



UANL

UNIVERSIDAD AUTÓNOMA DE NUEVO LEÓN



FCFM

FACULTAD DE CIENCIAS FÍSICO MATEMÁTICAS

Universidad Autónoma de Nuevo León

Licenciatura en Ciencias Computacionales

A14: Programando K-Nearest-Neighbor en Python

Materia: Inteligencia Artificial

Maestro: Luis Ángel Gutiérrez Rodríguez

Rebeca Jaramillo Camarillo

Matrícula: 2132988

Grupo: 031

Fecha: 30 de marzo de 2025

1. Introducción

El algoritmo de k vecinos más cercanos, también conocido como KNN o k-NN, es un clasificador de aprendizaje supervisado no paramétrico, que utiliza la proximidad para hacer clasificaciones o predicciones sobre la agrupación de un punto de datos individual.

Si bien se puede usar para problemas de regresión o clasificación, generalmente se usa como un algoritmo de clasificación, partiendo de la suposición de que se pueden encontrar puntos similares cerca uno del otro.

Para los problemas de clasificación, se asigna una etiqueta de clase sobre la base de un voto mayoritario, es decir, se utiliza la etiqueta que se representa con más frecuencia alrededor de un punto de datos determinado. Si bien esto técnicamente se considera "voto por mayoría", el término "voto por mayoría" se usa más comúnmente en la literatura. La distinción entre estas terminologías es que "voto mayoritario" técnicamente requiere una mayoría superior al 50 %, lo que funciona principalmente cuando solo hay dos categorías. Cuando tiene varias clases, por ejemplo, cuatro categorías, no necesita necesariamente el 50 % de los votos para llegar a una conclusión sobre una clase; puede asignar una etiqueta de clase con un voto superior al 25 %.

2. Metodología

A continuación, se muestran los pasos para resolver el ejercicio:

Importar librerías necesarias

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from matplotlib.colors import ListedColormap
5 import matplotlib.patches as mpatches
6 import seaborn as sb
7 from sklearn.model_selection import train_test_split
8 from sklearn.preprocessing import MinMaxScaler
9 from sklearn.neighbors import KNeighborsClassifier
10 from sklearn.metrics import classification_report
11 from sklearn.metrics import confusion_matrix
```

2.1. Leer el archivo csv y cargarlo como un dataset de Pandas

```
1 dataframe = pd.read_csv(r"reviews_sentiment.csv", sep=';')
2 print(dataframe.head(10))
3 print(dataframe.describe())
```

Visualización del archivo

```
1 dataframe.hist()
2 plt.show()
3
4 print(dataframe.groupby('Star Rating').size())
5
6 sb.countplot(x='Star Rating', data=dataframe)
7 plt.show()
8 sb.countplot(x='wordcount', data=dataframe)
9 plt.show()
```

2.2. Preparar el dataset

```
1 X = dataframe[['wordcount', 'sentimentValue']].values
2 y = dataframe['Star Rating'].values
3
4 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)
5 scaler = MinMaxScaler()
6 X_train = scaler.fit_transform(X_train)
7 X_test = scaler.transform(X_test)
```

2.3. Crear el Modelo

```
1 n_neighbors = 7
2
3 knn = KNeighborsClassifier(n_neighbors)
4 knn.fit(X_train, y_train)
5 print('Accuracy of K-NN classifier on training set: {:.2f}'
6       .format(knn.score(X_train, y_train)))
7 print('Accuracy of K-NN classifier on test set: {:.2f}'
8       .format(knn.score(X_test, y_test)))
9
10 ## Precision del modelo
11 pred = knn.predict(X_test)
12 print(confusion_matrix(y_test, pred))
13 print(classification_report(y_test, pred))
```

2.4. Gráfica de la Clasificación Obtenida

