# Unsupervised Learning

## Machine Learning for Data Science 1

Draws inference from data without labelled responses.

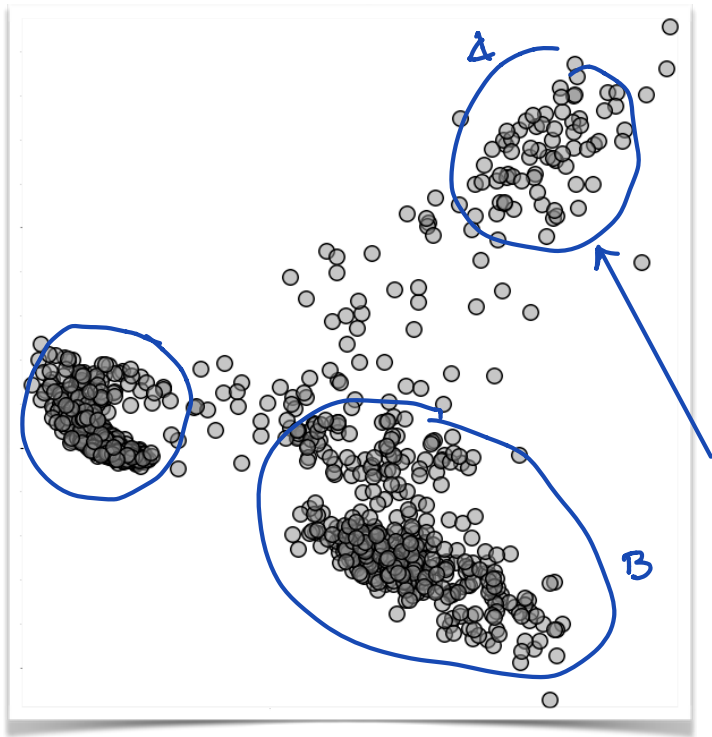Wikipedia: type of ML to find previously undetected patterns in a data with no prexisting labels and with minimal human supervision.

not true

interpretation is central
and ore is the known guidance
in these approach
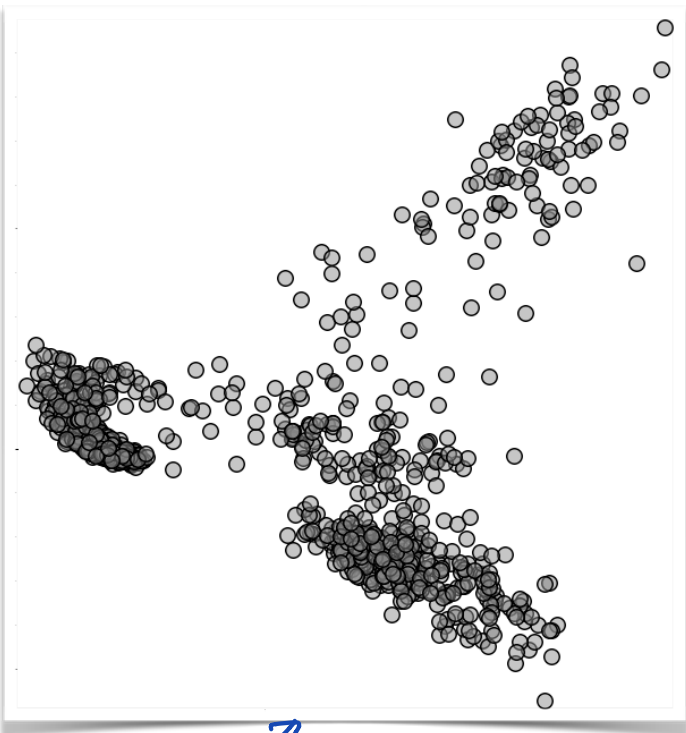
choice of parameters
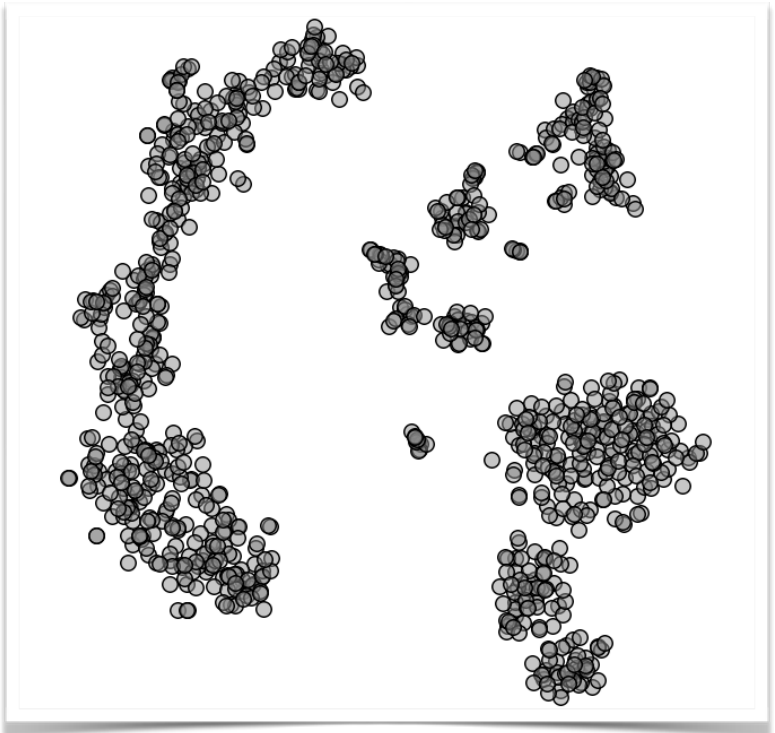approach

