

Capstone Project

Machine Learning Engineer Nanodegree

Ramkumar Singh

August 22, 2017

Contents

I.	Definition	2
A.	Project Overview	2
B.	Problem Statement.....	2
C.	Metric	2
II.	Analysis	3
A.	Data Exploration	3
B.	Exploratory Visualization	3
C.	Algorithms and Techniques	6
D.	Benchmark.....	6
III.	Methodology.....	7
A.	Data Pre-processing	7
B.	Implementation.....	7
C.	Refinement	11
IV.	Result	13
A.	Model Evaluation and Validation	13
B.	Justification.....	14
V.	Conclusion.....	16
A.	Free-Form Visualization	16
B.	Reflection	16
C.	Improvement.....	17

I. Definition

A. Project Overview

Spam filtering is a binary classification task familiar to any user of email services. We will use machine learning classifiers to implement a similar spam filter.

The task is to distinguish between two types of emails, "spam" and "non-spam" often called "ham". The machine learning classifier will detect that an email is spam if it is characterised by certain features. The textual content of the email – words like "Viagra" or "lottery" or phrases like "You've won 100,000,000 dollars! Click here!", "Join now!" – is crucial in spam detection and offers some of the strongest clues.

To train the classifier, we need a representative dataset with both spam and ham emails. In this project, we will use Apache SpamAssassin public corpus, which contains about 6047 message with approximate 30% spam ratio.

B. Problem Statement

Implement a Spam Filter with the help of a ML classifier which would classify a given mail as spam or ham

C. Metric

In real world scenario volume of ham mails are generally higher than the spam mails. In the given dataset, also the spam ratio is 30% only.

In this case Classification accuracy not enough since, it is the number of correct predictions made divided by the total number of predictions. So, a model that only predicted Ham mail in all case, would achieve an accuracy of approximately 70%. This is a high accuracy, but model is literally doing nothing and blindly classifying all mails as Ham.

For such problem f1score is a better metric to analyse performance of a classifier. It considers both the precision and the recall of the test to compute the score. The solution will evaluate f1score of ML model against the benchmark f1score

II. Analysis

A. Data Exploration

Apache SpamAssassin public corpus has spam and ham mails data. Below is the overview of the corpus data

URL:

<http://spamassassin.apache.org/old/publiccorpus/>

Content Folder:

- spam: 500 spam messages, all received from non-spam-trap sources.
- easy_ham: 2500 non-spam messages. These are typically quite easy to differentiate from spam, since they frequently do not contain any spammish signatures (like HTML etc).
- hard_ham: 250 non-spam messages which are closer in many respects to typical spam: use of HTML, unusual HTML mark-up, coloured text, "spammish-sounding" phrases etc.
- easy_ham_2: 1400 non-spam messages. A more recent addition to the set.
- spam_2: 1397 spam messages. Again, more recent.

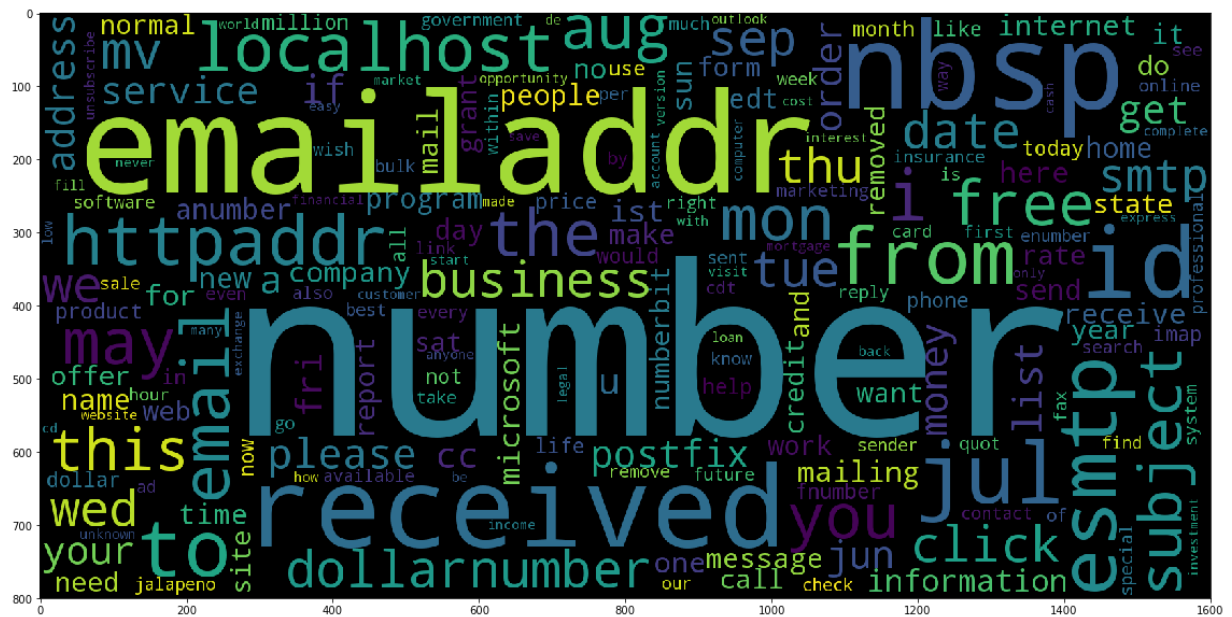
Total count: **6047** messages, with about a **31%** spam ratio.

