


```
In [135]: #convert school state to lower case
project_data['school_state'] = project_data['school_state'].str.lower()
project_data['school_state'].value_counts()

Out[135]:
ca    14990
tx     6764
de     6152
fl     5693
nc     4651
ar     3983
ga     3638
sc     3594
mi     2886
pa     2847
in     2499
mo     2368
al     2278
ma     2183
la     2164
de     2130
ok     2111
nj     2042
az     1969
mi     1874
wi     1668
al     1606
ms     1551
tn     1545
ct     1511
md     1384
nv     1252
ms     1198
ky     1197
mn     1190
co     1024
id     633
ia     606
ks     561
ar     498
dc     477
wv     472
hi     468
me     464
ak     321
sh     313
ph     312
ne     289
sd     268
ri     187
nd     216
mt     134
al     119
vt     75
Name: school_state, dtype: int64
```

```
In [136]: # we get the cost of the project using resource.csv file
resource_data.head(2)

Out[136]:
   id                description  quantity  price
0  0233245  L0552-Lakeshore Double-Space Mobile Drying Rack           1    140.00
1  p005062  Bouncy Bands for Questions (Blue support pipes)           3    14.95

In [137]: #https://stackoverflow.com/questions/22497798/how-to-reset-a-dataframes-indexes-for-all-groups-in-o
n-step
price_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_index()
price_data.head(2)

Out[137]:
   id  price  quantity
0  p000001    459.95           7
1  p000002   515.89          21
```

```
In [138]: #get two dataframes in python:
project_data = pd.merge(project_data, price_data, on='id', how='left')

In [139]: project_data.columns

Out[139]:
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
      'date', 'project_grade_category', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'essay',
      'preprocessed_essays', 'preprocessed_title', 'price', 'quantity'],
      dtype='object')

we are going to consider
- school_state : categorical data
- clean_categories : categorical data
- clean_subcategories : categorical data
- project_grade_category : categorical data
- teacher_prefix : categorical data
- project_title : text data
- text : text data
- project_resource_summary : text data (optional)
- quantity : numerical (optional)
- teacher_number_of_previously_posted_projects : numerical
- price : numerical
```

```
In [140]: project_data.head(2)

Out[140]:
   Unnamed: 0    id  teacher_id  teacher_prefix  school_state  Date  project_grade_category  project_title  price
0            0    8393  p005479   20H07h008965e58b2a3269b2b3c4e5    ms             ca      2016-04-27 00:27:36  grades_prek_2  Engineering STEAM into the Primary Classroom  14.95
1            1   37728  p043609   3B60494c61921b3b43b61dd0dc2904df    ms             ut      2018-04-27 00:31:25  grades_3_5  Solos Tools for Focus  515.89
```

2 rows x 22 columns

2. Naïve Bayes

2.1 Splitting data into Train and cross validation(or test): Stratified Sampling

```
In [145]: y = project_data['project_is_approved'].values
X = project_data.drop(['project_is_approved'], axis=1)
X.head(2)

# train test split
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, stratify=y)
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33, stratify=y_train)

print("=="100)

=====
```

2.2 Make Data Model Ready: encoding numerical, categorical features

```
In [146]: print(X_train.shape, y_train.shape)
print(X_test.shape, y_test.shape)
print(X_cv.shape, y_cv.shape)

print("=="100)

(44890, 21) (44890,)
(22110, 21) (22110,)
(33000, 21) (33000,)

=====
```

Vectorizing Categorical features

Vectorizing : clean_categories

```
In [147]: # we use count vectorizer to convert the values into one hot encoded features
from sklearn.feature_extraction.text import CountVectorizer
vectorizer_clean_categories = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
vectorizer_clean_categories.fit(X_train['clean_categories'].values)
print(vectorizer_clean_categories.get_feature_names())

X_train_categories_one_hot = vectorizer_clean_categories.transform(X_train['clean_categories']).value
X_cv_categories_one_hot = vectorizer_clean_categories.transform(X_cv['clean_categories'].values)
X_test_categories_one_hot = vectorizer_clean_categories.transform(X_test['clean_categories'].values)

print("Shape of matrix after one hot encoding ")
print(X_train_categories_one_hot.shape, y_train.shape)
print(X_cv_categories_one_hot.shape, y_cv.shape)
print(X_test_categories_one_hot.shape, y_test.shape)

['Arthritis', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'Applied_Learning', 'SpecialNeeds', 'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encoding
(44890, 9) (44890,)
(22110, 9) (22110,)
(33000, 9) (33000,)
```

