

CSCI-567: Machine Learning

Prof. Victor Adamchik

U of Southern California

July 23, 2020

Your model is only as good as your data.

July 23, 2020 1 / 47

Outline

- 1 Clustering
- 2 Gaussian mixture models
- 3 Problem Solving

July 23, 2020 2 / 47

Outline

- 1 Clustering
 - Problem setup
 - K-means algorithm
- 2 Gaussian mixture models
- 3 Problem Solving

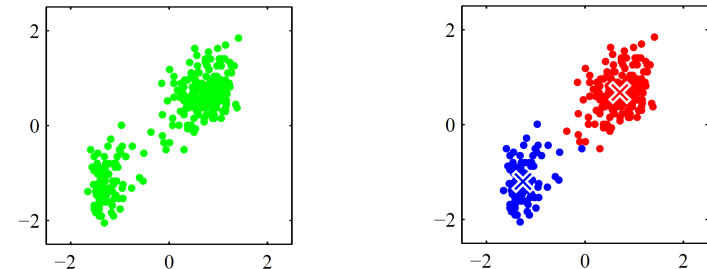
July 23, 2020 3 / 47

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



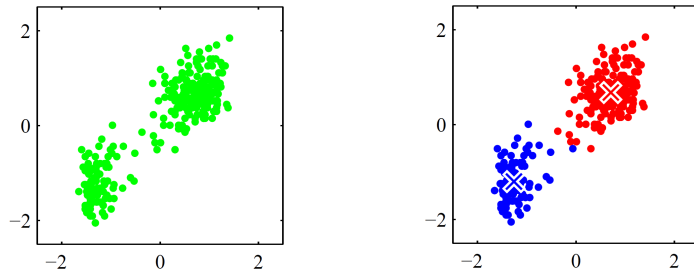
July 23, 2020 4 / 47

Clustering: formal definition

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

Output: group the data into K clusters, which means

- find an **assignment** $\gamma_{nk} \in \{0, 1\}$ s.t. if a data point $n \in [N]$ belongs to a cluster $k \in [K]$ then $\gamma_{nk} = 1$ and $\sum_{k \in [K]} \gamma_{nk} = 1$.
- find the cluster **centers** $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K \in \mathbb{R}^D$



July 23, 2020 5 / 47

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 γ_{nk} and $\boldsymbol{\mu}_k$ to minimize

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

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

Unfortunately, finding the exact minimizer is **NP-hard!**

July 23, 2020 7 / 47

Many applications

One example: **image compression** (vector quantization)

- 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 $K \rightarrow$ Small K

July 23, 2020 6 / 47

Alternating minimization

Instead, use a heuristic that **alternatively minimizes over $\{\gamma_{nk}\}$ and $\{\boldsymbol{\mu}_k\}$** :

Initialize $\{\gamma_{nk}^{(1)}\}$ and $\{\boldsymbol{\mu}_k^{(1)}\}$

For $t = 1, 2, \dots$

- fix centers $\{\boldsymbol{\mu}_k^{(t)}\}$, find assignments $\{\gamma_{nk}^{(t+1)}\}$

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

- fix assignments $\{\gamma_{nk}^{(t+1)}\}$, find new centers $\{\boldsymbol{\mu}_k^{(t+1)}\}$

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

July 23, 2020 8 / 47

A closer look

The first step (fixed centers, find assignments)

$$\operatorname{argmin}_{\{\gamma_{nk}\}} F(\{\gamma_{nk}\}, \{\boldsymbol{\mu}_k\}) = \operatorname{argmin}_{\{\gamma_{nk}\}} \sum_n \sum_k \gamma_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|_2^2$$

is simply to **assign each \mathbf{x}_n to the closest $\boldsymbol{\mu}_k$** , i.e.

$$\gamma_{nk} = \mathbb{I} \left[k == \operatorname{argmin}_c \|\mathbf{x}_n - \boldsymbol{\mu}_c\|_2^2 \right]$$

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

A closer look

The second step (fixed assignments, find centers)

$$\operatorname{argmin}_{\{\boldsymbol{\mu}_k\}} F(\{\gamma_{nk}\}, \{\boldsymbol{\mu}_k\}) = \operatorname{argmin}_{\{\boldsymbol{\mu}_k\}} \sum_n \sum_k \gamma_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|_2^2$$

We will do it for each cluster.

