# Data Modelling Methods-III

CS771: Introduction to Machine Learning

Purushottam Kar

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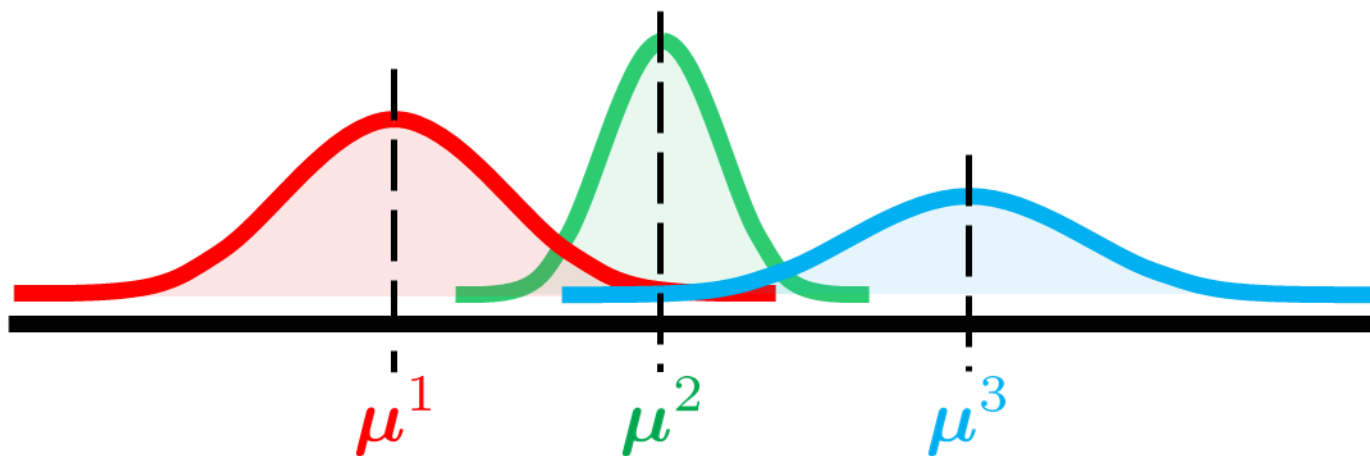
- September 21st, 2017 (Thursday) 1300–1500 hrs
- Venue L18, 19, L20 (all OROS)
- Syllabus: till whatever we cover today
- Open notes (handwritten only)
- No printed/photocopied material
- No laptops, i-pads, mobile phones (switched off)
- Please bring a notepad with you for rough work
- Please bring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

# Recap

# The generative story for labelled data



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

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

$$\hat{\boldsymbol{\Theta}}_{\mathrm{MLE}} = \arg\max_{\boldsymbol{\Theta}} \mathbb{P}\left[X, \left\{z^i\right\} \mid \boldsymbol{\Theta}\right]$$

$$\boldsymbol{\pi}_1^{\mathrm{MLE}} = \frac{\text{\# red emails}}{\text{\# total emails}} = \frac{n_r}{n}$$

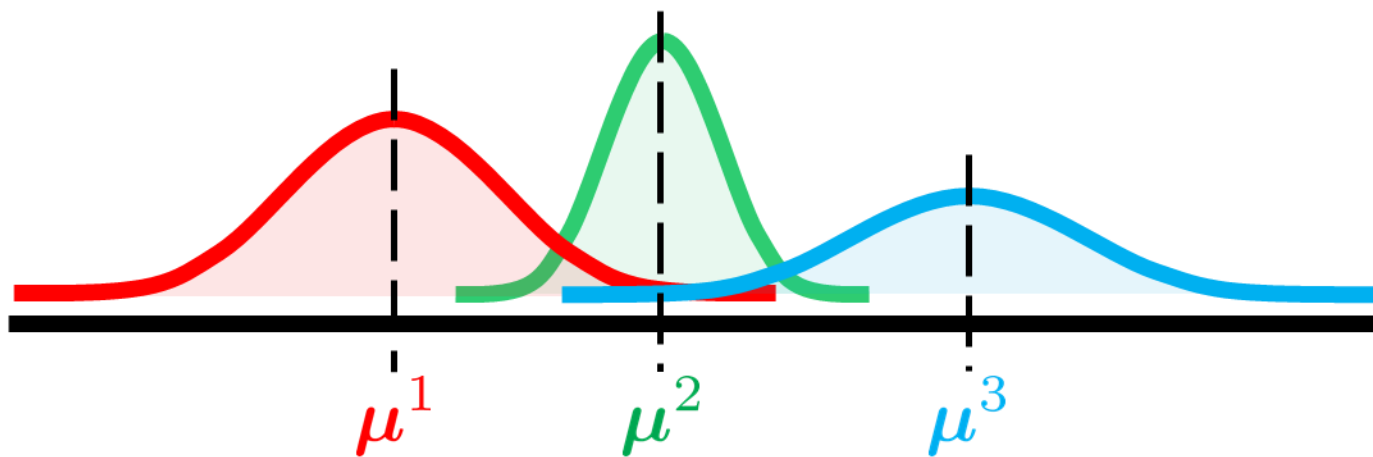$$\boldsymbol{\mu}_{\mathrm{MLE}}^1 = \frac{1}{n_r} \sum_{i:z^i=\bullet} \mathbf{x}^i$$

$$\Sigma_{\mathrm{MLE}}^1 = \frac{1}{n_r} \sum_{i:z^i=\bullet} \left(\mathbf{x}^i - \boldsymbol{\mu}_{MLE}^1\right)\left(\mathbf{x}^i - \boldsymbol{\mu}_{MLE}^1\right)^{\mathsf{T}}$$

Take log and apply 1st order optimality

Read [DAU] Sections 9.1-9.5

12

# Recap

# The generative story for unlabelled data



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

$$\mathbb{P}\left[\mathbf{x}^i \mid \boldsymbol{\Theta}\right] = \sum_{k=1}^{K} \mathbb{P}\left[\mathbf{x}^i, z^i = k \mid \boldsymbol{\Theta}\right] = \sum_{k=1}^{K} \boldsymbol{\pi}_k \cdot \mathbb{P}\left[\mathbf{x}^i \mid z^i = k, \boldsymbol{\Theta}\right]$$

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

$$= \sum_{k=1}^{K} \boldsymbol{\pi}_k \cdot \mathcal{N}(\mathbf{x}^i \mid \boldsymbol{\mu}^k, \Sigma^k)$$

Gaussian Mixture Model (GMM) with $K$ components

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

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

# Recap

# A Ray of Hope

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

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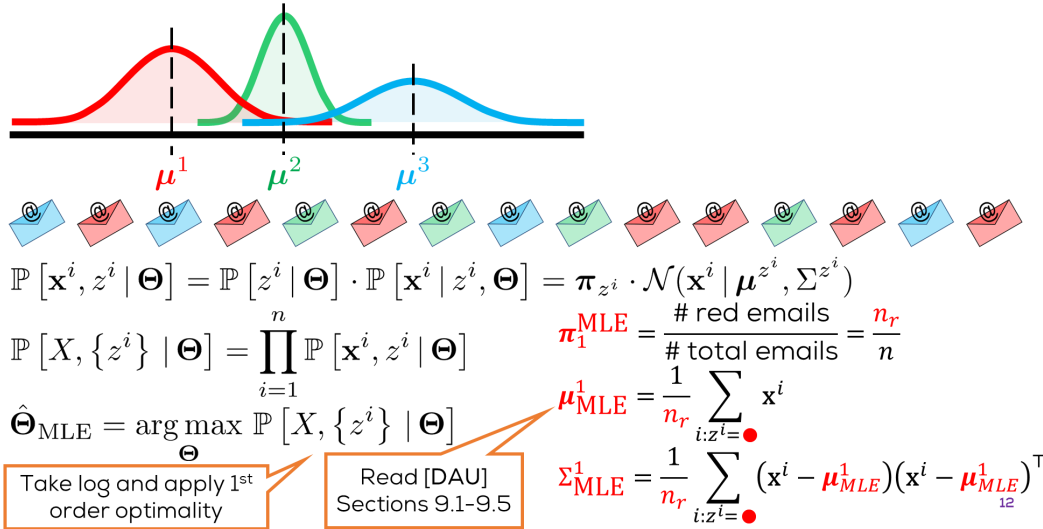
**ALTERNATING OPTIMIZATION**

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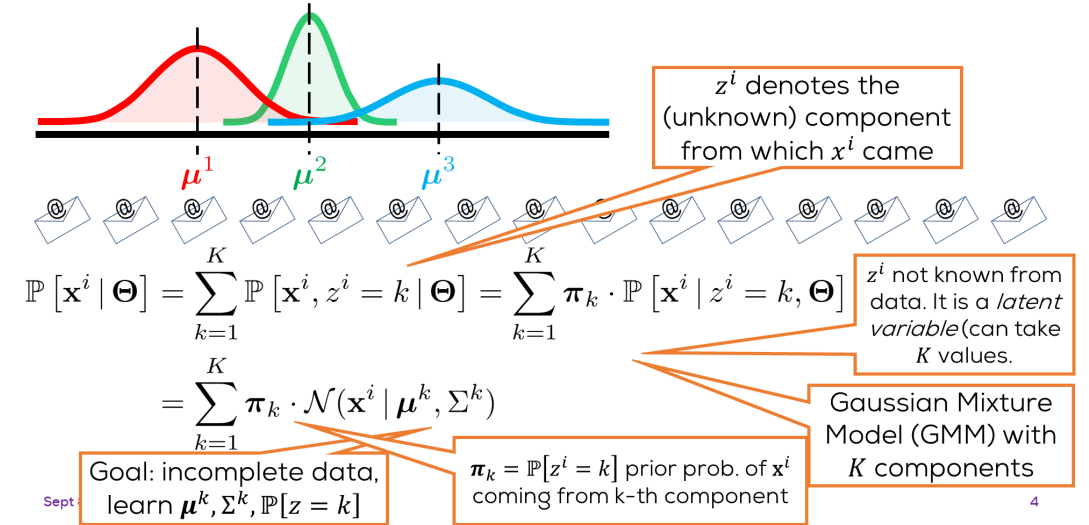4. Repeat until convergence

Various ways of updating $z^i$

# Recap

## The generative story for labelled data



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

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

$$\hat{\boldsymbol{\Theta}}_{\mathrm{MLE}} = \arg\max_{\boldsymbol{\Theta}} \mathbb{P}\left[X, \{z^i\} \mid \boldsymbol{\Theta}\right]$$

$$\boldsymbol{\pi}_1^{\mathrm{MLE}} = \frac{\text{\# red emails}}{\text{\# total emails}} = \frac{n_r}{n}$$

$$\boldsymbol{\mu}_{\mathrm{MLE}}^1 = \frac{1}{n_r} \sum_{i:z^i=\bullet} \mathbf{x}^i$$

$$\Sigma_{\mathrm{MLE}}^1 = \frac{1}{n_r} \sum_{i:z^i=\bullet} (\mathbf{x}^i - \boldsymbol{\mu}_{\mathrm{MLE}}^1)(\mathbf{x}^i - \boldsymbol{\mu}_{\mathrm{MLE}}^1)^\top$$

Take log and apply 1st order optimality

Read [DAU] Sections 9.1-9.5

12

## The generative story for unlabelled data



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

$$\mathbb{P}\left[\mathbf{x}^i \mid \boldsymbol{\Theta}\right] = \sum_{k=1}^{K} \mathbb{P}\left[\mathbf{x}^i, z^i = k \mid \boldsymbol{\Theta}\right] = \sum_{k=1}^{K} \boldsymbol{\pi}_k \cdot \mathbb{P}\left[\mathbf{x}^i \mid z^i = k, \boldsymbol{\Theta}\right]$$

$$= \sum_{k=1}^{K} \boldsymbol{\pi}_k \cdot \mathcal{N}(\mathbf{x}^i \mid \boldsymbol{\mu}^k, \Sigma^k)$$

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

Gaussian Mixture Model (GMM) with $K$ components

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

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

Sept

4

## A Ray of Hope

$$\hat{\boldsymbol{\Theta}}_{\mathrm{MLE}} = \arg\max_{\boldsymbol{\Theta}} \mathbb{P}\left[X \mid \boldsymbol{\Theta}\right]$$

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

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

Various ways of updating $z^i$

# Hard Assignment

The K-means algorithm

# A Ray of Hope

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

## ALTERNATING OPTIMIZATION

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

Various ways of updating $z^i$

# A Ray of Hope

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

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

# A Ray of Hope

$$\hat{\mathbf{\Theta}}_{\mathrm{MLE}} = \arg\max_{\mathbf{\Theta}} \mathbb{P}\left[X \mid \mathbf{\Theta}\right]$$

**ALTERNATING OPTIMIZATION**

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

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

# A Ray of Hope

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

## ALTERNATING OPTIMIZATION

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

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

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

# A Ray of Hope

$$\hat{\boldsymbol{\Theta}}_{\mathrm{MLE}} = \arg\max_{\boldsymbol{\Theta}} \mathbb{P}\left[X \mid \boldsymbol{\Theta}\right]$$

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

Bayes Rule!

