# Administration

# CSCI567 Machine Learning (Fall 2025)

Haipeng Luo

University of Southern California

Oct 24, 2025

Exam 1 grading still ongoing.

HW3 will be available after today's lecture.

1 / 50

Clusterin

#### Outline

- Clustering
- Question mixture models

Outline

- Clustering
  - Problem setup
  - K-means algorithm
  - Initialization and Convergence
- Question mixture models

3 / 50

4 / 50

# Unsupervised learning

Recall there are different types of machine learning problems

- supervised learning (what we have discussed so far)
  Aim to predict accurately, e.g., classification and regression
- unsupervised learning (this and next few lectures)
   Aim to discover hidden/latent patterns and explore data
- decision making (last two lectures)
   Aim to act optimally under uncertainty

Today's focus: clustering, a canonical unsupervised learning problem

5 / 50

To a section of

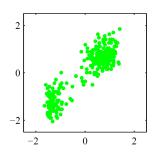
Problem setup

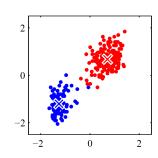
# Clustering: formal definition

**Given**: data points  $oldsymbol{x}_1,\dots,oldsymbol{x}_N\in\mathbb{R}^\mathsf{D}$  and  $\#\mathsf{clusters}\ K$  we want

**Output:** group the data into K clusters, which means

- find assignment  $\gamma_{nk} \in \{0,1\}$  for each data point  $n \in [N]$  and  $k \in [K]$  s.t.  $\sum_{k \in [K]} \gamma_{nk} = 1$  for any fixed n
- ullet find the cluster centers  $oldsymbol{\mu}_1,\ldots,oldsymbol{\mu}_K\in\mathbb{R}^{\mathsf{D}}$



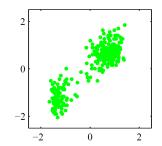


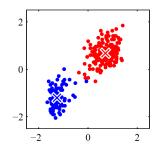
# Clustering: informal definition

Given: a set of data points (feature vectors), without labels

Output: group the data into some clusters, which means

- assign each point to a specific cluster
- find the center (representative/prototype/...) of each cluster





Clustering

Problem setup

# Many applications

- recognize communities in a social network
- group similar customers in market research
- grouping documents into different topics
- image segmentation
- identifying groups of genes with similar expression patterns
- accelerate other algorithms (e.g. NNC as in programing projects)
- . . .

- /

# One example

#### image compression:

- each pixel is a point
- perform clustering over these points
- replace each point by the center of the cluster it belongs to









Original image

Large  $K \longrightarrow \mathsf{Small}\ K$ 

9 / 50

K-means algorithm

# Alternating minimization

Instead, use a heuristic that alternatingly minimizes over  $\{\gamma_{nk}\}$  and  $\{\mu_k\}$ :

Clustering

Initialize  $\{\boldsymbol{\mu}_{k}^{(1)}\}$ 

For t = 1, 2, ...

find

$$\{\gamma_{nk}^{(t+1)}\} = \underset{\{\gamma_{nk}\}}{\operatorname{argmin}} F\left(\{\gamma_{nk}\}, \{\boldsymbol{\mu}_k^{(t)}\}\right)$$

find

$$\{\boldsymbol{\mu}_k^{(t+1)}\} = \operatorname*{argmin}_{\{\boldsymbol{\mu}_k\}} F\left(\{\gamma_{nk}^{(t+1)}\}, \{\boldsymbol{\mu}_k\}\right)$$

# Formal Objective

Key difference from supervised learning problems: no labels given, which means no ground-truth to even measure the quality of your answer!

Still, we can turn it into an optimization problem, e.g. through the popular "K-means" objective: find  $\gamma_{nk}$  and  $\mu_k$  to minimize

$$F(\{\gamma_{nk}\}, \{\boldsymbol{\mu}_k\}) = \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma_{nk} \|\boldsymbol{x}_n - \boldsymbol{\mu}_k\|_2^2$$

i.e. the sum of squared distances of each point to its center.

