

Twitter Sentiment Analysis

Abstract

In this report, we address the problem of sentiment classification on twitter dataset. We use a number of machine learning and deep learning methods to perform sentiment analysis. In the end, we use a majority vote ensemble method with 5 of our best models to achieve the classification accuracy of 83.58% on kaggle public leaderboard.

1 Problem Statement

Twitter is a popular social networking website where members create and interact with messages known as “tweets”. This serves as a mean for individuals to express their thoughts or feelings about different subjects. Various different parties such as consumers and marketers have done sentiment analysis on such tweets to gather insights into products or to conduct market analysis. Furthermore, with the recent advancements in machine learning algorithms, we are able improve the accuracy of our sentiment analysis predictions.

In this report, we will attempt to conduct sentiment analysis on “tweets” using various different machine learning algorithms. We attempt to classify the polarity of the tweet where it is either positive or negative. If the tweet has both positive and negative elements, the more dominant sentiment should be picked as the final label.

We use the dataset from [Kaggle](#) which was crawled and labeled positive/negative. The data provided comes with emoticons, usernames and hashtags which are required to be processed and converted into a standard form. We also need to extract useful features from the text such uni-grams and bigrams which is a form of representation of the “tweet”.

We use various machine learning algorithms to conduct sentiment analysis using the extracted features. However, just relying on individual models did not give a high accuracy so we pick the top few models to generate a model ensemble. Ensembling is a form of meta learning algorithm technique where we combine different classifiers in order to improve the prediction accuracy. Finally, we report our experimental results and findings at the end.

2 Data Description

The data given is in the form of a comma-separated values files with tweets and their corresponding sentiments. The training dataset is a csv file of type `tweet_id,sentiment,tweet` where the `tweet_id` is a unique integer identifying the tweet, `sentiment` is either 1 (positive) or 0 (negative), and `tweet` is the tweet enclosed in ". Similarly, the test dataset is a csv file of type `tweet_id,tweet`.

The dataset is a mixture of words, emoticons, symbols, URLs and references to people. Words

	Total	Unique	Average	Max	Positive	Negative
Tweets	800000	-	-	-	400312	399688
User Mentions	393392	-	0.4917	12	-	-
Emoticons	6797	-	0.0085	5	5807	990
URLs	38698	-	0.0484	5	-	-
Unigrams	9823554	181232	12.279	40	-	-
Bigrams	9025707	1954953	11.28	-	-	-

Table 1: Statistics of preprocessed train dataset

	Total	Unique	Average	Max	Positive	Negative
Tweets	200000	-	-	-	-	-
User Mentions	97887	-	0.4894	11	-	-
Emoticons	1700	-	0.0085	10	1472	228
URLs	9553	-	0.0478	5	-	-
Unigrams	2457216	78282	12.286	36	-	-
Bigrams	2257751	686530	11.29	-	-	-

Table 2: Statistics of preprocessed test dataset

and emoticons contribute to predicting the sentiment, but URLs and references to people don't. Therefore, URLs and references can be ignored. The words are also a mixture of misspelled words, extra punctuations, and words with many repeated letters. The tweets, therefore, have to be preprocessed to standardize the dataset.

The provided training and test dataset have 800000 and 200000 tweets respectively. Preliminary statistical analysis of the contents of datasets, after preprocessing as described in section 3.1, is shown in tables 1 and 2.

3 Methodology and Implementation

3.1 Pre-processing

Raw tweets scraped from twitter generally result in a noisy dataset. This is due to the casual nature of people's usage of social media. Tweets have certain special characteristics such as re-tweets, emoticons, user mentions, etc. which have to be suitably extracted. Therefore, raw twitter data has to be normalized to create a dataset which can be easily learned by various classifiers. We have applied an extensive number of pre-processing steps to standardize the dataset and reduce its size. We first do some general pre-processing on tweets which is as follows.

- Convert the tweet to lower case.
- Replace 2 or more dots (.) with space.
- Strip spaces and quotes (" and ') from the ends of tweet.
- Replace 2 or more spaces with a single space.

We handle special twitter features as follows.

3.1.1 URL

Users often share hyperlinks to other webpages in their tweets. Any particular URL is not important for text classification as it would lead to very sparse features. Therefore, we replace all the URLs in tweets with the word URL. The regular expression used to match URLs is `((www\.[\S]+)|(https?://[\S]+))`.

3.1.2 User Mention

Every twitter user has a handle associated with them. Users often mention other users in their tweets by `@handle`. We replace all user mentions with the word `USER_MENTION`. The regular expression used to match user mention is `@[\S]+`.

Emoticon(s)	Type	Regex	Replacement
:), :), :-), (:, (:, (-:, :')	Smile	(:\s?\) :-\) \(\s?: \(-: :\'\))	EMO_POS
:D, : D, :-D, xD, x-D, XD, X-D	Laugh	(:\s?D :-D x-?D X-?D)	EMO_POS
;-), ;), ;-D, ;D, (;, (-;	Wink	(:\s?\(:-\(\)\s?: \)\-:)	EMO_POS
<3, :*	Love	(<3 :*)	EMO_POS
:-(:, : (, :(:,):,)-:	Sad	(:\s?\(:-\(\)\s?: \)\-:)	EMO_NEG
:(, :'(, :"(Cry	(:,\(:\'\(:"\()	EMO_NEG

Table 3: List of emoticons matched by our method

3.1.3 Emoticon

Users often use a number of different emoticons in their tweet to convey different emotions. It is impossible to exhaustively match all the different emoticons used on social media as the number is ever increasing. However, we match some common emoticons which are used very frequently. We replace the matched emoticons with either `EMO_POS` or `EMO_NEG` depending on whether it is conveying a positive or a negative emotion. A list of all emoticons matched by our method is given in table 3.

