

Data Mining in Vector Space or Metric One?

Jaromír Kukal

DSEE FNSPE CTU Prague

May 13th 2010

Basic ideas

- Pattern set in vector space
- Pattern set in metric space
- Conversions
- Sammon's Mapping for metric spaces
- Cluster Analysis in both cases
- SOM in both cases
- PCA in both cases

Pattern in vector space

$$m, n \in \mathbf{N}$$

$$\mathbf{x}_i = (x_{i,1}, \dots, x_{i,n}) \in \mathbf{R}^n$$

$$\mathbf{S} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subset \mathbf{R}^n$$

$$\mathbf{X} = \begin{pmatrix} x_{1,1} & \cdots & x_{1,n} \\ \vdots & \ddots & \vdots \\ x_{m,1} & \cdots & x_{m,n} \end{pmatrix}$$

Normalization in vector space

$$u_j = \max_i x_{i,j}$$

$$l_j = \min_i x_{i,j}$$

$$x_{i,j}^* = \begin{cases} 2 \frac{x_{i,j} - l_j}{u_j - l_j} - 1 & \text{for } u_j > l_j \\ 0 & \text{for } u_j = l_j \\ 0 & \text{for empty } x_{i,j} \end{cases}$$

Standardization in vector space

$$\mu_j = \text{mean}_i x_{i,j}$$

$$\sigma_j = \text{std}_i x_{i,j}$$

$$x_{i,j}^* = \begin{cases} \frac{x_{i,j} - \mu_j}{\sigma_j} & \text{for } \sigma_j > 0 \\ 0 & \text{for } \sigma_j = 0 \\ 0 & \text{for empty } x_{i,j} \end{cases}$$

Robust transform in vector space

$$\eta_j = \operatorname{median}_i x_{i,j}$$

$$\lambda_j = \operatorname{median}_i |x_{i,j} - \eta_j|$$

$$x_{i,j}^* = \begin{cases} \frac{x_{i,j} - \eta_j}{\lambda_j} & \text{for } \lambda_j > 0 \\ 0 & \text{for } \lambda_j = 0 \\ 0 & \text{for empty } x_{i,j} \end{cases}$$

Minkowski distance in \mathbb{R}^n

$$d_p(\mathbf{x}, \mathbf{y}) = \left(\sum_k |x_k - y_k|^p \right)^{1/p}, \quad p \geq 1$$

$$\lim_{p \rightarrow \infty} d_p(\mathbf{x}, \mathbf{y}) = \max_k |x_k - y_k|$$

- Manhattan distance $p = 1$
- Compromise distance $p = 3/2$
- Euclidean distance $p = 2$
- Maximum distance $p \rightarrow \infty$

Distance representation

$$d_{i,j} = d_p(\mathbf{x}_i, \mathbf{x}_j)$$

$$\mathbf{D} = \begin{pmatrix} d_{1,1} & \cdots & d_{1,m} \\ \vdots & \ddots & \vdots \\ d_{m,1} & \cdots & d_{m,m} \end{pmatrix}$$

special case $mn > m(m-1)/2$

Metric space

$$\text{card}(\mathbf{U}) > 0$$

$$\langle \mathbf{U}, d \rangle$$

$$d: \mathbf{U} \times \mathbf{U} \rightarrow \mathbf{R}_0^+$$

$$\forall x, y \in \mathbf{U} \quad d(x, y) = d(y, x)$$

$$\forall x, y \in \mathbf{U} \quad d(x, y) = 0 \Leftrightarrow x = y$$

$$\forall x, y, z \in \mathbf{U} \quad d(x, y) \leq d(x, z) + d(z, y)$$

Pattern in metric space

$$m \in \mathbf{N}$$

$$x_i \in \mathbf{U}$$

$$\mathbf{S} = \{x_1, \dots, x_m\} \subseteq \mathbf{U}$$

$$d_{i,j} = d(x_i, x_j)$$

$$\mathbf{D} = \begin{pmatrix} d_{1,1} & \cdots & d_{1,m} \\ \vdots & \ddots & \vdots \\ d_{m,1} & \cdots & d_{m,m} \end{pmatrix}$$

