

1

Importing Data (Same in All notebooks)

In [1]:

```
1 ## Importing library
2 import numpy as np
3 import pandas as pd
4 np.random.seed(100)
5
6 data = pd.read_csv('/users/rohanchitte/downloads/Dataset_lyrics.csv_lyrics.csv')
```

In [2]:

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

Out[2]:

	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

1

Data Pre-processing (Same in All notebooks)

In [3]:

```

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

```

30/09/2021, 14:27Text data augmentation-Copy2- BERT with data augmented-Copy1 - Jupyter Notebook

```
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

Data Visualization - Histogram

In [4]:

```
1 df
```

Out[4]:

	genre	lyrics
0	Pop	Oh baby You know I'm gonna cut right chase Som...
1	Pop	playin everything easy like seem sure still wa...
2	Pop	If search For tenderness It hard find You love...
3	Pop	Oh oh oh I oh oh oh I If I wrote book stand Th...
4	Pop	Party people people party popping sitting arou...
...
185488	Country	I gotta say Boy couple date You're hand outrig...
185489	Country	I helped find diamond ring You made try everyt...
185490	Country	Look couple corner booth Looks lot like She's ...
185491	Country	When I fly mortal earth And I'm measured depth...
185492	Country	I heard friend friend friend You finally got r...

185493 rows × 2 columns

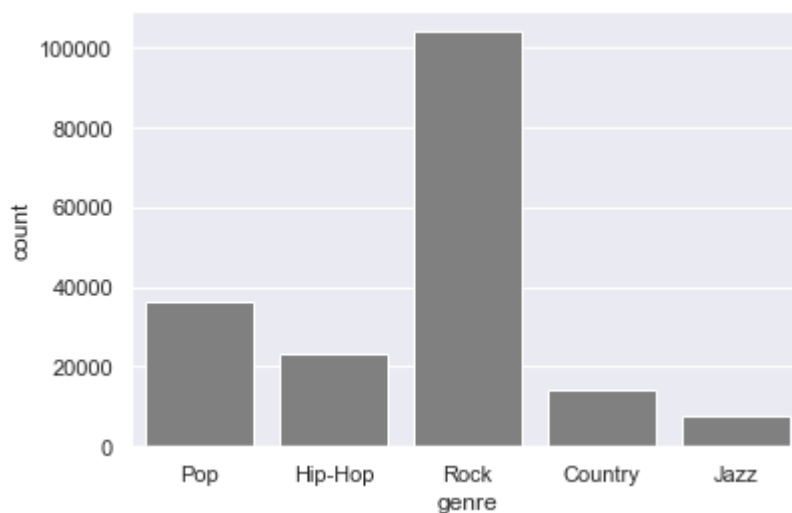
In [5]:

```
1 import seaborn as sns
2 sns.set()
3
4 sns.countplot(df['genre'], color='gray')
```

/opt/anaconda3/lib/python3.8/site-packages/seaborn/_decorators.py:36:
FutureWarning: Pass the following variable as a keyword arg: x. From v
ersion 0.12, the only valid positional argument will be `data`, and pa
ssing other arguments without an explicit keyword will result in an er
ror or misinterpretation.
warnings.warn(

Out[5]:

<AxesSubplot:xlabel='genre', ylabel='count'>



In [6]:

```
1 from sklearn.preprocessing import LabelEncoder
2 Y = df["genre"]
3 Y = LabelEncoder().fit_transform(Y)
```

In [7]:

```
1 df['Y'] = Y.tolist()
```

In [8]:

```
1 df["Y"]
```

Out[8]:

```
0      3
1      3
2      3
3      3
4      3
..
185488  0
185489  0
185490  0
185491  0
185492  0
```

Name: Y, Length: 185493, dtype: int64

```
1 # Splitting Data into Train and Test Set
```

In [9]:

```
1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(df["lyrics"], df["Y"], test_s
```

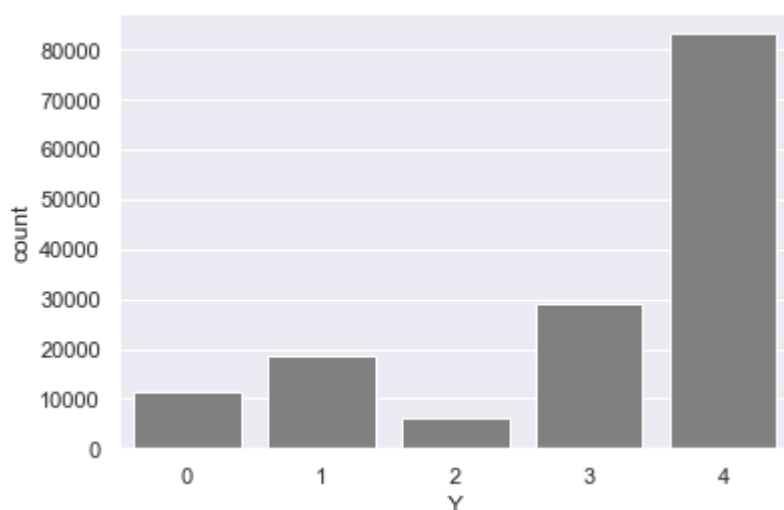
In [10]:

```
1 #Visualizing Y - Genres of training set
2 import seaborn as sns
3 sns.set()
4
5 sns.countplot(y_train, color='gray')
```

/opt/anaconda3/lib/python3.8/site-packages/seaborn/_decorators.py:36:
FutureWarning: Pass the following variable as a keyword arg: x. From v
ersion 0.12, the only valid positional argument will be `data`, and pa
ssing other arguments without an explicit keyword will result in an er
ror or misinterpretation.
warnings.warn(

Out[10]:

<AxesSubplot:xlabel='Y', ylabel='count'>



```
1 # Data Augmentation
```

In [11]:

```
1 #Creating artificial data to create more training data for Y Genre: Jazz (2)
2 import nlpaug.augmenter.char as nac
3 import nlpaug.augmenter.word as naw
4 import nlpaug.augmenter.sentence as nas
5 import nlpaug.flow as naflc
6
7 aug = naw.ContextualWordEmbsAug(model_path='bert-base-uncased', action="substitu
8
9 augmented_sentences=[]
10 augmented_sentences_labels=[]
11 jazz_index = []
12 for i in X_train.index:
13     if y_train[i]==2:
14         jazz_index.append(i)
15         temps=aug.augment(X_train[i],n=2)
16         for sent in temps:
17             augmented_sentences.append(sent)
18             augmented_sentences_labels.append(2)
19
20 X_train=X_train.append(pd.Series(augmented_sentences),ignore_index=True)
21 y_train=y_train.append(pd.Series(augmented_sentences_labels),ignore_index=True)
22
23
24 print(X_train.shape)
25 print(y_train.shape)
```

(160368,)

(160368,)

In []:

```
1 # Increase in training set from 148395 lyrics to 160368. So, Nearly 12000 ariti
```

In [12]:

```

1 #Visualizing the increase in data using histogram
2 import seaborn as sns
3 sns.set()
4
5 sns.countplot(y_train, color='gray')

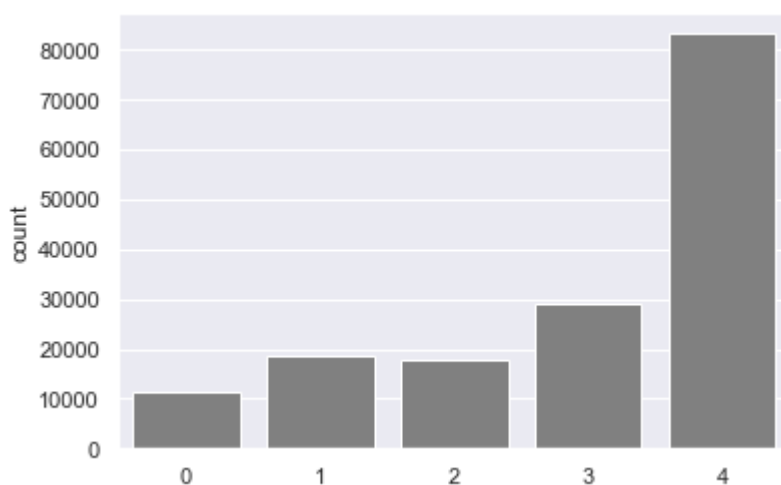
```

/opt/anaconda3/lib/python3.8/site-packages/seaborn/_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
warnings.warn(
```