Vectorizing : clean_subcategories

```
In [148]: # Vectorizing clean_subcategories
vectorizer_clean_subcategories = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
vectorizer_clean_subcategories.fit(X_train['clean_subcategories'].values)

X_train_sub_categories_one_hot = vectorizer_clean_subcategories.transform(X_train['clean_subcategories'].values)
X_cv_sub_categories_one_hot = vectorizer_clean_subcategories.transform(X_cv['clean_subcategories'].values)
X_test_sub_categories_one_hot = vectorizer_clean_subcategories.transform(X_test['clean_subcategories'].values)

print("Shape of matrix after one hot encoding ")
print(X_train_sub_categories_one_hot.shape, y_train.shape)
print(X_cv_sub_categories_one_hot.shape, y_cv.shape)
print(X_test_sub_categories_one_hot.shape, y_test.shape)

print("=="100)

Shape of matrix after one hot encoding
(44890, 30) (44890,)
(22110, 30) (22110,)
(33000, 30) (33000,)
```

Vectorizing : school_state

```
In [149]: # Vectorizing school_state
#code ref: https://www.youtube.com/watch?v=0H0qDcIn3Z4&t=5308
#provided did not work
vectorizer_school_state = CountVectorizer(binary=True)
vectorizer_school_state.fit(X_train['school_state'].values)

X_train_school_state_one_hot = vectorizer_school_state.transform(X_train['school_state'].values)
X_cv_school_state_one_hot = vectorizer_school_state.transform(X_cv['school_state'].values)
X_test_school_state_one_hot = vectorizer_school_state.transform(X_test['school_state'].values)

print("Shape of matrix after one hot encoding ")
print(X_train_school_state_one_hot.shape, y_train.shape)
print(X_cv_school_state_one_hot.shape, y_cv.shape)
print(X_test_school_state_one_hot.shape, y_test.shape)

print("=="100)

Shape of matrix after one hot encoding
(44890, 51) (44890,)
(22110, 51) (22110,)
(33000, 51) (33000,)
```

Vectorizing : teacher_prefix

```
In [150]: # Vectorizing teacher_prefix
vectorizer_teacher_prefix = CountVectorizer(binary=True)
vectorizer_teacher_prefix.fit(X_train['teacher_prefix'].values)

# we use the fitted CountVectorizer to convert the text to vector
X_train_teacher_prefix_one_hot = vectorizer_teacher_prefix.transform(X_train['teacher_prefix']).value
X_cv_teacher_prefix_one_hot = vectorizer_teacher_prefix.transform(X_cv['teacher_prefix'].values)
X_test_teacher_prefix_one_hot = vectorizer_teacher_prefix.transform(X_test['teacher_prefix'].values)

print("Shape of matrix after one hot encoding ")
print(X_train_teacher_prefix_one_hot.shape, y_train.shape)
print(X_cv_teacher_prefix_one_hot.shape, y_cv.shape)
print(X_test_teacher_prefix_one_hot.shape, y_test.shape)

print("=="100)

Shape of matrix after one hot encoding
(44890, 5) (44890,)
(22110, 5) (22110,)
(33000, 5) (33000,)
```

