

**Tratamiento de datos  
con R, STATISTICA y SPSS**

**CÁSTOR GUISANDE GONZÁLEZ**

Catedrático del área de Ecología. Universidad de Vigo (España)

**ANTONIO VAAMONDE LISTE**

Catedrático del área de Estadística e Investigación Operativa.  
Universidad de Vigo (España)

**ALDO BARREIRO FELPETO**

Investigador del área de Ecología. Universidad de Vigo (España)

# **Tratamiento de datos con R, STATISTICA y SPSS**



© Cástor Guisande González, Antonio Vaamonde Liste, Aldo Barreiro Felpeto, 2011  
Reservados todos los derechos.

Queda prohibida, salvo excepción prevista en la ley, cualquier forma de reproducción, distribución, comunicación pública y transformación de esta obra sin contar con la autorización de los titulares de propiedad intelectual. La infracción de los derechos mencionados puede ser constitutiva de delito contra la propiedad intelectual (art. 270 y siguientes del Código Penal). El Centro Español de Derechos Reprográficos (CEDRO) vela por el respeto de los citados derechos.

Ediciones Díaz de Santos  
E-mail: ediciones@diazdesantos.es  
Internet: <http://www.diazdesantos.es/ediciones>

ISBN: 978-84-7978-998-5  
Depósito legal: M. 7.321-2011

Diseño de cubierta: P55 Servicios Culturales  
Impresión: FER Fotocomposición  
Impreso en España

# ÍNDICE

---

<b>Prólogo y agradecimientos.....</b>	<b>XV</b>
<b>I. REPRESENTACIÓN DE DATOS.....</b>	<b>1</b>
I.1. COORDENADAS POLARES.....	1
I.1.1. Estandarización de los datos.....	2
I.1.2. Asignación de ángulos a las variables.....	2
I.1.3. Representación de las coordenadas polares.....	7
<b>II. ESTADÍSTICA DESCRIPTIVA.....</b>	<b>21</b>
II.1. MEDIDAS DE POSICIÓN.....	21
II.1.1. Medidas de posición central.....	21
II.1.1.1. <i>Media aritmética</i> .....	21
II.1.1.2. <i>Media geométrica</i> .....	22
II.1.1.3. <i>Media armónica</i> .....	22
II.1.1.4. <i>Moda</i> .....	23
II.1.1.5. <i>Mediana</i> .....	23
II.1.1.6. <i>Media acotada o recortada</i> .....	25
II.1.1.7. <i>Media ponderada</i> .....	25
II.1.2. Otras medidas de posición.....	27
II.2. MEDIDAS DE DISPERSIÓN.....	29
II.2.1. Amplitud.....	30
II.2.2. Varianza y cuasivarianza.....	30
II.2.3. Desviación típica y cuasidesviación típica.....	30
II.2.4. Desviación absoluta respecto a la media y la mediana.....	31
II.2.5. Coeficiente de variación.....	33
II.2.6. Error estándar de la media.....	33
II.2.7. Recorrido intercuartílico.....	33
II.3. ESTADÍSTICA DESCRIPTIVA CON R.....	33
II.4. ESTADÍSTICA DESCRIPTIVA CON SPSS.....	37
II.5. ESTADÍSTICA DESCRIPTIVA CON STATISTICA.....	39
<b>III. DISTRIBUCIÓN.....</b>	<b>47</b>
III.1. TIPOS DE VARIABLES.....	47
III.2. DISTRIBUCIONES PARA VARIABLES CONTINUAS.....	48

III.2.1. Normal.....	48
III.2.1.1. Aplicaciones de la distribución Normal....	49
III.2.1.2. Asimetría.....	56
III.2.1.3. Apuntamiento o curtosis.....	57
III.2.2. <i>t</i> de Student.....	58
III.2.3. Chi-cuadrado.....	60
III.2.4. <i>F</i> de Fisher-Snedecor.....	61
III.3. DISTRIBUCIONES PARA VARIABLES DISCRETAS.....	62
III.3.1. Binomial.....	62
III.3.2. Hipergeométrica.....	64
III.3.3. Poisson.....	66
<b>IV. INTERVALOS DE CONFIANZA.....</b>	<b>71</b>
IV.1. INTERVALO DE CONFIANZA DE LA MEDIA DE UNA POBLACIÓN NORMAL.....	71
IV.1.1. Desviación típica conocida.....	71
IV.1.2. Desviación típica desconocida.....	71
IV.1.2.1. Tamaño de muestra pequeño ( $< 30$ ).....	72
IV.1.2.2. Tamaño de muestra grande ( $\geq 30$ ).....	72
IV.2. INTERVALO DE CONFIANZA DE LA DIFERENCIA DE MEDIAS DE DOS POBLACIONES NORMALES.....	83
IV.2.1. Varianzas conocidas.....	83
IV.2.2. Varianzas desconocidas.....	83
IV.2.2.1. Tamaño de muestra grande ( $\geq 30$ ).....	83
IV.2.2.2. Tamaño de muestra pequeño ( $< 30$ ).....	83
IV.2.2.2.1. Varianzas iguales.....	84
IV.2.2.2.2. Varianzas diferentes.....	84
IV.3. INTERVALO DE CONFIANZA DE LA VARIANZA Y DESVIACIÓN TÍPICA DE UNA POBLACIÓN NORMAL.....	91
IV.4. INTERVALO DE CONFIANZA DE LA RAZÓN DE VARIANZAS DE DOS POBLACIONES NORMALES.....	94
<b>V. CONTRASTE DE HIPÓTESIS.....</b>	<b>103</b>
V.1. TIPOS DE HIPÓTESIS.....	103
V.2. ESTADÍSTICO DE CONTRASTE.....	103
V.3. TIPOS DE ERRORES.....	103
V.4. REGIONES CRÍTICAS Y DE ACEPTACIÓN.....	104

