

MILESTONE # 5: Final Report
IS 567 - Fall 2022

Muhammad Ibrahim Mian
UIN: 673846900

**Email Spam Classification using Machine
Leaning**

Introduction

Of all the different communication methods, email is one of the most extremely important form of communication. It has replaced the traditional snail mail and can be accessed from any part of the world just with the help of internet connectivity. With the ever growing number of users, emails are set to increase to a record high, but the worst part of that is, out of that approximately 57 % emails are of no use as they are spam emails.

In my project, I will be using Enron Email database to identify and classify spam emails. Enron corporation was an energy company which came into the spotlight in 2001 due to accounting frauds and emails of many executives and was made publicly available.

Descriptive Statistics

1. We will first look at the dataset itself.
 - a. Total emails (Testing & Training Set): 33,716
 - Spam Emails: 17,171
 - Ham Emails: 16,545
 - b. Total number of columns: 5
2. We remove the null values in column "Message"
 - a. Total emails (Testing & Training Set): 33,345
 - Spam Emails: 16,852
 - Ham Emails: 16,493
 - b. Total number of columns: 5
3. Train/Test Split
 - a. Total emails (Testing & Training Set): 33,345
 - b. Train Set
 - Spam Emails: 11,813
 - Ham Emails: 11,528
 - c. Test Set
 - Spam Emails: 5,039
 - Ham Emails: 4,965

4. Data extracted and in the data frame

Data present in data frame - (See Figure 1)

5. Email examples

Here is an example of the message. Initially we have it in raw form with all the numbers and characters. After cleaning it by removing non essential spaces, punctuation marks, special characters and opening contractions, we see the following result

- a. Raw email - (See Figure 2)
- b. Cleaned email - (See Figure 3)

6. Tokenized Words most common words (Plots)

Here, we can see the plots for the most common occurring words. We can see many stop words are present so we need to remove those first - (See Figure 4)

7. Stopwords removed most common words (Plots)

Here we get a better sense of what sort of words are common in the emails. This tells us some important information of the content of the emails - (See Figure 5)

8. Wordcloud (Messages column)

Similar to the plots, we can see in the word cloud the most occurring words. This is before any sort of filters on the data - (See Figure 6)

9. Wordcloud (ProcessedMessage column)

Similar to the plots, we can see in the word cloud the most occurring words. This is after filters on the data - (See Figure 7)

10. NLP Processing

When we performed some analysis on this data, we noticed that there were a lot of mentions regarding names and organizations. These connections help us analyze what sort of data was present in the emails.

- Email - (See Figure 8)
- POS Tagging - (See Figure 9)
- Dependency parsing - (See Figure 10)
- Name Entity Recognition - (See Figure 11)

Pre-processing / Transformation Steps

In this part, I will be discussing how I used my data in raw form and performed pre processing on it

1. We start with data in txt files labeled Ham and Spam. We need to join the files together and create a data frame from it
 - We start with downloading the dataset
 - We separate the emails from the subject
 - We check and add date to the row for the email.
2. Next, we will load the data in our data frame and perform checks
 - Label wise counts
 - All email counts
 - Null values
3. After this, we will perform steps on the text data
 - Remove special characters
 - Reduce spaces with punctuation to improve data quality (contractions)
 - All lowercase
 - Remove links
 - Open contractions
 - Tokenization
 - Stopwords Removal
 - Lemmatization

The reason to conduct these steps was to clean the data of any elements that would affect our models performance. In order to create better sense of the data present, we used the following steps above.

The remaining features after cleaning the data are important as they build to the context of the email. If there are words such as money, bank account, private emails which raise a red flag and set as SPAM email. The model will learn from the features and see if the same sort of words occur, they will mark it as SPAM or pass it as HAM

Feature Extraction / Selection

For this task, we will be using 5 methods for feature extraction

1. Count Vectorizer
2. TF-IDF
3. K best feature using Chi2
4. Variance threshold
5. Word2Vec

Using these, we will now apply models on different combinations

Preliminary models, parameters, evaluation results, and error analysis

For all results combined - (See Table 1)

1. Naive Bayes

I started with the simplest model first and applied both count vectorization and TF-IDF features on it

- Count Vectorizer: We started with top 1000 features of the data and using that to create the vectors for the model. Here are the results of the model - (See Figure 12)
- TF-IDF: We started with top 1000 features of the data and using that to create the vectors for the model. Here are the results of the model - (See Figure 13)

Comparing both features, we see that as we increase the number of features, we see that the false negatives decrease and false positives increase meaning HAM emails are also being classified as SPAM which makes emails unnecessarily been marked wrong

2. Support vector machines

Now we look at support vector machines and their performance

- Count Vectorizer: We started with top 1000 features of the data and using that to create the vectors for the model. Here are the results of the model - (See Figure 14)
- TF-IDF: We started with top 1000 features of the data and using that to create the vectors for the model. Here are the results of the model - (See Figure 15)

Comparing Naive bayes with Support vector machines, we see that it gives us a better result in terms of accuracy and with TF-IDF, it gives a better result compared to count vectorizer

3. K Nearest Neighbor (KNN)

I also took a look at unsupervised learning compared to supervised learning to see how it performs. Using elbow method, we found that K=4 is the best parameter - (See Figure 16)

- Count Vectorizer: We started with top 1000 features of the data and using that to create the vectors for the model. Here are the results of the model - (See Figure 17)
- TF-IDF: We started with top 1000 features of the data and using that to create the vectors for the model. Here are the results of the model - (See Figure 18)

Comparing Naive bayes and Support vector machines, we see that it performs very poorly and thus is not a good idea to implement unsupervised learning in our dataset

4. K Best features using Chi2

Now, we have established that unigram and TF-IDF gives us the best features, we will now use this combination with Chi2 to select our features and see how the number of features effects the models performance - (See Figure 19 & 20)

Looking at both the models, we see that SVM is giving us the best performance compared to Naive Bayes. It is also interesting to note that as the number of features increase, the better the results so it will be good to see how more features can improve the results as well as looking at the computational cost of running the model.

5. Variance Threshold

Taking the previous assumptions here as well, now the features are filtered with variance threshold and see how the features effect the model.

For this, TF-IDF was giving very poor results so we used count vectorizer with variance threshold - (See Figure 21 & 22)

Here, we see that both models are performing pretty much similar with variance threshold even with a huge number of features. K best features using chi2 still performs better than variance threshold

6. Decision Tree

Here, I decided to use decision tree by trying to limit the features even more using depth of the decision trees. - (See Figure 23)

Notice that the accuracy has dropped severely compared to our other models so it is not recommended

7. SVM - Word2Vec

For our final task, I used the Word2Vec word embedding to create features.

In Word2Vec, we will take a reference word embedding which already have the vectors for words. The reference word embeddings I am using is of Google which is publicly available. Google has used its search engine to create references for words and their connections.

The data will match the library and copy the vector representation for the available words. For the remaining words not found, we have couple of options. First is to ignore those words from the features which we are using here. Second is we can create our own vector for it and include in the dictionary - (See Figure 24)

Conclusion and Insights

After analyzing the results, following observations can be made

- Models perform best when only unigram data is present and accuracy starts to decrease as more pairs are made
- TF-IDF performed better compared to Count Vectorizer
- SVM performed better compared to Naïve Bayes
- We are currently working with 1000 sample size and looking at K best model, we can see that adding more features will help us achieve better results even with simple models
- Using only one trained word embedding vector on the Word2Vec, it shows promising trend and it would be interesting to see in the future that how other word embeddings can improve the model results

Some future improvements we can make to this project

- Data can be trained on neural networks to see how the model performs
- Different word embedding vector model for Word2Vec
- Testing model on different dataset

References

1. Singh, V. K., & Bhardwaj, S. (2018). Spam mail detection using classification techniques and global training set. *Intelligent Computing and Information and Communication*, 623–632. https://doi.org/10.1007/978-981-10-7245-1_61
2. Peng, W., Huang, L., Jia, J., & Ingram, E. (2018). Enhancing the naive Bayes spam filter through intelligent text modification detection. 2018 17th IEEE International Conference On Trust, Security And Privacy In Computing And Communications/ 12th IEEE International Conference On Big Data Science And Engineering (TrustCom/BigDataSE). <https://doi.org/10.1109/trustcom/bigdatase.2018.00122>
3. Singh, M., Pamula, R., & shekhar, S. kumar. (2018). Email spam classification by support vector machine. 2018 International Conference on Computing, Power and Communication Technologies (GUCON). <https://doi.org/10.1109/gucon.2018.8674973>
4. Muhammad Abdulhamid, S., Shuaib, M., Osho, O., Ismaila, I., & K. Alhassan, J. (2018). Comparative analysis of classification algorithms for email spam detection. *International Journal of Computer Network and Information Security*, 10(1), 60–67. <https://doi.org/10.5815/ijcnis.2018.01.07>
5. Kulwinder Kaur, & Dr. Mukesh Kumar. (2015). Spam detection using KNN, back propagation and Recurrent Neural Network. *International Journal of Engineering Research And*, V4(09). <https://doi.org/10.17577/ijertv4is090492>

