

# Práctica 3

## Clasificadores k-NN y regresión logística

Mario Emilio Jiménez Vizcaíno  
A01173359@itesm.mx  
Tecnológico de Monterrey  
Ingeniería en Tecnologías Computacionales  
Monterrey, N.L., México

Jesus Abraham Haros Madrid  
A01252642@itesm.mx  
Tecnológico de Monterrey  
Ingeniería en Tecnologías Computacionales  
Monterrey, N.L., México

### ABSTRACT

TODO

### 1 INTRODUCCIÓN

TODO

### 2 CONCEPTOS PREVIOS

- Programación básica en los lenguajes R y Python
- Conocimiento de las librerías scikit-learn, pandas y numpy
- Conocimientos de estadística y de regresión logística

### 3 METODOLOGÍA

#### 3.1 Datasets

Para comparar ambas implementaciones utilizamos dos datasets, expuestos a continuación:

##### 3.1.1 Dataset DEFAULT.

Este dataset está compuesto por 10,000 filas, cada una representa un cliente de un banco que puede o no cumplir con los pagos de su tarjeta de crédito (columna "default"). De cada cliente tenemos la siguiente información:

- Columna "default": Tiene los valores "Yes"/"No", representa si la persona realizó el pago mínimo a su tarjeta de crédito.
- Columna "student": Valores "Yes"/"No", representa si el cliente es un estudiante en ese momento.
- Columna "balance": Número decimal positivo que representa el balance de la tarjeta de crédito del cliente. Promedio de 835.4, números en el rango [0, 2654.3].
- Columna "income": Número decimal positivo que representa los ingresos que tiene el cliente. Promedio de 33517, números en el rango [772, 73554]

De este dataset, nuestro objetivo es predecir la columna "default" a partir de los otros tres parámetros, y como preparación cambiamos los valores de "student" ("Yes"/"No") a valores 1 y 0 respectivamente.

##### 3.1.2 Dataset GENERO.

Este dataset representa las mediciones de peso y altura de 10,000 personas, en conjunto con el género de la persona a la que se realizaron las medidas. Las columnas son:

- "Gender": Valores "Male"/"Female", el género de la persona a la que le corresponde esta fila de mediciones.
- "Height": La altura de la persona en pulgadas, promedio de 66.37, en el rango [54.26 y 79.00].

- "Weight": El peso de la persona en libras, promedio de 161.4, en el rango [64.7, 270.0].

La columna objetivo seleccionada de este dataset fue el género ya que tiene dos clasificaciones.

#### 3.2 Clasificación con k-NN

##### 3.2.1 Dataset DEFAULT.

TODO

El código fuente de este ejemplo se encuentra en el apéndice A.

##### 3.2.2 Dataset GENERO.

TODO

El código fuente puede ser encontrado en el apéndice B.

#### 3.3 Regresión logística

Para la primera parte de la práctica, en la que utilizamos la implementación de *sci-kit learn*, seleccionamos la clase *sklearn.linear\_model.LogisticRegression*[1] para nuestros scripts.

##### 3.3.1 Dataset DEFAULT.

Para este dataset primero leemos el archivo CSV a un *Dataframe* de *pandas*, transformamos las columnas "default" y "student" para que contengan valores booleanos y enteros respectivamente, seleccionamos las columnas que nos servirán como variables independientes (columnas "student", "balance" e "income") y variable dependiente (columna "default"). Después partimos las filas del dataset en una porción del 80% que usaremos para entrenar el modelo, y otra porción del 20% para probarlo.

Instanciamos el modelo de regresión de *sklearn*, lo entrenamos con los datos y después predecimos la variable dependiente con el modelo para así compararlo con los datos reales de prueba, usando una medida de tasa de precisión y la matriz de confusión.

El código de este ejemplo puede se encuentra en el apéndice C.

##### 3.3.2 Dataset GENERO.

Para este dataset realizamos un procedimiento similar: leer el dataset para crear un *Dataframe*, seleccionar las columnas de variables independientes ("Height" y "Weight") y la dependiente ("Gender"), dividir el dataset en 80%/20%, entrenar el modelo, y predecir la variable dependiente para los datos de prueba, para así comparar estos con los datos reales.

El código para esta sección se encuentra en el apéndice D.

