

Major Project : Music Mood Prediction

Krishni Agrahari, Rajasi Kesharwani, Nikhil Kamale, Kirri Mohitkar, Shazia Khan

```
In [1]: %matplotlib inline
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
import json
warnings.filterwarnings("ignore")
```

Part 2 : Data Pre-processing

We need to preprocess lyrics to use it for NLP models. The main preprocessing steps are - Tokenized words, Stemming, Stop-word removal, Lemmatization.

```
In [2]: import pandas as pd
import nltk.tokenize import word_tokenize
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer
from sklearn.preprocessing import LabelEncoder
from collections import defaultdict
from nltk.corpus import wordnet as wn
from sklearn.ensemble import RandomForestClassifier
import nltk
nltk.download('punkt')
nltk.download('wordnet')
nltk.download('averaged_perceptron_tagger')

[nltk_data] Error loading punkt: <urlopen error [Errno 11001]
[nltk_data]   getaddrinfo failed>
[nltk_data] Error loading wordnet: <urlopen error [Errno 11001]
[nltk_data]   getaddrinfo failed>
[nltk_data] Error loading averaged_perceptron_tagger: <urlopen error
[nltk_data]   [Errno 11001] getaddrinfo failed>

Out[2]: False

In [3]: df=pd.read_csv("training_backup.csv")
df_new=pd.read_csv("testing_backup.csv")
```

Tokenization of words in lyrics

```
In [4]: df_new['lyrics'] = [entry.lower() for entry in df_new['lyrics']]
df_new['lyrics']= [word_tokenize(entry) for entry in df_new['lyrics']]

In [5]: df['lyrics'] = [entry.lower() for entry in df['lyrics']]
df['lyrics']= [word_tokenize(entry) for entry in df['lyrics']]
```

Stemming, Stopword Removal, Lemmatization

```
In [6]: tag_map = defaultdict(lambda: wn.NOUN)
tag_map['V'] = wn.ADV
tag_map['V'] = wn.VERB
tag_map['R'] = wn.ADV
for index,entry in enumerate(df_new['lyrics']):
    Final_words = []
    word_lemmatized = WordNetLemmatizer()
    for word, tag in pos.tag(entry):
        Final_words.append(word_lemmatized(word).lower())
    if word not in stopwords.words('english') and word.isalpha():
        word_Final = word_lemmatized.lemmatize(word,tag_map[tag[0]])
        Final_words.append(word_Final)
    df_new.loc[index, 'text_Final'] = str(Final_words)

In [7]: df_tfidf()
```

```
Out[7]:
```

	lyrics	mood
1477	[yeah-yeah, yeah-yeah, yeah-yeah, yeah-yeah, y...	1
1478	[ip, them, down, hold, me, up, tek, them, th...	3
1479	[it, s, way, too, late, to, think, of, someo...	4
1480	[got, no, place, to, go, but, there, s, a, gi...	2
1481	[every, light, in, the, night, flicker, in, an...	1

```
In [8]: tag_map = defaultdict(lambda: wn.NOUN)
tag_map['V'] = wn.ADV
tag_map['V'] = wn.VERB
tag_map['R'] = wn.ADV
for index,entry in enumerate(df['lyrics']):
    Final_words = []
    word_lemmatized = WordNetLemmatizer()
    for word, tag in pos.tag(entry):
        Final_words.append(word_lemmatized(word).lower())
    if word not in stopwords.words('english') and word.isalpha():
        word_Final = word_lemmatized.lemmatize(word,tag_map[tag[0]])
        Final_words.append(word_Final)
    df.loc[index, 'text_Final'] = str(Final_words)

In [9]: train_x = df['text_Final']
valid_x = df_new['text_Final'][:233]
train_y = df['mood']
valid_y = df_new['mood'][:233]
```

Part 3 : Feature Engineering

We cannot give input as a list of words in Machine Learning models. Feature Engineering needs to be done for the lyrics column. We will use 3 types of NLP models for this -

