

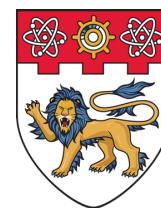
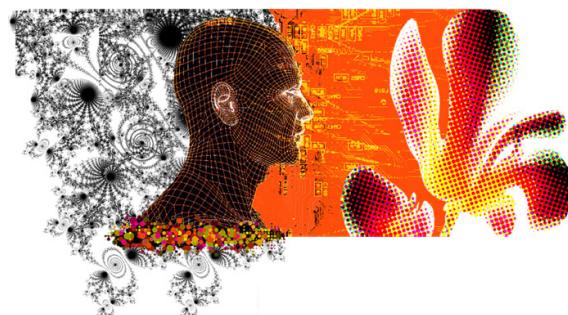
CE9010: Introduction to Data Science

Lecture 8: Unsupervised Learning

Semester 2 2017/18

Xavier Bresson

School of Computer Science and Engineering
Data Science and AI Research Centre
Nanyang Technological University (NTU), Singapore



NANYANG
TECHNOLOGICAL
UNIVERSITY
SINGAPORE

Outline

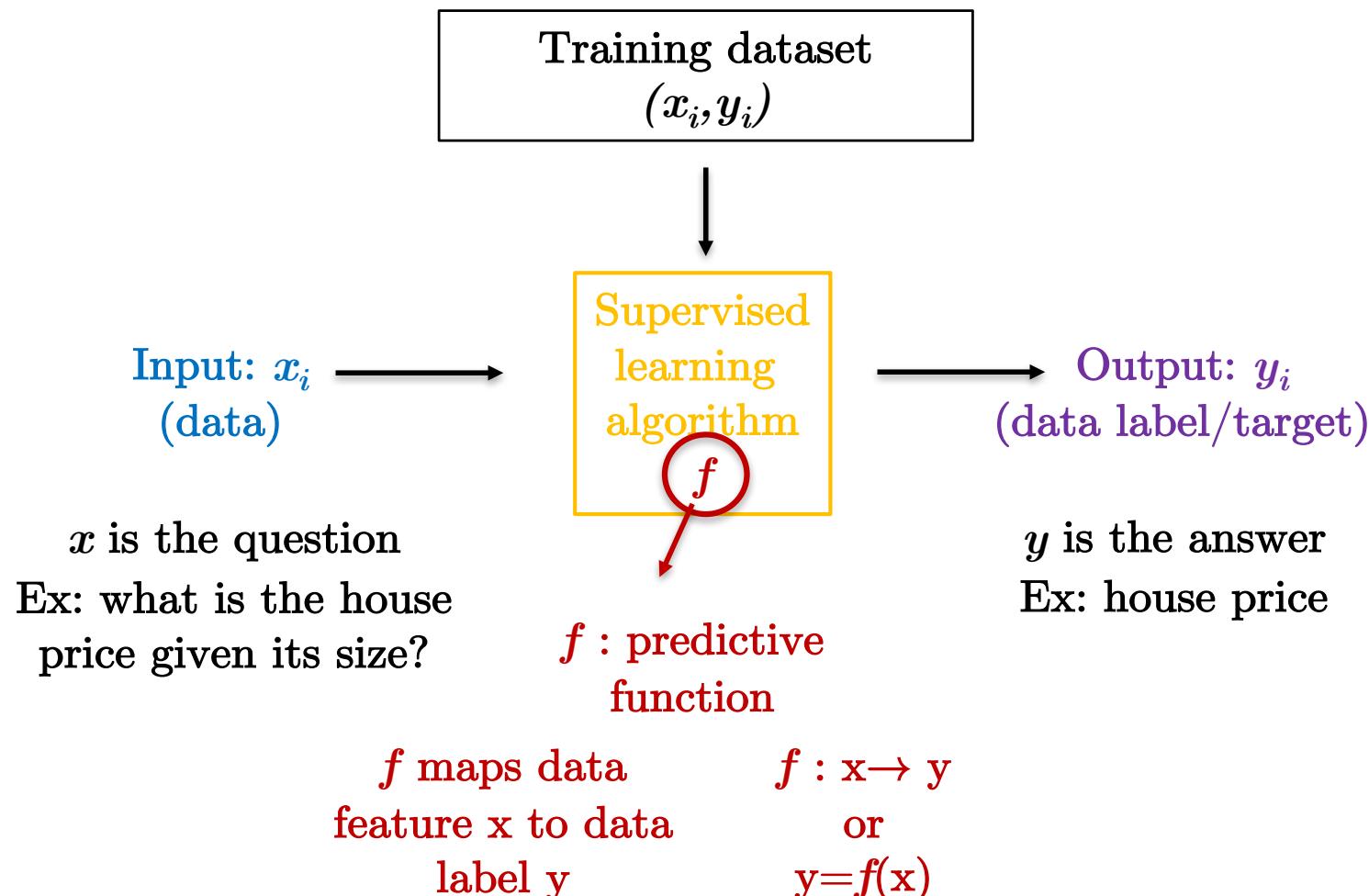
- Supervised vs unsupervised learning
- Unsupervised clustering with k-means
 - Clustering
 - k-means algorithm
 - Loss
 - k-means properties
- Unsupervised representation with PCA
 - Introduction and motivations
 - Principal directions and components
 - Formalization
 - Algorithm
 - PCA properties
- Conclusion

Outline

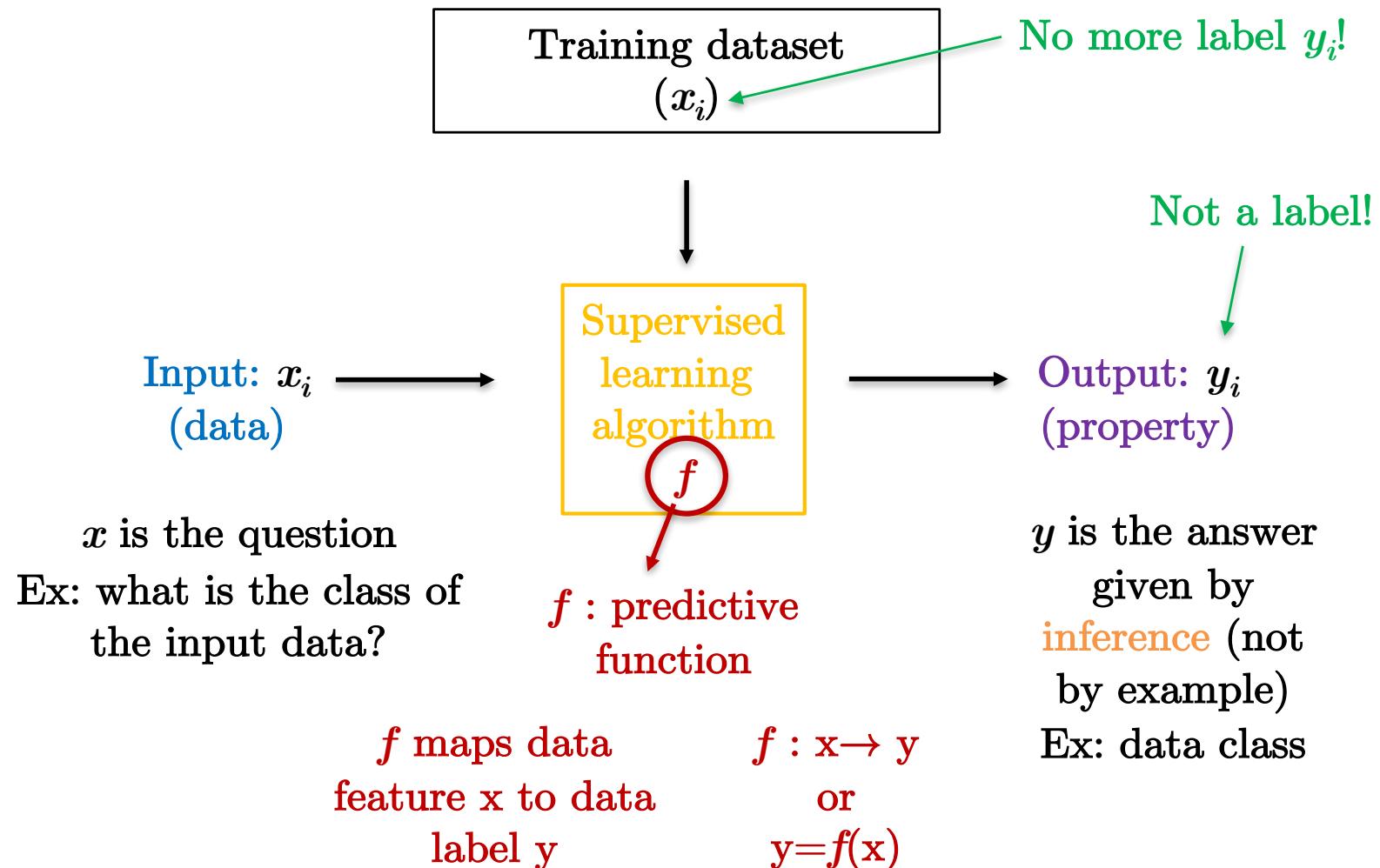
- **Supervised vs unsupervised learning**
- Unsupervised clustering with k-means
 - Clustering
 - k-means algorithm
 - Loss
 - k-means properties
- Unsupervised representation with PCA
 - Introduction and motivations
 - Principal directions and components
 - Formalization
 - Algorithm
 - PCA properties
- Conclusion

Supervised learning

- **Supervised learning:** Design a predictive function given a training set composed of data points (x_i, y_i) with x_i the data features and y_i the data label.

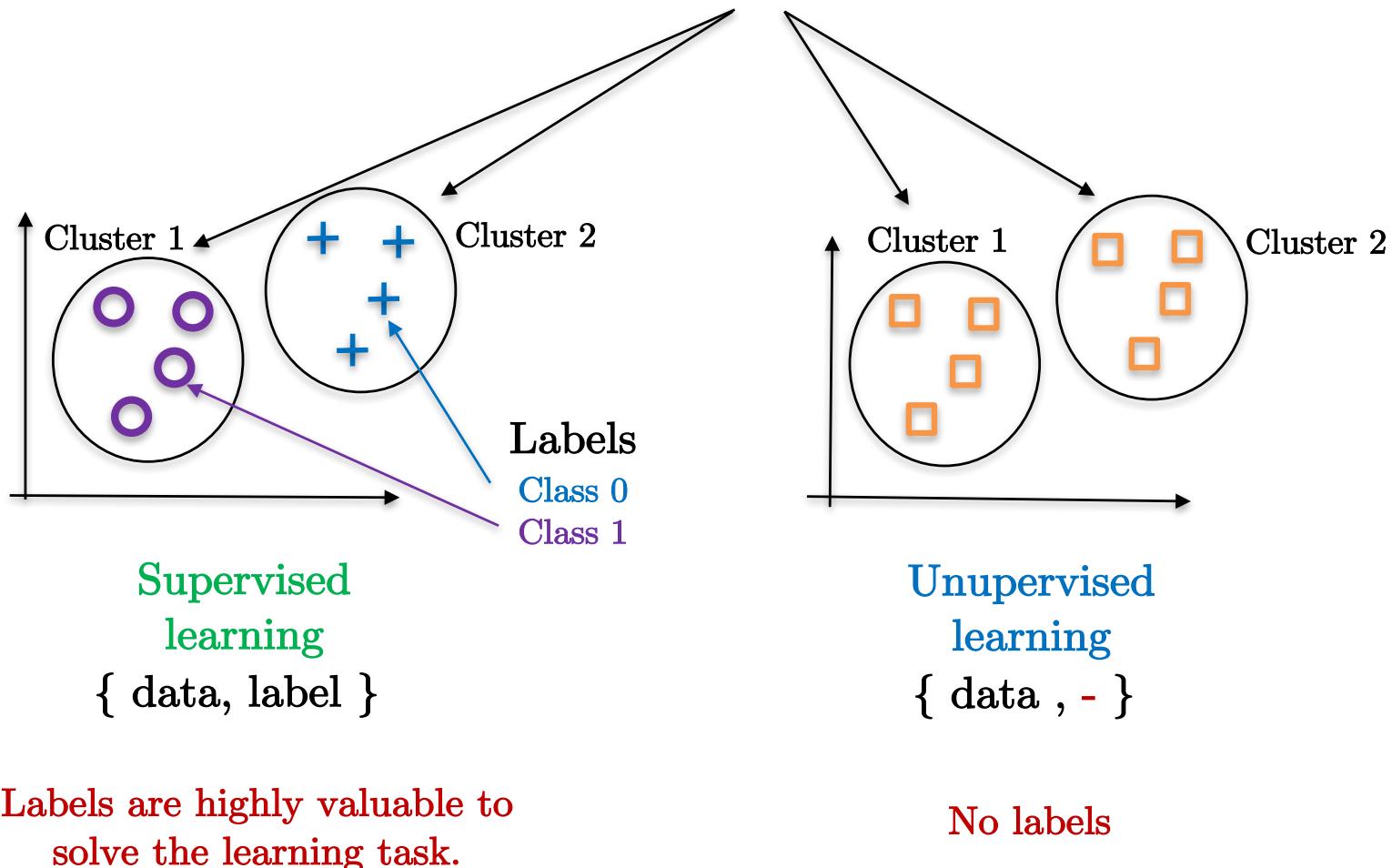


Unsupervised learning



Supervised vs unsupervised

- Difference between supervised and unsupervised learning:
 - Consider the task of finding clusters of data:



Unsupervised learning

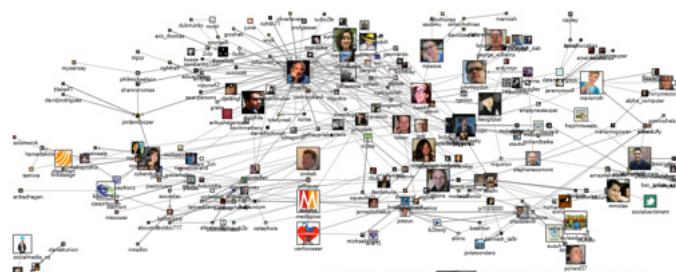
- **Unsupervised learning:** Design a predictive function given a training set which has no label. The training set is simply composed of the data points (x_i) with the data features (no y_i).
- **Main idea:** Find structures in data that can solve data analysis tasks.
- **Common tasks:**
 - Unsupervised clustering: k-means (structures are the means)
 - Unsupervised representation: PCA (structures are the variances)

Outline

- Supervised vs unsupervised learning
- **Unsupervised clustering with k-means**
 - Clustering
 - k-means algorithm
 - Loss
 - k-means properties
- Unsupervised representation with PCA
 - Introduction and motivations
 - Principal directions and components
 - Formalization
 - Algorithm
 - PCA properties
- Conclusion

Unsupervised clustering

- Definition: Find groups of data that are similar.
- Applications:



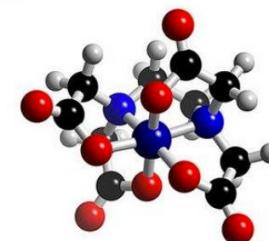
Social network analysis
(ads for groups of similar people)



Market analysis
(segment customers or products in homogeneous groups)



Image segmentation
(group pixels of same objects)



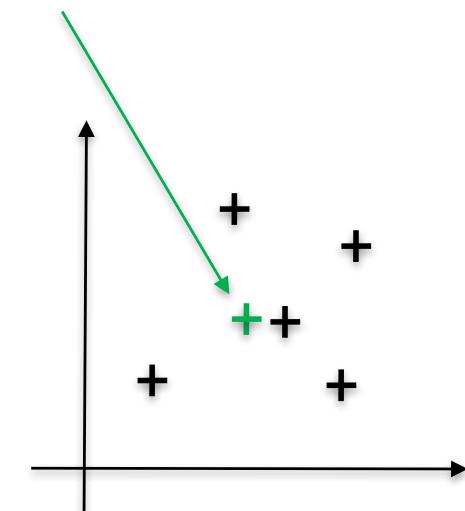
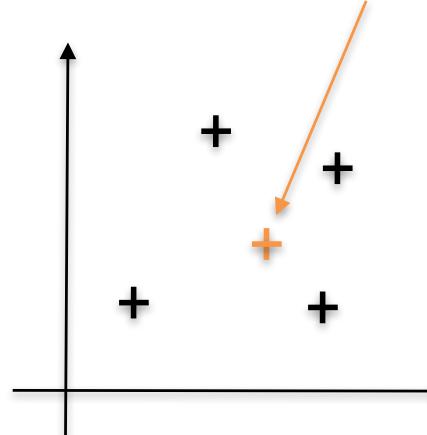
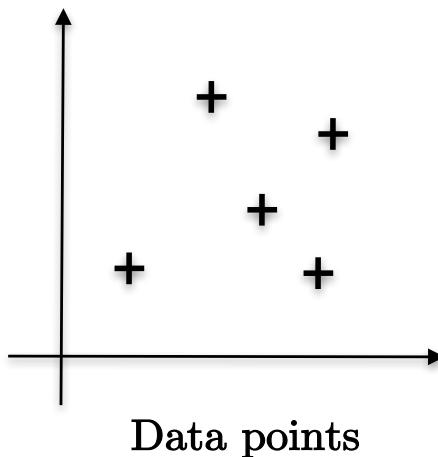
Molecular analysis
(find molecules related to drugs)

Outline

- Supervised vs unsupervised learning
- **Unsupervised clustering with k-means**
 - Clustering
 - **k-means algorithm**
 - Loss
 - k-means properties
- Unsupervised representation with PCA
 - Introduction and motivations
 - Principal directions and components
 - Formalization
 - Algorithm
 - PCA properties
- Conclusion

k-means

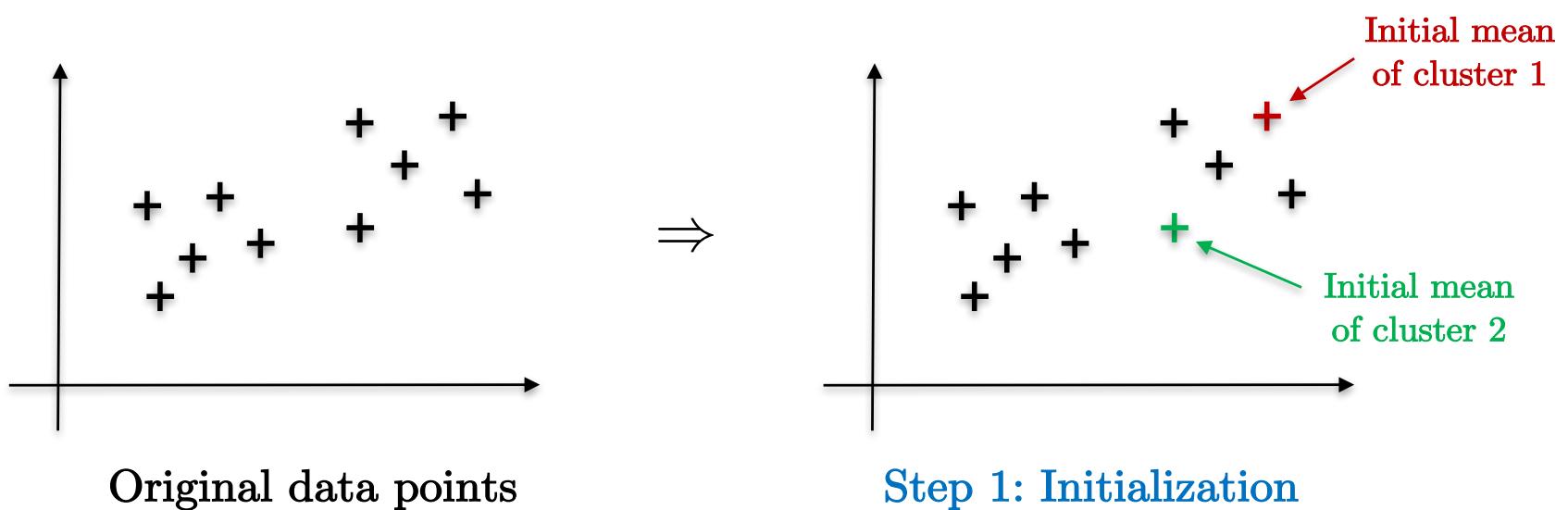
- k-means algorithm is the **most popular** unsupervised clustering algorithm.
- **Definition:** The mean (a.k.a. cluster center, centroid) of a cluster is the center of this cluster.
- **Center definition:**
 - The center is a point not necessarily belonging to the set of data.
 - The center is a point from the set of data.



- In this lecture, we choose the **center** to be a point from the set of data.

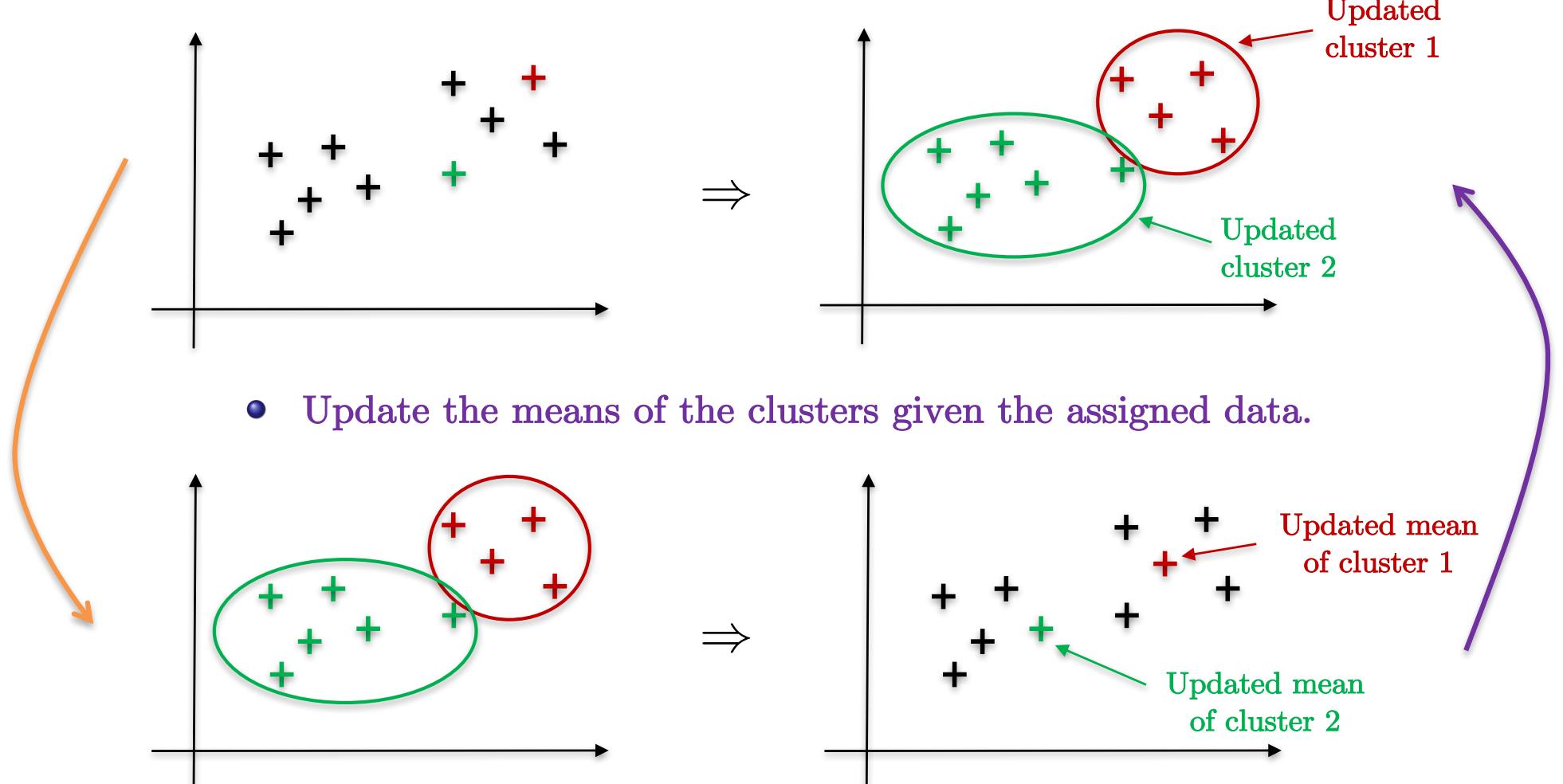
k-means

- k-means algorithm steps:
 - Step 1: Initialization. Select randomly k points called means:



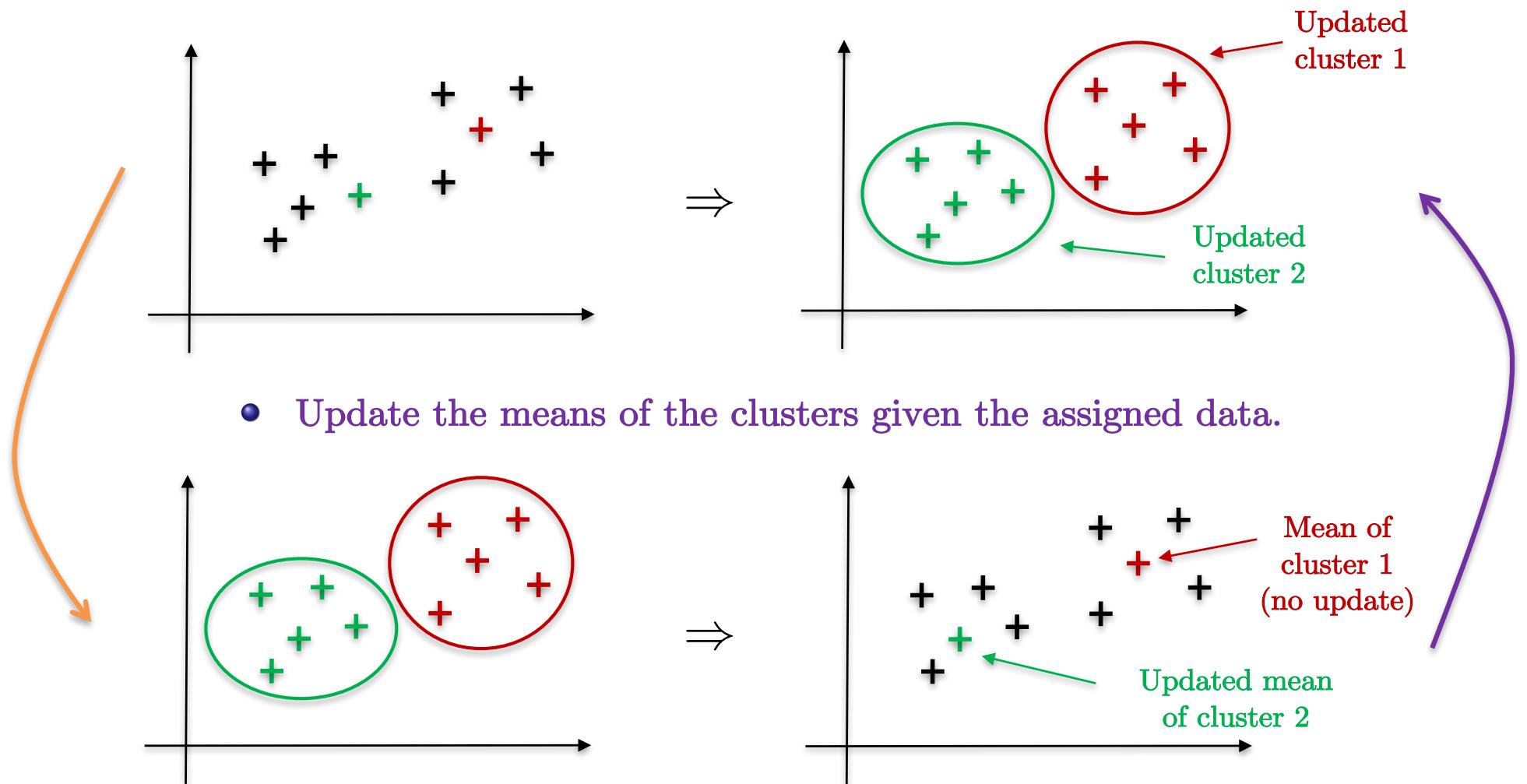
k-means

- Step 2: Loop until convergence.
 - Assign each data to one of the 2 means.