over fitting
overinterpretation

original space of >1000 features
principal component analysis
how are groups do we see
interpretation
Is the data that I have enough
knowledge.

PCA

t-SNE

principled
approaches to
unsupervised learning

dimensionality
reduction

- projections (change of coord. sys)
- embedding (new latent space)

2D, 1D

Clustering

grouping of data instances
similarity within
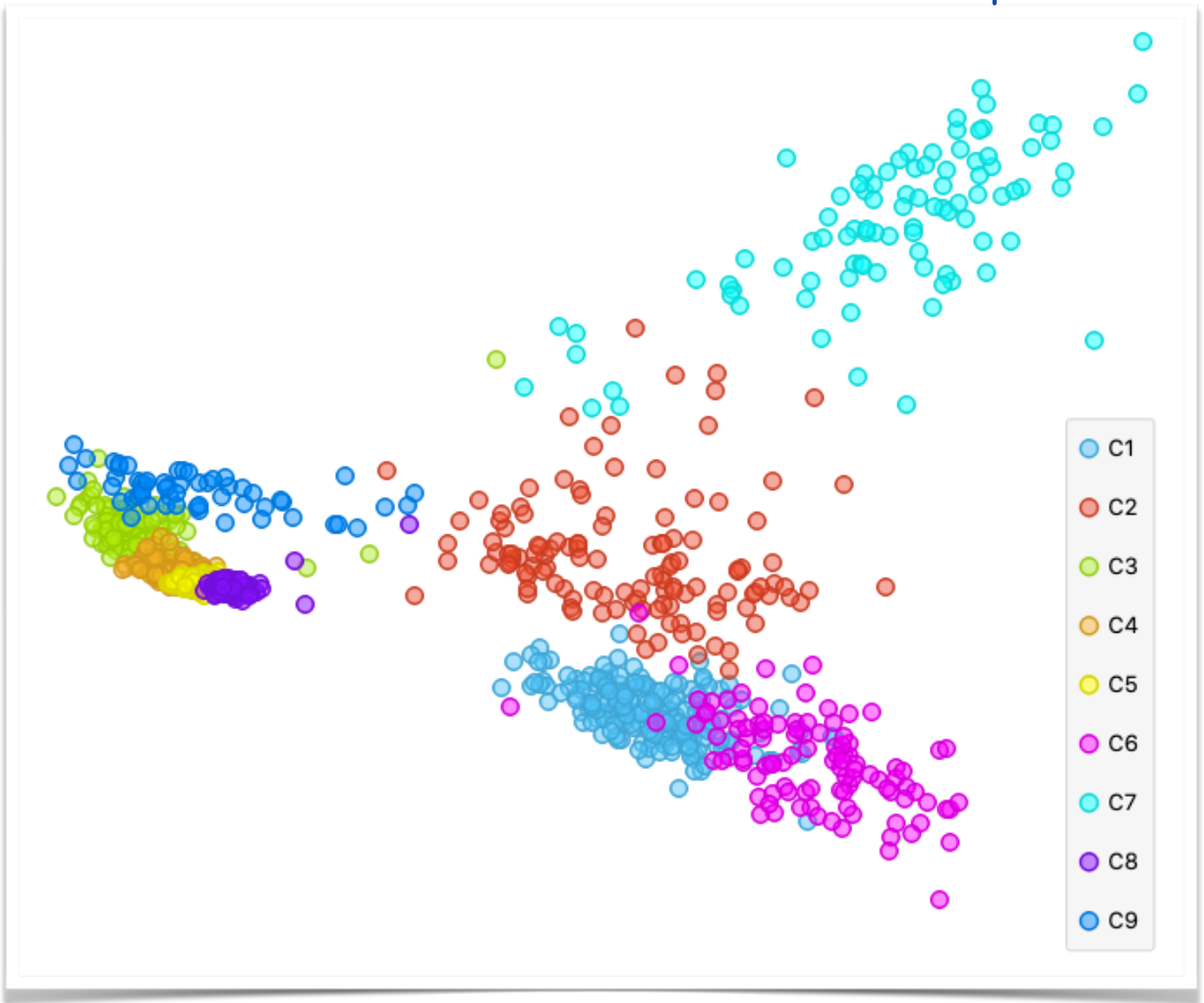dissimilarity betw. groups