## 4 RESULTADOS

### 4.1 Clasificación con k-NN

#### 4.1.1 Dataset DEFAULT.

TODO

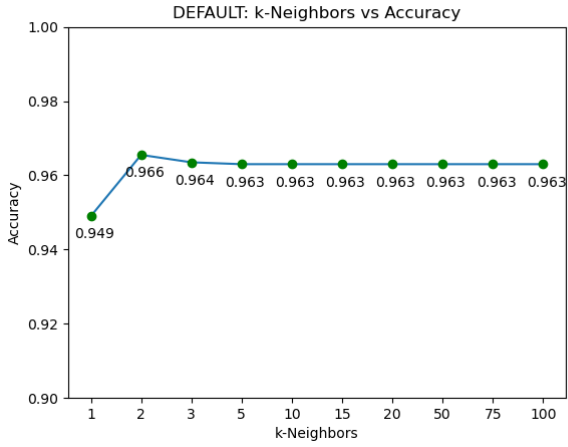


Figure 1: Variación del número de vecinos elegidos contra la precisión en el dataset DEFAULT

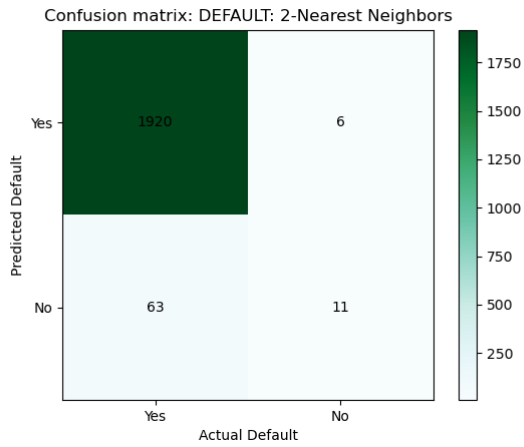


Figure 2: Matriz de confusión del dataset DEFAULT con 2-Nearest Neighbors

#### 4.1.2 Dataset GENERO.

TODO

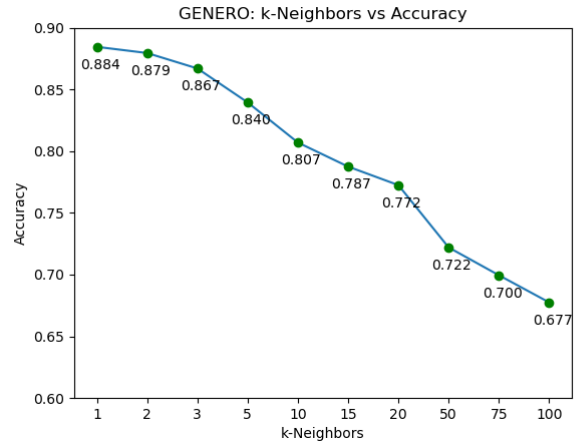


Figure 3: Variación del número de vecinos elegidos contra la precisión en el dataset GENERO