3.1.4 Hashtag

Hashtags are unspaced phrases prefixed by the hash symbol (`#`) which is frequently used by users to mention a trending topic on twitter. We replace all the hashtags with the words with the hash symbol. For example, `#hello` is replaced by `hello`. The regular expression used to match hashtags is `#(\S+)`.

3.1.5 Retweet

Retweets are tweets which have already been sent by someone else and are shared by other users. Retweets begin with the letters RT. We remove RT from the tweets as it is not an important feature for text classification. The regular expression used to match retweets is `\brt\b`.

After applying tweet level pre-processing, we processed individual words of tweets as follows.

- Strip any punctuation [`"?!,.():;`] from the word.
- Convert 2 or more letter repetitions to 2 letters. Some people send tweets like *I am soooooo happy* adding multiple characters to emphasize on certain words. This is done to handle such tweets by converting them to *I am soo happy*.
- Remove - and '. This is done to handle words like t-shirt and their's by converting them to the more general form tshirt and theirs.
- Check if the word is valid and accept it only if it is. We define a valid word as a word which begins with an alphabet with successive characters being alphabets, numbers or one of dot (.) and underscore(_).

Some example tweets from the training dataset and their normalized versions are shown in table 4.

3.2 Feature Extraction

We extract two types of features from our dataset, namely unigrams and bigrams. We create a frequency distribution of the unigrams and bigrams present in the dataset and choose top N unigrams and bigrams for our analysis.

3.2.1 Unigrams

Probably the simplest and the most commonly used features for text classification is the presence of single words or tokens in the the text. We extract single words from the training dataset and create a frequency distribution of these words. A total of 181232 unique words are extracted from

Raw	misses Swimming Class. http://plurk.com/p/12nt0b
Normalized	misses swimming class URL
Raw	@98PXYRochester HEYYYYYYYY!! its Fer from Chile again
Normalized	USER_MENTION hey its fer from chile again
Raw	Sometimes, You gotta hate #Windows updates.
Normalized	sometimes you gotta hate windows updates
Raw	@Santiago_Steph hii come talk to me i got candy :)
Normalized	USER_MENTION hii come talk to me i got candy EMO_POS
Raw	@bolly47 oh no :(r.i.p. your bella
Normalized	USER_MENTION oh no EMO_NEG r.i.p your bella

Table 4: Example tweets from the dataset and their normalized versions.

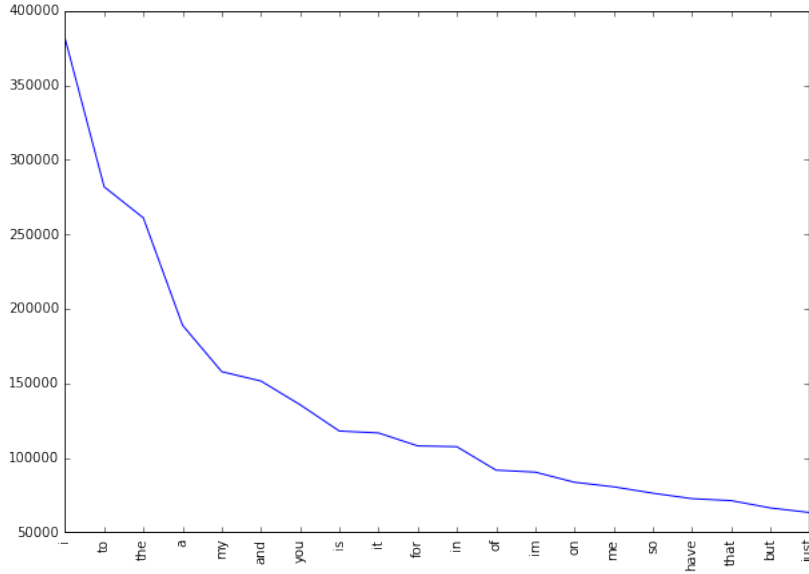


Figure 1: Frequencies of top 20 unigrams.

the dataset. Out of these words, most of the words at end of frequency spectrum are noise and occur very few times to influence classification. We, therefore, only use top N words from these to create our vocabulary where N is 15000 for sparse vector classification and 90000 for dense vector classification. The frequency distribution of top 20 words in our vocabulary is shown in figure 1. We can observe in figure 2 that the frequency distribution follows Zipf’s law which states that in a large sample of words, the frequency of a word is inversely proportional to its rank in the frequency table. This can be seen by the fact that a linear trendline with a negative slope fits the plot of $\log(\text{Frequency})$ vs. $\log(\text{Rank})$. The equation of the trendline shown in figure 2 is $\log(\text{Frequency}) = -0.78 \log(\text{Rank}) + 13.31$.

3.2.2 Bigrams

Bigrams are word pairs in the dataset which occur in succession in the corpus. These features are a good way to model negation in natural language like in the phrase – *This is not good*. A total of 1954953 unique bigrams were extracted from the dataset. Out of these, most of the bigrams at end of frequency spectrum are noise and occur very few times to influence classification. We therefore use only top 10000 bigrams from these to create our vocabulary. The frequency distribution of top 20 bigrams in our vocabulary is shown in figure 3.

3.3 Feature Representation

After extracting the unigrams and bigrams, we represent each tweet as a feature vector in either sparse vector representation or dense vector representation depending on the classification method.

