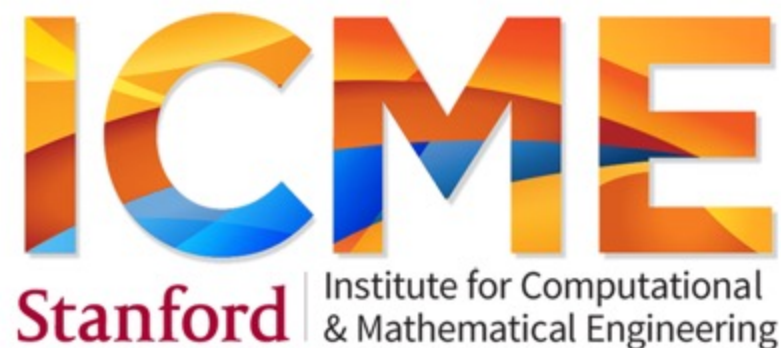


Welcome to CME 250 Introduction to Machine Learning!

Spring 2020 – Online version
April 16th 2020

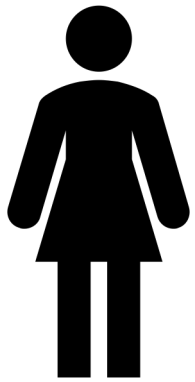


Today's schedule

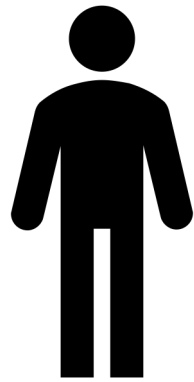
- Dimensionality reduction: Why is it useful?
- Reduce number of features
 - PCA as Maximal Variance Projection
 - Other (really useful) matrix decompositions
- Creating features out of similarity
 - Spectral clustering

Let's get to know each other...

Breakout room



You



Another
student

Name

Location

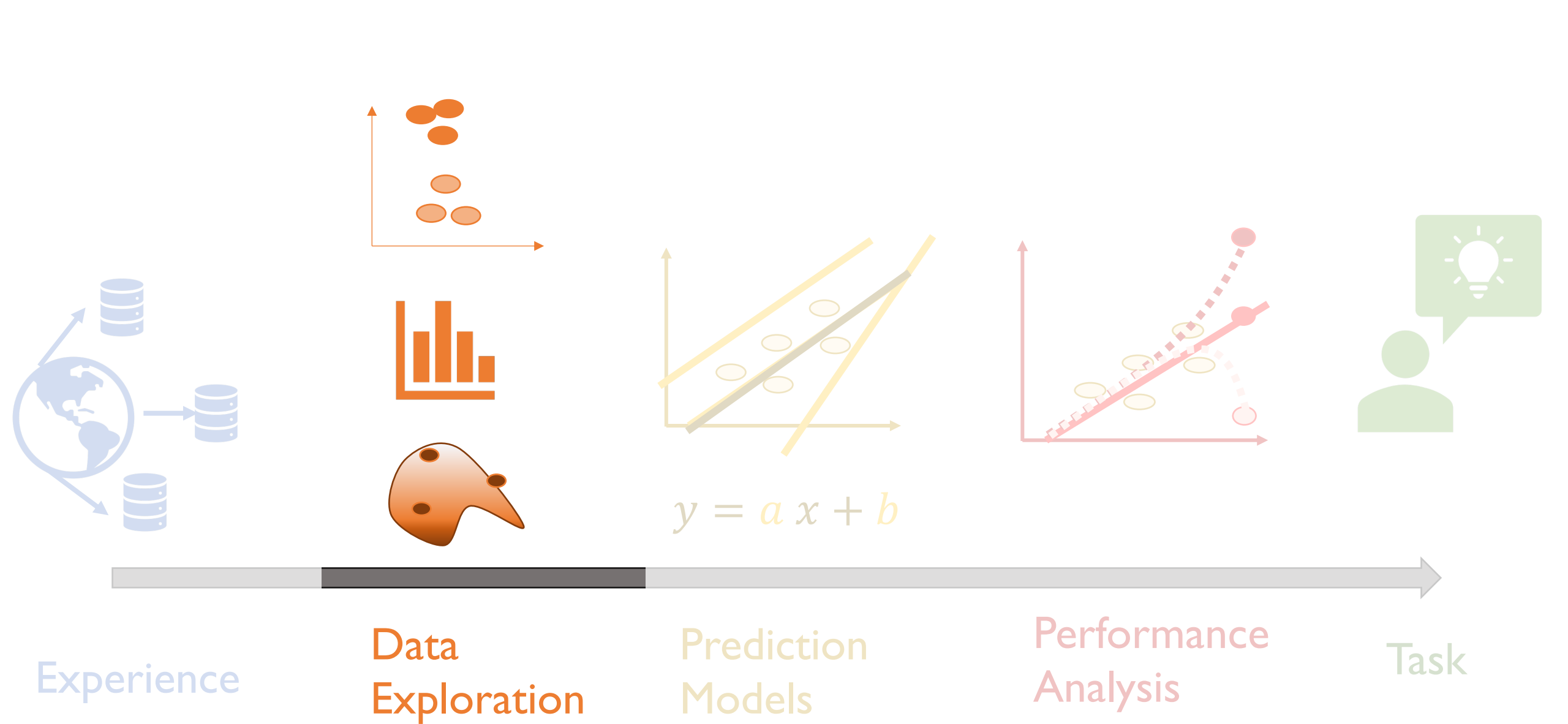
Department

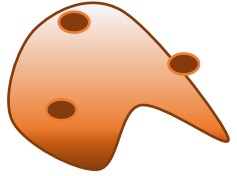
Year

Describe where you are:
kitchen, car, coffee shop ...

3 mins

Chat/Audio/Video





Unsupervised Learning Part II: Dimensionality Reduction



Data
Exploration

Introduction to Statistical Learning

Chapter 10.2: Principal Component Analysis,
10.4: Practical Lab in R

Elements Statistical Learning

Chapter 14.5 – 14.9: Different methods of DR

Recommended:

Ten quick tips for effective dimensionality reduction
Lan Huong Nguyen, Susan Holmes.(2019)

