

# PREDICCIÓN DE LA DIABETES EN LA POBLACIÓN FEMENINA DE LOS INDIOS PIMA

Los pima (o akimel o'odham, también deletreado akimel o'otham, "gente del río", antes conocido como pima) son un grupo de nativos americanos que viven en un área que consiste en lo que ahora es el centro y el sur de Arizona.

La población mayoritaria de las dos bandas supervivientes de los akimel o'odham se asienta en dos reservas: los akimel o'odham en la Comunidad India del Río Gila (GILIC) y los on'k akimel o'odham en la Comunidad India Pima-Maricopa del Río Salado (SRPMIC).

Referente a la salud los pimas consideran la intervención divina como causa de la enfermedad o curación de la enfermedad, por lo que continuamente hacen visitas y mandas a San Francisco (héroas de los pimas) para pedir salud.

También recurren a curanderos y mandas para curarse con hierbas medicinales (hierba inmortal y pezuña de vaca para el empacho, ajo y manzanilla para los dolores de parto, ocotillo, ajo, canela y ruda para los sustos, cola de caballo para el "mal de orín", torote prieto para picaduras de alacrán, etc)

Cuando la enfermedad persiste o es grave, acuden a pequeños centros médicos rurales que existen en la región (en Maycoba y el Kipor, bajo la responsabilidad de la Secretaría de Salud del Estado). En los últimos años viajan a los grandes centros de población como Ciudad Obregón y Hermosillo.

Los Pima presentan problemas de enfermedades bronco-respiratorias, diarreas, parásitos intestinales y anemia, principalmente en la población infantil. Existe baja mortalidad infantil y es raro que el promedio de vida supere los 60 años, predominando en un alto porcentaje la población menor de los 19 años.

Este conjunto de datos procede del Instituto Nacional de Diabetes y Enfermedades Digestivas y Renales. El objetivo del conjunto de datos es predecir de forma diagnóstica si un paciente tiene o no diabetes, basándose en determinadas mediciones diagnósticas incluidas en el conjunto de datos.

La selección de estas instancias de una base de datos más amplia está sujeta a varias restricciones. En particular, todos los pacientes son mujeres de al menos 21 años de edad de origen indio Pima.

## Diabetes

La diabetes es una enfermedad que se produce cuando la glucosa en sangre, también llamada azúcar en sangre, es demasiado alta. La glucosa en sangre es el principal fuente de energía y procede de los alimentos que se ingieren. La insulina, una hormona producida por el páncreas, ayuda a que la glucosa de los alimentos llegue a las células para ser utilizada como energía. A veces, el cuerpo no produce suficiente insulina, o no la utiliza bien. La glucosa se queda en la sangre y no llega a las células.

Con el tiempo, tener demasiada glucosa en la sangre puede causar problemas de salud. Aunque la diabetes no tiene cura, puedes tomar medidas para controlarla y mantenerte sano.

¿Cuáles son los distintos tipos de diabetes? Los tipos más comunes de diabetes son el tipo 1, el tipo 2 y la diabetes gestacional.

1. Diabetes de tipo 1: Si tiene diabetes de tipo 1, su cuerpo no produce insulina. Su sistema inmunitario ataca y destruye las células del páncreas que producen la insulina. La diabetes de tipo 1 suele diagnosticarse en niños y adultos jóvenes, aunque puede aparecer a cualquier edad. Las personas con diabetes de tipo 1 necesitan administrarse insulina todos los días para mantenerse con vida.

2. Diabetes de tipo 2: Si tiene diabetes de tipo 2, su cuerpo no fabrica ni utiliza de forma correcta la insulina. La diabetes de tipo 2 puede aparecer a cualquier edad, incluso durante la infancia. Sin embargo, este tipo de diabetes se da con mayor frecuencia en personas de mediana edad y mayores. El tipo 2 es el más común de los tipos de diabetes.

3. Diabetes gestacional La diabetes gestacional se desarrolla en algunas mujeres cuando están embarazadas. La mayoría de las veces, este tipo de diabetes desaparece después del nacimiento del bebé. Sin embargo, si has tenido diabetes gestacional, tienes más posibilidades de desarrollar diabetes de tipo 2 más adelante. A veces, la diabetes diagnosticada durante el embarazo es en realidad una diabetes de tipo 2.

Otros tipos de diabetes Los tipos menos comunes son la diabetes monogénica, que es una forma hereditaria de diabetes, y la diabetes relacionada con la fibrosis quística.

