# 1.2.1 : EDA: Advanced Feature Extraction.

```python
1   import warnings
2   warnings.filterwarnings("ignore")
3   import numpy as np
4   import pandas as pd
5   import seaborn as sns
6   import matplotlib.pyplot as plt
7   from subprocess import check_output
8   %matplotlib inline
9   import plotly.offline as py
10  py.init_notebook_mode(connected=True)
11  import plotly.graph_objs as go
12  import plotly.tools as tls
13  import os
14  import gc
15
16  import re
17  from nltk.corpus import stopwords
18  import distance
19  from nltk.stem import PorterStemmer
20  from bs4 import BeautifulSoup
21  import re
22  from nltk.corpus import stopwords
23  # This package is used for finding longest common subsequence between two strings
24  # you can write your own dp code for this
25  import distance
26  from nltk.stem import PorterStemmer
27  from bs4 import BeautifulSoup
28  from fuzzywuzzy import fuzz
29  from sklearn.manifold import TSNE
30  # Import the Required lib packages for WORD-Cloud generation
31  # https://stackoverflow.com/questions/45625434/how-to-install-wordcloud-in-python3-6
32  from wordcloud import WordCloud, STOPWORDS
33  from os import path
34  from PIL import Image
```

```python
#https://stackoverflow.com/questions/12468179/unicodedecodeerror-utf8-codec-cant-decode-byte-0x9c
if os.path.isfile('df_fe_without_preprocessing_train.csv'):
    df = pd.read_csv("df_fe_without_preprocessing_train.csv",encoding='latin-1')
    df = df.fillna('')
    df.head()
else:
    print("get df_fe_without_preprocessing_train.csv from drive or run the previous notebook")
```

```python
df.head(2)
```

| | id | qid1 | qid2 | question1 | question2 | is_duplicate | freq_qid1 | freq_qid2 | q1len | q2len | q1_n_words | q2_n_ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 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 | 66 | 57 | 14 | 12 |
| 1 | 1 | 3 | 4 | What is the story of Kohinoor (Koh-i-Noor) Dia... | What would happen if the Indian government sto... | 0 | 4 | 1 | 51 | 88 | 8 | 13 |

## 3.4 Preprocessing of Text

- Preprocessing:
  - Removing html tags
  - Removing Punctuations
  - Performing stemming
  - Removing Stopwords
  - Expanding contractions etc.

```python
 1   # To get the results in 4 decemal points
 2   SAFE_DIV = 0.0001
 3
 4   STOP_WORDS = stopwords.words("english")
 5
 6
 7   def preprocess(x):
 8       x = str(x).lower()
 9       x = x.replace(",000,000", "m").replace(",000", "k").replace("′", "'").replace("'", "'")\
10                           .replace("won't", "will not").replace("cannot", "can not").replace("can't", '
11                           .replace("n't", " not").replace("what's", "what is").replace("it's", "it is"
12                           .replace("'ve", " have").replace("i'm", "i am").replace("'re", " are")\
13                           .replace("he's", "he is").replace("she's", "she is").replace("'s", " own")\
14                           .replace("%", " percent ").replace("₹", " rupee ").replace("$", " dollar ")\
15                           .replace("€", " euro ").replace("'ll", " will")
16       x = re.sub(r"([0-9]+)000000", r"\1m", x)
17       x = re.sub(r"([0-9]+)000", r"\1k", x)
18
19
20       porter = PorterStemmer()
21       pattern = re.compile('\W')
22
23       if type(x) == type(''):
24           x = re.sub(pattern, ' ', x)
25
26
27       if type(x) == type(''):
28           x = porter.stem(x)
29           example1 = BeautifulSoup(x)
30           x = example1.get_text()
31
32
33       return x
34
```

- Function to Compute and get the features : With 2 parameters of Question 1 and Question 2

# 3.5 Advanced Feature Extraction (NLP and Fuzzy Features)

Definition:

- **Token**: You get a token by splitting sentence a space
- **Stop_Word** : stop words as per NLTK.
- **Word** : A token that is not a stop_word

Features:

- **cwc_min** : Ratio of common_word_count to min lenght of word count of Q1 and Q2
  cwc_min = common_word_count / (min(len(q1_words), len(q2_words))

