Estadística Descriptiva de una variable con STATGRAPHICS

Ficheros empleados: AlumnosIndustriales.sf3,

1. Introducción

El objetivo de este documento es la utilización de las técnicas de estadística descriptiva más importantes para resumir la información de un conjunto de datos. Se usará el programa Stagraphics (v 4.1). Los datos corresponden a 91 estudiantes de Ingeniería Industrial, a los que se les ha preguntado sobre variables tales como estatura, peso, número de hermanos, etc. Emplearemos así una base de datos muy sencilla que nos ayude a entender las posibilidades del programa Stagraphics.

En primer lugar leemos el fichero de datos.

STATGRAPHICS Plus - Untitled StatFolio	
File Edit Plot Describe Compare Relate Special View Window He	lp
Open Den StatFolio	Ctrl+F11
Close Den Data File	Ctrl+F12
Save Open StatGallery	
Save As • Open StatReporter	
StatLink Query Database (ODBC Read Clipboard	I)
Open Data File Buscar Buscar The construction of the	
Tipo de SG PLUS Files (*.sf3;*.sfx,*.sf) Cancela	91

O bien pulsamos la tecla 🛄 y seleccionamos el fichero.

2. Descripción de variables cualitativas

La variable residencia corresponde al lugar de residencia de los alumnos. Su codificación es

- 1-Madrid Sur
- 2-Madrid Centro
- 3-Madrid-otros
- 4-Fuera de Madrid

Para describir esta variable haremos primeramente un análisis gráfico y luego una tabla de frecuencias. Todas estas opciones están en Describe/Categorical Data/Tabulation...

🤌 S	🚀 STATGRAPHICS Plus - Untitled StatFolio									
File	Edit	Plot	Describe	Compare	Relate	Specia	l View	Window	Help	
a	👝 💼 🙀 Numeric Data					→ [जिल्ली ह	-	🥶 I 🖂 I 🕫	
			Categorical Data			•	Tabulation			
	Distributions				•	Crosstabulation				
Life Data				•	Contingency Tables					

donde nos preguntan por la variable que deseamos analizar

altura dinero	Data:	idencia		
locomocion nacimiento				
residencia sexo				
Variables	_			
Sort				
OK	Cancel	Delete	Transform	Help

2.1 Gráfico de barras

En las opciones gráficas seleccionamos barchart



donde puede verse que la población más grande de alumnos son los procedentes de Madrid Sur, con casi 50 alumnos.

Si queremos cambiar el aspecto del gráfico nos colocamos sobre él, pulsamos el botón derecho del ratón y seleccionamos Pane Options. Seleccionamos entonces que las barras sean verticales y que las frecuencias sean en porcentajes.

Pane Options		Barchart Options	×
Analysis Options		Chart Type	ОК
Graphics Options Undo		Clustered Stacked	Cancel
Select Locate Zoom In Undo Zoom		C Frequencies	Help
Reset Scaling/Viewpoint			
Сору	Ctrl+C	 Horizontal 	
Print Print Preview	F4 Shift+F3	Vertical	
Copy Pane to Gallery Copy Analysis to StatReporte	r	Baseline: 0,	

También vamos a cambiar la estética de los rectángulos. En lugar de Pane Options seleccionamos Graphics Options. Cambiaremos el relleno de las barras (fills)

		Graphics Op	tions	Barchart for
		Layout Gr	id Fills	Top Title X-Axis Y-Axis Profile
		Fill 1	C 11	
Pane Options		0 2	C 12	
Analysis Options		O 3	C 13	
Graphics Options	_	─ ▶ ○ 4	O 14	
Undo		C 5	O 15	Color
Select		0.6	0 16	
Locate		07	0 17	
Zoom In		0.8	O 18	
Undo Zoom		0.9	O 19	
Reset Scaling/Viewpoint		0 10	O 20	
Сору	Ctrl+C			
D	F (
Print Preview	5hift+F3	Aceptar	Ca	ancelar Aplicar
Conv Page to Gallery		1		Definir colores personalizado
Copy Analysis to StatRepo	orter		_	Aceptar Cancelar

y obtenemos

y obtenemos

Barchart for residencia



Al estar la altura de las barras en % podemos ver que casi el 50% de los alumnos proceden del sur de Madrid