Appendix

1. Figure 1

| Message ID | | Subject | | Message | Spam/Ham | Date |
|------------|---|--|---|---------|----------|------------|
| 0 | 0 | ena sales on hpl | just to update you on this project ' s status ... | | ham | 2000-05-10 |
| 1 | 1 | 98 - 6736 & 98 - 9638 for 1997 (ua 4 issues) | the above referenced meters need to be placed ... | | ham | 2000-02-18 |
| 2 | 2 | hpl nominations for december 28 , 1999 | (see attached file : hpl 228 . xls)\n- hpl... | | ham | 1999-12-27 |
| 3 | 3 | revised nom - kcs resources | daren ,\nit ' s in .\nbob\n- -----... | | ham | 2000-06-29 |
| 4 | 4 | new production - sitara deals needed | daren ,\nfyi .\nbob\n- -----... | | ham | 2000-07-28 |

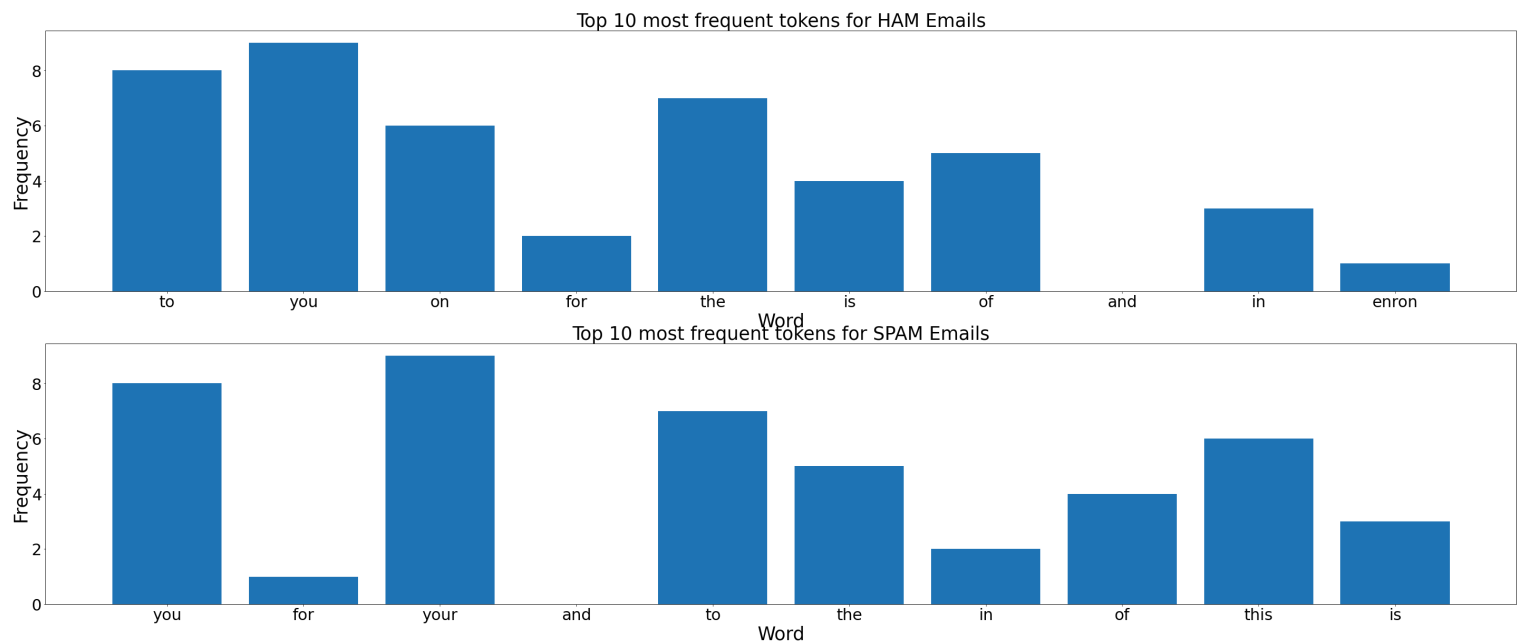
2. Figure 2

gentlemen :
please review and let me know if you have any questions .
for scheduling purposes , we will show a receipt from hpl (transportation
agreement # 4047) at aqua dulce of 45 , 000 , with deliveries to :
air products - la porte 5 , 000
oxy battleground 10 , 000
rohms & haas dp 20 , 000
dupont dp 10 , 000
dwright : you will need to coordinate these flow changes with the facilities .
thanks ,
>

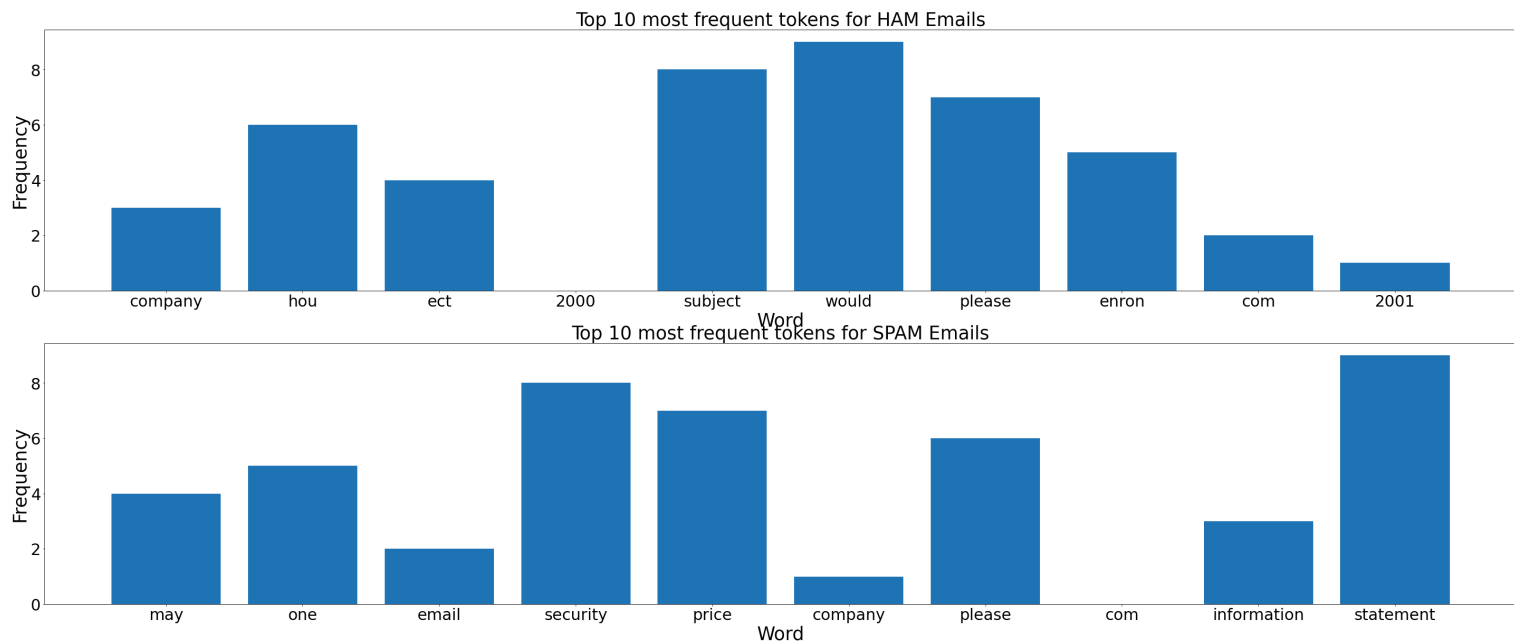
3. Figure 3

gentlemen please review and let me know if you have any questions for scheduling purposes we will show a receipt from
hpl transportation agreement 4047 at aqua dulce of 45 000 with deliveries to air products la porte 5 000 oxy battlegr
ound 10 000 rohms haas dp 20 000 dupont dp 10 000 dwright you will need to coordinate these flow changes with the facil
ities thanks this email and any files transmitted with it from the elpaso corporation are confidential and intended s
olely for the use of the individual or entity to whom they are addressed if you have received this email in error ple
ase notify the sender smartpigging 812 lomax market takes plan final xls

4. Figure 4



5. Figure 5



8. Figure 8

can you help me out on this darren ? mjj
 ----- forwarded by mary jo johnson / hou / ect on 11 / 09 / 2000
 10 : 04 am -----
 " john daugherty " on 11 / 08 / 2000 04 : 38 : 37 pm
 to :
 cc :
 subject : re : driscoll ranch # 3 gas pricing and interconnect estimate
 mary jo ,
 thanks for the update . regarding the notice provision of 6 business days
 prior to the close of business on the last business day of the month prior
 to selected month , does that mean we need to give you notice for december
 by tuesday , november 21 st at 5 : 00 pm or monday , november 20 th at 5 : 00 pm
 assuming the 23 rd and 24 th are holidays ?
 john daugherty
 ----- original message -----
 from :
 to :
 cc : ; ;
 ; ;
 ; ;
 ;
 sent : wednesday , november 08 , 2000 5 : 12 pm
 subject : re : driscoll ranch # 3 gas pricing and interconnect estimate

