + 1, 1):
                self.attribute_names.append(f'column_{i}')

            for i in range(0, len(parts), 1):
                if i == self.class_col:
                    self.prior_counts[parts[i]] += 1
                    continue

            self.likelihood_counts[
                f'{self.attribute_names[i]}|{parts[i]}|{parts[self.class_col]}'] += 1
```

```
"""
Parse the full file and count all data
"""
for line in fin:
    self.total_instances += 1
    parts = line.strip().split(delimiter)

    for i in range(0, len(parts), 1):
        if i == self.class_col:
            self.prior_counts[parts[i]] += 1
            continue

        self.likelihood_counts[
            f'{self.attribute_names[i]}|{parts[i]}|{parts[self.class_col]}'] += 1

    for key in self.prior_counts:
        self.priors[key] = self.prior_counts[key] / self.total_instances

    for key in self.likelihood_counts:
        cls = key.split('|')[2]
        self.likelihoods[key] = self.likelihood_counts[key] / self.prior_counts[cls]
```

Prepared by: Abu Naweem Khan & Sayed Muhammad Saifuddin

Visualization: [Display the Output]

Output generated using Implemented Algorithm:

```

2 Enter training dataset name: fishing.data
3 Enter new Instance, separated by comma: Strong,WarmAir,Cold,Sunny
4
5
6 *****
7
8 Training Instances: 14
9
10 Prior Probabilities
11 #Yes: 8      P(Yes): 0.5714
12 #No: 6      P(No): 0.4286
13
14 *****
15
16
17 New instance: ['Strong', 'WarmAir', 'Cold', 'Sunny']
18
19 P(Strong|Yes): 0.75 P(WarmAir|Yes): 0.62P(Cold|Yes): 0.12 P(Sunny|Yes): 0.75
20 P(Strong|No): 0.33 P(WarmAir|No): 0.33 P(Cold|No): 0.50 P(Sunny|No): 0.33
21
22 Class probabilities
23
24 Yes: 0.0251
25 No: 0.0079
26
27 Classify: Yes
28 Conditional Probability for class "Yes" : 75.98%
29
30
31 Continue?(Y/n): Y
32 Enter new Instance, separated by comma: Weak,WarmAir,Moderate,Rainy
33
34
35 *****
36
37 Training Instances: 14
38
39 Prior Probabilities
40 #Yes: 8      P(Yes): 0.5714
41 #No: 6      P(No): 0.4286
42
43 *****
44
45
46 New instance: ['Weak', 'WarmAir', 'Moderate', 'Rainy']
47
48 P(Weak|Yes): 0.25 P(WarmAir|Yes): 0.62P(Moderate|Yes): 0.50 P(Rainy|Yes): 0.12
49 P(Weak|No): 0.67 P(WarmAir|No): 0.33 P(Moderate|No): 0.33P(Rainy|No): 0.50
50
51 Class probabilities
52
53 Yes: 0.0056
54 No: 0.0159
55
56 Classify: No
57 Conditional Probability for class "No" : 73.99%
58
59
60 Continue?(Y/n): n

```

Output generated using Scikit-Learn Gaussian NB:

```

61 Enter new Instance, separated by comma: Strong,WarmAir,Cold,Sunny
62 Classify: Yes
63 Continue?(Y/n): Y
64 Enter new Instance, separated by comma: Weak,WarmAir,Moderate,Rainy
65 Classify: No
66 Continue?(Y/n): n
67
68 Process finished with exit code 0

```

Prepared by: Abu Naweem Khan & Sayed Muhammad Saifuddin

Analysis and Discussion:

We can see from the above visualization part that the output generated using the implemented algorithm and the output generated using the Naïve-Bayes module in Scikit-learn are same. So, we can say that our implemented algorithm is validated by a library package. To implement the algorithm, we have used the Python dictionary and there were few problems encountered during the whole implementation process like the dataset had no header whereas there might be header in average number of datasets. In that case we had to consider this to create the generalized function which can be used for any kind of dataset. Another problem was that the class was in the first column in our dataset but in general the class can be in the last column of the average number of datasets. To create the generalized function, we had to consider this issue too.

Future Work:

We can see from our sample dataset that it has only 14 instances which is too small. But if we use any larger dataset for this algorithm, we may get many small likelihoods. And if we multiply those likelihoods together as per our classifier model, arithmetic underflow can happen. Because, in that case, the result of the calculation will be a number smaller than absolute value than the computer can represent as a fixed length (fixed precision) binary digit. As a result, instead of returning the actual result of the calculation, the computer returns a zero. So, if we take the log of each likelihood and then add those instead of multiplying those together, it will increase the absolute value which will reduce the risk of arithmetic underflow.

[Note: You will find the source-code in the last part of this report]

Prepared by: Abu Naweem Khan & Sayed Muhammad Saifuddin

Part II: Naïve Bayes for Spam Classification

Objective:

The main goal of this part is to implement the Naïve Bayes algorithm to create a spam classifier (filter).

Background:

In this part of the project, we have used a sample dataset (textMsgs.data) which consists of a collection of 5574 labeled SMS text messages where about 13% of the messages are spam. Basically, we have implemented a Naïve Bayes algorithm that uses probabilities to perform classification. In part 1 of this project, we used multiplication technique to compute the probability of each class, given the probabilities of the observed data. But, as the data set, we have used is larger in this part, we have used the logarithm technique to compute those probabilities to avoid arithmetic underflow. The probabilities are estimated based on training data that consists of labeled (Spam | Ham) text messages. Then the target is to learn how to classify (or tag) new messages correctly using this Naïve Bayes classifier. To do this, we have used 10-fold cross validation where 90% data will be training data and rest 10% will be test data. In addition, the sample data was not pre-processed. So, we have observed the prediction both before and after data pre-processing. Moreover, we have tried to characterize misclassification by using precision and recall analysis.

Data pre-processing:

To pre-process the dataset, we have converted all words to lowercase, removed all numbers, punctuations and all stopwords such as articles, adjectives, and pronouns.

Visualization: [Display the Output]

Output generated before data pre-processing:

Confusion Matrix		Predict	
		Positive (P)	Negative (N)
Actual	Positive (P)	4827	0
	Negative (N)	332	415

Accuracy	Precision	Recall	Misclassification Rate
94.044%	93.56%	100%	5.96%

Prepared by: Abu Naweem Khan & Sayed Muhammad Saifuddin

```

1 C:\Users\shuvo\AppData\Local\Microsoft\WindowsApps\python3.9.exe "C:\Program Files\JetBrains\PyCharm
  2021.3\plugins\python\helpers\pydev\pydevd.py" --multiproc --qt-support=auto --client 127.0.0.1 --port
  64509 --file C:/Users/shuvo/Workspace/PyCharm/CIS678/main.py
2 Connected to pydev debugger (build 213.5744.248)
3 TP=4827, TN=415, FP=332, FN=0
4
5 Process finished with exit code 0
6

```

Prepared by: Abu Naweem Khan & Sayed Muhammad Saifuddin

Output generated after data pre-processing:

Confusion Matrix		Predict	
		Positive (P)	Negative (N)
Actual	Positive (P)	4825	2
	Negative (N)	181	566

Accuracy	Precision	Recall	Misclassification Rate
96.72%	96.38%	99.96%	10.19%

Prepared by: Abu Naweem Khan & Sayed Muhammad Saifuddin

```

1 C:\Users\shuvo\AppData\Local\Microsoft\WindowsApps\python3.9.exe C:/Users/shuvo/Workspace/PyCharm/CIS678
  /main.py
2 TP=4825, TN=566, FP=181, FN=2
3
4 Process finished with exit code 0
5

```

Analysis and Discussion:

We can see from the above visualization part that the output generated before data pre-processing and after data pre-processing are not the same. Though the accuracy of our spam classifier is not 100% but it has improved after data pre-processing and same goes for precision too. But, the interesting thing is the recall and misclassification rate have declined after data pre-processing.

Future Work:

There are a few other criteria to pre-process the given dataset. For instance, we can pre-process this dataset by substituting the abbreviations with their full forms, data stemming etc. If we can do so, the accuracy of the implemented classifier will definitely improve.

[Please turn the next page for the source-codes of both parts]

