

```
[1]: import pandas as pd
import numpy as np
import cufflinks as cf
from plotly.offline import download_plotlyjs, init_notebook_mode, plot, iplot
init_notebook_mode(connected=True) # to get the connection
cf.go_offline() #offline mode
from sklearn.model_selection import train_test_split
from scipy.sparse import hstack,csr_matrix
import warnings
warnings.filterwarnings("ignore")
from sklearn.metrics import roc_auc_score, confusion_matrix, roc_curve, auc, accuracy_score
import matplotlib.pyplot as plt
from sklearn.feature_selection import SelectKBest, chi2
from sklearn import naive_bayes
from tqdm import tqdm
import seaborn as sns
matplotlib inline
```

```
In [2]: data=pd.read_csv('/home/sathish/Desktop/appliedai/Donore_choose_preprocessed_data.csv',na_filter=False)

In [3]: y=data['project_is_approved'].values
x=data.drop(['project_is_approved'],axis=1)
```

# 1. Splitting Data

```
In [4]: x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3,shuffle=True)
len_train,x_train.shape[0],len_val,x_val.shape[0],len_test,x_test.shape[0]
```

# 2. Creating Data Matrix

## 2.1 converting text to vectors

### 2.1.1 Bag of words

```
In [5]: from sklearn.feature_extraction.text import CountVecorizer
vectorizer=CountVecorizer()
bow_features=[]
x_train_essay_bow=vectorizer.fit_transform(x_train['essay'].values)
x_val_essay_bow=vectorizer.transform(x_val['essay'].values)
x_test_essay_bow=vectorizer.transform(x_test['essay'].values)

bow_features.extend(vectorizer.get_feature_names())

x_train_summary_bow=vectorizer.fit_transform(x_train['project_resource_summary'].values)
x_val_summary_bow=vectorizer.transform(x_val['project_resource_summary'].values)
x_test_summary_bow=vectorizer.transform(x_test['project_resource_summary'].values)

bow_features.extend(vectorizer.get_feature_names())

x_train_title_bow=vectorizer.fit_transform(x_train['project_title'].values)
x_val_title_bow=vectorizer.transform(x_val['project_title'].values)
x_test_title_bow=vectorizer.transform(x_test['project_title'].values)

bow_features.extend(vectorizer.get_feature_names())

In [6]: len(bow_features)

Out[6]: 68536
```

### 2.1.2 TFIDF

```
In [7]: from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer=TfidfVectorizer()

tfidf_features=[]
x_train_essay_tfidf=vectorizer.fit_transform(x_train['essay'].values)
x_val_essay_tfidf=vectorizer.transform(x_val['essay'].values)
x_test_essay_tfidf=vectorizer.transform(x_test['essay'].values)

tfidf_features.extend(vectorizer.get_feature_names())

x_train_summary_tfidf=vectorizer.fit_transform(x_train['project_resource_summary'].values)
x_val_summary_tfidf=vectorizer.transform(x_val['project_resource_summary'].values)
x_test_summary_tfidf=vectorizer.transform(x_test['project_resource_summary'].values)

tfidf_features.extend(vectorizer.get_feature_names())

x_train_title_tfidf=vectorizer.fit_transform(x_train['project_title'].values)
x_val_title_tfidf=vectorizer.transform(x_val['project_title'].values)
x_test_title_tfidf=vectorizer.transform(x_test['project_title'].values)

tfidf_features.extend(vectorizer.get_feature_names())
```