<https://doi.org/10.1371/journal.pcbi.1006907>

Last class recap

Unsupervised Learning

Patterns + Properties in Data

Clustering

Subgroups of samples

Hard

Soft

K-means
Prototypes

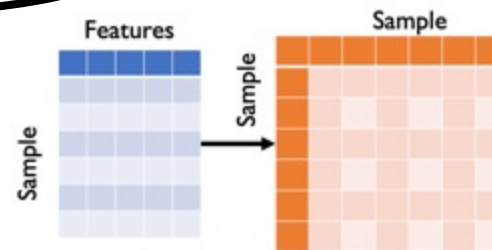
Hierarchical
Dendrograms

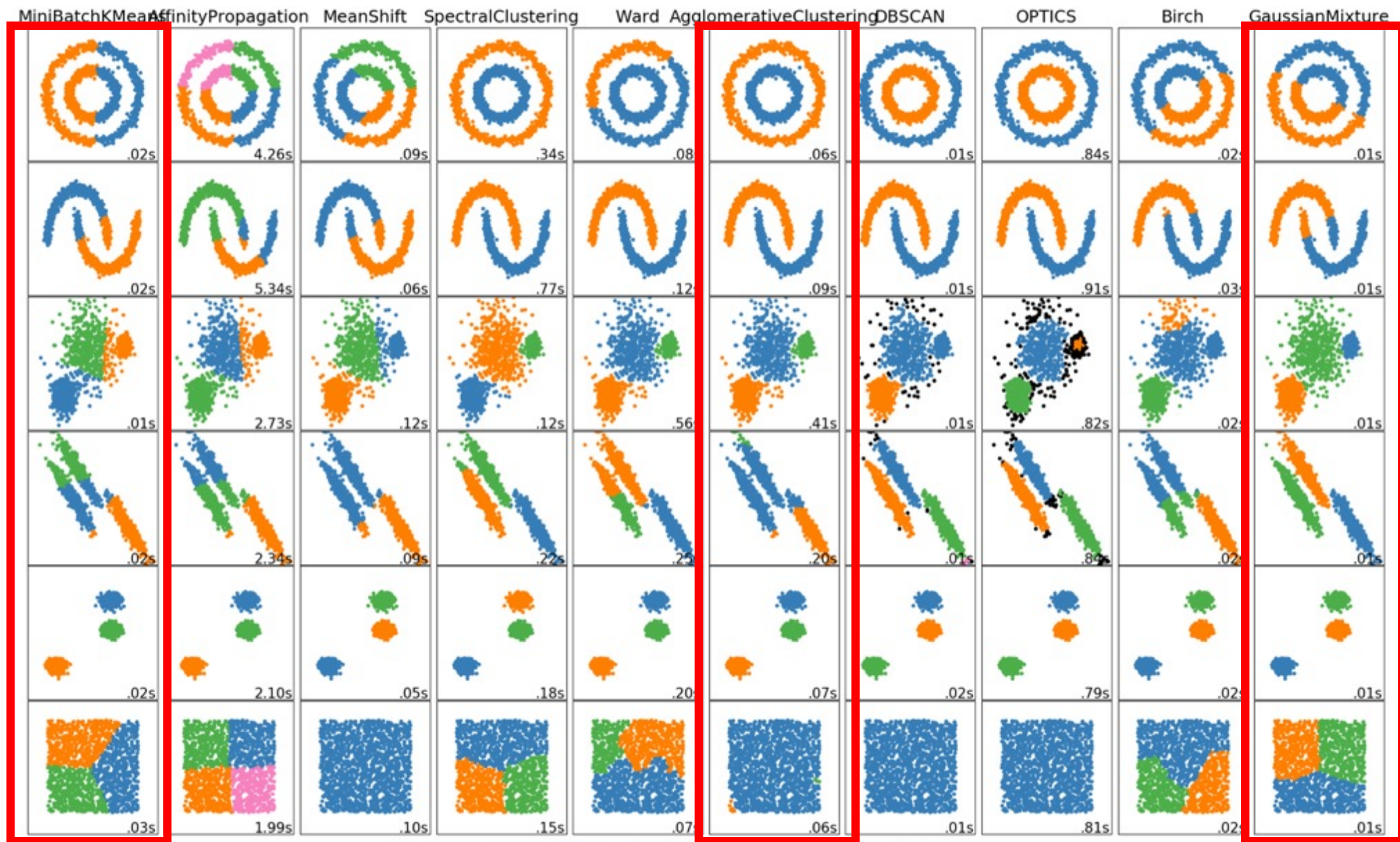
GMMs
Mixture Distribution

Dissimilarity or Similarity

Dimensionality Reduction

Reduce # variables

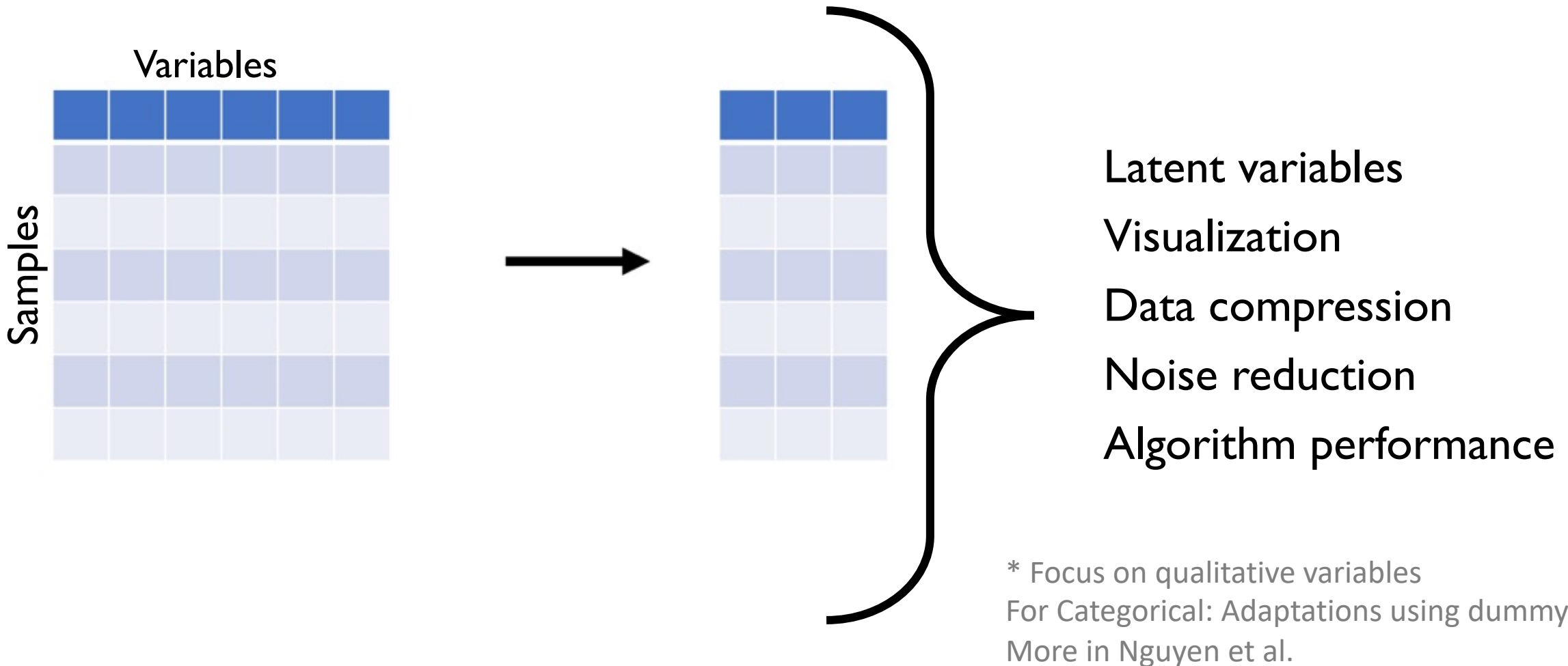




A comparison of the clustering algorithms in scikit-learn

What is dimensionality reduction?

Reduce # of variables **preserving** most of the **information**



PCA: Principal Component Analysis

The most common DR method

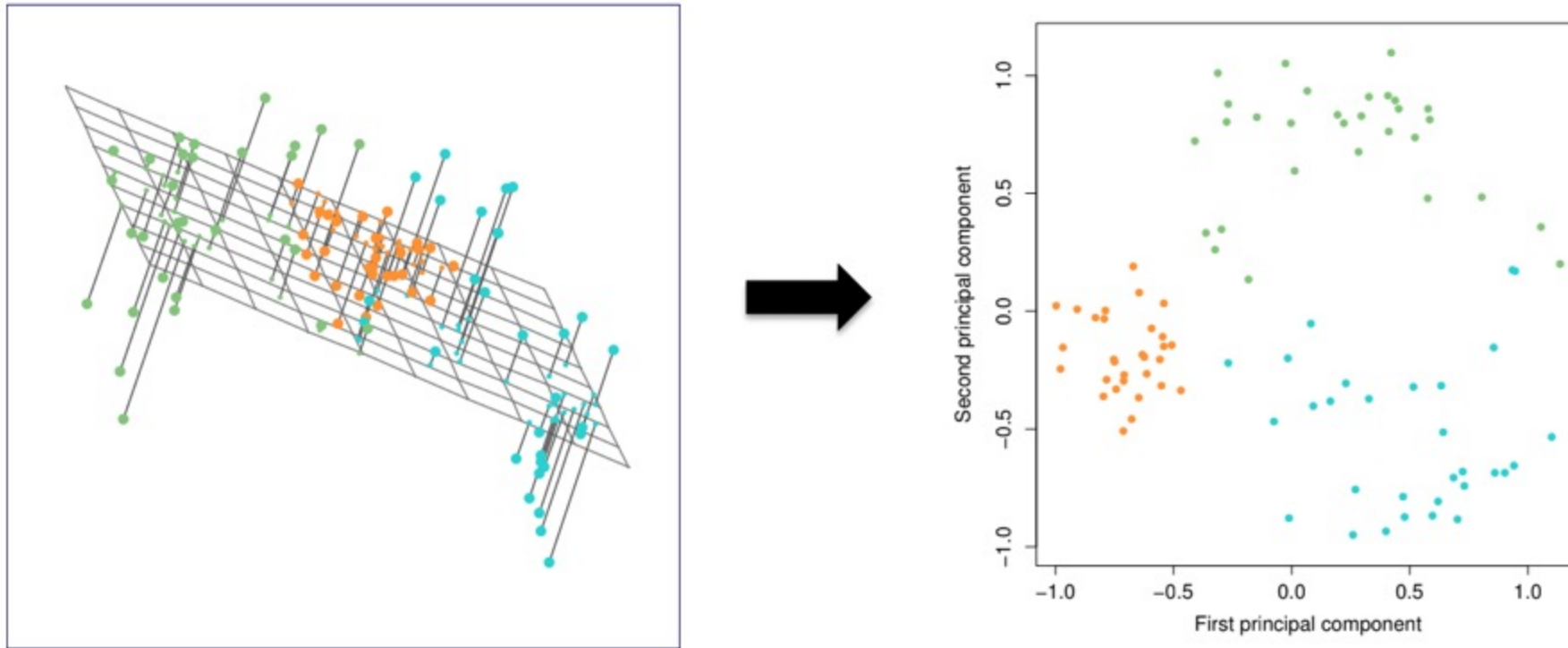


Figure 10.2 ISL (2013)

Goal: Find the projection that maximizes variation

PCA: What is 1st principal component?

Variables: X_1, X_2, \dots, X_p

1st principal
component

$$Z_1 = \phi_{11}X_1 + \phi_{21}X_2 + \dots + \phi_{p1}X_p$$

the largest
variance

PCA: Computing 1st principal component

For each sample: $x_1^{(i)}, x_2^{(i)}, \dots, x_p^{(i)}$

1) **Center** observations $\tilde{x}_j^{(i)} = x_j^{(i)} - \bar{x}_j$

2) We look for $z_1^{(i)} = \phi_{11}\tilde{x}_1^{(i)} + \phi_{21}\tilde{x}_2^{(i)} + \dots + \phi_{p1}\tilde{x}_p^{(i)}$

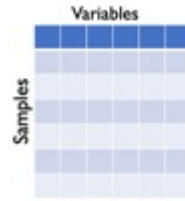
by solving

$$\max_{\phi_{11}, \dots, \phi_{p1}} \frac{1}{N} \sum_{i=1}^N \left(z_1^{(i)} \right)^2 \quad \text{s.t.} \quad \sum_{j=1}^p \phi_{p1}^2 = 1 \quad \text{normalized}$$

variance

PCA: Computing 1st principal component (In matrix form)

Observation vs. feature X



Variables			
Samples			

1) **Center** observations $\tilde{X} = \left(I - \frac{1}{n} \mathbf{1}\mathbf{1}^T \right) X$

2) We look for $z_1 = \tilde{X}\phi_1$

by solving

$$\max_{\phi_1} \phi_1^T (\tilde{X}^T \tilde{X}) \phi_1 \quad \text{s.t.} \quad \|\phi_1\|_2 = 1 \quad \text{normalized}$$

variance

Solution = Largest **eigenvector** of $\tilde{X}^T \tilde{X}$ **normalized**

PCA: What are first k principal components?

Variables: X_1, X_2, \dots, X_p

k principal
components

$$Z_1 = \phi_{11}X_1 + \phi_{21}X_2 + \dots + \phi_{p1}X_p$$

...

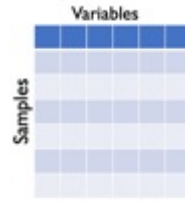
$$Z_k = \phi_{1k}X_1 + \phi_{2k}X_2 + \dots + \phi_{pk}X_p$$

the largest
variance
and

uncorrelated

PCA: Computing k principal components (In matrix form)

Observation vs. feature X



1) **Center** observations $\tilde{X} = \left(I - \frac{1}{n} \mathbf{1}\mathbf{1}^T \right) X$

2) We look for $Z_k = \tilde{X}V_k$

by solving

$$\max_{V_k} \text{tr} (V_k^T \tilde{X}^T \tilde{X} V_k)$$

variance

s.t.

$$V_k^T V_k = I$$

normalized
and
uncorrelated

Solution = **orthogonal** projection into k largest **eigenvectors** of $\tilde{X}^T \tilde{X}$

PCA: How to compute eigenvectors?

Eigenvalue Decomposition

$$\mathbf{A} = \mathbf{B}\mathbf{\Lambda}\mathbf{B}^{-1} \quad \mathbf{\Lambda}: \text{diagonal}$$

For $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$: $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} = \mathbf{V}\mathbf{D}^2\mathbf{V}^T$ s.p.d: symmetric positive definite

Singular Value Decomposition

$$\tilde{\mathbf{X}} = \mathbf{U}\mathbf{D}\mathbf{V}^T$$

$$d_{11} \geq d_{22} \geq \dots \geq d_{NN} \geq 0$$

$$\mathbf{U}^T \mathbf{U} = \mathbf{I}, \mathbf{V}^T \mathbf{V} = \mathbf{I}$$

PCA: The SVD and the principal components

We look for $\mathbf{Z}_k = \tilde{\mathbf{X}}\mathbf{V}_k$ uncorrelated and with largest variance

SVD: $\tilde{\mathbf{X}} = \mathbf{U}\mathbf{D}\mathbf{V}^T$

K principal components $\left\{ \begin{array}{l} \text{Take } \mathbf{V}_k = \text{first } k \text{ columns of } \mathbf{V} \\ \mathbf{Z}_k = \mathbf{U}_k\mathbf{D}_k \end{array} \right.$ Truncated SVD

PCA = SVD applied to centered matrix $\tilde{\mathbf{X}} = \left(\mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^T \right) \mathbf{X}$

The SVD and data compression

For any A , **truncated SVD** minimizes

$$\min_{U_k, D_k, V_k} \|A - U_k D_k V_k^T\|_2 \quad \begin{array}{ccc} U_k \in R^{N \times k} & D_k \in R^{k \times k} & V_k \in R^{p \times k} \\ \text{Orthogonal} & \text{Diagonal} & \text{Orthogonal} \end{array}$$



