

1. Conceptos básicos de matrices en R

Introducción.....	1
1.0 Resumen estadístico de datos multivariados.	2
1.1 Correlaciones.....	4
1.2 Distancias.....	5
1.3 Valores y Vectores propios.....	6
1.4 Ejercicios propuestos	7

Introducción

Este capítulo no pretende hacer una revisión detallada del álgebra lineal, por lo que se requiere de la revisión de otros documentos para aquellos que no estén familiarizados con esta temática. El énfasis de este capítulo se orienta a la aplicación del programa R, para la construcción de algunas matrices de uso común en el Análisis Multivariado (A.M.), las operaciones matriciales deberán ser consultadas para dar solución a las interrogantes que se consignan en el cuestionario de este capítulo.

La mayoría de los conjuntos de datos multivariantes tienen una forma común, y consisten en una matriz de datos, las filas de los que contienen corresponden a las unidades de la muestra, y las columnas se refieren a las variables medidas. Las letras mayúsculas suelen representar a las matrices (cuadradas o rectangulares) y en este documento, los valores relacionados en estas matrices corresponderán a números reales ubicados en medio de paréntesis. Simbólicamente un conjunto de datos multivariantes puede ser representados por una matriz A, dada por la siguiente notación:

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1q} \\ \vdots & & \ddots & \\ a_{n1} & a_{n2} & \cdots & a_{nq} \end{bmatrix}$$

Donde, n es el número de observaciones de la muestra, q es el número de variables medidas en la muestra y a_{nq} indica el valor de la variable q -ésima de la unidad n -ésima.

Las unidades en un conjunto de datos multivariantes con frecuencia son observaciones individuales, por ejemplo, cada uno de los peces en una investigación ictiológica, o de plantas en un estudio de investigación de vegetación. También pueden ser bahías, localidades, municipios o departamentos, por citar sólo cuatro posibilidades de unidades. En todos los casos las unidades son a menudo conocidas simplemente como "observaciones", un término que generalmente adoptado en este texto.

Un ejemplo de una matriz de datos multivariados es dado en la tabla 1. Aquí $n = 10$, $q = 7$ y por ejemplo $x_{44} = 1$. En esta base de datos las variables son de diferente naturaleza. Cuatro niveles de medición son usualmente distinguidos en la tipología de variables:

1. Nominal: Variables categóricas sin ordenar. Ejemplos de estas variables corresponden al tipo de ecosistema (tabla 1: E1, E2), sexo o color de los ojos de los individuos, presencia o ausencia de depredadores, entre otros.

2. Ordinal: Se otorga un orden, sin implicar la misma distancia de los puntos en una escala prediseñada por el usuario. Se incluyen ejemplos como las capas del suelo, niveles de contaminación de un ecosistema, o niveles de temperatura en diferentes comunidades.
3. Intervalo: Las distancias entre los puntos de una escala son iguales, pero la posición del cero en la escala es arbitraria. El ejemplo clásico es la escala de temperatura en grados Celsius o Fahrenheit.
4. Razón: Se analiza la magnitud relativa de las mediciones, así como las diferencias entre estas. La posición del cero es fija. Como ejemplos se incluye a la medida absoluta de la variable temperatura en grados Celsius, la abundancia, el peso o la edad de individuos en una o varias poblaciones biológicas.

La información cualitativa de la tabla xxx1, puede ser presentada en términos de códigos numéricos (requeridos por diferentes programas estadísticos), como el ecosistema E1 = 1 y el ecosistema E2 = 2, o las localidades S1, ..., S10, como 1, ..., 10. El análisis de datos *nominales* se puede limitar al resumen de resultados estadísticos como la moda. En el análisis de datos ordinales, la media y la desviación estándar no son recomendables.

Tabla 1. Datos hipotéticos de abundancia de órdenes de insectos en diferentes localidades ubicadas en dos tipos de ecosistemas.

Localidades	Ecosistema	Coleop	Dipter	Himeno	Hemip	Lepid
S1	E1	3	4	4	6	1
S2	E1	5	1	1	7	3
S3	E1	6	2	0	2	6
S4	E1	1	1	1	0	3
S5	E1	4	7	3	6	2
S6	E2	2	2	5	1	0
S7	E2	0	4	1	1	1
S8	E2	0	6	4	3	5
S9	E2	7	6	5	1	4
S10	E2	2	1	4	3	1

1.0 Resumen estadístico de datos multivariados.

El objeto de cualquier análisis multivariado consiste en iniciar con un resumen de cada una de las variables por separado, para poder entender su comportamiento general en el grupo de datos. Para este propósito se suele utilizar medidas de tendencia central y de variación (asumiendo que las variables son continuas o discretas). Adicionalmente se resume la información del conjunto de variables mediante sus correlaciones o covarianzas.