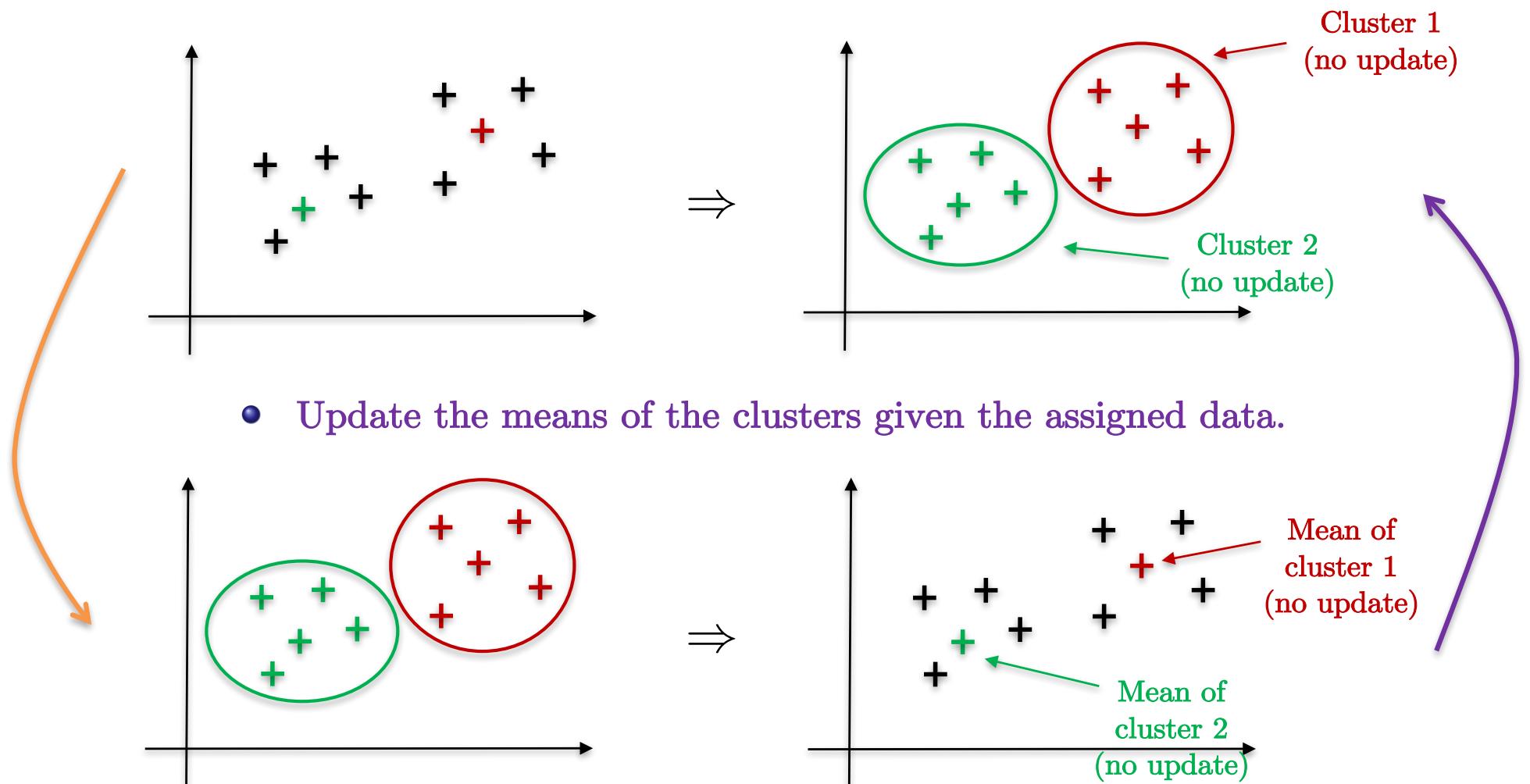
k-means

- Step 2: Loop until convergence.
 - Assign each data to one of the 2 means.



k-means

- Step 2: Loop until convergence.
 - Assign each data to one of the 2 means.



Formalization

- **Inputs:**
 - K = number of clusters
 - Training set $\{x_1, \dots, x_n\}$, x_i in R^d , d =number of features/dim
- **Outputs:**
 - Clusters = $\{C_1, \dots, C_K\}$ and means = $\{\mu_1, \dots, \mu_K\}$.
- **K-means algorithm:**
 - Initialization: Randomly initialize K means: $\{\mu_1, \dots, \mu_K\} \in R^d$.
 - Repeat until convergence:
 1. Clusters assignment $\{C_1, \dots, C_K\}$
 2. Means update $\{\mu_1, \dots, \mu_K\}$.

Formalization

1. Cluster assignment

For all $i = 1$ to n :

a_i gives the index of the cluster the closest to x_i .

$$a_i = \arg \min_k \|x_i - \mu_k\|_2^2$$

Take the k that minimizes the distance between x_i and the mean μ_k .

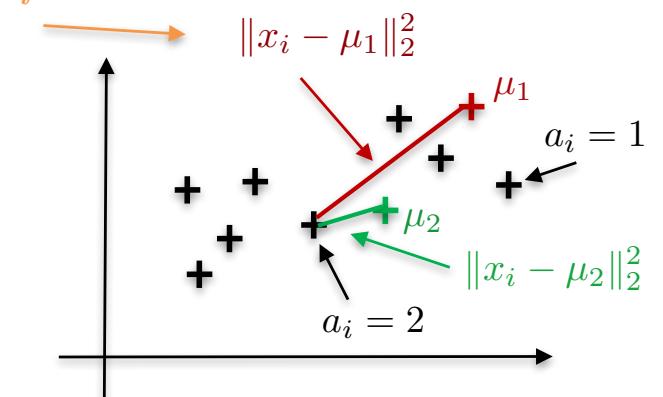
with $\|x\|_2^2 = \sum_{j=1}^d x_j^2$

Euclidean distance

$$a = \begin{bmatrix} 2 \\ 1 \\ 1 \\ \vdots \\ 2 \end{bmatrix}$$

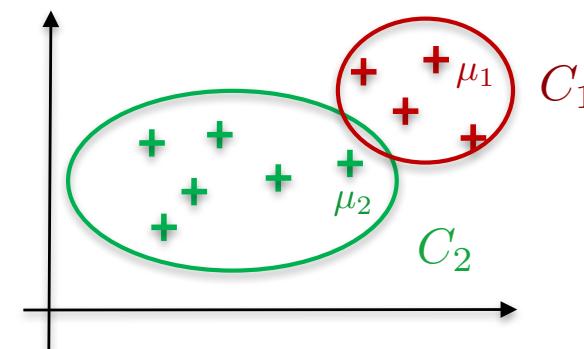
1st data in C_2

3rd data in C_1



For all $k = 1$ to K :

$$C_k = \{x_i : a_i = k\}$$



Formalization

2. Mean update

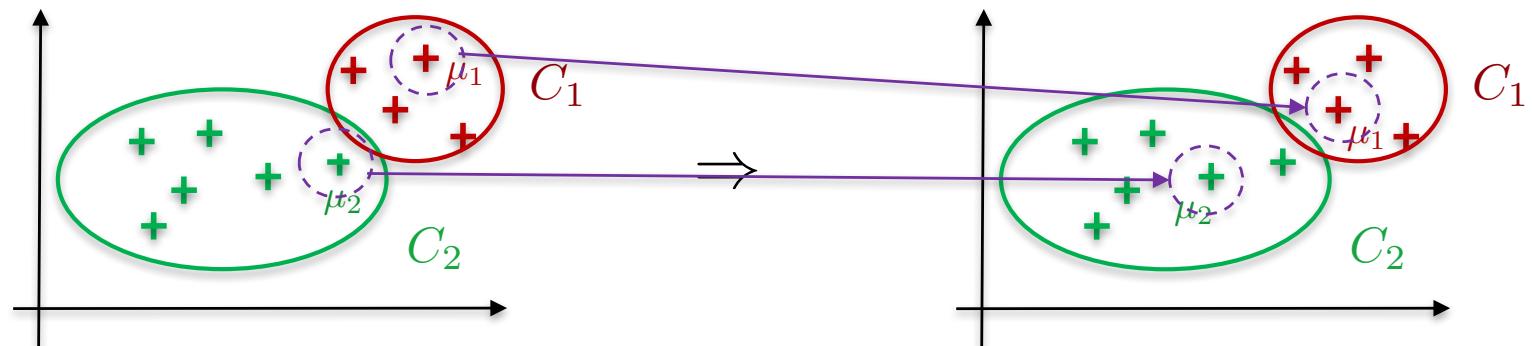
For all $k = 1$ to K :

$$\mu_k = \frac{\sum_{x_i \in C_k} x_i}{|C_k|}$$

μ_k is the mean of all data points in cluster C_k

nb of data in cluster C_k

$d \times 1$



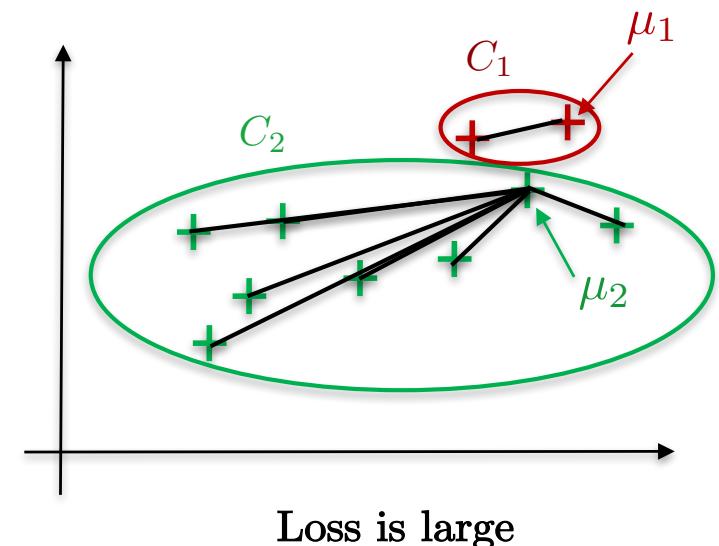
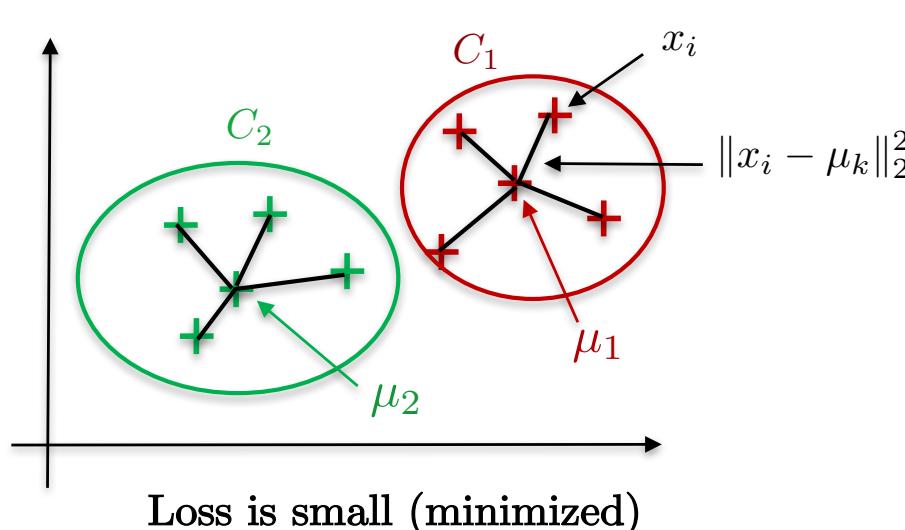
Outline

- Supervised vs unsupervised learning
- **Unsupervised clustering with k-means**
 - Clustering
 - k-means algorithm
 - **Loss**
 - k-means properties
- Unsupervised representation with PCA
 - Introduction and motivations
 - Principal directions and components
 - Formalization
 - Algorithm
 - PCA properties
- Conclusion

k-means loss

- Measure of fitness/accuracy between the k-means model and the data points:

$$\begin{aligned} L(\underbrace{C_1, \dots, C_K}_{\text{Clusters}}, \underbrace{\mu_1, \dots, \mu_K}_{\text{Means}}) &= \frac{1}{n} \sum_{i=1}^n \|x_i - \mu_{a_i}\|_2^2 \\ &= \frac{1}{n} \sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - \mu_k\|_2^2 \end{aligned}$$



Re-visiting k-means algorithm

- Optimize the k-means loss:

$$\min_{\substack{C_1, \dots, C_K \\ \mu_1, \dots, \mu_K}} L(C_1, \dots, C_K, \mu_1, \dots, \mu_K)$$

Find C_k and μ_k that minimize the loss function.

Re-visiting k-means algorithm

- Repeat:

1. Cluster assignment

$$\min_{C_1, \dots, C_K} L(C_1, \dots, C_K) \quad \text{for } \mu_1, \dots, \mu_K \text{ fixed}$$



$$C_k = \{x_i : a_i = k\}$$

$$a_i = \arg \min_k \|x_i - \mu_k\|_2^2$$

2. Mean update

$$\min_{\mu_1, \dots, \mu_K} L(\mu_1, \dots, \mu_K) \quad \text{for } C_1, \dots, C_K \text{ fixed}$$



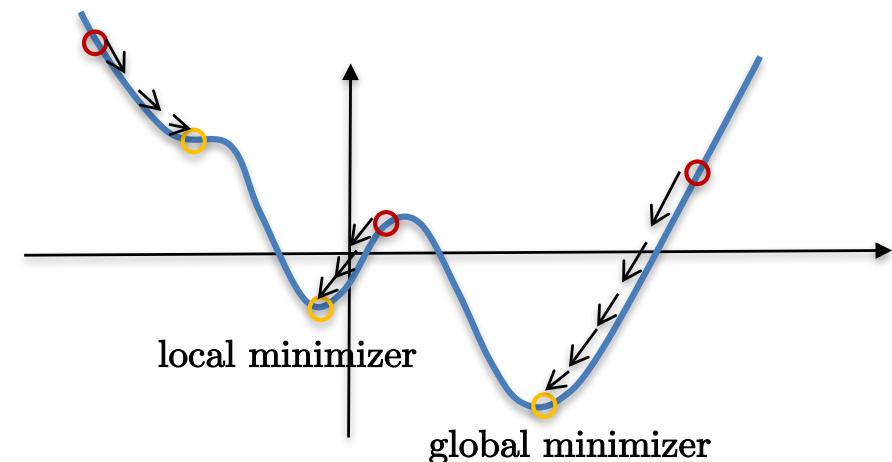
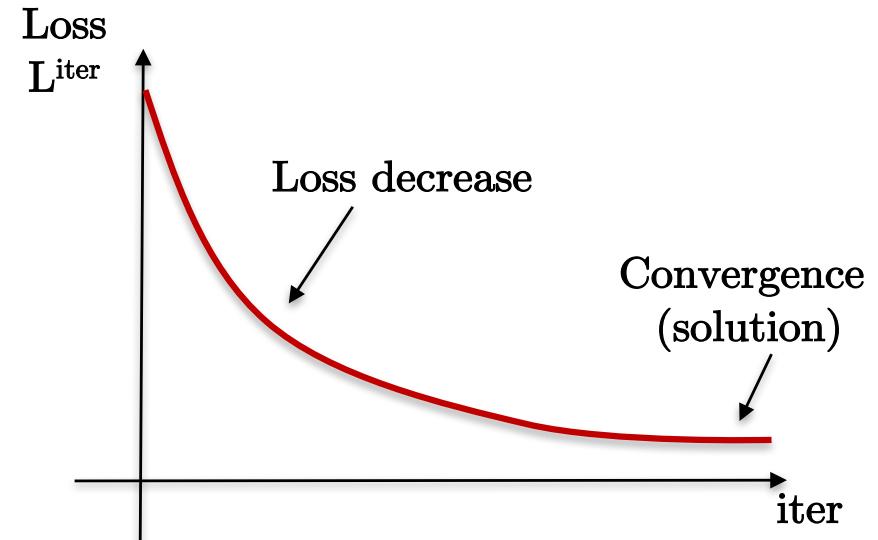
$$\mu_k = \frac{\sum_{x_i \in C_k} x_i}{|C_k|}$$

