

# Análisis de Componentes Principales (Principal Component Analysis, PCA) y t-SNE

Joaquín Amat Rodrigo [j.amatrodrigo@gmail.com](mailto:j.amatrodrigo@gmail.com)

Junio, 2017

## Tabla de contenidos

Introducción.....	2
Álgebra Matricial.....	2
Interpretación geométrica de las componentes principales .....	4
Cálculo de las componentes principales.....	6
Escalado de las variables .....	7
Reproducibilidad de las componentes .....	7
Influencia de outliers.....	7
Proporción de varianza explicada .....	8
Número óptimo de componentes principales.....	9
Ejemplo cálculo <i>eigenvectors</i> y <i>eigenvalues</i> .....	9
Ejemplo cálculo directo de PCA con R.....	13
Ejemplo PCA aplicado a genómica .....	19
PCR: PCA aplicado a regresión lineal.....	25
Ejemplo.....	25
PLS: PCA aplicado a regresión lineal.....	32
Ejemplo.....	33
t-SNE .....	35
Idea intuitiva .....	35
Algoritmo de t-SNE.....	35
Interpretación de los resultados de un t-SNE.....	36
Limitaciones y desventajas .....	37
Ejemplo con tsne .....	38
Ejemplo con Rtsne.....	39
Otros métodos de reducción de dimensionalidad.....	43
Bibliografía.....	43

Formato PDF: <https://github.com/JoaquinAmatRodrigo/Estadistica-con-R>

## Introducción

*Principal Component Analysis (PCA)* es un método estadístico que permite simplificar la complejidad de espacios muestrales con muchas dimensiones a la vez que conserva su información. Supóngase que existe una muestra con  $n$  individuos cada uno con  $p$  variables ( $X_1, X_2, \dots, X_p$ ), es decir, el espacio muestral tiene  $p$  dimensiones. PCA permite encontrar un número de factores subyacentes ( $z < p$ ) que explican aproximadamente lo mismo que las  $p$  variables originales. Donde antes se necesitaban  $p$  valores para caracterizar a cada individuo, ahora bastan  $z$  valores. Cada una de estas  $z$  nuevas variables recibe el nombre de componente principal.

*Principal Component Analysis* pertenece a la familia de técnicas conocida como *unsupervised learning*. Los métodos de *supervised learning* descritos en capítulos anteriores tienen el objetivo de predecir una variable respuesta  $Y$  a partir de una serie de predictores. Para ello, se dispone de  $p$  características ( $X_1, X_2, \dots, X_p$ ) y de la variable respuesta  $Y$  medidas en  $n$  observaciones. En el caso de *unsupervised learning*, la variable respuesta  $Y$  no se tiene en cuenta ya que el objetivo no es predecir  $Y$  sino extraer información empleando los predictores, por ejemplo, para identificar subgrupos. El principal problema al que se enfrentan los métodos de *unsupervised learning* es la dificultad para validar los resultados dado que no se dispone de una variable respuesta que permita contrastarlos.

El método de PCA permite por lo tanto "condensar" la información aportada por múltiples variables en solo unas pocas componentes. Esto lo convierte en un método muy útil de aplicar previa utilización de otras técnicas estadísticas tales como *regresión*, *clustering*... Aun así no hay que olvidar que sigue siendo necesario disponer del valor de las variables originales para calcular las componentes.

## Álgebra Matricial

En esta sección se describen dos de los conceptos matemáticos que se aplican en el PCA: *eigenvectors* y *eigenvalues*. Se trata simplemente de una descripción intuitiva con la única finalidad de facilitar el entendimiento del cálculo de componentes principales.

## Eigenvectors

Los *eigenvectors* son un caso particular de multiplicación entre una matriz y un vector. Obsérvese la siguiente multiplicación:

$$\begin{pmatrix} 2 & 3 \\ 2 & 1 \end{pmatrix} \times \begin{pmatrix} 3 \\ 2 \end{pmatrix} = \begin{pmatrix} 12 \\ 8 \end{pmatrix} = 4 \times \begin{pmatrix} 3 \\ 2 \end{pmatrix}$$

El vector resultante de la multiplicación es un múltiplo entero del vector original. Los *eigenvectors* de una matriz son todos aquellos vectores que, al multiplicarlos por dicha matriz, resultan en el mismo vector o en un múltiplo entero del mismo. Los *eigenvectors* tienen una serie de propiedades matemáticas específicas:

- Los *eigenvectors* solo existen para matrices cuadradas y no para todas. En el caso de que una matriz  $n \times n$  tenga *eigenvectors*, el número de ellos es  $n$ .
- Si se escala un *eigenvector* antes de multiplicarlo por la matriz, se obtiene un múltiplo del mismo *eigenvector*. Esto se debe a que si se escala un vector multiplicándolo por cierta cantidad, lo único que se consigue es cambiar su longitud pero la dirección es la misma.
- Todos los *eigenvectors* de una matriz son perpendiculares (ortogonales) entre ellos, independientemente de las dimensiones que tengan.

Dada la propiedad de que multiplicar un *eigenvector* solo cambia su longitud pero no su naturaleza de *eigenvector*, es frecuente escalarlos de tal forma que su longitud sea 1. De este modo se consigue que todos ellos estén estandarizados. A continuación se muestra un ejemplo:

El *eigenvector*  $\begin{pmatrix} 3 \\ 2 \end{pmatrix}$  tiene una longitud de  $\sqrt{3^2 + 2^2} = \sqrt{13}$ . Si se divide cada dimensión entre la longitud del vector, se obtiene el *eigenvector* estandarizado con longitud 1.

$$\begin{pmatrix} 3 \\ 2 \end{pmatrix} \div \sqrt{13} = \begin{pmatrix} 3/\sqrt{13} \\ 2/\sqrt{13} \end{pmatrix}$$

## Eigenvalue

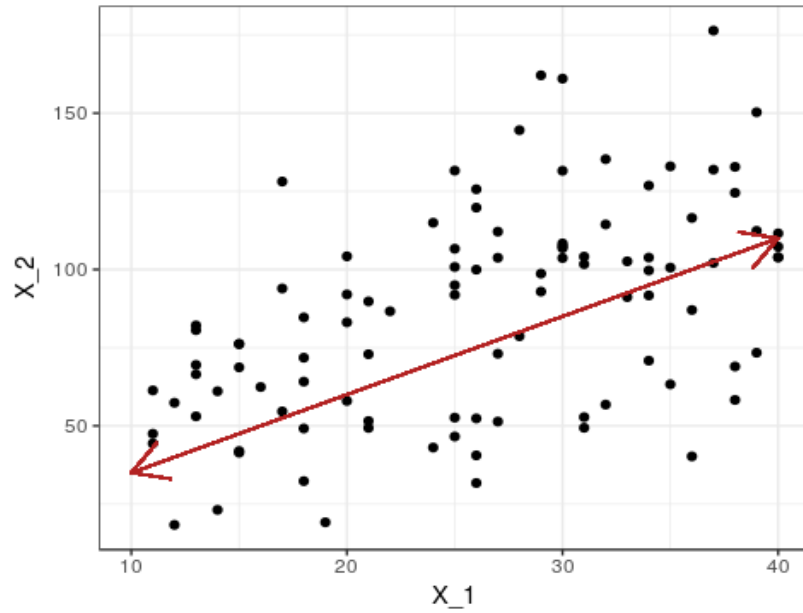
Cuando se multiplica una matriz por alguno de sus *eigenvectors* se obtiene un múltiplo del vector original, es decir, el resultado es ese mismo vector multiplicado por un número. Al valor por el que se multiplica el *eigenvector* resultante se le conoce como *eigenvalue*. A todo *eigenvector* le corresponde un *eigenvalue* y viceversa.

En el método PCA, cada una de las componentes se corresponde con un *eigenvector*, y el orden de componente se establece por orden decreciente de *eigenvalue*. Así pues, la primera componente es el *eigenvector* con el *eigenvalue* asociado más alto.

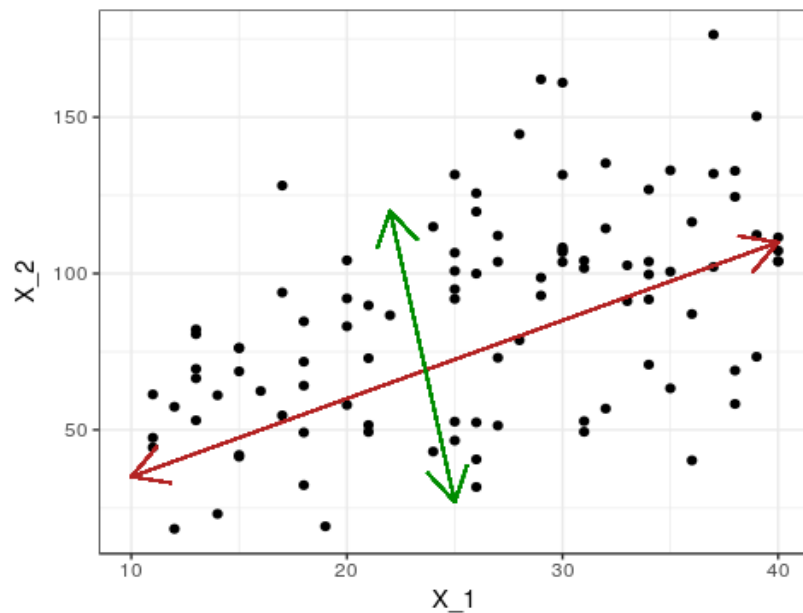
## Interpretación geométrica de las componentes principales

Una forma intuitiva de entender el proceso de PCA consiste en interpretar las componentes principales desde un punto de vista geométrico. Supóngase un conjunto de observaciones para las que se dispone de dos variables ( $X_1$ ,  $X_2$ ). El vector que define la primera componente principal ( $Z_1$ ) sigue la dirección en la que las observaciones varían más (línea roja). La proyección de cada observación sobre esa dirección equivale al valor de la primera componente para dicha observación (*principal component scores*,  $z_{i1}$ ).

```
library(ggplot2)
set.seed(435)
X_1 <- sample(x = 10:40, size = 100, replace = TRUE)
X_2 <- 2.5 * X_1 + 10
X_2 <- X_2 + rnorm(n = 100, mean = 10, sd = 30)
datos <- data.frame(X_1, X_2)
ggplot(data = datos, aes(x = X_1, y = X_2)) +
  geom_point() +
  geom_segment(aes(x = 10, y = 2.5 * 10 + 10, xend = 40, yend = 2.5 * 40 + 10),
               colour = "firebrick", arrow = arrow(ends = "both")) +
  theme_bw()
```



La segunda componente ( $Z_2$ ) sigue la segunda dirección en la que los datos muestran mayor varianza y que no está correlacionada con la primera componente. La condición de no correlación entre componentes principales equivale a decir que sus direcciones son perpendiculares/ortogonales.



## Cálculo de las componentes principales

Cada componente principal ( $Z_i$ ) se obtiene por combinación lineal de las variables originales. Se pueden entender como nuevas variables obtenidas al combinar de una determinada forma las variables originales. La primera componente principal de un grupo de variables ( $X_1, X_2, \dots, X_p$ ) es la combinación lineal normalizada de dichas variables que tiene mayor varianza:

$$Z_1 = \phi_{11}X_1 + \phi_{21}X_2 + \dots + \phi_{p1}X_p$$

Que la combinación lineal sea normalizada implica que:

$$\sum_{j=1}^p \phi_{j1}^2 = 1$$

Los términos  $\phi_{11}, \dots, \phi_{p1}$  reciben en el nombre de *loadings* y son los que definen a la componente.  $\phi_{11}$  es el *loading* de la variable  $X_1$  de la primera componente principal. Los *loadings* pueden interpretarse como el peso/importancia que tiene cada variable en cada componente y, por lo tanto, ayudan a conocer que tipo de información recoge cada una de las componentes.

Dado un set de datos  $X$  con  $n$  observaciones y  $p$  variables, el proceso a seguir para calcular la primera componente principal es:

- Centralización de las variables: se resta a cada valor la media de la variable a la que pertenece. Con esto se consigue que todas las variables tengan media cero.
- Se resuelve un problema de optimización para encontrar el valor de los *loadings* con los que se maximiza la varianza. Una forma de resolver esta optimización es mediante el cálculo de *eigenvector-eigenvalue* de la matriz de covarianzas.

Una vez calculada la primera componente ( $Z_1$ ) se calcula la segunda ( $Z_2$ ) repitiendo el mismo proceso, pero añadiendo la condición de que la combinación lineal no puede estar correlacionada con la primera componente. Esto equivale a decir que  $Z_1$  y  $Z_2$  tienen que ser perpendiculares. EL proceso se repite de forma iterativa hasta calcular todas las posibles componentes ( $\min(n-1, p)$ ) o hasta que se decida detener el proceso. El orden de importancia de las componentes viene dado por la magnitud del *eigenvalue* asociado a cada *eigenvector*.