The center is simply **an average of the points in that cluster** (hence the name)

$$\boldsymbol{\mu}_k = \frac{\sum_n \gamma_{nk} \mathbf{x}_n}{\sum_n \gamma_{nk}}$$

for each $k \in [K]$.

The K-means algorithm, S. Lloyd (1957)

Step 0 Initialization (choose K centers)

Step 1 Fix the centers $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$, **assign each point to the closest center**:

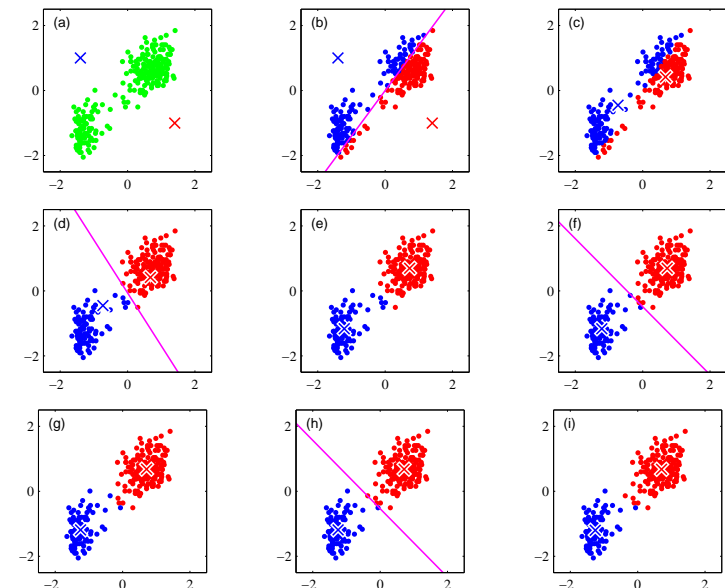
$$\gamma_{nk} = \mathbb{I} \left[k == \operatorname{argmin}_c \|\mathbf{x}_n - \boldsymbol{\mu}_c\|_2^2 \right]$$

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

$$\boldsymbol{\mu}_k = \frac{\sum_n \gamma_{nk} \mathbf{x}_n}{\sum_n \gamma_{nk}}$$

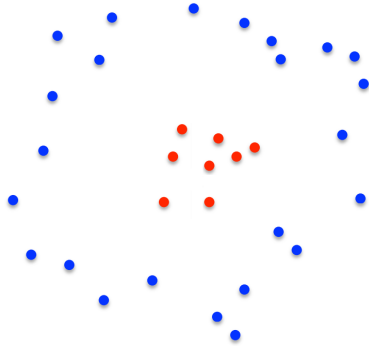
Step 3 Repeat Steps 1 and 2 until the centers no longer change.

An example



K-means algorithm is a heuristic!

K-means is not always able to properly cluster:

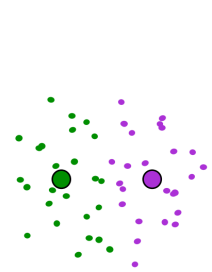


July 23, 2020 13 / 47

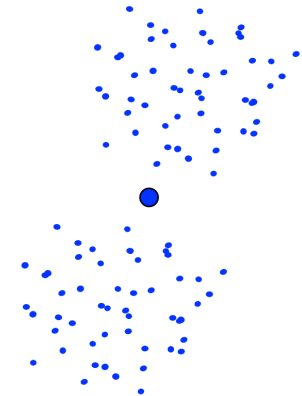
K-means algorithm is a heuristic!

