CdL in Scienze Statistiche ed Economiche - Università degli Studi di Milano-Bicocca

Lezione: Distanze

Docente: Aldo Solari

1 Distanze

L'analisi di raggruppamento (*cluster analysis*) ha per scopo far emergere dall'insieme dei dati a disposizione gruppi di unità statistiche "simili" tra loro e "dissimili" da quelle degli altri gruppi. Che cosa si intende per unità statistiche "simili", o equivalentemente, "dissimili"? Dobbiamo quantificare con un numero la "diversità" tra due unità statistiche

• Variabili Quantitative: Diversità = Distanza (Metrica e Indice di Distanza)

• Variabili Qualitative: Diversità = Indice di Dissimilarità

1.1 La funzione distanza

Consideriamo misurazioni su p variabili tutte quantitative. Quanto sono "distanti" due unità statistiche u_i e u_l ? Dipende da come definiamo la "distanza".

In generale, una distanza è una funzione

$$d: \mathbb{R}^p \times \mathbb{R}^p \to \mathbb{R}$$

che associa ad una coppia di unità statistiche u_i e u_l un numero reale Proprietà di una distanza:

- (D1) Non negatività
- (D2) Identità $d(u_i, u_l) = 0 \Leftrightarrow u_i = u_l$

 $d(u_i, u_l) > 0$

- (D3) Simmetria $d(u_i, u_l) = d(u_l, u_i)$
- (D4) Disuguaglianza triangolare $d(u_i, u_l) \leq d(u_i, u_t) + d(u_t, u_l)$
- METRICA: valgono (D1), (D2), (D3) e (D4)
- INDICE DI DISTANZA: valgono (D1), (D2) e (D3)

 $\bullet \;$ Distanza Euclidea d_2 tra due unità statistiche $u_i' \;$ e $\; u_l' \;$ e $\; u_l' \;$ e $\; u_l' \;$

$$d_2(u_i, u_l) = \sqrt{\sum_{j=1}^{p} (x_{ij} - x_{lj})^2}$$

 d_2 soddisfa le proprietà (D1), (D2), (D3) e (D4), quindi è una metrica

• Distanza di Manhattan (o della città a blocchi) d_1 tra due unità statistiche u_i' e u_l' u_l' e u_l' u_l' e u_l' e u_l'

$$d_1(u_i, u_l) = \sum_{i=1}^{p} |x_{ij} - x_{lj}|$$

 d_1 soddisfa le proprietà (D1), (D2), (D3) e (D4), quindi è una metrica

 $\bullet\,$ Distanza di Lagrange d_{∞} tra due unità statistiche u_i' e u_l' e u_l' $_{1\times p}$ $_{1\times p}$

$$d_{\infty}(u_i, u_l) = \max_{j \in \{1...,p\}} |x_{ij} - x_{lj}|$$

 d_{∞} soddisfa le proprietà (D1), (D2), (D3) e (D4), quindi è una metrica

Example 1.1.
$$u_1' = \begin{bmatrix} 1 & 1 \end{bmatrix}, \ u_2' = \begin{bmatrix} 2 & 3 \end{bmatrix}$$
 $d_2(u_1, u_2) = \sqrt{(1-2)^2 + (1-3)^2} = \sqrt{5}$ $d_1(u_1, u_2) = |1-2| + |1-3| = 3$ $d_{\infty}(u_1, u_2) = \max\{|1-2|, |1-3|\} = 2$

1.2 Distanza di Minkowski

Distanza di Minkowski d_m (di ordine $m \geq 1$) tra due unità statistiche u_i' e u_l' 1 $\times p$ 1 1 $\times p$ 1

$$d_m(u_i, u_l) = \left[\sum_{j=1}^{p} |x_{ij} - x_{lj}|^m\right]^{1/m}$$

- Per $m \ge 1$, d_m soddisfa le proprietà (D1), (D2), (D3) e (D4), quindi è una metrica
- $\bullet\,$ Casi particolari: distanza Euclidea (m=2), di Manhattan (m=1) e di Lagrange ($m=\infty$)
- E' funzione non crescente dell'indice m

$$d_{m'}(u_i, u_l) \ge d_{m''}(u_i, u_l) \quad 1 \le m' < m''$$

1.3 Distanza di Canberra

Distanza di Canberra d_C tra due unità statistiche u_i' e u_l'