## EDA

### Carga de las librerías

```
In [ ]: import pandas as pd
import numpy as np

import seaborn as sns
import matplotlib.pyplot as plt

import matplotlib inline

import scipy.stats as st
from scipy import interp
from scipy.stats import randint
from scipy.stats import uniform
import warnings
warnings.filterwarnings('ignore')
```

### Carga de los datos

```
In [ ]: from google.colab import drive
drive.mount('/content/drive')

In [ ]: data = pd.read_csv('/content/drive/MyDrive/O1_01col/O1_Machine_Learning/00_PROYECTO04_00_Proyectos_Machine_Lear
```

```
In [ ]: data.info()

Out [ ]: <class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
#   Column              Non-Null Count  Dtype
---  --
0   Pregnancies          768 non-null    int64
1   Glucose              768 non-null    int64
2   BloodPressure        768 non-null    int64
3   SkinThickness        768 non-null    int64
4   Insulin              768 non-null    float64
5   BMI                  768 non-null    float64
6   DiabetesPedigreeFunction  768 non-null    float64
7   Age                  768 non-null    int64
8   Outcome              768 non-null    int64
dtypes: float64(2), int64(7)
memory usage: 54.1 KB
```

```
In [ ]: data.head()

Out [ ]:   Pregnancies  Glucose  BloodPressure  SkinThickness  Insulin   BMI  DiabetesPedigreeFunction  Age  Outcome
0             6      148             72             35     0.336              0.627    50      1
1             1       85              66             29     0.266              0.351    31      0
2             8      183              64             0     0.233              0.672    32      1
3             1       89              66             23     94.281             0.167    21      0
4             0      137              40             35     168.431             2.288    33      1
```

```
In [ ]: data.shape

Out [ ]: (768, 9)
```

```
In [ ]: data.describe()

Out [ ]:   Pregnancies  Glucose  BloodPressure  SkinThickness  Insulin   BMI  DiabetesPedigreeFunction  Age  Outcome
count  768.000000  768.000000  768.000000  768.000000  768.000000  768.000000  768.000000  768.000000  768.000000
mean      3.845052   120.845311   69.105469      20.536458   79.799479   31.992578              0.471876   33.240885   0.349598
std       3.369578   31.972618   19.355807   15.952218   111.244002      7.884160              0.331329   11.760232   0.476951
min        0.000000   0.000000      0.000000      0.000000      0.000000      0.000000              0.078000   21.000000   0.000000
25%        1.000000   99.000000      62.000000      0.000000      0.000000      27.300000              0.247500   24.000000   0.000000
50%        3.000000  117.000000     72.000000     23.000000     30.500000     32.000000              0.323500   29.000000   0.000000
75%        6.000000  140.250000     80.000000     32.000000     127.250000     36.600000              0.626250   41.000000   1.000000
max       17.000000  199.000000    122.000000     99.000000     846.000000     67.100000              2.420000   81.000000   1.000000
```

El conjunto de datos proporciona la información de las pacientes. Incluye 768 registros y 8 atributos mas la variable objetivo. Cada atributo es un factor de riesgo potencial.

El conjunto de datos consta de varias variables médicas predictoras (independientes) y una variable objetivo (dependiente), el resultado. Las variables independientes incluyen:

- 1- Embarazos: número de embarazos
- 2- Glucosa: concentración de glucosa en plasma durante 2 horas en una prueba de tolerancia a la glucosa oral
- 3- Presión arterial: presión arterial diastólica (mm Hg)
- 4- SkinThickness: Espesor del pliegue cutáneo del tríceps (mm)
- 5- Insulina: insulina sérica de 2 horas (mu U / ml)
- 6- IMC: índice de masa corporal (peso en kg / (altura en m) 2)
- 7- DiabetesPedigreeFunction: función del pedigrí de la diabetes (una función que puntúa la probabilidad de diabetes según los antecedentes familiares)
- 8- Edad: Edad (años)
- 9- Resultado: variable de clase (0 si no es diabético, 1 si es diabético)

