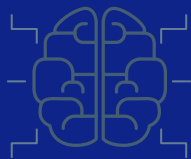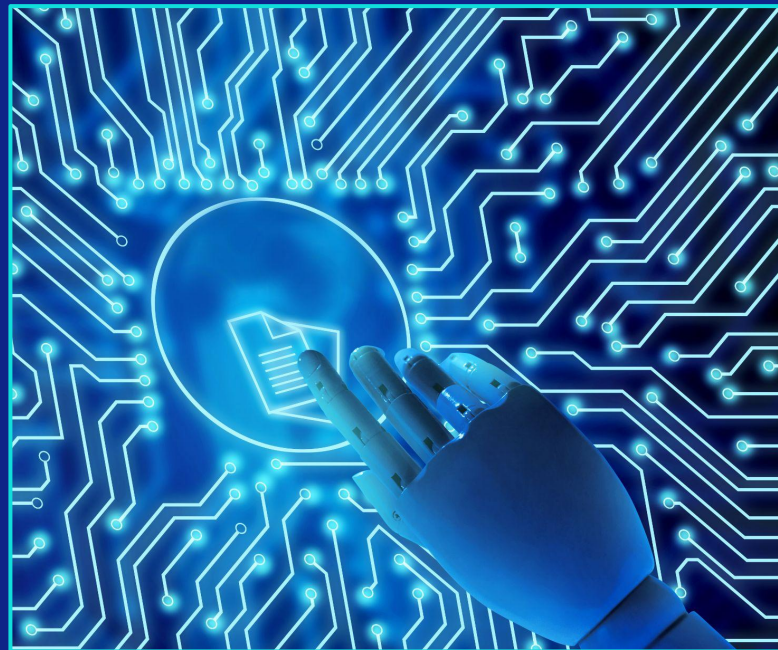# Squad Rapidminer

Francisco Ávila

Rubén Tenreiro

Marco Hernani

# INDICE

**UCI**
**Machine Learning Repository**
Center for Machine Learning and Intelligent Systems
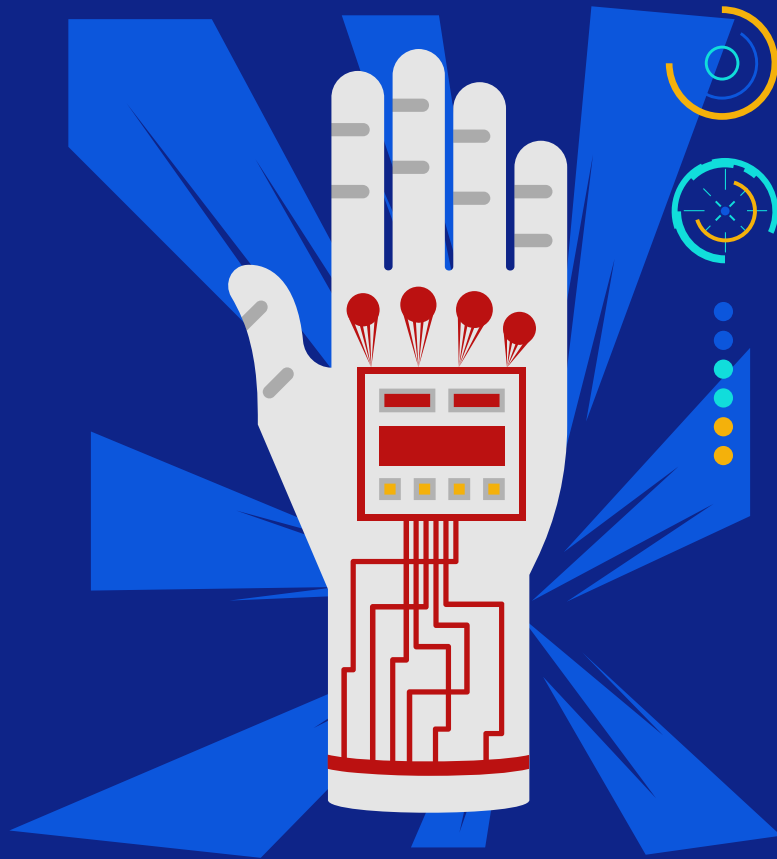
**01** Introducción