$$d_C(u_i, u_l) = \sum_{i=1}^{p} \frac{|x_{ij} - x_{lj}|}{|x_{ij} + x_{lj}|}$$

- Sono esclusi dalla somma i termini in cui il denominatore si annulla
- E' utilizzata per variabili quantitative non negative
- E' una versione pesata della distanza di Manhattan

$$d_C(u_i, u_l) = \sum_{i=1}^{p} w_j |x_{ij} - x_{lj}|$$

con pesi
$$w_j = \frac{1}{|x_{ij} + x_{lj}|}, j = 1, \dots, p$$

1.4 Distanza Euclidea al quadrato

• Distanza Euclidea al quadrato

$$d_2^2(u_i, u_l) = \sum_{j=1}^p (x_{ij} - x_{lj})^2$$

- La distanza Euclidea al quadrato d_2^2 soddisfa (D1), (D2) e (D3) ma non soddisfa (D4), quindi è un indice di distanza
- $\bullet\,$ Tuttavia, d_2^2 gode della proprietà di addittività (mentre la distanza Euclidea d_2 no):

per due insiemi K e Q tali che $K \cup Q = \{1, \dots, p\}$ e $K \cap Q = \emptyset$

$$d_2^2(u_i, u_l) = \sum_{j=1}^p (x_{ij} - x_{lj})^2 = \sum_{k \in K} (x_{ik} - x_{lk})^2 + \sum_{q \in Q} (x_{iq} - x_{lq})^2$$

Example 1.2.
$$u'_1 = [10 \ 5]$$
 $u'_2 = [13 \ 9]$

$$u'_{2} = \begin{bmatrix} 13 & 9 \end{bmatrix}^{1 \times 2}$$

$$u'_{3} = \begin{bmatrix} 11 & 7 \end{bmatrix}$$

$$u'_{3} = \begin{bmatrix} 11 & 7 \end{bmatrix}$$

$$d_{2}^{2}(u_{1}, u_{2}) = (10 - 13)^{2} + (5 - 9)^{2} = 25$$

$$d_{2}^{2}(u_{1}, u_{3}) = (10 - 11)^{2} + (5 - 7)^{2} = 5$$

$$d_{2}^{2}(u_{3}, u_{2}) = (11 - 13)^{2} + (7 - 9)^{2} = 8$$

$$25 = d_{2}^{2}(u_{1}, u_{2}) > d_{2}^{2}(u_{1}, u_{3}) + d_{2}^{2}(u_{3}, u_{2}) = 5 + 8 = 13$$

2 Distanza di Mahalanobis

Tra due unità statistiche u_i' e u_l' u_l' e u_l'

$$d_M(u_i, u_l) = \sqrt{\frac{(u_i - u_l)' S^{-1}(u_i - u_l)}{\sum_{p > p} \sum_{p > 1} u_i}}$$

Tra l'i-sima unità statistica u_i' e il baricentro \bar{x}' : $_{1 \times p}$

$$d_M(u_i, \bar{x}) = \sqrt{\frac{(u_i - \bar{x})' S^{-1}(u_i - \bar{x})}{1 \times p}}$$

L'insieme dei punti p-dimensionali $u'_{1\times p}$ con distanza di Mahalanobis costante c>0 dal baricentro $\bar x'_{1\times p}$ soddisfano l'equazione

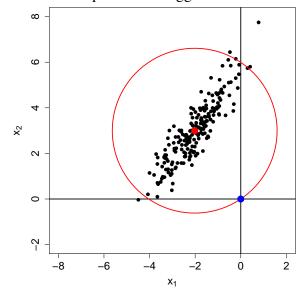
$$(u - \bar{x})' S^{-1}(u - \bar{x}) = c^2$$

che definisce un iper-ellissoide

Si noti che l'insieme di punti p-dimensionali $u'_{1\times p}$ con distanza Euclidea costante c>0 dal baricentro $\bar{x}'_{1\times p}$ soddisfano l'equazione

$$(u - \bar{x})'(u - \bar{x}) = c^2$$
_{1×p}
_{p×1}

che definisce una ipersfera di raggio c dal baricentro



2.1 Distanza di Mahalanobis e valori anomali

Se si può assumere che le righe della matrice X sono realizzazioni indipendenti generate dalla medesima distribuzione Normale p-variata, possiamo definire l'i-sima unità statistica u_i' un outlier $1 \times p$

