



UNIVERSIDAD AUTÓNOMA DE QUERÉTARO
FACULTAD DE INGENIERÍA
DIVISIÓN DE INVESTIGACIÓN Y POSGRADO

Machine Learning: Práctica 2. Imputación de datos.

Alumno:

Ing. Enrique Mena Camilo

Profesor:

Dr. Marco Antonio Aceves Fernández

Marzo 2023



Índice

1	Objetivo	1
2	Introducción	2
3	Marco teórico	3
3.1	Técnicas de imputación por información externa o deductiva	3
3.2	Técnicas deterministas de imputación	3
3.3	Técnicas estocásticas de imputación	3
4	Materiales y métodos	4
4.1	Conjunto de datos utilizado	4
4.2	Generación de datos faltantes	4
4.3	Procedimiento para imputación de datos	4
5	Resultados	5
5.1	Porcentaje de datos faltantes	5
5.2	Histogramas con datos originales	5
5.3	Mapa de calor de correlación	7
5.4	Histogramas con datos imputados	8
6	Conclusiones	9
	Referencias bibliográficas	11
A	Código documentado	12
A.1	DeleteData.py	12
A.2	Practica2_EnriqueMenaCamilo.py	13



1. Objetivo

Desarrollar un conjunto de herramientas propias que sirvan para realizar imputación de datos a atributos categóricos y continuos.

Se deberá implementar como mínimo dos métodos de imputación, así como procurar usar por lo menos 1 atributo categórico y 1 atributo continuo.



2. Introducción

Es de conocimiento general que todo modelo de inteligencia artificial es tan bueno como los datos con los que fue entrenado.

Dentro de los problemas más comunes que pueden presentarse en los datos se destacan: datos faltantes, problemas de cardinalidad irregular y valores atípicos [1].

Hablando específicamente del problema de datos faltantes, podemos encontrarnos con diversas causas de dicho problema, desde errores en el sistema de adquisición de datos hasta errores humanos [1].

Para solucionar el problema de datos faltantes existen diversos métodos, desde el proceso de eliminar las instancias que cuenten con datos faltantes o bien eliminar los atributos, sin embargo, cuando la cantidad de datos faltantes representa menos del 30 % del total de instancias del atributo, suele ser prudente utilizar algún método de imputación de datos [1].



3. Marco teórico

3.1. Técnicas de imputación por información externa o deductiva

Consisten en deducir los datos faltantes mediante reglas pre-establecidas que contemplan a instancias completas [1].

3.2. Técnicas deterministas de imputación

Útiles cuando de acuerdo a las mismas condiciones de los datos se producen las mismas respuestas.

- **Imputación por regresión.** Se ajusta un modelo lineal o polinomial usando atributos sin datos faltantes. Los datos faltantes se toman del ajuste resultante de la regresión [1].
- **Imputación de la media (o moda).** Los datos faltantes se llenan con la media (o moda en el caso de atributos cualitativos) de las instancias no faltantes [1].
- **Imputación por media de clases.** Se calcula la media (o moda) de las instancias que tienen valor por cada clase y se llena el valor faltante para cada una de las clases [1].
- **Imputación por vecino más cercano.** Se calcula la distancia entre la instancia a imputar y los datos que tienen valor establecido. El dato más cercano será el utilizado para imputar la instancia faltante.

3.3. Técnicas estocásticas de imputación

Se definen como aquellas técnicas de imputación que al repetirse bajo las mismas condiciones producen resultados diferentes.

- **Imputación aleatoria.** Se toman las posibles observaciones del atributo con datos faltantes, se selecciona un valor dentro del rango existente y se llena el dato faltante con dicha elección aleatoria [1].
- **Imputación secuencial.** Consiste en tomar de manera secuencial los datos existentes para reemplazar los datos faltantes. Se toma de manera aleatoria un dato existente y se utiliza dicho valor para reemplazar el primer dato faltante, posteriormente se toma el siguiente valor existente y se utiliza para reemplazar el siguiente dato faltante, el proceso se repite hasta haber llenado todos los datos faltantes [1].



4. Materiales y métodos

Para el desarrollo de esta práctica se utilizó el lenguaje de programación Python en su versión 3.10, con el que se diseñó un script para cumplir los objetivos de la práctica.

4.1. Conjunto de datos utilizado

El conjunto de datos utilizado para esta práctica consta de una colección cuyo objetivo es la predicción de un derrame cerebral (stroke, en inglés). Dicho conjunto de datos fue obtenido de la plataforma Kaggle, y consta de 10 atributos de diversos tipos y 1 variable objetivo, teniendo un total de *40,910* instancias.