## Escalado de las variables

El proceso de PCA identifica aquellas direcciones en las que la varianza es mayor. Como la varianza de una variable se mide en su misma escala elevada al cuadrado, si antes de calcular las componentes no se estandarizan todas las variables para que tengan media 0 y desviación estándar 1, aquellas variables cuya escala sea mayor dominarán al resto. De ahí que sea recomendable estandarizar siempre los datos.

## Reproducibilidad de las componentes

El proceso de PCA genera siempre las mismas componentes principales independientemente del software utilizado, es decir, el valor de los *loadings* resultantes es el mismo. La única diferencia que puede darse es que el signo de todos los *loadings* esté invertido. Esto es así porque el vector de *loadings* determina la dirección de la componente, y dicha dirección es la misma independientemente del signo (la componente sigue una línea que se extiende en ambas direcciones). Del mismo modo, el valor específico de las componentes obtenido para cada observación (*principal component scores*) es siempre el mismo, a excepción del signo.

## Influencia de outliers

Al trabajar con varianzas, el método PCA es altamente sensible a *outliers*, por lo que es altamente recomendable estudiar si los hay. La detección de valores atípicos con respecto a una determinada dimensión es algo relativamente sencillo de hacer mediante comprobaciones gráficas. Sin embargo, cuando se trata con múltiples dimensiones el proceso se complica. Por ejemplo, considérese un hombre que mide 2 metros y pesa 50 kg. Ninguno de los dos valores es atípico de forma individual, pero en conjunto se trataría de un caso muy excepcional. La distancia de *Mahalanobis* es una medida de distancia entre un punto y la media que se ajusta en función de la correlación entre dimensiones y que permite encontrar potenciales *outliers* en distribuciones multivariante.

## Proporción de varianza explicada

Una de las preguntas más frecuentes que surge tras realizar un PCA es: ¿Cuánta información presente en el set de datos original se pierde al proyectar las observaciones en un espacio de menor dimensión? o lo que es lo mismo ¿Cuánta información es capaz de capturar cada una de las componentes principales obtenidas? Para contestar a estas preguntas se recurre a la proporción de varianza explicada por cada componente principal.

Asumiendo que las variables se han normalizado para tener media cero, la varianza total presente en el set de datos se define como

$$\sum_{j=1}^p \text{Var}(X_j) = \sum_{j=1}^p \frac{1}{n} \sum_{i=1}^n x_{ij}^2$$

y la varianza explicada por la componente  $m$  es

$$\frac{1}{n} \sum_{i=1}^n z_{im}^2 = \frac{1}{n} \sum_{i=1}^n \left( \sum_{j=1}^p \phi_{jm} x_{ij} \right)^2$$

Por lo tanto, la proporción de varianza explicada por la componente  $m$  viene dada por el ratio

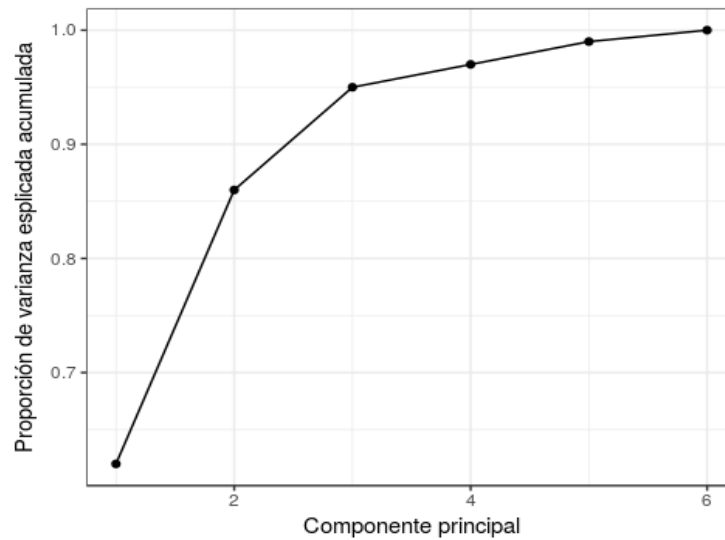
$$\frac{\sum_{i=1}^n \left( \sum_{j=1}^p \phi_{jm} x_{ij} \right)^2}{\sum_{j=1}^p \sum_{i=1}^n x_{ij}^2}$$

Tanto la proporción de varianza explicada como la proporción de varianza explicada acumulada son dos valores de gran utilidad a la hora de decidir el número de componentes principales a utilizar en los análisis posteriores. Si se calculan todas las componentes principales de un set de datos, entonces, aunque transformada, se está almacenando toda la información presente en los datos originales. El sumatorio de la proporción de varianza explicada acumulada de todas las componentes es siempre 1.



## Número óptimo de componentes principales

Por lo general, dada una matriz de datos de dimensiones  $n \times p$ , el número de componentes principales que se pueden calcular es como máximo de  $n-1$  o  $p$  (el menor de los dos valores es el limitante). Sin embargo, siendo el objetivo del PCA reducir la dimensionalidad, suelen ser de interés utilizar el número mínimo de componentes que resultan suficientes para explicar los datos. No existe una respuesta o método único que permita identificar cual es el número óptimo de componentes principales a utilizar. Una forma de proceder muy extendida consiste en evaluar la proporción de varianza explicada acumulada y seleccionar el número de componentes mínimo a partir del cual el incremento deja de ser sustancial.



## Ejemplo cálculo eigenvectors y eigenvalues

*El siguiente es un ejemplo que muestra cómo se pueden calcular las componentes principales de un set de datos mediante la identificación de eigenvectors y eigenvalues por el método de covarianza. Con el objetivo de ser un ejemplo sencillo, el set de datos empleado tiene dos dimensiones (variables). La mayoría de programas de análisis, entre ellos R, disponen de funciones que devuelven directamente el valor de las componentes principales (ver ejemplo siguiente).*

```
datos <- data.frame(X1 = c(2.5, 0.5, 2.2, 1.9, 3.1, 2.3, 2, 1, 1.5, 1.1),
                    X2 = c(2.4, 0.7, 2.9, 2.2, 3, 2.7, 1.6, 1.1, 1.6, 0.9))
datos
```

```
##      X1  X2
## 1  2.5 2.4
## 2  0.5 0.7
## 3  2.2 2.9
## 4  1.9 2.2
## 5  3.1 3.0
## 6  2.3 2.7
## 7  2.0 1.6
## 8  1.0 1.1
## 9  1.5 1.6
## 10 1.1 0.9
```

En primer lugar se resta a cada valor la media de la variable a la que pertenece. Con esto se consigue centralizar las variables y que su media sea 0.

```
datos_centrados <- datos
datos_centrados$X1 <- datos$X1 - mean(datos$X1)
datos_centrados$X2 <- datos$X2 - mean(datos$X2)
datos_centrados
```

```
##      X1  X2
## 1  0.69 0.49
## 2 -1.31 -1.21
## 3  0.39 0.99
## 4  0.09 0.29
## 5  1.29 1.09
## 6  0.49 0.79
## 7  0.19 -0.31
## 8 -0.81 -0.81
## 9 -0.31 -0.31
## 10 -0.71 -1.01
```

A continuación se calcula la matriz de correlaciones entre cada par de variables. Como en este ejemplo hay dos variables, el resultado es una matriz simétrica 2x2.

```
matriz_cov <- cov(datos_centrados)
matriz_cov
```

```
##      X1      X2
## X1 0.6165556 0.6154444
## X2 0.6154444 0.7165556
```

Dado que la matriz de covarianzas es cuadrada, se pueden obtener sus correspondientes *eigenvectors* y *eigenvalues*. La función `eigen()` calcula ambos y los almacena en una lista bajo el nombre de `vectors` y `values`. Los *eigenvalues* se devuelven en orden decreciente y los

*eigenvectors* (estandarizados) se ordenan de izquierda a derecha acorde a sus *eigenvalues* asociados.

```
eigen <- eigen(matriz_cov)
eigen$values
```

```
## [1] 1.2840277 0.0490834
```

```
eigen$vectors
```

```
##           [,1]      [,2]
## [1,] 0.6778734 -0.7351787
## [2,] 0.7351787  0.6778734
```

Los *eigenvectors* ordenados de mayor a menor *eigenvalues* se corresponden con las componentes principales.

Una vez obtenidos los *eigenvectors* (componentes principales) se calcula el valor que toma cada componente para cada observación en función de las variables originales (*principal component scores*). Para ello, simplemente se tienen que multiplicar los *eigenvectors* transpuestos por los datos originales centrados y también transpuestos.

```
t_eigenvectors <- t(eigen$vectors)
t_eigenvectors
```

```
##           [,1]      [,2]
## [1,] 0.6778734 0.7351787
## [2,] -0.7351787 0.6778734
```

```
t_datos_centrados <- t(datos_centrados)
t_datos_centrados
```

```
##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10]
## X1 0.69 -1.31 0.39 0.09 1.29 0.49 0.19 -0.81 -0.31 -0.71
## X2 0.49 -1.21 0.99 0.29 1.09 0.79 -0.31 -0.81 -0.31 -1.01
```

```
# Producto matricial
pc_scores <- t_eigenvectors %*% t_datos_centrados
rownames(pc_scores) <- c("PC1", "PC2")

# Se vuelve a transponer para que los datos estén en modo tabla
t(pc_scores)
```

```
##           PC1           PC2
## [1,]  0.82797019 -0.17511531
## [2,] -1.77758033  0.14285723
## [3,]  0.99219749  0.38437499
## [4,]  0.27421042  0.13041721
## [5,]  1.67580142 -0.20949846
## [6,]  0.91294910  0.17528244
## [7,] -0.09910944 -0.34982470
## [8,] -1.14457216  0.04641726
## [9,] -0.43804614  0.01776463
## [10,] -1.22382056 -0.16267529
```

En este ejemplo, al partir inicialmente de un espacio muestral con 2 dimensiones y haber calculado sus 2 componentes principales, no se ha reducido la dimensionalidad ni se ha perdido nada de información presente en los datos originales.

Cuando se calculan y conservan todos los *eigenvectors* (componentes principales), es posible recuperar de nuevo los valores iniciales. Solo es necesario invertir el proceso:

$$\text{datos originales} = (\text{eigenvectors}^{-1}) \times (\text{principal component scores}) + \text{medias originales}$$

Al estar incluidos todos los *eigenvectors*:

$$\text{eigenvectors}^{-1} = \text{eigenvectors}^T$$

```
datos_recuperados <- t(eigen$vectors %*% pc_scores)
datos_recuperados[, 1] <- datos_recuperados[, 1] + mean(datos$X1)
datos_recuperados[, 2] <- datos_recuperados[, 2] + mean(datos$X2)

datos_recuperados
```

```
##      [,1] [,2]
## [1,]  2.5  2.4
## [2,]  0.5  0.7
## [3,]  2.2  2.9
## [4,]  1.9  2.2
## [5,]  3.1  3.0
## [6,]  2.3  2.7
## [7,]  2.0  1.6
## [8,]  1.0  1.1
## [9,]  1.5  1.6
## [10,]  1.1  0.9
```

```
datos
```

```
##      X1  X2
## 1  2.5 2.4
## 2  0.5 0.7
## 3  2.2 2.9
## 4  1.9 2.2
## 5  3.1 3.0
## 6  2.3 2.7
## 7  2.0 1.6
## 8  1.0 1.1
## 9  1.5 1.6
## 10 1.1 0.9
```

En la mayoría de programas estadísticos, entre ellos `R`, existen funciones que calculan directamente las componentes principales y los *principal component scores* de cada observación sin tener que ir paso por paso.

## Ejemplo cálculo directo de PCA con R

El set de datos `USArrests` del paquete básico de `R` contiene el porcentaje de asaltos (Assault), asesinatos (Murder) y secuestros (Rape) por cada 100,000 habitantes para cada uno de los 50 estados de USA (1973). Además, también incluye el porcentaje de la población de cada estado que vive en zonas rurales (UrbanPop).

```
data("USArrests")
head(USArrests)
```

```
##      Murder Assault UrbanPop Rape
## Alabama    13.2    236      58 21.2
## Alaska     10.0    263      48 44.5
## Arizona      8.1    294      80 31.0
## Arkansas     8.8    190      50 19.5
## California   9.0    276      91 40.6
## Colorado     7.9    204      78 38.7
```

El promedio de los datos muestra que hay tres veces más secuestros que asesinatos y 8 veces más asaltos que secuestros.

```
apply(X = USArrests, MARGIN = 2, FUN = mean)
```

```
## Murder Assault UrbanPop Rape
## 7.788 170.760 65.540 21.232
```

La varianza es muy distinta entre las variables, en el caso de *Assault*, la varianza es varios órdenes de magnitud superior al resto.

```
apply(X = USArrests, MARGIN = 2, FUN = var)
```

```
## Murder Assault UrbanPop Rape
## 18.97047 6945.16571 209.51878 87.72916
```

Si no se estandarizan las variables para que tengan media cero y desviación estándar 1 antes de realizar el estudio PCA, la variable *Assault* dominará la mayoría de las componentes principales.

