

Machine Learning I
80-629A

Apprentissage Automatique I
80-629

Unsupervised Learning
— Week #7

Today

- **Unsupervised learning**
 - **Definition and properties**
 - **A few clustering models**
 - **Other instantiations of unsupervised learning**

Unsupervised Learning

Experience (E)

- What data does f experience?
 - (Focus on algorithms that experience whole datasets)
 - **Unsupervised.** Examples alone $\{\mathbf{x}_i\}_{i=0}^n$
 - **Supervised.** Examples come with labels $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=0}^n$

1. Unsupervised

$$\{\mathbf{x}_i\}_{i=0}^n$$

- Experience examples alone
- Learn “useful properties of the structure of the data”
 - E.g., clustering, density modeling ($p(\mathbf{x})$), PCA, FA.

Different tasks

- Finding patterns
 - Clustering $f : X \rightarrow \{1, 2, \dots, K\}$ (K clusters)
 - Dimensionality reduction $f : X^p \rightarrow X^k, k \ll p$
 - Density modelling $f : X \rightarrow [0, 1]$
 - ...

Why unsupervised learning?

1. Understand properties of the data
2. Learn useful representations
3. Use the results in a downstream application
 - Opportunity: There are lots of unlabelled data

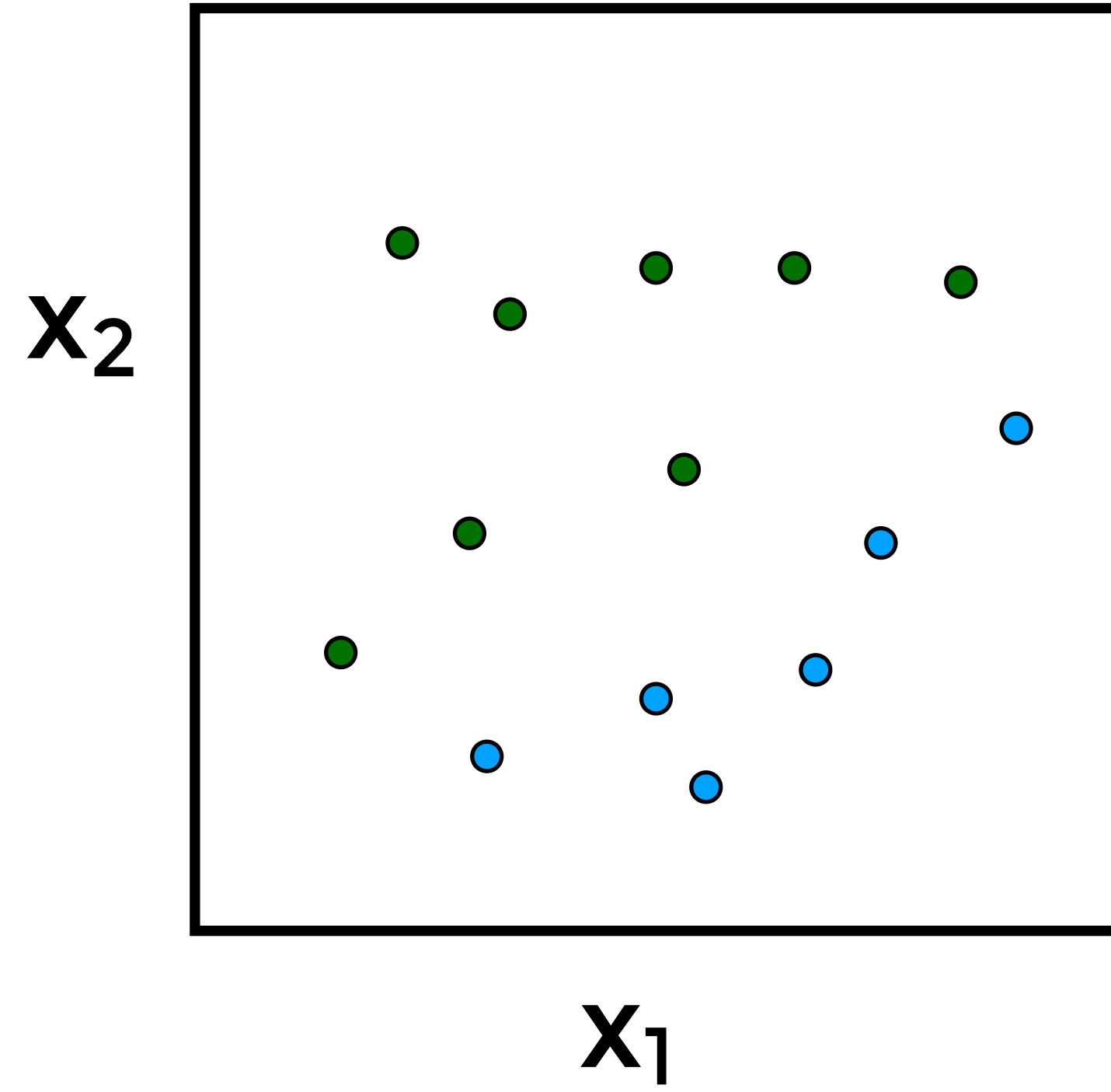
k-means clustering

Clustering

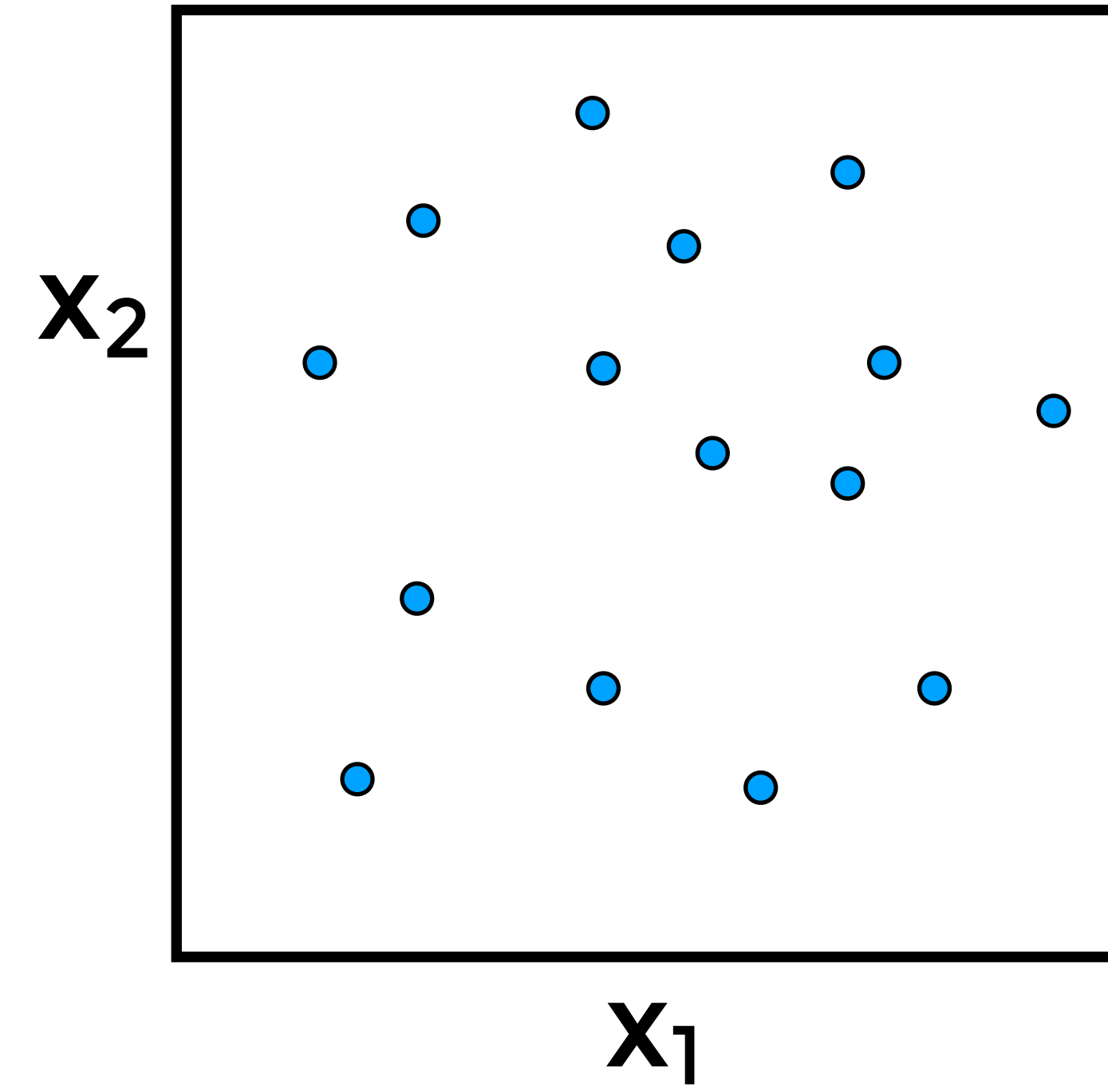
$$f : X \rightarrow \{1, \dots, K\}$$

- **Task:** Assign each point to one of K clusters
- **Cluster:** a set of similar points (a group)
- **Alternatively:** Divide the space into K regions. Assign points to a cluster based on the region they lie in
- **Similar to classification.**

