

Multinomial Naive Bayes

Índice

- 1 Entrenamiento del modelo
- 2 Visualización de resultados
- 3 Demo

In [1]:

```
import numpy as np
import matplotlib.pyplot as plt
import scipy.stats as st

from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import confusion_matrix, classification_report, roc_curve
from sklearn.feature_extraction.text import TfidfVectorizer

import pickle
X_train_smt = pickle.load(open("saved_feats/X_train_smt", "rb"))
X_test = pickle.load(open("saved_feats/X_test", "rb"))
y_train_smt = pickle.load(open("saved_feats/y_train_smt", "rb"))
y_test = pickle.load(open("saved_feats/y_test", "rb"))
data_text = pickle.load(open("saved_feats/data_text", "rb"))

# FORMATO DE PLOTS
plt.style.use('bmh')
```

Entrenamiento del modelo

Para ver el estudio realizado de la elección de mejores hiperparametros, consulta el siguiente notebook: [Parameters_NBayes](#)

In [2]:

```
clf_mnb = (MultinomialNB(alpha=1.0, fit_prior=True,
                        class_prior=[0.2, 0.5, 0.93]))
clf_mnb.fit(X_train_smt, y_train_smt)
pred_mnb = clf_mnb.predict(X_test)
pred_prob_mnb = clf_mnb.predict_proba(X_test)
```

Visualización de resultados

In [3]:

```
print('Predicted classes:', clf_mnb.classes_)
print('Average accuracy :', np.mean(pred_mnb == y_test)*100)
print('Train accuracy :', (clf_mnb.score(
    X_train_smt, y_train_smt))*100)
print('Test accuracy :', (clf_mnb.score(
    X_test, y_test))*100)

print('\n CONFUSION MATRIX')
print(confusion_matrix(y_test, pred_mnb))
print('\n CLASSIFICATION REPORT')
print(classification_report(y_test, pred_mnb))
```

Predicted classes: ['Negative' 'Neutral' 'Positive']
Average accuracy : 90.63016718721293
Train accuracy : 87.72431473082233
Test accuracy : 90.63016718721293

CONFUSION MATRIX

```
[[ 34  28  49]
 [ 20  88 142]
 [ 23 248 481]]
```

CLASSIFICATION REPORT

	precision	recall	f1-score	support
Negative	0.44	0.31	0.36	111
Neutral	0.24	0.35	0.29	250
Positive	0.96	0.95	0.95	5082
accuracy			0.91	5443
macro avg	0.55	0.53	0.53	5443
weighted avg	0.92	0.91	0.91	5443

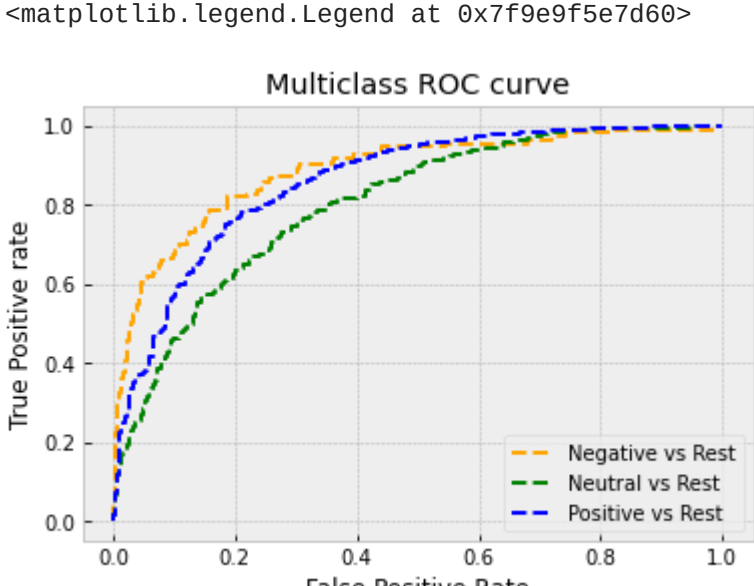
In [4]:

```
fpr = {}
tpr = {}
thresh = {}

fpr[0], tpr[0], thresh[0] = roc_curve(
    y_test, pred_prob_mnb[:, 0], pos_label='Negative')
fpr[1], tpr[1], thresh[1] = roc_curve(
    y_test, pred_prob_mnb[:, 1], pos_label='Neutral')
fpr[2], tpr[2], thresh[2] = roc_curve(
    y_test, pred_prob_mnb[:, 2], pos_label='Positive')

plt.plot(fpr[0], tpr[0], linestyle='--',
         color='orange', label='Negative vs Rest')
plt.plot(fpr[1], tpr[1], linestyle='--',
         color='green', label='Neutral vs Rest')
plt.plot(fpr[2], tpr[2], linestyle='--',
         color='blue', label='Positive vs Rest')
plt.title('Multiclass ROC curve')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive rate')
plt.legend(loc='best')
```