```
1 h = .02 # step size in the mesh
2
3 # Create color maps
4 cmap_light = ListedColormap(['#FFAAAA', '#ffcc99', '#ffffb3', '#b3ffff', '#
   c2f0c2'])
5 cmap_bold = ListedColormap(['#FF0000', '#ff9933', '#FFFF00', '#00ffff', '#00
   FF00'])
6
7 # we create an instance of Neighbours Classifier and fit the data.
8 clf = KNeighborsClassifier(n_neighbors, weights='distance')
9 clf.fit(X, y)
10
11 # Plot the decision boundary. For that, we will assign a color to each
12 # point in the mesh [x_min, x_max]x[y_min, y_max].
13 x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
14 y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
15 xx, yy = np.meshgrid(np.arange(x_min, x_max, h),
16                       np.arange(y_min, y_max, h))
17 Z = clf.predict(np.c_[xx.ravel(), yy.ravel()])
18
19 # Put the result into a color plot
20 Z = Z.reshape(xx.shape)
21 plt.figure()
22 plt.pcolormesh(xx, yy, Z, cmap=cmap_light)
23
24 # Plot also the training points
25 plt.scatter(X[:, 0], X[:, 1], c=y, cmap=cmap_bold,
26            edgecolor='k', s=20)
27 plt.xlim(xx.min(), xx.max())
28 plt.ylim(yy.min(), yy.max())
29
30 patch0 = mpatches.Patch(color='#FF0000', label='1')
31 patch1 = mpatches.Patch(color='#ff9933', label='2')
32 patch2 = mpatches.Patch(color='#FFFF00', label='3')
33 patch3 = mpatches.Patch(color='#00ffff', label='4')
34 patch4 = mpatches.Patch(color='#00FF00', label='5')
35 plt.legend(handles=[patch0, patch1, patch2, patch3, patch4])
36
37
38 plt.title("5-Class classification (k = %i, weights = '%s') "
39          % (n_neighbors, 'distance'))
40
41 plt.show()
```

2.5. Elegir el mejor valor de k

```
1 k_range = range(1, 20)
2 scores = []
3 for k in k_range:
4     knn = KNeighborsClassifier(n_neighbors = k)
5     knn.fit(X_train, y_train)
6     scores.append(knn.score(X_test, y_test))
7 plt.figure()
8 plt.xlabel('k')
9 plt.ylabel('accuracy')
10 plt.scatter(k_range, scores)
11 plt.xticks([0,5,10,15,20])
12 plt.show()
```

2.6. Clasificar ó Predecir nuevas muestras