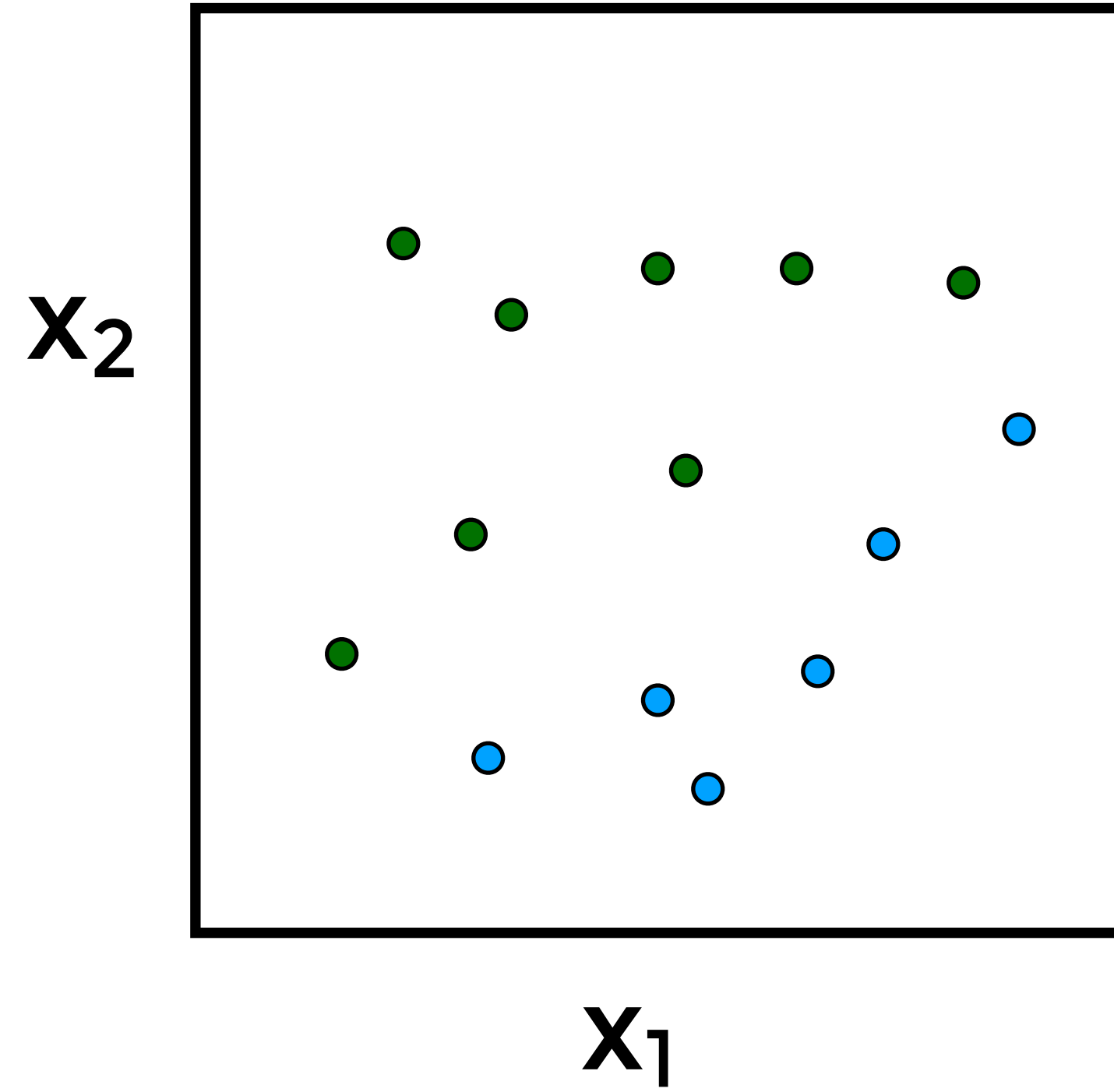
Supervised



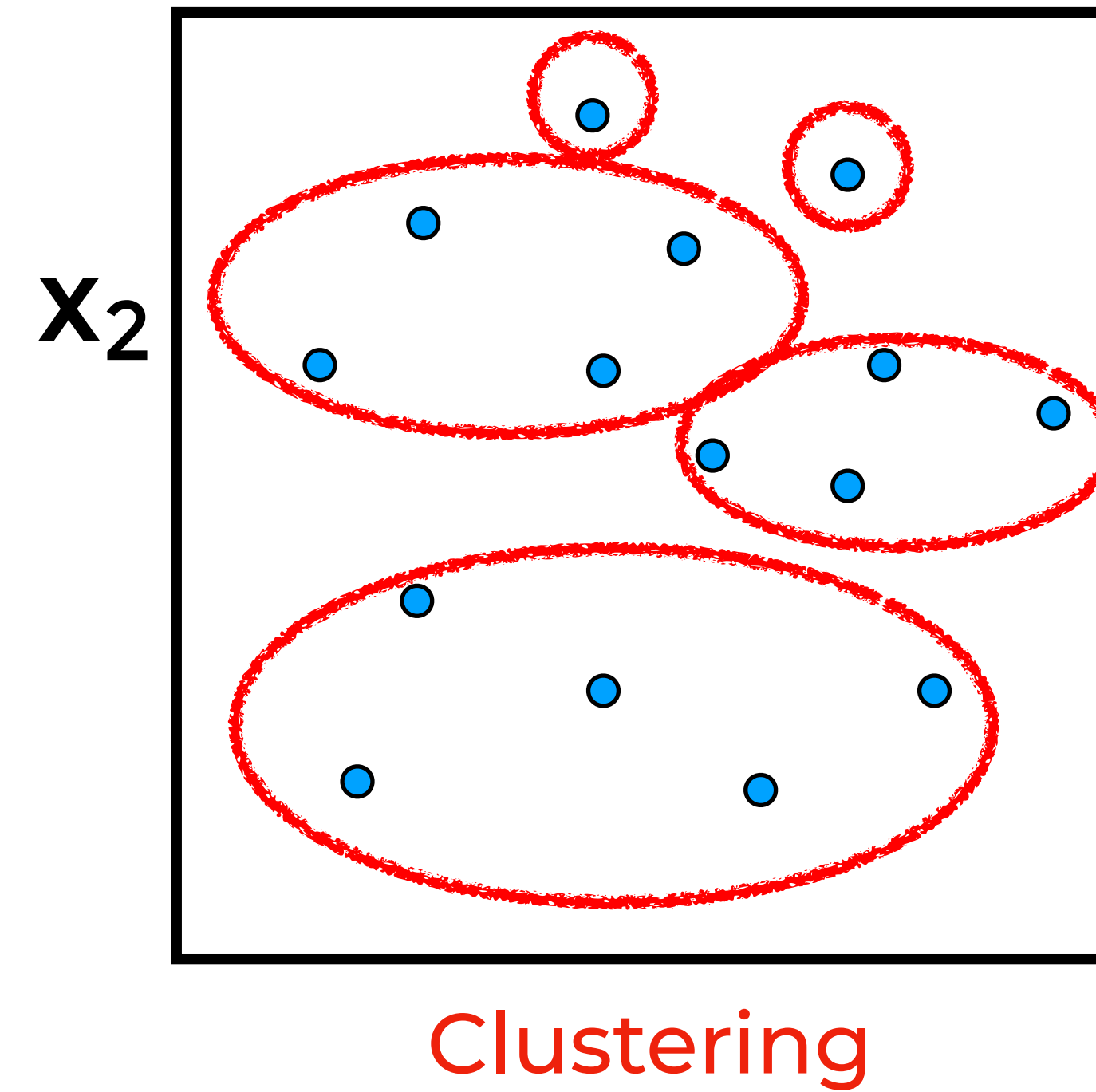
Unsupervised



Supervised



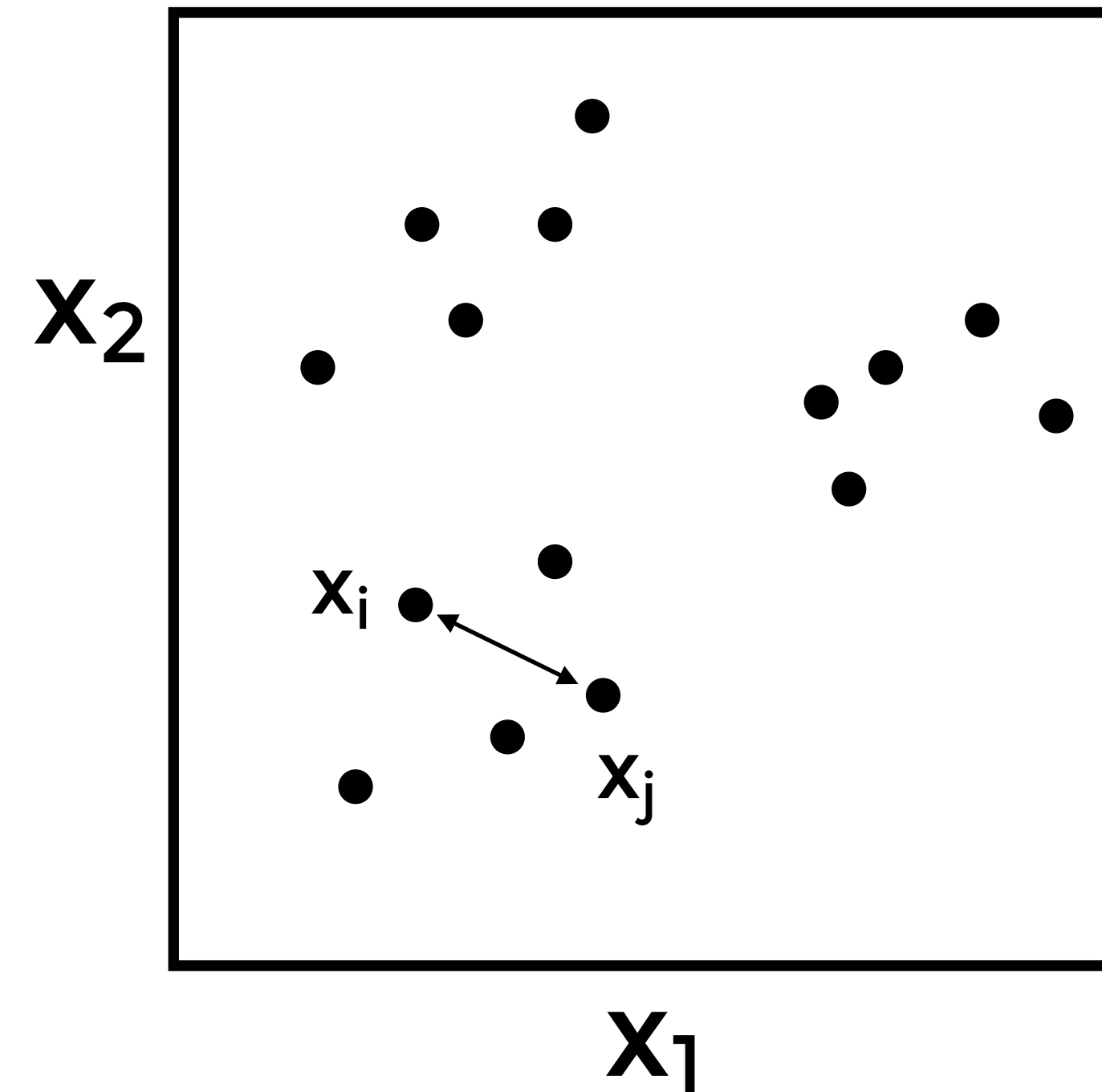
Unsupervised



Clustering

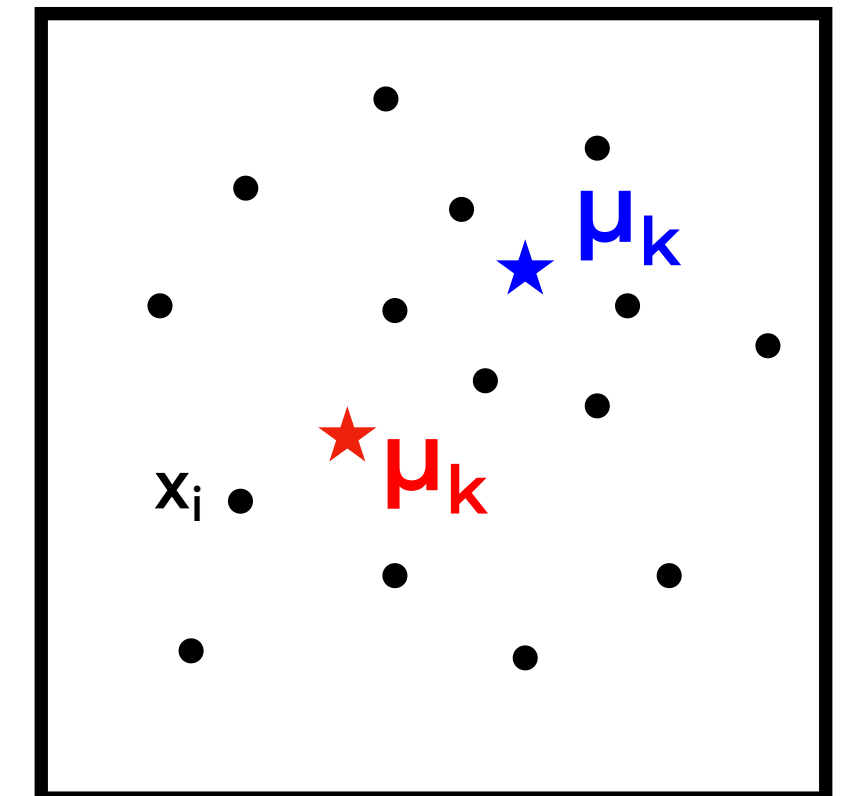
- Desideratum: group similar points together
- Similarity is often defined as being close according to some similarity or (inverse) distance function
- E.g., in Euclidean space

$$\text{distance}(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$$



K-means clustering

- A particular clustering model (and accompanying algorithm)
 - There are K clusters. Each point belongs to a cluster. Clusters have centers: μ
- Objective: Find cluster centers μ_k that minimize the within cluster distance

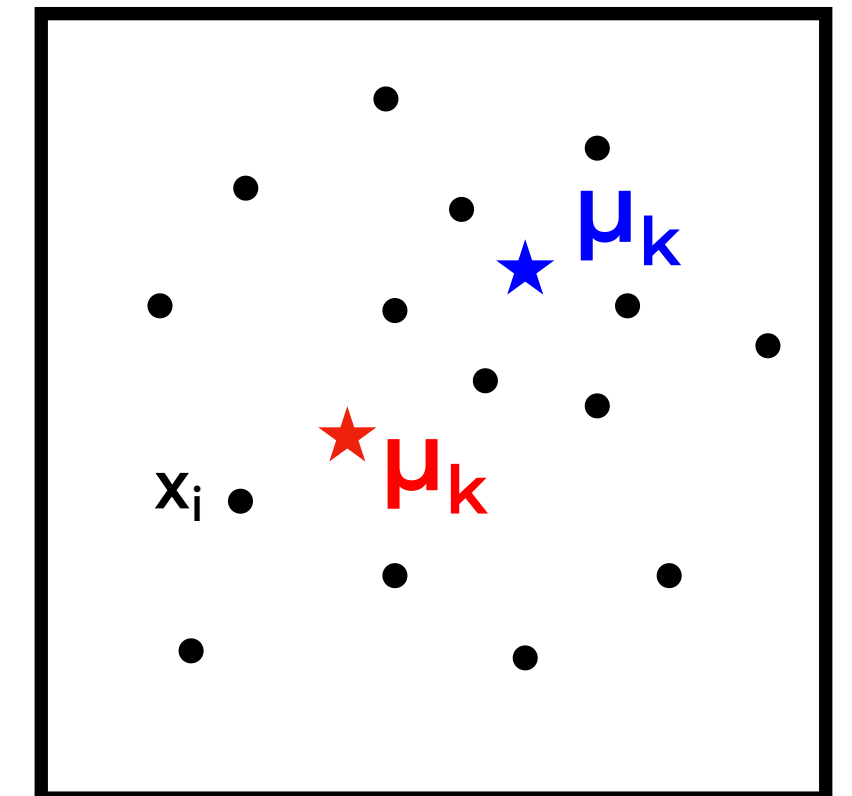


K-means clustering

- A particular clustering model (and accompanying algorithm)
 - There are K clusters. Each point belongs to a cluster. Clusters have centers: μ
- Objective: Find cluster centers μ_k that minimize the within cluster distance

$$\text{Objective} := \sum_{i=1}^N \sum_{k=1}^K r_{ik} \|x_i - \mu_k\|^2$$

$$r = \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ \vdots & \vdots \\ 0 & 1 \end{bmatrix}_{N \times 2}$$



K-means clustering

- A particular clustering model (and accompanying algorithm)
 - There are K clusters. Each point belongs to a cluster. Clusters have centers: μ
- Objective: Find cluster centers μ_k that minimize the within cluster distance

$$\text{Objective} := \sum_{i=1}^N \sum_{k=1}^K r_{ik} \|x_i - \mu_k\|^2$$

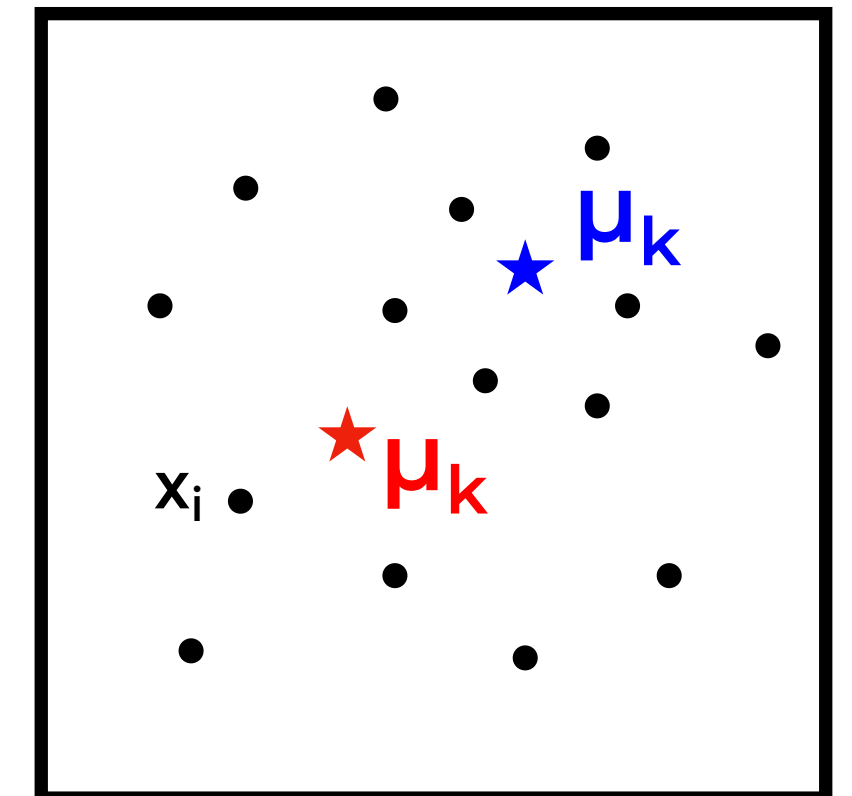
- Algorithm to minimize the objective:

- Initialize the cluster centers
- Until convergence:

1. Update responsibilities: r

2. Update cluster centers: $\mu_k \quad \forall k$

$$r = \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ \vdots & \vdots \\ 0 & 1 \end{bmatrix}_{N \times 2}$$



K-means clustering

- A particular clustering model (and accompanying algorithm)
 - There are K clusters. Each point belongs to a cluster. Clusters have centers: μ
- Objective: Find cluster centers μ_k that minimize the within cluster distance

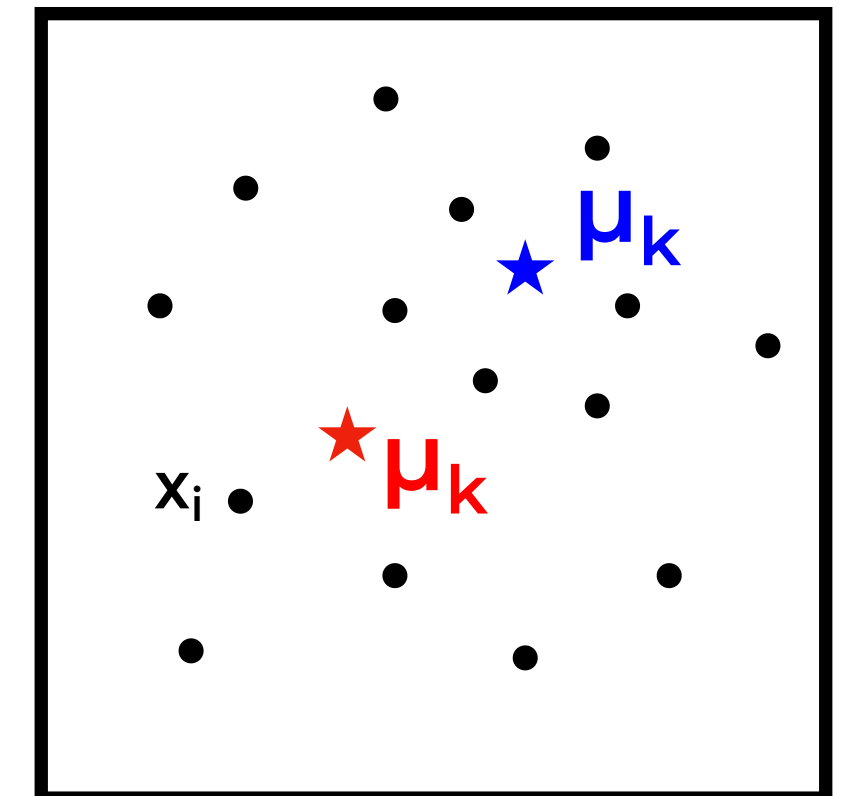
$$\text{Objective} := \sum_{i=1}^N \sum_{k=1}^K r_{ik} \|x_i - \mu_k\|^2$$

- Algorithm to minimize the objective:
 - Initialize the cluster centers
 - Until convergence:

1. Update responsibilities: r

2. Update cluster centers: $\mu_k \quad \forall k$

$$r = \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ \vdots & \vdots \\ 0 & 1 \end{bmatrix}_{N \times 2}$$



