

Data Modelling Methods-III

CS771: Introduction to Machine Learning
Purushottam Kar



Mid Semester Examination

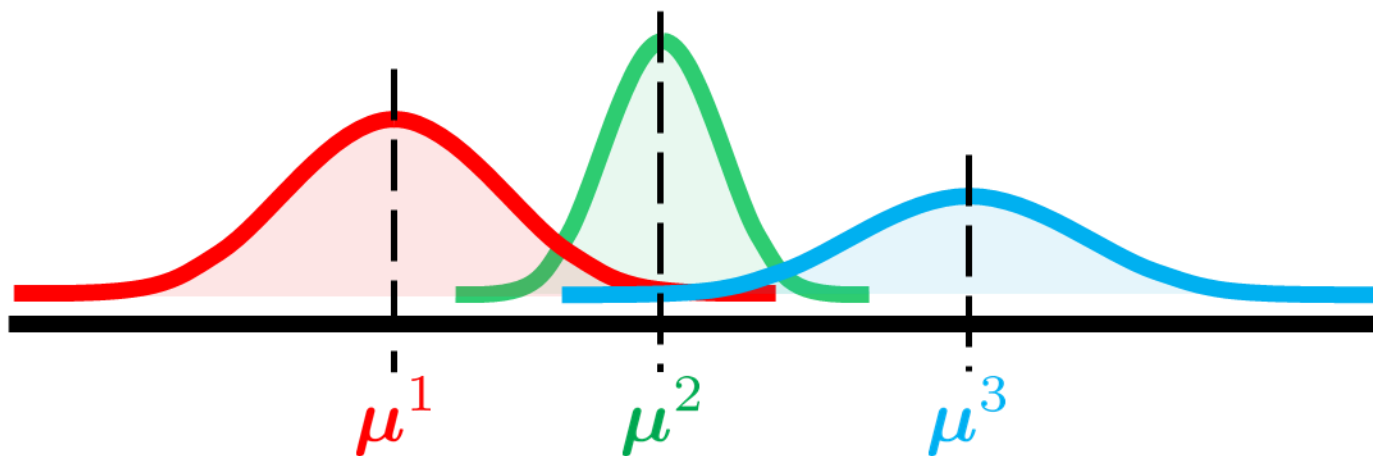
- September 21st, 2017 (Thursday) 1300–1500 hrs
- Venue L18, 19, L20 (all OROS)
- Syllabus: till whatever we cover today
- Open notes (handwritten only)
- No printed/photocopied material
- No laptops, i-pads, mobile phones (switched off)
- Please bring a notepad with you for rough work
- Please bring a pencil/eraser with you – we will not provide these
- Answers will have to be written on the question paper itself

Recap

Sept 13, 2017



The generative story for labelled data



$$\mathbb{P}[\mathbf{x}^i, z^i | \Theta] = \mathbb{P}[z^i | \Theta] \cdot \mathbb{P}[\mathbf{x}^i | z^i, \Theta] = \pi_{z^i} \cdot \mathcal{N}(\mathbf{x}^i | \mu^{z^i}, \Sigma^{z^i})$$

$$\mathbb{P}[X, \{z^i\} | \Theta] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i, z^i | \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X, \{z^i\} | \Theta]$$

Take log and apply 1st order optimality

Read [DAU]
Sections 9.1-9.5

$$\pi_1^{\text{MLE}} = \frac{\# \text{ red emails}}{\# \text{ total emails}} = \frac{n_r}{n}$$

$$\mu_{\text{MLE}}^1 = \frac{1}{n_r} \sum_{i: z^i = \bullet} \mathbf{x}^i$$

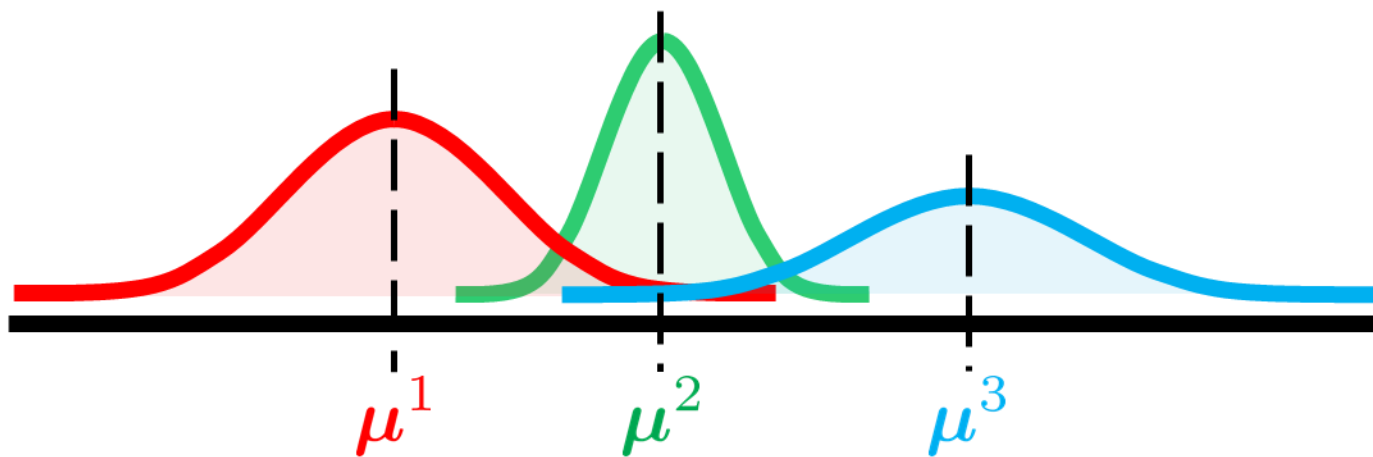
$$\Sigma_{\text{MLE}}^1 = \frac{1}{n_r} \sum_{i: z^i = \bullet} (\mathbf{x}^i - \mu_{\text{MLE}}^1)(\mathbf{x}^i - \mu_{\text{MLE}}^1)^{\top}$$

Recap

Sept 13, 2017



The generative story for unlabelled data



z^i denotes the (unknown) component from which x^i came

z^i not known from data. It is a *latent variable* (can take K values).

Gaussian Mixture Model (GMM) with K components

$\pi_k = \mathbb{P}[z^i = k]$ prior prob. of x^i coming from k -th component

Goal: incomplete data, learn $\mu^k, \Sigma^k, \mathbb{P}[z = k]$

$$\begin{aligned} \mathbb{P}[\mathbf{x}^i | \Theta] &= \sum_{k=1}^K \mathbb{P}[\mathbf{x}^i, z^i = k | \Theta] = \sum_{k=1}^K \pi_k \cdot \mathbb{P}[\mathbf{x}^i | z^i = k, \Theta] \\ &= \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}^i | \mu^k, \Sigma^k) \end{aligned}$$

Recap

Sept 13, 2017



A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X \mid \Theta]$$

Looks like block coordinate descent with $\Theta, \{z^i\}$ being two blocks of "coordinates"

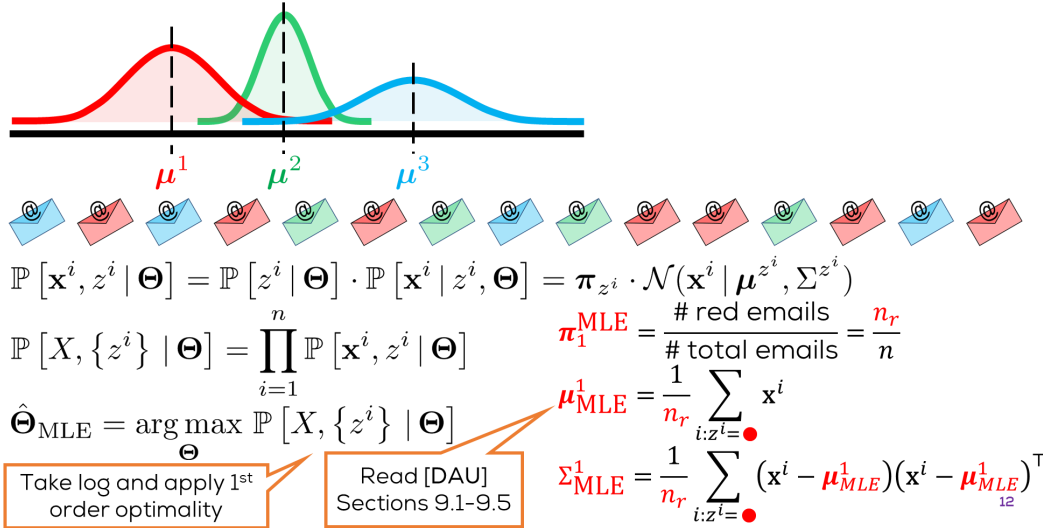
ALTERNATING OPTIMIZATION

1. Initialize Θ^0
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
3. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} \mid \Theta]$
4. Repeat until convergence

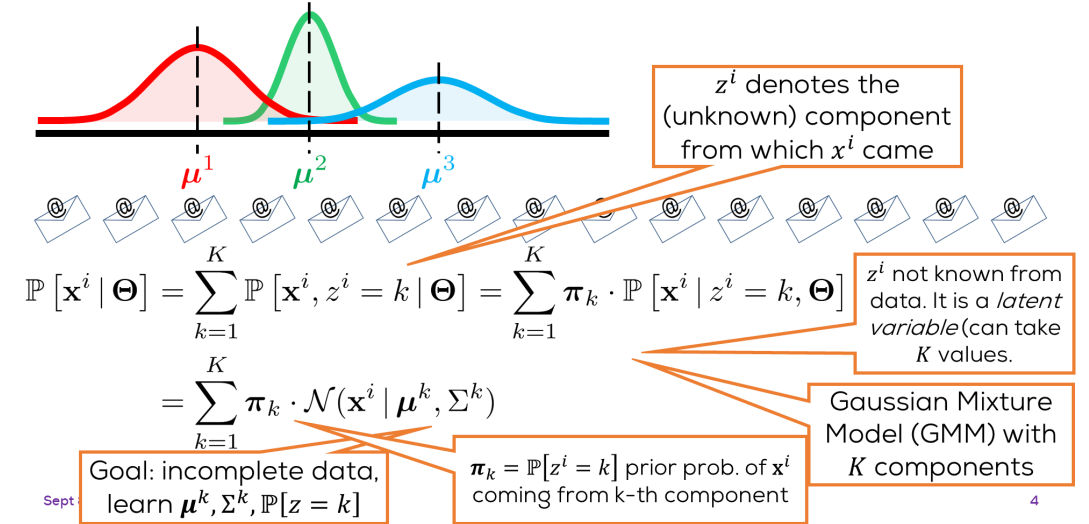
Various ways of updating z^i

Recap

The generative story for labelled data



The generative story for unlabelled data



A Ray of Hope

$$\hat{\Theta}_{MLE} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

ALTERNATING OPTIMIZATION

1. Initialize Θ^0
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
3. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

Various ways of updating z^i

Looks like block coordinate descent with $\Theta, \{z^i\}$ being two blocks of "coordinates"

Hard Assignment

The K-means algorithm

Sept 8, 2017



A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

Looks like block coordinate descent with $\Theta, \{z^i\}$ being two blocks of "coordinates"

ALTERNATING OPTIMIZATION

1. Initialize Θ^0
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
3. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

Various ways of updating z^i

A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

ALTERNATING OPTIMIZATION

1. Initialize Θ^0
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
3. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

ALTERNATING OPTIMIZATION

1. Initialize Θ^0
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
3. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t]$$

A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

ALTERNATING OPTIMIZATION

1. Initialize Θ^0
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
3. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t]$$

$$\Theta^t = \left\{ \pi^t, \left\{ \mu^{1,t}, \mu^{2,t}, \mu^{3,t} \right\}, \left\{ \Sigma^{1,t}, \Sigma^{2,t}, \Sigma^{3,t} \right\} \right\}$$

A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

ALTERNATING OPTIMIZATION

1. Initialize Θ^0
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
3. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

Bayes Rule!

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t] = \arg \max_{k \in [K]} \mathbb{P}[k | \Theta^t] \cdot \mathbb{P}[\mathbf{x}^i | k, \Theta^t]$$

$$\Theta^t = \left\{ \pi^t, \left\{ \mu^{1,t}, \mu^{2,t}, \mu^{3,t} \right\}, \left\{ \Sigma^{1,t}, \Sigma^{2,t}, \Sigma^{3,t} \right\} \right\}$$

A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

ALTERNATING OPTIMIZATION

1. Initialize Θ^0
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
3. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

Bayes
Rule!

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t] = \arg \max_{k \in [K]} \pi_k^t \cdot \mathbb{P}[\mathbf{x}^i | k, \Theta^t]$$

$$\Theta^t = \left\{ \pi^t, \left\{ \mu^{1,t}, \mu^{2,t}, \mu^{3,t} \right\}, \left\{ \Sigma^{1,t}, \Sigma^{2,t}, \Sigma^{3,t} \right\} \right\}$$

A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

ALTERNATING OPTIMIZATION

1. Initialize Θ^0
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
3. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

Bayes Rule!

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t] = \arg \max_{k \in [K]} \pi_k^t \cdot \mathcal{N}(\mathbf{x}^i | \mu^{k,t}, \Sigma^{k,t})$$

$$\Theta^t = \left\{ \pi^t, \left\{ \mu^{1,t}, \mu^{2,t}, \mu^{3,t} \right\}, \left\{ \Sigma^{1,t}, \Sigma^{2,t}, \Sigma^{3,t} \right\} \right\}$$

Towards the K-Means Algorithm

ALTERNATING OPTIMIZATION

1. Initialize Θ^0
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
 1. Let $z^{i,t} = \arg \max_k \pi_k^t \cdot \mathcal{N}(\mathbf{x}^i \mid \boldsymbol{\mu}^{k,t}, \Sigma^{k,t})$
3. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} \mid \Theta]$
 1. Let $\pi_k^{t+1} = \frac{n_k^t}{n}$, where $n_k^t = |\{i: z^{i,t} = k\}|$
 2. Let $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
 3. Let $\Sigma_k^{t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} (\mathbf{x}^i - \boldsymbol{\mu}^{k,t+1})(\mathbf{x}^i - \boldsymbol{\mu}^{k,t+1})^\top$
4. Repeat until convergence



A few simplifications

- Fix $\boldsymbol{\pi}_k^t = \frac{1}{K}$ for all iterations. Don't update it.
- Fix $\boldsymbol{\Sigma}^{k,t} = I$ for all iterations. Don't update it.

K-MEANS/LLOYD'S ALGORITHM

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1,\dots,K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
 1. Let $z^{i,t} = \arg \max_k \mathcal{N}(\mathbf{x}^i \mid \boldsymbol{\mu}^{k,t}, I)$
3. Update $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

A few simplifications

- Fix $\boldsymbol{\pi}_k^t = \frac{1}{K}$ for all iterations. Don't update it.
- Fix $\boldsymbol{\Sigma}^{k,t} = I$ for all iterations. Don't update it.

K-MEANS/LLOYD'S ALGORITHM

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1,\dots,K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
 1. Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\|_2$
3. Update $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

A few simplifications

- Fix $\boldsymbol{\pi}_k^t = \frac{1}{K}$ for all iterations. Don't update it.
- Fix $\boldsymbol{\Sigma}^{k,t} = I$ for all iterations. Don't update it.

K-MEANS/LLOYD'S ALGORITHM

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1,\dots,K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
 1. Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\|_2^2$
3. Update $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

The K-Means Objective

The K-Means Objective

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\boldsymbol{\mu}^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i: z^i=k} \|\mathbf{x}^i - \boldsymbol{\mu}^k\|_2^2$$

The K-Means Objective

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\boldsymbol{\mu}^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i: z^i=k} \|\mathbf{x}^i - \boldsymbol{\mu}^k\|_2^2$$

K-MEANS/LLOYD'S ALGORITHM

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
 1. Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\|_2^2$
3. Update $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

The K-Means Objective

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\boldsymbol{\mu}^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i: z^i=k} \|\mathbf{x}^i - \boldsymbol{\mu}^k\|_2^2$$

Alternates between updating $\{z^i\}$ and $\{\boldsymbol{\mu}^k\}$

K-MEANS/LLOYD'S ALGORITHM

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
 1. Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\|_2^2$
3. Update $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

The K-Means Objective

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\mu^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i: z^i=k} \|\mathbf{x}^i - \mu^k\|_2^2$$

An FA approach to solving a data modelling task!

Alternates between updating $\{z^i\}$ and $\{\mu^k\}$

K-MEANS/LLOYD'S ALGORITHM

1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
 1. Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

The K-Means Objective

NP-hard problem!

An FA approach to solving a data modelling task!

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\mu^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i: z^i = k} \|\mathbf{x}^i - \mu^k\|_2^2$$

Alternates between updating $\{z^i\}$ and $\{\mu^k\}$

K-MEANS/LLOYD'S ALGORITHM

1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
 1. Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

The K-Means Objective

NP-hard problem!

An FA approach to solving a data modelling task!

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\boldsymbol{\mu}^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i: z^i=k} \|\mathbf{x}^i - \boldsymbol{\mu}^k\|_2^2$$

Alternates between updating $\{z^i\}$ and $\{\boldsymbol{\mu}^k\}$

Very scalable but sensitive to initialization!

K-MEANS/LLOYD'S ALGORITHM

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
 1. Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\|_2^2$
3. Update $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

The K-Means Objective

NP-hard problem!

An FA approach to solving a data modelling task!

$$\hat{\Theta}_{\text{km}} = \arg \min_{\substack{\{\mu^k\}_{k=1,\dots,K} \\ \{z^i\}_{i=1,\dots,n}}} \sum_{k=1}^K \sum_{i: z^i=k} \|\mathbf{x}^i - \mu^k\|_2^2$$

Alternates between updating $\{z^i\}$ and $\{\mu^k\}$

Very scalable but sensitive to initialization!

k-means++ initialization

1. Sample $i_1 \sim [n]$, let $\mu^{1,0} = \mathbf{x}^{i_1}$
2. For $k = 2, \dots, K$
 - Sample $i_k \propto \text{min distance from } \{\mu^{1,0}, \dots, \mu^{k-1,0}\}$
 - Let $\mu^{k,0} = \mathbf{x}^{i_k}$

K-MEANS/LLOYD'S ALGORITHM

1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
 1. Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

K-Means Algorithm in action!

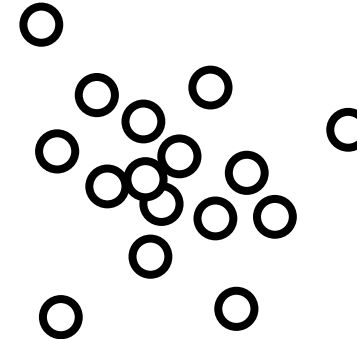
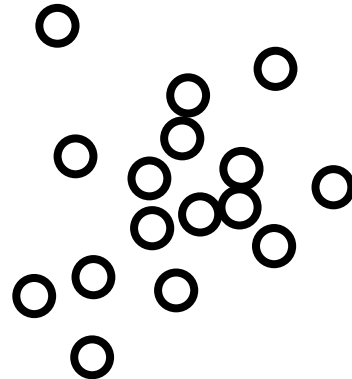
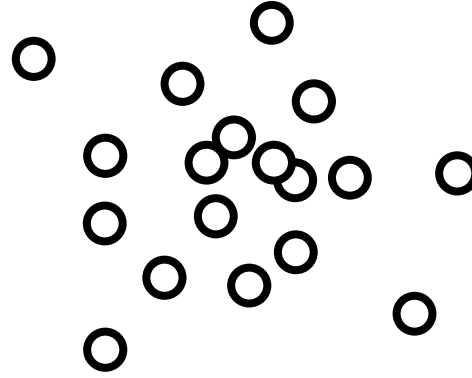
K-MEANS/LLOYD'S ALGORITHM

1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

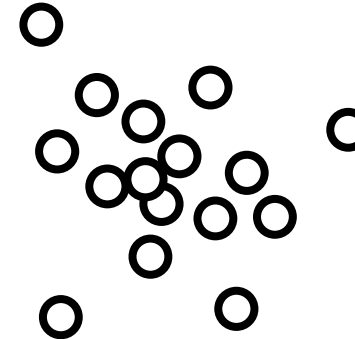
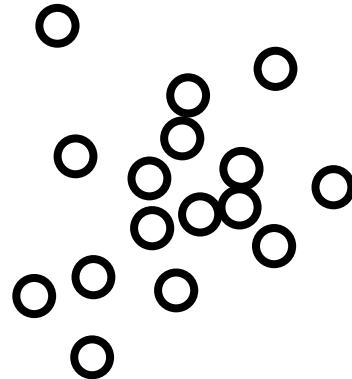
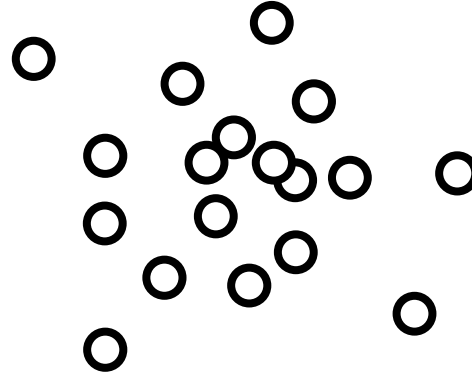
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

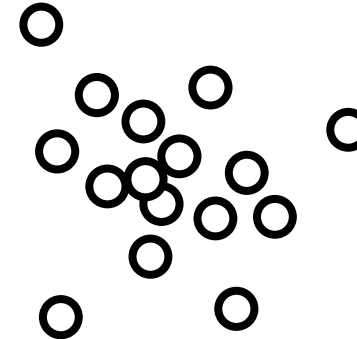
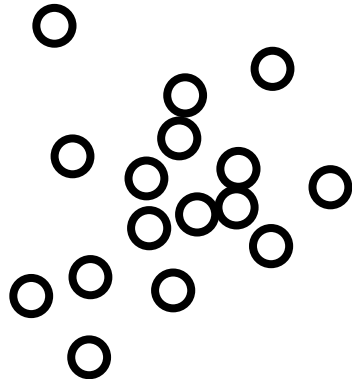
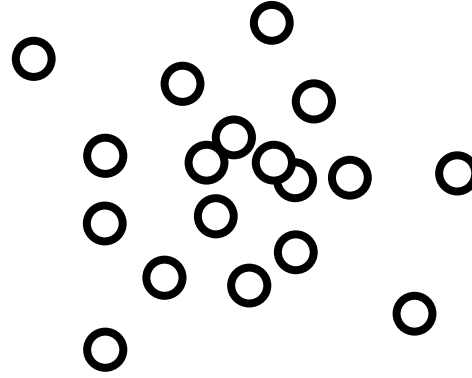
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

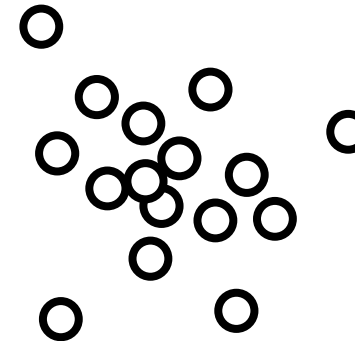
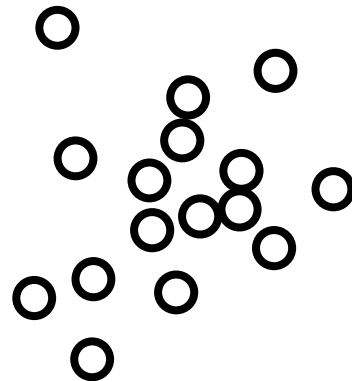
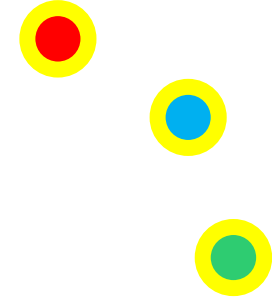
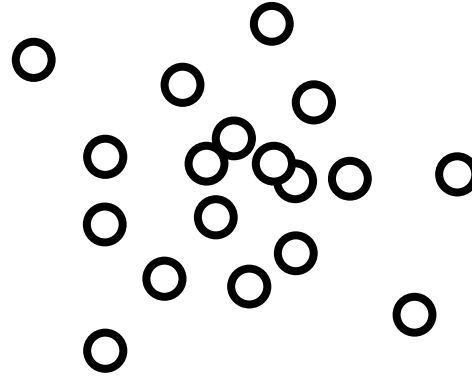
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

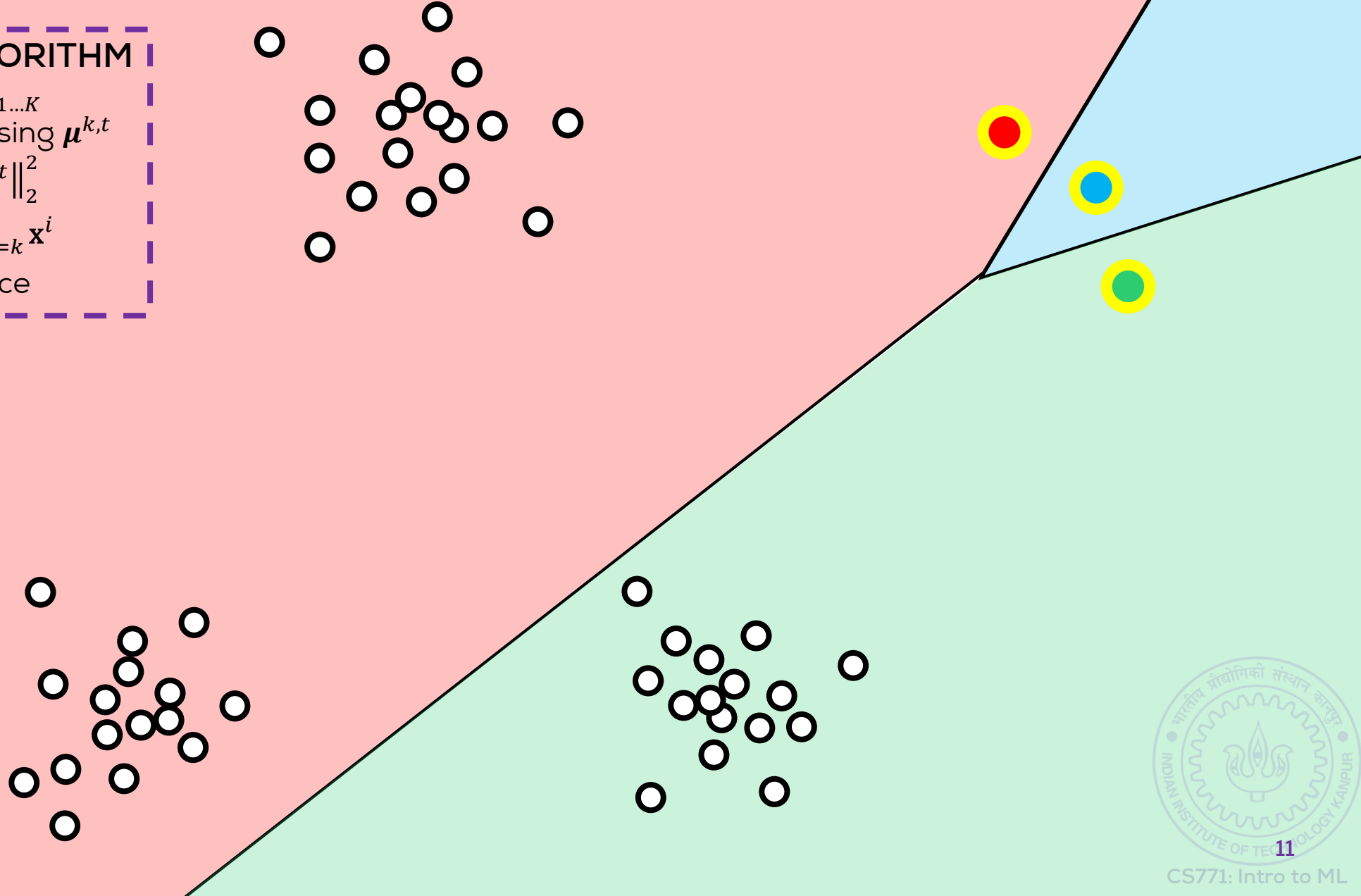
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