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

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

# A Ray of Hope

$$\hat{\mathbf{\Theta}}_{\mathrm{MLE}} = \arg\max_{\mathbf{\Theta}} \mathbb{P}\left[X \mid \mathbf{\Theta}\right]$$

## ALTERNATING OPTIMIZATION

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

Bayes Rule!

$$z^{i,t} = \arg\max_{k \in [K]} \mathbb{P}\left[k \mid \mathbf{x}^i, \mathbf{\Theta}^t\right] = \arg\max_{k \in [K]} \quad \boldsymbol{\pi}_k^t \quad \cdot \quad \mathbb{P}\left[\mathbf{x}^i \mid k, \mathbf{\Theta}^t\right]$$

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

# A Ray of Hope

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

## ALTERNATING OPTIMIZATION

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

Bayes Rule!

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

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

# Towards the K-Means Algorithm

1. Initialize $\boldsymbol{\Theta}^0$

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

   1. Let $z^{i,t} = \arg\max_k \boldsymbol{\pi}_k^t \cdot \mathcal{N}\left(\mathbf{x}^i \mid \boldsymbol{\mu}^{k,t}, \Sigma^{k,t}\right)$

3. Update $\boldsymbol{\Theta}^{t+1} = \arg\max_{\boldsymbol{\Theta}} \mathbb{P}\left[X, \left\{z^{i,t}\right\} \mid \boldsymbol{\Theta}\right]$

   1. Let $\boldsymbol{\pi}_k^{t+1} = \frac{n_k^t}{n}$, where $n_k^t = \left|\left\{i: \ z^{i,t} = k\right\}\right|$

   2. Let $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t}\sum_{i:z^{i,t}=k}\mathbf{x}^i$

   3. Let $\Sigma_k^{t+1} = \frac{1}{n_k^t}\sum_{i:z^{i,t}=k}\left(\mathbf{x}^i - \mu^{k,t}\right)\left(\mathbf{x}^i - \mu^{k,t}\right)^\top$

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 $\{\boldsymbol{\mu}^{k,0}\}_{k=1,\ldots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
   1. Let $z^{i,t} = \arg\max_k \mathcal{N}\left(\mathbf{x}^i \mid \boldsymbol{\mu}^{k,t}, I\right)$
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.
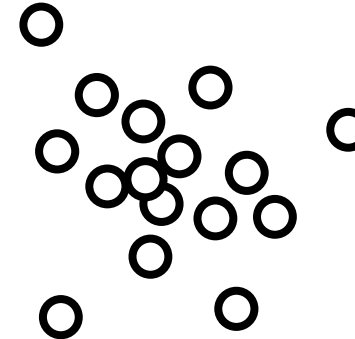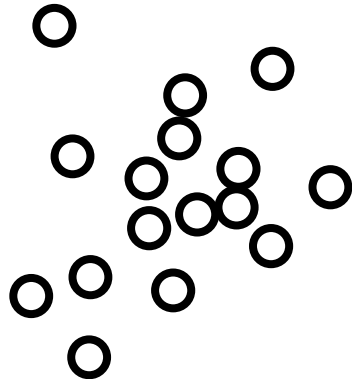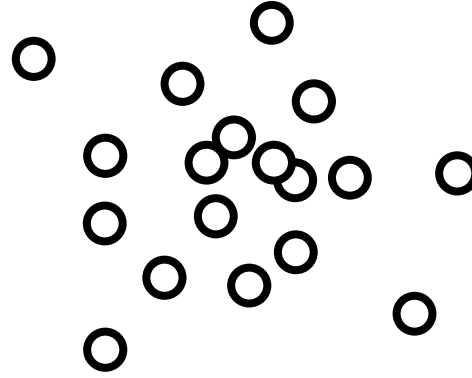
## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1,\ldots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
   1. Let $z^{i,t} = \arg\min_k \left\| \mathbf{x}^i - \boldsymbol{\mu}^{k,t} \right\|_2$
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 $\{\boldsymbol{\mu}^{k,0}\}_{k=1,\ldots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
    1. Let $z^{i,t} = \arg\min_k \left\| \mathbf{x}^i - \boldsymbol{\mu}^{k,t} \right\|_2^2$
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

$$\hat{\boldsymbol{\Theta}}_{\mathrm{km}} = \underset{\substack{\left\{\boldsymbol{\mu}^k\right\}_{k=1,\ldots,K} \\ \left\{z^i\right\}_{i=1,\ldots,n}}}{\arg\min} \sum_{k=1}^{K} \sum_{i:z^i=k} \left\|\mathbf{x}^i - \boldsymbol{\mu}^k\right\|_2^2$$

# The K-Means Objective

$$\hat{\mathbf{\Theta}}_{\mathrm{km}} = \underset{\substack{\{\boldsymbol{\mu}^k\}_{k=1,\ldots,K} \\ \{z^i\}_{i=1,\ldots,n}}}{\arg\min} \sum_{k=1}^{K} \sum_{i:z^i=k} \left\| \mathbf{x}^i - \boldsymbol{\mu}^k \right\|_2^2$$

**K-MEANS/LLOYD'S ALGORITHM**

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1\ldots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
   1. Let $z^{i,t} = \arg\min_k \left\| \mathbf{x}^i - \boldsymbol{\mu}^{k,t} \right\|_2^2$
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

$$\hat{\boldsymbol{\Theta}}_{\text{km}} = \operatorname*{arg\,min}_{\substack{\{\boldsymbol{\mu}^k\}_{k=1,\ldots,K} \\ \{z^i\}_{i=1,\ldots,n}}} \sum_{k=1}^{K} \sum_{i:z^i=k} \left\| \mathbf{x}^i - \boldsymbol{\mu}^k \right\|_2^2$$

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

## K-MEANS/LLOYD'S ALGORITHM

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1\ldots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
    1. Let $z^{i,t} = \operatorname*{arg\,min}_{k} \left\| \mathbf{x}^i - \boldsymbol{\mu}^{k,t} \right\|_2^2$
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

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

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 $\{\boldsymbol{\mu}^{k,0}\}_{k=1\ldots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
   1. Let $z^{i,t} = \arg\min_k \left\| \mathbf{x}^i - \boldsymbol{\mu}^{k,t} \right\|_2^2$
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!

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

$$\hat{\Theta}_{\mathrm{km}} = \underset{\substack{\{\boldsymbol{\mu}^k\}_{k=1,\ldots,K} \\ \{z^i\}_{i=1,\ldots,n}}}{\arg\min} \sum_{k=1}^{K} \sum_{i:z^i=k} \left\| \mathbf{x}^i - \boldsymbol{\mu}^k \right\|_2^2$$

**K-MEANS/LLOYD'S ALGORITHM**

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1\ldots K}$

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

   1. Let $z^{i,t} = \arg\min_k \left\| \mathbf{x}^i - \boldsymbol{\mu}^{k,t} \right\|_2^2$

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!

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

$$\hat{\boldsymbol{\Theta}}_{\mathrm{km}} = \underset{\substack{\{\boldsymbol{\mu}^k\}_{k=1,\ldots,K} \\ \{z^i\}_{i=1,\ldots,n}}}{\arg\min} \sum_{k=1}^{K} \sum_{i:z^i=k} \left\| \mathbf{x}^i - \boldsymbol{\mu}^k \right\|_2^2$$

Very scalable but sensitive to initialization!

**K-MEANS/LLOYD'S ALGORITHM**

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1\ldots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
    1. Let $z^{i,t} = \arg\min_k \left\| \mathbf{x}^i - \boldsymbol{\mu}^{k,t} \right\|_2^2$
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

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

NP-hard problem!

An FA approach to solving a data modelling task!

Alternates between updating $\{z^i\}$ and $\{\mu^k\}$
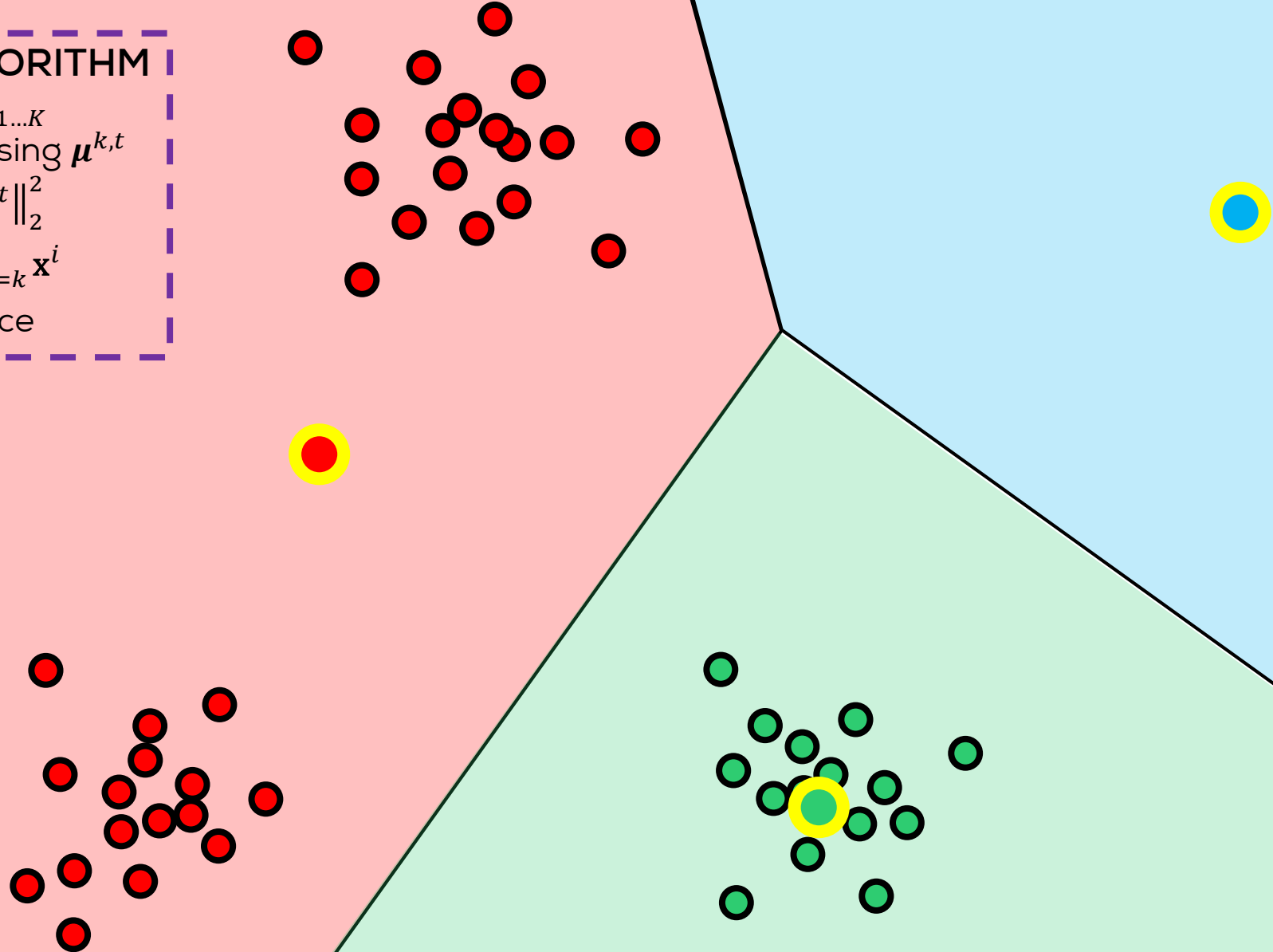
Very scalable but sensitive to initialization!

