

# PROGRAMACIÓN PARA CIENCIAS DE DATOS 2

Juan Fernando Rodríguez Medina

[https://media.istockphoto.com/id/1364317541/es/foto/Cien%C3%A1dicos-de-datos-marco-del-programador-tocando-y-analizando-el-desarrollo-en-diversas-pyg?si=1024x1024&w=8&s=20&crdtdt\\_7TNJUnw6b5sDQJyUy59ABU7JnD99NgU-YdQ=](https://media.istockphoto.com/id/1364317541/es/foto/Cien%C3%A1dicos-de-datos-marco-del-programador-tocando-y-analizando-el-desarrollo-en-diversas-pyg?si=1024x1024&w=8&s=20&crdtdt_7TNJUnw6b5sDQJyUy59ABU7JnD99NgU-YdQ=)

1. Define un problema específico dentro del campo elegido que pueda abordarse mediante el análisis de datos. Explica el problema y por qué vale la pena investigarlo.

- En el ámbito del sector salud hay muchas problemáticas que afectan a la población, uno de los asuntos más concernientes es que cada año hay una incidencia mayor de enfermedades crónicas que aumentan en gran medida el gasto público, sin embargo, según muchas revistas médicas en casi todos los casos de las enfermedades crónicas estas de habense tratado a tiempo se habría podido evitar la enfermedad final o un peor desenlace, en este trabajo se realizará la comparativa de una serie de relaciones con respecto a la diabetes en una población de un hospital entre algunas de las comparaciones que se harán están serán las siguientes
- ¿Qué relación hay entre el ser diabético (outcome 1 o 0) y el BMI (índice de masa corporal)
- ¿Qué relación hay entre el BMI y la edad?
- ¿Qué relación hay entre el BMI y la diabetes pedigre function (tener diabetes según la historia familiar) Una vez establecida la relación anterior, estudiar qué tan probable es padecer de diabetes dada esta relación puede ser crucial para tener una serie de controles médicos en estas personas para que el riesgo de diabetes sea menor.

1. ¿Cómo el análisis de datos podría contribuir a comprender o resolver el problema definido? ¿Qué ideas espera obtener?

El análisis de datos permite através de herramientas estadísticas y lenguajes de programación establecer una serie de evidencias reproducibles y contextualizadas, que ayuden a resolver problemas mediante hipótesis y procedimientos científicos matemáticos. De esta manera no solo se resuelve el problema, se puede llegar a realizar predicciones que orienten directrices preventivas hacia lo cual, este proyecto está dirigido.

En el análisis a las preguntas (hipótesis) que he propuesto, espero poder hacer los primeros "scratch" a una problemática interesante en sector salud que pueda ser de ayuda para evitar posibles casos de diabetes que se podrían haber evitado, siendo así una reducción de costos en el sector salud y mejoramiento de calidad de vida de muchos posibles pacientes diabéticos.

```
In [68]: # Importamos la librerías que pueden ser útiles en el trabajo y las nombramos
# Importamos la base de datos que está en formato .csv

import pandas as pd
import os
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split, cross_val_score, KFold
from sklearn.metrics import confusion_matrix
from sklearn import tree
diabetes= pd.read_csv(filepath_or_buffer="0:/Bases de datos/input/diabetes.csv")

In [69]: diabetes.head(5)

Out[69]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148	72	35	0	33.6	0.627	50	1
1	1	85	66	29	0	26.6	0.351	31	0
2	8	183	64	0	0	23.3	0.672	32	1
3	1	89	66	23	94	28.1	0.167	21	0
4	0	137	40	35	168	43.1	2.288	33	1

```
In [70]: diabetes.bloodPressure

Out[70]:
```

0	72
1	66
2	64
3	66
4	40
763	76
764	70
765	72
766	60
767	70

Name: BloodPressure, Length: 768, dtype: int64

```
In [71]: diabetes.info()

<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    int64
 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 [72]: diabetes.dropna(inplace=True)
diabetes.info()