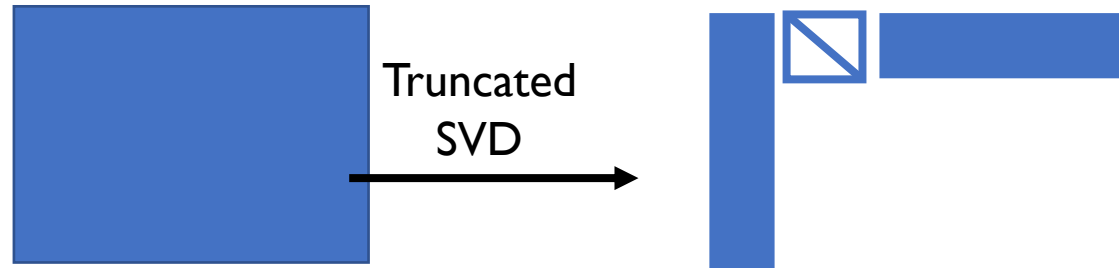
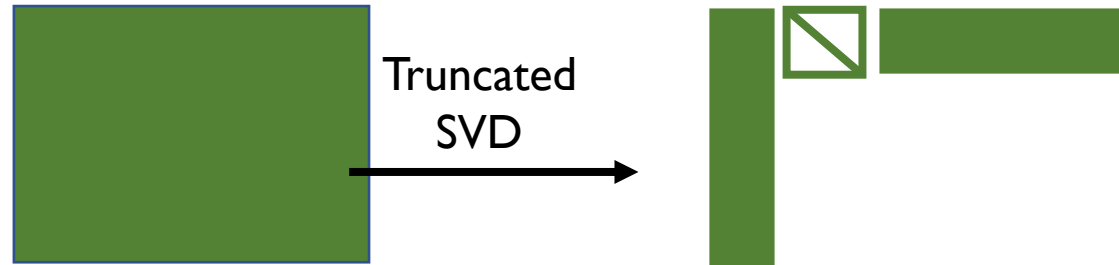
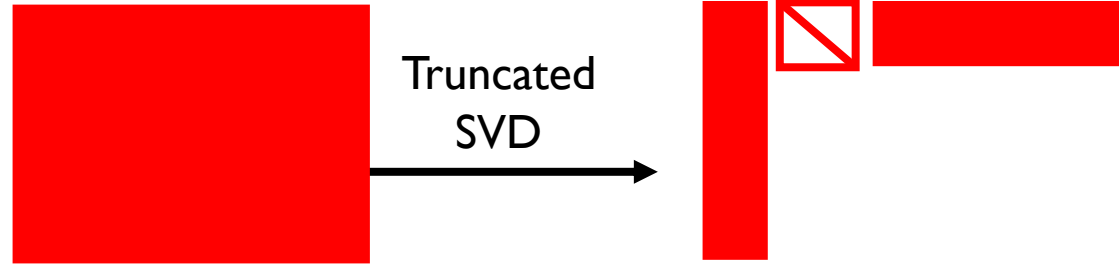
Truncated
SVD



12 times
smaller

<http://timbaumann.info/svd-image-compression-demo/>

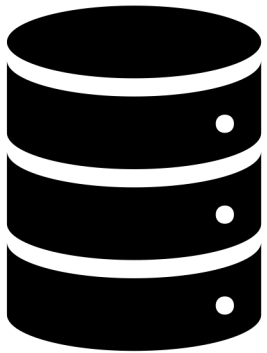
The SVD and data compression



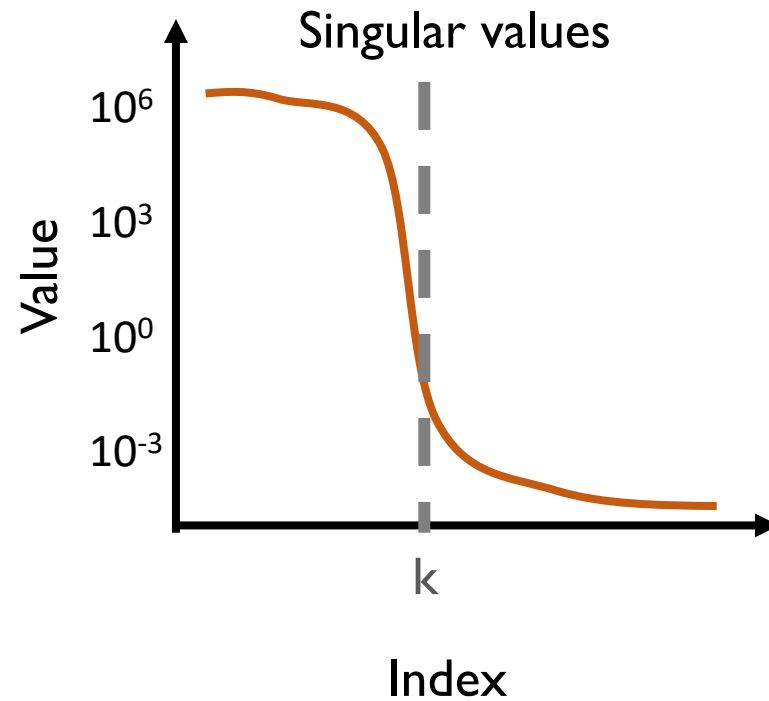
12 times smaller

PCA: How many k do we pick?

Space available



“Clear” cut

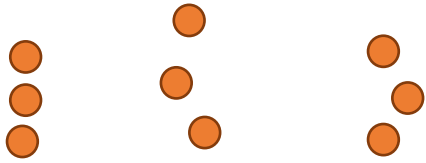


% Variance

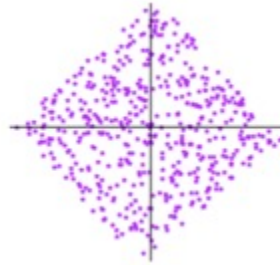
$$\frac{\sum_{i=1}^k d_i^2}{\sum_{i=1}^p d_i^2} = \frac{\text{Variance Explained}}{\text{Total Variance}}$$

PCA: SVD Challenges

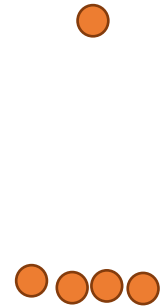
Not all variables have the same influence



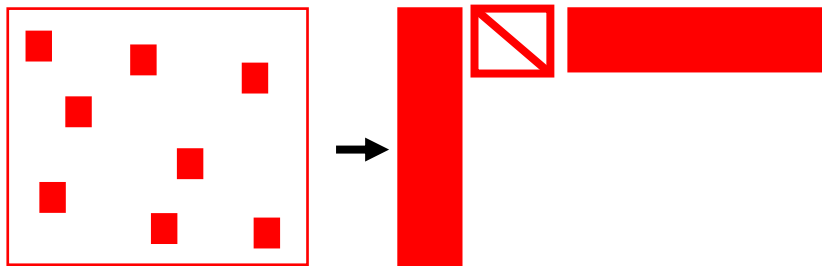
Uncorrelation, not independence



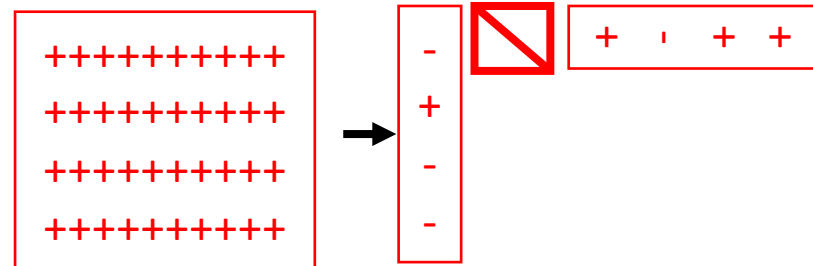
Non robust



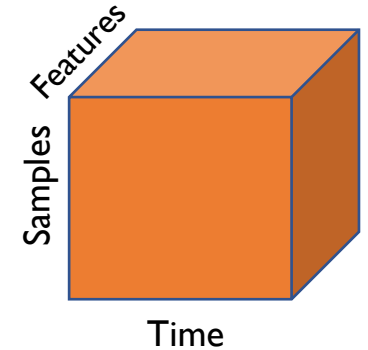
Non suitable for sparse data, expensive to compute



Does not preserve non negativity



Only 2D arrays



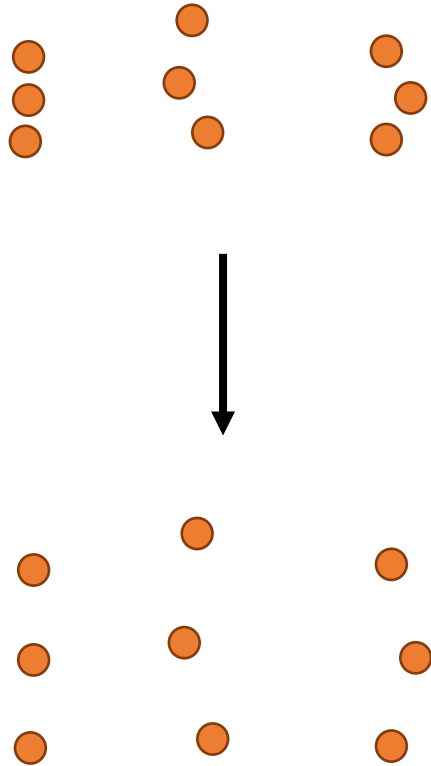
PCA: SVD Challenges

$$\min_{U_k, D_k, V_k} \|A - U_k D_k V_k^T\|_2 \quad \begin{matrix} U_k \in R^{N \times k}, & D_k \in R^{k \times k}, & V_k \in R^{p \times k} \\ \text{Orthogonal} & \text{Diagonal} & \text{Orthogonal} \end{matrix}$$



