

Normalizzazione

TRASFORM. TRA 0 e 1

$(0, 1)$
↓ ↓
c d

$$x' = \frac{x-a}{b-a}$$

TRASF. TRA -1 e 1

$(-1, 1)$
↓ ↓
c d

$$x' = -1 + \frac{x-a}{b-a} \cdot 2 = 2 \frac{x-a}{b-a} - 1$$

scalatura + traslazione del caso (0, 1)

STANDARDIZZAZIONE

$$x' = \frac{x - \bar{x}}{s_x}$$

$$\bar{x}' = \frac{\bar{x} - \bar{x}}{s_x} = 0$$

$$\left. \begin{aligned} s_{x'}^2 &= \frac{1}{s_x^2} \quad s_x^2 = 1 \\ s_{x'} &= 1 \end{aligned} \right\} \begin{array}{l} \text{la traslazione} \\ \text{di } -\bar{x} \text{ non} \\ \text{ha effetto} \end{array}$$

Non è tanto diversa dalla normalizzazione

$x' = \log x$ Si passa da dei dati che aumentano di tipo esponenziale a un qualcosa di lineare

(Scala logaritmica)

LATEX TBD chd

LEZIONE 08 - 25/03/2025

ANOVA (ANALISI DELLA VARIANZA)

Tecnica che si utilizza su popolazione ~~suddivisa~~ suddivisa in sottopopolazioni.

Si abbiano n individui suddivisi in G gruppi. Vengano queste osservazioni ordinate in base ai gruppi: (Si lavora in G gruppi indipendentemente dalla loro composizione)

x_i → gruppo

x_j → osserv. all'interno del gruppo

n_i ampiezza del gruppo i

$x_1^1, x_2^1, \dots, x_{n_1}^1, x_1^2, \dots, x_{n_2}^2, \dots, x_1^G, \dots, x_{n_G}^G$ → G ultimo gruppo

Media campionaria

$$\frac{1}{n_i} \sum_{j=1}^{n_i} x_j^i = \bar{x}^i \quad \text{Media camp. gruppo } i$$

$$n_i \bar{x}^i = \sum_{j=1}^{n_i} x_j^i$$

$$SS_T = \sum_{i=1}^G \sum_{j=1}^{n_i} (x_j^i - \bar{x})^2$$

Varianza rispetto alla media totale

$$\Rightarrow \sum_{i=1}^G n_i \bar{x}^i = \sum_{i=1}^G \sum_{j=1}^{n_i} x_j^i$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^G n_i \bar{x}^i \quad \text{Media totale}$$

SUM OF SQUARE TOTAL

$$\frac{1}{n-1} SS_T \quad \text{VAR. TOTALE}$$

$$SS_W = \sum_{i=1}^G \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2 \rightarrow \frac{1}{n-G} SS_W \quad \text{VARIANZA ENTRO I GRUPPI}$$

SUM OF SQUARES WITHIN

$$SS_B = \sum_{i=1}^G n_i (\bar{x}_i - \bar{x})^2$$

S. OF S. BETWEEN

Si calcola la varianza tra le medie campionarie dei gruppi. Gruppi con dentro + individui hanno un po' maggiore (x questo si moltiplica per n_i lo scarto)

$$\frac{1}{G-1} SS_B \quad \text{VARIANZA TRA I GRUPPI}$$

$$SS_T = SS_W + SS_B \quad \text{è dimostrabile, non lo faremo. (c'è nelle note)}$$

Vale sempre

$$\frac{SS_T}{n-1} = \frac{SS_W}{n-1} + \frac{SS_B}{n-1} = \frac{n-G}{n-1} \frac{SS_W}{n-G} + \frac{G-1}{n-1} \frac{SS_B}{G-1}$$

$$\underbrace{\frac{1}{n-1} \cdot SS_T}_{\text{VAR. TOTALE}} = \frac{n-G}{n-1} \cdot \underbrace{\frac{1}{n-G} SS_W}_{\text{VAR. ENTRO I GRUPPI}} + \frac{G-1}{n-1} \cdot \underbrace{\frac{1}{G-1} SS_B}_{\text{VAR. TRA I GRUPPI}}$$

LA VAR. TRA I GRUPPI, se è trascurabile, significa che VAR. TOT. e VAR. ENTRO I GRUPPI sono approssimativamente simili:

Se non è trascurabile, in uno dei gruppi, considerando quella risorsa, c'è una differenza significativa...

La divisione in gruppi è fornita a priori (non bisogna trovare una combinazione che rispetti le nostre aspettative, abbiamo i dati divisi in gruppi e li analizziamo)

(Test statistici)

CLASSIFICATORI

Automatismo che ci permette di classificare degli oggetti (es. ALBERI DI DECISIONE)

Bisogna valutare se questi classificatori siano buoni o meno. Sono dei predittori.

Si parlerà di classificatori binari x comodità.

Classe positiva / classe negativa.