still leaves room for interpretation
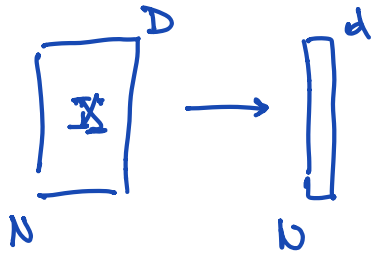
every choices of parameters

in original space

faithfulness

combination

cluster +
project, visualize

PCA, S1000

C1
C2
C3
C4
C5
C6
C7
C8
C9

# Principal Component Analysis

dimensionality reduction



$D \gg d$

$d = 2$ ← visualisation

$\underline{d = 1}$ ← sense of time
progression
development

↓ **maximise the**
**variance of projections**



$\mu_1$ ← direction of projection

$\mu_1^T \mu_1 = 1$
    or unit vector

$x \in \underline{X}$

$\mu_1^T x \in \mathbb{R}$

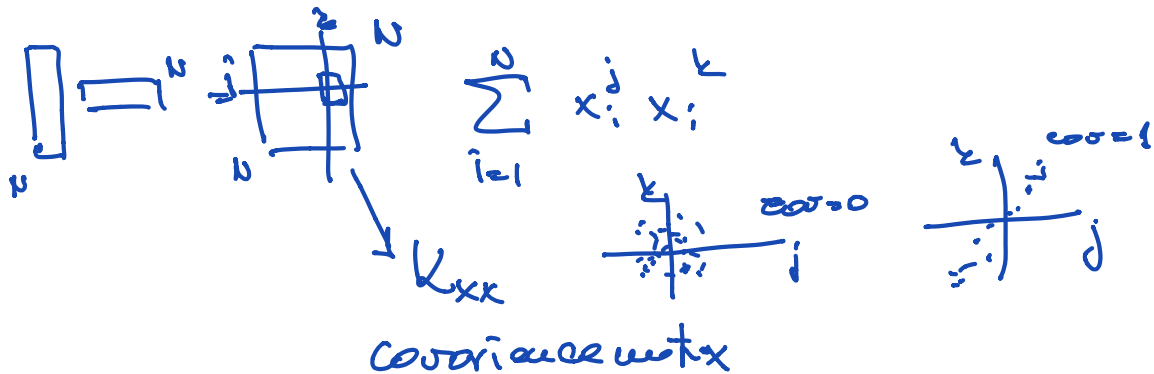$\mu_1^T \overline{x}$ , $\overline{x} = \frac{1}{m} \sum_{i=1}^{N} x_i$