La función `prcomp()` es una de las múltiples funciones en R que realizan PCA. Por defecto, `prcomp()` centra las variables para que tengan media cero, pero si se quiere además que su desviación estándar sea de uno, hay que indicar `scale = TRUE`.

```
pca <- prcomp(USArrests, scale = TRUE)
names(pca)
```

```
## [1] "sdev" "rotation" "center" "scale" "x"
```

Los elementos `center` y `scale` almacenados en el objeto `pca` contienen la media y desviación típica de las variables previa estandarización (en la escala original).

```
pca$center
```

```
## Murder Assault UrbanPop Rape
## 7.788 170.760 65.540 21.232
```

```
pca$scale
```

```
## Murder Assault UrbanPop Rape
## 4.355510 83.337661 14.474763 9.366385
```

`rotation` contiene el valor de los *loadings*  $\phi$  para cada componente (*eigenvector*). El número máximo de componentes principales se corresponde con el mínimo( $n-1, p$ ), que en este caso es  $\min(49, 4) = 4$ .

```
pca$rotation
```

##		PC1	PC2	PC3	PC4
## Murder		-0.5358995	0.4181809	-0.3412327	0.64922780
## Assault		-0.5831836	0.1879856	-0.2681484	-0.74340748
## UrbanPop		-0.2781909	-0.8728062	-0.3780158	0.13387773
## Rape		-0.5434321	-0.1673186	0.8177779	0.08902432

Analizar con detalle el vector de *loadings* que forma cada componente puede ayudar a interpretar que tipo de información recoge cada una de ellas. Por ejemplo, la primera componente es el resultado de la siguiente combinación lineal de las variables originales:

$$PC1 = -0.5358995 \text{ Murder} - 0.5831836 \text{ Assault} - 0.2781909 \text{ UrbanPop} - 0.5434321 \text{ Rape}$$

Los pesos asignados en la primera componente a las variables *Assault*, *Murder* y *Rape* son aproximadamente iguales entre ellos y bastante superiores al asignado a *UrbanPop*, esto significa que la primera componente recoge mayoritariamente la información correspondiente a los delitos. En la segunda componente, es la variable *UrbanPop* la que tiene con diferencia mayor peso, por lo que se corresponde principalmente con el nivel de urbanización del estado. Si bien en este ejemplo la interpretación de las componentes es bastante clara, no en todos los casos ocurre lo mismo.

La función `prcomp()` calcula automáticamente el valor de las componentes principales para cada observación (*principal component scores*) multiplicando los datos por los vectores de *loadings*. El resultado se almacena en la matriz `x`.

```
head(pca$x)
```

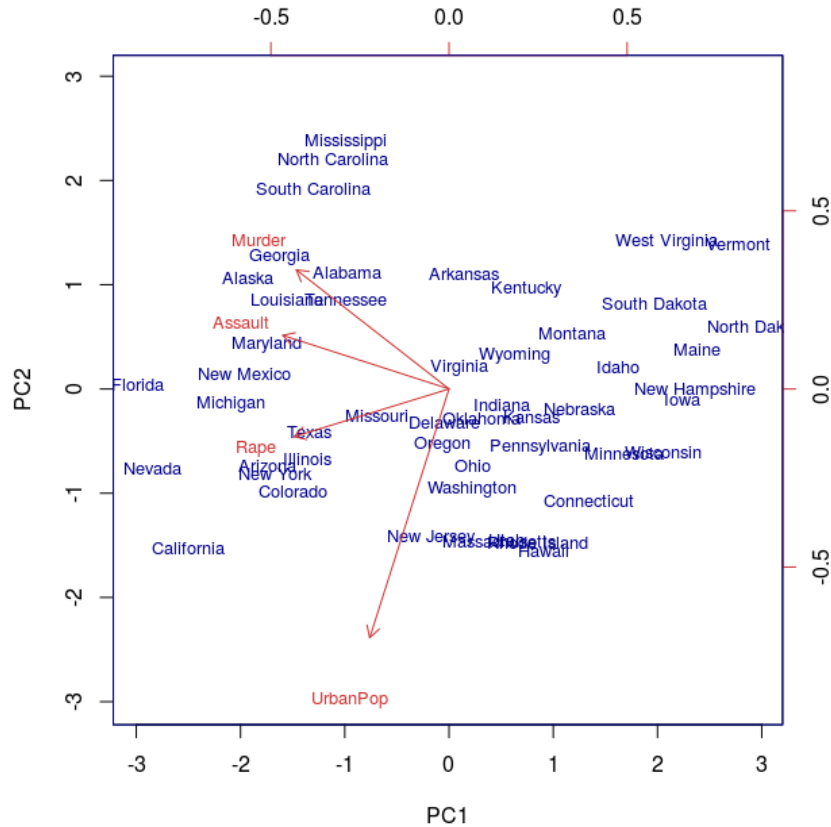
##		PC1	PC2	PC3	PC4
## Alabama		-0.9756604	1.1220012	-0.43980366	0.154696581
## Alaska		-1.9305379	1.0624269	2.01950027	-0.434175454
## Arizona		-1.7454429	-0.7384595	0.05423025	-0.826264240
## Arkansas		0.1399989	1.1085423	0.11342217	-0.180973554
## California		-2.4986128	-1.5274267	0.59254100	-0.338559240
## Colorado		-1.4993407	-0.9776297	1.08400162	0.001450164

```
dim(pca$x)
```

```
## [1] 50 4
```

Mediante la función `biplot()` se puede obtener una representación bidimensional de las dos primeras componentes. Es recomendable indicar el argumento `scale = 0` para que las flechas estén en la misma escala que las componentes.

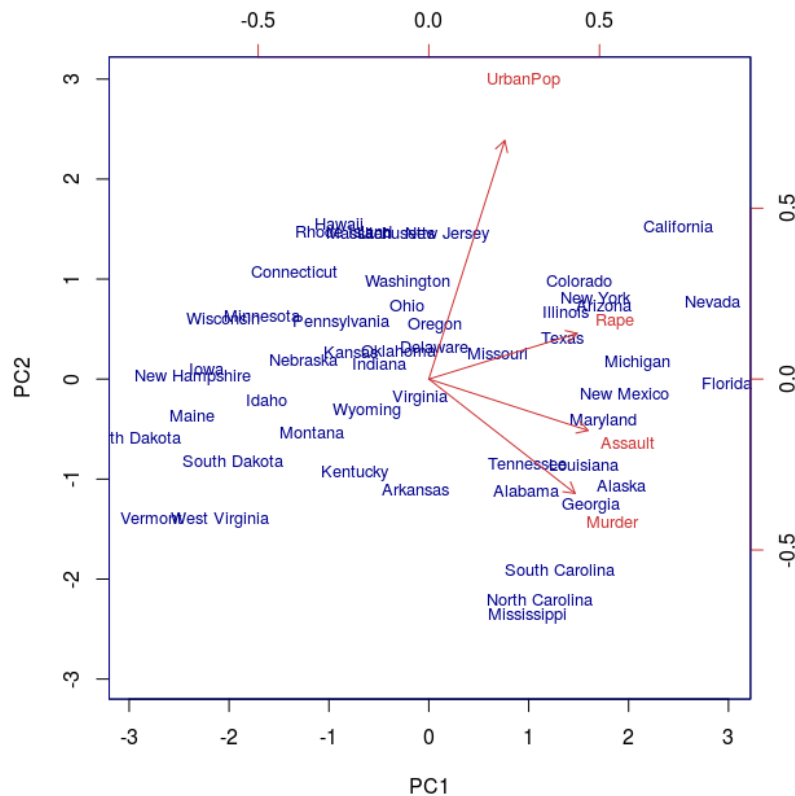
```
biplot(x = pca, scale = 0, cex = 0.8, col = c("blue4", "brown3"))
```



La imagen espejular, cuya interpretación es equivalente, se puede obtener invirtiendo el signo de los *loadings* y de los *principal component scores*.

```
pca$rotation <- -pca$rotation
pca$x <- -pca$x
biplot(x = pca, scale = 0, cex = 0.8, col = c("blue4", "brown3"))
```





Una vez calculadas las componentes principales, se puede conocer la varianza explicada por cada una de ellas, la proporción respecto al total y la proporción de varianza acumulada.

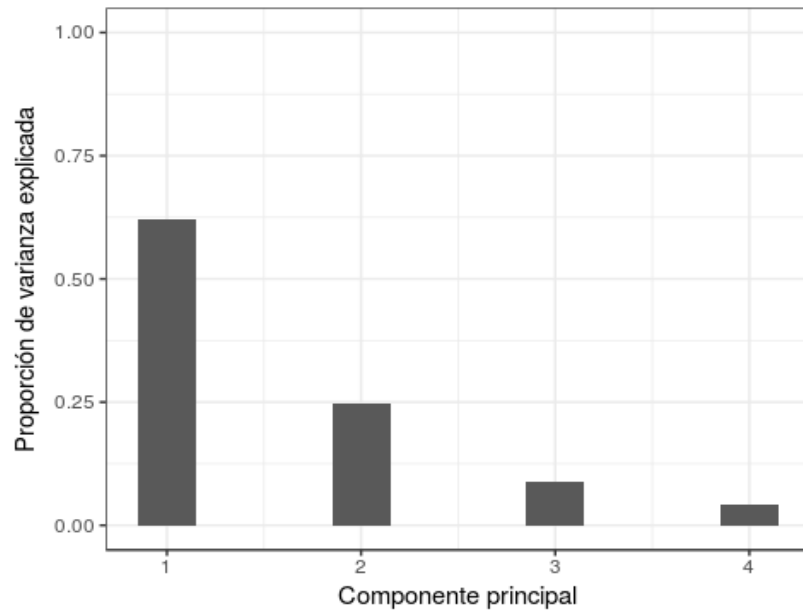
```
library(ggplot2)
pca$sdev^2
```

```
## [1] 2.4802416 0.9897652 0.3565632 0.1734301
```

```
prop_varianza <- pca$sdev^2/sum(pca$sdev^2)
prop_varianza
```

```
## [1] 0.62006039 0.24744129 0.08914080 0.04335752
```

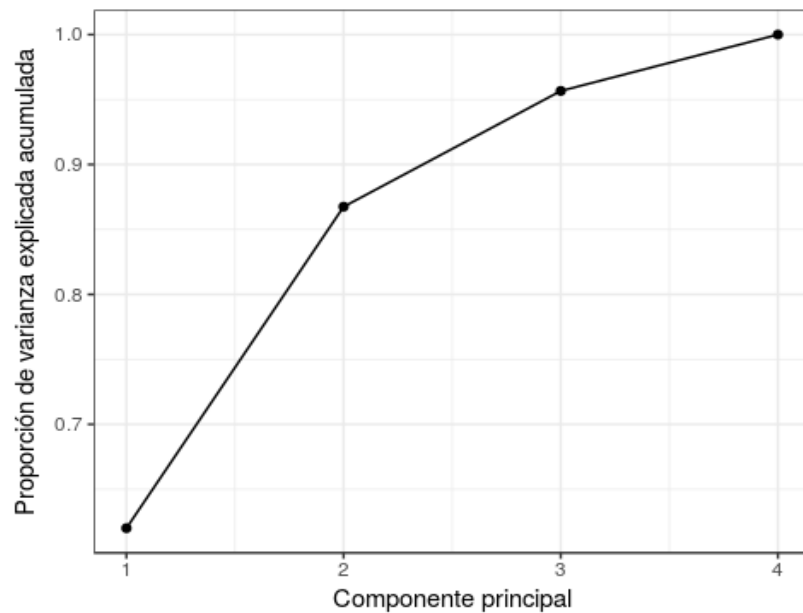
```
ggplot(data = data.frame(prop_varianza, pc = 1:4), aes(x = pc, y = prop_varianza)) +
  geom_col(width = 0.3) +
  scale_y_continuous(limits = c(0, 1)) +
  theme_bw() +
  labs(x = "Componente principal", y = "Proporción de varianza explicada")
```



```
prop_varianza_acum <- cumsum(prop_varianza)
prop_varianza_acum
```

```
## [1] 0.6200604 0.8675017 0.9566425 1.0000000
```

```
ggplot(data = data.frame(prop_varianza_acum, pc = 1:4),
       aes(x = pc, y = prop_varianza_acum, group = 1)) +
  geom_point() +
  geom_line() +
  theme_bw() +
  labs(x = "Componente principal", y = "Proporción de varianza explicada acumulada")
```



En este caso, la primera componente explica el 62% de la varianza observada en los datos y la segunda el 24.7%. Las dos últimas componentes no superan por separado el 1% de varianza explicada. Si se empleasen únicamente las dos primeras componentes se conseguiría explicar el 86.75% de la varianza observada.

