

Machine Intelligence 2

4.3 Self-Organising Maps

Prof. Dr. Klaus Obermayer

Fachgebiet Neuronale Informationsverarbeitung (NI)

SS 2018

Self-Organising Maps

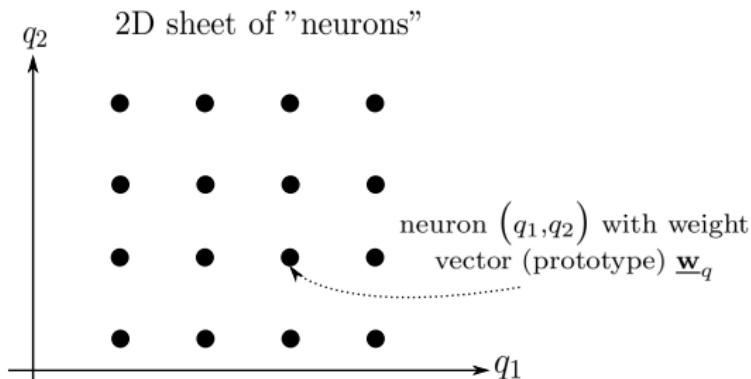
Self-organising maps (SOM)

- clustering & local embedding of vectorial data: $\underline{x}^{(\alpha)}, \alpha = 1, \dots, p$
- clustering of data based on similarity: squared Euclidean distance
- low-dimensional and **neighborhood preserving** representation for visualization

Model class

clustering with multidimensional cluster index \underline{q} , arranged in a grid:

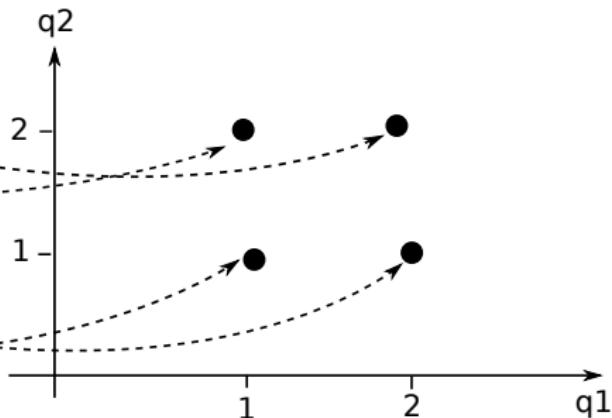
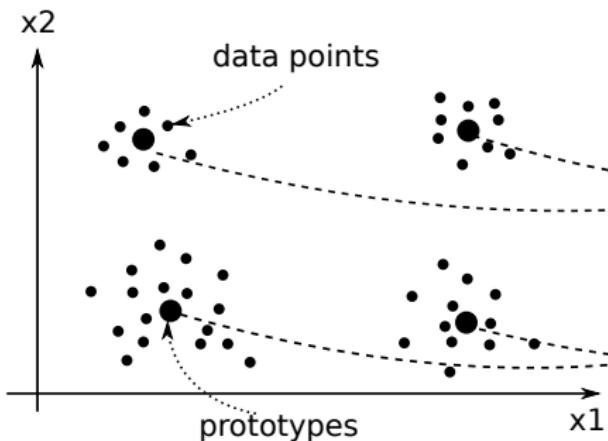
$$m_{\underline{q}}^{(\alpha)} = \begin{cases} 1, & \text{if } \underline{q} = \operatorname{argmin}_{\underline{r}} |\underline{x}^{(\alpha)} - \underline{w}_{\underline{r}}| \\ 0, & \text{else} \end{cases}$$



"Binary neurons" with position $\underline{q} = (q_1, q_2)^T$ are assigned to prototypes $\underline{w}_{\underline{q}}$.

"Topographic" maps

Following algorithm's aim by self-organization:



(low-dim.) "map" of (high-dim.) data space for $\left\{ \begin{array}{l} \text{visualization} \\ \text{dimension reduction} \\ \text{preprocessing for prediction} \end{array} \right.$

Algorithm 1: On-line learning for SOM

Initialization:

- choose no. M of partitions ("neurons")
- choose annealing schedule for ε and σ
- initialize M prototypes: $\underline{\mathbf{w}}_q = 1/p \sum_{\alpha} \underline{\mathbf{x}}^{(\alpha)} + \underline{\eta}$, $\underline{\eta}$ small noise vector

beginchoose a random data point $\underline{\mathbf{x}}^{(\alpha)}$

determine the closest prototype:

$$\underline{\mathbf{p}} = \operatorname{argmin}_{\underline{\mathbf{r}}} |\underline{\mathbf{x}}^{(\alpha)} - \underline{\mathbf{w}}_{\underline{\mathbf{r}}}|$$

change all prototypes $\underline{\mathbf{w}}_q$ according to:

$$\Delta \underline{\mathbf{w}}_q = \varepsilon \cdot h_{\underline{\mathbf{q}}\underline{\mathbf{p}}} \cdot (\underline{\mathbf{x}}^{(\alpha)} - \underline{\mathbf{w}}_q) \quad \text{for all } \underline{\mathbf{q}}$$

lower the learning rate ε and neighborhood width σ **end**

Neighborhood function $h_{\underline{q}\underline{p}}$

$$\Delta \underline{\mathbf{w}}_{\underline{\mathbf{q}}} = \varepsilon \cdot h_{\underline{\mathbf{q}}\underline{\mathbf{p}}} \cdot (\underline{\mathbf{x}}^{(\alpha)} - \underline{\mathbf{w}}_{\underline{\mathbf{q}}}^{\text{old}}) \quad \text{for all } \underline{\mathbf{q}}$$