Vectorizing : project_grade_category

```
In [151]: #encoding project_grade_category
vectorizer_project_grade_category = CountVectorizer(binary=True)
vectorizer_project_grade_category.fit(X_train['project_grade_category'].values)

# we use the fitted CountVectorizer to convert the text to vector
X_train_project_grade_category_one_hot = vectorizer_project_grade_category.transform(X_train['project_grade_category'].values)
X_cv_project_grade_category_one_hot = vectorizer_project_grade_category.transform(X_cv['project_grade_category'].values)
X_test_project_grade_category_one_hot = vectorizer_project_grade_category.transform(X_test['project_grade_category'].values)

print("Shape of matrix after one hot encoding ")
print(X_train_project_grade_category_one_hot.shape, y_train.shape)
print(X_cv_project_grade_category_one_hot.shape, y_cv.shape)
print(X_test_project_grade_category_one_hot.shape, y_test.shape)

print("=="100)

Shape of matrix after one hot encoding
(44890, 4) (44890,)
(22110, 4) (22110,)
(33000, 4) (33000,)
```

Encoding numerical features

Encoding numerical feature: Price

```
In [152]: # Encoding price feature
# check this one: https://www.youtube.com/watch?v=0H0qDcIn3Z4&t=5308
# standardization link: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html

"""from sklearn.preprocessing import StandardScaler

# price_standardized = StandardScaler.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ... 399. 287.73 5.5...])
# Reshape your data either using array.reshape(-1, 1)

price_scalar = StandardScaler()
price_scalar.fit(X_train['price'].values.reshape(-1,1)) # finding the mean and standard deviation of this data
print(f"Mean : {price_scalar.mean_[0]}, Standard deviation : {np.sqrt(price_scalar.var_[0])}")

# Now standardize the data with above mean and variance.
X_train_price_standardized = price_scalar.transform(X_train['price'].values.reshape(-1, 1))
X_cv_price_standardized = price_scalar.transform(X_cv['price'].values.reshape(-1, 1))
X_test_price_standardized = price_scalar.transform(X_test['price'].values.reshape(-1, 1))

X_train_price_standardized
X_cv_price_standardized
X_test_price_standardized

print("Shape of matrix after encoding ")
print(X_train_price_standardized.shape, y_train.shape)
print(X_cv_price_standardized.shape, y_cv.shape)
print(X_test_price_standardized.shape, y_test.shape)

print("=="100)

(44890, 1) (44890,)
(22110, 1) (22110,)
(33000, 1) (33000,)
```

Vectorizing : teacher_number_of_previously_posted_projects : numerical

```
In [154]: # Encoding teacher_number_of_previously_posted_projects : numerical feature
# check this one: https://www.youtube.com/watch?v=0H0qDcIn3Z4&t=5308
# standardization link: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html

"""from sklearn.preprocessing import StandardScaler

# price_standardized = StandardScaler.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ... 399. 287.73 5.5...])
# Reshape your data either using array.reshape(-1, 1)

teacher_number_of_previously_posted_projects_scalar = StandardScaler()
teacher_number_of_previously_posted_projects_scalar.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
print(f"Mean : {teacher_number_of_previously_posted_projects_scalar.mean_[0]}, Standard deviation : {np.sqrt(teacher_number_of_previously_posted_projects_scalar.var_[0])}")

# Now standardize the data with above mean and variance.
X_train_teacher_number_of_previously_posted_projects_standardized = teacher_number_of_previously_posted_projects_scalar.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
X_cv_teacher_number_of_previously_posted_projects_standardized = teacher_number_of_previously_posted_projects_scalar.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
X_test_teacher_number_of_previously_posted_projects_standardized = teacher_number_of_previously_posted_projects_scalar.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))

X_train_teacher_number_of_previously_posted_projects_standardized
X_cv_teacher_number_of_previously_posted_projects_standardized
X_test_teacher_number_of_previously_posted_projects_standardized

print("Shape of the matrix after vectorization:")
print(X_train_teacher_number_of_previously_posted_projects_standardized.shape, y_train.shape)
print(X_cv_teacher_number_of_previously_posted_projects_standardized.shape, y_cv.shape)
print(X_test_teacher_number_of_previously_posted_projects_standardized.shape, y_test.shape)

print("=="100)

(44890, 1) (44890,)
(22110, 1) (22110,)
(33000, 1) (33000,)
```