Makes optimization
more challenging

PCA: Not all variables have the same influence

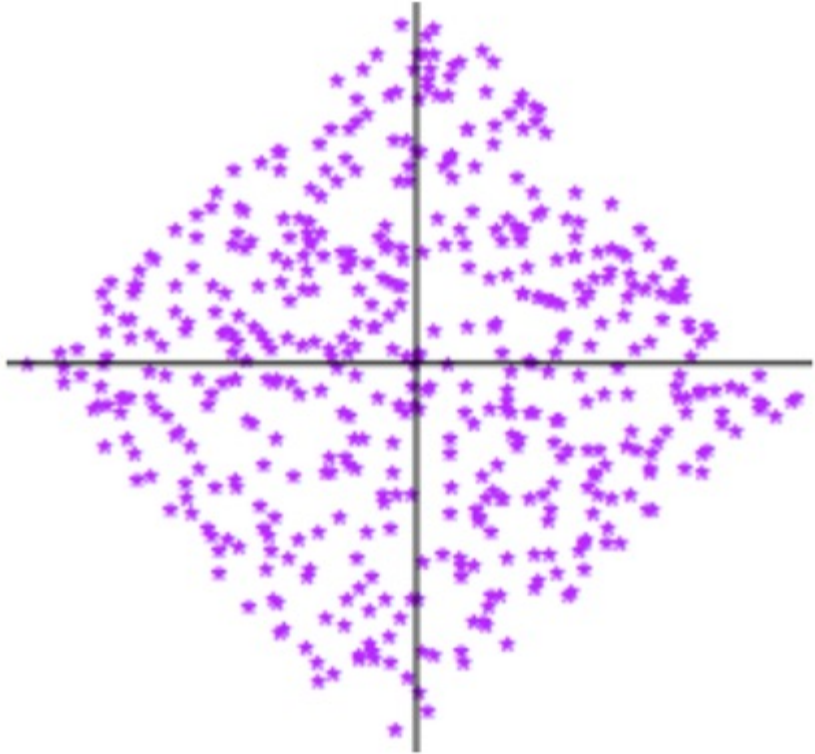


Weighted PCA

$$\tilde{X} \mathbf{W} \approx \mathbf{U}_k \mathbf{D}_k \mathbf{V}_k^T$$

e.g. \mathbf{W} diagonal $w_i = \frac{1}{\sqrt{\text{var}(x_i)}}$

PCA: Components are uncorrelated, but not independent

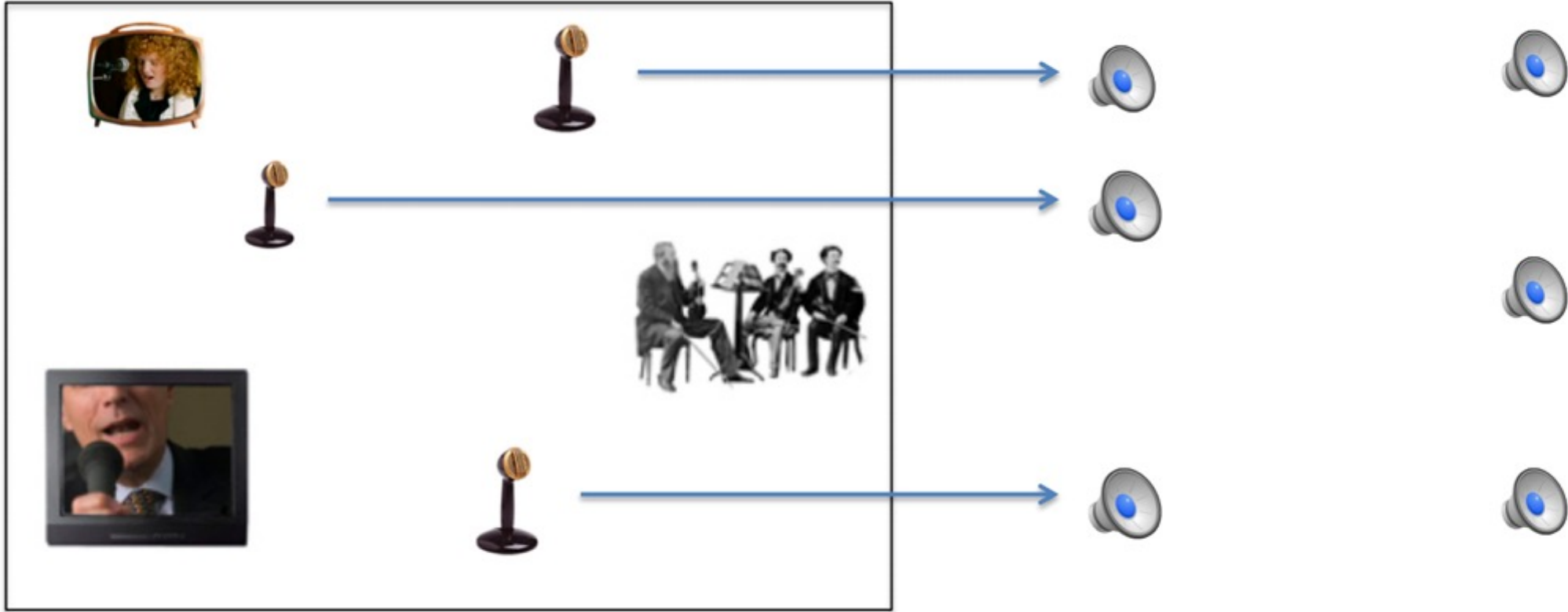


Independent
Component
Analysis ICA

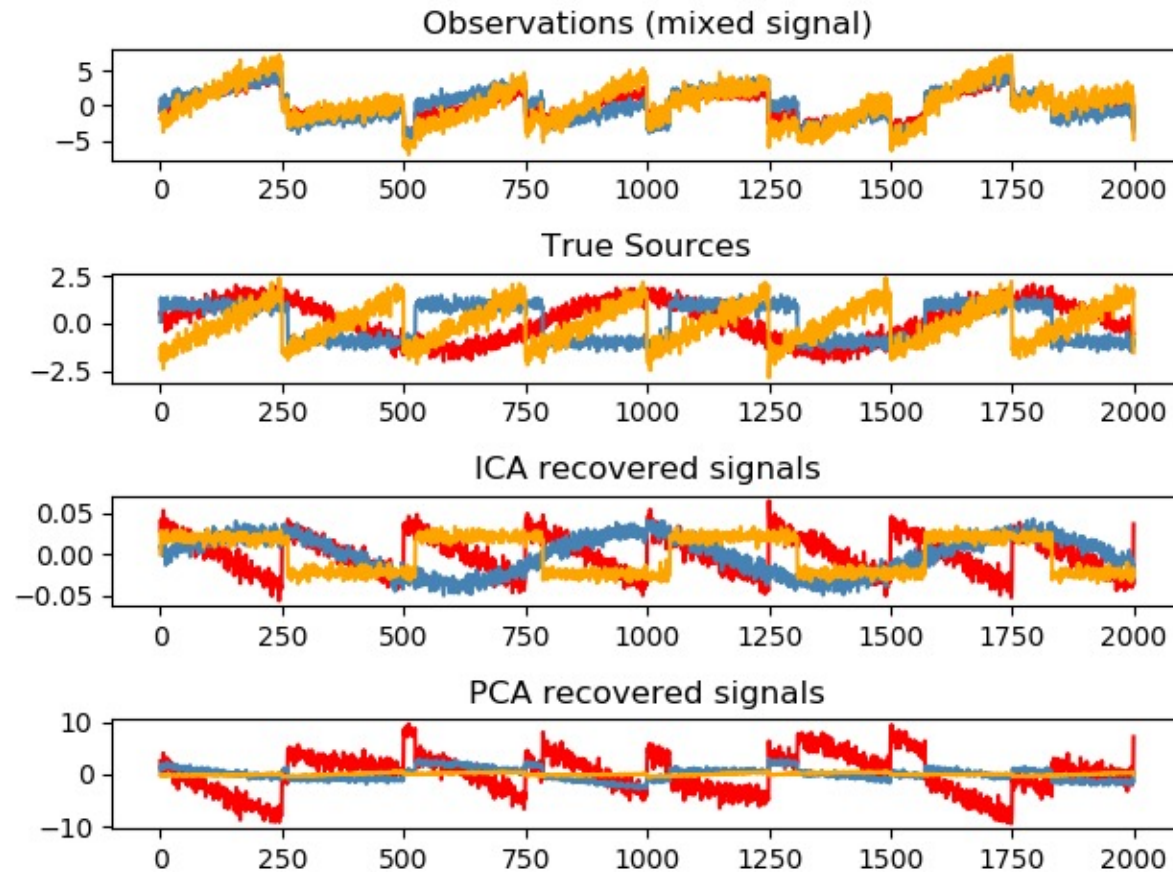
$$\tilde{\mathbf{X}} \Sigma^{-1/2} \approx \mathbf{U} \mathbf{A}$$

\mathbf{U} low entropy =
non-Gaussian
projection

PCA: Components are uncorrelated, but not independent



PCA: Components are uncorrelated, but not independent



https://scikit-learn.org/stable/auto_examples/decomposition/plot_ica_blind_source_separation.html

PCA: Non robust to outliers

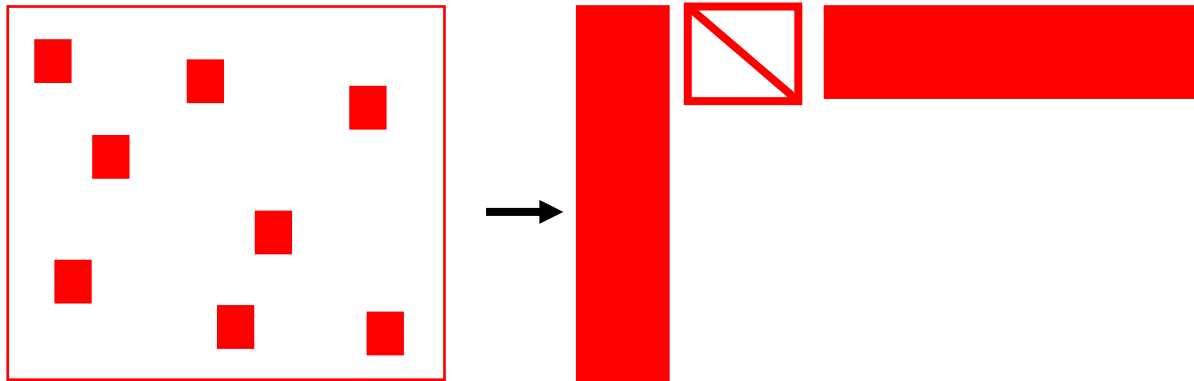


Robust PCA

$$\tilde{X} \approx U_k D_k V_k^T + S$$

S sparse

PCA: Non suitable for sparse data, expensive to compute



Sparse PCA

$$\tilde{X} \approx U_k D_k V_k^T$$

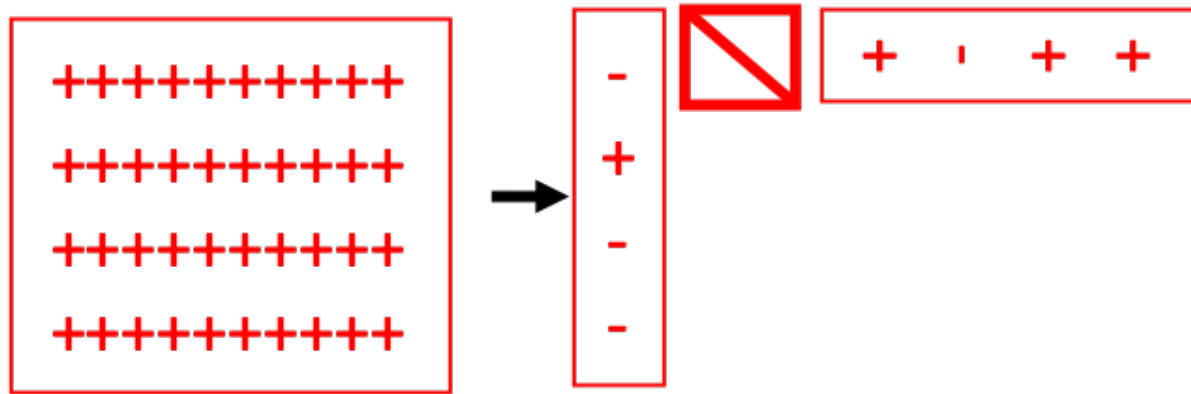
V_k sparse

CUR

$$\tilde{X} \approx CUR$$

C columns
 R rows

PCA: Non suitable for sparse data, expensive to compute



Non-Negative Factorization

$$\tilde{X} \approx WH$$

$$w_{i,k}, h_{k,j} \geq 0$$

PCA: Non suitable for sparse data, expensive to compute

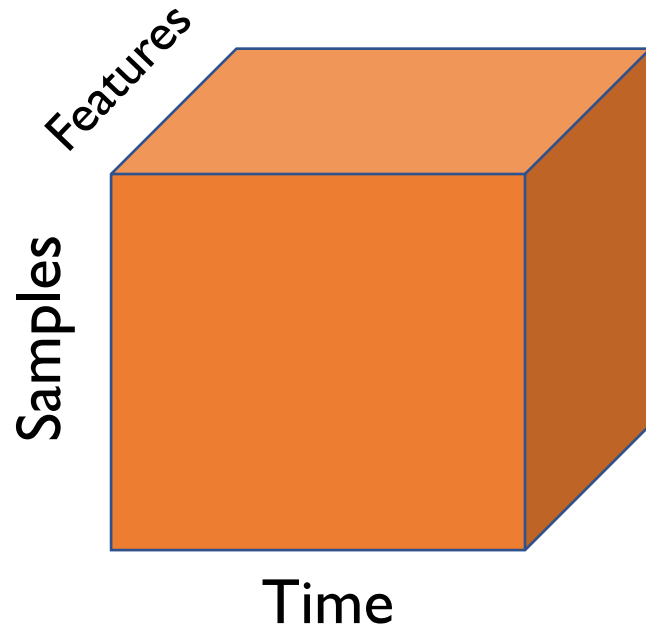
genfaces - PCA using randomized SVD - Train time 0.0



