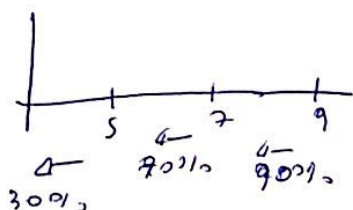


Javier Gómez López



1. Indicar si son verdaderas o falsas las siguientes afirmaciones, justificando la respuesta:

1.1a) Podemos afirmar que $P_{30} = 5$, $P_{70} = 7$, $P_{90} = 9$



observamos que el porcentaje entre el percentil 30 y 70 es del 40% y coincide con el intervalo $(5, 7)$, mientras que entre el percentil 70 y 90 hay un 20% y corresponde al intervalo $(7, 9)$. Por tanto es verdadera.

b)

1.2 Si x e y son estadísticamente independientes, $m_{33} = m_{30} \cdot m_{03}$

Así, $m_{33} = m_{30} \cdot m_{03} = 4 \cdot 5 = 20 \neq 18$. Falsa.

1.3. $y = ax + b$ si pasa por el origen, $b = 0$

$$b = \bar{y} - a\bar{x}; \quad \bar{y} = a\bar{x}; \quad a = \frac{\bar{y}}{\bar{x}} = \frac{m_{01}}{m_{10}}. \quad \underline{\text{Falsa}}$$

$$1.4 \quad r^2 = \frac{\sigma_{xy}^2}{\sigma_x^2 \sigma_y^2}$$

2. Sea $\{x, y\}$ una variable estadística bidimensional con distribución de frecuencias $\{(x_i, y_j); n_{ij}\}$. Obtener la media y la varianza de la distribución marginal de x .

$$\begin{aligned} \bar{x} &= \sum_{i=1}^k j_i \cdot x_i = \frac{1}{n} \cdot \sum_{i=1}^k n_{i.} \cdot x_i = \frac{1}{n} \sum_{j=1}^p n_{ij} \cdot \sum_{i=1}^k x_i = \sum_{j=1}^p \sum_{i=1}^k \frac{n_{ij}}{n} x_i = \\ &= \sum_{j=1}^p \sum_{i=1}^k \frac{n_{ij}}{n} \cdot \frac{n_{ij}}{n_{ij}} x_i = \sum_{j=1}^p \sum_{i=1}^k j_{.j} \cdot j_i^j x_i = \sum_{j=1}^p j_{.j} \sum_{i=1}^k j_i^j x_i = \sum_{j=1}^p j_{.j} \bar{x}_j \end{aligned}$$

$$\begin{aligned}
\sigma_x^2 &= \sum_{i=1}^k d_{i.} (x_i - \bar{x})^2 = \sum_{i=1}^k \frac{n_{i.}}{n} (x_i - \bar{x})^2 = \sum_{j=1}^p \sum_{i=1}^k \frac{n_{ij}}{n} (x_i - \bar{x})^2 = \\
&= \sum_{j=1}^p \sum_{i=1}^k \frac{n_{ij}}{n} \cdot \frac{n_{ij}}{n_{.j}} (x_i - \bar{x})^2 = \sum_{j=1}^p \sum_{i=1}^k d_{.j} d_{i.}^j (x_i - \bar{x})^2 = \\
&= \sum_{j=1}^p d_{.j} \left[\sum_{i=1}^k d_{i.}^j (x_i - \bar{x})^2 \right] = \sum_{j=1}^p d_{.j} \left[\sum_{i=1}^k d_{i.}^j (x_i - \bar{x}_j + \bar{x}_j - \bar{x})^2 \right] = \\
&= \sum_{j=1}^p d_{.j} \left[\sum_{i=1}^k d_{i.}^j (x_i - \bar{x}_j)^2 + \sum_{i=1}^k d_{i.}^j (\bar{x}_j - \bar{x})^2 + 2 \cdot \sum_{i=1}^k d_{i.}^j (x_i - \bar{x}_j) \cdot (\bar{x}_j - \bar{x}) \right] = \\
&= \sum_{j=1}^p d_{.j} \left[\sum_{i=1}^k d_{i.}^j (x_i - \bar{x}_j)^2 + \sum_{i=1}^k d_{i.}^j (\bar{x}_j - \bar{x})^2 \right] = \\
&= \sum_{j=1}^p d_{.j} \left[\sigma_{x_j}^2 + (\bar{x}_j - \bar{x})^2 \sum_{i=1}^k d_{i.}^j \right] = \sum_{j=1}^p d_{.j} \left[\sigma_{x_j}^2 + (\bar{x}_j - \bar{x})^2 \right] = \\
&= \sum_{j=1}^p d_{.j} \sigma_{x_j}^2 + \sum_{j=1}^p d_{.j} (\bar{x}_j - \bar{x})^2
\end{aligned}$$