- **cwc_max** : Ratio of common_word_count to max lenght of word count of Q1 and Q2
  cwc_max = common_word_count / (max(len(q1_words), len(q2_words))

- **csc_min** : Ratio of common_stop_count to min lenght of stop count of Q1 and Q2
  csc_min = common_stop_count / (min(len(q1_stops), len(q2_stops))

- **csc_max** : Ratio of common_stop_count to max lenght of stop count of Q1 and Q2
  csc_max = common_stop_count / (max(len(q1_stops), len(q2_stops))

- **ctc_min** : Ratio of common_token_count to min lenght of token count of Q1 and Q2
  ctc_min = common_token_count / (min(len(q1_tokens), len(q2_tokens))

- **ctc_max** : Ratio of common_token_count to max lenght of token count of Q1 and Q2
  ctc_max = common_token_count / (max(len(q1_tokens), len(q2_tokens))

- **last_word_eq** : Check if First word of both questions is equal or not
  last_word_eq = int(q1_tokens[-1] == q2_tokens[-1])

- **first_word_eq** : Check if First word of both questions is equal or not
  first_word_eq = int(q1_tokens[0] == q2_tokens[0])

- **abs_len_diff** : Abs. length difference
  abs_len_diff = abs(len(q1_tokens) - len(q2_tokens))

- **mean_len** : Average Token Length of both Questions
  mean_len = (len(q1_tokens) + len(q2_tokens))/2

- **fuzz_ratio** : https://github.com/seatgeek/fuzzywuzzy#usage (https://github.com/seatgeek/fuzzywuzzy#usage)
  http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/
  (http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/)

- **fuzz_partial_ratio** : https://github.com/seatgeek/fuzzywuzzy#usage
  (https://github.com/seatgeek/fuzzywuzzy#usage) http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-
  in-python/ (http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/)

- **token_sort_ratio** : https://github.com/seatgeek/fuzzywuzzy#usage
  (https://github.com/seatgeek/fuzzywuzzy#usage) http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-
  in-python/ (http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/)

- **token_set_ratio** : https://github.com/seatgeek/fuzzywuzzy#usage
  (https://github.com/seatgeek/fuzzywuzzy#usage) http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-
  in-python/ (http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/)