CountVectorizer

TfidfVectorizer

TfidfNGramModel

At first, Label Encoding -

```
In [10]: Encoder = LabelEncoder()
train_y = Encoder.fit_transform(train_y.ravel())
valid_y = Encoder.fit_transform(valid_y.ravel())
```

```
In [11]: all_texts = []
for stems in train_x:
    all_texts.append(stems)
for stems in valid_x:
    all_texts.append(stems)
print(all_texts[0])

['god', 'need', 'friend', 'god', 'come', 'end', 'god', 'lose', 'mind',
'god', 'find', 'wan', 'na', 'wan', 'na', 'wan', 'na', 'wan', 'na', 'lov
e', 'fell', 'like', 'hate', 'hate', 'felt', 'like', 'love', 'say', 'rea
r, 'disguist', 'wa', 'lake', 'tell', 'na', 'wan', 'na', 'wan', 'na', '
tell', 'believe', 'tell', 'see', 'cause', 'know', 'trust', 'heart', 'ri
ll', 'disguist', 'wa', 'lake', 'tell', 'believe', 'tell', 'see', 'cau
', 'know', 'trust', 'heart', 'fill', 'disguist', 'tell', 'believe', 'tel
', 'believe', 'lady', 'gentleman', 'may', 'attention', 'ready', 'joke', 're
ady', 'great', 'deception', 'tell', 'believe', 'tell', 'believe', 'tel
l', 'believe', 'tell', 'believe', 'tell', 'believe', 'tell', 'believe']
```

```
In [12]: from sklearn import model_selection, preprocessing, linear_model, naive_
bayes, metrics, svm
from sklearn.feature_extraction.text import TfidfVectorizer, CountVecto
rizer
from sklearn import decomposition, ensemble
import pandas, numpy, xgboost, textblob, string
from keras.preprocessing import text, sequence
from keras import layers, models, optimizers
```

CountVectorizer Model

```
In [13]: count_vect = CountVectorizer(analyzer='word')
count_vect.fit(all_texts)
xtrain_count = count_vect.transform(train_x)
xvalid_count = count_vect.transform(valid_x)
```

Tfidf/Vectorizer Model

```
In [14]: import nltk
import string
import re

porter_stemmer = nltk.stem.PorterStemmer()
def porter_tokenizer(text, stemmer=porter_stemmer):
    lower_txt = text.lower()
    tokens = nltk.wordpunct_tokenize(lower_txt)
    stems = [porter_stemmer.stem(t) for t in tokens]
    no_punct = [s for s in stems if re.match('[a-zA-Z]+$', s)]
    no_punct = no_punct
    return no_punct
from sklearn.feature_extraction.text import TfidfVectorizer
tfidf_vect = TfidfVectorizer(
    encoding='utf-8',
    decode_error='replace',
    strip_accents='unicode',
    analyzer='word',
    binary=False,
    stop_words='english',
    tokenizer=porter_tokenizer
)
```

Tfidf-NGram Model

```
In [15]: #tfidf_vect = TfidfVectorizer(analyzer='word',max_features=7000)
tfidf_vect.fit(all_texts)
xtrain_tfidf = tfidf_vect.transform(train_x)
xvalid_tfidf = tfidf_vect.transform(valid_x)

tfidf_vect_ngram = TfidfVectorizer(analyzer='word', ngram_range=(2,2), m
ax_features=7000)
tfidf_vect_ngram.fit(all_texts)
xtrain_tfidf_ngram = tfidf_vect_ngram.transform(train_x)
xvalid_tfidf_ngram = tfidf_vect_ngram.transform(valid_x)
```

Now 3 Training Dataset, 3 Test Dataset are ready to fit into Machine Learning models, one dataset per model defined above.

