Administration

CSCI567 Machine Learning (Fall 2020)

Prof. Haipeng Luo

U of Southern California

Oct 22, 2020

Quiz 1: 20% of total grade, everyone gets it

HW3: discuss solutions today

HW4: to be released, due on Sat, 10/31 (note the shorter time)

1/	46 2 / 46
	Clustering
Outline	Outline
Clustering	 Clustering Problem setup K-means algorithm Initialization and Convergence
2 Gaussian mixture models	2 Gaussian mixture models

Clustering Problem setup

Supervised learning v.s unsupervised learning

Clustering

Recall there are different types of machine learning problems

- **supervised learning** (what we have discussed so far) Aim to predict, e.g. classification and regression
- **unsupervised learning** (main focus from now on) Aim to discover hidden/latent patterns and explore data

Today's focus: clustering, an important unsupervised learning problem

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



5 / 46

Clustering Problem setup

Clustering: formal definition

Given: data points $x_1, \ldots, x_N \in \mathbb{R}^{\mathsf{D}}$ and #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
- find the cluster centers $\mu_1, \ldots, \mu_K \in \mathbb{R}^{\mathsf{D}}$





Clustering Problem setup

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

Clustering Problem setup

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 μ_k to minimize

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

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

Unfortunately, finding the exact minimizer is NP-hard!

Alternating minimization

Instead, use a heuristic that alternatingly minimizes over $\{\gamma_{nk}\}$ and $\{\mu_k\}$: Initialize $\{\mu_k^{(1)}\}$ For t = 1, 2, ...• find $\{\gamma_{nk}^{(t+1)}\} = \underset{\{\gamma_{nk}\}}{\operatorname{argmin}} F\left(\{\gamma_{nk}\}, \{\mu_k^{(t)}\}\right)$

find

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

9 / 46
Clustering K-means algorithm
Clustering K-means algorithm
A closer look

The first step

$$\min_{\{\gamma_{nk}\}} F\left(\{\gamma_{nk}\}, \{\boldsymbol{\mu}_k\}\right) = \min_{\{\gamma_{nk}\}} \sum_n \sum_k \gamma_{nk} \|\boldsymbol{x}_n - \boldsymbol{\mu}_k\|_2^2$$
$$= \sum_n \min_{\{\gamma_{nk}\}} \sum_k \gamma_{nk} \|\boldsymbol{x}_n - \boldsymbol{\mu}_k\|_2^2$$

is simply to assign each x_n to the closest μ_k , i.e.

$$\gamma_{nk} = \mathbb{I}\left[k = \operatorname*{argmin}_{c} \|\boldsymbol{x}_{n} - \boldsymbol{\mu}_{c}\|_{2}^{2}\right]$$

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



The second step

$$\min_{\{\boldsymbol{\mu}_k\}} F\left(\{\gamma_{nk}\}, \{\boldsymbol{\mu}_k\}\right) = \min_{\{\boldsymbol{\mu}_k\}} \sum_n \sum_k \gamma_{nk} \|\boldsymbol{x}_n - \boldsymbol{\mu}_k\|_2^2$$
$$= \sum_k \min_{\boldsymbol{\mu}_k} \sum_{n:\gamma_{nk}=1} \|\boldsymbol{x}_n - \boldsymbol{\mu}_k\|_2^2$$

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

$$\boldsymbol{\mu}_{k} = \frac{\sum_{n:\gamma_{nk}=1} \boldsymbol{x}_{n}}{|\{n:\gamma_{nk}=1\}|} = \frac{\sum_{n} \gamma_{nk} \boldsymbol{x}_{n}}{\sum_{n} \gamma_{nk}}$$

for each $k \in [K]$.

K-means algorithm Clustering

The K-means algorithm

Step 0 Initialize μ_1, \ldots, μ_K

Step 1 Fix the centers μ_1, \ldots, μ_K , assign each point to the closest center:

$$\gamma_{nk} = \mathbb{I}\left[k = \operatorname*{argmin}_{c} \|oldsymbol{x}_n - oldsymbol{\mu}_c\|_2^2
ight]$$

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

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

Step 3 Return to Step 1 if not converged

	Clustering	Initialization and Convergence	
How to initialize?			Со

There are different ways to initialize:

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

Initialization matters for convergence.



