

# DUPLICATE QUESTION IDENTIFICATION

Capstone Project 2: Final Report

Siri Surabathula

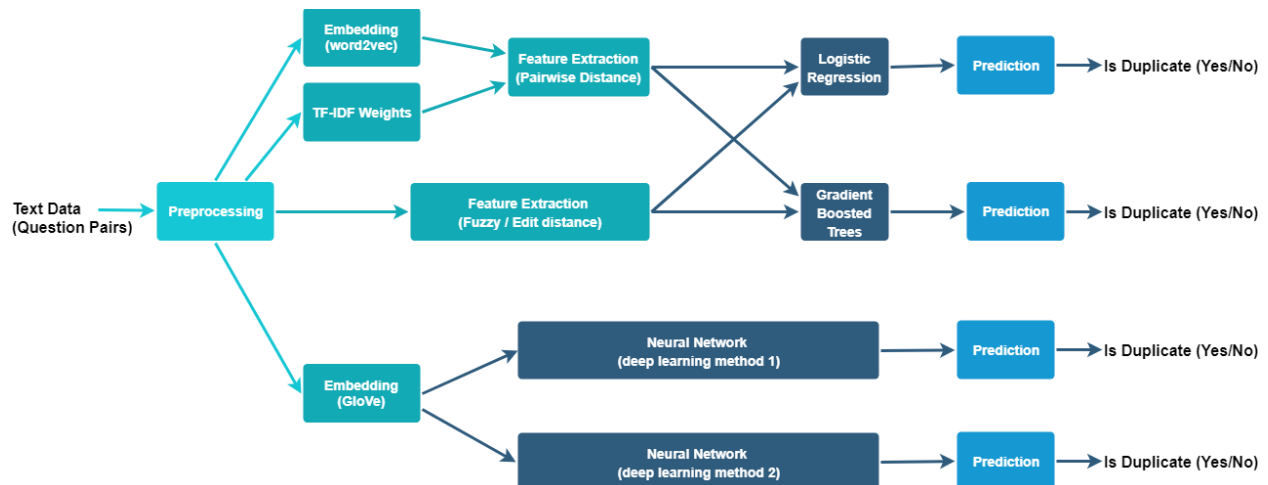
## Problem Statement

Question-answering and knowledge-sharing platforms like Quora and Stack-Exchange require mechanisms to group questions in their database based on similar intent. According to Quora about 100 million users visit their website every month, and a lot of the questions end up being duplicates (duplicate questions could be differently worded but they have the same intent)

For example, the questions 'What is the best way to travel from Houston to Atlanta ?' and 'Should I drive, fly or take a bus from Houston to Atlanta ?', are duplicates.

User experience is greatly improved if Quora is able to identify duplicate questions as this enables them to find questions that have already been answered and also avoids the need to answer the same question multiple times.

This project explores different classification methods to solve the problem of duplicate identification. These methods include logistic regression, gradient boosted trees and neural networks. The text data is preprocessed and embedded in vector space before it is fed to the classification models. The general approach of text preprocessing, embedding, modeling and classification is outlined in the high-level diagram below,



# Data

Source : [Kaggle \(Quora Question Pairs\)](#)

Description : The dataset contains a human-labeled training set and a test set. Each record in the training set represents a pair of questions and a binary label which indicates whether the question-pair is a duplicate or not.

## Pre-processing

Text preprocessing was performed on the text data using [gensim API](#)

The code for this is described in the ipython notebook [tfidf\\_word2vec](#)

The steps involved in preprocessing are as follows :

1. Tag Removal - Removed tags. For eg. "<i>Hello</i> <b>World</b>!" was converted to "Hello World!"
2. Repeating Whitespace Removal - Removed repeating whitespace characters (spaces, tabs, line breaks) and converted tabs & line breaks into spaces
3. Stopwords Removal - Removed stopwords. (the default stopwords list from gensim was used)
4. Case conversion and Stemming - Transformed text to lowercase and performed porter stemming.

## Word Embedding

### Word2vec

[Mikolov et al.](#) proposed this model based on two architectures, Continuous Bag of Words (CBOW) and Skip-Gram. Both methods take into account the context of the words while computing embeddings.

CBOW - This model takes into account a dynamic window of  $n$  words around (both before and after) the target word  $w_t$ . The order of the words is of no consequence.

The objective function is given by (here  $T$  is the total number of words in the training corpus)

$$J_{\theta} = \frac{1}{T} \sum_{t=1}^T \log p(w_t | w_{t-n}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+n})$$

Skip-gram - This model does the opposite of CBOW, i.e. it computes the probabilities of  $n$  words around the target word  $w_t$ , given the probability of  $w_t$ .

The objective function is given by (here  $T$  is the total number of words in the training corpus)

$$J_{\theta} = \frac{1}{T} \sum_{t=1}^T \sum_{-n < j < n, j \neq 0} \log p(w_{t+j} | w_t)$$

## GloVe

[Pennington et al.](#) proposed this model based on the premise that the ratio of co-occurrence probabilities of two words (and not the co-occurrence probabilities themselves) carries significant information about the context-specific probability distribution of the words. They came up with a weighted least squares objective function as follows,

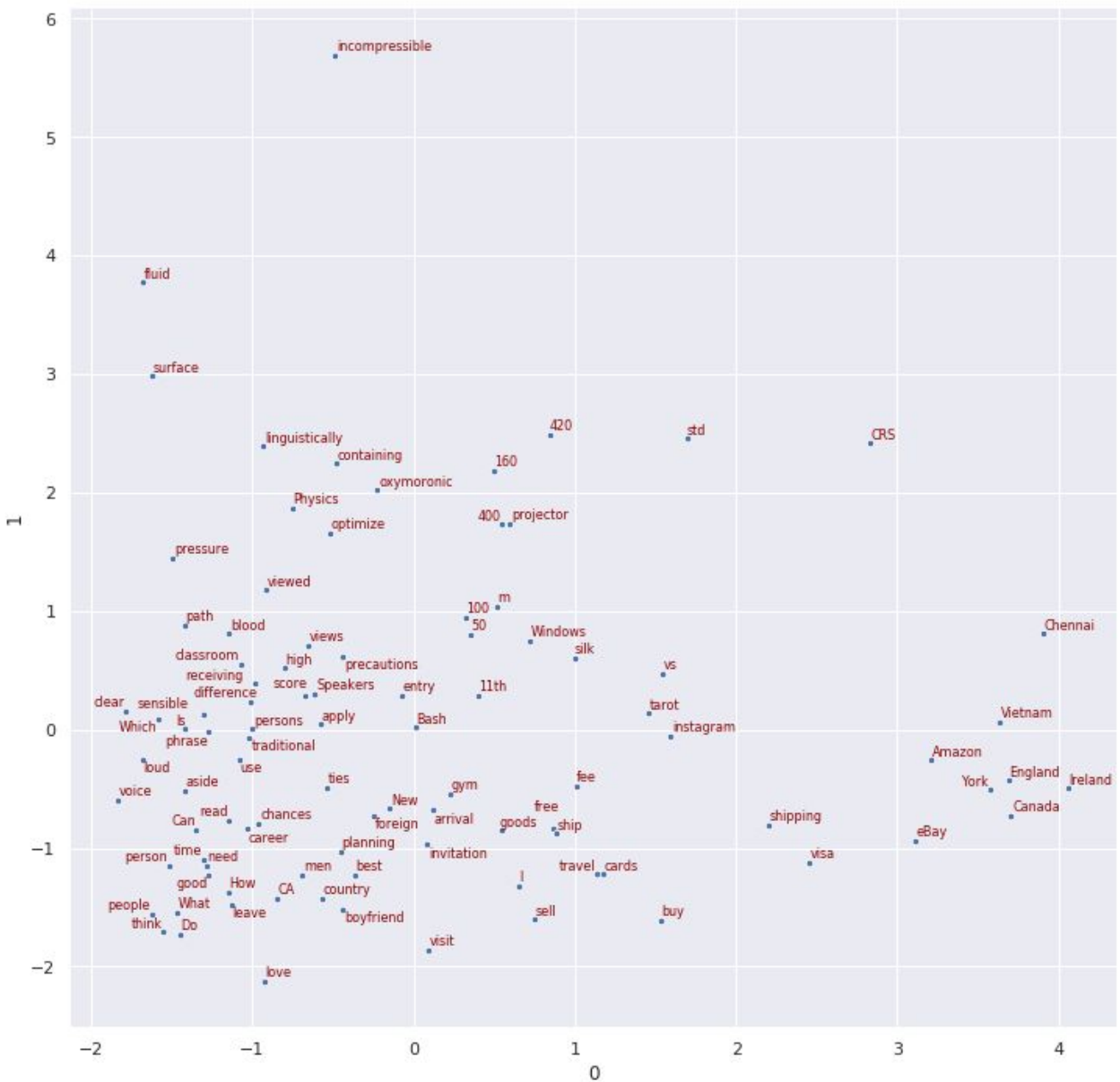
$$J_{\theta} = \sum_{i,j=1}^V f(X_{ij}) (w_i^T \bar{w}_j + b_i + \bar{b}_j - \log X_{ij})^2$$

Where  $w_i$  and  $b_i$  are the word vector and bias of the  $i$ th word,  $\bar{w}_j$  and  $\bar{b}_j$  are the word vector and bias of the  $j$ th context word,  $X_{ij}$  is the number of co-occurrences of the  $i$ th and  $j$ th words and  $f$  is the weighting function that assigns lower importance to co-occurrences involving rare and frequent words.