## Ejemplo PCA aplicado a genómica

El siguiente es un ejemplo de como PCA puede emplearse para encontrar patrones cuando se dispone de muchos de predictores.

*Supóngase un equipo de investigación que se encarga de clasificar los tumores de pacientes en 3 subtipos. Dependiendo del subtipo, el paciente recibe una medicación diferente. El proceso de caracterización se hace mediante tinciones y observaciones al microscopio. Este proceso es muy laborioso y lento, lo que incrementa mucho el tiempo de respuesta de los médicos. Nuevos estudios apuntan a que cuantificando la expresión de un grupo de 9 genes se podría clasificar los tumores con una alta precisión. Se quiere determinar si tal patrón existe dentro de los datos.*

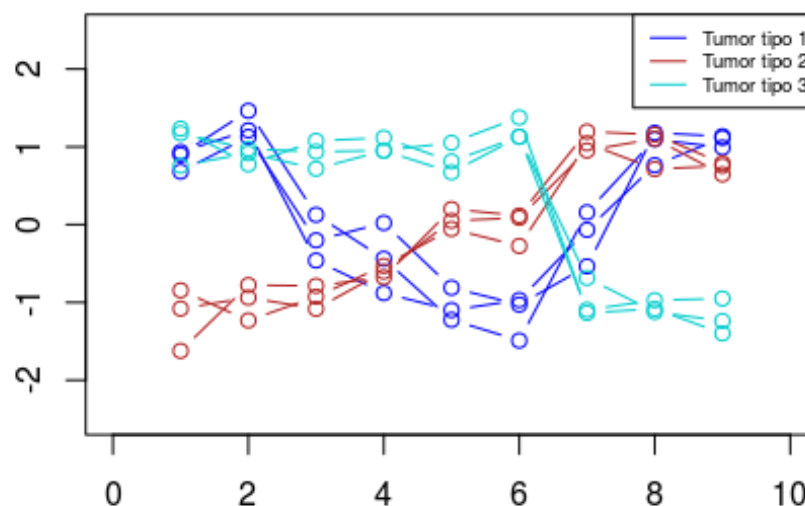
*Se trata de un ejemplo muy simplificado (en la realidad suele disponerse de varios miles de genes). El método de PCA puede combinarse con técnicas de clustering para crear agrupaciones automáticamente. En este caso, simplemente se muestra como PCA puede ayudar a la identificación de patrones que facilitarán posibles agrupamientos posteriores.*

```
#Se crean 3 perfiles modelo de expresión génica que pueden tener los tumores
tumor_1 <- c(1,1,0,-0.5,-1,-1,0,1,1)
tumor_2 <- c(-1,-1,-1,-0.5,0,0,1,1,1)
tumor_3 <- c(1,1,1,1,1,1,-1,-1,-1)

#Añadiendo ruido aleatorio se generan 3 muestras a partir de cada perfil modelo
set.seed(755)
tumor_a <- tumor_1 + rnorm(n = 9, mean = 0, sd = 0.2)
tumor_b <- tumor_1 + rnorm(n = 9, mean = 0, sd = 0.2)
tumor_c <- tumor_1 + rnorm(n = 9, mean = 0, sd = 0.2)
tumor_d <- tumor_2 + rnorm(n = 9, mean = 0, sd = 0.2)
tumor_e <- tumor_2 + rnorm(n = 9, mean = 0, sd = 0.2)
tumor_f <- tumor_2 + rnorm(n = 9, mean = 0, sd = 0.2)
tumor_g <- tumor_3 + rnorm(n = 9, mean = 0, sd = 0.2)
tumor_h <- tumor_3 + rnorm(n = 9, mean = 0, sd = 0.2)
tumor_i <- tumor_3 + rnorm(n = 9, mean = 0, sd = 0.2)
```

```
# Representación de los perfiles
plot(0,0, xlim = c(0,10), ylim = c(-2.5, 2.5), type = "n",
     main = "Perfil de expresión génica de 9 tumores")
lines(x = 1:9, y = tumor_a, type = "b", col = "blue")
lines(x = 1:9, y = tumor_b, type = "b", col = "blue")
lines(x = 1:9, y = tumor_c, type = "b", col = "blue")
lines(x = 1:9, y = tumor_d, type = "b", col = "firebrick")
lines(x = 1:9, y = tumor_e, type = "b", col = "firebrick")
lines(x = 1:9, y = tumor_f, type = "b", col = "firebrick")
lines(x = 1:9, y = tumor_g, type = "b", col = "cyan3")
lines(x = 1:9, y = tumor_h, type = "b", col = "cyan3")
lines(x = 1:9, y = tumor_i, type = "b", col = "cyan3")
legend("topright", legend = c("Tumor tipo 1", "Tumor tipo 2", "Tumor tipo 3"),
      col = c("blue", "firebrick", "cyan3"), lty = 1, cex = 0.6)
```

**Perfil de expresión génica de 9 tumores**

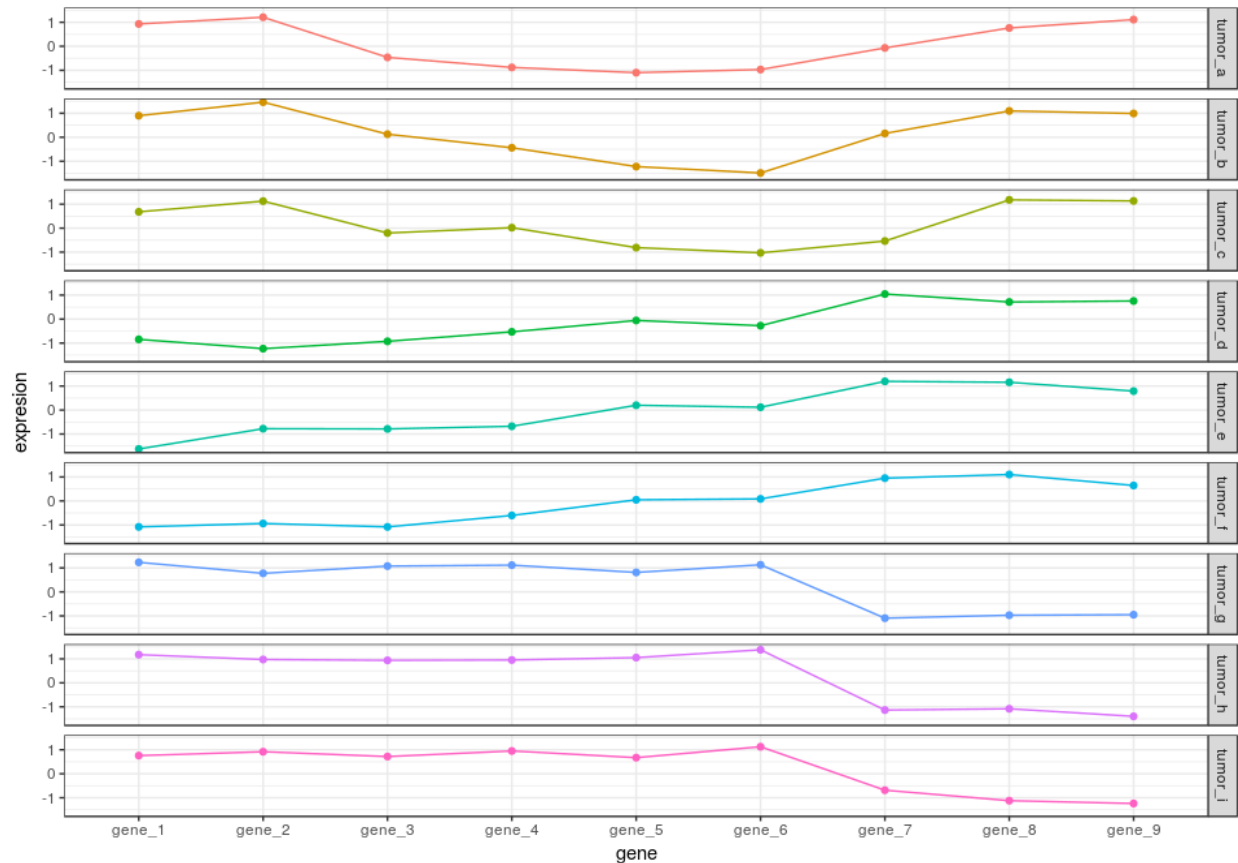


La representación del perfil de expresión muestra que, aunque con variabilidad, cada grupo tiene su propio perfil característico.

Se crea un *data.frame* con los datos para poder empelar `ggplot2`.

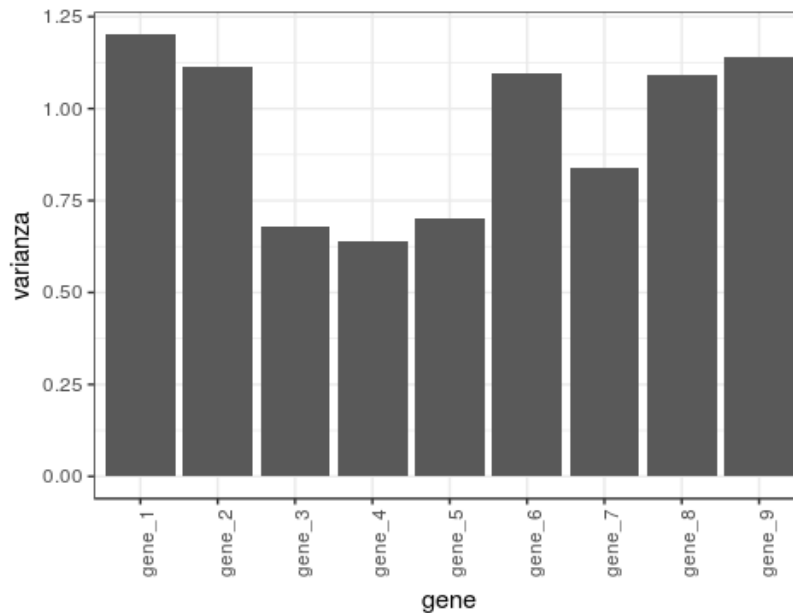
```
library(ggplot2)
library(tidyr)
library(dplyr)
datos <- data.frame(gene = paste("gene_", 1:9, sep = ""), tumor_a, tumor_b,
                    tumor_c, tumor_d, tumor_e, tumor_f, tumor_g,
                    tumor_h, tumor_i)
datos_tidy <- gather(data = datos, key = tumor, value = expresion, -1)
```

```
ggplot(data = datos_tidy, aes(x = gene, y = expresion, color = tumor)) +
  geom_path(aes(group = tumor)) +
  geom_point() +
  theme_bw() +
  facet_grid(tumor~.) +
  theme(legend.position = "none")
```



PCA es dependiente de la escala y varianza de los predictores. En este caso se va a suponer que los genes se han medido con el mismo instrumento y la misma escala. Solo queda entonces comprobar que ninguno de ellos tiene una varianza mucho mayor al resto.

```
datos_tidy %>% group_by(gene) %>% summarise(varianza = var(expresion)) %>%
  ggplot(aes(x = gene, y = varianza)) +
  geom_col() +
  theme_bw() +
  theme(axis.text.x = element_text(angle = 90, hjust = 1))
```



No hay ningún gen cuya varianza sea mucho mayor que la del resto, lo que significa que ninguno de ellos va dirigir el patrón resultante más que los demás.

