

VADEMECUM . ESTADISTIKA ETA DATUEN ANALISIA . II gn ZATIA

Asoziazio: aldagai kualitatiboen arteko erlazioz

abiapuntua: kontingentzia-taula

maizt. empiniko	a	b	c	a+b+c	} maizt. maizt
	d	e	f	d+e+f	
	a+d	b+e	c+f	n	

bastez maizt

independentzia: asoziazio-erlazio berdirak
% berdirak errenkadetan eta utabeetan

portzentajeak: asoziazio esploratzaileko

$$\left[\frac{a}{a+d}, \frac{d}{a+d} \right], \left[\frac{b}{b+e}, \frac{e}{b+e} \right], \dots$$

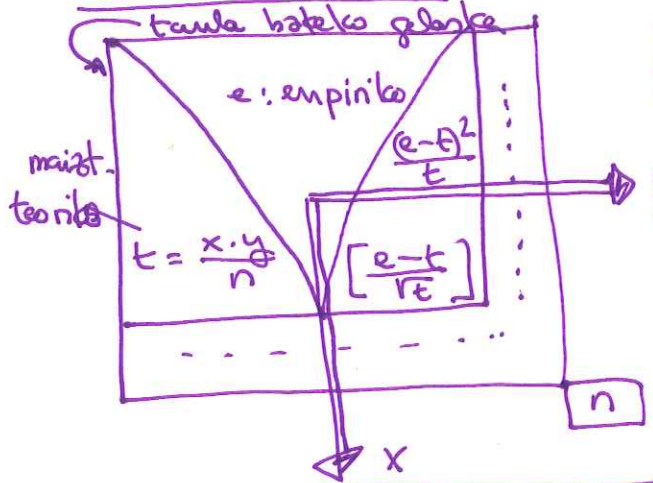
100

errenkadako aldagaiarenak

$$\left[\frac{a}{a+b+c}, \frac{b}{a+b+c}, \frac{c}{a+b+c} \right], \dots$$

utabeako aldagaiarenak.

Asoziatio-neurriak



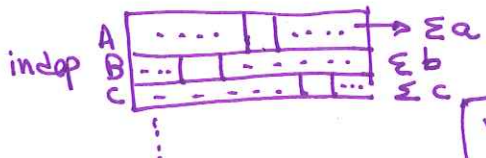
$$X^2 = \sum \frac{(e-t)^2}{t} \quad [0, \infty)$$

khi karratu gaitza gutxiak b ein da interpretatu

$\phi = \sqrt{\frac{X^2}{n}}$	0-1 balioak 2x2 taulak soilik
$C = \sqrt{\frac{X^2}{X^2+n}}$ kontingentzia koef	0-1 taula karratuak $C_{max} = \sqrt{\frac{m-1}{m}}$ taula es karratuak C/Cmax
$V = \sqrt{\frac{X^2}{n(m-1)}}$ Cramerren V	taula gutxiak

λ (lambda) (PRE neurria)

indep? indep. kontuan hartuta, zenbat errore? m (erroreok gutxiak)
dep? indep. kontuan hartuta, indep. A bada, zenbat errore? a
B bada, zenbat errore? b
C bada, zenbat errore? c



$$\lambda = \frac{m - (a+b+c+\dots)}{m}$$

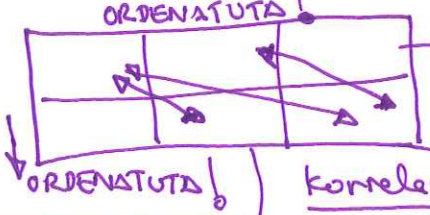
[zenbat errore gutxiago]

Goodman eta Kruskalen gamma

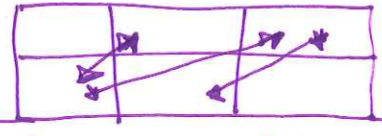
$$\gamma = \frac{k-d}{k+d}$$

int. $< 0 \rightarrow x \uparrow y \downarrow$
 $> 0 \rightarrow x \uparrow y \uparrow$

bi aldagai ordinal



k konkordantziak: gezi puntuak biertu eta bati



d diskordantziak: gezi puntuak biertu eta bati

Korrelezioz: aldagai kuantitatiboen arteko erlazioz

$$S_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n} = \frac{\sum xy}{n} - \bar{x}\bar{y}$$

kobariantza

$S_{xy} \in [-\infty, \infty]$
 $< 0 \rightarrow x \uparrow y \downarrow$
 $> 0 \rightarrow x \uparrow y \uparrow$
L sendatasuna ez.

$$r_{xy} = \frac{S_{xy}}{S_x S_y}$$

korelatio koef lineala

x	y	xy	x ²	y ²
...

