NLTK Classifier

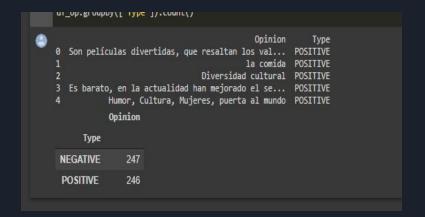
Daniel Steven Vargas Guzman

A	В	С	D	Е
Marca temporal	¿En promedio, cuán	Mencione los aspect	Mencione los aspect	Asigne una punt
22/11/2018 15:07:24	4	Son películas diverti	Siempre muestran la	
22/11/2018 18:19:23	2	la comida	los costos	
22/11/2018 18:19:24	2	Diversidad cultural	películas bélicas de	
22/11/2018 18:20:47	2	Es barato, en la actu	cuando se va la luz r	
22/11/2018 18:21:01	1	Humor, Cultura, Muje	Groserías y más Gros	
22/11/2018 18:21:02	7	El cine en Colombia	La mayor debilidad (
22/11/2018 18:21:20	1	Que es muy comico,	Falta d calidad en p	
22/11/2018 18:21:50	1	Da a Conocer la real	Muchas veces están	
22/11/2018 18:21:57	4	Con pocos recursos s	La falta de recursos	
22/11/2018 18:22:04	2	comedia, diversidad	baja producción de :	
22/11/2018 18:22:18	4	Es creativo, divertido	Con frecuencia resal	
9/8/2019 11:19:01	4	Me gustan las última	Hay bajo presupuesto	
9/8/2019 11:24:29	9	El cine mejora día a	Falta realizar más in	
1/8/2019 19:25:23	2	la forma como mues	al salir al exterior se	
21/8/2019 19:25:29	2	No tiene ningún aspo	Solo muestra aspecti	
21/8/2019 19:25:38	1	Esta bien catalogado	no son de un gran pr	
21/8/2019 19:25:41	1	Enfoque ambiental	Drogas y putas	
21/8/2019 19:25:43	1	La comodidad de la:	La rapidez en iniciar	
21/8/2019 19:25:46	1	La representacion de	La falta de tramas de	
21/8/2019 19:25:46	4	Bueno, Mejoro la tra	Experimentar nuevas	
21/8/2019 19:25:52	4	temas cotidianos de	se resalta mucho el f	
21/8/2019 19:25:58	2	La producción y la tr	Que las grandes proc	
21/8/2019 19:26:01	2	El buen humor	Poco Interesante	
21/8/2019 19:26:07	1	Ultimamente hay pro	Referencias a historia	
21/8/2019 19:26:12	1	Argumento, comedia	efectos especiales y	
21/8/2019 19:26:24	3	He visto algunas peli	Hay algunas que glo	
21/8/2019 19:26:28	1		Gran parte de las pro	
21/8/2019 19:26:43	3	la mayoría son hech	hacen muchas pelic	
21/8/2019 19:26:47		Las expresiones de c		
21/8/2019 19:26:47		La forma de criticar		

Data information

We have a large amount of data represented as opinions about Colombian cinema, at this point, aspects such as good and bad characteristics of this topic in question remain, finally the information that is issued regarding Colombian cinema (quality)

Methodology



we capture the information required in the tcv to obtain the good and bad opinions with the classification if it is a good or bad perception

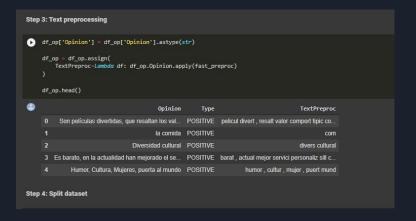
1. Define a text preprocessing data

```
Step 2: Define a text preprocessing function
          n nltk.tokenize import word_tokenize
n nltk.stem import SnowballStemmer
        rom nltk.corpus import stopwords
     nltk.download('stopwords')
     stemmer = SnowballStemmer('spanish')
     nltk.download('punkt')
     stop words = set(stopwords.words('spanish'))
     stop words = stop words.union(set(['me', 'le', 'da', 'mi', 'su', 'ha', 'he', 'ya', 'un', 'una', 'es', 'del', 'las', 'los', 'en', 'que', 'y', 'la', 'de']))
      def fast_preproc(text):
       words = word tokenize(text.lower())
       words = [word for word in words if not word in stop_words]
words = [stemmer.stem(word) for word in words]
         text = " ".join(str(word) for word in words)
              pt Exception as e:
     print(fast preproc("La que me gusta"))
 [nltk_data] Downloading package stopwords to /root/nltk_data...
      [nltk_data] Unzipping corpora/stopwords.zip.
     [nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Unzipping tokenizers/punkt.zip.
```

We define the function of word processing

We obtain all the "stop words" to be able to make the union with some words defined for the text processing, for this point the information is tokenized, allowing to have a perception of which are the words that should be classified with NLTK library

2. Text Preprocessing



At this point we transform the information in such a way that it is recognized as a string and the data is preprocessed with the functions that are expressed in the code.