Clustering Initialization and Convergence

nvergence

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

Clustering Initialization and Convergence

Local minimum v.s global minimum

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



K-means converges immediately in both cases, but

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



- 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: K-means++ guarantees to find a solution that in expectation is at most O(log K) times of the optimal

Local minimum v.s global minimum



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



K-means++ is K-means with a better initialization procedure:

Start with a random data point as the first center μ_1 For $k = 2, \dots, K$

• randomly pick the k-th center μ_k such that

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

Intuitively this *spreads out the initial centers*.

Clustering Initialization and Convergence

K-means++ on the same example



Suppose we pick top left as μ_1 , then

• $\Pr[\mu_2 = \text{bottom left}] \propto W^2$, $\Pr[\mu_2 = \text{top right}] \propto L^2$

Gaussian mixture models

• $\Pr[\mu_2 = \text{bottom right}] \propto W^2 + L^2$

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.

21 / 46

Outline



2 Gaussian mixture models

- Motivation and Model
- EM algorithm
- EM applied to GMMs

Summary for K-means

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.

Gaussian mixture models Motivation and Model

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**

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.



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".



26 / 46

Gaussian mixture models Motivation and Model

GMM: formal definition

A GMM has the following density function:

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

where

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

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

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

- \boldsymbol{x} and \boldsymbol{z} are both random variables drawn from the model
 - x is observed
 - z is unobserved/latent

An example



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



The marginal distribution is

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

Motivation and Model

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:

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

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

- both learn the cluster centers μ_k 's
- in addition, GMM learns cluster weight ω_k and covariance Σ_k , thus
 - we can predict probability of seeing a new point
 - we can *generate synthetic data*

30 / 46

29 / 46

Gaussian mixture models

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**.

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)

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

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

$$\begin{split} \omega_k &= \frac{\sum_n \gamma_{nk}}{N} \qquad \boldsymbol{\mu}_k = \frac{\sum_n \gamma_{nk} \boldsymbol{x}_n}{\sum_n \gamma_{nk}} \\ \boldsymbol{\Sigma}_k &= \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (\boldsymbol{x}_n - \boldsymbol{\mu}_k) (\boldsymbol{x}_n - \boldsymbol{\mu}_k)^{\mathrm{T}} \end{split}$$

Step 3 return to Step 1 if not converged

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

Gaussian mixture models Motivation and Model

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



 $\mathsf{EM}_\mathsf{demo.pdf}$ shows how the blue curve moves towards red curve quickly via EM

Gaussian mixture models EM algorithm

High level idea







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

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

- heta 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.

Gaussian mixture models EM algorithm

Derivation of EM

Finding the lower bound of *P*:

$$\ln p(\boldsymbol{x} ; \boldsymbol{\theta}) = \ln \frac{p(\boldsymbol{x}, z ; \boldsymbol{\theta})}{p(z|\boldsymbol{x} ; \boldsymbol{\theta})} \qquad (\text{true for any } z)$$
$$= \mathbb{E}_{z \sim q} \left[\ln \frac{p(\boldsymbol{x}, z ; \boldsymbol{\theta})}{p(z|\boldsymbol{x} ; \boldsymbol{\theta})} \right] \qquad (\text{true for any dist. } q)$$
$$= \mathbb{E}_{z \sim q} \left[\ln p(\boldsymbol{x}, z ; \boldsymbol{\theta}) \right] - \mathbb{E}_{z \sim q} \left[\ln q(z) \right] - \mathbb{E}_{z \sim q} \left[\ln \frac{p(z|\boldsymbol{x} ; \boldsymbol{\theta})}{q(z)} \right]$$
$$= \mathbb{E}_{z \sim q} \left[\ln p(\boldsymbol{x}, z ; \boldsymbol{\theta}) \right] + H(q) - \mathbb{E}_{z \sim q} \left[\ln \frac{p(z|\boldsymbol{x} ; \boldsymbol{\theta})}{q(z)} \right] \qquad (H \text{ is entropy})$$
$$\geq \mathbb{E}_{z \sim q} \left[\ln p(\boldsymbol{x}, z ; \boldsymbol{\theta}) \right] + H(q) - \ln \mathbb{E}_{z \sim q} \left[\frac{p(z|\boldsymbol{x} ; \boldsymbol{\theta})}{q(z)} \right] \qquad (Jensen's inequality)$$

