

Taller Preprocesamiento - Minería de Datos - 2025 - I

Facultad de Ingeniería - Departamento de Ingeniería de Sistemas e Industrial

Estudiantes: Francisco José Salamanca Rivera, fsalamancar@unal.edu.co,

Daniel Mauricio Bonilla Muñoz, dabonilla@unal.edu.co,

David Camilo Gómez Medina, dcgomezme@unal.edu.co

Profesora: Elizabeth León Guzmán, eleonguz@unal.edu.co

ACTIVIDADES

1. Dado un conjunto de datos de una dimensión $X = \{-5.0, 23.0, 17.6, 7.23, 1.11\}$, normalizar usando:

- Min-max normalización en el intervalo $[0, 1]$.
- Min-max normalización en el intervalo $[-1, 1]$.
- Desviación estándar.
- Escala decimal en el intervalo $[-1, 1]$.
- Comparar los resultados y discutir ventajas y desventajas.

Solución

```
1 import numpy as np
2 from sklearn.preprocessing import MinMaxScaler
3
4 v = np.array([-5.0, 23.0, 17.6, 7.23, 1.11])
5
6 # a) Normalizacion min-max con el intervalo [0,1]
7 v_n = (v - v.min()) / (v.max() - v.min())
8
9 print('Datos normalizados [0,1]:\n', v_n)
10 print()
11
12 # b) min-max intervalo [-1, 1]
13 v_n_sk = MinMaxScaler(feature_range=(-1, 1))
14 v_n_sk = v_n_sk.fit_transform(v.reshape(-1, 1))
15 v_n_sk = v_n_sk.T
16
17 print('Datos normalizados [-1, 1]:\n', v_n_sk)
18 print()
19
20 # c) Desviacion estandar
21 v_n_std = (v - v.mean()) / v.std()
22
23 print('Datos normalizados con desviacion estandar:\n', v_n_std)
24 print()
25
26 # d) Escala decimal en el intervalo [-1,1]
27 # Valor absoluto maximo
28 m_v = np.max(np.abs(v))
```

```

29
30 # j = numero de digitos
31 j = int(np.ceil(np.log10(m_v)))
32
33 v_n_dec = v / (10 ** j)
34
35 print('Datos normalizados con escala decimal:\n', v_n_dec)
36 print()

```

Listing 1: Solución Punto 1.

Salida del código:

```

Datos normalizados [0,1]: [0. 1. 0.80714286 0.43678571 0.21821429]
Datos normalizados [-1, 1]: [[-1. 1. 0.61428571 -0.12642857 -0.56357143]]
Datos normalizados con desviacion estandar: [-1.33779582 1.37893488
0.85499396 -0.15116666 -0.74496637]
Datos normalizados con escala decimal: [-0.05 0.23 0.176 0.0723 0.0111]

```

Ventajas y Desventajas

1. Min-Max [0,1]

Ventajas

- Conserva la forma de la distribución original (solo escala y traslada).
- Todos los valores quedan en un rango fijo [0,1], útil para redes neuronales con activaciones en ese intervalo.
- Fácil de interpretar: 0 = mínimo, 1 = máximo.

Desventajas

- Muy sensible a *outliers*: un valor extremo “estira” el rango y comprime el resto.

2. Min-Max [-1,1]

Ventajas

- Igual que [0,1], pero en un rango centrado en cero, lo cual favorece algunos algoritmos (p.ej. redes con activación *tanh*, modelos basados en distancias).
- Todos los valores quedan en un rango fijo [0,1], útil para redes neuronales con activaciones en ese intervalo.
- Fácil de interpretar: 0 = mínimo, 1 = máximo.

Desventajas

- Comparte la misma sensibilidad a *outliers* que el min-max [0,1].

3. Z-score (normalización por desviación estándar)

Ventajas

- Centra los datos en media 0 y varianza 1; los valores quedan en unidades de σ , lo que ayuda a comparar la “lejanía” de la media..
- Menos afectado por outliers extremos que el min-max, pues esos solo influyen en σ y μ .
- Los algoritmos basados en supuestos gaussianos (p.ej. regresión lineal, PCA, SVM) suelen mejorar su rendimiento.

Desventajas

- No acota los valores a un rango fijo: pueden surgir valores muy grandes o pequeños si hay *outliers*.
- La interpretación (unidades de desviación estándar) puede ser menos intuitiva.

4. Escala decimal)

Ventajas

- Normalización sencilla de calcular.
- Se garantiza que $\max(|v'|) < 1$, donde v' es la normalización decimal. Esto, sin necesidad de conocer distribución ni estadísticos.

Desventajas

- No tiene en cuenta la dispersión interna de los datos ni la forma de la distribución.
- Rango muy amplio de posibles escalas (p.ej. pasar de 0.07 a 0.001 si aumenta el máximo), poco control sobre cómo quedan los valores intermedios.
- Menos frecuente en algoritmos de ML modernos, que suelen preferir min-max o *z-score*.

¿Cuál elegir?

- Si se necesita un rango fijo y el dataset no tiene *outliers* fuertes, se puede utilizar **Min-Max**.
- Si el modelo se beneficia de datos centrados en 0, se puede hacer uso de **Min-Max [-1,1]**.
- Para algoritmos que asumen distribución gaussiana o requieren características con varianza unitaria, se puede utilizar *z-score*.
- Para casos muy rápidos y sencillos donde solo importa que $|v'| < 1$, se puede utilizar escala decimal.

2. Dado un conjunto de datos de 4 dimensiones con valores perdidos:

<i>I1</i>	<i>I2</i>	<i>I3</i>	<i>I4</i>
0	1	1	2
2	1	?	1
1	?	?	0
?	2	1	?
2	2	1	0

- a) Dado que el dominio para todos los atributos es $[0, 1, 2]$, ¿cuál debe ser el número de ejemplos “artificiales” si los valores perdidos son interpretados como “no importa el valor” y ellos son reemplazados con todos los posibles valores para su dominio?
- b) ¿Cuál otro método utilizaría para reemplazar los valores perdidos?

Solución

Solución a) Para la columna **I1**, el número de ejemplos artificiales es **3**, y son los siguientes:

- Ejemplo 1 de la columna I1: [0, 2, 1, 0, 2]
- Ejemplo 2 de la columna I1: [0, 2, 1, 1, 2]
- Ejemplo 3 de la columna I1: [0, 2, 1, 2, 2]

Para la columna **I2**, el número de ejemplos artificiales es **3**, y son los siguientes:

- Ejemplo 1 de la columna I2: [1, 1, 0, 2, 2]
- Ejemplo 2 de la columna I2: [1, 1, 1, 2, 2]
- Ejemplo 3 de la columna I2: [1, 1, 2, 2, 2]

Para la columna **I3**, el número de ejemplos artificiales es **9**, y son los siguientes:

- Ejemplo 1 de la columna I3: [1, 0, 0, 1, 1]
- Ejemplo 2 de la columna I3: [1, 0, 1, 1, 1]
- Ejemplo 3 de la columna I3: [1, 0, 2, 1, 1]
- Ejemplo 4 de la columna I3: [1, 1, 0, 1, 1]
- Ejemplo 5 de la columna I3: [1, 1, 1, 1, 1]
- Ejemplo 6 de la columna I3: [1, 1, 2, 1, 1]
- Ejemplo 7 de la columna I3: [1, 2, 0, 1, 1]
- Ejemplo 8 de la columna I3: [1, 2, 1, 1, 1]
- Ejemplo 9 de la columna I3: [1, 2, 2, 1, 1]

Para la columna **I4**, el número de ejemplos artificiales es **3**, y son los siguientes:

- Ejemplo 1 de la columna I4: [2, 1, 0, 0, 0]
- Ejemplo 2 de la columna I4: [2, 1, 0, 1, 0]
- Ejemplo 3 de la columna I4: [2, 1, 0, 2, 0]

En total, el número de ejemplos artificiales para todas las columnas es **18**.