V.5. TIPOS DE CONTRASTES.....	106
<b>VI. CONTRASTES DE BONDAD DE AJUSTE.....</b>	<b>107</b>
VI.1. CHI-CUADRADO.....	107
VI.2. TEST DE KOLMOGOROV-SMIRNOV.....	108
VI.3. TEST DE SHAPIRO-WILK .....	108
<b>VII. CONTRASTES DE HOMOGENEIDAD EN VARIABLES CUALITATIVAS.....</b>	<b>135</b>
VII.1. MUESTRAS INDEPENDIENTES.....	135
VII.1.1. Variables politómicas.....	135
<i>VII.1.1.1. Chi-cuadrado</i> .....	136
<i>VII.1.1.2. Razón de verosimilitud (test G)</i> .....	136
VII.1.2. Variables dicotómicas.....	150
<i>VII.1.2.1. Corrección de Yates</i> .....	150
<i>VII.1.2.2. Prueba de Fisher</i> .....	150
VII.2. MUESTRAS RELACIONADAS.....	161
VII.2.1. Variables politómicas.....	161
<i>VII.2.1.1. ANOVA de Friedman</i> .....	161
VII.2.2. Variables dicotómicas.....	167
<i>VII.2.2.1. Prueba de McNemar</i> .....	167
<i>VII.2.2.2. Q de Cochran</i> .....	179
<b>VIII. CONTRASTES DE INDEPENDENCIA Y ASOCIACIÓN EN VARIABLES CUALITATIVAS.....</b>	<b>189</b>
VIII.1. VARIABLES POLITÓMICAS.....	189
VIII.1.1. Pruebas de independencia.....	189
<i>VIII.1.1.1. Chi-cuadrado de Pearson y test G de razón de verosimilitud</i> .....	189
VIII.1.2. Medidas de asociación.....	189
<i>VIII.1.2.1. V de Cramer</i> .....	189
<i>VIII.1.2.2. Coeficiente de Contingencia</i> .....	190
<i>VIII.1.2.3. Coeficiente de Incertidumbre</i> .....	190
VIII.2. VARIABLES DICOTÓMICAS.....	202
VIII.2.1. Pruebas de independencia.....	202
<i>VIII.2.1.1. Corrección de Yates y prueba de Fisher</i> ...	202
<i>VIII.2.1.2. Pruebas de Cochran y Mantel-Haenszel</i> ...	202
VIII.2.2. Medidas de asociación.....	203
<i>VIII.2.2.1. Phi (<math>\phi</math>)</i> .....	203
<i>VIII.2.2.2. Riesgo relativo, razón de ventajas (odds ratio) y las pruebas de Breslow-Day y de Tarone</i> ..	203

<b>IX. CONTRASTES DE HOMOGENEIDAD EN VARIABLES CUANTITATIVAS.....</b>	217
IX.1. PRUEBAS PARAMÉTRICAS.....	217
IX.1.1. Requisitos.....	217
IX.1.2. Transformaciones.....	218
IX.1.3. t-test.....	218
IX.1.3.1. Muestras independientes.....	218
IX.1.3.2. Muestras dependientes.....	240
IX.1.4. Análisis de varianza.....	251
IX.1.5. Análisis de covarianza.....	289
IX.1.6. Análisis de varianza y covarianza con medidas repetidas.....	318
IX.1.7. Análisis de varianza anidado.....	363
IX.2. PRUEBAS NO PARAMÉTRICAS.....	386
IX.2.1. Contrastes para dos muestras independientes....	386
IX.2.1.1. Prueba U de Mann-Whitney.....	386
IX.2.1.2 Test de rachas de Wald-Wolfowitz.....	387
IX.2.1.3. Prueba de Kolmogorov-Smirnov para dos muestras.....	388
IX.2.2. Contrastes para k-muestras independientes.....	400
IX.2.2.1. Contraste de la mediana.....	400
IX.2.2.2 ANOVA de Kruskal-Wallis.....	400
IX.2.3. Contrastes para dos muestras dependientes.....	408
IX.2.3.1. Contraste de los signos.....	408
IX.2.3.2. Prueba de Wilcoxon para pares relacionados.....	408
IX.2.4. Contrastes para k-muestras dependientes.....	415
<b>X. CONTRASTES DE INDEPENDENCIA Y ASOCIACIÓN EN VARIABLES CUANTITATIVAS.....</b>	421
X.1. CORRELACIONES BIVARIADAS SIMPLES.....	421
X.1.1. Coeficiente de correlación de Pearson ( $r$ ).....	421
X.1.2. Coeficiente de correlación de Spearman ( $\rho$ ).....	422
X.1.3. Coeficiente $\tau$ de Kendall.....	423
X.1.4. Coeficiente Gamma ( $\gamma$ ).....	424
X.2. CORRELACIONES BIVARIADAS PARCIALES....	439
X.2.1. Coeficiente de correlación parcial de Pearson....	439
X.2.2. Coeficiente de correlación parcial de Kendall ( $T_{xy.z}$ ).....	440
X.3. CORRELACIÓN MÚLTIPLE.....	449
X.3.1. Coeficiente de correlación múltiple de Pearson....	449