K-means clustering

- A particular clustering model (and accompanying algorithm)
 - There are K clusters. Each point belongs to a cluster. Clusters have centers: μ
- Objective: Find cluster centers μ_k that minimize the within cluster distance

$$\text{Objective} := \sum_{i=1}^N \sum_{k=1}^K r_{ik} \|x_i - \mu_k\|^2$$

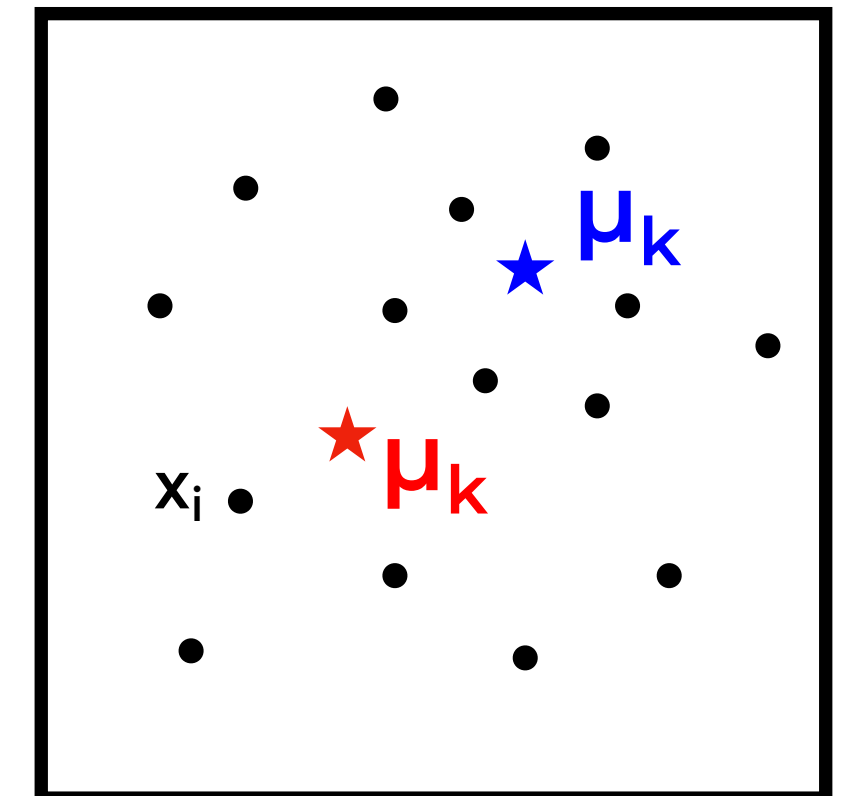
- Algorithm to minimize the objective:

- Initialize the cluster centers
- Until convergence:

1. Update responsibilities: r

2. Update cluster centers: $\mu_k \quad \forall k$

$$r = \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ \vdots & \vdots \\ 0 & 1 \end{bmatrix}_{N \times 2}$$



K-means clustering

- A particular clustering model (and accompanying algorithm)
 - There are K clusters. Each point belongs to a cluster. Clusters have centers: μ
- Objective: Find cluster centers μ_k that minimize the within cluster distance

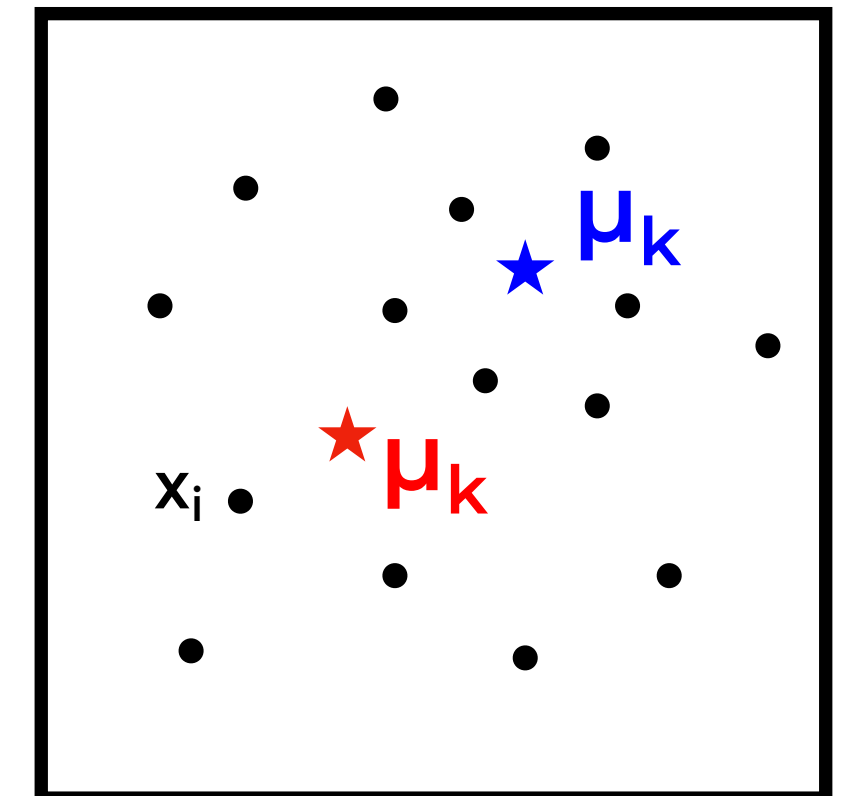
$$\text{Objective} := \sum_{i=1}^N \sum_{k=1}^K r_{ik} \|x_i - \mu_k\|^2$$

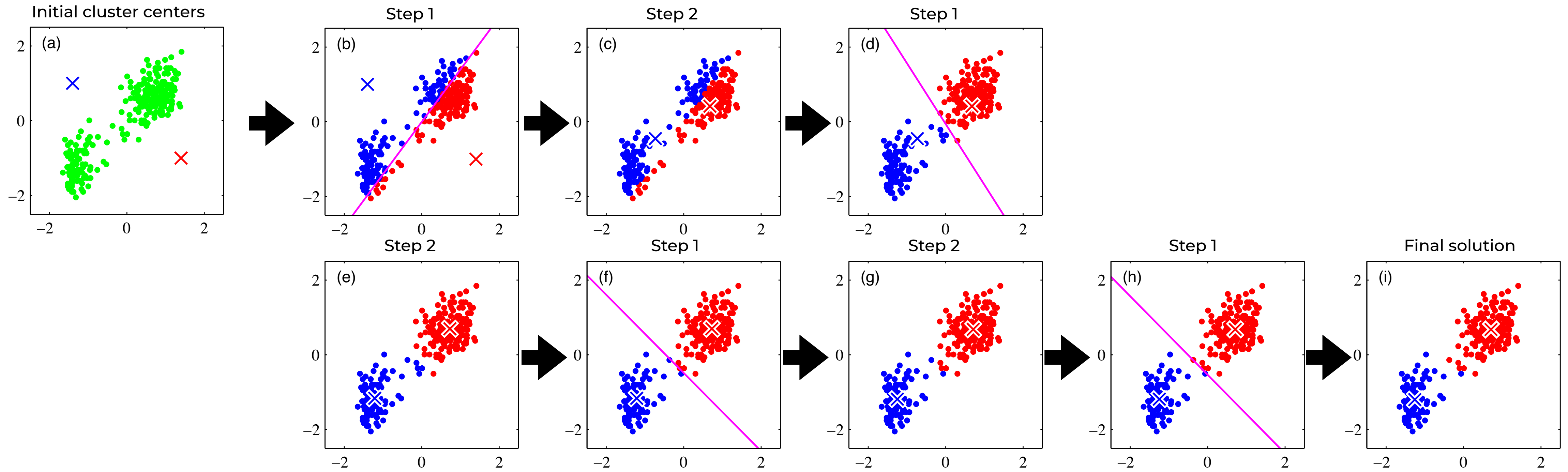
- Algorithm to minimize the objective:
 - Initialize the cluster centers
 - Until convergence:

1. Update responsibilities: r

2. Update cluster centers: $\mu_k \forall k$

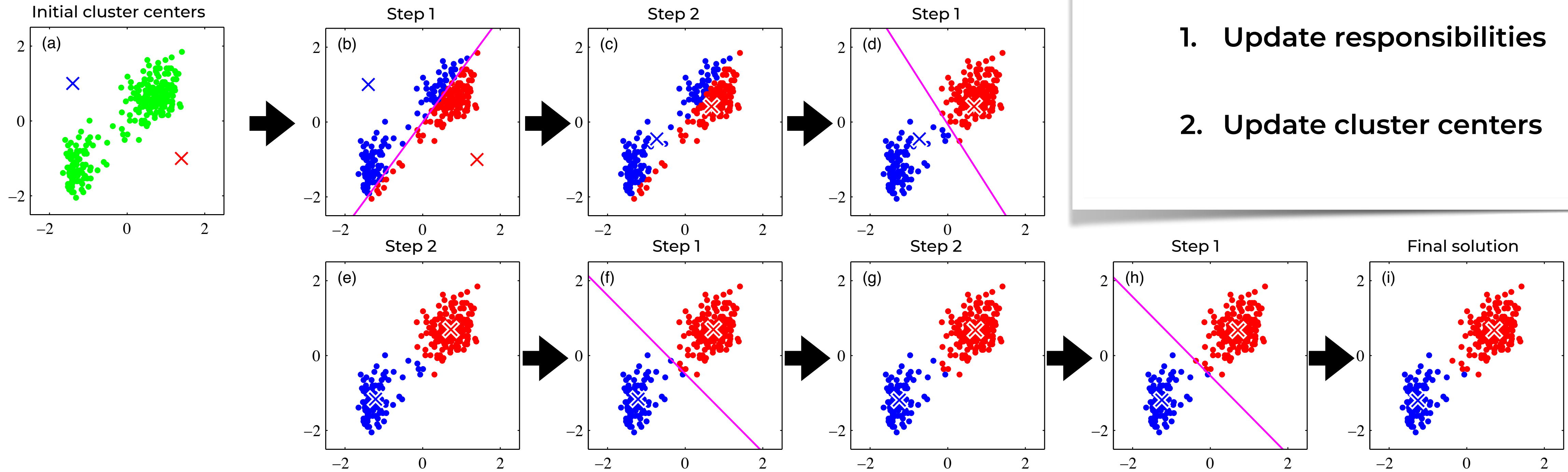
$$r = \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ \vdots & \vdots \\ 0 & 1 \end{bmatrix}_{N \times 2}$$





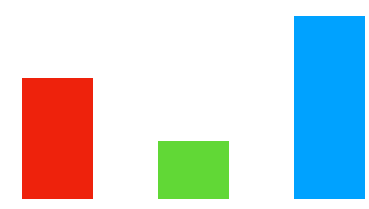
Algorithm

- Initialize the cluster centers
- Until convergence:
 1. Update responsibilities
 2. Update cluster centers



K-means on images

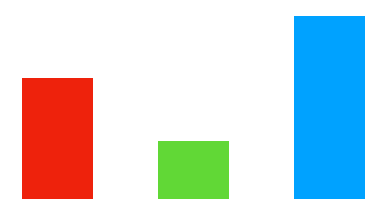
Each pixel is a datum

$$\mathbf{x}_i = (r_i, g_i, b_i)$$




K-means on images

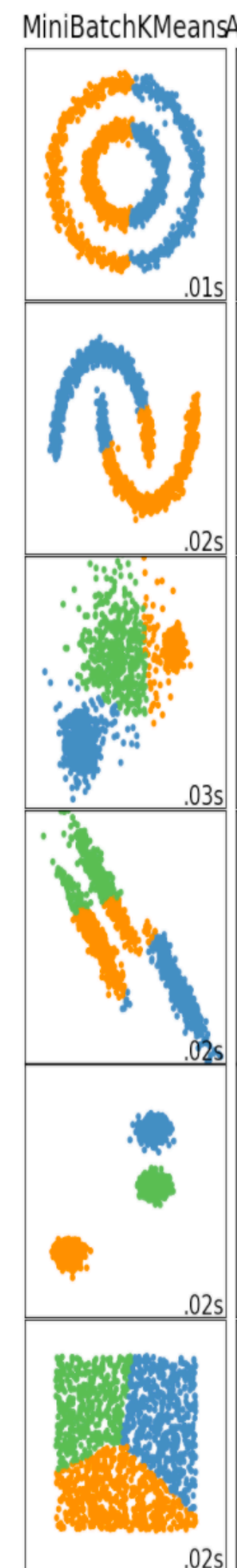
Each pixel is a datum

$$\mathbf{x}_i = (r_i, g_i, b_i)$$




- We cluster pixels by color
 - I.e., pixels of similar color will belong to the same cluster
- We re-draw the image by replacing the color of each pixel by the color of its cluster

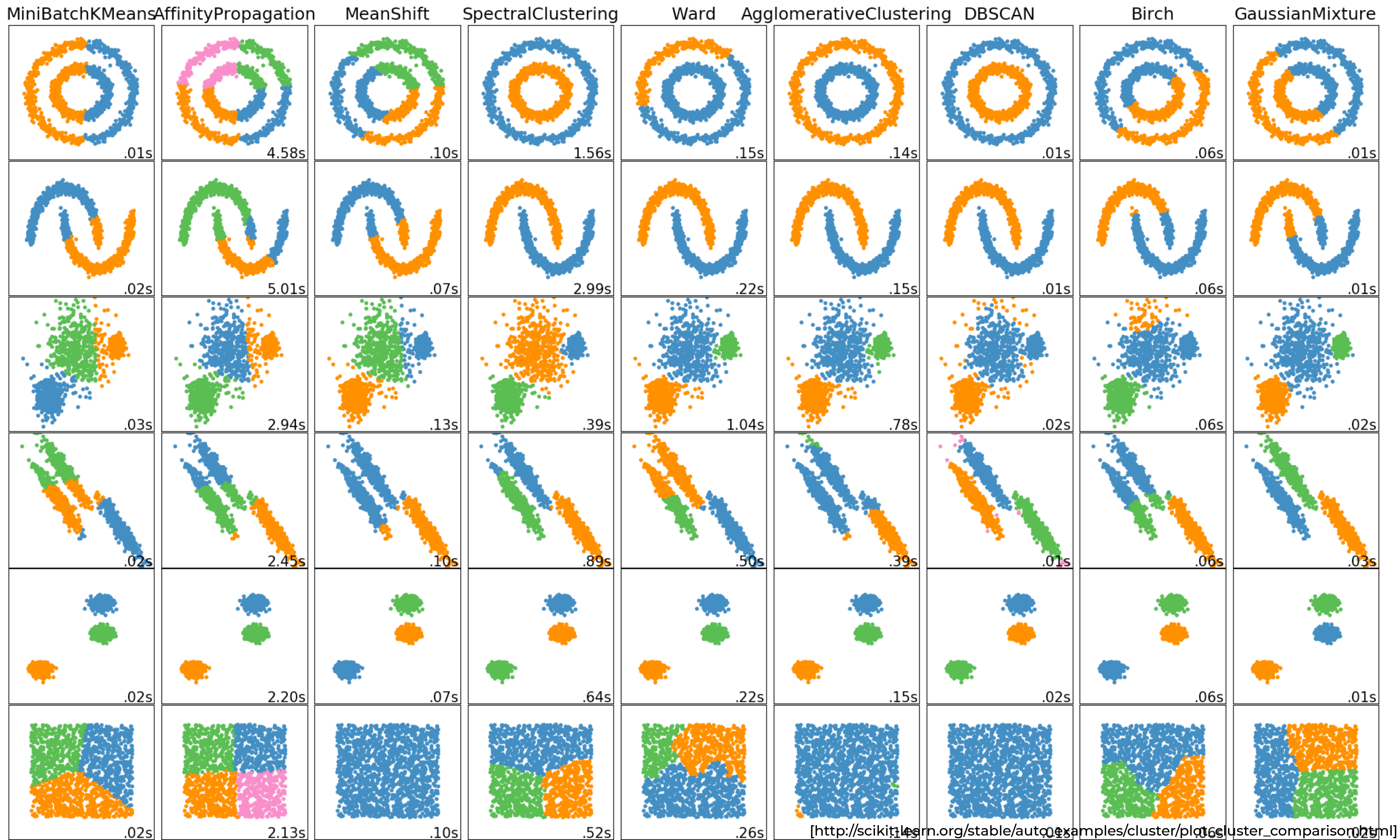
A few failure cases of k-means



- Only works with squared Euclidean distance
- What about non-continuous data?

$$\text{K-medoids Objective} := \sum_{i=1}^N \sum_{k=1}^K r_{ik} d(\mathbf{x}_i, \mu_k)$$

- Not robust to outliers
- Tends to result in relatively uniform cluster sizes
- Hard assignments