- **longest_substr_ratio** : Ratio of length longest common substring to min lenghth of token count of Q1 and Q2
  longest_substr_ratio = len(longest common substring) / (min(len(q1_tokens), len(q2_tokens))

```python
def get_token_features(q1, q2):
    token_features = [0.0]*10

    # Converting the Sentence into Tokens:
    q1_tokens = q1.split()
    q2_tokens = q2.split()

    if len(q1_tokens) == 0 or len(q2_tokens) == 0:
        return token_features
    # Get the non-stopwords in Questions
    q1_words = set([word for word in q1_tokens if word not in STOP_WORDS])
    q2_words = set([word for word in q2_tokens if word not in STOP_WORDS])

    #Get the stopwords in Questions
    q1_stops = set([word for word in q1_tokens if word in STOP_WORDS])
    q2_stops = set([word for word in q2_tokens if word in STOP_WORDS])

    # Get the common non-stopwords from Question pair
    common_word_count = len(q1_words.intersection(q2_words))

    # Get the common stopwords from Question pair
    common_stop_count = len(q1_stops.intersection(q2_stops))

    # Get the common Tokens from Question pair
    common_token_count = len(set(q1_tokens).intersection(set(q2_tokens)))


    token_features[0] = common_word_count / (min(len(q1_words), len(q2_words)) + SAFE_DIV)
    token_features[1] = common_word_count / (max(len(q1_words), len(q2_words)) + SAFE_DIV)
    token_features[2] = common_stop_count / (min(len(q1_stops), len(q2_stops)) + SAFE_DIV)
    token_features[3] = common_stop_count / (max(len(q1_stops), len(q2_stops)) + SAFE_DIV)
    token_features[4] = common_token_count / (min(len(q1_tokens), len(q2_tokens)) + SAFE_DIV)
    token_features[5] = common_token_count / (max(len(q1_tokens), len(q2_tokens)) + SAFE_DIV)
```

```
35        # Last word of both question is same or not
36        token_features[6] = int(q1_tokens[-1] == q2_tokens[-1])
37
38        # First word of both question is same or not
39        token_features[7] = int(q1_tokens[0] == q2_tokens[0])
40
41        token_features[8] = abs(len(q1_tokens) - len(q2_tokens))
42
43        #Average Token Length of both Questions
44        token_features[9] = (len(q1_tokens) + len(q2_tokens))/2
45        return token_features
46
47  # get the Longest Common sub string
48
49  def get_longest_substr_ratio(a, b):
50      strs = list(distance.lcsubstrings(a, b))
51      if len(strs) == 0:
52          return 0
53      else:
54          return len(strs[0]) / (min(len(a), len(b)) + 1)
55
56  def extract_features(df):
57      # preprocessing each question
58      df["question1"] = df["question1"].fillna("").apply(preprocess)
59      df["question2"] = df["question2"].fillna("").apply(preprocess)
60
61      print("token features...")
62
63      # Merging Features with dataset
64
65      token_features = df.apply(lambda x: get_token_features(x["question1"], x["question2"]), axis=1)
66
67      df["cwc_min"]       = list(map(lambda x: x[0], token_features))
68      df["cwc_max"]       = list(map(lambda x: x[1], token_features))
69      df["csc_min"]       = list(map(lambda x: x[2], token_features))
```

```python
70        df["csc_max"]        = list(map(lambda x: x[3], token_features))
71        df["ctc_min"]        = list(map(lambda x: x[4], token_features))
72        df["ctc_max"]        = list(map(lambda x: x[5], token_features))
73        df["last_word_eq"]   = list(map(lambda x: x[6], token_features))
74        df["first_word_eq"]  = list(map(lambda x: x[7], token_features))
75        df["abs_len_diff"]   = list(map(lambda x: x[8], token_features))
76        df["mean_len"]       = list(map(lambda x: x[9], token_features))
77
78        #Computing Fuzzy Features and Merging with Dataset
79
80        # do read this blog: http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/
81        # https://stackoverflow.com/questions/31806695/when-to-use-which-fuzz-function-to-compare-2-strings
82        # https://github.com/seatgeek/fuzzywuzzy
83        print("fuzzy features..")
84
85        df["token_set_ratio"]       = df.apply(lambda x: fuzz.token_set_ratio(x["question1"], x["question2"]
86        # The token sort approach involves tokenizing the string in question, sorting the tokens alphabetic
87        # then joining them back into a string We then compare the transformed strings with a simple ratio(
88        df["token_sort_ratio"]      = df.apply(lambda x: fuzz.token_sort_ratio(x["question1"], x["question2"
89        df["fuzz_ratio"]            = df.apply(lambda x: fuzz.QRatio(x["question1"], x["question2"]), axis=
90        df["fuzz_partial_ratio"]    = df.apply(lambda x: fuzz.partial_ratio(x["question1"], x["question2"])
91        df["longest_substr_ratio"]  = df.apply(lambda x: get_longest_substr_ratio(x["question1"], x["questi
92        return df
```

```
1  if os.path.isfile('nlp_features_train.csv'):
2      df = pd.read_csv("nlp_features_train.csv",encoding='latin-1')
3      df.fillna('')
4  else:
5      print("Extracting features for train:")
6      df = pd.read_csv("train.csv")
7      df = extract_features(df)
8      df.to_csv("nlp_features_train.csv", index=False)
9  df.head(2)
```

| | id | qid1 | qid2 | question1 | question2 | is_duplicate | cwc_min | cwc_max | csc_min | csc_max | ... | ctc_max |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **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 | 0.999980 | 0.833319 | 0.999983 | 0.999983 | ... | 0.785709 |
| **1** | 1 | 3 | 4 | what is the story of kohinoor koh i noor dia... | what would happen if the indian government sto... | 0 | 0.799984 | 0.399996 | 0.749981 | 0.599988 | ... | 0.466664 |