Normalization in metric space

$$M = \begin{cases} \max_{i < j} d_{i,j} \\ \text{mean}_{i < j} d_{i,j} \\ \text{median}_{i < j} d_{i,j} \end{cases}$$

$$d_{i,j}^* = d_{i,j} / M$$

$$\mathbf{D}^* = \mathbf{D} / M$$

Conversion from U to \mathbf{R}^n

$$n \in \mathbf{N}$$

$$e_k \in U$$

$$\mathbf{E} = \{e_1, \dots, e_n\} \subseteq U$$

$$\mathbf{x}_i = (d(x_i, e_1), \dots, d(x_i, e_n))$$

$$\mathbf{x}_i \in \mathbf{R}^n$$

$$\mathbf{E} = \mathbf{S} \Rightarrow \mathbf{X} = \mathbf{D}$$

PCA preliminaries

PCA : $\mathbf{R}^n \rightarrow \mathbf{R}^H$, where $H < n$

$$\mathbf{y} = \text{PCA}(\mathbf{x}) = \mathbf{x}\mathbf{W} \quad \mathbf{W} \in \mathbf{R}^{n \times H}$$

$$\mathbf{X} \in \mathbf{R}^{m \times n} \quad \sum_i x_{i,j} = 0$$

$$H \leq \text{rank}(\mathbf{X}) \leq \min(m-1, n)$$

$$\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^T$$

$$\mathbf{U}^T \mathbf{U} = \mathbf{I}_m \quad \mathbf{V}^T \mathbf{V} = \mathbf{I}_n$$

$$\mathbf{S} = \text{diag}(\mathbf{s}) \geq \mathbf{0}$$

Traditional PCA

$$\mathbf{A} = \mathbf{X}^T \mathbf{X} \geq \mathbf{0}$$

$$(\mathbf{A} - \lambda \mathbf{I}_n) \mathbf{v} = \mathbf{0} \quad \|\mathbf{v}\| = 1$$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$$

$$\mathbf{W} = \left(\begin{array}{ccc} \frac{\mathbf{v}_1}{\sqrt{\lambda_1}} & \dots & \frac{\mathbf{v}_H}{\sqrt{\lambda_H}} \end{array} \right)$$

$$\mathbf{Y} = \mathbf{XW} = (\mathbf{u}_1 \quad \dots \quad \mathbf{u}_H)$$

PCA for large n

$$\mathbf{B} = \mathbf{X}\mathbf{X}^T \geq \mathbf{0}$$

$$(\mathbf{B} - \lambda \mathbf{I}_m) \mathbf{u} = \mathbf{0} \quad \|\mathbf{u}\| = 1$$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m \geq 0$$

$$\mathbf{W} = \mathbf{X}^T \begin{pmatrix} \frac{\mathbf{u}_1}{\lambda_1} & \dots & \frac{\mathbf{u}_H}{\lambda_H} \end{pmatrix}$$

$$\mathbf{Y} = \mathbf{X}\mathbf{W} = (\mathbf{u}_1 \quad \dots \quad \mathbf{u}_H)$$

Sammon's mapping

$$f : \mathbf{U} \rightarrow \mathbf{V} \quad y = f(x)$$

$$\langle \mathbf{U}, d \rangle \quad d_{i,j} = d(x_i, x_j)$$

$$\langle \mathbf{V}, \delta \rangle \quad \delta_{i,j} = \delta(y_i, y_j)$$

$$F = \frac{\sum_{i < j} \frac{(d_{i,j} - \delta_{i,j})^2}{d_{i,j}}}{\sum_{i < j} d_{i,j}} = \min_f$$

$$\text{postulate : } d_{i,j} = \text{undefined} \Rightarrow \frac{(d_{i,j} - \delta_{i,j})^2}{d_{i,j}} = 0$$