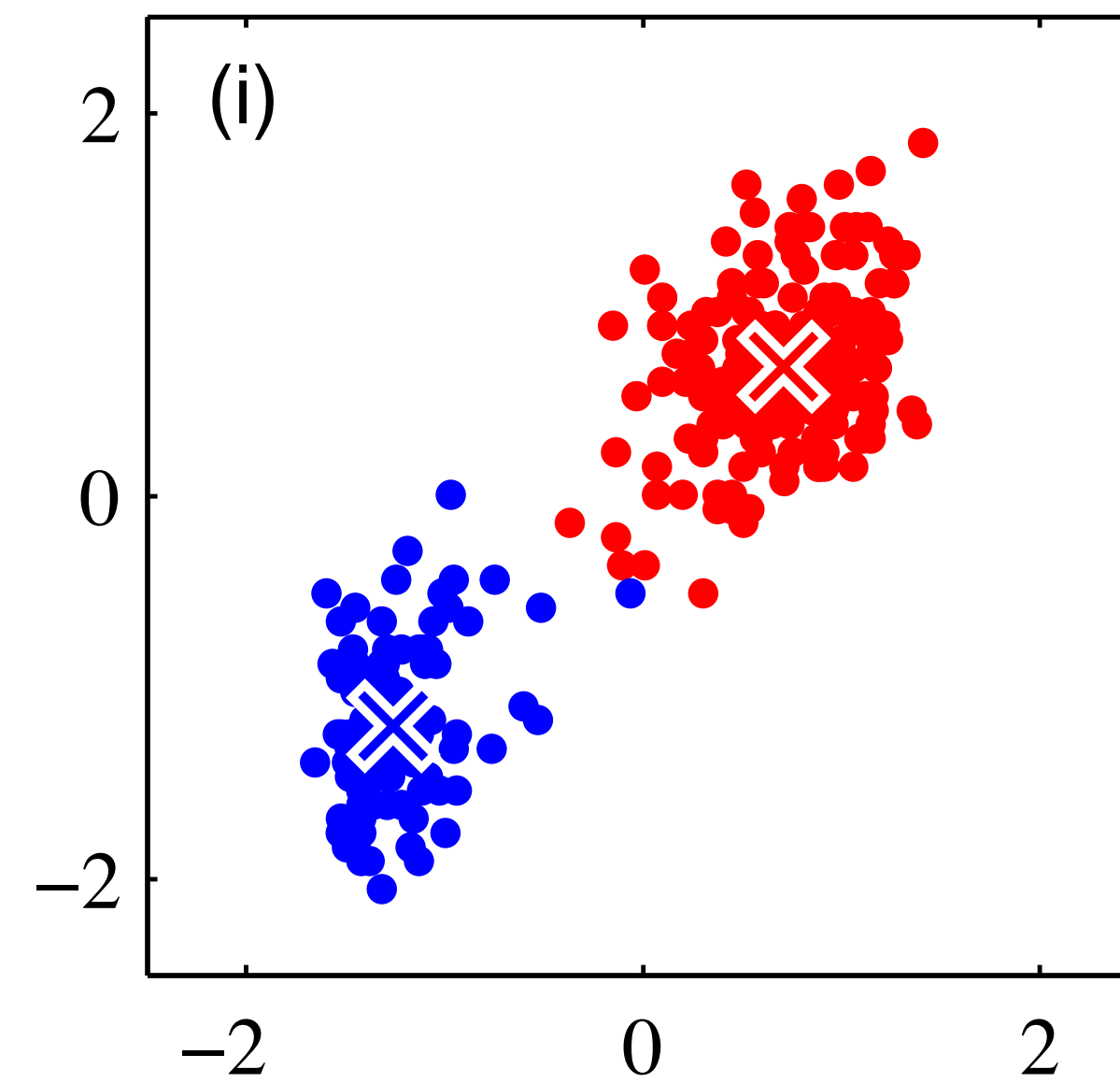
Unsupervised learning

- In supervised learning you have:
 - Clear metric (e.g., mean squared error, accuracy, etc.)
 - Procedures for comparing different models
- Unsupervised learning
 - Unclear what the right metric is, domain dependant
 - How to compare different models?
- Data dimensions can have a big impact on solution
 - You should carefully select the input

Gaussian Mixture for Probabilistic Clustering

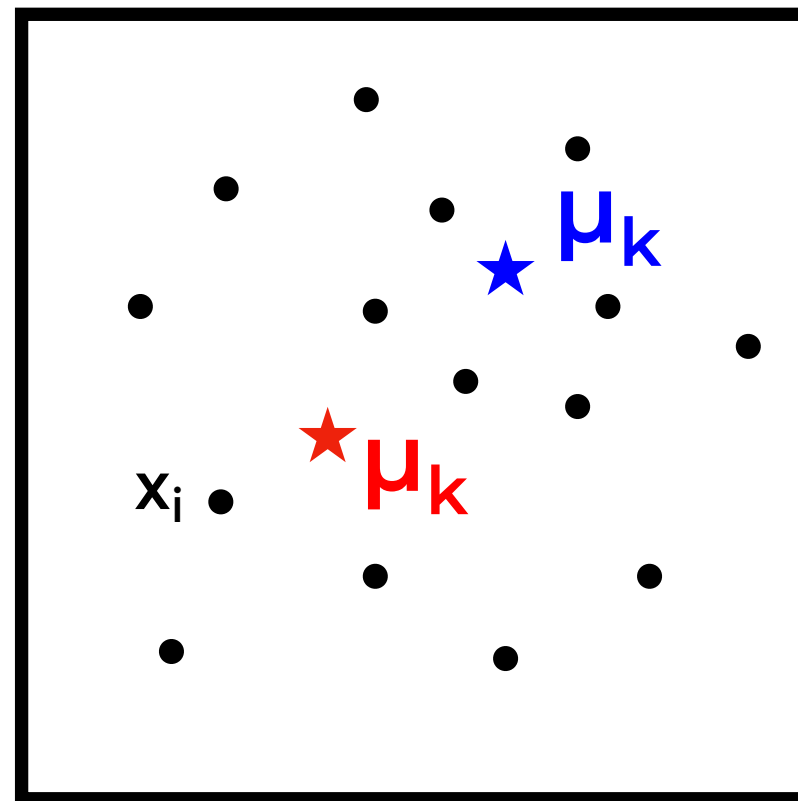
Motivation

Clustering found
using K-means



A probabilistic approach to k-means clustering

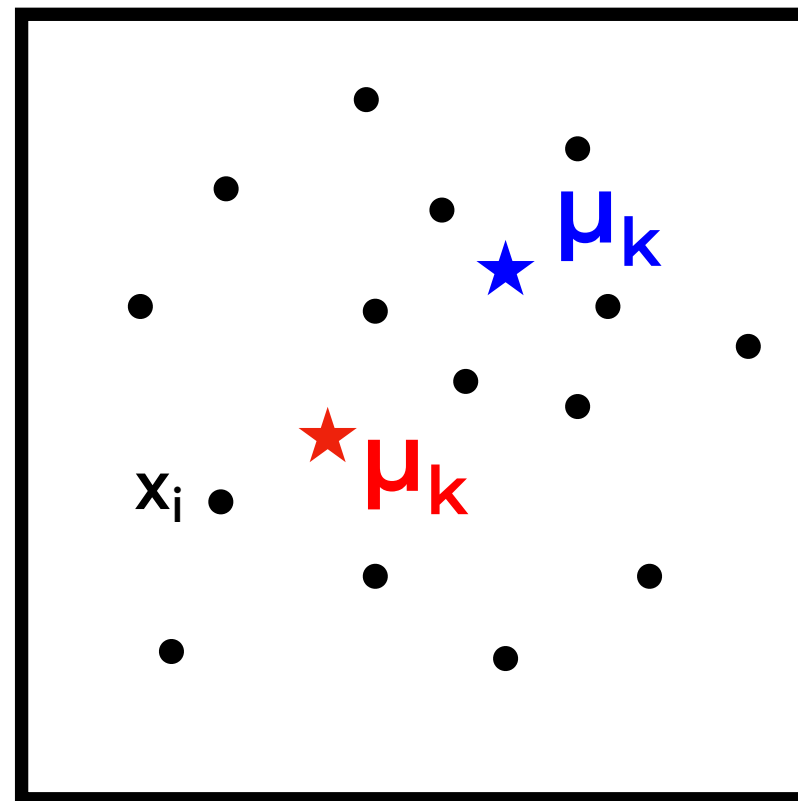
K-means Clustering



$$\text{Objective} := \sum_{i=1}^N \sum_{k=1}^K r_{ik} \| \mathbf{x}_i - \mu_k \|^2$$

A probabilistic approach to k-means clustering

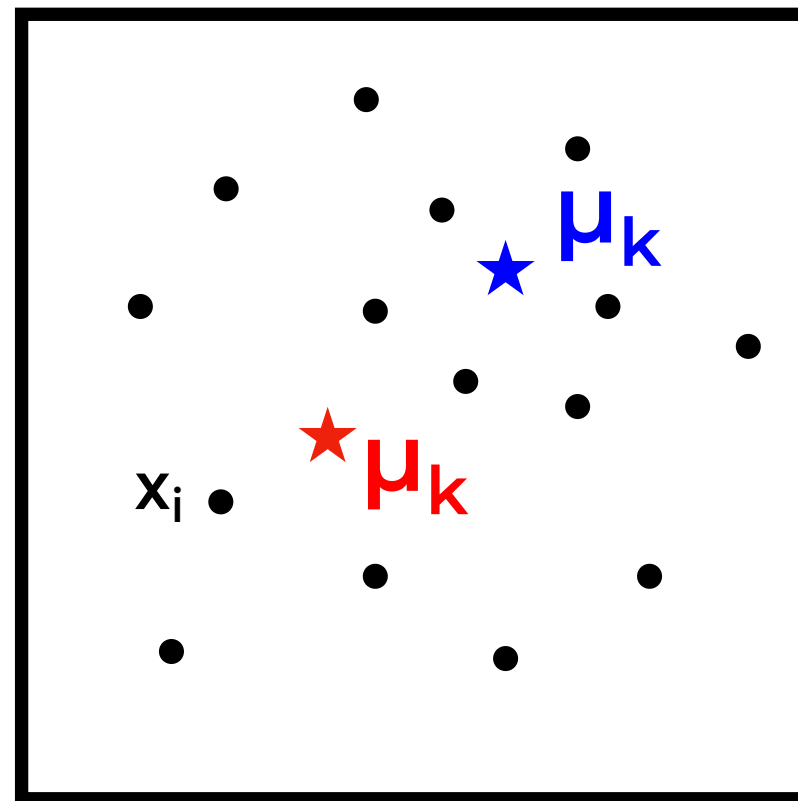
K-means Clustering



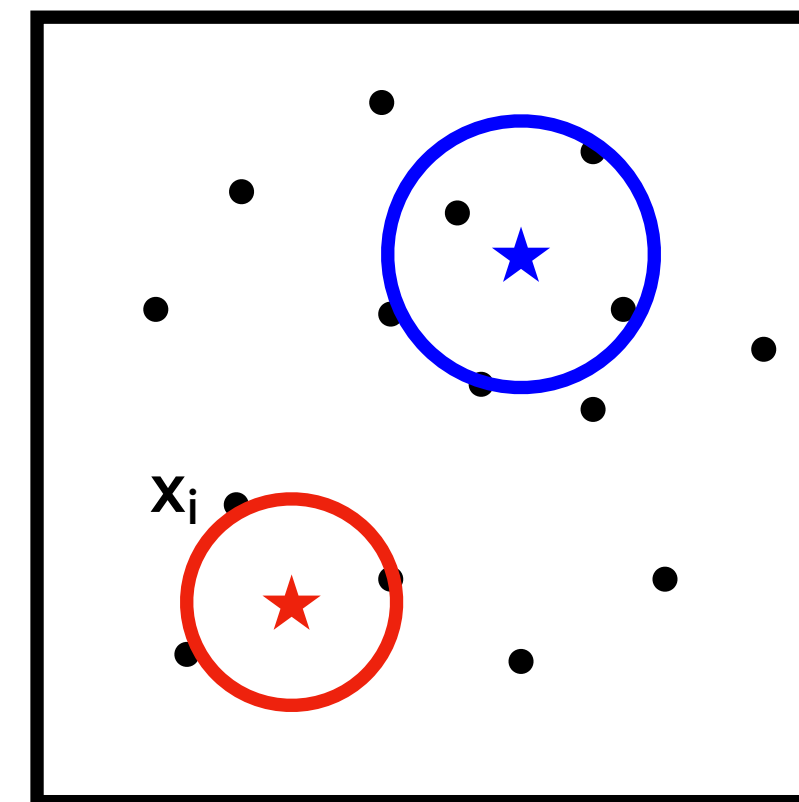
$$\text{Objective} := \sum_{i=1}^N \sum_{k=1}^K \overset{\text{Responsibility}}{r_{ik}} \left\| \mathbf{x}_i - \overset{\text{center}}{\mu_k} \right\|^2$$

A probabilistic approach to k-means clustering

K-means Clustering



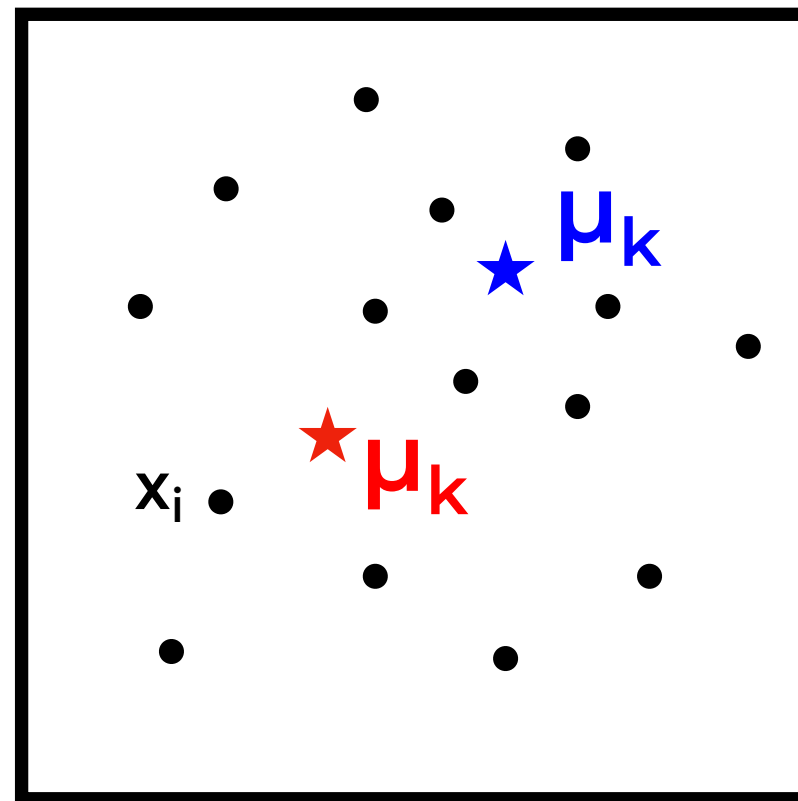
Soft K-means Clustering



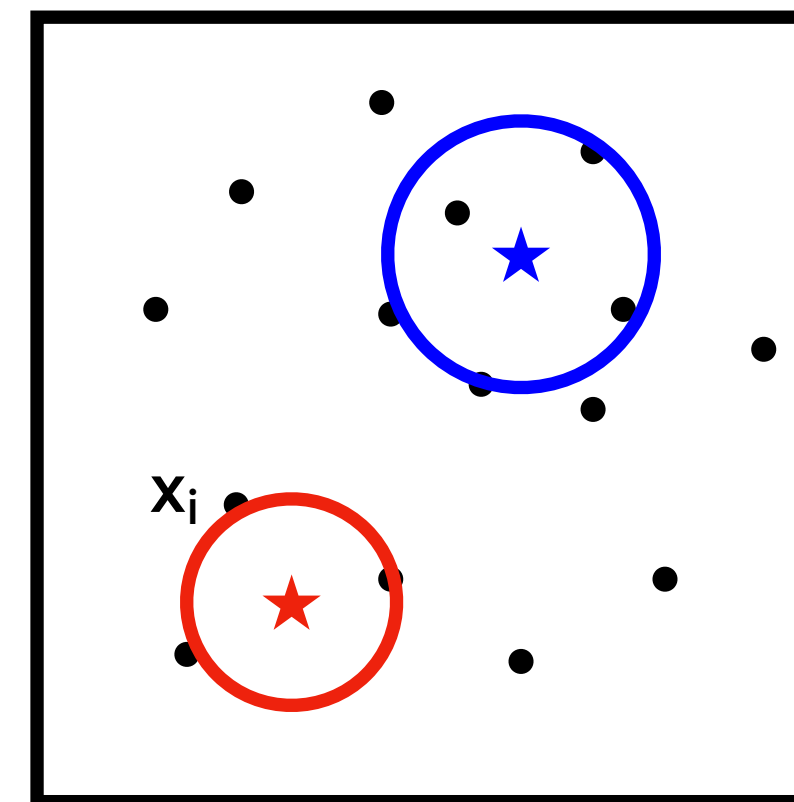
$$\text{Objective} := \sum_{i=1}^N \sum_{k=1}^K \overset{\text{Responsibility}}{r_{ik}} \left\| \mathbf{x}_i - \overset{\text{center}}{\mu_k} \right\|^2$$