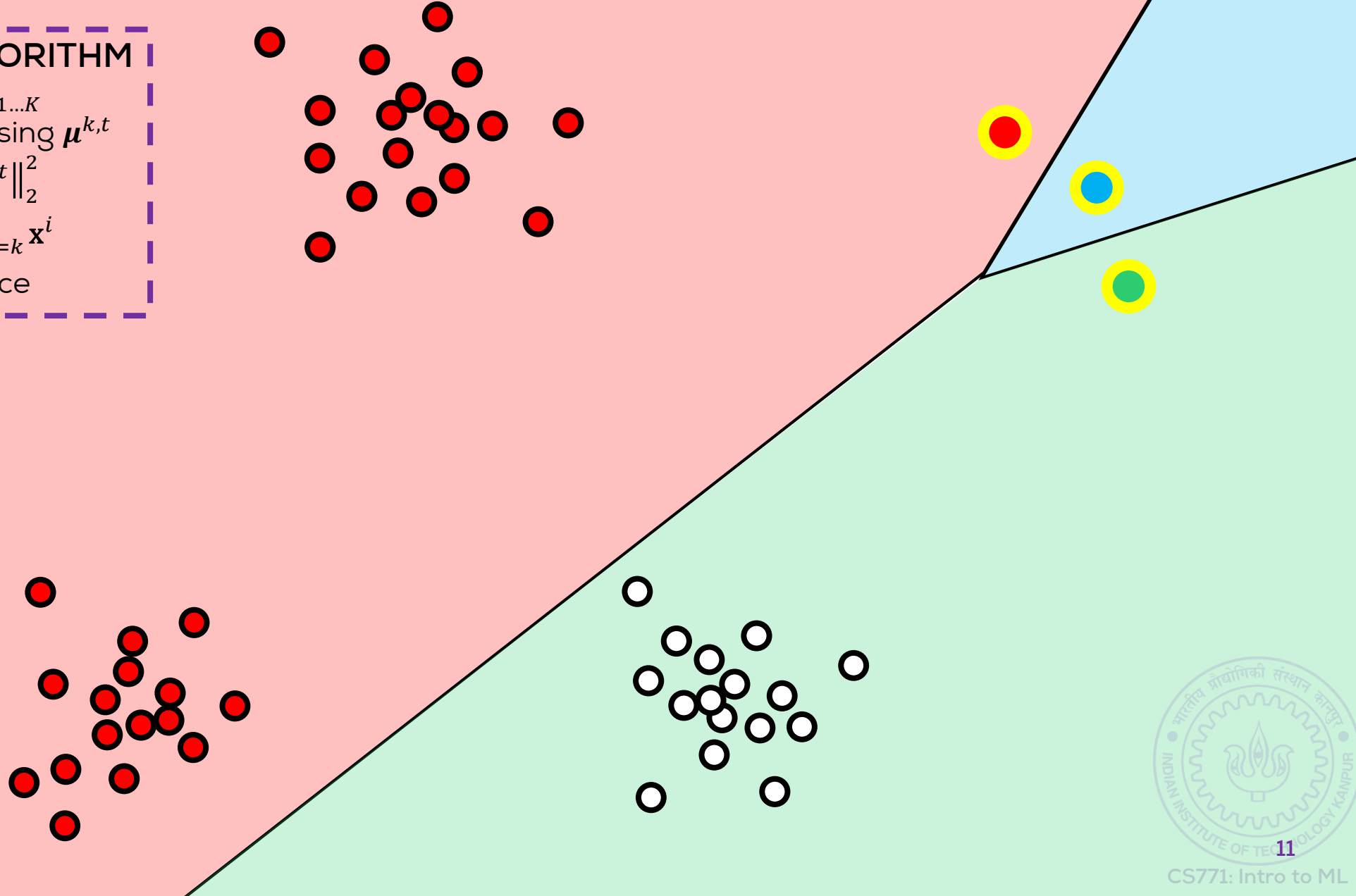
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

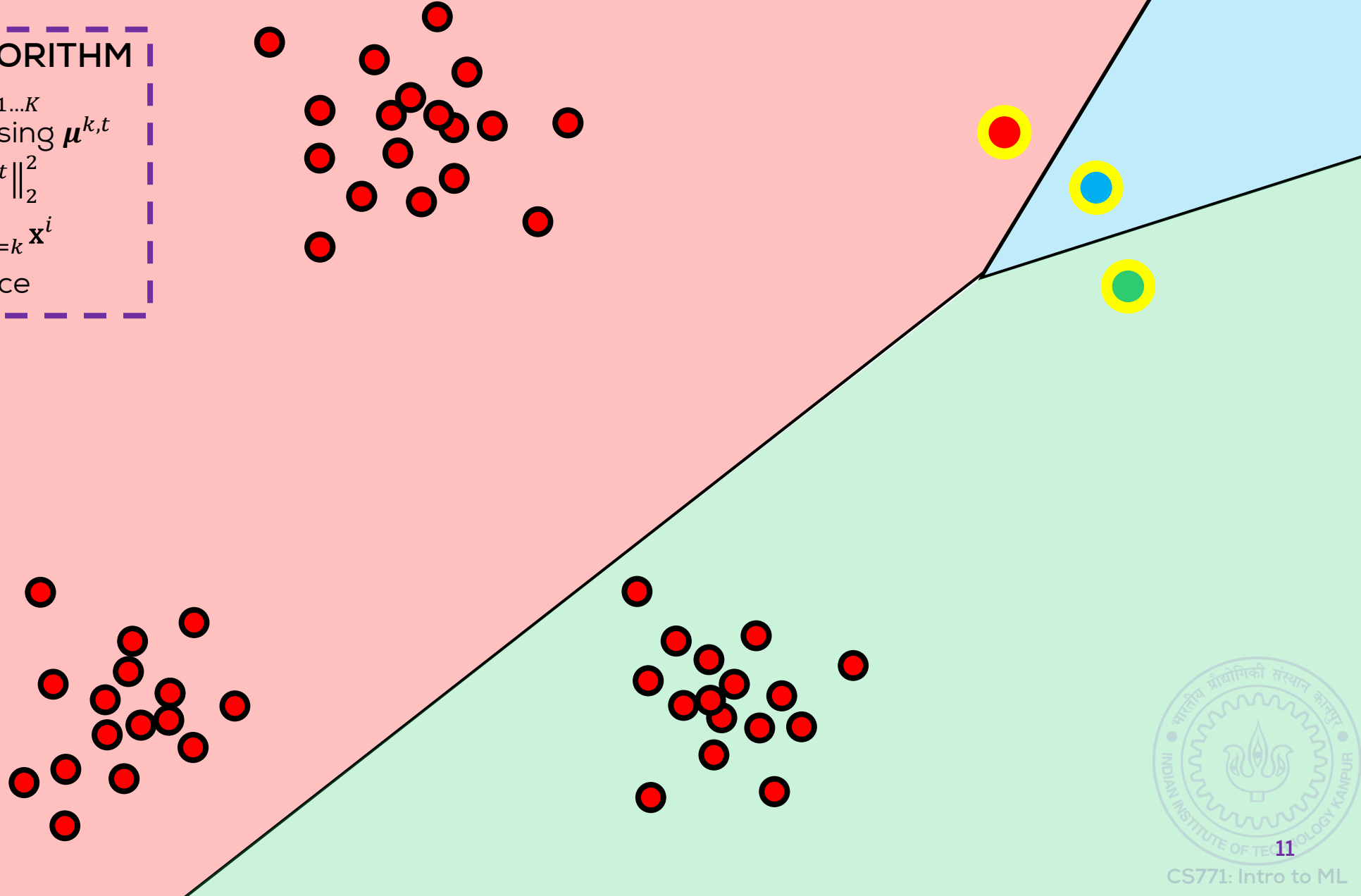
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

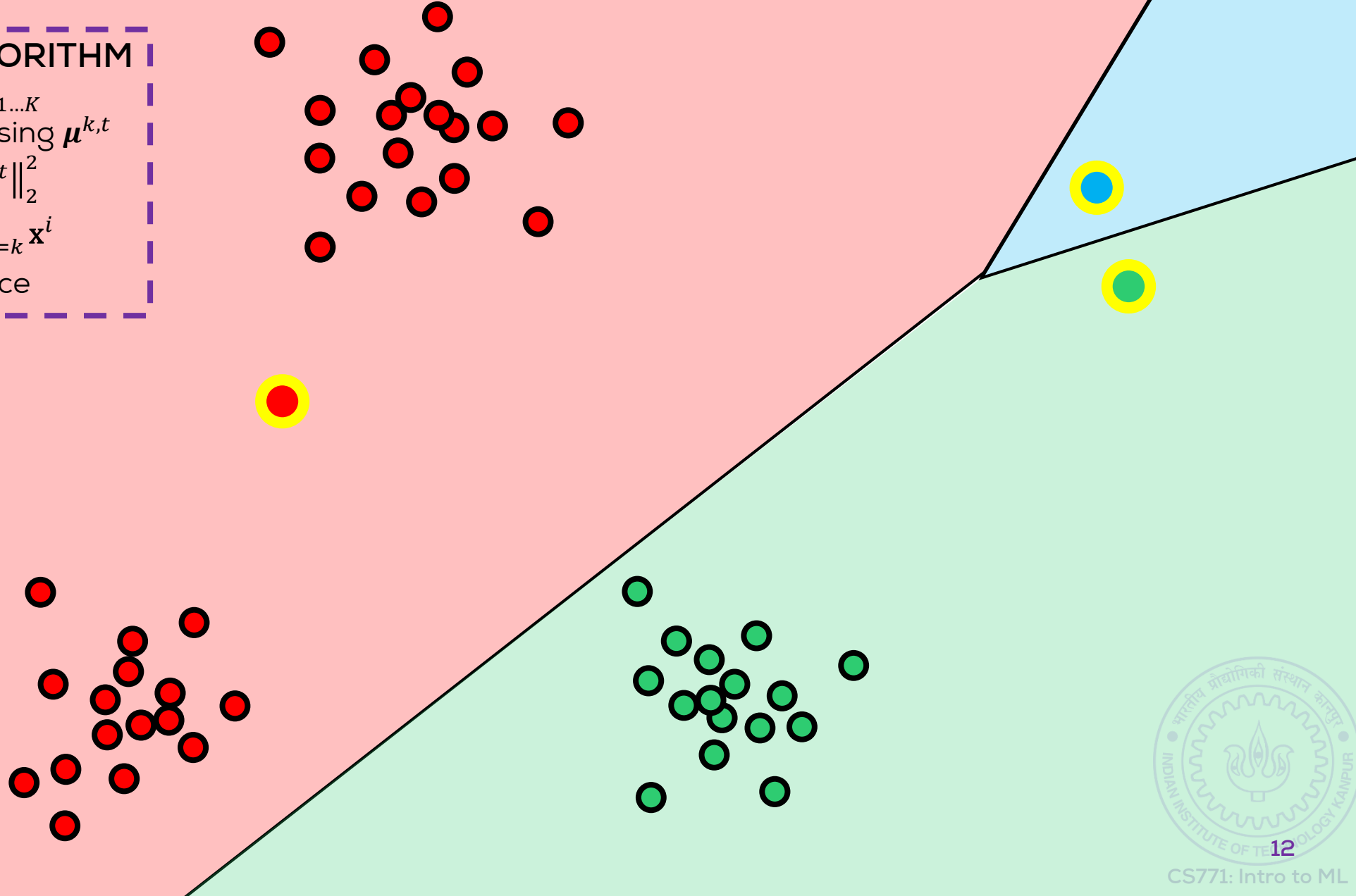
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

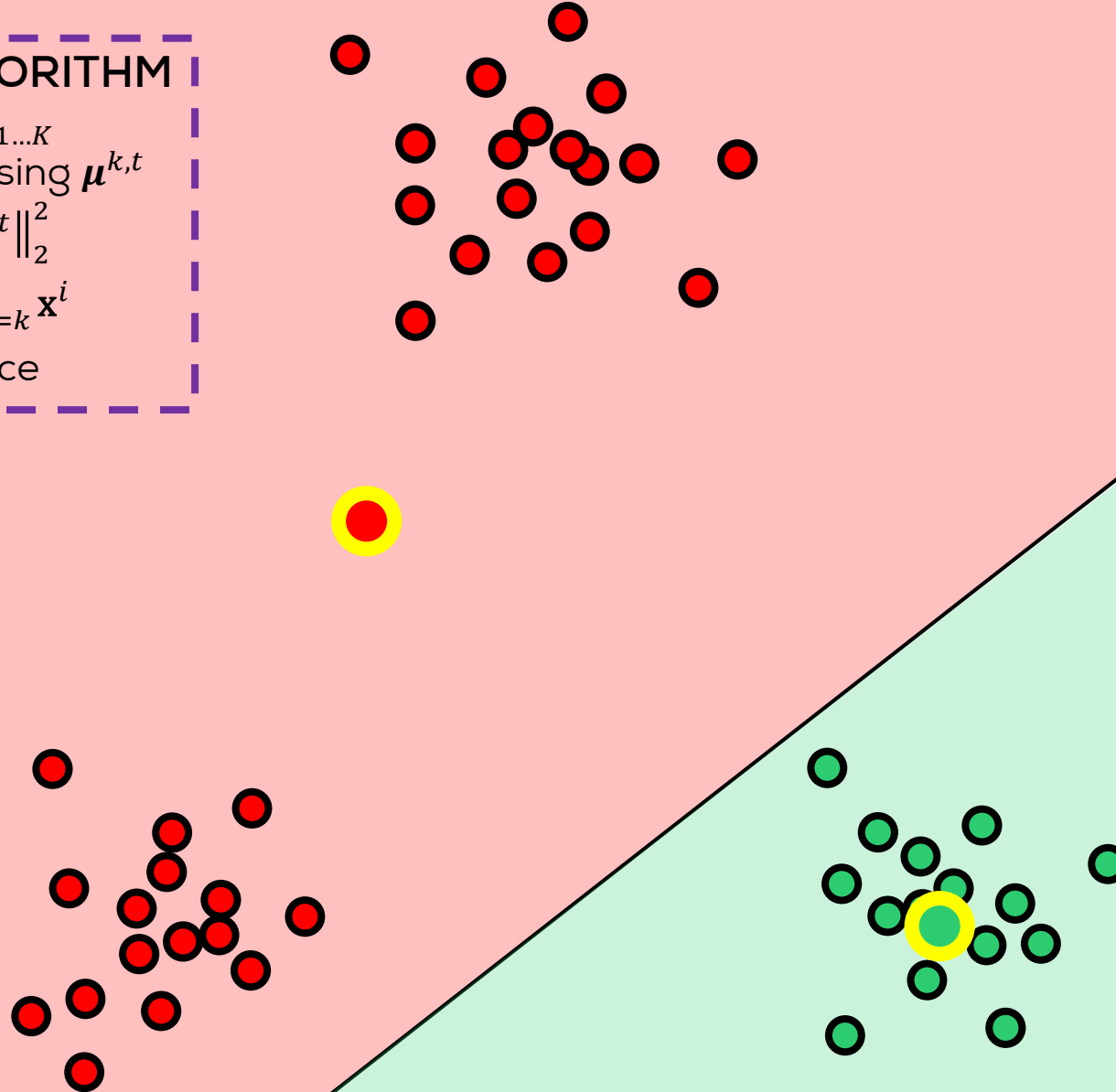
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

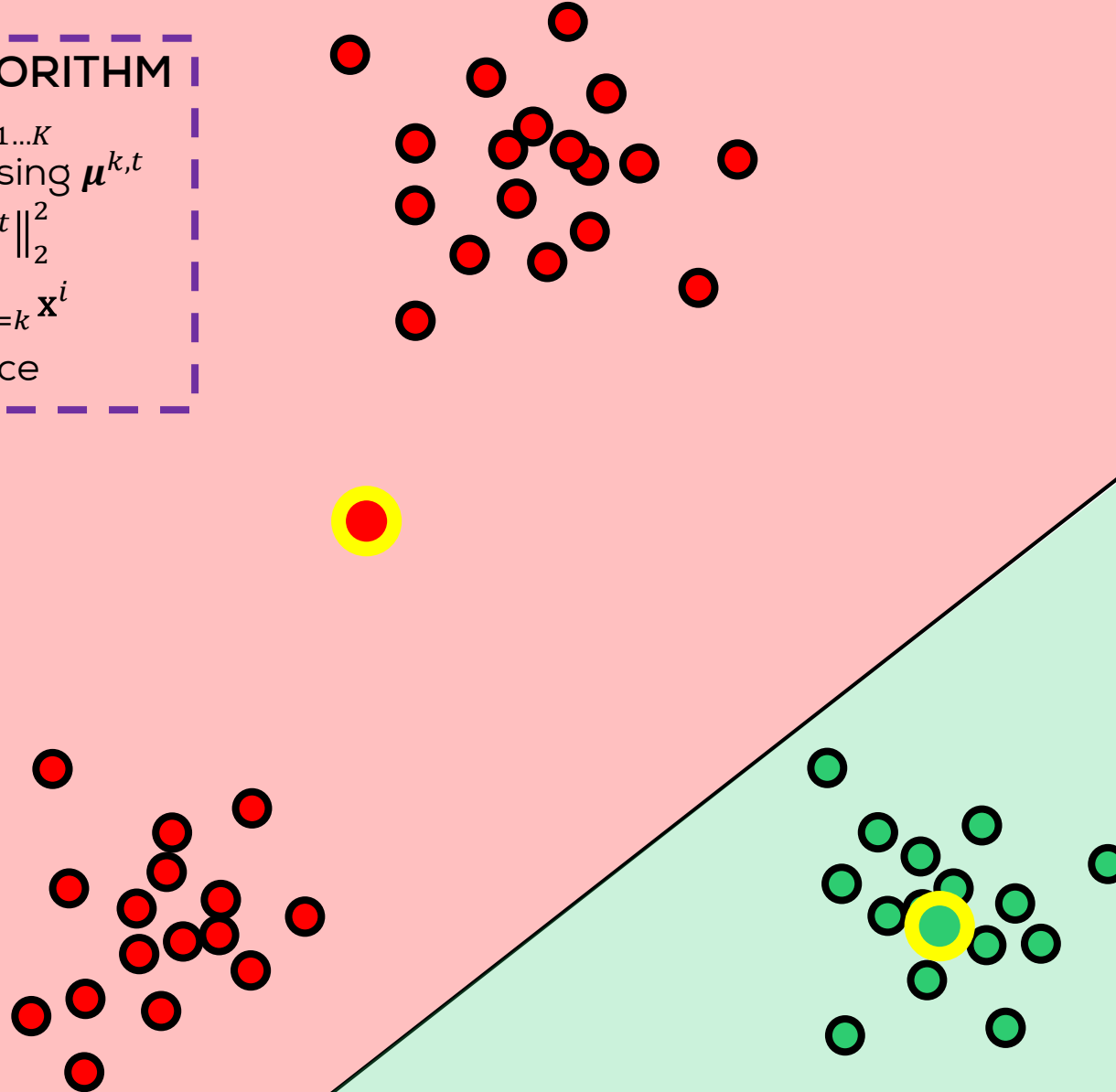
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

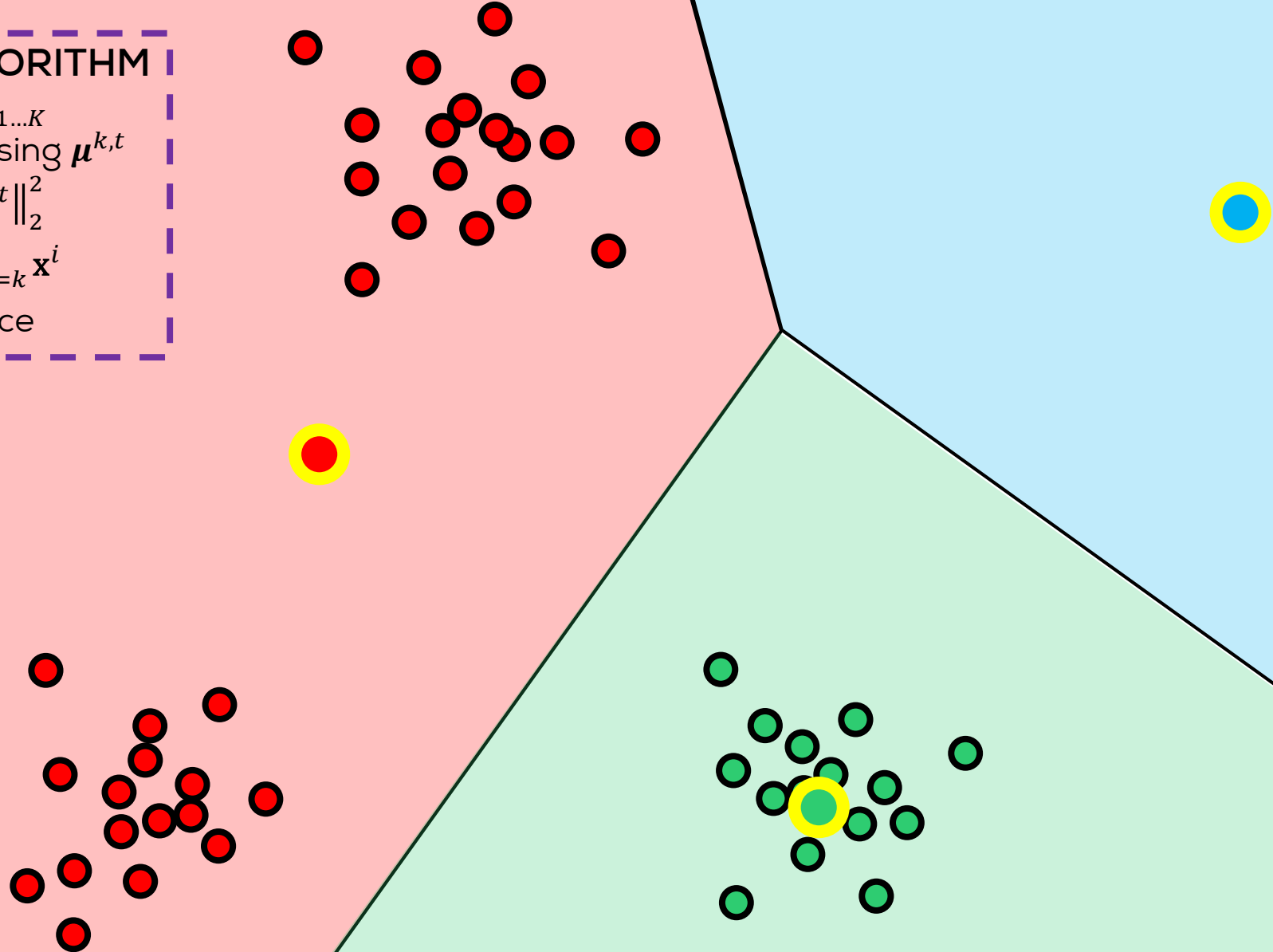
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

Stuck!!!

K-Means Algorithm in action!

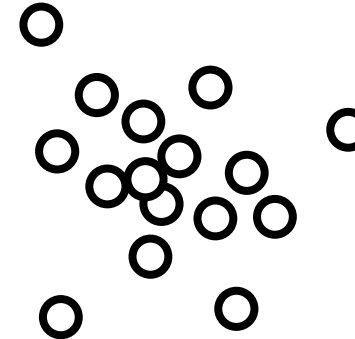
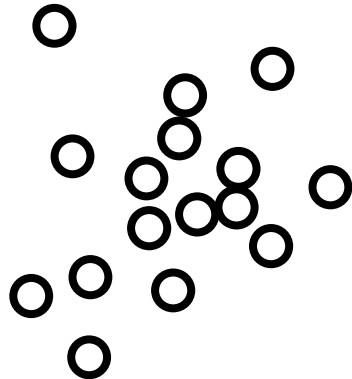
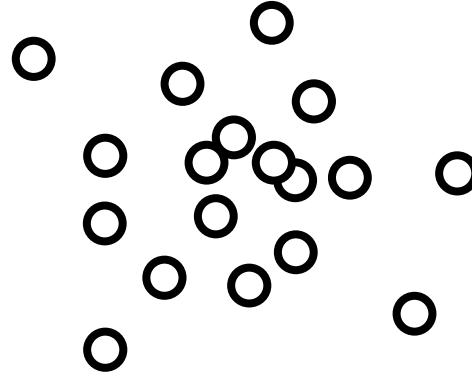
K-MEANS/LLOYD'S ALGORITHM

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\|_2^2$
3. Update $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

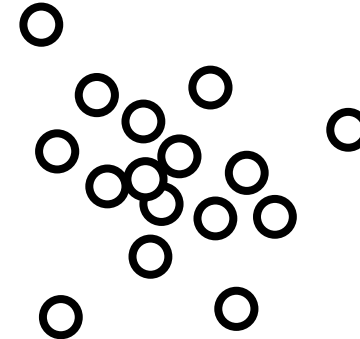
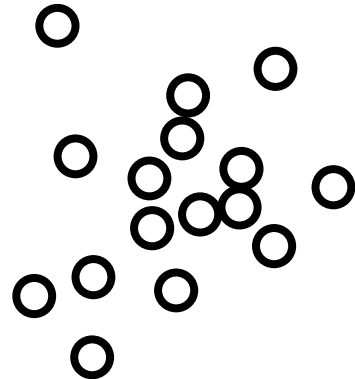
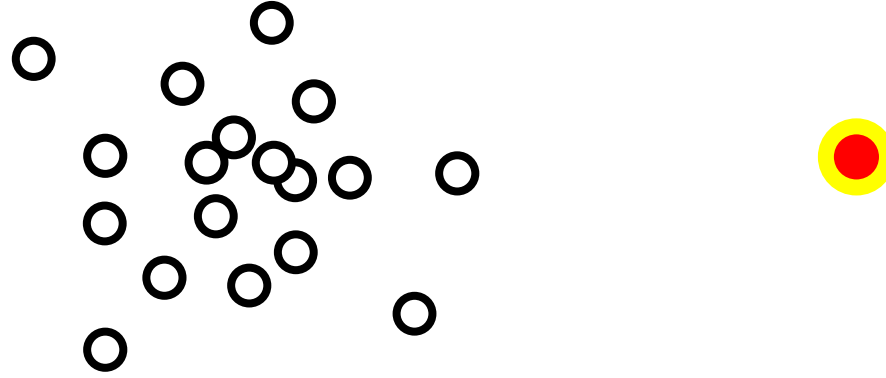
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

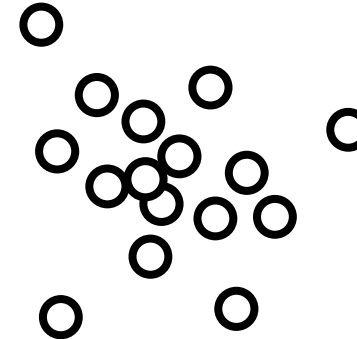
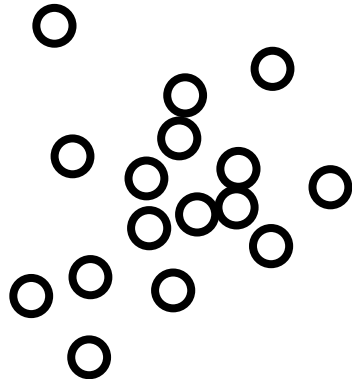
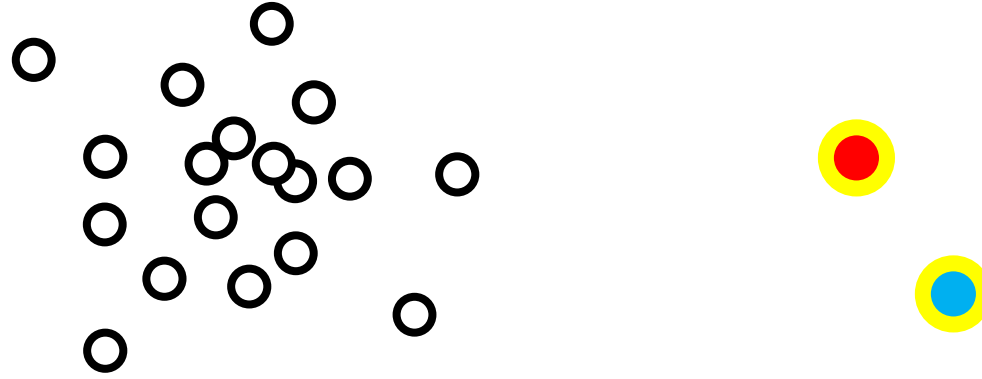
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

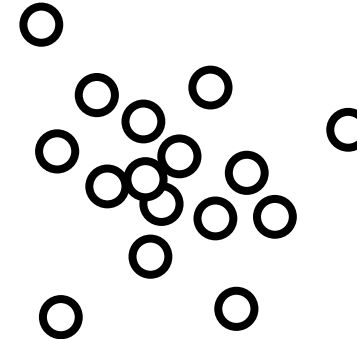
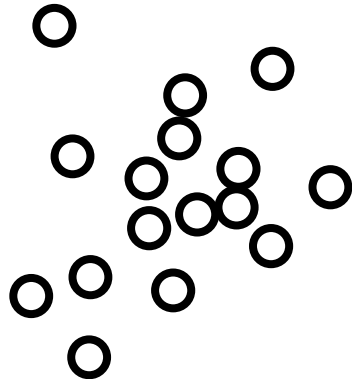
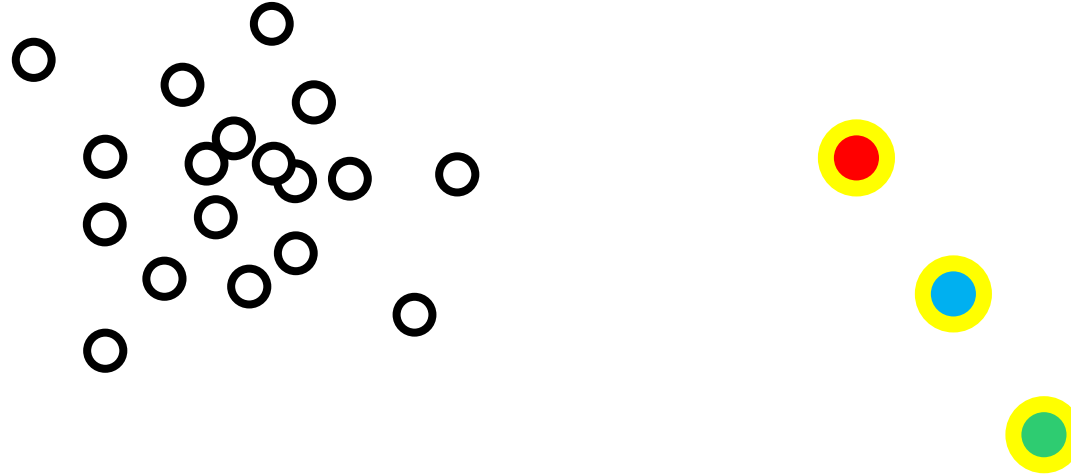
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