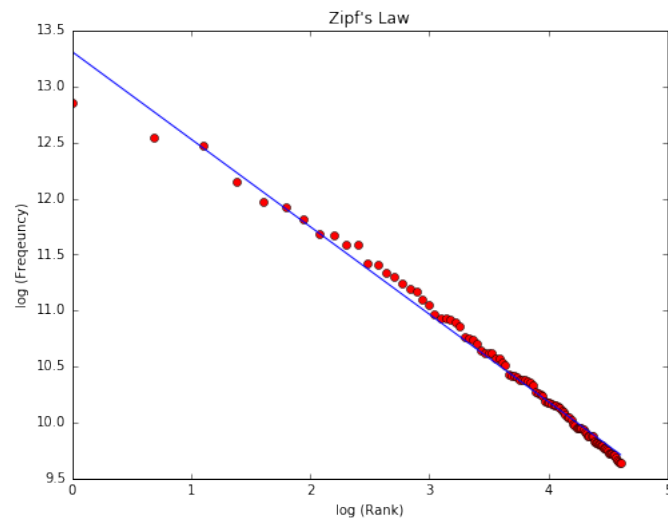


Figure 2: Unigrams frequencies follow Zipf's Law.

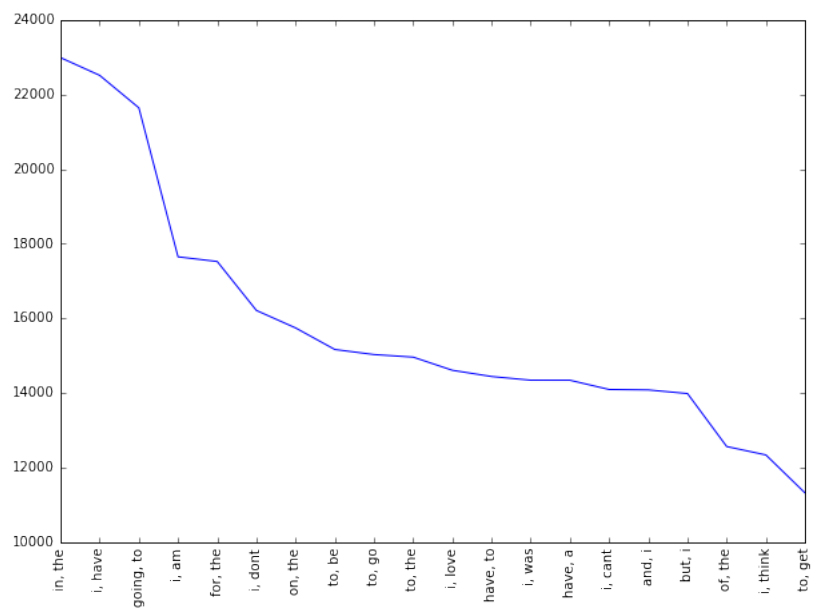


Figure 3: Frequencies of top 20 bigrams.

3.3.1 Sparse Vector Representation

Depending on whether or not we are using bigram features, the sparse vector representation of each tweet is either of length 15000 (when considering only unigrams) or 25000 (when considering unigrams and bigrams). Each unigram (and bigram) is given a unique index depending on its rank. The feature vector for a tweet has a positive value at the indices of unigrams (and bigrams) which are present in that tweet and zero elsewhere which is why the vector is sparse. The positive value at the indices of unigrams (and bigrams) depends on the feature type we specify which is one of *presence* and *frequency*.

- **presence** In the case of *presence* feature type, the feature vector has a 1 at indices of unigrams (and bigrams) present in a tweet and 0 elsewhere.
- **frequency** In the case of *frequency* feature type, the feature vector has a positive integer at indices of unigrams (and bigrams) which is the frequency of that unigram (or bigram) in the tweet and 0 elsewhere. A matrix of such term-frequency vectors is constructed for the entire training dataset and then each term frequency is scaled by the inverse-document-frequency of the term (idf) to assign higher values to important terms. The inverse-document-frequency of a term t is defined as.

$$idf(t) = \log \left(\frac{1 + n_d}{1 + df(d, t)} \right) + 1$$

where n_d is the total number of documents and $df(d, t)$ is the number of documents in which the term t occurs.

Handling Memory Issues Which dealing with sparse vector representations, the feature vector for each tweet is of length 25000 and the total number of tweets in the training set is 800000 which means allocation of memory for a matrix of size 800000×25000 . Assuming 4 bytes are required to represent each float value in the matrix, this matrix needs a memory of 8×10^{10} bytes (≈ 75 GB) which is far greater than the memory available in common notebooks. To tackle this issue, we used `scipy.sparse.lil_matrix` data structure provided by Scipy which is a memory efficient linked list based implementation of sparse matrices. In addition to that, we used Python generators wherever possible instead of keeping the entire dataset in memory.

3.3.2 Dense Vector Representation

For dense vector representation we use a vocabulary of unigrams of size 90000 i.e. the top 90000 words in the dataset. We assign an integer index to each word depending on its rank (starting from 1) which means that the most common word is assigned the number 1, the second most common word is assigned the number 2 and so on. Each tweet is then represented by a vector of these indices which is a dense vector.

3.4 Classifiers

3.4.1 Naive Bayes

Naive Bayes is a simple model which can be used for text classification. In this model, the class \hat{c} is assigned to a tweet t , where

$$\hat{c} = \underset{c}{\operatorname{argmax}} P(c|t)$$
$$P(c|t) \propto P(c) \prod_{i=1}^n P(f_i|c)$$

In the formula above, f_i represents the i -th feature of total n features. $P(c)$ and $P(f_i|c)$ can be obtained through maximum likelihood estimates.

3.4.2 Maximum Entropy

Maximum Entropy Classifier model is based on the Principle of Maximum Entropy. The main idea behind it is to choose the most uniform probabilistic model that maximizes the entropy, with given constraints. Unlike Naive Bayes, it does not assume that features are conditionally independent of each other. So, we can add features like bigrams without worrying about feature overlap. In a binary classification problem like the one we are addressing, it is the same as using Logistic Regression to find a distribution over the classes. The model is represented by

$$P_{ME}(c|d, \lambda) = \frac{\exp[\sum_i \lambda_i f_i(c, d)]}{\sum_{c'} \exp[\sum_i \lambda_i f_i(c', d)]}$$

Here, c is the class, d is the tweet and λ is the weight vector. The weight vector is found by numerical optimization of the lambdas so as to maximize the conditional probability.

