

Taller 1 – Python y datos: Analítica Computacional de Datos

Nombre: Isabela Castillo Mercado

Código: 201813093

1. Se realizo la Instalación de Python 3.11

```
Microsoft Windows [Versión 10.0.22631.3880]
(c) Microsoft Corporation. Todos los derechos reservados.

C:\Users\fmerc>where python
C:\Users\fmerc\AppData\Local\Microsoft\WindowsApps\python.exe

C:\Users\fmerc>_
```

(*) Aclaración: Mi usuario en este pc es fmerc.

2. Dentro de Visual Studio Code realizamos todo lo mencionado em el enunciado del taller y la instalación de los paquetes correspondientes.

```
PS C:\Users\fmerc\OneDrive\Escritorio\Analitica Computacional> py -3.11 -m pip install numpy scipy pandas matplotlib seaborn
Collecting numpy
  Downloading numpy-2.0.1-cp311-cp311-win_and64.whl.metadata (60 kB)
    Downloading numpy-2.0.1-cp311-cp311-win_and64.whl (399.3 kB)
    etas 0:00:00
Collecting scipy
  Downloading scipy-1.14.0-cp311-cp311-win_and64.whl.metadata (60 kB)
    Downloading scipy-1.14.0-cp311-cp311-win_and64.whl (16.6 MB)
    etas 0:00:00
Collecting pandas
  Downloading pandas-2.2.2-cp311-cp311-win_and64.whl.metadata (19 kB)
Collecting matplotlib
  Downloading matplotlib-3.9.1.post1-cp311-cp311-win_and64.whl.metadata (11 kB)
Collecting seaborn
  Downloading seaborn-0.13.2-py3-none-any.whl.metadata (5.4 kB)
Requirement already satisfied: python-dateutil>=2.8.2 in c:\users\fmerc\appdata\roaming\python\python311\site-packages (from pandas) (2.9.0.post0)
Collecting pytz>=2020.1 (from pandas)
  Downloading pytz-2024.1-py2.py3-none-any.whl.metadata (22 kB)
Collecting tzdata>=2022.7 (from pandas)
  Downloading tzdata-2024.1-py2.py3-none-any.whl.metadata (1.4 kB)
Collecting contourpy>=1.0.1 (from matplotlib)
  Downloading contourpy-1.2.1-cp311-cp311-win_and64.whl.metadata (5.8 kB)
Collecting cycler>=0.10 (from matplotlib)
  Downloading cycler-0.12.1-py3-none-any.whl.metadata (3.8 kB)
Collecting fonttools>=4.22.0 (from matplotlib)
  Downloading fonttools-4.53.1-cp311-cp311-win_and64.whl.metadata (165 kB)
    Downloading fonttools-4.53.1-cp311-cp311-win_and64.whl (2.5 MB)
    etas 0:00:00
Collecting kdisolver>=1.3.1 (from matplotlib)
  Downloading kdisolver-1.4.5-cp311-cp311-win_and64.whl.metadata (6.5 kB)
Requirement already satisfied: pillow>=20.0 in c:\users\fmerc\appdata\roaming\python\python311\site-packages (from matplotlib) (24.1)
Collecting pillow>=8 (from matplotlib)
  Downloading pillow-10.4.0-cp311-cp311-win_and64.whl.metadata (9.3 kB)
Collecting pyparsing>=2.3.1 (from matplotlib)
  Downloading pyparsing-3.1.2-py3-none-any.whl.metadata (5.1 kB)
Requirement already satisfied: six>=1.5 in c:\users\fmerc\appdata\roaming\python\python311\site-packages (from python-dateutil>=2.8.2->pandas) (1.16.0)
Downloaded numpy-2.0.1-cp311-cp311-win_and64.whl (16.6 MB)
    Downloading scipy-1.14.0-cp311-cp311-win_and64.whl (44.7 MB)
    Downloading pandas-2.2.2-cp311-cp311-win_and64.whl (11.0 MB)
    Downloading matplotlib-3.9.1.post1-cp311-cp311-win_and64.whl (8.0 MB)
    Downloading seaborn-0.13.2-py3-none-any.whl (294 kB)
    Downloading contourpy-1.2.1-cp311-cp311-win_and64.whl (188 kB)
    Downloading cycler-0.12.1-py3-none-any.whl (2.2 MB)
    Downloading fonttools-4.53.1-cp311-cp311-win_and64.whl (28.1 MB)
    Downloading kdisolver-1.4.5-cp311-cp311-win_and64.whl (56 kB)
    Downloading pillow-10.4.0-cp311-cp311-win_and64.whl (183 kB)
    Downloading pyparsing-3.1.2-py3-none-any.whl (103.2 kB)
    Downloading pytz-2024.1-py2.py3-none-any.whl (505 kB)
    Downloading tzdata-2024.1-py2.py3-none-any.whl (345 kB)
Installing collected packages: pytz, tzdata, pyparsing, pillow, numpy, kdisolver, fonttools, cycler, scipy, pandas, contourpy, matplotlib, seaborn
```

