

```

In [1]: # Do ignore warnings
import warnings
warnings.filterwarnings('ignore')
# Do use sqllite3 database
import sqllite3
import numpy as np
import pandas as pd
import string
import nltk
import matplotlib.pyplot as plt
from nltk.stem.porter import PorterStemmer
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer
import re
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import cross_validation
from sklearn.metrics import accuracy_score
from sklearn.cross_validation import cross_val_score
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

In [2]: con = sqllite3.connect('database_sqllite')

filtered_data = pd.read_sqllite_query(""" SELECT * FROM Reviews WHERE Score != 3 """, con)

# Give reviews with Score=3 a positive rating, and reviews with a score=3 a negative rating
def partition(N):
    if N < 3:
        return 'negative'
    return 'positive'

# changing reviews with score less than 3 to be positive and vice-versa
actualScore = filtered_data['Score']
positiveNegative = actualScore.map(partition)
filtered_data['Score'] = positiveNegative



## Text Preprocessing on all data points



In [3]: stop = set(stopwords.words('english')) #set of stopwords
snow = nltk.Stem.SnowballStemmer('english') #initialising the snowball stemmer

def cleanhtml(sentence): #function to clean the word of any html-tags
    return re.compile('<.*>?').sub('', sentence)

cleantext = re.sub(cleanhtml, '', sentence)

return cleantext

def cleanup(sentence): #function to clean the word of any punctuation or special character
    cleaned = re.sub('[!@#$%^&*(){}~`~\']+', '', sentence)
    cleaned = re.sub('[\s|!@#$%^&*(){}~`~\']+', '', cleaned)

    return cleaned

In [4]: #Code for implementing step-by-step the checks mentioned in the pre-processing phase
# This code takes a while to run as it needs to run on 500k sentences
i=0
str=""
final string=[]
all_positive_words=[] # store words from +ve reviews here

```

```
for sent in filtered_data['text'].values:
    filtered_sentence=[]
    #print(sent);
```

```

for w in sent.split():
    for cleaned_words in cleanup(w).split():
        if(cleaned_words.isalpha()) & (len(cleaned_words)>1)):
            if(cleaned_words.lower() not in stop):
                s=(sno.stem(cleaned_words.lower().lower())).encode('utf8')
                filtered_sentence.append(s)
                if filtered_data['Score'].values[i] == 'positive':
                    all_positive_words.append(s) #list of all words used to describe positive text
views

views reviews

            if(filtered_data['Score'].values[i] == 'negative':
                all_negative_words.append(s) #list of all words used to describe negative text

            else:
                continue

        else:
            continue

    #print(filtered_sentence)
    str1 = b" ".join(filtered_sentence) #final string of cleaned words
    #print("*****")
    final_string.append(str1)
    i+=1

```

In [5]: `filtered_data['CleanedText'] = final_string #adding a column of CleanedText which displays the data after pre-processing of review`  
`filtered_data['CleanedText'] = filtered_data['CleanedText'].str.decode('utf-8')`

## Sort the datapoints according to time and take first 100000 points

In [6]: `sorted_data = filtered_data.sort_values(by=['Time'])`  
`3 = sorted_data['Score']`  
`Score = 3`

## Splitting

In [51]: `# HERE WE ARE SPLITTING THE DATA POINTS IN TO 80% TRAIN AND 20% FOR TEST`  
`X_1, X_test, y_1, y_test = cross_validation.train_test_split(sorted_data, Score, test_size=0.2, random_state=0)`  
`# HERE WE ARE AGAIN SPLITTING THE TRAIN DATA IN EARLIER LINE X_1 IN TO 75% TRAINING AND 25% CROSS VALIDATION DATA so we have 60% 20% 20% ratio`  
`X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X_1, y_1, test_size=0.25)`

## BAG OF WORDS

In [52]: `count_vect = CountVectorizer() #in scikit-learn`  
`vec = count_vect.fit(X_tr['CleanedText'].values)`

In [53]: `X_trvec = vec.transform(X_tr['CleanedText'].values)`  
`X_cvvec = vec.transform(X_cv['CleanedText'].values)`  
`X_testvec = vec.transform(X_test['CleanedText'].values)`