X.3.2. Coeficiente de concordancia de Kendall ( $W$ ).....	450
<b>XI. REGRESIONES .....</b>	<b>453</b>
<b>XI.1 MODELOS DE REGRESIÓN PARA VARIABLES     DEPENDIENTES CUANTITATIVAS.....</b>	<b>453</b>
XI.1.1. Requisitos.....	453
XI.1.2. Regresión simple.....	455
XI.1.3. Regresión múltiple.....	486
XI.1.4. Otras regresiones simples o múltiples no lineales.....	519
<i>XI.1.4.1. Curva logística.....</i>	<i>519</i>
<i>XI.1.4.2. Curva de crecimiento de von Bertalanffy...</i>	<i>520</i>
<i>XI.1.4.3. Curva de crecimiento de Gompertz.....</i>	<i>521</i>
<i>XI.1.4.4. Relación entre tasas y variables.....</i>	<i>521</i>
<b>XI.2 MODELOS DE REGRESIÓN PARA VARIABLES     DEPENDIENTES CUALITATIVAS.....</b>	<b>537</b>
XI.2.1. Regresión logística binomial.....	537
XI.2.2. Regresión logística multinomial.....	557
<b>XII. SERIES TEMPORALES.....</b>	<b>585</b>
<b>XII.1. MODELO ARIMA.....</b>	<b>585</b>
XII.1.1. Autocorrelación.....	585
XII.1.2. Series estacionarias o no estacionarias.....	586
XII.1.3. Estacionalidad.....	586
XII.1.4. Covariables.....	586
XII.1.5. Construcción del modelo.....	586
<b>XII.2. CICLOS.....</b>	<b>637</b>
<b>XIII. ANÁLISIS MULTIVARIANTE: MÉTODOS FAC- TORIALES.....</b>	<b>667</b>
XIII.1. ANÁLISIS DE COMPONENTES PRINCIPALES.....	667
XIII.2. ANÁLISIS FACTORIAL.....	697
XIII.3. ANÁLISIS DE CORRESPONDENCIAS.....	719
XIII.4. ESCALAMIENTO MULTIDIMENSIONAL.....	743
<b>XIV. ANÁLISIS MULTIVARIANTE: MÉTODOS DE CLASIFICACIÓN.....</b>	<b>759</b>
<b>XIV.1. CLASIFICACIÓN JERÁRQUICA.....</b>	<b>759</b>
XIV.1.1. Criterio de distancia.....	760
XIV.1.2. Criterio de aglomeración.....	760
XIV.1.3. Dendrograma.....	761



XIV.2. CLASIFICACIÓN DE <i>K</i> -MEDIAS.....	773
XIV.3. ANÁLISIS DISCRIMINANTE.....	787
XIV.3.1. Hipótesis.....	787
XIV.3.2. Método paso a paso.....	787
XIV.3.3. Influencia de las distintas variables independientes.....	788
XIV.3.4. Validación.....	788
<b>XV. MODELOS DE SIMULACIÓN.....</b>	<b>811</b>
XV.1. EL USO DE MODELOS.....	811
XV.2. PASOS A CONSIDERAR PARA EL DESARROLLO DE UN MODELO.....	812
XV.3. INTRODUCCIÓN AL MODELADO CON EL PROGRAMA STELLA.....	813
XV.3.1. Significado de los iconos y menús específicos de la barra de herramientas en niveles <i>Model</i> , <i>Map</i> e <i>Interface</i> .....	814
XV.3.2. Ejemplo de funcionamiento de los iconos y menús básicos.....	819
XV.4. EJEMPLO DE MODELADO CON EL PROGRAMA STELLA: DESARROLLO CONCEPTUAL Y MANEJO PRÁCTICO.....	832
XV.4.1. Dinámica poblacional de la presa.....	833
XV.4.2. Dinámica poblacional del depredador y su influencia sobre la presa.....	834
XV.4.3. Influencia de la explotación sobre la población del depredador.....	836
<b>BIBLIOGRAFÍA.....</b>	<b>839</b>
<b>GUÍA RESUMEN.....</b>	<b>845</b>
<b>APÉNDICE I. TABLAS ESTADÍSTICAS.....</b>	<b>851</b>
Tabla 1. Áreas bajo la curva Normal estándar.....	851
Tabla 2. Valores críticos de la distribución <i>t</i> de Student.....	852
Tabla 3. Valores críticos de la distribución $\chi^2$ .....	853
Tabla 4. Valores críticos de la distribución <i>F</i> de Fisher-Snedecor.....	861
Tabla 5. Valores críticos del estadístico de Kolmogorov-Smirnov.....	873
Tabla 6. Valores críticos del estadístico de Lilliefors.....	874

Tabla 7. Distribución del estadístico de contraste de Durbin-Watson.....	875
--	-----