4.2. Generación de datos faltantes

Originalmente, el conjunto de datos utilizado solamente contaba con 3 instancias faltantes de un único atributo, por lo que se diseñó un script en Python que borra de forma aleatoria un porcentaje dado de datos del conjunto de datos.

4.3. Procedimiento para imputación de datos

Teniendo el conjunto de datos listo para aplicar los métodos de imputación, se procedió a realizar la siguiente lista de pasos:

1. Estimar el porcentaje de datos faltantes.
2. Generar histogramas de todas los atributos del conjunto de datos.
3. Generar mapa de calor de correlación del conjunto de datos.
4. Elegir técnica de imputación adecuada para cada atributo.
5. Realizar imputación de datos faltantes.
6. Generar histogramas de todos los atributos del conjunto de datos con los datos imputados.

Dichos pasos fueron implementados dentro de un script de Python que implementa los métodos de imputación elegidos, así como generar gráficos y mostrar resultados en terminal.

5. Resultados

5.1. Porcentaje de datos faltantes

Se determinó que el conjunto de datos utilizado cuenta con 5 atributos con datos faltantes, todos con un porcentaje menor al 30 %, por lo que se puede aplicar técnicas de imputación a todos los atributos con datos faltantes. La Tabla 1 muestra el detalle de los datos faltantes encontrados.

Tabla 1: Porcentaje de datos faltantes por atributo.

Columna	Total de instancias	Instancias con datos faltantes	Porcentaje de datos faltantes
sex	40910	3	0.01 %
age	40910	0	0.00 %
hypertension	40910	5676	16.11 %
heart_disease	40910	0	0.00 %
ever_married	40910	0	0.00 %
work_type	40910	5675	16.11 %
Residence_type	40910	5704	16.20 %
avg_glucose_level	40910	0	0.00 %
bmi	40910	5720	16.25 %
smoking_status	40910	0	0.00 %
stroke	40910	0	0.00 %

5.2. Histogramas con datos originales

Observando los histogramas de todas las columnas del conjunto de datos (mostrados en la Figura 1 podemos observar lo siguiente de los atributos con datos faltantes.

- **sex, hypertension, Residence_type**
 - Los datos faltantes son categóricos.
 - Los datos faltantes sólo pueden tomar los valores 0 o 1.
- **work_type**
 - Los datos faltantes son categóricos.
 - Los datos faltantes sólo pueden tomar los valores 1, 2, 3, o 4.

- **bmi**
- Los datos faltantes son continuos.