In [10]: `#(0.0001, 0.0001, 0.001, 0.01, 0.1, 1, 10)`  
`from sklearn.naive_bayes import MultinomialNB`  
`listacc=[]`  
`alpha = [0.0001, 0.0001, 0.001, 0.001, 0.01, 0.1, 1, 10]`  
`for i in alpha:`  
 `b = MultinomialNB(alpha=float(i))`  
 `# fitting the model on crossvalidation train`  
 `b.fit(X_trvec, y_tr)`  
 `# predict the response on the crossvalidation train`  
 `pred = b.predict(X_cvvec)`  
 `# evaluate CV accuracy`

```
listacc.append(acc)
MSE = [100 - x for x in listacc]
```

```

optimal_k = alpha[MSE.index(min(MSE))]
print('\nThe best alpha is %f.' % optimal_k)

plt.plot(alpha, MSE)

for xy in zip(alpha, np.round(MSE,3)):
    plt.annotate('%s, %s)' % xy, xy=xy, textcoords='data')

plt.xlabel('Number of Neighbors K')
plt.ylabel('Misclassification Error')
plt.show()

print('The misclassification error for each alpha value is : ', np.round(MSE,3))

```

---

```

CV accuracy for alpha = 0.000010 is 90%
CV accuracy for alpha = 0.000100 is 90%
CV accuracy for alpha = 0.001000 is 90%
CV accuracy for alpha = 0.010000 is 90%
CV accuracy for alpha = 0.100000 is 90%
CV accuracy for alpha = 1.000000 is 90%
CV accuracy for alpha = 10.000000 is 89%

The best alpha is 0.010000.

```

the misclassification error for each alpha value is : [ 9.448 9.424 9.379 9.345 9.398 9.36 1 10.676]

```

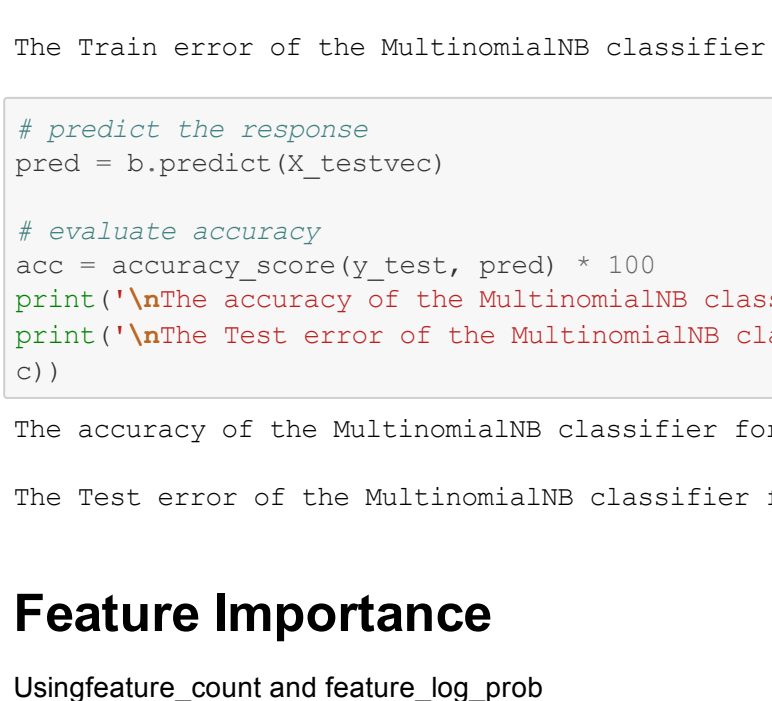
In [11]: b = MultinomialNB(alpha=float(optimal_k))#Using Multinomial nb for best alpha
          s = fit(X_train, y_train)#fitting this for training data
          b.get_params(deep=True)#function for getting parameters

Out[11]: {'alpha': 0.01, 'class_prior': None, 'fit_prior': True}

In [12]: pred = b.predict(X_test)
          acc = accuracy_score(y_test, pred) * 100
          print('\nThe accuracy of the MultinomialNB classifier for alpha = %f is %f%%' % (optimal_k, acc))
          print('\nThe Train error of the MultinomialNB classifier for alpha = %f is %f%%' % (optimal_k, 100-acc))