<b>APÉNDICE II. INSTALACIÓN Y CONCEPTOS BÁSICOS SOBRE R .....</b>	<b>877</b>
<b>AII.1. SOBRE R.....</b>	<b>877</b>
AII.1.1. ¿Qué es R?.....	877
AII.1.2. Instalación.....	877
AII.1.3. Interfaz.....	878
<b>AII.2. PRIMEROS PASOS.....</b>	<b>879</b>
AII.2.1. Entrar/salir de R y seleccionar directorio de trabajo.....	879
AII.2.2. Crear un archivo script.....	880
AII.2.3. Instalar y cargar paquetes y conjuntos de datos.....	882
AII.2.4. Actualización de R.....	884
AII.2.5. Ayuda, ejemplos y demos.....	887
AII.2.6. Documentación sobre R .....	890
<b>AII.3. R COMMANDER .....</b>	<b>890</b>
AII.3.1. Características de R Commander.....	890
AII.3.2. El menú de R Commander.....	893
<i>AII.3.2.1. Fichero.....</i>	<i>893</i>
<i>AII.3.2.2. Editar.....</i>	<i>893</i>
<i>AII.3.2.3. Datos.....</i>	<i>894</i>
<i>AII.3.2.4. Estadísticos.....</i>	<i>895</i>
<i>AII.3.2.5. Gráficas.....</i>	<i>895</i>
<i>AII.3.2.6. Modelos.....</i>	<i>896</i>
<i>AII.3.2.7. Distribuciones.....</i>	<i>896</i>
<i>AII.3.2.8. Herramientas.....</i>	<i>897</i>
<i>AII.3.2.9. Ayuda.....</i>	<i>897</i>
<b>AII.4. EL LENGUAJE R.....</b>	<b>897</b>
AII.4.1. Apuntes preliminares.....	897
AII.4.2. Importación de datos.....	900
AII.4.3. Guardar datos y generar archivos de salida.....	903
AII.4.4. Objetos.....	906
AII.4.5. Operadores.....	909
AII.4.6. Funciones.....	909
<i>AII.4.6.1. Funciones built in.....</i>	<i>909</i>
<i>AII.4.6.2. Funciones definidas por el usuario.....</i>	<i>914</i>
AII.4.7. Bucles, ciclos o <i>loops</i> .....	918
AII.4.8. Vectores.....	920
<i>AII.4.8.1. Creación de vectores.....</i>	<i>920</i>

AII.4.8.2. <i>Trabajando con vectores</i> .....	922
AII.4.9. <i>Matrices y Arrays</i> .....	924
AII.4.9.1. <i>Creación de matrices y Arrays</i> .....	924
AII.4.9.2. <i>Trabajando con matrices y Arrays</i> .....	929
AII.4.10. <i>Listas</i> .....	933
AII.4.10.1. <i>Creación de listas</i> .....	933
AII.4.10.2. <i>Trabajando con listas</i> .....	934
AII.4.11. <i>Data frames</i> .....	936
AII.4.11.1. <i>Creación de data frames</i> .....	936
AII.4.11.2. <i>Seleccionando variables y casos en data frames</i> .....	939
AII.4.11.3. <i>Funciones importantes en data frames</i> ....	944
AII.4.12. <i>Editor de datos</i> .....	948
AII.4.13. <i>Texto</i> .....	948
AII.4.14. <i>Gráficos</i> .....	950
AII.4.14.1. <i>Aspectos generales</i> .....	950
AII.4.14.2. <i>Gráfico de dispersión</i> .....	952
AII.4.14.3. <i>Representar modelos y gráficos en paneles</i> .....	957
AII.4.14.4. <i>Diagrama de cajas</i> .....	958
AII.4.14.5. <i>Graficar curvas de ajuste y caracteres especiales</i> .....	960
AII.4.14.6. <i>Histograma</i> .....	961
AII.4.14.7. <i>Gráficos de barras</i> .....	965
AII.4.14.8. <i>Gráficos de contornos y superficies</i> .....	969
AII.4.14.9. <i>Configuración de márgenes y ejes adicionales</i> .....	971
AII.4.14.10. <i>Formas especiales</i> .....	972
<b>ÍNDICE DE CONCEPTOS</b> .....	975

# PRÓLOGO Y AGRADECIMIENTOS

---

Los investigadores necesitan la Estadística para obtener conclusiones válidas a partir de los datos. El método estadístico se ha convertido en parte esencial de la generación de conocimiento científico, y en prácticamente todas las publicaciones especializadas –en casi cualquier campo de conocimiento– las técnicas estadísticas tienen un papel muy importante, hasta el punto de que las revistas científicas suelen disponer de revisores con amplios conocimientos en este campo.

Aunque los usuarios de los métodos de tratamiento y análisis de datos tienen en general conocimientos de Estadística, no suelen ser expertos; para resolver problemas específicos tienen dos opciones: consultar con un especialista –siempre recomendable si el problema es complejo– o recurrir a algún manual de consulta. Frecuentemente los libros de Estadística son excesivamente teóricos, prestando atención únicamente a los aspectos matemáticos, demostraciones, teoremas y propiedades (sin duda de gran importancia para conocer en profundidad las técnicas estadísticas), o bien libros que omiten totalmente los principios básicos y el esquema conceptual que permite comprender el funcionamiento de cada técnica estadística, para centrar todos sus esfuerzos en la aplicación en sí misma, presentada como receta al alcance del lector profano.

