

# Introduction to Data Mining

## 05 - Clustering

Benjamin Paaßen

WS 2023/2024, Bielefeld University

- ▶ Sheet 02 is out today

- ▶ Sheet 02 is out today
- ▶ We noticed some difficulties in assigning e-mail submissions to groups

- ▶ Sheet 02 is out today
- ▶ We noticed some difficulties in assigning e-mail submissions to groups
- ⇒ Please use moodle exclusively for your submissions, from now on

1. K-Means algorithm

2. Gaussian Mixture Models

## Motivation: Education Example

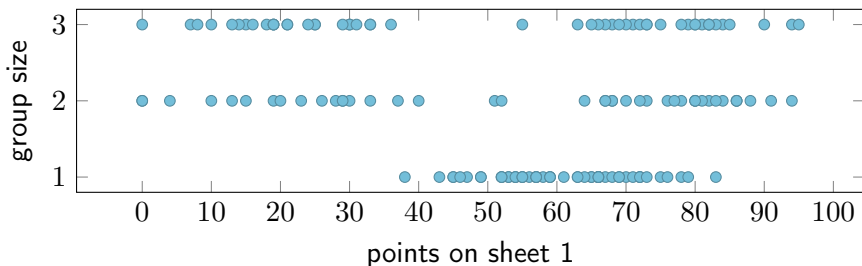
- Say we want to personalize a course, but have limited teaching capacity

## Motivation: Education Example

- ▶ Say we want to personalize a course, but have limited teaching capacity
- ⇒ Find **clusters** of students who have similar support needs.

## Motivation: Education Example

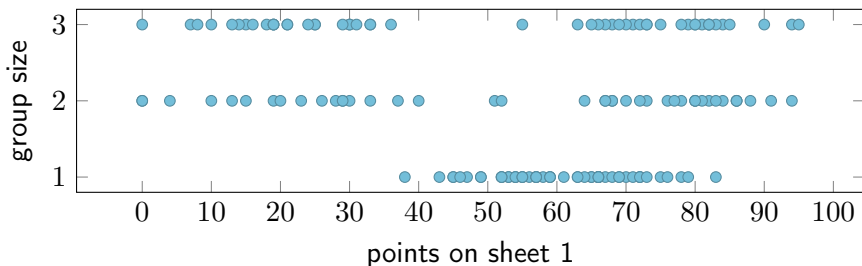
- Say we want to personalize a course, but have limited teaching capacity
- ⇒ Find **clusters** of students who have similar support needs.





## Motivation: Education Example

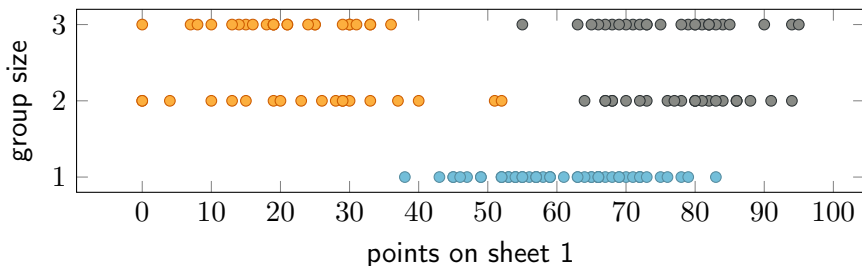
- Say we want to personalize a course, but have limited teaching capacity
- ⇒ Find **clusters** of students who have similar support needs.



- Which patterns do you notice?

## Motivation: Education Example

- Say we want to personalize a course, but have limited teaching capacity
- ⇒ Find **clusters** of students who have similar support needs.



- Which patterns do you notice?

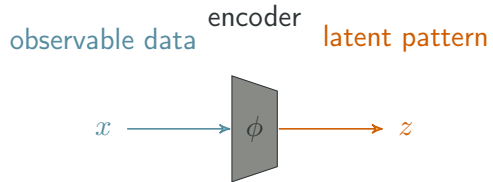
## K-Means algorithm

# Autoencoding Setup

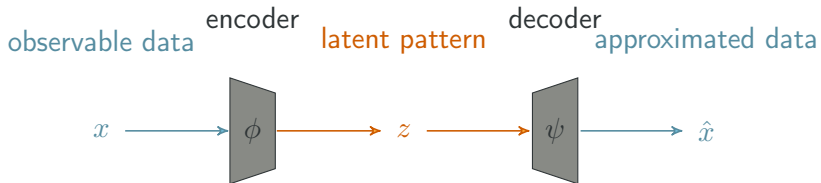
observable data

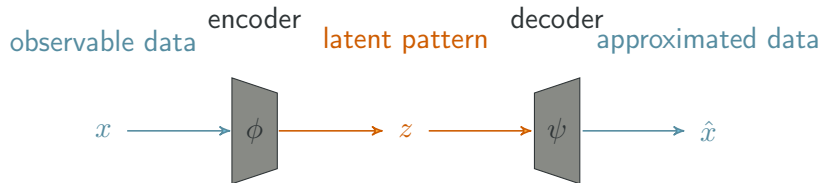
$x$

# Autoencoding Setup

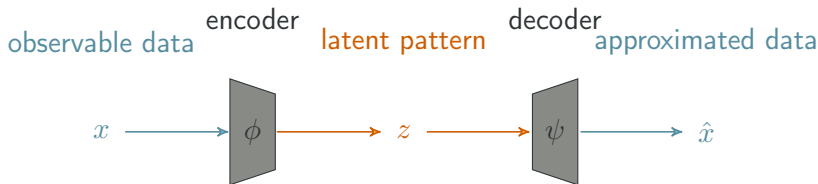


# Autoencoding Setup





- Clustering means that  $z \in \{1, \dots, K\}$



► Clustering means that  $z \in \{1, \dots, K\}$

⇒ decoder:  $\psi(z) = \mu_z \in \mathbb{R}^m$  (discrete catalogue of “prototypes”)



## Optimal encoder

- For each  $x$ , directly optimize reconstruction error

$$\phi(x) = \arg \min_k \|\psi(k) - x\|$$

- For each  $x$ , directly optimize reconstruction error

$$\phi(x) = \arg \min_k \|\psi(k) - x\| = \arg \min_k \|\mu_k - x\|$$

## Optimal encoder

- For each  $x$ , directly optimize reconstruction error

$$\phi(x) = \arg \min_k \|\psi(k) - x\| = \arg \min_k \|\mu_k - x\|$$

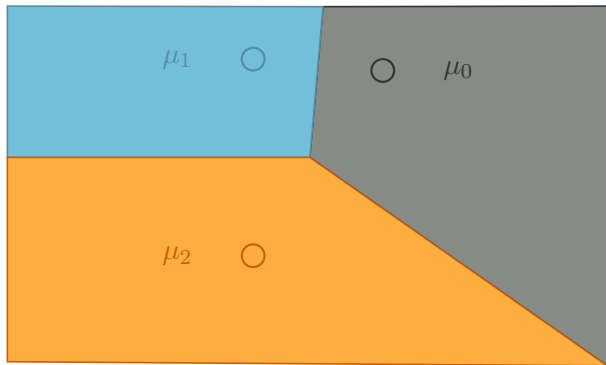
$\mu_1$     $\mu_0$

$\mu_2$  

## Optimal encoder

- For each  $x$ , directly optimize reconstruction error

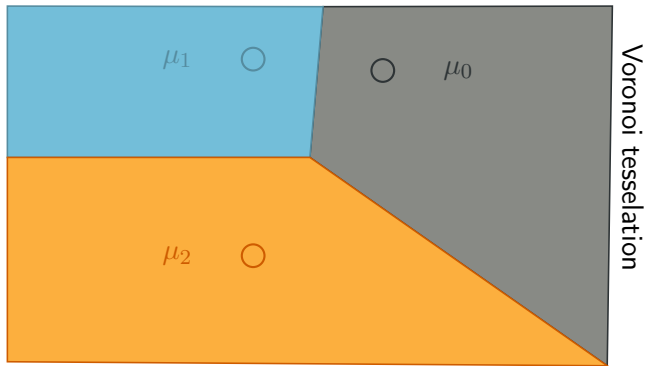
$$\phi(x) = \arg \min_k \|\psi(k) - x\| = \arg \min_k \|\mu_k - x\|$$