It does matter how you initialize the centers!

In the following example $K = 3$:



Would be better to have
one cluster here



... and two clusters here

July 23, 2020 14 / 47

How to initialize?

A bad selection for the initial centers can lead to a very poor clustering of data.

It also may lead a very long to converge.

There are [different ways to initialize](#):

- randomly pick K points as initial centers
- as it turns out, good initial centers are ones that aren't close to each other. (e.g. **K-means++**, 2007)

July 23, 2020 15 / 47

How to initialize?

The K-means++ algorithm.

The algorithm selects initial centers that aren't close to each other, then uses K-means algorithm for clustering.

The high-level pseudo-code for the K-means++:

- select a data point at random as the first center
- loop $K-1$ times
 - ▶ compute distance squared $d(x)^2$ from each point to the nearest cluster center
 - ▶ select a point that has largest probability $\frac{d(x)^2}{\sum_x d(x)^2}$ as the next center

July 23, 2020 16 / 47

Convergence

It will **converge in a finite number of iterations** to a **local** minimum.

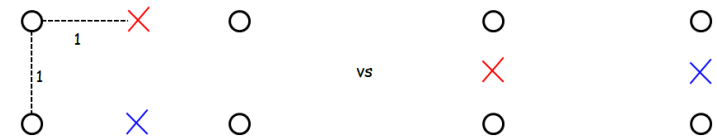
- objective decreases at each step
- objective is lower bounded by 0
- #possible_assignments is finite (K^N , exponentially large though)
- it may take *exponentially* many iterations to converge
- it might not converge to the *global* minimum

July 23, 2020 17 / 47

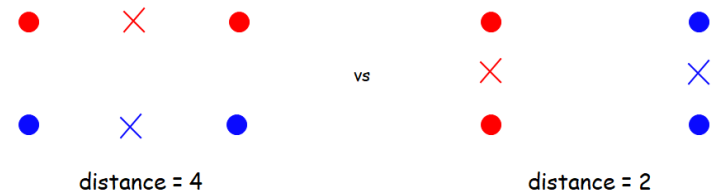
Local minimum v.s global minimum

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

We initialize the centers by the mean of two points.



K-means converges immediately in both cases.

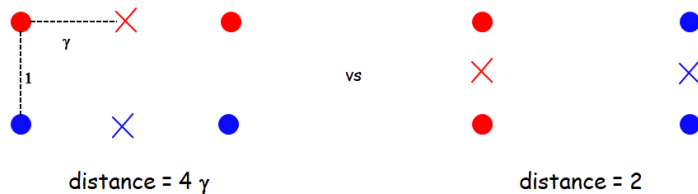


July 23, 2020 18 / 47

Local minimum v.s global minimum

In the left picture we get a local minimum, but in the right - a global minimum!

Moreover, local minimum can be *arbitrarily worse* if we increase the width of this “rectangle” to 2γ .



So, we get stuck at a local minimum.

Initialization matters a lot!

July 23, 2020 19 / 47

Cluster Quality Measures

We need to define a measure of cluster quality Q and then try different values of K until we get an optimal value for Q

There are different metrics for evaluating clustering algorithms, depending on what types of clusters we want

K-means emphasizes similarity of data within clusters:

$$Q = \sum_{k=1}^K \frac{1}{C_k} \sum_{x \in C_k} \|x - \mu_k\|_2^2$$

where C_k is the number of data points in cluster k .

July 23, 2020 20 / 47

Cluster Quality Measures

Other Quality measures:

The aim is to identify sets of clusters that are compact and at the same time are well separated

- Dunn Index
- Davies-Bouldin Index
- Silhouette Index

Outline

- 1 Clustering
- 2 Gaussian mixture models
 - Motivation and Model
 - EM algorithm
- 3 Problem Solving

Taxonomy of ML Models

There are two kinds of classification models in machine learning — [generative](#) models and [discriminative](#) models.