2 rows × 21 columns

## 3.5.1 Analysis of extracted features

### 3.5.1.1 Plotting Word clouds

- Creating Word Cloud of Duplicates and Non-Duplicates Question pairs
- We can observe the most frequent occuring words

```
1   df_duplicate = df[df['is_duplicate'] == 1]
2   dfp_nonduplicate = df[df['is_duplicate'] == 0]
3
4   # Converting 2d array of q1 and q2 and flatten the array: like {{1,2},{3,4}} to {1,2,3,4}
5   p = np.dstack([df_duplicate["question1"], df_duplicate["question2"]]).flatten()
6   n = np.dstack([dfp_nonduplicate["question1"], dfp_nonduplicate["question2"]]).flatten()
7
8   print ("Number of data points in class 1 (duplicate pairs) :",len(p))
9   print ("Number of data points in class 0 (non duplicate pairs) :",len(n))
10
11  #Saving the np array into a text file
12  np.savetxt('train_p.txt', p, delimiter=' ', fmt='%s')
13  np.savetxt('train_n.txt', n, delimiter=' ', fmt='%s')
```

```
Number of data points in class 1 (duplicate pairs) : 298526
Number of data points in class 0 (non duplicate pairs) : 510054
```
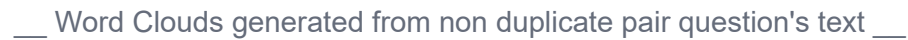
```
 1   # reading the text files and removing the Stop Words:
 2   d = path.dirname('.')
 3
 4   textp_w = open(path.join(d, 'train_p.txt')).read()
 5   textn_w = open(path.join(d, 'train_n.txt')).read()
 6   stopwords = set(STOPWORDS)
 7   stopwords.add("said")
 8   stopwords.add("br")
 9   stopwords.add(" ")
10   stopwords.remove("not")
11
12   stopwords.remove("no")
13   #stopwords.remove("good")
14   #stopwords.remove("Love")
15   stopwords.remove("like")
16   #stopwords.remove("best")
17   #stopwords.remove("!")
18   print ("Total number of words in duplicate pair questions :",len(textp_w))
19   print ("Total number of words in non duplicate pair questions :",len(textn_w))
```
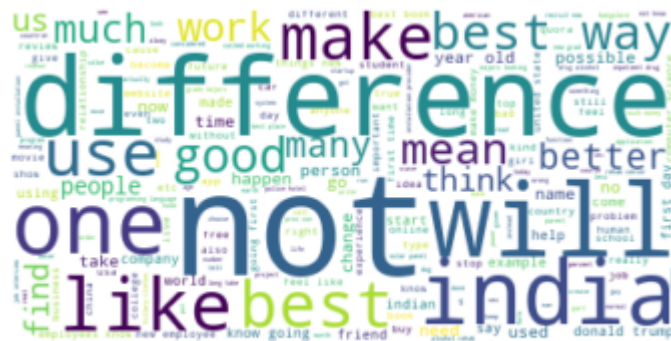
```
Total number of words in duplicate pair questions : 16109886
Total number of words in non duplicate pair questions : 33193130
```

__ Word Clouds generated from duplicate pair question's text __

```
1  wc = WordCloud(background_color="white", max_words=len(textp_w), stopwords=stopwords)
2  wc.generate(textp_w)
3  print ("Word Cloud for Duplicate Question pairs")
4  plt.imshow(wc, interpolation='bilinear')
5  plt.axis("off")
6  plt.show()
```

Word Cloud for Duplicate Question pairs



__ Word Clouds generated from non duplicate pair question's text __

```
1   wc = WordCloud(background_color="white", max_words=len(textn_w),stopwords=stopwords)
2   # generate word cloud
3   wc.generate(textn_w)
4   print ("Word Cloud for non-Duplicate Question pairs:")
5   plt.imshow(wc, interpolation='bilinear')
6   plt.axis("off")
7   plt.show()
```

Word Cloud for non-Duplicate Question pairs:



### 3.5.1.2 Pair plot of features ['ctc_min', 'cwc_min', 'csc_min', 'token_sort_ratio']

```
1  n = df.shape[0]
2  sns.pairplot(df[['ctc_min', 'cwc_min', 'csc_min', 'token_sort_ratio', 'is_duplicate']][0:n], hue='is_dup
3  plt.show()
```
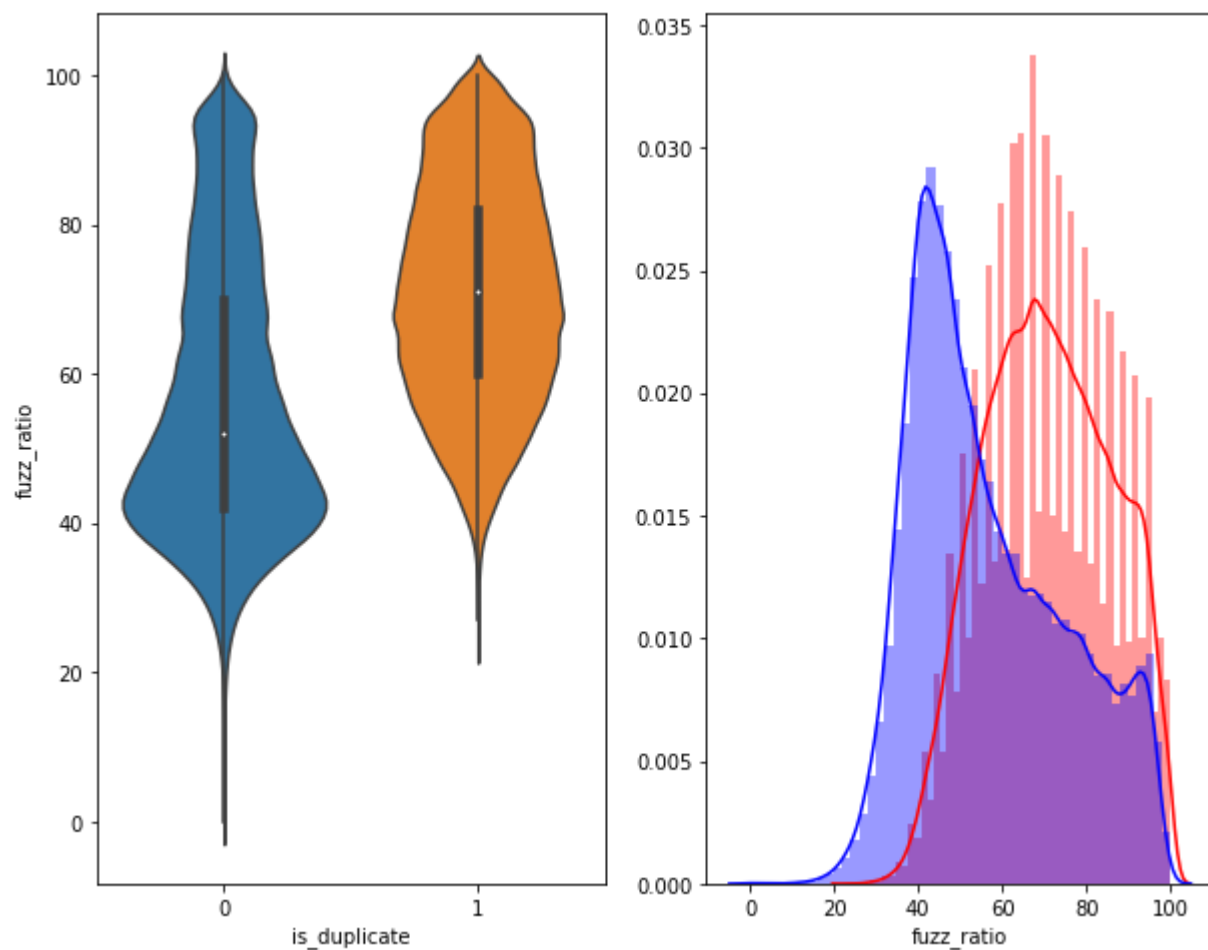
```
1   # Distribution of the token_sort_ratio
2   plt.figure(figsize=(10, 8))
3
4   plt.subplot(1,2,1)
5   sns.violinplot(x = 'is_duplicate', y = 'token_sort_ratio', data = df[0:] , )
6
7   plt.subplot(1,2,2)
8   sns.distplot(df[df['is_duplicate'] == 1.0]['token_sort_ratio'][0:] , label = "1", color = 'red')
9   sns.distplot(df[df['is_duplicate'] == 0.0]['token_sort_ratio'][0:] , label = "0" , color = 'blue' )
10  plt.show()
```

```
1  plt.figure(figsize=(10, 8))
2
3  plt.subplot(1,2,1)
4  sns.violinplot(x = 'is_duplicate', y = 'fuzz_ratio', data = df[0:] , )
5
6  plt.subplot(1,2,2)
7  sns.distplot(df[df['is_duplicate'] == 1.0]['fuzz_ratio'][0:] , label = "1", color = 'red')
8  sns.distplot(df[df['is_duplicate'] == 0.0]['fuzz_ratio'][0:] , label = "0" , color = 'blue' )
9  plt.show()
```
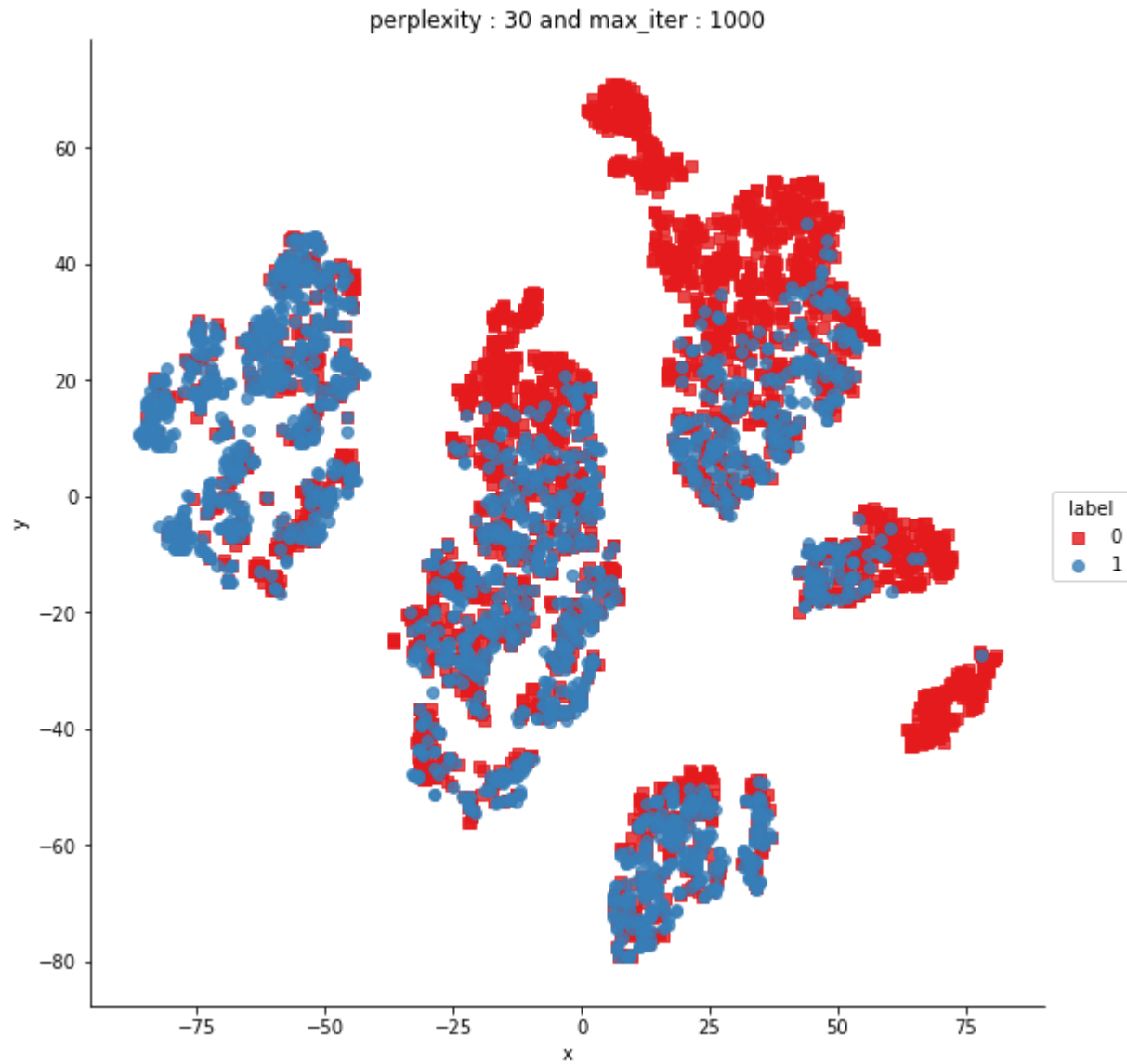
### 3.5.2 Visualization

```python
# Using TSNE for Dimentionality reduction for 15 Features(Generated after cleaning the data) to 3 dimen

from sklearn.preprocessing import MinMaxScaler

dfp_subsampled = df[0:5000]
X = MinMaxScaler().fit_transform(dfp_subsampled[['cwc_min', 'cwc_max', 'csc_min', 'csc_max' , 'ctc_min'
y = dfp_subsampled['is_duplicate'].values
```

```python
1  tsne2d = TSNE(
2      n_components=2,