<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    int64
 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 [73]: # Revisamos si la base de datos tiene subcategorías que sean datos con solo 1 variable, lo que significaría que esa variable no es de interés

In [74]: cols_cat = ['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
                'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome']

for col in cols_cat:
    print(f"Columna {col}: {diabetes[col].nunique()} subniveles")

Columna Pregnancies: 17 subniveles
Columna Glucose: 136 subniveles
Columna BloodPressure: 47 subniveles
Columna SkinThickness: 51 subniveles
Columna Insulin: 186 subniveles
Columna BMI: 248 subniveles
Columna DiabetesPedigreeFunction: 517 subniveles
Columna Age: 52 subniveles
Columna Outcome: 2 subniveles
```

3. Realiza un cálculo de al menos dos medidas estadísticas descriptivas clave (por ejemplo, media, mediana, desviación estándar) utilizando NumPy. ¿Qué información te proporcionan estas medidas sobre los datos?

```
In [75]: # media mediana y moda con numpy

print("Media: ", np.mean(diabetes))
print("Mediana: ", np.median(diabetes))
# numpy parece que no tiene el módulo para sacar la moda de un df "print("Moda: ", np.mode(diabetes))"
print("Desviación Estándar: ", np.std(diabetes))
print("Varianza: ", np.var(diabetes))

Media: Pregnancies          3.845052
        Glucose            129.884531
        BloodPressure      69.185469
        SkinThickness       29.536458
        Insulin            79.799479
        BMI                 31.992578
        DiabetesPedigreeFunction 0.471876
        Age                 33.248885
        Outcome             0.348958
dtype: float64
Mediana: 24.3
Desviación Estándar: Pregnancies          3.367384
                     Glucose            31.951796
                     BloodPressure      19.343202
                     SkinThickness       15.941229
                     Insulin            115.168949
                     BMI                 7.879926
                     DiabetesPedigreeFunction 0.351113
                     Age                 11.752573
                     Outcome             0.476841
dtype: float64
Varianza: Pregnancies          11.339272
          Glucose            1020.917283
          BloodPressure      374.159449
          SkinThickness       132.63.886875
          Insulin            62.079846
          BMI                 6.109639
          DiabetesPedigreeFunction 0.138.122944
          Age                 0.227186
dtype: float64

<UserWarning: \anaconda3\lib\site-packages\numpy\core\fromnumeric.py:3438: FutureWarning: In a future version, DataFrame.mean(axis=None) will
1 return a scalar mean over the entire DataFrame. To retain the old behavior, use 'frame.mean(axis=0)' or just 'frame.mean()'
return mean(axis=None, dtype=None, out=None, **kwargs)

La información que me proporciona y que resulta más interesante la desviación estándar la cual para Diabetes Pedigree function es mínima se explica porque
los valores son muy pequeños y en los extremos la variación es mínima, y por otra parte vemos que en datos como la glucosa la desviación estándar es
muy alta, lo cual se explica porque los no diabéticos y los diabéticos tienen una diferencia marcada, sin embargo hay que mencionar que en estos valores están
valores incoherentes como podrían ser los 0 valores por encima de los 900, cuando en la teoría se describen hasta 200mg/dl de glucosa en sangre, es decir
existen valores muy discordantes para el dataframe conseguido
```

```
In [76]: # Sacando la mediana, media, moda, mínimo, máximo desviación estándar y los cuantiles usando simplemente la función .describe de pandas
diabetes.describe()

Out[76]:
```

	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	129.884531	69.185469	29.536458	79.799479	31.992578	0.471876	33.248885	0.348958
std	3.369578	31.972618	19.356807	15.952218	115.244002	7.884160	0.331329	11.760232	0.476961
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	0.243750	29.000000	0.000000
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	0.372500	29.000000	0.000000
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.000000	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

```
In [77]: # la desviación estándar nos indica que hay valores dispersos en las columnas