```
1 from collections import defaultdict
2 from nltk import corpus
3 import math
4
5
6 class NB:
7     __class_col = None
8     __attribute_names = None
9     __prior_counts = None
10    __priors = None
11    __likelihood_counts = None
12    __likelihoods = None
13    __total_instances = None
14
15    def __init__(self):
16        self.__class_col = None
17        self.__attribute_names = []
18
19        self.__prior_counts = defaultdict(int)
20        self.__priors = defaultdict(float)
21
22        self.__likelihood_counts = defaultdict(float)
23        self.__likelihoods = defaultdict(float)
24
25        self.__total_instances = 0
26
27    def train(self, class_at_column=1, header=True, filepath='', delimiter=','):
28        """
29        Trying to create a Generalized function for training with any
30        type of textual dataset. .csv/.tsv and with/without header in
31        file support.
32
33        :param class_at_column: The column that contains the class.
34        All other columns are considered attributes.
35        Default class at 1st column.
36        :param header: If the attribute names present at the first row of the file.
37        Default is True.
38        :param filepath: The full/relative to .py file location of data set path.
39        :param delimiter: The attribute separator for each line.
40        :return: None
41        """
42        self.__class_col = class_at_column - 1
43
44        with open(filepath, mode='r', encoding='utf8') as fin:
45            """
46            Take the first line and based on header(True/False)
47            Create Prior, Likelihood and Variable name lists
48            """
49
50            lines = fin.readlines()
51            parts = lines[0].strip().split(delimiter)
52
53            if header:
54                i = 1
55                for col in parts:
56                    if col in self.__attribute_names:
57                        raise IOError("Column Names not Unique.")
58                    self.__attribute_names.append(col.strip())
59                    i += 1
60
61            """
62            Parse the full file and count all data
63            """
64            for line in lines:
65                self.__total_instances += 1
66                parts = line.strip().split(delimiter)
67
68                cls = parts[self.__class_col].strip()
69
70                for i in range(0, len(parts), 1):
```

```

71         if i == self.__class_col:
72             self.__prior_counts[cls] += 1
73             continue
74
75         self.__likelihood_counts[f'{parts[i].strip()}|{cls}'] += 1
76
77     for cls in self.__prior_counts:
78         self.__priors[cls] = self.__prior_counts[cls] / self.__total_instances
79
80     for key in self.__likelihood_counts:
81         cls = key.split('|')[-1]
82         self.__likelihoods[key] = self.__likelihood_counts[key] / self.__prior_counts[cls]
83
84     """
85     For tabular data, if any attribute is previously unseen, we assign probability to 0.
86     For textual data, if any word is previously unseen, we assign a
87     probability using - 1/(total #word + size of vocab)
88
89     This helps keep the test code same for both tabular and textual data.
90     """
91     self.__likelihoods['_any_|_any_'] = 0
92
93 def test(self, new_instance=None, no_print=False):
94     if new_instance is None:
95         raise AttributeError(f"Invalid New Instance.")
96
97     if not no_print: print('\n\n*****\n')
98     if not no_print: print(f'Training Instances: {self.__total_instances}\n')
99     if not no_print: print('Prior Probabilities')
100
101     for cls in self.__priors:
102         if not no_print: print(f'#{cls}: {self.__prior_counts[cls]}\t\t', end='')
103         if not no_print: print(f'P({cls}): {self.__priors[cls]}\t\t', end='')
104         if not no_print: print(f'log(P({cls})): {math.log(self.__priors[cls])}\t\t')
105
106     if not no_print: print('\n*****\n\n')
107     if not no_print: print(f'New instance: {str(new_instance)}\n')
108
109     posteriors = {}
110
111     for cls in self.__priors:
112         posteriors[cls] = self.__priors[cls]
113         posteriors[f'log({cls})'] = math.log(self.__priors[cls])
114
115         for attr in new_instance:
116             key = f'{attr}|{cls}'
117
118             if key in self.__likelihoods:
119                 lkh = self.__likelihoods[key]
120             else:
121                 lkh = self.__likelihoods['_any_|_any_']
122
123             posteriors[cls] *= lkh
124             posteriors[f'log({cls})'] += math.log(lkh)
125             if not no_print: print(f'P({attr}|{cls}): {lkh}\t\t', end='')
126             if not no_print: print(f'log(P({attr}|{cls})): {math.log(lkh)}\t\t', end='')
127
128         if not no_print: print('')
129
130     if not no_print: print('\nClass probabilities\n')
131     max_prob = -99999999
132     verdict = ''
133     total = 0
134     for cls in posteriors:
135         if not cls.startswith('log'):
136             total += posteriors[cls]
137             if not no_print: print(f'{cls}: {posteriors[cls]}\t\t', end='')
138             continue
139
140     if not no_print: print(f'{cls}: {posteriors[cls]}\t\t')

```