We are looking for $\mu_1$ so that projected data points are maximally dispersed

$$\frac{1}{N} \sum_{i=1}^{N} \left( \mu_1^T x_i - \mu_1^T \bar{x} \right)^2 = \underline{Var\left( \mu_1^? X^T \right)}$$

$$Var\left( \mu_1^T X^T \right) = \frac{1}{N} \sum \left( \underline{\mu_1^T x_i} \ \underline{\mu_1^T x_i} - 2 \underline{\mu_1^T x_i} \underline{\mu_1^T \bar{x}} + \underline{\mu_1^T \bar{x}} \ \underline{\mu_1^T \bar{x}} \right)$$

$$\underbrace{\left( \mu_1^T x_i \right)^T = x_i^T \mu_1^T}$$

$$= \mu_1^T \left( \frac{1}{N} \sum \left( x_i x_i^T - 2 x_i \bar{x}^T + \bar{x} \bar{x}^T \right) \right) \mu_1$$

$$= \mu_1^T \left( \underline{\frac{1}{N} \sum_{i=1}^{N} \left( x_i - \bar{x} \right) \left( x_i - \bar{x} \right)^T} \right) \mu_1$$

$$\sum_{i=1}^{N} x_i^j x_i^k$$

$K_{xx}$

covariance matrix

$$\text{Var}\left(\mu_1^T X^T\right) = \underline{\mu_1^T K \mu_1} \longrightarrow \text{For PCA} \\ \text{maximize}$$

$$\text{constrain: } \underline{\mu_1^T \mu_1 = 1}$$

Lagrange:

$$f(\mu_1) = \mu_1^T K \mu_1 + \lambda_1 \left(1 - \underline{\mu_1^T \mu_1}\right)$$

$$\nabla f(\mu_1) = \underline{K \mu_1 - \lambda_1 \mu_1 = 0}$$
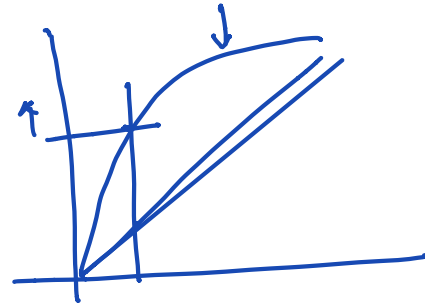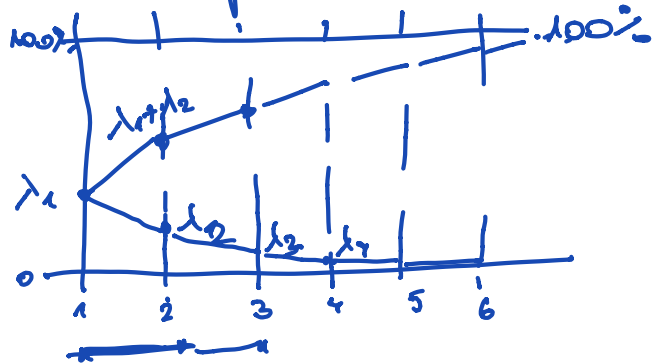
$$\underline{K \mu_1 = \lambda_1 \mu_1}$$

eigenvector

eigenvalue

$$\mu_1^T K \mu_1 = \text{Var}\left(\underline{\mu_1^T X^T}\right) = \mu_1^T \underline{\lambda_1} \mu_1 = \underline{\lambda_1}$$

$$\text{PCA: 1st component: } \mu_1 \text{ of } K, \lambda_1$$

# Scree diagram

variance explained



$\lambda_1$

$\lambda_1 + \lambda_2$

100%

$\lambda_2$ $\lambda_3$ $\lambda_4$

0

1 2 3 4 5 6

.100%

$\underline{80\% , 90\%}$

— how many dimensions?
— are the first 2 dim. "information enough"