```

The accuracy of the MultinomialNB classifier for alpha = 0.010000 is 91.03566%



```
pos_sort=sorted(b.feature_count_[1],reverse=True)
neg_sort=sorted(b.feature_count_[0],reverse=True)
```

```
[361]: print("pos img features", pos_sort[0:20])
print("neg img features", neg_sort[0:20])

pos img features: [125452, 103996, 99238, 99167, 97297, 90788, 86462, 83395, 0, 789
16770, 75938, 73907, 69339, 64100, 61227, 48954, 48474, 48174, 47495,
46075,
94075]
neg img features: [29267, 28790, 23962, 18308, 17614, 16123, 15779, 15416, 1392
5, 13537, 13093, 12263, 12226, 11044, 10661, 10652, 10619, 10144, 9
443, 0]

In [362]: pos_sort=sorted(b.feature_log_prob_[1],reverse=True)
neg_sort=sorted(b.feature_log_prob_[0],reverse=True)

In [363]: print("pos img features", pos_sort[0:20])
print("neg img features", neg_sort[0:20])

pos img features: [-4.398344577443124, -4.5188914942563, -6.62546143052712, -6.6261711363818
0, -6.62617113638180, -6.62617113638180, -6.763279394789964, -6.799320148637095, -6.8548583
528926, -6.8821655894065, -6.8943789203764, -8.920174187506145, -8.9837955496643, -5.06
2532308145595, -5.1083788766046, -5.320090964674855, -5.34194294214805, -5.3481522894303,
-5.36343499275147, -5.387703081888]
neg img features: [-4.398344577443124, -4.319917401981268, -4.50436962194316, -6.7725814805440
7, -6.81122942491666, -6.839871395029948, -6.91973000251947, -6.943494943046666, -6.9642262437
447, -6.987849478398, -6.99404611449058, -7.17331208871335, -7.17633361812845, -5.2927
45622995, -5.313242241884285, -5.31413897655807, -5.31742122447076, -5.321392975155365,
-5.342399197070905, -5.43440601261602]

In [30]: def important_features(vectORIZER, classifier, n=20):
class_labels = classifier.classes_
pos_img_name = vectORIZER.get_feature_names()
topn_class1 = sorted(zip(classifier.feature_log_prob_[1], feature_names), reverse=True)[:n]
topn_class2 = sorted(zip(classifier.feature_log_prob_[0], feature_names), reverse=True)[:n]
print("Important words in negative review")
for coef, feat in topn_class1:
print(class_labels[0], coef, feat)
print("\n")
print("Important words in positive review")
for coef, feat in topn_class2:
print(class_labels[1], coef, feat)
```

negative -4.27556236346887 ta  
negative -4.29421328285664 li  
negative -4.48669351695775 pr

negative -4.778228473412305 flavor  
negative -4.67943960113963 would  
negative -4.484076040273313 tri  
negative -4.91821242723155 food  
negative -5.00124308493559 coffee  
negative -5.06531471388948 good  
negative -5.033646087087801 use  
negative -5.140613595700686 buy  
negative -5.15669595760521 get  
negative -5.25429156316126 order  
negative -5.70234913544413 dog  
negative -5.276582935636738 drink  
negative -5.286538943615864 eat  
negative -5.304902179395902 eat  
negative -5.343920645972659 even  
negative -5.40656890595022 bag

---

Important words in positive reviews  
positive -4.39361154973788 like  
positive -4.510193642682598 eat  
positive -4.6181432939678525 love  
positive -4.6197777813836485 flavor  
positive -4.638747031653936 good  
positive -4.710455799791102 great  
positive -4.725020212887623 one  
positive -4.7620013162302845 use  
positive -4.851113668895646 coffee  
positive -4.87813623644986 tri  
positive -4.880398517553657 product  
positive -4.918299272041327 tea  
positive -4.9778832425218 food  
positive -5.05246286141371 get  
positive -5.106293281507648 make  
positive -5.24243972434669 would  
positive -5.33935324131607 dog  
positive -5.533947766842575 eat  
positive -5.53382763742556 wine  
positive -5.5813216407983 buy

consists 4 things • True Positive (TP): Observation is positive, and is predicted to be positive [0,0] • False Negative (FN): Observation is positive, but is predicted negative [0,1] • True Negative (TN): Observation is negative, and is predicted to be negative [1,1] • False Positive (FP): Observation is negative, but is predicted positive [0,1] Index if the data is imbalanced, then we will use these things as accuracy =

```
score = 2*recall/precision+recall + precision These can also be calculated directly from sklearn...
```

```
In [32]: pred = b.predict(X_testvec)
from sklearn.metrics import confusion_matrix
import seaborn as sn
print(confusion_matrix(y_test, pred))
CFM = confusion_matrix(y_test, pred)
df_cm = pd.DataFrame(CFM, range(2), range(2))
plt.figure(figsize=(10,7))
```

```

In [32]:
sns.set(font_scale=1.4)
fig, axs = plt.subplots(2, 2)
sns.heatmap(df_cm, annot=True, cbar=True, annot_kws={"size": 16})

Out[32]:
<matplotlib.axes._subplots.AxesSubplot at 0x210b5a39b0>