Vectorizing Text data using BOW

Encoding : preprocessed_essays

```
In [155]: print(X_train.shape, y_train.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("=="100)

#Encoding preprocessed essays using BOW
vectorizer_preprocessed_essays = CountVectorizer(min_df=10, ngram_range=(1,4), max_features=5000)
vectorizer_preprocessed_essays.fit(X_train['preprocessed_essays'].values)

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_bow = vectorizer_preprocessed_essays.transform(X_train['preprocessed_essays'].values)
X_cv_essay_bow = vectorizer_preprocessed_essays.transform(X_cv['preprocessed_essays'].values)
X_test_essay_bow = vectorizer_preprocessed_essays.transform(X_test['preprocessed_essays'].values)

print("After vectorizations")
print(X_train_essay_bow.shape, y_train.shape)
print(X_cv_essay_bow.shape, y_cv.shape)
print(X_test_essay_bow.shape, y_test.shape)

print("=="100)

(44890, 21) (44890,)
(22110, 21) (22110,)
(33000, 21) (33000,)

After vectorizations
(44890, 5000) (44890,)
(22110, 5000) (22110,)
(33000, 5000) (33000,)
```

Encoding : preprocessed_title

```
In [157]: # encoding preprocessed title using BOW
print("=="100)

vectorizer_preprocessed_title = CountVectorizer(min_df=10, ngram_range=(1,4), max_features=5000)
vectorizer_preprocessed_title.fit(X_train['preprocessed_title'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_title_bow = vectorizer_preprocessed_title.transform(X_train['preprocessed_title'].values)
X_cv_title_bow = vectorizer_preprocessed_title.transform(X_cv['preprocessed_title'].values)
X_test_title_bow = vectorizer_preprocessed_title.transform(X_test['preprocessed_title'].values)

print("After vectorizations")
print(X_train_title_bow.shape, y_train.shape)
print(X_cv_title_bow.shape, y_cv.shape)
print(X_test_title_bow.shape, y_test.shape)

print("=="100)

(44890, 3779) (44890,)
(22110, 3779) (22110,)
(33000, 3779) (33000,)
```

2.3 Merging all the above features

Preparing Set 1: categorical, numerical features + project_title(BOW) + preprocessed_essay (BOW)

```
In [158]: #we need to merge all the numerical features i.e categorical, text, numerical vectors
#https://stackoverflow.com/questions/59526898/how-to-get-feature-importance-in-naive-bayes#59530697
#https://www.kaggle.com/vsundar/amazon-review-log-nb""
X_train_set1 = hstack((X_train_title_bow, X_train_essay_bow, X_train_categories_one_hot, X_train_sub_categories_one_hot, X_train_school_state_one_hot, X_train_teacher_prefix_one_hot, X_train_project_grade_category_one_hot, X_cv_price_norm, X_cv_clean_category, X_train_price_norm, X_train_teacher_number_of_previously_posted_projects_norm, X_test_set1 = hstack((X_cv_title_bow, X_cv_essay_bow, X_cv_sub_categories_one_hot, X_cv_sub_categories_norm, X_cv_school_state_one_hot, X_cv_teacher_prefix_one_hot, X_cv_project_grade_category_one_hot, X_cv_price_norm, X_cv_clean_category, X_test_price_norm, X_test_teacher_number_of_previously_posted_projects_norm), tocsr()

print("Final Data matrix")
print(X_train_set1.shape, y_train.shape)
print(X_cv_set1.shape, y_cv.shape)
print(X_test_set1.shape, y_test.shape)

print("=="100)

Final Data matrix
(44890, 8880) (44890,)
(22110, 8880) (22110,)
(33000, 8880) (33000,)
```

2.4 Applying NB on different kind of featurization as mentioned in the instructions

Apply NB on different kind of featurization as mentioned in the instructions
For Every model that you work on make sure you do the step 2 and step 3 of instructions