Unfortunately, finding the exact minimizer is NP-hard!

10 / 50

Clustering

K-means algorithm

#### A closer look

The first step

$$\min_{\{\gamma_{nk}\}} F(\{\gamma_{nk}\}, \{\mu_k\}) = \min_{\{\gamma_{nk}\}} \sum_{n} \sum_{k} \gamma_{nk} \|x_n - \mu_k\|_2^2 
= \sum_{n} \min_{\{\gamma_{nk}\}} \sum_{k} \gamma_{nk} \|x_n - \mu_k\|_2^2$$

is simply to assign each  $x_n$  to the closest  $\mu_k$ , i.e.

$$\gamma_{nk} = egin{cases} 1, & ext{if } k = \operatorname{argmin}_c \|oldsymbol{x}_n - oldsymbol{\mu}_c\|_2^2 \ 0, & ext{else} \end{cases}$$

for all  $k \in [K]$  and  $n \in [N]$ .

#### A closer look

The second step

$$\min_{\{\mu_k\}} F(\{\gamma_{nk}\}, \{\mu_k\}) = \min_{\{\mu_k\}} \sum_n \sum_k \gamma_{nk} ||x_n - \mu_k||_2^2 
= \sum_k \min_{\mu_k} \sum_{n:\gamma_{nk}=1} ||x_n - \mu_k||_2^2$$

is simply to average the points of each cluster (hence the name)

Clustering

K-means algorithm

$$oldsymbol{\mu}_k = rac{\sum_{n:\gamma_{nk}=1} oldsymbol{x}_n}{|\{n:\gamma_{nk}=1\}|} = rac{\sum_n \gamma_{nk} oldsymbol{x}_n}{\sum_n \gamma_{nk}}$$

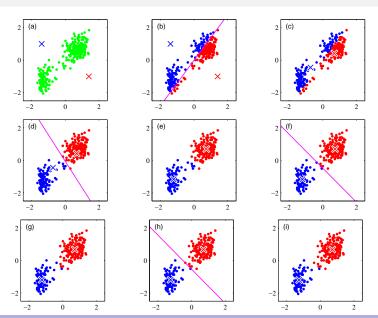
for each  $k \in [K]$ .

13 / 50

Clustering

Initialization and Convergence

# An example



# The K-means algorithm

**Step 0** Initialize  $\mu_1, \dots, \mu_K$ 

**Step 1** Fix the centers  $\mu_1, \ldots, \mu_K$ , assign each point to the closest center:

$$\gamma_{nk} = egin{cases} 1, & \text{if } k = \operatorname{argmin}_c \|m{x}_n - m{\mu}_c\|_2^2 \ 0, & \text{else} \end{cases}$$

Step 2 Fix the assignment  $\{\gamma_{nk}\}$ , update the centers

$$oldsymbol{\mu}_k = rac{\sum_n \gamma_{nk} oldsymbol{x}_n}{\sum_n \gamma_{nk}}$$

Step 3 Return to Step 1 if not converged

How to initialize?

There are different ways to initialize:

- ullet randomly pick K points as initial centers
- or randomly assign each point to a cluster, then average
- or more sophisticated approaches (e.g. greedy or K-means++)

Initialization matters for convergence.

#### lustering Initialization and Convergence

# Convergence

K-means will converge in a finite number of iterations, why?

- objective decreases at each step
- objective is lower bounded by 0
- #possible\_assignments is finite ( $K^N$ , exponentially large though)

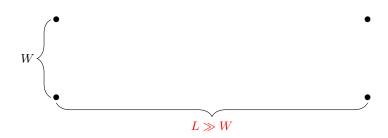
#### However

- it could take exponentially many iterations to converge
- and it *might not converge to the global minimum* of the K-means objective

17 / 50

Clustering Initialization and Convergence

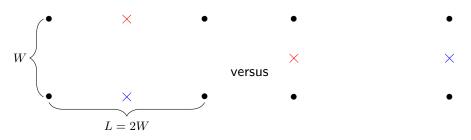
## Local minimum v.s global minimum



- ullet moreover, local minimum can be arbitrarily worse if we increase L
- so initialization matters a lot for K-means