**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 $z^{i,t} = \underset{k}{\arg\min}\left\|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\right\|_2^2$

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++ initialization
1.　Sample $i_1 \sim [n]$, let $\boldsymbol{\mu}^{1,0} = \mathbf{x}^{i_1}$
2.　For k = 2,...K
　　- Sample $i_k \propto$ min distance from $\{\boldsymbol{\mu}^{1,0}, \dots, \boldsymbol{\mu}^{k-1,0}\}$
　　- Let $\boldsymbol{\mu}^{k,0} = \mathbf{x}^{i_k}$
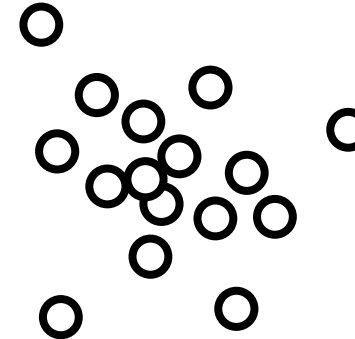
# K-Means Algorithm in action!

**K-MEANS/LLOYD'S ALGORITHM**

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1\ldots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$

   Let $z^{i,t} = \arg\min_k \left\|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\right\|_2^2$

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!

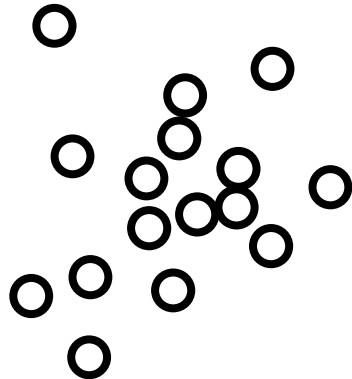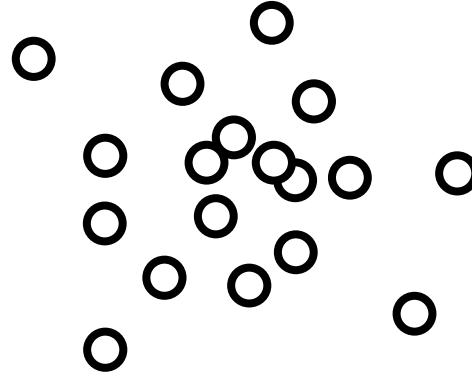# K-Means Algorithm in action!

# K-Means Algorithm in action!

**K-MEANS/LLOYD'S ALGORITHM**
1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1\ldots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
   Let $z^{i,t} = \arg\min_k \left\| \mathbf{x}^i - \boldsymbol{\mu}^{k,t} \right\|_2^2$
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

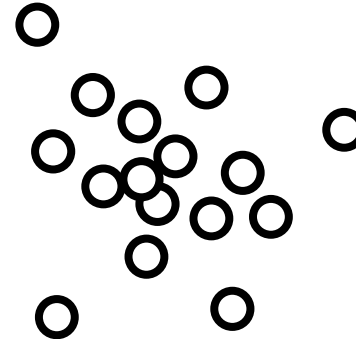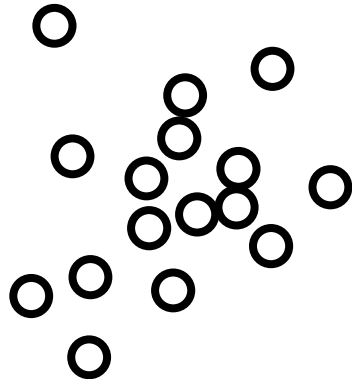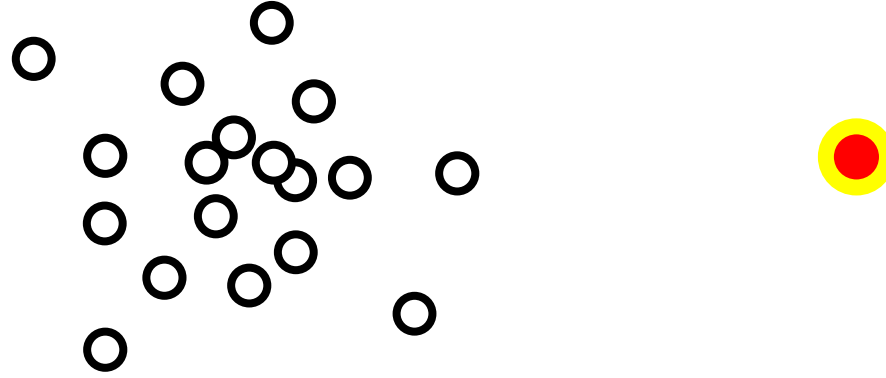# K-Means Algorithm in action!

**K-MEANS/LLOYD'S ALGORITHM**
1.	Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1\ldots K}$
2.	For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$

	Let $z^{i,t} = \arg\min_k \left\|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\right\|_2^2$

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 Algorithm in action!

# K-Means Algorithm in action!

**K-MEANS/LLOYD'S ALGORITHM**
1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1...K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$

   Let $z^{i,t} = \arg\min_{k}\left\|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\right\|_2^2$
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
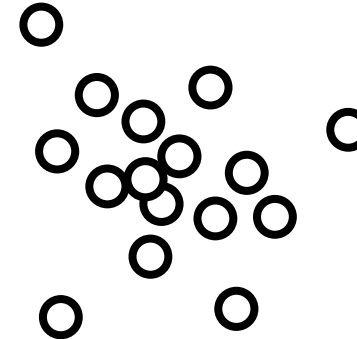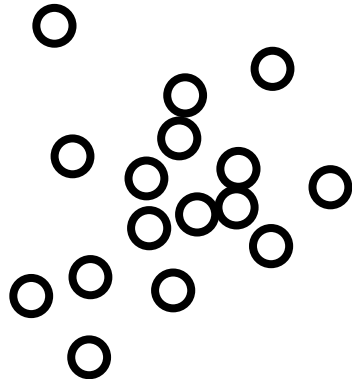
# K-Means Algorithm in action!

# K-Means Algorithm in action!
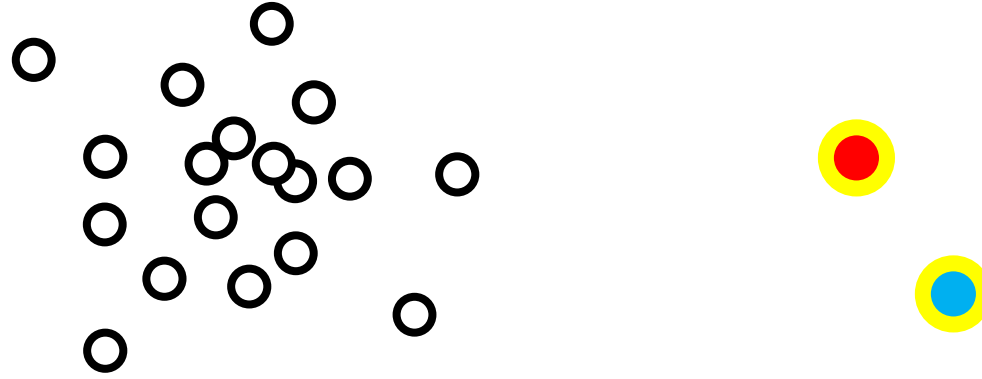
# K-Means Algorithm in action!

# K-Means Algorithm in action!

# K-Means Algorithm in action!

Stuck!!!

# K-Means Algorithm in action!

**K-MEANS/LLOYD'S ALGORITHM**
1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1\ldots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$

   Let $z^{i,t} = \arg\min_k \left\|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\right\|_2^2$
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

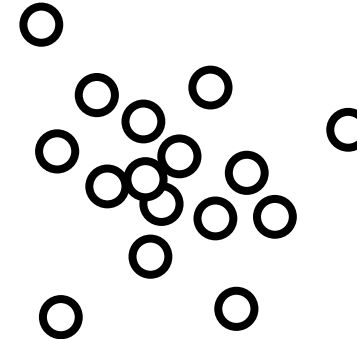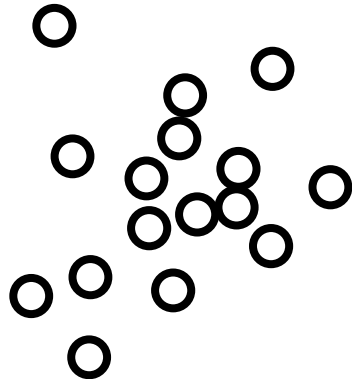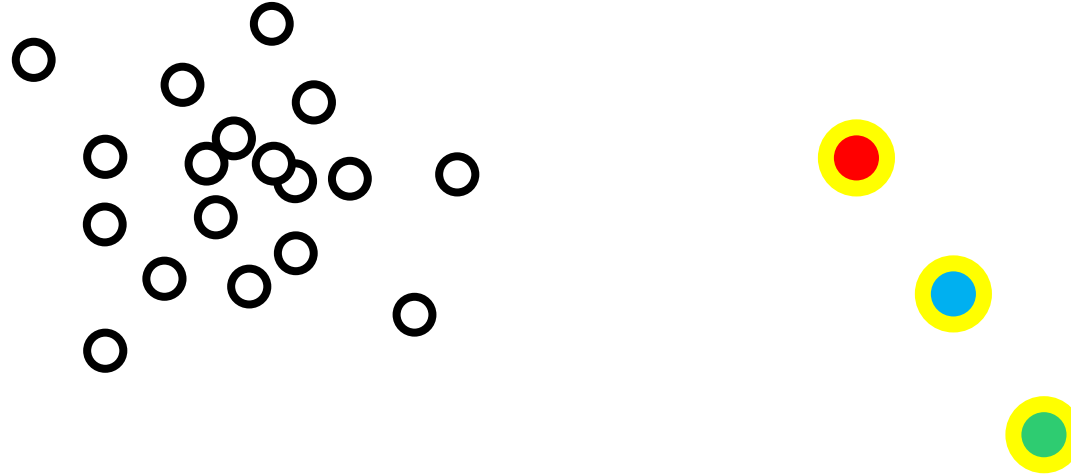# K-Means Algorithm in action!

