

Introduction to Machine Learning

Lecture 9 Representation and Clustering - Gaussian Mixture Models and EM Algorithm

Hongteng Xu



中國人民大學
RENMIN UNIVERSITY OF CHINA

高瓴人工智能學院
Gaoling School of Artificial Intelligence

Outline

Review

- ▶ Kmeans
- ▶ Spectral clustering and its connections to Kmeans and manifold learning
- ▶ Evaluation of clustering results

Outline

Review

- ▶ Kmeans
- ▶ Spectral clustering and its connections to Kmeans and manifold learning
- ▶ Evaluation of clustering results

Today

- ▶ Generative modeling and Gaussian mixture model
- ▶ EM algorithm

Generative Model v.s. Discriminative Model

Denote X as the random variable of samples, optionally Y as the random variable of labels.

- ▶ **Generative Model:**

- ▶ **Functionality:** Capture the mechanism of generating data

Generative Model v.s. Discriminative Model

Denote X as the random variable of samples, optionally Y as the random variable of labels.

- ▶ **Generative Model:**

- ▶ **Functionality:** Capture the mechanism of generating data
- ▶ **Principle:** Model the data distribution $P(X)$ or the joint distribution $P(X, Y)$

Generative Model v.s. Discriminative Model

Denote X as the random variable of samples, optionally Y as the random variable of labels.

- ▶ **Generative Model:**

- ▶ **Functionality:** Capture the mechanism of generating data
- ▶ **Principle:** Model the data distribution $P(X)$ or the joint distribution $P(X, Y)$

- ▶ **Discriminative Model:**

- ▶ **Functionality:** Capture the difference between different data points

Generative Model v.s. Discriminative Model

Denote X as the random variable of samples, optionally Y as the random variable of labels.

- ▶ **Generative Model:**

- ▶ **Functionality:** Capture the mechanism of generating data
- ▶ **Principle:** Model the data distribution $P(X)$ or the joint distribution $P(X, Y)$

- ▶ **Discriminative Model:**

- ▶ **Functionality:** Capture the difference between different data points
- ▶ **Principle:** Model the conditional probability $P(Y|X)$

Generative Model v.s. Discriminative Model

Denote X as the random variable of samples, optionally Y as the random variable of labels.

- ▶ **Generative Model:**

- ▶ **Functionality:** Capture the mechanism of generating data
- ▶ **Principle:** Model the data distribution $P(X)$ or the joint distribution $P(X, Y)$

- ▶ **Discriminative Model:**

- ▶ **Functionality:** Capture the difference between different data points
- ▶ **Principle:** Model the conditional probability $P(Y|X)$

- ▶ How to categorize the models we learned before?

Generative Model v.s. Discriminative Model

Denote X as the random variable of samples, optionally Y as the random variable of labels.

- ▶ **Generative Model:**

- ▶ **Functionality:** Capture the mechanism of generating data
- ▶ **Principle:** Model the data distribution $P(X)$ or the joint distribution $P(X, Y)$

- ▶ **Discriminative Model:**

- ▶ **Functionality:** Capture the difference between different data points
- ▶ **Principle:** Model the conditional probability $P(Y|X)$

- ▶ How to categorize the models we learned before?

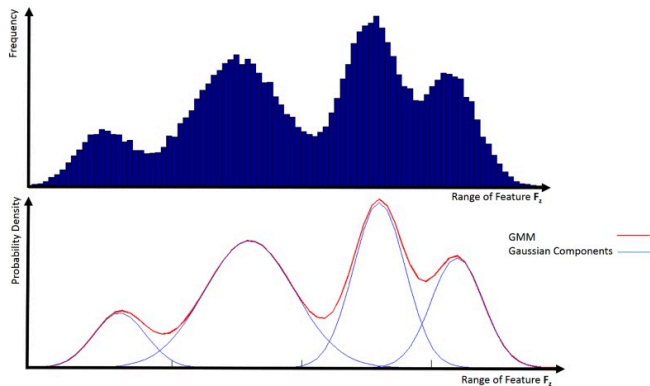
- ▶ Which one is harder?

Let's design the two modeling strategies in practice

- ▶ **Data:** Given $\{x_n, y_n\}_{n=1}^N$, where x 's represent education years, and y 's represent yearly incomes.
- ▶ **Task:** Learn an estimator predicting yearly incomes based on education years.

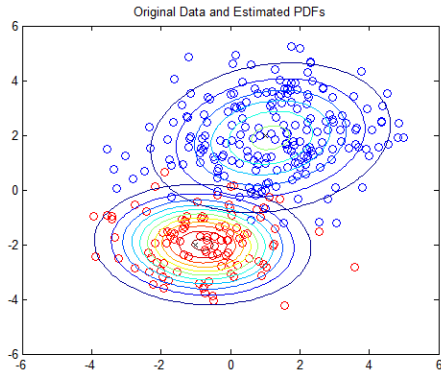
Gaussian Mixture Models

A generative model for the data with clustering structures



Gaussian Mixture Models

A generative model for the data with clustering structures



GMM: Generative Mechanism

Suppose that there are K Gaussian distributions defined on the sample space $\mathcal{X} \subset \mathbb{R}^D$.