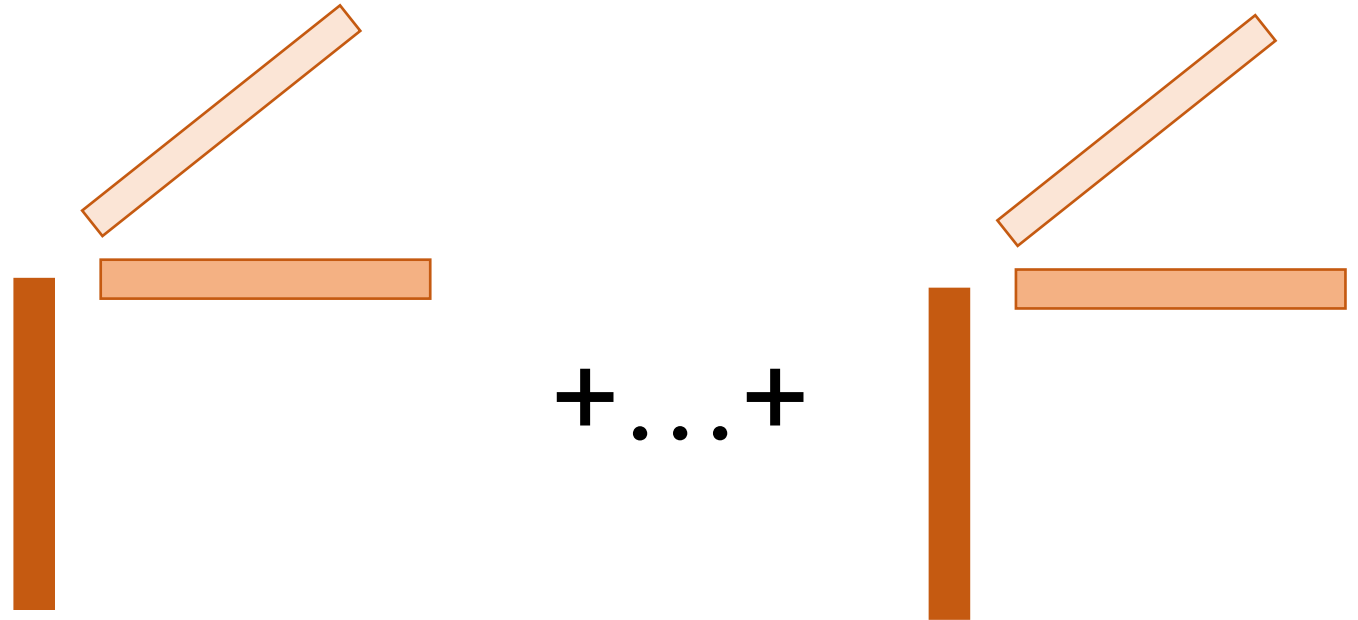
Non-negative components - NMF - Train time 0.2s



PCA: Only 2D arrays



Tensor Decomposition
CP decomposition



PCA: Only 2D arrays

Neuron

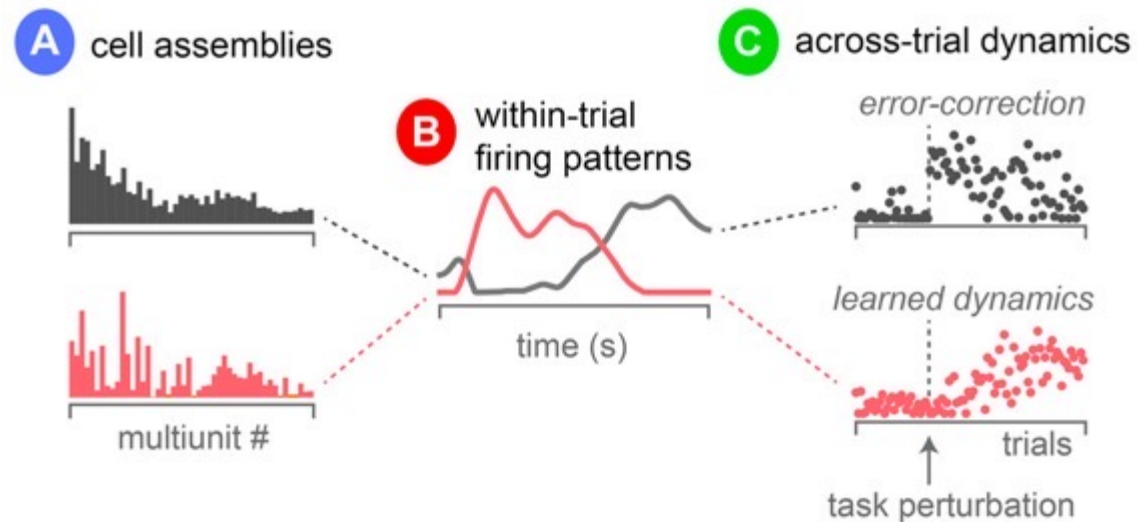
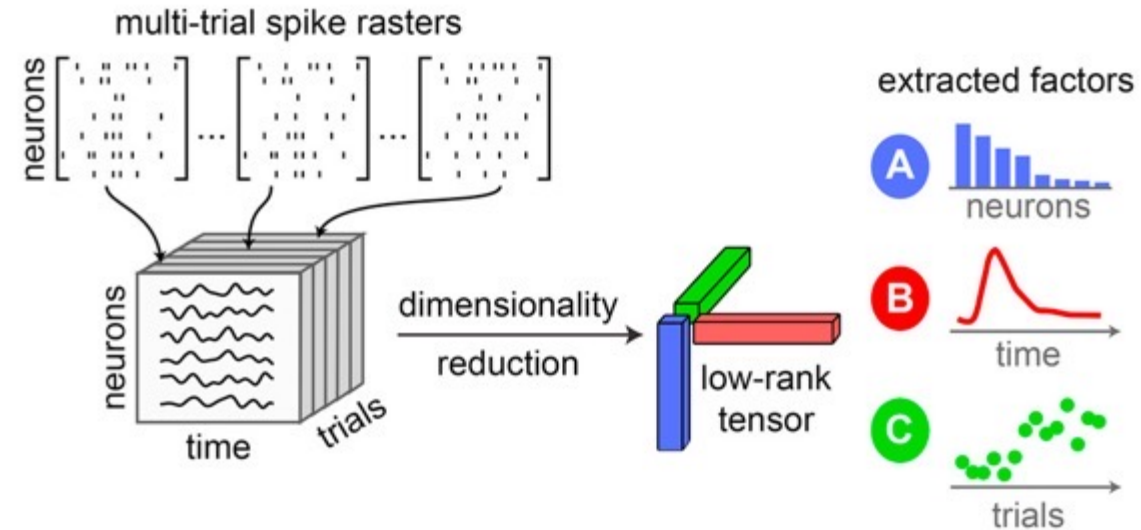
Volume 98, Issue 6, 27 June 2018, Pages 1099-1115.e8



NeuroResource

Unsupervised Discovery of Demixed, Low-Dimensional Neural Dynamics across Multiple Timescales through Tensor Component Analysis

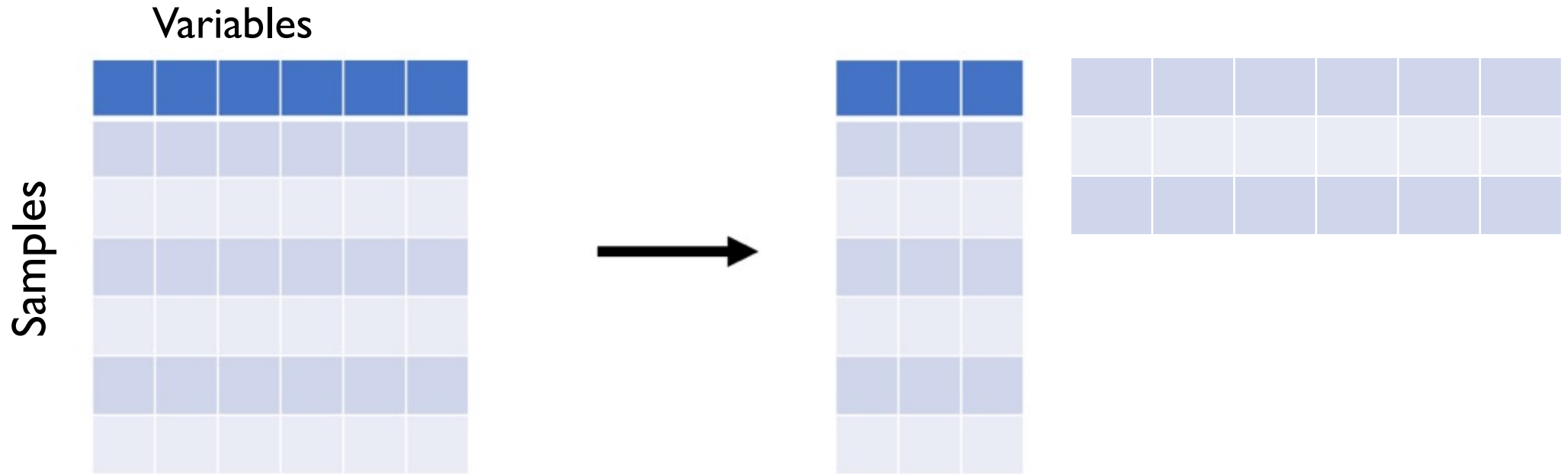
Alex H. Williams^{1, 13}, Tony Hyun Kim², Forea Wang¹, Saurabh Vyas^{2, 3}, Stephen I. Ryu^{2, 11}, Krishna V. Shenoy^{2, 3, 6, 7, 8, 9}, Mark Schnitzer^{4, 5, 7, 9, 10}, Tamara G. Kolda¹², Surya Ganguli^{4, 6, 7, 8}



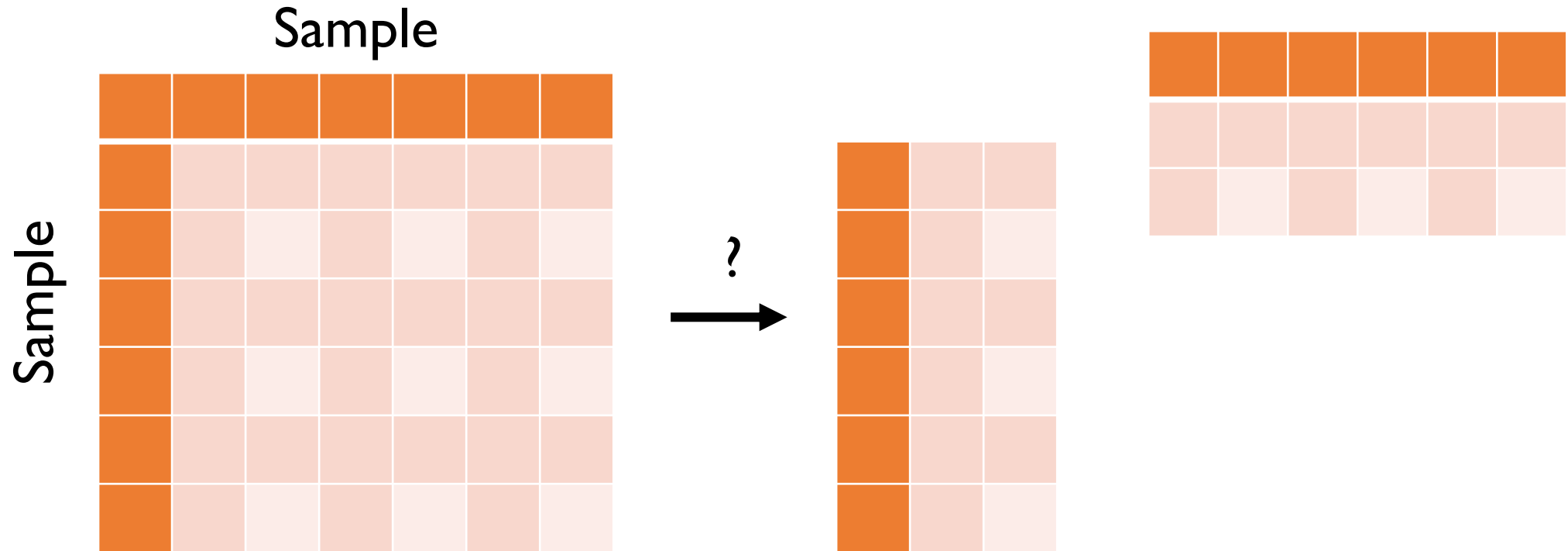
<https://doi.org/10.1016/j.neuron.2018.05.015>

PCA + Other matrix factorizations

$$\tilde{X} \approx UA$$



What about similarity matrix?