# Local minimum v.s global minimum

Simple example: 4 data points, 2 clusters, 2 different initializations



K-means converges immediately in both cases, but

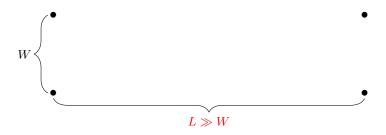
- ullet left has K-means objective  $L^2=4W^2$
- ullet right has K-means objective  $W^2$ , 4 times better than left!
- in fact, left is **local minimum**, and right is **global minimum**.

18 / 50

Clustering Initia

Initialization and Convergence

# How common initialization methods perform?



- ullet randomly pick K points as initial centers: fails with 1/3 probability
- or randomly assign each point to a cluster, then average: similarly fail with a constant probability
- or more sophisticated approaches?

#### Initialization and Convergence

# An Intuitive and Greedy Approach

Idea: spread out the initial centers

Start with a random data point as the first center  $\mu_1$ 

For  $k = 2, \ldots, K$ 

ullet let the k-th center  $\mu_k$  be the farthest from the chosen centers:

$$oldsymbol{\mu}_k = \operatorname*{argmax} \min_{oldsymbol{x} \in \{oldsymbol{x}_1, ..., oldsymbol{x}_N\}} \|oldsymbol{x} - oldsymbol{\mu}_j\|_2^2$$

21 / 50

Chartenia

Initialization and Convergence

# Does greedy initialization always work?

Not really; it is too sensitive to outliers!

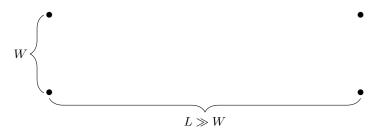


outlier



In HW3, you will verify that greedy initialization can also lead to *arbitrarily bad* local minimums.

# Greedy initialization on the same example



Suppose we pick top left as  $\mu_1$ , then  $\mu_2$  is the bottom right

• K-means converges to the **global minimum** after one iteration!

See demo at https://www.naftaliharris.com/blog/visualizing-k-means-clustering/.

22 / 50

Clustering

Initialization and Convergence

# Solution: K-means++

K-means++: robustify the greedy approach via randomness

Start with a random data point as the first center  $\mu_1$ 

For 
$$k = 2, \dots, K$$

ullet randomly pick the k-th center  $\mu_k$  such that

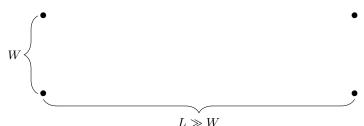
$$\mu_k = \underset{\boldsymbol{x} \in \{\boldsymbol{x}_1, \dots, \boldsymbol{x}_N\}}{\operatorname{argmax}} \min_{j=1,\dots,k-1} \|\boldsymbol{x} - \boldsymbol{\mu}_j\|_2^2$$

$$\Pr[oldsymbol{\mu}_k = oldsymbol{x}_n] \propto \min_{j=1,...,k-1} \|oldsymbol{x}_n - oldsymbol{\mu}_j\|_2^2$$

**K-means++** *guarantees* to find a solution that in expectation is at most  $O(\log K)$  times of the global optimal.

#### ering Initialization and Convergence

# K-means++ on the same example



Suppose we pick top left as  $\mu_1$ , then

- ullet  $\Pr[oldsymbol{\mu}_2 = \mathsf{bottom} \; \mathsf{left}] \propto W^2$ ,  $\Pr[oldsymbol{\mu}_2 = \mathsf{top} \; \mathsf{right}] \propto L^2$
- $\Pr[\mu_2 = \text{bottom right}] \propto W^2 + L^2$

So the expected K-means objective is

$$\frac{W^2}{2(W^2 + L^2)} \cdot L^2 + \left(\frac{L^2}{2(W^2 + L^2)} + \frac{1}{2}\right) \cdot W^2 \le \frac{3}{2}W^2,$$

that is, at most 1.5 times of the optimal.