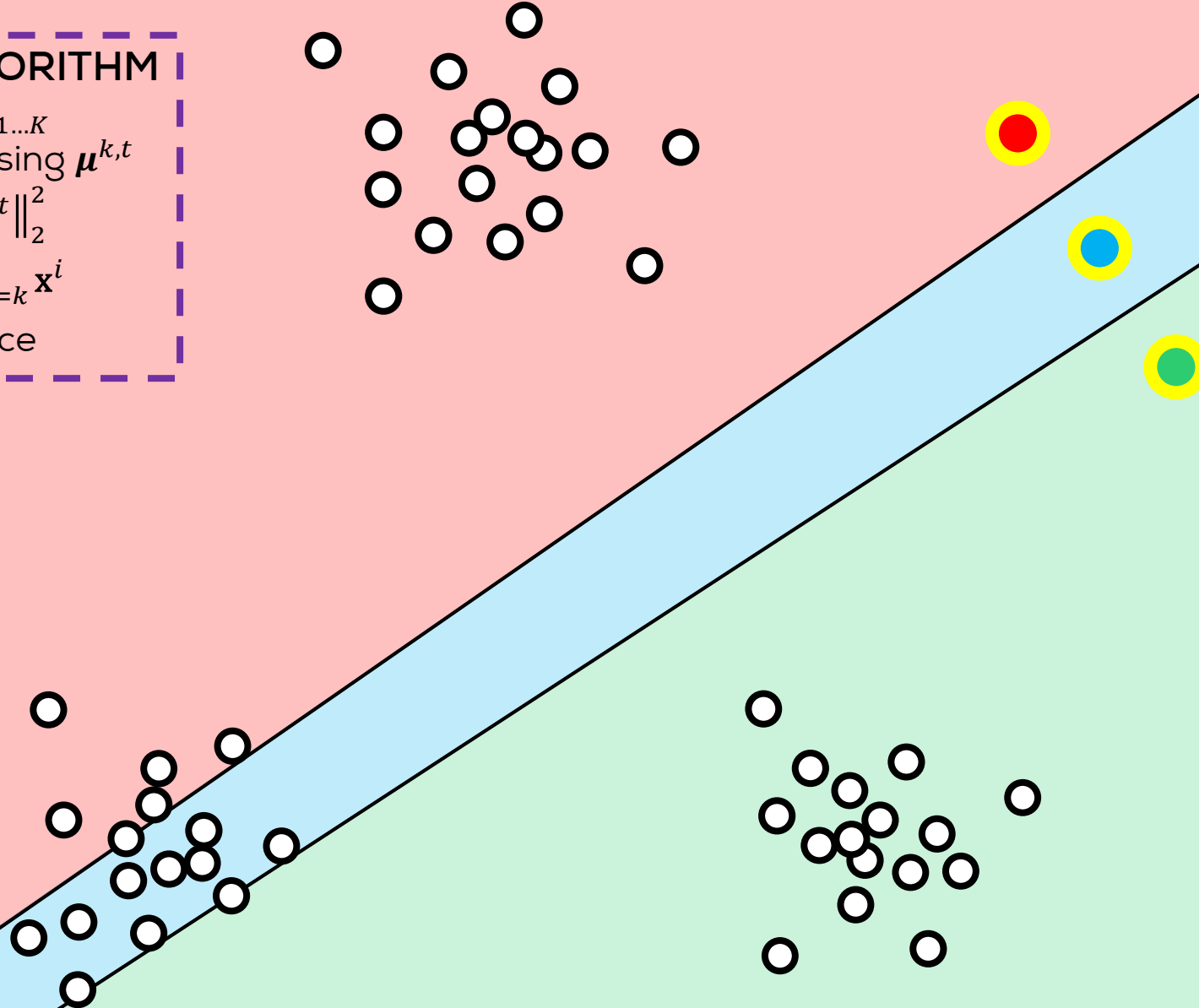
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

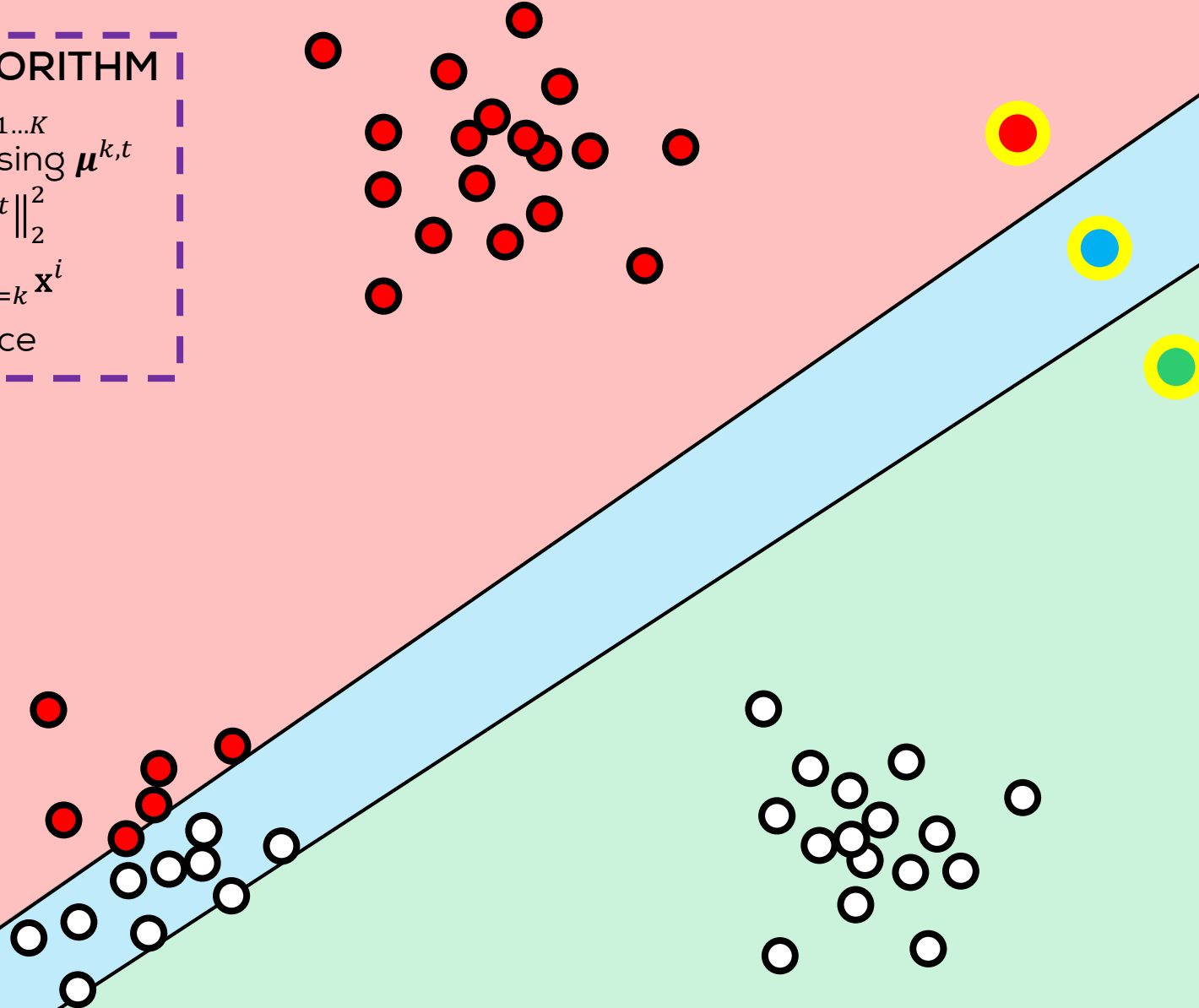
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

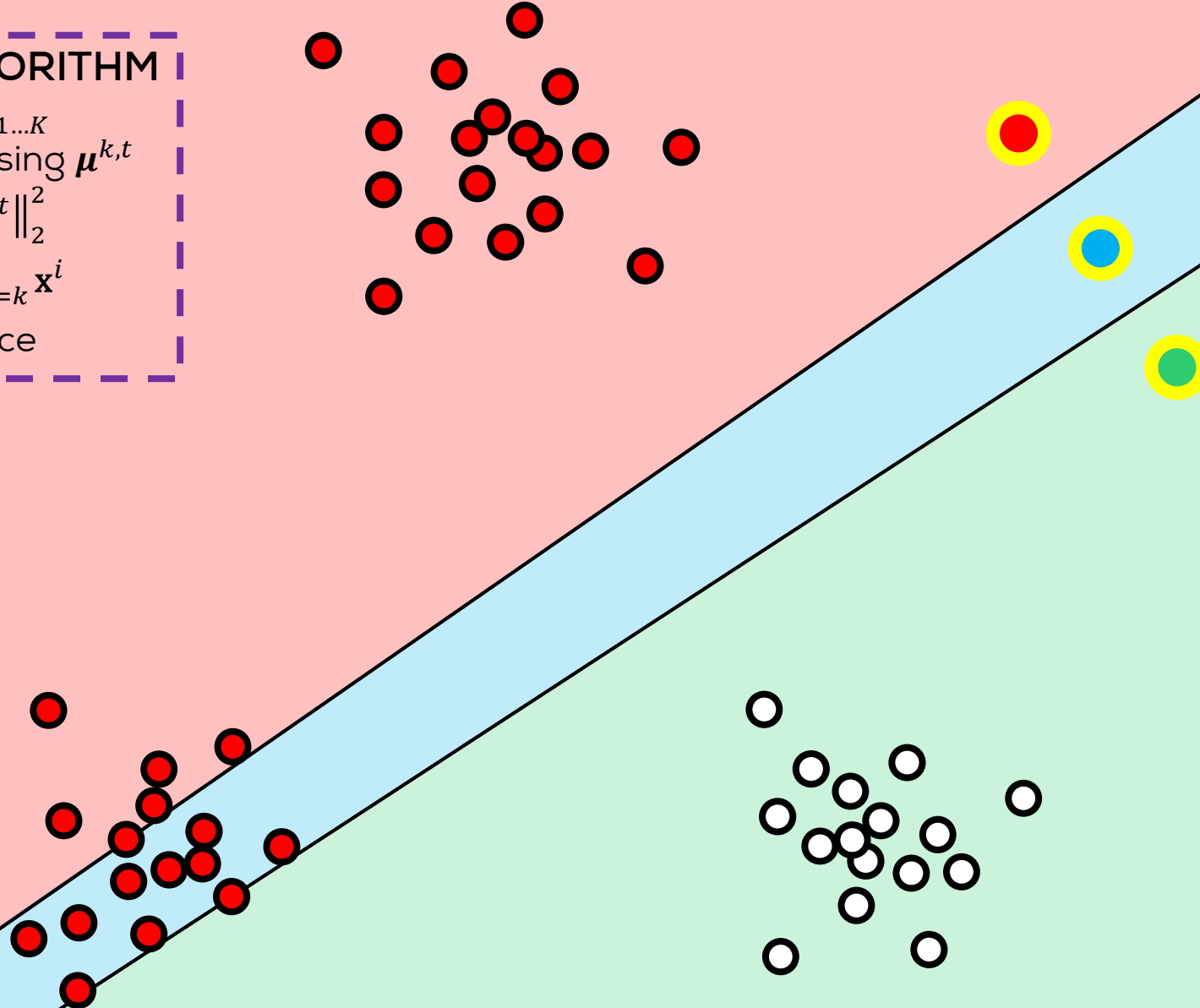
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

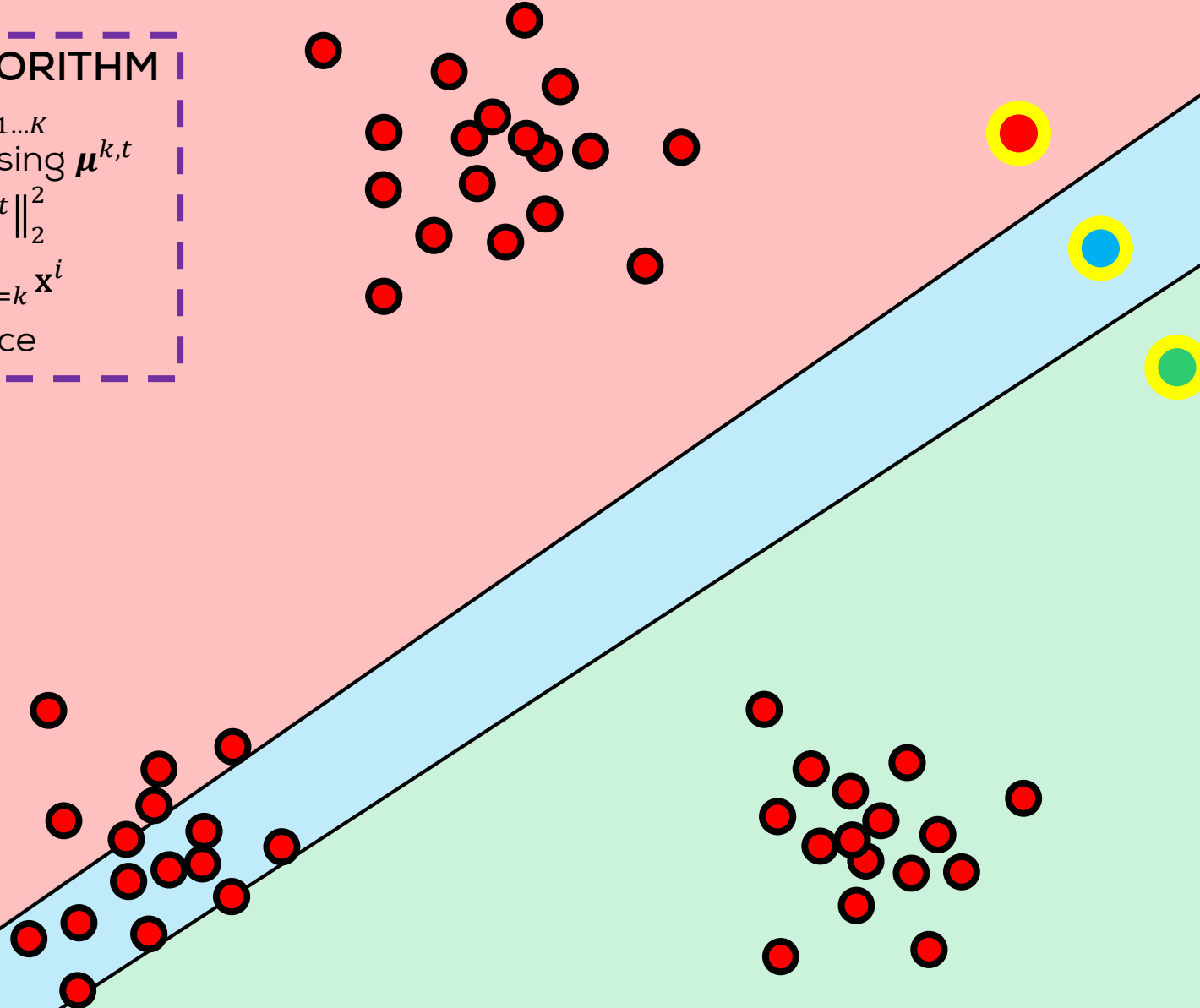
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

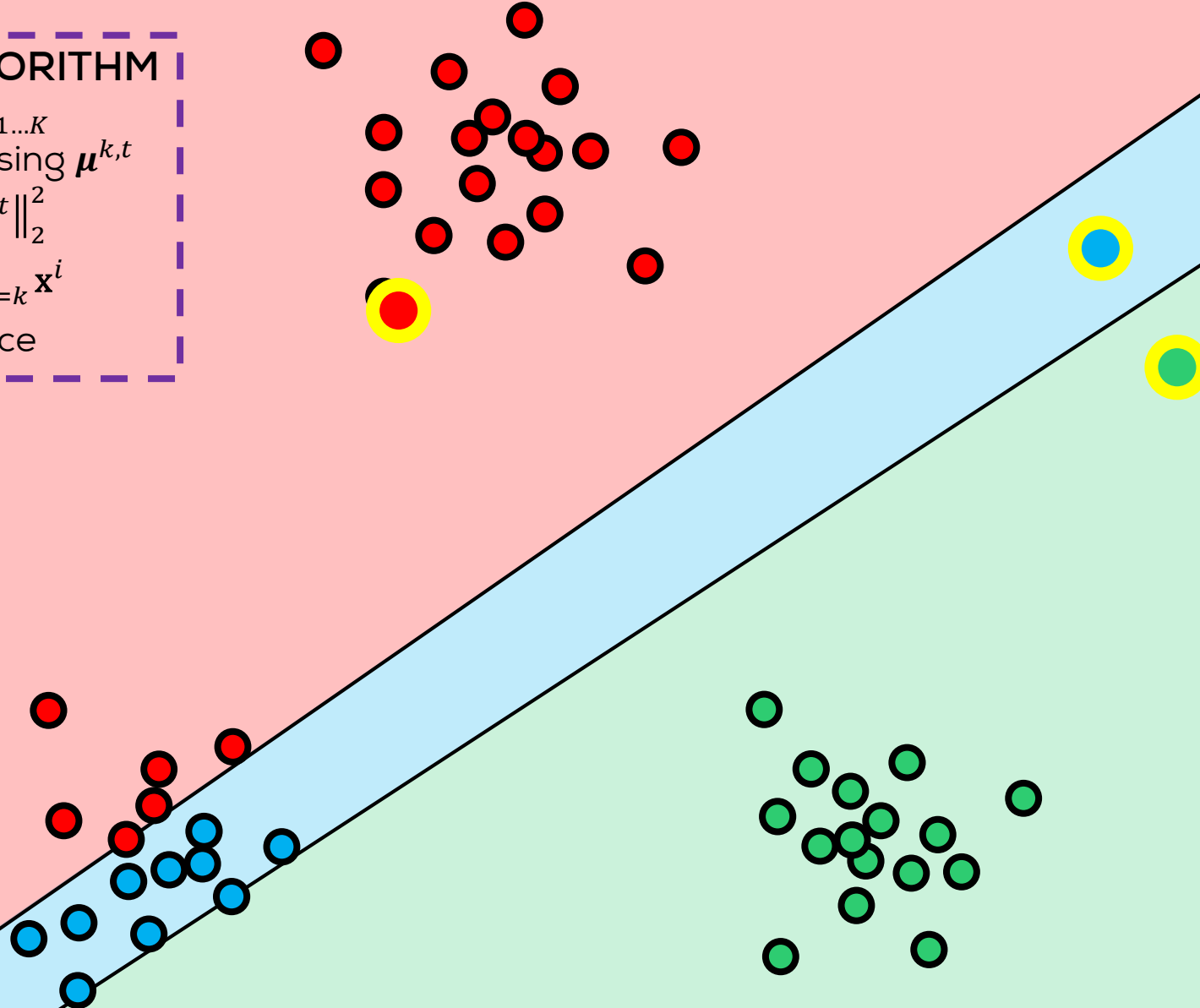
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

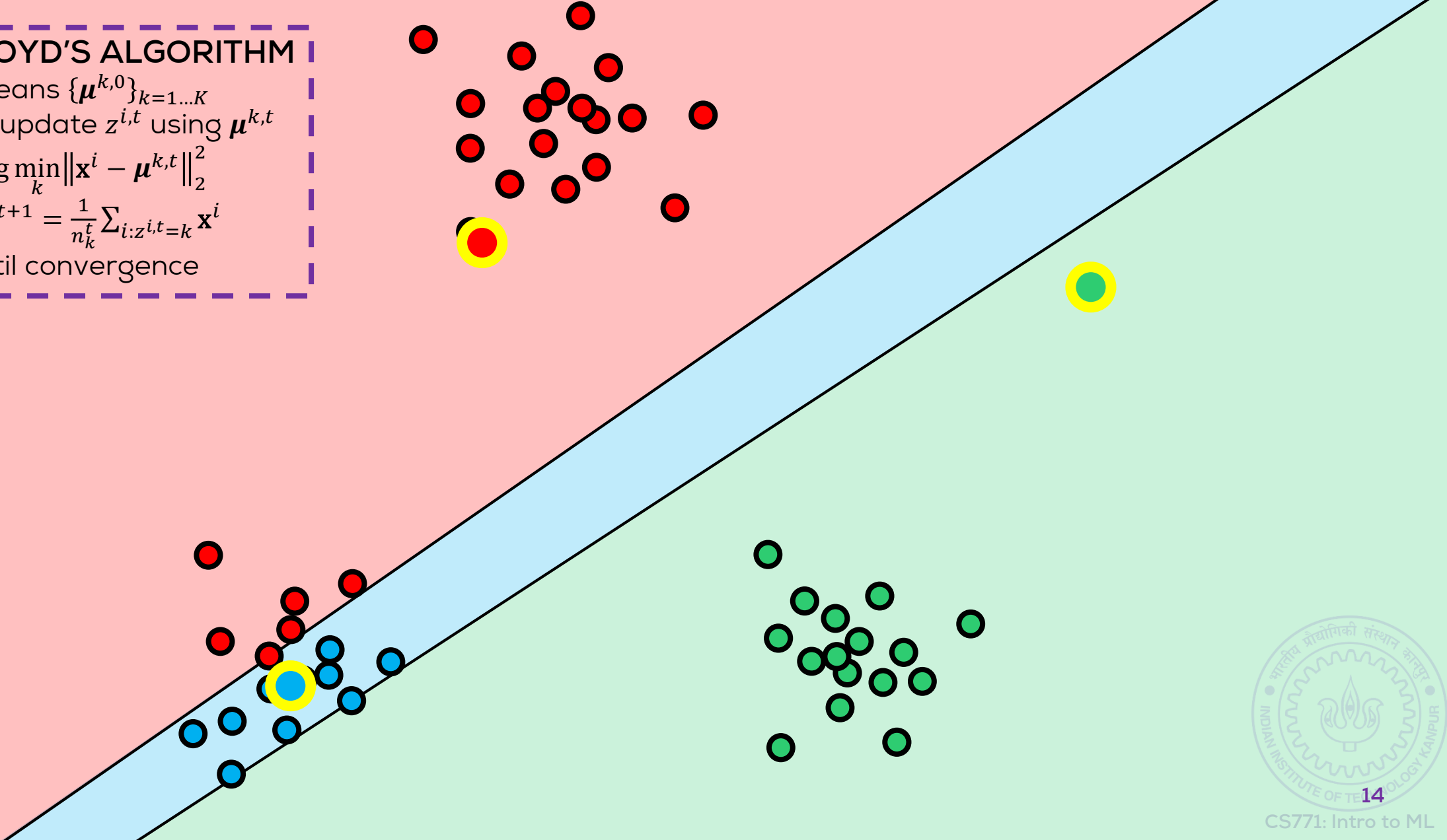
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

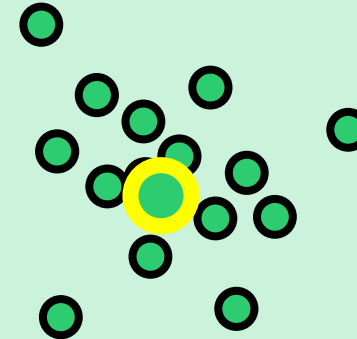
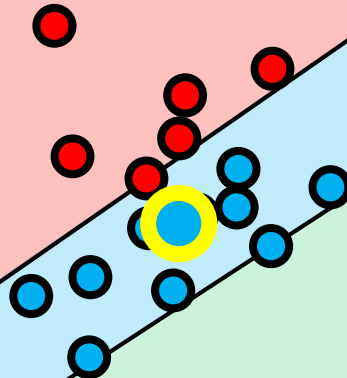
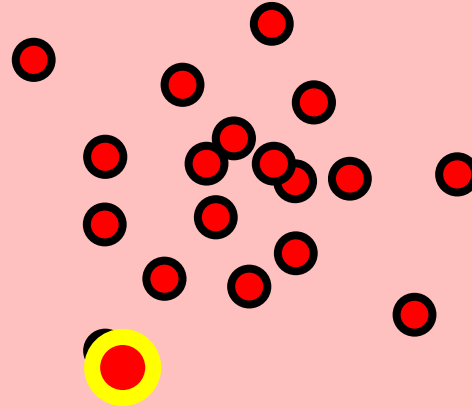
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

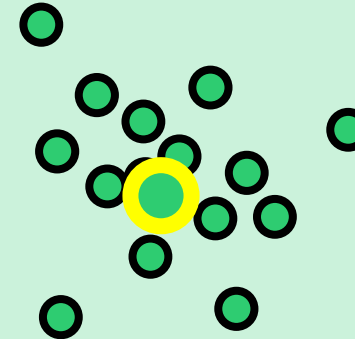
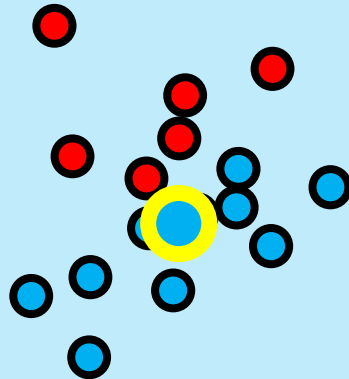
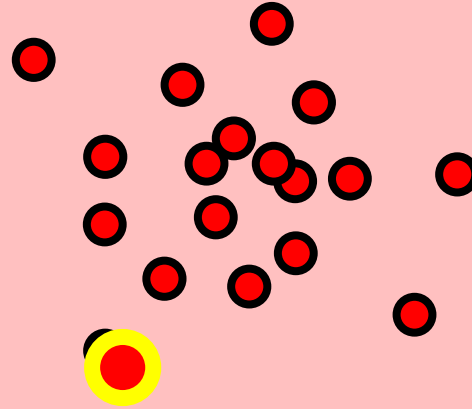
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

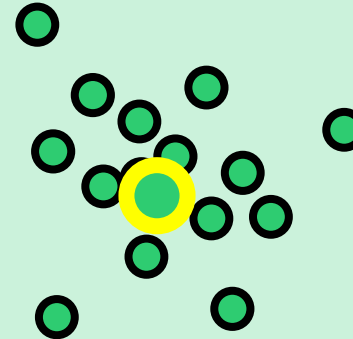
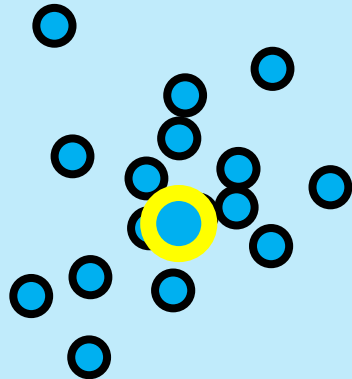
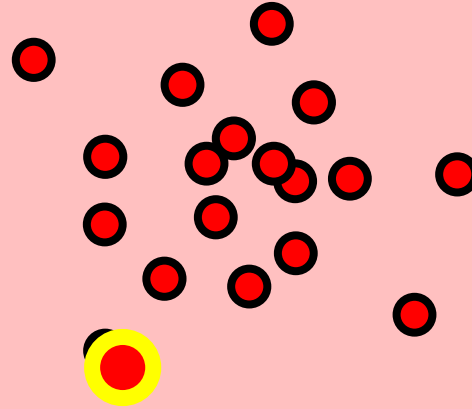
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

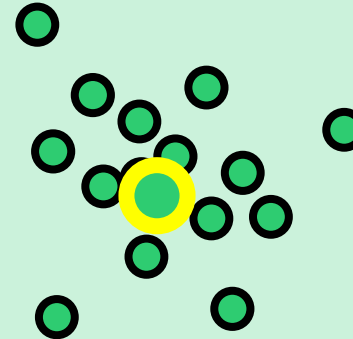
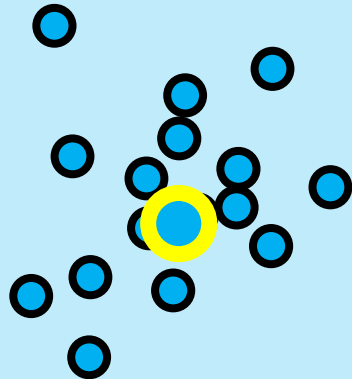
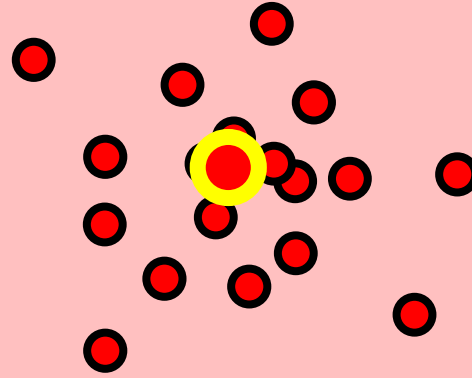
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