A probabilistic approach to k-means clustering

K-means Clustering



Soft K-means Clustering



Objective := $\sum_{i=1}^N \sum_{k=1}^K \boxed{r_{ik}} \|\mathbf{x}_i - \boxed{\mu_k}\|^2$

Responsibility center

- Responsibilities are continuous [0, 1]
- Each cluster has a responsibility: π_k
- Each cluster models data using a Gaussian: $\mathcal{N}(\mathbf{x}_i | \mu_k, \Sigma_k)$

Soft K-means Clustering

- Each cluster gives a probability to each datum: $\mathcal{N}(\mathbf{x}_i \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

Soft K-means Clustering

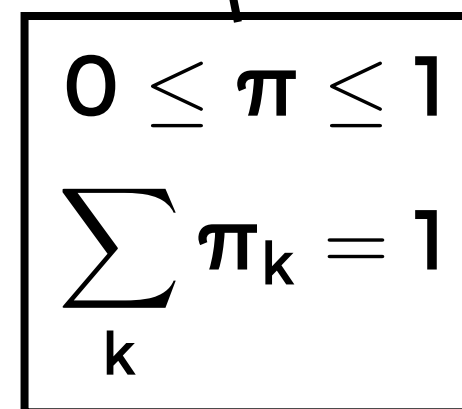
- Each cluster gives a probability to each datum: $\mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
- We can write the complete model as:

$$\mathbf{P}(\mathbf{x}_i | \text{parameters}) = \sum_{k=1}^K \pi_k \mathbf{P}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Soft K-means Clustering

- Each cluster gives a probability to each datum: $\mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
- We can write the complete model as:

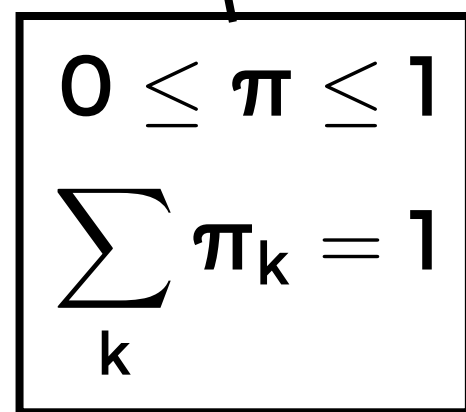
$$\mathbf{P}(\mathbf{x}_i | \text{parameters}) = \sum_{k=1}^K \pi_k \mathbf{P}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$


$$\begin{array}{l} 0 \leq \pi \leq 1 \\ \sum_k \pi_k = 1 \end{array}$$

Soft K-means Clustering

- Each cluster gives a probability to each datum: $\mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
- We can write the complete model as:

$$\mathbf{P}(\mathbf{x}_i | \text{parameters}) = \sum_{k=1}^K \pi_k \mathbf{P}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

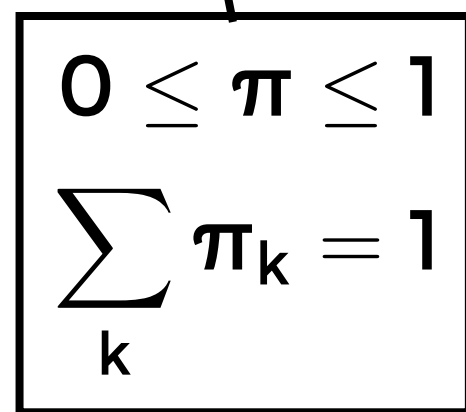

$$\begin{array}{l} 0 \leq \pi \leq 1 \\ \sum_k \pi_k = 1 \end{array}$$

$$\mathbf{P}(\mathbf{x}_i | \text{parameters}) = \sum_{k=1}^K \mathbf{P}(\mathbf{z}_i = \mathbf{k}) \mathbf{P}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Soft K-means Clustering

- Each cluster gives a probability to each datum: $\mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
- We can write the complete model as:

$$\mathbf{P}(\mathbf{x}_i | \text{parameters}) = \sum_{k=1}^K \pi_k \mathbf{P}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$


$$\begin{array}{l} 0 \leq \pi \leq 1 \\ \sum_k \pi_k = 1 \end{array}$$

$$\mathbf{P}(\mathbf{x}_i | \text{parameters}) = \sum_{k=1}^K \mathbf{P}(z_i = k) \mathbf{P}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$\text{Posterior: } \mathbf{P}(z_i = k | \mathbf{x}_i) \propto \mathbf{P}(z_i = k) \mathbf{P}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

Soft K-means Clustering

- Each cluster gives a probability to each datum: $\mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
- We can write the complete model as:
- Model is known as a **Gaussian Mixture Model**

$$P(\mathbf{x}_i | \text{parameters}) = \sum_{k=1}^K \pi_k P(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$0 \leq \pi_k \leq 1$$
$$\sum_k \pi_k = 1$$

$$P(\mathbf{x}_i | \text{parameters}) = \sum_{k=1}^K P(\mathbf{z}_i = \mathbf{k}) P(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Posterior: $P(\mathbf{z}_i = \mathbf{k} | \mathbf{x}_i) \propto P(\mathbf{z}_i = \mathbf{k}) P(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

Soft K-means Clustering

- Each cluster gives a probability to each datum: $\mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
- We can write the complete model as:

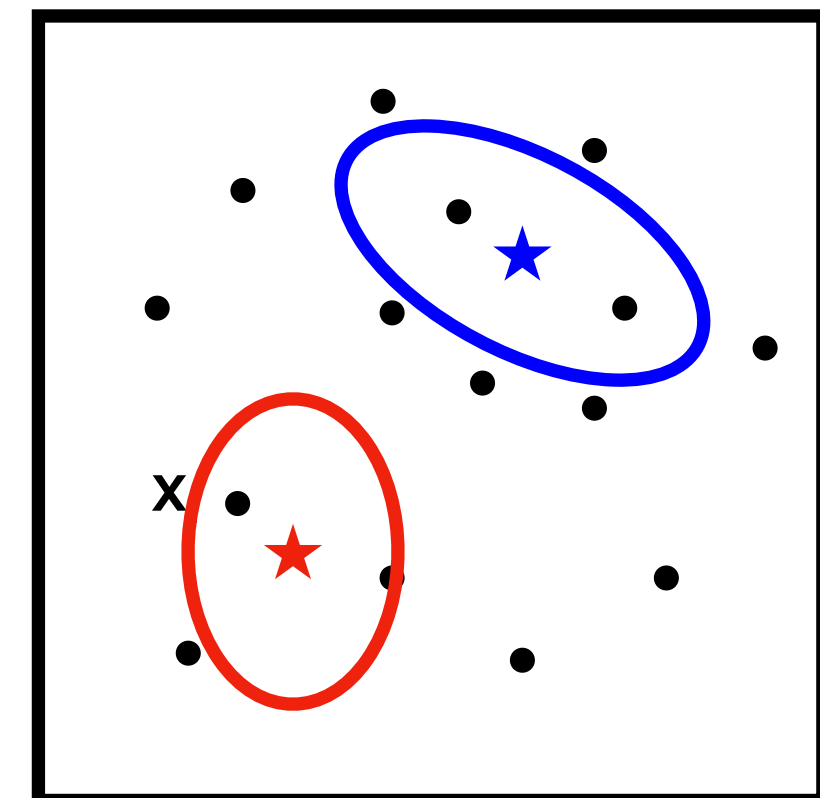
$$\mathbf{P}(\mathbf{x}_i | \text{parameters}) = \sum_{k=1}^K \pi_k \mathbf{P}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$\begin{aligned} 0 &\leq \pi_k \leq 1 \\ \sum_k \pi_k &= 1 \end{aligned}$$

$$\mathbf{P}(\mathbf{x}_i | \text{parameters}) = \sum_{k=1}^K \mathbf{P}(z_i = k) \mathbf{P}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$\text{Posterior: } \mathbf{P}(z_i = k | \mathbf{x}_i) \propto \mathbf{P}(z_i = k) \mathbf{P}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

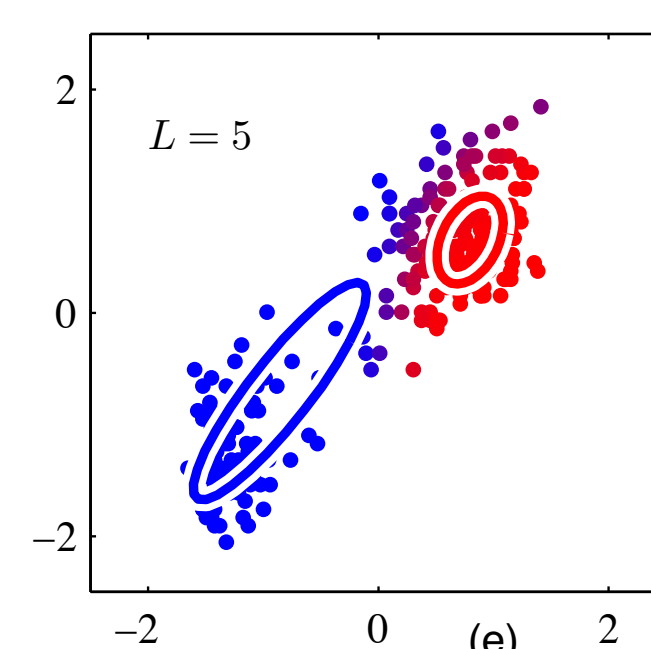
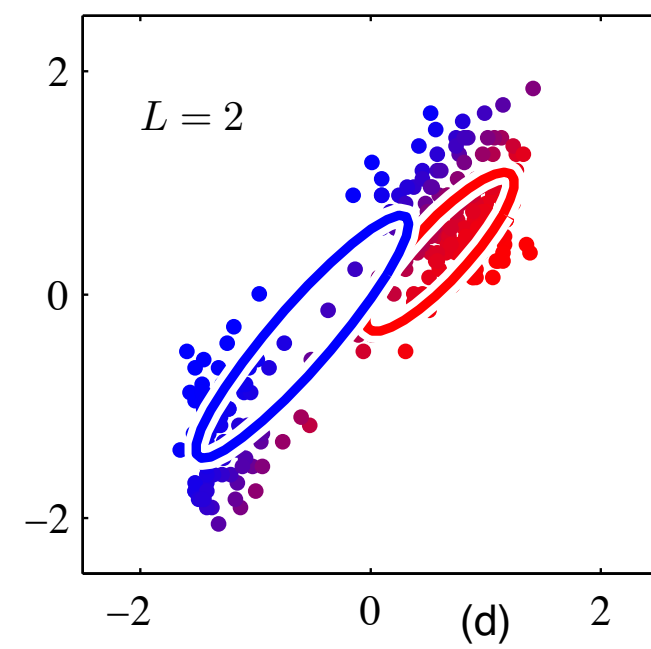
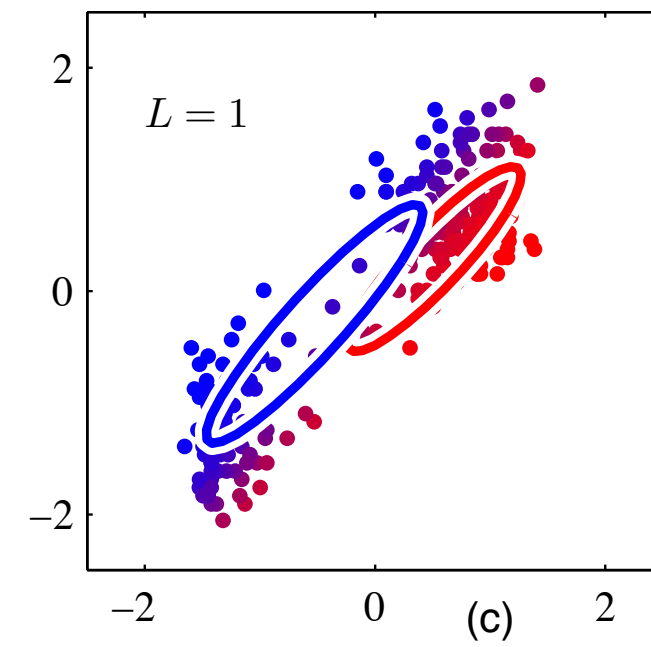
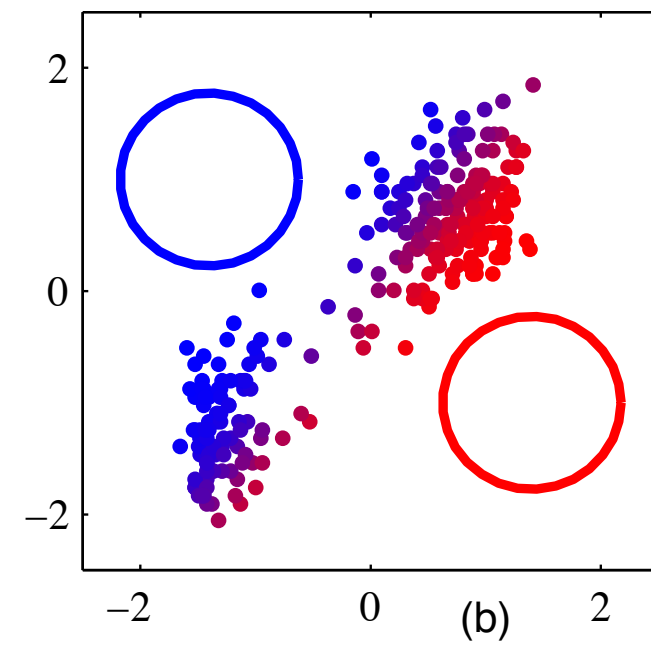
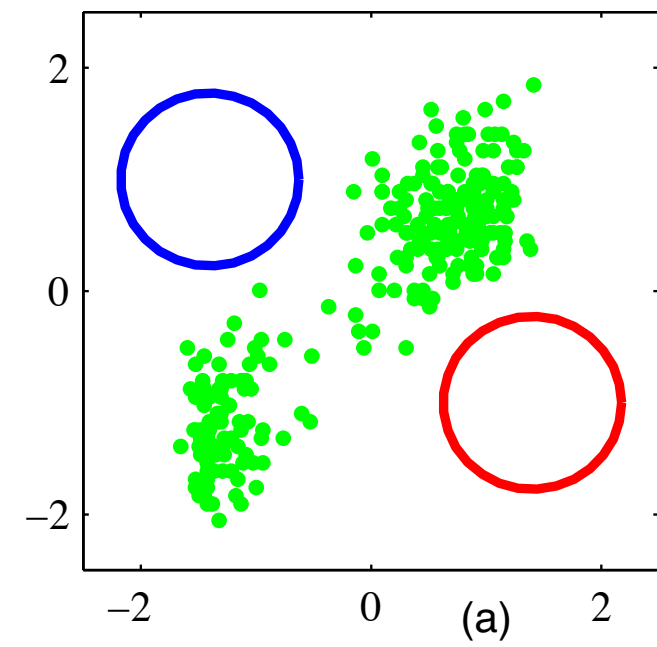
- Model is known as a **Gaussian Mixture Model**
- Advantages over k-means:
 - Soft Clustering
 - Clusters don't have to be spherical