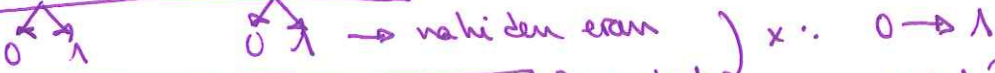
$r_{xy} \in [-1, 1]$

asoc. neurriak, $|\gamma|$ eta $|r_{xy}|$ nola interpretatu?
0-0.3 \rightarrow asoz edo kor gutxiak
0.3-0.6 \rightarrow asoz edo kor ertainak
 $> 0.6 \rightarrow$ asoz bir sendoa

r_{xy} dikotomiko - kuantitatibo



r_{xy} dikotomiko - dikotomiko



Korrelazio partziala:

$$r_{12.3} = \frac{r_{12} - r_{13} r_{23}}{\sqrt{1 - r_{13}^2} \sqrt{1 - r_{23}^2}}$$

bastertzen den aldagaia.

r_{xy} int } x: 0 to 1
y: 0 to 1? x up y up
1 to 0? x up y down

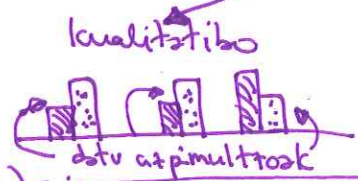
aldagaiak korrelazioa:
sendas baino fundamentu gabes

Kuantitatibo - kualitatibo erlazio estatistikoa

independentearren arabera ~~multzoak~~ erin eta horien arabera berstenen datu-azpimultzoak sakertu.

Korrelazio matritza

- simetrikoa
- diagonal 1
- > $x = c(\dots)$
- > $y = c(\dots)$
- > $\text{data frame}(x, y) = \text{itera}$
- > $\text{cor}(\text{itera}, \text{itera})$



kuantitatibo } histograma
relakoa den: puntu dispersio
karr - diagrama
 \bar{x}

ERREGRESIOA

1) datuak: x, y
zuzena: $\hat{y} = a + bx$

x	y	xy	x ²
Σ	Σ	Σ	Σ

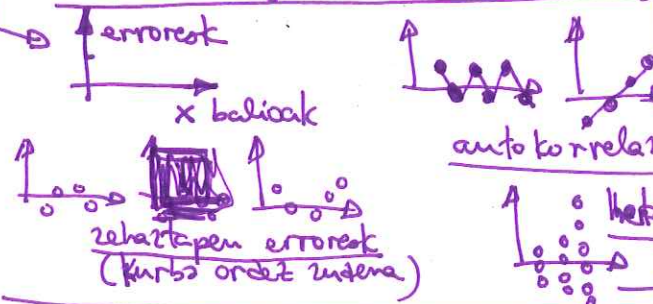
2) $e = y - \hat{y}$ erroreak
PROP: $\Sigma e = 0$
 $\Sigma y = \Sigma \hat{y}$
 $R: \text{zuzena} = \text{lm}(y \sim x)$

poikantziaren egokitasuna } zuzena \$ residuals
zuzena \$ fitted values

3) PROP: $S_y^2 = S_{\hat{y}}^2 + S_e^2$
Totala } azaldu
azaldu } azaldu gabes
 $R^2 = \frac{S_{\hat{y}}^2}{S_y^2} = 1 - \frac{S_e^2}{S_y^2}$
mugakoeff [0,1] } > 0,6 -> egokitasun handia

xy	xy	x ²	y ²
Σ	Σ	Σ	Σ

4) Erroren diagramak (diagnostikoa)



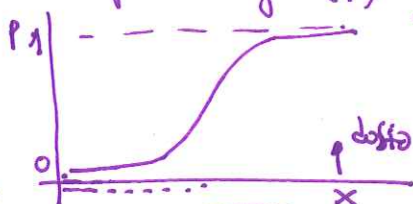
heteroskedastikotasuna } σ^2 berdina
 σ^2 berdina } homoskedastikotasuna

logit eredu (erregresio logistikoa)

aldagi dep: probabilitates edo porzentajes (p)
indep: doria

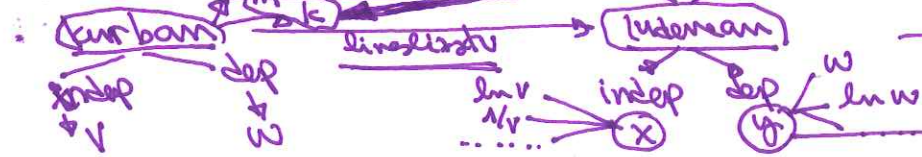
$$\ln \frac{p}{1-p} = a + bx$$

↓ doria



datuen joera alderantzizkoa bada eztekoaren edo aurtekoaren probabilitatea erin.

datuak es linealak baldin a b



R^2 beti x eta y-rekin db