se

$$d_M^2(u_i, \bar{x}) > q_{0.95}$$

dove $q_{0.95}$ è il 0.95 quantile di una distribuzione χ^2 con p gradi di libertà

p	$q_{0.95}$
1	3.8415
2	5.9915
3	7.8147
5	11.0705
10	18.3070
100	124.3421

Se le n unità statistiche sono realizzazioni indipendenti generate dalla medesima distribuzione Normale p-variata

valore atteso di outliers = $n \times 0.05$

Qualora osserviamo un numero sostanzialmente più elevato di quello atteso, abbiamo un eccesso di *outliers*

Example 2.1. *Dati Animals*

i		$d_M^2(u_i,\bar{x})$	i		$d_M^2(u_i,\bar{x})$
1	Mountain beaver	3.339	15	African elephant	5.445
2	Cow	1.529	16	Triceratops	22.392
3	Grey wolf	0.134	17	Rhesus monkey	4.12
4	Goat	0.328	18	Kangaroo	0.014
5	Guinea pig	3.533	19	Golden hamster	8.51
6	Dipliodocus	26.092	20	Mouse	14.908
7	Asian elephant	2.963	21	Rabbit	2.243
8	Donkey	0.343	22	Sheep	0.082
9	Horse	1.348	23	Jaguar	0.169
10	Potar monkey	2.087	24	Chimpanzee	0.795
11	Cat	2.573	25	Rat	6.249
12	Giraffe	1.346	26	Brachiosaurus	38.629
13	Gorilla	0.423	27	Mole	11.303
14	Human	2.084	28	Pig	0.782

3 Distanze e trasformazioni lineari

3.1 Invarianza di d_M rispetto alle trasformazioni lineari

La trasformazione lineare dell' i-sima unità statistica $u_i' = x_i'$ $1 \times p$ $1 \times p$

$$y_i' = x_i' A' + b'_{1 \times p}$$

è definita da

- la matrice $A_{p \times p}$
- il vettore $b \atop p \times 1$

La matrice dei dati linearmente trasformati risulta

$$Y_{n \times p} = X_{n \times pp \times p} A' + 1_{n \times 11 \times p} b'$$

con con vettore delle medie e matrice di varianze/covarianze

$$\bar{y}_{p \times 1} = A \bar{x}_1 + b_1, \qquad S^Y_{p \times p} = A S^X A'_{p \times pp \times pp \times p}$$

Si noti che ci stiamo limitando a considerare trasformazioni lineari delle unità statistiche da \mathbb{R}^p a \mathbb{R}^p , ovvero trasformazioni lineari che non riducono la dimensionalità originale p.

Proposition 3.1. Siano $y_i' = x_i' A' + b' e y_l' = x_l' A' + b' con A non singolare. La distanza di Mahalanobis <math>d_m$ è invariate rispetto alle trasformazioni lineari (non singolari):

$$d_{M}(y_{i}, y_{l}) = \sqrt{(y_{i} - y_{l})'[S^{Y}]^{-1}(y_{i} - y_{l})}$$

$$= \sqrt{(Ax_{i} + b - Ax_{l} - b)'[AS^{X}A']^{-1}(Ax_{i} + b - Ax_{l} - b)}$$

$$= \sqrt{[A(x_{i} - x_{l})]'[AS^{X}A']^{-1}[A(x_{i} - x_{l})]}$$

$$= \sqrt{(x_{i} - x_{l})'A'[A']^{-1}[S^{X}]^{-1}[A]^{-1}A(x_{i} - x_{l})}$$

$$= d_{M}(u_{i}, u_{l})$$

 $ricordando\ che\ (AB)^{-1} = B^{-1}A^{-1}$

3.2 Invarianza di d_m rispetto alle traslazioni

Una traslazione è un caso particolare di trasformazione lineare con

$$\bullet \ \ \underset{p \times p}{A} = \underset{p \times p}{I}$$

• $b'_{1 \times p}$ arbitraria

Traslazione della matrice dei dati X

$$Y_{n \times p} = X_{n \times p} + \underset{n \times 11 \times p}{1}b'$$

con vettore delle medie e matrice di varianze/covarianze

$$\bar{y}_{p \times 1} = \bar{x}_{p \times 1} + b_{p \times 1}, \qquad S^{Y}_{p \times p} = S^{X}_{p \times p}$$

