

Udacity Machine Learning Nanodegree Capstone  
**Quora Duplicate Question Detection**

*Raahul Seshadri*  
April, 2017

# Contents

<b>1</b>	<b>Definition</b>	<b>2</b>
1.1	Project Overview . . . . .	2
1.2	Problem Statement . . . . .	2
1.3	Metrics . . . . .	3
<b>2</b>	<b>Analysis</b>	<b>4</b>
2.1	Data Exploration . . . . .	4
2.1.1	Why do we have question IDs for each pair? . . . . .	4
2.1.2	Summary statistics . . . . .	5
2.2	Exploratory Visualization . . . . .	5
2.2.1	Summary visualization . . . . .	5
2.2.2	Word count differences . . . . .	5
2.2.3	Jaccard index . . . . .	7
2.3	Algorithms & Techniques . . . . .	8
2.3.1	Feature extraction . . . . .	8
2.3.2	Pooling . . . . .	9
2.3.3	Neural Networks . . . . .	10
2.3.4	Improving Neural Network Accuracy . . . . .	10

# 1 Definition

## 1.1 Project Overview

Websites like Quora<sup>1</sup>, StackExchange<sup>2</sup> network (includes the popular programming website StackOverflow<sup>3</sup>) allow the users to post questions, and the entire community can answer those questions. Ideally, every unique question should be present just once in the system, so that every answer to those questions are present at once place.

However, users are prone to posting duplicate questions, because not all of them check if the question they're asking has already been asked by someone else. This necessitates having an automated duplicate question checker that can check if a question is a duplicate. Whenever the user tries to post a new question, the system can suggest an existing one for perusal.

This project was inspired by Quora's Kaggle challenge<sup>4</sup>. A dataset of question pairs, manually tagged as duplicate or not by human reviewers, has been provided<sup>5</sup> by Quora to train on.

**For the purpose of this project, the inputs are strictly assumed to be in English.**

## 1.2 Problem Statement

The problem can be stated as follows:

**Given a pair of question texts, detect if they are duplicate or not.**

For example, given the following two questions:

*Q1: What is the average salary in India?*  
*Q2: What is the average salary in the United States?*

The system could flag them as either duplicates or not-duplicates. The steps to achieving our duplicate detector are as follows:

1. Download and pre-process the Quora training dataset
2. Split the Quora training dataset into 90% training and 10% test set.
3. In the 90% training set, further 10% will be used as a validation set.
4. Extract usable features from the dataset.
5. Train a binary classifier to differentiate between duplicates/non-duplicates.
6. Provide a command-line interface so that users can check for duplicates from Quora's test set, or provide his own questions.

---

<sup>1</sup><https://www.quora.com>

<sup>2</sup><https://stackexchange.com>

<sup>3</sup><https://stackoverflow.com>

<sup>4</sup><https://www.kaggle.com/c/quora-question-pairs>

<sup>5</sup><https://www.kaggle.com/c/quora-question-pairs>

### 1.3 Metrics

Being a binary classification problem, the metric is going to be accuracy:

$$\text{accuracy} = \frac{\text{true positives} + \text{true negatives}}{\text{total samples}}$$

This accuracy will be measured on the test set (10% of the total dataset).

1. **False negatives**, where a duplicate pair is not detected as such, will lead to degraded user experiences. However, it is also rather harmless, since no information is lost. There is still an opportunity to manually flag them as duplicates, or build such a feature on the platform that is using the duplicate detector.
2. **False positives**, where a question pair that's not a duplicate, but is marked as such, is more dangerous. It can lead to very degraded user experience, since it directly hinders the user trying to perform an action. It is better for the system to err on the side of false negatives than false positives.

