

Ahmed Awadallah

Text Classification using Clustering and Naive Bayes

Introduction

Employees often use up a large portion of their work time parsing through and responding to important emails. Usually inboxes are also filled with another type of email called spam. The spam emails are generally emails that are sent to a large group of people to solicit something. Spam can be parsed through but doing this manually takes precious time and concentration out of someone's day. Occasionally spam emails can be much more malicious and contain code or links that can cause damage to the computer or company as a whole. Thus there is a vested interest in developing tools that can sift through emails automatically and discard those that are spam while letting the real emails through.

For this case study we are using the spam assassin data. It is a collection of 9348 emails where some are spam and some are not. The goal is to develop a classification model that can distinguish these two classes.

Methods

The spam assassin dataset is a collection of 9348 that are distributed among various folders. To collect that data the files all must be read and recorded within the program using an os walk. Some emails are not simply text but contain a collection of email components. These are a variety of these email types but they all fall under multipart email. As these multipart emails are connected the decision we made was to consider multipart emails as one larger email. This means that for multipart emails we simply combine all the body text of each partner email into one string.

No text preprocessing was performed on any email body. This may look lazy but it appears the tools from sklearn (the countvectorizer) seem to automatically do this. If we are somehow wrong about this an additional reason is that the model performance later is impressive. Because the model performance is good I'd argue against over engineering the problem. If later with new data the model performs significantly worse then this decision should be revisited.

The text for each email contains information about each email. For this project we will be using a bag of words approach to derive meaning from the text. The approach we chose was a TfidfVectorizer as this will help account for words that are frequent but don't help identify a document. For example the word "the" is already used quite a few times in this case study, it is not likely the word "the" would be helpful in identifying this text as belonging to a case study. The same logic is being applied to the emails.

In addition to using the TfidfVectorizer we also set the hyperparameter min_df to 0.002. The min_df hyperparameter lets the vectorizer ignore words below a certain frequency threshold. We set this hyperparameter as incredibly rare words are not generalizable to new emails and thus shouldn't be included. For example the word "thisiscasestudy3forsmu" is unlikely to appear in any other text and could be used to identify this document as a case study. At the same time since no other document would have it would only be useful for this specific case study and thus not generalizable. Our choices lead to the bag of words to contain 8660 instead of 123781.

We chose 0.002 through an inspection method. We would lower the threshold and inspect the words that were recorded. As soon as we saw the appearance of words that were absurd or nonsensical we would revert to the last value we tried. An example of a word that lead to us deciding to not use 0.001 was

"dqoncjxozwfkpg0kpg1ldgegahr0cc1lcxvpdj0iq29udgvudc1myw5ndwfnzsigy29u".

Once we had our word frequency feature matrix we used this as input into a KMeans model that expected two clusters. Two clusters was specified as this problem is a binary problem where we are predicting spam and not spam.

The clusters outputted by the KMeans model were added to the training set as an additional feature now leading us to have 8661 features.

This larger feature set was then used in a MultinomialNB model. Spam will be considered the positive class. We optimized the alpha hyperparameter using 5 fold cross validation and f1 score as an optimizing metric.

F1 score was chosen as we wanted a balanced performance of predicting an email as spam when it's not and predicting an email as spam when it is. It was not specified what the business wanted and there is no clear reason to pick one over the other as the cost is about equal. Misclassifying a couple real emails leads to confusion and frustration at a lack of a response. Misclassifying a couple spam emails leads people to manually check the emails for spam occasionally. This has an added risk of malware or social engineering to occur. Both situations are not ideal but it's hard to say in general for all cases we prefer one over the other. This is why we opted for the balanced approach.

Finally when the optimal alpha was chosen a final model was made. Analysis of the models performance was done through a classification report, confusion matrix, and roc curve.