Proposition 3.2. Siano $y_i' = x_i' + b'_{1 \times p} e y_l' = x_l' + b'_{1 \times p} La distanza di Minkowski <math>d_m$ è invariate rispetto alle traslazioni:

$$d_{m}(y_{i}, y_{l}) = \left[\sum_{j=1}^{p} |y_{ij} - y_{lj}|^{m}\right]^{1/m}$$

$$= \left[\sum_{j=1}^{p} |(x_{ij} + b_{j}) - (x_{lj} + b_{j})|^{m}\right]^{1/m}$$

$$= \left[\sum_{j=1}^{p} |x_{ij} - x_{lj}|^{m}\right]^{1/m}$$

$$= d_{m}(u_{i}, u_{l})$$

3.3 Invarianza di d_2 rispetto alle trasformazioni ortogonali

- $\bullet \quad b'_{1\times p} = 0_{1\times p}$

Trasformazione ortogonale della matrice dei dati X

$$\underset{n \times p}{Y} = \underset{n \times pp \times p}{X} A'$$

con vettore delle medie e matrice di varianze/covarianze

$$\bar{y}_{p\times 1} = \underset{p\times pp\times 1}{A}\bar{x}, \quad S^Y_{p\times p} = \underset{p\times pp\times pp\times p}{A}S^XA'$$

Esempi di trasformazioni ortogonali sono:

• Trasformazione identità: $A = I_{p \times p} = I_{p \times p}$

• Permutazione: A è una matrice di permutazione che si ottiene permutando le righe (o le colonne) della matrice identità I. In due dimensioni, la seguente matrice di permutazione comporta scambiare l'ordine due delle colonne di X:

$$A_{2\times 2} = \left[\begin{array}{cc} 0 & 1\\ 1 & 0 \end{array} \right]$$

• Rotazione: A è una matrice di rotazione, ovvero A ortogonale con $\det(A) = 1$ o -1. In due dimensioni, la seguente matrice di rotazione comporta una rotazione antioraria di angolo θ radianti intorno all'origine:

$$A_{2\times 2} = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}$$

Proposition 3.3. Siano $y'_i = x'_i A'_i e y'_l = x'_l A'_i con A_{p \times p}$ matrice ortogonale. La distanza Euclidea d_2 è invariate rispetto alle trasformazioni ortogonali:

$$d_{2}(y_{i}, y_{l}) = \sqrt{(y_{i} - y_{l})'(y_{i} - y_{l})}$$

$$= \sqrt{(Ax_{i} - Ax_{l})'(Ax_{i} - Ax_{l})}$$

$$= \sqrt{[A(x_{i} - x_{l})]'[A(x_{i} - x_{l})]}$$

$$= \sqrt{(x_{i} - x_{l})'A'A(x_{i} - x_{l})}$$

$$= \sqrt{(x_{i} - x_{l})'A^{-1}A(x_{i} - x_{l})}$$

$$= \sqrt{(x_{i} - x_{l})'(x_{i} - x_{l})}$$

$$= d_{2}(x_{i}, x_{l})$$

3.4 Distanza Euclidea calcolata su \tilde{X} , Z e \tilde{Z}

• $\tilde{x}'_i = (u_i - \bar{x})'$ è l'i-sima riga di $\tilde{X}_{n \times p} = \underset{n \times nn \times p}{HX}$

$$d_2(\tilde{x}_i, \tilde{x}_l) = \sqrt{\frac{(u_i - u_l)'(u_i - u_l)}{\sum_{l \in P} u_l + \sum_{l \in P} u_l}} = d_2(u_i, u_l)$$

• $z_i' = (u_i - \bar{x})' D^{-\frac{1}{2}}_{p \times p}$ è l'i-sima riga di $Z_{n \times p} = H X_n D^{-\frac{1}{2}}_{p \times p}$

$$d_2(z_i, z_l) = \sqrt{\frac{(u_i - u_l)' D^{-1}(u_i - u_l)}{\sum_{p \times p}^{p} (u_i - u_l)}} = \sqrt{\sum_{j=1}^{p} \frac{1}{s_{jj}} (x_{ij} - x_{lj})^2}$$