25 / 50

K-means is alternating minimization for the K-means objective.

The initialization matters a lot for the convergence.

K-means++ uses a theoretically (and often empirically) better initialization.

26 / 50

Gaussian mixture models

Motivation and Model

#### Gaussian mixture models

Summary for K-means

Gaussian mixture models (GMM) is a probabilistic approach for clustering

- more explanatory than minimizing the K-means objective
- can be seen as a soft version of K-means

To solve GMM, we will introduce a powerful method for learning probabilistic model: **Expectation–Maximization (EM) algorithm** 

Gaussian mixture models

#### Outline

- Clustering
- Question mixture models
  - Motivation and Model
  - EM algorithm
  - EM applied to GMMs

27 / 50

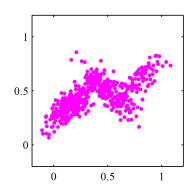
# A generative model

For classification, we discussed the sigmoid model to "explain" how the labels are generated.

Similarly, for clustering, we want to come up with a probabilistic model p to "explain" how the data is generated.

That is, each point is an independent sample of  ${m x} \sim p$ .

What probabilistic model generates data like this?



29 / 50

Gaussian mixture models

Motivation and Model

#### GMM: formal definition

A GMM has the following density function:

$$p(oldsymbol{x}) = \sum_{k=1}^K \omega_k N(oldsymbol{x} \mid oldsymbol{\mu}_k, oldsymbol{\Sigma}_k)$$

where

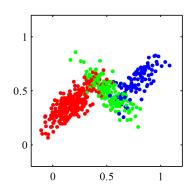
- K: the number of Gaussian components (same as #clusters we want)
- $\omega_1, \ldots, \omega_K$ : mixture weights, a distribution over K components
- ullet  $\mu_k$  and  $\Sigma_k$ : mean and covariance matrix of the k-th Gaussian
- ullet N: the density function for a Gaussian

#### **GMM**: intuition

GMM is a natural model to explain such data

Assume there are 3 ground-truth Gaussian models. To generate a point, we

- first randomly pick one of the Gaussian models,
- then draw a point according this Gaussian.



Hence the name "Gaussian mixture model".

,

Gaussian mixture models

Motivation and Model

#### Another view

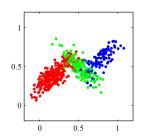
By introducing a latent variable  $z \in [K]$ , which indicates cluster membership, we can see p as a marginal distribution

$$p(\boldsymbol{x}) = \sum_{k=1}^{K} p(\boldsymbol{x}, z = k) = \sum_{k=1}^{K} p(z = k) p(\boldsymbol{x} | z = k) = \sum_{k=1}^{K} \omega_k N(\boldsymbol{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

 $oldsymbol{x}$  and z are both random variables drawn from the model

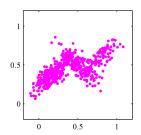
- ullet x is observed
- z is unobserved/latent

## An example



The conditional distributions are

$$\begin{split} p(\boldsymbol{x} \mid z = \mathsf{red}) &= N(\boldsymbol{x} \mid \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \\ p(\boldsymbol{x} \mid z = \mathsf{blue}) &= N(\boldsymbol{x} \mid \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \\ p(\boldsymbol{x} \mid z = \mathsf{green}) &= N(\boldsymbol{x} \mid \boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3) \end{split}$$



The marginal distribution is

$$\begin{split} p(\boldsymbol{x}) &= p(\text{red}) N(\boldsymbol{x} \mid \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + p(\text{blue}) N(\boldsymbol{x} \mid \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \\ &+ p(\text{green}) N(\boldsymbol{x} \mid \boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3) \end{split}$$

33 / 5

Gaussian mixture models

Motivation and Model

## How to learn these parameters?

An obvious attempt is maximum-likelihood estimation (MLE): find

$$\underset{\boldsymbol{\theta}}{\operatorname{argmax}} \ln \prod_{n=1}^{N} p(\boldsymbol{x}_{n}; \boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_{n=1}^{N} \ln p(\boldsymbol{x}_{n}; \boldsymbol{\theta}) \triangleq \underset{\boldsymbol{\theta}}{\operatorname{argmax}} P(\boldsymbol{\theta})$$

This is called incomplete log-likelihood (since  $z_n$ 's are unobserved), and is intractable in general (non-concave problem).