- ▶ $\mathbf{w} = [w_k] \in \Delta^{K-1}$
- ▶ $\{\mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\}_{k=1}^K$

GMM: Generative Mechanism

Suppose that there are K Gaussian distributions defined on the sample space $\mathcal{X} \subset \mathbb{R}^D$.

- ▶ $\mathbf{w} = [w_k] \in \Delta^{K-1}$
- ▶ $\{\mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\}_{k=1}^K$

Generative process:

- 1 Determine the cluster: $k \sim \text{Categorical}(\mathbf{w})$
- 2 Determine the sample based on the cluster: $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$.

GMM: Generative Mechanism

Suppose that there are K Gaussian distributions defined on the sample space $\mathcal{X} \subset \mathbb{R}^D$.

- ▶ $\mathbf{w} = [w_k] \in \Delta^{K-1}$
- ▶ $\{\mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\}_{k=1}^K$

Generative process:

- 1 Determine the cluster: $k \sim \text{Categorical}(\mathbf{w})$
- 2 Determine the sample based on the cluster: $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$.

How to achieve the above sampling processes?

The MLE Framework of GMM

What is the probability of \mathbf{x} ? (Derive it)

The MLE Framework of GMM

What is the probability of \mathbf{x} ? (Derive it)

$$p(\mathbf{x}) = \sum_{k=1}^K p(k)p(\mathbf{x}|k)$$

The MLE Framework of GMM

What is the probability of \mathbf{x} ? (Derive it)

$$p(\mathbf{x}) = \sum_{k=1}^K p(k)p(\mathbf{x}|k) = \sum_{k=1}^K w_k \det(2\pi \Sigma_k)^{-\frac{1}{2}} \exp \left(-\frac{1}{2}(\mathbf{x} - \mu_k) \Sigma_k^{-1} (\mathbf{x} - \mu_k) \right) \quad (1)$$

The MLE Framework of GMM

What is the probability of \mathbf{x} ? (Derive it)

$$p(\mathbf{x}) = \sum_{k=1}^K p(k)p(\mathbf{x}|k) = \sum_{k=1}^K w_k \det(2\pi \Sigma_k)^{-\frac{1}{2}} \exp \left(-\frac{1}{2}(\mathbf{x} - \mu_k) \Sigma_k^{-1} (\mathbf{x} - \mu_k) \right) \quad (1)$$

Given $\{\mathbf{x}_n\}_{n=1}^N$, what is the likelihood function of model parameters?

The MLE Framework of GMM

What is the probability of \mathbf{x} ? (Derive it)

$$p(\mathbf{x}) = \sum_{k=1}^K p(k)p(\mathbf{x}|k) = \sum_{k=1}^K w_k \det(2\pi \Sigma_k)^{-\frac{1}{2}} \exp \left(-\frac{1}{2}(\mathbf{x} - \mu_k) \Sigma_k^{-1} (\mathbf{x} - \mu_k) \right) \quad (1)$$

Given $\{\mathbf{x}_n\}_{n=1}^N$, what is the likelihood function of model parameters?

$$L(\theta) = \prod_{n=1}^N p(\mathbf{x}_n; \theta) \quad (2)$$

The Challenges of GMM's MLE

Phenomenon in formulation:

- ▶ $\log L(\theta)$ falls into a log-sum-exp formulation.

The Challenges of GMM's MLE

Phenomenon in formulation:

- ▶ $\log L(\theta)$ falls into a log-sum-exp formulation.
- ▶ Hard to optimize (Could you derive its gradient?)

The Challenges of GMM's MLE

Phenomenon in formulation:

- ▶ $\log L(\theta)$ falls into a log-sum-exp formulation.
- ▶ Hard to optimize (Could you derive its gradient?)

Reason in Statistics:

- ▶ Hidden/unobserved generative process — we don't know which cluster is each sample from?

The Challenges of GMM's MLE

Phenomenon in formulation:

- ▶ $\log L(\theta)$ falls into a log-sum-exp formulation.
- ▶ Hard to optimize (Could you derive its gradient?)

Reason in Statistics:

- ▶ Hidden/unobserved generative process — we don't know which cluster is each sample from?
- ▶ What if we know it?

Learning GMM: An EM Algorithm

In practice, we often apply an **expectation-maximization (EM)** algorithm to learn GMM.