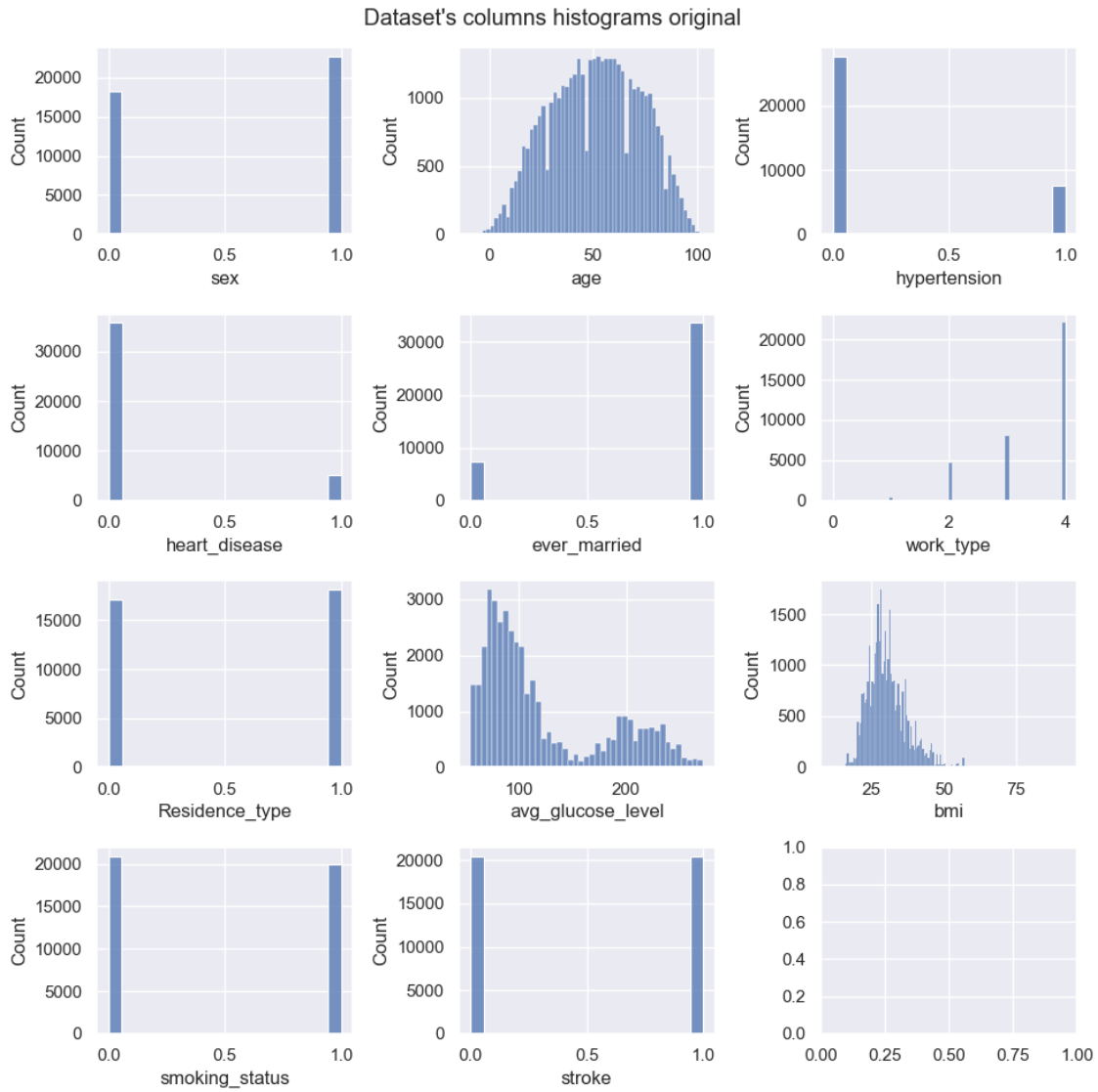


Figura 1: Histogramas de datos originales.

5.3. Mapa de calor de correlación

Al analizar el mapa de calor de correlación (Figura 2) se puede observar lo siguiente:

- Los atributos *work_type* y *Residence_type* presentan índices de correlación cercanos a cero, por lo que la *imputación aleatoria* podría aplicarse correctamente.
- El atributo *hypertension* consigue su mayor índice de correlación con la variable objetivo (*stroke*), por lo que una imputación por *moda de clases* es una elección apropiada.
- El atributo *bmi* consigue su mayor índice de correlación con el atributo *avg_glucose_level*, y dado que ambos son atributos de tipo continuo la *imputación por correlación* presenta una alternativa tentadora.

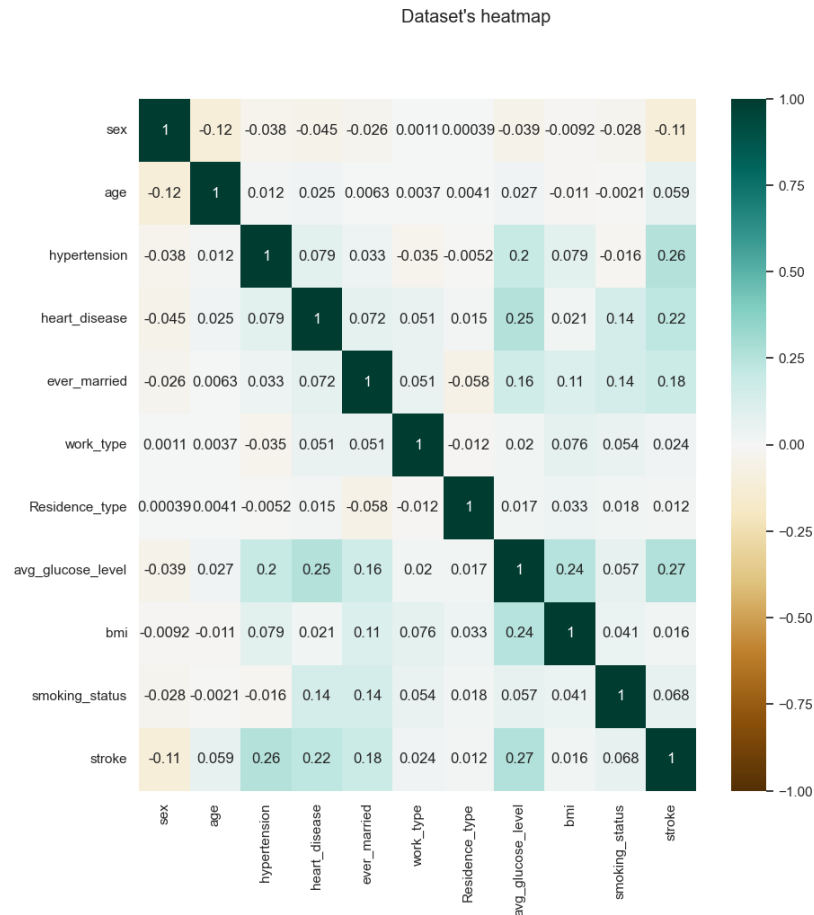


Figura 2: Mapa de calor de correlación entre columnas del conjunto de datos.