Nonlinear alternative to PCA

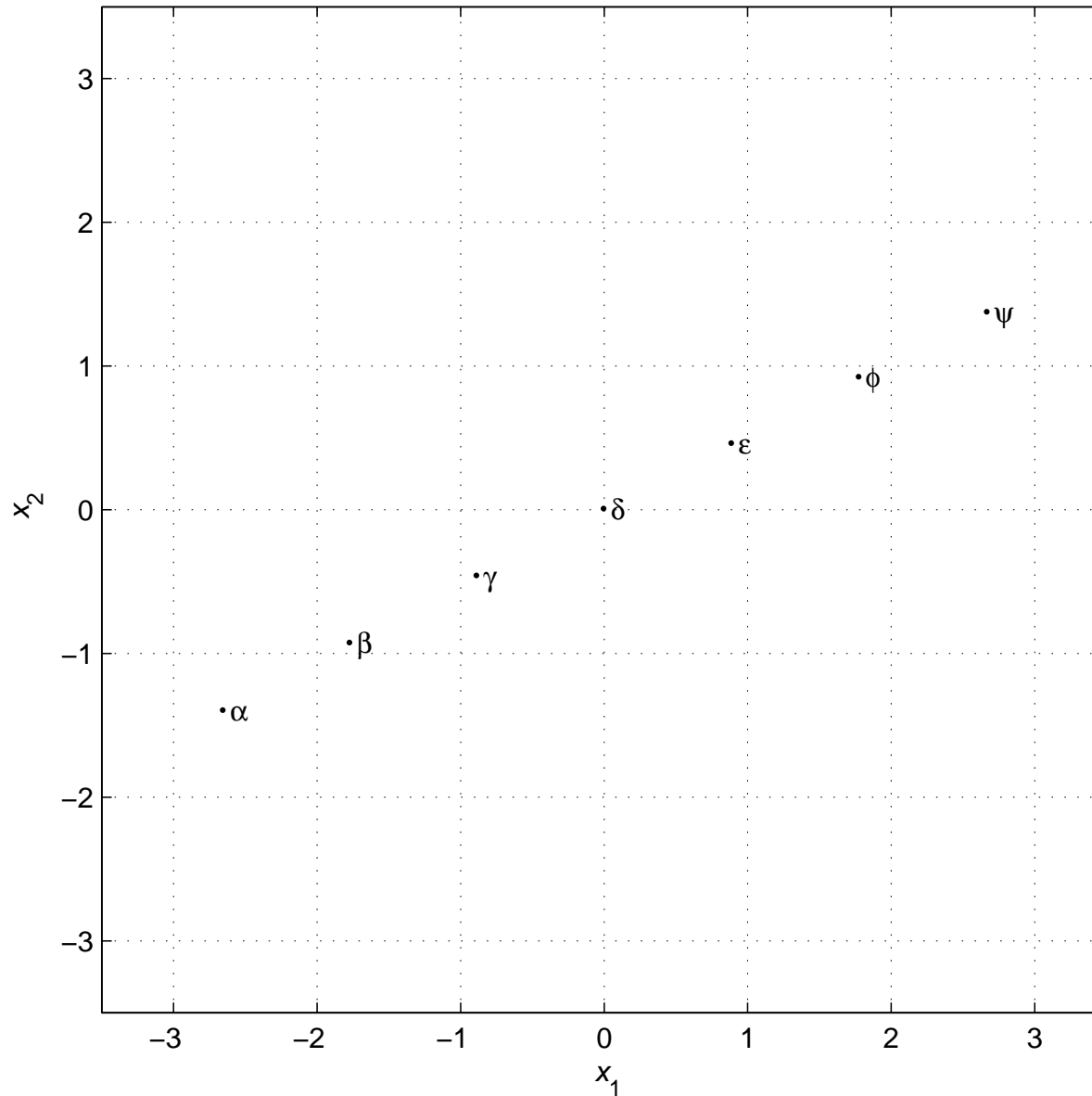
$$H \in \mathbf{N}$$

$$\mathbf{V} = \mathbf{R}^H$$

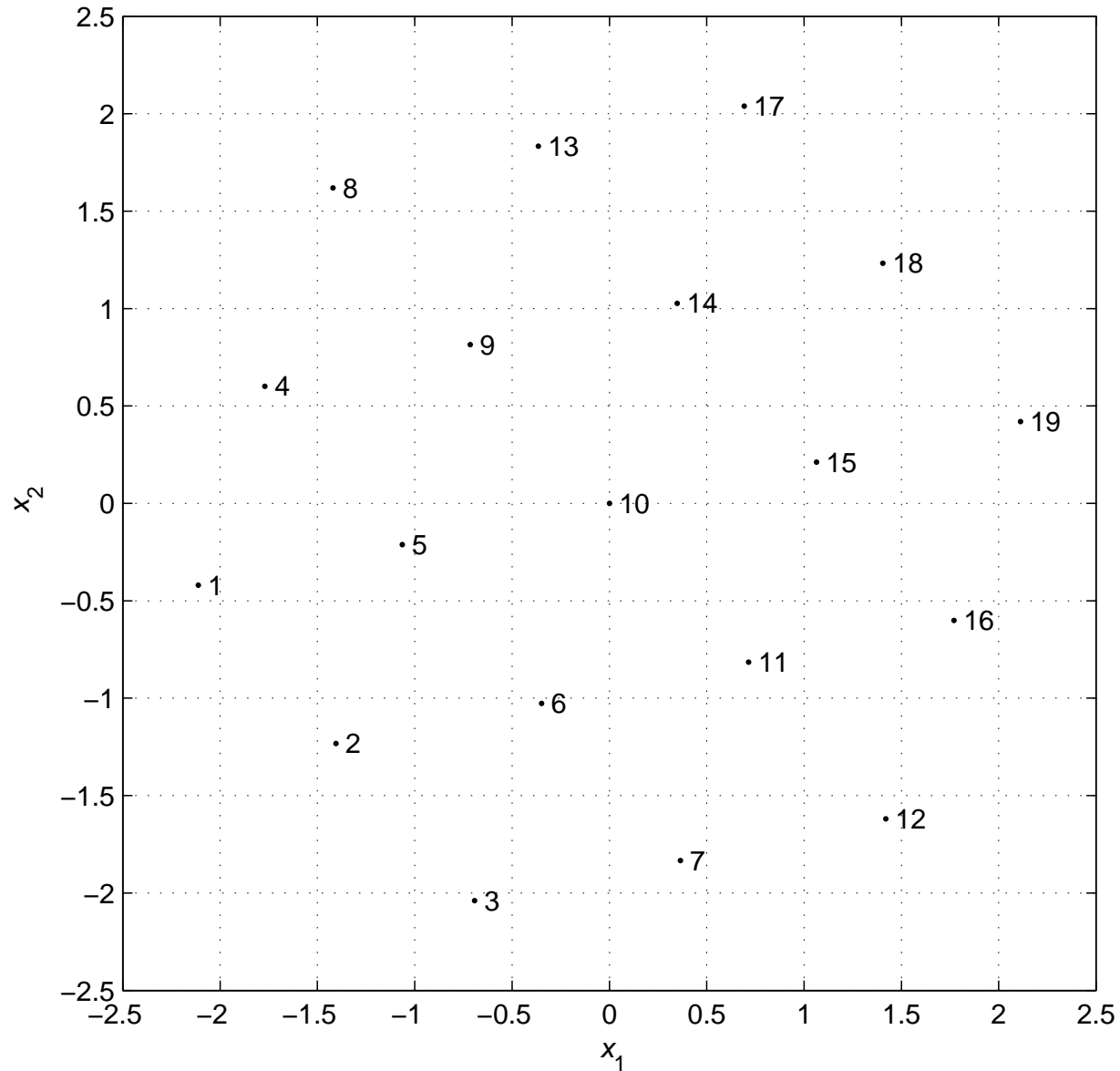
$$\delta(\mathbf{a}, \mathbf{b}) = \left(\sum_k (a_k - b_k)^2 \right)^{1/2}$$

$$\Phi_1 = \sum_{i < j} \frac{(d_{i,j} - \delta(\mathbf{y}_i, \mathbf{y}_j))^2}{d_{i,j}} = \min_{\mathbf{y}_1, \dots, \mathbf{y}_m}$$