**02** EDA

**03** Preprocesado

**04** Modelos

https://archive.ics.uci.edu/ml/datasets/Census+Income

# 01

# Introducción

# Variables cuantitativas

| age | Edad |
|---|---|
| fnlwgt | Número que el censo cree que la entrada representa |
| education_num | Un número relacionado con el nivel de educación (sin un orden específico) |
| capital_gain | Capital ganado con inversiones |
| capital_loss | Capital perdido en inversiones |
| hours_per_week | Horas trabajadas por semanas |

# Variables cualitativas

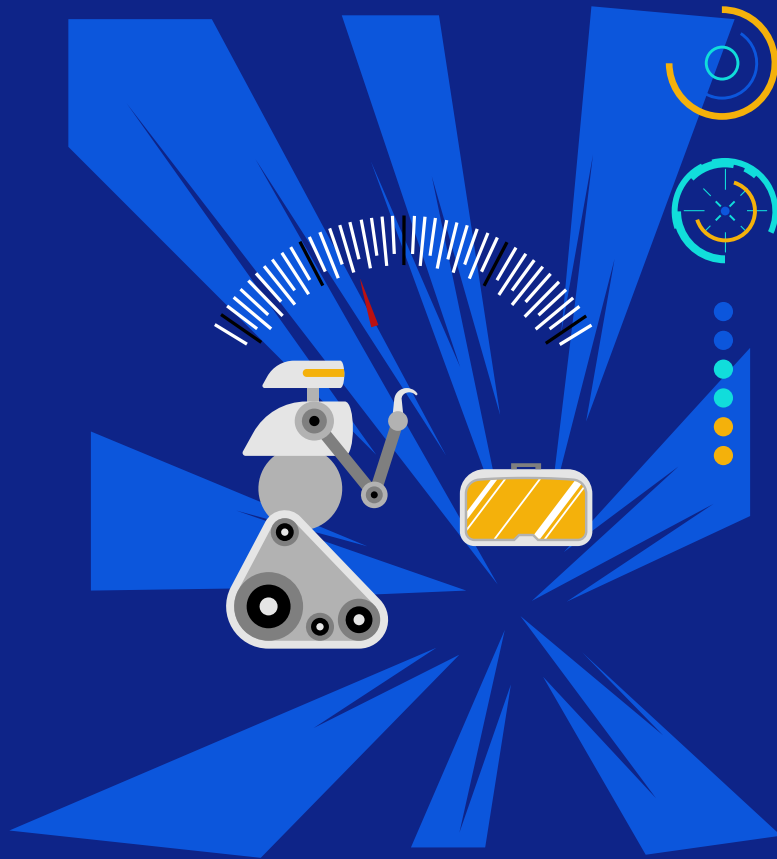| | |
|---|---|
| **work_class** | Clase en la que se categoriza el trabajo de la persona |
| **education** | Nivel más alto de educación de la persona |
| **marital_status** | Estado civil de la persona |
| **occupation** | Trabajo de la persona |
| **relationship** | Tipo de relación que tiene esa persona |
| **race** | Raza de la persona |
| **native_country** | País de nacimiento de la persona |
| **sex** | Sexo de la persona |

# Variable objetivo

| target | Variable objetivo que contiene un valor binario en función de si una persona gana o no más de 50K |
|---|---|

# 02

## EDA

# Dataframe

```
1  df = pd.read_csv("adult.data", header = None, names = columnas)
2  df.head()
```

| | age | work_class | fnlwgt | education | education_num | marital_status | occupation | relationship | race | sex | capital_gain | capital_loss | hours_per_week | native_country | target |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 39 | State-gov | 77516 | Bachelors | 13 | Never-married | Adm-clerical | Not-in-family | White | Male | 2174 | 0 | 40 | United-States | <=50K |
| 1 | 50 | Self-emp-not-inc | 83311 | Bachelors | 13 | Married-civ-spouse | Exec-managerial | Husband | White | Male | 0 | 0 | 13 | United-States | <=50K |
| 2 | 38 | Private | 215646 | HS-grad | 9 | Divorced | Handlers-cleaners | Not-in-family | White | Male | 0 | 0 | 40 | United-States | <=50K |
| 3 | 53 | Private | 234721 | 11th | 7 | Married-civ-spouse | Handlers-cleaners | Husband | Black | Male | 0 | 0 | 40 | United-States | <=50K |
| 4 | 28 | Private | 338409 | Bachelors | 13 | Married-civ-spouse | Prof-specialty | Wife | Black | Female | 0 | 0 | 40 | Cuba | <=50K |

