

# # Importing Data and Preprocessing (Same in All notebooks)

In [1]:

```
1 import numpy as np
2 import pandas as pd
3 import string
```

In [2]:

```
1 data = pd.read_csv('/users/rohanchitte/downloads/Dataset_lyrics.csv_lyrics.csv')
```

In [3]:

```
1 filtered = data[data['lyrics'].notnull()]
2 filtered
```

Out[3]:

	index	song	year	artist	genre	lyrics
0	0	ego-remix	2009	beyonce-knowles	Pop	Oh baby, how you doing?\nYou know I'm gonna cu...
1	1	then-tell-me	2009	beyonce-knowles	Pop	playin' everything so easy,\nit's like you see...
2	2	honesty	2009	beyonce-knowles	Pop	If you search\nFor tenderness\nIt isn't hard t...
3	3	you-are-my-rock	2009	beyonce-knowles	Pop	Oh oh oh I, oh oh oh I\n[Verse 1:]\nIf I wrote...
4	4	black-culture	2009	beyonce-knowles	Pop	Party the people, the people the party it's po...
...	...	...	...	...	...	...
362232	362232	who-am-i-drinking-tonight	2012	edens-edge	Country	I gotta say\nBoy, after only just a couple of ...
362233	362233	liar	2012	edens-edge	Country	I helped you find her diamond ring\nYou made m...
362234	362234	last-supper	2012	edens-edge	Country	Look at the couple in the corner booth\nLooks ...
362235	362235	christ-alone-live-in-studio	2012	edens-edge	Country	When I fly off this mortal earth\nAnd I'm meas...
362236	362236	amen	2012	edens-edge	Country	I heard from a friend of a friend of a friend ...

266557 rows × 6 columns

In [4]:

```

1 import nltk
2 from nltk.corpus import stopwords
3
4 cleaned = filtered.copy()
5
6 # Remove punctuation
7 cleaned['lyrics'] = cleaned['lyrics'].str.replace("[-\?.,\/#!\$%\^&\*;:{ }=\_~()]",
8
9 # Remove song-related identifiers like [Chorus] or [Verse]
10 cleaned['lyrics'] = cleaned['lyrics'].str.replace("\[(.*?)\]", ' ')
11 cleaned['lyrics'] = cleaned['lyrics'].str.replace("' | '", ' ')
12 cleaned['lyrics'] = cleaned['lyrics'].str.replace('x[0-9]+', ' ')
13
14 # Remove all songs without lyrics (e.g. instrumental pieces)
15 cleaned = cleaned[cleaned['lyrics'].str.strip().str.lower() != 'instrumental']
16
17 # Remove any songs with corrupted/non-ASCII characters, unavailable lyrics
18 cleaned = cleaned[~cleaned['lyrics'].str.contains(r'[\x00-\x7F]+')]
19 cleaned = cleaned[cleaned['lyrics'].str.strip() != '']
20 cleaned = cleaned[cleaned['genre'].str.lower() != 'not available']
21
22 #Selecting Pop, Rock, Country, Jazz
23 cleaned = cleaned.loc[(cleaned['genre'] == 'Pop') |
24                       (cleaned['genre'] == 'Country') |
25                       (cleaned['genre'] == 'Rock') |
26                       (cleaned['genre'] == 'Hip-Hop') |
27                       (cleaned['genre'] == 'Jazz') ]
28 cleaned.reset_index(inplace = True)
29
30 cleaned
31 print(len(cleaned))
32
33 from nltk.corpus import stopwords
34 stop = stopwords.words('english')
35 #removing stop words from lyrics
36
37 cleaned['lyrics'] = cleaned['lyrics'].apply(lambda x: ' '.join([word for word in
38
39 #lemmatizing lyrics
40 import nltk
41
42 w_tokenizer = nltk.tokenize.WhitespaceTokenizer()
43 lemmatizer = nltk.stem.WordNetLemmatizer()
44
45 def lemmatize_text(text, flg_lemm=True):
46     #Convert string to list (tokenize)
47     lst_text = text.split()
48
49     ## Lemmatisation (convert the word into root word)
50     if flg_lemm == True:
51         lem = nltk.stem.wordnet.WordNetLemmatizer()
52         lst_text = [lem.lemmatize(word) for word in lst_text]
53
54     ## back to string from list
55     text = " ".join(lst_text)
56     return text
57
58 #cleaned["lyrics"] = cleaned["lyrics"].apply(lemmatize_text)
59