Part 4 : Mood Prediction Model

```
In [16]: def train_model(classifier, feature_vector_train, label, feature_vector_
valid, is_neural_net=False):
    classifier.fit(feature_vector_train, label)
    predict_the_model=classifier.predict(feature_vector_valid)
    return metrics.accuracy_score(predictions, valid_y)
```

Multinomial Naive Bayes Model

```
In [17]: accuracy_count_nb = train_model(naive_bayes.MultinomialNB(), xtrain_count,
train_y, xvalid_count)
print ("NB, Count Vectors: ", accuracy_count_nb)

accuracy_word_nb = train_model(naive_bayes.MultinomialNB(), xtrain_tfidf,
train_y, xvalid_tfidf)
print ("NB, WordLevel TF-IDF: ", accuracy_word_nb)

accuracy_ngram_nb = train_model(naive_bayes.MultinomialNB(), xtrain_tfidf_ngram,
train_y, xvalid_tfidf_ngram)
print ("NB, N-Gram Vectors: ", accuracy_ngram_nb)

NB, Count Vectors: 0.5836989871244635
NB, Wordlevel TF-IDF: 0.459227467811588
NB, N-Gram Vectors: 0.5793991416399013
```

Random Forest Classifier

```
In [18]: accuracy_count_rf = train_model(RandomForestClassifier(n_estimators=100
),xtrain_count,train_y,xvalid_count)
print("RF, Count Vectors: ",accuracy_count_rf)

accuracy_word_rf = train_model(RandomForestClassifier(n_estimators=100),
xtrain_tfidf,train_y,xvalid_tfidf)
print("RF, WordLevel TF-IDF Vectors: ",accuracy_word_rf)

accuracy_ngram_rf = train_model(RandomForestClassifier(n_estimators=100
),xtrain_tfidf_ngram,train_y,xvalid_tfidf_ngram)
print("RF, N-Gram Vectors: ",accuracy_ngram_rf)

RF, Count Vectors: 0.639981287553648
RF, Wordlevel TF-IDF Vectors: 0.643776824034348
RF, N-Gram Vectors: 0.6180257510729614
```

Logistic Regression Model

```
In [19]: accuracy_count_lr = train_model(linear_model.LogisticRegression(), xtra
in_count, train_y, xvalid_count)
print ("LR, Count Vectors: ", accuracy_count_lr)

accuracy_word_lr = train_model(linear_model.LogisticRegression(), xtrain_
tfidf, train_y, xvalid_tfidf)
print ("LR, Wordlevel TF-IDF: ", accuracy_word_lr)

accuracy_ngram_lr = train_model(linear_model.LogisticRegression(), xtrain_
tfidf_ngram, train_y, xvalid_tfidf_ngram)
print ("LR, N-Gram Vectors: ", accuracy_ngram_lr)

LR, Count Vectors: 0.6351931330472103
LR, Wordlevel TF-IDF: 0.6266994420608058
LR, N-Gram Vectors: 0.5836989871244635
```

XGBoost Classifier

```
In [20]: accuracy_count_xgb = train_model(xgboost.XGBClassifier(), xtrain_count.to
csc(), train_y, xvalid_count.tocsc())
print ("Xgb, Count Vectors: ", accuracy_count_xgb)

accuracy_word_xgb = train_model(xgboost.XGBClassifier(), xtrain_tfidf.toc
sc(), train_y, xvalid_tfidf.tocsc())
print ("Xgb, WordLevel TF-IDF: ", accuracy_word_xgb)

accuracy_ngram_xgb = train_model(xgboost.XGBClassifier(), xtrain_tfidf_ng
ram.tocsc(), train_y, xvalid_tfidf_ngram.tocsc())
print ("Xgb, Ngram Level Vectors: ", accuracy_ngram_xgb)

[14:36:18] WARNING: C:/Users/Administrator/WorkSpace/xgboost-win64_relea
se_1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default eval
uation metric used with the objective 'multi:softprob' was changed from
'merror' to 'mlogloss'. Explicitly set eval_metric if you'd like to rest
ore the old behavior.
Xgb, Count Vectors: 0.686695278999571
[14:36:12] WARNING: C:/Users/Administrator/WorkSpace/xgboost-win64_relea
se_1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default eval
uation metric used with the objective 'multi:softprob' was changed from
'merror' to 'mlogloss'. Explicitly set eval_metric if you'd like to rest
ore the old behavior.
Xgb, Wordlevel TF-IDF: 0.639981287553648
[14:36:10] WARNING: C:/Users/Administrator/WorkSpace/xgboost-win64_relea
se_1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default eval
uation metric used with the objective 'multi:softprob' was changed from
'merror' to 'mlogloss'. Explicitly set eval_metric if you'd like to rest
ore the old behavior.
Xgb, Ngram Level Vectors: 0.5665236651562146
```