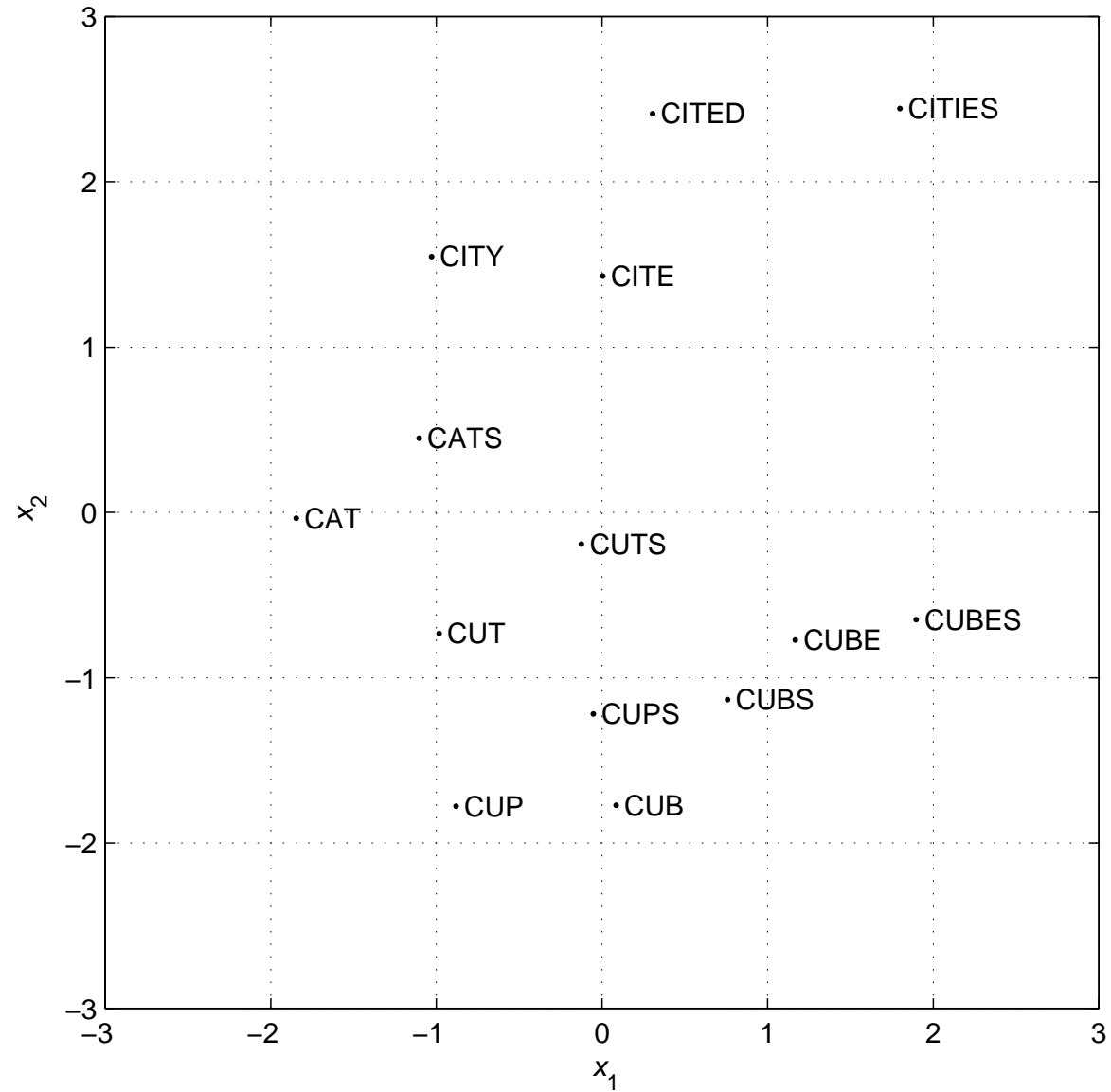
Are they sitting on the bench?



Plot of hexagonal graph



Levensthein distance in \mathbb{U}



Modified PCA

$$\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_H) \in (\mathbf{R}^+)^H$$

$$\mathbf{y} = (PCA_1, \dots, PCA_H) \in \mathbf{R}^H$$

$$\mathbf{y}_{\text{modi}} = \boldsymbol{\varepsilon} \otimes \mathbf{y} = (\varepsilon_1 PCA_1, \dots, \varepsilon_H PCA_H) \in \mathbf{R}^H$$

$$\delta(\mathbf{a}, \mathbf{b}) = \left(\sum_k (a_k - b_k)^2 \right)^{1/2}$$

$$\Phi_2 = \sum_{i < j} \frac{(d_{i,j} - \delta(\boldsymbol{\varepsilon} \otimes \mathbf{y}_i, \boldsymbol{\varepsilon} \otimes \mathbf{y}_j))^2}{d_{i,j}} = \min_{\boldsymbol{\varepsilon}}$$

Alternative to cluster analysis

$$H > 1 \quad \mathbf{V} = \{1, \dots, H\}$$

$$\delta(a, b) = \begin{cases} 0 & \text{for } a = b \\ \text{undefined} & \text{for } a \neq b \end{cases}$$

$$\mathbf{p} = (p_1, \dots, p_m) \in \mathbf{V}^m$$

$$\text{postulate : } \delta_{i,j} = \text{undefined} \Rightarrow \frac{(d_{i,j} - \delta_{i,j})^2}{d_{i,j}} = 0$$

$$\Phi_3 = \sum_{\substack{i < j \\ \delta_{i,j}=0}} \frac{(d_{i,j} - \delta_{i,j})^2}{d_{i,j}} = \sum_{\substack{i < j \\ \delta_{i,j}=0}} d_{i,j} = \sum_{\substack{i < j \\ p_i=p_j}} d_{i,j} = \min_{\mathbf{p}} \sum_{\substack{i < j \\ p_i=p_j}} d_{i,j}$$