# Resumen variables numéricas
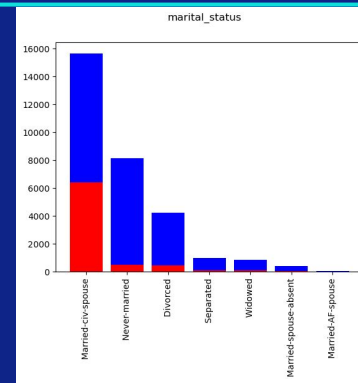
```
1  df.describe().T
```

|                | count    | mean          | std           | min     | 25%      | 50%      | 75%      | max       |
|----------------|----------|---------------|---------------|---------|----------|----------|----------|-----------|
| age            | 32561.0  | 38.581647     | 13.640433     | 17.0    | 28.0     | 37.0     | 48.0     | 90.0      |
| fnlwgt         | 32561.0  | 189778.366512 | 105549.977697 | 12285.0 | 117827.0 | 178356.0 | 237051.0 | 1484705.0 |
| education_num  | 32561.0  | 10.080679     | 2.572720      | 1.0     | 9.0      | 10.0     | 12.0     | 16.0      |
| capital_gain   | 32561.0  | 1077.648844   | 7385.292085   | 0.0     | 0.0      | 0.0      | 0.0      | 99999.0   |
| capital_loss   | 32561.0  | 87.303830     | 402.960219    | 0.0     | 0.0      | 0.0      | 0.0      | 4356.0    |
| hours_per_week | 32561.0  | 40.437456     | 12.347429     | 1.0     | 40.0     | 40.0     | 45.0     | 99.0      |

# Scatter Matrix

EDA 02

# 03

# Preprocesado

# Agrupamos



## work_class



"Private", "Other"

## race



"White", "Other"

## native_country



"United-States", "Other"
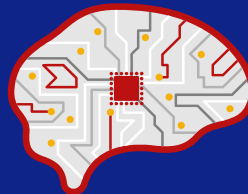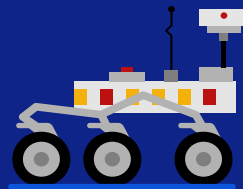
# OneHotEncodeamos

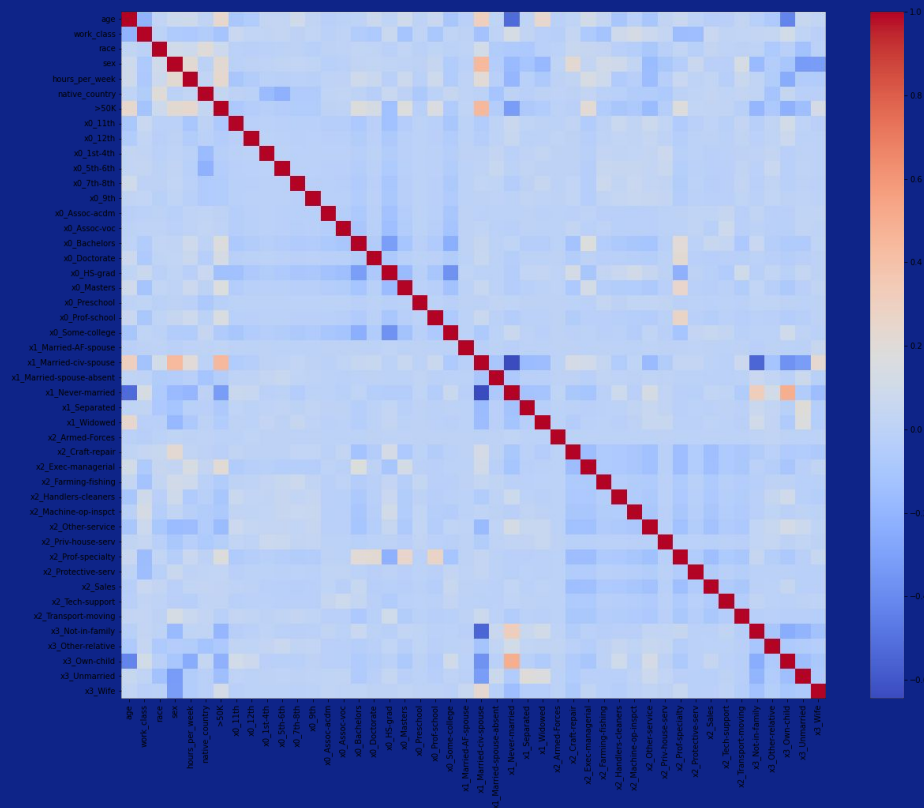**education**   **marital_status**   **occupation**   **relationship**

```
1  jj = OneHotEncoder(drop= 'first')
```

```
1  encoded_df = pd.DataFrame(jj.fit_transform(df_semiclean[['education',
2                          'marital_status', 'occupation', 'relationship']]).toarray())
```

