



**Universidad de Buenos Aires - Facultad de Ingeniería  
Organización de Datos 75.06/95.58**

## **Informe trabajo práctico N°2**

**Lucas Risaro - Javier Nuñez Leyes**

94335 - 94455

# Índice

|   |          |
|---|----------|
| <b>Introducción</b>                         | <b>3</b> |
| <b>Tabla de métodos de preprocesamiento</b> | <b>3</b> |
| <b>Tabla de modelos</b>                     | <b>4</b> |
| <b>Matrices de confusión</b>                | <b>5</b> |
| <b>Conclusión</b>                           | <b>6</b> |
| Mejor modelo                                | 6        |
| Análisis de falsos positivos                | 6        |

## Informe TP2

### Introducción

A continuación se enuncian los detalles de los modelos utilizados para realizar las predicciones solicitadas sobre el set de datos entregado.

Se listan los nombres, métodos utilizados y resultados obtenidos de cada uno junto con las métricas que se tuvieron en cuenta para decidir cuál de los modelos fue el mejor.

La información se mostrará en dos tablas, una que listara los métodos de preprocesamiento y otra con la información pertinente de cada modelo.

### Tabla de métodos de preprocesamiento

| Nombre del preprocesamiento  | Explicación   | Función de python                   |
|--|---|-------------------------------------|
| Preprocesamiento General de los modelos XGBOOST y Random Forest                  | Función donde se realiza todo el feature engineering sobre el set de datos. Se unifican y eliminan columnas y se aplica One Hot Encoding a las definitivas.             | feature_engineering_original_set()  |
| Preprocesamiento general más rango de edad para XGB y RF                         | Preprocesamiento general para XGB y RF y además crea columnas para encasillar a cada usuario en un rango de edad  | feature_engineering_add_age_range() |
| One Hot Encoding   | Aplica One Hot Encoding a las columnas y set que se le especifique  | apply_one_hot_encoding()            |
| Codificación de combinación de valores de "trabajo" y "rol_familiar_registrado". | Función que devuelve un 1 si esa row posee la combinación rol_familiar_registrado = "casado" y (trabajo = "profesional_especializado" o trabajo = "directivo_gerente" ) | set_value_row_casado_trabajo()      |
| Unifica valores Casado y Casada  | Unifica los valores casado y casada de la columna   | unificar_valores_casado_casada()    |

|   |   |                                      |
|---|---|--------------------------------------|
|   | rol_familiar_registrado en casado   |                                      |
| Obtener columnas por índice                             | Devuelve las columnas del set de datos indicadas por sus índices  | get_columns_by_index                 |
| Preprocesamiento de todas las columnas                  | Convierte las columnas a valores numéricos, unifica valores categóricos y aplica one hot encoding.                | preprocessing_con_todos_los_features |
| Preprocesamiento determinado en la primera parte del TP | Unifica valores, aplica one hot encoding y elimina las columnas innecesarias (según análisis de la primera parte) | feature_engineering_TP_primera_parte |

## Tabla de modelos

| Nombre de modelo | Preprocesamiento   | AUC-ROC | Accuracy | Precisión | Recall | F1-Score |
|------------------|--|---------|----------|-----------|--------|----------|
| Random Forest    | Preprocesamiento General de los modelos XGBOOST y Random Forest                                    | 0.890   | 0.838    | 0.692     | 0.916  | 0.638    |
| 1-XGBoost        | Preprocesamiento General de los modelos XGBOOST y Random Forest                                    | 0.892   | 0.840    | 0.692     | 0.915  | 0.645    |
| 2-XGBoost        | Preprocesamiento general más rango de edad para XGB y RF   | 0.898   | 0.846    | 0.711     | 0.921  | 0.657    |
| 3-XGBoost        | Preprocesamiento General de los modelos XGBOOST y Random Forest<br><br>Obtener columnas por índice | 0.782   | 0.794    | 0.694     | 0.963  | 0.380    |