## 2.2 one hot encoding the categorical features

```
In [8]: cat_features=[]
vectorizer=CountVecorizer()
x_train_state_onehot=vectorizer.fit_transform(x_train['school_state'].values)
x_val_state_onehot=vectorizer.transform(x_val['school_state'].values)
x_test_state_onehot=vectorizer.transform(x_test['school_state'].values)

cat_features.extend(vectorizer.get_feature_names())

x_train_tprefix_onehot=vectorizer.fit_transform(x_train['teacher_prefix'].values)
x_val_tprefix_onehot=vectorizer.transform(x_val['teacher_prefix'].values)
x_test_tprefix_onehot=vectorizer.transform(x_test['teacher_prefix'].values)

cat_features.extend(vectorizer.get_feature_names())

x_train_grade_onehot=vectorizer.fit_transform(x_train['project_grade_category'].values)
x_val_grade_onehot=vectorizer.transform(x_val['project_grade_category'].values)
x_test_grade_onehot=vectorizer.transform(x_test['project_grade_category'].values)

cat_features.extend(vectorizer.get_feature_names())

x_train_cat_onehot=vectorizer.fit_transform(x_train['project_subject_categories'].values)
x_val_cat_onehot=vectorizer.transform(x_val['project_subject_categories'].values)
x_test_cat_onehot=vectorizer.transform(x_test['project_subject_categories'].values)

cat_features.extend(vectorizer.get_feature_names())

x_train_subcat_onehot=vectorizer.fit_transform(x_train['project_subject_subcategories'].values)
x_val_subcat_onehot=vectorizer.transform(x_val['project_subject_subcategories'].values)
x_test_subcat_onehot=vectorizer.transform(x_test['project_subject_subcategories'].values)

cat_features.extend(vectorizer.get_feature_names())
```

## 2.3 Normalizing the numerical features

```
In [9]: from sklearn.preprocessing import Normalizer
import scipy
normalizer=Normalizer()
x_train_price_norm=normalizer.fit_transform(x_train['price'].values.reshape(-1,1))
x_val_price_norm=normalizer.transform(x_val['price'].values.reshape(-1,1))
x_test_price_norm=normalizer.transform(x_test['price'].values.reshape(-1,1))

x_train_quantity_norm=normalizer.fit_transform(x_train['quantity'].values.reshape(-1,1))
x_val_quantity_norm=normalizer.transform(x_val['quantity'].values.reshape(-1,1))
x_test_quantity_norm=normalizer.transform(x_test['quantity'].values.reshape(-1,1))
```

## Combining all the feature names

```
In [10]: bow_features.extend(cat_features)
bow_features.extend(['price','quantity'])

tfidf_features.extend(cat_features)
tfidf_features.extend(['price','quantity'])
```

## 2.4 Building data matrix

### 2.4.1 Set 1: Categorical, numerical, essay(BOW), summary(BOW), title(BOW)

```
In [11]: #text features
bow_train=hstack((x_train_essay_bow,x_train_summary_bow,x_train_title_bow))
bow_valid=hstack((x_val_essay_bow,x_val_summary_bow,x_val_title_bow))
bow_test=hstack((x_test_essay_bow,x_test_summary_bow,x_test_title_bow))

#all categorical features
x_train_cat=hstack((x_train_state_onehot,x_train_tprefix_onehot,x_train_grade_onehot,x_train_cat_onehot,x_train_subcat_onehot))
x_val_cat=hstack((x_val_state_onehot,x_val_tprefix_onehot,x_val_grade_onehot,x_val_cat_onehot,x_val_subcat_onehot))
x_test_cat=hstack((x_test_state_onehot,x_test_tprefix_onehot,x_test_grade_onehot,x_test_cat_onehot,x_test_subcat_onehot))

#numerical features
x_train_num=hstack((csr_matrix(x_train_price_norm),x_train_quantity_norm))
x_val_num=hstack((csr_matrix(x_val_price_norm),x_val_quantity_norm))
x_test_num=hstack((csr_matrix(x_test_price_norm),x_test_quantity_norm))

#final matrices
final_bow_train=hstack((bow_train,x_train_cat,x_train_num)).tocsr()
final_bow_val=hstack((bow_val,x_val_cat,x_val_num)).tocsr()
final_bow_test=hstack((bow_test,x_test_cat,x_test_num)).tocsr()
```