### Distribución de la variable objetivo

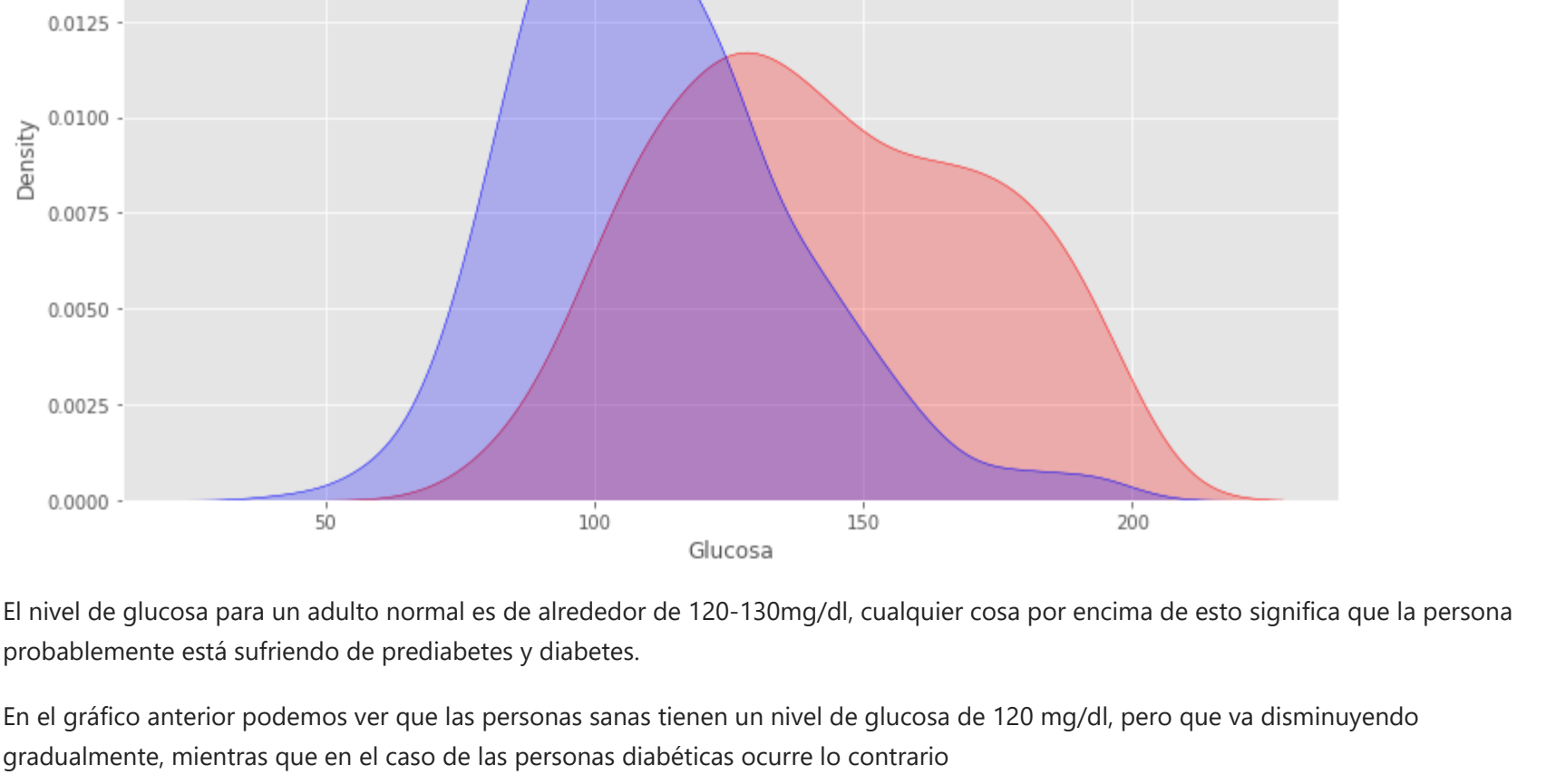
```
In [ ]: data['Outcome'].value_counts()

Out [ ]: 0      500
        1      268
        Name: Outcome, dtype: int64

Analizando y viendo los atributos observamos que disponemos de un dataset desbalanceado, ya que tan solo disponemos de 268 casos con diabetes y 500 casos sin diabetes
```

```
In [ ]: sns.countplot(data['Outcome'], label='count')

Out [ ]: <matplotlib.axes._subplots.AxesSubplot at 0x7f8c0e6361d0>
```



### Datos Faltantes

A la hora de analizar el head() del dataset se puede observar que algunos registros contienen 0 en alguna de sus características. Este hecho nos indica que que son datos faltantes ya que un 0 no tendría sentido en los atributos de este proyecto.

```
In [ ]: data[['Glucose','BloodPressure','SkinThickness','Insulin','BMI']] = data[['Glucose','BloodPressure','SkinThickn
```

```
In [ ]: missing_data = data.isnull().sum()
total_percentage = (missing_data.sum()/data.shape[0]) * 100
print(f'El porcentaje total de datos faltantes es del {round(total_percentage,2)}%')

El porcentaje total de datos faltantes es del 84.9%
```

```
In [ ]: # Porcentaje de datos faltantes por categoria

total = data.isnull().sum().sort_values(ascending=False)
percent_total = (data.isnull().sum()/data.isnull().count()).sort_values(ascending=False)*100
missing = pd.concat([total, percent_total, axis=1, keys=['Total', 'Percentage']]
missing_data = missing[missing['Total']>0]
missing_data

Out [ ]:   Total  Percentage
SkinThickness  227   29.557292
BloodPressure  35   4.557292
BMI           11   1.432292
Glucose        5    0.651042
```



Podemos observar que todos los valores faltantes se han transformado en valores NaN.