- $h_{\underline{\mathbf{q}}\underline{\mathbf{p}}}$ enforces similar learning steps for neighboring neurons
- common choice:

$$h_{\underline{\mathbf{q}}\underline{\mathbf{p}}} = \exp \left\{ - \frac{(\underline{\mathbf{q}} - \underline{\mathbf{p}})^2}{2\sigma^2} \right\}$$

- $\sigma = 0$: standard on-line K-means (only update $\underline{\mathbf{q}} = \underline{\mathbf{p}}$)

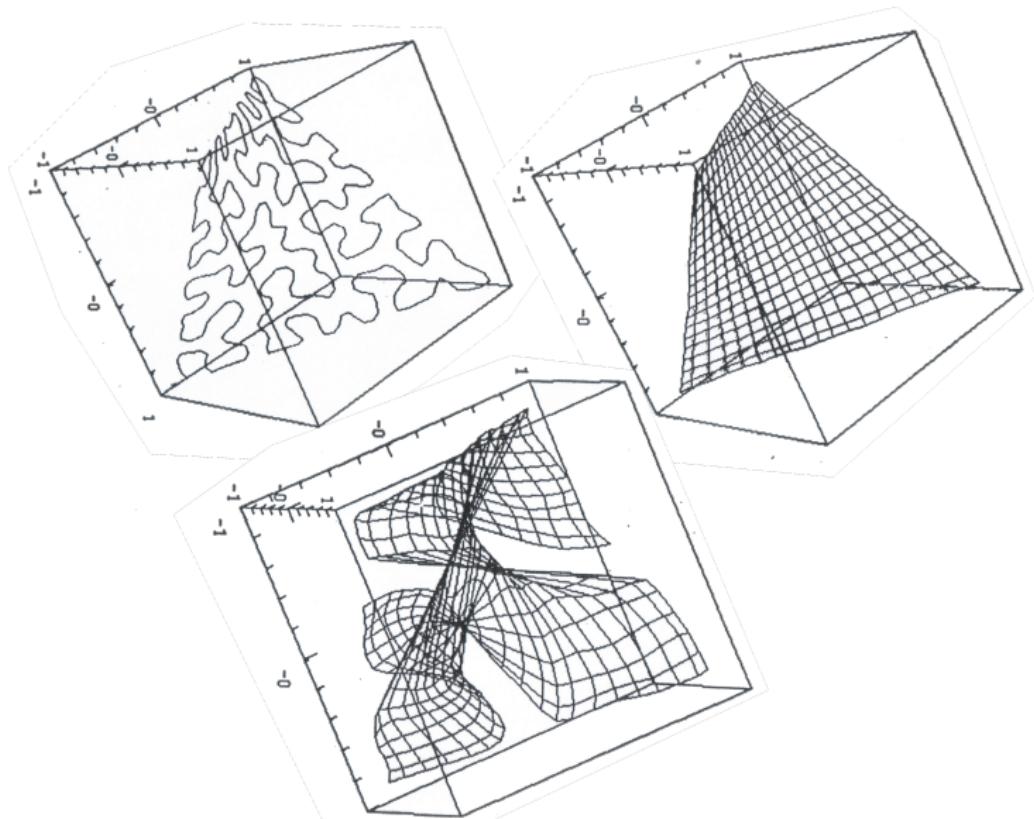
Neighborhood function $h_{\underline{\mathbf{q}}\underline{\mathbf{p}}}$

$$h_{\underline{\mathbf{q}}\underline{\mathbf{p}}} = \exp \left\{ -\frac{(\underline{\mathbf{q}} - \underline{\mathbf{p}})^2}{2\sigma^2} \right\}$$

σ annealing

- start with large σ (\leadsto neighborhood function convex over its support)
- decrease linearly or exponentially (but "slow") during learning.
- solution depends on final value of σ
- $\sigma = 0$: minimum of the K-means clustering cost function but "*neighborhood preserving*"
- σ small : better representation capabilities at the expense of a non-optimal clustering cost

2d manifold example



Measurements of Leptograpsus variegatus

Dead crabs loose their color and their sexual features –
 ⇒ Can we infer species & sex from the shell size and shape?