```
1 print(clf.predict([[5, 1.0]]))
2 print(clf.predict_proba([[20, 0.0]]))
```

3. Resultados

Visualización del archivo

Se observa que la distribución de “estrellas” no está balanceada, la gráfica de Valores de Sentimientos parece una campana movida levemente hacia la derecha del cero y la cantidad de palabras se centra sobre todo de 0 a 10.

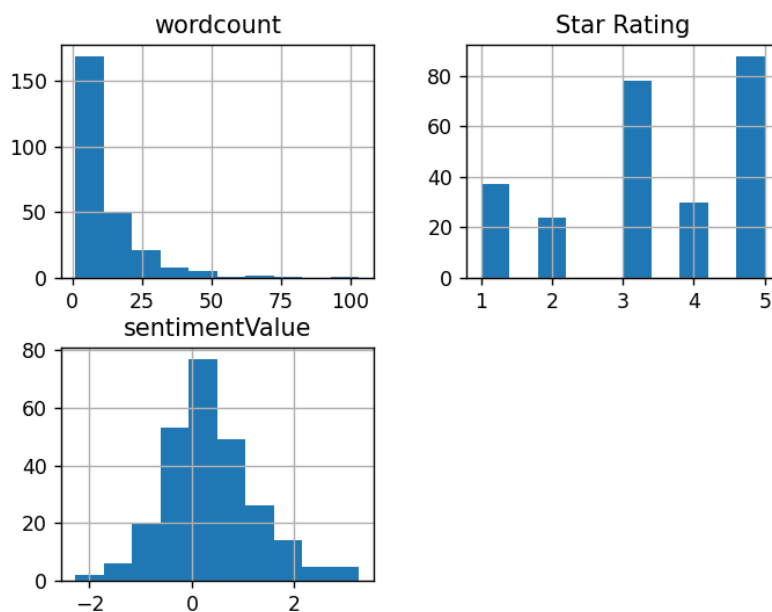


Figura 1: Gráficas de datos de entrada

Después se observa cuantas Valoraciones de Estrellas realmente se tienen, y se concluye que hay una mayoría de valoración de 3 y 5 estrellas.

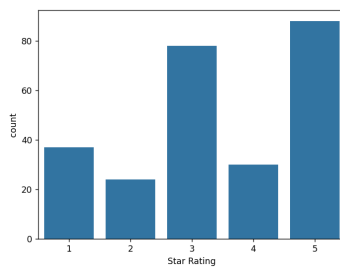


Figura 2: Gráficas de datos de entrada

También se grafica la mayoría de palabras y se observa que la mayoría de reviews tienen entre 1 y 10 palabras.

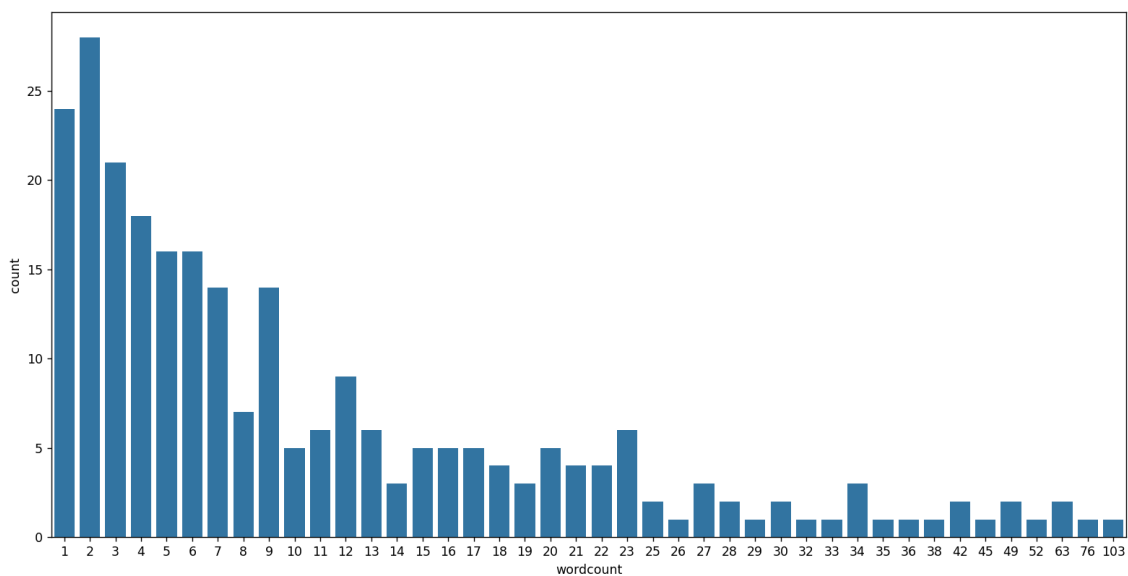


Figura 3: Gráficas de datos de entrada

Gráfica de la Clasificación Obtenida

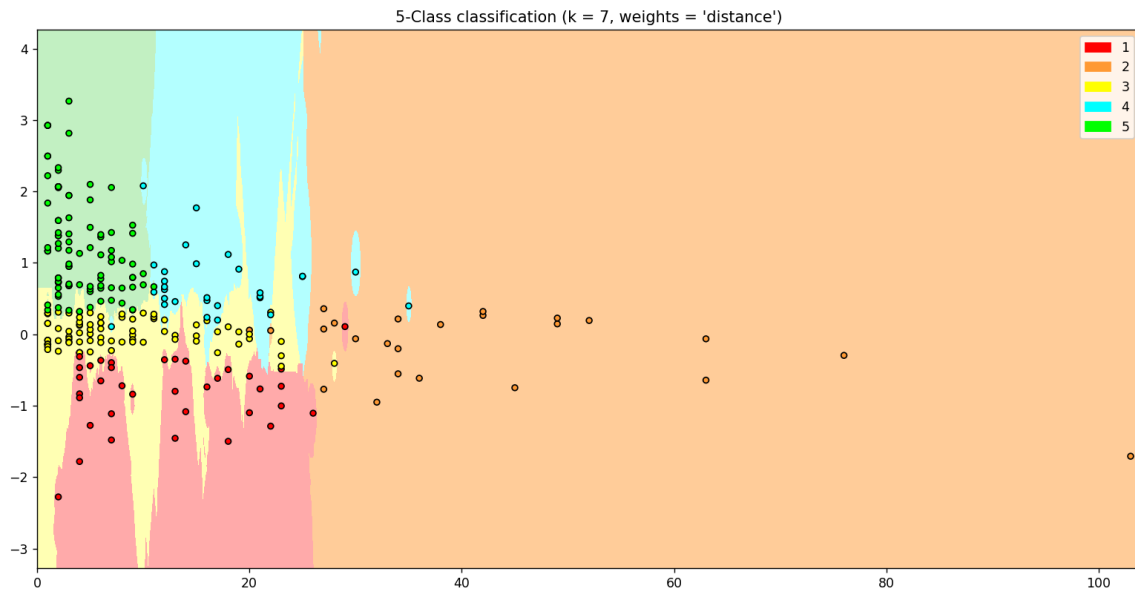


Figura 4: Gráficas de la Clasificación Obtenida

Algunos datos que se podrían "leer" de la gráfica son los siguientes:

- Los usuarios que ponen 1 estrella tienen sentimiento negativo y hasta 25 palabras.
- Los usuarios que ponen 2 estrellas dan muchas explicaciones (hasta 100 palabras) y su sentimiento puede variar entre negativo y algo positivo.
- Los usuarios que ponen 3 estrellas son bastante neutrales en sentimientos, puesto que están en torno al cero y hasta unas 25 palabras.