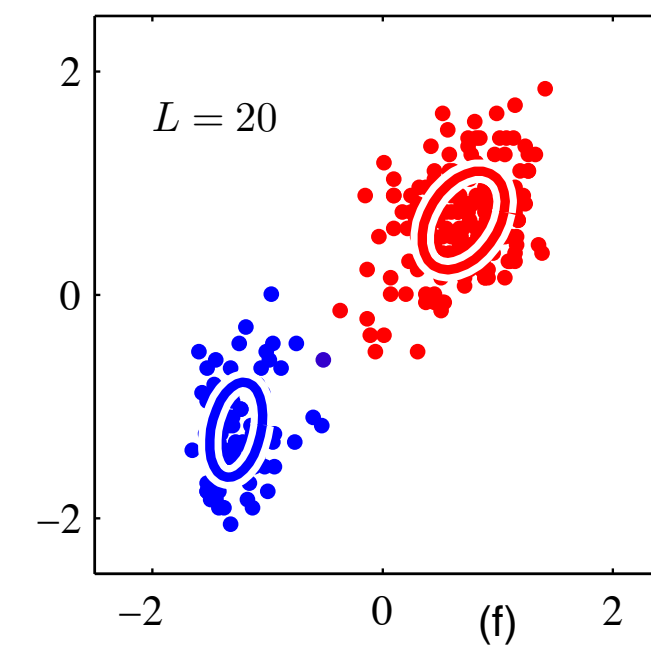
A probabilistic approach to k-means clustering

- Estimation of the parameters: Maximum likelihood estimate
 - Could do it jointly (“a la” neural network)
 - Often in two separate steps (similar as the non-probabilistic version)
 - This leads to the Expectation-Maximization (EM) algorithm

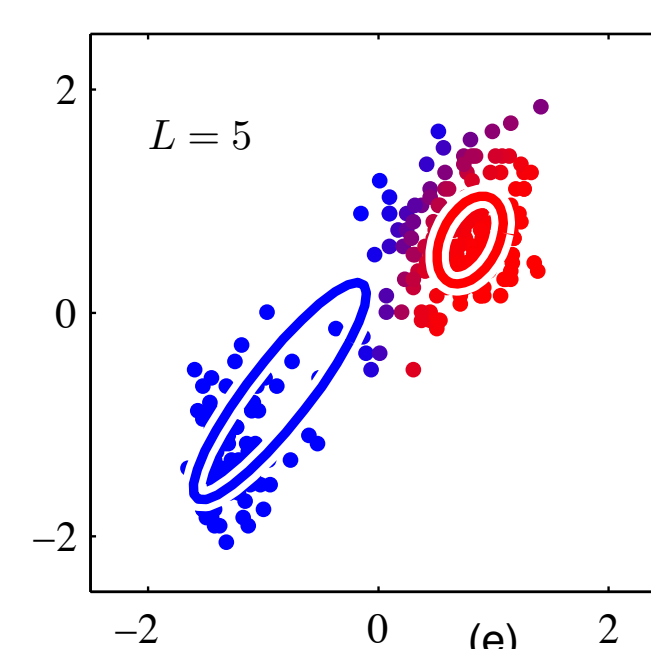
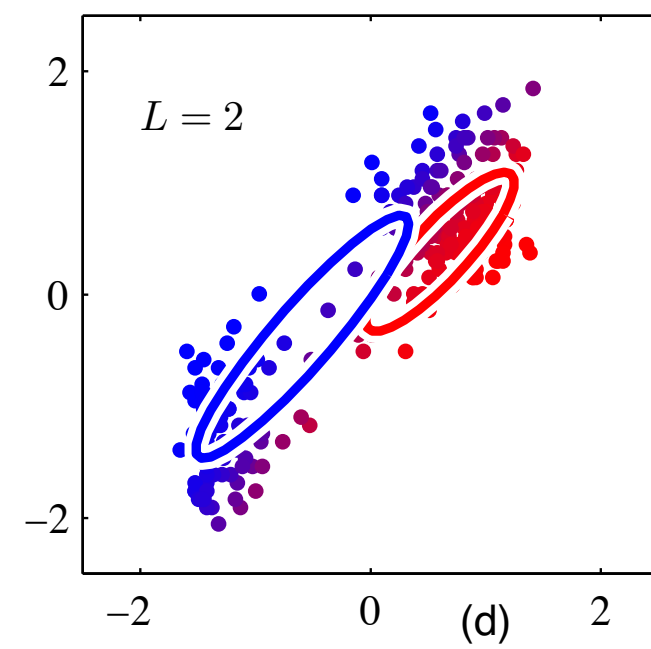
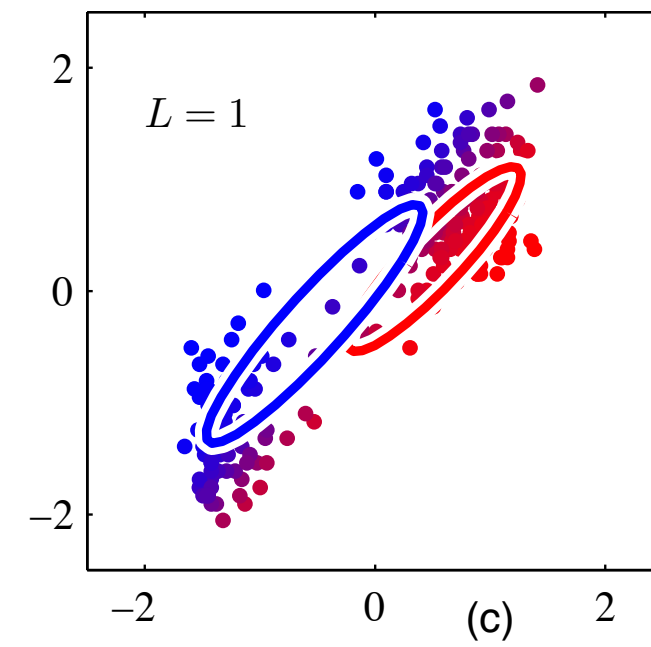
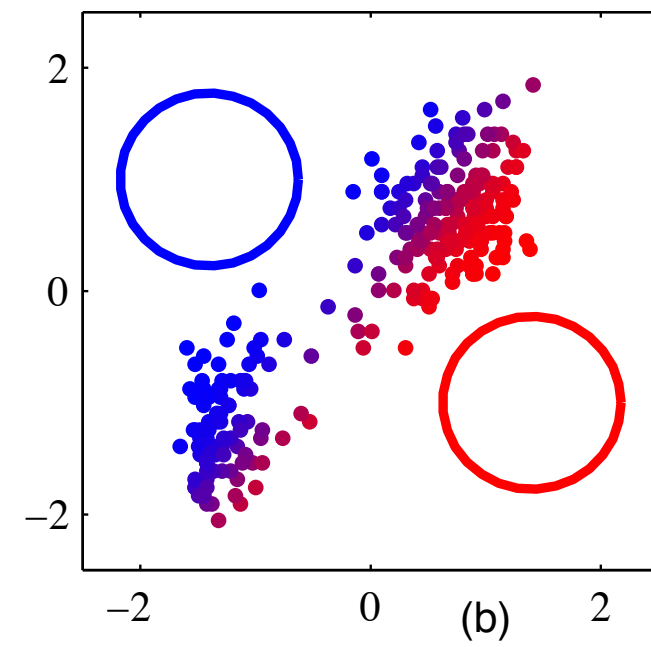
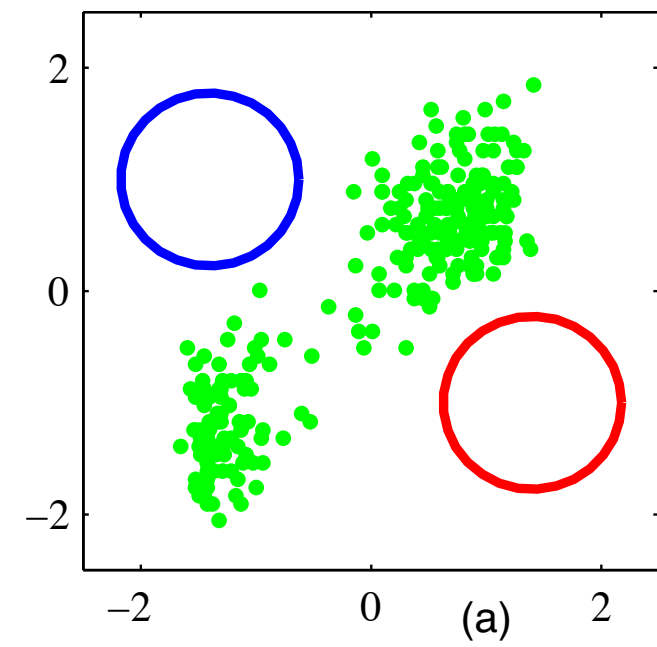
Initial cluster centers