2.4.1 Applying NB on BOW featurization, SET 1

Method 1: Simple for loop (if you are having memory limitations use this)

```
In [159]: def batch_predict(clf, data):
# clf = roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_data_pred = []
tr_loop = data.shape[0] - data.shape[0]*10000
for i in range(0, tr_loop, 10000):
    y_data_pred.extend(clf.predict_proba(data[i:i+10000])[1,:])
# we will be predicting for the last data points
if data.shape[0]*10000 != 0:
    y_data_pred.extend(clf.predict_proba(data[tr_loop:])[1,:])

return y_data_pred

In [160]: import matplotlib.pyplot as plt
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import roc_auc_score

y_true = array, shape = [n_samples] or [n_samples, n_classes]
y_predicted = array, shape = [n_samples] or [n_samples, n_classes]
y_score = array, shape = [n_samples] or [n_samples, n_classes]
Target scores, can either be probability estimates of the positive class, confidence values, or non-threshold measure of decisions (as returned by "decision_function" on some classifiers).
For binary y_true, y_score is supposed to be the score of the class with greater label.

train_auc = []
alpha = [0.001, 0.01, 0.1, 0.5, 1.0, 10.0,]
for i in tqdm(alpha):
    model = MultinomialNB(alpha=i)
    y_train_pred = batch_predict(model, X_train_set1)
    y_cv_pred = batch_predict(model, X_cv_set1)

# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = batch_predict(model, X_train_set1)
y_test_pred = batch_predict(model, X_test_set1)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

plt.plot(train_fpr, train_tpr, label='train AUC')
plt.plot(test_fpr, test_tpr, label='test AUC')
plt.legend()
plt.xlabel('False Positive Rate(FPR)')
plt.ylabel('True Positive Rate(TPR)')
plt.title('ROC CURVE')
plt.grid()
plt.show()

100%|#####| 6/6 [00:01:00:00, 3.58it/s]
```

Testing the performance of the model on test data, plotting ROC Curves

```
In [161]: #here we are choosing best value of alpha based on for loop results
best_alpha = 0.1

In [162]: #https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

model = MultinomialNB(alpha=best_alpha)
model.fit(X_train_set1, y_train)
y_train_pred = batch_predict(model, X_train_set1)
y_test_pred = batch_predict(model, X_test_set1)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

plt.plot(train_fpr, train_tpr, label='train AUC')
plt.plot(test_fpr, test_tpr, label='test AUC')
plt.legend()
plt.xlabel('False Positive Rate(FPR)')
plt.ylabel('True Positive Rate(TPR)')
plt.title('ROC CURVE')
plt.grid()
plt.show()
```