The text data was embedded in 300 dimensional space using two methods:

1. Word2Vec - For the classification methods involving logistic regression and gradient boosted trees, the [Fasttext pretrained model](#) for English (trained on Common Crawl and Wikipedia) was used for embedding. The pretrained fasttext model was transfer-trained with the questions in training data using the gensim method [intersect\\_word2vec\\_format](#). This method initializes a word2vec model with a vocabulary of the training data, then intersects this vocabulary with the pretrained model. No words are added to the existing vocabulary, but intersecting words adopt the weights of the pretrained model, while non-intersecting words are left alone.
2. GloVe - For the classification methods involving neural networks, the [Spacy pretrained model](#) for English (trained on Common Crawl using GloVe) was used for embedding. Provision was made for Out Of Vocabulary (OOV) words by randomly mapping each OOV word to one of 50 randomly generated vectors.

The plot below shows GloVe (Spacy) word-vectors for 20 questions picked from the Quora dataset. Dimension reduction (PCA) was used to reduce from 300 to 2 dimensions.



# Feature Extraction

## Feature Set 1 - Similarity using word-to-word (pairwise) Distance

### Pairwise Distance Metrics

The following pairwise distance and similarity metrics were computed. Each metric was computed for every combination of token-pairs formed by picking the first token ( $t_1$ ) from question 1 and the second token ( $t_2$ ) from question 2:

1. Chebyshev Distance - which is defined as,

$$\max|t_1 - t_2|$$

2. Bray-Curtis Distance - which is defined as,

$$\sum |t_1 - t_2| / \sum |t_1 + t_2|$$

3. Cosine Distance - which is defined as,

$$1 - \frac{t_1 \cdot t_2}{\|t_1\|_2 \|t_2\|_2}$$

4. Correlation Distance - which is defined as,

$$1 - \frac{(t_1 - \bar{t}_1)(t_2 - \bar{t}_2)}{\|t_1 - \bar{t}_1\|_2 \|t_2 - \bar{t}_2\|_2}$$

5. Canberra Distance - which is defined as,

$$\sum \frac{|t_1 - t_2|}{|t_1| + |t_2|}$$

6. Cityblock Distance - which is defined as,

$$\sum |t_1 - t_2|$$

7. Euclidean Distance - which is defined as,

$$\|t_1 - t_2\|_2$$

8. L1 Distance - which is defined as,

$$\|t_1 - t_2\|_1$$

9. Minkowski Distance - which is defined as,

$$\sqrt[p]{\sum |t_1 - t_2|^p}$$

10. Squared Euclidean Distance - which is defined as,

$$\|t_1 - t_2\|_2^2$$

## TF-IDF

TF-IDF or Term Frequency - Inverse Document Frequency, is a weighting statistic used in many NLP applications. The value of this statistic increases proportionally with the number of times a word appears in the document and is offset by the number of documents in the corpus that contain the word.

Term frequency  $tf_{(t,d)}$  of a term  $t$  occurring in a document  $d$  is given by,

$$tf_{(t,d)} = \frac{f_{t,d}}{\sum_{t' \in d} f_{t',d}} \quad \text{where } f_{t,d} \text{ is the raw count of term } t \text{ in document } d$$

Inverse document frequency  $idf(t, D)$  of a term  $t$  occurring in a document corpus  $D$  is given by,

$$idf(t, D) = \log \frac{N}{|\{d \in D : t \in d\}|}$$

where  $N$  is the total number of documents in the corpus  $D$  and  $|\{d \in D : t \in d\}|$  is the number of times the term  $t$  occurs in the corpus  $D$

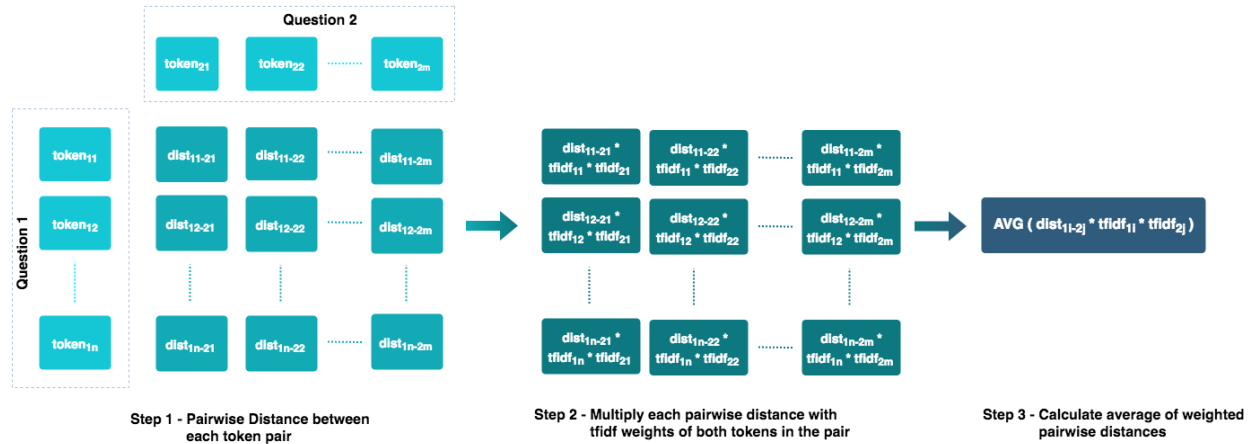
Then the TF-IDF of term  $t$  in the corpus  $D$  is given by,

$$tfidf(t, D) = tf_{(t,d)} \cdot idf(t, D)$$

## Avg. Weighted Distances between the Questions using TF-IDF Weights

The term-frequency-inverse-document-frequency (TF-IDF) weight was computed for each token in the entire dataset using the scikit-learn API ([TFIDFVectorizer](#)). The Vectorizer was trained using the training set only. These weights were then used to compute weighted averages of the pairwise distances.

The process is illustrated in the diagram below,



## Feature Set 2 - Similarity using Levenshtein (Edit) Distance

### Levenshtein Distance

This is a method for measuring the difference between two text sequences based on the number of single character edits (insertions, deletions and substitutions) it takes to change one sequence to another. It is also known as 'edit distance'.

It is computed for two sequences  $a$  and  $b$  (of length  $|a|$  and  $|b|$  respectively) as follows,

$$lev_{a,b}(i, j) = \begin{cases} \max(i, j) & , \text{if } \min(i, j) = 0 \\ \min(lev_{a,b}(i-1, j) + 1, lev_{a,b}(i, j-1) + 1, lev_{a,b}(i-1, j-1) + 1_{a_i \neq b_j}) & , \text{otherwise} \end{cases}$$

### Fuzzy-Wuzzy Distance Metrics

The python library [fuzzy-wuzzy](https://pypi.org/project/fuzzy-wuzzy/) was used to compute the following metrics:

1. Simple Ratio - which computes the similarity between two strings (in this case the two questions) using the simple edit distance (using Levenshtein distance) between the two strings
2. Partial Ratio - which improves on the simple ratio method above using a heuristic called "best partial" which is useful when the two strings are of noticeably different lengths. If the shorter string is length  $m$ , and the longer string is length  $n$ , the score of the best matching length- $m$  substring is taken into account.
3. Token Sort Ratio - which involves tokenizing each string in question, sorting the tokens alphabetically, and then joining them back. These new strings are then compared using the simple ratio method.
4. Token Set Ratio - which involves tokenizing both the strings in question, but instead of immediately sorting and comparing, the tokens in the strings are split into three groups: the intersection component common to both strings and the two remainder components from each

string. These sets are then used to perform the comparison. The scores increase when the intersection component makes up a larger percentage of the full string, and when the string remainders are more similar.