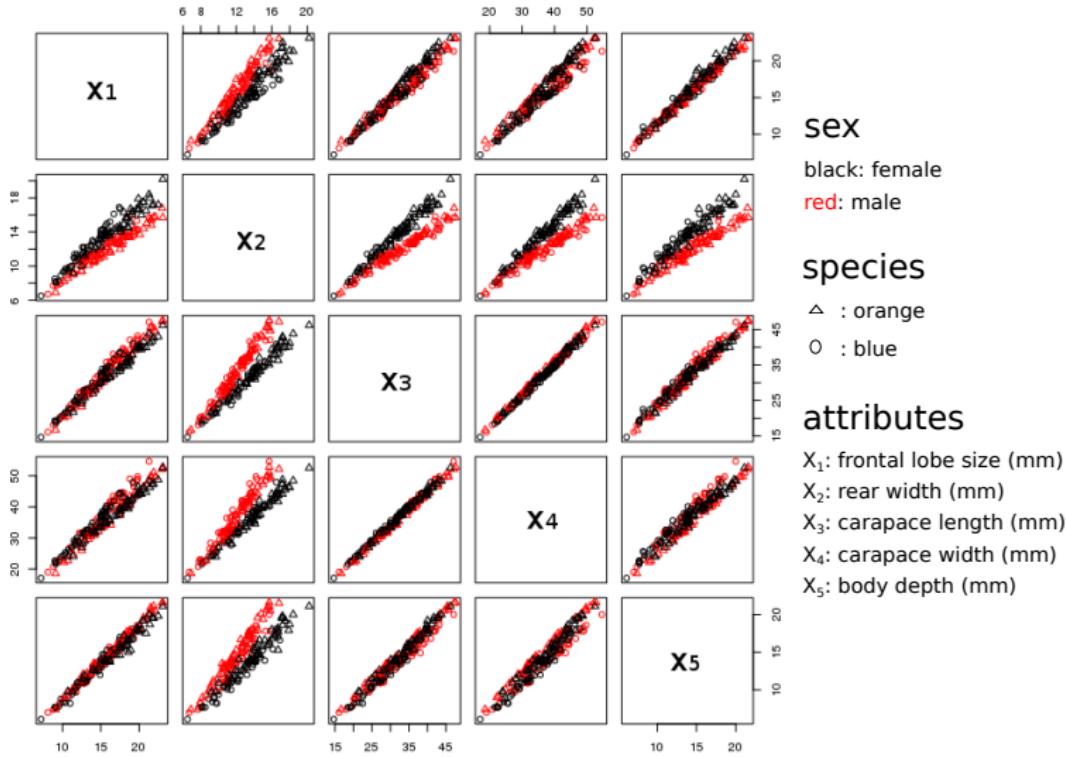
crabs: L. variegatus $\left\{ \begin{array}{l} \text{orange} \\ \text{blue} \end{array} \right\}$ two (sub-)species
 → male and female crabs

- *elementary features:*

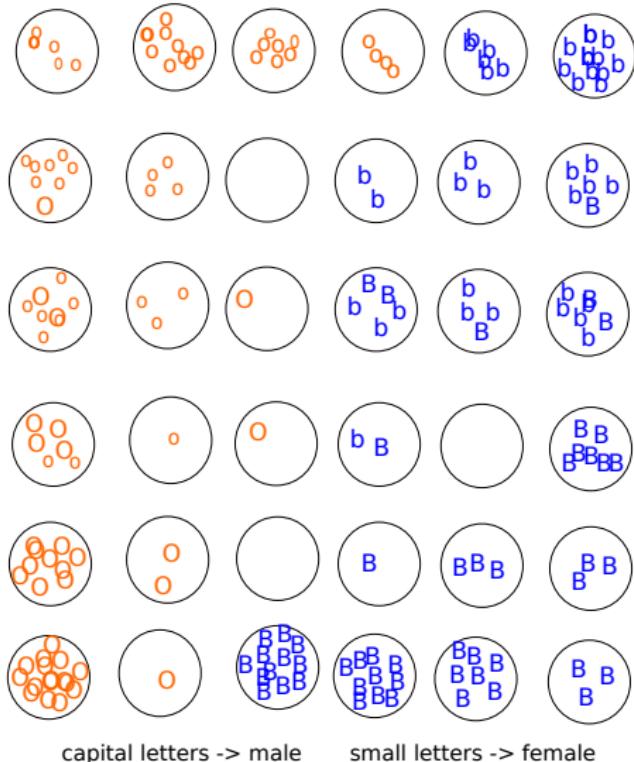
frontal lobe size:	x_1	5-dim. feature vector
rear width:	x_2	
carapace length:	x_3	
carapace width:	x_4	
body depth:	x_5	

- *complex features:* linear combinations of elementary features
 ⇒ directions in feature space which could identify color and/or sex.

The Leptograpsus data: scatter plot



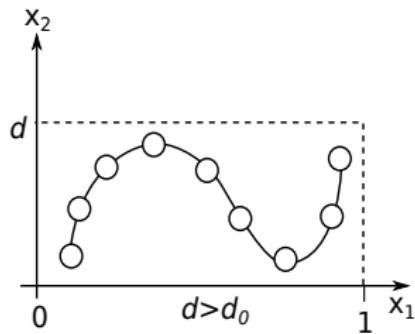
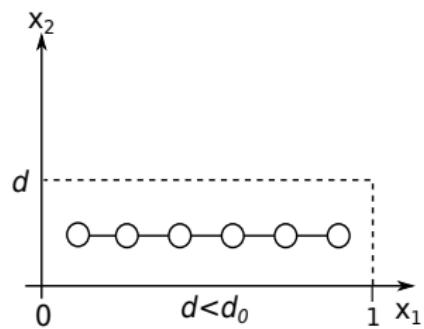
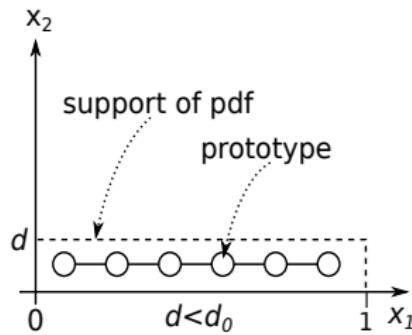
Leptograpsus example