**K-MEANS/LLOYD'S ALGORITHM**
1.  Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1\ldots K}$
2.  For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$

Let $z^{i,t} = \arg\min_{k}\left\|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\right\|_2^2$

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 Algorithm in action!

# 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 $z^{i,t} = \arg\min_{k}\left\|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\right\|_2^2$

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

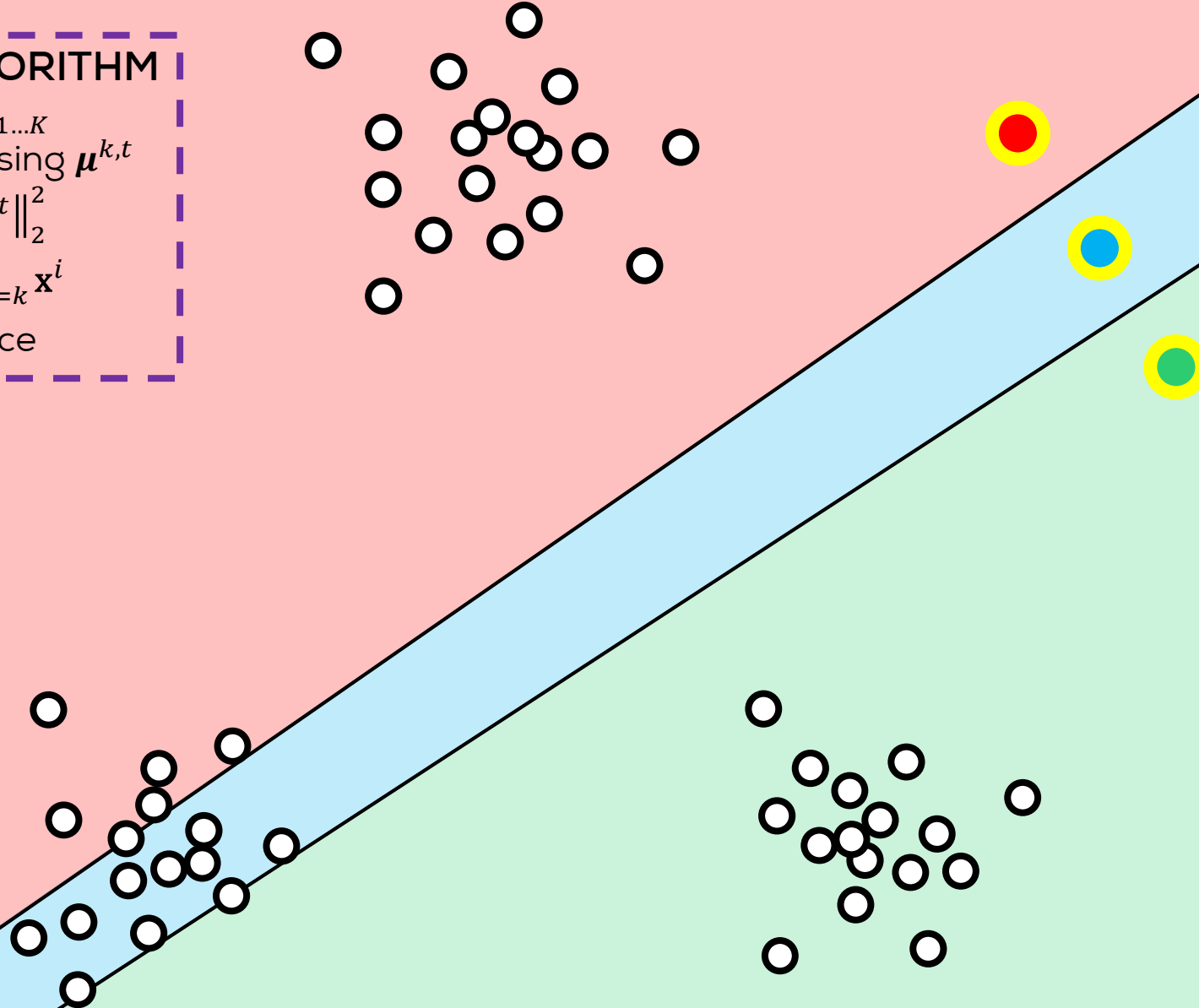# K-Means Algorithm in action!

**K-MEANS/LLOYD'S ALGORITHM**
1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1...K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
   Let $z^{i,t} = \arg\min_{k}\left\|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\right\|_2^2$
3. Update $\boldsymbol{\mu}^{k,t+1} = \frac{1}{n_k^t}\sum_{i:z^{i,t}=k}\mathbf{x}^i$
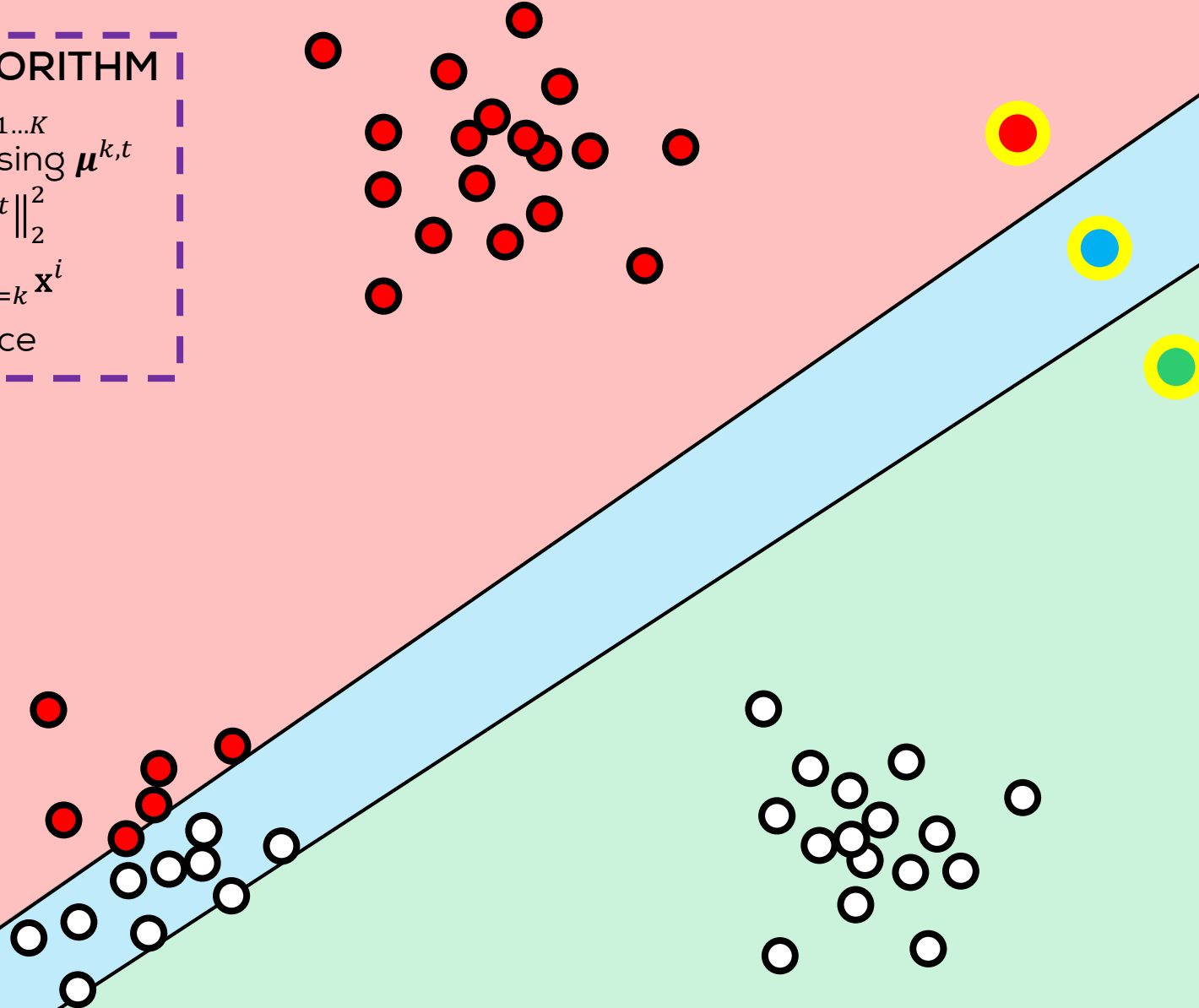4. Repeat until convergence

# K-Means Algorithm in action!

**K-MEANS/LLOYD'S ALGORITHM**

1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1...K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$

   Let $z^{i,t} = \arg\min_{k}\left\|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\right\|_2^2$

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
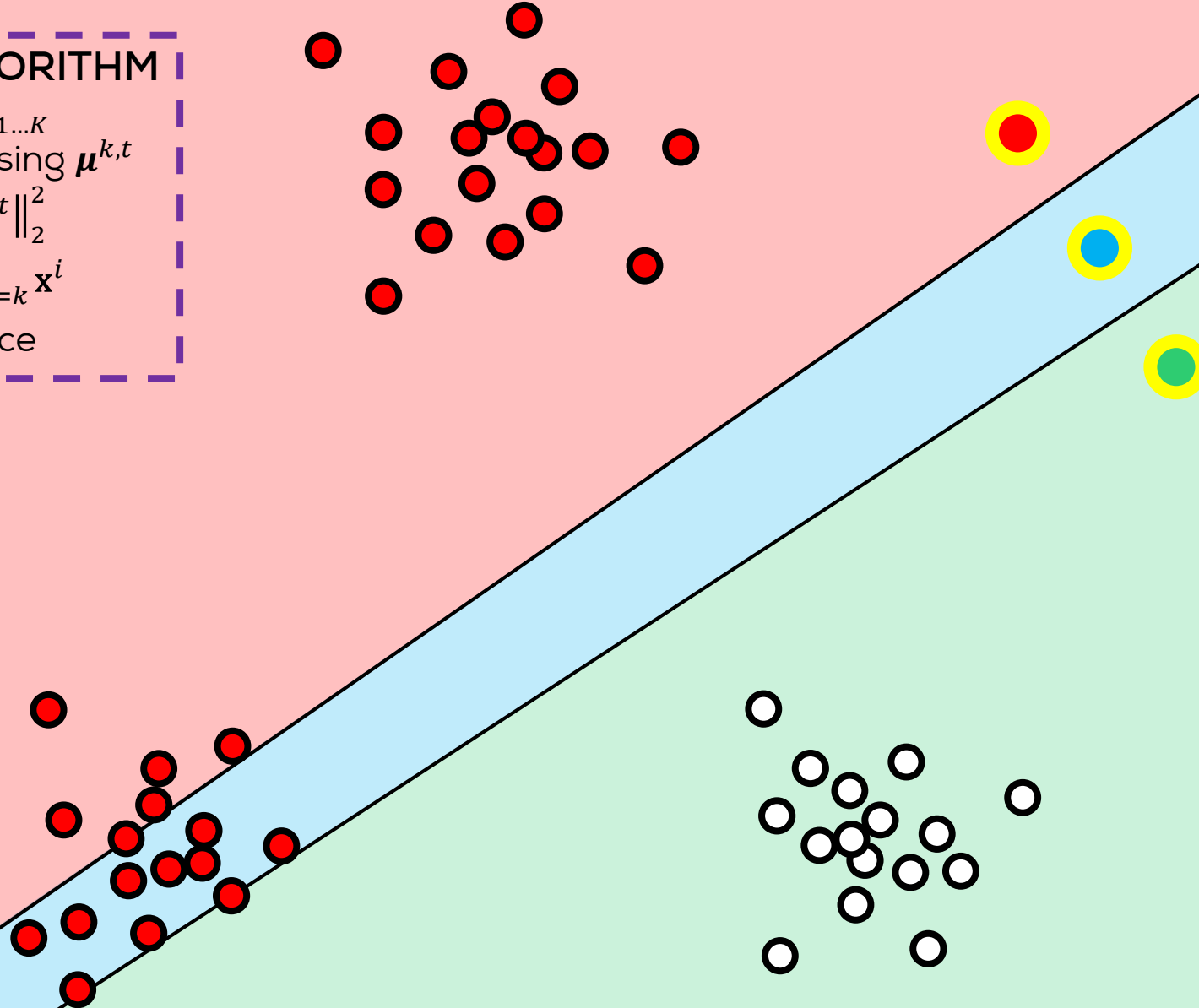
# K-Means Algorithm in action!