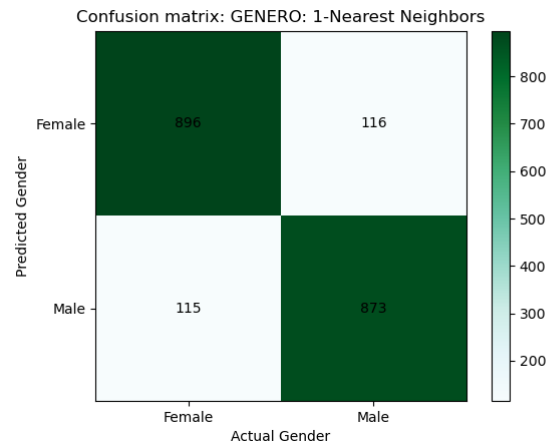


Figure 4: Matriz de confusión del dataset GENERO con 1-Nearest Neighbors

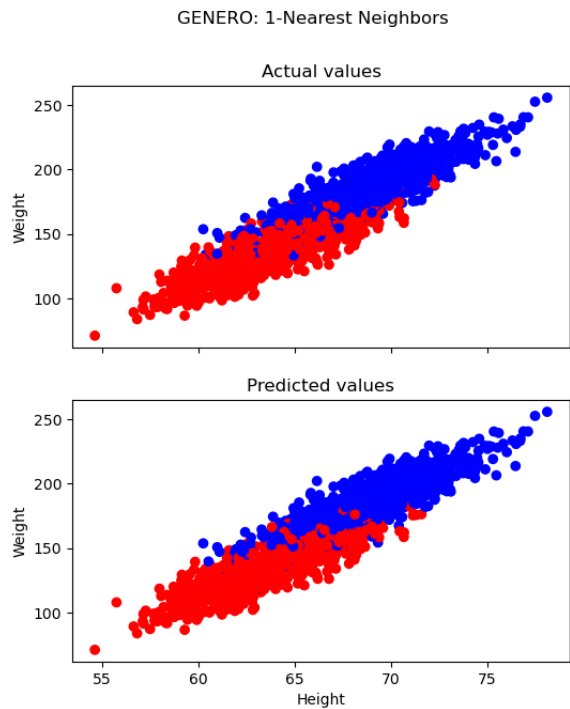


Figure 5: Gráfica de dispersión del dataset GENERO con 1-Nearest Neighbors

4.2 Regresión logística

4.2.1 Dataset DEFAULT.

TODO

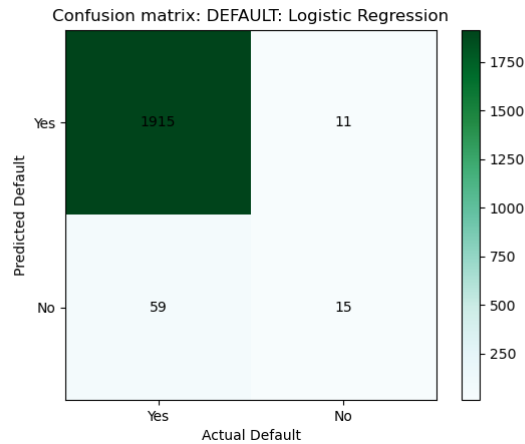


Figure 6: Matriz de confusión del dataset DEFAULT con regresión logística

4.2.2 Dataset GENERO.

TODO



Figure 7: Matriz de confusión del dataset GENERO con regresión logística

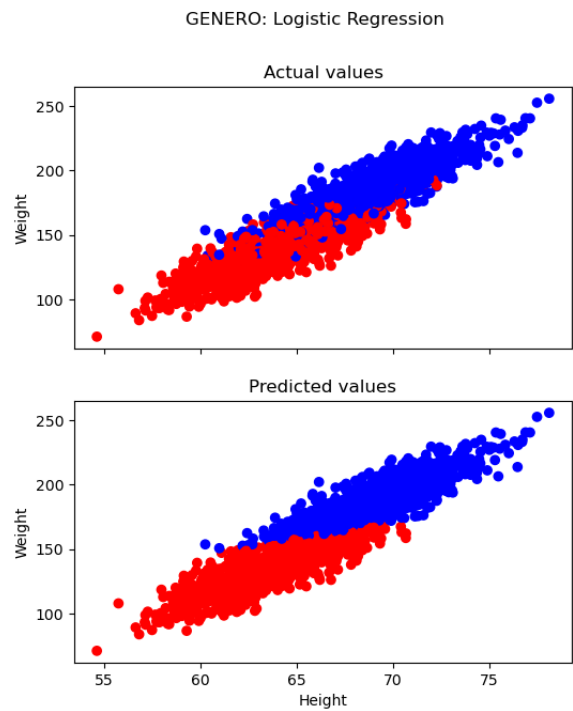


Figure 8: Gráfica de dispersión del dataset GENERO con regresión logística

5 CONCLUSIONES Y REFLEXIONES

TODO

5.1 Refrexión de Abraham

TODO

## 5.2 Reflexión de Mario

TODO

### REFERENCES

- [1] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12 (2011), 2825–2830.