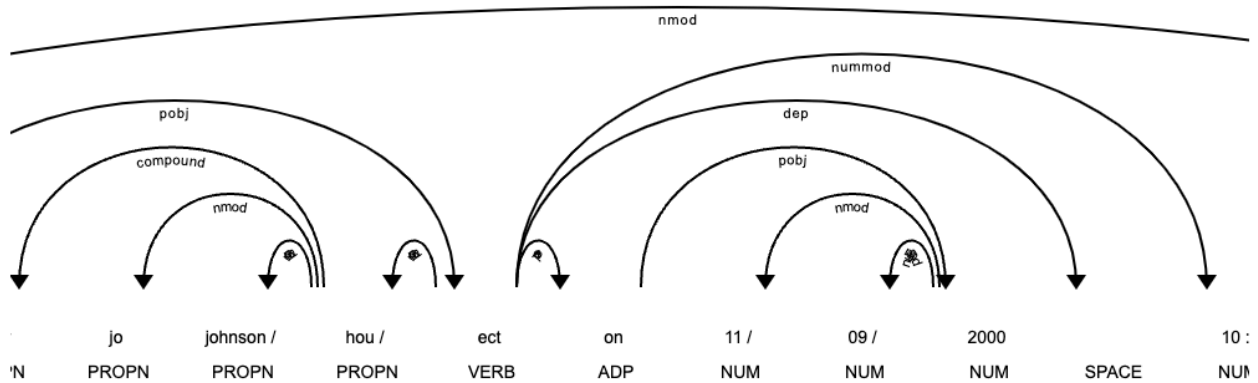
9. Figure 9

```
In [19]: #POS Tagging
```

```
nlp = spacy.load("en_core_web_sm")
doc = nlp(message)
for token in doc:
    print(f'{token.text:{10}} {token.pos_:{10}} {token.tag_:{10}}')
```

[illegible]

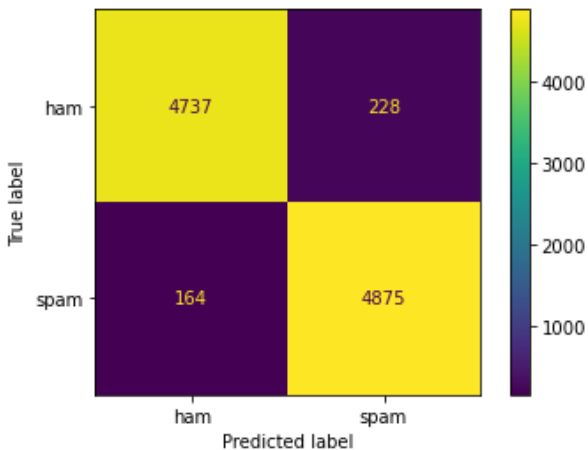
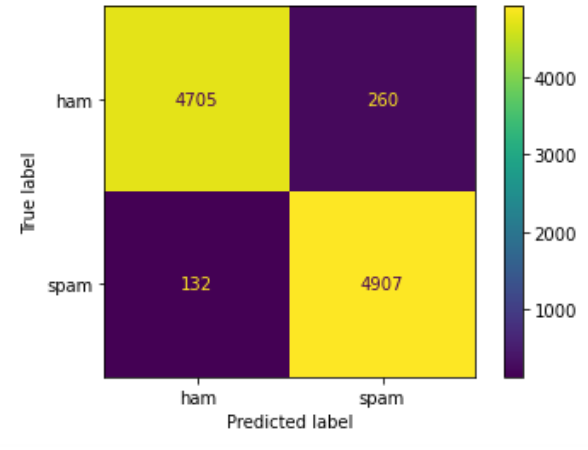
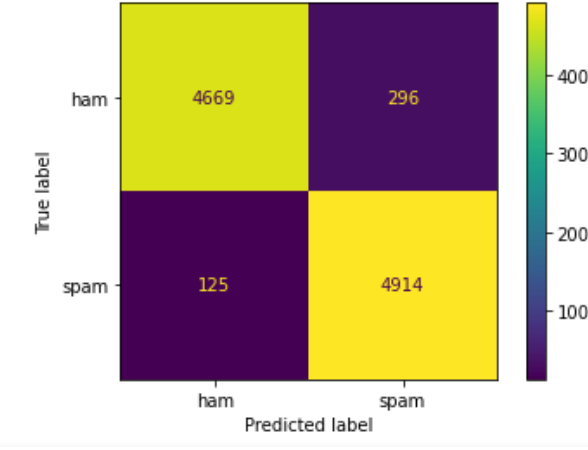
10. Figure10



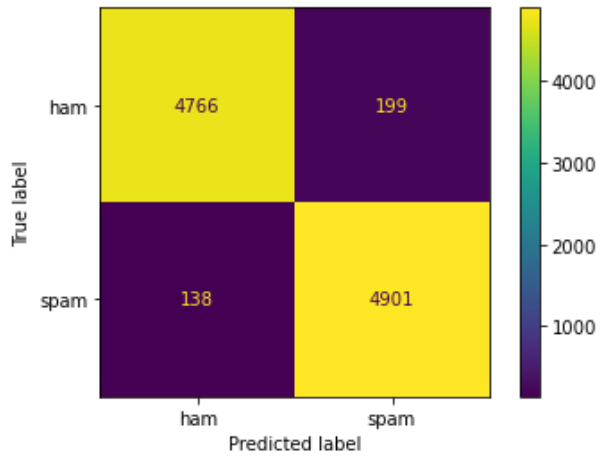
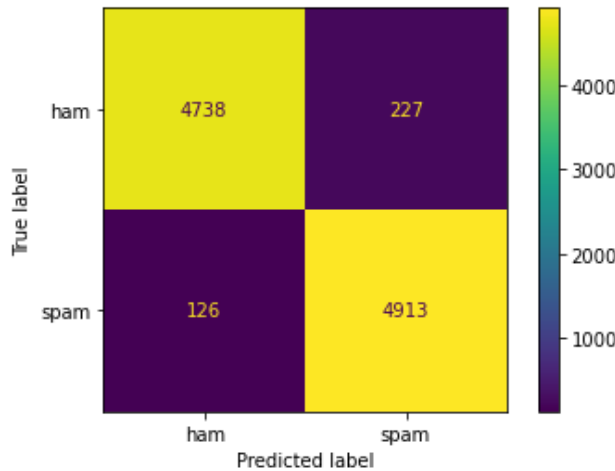
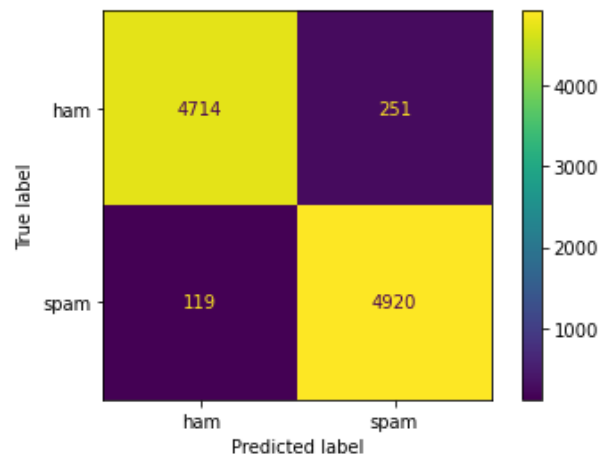
11. Figure 11

```
[('mary jo johnson / hou', 'PERSON'),
 ('11 / 09 / 2000', 'DATE'),
 ('04', 'CARDINAL'),
 ('john daugherty', 'PERSON'),
 ('11 / 08 / 2000 04', 'DATE'),
 ('38', 'CARDINAL'),
 ('37 pm', 'QUANTITY'),
 ('6', 'CARDINAL'),
 ('the last business day of the month', 'DATE'),
 ('month', 'DATE'),
 ('december', 'DATE'),
 ('tuesday', 'DATE'),
 ('5', 'CARDINAL'),
 ('00 pm', 'TIME'),
 ('monday , november 20', 'DATE'),
 ('5', 'CARDINAL'),
 ('00', 'CARDINAL'),
 ('23', 'CARDINAL'),
 ('24', 'CARDINAL'),
 ('john daugherty\n- - - - -', 'PERSON'),
 ('wednesday , november 08 , 2000 5', 'DATE'),
 ('12 pm', 'QUANTITY')]
```