## 2 Analysis

### 2.1 Data Exploration

Below is the head (first 20 entries) of the Quora dataset:

id	qid1	qid2	question1	question2	is_duplicate
0	0	1	2	What is the step by step guide to invest in sh... What is the step by step guide to invest in sh...	0
1	1	3	4	What is the story of Kohinoor (Koh-i-Noor) Dia... What would happen if the Indian government sto...	0
2	2	5	6	How can I increase the speed of my internet co... How can Internet speed be increased by hacking...	0
3	3	7	8	Why am I mentally very lonely? How can I solve... Find the remainder when $23^{24}$ is divided by 100...	0
4	4	9	10	Which one dissolve in water quikly sugar, salt... Which fish would survive in salt water?	0
5	5	11	12	Astrology: I am a Capricorn Sun Cap moon and c... I'm a triple Capricorn (Sun, Moon and ascendan...	1
6	6	13	14	Should I buy tiago? What keeps childern active and far from phone ...	0
7	7	15	16	How can I be a good geologist? What should I do to be a great geologist?	1
8	8	17	18	When do you use $\ni$ instead of $\hookleftarrow$ ? When do you use "&" instead of "and"?	0
9	9	19	20	Motorola (company): Can I hack my Charter Moto... How do I hack Motorola DCX3400 for free internet?	0
10	10	21	22	Method to find separation of slits using fresn... What are some of the things technicians can te...	0
11	11	23	24	How do I read and find my YouTube comments? How can I see all my Youtube comments?	1
12	12	25	26	What can make Physics easy to learn? How can you make physics easy to learn?	1
13	13	27	28	What was your first sexual experience like? What was your first sexual experience?	1

The columns correspond to the following features:

1. **id**: ID to identify the question pair uniquely
2. **qid1**: ID of the first question in the question pair
3. **qid2**: ID of the second question in the question pair
4. **question1**: Text of the first question in the question pair
5. **question2**: Text of the second question in the question pair
6. **is\_duplicate**: Quora reviewer's decision on whether the question pair is duplicate (1) or not (0)

#### 2.1.1 Why do we have question IDs for each pair?

There are two columns of interest, "qid1" and "qid2" that are of interest. The algorithm that we're going to design will only tell us if a given pair of question is a duplicate. However, given a question, we'd also like to find out all the duplicates of it in the system.

For example, if question "1" is a duplicate of question "2", and question "2" is a duplicate of "3", then a system should also know that "1" and "3" are duplicates. Data structures like "union find" allow us to do that. However, this is not in scope of the algorithm that we'll design, but the responsibility of the higher system that will use this duplicate detector.

Which is why the question IDs won't be used as a feature for the duplicate detector, but is still important for the system using the duplicate detector.

### 2.1.2 Summary statistics

Below are some basic summary statistics:

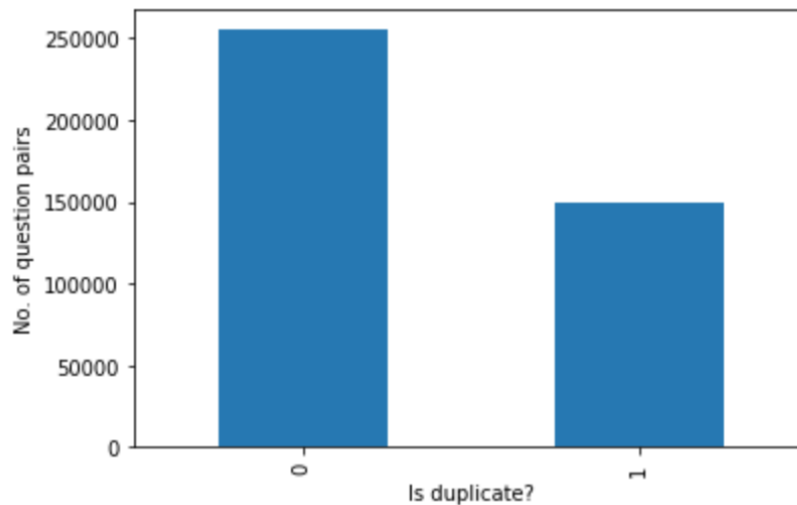
1. **Total entries:** 4,04,290
2. **Total positive entries:** 1,49,263
3. **Total negative entries:** 2,55,027
4. **Percent positive entries:** 36.91%
5. **Percent negative entries:** 63.08%

We have a sample with more negative than positive examples.