3      init='random', # pca
4      random_state=101,
5      method='barnes_hut',
6      n_iter=1000,
7      verbose=2,
8      angle=0.5
9  ).fit_transform(X)
```

```
[t-SNE] Computing 91 nearest neighbors...
[t-SNE] Indexed 5000 samples in 0.011s...
[t-SNE] Computed neighbors for 5000 samples in 0.912s...
[t-SNE] Computed conditional probabilities for sample 1000 / 5000
[t-SNE] Computed conditional probabilities for sample 2000 / 5000
[t-SNE] Computed conditional probabilities for sample 3000 / 5000
[t-SNE] Computed conditional probabilities for sample 4000 / 5000
[t-SNE] Computed conditional probabilities for sample 5000 / 5000
[t-SNE] Mean sigma: 0.116557
[t-SNE] Computed conditional probabilities in 0.433s
[t-SNE] Iteration 50: error = 80.9244080, gradient norm = 0.0428133 (50 iterations in 13.099s)
[t-SNE] Iteration 100: error = 70.3858795, gradient norm = 0.0100968 (50 iterations in 9.067s)
[t-SNE] Iteration 150: error = 68.6138382, gradient norm = 0.0058392 (50 iterations in 9.602s)
[t-SNE] Iteration 200: error = 67.7700119, gradient norm = 0.0036596 (50 iterations in 9.121s)
[t-SNE] Iteration 250: error = 67.2725067, gradient norm = 0.0034962 (50 iterations in 11.305s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 67.272507
[t-SNE] Iteration 300: error = 1.7737305, gradient norm = 0.0011918 (50 iterations in 8.289s)
[t-SNE] Iteration 350: error = 1.3720417, gradient norm = 0.0004822 (50 iterations in 10.526s)
[t-SNE] Iteration 400: error = 1.2039998, gradient norm = 0.0002768 (50 iterations in 9.600s)
