

Preconceptions

Dataset = Started with Enron corpus dataset, approximately 500k emails after cleanup ended with 16,769 emails.

Examples: $X := \mathbb{R}^d$ where d holds for the dimension of the features.

Features

1. Content-Length: Length of the email body, for example:

"Hello Dan, how are you?" the content length is 25.

2. Number Of Recipients: number of recipients of the email.

The following two features is word embedding features:

- words are converted into d -dimensional vector.

- The average of the words vector creates single d dimensional vector.

3. Email Subject Embedding, Example: Subject: "Meeting Tomorrow @ 10AM"

Word2Vec vectors: $V_{meeting}, V_{Tomorrow}, V_{@}, V_{10AM} \in \mathbb{R}^d$, the final subject vec:

$$V_{subject} = \frac{V_{meeting} + V_{Tomorrow} + V_{@} + V_{10AM}}{4} \in \mathbb{R}^d$$

4. Email body Embedding: Same approach as Subject Embedding:

$$V_{Body} = \sum_{i=1}^n \frac{V_{word_i}}{n} \in \mathbb{R}^d \quad | \text{where } |Body| = n$$

Labels

Created 7 labels in order to categorize:

Personal

hr

meeting and scheduling

Operations and logistics

corporate and legal

Projects

Finance

note: We started off with 9 labels as w/ progress we remove two features (all details and justification in the code).

Hypothesis class H

- consists of all possible mapping from X to the selected labels.

$x = \text{features}, x \in \mathbb{R}^d$

$y = \text{Labels}$

$h \in H = \text{trained classifier, defined as } h: x \rightarrow y$

The process ahead

Dataset

Enron dataset is a huge dataset of emails, that requires gentle and strict cleanup.

Owing to the fact we cannot inspect the emails manually, much noise and outliers is expected.

Labels

Given that the emails are organized as folders, after the cleanup we should remain with properly categorized emails to the respective folders.

Eventually the folders rule is to be our "ground truth" A.K.A Labels.

note: we will create new folders by merging semantic related folders (the process and details in the code)

Features

Crucial part is the analysis of the email text, hence much consideration given to the NLP model that handle vast amount of data while maintaining the semantic relation of the text.

Models

Given we have labels, supervised models of course will be used.

Since our objective is to categorize emails, classification models naturally comes to mind. we have chosen to analyze the text with the following:

- Logistic regression
- Support vector machine (svm)
- Random forest

Multi class Logistic Regression

In this model, we demand multi-class classification and not binary one. Two approaches were considered:

1. One-vs-All

- Train 7 separate classifiers. Very expensive on high dimensional data as each embedding.
- choose the class with highest probability.

For each class k the probability of an email belong to class k :

$$\text{let } z = w_k^T x + b_k \text{ where:}$$

- $x \in \mathbb{R}^d$ feature vector,
- $w_k \in \mathbb{R}^d$ weight vector for class $k \in \{1, 2, \dots, 9\}$
- b_k is the bias for class k
- recall sigmoid function: $\sigma(z) = \frac{1}{1 + e^{-z}}$

$$P(y=k | x) = \sigma(z) = \frac{1}{1 + \exp(-(w_k^T x + b_k))}$$

2. Softmax Regression

- Train single classifier for all classes simultaneously
- calculate the probabilities for $x \in \mathbb{R}^d$ that input x belong to each class $k \in C$ as C holds for the set of classes.
- Probabilities normalized by softmax, sum to 1.

- Softmax generalization the sigmoid function for multi class classification

$$P(y=k|x) = \frac{\exp(w_k^T x + b_k)}{\sum_{j=1}^{|C|} \exp(w_j^T x + b_j)}$$

In conclusion

as Softmax train one classifier instead $|C|$ classifiers,

the computational cost of softmax is much faster.

Thanks to word embeddings, which increase the feature space,

computational speed takes place as primary factor,

given that we have selected the softmax approach.

Random forest

Due to the noise and imbalance nature of the dataset

linear predictors struggle to capture meaningful patterns here. Owing to the Random forest randomness:

- Rf is more robust to noisy dataset compared to linear predictors.
- Rf is more robust to overfitting than using K decision trees with OVE (one vs all) approach.

On the other hand, emails are represented using word embedding,

leading to high number of features that Random forest struggles with.

We keep in mind that text-based classification problems often result in complex non-linear decision boundaries.

In conclusion

We expect Rf to deliver reasonable results, however we keep in mind that

due to the curse of dimensionality its performance may be overshadowed by other models better suited for high-dimensional analysis.

Support Vector Machine

When considering SVM for this dataset, a clear approach emerges:

Leveraging Soft SVM take advantage of its qualities such as:

- Soft SVM more robust to noisy data.

- In real world email classification, Categories often overlap, Hard SVM fails in such cases while Soft SVM more flexible and robust.

- SVM is well-equipped to handle complexity of high-dimensional dataset, such as those found in Enron dataset.

In conclusion

We expect Soft-SVM will perform well due its clear advantages regarding Enron dataset.