attributes

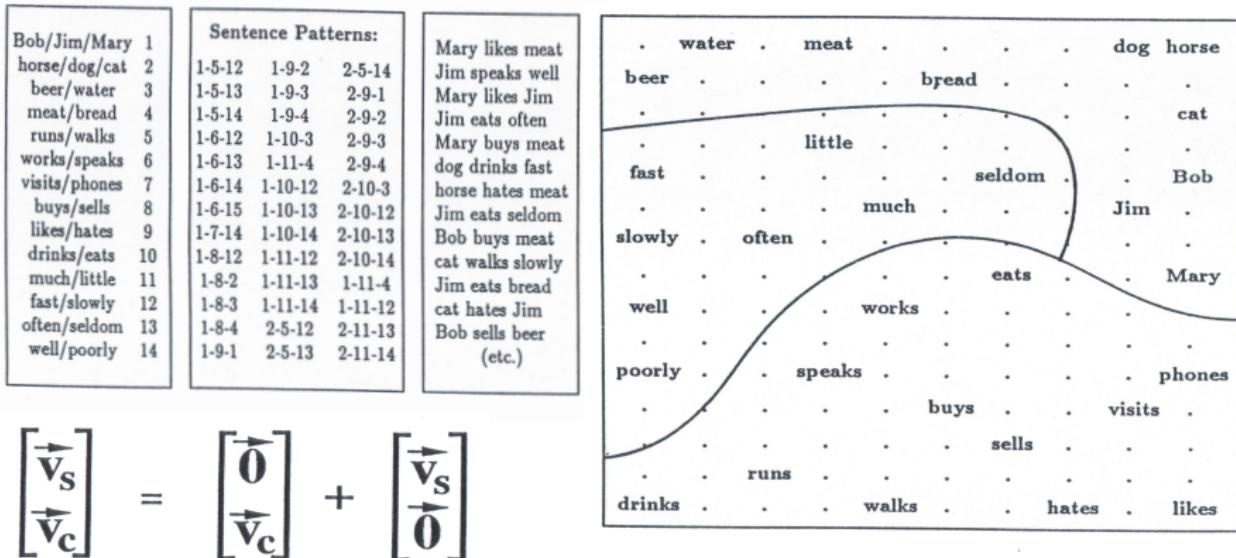
- X_1 : frontal lobe size (mm)
- X_2 : rear width (mm)
- X_3 : carapace length (mm)
- X_4 : carapace width (mm)
- X_5 : body depth (mm)