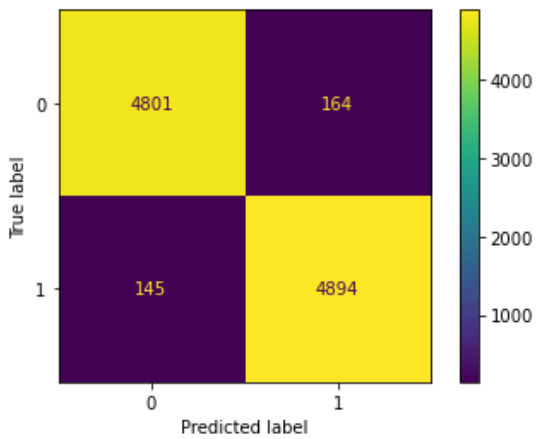
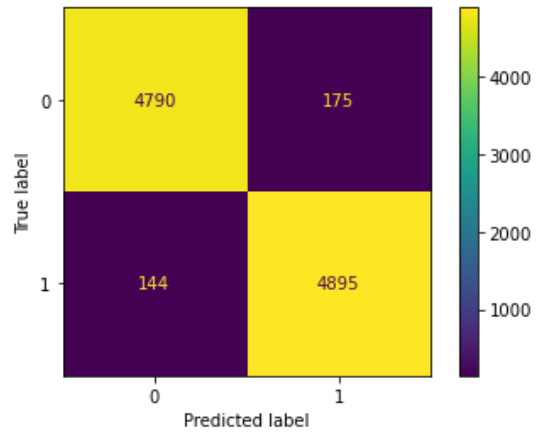
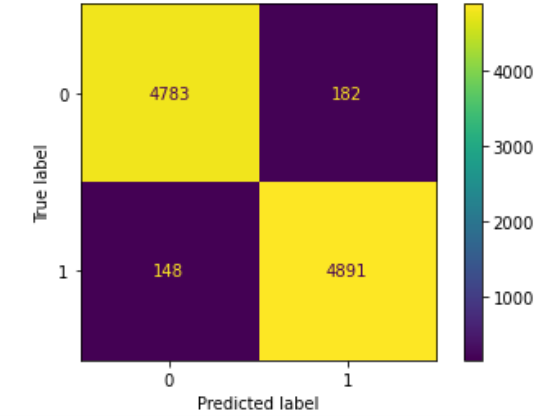
12. Figure 12

| Model | Sub Type | Feature space | Accuracy | Confusion Matrix | | | | | | | | | |
|--------------------------------|------------------------------|--|------------------------------|---|------------------------------|-----|------|-----|------|-----|------|-----|------|
| Naive Bayes - Count Vectorizer | Unigram | 1000 | 0.961 |  <table><tr><th>True label \ Predicted label</th><th>ham</th><th>spam</th></tr><tr><th>ham</th><td>4737</td><td>228</td></tr><tr><th>spam</th><td>164</td><td>4875</td></tr></table> | True label \ Predicted label | ham | spam | ham | 4737 | 228 | spam | 164 | 4875 |
| | True label \ Predicted label | | ham | spam | | | | | | | | | |
| | ham | | 4737 | 228 | | | | | | | | | |
| spam | 164 | 4875 | | | | | | | | | | | |
| Unigram+Bi Gram | 0.961 |  <table><tr><th>True label \ Predicted label</th><th>ham</th><th>spam</th></tr><tr><th>ham</th><td>4705</td><td>260</td></tr><tr><th>spam</th><td>132</td><td>4907</td></tr></table> | True label \ Predicted label | ham | spam | ham | 4705 | 260 | spam | 132 | 4907 | | |
| True label \ Predicted label | ham | spam | | | | | | | | | | | |
| ham | 4705 | 260 | | | | | | | | | | | |
| spam | 132 | 4907 | | | | | | | | | | | |
| Unigram+Bi gram+Trigram | 0.958 |  <table><tr><th>True label \ Predicted label</th><th>ham</th><th>spam</th></tr><tr><th>ham</th><td>4669</td><td>296</td></tr><tr><th>spam</th><td>125</td><td>4914</td></tr></table> | True label \ Predicted label | ham | spam | ham | 4669 | 296 | spam | 125 | 4914 | | |
| True label \ Predicted label | ham | spam | | | | | | | | | | | |
| ham | 4669 | 296 | | | | | | | | | | | |
| spam | 125 | 4914 | | | | | | | | | | | |

13. Figure 13

| Model | Sub Type | Feature space | Accuracy | Confusion Matrix | | | | | | | | | |
|------------------------------|------------------------------|--|------------------------------|--|------------------------------|-----|------|-----|------|-----|------|-----|------|
| Naive Bayes - TF-IDF | Unigram | 1000 | 0.966 |  <table><tr><th>True label \ Predicted label</th><th>ham</th><th>spam</th></tr><tr><th>ham</th><td>4766</td><td>199</td></tr><tr><th>spam</th><td>138</td><td>4901</td></tr></table> | True label \ Predicted label | ham | spam | ham | 4766 | 199 | spam | 138 | 4901 |
| | True label \ Predicted label | | ham | spam | | | | | | | | | |
| | ham | | 4766 | 199 | | | | | | | | | |
| spam | 138 | 4901 | | | | | | | | | | | |
| Unigram+Bi Gram | 0.965 |  <table><tr><th>True label \ Predicted label</th><th>ham</th><th>spam</th></tr><tr><th>ham</th><td>4738</td><td>227</td></tr><tr><th>spam</th><td>126</td><td>4913</td></tr></table> | True label \ Predicted label | ham | spam | ham | 4738 | 227 | spam | 126 | 4913 | | |
| True label \ Predicted label | ham | spam | | | | | | | | | | | |
| ham | 4738 | 227 | | | | | | | | | | | |
| spam | 126 | 4913 | | | | | | | | | | | |
| Unigram+Bigram+Trigram | 0.963 |  <table><tr><th>True label \ Predicted label</th><th>ham</th><th>spam</th></tr><tr><th>ham</th><td>4714</td><td>251</td></tr><tr><th>spam</th><td>119</td><td>4920</td></tr></table> | True label \ Predicted label | ham | spam | ham | 4714 | 251 | spam | 119 | 4920 | | |
| True label \ Predicted label | ham | spam | | | | | | | | | | | |
| ham | 4714 | 251 | | | | | | | | | | | |
| spam | 119 | 4920 | | | | | | | | | | | |

