

# Informes automatizados de estadística univariante y bivalente

III Jornadas de Usuarios de R

Tania Iglesias, Patricia Díaz, Alexandra González

Unidad de Consultoría Estadística  
Servicios Científico-Técnicos  
Universidad de Oviedo

Madrid, 18 de noviembre de 2011

## 1 Introducción

## 2 Análisis estadístico

- Análisis Univariante
  - Variables continuas
  - Variables nominales
- Análisis Bivariante
  - Análisis bivariante nominal nominal
  - Análisis bivariante nominal continua
  - Análisis bivariante continua continua

- 1 Petición por parte del usuario
- 2 Reunión conjunta para establecer objetivos
- 3 Realización del análisis estadístico
- 4 Entrega de informe de resultados



- 1 Petición por parte del usuario
- 2 Reunión conjunta para establecer objetivos
- 3 Realización del análisis estadístico
- 4 Entrega de informe de resultados



- 1 Petición por parte del usuario
- 2 Reunión conjunta para establecer objetivos
- 3 Realización del análisis estadístico
- 4 Entrega de informe de resultados



- 1 Petición por parte del usuario
- 2 Reunión conjunta para establecer objetivos
- 3 Realización del análisis estadístico
- 4 Entrega de informe de resultados



## Pasos previos

- 1 Análisis exploratorio o descriptivo de datos
- 2 Representaciones gráficas
- 3 Estudio de relaciones entre variables

## Necesidad y realización de informes

- Automatización mediante la definición de funciones en R
- Utilización del paquete Sweave: integración con  $\text{\LaTeX}$

## Pasos previos

- 1 Análisis exploratorio o descriptivo de datos
- 2 Representaciones gráficas
- 3 Estudio de relaciones entre variables

## Necesidad y realización de informes

- Automatización mediante la definición de funciones en R
- Utilización del paquete Sweave: integración con  $\text{\LaTeX}$



## Estadística Univariante

- Descripción de variables continuas
- Descripción de variables categóricas o nominales

## Estadística Bivariante

- Relación entre variables nominales
- Relación entre variables continuas
- Relación entre una variable continua y otra nominal

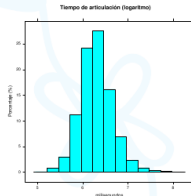
## Variables continuas

```
univariantecontinualatexparagraph <- function
( x
, labelx = deparse(substitute(x))
, namex = deparse(substitute(x))
, printoutliers = TRUE
, printtable = TRUE
, printgraph = TRUE
, printtestdistribution = TRUE
, printblock = TRUE
, h0mu = mean(x, na.rm=TRUE )
, latexwidthblock = 0.485
, dirgraph = "Dibujos"
, latexwidthgraph = 0.485
, significancelevelalpha = 0.05
, ... )
```

# Análisis univariante variable numérica

- Descriptivos
  - Media, mediana, cuartiles
  - Número de casos válidos
  - Número de casos perdidos
  - Mínimo y máximo
- Valores atípicos
- Normalidad: Shapiro Wilk
- Simetría: Test de Mira y Wilcoxon
- Igualdad de medias: Student o Bootstrap

Respecto a la variable *Tiempo de articulación (logaritmo)*, resulta que se dispone de 3561 casos válidos, ya que se produce un 6.71% de casos perdidos en esta magnitud. El valor medio se alcanza en 6.3 milisegundos, con una desviación típica de 0.37, mientras que la mediana disminuye hasta 6.27 milisegundos. El valor mínimo se alcanza en 5.08 y el máximo asciende a 8.1. El 50% de las observaciones centrales se encuentran entre 6.07 y 6.51, estando estos datos agrupados en un intervalo de amplitud 0.44 milisegundos, si bien el 25% de los datos inferiores se encuentran en una amplitud de 0.99 y el 25% de los valores más grandes se ubican en una distancia de 1.6 milisegundos. Se detectan 19 y 65 observaciones extremas en la tramo inferior y superior de la distribución, respectivamente. La presencia de estas observaciones atípicas no influye significativamente en los estadísticos calculados. En la tabla aparece como NA los valores perdidos, y %(NA+) y %(NA-) representan la distribución porcentual incluyendo o no los casos perdidos, respectivamente.