Dimension reducing mappings



Semantic maps

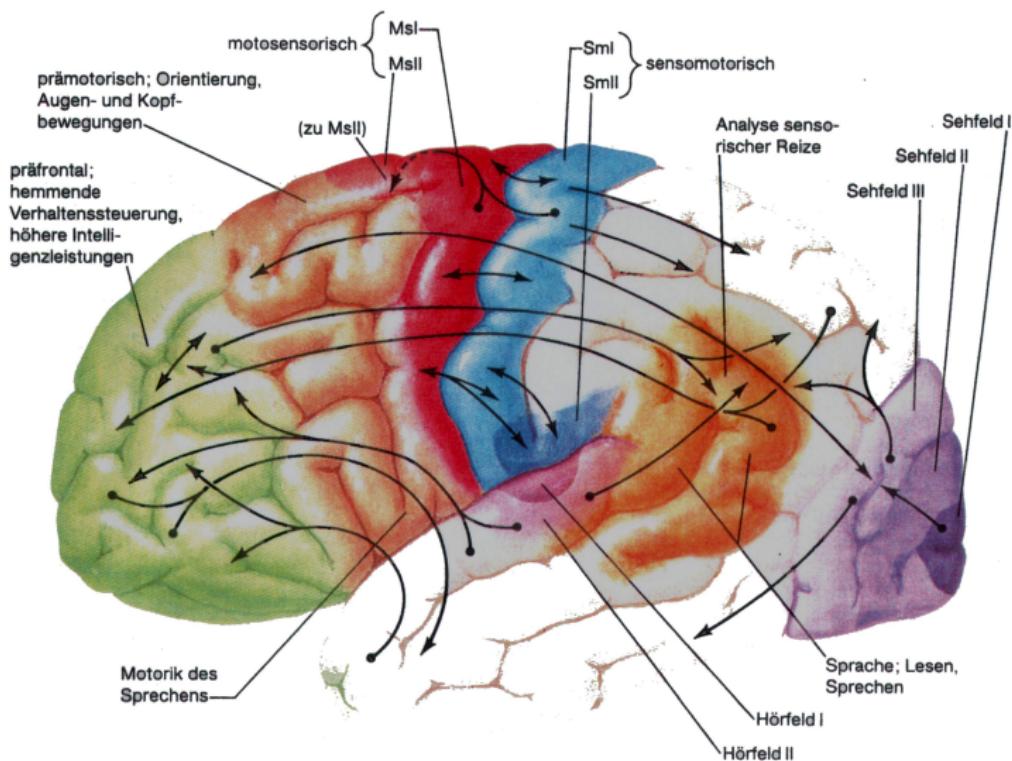
Abstract Data: Verbal Statements



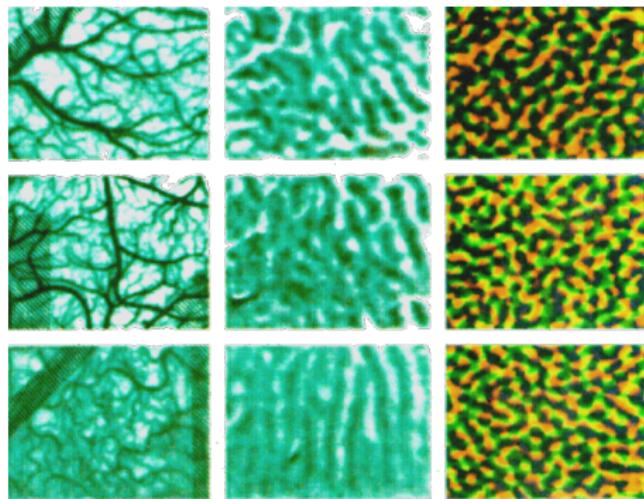
$$\begin{bmatrix} \vec{v}_s \\ \vec{v}_c \end{bmatrix} = \begin{bmatrix} \vec{0} \\ \vec{0} \end{bmatrix} + \begin{bmatrix} \vec{v}_s \\ \vec{v}_c \end{bmatrix}$$

Ritter & Kohnen 1989

Orientation & ocular dominance maps

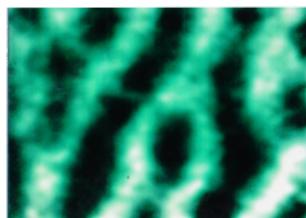
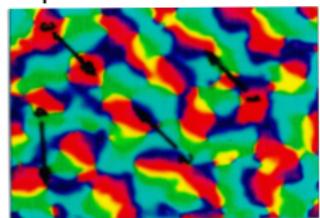


Orientation & ocular dominance maps

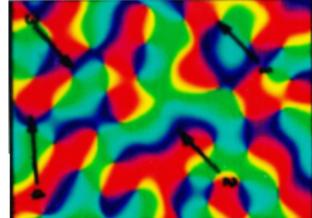


Orientation & ocular dominance maps

experimental data:



mapping results:



5d feature vector

position: $x, y \in [0, d]$

orientation preference & specificity:

$$r \cos(2\phi), \phi \in [0, \pi[$$

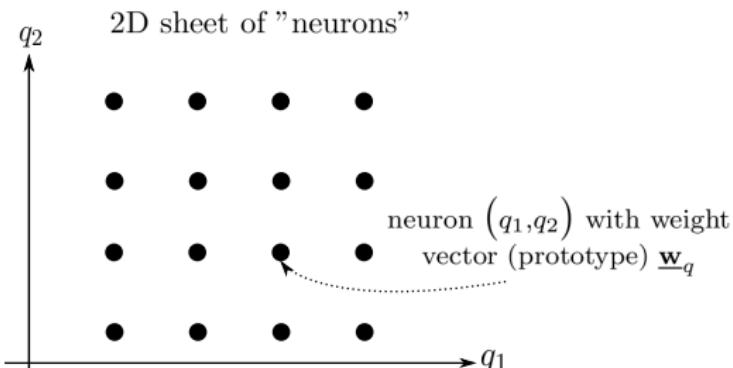
$$r \sin(2\phi), \phi \in [0, \pi[$$

ocular dominance: $z \in \{-1, 1\}$

Self-organizing maps for pairwise data

- distance matrix \mathbf{D} specifying the distances or 'dissimilarities' $d_{\alpha\alpha'}$ between p "objects" $\alpha = 1, \dots, p$
- set of M clusters (partitions) $\underline{\mathbf{q}}$ with a geometrical structure (e.g. 1-d line or 2-d grid).
- binary assignment variables (normalized):

$$m_{\underline{\mathbf{q}}}^{(\alpha)} = \begin{cases} 1, & \text{if object } \alpha \text{ belongs to cluster } \underline{\mathbf{q}} \\ 0, & \text{else} \end{cases}$$



	1	2	3	\dots	p
1	0	1.7	0.99	\dots	3.0
2	1.7	0	0.3	\dots	0.1
3	0.9	0.3	0	\dots	0.2
\vdots	\vdots	\vdots	\ddots	\ddots	\vdots
p	3.0	0.1	0.2	\dots	0

relational representation
"pairwise data"

Cost function & model selection

$$E\left[\{m_{\underline{\mathbf{q}}}^{(\alpha)}\}\right] = \frac{1}{p} \sum_{\underline{\mathbf{r}}} \frac{\sum_{\alpha, \alpha'} \left(\sum_{\underline{\mathbf{q}}} h_{\underline{\mathbf{r}}\underline{\mathbf{q}}} m_{\underline{\mathbf{q}}}^{(\alpha)} \right) \left(\sum_{\underline{\mathbf{q}}} h_{\underline{\mathbf{r}}\underline{\mathbf{q}}} m_{\underline{\mathbf{q}}}^{(\alpha')} \right) d_{\alpha\alpha'}}{\sum_{\alpha} \left(\sum_{\underline{\mathbf{q}}} h_{\underline{\mathbf{r}}\underline{\mathbf{q}}} m_{\underline{\mathbf{q}}}^{(\alpha)} \right)} \stackrel{!}{=} \min$$

- $m_q^{(\alpha)} \rightarrow \sum_{\underline{\mathbf{q}}} h_{\underline{\mathbf{r}}\underline{\mathbf{q}}} m_{\underline{\mathbf{q}}}^{(\alpha)}$
- "neighboring" clusters (w.r.t. $h_{\underline{\mathbf{r}}\underline{\mathbf{q}}}$) contribute to the total average distance
- "neighborhood preserving maps" induce lower cost