— power method
— gram-schmidt orthog. |

**PCA**

Components Selection

Components: 2

Explained variance: 47%

Options

☑ Normalize variables

Show only first 20

☑ Apply Automatically

cumulative variance

0.474

0.195

component variance

Principal Components

Proportion of variance

# Singular Value Decomposition

|  |  | Beer | | | SF | | | | |
|---|---|---|---|---|---|---|---|---|---|
|  |  | A star is Born | Titanic | When Sally met Harry | Matrix | Fractle. |  | σ |  |
| women | Rory | 1 | 2 | 1 | 0 | 0 | — | 1 | 0 |
|  | Eoc | 5 | 6 | 7 | 0 | 0 | — | 1 | 0 |
|  | JoAnn | 2 | 3 | 1 | 0 | 0 | — | 1 | 0 |
|  | Jim | 1 | 2 | 1 | 1 | 0 | — | 3 | 0 |
| men | Fritz | 0 | 0 | 0 | 5 | 7 | — | 0 | 1 |
|  | Bill | 6 | 0 | 0 | 4 | 2 | — | 0 | 1 |

$$
\rightarrow \quad \underline{1 \quad \underline{1} \quad 1 \qquad \qquad 0 \quad 0}
$$

$$
\underline{0 \quad 0 \quad 0 \qquad \qquad \underline{1 \quad 1}}
$$

$\nearrow$ S $\overline{\text{truncated}}$ decomposition

$$\underline{X} = U \, \Sigma \, V^T$$

$N \times D$  $N \times R$  $R \times R$  $R \times D$

$U^T U = I$

$V^T V = I$

$\Sigma : \text{diagonal}$

$$\hat{X} =$$

coeff. rows

mean, var

$U$

$$\underline{X} = U \Sigma V^T_{\cdot}$$

$$\underline{X}^T = V \Sigma^T U^T = V \Sigma U^T$$

$$X^T X = V \Sigma^T U^T U \Sigma V^T = U \Sigma^2 V^T$$

$$\underbrace{\qquad}_{I}$$

$$\boxed{X^T X} \, \boxed{V} = V \Sigma^2 V^T V = V \boxed{\Sigma^2}$$

PCA : projection

embedding : Multidimensional Scaling MDS



$$p_{ij} = \| x_i - x_j \|$$

Any kind of distance

Preserve the distances !

$$\equiv \quad J(\Theta) = \sum_{i \neq j}^{N} \left( p_{ij} - g_{ij} \right)^2 \quad : \text{minimize}$$

$$x_i \longmapsto \Theta_i$$

$$g_{ij} = \| \Theta_i - \Theta_j \|$$

$$\begin{bmatrix} \chi \end{bmatrix}^D \rightarrow \begin{bmatrix} \theta \end{bmatrix}^2$$

$N = 100$          $200$

$$\dfrac{200}{\phantom{x}}$$

$$\frac{\partial J(\theta)}{\partial \phi_1^{(i)}} = \sum_{i \neq j} \frac{\partial J(\theta)}{\partial g_{ij}} \frac{\partial g_{ij}}{\partial \phi_1^{(i)}}$$

$$= 2 \sum_{i \neq j} \underbrace{\left( g_{ij} - p_{ij} \right)}_{error} \quad \underbrace{\frac{\phi_1^{(i)} - \phi_2^{(j)}}{g_{ij}}}_{} \quad \bigg| \quad \text{direction of change}$$

Solution: majoisation, SMACOF
                    ↳ speedup
                    convergence

$\phi_2$

$\phi_1$

seasnake

porpoise
dolphin

platypus

- amphibian
- bird
- fish
- insect
- invertebrate
- mammal
- reptile

# Stochastic Neighbor Embedding — SNE
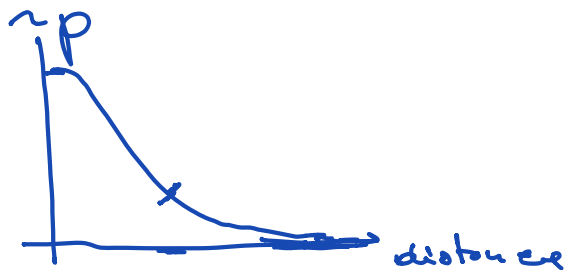
2002    Hinton & Roweis

2008    von der Maaten     t-SNE

## Idea

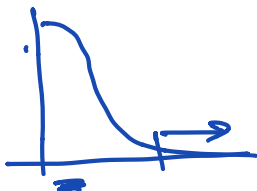data instances close to each other
in the original space
Should be close in embedding

$$p_{j|i} = \frac{e^{-\frac{\|x_i - x_j\|^2}{2\sigma_i^2}}}{\sum\limits_{k \neq i} e^{-\frac{\|x_i - x_k\|^2}{2\sigma_i^2}}}$$

perplexity

original space

distance

latent space $\phi_2$



$$g_{j|i} = \frac{e^{-\|\phi_i - \phi_j\|^2}}{\sum\limits_{k \neq i} e^{-\|\phi_i - \phi_k\|^2}}$$

