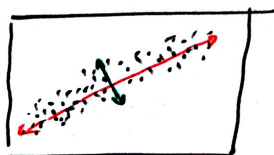
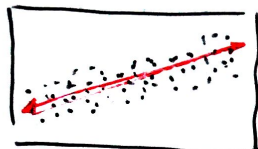


## Interpretación Geométrica

Supongase un conjunto de datos definido por dos variables  $(x_1, x_2)$  el vector que define la primera componente principal ( $z_1$ ) es aquel en la dirección de mayor varianza, la dirección en la que las observaciones presentan mayor varianza, la proyección de cada observación sobre esa dirección es el valor de la primera componente



la segunda componentes es la de mayor varianza también y que no tenga correlación con la primera, es decir, ortogonal

## Cálculo de componentes

Cada  $z$  se obtiene por combinación lineal de las variables originales

$$(x_1, x_2, x_3, \dots, x_p)$$

$$z_1 = \phi_{11} x_1 + \phi_{12} x_2 + \dots + \phi_{1p} x_p$$

$$\sum_{j=1}^p \phi_{1j}^2 = 1 \quad \text{combinación lineal normalizada}$$

↓ loadings

Proceso:

- Centrar las variables. Se resta a cada valor la media de la variable  $\Rightarrow$  media  $= 0$  . std  $= 1$
  - Resolver problema optimización para loadings para maximizar varianza  
cálculo de eigenvalues y eigenvectors de la matriz de covarianza
- $\Rightarrow$  Una vez se calcula la primera  $z_1$ , se calcula la segunda agregando la condición de que  $z_2$  y  $z_1$  no pueden estar correlacionados - ortogonales, el orden de importancia lo da el eigenvalue

## ## PCA\_ALGORITHM ##

Este sirve para reducir la dimension del dataset basado en la relacion existente entre los datos

Evaluacion previa

correlograma 
$$S(x, y) = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x) \text{Var}(y)}}$$

Varianza

$$\text{Var}(x) = \sum_{x \in R_x} (x - \mu)^2$$

covarianza

$$\text{Cov}(x, y) = \frac{1}{2n^2} \sum_{i=2}^n \sum_{j=1}^n (x_i - \bar{x}_j)(y_i - \bar{y}_j)$$

- Permite la reduccion de dimensiones bajando la complejidad del dataset

Suponga un espacio muestral de  $n$  muestras y  $p$  variables

PCA permite encontrar unos factores subyacentes  $z$  ( $z < p$ )

que explican aproximadamente lo mismo que  $p$ . estas  $z$  variables reciben el nombre de componentes principales

Cabe resaltar que es importante contar con las variables originales  
PCA se aplica para visualizacion y el preprocesado de predictores  
previo al ajuste de modelos supervisados

Conceptos base (Algebra lineal)

- eigenvalues

- eigenvectors

es el factor de multiplicacion del eigenvalor.

vectores que al multiplicarles por la matriz dan el mismo vector o un multiplo entero

• los eigenvalores son ortogonales

- Se recomienda remover outliers

Varianza explicada

Cuanta información se pierde al realizar PCA

Varianza total en el set de datos  $\sum_{j=1}^p \text{Var}(X_j) = \sum_{j=1}^p \frac{1}{n} \sum_{i=1}^n x_{ij}^2 \quad (1)$

Varianza explicada por la componente  $m$   $\frac{1}{n} \sum_{i=1}^n z_{i,m}^2 = \frac{1}{n} \sum_{i=1}^n \left( \sum_{j=1}^p \phi_{jm} x_{ij} \right)^2 \quad (2)$

la proporción de varianza explicada es el ratio  $(2)/(1)$

Número de componentes principales

Se estima con la varianza explicada acumulada mirando a partir de que componente deja de aumentar sustancialmente.