Algorithm 2: SOM for pairwise data (mean-field approximation)**Initialization:**

- choose no. M of partitions, initial (β_0) and final (β_f) noise parameters, annealing factor η , width σ of neighborhood function $h_{\underline{sq}}$, tolerance θ
- initialize mean-fields $e_{\underline{q}}^{(\alpha)}$: random numbers $\in [0, 1]$

$$\beta \leftarrow \beta_0$$

while $\beta < \beta_f$ **do** annealing

repeat EM

$$\text{compute assignment probabilities } \langle m_{\underline{q}}^{(\alpha)} \rangle_Q = \frac{\exp \left\{ -\beta \left(e_{\underline{q}}^{(\alpha)} \right)_{\text{old}} \right\}}{\sum_{\underline{r}} \exp \left\{ -\beta \left(e_{\underline{r}}^{(\alpha)} \right)_{\text{old}} \right\}} \quad \forall \underline{q}, \alpha$$

compute new mean-fields

$$(e_{\underline{q}}^{(\alpha)})_{\text{new}} = \frac{1}{p} \sum_{\underline{s}} h_{\underline{sq}} \left[\frac{1}{\sum_{\gamma} \left(\sum_{\underline{r}} h_{\underline{sr}} \langle m_{\underline{r}}^{(\gamma)} \rangle_Q \right)} \sum_{\delta} \left(\sum_{\underline{r}} h_{\underline{sr}} \langle m_{\underline{r}}^{(\delta)} \rangle_Q \right) \right]$$

$$\cdot \left\{ d_{\delta\alpha} - \frac{1}{2} \frac{1}{\sum_{\gamma} \left(\sum_{\underline{r}} h_{\underline{sr}} \langle m_{\underline{r}}^{(\gamma)} \rangle_Q \right)} \sum_{\varepsilon} \left(\sum_{\underline{r}} h_{\underline{sr}} \langle m_{\underline{r}}^{(\varepsilon)} \rangle_Q \right) d_{\varepsilon\delta} \right\} \quad \forall \underline{q}, \alpha$$

until $| (e_{\underline{q}}^{(\alpha)})_{\text{new}} - (e_{\underline{q}}^{(\alpha)})_{\text{old}} | < \theta \quad \forall \underline{q}, \alpha$

$$\beta \leftarrow \eta \beta$$

end

Comments

- replacing $h_{\underline{\text{sp}}}$ by $\delta_{\underline{\text{sp}}}$ recovers standard pairwise clustering
- "Kohonen-approximation": replace $h_{\underline{\text{sp}}}$ by $\delta_{\underline{\text{sp}}}$ for the neighborhood function (only) at \circledast in algorithm 2
 - reduction of computational cost
 - visualization properties remain
 - original algorithm suggested by T. Kohonen is recovered for squared Euclidean distances $d_{\alpha\alpha'}$ and $\beta \rightarrow \infty$
- no need for σ annealing

Toy example: noisy spiral

1000 data points:

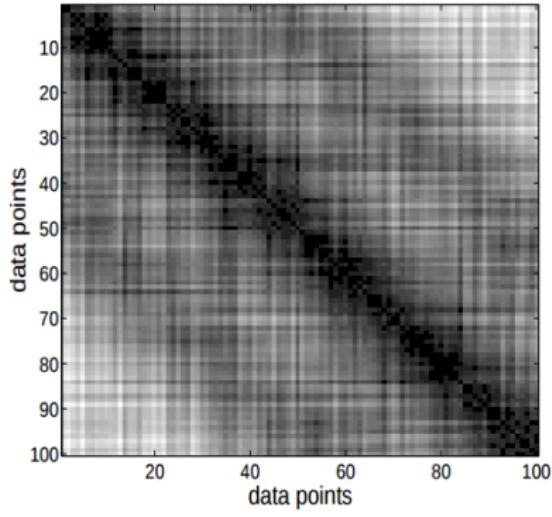
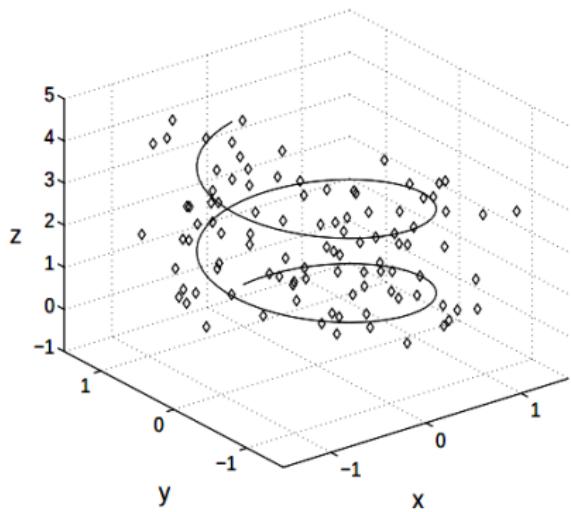
$$x = \sin \theta + \eta_x$$

