

(1)

Solución del examen 7/02/06

① Llamamos X a la variable "nº de días de baja" y construimos la tabla auxiliar:

x_i	n_i	N_i	$x_i \cdot n_i$	$x_i^2 \cdot n_i$
7	1	1	7	49
8	2	3	16	128
12	1	4	12	144
14	1	5	14	196
15	1	6	15	225
19	1	7	19	361
22	1	8	22	484
85	1	9	85	7225
total	9		190	8812

a) media aritmética: $\bar{x} = \frac{190}{9} = \underline{\underline{21.11}}$

Mediana: $n=9$ impar $\Rightarrow M_E = x_{\left(\frac{n+1}{2}\right)} = x_{(5)} = \underline{\underline{14}}$

desviación típica:

$$s_n^2 = \bar{x}^2 - \bar{x}^2 = \frac{8812}{9} - \left(\frac{190}{9}\right)^2 = 533.432$$

$$s_n = \sqrt{533.432} = 23.096 \approx \underline{\underline{23.10}}$$

- b) Para dibujar el diagrama de caja necesitamos calcular Q_1 , Q_3 y las barreras exteriores.

Para calcular Q_1 y Q_3 podemos

- utilizar la columna de frecuencias acumuladas:

$$\frac{n}{4} = 2.75 \Rightarrow Q_1 = 8$$

$$\frac{3n}{4} = 6.75 \approx 7 \Rightarrow Q_3 = 19$$

- O bien, como hay pocos datos, calcularlos directamente:

$$7 \quad \boxed{8 \quad 8} \quad 12 \quad (\textcircled{14}) \quad 15 \quad \boxed{19 \quad 22} \quad 8.5$$

↑ ↑ ↑
 $Q_1 = \frac{8+8}{2} = 8$ $Q_2 = \text{Me} = 14$ $Q_3 = \frac{19+22}{2} = 20.5$

Observación: Q_3 es distinto según el método usado.
Se han elegido por buenos ambos valores.

Para seguir con los cálculos voy a tomar $Q_3 = 20.5$

Rango intercuartílico: $RI = Q_3 - Q_1 = 12.5$

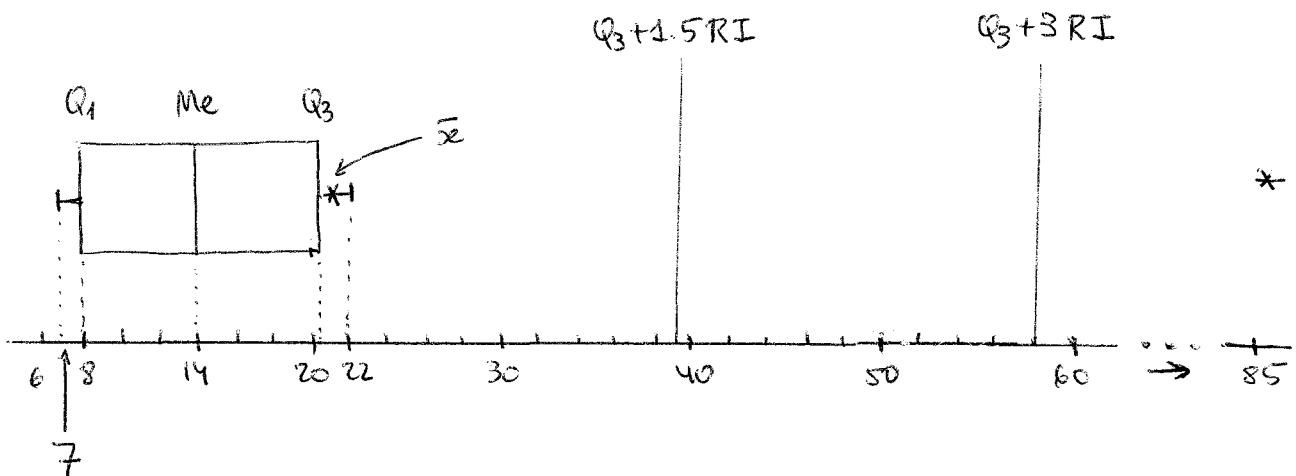
Barreras exteriores:

$$\begin{aligned} & \parallel Q_3 + 1.5 RI = 39.25 \quad (\text{primera B.E. superior}) \\ & \parallel Q_3 + 3 RI = 58 \quad (\text{segunda B.E. superior}) \end{aligned}$$

$$Q_1 - 1.5 RI = -10.75 \leftarrow \text{esta fuera del rango de los datos.}$$

Con esta información dibujamos el diagrama de caja:

(3)



$x_i = 85$ es un valor atípico extremo, puesto que está fuera de la segunda B.E. superior.