## Feature Set 3 - Word Vectors

The word-vectors produced by the embedding process described earlier form this feature set. These features contain substantial information about the question pairs but are not as interpretable as features in sets 1 & 2.

## Parallel Processing for Feature Computation

[Dask Delayed](#) was used to compute the features on the train and test datasets as all these computations involved embarrassingly parallel operations.

Dask Delayed is useful when there clearly exists scope for parallelism in a problem (for example one or more functions can be evaluated independently), but it's not possible to convert the data structures involved into a big array or big DataFrame for computation.

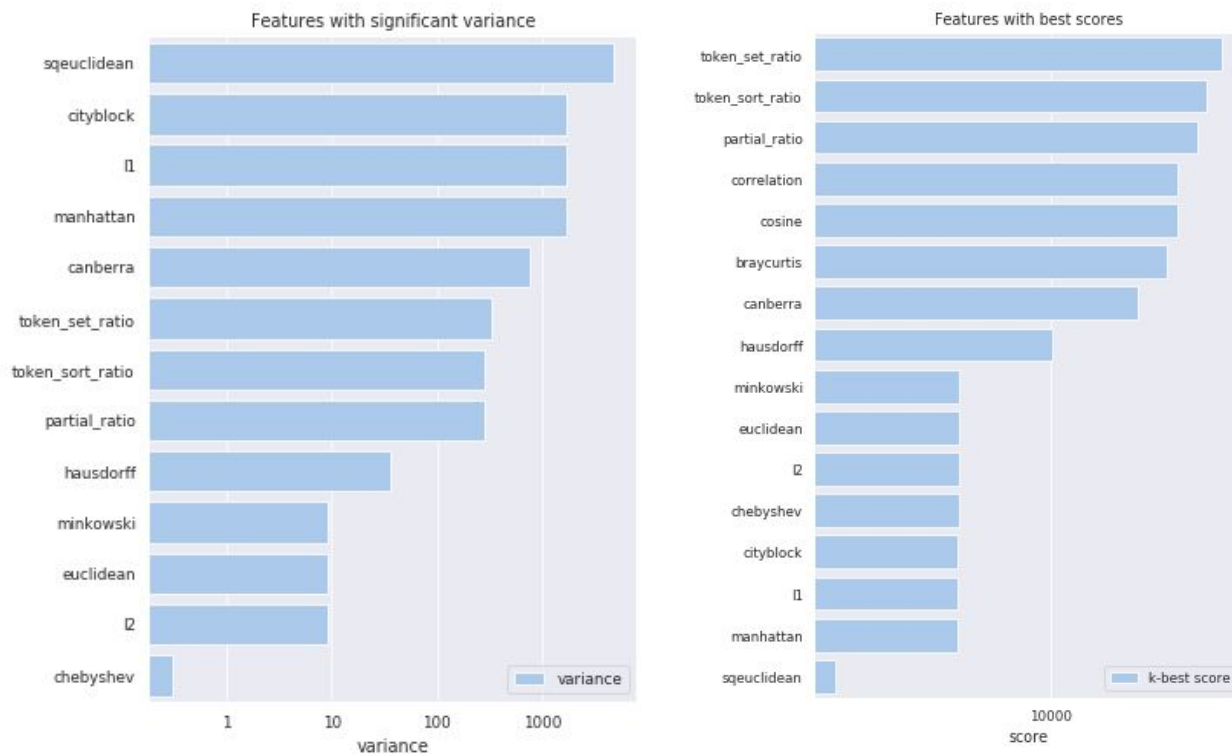
In this case, the data structures involved (the set of arrays corresponding to question 1 as well as the set of arrays corresponding to question 2) were each numpy arrays where each row again contained a numpy array of varying size. These could not be converted to numpy arrays with fixed number of dimensions or to DataFrames with fixed number of columns. (this was ultimately because of the fact that the number of tokens varied from question to question)

Dask Delayed allows functions to be operated lazily. Rather than executing a function immediately, it defers execution, placing the function and its arguments into a task graph. Once the task graph of the entire computation is constructed, Dask schedulers execute the tasks in way that exploits all possibility of parallelism, which in turn leads to improvement in performance.

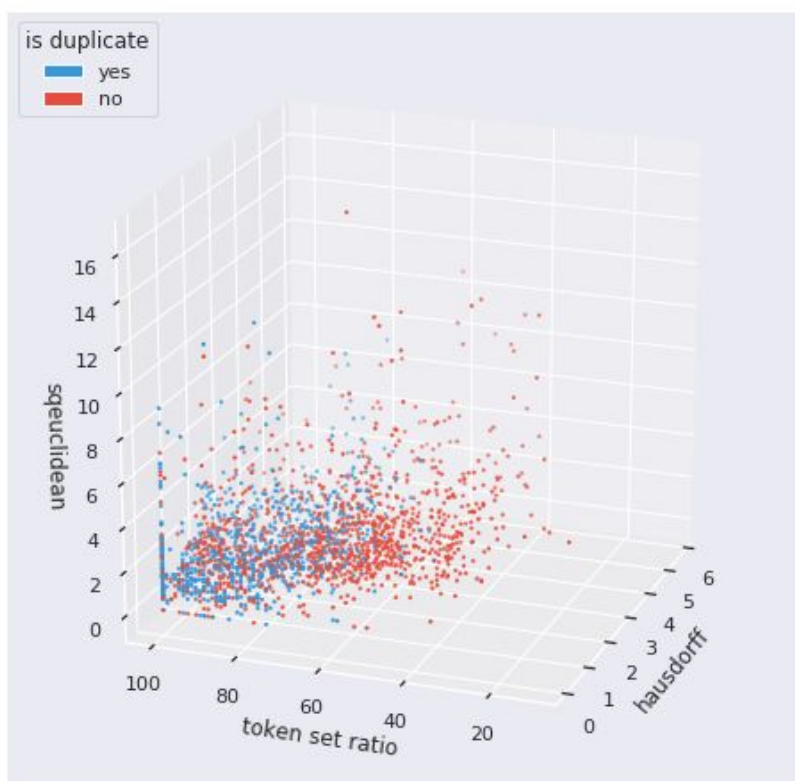


## Feature Importance

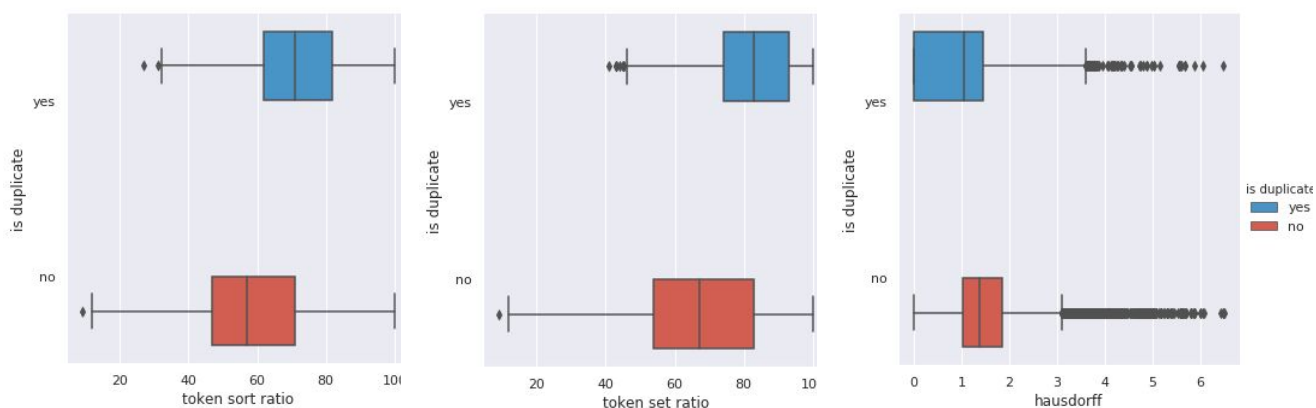
The barplots below show a comparison of features based on variance and best score (score is based on ANOVA F-value). The pairwise distance features like square euclidean, cityblock, l1 and manhattan have the highest variance over the training set, while the fuzzy (levenshtein distance based) features and some pairwise distance features like correlation and cosine have the highest F-value based scores.



The 3-D scatter plot below shows how duplicate(blue) and non-duplicate(red) question-pairs vary with square euclidean distance, token sort ratio and hausdorff distance.



