

Similarità, distanza, associazione

Misure di similarità

- Variano da 1 (massima somiglianza, osservazioni identiche) a 0 (nessuna somiglianza, osservazioni completamente diverse)
- Possono essere simmetriche (l'assenza di una specie è considerata informativa) o asimmetriche (l'assenza non è un dato certo)
- Se trasformate in dissimilarità ($D=1-S$), possono godere di proprietà metriche o meno

| | | | | | | |
|-----------------------|---|-----------------------|----------|--------|-------|-------|
| | | Osservazione <i>j</i> | | | St. A | St. B |
| | | 1 | 0 | Sp. 1 | 3 | 0 |
| Osservazione <i>k</i> | 1 | <i>a</i> | <i>b</i> | Sp. 2 | 4 | 2 |
| | 0 | <i>c</i> | <i>d</i> | Sp. 3 | 0 | 0 |
| | | $p = a + b + c + d$ | | Sp. 4 | 2 | 5 |
| | | | | Sp. 5 | 1 | 16 |
| | | | | Sp. 6 | 0 | 4 |
| | | | | Sp. 7 | 12 | 5 |
| | | | | Sp. 8 | 0 | 1 |
| | | | | Sp. 9 | 0 | 4 |
| | | | | Sp. 10 | 1 | 0 |

| | |
|---------|---------|
| $a = 4$ | $b = 3$ |
| $c = 2$ | $d = 1$ |

Alcune misure di similarità

simmetriche

concordanza semplice

$$S_{jk} = \frac{a+d}{p}$$

Rogers & Tanimoto

$$S_{jk} = \frac{a+d}{a+2b+2c+d}$$

asimmetriche

Jaccard

$$S_{jk} = \frac{a}{a+b+c}$$

Sørensen

$$S_{jk} = \frac{2a}{2a+b+c}$$

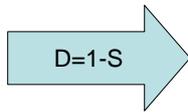
Gower

$$S_{jk} = \frac{\sum_{i=1}^p w_i S_i}{\sum_{i=1}^p w_i}$$

- per descrittori binari $s_i=1$ nei casi di concordanza e $s_i=0$ altrimenti (la concordanza da doppio zero viene trattata in accordo con il significato dello zero)
- per descrittori semi-quantitativi ordinali e quantitativi:
 $s_i=1-|x_{ij}-x_{ik}| R_i^{-1}$ (dove R_i è l'intervallo di variazione dell' i -mo descrittore)

Steinhaus

$$S_{jk} = \frac{2 \sum_{i=1}^p \min(x_{ij}, x_{ik})}{\sum_{i=1}^p x_{ij} + x_{ik}}$$



Bray-Curtis

$$D_{jk} = \frac{\sum_{i=1}^p |x_{ij} - x_{ik}|}{\sum_{i=1}^p x_{ij} + x_{ik}}$$

Dissimilarità metriche se...

simmetrica

1. $D_{jk}=0$ se $j=k$

2. $D_{jk}>0$ se $j \neq k$

3. $D_{jk}=D_{kj}$

4. $D_{jk}+D_{kh} \geq D_{jh}$ (assioma della disuguaglianza triangolare)

Misure di distanza

euclidea

$$D_{jk} = \sqrt{\sum_{i=1}^p (x_{ij} - x_{ik})^2}$$

Manhattan

$$D_{jk} = \sum_{i=1}^p |x_{ij} - x_{ik}|$$

Minkowski

$$D_{jk} = r \sqrt{\left(\sum_{i=1}^p |x_{ij} - x_{ik}|^r \right)}$$

Czekanowski

$$D_{jk} = \frac{1}{p} \sum_{i=1}^p |x_{ij} - x_{ik}|$$

Canberra

$$D_{ij} = \sum_{i=1}^p \frac{|x_{ij} - x_{ik}|}{(x_{ij} + x_{ik})}$$

Bray-Curtis

$$D_{ij} = \frac{\sum_{i=1}^s |x_{ij} - x_{ik}|}{\sum_{i=1}^s (x_{ij} + x_{ik})}$$

corda

$$D_{jk} = \sqrt{2 \left(1 - \frac{\sum_{i=1}^p x_{ij} x_{ik}}{\sqrt{\sum_{i=1}^p x_{ij}^2 \sum_{i=1}^p x_{ik}^2}} \right)}$$

Misure di associazione

Fager & McGowan

$$S_{jk} = \frac{a}{\sqrt{(a+b)(a+c)}} - \frac{1}{2 \cdot \sqrt{a+c}} \quad (c \geq b)$$

...ma possono essere utilizzati anche i coefficienti di correlazione.