Se procede a instalar la librería 'verstack' con el objetivo de poder rellenar los NaN por valores nuevos según correlación entre las características mediante la aplicación de un modelo XGBoost.

```
In [ ]: pip install verstack

Collecting verstack
  Downloading https://files.pythonhosted.org/packages/bc/7e/6319afad95521175557db0f30c31a6edd6cfc795fec27babb1be31/verstack-0.1.1.tar.gz
Requirement already satisfied: pandas in /usr/local/lib/python3.7/dist-packages (from verstack) (1.1.5)
Requirement already satisfied: xgboost in /usr/local/lib/python3.7/dist-packages (from verstack) (0.90)
Requirement already satisfied: sklearn in /usr/local/lib/python3.7/dist-packages (from verstack) (0.90)
Requirement already satisfied: python-dateutil<=2.7.3 in /usr/local/lib/python3.7/dist-packages (from pandas->verstack) (2.8.3)
Requirement already satisfied: pytz<=2017.2 in /usr/local/lib/python3.7/dist-packages (from pandas->verstack) (2018.9)
Requirement already satisfied: scipy in /usr/local/lib/python3.7/dist-packages (from xgboost->verstack) (1.4.1)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.7/dist-packages (from sklearn->verstack) (0.90)
Requirement already satisfied: six<=1.5 in /usr/local/lib/python3.7/dist-packages (from python-dateutil->verstack) (1.15.0)
Requirement already satisfied: joblib<=0.11 in /usr/local/lib/python3.7/dist-packages (from scikit-learn->sklearn->verstack) (1.0.1)
Building wheels for collected packages: verstack
  Building wheel for verstack (setup.py) ... done
  Created wheel for verstack: filename=verstack-0.3.1-cp37-none-any.whl size=14342 sha256=8ae0c94577eb3aa61bfbcad1039b0f6d66369d5428d0b3a6ae216f
  Stored in directory: /root/.cache/pip/wheels/15/1b/58/1de59516150ca4d9b1dbaceb3bfcc0cf42d166fabec2f6
Successfully built verstack
Installing collected packages: verstack
Successfully installed verstack-0.3.1
```

```
In [ ]: from verstack import NaNImputer
imputer = NaNImputer()
new_data = imputer.impute(data)

Dataset dimensions:
- columns: 9
- rows: 768
- mb in memory: 0.05
- NaN cols num: 5

Deploy multiprocessing with 2 parallel processes
- Glucose: imputed 5 NaNs
- SkinThickness: imputed 227 NaNs
- BloodPressure: imputed 35 NaNs
- BMI: imputed 11 NaNs
- Insulin: imputed 374 NaNs

NaNs imputation time: 0.88 minutes
```

```
In [ ]: new_data.head()

Out [ ]:   Pregnancies  Glucose  BloodPressure  SkinThickness  Insulin   BMI  DiabetesPedigreeFunction  Age  Outcome
0             6      148.0             72.0             35.0   79.799479   31.992578              0.471876    50      1
1             1       85.0              66.0             29.0   59.483383   23.6              0.351    31      0
2             8      183.0             64.0             0.0   197.233948   23.3              0.672    32      1
3             1       89.0              66.0             23.0   94.000000   28.1              0.167    21      0
4             0      137.0             40.0             35.0   168.000000   43.1              2.288    33      1
```

```
In [ ]: # Comprobamos que se ha realizado correctamente la asignación de valores faltantes
missing_data = new_data.isnull().sum()
total_percentage = (missing_data.sum()/new_data.shape[0]) * 100
print(f'El porcentaje total de datos faltantes es del {round(total_percentage,2)}%')

El porcentaje total de datos faltantes es del 0.0%
```

### Distribución de las variables

```
In [ ]: new_data_diabetes = new_data[new_data['Outcome'] == 1]
new_data_no_diabetes = new_data[new_data['Outcome'] == 0]
```