One solution is to still apply GD/SGD, but a much more effective approach is the **Expectation–Maximization (EM) algorithm**.

# Learning GMMs

Learning a GMM means finding all the parameters  $\theta = \{\omega_k, \mu_k, \Sigma_k\}_{k=1}^K$ . In the process, we will learn the latent variable  $z_n$  as well:

$$p(z_n = k \mid \boldsymbol{x}_n) \triangleq \gamma_{nk} \in [0, 1]$$

i.e. "soft assignment" of each point to each cluster, as opposed to "hard assignment" by K-means.

GMM is more explanatory than K-means

- ullet both learn the cluster centers  $oldsymbol{\mu}_k$ 's
- ullet in addition, GMM learns cluster weight  $\omega_k$  and covariance  $\Sigma_k$ , thus
  - we can predict probability of seeing a new point
  - we can generate synthetic data

,

34 / 50

Gaussian mixture models

Motivation and Model

# Preview of EM for learning GMMs

**Step 0** Initialize  $\omega_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$  for each  $k \in [K]$ 

Step 1 (E-Step) update the "soft assignment" (fixing parameters)

$$\gamma_{nk} = p(z_n = k \mid \boldsymbol{x}_n) \propto \omega_k N\left(\boldsymbol{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\right)$$

Step 2 (M-Step) update the model parameter (fixing assignments)

$$\omega_k = rac{\sum_n \gamma_{nk}}{N}$$
  $oldsymbol{\mu}_k = rac{\sum_n \gamma_{nk} oldsymbol{x}_n}{\sum_n \gamma_{nk}}$ 

$$oldsymbol{\Sigma}_k = rac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (oldsymbol{x}_n - oldsymbol{\mu}_k) (oldsymbol{x}_n - oldsymbol{\mu}_k)^{ ext{T}}$$

Step 3 return to Step 1 if not converged

We will see how this is a special case of EM.

#### Demo

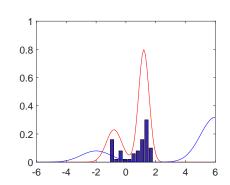
Generate 50 data points from a mixture of 2 Gaussians with

- $\omega_1 = 0.3, \mu_1 = -0.8, \Sigma_1 = 0.52$
- $\omega_2 = 0.7, \mu_2 = 1.2, \Sigma_2 = 0.35$

## histogram represents the data

red curve represents the ground-truth density  $p(\boldsymbol{x}) = \sum_{k=1}^K \omega_k N(\boldsymbol{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ 

blue curve represents the learned density for a specific round



 ${\rm EM\_demo\_1D.pdf}$  shows how the blue curve moves towards red curve quickly via  ${\rm EM}$ 

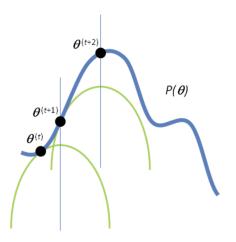
37 / 50

Gaussian mixture models

EM algorithm

## High level idea

Keep maximizing a lower bound of P that is more manageable



# EM algorithm

In general EM is a heuristic to solve MLE with latent variables (not just GMM), i.e. find the maximizer of

$$P(\boldsymbol{\theta}) = \sum_{n=1}^{N} \ln p(\boldsymbol{x}_n ; \boldsymbol{\theta}) = \sum_{n=1}^{N} \ln \int_{z_n} p(\boldsymbol{x}_n, z_n ; \boldsymbol{\theta}) dz_n$$

- $oldsymbol{\theta}$  is the parameters for a general probabilistic model
- $x_n$ 's are observed random variables
- $z_n$ 's are latent variables

Again, directly solving the objective is intractable.