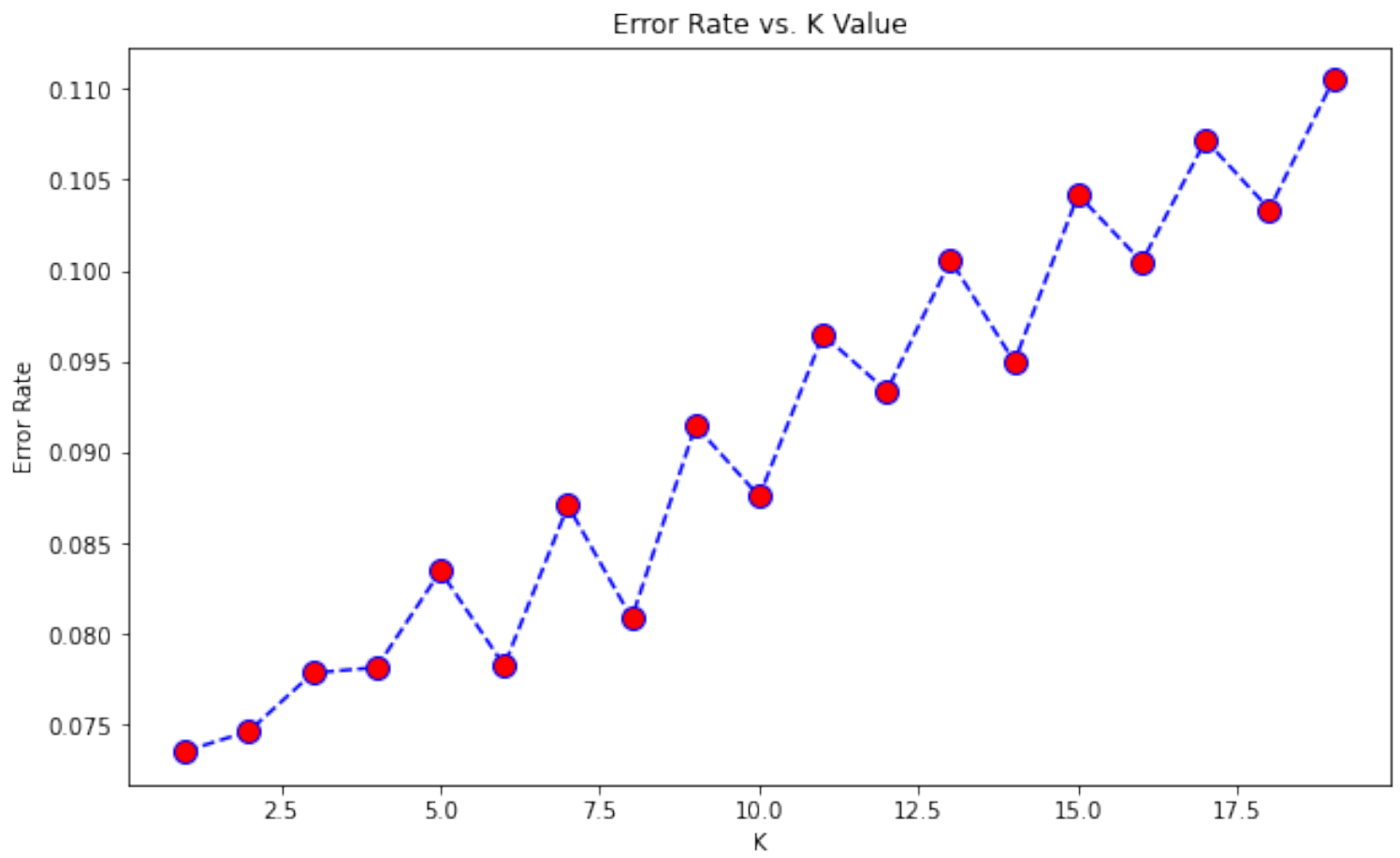
14. Figure 14

| Model | Sub Type | Feature space | Accuracy | Confusion Matrix | | | | | | | | | |
|--------------------------------|------------------------------|--|------------------------------|--|------------------------------|---|------|-----|------|-----|------|-----|------|
| SVM - Count Vectorizer | Unigram | 1000 | 0.969 |  <table><tr><th>True label \ Predicted label</th><th>0</th><th>1</th></tr><tr><th>0</th><td>4801</td><td>164</td></tr><tr><th>1</th><td>145</td><td>4894</td></tr></table> | True label \ Predicted label | 0 | 1 | 0 | 4801 | 164 | 1 | 145 | 4894 |
| | True label \ Predicted label | | 0 | 1 | | | | | | | | | |
| | 0 | | 4801 | 164 | | | | | | | | | |
| 1 | 145 | 4894 | | | | | | | | | | | |
| Unigram+ Bi Gram | 0.968 |  <table><tr><th>True label \ Predicted label</th><th>0</th><th>1</th></tr><tr><th>0</th><td>4790</td><td>175</td></tr><tr><th>1</th><td>144</td><td>4895</td></tr></table> | True label \ Predicted label | 0 | 1 | 0 | 4790 | 175 | 1 | 144 | 4895 | | |
| True label \ Predicted label | 0 | 1 | | | | | | | | | | | |
| 0 | 4790 | 175 | | | | | | | | | | | |
| 1 | 144 | 4895 | | | | | | | | | | | |
| Unigram+ Bigram+T rigram | 0.967 |  <table><tr><th>True label \ Predicted label</th><th>0</th><th>1</th></tr><tr><th>0</th><td>4783</td><td>182</td></tr><tr><th>1</th><td>148</td><td>4891</td></tr></table> | True label \ Predicted label | 0 | 1 | 0 | 4783 | 182 | 1 | 148 | 4891 | | |
| True label \ Predicted label | 0 | 1 | | | | | | | | | | | |
| 0 | 4783 | 182 | | | | | | | | | | | |
| 1 | 148 | 4891 | | | | | | | | | | | |

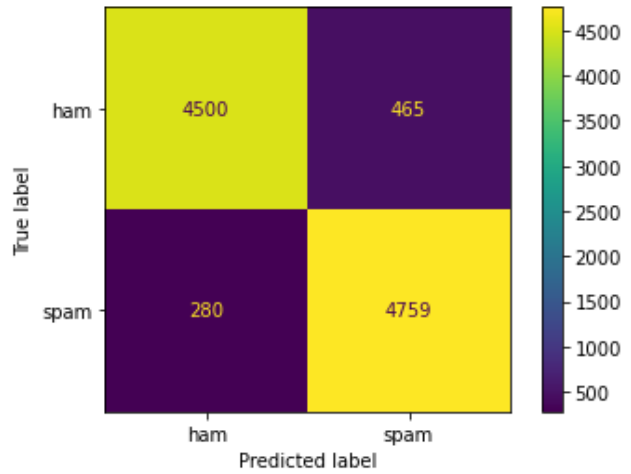
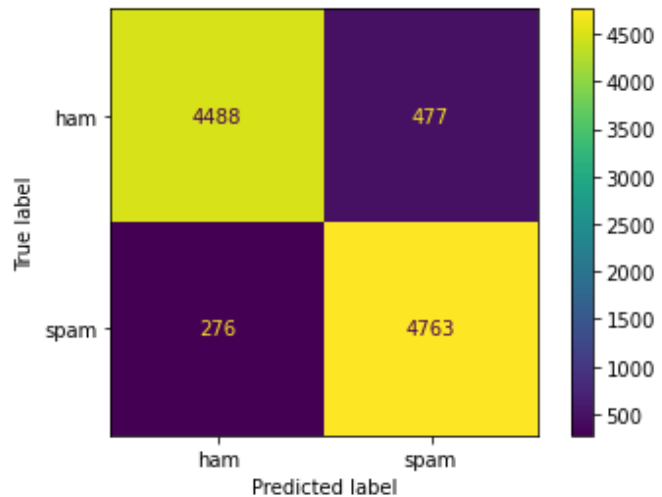
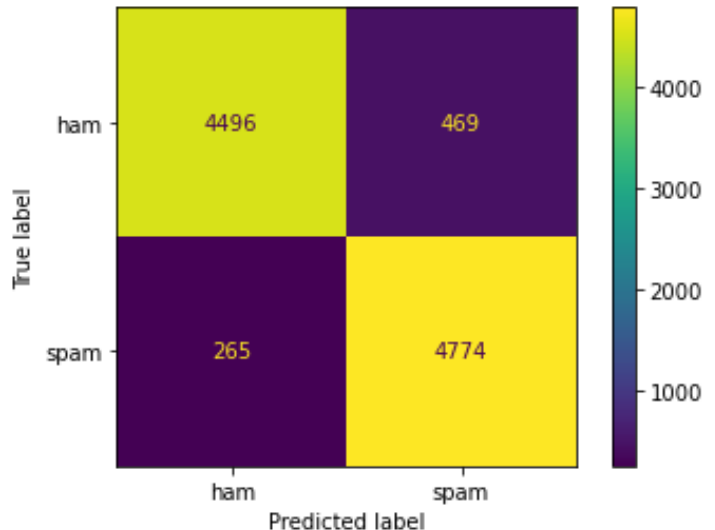
15. Figure 15

| Model | Sub Type | Feature space | Accuracy | Confusion Matrix | | | | | | | | | |
|------------------------------|------------------------------|---|------------------------------|---|------------------------------|---|------|-----|------|-----|------|-----|------|
| SVM - TF-IDF | Unigram | 1000 | 0.969 | <table><tr><th>True label \ Predicted label</th><th>0</th><th>1</th></tr><tr><th>0</th><td>4821</td><td>144</td></tr><tr><th>1</th><td>163</td><td>4876</td></tr></table> | True label \ Predicted label | 0 | 1 | 0 | 4821 | 144 | 1 | 163 | 4876 |
| | True label \ Predicted label | | 0 | 1 | | | | | | | | | |
| | 0 | | 4821 | 144 | | | | | | | | | |
| 1 | 163 | 4876 | | | | | | | | | | | |
| Unigram+ Bi Gram | 0.968 | <table><tr><th>True label \ Predicted label</th><th>0</th><th>1</th></tr><tr><th>0</th><td>4827</td><td>138</td></tr><tr><th>1</th><td>176</td><td>4863</td></tr></table> | True label \ Predicted label | 0 | 1 | 0 | 4827 | 138 | 1 | 176 | 4863 | | |
| True label \ Predicted label | 0 | 1 | | | | | | | | | | | |
| 0 | 4827 | 138 | | | | | | | | | | | |
| 1 | 176 | 4863 | | | | | | | | | | | |
| Unigram+ Bigram+Trigram | 0.968 | <table><tr><th>True label \ Predicted label</th><th>0</th><th>1</th></tr><tr><th>0</th><td>4810</td><td>155</td></tr><tr><th>1</th><td>162</td><td>4877</td></tr></table> | True label \ Predicted label | 0 | 1 | 0 | 4810 | 155 | 1 | 162 | 4877 | | |
| True label \ Predicted label | 0 | 1 | | | | | | | | | | | |
| 0 | 4810 | 155 | | | | | | | | | | | |
| 1 | 162 | 4877 | | | | | | | | | | | |