The Box-and-whisker plots below show the 3 quantiles and extreme values for the levenshtein features token-sort-ratio and token-set-ratio and pairwise distance feature hausdorff.



The plots above give a qualitative idea about how hand-engineered features (sets 1 & 2) might contribute to the classification. It is obvious from these plots that a few hand-engineered features are not sufficient, and it is necessary to include feature-set 3 (embedding vectors)

# Modeling

## Train-Validation-Test Split

The data was split into training and hold-out sets using [sklearn.model\\_selection.train\\_test\\_split](#) with a train-test split of 67-33.

The non deep learning methods used a 5-fold cross validation on the training set. The deep learning methods split the training set using a training-validation split of 80-20 during hyperparameter optimization and 90-10 during training of the tuned model.

## Model 1 - Logistic Regression

The code for this is described in the notebook [tfidf\\_word2vec](#) (section Modeling)

The scikit-learn API ([sklearn.linear\\_model.LogisticRegression](#)) was used.

Logistic regression is a linear model for classification. In this model, the probabilities describing the possible outcomes of a single trial are modeled using a logistic function.

As an optimization problem, binary class L2 penalized logistic regression minimizes the following cost function,

$$\min_{\omega, c} \frac{1}{2} \omega^T \omega + C \sum_{i=1}^n \log \left( \exp \left( -y_i \left( X_i^T \omega + c \right) \right) + 1 \right)$$

The hyperparameters of the model are:

1. C (Inverse of regularization strength)
2. tolerance (Tolerance for stopping criteria)

## Model Selection

5-fold cross-validation was used with scikit-learn API ([sklearn.model\\_selection.RandomizedSearchCV](#)) to find the optimum combination of hyperparameters.

The optimum parameters are as follows,

C (Inverse of regularization strength)	100
tolerance (Tolerance for stopping criteria)	1e-05

## Prediction

The following is the classification report,

	<b>precision</b>	<b>recall</b>	<b>f1-score</b>	<b>support</b>
<b>not duplicate</b>	0.79	0.83	0.81	84248
<b>duplicate</b>	0.68	0.61	0.65	49146
<b>micro avg</b>	0.75	0.75	0.75	133394
<b>macro avg</b>	0.74	0.72	0.73	133394
<b>weighted avg</b>	0.75	0.75	0.75	133394

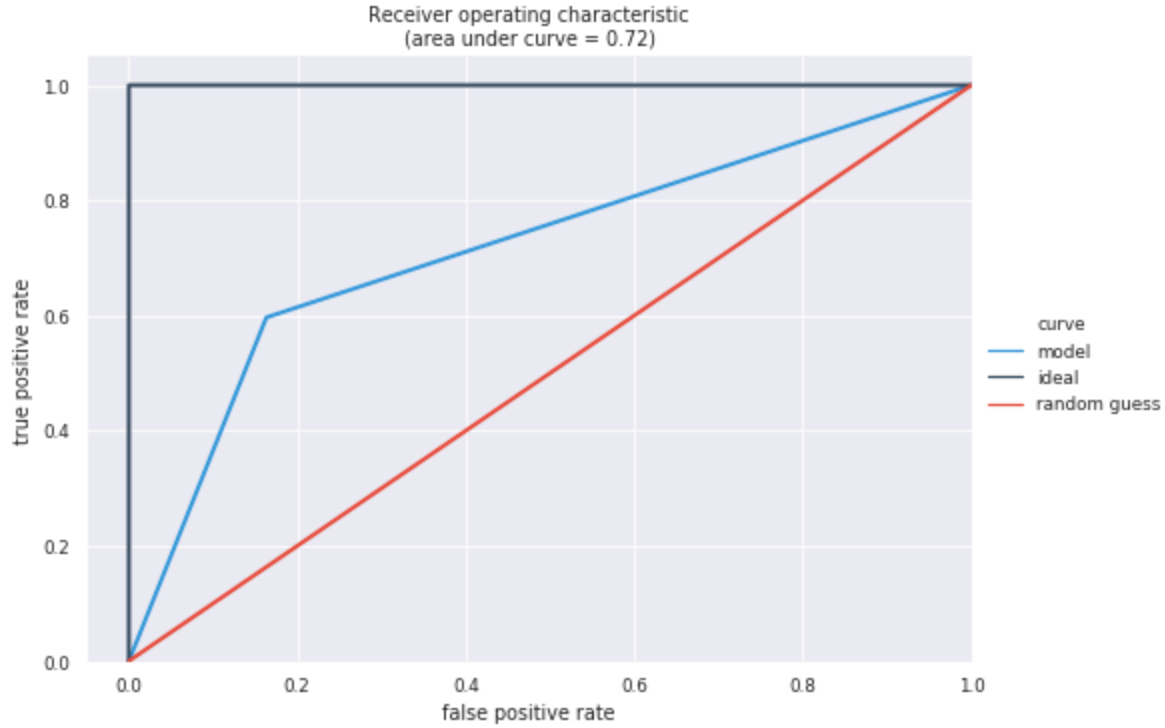
## Model Accuracy Score

The accuracy, recall and prediction scores for the model were computed as follows:

accuracy score	0.7526
precision score	0.6836
recall score	0.6118

## ROC Curve and AUC

The Receiver Operating Characteristic (ROC) curve is shown in light blue below. The area under curve (AUC) is 0.72



## Model 2 - XGBoost

The code for this is described in the notebook [tfidf\\_word2vec](#) (section Modeling)

The [XGBoost API](#) was used.

XGBoost (Extreme Gradient Boosting), is an implementation of the Gradient Boosting technique introduced in the paper [Greedy Function Approximation: A Gradient Boosting Machine, by Friedman](#).

XGBoost can be considered a special case of gradient descent in functional space,

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

where  $f_k \in F$  and  $F$  is the space of functions of all trees

The objective function to be minimized is given by,

$$Obj = \sum_{i=1}^n l(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k)$$

where  $l$  is the training loss and  $\Omega$  the tree complexity

The prediction at time-step  $t$  is given by,

$$\hat{y}_i^{(t)} = \hat{y}_i^{(t-1)} + f_t(x_i)$$

The objective at time-step  $t$  is given by,

$$Obj^{(t)} = \sum_{i=1}^n l(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \sum_{i=1}^t \Omega(f_i) + constant \approx \sum_{i=1}^n \left[ l(y_i, \hat{y}_i^{(t-1)}) + g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i) \right] + \Omega(f_t) + constant$$

where  $g_i = \partial_{\hat{y}_i^{(t-1)}} l(y_i, \hat{y}_i^{(t-1)})$ ,  $h_i = \partial_{\hat{y}_i^{(t-1)}}^2 l(y_i, \hat{y}_i^{(t-1)})$

If the tree function is defined as follows,

$$f_t(x) = w_{q(x)}$$

where  $w$  is the weight of the  $j$ th leaf and  $q(x)$  the mapping from  $x$  to leaf  $j$

And the tree complexity is defined as follows,

$$\Omega(f_t) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

where  $T$  is the total number of leaves and  $w_j$  is the weight of the  $j$ th leaf

Then the minimum value of the objective at  $t$  becomes,

$$Obj^{(t)} = -\frac{1}{2} \sum_{j=1}^T \frac{G_j^2}{H_j + \lambda} + \gamma T \quad \text{where the optimal leaf weight } w_j^* = -\frac{G_j}{H_j + \lambda}$$

where  $G_j = \sum_{(over\ jth\ leaf)} g_i$  and  $H_j = \sum_{(over\ jth\ leaf)} h_i$

The algorithm now finds the best possible tree using the above optimum value of the objective which can also be considered the score of the tree (analogous to the gini-value score in decision trees). However, there could be infinite number of possible trees at any time-step  $t$ . To get over this, the algorithm adopts a greedy approach and calculates the gain in objective at the split level (instead of the entire tree) and thus finds the next optimum split that minimizes the objective. It continues to make splits until the complexity or regularization factor for the total number of leaves in the tree,  $\gamma$ , makes it infeasible for making any more splits.

The hyperparameters of the model are:

1. n\_estimators(Number of boosted trees to fit)
2. reg\_lambda(L2 regularization term on weights)
3. max\_depth(Maximum tree depth for base learners)
4. learning\_rate(Boosting learning rate (xgb's "eta"))
5. gamma(Min. loss to make a partition on a leaf)

## Model Selection

5-fold cross-validation was used with scikit-learn API ([sklearn.model\\_selection.RandomizedSearchCV](#)) to find the optimum combination of hyperparameters.