- **E-step:** Given current parameters $\{w_k^{(t)}, \mu_k^{(t)}, \Sigma_k^{(t)}\}_{k=1}^K$, calculate **responsibility** (the posterior distribution of clusters given a sample.)

$$p^{(t)}(k|\mathbf{x}_n) = \frac{w_k^{(t)} p(\mathbf{x}_n; \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{i=1}^K w_i^{(t)} p(\mathbf{x}_n; \mu_i^{(t)}, \Sigma_i^{(t)})}, \quad \forall n = 1, \dots, N, k = 1, \dots, K$$

Learning GMM: An EM Algorithm

In practice, we often apply an **expectation-maximization (EM)** algorithm to learn GMM.

- **E-step:** Given current parameters $\{w_k^{(t)}, \mu_k^{(t)}, \Sigma_k^{(t)}\}_{k=1}^K$, calculate **responsibility** (the posterior distribution of clusters given a sample.)

$$p^{(t)}(k|\mathbf{x}_n) = \frac{w_k^{(t)} p(\mathbf{x}_n; \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{i=1}^K w_i^{(t)} p(\mathbf{x}_n; \mu_i^{(t)}, \Sigma_i^{(t)})}, \quad \forall n = 1, \dots, N, k = 1, \dots, K \quad (3)$$

How to explain the denominator and the numerator?

Learning GMM: An EM Algorithm

In practice, we often apply an **expectation-maximization (EM)** algorithm to learn GMM.

- **E-step:** Given current parameters $\{w_k^{(t)}, \mu_k^{(t)}, \Sigma_k^{(t)}\}_{k=1}^K$, calculate **responsibility** (the posterior distribution of clusters given a sample.)

$$p^{(t)}(k|\mathbf{x}_n) = \frac{w_k^{(t)} p(\mathbf{x}_n; \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{i=1}^K w_i^{(t)} p(\mathbf{x}_n; \mu_i^{(t)}, \Sigma_i^{(t)})}, \quad \forall n = 1, \dots, N, k = 1, \dots, K \quad (3)$$

How to explain the denominator and the numerator?

- **M-step:** Update model parameters based on the responsibilities

$$w_k^{(t+1)} = \frac{1}{N} \sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n),$$

Learning GMM: An EM Algorithm

In practice, we often apply an **expectation-maximization (EM)** algorithm to learn GMM.

- **E-step:** Given current parameters $\{w_k^{(t)}, \mu_k^{(t)}, \Sigma_k^{(t)}\}_{k=1}^K$, calculate **responsibility** (the posterior distribution of clusters given a sample.)

$$p^{(t)}(k|\mathbf{x}_n) = \frac{w_k^{(t)} p(\mathbf{x}_n; \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{i=1}^K w_i^{(t)} p(\mathbf{x}_n; \mu_i^{(t)}, \Sigma_i^{(t)})}, \quad \forall n = 1, \dots, N, k = 1, \dots, K \quad (3)$$

How to explain the denominator and the numerator?

- **M-step:** Update model parameters based on the responsibilities

$$w_k^{(t+1)} = \frac{1}{N} \sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n), \quad \mu_k^{(t+1)} = \frac{\sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n) \mathbf{x}_n}{\sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n)}$$

Learning GMM: An EM Algorithm

In practice, we often apply an **expectation-maximization (EM)** algorithm to learn GMM.

- **E-step:** Given current parameters $\{w_k^{(t)}, \mu_k^{(t)}, \Sigma_k^{(t)}\}_{k=1}^K$, calculate **responsibility** (the posterior distribution of clusters given a sample.)

$$p^{(t)}(k|\mathbf{x}_n) = \frac{w_k^{(t)} p(\mathbf{x}_n; \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{i=1}^K w_i^{(t)} p(\mathbf{x}_n; \mu_i^{(t)}, \Sigma_i^{(t)})}, \quad \forall n = 1, \dots, N, k = 1, \dots, K \quad (3)$$

How to explain the denominator and the numerator?

- **M-step:** Update model parameters based on the responsibilities

$$\begin{aligned} w_k^{(t+1)} &= \frac{1}{N} \sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n), & \mu_k^{(t+1)} &= \frac{\sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n) \mathbf{x}_n}{\sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n)} \\ \Sigma_k^{(t+1)} &= \frac{\sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n) (\mathbf{x}_n - \mu_k^{(t+1)}) (\mathbf{x}_n - \mu_k^{(t+1)})^T}{\sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n)} \end{aligned} \quad (4)$$

An Optimization Viewpoint of The EM Algorithm

Original MLE: Denote $\theta = \{w_k, \mu_k, \Sigma_k\}_{k=1}^K$, the MLE of GMM corresponds to

$$\min_{\theta} \sum_{n=1}^N -\log p(\mathbf{x}_n; \theta) = \min_{\theta} \underbrace{\sum_{n=1}^N -\log \left(\sum_{k=1}^K w_k p(\mathbf{x}_n; \mu_k, \Sigma_k) \right)}_{L(\theta)}.$$

An Optimization Viewpoint of The EM Algorithm

Original MLE: Denote $\theta = \{w_k, \mu_k, \Sigma_k\}_{k=1}^K$, the MLE of GMM corresponds to

$$\min_{\theta} \sum_{n=1}^N -\log p(\mathbf{x}_n; \theta) = \min_{\theta} \underbrace{\sum_{n=1}^N -\log \left(\sum_{k=1}^K w_k p(\mathbf{x}_n; \mu_k, \Sigma_k) \right)}_{L(\theta)}. \quad (5)$$

- $L(\theta)$ is a convex function.