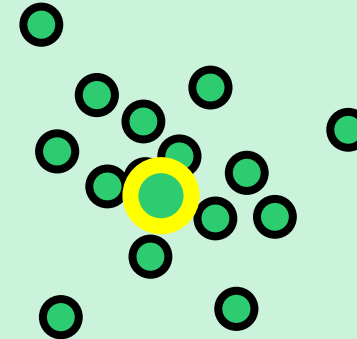
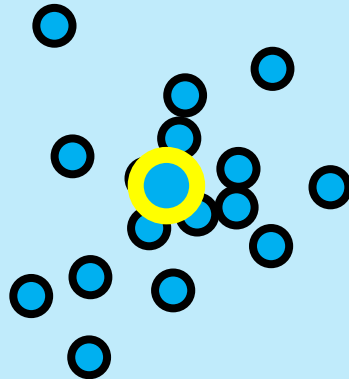
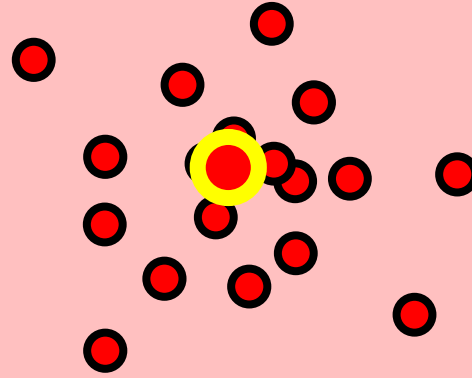
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

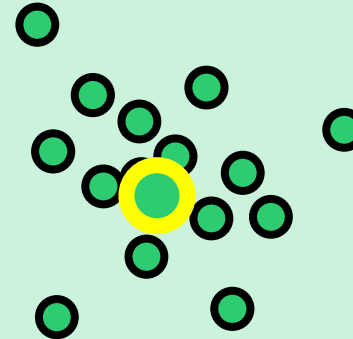
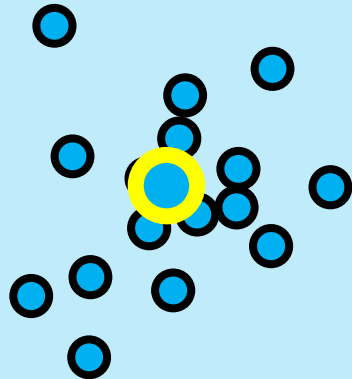
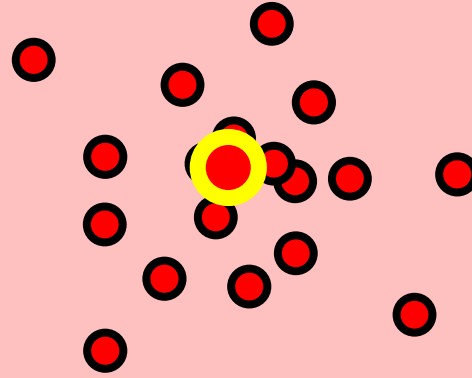
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k^t} \sum_{i:z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

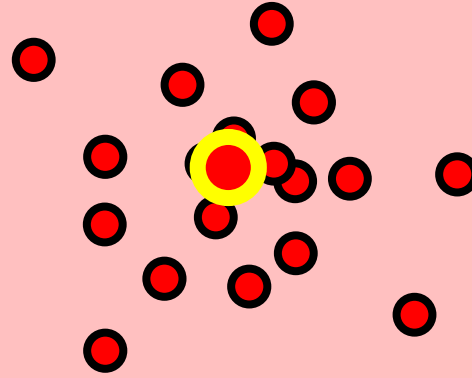
1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence



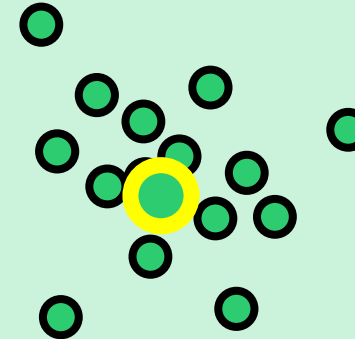
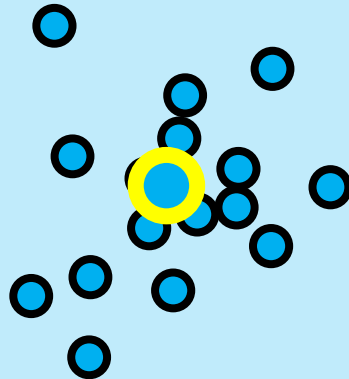
K-Means Algorithm in action!

K-MEANS/LLOYD'S ALGORITHM

1. Initialize means $\{\mu^{k,0}\}_{k=1\dots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mu^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \mu^{k,t}\|_2^2$
3. Update $\mu^{k,t+1} = \frac{1}{n_k} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
4. Repeat until convergence

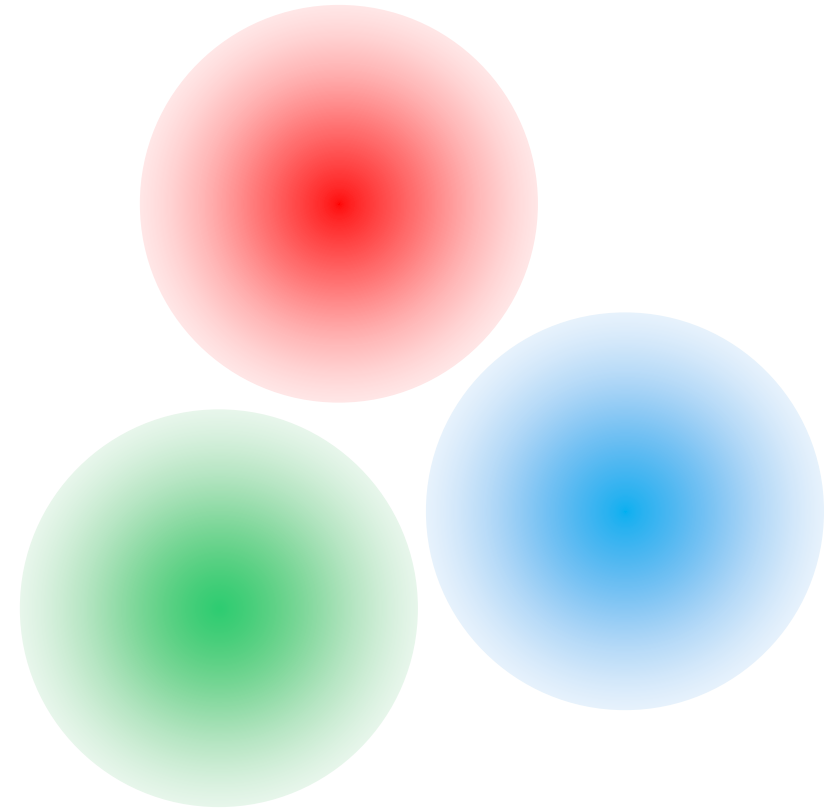


Stuck!!! ... but
at the global
optimum 😊

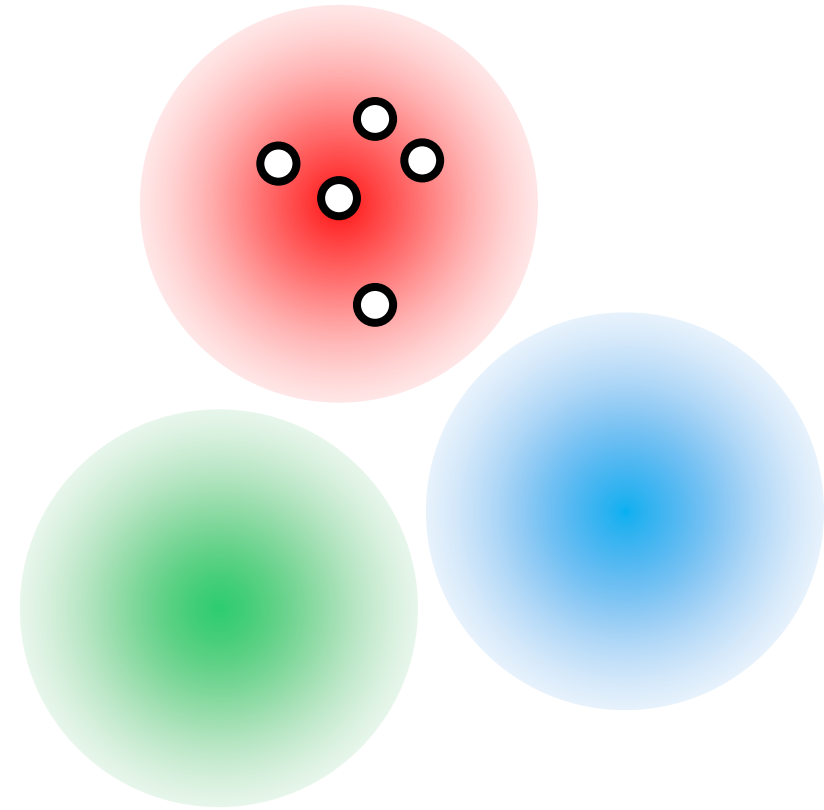


The Magic is not Always very Useful!

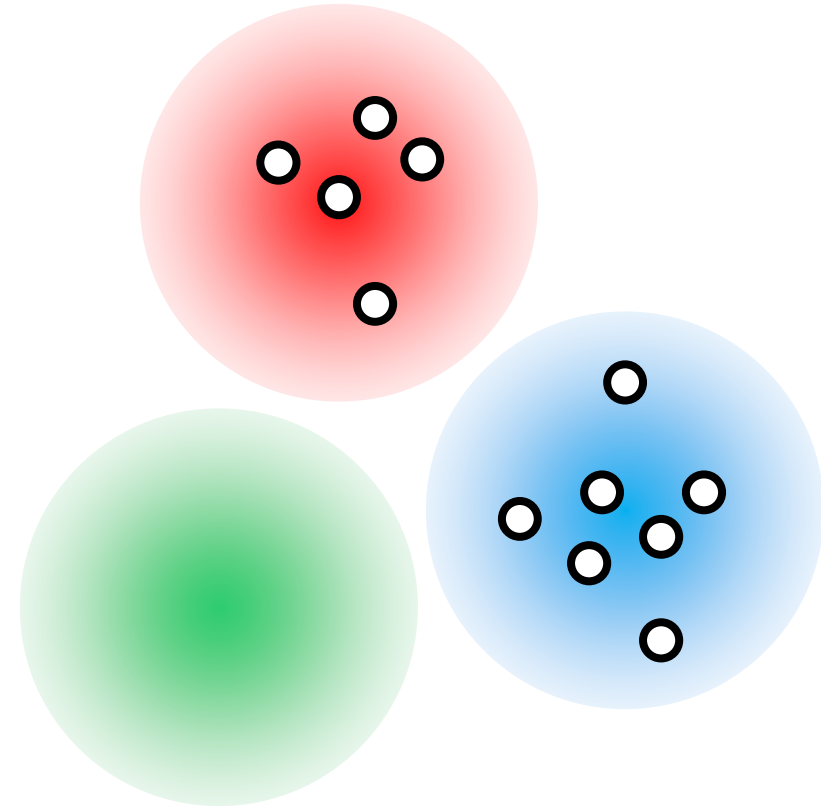
The Magic is not Always very Useful!



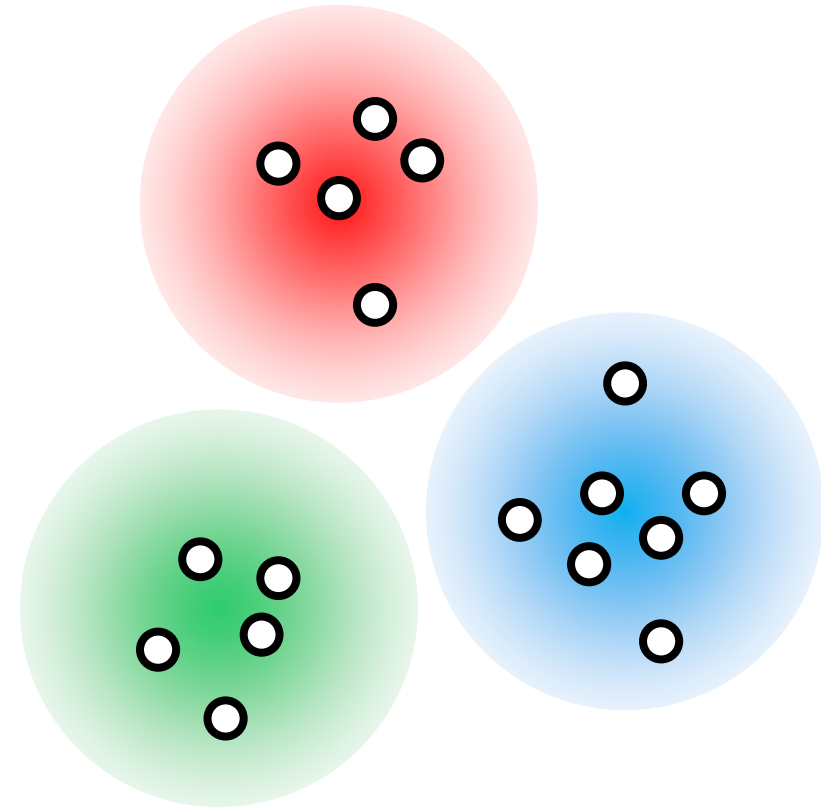
The Magic is not Always very Useful!



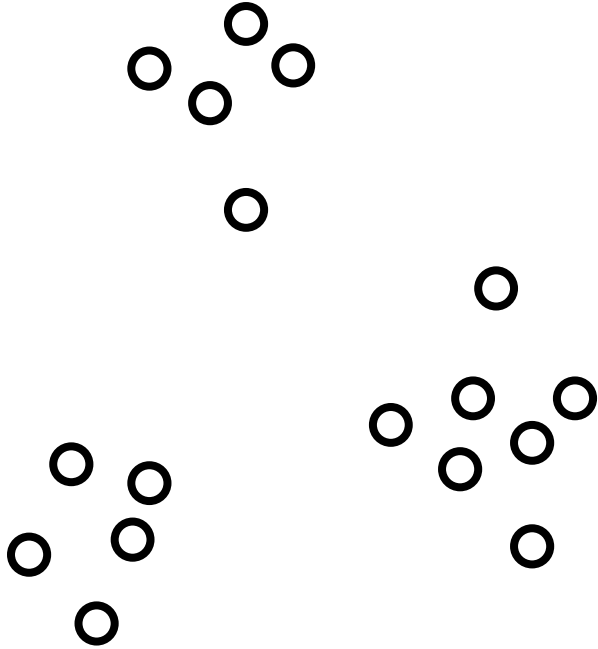
The Magic is not Always very Useful!



The Magic is not Always very Useful!

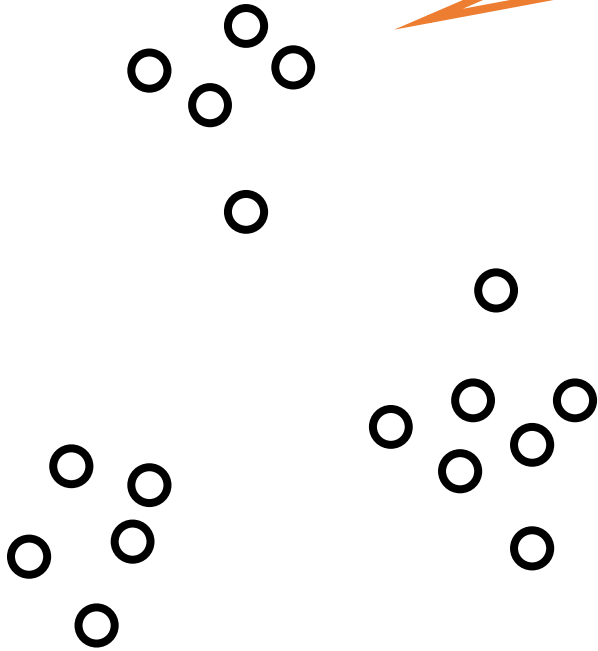


The Magic is not Always very Useful!



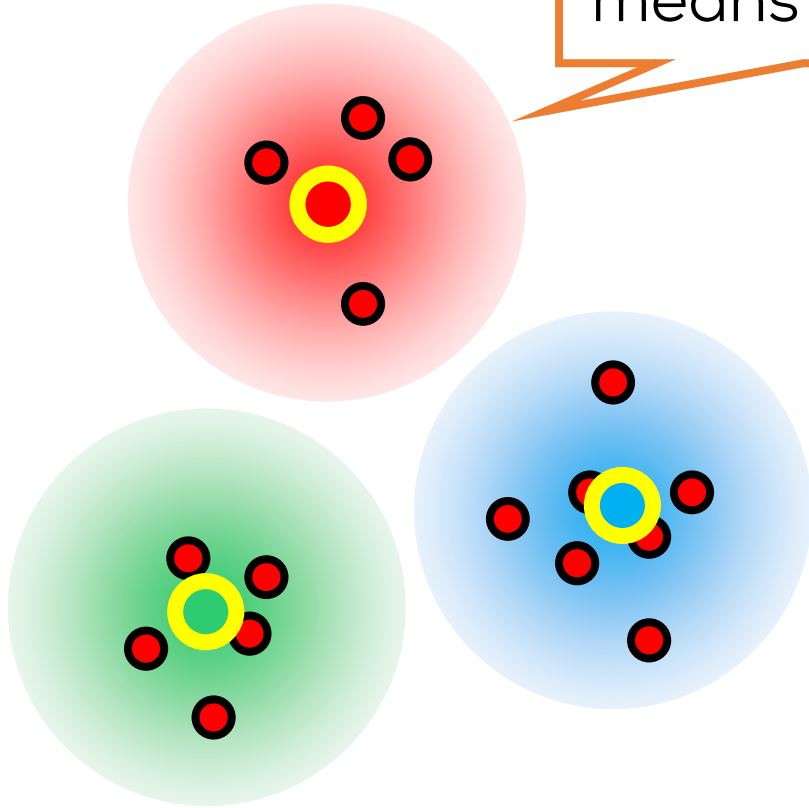
The Magic is not Always very Useful!

Apply the k-means algorithm



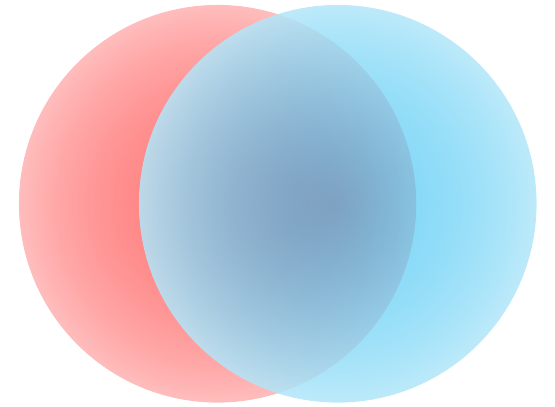
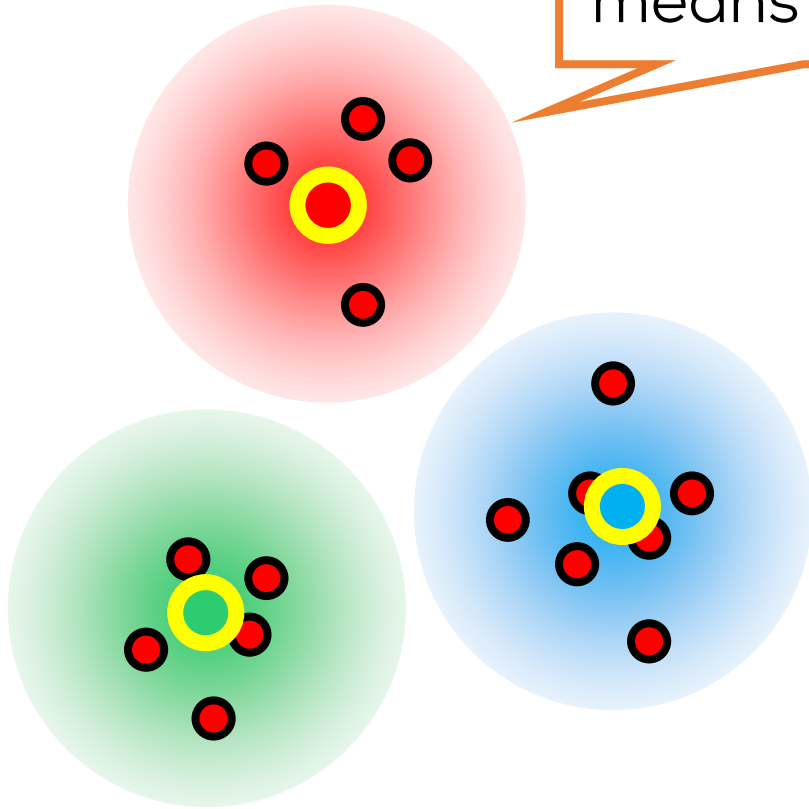
The Magic is not Always very Useful!

Apply the k-means algorithm



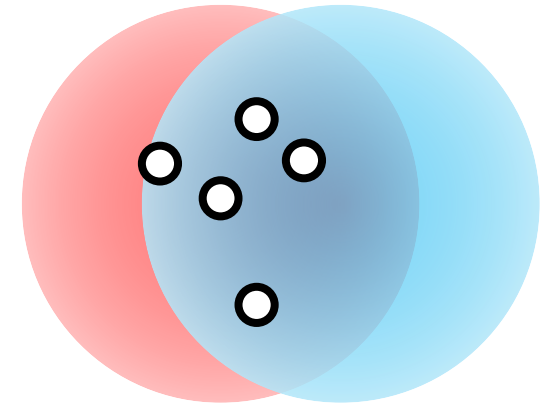
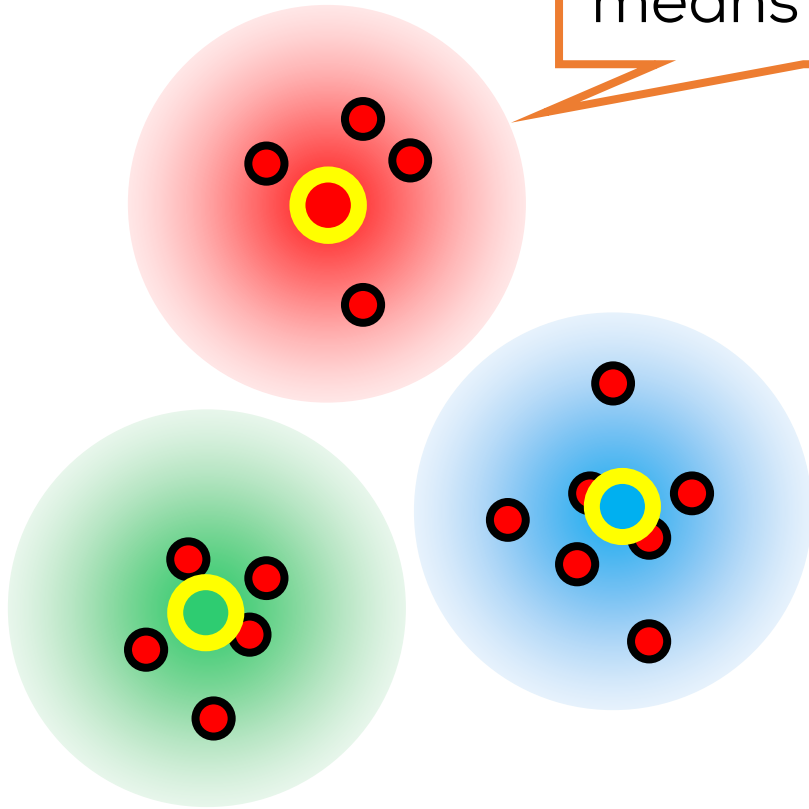
The Magic is not Always very Useful!

Apply the k-means algorithm



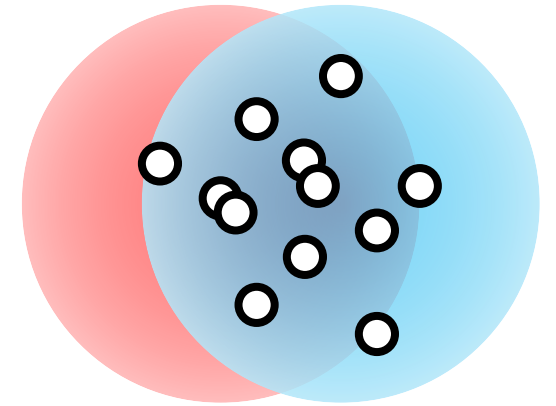
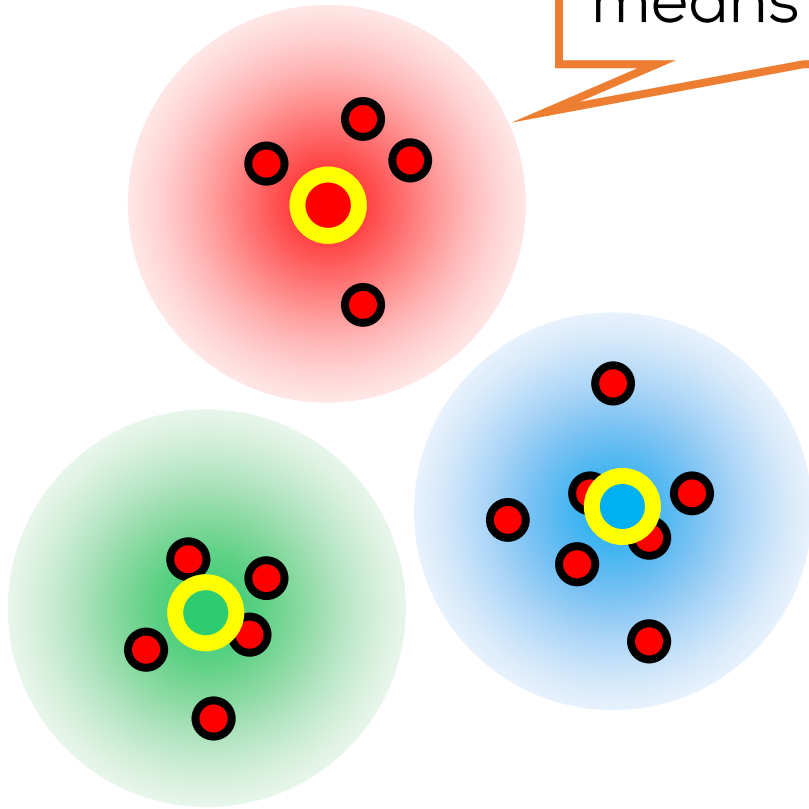
The Magic is not Always very Useful!

Apply the k-means algorithm



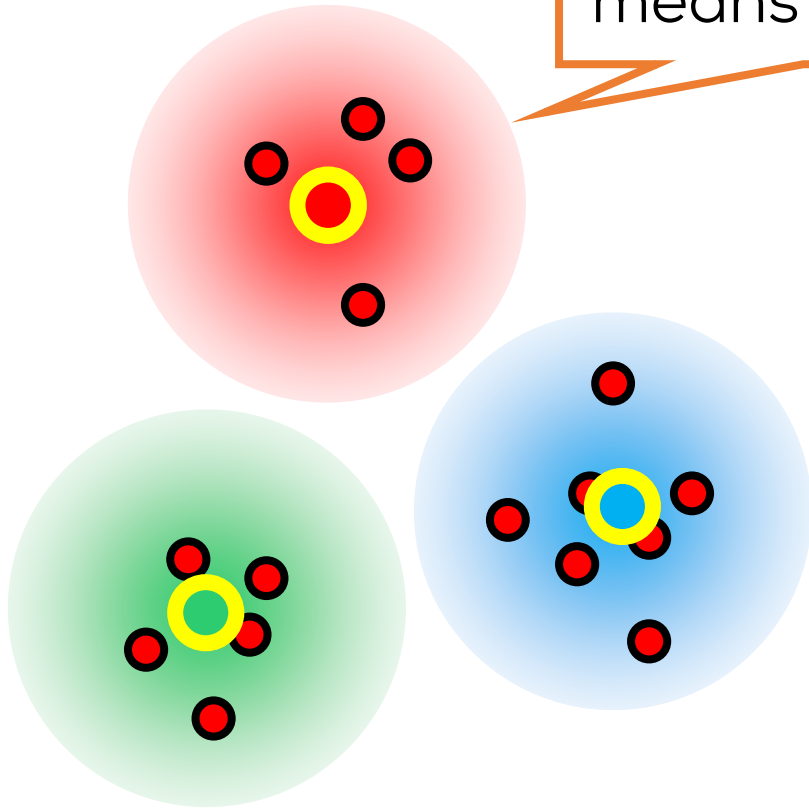
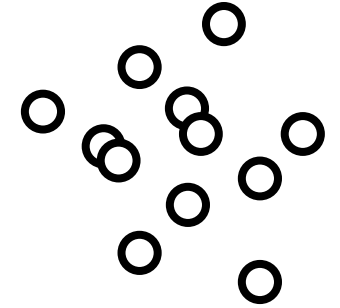
The Magic is not Always very Useful!