5.4. Histogramas con datos imputados

Prestando especial atención a los atributos que se sometieron a métodos de imputación, se puede observar que el procedimiento no ha cambiado la forma de la distribución de los datos, lo cual es un indicador de una correcta implementación de los métodos de imputación.

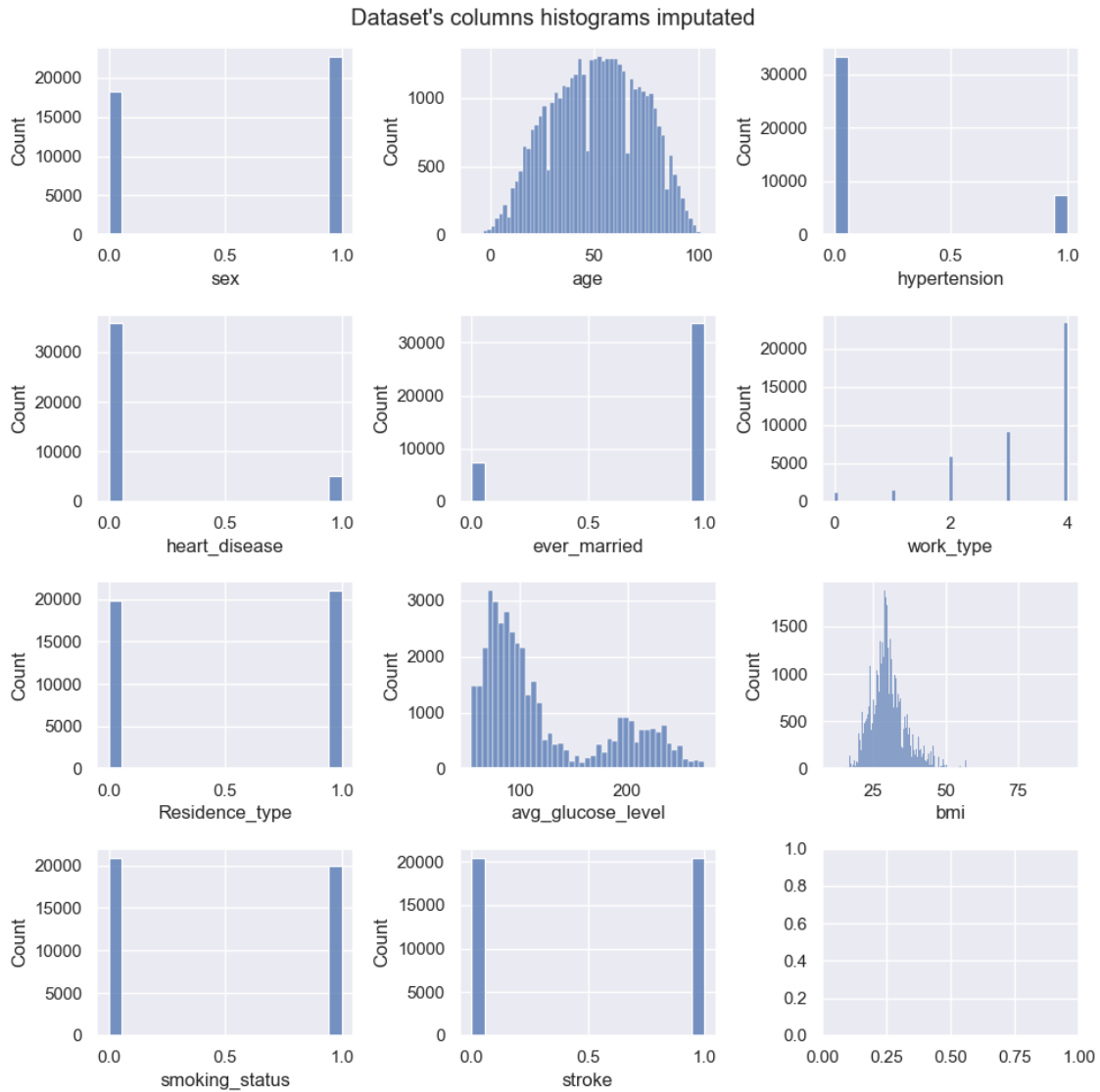


Figura 3: Histogramas de datos originales con datos imputados.