B. Exploratory Visualization

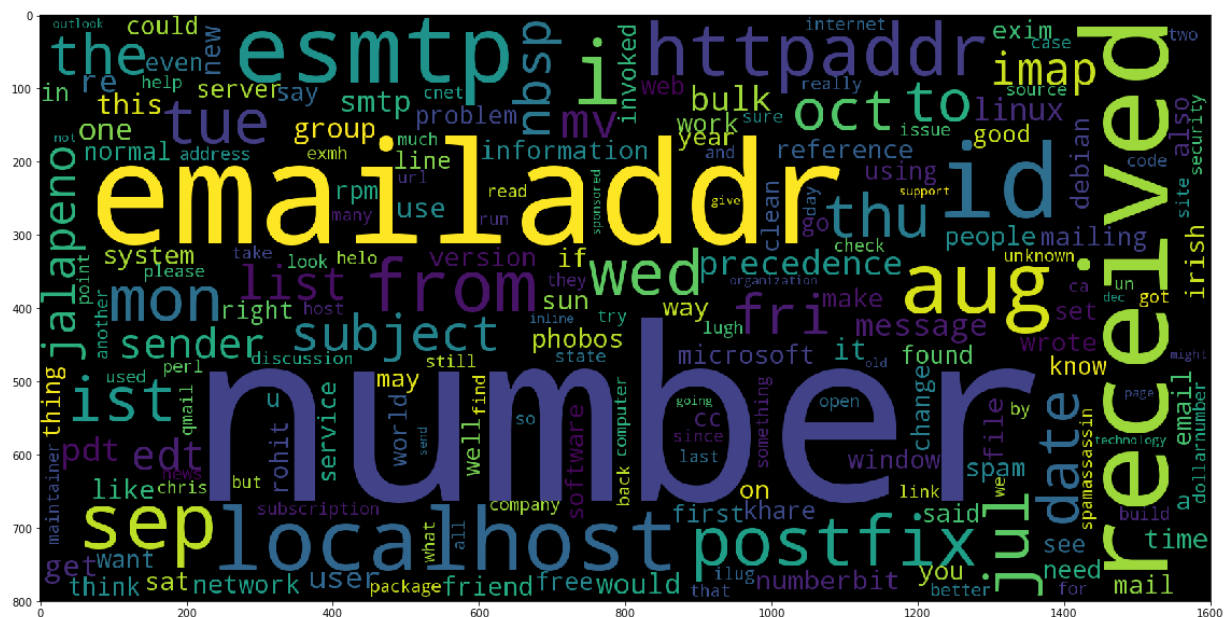
Emails are pre-processed to normalize, tokenize and lemmatize. After that, for each word, its frequency of appearance is calculated. This word and frequency data is used to create two type of plot for visualization.

1. Word cloud (tag cloud, or weighted list in visual design) is a visual representation of text data. Tags are usually single words, and the importance of each tag is shown with font size or color, which gives greater prominence to words that appear more frequently. Two separate word cloud map is created for spam and ham words. This visualization makes it easier to see relative frequency of words in spam and ham mails.