## Optimal encoder

- For each  $x$ , directly optimize reconstruction error

$$\phi(x) = \arg \min_k \|\psi(k) - x\| = \arg \min_k \|\mu_k - x\|$$



# Optimize prototype positions

- ▶ Let  $x_1, \dots, x_N \in \mathbb{R}^m$  be example data

- ▶ Let  $x_1, \dots, x_N \in \mathbb{R}^m$  be example data
- ▶ Try to find prototypes  $\mu_1, \dots, \mu_K \in \mathbb{R}^m$  that minimize reconstruction error

$$\min_{\mu_1, \dots, \mu_K} \sum_{i=1}^N \|\mu_{z_i} - x_i\|^2$$

such that  $z_i = \phi(x_i) = \arg \min_k \|\mu_k - x_i\|$

- ▶ Let  $x_1, \dots, x_N \in \mathbb{R}^m$  be example data
- ▶ Try to find prototypes  $\mu_1, \dots, \mu_K \in \mathbb{R}^m$  that minimize reconstruction error

$$\min_{\mu_1, \dots, \mu_K} \sum_{i=1}^N \|\mu_{z_i} - x_i\|^2$$

such that  $z_i = \phi(x_i) = \arg \min_k \|\mu_k - x_i\|$

- ▶ Unfortunately, this is NP-hard :(



## Optimize prototype positions (continued)

- ▶ Idea: Optimization becomes really simple, if  $z_i$  is fixed :)

## Optimize prototype positions (continued)

- ▶ Idea: Optimization becomes really simple, if  $z_i$  is fixed :)
- ▶ Let  $\mathcal{P}_k : \{i | z_i = k\}$ ; the “receptive field” of prototype  $k$

## Optimize prototype positions (continued)

- ▶ Idea: Optimization becomes really simple, if  $z_i$  is fixed :)
- ▶ Let  $\mathcal{P}_k : \{i | z_i = k\}$ ; the “receptive field” of prototype  $k$
- ▶ Calculate gradient:

$$\nabla_{\mu_k} \sum_{i=1}^N \|\mu_{z_i} - x_i\|^2$$

- ▶ Idea: Optimization becomes really simple, if  $z_i$  is fixed :)
- ▶ Let  $\mathcal{P}_k : \{i | z_i = k\}$ ; the “receptive field” of prototype  $k$
- ▶ Calculate gradient:

$$\nabla_{\mu_k} \sum_{i=1}^N \|\mu_{z_i} - x_i\|^2 = \nabla_{\mu_k} \sum_{i \in \mathcal{P}_k} \|\mu_k - x_i\|^2$$

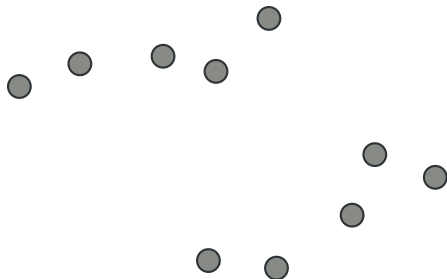
- ▶ Idea: Optimization becomes really simple, if  $z_i$  is fixed :)
- ▶ Let  $\mathcal{P}_k : \{i | z_i = k\}$ ; the “receptive field” of prototype  $k$
- ▶ Calculate gradient:

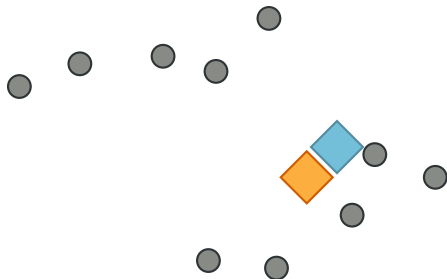
$$\begin{aligned}\nabla_{\mu_k} \sum_{i=1}^N \|\mu_{z_i} - x_i\|^2 &= \nabla_{\mu_k} \sum_{i \in \mathcal{P}_k} \|\mu_k - x_i\|^2 \\ &= \sum_{i \in \mathcal{P}_k} 2(\mu_k - x_i)\end{aligned}$$

- ▶ Idea: Optimization becomes really simple, if  $z_i$  is fixed :)
- ▶ Let  $\mathcal{P}_k : \{i | z_i = k\}$ ; the “receptive field” of prototype  $k$
- ▶ Calculate gradient:

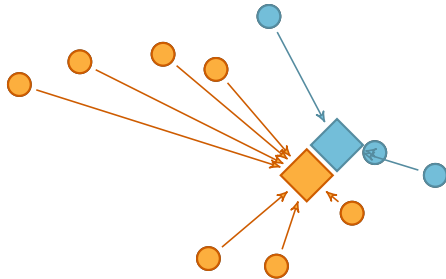
$$\begin{aligned}\nabla_{\mu_k} \sum_{i=1}^N \|\mu_{z_i} - x_i\|^2 &= \nabla_{\mu_k} \sum_{i \in \mathcal{P}_k} \|\mu_k - x_i\|^2 \\ &= \sum_{i \in \mathcal{P}_k} 2(\mu_k - x_i)\end{aligned}$$

- ▶ Setting gradient to zero yields  $\mu_k = \frac{1}{|\mathcal{P}_k|} \sum_{i \in \mathcal{P}_k} x_i$ .

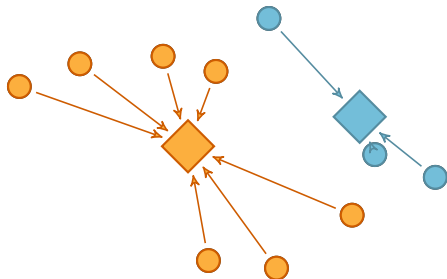




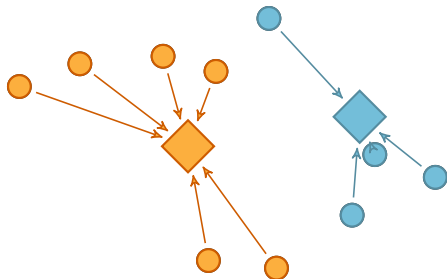




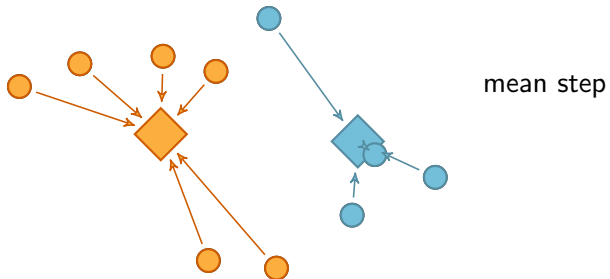
assignment step

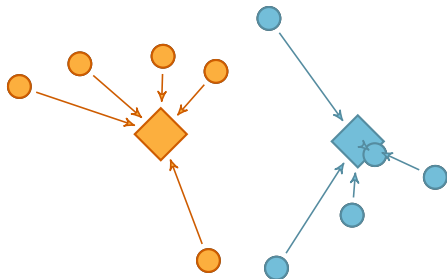


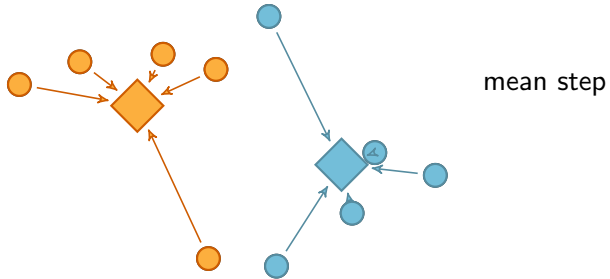
mean step

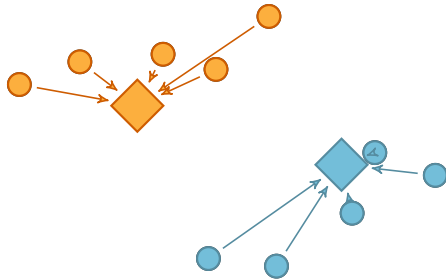


assignment step

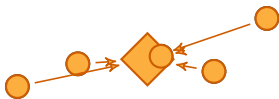




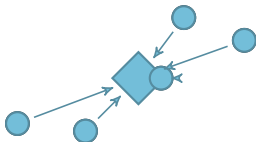




assignment step



mean step





**function** KMeans(data matrix  $\mathbf{X}$  with  $N$  rows and  $m$  columns, desired number of clusters  $K$ )

**end function**

**function** KMeans(data matrix  $\mathbf{X}$  with  $N$  rows and  $m$  columns, desired number of clusters  $K$ )

Randomly initialize  $\mu_1, \dots, \mu_K \in \mathbb{R}^m$  (in the convex hull of the data).