3.4.3 Decision Tree

Decision trees are a classifier model in which each node of the tree represents a test on the attribute of the data set, and its children represent the outcomes. The leaf nodes represents the final classes of the data points. It is a supervised classifier model which uses data with known labels to form the decision tree and then the model is applied on the test data. For each node in the tree the best test condition or decision has to be taken. We use the GINI factor to decide the best split. For a given node t , $GINI(t) = 1 - \sum_j [p(j|t)]^2$, where $p(j|t)$ is the relative frequency of class j at node t , and $GINI_{split} = \sum_{i=1}^k \frac{n_i}{n} GINI(i)$ (n_i = number of records at child i , n = number of records at node p) indicates the quality of the split. We choose a split that minimizes the GINI factor.

3.4.4 Random Forest

Random Forest is an ensemble learning algorithm for classification and regression. Random Forest generates a multitude of decision trees classifies based on the aggregated decision of those trees. For a set of tweets x_1, x_2, \dots, x_n and their respective sentiment labels y_1, y_2, \dots, y_n bagging repeatedly selects a random sample (X_b, Y_b) with replacement. Each classification tree f_b is trained using a different random sample (X_b, Y_b) where b ranges from $1 \dots B$. Finally, a majority vote is taken of predictions of these B trees.

3.4.5 XGBoost

Xgboost is a form of gradient boosting algorithm which produces a prediction model that is an ensemble of weak prediction decision trees. We use the ensemble of K models by adding their outputs in the following manner

$$\hat{y}_i = \sum_{k=1}^K f_k(x_i), f_k \in F$$

where F is the space of trees, x_i is the input and \hat{y}_i is the final output. We attempt to minimize the following loss function

$$L(\Phi) = \sum_i l(\hat{y}_i, y_i) + \sum \Omega(f_k)$$

$$\text{where } \Omega(f) = \gamma T + \frac{1}{2} \lambda \|w\|^2$$

where Ω is the regularisation term.

3.4.6 SVM

SVM, also known as support vector machines, is a non-probabilistic binary linear classifier. For a training set of points (x_i, y_i) where x is the feature vector and y is the class, we want to find the

maximum-margin hyperplane that divides the points with $y_i = 1$ and $y_i = -1$. The equation of the hyperplane is as follow

$$w \cdot x - b = 0$$

We want to maximize the margin, denoted by γ , as follows

$$\max_{w, \gamma} \gamma, s.t. \forall i, \gamma \leq y_i(w \cdot x_i + b)$$

in order to separate the points well.

3.4.7 Multi-Layer Perceptron

MLP or Multilayer perceptron is a class of feed-forward neural networks, which has atleast three layers of neurons. Each neuron uses a non-linear activation function, and learns with supervision using backpropagation algorithm. It performs well in complex classification problems such as sentiment analysis by learning non-linear models.

3.4.8 Convolutional Neural Networks

Convolutional Neural Networks or CNNs are a type of neural networks which involve layers called convolution layers which can interpret spacial data. A convolution layers has a number of filters or kernels which it learns to extract specific types of features from the data. The kernel is a 2D window which is slid over the input data performing the convolution operation. We use temporal convolution in our experiments which is suitable for analyzing sequential data like tweets.

3.4.9 Recurrent Neural Networks

Recurrent Neural Network are a network of neuron-like nodes, each with a directed (one-way) connection to every other node. In RNN, hidden state denoted by h_t acts as memory of the network and learns contextual information which is important for classification of natural language. The output at each step is calculated based on the memory h_t at time t and current input x_t . The main feature of an RNN is its hidden state, which captures sequential dependence in information. We used Long Term Short Memory (LSTM) networks in our experiments which is a special kind of RNN capable of remembering information over a long period of time.

4 Experiments

We perform experiments using various different classifiers. Unless otherwise specified, we use 10% of the training dataset for validation of our models to check against overfitting i.e. we use 720000 tweets for training and 80000 tweets for validation. For Naive Bayes, Maximum Entropy, Decision Tree, Random Forest, XGBoost, SVM and Multi-Layer Perceptron we use sparse vector representation of tweets. For Recurrent Neural Networks and Convolutional Neural Networks we use the dense vector representation.

4.1 Baseline

For a baseline, we use a simple positive and negative word counting method to assign sentiment to a given tweet. We use the [Opinion Dataset](#) of positive and negative words to classify tweets. In cases when the number of positive and negative words are equal, we assign positive sentiment. Using this baseline model, we achieve a classification accuracy of 63.48% on Kaggle public leaderboard.

4.2 Naive Bayes

We used MultinomialNB from `sklearn.naive_bayes` package of *scikit-learn* for Naive Bayes classification. We used Laplace smoothed version of Naive Bayes with the smoothing parameter α set to its default value of 1. We used sparse vector representation for classification and ran experiments using both *presence* and *frequency* feature types. We found that *presence* features outperform *frequency* features because Naive Bayes is essentially built to work better on integer features rather

than floats. We also observed that addition of bigram features improves the accuracy. We obtain a best validation accuracy of 79.68% using Naive Bayes with *presence* of unigrams and bigrams. A comparison of accuracies obtained on the validation set using different features is shown in table 5.