An Optimization Viewpoint of The EM Algorithm

Original MLE: Denote $\theta = \{w_k, \mu_k, \Sigma_k\}_{k=1}^K$, the MLE of GMM corresponds to

$$\min_{\theta} \sum_{n=1}^N -\log p(\mathbf{x}_n; \theta) = \min_{\theta} \underbrace{\sum_{n=1}^N -\log \left(\sum_{k=1}^K w_k p(\mathbf{x}_n; \mu_k, \Sigma_k) \right)}_{L(\theta)}. \quad (5)$$

- ▶ $L(\theta)$ is a convex function.
- ▶ Apply **Jensen's inequality** — instead of minimizing $L(\theta)$, we minimize its upper bound $Q(\theta; \theta^{(t)})$ instead.

An Optimization Viewpoint of The EM Algorithm

Original MLE: Denote $\theta = \{w_k, \mu_k, \Sigma_k\}_{k=1}^K$, the MLE of GMM corresponds to

$$\min_{\theta} \sum_{n=1}^N -\log p(\mathbf{x}_n; \theta) = \min_{\theta} \underbrace{\sum_{n=1}^N -\log \left(\sum_{k=1}^K w_k p(\mathbf{x}_n; \mu_k, \Sigma_k) \right)}_{L(\theta)}. \quad (5)$$

- ▶ $L(\theta)$ is a convex function.
- ▶ Apply **Jensen's inequality** — instead of minimizing $L(\theta)$, we minimize its upper bound $Q(\theta; \theta^{(t)})$ instead.

$$L(\theta) = - \sum_{n=1}^N \log \left(\sum_{k=1}^K w_k p(\mathbf{x}_n; \mu_k, \Sigma_k) \right)$$

An Optimization Viewpoint of The EM Algorithm

Original MLE: Denote $\theta = \{w_k, \mu_k, \Sigma_k\}_{k=1}^K$, the MLE of GMM corresponds to

$$\min_{\theta} \sum_{n=1}^N -\log p(\mathbf{x}_n; \theta) = \min_{\theta} \underbrace{\sum_{n=1}^N -\log \left(\sum_{k=1}^K w_k p(\mathbf{x}_n; \mu_k, \Sigma_k) \right)}_{L(\theta)}. \quad (5)$$

- ▶ $L(\theta)$ is a convex function.
- ▶ Apply **Jensen's inequality** — instead of minimizing $L(\theta)$, we minimize its upper bound $Q(\theta; \theta^{(t)})$ instead.

$$\begin{aligned} L(\theta) &= - \sum_{n=1}^N \log \left(\sum_{k=1}^K w_k p(\mathbf{x}_n; \mu_k, \Sigma_k) \right) \\ &= - \sum_{n=1}^N \log \left(\sum_{k=1}^K \frac{p^{(t)}(k|\mathbf{x}_n)}{p^{(t)}(k|\mathbf{x}_n)} w_k p(\mathbf{x}_n; \mu_k, \Sigma_k) \right) \end{aligned}$$

An Optimization Viewpoint of The EM Algorithm

Original MLE: Denote $\theta = \{w_k, \mu_k, \Sigma_k\}_{k=1}^K$, the MLE of GMM corresponds to

$$\min_{\theta} \sum_{n=1}^N -\log p(\mathbf{x}_n; \theta) = \min_{\theta} \underbrace{\sum_{n=1}^N -\log \left(\sum_{k=1}^K w_k p(\mathbf{x}_n; \mu_k, \Sigma_k) \right)}_{L(\theta)}. \quad (5)$$

- ▶ $L(\theta)$ is a convex function.
- ▶ Apply **Jensen's inequality** — instead of minimizing $L(\theta)$, we minimize its upper bound $Q(\theta; \theta^{(t)})$ instead.

$$\begin{aligned} L(\theta) &= - \sum_{n=1}^N \log \left(\sum_{k=1}^K w_k p(\mathbf{x}_n; \mu_k, \Sigma_k) \right) \\ &= - \sum_{n=1}^N \log \left(\sum_{k=1}^K \frac{p^{(t)}(k|\mathbf{x}_n)}{p^{(t)}(k|\mathbf{x}_n)} w_k p(\mathbf{x}_n; \mu_k, \Sigma_k) \right) \\ &\leq - \sum_{n=1}^N \sum_{k=1}^K p^{(t)}(k|\mathbf{x}_n) \log \frac{w_k p(\mathbf{x}_n; \mu_k, \Sigma_k)}{p^{(t)}(k|\mathbf{x}_n)} = Q(\theta; \theta^{(t)}) \end{aligned} \quad (6)$$