$$y = \cos \theta + \eta_y$$

$$z = \frac{\theta}{\pi} + \eta_z$$

$$\theta \in [0, 4\pi]; \quad \eta \sim \mathcal{N}_{(0,0.3)}$$

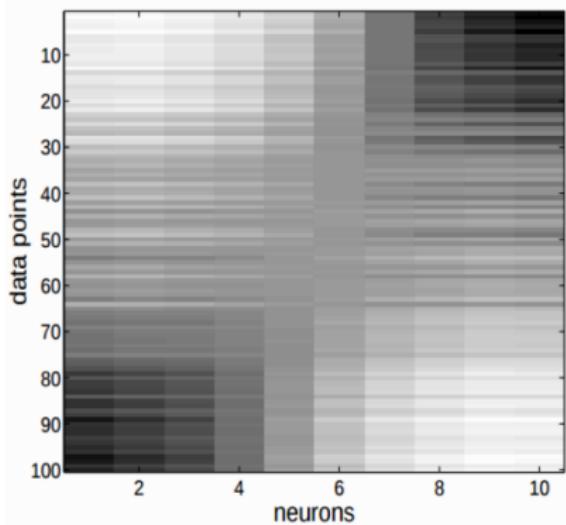
Euclidean squared distance



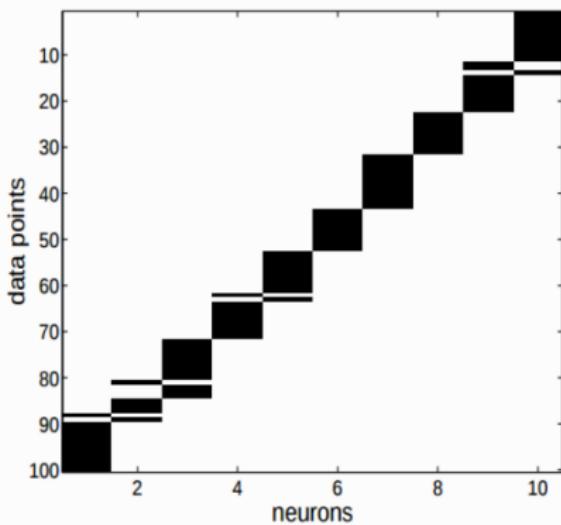
Toy example: noisy spiral

- 10 neurons in 1D
- Gaussian neighborhood function, $\sigma = 0.5$

$\beta=0.8$



$\beta=186.21$

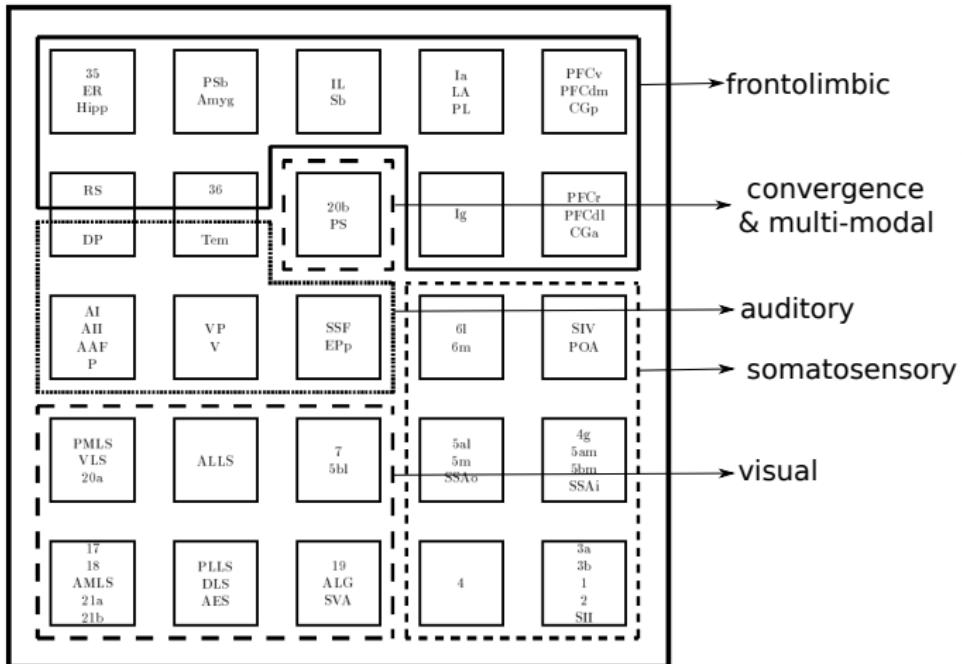


Mapping cat's cerebral cortex

AREA	DIFFERENT		SIMILAR	
	DIFF	DIFF	SIM	SIM
PMLs	17	18	PMLs	PMLs
PFLs	3	3	PFLs	PFLs
AMLS	3	3	AMLS	AMLS
VLS	1	1	VLS	VLS
DEs	0	0	DEs	DEs
DLS	0	0	DLS	DLS
21a	2	2	2	2
21b	2	2	2	2
20b	0	0	0	0
ALG	0	0	0	0
7	0	0	0	0
AES	0	0	0	0
SVA	0	0	0	0
PS	0	0	0	0
All	0	0	0	0
AAF	0	0	0	0
DP	0	0	0	0
VP	0	0	0	0
V	0	0	0	0
S&F	0	0	0	0
EPr	0	0	0	0
Tan	0	0	0	0
3a	0	0	0	0
3b	0	0	0	0
2	0	0	0	0
SII	0	0	0	0
SIV	0	0	0	0
4S	0	0	0	0
6I	0	0	0	0
6m	0	0	0	0
POA	0	0	0	0
Sam	0	0	0	0
Sal	0	0	0	0
Sai	0	0	0	0
Ssi	0	1	2	2
Ssa	0	0	0	0
SSAAo	0	0	0	0
SSAAi	0	0	0	0
PPCf	0	0	0	0
PPCfi	0	0	0	0
PPCdm	0	0	0	0
is	0	0	0	0
CGs	0	0	0	0
CGp	0	0	0	0
LA	0	0	0	0
RS	0	0	0	0
PL	0	0	0	0
IL	0	0	0	0
3S	0	0	0	0
SPsh	0	0	0	0
Sb	0	0	0	0
EP	0	0	0	0
Hipp	0	0	0	0
Amyg	0	0	0	0
17	18	18	18	18
PMLs	18	18	PMLs	PMLs
AMLS	18	18	AMLS	AMLS
VLS	18	18	VLS	VLS
DEs	18	18	DEs	DEs
DLS	18	18	DLS	DLS
21a	18	18	21a	21a
21b	18	18	21b	21b
20b	18	18	20b	20b
ALG	18	18	ALG	ALG
7	18	18	7	7
AES	18	18	AES	AES
SVA	18	18	SVA	SVA
PS	18	18	PS	PS
All	18	18	All	All
AAF	18	18	AAF	AAF
DP	18	18	DP	DP
VP	18	18	VP	VP
V	18	18	V	V
S&F	18	18	S&F	S&F
EPr	18	18	EPr	EPr
Tan	18	18	Tan	Tan
3a	18	18	3a	3a
3b	18	18	3b	3b
2	18	18	2	2
SII	18	18	SII	SII
SIV	18	18	SIV	SIV
4S	18	18	4S	4S
6I	18	18	6I	6I
6m	18	18	6m	6m
POA	18	18	POA	POA
Sam	18	18	Sam	Sam
Sal	18	18	Sal	Sal
Sai	18	18	Sai	Sai
Ssi	18	18	Ssi	Ssi
Ssa	18	18	Ssa	Ssa
SSAAo	18	18	SSAAo	SSAAo
SSAAi	18	18	SSAAi	SSAAi
PPCf	18	18	PPCf	PPCf
PPCfi	18	18	PPCfi	PPCfi
PPCdm	18	18	PPCdm	PPCdm
is	18	18	is	is
CGs	18	18	CGs	CGs
CGp	18	18	CGp	CGp
LA	18	18	LA	LA
RS	18	18	RS	RS
PL	18	18	PL	PL
IL	18	18	IL	IL
3S	18	18	3S	3S
SPsh	18	18	SPsh	SPsh
Sb	18	18	Sb	Sb
EP	18	18	EP	EP
Hipp	18	18	Hipp	Hipp
Amyg	18	18	Amyg	Amyg