Solución b) Otro método para reemplazar valores perdidos es asignarles una **constante global**, como un valor fijo previamente definido. Sin embargo, esta estrategia depende del contexto y de la aplicación que el usuario le dará a los datos. Otra técnica común es reemplazar los valores perdidos con la **media** (si el atributo es numérico) o la **moda** (si es categórico) del atributo correspondiente. Además, existen métodos más avanzados, como utilizar modelos de **predicción** (por ejemplo, regresión o árboles de decisión) que estiman los valores faltantes a partir de la información de otros atributos.

3. El número de hijos de diferentes pacientes es dado por el siguiente vector:

$C = \{3, 1, 0, 2, 7, 3, 6, 4, -2, 0, 0, 10, 15, 6\}$.

a) Encontrar “outliers” usando parámetros estadísticos estándar: media y varianza.

- 1) Si el umbral cambia de ± 3 desviaciones estándar a ± 2 desviaciones estándar, ¿cuál “outlier adicional” se encuentra?

Solución

```

1  import numpy as np
2
3  # Datos
4  C = np.array([3, 1, 0, 2, 7, 3, 6, 4, -2, 0, 0, 10, 15, 6])
5
6  # Cálculo de media y desviacion estandar
7  mean = C.mean()
8  std = C.std(ddof=1)
9  print(f'Media: {mean:.4f}')
10 print(f'Desviacion estandar: {std:.4f}\n')
11
12 # Umbral +- 3 desv estandar
13 lower_3 = mean - 3 * std
14 upper_3 = mean + 3 * std
15 outliers_3 = C[(C < lower_3) | (C > upper_3)]
16 print('Outliers con umbral +- 3 desv est:', outliers_3)
17
18 # Umbral +- 2 desv estandar
19 lower_2 = mean - 2 * std
20 upper_2 = mean + 2 * std
21 outliers_2 = C[(C < lower_2) | (C > upper_2)]
22 print('Outliers con umbral +- 2 desv est:', outliers_2)
23
24 # Outlier adicional al cambiar de +- 3 desv est a +- 2 desv est

```

```

25 additional = np.setdiff1d(outliers_2, outliers_3)
26 print('\n Outlier adicional encontrado al usar +- 2 desv est:', additional)

```

Salida del código:

Media: 3.9286

Desviacion estandar: 4.5820

Outliers con umbral +- 3 desv est: []

Outliers con umbral +- 2 desv est: [15]

Outlier adicional encontrado al usar +- 2 desv est: [15]

4. Dado un conjunto de tres dimensiones,

$X = \{ \{1, 2, 0\}, \{3, 1, 4\}, \{2, 1, 5\}, \{0, 1, 6\}, \{2, 4, 3\}, \{4, 4, 2\}, \{5, 2, 1\}, \{7, 7, 7\}, \{0, 0, 0\}, \{3, 3, 3\} \}$

- a) Describir el procedimiento e interpretar los resultados de detección de *outliers* basado en la media y varianza.

Solución

Siguiendo el planteamiento del ejercicio anterior, se decide utilizar 2 veces la desviación estándar (2σ) para el análisis de los *outliers*. Esto, con el fin de poder conseguir por lo menos, un dato atípico.

```

1  import numpy as np
2
3  # Datos proporcionados
4  X = [
5      [1, 2, 0],
6      [3, 1, 4],
7      [2, 1, 5],
8      [0, 1, 6],
9      [2, 4, 3],
10     [4, 4, 2],
11     [5, 2, 1],
12     [7, 7, 7],
13     [0, 0, 0],
14     [3, 3, 3]
15 ]
16
17 # Se convierte a array de numpy para facilitar calculos
18 data = np.array(X)
19
20 # Calculo de la media y la desv estandar
21 means = np.mean(data, axis = 0)
22 stds = np.std(data, axis = 0)
23
24 # Limite sup e inf desv estandar

```

```

25 low_b = means - 2 * stds
26 upp_b = means + 2 *stds
27
28 # Identificar outliers
29 outliers = []
30 for i, point in enumerate(data):
31     outlier_dims = []
32     for dim in range(0,3):
33         if not (low_b[dim] <= point[dim] <= upp_b[dim]):
34             outlier_dims.append(dim + 1)
35     if outlier_dims:
36         outliers.append({"punto": point.tolist(),
37                         "indice": i,
38                         "dimensiones": outlier_dims,
39                         "motivo": f"Fuera de rango en dimensión(es)
40                             {outlier_dims}"})
41 # Mostrar resultados
42 print("=== Estadísticas por dimensión ===")
43 for dim in range(3):
44     print(f"Dimensión {dim + 1}:")
45     print(f"  Media: {means[dim]:.2f}")
46     print(f"  Desviación: {stds[dim]:.2f}")
47     print(f"  Rango aceptable: [{low_b[dim]:.2f}, {upp_b[dim]:.2f}]\n")
48
49 print("\n=== Resultados de detección ===")
50 if outliers:
51     for outlier in outliers:
52         print(f"Punto {outlier['punto']} (índice {outlier['indice']}):")
53         print(f"  - {outlier['motivo']}")
54 else:
55     print("No se encontraron outliers")

```

Salida del código:

```
=== Estadísticas por dimensión ===
```

```
Dimensión 1:
```

```
Media: 2.70
```

```
Desviación: 2.10
```

```
Rango aceptable: [-1.50, 6.90]
```

```
Dimensión 2:
```

```
Media: 2.50
```

```
Desviación: 1.96
```

```
Rango aceptable: [-1.42, 6.42]
```

```
Dimensión 3:
```

```
Media: 3.10
```

```
Desviación: 2.30
```

Rango aceptable: [-1.50, 7.70]

=== Resultados de detección ===

Punto [7, 7, 7] (índice 7):

- Fuera de rango en la dimensión(es) [1, 2]

Esto, refleja que, los datos atípicos son el 7, para las dimensiones 1 y 2.

5. En Weka cargar el conjunto de datos iris

- Eliminar manualmente valores (15%) en sus atributos, para simular valores perdidos. Luego aplicar varios métodos que están en weka para remplazar esos valores perdidos. Discutir las diferencias entre el valor real y el que valor que lo reemplaza, y las diferencias entre los métodos.
- Normalizar usando varios métodos.
- Discretizar usando varios métodos.

Solución

Solución a)

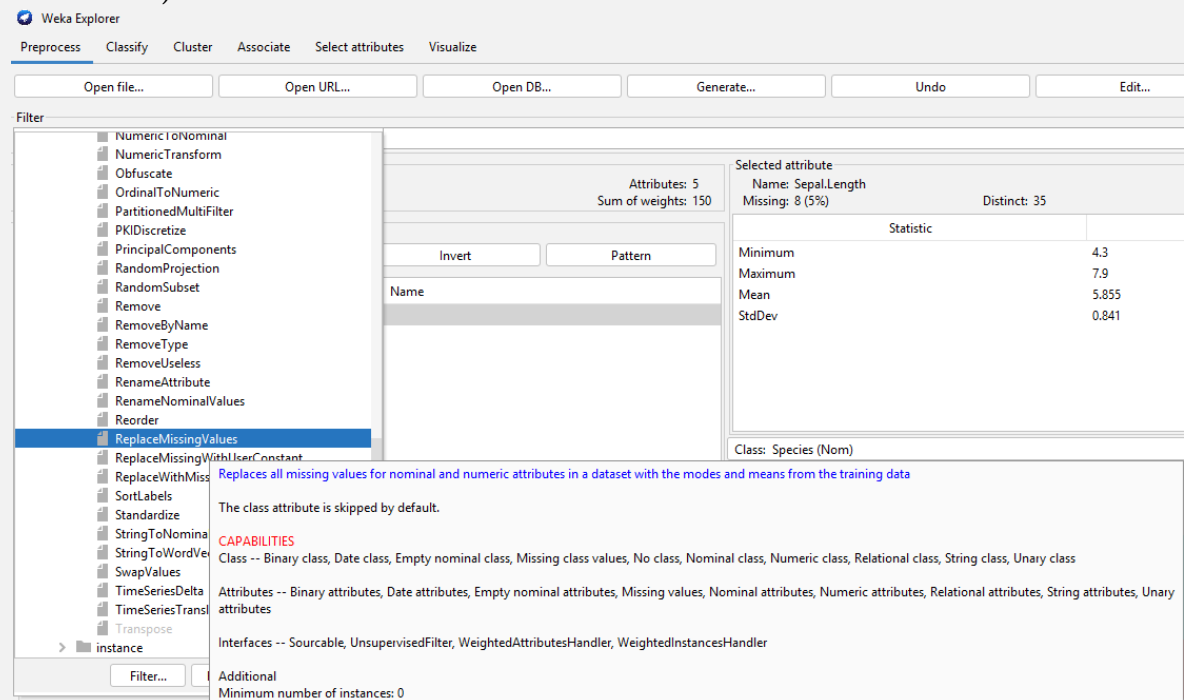


Figura 1: Captura de pantalla usando el método ReplaceMissingValues.