An Optimization Viewpoint of The EM Algorithm

The optimization problem in the t -th step:

$$\min_{\theta} Q(\theta; \theta^{(t)}) = \min_{\theta} - \sum_{n=1}^N \sum_{k=1}^K p^{(t)}(k|\mathbf{x}_n) \log \frac{w_k p(\mathbf{x}_n; \mu_k, \Sigma_k)}{p^{(t)}(k|\mathbf{x}_n)}.$$

An Optimization Viewpoint of The EM Algorithm

The optimization problem in the t -th step:

$$\min_{\theta} Q(\theta; \theta^{(t)}) = \min_{\theta} - \sum_{n=1}^N \sum_{k=1}^K p^{(t)}(k|\mathbf{x}_n) \log \frac{w_k p(\mathbf{x}_n; \mu_k, \Sigma_k)}{p^{(t)}(k|\mathbf{x}_n)}. \quad (7)$$

Then, in the M-step:

$$\begin{aligned} & \min_{\mathbf{w} \in \Delta^{K-1}} - \sum_{n=1}^N \sum_{k=1}^K p^{(t)}(k|\mathbf{x}_n) \log w_k \\ &= \min_{\mathbf{w} \in \Delta^{K-1}} - \sum_{k=1}^K \underbrace{\frac{\sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n)}{N}}_{p_k^{(t)}} \log w_k \end{aligned}$$

An Optimization Viewpoint of The EM Algorithm

The optimization problem in the t -th step:

$$\min_{\theta} Q(\theta; \theta^{(t)}) = \min_{\theta} - \sum_{n=1}^N \sum_{k=1}^K p^{(t)}(k|\mathbf{x}_n) \log \frac{w_k p(\mathbf{x}_n; \mu_k, \Sigma_k)}{p^{(t)}(k|\mathbf{x}_n)}. \quad (7)$$

Then, in the M-step:

$$\begin{aligned} & \min_{\mathbf{w} \in \Delta^{K-1}} - \sum_{n=1}^N \sum_{k=1}^K p^{(t)}(k|\mathbf{x}_n) \log w_k \\ &= \min_{\mathbf{w} \in \Delta^{K-1}} - \sum_{k=1}^K \underbrace{\frac{\sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n)}{N}}_{p_k^{(t)}} \log w_k = \min_{\mathbf{w}} \text{KL}(\mathbf{p}^{(t)} \parallel \mathbf{w}) \end{aligned}$$

An Optimization Viewpoint of The EM Algorithm

The optimization problem in the t -th step:

$$\min_{\theta} Q(\theta; \theta^{(t)}) = \min_{\theta} - \sum_{n=1}^N \sum_{k=1}^K p^{(t)}(k|\mathbf{x}_n) \log \frac{w_k p(\mathbf{x}_n; \mu_k, \Sigma_k)}{p^{(t)}(k|\mathbf{x}_n)}. \quad (7)$$

Then, in the M-step:

$$\begin{aligned} & \min_{\mathbf{w} \in \Delta^{K-1}} - \sum_{n=1}^N \sum_{k=1}^K p^{(t)}(k|\mathbf{x}_n) \log w_k \\ &= \min_{\mathbf{w} \in \Delta^{K-1}} - \sum_{k=1}^K \underbrace{\frac{\sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n)}{N}}_{p_k^{(t)}} \log w_k = \min_{\mathbf{w}} \text{KL}(\mathbf{p}^{(t)} \parallel \mathbf{w}) \quad \Rightarrow \quad \mathbf{w}^{(t+1)} = \mathbf{p}^{(t)} \quad (8) \end{aligned}$$

How to update μ_k and Σ_k ?

An Optimization Viewpoint of The EM Algorithm

The optimization problem in the t -th step:

$$\min_{\theta} Q(\theta; \theta^{(t)}) = \min_{\theta} - \sum_{n=1}^N \sum_{k=1}^K p^{(t)}(k|\mathbf{x}_n) \log \frac{w_k p(\mathbf{x}_n; \mu_k, \Sigma_k)}{p^{(t)}(k|\mathbf{x}_n)}. \quad (7)$$

Then, in the M-step:

$$\begin{aligned} & \min_{\mathbf{w} \in \Delta^{K-1}} - \sum_{n=1}^N \sum_{k=1}^K p^{(t)}(k|\mathbf{x}_n) \log w_k \\ &= \min_{\mathbf{w} \in \Delta^{K-1}} - \sum_{k=1}^K \underbrace{\frac{\sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n)}{N}}_{p_k^{(t)}} \log w_k = \min_{\mathbf{w}} \text{KL}(\mathbf{p}^{(t)} \parallel \mathbf{w}) \quad \Rightarrow \quad \mathbf{w}^{(t+1)} = \mathbf{p}^{(t)} \quad (8) \end{aligned}$$