Como se puede observar en la salida del terminal, tuve que especificar la versión python la que instalaba los paquetes, debido a que ya hacia uso de pyehon en Visual Studio Code y tenia la versión 3.92

```

Location: C:\Users\fmmerc\AppData\Local\Programs\Python\Python311\Lib\site-packages
Requires: numpy
Required-by: statsmodels
---
Name: pandas
Version: 2.2.2
Summary: Powerful data structures for data analysis, time series, and statistics
Home-page: https://pandas.pydata.org
Author:
Author-email: The Pandas Development Team <pandas-dev@python.org>
License: BSD 3-Clause License

Copyright (c) 2008-2011, AQR Capital Management, LLC, Lambda Foundry, Inc. and PyData Development Team
All rights reserved.

---
Name: matplotlib
Version: 3.9.1.post1
Summary: Python plotting package
Home-page:
Author: John D. Hunter, Michael Droettboom
Author-email: Unknown <matplotlib-users@python.org>
License: license agreement for matplotlib versions 1.3.0 and later
=====
Name: seaborn
Version: 0.13.2
Summary: Statistical data visualization
Home-page:
Author:
Author-email: Michael Waskom <mwaskom@gmail.com>
License:
Location: C:\Users\fmmerc\AppData\Local\Programs\Python\Python311\Lib\site-packages
Requires: matplotlib, numpy, pandas
Required-by:
PS C:\Users\fmmerc\OneDrive\Escritorio\Analitica Computacional>

Name: pandas
Version: 2.2.2
Summary: Powerful data structures for data analysis, time series, and statistics
Home-page: https://pandas.pydata.org
Author:
Author-email: The Pandas Development Team <pandas-dev@python.org>
License: BSD 3-Clause License

Copyright (c) 2008-2011, AQR Capital Management, LLC, Lambda Foundry, Inc. and PyData Development Team
All rights reserved.

Name: scipy
Version: 1.11.4
Summary: SciPy: Scientific Library for Python
Home-page: https://www.scipy.org
Author:
Author-email:
License: BSD
Location: c:\users\fmmerc\anaconda3\lib\site-packages
Requires: numpy
Required-by: cupy, datashader, ecos, gensim, osp, pyportfolioopt, qidi, schitt-image, schitt-learn, schitt-plot, scs, shap, statsmodels, xgboost

```

2. Exploración de Datos en Python

```
import numpy as np
import pandas as pd
```

Se importan las librerías numpy como np y pandas como pd para su uso en el resto del cuaderno.

```
df = pd.read_csv('..\\BankChurn.csv')
```

Se lee e importa el archivo CSV como un DataFrame utilizando la función read_csv de la librería pandas (pd).

```
df.shape
```

Muestra la dimensión de la data frame

```
df.head()
```

	Attrition_Flag	Customer_Age	Gender	Dependent_count	Education_Level	Marital_Status	Income_Category	Card_Category	Months_on_book	Total_Relationship_Count	Months_Inactive_12_mon	Contacts_Count_12_mon	Credit_Limit	Total_Revolving_Bal	Avg_Credit_Utilization
0	Existing Customer	45	M	3	High School	Married	\$0K - \$10K	Blue	39	5	1	3	12691.0	777	0.000000
1	Existing Customer	49	F	5	Graduate	Single	Less than \$40K	Blue	44	6	1	2	8256.0	864	0.000000
2	Existing Customer	51	M	3	Graduate	Married	\$10K - \$120K	Blue	36	4	1	0	3418.0	0	0.000000
3	Existing Customer	40	F	4	High School	Unknown	Less than \$40K	Blue	34	3	4	1	3313.0	2517	0.000000
4	Existing Customer	40	M	3	Uneducated	Married	\$0K - \$10K	Blue	21	5	1	0	4716.0	0	0.000000