Confusion Matrix

```
In [163]: # we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def find_best_threshold(threshold, fpr, tpr):
    t = threshold
    # tpr(1-fpr) will be maximum if your fpr is very low and tpr is very high
    return t, max(tpr(1-fpr)), max(tpr(1-fpr)), "for threshold", np.round(t,3))

def predict_with_best_t(proba, threshold):
    predictions = []
    for i in range(0, proba.shape[0]):
        if proba[i,1] > threshold:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions

In [164]: #ref link: https://zeelcode.com/python/prettytable/
from prettytable import PrettyTable

print("=="100)

from sklearn.metrics import confusion_matrix
def find_best_threshold(threshold, fpr, tpr):
    t = threshold
    # tpr(1-fpr) will be maximum if your fpr is very low and tpr is very high
    return t, max(tpr(1-fpr)), max(tpr(1-fpr)), "for threshold", np.round(t,3))

tb = PrettyTable()
tb.title = "Train confusion matrix"
tb.field_names = ("Predicted: NO", "Predicted: YES")
tb.add_row(["Actual: NO", tn, fp])
tb.add_row(["Actual: YES", fn, tp])
print(tb)

print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
tb = PrettyTable()
tb.title = "Test confusion matrix"
tb.field_names = ("Predicted: NO", "Predicted: YES")
tb.add_row(["Actual: NO", tn, fp])
tb.add_row(["Actual: YES", fn, tp])
print(tb)

=====
Train confusion matrix
[[ 4603 2200]
 [11570 2651]]
-----+-----+
| Train confusion matrix |
|-----+-----+
| Actual: NO | 4603 | 2200 |
| Actual: YES | 11570 | 2651 |
|-----+-----+
Test confusion matrix
[[ 2943 2058]
 [ 8796 8293]]
-----+-----+
| Test confusion matrix |
|-----+-----+
| Predicted: NO | Predicted: YES |
|-----+-----+
| Actual: NO | 2943 | 2058 |
| Actual: YES | 8796 | 10293 |
|-----+-----+

Feature importance: Top 20 features from Set1

In [165]: # Get feature names from Set1
set1_feature_names = list(project_data.columns)
for i in vectorizer_preprocessed_title.get_feature_names():
    set1_feature_names.append(i)
#Get features corresponding to preprocessed_essays
for i in vectorizer_preprocessed_essays.get_feature_names():
    set1_feature_names.append(i)
#Get features corresponding to clean_categories
for i in vectorizer_clean_categories.get_feature_names():
    set1_feature_names.append(i)
#Get features corresponding to clean_subcategories
for i in vectorizer_clean_subcategories.get_feature_names():
    set1_feature_names.append(i)
#Get features corresponding to school_state
for i in vectorizer_school_state.get_feature_names():
    set1_feature_names.append(i)
#Get features corresponding to teacher_prefix
for i in vectorizer_teacher_prefix.get_feature_names():
    set1_feature_names.append(i)
#Get features corresponding to project_grade_category
for i in vectorizer_project_grade_category.get_feature_names():
    set1_feature_names.append(i)
set1_feature_names.append('price')
set1_feature_names.append('teacher_number_of_previously_posted_projects')

#print(set1_feature_names)
print(len(set1_feature_names))

8880

In [166]: # Get top two features displayed for both the classes
#Code Ref:https://www.kaggle.com/vsundar/amazon-review-log-nb""
neg_class_prob_sorted = (-model.feature_log_prob_[0,:]).argsort()
pos_class_prob_sorted = (-model.feature_log_prob_[1,:]).argsort()

neg_class_features = np.take(set1_feature_names, neg_class_prob_sorted[:20])
pos_class_features = np.take(set1_feature_names, pos_class_prob_sorted[:20])

print("The top 20 features from the positive class are:\n")
print(neg_class_features, neg_class_prob_sorted[:20])

print("The top 20 features from the negative class are:\n")
print(pos_class_features, pos_class_prob_sorted[:20])

The top 20 features from the positive class are:
['students', 'school', 'learning', 'classroom', 'not', 'learn', 'help', 'price', 'many', 'nannan', 'need', 'work', 'reading', 'love', 'able', 'reading', 'teacher_number_of_previously_posted_projects', 'day', 'able', 'come', 'class'] [789 7451 6154 4472 6758 8697 5719 8878 6655 6443 8669 8672 4566 8879 6333 3875 7259 8479 4803 6489]

The top 20 features from the negative class are:
['students', 'school', 'learning', 'classroom', 'not', 'learn', 'help', 'price', 'nannan', 'many', 'need', 'work', 'come', 'teacher_number_of_previously_posted_projects', 'love', 'able', 'reading', 'class'] [789 7451 6154 4472 6758 8697 5719 8878 6655 6443 8669 8672 4566 8879 6333 3875 7259 8479 4803 6489]
```