4.3 Maximum Entropy

The `nltk` library provides several text analysis tools. We use the `MaxentClassifier` to perform sentiment analysis on the given tweets. Unigrams, bigrams and a combination of both were given as input features to the classifier. The Improved Iterative Scaling algorithm for training provided better results than Generalised Iterative Scaling. Feature combination of unigrams and bigrams, gave better accuracy of 80.98% compared to just unigrams (79.34%) and just bigrams (79.2%).

For a binary classification problem, Logistic Regression is essentially the same as Maximum Entropy. So, we implemented a sequential Logistic Regression model using `keras`, with sigmoid activation function, binary cross-entropy loss and Adam's optimizer achieving better performance than `nltk`. Using frequency and presence features we get almost the same accuracies, but the performance is slightly better when we use unigrams and bigrams together. The best accuracy achieved was 81.52%. A comparison of accuracies obtained on the validation set using different features is shown in table 5.

4.4 Decision Tree

We use the `DecisionTreeClassifier` from `sklearn.tree` package provided by *scikit-learn* to build our model. GINI is used to evaluate the split at every node and the best split is chosen always. The model performed slightly better using the presence feature compared to frequency. Also using unigrams with or without bigrams didn't make any significant improvements. The best accuracy achieved using decision trees was 68.1%. A comparison of accuracies obtained on the validation set using different features is shown in table 5.

4.5 Random Forest

We implemented random forest algorithm by using `RandomForestClassifier` from `sklearn.ensemble` provided by *scikit-learn*. We experimented using 10 estimators (trees) using both *presence* and *frequency* features. *presence* features performed better than *frequency* though the improvement was not substantial. A comparison of accuracies obtained on the validation set using different features is shown in table 5.

4.6 XGBoost

We also attempted tackling the problem with XGboost classifier. We set max tree depth to 25 where it refers to the maximum depth of a tree and is used to control over-fitting as a high value might result in the model learning relations that are tied to the training data. Since XGboost is an algorithm that utilises an ensemble of weaker trees, it is important to tune the number of estimators that is used. We realised that setting this value to 400 gave the best result. The best result was 0.78.72 which came from the configuration of presence with Unigrams + Bigrams.

4.7 SVM

We utilise the SVM classifier available in `sklearn`. We set the C term to be 0.1. C term is the penalty parameter of the error term. In other words, this influences the misclassification on the objective function. We run SVM with both Unigram as well Unigram + Bigram. We also run the configurations with frequency and presence. The best result was 81.55 which came the configuration of frequency and Unigram + Bigram.

4.8 Multi-Layer Perceptron

We used `keras` with `TensorFlow` backend to implement the Multi-Layer Perceptron model. We used a 1-hidden layer neural network with 500 hidden units. The output from the neural network

Algorithms	Presence		Frequency	
	Unigrams	Unigrams+Bigrams	Unigrams	Unigrams+Bigrams
Naive Bayes	78.16	79.68	77.52	79.38
Max Entropy	79.96	81.52	79.7	81.5
Decision Tree	68.1	68.01	67.82	67.78
Random Forest	76.54	77.21	76.16	77.14
XGBoost	77.56	78.72	77.42	78.32
SVM	79.54	81.11	79.83	81.55
MLP	80.1	81.7	80.15	81.35

Table 5: Comparison of various classifiers which use sparse vector representation

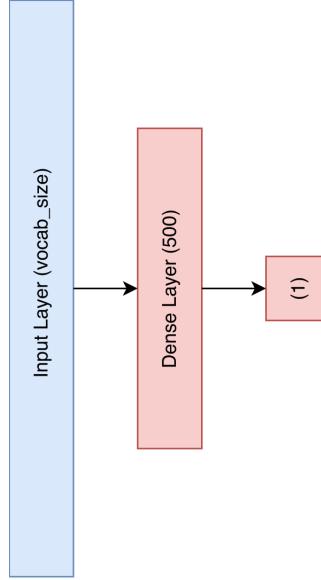


Figure 4: Architecture of the MLP Model.

is a single value which we pass through the sigmoid non-linearity to squish it in the range $[0, 1]$. The sigmoid function is defined by $f(z) = \frac{1}{1+\exp^{-z}}$. The output from the neural network gives the probability $\Pr(\text{positive}|\text{tweet})$ i.e. the probability of the tweets sentiment being positive. At the prediction step, we round off the probability values to convert them to class labels 0 (negative) and 1 (positive). The architecture of the model is shown in figure . Red hidden layers represent layers with sigmoid non-linearity. We trained our model using binary cross entropy loss with the weight update scheme being the one defined by Adam et. al. We also conducted experiments using SGD + Momentum weight updates and found out that it takes too long to converge. We ran our model upto 20 epochs after which it began to overfit. We used sparse vector representation of tweets for training. We found that the presence of bigrams features significantly improved the accuracy.

4.9 Convolutional Neural Networks

We used `keras` with `TensorFlow` backend to implement the Convolutional Neural Network model. We used the dense vector representation of the tweets to train our CNN models. We used a vocabulary of top 90000 words from the training dataset. We represent each word in our vocabulary with an integer index from $1 \dots 90000$ where the integer index represents the rank of the word in the dataset. The integer index 0 is reserved for the special padding word. Further each of these $90000+1$ words is represented by a 200 dimensional vector. The first layer of our models is the **Embedding** layer which is a matrix of shape $(v+1) \times d$ where v is vocabulary size ($=90000$) and d is the dimension of each word vector ($=200$). We initialize the embedding layer with random weights from $\mathcal{N}(0, 0.01)$. Each row of this embedding matrix represents represents the 200 dimensional word vector for a word in the vocabulary. For words in our vocabulary which match GloVe word vectors provided by the [StanfordNLP](#) group, we seed the corresponding row of the embedding matrix from GloVe vectors. Each tweet i.e. its dense vector representation is padded with 0s at

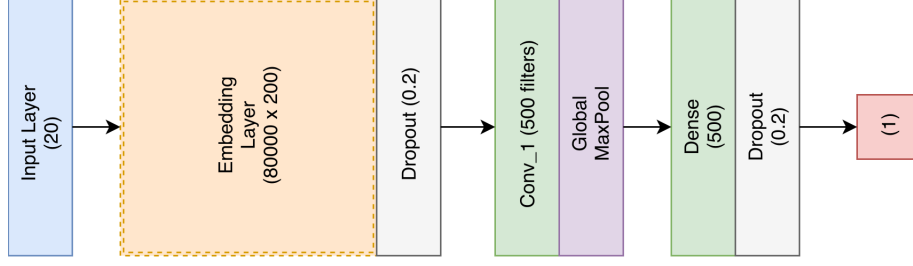


Figure 5: Neural Network Architecture with 1 Conv Layer.