La función head muestra las primeras 5 filas con todas las columnas del DataFrame, comenzando con el índice cero, ya que Python utiliza indexación basada en cero.

```
df["Attrition_Flag"].unique()
```

```
array(['Existing Customer', 'Attrited Customer'], dtype=object)
```

La función unique en Python muestra los valores únicos de una columna, eliminando duplicados. Dado que la columna Attrition_Flag son 2 categorías, te devolverá solo los valores únicos.

```
df.groupby(["Attrition_Flag"]).count()
```

	Customer_Age	Gender	Dependent_count	Education_Level	Marital_Status	Income_Category	Card_Category	Months_on_book	Total_Relationship_Count	Months_Inactive_12_mon	Contacts_Count_12_mon	Credit_Limit	Total_Revolving_Bal	Avg_Credit_Utilization
Attrited Customer	1627	1627	1627	1627	1627	1627	1627	1627	1627	1627	1627	1627	1627	1627
Existing Customer	8500	8500	8500	8500	8500	8500	8500	8500	8500	8500	8500	8500	8500	8500

La función groupby agrupa por categoría y cuenta la cantidad de filas que tiene cada categoría en la columna Attrition_Flag.

```
df["Attrition_Flag"].value_counts()
```

Attrition_Flag	count
Existing Customer	8500
Attrited Customer	1627

Name: count, dtype: int64

La función value_counts muestra el número de observaciones por categoría en Attrition_Flag, asociando cada una con su respectiva frecuencia.

```
df.describe()
```

	Customer_Age	Dependent_count	Months_on_book	Total_Relationship_Count	Months_Inactive_12_mon	Contacts_Count_12_mon	Credit_Limit	Total_Revolving_Bal	Avg_Open_To_Buy	Total_Amt_Chng_Q4_Q1	Total_Trans_Amt	Total_Trans_Ct	Total_Ct_Chng
count	10127.000000	10127.000000	10127.000000	10127.000000	10127.000000	10127.000000	10127.000000	10127.000000	10127.000000	10127.000000	10127.000000	10127.000000	10127.000000
mean	46.325960	2.346203	35.929409	3.912590	2.341167	2.455317	8631.953698	1162.814061	7469.139637	0.759941	4404.086104	64.858695	0.000000
std	8.016014	1.298908	7.986416	1.554408	1.010622	1.106225	9088.776650	814.987335	9090.695324	0.219207	3197.129254	23.472570	0.000000
min	26.000000	0.000000	13.000000	1.000000	0.000000	0.000000	1438.100000	0.000000	3.000000	0.000000	510.000000	10.000000	0.000000
25%	41.000000	1.000000	31.000000	3.000000	2.000000	2.000000	2555.000000	359.000000	1124.500000	0.631000	2155.500000	45.000000	0.000000
50%	46.000000	2.000000	36.000000	4.000000	2.000000	2.000000	4549.000000	1276.000000	3434.000000	0.736000	3899.000000	67.000000	0.000000
75%	52.000000	3.000000	40.000000	5.000000	3.000000	3.000000	11067.500000	1784.000000	9859.000000	0.859000	4741.000000	81.000000	0.000000
max	79.000000	5.000000	56.000000	6.000000	6.000000	6.000000	34516.000000	2517.000000	34516.000000	3.397000	18494.000000	139.000000	0.000000

La función describe proporciona un análisis descriptivo de cada columna con atributos cuantitativos, mostrando el número de observaciones (count), el promedio (mean), la desviación estándar (std), el valor mínimo (min), los cuartiles del 25%, 50% y 75%, y el valor máximo (max) dentro de la columna.

```
# número de clientes perdido y no perdidos
counts = df.Attrition_Flag.value_counts()
perc_churn = (counts.iloc[1]) / (counts.iloc[0] + counts.iloc[1]) * 100
print(f"Churn_Rate = {perc_churn:1f}%")
```

La función calcula la tasa de abandono como el porcentaje de Attrited Customer sobre el total de clientes (la suma de Existing Customer y Attrited Customer), multiplicado por 100. Utiliza iloc para seleccionar valores específicos en el DataFrame count (almacena el conteo por categoría de la columna Attrition_Flag utilizando la función value_counts), donde Existing Customer está en la fila 1 y Attrited Customer en la fila 0. Finalmente, imprime el valor

calculado en `perc_churn` junto con el texto "Churn Rate", formateado como porcentaje con un decimal.

```
# número de duplicados
duplicates = len(df[df.duplicated()])
print("Number of Duplicate Entries: {duplicates}")

[18]: Number of Duplicate Entries: 0
```