Results

The following is the results of the KMeans clustering on the data. While this is not a model in and of itself it's interesting to see its performance in classifying the two classes

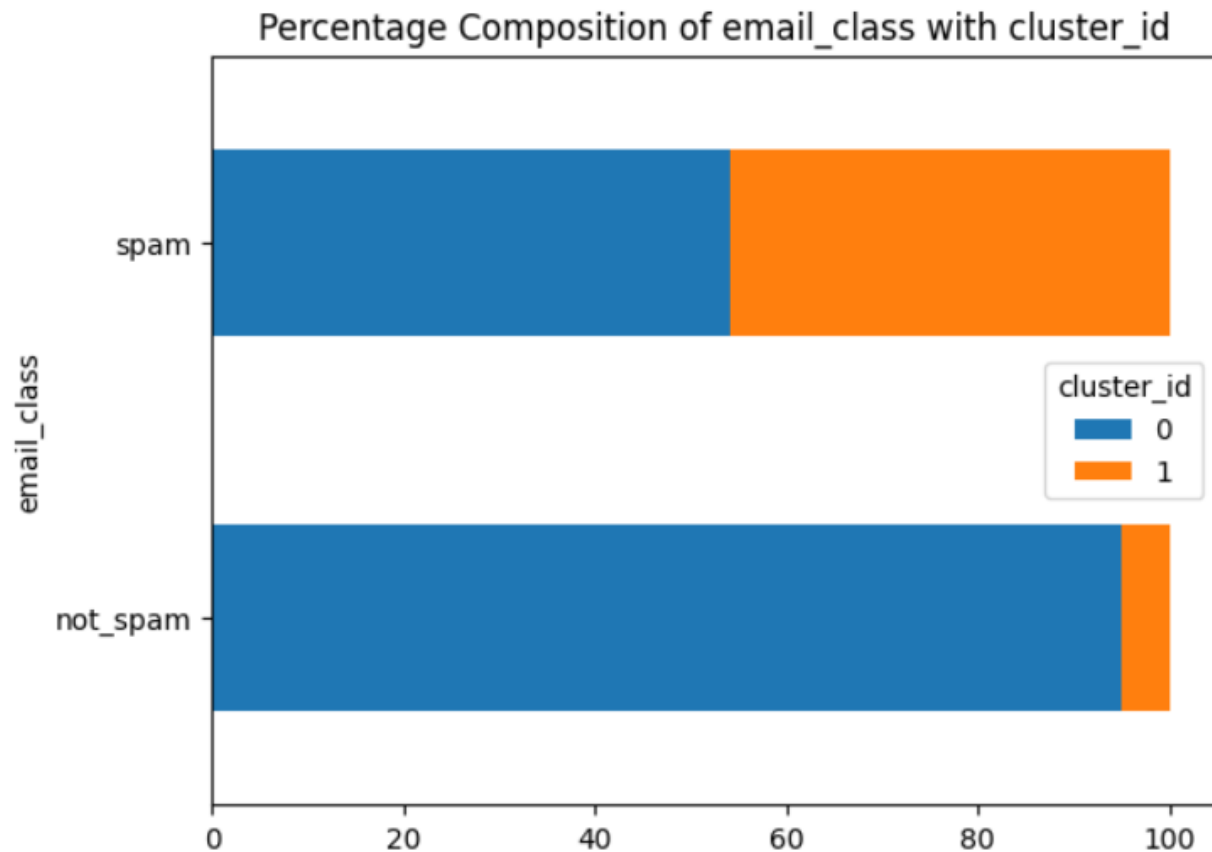


Figure 1: Plot showing how well clusters associated with the spam and not spam classes

Generally the clustering of KMeans did not do a great job separating spam and not spam as the two bars are mixed between the two cluster ids. What can be said is that cluster id 0 is likely associated with not spam and therefore cluster 1 is associated with spam. While we don't need to do this analysis to feed the clustering results into our next model it is still good to see some information about what it means to be spam versus not spam is contained within cluster id.

Remember that our final model is a Multinomial Naive Bayes model. The classification report of our final version of this model is displayed below.