Spam Mail Word Cloud:

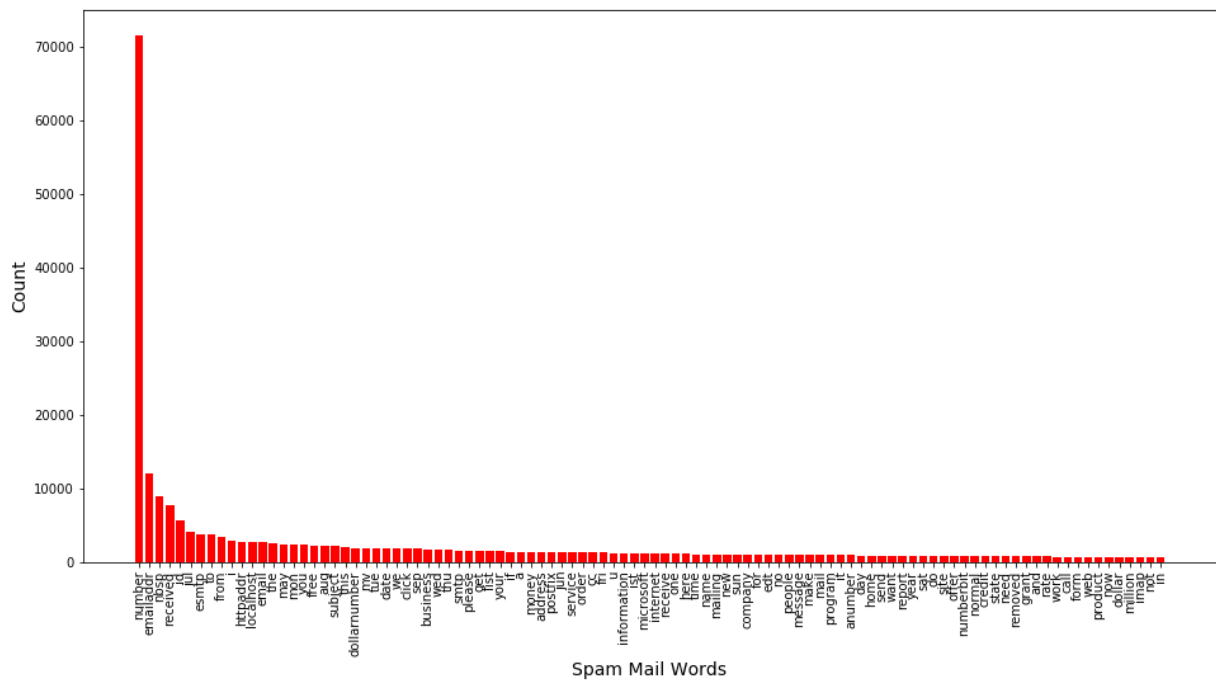


Ham Mail Word Cloud:

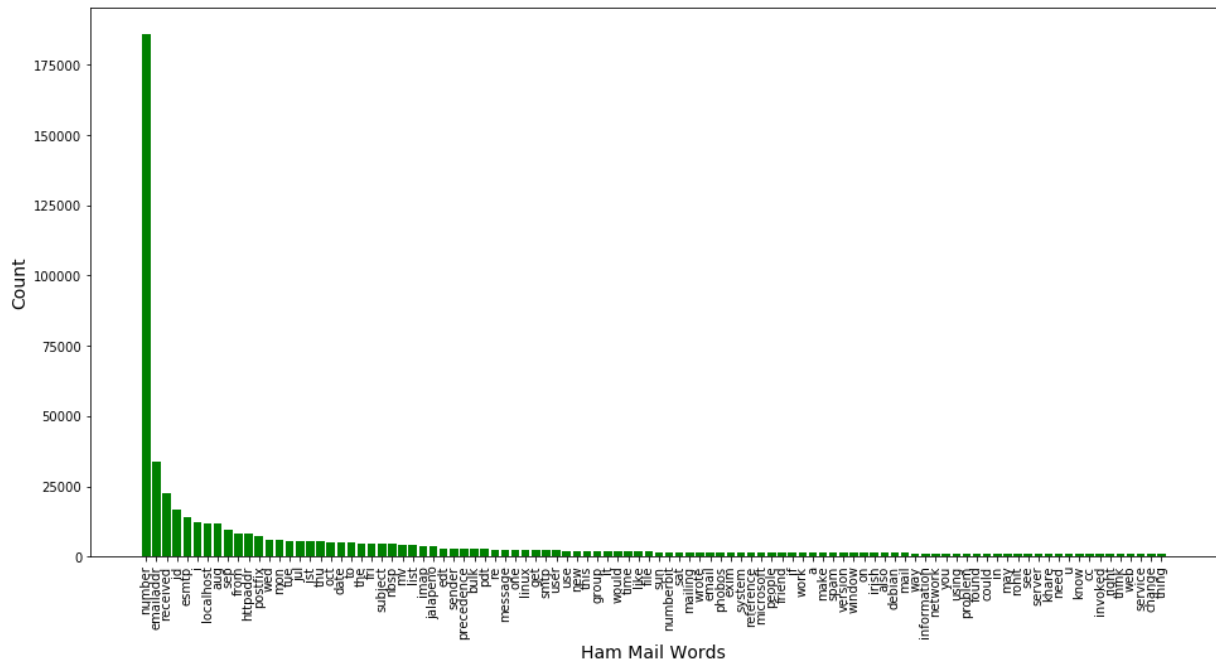


2. Frequency Histogram: 100 most frequent word and its frequency is plotted for spam email words and ham email words. Y axis has frequency count of top 100 high frequency words and X-axis has the corresponding word. This plot helps in visualizing the actual frequency of words in spam and ham mails.

Spam Mail Word Frequency Histogram (Top 100):



Ham Mail Word Frequency Histogram (Top 100):



C. Algorithms and Techniques

To create solution of the problem below are the high level of steps:

- Pre-processing the mail: - In this step we will process the mail content with following changes:
 - Remove html tags
 - Normalize numbers, Urls, email address and Dollar sign (Replace 0-9 with 'number, hyperlinks with 'httpaddr', any email address with 'emailaddr' and \$ with 'dollar')
 - Tokenize the words
 - Take only alphanumeric words
 - Lemmatize and change case to lower to reduce the words to their stemmed form.
 - Filter stop words (words which do not have significance like : to, the, a etc)
- Train test Split: - Create separate set for training and testing purpose..
- Extracting the features: - Use the training set data calculate frequency of each word and select high frequency word as feature to be used for learning
- Training different classifier :- Train the classifiers like Logistic Regression, K-Neighbour Classifier, SVM , Naïve Bayes etc
- Evaluating the classifiers: - The above classifiers will be evaluated against the dummy classifier, by comparing their prediction time and f1score.
- Fine Tune the best classifier :- Use GridSearch technique to fine tune the best classifier using different parameters like C, gamma and other applicable parameters

D. Benchmark

In this project, we will use the sklearn dummy classifier as a simple baseline to compare with. DummyClassifier is a classifier that makes predictions using simple rules, so the goal is to perform better than this.

III. Methodology

A. Data Pre-processing

Reading Data:

While reading the email from corpus file, opening the file with utf8 encoding. If UnicodeDecodeError occurs due to not supported character, we catch the exception and continue reading next file.

Pre-processing:

In data pre-processing step we normalize, tokenize, lemmatize and filter out stop words and non-alphanumeric words. Below are the details of pre-processing.

- Clean the mail text by removing html tags
- Normalize numbers, Urls, email address and Dollar sign (Replace numeric values with 'number' , hyperlinks with 'httpaddr', any email address with 'emailaddr' and \$ with 'dollar')
- Tokenize the words and take only alphanumeric words
- Filter stop words (words which do not have significance like : to, the, a etc)
- Lemmatize and change case to lower to reduce the words to their stemmed form.

Feature Extraction:

We split the data in training and testing sets. To generate NLP features, we use ONLY the training data. While working on GridSearchCV and cross validation, to easily and correctly implement the feature extraction, we are using Sklearn pipelines.

Some of the classifier need dense matrix but the Pipeline produces sparse matrix, so sklearn.preprocessing FunctionTransformer is used to generate dense matrix. This option influences memory footprint.

B. Implementation

Pre-processing helper methods:

- ***preprocess_normalize*** :

This method is implemented to normalize the email data using regex rules, which removes html tags and replace numeric values with 'number' , hyperlinks with 'httpaddr', any email address with 'emailaddr' and \$ with 'dollar'

- ***preprocess_lemmatize_tokenize:***

This method is using ***nltk word_tokenize*** to tokenize and ***nltk WordNetLemmatizer*** to lemmatize the email words. It also filters stopwords using ***nltk.corpus stopwords***. This method also converts email words to lower case and takes only alpha-numeric words to return at the end.