Calcula el número de duplicados (Calcula el largo de un arreglo) que se encuentran en toda la data frame de igual manera imprimimos el resultado con la función `print`

```
# número de valores perdidos
missing_values = df.isnull().sum().sum()
print("Number of Missing Values: {missing_values}")

[11]: Number of Missing Values: 0
```

La función `is_null()` identifica los valores nulos en cada celda del DataFrame, asignando un valor booleano: True para nulos (representado internamente como 1) y False para no nulos (representado como 0). Primero, se suman los valores nulos por columna, y luego se suman estas sumas para obtener el total de valores faltantes en todo el DataFrame. Finalmente, este total se imprime junto con un texto descriptivo usando `print`.

```
# Tipos de datos en el dataset
types = df.dtypes.value_counts()

print("Number of Features: %d"%(df.shape[1]))
print("Number of Customers: %d"%(df.shape[0]))
print("Data Types and Frequency in Dataset:")
print(types)

[11]: Number of Features: 20
      Number of Customers: 10127
      Data Types and Frequency in Dataset:
      int64      9
      object      6
      float64     5
      Name: count, dtype: int64
```

Esta función cuenta y muestra la cantidad de columnas, filas y la frecuencia de cada tipo de dato en el DataFrame `df`.

```
# Conversión de características
df['gender'] = df['gender'].map({'M': 1, 'F': 0})
df['Attrition_Flag'] = df['Attrition_Flag'].map({'Attrited Customer': 1, 'Existing Customer': 0})

[19]: 0/0
```

Esta función convierte los valores categóricos de las columnas `Gender` y `Attrition_Flag` en valores numéricos (1 y 0) en el DataFrame `df`.

```
catcols = df.select_dtypes(exclude = ['int64', 'float64']).columns
intcols = df.select_dtypes(include = ['int64']).columns
floatcols = df.select_dtypes(include = ['float64']).columns

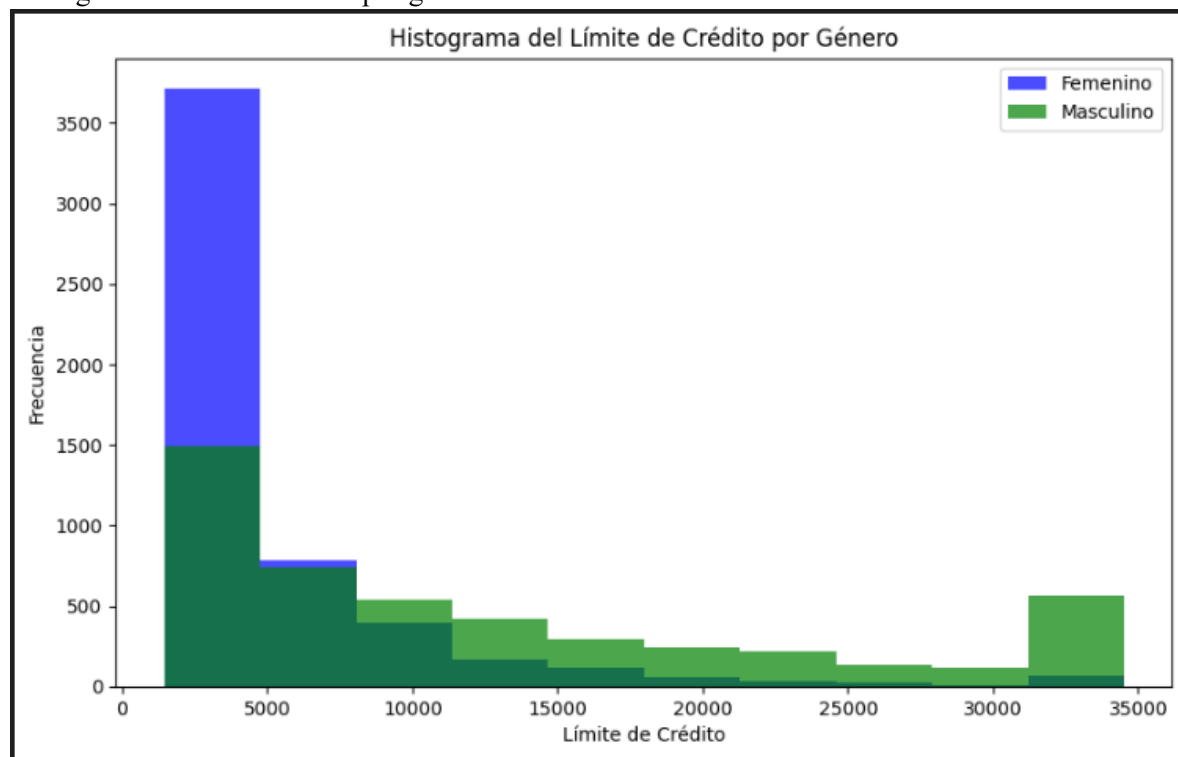
# codificación
df = pd.get_dummies(df, columns = catcols)

print("New Number of Features: %d"%(df.shape[1]))

New Number of Features: 37
```