[t-SNE] Iteration 450: error = 1.1133438, gradient norm = 0.0001881 (50 iterations in 11.827s)
[t-SNE] Iteration 500: error = 1.0579143, gradient norm = 0.0001434 (50 iterations in 8.941s)
[t-SNE] Iteration 550: error = 1.0221983, gradient norm = 0.0001164 (50 iterations in 11.092s)
[t-SNE] Iteration 600: error = 0.9987167, gradient norm = 0.0001039 (50 iterations in 11.467s)
[t-SNE] Iteration 650: error = 0.9831534, gradient norm = 0.0000938 (50 iterations in 11.799s)
[t-SNE] Iteration 700: error = 0.9722011, gradient norm = 0.0000858 (50 iterations in 12.028s)
[t-SNE] Iteration 750: error = 0.9643636, gradient norm = 0.0000799 (50 iterations in 12.120s)
[t-SNE] Iteration 800: error = 0.9584482, gradient norm = 0.0000785 (50 iterations in 11.867s)
[t-SNE] Iteration 850: error = 0.9538348, gradient norm = 0.0000739 (50 iterations in 11.461s)
[t-SNE] Iteration 900: error = 0.9496906, gradient norm = 0.0000712 (50 iterations in 11.023s)
[t-SNE] Iteration 950: error = 0.9463405, gradient norm = 0.0000673 (50 iterations in 11.755s)
```

```
[t-SNE] Iteration 1000: error = 0.9432716, gradient norm = 0.0000662 (50 iterations in 11.493s)
[t-SNE] Error after 1000 iterations: 0.943272
```

```
1  df = pd.DataFrame({'x':tsne2d[:,0], 'y':tsne2d[:,1] ,'label':y})
2
3  # draw the plot in appropriate place in the grid
4  sns.lmplot(data=df, x='x', y='y', hue='label', fit_reg=False, size=8,palette="Set1",markers=['s','o'])
5  plt.title("perplexity : {} and max_iter : {}".format(30, 1000))
6  plt.show()
```