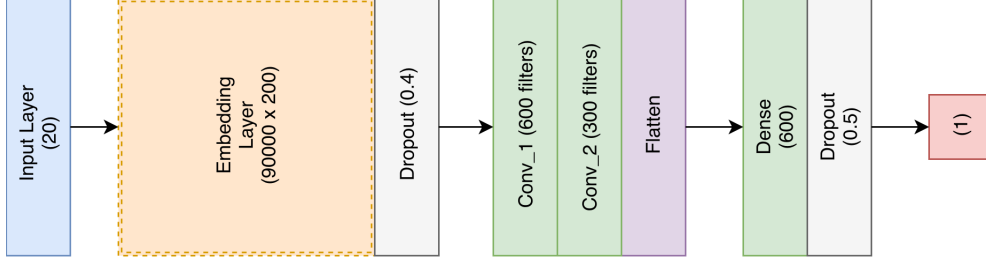


Figure 6: Neural Network Architecture with 2 Conv Layers.

the end until its length is equal to `max_length` which is a parameter we tweak in our experiments. We trained our model using binary cross entropy loss with the weight update scheme being the one defined by Adam et. al. We also conducted experiments using SGD + Momentum weight updates and found out that it takes longer (≈ 100 epochs) to converge compared to validation accuracy equivalent to Adam. We ran our model upto 10 epochs. Using the Adam weight update scheme, the model converges very fast (≈ 4 epochs) and begins to overfit badly after that. We, therefore, use models from 3rd or 4th epoch for our results. We tried four different CNN architectures which are as follows.

- **1-Conv-NN:** As the name suggests, this is an architecture with 1 convolution layer. We perform temporal convolution with a kernel size of 3 and zero padding. After the convolution layer, we apply *relu* activation function (which is defined as $f(x) = \max(0, x)$) and then perform Global Max Pooling over time to reduce the dimensionality of the data. We pass the output of the Global Max Pool layer to a fully-connected layer which then outputs a single value which is passed through sigmoid activation function to convert it into a probability value. We also added dropout layers after the embedding layer and the fully-connected layer to regularize our network and prevent it from overfitting. We use a tweet `max_length` of 20 in this network with a vocabulary of 80000 words. The complete architecture of the network is `embedding_layer (800001x200) → dropout(0.2) → conv_1 (500 filters) → relu → global_maxpool → dense(500) → relu → dropout(0.2) → dense(1) → sigmoid` as shown in figure 5. Green layers indicate *relu* activation while red indicates *sigmoid*.
- **2-Conv-NN:** In this architecture we increased the vocabulary from 80000 to 90000. We also increased the dropout after embedding layer to 0.4 and that after the fully connected layer to 0.5 to further regularize the network and thus prevent overfitting. We changed the number of filters in the first convolution layer to 600 and added another convolution layer with 300 filters after the first convolution layer. We also replaced the Global MaxPool layer with a Flatten layer as we believed some features of the input tweets got lost while max pooling. We also increased the number of units in the fully-connected layer to 600. All of these changes allowed the network to learn and regularize better thereby improving the validation accuracy. The complete architecture of the network is `embedding_layer (900001x200) → dropout(0.4) → conv_1 (600 filters) → relu → conv_2 (300 filters) → relu → flatten → dense(600) → relu → dropout(0.5) → dense(1) → sigmoid` as shown in figure 6.
- **3-Conv-NN:** In this architecture we added another convolution layer with 150 filters after the second convolution layer. The complete architecture of the network is `embedding_layer`

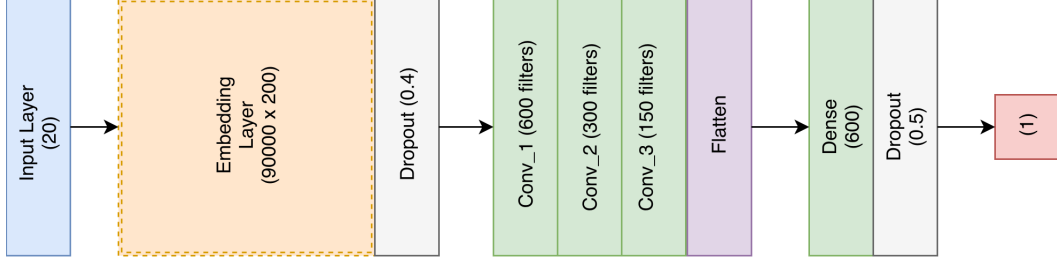


Figure 7: Neural Network Architecture with 3 Conv Layers.