Classification Report					
		precision	recall	f1-score	support
	not_spam	0.98	0.96	0.97	6951
	spam	0.89	0.94	0.91	2397
	accuracy			0.95	9348
	macro avg	0.94	0.95	0.94	9348
	weighted avg	0.96	0.95	0.96	9348

Figure 2: Classification Report of MultinomialNB Model

As can be seen in the report the overall performance of the model is impressive. The two classes can easily be distinguished with almost 100% accuracy.

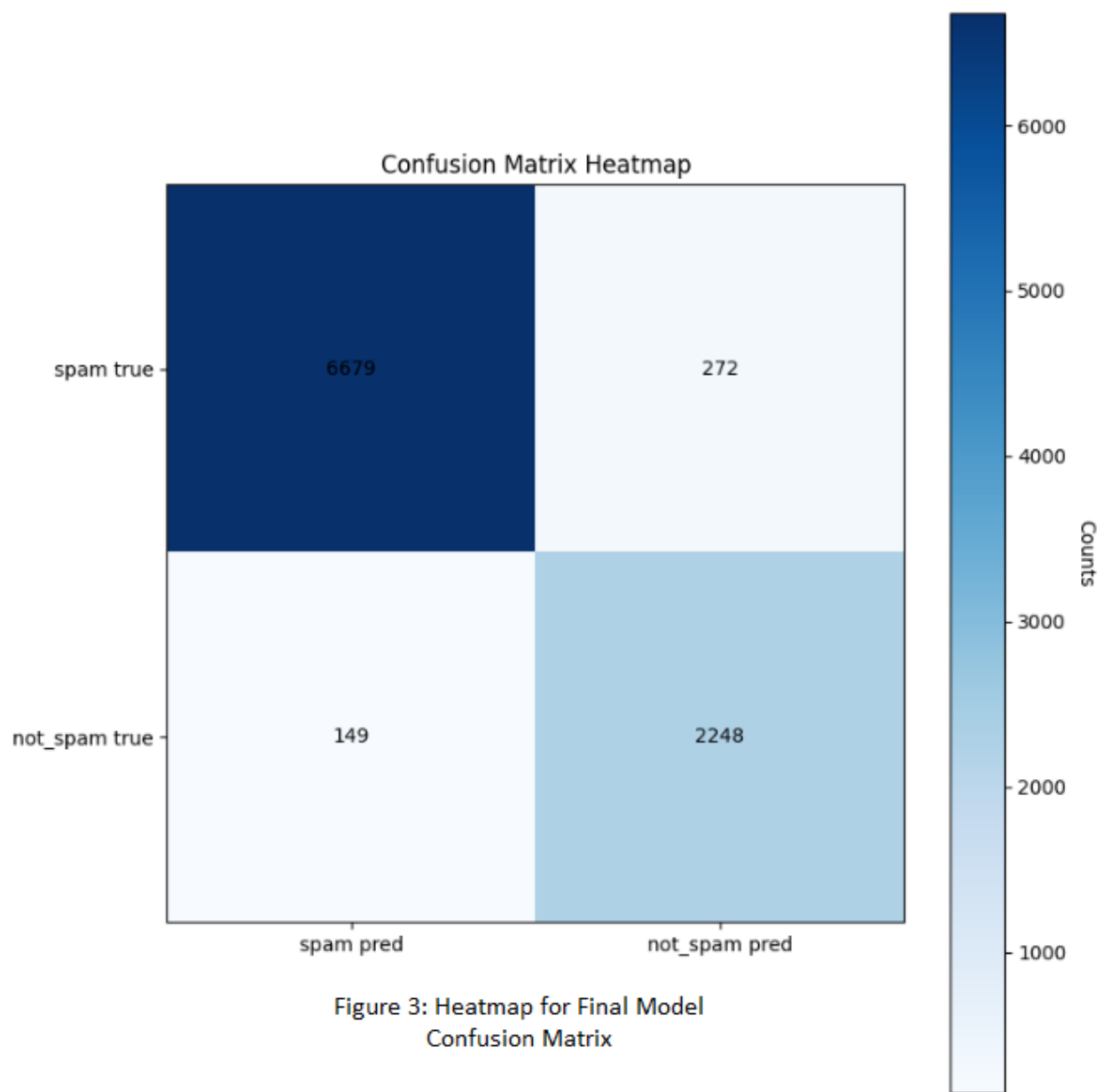


Figure 3: Heatmap for Final Model
Confusion Matrix

This confusion matrix provides further evidence of the strength of our model but from a different perspective. The truth labels and the predicted labels align well with each other. Only very few cases fail to be predicted properly for both classes. Since we decided a balanced optimizing approach these results are in line with that decision.