Esta función identifica las columnas categóricas, enteras y de punto flotante en `df`, aplica codificación one-hot (`get_dummies`) a las categóricas, y luego imprime el nuevo número de columnas del DataFrame.

Histograma limite de Credito por genero



Realice un histograma, segmentado por genero para tener una noción de si la concentración del límite de créditos se ve impactado por el género. Si bien podemos observar que en los valores más menores de crédito se concentran mas hombres que mujeres, de lo contrario en valores más altos de crédito se concentran más mujeres que hombres.

El análisis, no tiene un sustento estadístico, por lo que simplemente es una noción cualitativa de la grafica

El análisis de la BBDD BikePrices.csv se encuentra como archivo de soporte BikePrices.ipynb

3. Números Aleatorios y Bondad de Ajuste en Python

Explicación de cada celda de codigo del Notebook de data-gen.ipynb

```
Generación de muestras de una variable aleatoria normal usando numpy.

import numpy as np

mu = 3
sigma = 0.5
n=1000
vals = np.random.normal(loc=mu, scale=sigma, size=n)
print(vals)
```

La función de numpy genera un arreglo de 1000 números (n), ajustado a una normal con media 3 y desviación estándar de 0.5. y después se imprime el arreglo

```
import pandas as pd

df = pd.DataFrame(vals)
print(df.describe())
```

count	1000.000000
mean	3.022082
std	0.422086
min	1.281704
25%	2.688698
50%	3.022621
75%	3.351697
max	4.747871

Realiza un analisis descriptivo de el arreglo vals que se genero

```
import matplotlib.pyplot as plt
count, bins, ignored = plt.hist(x=vals, bins=30)
plt.title('Histograma de tiempos de servicio')
plt.xlabel('Tiempos de servicio')
plt.ylabel('Frecuencia')
plt.show()

print(bins)
print(count)
```

La función realiza un histograma, con particiones de 30, y el comando plt llama las función de matplotlib para asignar titulo, eje x y eje y.El Histograma replica la forma de una Normal.

```
from scipy.stats import norm
import numpy as np
import matplotlib.pyplot as plt

# Normal (mu=0, sigma=1)
x = np.arange(-4,4,0.001)
plt.plot(x, norm.pdf(x))
plt.title('Función de densidad normal estándar')
plt.xlabel('Valores')
plt.ylabel('Densidad')
plt.show()
```

Esta función genera y visualiza la función de densidad de probabilidad de una distribución normal estándar (media 0, desviación estándar 1) usando matplotlib.

```
from scipy.stats import norm
import numpy as np
import matplotlib.pyplot as plt

# Normal (mu=3, sigma=1)
x = np.arange(0.5, 6, 0.001)
plt.plot(x, norm.pdf(x, loc=3))
plt.title('Función de densidad normal(mu=3, sigma=1)')
plt.xlabel('Valores')
plt.ylabel('Densidad')
plt.show()
```

Esta función genera y visualiza la función de densidad de probabilidad de una distribución normal con media 3 y desviación estándar 1 usando matplotlib.

```

from scipy.stats import norm
import numpy as np
import matplotlib.pyplot as plt

# Normal (mu=0, sigma=1)
x = np.arange(0,4,0.001)
plt.plot(x, norm.pdf(x, loc=3, scale=0.5))
plt.title('Función de densidad normal(mu=3, sigma=0.5)')
plt.xlabel('Valores')
plt.ylabel('Densidad')
plt.show()

```

Esta función genera y visualiza la función de densidad de probabilidad de una distribución normal con media 3 y desviación estándar 0.5 usando matplotlib.