Respecto a la distribución de los datos, se rechaza la hipótesis de normalidad (test de Shapiro-Wilk,  $p$ -valor=0).

#### Valoración

El valor medio obtenido es consistente.

	Frec.	%(NA+)	%(NA-)
[5.08,5.38)	17	0.4	0.5
[5.38,5.68)	108	2.8	3.0
[5.68,5.99)	511	13.4	14.3
[5.99,6.29)	1220	32.1	34.3
[6.29,6.59)	1041	27.4	29.2
[6.59,6.89)	471	12.4	13.2
[6.89,7.2)	135	3.6	3.8
[7.2,7.5)	37	1.0	1.0
[7.5,7.8)	16	0.4	0.4
[7.8,8.1]	5	0.1	0.1
NA's	239	6.3	0.0
Total	3800	100.0	100.0

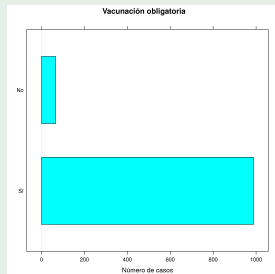
## Variables nominales

```
univariantenominallatexparagraph <- function
  ( x
  , labelx = "etiqueta x"
  , namex = "x"
  , isordered = FALSE
  , nombrefiltro = ""
  , unitsex = "unidades"
  , ndigits = 2
  , printtable = TRUE
  , printgraph = TRUE
  , latexwidthblock = 0.485
  , dirgraph = "Dibujos"
  , latexwidthgraph = 0.485
  , significancelevelalpha = 0.05
  , ... )
```

- Frecuencias
- Porcentajes

## Vacunación obligatoria

	<b>Frec.</b>	<b>%</b>	<b>Acum. %</b>
No	66	6.3	6.3
Sí	988	93.7	100.0
Total	1054	100.0	100.0



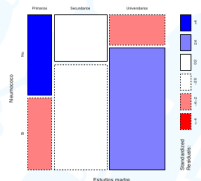
## Tipos de análisis

- Nominal- Nominal
- Nominal- Continua
- Continua- Continua

## Test implementados

- Número suficiente de casos:
  - Test de Pearson
- Número insuficiente de casos:
  - Test de Fisher
  - Test de Barnard (2x2)

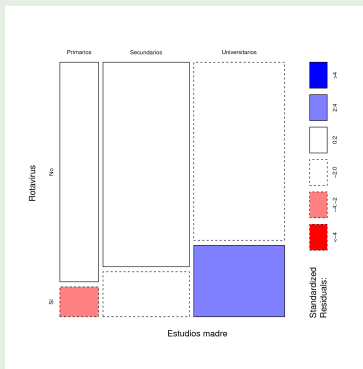
Se realizó el análisis para estudiar la relación entre *Neumococo* y *Estudios madre*, obteniéndose que se rechaza la hipótesis de independencia (test de Pearson, p-valor=0). Se detectan relaciones en las siguientes modalidades: *No* con *Universitarios* (residuo=-3.93), *Si* con *Primarios* (residuo=3.69), *Si* con *Universitarios* (residuo=2.57) y *No* con *Primarios* (residuo=-5.64).



	Primarios			Secundarios			Universitarios				
	n	%Fila	%Col.	n	%Fila	%Col.	n	%Fila	%Col.		
No	97	52.72	31.49	126	30.66	43.91	0.28	85	19.63	27.63	
Si	87	47.28	12.98	-5.64	288	69.34	39.68	-0.17	349	80.37	49.33

## Relación tipo de estudios y vacunación

	Primarios				Secundarios				Universitarios			
	n	%Fila	%Col.	Resid.	n	%Fila	%Col.	Resid.	n	%Fila	%Col.	Resid.
No	162	88.04	20.05	1.44	337	82.00	41.71	0.78	309	71.36	38.24	-1.7
Sí	22	11.96	10.00	<b>-2.77</b>	74	18.00	33.64	-1.49	124	28.64	56.36	<b>3.26</b>