**end function**

**function** KMeans(data matrix  $\mathbf{X}$  with  $N$  rows and  $m$  columns, desired number of clusters  $K$ )

Randomly initialize  $\mu_1, \dots, \mu_K \in \mathbb{R}^m$  (in the convex hull of the data).

**while** at least one  $z_i$  changed in the last iteration **do**

**end while**

**end function**

**function** KMeans(data matrix  $\mathbf{X}$  with  $N$  rows and  $m$  columns, desired number of clusters  $K$ )

Randomly initialize  $\mu_1, \dots, \mu_K \in \mathbb{R}^m$  (in the convex hull of the data).

**while** at least one  $z_i$  changed in the last iteration **do**

    Compute  $z_i \leftarrow \arg \min_k \|\mu_k - x_i\|$  for all  $i \in \{1, \dots, N\}$ .

**end while**

**end function**

**function** KMeans(data matrix  $\mathbf{X}$  with  $N$  rows and  $m$  columns, desired number of clusters  $K$ )

Randomly initialize  $\mu_1, \dots, \mu_K \in \mathbb{R}^m$  (in the convex hull of the data).

**while** at least one  $z_i$  changed in the last iteration **do**

    Compute  $z_i \leftarrow \arg \min_k \|\mu_k - x_i\|$  for all  $i \in \{1, \dots, N\}$ .

    Compute  $\mathcal{P}_k \leftarrow \{i | z_i = k\}$  for all  $k \in \{1, \dots, K\}$ .

**end while**

**end function**

**function** KMeans(data matrix  $\mathbf{X}$  with  $N$  rows and  $m$  columns, desired number of clusters  $K$ )

Randomly initialize  $\mu_1, \dots, \mu_K \in \mathbb{R}^m$  (in the convex hull of the data).

**while** at least one  $z_i$  changed in the last iteration **do**

    Compute  $z_i \leftarrow \arg \min_k \|\mu_k - x_i\|$  for all  $i \in \{1, \dots, N\}$ .

    Compute  $\mathcal{P}_k \leftarrow \{i | z_i = k\}$  for all  $k \in \{1, \dots, K\}$ .

    Compute  $\mu_k \leftarrow \frac{1}{|\mathcal{P}_k|} \sum_{i \in \mathcal{P}_k} x_i$  for all  $k \in \{1, \dots, K\}$ .

**end while**

**end function**

**function** KMeans(data matrix  $\mathbf{X}$  with  $N$  rows and  $m$  columns, desired number of clusters  $K$ )

Randomly initialize  $\mu_1, \dots, \mu_K \in \mathbb{R}^m$  (in the convex hull of the data).

**while** at least one  $z_i$  changed in the last iteration **do**

    Compute  $z_i \leftarrow \arg \min_k \|\mu_k - x_i\|$  for all  $i \in \{1, \dots, N\}$ .

    Compute  $\mathcal{P}_k \leftarrow \{i | z_i = k\}$  for all  $k \in \{1, \dots, K\}$ .

    Compute  $\mu_k \leftarrow \frac{1}{|\mathcal{P}_k|} \sum_{i \in \mathcal{P}_k} x_i$  for all  $k \in \{1, \dots, K\}$ .

**end while**

**return**  $\mu_1, \dots, \mu_K$ .

**end function**

## Proof of stopping (Sketch)

- ▶ Let  $z_1, \dots, z_t$  be the series of assignment vectors for subsequent iterations of  $K$ -means



## Proof of stopping (Sketch)

- ▶ Let  $z_1, \dots, z_t$  be the series of assignment vectors for subsequent iterations of  $K$ -means
- ▶ There are only finitely many possible assignments – so for sufficiently high  $t$ , one must repeat

## Proof of stopping (Sketch)

- ▶ Let  $z_1, \dots, z_t$  be the series of assignment vectors for subsequent iterations of  $K$ -means
- ▶ There are only finitely many possible assignments – so for sufficiently high  $t$ , one must repeat
- ▶ An immediate repeat stops the algorithm

## Proof of stopping (Sketch)

- ▶ Let  $z_1, \dots, z_t$  be the series of assignment vectors for subsequent iterations of  $K$ -means
- ▶ There are only finitely many possible assignments – so for sufficiently high  $t$ , one must repeat
- ▶ An immediate repeat stops the algorithm
- ▶ Otherwise, note that each assignment vector  $z$  corresponds to **one** reconstruction error (after re-computing the means)

## Proof of stopping (Sketch)

- ▶ Let  $z_1, \dots, z_t$  be the series of assignment vectors for subsequent iterations of  $K$ -means
- ▶ There are only finitely many possible assignments – so for sufficiently high  $t$ , one must repeat
- ▶ An immediate repeat stops the algorithm
- ▶ Otherwise, note that each assignment vector  $z$  corresponds to **one** reconstruction error (after re-computing the means)
- ▶  $K$ -means only changes the assignments to **reduce** the reconstruction error

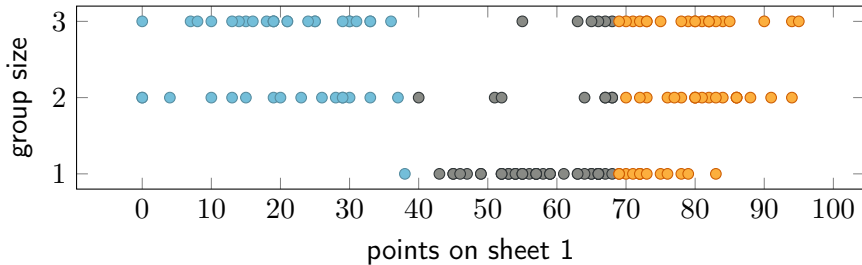
# Proof of stopping (Sketch)

- ▶ Let  $z_1, \dots, z_t$  be the series of assignment vectors for subsequent iterations of  $K$ -means
- ▶ There are only finitely many possible assignments – so for sufficiently high  $t$ , one must repeat
- ▶ An immediate repeat stops the algorithm
- ▶ Otherwise, note that each assignment vector  $z$  corresponds to **one** reconstruction error (after re-computing the means)
- ▶  $K$ -means only changes the assignments to **reduce** the reconstruction error
- ▶ Hence, repeating an assignment vector is impossible (this would increase the loss)

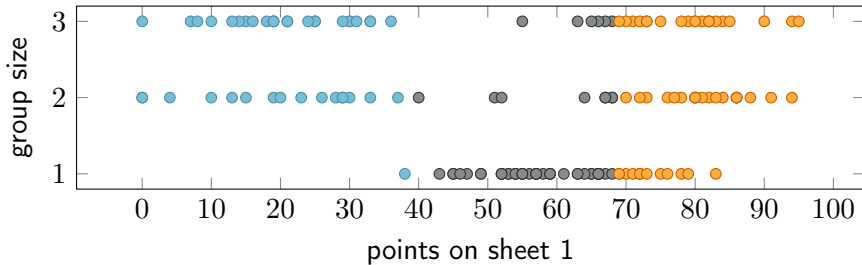
## Proof of stopping (Sketch)

- ▶ Let  $z_1, \dots, z_t$  be the series of assignment vectors for subsequent iterations of  $K$ -means
- ▶ There are only finitely many possible assignments – so for sufficiently high  $t$ , one must repeat
- ▶ An immediate repeat stops the algorithm
- ▶ Otherwise, note that each assignment vector  $z$  corresponds to **one** reconstruction error (after re-computing the means)
- ▶  $K$ -means only changes the assignments to **reduce** the reconstruction error
- ▶ Hence, repeating an assignment vector is impossible (this would increase the loss)
- ▶ **note!** in practice,  $K$ -means mostly stops already after ca. 10-30 iterations

## Example: non-normalized data



## Example: non-normalized data



- What is the problem here?