```

141
142         if max_prob < posteriors[cls]:
143             max_prob = posteriors[cls]
144             verdict = cls
145
146         if not no_print: print(f'\nClassify: {verdict[4:-1]}')
147
148         if total > 0:
149             if not no_print: print(f'Conditional Probability for class "{verdict[4:-1]}" : {posteriors[
verdict[4:-1]] * 100 / total:.2f}%\n\n')
150
151         return verdict[4:-1]
152
153     def __fix_likelihoods_for_text(self):
154         size_of_vocab = len(self.__likelihood_counts)
155         total_word_count = 0
156         for value in self.__likelihood_counts.values():
157             total_word_count += value
158
159         for key in self.__likelihoods:
160             self.__likelihoods[key] = (self.__likelihood_counts[key] + 1) / (total_word_count +
size_of_vocab)
161
162         self.__likelihoods['_any_|_any_'] = 1 / (total_word_count + size_of_vocab)
163
164     def train_for_text(self, datafile):
165         lines = []
166         max_len = 0
167         # total_instances_by_cls = defaultdict(int)
168         class_at_column = 1
169         stop_words = set(corpus.stopwords.words('english'))
170
171         with open(datafile, mode='r', encoding='utf8') as fin:
172             for line in fin:
173                 # case-insensitive prediction
174                 line = line.strip().lower()
175
176                 # removing stopwords
177                 # words = word_tokenize(line)
178                 # valid_words = []
179                 # for w in words:
180                 #     if w not in stop_words:
181                 #         valid_words.append(w)
182                 #
183                 # if max_len < len(valid_words):
184                 #     max_len = len(valid_words)
185
186                 line = line.replace(',', ' ')
187                 valid_words = line.split()
188
189                 # total_instances_by_cls[valid_words[class_at_column - 1]] += 1
190
191                 lines.append(valid_words)
192
193         newfile = 'processed_texts.csv'
194         with open(newfile, mode='w', encoding='utf8') as fout:
195
196             # for cls in total_instances_by_cls:
197             #     fout.write(f'#{cls}_{total_instances_by_cls[cls]}\n')
198
199             for line in lines:
200                 length = len(line)
201                 for w in line:
202                     fout.write(f'{w},')
203
204                 for i in range(0, max_len - length, 1):
205                     fout.write(',')
206
207                 fout.write('\n')
208

```


Prepared by: Abu Naweem Khan & Sayed Muhammad Saifuddin

```
209         self.train(class_at_column=class_at_column, header=False, filepath=newfile, delimiter=',')
210         self.__fix_likeliheids_for_text()
211
```

```
1 from sklearn.naive_bayes import GaussianNB
2 import numpy as np
3 from nltk import word_tokenize
4 from nltk import corpus
5 from project2.naive_bayes import NB
6 from collections import defaultdict
7
8 def sampling(filename):
9     """
10    10-fold cross-validation files. each training set will have 90% of each
11    class.
12    :param filename: dictionary of lists containing filenames for each class
13    :return:
14    """
15    hams = []
16    spams = []
17    with open(filename, mode='r', encoding='utf8') as fin:
18        for line in fin:
19            if line.startswith('ham'):
20                hams.append(line.strip())
21            else:
22                spams.append(line.strip())
23
24    filenames = {
25        'ham': [],
26        'spam': []
27    }
28
29    len_train_ham = int(len(hams) * 0.9)
30    len_test_ham = len(hams) - len_train_ham
31    index = 0
32    while index < len(hams):
33        file_train = f'ham_train_{index}.txt'
34        file_test = f'ham_test_{index}.txt'
35        filenames['ham'].append((file_train, file_test))
36
37        upto = index + len_test_ham
38        if upto > len(hams):
39            upto = len(hams)
40
41        with open(file_train, mode='w', encoding='utf8') as fout_train, open(
42            file_test, mode='w', encoding='utf8') as fout_test:
43            for i in range(0, index, 1):
44                fout_train.write(f'{hams[i]}\n')
45
46            for i in range(index, upto, 1):
47                fout_test.write(f'{hams[i]}\n')
48
49            for i in range(upto, len(hams), 1):
50                fout_train.write(f'{hams[i]}\n')
51
52        index += len_test_ham
53
54    len_train_spam = int(len(spams) * 0.9)
55    len_test_spam = len(spams) - len_train_spam
56    index = 0
57    while index < len(spams):
58        file_train = f'spam_train_{index}.txt'
59        file_test = f'spam_test_{index}.txt'
60        filenames['spam'].append((file_train, file_test))
61
62        upto = index + len_test_spam
63        if upto > len(spams):
64            upto = len(spams)
65
66        with open(file_train, mode='w', encoding='utf8') as fout_train, open(
67            file_test, mode='w', encoding='utf8') as fout_test:
68            for i in range(0, index, 1):
69                fout_train.write(f'{spams[i]}\n')
70
```