Matriz de correlación

# Train y Test split

```python
1  y = df_final['>50K']
2  X = df_final.drop('>50K', axis= 1)
3  X_train_sn,  X_test_sn, y_train_sn, y_test_sn = train_test_split(X, y, train_size=0.8, random_state=1729)
```

# StandardScaler()

```python
1  X_train_std,  X_test_std, y_train_std, y_test_std = train_test_split(X_std, y_std, train_size=0.8, random_state=1729)
2
3  scaler = preprocessing.StandardScaler().fit(X_train_std[columns_name])
4  X_train_std[columns_name] = scaler.transform(X_train_std[columns_name])
5
6  scaler = preprocessing.StandardScaler().fit(X_test_std[columns_name])
7  X_test_std[columns_name] = scaler.transform(X_test_std[columns_name])
```

# Dos datasets



Sin estandarizar

Estandarizado

# Regresión Logística

Regresión logística SN train



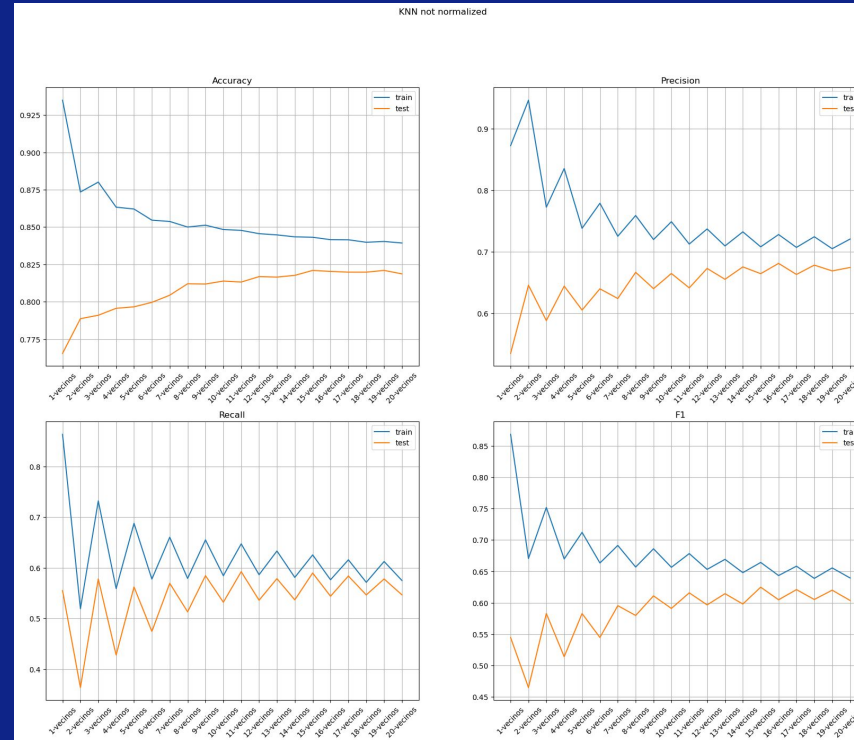Regresión logística SN test
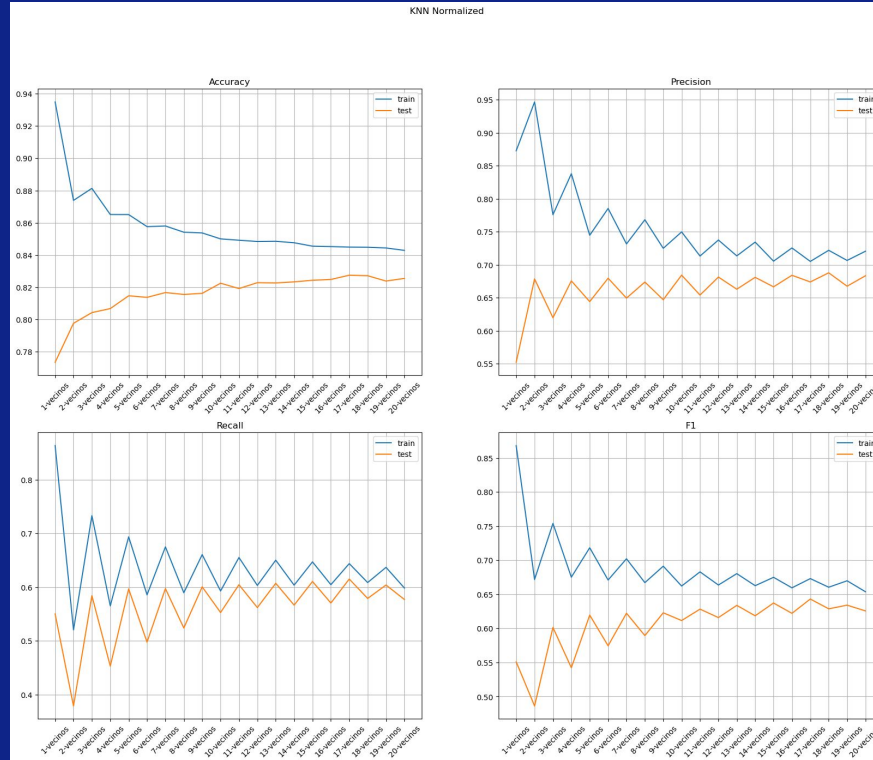
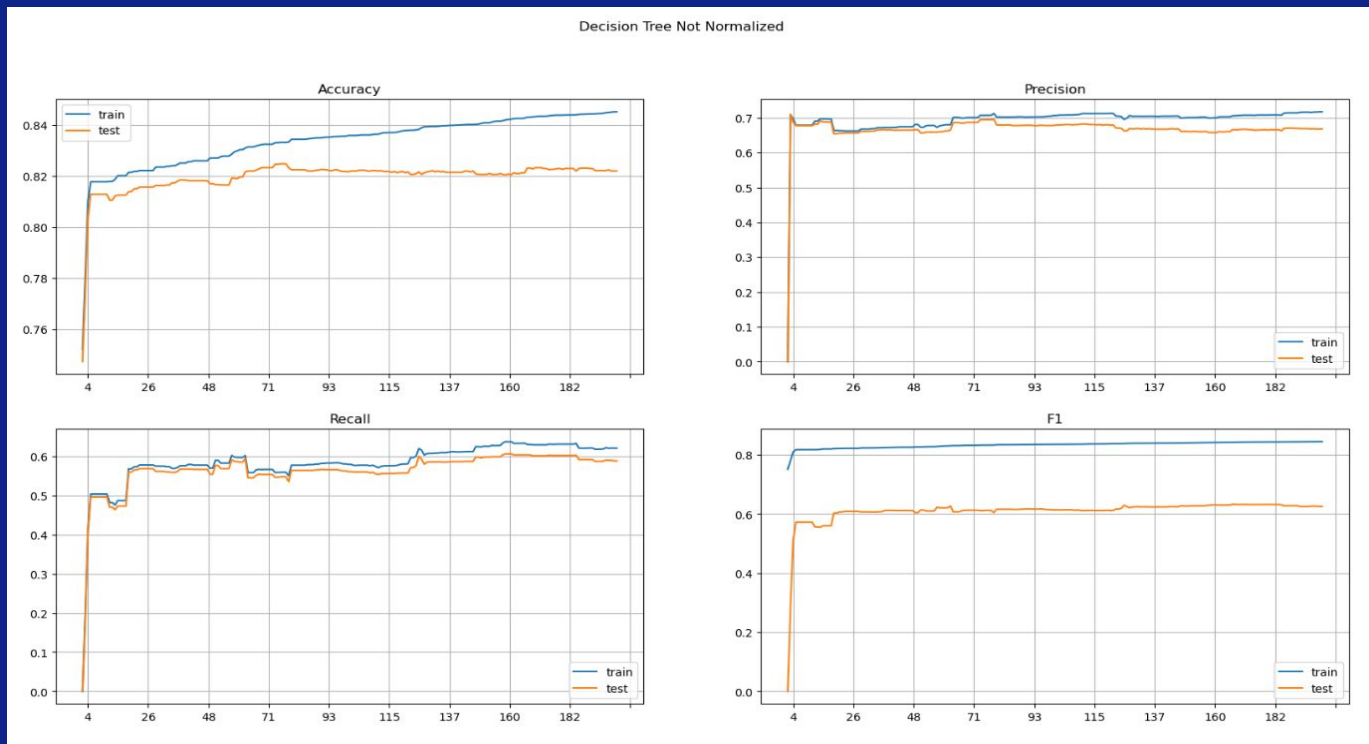# Regresión Logística