## A CÓDIGO DE REGRESIÓN LOGÍSTICA DEL DATASET DEFAULT

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.neighbors import KNeighborsRegressor
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import accuracy_score, confusion_matrix
6
7 from graphs import graphConfusionMatrix, graphNeighborsAccuracy
8
9 print(" ~ Reading default.txt and generating train and test sets")
10 data = pd.read_csv('default.txt', sep=" ")
11 # Transform 'default' and 'student' columns from Yes/No to integers (1/0)
12 data["default"] = (data["default"] == "Yes").astype(int)
13 data["student"] = (data["student"] == "Yes").astype(int)
14 x = data.iloc[:, 1:4].values.reshape(-1, 3)
15 y = data.iloc[:, 0].values.reshape(-1, 1)
16 xTrain, xTest, yTrain, yTest = train_test_split(
17     x, y, test_size=0.2, random_state=0)
18
19 print(" ~ Creating and testing the k-NN models for different values")
20 print(" → neighbors accuracy confusion matrix")
21 k_neighbors_vals = [1, 2, 3, 5, 10, 15, 20, 50, 75, 100]
22 accuracy_vals = []
23 best_k_val = (-1, 0, []) # k_neighbors, accuracy, conf_matrix
24 for val in k_neighbors_vals:
25     nn = KNeighborsRegressor(n_neighbors=val, algorithm="brute")
26     nn.fit(xTrain, yTrain)
27     yPredicted = nn.predict(xTest).astype(int)
28     accuracy = accuracy_score(yTest, yPredicted)
29     accuracy_vals.append(accuracy)
30     cm = confusion_matrix(yTest, yPredicted)
31     print(" → {:^9} {:^8} {}".format(val, accuracy, cm.tolist()))
32     if best_k_val[1] < accuracy:
33         best_k_val = (val, accuracy, cm)
34
35 graphNeighborsAccuracy(k_neighbors_vals, accuracy_vals,
36     "DEFAULT", "default_neighbors_acc.png")
37 graphConfusionMatrix(best_k_val[2], ["Yes", "No"], "Default",
38     "DEFAULT: {}-Nearest Neighbors".format(best_k_val[0]), "default_neighbors_cm.png")
```

## B CÓDIGO DE CLASIFICACIÓN K-NN DEL DATASET GENERO

```

1 import numpy as np
2 import pandas as pd
3 from sklearn.neighbors import KNeighborsRegressor
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import accuracy_score, confusion_matrix
6
7 from graphs import graphConfusionMatrix, graphNeighborsAccuracy, graphGENERO
8
9 print(" ~ Reading genero.txt and generating train and test sets")
10 data = pd.read_csv('genero.txt')
11 x = data.iloc[:, 1:3].values.reshape(-1, 2)
12 y = data.iloc[:, 0].values.reshape(-1, 1)
13 xTrain, xTest, yTrain, yTest = train_test_split(
14     x, y, test_size=0.2, random_state=0)
15
16
17 yIntToStr = np.vectorize(lambda x: ["Female", "Male"][int(x)])
18
19 print(" ~ Creating and testing the k-NN models for different values")
20 print(" → neighbors accuracy confusion matrix")
21 k_neighbors_vals = [1, 2, 3, 5, 10, 15, 20, 50, 75, 100]
22 accuracy_vals = []
23 best_k_val = (-1, 0, [], []) # k_neighbors, accuracy, conf_matrix, yPredicted
24 for val in k_neighbors_vals:
25     nn = KNeighborsRegressor(n_neighbors=val, algorithm="brute")
26     # Convert String Male/Female to integers 1/0
27     yTrainInt = (yTrain.ravel() == "Male").astype(int)
28     nn.fit(xTrain, yTrainInt)
29     yPredicted = nn.predict(xTest)
30     yPredicted = yIntToStr(yPredicted)
31     accuracy = accuracy_score(yTest, yPredicted)
32     accuracy_vals.append(accuracy)
33     cm = confusion_matrix(yTest, yPredicted)
34     print(" → {:^9} {:^8} {}".format(val, accuracy, cm.tolist()))
35     if best_k_val[1] < accuracy:
36         best_k_val = (val, accuracy, cm, yPredicted)
37
38 graphNeighborsAccuracy(k_neighbors_vals, accuracy_vals,
39     "GENERO", "genero_neighbors_acc.png")
40 graphGENERO(xTest, yTest, best_k_val[3], "{} - Nearest Neighbors".format(best_k_val[0]),
41     "genero_neighbors_pred.png")
42 graphConfusionMatrix(best_k_val[2], ["Female", "Male"], "Gender",
43     "GENERO: {} - Nearest Neighbors".format(best_k_val[0]),
44     "genero_neighbors_cm.png")

```

## C CÓDIGO DE REGRESIÓN LOGÍSTICA DEL DATASET DEFAULT

```
1 import pandas as pd
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.model_selection import train_test_split
4 from sklearn.metrics import accuracy_score, confusion_matrix
5
6 from graphs import graphConfusionMatrix
7
8 print(" ~ Reading default.txt and generating train and test sets")
9 data = pd.read_csv('default.txt', sep=" ")
10 # Transform 'default' column from Yes/No to a boolean
11 data["default"] = (data["default"] == "Yes").astype(bool)
12 # Transform 'student' column from Yes/No to an integer
13 data["student"] = (data["student"] == "Yes").astype(int)
14 x = data.iloc[:, 1:4].values.reshape(-1, 3)
15 y = data.iloc[:, 0].values.reshape(-1, 1)
16 xTrain, xTest, yTrain, yTest = train_test_split(
17     x, y, test_size=0.2, random_state=0)
18
19 print(" ~ Creating the logistic regression model")
20 regressor = LogisticRegression()
21 regressor.fit(xTrain, yTrain.ravel())
22
23 print(" ~ Testing the logistic regression model")
24 yPredicted = regressor.predict(xTest)
25 accuracy = accuracy_score(yTest, yPredicted)
26 print(" → Accuracy:", accuracy)
27 cm = confusion_matrix(yTest, yPredicted)
28 print(" → Confusion matrix:", cm.tolist())
29
30 graphConfusionMatrix(cm, ["Yes", "No"], "Default",
31     "DEFAULT: Logistic Regression", "default_regression_cm.png")
```

## D CÓDIGO DE REGRESIÓN LOGÍSTICA DEL DATASET GENERO