$\phi_1$

SNE

$$J(\theta) = \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{g_{j|i}} =$$
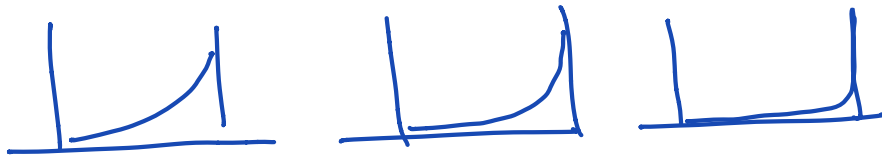
$$= \sum_i KL(P_i \| Q_i)$$

$$\frac{\partial J}{\partial \theta^{(i)}} = 2 \sum_j \left( p_{j|i} - g_{j|i} + p_{ij} - g_{i|j} \right) \left( \theta^{(i)} - \theta^{(j)} \right)$$
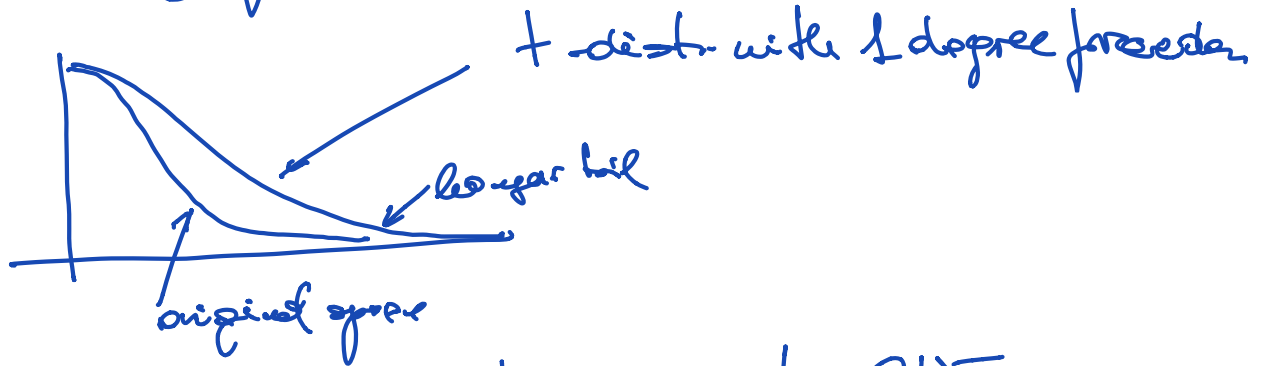
SNE

