

University of Waterloo

CS 486, Winter 2024

Assignment 1

Part a

1) Printout of the code:

```
1  from math import log2
2  from queue import PriorityQueue
3  import matplotlib.pyplot as plt
4  import decimal

6  DEBUG = False

8  ##### Variables #####
9  NUM_FEATURES = 0 # Define later
10 NUM_SAMPLE = 1500
11 ATHEISM_ID = 1
12 BOOKS_ID = 2
13 LABEL_STR = {ATHEISM_ID: "Atheism", BOOKS_ID: "Books" }

15 # We define  $P(x)$  as # belongs to atheism over total.

17 ##### READING FILE #####

19 # Function to read label data, returning a dictionary mapping document ID to label
20 def read_label_data(file_name):
21     label_data = {}
22     with open(file_name, 'r') as file:
23         lineNum = 0
24         for line in file:
25             lineNum += 1
26             label_data[lineNum] = int(line.strip())
27             # docId = line number
28     file.close()
29     return label_data

31 # Function to read word data, returning a dictionary mapping word ID to word
32 def read_word_data(file_name):
33     word_data = {}
34     with open(file_name, 'r') as file:
35         lineNum = 0
36         for line in file:
37             lineNum += 1
38             word_data[int(lineNum)] = line.strip()
39             # wordId = line number
40     # Define the Number of features
41     global NUM_FEATURES
42     NUM_FEATURES = lineNum
43     file.close()
44     return word_data

46 # Function to read train or test data
47 # document_data[n] = array of word_id that n+1'th doc have.
48 def read_document_data(file_name):
49     document_data = {i: [] for i in range(1, NUM_SAMPLE + 1)} # 1 to 1500
50     with open(file_name, 'r') as file:
51         for line in file:
52             doc_id, word_id = line.strip().split()
53             document_data[int(doc_id)].append(int(word_id))
54     file.close()
55     return document_data

57 ##### TESTING #####

59 # 0.5 as we are ML
60 def predict_class(document, Theta_i_atheism, Theta_i_books, theta_atheism=0.5, theta_books=0.5):
61     # Start with the theta_c
62     prob_atheism = decimal.Decimal(theta_atheism)
63     prob_books = decimal.Decimal(theta_books)

65     # multiply all the likelihood
66     for word_id in range(1, NUM_FEATURES + 1):
67         if word_id in document:
68             prob_atheism *= decimal.Decimal(Theta_i_atheism[word_id])
69         else:
```

```

70         prob_atheism == decimal.Decimal(1 - Theta_i_atheism[word_id])

72     if word_id in document:
73         prob_books == decimal.Decimal(Theta_i_books[word_id])
74     else:
75         prob_books == decimal.Decimal(1 - Theta_i_books[word_id])

77     normalized_prob_atheism = decimal.Decimal(prob_atheism / decimal.Decimal(prob_atheism + prob_books))
78     # Return the class with the highest posterior probability
79     if normalized_prob_atheism >= 0.5:
80         return ATHEISM_ID
81     else:
82         return BOOKS_ID

84 ##### MAIN #####

86 if __name__ == '__main__':

88     # Reading data from files
89     words = read_word_data('./dataset/words.txt')
90     train_data = read_document_data('./dataset/trainData.txt')
91     train_labels = read_label_data('./dataset/trainLabel.txt')
92     test_data = read_document_data('./dataset/testData.txt')
93     test_labels = read_label_data('./dataset/testLabel.txt')

95     ### First calculate relative frequency of books belong to subreddit 'Atheism' or 'Books' ###
96     total_count = 0
97     atheism_count = 0
98     books_count = 0

100     for _, label in train_labels.items():
101         total_count += 1
102         if label == ATHEISM_ID:
103             atheism_count += 1
104         elif label == BOOKS_ID:
105             books_count += 1

107     theta_atheism = atheism_count / total_count
108     theta_books = books_count / total_count

110     if DEBUG:
111         print(f"Total_number_of_works: {NUM_FEATURES}")
112         print(f"Total_number_of_documents: {total_count}")
113         print(f"Number_of_documents_labeled_as 'Atheism': {atheism_count}")
114         print(f"Number_of_documents_labeled_as 'Books': {books_count}")
115         print(f"Prior_probability_of 'Atheism' class: {theta_atheism:.4f}")
116         print(f"Prior_probability_of 'Books' class: {theta_books:.4f}")

118     ### Split train data for each label ###
119     atheism_train_data = {}
120     books_train_data = {}

122     # Split the train_data based on the labels
123     for doc_id, label in train_labels.items():
124         if label == ATHEISM_ID:
125             atheism_train_data[doc_id] = train_data[doc_id]
126         elif label == BOOKS_ID:
127             books_train_data[doc_id] = train_data[doc_id]

129     ### calculate number of document (value) that have feature word_id (key) ###

131     atheism_word_counts = {i: 0 for i in range(1, NUM_FEATURES + 1)}
132     for i in range(1, NUM_FEATURES + 1):
133         # count the occurrence of word_id = i
134         for _, word_ids in atheism_train_data.items():
135             if i in word_ids:
136                 atheism_word_counts[i] += 1

138     books_word_counts = {i: 0 for i in range(1, NUM_FEATURES + 1)}
139     for i in range(1, NUM_FEATURES + 1):
140         # count the occurrence of word_id = i
141         for _, word_ids in books_train_data.items():
142             if i in word_ids:
143                 books_word_counts[i] += 1

145     ### Account for Laplace correction when calculating the actual theta_i/0 ###

147     Theta_i_atheism = {i: 0 for i in range(1, NUM_FEATURES + 1)}
148     Theta_i_books = {i: 0 for i in range(1, NUM_FEATURES + 1)}

150     for i in range(1, NUM_FEATURES + 1):
151         Theta_i_atheism[i] = (atheism_word_counts[i] + 1) / (atheism_count + 2)
152         Theta_i_books[i] = (books_word_counts[i] + 1) / (books_count + 2)

154     ##### PROCESSING FINISHED #####

156     print("")