## PCA

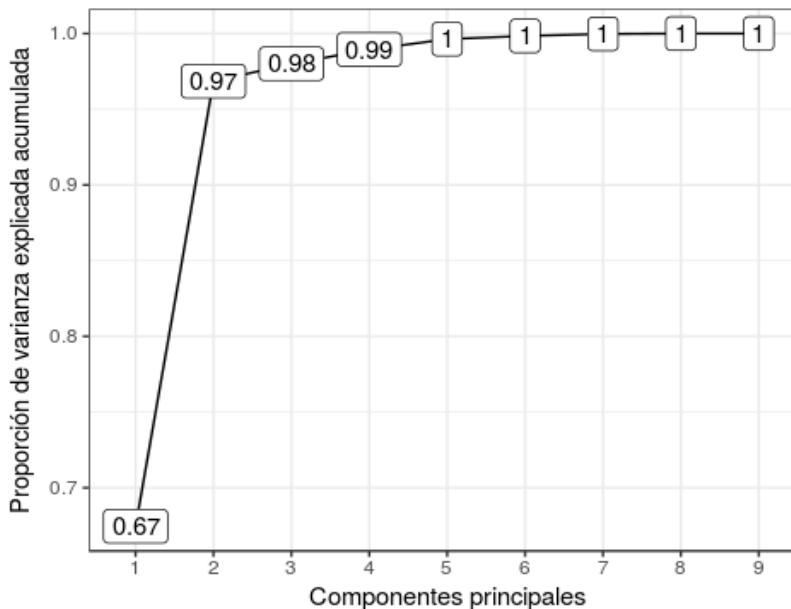
# Dado que se quieren estudiar patrones empleando la expresión de los genes, cada gen debe de estar en una columna (predictores)

```
datos_trans <- datos_tidy %>% spread(key = gene, value = expresion)
rownames(datos_trans) <- datos_trans$tumor
datos_trans <- datos_trans[, -1]
head(datos_trans)
```

```
##           gene_1    gene_2    gene_3    gene_4    gene_5
## tumor_a  0.9323190  1.2113278 -0.4616314 -0.88042240 -1.09941170
## tumor_b  0.8975893  1.4624217  0.1259876 -0.43860412 -1.22377397
## tumor_c  0.6811663  1.1275117 -0.2017863  0.02108253 -0.81164718
## tumor_d -0.8451067 -1.2319309 -0.9265317 -0.53071756 -0.05634140
## tumor_e -1.6228992 -0.7777483 -0.7868679 -0.67712937  0.19922605
## tumor_f -1.0805006 -0.9383486 -1.0821545 -0.60355858  0.04880287
##           gene_6    gene_7    gene_8    gene_9
## tumor_a -0.96955658 -0.06758821 0.7638985 1.1145246
## tumor_b -1.49015370  0.15747862 1.0964449 0.9901505
## tumor_c -1.02881566 -0.53623375 1.1793956 1.1355446
## tumor_d -0.27472142  1.04447944 0.7132720 0.7550505
## tumor_e  0.11591201  1.19523024 1.1565408 0.7884622
## tumor_f  0.08683008  0.95150260 1.1025851 0.6450739
```

```
pca <- prcomp(datos_trans)

# Cálculo de la varianza explicada acumulada
prop_varianza <- pca$sdev^2/sum(pca$sdev^2)
prop_varianza_acum <- cumsum(prop_varianza)
ggplot(data = data.frame(prop_varianza_acum, pc = factor(1:9)),
      aes(x = pc, y = prop_varianza_acum, group = 1)) +
  geom_point() +
  geom_line() +
  geom_label(aes(label = round(prop_varianza_acum,2))) +
  theme_bw() +
  labs(x = "Componentes principales",
       y = "Proporción de varianza explicada acumulada")
```



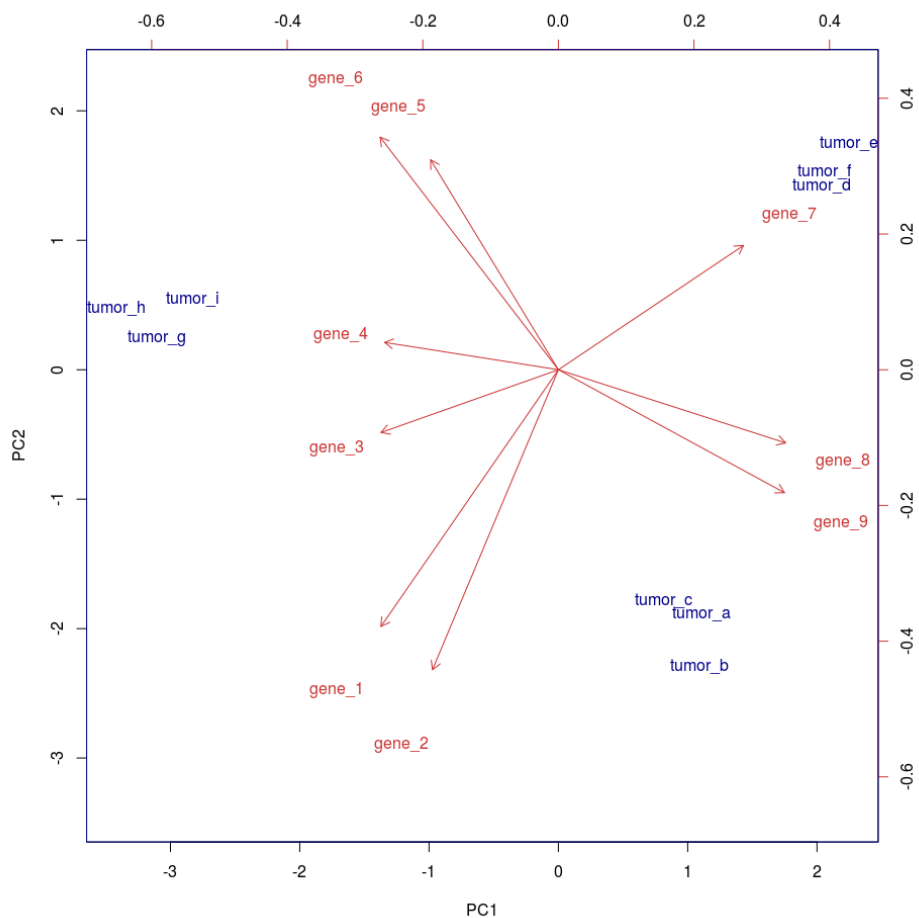
El estudio de la varianza explicada acumulada muestra que con solo dos componentes se puede capturar el 97% de la varianza, prácticamente la totalidad de la información contenida en la expresión de los 9 genes. Viendo en detalle el valor de los *loadings* de las dos primeras componentes se observa que en la primera los 9 genes tienen la misma importancia y lo mismo ocurre en la segunda, a excepción del gen 4 que tiene un peso mucho menor.

```
pca$rotation[, 1:2]
```

```
##          PC1          PC2
## gene_1 -0.3271796 -0.47294733
## gene_2 -0.2321769 -0.55256710
## gene_3 -0.3270151 -0.11554638
## gene_4 -0.3203507  0.05025632
## gene_5 -0.2356233  0.38672726
## gene_6 -0.3286415  0.42850939
## gene_7  0.3414044  0.22898673
## gene_8  0.4191123 -0.13487735
## gene_9  0.4168313 -0.22651416
```

Finalmente, la proyección de los datos en las dos primeras componentes muestra que el PCA ha sido capaz de encontrar un claro patrón que diferencia los 3 tipos de tumores.

```
biplot(x = pca, scale = 0, cex = 1, col = c("blue4", "brown3"))
```





## PCR: PCA aplicado a regresión lineal

El método *Principal Components Regression PCR* consiste en ajustar un modelo de regresión lineal por mínimos cuadrados empleando como predictores las componentes generadas a partir de un *Principal Component Analysis (PCA)*. De esta forma, con un número reducido de componentes se puede explicar la mayor parte de la varianza de los datos.

En los estudios observacionales, es frecuente disponer de un número elevado de variables que se pueden emplear como predictores. A pesar de ello, un alto número de predictores no implica necesariamente mucha información. Si las variables están correlacionadas entre ellas, la información que aportan es redundante y además viola la condición de no colinealidad necesaria en la regresión por mínimos cuadrados. Dado que el PCA es útil eliminando información redundante, si se emplean como predictores las componentes principales se puede mejorar el modelo de regresión. Es importante tener en cuenta que, si bien el *Principal Components Regression* reduce el número de predictores del modelo, no se puede considerar como un método de selección de variables ya que todas ellas se necesitan para el cálculo de las componentes. La identificación del número óptimo de componentes principales que se emplean como predictores en *PCR* puede identificarse por *cross validation*.

El método PCR no deja de ser un ajuste lineal por mínimos cuadrados que emplea componentes principales como predictores, para que sea válido se tienen que cumplir las condiciones requeridas para la regresión por mínimos cuadrados.

## Ejemplo

*La cuantificación del contenido en grasa de la carne puede hacerse mediante técnicas de analítica química, sin embargo, este proceso es costoso en tiempo y recursos. Una posible alternativa para reducir costes y optimizar tiempo es emplear un espectrofotómetro (instrumento capaz de detectar la absorbancia que tiene un material a diferentes tipos de luz en función de sus características). Para comprobar su efectividad se mide el espectro de absorbancia de 100 longitudes de onda en 215 muestras de carne, cuyo contenido en grasa se obtiene también por análisis químico para poder comparar los resultados. El set de datos `meatspec` del paquete `faraway` contiene toda la información.*

```
library(faraway)
data(meatspec)
dim(meatspec)
```

```
## [1] 215 101
```

El set de datos contiene 101 columnas. Las 100 primeras, nombradas como  $V1, \dots, V100$  recogen el valor de absorbancia para cada una de las 100 longitudes de onda analizadas, y la columna *fat* el contenido en grasa medido por técnicas químicas.

Para poder evaluar la capacidad predictiva del modelo, se dividen las observaciones disponibles en dos grupos: uno de entrenamiento para ajustar el modelo (80% de los datos) y uno de test (20% de los datos).

```
training <- meatspec[1:172, ]
test <- meatspec[173:215, ]
```

En primer lugar se ajusta un modelo incluyendo todas las longitudes de onda como predictores.

```
modelo <- lm(fat ~ ., data = training)
summary(modelo)
```

```
##
## Call:
## lm(formula = fat ~ ., data = training)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.09837 -0.35779  0.04555  0.38080  2.33860
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      6.324       2.012   3.143 0.002439 **
## V1             12134.077    3659.798   3.316 0.001443 **
## V2             -12585.857    5971.891  -2.108 0.038605 *
## V3              -5107.556    9390.265  -0.544 0.588200
## V4              23880.493   17143.644   1.393 0.167977
## V5             -40509.555   22129.359  -1.831 0.071360 .
## V6              28469.416   19569.400   1.455 0.150134
## V7             -20901.082   12501.639  -1.672 0.098952 .
## V8               8369.465    7515.467   1.114 0.269193
## V9             -1539.328    5397.505  -0.285 0.776327
## V10             4706.267    7406.895   0.635 0.527217
## V11             7012.943   11720.620   0.598 0.551516
```