Si organizzano come dalle matrici:

- RIGHE: VALORE EFFETTIVO Positivo Negativo
- COLONNE: PREDEZIONE Positivo Negativo

VALORE EFFETTIVO

	P	N
P	True Positive	False Positive
N	False Negative	True Negative

- Se si somma sulla prima colonna si ha la cardinalità degli elementi realmente positivi
- Idem x la seconda (negativi)
- Se si somma tutto si ottiene la cardinalità del campione

ACCURATEZZA: $\frac{TP + TN}{TP + FP + FN + TN}$ $\frac{\# \text{ PREDEZIONI CORRETTE}}{\# \text{ CASI TOTALI}}$

→ 1 BUONO
→ 0 NON ACCURATO
TENDENZA

Metrica naturale e facile da interpretare. In alcuni casi è però poco significativa. In una situazione con uno sbilanciamento forte nei valori effettivi dei casi, l'accuratezza perde di utilità (esempio se si ha 1 solo positivo ogni 10k) ⊕

Bisogna vedere quanto si è bravi a

- prevedere i positivi guardando solo i positivi
 - " i negativi " " negativi "
- } vanno fatte entrambe

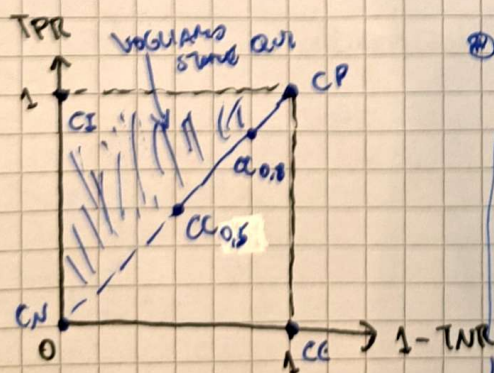
Si fanno 2 verifiche separate:

True Positive Rate = $\frac{TP}{TOT P}$ DA 0 A 1 → BUONO A PREDIRE I POSITIVI

True Negative Rate = $\frac{TN}{TOT N}$

[TPR detto anche recall o sensitivity
TNR detto anche specificity]

Considerando entrambe TPR e TNR si possono fare delle considerazioni sul nostro classificatore.



① CP classificatore costante positivo

	P	N
P	TOT P	0
N	0	0

1 - TNR = 1 - 0 = 1
TPR = 1

CN

	P	N
P	0	0
N	TOT P	TOT N

1 - TNR = 1 - 1 = 0
TPR = 0

② CI classificatore ideale

	P	N
P	0	0
N	0	TOT N

1 - TNR = 1 - 1 = 0
TPR = 1

Sulla diagonale

CC_{0,5} classifikatore casuale con 0,5 di probabilità

$\frac{1}{2} \text{TOTP}$	$\frac{1}{2} \text{TOTN}$
$\frac{1}{2} \text{TOTP}$	$\frac{1}{2} \text{TOTN}$

$$1 - \text{TNR} = \frac{1}{2}$$

$$\text{TPR} = \frac{1}{2}$$

Se si è sulla diagonale si ha un comportamento ~~equo~~ equiparabile a quello del lancio della moneta
→ poco significativo

CC_{0,8}

0,8TOTP	0,8TOTN
0,2TOTP	0,2TOTN

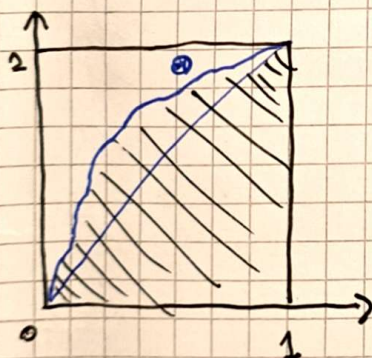
$$1 - \text{TNR} = 0,8$$

$$\text{TPR} = 0,8$$

• Se si ha un classifikatore che sta nella parte bassa, si inverte la predizione.

CLASSIFICATORE A SOGLIA Se sta sotto la soglia negativo, sopra positivo (esempio)

Come si fissa la soglia?



$$C(x) = \begin{cases} P & \text{se } x > T(\text{tau}) \\ N & \text{se } x \leq T \end{cases}$$

Se si prende come $T =$ valore minimo possibile nel campione si è nel caso CP

Se $T =$ valore max possibile nel campione si è nel caso CN

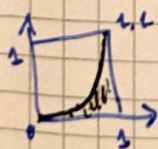
se T cresce gradualmente si ottiene ^{la curva}

Più l'area III è vicino a 1 (si tende a essere nel caso migliore), più si ottiene una curva detta ROC. L'AREA SOTTO LA CURVA ROC PRENDE IL NOME DI AUROC (AREA UNDER ROC)

↳ Serve a valutare la bontà di un classifikatore indipendentemente dalla soglia.
➢ AREA: MIGLIOR CLASSIFICAZIONE INDIPEND. DALLA SOGLIA SCELTA

Si prendano due classifikatori. Se la curva ROC di uno domina ~~l'altra~~ l'altra, allora quel classifikatore ha risultati migliori

SI PUÒ FARE IL DISCORSO INVERSO & COMPLEMENTARE



Se lo si specchia è buono

$B = \cancel{A} \cup (B \setminus A)$ sono disgiunti:

$$P(B) = P(A) + P(B \setminus A) \geq 0 \quad \text{ecc...}$$

ESERCIZIO LEZIONE 13

$$P(E|M) = 1 - P(\bar{E}|M) \quad \text{Dimostrazione}$$

$$P(A|M) = \frac{P(A \cap M)}{P(M)} \quad P(\bar{A}|M) = \frac{P(\bar{A} \cap M)}{P(M)}$$

Dato che $(A \cap M) \cup (\bar{A} \cap M) = M$ e i due eventi sono disgiunti:

$$P(M) = P((A \cap M) \cup (\bar{A} \cap M)) = P(A \cap M) + P(\bar{A} \cap M)$$

Dividendo entrambi i membri per $P(M)$:

$$\frac{P(A \cap M)}{P(M)} + \frac{P(\bar{A} \cap M)}{P(M)} = 1 \quad \Rightarrow \quad P(A|M) = 1 - P(\bar{A}|M)$$