Discriminative models:

- nearest neighbor, k-means clustering, traditional neural networks, SVM.
- we learn $f()$ on data set (x_i, y_i) to output the most likely y on unseen x .
- having $f()$ we know how to discriminate unseen x 's from different classes.
- we learn the decision boundary between the classes.
- we have no idea how the data is generated.

Taxonomy of ML Models

There are two kinds of classification models in machine learning — [generative](#) models and [discriminative](#) models.

Generative models:

- Naïve Bayes, Gaussian mixture model, Hidden Markov model, Adversarial Network (GAN).
- it's used widely in unsupervised machine learning.
- it's a probabilistic way to think about how the data might have been generated.
- learn the joint probability distribution $P(x, y)$ and predict $P(y|x)$ with the help of Bayes Theorem.

Gaussian mixture models

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**

July 23, 2020 25 / 47

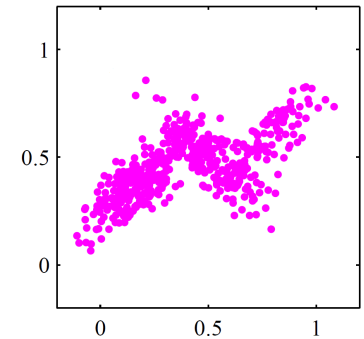
A generative model

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**.

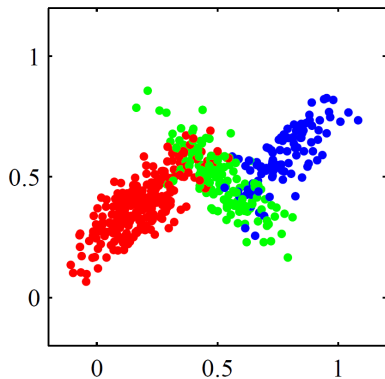
That is, each point is **an independent sample** of $x \sim p$.

What probabilistic model generates data like this?



July 23, 2020 26 / 47

Gaussian mixture models: intuition



We will model each region with a Gaussian distribution. This leads to the idea of Gaussian **mixture** models (GMMs).

The problem we are now facing is that i) we do not know which (color) region a data point comes from; ii) the parameters of Gaussian distributions in each region. We need to find all of them from *unsupervised* data $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N$.

July 23, 2020 27 / 47

GMM: formal definition

A GMM has the following density function:

$$p(\mathbf{x}) = \sum_{k=1}^K \omega_k N(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \sum_{k=1}^K \omega_k \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}_k|}} e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)}$$

where

- K : the number of **Gaussian components** (same as #clusters we want)
- $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$: **mean and covariance matrix** of the k -th Gaussian
- $\omega_1, \dots, \omega_K$: **mixture weights**, they represent how much each component contributes to the final distribution. It satisfies two properties:

$$\forall k, \omega_k > 0, \quad \text{and} \quad \sum_k \omega_k = 1$$

July 23, 2020 28 / 47

Another view

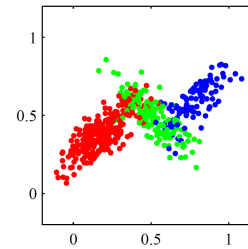
By introducing a **latent variable** $z \in [K]$, which indicates cluster membership, we can see p as a **marginal distribution**

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

\mathbf{x} and z are both random variables drawn from the model

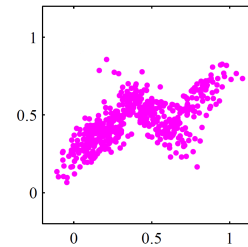
- \mathbf{x} is **observed**
- z is **unobserved/latent**

An example



The conditional distributions are

$$\begin{aligned} p(\mathbf{x} | z = \text{red}) &= N(\mathbf{x} | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \\ p(\mathbf{x} | z = \text{blue}) &= N(\mathbf{x} | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \\ p(\mathbf{x} | z = \text{green}) &= N(\mathbf{x} | \boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3) \end{aligned}$$



The marginal distribution is

$$p(\mathbf{x}) = p(\text{red}) N(\mathbf{x} | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + p(\text{blue}) N(\mathbf{x} | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) + p(\text{green}) N(\mathbf{x} | \boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3)$$

Learning GMMs

Learning a GMM means **finding all the parameters** $\boldsymbol{\theta} = \{\omega_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$.