Outline

- Supervised vs unsupervised learning
- **Unsupervised clustering with k-means**
 - Clustering
 - k-means algorithm
 - Loss
 - **k-means properties**
- Unsupervised representation with PCA
 - Introduction and motivations
 - Principal directions and components
 - Formalization
 - Algorithm
 - PCA properties
- Conclusion

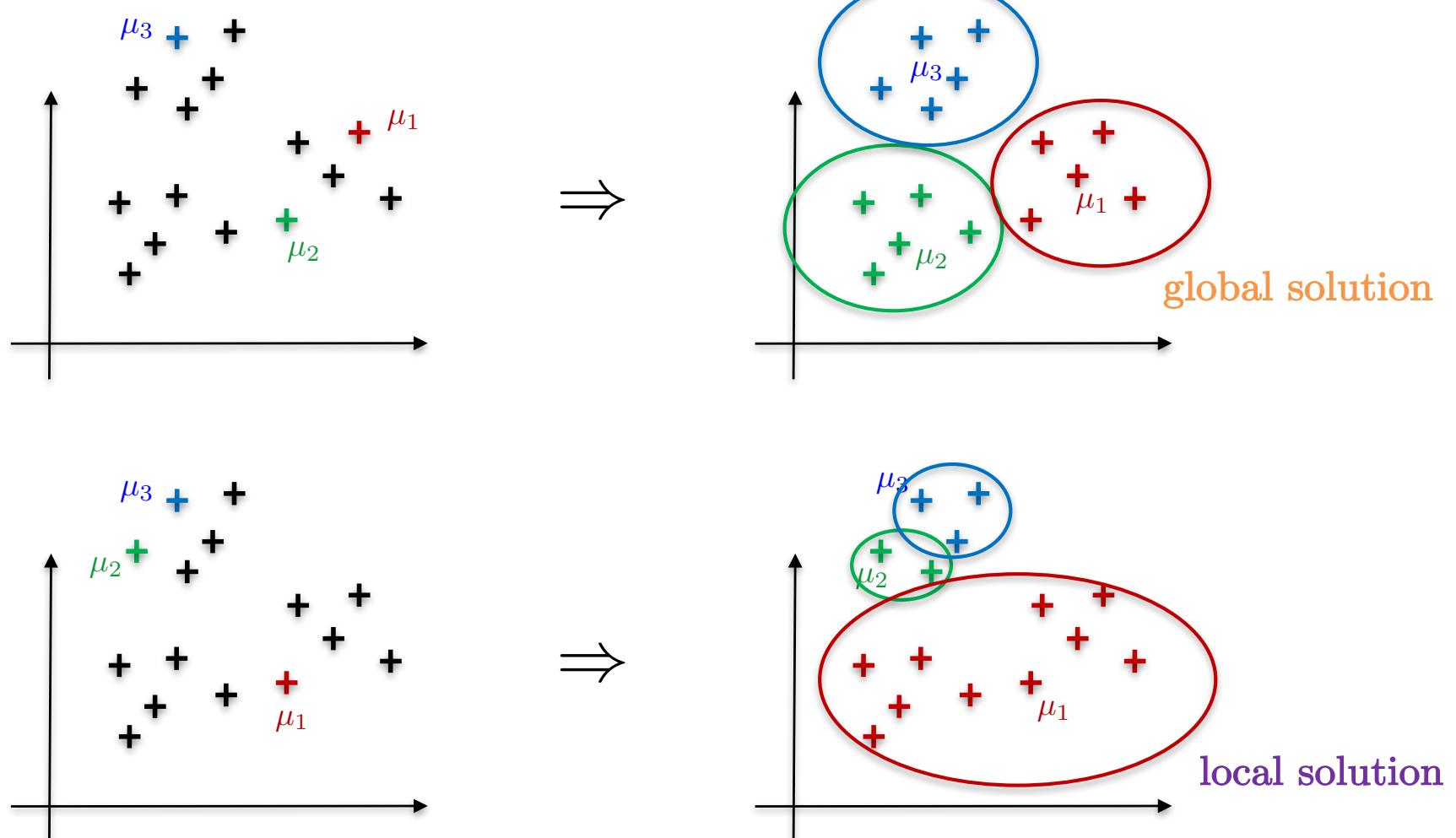
Monotonicity and non-convexity

- **Monotonicity:** k-means is guaranteed to decrease the loss at each iteration.
- **Non-convex:** k-means loss is non-convex
 - ⇒ No guarantee to find the global minimizer of the loss.
 - ⇒ Initial condition is important.



Initialization

- k-means has no guarantee to find a global solution:

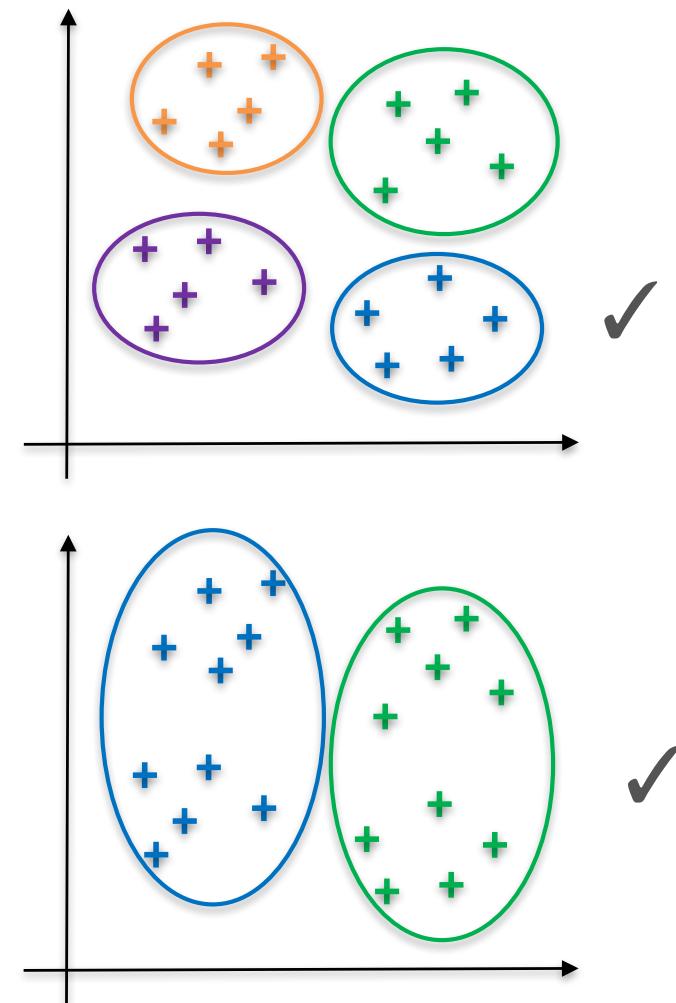
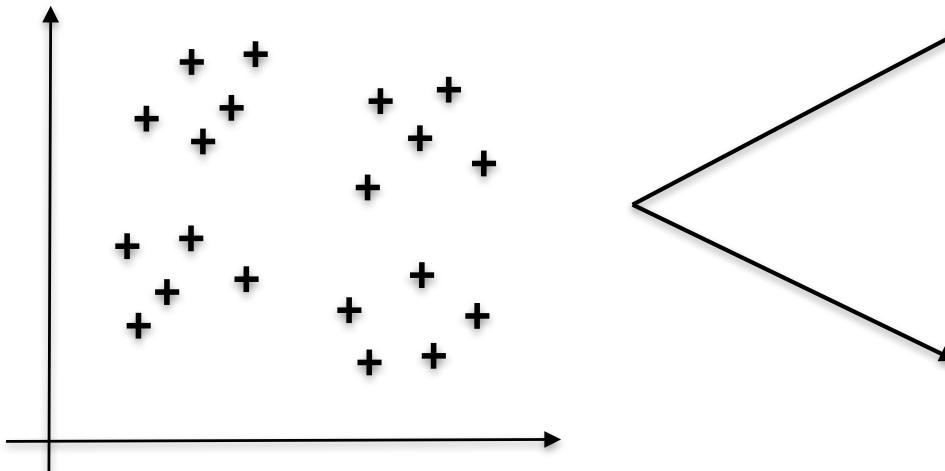


Initialization

- Reduce the dependence w.r.t. initialization:
 - Run the k-means algorithm 100 times with random initialization.
 - Save clusters and the loss value for each run.
 - Pick the solution with the smallest loss value.

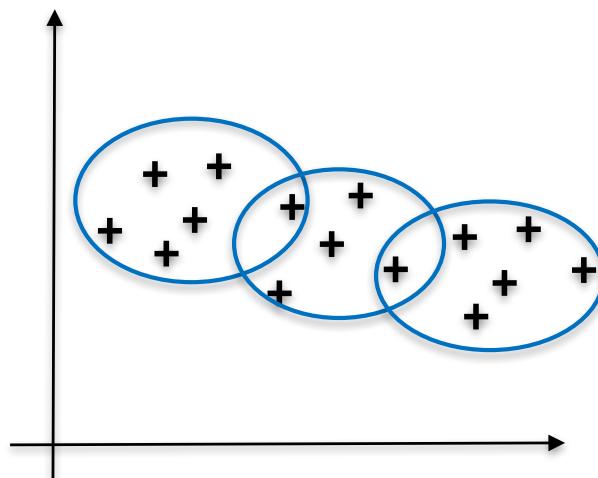
How to choose the number k of clusters?

- No obvious answer. It depends on the data and the next task that will make use of the clusters (sometimes, post-processing selection after k-means ran).
- Example: Ambiguous “optimal” choice of k?

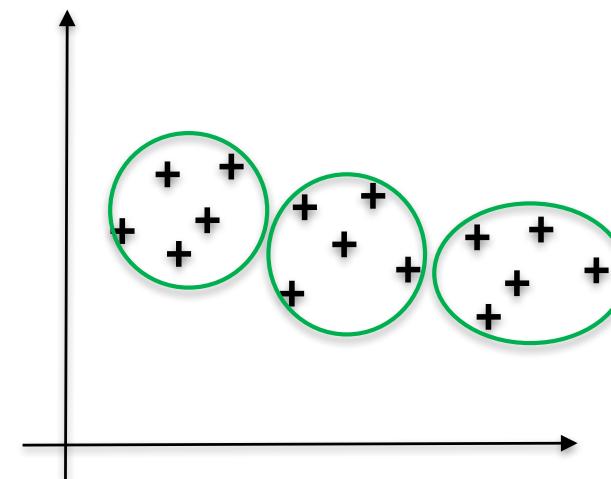


Separated clusters

- K-means algorithm **cannot be applied to non-separated/overlapping clusters**. It is designed to find **separated clusters**.



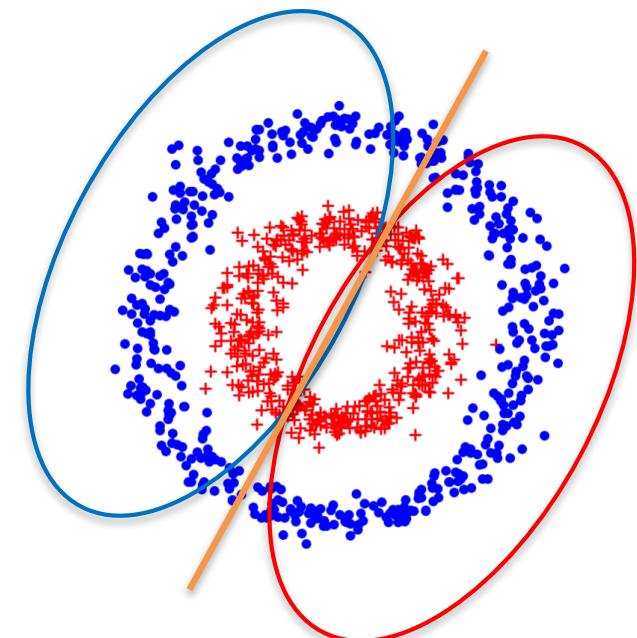
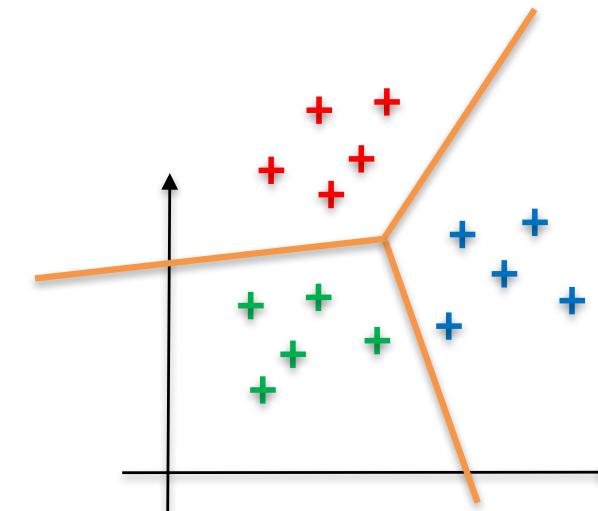
Overlapping clusters



Separated clusters

Limitation

- k-means algorithm assumes the **data are linearly separable** \Rightarrow Data can be separated by **straight lines**.
- k-means is **not** supposed to work for **data with more complex distributions** (solution is given by kernel k-means).



Summary

- k-means is the oldest unsupervised clustering algorithm (by physicist Lloyd in 1957).
- k-means is sound: Loss decrease and convergence guaranteed, speed is linear w.r.t. the number of data, the number of features and the number of clusters, easy to implement, and multiple extensions exist (k-medians, k-plans)
- k-means is limited: linear data and existence of local minimizers.