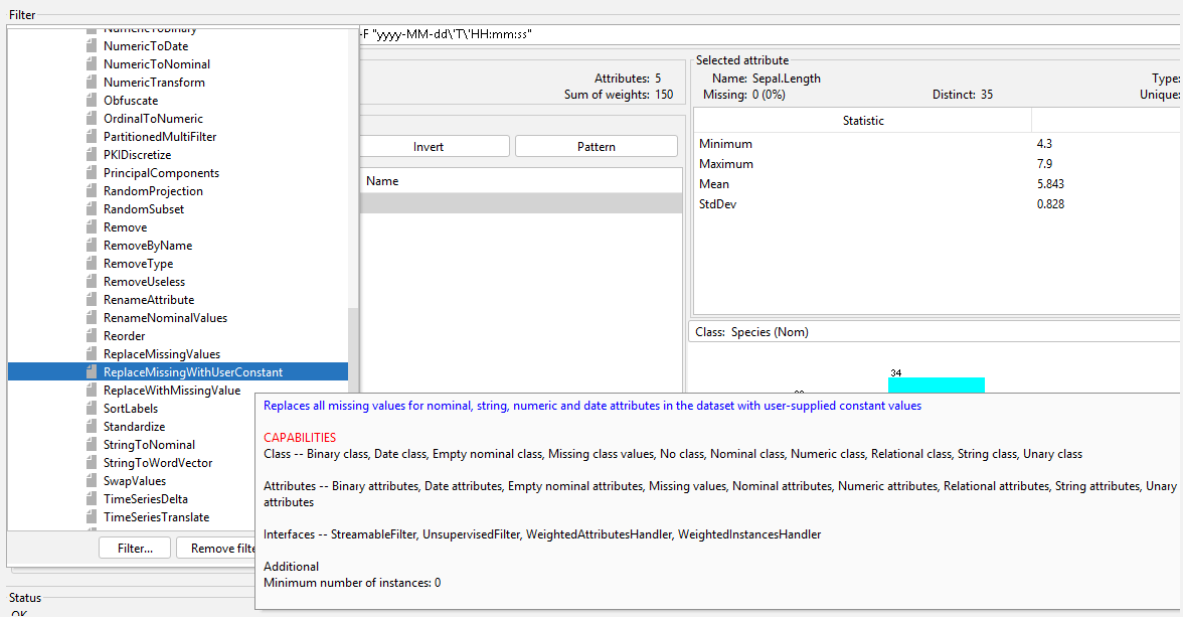


Figura 2: Captura de pantalla usando el método ReplaceMissingWithUserConstant.

Al comparar los valores reales con los valores que se reemplazan mediante los métodos de Weka, se observa que rara vez coinciden de forma exacta.

El método ReplaceMissingValues sustituye los valores perdidos por la media (en atributos numéricos) o la moda (en atributos nominales), lo que permite mantener cierta coherencia con la distribución general del conjunto de datos, aunque puede reducir la variabilidad y alejarse del valor específico que se eliminó.

Por otro lado, el método ReplaceMissingWithUserConstant reemplaza todos los valores perdidos con un valor fijo definido por el usuario, lo que resulta en una solución simple pero menos representativa, ya que no tiene en cuenta la distribución real de los datos. La principal diferencia entre ambos métodos es que ReplaceMissingValues adapta los valores de forma dinámica según el contenido del dataset, mientras que ReplaceMissingWithUserConstant aplica un único valor igual para todos los casos, lo cual puede introducir un sesgo mayor si no se elige cuidadosamente el valor constante.

Solución b)

Método Normalize

El método Normalize se utiliza para escalar los valores numéricos de los atributos dentro de un rango entre 0 y 1. Este proceso es importante porque permite que todos los atributos tengan la misma escala, evitando que aquellos con valores más grandes dominen sobre los más pequeños durante el análisis.

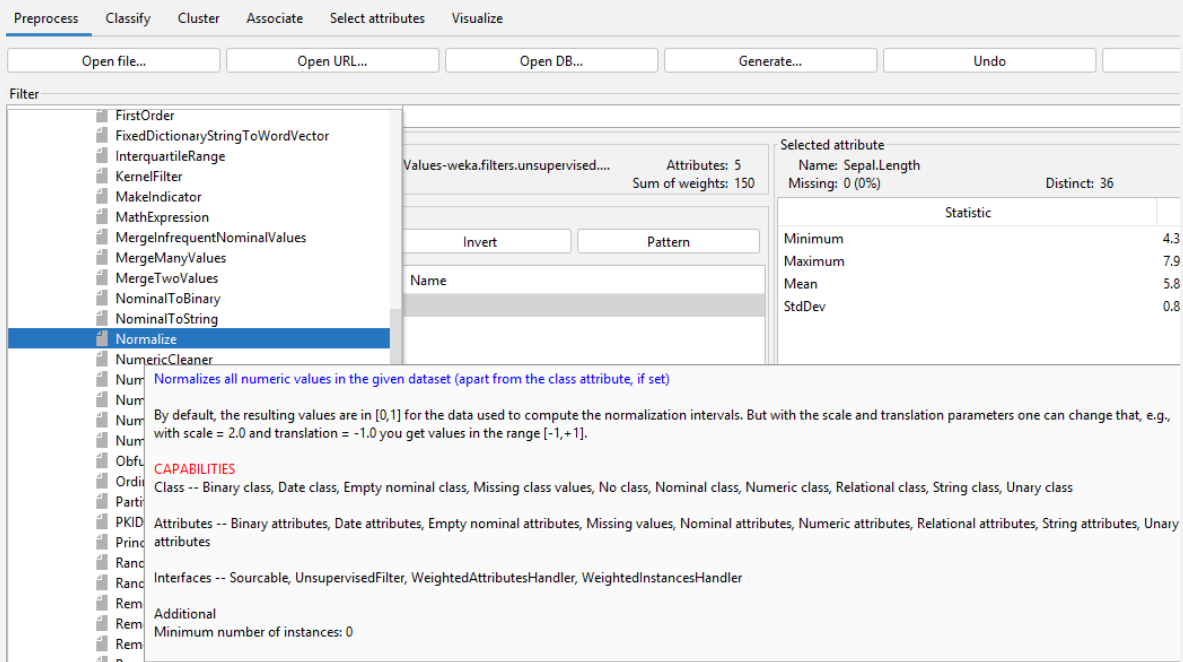


Figura 3: Captura de pantalla usando el método normalize.

Para nuestro conjunto de datos *Iris*, los nuevos valores de la media y la desviación estándar de las variables, después de aplicar los métodos, son los siguientes:

- **Sepal length:** media = 0.429, desviación estándar = 0.23
- **Sepal width:** media = 0.439, desviación estándar = 0.181
- **Petal length:** media = 0.468, desviación estándar = 0.299
- **Petal width:** media = 0.458, desviación estándar = 0.318

Método Standardize

El método Standardize transforma los valores numéricos de los atributos para que tengan una media de 0 y una desviación estándar de 1. Este proceso, se conoce como estandarización o normalización Z-score.

