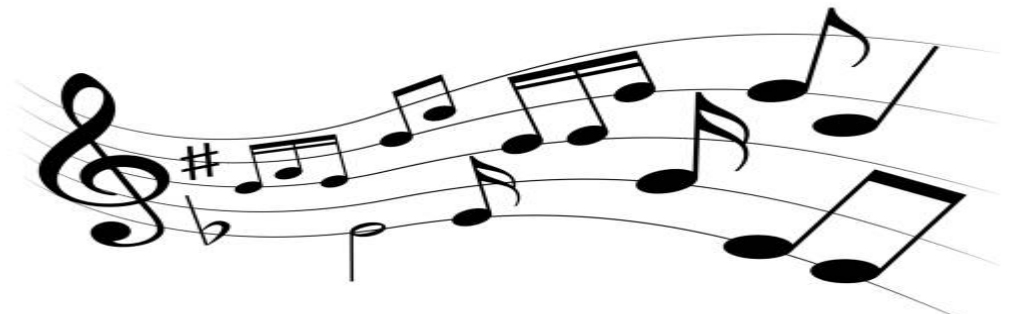


# Text Classification

Guess Artist from Lyrics



# 1- Get Data: LyricsGenius



- LyricsGenius: a Python client for the Genius.com API
- lyricsgenius provides a simple interface to the song, artist, and lyrics data stored on Genius.com.
- **Installation:** `pip install lyricsgenius`
- **Usage:** Import the package and initiate Genius:  

```
import lyricsgenius  
genius = lyricsgenius.Genius(token)
```

# API Client



← → ↺ 🔒 genius.com/api-clients#

Documentation


Support Forum

TOS

All API Clients

New API Client

## API Clients

 Lyrics-Artist Edit

APP WEBSITE URL  
<https://github.com/spicedacademy/tensor-tarragon-student-code/tree/Mahmoud>

CLIENT ID  

n99AGo0JB4cpSRjAlh15F1AQbERD9YSMN5DureP25eSsp6UJSay6N3j

CLIENT SECRET  

This key is secret! Hover to view and copy

CLIENT ACCESS TOKEN  

oIA4K9BY5VHZEnJTHOIGmA\_uWwcbN71sob1h9MzcVmWhzPT\_laG>

# Collect songs



```
39
40 #define list of artists
41 artists = ['Frank Sinatra','Ed Sheeran','Taylor Swift']
42 max_songs=150
43
44
45
46 #search for songs for each artist in the list and save the songs to the artist folder
47 def collect_songs(artists,max_Songs):
48
49     api=genius.Genius('r2J0C0T6SWh2hLSP5yWaTu2f5LkXbJXc-mMvNg1rE7VDYGoBt-jNcFxRnYbKUjGS',
50                       excluded_terms = ["(Remix)", "(Live)"] ,
51                       skip_non_songs=True,
52                       remove_section_headers=True)
53
54     for artist in artists:
55         songs = (api.search_artist(artist, max_songs=max_Songs, sort='popularity')).songs
56         c=0
57         for song in songs:
58             fileName = os.path.join('lyrics/'+artist,"songnumber"+str(c)+".txt")
59             file = open(fileName, "w")
60             file.write(song.lyrics)
61             file.close()
62             c+=1
63
```

# Collect Songs



```
In [9]: collect_songs(artists,max_songs)
```

```
Searching for songs by Frank Sinatra...
```

```
Song 1: "Fly Me to the Moon"
```

```
Song 2: "My Way"
```

```
Song 3: "That's Life"
```

```
Song 4: "New York, New York"
```

```
Song 5: "Somethin' Stupid"
```

```
Song 6: "The Girl From Ipanema (Garota De Ipanema)"
```

```
Song 7: "The Way You Look Tonight"
```

```
Song 8: "Come Fly with Me"
```

```
Song 9: "Have Yourself a Merry Little Christmas"
```

```
Song 10: "I've Got You Under My Skin"
```

```
Song 11: "My Funny Valentine"
```

```
Song 12: "It Was a Very Good Year"
```

```
Song 13: "Strangers in the Night"
```

```
Song 14: "Blue Moon"
```

```
Song 15: "Luck Be a Lady"
```

```
Song 16: "You Make Me Feel So Young"
```

```
Song 17: "Let It Snow! Let It Snow! Let It Snow!"
```

```
Song 18: "It Had To Be You"
```