perplexity : 30 and max_iter : 1000

```python
from sklearn.manifold import TSNE
tsne3d = TSNE(
    n_components=3,
    init='random', # pca
    random_state=101,
    method='barnes_hut',
    n_iter=1000,
    verbose=2,
    angle=0.5
).fit_transform(X)
```

```
[t-SNE] Computing 91 nearest neighbors...
[t-SNE] Indexed 5000 samples in 0.010s...
[t-SNE] Computed neighbors for 5000 samples in 0.935s...
[t-SNE] Computed conditional probabilities for sample 1000 / 5000
[t-SNE] Computed conditional probabilities for sample 2000 / 5000
[t-SNE] Computed conditional probabilities for sample 3000 / 5000
[t-SNE] Computed conditional probabilities for sample 4000 / 5000
[t-SNE] Computed conditional probabilities for sample 5000 / 5000
[t-SNE] Mean sigma: 0.116557
[t-SNE] Computed conditional probabilities in 0.363s
[t-SNE] Iteration 50: error = 77.7944183, gradient norm = 0.1014017 (50 iterations in 34.931s)
[t-SNE] Iteration 100: error = 69.2682266, gradient norm = 0.0248657 (50 iterations in 15.147s)
[t-SNE] Iteration 150: error = 67.7877655, gradient norm = 0.0150941 (50 iterations in 13.761s)
[t-SNE] Iteration 200: error = 67.1991119, gradient norm = 0.0126559 (50 iterations in 13.425s)
[t-SNE] Iteration 250: error = 66.8560715, gradient norm = 0.0074975 (50 iterations in 12.904s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 66.856071
[t-SNE] Iteration 300: error = 1.2356015, gradient norm = 0.0007033 (50 iterations in 13.302s)
[t-SNE] Iteration 350: error = 0.9948602, gradient norm = 0.0001997 (50 iterations in 18.898s)
[t-SNE] Iteration 400: error = 0.9168936, gradient norm = 0.0001430 (50 iterations in 13.397s)
[t-SNE] Iteration 450: error = 0.8863022, gradient norm = 0.0000975 (50 iterations in 16.379s)
[t-SNE] Iteration 500: error = 0.8681002, gradient norm = 0.0000854 (50 iterations in 17.791s)
[t-SNE] Iteration 550: error = 0.8564141, gradient norm = 0.0000694 (50 iterations in 17.060s)
[t-SNE] Iteration 600: error = 0.8470711, gradient norm = 0.0000640 (50 iterations in 15.454s)
[t-SNE] Iteration 650: error = 0.8389117, gradient norm = 0.0000561 (50 iterations in 17.562s)
[t-SNE] Iteration 700: error = 0.8325295, gradient norm = 0.0000529 (50 iterations in 13.443s)
[t-SNE] Iteration 750: error = 0.8268463, gradient norm = 0.0000528 (50 iterations in 17.981s)
[t-SNE] Iteration 800: error = 0.8219477, gradient norm = 0.0000477 (50 iterations in 17.448s)
[t-SNE] Iteration 850: error = 0.8180174, gradient norm = 0.0000490 (50 iterations in 18.376s)
[t-SNE] Iteration 900: error = 0.8150476, gradient norm = 0.0000456 (50 iterations in 17.778s)
```

```
[t-SNE] Iteration 950: error = 0.8122067, gradient norm = 0.0000472 (50 iterations in 16.983s)
[t-SNE] Iteration 1000: error = 0.8095787, gradient norm = 0.0000489 (50 iterations in 18.581s)
[t-SNE] Error after 1000 iterations: 0.809579
```

```
 1  trace1 = go.Scatter3d(
 2      x=tsne3d[:,0],
 3      y=tsne3d[:,1],
 4      z=tsne3d[:,2],
 5      mode='markers',
 6      marker=dict(
 7          sizemode='diameter',
 8          color = y,
 9          colorscale = 'Portland',
10          colorbar = dict(title = 'duplicate'),
11          line=dict(color='rgb(255, 255, 255)'),
12          opacity=0.75
13      )
14  )
15
16  data=[trace1]
17  layout=dict(height=800, width=800, title='3d embedding with engineered features')
18  fig=dict(data=data, layout=layout)
19  py.iplot(fig, filename='3DBubble')
```