Como medida de tendencia central sería preferible usar la mediana en lugar de la media, puesto que la mediana es robusta y, por tanto, no se ve afectada por la presencia de atípicos.

- c) Se observa una asimetría hacia la derecha, puesto que $Me < \bar{x}$. Esto significa que hay más del 50% de los datos por debajo del punto medio de la distribución (\bar{x}). De hecho, mirando el diagrama de caja, vemos que hay más del 75% de los datos por debajo de \bar{x} . Esto es causado en parte por el valor atípico $x_i = 85$, que "tira de la media".

Puesto que la distribución no tiene forma de campana, calcularemos el coeficiente de asimetría de Fisher:

$$As_F = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{S_n^3} = \frac{252100.69/9}{(23.096)^3} = 2.27 > 0$$

↓
asimetría positiva.

- d) Para simetrizar estos datos, podemos utilizar las transformaciones siguientes: \sqrt{x} , $\log x$, $1/x$, ordenadas de menor efecto a mayor.
 Las transformaciones del tipo x^p , $p > 1$, no son adecuadas para este tipo de asimetría.
- e) La relación entre los diagramas A,B,C,D y las transformaciones x^2 , $\log x$, \sqrt{x} , $1/x$ es:

$$A \rightarrow \log x$$

$$B \rightarrow \sqrt{x}$$

$$C \rightarrow x^2$$

$$D \rightarrow 1/x$$

② Llamemos $X = \text{"nº de horas"}$, $Y = \text{"edad"}$, $Z = \text{"sexo"}$.

- a) X es cuantitativa continua, agrupada en intervalos.
 Y es cuantitativa discreta (o continua), también agrupada en intervalos.
 Z es categórica nominal.

Distribuciones marginales de las 3 variables:

Para X	x_{ij}	n_{ij}	$x_{i\text{inic}}$	$x_{ij}^2 n_{ij}$	
[0, 20]	10	33	330	3300	
(20, 30]	25	24	600	15000	
(30, 40]	35	43	1505	52675	
		100	2435	70975	

tabla auxiliar para el apartado f)

Para Y	y_{ij}	n_{ij}	N_{ej}	$y_{ij} n_{ij}$	$y_{ij}^2 \cdot n_{ij}$
[15, 25]	20	31	31	620	12400
(25, 35]	35	44	75	1540	53900
(35, 45]	55	25	100	1375	75625
		100		141925	

Para Z :

Z_K	n_K
Hombre	46
Mujer	54
	100

tabla auxiliar para el apartado b) y para el apartado f)

b) Mediana de edad:

$\frac{n}{2} = 50 \Rightarrow$ intervalo mediano $(25, 45]$.

$$Me = 25 + (45-25) \frac{50-31}{75-31} = \underline{\underline{33.64}}$$

c) $\bar{x}_H = \frac{10 \cdot 13 + 25 \cdot 11 + 35 \cdot 22}{13+11+22} = \frac{1175}{46} = 25.54$

$$\bar{x}_M = \frac{10 \cdot 20 + 25 \cdot 13 + 35 \cdot 21}{20+13+21} = \frac{1260}{54} = 23.33$$

$\bar{x}_H > \bar{x}_M \Rightarrow$ En media, los hombres miran más la TV.

d) El % de hombres que mira la TV más de 20 horas es $\frac{11+22}{46} \cdot 100\% = \underline{\underline{71.74\%}}$

e) Nos piden la moda de las mujeres que miran la TV menos de 30 horas:

$y = \text{"edad"}$	$[0, 30]$	n_i / L_i
$[15, 25]$	7	0.7
$(25, 45]$	24	1.2
$(45, 65]$	2	0.1

← intervalo modal
 $(25, 45]$

$$Mo = 25 + (45-25) \frac{0.1}{0.7+0.1} = \underline{\underline{27.5}}$$

f) Para medir la dispersión utilizaremos el coeficiente de variación, puesto que se trata de variables expresadas en diferentes unidades de medida (una son horas y la otra años). Por tanto, no se pueden comparar sólo utilizando las desviaciones estándares o sólo las varianzas.