Alternative to SOM

$$G = \langle \mathbf{V}, \mathbf{E} \rangle \quad \mathbf{V} = \{1, \dots, H\} = V(G)$$

$$\mathbf{E} = E(G) \subset \binom{\mathbf{V}}{2} \quad \delta(a, b) \text{ is vertex distance in } G$$

$$\mathbf{p} = (p_1, \dots, p_m) \in \mathbf{V}^m \quad a > 0$$

$$\Phi_4 = \sum_{i < j} \frac{(d_{i,j} - a \delta(p_i, p_j))^2}{d_{i,j}} = \min_{\mathbf{p}, a}$$

Cluster and partition

$$H > 1 \quad \mathbf{V} = \{1, \dots, H\}$$

$$\mathbf{p} = (p_1, \dots, p_m) \in \mathbf{V}^m$$

$$\mathbf{C}_k = \{x_i \in \mathbf{S} \mid p_i = k\}$$

$$\mathbf{S} = \bigcup_{k \in \mathbf{V}} \mathbf{C}_k$$

Momentum in metric space

$$F_k(y) = \sum_{x_i \in \mathbf{C}_k} d^2(x_i, y)$$

$$\mathbf{S} \subseteq \mathbf{Q} \subseteq \mathbf{U}$$

$$t_k \in \arg \min_{y \in \mathbf{U}} F_k(y)$$

$$t_k^+ \in \arg \min_{y \in \mathbf{Q}} F_k(y)$$

$$F_k^* = F_k(t_k)$$

$$F_k^* \leq F_k^+ = F_k(t_k^+)$$

$\mathbf{U} = \mathbf{R}^n$ with Euclidean metrics $\Rightarrow \mathbf{t}_k = \text{mean}_{\mathbf{x}_i \in \mathbf{C}_k} \mathbf{x}_i$

Center in metric space

$$G_k(y) = \sum_{x_i \in \mathbf{C}_k} d(x_i, y)$$

$$\mathbf{S} \subseteq \mathbf{Q} \subseteq \mathbf{U}$$

$$s_k \in \arg \min_{y \in \mathbf{U}} G_k(y)$$

$$s_k^+ \in \arg \min_{y \in \mathbf{Q}} G_k(y)$$

$$G_k^* = G_k(s_k)$$

$$G_k^* \leq G_k^+ = G_k(s_k^+)$$

$\mathbf{U} = \mathbf{R}^n$ with Manhattan metrics $\Rightarrow \mathbf{s}_k = \operatorname{median}_{\mathbf{x}_i \in \mathbf{C}_k} \mathbf{x}_i$

ISODATA in metric space

Random initialization $\mathbf{p} = (p_1, \dots, p_m) \in \mathbf{V}^m$ $\text{card}(\mathbf{C}_k) > 0$

Momentum revision $t_k^+ \in \arg \min_{y \in \mathbf{Q}} F_k(y)$

Cluster revision $x_i \in \mathbf{C}_k \Rightarrow k \in \arg \min_{1 \leq j \leq H} d(x_i, t_j^+)$

Objective function $\Phi_5 = \sum_k F_k^+ = \min_{\mathbf{p}}$

ISODATA in vector space

Random initialization $\mathbf{p} = (p_1, \dots, p_m) \in \mathbf{V}^m$ $\text{card}(\mathbf{C}_k) > 0$

Momentum revision $\mathbf{t}_k = \text{mean}_{\mathbf{x}_i \in \mathbf{C}_k} \mathbf{x}_i$

Cluster revision $\mathbf{x}_i \in \mathbf{C}_k \Rightarrow k \in \arg \min_{1 \leq j \leq H} d(\mathbf{x}_i, \mathbf{t}_j)$

Objective function $\Phi_6 = \sum_k F_k^* = \min_{\mathbf{p}}$

SOM preliminaries

$$G = \langle \mathbf{V}, \mathbf{E} \rangle$$

$$\mathbf{V} = \{1, \dots, H\} = V(G)$$

$$\mathbf{E} = E(G) \subset \binom{\mathbf{V}}{2}$$

$\delta(a, b)$ is vertex distance in G

$R > 0$ is learning radius

$$\mathbf{p} = (p_1, \dots, p_m) \in \mathbf{V}^m \quad \mathbf{w} = (w_1, \dots, w_H) \in \mathbf{U}^H$$

$$\chi: \mathbf{N}_0 \rightarrow [0;1]$$

$$\forall k \in \mathbf{N}_0 : \chi(k) \geq \chi(k+1)$$

$$\chi(0) = 1 \quad \chi(M) < 1 \quad M = \max_{a, b \in \mathbf{V}} \delta(a, b)$$

