

*Universidad de La Laguna*

*Escuela de Doctorado y Estudios de Posgrado*

*Departamento de Ingeniería Informática y de Sistemas*

# Extracción de conocimiento en bases de datos

PRÁCTICA: CLASIFICACIÓN

*J. Marcos Moreno Vega*

Copyright © 2023 J. Marcos Moreno Vega

UNIVERSIDAD DE LA LAGUNA

TUFTE-LATEX.GOOGLECODE.COM

Licensed under the Apache License, Version 2.0 (the “License”); you may not use this file except in compliance with the License. You may obtain a copy of the License at <http://www.apache.org/licenses/LICENSE-2.0>. Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an “AS IS” BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

---

**Objetivo:**

Construir y evaluar modelos de clasificación.

**Base de datos:**

homeLoanAproval.csv

**Lenguaje de programación:**

R o Python.

---

### *Concesión de hipotecas*

El mercado hipotecario mueve anualmente una cifra de negocios que alcanza los miles de millones de euros. Las entidades bancarias deben gestionar miles de solicitudes anuales y decidir qué hipotecas conceder y cuáles denegar. Conceder una hipoteca que luego no será devuelta supone un perjuicio económico para la entidad. Del mismo modo, denegar una hipoteca a un solicitante con capacidad para afrontar los pagos es perder una oportunidad de negocio.

Se desea construir un sistema basado en Analítica de datos que ayude en la tarea de concesión de los créditos. Para ello, se dispone de información detallada de los créditos gestionados en el pasado, así como de la decisión que se adoptó con cada solicitud.

### *Base de datos*

La base de datos homeLoanAproval.csv almacena información sobre 614 solicitudes de crédito hipotecario. El fichero consta de 12 variables, 11 de las cuáles son predictivas. En la tabla 1.1 se enumeran y describen todas las variables.

### *Tareas*

- a) Preparar los datos. Desarrollar las tareas de preparación de los datos necesarias para disponer de una base de datos apropiada para la tarea de clasificación. Entre las labores a realizar se encuentran las de visualización y comprensión de las variables, cambio de tipos, transformación e identificación de outliers y datos faltantes. También debe quedar claro si nos enfrentamos a una clasificación con clases balanceadas o no.
- b) Analizar los datos. Construir un clasificador  $k$ -NN, un árbol de clasificación y un clasificador naive Bayes. Se propone considerar varios escenarios. En el primero, independientemente de que existan o no datos faltantes o outliers, se aplicará el correspondiente algoritmo de inducción sin tener en cuenta estos hechos. En el segundo escenario, se aplicarán técnicas para tratar con los outliers y los datos faltantes y se construirán los anteriores clasificadores. Debe evaluarse también el uso de técnicas para tratar con clases

Variable	Descripción	Valores
LoanID	identificador de la solicitud	Alfanumérico
Gender	género del solicitante	Male/Female
Married	variable que indica si el solicitante es una persona casada	Yes/No
Dependents	número de personas a cargo del solicitante	Numérico
Education	nivel educativo del solicitante	Graduate/Not Graduate
SelfEmployed	variable que indica si el solicitante es autónomo o no	Yes/No
ApplicationIncome	ingresos del solicitante	Numérico
CoapplicationIncome	ingresos del segundo solicitante del crédito	Numérico
LoanAmount	importe del crédito hipotecario	Numérico
LoanAmountTerm	periodo del crédito (en meses)	Numérico
PropertyArea	lugar en que se ubica la propiedad hipotecada	Rural/Urban/Semiurban
LoanStatus	variable que indica si el crédito fue aprobado	Yes/No

desbalanceadas y el efecto que tienen estas técnicas en la calidad de los clasificadores obtenidos.

Tabla 1.1: Base de datos home-LoanApproval.csv

### *Qué debe presentar el alumnado*

- Cuaderno Jupyter, debidamente comentado, usado para preparar y analizar los datos.
- Una memoria en formato pdf en la que se describan brevemente el análisis realizado y el modelo de clasificación obtenido para ayudar a los responsables de las entidades bancarias.
- La memoria debe incluir también gráficas y tablas que muestren la calidad de los clasificadores obtenidos y que permitan identificar cuál es el más apropiado.