$$\text{Edad: } \bar{y} = \frac{3535}{150} = 35.35$$

$$s_y^2 = \bar{y}^2 - \bar{y}^2 = \frac{141925}{100} - (35.35)^2 = 169.6275$$

$$\text{Nº horas: } \bar{x} = \frac{2435}{100} = 24.35$$

$$s_x^2 = \bar{x}^2 - \bar{x}^2 = \frac{70975}{100} - (24.35)^2 = 116.8275$$

$$CV_X = \frac{s_x}{\bar{x}} = \frac{\sqrt{116.8275}}{24.35} = \underline{\underline{0.44}}$$

$$CV_Y = \frac{s_y}{\bar{y}} = \frac{\sqrt{169.6275}}{35.35} = \underline{\underline{0.37}}$$

\Rightarrow X presenta más dispersión que Y.

g) Para ver si existe alguna relación entre la edad y el nº de horas que miran la TV podemos utilizar el coef. de correlación lineal de Pearson, puesto que las otras variables son cuantitativas. El coef. de correlación de Spearman no es calculado porque las variables de nuestro estudio no son puntuaciones dadas por un observador externo.
 Las medidas χ^2 , χ^4 y χ^5 son para tablas de contingencia, no para tablas de doble entrada, como es nuestro caso.

Nº horas (x_i)

Edad (y_j)	10	25	35	$(\sum_i x_i n_{ij}) y_j$
20	2	10	19	18700
35	28	12	4	25200
55	3	2	20	42900
				86800

$$s_{xy} = \overline{xy} - \bar{x}\bar{y} = \frac{86800}{100} - (24.35)(35.35) = 7.2275$$

$$r_{xy} = \frac{s_{xy}}{s_x s_y} = \frac{7.2275}{\sqrt{(169.6275)(116.8275)}} = \underline{\underline{0.05}}$$

r_{xy} es muy cercano a cero, por lo que deducimos que la relación lineal es extremadamente débil. De hecho, podríamos decir que no existe. Aunque puede existir algún otro tipo de relación.

h) La variable "sexo" es categórica nominal, por tanto para estudiar la relación entre esta variable y el "nº de horas" deberemos utilizar medidas para tablas de contingencia, y en concreto, para variables nominales.

Quedan excluidas h1, h3.

Las medidas que se pueden utilizar son: h2, h4, h5.

El procedimiento a seguir es:

1º Realizamos un contraste de independencia basado en el estadístico χ^2 para saber si hay o no independencia. Si las variables son independientes, bemos firmar esto.
Si las variables son dependientes:

2º Medimos el grado de dependencia mediante el coef. de contingencia de Pearson o la V de Cramér. Ambas toman valores entre 0 y 1. Cuanto más cercanas a 1 mayor grado de asociación (dependencia) entre las variables.

- ③ Llamamos Y = "nº de televisores", X = "ingresos".
 $\bar{y} = 2$, $s_y^2 = 0.04$, $\bar{x} = 930$, $s_x^2 = 2410$, $r_{xy} = 0.9$.

a) Encuentra la recta de regresión de Y sobre X , nº de televisores en función de los ingresos mensuales.

$$Y = a + bX,$$

$$b = \frac{s_{xy}}{s_x^2} = \underbrace{\frac{r_{xy} s_x s_y}{s_x^2}}_{r_{xy} = \frac{s_{xy}}{s_x s_y}, s_{xy} = r_{xy} s_x s_y} = r_{xy} \frac{s_y}{s_x} = 0.9 \sqrt{\frac{0.04}{2410}} = \underline{\underline{0.0037}}$$

$$a = \bar{y} - b\bar{x} = 2 - 0.0037(930) = \underline{\underline{-1.41}}$$

b) Para valorar la bondad de ajuste del modelo utilizamos el coef. de determinación R^2 , que en este caso coincide con r_{xy}^2 .

$$R^2 = r_{xy}^2 = 0.9^2 = 0.81 \rightarrow 81\% \text{ de la variabilidad de } Y \text{ es explicada por el modelo.}$$

Queda sin explicar el $100 - 81 = 19\%$ de la variabilidad.

El modelo \Rightarrow bastante fiable.

c) $x_{ei} = 900$

$$\hat{y}_i = -1.41 + (0.0037)900 = 1.92 \approx 2 \rightarrow 2 \text{ televisores.}$$

$x_{ei} = 900$ es muy cercano a $\bar{x} = 930$.

(Recordemos que las mejores predicciones se obtienen para valores cercanos a la media)

d) $b > 0 \Rightarrow$ pendiente positiva \Rightarrow descartamos D

$r_{xy} = 0.9 \Rightarrow$ relación lineal fuerte \Rightarrow descartamos B

$R^2 = 0.81 \Rightarrow 81\% \text{ de ajuste} \Rightarrow$ descartamos C.

El gráfico A \Rightarrow el que corresponde a nuestra nube de puntos.