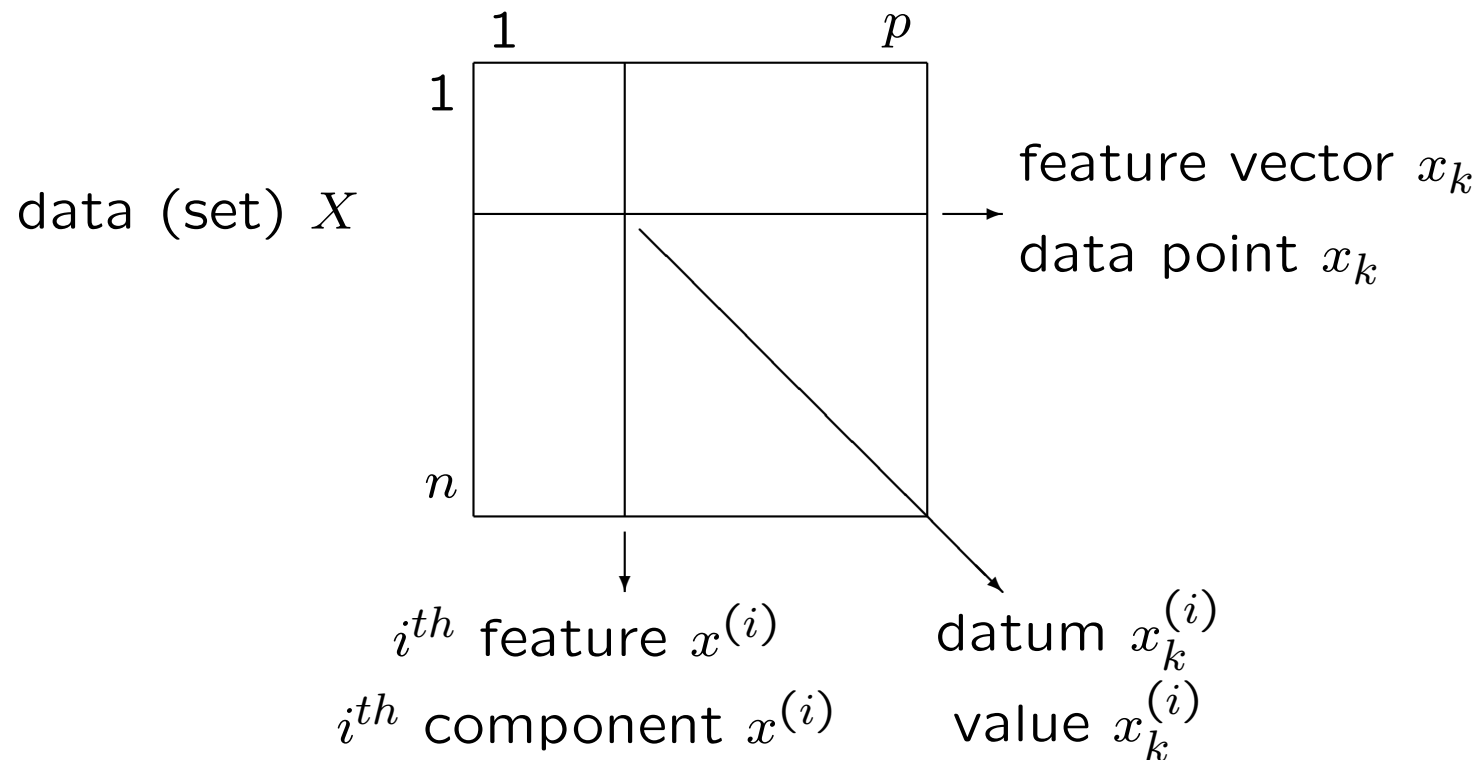


Scales

scale	operations		example	statistics
ratio	·	/	21 years, $273^{\circ}K$	generalized mean
interval	+	–	2015 A.D., $20^{\circ}C$	mean
ordinal	>	<	A, B, C, D, F	median
nominal	=	≠	Alice, Bob, Carol	mode

Matrix Representation of Numerical Data

- numerical data set $X = \{x_1, \dots, x_n\} \subset \mathbb{R}^p$



Relations on Object Data

- object data $O = \{o_1, \dots, o_n\}$
- relational data $R = \begin{pmatrix} r_{11} & \cdots & r_{1n} \\ \vdots & \ddots & \vdots \\ r_{n1} & \cdots & r_{nn} \end{pmatrix} \in \mathbb{R}^{n \times n}$
- semantics r_{ij} : similarity / dissimilarity
- symmetric relation: $r_{ij} = r_{ji} \ \forall i, j \in \{1, \dots, n\}$
- for $O \subset \mathbb{R}^p$: R can be obtained by norm on \mathbb{R}^p