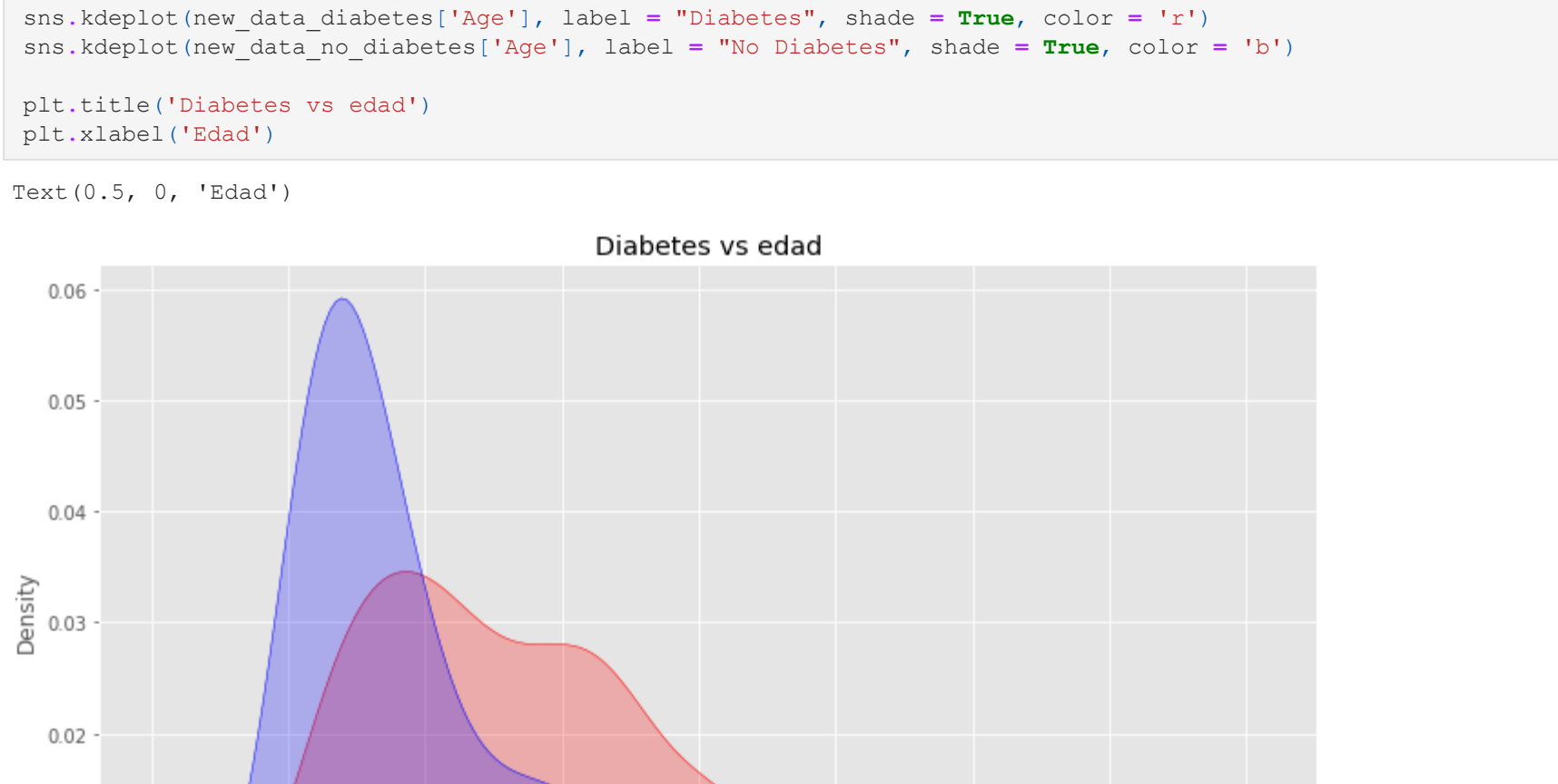
### Diabetes y Glucosa en sangre

```
In [ ]: plt.figure(figsize=(12, 7))

sns.kdeplot(new_data_diabetes['Glucose'], label = "Diabetes", shade = True, color = 'r')
sns.kdeplot(new_data_no_diabetes['Glucose'], label = "No Diabetes", shade = True, color = 'b')

plt.title('Diabetes vs Niveles de Glucosa')
plt.xlabel('Glucose')

Out [ ]: Text(0.5, 0, 'Glucose')
```



El nivel de glucosa para un adulto normal es de alrededor de 120-130mg/dl, cualquier cosa por encima de esto significa que la persona probablemente está sufriendo de prediabetes y diabetes.

En el gráfico anterior podemos ver que las personas sanas tienen un nivel de glucosa de 120 mg/dl, pero que va disminuyendo gradualmente, mientras que en el caso de las personas diabéticas ocurre lo contrario

### Diabetes y embarazos

```
In [ ]: plt.figure(figsize=(12, 7))

sns.kdeplot(new_data_diabetes['Pregnancies'], label = "Diabetes", shade = True, color = 'r')
sns.kdeplot(new_data_no_diabetes['Pregnancies'], label = "No Diabetes", shade = True, color = 'b')

plt.title('Diabetes vs embarazos')
plt.xlabel('embarazos')

Out [ ]: Text(0.5, 0, 'embarazos')
```



Se observa una distribución uniforme entre los dos grupos con frecuencias diferentes. En este sentido se puede afirmar que los embarazos no son la causa de la aparición de la enfermedad de diabetes

### Diabetes y presión sanguínea

```
In [ ]: plt.figure(figsize=(12, 7))

sns.kdeplot(new_data_diabetes['BloodPressure'], label = "Diabetes", shade = True, color = 'r')
sns.kdeplot(new_data_no_diabetes['BloodPressure'], label = "No Diabetes", shade = True, color = 'b')

plt.title('Diabetes vs Presión sanguínea')
plt.xlabel('Presión sanguínea')

Out [ ]: Text(0.5, 0, 'Presión sanguínea')
```



La presión arterial alta (también conocida como "hipertensión") es muy común en personas con diabetes. De hecho, las dos afecciones suelen ir de la mano porque ambas pueden ser consecuencia de los mismos factores de estilo de vida.