### 2.4.2 Set 2: Categorical, numerical, essay(TFIDF), summary(TFIDF), title(TFIDF)

```
In [12]: #text features
tfidf_train=hstack((x_train_essay_tfidf,x_train_summary_tfidf,x_train_title_tfidf))
tfidf_valid=hstack((x_val_essay_tfidf,x_val_summary_tfidf,x_val_title_tfidf))
tfidf_test=hstack((x_test_essay_tfidf,x_test_summary_tfidf,x_test_title_tfidf))

#all categorical features
x_train_cat=hstack((x_train_state_onehot,x_train_tprefix_onehot,x_train_grade_onehot,x_train_cat_onehot,x_train_subcat_onehot))
x_val_cat=hstack((x_val_state_onehot,x_val_tprefix_onehot,x_val_grade_onehot,x_val_cat_onehot,x_val_subcat_onehot))
x_test_cat=hstack((x_test_state_onehot,x_test_tprefix_onehot,x_test_grade_onehot,x_test_cat_onehot,x_test_subcat_onehot))

#numerical features
x_train_num=hstack((csr_matrix(x_train_price_norm),x_train_quantity_norm))
x_val_num=hstack((csr_matrix(x_val_price_norm),x_val_quantity_norm))
x_test_num=hstack((csr_matrix(x_test_price_norm),x_test_quantity_norm))

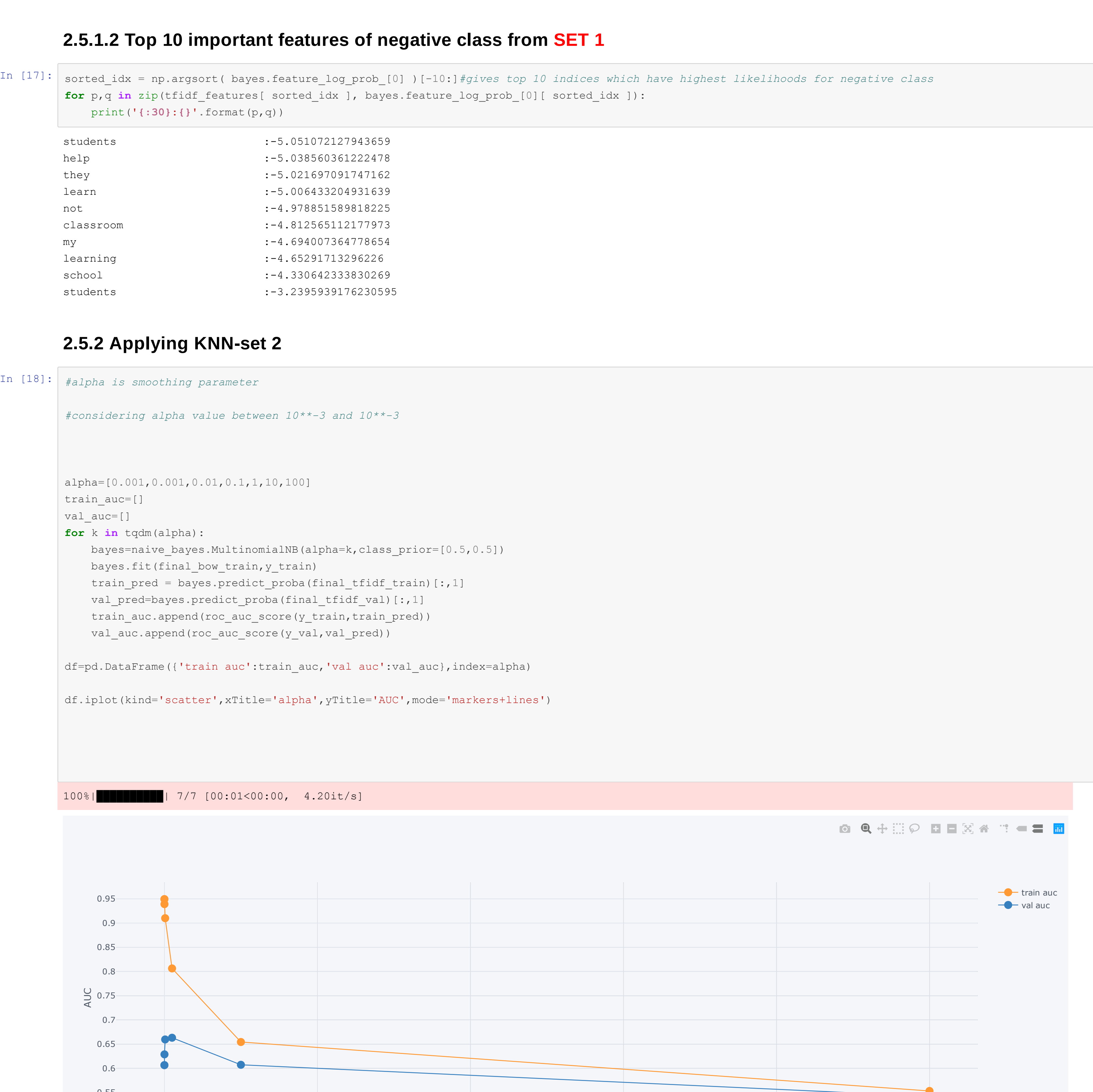
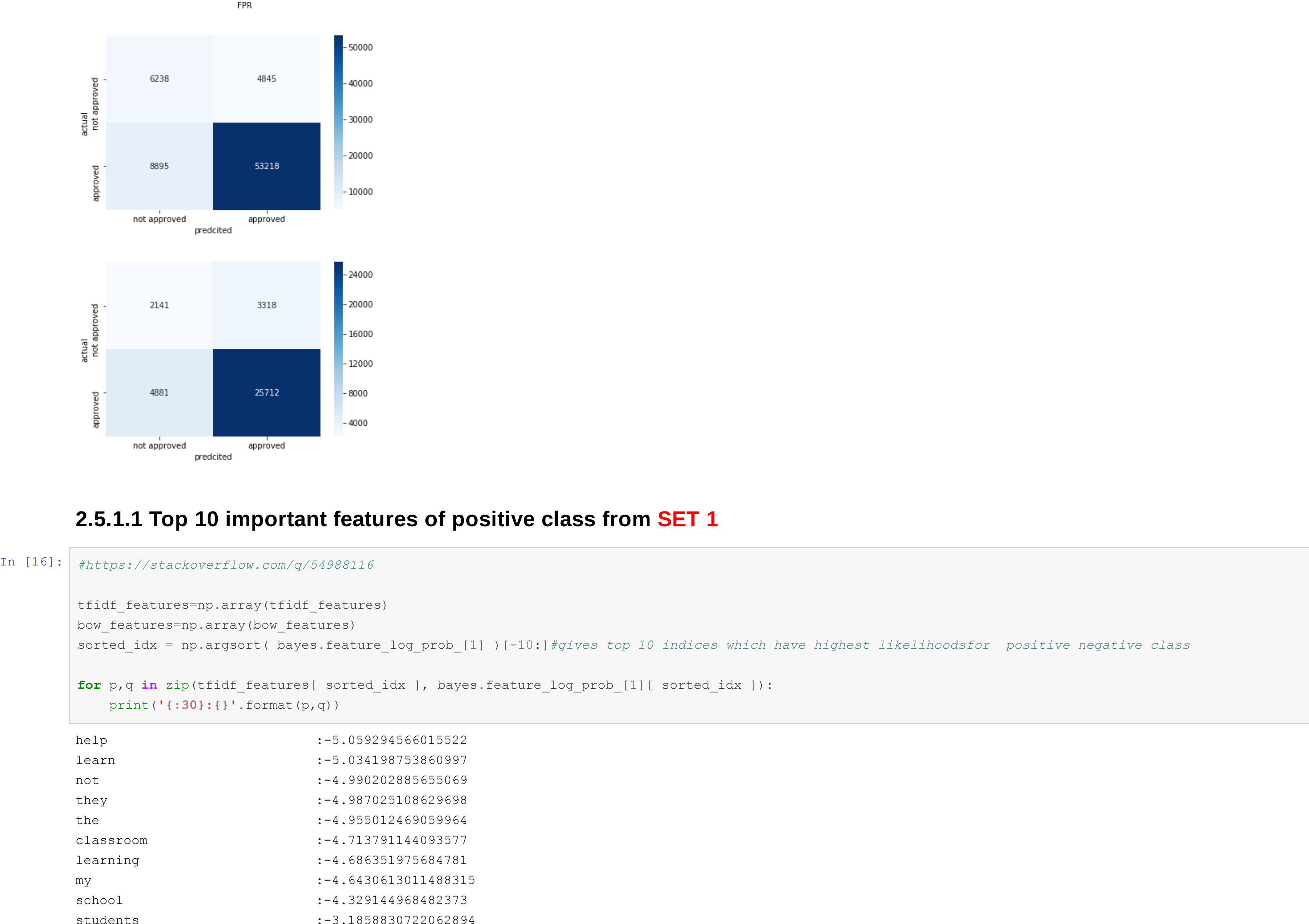
#final matrices
final_tfidf_train=hstack((tfidf_train,x_train_cat,x_train_num)).tocsr()
final_tfidf_val=hstack((tfidf_valid,x_val_cat,x_val_num)).tocsr()
final_tfidf_test=hstack((tfidf_test,x_test_cat,x_test_num)).tocsr()

In [13]: #making sure that everything is ok
print(final_bow_train.shape,final_bow_val.shape,final_bow_test.shape)
print(final_tfidf_train.shape,final_tfidf_val.shape,final_tfidf_test.shape)

(49041, 68664) (24155, 68664) (36052, 68664)
(49041, 68664) (24155, 68664) (36052, 68664)
```