SOM characteristic

$$\chi(\delta) = \begin{cases} 1 & \text{for } \delta \leq R \\ 0 & \text{for } \delta > R \end{cases}$$

$$\chi(\delta) = \exp\left(-\frac{\delta^2}{2R^2}\right)$$

$$\chi(\delta) = \left(1 + \frac{\delta^2}{R^2}\right)^{-1}$$

Batch SOM in metric space

Random initialization $\mathbf{p} = (p_1, \dots, p_m) \in \mathbf{V}^m$ $\text{card}(\mathbf{C}_k) > 0$

Weight revision $w_i \in \arg \min_{y \in \mathbf{Q}} \sum_{k=1}^H \sum_{x_j \in \mathbf{C}_k} d^2(x_j, y) \chi(\delta(i, k))$

Vertex revision $x_i \in \mathbf{C}_k \Rightarrow k \in \arg \min_{1 \leq k \leq H} \sum_{i=1}^H d^2(x_j, w_i) \chi(\delta(i, k))$

Objective function $\Phi_\gamma = \sum_{k=1}^H \sum_{i=1}^H \sum_{x_j \in \mathbf{C}_k} d^2(x_j, w_i) \chi(\delta(i, k)) = \min_{\mathbf{p}, \mathbf{w}}$

Batch SOM in vector space

Random initialization $\mathbf{p} = (p_1, \dots, p_m) \in \mathbf{V}^m$ $\text{card}(\mathbf{C}_k) > 0$

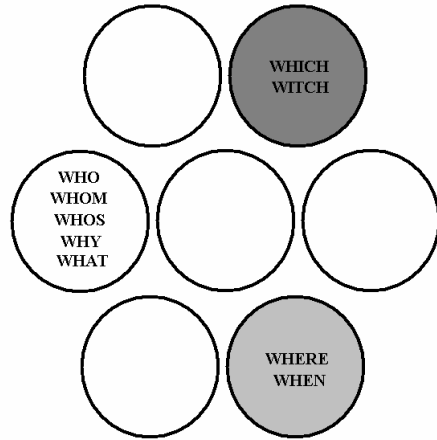
Weight revision
$$\mathbf{w}_i = \frac{\sum_{k=1}^H \sum_{\mathbf{x}_j \in \mathbf{C}_k} \mathbf{x}_j \chi(\delta(i, k))}{\sum_{k=1}^H \sum_{\mathbf{x}_j \in \mathbf{C}_k} \chi(\delta(i, k))}$$

Vertex revision $\mathbf{x}_i \in \mathbf{C}_k \Rightarrow k \in \arg \min_{1 \leq k \leq H} \sum_{i=1}^H d^2(\mathbf{x}_j, \mathbf{w}_i) \chi(\delta(i, k))$

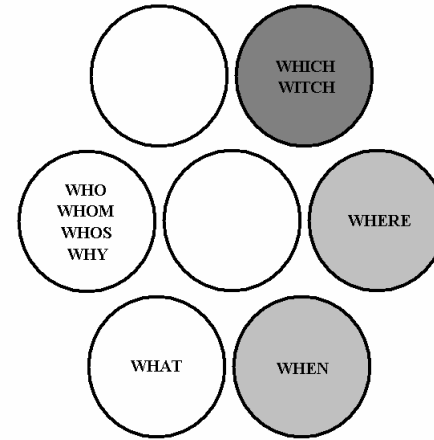
Objective function
$$\Phi_8 = \sum_{k=1}^H \sum_{i=1}^H \sum_{\mathbf{x}_j \in \mathbf{C}_k} d^2(\mathbf{x}_j, \mathbf{w}_i) \chi(\delta(i, k)) = \min_{\mathbf{p}, \mathbf{w}_1, \dots, \mathbf{w}_H}$$