## Example: normalized data

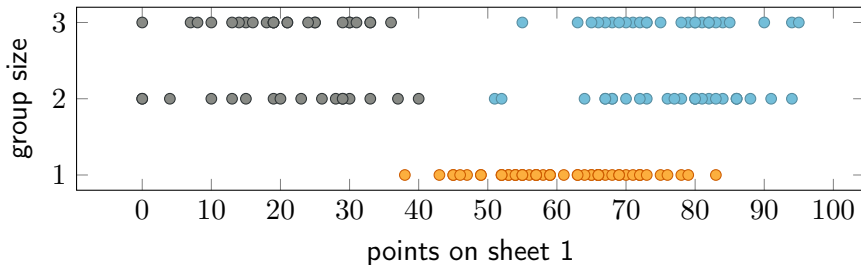
- ▶ Trick:  $z$ -normalize data before using  $K$ -means to avoid scaling issues:

$$x_{i,j} \leftarrow \frac{x_{i,j} - \mu_j}{\sigma_j}$$

## Example: normalized data

- Trick:  $z$ -normalize data before using  $K$ -means to avoid scaling issues:

$$x_{i,j} \leftarrow \frac{x_{i,j} - \mu_j}{\sigma_j}$$



## How to choose $K$ ?

- ▶ Idea: Compute “quality indices” for clusterings and take  $K$  that yields the best quality index

## How to choose $K$ ?

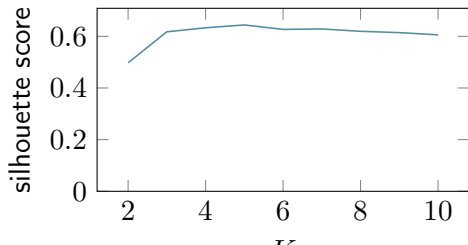
- ▶ Idea: Compute “quality indices” for clusterings and take  $K$  that yields the best quality index
- ▶ Silhouette Score: Let  $d_i^+$  be average distance to other samples in the same cluster as  $i$  and let  $d_i^-$  be the distance of  $i$  to the closest, other cluster

- ▶ Idea: Compute “quality indices” for clusterings and take  $K$  that yields the best quality index
- ▶ Silhouette Score: Let  $d_i^+$  be average distance to other samples in the same cluster as  $i$  and let  $d_i^-$  be the distance of  $i$  to the closest, other cluster  $\Rightarrow$   
$$\frac{1}{N} \sum_{i=1}^N \frac{d_i^- - d_i^+}{\max\{d_i^-, d_i^+\}}$$

- ▶ Idea: Compute “quality indices” for clusterings and take  $K$  that yields the best quality index
- ▶ Silhouette Score: Let  $d_i^+$  be average distance to other samples in the same cluster as  $i$  and let  $d_i^-$  be the distance of  $i$  to the closest, other cluster  $\Rightarrow$ 
$$\frac{1}{N} \sum_{i=1}^N \frac{d_i^- - d_i^+}{\max\{d_i^-, d_i^+\}}$$
- ▶ Bayesian information criterion (BIC): reconstruction loss plus number of model parameters:  $K \cdot m$

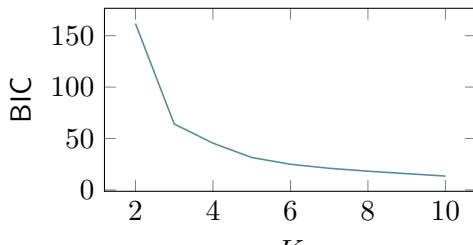
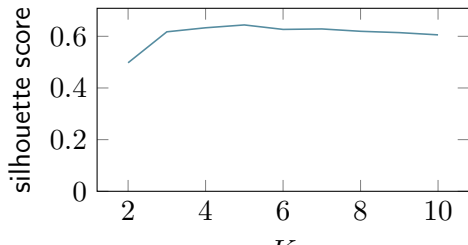
# How to choose $K$ ?

- ▶ Idea: Compute “quality indices” for clusterings and take  $K$  that yields the best quality index
- ▶ Silhouette Score: Let  $d_i^+$  be average distance to other samples in the same cluster as  $i$  and let  $d_i^-$  be the distance of  $i$  to the closest, other cluster  $\Rightarrow$   
$$\frac{1}{N} \sum_{i=1}^N \frac{d_i^- - d_i^+}{\max\{d_i^-, d_i^+\}}$$
- ▶ Bayesian information criterion (BIC): reconstruction loss plus number of model parameters:  $K \cdot m$



# How to choose $K$ ?

- ▶ Idea: Compute “quality indices” for clusterings and take  $K$  that yields the best quality index
- ▶ Silhouette Score: Let  $d_i^+$  be average distance to other samples in the same cluster as  $i$  and let  $d_i^-$  be the distance of  $i$  to the closest, other cluster  $\Rightarrow$   
$$\frac{1}{N} \sum_{i=1}^N \frac{d_i^- - d_i^+}{\max\{d_i^-, d_i^+\}}$$
- ▶ Bayesian information criterion (BIC): reconstruction loss plus number of model parameters:  $K \cdot m$



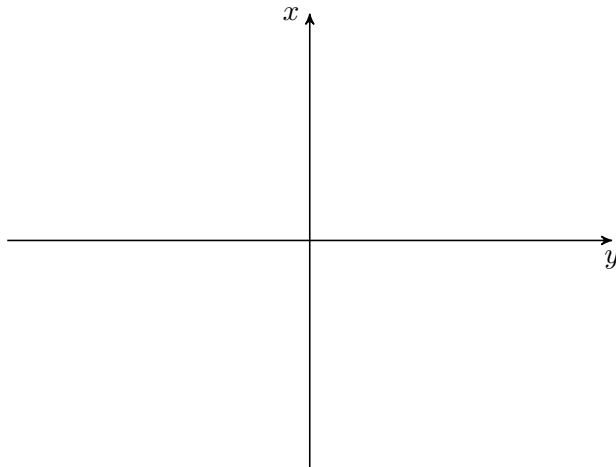


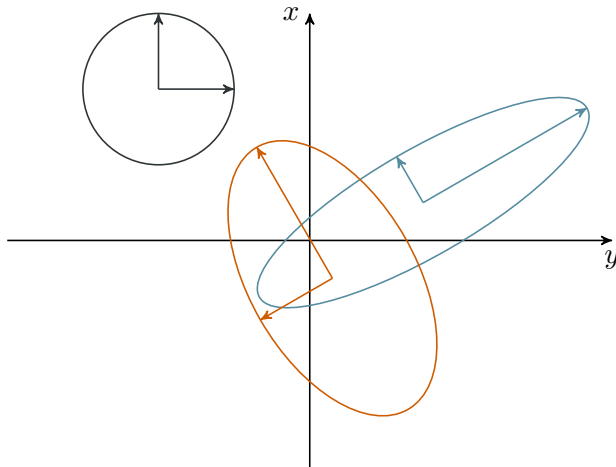
## Gaussian Mixture Models

- ▶ Assume data is generated by a mixture of  $K$  Gaussians as follows:

- ▶ Assume data is generated by a mixture of  $K$  Gaussians as follows:
- ▶ Step 1: sample latent cluster index  $k$  via marginal probability  $p_Z(k)$

- ▶ Assume data is generated by a mixture of  $K$  Gaussians as follows:
- ▶ Step 1: sample latent cluster index  $k$  via marginal probability  $p_Z(k)$
- ▶ Step 2: Sample observable data point  $x$  from Gaussian  $p_{X|Z}(x|k)$  with mean  $\mu_k$  and covariance matrix  $\Sigma_k$





$$p_X(x) = \sum_{k=1}^K p_{X|Z}(x|k) \cdot p_Z(k)$$

$$p_X(x) = \sum_{k=1}^K p_{X|Z}(x|k) \cdot p_Z(k)$$
$$p_{Z|X}(k|x) = \frac{p_{X|Z}(x|k) \cdot p_Z(k)}{\sum_{z=1}^K p_{X|Z}(x|z) \cdot p_Z(z)}$$



- Negative log likelihood:

$$\ell = \sum_{i=1}^N -\log \left[ \sum_{k=1}^K p_{X|Z}(x_i|k) \cdot p_Z(k) \right]$$

- Negative log likelihood:

$$\ell = \sum_{i=1}^N -\log \left[ \sum_{k=1}^K p_{X|Z}(x_i|k) \cdot p_Z(k) \right]$$

⇒ Nasty to optimize (log of sum, non-convex) :(

- Negative log likelihood:

$$\ell = \sum_{i=1}^N -\log \left[ \sum_{k=1}^K p_{X|Z}(x_i|k) \cdot p_Z(k) \right]$$