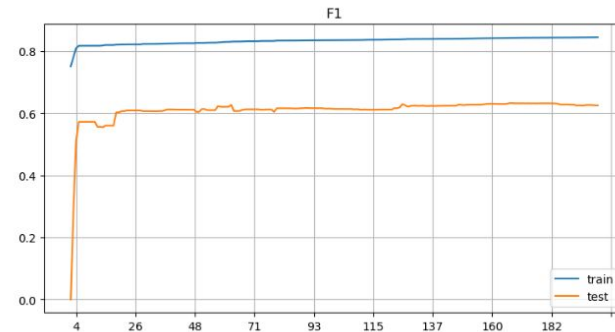Regresión logística STD train



Regresión logística STD test

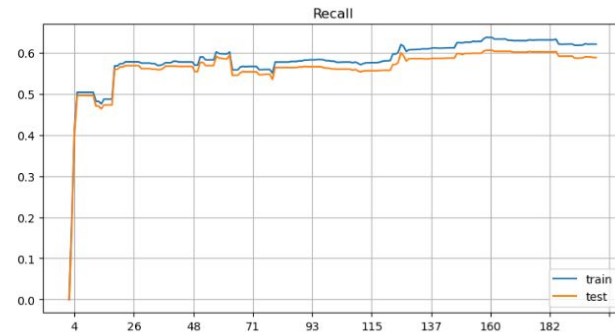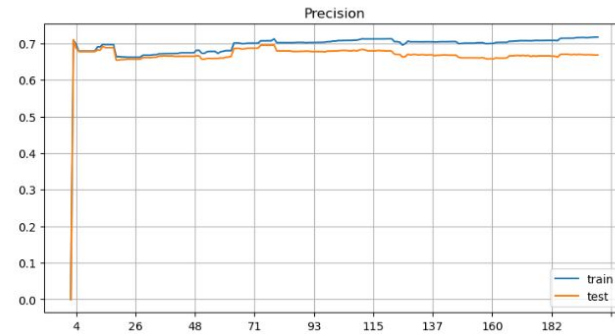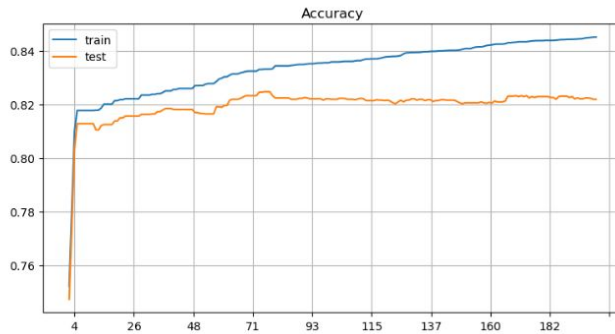# KNN not normalized