# K-Means Algorithm in action!
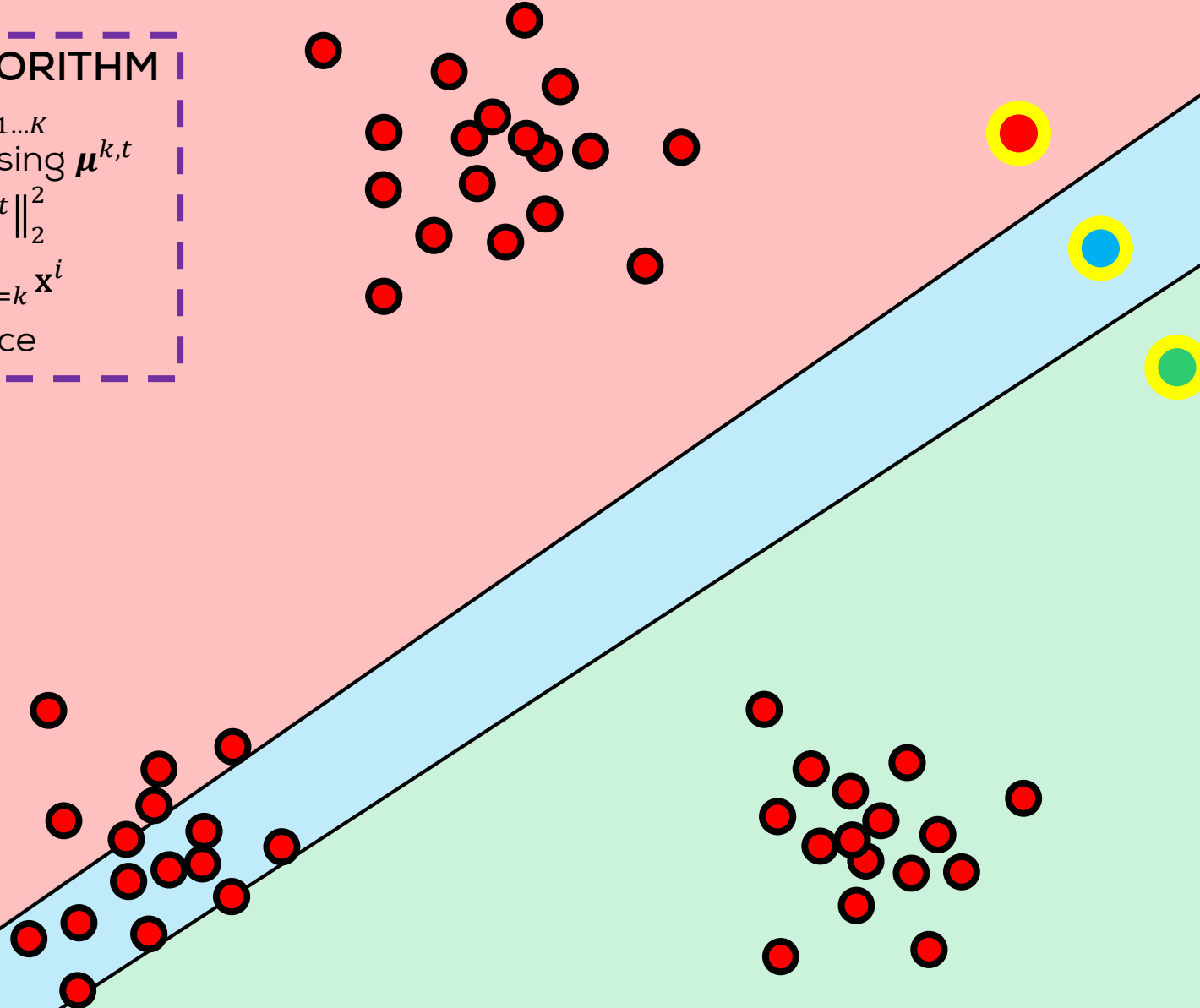
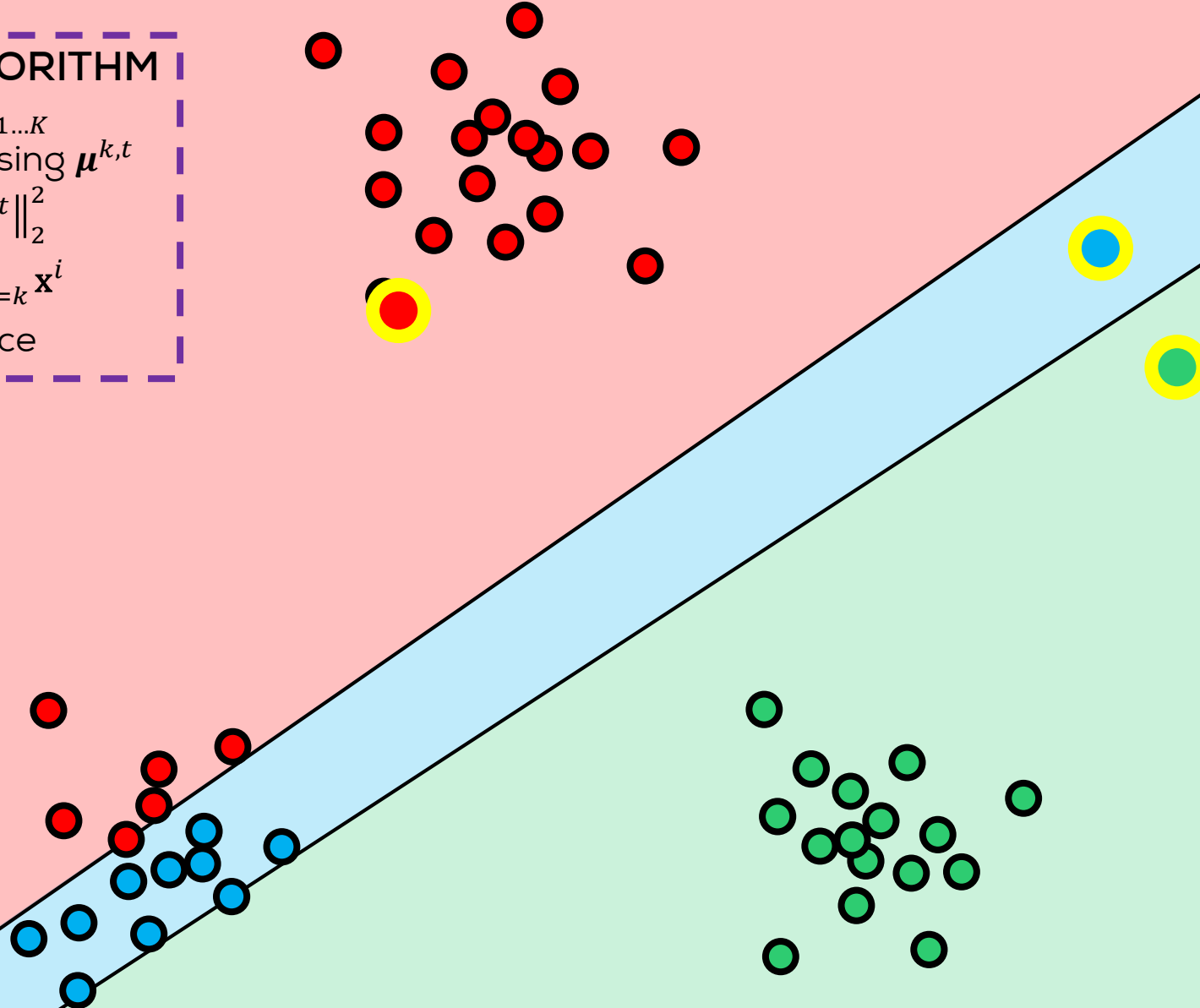# K-Means Algorithm in action!

# K-Means Algorithm in action!

# K-Means Algorithm in action!

# K-Means Algorithm in action!

# K-Means Algorithm in action!

# K-Means Algorithm in action!

# K-Means Algorithm in action!

# K-Means Algorithm in action!

# K-Means Algorithm in action!

**K-MEANS/LLOYD'S ALGORITHM**
1. Initialize means $\{\boldsymbol{\mu}^{k,0}\}_{k=1\ldots K}$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\mu}^{k,t}$
   Let $z^{i,t} = \arg\min_k \left\|\mathbf{x}^i - \boldsymbol{\mu}^{k,t}\right\|_2^2$
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

Stuck!!! … but at the global optimum ☺

# The Magic is not Always very Useful!

lotr.wikia.com

CS771: Intro to ML

# The Magic is not Always very Useful!

lotr.wikia.com

CS771: Intro to ML

# The Magic is not Always very Useful!

lotr.wikia.com

# The Magic is not Always very Useful!

lotr.wikia.com

CS771: Intro to ML

# The Magic is not Always very Useful!

lotr.wikia.com

# The Magic is not Always very Useful!

lotr.wikia.com

CS771: Intro to ML

# The Magic is not Always very Useful!

lotr.wikia.com

CS771: Intro to ML

# The Magic is not Always very Useful!

lotr.wikia.com

# The Magic is not Always very Useful!

# The Magic is not Always very Useful!



Apply the k-means algortihm

# The Magic is not Always very Useful!

# The Magic is not Always very Useful!

lotr.wikia.com

# The Magic is not Always very Useful!

lotr.wikia.com

CS771: Intro to ML

# The Magic is not Always very Useful!

Apply the k-means algortihm

?

# The Magic is not Always very Useful!

lotr.wikia.com

# The Magic is not Always very Useful!



Apply the k-means algortihm

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

?

The problem is NP hard for a reason!

lotr.wikia.com

# Soft Assignment

The EM algorithm

# A Ray of Hope

$$\hat{\mathbf{\Theta}}_{\mathrm{MLE}} = \arg\max_{\mathbf{\Theta}} \mathbb{P}\left[X \mid \mathbf{\Theta}\right]$$

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

# A Ray of Hope

$$\hat{\boldsymbol{\Theta}}_{\mathrm{MLE}} = \arg\max_{\boldsymbol{\Theta}} \mathbb{P}\left[X \mid \boldsymbol{\Theta}\right]$$

**ALTERNATING OPTIMIZATION**

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

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

# A Ray of Hope

$$\hat{\boldsymbol{\Theta}}_{\mathrm{MLE}} = \arg\max_{\boldsymbol{\Theta}} \mathbb{P}\left[X \mid \boldsymbol{\Theta}\right]$$

**ALTERNATING OPTIMIZATION**

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

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

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

# A Ray of Hope

$$\hat{\mathbf{\Theta}}_{\mathrm{MLE}} = \arg\max_{\mathbf{\Theta}} \mathbb{P}\left[X \mid \mathbf{\Theta}\right]$$

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

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

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

# A Ray of Hope

$$\hat{\boldsymbol{\Theta}}_{\mathrm{MLE}} = \arg\max_{\boldsymbol{\Theta}} \mathbb{P}\left[X \mid \boldsymbol{\Theta}\right]$$

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

$$z^{i,t} = \arg\max_{k \in [K]} \mathbb{P}\left[k \mid \mathbf{x}^i, \boldsymbol{\Theta}^t\right] = \arg\max_{k \in [K]} \quad \boldsymbol{\pi}_k^t \quad \cdot \quad \mathbb{P}\left[\mathbf{x}^i \mid k, \boldsymbol{\Theta}^t\right]$$

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

# A Ray of Hope

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

**ALTERNATING OPTIMIZATION**

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

Bayes Rule!