- ***LemmatizerTokenizer:***

A custom tokenizer class is implemented which Lemmetize, Tokenize and perform other optimizations using the ***preprocess_lemmatize_tokenize*** method mentioned above

Feature Extraction:

- **TF-IDF features:**

For feature extraction, word occurrence count has an issue; longer documents will have higher average count values than shorter documents, even though they might talk about the same topics.

To avoid these potential discrepancies, it suffices to divide the number of occurrences of each word in a document by the total number of words in the document: these features are called tf (Term Frequencies).

Another refinement on top of tf is to downscale weights for words that occur in many documents in the corpus and are therefore less informative than those that occur only in a smaller portion of the corpus. This downscaling is called tf-idf (Term Frequency Times Inverse Document Frequency).

To convert a collection of raw documents to a matrix of TF-IDF features, I am using `sklearn.feature_extraction.text TfidfVectorizer`.

- **Custom Tokenizer:**

LemmatizerTokenizer, the custom tokenizer class is passed to *TfidfVectorizer* to be used during feature extraction.

- **Pipeline:**

We need to generate the NLP features on training data only. This is important because if we generate NLP feature on all of the data and then split the data in train and test set, we have contaminated the test.

Since the NLP features will have already TF-IDF/stats of testing data as well. This separation of training data and testing data for feature generation and testing becomes complicated with cross validation methods like GridSearchCV, since each new fold need a new set of features. To achieve this separation we are using ***sklearn.pipeline***

- **Function Transformer:**

Sklearn Pipeline produces sparse matrix, but some of the classifier need dense matrix due to which the pipeline.fit method throws error "*A sparse matrix was passed, but dense data is required*". To overcome this issue ***sklearn.preprocessing.FunctionTransformer*** is passed as an argument to ***sklearn.pipeline.make_pipeline***. This Function Transformer turns a Python function into a Pipeline-compatible transformer object and hence used to convert sparse matrix to dense

Sklearn Classifiers

Following *scikit-learn* supervised learning models are used to train the classifier.

- Gaussian Naive Bayes (GaussianNB)
- Logistic Regression
- K-Nearest Neighbors (KNeighbors)
- Support Vector Machines (SVM)

- **Gaussian Naïve Bayes**

Why: Naive Bayes is a simple technique for constructing classifiers that assign class labels to problem instances, represented as vectors of feature values. Despite their naive design and apparently oversimplified assumptions, naive Bayes classifiers have worked quite well in many complex real-world situations. It only requires a small number of training data to estimate the parameters necessary for classification. Also, the Naive Bayes has been used real Spam classification system.

Advantage: Naive Bayes classifiers are highly scalable, requiring parameters linear in the number of variables (features/predictors) in a learning problem. Maximum-likelihood training can be done by evaluating a closed-form expression, which takes linear time, rather than by expensive iterative approximation as used for many other types of classifiers

Disadvantage: Due the feature independence assumption, it loses the ability to exploit interaction between features, however for classification task this is often is not a problem.

Source: https://en.wikipedia.org/wiki/Naive_Bayes_classifier

- **Logistic Regression**

Why: Logistic Regression classification model is a probabilistic model which maximizes the posterior class probability. It is also suitable for binomial classification problem like this. Also, as it is a high bias/low variance model, it works well when we have a large number of features in comparison to the number of data/sample. So, it is a good choice for the problem at hand where the number of data is low but the feature space is relatively large.

Advantage: If there's a lot of noise, logistic regression can handle it better. Logistic regression is intrinsically simple, it has low variance and so is less prone to overfitting.

Disadvantage: Unlike SVM, it considers all the points in the data set. This may not be a preferred approach for certain problems.

Source: https://en.wikipedia.org/wiki/Logistic_regression

- **K-Nearest Neighbours (KNeighbors)**

Why: Because the training set size is small so implementing instance based learning such as KNN classification can be a good candidate for such scenario. Also, we do not have to worry about the linear separability of the data and implementing the model is very easy as well.

Advantage: Makes no assumption about the data distribution. Learning time is very low. Immediately adapts to the new data not affected by linear separability of the data. An easy algorithm to explain and implement.