# Correct Artist name



```
In [4]: collect_songs(artists,max_songs)
```

```
Searching for songs by ed sheran...
```

```
Changing artist name to 'Ed Sheeran'
```

```
Song 1: "Shape of You"
```

```
Song 2: "Perfect"
```

```
Song 3: "Castle on the Hill"
```

## 2- Data Cleaning



- 1- Natural language toolkit (NLTK)
  - 2- SpaCy
- 
- NLTK is essentially a string processing library, where each function takes strings as input and returns a processed string. But spaCy takes an object-oriented approach.
  - **Which is better?**  
spaCy provides the best way to do it. It provides the fastest and most accurate syntactic analysis of any NLP library released to date.

# Data Cleaning with SpaCy



```
120 # Apply spacy to the text
121 doc=nlp(text)
122 # Lemmatization,remove noise (stopwords, digit, punctuation and single characters)
123 tokens=[token.lemma_.strip() for token in doc if
124         not token.is_stop and not nlp.vocab[token.lemma_].is_stop # remove StopWords
125         and not token.is_punct # Remove punctuation
126         and not token.is_digit # Remove digit
127         and not token.is_space
128         and not token.is_quote
129         and not token.is_bracket
130         and not token.like_num
131         and not token.is_currency
132     ]
133 # Remove empty tokens and one letter tokens
134 tokens = [token for token in tokens if token != "" and len(token)>1]
135 # Recreation of the text
136 new_text=" ".join(tokens)
137 # Remove non alphabetic characters
138 new_text = re.sub(r"^[a-zA-Z]", " ", new_text)
139 # remove non-Unicode characters
140 new_text = re.sub(r"^\x00-\x7F+", "", new_text)
141
142 new_text=new_text.lower()
```



# Data Cleaning with SpaCy



```
1 Things were all good yesterday
2 And then the devil took your memory
3 And if you fell to your death today
4 I hope that heaven is your resting place
5 I heard the doctors put your chest in pain
6 But then that could've been the medicine
7 And now you're lying in the bed again
8 Either way I'll cry with the rest of them
9
10 And my father told me "son
11 It's not his fault he doesn't know your face
12 And you're not the only one"
13 Although my grandma used to say that he used to sing
14
15 Darling, hold me in your arms the way you did last night
16 And we'll lie inside for a little while here, oh
17 I could look into your eyes until the sun comes up
18 And we're wrapped in light and life and love
19 Put your open lips on mine and slowly let them shut
20 For they're designed to be together, oh
21 With your body next to mine, our hearts will beat as one
22 And we're set alight, we're afire love
23 Oh love
```

songnumber36.txt



song\_cleaned\_number36.txt



```
1 thing good yesterday devil memory fall death today hope heaven rest place hear doctor chest pain medicine lie bed way cry rest father tell
son fault know face grandma sing darling hold arm way night lie inside little oh look eye sun come wrap light life love open lip slowly let
shut design oh body heart beat set alight afire love oh love thing good yesterday devil breath away leave pain black suit black tie stand
rain family staple stranger friend come mind paint pen year old remember father tell son fault know face grandma sing darling hold arm way
night lie inside little oh look eye sun come wrap light life love open lip slowly let shut design oh body heart beat set alight afire love
oh love father family rise seat sing hallelujah brother family rise seat sing hallelujah brother family rise seat sing hallelujah brother
sister father family rise seat sing hallelujah
```

# Train a Model



- Create CORPUS
- Vectorize using TfidfVectorizer
- train with:
  - ❖ - Random Forest Classifier
  - ❖ - Multinomial Naive Bayes
- Use GridSearch for Hyperparameters tuning