La diabetes daña las arterias y las convierte en objetivos de endurecimiento, lo que se denomina aterosclerosis. Esto puede provocar una presión arterial alta que, si no se trata, puede causar problemas como daños en los vasos sanguíneos, infarto de miocardio e insuficiencia renal.

Para una persona normal la PA debe estar en o por debajo de 120/80 mm Hg, la persona con hipertensión puede estar por encima de 139/89 mm Hg.

Del gráfico anterior, podemos decir que, los diabéticos y las personas sanas se distribuyen uniformemente con PA baja y normal pero, hay menos personas sanas que tienen PA alta.

### Diabetes y pliegue de piel

```
In [ ]: plt.figure(figsize=(12, 7))

sns.kdeplot(new_data_diabetes['SkinThickness'], label = "Diabetes", shade = True, color = 'r')
sns.kdeplot(new_data_no_diabetes['SkinThickness'], label = "No Diabetes", shade = True, color = 'b')

plt.title('Diabetes vs pliegue de la piel')
plt.xlabel('pliegue de la piel en tríceps (mm)')

Out [ ]: Text(0.5, 0, 'pliegue de la piel en tríceps (mm)')
```



Los cambios en los vasos sanguíneos debidos a la diabetes pueden provocar una afección en la piel denominada dermatopatía diabética. La dermatopatía aparece en forma de manchas escamosas de color marrón claro o rojo, a menudo en la parte delantera de las piernas. Las manchas no duelen, no tienen ampollas ni pican, y generalmente no es necesario un tratamiento.

En el gráfico anterior, la distribución entre personas sanas y diabéticas es más o menos la misma en cuanto al grosor de la piel.

### Diabetes e insulina

```
In [ ]: plt.figure(figsize=(12, 7))

sns.kdeplot(new_data_diabetes['Insulin'], label = "Diabetes", shade = True, color = 'r')
sns.kdeplot(new_data_no_diabetes['Insulin'], label = "No Diabetes", shade = True, color = 'b')

plt.title('Diabetes vs niveles de insulina')
plt.xlabel('Nivel de insulina en sangre')

Out [ ]: Text(0.5, 0, 'Nivel de insulina en sangre')
```



La insulina es una hormona que produce el páncreas para que las células puedan utilizar la glucosa. Cuando el cuerpo no produce o utiliza la insulina correctamente, se puede tomar insulina artificial para ayudar a controlar el azúcar en sangre. Se pueden utilizar muchos tipos para tratar la diabetes.

La insulina ayuda a controlar los niveles de glucosa en sangre indicando al hígado y a las células musculares y grasas que tomen la glucosa de la sangre. Por lo tanto, la insulina ayuda a las células a tomar la glucosa para utilizarla como energía. Si el cuerpo tiene suficiente energía, la insulina indica al hígado que tome la glucosa y la almacene como glucógeno.

De la gráfica anterior, podemos ver que hay personas diabéticas que aumentan a medida que los niveles de insulina aumentan gradualmente. Hay más personas sanas en torno a los niveles de insulina 0-100.

### Diabetes y IMC

```
In [ ]: plt.figure(figsize=(12, 7))

sns.kdeplot(new_data_diabetes['BMI'], label = "Diabetes", shade = True, color = 'r')
sns.kdeplot(new_data_no_diabetes['BMI'], label = "No Diabetes", shade = True, color = 'b')

plt.title('Diabetes vs IMC')
plt.xlabel('IMC')

Out [ ]: Text(0.5, 0, 'IMC')
```



El sobrepeso (IMC de 25-29.9) o la obesidad (IMC de 30-39.9) o la obesidad mórbida (IMC de 40 o más), aumentan considerablemente el riesgo de desarrollar diabetes de tipo 2. Cuanto más exceso de peso tengas, más resistentes se vuelven las células de tus músculos y tejidos a tu propia hormona de la insulina.

De la gráfica anterior podemos determinar que, a medida que aumenta el IMC disminuye la probabilidad de ser saludable y aumenta la de ser diabético.

### Diabetes función de tipo

```
In [ ]: plt.figure(figsize=(12, 7))

sns.kdeplot(new_data_diabetes['DiabetesPedigreeFunction'], label = "Diabetes", shade = True, color = 'r')
sns.kdeplot(new_data_no_diabetes['DiabetesPedigreeFunction'], label = "No Diabetes", shade = True, color = 'b')

plt.title('Diabetes vs Diabetes Pedigree Function')
plt.xlabel('Diabetes Pedigree Function')

Out [ ]: Text(0.5, 0, 'Diabetes Pedigree Function')
```



La función pedigrí de la diabetes es una función que puntúa la probabilidad de padecer diabetes basándose en los antecedentes familiares. Proporciona algunos datos sobre la historia de la diabetes mellitus en los familiares y la relación genética de esos familiares con el paciente.

De la gráfica anterior, a medida que la función aumenta las personas diabéticas aumentan, mostrando que la diabetes podría ser hereditaria para ese individuo.