1.0.1 Medias. Para las q variables, el vector de medias de una población es usualmente representado por $\mu' = [\mu_1, \mu_2, \dots, \mu_q]$, donde: $\mu_i = E(x_i)$

μ_i es la media de la población (o el valor esperado, denotado por el operador E) de la i -ésima variable. Un estimador de μ' , se basa en n , de las observaciones con q -dimensiones. De esta manera $\bar{x}' = [\bar{x}_1, \bar{x}_2, \dots, \bar{x}_q]$, donde \bar{x}' es la media de la variable x_i . Para ilustrar el cálculo de un vector, tomaremos las columnas de la tabla xxx1 (archivo: Insectos.csv), en donde se exponen las abundancias de 5 órdenes de insectos acuáticos de una muestra con 10 localidades. El vector de medias se puede extraer directamente de R, mediante la siguiente función:

```
insecto<-read.csv2("Insectos.csv", row.names=1)
insecto <- (insecto [,c(2:6)])
mean(insecto)
```

Los valores resultantes son:

<i>Coleop</i>	<i>Dipter</i>	<i>Himeno</i>	<i>Hemip</i>	<i>Lepid</i>
3,0	3,4	2,8	3,0	2,6

1.0.2 Varianzas. El vector de la varianza de una población puede representarse por $\sigma' = [\sigma_1^2, \sigma_2^2, \dots, \sigma_q^2]$, donde: $\sigma_i^2 = E(x_i - \mu_i)^2$

Un estimador de σ' basado en n , de las observaciones con q -dimensiones es $s_i^2 s' = [s_1^2, s_2^2, \dots, s_q^2]$, donde s_i^2 es la varianza muestral de x_i . El cálculo de las varianzas de la tabla xxx1, es realizado en el programa R mediante la función *sd*, de la siguiente forma:

$sd(insecto)^2$

Los valores resultantes de la varianza son:

<i>Coleop</i>	<i>Dipter</i>	<i>Himeno</i>	<i>Hemip</i>	<i>Lepid</i>
6,00000	5,37778	3,51111	6,22222	3,82222

Los valores resultantes de la desviación estándar son:

$sd(insecto)$

<i>Coleop</i>	<i>Dipter</i>	<i>Himeno</i>	<i>Hemip</i>	<i>Lepid</i>
2,44949	2,31900	1,87380	2,49444	1,95505

1.0.3 Covarianzas. La covarianza para dos variables: x_i y x_j , es definida por la siguiente función $\Sigma \sigma_i^2 Cov(x_i, x_j) = (x_i - \mu_i)(x_j - \mu_j)$, cuando $i = j$, la covarianza de la variable es su varianza. La covarianza de x_i y x_j usualmente se denota por σ_{ij} (la varianza de x_i se denota como σ_{ii}). La matriz de covarianza es simétrica y de dimensión $q \times q$, se denota como Σ , donde:

$$\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1q} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2q} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{q1} & \sigma_{q2} & \cdots & \sigma_{qq} \end{pmatrix}$$

Nótese que $\sigma_{ij} = \sigma_{ji}$. Esta matriz es conocida como varianza-covarianza o simplemente covarianza. La matriz Σ es estimada por la matriz S de la siguiente manera:

$$S = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_j - \bar{x})'$$

Donde $x_i = [x_{i1}, x_{i2}, \dots, x_{iq}]$ es el vector de observaciones para el i-ésimo objeto. La diagonal principal de S contiene a la varianza de cada variable. La matriz de covarianza se obtiene mediante la función *var* de R.

var(insecto)

Se genera la matriz de varianzas (en la diagonal principal) y de covarianzas (elementos por fuera de la diagonal).