|               |   |       |       |       |       |       |
|---------------|---|-------|-------|-------|-------|-------|
| 1-KNN         | Preprocesamiento de todas las columnas                  | 0.892 | 0.842 | 0.704 | 0.920 | 0.640 |
| 2-KNN         | Preprocesamiento determinado en la primera parte del TP | 0.847 | 0.810 | 0.652 | 0.923 | 0.533 |
| 1-Naive Bayes | Preprocesamiento de todas las columnas                  | 0.859 | 0.785 | 0.540 | 0.806 | 0.616 |
| 2-Naive Bayes | Preprocesamiento determinado en la primera parte del TP | 0.836 | 0.806 | 0.683 | 0.946 | 0.475 |
| 1-SVM         | Preprocesamiento de todas las columnas                  | 0.889 | 0.831 | 0.756 | 0.954 | 0.559 |
| 2-SVM         | Preprocesamiento determinado en la primera parte del TP | 0.719 | 0.821 | 0.652 | 0.907 | 0.596 |

## Matrices de confusión

| Modelo        | Matriz                   |
|---------------|--------------------------|
| Random Forest | [4544 403]<br>[ 647 919] |
| 1-XGBoost     | [4528 419]<br>[ 621 945] |
| 2-XGBoost     | [4559 388]<br>[ 610 956] |
| 3-XGBoost     | [4766 181]<br>[1155 411] |
| 1-KNN         | [4554 393]<br>[ 631 935] |
| 2-KNN         | [4571 376]<br>[ 860 706] |

|               |                           |
|---------------|---------------------------|
| 1-Naive Bayes | [3988 959]<br>[ 440 1126] |
| 2-Naive Bayes | [4682 265]<br>[ 995 571]  |
| 1-SVM         | [4723 224]<br>[ 871 695]  |
| 2-SVM         | [4487 460]<br>[ 704 862]  |

## Conclusión

### Mejor modelo

En base al área bajo la curva AUC-ROC (0.898), el mejor modelo es XGBoost version 2. Este número nos indica que tan bueno es nuestro modelo para distinguir entre las clases. Cuanto más alto es el valor de AUC, mejor el modelo está prediciendo los 0 como 0 y los 1 como 1.

Con respecto al baseline del TP 1 este está fuertemente ligado a los datos de “entrenamiento”, es decir que el “modelo” con “IFs” tiene problemas de overfitting, no es capaz de predecir de manera correcta nuevos sets de datos ya que memoriza los datos que uso para su entrenamiento.

### Análisis de falsos positivos

Para recomendar un modelo basado en el manejo de los falsos positivos analizaremos el recall y la precisión de cada uno.

Si no nos importa tener entre nuestras predicciones un número alto de falsos positivos nos fijaremos en el recall del modelo en cambio si lo que queremos es un modelo que posea pocos falsos positivos entre sus predicciones nos concentraremos en la precisión del modelo.

Al tener una alta precision nos aseguramos que la gran mayoría de resultados que el modelo dice que son positivos efectivamente lo son, quiere decir que el modelo devuelve pocos falsos positivos. Pero esto hace que en el esfuerzo de asegurar esos resultados como positivos se le escapen otros y lo predigan como negativos, es decir tendría más falsos negativos. Lo que podría ocasionar que su efectividad en la predicción general disminuya.

Al tener un alto recall el modelo devuelve mucho más resultados positivos pero esto hace que tenga muchos falsos positivos entre ellos.

Entonces viendo los resultados detallados en las tablas anteriores si lo que el usuario necesita es un modelo que posea la menor cantidad de falsos positivos posibles, corriendo

el riesgo de tener muchos falsos negativos, entre sus predicciones le recomendamos usar SVM, que de los modelos probados es el que mayor precisión tiene (0.756).

En cambio si lo que queremos es obtener la mayor cantidad de predicciones positivas sin preocuparnos por que un gran número de ellas en realidad no seas positivas recomendamos el modelo con mayor Recall, en nuestro caso XGBoost version 3 que posee un recall de 0.963.