



Universidad de Alcalá

Escuela Politécnica Superior

Universidad de Alcalá

Fundamentos de Ciencia de Datos

Juan J. Cuadrado Gallego

Martes 10:00 – 12:00

Grado en Ingeniería Informática – Curso 2019/2020

Marcos Barranquero Fernández – 51129104N

TEMA 1 – PROBABILIDAD Y ESTADÍSTICA

FRECUENCIA

F. Absoluta (f_i)	F. Relativa (f_{ri})	F. Abs/Rel Acum.
N° Apariciones elemento	$\frac{N^{\circ} \text{ apariciones elemento}}{\text{Elementos totales}}$	$\sum_{i=0}^{\text{indice elem.}} f_i/f_{ri}(\text{elemento}_i)$

MEDIA

Aritmética (\bar{x}_a)	Geométrica (\bar{x}_g)	Armónica ($(\bar{x}_a)^{-1}$):
$\frac{\sum_{i=0}^n x_i}{n}$	$\bar{x}_g = \left(\prod_{i=1}^n x_i \right)^{\frac{1}{n}}$	$\frac{\sum_{i=0}^n \frac{1}{x_i}}{n}$

CLASES DE EQUIVALENCIA

Rango	Límite	Amplitud	Marca
$v_{sup} - v_{inf}$	<ul style="list-style-type: none"> Superior: l_{sup} Inferior: l_{inf} 	$l_{sup} - l_{inf}$	$(l_{sup} + l_{inf})/2$

MODA

La **moda** es el **dato que aparece más veces** en un conjunto.

MEDIDAS DE DISPERSIÓN

Desviación típica o estándar (σ)	Desviación media (s)	Varianza (σ^2)	CV. de Pearson
$\sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}}$	$\frac{\sum_{i=1}^n x_i - \bar{x} }{n}$	$\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$	$cv_{\bar{x}} = \frac{\sigma}{\bar{x}}$

TEOREMA DE TCHEBYCHEV

Pertenencia	Amplitud	Intervalo
$pertenencia = 1 - \frac{1}{k^2}$	$amplitud = k \cdot \sigma$	$[\bar{x} - amplitud, \bar{x} + amplitud]$

El porcentaje de datos *pertenencia* se encontrará representado en el *intervalo*.

MEDIDAS DE ORDENACIÓN

Cuartil	Decil	Percentil	Centil
$n_c = \frac{1}{4}n_t$	$n_c = \frac{1}{10}n_t$	$n_c = \frac{1}{100}n_t$	$n_c = \frac{1}{n}n_t$

- Si n_c es un número natural:
 - Si n_t es un número impar: $Q_i = x_{nc}$
 - Si n_t es un número par: $Q_i = \frac{x_{nc\downarrow} + x_{nc\uparrow}}{2}$
- Si n_c es un número real: $Q_i = \frac{x_{nc\downarrow} + x_{nc\uparrow}}{2}$

REGLAS DE PROBABILIDAD

- $0 \leq P(A) \leq 1$
- $P(\bar{A}) = 1 - P(A)$
- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
- Si $A \subset B \rightarrow P(A) \leq P(B)$
- $P(E) = 1$
- $P(\emptyset) = 0$

REGLAS DE CONJUNTOS

- $B \cup C = \{\text{elementos}(B) + \text{elementos}(C)\}$
- $B \cap C = \{\text{elementos compartidos}(B, C)\}$
- $B - C = \{\text{elementos}(B) - \text{elementos compartidos}(B, C)\}$
- $B \Delta C = \{\text{elementos}(B) + \text{elementos}(C) - \text{elementos compartidos}(B, C)\}$

TEMA 2 – ASOCIACIÓN

SOPORTE Y CONFIANZA

Suporte	Confianza
$S(C) = \frac{\text{numero de apariciones}(C)}{\text{tamaño del espacio muestral}}$	$C(C_1 \rightarrow C_2) = \frac{\text{numero de apariciones}(C_1 \cup C_2)}{\text{numero de apariciones}(C_1)}$

APRIORI

1. Se calcula soporte de los sucesos elementales de ese espacio muestral, y se descartan sucesos elementales que no superen el umbral de soporte.
2. Se generan sucesos candidatos empleando Apriori-Gen. Partiendo de los sucesos elementales, para generar casos de tamaño $k+1$, se toman sucesos de tamaño k que compartan $k-1$ elementos.
3. Se calcula soporte para los sucesos candidatos generados, y se descartan aquellos que no superen el umbral de soporte.
4. Se generan todas las asociaciones (permutaciones) posibles partiendo de los sucesos candidatos. Para cada asociación, se calcula la confianza y se descartan las asociaciones que no superen el umbral de confianza.
5. Nos quedamos con los casos no descartados.