## V12	14891.444	20169.170	0.738	0.462749	
## V13	-30963.902	26186.839	-1.182	0.240983	
## V14	34338.612	22323.830	1.538	0.128444	
## V15	-22235.237	13842.268	-1.606	0.112640	
## V16	-7466.797	8558.172	-0.872	0.385890	
## V17	6716.653	6561.805	1.024	0.309500	
## V18	-2033.071	6741.330	-0.302	0.763851	
## V19	8541.212	9419.998	0.907	0.367627	
## V20	-1667.207	17300.433	-0.096	0.923500	
## V21	-31972.494	24622.615	-1.299	0.198317	
## V22	59526.389	27730.712	2.147	0.035244	*
## V23	-49241.388	23117.226	-2.130	0.036632	*
## V24	16184.597	16679.609	0.970	0.335180	
## V25	12077.951	10751.912	1.123	0.265081	
## V26	-12632.330	6774.573	-1.865	0.066361	.
## V27	-6298.837	7032.334	-0.896	0.373442	
## V28	29625.988	9011.227	3.288	0.001573	**
## V29	-39374.835	13561.228	-2.903	0.004914	**
## V30	31251.427	18742.000	1.667	0.099829	.
## V31	-27238.189	21335.756	-1.277	0.205887	
## V32	23009.543	19776.156	1.163	0.248522	
## V33	-4584.373	14572.471	-0.315	0.753995	
## V34	-5437.943	10344.728	-0.526	0.600754	
## V35	-6128.931	8762.663	-0.699	0.486564	
## V36	5599.605	6652.640	0.842	0.402776	
## V37	-5569.160	6670.198	-0.835	0.406557	
## V38	97.451	9291.480	0.010	0.991661	
## V39	36021.407	12574.711	2.865	0.005488	**
## V40	-54273.400	17144.384	-3.166	0.002280	**
## V41	52084.876	21758.024	2.394	0.019318	*
## V42	-48458.089	23950.549	-2.023	0.046813	*
## V43	29334.488	20232.617	1.450	0.151500	
## V44	-18282.834	13508.157	-1.353	0.180200	
## V45	22110.934	9725.348	2.274	0.026020	*
## V46	-11735.692	6631.245	-1.770	0.081061	.
## V47	-514.521	3800.612	-0.135	0.892696	
## V48	2551.480	6131.893	0.416	0.678592	
## V49	3707.639	8970.401	0.413	0.680618	
## V50	-25762.703	10934.783	-2.356	0.021236	*
## V51	46844.468	15367.852	3.048	0.003233	**
## V52	-47783.626	18069.344	-2.644	0.010065	*
## V53	26233.604	18822.491	1.394	0.167744	
## V54	87.825	17403.836	0.005	0.995988	
## V55	-8475.119	13232.005	-0.641	0.523908	
## V56	3488.507	7228.428	0.483	0.630858	
## V57	-1520.733	4988.093	-0.305	0.761355	
## V58	2275.175	5495.630	0.414	0.680124	
## V59	-5415.427	5721.475	-0.947	0.347099	
## V60	7152.015	4754.317	1.504	0.136935	
## V61	-4494.234	4512.937	-0.996	0.322702	

```
## V62          3662.045    4811.634    0.761 0.449129
## V63          13993.987    7098.106    1.972 0.052563 .
## V64         -23252.133    8973.839   -2.591 0.011604 *
## V65          4373.731   10048.591    0.435 0.664695
## V66          4580.913   10146.146    0.451 0.653011
## V67          -837.676   10747.974   -0.078 0.938097
## V68         -7074.425   10852.430   -0.652 0.516587
## V69          9506.571    9739.256    0.976 0.332325
## V70         -2765.100    9519.031   -0.290 0.772295
## V71         -1125.135    8586.061   -0.131 0.896113
## V72         -7295.096    7489.488   -0.974 0.333341
## V73          17059.811    6522.093    2.616 0.010870 *
## V74         -9889.553    6543.945   -1.511 0.135162
## V75          -325.615    6125.973   -0.053 0.957759
## V76           782.219    5421.002    0.144 0.885677
## V77          8058.935    5793.416    1.391 0.168554
## V78        -15869.978    6448.208   -2.461 0.016282 *
## V79          21768.619    6435.678    3.382 0.001172 **
## V80        -28338.145    8180.874   -3.464 0.000906 ***
## V81           8523.317   10053.153    0.848 0.399384
## V82          22319.451   12098.046    1.845 0.069226 .
## V83        -17244.722   13991.685   -1.232 0.221829
## V84        -18325.836   14959.964   -1.225 0.224627
## V85          33345.457   13868.197    2.404 0.018808 *
## V86         -7955.157   14571.278   -0.546 0.586813
## V87         -7837.966   16141.553   -0.486 0.628762
## V88         -1815.552   17261.928   -0.105 0.916532
## V89           631.595   15684.751    0.040 0.967992
## V90        -2701.955   16187.612   -0.167 0.867911
## V91          4375.678   19400.005    0.226 0.822199
## V92          12925.188   16456.244    0.785 0.434816
## V93         -7441.235   12417.883   -0.599 0.550923
## V94         -2464.532   11815.234   -0.209 0.835366
## V95         -2090.635    9666.576   -0.216 0.829394
## V96          10912.352    9950.716    1.097 0.276505
## V97        -20331.405   11022.234   -1.845 0.069270 .
## V98          3948.443    8227.133    0.480 0.632753
## V99          6358.930    8652.372    0.735 0.464800
## V100         -263.365    4104.463   -0.064 0.949019
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.074 on 71 degrees of freedom
## Multiple R-squared:  0.997, Adjusted R-squared:  0.9928
## F-statistic: 237.5 on 100 and 71 DF, p-value: < 2.2e-16
```

El valor  $R^2_{ajustado}$  obtenido es muy alto (0.9928) lo que indica que el modelo es capaz de predecir con gran exactitud el contenido en grasa de las observaciones con las que se ha

entrenado. El hecho de que el modelo en conjunto sea significativo (p-value: < 2.2e-16), pero que muy pocos de los predictores lo sean a nivel individual, es un indicativo de una posible redundancia entre los predictores (colinealidad).

¿Cómo de bueno es el modelo prediciendo nuevas observaciones que no han participado en ajuste? Al tratarse de un modelo de regresión, la estimación del error de predicción se obtiene mediante el *Mean Square Error (MSE)*.

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

```
# MSE empleando las observaciones de entrenamiento
training_mse <- mean((modelo$fitted.values - training$fat)^2)
training_mse
```

```
## [1] 0.4765372
```

```
# MSE empleando nuevas observaciones
predicciones <- predict(modelo, newdata = test)
test_mse <- mean((predicciones - test$fat)^2)
test_mse
```

```
## [1] 14.54659
```

Se observa que el modelo tiene un MSE muy bajo (0.48) cuando predice las mismas observaciones con las que se ha entrenado, pero 30 veces más alto (14.54) al predecir nuevas observaciones. Esto significa que el modelo no es útil, ya que el objetivo es aplicarlo para predecir el contenido en grasa de futuras muestras de carne. A este problema se le conoce como *overfitting*. Una de las causas por las que un modelo puede sufrir *overfitting* es la incorporación de predictores innecesarios, que no aportan información o que la información que aportan es redundante.

Se recurre en primer lugar a la selección de predictores mediante *stepwise selection* empleando el AIC como criterio de evaluación:

```
modelo_step_selection <- step(object = modelo, trace = FALSE)
# Número de predictores del modelo resultante
length(modelo_step_selection$coefficients)
```

```
## [1] 73
```

```
# Training-MSE
training_mse <- mean((modelo_step_selection$fitted.values - training$fat)^2)
training_mse
```

```
## [1] 0.5034001
```

```
# Test-MSE
predicciones <- predict(modelo_step_selection, newdata = test)
test_mse <- mean((predicciones - test$fat)^2)
test_mse
```

```
## [1] 12.88986
```

El proceso de *stepwise selection* devuelve como mejor modelo el formado por 73 de los 100 predictores disponibles. Al haber eliminado predictores del modelo, el *training MSE* siempre aumenta, en este caso de 0.48 a 0.05, pero el *test-MSE* se ha reducido a 12.88986.

Véase ahora el resultado si se ajusta el modelo empleando las componentes principales:

```
# Cálculo de componentes principales. Se excluye la columna con la variable
# respuesta *fat*
pca <- prcomp(training[, -101], scale. = TRUE)

# Se muestra la proporción de varianza explicada y acumulada de las 9
# primeras componentes
summary(pca)$importance[, 1:9]
```

##	PC1	PC2	PC3	PC4	PC5
## Standard deviation	9.92492	1.043606	0.5357885	0.3312792	0.07898436
## Proportion of Variance	0.98504	0.010890	0.0028700	0.0011000	0.00006000
## Cumulative Proportion	0.98504	0.995930	0.9988000	0.9999000	0.99996000
##	PC6	PC7	PC8	PC9	
## Standard deviation	0.04974461	0.02700194	0.02059129	0.008603878	
## Proportion of Variance	0.00002000	0.00001000	0.00000000	0.00000000	
## Cumulative Proportion	0.99999000	0.99999000	1.00000000	1.00000000	

El estudio de la proporción de varianza explicada muestra que la primera componente recoge la mayor parte de la información (98.5%), decayendo drásticamente la varianza en las sucesivas componentes.

Una vez obtenido el valor de las componentes para cada observación (*principal component scores*), puede ajustarse el modelo lineal empleando dichos valores junto con la

variable respuesta que le corresponde a cada observación. Con la función `pcr()` del paquete `pls` se evita tener que codificar cada uno de los pasos intermedios.

Acorde a la proporción de varianza acumulada, emplear las 4 primeras componentes podría ser una buena elección, ya que en conjunto explican el 99.99100 % de varianza.

```
library(pls)
modelo_pcr <- pcr(formula = fat ~ ., data = training, scale. = TRUE, ncomp = 4)

# Test-MSE
predicciones <- predict(modelo_pcr, newdata = test, ncomp = 4)
test_mse <- mean((predicciones - test$fat)^2)
test_mse
```

```
## [1] 20.55699
```

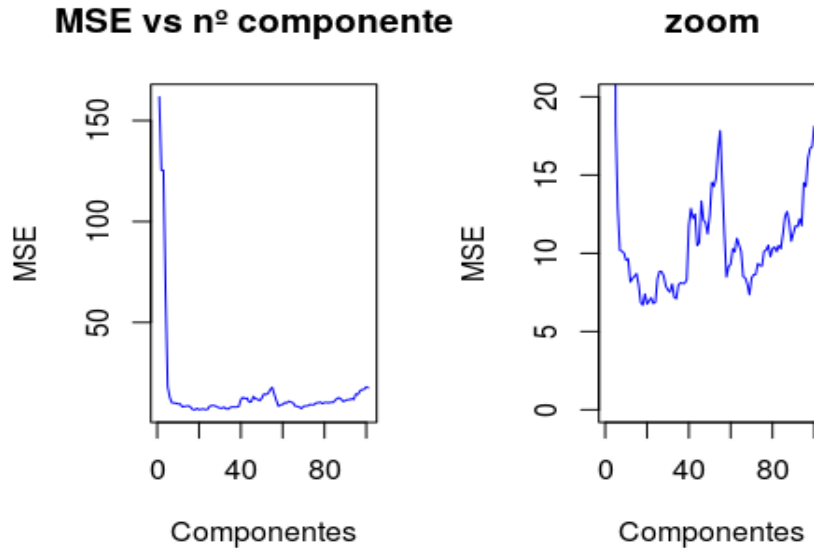
El *test-MSE* obtenido (20.56) para el modelo que emplea como predictores las 4 primeras componentes es mucho mayor que el obtenido con el modelo generado por *stepwise selection* (12.89) e incluso que el obtenido incluyendo todos los predictores (14.54659). Esto significa que, o bien el hecho de emplear componentes principales como predictores no es útil para este caso, o que el número de componentes incluido no es el adecuado.

La función `pcr()` incluye la posibilidad de recurrir a *cross validation* para identificar el número óptimo de componentes con el que se minimiza el *MSE*.

```
set.seed(123)
modelo_pcr <- pcr(formula = fat ~ ., data = training, scale. = TRUE, validation =
"CV")
modelo_pcr_CV <- MSEP(modelo_pcr, estimate = "CV")
which.min(modelo_pcr_CV$val)
```

```
## [1] 18
```

```
par(mfrow = c(1, 2))
plot(modelo_pcr_CV$val, main = "MSE vs nº componentes", type = "l", ylab = "MSE",
col = "blue", xlab = "Componentes")
plot(modelo_pcr_CV$val, main = "zoom", type = "l", ylab = "MSE",
xlab = "Componentes", col = "blue", ylim = c(0, 20))
```



```
# Test-MSE
predicciones <- predict(modelo_pcr, newdata = test, ncomp = 18)
test_mse <- mean((predicciones - test$fat)^2)
test_mse
```

```
## [1] 4.524698
```

El número óptimo de componentes principales identificado por *cross validation* es de 18. Empleando este número en la *PCR* se consigue reducir el *test-MSE* a 4.52, un valor muy por debajo del conseguido con los otros modelos.

## PLS: PCA aplicado a regresión lineal

El método *Partial Least Squares (PLS)* es muy similar al *PCR* en cuanto que ambos emplean las componentes principales resultantes de un análisis *PCA* como predictores. La diferencia reside en que, mientras *PCR* ignora la variable respuesta *Y* para determinar las combinaciones lineales, *PLS* busca aquellas que, además de explicar la varianza observada, predicen *Y* lo mejor posible. Puede considerarse como una versión *supervised* de *PCR*.



## Ejemplo

La cuantificación del contenido en grasa de la carne puede hacerse mediante técnicas de analítica química, sin embargo, este proceso es costoso en tiempo y recursos. Una posible alternativa para reducir costes y optimizar tiempo es emplear un espectrofotómetro (instrumento capaz de detectar la absorbancia que tiene un material a diferentes tipos de luz en función de sus características). Para comprobar su efectividad se mide el espectro de absorbancia de 100 longitudes de onda en 215 muestras de carne, cuyo contenido en grasa se obtiene también por análisis químico para poder comparar los resultados. El set de datos `meatspec` del paquete `faraway` contiene toda la información.

```
library(faraway)
data(meatspec)
```

Igual que en el ejemplo anterior de *PCR*, para poder evaluar la capacidad predictiva del modelo, se dividen las observaciones disponibles en dos grupos: uno de entrenamiento para ajustar el modelo (80% de los datos) y uno de test (20% de los datos).