Apply the k-means algorithm



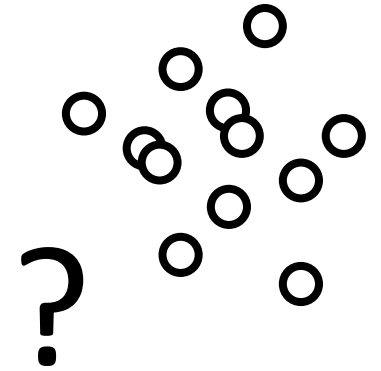
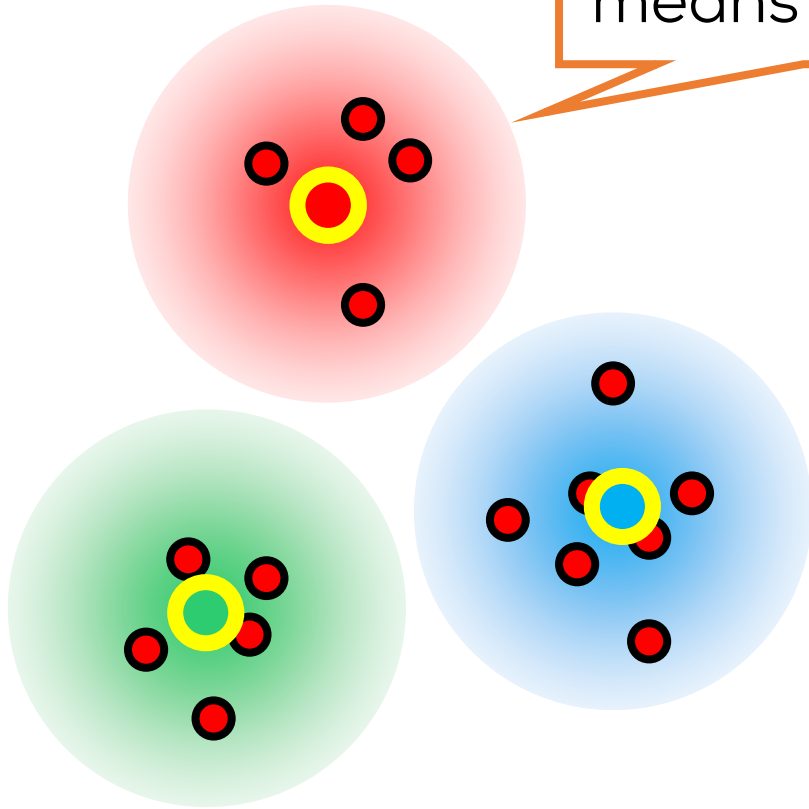
The Magic is not Always very Useful!

Apply the k-means algorithm



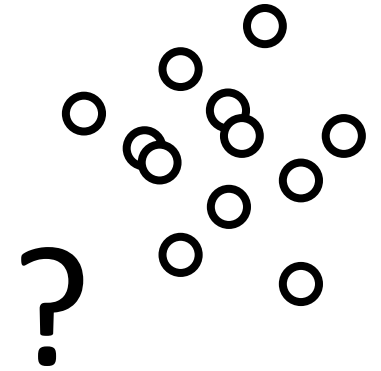
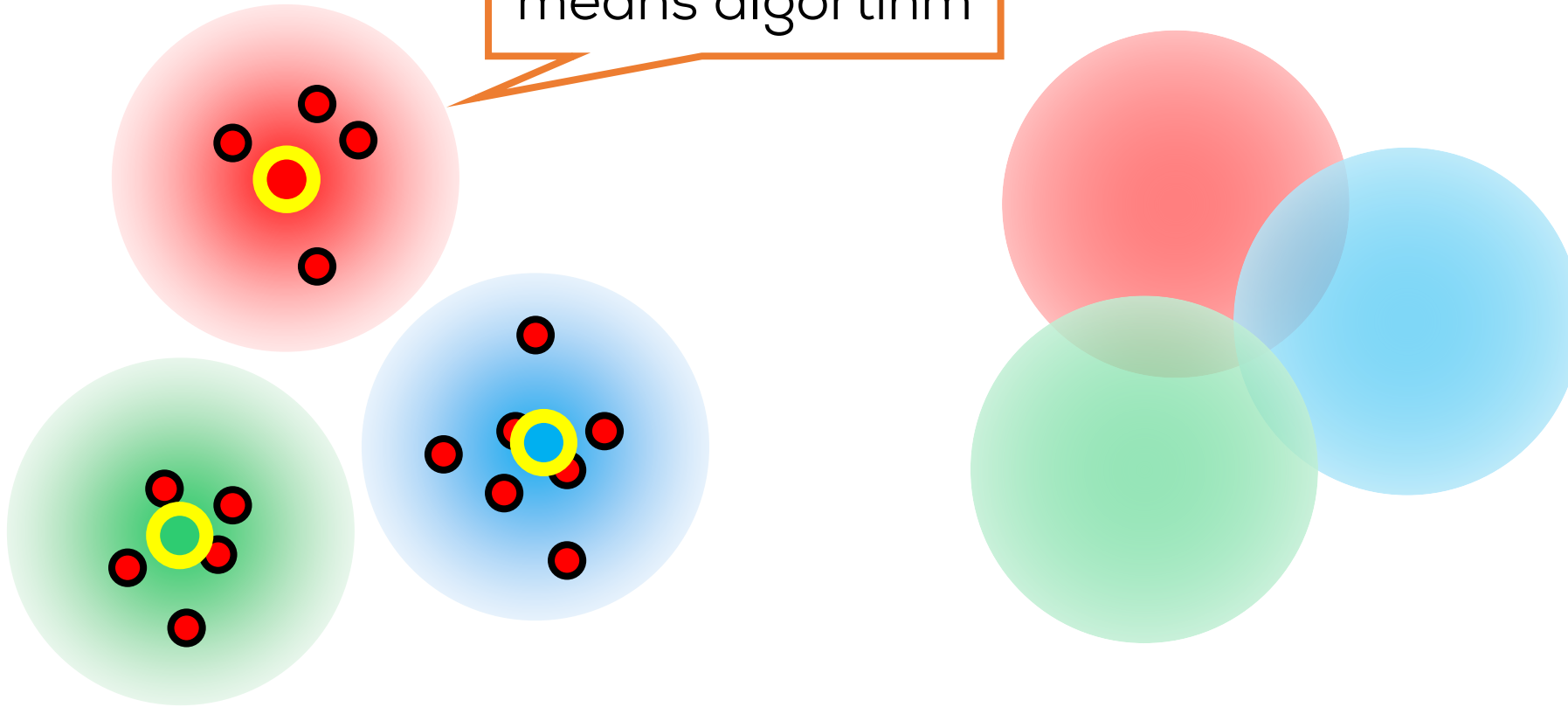
The Magic is not Always very Useful!

Apply the k-means algorithm



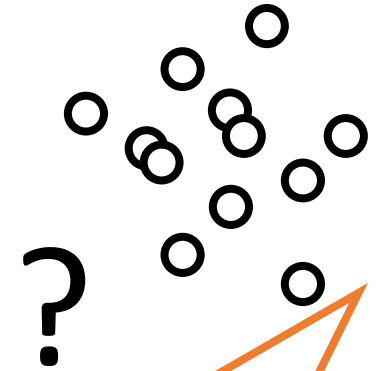
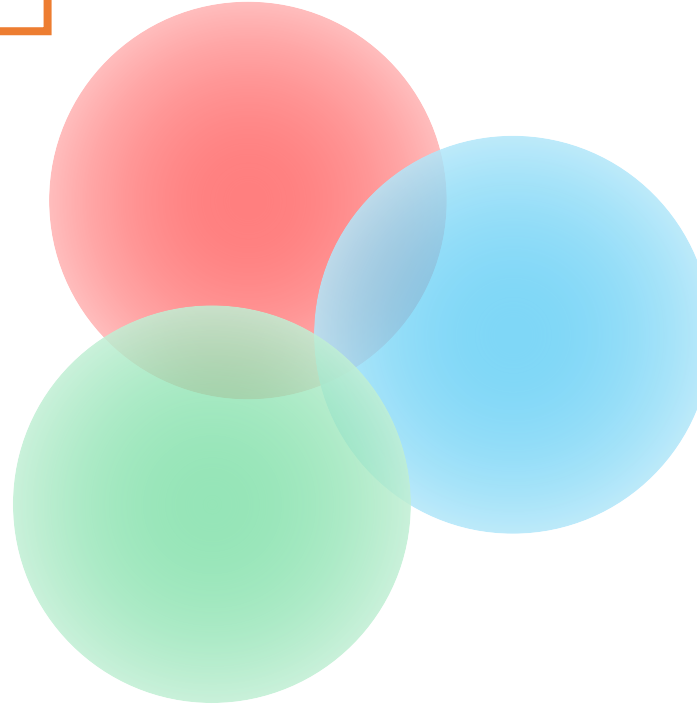
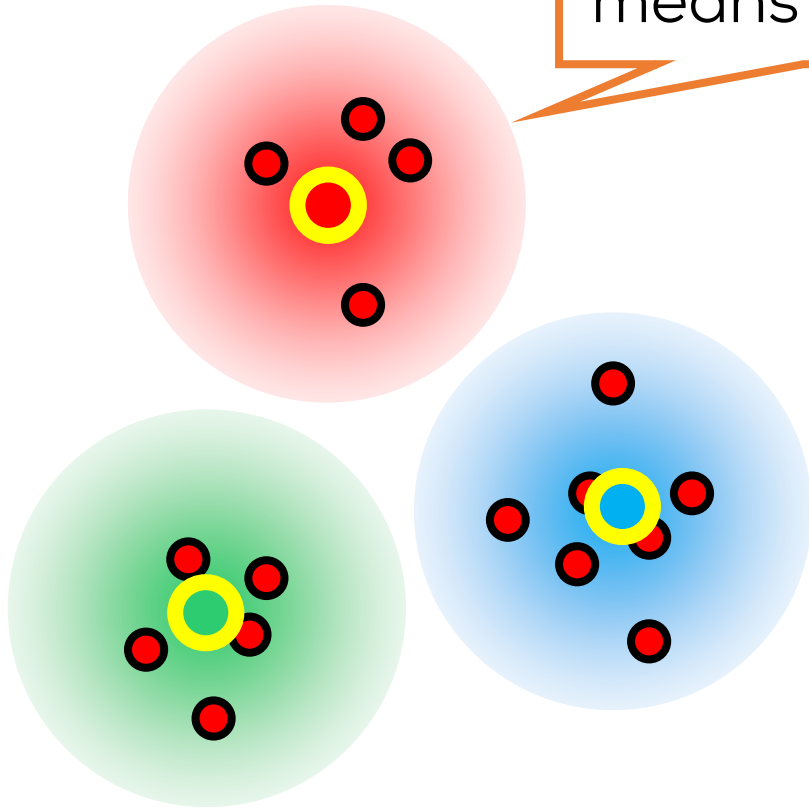
The Magic is not Always very Useful!

Apply the k-means algorithm



The Magic is not Always very Useful!

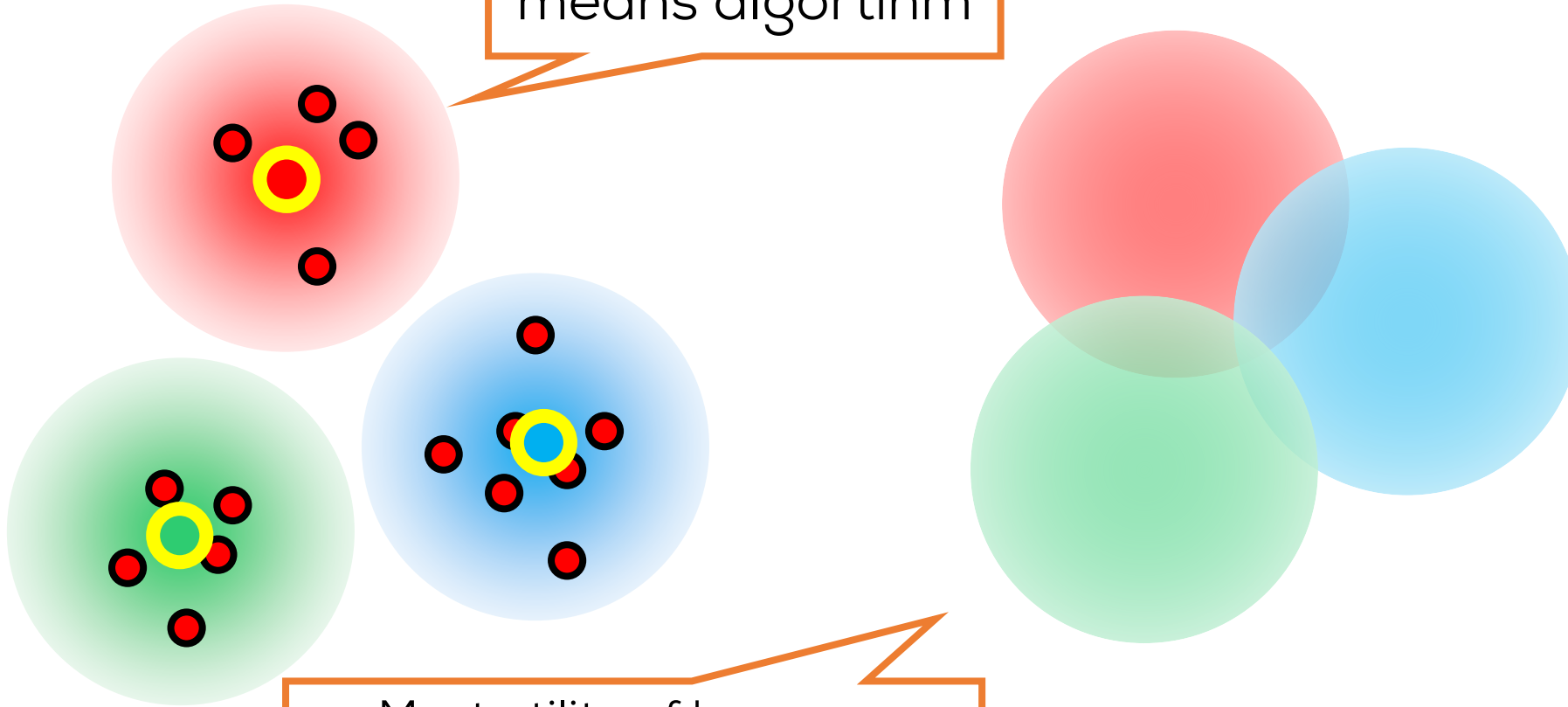
Apply the k-means algorithm



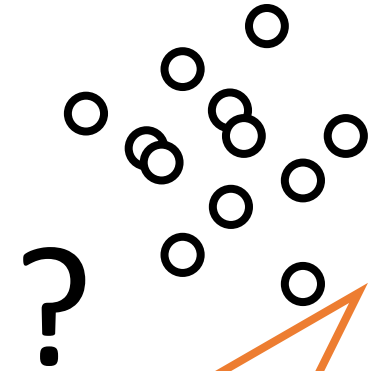
The problem is NP hard for a reason!

The Magic is not Always very Useful!

Apply the k-means algorithm



Most utility of k-means comes in problems that have some apparent structure



The problem is NP hard for a reason!

Soft Assignment

The EM algorithm

A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

ALTERNATING OPTIMIZATION

1. Initialize Θ^0
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
3. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

ALTERNATING OPTIMIZATION

1. Initialize Θ^0
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
3. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t]$$

A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

ALTERNATING OPTIMIZATION

1. Initialize Θ^0
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
3. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t]$$

$$\Theta^t = \left\{ \pi^t, \left\{ \mu^{1,t}, \mu^{2,t}, \mu^{3,t} \right\}, \left\{ \Sigma^{1,t}, \Sigma^{2,t}, \Sigma^{3,t} \right\} \right\}$$

A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

ALTERNATING OPTIMIZATION

1. Initialize Θ^0
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
3. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t] = \arg \max_{k \in [K]} \mathbb{P}[k | \Theta^t] \cdot \mathbb{P}[\mathbf{x}^i | k, \Theta^t]$$

$$\Theta^t = \left\{ \pi^t, \left\{ \mu^{1,t}, \mu^{2,t}, \mu^{3,t} \right\}, \left\{ \Sigma^{1,t}, \Sigma^{2,t}, \Sigma^{3,t} \right\} \right\}$$

A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

ALTERNATING OPTIMIZATION

1. Initialize Θ^0
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
3. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t] = \arg \max_{k \in [K]} \pi_k^t \cdot \mathbb{P}[\mathbf{x}^i | k, \Theta^t]$$

$$\Theta^t = \left\{ \pi^t, \left\{ \mu^{1,t}, \mu^{2,t}, \mu^{3,t} \right\}, \left\{ \Sigma^{1,t}, \Sigma^{2,t}, \Sigma^{3,t} \right\} \right\}$$

A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

ALTERNATING OPTIMIZATION

1. Initialize Θ^0
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
3. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$
4. Repeat until convergence

Bayes
Rule!

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t] = \arg \max_{k \in [K]} \pi_k^t \cdot \mathcal{N}(\mathbf{x}^i | \mu^{k,t}, \Sigma^{k,t})$$

$$\Theta^t = \left\{ \pi^t, \left\{ \mu^{1,t}, \mu^{2,t}, \mu^{3,t} \right\}, \left\{ \Sigma^{1,t}, \Sigma^{2,t}, \Sigma^{3,t} \right\} \right\}$$

A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

ALTERNATING OPTIMIZATION

1. Initialize Θ^0

2. For $i \in [n]$ update $z^{i,t}$ using Θ^t

May be throwing away a lot of information!

$$z^{i,t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$$

Repeat until convergence

Bayes Rule!

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t] = \arg \max_{k \in [K]} \pi_k^t \cdot \mathcal{N}(\mathbf{x}^i | \mu^{k,t}, \Sigma^{k,t})$$

$$\Theta^t = \left\{ \pi^t, \left\{ \mu^{1,t}, \mu^{2,t}, \mu^{3,t} \right\}, \left\{ \Sigma^{1,t}, \Sigma^{2,t}, \Sigma^{3,t} \right\} \right\}$$

A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P} [X \mid \Theta]$$

$$\begin{aligned} \text{E.g. } \mathbb{P}[\text{red} | \mathbf{x}^i, \Theta^t] &= 0.5, \\ \mathbb{P}[\text{blue} | \mathbf{x}^i, \Theta^t] &= 0.4, \\ \mathbb{P}[\text{green} | \mathbf{x}^i, \Theta^t] &= 0.1, \end{aligned}$$

ERNATING OPTIMIZATION

Realize Θ^0

2 For $i \in [n]$ update $z^{i,t}$ using Θ^t

May be throwing away a lot of information!

$$z^{i,t} = \arg \max_{\Theta} \mathbb{P} [X, \{z^{i,t}\} \mid \Theta]$$

Bayes Rule!

information! \rightarrow still convergence

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P} [k \mid \mathbf{x}^i, \Theta^t] = \arg \max_{k \in [K]} \pi_k^t \cdot \mathcal{N}(\mathbf{x}^i \mid \boldsymbol{\mu}^{k,t}, \Sigma^{k,t})$$

$$\Theta^t = \left\{ \pi^t, \left\{ \mu^{1,t}, \mu^{2,t}, \mu^{3,t} \right\}, \left\{ \Sigma^{1,t}, \Sigma^{2,t}, \Sigma^{3,t} \right\} \right\}$$

A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

E.g. $\mathbb{P}[\text{red} | \mathbf{x}^i, \Theta^t] = 0.5,$
 $\mathbb{P}[\text{blue} | \mathbf{x}^i, \Theta^t] = 0.4,$
 $\mathbb{P}[\text{green} | \mathbf{x}^i, \Theta^t] = 0.1,$

Can we use $\mathbb{P}[k | \mathbf{x}^i, \Theta^t]$ as weights instead?

For $i \in [n]$ update $z^{i,t}$ using Θ^t

May be throwing away a lot of information!

$$z^{i,t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$$

Bayes Rule!

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t] = \arg \max_{k \in [K]} \pi_k^t \cdot \mathcal{N}(\mathbf{x}^i | \mu^{k,t}, \Sigma^{k,t})$$

$$\Theta^t = \{ \pi^t, \{ \mu^{1,t}, \mu^{2,t}, \mu^{3,t} \}, \{ \Sigma^{1,t}, \Sigma^{2,t}, \Sigma^{3,t} \} \}$$

A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

Assign point \mathbf{x}^i to cluster k with weight $\propto \mathbb{P}[k | \mathbf{x}^i, \Theta^t]$

E.g. $\mathbb{P}[\text{red} | \mathbf{x}^i, \Theta^t] = 0.5$,
 $\mathbb{P}[\text{blue} | \mathbf{x}^i, \Theta^t] = 0.4$,
 $\mathbb{P}[\text{green} | \mathbf{x}^i, \Theta^t] = 0.1$,

Can we use $\mathbb{P}[k | \mathbf{x}^i, \Theta^t]$ as weights instead?

EMERNA

ATION

For $i \in [n]$ update $z^{i,t}$ using Θ^t

May be throwing away a lot of information!

$$z^{i,t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$$

Bayes Rule!

$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t] = \arg \max_{k \in [K]} \pi_k^t \cdot \mathcal{N}(\mathbf{x}^i | \mu^{k,t}, \Sigma^{k,t})$$

$$\Theta^t = \left\{ \pi^t, \left\{ \mu^{1,t}, \mu^{2,t}, \mu^{3,t} \right\}, \left\{ \Sigma^{1,t}, \Sigma^{2,t}, \Sigma^{3,t} \right\} \right\}$$

A Ray of Hope

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X | \Theta]$$

E.g. $\mathbb{P}[\text{red} | \mathbf{x}^i, \Theta^t] = 0.5$,
 $\mathbb{P}[\text{blue} | \mathbf{x}^i, \Theta^t] = 0.4$,
 $\mathbb{P}[\text{green} | \mathbf{x}^i, \Theta^t] = 0.1$,

Assign point \mathbf{x}^i to cluster k with weight $\propto \mathbb{P}[k | \mathbf{x}^i, \Theta^t]$

Can we use $\mathbb{P}[k | \mathbf{x}^i, \Theta^t]$ as weights instead?

Has a "Bayesian" feel to it – use all available posterior information

For $i \in [n]$ update $z^{i,t}$ using Θ^t

May be throwing away a lot of information!

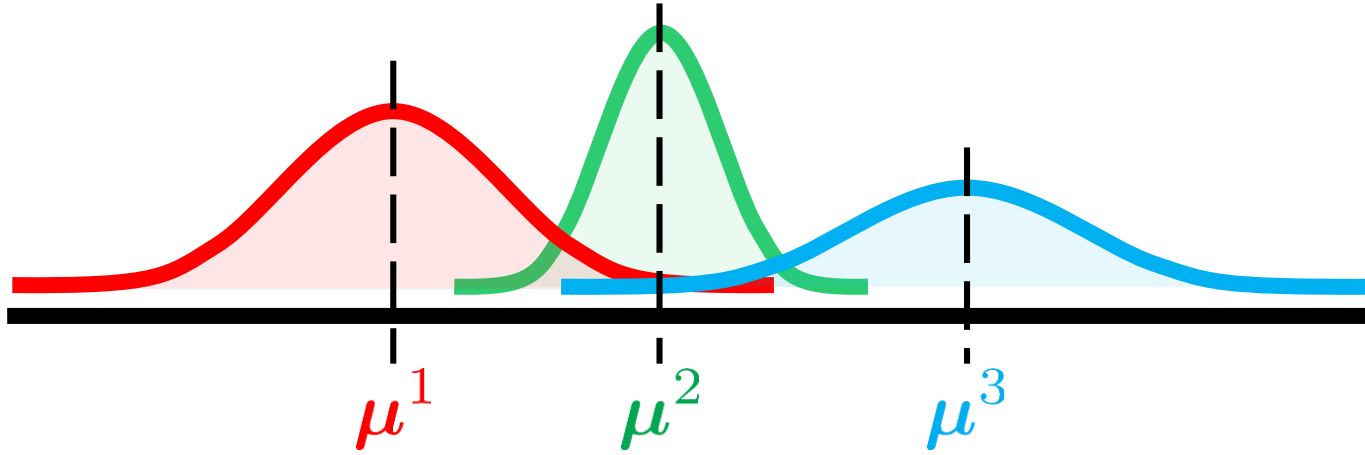
$$z^{i,t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} | \Theta]$$

Bayes Rule!