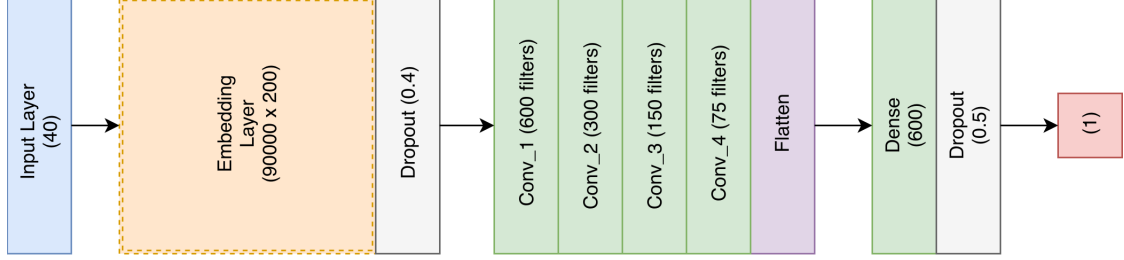


Figure 8: Neural Network Architecture with 4 Conv Layers.

$(900001 \times 200) \rightarrow \text{dropout}(0.4) \rightarrow \text{conv}_1$ (600 filters) $\rightarrow \text{relu} \rightarrow \text{conv}_2$ (300 filters) $\rightarrow \text{relu} \rightarrow \text{conv}_3$ (150 filters) $\rightarrow \text{relu} \rightarrow \text{flatten} \rightarrow \text{dense}(600) \rightarrow \text{relu} \rightarrow \text{dropout}(0.5) \rightarrow \text{dense}(1) \rightarrow \text{sigmoid}$ as shown in figure 7.

- **4-Conv-NN:** In this architecture we added another convolution layer with 75 filters after the third convolution layer. We also increased `max_length` of the tweet to 40 going by the fact that the length of largest tweet in our pre-processed dataset is about 40 words. The complete architecture of the network is `embedding_layer` $(900001 \times 200) \rightarrow \text{dropout}(0.4) \rightarrow \text{conv}_1$ (600 filters) $\rightarrow \text{relu} \rightarrow \text{conv}_2$ (300 filters) $\rightarrow \text{relu} \rightarrow \text{conv}_3$ (150 filters) $\rightarrow \text{relu} \rightarrow \text{conv}_4$ (75 filters) $\rightarrow \text{relu} \rightarrow \text{flatten} \rightarrow \text{dense}(600) \rightarrow \text{relu} \rightarrow \text{dropout}(0.5) \rightarrow \text{dense}(1) \rightarrow \text{sigmoid}$ as shown in figure 8.

We notice that each successive CNN model is better than the previous one with 1-Conv-NN, 2-Conv-NN, 3-Conv-NN and 4-Conv-NN achieving accuracies of 82.40, 82.76, 82.95 and 83.34 respectively on Kaggle public leaderboard.

4.10 Recurrent Neural Networks

We used neural networks with LSTM layers in our experiments. We used a vocabulary of top 20000 words from the training dataset. We used the dense vector representation for training our models. We pad or truncate each dense vector representation to make it equal to `max_length` which is a parameter we tweak in our experiments. The first layer of our network is the Embedding layer which as described in section 4.9 We test two different types of LSTM models.

- *Random Embedding Initialization:* In these models, we use a word embedding dimension of 32 and train the embeddings from scratch. The embedding layer is followed by an LSTM layer where we experimented with different number of LSTM units. The LSTM layer is followed by a fully-connected layer with 32 units and *relu* activation. Finally, the output is a single value with *sigmoid* activation. We also add dropouts of 0.2 after embeddings layer and the penultimate layer to regularize the network.
- *Embeddings Seeded with GloVe:* In these models, we use a word vector dimension of 200 instead and seed it with GloVe word vectors provided by the StanfordNLP group. The word embeddings are fine tuned during the course of training. We follow the embeddings layer with an LSTM layer which is followed by a fully-connected layer with *relu* activation. Finally, the output is a single value with *sigmoid* activation. We add dropouts of 0.4 and 0.5 after embeddings layer and the penultimate layer respectively to further regularize the network.

LSTM Units	Dense Units	max_length	Loss	Embedding Initialization	Accuracy
100	32	40	MSE	Random	79.8%
100	32	40	BCE	Random	82.2%
50	32	40	MSE	Random	78.96%
50	32	40	BCE	Random	81.97%
100	600	20	BCE	GloVe	82.7%
128	64	40	BCE	GloVe	83.0%

Table 6: Comparison of different LSTM models. MSE is mean squared error and BCE is binary cross entropy.

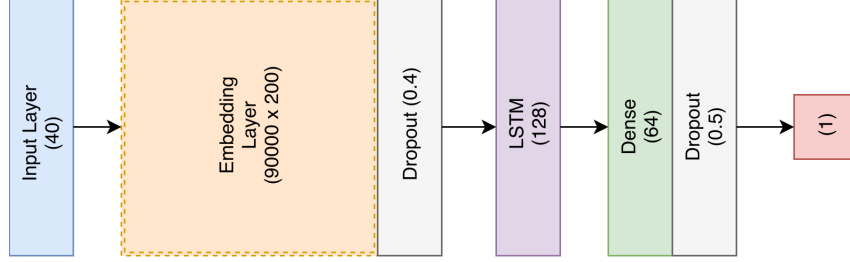


Figure 9: Architecture of best performing LSTM-NN

We experimented with different values of LSTM and fully-connected units and the results are summarized in table 6. The architecture of our best performing LSTM-NN is shown in figure 9.

We experimented with both Adam optimizer and SGD with momentum for training our networks and find the Adam worked better and converges faster. We trained our model using `mean_squared_error` and `binary_cross_entropy` loss. We found that `binary_cross_entropy` worked better than `mean_squared_error` which is expected given our binary classification problem. The results from various different LSTM models are summarized in table 6. We obtain best accuracy of 83.0% among the different LSTM models.