Outline

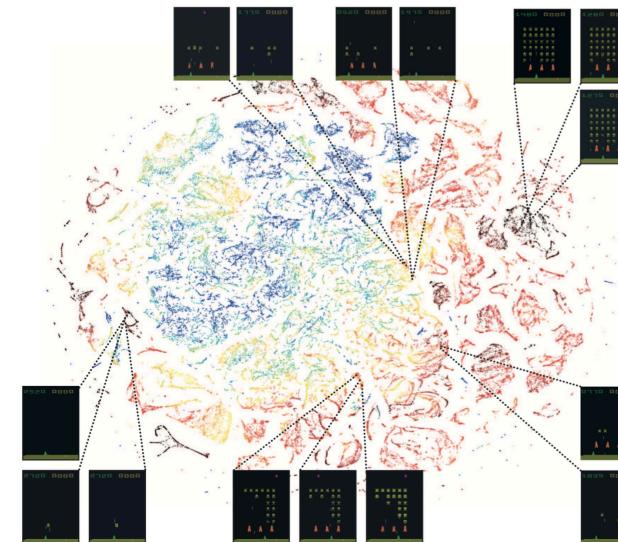
- Supervised vs unsupervised learning
- Unsupervised clustering with k-means
 - Clustering
 - k-means algorithm
 - Loss
 - k-means properties
- **Unsupervised representation with PCA**
 - **Introduction and motivations**
 - Principal directions and components
 - Formalization
 - Algorithm
 - PCA properties
- Conclusion

Dimensionality reduction

- Principal component analysis (PCA) is a dimensionality reduction technique (introduced by statistician Karl Pearson in 1901).
- Why dimensionality reduction is useful?
 - Data compression (compact representation)
 - Data visualization (2D and 3D)



Data compression
 $64 \times 64 \approx 4k$ reduced to
10 variables

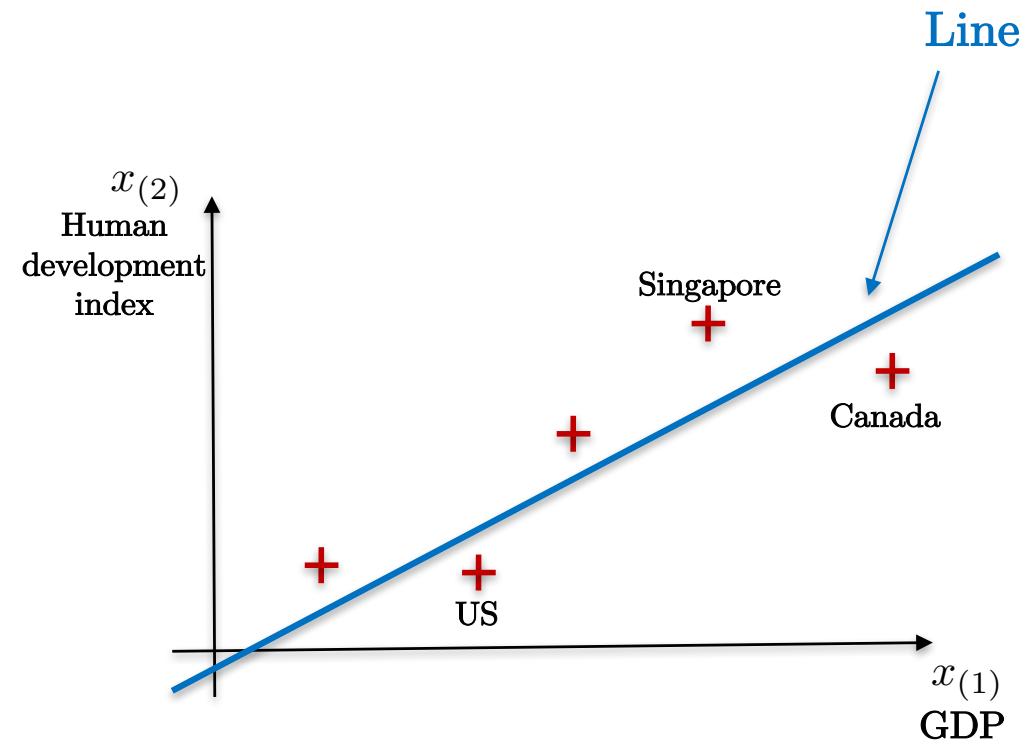


Data visualization
Video games

Data compression motivation

- Example: Let's consider some country properties

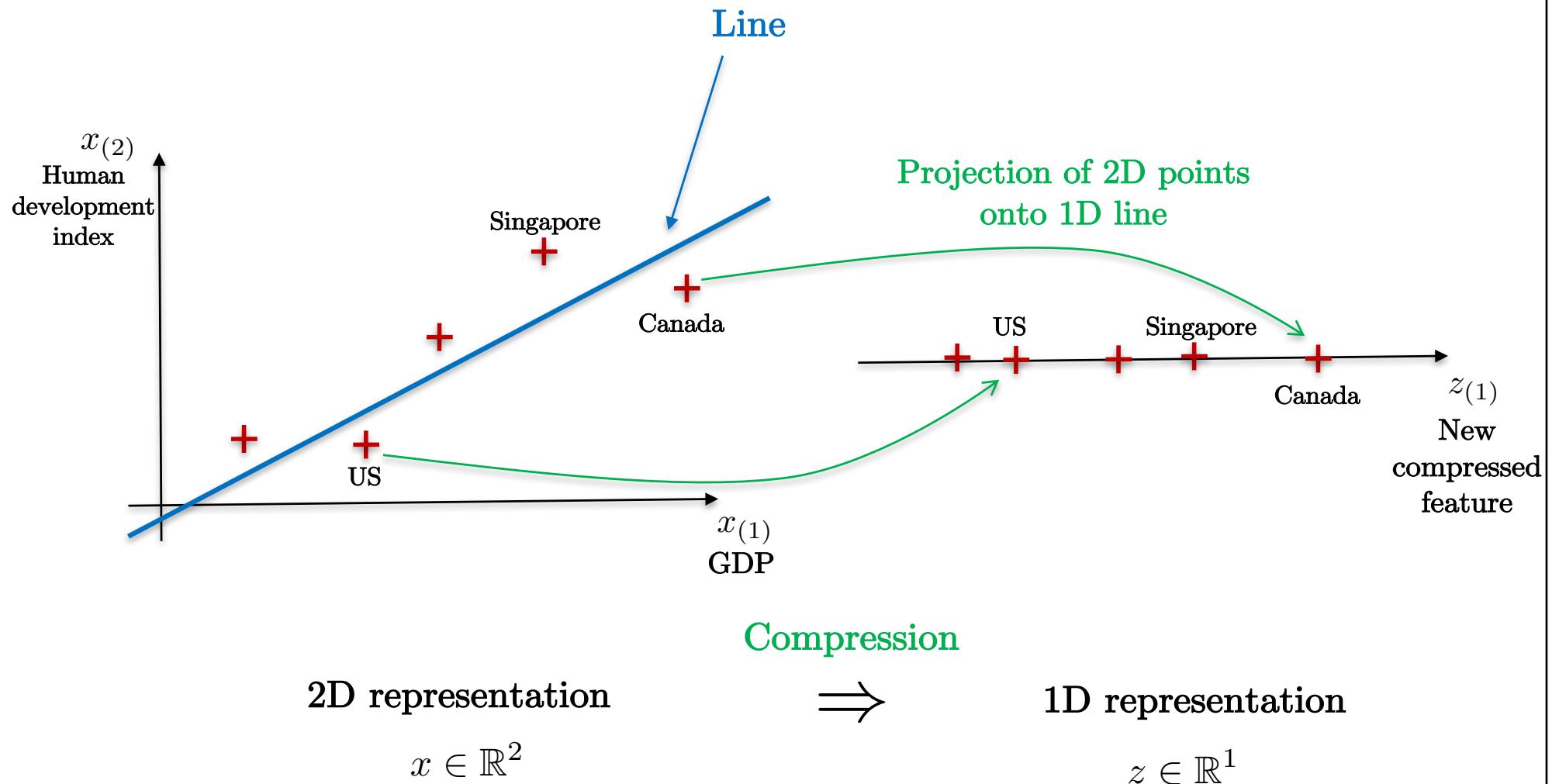
Country	GDP (gross domestic product in trillions of US\$) $x_{(1)}$	Human development index $x_{(2)}$
US	14.3	0.91
Canada	1.5	0.90
China	5.8	0.68
India	1.6	0.5
Singapore	0.2	0.8
Russia	1.4	0.7



Data are close to a line. There is some redundancy between the 2 features representing the country \Rightarrow Compression is possible.

Data compression motivation

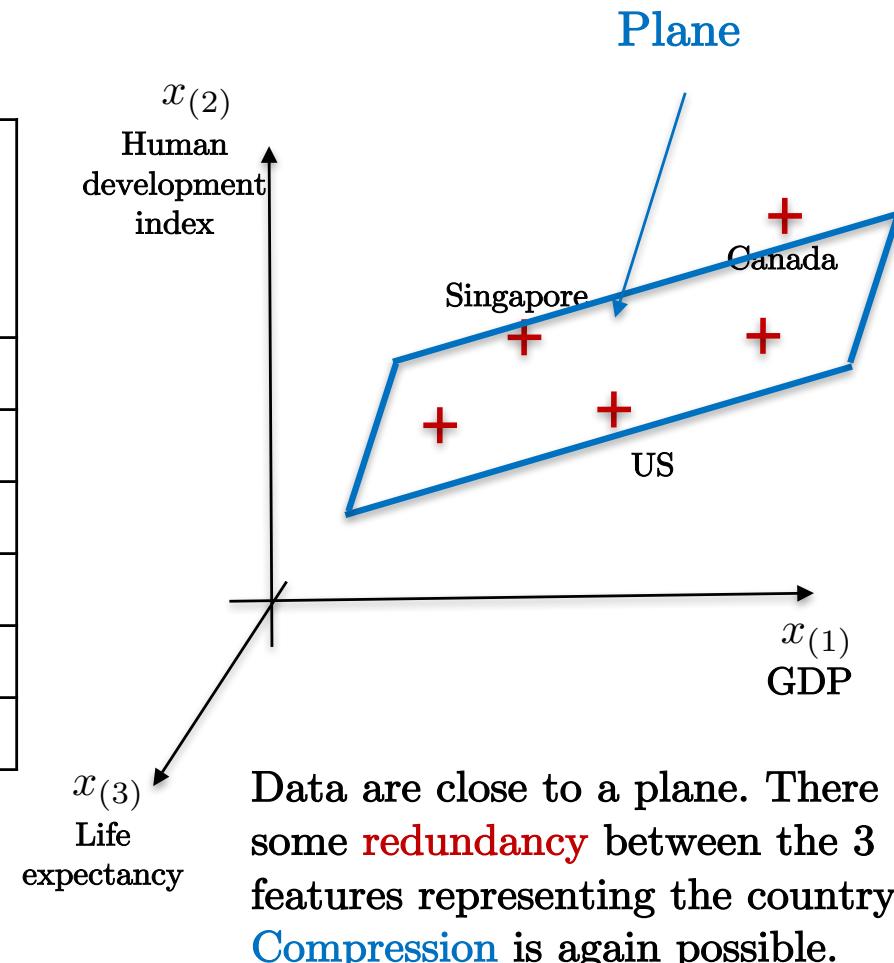
- We can **reduce** the data **from 2 features** (2D representation) to **1 feature** (1D representation):



Data compression motivation

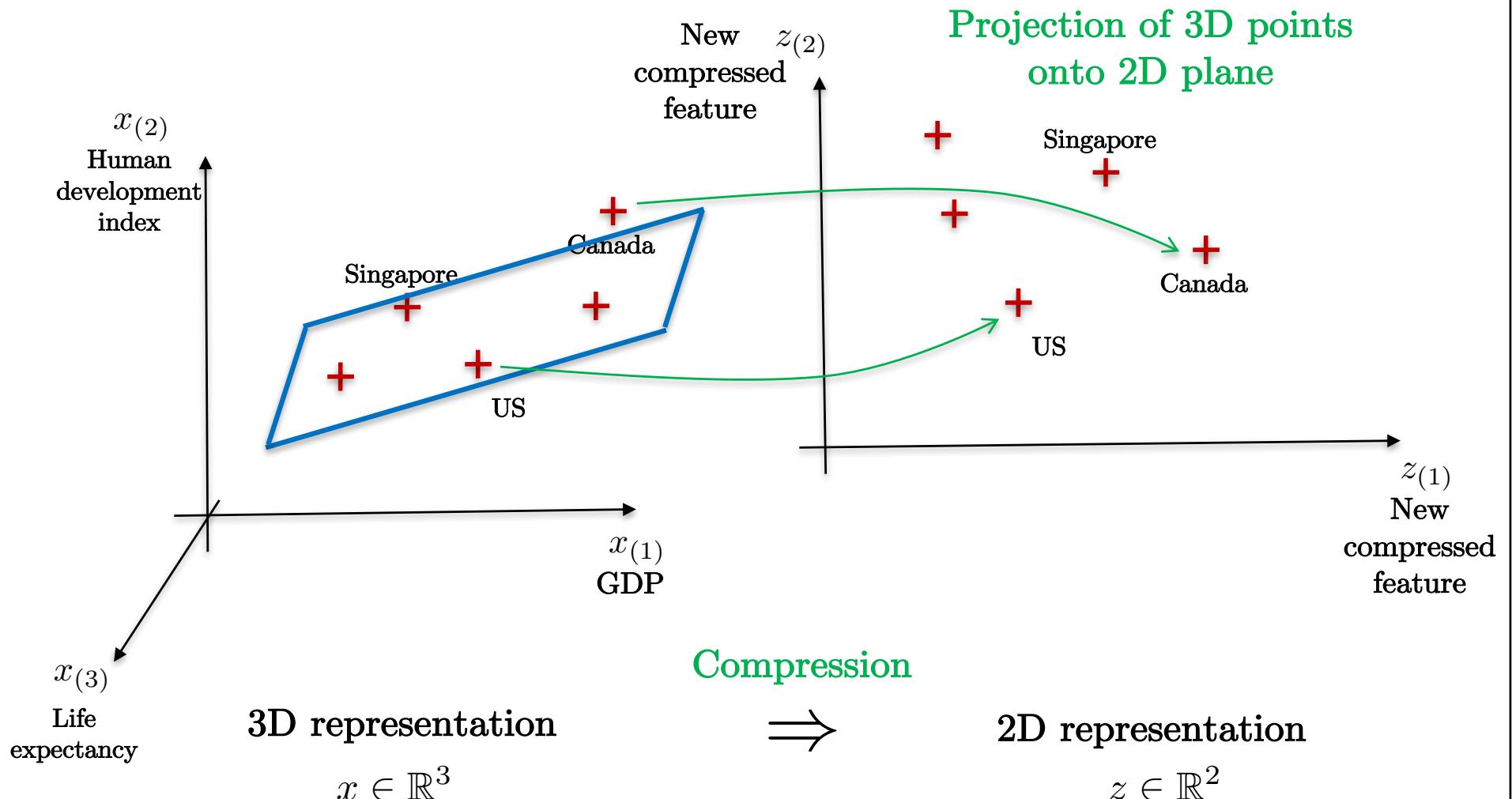
- Example: Let's consider three country properties

Country	GDP (gross domestic product in trillions of US\$) $x(1)$	Human development index $x(2)$	Life expectancy $x(3)$
US	14.3	0.91	78.3
Canada	1.5	0.90	80.7
China	5.8	0.68	73
India	1.6	0.5	64.7
Singapore	0.2	0.8	80
Russia	1.4	0.7	65.5



Data compression motivation

- We can also reduce the data from 3 features (3D representation) to 2 features (2D representation):



Data compression motivation

- General case: Data features can be compressed from 10,000-D compressed to 100-D.

Dim reduction/
Compression

10,0000-D representation \Rightarrow 100-D representation

$$x = \begin{bmatrix} x_{(1)} \\ x_{(2)} \\ \vdots \\ x_{(10,000)} \end{bmatrix} \in \mathbb{R}^{10,000}$$
$$z = \begin{bmatrix} z_{(1)} \\ z_{(2)} \\ \vdots \\ z_{(100)} \end{bmatrix} \in \mathbb{R}^{100}$$

- Main advantage of compression: Speeding up learning algorithm and less memory consuming.