```

	0	1
0	1.1e+04	4.8e+03
1	5e+03	8.4e+04

```

In [33]:
TP = CFM[0][0]
FP = CFM[1][0]
FN = CFM[0][1]
TN = CFM[1][1]
P = TP/FP

```

```
print('TP Value is', FN)
print('FP Value is', FN)
print('TN Value is', FN)
print('FN Value is', FN)

TPR = float(TP/P)
FPR = float(FP/P)
TNR = float(TN/N)
FNR = float(FN/N)
print('TPR Value is', TPR)
print('FPR Value is', FPR)
print('TNR Value is', TNR)
print('FNR Value is', FNR)

TF Value is 11345
TF Value is 4999
TN Value is 84010
FN Value is 4809
TFR Value is 0.6941385217816936
FNR Value is 0.3058614782183064
TFR Value is 0.9458561794210698
TNR Value is 0.95413820578930184

In [34]: from sklearn.metrics import accuracy_score, f1_score, precision_score, recall_score
print(accuracy_score(y_test, pred))
print(f1_score(y_test, pred, average="macro"))
print(precision_score(y_test, pred, average="macro"))
print(recall_score(y_test, pred, average="macro"))

0.9047352569895565
0.8125313466176236
0.8199973506013817
0.8230699876377905

TF_IDF

In [35]: tfidf_vect = TfidfVectorizer(ngram_range=(1,3))
final_tf_idf = tfidf_vect.fit(X_train['CleanedText']).values

In [36]: X_train_tf_idf = final_tf_idf.transform(X_train['CleanedText']).values
```

```
X_test_tf_idf = final_tf_idf.transform(X_test['CleanedText'].values)

In [37]: #[0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10]
from sklearn.naive_bayes import MultinomialNB
listacc=[]

alpha= [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10]
for i in alpha:
    nb=MultinomialNB(alpha=alpha)
    listacc.append(nb.fit(X_train_tf_idf, y_train).score(X_test_tf_idf, y_test))
```

```
# fitting the model on crossvalidation train
b.fit(X_tr_tf_idf, y_tr)

# predict the response on the crossvalidation train
pred = b.predict(X_cv_tf_idf)

# evaluate accuracy
acc = accuracy_score(y_cv, pred, normalize=True) * float(100)
print('NaCV accuracy for alpha = %f is %d%%' % (l, acc))
listacc.append(acc)

MSE = [100 - x * acc for x in listacc]

# determining best k
```

```
plt.plot(a
```

```
plt.annotate('(k0, ss)' % xy, xy=xy, textcoords='data')
plt.xlabel('Number of Neighbors K')
plt.ylabel('Misclassification Error')
plt.show()

print("the misclassification error for each alpha is : ", np.round(MSE, 3))
```

CV accuracy for alpha = 0.000010 is 91%

CV accuracy for alpha = 0.000100 is 92%

CV accuracy for alpha = 0.001000 is 92%

CV accuracy for alpha = 0.010000 is 93%

CV accuracy for alpha = 0.100000 is 92%

CV accuracy for alpha = 1.000000 is 84%

CV accuracy for alpha = 10.000000 is 84%

The best alpha is 0.010000.

Number of Neighbors K	Misclassification Error
1	15.563
0.05	8.003
0.001	7.108
0.01	6.308

the misclassification error for each alpha is : [ 8.003 7.663 7.108 6.308 7.511 15.159 15.563]

```
In [38]: b = MultinomialNB(alpha=float(optimal_k))
         s = b.fit(X_tr_tf_idf, y_tr)

In [39]: pred = b.predict(X_tr_tf_idf)
         acc = accuracy_score(y_tr, pred) + 100
```

```
print('The accuracy of the MultinomialNB classifier for alpha = %f is %f' % (alpha, (tpCm1_A, auc)))
print('\nThe Train error of the MultinomialNB classifier for alpha = %f is %f' % (optimal_k, 100-a_c))
```