The optimum parameters are as follows,

N_estimators (Number of boosted trees to fit)	200
Reg_lambda (L2 regularization term on weights)	6.31947
Max_depth (Maximum tree depth for base learners)	29
Learning_rate (Boosting learning rate (xgb's "eta"))	0.23
Gamma (Min. loss to make a partition on a leaf)	0.23

## Prediction

The following is the classification report,

	precision	recall	f1-score	support
<b>not duplicate</b>	0.87	0.88	0.87	84248
<b>duplicate</b>	0.79	0.76	0.78	49146
<b>micro avg</b>	0.84	0.84	0.84	133394
<b>macro avg</b>	0.83	0.82	0.83	133394
<b>weighted avg</b>	0.84	0.84	0.84	133394

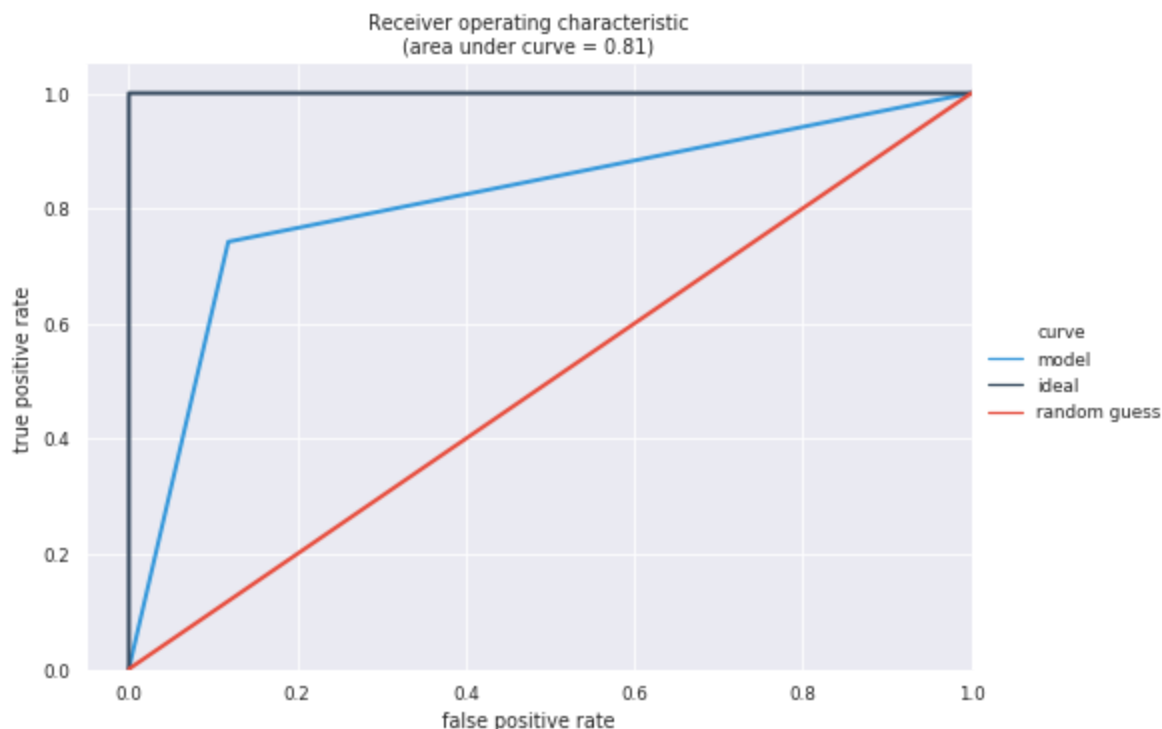
## Model Accuracy Score

The accuracy, recall and prediction scores for the model were computed as follows:

accuracy score	0.8383
precision score	0.7647
recall score	0.7898

## ROC Curve and AUC

The Receiver Operating Characteristic (ROC) curve is shown in light blue below. The area under curve (AUC) is 0.81



## Model 3 - Deep Learning with Decomposable Attention Model

The code for this is described in the notebook [DL\\_encode\\_attend](#) (section Method 1)

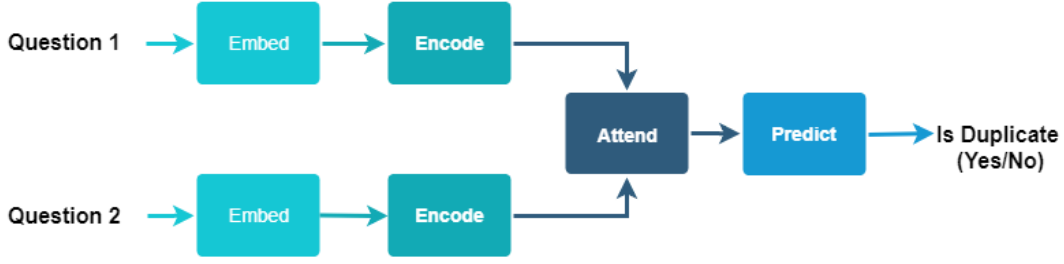
This method was adapted from the section [Example 1: A Decomposable Attention Model for Natural Language Inference](#) from a blog by Matthew Honnibal.

The blog example itself is adapted from [A Decomposable Attention Model for Natural Language Inference by Parikh et al.](#) where the problem of predicting a class label over a pair of sentences, where the class represents the logical relationship between them, is tackled using attention mechanism. A crucial advantage is provided by the fact that sentence-to-vector reduction operates over the pair of sentences jointly, as opposed to encoding the sentences into vectors independently.

The [Keras API](#) with [TensorFlow backend](#) was used. A single GPU (Nvidia GeForce GTX 980 Ti) with 5.2 compute capability and 6GB memory was used for computation.

Keras provides a high-level Neural Network Python API that runs on top of high performance numerical computation libraries like TensorFlow and Theano. It provides a modular and extensible api with standalone and customizable modules for neural layers, cost functions, optimizers, initialization schemes, activation functions and regularization schemes which can be combined in different configurations to create new models.





The layers of the model are:

1. Embed - Convert tokens from text to vector space. A map for token ID to word vector (based on [spacy's pretrained 300-dimensional word vectors trained on Common Crawl with GloVe](#)) is created before-hand, and each question is converted to a sequence of token IDs outside the model. The model then takes the questions as sequences of token IDs and the ID-vector map as input.
2. Encode - Normally, an encode layer would have aggregated the sequence of vectors in each question into a single question matrix. However, this encode layer performs [soft alignment \(Parikh et al. sec 3.1\)](#) with tokens in the other question to compute attention weights for the subsequent attention layer. The weights are computed using the dot product (pairwise cosine distance) of the vectors of question a with the vectors of question b as follows (here  $\bar{a}_i$  is the  $i$ th token vector of question a and  $\bar{b}_j$  is the  $j$ th token vector of question b,  $e_{ij}$  is the unnormalized attention weight between  $\bar{a}_i$  and  $\bar{b}_j$ ,  $F$  is the attention-based encoding function and  $F'$  is the dot product of this function):

$$e_{ij} = F'(\bar{a}_i, \bar{b}_j) = F(\bar{a}_i)^T F(\bar{b}_j)$$

Here  $F$  is a [feed-forward neural network with ReLU activations \(Parikh et al. sec 3.1\)](#)

Next, these weights are normalized as follows (here  $\beta_i$  is the sub-phrase of question b that is soft-aligned with the  $i$ th token  $\bar{a}_i$  of question a and  $\alpha_j$  is the sub-phrase of question a that is soft-aligned with the  $j$ th token  $\bar{b}_j$  of question b,  $l_b$  and  $l_a$  are the total tokens in question b and question a respectively)

$$\beta_i = \sum_{j=1}^{l_b} \frac{\exp(e_{ij})}{\sum_{k=1}^{l_b} \exp(e_{ik})} \bar{b}_j$$

$$\alpha_j = \sum_{i=1}^{l_a} \frac{\exp(e_{ij})}{\sum_{k=1}^{l_a} \exp(e_{kj})} \bar{a}_i$$

3. Attend - This step combines the above encoded matrices and reduces them to a single vector. This first involves a comparison step as follows,

$$v_{1i} = G([\bar{a}_i, \beta_i])$$

$$v_{2j} = G([\bar{b}_j, \alpha_j])$$

Here  $G$  is again a feed-forward neural network with ReLU activations and  $[.,.]$  represents the concatenation of two vectors.