Disadvantage: Slow during prediction. So difficult to use this for prediction in real time. Storage space can be a challenge if working on more data. Performance impact is high if more features are added (Curse of Dimensionality).

Source: <http://www.dummies.com/programming/big-data/data-science/solving-real-world-problems-with-nearest-neighbor-algorithms/>

- **SVM**

Why: SVM performs well for most of the cases. It can be used with various Kernel to fit the nonlinear decision boundary as well. Various parameters can be tuned to get a high level of accuracy.

Advantage High accuracy. By implementing Kernel Trick, we can inject the domain knowledge in the classifier. With appropriate kernel, it can work well even if the data is not linearly separable

Disadvantage. Tuning the model for optimal parameters can be difficult. Memory-intensive. Hard to interpret

Source: https://en.wikipedia.org/wiki/Support_vector_machine

Training

A helper function `train_predict` is implemented for training and testing the four supervised learning models we have chosen above.

The function takes as input a classifier, and the training and testing data. It creates a sklearn pipeline using the vectorizer created in data processing step. Using this pipeline, the model is trained.

This function will report the training time, prediction time and `f1score` for both the training and testing data separately.

We import the dummy classifier and four other supervised learning models mentioned earlier, and run the ***train_predict*** function for each one. We will perform following action:

- From the `X_train` data set create four different training set of sizes: 25%, 50%, 75% and 100% of data.
- Import the supervised learning models discussed in the previous section.
- Initialize the three models and store them in `clf_A`, `clf_B`, `clf_C` and `clf_D`
 - Use a `random_state` for each model.
 - Use the default settings for each model — we will tune one specific model in a later section.
- Fit each model with each training set size and make predictions on the test set.

C. Refinement

Initially I trained the classifier on sklearn.svm SVC which is not giving a meaningful result. `F1_score` is quite low in comparison to other classifiers. Below is the performance matrix for the same

A regular SVC with default values uses a radial basis function as the kernel. This nonlinear basis function has higher variance. We can add regularization to the nonlinear model and we will probably see much better results. (This is the `C` parameter in scikit learn SVMs), however I decided to use a LinearSVC instead since we are training on default parameters

and linear kernel has lower variance by default. LinearSVC seems a better choice because we have large feature space in comparison to the training dataset.

Model Tuning:

In this step, we fine tune the chosen model. We have used grid search (GridSearchCV) with parameters tuned with at least 3 different values. We have to use the entire training set for this.

Following are the steps implemented for Model Tuning

- Import `sklearn.grid_search.GridSearchCV` and `sklearn.metrics.make_scorer`.
- Initialize the classifier chosen for fine tuning
- Create the parameters list to tune
- Make an f1 scoring function using 'make_scorer' and store it in `f1_scorer`
- Initialize Pipeline using the vectorizer, FunctionTransformer and classifier
- Perform grid search on the classifier using the pipeline, parameters and `f1_scorer` as the scoring method and store it in `grid_obj`
- Fit the `grid_object` to the training data and find the optimal parameters,
- Get the best estimator from the `grid_object`
- Report the final F1 score for training and testing using the best estimator

IV. Result

A. Model Evaluation and Validation

Below are the 16 different outputs— 4 for each model using the varying training set sizes.

Training Classifier: GaussianNB				
Training Set Size	Training Time	Prediction Time (test)	f1 score (training)	f1 score (testing)
1107	12.021	10.775	0.9985	0.7766
2214	19.214	10.936	0.9992	0.8402
3321	29.743	11.136	0.9985	0.8529
4429	40.91	11.493	0.9989	0.8677

Training Classifier: LogisticRegression				
Training Set Size	Training Time	Prediction Time (test)	f1 score (training)	f1 score (testing)
1107	9.893	9.509	0.894	0.8306
2214	18.691	10.193	0.9298	0.904
3321	28.036	10.767	0.9453	0.9219
4429	37.275	9.884	0.9573	0.9329

Training Classifier: KNeighborsClassifier				
Training Set Size	Training Time	Prediction Time (test)	f1 score (training)	f1 score (testing)
1107	9.299	60.775	0.8486	0.7885
2214	20.927	166.08	0.8993	0.8311
3321	33.955	274.445	0.9123	0.8576
4429	48.149	493.512	0.9333	0.8826