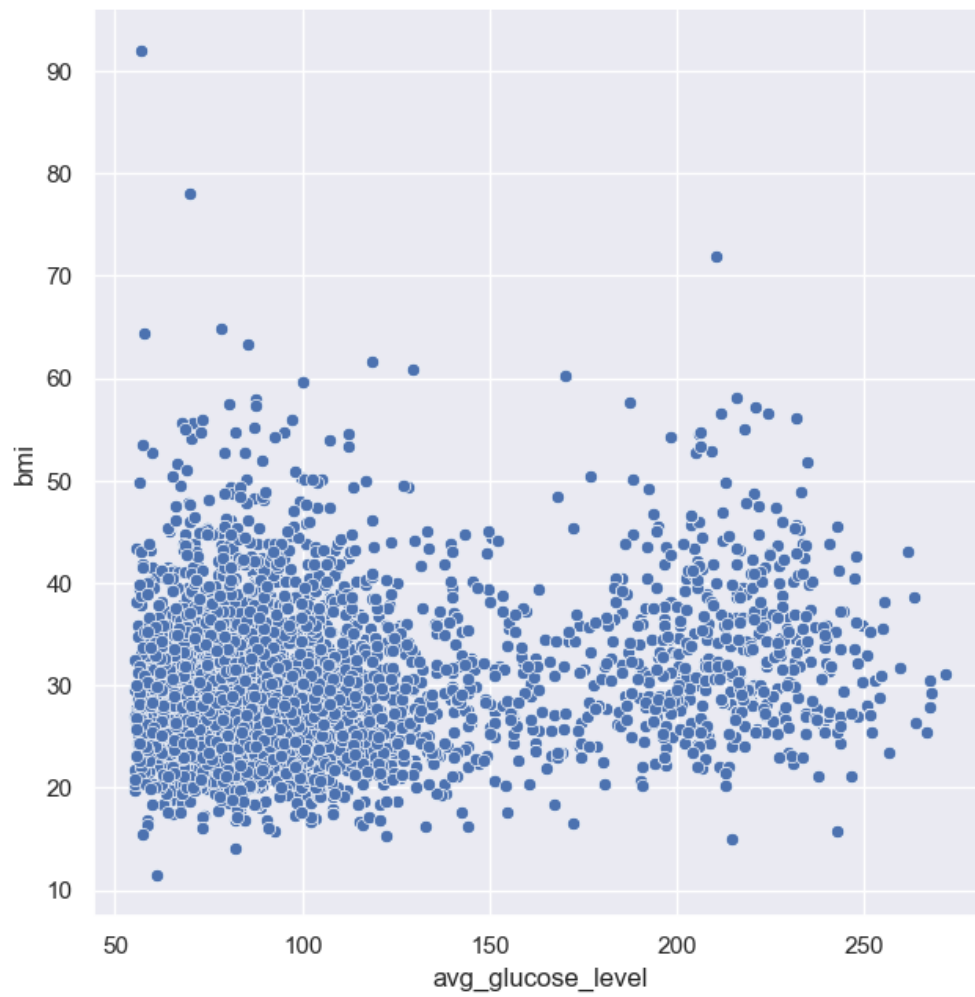


6. Conclusiones

El manejo de datos faltantes siempre presenta un gran reto dentro del área de inteligencia artificial. Los métodos analizados e implementados en esta práctica han demostrado ser útiles para llenar estos espacios vacíos y no cambiar la distribución original de los datos.

Hablando específicamente de cada uno de los atributos trabajados tenemos:

- **sex.** Dado a que solamente había 3 datos faltantes, fue bastante claro el implementar un método como la imputación por moda. Habría que investigar si dicho método es recomendable solamente ante ciertos rangos de datos faltantes.
- **hypertension.** Considerando que el índice de correlación máximo obtenido coincide con la variable objetivo, fue fácil imaginar que la imputación por media/moda de clases sería ideal, sin embargo al prestar atención al método podemos observar que la distribución de los posibles valores de este atributo es aproximadamente de 3:1, por lo que el método de moda de clase relleno todos los faltantes con la clase 0, mismo resultado se hubiera obtenido con una imputación por moda y hubiera disminuido la complejidad computacional.
- **workt_type y Residence_type.** El método de imputación aleatoria logró mantener de forma adecuada la distribución de estas variables, comportamiento que se atribuye a la claras diferencias en la frecuencia de cada una de las observaciones de cada atributo, es muy probable que utilizar este método en atributos con mayor resolución no será eficaz.
- **bmi.** Considerar que el mayor índice de correlación fue con el atributo *avg_glucose_level* fue de ayuda para lograr estimar una línea de tendencia entre ambos atributos, sin embargo, dicho valor de correlación apenas fue de 0,24, lo cuál es un claro indicador de que no existe una relación lineal entre ambos atributos (afirmación que se confirma al generar un gráfico de dispersión entre ambos atributos (Figura ??)). Existe la posibilidad de que utilizar algún tipo de regresión no lineal pueda mejorar el rendimiento de este método. A pesar de sus limitantes, es destacable el hecho de que se logró conservar la forma de la distribución de los datos originales.