In [78]: print(f"¿Tamaño del set antes de eliminar las filas repetidas: {diabetes.shape}")
# están en rangos diferentes
cols_num = ['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Age']
print(f"¿Tamaño del set después de eliminar las filas repetidas: {diabetes.shape}")

Tamaño del set antes de eliminar las filas repetidas: (768, 9)
Tamaño del set después de eliminar las filas repetidas: (768, 9)
```

4. Cree una visualización gráfica adecuada para explorar los datos. Puede ser un histograma, un gráfico de dispersión u otro tipo relevante. Adjunta la visualización y explica cómo esta visualización puede ayudarte a abordar el problema inicial.

```
In [79]: # Generar diagramas de caja individuales para las variables numéricas
# están en rangos diferentes
cols_num = ['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
            'BMI', 'DiabetesPedigreeFunction', 'Age']

fig, ax = plt.subplots(nrows=8, ncols=1, figsize=(8,38))
fig.subplots_adjust(hspace=0.5)

for i, col in enumerate(cols_num):
    sns.boxplot(x=col, data=diabetes, ax=ax[i])
    ax[i].set_title(col)

Pregnancies

Glucose

BloodPressure

SkinThickness

Insulin

BMI

DiabetesPedigreeFunction

Age

Los gráficos y precisamente para el caso propuesto, el diagrama de caja ayuda a entender fácilmente en qué valores se encuentran los datos según el campo, y
por lo tanto permite hacer un análisis general de cómo están distribuidos los datos y qué tantos datos atípicos existen, es decir todos los datos que salen de los
bigotes. Por otra parte también indicará los datos que se encuentran dentro de la media, y qué tan homogéneos o no, son los datos.
```

```
In [80]: # Graficar los subniveles de cada variable categórica en forma de histograma
cols_cat = ['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Age']
fig, subplots = plt.subplots(nrows=5, ncols=1, figsize=(25,38))
fig.subplots_adjust(hspace=1)

for i, col in enumerate(cols_cat):
    sns.countplot(x=col, data=diabetes, ax=ax[i])
    ax[i].set_title(col)
    ax[i].set_xticklabels(ax[i].get_xticklabels(), rotation=30)

Pregnancies

Glucose

BloodPressure

SkinThickness

Age
```

```
In [81]: # ELIMINANDO LOS REGISTROS INÚTILES O MUY FUERA DE LOS PARÁMETROS DE INTERÉS CUANTITATIVAS

In [82]: print(f"¿Tamaño del set antes de eliminar registros de BMI: {diabetes.shape}")
data = diabetes[diabetes['BMI']>10]
print(f"¿Tamaño del set después de eliminar registros de BMI: {data.shape}")

Tamaño del set antes de eliminar registros de BMI: (768, 9)
Tamaño del set después de eliminar registros de BMI: (767, 9)

In [83]: sns.boxplot(y="BMI", data=data)

<AxesSubplot: ylabel='BMI'>

BMI
```

```
In [84]: print(f"¿Tamaño del set antes de eliminar registros de Presión Sanguínea: {diabetes.shape}")
data2 = diabetes[diabetes['BloodPressure']>24]
print(f"¿Tamaño del set después de eliminar registros de Presión Sanguínea: {data2.shape}")

Tamaño del set antes de eliminar registros de Presión Sanguínea: (768, 9)
Tamaño del set después de eliminar registros de Presión Sanguínea: (763, 9)

In [85]: sns.boxplot(y="BloodPressure", data=data2)

<AxesSubplot: ylabel='BloodPressure'>

BloodPressure
```

```
In [86]: print(f"¿Tamaño del set antes de eliminar registros de Insulina: {diabetes.shape}")
data2 = diabetes[diabetes['Insulin']<300]
print(f"¿Tamaño del set después de eliminar registros de Insulina: {data2.shape}")

Tamaño del set antes de eliminar registros de Insulina: (768, 9)
Tamaño del set después de eliminar registros de Insulina: (731, 9)