Training Classifier: LinearSVC				
Training Set Size	Training Time	Prediction Time (test)	f1 score (training)	f1 score (testing)
1107	11.533	9.919	0.997	0.9601
2214	18.27	9.88	0.9985	0.9721
3321	29.858	9.749	0.998	0.9736
4429	38.874	10.208	0.9985	0.9796

In the above table gives we can see that KNN classification is very slow as per expected and discussed in the Sklearn Classifiers section of implementation. Gaussian Naïve Bayes and Logistic Regression Classifier are having similar training time and prediction time, but Logistic Regression Classifiers has higher score. Overall Linear SVC model seems to be the best classifier, because the model is giving the highest testing F1 score among the other models. Training and prediction time is also comparable to Logistic Regression and Gaussian NB models in the experiment.

Tuning the Model

LinearSVC, which is the best classifier as per the experiment is tuned further. This results a final model having following F1 Scores:

- Training F1 score of 0.9993.
- Testing F1 score of 0.9855.

Best Classifier Parameters

Classifier Best Parameters: {'linearsvc__C': 10, 'tfidfvectorizer__ngram_range': (1, 1)}

Confusion Matrix:

	Predicted		
		Ham	Spam
	Actual		
	Ham	3841	4
	Spam	8	1684

Sensitivity Analysis

We are analyzing performance of the model while working on different set of data. The final classifier gives following scores on different set of data:

- F1 score on 25% of training data 0.9985.
- F1 score on 50% of training data 0.9992.
- F1 score on 75% of training data 0.9990.
- F1 score on entire data 0.9964

The above data shows that the final classifier is robust and performance is not affected much by the changes in input data.

B. Justification

The dummy classifier offers performance is as below. The F1 Score for testing and training is quite low. This score represents performance of a model which do not learn anything.

Training Classifier: DummyClassifier , Time in Seconds				
Training Set Size	Training Time	Prediction Time (test)	f1 score (training)	f1 score (testing)
1107	9.316	10.508	0.2888	0.3127
2214	18.956	11.844	0.298	0.3014
3321	28.462	10.744	0.2981	0.2917
4429	39.874	9.797	0.2996	0.315

The final selected LinearSV model has below performance matrix over the same set of data.

Training Classifier: LinearSVC				
Training Set Size	Training Time	Prediction Time (test)	f1 score (training)	f1 score (testing)
1107	11.533	9.919	0.997	0.9601
2214	18.27	9.88	0.9985	0.9721
3321	29.858	9.749	0.998	0.9736
4429	38.874	10.208	0.9985	0.9796

The result shows that the LinearSVC model has learnt significantly from the training data set and able to predict lot better on the test data set. The fine-tuned Linear SVC model perform marginally better that the LinearSVC with default parameters. The fine-tuned classifier is giving training F1 score of 0.9993 and Testing score of 0.9855.