Los libros más teóricos, pese a su corrección formal, resultan de escasa utilidad para el usuario normal de las técnicas estadísticas, ya que se ve obligado a convertirse en un experto, o al menos a intentarlo, dedicando a ello más tiempo del que dispone, para poder llegar a resolver sus problemas inmediatos de tratamiento de datos, a menudo sin lograrlo. Los libros más aplicados, por el contrario, parecen muy adecuados para el lector que tiene un problema estadístico similar al descrito en el capítulo correspondiente, pero pronto averigua que lamentablemente los problemas son siempre distintos, que existen cuestiones que la receta no trata o no resuelve, y que al utilizarla nunca está a salvo de cometer errores insalvables, que invalidarán los resultados de su investigación.

El principal objetivo de este libro es servir de ayuda a los investigadores, cualquiera que sea su ámbito de conocimiento, que necesitan utilizar técnicas estadísticas para tratar y analizar sus datos. Con un enfoque centrado en la aplicación, se hace énfasis en la explicación de cada método estadístico con un lenguaje accesible, tratando de aclarar de forma breve pero razonada los principios básicos de cada técnica. Numerosas aplicaciones con casos reales son utilizadas para fijar los conceptos y explicar los pasos a seguir en el tratamiento de los datos, desde la introducción de datos hasta la interpretación de los resultados obtenidos.

Se utilizan tres programas estadísticos: SPSS, STATISTICA, y R. Todos ellos son programas de ámbito general muy conocidos y utilizados para el tratamiento estadístico de datos en diversos campos. R es un programa estadístico de uso libre y gratuito, que recientemente ha alcanzado una gran difusión. A la ventaja –nada desdeñable– de su utilización gratuita, une una potencia de análisis superior a la de cualquier otro programa de carácter general, junto con un ritmo de desarrollo de aplicaciones que permite adivinar que pronto desplazará del mercado a los programas más costosos. R es un proyecto abierto y participativo en el

que colaboran –de forma libre y espontánea– miles de investigadores estadísticos, lo que asegura el debate, la crítica, y la actualización permanente de los procedimientos; su único inconveniente, por el momento, es el entorno o interfase de aplicación, con menús de usuario poco (o nada) desarrollados en algunas técnicas, que utiliza un lenguaje de programación que puede parecer complejo a los no iniciados.

Todas las técnicas son aplicadas con los tres programas. Para ello se utilizan datos distintos, aunque a veces el mismo ejemplo es resuelto con dos programas diferentes para que el lector pueda comparar y deducir sus ventajas y limitaciones. Se trata de conseguir que el lector pueda reproducir con sus propios datos cualquier aplicación y, por ello, se sigue siempre un esquema de pasos sucesivos, con el que debe resultar razonablemente sencillo completar un nuevo análisis mediante cualquiera de los tres programas estadísticos mencionados.

Todos los ejemplos se explican detalladamente, para que puedan ser repetidos incluso por el usuario que no tiene conocimientos previos de Estadística ni del programa. Las distintas técnicas estadísticas son presentadas de forma independiente, de manera que el lector interesado en cualquiera de ellas pueda empezar directamente por el capítulo correspondiente.

Los supuestos o hipótesis que cada método requiere son verificados cuidadosamente: es necesario que el lector compruebe también todas las hipótesis en cualquier aplicación de la misma técnica, y lo haga constar así junto a la presentación de los resultados. Los métodos estadísticos solo son válidos cuando se aplican correctamente en su contexto adecuado, y es muy frecuente encontrar aplicaciones desafortunadas en las que obviamente no se cumplen las condiciones necesarias, o en las que simplemente no se han comprobado.

Hemos prestado también especial atención a la interpretación correcta de los resultados obtenidos. Es bastante habitual observar –en la publicación de resultados– cómo se obtienen conclusiones que no pueden deducirse de los datos, generalmente porque se interpreta de forma laxa o generalizada el resultado de una prueba estadística, porque se ignoran las características de los datos utilizados, o porque se desconocen las posibilidades y limitaciones del método estadístico. Aunque en este libro se intenta explicar con claridad cómo deben interpretarse correctamente los distintos resultados, lo cierto es que éste es el aspecto más difícil para el usuario no experto. El único antídoto frente al mal uso de la Estadística consiste en un conocimiento profundo de las técnicas, por lo que recomendamos vivamente que el usuario consulte a un especialista cuando tenga dudas de interpretación en la aplicación de métodos estadísticos complejos.

El libro se acompaña de un CD con los archivos de datos de todos los ejemplos utilizados, así como de los programas de R necesarios para su análisis, con el fin de facilitar su utilización. Una breve guía resumen al final del libro puede utilizarse para una primera orientación sobre la técnica estadística que debe ser utilizada, en función del tipo de datos del que se dispone y de los objetivos buscados con el análisis.

Esta obra no tiene la pretensión de ser exhaustiva. Muchas técnicas simples de Estadística Descriptiva y algunos métodos estadísticos de carácter más específico no son tratados para no aumentar excesivamente su extensión. Hemos tratado de incluir la mayoría de las técnicas de tratamiento estadístico de datos que los investigadores necesitan en su trabajo habitual, añadiendo la resolución de problemas frecuentes que no tienen en otros textos una consideración suficiente, y hemos procurado también situarnos en el lugar del usuario para dar una respuesta satisfactoria a sus demandas, de modo que pueda superar sin dificult-

tades las barreras que a menudo impiden la aplicación de esta herramienta imprescindible que es la Estadística.

