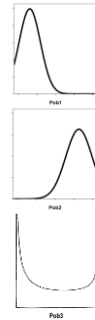


Las distribuciones de probabilidad en el análisis de los datos

1

Ejercicio: Distribuciones

Las gráficas que se presentan a continuación muestran la distribución de probabilidad de los salarios de 3 universidades de un país AAAAA



2

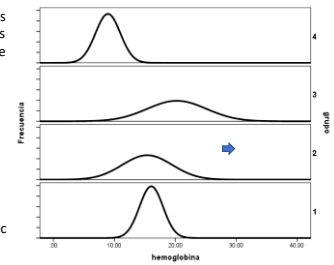
Problema: En un ensayo clínico con cuatro ramas se comparan tres productos que se espera normalicen los niveles de hemoglobina a los hombres anémicos

Valores normales:
14 a 18 g/100ml

3

Problema: En un ensayo clínico con cuatro ramas se comparan tres productos que se espera normalicen los niveles de hemoglobina a los hombres anémicos

Valores normales:
14 a 18 g/100ml



¿Qué opina usted de los resultados?

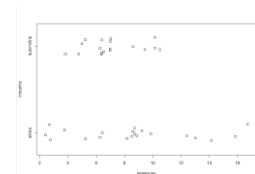
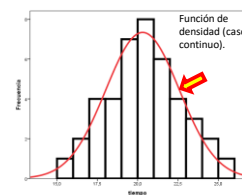
4

Representación gráfica
de la distribución de
las variables continuas

Histograma y
diagramas de
puntos

5

Representación de las distribuciones en el caso continuo: Histogramas y diagramas de puntos



6

Problema

Se sospecha que el si las adolescentes practican un ejercicio físico de alta intensidad esto puede afectar la edad de la menarquia (primera menstruación). Para analizar esta preocupación se realizó un estudio en el que se analizó la edad de la menarquia a adolescentes que practicaron el deporte antes o después de su primera menstruación

Represente la distribución (utilice el histograma y el diagrama de puntos)
Comente sobre el comportamiento de las edades y analice si se observa o no efecto del deporte sobre la edad

Base de datos
menarquia

7

8

Una agencia de empleos hizo evaluaciones para caracterizar las habilidades para la captura de datos de cuatro personas. El número de datos digitados correctamente durante un minuto por cada una de las personas es:

Persona	datos					
A	63	66	68	64	69	72
B	68	67	66	67	68	
C	50	79	75	59	72	20
D	64	68	50	57	59	

Represente gráficamente los valores de cada persona y utilice el gráfico para describir su desempeño

Medidas para representar aspectos de la distribución de probabilidad

9

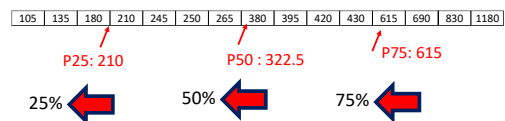
10

Medidas de posición

Posición central

Posición no central

Medidas de posición: No Centrales. Los percentiles



Percentiles, Cuartiles, Deciles

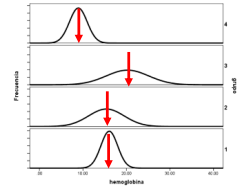
11

12

Los diagramas de cajas

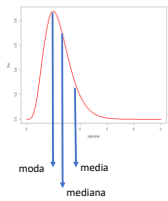
13

Medidas de posición

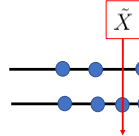


14

Las medidas de posición central


 $X_1 \ X_2 \ \dots \dots \dots \ X_n$

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

 $X_{(1)} \ X_{(2)} \ \dots \dots \dots \ X_{(n)}$


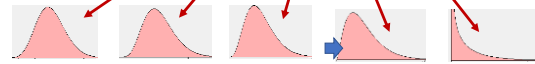
n par

n impar

15

Medidas de posición: ¿Media o mediana?

	23	23	23	23	23
	24	24	24	24	24
	22	22	22	22	22
	25	25	25	25	25
	28	28	28	28	28
	29	29	29	29	29
	31	31	31	31	31
promedios	26.5	27.1	29	320.9	3008.0



Cuando las distribuciones son muy asimétricas la media aritmética no identifica el puntaje "típico" por eso se debe usar otra medida de posición: LA MEDIANA

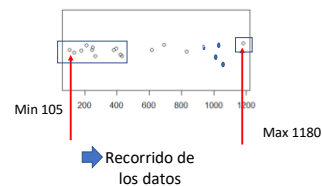
16

Medidas de dispersión (de escala)

17

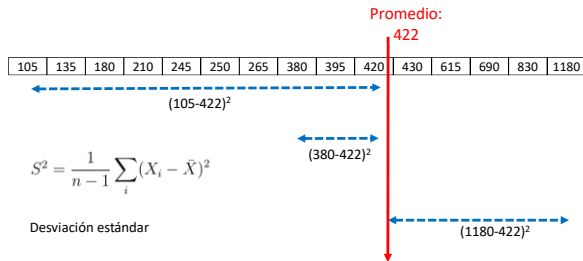
Medidas de dispersión:

Cuantifican homogeneidad / heterogeneidad



18

Medidas de dispersión: Varianza y Desviación estándar



19

¿Por qué no puede ser de una forma más "racional"?