-- /

38 / 50

Gaussian mixture models

EM algorithm

# Derivation of EM (not required)

Finding the lower bound of P:

$$\ln p(\boldsymbol{x}\;;\boldsymbol{\theta}) = \ln \frac{p(\boldsymbol{x},z\;;\boldsymbol{\theta})}{p(z|\boldsymbol{x}\;;\boldsymbol{\theta})} \qquad \text{(true for any } z\text{)}$$

$$= \mathbb{E}_{z\sim q} \left[ \ln \frac{p(\boldsymbol{x},z\;;\boldsymbol{\theta})}{p(z|\boldsymbol{x}\;;\boldsymbol{\theta})} \right] \qquad \text{(true for any dist. } q\text{)}$$

$$= \mathbb{E}_{z\sim q} \left[ \ln p(\boldsymbol{x},z\;;\boldsymbol{\theta}) \right] - \mathbb{E}_{z\sim q} \left[ \ln q(z) \right] - \mathbb{E}_{z\sim q} \left[ \ln \frac{p(z|\boldsymbol{x}\;;\boldsymbol{\theta})}{q(z)} \right]$$

$$= \mathbb{E}_{z\sim q} \left[ \ln p(\boldsymbol{x},z\;;\boldsymbol{\theta}) \right] + \frac{H(q)}{p(z)} - \mathbb{E}_{z\sim q} \left[ \ln \frac{p(z|\boldsymbol{x}\;;\boldsymbol{\theta})}{q(z)} \right] \qquad \text{($H$ is entropy)}$$

$$\geq \mathbb{E}_{z \sim q} \left[ \ln p(oldsymbol{x}, z \; ; oldsymbol{ heta}) 
ight] + H(q) - \ln \mathbb{E}_{z \sim q} \left[ rac{p(z|oldsymbol{x} \; ; oldsymbol{ heta})}{q(z)} 
ight]$$

(Jensen's inequality)

$$= \mathbb{E}_{z \sim q} \left[ \ln p(\boldsymbol{x}, z ; \boldsymbol{\theta}) \right] + H(q)$$

39 / 50

#### 

# Alternatively maximize the lower bound

Therefore, we obtain a lower bound for the log-likelihood function

$$egin{aligned} P(oldsymbol{ heta}) &= \sum_{n=1}^N \ln p(oldsymbol{x}_n \ ; oldsymbol{ heta}) \ &\geq \sum_{n=1}^N \left( \mathbb{E}_{z_n \sim q_n} \left[ \ln p(oldsymbol{x}_n, z_n \ ; oldsymbol{ heta}) 
ight] + H(q_n) 
ight) = oldsymbol{F}(oldsymbol{ heta}, \{q_n\}) \end{aligned}$$

This holds for any  $\{q_n\}$ , so how do we choose? Naturally, the one that maximizes the lower bound (i.e. the tightest lower bound)!

Equivalently, this is the same as alternatingly maximizing F over  $\{q_n\}$  and  $\theta$  (similar to K-means).

41 / 50

Gaussian mixture models

EM algorithm

# Maximizing over heta

Fix  $\{q_n^{(t)}\}$ , maximize over  $\boldsymbol{\theta}$ :

$$\begin{split} & \underset{\boldsymbol{\theta}}{\operatorname{argmax}} F\left(\boldsymbol{\theta}, \{q_n^{(t)}\}\right) \\ &= \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_{n=1}^N \mathbb{E}_{z_n \sim q_n^{(t)}} \left[\ln p(\boldsymbol{x}_n, z_n \ ; \boldsymbol{\theta})\right] \quad \left(H(q_n^{(t)}) \text{ is independent of } \boldsymbol{\theta}\right) \\ &\triangleq \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \ Q(\boldsymbol{\theta} \ ; \boldsymbol{\theta}^{(t)}) & \left(\{q_n^{(t)}\} \text{ are computed via } \boldsymbol{\theta}^{(t)}\right) \end{split}$$

Q is the (expected) **complete likelihood** and is usually more tractable.

