

# 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 C.

#### 3.2.2 Dataset GENERO.

TODO

El código fuente puede ser encontrado en el apéndice D.

### 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 A.

#### 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 B.

## 4 RESULTADOS

### 4.1 Clasificación con k-NN

4.1.1 *Dataset DEFAULT.*

TODO

4.1.2 *Dataset GENERO.*

TODO

### 4.2 Regresión logística

4.2.1 *Dataset DEFAULT.*

TODO

4.2.2 *Dataset GENERO.*

TODO

## 5 CONCLUSIONES Y REFLEXIONES

TODO

### 5.1 Refrexi3n de Abraham

TODO

### 5.2 Reflexi3n 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**
- B CÓDIGO DE REGRESIÓN LOGÍSTICA DEL DATASET GENERO**
- C CÓDIGO DE CLASIFICACIÓN K-NN DEL DATASET DEFAULT**
- D CÓDIGO DE CLASIFICACIÓN K-NN DEL DATASET GENERO**
- E CÓDIGO DE GENERACIÓN DE GRÁFICAS**