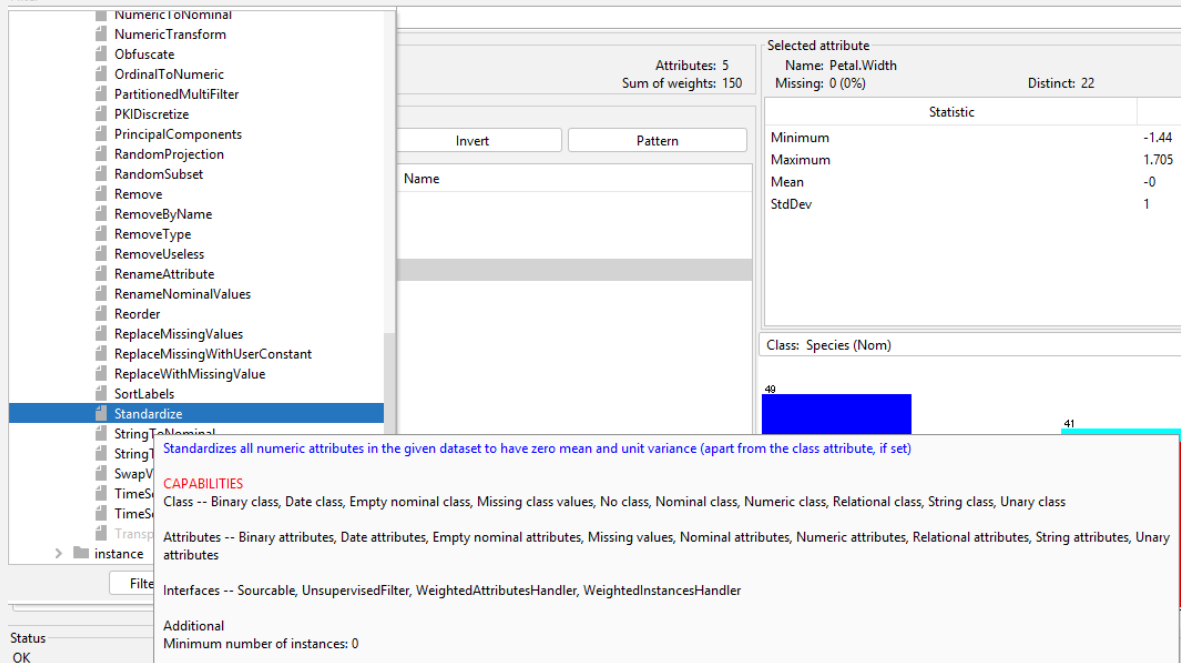


Figura 4: Captura de pantalla usando el método standardize.

Después de aplicar el método los valores de la media y desviación estandar para los cuatro atributos es de 0 y 1 respectivamente.

Solución c)

Método Discretize

El método Discretize en Weka convierte atributos numéricos en valores categóricos o intervalos, dividiendo el rango de datos en un número específico de grupos llamados bins. Cada bin representa un rango de valores y es tratado como una categoría

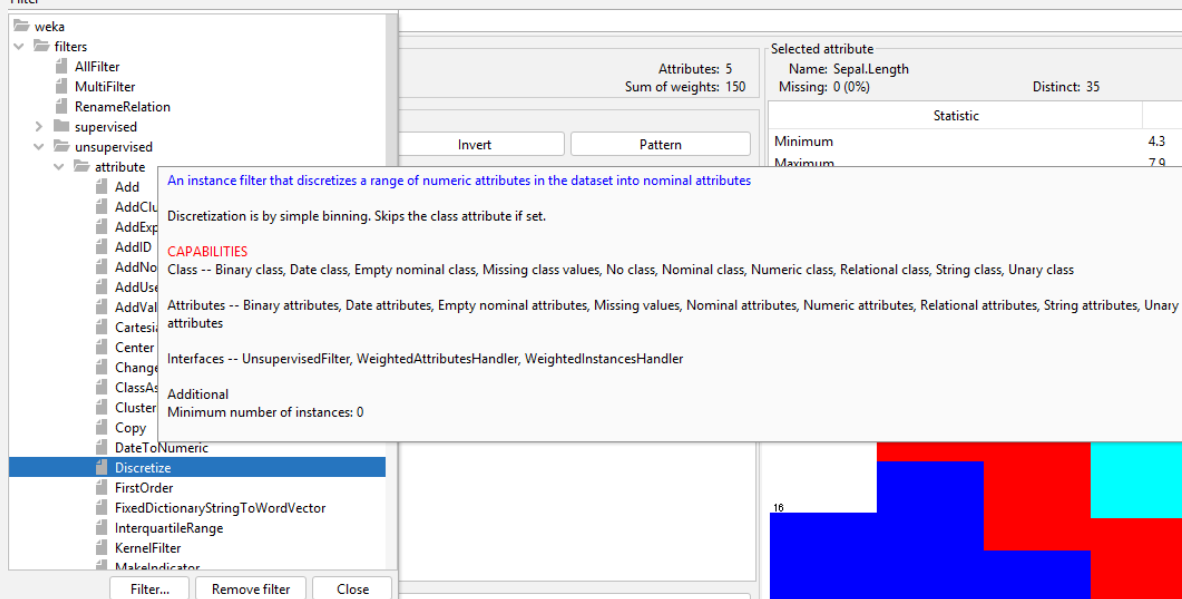


Figura 5: Captura de pantalla usando el método Discretize.

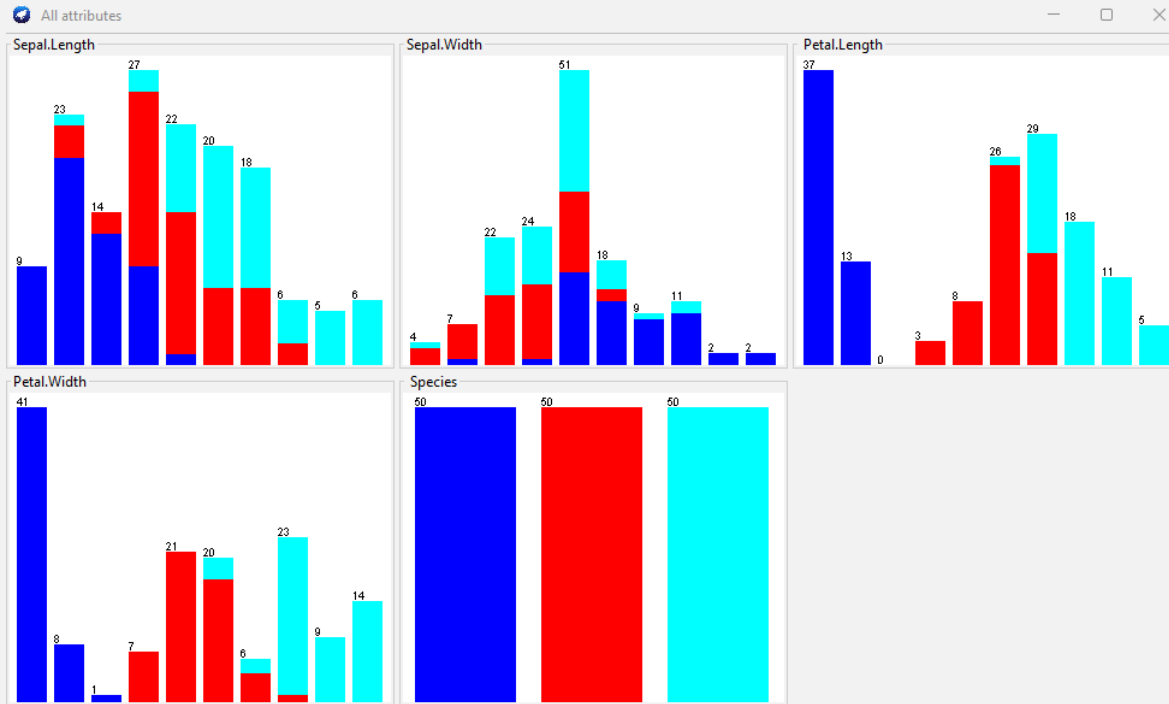


Figura 6: Gráfico de barras con las nuevas categorías después de aplicar el método discretize. Por defecto, el método Discretize utiliza 10 bins para crear las nuevas categorías, en la siguiente figura se puede observar un gráfico de barras para las nuevas categorías en cada atributo.