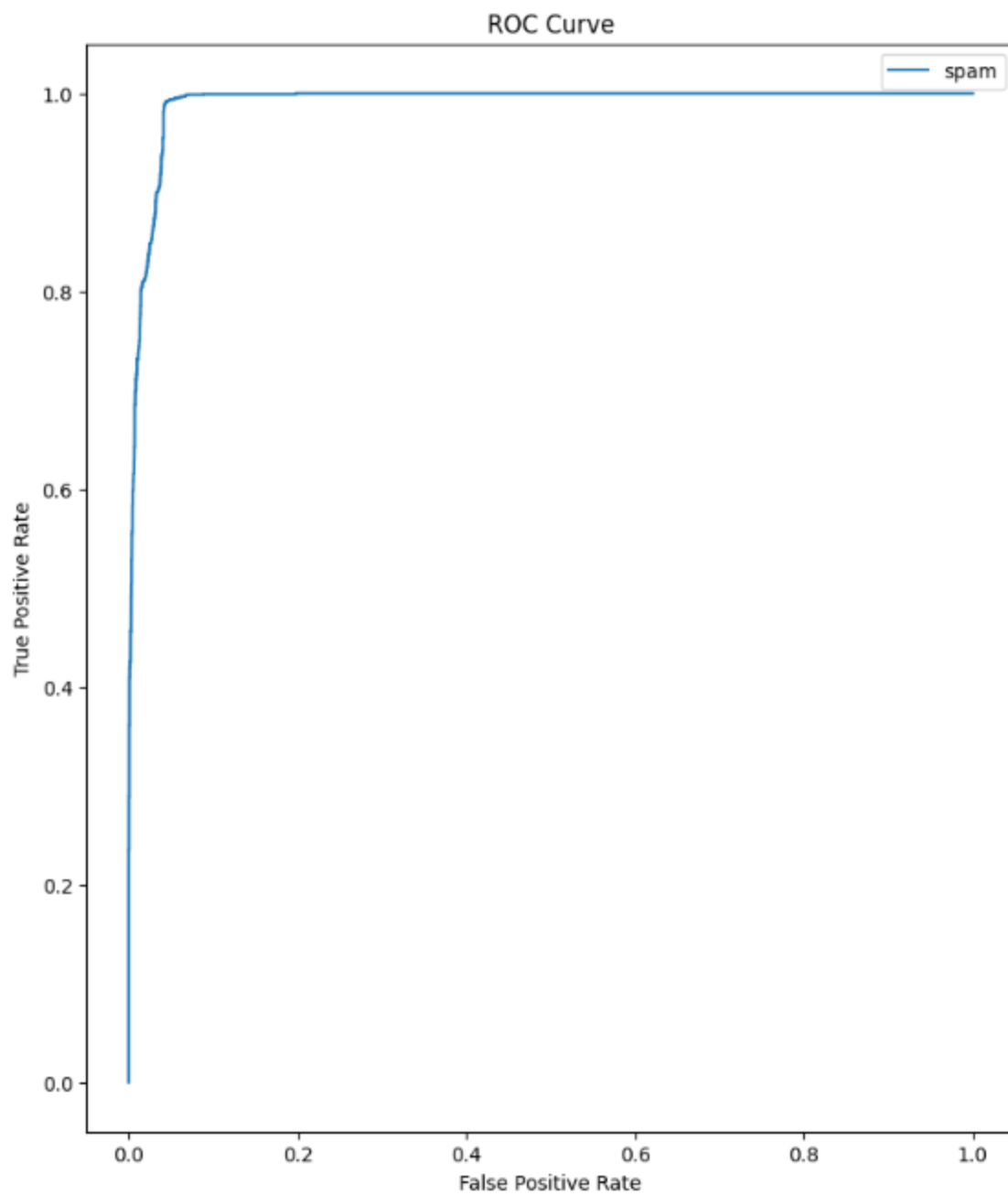


Figure 4: ROC of Final Model

```
AUC Report
~~~~~
spam AUC: 0.9906
```

Figure 5: AUC for Final Model ROC Curve

The ROC Curve provides further confidence in the performance of our model. Visually this is the curve of a good classifier but the numeric metric is also supplied for those who prefer that. These results are here to again show that this model does well but it is also here to provide some new context for an old question. Earlier we decided to use spam as the positive class and use f1 score as our optimizing metric. This was based on the idea that we were okay with both False Positives and False Negatives occurring at a close to equal rate. The above plot shows that there exists a threshold where we could in fact maximize the TPR while incurring a False Positive Rate Cost. Before it may not have been worth considering but now that we can show that it can be possible to choose one with minimal cost to the other may make this an attractive option to choose. Food for thought.

Conclusion

Overall through the use of bags of words, clusters, and naive bayes we were able to create a strong classifier for the spam assassin dataset. Further avenues of exploration remain in regards to what category we care about more and the costs of those decisions.

There is one loose end though. At the start of the analysis we decided to vectorize all the text data first and then perform analysis. This vectorized data formed our feature set. The question is was this the right choice? Technically by looking at all the text input information about the hold out set is present within the data even if the naive bayes model did not train on them.

Code

Appendix A

```
1  import os
2  import sys
3  import email
4  import pandas as pd
5  import numpy as np
6  import matplotlib.pyplot as plt
7  import time
8  from sklearn.model_selection import cross_val_score
9  from sklearn.metrics import confusion_matrix
10 from sklearn import metrics as mt
11 from sklearn.feature_extraction.text import TfidfVectorizer
12 from sklearn.cluster import KMeans
13 from scipy.sparse import hstack
14 from sklearn.naive_bayes import MultinomialNB
15 from sklearn.metrics import classification_report
16
17
18 # Build a spam classifier using naive Bayes and clustering.