• $\tilde{z}_i' = (u_i - \bar{x})' S^{-\frac{1}{2}}$ è l'i-sima riga di $\tilde{Z}_n = H X S^{-\frac{1}{2}} X S^{-\frac{1}{2}}$

$$d_2(\tilde{z}_i, \tilde{z}_l) = \sqrt{\frac{(u_i - u_l)' S^{-1}(u_i - u_l)}{\sum_{\substack{1 \le p \ p \le p}} (u_i - u_l)}} = d_M(u_i, u_l)$$

4 Indici di similarità

Consideriamo misurazioni su p variabili, qualitative e/o quantitative. Ciascuna unità statistica presenta misurazioni appartenenti allo spazio campionario $\mathcal{X} = \mathcal{X}_1 \times \ldots \times \mathcal{X}_p$. Ad esempio, se abbiamo p = 2 variabili, Sesso e Posizione geografica, lo spazio campionario è:

$$\mathcal{X} = \mathcal{X}_{Sesso} \times \mathcal{X}_{Pos.Geog.} = \{M, F\} \times \{Nord, Centro, Sud\} = \{(M, Nord), (F, Nord), (M, Centro), (F, Centro), (M, Sud), (F, Sud)\}$$

4.1 Indice di similarità

In generale, un indice di similarità è una funzione

$$s: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$$

che associa ad una coppia di unità statistiche u_i' e u_l' un numero reale. Un indice di similarità soddisfa le seguenti proprietà:

(S1) Non negatività

$$s(u_i, u_l) > 0$$

(S2) Normalizzazione

$$u_i = u_l \Rightarrow s(u_i, u_l) = 1$$

(S3) Simmetria

$$s(u_i, u_l) = s(u_l, u_i)$$

dove 1 è il massimo valore assumibile dall'indice di similarità. Un indice di dissimilarità è definito come

$$d(u_i, u_j) = 1 - s(u_i, u_j)$$

e soddisfa (D1) e (D3)

4.2 Variabili binarie

Supponiamo che il profilo dell'i-esima unità statistica u'_i sia composto di sole variabili binarie (o dicotomiche), codificate per comodità come 0 e 1

$$X_{n \times p} = \begin{bmatrix} 0 & 0 & \cdots & 1 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 1 & 0 & \cdots & 1 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 1 & 0 & \cdots & 0 & \cdots & 1 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 1 & 1 & \cdots & 1 & \cdots & 1 \end{bmatrix} = \begin{bmatrix} u'_1 \\ \cdots \\ u'_i \\ \cdots \\ u'_l \\ \cdots \\ u'_n \end{bmatrix}$$

Possiamo costruire, per ciascuna coppia u_i' e u_l' , la seguente tabella di contingenza

	uni	t l	
unit i	1	0	
1	a	b	a+b
0	c	d	c+d
	a+c	b+d	p = a + b + c + d

dove

- a è la frequenza di variabili binarie con valore 1 per l'unità i e valore 1 per l'unità l
- b è la frequenza di variabili binarie con valore 1 per l'unità i e valore 0 per l'unità l
- etc.

Example 4.1.

$$\begin{bmatrix} 1 & 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} u'_1 \\ u'_2 \end{bmatrix}$$

$$\begin{array}{c|c} u_2 \\ \hline 1 & 2 & 1 & 3 \\ \hline 0 & 1 & 1 & 2 \\ \hline 3 & 2 & p = 5 \end{array}$$

4.2.1 Indice di corrispondenza e di Jaccard

Consideriamo 1 come 'presenza' e 0 come 'assenza'. Non è ovvio se la contemporanea presenza 1-1 o la contemporanea assenza 0-0 siano egualmente indicativi di somiglianza. Ad esempio, se le unità sono individui e la variabile binaria è "capelli castani (1)/capelli non castani (0)" la contemporanea presenza 1-1 è indubbiamente indicativa di somiglianza, non così la contemporanea assenza 0-0. Si parla in questo caso di variabile binaria *asimmetrica*. Per contro se la variabile binaria è "maschio (1)/femmina (0)" la contemporanea assenza 0-0 ha lo stesso valore della contemporanea presenza 1-1. Si parla in questo caso di variabile binaria *simmetrica*.