Método NumericToNominal

El método NumericToNominal en Weka convierte atributos numéricos en categorías o etiquetas. Esto es útil cuando los números no tienen un valor continuo, sino que representan diferentes grupos o clases

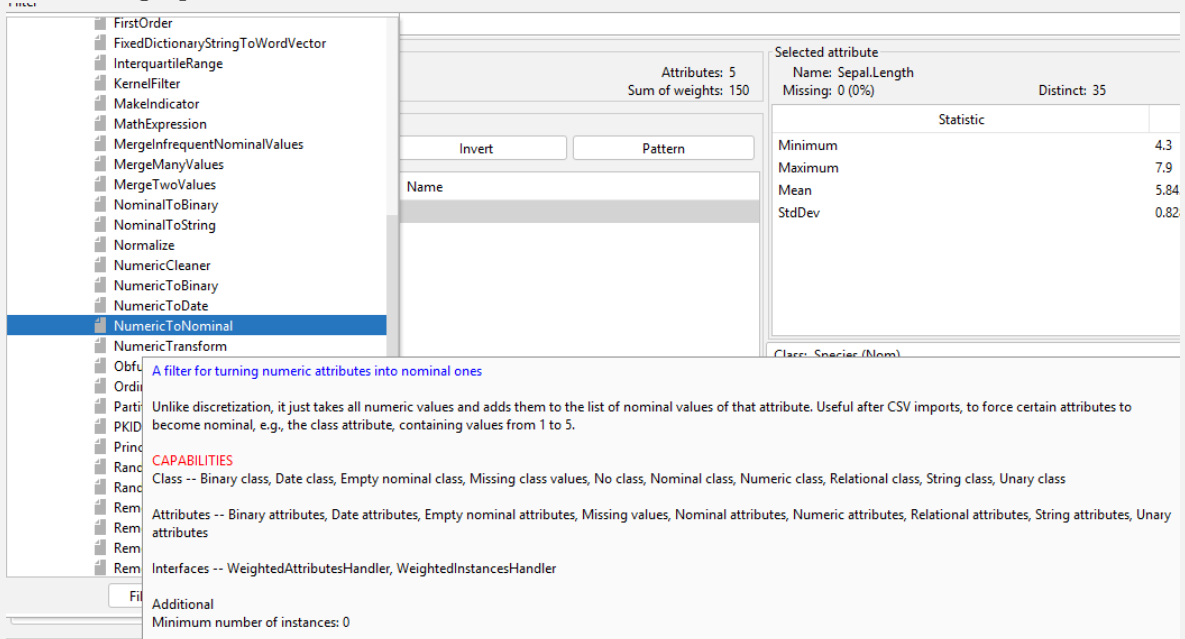


Figura 7: Captura de pantalla usando el método NumericToNominal.

6. Dado el conjunto de datos: Realizar reducción de valores basado en la técnica de BIN con el mejor corte para lo siguiente (mostrar pasos):

$I1$	$I2$	$I3$
1	5.9	3.4
2	2.1	6.2
1	1.6	2.8
2	6.8	5.8
1	3.1	3.1
1	8.3	4.1
2	2.4	5.0

a) Dimensión $I2$ usando la media como representantes de 2 BINS

b) Dimensión $I3$ usando el limite más cercano como representante de 2 BINS

```

1 import numpy as np
2 punto6= np.array([["I1","I2","I3"],
3                   [1,5.9,3.4],
4                   [2,2.1,6.2],
5                   [1,1.6,2.8],
6                   [2,6.8,5.8],
7                   [1,3.1,3.1],
8                   [1,8.3,4.1],
9                   [2,2.4,5.0]])
10
11
12
13 #a. Dimensión I2 usando la media como representantes de 2 BINS (Separar dos bins
    mediante la media)
14
15 punto6a=punto6[1:,1]
16 punto6a=np.sort(punto6a)
17 punto6a = np.array(punto6a, dtype=float)
18 mean=np.mean(punto6a)
19
20 #Separar bins segun las medias
21 bin1=punto6a[punto6a <= mean]
22 bin2=punto6a[punto6a > mean]
23
24 mean_bin1=np.mean(bin1)
25 mean_bin2=np.mean(bin2)
26
27 #reemplaza aquellos menores a la media por el valor de bin1, y el resto con bin2
28 I2=np.where(punto6a <= mean, mean_bin1, mean_bin2)
29
30 print(I2)
31
32 #otra manera usando numpy
33 bins=np.digitize(punto6a,bins=[mean])
34 print(bins)
35
36
37 #b.Dimension I3 usando el limite más cercano como representante de 2 BINS
38

```

```

39 #separar en 2 bins, pero el valor se reemplazara por el limite mas cercano, ya sea
    el superior o el inferior del bin
40
41 punto6b=punto6[1:,1]
42 punto6b=np.sort(punto6b)
43 punto6b=np.array(punto6b, dtype=float)
44 meanb=np.mean(punto6b)
45
46 bin1b=punto6b[punto6b <= meanb]
47 bin2b=punto6b[punto6b > meanb ]
48
49 #bin1
50 lim_sup1=np.max(bin1b)
51 lim_inf1=np.min(bin1b)
52
53 #bin2
54 lim_sup2=np.max(bin2b)
55 lim_inf2=np.min(bin2b)
56
57 #Funcion reemplazo
58 def reemplazo_limite(valor):
59     if valor <= mean:
60         return lim_inf1 if abs(valor - lim_inf1) <= abs(valor - lim_sup1) else lim
        _sup1
61     else:
62         return lim_inf2 if abs(valor - lim_inf2) <= abs(valor - lim_sup2) else lim
        _sup2
63
64 I3 = np.array([reemplazo_limite(v) for v in punto6b])
65 print(I3)

```

Listing 2: Punto 6.

Salida del código:

- `print(I2)`: Imprime los valores de la dimensión I2 una vez reemplazados por la media del bin correspondiente. Los valores menores o iguales a la media general se reemplazan por la media de su grupo (bin1) y los mayores por la media del otro grupo (bin2).
 $I2 = [2.3, 2.3, 2.3, 2.3, 7., 7., 7.]$
- `print(I3)`: Imprime los valores de la dimensión I3 reemplazados por el límite más cercano (inferior o superior) del bin correspondiente, en lugar de por la media.
 $I3 = [1.6, 1.6, 3.1, 3.1, 5.9, 5.9, 8.3]$

7. Dado el conjunto de datos con tres dimensiones de entrada y una dimension representando la clase:

<i>I1</i>	<i>I2</i>	<i>I3</i>	<i>C</i>
2.5	1.6	5.9	0
7.2	4.3	2.1	1
3.4	5.8	1.6	1
5.6	3.6	6.8	0
4.8	7.2	3.1	1
8.1	4.9	8.3	0
6.3	4.8	2.4	1

Hacer el ranking de las dimensiones realizando comparación de medias y varianzas.

```
1 import numpy as np
2 punto7= np.array([["I1","I2","I3", "C"],
3                   [2.5,1.6,5.9,0],
4                   [7.2,4.3,2.1,1],
5                   [3.4,5.8,1.6,1],
6                   [5.6,3.6,6.8,0],
7                   [4.8,7.2,3.1,1],
8                   [8.1,4.9,8.3,0],
9                   [6.3,4.8,2.4,1]])
10
11
12 I1= punto7[1:,0]
13 I1=np.array(I1, dtype=float)
14 I2= punto7[1:,1]
15 I2=np.array(I2, dtype=float)
16 I3= punto7[1:,2]
17 I3=np.array(I3, dtype=float)
18 I4= punto7[1:,3]
19
20
21 print(I1)
22 print(I2)
23 print(I3)
24
25
26 #Varianzas
27 var=[]
28 I1_var= np.var(I1)
29 var.append(I1_var)
30 I2_var= np.var(I2)
31 var.append(I2_var)
32 I3_var= np.var(I3)
33 var.append(I3_var)
34
35 print(var)
36 sorted_var=np.sort(var)
37 print(sorted_var)
38
39 print(var)
40
41
42 print(f"El ranking de las clases segun varianzas es I2, I1, I3")
43
44 #medias
45 med=[]
46 I1_var= np.median(I1)
47 med.append(I1_var)
48 I2_var= np.median(I2)
49 med.append(I2_var)
50 I3_var= np.median(I3)
51 med.append(I3_var)
52
53 print(med)
54 sorted_med=np.sort(med)
55 print(sorted_med)
56
```

```
57 print(f"El ranking de las clases segun medias es I1, I2, I3")
```

Listing 3: Punto 7.