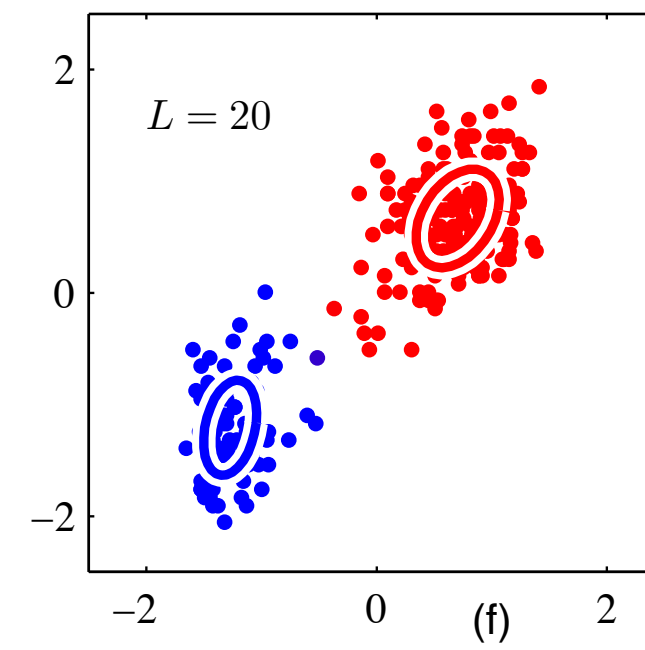
Final solution



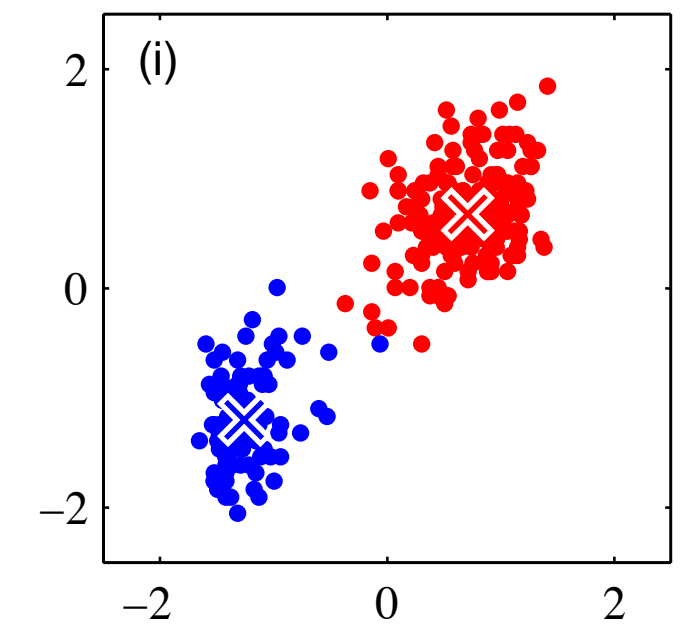
Initial cluster centers



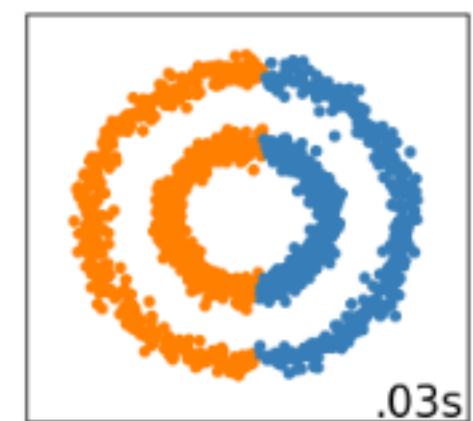
Final solution



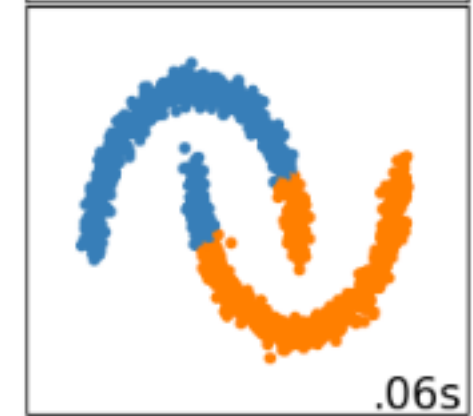
k-means
Final solution



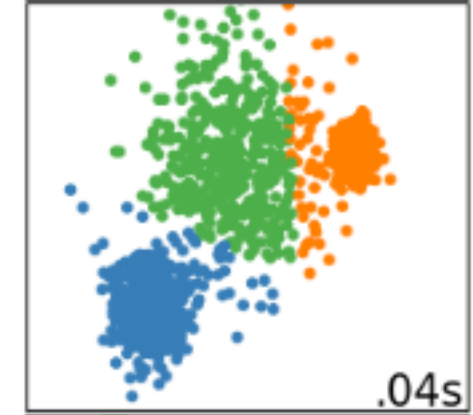
K-means



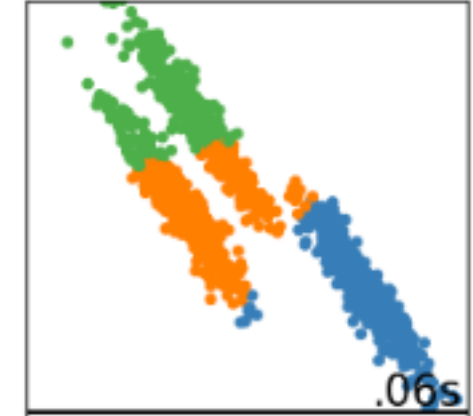
Similar



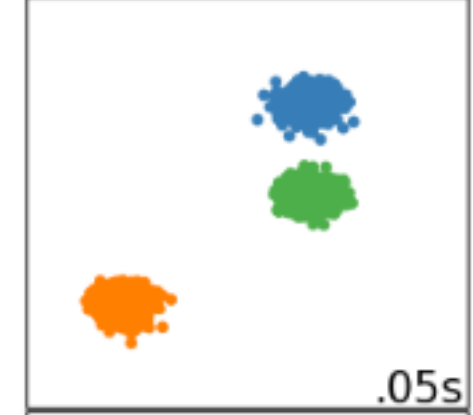
Similar



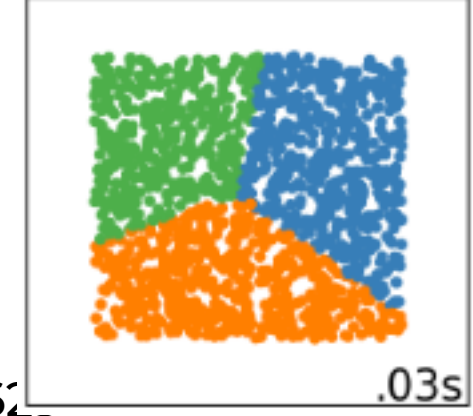
GMM better



GMM better

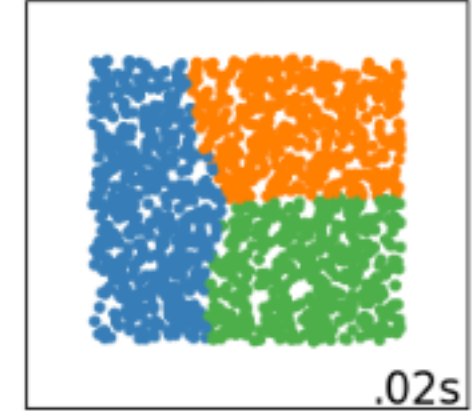
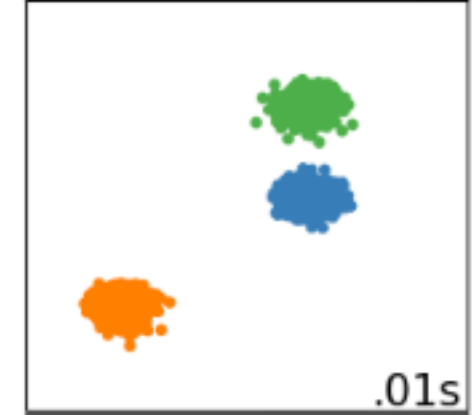
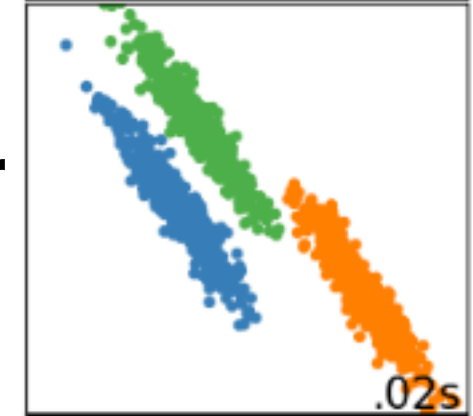
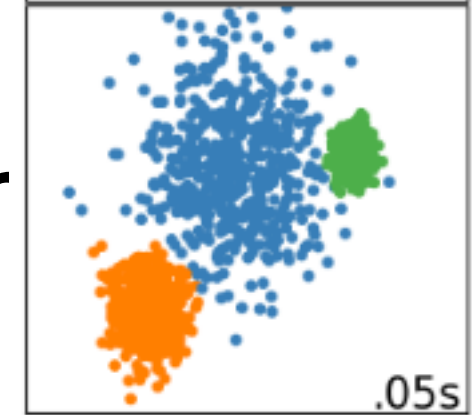
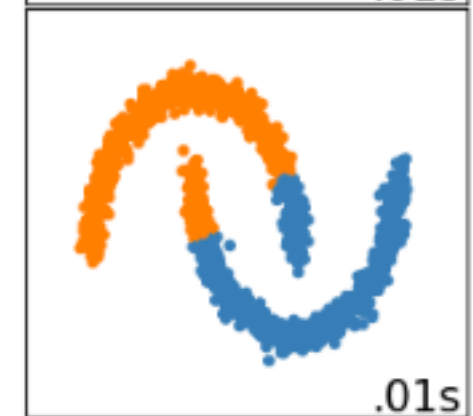
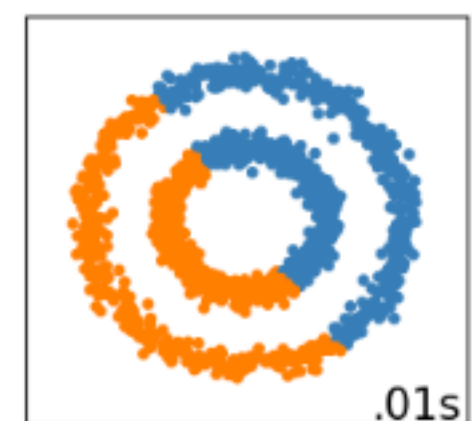


Similar

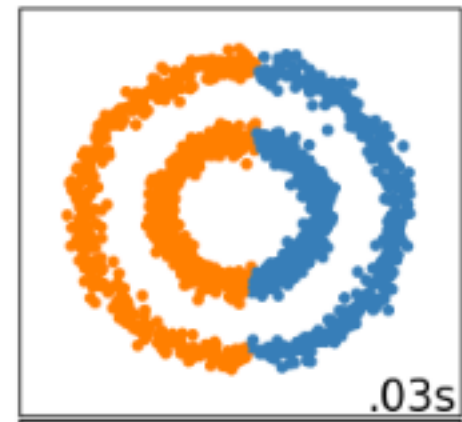


Similar

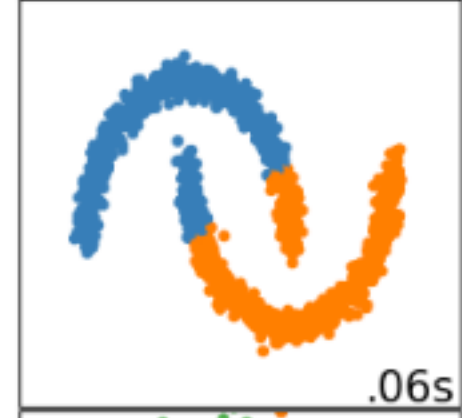
GMMs



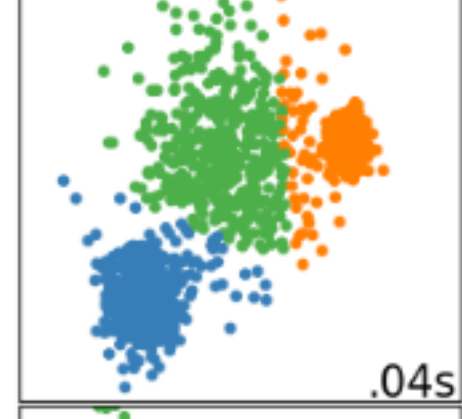
K-means



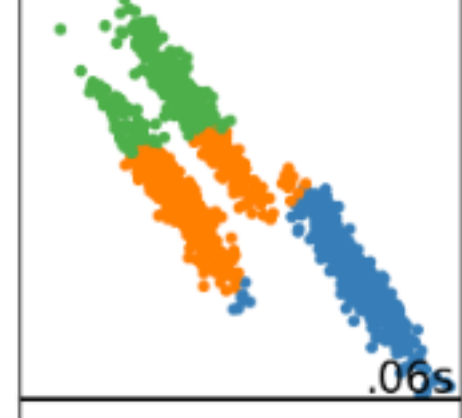
Similar



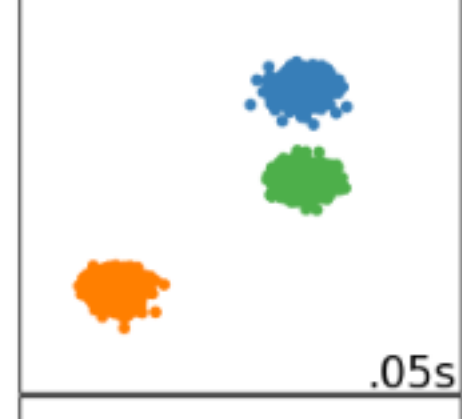
Similar



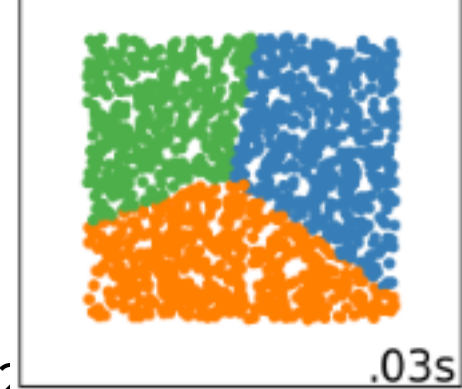
GMM better



GMM better

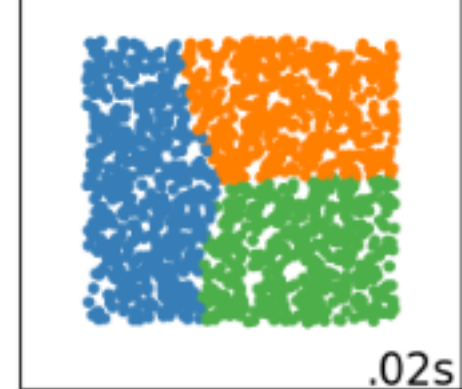
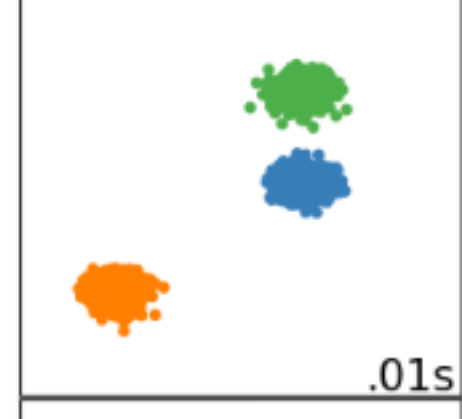
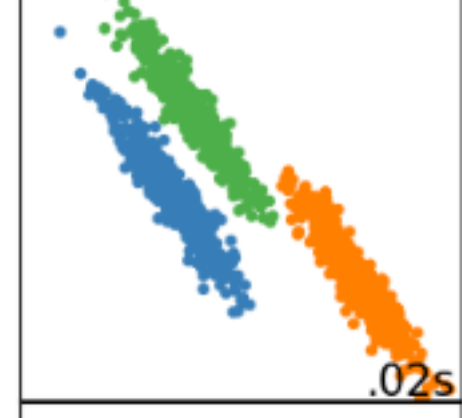
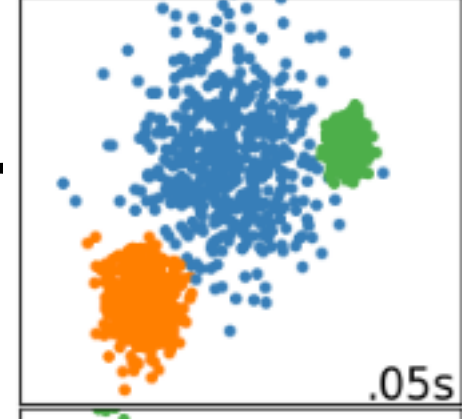
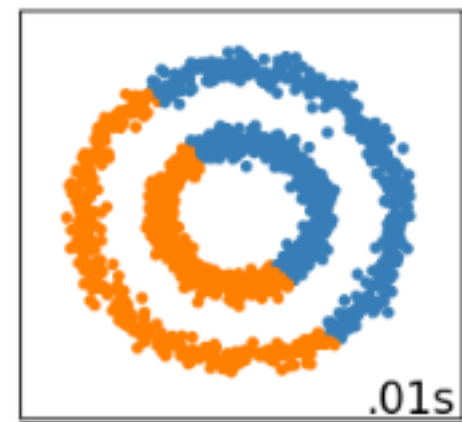


Similar



Similar

GMMs



Comparing K-means to GMMs

- GMMs learns covariance matrix
 - Per cluster variance
 - Covariance terms
- GMMs has many more parameters
 - Covariance matrix (MxM)

Unsupervised Learning Beyond Clustering

Different tasks

- Finding patterns

- Clustering

$$f : X \rightarrow \{1, 2, \dots, K\} \text{ (K clusters)}$$

- Dimensionality reduction

$$f : X^p \rightarrow X^k, k \ll p$$

- Density modelling

$$f : X \rightarrow [0, 1]$$

- ...

Different tasks

- Finding patterns

- Clustering

$$f : X \rightarrow \{1, 2, \dots, K\} \text{ (K clusters)}$$

- Dimensionality reduction

$$f : X^p \rightarrow X^k, k \ll p$$

- Density modelling

$$f : X \rightarrow [0, 1]$$

- ...

Autoencoders \longrightarrow

Autoencoders

- A type of neural network for dimensionality reduction
 - Non-linear PCA

Autoencoders

- A type of neural network for dimensionality reduction
 - Non-linear PCA
 - Intuition: let's learn to copy the data

$$\mathbf{x} = \mathbf{f}(\mathbf{x})$$

Autoencoders

- A type of neural network for dimensionality reduction
 - Non-linear PCA
 - Intuition: let's learn to copy the data

$$\mathbf{x} = \mathbf{f}(\mathbf{x})$$

- We force a “bottleneck”

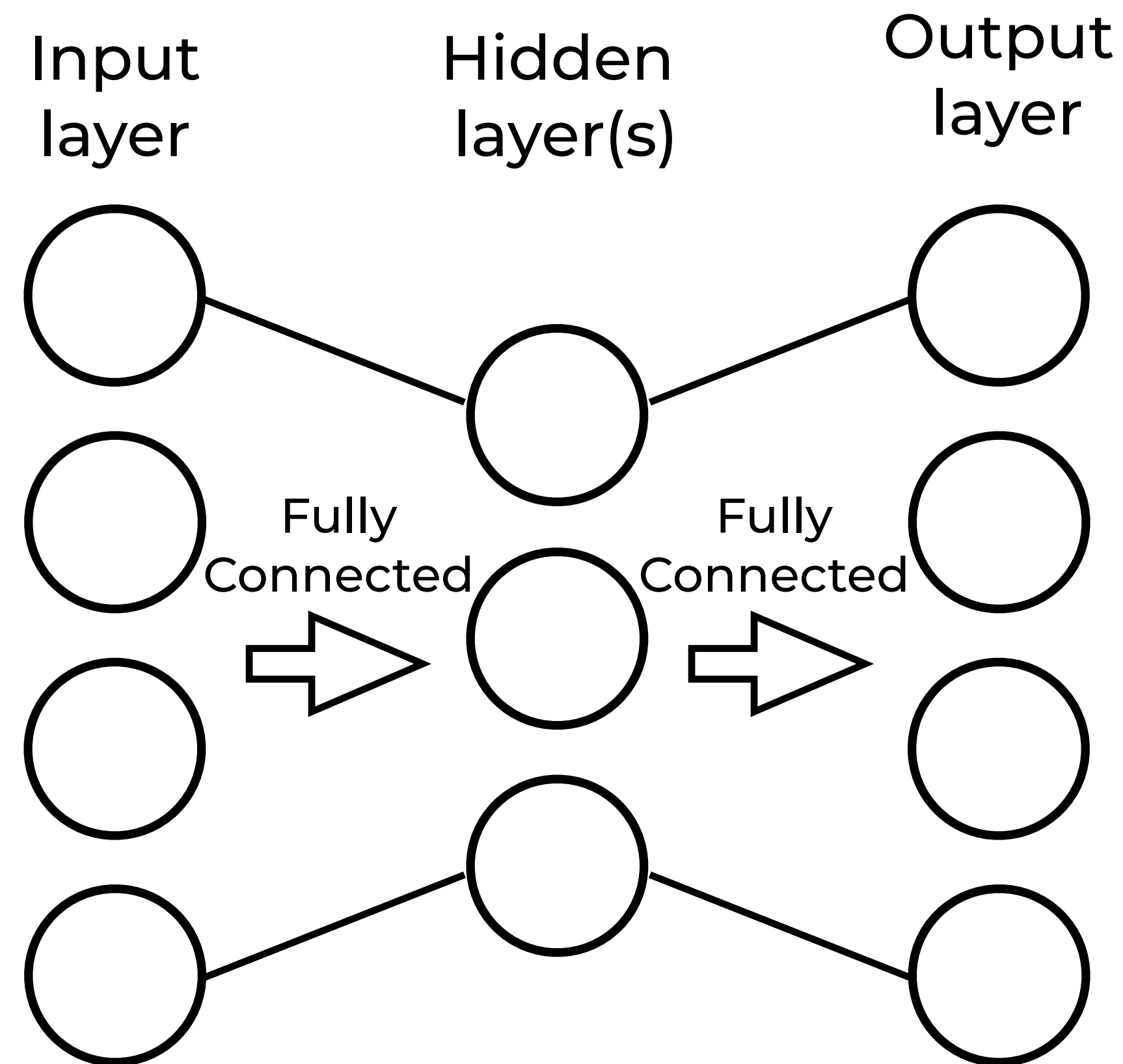
$$\mathbf{z} = \mathbf{f}_1(\mathbf{x})$$