In the process, we will **learn the latent variable z_n as well**:

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

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

- both learn the cluster centers $\boldsymbol{\mu}_k$'s
- in addition, GMM learns cluster weight ω_k and covariance $\boldsymbol{\Sigma}_k$, thus
 - ▶ we can **predict probability of seeing a new point**
 - ▶ we can **generate synthetic data**

How to learn these parameters?

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

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

This is called **incomplete 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**.

Preview of EM for learning GMMs

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

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

$$\gamma_{nk} = p(z_n = k \mid \mathbf{x}_n) \propto \omega_k N(\mathbf{x}_n \mid \mu_k, \Sigma_k)$$

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

$$\omega_k = \frac{\sum_n \gamma_{nk}}{N} \quad \mu_k = \frac{\sum_n \gamma_{nk} \mathbf{x}_n}{\sum_n \gamma_{nk}}$$
$$\Sigma_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T$$

Step 3 return to Step 1 if not converged

July 23, 2020 33 / 47

EM algorithm

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

$$P(\theta) = \sum_{n=1}^N \ln p(\mathbf{x}_n; \theta)$$

- θ is the **parameters** for a general probabilistic model
- \mathbf{x}_n 's are **observed random variables**
- z_n 's are **latent variables**

Again, directly solving the objective is intractable.

July 23, 2020 34 / 47

EM algorithm

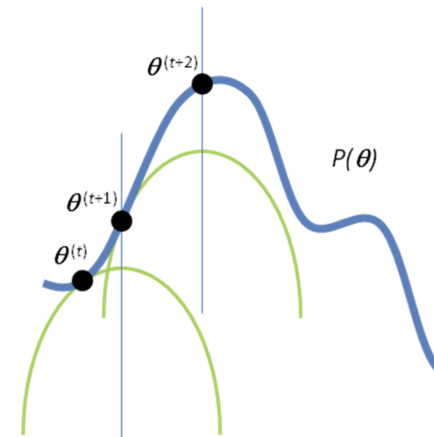
A general algorithm for dealing with hidden data.

- EM is an optimization strategy for objective functions that can be interpreted as likelihoods in the presence of missing data.
- EM is much simpler than gradient methods: no need to choose step size.
- EM is an iterative algorithm with two steps:
 - ▶ E-step: fill-in hidden values using inference
 - ▶ M-step: apply standard MLE method to completed data
- We will prove that EM always converges to a local optimum of the likelihood.

July 23, 2020 35 / 47

High level idea

Keep maximizing **a lower bound of P** that is more manageable



July 23, 2020 36 / 47

Derivation of EM

Finding the lower bound of P :

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

Let us recall the definition of expectation

$$\mathbb{E}_{z \sim q} [f(z)] = \sum_z q(z) f(z)$$

Jensen's inequality

Claim: $\mathbb{E}[\ln X] \leq \ln(\mathbb{E}[X])$

Proof. By the definition of $\mathbb{E}[X] = \frac{1}{N} (x_1 + x_2 + \dots + x_n)$, then

$$\mathbb{E}[\ln X] = \frac{1}{N} (\ln x_1 + \ln x_2 + \dots + \ln x_n) = \frac{1}{N} \ln \prod_{n=1}^N x_n$$

It follows, that the above claim can be rewritten as

$$\frac{1}{N} \ln \prod_{n=1}^N x_n \leq \ln \frac{1}{N} \sum_{n=1}^N x_n$$

$$\sqrt[N]{\prod_{n=1}^N x_n} \leq \frac{1}{N} \sum_{n=1}^N x_n$$

This is the AGM inequality. For $N = 2$, it is just $(x_1 - x_2)^2 \geq 0$.