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

33 / 46

aussian mixture models EM algorithm

Alternatively maximize the lower bound

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

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

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

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)!

Gaussian mixture models

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

EM algorithm

37 / 46

Maximizing over θ

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

$$\underset{\boldsymbol{\theta}}{\operatorname{argmax}} F\left(\boldsymbol{\theta}, \{q_n^{(t)}\}\right)$$

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

$$\triangleq \underset{\boldsymbol{\theta}}{\operatorname{argmax}} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(t)}) \qquad \left(\{q_n^{(t)}\} \text{ are computed via } \boldsymbol{\theta}^{(t)}\right)$$

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

Maximizing over $\{q_n\}$

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

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

is $q_n^{(t)}$ 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 $\theta^{(t)}$. (Verified in HW4)

So at
$$\theta^{(t)}$$
, we found the tightest lower bound $F\left(\theta, \{q_n^{(t)}\}\right)$:
• $F\left(\theta, \{q_n^{(t)}\}\right) \leq P(\theta)$ for all θ .
• $F\left(\theta^{(t)}, \{q_n^{(t)}\}\right) = P(\theta^{(t)})$ (verify yourself by going through Slide 36

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



 $P(\theta)$ is non-concave, but $Q(\theta; \theta^{(t)})$ 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.

41 / 46

Gaussian mixture models EM applied to GMMs

Apply EM to learn GMMs

M-Step:

$$\operatorname{argmax}_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) = \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{n=1}^{N} \mathbb{E}_{z_n \sim q_n^{(t)}} \left[\ln p(\boldsymbol{x}_n, z_n ; \boldsymbol{\theta}) \right]$$
$$= \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{n=1}^{N} \mathbb{E}_{z_n \sim q_n^{(t)}} \left[\ln p(z_n ; \boldsymbol{\theta}) + \ln p(\boldsymbol{x}_n | z_n ; \boldsymbol{\theta}) \right]$$
$$= \operatorname{argmax}_{\{\omega_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}} \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma_{nk} \left(\ln \omega_k + \ln N(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)$$

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

To find each μ_k, Σ_k , solve

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

Gaussian mixture models EM applied to GMMs

M-Step (continued)

Solutions to previous two problems are very natural, for each \boldsymbol{k}

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

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

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

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

$$\boldsymbol{\Sigma}_{k} = rac{1}{\sum_{n} \gamma_{nk}} \sum_{n} \gamma_{nk} (\boldsymbol{x}_{n} - \boldsymbol{\mu}_{k}) (\boldsymbol{x}_{n} - \boldsymbol{\mu}_{k})^{\mathrm{T}}$$

i.e (weighted) covariance of examples belonging to cluster \boldsymbol{k}

You will verify some of these in HW4.

43 / 46

Gaussian mixture models EM applied to GMMs

Putting it together

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)

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

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

$$\omega_k = \frac{\sum_n \gamma_{nk}}{N} \qquad \boldsymbol{\mu}_k = \frac{\sum_n \gamma_{nk} \boldsymbol{x}_n}{\sum_n \gamma_{nk}}$$
$$\boldsymbol{\Sigma}_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (\boldsymbol{x}_n - \boldsymbol{\mu}_k) (\boldsymbol{x}_n - \boldsymbol{\mu}_k)^{\mathrm{T}}$$

Step 3 return to Step 1 if not converged

45 / 46

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 σ so only ω_k and μ_k are parameters
- when $\sigma \rightarrow 0$, 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.