

Actividad 2.1 – Análisis de observaciones influyentes

1.

A)

```
Ejercicio: Calcular la media y la mediana antes de realizar la modificación de incluir unos ingresos de 500.000€  
+ Código + Markdown  
1 media= np.mean(datos)  
2 mediana= np.median(datos) # Es el valor central cuando los datos se ordenan  
3 print(f" Media= {media} y Mediana={mediana}")  
Media= 20364.58936071857 y Mediana=21111.15580445307
```

B)

```
Ejercicio: Aplicar el método de Probabilidad global, para detectar los outliers utilizado en el ejemplo 2_3_Outliers  
1 import scipy.stats as st  
2  
3 # CRITERIO 1: PROBABILIDAD GLOBAL - Explicado en el capítulo 5 - Preprocesamiento (Semana 2), en la pág 10  
4 # Asumiendo que las variables tienen una distribución normal.  
5 # Probabilidad de la muestra de estar dentro de las bandas  
6 p_g=0.95  
7 # probabilidad global  
8 alfa_g=(1-p_g)/2  
9 # probabilidad para un solo dato  
10 alfa= 1-(1-alfa_g)**(1/len(datos)) # Se realiza este ajuste para ser más precisos.  
11 # alfa = alfa_g # Si no realiza la corrección de la línea anterior entonces hubiese detectado más valores que son outliers cuando re  
12 # con esta corrección lo que se pretende es obtener bandas lo suficientemente anchas.  
13 # alfa=1/(2*len(datos))  
14 Z_alfa=st.norm.ppf(1-alfa/2)  
15  
16 # Impresión de resultados  
17 alfa=round(alfa,5)  
18 Z_alfa=round(Z_alfa,5)  
19 print(f" Alfa = {alfa}")  
20 print(f" Z_alfa = {Z_alfa}")  
Alfa =0.00025  
Z_alfa =3.65906
```

```
1 # Utilizaremos el Criterio 1 ó el Criterio 2 para calcular alfa y Z-alfa en función de la naturaleza de los datos que estamos proc  
2 xL= round(np.mean(datos)-Z_alfa* np.std(datos),4)  
3 xU= round(np.mean(datos)+Z_alfa* np.std(datos),4)  
4 print(f" Banda= [ {xL},{xU}]")  
Banda= [ -149826.1041,200131.9466]
```

Python

C)

Python

```
1 # Aplicamos el método de la distancia entre cuartiles, en este caso, al aplicar el método Jackknife no obtenemos
2 # el valor del elemento que es considerado Outlier, sino la posición del mismo.
3 Q1 = np.quantile(phi,0.25)
4 Q3 = np.quantile(phi,0.75)
5 IQR = Q3 - Q1
6 xL=Q1 - 1.5 * IQR
7 xU=Q3 + 1.5 * IQR
8 for i in range(len(datos)):
9     if phi[i] < xL or phi[i]>xU:
10         print(f" El dato {i} es una observación influyente para la media")
```

```
1 pd.DataFrame(phi).describe()
```

	0
count	100.000000
mean	21131.484935
std	76.369084
min	21055.498657
25%	21055.498657
50%	21131.484935
75%	21207.471214
max	21207.471214

En este caso no encuentra el dato influyente.

2.

A)

```
1. ¿Cuánto vale la media, mediana, la desviación estándar muestral, la
varianza muestral y el rango de la variable X?
```

```
1 X = np.array([1, 2, 3, 4, 5, 6, 7, 8, 9, 10])
2
3 # Media
4 media = np.mean(X)
5
6 # Mediana
7 mediana = np.median(X)
8
9 # Desviación estándar muestral
10 desviacion_estandar = np.std(X)
11
12 # Varianza muestral
13 varianza = np.var(X)
14
15 # Rango
16 rango = np.max(X) - np.min(X)
17
18 print("Media:", media)
19 print("Mediana:", mediana)
20 print("Desviación estándar muestral:", desviacion_estandar)
21 print("Varianza muestral:", varianza)
22 print("Rango:", rango)
```

Media: 5.5
Mediana: 5.5
Desviación estándar muestral: 2.8722813232690143
Varianza muestral: 8.25
Rango: 9

B)

2. Utilizar la función describe() de Panda, para obtener la media, desviación estándar, etc...

```
1 X = pd.Series([1, 2, 3, 4, 5, 6, 7, 8, 9, 10])
2
3 print(X.describe());
```

```
count    10.00000
mean      5.50000
std       3.02765
min       1.00000
25%       3.25000
50%       5.50000
75%       7.75000
max      10.00000
dtype: float64
```

C)

3. ¿Por qué el resultado de calcular la desviación estándar con Numpy es diferente a la calculada por describe de Panda? ¿Qué ajuste sería necesario realizar para que los resultados fuesen similares/iguales?

Por defecto, en NumPy, ddof es 0, lo que significa que se está calculando la desviación estándar y la varianza de la población. Sin embargo, en estadísticas, cuando trabajamos con muestras en lugar de poblaciones completas, a menudo queremos ajustar los cálculos para tener en cuenta los grados de libertad perdidos al estimar parámetros poblacionales a partir de la muestra. En este caso, se utiliza ddof=1.

Para que los datos fuesen iguales o parecidos habría que añadir el parametro ddof al metodo de NumPy:

```
1 X = np.array([1, 2, 3, 4, 5, 6, 7, 8, 9, 10])
2
3 # Desviación estándar muestral
4 desviacion_estandar = np.std(X, ddof=1)
5
6 print("Desviación estándar muestral:", desviacion_estandar)
```

```
Desviación estándar muestral: 3.0276583540974917
```

Python

D)

4. Estandarizar la variable (escalamiento) mediante rangos y a continuación calcular la media y la mediana de la variable escalada

$$X_{\text{escalado}} = \frac{X - \min(X)}{\max(X) - \min(X)}$$

```
1 X = np.array([1, 2, 3, 4, 5, 6, 7, 8, 9, 10])
2
3 # Rango
4 rango = np.max(X) - np.min(X)
5
6 # Estandarizar mediante rangos
7 X_escalado = (X - np.min(X)) / rango
8
9 # Calcular la media y la mediana de la variable escalada
10 media_escalada = np.mean(X_escalado)
11 mediana_escalada = np.median(X_escalado)
12
13 print("Media de la variable escalada:", media_escalada)
14 print("Mediana de la variable escalada:", mediana_escalada)
```

```
Media de la variable escalada: 0.5
Mediana de la variable escalada: 0.5
```

E)

5. Repetir el apartado anterior con el escalamiento Z - score

$$Z = \frac{X - \text{media}(X)}{\text{desviación estándar}(X)}$$

```
1 X = pd.Series([1, 2, 3, 4, 5, 6, 7, 8, 9, 10])
2
3 # Estandarizar mediante Z-score
4 X_zscore = (X - X.mean()) / X.std()
5
6 # Calcular la media y la mediana de la variable estandarizada
7 media_zscore = np.mean(X_zscore)
8 mediana_zscore = np.median(X_zscore)
9
10 print("Media de la variable estandarizada (Z-score):", media_zscore)
11 print("Mediana de la variable estandarizada (Z-score):", mediana_zscore)
```

```
Media de la variable estandarizada (Z-score): 4.4408920985006264e-17
Mediana de la variable estandarizada (Z-score): 0.0
```

[GITHUB](#)