

# K-Vecinos más cercanos

Versión : 1.2, Para uso educativo  
Autor : Luis Beltran Palma Ttito  
Lugar : Cusco, Perú, 2023.  
Proposito : K-NN en la clasificación de dígitos manuscritos

## 1. PASO PREVIOS

### Importanci3n de librerías

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 import seaborn as sns
5 import random
6
7 # Librería: scikit learn, "sklearn"
8
9 # Módulo train_test_split: separaci3n de datos para entrenamiento y test
10 from sklearn.model_selection import train_test_split
11 # Módulo KNeighborsClassifier: Implementa K-NN para clasificaci3n
12 from sklearn.neighbors import KNeighborsClassifier
13 # Módulo classification_report: Métricas
14 from sklearn.metrics import classification_report
15 # Módulo accuracy_score: métrica de exactitud
16 from sklearn.metrics import accuracy_score
17 # Módulo confusion_matrix: matriz de confusi3n
18 from sklearn.metrics import confusion_matrix
```

### Importaci3n de datos

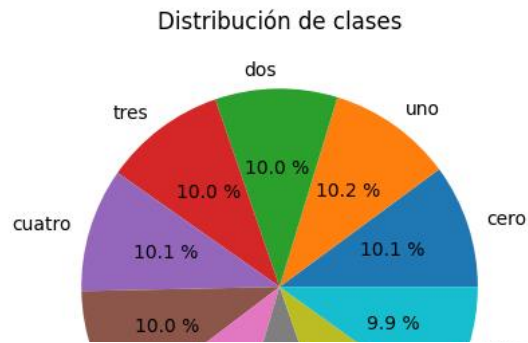
```
1 # leer el dataset
2 datos = pd.read_csv("digito16x16.csv", sep=',', header=0)
3 # muestras los primero 5 valores
4 datos.head()
```

	p1	p2	p3	p4	p5	p6	p7	p8	p9	p10	...	p248	p249	p250	p251	p252	p253	p254	p255	p256	digito
0	0	0	0	0	0	0	1	1	1	1	...	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	1	1	1	1	1	...	1	1	1	1	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	1	...	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	1	1	1	1	...	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	1	...	0	0	0	0	0	0	0	0	0	0

5 rows × 257 columns

### Distribuci3n de clases

```
1 # mostrar distribuci3n de clases mediante un gráfico de pye
2 distribucion = datos.groupby('digito').size()
3 plt.pie(distribucion, labels = ['cero', 'uno', 'dos', 'tres', 'cuatro', 'cinco', 'seis', 'siete', 'ocho', 'nueve'], autopct="%0.1f %%")
4 plt.title('Distribuci3n de clases')
5 plt.show()
6
```



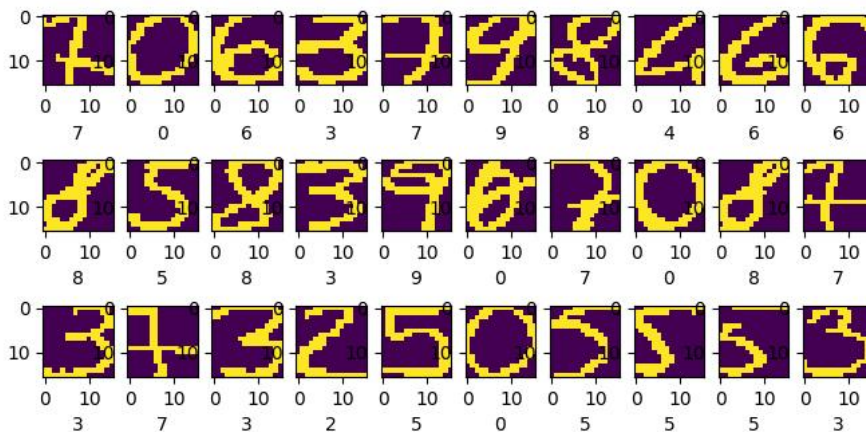
## Visualización de datos



```

1 #visualizar algunos datos aleatoriamente
2 X = datos.drop(['digito'], axis=1)
3 y = datos['digito']
4
5 ran = datos.shape[0]-1
6 fig, ax = plt.subplots(3, 10, figsize=(8,4))
7 for k in range(3):
8     for j in range(10):
9         azar = random.sample(range(ran),1)
10        img = np.array(X.iloc[azar].values)
11        img = np.array(img).reshape(16,16)
12        ax[k,j].imshow(img)
13        ax[k,j].set_xlabel(y[azar[0]])
14 plt.show()

```



## 2. SEPARACIÓN DE DATOS PARA TRAINIG Y TESTING

```

1 # Separación de datos en entrada y salida
2 X = datos.drop(['digito'], axis=1)
3 y = datos['digito']
4
5 # Separación de datos para entrenamiento y test
6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.30, random_state=5)

```

## 3. Entrenamiento y Test de K-NN