How to update μ_k and Σ_k ?

$$\min_{\mu_k} - \sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n) \log p(\mathbf{x}_n; \mu_k, \Sigma_k^{(t)})$$

An Optimization Viewpoint of The EM Algorithm

The optimization problem in the t -th step:

$$\min_{\theta} Q(\theta; \theta^{(t)}) = \min_{\theta} - \sum_{n=1}^N \sum_{k=1}^K p^{(t)}(k|\mathbf{x}_n) \log \frac{w_k p(\mathbf{x}_n; \mu_k, \Sigma_k)}{p^{(t)}(k|\mathbf{x}_n)}. \quad (7)$$

Then, in the M-step:

$$\begin{aligned} & \min_{\mathbf{w} \in \Delta^{K-1}} - \sum_{n=1}^N \sum_{k=1}^K p^{(t)}(k|\mathbf{x}_n) \log w_k \\ &= \min_{\mathbf{w} \in \Delta^{K-1}} - \sum_{k=1}^K \underbrace{\frac{\sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n)}{N}}_{p_k^{(t)}} \log w_k = \min_{\mathbf{w}} \text{KL}(\mathbf{p}^{(t)} \parallel \mathbf{w}) \quad \Rightarrow \quad \mathbf{w}^{(t+1)} = \mathbf{p}^{(t)} \end{aligned} \quad (8)$$

How to update μ_k and Σ_k ?

$$\begin{aligned} & \min_{\mu_k} - \sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n) \log p(\mathbf{x}_n; \mu_k, \Sigma_k^{(t)}) \\ & \frac{\partial Q}{\partial \mu_k} = 0 \quad \Rightarrow \quad \mu_k^{(t+1)} = \frac{\sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n) \mathbf{x}_n}{\sum_{n=1}^N p^{(t)}(k|\mathbf{x}_n)} \quad (\text{Derive It.}) \end{aligned} \quad (9)$$

Revisit Kmeans from The Viewpoint of EM Algorithm

Kmeans learns K centroids, which also provides us with a **generative model**.

- 1 Initialize K centroids randomly from the observed data, e.g., $\{\mathbf{c}_k\}_{k=1}^K$.

Revisit Kmeans from The Viewpoint of EM Algorithm

Kmeans learns K centroids, which also provides us with a **generative model**.

- 1 Initialize K centroids randomly from the observed data, e.g., $\{\mathbf{c}_k\}_{k=1}^K$.
- 2 Repeat the following steps till converge
 - a Assign each data to the nearest cluster: $\forall \mathbf{x}_n$

$$\mathbf{x}_n \in \mathcal{C}_k, \quad \text{if } k = \arg \min_{k \in \{1, \dots, K\}} d(\mathbf{x}_n, \mathbf{c}_k).$$

Revisit Kmeans from The Viewpoint of EM Algorithm

Kmeans learns K centroids, which also provides us with a **generative model**.

- 1 Initialize K centroids randomly from the observed data, e.g., $\{\mathbf{c}_k\}_{k=1}^K$.
- 2 Repeat the following steps till converge
 - a Assign each data to the nearest cluster: $\forall \mathbf{x}_n$

$$\mathbf{x}_n \in \mathcal{C}_k, \quad \text{if } k = \arg \min_{k \in \{1, \dots, K\}} d(\mathbf{x}_n, \mathbf{c}_k). \quad (10)$$

- b Update the centriods:

$$\mathbf{c}_k = \frac{1}{|\mathcal{C}_k|} \sum_{\mathbf{x}_n \in \mathcal{C}_k} \mathbf{x}_n.$$

Revisit Kmeans from The Viewpoint of EM Algorithm

Kmeans learns K centroids, which also provides us with a **generative model**.

- 1 Initialize K centroids randomly from the observed data, e.g., $\{\mathbf{c}_k\}_{k=1}^K$.
- 2 Repeat the following steps till converge
 - a Assign each data to the nearest cluster: $\forall \mathbf{x}_n$

$$\mathbf{x}_n \in \mathcal{C}_k, \quad \text{if } k = \arg \min_{k \in \{1, \dots, K\}} d(\mathbf{x}_n, \mathbf{c}_k). \quad (10)$$

- b Update the centriods:

$$\mathbf{c}_k = \frac{1}{|\mathcal{C}_k|} \sum_{\mathbf{x}_n \in \mathcal{C}_k} \mathbf{x}_n. \quad (11)$$

- $w_k = p(k) = \frac{|\mathcal{C}_k|}{N}$, and

$$p(\mathbf{x}; \mathbf{c}_k) = \begin{cases} \frac{\Gamma(\frac{D}{2}+1)}{\pi^{D/2} r}, & \|\mathbf{x} - \mathbf{c}_k\|_2 \leq r, \\ 0, & \text{Otherwise.} \end{cases}$$

Revisit Kmeans from The Viewpoint of EM Algorithm

Kmeans learns K centroids, which also provides us with a **generative model**.