# KNN Normalized

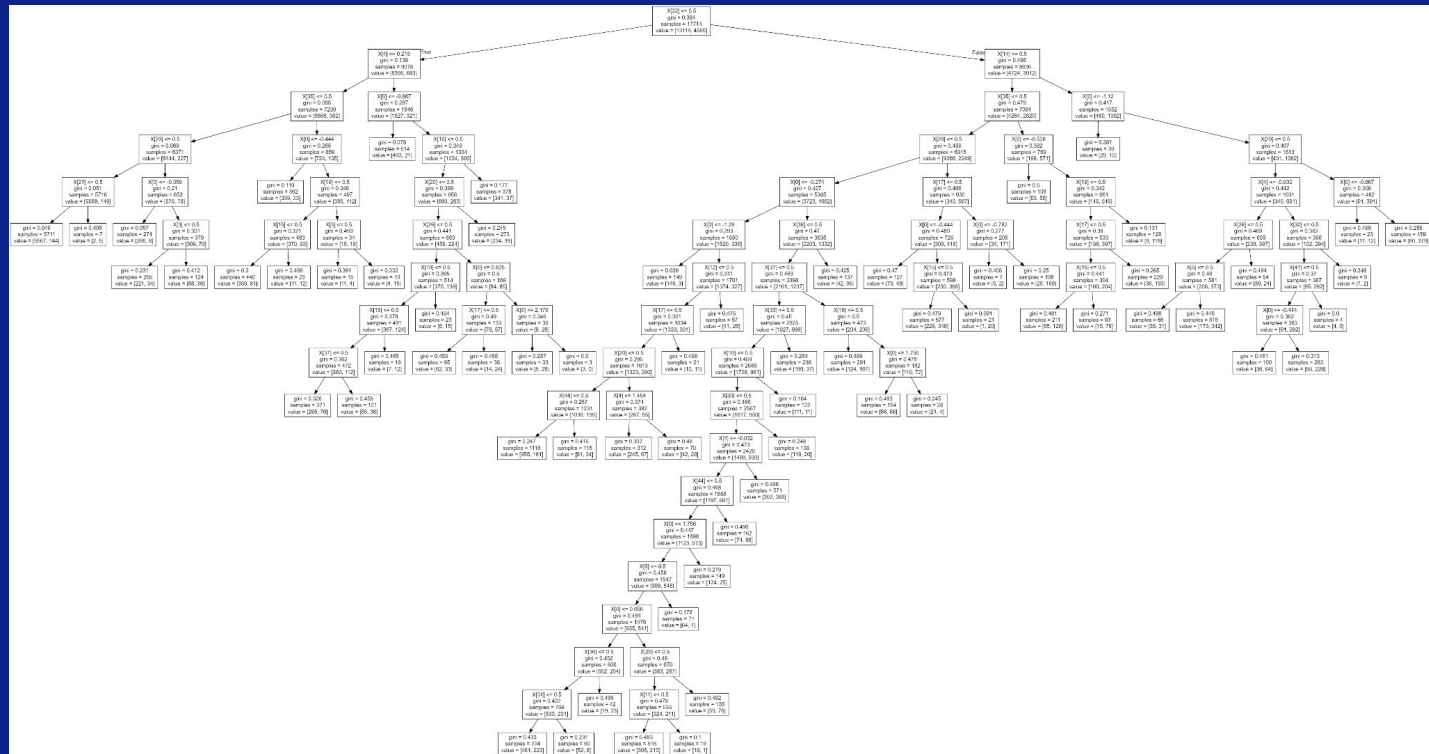# Decision Tree not normalized
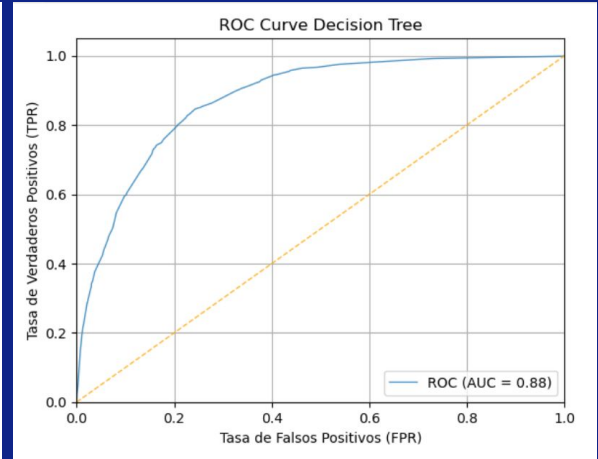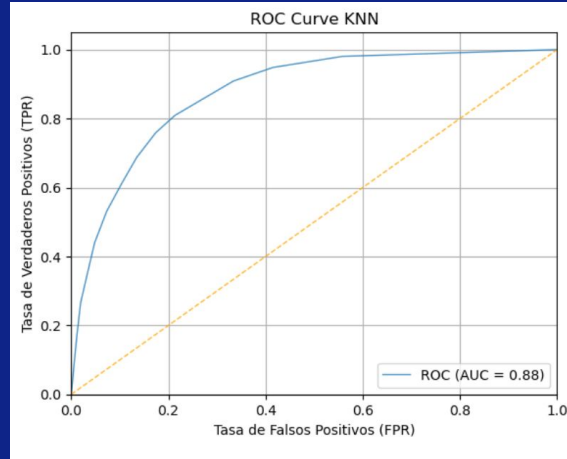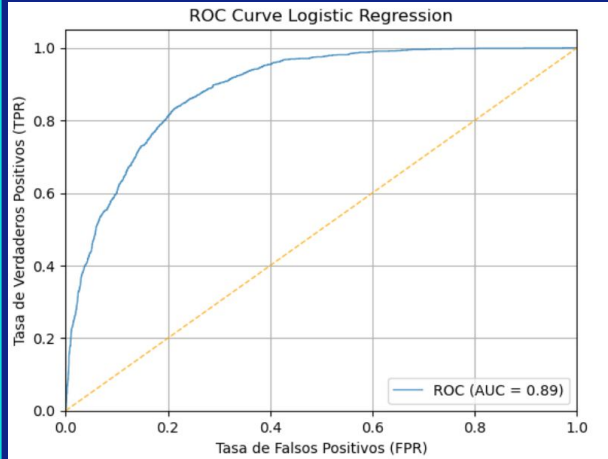