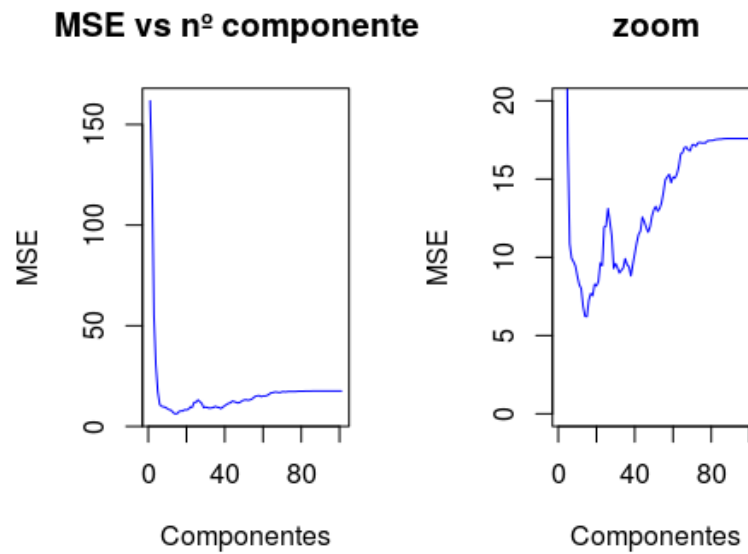
```
training <- meatspec[1:172, ]
test <- meatspec[173:215, ]
```

La función `plsr()` del paquete `pls` ajusta modelos por *Partial Least Squares* e incluye la posibilidad de recurrir a *cross validation* para identificar el número óptimo de componentes con el que se minimiza el *MSE*.

```
library(pls)
set.seed(123)
modelo_pls <- plsr(formula = fat ~ ., data = training, scale. = TRUE,
                    validation = "CV")
modelo_pls_CV <- MSEP(modelo_pls, estimate = "CV")
which.min(modelo_pls_CV$val)
```

```
## [1] 15
```

```
par(mfrow = c(1, 2))
plot(modelo_pls_CV$val, main = "MSE vs nº componentes", type = "l", ylab = "MSE",
     col = "blue", xlab = "Componentes")
plot(modelo_pls_CV$val, main = "zoom", type = "l", ylab = "MSE",
     xlab = "Componentes", col = "blue", ylim = c(0, 20))
```



```
# Test-MSE
predicciones <- predict(modelo_pls, newdata = test, ncomp = 15)
test_mse <- mean((predicciones - test$fat)^2)
test_mse
```

```
## [1] 3.888104
```

Si se comparan los resultados obtenidos por *PCR* y *PLS* se observa que el número de componentes óptimo es inferior en *PLS*. Esto suele ser así ya que en el proceso de *PLS* se está incluyendo información adicional a través de la variable respuesta. Para este ejemplo, el método *PLS* consigue un *test-MSE* ligeramente inferior al obtenido por *PCR*.

## t-SNE

### Idea intuitiva

Cuando se desea reducir la dimensionalidad de un conjunto de datos, es decir, condensar la información contenida en muchas variables en solo unas pocas a la vez que se mantiene la "estructura" original de los datos, la técnica de *Principal Component Analysis (PCA)* es uno de los estándares. A pesar de las muy buenas propiedades que tiene el *PCA*, sufre de algunas limitaciones, por ejemplo, solo tiene en cuenta combinaciones lineales de las variables originales. En determinados escenarios, el no poder considerar otro tipo de combinaciones supone perder mucha información. En 2008, *Geoffrey Hinton* y *Laurens van der Maaten* publicaron el método no lineal *t-distributed stochastic neighbor embedding (t-SNE)* que consigue superar esta limitación y por lo tanto supera al *PCA* en muchos escenarios.

### Algoritmo de t-SNE

*t-SNE* surge como una extensión del algoritmo *stochastic neighbor embedding (SNE)*. Los dos primeros pasos de ambos algoritmos son idénticos:

1. Se convierten las distancias euclídeas multidimensionales entre pares de observaciones en probabilidades condicionales. En lugar de medir la similitud entre dos observaciones  $x_i$  y  $x_j$  mediante la distancia euclídea, se mide como la probabilidad condicional ( $p_{j|i}$ ) de que  $x_j$  fuese seleccionada como observación vecina de  $x_i$  si las observaciones perteneciesen a una distribución de densidad gaussiana centrada en  $x_i$ . Para observaciones cercanas,  $p_{j|i}$  es alta, mientras que para observaciones alejadas la probabilidad es mínima.

$$p_{j|i} = \frac{\exp(-||x_i - x_j||^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-||x_i - x_k||^2 / 2\sigma_i^2)}$$

donde  $\sigma_i$  es la varianza de la distribución gaussiana centrada en  $x_i$ .

2. Se crean dos observaciones  $y_i, y_j$  homólogas a  $x_i$  y  $x_j$  pero en una dimensionalidad menor, cuya distancia se define como la probabilidad condicional  $q_{j|i}$ .

$$q_{j|i} = \frac{\exp(-||y_i - y_j||^2)}{\sum_{k \neq i} \exp(-||y_i - y_k||^2)}$$

Como resultado de los pasos 1 y 2 se obtiene la probabilidad condicional de similitud entre pares de observaciones en el espacio multidimensional original y en un espacio de menor dimensión. Para que las observaciones  $y_i$ ,  $y_j$  creadas en el espacio de menor dimensión sean files representantes de las observaciones  $x_i$  y  $x_j$  del espacio superior, las probabilidades condicionales  $p_{j|i}$  y  $q_{j|i}$  deben ser exactamente las mismas, o lo que es lo mismo, su diferencia tiene que ser cero. Es el modo en el que *t-SNE* minimiza las diferencias de las probabilidades condicionales lo que lo distingue del método *SNE* (el método matemático empleado queda fuera del objetivo de este ensayo).

Para poder obtener las probabilidades condicionales tal como se ha descrito en los pasos 1 y 2, es necesario seleccionar un valor de  $\sigma_i$ . Es muy poco probable que exista un valor de varianza  $\sigma_i$  óptimo para todas las observaciones del set de datos, ya que la densidad de los datos suele variar de una región a otra. Para regiones densas, un valor pequeño de  $\sigma_i$  es mucho más apropiado que para regiones dispersas. *t-SNE* realiza una búsqueda binaria para encontrar el valor óptimo  $\sigma_i$  con una *perplexity* fijada por el usuario. La *perplexity* puede entenderse como una medida del número de observaciones vecinas que tienen que emplearse en cada estimación local. Suele ser recomendable seleccionar valores entre 5 y 50.

El algoritmo de *t-SNE* calcula la probabilidad condicional de cada par de observaciones y trata de minimizar la suma de las diferencias entre las probabilidades de la dimensión superior e inferior. Esto implica una cantidad enorme de cálculos, lo que se traduce en un proceso muy lento cuando se trabaja con varios miles de observaciones. Existe una aproximación más rápida del algoritmo conocida como *Barnes-Hut*.

## Interpretación de los resultados de un t-SNE

A continuación, se enumeran algunos puntos clave para interpretar correctamente los resultados obtenidos por el método *t-SNE*. En la página <https://distill.pub/2016/misread-tsne/> pueden encontrarse ejemplos muy detallados e intuitivos.

- El valor de *perplexity* debe de ser menor que el número de observaciones del set de datos. Por lo general se recomienda trabajar dentro del rango 5-50.
- Dependiendo del valor de *perplexity* empleado, tanto la distancia entre *clusters* como su forma puede variar en gran medida.
- El proceso no es determinista, puede ocurrir que, ejecutando el mismo algoritmo, con los mismos parámetros, en los mismos datos, se obtengan resultados ligeramente distintos.

- Idealmente, al representar los resultados de un *t-SNE*, se observa un patrón de agrupaciones. Sin embargo, su interpretación no debe hacerse de la misma forma que en un proceso de *clustering*. En *t-SNE*, ni el tamaño de los *clusters*, ni su densidad, ni su separación deben entenderse como indicadores de dispersión u otras características de la población con la que se está trabajando. El algoritmo de *t-SNE* actúa de tal forma que puede expandir *clusters* densos y compactar *clusters* dispersos, de ahí lo buenas que suelen ser sus visualizaciones.

## Limitaciones y desventajas

La reducción de dimensionalidad que consigue el *t-SNE* está basada en el concepto matemático de *manifold*. Un *manifold* se define como una superficie  $d$ -dimensional que reside dentro de un espacio  $D$ -dimensional, siendo  $d < D$ . Un ejemplo en 3 dimensiones sería el siguiente: una hoja de papel tiene dos dimensiones, pero, si se arruga en forma de pelota, pasa a estar en un espacio de 3 dimensiones. Por mucho que se moldee, es posible desplegarla recuperando la estructura original de dos dimensiones. La hoja de papel representa un *manifold* 2D en un espacio 3D. Si en lugar de una hoja, se tratase de un hilo enredado formando de ovillo, sería un *manifold* 1D en un espacio 3D. En aquellos casos en los que se cumple la hipótesis o asunción de *manifold* (los datos forman un *manifold* de menor dimensión que el espacio en el que se están observando), entonces, el *t-SNE* es capaz de proyectar correctamente los datos en esa dimensionalidad inferior. En la realidad, la asunción de *manifold* no puede cumplirse siempre, de hecho, si así fuese, cualquier set de datos con  $N$ -dimensiones, podría considerarse como un *manifold* de dimensión  $M$ , siendo  $M < N$ . Una vez proyectados los datos en ese espacio  $M$ -dimensional, se podría repetir el proceso sucesivamente hasta concluir que los datos originales de dimensión  $N$  pertenecen a un *manifold* de dimensión 1. Dado que esto no es cierto para todos los sets de datos, la asunción de *manifold* no siempre se cumple. Es en estos casos, en los que *t-SNE* no resulta útil reduciendo la dimensionalidad.

En cuanto a las desventajas del *t-SNE*, pueden destacarse:

- Está diseñado para reducir los datos a 2 o 3 dimensiones. No debe aplicarse si la reducción de dimensionalidad se hace a espacios de  $d > 3$ .
- La baja interpretación de su algoritmo (*black box*), es decir, permite visualizar muy bien datos, pero no genera, por ejemplo, una serie de ecuaciones fácilmente interpretables como hace el *PCA*.
- No es un proceso totalmente determinista, por lo que, a pesar de emplear los mismos datos, los resultados pueden variar.

- Su algoritmo no es incremental, en otras palabras, no puede aplicarse sobre un set de datos y después actualizarlo con unas pocas observaciones nuevas. Se tiene que ejecutar de nuevo todo el algoritmo incluyendo todas las observaciones (las viejas y las nuevas).

## Ejemplo con tsne

El paquete `tsne` contiene una implementación del algoritmo *t-SNE* original, no la aproximación de *Barnes-Hut*.

*El UC Irvine Machine Learning Repository contiene multitud de sets de datos muy adecuados para el estudio de métodos estadísticos y de machine learning. En concreto, el [Optical Recognition of Handwritten Digits Data Set](#), contiene la imagen digitalizada de 3822 números (del 0 al 9) escritos a mano. La información digitalizada de cada número genera un espacio de 64 dimensiones, es decir, digitalmente, cada número está descrito por 64 variables. Se pretende evaluar si t-SNE es capaz de reducir todas esas variables a solo 2, sin perder el patrón/estructura de los datos.*

```
library(readr)
library(dplyr)

# Carga de datos
datos <- read_csv(paste0("http://archive.ics.uci.edu/ml/machine-learning-",
                        "databases/optdigits/optdigits.tra"))
dim(datos)
```

```
## [1] 3822 65
```

```
# La última columna contiene el número real al que se corresponde la observación.
# Se renombra como "numero"
datos <- datos %>% rename(numero = `0_26`)

# La función tsne() recibe como argumento una matriz, no un data.frames
datos <- data.matrix(datos)

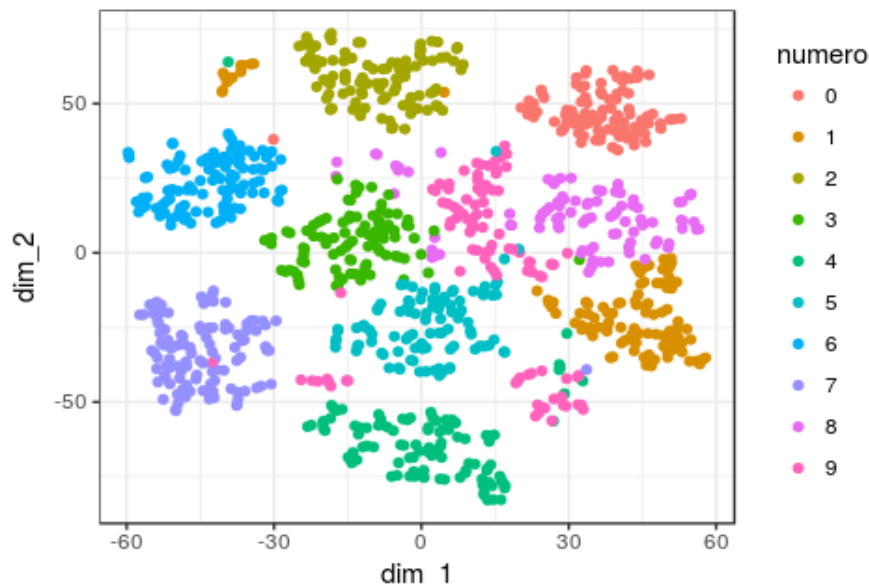
# Debido a los requerimientos computacionales del t-SNE, se limita este ejemplo
# únicamente a 1000 observaciones.
datos <- datos[1:1000,]
```