## 2.5 Naive Bayes

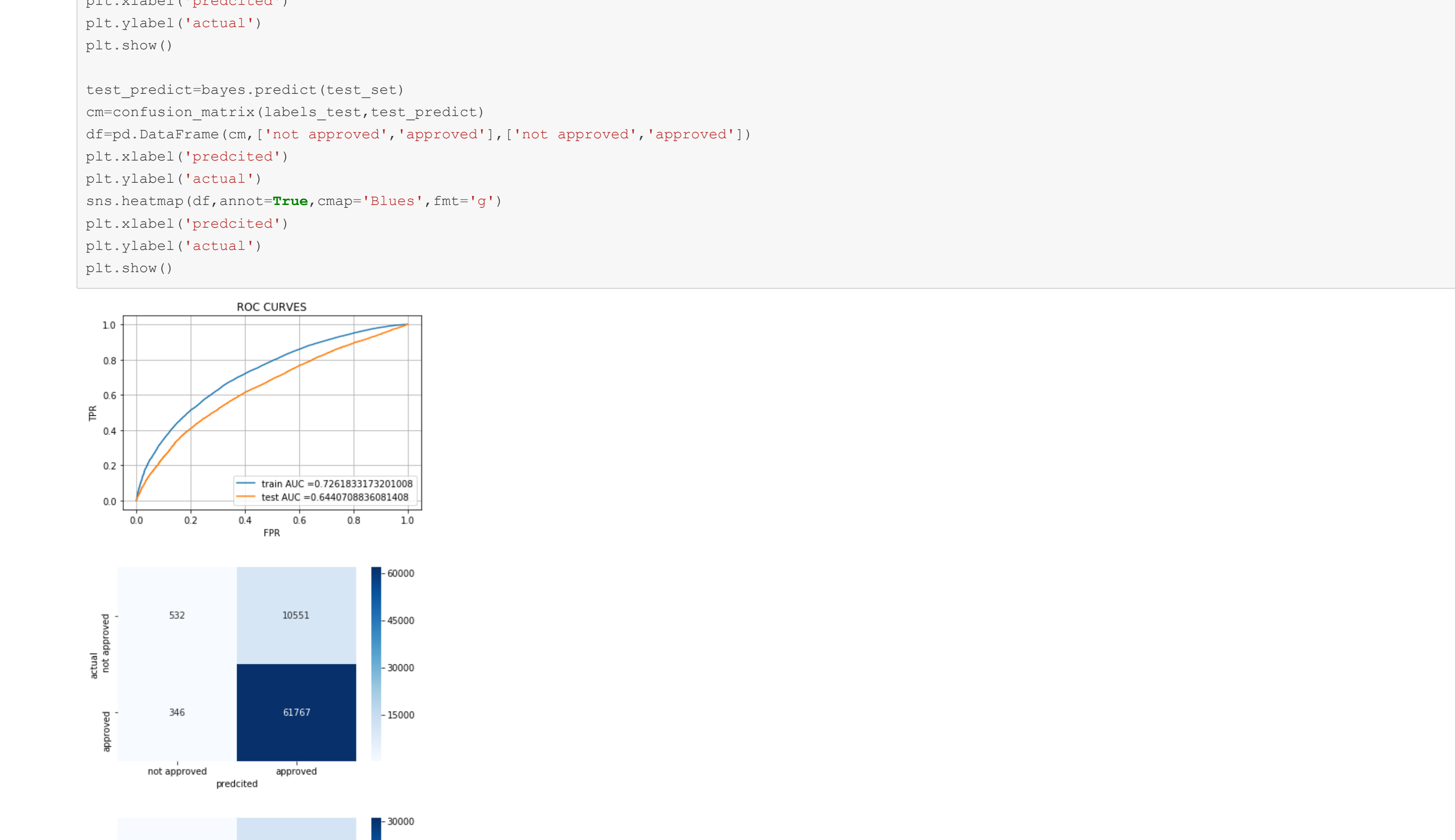
### 2.5.1 Applying Naive Bayes -Set 1