4.11 Ensemble

In a quest to further improve accuracy, we developed a simple ensemble model. We first extract 600 dimensional feature vectors for each tweet from the penultimate layer of our best performing 4-Conv-NN model. Each tweet is now represented by a 600 dimensional feature vector. We use these features to classify the tweets using a linear SVM model with $C=1$. We classify the tweets using this SVM model. We then take the majority vote of predictions from the following 5 models.

1. LSTM-NN
2. 4-Conv-NN
3. 4-Conv-NN features + SVM
4. 4-Conv-NN with `max_length` = 20
5. 3-Conv-NN

The accuracies from each of these individual models and their majority voting ensemble are shown in table 7. The flowchart of ensemble is shown in figure 10.

5 Conclusion

5.1 Summary of achievements

The provided tweets were a mixture of words, emoticons, URLs, hastags, user mentions, and symbols. Before training the we pre-process the tweets to make it suitable for feeding into models. We

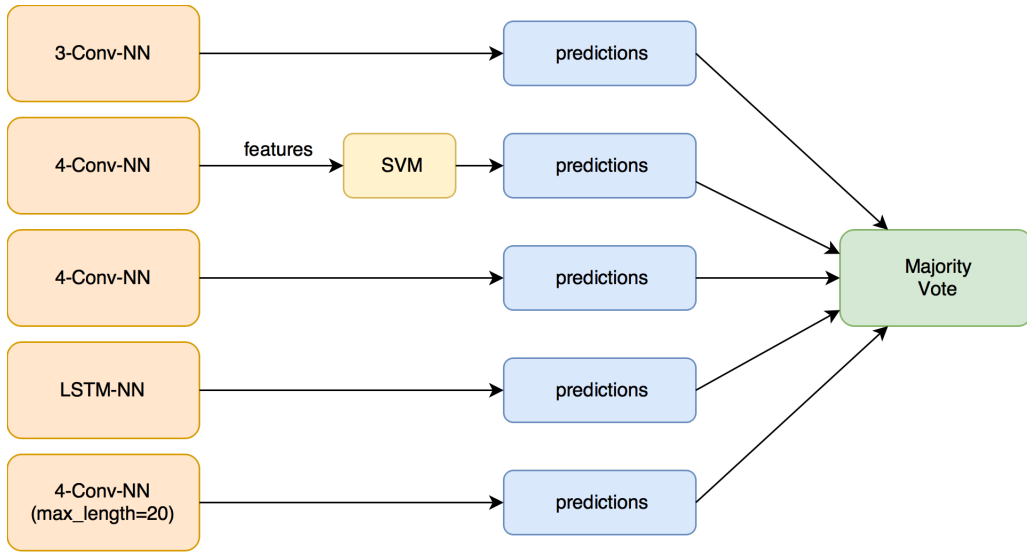


Figure 10: Flowchart of Majority Voting Ensemble.

Model	Accuracy
LSTM-NN	83.00
4-Conv-NN	83.34
4-Conv-NN features + SVM	83.39
4-Conv-NN with max_length = 20	82.85
3-Conv-NN	82.95
Majority Vote Ensemble	83.58

Table 7: Models used for ensemble and their accuracies on Kaggle public leaderboard

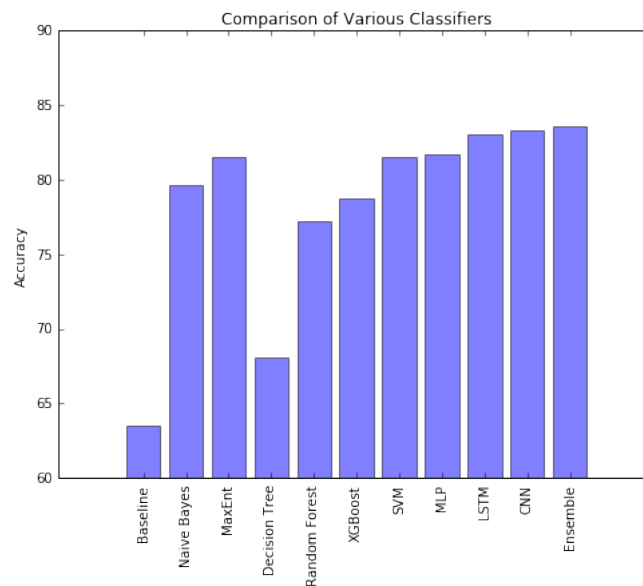


Figure 11: Comparison of accuracies of various models

implemented several machine learning algorithms like Naive Bayes, Maximum Entropy, Decision Tree, Random Forest, XGBoost, SVM, Multi-Layer Perceptron, Recurrent Neural networks and Convolutional Neural Networks to classify the polarity of the tweet. We used two types of features namely unigrams and bigrams for classification and observes that augmenting the feature vector with bigrams improved the accuracy. Once the feature has been extracted it was represented as either a sparse vector or a dense vector. It has been observed that *presence* in the sparse vector representation recorded a better performance than *frequency*.

Neural methods performed better than other classifiers in general. Our best LSTM model achieved an accuracy of 83.0% on Kaggle while the best CNN model achieved 83.34%. The model which used features from our best CNN model and classifies using SVM performed slightly better than only CNN. We finally used an ensemble method taking a majority vote over the predictions of 5 of our best models achieving an accuracy of 83.58%.

5.2 Future directions

- *Handling emotion ranges*: We can improve and train our models to handle a range of sentiments. Tweets don't always have positive or negative sentiment. At times they may have no sentiment i.e. neutral. Sentiment can also have gradations like the sentence, *This is good*, is positive but the sentence, *This is extraordinary*, is somewhat more positive than the first. We can therefore classify the sentiment in ranges, say from -2 to +2.
- *Using symbols*: During our pre-processing, we discard most of the symbols like commas, full-stops, and exclamation mark. These symbols may be helpful in assigning sentiment to a sentence.