A diferencia del proceso de PCA, en el que se suelen calcular todas o gran parte de las componentes principales y luego se estudia qué porcentaje de información contiene cada una, en *tSNE* se especifica de antemano a cuántas dimensiones se tiene que reducir el espacio original.

```
library(tsne)
library(ggplot2)

# También se limita el número de iteraciones (epoch) a 100, aunque los
# resultados podrían mejorar si se aumentara
set.seed(321)
tsne_reduction <- tsne(datos, k = 2, perplexity = 30, epoch = 100)

# Para poder representar el verdadero número al que corresponde cada imagen,
# se adjunta la variable "numero" del set de datos
resultados <- as.data.frame(tsne_reduction)
colnames(resultados) <- c("dim_1", "dim_2")
resultados$numero <- as.character(datos[, "numero"])
ggplot(data = resultados, aes(x = dim_1, y = dim_2)) +
  geom_point(aes(color = numero)) +
  theme_bw()
```



## Ejemplo con Rtsne

El paquete `Rtsne` contiene una implementación del algoritmo *t-SNE* mediante la aproximación de *Barnes-Hut*. Se trata de un algoritmo ligeramente distinto al *t-SNE* original, pero más veloz, que permite aplicar *t-SNE* a sets de datos con miles de observaciones.

*Utilizando el mismo set de datos que en el ejemplo anterior, pero incluyendo las 3822 observaciones, se reduce la dimensionalidad, esta vez empleando la aproximación de Barnes-Hut.*

```
library(readr)
library(dplyr)

# Carga de datos
datos <- read_csv(paste0("http://archive.ics.uci.edu/ml/machine-learning-",
                        "databases/optdigits/optdigits.tra"))
dim(datos)
```

```
## [1] 3822 65
```

```
# La última columna contiene el número real al que se corresponde la observación.
# Se renombra como "numero"
datos <- datos %>% rename(numero = `0_26`)

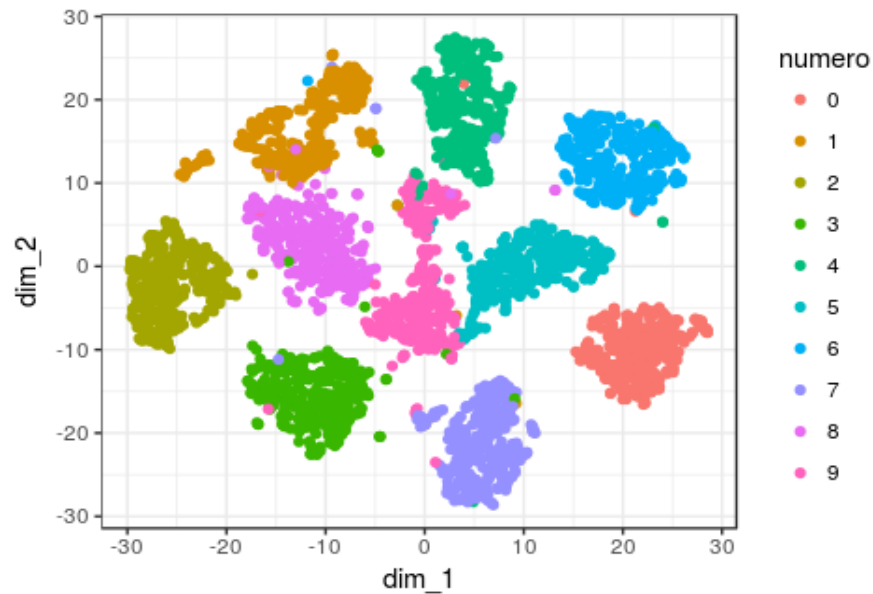
library(Rtsne)
library(ggplot2)
tsne <- Rtsne(X = datos[, -65], is_distance = FALSE, dims = 2, perplexity = 30,
             theta = 0.5, max_iter = 500)

# El objeto devuelto por Rtsne() almacena los valores de las dimensiones en el
# elemento Y. Como en este caso se ha especificado que la reducción se haga
# a dos dimensiones (k=2), Y tiene solo dos columnas.
head(tsne$Y)
```

```
##           [,1]      [,2]
## [1,] 19.3153643 -8.000100
## [2,]  7.1142258 -24.484858
## [3,] -0.2274412  9.207366
## [4,] 21.8018590  7.262071
## [5,] -26.6238324 -7.480648
## [6,]  0.8084161  5.393628
```

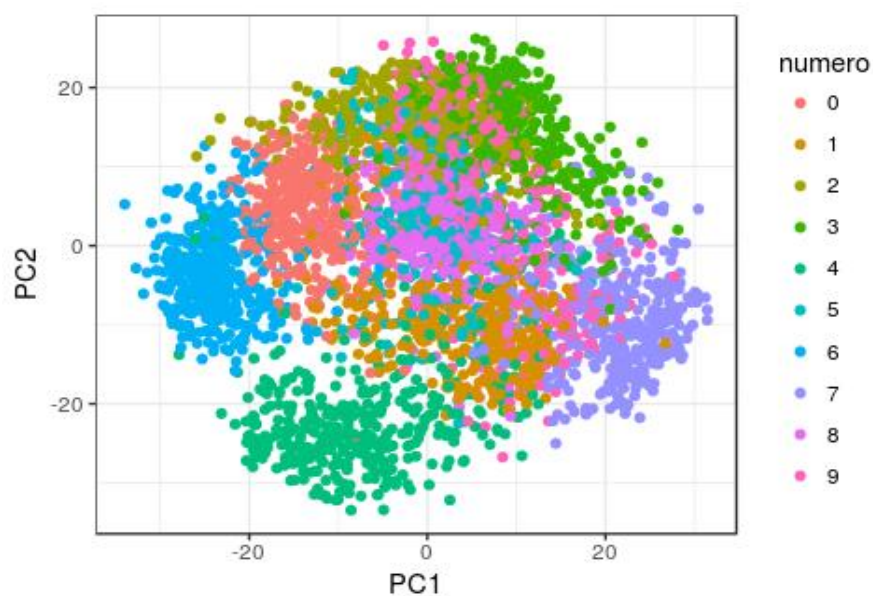
```
# Para poder representar el verdadero número al que corresponde cada imagen,
# se adjunta la variable "numero" del set de datos
resultados <- as.data.frame(tsne$Y)
colnames(resultados) <- c("dim_1", "dim_2")
resultados$numero <- as.character(datos$numero)
ggplot(data = resultados, aes(x = dim_1, y = dim_2)) +
  geom_point(aes(color = numero)) +
  theme_bw()
```





Se reduce de nuevo la dimensionalidad, pero esta vez empleando *PCA* y representando las dos primeras componentes.

```
pca <- prcomp(x = datos[, -65])
resultados <- as.data.frame(pca$x[, 1:2])
resultados$numero <- as.character(datos$numero)
ggplot(data = resultados, aes(x = PC1, y = PC2)) +
  geom_point(aes(color = numero)) +
  theme_bw()
```

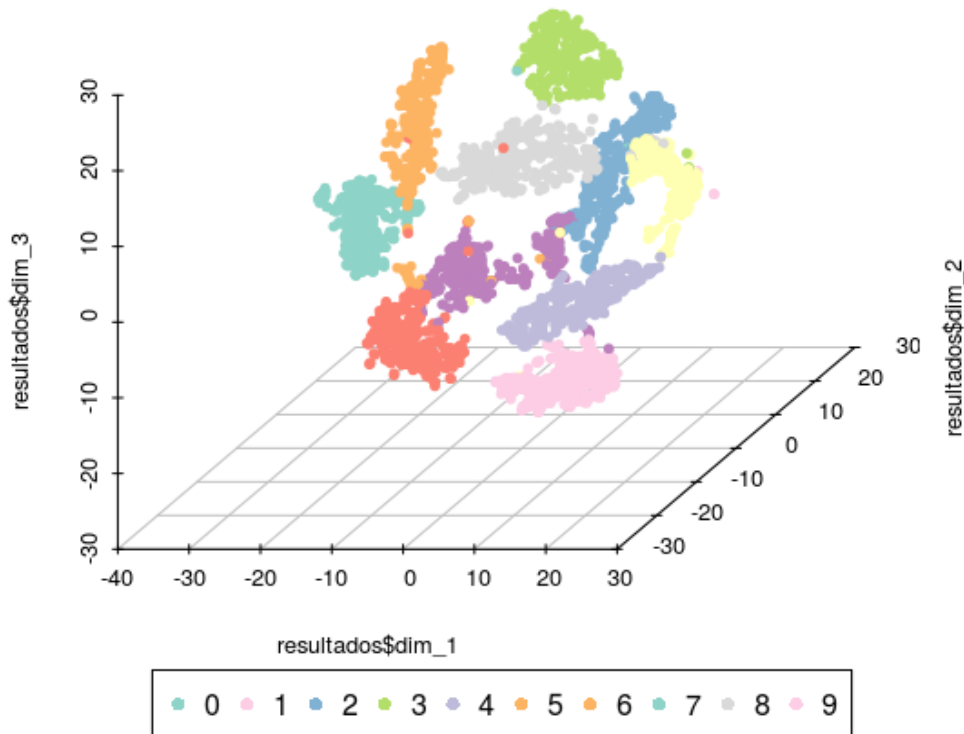


En este escenario, el método *tSNE* consigue que, al reducir la dimensionalidad, la separación entre observaciones sea sustancialmente mejor que cuando se emplea *PCA*.

Véase ahora la reducción a un espacio de 3 dimensiones.

```
library(scatterplot3d)
library(RColorBrewer)
tsne <- Rtsne(X = datos[, -65], is_distance = FALSE, dims = 3, perplexity = 30,
              theta = 0.5, max_iter = 500)
resultados <- as.data.frame(tsne$Y)
colnames(resultados) <- c("dim_1", "dim_2", "dim_3")
resultados$numero <- as.factor(datos$numero)

colores <- brewer.pal(n = 10, name = "Set3")
colores <- colores[as.numeric(resultados$numero)]
scatterplot3d(x = resultados$dim_1,
              y = resultados$dim_2,
              z = resultados$dim_3,
              pch = 20, color = colores, cex.lab = 0.8,
              grid = TRUE, box = FALSE)
legend("bottom", legend = levels(resultados$numero),
      col = colores, pch = 16,
      inset = -0.23, xpd = TRUE, horiz = TRUE)
```



## Otros métodos de reducción de dimensionalidad

El siguiente listado contiene otros métodos desarrollados para reducir la dimensionalidad de los datos.

- PCA (linear)
- t-SNE (non-parametric/ nonlinear)
- Sammon mapping (nonlinear)
- Isomap (nonlinear)
- LLE (nonlinear)
- CCA (nonlinear)
- SNE (nonlinear)
- MVU (nonlinear)
- Laplacian Eigenmaps (nonlinear)

## Bibliografía

*Introduction to Statistical Learning*

*A tutorial on Principal Components Analysis, Lindsay I Smith February 2002*

*Linear Models with R, Julian J. Faraway*

[https://en.wikipedia.org/wiki/Principal\\_component\\_analysis](https://en.wikipedia.org/wiki/Principal_component_analysis)

*Points of Significance: Principal component analysis, Nature Methods*

*What is principal component analysis?, Markus Ringnér*

*Visualizing Data using t-SNE, Laurens van der Maaten and Geoffrey Hinton*

<https://www.analyticsvidhya.com/blog/2017/01/t-sne-implementation-r-python/>

<https://datascienceplus.com/multi-dimensional-reduction-and-visualisation-with-t-sne/>

*An illustrated introduction to the t-SNE algorithm by Cyrille Rossant March 3, 2015*

<https://distill.pub/2016/misread-tsne/>

<http://lvdmaaten.github.io/tsne/>



This work by Joaquín Amat Rodrigo is licensed under a [Creative Commons Attribution-ShareAlike 4.0 International License](https://creativecommons.org/licenses/by-sa/4.0/).