```

```

158     #### 10 most discriminative word features ####
160     # Compute the log probabilities for each word
161     log_prob_diffs = {}
162     for word_id in range(1, NUM_FEATURES + 1):
163         log_prob_atheism = log2(Theta_i_atheism[word_id])
164         log_prob_books = log2(Theta_i_books[word_id])
165         log_prob_diffs[word_id] = abs(log_prob_atheism - log_prob_books)
167     # Sort the words by the most discriminative features
168     most_discriminative = sorted(log_prob_diffs, key=log_prob_diffs.get, reverse=True)[:10]
170     # Print the 10 most discriminative words
171     print("10_Most_Discriminative_Word_Features:")
172     for word_id in most_discriminative:
173         word = words[word_id]
174         if DEBUG:
175             print(f"Word:_{word},_Difference_in_Log_Probabilities:_{log_prob_diffs[word_id]:.4f}")
176         else:
177             print(f"Word:_{word}")
179     print("")
181     #### Predict the class for each document in the training set ####
182     predicted_labels_train = {}
183     for doc_id, document in train_data.items():
184         predicted_labels_train[doc_id] = predict_class(document, Theta_i_atheism, Theta_i_books)
186     # Calculate the accuracy
187     correct_predictions_train = 0
188     for doc_id, prediction in predicted_labels_train.items():
189         if train_labels[doc_id] == prediction:
190             correct_predictions_train += 1
192     accuracy_train = correct_predictions_train / len(train_labels)
194     print(f'Training_accuracy_of_the_Naive_Bayes_classifier():_{accuracy_train*100:.2f}%')
196     #### Predict the class for each document in the test set ####
197     predicted_labels_test = {}
198     for doc_id, document in test_data.items():
199         predicted_labels_test[doc_id] = predict_class(document, Theta_i_atheism, Theta_i_books)
201     # Calculate the accuracy
202     correct_predictions_test = 0
203     for doc_id, prediction in predicted_labels_test.items():
204         if train_labels[doc_id] == prediction:
205             correct_predictions_test += 1
207     accuracy_test = correct_predictions_test / len(test_labels)
209     print(f'Testing_accuracy_of_the_Naive_Bayes_classifier():_{accuracy_test*100:.2f}%')
211     print("")

```

2) 10 most discriminative word features (from most to least, break ties randomly):

christian, religion, atheism, books, christians,
library, religious, libraries, novel, beliefs

In my opinion, these words are indeed very good word features to classify the originated sub-reddit for each posts. All of these words are strongly associated to one of the topic "atheism" (like the words christian, religion, atheism, christians, religious and beliefs) or "books" (the words books, library, libraries and novel).

3) Training accuracy of the Naive Bayes classifier (): 90.27%
Testing accuracy of the Naive Bayes classifier (): 74.47%

4) The naive Bayes model's assumption on the independency of words is not reasonable. Because for natural languages that we use, it has a structure with syntax and semantics where certain words are more likely to appear together.

For example, if a text has word **example**, then the likelihood of the word **for** appears in the same text is a lot higher, because **for example** is a popular phrase that people use.

Also, the meaning of a word can be contextdependent in a lot cases. Words can have different meanings in different contexts and even the presence of certain words can influence the interpretation of other words.

These correlation between words means that the naive Bayes model's assumption on the words are not exactly reasonable, but it's a good simplification that works quite well in practice for certain applications such as spam detection and preliminary topic classification.

- 5) There are several things we can do to extend the Naive Bayes model to take into account of dependencies between words, they include but not limit to:
 - Extend the naive Bayes model into a more complex Bayesian Network that have edges between some word feature nodes will allow a word feature to be dependent to some other word features, which would allow the model to account for some dependencies between word features.
 - Instead of use individual words as features with Naive Bayes model, we can instead use a sequence of few words (n -gram, note that our word features is where $n = 1$) as features instead. This will allow the model to take into account some dependencies between words that often appear together.
- 6) If I were to use MAP learning instead of ML learning under the setting of continuous parameter θ , then I will need to calculate the most likely posterior hypothesis $\theta^* = \operatorname{argmax}_{\theta} P(d|\theta)P(\theta)$, where $P(\theta)$ is the prior of the hypothesis (every parameter has a prior), instead of what ML required, which is the hypothesis that makes the data most likely $\theta = \operatorname{argmax}_{\theta} P(d|\theta)$ (where we assumed a uniform prior, similar as what I did in my function `predict_class()`, so we just omitted it), then take derivative of the RHS making it equal 0. After this, we find the most likely posterior hypothesis θ , then make prediction use this value along with Laplace correction in a similar manner.