Visualization motivation

- Data visualization is an exploratory tool to get insights about the data, and therefore better understand them.
- The problem is that data with more than 3 features cannot be visualized. How to visualize 10,000 dimensions?
⇒ Dimensionality reduction

Country	GDP (gross domestic product in trillions of US\$) $x_{(1)}$	Per capita GDP $x_{(2)}$	Human developme nt index $x_{(3)}$	Life expectanc y $x_{(4)}$	Poverty index (Gini) $x_{(5)}$	Mean household income (thousand s of US\$) $x_{(6)}$
US	14.3	46.7	0.91	78.3	40.8	84.3
Canada	1.5	39.1	0.90	80.7	32.6	67.2
China	5.8	7.5	0.68	73	46.9	10.22
India	1.6	3.4	0.5	64.7	36.8	0.73
Singapore	0.2	56.6	0.8	80	42.5	67.1
Russia	1.4	19.8	0.7	65.5	39.9	0.72

Each country
has 6 features

Visualization motivation

- Dimensionality reduction interpretation:

Country	GDP (gross domestic product in trillions of US\$) $x_{(1)}$	Per capita GDP $x_{(2)}$	Human development index $x_{(3)}$	Life expectancy $x_{(4)}$	Poverty index (Gini) $x_{(5)}$	Mean household income (thousands of US\$) $x_{(6)}$
US	14.3	46.7	0.91	78.3	40.8	84.3
Canada	1.5	39.1	0.90	80.7	32.6	67.2
China	5.8	7.5	0.68	73	46.9	10.22
India	1.6	3.4	0.5	64.7	36.8	0.73
Singapore	0.2	56.6	0.8	80	42.5	67.1
Russia	1.4	19.8	0.7	65.5	39.9	0.72

The meaning of hand-crafted features are **easy** to interpret.

Dim reduction/
Compression



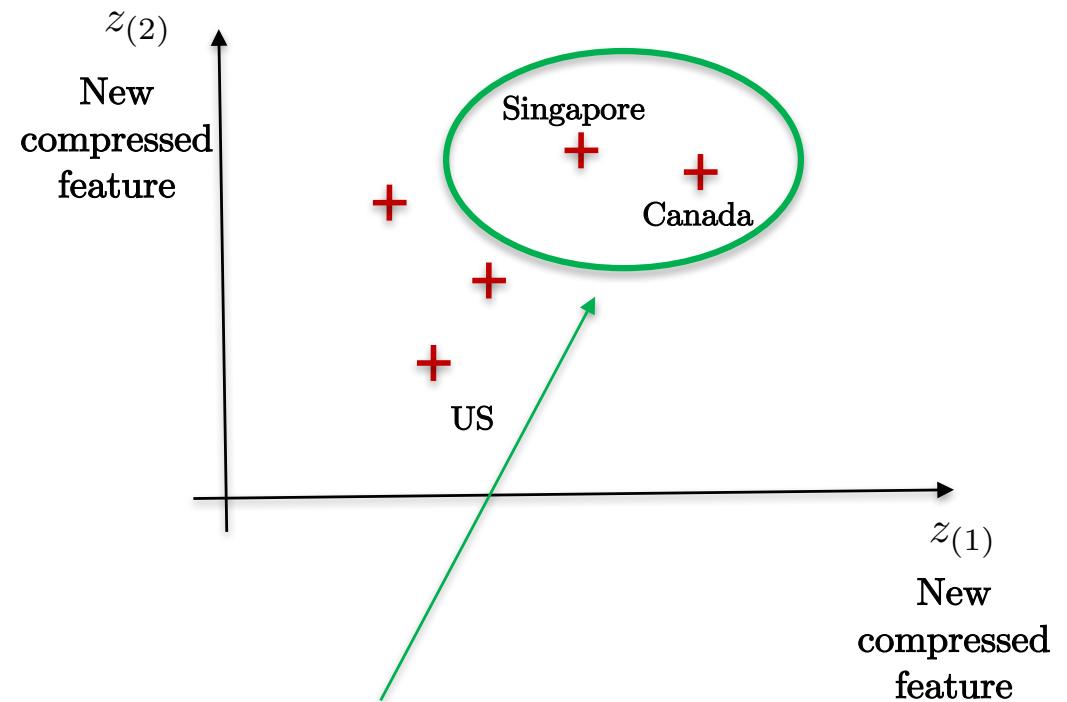
Country	Learned feature 1 $z_{(1)}$	Learned feature 2 $z_{(2)}$
US	-4.5	46.7
Canada	1.6	6.7
China	-5.8	7.5
India	1.6	-3.4
Singapore	0.2	3.2
Russia	2.4	-7.0

The meaning of learned features can be **hard** to interpret.

Visualization motivation

- Plot/visualization:

Country	Learned feature 1 $z(1)$	Learned feature 2 $z(2)$
US	-4.5	46.7
Canada	1.6	6.7
China	-5.8	7.5
India	1.6	-3.4
Singapore	0.2	3.2
Russia	2.4	-7.0



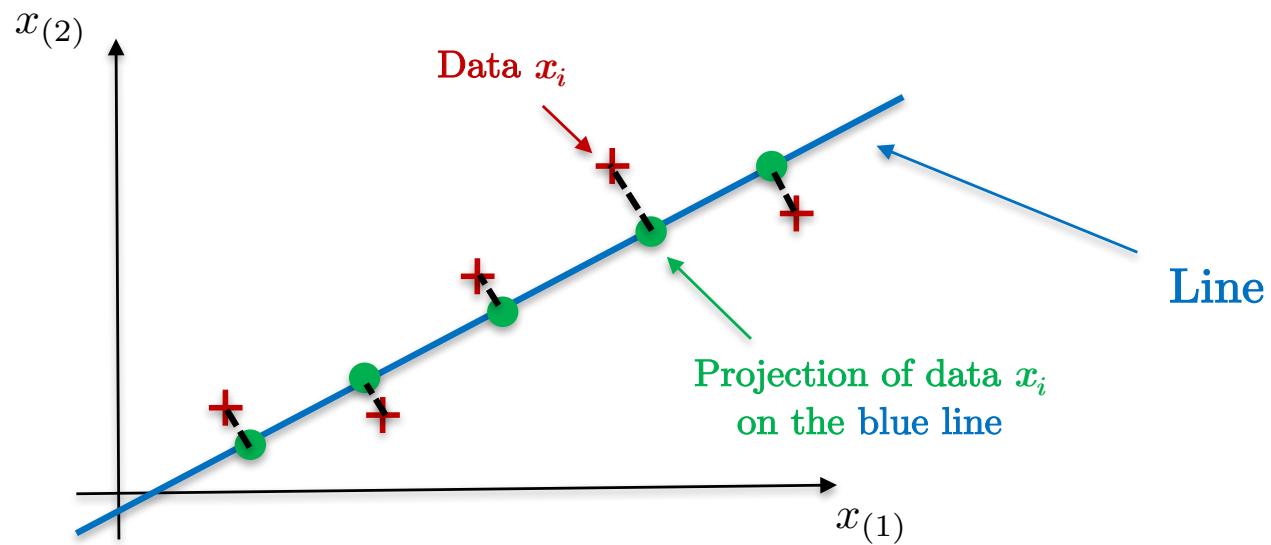
Although interpreting the new features may be hard,
the dim reduction technique always put together
data that are similar \Rightarrow Clustering property of dim
reduction. It makes easier to analyze data.

Outline

- Supervised vs unsupervised learning
- Unsupervised clustering with k-means
 - Clustering
 - k-means algorithm
 - Loss
 - k-means properties
- **Unsupervised representation with PCA**
 - Introduction and motivations
 - **Principal directions and components**
 - Formalization
 - Algorithm
 - PCA properties
- Conclusion

Basic idea

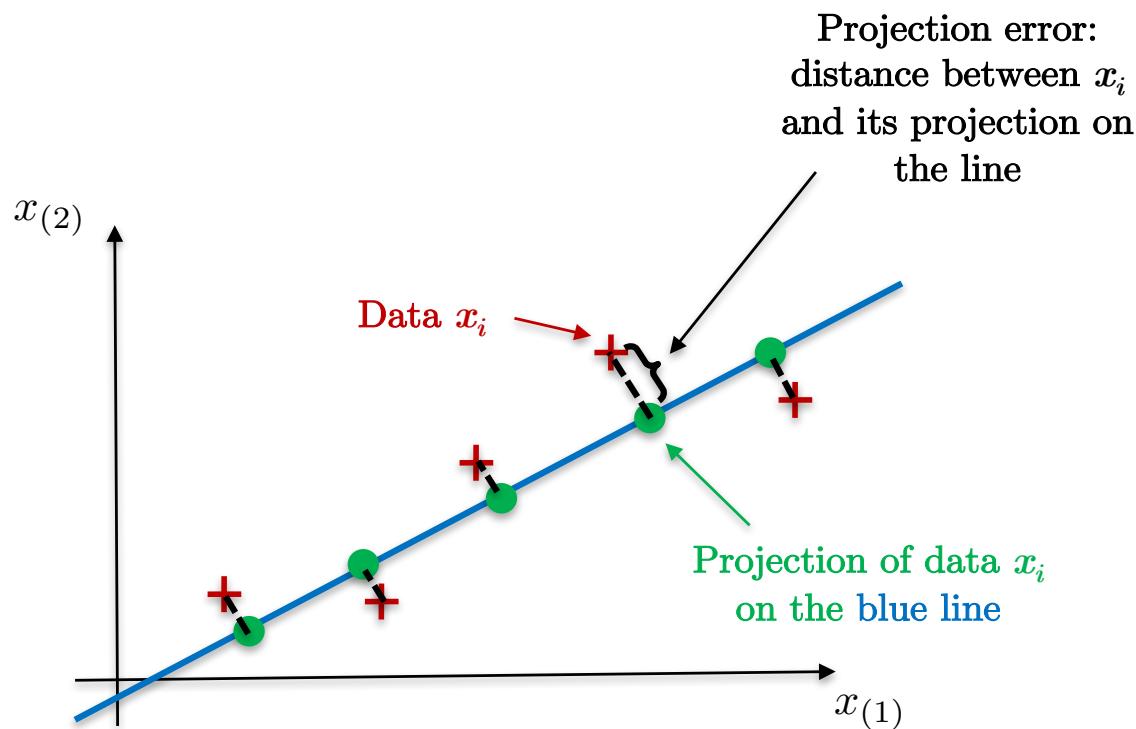
- Reduce the dimensionality of data by projecting the data on the closest line for 1D reduction, the closest plane for 2D reduction, etc.



- How to define a “good” line, a “good” plane, etc?

Basic idea

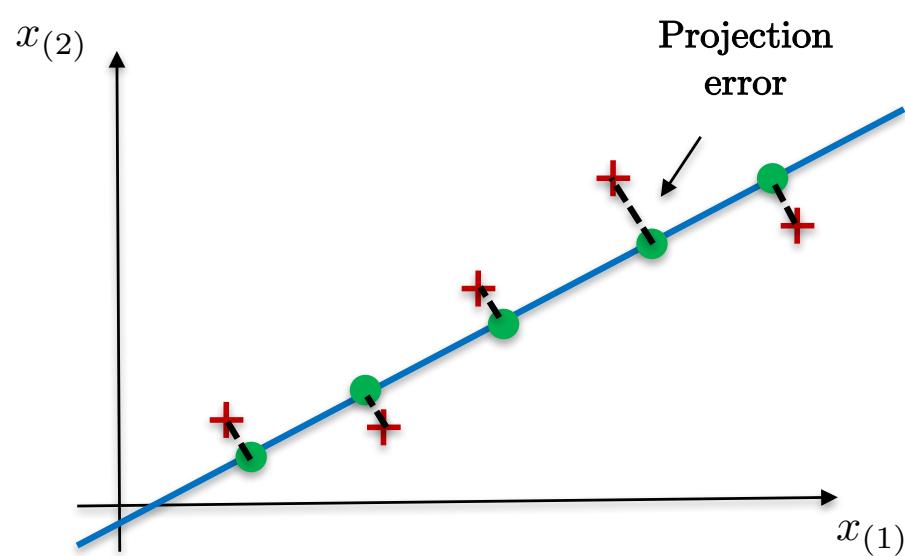
- What is a good line? The line that minimizes as much as possible the projection errors.



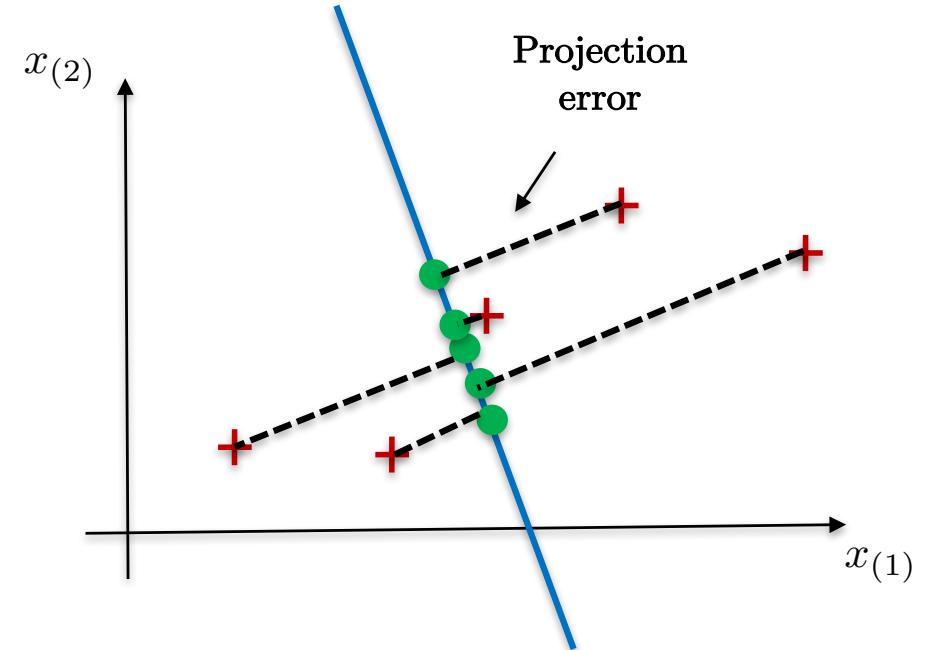
The best blue line if the line that minimizes the sum of all projection errors.