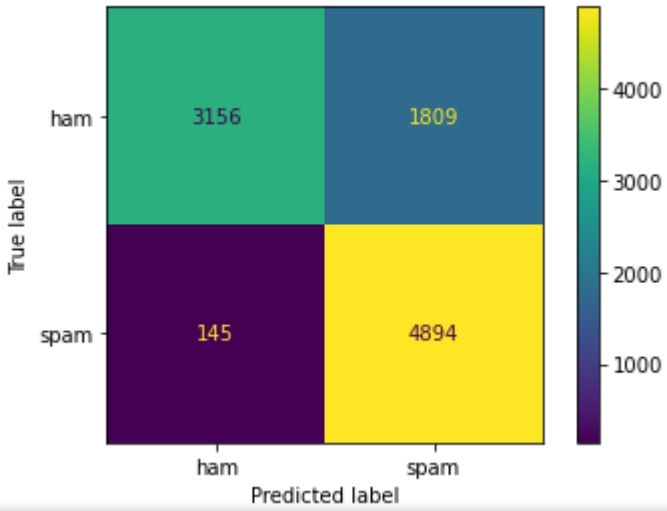
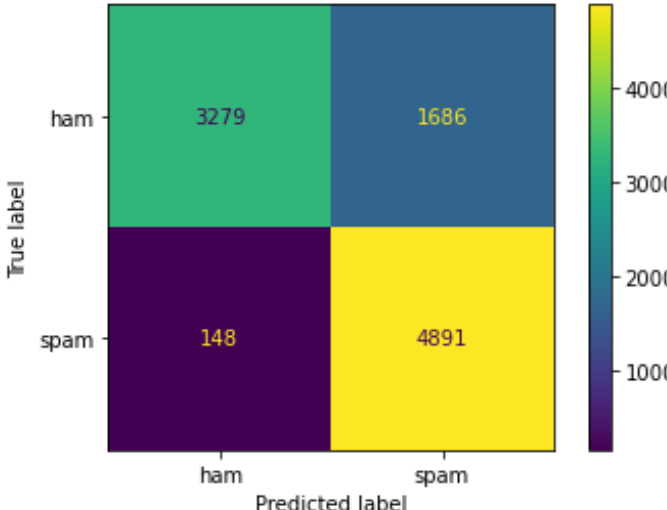
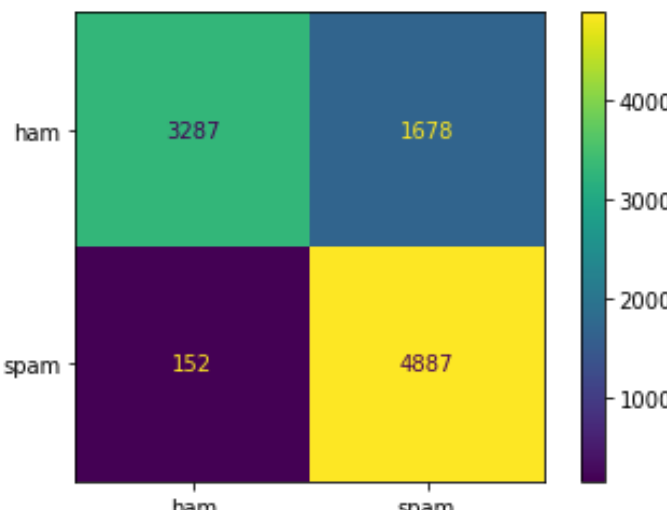
16. Figure 16



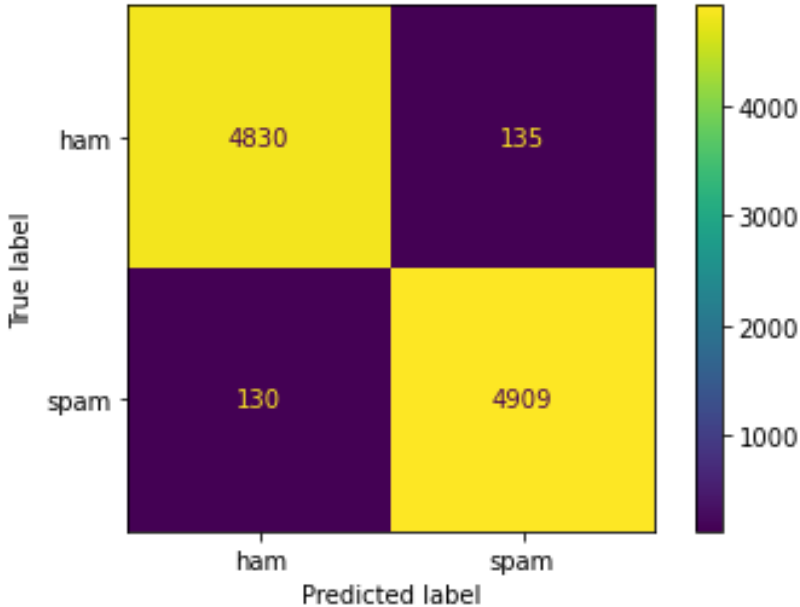
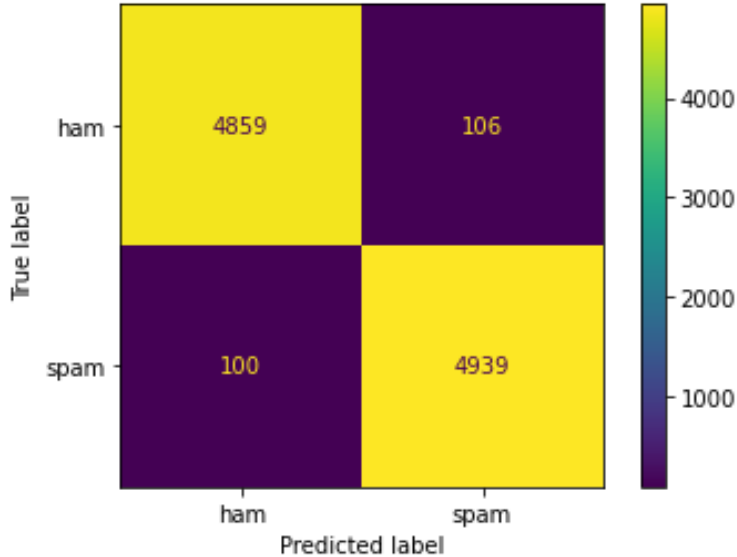
17. Figure 17

| Model | Sub Type | Feature space | Accuracy | Confusion Matrix | | | | | | | | | |
|------------------------------|------------------------------|--|------------------------------|--|------------------------------|-----|------|-----|------|-----|------|-----|------|
| KNN - Count Vectorizer (K=4) | Unigram | 1000 | 0.926 |  <table><tr><th>True label \ Predicted label</th><th>ham</th><th>spam</th></tr><tr><th>ham</th><td>4500</td><td>465</td></tr><tr><th>spam</th><td>280</td><td>4759</td></tr></table> | True label \ Predicted label | ham | spam | ham | 4500 | 465 | spam | 280 | 4759 |
| | True label \ Predicted label | | ham | spam | | | | | | | | | |
| | ham | | 4500 | 465 | | | | | | | | | |
| spam | 280 | 4759 | | | | | | | | | | | |
| Unigram+ Bi Gram | 0.925 |  <table><tr><th>True label \ Predicted label</th><th>ham</th><th>spam</th></tr><tr><th>ham</th><td>4488</td><td>477</td></tr><tr><th>spam</th><td>276</td><td>4763</td></tr></table> | True label \ Predicted label | ham | spam | ham | 4488 | 477 | spam | 276 | 4763 | | |
| True label \ Predicted label | ham | spam | | | | | | | | | | | |
| ham | 4488 | 477 | | | | | | | | | | | |
| spam | 276 | 4763 | | | | | | | | | | | |
| Unigram+ Bigram+Tri gram | 0.927 |  <table><tr><th>True label \ Predicted label</th><th>ham</th><th>spam</th></tr><tr><th>ham</th><td>4496</td><td>469</td></tr><tr><th>spam</th><td>265</td><td>4774</td></tr></table> | True label \ Predicted label | ham | spam | ham | 4496 | 469 | spam | 265 | 4774 | | |
| True label \ Predicted label | ham | spam | | | | | | | | | | | |
| ham | 4496 | 469 | | | | | | | | | | | |
| spam | 265 | 4774 | | | | | | | | | | | |

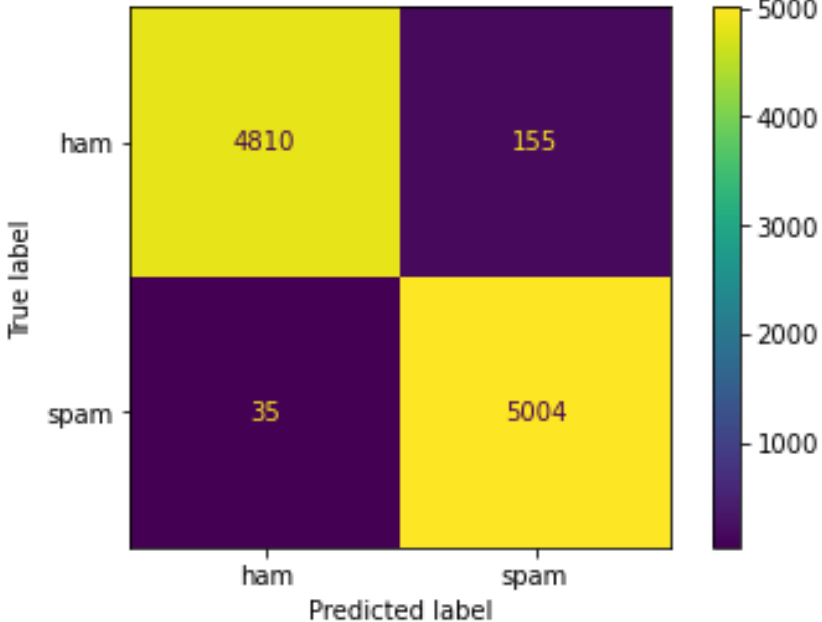
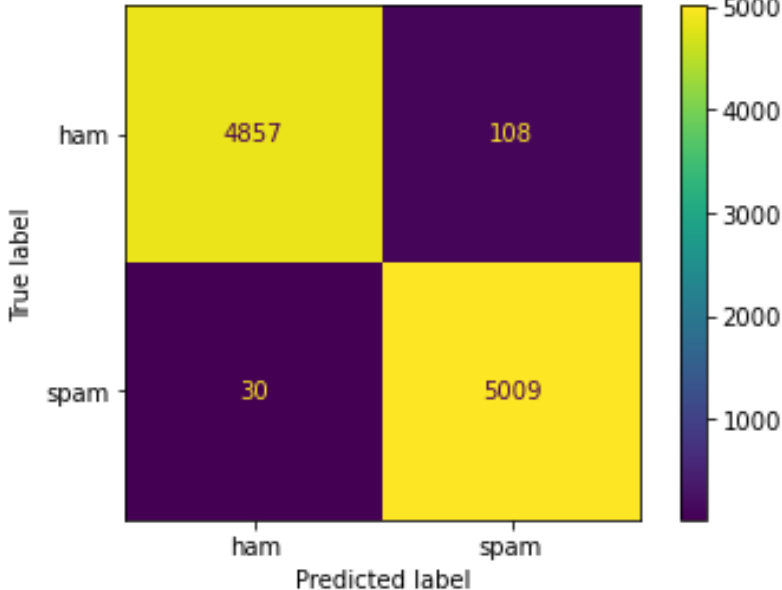
18. Figure 18