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

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

# A Ray of Hope

$$\hat{\boldsymbol{\Theta}}_{\mathrm{MLE}} = \arg\max_{\boldsymbol{\Theta}} \mathbb{P}\left[X \mid \boldsymbol{\Theta}\right]$$

**ALTERNATING OPTIMIZATION**

1. Initialize $\boldsymbol{\Theta}^0$
2. For $i \in [n]$, update $z^{i,t}$ using $\boldsymbol{\Theta}^t$
   $$\phantom{xxx}^{-1} = \arg\max_{\boldsymbol{\Theta}} \mathbb{P}\left[X, \{z^{i,t}\} \mid \boldsymbol{\Theta}\right]$$
   until convergence

May be throwing away a lot of information!

Bayes Rule!

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

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

# A Ray of Hope

$$\hat{\boldsymbol{\Theta}}_{\mathrm{MLE}} = \arg\max_{\boldsymbol{\Theta}} \mathbb{P}\left[X \mid \boldsymbol{\Theta}\right]$$

E.g. $\mathbb{P}\left[\ \textcolor{red}{\bullet}\ |\mathbf{x}^i,\boldsymbol{\Theta}^t\right] = 0.5,$
$\mathbb{P}\left[\ \textcolor{cyan}{\bullet}\ |\mathbf{x}^i,\boldsymbol{\Theta}^t\right] = 0.4,$
$\mathbb{P}\left[\ \textcolor{green}{\bullet}\ |\mathbf{x}^i,\boldsymbol{\Theta}^t\right] = 0.1,$

**ALTERNATING OPTIMIZATION**

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

May be throwing away a lot of information!

Bayes Rule!

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

$$\boldsymbol{\Theta}^t = \left\{ \boldsymbol{\pi}^t, \left\{ \textcolor{red}{\boldsymbol{\mu}^{1,t}}, \textcolor{green}{\boldsymbol{\mu}^{2,t}}, \textcolor{cyan}{\boldsymbol{\mu}^{3,t}} \right\}, \left\{ \textcolor{red}{\Sigma^{1,t}}, \textcolor{green}{\Sigma^{2,t}}, \textcolor{cyan}{\Sigma^{3,t}} \right\} \right\}$$

# A Ray of Hope

$$\hat{\boldsymbol{\Theta}}_{\mathrm{MLE}} = \arg\max_{\boldsymbol{\Theta}} \mathbb{P}\left[X \mid \boldsymbol{\Theta}\right]$$

E.g. $\mathbb{P}\left[\,\textcolor{red}{\bullet}\,|\mathbf{x}^i,\boldsymbol{\Theta}^t\right] = 0.5,$
$\mathbb{P}\left[\,\textcolor{cyan}{\bullet}\,|\mathbf{x}^i,\boldsymbol{\Theta}^t\right] = 0.4,$
$\mathbb{P}\left[\,\textcolor{green}{\bullet}\,|\mathbf{x}^i,\boldsymbol{\Theta}^t\right] = 0.1,$

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

**ALTERNATING ... ATION**

... alize ...

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

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

... til convergence

May be throwing away a lot of information!

Bayes Rule!

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

$$\boldsymbol{\Theta}^t = \left\{ \boldsymbol{\pi}^t, \left\{ \textcolor{red}{\boldsymbol{\mu}^{1,t}}, \textcolor{green}{\boldsymbol{\mu}^{2,t}}, \textcolor{cyan}{\boldsymbol{\mu}^{3,t}} \right\}, \left\{ \textcolor{red}{\Sigma^{1,t}}, \textcolor{green}{\Sigma^{2,t}}, \textcolor{cyan}{\Sigma^{3,t}} \right\} \right\}$$

# A Ray of Hope

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

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

E.g. $\mathbb{P}\left[\,\textcolor{red}{\bullet}\,|\mathbf{x}^i, \boldsymbol{\Theta}^t\right] = 0.5,$
$\mathbb{P}\left[\,\textcolor{cyan}{\bullet}\,|\mathbf{x}^i, \boldsymbol{\Theta}^t\right] = 0.4,$
$\mathbb{P}\left[\,\textcolor{green}{\bullet}\,|\mathbf{x}^i, \boldsymbol{\Theta}^t\right] = 0.1,$

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

**ERNA**ATION

lize

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

$^{-1} = \arg\max_{\boldsymbol{\Theta}} \mathbb{P}\left[X, \{z^{i,t}\} \mid \boldsymbol{\Theta}\right]$

May be throwing away a lot of information!

til convergence

Bayes Rule!

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

$$\boldsymbol{\Theta}^t = \left\{\boldsymbol{\pi}^t, \left\{\textcolor{red}{\boldsymbol{\mu}^{1,t}}, \textcolor{green}{\boldsymbol{\mu}^{2,t}}, \textcolor{cyan}{\boldsymbol{\mu}^{3,t}}\right\}, \left\{\textcolor{red}{\Sigma^{1,t}}, \textcolor{green}{\Sigma^{2,t}}, \textcolor{cyan}{\Sigma^{3,t}}\right\}\right\}$$

# A Ray of Hope

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

**ALTERNATING OPTIMIZATION**

1. Initialize $\Theta^0$
2. For $i \in [n]$ update $z^{i,t}$ using $\Theta^t$
3. $\Theta^{t+1} = \arg\max_{\Theta} \mathbb{P}\left[X, \{z^{i,t}\} \mid \Theta\right]$

Repeat until convergence

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

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

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

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

May be throwing away a lot of information!

Bayes Rule!

# Weighted MLE

# Weighted MLE



Data points have weights

$\mu^2$    $\mu^3$

$\gamma^1$    $\gamma^2$    $\gamma^3$    $\gamma^4$    $\gamma^n$

# Weighted MLE



Data points have weights

$\mu^2$  $\mu^3$

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

# Weighted MLE



Data points have weights

$\boldsymbol{\mu}^2 \qquad \boldsymbol{\mu}^3$

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

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

# Weighted MLE



Data points have weights

$\boldsymbol{\mu}^2$    $\boldsymbol{\mu}^3$

$\gamma^1$    $\gamma^2$    $\gamma^3$    $\gamma^4$    $\gamma^n$

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

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

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

# Weighted MLE



Data points have weights

$\boldsymbol{\mu}^2$   $\boldsymbol{\mu}^3$

$\gamma^1$ $\gamma^2$ $\gamma^3$ $\gamma^4$ $\gamma^n$

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

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

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

Take log and apply 1st order optimality

# Weighted MLE



Data points have weights

$\mu^2$   $\mu^3$

$\gamma^1$  $\gamma^2$  $\gamma^3$  $\gamma^4$  $\gamma^n$

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

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

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

Take log and apply 1st order optimality

$$\boldsymbol{\pi}_1^{\mathrm{MLE}} = \frac{\text{\# eff. red emails}}{\text{\# eff. total emails}} = \frac{\tilde{n}_r}{n} = \frac{\sum_{i:z^i=\bullet} \gamma^i}{\sum_{j=1}^{n} \gamma^j}$$

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

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

# Weighted MLE



Data points have weights

$\boldsymbol{\mu}^2$    $\boldsymbol{\mu}^3$

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

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

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

$$\boldsymbol{\pi}_1^{\mathrm{MLE}} = \frac{\text{\# eff. red emails}}{\text{\# eff. total emails}} = \frac{\tilde{n}_r}{n} = \frac{\sum_{i:z^i=\bullet} \gamma^i}{\sum_{j=1}^{n} \gamma^j}$$

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

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

Take log and apply 1st order optimality

Exercise!

# Weighted MLE



$\boldsymbol{\mu}^2$      $\boldsymbol{\mu}^3$

Data points have weights

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

$\gamma^1 \quad \gamma^2 \quad \gamma^3 \quad \gamma^4 \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \gamma^n$

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

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

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

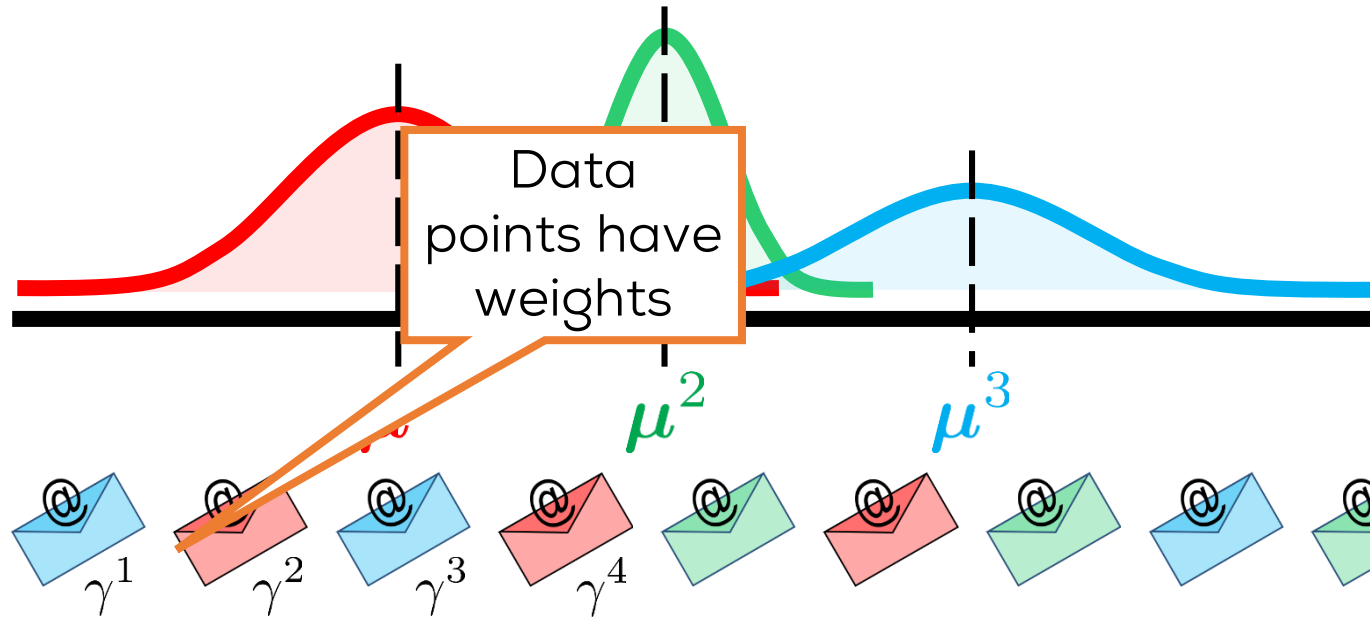$$\boldsymbol{\pi}_1^{\mathrm{MLE}} = \frac{\text{\# eff. red emails}}{\text{\# eff. total emails}} = \frac{\tilde{n}_r}{n} = \frac{\sum_{i:z^i=\bullet} \gamma^i}{\sum_{j=1}^{n} \gamma^j}$$

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

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

Take log and apply 1ˢᵗ order optimality

Exercise!