```
71         for i in range(index, upto, 1):
72             fout_test.write(f'{spams[i]}\n')
73
74         for i in range(upto, len(spams), 1):
75             fout_train.write(f'{spams[i]}\n')
76
77         index += len_test_spam
78
79     return filenames
80
81
82 def textual_data(filename):
83     train_test_files = sampling(filename)
84
85     TP = TN = FP = FN = 0
86
87     for i in range(0, len(train_test_files['ham']), 1):
88         train_ham, test_ham = train_test_files['ham'][i]
89         train_spam, test_spam = train_test_files['spam'][i]
90         train_filename = 'training_set.txt'
91         test_filename = 'test_set.txt'
92         with open(train_filename, mode='w', encoding='utf8') as fout_train, open(
93             test_filename, mode='w', encoding='utf8') as fout_test, open(
94                 train_ham, mode='r', encoding='utf8') as fin_train_ham, open(
95                     train_spam, mode='r', encoding='utf8') as fin_train_spam, open(
96                         test_ham, mode='r', encoding='utf8') as fin_test_ham, open(
97                             test_spam, mode='r', encoding='utf8') as fin_test_spam:
98
99             trains = fin_train_ham.readlines()
100             trains.extend(fin_train_spam.readlines())
101             test = fin_test_ham.readlines()
102             test.extend(fin_test_spam.readlines())
103
104             for line in trains:
105                 fout_train.write(f'{line.strip()}\n')
106
107             for line in test:
108                 fout_test.write(f'{line.strip()}\n')
109
110     nb = NB()
111     nb.train_for_text(datafile=train_filename)
112
113     with open(test_filename, mode='r', encoding='utf8') as fin:
114         for line in fin:
115             new_instance = line.strip()
116
117             if len(new_instance) == 0:
118                 continue
119
120             # lowercase trained model, so lowercase test
121             new_instance = new_instance.lower()
122
123             words = new_instance.split()
124             # do pre-processing
125
126             actual_cls = words[0]
127             prediction_cls = nb.test(words[1:], no_print=True)
128
129             if actual_cls == 'ham' and prediction_cls == 'ham':
130                 TP += 1
131             elif actual_cls == 'spam' and prediction_cls == 'spam':
132                 TN += 1
133             elif actual_cls == 'ham' and prediction_cls == 'spam':
134                 FN += 1
135             elif actual_cls == 'spam' and prediction_cls == 'ham':
136                 FP += 1
137
138     print(f'TP={TP}, TN={TN}, FP={FP}, FN={FN}')
```

```
141
142 def tabular_data():
143     filename = input("Enter training dataset name: ")
144     nb = NB()
145     nb.train(class_at_column=1, header=False, filepath=filename, delimiter=' ')
146
147     while True:
148         new_instance = input("Enter new Instance, separated by comma: ")
149         nb.test(new_instance.split(','))
150
151         s = input("Continue?(Y/n): ")
152         if s == 'n':
153             break
154
155
156 def gaussian_nb():
157     X = []
158     Y = []
159     classes = {
160         'Yes': 1,
161         'No': 0
162     }
163     wind = {
164         'Strong': 0,
165         'Weak': 1
166     }
167     air = {
168         'WarmAir': 0,
169         'ColdAir': 1
170     }
171     water = {
172         'Warm': 0,
173         'Moderate': 1,
174         'Cold': 2
175     }
176     sky = {
177         'Sunny': 0,
178         'Cloudy': 1,
179         'Rainy': 2
180     }
181     with open('fishing.data', mode='r') as fin:
182         for line in fin:
183             parts = line.strip().split(' ')
184             Y.append(classes[parts[0]])
185
186             X.append([wind[parts[1]], air[parts[2]], water[parts[3]], sky[parts[4]]])
187
188     classifier = GaussianNB()
189     classifier.fit(nb.array(X), nb.array(Y))
190
191     while True:
192         new_instance = input("Enter new Instance, separated by comma: ")
193         parts = new_instance.split(',')
194         new_instance = [[wind[parts[0]], air[parts[1]], water[parts[2]], sky[parts[3]]]]
195         c = classifier.predict(new_instance)
196
197         if c == classes['Yes']:
198             print(f'Classify: Yes')
199         else:
200             print(f'Classify: No')
201
202         s = input("Continue?(Y/n): ")
203         if s == 'n':
204             break
205
```