Derivation of EM

Finding the lower bound of P :

$$\begin{aligned}\ln p(\mathbf{x}; \boldsymbol{\theta}) &= \ln \frac{p(\mathbf{x}, z; \boldsymbol{\theta})}{p(z|\mathbf{x}; \boldsymbol{\theta})} && \text{(true for any } z) \\ &= \mathbb{E}_{z \sim q} \left[\ln \frac{p(\mathbf{x}, z; \boldsymbol{\theta})}{p(z|\mathbf{x}; \boldsymbol{\theta})} \right] && \text{(true for any dist. } q) \\ &= \mathbb{E}_{z \sim q} [\ln p(\mathbf{x}, z; \boldsymbol{\theta})] - \mathbb{E}_{z \sim q} [\ln q(z)] - \mathbb{E}_{z \sim q} \left[\ln \frac{p(z|\mathbf{x}; \boldsymbol{\theta})}{q(z)} \right] \\ &\geq \mathbb{E}_{z \sim q} [\ln p(\mathbf{x}, z; \boldsymbol{\theta})] - \mathbb{E}_{z \sim q} [\ln q(z)] - \ln \mathbb{E}_{z \sim q} \left[\frac{p(z|\mathbf{x}; \boldsymbol{\theta})}{q(z)} \right] \\ &&& \text{(Jensen's inequality)}\end{aligned}$$

Derivation of EM

After applying Jensen's inequality, we obtain

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

Next, we observe that

$$\mathbb{E}_{z \sim q} \left[\frac{p(z|\mathbf{x}; \boldsymbol{\theta})}{q(z)} \right] = \sum_z q(z) \left(\frac{p(z|\mathbf{x}; \boldsymbol{\theta})}{q(z)} \right) = \sum_z p(z|\mathbf{x}; \boldsymbol{\theta}) = 1$$

It follows,

$$\ln p(\mathbf{x}; \boldsymbol{\theta}) \geq \mathbb{E}_{z \sim q} [\ln p(\mathbf{x}, z; \boldsymbol{\theta})] - \mathbb{E}_{z \sim q} [\ln q(z)]$$

Alternatively maximize the lower bound

We have found a lower bound for the log-likelihood function

$$\begin{aligned} P(\boldsymbol{\theta}) &= \sum_{n=1}^N \ln p(\mathbf{x}_n; \boldsymbol{\theta}) \\ &\geq \sum_{n=1}^N \left(\mathbb{E}_{z_n \sim q_n} [\ln p(\mathbf{x}_n, z_n; \boldsymbol{\theta})] - \mathbb{E}_{z_n \sim q_n} [\ln q_n(z_n)] \right) = F(\boldsymbol{\theta}, \{q_n\}) \end{aligned}$$

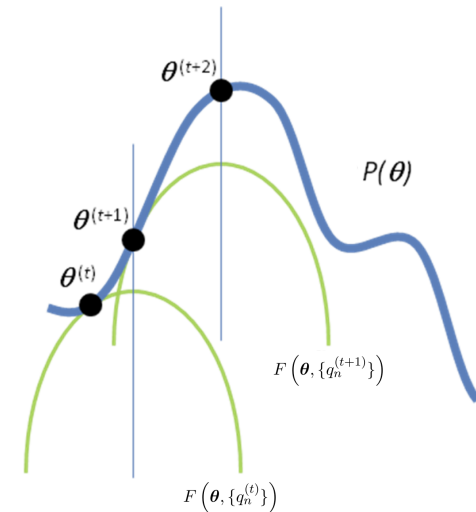
This holds for **any** $\{q_n\}$, so how do we choose?

