

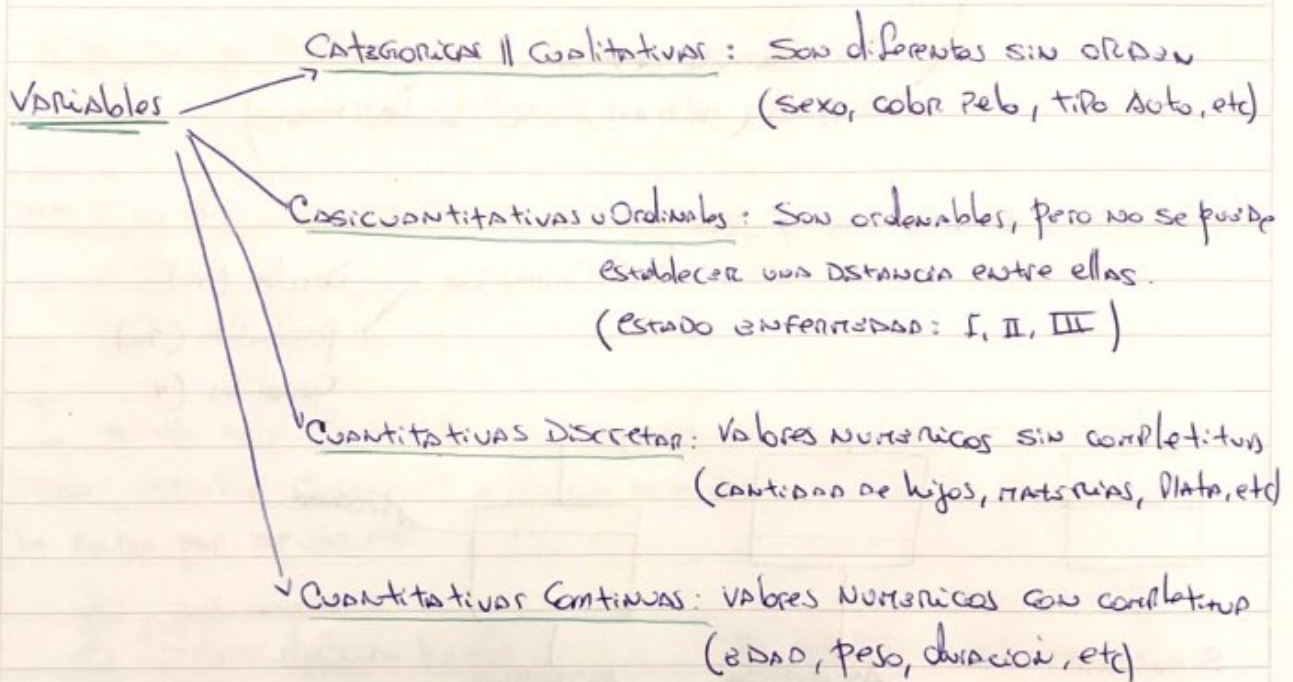
U₁/U₂

Estadística

- Procedimiento hipotético Deductivo
- Supuestos Iniciales
- Sin Herramientas Informáticas

Data Mining

- Procedimiento Inductivo
- Sin Supuestos Iniciales
- Amplia difusión entre computadores



Frec. Absoluta: # observaciones de cada modalidad/intervalo

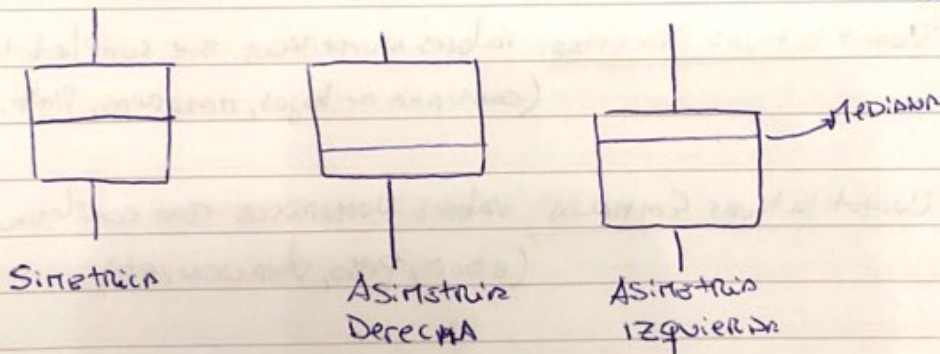
Frec. Relativa: $\frac{\text{Frec. Absoluta}}{\text{\# total observaciones}}$

11. Tendencias Central:
- Media / Promedio
 - Media α -Podada
 - Mediana
 - MODA

12. Dispersion:
- Rango muestral (MAX-MIN)
 - Varianza muestral: Promedio de las cuadradas de las distancias de las obs. a media
 - Coeficiente de Variancia
 - MAD (Median Absolute Deviation)