Next, the two sets of comparison vectors are aggregated by summation,

$$v_1 = \sum_{i=1}^{l_a} v_{1i}$$

$$v_2 = \sum_{j=1}^{l_b} v_{2j}$$

4. Predict / Classify - This involves the final feed-forward network classifier H as follows

$$\hat{y} = H([v_1, v_2])$$

the predicted class is given by  $\hat{y} = \operatorname{argmax}_i (\hat{y}_i)$

The cross-entropy loss for the model is given by,

$$L(\theta_F, \theta_G, \theta_H) = \frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \log \frac{\exp(\hat{y}_c)}{\sum_{c'=1}^C \exp(\hat{y}_{c'})}$$

The learnable parameters of the model are:

$\theta_F, \theta_G, \theta_H$  the learnable parameters of the feed-forward networks F, G and H respectively

The hyperparameters of the model are:

1. Dropout regularization rate for network F
2. Dropout regularization rate for network G
3. Dropout regularization rate for network H
4. Optimization method (Adam / RMSProp / Nadam / SGD )
5. Learning rate for the optimization method
6. Number of samples (batch size) per gradient update
7. The numbers of epochs used for training

## Model Selection

[Hyperas](#) was used to perform automatic search for the set of optimum values in the model's multidimensional hyperparameter space. Six of the seven listed hyperparameters (above) were optimized using hyperas. The last parameter, number of epochs, was tuned manually.

Two search runs, one coarse-grained and the other fine-grained, were performed. The coarse-grained search determined the batch-size and the fine-grained search determined optimum values for the remaining parameters.

The optimum parameters are as follows,

Dropout regularization rate for network F	0.8935
Dropout regularization rate for network G	0.0716
Dropout regularization rate for network H	0.189
Optimization method (Adam / RMSProp / Nadam / SGD )	SGD
Initial Learning rate (an exponential decay of 0.1 was used)	0.0244
Number of samples (batch size) per gradient update	64
The numbers of epochs used for training	50

The variation of training and validation accuracy during the fine-grained search (50 evaluations) is shown below,



## Hyperas

Hyperas is a high-level interface to the Hyperparameter optimization library [Hyperopt](#), which makes it very easy to be used with Keras.

Hyperopt offers the Tree-structured Parzen Estimator (TPE) algorithm for efficiently searching in high-dimensional spaces. This is a Sequential Model Based Optimization (SMBO) method, which builds and evaluates models sequentially, based on historic values of hyperparameters.

## Early Stopping

The Keras implementation of [early stopping](#) was used during the hyperparameter search process to prevent the search algorithm from wasting time exploring local minima. This is implemented in the form of a callback which stops training if there is no improvement in the objective function.

## Prediction

The following is the classification report,

	<b>precision</b>	<b>recall</b>	<b>f1-score</b>	<b>support</b>
<b>not duplicate</b>	0.85	0.82	0.84	84267
<b>duplicate</b>	0.72	0.76	0.74	49148
<b>micro avg</b>	0.80	0.80	0.80	133415
<b>macro avg</b>	0.79	0.79	0.79	133415
<b>weighted avg</b>	0.80	0.80	0.80	133415

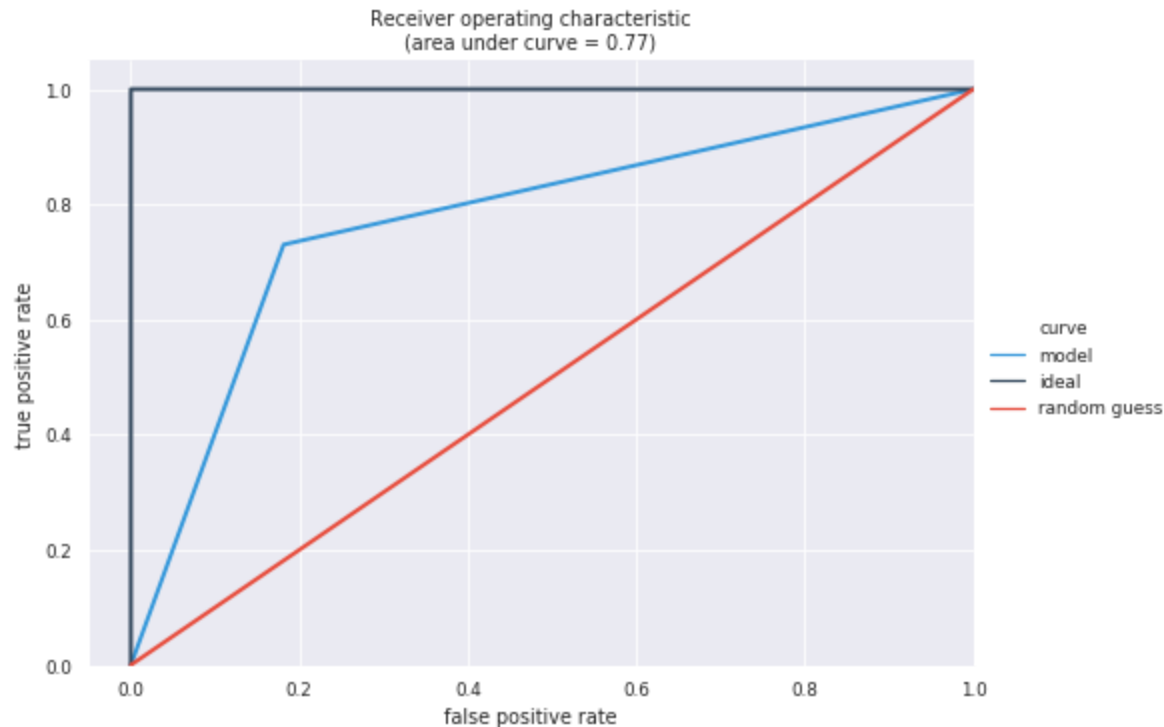
## Model Accuracy Score

The accuracy, recall and prediction scores for the model were computed as follows:

accuracy score	0.8004
precision score	0.7604
recall score	0.7156

## ROC Curve and AUC

The Receiver Operating Characteristic (ROC) curve is shown in light blue below. The area under curve (AUC) is 0.77



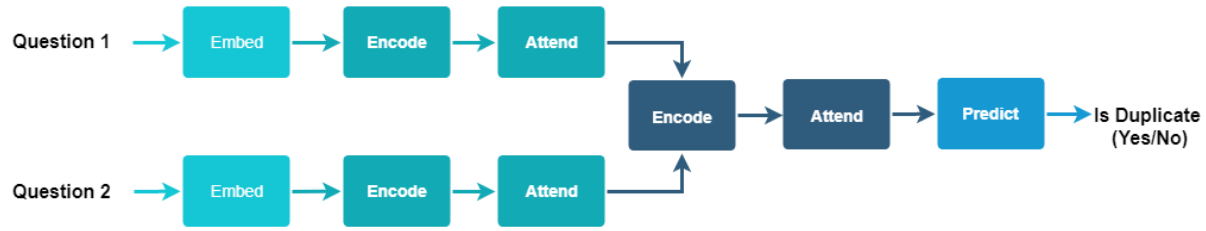
## Model 4 - Deep Learning with Hierarchical Attention Model

The code for this is described in the notebook [DL\\_encode\\_attend](#) (section Method 2)

This method was adapted from the section [Example 2: Hierarchical Attention Networks for Document Classification](#) from the same blog as the previous method, by Matthew Honnibal.

The blog example itself is adapted from [Hierarchical Attention Networks for Document Classification by Yang et al.](#) where the problem of classification of documents is tackled using attention mechanism. The method uses attention to reduce a matrix to a vector. It does this by learning context vectors for the two attention steps (one for words and one for sentences), which is analogous to learning words or sentences from a document that summarize it. Applying the word and sentence context-vector analogy to a generic NLP problem, the attention steps can be seen as an effective feature extraction process, which in the case of duplicate question identification, can be employed to learn similarity between words and sentences belonging to a pair of questions.

Again [Keras API](#) with [TensorFlow backend](#) was used.