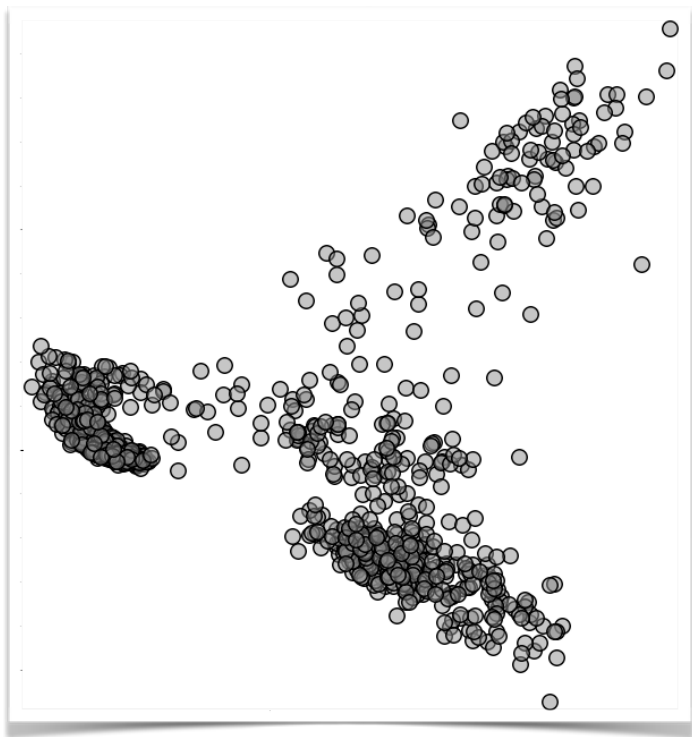multi dimension
space

1000 features

→

two dimension
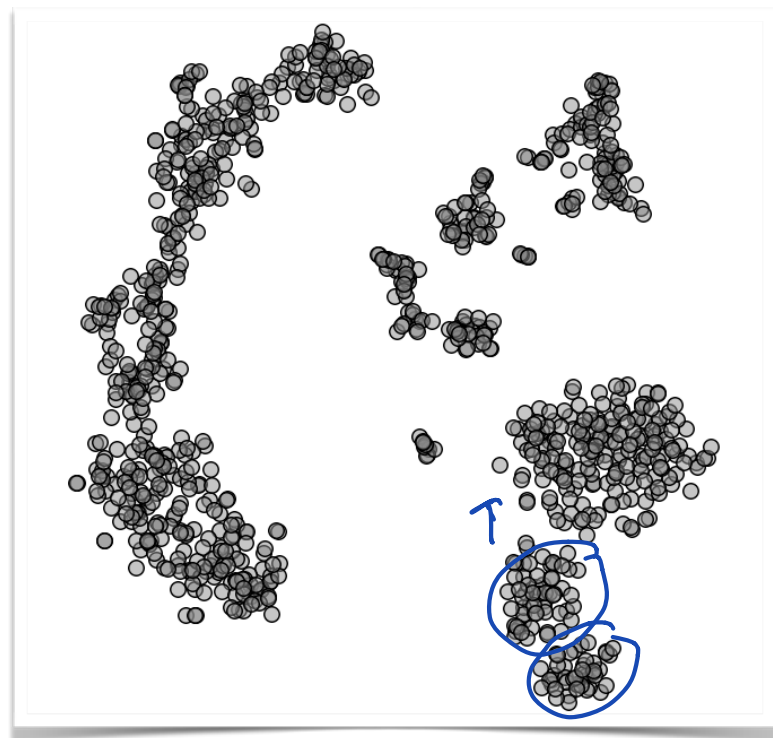space

2

The crossing problem

t distr with 1 degree freedom

longer tail

original space

$$g_{ij} = \frac{(1 + \|x_i - x_j\|^2)^{-1}}{\sum_{k \neq i} (1 + \|x_i - x_k\|^2)^{-1}}$$