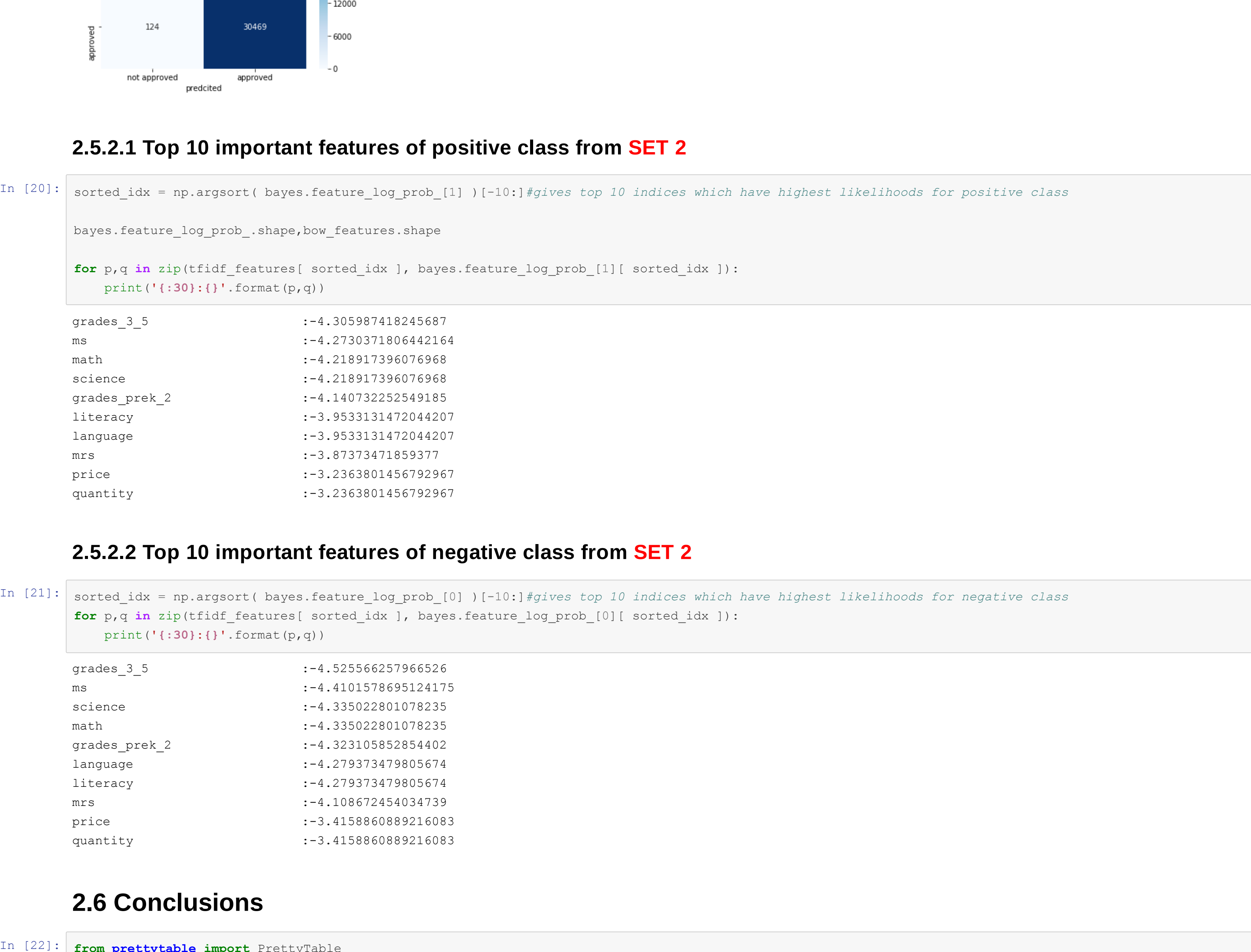
### 2.5.1.1 Top 10 important features of positive class from SET 1