Salida del código:

- `print(var): [3.449795918367347, 2.6314285714285712, 5.9983673469387755]`.
- `print(var): El ranking de las clases segun las varianzas es I2,I1,I3.`
- `print(med): [np.float64(5.6), np.float64(4.8), np.float64(3.1)]`
- `print(b): El ranking de las clases segun las varianzas es I1,I2,I3.`

8. Dado el conjunto de datos X, donde X1 y X2 son dimensiones numericas, X3 y X4 son dimensiones con datos categoricos.

X1	X2	X3	X4
2.7	3.4	1	A
3.1	6.2	2	A
4.5	2.8	1	B
5.3	5.8	2	B
6.6	3.1	1	A
5.0	4.1	2	B

- a) Aplicar el método de selección de características basado en la entropía para reducir una dimensión (mostrar pasos).
- b) Implementar un programa para realizar el “ranking” de dimensiones usando entropía.

Solución

a) Similaridades:

	R1	R2	R3	R4	R5	R6
R1		0.25	0.25	0	0.5	0
R2			0	0.25	0.25	0.25
R3				0.25	0.25	0.25
R4					0	0.5
R5						0
R6						

Calculamos la entropía de todo el conjunto de datos con la siguiente formula:

$$E = - \sum_{i=1}^{N-1} \sum_{j=i+1}^N (S_{ij} \cdot \log(S_{ij}) + (1 - S_{ij}) \cdot \log(1 - S_{ij}))$$

$$E = 5.88$$

Descartamos la siguiente variable X1, y calculamos nuevamente la entropía:

	$R1$	$R2$	$R3$	$R4$	$R5$	$R6$
$R1$		0.33	0.33	0	0.66	0
$R2$			0	0.33	0.33	0.33
$R3$				0.33	0.33	0.33
$R4$					0	0.66
$R5$						0
$R6$						

$$E_{f1} = 6.36$$

$$E - E_{f1} = 5.88 - 6.36 = -0.48$$

Descartamos X2:

	$R1$	$R2$	$R3$	$R4$	$R5$	$R6$
$R1$		0.33	0.33	0	0.66	0
$R2$			0	0.33	0.33	0.33
$R3$				0.33	0.33	0.33
$R4$					0	0.66
$R5$						0
$R6$						

$$E_{f2} = 6.36$$

$$E - E_{f2} = 5.88 - 6.36 = -0.48$$

Descartamos X3:

	$R1$	$R2$	$R3$	$R4$	$R5$	$R6$
$R1$		0.33	0	0	0.33	0
$R2$			0	0	0.33	0
$R3$				0.33	0	0.33
$R4$					0	0.33
$R5$						0
$R6$						

$$E_{f3} = 3.81 \quad E - E_{f3} = 5.88 - 3.81 = 2.06$$

Descartamos X4:

	$R1$	$R2$	$R3$	$R4$	$R5$	$R6$
$R1$		0	0.33	0	0.33	0
$R2$			0	0.33	0	0.33
$R3$				0	0.33	0
$R4$					0	0.33
$R5$						0
$R6$						

$$E_{f4} = 3.81 \quad E - E_{f4} = 5.88 - 3.81 = 2.06$$

Ranking

- a) Descartando x_1 ($\Delta E = -0.48$)
- b) Descartando x_2 ($\Delta E = -0.48$)
- c) Descartando x_4 ($\Delta E = 2.06$)

d) Descartando x_3

($\Delta E = 2.06$)

Se puede reducir la dimensionalidad del conjunto de datos eliminando la variable X1 o X2 porque, al hacerlo, la entropía total del sistema no disminuye significativamente. De hecho, la entropía aumenta ligeramente, lo que indica que estas variables no están aportando información relevante al conjunto de datos.

b)

```
1     import numpy as np
2
3     def hamming_similarity(matrix):
4         # Obtener el número de filas
5         n_rows = matrix.shape[0]
6
7         # Inicializamos una matriz para almacenar las similitudes de Hamming
8         similarity_matrix = np.zeros((n_rows, n_rows))
9
10        # Comparamos todas las filas
11        for i in range(n_rows):
12            for j in range(i, n_rows): # Solo comparar una vez para evitar duplicados
13                # Calculamos la cantidad de coincidencias en las filas
14                hamming_distance = np.sum(matrix[i] != matrix[j])
15                # La similitud de Hamming es el complemento de la distancia de Hamming
16                similarity = 1 - (hamming_distance / len(matrix[i]))
17                similarity_matrix[i, j] = similarity
18                similarity_matrix[j, i] = similarity # Simetría de la matriz
19
20        return similarity_matrix
21
22    def binary_entropy_matrix(S):
23        """Calcula la entropía total de una matriz S con la fórmula dada."""
24        E = 0.0
25        n = S.shape[0]
26        for i in range(n - 1):
27            for j in range(i + 1, n):
28                sij = S[i][j]
29
30                if sij in [0, 1]:
31                    term = 0 # Evita log(0) que es indefinido
32                else:
33                    term = sij * np.log(sij) + (1 - sij) * np.log(1 - sij)
34
35                E += term
36        return -E
37
38    def remove_column_and_recalculate(S, col_index):
39        """Elimina una columna y su fila correspondiente de la matriz de similitud."""
```

```

40 S_reduced = np.delete(S, col_index, axis=1)
41     # elimina columna
42 S_reduced = hamming_similarity(S_reduced)
43 return S_reduced
44
45 def rank_features_by_entropy(S):
46     s= hamming_similarity(S)
47     total_entropy = binary_entropy_matrix(s)
48     print(f"Entropía total del conjunto original: {total_entropy:.4f}\n")
49     rankings = []
50
51
52     for i in range(S.shape[1]):
53         S_reduced = remove_column_and_recalculate(S, i)
54         entropy_reduced = binary_entropy_matrix(S_reduced)
55         delta_entropy = total_entropy - entropy_reduced
56         rankings.append((f"x{i+1}", delta_entropy))
57
58
59     # Ordenamos de menor a mayor contribución (menos importante primero)
60     rankings.sort(key=lambda x: x[1])
61
62     print("Ranking de variables (de menor a mayor aporte de información):\n")
63     for var, delta in rankings:
64         print(f"{var}: E = {delta:.4f}")
65
66
67 X= np.array([
68     [2.7, 3.4, 1, 'A'],
69     [3.1, 6.2, 2, 'A'],
70     [4.5, 3.8, 1, 'B'],
71     [5.3, 5.8, 2, 'B'],
72     [6.6, 3.1, 1, 'A'],
73     [5, 4.1, 2, 'B'],
74 ])
75
76
77
78 rank_features_by_entropy(X)

```

```

Entropía total del conjunto original: 5.8850

Ranking de variables (de menor a mayor aporte de información):

x1: ΔE = -0.4802
x2: ΔE = -0.4802
x3: ΔE = 2.0659
x4: ΔE = 2.0659

```

Figura 8: Salida `rank_features_by_entropy(X)`.