13. Posición estadísticas y Orden

- Cuantiles
 - Deciles (10)
 - Percentiles (100)
 - Quartiles (4)



INFO

- MULTIVARIADA

	VAR-1	VAR-2	...	VAR-N
IND-1	:	:	...	:
IND-2	:	:	...	:
:	:	:	...	:
IND-N	:	:	...	:

N Rows = # INDIVIDUOS

P-variables = # variables a ANALIZAR

- GRAFICO DE PLOTTING: Distr. Conjuntas MULTIVARIADAS

Singular = $\text{Det}(C) = 0$ = l. dependiente

DISPERSOGRAMA: variables Cuantitativas, vincula de a pares de variables

G. Estrellas: Cuantitativas para detectar similitudes

C. Chernoff: \nearrow

Posición y Dispersión Multivariadas

Vector medias muestral: $\bar{X} = (\bar{X}_1, \bar{X}_2, \dots, \bar{X}_p)$ donde \bar{X}_i promedio de la i -ésima variable.

Matriz de Varianzas y Covarianzas Muestral

$$\hat{\Sigma} = \frac{1}{n} (X - \bar{X})^t (X - \bar{X})$$

$\hat{\Sigma}$ tamaño $p \times p$, Simétrica, La diagonal son las varianzas muestrales de cada variable observada y fuera de la diagonal esta la covarianza de cada par de variables

$\hat{\Sigma}$ $\left\{ \begin{array}{l} \text{Simétrica} \\ \text{Semi-definida positiva} \\ \text{Autovalores} \geq 0 \end{array} \right.$

$\hat{\Sigma}_{ij} \left\{ \begin{array}{l} > 0 \text{ Asoc. lineal } \oplus \\ < 0 \text{ Asoc. lineal } \ominus \\ = 0 \text{ No hay Asoc. lineal} \end{array} \right.$

~~Matriz~~ Matriz de Correlación

\hookrightarrow Mat. de VAR y COV Normalizada

$$\frac{\hat{\Sigma}_{ij} - \bar{X}_i \bar{X}_j}{\sigma_i \sigma_j}$$

Normalizo si \bar{X}_i y \bar{X}_j son buenos estimadores

[La COV Detecta Asociación lineal únicamente]

r_{ik} Sensibles a
 S_{ik} outliers

$$r_{ik} = \frac{S_{ik}}{\sqrt{S_{ii}} \sqrt{S_{kk}}}$$

$$R = \begin{bmatrix} 1 & r_{12} & \dots & r_{1n} \\ r_{21} & 1 & & \\ \vdots & & \ddots & \\ r_{n1} & r_{n2} & \dots & 1 \end{bmatrix}$$

traza de Σ o R es $(+)$

Correlograma visualiza la fuerza y sentido de la correlación entre un conjunto de variables

- Azul : correlación $(+)$
- Rojo : correlación $(-)$

Alternativas Robustas

↳ Alternativa a las clásicas cuando hay presencia de outliers.

Efecto de Outliers → ENTASCARAMIENTO : un grupo esconde otro grupo

→ INUNDACIÓN : una observación es outlier solo si hay presencia de otra.

/// Dist Mahalanobis : considera la correlación entre variables

/// Vector D/mediana : Sustituir el μ de medias por el de medianas y calcular la mat. de covarianzas computando la menor Dist. Mahalanobis

/// MVE (Minimum Volume Ellipsoid) : Buscar el elipsoide de menor volumen que cubra m de las n observaciones.

/// MCD (Minimum Covariance Determinant) : minimiza el Det de la matriz de covarianzas de m ~~medias~~ observaciones de las n que hay.

U3//: ACP

$A \in \mathbb{R}^{m \times m}$ con autovalores $\lambda_1, \lambda_2, \dots, \lambda_m$

$$\text{Traza}(A) = \sum_{i=1}^m \lambda_i$$

$$\text{Det}(A) = \prod_{i=1}^m \lambda_i$$

Reducción de Dimensión \rightarrow ACP

- TRANSFORMA UN CONJUNTO DE VARIABLES CORRELACIONADAS EN UN GTO DE VARIABLES NO CORRELACIONADAS, de menor dimensión que se obtiene a partir de combinaciones lineales de las variables originales.
- Técnica Descriptiva, exploratoria, no tiene supuestos \Rightarrow "SIEMPRE PUEDE APLICARSE"
- Elige comb. lineales de las variables que maximizan la Varianza. (minimizan la pérdida de info inicial)
- Todas las variables juegan el mismo papel, no existen indep & dep.
- Descarta información redundante
- Alternativa para visualizar info multidimensional
- Solo tiene sentido si las variables originales están fuertemente correlacionadas.
- Es recomendable el ACP con la M. de Correlaciones (R) ya que Σ se ve influenciada por la variabilidad del conjunto.

$$Y_1 = Q_{11}X_1 + Q_{12}X_2 + \dots + Q_{1N}X_N$$

$$Y_2 = Q_{21}X_1 + Q_{22}X_2 + \dots + Q_{2N}X_N$$

donde $\underline{Q_i} = (Q_{i1}, Q_{i2}, \dots, Q_{iN}) \in \mathbb{R}^N$
es el vector de CARGAS/
LOADINGS

Buscan comb. lineal de variables que maximice variabilidad.