# Hard Alternating Minimization

# Soft Alternating Minimization

# Soft Alternating Minimization

## ALTERNATING OPTIMIZATION

1. For $i \in [n]$, create $k$ copies of the data point
   1. Let $\mathbf{x}^i \rightarrow \left\{ \mathbf{x}^{\{i,1\}}, \mathbf{x}^{\{i,2\}}, \ldots, \mathbf{x}^{\{i,k\}} \right\}$
   2. Assign the $k$-th copy label $k$ i.e. $z^{\{i,k\}} = k$
2. Initialize $\boldsymbol{\Theta}^0$
3. Update weights $\gamma^{i,k,t}$ using $\boldsymbol{\Theta}^t$

   1. Let $\gamma^{i,k,t} = \mathbb{P}[k \mid \mathbf{x}^i, \boldsymbol{\Theta}^t] = \dfrac{\boldsymbol{\pi}_k^t \cdot \mathcal{N}(\ma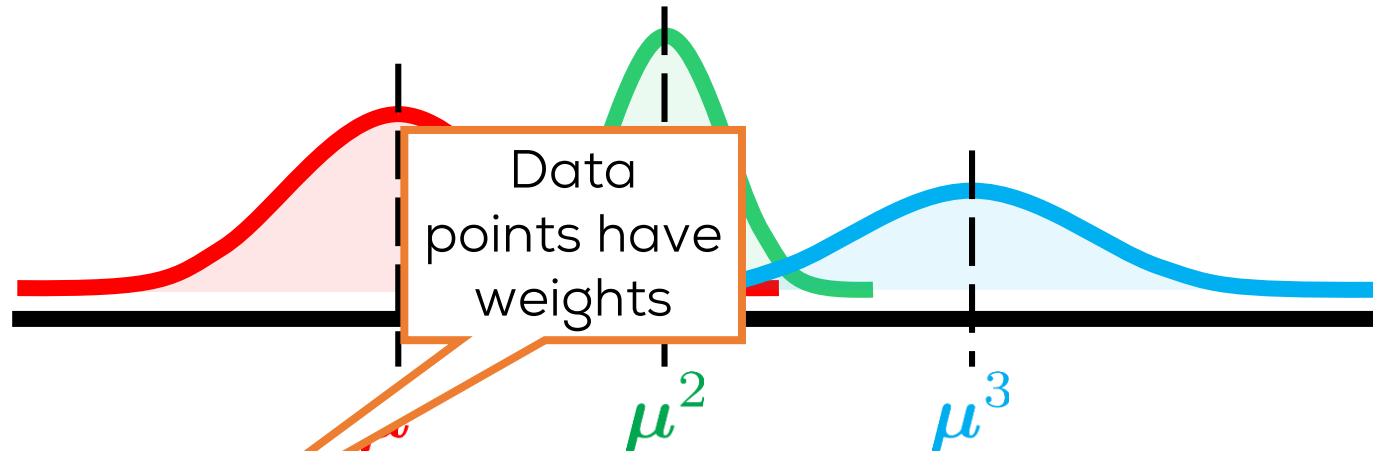thbf{X}^i \mid \boldsymbol{\mu}^{k,t}, \Sigma^{k,t})}{\sum_j \boldsymbol{\pi}_j^t \cdot \mathcal{N}(\mathbf{X}^i \mid \boldsymbol{\mu}^{j,t}, \Sigma^{j,t})}$

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

4. Update $\boldsymbol{\Theta}^{t+1} = \arg\max_{\boldsymbol{\Theta}} \mathbb{P}\left[ \left\{ x^{\{i,k\}} \right\}, \left\{ z^{\{i,k\}} \right\}, \left\{ \gamma^{\{i,k,t\}} \right\} \mid \boldsymbol{\Theta} \right]$
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 $\mathbf{\Theta}^0$
3. Update weights $\gamma^{i,k,t}$ using $\mathbf{\Theta}^t$

    1. Let $\gamma^{i,k,t} = \mathbb{P}[k \mid \mathbf{x}^i, \mathbf{\Theta}^t] = \dfrac{\boldsymbol{\pi}_k^t \cdot \mathcal{N}(\mathbf{X}^i \mid \boldsymbol{\mu}^{k,t}, \Sigma^{k,t})}{\sum_j \boldsymbol{\pi}_j^t \cdot \mathcal{N}(\mathbf{X}^i \mid \boldsymbol{\mu}^{j,t}, \Sigma^{j,t})}$

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

**ALTERNATING OPTIMIZATION**

1. For $i \in [n]$, create $k$ copies of the data point
   1. Let $\mathbf{x}^i \rightarrow \left\{\mathbf{x}^{\{i,1\}}, \mathbf{x}^{\{i,2\}}, \dots, \mathbf{x}^{\{i,k\}}\right\}$
   2. Assign the $k$-th copy label $k$ i.e. $z^{\{i,k\}} = k$
2. Initialize $\mathbf{\Theta}^0$
3. Update weights $\gamma^{i,k,t}$ using $\mathbf{\Theta}^t$

   1. Let $\gamma^{i,k,t} = \mathbb{P}[k \mid \mathbf{x}^i, \mathbf{\Theta}^t] = \dfrac{\boldsymbol{\pi}_k^t \cdot \mathcal{N}(\mathbf{X}^i \mid \boldsymbol{\mu}^{k,t}, \Sigma^{k,t})}{\sum_j \boldsymbol{\pi}_j^t \cdot \mathcal{N}(\mathbf{X}^i \mid \boldsymbol{\mu}^{j,t}, \Sigma^{j,t})}$

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

5. Repeat until convergence

Hard–AM: all weight was on a single cluster

Distribute a unit weight across clusters

$\sum_j \gamma^{\{i,j,t\}} = 1$ for all $i$ 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 $\{\boldsymbol{\mu}^{k,0}\}_{k=1,\ldots K}$
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 $\{\boldsymbol{\mu}^{k,0}\}_{k=1,\ldots K}$
2. For all $i$, update $\gamma^{i,k,t}$ using $\boldsymbol{\mu}^{k,t}$

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

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 $\boldsymbol{\mu}^{k,t}$

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

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

# Mixed Regression

# Mixed Regression

# Mixed Regression

# Mixed Regression

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

$$= (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1} \mathbf{X}\mathbf{y}$$

# Mixed Regression



$$\hat{\mathbf{w}}_{\mathrm{MAP}} = \arg\min \sum_{i=1}^{n} \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle\right)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$= (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1}\mathbf{X}\mathbf{y}$$

# Mixed Regression



$$\hat{\mathbf{w}}_{\mathrm{MAP}} = \arg\min \sum_{i=1}^{n} \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle\right)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$= (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1}\mathbf{X}\mathbf{y}$$

# Mixed Regression

$$\hat{\mathbf{w}}_{\mathrm{MAP}} = \arg\min \sum_{i=1}^{n} \left( y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$= (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1}\mathbf{X}\mathbf{y}$$



$y$

$x$

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

$$\hat{\mathbf{w}}_{\mathrm{MAP}} = \arg\min \sum_{i=1}^{n} \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle\right)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$= (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1}\mathbf{X}\mathbf{y}$$

$y$

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

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

- 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

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

$z^i \in \{0,1\}$ 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

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

$z^i \in \{0,1\}$ 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 $\mathbf{\Theta}^0 = \left\{\mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}}\right\}$

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

   1. Let $z^{i,t} = \arg\max_{k \in \{0,1\}} \mathcal{N}\left(y^i \mid \mathbf{x}^i, \boldsymbol{w}^{k,t}\right)$

3. Update $\mathbf{w}^{\{t+1,k\}} = \arg\min_{\mathbf{W}} \sum_{i:z^{\{i,t\}}=k} \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle\right)^2 + \frac{1}{2}\|\mathbf{w}\|_2^2$

4. Set $\mathbf{\Theta}^{t+1} = \left\{\mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}}\right\}$

5. Repeat until convergence

# Hard Mixed Regression

ALTERNATING OPTIMIZATION

1. Initialize $\boldsymbol{\Theta}^0 = \{\mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}}\}$

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

   1. Let $z^{i,t} = \arg \max_{k \in \{0,1\}} \mathcal{N}\left(y^i \mid \mathbf{x}^i, \boldsymbol{w}^{k,t}\right)$

3. Update $\mathbf{w}^{\{t+1,k\}} = \arg \min_{\mathbf{W}} \sum_{i:z^{\{i,t\}}=k} \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle\right)^2 + \frac{1}{2}\|\mathbf{w}\|_2^2$

4. Set $\boldsymbol{\Theta}^{t+1} = \{\mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}}\}$

5. Repeat until convergence

*Exercise: verify these updates*

# Hard Mixed Regression

## ALTERNATING OPTIMIZATION

1. Initialize $\mathbf{\Theta}^0 = \left\{ \mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}} \right\}$

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

   1. Let $z^{i,t} = \arg \min_{k \in \{0,1\}} \left| y^i - \left\langle \mathbf{w}^{t,k}, \mathbf{x}^i \right\rangle \right|$

3. Update $\mathbf{w}^{\{t+1,k\}} = \arg \min_{\mathbf{W}} \sum_{i: z^{\{i,t\}}=k} \left( y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \frac{1}{2} \|\mathbf{w}\|_2^2$

4. Set $\mathbf{\Theta}^{t+1} = \left\{ \mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}} \right\}$

5. Repeat until convergence

*Exercise: verify these updates*

# Hard Mixed Regression

Assign to the "closest" line!

## ALTERNATING OPTIMIZATION

1. Initialize $\mathbf{\Theta}^0 = \{\mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}}\}$
2. For $i \in [n]$, update $z^{i,t}$ using $\mathbf{\Theta}^t$
   1. Let $z^{i,t} = \arg \min_{k \in \{0,1\}} |y^i - \langle \mathbf{w}^{t,k}, \mathbf{x}^i \rangle|$

3. Update $\mathbf{w}^{\{t+1,k\}} = \arg \min_{\mathbf{W}} \sum_{i:z^{\{i,t\}}=k} \left( y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \frac{1}{2} \|\mathbf{w}\|_2^2$
4. Set $\mathbf{\Theta}^{t+1} = \{\mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}}\}$
5. Repeat until convergence

*Exercise: verify these updates*

# Hard Mixed Regression

## ALTERNATING OPTIMIZATION

1. Initialize $\boldsymbol{\Theta}^0 = \{\mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}}\}$

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

   1. Let $z^{i,t} = \arg\min_{k \in \{0,1\}} \left| y^i - \langle \mathbf{w}^{t,k}, \mathbf{x}^i \rangle \right|$

3. Update $\mathbf{w}^{\{t+1,k\}} = \arg\min_{\mathbf{W}} \sum_{i:z^{\{i,t\}}=k} \left( y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \frac{1}{2} \|\mathbf{w}\|_2^2$

4. Set $\boldsymbol{\Theta}^{t+1} = \{\mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}}\}$

5. Repeat until convergence

Assign to the "closest" line!