### Diabetes y edad

```
In [ ]: plt.figure(figsize=(12, 7))

sns.kdeplot(new_data_diabetes['Age'], label = "Diabetes", shade = True, color = 'r')
sns.kdeplot(new_data_no_diabetes['Age'], label = "No Diabetes", shade = True, color = 'b')

plt.title('Diabetes vs edad')
plt.xlabel('Edad')

Out [ ]: Text(0.5, 0, 'Edad')
```



A medida que la persona envejece, corre un alto riesgo de desarrollar diabetes de tipo 2 debido a los efectos combinados del aumento de la resistencia a la insulina y el deterioro de la función de los islotes pancreáticos con el envejecimiento.

En el gráfico anterior, podemos ver que hay más personas sanas en torno a los 20-25 años, pero a medida que la edad aumenta gradualmente también lo hace la gente que es diabética, lo que demuestra que la edad y la diabetes van de la mano.



### Outliers

Un valor atípico o 'outlier' es una observación que se encuentra a una distancia anormal de otros valores en una muestra aleatoria de una población.

En este cuaderno, estamos utilizando el Box Plot para detectar los valores atípicos de cada característica en personas conjunto de datos, donde cualquier punto por encima o por debajo de los bigotes representa un valor atípico. Esto también se conoce como "método univariante", ya que aquí estamos utilizando un análisis de valores atípicos de una variable.

Se representa con la fórmula IQR = Q3 - Q1. Las líneas de código siguientes calculan e imprimen el rango intercuartil para cada una de las variables del conjunto de datos. La salida anterior tiene las puntuaciones IQR, que pueden utilizarse para detectar valores atípicos.

Después de la detección, utilizamos la Imputación de la Mediana para cuidar de los valores atípicos. En esta técnica, reemplazamos los valores extremos con los valores de la mediana. Se aconseja no utilizar los valores medios, ya que se ven afectados por los valores atípicos.

```
In [ ]: def remove_outliers(data):
arr=[]

q1=np.percentile(data,25)
q3=np.percentile(data,75)
iqr=q3-q1
m=1*(1.5*iqr)
ma=q3+(1.5*iqr)

for i in list(data):
    if i<m1:
        i=mi
        arr.append(i)
    elif i>ma:
        i=ma
        arr.append(i)
    else:
        arr.append(i)
    #print(max(arr))
return arr
```

```
In [ ]: new_data['Glucose'] = remove_outliers(new_data['Glucose'])
new_data['BloodPressure'] = remove_outliers(new_data['BloodPressure'])
new_data['SkinThickness'] = remove_outliers(new_data['SkinThickness'])
new_data['Insulin'] = remove_outliers(new_data['Insulin'])
new_data['BMI'] = remove_outliers(new_data['BMI'])
new_data['Pregnancies'] = remove_outliers(new_data['Pregnancies'])
new_data['Age'] = remove_outliers(new_data['Age'])
new_data['DiabetesPedigreeFunction'] = remove_outliers(new_data['DiabetesPedigreeFunction'])

print('Outliers successfully removed')

Outliers successfully removed
```

```
In [ ]: new_data.shape

Out [ ]: (768, 9)
```