PROPIEDADES

$\text{soporte}(\text{caso}) \geq \text{soporte}(\text{subcaso que compone el caso})$
$\begin{aligned} \text{Si } C(A \rightarrow (B - A)) \leq \text{umbral de confianza, entonces} \\ C(A' \rightarrow (B - A')) \leq \text{umbral de confianza} \\ \text{donde } A' \subseteq A \end{aligned}$

TEMA 3 – CLASIFICACIÓN SUPERVISADA

ÁRBOLES DE DECISIÓN – ALGORITMO DE HUNT

1. Obtener suceso elemental en nodo inicial (no puede ser el clasificador).
2. Clasificar sucesos en base a nodo elegido.
 - a. Si podemos clasificar todos, hemos finalizado.
 - b. Si no, obtener suceso elemental en nodo intermedio y ejecutar paso 2.

GANANCIA DE INFORMACIÓN

Cálculo de la impureza		
$\Delta Impureza = I_{padre} \cdot freq_{padre} - I_{hijos} \cdot freq_{hijos}$		
$I_{hijos} = \sum_{j=1}^k \frac{N_{sucesos}(Nodo_j)}{N_{total\ sucesos}} I(Nodo_j)$		
Entropía	Error	Gini
$-\sum(freq_i(nodo_i) \cdot \log_2(freq_i(nodo_i)))$	$1 - \max(freq_i(nodo_i))$	$1 - \sum(freq_i(nodo_i))^2$

REGRESIÓN – COVARIANZA Y CORRELACIÓN

Covarianza	Correlación [0,1]
$S_{xy} = \frac{(\sum_{i=1}^n x_i \cdot y_i)}{n} - (\bar{x} \cdot \bar{y})$	$r_{xy} = \frac{S_{xy}}{s_x \cdot s_y}$
<ul style="list-style-type: none"> $S_{xy} > 0 \rightarrow$ Dependencia directa positiva. $S_{xy} = 0 \rightarrow$ No existe relación lineal. $S_{xy} < 0 \rightarrow$ Dependencia inversa negativa. 	<ul style="list-style-type: none"> $[0,1] \rightarrow$ relación ascendente. $0 \rightarrow$ no hay correlación. $[-1,0) \rightarrow$ relación descendente. <p>Se consideran fuertemente correlacionadas a partir de 0'8.</p>

REGRESIÓN POR MÍNIMOS CUADRADOS

$y = a + bx$	<ul style="list-style-type: none"> $a = \bar{y} - b \cdot \bar{x}$ $b = \frac{s_{xy}}{(s_x)^2}$
--------------	---

ANÁLISIS REGRESIÓN

$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})$	<ul style="list-style-type: none"> \hat{y}_i = valor calculado con la recta de regresión, tomando su x_i asociado. \bar{y} = media de y.
$SSy = \sum_{i=1}^n (y_i - \bar{y})$	<ul style="list-style-type: none"> y_i = valor original del dato. \bar{y} = media de y.
$r^2 = \frac{SSR}{SSy}$	<ul style="list-style-type: none"> $0 \leq r \leq 1$; Conforme más se aproxime a 1, más correcta es la recta de regresión.

ANÁLISIS DEL ERROR

$Sr = \sqrt{\frac{\sum_{i=1}^n ((y_i - \hat{y}_i)^2)}{n}}$	<ul style="list-style-type: none"> \hat{y}_i = valor calculado con la recta de regresión, tomando su x_i asociado. y_i = valor original del dato.
--	--

Cuanto más cercano a 0, mejor.