t -SNE

$$\frac{\partial \mathfrak{I}}{\partial \theta^{(i)}} = 4 \sum \left( \underbrace{p_{ij} - q_{ij}}_{\text{error}} \right) \left( \underbrace{\phi_i - \phi_j}_{\text{direction}} \right) \left( \underbrace{1 + \| \phi_i - \phi_j \|^2}_{\text{weight}} \right)^{-1}$$
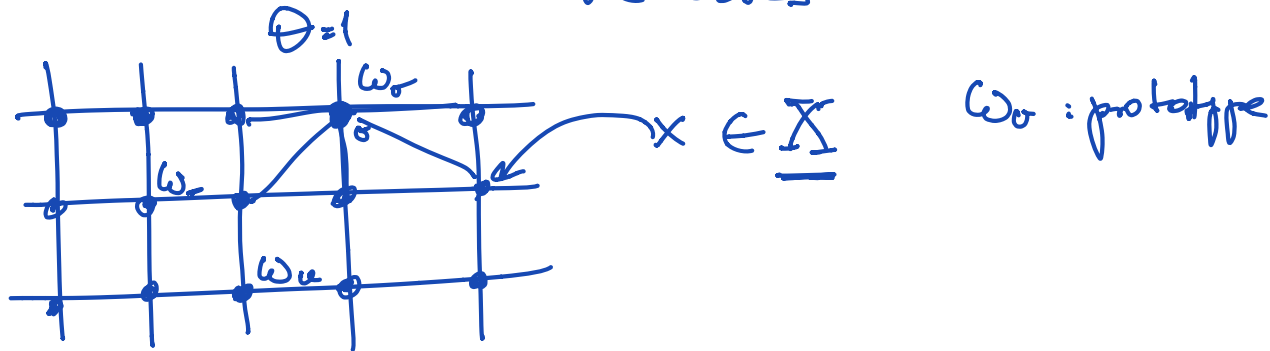
PCA
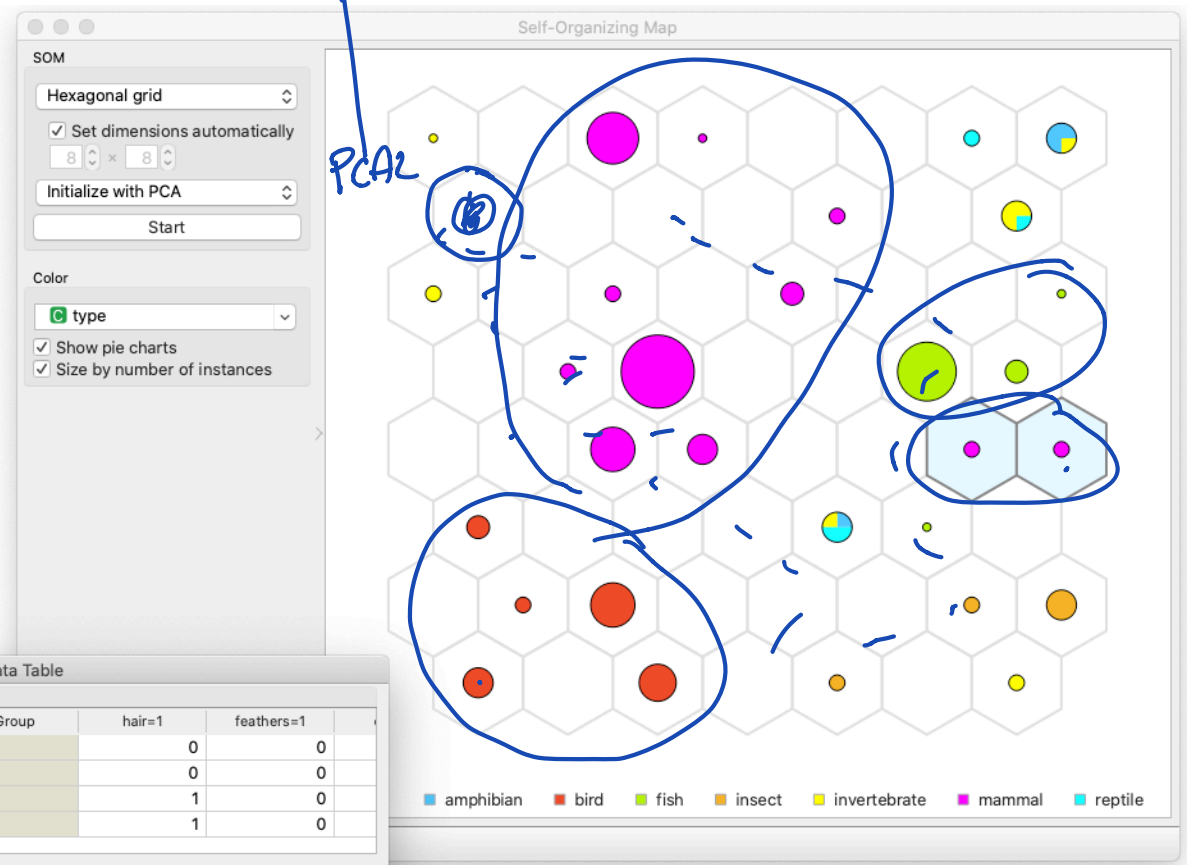
t-SNE

# Self-Growing Maps

## Kohonen maps
## networks

$\Theta = 1$



$x \in \underline{X}$

$\omega_\sigma$ : prototype

1. initialize $\omega_\sigma$
2. randomly pick $x \in \underline{X}$
3. find noble $\mu$.

$$\mu = \arg\max_\sigma \| x - \omega_\sigma \| \quad , \text{ we place } x \text{ in node } \mu$$

4. $\omega_\sigma \leftarrow \omega_\sigma + \Theta(\underline{\mu}, \underline{\sigma}) \cdot \alpha \left( \underline{\omega}_\sigma - \underline{x} \right)$

5. if stopping condition not met
goto 2