Naturally, **the one that maximizes the lower bound** (i.e. the tightest lower bound)!

This is similar to K-means: we will alternatively maximizing F over $\{q_n\}$ and $\boldsymbol{\theta}$.

Pictorial explanation

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



Maximizing over $\{q_n\}$

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

$$\begin{aligned} \operatorname{argmax}_{q_n} F(\boldsymbol{\theta}, \{q_n\}) &= \operatorname{argmax}_{q_n} \left(\mathbb{E}_{z_n \sim q_n} [\ln p(\mathbf{x}_n, z_n; \boldsymbol{\theta}^{(t)})] - \mathbb{E}_{z_n \sim q_n} [\ln q_n(z_n)] \right) \\ &= \operatorname{argmax}_{q_n} \sum_{k=1}^K \left(q_n(k) \ln p(\mathbf{x}_n, z_n = k; \boldsymbol{\theta}^{(t)}) - q_n(k) \ln q_n(k) \right) \end{aligned}$$

subject to conditions:

$$q_n(k) \geq 0 \quad \text{and} \quad \sum_k q_n(k) = 1$$

Next, write down the Lagrangian and then apply KKT conditions.

Maximizing over $\{q_n\}$

The solution to

$$\operatorname{argmax}_{q_n} F(\boldsymbol{\theta}, \{q_n\})$$

is (we will solve it in discussions)

$$q_n^{(t)}(z_n) = p(z_n = k | \mathbf{x}_n; \boldsymbol{\theta}^{(t)})$$

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

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

- $F(\boldsymbol{\theta}, \{q_n^{(t)}\}) \leq P(\boldsymbol{\theta})$ for all $\boldsymbol{\theta}$.
- $F(\boldsymbol{\theta}^{(t)}, \{q_n^{(t)}\}) = P(\boldsymbol{\theta}^{(t)})$

Maximizing over θ

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

$$\begin{aligned} & \operatorname{argmax}_{\theta} F(\theta, \{q_n^{(t)}\}) \\ &= \operatorname{argmax}_{\theta} \sum_{n=1}^N \mathbb{E}_{z_n \sim q_n^{(t)}} [\ln p(\mathbf{x}_n, z_n; \theta)] \\ &\triangleq \operatorname{argmax}_{\theta} Q(\theta; \theta^{(t)}) \quad (\{q_n^{(t)}\} \text{ are computed via } \theta^{(t)}) \end{aligned}$$

(we will solve it in discussions)

Q is called a **complete likelihood** and is usually more tractable, since z_n are not latent variables anymore.

July 23, 2020 45 / 47

Outline

- 1 Clustering
- 2 Gaussian mixture models
- 3 Problem Solving

July 23, 2020 47 / 47

Summary

EM is an algorithm to solve MLE with latent variables (not just GMM), i.e. find the maximizer of

$$P(\theta) = \sum_{n=1}^N \ln p(\mathbf{x}_n; \theta)$$

Directly solving the objective is intractable. Instead we optimize the lower bound

$$P(\theta) \geq F(\theta, \{q_n^{(t)}\})$$

where

$$F(\theta, \{q_n^{(t)}\}) = \sum_{n=1}^N \sum_{k=1}^K (q_n(k) \ln p(\mathbf{x}_n, z_n = k; \theta^{(t)}) - q_n(k) \ln q_n(k))$$

July 23, 2020 46 / 47

Problem 1

Maximize the lower bound $F(\theta, \{q_n\})$ over q_n assuming that θ is fixed. See slide 43.

Solution

Problem 2

On slide 45 we defined a complete likelihood $Q(\theta; \theta^{(t)})$.
Maximize Q over μ_k to get

$$\mu_k = \frac{\sum_n \gamma_{nk} x_n}{\sum_n \gamma_{nk}}$$

Solution

Solution