| Model | Sub Type | Feature space | Accuracy | Confusion Matrix |
|--------------------------|------------------------|---------------|----------|--|
| KNN - TF-IDF (K=4) | Unigram | 1000 | 0.805 |  |
| | Unigram+Bi Gram | | 0.816 |  |
| | Unigram+Bigram+Trigram | | 0.817 |  |

19. Figure 19

| Model | Sub Type | Feature space | Accuracy | Confusion Matrix | | | | | | | | | |
|------------------------------------|----------|------------------------------|----------|--|------------------------------|-----|------|-----|------|-----|------|-----|------|
| Naive Bayes - K Best TF-IDF | Unigram | 1000 | 0.974 |  <table><tr><th>True label \ Predicted label</th><th>ham</th><th>spam</th></tr><tr><th>ham</th><td>4830</td><td>135</td></tr><tr><th>spam</th><td>130</td><td>4909</td></tr></table> | True label \ Predicted label | ham | spam | ham | 4830 | 135 | spam | 130 | 4909 |
| | | True label \ Predicted label | ham | spam | | | | | | | | | |
| ham | 4830 | 135 | | | | | | | | | | | |
| spam | 130 | 4909 | | | | | | | | | | | |
| | | 5000 | 0.979 |  <table><tr><th>True label \ Predicted label</th><th>ham</th><th>spam</th></tr><tr><th>ham</th><td>4859</td><td>106</td></tr><tr><th>spam</th><td>100</td><td>4939</td></tr></table> | True label \ Predicted label | ham | spam | ham | 4859 | 106 | spam | 100 | 4939 |
| True label \ Predicted label | ham | spam | | | | | | | | | | | |
| ham | 4859 | 106 | | | | | | | | | | | |
| spam | 100 | 4939 | | | | | | | | | | | |

20. Figure 20

| Model | Sub Type | Feature space | Accuracy | Confusion Matrix | | | | | | | | | |
|------------------------------|----------|------------------------------|----------|---|------------------------------|-----|------|-----|------|-----|------|----|------|
| SVM - K Best TF-IDF | Unigram | 1000 | 0.981 |  <table><tr><td>True label \ Predicted label</td><td>ham</td><td>spam</td></tr><tr><td>ham</td><td>4810</td><td>155</td></tr><tr><td>spam</td><td>35</td><td>5004</td></tr></table> | True label \ Predicted label | ham | spam | ham | 4810 | 155 | spam | 35 | 5004 |
| | | True label \ Predicted label | ham | spam | | | | | | | | | |
| ham | 4810 | 155 | | | | | | | | | | | |
| spam | 35 | 5004 | | | | | | | | | | | |
| | | 5000 | 0.986 |  <table><tr><td>True label \ Predicted label</td><td>ham</td><td>spam</td></tr><tr><td>ham</td><td>4857</td><td>108</td></tr><tr><td>spam</td><td>30</td><td>5009</td></tr></table> | True label \ Predicted label | ham | spam | ham | 4857 | 108 | spam | 30 | 5009 |
| True label \ Predicted label | ham | spam | | | | | | | | | | | |
| ham | 4857 | 108 | | | | | | | | | | | |
| spam | 30 | 5009 | | | | | | | | | | | |

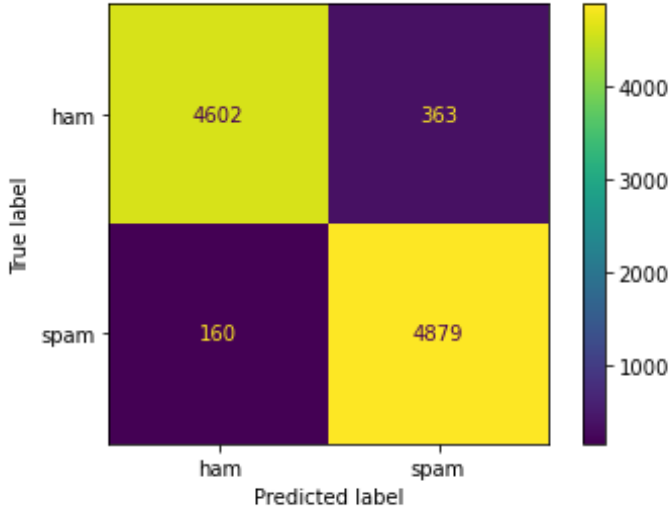
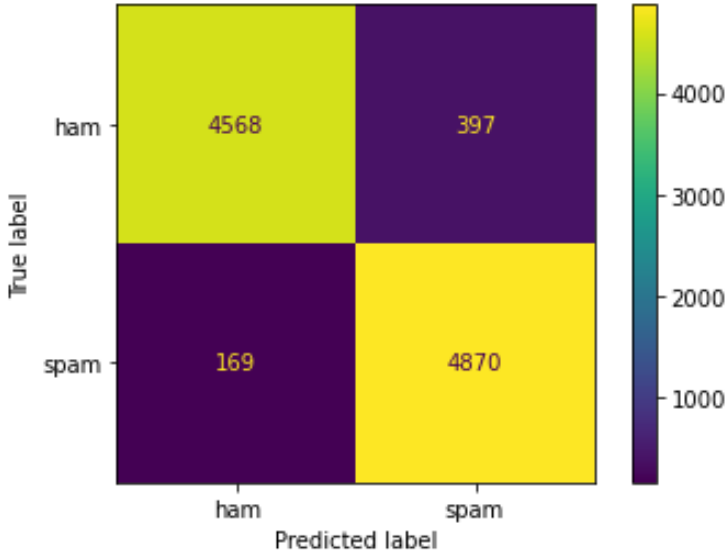
21. Figure 21

| Model | Sub Type | Feature space | Accuracy | Confusion Matrix | | | | | | | | | |
|--|------------|---------------|----------|--|------------------|-----|------|-----|------|-----|------|-----|------|
| Naive Bayes - Variance Threshold (0.001) | Unigram CV | 14,910 | 0.980 | <p>A confusion matrix heatmap for a binary classification task. The y-axis is labeled 'True label' with categories 'ham' and 'spam'. The x-axis is labeled 'Predicted label' with categories 'ham' and 'spam'. The matrix values are: True ham predicted ham: 4869; True ham predicted spam: 96; True spam predicted ham: 99; True spam predicted spam: 4940. A color bar on the right indicates counts from 0 to 4000, with yellow representing higher counts and purple representing lower counts.</p> <table><tr><th>True \ Predicted</th><th>ham</th><th>spam</th></tr><tr><th>ham</th><td>4869</td><td>96</td></tr><tr><th>spam</th><td>99</td><td>4940</td></tr></table> | True \ Predicted | ham | spam | ham | 4869 | 96 | spam | 99 | 4940 |
| True \ Predicted | ham | spam | | | | | | | | | | | |
| ham | 4869 | 96 | | | | | | | | | | | |
| spam | 99 | 4940 | | | | | | | | | | | |
| Naive Bayes - Variance Threshold (0.005) | Unigram CV | 6,380 | 0.976 | <p>A confusion matrix heatmap for a binary classification task. The y-axis is labeled 'True label' with categories 'ham' and 'spam'. The x-axis is labeled 'Predicted label' with categories 'ham' and 'spam'. The matrix values are: True ham predicted ham: 4852; True ham predicted spam: 113; True spam predicted ham: 122; True spam predicted spam: 4917. A color bar on the right indicates counts from 0 to 4000, with yellow representing higher counts and purple representing lower counts.</p> <table><tr><th>True \ Predicted</th><th>ham</th><th>spam</th></tr><tr><th>ham</th><td>4852</td><td>113</td></tr><tr><th>spam</th><td>122</td><td>4917</td></tr></table> | True \ Predicted | ham | spam | ham | 4852 | 113 | spam | 122 | 4917 |
| True \ Predicted | ham | spam | | | | | | | | | | | |
| ham | 4852 | 113 | | | | | | | | | | | |
| spam | 122 | 4917 | | | | | | | | | | | |