19 # You will have to create your own dataset from the input messages.
20 # Be sure to document how you created your dataset.
21
22 INVALID_INPUT = 1
23
24
25 class UnexpectedEmailFileStructureError(Exception):
26     def __init__(self):
27         msg = "This script expects a file path to a directory that contains "
28         msg += "the spam assassin dataset. The spam assassin dataset is a "
29         msg += "group of five folders with the names, easy_ham, easy_ham_2, "
30         msg += "hard_ham, spam, spam_2. Did you specify a folder path that "
31         msg += "contains these folders?"
32         super().__init__(msg)
33
```



```

35 class MissingInputError(Exception):
36     def __init__(self):
37         msg = "This script expects exactly one input, a file directory that "
38         msg += "contains the spam assassin dataset. You input nothing. "
39         msg += "Did you forget to add an argument with the script?"
40         super().__init__(msg)
41
42
43 class ExcessiveInputError(Exception):
44     def __init__(self):
45         msg = "This script expects atleast one input, a file directory that "
46         msg += "contains the spam assassin dataset. You called this script "
47         msg += "with more than one argument. Only use one and try again. "
48         super().__init__(msg)
49
50
51 def recursive_payload_retrieval(email_payload):
52     email_body = ""
53     for sub_email in email_payload:
54         sub_email_payload = sub_email.get_payload()
55         if type(sub_email_payload) is list:
56             email_body += recursive_payload_retrieval(sub_email_payload)
57         elif type(sub_email_payload) is str:
58             email_body += sub_email_payload
59     return email_body
60
61
62 def check_correct_file_structure(inputs):
63     email_direc = inputs[1]
64     expected_dirs = ["easy_ham", "easy_ham_2", "hard_ham", "spam", "spam_2"]
65     first = True
66     for root, dirs, files in os.walk(email_direc):
67         if first:
68             first = False
69             for expected_direc in expected_dirs:
70                 if expected_direc not in dirs:
71                     print(expected_direc)
72                     print(dirs)
73                     raise UnexpectedEmailFileStructureError
74

```