In k-means, we assigned to the closest mean

Exercise: verify these updates

# Weighted Regression

# Weighted Regression

$$\hat{\mathbf{w}}_{\mathrm{MAP}} = \arg\min \sum_{i=1}^{n} \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle\right)^2 + \lambda \cdot \|\mathbf{w}\|_2^2$$

$$= (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1}\mathbf{X}\mathbf{y}$$

# Weighted Regression

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

$$= (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1}\mathbf{X}\mathbf{y}$$

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

$$= (M + \lambda \cdot I)^{-1}\mathbf{b}$$

$$M = \sum_{i=1}^{n} \gamma_i \mathbf{x}^i (\mathbf{x}^i)^\top$$

$$\mathbf{b} = \sum_{i=1}^{n} \gamma_i y_i \cdot \mathbf{x}^i$$

# Weighted Regression

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

$$= (\mathbf{X}\mathbf{X}^\top + \lambda \cdot I)^{-1}\mathbf{X}\mathbf{y}$$

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

$$= (M + \lambda \cdot I)^{-1}\mathbf{b}$$

$$M = \sum_{i=1}^{n} \gamma_i \mathbf{x}^i (\mathbf{x}^i)^\top$$

$$\mathbf{b} = \sum_{i=1}^{n} \gamma_i y_i \cdot \mathbf{x}^i$$

Exercise

# Soft Mixed Regression

## SOFT ALTERNATING OPTIMIZATION

1. Initialize $\mathbf{\Theta}^0 = \left\{ \mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}} \right\}$

2. For $i \in [n]$, update $\gamma^{\{i,k,t\}}$ using $\mathbf{\Theta}^t$

   1. Let $c^{\{i,k,t\}} = \exp\left( -\frac{\left(y^i - \langle \mathbf{w}^{\{t,k\}}, \mathbf{x}^i \rangle \right)^2}{2} \right)$

   2. Let $\gamma^{\{i,k,t\}} = \frac{c^{\{i,k,t\}}}{c^{\{i,0,t\}} + c^{\{i,1,t\}}}$

3. Update $\mathbf{w}^{\{t+1,k\}} = \arg\min_{\mathbf{w}} \sum_i \gamma^{\{i,k,t\}} \cdot \left( y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \frac{1}{2} \|\mathbf{w}\|_2^2$

4. Set $\mathbf{\Theta}^{t+1} = \left\{ \mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}} \right\}$

5. Repeat until convergence

# Soft Mixed Regression

## SOFT ALTERNATING OPTIMIZATION

1. Initialize $\boldsymbol{\Theta}^0 = \{\mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}}\}$

2. For $i \in [n]$, update $\gamma^{\{i,k,t\}}$ using $\boldsymbol{\Theta}^t$

    1. Let $c^{\{i,k,t\}} = \exp\left(-\dfrac{\left(y^i - \langle \mathbf{w}^{\{t,k\}}, \mathbf{x}^i \rangle\right)^2}{2}\right)$

    2. Let $\gamma^{\{i,k,t\}} = \dfrac{c^{\{i,k,t\}}}{c^{\{i,0,t\}} + c^{\{i,1,t\}}}$

3. Update $\mathbf{w}^{\{t+1,k\}} = \arg\min_{\mathbf{w}} \sum_i \gamma^{\{i,k,t\}} \cdot \left(y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle\right)^2 + \dfrac{1}{2} \|\mathbf{w}\|_2^2$

4. Set $\boldsymbol{\Theta}^{t+1} = \{\mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}}\}$

5. Repeat until convergence

*Exercise: derive these updates*

# Soft Mixed Regression

## SOFT ALTERNATING OPTIMIZATION

1. Initialize $\boldsymbol{\Theta}^0 = \left\{ \mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}} \right\}$
2. For $i \in [n]$, update $\gamma^{\{i,k,t\}}$ using $\boldsymbol{\Theta}^t$

   $\boxed{\propto \mathbb{P}\left[ k \mid y^i, \mathbf{x}^i, \boldsymbol{\Theta}^t \right]}$

   1. Let $c^{\{i,k,t\}} = \exp\left( -\dfrac{\left( y^i - \langle \mathbf{w}^{\{t,k\}}, \mathbf{x}^i \rangle \right)^2}{2} \right)$
   2. Let $\gamma^{\{i,k,t\}} = \dfrac{c^{\{i,k,t\}}}{c^{\{i,0,t\}} + c^{\{i,1,t\}}}$
3. Update $\mathbf{w}^{\{t+1,k\}} = \arg\min_{\mathbf{w}} \sum_i \gamma^{\{i,k,t\}} \cdot \left( y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \dfrac{1}{2} \|\mathbf{w}\|_2^2$
4. Set $\boldsymbol{\Theta}^{t+1} = \left\{ \mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}} \right\}$
5. Repeat until convergence

*Exercise: derive these updates*

# Soft Mixed Regression
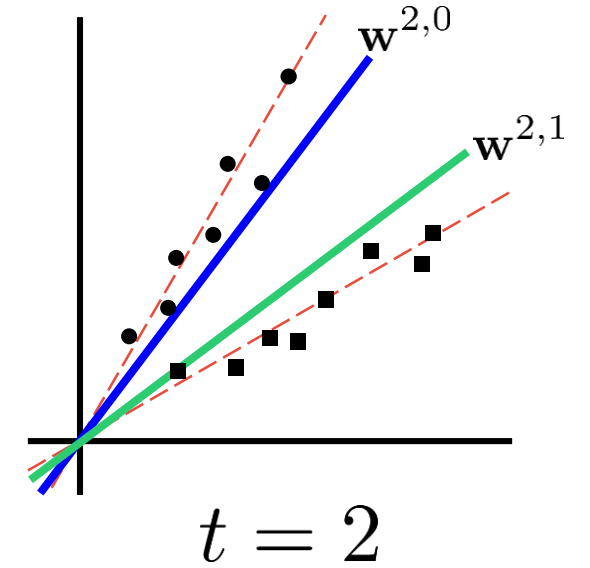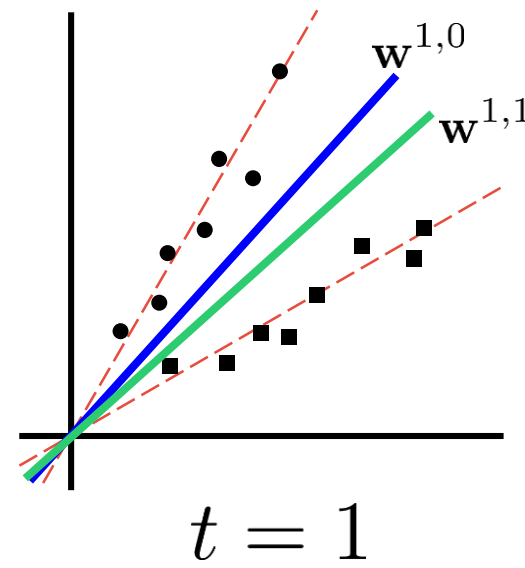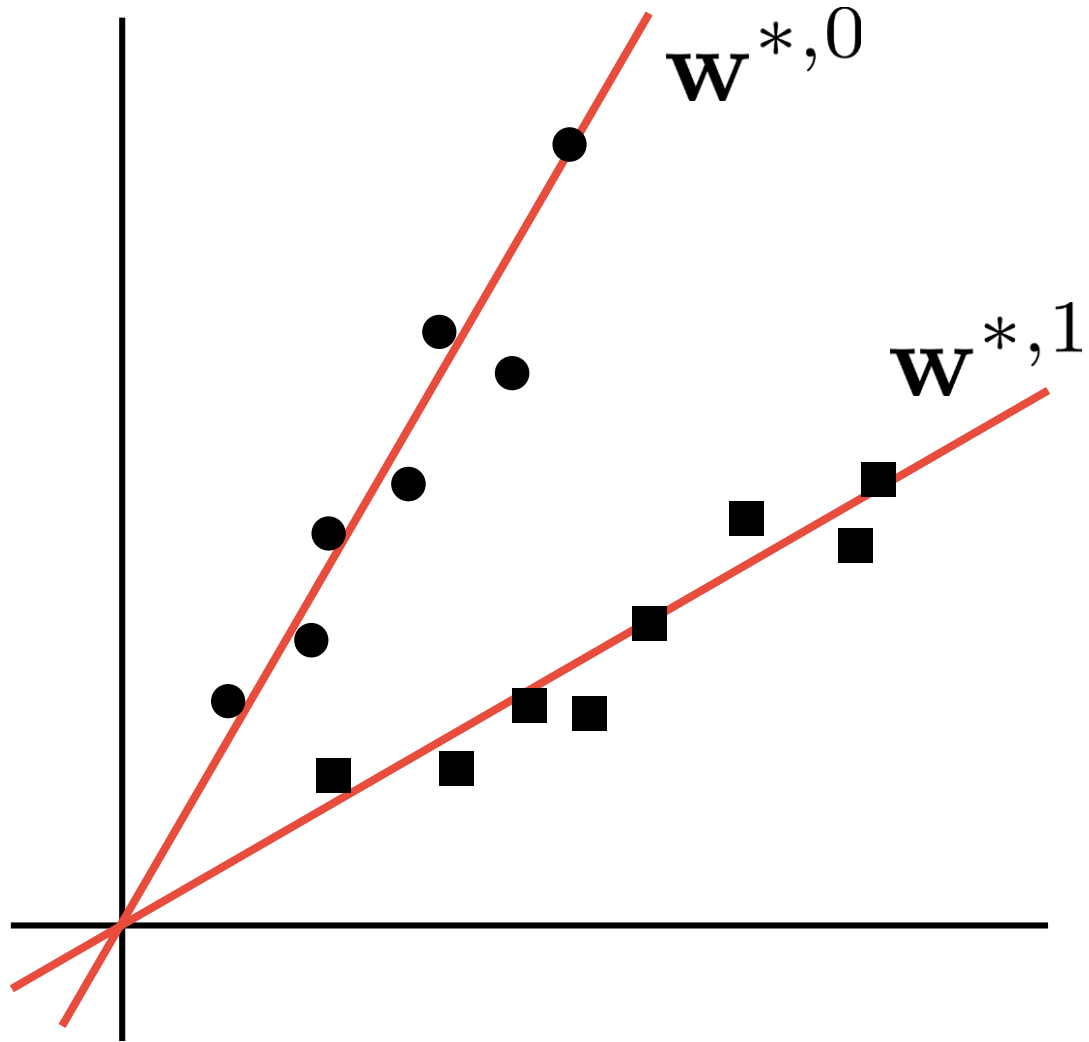
## SOFT ALTERNATING OPTIMIZATION

1. Initialize $\mathbf{\Theta}^0 = \left\{ \mathbf{w}^{\{0,0\}}, \mathbf{w}^{\{0,1\}} \right\}$
2. For $i \in [n]$, update $\gamma^{\{i,k,t\}}$ using $\mathbf{\Theta}^t$

    $\boxed{\propto \mathbb{P}\left[ k \mid y^i, \mathbf{x}^i, \mathbf{\Theta}^t \right]}$

    1. Let $c^{\{i,k,t\}} = \exp\left( -\dfrac{\left( y^i - \langle \mathbf{w}^{\{t,k\}}, \mathbf{x}^i \rangle \right)^2}{2} \right)$

    2. Let $\gamma^{\{i,k,t\}} = \dfrac{c^{\{i,k,t\}}}{c^{\{i,0,t\}} + c^{\{i,1,t\}}}$ $\boxed{= \mathbb{P}\left[ k \mid y^i, \mathbf{x}^i, \mathbf{\Theta}^t \right]}$

3. Update $\mathbf{w}^{\{t+1,k\}} = \arg\min_{\mathbf{w}} \sum_i \gamma^{\{i,k,t\}} \cdot \left( y^i - \langle \mathbf{w}, \mathbf{x}^i \rangle \right)^2 + \dfrac{1}{2} \|\mathbf{w}\|_2^2$

4. Set $\mathbf{\Theta}^{t+1} = \left\{ \mathbf{w}^{\{t+1,0\}}, \mathbf{w}^{\{t+1,1\}} \right\}$

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