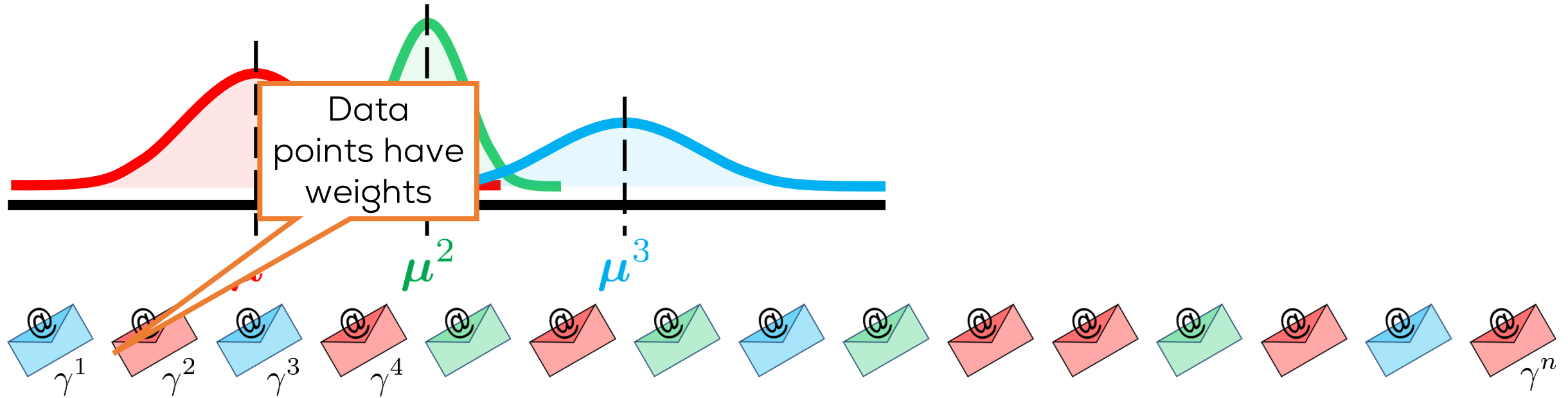
$$z^{i,t} = \arg \max_{k \in [K]} \mathbb{P}[k | \mathbf{x}^i, \Theta^t] = \arg \max_{k \in [K]} \pi_k^t \cdot \mathcal{N}(\mathbf{x}^i | \mu^{k,t}, \Sigma^{k,t})$$

$$\Theta^t = \{ \pi^t, \{ \mu^{1,t}, \mu^{2,t}, \mu^{3,t} \}, \{ \Sigma^{1,t}, \Sigma^{2,t}, \Sigma^{3,t} \} \}$$

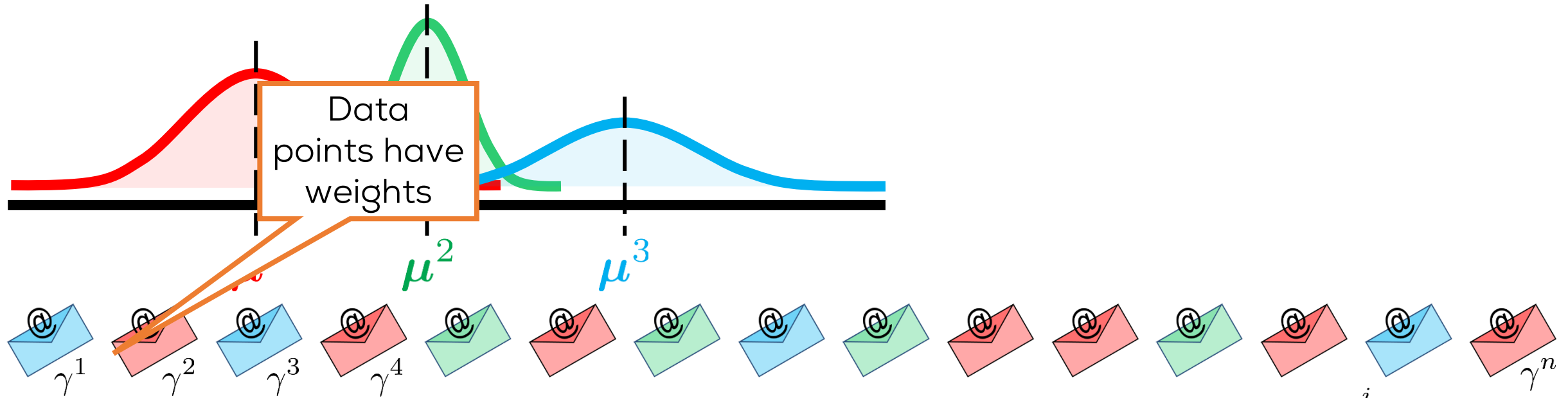
Weighted MLE



Weighted MLE

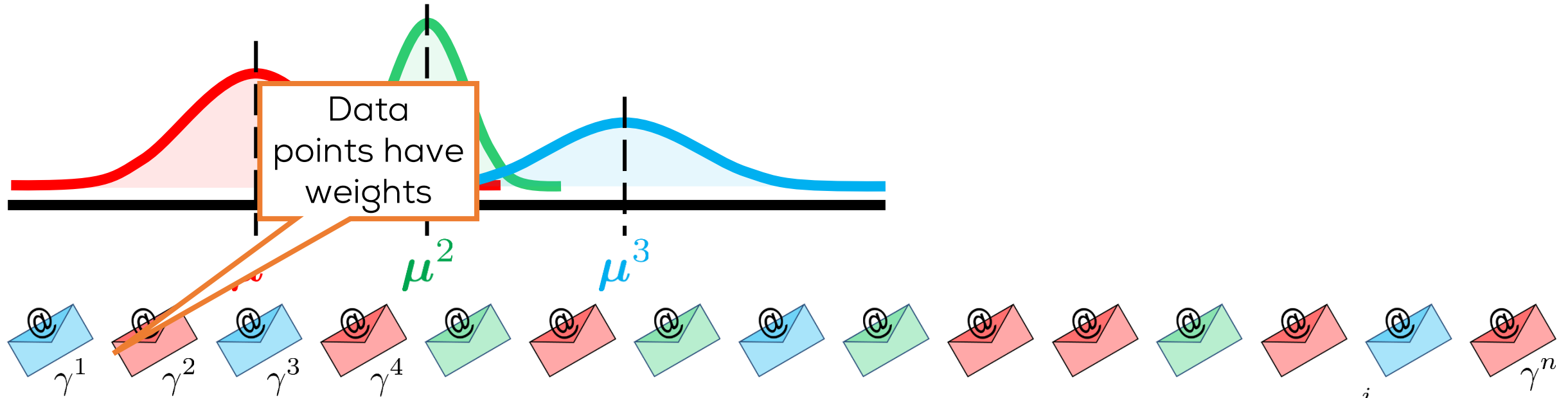


Weighted MLE



$$\mathbb{P} [\mathbf{x}^i, z^i, \gamma^i \mid \Theta] = \left(\mathbb{P} [z^i \mid \Theta] \cdot \mathbb{P} [\mathbf{x}^i \mid z^i, \Theta] \right)^{\gamma^i} = \left(\pi_{z^i} \cdot \mathcal{N}(\mathbf{x}^i \mid \boldsymbol{\mu}^{z^i}, \Sigma^{z^i}) \right)^{\gamma^i}$$

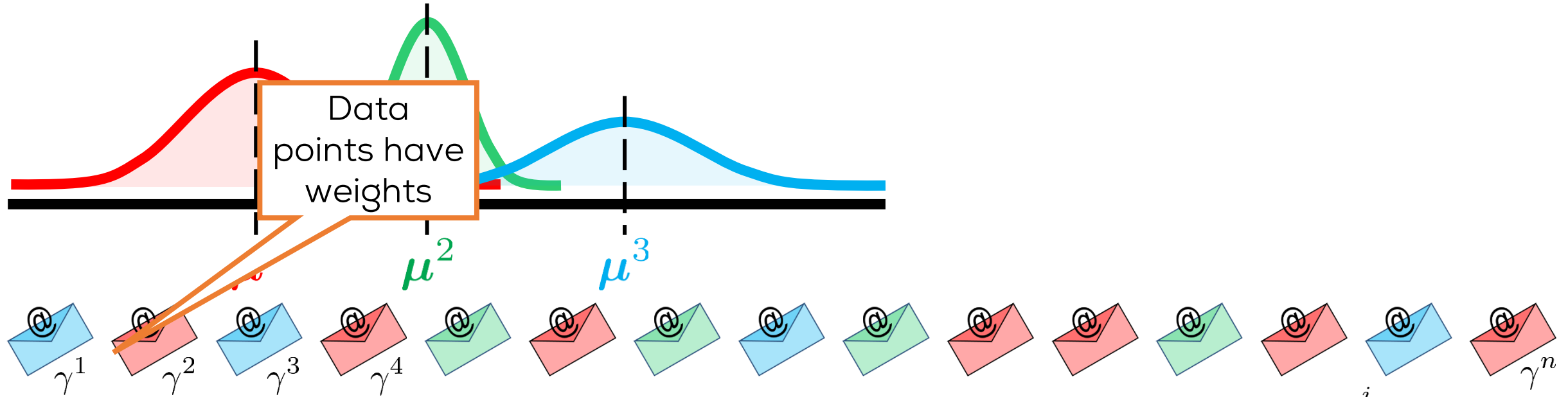
Weighted MLE



$$\mathbb{P} [\mathbf{x}^i, z^i, \gamma^i \mid \Theta] = \left(\mathbb{P} [z^i \mid \Theta] \cdot \mathbb{P} [\mathbf{x}^i \mid z^i, \Theta] \right)^{\gamma^i} = \left(\pi_{z^i} \cdot \mathcal{N}(\mathbf{x}^i \mid \boldsymbol{\mu}^{z^i}, \Sigma^{z^i}) \right)^{\gamma^i}$$

$$\mathbb{P} [X, \{z^i\}, \{\gamma^i\} \mid \Theta] = \prod_{i=1}^n \mathbb{P} [\mathbf{x}^i, z^i, \gamma^i \mid \Theta]$$

Weighted MLE

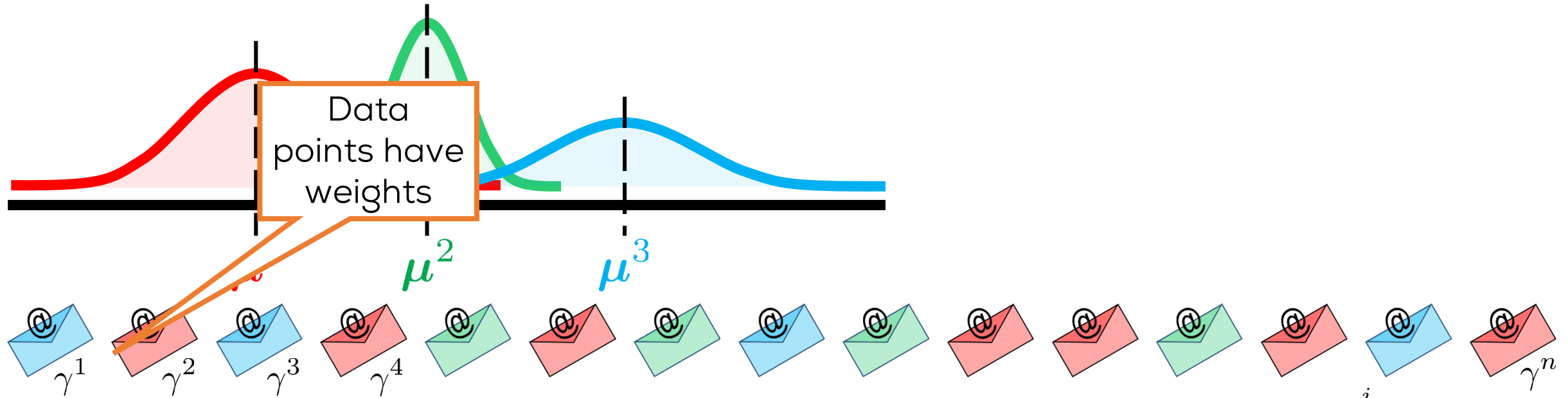


$$\mathbb{P} [\mathbf{x}^i, z^i, \gamma^i \mid \Theta] = \left(\mathbb{P} [z^i \mid \Theta] \cdot \mathbb{P} [\mathbf{x}^i \mid z^i, \Theta] \right)^{\gamma^i} = \left(\pi_{z^i} \cdot \mathcal{N}(\mathbf{x}^i \mid \mu^{z^i}, \Sigma^{z^i}) \right)^{\gamma^i}$$

$$\mathbb{P} [X, \{z^i\}, \{\gamma^i\} \mid \Theta] = \prod_{i=1}^n \mathbb{P} [\mathbf{x}^i, z^i, \gamma^i \mid \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P} [X, \{z^i\}, \{\gamma^i\} \mid \Theta]$$

Weighted MLE



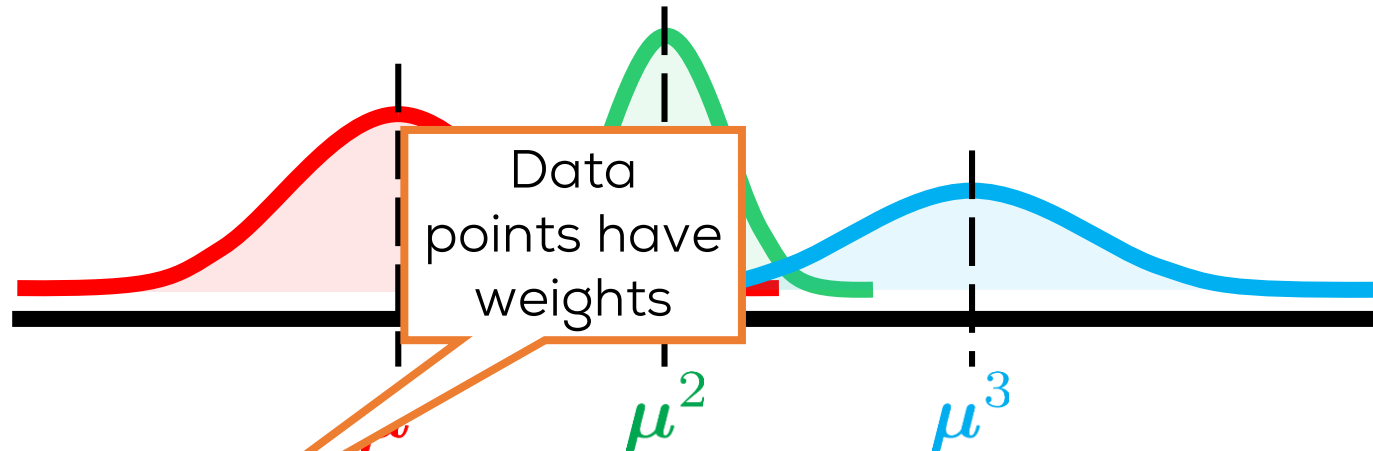
$$\mathbb{P} [\mathbf{x}^i, z^i, \gamma^i \mid \Theta] = \left(\mathbb{P} [z^i \mid \Theta] \cdot \mathbb{P} [\mathbf{x}^i \mid z^i, \Theta] \right)^{\gamma^i} = \left(\pi_{z^i} \cdot \mathcal{N}(\mathbf{x}^i \mid \boldsymbol{\mu}^{z^i}, \Sigma^{z^i}) \right)^{\gamma^i}$$

$$\mathbb{P} [X, \{z^i\}, \{\gamma^i\} \mid \Theta] = \prod_{i=1}^n \mathbb{P} [\mathbf{x}^i, z^i, \gamma^i \mid \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P} [X, \{z^i\}, \{\gamma^i\} \mid \Theta]$$

Take log and apply 1st order optimality

Weighted MLE



$$\mathbb{P} [\mathbf{x}^i, z^i, \gamma^i \mid \Theta] = \left(\mathbb{P} [z^i \mid \Theta] \cdot \mathbb{P} [\mathbf{x}^i \mid z^i, \Theta] \right)^{\gamma^i} = \left(\pi_{z^i} \cdot \mathcal{N}(\mathbf{x}^i \mid \boldsymbol{\mu}^{z^i}, \Sigma^{z^i}) \right)^{\gamma^i}$$

$$\mathbb{P} [X, \{z^i\}, \{\gamma^i\} \mid \Theta] = \prod_{i=1}^n \mathbb{P} [\mathbf{x}^i, z^i, \gamma^i \mid \Theta]$$

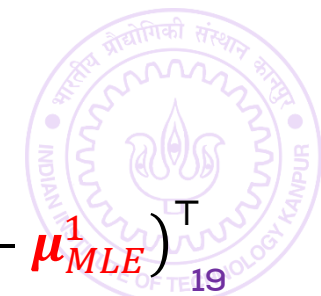
$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P} [X, \{z^i\}, \{\gamma^i\} \mid \Theta]$$

Take log and apply 1st order optimality

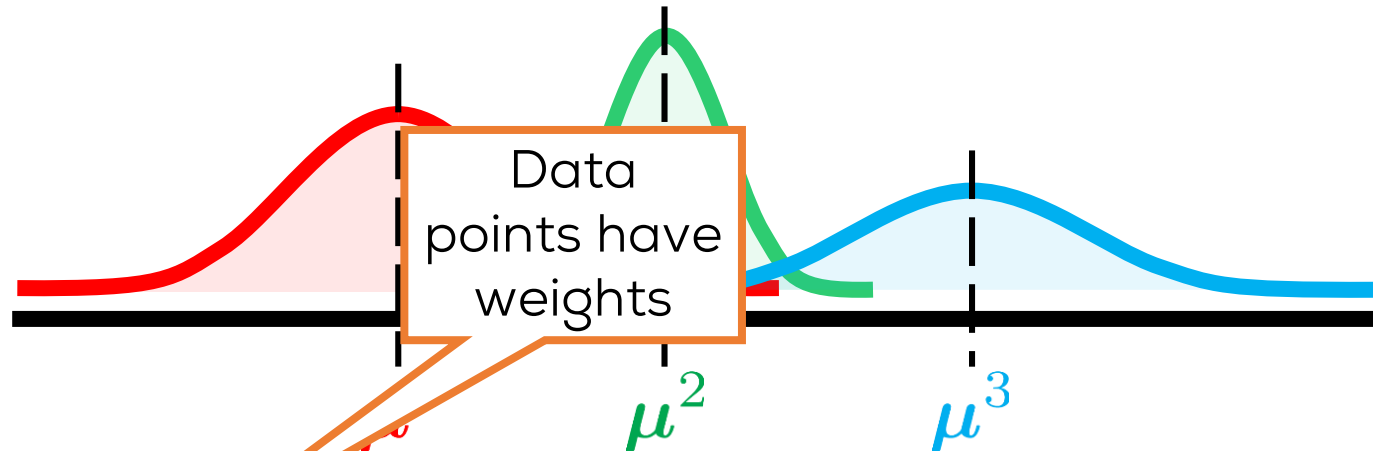
$$\pi_1^{\text{MLE}} = \frac{\# \text{ eff. red emails}}{\# \text{ eff. total emails}} = \frac{\tilde{n}_r}{n} = \frac{\sum_{i: z^i = \bullet} \gamma^i}{\sum_{j=1}^n \gamma^j}$$

$$\boldsymbol{\mu}_{\text{MLE}}^1 = \frac{1}{\tilde{n}_r} \sum_{i: z^i = \bullet} \gamma^i \cdot \mathbf{x}^i$$

$$\Sigma_{\text{MLE}}^1 = \frac{1}{\tilde{n}_r} \sum_{i: z^i = \bullet} \gamma^i \cdot (\mathbf{x}^i - \boldsymbol{\mu}_{\text{MLE}}^1)(\mathbf{x}^i - \boldsymbol{\mu}_{\text{MLE}}^1)^{\top}$$



Weighted MLE



$$\mathbb{P}[\mathbf{x}^i, z^i, \gamma^i | \Theta] = (\mathbb{P}[z^i | \Theta] \cdot \mathbb{P}[\mathbf{x}^i | z^i, \Theta])^{\gamma^i} = \left(\pi_{z^i} \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^{z^i}, \Sigma^{z^i}) \right)^{\gamma^i}$$

$$\mathbb{P}[X, \{z^i\}, \{\gamma^i\} | \Theta] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i, z^i, \gamma^i | \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X, \{z^i\}, \{\gamma^i\} | \Theta]$$

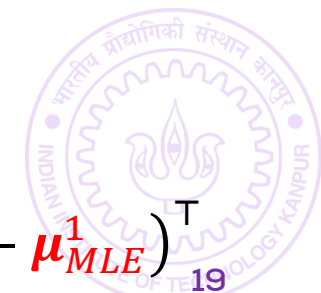
Take log and apply 1st order optimality

Exercise!

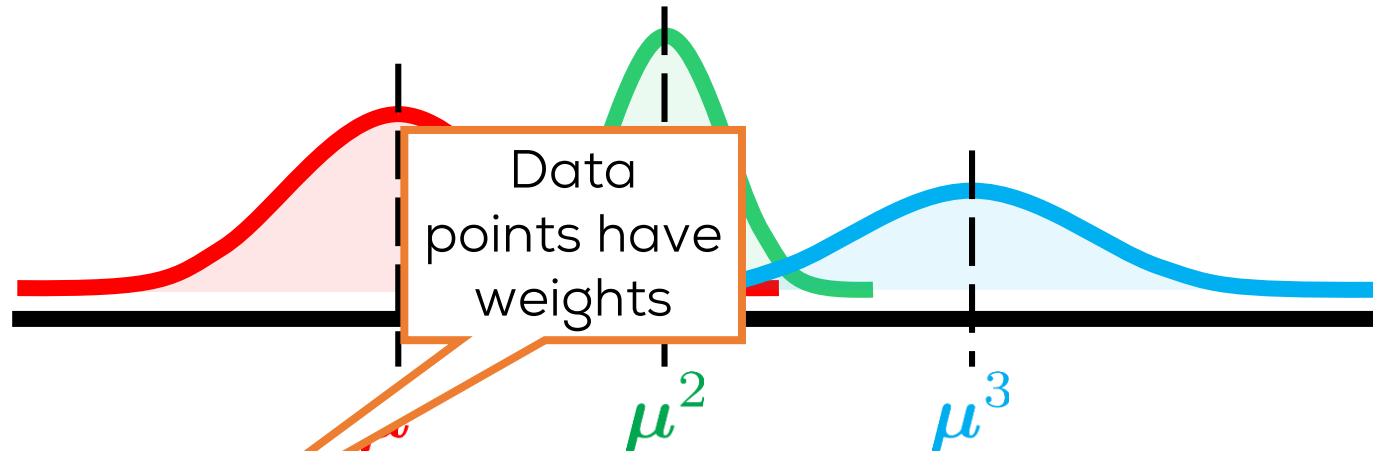
$$\pi_1^{\text{MLE}} = \frac{\# \text{ eff. red emails}}{\# \text{ eff. total emails}} = \frac{\tilde{n}_r}{n} = \frac{\sum_{i: z^i = \bullet} \gamma^i}{\sum_{j=1}^n \gamma^j}$$

$$\boldsymbol{\mu}_1^{\text{MLE}} = \frac{1}{\tilde{n}_r} \sum_{i: z^i = \bullet} \gamma^i \cdot \mathbf{x}^i$$

$$\Sigma_1^{\text{MLE}} = \frac{1}{\tilde{n}_r} \sum_{i: z^i = \bullet} \gamma^i \cdot (\mathbf{x}^i - \boldsymbol{\mu}_1^{\text{MLE}})(\mathbf{x}^i - \boldsymbol{\mu}_1^{\text{MLE}})^{\top}$$



Weighted MLE



Reduces to normal MLE if $\gamma^i = 1$ for all i



$$\mathbb{P}[\mathbf{x}^i, z^i, \gamma^i | \Theta] = (\mathbb{P}[z^i | \Theta] \cdot \mathbb{P}[\mathbf{x}^i | z^i, \Theta])^{\gamma^i} = \left(\pi_{z^i} \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^{z^i}, \Sigma^{z^i}) \right)^{\gamma^i}$$

$$\mathbb{P}[X, \{z^i\}, \{\gamma^i\} | \Theta] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i, z^i, \gamma^i | \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X, \{z^i\}, \{\gamma^i\} | \Theta]$$

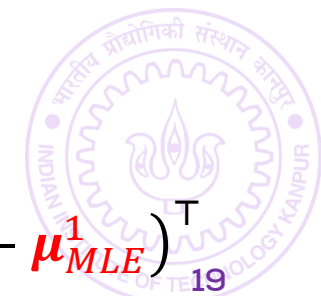
Take log and apply 1st order optimality

Exercise!

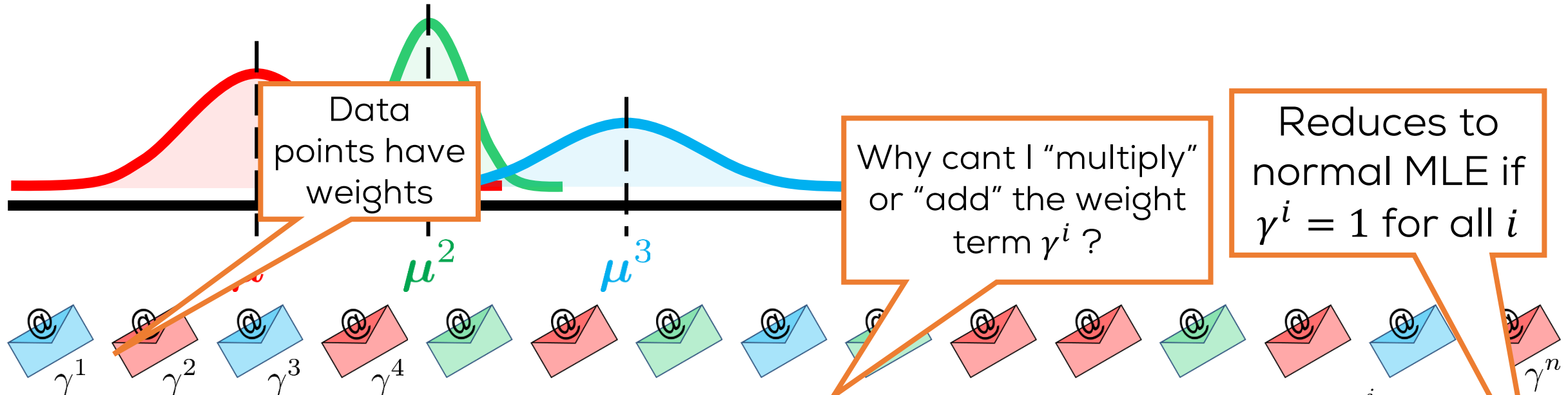
$$\pi_1^{\text{MLE}} = \frac{\# \text{ eff. red emails}}{\# \text{ eff. total emails}} = \frac{\tilde{n}_r}{n} = \frac{\sum_{i: z^i = \bullet} \gamma^i}{\sum_{j=1}^n \gamma^j}$$

$$\boldsymbol{\mu}_{\text{MLE}}^1 = \frac{1}{\tilde{n}_r} \sum_{i: z^i = \bullet} \gamma^i \cdot \mathbf{x}^i$$

$$\Sigma_{\text{MLE}}^1 = \frac{1}{\tilde{n}_r} \sum_{i: z^i = \bullet} \gamma^i \cdot (\mathbf{x}^i - \boldsymbol{\mu}_{\text{MLE}}^1)(\mathbf{x}^i - \boldsymbol{\mu}_{\text{MLE}}^1)^{\top}$$



Weighted MLE



$$\mathbb{P}[\mathbf{x}^i, z^i, \gamma^i | \Theta] = (\mathbb{P}[z^i | \Theta] \cdot \mathbb{P}[\mathbf{x}^i | z^i, \Theta])^{\gamma^i} = \left(\pi_{z^i} \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^{z^i}, \Sigma^{z^i}) \right)^{\gamma^i}$$

$$\mathbb{P}[X, \{z^i\}, \{\gamma^i\} | \Theta] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i, z^i, \gamma^i | \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X, \{z^i\}, \{\gamma^i\} | \Theta]$$

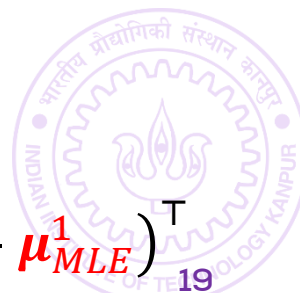
Take log and apply 1st order optimality

Exercise!

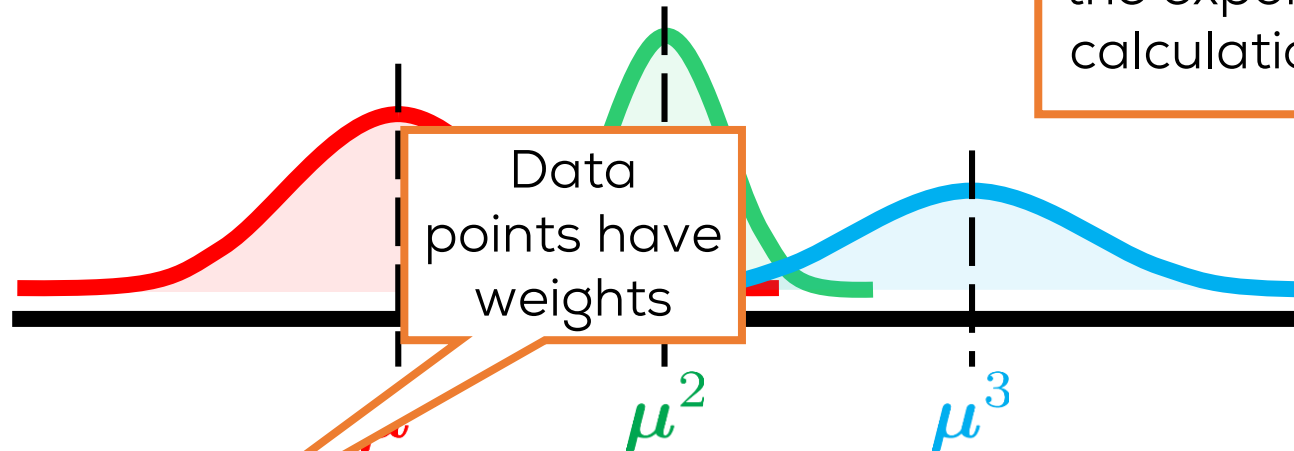
$$\pi_1^{\text{MLE}} = \frac{\# \text{ eff. red emails}}{\# \text{ eff. total emails}} = \frac{\tilde{n}_r}{n} = \frac{\sum_{i: z^i = \bullet} \gamma^i}{\sum_{j=1}^n \gamma^j}$$

$$\boldsymbol{\mu}_{\text{MLE}}^1 = \frac{1}{\tilde{n}_r} \sum_{i: z^i = \bullet} \gamma^i \cdot \mathbf{x}^i$$

$$\Sigma_{\text{MLE}}^1 = \frac{1}{\tilde{n}_r} \sum_{i: z^i = \bullet} \gamma^i \cdot (\mathbf{x}^i - \boldsymbol{\mu}_{\text{MLE}}^1)(\mathbf{x}^i - \boldsymbol{\mu}_{\text{MLE}}^1)^{\top}$$



Weighted MLE



$$\mathbb{P}[\mathbf{x}^i, z^i, \gamma^i | \Theta] = (\mathbb{P}[z^i | \Theta] \cdot \mathbb{P}[\mathbf{x}^i | z^i, \Theta])^{\gamma^i} = \left(\pi_{z^i} \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^{z^i}, \Sigma^{z^i}) \right)^{\gamma^i}$$

$$\mathbb{P}[X, \{z^i\}, \{\gamma^i\} | \Theta] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i, z^i, \gamma^i | \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X, \{z^i\}, \{\gamma^i\} | \Theta]$$

Take log and apply 1st order optimality

Exercise!

$$\pi_1^{\text{MLE}} = \frac{\# \text{ eff. red emails}}{\# \text{ eff. total emails}} = \frac{\tilde{n}_r}{n} = \frac{\sum_{i: z^i = \bullet} \gamma^i}{\sum_{j=1}^n \gamma^j}$$

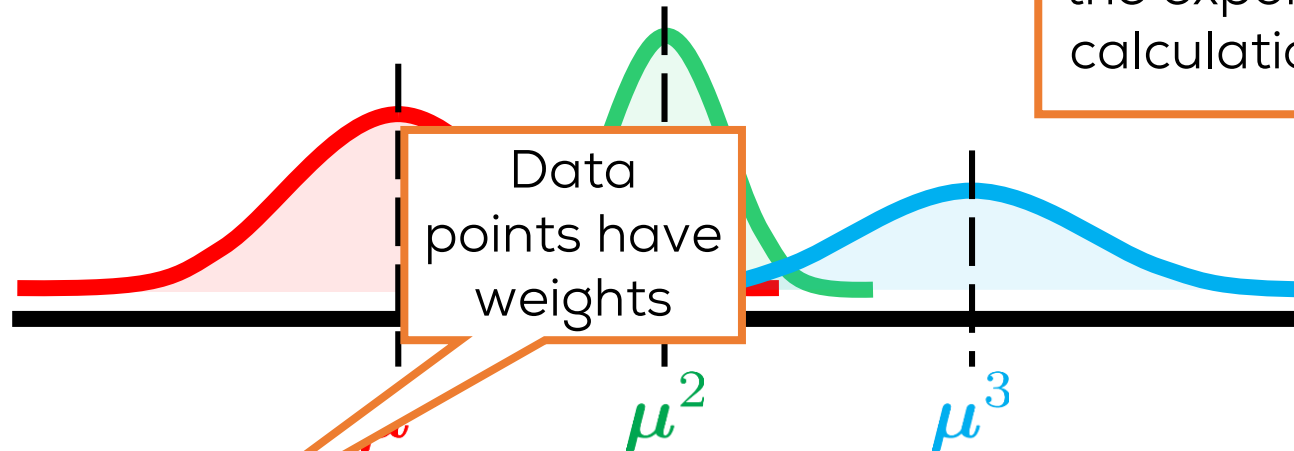
$$\boldsymbol{\mu}_{\text{MLE}}^1 = \frac{1}{\tilde{n}_r} \sum_{i: z^i = \bullet} \gamma^i \cdot \mathbf{x}^i$$

$$\Sigma_{\text{MLE}}^1 = \frac{1}{\tilde{n}_r} \sum_{i: z^i = \bullet} \gamma^i \cdot (\mathbf{x}^i - \boldsymbol{\mu}_{\text{MLE}}^1)(\mathbf{x}^i - \boldsymbol{\mu}_{\text{MLE}}^1)^{\top}$$

Reduces to normal MLE if $\gamma^i = 1$ for all i



Weighted MLE



Incorporating γ^i into the exponent makes calculations simpler

Also makes more sense (when we study the EM algo)

Why cant I "multiply" or "add" the weight term γ^i ?