## 2.2 Exploratory Visualization

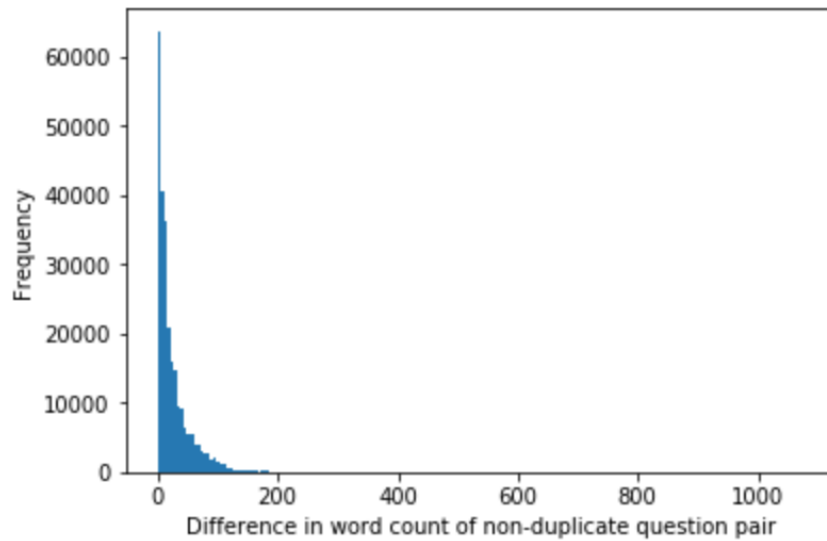
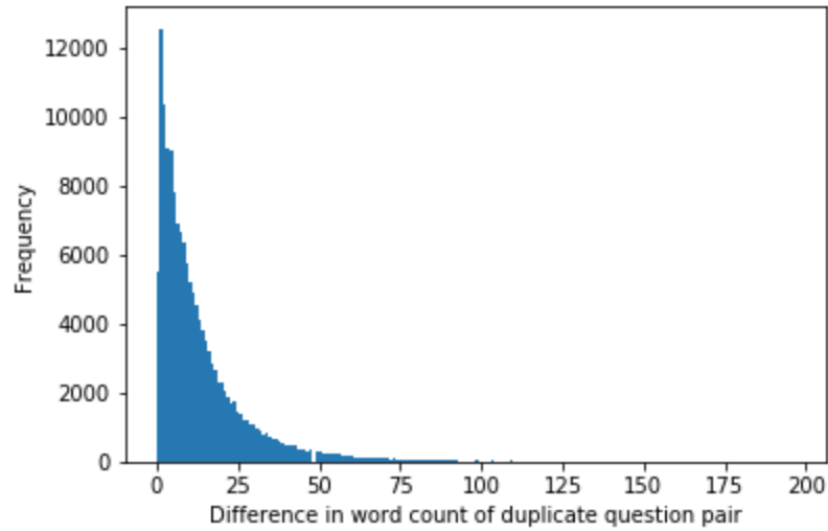
### 2.2.1 Summary visualization

Let's look at the class distribution (duplicate or not). This is the same from our summary statistics.



### 2.2.2 Word count differences

Let's look at if word count has any bearing on whether a question is marked as duplicate or not, or can we discern any important facts/patterns.



1. For duplicates, the word count difference tapers smoothly
2. For non-duplicates, there is a sharp decrease in the 2nd bin
3. There are question pairs that have word count difference 75 and above, and are still duplicates. So simple bag of words model will not work. *This is the most important conclusion.*

### 2.2.3 Jaccard index

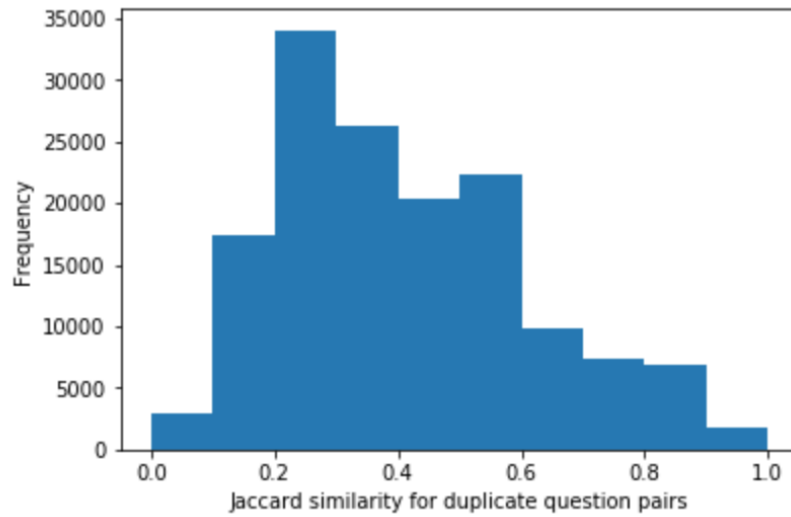
To explore a baseline model, let's look at the bag-of-words similarity of duplicate and non-duplicate question pairs using the Jaccard Index<sup>6</sup>. The relevant formula is:

$$J(A, B) = \frac{A \cap B}{A \cup B}$$

In our case, A and B refer to individual words in the question pairs. So, the formula basically is:

$$J(A, B) = \frac{\text{No. of common words in question1 and question2}}{\text{No. of unique words in question1 and question2 combined}}$$

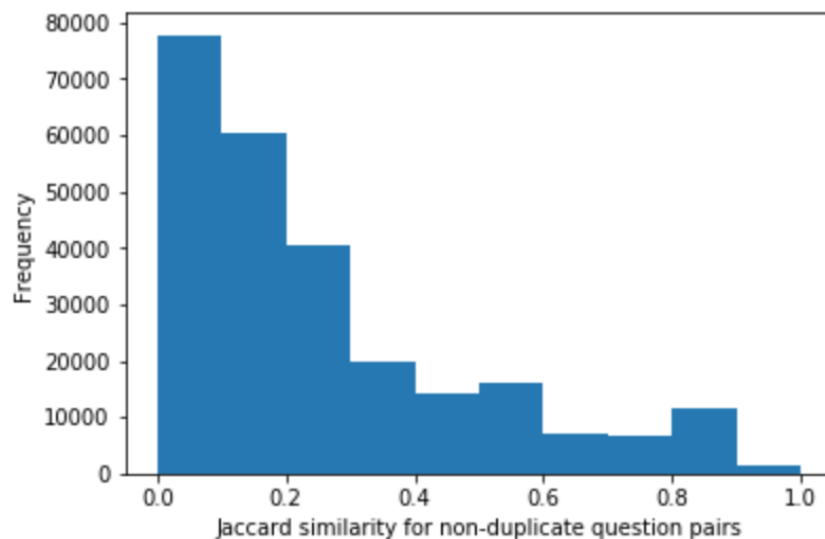
Let's look at the Jaccard similarity of duplicate and non-duplicate question pairs separately.



---

<sup>6</sup>[https://en.wikipedia.org/wiki/Jaccard\\_index](https://en.wikipedia.org/wiki/Jaccard_index)





1. For duplicates, Jaccard similarity is like a gamma distribution. Jaccard is bad at predicting a pair is a duplicate.
2. For non-duplicates, the distribution is again gamma but with a right skew. So Jaccard positively gives low scores for most of the non-duplicates.

We had previously concluded from the word counts that bag-of-words similarity will not be a good method of differentiating between duplicates and non-duplicates, and the two graphs above confirm the fact.

## 2.3 Algorithms & Techniques

### 2.3.1 Feature extraction

The inputs to our algorithm are two variable length strings. In NLP (Natural Language Processing), there are two primary ways of extracting features from text, after the input has been word-tokenized, that is being split into individual words:

#### **Bag of words:**

Here, every word is assigned a unique number. The input to any ML algorithm then becomes a one-hot encoded array of numbers that correspond to these words.

For example, if the input corpus has two words: "hello" and "world". And we assign 0 to "hello" and 1 to "world". Then, the word "hello" would be represented as [1 0] while "world" would be represented as [0 1]. (one hot encoded).

The number of features become equal to the number of unique words in the corpus (training set).

### Word vectors:

The problem with one hot encoding is that we come up with very large feature vectors, since they're one hot encoded. This does not scale to large corpora.

Instead of a sparse array, word vectors convert input words into a fixed-length vector. A common vector size is 300. Thus, every word in the input corpus gets converted into a dense feature vector of length, say, 300, regardless of how many words there are in the corpus.

As opposed to bag of words, word vectors have to be trained on a corpus, before they can be used to extract features from different corpora. Thankfully, there are a lot of pre-trained models for word vectors available.