2.2 Gráfico de tarta o porciones

Dentro de las opciones gráficas, 🔛 seleccionamos Piechart

2.3 Tabla de frecuencias

En las opciones numéricas (Tabular Options) encontramos la tabla de frecuencias, que nos da la información numérica que antes hemos representado en gráficos.



y la tabla de frecuencias resultante es

Frequency Table for residencia

Class	Value	Frequency	Relative Frequency	Cumulative Frequency	Cum. Rel. Frequency
1	1	46	0,4842	46	0,4842
2	2	36	0,3789	82	0,8632
з	3	12	0,1263	94	0,9895
4	4	1	0,0105	95	1,0000

donde puede verse que entre el sur (48.4%) y el centro (37.9%) de Madrid, se abarca al 86.3% de los estudiantes

3. Descripción de variables cuantitativas

3.1 Análisis gráfico de variables discretas con pocos valores

En el caso de variables cuantitativas discretas con pocos valores, <u>el análisis gráfico es igual que para variables</u> <u>cualitativas</u>. Podemos suministrar un gráfico de barras. También la tabla de frecuencias sería igual que para el caso de variables cualitativas.

Por ejemplo, la variable hermanos proporciona el número de hermanos que tiene cada alumno. Para selecionar

una variable nueva dentro de un mismo análisis, basta con pulsar la tecla denominada Input Dialog



donde puede verse que las familias más frecuentes (entre las que tienen hijos cursando Ingeniería Industrial) son las de 2 y 3 hijos (1 y 2 hermanos)

3.2 Análisis gráfico de variables cuantitativas

El análisis gráfico y numérico de variables cuantitativas se hace en Describe/Numeric Data/One-Variable Analisys

🥕 S	🥕 STATGRAPHICS Plus - Untitled StatFolio										
File	Edit	Plot	Describe	Compare	Relate	Spec	ial	View	Window	Help	
🚗 📩 Numeric Data						•		One-Va	ariable Ana	alysis	
			Catego	Categorical Data				Multiple-Variable Analysis			
	Fabul	ation	Distribu	utions		•		Subset Analysis			
Life Data				•		Row-W	/ise Statist	ics			
Hypothesis Tests							Power	Transform	ations		
Fre	quen	cv Ta	T ₁ Sample-Size Determination			n Statistical Tolerance L				nce Limits	

La variable altura tiene las alturas de los alumnos. Seleccionamos esta variable

One-¥ariable Ana	ysis	X
altura dinero hermanos locomocion nacimiento peso residencia sexo tiempo Variables	Data: altura (Select.)	

Vamos a hacer su histograma y su diagrama box-plot. Seleccionamos las opciones gráficas que queremos.



Vamos a añadir en el Box-plot la marca de la posición de la media muestral. Nos colocamos en ese gráfico, pulsamos el botón derecho del ratón, y seleccionamos Pane Options



El Box-plot nos muestra que la distribución de las alturas es algo asimétrica. La caja central muestra una asimetría negativa, si bien las colas de la distribución no son muy largas. Este efecto se ve tambi⁻çen en el histograma. Vamos a cambiar el número de clases del histograma. Como tenemos 95 observaciones, tomaremos $\sqrt{95} \approx 10$ clases. Nos posicionamos en el histograma y con el botón derecho del ratón seleccionamos Pane Options



Este mayor número de clases nos muestra una bimodalidad que es imposible de visualizar en un boxplot. Hay una moda en torno a 165cm. Y otra en torno a 178 cm. Esas dos modas sugieren que la población no es homogénea. Es muy posible que se deba a las alturas de chicos y chicas (ver guión de análisis de varias variables).

La variable sexo tiene el sexo de los alumnos (1=chico, 0=chica). Vamos a emplear esa variable para seleccionar la altura de los chicos o de las chicas.