```

```

60 cleaned["lyrics"] = cleaned["lyrics"].apply(lambda x: lemmatize_text(x))
61
62 df = cleaned.drop(labels=["level_0", "index", "song", "year", "artist"], axis=1)

```

185493

## 1 # Splitting Data into training and test set

In [5]:

```
1 from sklearn.model_selection import train_test_split
```

In [6]:

```
1 df_train, df_test = train_test_split(df, test_size=0.33, random_state=42)
```

In [26]:

```

1 df_train.reset_index()
2 df_test.reset_index()

```

Out[26]:

	index	genre	lyrics
0	35835	Jazz	I dance ask I dance ask I dance madame My hear...
1	2538	Hip-Hop	Sonic boom head dread cause he's tread Upon Fl...
2	63159	Rock	If I could turn page In time I'd rearrange Jus...
3	6483	Rock	record stop stop skipping equipped stor ear fu...
4	15496	Hip-Hop	Hey yeah ya know I like playersNo Diggity No d...
...	...	...	...
61208	10254	Hip-Hop	We're never done found place belong Don't stan...
61209	31630	Country	It's fake hoax nowhere road one go anywhere an...
61210	107267	Rock	I've spent much time throwing rock window That...
61211	67806	Rock	You're lookin fine long time I still remember ...
61212	23935	Pop	I I get creepin feelin' That might start belie...

61213 rows × 3 columns

In [27]:

```

1 #train_test split
2 x_tr, y_tr = df_train['lyrics'].values, df_train['genre'].values
3 x_val, y_val = df_test['lyrics'].values, df_test['genre'].values

```

## 1 # Function to calculcate max length of the sequence in the corpus

In [28]:

```
1 def get_max_length(df_train):
2     """
3     get max token counts from train data,
4     so we use this number as fixed length input to LSTM cell
5     """
6     max_length = 0
7     for row in df_train['lyrics']:
8         if len(row.split(" ")) > max_length:
9             max_length = len(row.split(" "))
10    return max_length
11
```

In [30]:

```
1 maximumlen = get_max_length(df_train)
2 maximumlen
```

Out[30]:

3666

```
1 # One hot encoding genres
```

In [31]:

```
1 def genre_encode(genre):
2     """
3     return one hot encoding for Y value
4     """
5     if genre == 'Pop':
6         return [1,0,0,0,0]
7     elif genre == 'Rock':
8         return [0,1,0,0,0]
9     elif genre == 'Country':
10        return [0,0,1,0,0]
11    elif genre == 'Hip-Hop':
12        return [0,0,0,1,0]
13    else:
14        return [0,0,0,0,1]
```

In [32]:

```
1 genres = df_train['genre'].tolist()
2 y_tr = [genre_encode(genre) for genre in genres]
```

In [33]:

```
1 genres = df_test['genre'].tolist()
2 y_val = [genre_encode(genre) for genre in genres]
```

```
1 # Tokenization and Padding of the sequences to
  make their length same
```

In [34]:

```

1 from keras.preprocessing.text import Tokenizer
2 from keras.preprocessing.sequence import pad_sequences
3
4 #Tokenize the sentences
5 tokenizer = Tokenizer()
6
7 #preparing vocabulary
8 tokenizer.fit_on_texts(list(x_tr))
9
10 #converting text into integer sequences
11 x_tr_seq = tokenizer.texts_to_sequences(x_tr)
12 x_val_seq = tokenizer.texts_to_sequences(x_val)
13
14 #padding to prepare sequences of same length
15 x_tr_seq = pad_sequences(x_tr_seq, maxlen=maximumlen)
16 x_val_seq = pad_sequences(x_val_seq, maxlen=maximumlen)

```

## 1 # Loading Pretrained Glove Word Embedding

In [35]:

```

1 # load the whole embedding into memory
2 embeddings_index = dict()
3 f = open('/users/rohanchitte/glove.6B.300d.txt')
4
5 for line in f:
6     values = line.split()
7     word = values[0]
8     coefs = np.asarray(values[1:], dtype='float32')
9     embeddings_index[word] = coefs
10
11 f.close()
12 print('Loaded %s word vectors.' % len(embeddings_index))

```

Loaded 400000 word vectors.