ullet versus the incomplete likelihood:  $P(oldsymbol{ heta}) = \sum_{n=1}^N \ln p(oldsymbol{x}_n \ ; oldsymbol{ heta})$ 

# Maximizing over $\{q_n\}$

Fix  $\theta^{(t)}$ , the solution to

$$\operatorname*{argmax}_{q_n} \mathbb{E}_{z_n \sim q_n} \left[ \ln p(\boldsymbol{x}_n, z_n ; \boldsymbol{\theta}^{(t)}) \right] + H(q_n)$$

turns out to be  $q_n^{(t)}$  s.t.

$$q_n^{(t)}(z_n) = p(z_n \mid \boldsymbol{x}_n ; \boldsymbol{\theta}^{(t)}) \propto p(\boldsymbol{x}_n, z_n ; \boldsymbol{\theta}^{(t)})$$

i.e., the *posterior distribution of*  $z_n$  given  $x_n$  and  $\theta^{(t)}$ .

So at  $\theta^{(t)}$ , we found the tightest lower bound  $F\left(\boldsymbol{\theta},\{q_n^{(t)}\}\right)$ :

- $F\left(\boldsymbol{\theta}, \{q_n^{(t)}\}\right) \leq P(\boldsymbol{\theta})$  for all  $\boldsymbol{\theta}$ .
- ullet  $F\left(m{ heta}^{(t)},\{q_n^{(t)}\}
  ight)=P(m{ heta}^{(t)})$  (can be verified through Slide 40)

42 / 50

Gaussian mixture models

EM algorithm

# General EM algorithm

**Step 0** Initialize  $\theta^{(1)}$ , t=1

Step 1 (E-Step) update the posterior of latent variables

$$q_n^{(t)}(\cdot) = p(\cdot \mid \boldsymbol{x}_n ; \boldsymbol{\theta}^{(t)})$$

and obtain Expectation of complete likelihood

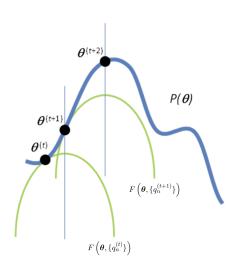
$$Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)}) = \sum_{n=1}^{N} \mathbb{E}_{z_n \sim q_n^{(t)}} \left[ \ln p(\boldsymbol{x}_n, z_n; \boldsymbol{\theta}) \right]$$

Step 2 (M-Step) update the model parameter via Maximization

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \operatorname*{argmax}_{\boldsymbol{\theta}} Q(\boldsymbol{\theta} ; \boldsymbol{\theta}^{(t)})$$

**Step 3**  $t \leftarrow t + 1$  and return to Step 1 if not converged

# Pictorial explanation



 $P(\boldsymbol{\theta})$  is non-concave, but  $Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)})$ often is concave and easy to maximize.

$$P(\boldsymbol{\theta}^{(t+1)}) \ge F\left(\boldsymbol{\theta}^{(t+1)}; \{q_n^{(t)}\}\right)$$
$$\ge F\left(\boldsymbol{\theta}^{(t)}; \{q_n^{(t)}\}\right)$$
$$= P(\boldsymbol{\theta}^{(t)})$$

So EM always increases the objective value and will converge to some local maximum (similar to K-means).

# Apply EM to learn GMMs

E-Step:

$$q_n^{(t)}(z_n = k) = p\left(z_n = k \mid \boldsymbol{x}_n ; \boldsymbol{\theta}^{(t)}\right)$$

$$\propto p\left(\boldsymbol{x}_n, z_n = k ; \boldsymbol{\theta}^{(t)}\right)$$

$$= p\left(z_n = k ; \boldsymbol{\theta}^{(t)}\right) p(\boldsymbol{x}_n \mid z_n = k ; \boldsymbol{\theta}^{(t)})$$

$$= \omega_k^{(t)} N\left(\boldsymbol{x}_n \mid \boldsymbol{\mu}_k^{(t)}, \boldsymbol{\Sigma}_k^{(t)}\right)$$

This computes the "soft assignment"  $\gamma_{nk} = q_n^{(t)}(z_n = k)$ , i.e. conditional probability of  $x_n$  belonging to cluster k.