Bagging - Random Forest

```
In [21]: accuracy_count_bg = train_model(ensemble.RandomForestClassifier(), xtra
in_count, train_y, xvalid_count)
print ("RF, Count Vectors: ", accuracy_count_bg)

accuracy_word_bg = train_model(ensemble.RandomForestClassifier(), xtrain_
tfidf, train_y, xvalid_tfidf)
print ("RF, WordLevel TF-IDF: ", accuracy_word_bg)

accuracy_ngram_bg = train_model(ensemble.RandomForestClassifier(), xtra
in_tfidf_ngram, train_y, xvalid_tfidf_ngram)
print ("RF, Ngram Vectors: ", accuracy_ngram_bg)

RF, Count Vectors: 0.648068689527897
RF, Wordlevel TF-IDF: 0.644068689527897
RF, Ngram Vectors: 0.5879822822618258
```

```
In [22]: from sklearn.svm import LinearSVC
accuracy_count_svm = train_model(LinearSVC(multi_class='ovr'), xtrain_co
unt, train_y, xvalid_count)
print ("SVM, Count Vectors: ", accuracy_count_svm)

accuracy_word_svm = train_model(LinearSVC(multi_class='ovr'), xtrain_tfi
df, train_y, xvalid_tfidf)
print ("SVM, WordLevel TF-IDF: ", accuracy_word_svm)

accuracy_ngram_svm = train_model(LinearSVC(multi_class='ovr'), xtrain_tf
idf_ngram, train_y, xvalid_tfidf_ngram)
print ("SVM, Ngram Vectors: ", accuracy_ngram_svm)

SVM, Count Vectors: 0.639981287553648
SVM, Wordlevel TF-IDF: 0.652366524034348
SVM, Ngram Vectors: 0.622317595665236
```

```
In [23]: import matplotlib as mpl
import numpy as np
def plot_cm(X, y, clf, title):
    cm = metrics.confusion_matrix(y, clf.predict(X))
    print(metrics.classification_report(y, clf.predict(X)))
    np.set_printoptions(suppress=True)
    mpl.rc('figure', figsize=(10,10))

    hm = sns.heatmap(cm,
        char=False,
        annot=True,
        square=True,
        fmt='.1',
        yticklabels=['happy','sad','angry','relaxed'],
        xticklabels=['happy','sad','angry','relaxed'],
        cmap='Blues')

    plt.title(title)
    plt.ylabel('actual class')
    plt.xlabel('predicted class')
    plt.tight_layout()
    #fig.savefig('rf_cold.jpg')
    plt.show()
```

```
In [24]: # import seaborn as sns
import matplotlib.pyplot as plt

clf=ensemble.RandomForestClassifier()
clf.fit(xtrain_count,train_y)
plot_cm(xvalid_count, valid_y, clf, "Training CountVectorizer on LR")
```

	precision	recall	f1-score	support
0	0.78	0.63	0.67	97
1	0.59	0.85	0.63	68
2	0.69	0.33	0.59	12
3	0.79	0.41	0.54	56
accuracy				
macro avg	0.75	0.56	0.59	233
weighted avg	0.69	0.63	0.62	233

Training Word Level - TfidfVectorizer on LR

```
In [25]: # import seaborn as sns
import matplotlib.pyplot as plt

clf=linear_model.LogisticRegression()
clf.fit(xtrain_count,train_y)
plot_cm(xvalid_count, valid_y, clf, "Training CountVectorizer on LR")
```

	precision	recall	f1-score	support
0	0.78	0.66	0.72	97
1	0.58	0.66	0.62	68
2	0.44	0.33	0.38	12
3	0.54	0.82	0.58	56
accuracy				
macro avg	0.59	0.57	0.54	233
weighted avg	0.65	0.64	0.67	233