⇒ Nasty to optimize (log of sum, non-convex) :(

- Trick: Compute  $\gamma_{k,i} = p_{Z|X}(k|x_i)$ , then optimize the **expected** negative log likelihood instead:

$$Q = \sum_{i=1}^N \sum_{k=1}^K -\gamma_{k,i} \cdot \log \left[ p_{X|Z}(x_i|k) \cdot p_Z(k) \right]$$

- Negative log likelihood:

$$\ell = \sum_{i=1}^N -\log \left[ \sum_{k=1}^K p_{X|Z}(x_i|k) \cdot p_Z(k) \right]$$

⇒ Nasty to optimize (log of sum, non-convex) :(

- Trick: Compute  $\gamma_{k,i} = p_{Z|X}(k|x_i)$ , then optimize the **expected** negative log likelihood instead:

$$\begin{aligned} Q &= \sum_{i=1}^N \sum_{k=1}^K -\gamma_{k,i} \cdot \log \left[ p_{X|Z}(x_i|k) \cdot p_Z(k) \right] \\ &= \sum_{i=1}^N \sum_{k=1}^K \gamma_{k,i} \cdot \left( \frac{1}{2} \log[2\pi \det(\Sigma_k)] + \frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) - \log[p_Z(k)] \right) \end{aligned}$$

- Negative log likelihood:

$$\ell = \sum_{i=1}^N -\log \left[ \sum_{k=1}^K p_{X|Z}(x_i|k) \cdot p_Z(k) \right]$$

⇒ Nasty to optimize (log of sum, non-convex) :(

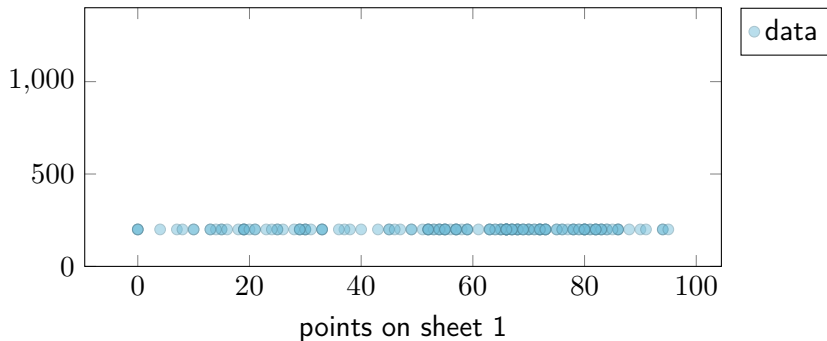
- Trick: Compute  $\gamma_{k,i} = p_{Z|X}(k|x_i)$ , then optimize the **expected** negative log likelihood instead:

$$\begin{aligned} Q &= \sum_{i=1}^N \sum_{k=1}^K -\gamma_{k,i} \cdot \log \left[ p_{X|Z}(x_i|k) \cdot p_Z(k) \right] \\ &= \sum_{i=1}^N \sum_{k=1}^K \gamma_{k,i} \cdot \left( \frac{1}{2} \log[2\pi \det(\Sigma_k)] + \frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) - \log[p_Z(k)] \right) \end{aligned}$$

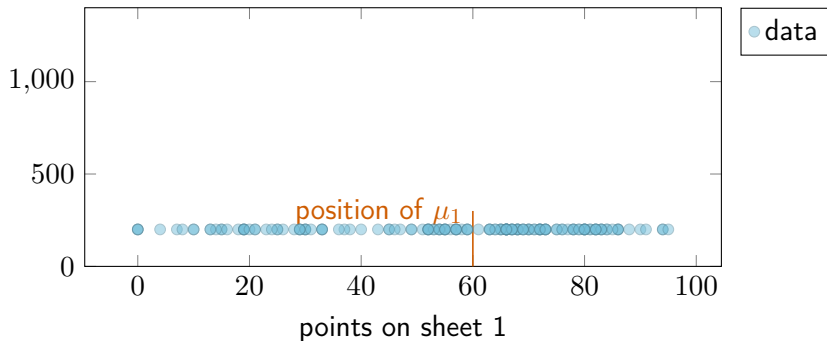
⇒ Convex and closed-form solution (see exercises) :)

- Scenario: only consider one dimension,  $K = 2$ ,  $\sigma_1 = \sigma_2 = 10$ ,  $\mu_1 = 60$ , only  $\mu_2$  varies

- Scenario: only consider one dimension,  $K = 2$ ,  $\sigma_1 = \sigma_2 = 10$ ,  $\mu_1 = 60$ , only  $\mu_2$  varies

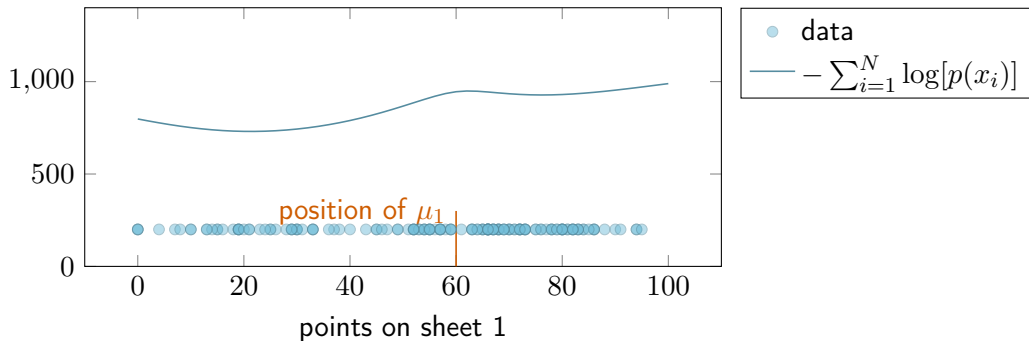


- Scenario: only consider one dimension,  $K = 2$ ,  $\sigma_1 = \sigma_2 = 10$ ,  $\mu_1 = 60$ , only  $\mu_2$  varies

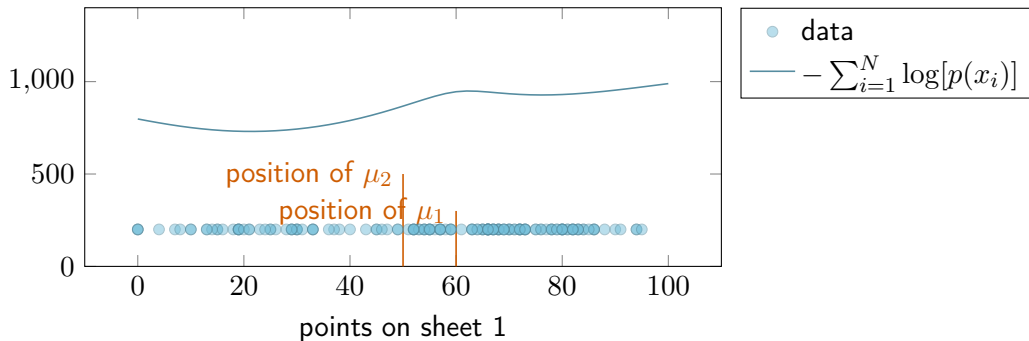




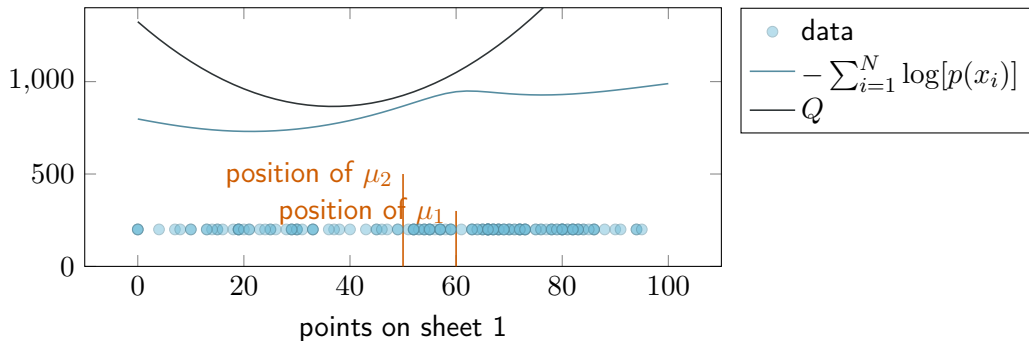
- Scenario: only consider one dimension,  $K = 2$ ,  $\sigma_1 = \sigma_2 = 10$ ,  $\mu_1 = 60$ , only  $\mu_2$  varies