2.4.2 Applying NB on TFIDF, SET 2

Encoding essay and project title using TFIDF


```
In [167]: # Please write all the code with proper documentation
from sklearn.feature_extraction.text import TfidfVectorizer

print(X_train.shape, y_train.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)

# Encoding preprocessed essays
vectorizer = TfidfVectorizer(min_df=10)
vectorizer.fit(X_train['preprocessed_essays'].values)

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_tfidf = vectorizer.transform(X_train['preprocessed_essays'].values)
X_cv_essay_tfidf = vectorizer.transform(X_cv['preprocessed_essays'].values)
X_test_essay_tfidf = vectorizer.transform(X_test['preprocessed_essays'].values)

print("After vectorizations")
print(X_train_essay_tfidf.shape, y_train.shape)
print(X_cv_essay_tfidf.shape, y_cv.shape)
print(X_test_essay_tfidf.shape, y_test.shape)
print("="*100)

#Encoding preprocessed title
vectorizer = TfidfVectorizer(min_df=10)
vectorizer.fit(X_train['preprocessed_title'].values)

# we use the fitted CountVectorizer to convert the text to vector
X_train_title_tfidf = vectorizer.transform(X_train['preprocessed_title'].values)
X_cv_title_tfidf = vectorizer.transform(X_cv['preprocessed_title'].values)
X_test_title_tfidf = vectorizer.transform(X_test['preprocessed_title'].values)

print("After vectorizations")
print(X_train_title_tfidf.shape, y_train.shape)
print(X_cv_title_tfidf.shape, y_cv.shape)
print(X_test_title_tfidf.shape, y_test.shape)
print("="*100)

(44899, 21) (44899,)
(22110, 21) (22110,)
(33000, 21) (33000,)

After vectorizations
(44899, 11632) (44899,)
(22110, 11632) (22110,)
(33000, 11632) (33000,)

After vectorizations
(44899, 2005) (44899,)
(22110, 2005) (22110,)
(33000, 2005) (33000,)

=====
```

Preparing Set 2: categorical, numerical features + project_title(TFIDF)+ preprocessed_essay (TFIDF)

```
In [168]: # preparing Set2
# we need to merge all the numerical vectors i.e categorical, text, numerical vectors
from scipy.sparse import hstack
X_train_set2 = hstack((X_train_title_tfidf, X_train_essay_tfidf, X_train_categories_one_hot, X_train_sub_categories_one_hot, X_train_school_state_one_hot, X_train_teacher_prefix_one_hot, X_train_project_grade_category_one_hot, X_train_price_norm, X_train_teacher_number_of_previously_posted_projects_norm)).tocsr()
X_cv_set2 = hstack((X_cv_title_tfidf, X_cv_essay_tfidf, X_cv_categories_one_hot, X_cv_sub_categories_one_hot, X_cv_price_norm, X_cv_teacher_number_of_previously_posted_projects_norm)).tocsr()
X_test_set2 = hstack((X_test_title_tfidf, X_test_essay_tfidf, X_test_categories_one_hot, X_test_sub_categories_one_hot, X_test_school_state_one_hot, X_test_teacher_prefix_one_hot, X_test_project_grade_category_one_hot, X_test_price_norm, X_test_teacher_number_of_previously_posted_projects_norm)).tocsr()

print("Final Data matrix")
print(X_train_set2.shape, y_train.shape)
print(X_cv_set2.shape, y_cv.shape)
print(X_test_set2.shape, y_test.shape)
print("="*100)

Final Data matrix
(44899, 13738) (44899,)
(22110, 13738) (22110,)
(33000, 13738) (33000,)

=====
```

Hyperparameter tuning

```
In [169]: import matplotlib.pyplot as plt
from sklearn.metrics import MultinomialNB
from sklearn.metrics import roc_auc_score

y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.

y_score : array, shape = [n_samples] or [n_samples, n_classes]
Target scores, can either be probability estimates of the positive class, confidence values, or non-thresholded measure of decisions (as returned by "decision_function" on some classifiers).
For binary y_true, y_score is supposed to be the score of the class with greater label.

train_auc = []
cv_auc = []
alpha = [0.001, 0.01, 0.1, 0.5, 1.0, 10.0]
for i in tqdm(alpha):
    model = MultinomialNB(alpha=i)
    model.fit(X_train_set2, y_train)

    y_train_pred = batch_predict(model, X_train_set2)
    y_cv_pred = batch_predict(model, X_cv_set2)

    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train, y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

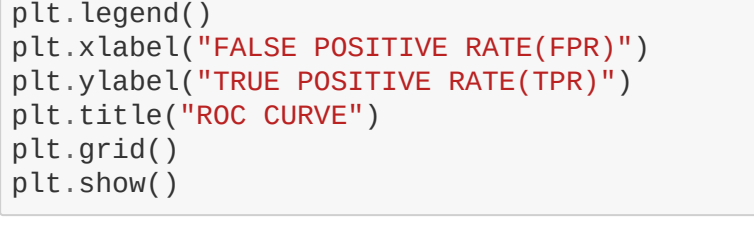
plt.plot(alpha, train_auc, label='Train AUC')
plt.plot(alpha, cv_auc, label='CV AUC')

plt.scatter(alpha, train_auc, label='Train AUC points')
plt.scatter(alpha, cv_auc, label='CV AUC points')

plt.legend()
plt.xlabel("alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```

100% | 8/6 [00:01<00:00, 4.34it/s]

ERROR PLOTS



Testing the performance of the model on test data, plotting ROC Curves