Training CountVectorizer on LR

```
In [26]: # import seaborn as sns
import matplotlib.pyplot as plt

clf=ensemble.RandomForestClassifier()
clf.fit(xtrain_tfidf_ngram,train_y)
plot_cm(xvalid_tfidf_ngram, valid_y, clf, "Training NGram Vectors on R
F")
```

	precision	recall	f1-score	support
0	0.84	0.47	0.61	97
1	0.52	0.79	0.63	68
2	0.80	0.33	0.47	12
3	0.51	0.62	0.56	56
accuracy				
macro avg	0.67	0.56	0.57	233
weighted avg	0.68	0.60	0.59	233

Training NGram Vectors on RF

```
In [27]: # import seaborn as sns
import matplotlib.pyplot as plt

clf=svm.LinearSVC(multi_class='ovr')
clf.fit(xtrain_tfidf,train_y)
plot_cm(xvalid_tfidf, valid_y, clf, "Training Word Level - TfidfVecto
r on SVM")
```

	precision	recall	f1-score	support
0	0.76	0.69	0.72	97
1	0.64	0.69	0.67	68
2	0.80	0.33	0.47	12
3	0.55	0.61	0.55	56
accuracy				
macro avg	0.68	0.58	0.65	233
weighted avg	0.67	0.65	0.65	233

Training Word Level - TfidfVectorizer on SVM

```
In [28]: def pred(lyrics):
    wt=word_tokenize(lyrics)
    tag_map = defaultdict(lambda: wn.NOUN)
    tag_map['J'] = wn.ADJ
    tag_map['V'] = wn.ADV
    tag_map['R'] = wn.VERB
    tag_map['N'] = wn.ADV
    Final_words = []
    word_lemmatized = WordNetLemmatizer()
    for word, tag in pos.tag(wt):
        if word not in stopwords.words('english') and word.isalpha():
            word_Final = word_lemmatized.lemmatize(word,tag_map[tag[0]])
            Final_words.append(word_Final)
    result = str(Final_words)
    df=pd.DataFrame(columns=['lyrics'])
    df=df.append({'lyrics':result,ignore_index=True})
    # xvalid_count = count_vect.transform(result)
    # res = Encoder.fit_transform(result)
    test_x=df['lyrics']
    # print(test_x.shape)
    xvalid_tfidf = tfidf_vect.transform(testx)
    y=clf.predict(xvalid_tfidf)
    print("Tag: ",end='')
    if y==0:
        print("Happy")
    elif(y==1):
        print("Sad")
    elif(y==2):
        print("Angry")
    elif(y==3):
        print("Relaxed")

In [29]: from googletrans import Translator
def pred1(lyrics):
    translator = Translator()
    l1=translator.translate(lyrics)
    sprint(l1)
    lyrics=l1.text
    wt=word_tokenize(lyrics)
    tag_map = defaultdict(lambda: wn.NOUN)
    tag_map['J'] = wn.ADJ
    tag_map['V'] = wn.VERB
    tag_map['R'] = wn.ADV
    Final_words = []
    word_lemmatized = WordNetLemmatizer()
    for word, tag in pos.tag(wt):
        if word not in stopwords.words('english') and word.isalpha():
            word_Final = word_lemmatized.lemmatize(word,tag_map[tag[0]])
            Final_words.append(word_Final)
    result = str(Final_words)
    df=pd.DataFrame(columns=['lyrics'])
    df=df.append({'lyrics':result,ignore_index=True})
    # xvalid_count = count_vect.transform(result)
    # res = Encoder.fit_transform(result)
    test_x=df['lyrics']
    # print(test_x.shape)
    xvalid_tfidf = tfidf_vect.transform(testx)
    y=clf.predict(xvalid_tfidf)
    print("Tag: ",end='')
    if y==0:
        print("Happy")
    elif(y==1):
        print("Sad")
    elif(y==2):
        print("Angry")
    elif(y==3):
        print("Relaxed")
```