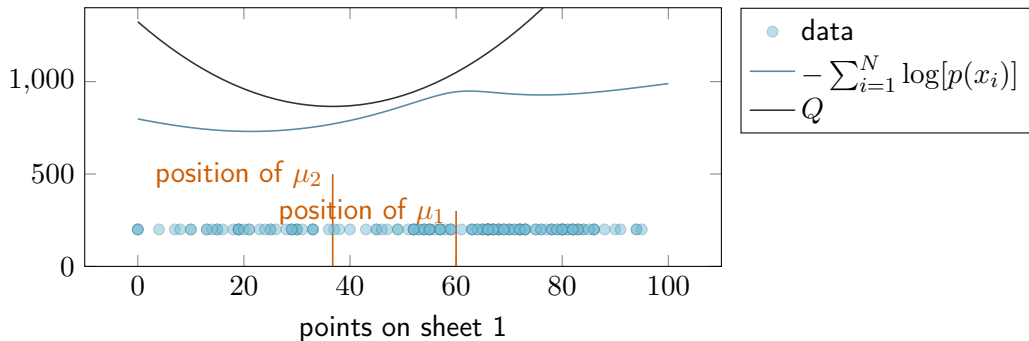
- Scenario: only consider one dimension,  $K = 2$ ,  $\sigma_1 = \sigma_2 = 10$ ,  $\mu_1 = 60$ , only  $\mu_2$  varies



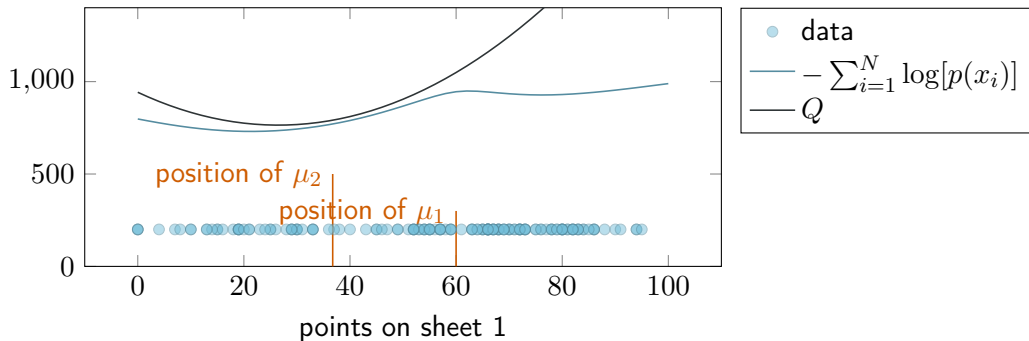
- Scenario: only consider one dimension,  $K = 2$ ,  $\sigma_1 = \sigma_2 = 10$ ,  $\mu_1 = 60$ , only  $\mu_2$  varies



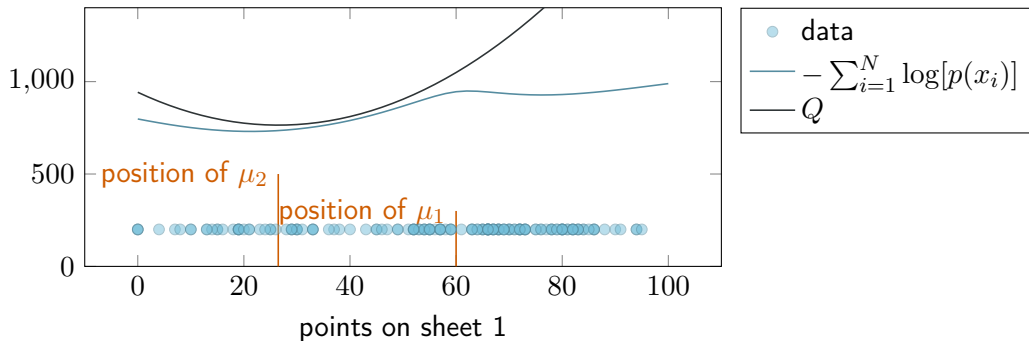
- Scenario: only consider one dimension,  $K = 2$ ,  $\sigma_1 = \sigma_2 = 10$ ,  $\mu_1 = 60$ , only  $\mu_2$  varies



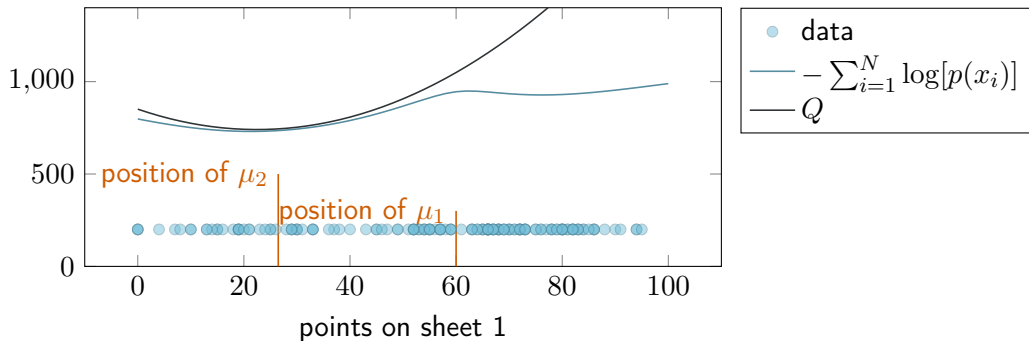
- Scenario: only consider one dimension,  $K = 2$ ,  $\sigma_1 = \sigma_2 = 10$ ,  $\mu_1 = 60$ , only  $\mu_2$  varies



- Scenario: only consider one dimension,  $K = 2$ ,  $\sigma_1 = \sigma_2 = 10$ ,  $\mu_1 = 60$ , only  $\mu_2$  varies



- Scenario: only consider one dimension,  $K = 2$ ,  $\sigma_1 = \sigma_2 = 10$ ,  $\mu_1 = 60$ , only  $\mu_2$  varies



Let  $\gamma_{i,k}$  be arbitrary real number s.t.  $\sum_{k=1}^K \gamma_{i,k} = 1$  and  $\gamma_{i,k} \geq 0$  for all  $i$  and all  $k$ .  
For every  $x_i$ , we have:



Let  $\gamma_{i,k}$  be arbitrary real number s.t.  $\sum_{k=1}^K \gamma_{i,k} = 1$  and  $\gamma_{i,k} \geq 0$  for all  $i$  and all  $k$ .  
For every  $x_i$ , we have:

$$\log[p_X(x_i)] = \sum_{k=1}^K \gamma_{i,k} \cdot \log[p_X(x_i)]$$

Let  $\gamma_{i,k}$  be arbitrary real number s.t.  $\sum_{k=1}^K \gamma_{i,k} = 1$  and  $\gamma_{i,k} \geq 0$  for all  $i$  and all  $k$ .  
For every  $x_i$ , we have:

$$\begin{aligned}\log[p_X(x_i)] &= \sum_{k=1}^K \gamma_{i,k} \cdot \log[p_X(x_i)] \\ &= \sum_{k=1}^K \gamma_{i,k} \cdot \log \left[ \frac{p_{X,Z}(x_i, k)}{p_{Z|X}(k|x_i)} \right]\end{aligned}$$

Let  $\gamma_{i,k}$  be arbitrary real number s.t.  $\sum_{k=1}^K \gamma_{i,k} = 1$  and  $\gamma_{i,k} \geq 0$  for all  $i$  and all  $k$ .  
For every  $x_i$ , we have:

$$\begin{aligned}\log[p_X(x_i)] &= \sum_{k=1}^K \gamma_{i,k} \cdot \log[p_X(x_i)] \\ &= \sum_{k=1}^K \gamma_{i,k} \cdot \log \left[ \frac{p_{X,Z}(x_i, k)}{p_{Z|X}(k|x_i)} \right] \\ &= \sum_{k=1}^K \gamma_{i,k} \cdot \log[p_{X,Z}(x_i, k)] - \sum_{k=1}^K \gamma_{i,k} \cdot \log[p_{Z|X}(k|x_i)]\end{aligned}$$

Let  $\gamma_{i,k}$  be arbitrary real number s.t.  $\sum_{k=1}^K \gamma_{i,k} = 1$  and  $\gamma_{i,k} \geq 0$  for all  $i$  and all  $k$ .