- Los usuarios que dan 5 estrellas son bastante positivos (de 0, 5 en adelante, aproximadamente) y ponen pocas palabras (hasta 10).

El mejor valor de k

En la siguiente gráfica se muestra de donde se obtiene el valor de "k" utilizado en el código. Como resultado, se puede ver que el mejor valor de "k" es entre 7 y 14.

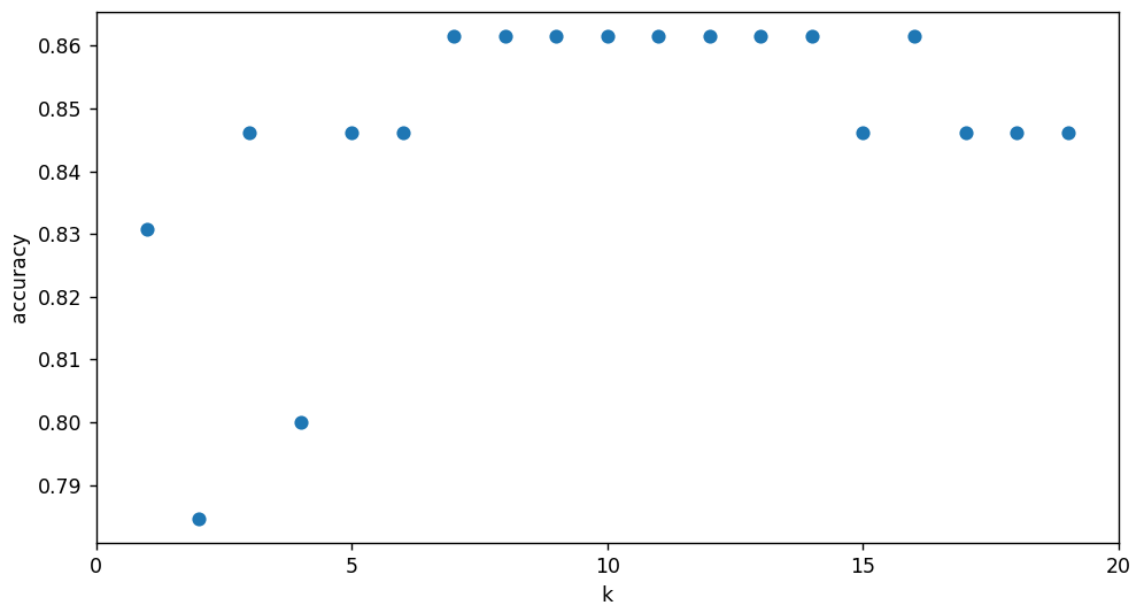


Figura 5: Gráficas de datos de entrada

Clasificar ó Predecir nuevas muestras

La siguiente imagen muestra la clasificación de una review de 5 palabras y sentimiento 1, es decir, 5 estrellas. Otra opción seria ver la probabilidad de obtener 1, 2, 3, 4 o 5 estrellas de una review de 20 palabras y sentimiento 0, que seria de 3 estrellas con una probabilidad del 97 %.

```
[5]  
[[0.00381998 0.02520212 0.97097789 0.         0.         ]]
```

Figura 6: Gráficas de datos de entrada

4. Conclusión

Este ejercicio me permitió aplicar el algoritmo KNN (K-Nearest Neighbors) para clasificar reseñas según su calificación de estrellas, utilizando como características el conteo de palabras (word-count) y el valor de sentimiento (sentimentValue). Desde mi punto de vista, el método KNN es requiere un bien análisis de datos, ya que es importante una buena elección del parámetro "k" y es necesario el uso de una mayor cantidad de datos y características si se quiere lograr un omdelo con una precisión.