Referencias bibliográficas

- [1] M. A. Aceves Fernández, *Inteligencia artificial para programadores con prisa*. Universo de Letras, 2021, ISBN: 9788418854613.



A. Código documentado

El código completo y funcional se puede encontrar anexo en el archivo *zip* compartido en conjunto con este reporte, así como en las secciones anexas.

A.1. DeleteData.py

Script implementado para generar datos faltantes.

```
1 #!/usr/bin/env python3
2 """Script to delete data from data set.
3 """
4
5 import numpy as np
6 import pandas as pd
7
8 from typing import Union
9
10 __author__ = "Enrique Mena Camilo"
11 __email__ = "enriquemece97@gmail.com"
12
13
14 def delete_data(dataset: pd.DataFrame, percentage: float = 0.15, to_delete
15 : Union[None, list] = None):
16     """Performs random data deletion.
17
18     :param dataset: Data set to perform data deletion.
19     :param percentage: Percentage of deletion
20     :param to_delete: Columns to delete.
21     :return pd.DataFrame: Data set with deleted data.
22     """
23     if percentage <= 0 or percentage >= 1:
24         raise ValueError("'percentage' must be more than 0 and less than 1
25 ")
26
27     dataset = dataset.copy()
28     total_samples = dataset.shape[0]
29     to_delete_samples = int(total_samples*percentage)
30
31     if to_delete is None:
32         to_delete = dataset.columns.to_list()
33
34     for feature in to_delete:
35         delete_samples = np.random.randint(0, total_samples, size=
36 to_delete_samples)
37         dataset.loc[delete_samples, feature] = np.nan
```



```
36     return dataset
37
38
39 if __name__ == "__main__":
40     data = pd.read_csv("../data/stroke.csv")
41     features_to_delete = data.drop(columns=["sex", "age", "stroke", "
42     heart_disease", "avg_glucose_level",
43                                     "ever_married", "
44     smoking_status"]).columns.to_list()
45     data_deleted = delete_data(data, to_delete=features_to_delete)
46     data_deleted.to_csv("../data/stroke_nan.csv", index=False)
```