SOM of written words

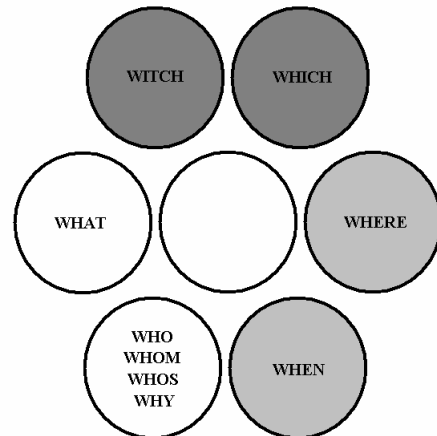
$R = 1.0$



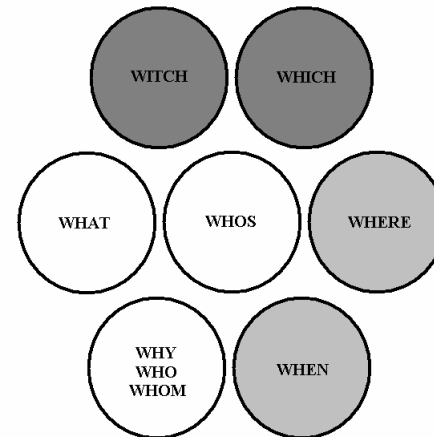
$R = 0.5$



$R = 0.2$

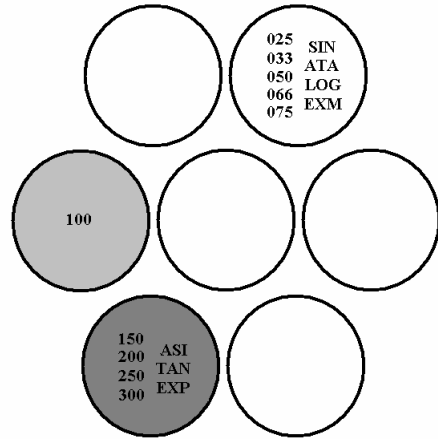


$R = 0.1$

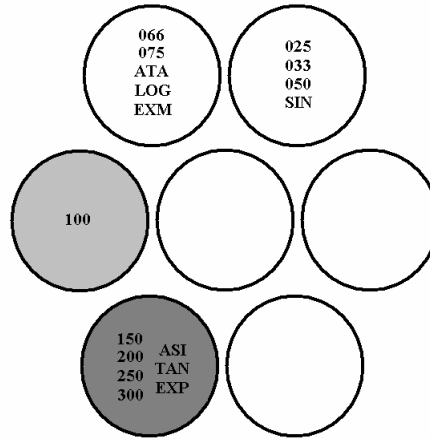


SOM of functions

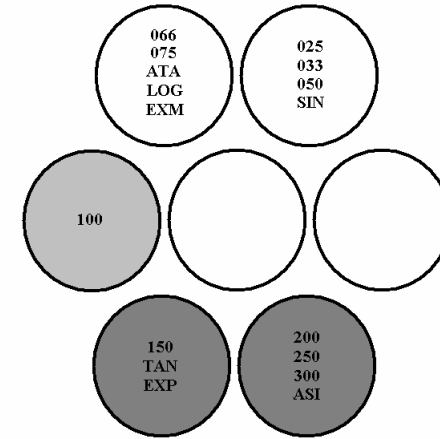
$R = 0.6$



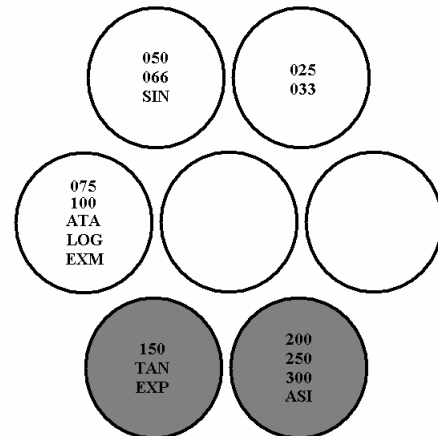
$R = 0.5$



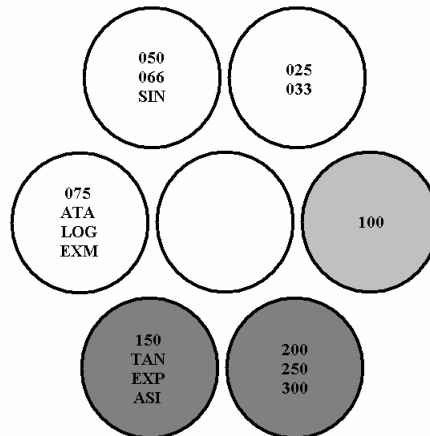
$R = 0.4$



$R = 0.3$



$R = 0.2$



$R = 0.1$