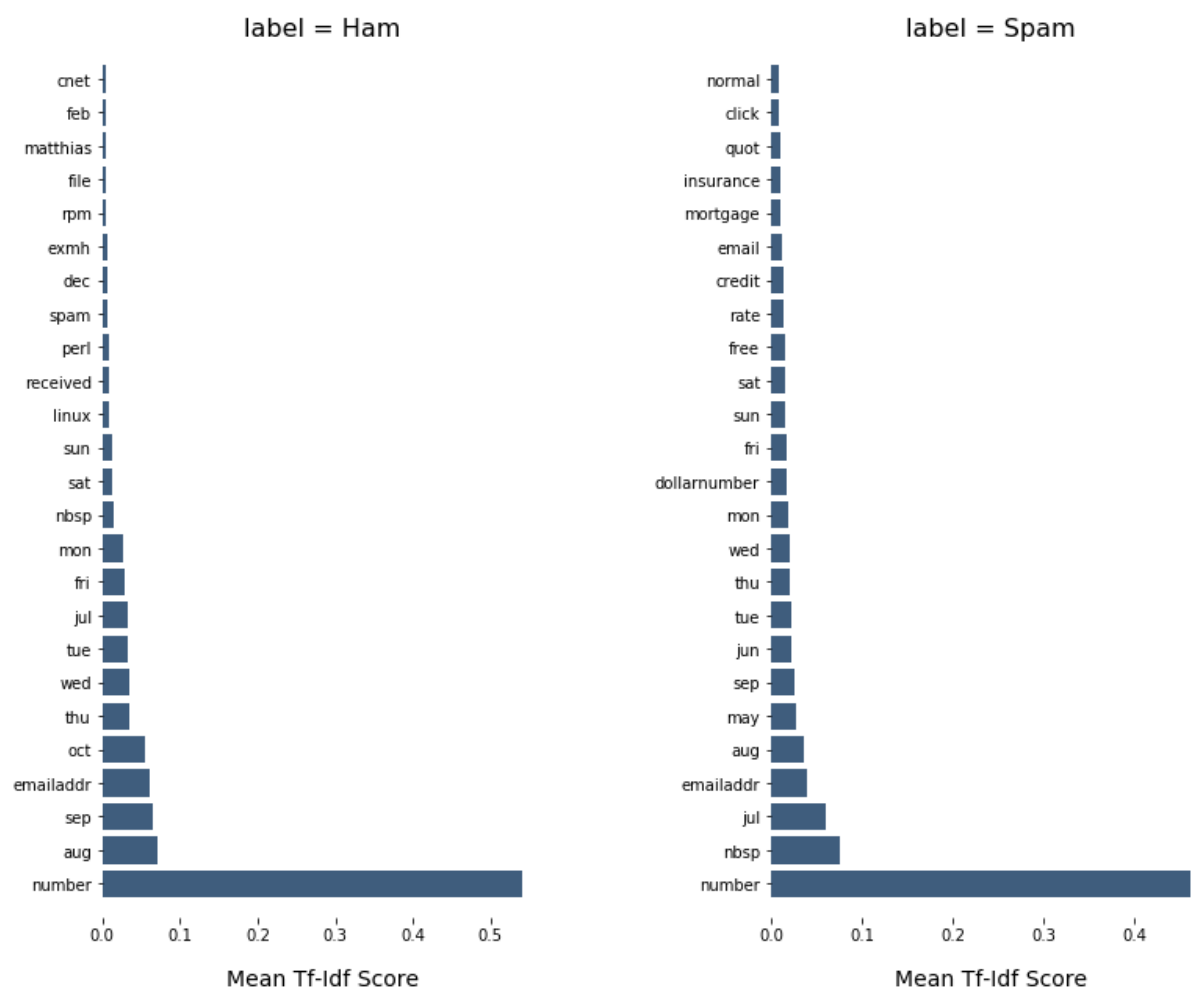
By looking at the confusion matrix, it is evident that out of 3845 ham mails, model is able to successfully classify 3841 emails and out of 1692 ham mails, the model is able to classify 1684 emails. The false positive count is 4 and false negative is 8. Cases of false positive and false negative are quite low.

V. Conclusion

A. Free-Form Visualization

TF-IDF Analysis:

TF-IDF vectorization of a corpus of text documents assigns each word in a document a number that is proportional to its frequency in the document and inversely proportional to the number of documents in which it occurs. Very common words, thereby receive heavily discounted tf-idf scores, in contrast to words that are very specific to the document in question. Visualizing such TF-IDF score for Spam and Ham mails below



B. Reflection

Cleaning the data and preprocessing step to normalize, tokenize and lemmatize the input data is quite interesting. There is lot of preprocessing needed to clean and transform the data, to make it usable to work with feature extraction and model training.

Extracting the feature using vectorization and implementing pipeline for gridsearchCV was not intuitive earlier. Got to know about these during proposal review and hence learned about the limitation of earlier approach.

Also, earlier I was relying on word frequency count for feature selection and relying on an empirical value to decide the selection cut-off. While learning about vectorization, came across the TF-IDF (Term Frequency Times Inverse Document Frequency), which made more sense to me and hence I have used the same in this project.

C. Improvement

While researching about Email Spam classifiers, I tried to find out how this is implemented in the email providers like Gmail or Hotmail. I came to know that in 2015, Google added Deep Learning capabilities to the Gmail spam filtering system and it has improved their accuracy in a great deal.

Deep Learning achieves great power and flexibility by learning to represent the world as nested hierarchy of concepts. It also automatically finds out the features which are important for classification, but deep learning algorithms need a large amount of data to understand it perfectly.

With large amount of data and implementing a Deep Learning classifier we may achieve a better result.

Also, the classifiers in this project just analyze the textual content of the email. Any information about sender's history and reputation can help in achieving better performance.