```

76 def check_valid_inputs(inputs):
77     if len(inputs) < 2:
78         raise MissingInputError
79     elif len(inputs) > 2:
80         raise ExcessiveInputError
81
82
83 def load_emails(inputs):
84     email_direct = inputs[1]
85     data = {"email_body": [], "email_type": [], "email_class": []}
86     for root, dirs, files in os.walk(email_direct):
87         for f in files:
88             file_path = os.path.join(root, f)
89             with open(file_path, "r", encoding="latin-1") as email_file:
90                 if "cmds" in file_path:
91                     continue
92                 body = ""
93                 msg = email.message_from_file(email_file)
94                 t = msg.get_content_type()
95                 payload = msg.get_payload()
96                 if type(payload) is list:
97                     body = recursive_payload_retrieval(payload)
98                 elif type(payload) == str:
99                     body = payload
100             else:
101                 print("Unexpected Condition")
102                 raise Exception()
103
104             data["email_type"].append(t)
105
106             data["email_body"].append(body)
107
108             if "spam" in root:
109                 data["email_class"].append("spam")
110             else:
111                 data["email_class"].append("not_spam")
112     return data
113

```

```
115 def categ_categ_compare(df, categ1, categ2, file_name="no_name.png"):
116     horizontal = 1
117     vertical = 0
118     df_count = df.groupby([categ1, categ2])[categ1].count().unstack().fillna(0)
119     df_relative_count = df_count.div(
120         df_count.sum(axis=horizontal),
121         axis=vertical)*100
122     title = "Percentage Composition of " + categ1 + " with " + categ2
123     df_relative_count.plot.barh(stacked=True, title=title)
124     plt.savefig(file_name, bbox_inches='tight')
125     plt.clf()
```

```

127
128 def multinom_nb_cv(x, y, print_progress=False, n_jobs=1):
129     best_alpha = 0
130     best_f1 = 0
131     first = True
132     iter_count = 20
133     report_freq = int(iter_count/20)
134     i = 0
135     t1 = time.time()
136     for param_alpha in np.logspace(-3, 3, iter_count):
137
138         model = MultinomialNB(alpha=param_alpha)
139
140         f1_cv = cross_val_score(model, x, y, cv=5,
141                                 scoring='f1')
142
143         f1 = sum(f1_cv)/len(f1_cv)
144
145         if first:
146             first = False
147             best_alpha = param_alpha
148             best_f1 = f1
149         else:
150             if (f1 > best_f1):
151                 best_alpha = param_alpha
152                 best_f1 = f1
153
154         if print_progress and ((i+1) % report_freq == 0):
155             print(f"Iteration {i+1}. "
156                   f"Percent done = {(i+1)/iter_count*100:.4f}%")
157
158         i += 1
159     t2 = time.time()
160     elapsed_time = t2-t1
161
162     report_str = f"Best f1={best_f1:.4f}, "
163     report_str += f"alpha={best_alpha}, time={elapsed_time:.4f}s\n"
164
165     return best_alpha, report_str

```

```

168 def log_section(title="No Title", content="No Content"):
169     log_file.write(title + "\n")
170     log_file.write("~~~~~\n")
171     log_file.write(content + "\n")
172     log_file.write("\n")
173
174
175 def plot_heatmap(data,
176                 xlab,
177                 ylab,
178                 size=9,
179                 title="untitled heatmap",
180                 cbar_label="untitled",
181                 save_name="untitled.png"):
182
183     fig, ax = plt.subplots(figsize=(size, size))
184     im = ax.imshow(data, cmap=plt.cm.Blues)
185
186     for i in range(len(ylab)):
187         for j in range(len(xlab)):
188             ax.text(j, i, data[i, j],
189                   ha="center", va="center", color="black")
190
191     # Show all ticks and label them with the respective list entries
192     ax.set_xticks(np.arange(len(xlab)), labels=xlab)
193     ax.set_yticks(np.arange(len(ylab)), labels=ylab)
194
195     cbar = ax.figure.colorbar(im, ax=ax)
196     cbar.ax.set_ylabel(cbar_label, rotation=-90, va="bottom")
197
198     ax.set_title(title)
199     fig.tight_layout()
200     plt.savefig(save_name, bbox_inches='tight')
201     plt.clf()

```