Si queremos seleccionar sólo a los chicos hacemos lo siguiente:



y vemos que con sólo los chicos, la distribución es muy simétrica, unimodal, con moda en los 180 cm. Si lo repetimos para las chicas tenemos



La distribución de las chicas es más uniforme. No es con forma de campana como la de los chicos. Tal vez sea porque hay menos datos (sólo 32) o porque las chicas de esa titulación sean realmente más heterogéneas.

3.3 Tabla de Frecuencias

La distribución de frecuencias mediante una tabla nos proporiocna la misma información que un histograma, pero nos permite ver los valores numéricos de las frecuencias de cada intervalo. Para hacer la tabla de

🖥 One-Variable Analysis - altura AN Row: Lbl: cy Tabula altura Tabular Options Lower Limit Upper Limit lative Midpoint Fre Frequency Analysis Summary or below 150,0 0,0000 153,125 🔲 Summary Statistic 150,0 156,25 0,0000 156,25 162,5 159,375 0,0947 Percentiles 162,5 168,75 16 0,1684 168.75 175.0 171.875 20 0.2105 Frequency Tabulation 173,0 181,25 187,5 193,75 178,125 184,375 190,625 196,875 175.0 32 16 2 0 193,0 181,25 187,5 193,75 200,0 0,1684 Stem-and-Leaf Display 0,0211 0,0000 0,0000 Confidence Intervals 200,0 Hypothesis Tests 174,621 Standard deviation = 8,22707 Cancel ПK All Help

frecuencias vamos a las opciones numéricas (Tabular Options 🗐)

La tabla de frecuencias tienen las mismas opciones que el histograma. Podemos cambiar el número de clases o limitar el rango de valores. Para acceder a estas opciones nos colocamos sobre la ventana de resultados y pulsamos el botón derecho del ratón. Selecccionamos entonces Pane Options. En la ventana que obtenemos seleccionamos 10 clases. Los cambios que propongamos para la tabla de frecuencias también afectan al histograma de frecuencias

🖥 One-V	ariable <i>i</i>	Analysis - a	ltura											
	2	IIII 幸· 容	<u>> 8? † (</u>	LbI:	A +B	Row:	<i>8</i> *8							
requenc	y Tabul	ation for	altura		Dolotivo	Cumulativa	Cum Bol			Histo	ogram	for a	ltura	
lass	Limit	Limit	Midpoint	Frequency	Frequency	Frequency	Fremiency							
at or 1 2 3 4 5 6 7 8 9	below 150,0 155,0 160,0 165,0 170,0 175,0 180,0 185,0 190,0	150,0 155,0 160,0 165,0 170,0 175,0 180,0 185,0 190,0 195,0	152,5 157,5 162,5 167,5 172,5 177,5 182,5 187,5 192,5	0 0 7 11 15 12 29 15 5 1	0,0000 0,0000 0,0737 0,1158 0,1579 0,1263 0,3053 0,1579 0,0526 0,0526	Frequency Number of Cl 10 Lower Limit: 150, Upper Limit:	Tabulation (asses:	Options OK Cancel Help						
10 howo	195,0	200,0	197,5	0	0,0000	200,		T Hold	ht.					
:ean = 1	74,621	Standard	l deviation	= 8,22707				<u>~</u>	150	160	170 alti	180 ura	190	200

Esta tabla muestra que los dos intervalos modales son alrededor de los valores (midpoint), 167.5 y 177.5 y que el intervalo más frecuente, el centrado en 177.5 contiene a más del 30% de los alumnos.

3.4 Medidas características de variables cuantitativas

Para calcular las medidas características de la variable altura vamos a las opciones numéricas (Tabular Options



Podemos selecionar las medidas que deseemos. Nos posicionamos en la ventana de resultados y seleccionamos Pane Options. Aparece una ventana con todos los estadísticos univariantes que calcula el Statgraphics. Si los selecionamos todos obtenemos los siguientes resultados:

Summary Statistics Options								
🔽 Average	🔽 Min.	🔽 Skewness						
🔽 Median	🔽 Max.	🔽 Std. Skewness						
🔽 Mode	🔽 Range	🔽 Kurtosis						
🔽 Geo. Mean	🔽 Lower Quartile	🔽 Std. Kurtosis						
🔽 Variance	🔽 Upper Quartile	🔽 Coeff. of Var.						
🔽 Std. Deviation	🔽 Interquartile Range	🔽 Sum						
🔽 Std. Error								
ОК	Cancel All	Help						



Summary Statistics for altura

```
Count = 95
Average = 174,621
Median = 177,0
Mode = 180,0
Geometric mean = 174,427
Variance = 67,6847
Standard deviation = 8,22707
Standard error = 0,844079
Minimum = 158,0
Maximum = 193,0
Range = 35,0
Lower quartile = 168,0
Upper quartile = 180,0
Interquartile range = 12,0
Skewness = -0,302876
Stnd. skewness = -1,20518
Kurtosis = -0,855173
Stnd. kurtosis = -1,70142
Coeff. of variation = 4,71138%
Sum = 16589,0
```

Es necesario hacer algunas puntualizaciones sobre estas medidas características:

- Al ser la altura una medida continua, LA MODA NO TIENE SENTIDO. La moda es el valor más frecuente, y en una variable continua podría suceder que no se repitiese ningún valor. En esos casos, el programa nos devolvería el primer valor que leyese en el fichero de datos. En este tipo de variables sólo tiene sentido hablar de intervalo modal de un histograma
- La varianza está calculada con la expresión

$$s^{2} = \frac{\sum_{i=1}^{n} (x_{i} - \overline{x})^{2}}{n-1}$$

en lugar de

$$s^{2} = \frac{\sum_{i=1}^{n} \left(x_{i} - \overline{x}\right)^{2}}{n}$$

La conveniencia de dividir por *n-1* en lugar de *n* no es inmediata, y su justificación teórica se verá en temas más avanzados.

- La desviación típica usa también esta misma formulación, dividiéndose por n-1.
- Las siguientes medidas
 - Standard Error
 - Stnd. Skewness
 - Stnd. Kurtosis

no son propiamente de estadística descriptiva, sino de inferencia. Por tanto no se cubren en este documento.

• El coeficiente de curtosis es

$$K = \frac{\sum_{i=1}^{n} (x_i - \overline{x})^4}{n \times s^4} - 3$$

por tanto, para una variable que tenga forma de campana, la curtosis es 0.

A continuación se muestra la comparación entre chicos y chicas mediante algunas medidas características



Puede verse ahora que los chicos son más altos que las chicas. La estaura media de los chicos es 179 mientras que para las chicas es sólo 165. Los chicos se concentran mucho alrededor de esa media. Puede verse que la media es casi igual que la moda, lo que se ve también fácilmente en la simetría del histograma y el bajo valor del coeficiente se asimetría. Esa concentración cerca de la media se ve también en el rango intercuartílico. El 50% de los chicos colocados en las posiciones centrales sólo se diferencian en un máximo de 5 centímetros, mientras que para las chicas ese rango intercuartílico es de 7 cm. Otra medida que resume esa mayor concentración es la curtosis. En los chicos es positiva y alta, mientras que en las chicas es negativa. En este caso, la curtosis es la medida que mejor resume la diferencia entre esos dos histogramas.

3.5 Percentiles



Para la obtención de los percentiles seleccionamos Tabular Options...Percentiles...OK.

Para seleccionar algún percentil concreto, nos ponemos sobre la ventana de resultados, pulsamos el botón derecho del ratón y selecionamos Pane Olptions. Vamos a calcular el percentil 20 y 80.

E I I I I I I I I I I I I I I I I I I I	Lbl:	##	Row:
Percentiles for altura			
20,0% = 167,0	Percentiles (Options	\mathbf{X}
80,0% = 181,0	Percentiles:		ОК
The StatAdvisor	20.	0,	Cancel
This pane shows sample percentil values below which specific percent;	0,	0,	Help
can see the percentiles graphically the list of Graphical Options.	0,	0,	
	U,	0,	

El 20% de los alumnos mide menos de 167 cm, y el 80% mide menos de 181 cm.