The accuracy of the MultinomialNB classifier for alpha = 0.010000 is 99.668133%

The Train error of the MultinomialNB classifier for alpha = 0.010000 is 0.331867%

```
in [40]: # predict the response
```

```
# evaluate accuracy
acc = accuracy_score(y_test, pred) + 100
print('The accuracy of the MultinomialNB classifier for alpha = %f is %f%%' % (optimal_k, acc))
print('The Test error of the MultinomialNB classifier for alpha = %f is %f%%' % (optimal_k, 100-acc))

The accuracy of the MultinomialNB classifier for alpha = 0.010000 is 93.7135684

The Test error of the MultinomialNB classifier for alpha = 0.010000 is 6.28643224
```

In [41]:

```
pos_sort=sorted(b.feature_count_[1],reverse=True)
neg_sort=sorted(b.feature_count_[0],reverse=True)
```

[illegible][illegible]

```
In [45]: #We can use the function which we made previously
import_features(final_tf_idf,bw=20)

Important words in negative reviews
negative -6.0471396641414142
negative -6.156938913491921 like
negative -6.225340811202625 product
negative -6.460384925465203 coffee
negative -6.5176613100550392 flavor
negative -6.5234238873987635 would
negative -6.565866495236424 one
negative -6.670588241915221 buy
negative -6.674801591692364 tri
```

negative -6.7379489358578884 tea  
negative -6.738163093222976 food  
negative -6.77931620404261 box  
negative -6.785926304391397 dog  
negative -6.832266512995674 dont  
negative -6.84253697105188 good  
negative -6.860893061176245 disappoint  
negative -6.882239450867597 got  
negative -6.9033438940485197 bag  
negative -6.916416084917896 use

-----  
Important words in positive reviews  
positive -6.178585967445461 love  
positive -6.20840839789137492 great

```
positive -6.2513077809303815 like
positive -6.2573455841619275 eat
positive -6.282535406206723a food
positive -6.2958985982952027 run
positive -6.296838402377595 flavor
positive -6.482133131930014 processor
positive -6.30330977468862 use
positive -6.508541616253829 one
positive -6.5837148505002602 uci
positive -6.622349985849658 dog
positive -6.68468303002641 dog
positive -6.7132183949276885 get
positive -6.734167404949971 make
positive -6.79792111366242 price
```

```
positive -6.81825670109113 buy
positive -6.84333598634462 sell
```

## Different Parameters in MultinomialNB

```
In [46]: #Estimated empirical log probability for each class.
          print(class_log_prior_)
          #Mirrors class_log_prior_ for interpreting MultinomialNB as a linear model.
          print(class_intercept_)
          #Empirical log probability of features given a class, P(x_i|y).
          print(class_feature_log_prob_)
          #Mirrors feature_log_prob_ for interpreting MultinomialNB as a linear model
```

```
print(b.class_count)
print(b.classes_)

[-1.85182048 -0.17072961]
[-0.17072961]
[[-17.47110119 -17.47110119 -17.47110119 ... -17.47110119 -17.47110119
 -17.47110119]
 [-15.85137947 -16.0313027 -17.70201943 ... -15.22673249 -15.79780411]
 [-15.79780411]]
[-15.95113747 -16.03113027 -17.70201943 ... -15.22673249 -15.79780411
 -15.79780411]
[ 49516. 265972.]
['negative' 'positive']
```

```
(47): pred = b.predict(X_test_tf_idf)
      from sklearn.metrics import confusion_matrix
      import seaborn as sn
      print(confusion_matrix(y_test, pred))
      CPM = confusion_matrix(y_test, pred)
      df_cm = pd.DataFrame(CPM, range(2), range(2))
      #plt.figure(figsize=(10,7))
      sn.set(font_scale=1.4) #for label size
      sn.heatmap(df_cm, annot=True, annot_kws={"size": 16})

[10196 5958]
[ 635 88356]

Out[47]: <matplotlib.axes._subplots.AxesSubplot at 0x2111f036d30>
```

	0	1
0	1e+04	6e+03
1	6.5e+02	8.8e+04

```
[48]: TP = CPM[0][0]
      FP = CPM[1][0]
      FN = CPM[0][1]
      TN = CPM[1][1]
      P = TP/FP
      N = TN/FN
      print('TP Value is', TP)
      print('FP Value is', FP)
      print('TN Value is', TN)
      print('FN Value is', FN)

      TPR = float(TP/P)
      FPR = float(FP/P)
      TNR = float(TN/N)
```

```
print('TFR Value is',TFR )
print('FFR Value is',FFR )
print('TN Value is',TN )
print('FN Value is',FN )

TF Value is 10196
FF Value is 653
TN Value is 88356
FN Value is 5958
TFR Value is 0.393810120748456
```