- 1 Initialize K centroids randomly from the observed data, e.g., $\{\mathbf{c}_k\}_{k=1}^K$.
- 2 Repeat the following steps till converge
 - a Assign each data to the nearest cluster: $\forall \mathbf{x}_n$

$$\mathbf{x}_n \in \mathcal{C}_k, \quad \text{if } k = \arg \min_{k \in \{1, \dots, K\}} d(\mathbf{x}_n, \mathbf{c}_k). \quad (10)$$

- b Update the centriods:

$$\mathbf{c}_k = \frac{1}{|\mathcal{C}_k|} \sum_{\mathbf{x}_n \in \mathcal{C}_k} \mathbf{x}_n. \quad (11)$$

- $w_k = p(k) = \frac{|\mathcal{C}_k|}{N}$, and

$$p(\mathbf{x}; \mathbf{c}_k) = \begin{cases} \frac{\Gamma(\frac{D}{2}+1)}{\pi^{D/2} r}, & \|\mathbf{x} - \mathbf{c}_k\|_2 \leq r, \\ 0, & \text{Otherwise.} \end{cases} \quad (12)$$

Recall nonparametric kernel functions:)

Revisit Kmeans from The Viewpoint of EM Algorithm

Original MLE: Denote $\theta = \{w_k, \mathbf{c}_k, r_k\}_{k=1}^K$, the MLE of GMM corresponds to

$$\min_{\theta} \sum_{n=1}^N -\log p(\mathbf{x}_n; \theta) = \min_{\theta} \underbrace{\sum_{n=1}^N -\log \left(\sum_{k=1}^K w_k p(\mathbf{x}_n; \mathbf{c}_k, r_k) \right)}_{L(\theta)}.$$

Revisit Kmeans from The Viewpoint of EM Algorithm

Original MLE: Denote $\theta = \{w_k, \mathbf{c}_k, r_k\}_{k=1}^K$, the MLE of GMM corresponds to

$$\min_{\theta} \sum_{n=1}^N -\log p(\mathbf{x}_n; \theta) = \min_{\theta} \underbrace{\sum_{n=1}^N -\log \left(\sum_{k=1}^K w_k p(\mathbf{x}_n; \mathbf{c}_k, r_k) \right)}_{L(\theta)}. \quad (13)$$

EM Algorithm:

► **E-step:** $\exists r > 0$ for \mathbf{x}_n and $k = 1, \dots, K$

$$p^{(t)}(k|\mathbf{x}_n) = \frac{w_k^{(t)} p(\mathbf{x}_n; \mathbf{c}_k^{(t)}, r_k^{(t)})}{\sum_{i=1}^K w_i^{(t)} p(\mathbf{x}_n; \mathbf{c}_i^{(t)}, r_i^{(t)})} = \begin{cases} 1 & \|\mathbf{x}_n - \mathbf{c}_k^{(t)}\|_2 \leq r \\ 0 & \text{Otherwise} \end{cases}$$

Revisit Kmeans from The Viewpoint of EM Algorithm

Original MLE: Denote $\theta = \{w_k, \mathbf{c}_k, r_k\}_{k=1}^K$, the MLE of GMM corresponds to

$$\min_{\theta} \sum_{n=1}^N -\log p(\mathbf{x}_n; \theta) = \min_{\theta} \underbrace{\sum_{n=1}^N -\log \left(\sum_{k=1}^K w_k p(\mathbf{x}_n; \mathbf{c}_k, r_k) \right)}_{L(\theta)}. \quad (13)$$

EM Algorithm:

► **E-step:** $\exists r > 0$ for \mathbf{x}_n and $k = 1, \dots, K$

$$p^{(t)}(k|\mathbf{x}_n) = \frac{w_k^{(t)} p(\mathbf{x}_n; \mathbf{c}_k^{(t)}, r_k^{(t)})}{\sum_{i=1}^K w_i^{(t)} p(\mathbf{x}_n; \mathbf{c}_i^{(t)}, r_i^{(t)})} = \begin{cases} 1 & \|\mathbf{x}_n - \mathbf{c}_k^{(t)}\|_2 \leq r \\ 0 & \text{Otherwise} \end{cases} \quad (14)$$

► **M-step:**

$$\min_{\theta} Q(\theta; \theta^{(t)}) = \min_{\theta} - \sum_{n=1}^N \sum_{k=1}^K p^{(t)}(k|\mathbf{x}_n) \log \frac{w_k p(\mathbf{x}_n; \mathbf{c}_k)}{p^{(t)}(k|\mathbf{x}_n)}. \quad (15)$$

A Generalized Formulation of The EM Algorithm

- ▶ Essentially, the EM algorithm is used to find **local** maximum likelihood parameters of a statistical model that involve **unobserved latent variables**

A Generalized Formulation of The EM Algorithm

- ▶ Essentially, the EM algorithm is used to find **local** maximum likelihood parameters of a statistical model that involve **unobserved latent variables**
- ▶ In GMMs, the unobserved latent variable, its the cluster ID of each sample.