9. Al conjunto de datos Adult del repositorio de Machine Learning:

- a) Convertir todos los atributos numéricos a categóricos utilizando dos estrategias diferentes.
- a) Transformar el conjunto de datos de manera que todos los atributos sean numéricos.

Solución

a) La primer estrategia a implementar, es el uso del paquete `sklearn`. Se hizo uso del módulo `preprocessing` y de la clase `KBinsDiscretizer`. El código se detalla a continuación:

```
1  pip install ucimlrepo
2
3  import numpy as np
4  import pandas as pd
5  from sklearn.preprocessing import LabelEncoder
6  from ucimlrepo import fetch_ucirepo
7
8  # fetch dataset
9  adult = fetch_ucirepo(id=2)
10
11 # data (as pandas dataframes)
12 X = adult.data.features
13 y = adult.data.targets
14
15 # metadata
16 print(adult.metadata)
17
18 # variable information
19 print(adult.variables)
20
21 from sklearn.preprocessing import KBinsDiscretizer
22
23 numeric_cols = ['age', 'fnlwgt', 'education-num', 'capital-gain',
24 'capital-loss', 'hours-per-week']
25
26 # Crea copia del DataFrame original para no modificarlo
27 X_cat_strat1 = X.copy()
28
29 # Aplica discretización
30 discretizer = KBinsDiscretizer(n_bins=3, encode='ordinal', strategy='quantile')
31 X_cat_strat1[numeric_cols] = discretizer.fit_transform(X[numeric_cols])
32
33 # Renombrar categorías (opcional)
34 for col in numeric_cols:
35     X_cat_strat1[col] = X_cat_strat1[col].map({0: 'low', 1: 'medium',
36     2: 'high'})
```

Además, se realizó el mismo método de *Binning* pero sin el apoyo de la clase

KBinsDiscretizer. Esto, con la finalidad de comprender mejor la metodología.

```
1  # columna 'age'
2  bins = [0, 30, 60, 100]
3  labels = ['young', 'adult', 'senior']
4  X_cat_strat2 = X.copy()
5  X_cat_strat2['age'] = pd.cut(X['age'], bins=bins, labels=labels)
6
7  # columna 'education num'
8  bins = [0, 8, 12, X['education-num'].max()]
9  labels = ['elementary', 'middle', 'high']
10 X_cat_strat2['education-num'] = pd.cut(X['education-num'], bins=bins,
11 labels=labels)
12
13 #columna 'capital gain'
14 bins = [-1, 10000, 50000, X['capital-gain'].max()]
15 labels = ['low', 'medium', 'high']
16 X_cat_strat2['capital-gain'] = pd.cut(X['capital-gain'], bins=bins,
17 labels=labels)
18
19 #columna 'capital loss'
20 bins = [-1, 1000, 3000, X['capital-loss'].max()]
21 labels = ['low', 'medium', 'high']
22 X_cat_strat2['capital-loss'] = pd.cut(X['capital-loss'], bins=bins,
23 labels=labels)
24
25 #columna 'hours per week'
26 bins = [0, 30, 40, X['hours-per-week'].max()]
27 labels = ['part-time', 'full-time', 'overtime']
28 X_cat_strat2['hours-per-week'] = pd.cut(X['hours-per-week'], bins=bins,
29 labels=labels)
```

La segunda metodología para abordar el problema, fue la Binaria. A cotinuación, se detalla el código:

```
1  X_cat_strat3 = X.copy()
2
3  # 'age': binario basado en un umbral
4  X_cat_strat3['age'] = np.where(X['age'] < 40, 'young', 'adul')
5
6  # 'capital-gain': binario basado en la media
7  mean_gain = X['capital-gain'].mean()
8  X_cat_strat3['capital-gain'] = np.where(X['capital-gain'] <= mean_gain,
9  'low', 'high')
10
11 # Aplica a todas las columnas numéricas
12 for col in numeric_cols:
```

```

13     if col not in ['age', 'capital-gain']: # Evita columnas ya procesadas
14         threshold = X[col].median() # Usa mediana como umbral
15         X_cat_strat3[col] = np.where(X[col] <= threshold, 'low', 'high')
16
17     X_cat_strat3[numeric_cols] = X_cat_strat3[numeric_cols].astype('category')

```

b) La transformación de todos los atributos a numéricos, se hizo por medio de la metodología *One Hot Encoding*.

```

1  from sklearn.preprocessing import OneHotEncoder
2
3  # Codificar variables categóricas originales
4  encoder = OneHotEncoder(drop='first', sparse_output=False)
5  X_encoded = pd.DataFrame(
6      encoder.fit_transform(X.select_dtypes(include='object')),
7      columns=encoder.get_feature_names_out()
8
9  # Unir columnas numéricas originales y categóricas codificadas
10 X_numeric = pd.concat([X.select_dtypes(exclude='object'), X_encoded], axis=1)

```

Resultados:

```

1  # Para el inciso a)
2  print("Estrategia 1 (cuantiles):")
3  print(X_cat_strat1.dtypes) #Todas las columnas deben ser 'object' o 'category'
4
5  print("\nEstrategia 2 (Binaria):")
6  print(X_cat_strat3[numeric_cols].dtypes)
7
8  # Para el inciso b)
9  print("\nDataset totalmente numérico:")
10 print(X_numeric.dtypes) # Todas las columnas deben ser 'float64' o 'int64'

```

```

Estrategia 1 (cuantiles):
age          object
workclass    object
fnlwgt       object
education    object
education-num object
marital-status object
occupation   object
relationship object
race         object
sex          object
capital-gain object
capital-loss object
hours-per-week object
native-country object
dtype: object

Estrategia 2 (Binaria):
age          category
fnlwgt       category
education-num category
capital-gain category
capital-loss category
hours-per-week category
dtype: object

Dataset totalmente numérico:
age          int64
fnlwgt       int64
education-num int64
capital-gain int64
capital-loss int64
...
native-country_Trinidad&Tobago float64
native-country_United-States float64
native-country_Vietnam float64
native-country_Yugoslavia float64
native-country_nan float64
Length: 103, dtype: object

```

Figura 9: Verificación de los tipos de datos en las estrategias 1 y 2 del inciso a) y del inciso b).

	age	workclass	fnlwgt	education	education-num	marital-status	occupation	relationship	race	sex	capital-gain	capital-loss	hours-per-week	native-country
0	medium	State-gov	low	Bachelors	high	Never-married	Adm-clerical	Not-in-family	White	Male	low	low	medium	United-States
1	high	Self-emp-not-inc	low	Bachelors	high	Married-civ-spouse	Exec-managerial	Husband	White	Male	low	low	low	United-States
2	medium	Private	high	HS-grad	medium	Divorced	Handlers-cleaners	Not-in-family	White	Male	low	low	medium	United-States
3	high	Private	high	11th	low	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	low	low	medium	United-States
4	low	Private	high	Bachelors	high	Married-civ-spouse	Prof-specialty	Wife	Black	Female	low	low	medium	Cuba

Figura 10: Salida primera metodología inciso a) - KBinsDiscretizer.

	age	workclass	fnlwgt	education	education-num	marital-status	occupation	relationship	race	sex	capital-gain	capital-loss	hours-per-week	native-country
0	adult	State-gov	77516	Bachelors	high	Never-married	Adm-clerical	Not-in-family	White	Male	low	low	full-time	United-States
1	adult	Self-emp-not-inc	83311	Bachelors	high	Married-civ-spouse	Exec-managerial	Husband	White	Male	low	low	part-time	United-States
2	adult	Private	215646	HS-grad	middle	Divorced	Handlers-cleaners	Not-in-family	White	Male	low	low	full-time	United-States
3	adult	Private	234721	11th	elementary	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	low	low	full-time	United-States
4	young	Private	338409	Bachelors	high	Married-civ-spouse	Prof-specialty	Wife	Black	Female	low	low	full-time	Cuba

Figura 11: Salida primera metodología inciso a) - Manual.

	age	workclass	fnlwgt	education	education-num	marital-status	occupation	relationship	race	sex	capital-gain	capital-loss	hours-per-week	native-country
0	young	State-gov	low	Bachelors	high	Never-married	Adm-clerical	Not-in-family	White	Male	high	low	low	United-States
1	adul	Self-emp-not-inc	low	Bachelors	high	Married-civ-spouse	Exec-managerial	Husband	White	Male	low	low	low	United-States
2	young	Private	high	HS-grad	low	Divorced	Handlers-cleaners	Not-in-family	White	Male	low	low	low	United-States
3	adul	Private	high	11th	low	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	low	low	low	United-States
4	young	Private	high	Bachelors	high	Married-civ-spouse	Prof-specialty	Wife	Black	Female	low	low	low	Cuba

Figura 12: Salida segunda metodología inciso a) - Binaria.