45 / 50

Gaussian mixture models

EM applied to GMMs

# Apply EM to learn GMMs

M-Step:

$$\begin{aligned} \operatorname*{argmax}_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) &= \operatorname*{argmax}_{\boldsymbol{\theta}} \sum_{n=1}^{N} \mathbb{E}_{z_{n} \sim q_{n}^{(t)}} \left[ \ln p(\boldsymbol{x}_{n}, z_{n} ; \boldsymbol{\theta}) \right] \\ &= \operatorname*{argmax}_{\boldsymbol{\theta}} \sum_{n=1}^{N} \mathbb{E}_{z_{n} \sim q_{n}^{(t)}} \left[ \ln p(z_{n} ; \boldsymbol{\theta}) + \ln p(\boldsymbol{x}_{n} | z_{n} ; \boldsymbol{\theta}) \right] \\ &= \operatorname*{argmax}_{\{\omega_{k}, \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}\}} \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma_{nk} \left( \ln \omega_{k} + \ln N(\boldsymbol{x}_{n} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) \right) \end{aligned}$$

To find  $\omega_1,\ldots,\omega_K$ , solve

To find each  $\mu_k, \Sigma_k$ , solve

$$\underset{\pmb{\omega}}{\operatorname{argmax}} \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma_{nk} \ln \omega_{k}$$

$$\underset{\boldsymbol{\omega}}{\operatorname{argmax}} \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma_{nk} \ln \omega_{k} \qquad \underset{\boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}}{\operatorname{argmax}} \sum_{n=1}^{N} \gamma_{nk} \ln N(\boldsymbol{x}_{n} \mid \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})$$

Gaussian mixture models

EM applied to GMMs

# M-Step (continued)

Solutions to previous two problems are very natural, for each k

$$\omega_k = \frac{\sum_n \gamma_{nk}}{N}$$

i.e. (weighted) fraction of examples belonging to cluster k

$$oldsymbol{\mu}_k = rac{\sum_n \gamma_{nk} oldsymbol{x}_n}{\sum_n \gamma_{nk}}$$

i.e. (weighted) average of examples belonging to cluster k

$$oldsymbol{\Sigma}_k = rac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (oldsymbol{x}_n - oldsymbol{\mu}_k) (oldsymbol{x}_n - oldsymbol{\mu}_k)^{ ext{T}}$$

i.e (weighted) covariance of examples belonging to cluster k

(You can try to verify these for the 1D case.)

# Putting it together

EM for learning GMMs:

(see 2D demo)

**Step 0** Initialize  $\omega_k, \mu_k, \Sigma_k$  for each  $k \in [K]$ 

Step 1 (E-Step) update the "soft assignment" (fixing parameters)

$$\gamma_{nk} = p(z_n = k \mid \boldsymbol{x}_n) \propto \omega_k N(\boldsymbol{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Step 2 (M-Step) update the model parameter (fixing assignments)

$$\omega_k = rac{\sum_n \gamma_{nk}}{N}$$
  $oldsymbol{\mu}_k = rac{\sum_n \gamma_{nk} oldsymbol{x}_n}{\sum_n \gamma_{nk}}$ 

$$oldsymbol{\Sigma}_k = rac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (oldsymbol{x}_n - oldsymbol{\mu}_k) (oldsymbol{x}_n - oldsymbol{\mu}_k)^{ ext{T}}$$

**Step 3** return to Step 1 if not converged

#### Connection to K-means

K-means is in fact a special case of EM for (a simplified) GMM:

- assume  $\Sigma_k = \sigma^2 I$  for some fixed  $\sigma$  so only  $\omega_k$  and  $\mu_k$  are parameters
- when  $\sigma \to 0$ , EM becomes K-means

GMM is a soft version of K-means and it provides a probabilistic interpretation of the data, which means we can predict and generate data after learning.