In [87]: sns.boxplot(y="Insulin", data=data2)

<AxesSubplot: ylabel='Insulin'>

Insulin
```

5. Teniendo en cuenta los resultados del análisis estadístico y la visualización. ¿Qué patrones o tendencias interesantes puedes observar? ¿Hay algún valor atípico o característica inesperada en los datos?

En el análisis exploratorio que se ha hecho con los datos através de los gráficos de caja de bigotes, vemos que hay en la mayoría de los parámetros muchos
datos exócticos que pueden hacer que futuros procedimientos estadísticos presenten sesgos, por lo tanto enumeraré los pasos a seguir

- Eliminar datos que están muy fuera de los bigotes o no tienen coherencia con los tipos de datos, por ejemplo, no tiene sentido pacientes con 0 de presión sanguínea si todos son pacientes con vida, lo mismo se podría decir con los valores que tienen 0 en grosor de piel entre otros.
- Hay una relación entre tener diabetes y el IMC (o BMI en inglés)
- Para algunos campos como la glucosa en sangre hay una cantidad de subniveles tan grande que valdría la pena categorizarlos en secciones

6. Seleccione dos variables que podrían estar relacionadas dentro del contexto del problema. Convierta las series correspondientes a estas variables en vectores de NumPy e implemente una regresión lineal utilizando numpy.polyfit().

```
In [88]: 1. Outcome + DPF
        2. Outcome + Age
        3. Outcome + BMI

In [89]: bi, b0 = np.polyfit(diabetes["Age"], diabetes["Outcome"], 1)
print(bi, b0)
0.00966634703521698 0.0276241886117164

In [90]: plt.scatter(diabetes["Age"], diabetes["Outcome"], c="orange")
plt.plot(diabetes["Age"], b0+bi*diabetes["Outcome"])
plt.title("Gráfica de dispersión Age vs Outcome")
plt.xlabel("Age")
plt.ylabel("Outcome")
plt.grid()
plt.show()

Out[90]: <function matplotlib.pyplot.show(close=None, block=None)>

In [91]: # Gráfica de dispersión Age vs Outcome

Outcome

Age

In [92]: from sklearn.linear_model import LogisticRegression

# definiendo input y output
X_train = np.array(diabetes["Age"]).reshape((-1, 1))
y_train = np.array(diabetes["Outcome"])

# creando modelo
model = LogisticRegression()
model.fit(X_train, y_train)

# imprimiendo parámetros
print(f"Intercepto (b): {model.intercept_}")
print(f"Pendiente (w): {model.coef_[0]}")

Intercepto (b): [-2.04744865]
Pendiente (w): [[0.04202466]]

In [93]: w = 0.04202466
b = -2.04744865
# puntos de la recta
x = np.linspace(0, diabetes["Age"].max(), 100)
y = 1/(1+np.exp(-(w*x+b)))

# grafica de la recta
diabetes.plot.scatter(x="Age", y="Outcome")
plt.plot(x, y, '-r')
plt.ylim(0, diabetes["Outcome"].max()+1.1)
plt.grid()
plt.show()

Outcome

Age

In [94]: bi, b0 = np.polyfit(diabetes["DiabetesPedigreeFunction"], diabetes["Outcome"], 1)
print(bi, b0)
0.28020566278803 0.238070243542176

In [95]: plt.scatter(diabetes["DiabetesPedigreeFunction"], diabetes["Outcome"], c="orange")
plt.plot(diabetes["DiabetesPedigreeFunction"], b0+bi*diabetes["Outcome"])
plt.title("Gráfica de dispersión Outcome vs BMI")
plt.xlabel("DiabetesPedigreeFunction")
plt.ylabel("Outcome")
plt.grid()
plt.show()