La variabilidad de x_1 es máxima cuando el vector de cargas es el autovector asociado al mayor autovalue de Σ (mat de Var + Cov)

x_2 : Hago lo mismo con las variables originales no correlacionadas con x_1 .

$$\text{Proporción de Variabilidad de } x_1 : \frac{\lambda_1}{\text{tr}(\Sigma)} = \frac{\lambda_1}{\lambda_1 + \lambda_2 + \dots + \lambda_n}$$

Comp Principales

Criterios De Parada.

- 1: Porcentaje de Variabilidad explicada.
- 2: Criterio de KAISER: Retener las m primeras componentes / sus autovectores Resultan $\gg 1$
 $\lambda_1 \gg \lambda_2 \gg \dots \gg \lambda_n \gg 1$
Algunos Autores recomiendan $\gg 0.7$
Se puede extender a Σ como $\gg \frac{\text{tr}(\Sigma)}{p}$
- 3: Criterio Del baston Pato: Se para cuando la variabilidad (Grafico de sedimentación) explicada se estabiliza a partir de un cierto valor.
- 4: Prueba de esfericidad: En un punto dga de haber direcciones de máxima variabilidad \Rightarrow La distribución tiene forma / empieza a tener forma de esferas.

Estimación, se usa $\sum^A = S$

Autovectores corresponden a la varianza de la componente $\Rightarrow (\text{Desvio estándar})^2$

Comunes usar R por sobre Σ

CARGAS

- \hookrightarrow De y_i es $\oplus \Rightarrow$ La variable y la componente tienen correlación \oplus
- \hookrightarrow Si es $\oplus \Rightarrow$ un individuo que tenga puntuación alta en esta variable tendrá valores más alto de esta variable que otros individuos
- \hookrightarrow y_i es $\ominus \Rightarrow$ correlación \ominus
- \hookrightarrow Dos individuos, con puntuaciones similares en las restantes variables, el que tenga más alta en esta se ubicará en un valor menor de la componente.

Primera Componente tiene todas las cargas positivas (+) y negativas (-)

\Rightarrow Componente de Tamaño: \Rightarrow Individuo con alto en esta tendrá en todas las tallas (Rapidez)

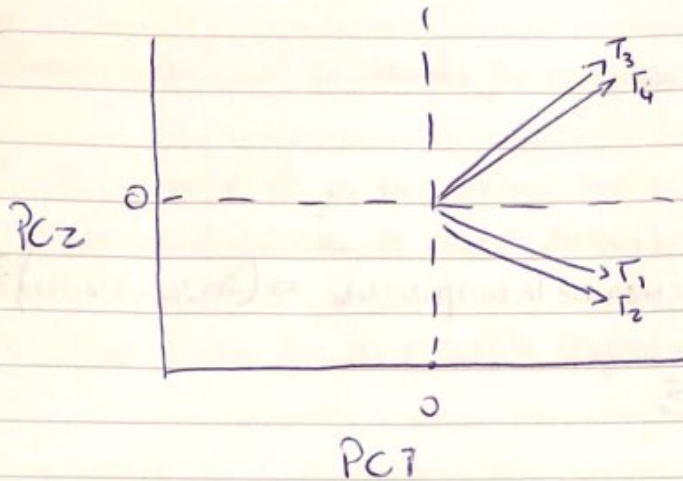
Segunda componente con valores intercalados (\oplus y \ominus)

\Rightarrow Componente de Forma contrasta los primeros 2 tramos con los últimos dos.

• Si tengo pocos individuos los grupo todos, con muchos agrupo

b. Plot: Las proyecciones sobre las componentes nos dan los Loadings

- \times chico \Rightarrow mucha correlación
- variables (T_1, T_2, T_3, \dots) muy o poco correlacionadas.



- / Las 4 proyectadas sobre PC1 son \oplus "Tamaño"
- / T3 y T4 proyectadas sobre PC2 son \oplus y T1 y T2 son \ominus "Forma"
- / Ptos cerca del origen son "Promedio"

Tamaño y Forma son no correlacionadas

Sino tengo la muestra y solo tengo Σ y R no puedo obtener los scores (Puntos en la PC)

ACP-Robustas

- Outliers pueden distorsionar la Mat. De Cov Muestral
- Métodos con menos supuestos pero computacionalmente costosos.

MCD (Minimum Covariance Det)
 Estimador Stahel-Duchu
 MVE (Minimum Variance Ellipse)

} Opciones para Reemplazar el vector de medias y la matriz de Cov.