```

1 # Definición del modelo
2 Modelo = KNeighborsClassifier(n_neighbors = 5)
3 # Entrenamiento del modelo
4 Modelo.fit(X_train, y_train)
5 # Test del modelo
6 Exactitud = Modelo.score(X_test, y_test)
7 print('Exactitud de K-NN, k=3, :', Exactitud)

```

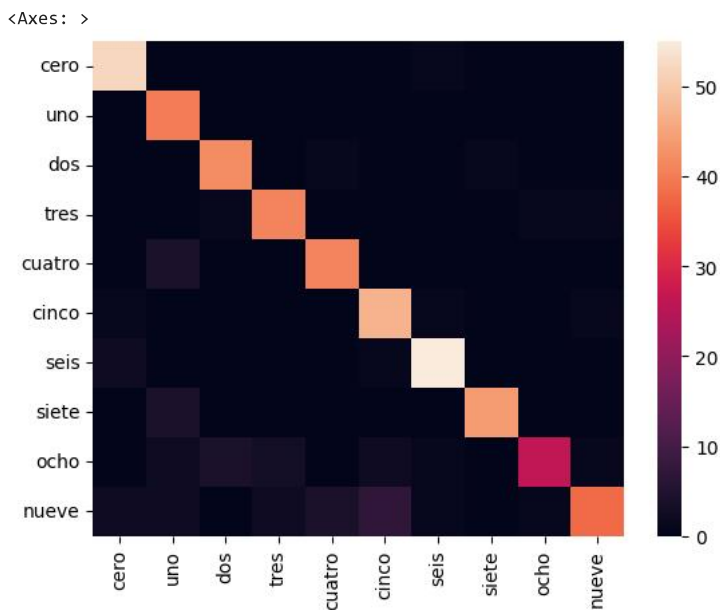
Exactitud de K-NN, k=3, : 0.891213389121339

## ▼ Matriz de confusión y exactitud

```
1 # Test del modelo
2 pred = Modelo.predict(X_test)
3 # Matriz de confusión
4 matriz = confusion_matrix(y_test, pred)
5 print(matriz)
```

```
[[52  0  0  0  0  0  1  0  0  0]
 [ 0 40  0  0  0  0  0  0  0  0]
 [ 0  0 42  0  1  0  0  1  0  0]
 [ 0  0  1 41  0  0  0  0  1  1]
 [ 0  4  0  0 41  0  0  0  0  0]
 [ 1  0  0  0  0 47  1  0  0  1]
 [ 2  0  0  0  0  1 55  0  0  0]
 [ 0  4  0  0  0  0  0 44  0  0]
 [ 0  2  4  3  0  2  1  0 26  1]
 [ 2  2  0  2  4  7  1  0  1 38]]
```

```
1 #visualización de la matriz de confusión como un mapa de calor
2 categorias = ['cero','uno','dos','tres','cuatro','cinco','seis','siete','ocho','nueve']
3 sns.heatmap(matriz, yticklabels=categorias, xticklabels=categorias)
```



```
1 from sklearn.metrics import accuracy_score
2 # Cálculo de accuracy (exactitud)
3 accuracy = accuracy_score(y_test, pred)
4 print('Exactitud: ',accuracy)
```

Exactitud: 0.891213389121339

## ▼ 4. PREDICCIÓN DE NUEVOS CASOS

```
1 nuevo = np.array([[0,0,0,0,0,0,0,0,1,1,1,1,1,1,1,1,
2                   0,0,0,0,0,0,0,1,1,1,1,1,1,0,1,1,
3                   0,0,0,0,0,0,1,1,1,0,0,0,0,0,1,1,
4                   0,0,0,0,0,1,1,1,1,0,0,0,0,0,1,0,
5                   0,0,0,0,1,1,1,0,0,0,0,0,0,0,0,0,
6                   0,0,0,1,1,1,0,1,0,1,0,1,0,1,0,0,0,
7                   0,0,1,1,1,0,0,0,0,0,0,0,0,0,0,0,
8                   0,0,1,1,0,0,0,1,1,1,0,0,0,0,0,0,
9                   0,1,1,1,0,1,1,1,1,1,1,1,1,1,1,0,
10                  1,1,1,1,1,1,1,1,0,0,0,0,0,0,1,1,0,
11                  1,1,1,1,1,0,0,0,0,0,0,0,0,0,1,1,1,
12                  1,1,1,0,0,0,0,0,0,0,0,0,0,0,1,1,0,
13                  1,1,1,0,0,0,0,0,0,0,0,0,0,0,1,1,0,
14                  1,1,1,1,0,0,0,0,0,0,1,1,1,1,0,0,
15                  1,1,1,0,1,1,1,0,1,0,1,1,1,1,0,0,0,
16                  0,1,1,1,1,1,1,0,0,0,1,0,0,0,0,0,0]])
```

```
1 y_pred = Modelo.predict(nuevo)
2 # mostrar el resultado de la clasificación del dígito
3 print('El dígito es: ',y_pred[0])
```

```
El dígito es: 6
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but KNeighborsClassi
warnings.warn(
```

## 4. PRÁCTICA

Aplicar el algoritmo k-NN para los siguientes dataset

```
> winequality
> dureza.csv
```

De igual manera experimente con los datos vectorizados de tubérculos y corte de papas nativas del proyecto.

...

✓ 0 s completado a las 8:32