3. Se realiza un estudio para observar el tiempo que tardan en resolver unos escolares que han seguido un curso de formación por módulos. Se observa el número de módulos que han superado (X) junto con el tiempo en minutos que tardan en resolver el problema (Y)

$X \setminus Y$	[1-9]	[9-21]	[21-39]	$n_{i.}$	$n_{i.} x_i$	$n_{i.} x_i^2$
2	0	1	5	6	12	24
4	0	5	5	10	40	160
5	5	3	0	8	40	200
8	15	1	0	16	128	1024
$n_{.j}$	20	10	10	40	220	1408
$n_{.j} y_j$	100	150	300	550		
$n_{.j} y_j^2$						

a) Sabiendo que $\sigma_y^2 = 104'6875$, ¿qué valor medio es más representativo, el de X o el de Y ?

Primer calculamos la media de X e Y :

$$\bar{X} = \frac{1}{n} \sum_{i=1}^4 n_{i.} x_i = 5'5 \text{ módulos superados}$$

$$\bar{Y} = \frac{1}{n} \sum_{j=1}^3 n_{.j} y_j = \frac{550}{40} = 13'75 \text{ minutos}$$

Ahora calculamos sus varianzas y sus desviaciones típicas:

$$\sigma_y^2 = 104'6875 \text{ minutos}^2$$

$$\sigma_x^2 = m_2 - m_1^2 = \frac{1}{n} \sum_{i=1}^4 n_{i.} x_i^2 - (5'5)^2 = \frac{1408}{40} - 5'5^2 = 4'95 \text{ módulos}^2$$

$$\sigma_x = 2'225 \text{ módulos}$$

$$\sigma_y = 10'232 \text{ minutos}$$

Ahora calculamos sus coeficientes de variación de Pearson

$$C.V(X) = \frac{\sigma_x}{\bar{X}} = \frac{2'225}{5'5} = 0'405$$

$$C.V(Y) = \frac{\sigma_y}{\bar{Y}} = \frac{10'232}{13'75} = 0'744$$

[Por lo tanto el valor medio más representativo es el de X]

Javier Gómez López

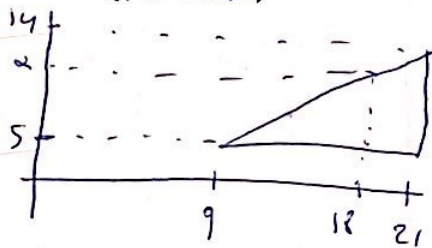
4/8

5) Para los estudiantes que superan menos de 6 módulos, ¿qué porcentaje tarda menos de 18 minutos en resolver el problema? ¿Cuál es el tiempo de respuesta más frecuente?

$X \backslash Y$	[1-9]	(9-21]	(21-39]
2	0	1	5
4	0	5	5
5	5	3	0
$n_{.j}$	5	9	10
$n_{i.}$	5,18	9,12	10,18

→ (Marginal de $Y, x < 6$)

Calculamos cuánto tardan menos de 18 min
(No es la)



observamos que $18 \in (9-21]$

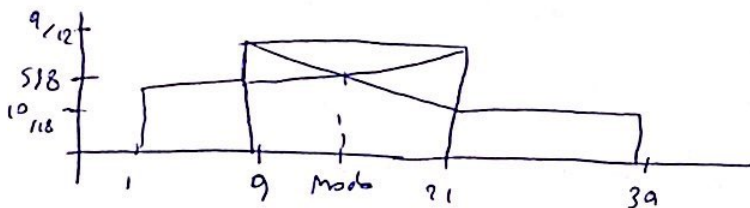
Hacemos semejanza de triángulos

$$\frac{18-9}{21-9} = \frac{2-5}{14-5} \quad ; \quad 2 = 11,75$$

Es decir, 11,75 alumnos tardan menos de 18 min en superar 6 o menos módulos

$\frac{11,75}{24} = 0,4896$. Entonces, podemos afirmar que el 48,96% tardan menos de 18 min en resolver menos de 6 módulos.