Correlation as similarity

Pattern

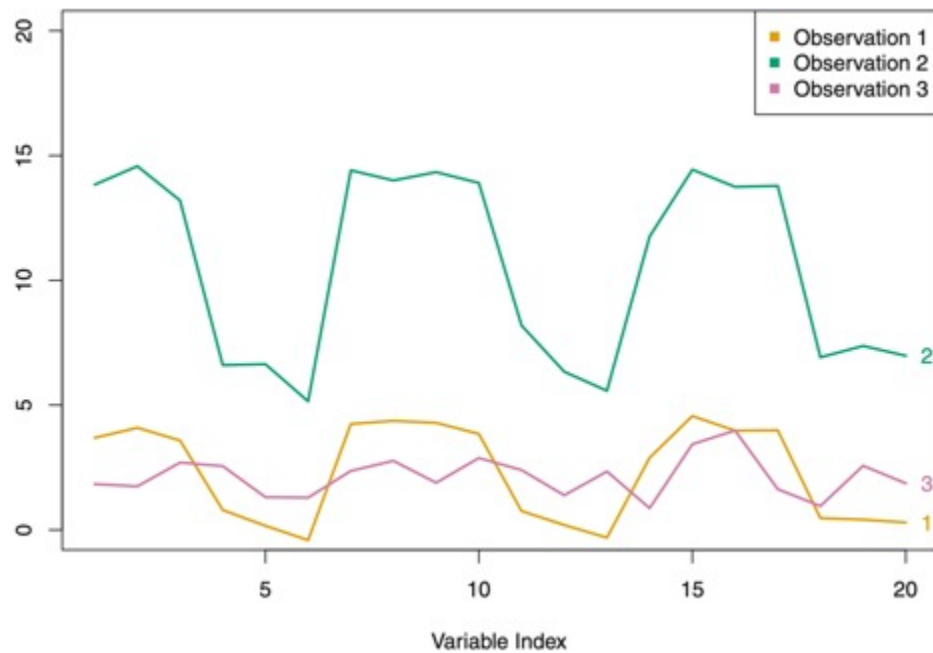


Figure 10.13 ISL (2013)

Correlation

$$d(x^{(1)}, x^{(2)}) = \sum_{k=1}^P \frac{(x_k^{(1)} - \bar{x}^{(1)})(x_k^{(2)} - \bar{x}^{(2)})}{\sqrt{\sum_{k=1}^P (x_k^{(1)} - \bar{x}^{(1)})^2} \sqrt{\sum_{k=1}^P (x_k^{(2)} - \bar{x}^{(2)})^2}}$$

In matrix notation, \tilde{K} similarity matrix

$$\tilde{K} = \tilde{X}\tilde{X}^T$$

SVD and Similarity matrix

In matrix notation, K similarity matrix

$$\tilde{K} = \tilde{X}\tilde{X}^T$$

In terms of cosine similarity

$$\tilde{K} = \left(I - \frac{1}{n}\mathbf{1}\mathbf{1}^T\right)XX^T\left(I - \frac{1}{n}\mathbf{1}\mathbf{1}^T\right)$$

Using SVD of \tilde{X} (Principal Components)

$$\tilde{K} = UD^2U^T = ZZ^T$$

What if we prefer another similarity measure?

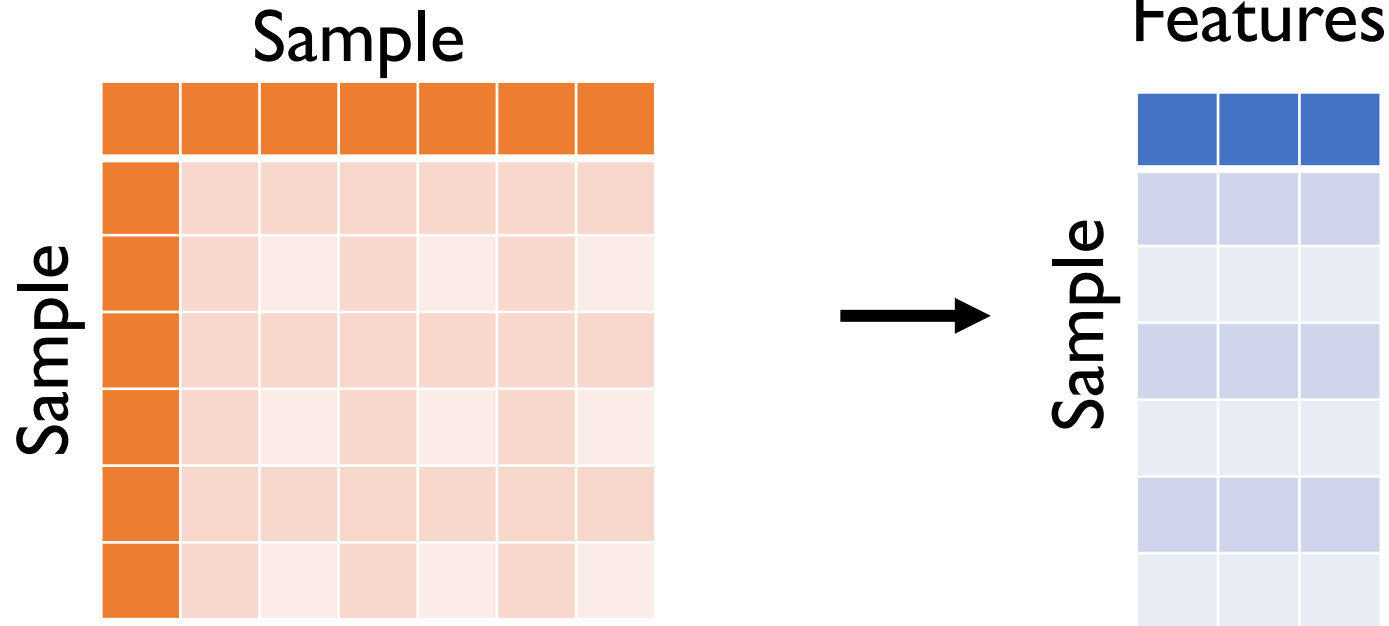
Kernel PCA

In terms of **any similarity**

$$\tilde{K} = \left(I - \frac{1}{n} \mathbf{1}\mathbf{1}^T\right) \mathbf{K} \left(I - \frac{1}{n} \mathbf{1}\mathbf{1}^T\right)$$

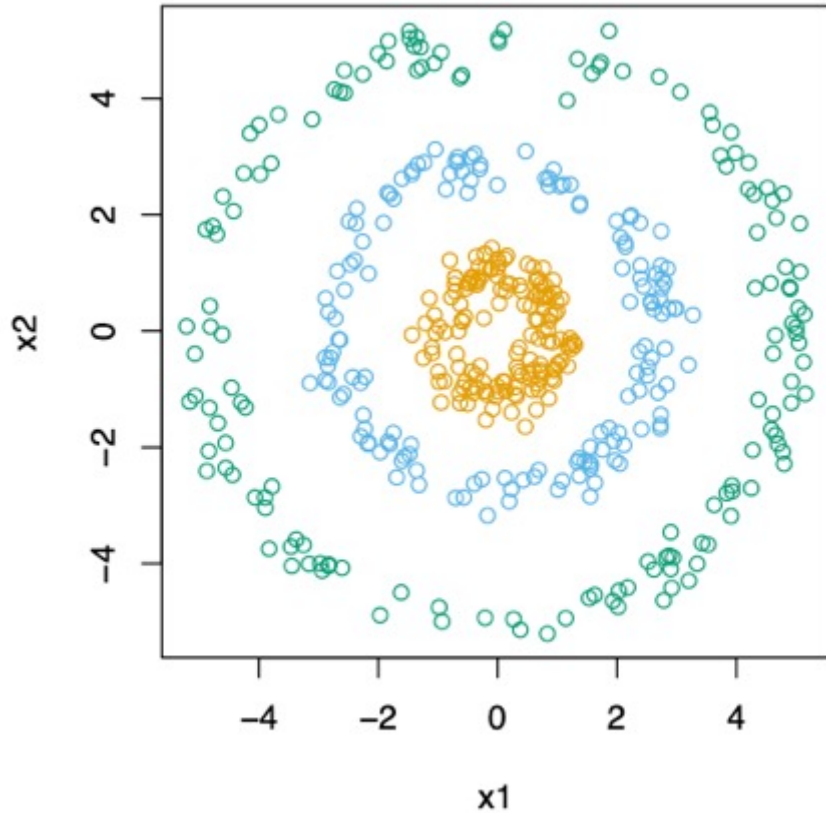
computing **SVD (Principal Components)**

$$\tilde{K} = \mathbf{U}\mathbf{D}^2\mathbf{U}^T = \mathbf{Z}\mathbf{Z}^T$$



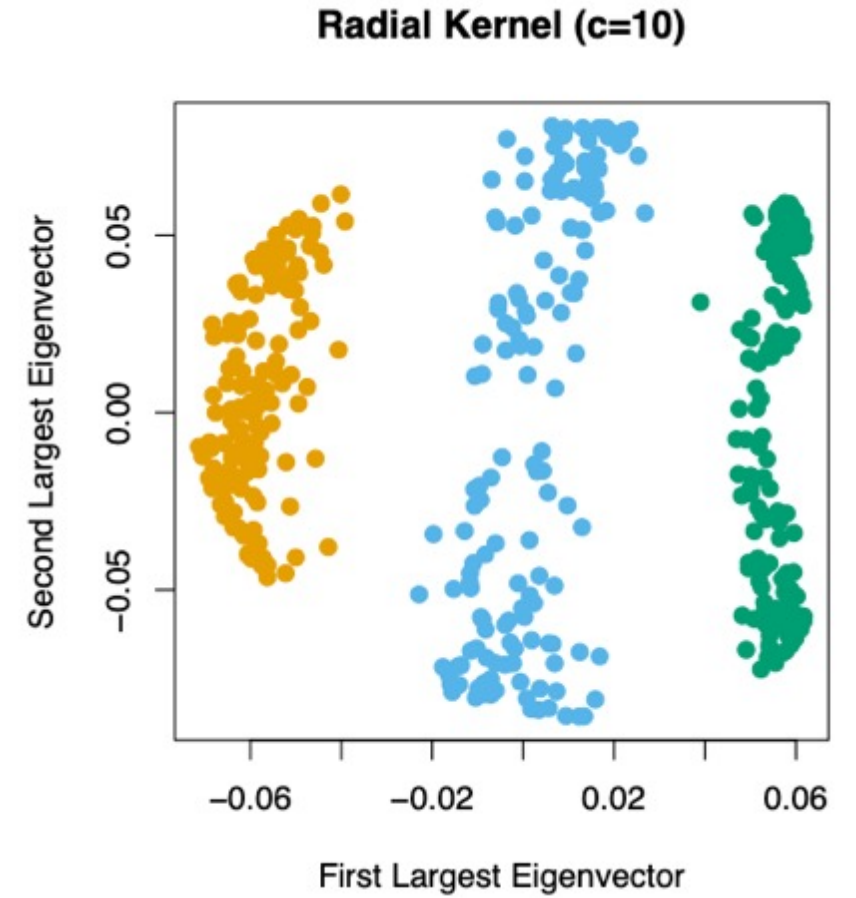
Create features that preserve similarity