A.2. Practica2_EnriqueMenaCamilo.py

Script implementado para el desarrollo de la práctica.

```
1  #!/usr/bin/env python3
2  """Machine learning: Practice 2.
3  Implement data imputation methods for categorical and continuous
4  attributes.
5  By: Enrique Mena Camilo.
6  """
7
8  import numpy as np
9  import pandas as pd
10 import seaborn as sns
11 import matplotlib.pyplot as plt
12
13 from tqdm import tqdm
14 from sklearn.preprocessing import MinMaxScaler
15
16 sns.set_theme(style="darkgrid")
17
18 __author__ = "Enrique Mena Camilo"
19 __email__ = "enriquemece97@gmail.com"
20
21 def predict_linear_regression(x_data: np.ndarray, w_data: np.ndarray) ->
22     np.ndarray:
23     """Gets predictions using the linear regression method.
24
25     :param np.ndarray x_data: Data with which the prediction will be made.
26     :param np.ndarray w_data: Weights obtained from linear regression
27     training.
28     :return np.ndarray: Results obtained from the prediction.
29     """
30     x_data = np.append(np.ones((len(x_data), 1)), x_data, axis=1)
31     return np.dot(x_data, w_data)
```



```
30
31
32 def train_linear_regression(x_data: np.ndarray, y_data: np.ndarray, lr:
    float = 0.1, it: int = 1000) -> np.ndarray:
33     """Fit linear regression model.
34
35     :param np.ndarray x_data: Training data.
36     :param np.ndarray y_data: Target values.
37     :param float lr: Learning rate. Defaults to 0.1.
38     :param int it: Iterations. Defaults to 1000.
39     :return np.ndarray: Weights from linear regression.
40     """
41     def cost(x_: np.ndarray, y_: np.ndarray, w_: np.ndarray) -> float:
42         """MSE as cost function.
43
44         :param np.ndarray x_: Attributes values.
45         :param np.ndarray y_: Target values
46         :param np.ndarray w_: Weights values.
47         :return float:
48         """
49         return 1/(2 * len(y_)) * np.dot((np.dot(x_, w_) - y_).T, (np.dot(
            x_, w_) - y_))
50
51     x_data, y_data = x_data.copy(), y_data.copy()
52
53     x_data = np.append(np.ones((len(x_data), 1)), x_data, axis=1)
54     x_len = len(x_data)
55     theta = np.ones((x_data.shape[1], 1))
56
57     for _ in tqdm(range(it)):
58         theta = theta - lr * (1/x_len) * np.dot(x_data.T, np.dot(x_data,
            theta) - y_data)
59     print(f"Final cost: {cost(x_data, y_data, theta)}")
60
61     return theta
62
63
64 def target_mean_imputation(dataset: pd.DataFrame, feature: str, target:
    str, feature_type: str = "quant") -> pd.DataFrame:
65     """Performs data imputation using the target mean method.
66
67     :param pd.DataFrame dataset: Full dataset.
68     :param str feature: Feature name to impute.
69     :param str target: Target name.
70     :param str feature_type: One of ['quant', 'qual'].
71     :return pd.DataFrame: Feature with imputed data.
72     """
73     to_fill = dataset[feature].copy()
```




```
74
75     for class_ in dataset[target].unique():
76         if feature_type == "quali":
77             estimated = to_fill.loc[dataset[target] == class_].mode()[0]
78         elif feature_type == "quanti":
79             estimated = to_fill.loc[dataset[target] == class_].mean()
80         else:
81             raise ValueError("'feature_type' must be one of ['quanti', 'quali']")
82
83         to_fill.loc[(dataset[target] == class_) & (to_fill.isna())] =
estimated
84
85     return to_fill
86
87
88 def random_imputation(dataset: pd.DataFrame, feature: str, feature_type:
str = "quanti") -> pd.DataFrame:
89     """Performs data imputation using the random method.
90
91     :param pd.DataFrame dataset: Full dataset.
92     :param str feature: Feature name to impute.
93     :param str feature_type: One of ['quanti', 'quali'].
94     :return pd.DataFrame: Feature with imputed data.
95     """
96     to_fill = dataset[feature].copy()
97     nan_count = to_fill.isna().sum()
98
99     if feature_type == "quanti":
100         estimated = np.random.uniform(low=to_fill.min(), high=to_fill.max
()+1, size=nan_count)
101     elif feature_type == "quali":
102         estimated = np.random.randint(low=to_fill.min(), high=to_fill.max
()+1, size=nan_count)
103     else:
104         raise ValueError("'feature_type' must be one of ['quanti', 'quali
']")
105
106     estimated = pd.Series(estimated, index=to_fill.loc[to_fill.isna()].
index.to_list())
107     return to_fill.fillna(estimated)
108
109
110 class DatasetNanHandler:
111     def __init__(self, dataset: pd.DataFrame, target: str):
112         self.dataset = dataset.copy()
113         self.columns = dataset.columns.to_list()
114         self.target = target
```



```
115     self.features = dataset.drop(columns=[target]).columns.to_list()
116
117     self.nan_count = None
118
119     def get_nan_count(self, rebuild: bool = False) -> pd.DataFrame:
120         """Build information about missing values in the data set.
121
122         :param bool rebuild: Flag to indicate whether to use previous data
123         or recalculate. Defaults to False.
124         :return pd.DataFrame: DataFrame with information about nan values
125         around full dataset.
126         """
127         if self.nan_count is None or rebuild:
128             nan_count = list()
129             for column in self.columns:
130                 nan_count.append({
131                     "column": column,
132                     "total_count": self.dataset[column].shape[0],
133                     "non_nan_count": self.dataset[column].count(),