Dissimilarity/Distance Measures

- Properties:

$$d(x, y) = d(y, x)$$

$$d(x, y) = 0 \Leftrightarrow x = y$$

$$d(x, z) \leq d(x, y) + d(y, z)$$

Norms

- $\|\cdot\| : \mathbb{R}^p \rightarrow \mathbb{R}^+$ is a norm, iff

1. $\|x\| = 0 \Leftrightarrow x = (0, \dots, 0)$

2. $\|a \cdot x\| = |a| \cdot \|x\| \quad \forall a \in \mathbb{R}, x \in \mathbb{R}^p$

3. $\|x + y\| \leq \|x\| + \|y\| \quad \forall x, y \in \mathbb{R}^p$

- “hyperbolic norm”

$$\|x\|_h = \prod_{i=1}^p x^{(i)}$$

ist not a norm!

Matrix Norms

$$\|x\|_A = \sqrt{x A x^T}$$

Examples:

- Euclidean:

$$A = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix}$$

- diagonal:

$$A = \begin{pmatrix} d_1 & 0 & \cdots & 0 \\ 0 & d_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & d_p \end{pmatrix}$$

- Mahalanobis:

$$A = \text{cov}^{-1}(X) = \left(\frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^T (x_k - \bar{x}) \right)^{-1}$$

Minkowski/Lebesgue Norms

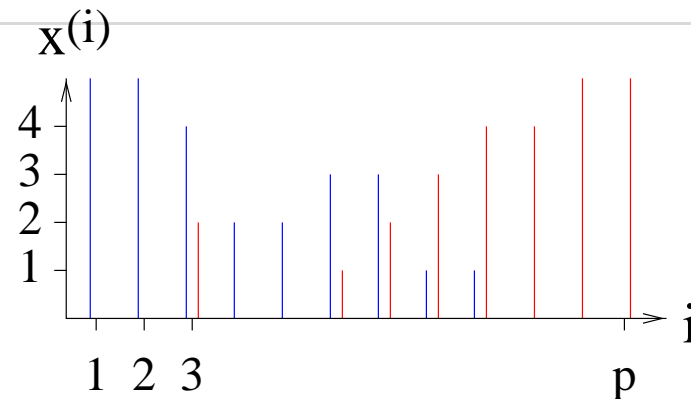
$$\|x\|_\alpha = \sqrt[\alpha]{\sum_{j=1}^p |x^{(j)}|^\alpha}$$

Examples:

- Manhattan or city block: $\alpha = 1$
- Euclidean: $\alpha = 2$
- sup or max: $\alpha \rightarrow \infty$

$$\lim_{\alpha \rightarrow \infty} \sqrt[\alpha]{\sum_{j=1}^p |x^{(j)}|^\alpha} = \max_{j=1, \dots, p} |x^{(j)}|$$

Similarity/Proximity Measures



- Properties:

$$s(x, y) = s(y, x)$$

$$s(x, y) \leq s(x, x)$$

$$s(x, y) \geq 0$$

- normalized similarity measures:

$$s(x, x) = 1$$

Similarity Measures $(x^{(i)}, y^{(i)} \geq 0)$

- cosine

$$s(x, y) = \frac{\sum_{i=1}^p x^{(i)} y^{(i)}}{\sqrt{\sum_{i=1}^p (x^{(i)})^2 \sum_{i=1}^p (y^{(i)})^2}}$$

- overlap

$$s(x, y) = \frac{\sum_{i=1}^p x^{(i)} y^{(i)}}{\min \left(\sum_{i=1}^p (x^{(i)})^2, \sum_{i=1}^p (y^{(i)})^2 \right)}$$

Similarity Measures $(x^{(i)}, y^{(i)} \geq 0)$

- Dice

$$s(x, y) = \frac{2 \sum_{i=1}^p x^{(i)} y^{(i)}}{\sum_{i=1}^p (x^{(i)})^2 + \sum_{i=1}^p (y^{(i)})^2}$$

- Jaccard/Tanimoto

$$s(x, y) = \frac{\sum_{i=1}^p x^{(i)} y^{(i)}}{\sum_{i=1}^p (x^{(i)})^2 + \sum_{i=1}^p (y^{(i)})^2 - \sum_{i=1}^p x^{(i)} y^{(i)}}$$