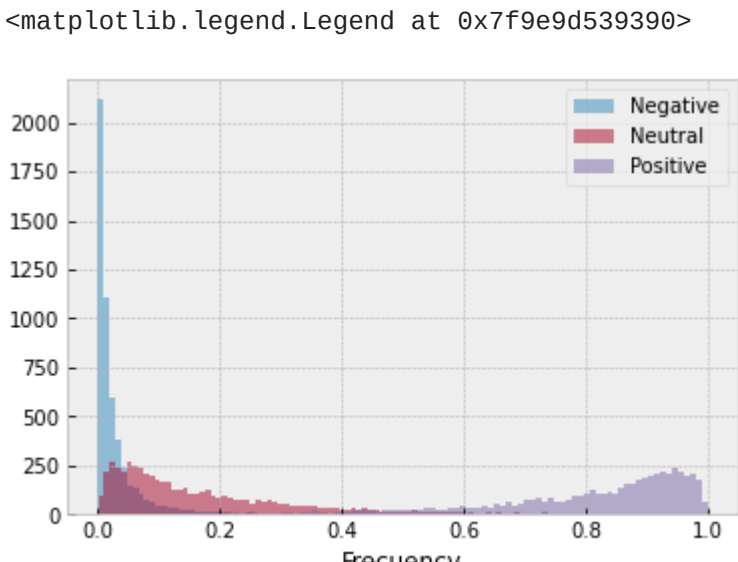
Out[4]:



In [5]:

```
fig = plt.figure()
ax = fig.add_subplot()
ax.hist(pred_prob_mnb[:, 0], 100, alpha=0.5, label="Negative")
ax.hist(pred_prob_mnb[:, 1], 100, alpha=0.5, label="Neutral")
ax.hist(pred_prob_mnb[:, 2], 100, alpha=0.5, label="Positive")
ax.set_xlabel('Probability')
ax.set_xlabel('Frequency')
ax.legend(prop={'size': 10})
```

Out[5]:



In [6]:

```
x = pred_prob_mnb[:, 0]
y = pred_prob_mnb[:, 1]
z = pred_prob_mnb[:, 2]
cmap = plt.get_cmap('magma')

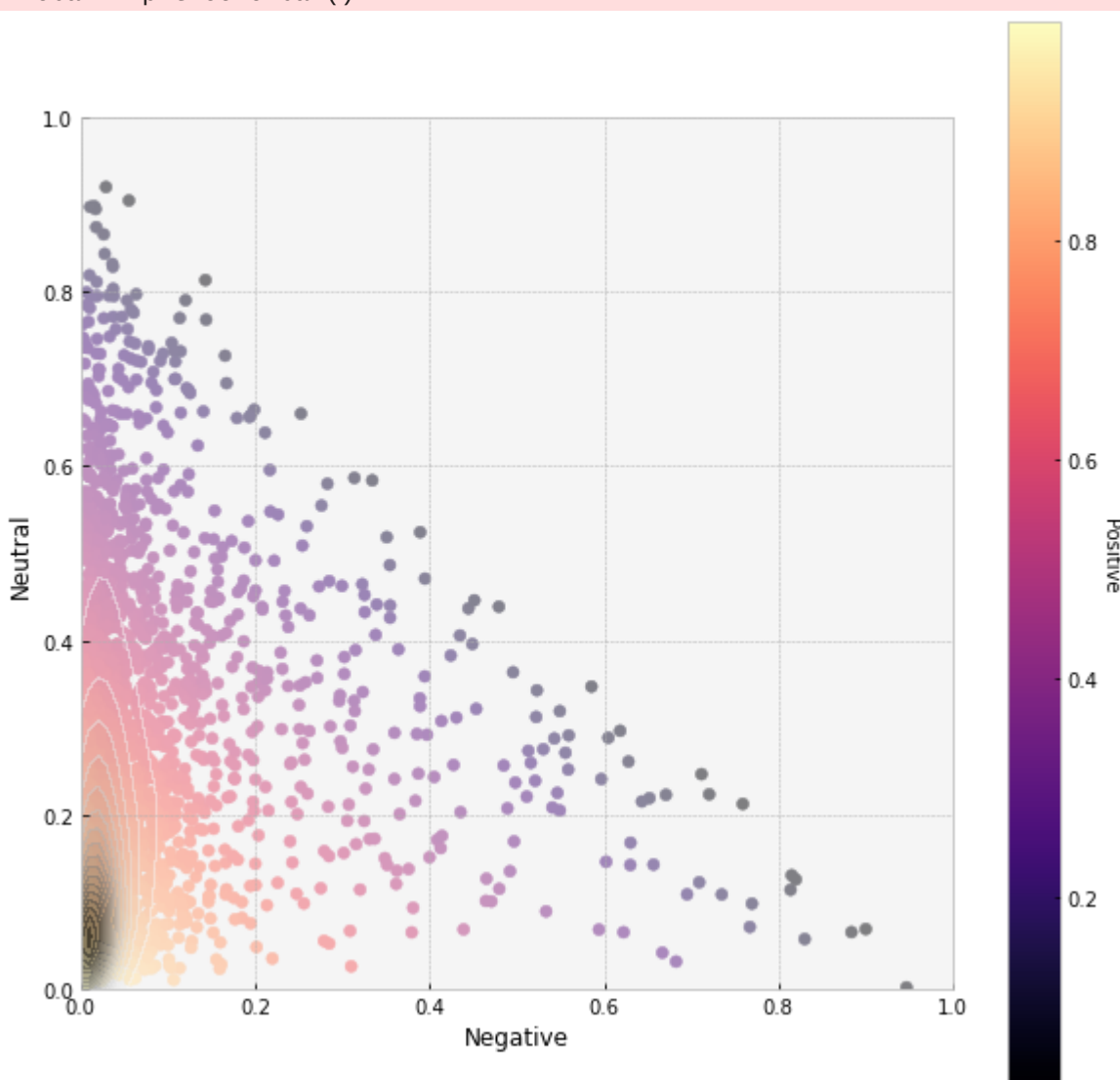
fig = plt.figure(figsize=(10, 10))
ax = fig.add_subplot()
ax.set_xlabel('Negative')
ax.set_ylabel('Neutral')
plt.scatter(x, y, c=z, cmap=cmap)
cbar = plt.colorbar()
cbar.set_label('Positive', rotation=270, fontsize=10, labelpad=10.0)
ax.set_aspect('equal', adjustable='box')

xx, yy = np.mgrid[0:1:100j, 0:1:100j]
positions = np.vstack([xx.ravel(), yy.ravel()])
values = np.vstack([x, y])
kernel = st.gaussian_kde(values)
f = np.reshape(kernel(positions).T, xx.shape)

plt.contourf(xx, yy, f, levels=20, cmap=plt.get_cmap('Greys'), alpha=0.5)
plt.show()
```

/tmp/ipykernel_27019/2290017891.py:11: MatplotlibDeprecationWarning: Auto-removal of grids by pcolor() and pcmesh() is deprecated since 3.5 and will be removed two minor releases later; please call grid(False) first.

```
cbar = plt.colorbar()
```



Demo

In [7]:

```
# CountVectorizer followed by TfidfTransformer
tfidf = TfidfVectorizer(use_idf=True, norm='l2',
                        smooth_idf=True, sublinear_tf=False, max_features=5000)
tfidf.fit_transform(data_text)

# Input sentence
text = input("Enter sentence: ")
arr_text = np.array([text])
tfidf_arr_text = tfidf.transform(arr_text)

# Output
print(clf_mnb.predict(tfidf_arr_text))
```

Enter sentence: Product description is accurate
['Neutral']