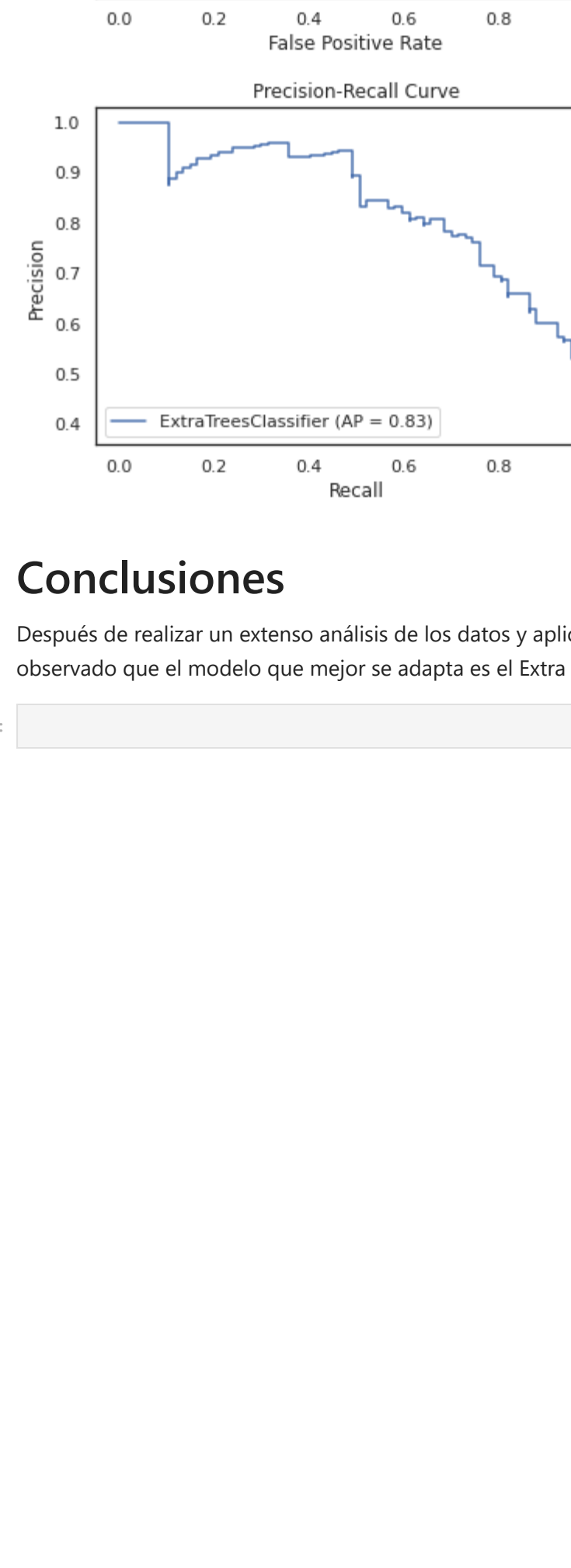
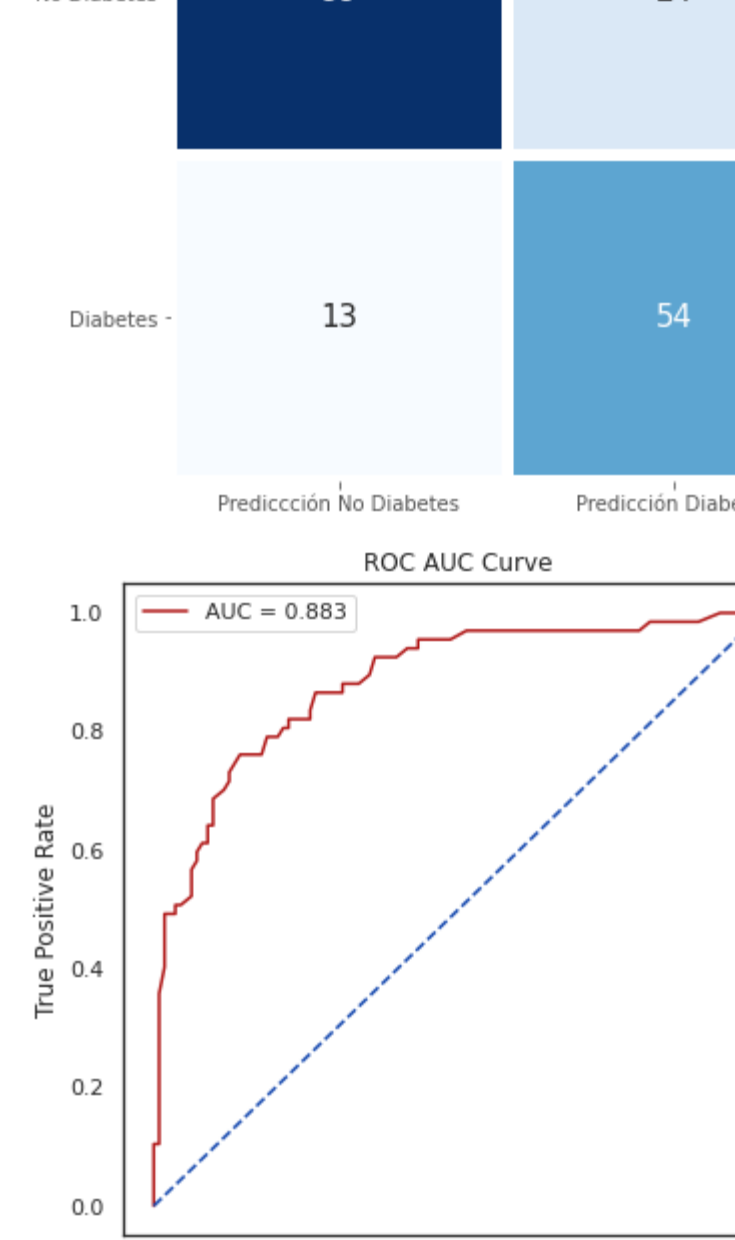






	precision	recall	f1-score	support
0	0.87	0.79	0.83	113
1	0.69	0.81	0.74	67
accuracy			0.79	180
macro avg	0.78	0.80	0.79	180
weighted avg	0.81	0.79	0.80	180

ROC AUC score: 0.8825782591445744  
Accuracy Score: 0.7944444444444444



## Conclusiones

Después de realizar un extenso análisis de los datos y aplicar diferentes modelos y tunear los mismos mediante Grid Search, hemos observado que el modelo que mejor se adapta es el Extra Tree Classifier.