Ilahi song by Anjii Singh

```
In [33]: pred("Evenings,as if they're mad,\nNights like some tunnel,\nI don't k
now why on rebel flights only,\nO God, my heart comes,\nO God, my heart is
interested,\nThere is a question about tomorrow,\nAs in, no one kno
ws what happens tomorrow,\nLiving life is now,\nI don't know why in roma
dic ways only,\nI find my heart interested, O Lord,\nMy heart is interest
ed,\nMy philosophy is : 'My bab is on my shoulders,\nI go wherever the p
ath takes me,\nI don't on the road,\nBut the sea of dross itself,\nO God...
O God...")
Tag: Happy
```

```
In [32]: pred("shaam mein milang se\nVandetiin surang se\nVBaghi udan pe hi na a
ane kyun\nIlahi mera jee aaye aye\nIlahi mera jee aaye aye\nDaba da
ang dung dung dung dang...,\nKai pe saawal hai\nJeena filhaal hai\nKhanaa
bahiyaon pe hi jana kyun\nIlahi mera jee aaye aye\nIlahi mera jee aaye aye\nMeraa
faisafa kandhe pe mera basteen\nChalaa main jahaan le
o halaa mujhe rasteen\nBodonon pe nahin\nBodonon ke samandar pe wo-o-wo-o-
\nIlahi mera jee aaye aye\nIlahi mera jee aaye aye...,\nShaam mein mil
ang se\nVandetiin surang se\nVBaghi udan pe hi na jana kyun\nIlahi mera
jee aaye aye\nIlahi mera jee aaye aye\nIlahi... Ilahi... Ilahi...")
Tag: Happy
```

Tera Hone Laga Hoon song by Anif Aslam

```
In [32]: pred("I'm becoming yours,\nGetting lost in my thoughts,\nSince I met yo
u...,\nShould I touch you with my sight,\nUnbecame my arns yearn for you,
\nMy heart has called out,\nWho come on now...,\nLook, the dewdrops fall as
 rain,\nEven the weather drops a hint,\nWho come on now...,\nFemale:\nMera i
n arn,\nLike an embrace,\nI admit, and I agree,\nI'm in love with you. A
I'm...,\nBecoming yours,\nGetting lost in my thoughts,\nSince I met yo
u...")
Tag: Happy
```

```
In [33]: pred("bahon mein dale bahin bahain\nBahon ka jaise haar hua\nMehaan main
a main ne maana maana\nMujhe bhi pyaar hua\nTera hone laga hoon\nKha
one Laga hoon\nJab se mila hoon\nTera hone laga hoon\nKhone Laga hoon\nJ
ab se mila hoon")
Tag: Happy
```

```
In [34]: pred("All the faces are story, even the hearts are made of stone\nThen w
hy do I remain lonely, and roam the streets like a destitute\nWho shall
not attain anything here, my heart\nI'm a dream breakin' a single moment
\nThe world seems lonely\nThe world seems lonely\nWhen no loved one rema
ins\nThe world seems lonely\nThe world seems lonely\nWho does this happ
en\nThen this heart cries out\nEven the wind seems as if it is crying?
\nThe world seems lonely\nI'm a dream breaks in a single moment\nThe worl
d seems lonely\nThe world seems lonely\nWhen no loved one remains\nThe w
orld seems lonely\nThe world seems lonely\nSeems lonely")
Tag: Sad
```

```
In [35]: pred("Naine roya tuje dhund Sisak Sisak Ke Hawayein\nAnjag Soona Laage
\nChhai Se Jo Tootte Koi ")
Tag: Happy
```

Maine Roy's song by Tanveer Evan

```
In [36]: pred("I've said I cried you in the mist of mist, I didn't get you to sleep
p in these nights, I've never thought of reha in your memories, I never
thought you just cried.")
Tag: Sad
```

```
In [37]: pred("Naine roya tuje dhund dhunde Roxytha main Tu na mila\nBehrissi in
raon main soyaa main tha main Rotathi Teri yaadon mein woh hi raat")
Tag: Sad
```

```
In [38]: pred("Jo dil mein bhara tune Dekhke uss zehar ko\nBhugtegi mere gham
ko\nKehi aah ke Kahaar ko\nApni Rotgarzi ka ab anzaam dekhegi")
Tag: Sad
```

```
In [ ]:
```