Decision Tree Not Normalized
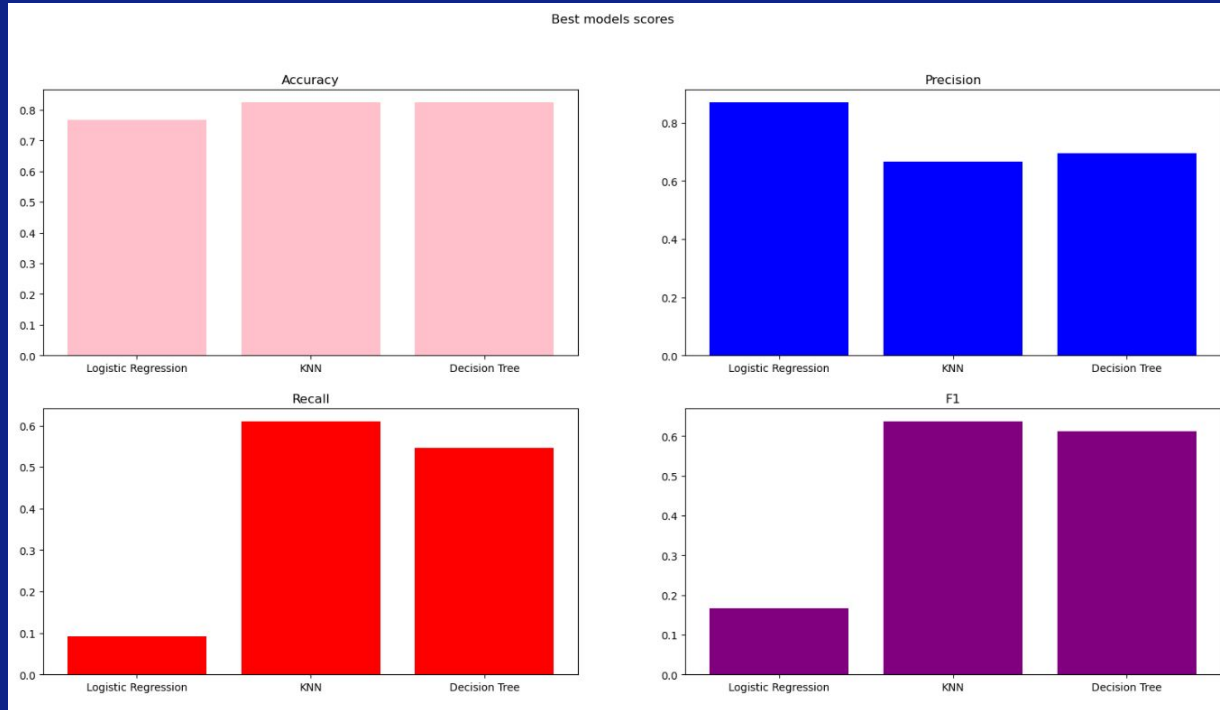
# Decision Tree normalized



Decision Tree Normalized

Modelos 04

# Decision Tree

# Curvas ROC

# Comparación de Resultados



Modelos **04**

# Preguntas