```
In [170]: #here we are choosing best value of alpha based on for loop results
best_alpha = 0.1

In [171]: # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

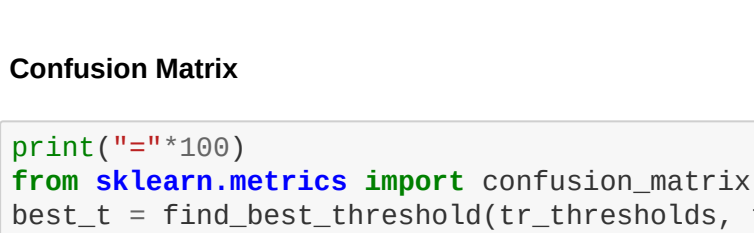
model = MultinomialNB(alpha=best_alpha)
model.fit(X_train_set2, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = batch_predict(model, X_train_set2)
y_test_pred = batch_predict(model, X_test_set2)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

plt.plot(train_fpr, train_tpr, label="train AUC")
plt.plot(test_fpr, test_tpr, label="test AUC")

plt.legend()
plt.xlabel("FALSE POSITIVE RATE(FPR)")
plt.ylabel("TRUE POSITIVE RATE(TPR)")
plt.title("ROC CURVE")
plt.grid()
plt.show()
```



Confusion Matrix

```
In [172]: print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)))
tn,fp,fn,tp=confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)).ravel()

tb = PrettyTable()
tb.title = 'Train confusion matrix'
tb.field_names = ("","Predicted: NO","Predicted: YES")
tb.add_row(["Actual : NO",tn,fp ])
tb.add_row(["Actual : YES",fn,tp])
print(tb)

print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
tn,fp,fn,tp=confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)).ravel()

tb = PrettyTable()
tb.title = 'Test confusion matrix'
tb.field_names = ("","Predicted: NO","Predicted: YES")
tb.add_row(["Actual : NO",tn,fp ])
tb.add_row(["Actual : YES",fn,tp])
print(tb)
```

the maximum value of tpr*(1-fpr) 0.5131340680403463 for threshold 0.847

Train confusion matrix				
[[4862 1941]				
[10741 27346]]				
-----+-----				
Train confusion matrix				
-----+-----				
Predicted: NO Predicted: YES				
-----+-----				
Actual : NO 4862 1941				
Actual : YES 10741 27346				
-----+-----				

Test confusion matrix				
[[2736 2205]				
[8583 19410]]				
-----+-----				
Test confusion matrix				
-----+-----				
Predicted: NO Predicted: YES				
-----+-----				
Actual : NO 2736 2205				
Actual : YES 8583 19410				
-----+-----				

2. Summary

Comparing all models performance using Prettytable library

```
In [174]: #ref link: http://zetcode.com/python/prettytable/
tb = PrettyTable()
tb.title = 'Summary'
tb.field_names = ("Vectorizer", "Model", "Hyperparameter", "AUC")
tb.add_row(["Bow", "Naive Bayes", 0.1, 0.6852])
tb.add_row(["TFIDF", "Naive Bayes", 0.1, 0.6648])
print(tb)

+-----+
| Summary |
+-----+
| Vectorizer | Model | Hyperparameter | AUC |
+-----+
| Bow | Naive Bayes | 0.1 | 0.6852 |
| TFIDF | Naive Bayes | 0.1 | 0.6648 |
+-----+
```

Observation:

- We observed better AUC value i.e 0.6852 using Bow vectorizer as compare to using TFIDF vectorizer.

```
In [ ] :
```