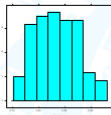
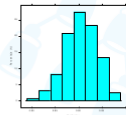


## Test estadísticos

- Univariante para cada uno de los niveles
- Normalidad
- Homocedasticidad:
  - Test F
  - Test de Ansary Bradley
  - Test de Barlett
  - Test de Flinger-Killen

	Frec. %	Acum. %
(0.707;0.816)	1	0.8
(0.816;0.850)	3	2.8
(0.850;0.850)	6	6.8
(0.850;0.871)	12	16.8
(0.871;0.880)	10	26.8
(0.880;0.888)	20	46.8
(0.888;0.907)	27	73.8
(0.907;0.940)	17	90.8
(0.940;0.940)	10	100.7
(0.940;0.952)	0	100.6
<b>Total</b>	<b>120</b>	<b>100.6</b>

	Frec. %	Acum. %
(0.707;0.774)	7	5.6
(0.774;0.774)	1	6.6
(0.774;0.807)	18	23.2
(0.807;0.809)	10	33.2
(0.809;0.840)	15	48.2
(0.840;0.871)	19	67.2
(0.871;0.871)	14	81.2
(0.871;0.880)	15	96.2
(0.880;0.880)	5	101.2
<b>Total</b>	<b>120</b>	<b>100.0</b>



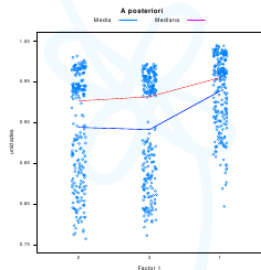
2. Respecto a la variable A posteriori (A), se disponen de 120 casos válidos. El valor medio se sitúa en 0.88 unidades, con una desviación típica de 0.04, mientras que la mediana disminuye hasta 0.83 unidades. El valor mínimo se sitúa en 0.70 y el máximo asciende a 0.95. El 50 % de las observaciones centrales se encuentran entre 0.8 y 0.87, mientras otros datos agrupados en un intervalo de amplitud 0.06 unidades, si bien el 25 % de los datos inferiores se encuentran en una amplitud de 0.05 y el 25 % de los valores más grandes se sitúan en una distancia de 0.06 unidades. No se detectan valores atípicos.

3. Respecto a la variable A posteriori (B), se disponen de 120 casos válidos. El valor medio se sitúa en 0.83 unidades, con una desviación típica de 0.03, mientras que la mediana aumenta hasta 0.83 unidades. El valor mínimo se sitúa en 0.70 y el máximo asciende a 0.94. El 50 % de las observaciones centrales se encuentran entre 0.8 y 0.86, mientras otros datos agrupados en un intervalo de amplitud 0.05 unidades, si bien el 25 % de los datos inferiores se encuentran en una amplitud de 0.04 y el 25 % de los valores más grandes se sitúan en una distancia de 0.06 unidades. No se detectan valores atípicos.

## Test Igualdad de medias

- Test paramétricos
  - Student
  - Welch
  - Anova
- Test no paramétricos
  - Wilcoxon
  - Kruskal-Wallis
  - Permutaciones

El test de permutaciones indica que en ciertos pares de niveles se producen diferencias significativas entre sí. Si los ordenados por orden de significatividad, resulta la siguiente prelación: 2 y 1 (p-valor=0) y 3 y 1 (p-valor=0), respectivamente. Por el contrario, no se producen diferencias entre los siguientes niveles: 3 y 2 (p-valor=0.75), respectivamente.



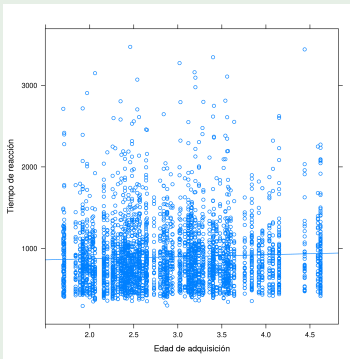
### Valoración

Se detectan diferencias significativas entre los grupos considerados.

Análisis

Regresión lineal

## Gráfico



# ¡Gracias por vuestra atención!