In [36]:

```

1 size_of_vocabulary=len(tokenizer.word_index) + 1 #+1 for padding
2 print(size_of_vocabulary)

```

220815

## 1 # Creating a weight matrix for words in training docs

In [37]:

```

1
2 embedding_matrix = np.zeros((size_of_vocabulary, 300))
3
4 for word, i in tokenizer.word_index.items():
5     embedding_vector = embeddings_index.get(word)
6     if embedding_vector is not None:
7         embedding_matrix[i] = embedding_vector

```

## 1 # LSTM Model

In [38]:

```

1 #deep learning library
2 from keras.models import *
3 from keras.layers import *
4 from keras.callbacks import *

```

In [39]:

```

1 model=Sequential()
2
3 #embedding layer
4 model.add(Embedding(size_of_vocabulary,300,weights=[embedding_matrix],input_length=
5
6 #lstm layer
7 model.add(LSTM(128,return_sequences=True,dropout=0.2))
8
9 #Global Maxpooling
10 model.add(GlobalMaxPooling1D())
11
12 #Dense Layer
13 model.add(Dense(64,activation='relu'))
14 model.add(Dense(5,activation='softmax'))
15
16 #Add loss function, metrics, optimizer
17 model.compile(optimizer='adam', loss='categorical_crossentropy',metrics=["accuracy"])
18
19 #Adding callbacks
20 #es = EarlyStopping(monitor='val_loss', mode='min', verbose=1,patience=3)
21
22 #mc=ModelCheckpoint('best_model.h5', monitor='val_acc', mode='max', save_best_only=True)
23
24 #Print summary of model
25 print(model.summary())

```

Model: "sequential\_1"

Layer (type)	Output Shape	Param #
=====		
embedding_1 (Embedding)	(None, 3666, 300)	66244500
=====		
lstm_1 (LSTM)	(None, 3666, 128)	219648
=====		
global_max_pooling1d_1 (GlobalMaxPooling1D)	(None, 128)	0
=====		
dense_2 (Dense)	(None, 64)	8256
=====		
dense_3 (Dense)	(None, 5)	325
=====		
Total params: 66,472,729		
Trainable params: 228,229		
Non-trainable params: 66,244,500		
=====		
None		

In [40]:

```
1 history = model.fit( np.array(x_tr_seq), np.array(y_tr), batch_size=128, epochs=
```

Epoch 1/5

971/971 [=====] - 57327s 59s/step - loss: 0.9

087 - accuracy: 0.6608 - val\_loss: 0.8444 - val\_accuracy: 0.6809

Epoch 2/5

971/971 [=====] - 54529s 56s/step - loss: 0.8

149 - accuracy: 0.6924 - val\_loss: 0.8185 - val\_accuracy: 0.6884

Epoch 3/5

971/971 [=====] - 74745s 77s/step - loss: 0.7

791 - accuracy: 0.7066 - val\_loss: 0.7951 - val\_accuracy: 0.7013

Epoch 4/5

971/971 [=====] - 60208s 62s/step - loss: 0.7

513 - accuracy: 0.7172 - val\_loss: 0.8008 - val\_accuracy: 0.7010

Epoch 5/5

971/971 [=====] - 60168s 62s/step - loss: 0.7

281 - accuracy: 0.7261 - val\_loss: 0.7911 - val\_accuracy: 0.7045

In [41]:

```
1 _,val_acc = model.evaluate(np.array(x_val_seq),np.array(y_val))
2 print(val_acc)
```

1913/1913 [=====] - 11597s 6s/step - loss: 0.

7911 - accuracy: 0.7045

0.7045398950576782

In [42]:

```
1 # serialize weights to HDF5
2 model.save_weights("lyrics-5-categories-model-glove.h5")
3 print("Saved model to disk")
```

Saved model to disk