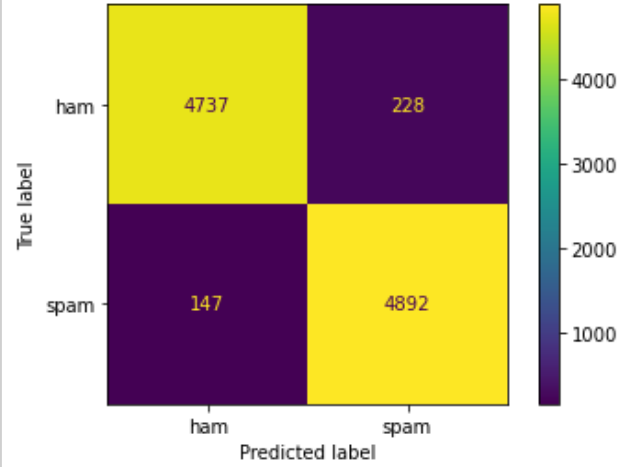
22. Figure 22

| Model | Sub Type | Feature space | Accuracy | Confusion Matrix | | | | | | | | | |
|----------------------------------|----------|---------------|----------|--|------------------------------|-----|------|-----|------|-----|------|----|------|
| SVM - Variance Threshold (0.001) | Unigram | 14,910 | 0.980 | <p>A 2x2 confusion matrix heatmap for the SVM model with a variance threshold of 0.001. The y-axis is labeled 'True label' with categories 'ham' and 'spam'. The x-axis is labeled 'Predicted label' with categories 'ham' and 'spam'. The matrix values are: True ham predicted ham: 4851; True ham predicted spam: 114; True spam predicted ham: 82; True spam predicted spam: 4957. A color bar on the right indicates counts from 0 to 4000, with yellow representing higher counts and purple representing lower counts.</p> <table><tr><td>True label \ Predicted label</td><td>ham</td><td>spam</td></tr><tr><td>ham</td><td>4851</td><td>114</td></tr><tr><td>spam</td><td>82</td><td>4957</td></tr></table> | True label \ Predicted label | ham | spam | ham | 4851 | 114 | spam | 82 | 4957 |
| True label \ Predicted label | ham | spam | | | | | | | | | | | |
| ham | 4851 | 114 | | | | | | | | | | | |
| spam | 82 | 4957 | | | | | | | | | | | |
| SVM - Variance Threshold (0.005) | Unigram | 6,380 | 0.978 | <p>A 2x2 confusion matrix heatmap for the SVM model with a variance threshold of 0.005. The y-axis is labeled 'True label' with categories 'ham' and 'spam'. The x-axis is labeled 'Predicted label' with categories 'ham' and 'spam'. The matrix values are: True ham predicted ham: 4833; True ham predicted spam: 132; True spam predicted ham: 82; True spam predicted spam: 4957. A color bar on the right indicates counts from 0 to 4000, with yellow representing higher counts and purple representing lower counts.</p> <table><tr><td>True label \ Predicted label</td><td>ham</td><td>spam</td></tr><tr><td>ham</td><td>4833</td><td>132</td></tr><tr><td>spam</td><td>82</td><td>4957</td></tr></table> | True label \ Predicted label | ham | spam | ham | 4833 | 132 | spam | 82 | 4957 |
| True label \ Predicted label | ham | spam | | | | | | | | | | | |
| ham | 4833 | 132 | | | | | | | | | | | |
| spam | 82 | 4957 | | | | | | | | | | | |

23. Figure 23

| Model | Sub Type | Feature space | Accuracy | Confusion Matrix | | | | | | | | | |
|-----------------------------------|----------------|---------------|----------|---|------------------------------|-----|------|-----|------|-----|------|-----|------|
| Decision Tree (GINI) Depth =35 | Unigram TF-IDF | 1000 | 0.948 |  <table><tr><td>True label \ Predicted label</td><td>ham</td><td>spam</td></tr><tr><td>ham</td><td>4602</td><td>363</td></tr><tr><td>spam</td><td>160</td><td>4879</td></tr></table> | True label \ Predicted label | ham | spam | ham | 4602 | 363 | spam | 160 | 4879 |
| True label \ Predicted label | ham | spam | | | | | | | | | | | |
| ham | 4602 | 363 | | | | | | | | | | | |
| spam | 160 | 4879 | | | | | | | | | | | |
| Decision Tree (Entropy) Depth =35 | Unigram TF-IDF | 1000 | 0.943 |  <table><tr><td>True label \ Predicted label</td><td>ham</td><td>spam</td></tr><tr><td>ham</td><td>4568</td><td>397</td></tr><tr><td>spam</td><td>169</td><td>4870</td></tr></table> | True label \ Predicted label | ham | spam | ham | 4568 | 397 | spam | 169 | 4870 |
| True label \ Predicted label | ham | spam | | | | | | | | | | | |
| ham | 4568 | 397 | | | | | | | | | | | |
| spam | 169 | 4870 | | | | | | | | | | | |

24. Figure 24

| Model | Sub Type | Feature space | Accuracy | Confusion Matrix |
|----------------|----------|---------------|----------|---|
| SVM - Word2Vec | - | 23,341 | 0.963 |  |

25. Table 1

| Model | Sub Type | Feature space | Micro | Macro | | | Weighted | | |
|--------------------------------|--------------------------|---------------|-----------------------|-----------|--------|------|-----------|--------|------|
| | | | Precision/ Recall/ F1 | Precision | Recall | F1 | Precision | Recall | F1 |
| Naive Bayes - Count Vectorizer | Unigram | 1000 | 0.961 | 0.96 | 0.96 | 0.96 | 0.96 | 0.96 | 0.96 |
| | Unigram +Bi Gram | | 0.961 | 0.96 | 0.96 | 0.96 | 0.96 | 0.96 | 0.96 |
| | Unigram +Bigram +Trigram | | 0.958 | 0.96 | 0.96 | 0.96 | 0.96 | 0.96 | 0.96 |
| Naive Bayes - TF-IDF | Unigram | 1000 | 0.966 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 |
| | Unigram +Bi Gram | | 0.965 | 0.96 | 0.96 | 0.96 | 0.96 | 0.96 | 0.96 |
| | Unigram +Bigram +Trigram | | 0.963 | 0.96 | 0.96 | 0.96 | 0.96 | 0.96 | 0.96 |
| SVM - Count Vectorizer | Unigram | 1000 | 0.969 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 |
| | Unigram +Bi Gram | | 0.968 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 |
| | Unigram +Bigram +Trigram | | 0.967 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 |

| | | | | | | | | | |
|------------------------------|--------------------------|------|-------|------|------|------|------|------|------|
| SVM - TF-IDF | Unigram | 1000 | 0.969 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 |
| | Unigram +Bi Gram | | 0.968 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 |
| | Unigram +Bigram +Trigram | | 0.968 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 |
| KNN - Count Vectorizer (K=4) | Unigram | 1000 | 0.926 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 |
| | Unigram +Bi Gram | | 0.925 | 0.93 | 0.92 | 0.92 | 0.93 | 0.92 | 0.92 |
| | Unigram +Bigram +Trigram | | 0.927 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 |
| KNN - TF-IDF (K=4) | Unigram | 1000 | 0.805 | 0.84 | 0.80 | 0.80 | 0.84 | 0.80 | 0.80 |
| | Unigram +Bi Gram | | 0.816 | 0.85 | 0.82 | 0.81 | 0.85 | 0.82 | 0.81 |
| | Unigram +Bigram +Trigram | | 0.817 | 0.85 | 0.82 | 0.81 | 0.85 | 0.82 | 0.81 |
| Naive Bayes - K Best | Unigram | 1000 | 0.974 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 |

| | | | | | | | | | |
|--|----------------|--------|-------|------|------|------|------|------|------|
| TF-IDF | Unigram | 5000 | 0.979 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 |
| SVM - K Best | Unigram | 1000 | 0.981 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 |
| TF-IDF | | 5000 | 0.986 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 |
| Naive Bayes - Variance Threshold (0.001) | Unigram CV | 14,910 | 0.980 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 |
| Naive Bayes - Variance Threshold (0.005) | Unigram CV | 6,380 | 0.976 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 |
| SVM - Variance Threshold (0.001) | Unigram | 14,910 | 0.980 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 |
| SVM - Variance Threshold (0.005) | Unigram | 6,380 | 0.978 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 |
| Decision Tree (GINI) Depth =35 | Unigram TF-IDF | 1000 | 0.948 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 |
| Decision Tree (Entropy) Depth =35 | Unigram TF-IDF | 1000 | 0.943 | 0.94 | 0.94 | 0.94 | 0.94 | 0.94 | 0.94 |
| SVM - Word2Vec | - | 23,341 | 0.963 | 0.96 | 0.96 | 0.96 | 0.96 | 0.96 | 0.96 |