	<i>Coleop</i>	<i>Dipter</i>	<i>Himeno</i>	<i>Hemip</i>	<i>Lepid</i>
<i>Coleop</i>	6,000	0,444	-0,111	1,556	1,889
<i>Dipter</i>	0,444	5,378	1,756	0,889	0,733
<i>Himeno</i>	-0,111	1,756	3,511	0,000	-1,311
<i>Hemip</i>	1,556	0,889	0,000	6,222	-0,444
<i>Lepid</i>	1,889	0,733	-1,311	-0,444	3,822

1.0.4 Resumen estadístico. Una forma abreviada de reportar los resultados estadísticos de una base de datos se obtiene en el programa R mediante la función *summary*, de la siguiente manera:

summary(insecto)

	<i>Coleop</i>		<i>Dipter</i>		<i>Himeno</i>		<i>Hemip</i>		<i>Lepid</i>
<i>Min.</i>	:0.00	<i>Min.</i>	:1.00	<i>Min.</i>	:0.0	<i>Min.</i>	:0.00	<i>Min.</i>	:0.00
<i>1st Qu.</i>	:1.25	<i>1st Qu.</i>	:1.25	<i>1st Qu.</i>	:1.00	<i>1st Qu.</i>	:1.00	<i>1st Qu.</i>	:1.00
<i>Median</i>	:2.50	<i>Median</i>	:3.00	<i>Median</i>	:3.5	<i>Median</i>	:2.50	<i>Median</i>	:2.50
<i>Mean</i>	:3.00	<i>Mean</i>	:3.40	<i>Mean</i>	:2.8	<i>Mean</i>	:3.00	<i>Mean</i>	:2.60
<i>3rd Qu.</i>	:4.75	<i>3rd Qu.</i>	:5.50	<i>3rd Qu.</i>	:4.0	<i>3rd Qu.</i>	:5.25	<i>3rd Qu.</i>	:3.75
<i>Max.</i>	:7.00	<i>Max.</i>	:7.00	<i>Max.</i>	:5.0	<i>Max.</i>	:7.00	<i>Max.</i>	:6.00

Donde para cada variable, Min: mínimo valor, 1st Qu.: primer cuartil, Median: mediana, Mean: media o valor promedio, 3rd Qu.: tercer cuartil, Max.: valor máximo.

1.1 Correlaciones.

Es una forma de estandarizar a las covarianzas en variables que presentan unidades disímiles (ej. Variables físico-químicas). Esta estandarización se realiza dividiendo por el producto de las desviaciones estándar de dos variables, mediante un coeficiente de correlación de Pearson ρ_{ij} , donde:

$$\rho_{ij} = \frac{\sigma_{ij}}{\sqrt{\sigma_{ii}}\sqrt{\sigma_{jj}}}$$

El coeficiente de correlación ρ , presenta sus límites entre -1 y +1, dando una medida del nivel de asociación o correspondencia lineal entre dos variables x_i y x_j . Cuando la asociación es positiva, el coeficiente tiende a +1 y si la asociación es negativa se tiende a -1. Si no existe asociación entre las variables el coeficiente tiende a 0. La matriz de correlación se obtiene mediante la siguiente función:

cor(insecto)

De esta manera se genera es la siguiente matriz de correlación:

	Coleop	Dipter	Himeno	Hemip	Lepid
Coleop	1,00	0,078	-0,024	0,255	0,394
Dipter	0,078	1,00	0,404	0,154	0,162
Himeno	-0,024	0,404	1,00	0,000	-0,358
Hemip	0,255	0,154	0,000	1,00	-0,091
Lepid	0,394	0,162	-0,358	-0,091	1,00

1.2 Distancias.

Este concepto es detallado en el capítulo de clasificación de este libro, pues puede ser categorizado en múltiples tipos de distancias (Legendre y Legendre 1998). Una visión general demuestra que las distancias entre observaciones son atributos de mucha importancia para las técnicas multivariadas. La distancia de mayor uso es la Euclídea, la cual se deriva del Teorema de Pitágoras, en donde cada par observaciones corresponden a dos puntos del plano cartesiano y su distancia se determina mediante la ecuación de triángulos (distancia métrica que cumple con la desigualdad triangular). De esta manera la distancia euclídea para dos vectores (dos filas i e j , representadas por dos observaciones), se calcula mediante la siguiente función:

$$d_{ij} = \sqrt{\sum_{k=1}^q (x_{ik} - x_{jk})^2}$$

Para el cálculo de la Distancia Euclídea se utiliza la función *dist* de la siguiente manera:

dist <- dist(insecto)

Con base a la función anterior, se obtiene la siguiente matriz de distancias euclídeas la cual se resume en una matriz triangular, por ser simétrica.