$$\mathbf{x} = \mathbf{f}_2(\mathbf{z})$$

$$\dim(\mathbf{z}) < \dim(\mathbf{x})$$

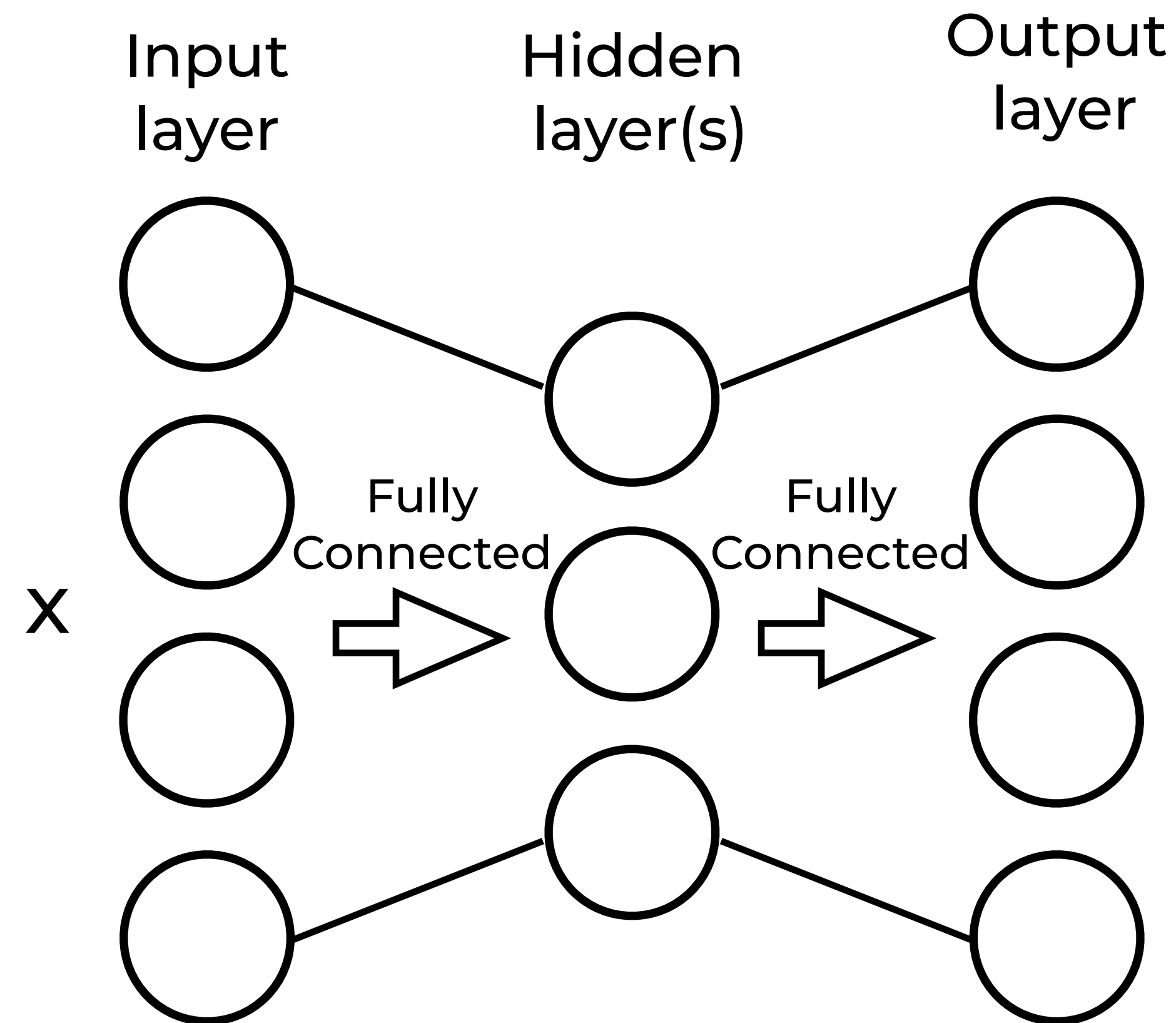
Autoencoders

- A neural network architecture for unsupervised learning



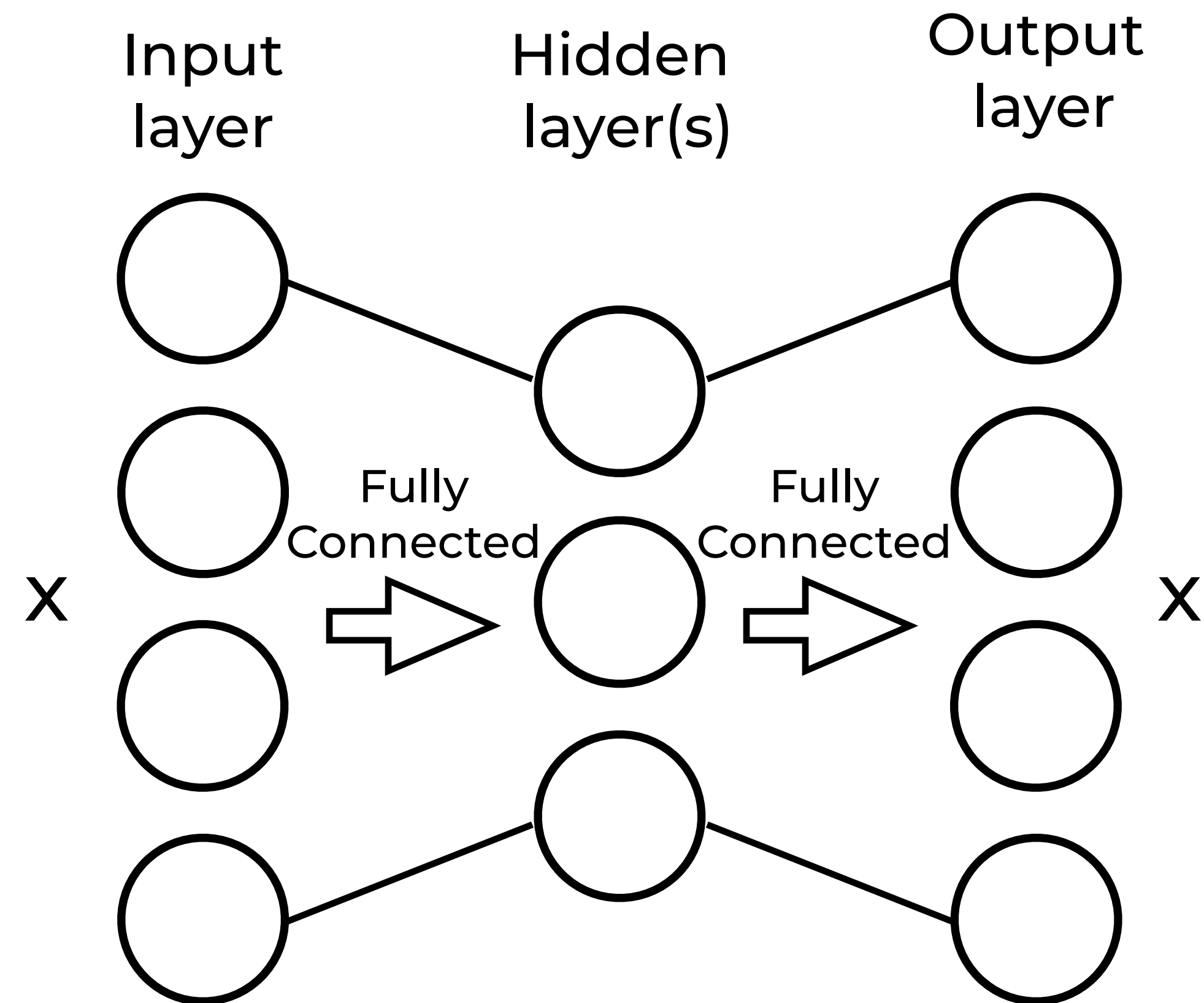
Autoencoders

- A neural network architecture for unsupervised learning



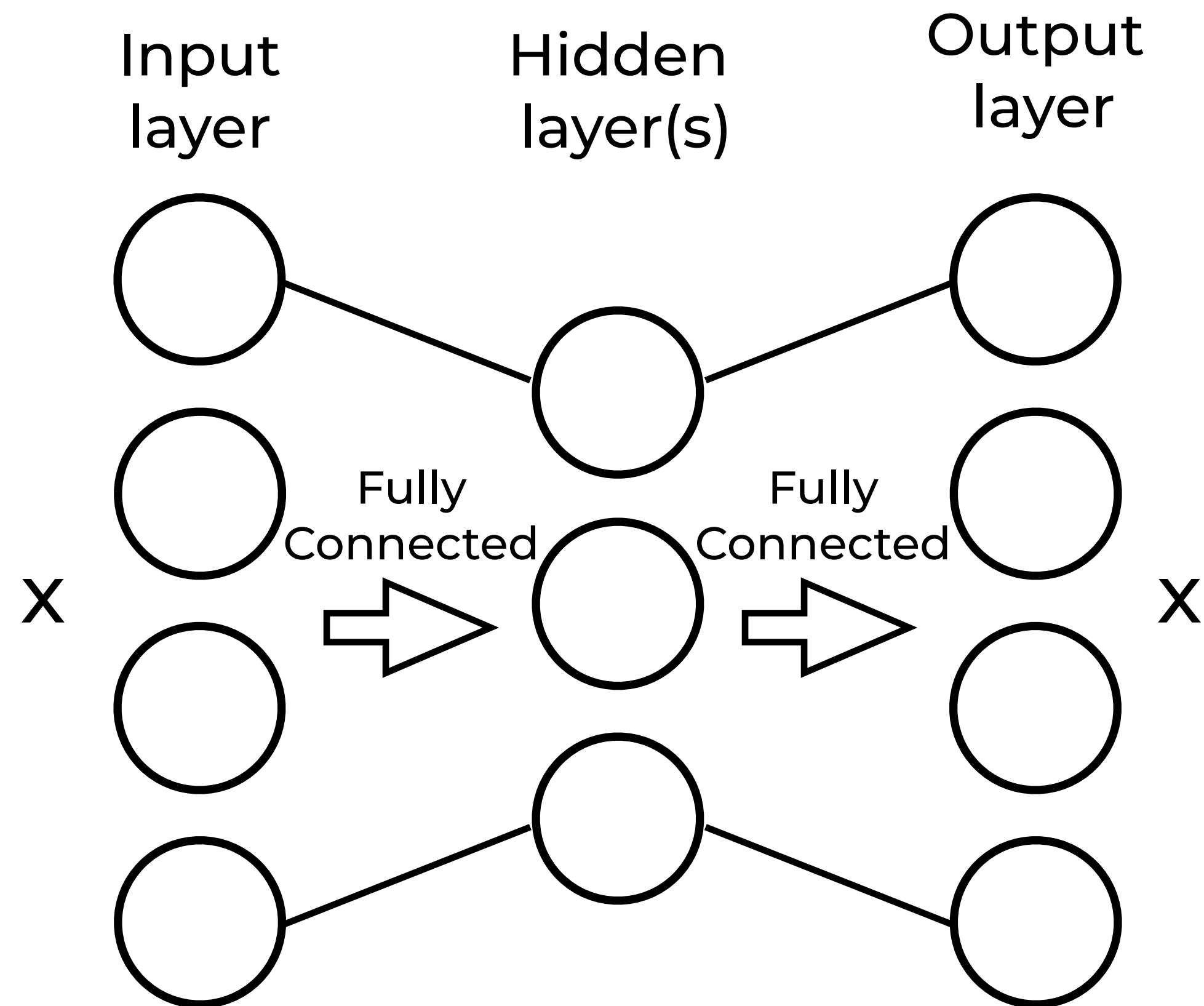
Autoencoders

- A neural network architecture for unsupervised learning



Autoencoders

- A neural network architecture for unsupervised learning



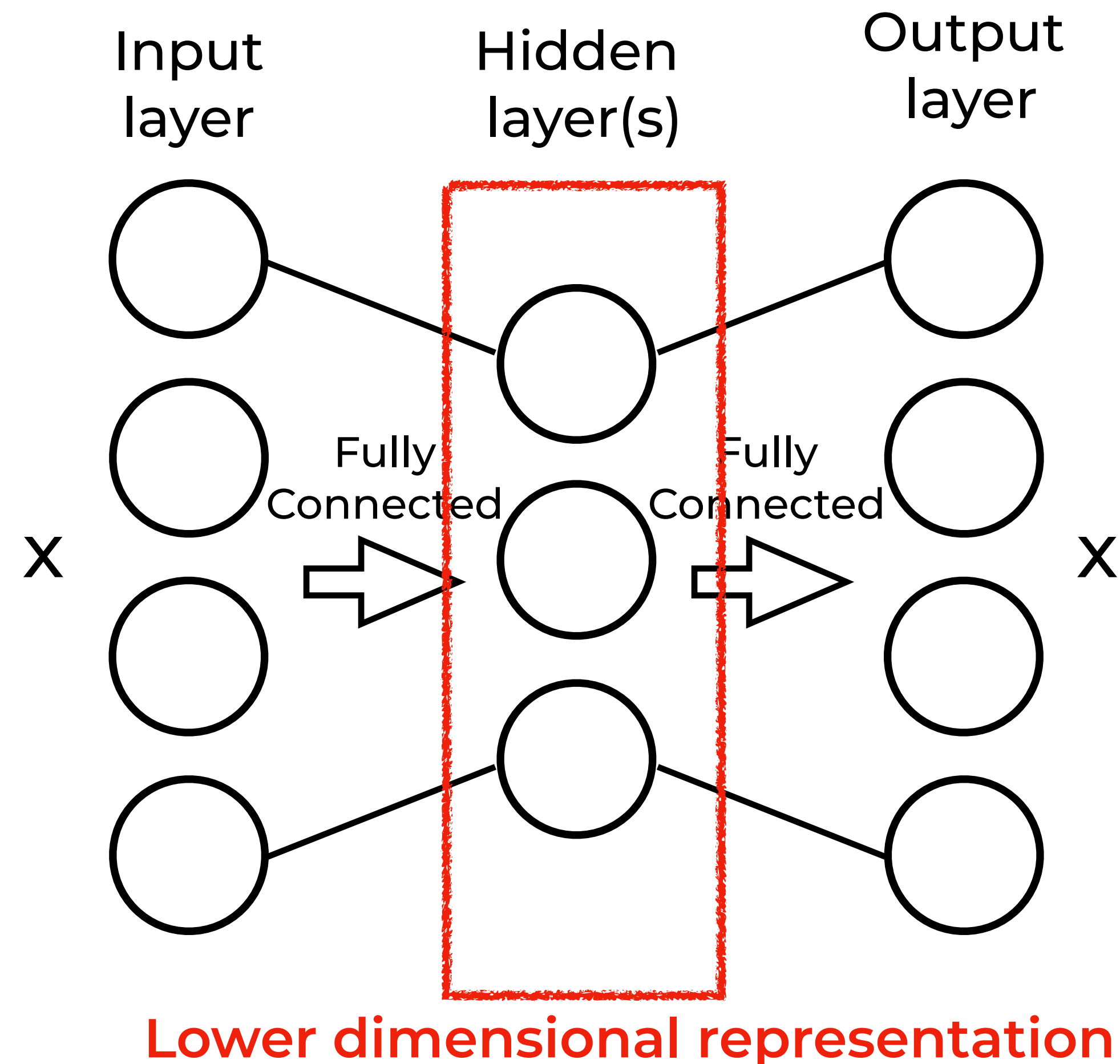
Objective:

How well the network predicts X ?

$$\begin{aligned} \text{Loss} &:= \sum_{i=1}^N (x_i - \hat{x}_i)^2 \\ &= \sum_{i=1}^N (x_i - f_2(f_1(x)))^2 \end{aligned}$$

Autoencoders

- A neural network architecture for unsupervised learning



Objective:
How well the network predicts X ?

$$\begin{aligned} \text{Loss} &:= \sum_{i=1}^N (x_i - \hat{x}_i)^2 \\ &= \sum_{i=1}^N (x_i - f_2(f_1(x)))^2 \end{aligned}$$

Unsupervised learning as supervised learning

- Examples:
 - Auto-encoders “predict” their inputs
 - Language models “predict” the next word

Unsupervised learning as supervised learning

- Examples:
 - Auto-encoders “predict” their inputs
 - Language models “predict” the next word
- Create a target from the x 's (or a subset)
 - Find a task to ensure that you learn something useful
 - *Self-supervised learning*

Unsupervised learning takeaways

- Most data (in the world) is unlabeled
- Useful tasks: clustering, density estimation, dimensionality reduction
- K-means and Gaussian mixtures (GMMs)
- Performance is harder to judge
 - Note that all our examples were in 2D
- Can be used in downstream tasks