# Train with Random Forest Classifier



```
182
183 """train model with Random Forest classifier + GridSearch"""
184 def train_model_RF_GridSearch(X_train,y_train):
185     pipeline = Pipeline([
186         ('Tfidf',TfidfVectorizer()),
187         ('RF',RandomForestClassifier())
188     ])
189
190     params = {
191         'Tfidf__max_features':[1000,2000],
192         'Tfidf__ngram_range': [(1, 1), (1, 2), (1, 3)],
193         'Tfidf__min_df':[5,7,10],
194         'Tfidf__max_df':[0.5,0.6,0.7],
195         'RF__n_estimators':[1000,2000],
196         'RF__max_depth':[3,5,7],
197     }
198
199     tfidf_gs = GridSearchCV(pipeline, param_grid=params, cv = 5, verbose = 1,scoring='accuracy',n_jobs=-1)
200     tfidf_gs.fit(X_train,y_train)
201     print('Best parameters:',tfidf_gs.best_params_)
202     best_model=tfidf_gs.best_estimator_
203     return best_model
204
```

# Evaluate RFC



Model accuracy: 0.81



confusion matrix

```
[[19  3  1]
 [ 2 33  0]
 [ 8  3 21]]
```



classification report

	precision	recall	f1-score	support
Ed Sheeran	0.66	0.83	0.73	23
Frank Sinatra	0.85	0.94	0.89	35
Taylor Swift	0.95	0.66	0.78	32
accuracy			0.81	90
macro avg	0.82	0.81	0.80	90
weighted avg	0.84	0.81	0.81	90

# Train with MultinomialNB



```
216
217 """ train model with Naive bias using GridSearch"""
218 def train_model_NB_GridSearch(X_train,y_train):
219
220     pipeline = Pipeline([
221
222         ('Tfidf',TfidfVectorizer(stop_words='english')),
223         ('NB',MultinomialNB())
224     ])
225     params = {
226         'Tfidf__max_features':[1000,2000,4000],
227         'Tfidf__ngram_range': [(1, 1), (1, 2), (1, 3)],
228         'Tfidf__min_df':[5,7,10],
229         'Tfidf__max_df':[0.5,0.6,0.7],
230         'NB__alpha':[1,0.5,0.1, 0.01, 0.001, 0.0001],
231     }
232
233     tfidf_gs = GridSearchCV(pipeline, param_grid=params, cv = 5, verbose = 1,scoring='accuracy',n_jobs=-1)
234     print('start training\n')
235     tfidf_gs.fit(X_train,y_train)
236     print('Best parameters:',tfidf_gs.best_params_)
237     best_model=tfidf_gs.best_estimator_
238     return best_model
239
```

# Evaluate MultinomialNB



```
start training
```

```
Fitting 5 folds for each of 486 candidates, totalling 2430 fits
```

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
```

```
[Parallel(n_jobs=-1)]: Done 42 tasks      | elapsed: 10.9s
```

```
[Parallel(n_jobs=-1)]: Done 192 tasks    | elapsed: 29.2s
```

```
[Parallel(n_jobs=-1)]: Done 442 tasks    | elapsed: 58.0s
```

```
[Parallel(n_jobs=-1)]: Done 792 tasks    | elapsed: 1.6min
```

```
[Parallel(n_jobs=-1)]: Done 1242 tasks   | elapsed: 2.5min
```

```
[Parallel(n_jobs=-1)]: Done 1792 tasks   | elapsed: 3.6min
```

```
[Parallel(n_jobs=-1)]: Done 2430 out of 2430 | elapsed: 4.8min finished
```

```
Best parameters: {'NB__alpha': 0.1, 'Tfidf__max_df': 0.5, 'Tfidf__max_features': 2000, 'Tfidf__min_df': 5, 'Tfidf__ngram_range': (1, 2)}
```

```
Model accuracy: 0.72
```