Kernel PCA



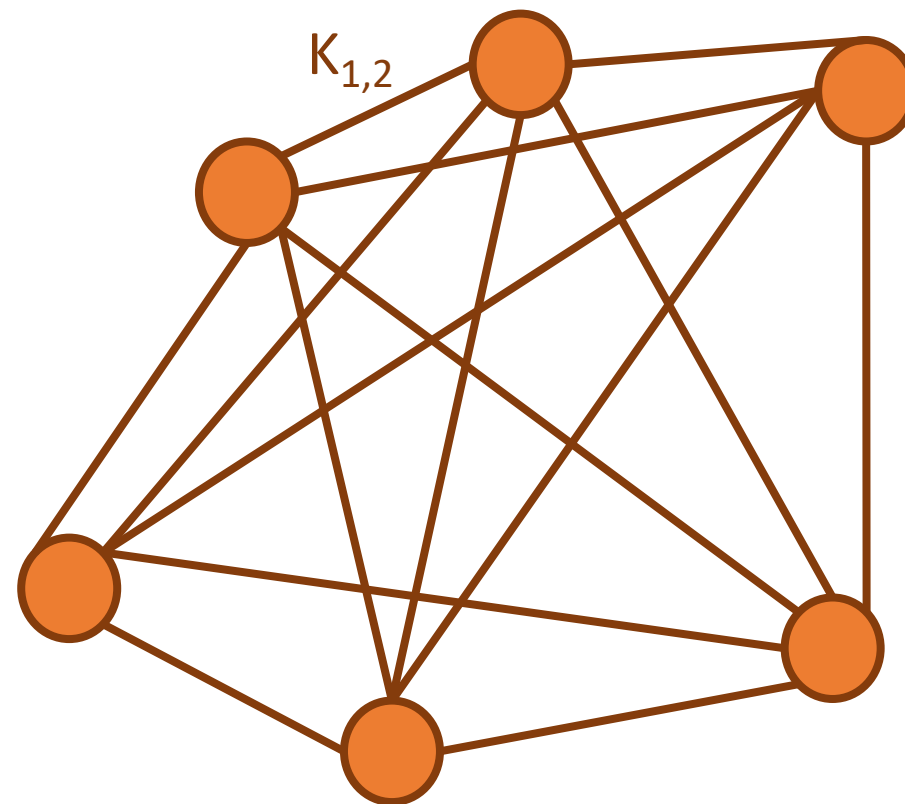
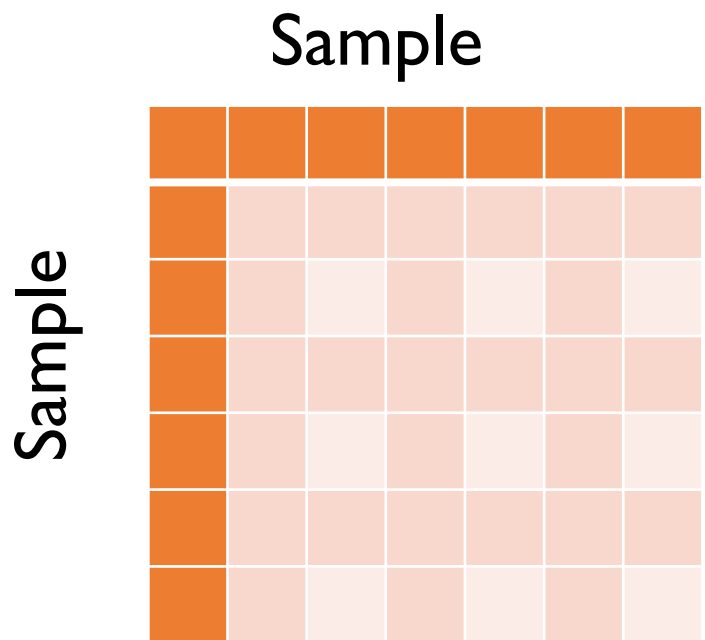
$$K_{i,j} = \exp\left(\frac{\|x^{(i)} - x^{(j)}\|^2}{c}\right)$$

A thick black arrow points from the equation to the right, indicating the transformation of the data into the kernel space.



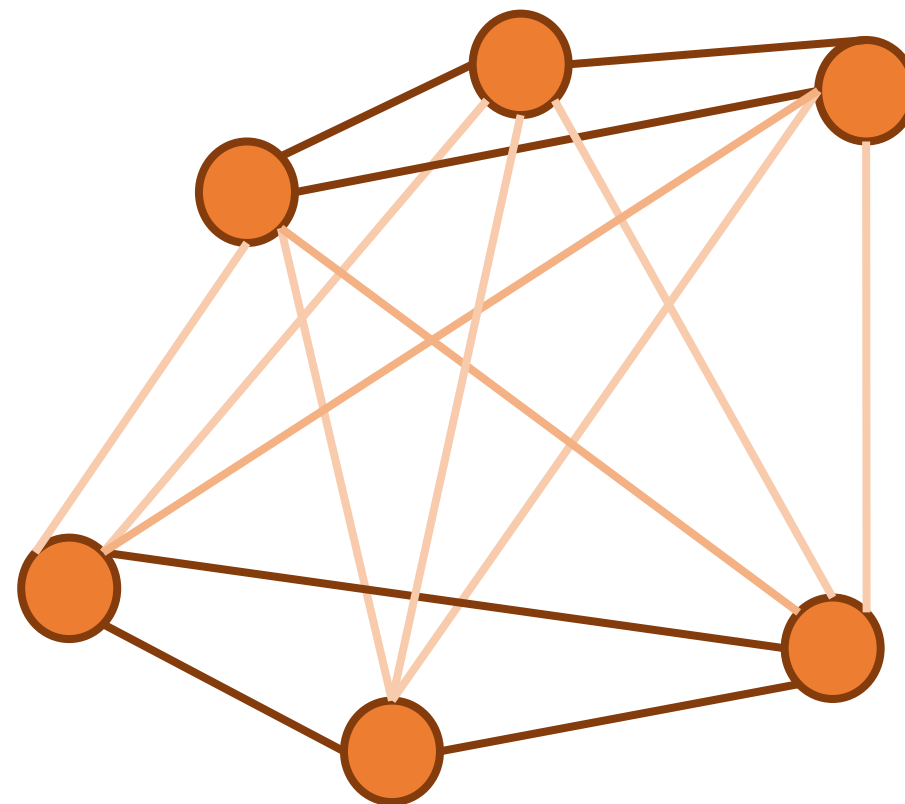
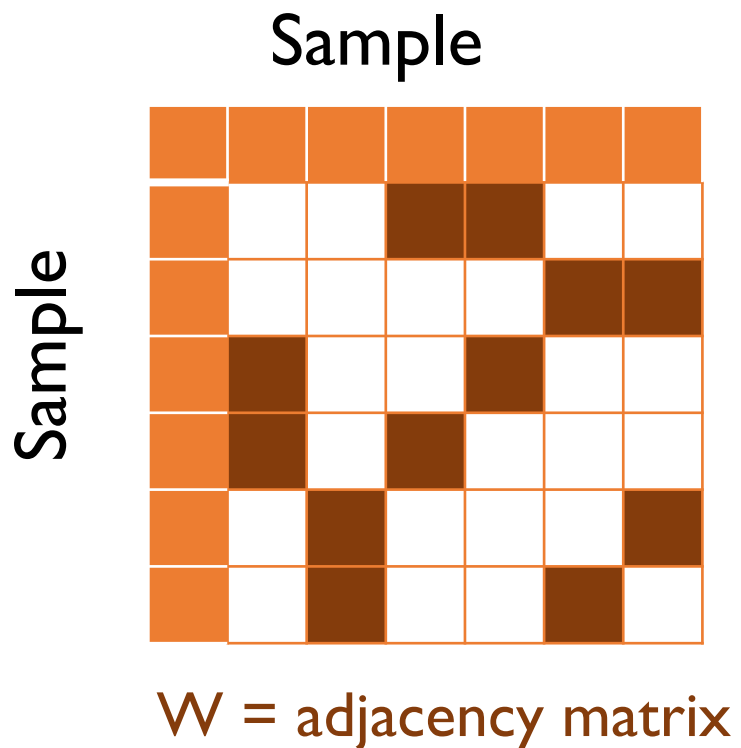
Fragment of Figure 14.30, ESL (2nd edition 2009)

How to interpret similarity matrix?