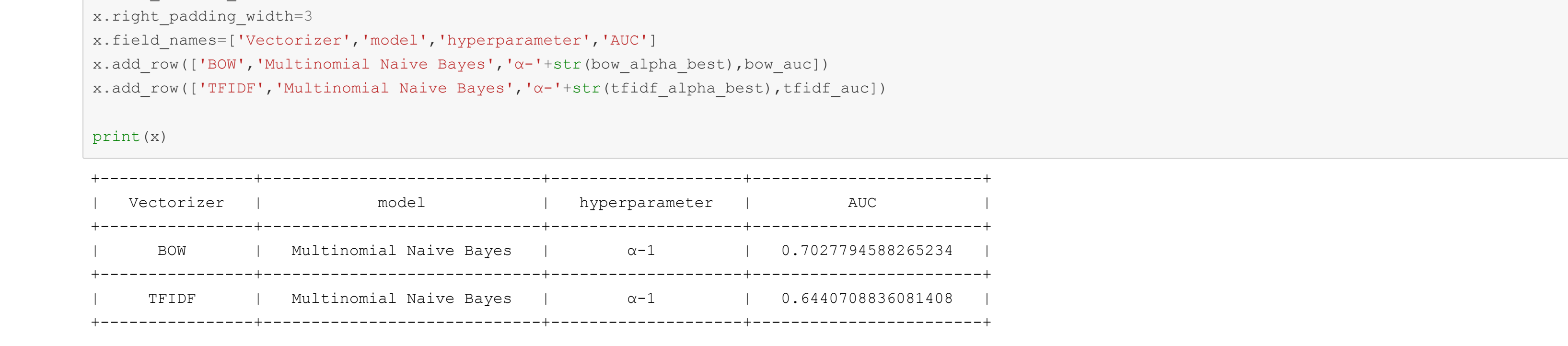
### 2.5.1.2 Top 10 important features of negative class from SET 1



### 2.5.2 Applying KNN-set 2



### 2.5.2.1 Top 10 important features of positive class from SET 2



### 2.5.2.2 Top 10 important features of negative class from SET 2



## 2.6 Conclusions

```
In [22]: from prettytable import PrettyTable
import prettytable
p=PrettyTable()
x.headers=prettytable.ALL
x.left_padding_width=3
x.right_padding_width=3
x.field_names=['Vectorizer','model','hyperparameter','AUC']
x.add_row(['BOW','Multinomial Naive Bayes','alpha=1','0.702779458902634'])
x.add_row(['TFIDF','Multinomial Naive Bayes','alpha=1','0.644070883081408'])
print(x)
```

Vectorizer	model	hyperparameter	AUC
BOW	Multinomial Naive Bayes	alpha=1	0.702779458902634
TFIDF	Multinomial Naive Bayes	alpha=1	0.644070883081408