Por último, queremos agradecer al profesor Kenneth Roy Cabrera Torres de la Universidad Nacional de Colombia Sede Medellín, su ayuda con los *scripts* de R en los primeros capítulos de este libro.



## REPRESENTACIÓN DE DATOS

---

### I.1. Coordenadas polares

Un paso obligado previo a la realización de cualquier tratamiento de datos es representar los datos gráficamente. Esto es necesario por muchos motivos: por ejemplo, para ver el tipo de relación que existe entre dos variables (lineal, logarítmica, exponencial, etc.); para identificar posibles "outliers", datos que son muy diferentes del resto y que se pueden deber simplemente a que hemos introducido mal los datos en el ordenador; para ver el tipo de distribución y variabilidad de los datos, etc. Existen muchos tipos de gráficas que permiten hacer representaciones, que son bien conocidas y de uso común por la mayoría de las personas que trabajan con datos. En este capítulo solo vamos a tratar un tipo de representación que no es tan conocido, el gráfico de coordenadas polares.

Un problema frecuente que surge a la hora de mostrar gráficamente los resultados obtenidos es que es necesario representar más de dos ejes de coordenadas. Sin embargo, en un plano bidimensional lo máximo que se puede dibujar son tres ejes. Las coordenadas polares permiten representar en un gráfico bidimensional cualquier número de ejes de coordenadas.

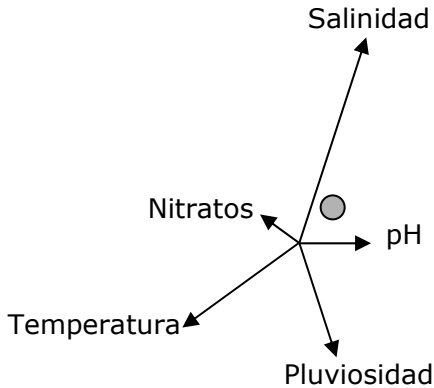
Supongamos que se quiere representar el nicho de varias especies de árboles considerando como variables del nicho el pH medio del suelo en el que aparece la especie, la temperatura media, la salinidad media del suelo, la pluviosidad media y la concentración media de nitratos en el suelo. Las coordenadas polares consideran las diferentes variables como vectores de fuerza, de tal forma que las coordenadas polares  $X$  e  $Y$  de cada especie serían el punto de equilibrio de todos los vectores es decir, de todas las variables (Figura I.1).

El módulo de cada vector sería el valor de la variable y el ángulo de cada vector tendría un valor convenido asignado a cada variable. Por lo tanto, cada especie tendría unas coordenadas  $X$  e  $Y$ , lo cual permite su representación en un plano bidimensional y, además, estas coordenadas  $X$  e  $Y$  vendrían determinadas por el efecto conjunto de todas las variables consideradas para definir el nicho de las especies. Es decir, cada especie ocuparía una posición en el plano en función de las variables del nicho, lo cual permitiría ver gráficamente que especies tienen un nicho más parecido.

Debido a que la representación polar se basa en vectores de fuerza, solo se puede utilizar con variables que tengan igual peso, es decir, que los valores de las variables puedan ser comparables. En el caso de que las variables no sean equiparables es necesario realizar una estandarización previa de los datos.

En un análisis multivariante de componentes principales, las puntuaciones obtenidas en cada eje para cada caso serían un ejemplo de variables que son en general equiparables. Los porcentajes de alimento consumido por varias especies también serían otro caso de variables equiparables, así como, por ejemplo, el porcentaje granulométrico del suelo.





**Figura I.1** Esquema de un sistema de coordenadas polares donde el círculo gris indica el centro de gravedad o punto de equilibrio de todas las variables. Las coordenadas de ese punto gris serían las coordenadas X e Y de la especie.

Por el contrario, si por ejemplo queremos una representación de varias lagunas, para ver gráficamente qué lagunas tienen unas características físico-químicas semejantes, considerando la temperatura, conductividad, pH, oxígeno, etc., estas variables no serían equiparables y sería necesario realizar una estandarización previa.

### I.1.1. ESTANDARIZACIÓN DE LOS DATOS

Como ejemplo de estandarización previa al cálculo de las coordenadas polares vamos a trabajar con datos físico-químicos de distintas lagunas, para las que, en algunas de ellas, existen varias estaciones de muestreo.

En la Tabla I.1 (Cuadro I.1) se muestran los datos y en el archivo **Cuadro I.1.xls**, además de los datos, se muestran todos los pasos a seguir para estandarizar los datos, los cuales también se describen en el Cuadro I.1.

En el Cuadro I.2 hay otro ejemplo a modo de práctica, del método para estandarizar datos.

### I.1.2. ASIGNACIÓN DE ÁNGULOS A LAS VARIABLES

Un paso previo al cálculo de las coordenadas es asignar ángulos a las variables consideradas. La asignación de ángulos a las variables es diferente si hay valores negativos (por ejemplo una salida de un análisis multivariante) o si, por el contrario, solo hay valores estandarizados en una escala de 0 a 1 y, por tanto, no hay valores negativos.