Mapping cat's cerebral cortex

"Kohonen"-map:
5x5 neurons,
Gaussian
neighborhood
 $\delta_n = 0.4$



Squared Euclidean distances

$$d_{\alpha\alpha'} = \frac{1}{2} (\underline{\mathbf{x}}^{(\alpha)} - \underline{\mathbf{x}}^{(\alpha')})^2$$

cost function:

$$E[\{m_{\underline{\mathbf{q}}}^{(\alpha)}\}] = \frac{1}{p} \sum_{\underline{\mathbf{q}}, \alpha} \left(\sum_{\underline{\mathbf{p}}} h_{\underline{\mathbf{q}}\underline{\mathbf{p}}} m_{\underline{\mathbf{p}}}^{(\alpha)} \right) (\underline{\mathbf{x}}^{(\alpha)} - \underline{\mathbf{w}}_q)^2 = \frac{1}{p} \sum_{\underline{\mathbf{p}}^\alpha} m_{\underline{\mathbf{p}}}^{(\alpha)} \sum_{\underline{\mathbf{q}}} h_{\underline{\mathbf{q}}\underline{\mathbf{p}}} (\underline{\mathbf{x}}^{(\alpha)} - \underline{\mathbf{w}}_{\underline{\mathbf{q}}})^2$$

$$\underline{\mathbf{w}}_{\underline{\mathbf{q}}} = \frac{\sum_{\alpha'} \left(\sum_{\underline{\mathbf{p}}} h_{\underline{\mathbf{q}}\underline{\mathbf{p}}} m_{\underline{\mathbf{p}}}^{(\alpha')} \right) \underline{\mathbf{x}}^{(\alpha')}}{\sum_{\alpha'} \left(\sum_{\underline{\mathbf{p}}} h_{\underline{\mathbf{q}}\underline{\mathbf{p}}} m_{\underline{\mathbf{p}}}^{(\alpha')} \right)}$$

$\underline{\mathbf{w}}_{\underline{\mathbf{q}}}$: center of mass of all data which belongs to cluster weighted by the neighborhood function $h_{\underline{\mathbf{q}}\underline{\mathbf{p}}}$

On-line minimization of E^T

Algorithm 3: Online learning for SOM with Euclidean distances

Initialization of prototypes

Select learning step ε

begin

choose a random data point $\underline{x}^{(\alpha)}$

assign data point to the prototype with minimum assignment cost

$$\underline{p} = \operatorname{argmin}_{\underline{r}} \sum_{\underline{q}} h_{\underline{r}\underline{q}} \underbrace{\left(\underline{x}^{(\alpha)} - \underline{w}_{\underline{q}}^{\text{old}} \right)^2}_{*}$$

change all prototypes according to

$$\Delta \underline{w}_{\underline{q}} = \varepsilon h_{\underline{p}\underline{q}} \cdot \left(\underline{x}^{(\alpha)} - \underline{w}_{\underline{q}} \right) \text{ for all } \underline{q}$$

end