	S1	S2	S3	S4	S5	S6	S7	S8	S9
S2	5,20								
S3	8,37	6,08							
S4	7,87	8,06	6,32						
S5	3,46	6,56	8,37	9,27					
S6	5,66	8,43	8,83	5,29	7,87				
S7	6,56	8,60	8,19	3,87	7,42	5,00			
S8	6,16	8,89	8,37	6,93	6,00	7,07	5,74		
S9	7,42	9,06	6,86	8,89	6,56	7,55	8,83	7,42	
S10	4,36	6,16	7,68	4,80	7,14	2,65	5,10	6,71	8,00

El cálculo de la matriz de distancias euclídeas ignora las diferencias en las unidades o escalas de las variables (ej. En el caso de variables fisicoquímicas), por lo que se sugiere hacer una estandarización, previa al cálculo de las distancias. Una buena opción para estandarizar las variables consiste en utilizar la función *scale* de R, la cual toma cada dato y le resta su media y le divide su desviación estándar. Luego se calcula la distancia euclídea sobre los datos estandarizados, de la siguiente manera:

```
insecto.estand <- scale (insecto)
dist(insecto.estand)
```

Con base a las funciones anteriores, se obtiene la siguiente matriz estandarizada de distancias euclídeas triangular y simétrica.

	S1	S2	S3	S4	S5	S6	S7	S8	S9
S2	2,47								
S3	3,99	2,65							
S4	3,43	3,25	2,76						
S5	1,54	2,90	3,83	3,92					
S6	2,34	3,79	4,40	2,72	3,39				
S7	2,84	3,56	3,71	1,75	3,12	2,50			
S8	2,81	3,87	3,74	3,14	2,63	3,33	2,85		
S9	3,17	3,99	3,39	3,95	2,81	3,37	3,98	3,06	
S10	1,81	2,77	3,76	2,29	3,06	1,17	2,36	3,08	3,48

1.3 Valores y Vectores propios.

También conocido como autovalores - λ_i (eigenvalues) y autovectores - μ_i (eigenvectors). Los valores propios de una matriz cuadrada A, son tales que para μ_i cumplen la siguiente propiedad:

$$A \mu_i = \lambda_i \mu_i$$

De esta manera se observa que los valores propios - λ_i son escalares que mediante un valor único resumen la información contenida en una matriz A, siempre y cuando cumplan con la igualdad de la función anterior. Para el cálculo de λ_i y μ_i se debe realizar el siguiente determinante, que hace parte de la ecuación característica:

$$|A - \lambda_i I| = 0$$

Donde I corresponde a la matriz identidad de A. de esta manera se obtiene la solución a la ecuación característica de autovalores y autovectores de A. revisemos el siguiente ejemplo, en donde tenemos una matriz cuadrada A de dimensión 2 x 2.

$$A = \begin{pmatrix} 1 & 2 \\ -1 & 4 \end{pmatrix}$$

El cálculo de λ_i y μ_i se realiza en R mediante la función *eigen*, luego se hace una descomposición mediante los siguientes comandos:

```
A <- matrix(c(1,2,-1,4),2,2,byrow=T)
A
```

Se obtiene la matriz cuadrada A:

	[,1]	[,2]
[1,]	1	2
[2,]	-1	4

Con el comando “*eigen*” se obtienen los autovectores (dos vectores) y autovalores (dos escalares):

eigen(A)

Los autovalores obtenidos son los siguientes:

<i>\$values</i>
[1] 3 2

Los autovectores obtenidos son los siguientes:

	[,1]	[,2]
[1,]	-0,71	-0,89
[2,]	-0,71	-0,45

Para probar la igualdad ($A \mu_i = \lambda_i \mu_i$) primero se extraen los autovalores y los autovectores mediante las dos funciones siguientes:

```
autovalor <- eigen(A)$values[1]
autovector <- eigen(A)$vectors[,1]
```

Posteriormente se prueba la igualdad (==) con la siguiente operación:

```
A %*% autovector == autovalor * autovector
```

R reporta que la igualdad es verdadera, mediante el siguiente vector:

	[,1]
[1,]	TRUE
[2,]	TRUE

Similar a lo realizado con la matriz A de 2 x 2, se puede obtener los valores y vectores propios de bases de datos más complejas siempre y cuando sean cuadradas. En el caso de la matriz de distancias euclídeas (dist) obtenidas de la tabla xxx1, esta operación puede ser realizada, debido a que es una matriz cuadrada y simétrica.

1.4 Ejercicios propuestos

Nota: Se requiere contar con la versión completa del libro *Análisis de Datos Ambientales y Ecológicos - ANDEA*