Basic idea

- Different lines:



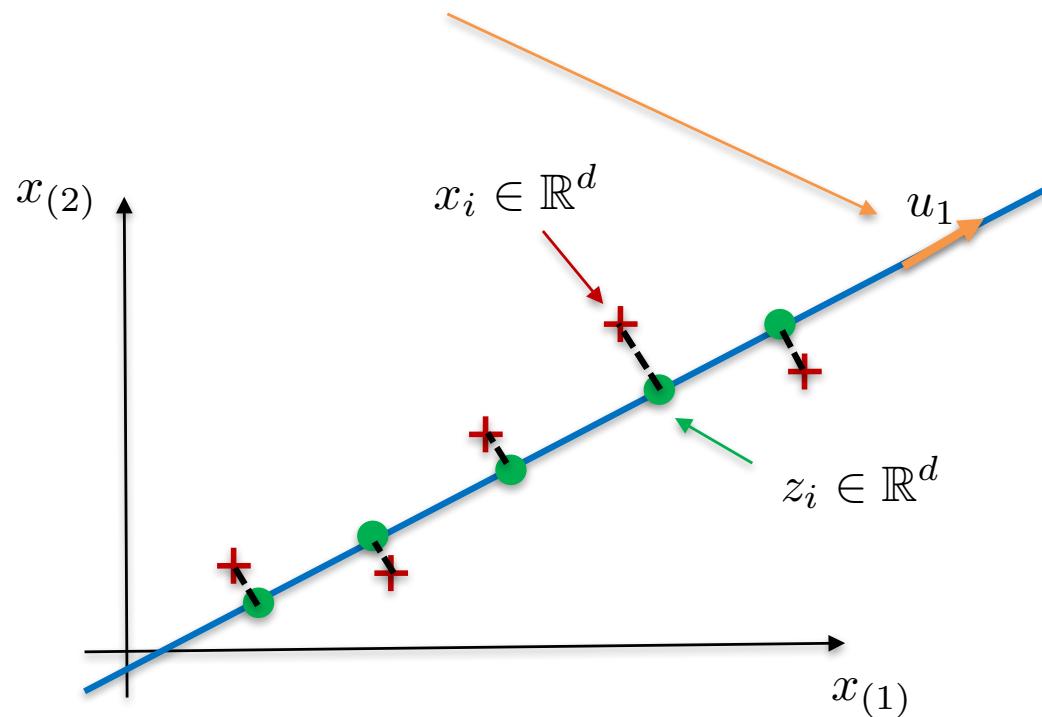
Sum of projection errors is **small** (and **minimized**).



Sum of projection errors is **high**.

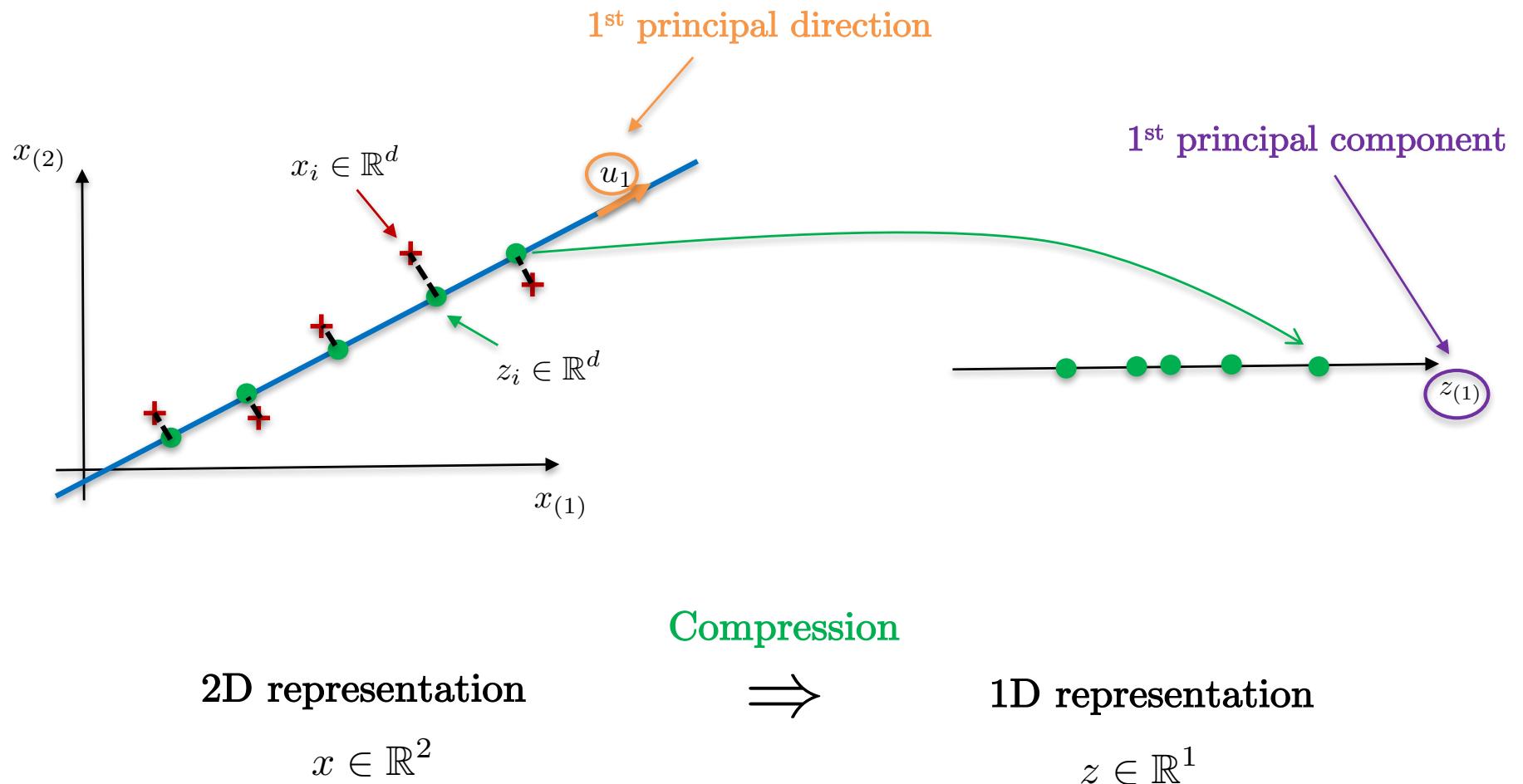
Notations

- x_i : i^{th} data point
- $x_{i(j)}$: j^{th} data feature of the i^{th} data point
- z_i : projection of the i^{th} data point
- $z_{i(j)}$: j^{th} learned feature of the i^{th} projected data point
- u_k : unit vector of the k^{th} principal direction (the line if $k=1$)



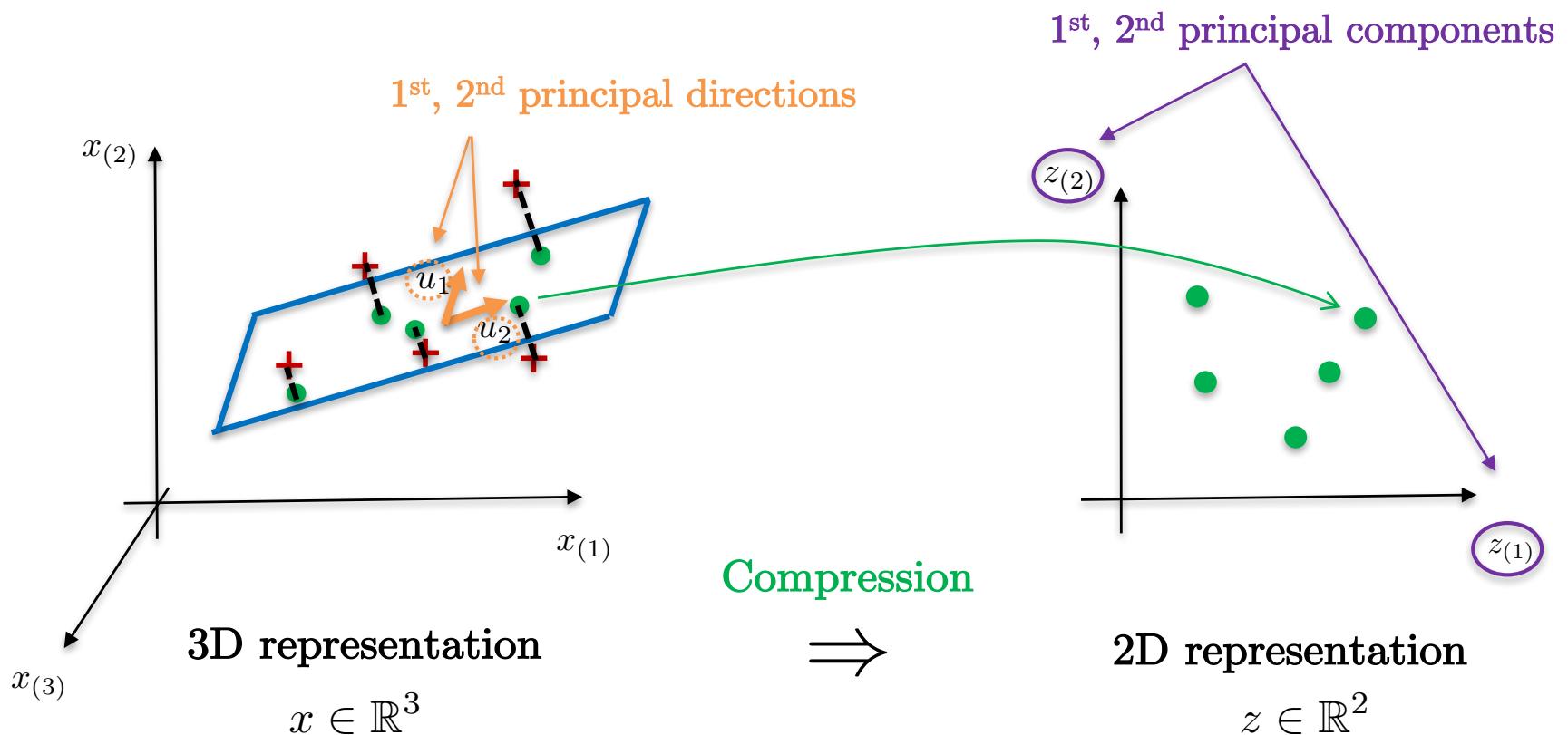
Principal directions and components

- **1st principal direction:** The direction u_1 in \mathbb{R}^2 onto which the projection of data is minimized and the coordinate $z_{(1)}$ in \mathbb{R}^1 on the projected direction u_1 .



Principal directions and components

- **1st, 2nd principal directions:** The directions u_1, u_2 in \mathbb{R}^3 onto which the projection of data is minimized and the coordinate $z_{(1)}, z_{(2)}$ in \mathbb{R}^2 on the projected directions u_1, u_2 .



- **1st, ..., kth principal directions:** The directions u_1, \dots, u_k in \mathbb{R}^d onto which the projection of data is minimized and the coordinate $z_{(1)}, \dots, z_{(k)}$ in \mathbb{R}^k on the projected directions u_1, \dots, u_k .

PCA pre-processing

- Required pre-processing: Data centering (zero mean)
 - Training set: x_1, \dots, x_n
 - Mean vector: $\mu = \frac{1}{n} \sum_{i=1}^n x_i \in \mathbb{R}^d$
 - Centering data: $x_i \leftarrow x_i - \mu$
- Other pre-processing: z-scoring may help or not.

$$x_i \leftarrow \frac{x_i - \mu}{\sigma}$$

Outline

- Supervised vs unsupervised learning
- Unsupervised clustering with k-means
 - Clustering
 - k-means algorithm
 - Loss
 - k-means properties
- **Unsupervised representation with PCA**
 - Introduction and motivations
 - Principal directions and components
 - **Formalization**
 - Algorithm
 - PCA properties
- Conclusion

Formalization

- PCA technique finds a linear space (line, plane, hyper-plane) that minimizes the sum of all projection errors.
- How to find this linear space?
 - Linear algebra !



Co-variance matrix

- Matrix of all variances and co-variances:

$$\underset{d \times d}{\Sigma} = \frac{1}{n} \sum_{i=1}^n \underbrace{\underbrace{x_i x_i^T}_{\substack{1 \times d \\ \underbrace{\quad\quad\quad}_{d \times d}}} \underbrace{\quad\quad\quad}_{d \times 1}}_{d \times d} = \frac{1}{n} X^T X \quad \text{with data matrix } X = \begin{bmatrix} -x_1^T- \\ \vdots \\ -x_n^T- \end{bmatrix}_{n \times d}$$

Co-variance
matrix
(encode all data
variations)

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} & \dots & \Sigma_{1d} \\ & \ddots & & \\ \Sigma_{d1} & \Sigma_{d2} & \dots & \Sigma_{dd} \end{bmatrix}$$

$$\left\{ \begin{array}{l} \Sigma_{jj} = \sum_{i=1}^n x_{i(j)}^2 \quad \text{data variance along } j\text{-dim} \\ \Sigma_{jl} = \sum_{i=1}^n x_{i(j)} x_{i(l)} \quad \text{data co-variance along } j\text{-dim and } l\text{-dim} \end{array} \right.$$

Principal directions

- The principal directions u_1, \dots, u_k in \mathbb{R}^d onto which the projection of data is minimized are given by the eigenvalue decomposition (EVD) of the co-variance matrix (no proof given):

$$\Sigma = \underbrace{U S U^T}_{\substack{d \times d \quad d \times d \quad d \times d \\ d \times d}}$$

Eigenvalue decomposition is a **matrix factorization** technique.

Eigenvector matrix: U

Eigenvalue matrix: S

$$U = \underbrace{\left[\begin{array}{ccc|c} | & & | & \\ u_1 & \dots & u_k & \\ | & & | & \\ \hline & & & \end{array} \right]}_{\substack{d \times d \\ d \times k}} \dots \underbrace{\left[\begin{array}{c|c} | & \\ u_d & | \end{array} \right]}_{d \times d}$$

Matrix of k principal directions
 u_1, \dots, u_k in \mathbb{R}^d

$$u_p^T u_q = \begin{cases} 1 & \text{if } p = q \\ 0 & \text{if } p \neq q \end{cases}$$

$$u_p^T u_p = \|u_p\|_2^2 = 1, \forall p \quad (\text{unit vector})$$

$$u_p^T u_q = 0, \forall p \neq q \quad (\text{orthonormal vectors})$$



Principal directions

- u_1 represents the direction of the largest data variance.
- u_2 represents the direction of the second largest data variance.
- ...
- S_{11} is the value of the data variance along the direction u_1 .
- S_{22} is the value of the data variance along the direction u_2 .
- ...