Out[95]: <function matplotlib.pyplot.show(close=None, block=None)>

In [96]: # Gráfica de dispersión Outcome vs BMI

Outcome

DiabetesPedigreeFunction

In [97]: from sklearn.linear_model import LogisticRegression

# definiendo input y output
X_train = np.array(diabetes["DiabetesPedigreeFunction"]).reshape((-1, 1))
y_train = np.array(diabetes["Outcome"])

# creando modelo
model = LogisticRegression()
model.fit(X_train, y_train)

# imprimiendo parámetros
print(f"Intercepto (b): {model.intercept_}")
print(f"Pendiente (w): {model.coef_[0]}")

Intercepto (b): [-1.12809097]
Pendiente (w): [[0.68596889]]

In [98]: w = 1.02697365
b = -1.12809097
# puntos de la recta
x = np.linspace(0, diabetes["BMI"].max(), 100)
y = 1/(1+np.exp(-(w*x+b)))

# grafica de la recta
diabetes.plot.scatter(x="BMI", y="Outcome")
plt.plot(x, y, '-r')
plt.ylim(0, diabetes["Outcome"].max()+1.1)
plt.grid()
plt.show()

Outcome

BMI

In [99]: bi, b0 = np.polyfit(diabetes["BMI"], diabetes["Outcome"], 1)
print(bi, b0)
0.0179853815593836 -0.2175192159743995

In [100]: plt.scatter(diabetes["BMI"], diabetes["Outcome"], c="orange")
plt.plot(diabetes["BMI"], b0+bi*diabetes["Outcome"])
plt.title("Gráfica de dispersión Outcome vs BMI")
plt.xlabel("BMI")
plt.ylabel("Outcome")
plt.grid()
plt.show()

Out[100]: <function matplotlib.pyplot.show(close=None, block=None)>

In [101]: # Gráfica de dispersión Outcome vs BMI

Outcome

BMI
```

8. Considerando la calidad del ajuste de la regresión lineal a los datos, ¿crees que este modelo es suficiente para describir completamente la relación entre las variables? ¿Qué podrías hacer para mejorar la precisión de la predicción?

No es suficiente la regresión Lineal hay que aplicar otro tipo de herramientas de estadísticas para poder realizar un análisis de datos suficiente para hallar una
relación entre los diferentes datos que se quieren estudiar. La otra herramienta estadística que utilice fue la de regresión logística o función sigmoide que permite
predecir clases binarias, de manera que la relación que existe entre dos variables es de sí o no, una vez que pasa o no el 0.5 en los casos del Outcome que me
dice si un paciente tiene o no diabetes. Este procedimiento de comparar el Outcome con las diferentes variables para predecir qué valores me determinen si un
paciente que tiene determinadas características tiene tendencias de padecer diabetes lo realicé con variables como la edad, el BMI y Diabetes Pedigree
Function, sin embargo se puede aplicar a cualquier otro parámetro de interés de estudio.

9. Basándote en el gráfico de los datos y la regresión lineal, ¿qué conclusiones puedes extraer sobre la relación entre las dos variables? ¿Observas una tendencia clara en los datos que respalda la línea de regresión?

En el caso de Outcome vs BMI, Outcome vs Age, Outcome vs DPF se observa una relación que no puede ser explicada a través de una regresión lineal,
significando que los datos no tienen una relación lineal o directamente proporcional aparente. En este orden de ideas, para poder hallar una relación entre las
variables, he realizado la regresión logística, y los resultados plasmados en la gráfica indican que sí hay una relación que podríamos denominar directamente
proporcional

- Outcome vs BMI: los resultados alrededor de entre 40 BMI podemos determinar que el Outcome de ser diabético o no empieza a incrementar bastante cuando este es un punto donde el aumento del 0.5 del Outcome.
- Outcome vs Age: para los resultados entre 40 y 50 años la probabilidad de tener diabetes aumenta considerablemente y a la misma tiene una tendencia exponencial
- Outcome vs DPF: para resultados donde el diabetes pedigree function es 1 o superior, hay una tendencia de considerar que esa persona tiene a tener diabetes.