Two popular algorithms to generate them are Word2Vec<sup>7</sup> and GloVe<sup>8</sup>. spaCy<sup>9</sup> comes with a GloVe model<sup>10</sup> of 300-dimensional vectors trained on the Common Crawl corpus. This is the Python library that we will use to do all our processing.

Naturally, I'll be using the word vectors approach, since it also allows me to train for words that don't appear in the Quora training corpus, but a similar word does.

### 2.3.2 Pooling

Given a pair of questions, the questions themselves can be of different sizes. If we get word vectors for each word in the question, then we get a matrix of size  $m \times 300$ , where we're assuming that 300 is the word vector size and  $m$  is the number of words in the question.

However, for a pair,  $m$  could be different. Thus, we need a way to normalize this. There are two ways<sup>11</sup>:

1. **Max pooling:** Take the maximum of each column in the  $m \times 300$  matrix to get a  $1 \times 300$  matrix.
2. **Mean pooling:** Take the mean of each column in the  $m \times 300$  matrix to get a  $1 \times 300$  matrix.

Max pooling is reported to perform better<sup>12</sup>. However, I've decided to concatenate mean and max pooling to get a matrix  $1 \times 600$  in dimension. Thus, every question is converted into 600-dimensional vector.

---

<sup>7</sup><https://en.wikipedia.org/wiki/Word2vec>

<sup>8</sup><https://nlp.stanford.edu/projects/glove/>

<sup>9</sup><https://spacy.io>

<sup>10</sup><https://spacy.io/docs/usage/word-vectors-similarities>

<sup>11</sup><https://explosion.ai/blog/quora-deep-text-pair-classification#example-neural-network>

<sup>12</sup><https://explosion.ai/blog/quora-deep-text-pair-classification#results>

### 2.3.3 Neural Networks

We will be using Neural Networks to create our model. Neural Networks, especially deep ones, have lately become popular in the field of NLP.<sup>13 14 15</sup>

A neural network learns to model functions, just like supervised algorithms. However, the deeper the neural network (deeper meaning more number of layers), the more complex characteristics it can detect in the input.

A downside to neural networks is that it is incredibly difficult to correctly design them. The network structure itself has many possibilities, and then we also have a number of hyperparameters like epochs, depending on the optimizer being used. They also take a long time to train.

### 2.3.4 Improving Neural Network Accuracy

There are several techniques in neural networks to improve the accuracy of the network, or to reduce overfitting.

#### Dropout:

Here, random units inside of the network are switched off in every epoch. Thus, the network is forced to learn redundant representations of the input. This, in totality, avoids overfitting. Dropout only happens during training, and is switched off when running the model in production.<sup>16 17</sup>

Due to the low amount of training data that we have, I'm not using dropouts in my network.

#### Batch Normalization:

A Batch Normalization layer shifts the inputs from the previous layer to have zero-mean and unit-variance. This prevents data flowing in the network to not become too big or too small. It is said to result in higher accuracy and faster learning (convergence).<sup>18</sup> To try this out, I trained variants of my network with Batch Normalization on and off.

## 2.4 Benchmark

In the field of Information Retrieval, there is a problem called Near Duplicate Detection.<sup>19</sup> Jaccard Similarity of bag-of-words ( $k = 1$  shinglings) is one way of finding near duplicates.<sup>20</sup>

---

<sup>13</sup><http://colah.github.io/posts/2014-07-NLP-RNNs-Representations>

<sup>14</sup><https://www.quora.com/How-are-neural-networks-used-in-Natural-Language-Processing>

<sup>15</sup><https://nlp.stanford.edu/projects/DeepLearningInNaturalLanguageProcessing.shtml>

<sup>16</sup><https://www.quora.com/How-does-the-dropout-method-work-in-deep-learning>

<sup>17</sup><https://www.cs.toronto.edu/~hinton/absps/JMLRdropout.pdf>

<sup>18</sup><https://www.quora.com/Why-does-batch-normalization-help>

<sup>19</sup><https://nlp.stanford.edu/IR-book/html/htmledition/near-duplicates-and-shingling-1.html>

<sup>20</sup><http://stackoverflow.com/questions/23053688/how-to-detect-duplicates-among-text-documents-and-return-the-duplicates-similar>

Since we split Quora dataset into 10% test set, this benchmark model would be run on the test set. This baseline model gives an accuracy of 65.153%.