For every  $x_i$ , we have:

$$\begin{aligned}\log[p_X(x_i)] &= \sum_{k=1}^K \gamma_{i,k} \cdot \log[p_X(x_i)] \\ &= \sum_{k=1}^K \gamma_{i,k} \cdot \log \left[ \frac{p_{X,Z}(x_i, k)}{p_{Z|X}(k|x_i)} \right] \\ &= \sum_{k=1}^K \gamma_{i,k} \cdot \log[p_{X,Z}(x_i, k)] - \sum_{k=1}^K \gamma_{i,k} \cdot \log[p_{Z|X}(k|x_i)] \\ &= \sum_{k=1}^K \gamma_{i,k} \cdot \log[p_{X,Z}(x_i, k)] - \sum_{k=1}^K \gamma_{i,k} \cdot \left( \log[p_{Z|X}(k|x_i)] - \log[\gamma_{i,k}] + \log[\gamma_{i,k}] \right)\end{aligned}$$

$$\log[p_X(x_i)] = \sum_{k=1}^K \gamma_{i,k} \cdot \log[p_{X,Z}(x_i, k)] - \sum_{k=1}^K \gamma_{i,k} \cdot \log \left[ \frac{p_{Z|X}(k|x_i)}{\gamma_{i,k}} \right] - \sum_{k=1}^K \gamma_{i,k} \cdot \log[\gamma_{i,k}]$$

$$\log[p_X(x_i)] = \sum_{k=1}^K \gamma_{i,k} \cdot \log[p_{X,Z}(x_i, k)] - \sum_{k=1}^K \gamma_{i,k} \cdot \log \left[ \frac{p_{Z|X}(k|x_i)}{\gamma_{i,k}} \right] - \sum_{k=1}^K \gamma_{i,k} \cdot \log[\gamma_{i,k}]$$

► We summarize:

$$Q = - \sum_{i=1}^N \sum_{k=1}^K \gamma_{i,k} \cdot \log[p_{X,Z}(x_i, k)] \quad (\text{exp. neg. log likelihood})$$

$$\mathcal{D}_{\text{KL}}(\gamma_i || p_{Z_i|X_i}) = - \sum_{k=1}^K \gamma_{i,k} \cdot \log \left[ \frac{p_{Z|X}(k|x_i)}{\gamma_{i,k}} \right] \quad (\text{Kullback-Leibler-Divergence})$$

$$\mathcal{H}(\gamma_i) = - \sum_{k=1}^K \gamma_{i,k} \cdot \log[\gamma_{i,k}] \quad (\text{entropy})$$

$$-\sum_{i=1}^N \log[p_X(x_i)] = Q - \sum_{i=1}^N \mathcal{D}_{\text{KL}}(\gamma_i || p_{Z_i|X_i}) + \mathcal{H}(\gamma_i)$$

$$-\sum_{i=1}^N \log[p_X(x_i)] = Q - \sum_{i=1}^N \mathcal{D}_{\text{KL}}(\gamma_i || p_{Z_i|X_i}) + \mathcal{H}(\gamma_i)$$

Interpretation:



$$-\sum_{i=1}^N \log[p_X(x_i)] = Q - \sum_{i=1}^N \mathcal{D}_{\text{KL}}(\gamma_i || p_{Z_i|X_i}) + \mathcal{H}(\gamma_i)$$

Interpretation:

- The true neg. log likelihood  $\mathcal{L}$  is equal to  $Q$  (our much nicer loss) minus KL divergences minus entropies

$$-\sum_{i=1}^N \log[p_X(x_i)] = Q - \sum_{i=1}^N \mathcal{D}_{\text{KL}}(\gamma_i || p_{Z_i|X_i}) + \mathcal{H}(\gamma_i)$$

Interpretation:

- ▶ The true neg. log likelihood  $\mathcal{L}$  is equal to  $Q$  (our much nicer loss) minus KL divergences minus entropies
- ▶ Both KL divergence and entropy are non-negative!

$$-\sum_{i=1}^N \log[p_X(x_i)] = Q - \sum_{i=1}^N \mathcal{D}_{\text{KL}}(\gamma_i || p_{Z_i|X_i}) + \mathcal{H}(\gamma_i)$$

Interpretation:

- ▶ The true neg. log likelihood  $\mathcal{L}$  is equal to  $Q$  (our much nicer loss) minus KL divergences minus entropies
  - ▶ Both KL divergence and entropy are non-negative!
- $\Rightarrow \mathcal{L}$  is upper-bounded by  $Q \Rightarrow$  optimizing  $Q$  improves our upper bound on  $\mathcal{L}$

$$-\sum_{i=1}^N \log[p_X(x_i)] = Q - \sum_{i=1}^N \mathcal{D}_{\text{KL}}(\gamma_i || p_{Z_i|X_i}) + \mathcal{H}(\gamma_i)$$

Interpretation:

- ▶ The true neg. log likelihood  $\mathcal{L}$  is equal to  $Q$  (our much nicer loss) minus KL divergences minus entropies
  - ▶ Both KL divergence and entropy are non-negative!
- $\Rightarrow \mathcal{L}$  is upper-bounded by  $Q \Rightarrow$  optimizing  $Q$  improves our upper bound on  $\mathcal{L}$
- ▶ Even better: KL divergence gets zero if  $\gamma_{i,k} = p_{Z_i|X_i}(k|x_i)$

$$-\sum_{i=1}^N \log[p_X(x_i)] = Q - \sum_{i=1}^N \mathcal{D}_{\text{KL}}(\gamma_i || p_{Z_i|X_i}) + \mathcal{H}(\gamma_i)$$

Interpretation:

- ▶ The true neg. log likelihood  $\mathcal{L}$  is equal to  $Q$  (our much nicer loss) minus KL divergences minus entropies
- ▶ Both KL divergence and entropy are non-negative!
- $\Rightarrow \mathcal{L}$  is upper-bounded by  $Q \Rightarrow$  optimizing  $Q$  improves our upper bound on  $\mathcal{L}$
- ▶ Even better: KL divergence gets zero if  $\gamma_{i,k} = p_{Z_i|X_i}(k|x_i)$
- $\Rightarrow$  Every time we re-compute  $\gamma_{i,k} = p_{Z_i|X_i}(k|x_i)$ , we improve how well  $Q$  approximates  $\mathcal{L}$

$$-\sum_{i=1}^N \log[p_X(x_i)] = Q - \sum_{i=1}^N \mathcal{D}_{\text{KL}}(\gamma_i || p_{Z_i|X_i}) + \mathcal{H}(\gamma_i)$$

Interpretation:

- ▶ The true neg. log likelihood  $\mathcal{L}$  is equal to  $Q$  (our much nicer loss) minus KL divergences minus entropies
- ▶ Both KL divergence and entropy are non-negative!
- $\Rightarrow \mathcal{L}$  is upper-bounded by  $Q \Rightarrow$  optimizing  $Q$  improves our upper bound on  $\mathcal{L}$
- ▶ Even better: KL divergence gets zero if  $\gamma_{i,k} = p_{Z_i|X_i}(k|x_i)$
- $\Rightarrow$  Every time we re-compute  $\gamma_{i,k} = p_{Z_i|X_i}(k|x_i)$ , we improve how well  $Q$  approximates  $\mathcal{L}$
- ▶ Entropy gets zero if every point belongs exactly to one cluster

$$-\sum_{i=1}^N \log[p_X(x_i)] = Q - \sum_{i=1}^N \mathcal{D}_{\text{KL}}(\gamma_i || p_{Z_i|X_i}) + \mathcal{H}(\gamma_i)$$

Interpretation:

- ▶ The true neg. log likelihood  $\mathcal{L}$  is equal to  $Q$  (our much nicer loss) minus KL divergences minus entropies
- ▶ Both KL divergence and entropy are non-negative!
- ⇒  $\mathcal{L}$  is upper-bounded by  $Q$  ⇒ optimizing  $Q$  improves our upper bound on  $\mathcal{L}$
- ▶ Even better: KL divergence gets zero if  $\gamma_{i,k} = p_{Z_i|X_i}(k|x_i)$
- ⇒ Every time we re-compute  $\gamma_{i,k} = p_{Z_i|X_i}(k|x_i)$ , we improve how well  $Q$  approximates  $\mathcal{L}$
- ▶ Entropy gets zero if every point belongs exactly to one cluster
- ⇒ The more clearly defined our clusters get over time, the better  $Q$  approximates  $\mathcal{L}$

## EM algorithm for GMMs

**function** GMM(data matrix  $\mathbf{X}$  with  $N$  rows and  $m$  columns, desired number of Gaussians  $K$ )

**end function**



**function** GMM(data matrix  $\mathbf{X}$  with  $N$  rows and  $m$  columns, desired number of Gaussians  $K$ )

Randomly initialize  $\mu_1, \dots, \mu_K \in \mathbb{R}^m$  (in the convex hull of the data).