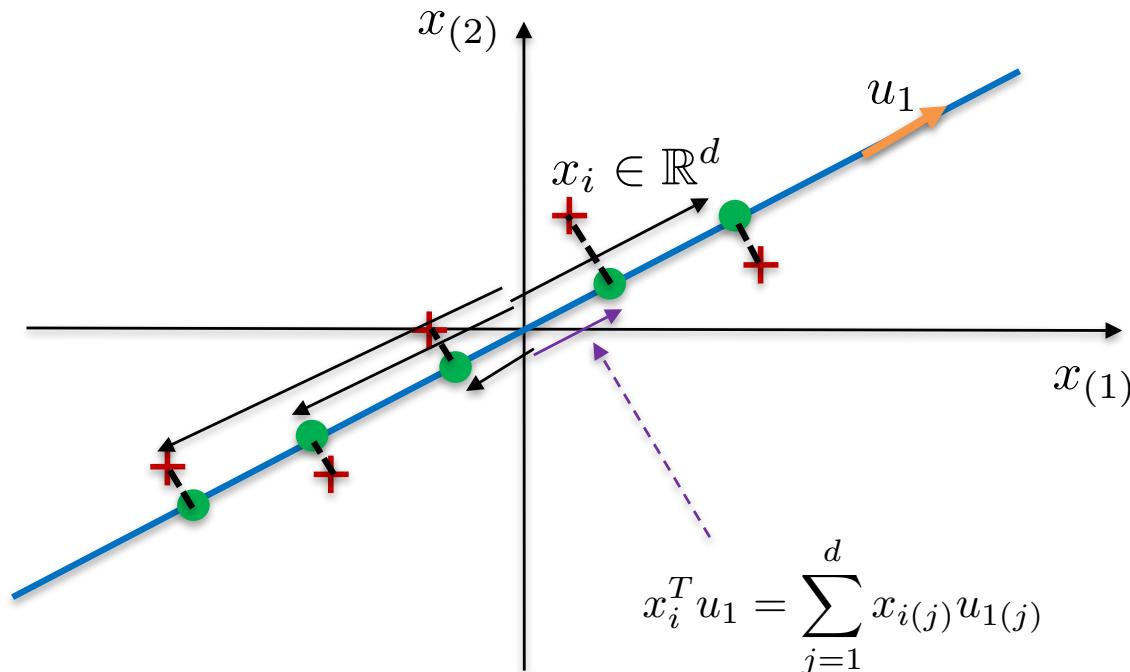
$$\Sigma = U S U^T \quad \text{with} \quad d \times d \quad S = \begin{bmatrix} S_{11} & & & 0 \\ & S_{22} & & \\ & & \ddots & \\ 0 & & & S_{dd} \end{bmatrix}$$

$$\text{and} \quad S_{ll} = \sum_{i=1}^n |x_i^T u_l|^2 \quad (\text{no proof given})$$

Principal directions capture the data variance

- S_{11} represents the variance of the data along the direction u_1 of the largest data variance :

$$S_{11} = \sum_{i=1}^n |x_i^T u_1|^2 \quad \text{Sum of square of projected data along } u_1$$



Outline

- Supervised vs unsupervised learning
- Unsupervised clustering with k-means
 - Clustering
 - k-means algorithm
 - Loss
 - k-means properties
- **Unsupervised representation with PCA**
 - Introduction and motivations
 - Principal directions and components
 - Formalization
 - **Algorithm**
 - PCA properties
- Conclusion

PCA algorithm

- Algorithm:

- Pre-process data: zero-mean

$$x_i \leftarrow x_i - \frac{1}{n} \sum_{i=1}^n x_i$$

- Construct the co-variance matrix:

$$\Sigma = \frac{1}{n} X^T X$$

- Compute the EVD of Σ (principal directions): $\Sigma = U S U^T$

- Compute the principal components (or projected data):

$$Z = U_k^T X^T$$

$k \times n$ $k \times d$ $d \times n$

Projected data
coordinates
on u_1, \dots, u_k

$$X^T = \begin{bmatrix} | & & | \\ x_1 & \dots & x_n \\ | & & | \end{bmatrix}$$

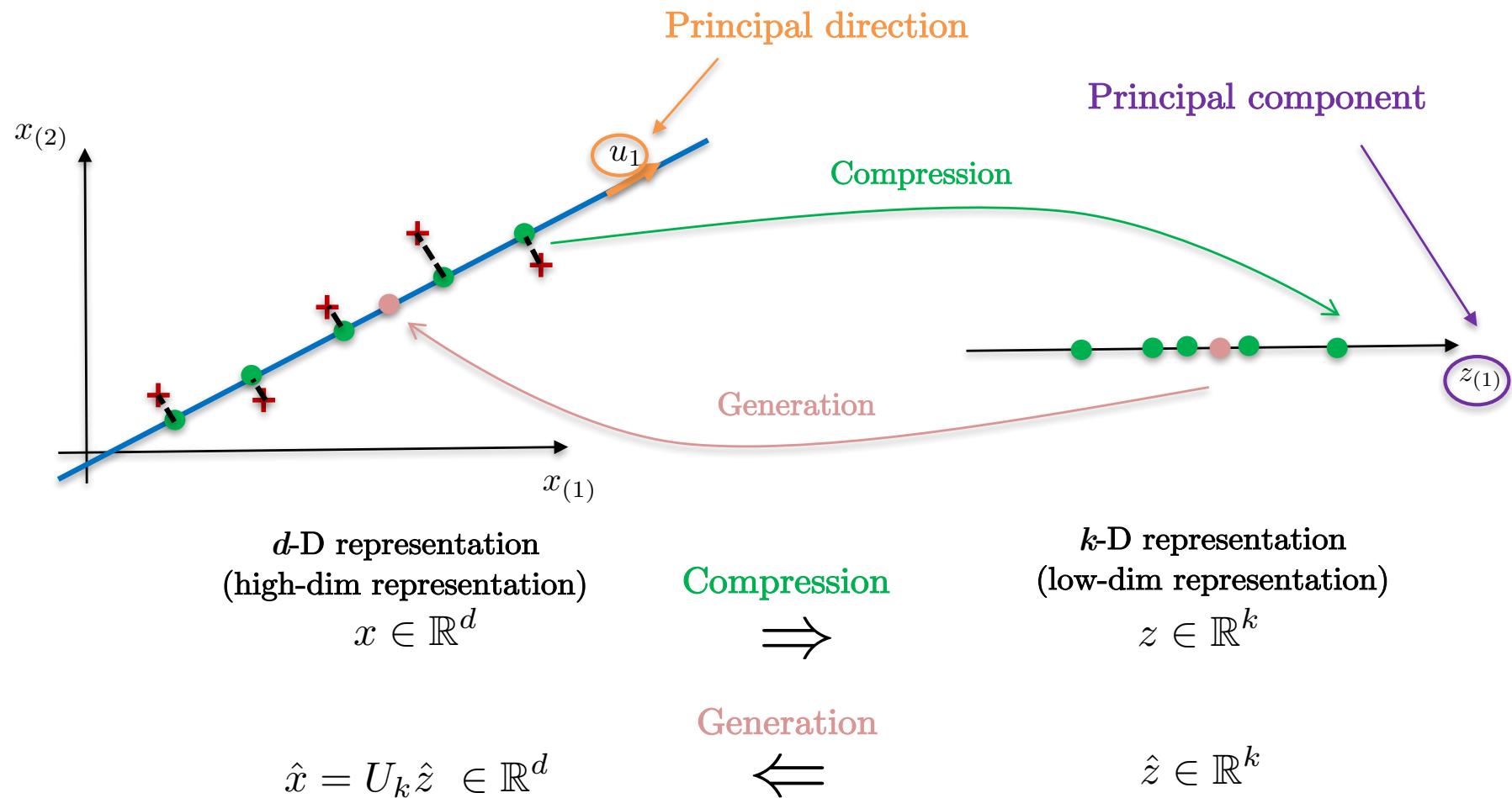
$$U_k^T = \begin{bmatrix} -u_1^T - \\ \vdots \\ -u_k^T - \end{bmatrix} \begin{matrix} Z \\ \downarrow \\ \downarrow \end{matrix}$$

Outline

- Supervised vs unsupervised learning
- Unsupervised clustering with k-means
 - Clustering
 - k-means algorithm
 - Loss
 - k-means properties
- **Unsupervised representation with PCA**
 - Introduction and motivations
 - Principal directions and components
 - Formalization
 - Algorithm
 - **PCA properties**
- Conclusion

Data generation

- We can construct new data from the compressed representation:

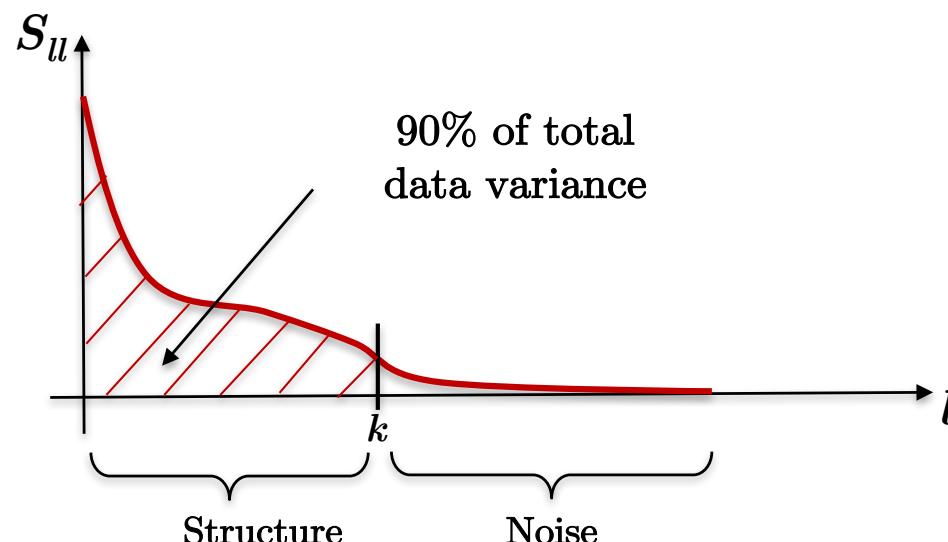


How to select k ?

- **Selection rule:** Data variations are captured by each principal direction \Rightarrow A natural rule is to **select the first k principal directions that capture e.g. 90% of total variance:**

$$\frac{\sum_{l=1}^k S_{ll}}{\sum_{l=1}^d S_{ll}} \geq 0.9$$

S_{ll} represents the variance of the data along the direction u_l



PCA in practice

- Applications of PCA:
 - Visualization:
 - Insights about data properties
 - Compression:
 - Reduce memory needed for data
 - Speed up learning speed

$$(x_1, y_1), \dots, (x_n, y_n), x_i \in \mathbb{R}^d$$

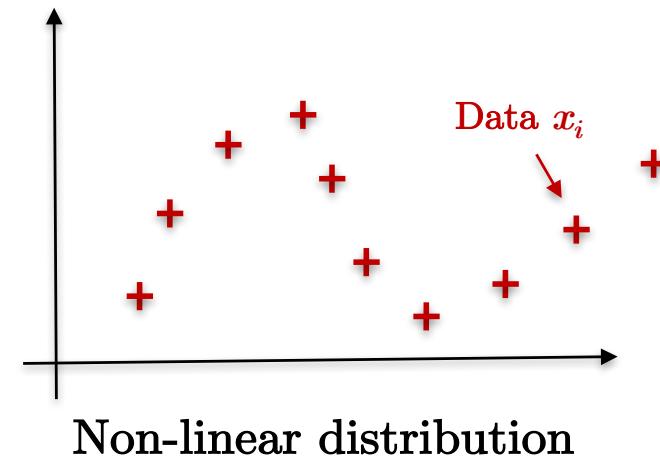
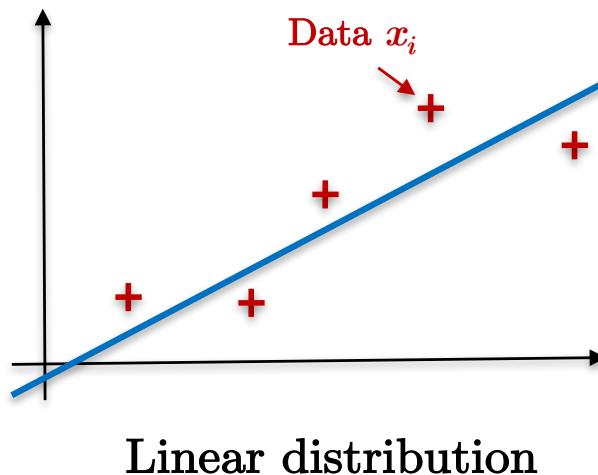


$$(z_1, y_1), \dots, (z_n, y_n), z_i \in \mathbb{R}^k$$

- Example: Images $d=256 \times 256 \Rightarrow k=128$
- Important: There is no guarantee that applying PCA will improve the performances of learning algorithm. Actually, **PCA usually decreases the results, but it is fine if speed is a priority.**

Limitation

- PCA is a linear dimensionality reduction technique \Rightarrow It means that it is guaranteed to work when the data follow a linear distribution, but no guarantee otherwise.



Summary

- PCA is the most common and the oldest linear dimensionality reduction technique.
- There exist several improvements of PCA: sparse PCA, robust PCA, graph PCA, non-linear PCA, etc.
- PCA algorithm is simple, sound, linear w.r.t. number of training data and number of principal components, but does not scale w.r.t. number of features d (complexity is $O(d^3)$).
- PCA is optimal for data with linear distribution.

Outline

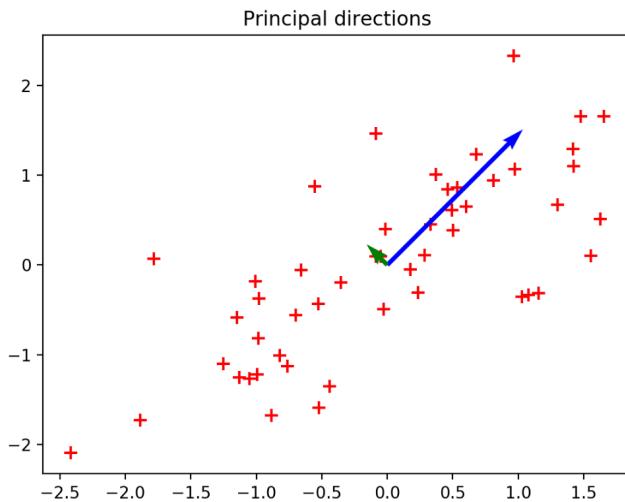
- Supervised vs unsupervised learning
- Unsupervised clustering with k-means
 - Clustering
 - k-means algorithm
 - Loss
 - k-means properties
- Unsupervised representation with PCA
 - Introduction and motivations
 - Principal directions and components
 - Formalization
 - Algorithm
 - PCA properties
- Conclusion

Conclusion

- Unsupervised learning techniques solve data analysis tasks **without any label information.**
- Unsupervised techniques can be applied to data clustering, data representation (beyond linear PCA s.a. **generative adversarial networks** (GAN)), data visualization(beyond PCA s.a. **t-SNE**), etc.
- **Unsupervised** learning problems are **much more challenging** to solve than **supervised** learning problems.
- The next AI revolution will be **unsupervised** [LeCun].

Coding exercise

- [tutorial07.ipynb](#)



1.5 Main principal directions

Data variations are captured by each principal direction. A natural rule is to select the first k principal directions that capture e.g. 85% of total variance:

$$\frac{\sum_{l=1}^k S_{ll}}{\sum_{l=1}^d S_{ll}} \geq 0.85$$

How many principal directions do we need to capture 85% of the total data variance?



Questions?