En el caso de que solo existan valores positivos en una escala de 0 a 1 (los valores están estandarizados), lo que se hace es restar a todos los valores 0,5 antes de comenzar con el proceso de asignar ángulos a las variables. Si el rango de los datos no es de 0 a 1 (los valores no están estandarizados) ya que, por ejemplo, es una salida de un análisis multivariante que tiene valores negativos, no se resta el valor de 0,5 a todos los datos.

### CUADRO I.1. Estandarización de datos

**EJEMPLO.** La estandarización de los datos físico-químicos de diferentes lagunas se muestra en la siguiente tabla. Los datos están en el archivo **Cuadro I.1.xls**.

**Tabla I.1.** Datos físico-químicos de distintas lagunas.

Laguna	Estación	Temperatura (°C)	Conductividad ( $\mu\text{S cm}^{-1}$ )	pH	Nitrato ( $\mu\text{M}$ )	Nitrito ( $\mu\text{M}$ )	Amonio ( $\mu\text{M}$ )	Fosfato ( $\mu\text{M}$ )	Silicato ( $\mu\text{M}$ )
1	1	29,1	495	6,08	0,00	1,090	1,140	0,320	62,78
2	1	28,1	1537	7,85	0,52	0,280	0,940	0,384	231,67
	2	28,6	1568	8,52	0,00	0,253	0,839	0,428	262,31
3	1	28,2	755	8,52	0,52	0,275	0,937	0,384	231,67
	2	28,5	739	8,24	0,00	0,348	0,339	0,761	172,69
4	3	28,6	748	8,60	0,13	0,386	0,419	0,867	189,41
	1	27,5	905	7,63	0,00	0,893	0,915	1,056	91,699
5	1	26,4	419	7,72	0,00	0,708	0,668	0,258	113,18
6	1	23,9	1034	7,64	0,42	0,301	1,229	0,263	189,04
7	1	32,6	217	7,42	0,30	0,843	2,076	0,618	118,00
8	1	31,3	371	7,32	0,00	0,697	0,635	0,489	99,225
9	1	32,3	1162	7,17	0,04	0,248	0,762	1,811	326,63

**Paso 1.** Cálculo de los valores máximos y mínimos de las variables.

	Temperatura (°C)	Conductividad ( $\mu\text{S cm}^{-1}$ )	pH	Nitrato ( $\mu\text{M}$ )	Nitrito ( $\mu\text{M}$ )	Amonio ( $\mu\text{M}$ )	Fosfato ( $\mu\text{M}$ )	Silicato ( $\mu\text{M}$ )
Máximo	32,6	1568	8,6	0,52	1,09	2,076	1,811	326,63
Mínimo	23,9	217	6,08	0	0,248	0,339	0,258	62,78

**Paso 2.** Estandarización a una escala de 0 a 1 de todas las variables.

A cada uno de los valores de las variables se aplica la siguiente fórmula:

$$VE = \frac{x - Min}{Max - Min}$$

donde *VE* es el valor estandarizado, *Max* y *Min* son los valores máximo y mínimo de cada variable, respectivamente, que se calcularon en el paso 1, y *x* es cada uno de los valores de cada variable. Los valores estandarizados se muestran en la Tabla I.2.

**Tabla I.2.** Valores estandarizados de los datos físico-químicos que se muestran en la Tabla I.1.

Laguna	Estación	Temperatura (°C)	Conductividad ( $\mu\text{S cm}^{-1}$ )	pH	Nitrato ( $\mu\text{M}$ )	Nitrito ( $\mu\text{M}$ )	Amonio ( $\mu\text{M}$ )	Fosfato ( $\mu\text{M}$ )	Silicato ( $\mu\text{M}$ )
1	1	0,598	0,206	0,000	0,000	1,000	0,461	0,040	0,000
2	1	0,483	0,977	0,702	1,000	0,038	0,346	0,081	0,640
	2	0,540	1,000	0,968	0,000	0,006	0,288	0,109	0,756
3	1	0,494	0,398	0,968	1,000	0,032	0,344	0,081	0,640
	2	0,529	0,386	0,857	0,000	0,119	0,000	0,324	0,417
4	3	0,540	0,393	1,000	0,250	0,164	0,046	0,392	0,480
	1	0,414	0,509	0,615	0,000	0,766	0,332	0,514	0,110
5	1	0,287	0,150	0,651	0,000	0,546	0,189	0,000	0,191
6	1	0,000	0,605	0,619	0,808	0,063	0,512	0,003	0,479
7	1	1,000	0,000	0,532	0,577	0,707	1,000	0,232	0,209
8	1	0,851	0,114	0,492	0,000	0,533	0,170	0,149	0,138
9	1	0,966	0,699	0,433	0,077	0,000	0,244	1,000	1,000

**CUADRO I.1.** (Continuación)**Estadística descriptiva con R**

Después de haber invocado al programa R y seleccionado el directorio, como se explica en el Apéndice II, abrimos el archivo **Cuadro I.1.R**, el cual es un archivo *script* que contiene las instrucciones para realizar la estadística descriptiva que se muestra en este ejemplo. Este archivo está en el CD que acompaña al libro en la carpeta Capítulo I. Para grabar los resultados, es mejor pasar el archivo **Cuadro I.1.R** a un directorio del disco duro y seleccionarlo como directorio de trabajo.