```

204 if __name__ == '__main__':
205     inputs = sys.argv
206     check_valid_inputs(inputs)
207     check_correct_file_structure(inputs)
208     # Log Start
209     log_file = open("log.txt", "w")
210     log_file.write("Case Study 3 Report\n\n")
211
212     # Use input to find and load emails
213     email_data_raw = load_emails(inputs)
214     email_df = pd.DataFrame(email_data_raw)
215
216     # Id to class maps
217     email_class_to_id = {"spam": 1, "not_spam": 0}
218     id_to_email_class = {1: "spam", 0: "not_spam"}
219
220     # TFIDF Vectorizer
221     email_bodies = email_df["email_body"]
222     vectorizer = TfidfVectorizer(min_df=0.002)
223     features = vectorizer.fit_transform(email_bodies)
224
225     # Cluster with Kmeans
226     cluster_model = KMeans(n_clusters=2)
227     cluster_model.fit(features)
228     email_df["cluster_id"] = cluster_model.labels_
229
230     # Plot to see which clusters refers to what email class
231     categ_categ_compare(email_df,
232                         "email_class",
233                         "cluster_id",
234                         "categorical_cluster_comparison.png")
235
236     # add cluster feature to vectorizer features
237     cluster_ids_raw = cluster_model.labels_
238     cluster_ids_resaped = cluster_ids_raw.reshape(-1, 1)
239     features_and_cluster_ids = hstack((features, cluster_ids_resaped)).A

```

```

# Optimizing Alpha
email_id = email_df["email_class"].map(email_class_to_id)
best_alpha, report = multinom_nb_cv(features_and_cluster_ids,
                                     email_id,
                                     n_jobs=5,
                                     print_progress=True)
log_section(title="Optimization Report",
            content=report)

# Best Model
final_clf = MultinomialNB(alpha=best_alpha)
final_clf.fit(features_and_cluster_ids, email_df["email_class"])

# Classification Report
y_pred_final = final_clf.predict(features_and_cluster_ids)
y_true = email_df["email_class"]
classif_rep = classification_report(
    y_true,
    y_pred_final
)
log_section(title="Classification Report",
            content=classif_rep)

```

```

# Confusion Matrix Heatmap
conf_mat = confusion_matrix(y_true, y_pred_final)

base_labels = list(id_to_email_class.values())
label_true = []
label_pred = []
for i in range(len(base_labels)):
    true_label = base_labels[i] + " true"
    pred_label = base_labels[i] + " pred"
    label_pred.append(pred_label)
    label_true.append(true_label)
plot_heatmap(conf_mat,
              label_pred,
              label_true,
              size=8,
              title="Confusion Matrix Heatmap",
              cbar_label="Counts",
              save_name="conf_heatmap.png")

# ROC Curve
y_pred_prob = final_clf.predict_proba(features_and_cluster_ids)
y_label = id_to_email_class[1]
preds = y_pred_prob[:, 1]
fpr, tpr, thresholds = mt.roc_curve(email_id, preds, pos_label=i)
plt.plot(fpr, tpr, label=y_label)
title_str = 'ROC Curve'
plt.title(title_str)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
auc_report = f"{y_label} AUC: {mt.auc(fpr, tpr):.4f}\n"
log_section(title="AUC Report",
            content=auc_report)
plt.legend()
plt.savefig("roc_plot.png", bbox_inches='tight')
plt.clf()
log_file.write("\n")
log_file.close()

```