4. Split information

Step 4: Split dataset

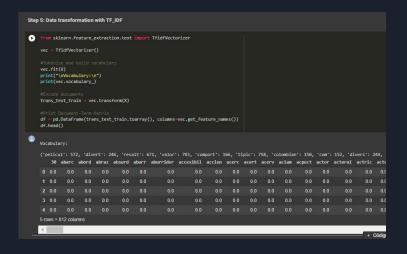


#Split dataset
X = df_op['TextPreproc']

Y = df_op['Type']

we divide the information with the split function

Data transformation with TF_IDF



We carry out the vectorization of the information to convert it into a vocabulary, this with the vectorized transformation implies the immersion of the information in a dataframe

6. Split test-train datasets

Step 6: Split Test-Train Datasets [] from sklearn.model_selection import train_test_split X_train, X_test, Y_train, Y_test = train_test_split(df, Y, test_size=0.15) print("Train and test shapes:\n") print("Train. X: " + str(X_train.shape) + " - Y:" + str(Y_train.shape)) print("Test. X: " + str(X_test.shape) + " - Y:" + str(Y_test.shape)) Train and test shapes: Train. X: (419, 812) - Y:(419,) Test. X: (75, 812) - Y:(75,)

We carry out the training of the previously vectorized dataframe, to adjust a correct interweaving of the model for the classification of the opinions that have been transformed

7. Built Classifier

```
Step 7. Build a Classifier
     from sklearn.linear_model import LogisticRegression
    classifier = LogisticRegression(solver='lbfgs')
    classifier.fit(X_train, Y_train)
    from sklearn.metrics import classification report, confusion matrix
    y pred = classifier.predict(X test)
    print("\nMatriz de confusión:\n")
    print(confusion_matrix(Y_test,y_pred))
    print("\nEstadisticas del clasificador:\n")
    print(classification_report(Y_test,y_pred))
    Matriz de confusión:
    [[28 10]
     ſ 2 3511
    Estadisticas del clasificador:
                  precision
                               recall f1-score support
        NEGATIVE
                       0.93
                                 0.74
                                           0.82
        POSITIVE
                                           0.85
        accuracy
                                           0.84
                       0.86
                                 0.84
                                           0.84
        macro avg
    weighted avg
                       0.86
```

At this point, the classifier is built where all the opinions are obtained together with the sklearn.metrics library, showing their statistical values for each of the trained points and their level of precision.

8. Predict new instance

```
Step 7. Predict New Instance

[ ] new_opinion = ["We parecen malas por la razon de que no llevan un hilo conductor claro"]

#new_opinion = ["We gustan porque son chistosas"]

Xt_new = [fast_preproc(str(new_opinion))]

trans_new_doc = vec.transform(Xt_new) #Use same TfIdfVectorizer

print("\nPredicted result: " + str(classifier.predict(trans_new_doc)))

Predicted result: ['NEGATIVE']

Step 8. Clustering Text
```

Previously trained the model with the pertinent libraries, we used the opinions as input, to visualize their classification according to the model, implying what the perception about Colombian cinema is, in this case the text in Spanish was used "Me parecen malas por la razon de que no llevan un hilo conductor claro" and the result of the perception was negative

9. Clustering

Step 8. Clustering Text

```
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.cluster import KMeans
from sklearn.decomposition import PCA
import matplotlib.pyplot as plt
import numpy as np
X = df_op['TextPreproc']
vec = TfidfVectorizer()
vec.fit(X)
trans_text_train = vec.transform(X)
cls = KMeans(n clusters=2, random state=0)
cls.fit(trans_text_train)
cls.predict(trans_text_train)
print(cls.labels ) #0: First cluster, 1: Second cluster, 2: Third cluster...
print(cls.labels .shape)
print(sum(cls.labels [0:246]))
print(sum(cls.labels [247:494]))
pca = PCA(n_components=2, random_state=0)
reduced_features = pca.fit_transform(trans_text_train.toarray())
reduced_cluster_centers = pca.transform(cls.cluster_centers_)
plt.scatter(reduced_features[:,0], reduced_features[:,1], c=cls.predict(trans_text_train))
plt.scatter(reduced cluster centers[:, 0], reduced cluster centers[:,1], marker='x', s=150, c='b')
```

With the training of the model and its previous validation, it is now able to interpret the text and classify it so that it can understand if it is a good or bad opinion according to the model that was built with the survey of perceptions of Colombian cinema

1. Clustering result

```
pit.scatter(reduced_cluster_centers[:, v], reduced_cluster_centers[:,1], marker=[X', S=150, C=100]
(724,)
<matplotlib.collections.PathCollection at 0x7f23f77d3190>
 0.6
 0.4
 -0.2
 -0.4
```

Conclusion

In conclusion, with the definition of a classification model (NLTK) for the preprocessing of data (text), we can classify them according to a training model that emits a predictive response about the possible result of the opinion of perception about Colombian cinema.