QQ plot usando statsmodels.

```

import statsmodels.api as sm
import scipy.stats as stats
fig = sm.qqplot(data = vals, dist=stats.norm, loc=3, scale=0.5, line='45')
plt.show()

```

Este código genera un gráfico Q-Q (cuantiles-cuantiles) para comparar la distribución de los datos en vals con una distribución normal con media 3 y desviación estándar 0.5, utilizando statsmodels.

```

import statsmodels.api as sm
import scipy.stats as stats
fig = sm.qqplot(data = vals, dist=stats.norm, loc=3, scale=1, line='45')
plt.show()

```

Este código genera un gráfico Q-Q (cuantiles-cuantiles) para comparar la distribución de los datos en vals con una distribución normal con media 3 y desviación estándar 1, utilizando statsmodels.

```

import statsmodels.api as sm
import scipy.stats as stats
fig = sm.qqplot(data = vals, dist=stats.norm, loc=0, scale=1, line='45')
plt.show()

```

Este código genera un gráfico Q-Q (cuantiles-cuantiles) para comparar la distribución de los datos en vals con una distribución normal con media 0 y desviación estándar 1, utilizando statsmodels.

```

import matplotlib.pyplot as plt
# Histograma
count, bins, ignored = plt.hist(x=vals, bins=30, density=True)

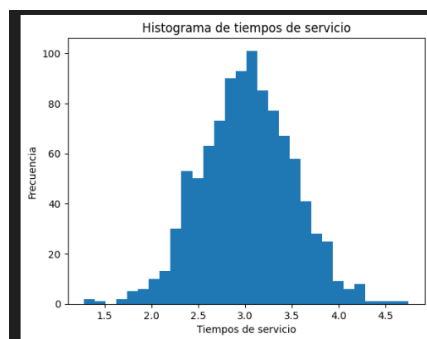
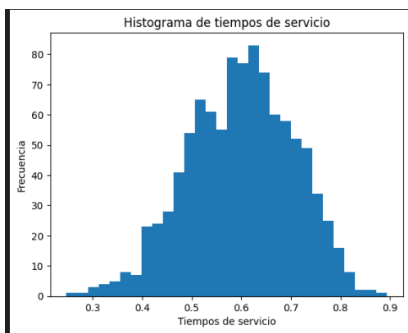
# densidad normal(3, 0.5)
x = np.arange(1,5,0.001)
plt.plot(x, norm.pdf(x, loc=3, scale=0.5))

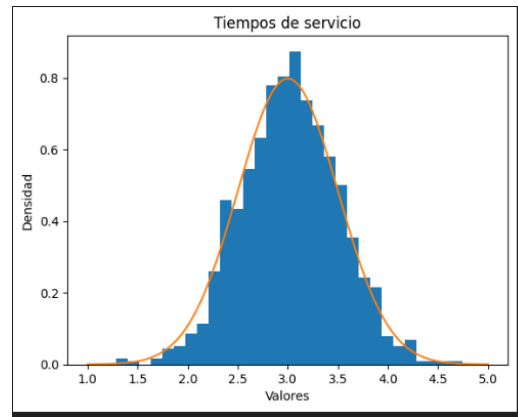
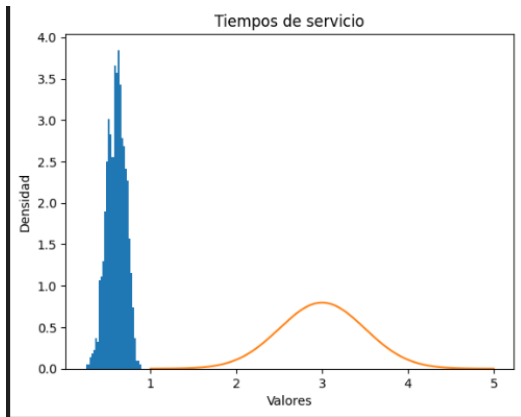
plt.title('Tiempo de servicio')
plt.xlabel('Valores')
plt.ylabel('Densidad')
plt.show()

```

Este código crea un histograma de los datos en vals con 30 intervalos y superpone la función de densidad de una distribución normal con media 3 y desviación estándar 0.5, y luego muestra el gráfico con matplotlib.

Cambio a distribución beta





Vemos que en el cambio de distribución la grafica un poco con tipo de rezago hacia la derecha, haciendo una contraste de estas graficas podemos observar que una grafica se ajusta a la normal y se observa como la distribución beta no se ajusta a la normal