Part b

The content of Q1 is submit to marmoset.

1) Code for k-fold cross validation with $k = 5$:

```
1  import numpy as np
2  import pandas as pd
3  import matplotlib.pyplot as plt

5  from neural_net import NeuralNetwork
6  from operations import *

8  from sklearn.model_selection import KFold

10 def load_dataset(csv_path, target_feature):
11     dataset = pd.read_csv(csv_path)
12     t = np.expand_dims(dataset[target_feature].to_numpy().astype(float), axis=1)
13     X = dataset.drop([target_feature], axis=1).to_numpy()
14     return X, t

16 X, y = load_dataset("data/wine_quality.csv", "quality")

18 n_features = X.shape[1]

20 # Hyperparameters and settings
21 n_splits = 5
22 epochs = 500
23 learning_rate = 0.001

25 # Initialize k-fold cross-validation
26 kf = KFold(n_splits=n_splits, shuffle=True, random_state=None)

28 # Initialize variables to store metrics
29 fold_losses = []
30 fold_maes = []
31 epoch_losses = np.zeros((n_splits, epochs))

33 # neural network architecture
34 layer_sizes = [32,32,16,1]
35 activations = [ReLU(), ReLU(), Sigmoid(), Identity()]

38 # K-Fold Cross Validation
39 for fold, (train_idx, val_idx) in enumerate(kf.split(X)):
40     print(f"Training on fold_{fold+1}/{n_splits}...")

42     # Split data into training and validation sets
43     X_train, X_val = X[train_idx], X[val_idx]
44     y_train, y_val = y[train_idx], y[val_idx]

46     # Initialize the neural network
47     nn = NeuralNetwork(n_features=n_features,
48                        layer_sizes=layer_sizes,
49                        activations=activations,
50                        loss=MeanSquaredError(),
51                        learning_rate=learning_rate)

53     # Train the neural network
54     W, epoch_loss = nn.train(X_train, y_train, epochs=epochs)
55     epoch_losses[fold, :] = epoch_loss

57     # Evaluate on validation set
58     mae = nn.evaluate(X_val, y_val, mean_absolute_error)
59     fold_maes.append(mae)
60     print(f"Fold_{fold+1}_MAE: {mae}")

62 # Average training loss over folds for each epoch
63 average_epoch_losses = epoch_losses.mean(axis=0)
64 std_deviation_maes = np.std(fold_maes)
65 average_mae = np.mean(fold_maes)

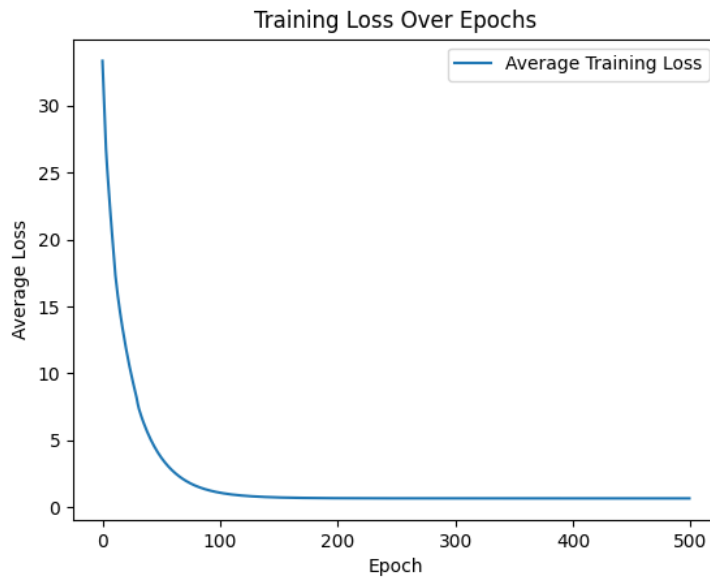
67 # Output the results
68 print(f"Average_MAE_over_all_folds: {average_mae}")
69 print(f"Standard_Deviation_of_MAE_over_all_folds: {std_deviation_maes}")

71 # Plot the results
72 plt.plot(range(epochs), average_epoch_losses, label='Average_Training_Loss')
73 plt.xlabel('Epoch')
74 plt.ylabel('Average_Loss')
75 plt.title('Training_Loss_Over_Epochs')
76 plt.legend()
77 plt.show()
```

2) Size of layers: [32,32,16,1]

Layer used in my neural network: [ReLU(), ReLU(), Sigmoid(), Identity()]
(Same as the one given in `neural_net.py`)

3) Plot:



4) Average MAE over all folds: 0.6851463920029073

Standard Deviation of MAE over all folds: 0.01271361092488064