Reduces to normal MLE if $\gamma^i = 1$ for all i

$$\mathbb{P}[\mathbf{x}^i, z^i, \gamma^i | \Theta] = (\mathbb{P}[z^i | \Theta] \cdot \mathbb{P}[\mathbf{x}^i | z^i, \Theta])^{\gamma^i} = \left(\pi_{z^i} \cdot \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}^{z^i}, \Sigma^{z^i}) \right)^{\gamma^i}$$

$$\mathbb{P}[X, \{z^i\}, \{\gamma^i\} | \Theta] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i, z^i, \gamma^i | \Theta]$$

$$\hat{\Theta}_{\text{MLE}} = \arg \max_{\Theta} \mathbb{P}[X, \{z^i\}, \{\gamma^i\} | \Theta]$$

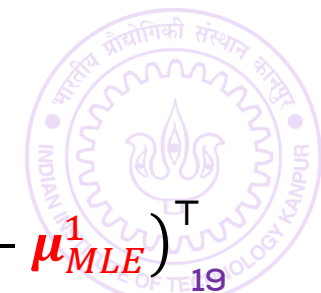
Take log and apply 1st order optimality

Exercise!

$$\pi_1^{\text{MLE}} = \frac{\# \text{ eff. red emails}}{\# \text{ eff. total emails}} = \frac{\tilde{n}_r}{n} = \frac{\sum_{i: z^i = \bullet} \gamma^i}{\sum_{j=1}^n \gamma^j}$$

$$\boldsymbol{\mu}_{\text{MLE}}^1 = \frac{1}{\tilde{n}_r} \sum_{i: z^i = \bullet} \gamma^i \cdot \mathbf{x}^i$$

$$\Sigma_{\text{MLE}}^1 = \frac{1}{\tilde{n}_r} \sum_{i: z^i = \bullet} \gamma^i \cdot (\mathbf{x}^i - \boldsymbol{\mu}_{\text{MLE}}^1)(\mathbf{x}^i - \boldsymbol{\mu}_{\text{MLE}}^1)^{\top}$$



Hard Alternating Minimization

ALTERNATING OPTIMIZATION

1. Initialize Θ^0
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
 1. Let $z^{i,t} = \arg \max_k \pi_k^t \cdot \mathcal{N}(\mathbf{x}^i \mid \boldsymbol{\mu}^{k,t}, \Sigma^{k,t})$
3. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[X, \{z^{i,t}\} \mid \Theta]$
4. Repeat until convergence

Soft Alternating Minimization

ALTERNATING OPTIMIZATION

1. For $i \in [n]$, create k copies of the data point
 1. Let $\mathbf{x}^i \rightarrow \{\mathbf{x}^{\{i,1\}}, \mathbf{x}^{\{i,2\}}, \dots, \mathbf{x}^{\{i,k\}}\}$
 2. Assign the k -th copy label k i.e. $z^{\{i,k\}} = k$
2. Initialize Θ^0
3. Update weights $\gamma^{i,k,t}$ using Θ^t
 1. Let $\gamma^{i,k,t} = \mathbb{P}[k \mid \mathbf{x}^i, \Theta^t] = \frac{\pi_k^t \cdot \mathcal{N}(\mathbf{X}^i \mid \mu^{k,t}, \Sigma^{k,t})}{\sum_j \pi_j^t \cdot \mathcal{N}(\mathbf{X}^i \mid \mu^{j,t}, \Sigma^{j,t})}$
4. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P} [\{x^{\{i,k\}}\}, \{z^{\{i,k\}}\}, \{\gamma^{\{i,k,t\}}\} \mid \Theta]$
5. Repeat until convergence



Soft Alternating Minimization

ALTERNATING OPTIMIZATION

1. For $i \in [n]$, create k copies of the data point

1. Let $\mathbf{x}^i \rightarrow \{\mathbf{x}^{\{i,1\}}, \mathbf{x}^{\{i,2\}}, \dots, \mathbf{x}^{\{i,k\}}\}$

2. Assign the k -th copy label k i.e. $z^{\{i,k\}} = k$

2. Initialize Θ^0

3. Update weights $\gamma^{i,k,t}$ using Θ^t

1. Let $\gamma^{i,k,t} = \mathbb{P}[k \mid \mathbf{x}^i, \Theta^t] = \frac{\pi_k^t \cdot \mathcal{N}(\mathbf{x}^i \mid \mu^{k,t}, \Sigma^{k,t})}{\sum_j \pi_j^t \cdot \mathcal{N}(\mathbf{x}^i \mid \mu^{j,t}, \Sigma^{j,t})}$

4. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[\{\mathbf{x}^{\{i,k\}}\}, \{z^{\{i,k\}}\}, \{\gamma^{\{i,k,t\}}\} \mid \Theta]$

5. Repeat until convergence

$$\sum_j \gamma^{\{i,j,t\}} = 1$$

for all i and
all t



Soft Alternating Minimization

ALTERNATING OPTIMIZATION

1. For $i \in [n]$, create k copies of the data point

1. Let $\mathbf{x}^i \rightarrow \{\mathbf{x}^{\{i,1\}}, \mathbf{x}^{\{i,2\}}, \dots, \mathbf{x}^{\{i,k\}}\}$

2. Assign the k -th copy label k i.e. $z^{\{i,k\}} = k$

2. Initialize Θ^0

3. Update weights $\gamma^{i,k,t}$ using Θ^t

1. Let $\gamma^{i,k,t} = \mathbb{P}[k \mid \mathbf{x}^i, \Theta^t] = \frac{\pi_k^t \cdot \mathcal{N}(\mathbf{x}^i \mid \mu^{k,t}, \Sigma^{k,t})}{\sum_j \pi_j^t \cdot \mathcal{N}(\mathbf{x}^i \mid \mu^{j,t}, \Sigma^{j,t})}$

4. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[\{\mathbf{x}^{\{i,k\}}\}, \{z^{\{i,k\}}\}, \{\gamma^{\{i,k,t\}}\} \mid \Theta]$

5. Repeat until convergence

Distribute a unit weight across clusters

$$\sum_j \gamma^{\{i,j,t\}} = 1$$

for all i and all t



Soft Alternating Minimization

Hard-AM: all weight was on a single cluster

ALTERNATING OPTIMIZATION

1. For $i \in [n]$, create k copies of the data point

1. Let $\mathbf{x}^i \rightarrow \{\mathbf{x}^{\{i,1\}}, \mathbf{x}^{\{i,2\}}, \dots, \mathbf{x}^{\{i,k\}}\}$

2. Assign the k -th copy label k i.e. $z^{\{i,k\}} = k$

2. Initialize Θ^0

3. Update weights $\gamma^{i,k,t}$ using Θ^t

1. Let $\gamma^{i,k,t} = \mathbb{P}[k | \mathbf{x}^i, \Theta^t] = \frac{\pi_k^t \cdot \mathcal{N}(\mathbf{x}^i | \mu^{k,t}, \Sigma^{k,t})}{\sum_j \pi_j^t \cdot \mathcal{N}(\mathbf{x}^i | \mu^{j,t}, \Sigma^{j,t})}$

4. Update $\Theta^{t+1} = \arg \max_{\Theta} \mathbb{P}[\{\mathbf{x}^{\{i,k\}}\}, \{z^{\{i,k\}}\}, \{\gamma^{\{i,k,t\}}\} | \Theta]$

5. Repeat until convergence

Distribute a unit weight across clusters

$$\sum_j \gamma^{\{i,j,t\}} = 1 \text{ for all } i \text{ and all } t$$



Soft k-means Algorithm

- Fix $\boldsymbol{\pi}_k^t = \frac{1}{K}$ for all iterations. Don't update it.
- Fix $\boldsymbol{\Sigma}^{k,t} = I$ for all iterations. Don't update it.

K-MEANS ALGO

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1,\dots,K}$
2. For all i , update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\|_2^2$
3. Let $n_k^t = |i: z^{i,t} = k|$
4. Update $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
5. Repeat until convergence

SOFT K-MEANS ALGO

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1,\dots,K}$
2. For all i , update $\gamma^{i,k,t}$ using $\boldsymbol{\mu}^{k,t}$

$$\text{Let } \gamma^{i,k,t} = \frac{\exp\left(-\frac{\|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\|_2^2}{2}\right)}{\sum_j \exp\left(-\frac{\|\mathbf{x}^i - \boldsymbol{\mu}^{j,t}\|_2^2}{2}\right)}$$

3. Let $\tilde{n}_k^t = \sum_i \gamma^{i,k,t}$
4. Update $\boldsymbol{\mu}^{k,t+1} = \frac{1}{\tilde{n}_k^t} \sum_i \gamma^{i,k,t} \cdot \mathbf{x}^i$
5. Repeat until convergence

Soft k-means Algorithm

- Fix $\boldsymbol{\pi}_k^t = \frac{1}{K}$ for all iterations. Don't update it.
- Fix $\boldsymbol{\Sigma}^{k,t} = I$ for all iterations. Don't update it.

K-MEANS ALGO

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1,\dots,K}$
2. For all i , update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
Let $z^{i,t} = \arg \min_k \|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\|_2^2$
3. Let $n_k^t = |i: z^{i,t} = k|$
4. Update $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t} \sum_{i: z^{i,t}=k} \mathbf{x}^i$
5. Repeat until convergence

SOFT K-MEANS ALGO

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1,\dots,K}$
2. For all i , update $\gamma^{i,k,t}$ using $\boldsymbol{\mu}^{k,t}$

$$\text{Let } \gamma^{i,k,t} = \frac{\exp\left(-\frac{\|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\|_2^2}{2}\right)}{\sum_j \exp\left(-\frac{\|\mathbf{x}^i - \boldsymbol{\mu}^{j,t}\|_2^2}{2}\right)}$$

3. Let $\tilde{n}_k^t = \sum_i \gamma^{i,k,t}$
4. Update $\boldsymbol{\mu}^{k,t+1} = \frac{1}{\tilde{n}_k^t} \sum_i \gamma^{i,k,t} \cdot \mathbf{x}^i$
5. Repeat until convergence

Mixed Regression

Mixed Regression

Sept 13, 2017



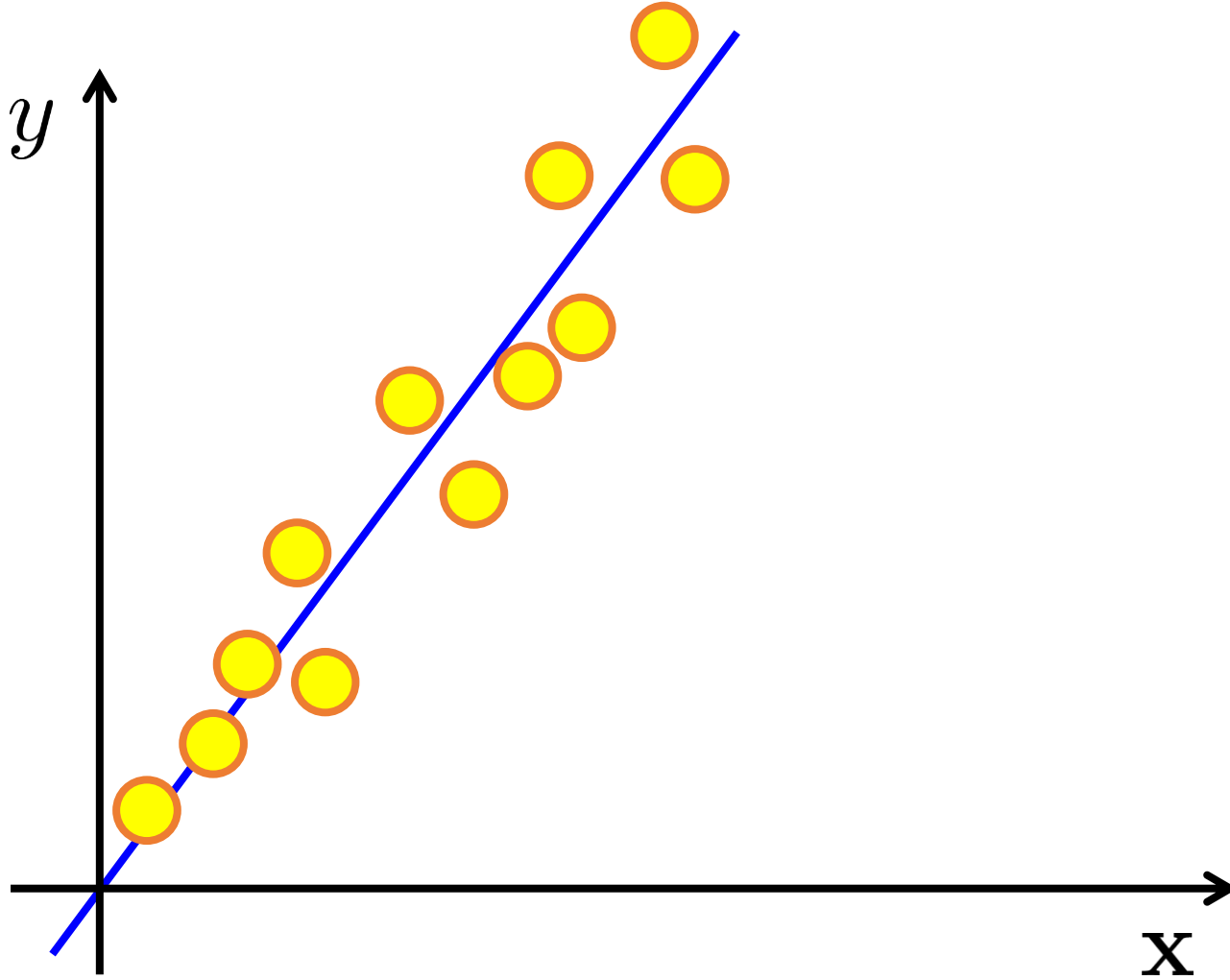
Mixed Regression



Sept 13, 2017



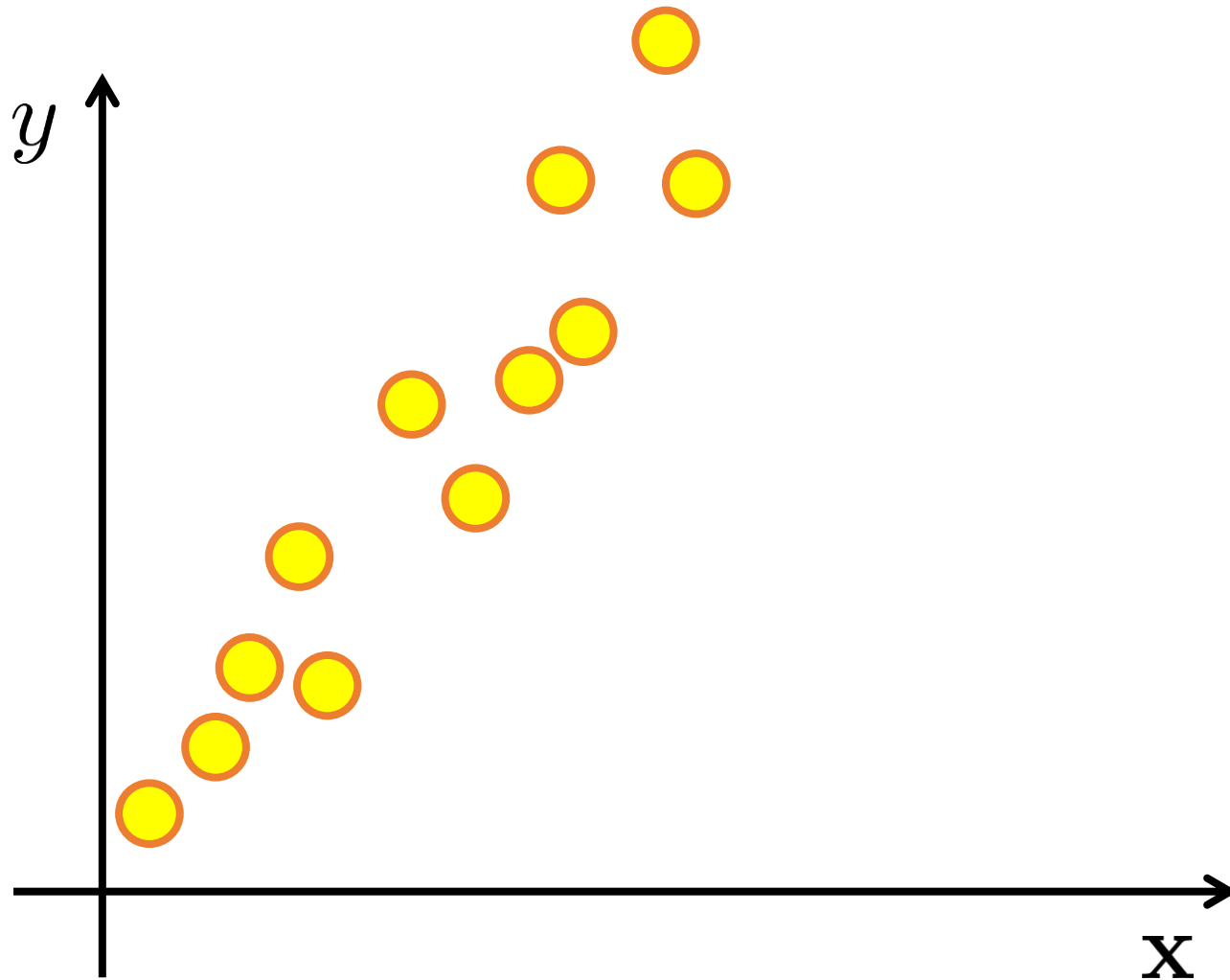
Mixed Regression



Sept 13, 2017



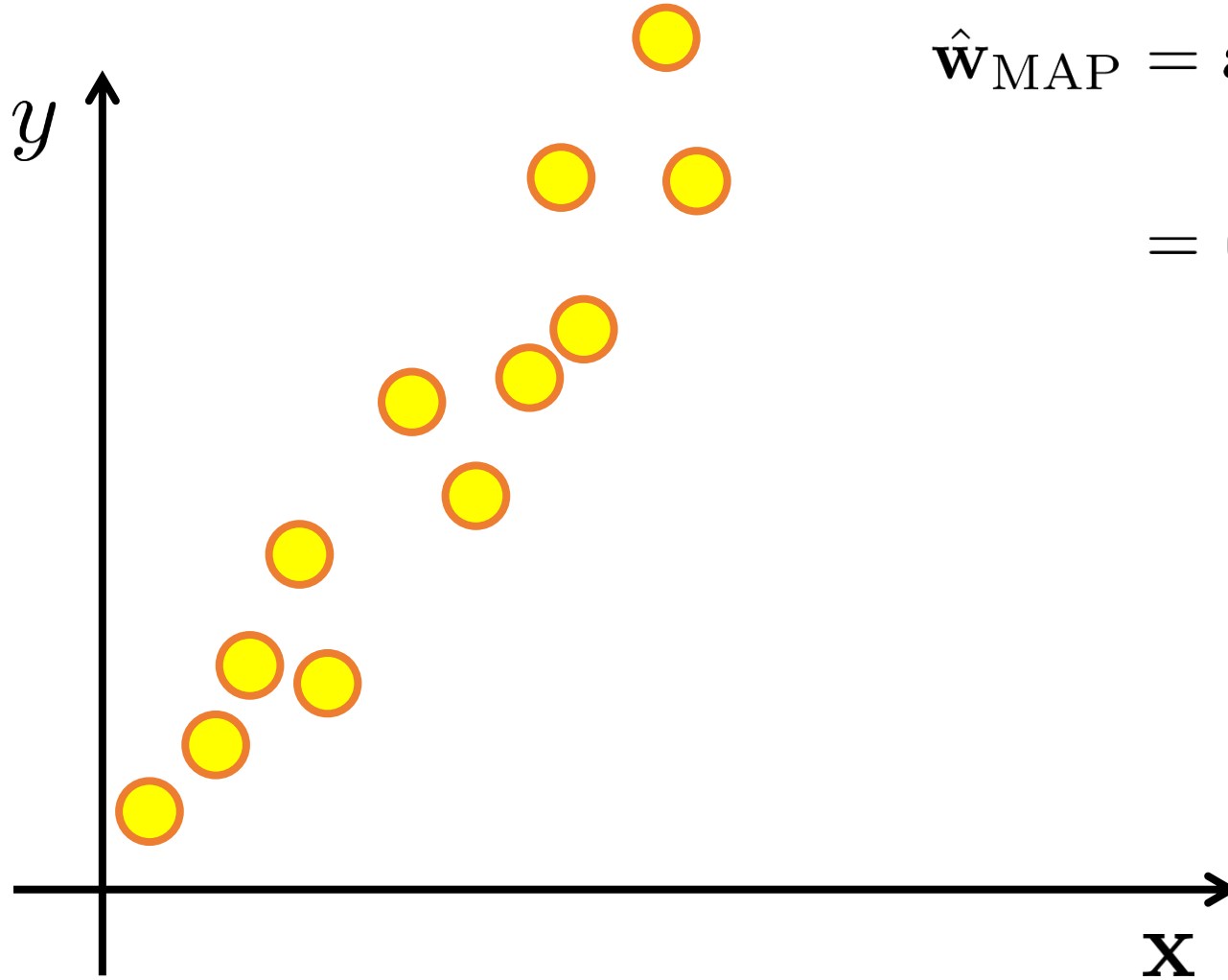
Mixed Regression



Sept 13, 2017

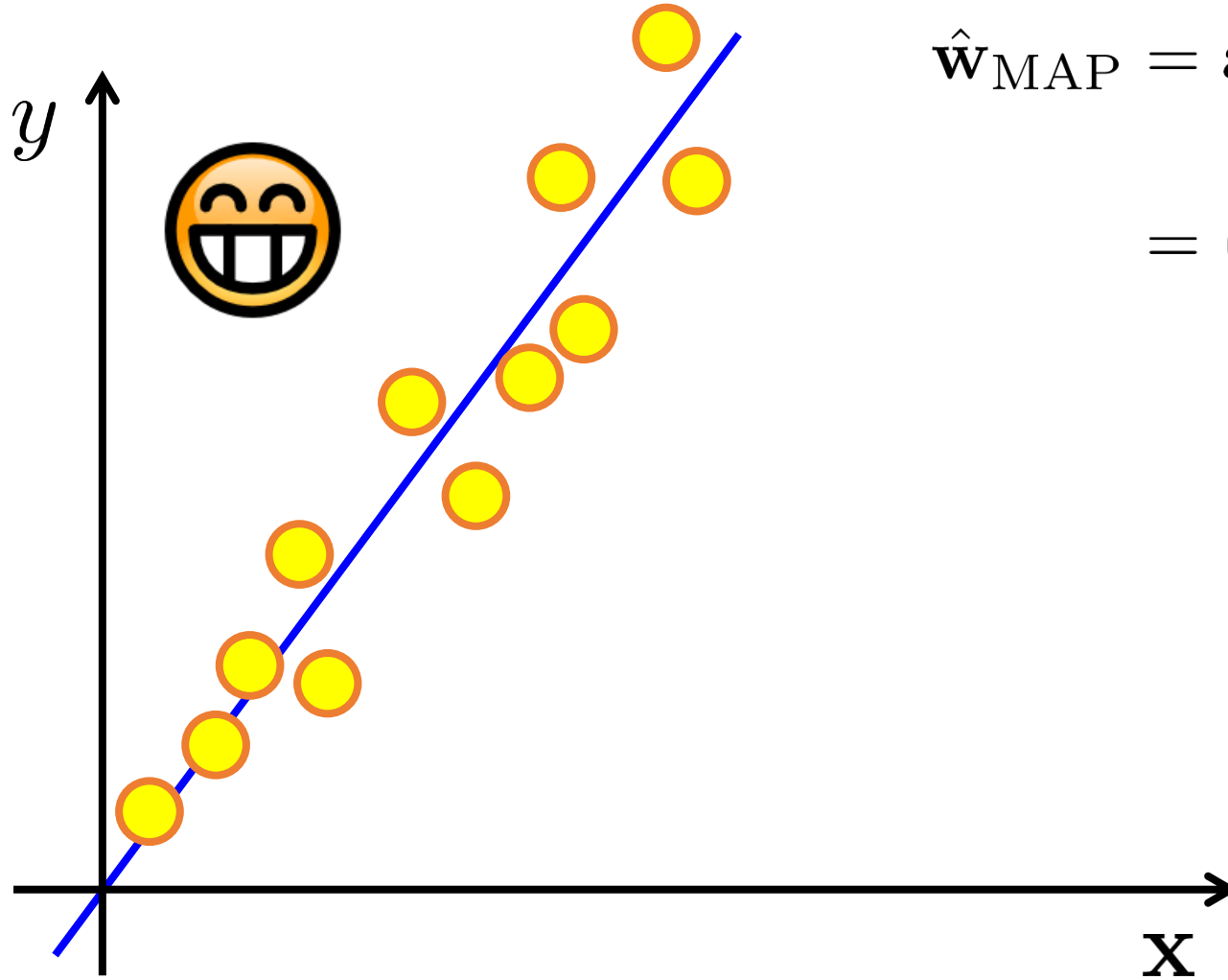


Mixed Regression



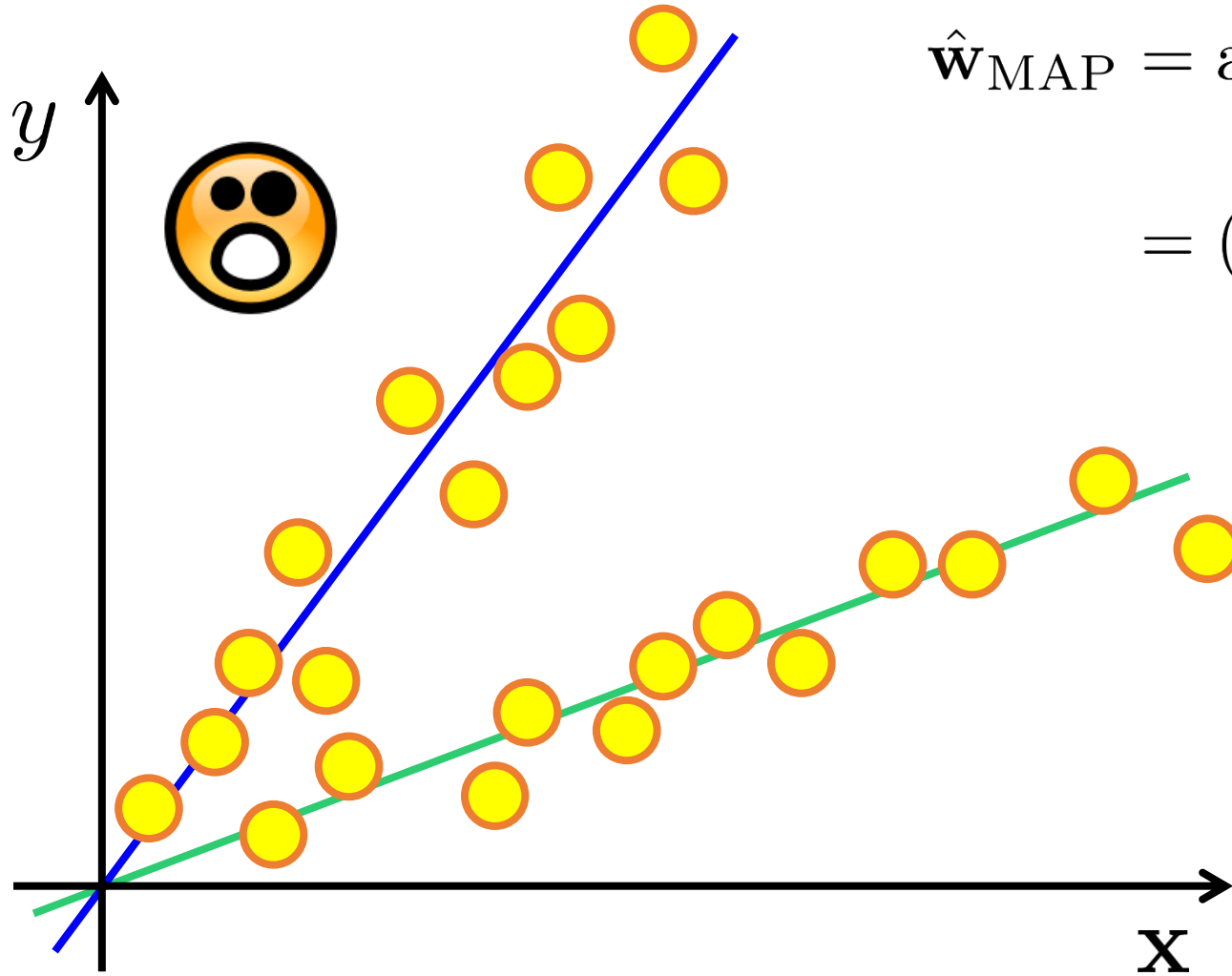
$$\begin{aligned}\hat{\mathbf{w}}_{\text{MAP}} &= \arg \min \sum_{i=1}^n \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \lambda \cdot \|\mathbf{w}\|_2^2 \\ &= (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}\end{aligned}$$