The layers of the model are:

1. Embed - Convert tokens from text to vector space. A map for token ID to word vector (based on [spacy's pretrained 300-dimensional word vectors trained on Common Crawl with GloVe](#)) is created before-hand, and each question is converted to a sequence of token IDs outside the model. The model then takes the questions as sequences of token IDs and the ID-vector map as input.
2. Encode Words - This uses a bidirectional GRU to encode each word into a context-specific vector. For every token vector  $x_{it}$  (embedded word vector for the word  $w_{it}$ ) a forward GRU encodes the forward hidden state  $h_{it}^{\rightarrow}$  and a backward GRU encodes the backward hidden state  $h_{it}^{\leftarrow}$ . A bidirectional wrapper layer then concatenates the forward and backward hidden states to summarize the information of the whole sentence centered around the word as  $h_{it}$ :

$$x_{it} = W_e w_{it} \quad (\text{here } W_e \text{ is the embedding matrix})$$

$$h_{it}^{\rightarrow} = GRU^{\rightarrow}(x_{it})$$

$$h_{it}^{\leftarrow} = GRU^{\leftarrow}(x_{it})$$

$$h_{it} = [h_{it}^{\rightarrow}, h_{it}^{\leftarrow}]$$

Here  $GRU$  is a [recurrent neural network that uses a gating mechanism to track the state of sequences without using separate memory cells \(sec 2.1 of Yang et al.\)](#)

3. Word Attention - This step extracts the important words that are representative of the meaning of the sentence (or question) by using attention mechanism and computes a single vector for each question by giving higher weights to the more important words. First  $h_{it}$  (the encoded vector of every word  $w_{it}$ ) is converted to a hidden representation  $u_{it}$  by passing it through a Multilayer Perceptron (MLP) with one hidden layer and the tanh activation function with learned parameters  $W_w$  and  $b_w$ ,

$$u_{it} = \tanh(W_w h_{it} + b_w)$$

Next, the normalized importance of the word  $w_{it}$  is measured as the normalized similarity of its hidden vector  $u_{it}$  with the context vector  $u_w$  using the softmax function (the context vector  $u_w$  is a randomly initialized learned parameter of the attention layer),

$$\alpha_{it} = \frac{\sum_t \frac{\exp(u_{it}^T u_w)}{\sum_t \exp(u_{it}^T u_w)}}$$

Now, a single sentence vector is computed for each question by aggregating the vectors  $h_{it}$  and the weights  $\alpha_{it}$  with summation as follows,

$$s_i = \sum_t \alpha_{it} h_{it}$$

5. Encode Sentences - This uses a bidirectional GRU to encode each sentence (question) into a context-specific vector. For every sentence vector  $s_i$  (single sentence vector computed in the previous step for every question) a forward GRU encodes the forward state  $h_i^{\rightarrow}$  and a backward GRU encodes the backward hidden state  $h_i^{\leftarrow}$ . A bidirectional wrapper layer then concatenates the forward and backward hidden states to summarize the information of the set of sentences centered around the sentence as  $h_i$ . Unlike the case of document classification where bidirectional GRU may deal with any number of sentence vectors at a time, the duplicate question problem deals with only two sentence (question) vectors at a time:

$$h_i^{\rightarrow} = GRU^{\rightarrow}(s_i)$$

$$h_i^{\leftarrow} = GRU^{\leftarrow}(s_i)$$

$$h_i = [h_i^{\rightarrow}, h_i^{\leftarrow}]$$

6. Sentence Attention - This step extracts the important features that are representative of the relationship between the two sentences (or questions) by using attention mechanism. It computes a single vector for both questions by giving higher weights to the more important features. First  $h_i$  (the encoded vector of every sentence  $s_i$ ) is converted to a hidden representation  $u_i$  by passing it through a Multilayer Perceptron (MLP) with one hidden layer and the tanh activation function with learned parameters  $W_s$  and  $b_s$ ,

$$u_i = \tanh(W_s h_i + b_s)$$

Next, the normalized importance of the sentence  $s_i$  is measured as the normalized similarity of its hidden vector  $u_i$  with the context vector  $u_s$  using the softmax function (the context vector  $u_s$  is a randomly initialized learned parameter of the attention layer),

$$\alpha_i = \frac{\exp(u_i^T u_s)}{\sum_i \exp(u_i^T u_s)}$$

Now, a single sentence vector is computed for both questions by aggregating the vectors  $h_i$  and the weights  $\alpha_i$  with summation as follows,

$$v = \sum_i \alpha_i h_i$$

7. Predict / Classify - This involves the final MLP with sigmoid activation as follows,

$$\hat{y} = \text{sigmoid}(W_c v + b_c)$$

the predicted class is given by  $\hat{y} = \text{argmax}_i(\hat{y}_i)$

The negative log-likelihood loss for the model is given by,

$$L = -\sum_d \log \hat{y}_d$$

The learnable parameters of the model are:

$W_w, b_w, u_w$  of the Word Attention Layer

$W_s$ ,  $b_s$ ,  $u_s$  of the Sentence Attention Layer

The hyperparameters of the model are:

1. Dropout regularization rate (word encoder)
2. Dropout regularization rate (sentence encoder)
3. Optimization method (Adam / RMSProp / Nadam / SGD )
4. Initial Learning rate (for the exponential decay)
5. Learning rate decay (for the exponential decay)
6. Number of samples (batch size) per gradient update
7. The numbers of epochs used for training

## Model Selection

[Hyperopt](#) was used to perform automatic search for the set of optimum values in the model's multidimensional hyperparameter space. Six of the seven listed hyperparameters (above) were optimized using hyperopt. The final parameter, number of epochs, was tuned manually.

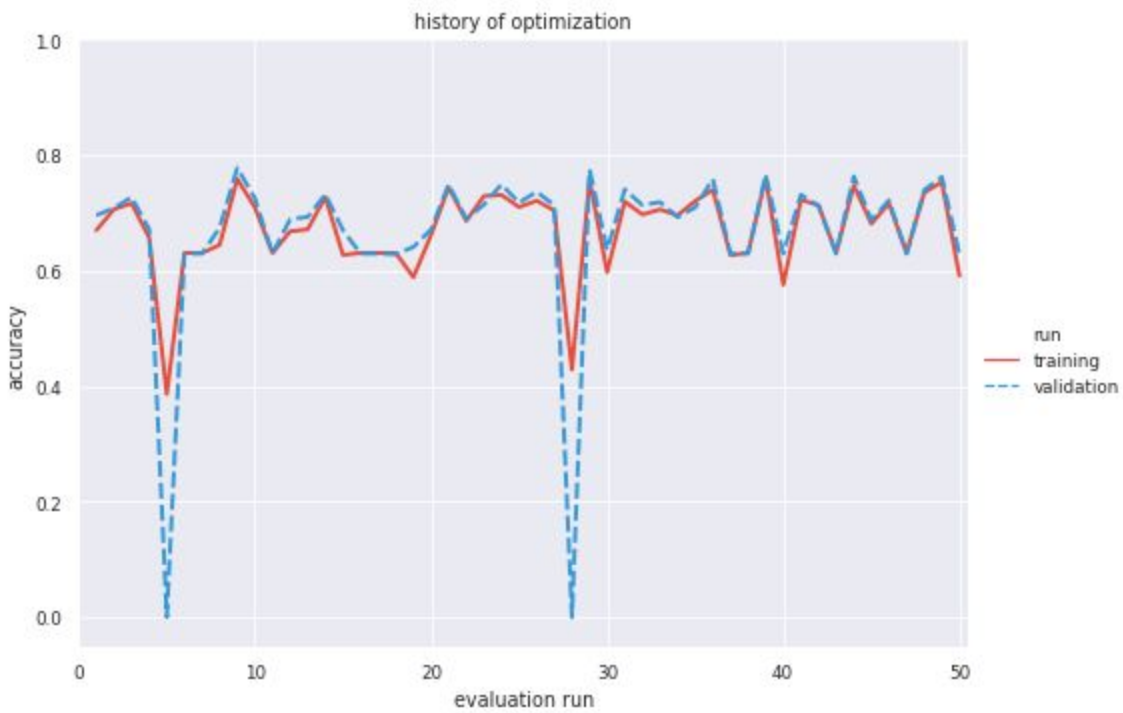
Two search runs, one coarse-grained and the other fine-grained, were performed. The coarse-grained search determined the batch-size and optimization method while the fine-grained search determined optimum values for the remaining parameters.