Out[12]:

```
<AxesSubplot:ylabel='count'>
```



```
1 # Bert Model
```

In [13]:

```

1 import tensorflow as tf
2 import tensorflow_hub as hub
3 import tensorflow_text as text
4 from official.nlp import optimization # to create AdamW optimizer
5 from official.nlp import bert
6 from tensorflow import keras

```

/opt/anaconda3/lib/python3.8/site-packages/tensorflow_addons/utils/ensure_tf_install.py:53: UserWarning: Tensorflow Addons supports using Python ops for all Tensorflow versions above or equal to 2.3.0 and strictly below 2.6.0 (nightly versions are not supported).

The versions of TensorFlow you are currently using is 2.6.0 and is not supported.

Some things might work, some things might not.

If you were to encounter a bug, do not file an issue.

If you want to make sure you're using a tested and supported configuration, either change the TensorFlow version or the TensorFlow Addons's version.

You can find the compatibility matrix in TensorFlow Addon's readme:

<https://github.com/tensorflow/addons> (<https://github.com/tensorflow/addons>)

```
warnings.warn(
```


In [14]:

```

1 import spacy
2 nlp = spacy.load("en_core_web_sm")
3
4 def remove_non_ascii(text):
5     doc = nlp(text)
6     to_return = " ".join([str(token) for token in doc if token.is_ascii])
7     return to_return
8
9 X_train = X_train.apply(remove_non_ascii)

```

In [15]:

```

1 X_test = X_test.apply(remove_non_ascii)

```

In [16]:

```

### Selecting BERT encoder having transformer layers(L) = 4, ####dimension of o/p = 3
f_bert_encoder = 'https://tfhub.dev/tensorflow/small_bert/bert_en_uncased_L-4_H-512_A-8/1'
3
###choosing pre-processor that is compatible with BERT encoder ####selected
f_bert_pre_process = 'https://tfhub.dev/tensorflow/bert_en_uncased_preprocess/3'
4
def build_classifier_model():
    text_input = tf.keras.layers.Input(shape=(), dtype=tf.string, name='text')
    preprocessing_layer = hub.KerasLayer(f_bert_pre_process, name='preprocessing')
    encoder_inputs = preprocessing_layer(text_input)
    encoder = hub.KerasLayer(f_bert_encoder, trainable=True, name='BERT_encoder')
    outputs = encoder(encoder_inputs)
    net = outputs['pooled_output']
    net = tf.keras.layers.Dropout(0.1)(net)
    net = tf.keras.layers.Dense(5, activation='softmax', name='classifier')(net)
    #net = tf.keras.layers.Dense(1, activation=None, name='classifier')(net)
    return tf.keras.Model(text_input, net)
17
classifier_model = build_classifier_model()

```

```

INFO:absl:Using /var/folders/fl/kwcrn5_93n55xjv_rvr1d1080000gn/T/tfhub
_modules to cache modules.
INFO:absl:Downloading TF-Hub Module 'https://tfhub.dev/tensorflow/bert
_en_uncased_preprocess/3'.
INFO:absl:Downloaded https://tfhub.dev/tensorflow/bert_en_uncased_prep
rocess/3, (https://tfhub.dev/tensorflow/bert_en_uncased_preprocess/3,)
Total size: 1.96MB
INFO:absl:Downloaded TF-Hub Module 'https://tfhub.dev/tensorflow/bert
_en_uncased_preprocess/3'.
INFO:absl:Downloading TF-Hub Module 'https://tfhub.dev/tensorflow/smal
l_bert/bert_en_uncased_L-4_H-512_A-8/1'.
INFO:absl:Downloading https://tfhub.dev/tensorflow/small_bert/bert_en
_uncased_L-4_H-512_A-8/1: (https://tfhub.dev/tensorflow/small_bert/bert
_en_uncased_L-4_H-512_A-8/1:) 70.00MB
INFO:absl:Downloaded https://tfhub.dev/tensorflow/small_bert/bert_en_u
ncased_L-4_H-512_A-8/1, (https://tfhub.dev/tensorflow/small_bert/bert
_en_uncased_L-4_H-512_A-8/1,) Total size: 115.55MB
INFO:absl:Downloaded TF-Hub Module 'https://tfhub.dev/tensorflow/small
_bert/bert_en_uncased_L-4_H-512_A-8/1'.