**Paso 1.** Para importar bases de datos desde Excel a R se recomienda guardar el archivo de Excel en formato ".csv" y luego leer desde R utilizando la función `read.csv` o `read.csv2` según el tipo de formato ".csv" usado por Excel para guardar los datos. Las instrucciones para realizar este tipo de importación de datos están en el Apéndice II en el apartado AII.4.2. En este ejemplo usaremos el archivo **Cuadro I.1.V.csv**, en el cual se encuentran los datos mostrados anteriormente en la Tabla I.1.

**Paso 2.** La siguiente ventana muestra las instrucciones del archivo *script* **Cuadro I.1.R**. Lo único que hay que hacer es cambiar la sección modificable por el usuario, para poder hacer la estandarización de cualquier tipo de archivo de datos.

```
# Cuadro I.1
# Estandarización de datos

#####
# Sección modificable por el usuario
#####
# Lectura de la base de datos
datos<-read.csv2("Cuadro I.1.V.csv",encoding="latin1")

# Selección de las variables de interés
varInteres<-c("Temperatura","Conductividad","pH","Nitrato",
              "Nitrito","Amonio","Fosfato","Silicato")

# Selección de las variables de identificación

varID<-c("Laguna","Estación")

# Nombre del archivo de salida con los resultados de la
# estandarización
nomSalida<-"Salida Cuadro I.1.V.csv"
```

**Paso 3.** En primer lugar hay que importar los datos y, por tanto, especificar el nombre del archivo de entrada, en este caso **Cuadro I.1.V.csv**.

```
datos<-read.csv2("Cuadro I.1.V.csv",encoding="latin1")
```

**CUADRO I.1.** (Continuación)

**Paso 4.** Una vez importados los datos hay que definir las variables que se quieren estandarizar.

```
varInteres<-c("Temperatura","Conductividad","pH","Nitrato",
"Nitrito","Amonio","Fosfato","Silicato")
```

**Paso 5.** A continuación se seleccionan las variables que se utilizan como códigos de identificación de los casos.

```
varID<-c("Laguna","Estación")
```

**Paso 6.** El último paso que tiene que hacer el usuario es definir el nombre del archivo donde se grabarán los valores estandarizados.

```
nomSalida<-"Salida Cuadro I.1.V.csv"
```

**Paso 7.** Para ejecutar el archivo *script* es posible posicionarse en cada línea del archivo *script* y pulsar *ctrl-R*, o seleccionar varias líneas o todo el archivo con *ctrl-A* y luego pulsar *ctrl-R*, para ejecutar varias líneas o todo el archivo a la vez. Los resultados se muestran en la siguiente ventana.

```
> resultado<-rbind(maximos,minimos)
>
> # Estandarización a una escala de 0 a 1 de todas las variables dada
> resultado2<-cbind(baseID,t((t(baseVar)-minimos)/(maximos-minimos)))
> if(length(varID)==1) colnames(resultado2)[1]<-varID
>
> #####
> # Sección que muestra los resultados
> #####
>
> # Muestra los resultados de mínimos y máximos
> resultado
      Temperatura Conductividad   pH Nitrato Nitrito Amonio Fosfato Silicato
maximos    32.6         1568 8.60    0.52    1.09    2.08    1.81    326.63
minimos     23.9         217 6.08    0.00    0.25    0.34    0.26    62.78
>
> # Muestra los resultados de la estandarización en la escala 0 a 1.
> resultado2
      Laguna Estación Temperatura Conductividad   pH  Nitrato  Nitrito
1         1         1    0.5977011    0.2057735 0.0000000 0.0000000 1.0000000
2         2         1    0.4827586    0.9770540 0.7023810 1.0000000 0.03571429
3         2         2    0.5402299    1.0000000 0.9682540 0.0000000 0.0000000
4         3         1    0.4942529    0.3982235 0.9682540 1.0000000 0.03571429
5         3         2    0.5287356    0.3863805 0.8571429 0.0000000 0.11904762
6         3         3    0.5402299    0.3930422 1.0000000 0.2500000 0.16666667
7         4         1    0.4137931    0.5092524 0.6150794 0.0000000 0.76190476
8         5         1    0.2873563    0.1495189 0.6507937 0.0000000 0.54761905
9         6         1    0.0000000    0.6047372 0.6190476 0.80769231 0.05952381
10        7         1    1.0000000    0.0000000 0.5317460 0.57692308 0.70238095
11        8         1    0.8505747    0.1139896 0.4920635 0.0000000 0.53571429
12        9         1    0.9655172    0.6994819 0.4325397 0.07692308 0.0000000
      Amonio  Fosfato  Silicato
1 0.45977011 0.03870968 0.0000000
2 0.34482759 0.07741935 0.6400985
3 0.28735632 0.10967742 0.7562251
4 0.34482759 0.07741935 0.6400985
5 0.00000000 0.32258065 0.4165624
6 0.04597701 0.39354839 0.4799318
7 0.33333333 0.51612903 0.1096077
8 0.18965517 0.00000000 0.1910176
9 0.51149425 0.00000000 0.4785295
10 1.00000000 0.23225806 0.2092856
11 0.17241379 0.14838710 0.1381467
12 0.24137931 1.00000000 1.0000000
```