10. Escoger un conjunto de datos del repositorio de Machine learning, que tenga varias dimensiones y que sean numéricas, y aplicar PCA. Describir el nuevo conjunto de datos.

Hacer el ranking de las dimensiones realizando comparación de medias y varianzas.

```

1 import pandas as pd
2 import numpy as np
3 import random as rd
4 from sklearn.decomposition import PCA
5 from sklearn import preprocessing
6 import matplotlib.pyplot as plt
7 from sklearn.datasets import load_wine
8
9
10 wine = load_wine()
11 wine_df = pd.DataFrame(wine.data, columns=wine.feature_names)
12
13 wine_df.head()
14
15 scaled_wine = preprocessing.scale(wine_df)
16 pca = PCA()
17 pca.fit(scaled_wine)
18 pca_data = pca.transform(scaled_wine)
19
20 #Calcular el porcentaje de variacion de cada PC
21 per_var = np.round(pca.explained_variance_ratio_* 100, decimals=1)
22 labels = ['PC' + str(x) for x in range(1, len(per_var)+1)] #Etiquetas
23
24
25 #Crear grafico
26 plt.bar(x=range(1,len(per_var)+1), height=per_var, tick_label=labels)
27 plt.ylabel('Porcentaje de varianza explicada')
28 plt.xlabel('Componentes principales')
29 plt.title('Porcentaje de varianza explicada por cada PC')
30 plt.show()
31
32 #Se observa que la variacion se condensa principalmente en los 3 primeros PC

```

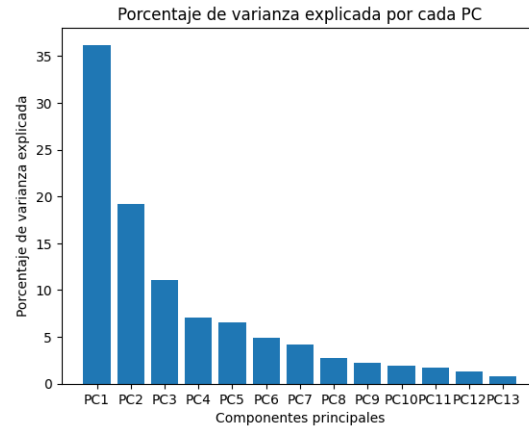



Figura 13: Porcentaje de varianza explicada por cada componente principal

```

1
2
3 import matplotlib.pyplot as plt
4
5 # Crear DataFrame
6 pca_df = pd.DataFrame(pca_data, columns=labels)
7 pca_df['target'] = wine.target
8
9 # Colores por clase
10 colors = ['red', 'green', 'blue']
11 labels_unique = wine.target_names
12
13 plt.figure(figsize=(8, 6))
14
15 for i, label in enumerate(np.unique(wine.target)):
16     subset = pca_df[pca_df['target'] == label]
17     plt.scatter(subset['PC1'], subset['PC2'], color=colors[i], label=labels_unique[i], s=80)
18
19     # Anotaciones opcionales
20     for idx in subset.index:
21         plt.annotate(wine.target[idx], (pca_df.PC1[idx], pca_df.PC2[idx]),
22                     fontsize=8, alpha=0.5)
23
24 plt.title('PCA de Wine Dataset')
25 plt.xlabel('PC1 - {:.2f}%'.format(per_var[0]))
26 plt.ylabel('PC2 - {:.2f}%'.format(per_var[1]))
27 plt.legend(title='Clase')
28 plt.show()

```

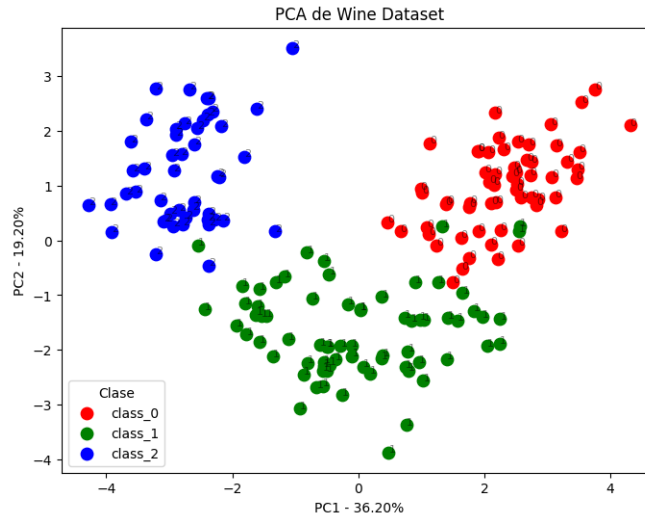


Figura 14: PCA del dataset Wine

```
1 loading_scores = pd.Series(pca.components_[0], index=wine_df.columns)
2 sorted_loading_scores = loading_scores.abs().sort_values(ascending=False)
3
4
5 #Top de las variables que mas contribuyeron a la separacion de las variables
6 top_10_variables = sorted_loading_scores.head(10)
7 print(top_10_variables)
```

```
flavanoids      0.422934
total_phenols   0.394661
od280/od315_of_diluted_wines 0.376167
proanthocyanins 0.313429
nonflavanoid_phenols 0.298533
hue             0.296715
proline         0.286752
malic_acid      0.245188
alkalinity_of_ash 0.239320
alcohol         0.144329
dtype: float64
```

Figura 15: Top variables que mas contribuyeron a la separación de variables

Solución:

- : Se escogió la base de datos **Wine**, proveniente del repositorio de *machine learning* de la Universidad de California en Irvine (UCI). Esta base de datos contiene el análisis químico de vinos cultivados en la región italiana de **Piamonte**, con el objetivo de clasificar los vinos según el tipo de uva (*cultivar*) del cual provienen.

La base consta de vinos pertenecientes a **tres cultivares diferentes**, y cada observación incluye **13 características químicas** medidas en una muestra del vino. Estas características incluyen, entre otras:

- Porcentaje de alcohol
- Ácido málico

- Cenizas
- Alcalinidad de las cenizas
- Magnesio
- Fenoles totales
- Flavonoides
- Intensidad del color

El objetivo principal es utilizar estas variables químicas para construir modelos que permitan **clasificar correctamente** a qué tipo de uva corresponde cada vino. En este punto, se imprimen los valores de la dimensión I2 una vez reemplazados por la media del bin correspondiente. Los valores menores o iguales a la media general se reemplazan por la media de su grupo (bin1) y los mayores por la media del otro grupo (bin2).

- **PCA:** Se observó que los dos primeros componentes principales explican en conjunto un **56 %** de la variabilidad total de los datos, lo cual permite visualizar una separación clara entre las clases de vino. Esta reducción de dimensionalidad facilita la agrupación y análisis visual de los vinos según su cultivar.
- **Importancia de variables:** Dentro de las variables más relevantes identificadas mediante el análisis de componentes principales, las tres con mayor contribución al primer componente fueron: **flavonoids** con un valor de **0.42**, **total phenols** con **0.39** y **od280/od315** con **0.37**. Estas variables son claves para la separación entre clases de vino.