Indice di corrispondenza semplice

$$s_c(u_i, u_l) = \frac{a+d}{p}$$

considera allo stesso modo co-presenze 1-1 e co-assenze 0-0, quindi è opportuno per variabili binarie simmetriche

Indice di Jaccard

$$s_J(u_i, u_j) = \frac{a}{a+b+c}$$

ignora le coassenze 0-0 (ed è indeterminato se d=p), quindi è opportuno per variabili binarie asimmetriche

Per l'esempio precedente abbiamo

$$s_c(u_1, u_2) = \frac{3}{5} = 0.6, \quad s_J(u_1, u_2) = \frac{2}{4} = 0.5$$

Example 4.2.

$$\begin{bmatrix} 1 & 0 & 1 & 1 & 1 \\ 1 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} u_1' \\ u_2' \\ u_3' \end{bmatrix}$$

Per ciascuna coppia di osservazioni calcoliamo la tabella di contingenza, ottenendo le tre tabelle

Si noti che u_1 è equi-somigliante a u_2 e u_3 secondo s_c , mentre è più somigliante a u_2 che a u_3 secondo s_J , questo poichè la co-assenza che lo accomuna a u_3 non ha peso nell'indice di Jaccard.

4.3 Variabili qualitative nominali

Se tutte le variabili sono qualitative nominali (factor in R), possiamo considerare come indice di corrispondenza semplice la proporzione di variabili in cui le due unità u_i' e u_j' assumono la stessa modalità

$$s_c(u_i, u_j) = \frac{\sum_{j=1}^p I\{x_{ij} = x_{lj}\}}{p}$$

dove $I\{\cdot\}$ rappresenta la funzione indicatrice

4.4 Variabili qualitative ordinali

Variabili qualitative ordinali (Ord.factor in R) con modalità ordinate, ad esempio, mai \prec qualche volta \prec spesso \prec sempre.

Trattare queste variabili come qualitative non ordinate, sebbene possibile, fa perdere l'informazione relativa all'ordinamento delle modalità (mai e qualche volta sono misurate egualmente 'distanti' di mai e sempre).

Se la j-sima variabile è qualitativa ordinale, una soluzione alternativa consiste nel trasformare le m_i modalità ordinate nei corrispondenti numeri interi da 1 a m_j normalizzando il risultato:

$$y_{ij} = \frac{\text{punteggio}(x_{ij}) - 1}{m_i - 1}$$

e trattare la *j*-sima variabile come quantitativa

In questo caso si assume che le 'distanze' tra le categorie ordinate sono le stesse Ad esempio

4.5 Variabili miste: indice di Gower

$$s_G(u_i, u_l) = \frac{\sum_{j=1}^{p} \delta_{il}(j) s_{il}(j)}{\sum_{j=1}^{p} \delta_{il}(j)}$$

dove

$$s_{il}(j) = \begin{cases} 1 - \frac{|x_{ij} - x_{lj}|}{\text{range } j\text{-sima variabile}} & \text{se } j\text{-sima variabile quantitativa} \\ I(x_{ij} = x_{lj}) & \text{se } j\text{-sima variabile binaria/nominale} \\ 1 - |y_{ij} - y_{lj}| & \text{se } j\text{-sima variabile ordinale} \end{cases}$$

$$\delta_{il}(j) = \begin{cases} 1 & i,l \text{ confrontabili rispetto } j\text{-sima variabile} \\ 0 & i,l \text{ non confrontabili rispetto } j\text{-sima variabile} \end{cases}$$

dove due unità sono non confrontabili rispetto alla j-sima variabile se c'è un valore mancante in almeno una delle due o se la j-sima variabile è binaria asimmetrica e si ha co-assenza 0-0.

4.6 Matrice delle distanze/dissimilarità

A $\underset{n \times p}{X}$ si associa una matrice $\underset{n \times n}{D}$ delle distanze/dissimilarità tra le n unità statistiche

$$D_{n \times n} = \begin{bmatrix} 0 & d_{12} & \cdots & d_{1i} & \cdots & d_{1n} \\ & 0 & \cdots & d_{2i} & \cdots & d_{2n} \\ & & \ddots & \vdots & & \vdots \\ & & 0 & \cdots & d_{in} \\ & & & \ddots & \vdots \\ & & & 0 \end{bmatrix}$$

dove

- $\bullet \ d_{il} = d(u_i, u_l)$
- $d_{il} = d_{li}$ (la matrice è simmetrica)
- $d_{ii} = 0$