Ahora calculamos la moda:
(No es la)



Usando semejanza de triángulos:

$$\frac{9,12 - 5,18}{(9,12 - 5,18) + (9,12 - 10,18)} = \frac{M_0 - 9}{(21 - M_0) + (M_0 - 9)} \quad ; \quad [M_0 = 13,696 \text{ minutos}]$$

Javier Gómez

5,8

c) Calcular el número máximo de módulos del 40% de los estudiantes que más módulos superan

Mostrar el percentil 60, ya que deja por encima al 40%.

$$\frac{n \cdot k}{100} = \frac{40 \cdot 60}{100} = 24 \quad 24 = N_3 \rightarrow x_3 = 5 \text{ módulos}$$

[Concluimos que, el número mínimo de módulos que superan el 40% en más módulos es 5]

d) Sabiendo que $m_{11} = 57'15$, estimar el valor de y cuando $x = 4$ mediante una recta de regresión mínima cuadrática y dar una medida de la bondad de la predicción.

$$y = ax + b \quad a = \frac{\sigma_{xy}}{\sigma_x^2} \quad b = \bar{y} - a\bar{x}$$

$$\sigma_{xy} = m_{11} - m_{10} \cdot m_{01} = m_{11} - \bar{x}\bar{y} = 57'15 - 5'5 \cdot 13'75 = -18'125$$

$$a = \frac{-18'125}{4'125} = -3'662 \quad b = 13'75 - (-3'662) \cdot 5'5 = 33'89$$

$$y = -3'662x + 33'89 \quad \text{Cuando } x = 4, [y = 19'242] \text{ minutos}$$

Ahora calculamos la razón de correlación. $r^2_{y,x} = \frac{\sigma_{xy}^2}{\sigma_x^2 \sigma_y^2} = \frac{(-18'125)^2}{4'125 \cdot 104'6875}$

$$[r^2_{y,x} = 0'634]$$

Javier Gómez

610

e) Ajustar a los datos un modelo de regresión hipérbolica para predecir el tiempo de respuesta conociendo el número de módulos superados y predecir el tiempo de respuesta de un estudiante que superó 4 módulos

$$a = \frac{\sigma_{zy}}{\sigma_z^2} \quad b = \bar{y} - a\bar{z}$$

$$y = a \frac{1}{x} + b; \quad z = \frac{1}{x} \Rightarrow y = az + b$$

Y	9	12	18
z	1/9	1/12	1/18

z \ Y	3	15	30	n _{i.}	z _{i.} · n _{i.}	z _{i.} ² · n _{i.}
1/2	0	1	5	6	3	3/2
1/4	0	5	5	10	2/5	0/625
1/5	5	3	0	8	1/6	0/32
1/8	15	1	0	16	2	0/25
				40	9/1	2/625

los medios de y los tenemos ya.

$$\bar{z} = \frac{1}{n} \sum_{i=1}^4 z_i \cdot n_{i.} = \frac{9/1}{40} = 0'2275 \text{ 1/módulo}$$

$$\sigma_z^2 = m_2 - m_1^2 = \frac{1}{n} \sum_{i=1}^4 z_i^2 \cdot n_{i.} - \bar{z}^2 = \frac{2/625}{40} - 0'2275^2 = 0'0156 \text{ 1/módulo}^2$$

$$\sigma_{zy} = m_{11} - m_{10} m_{01} = \frac{1}{n} \sum_{i=1}^4 \sum_{j=1}^4 z_i \cdot y_j \cdot n_{ij} - \bar{z} \bar{y} = \frac{1/64}{40} - 0'2275 \cdot 13'75 = 0'972 \text{ min/módulo}$$

$$a = \frac{0'972}{0'0156} = 62'13 \quad b = -0'422$$

$$y = 62'13z - 0'422 \Rightarrow y = \frac{62'13}{x} - 0'422$$

$$[Si \ x=4; \ y=15'153 \text{ minutos}]$$

f) Para predecir el tiempo de respuesta, ¿qué modelo de regresión es más adecuado, el lineal o el hipérbolico?

Calculamos la razón de correlación $r^2_{y,x}$ de la hipérbolico

$$f(2) = 30'728$$

$$f(4) = 15'153$$

$$f(5) = 12'038$$

$$f(8) = 7'1358$$

$$\sigma_{e_y}^2 = \frac{1}{n} \sum_{i=1}^4 \sum_{j=1}^3 n_{ij} (f(x_i) - \bar{y})^2 =$$

$$= \frac{1}{40} \cdot (1729'515 + 181'684 + 23'948 + 652'687)$$

$$= 60'610$$

$$r^2_{y,x} = \frac{\sigma_{e_y}^2}{\sigma_y^2} = \frac{60'61}{104'6075} = 0'5789$$

Como el valor de la razón de correlación en la recta es mayor, afirmamos que este ajuste es mejor que el hipérbolico.