```

In [17]:

```
1 loss = tf.keras.losses.CategoricalCrossentropy(from_logits=True)
2 metrics = tf.metrics.CategoricalAccuracy()
3 epochs = 7
4 steps_per_epoch = 11370
5 num_train_steps = steps_per_epoch * epochs
6 num_warmup_steps = int(0.1*num_train_steps)
7 init_lr = 3e-5
8 optimizer = optimization.create_optimizer(init_lr=init_lr, num_train_steps=num_t
9 classifier_model.compile(optimizer=optimizer,
10                           loss=loss,
11                           metrics=metrics)
```

INFO:absl:using Adamw optimizer

INFO:absl:gradient_clip_norm=1.000000

In [20]:

```
1 from sklearn.preprocessing import LabelBinarizer
2
3 def get_encoded_labels(topic_clusters):
4     encoder = LabelBinarizer()
5     encoded_labels = encoder.fit_transform(topic_clusters)
6     return encoded_labels
7 #http://localhost:8888/notebooks/Text%20data%20augmentation-Copy2-%20BERT%20with
8 labels = get_encoded_labels(y_train)
9
10 labelsval = get_encoded_labels(y_test)
```

```
1 # Fitting the model
```

In []:

```
1 history = classifier_model.fit(X_train, labels, epochs=epochs, verbose=1, validation
```

Epoch 1/7

```
/opt/anaconda3/lib/python3.8/site-packages/keras/backend.py:4846: User
Warning: "`categorical_crossentropy` received `from_logits=True`, but
the `output` argument was produced by a sigmoid or softmax activation
and thus does not represent logits. Was this intended?"
```

```
warnings.warn(
```

```
5012/5012 [=====] - 60383s 12s/step - loss:
1.0618 - categorical_accuracy: 0.5908 - val_loss: 0.8502 - val_categor
ical_accuracy: 0.6767
```

Epoch 2/7

```
5012/5012 [=====] - 68383s 14s/step - loss:
0.8581 - categorical_accuracy: 0.6745 - val_loss: 0.8028 - val_categor
ical_accuracy: 0.6989
```

Epoch 3/7

```
5012/5012 [=====] - 65567s 13s/step - loss:
0.7252 - categorical_accuracy: 0.7272 - val_loss: 0.8014 - val_categor
ical_accuracy: 0.7076
```

Epoch 4/7

```
5012/5012 [=====] - 65542s 13s/step - loss:
0.6106 - categorical_accuracy: 0.7691 - val_loss: 0.8198 - val_categor
ical_accuracy: 0.7158
```

Epoch 5/7

```
5012/5012 [=====] - 64529s 13s/step - loss:
0.5200 - categorical_accuracy: 0.8039 - val_loss: 0.8904 - val_categor
ical_accuracy: 0.7136
```

Epoch 6/7

```
5012/5012 [=====] - 68386s 14s/step - loss:
0.4462 - categorical_accuracy: 0.8327 - val_loss: 0.9507 - val_categor
ical_accuracy: 0.6939
```

Epoch 7/7

```
4152/5012 [=====>.....] - ETA: 2:55:01 - loss: 0.38
06 - categorical_accuracy: 0.8585
```

```
1 # Training interrupted as the validation accuracy was
not improving. Best accuracy of 71.58 % was
observed after training the model for 5 epochs.
```