Mixed Regression



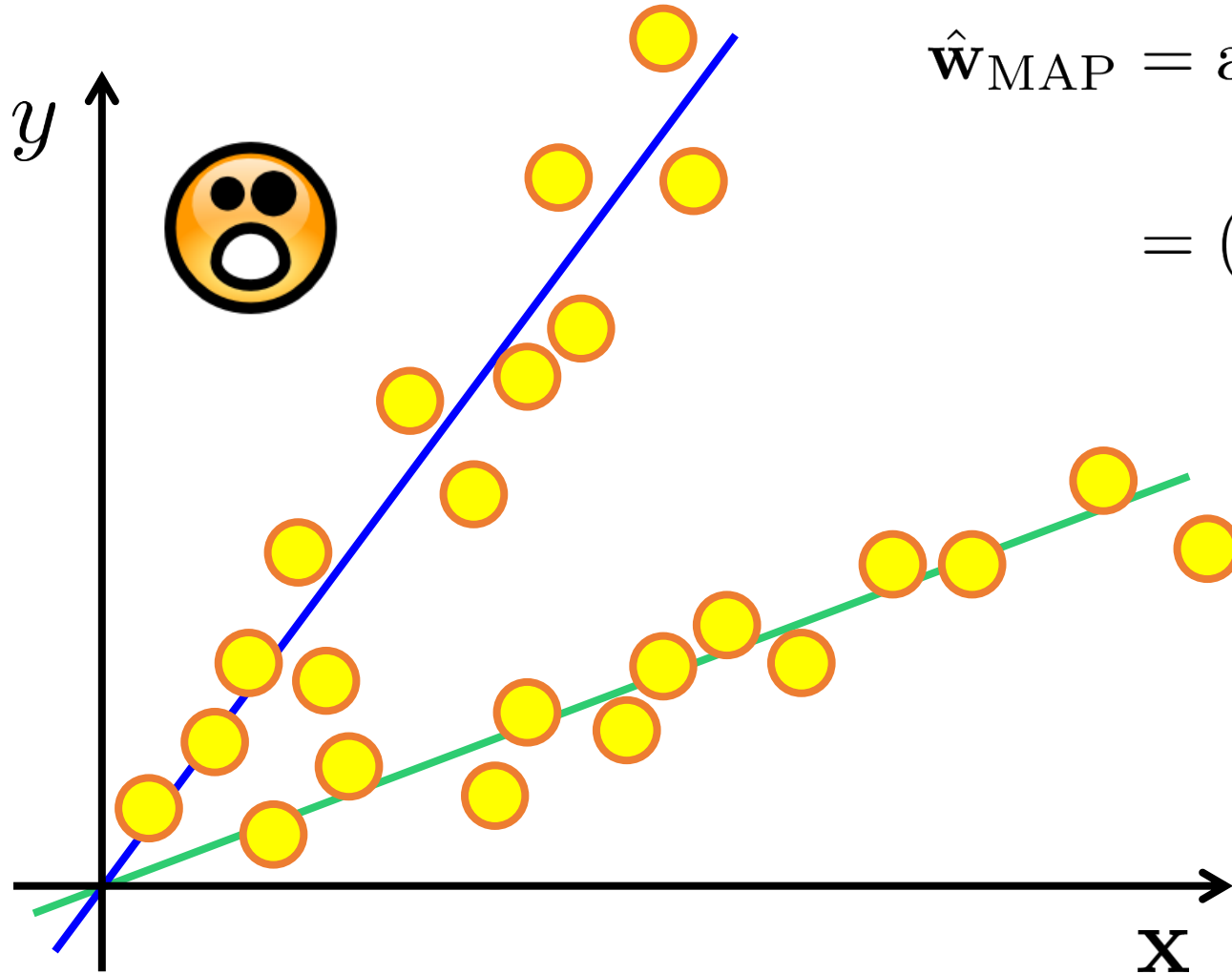
$$\begin{aligned}\hat{\mathbf{w}}_{\text{MAP}} &= \arg \min \sum_{i=1}^n \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \lambda \cdot \|\mathbf{w}\|_2^2 \\ &= (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}\end{aligned}$$

Mixed Regression



$$\begin{aligned}\hat{\mathbf{w}}_{\text{MAP}} &= \arg \min \sum_{i=1}^n \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \lambda \cdot \|\mathbf{w}\|_2^2 \\ &= (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}\end{aligned}$$

Mixed Regression

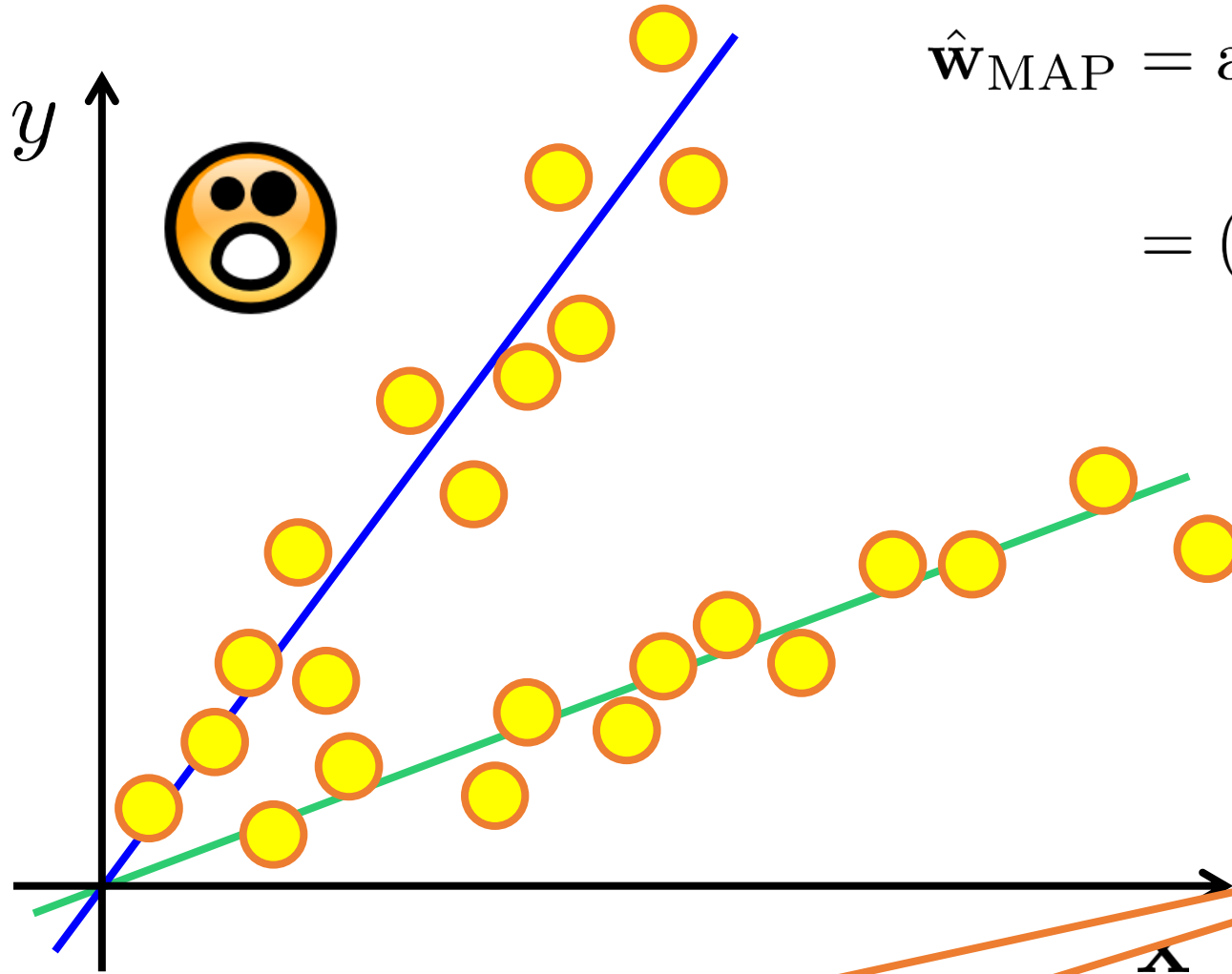


$$\begin{aligned}\hat{\mathbf{w}}_{\text{MAP}} &= \arg \min \sum_{i=1}^n \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \lambda \cdot \|\mathbf{w}\|_2^2 \\ &= (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}\end{aligned}$$

Regression with a
classification touch ☺

E.g. recommendation systems, a
different system might work for
young and old but we don't know
age of users

Mixed Regression



$$\begin{aligned}\hat{\mathbf{w}}_{\text{MAP}} &= \arg \min \sum_{i=1}^n \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \lambda \cdot \|\mathbf{w}\|_2^2 \\ &= (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}\end{aligned}$$

Regression with a
classification touch 😊

E.g. recommendation systems, a
different system might work for
young and old but we don't know
age of users

Mixed models, Mixture of Experts

Mixed Regression

$$\mathbb{P}[y \mid \mathbf{x}^i, z^i, \mathbf{w}] = \mathcal{N}(\langle \mathbf{w}^{z^i}, \mathbf{x}^i \rangle, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \langle \mathbf{w}^{z^i}, \mathbf{x}^i \rangle)^2}{2\sigma^2}\right)$$

- Assume for sake of simplicity that both models are equally sampled $\mathbb{P}[z = 0] = \mathbb{P}[z = 1] = 0.5$
- Assume for sake of simplicity that Gaussian noise $\sigma_1 = \sigma_2 = 1$
- Only unknowns are the two linear models $\mathbf{w}^1, \mathbf{w}^2$

Mixed Regression

$$\mathbb{P}[y \mid \mathbf{x}^i, z^i, \mathbf{w}] = \mathcal{N}(\langle \mathbf{w}^{z^i}, \mathbf{x}^i \rangle, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \langle \mathbf{w}^{z^i}, \mathbf{x}^i \rangle)^2}{2\sigma^2}\right)$$

$z^i \in \{0,1\}$ indicates
which line
generated the data

- Assume for sake of simplicity that both models are equally sampled $\mathbb{P}[z = 0] = \mathbb{P}[z = 1] = 0.5$
- Assume for sake of simplicity that Gaussian noise $\sigma_1 = \sigma_2 = 1$
- Only unknowns are the two linear models $\mathbf{w}^1, \mathbf{w}^2$

Mixed Regression

$$\mathbb{P}[y \mid \mathbf{x}^i, z^i, \mathbf{w}] = \mathcal{N}(\langle \mathbf{w}^{z^i}, \mathbf{x}^i \rangle, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \langle \mathbf{w}^{z^i}, \mathbf{x}^i \rangle)^2}{2\sigma^2}\right)$$

$z^i \in \{0,1\}$ indicates
which line
generated the data

z^i is a latent
variable!

- Assume for sake of simplicity that both models are equally sampled $\mathbb{P}[z = 0] = \mathbb{P}[z = 1] = 0.5$
- Assume for sake of simplicity that Gaussian noise $\sigma_1 = \sigma_2 = 1$
- Only unknowns are the two linear models $\mathbf{w}^1, \mathbf{w}^2$

Hard Mixed Regression

ALTERNATING OPTIMIZATION

1. Initialize $\Theta^0 = \{\mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}}\}$
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
 1. Let $z^{i,t} = \arg \max_{k \in \{0,1\}} \mathcal{N}(y^i \mid \mathbf{x}^i, \mathbf{w}^{k,t})$
3. Update $\mathbf{w}^{\{t+1,k\}} = \arg \min_{\mathbf{w}} \sum_{i: z^{i,t}=k} (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \frac{1}{2} \|\mathbf{w}\|_2^2$
4. Set $\Theta^{t+1} = \{\mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}}\}$
5. Repeat until convergence

Hard Mixed Regression

ALTERNATING OPTIMIZATION

1. Initialize $\Theta^0 = \{\mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}}\}$
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
 1. Let $z^{i,t} = \arg \max_{k \in \{0,1\}} \mathcal{N}(y^i \mid \mathbf{x}^i, \mathbf{w}^{k,t})$
3. Update $\mathbf{w}^{\{t+1,k\}} = \arg \min_{\mathbf{w}} \sum_{i: z^{i,t}=k} (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \frac{1}{2} \|\mathbf{w}\|_2^2$
4. Set $\Theta^{t+1} = \{\mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}}\}$
5. Repeat until convergence

Hard Mixed Regression

ALTERNATING OPTIMIZATION

1. Initialize $\Theta^0 = \{\mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}}\}$
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
 1. Let $z^{i,t} = \arg \min_{k \in \{0,1\}} |y^i - \langle \mathbf{w}^{t,k}, \mathbf{x}^i \rangle|$
3. Update $\mathbf{w}^{\{t+1,k\}} = \arg \min_{\mathbf{w}} \sum_{i: z^{i,t}=k} (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \frac{1}{2} \|\mathbf{w}\|_2^2$
4. Set $\Theta^{t+1} = \{\mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}}\}$
5. Repeat until convergence

Hard Mixed Regression

ALTERNATING OPTIMIZATION

Assign to
the “closest”
line!

1. Initialize $\Theta^0 = \{\mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}}\}$
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
 1. Let $z^{i,t} = \arg \min_{k \in \{0,1\}} |y^i - \langle \mathbf{w}^{t,k}, \mathbf{x}^i \rangle|$
3. Update $\mathbf{w}^{\{t+1,k\}} = \arg \min_{\mathbf{w}} \sum_{i: z^{i,t}=k} (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \frac{1}{2} \|\mathbf{w}\|_2^2$
4. Set $\Theta^{t+1} = \{\mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}}\}$
5. Repeat until convergence

Exercise: verify
these updates

Hard Mixed Regression

ALTERNATING OPTIMIZATION

1. Initialize $\Theta^0 = \{\mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}}\}$
2. For $i \in [n]$, update $z^{i,t}$ using Θ^t
 1. Let $z^{i,t} = \arg \min_{k \in \{0,1\}} |y^i - \langle \mathbf{w}^{t,k}, \mathbf{x}^i \rangle|$
3. Update $\mathbf{w}^{\{t+1,k\}} = \arg \min_{\mathbf{w}} \sum_{i: z^{i,t}=k} (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \frac{1}{2} \|\mathbf{w}\|_2^2$
4. Set $\Theta^{t+1} = \{\mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}}\}$
5. Repeat until convergence

Assign to the "closest" line!

In k-means, we assigned to the closest mean

Exercise: verify these updates

Weighted Regression

Sept 13, 2017



Weighted Regression

$$\begin{aligned}\hat{\mathbf{w}}_{\text{MAP}} &= \arg \min \sum_{i=1}^n \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \lambda \cdot \|\mathbf{w}\|_2^2 \\ &= (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}\end{aligned}$$

Weighted Regression

$$\begin{aligned}\hat{\mathbf{w}}_{\text{MAP}} &= \arg \min \sum_{i=1}^n \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \lambda \cdot \|\mathbf{w}\|_2^2 \\ &= (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y} \\ \hat{\mathbf{w}}_{\text{MAP}} &= \arg \min \sum_{i=1}^n \gamma_i \cdot \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \lambda \cdot \|\mathbf{w}\|_2^2 \\ &= (M + \lambda \cdot I)^{-1} \mathbf{b} \\ M &= \sum_{i=1}^n \gamma_i \mathbf{x}^i (\mathbf{x}^i)^\top \\ \mathbf{b} &= \sum_{i=1}^n \gamma_i y_i \cdot \mathbf{x}^i\end{aligned}$$

Weighted Regression

$$\begin{aligned}\hat{\mathbf{w}}_{\text{MAP}} &= \arg \min \sum_{i=1}^n \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \lambda \cdot \|\mathbf{w}\|_2^2 \\ &= (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y} \\ \hat{\mathbf{w}}_{\text{MAP}} &= \arg \min \sum_{i=1}^n \gamma_i \cdot \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \lambda \cdot \|\mathbf{w}\|_2^2 \\ &= (M + \lambda \cdot I)^{-1} \mathbf{b} \\ M &= \sum_{i=1}^n \gamma_i \mathbf{x}^i (\mathbf{x}^i)^\top \\ \mathbf{b} &= \sum_{i=1}^n \gamma_i y_i \cdot \mathbf{x}^i\end{aligned}$$

Soft Mixed Regression

SOFT ALTERNATING OPTIMIZATION

1. Initialize $\Theta^0 = \{\mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}}\}$
2. For $i \in [n]$, update $\gamma^{\{i,k,t\}}$ using Θ^t
 1. Let $c^{\{i,k,t\}} = \exp\left(-\frac{(y^i - \langle \mathbf{w}^{\{t,k\}}, \mathbf{x}^i \rangle)^2}{2}\right)$
 2. Let $\gamma^{\{i,k,t\}} = \frac{c^{\{i,k,t\}}}{c^{\{i,0,t\}} + c^{\{i,1,t\}}}$
3. Update $\mathbf{w}^{\{t+1,k\}} = \arg \min_{\mathbf{w}} \sum_i \gamma^{\{i,k,t\}} \cdot (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \frac{1}{2} \|\mathbf{w}\|_2^2$
4. Set $\Theta^{t+1} = \{\mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}}\}$
5. Repeat until convergence

Soft Mixed Regression

SOFT ALTERNATING OPTIMIZATION

1. Initialize $\Theta^0 = \{\mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}}\}$
2. For $i \in [n]$, update $\gamma^{\{i,k,t\}}$ using Θ^t
 1. Let $c^{\{i,k,t\}} = \exp\left(-\frac{(y^i - \langle \mathbf{w}^{\{t,k\}}, \mathbf{x}^i \rangle)^2}{2}\right)$
 2. Let $\gamma^{\{i,k,t\}} = \frac{c^{\{i,k,t\}}}{c^{\{i,0,t\}} + c^{\{i,1,t\}}}$
3. Update $\mathbf{w}^{\{t+1,k\}} = \arg \min_{\mathbf{w}} \sum_i \gamma^{\{i,k,t\}} \cdot (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \frac{1}{2} \|\mathbf{w}\|_2^2$
4. Set $\Theta^{t+1} = \{\mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}}\}$
5. Repeat until convergence

Exercise: derive these updates

Soft Mixed Regression

SOFT ALTERNATING OPTIMIZATION

1. Initialize $\Theta^0 = \{\mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}}\}$
2. For $i \in [n]$, update $\gamma^{\{i,k,t\}}$ using Θ^t
 1. Let $c^{\{i,k,t\}} = \exp\left(-\frac{(y^i - \langle \mathbf{w}^{\{t,k\}}, \mathbf{x}^i \rangle)^2}{2}\right)$
 2. Let $\gamma^{\{i,k,t\}} = \frac{c^{\{i,k,t\}}}{c^{\{i,0,t\}} + c^{\{i,1,t\}}}$
3. Update $\mathbf{w}^{\{t+1,k\}} = \arg \min_{\mathbf{w}} \sum_i \gamma^{\{i,k,t\}} \cdot (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \frac{1}{2} \|\mathbf{w}\|_2^2$
4. Set $\Theta^{t+1} = \{\mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}}\}$
5. Repeat until convergence

$$\propto \mathbb{P}[k | y^i, \mathbf{x}^i, \Theta^t]$$

Exercise: derive these updates

Soft Mixed Regression

SOFT ALTERNATING OPTIMIZATION

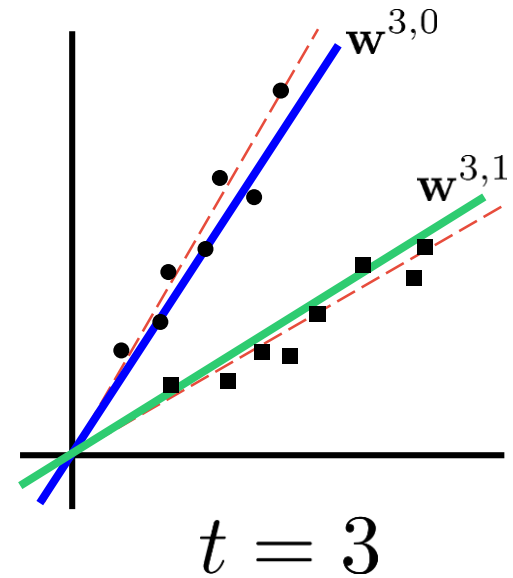
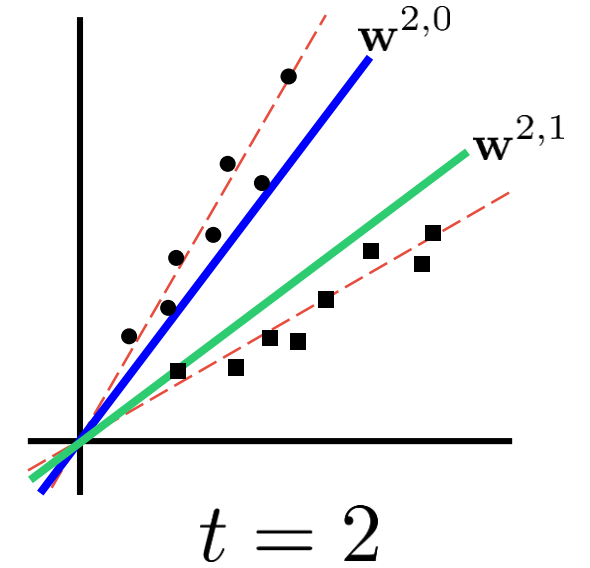
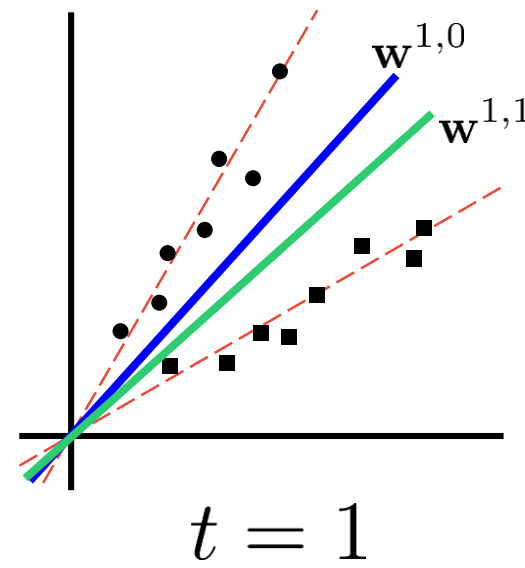
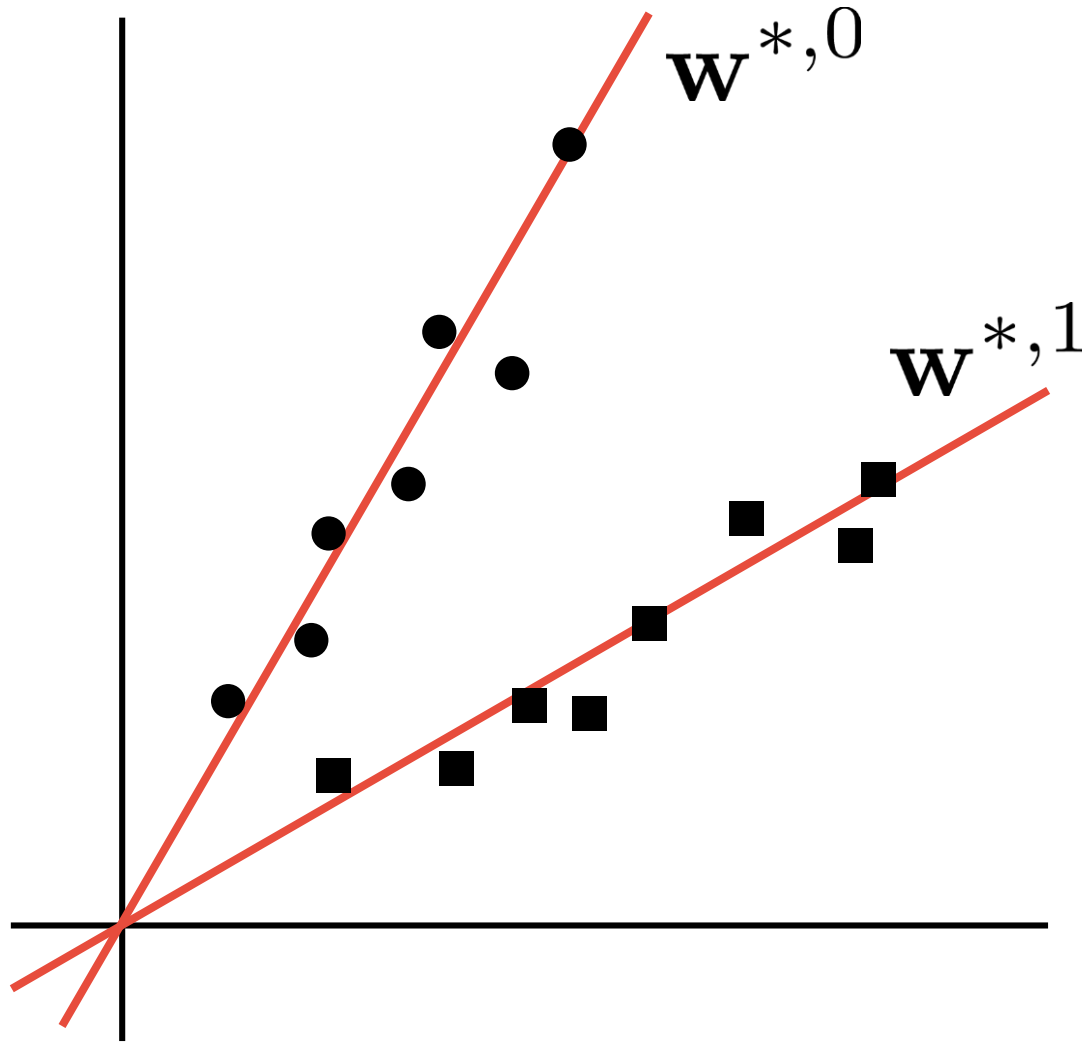
1. Initialize $\Theta^0 = \{\mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}}\}$
2. For $i \in [n]$, update $\gamma^{\{i,k,t\}}$ using Θ^t
 1. Let $c^{\{i,k,t\}} = \exp\left(-\frac{(y^i - \langle \mathbf{w}^{\{t,k\}}, \mathbf{x}^i \rangle)^2}{2}\right)$

$\propto \mathbb{P}[k | y^i, \mathbf{x}^i, \Theta^t]$
 2. Let $\gamma^{\{i,k,t\}} = \frac{c^{\{i,k,t\}}}{c^{\{i,0,t\}} + c^{\{i,1,t\}}}$

$= \mathbb{P}[k | y^i, \mathbf{x}^i, \Theta^t]$
3. Update $\mathbf{w}^{\{t+1,k\}} = \arg \min_{\mathbf{w}} \sum_i \gamma^{\{i,k,t\}} \cdot (y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle)^2 + \frac{1}{2} \|\mathbf{w}\|_2^2$
4. Set $\Theta^{t+1} = \{\mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}}\}$
5. Repeat until convergence

Exercise: derive these updates

Soft MR in action



The EM Algorithm

- Generalizes the notion of “soft” updates
- Very powerful algorithm
- Related to alternating minimization
- Will study this next time!

Please give your Feedback

<http://tinyurl.com/ml17-18afb>