Edge weight = similarity

K-Nearest Neighbors graph



Edge weight = similarity

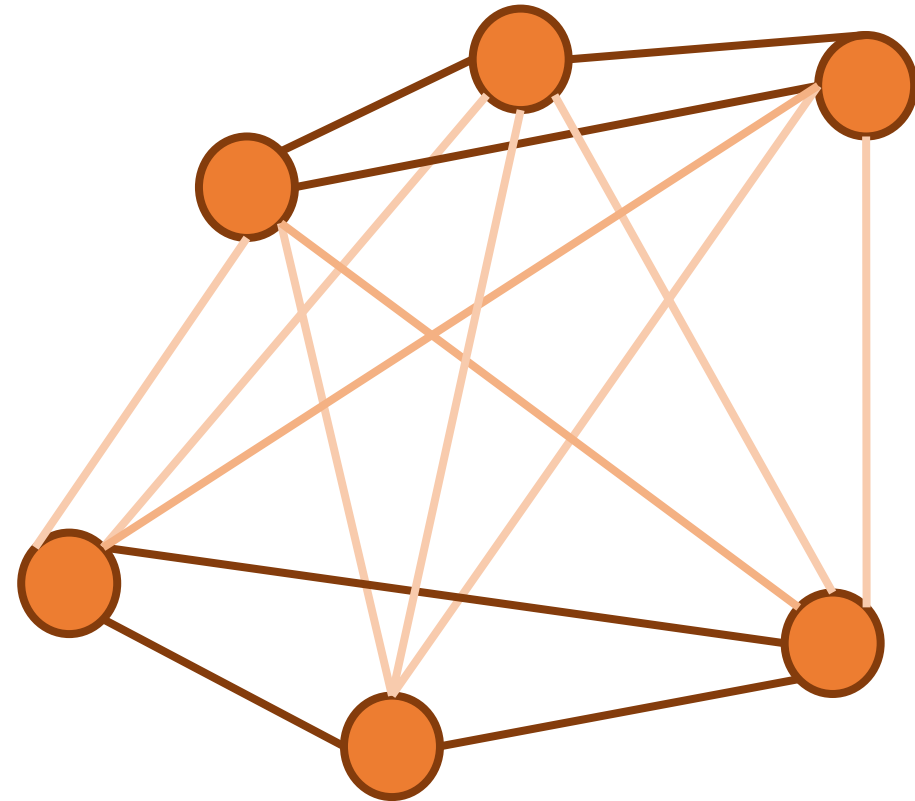
Graph Laplacian

W = Adjacency Matrix

Graph Laplacian

$$L = G - W$$

where $g_{ii} = \sum_{j \neq i} w_{ij}$



Edge weight = similarity

Spectral Clustering

Eigenvectors of L
corresponding to
smallest eigenvalue



New features



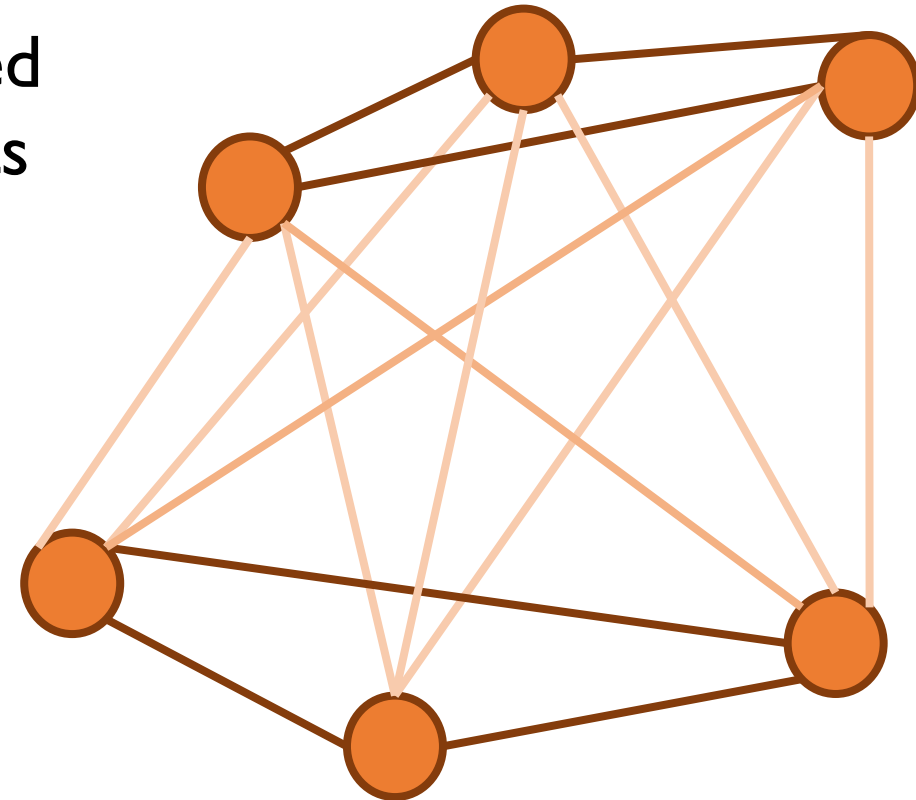
Apply clustering
(e.g. K-means)

\approx

Connected
components



Sample



Edge weight = similarity

Spectral Clustering

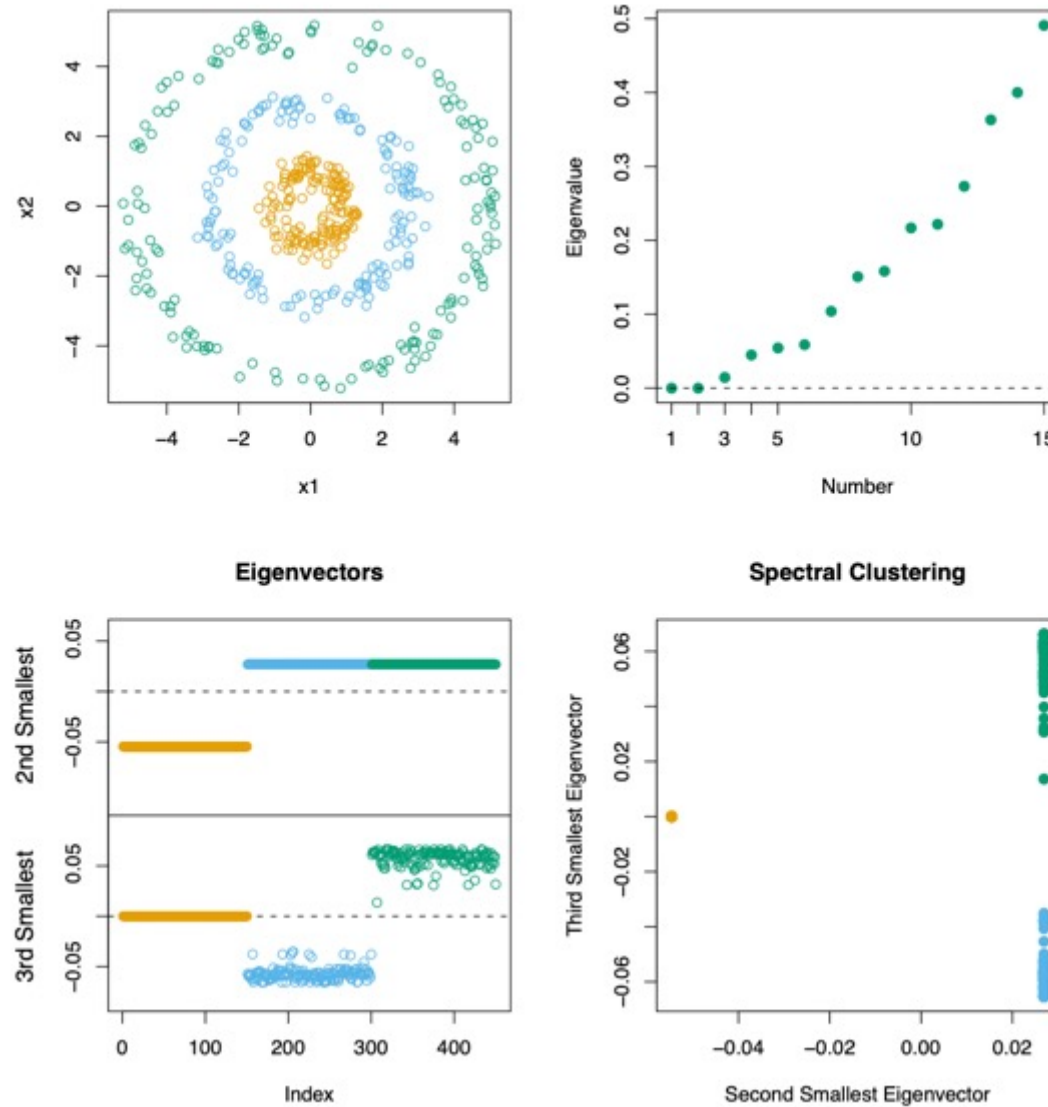
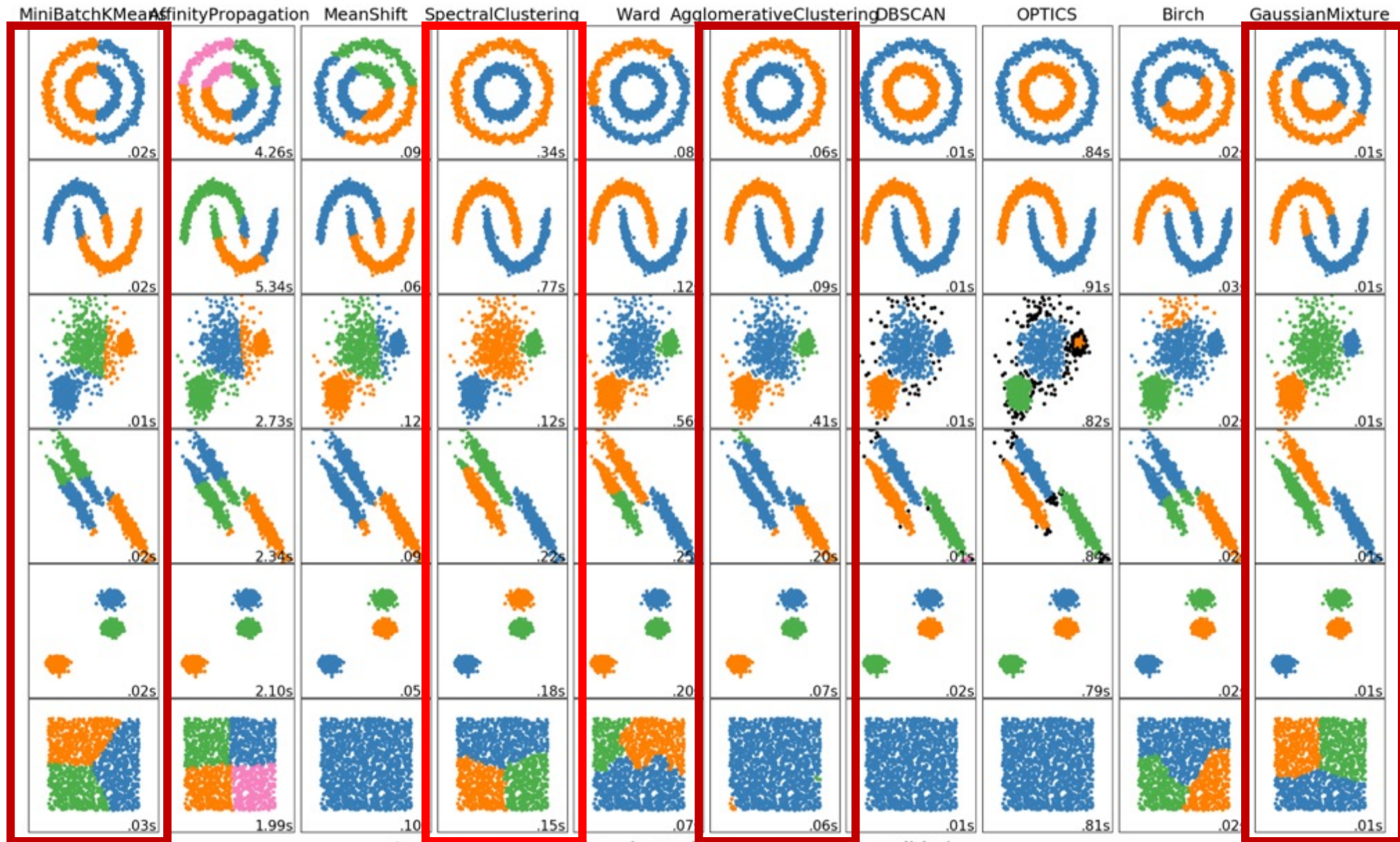


Figure 14.29, ESL (2nd edition 2009)



A comparison of the clustering algorithms in scikit-learn

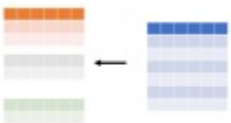
Today recap

Unsupervised Learning

Patterns + Properties in Data

Clustering

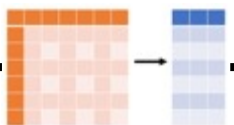
Subgroups of samples



Spectral
clustering



From similarities



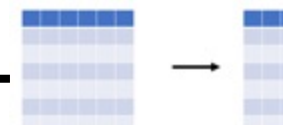
Kernel
PCA



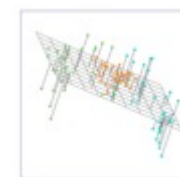
Dimensionality Reduction

Reduce # variables

From features



PCA



Weighted PCA
Robust PCA
ICA
Sparse PCA
CUR
NNMF
CP Decomp.
...

SVD



Truncated
SVD



Further reading

- What about categorical variables?
- What about DR for visualization?

Ten quick tips for effective dimensionality reduction

Lan Huong Nguyen, Susan Holmes.(2019)

<https://doi.org/10.1371/journal.pcbi.1006907>

Logistics

- Join slack channel cme250202.slack.com
- Office hours tomorrow Noon to 1pm (link on Canvas)
- Part I project deadline: April 26.