Lo lógico sería definir la dispersión respecto a la media como el promedio de las desviaciones, esto es:

$$\frac{1}{n} \sum_i (X_i - \bar{X})$$

¡Pero esto no funciona!,

Entonces para seguir esa idea debe ser:

$$S^2 = \frac{1}{n-1} \sum_i (X_i - \bar{X})^2$$

¿y qué le parece esta?

$$\frac{1}{n} \sum_i |X_i - \bar{X}|$$

¿Por qué complicarse la vida y no usar mejor?:

$$\sum_i (X_i - \bar{X})^2 \quad \text{o} \quad \sum_i |X_i - \bar{X}|$$

20

La desviación estándar para caracterizar variabilidad

datos	105	110	125	200
	100	100	100	100
	95	90	75	0
media	100	100	100	100
DE	5	10	25	100

21

Uso de la varianza. Descripción de la alimentación de dos personas

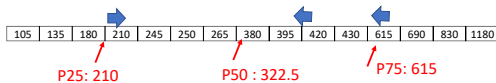
persona 1	persona 2
4500.0	4500.0
1200.0	1200.0
1000.0	1000.0
1100.0	1100.0
1200.0	1200.0
1300.0	1300.0
1400.0	1400.0
1500.0	1500.0
1600.0	1600.0
1700.0	1700.0
1800.0	1800.0
1900.0	1900.0
2000.0	2000.0
2100.0	2100.0
2200.0	2200.0
2300.0	2300.0
2400.0	2400.0
2500.0	2500.0
2600.0	2600.0
2700.0	2700.0
2800.0	2800.0
2900.0	2900.0
3000.0	3000.0
3100.0	3100.0
3200.0	3200.0
3300.0	3300.0
3400.0	3400.0
3500.0	3500.0
3600.0	3600.0
3700.0	3700.0
3800.0	3800.0
3900.0	3900.0
4000.0	4000.0
4100.0	4100.0
4200.0	4200.0
4300.0	4300.0
4400.0	4400.0
4500.0	4500.0

característica	persona 1	persona 2
media	1500.2	811.3
DE	524.1	150.3

Recomienda para una persona
ICBF: 2100Kcal/día

22

Medidas de dispersión: Recorrido intercuartílico



Recorrido intercuartílico:
 $P75 - P25 = 615 - 210$

23

Medidas de dispersión: Mediana de la desviación absoluta (MAD)

$X_1 \ X_2 \ \dots \dots \dots \ X_n$

\tilde{X}

Nueva variable

$$Y_i = |X_i - \tilde{X}|, i = 1, \dots, n$$

$$MAD(X) = \tilde{Y}$$

24

Problema: Tenemos 4 medidas de dispersión (escala). Comente las características de cada una y reflexiones sobre cuándo utilizaría o no utilizaría cada una

$$\text{Recorrido} = X_{(n)} - X_{(1)}$$

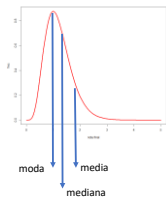
Recorrido intercuartílico: $P75 - P25$

$$\text{Varianza } S^2 = \frac{1}{n-1} \sum_i (X_i - \bar{X})^2$$

Mediana de la desviación absoluta:
mediana ($|X - \bar{X}|$)

Las medidas de forma

25

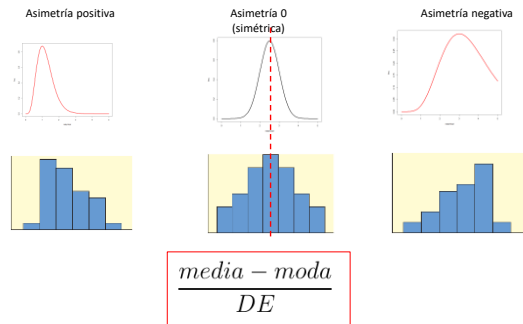


Medidas de forma: El coeficiente de asimetría (1er coeficiente de Pearson)

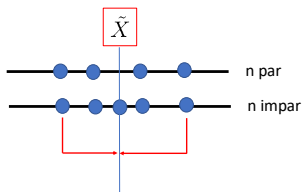
$$\frac{\text{media} - \text{moda}}{DE}$$

27

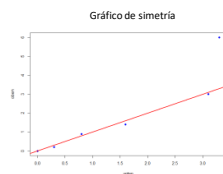
26



28



29



30