```

1 import pandas as pd
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.model_selection import train_test_split
4 from sklearn.metrics import accuracy_score, confusion_matrix
5
6 from graphs import graphConfusionMatrix, graphGENERO
7
8 print(" ~ Reading genero.txt and generating train and test sets")
9 data = pd.read_csv('genero.txt')
10 x = data.iloc[:, 1:3].values.reshape(-1, 2)
11 y = data.iloc[:, 0].values.reshape(-1, 1)
12 xTrain, xTest, yTrain, yTest = train_test_split(
13     x, y, test_size=0.2, random_state=0)
14
15 print(" ~ Creating the logistic regression model")
16 regressor = LogisticRegression()
17 regressor.fit(xTrain, yTrain.ravel())
18
19 print(" ~ Testing the logistic regression model")
20 yPredicted = regressor.predict(xTest)
21 accuracy = accuracy_score(yTest, yPredicted)
22 print(" → Accuracy:", accuracy)
23 cm = confusion_matrix(yTest, yPredicted)
24 print(" → Confusion matrix:", cm.tolist())
25
26 graphGENERO(xTest, yTest, yPredicted, "Logistic Regression",
27             "genero_regression_pred.png")
28 graphConfusionMatrix(cm, ["Male", "Female"], "Gender",
29                         "GENERO: Logistic Regression", "genero_regression_cm.png")

```



## E CÓDIGO DE GENERACIÓN DE GRÁFICAS

```
1 import math
2 import sys
3 import matplotlib.pyplot as plt
4 import numpy as np
5
6 shouldDisplay = "--display-graphs" in sys.argv
7 shouldSave = "--save-graphs" in sys.argv
8
9 mapColor = np.vectorize(lambda x: "b" if x == "Male" else "r")
10
11
12 def graphGENERO(xys, actualVals, predictedVals, title, filename):
13     fig, (ax1, ax2) = plt.subplots(2, sharex=True)
14     fig.set_size_inches(6, 7)
15     fig.suptitle("GENERO: " + title)
16
17     ax1.set_title("Actual values")
18     ax1.set_ylabel("Weight")
19     actualColors = mapColor(actualVals).flatten()
20     ax1.scatter(xys[:, 0], xys[:, 1], c=actualColors)
21
22     ax2.set_title("Predicted values")
23     ax2.set_xlabel("Height")
24     ax2.set_ylabel("Weight")
25     predictedColors = mapColor(predictedVals).flatten()
26     ax2.scatter(xys[:, 0], xys[:, 1], c=predictedColors)
27
28     if shouldSave:
29         fig.savefig(filename)
30     if shouldDisplay:
31         plt.show()
32
33
34 def graphConfusionMatrix(cm, labels, variableName, title, filename):
35     plt.figure()
36     plt.imshow(cm, interpolation="nearest", cmap="BuGn")
37     plt.title("Confusion matrix: " + title)
38     plt.colorbar()
39     plt.xticks(np.arange(2), labels, size=10)
40     plt.yticks(np.arange(2), labels, size=10)
41     plt.xlabel("Actual " + variableName)
42     plt.ylabel("Predicted " + variableName)
43
44     for x in range(2):
45         for y in range(2):
46             plt.annotate(cm[x][y], xy=(y, x),
47                         horizontalalignment="center",
48                         verticalalignment="center")
49
50     if shouldSave:
51         plt.savefig(filename)
52     if shouldDisplay:
53         plt.show()
54
55
56 def graphNeighborsAccuracy(k_neighbors, accuracy, title, filename):
57     x = list(map(str, k_neighbors))
58
59     plt.figure()
60     plt.title(title + ": k-Neighbors vs Accuracy")
61     plt.plot(x, accuracy)
62     plt.plot(x, accuracy, 'og')
```