**end function**

**function** GMM(data matrix  $\mathbf{X}$  with  $N$  rows and  $m$  columns, desired number of Gaussians  $K$ )

Randomly initialize  $\mu_1, \dots, \mu_K \in \mathbb{R}^m$  (in the convex hull of the data).

Initialize  $p_Z(k) \leftarrow \frac{1}{K}$ , initialize  $\Sigma_k \leftarrow \mathbf{I}$ .

**end function**

**function** GMM(data matrix  $\mathbf{X}$  with  $N$  rows and  $m$  columns, desired number of Gaussians  $K$ )

Randomly initialize  $\mu_1, \dots, \mu_K \in \mathbb{R}^m$  (in the convex hull of the data).

Initialize  $p_Z(k) \leftarrow \frac{1}{K}$ , initialize  $\Sigma_k \leftarrow \mathbf{I}$ .

**for** desired number of iterations **do**

end **for**

**end function**

**function** GMM(data matrix  $\mathbf{X}$  with  $N$  rows and  $m$  columns, desired number of Gaussians  $K$ )

Randomly initialize  $\mu_1, \dots, \mu_K \in \mathbb{R}^m$  (in the convex hull of the data).

Initialize  $p_Z(k) \leftarrow \frac{1}{K}$ , initialize  $\Sigma_k \leftarrow \mathbf{I}$ .

**for** desired number of iterations **do**

    Compute  $\gamma_{k,i} = p_{Z|X}(k|x_i)$  for all  $i \in \{1, \dots, N\}$  and all  $k \in \{1, \dots, K\}$ .

        ▷ Expectation step

**end for**

**end function**

**function** GMM(data matrix  $\mathbf{X}$  with  $N$  rows and  $m$  columns, desired number of Gaussians  $K$ )

Randomly initialize  $\mu_1, \dots, \mu_K \in \mathbb{R}^m$  (in the convex hull of the data).

Initialize  $p_Z(k) \leftarrow \frac{1}{K}$ , initialize  $\Sigma_k \leftarrow \mathbf{I}$ .

**for** desired number of iterations **do**

Compute  $\gamma_{k,i} = p_{Z|X}(k|x_i)$  for all  $i \in \{1, \dots, N\}$  and all  $k \in \{1, \dots, K\}$ .

▷ Expectation step

Set  $p_Z(k) \leftarrow \frac{1}{N} \sum_{i=1}^N \gamma_{i,k}$  for all  $k \in \{1, \dots, K\}$ .

Set  $\mu_k \leftarrow \frac{\sum_{i=1}^N \gamma_{i,k} \cdot x_i}{\sum_{i=1}^N \gamma_{i,k}}$  for all  $k \in \{1, \dots, K\}$ .

Set  $\Sigma_k \leftarrow \frac{\sum_{i=1}^N \gamma_{i,k} \cdot (x_i - \mu_k)(x_i - \mu_k)^T}{\sum_{i=1}^N \gamma_{i,k}}$  for all  $k \in \{1, \dots, K\}$ .

▷ Maximization step

**end for**

**end function**

**function** GMM(data matrix  $\mathbf{X}$  with  $N$  rows and  $m$  columns, desired number of Gaussians  $K$ )

Randomly initialize  $\mu_1, \dots, \mu_K \in \mathbb{R}^m$  (in the convex hull of the data).

Initialize  $p_Z(k) \leftarrow \frac{1}{K}$ , initialize  $\Sigma_k \leftarrow \mathbf{I}$ .

**for** desired number of iterations **do**

Compute  $\gamma_{k,i} = p_{Z|X}(k|x_i)$  for all  $i \in \{1, \dots, N\}$  and all  $k \in \{1, \dots, K\}$ .

▷ Expectation step

Set  $p_Z(k) \leftarrow \frac{1}{N} \sum_{i=1}^N \gamma_{i,k}$  for all  $k \in \{1, \dots, K\}$ .

Set  $\mu_k \leftarrow \frac{\sum_{i=1}^N \gamma_{i,k} \cdot x_i}{\sum_{i=1}^N \gamma_{i,k}}$  for all  $k \in \{1, \dots, K\}$ .

Set  $\Sigma_k \leftarrow \frac{\sum_{i=1}^N \gamma_{i,k} \cdot (x_i - \mu_k)(x_i - \mu_k)^T}{\sum_{i=1}^N \gamma_{i,k}}$  for all  $k \in \{1, \dots, K\}$ .

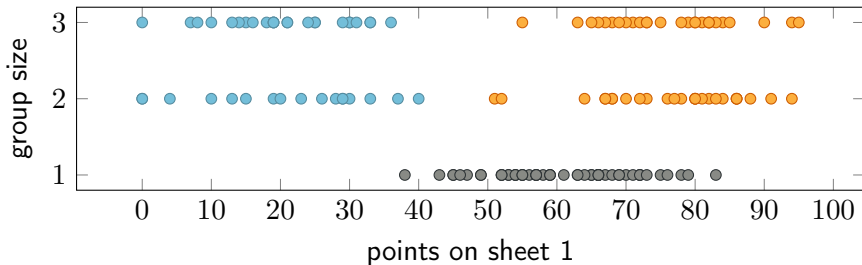
▷ Maximization step

**end for**

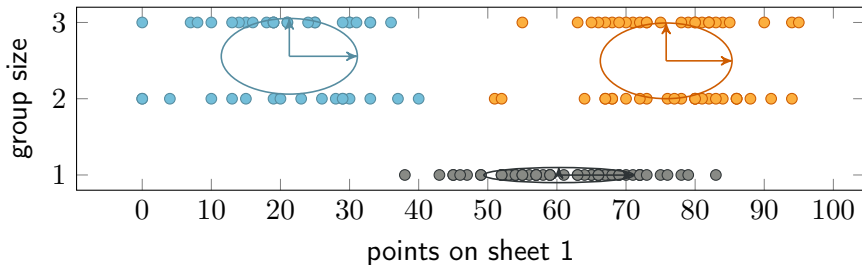
**return**  $p_Z, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K$ .

**end function**

## GMM on example data



## GMM on example data





- ▶  $K$ -Means discovers clusters by assigning data to closest prototype and prototypes to cluster means

- ▶  $K$ -Means discovers clusters by assigning data to closest prototype and prototypes to cluster means
- ▶ Sensitive to data scaling and “unlucky” initializations

- ▶  $K$ -Means discovers clusters by assigning data to closest prototype and prototypes to cluster means
- ▶ Sensitive to data scaling and “unlucky” initializations
- ▶ Selection of  $K$ : Silhouette score and/or BIC

- ▶  $K$ -Means discovers clusters by assigning data to closest prototype and prototypes to cluster means
- ▶ Sensitive to data scaling and “unlucky” initializations
- ▶ Selection of  $K$ : Silhouette score and/or BIC
- ▶ Gaussian mixture models can handle scaling and elliptical cluster shapes – but need much more computation

- Barber, David (2012). **Bayesian Reasoning and Machine Learning**. Cambridge, UK: Cambridge University Press. url:  
<http://www.cs.ucl.ac.uk/staff/d.barber/brml/>.
- Jain, Anil K. (2010). “Data clustering: 50 years beyond K-means”. In: **Pattern Recognition Letters** 31.8. Award winning papers from the 19th International Conference on Pattern Recognition (ICPR), pp. 651–666. doi:  
[10.1016/j.patrec.2009.09.011](https://doi.org/10.1016/j.patrec.2009.09.011).