The optimum parameters are as follows,

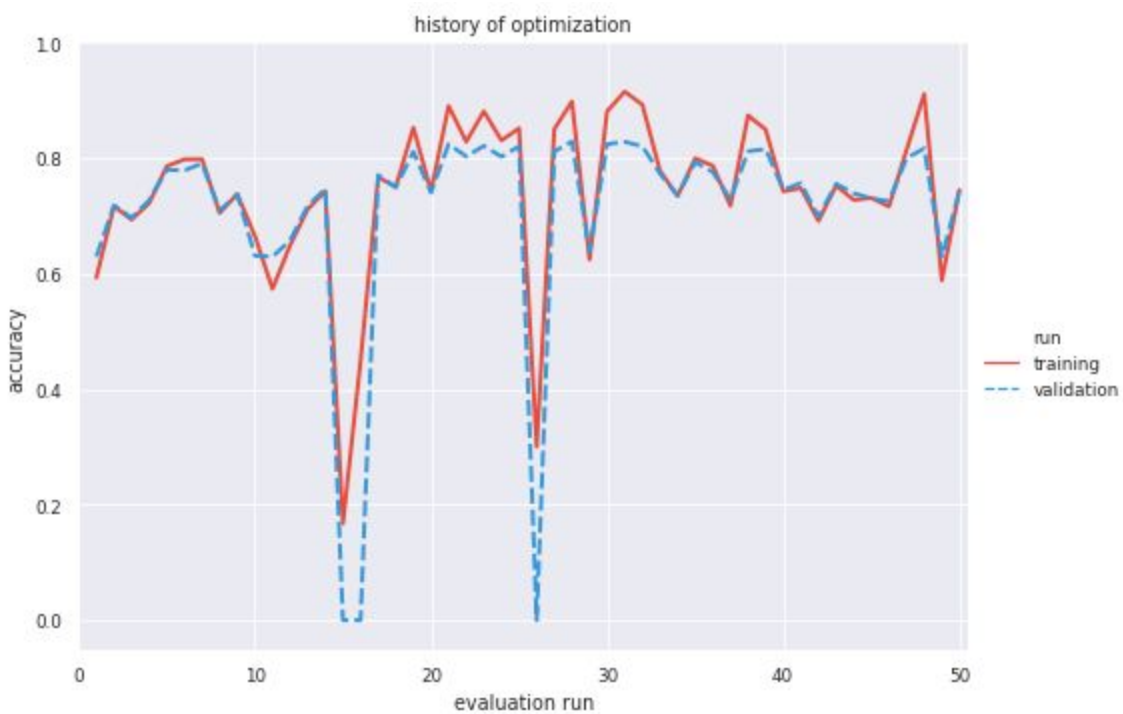
Dropout regularization rate (word encoder)	0.3503
Dropout regularization rate (sentence encoder)	0.0074
Optimization method (Adam / RMSProp / Nadam / SGD )	RMSProp
Initial Learning rate (for the exponential decay)	0.0028
Learning rate decay (for the exponential decay)	0.1875
Number of samples (batch size) per gradient update	128
The numbers of epochs used for training	40



The variation of training and validation accuracy during the coarse-grained search (50 evaluations) is shown below,



The variation of training and validation accuracy during the fine-grained search (50 evaluations) is shown below,



## Hyperopt

For the sake of experimentation and comparison, the lower level API Hyperopt was used in this case (instead of Hyperas). Hyperopt proved to be fairly easy to use with a less limiting interface than Hyperas.

## Early Stopping

The Keras implementation of [early stopping](#) was used during the hyperparameter search process to prevent the search algorithm from wasting time exploring local minima. This is implemented in the form of a callback which stops training if there is no improvement in the objective function.

## Learning Rate Decay

An exponential decay was used for the learning rate with help of the Keras callback [LearningRateScheduler](#). This creates a fast gradient descent in the beginning few epochs and then slows down quickly to gradually approach the minimum (without 'jumping' over it). Both the initial learning rate and the rate of exponential decay were tuned using Hyperopt.

## Prediction

The following is the classification report,

	precision	recall	f1-score	support
not duplicate	0.88	0.86	0.87	84267
duplicate	0.76	0.79	0.78	49148
micro avg	0.83	0.83	0.83	133415
macro avg	0.82	0.83	0.82	133415
weighted avg	0.84	0.83	0.83	133415

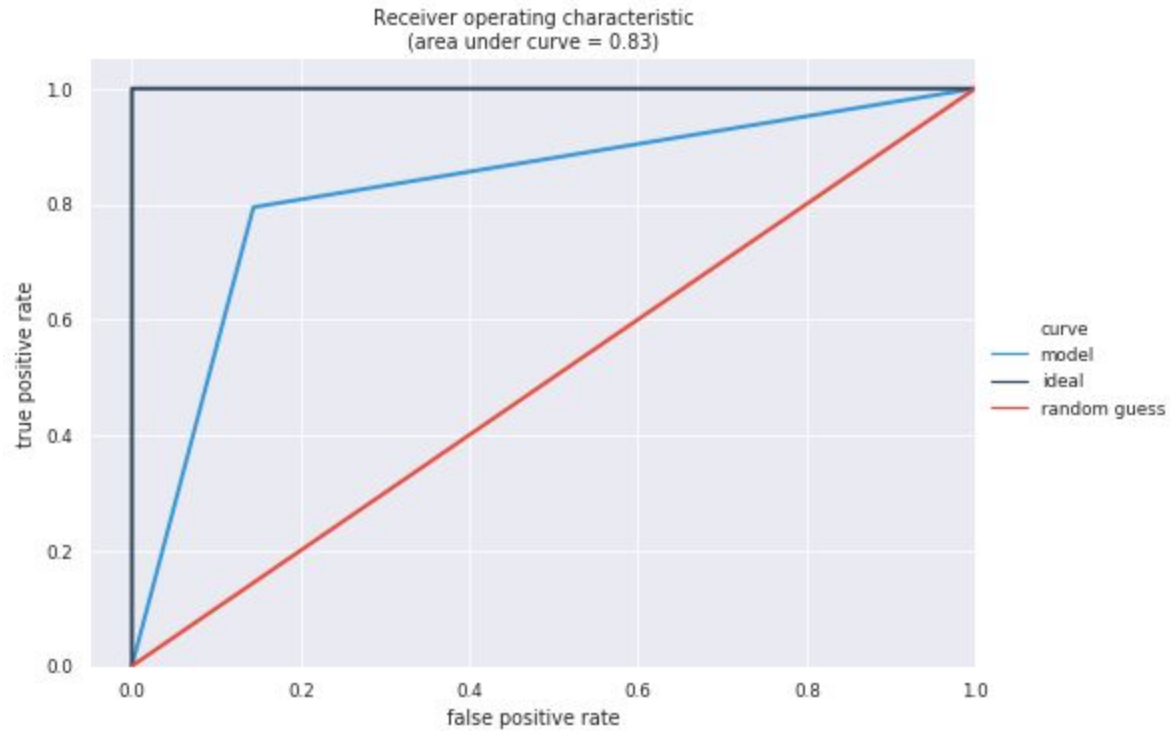
## Model Accuracy Score

The accuracy, recall and prediction scores for the model were computed as follows:

accuracy score           0.8332  
precision score         0.7627  
recall score             0.7945

## ROC Curve and AUC

The Receiver Operating Characteristic (ROC) curve is shown in light blue below. The area under curve (AUC) is 0.83



## Conclusion

A high precision is more important in this problem, as it is necessary to have as low a number of false positives as possible (i.e. question pairs which are not duplicate but which the model falsely classifies as duplicate). A high recall (low number of false negatives i.e. question pairs which are duplicate but which the model falsely classifies as not-duplicate) is also good to have as this reduces the number of redundant answers and also avoids the unnecessary overhead of answering the same question more than once.

It is obvious from the model evaluation results for the various approaches used, that XGBoost as well as Deep Learning with Hierarchical Attention are effective ways to solve the problem of duplicate identification.

# Further Exploration

Further exploration involves adding more layers to the deep learning models and also exploring alternative neural network architectures. In addition to this, exploring alternative hyperparameter tuning methods could also improve the accuracy and precision.

## References

### NLP

1. [On word embeddings - Part 1 by Sebastian Ruder](#)
2. [GloVe - On word embeddings - Part 3 by Sebastian Ruder](#)
3. [Spacy](#)
4. [Gensim](#)
5. [Fasttext](#)
6. [FuzzyWuzzy](#)

### Modeling

1. [Logistic Regression](#)
2. [XGBoost](#)
3. [Hyper-parameter Tuning](#)
4. [Embed, encode, attend, predict: The new deep learning formula for state-of-the-art NLP models, by Matthew Honnibal](#)
5. [Hierarchical Attention Networks for Document Classification by Yang et al.](#)
6. [A Decomposable Attention Model for Natural Language Inference by Parikh et al.](#)

### Deep Learning

1. [Keras](#)
2. [TensorFlow](#)
3. [Hyperas](#)
4. [Hyperopt tutorial for Optimizing Neural Networks' Hyperparameters](#)
5. [Hyperopt](#)