134                     "nan_count": self.dataset[column].isna().sum(),
135                     "nan_percentage": round(self.dataset[column].isna().
136 sum()/self.dataset[column].count(), 4)
137                 })
138
139             self.nan_count = pd.DataFrame(nan_count).set_index("column")
140         return self.nan_count
141
142     def get_histograms(self, shape: tuple, suffix: str = None):
143         """Build column's histograms in subplots matrix.
144
145         :param tuple shape: Subplots shape.
146         :param str suffix: Suffix in the png file. Defaults to None.
147         """
148         fig, axes = plt.subplots(shape[0], shape[1])
149         fig.suptitle(f"Dataset's columns histograms {suffix}")
150         fig.set_size_inches(10, 10)
151         suffix = f"_{suffix}" if suffix else ""
152
153         for ax, column in zip(axes.flatten(), self.columns):
154             sns.histplot(data=self.dataset, x=column, ax=ax)
155             ax.grid(True)
156             plt.tight_layout()
157             plt.savefig(f"./figures/dataset_histogram{suffix}.png", bbox_inches
158 = 'tight')
159             plt.show()
160
161     def get_heatmap(self):
162         """Build dataset correlation heatmap.
```



```
159     """
160     fig, ax = plt.subplots(1)
161     fig.suptitle(f"Dataset's heatmap")
162     fig.set_size_inches(10, 10)
163     sns.heatmap(data=self.dataset.corr(), annot=True, cmap='BrBG',
164 vmin=-1, vmax=1)
165     plt.savefig(f"./figures/dataset_heatmap.png", bbox_inches='tight')
166     plt.show()
167
168     def fill_nan(self, method: str, feature: str, **kwargs):
169         """Perform data imputation according to the choiced method.
170
171         :param str method: Once of ['mean', 'mode', 'target_mean', 'random
172 ', 'regression']
173         :param str feature: Feature name to impute.
174         :param kwargs:
175         :return:
176         """
177         if method == "mean":
178             estimated = self.dataset[feature].mean()
179             self.dataset[feature].fillna(estimated, inplace=True)
180             return
181
182         if method == "mode":
183             estimated = self.dataset[feature].mode()[0]
184             self.dataset[feature].fillna(estimated, inplace=True)
185             return
186
187         if method == "target_mean":
188             estimated = target_mean_imputation(self.dataset, feature, self
189 .target, kwargs.get("feature_type", "quali"))
190             self.dataset[feature] = estimated
191             return
192
193         if method == "random":
194             estimated = random_imputation(self.dataset, feature, kwargs.
195 get("feature_type", "quali"))
196             self.dataset[feature] = estimated
197             return
198
199         if method == "regression":
200             if "predictors" in kwargs:
201                 predictors = kwargs["predictors"]
202             else:
203                 raise TypeError("Missing 1 required keyword-only argument:
204 'predictors'")
205
206             # Prepare predictors and target
```



```
202         x = self.dataset[self.dataset[feature].notnull()][predictors].
to_numpy()
203         y = self.dataset[self.dataset[feature].notnull()][feature].
to_numpy()
204         y = y.reshape(y.shape + (1,))
205
206         # Normalize predictors
207         scaler = MinMaxScaler()
208         scaler.fit(x)
209         x = scaler.transform(x)
210
211         # Train regression
212         weights = train_linear_regression(x, y)
213
214         # Fill NaN values
215         x_nan = self.dataset[self.dataset[feature].isnull()][
predictors]
216         to_fill_index = x_nan.index.to_list()
217         x_nan = scaler.transform(x_nan.to_numpy())
218         estimated = predict_linear_regression(x_nan, weights)[: , 0]
219         self.dataset.loc[to_fill_index, feature] = estimated
220         return
221
222         raise ValueError("Invalid 'method'")
223
224
225 if __name__ == "__main__":
226     data = pd.read_csv("../data/stroke_nan.csv")
227     data_handler = DatasetNanHandler(data, target="stroke")
228
229     print("====NaN values count (original)====")
230     print(data_handler.get_nan_count())
231
232     print("====Histograms (original)====")
233     data_handler.get_histograms((4, 3), suffix="original")
234
235     print("====Heatmap====")
236     data_handler.get_heatmap()
237
238     print("====Mode imputation====")
239     data_handler.fill_nan("mode", "sex")
240
241     print("====Target mean imputation====")
242     data_handler.fill_nan("target_mean", "hypertension", feature_type="
quali")
243
244     print("====Random imputation====")
```



```
245     data_handler.fill_nan("random", "Residence_type", feature_type="quali"  
246 )  
247     data_handler.fill_nan("random", "work_type", feature_type="quali")  
248  
249     print("====Regression imputation====")  
250     data_handler.fill_nan("regression", "bmi", predictors=["  
251     avg_glucose_level"])  
252  
253     # bmi vs glucose dispersion  
254     ax = sns.pairplot(data, y_vars=["bmi"], x_vars=["avg_glucose_level"])  
255     ax.fig.set_size_inches(8, 8)  
256     plt.savefig("./figures/bmi_glucose.png", bbox_inches='tight')  
257     plt.show()  
258  
259     print("====NaN values count (after imputation)====")  
260     print(data_handler.get_nan_count(rebuild=True))  
261  
262     print("====Histograms (imputed)====")  
263     data_handler.get_histograms((4, 3), suffix="imputed")
```