A Generalized Formulation of The EM Algorithm

- ▶ Essentially, the EM algorithm is used to find **local** maximum likelihood parameters of a statistical model that involve **unobserved latent variables**
- ▶ In GMMs, the unobserved latent variable, its the cluster ID of each sample.
- ▶ Denote \mathbf{X} as observed data, and \mathbf{Z} as latent variables of the data, θ as the model parameters.

A Generalized Formulation of The EM Algorithm

- ▶ Essentially, the EM algorithm is used to find **local** maximum likelihood parameters of a statistical model that involve **unobserved latent variables**
- ▶ In GMMs, the unobserved latent variable, its the cluster ID of each sample.
- ▶ Denote \mathbf{X} as observed data, and \mathbf{Z} as latent variables of the data, θ as the model parameters.
- ▶ $L(\theta; \mathbf{X}, \mathbf{Z}) = p(\mathbf{X}, \mathbf{Z} | \theta)$: a likelihood function.

A Generalized Formulation of The EM Algorithm

- ▶ Essentially, the EM algorithm is used to find **local** maximum likelihood parameters of a statistical model that involve **unobserved latent variables**
- ▶ In GMMs, the unobserved latent variable, its the cluster ID of each sample.
- ▶ Denote \mathbf{X} as observed data, and \mathbf{Z} as latent variables of the data, θ as the model parameters.
- ▶ $L(\theta; \mathbf{X}, \mathbf{Z}) = p(\mathbf{X}, \mathbf{Z} | \theta)$: a likelihood function.
- ▶ The MLE of θ given observed \mathbf{X} : maximizing the marginal probability of \mathbf{X} :

$$L(\theta; \mathbf{X}) = p(\mathbf{X} | \theta)$$

A Generalized Formulation of The EM Algorithm

- ▶ Essentially, the EM algorithm is used to find **local** maximum likelihood parameters of a statistical model that involve **unobserved latent variables**
- ▶ In GMMs, the unobserved latent variable, its the cluster ID of each sample.
- ▶ Denote \mathbf{X} as observed data, and \mathbf{Z} as latent variables of the data, θ as the model parameters.
- ▶ $L(\theta; \mathbf{X}, \mathbf{Z}) = p(\mathbf{X}, \mathbf{Z} | \theta)$: a likelihood function.
- ▶ The MLE of θ given observed \mathbf{X} : maximizing the marginal probability of \mathbf{X} :

$$L(\theta; \mathbf{X}) = p(\mathbf{X} | \theta) = \int p(\mathbf{X}, \mathbf{Z} | \theta) d\mathbf{Z}$$

A Generalized Formulation of The EM Algorithm

- ▶ Essentially, the EM algorithm is used to find **local** maximum likelihood parameters of a statistical model that involve **unobserved latent variables**
- ▶ In GMMs, the unobserved latent variable, its the cluster ID of each sample.
- ▶ Denote \mathbf{X} as observed data, and \mathbf{Z} as latent variables of the data, θ as the model parameters.
- ▶ $L(\theta; \mathbf{X}, \mathbf{Z}) = p(\mathbf{X}, \mathbf{Z} | \theta)$: a likelihood function.
- ▶ The MLE of θ given observed \mathbf{X} : maximizing the marginal probability of \mathbf{X} :

$$\begin{aligned} L(\theta; \mathbf{X}) &= p(\mathbf{X} | \theta) = \int p(\mathbf{X}, \mathbf{Z} | \theta) d\mathbf{Z} \\ &= \int \underbrace{p(\mathbf{X} | \mathbf{Z}, \theta)}_{L(\theta; \mathbf{X}, \mathbf{Z})} p(\mathbf{Z} | \theta) d\mathbf{Z} \end{aligned}$$

A Generalized Formulation of The EM Algorithm

- ▶ Essentially, the EM algorithm is used to find **local** maximum likelihood parameters of a statistical model that involve **unobserved latent variables**
- ▶ In GMMs, the unobserved latent variable, its the cluster ID of each sample.
- ▶ Denote \mathbf{X} as observed data, and \mathbf{Z} as latent variables of the data, θ as the model parameters.
- ▶ $L(\theta; \mathbf{X}, \mathbf{Z}) = p(\mathbf{X}, \mathbf{Z} | \theta)$: a likelihood function.
- ▶ The MLE of θ given observed \mathbf{X} : maximizing the marginal probability of \mathbf{X} :

$$\begin{aligned} L(\theta; \mathbf{X}) &= p(\mathbf{X} | \theta) = \int p(\mathbf{X}, \mathbf{Z} | \theta) d\mathbf{Z} \\ &= \int \underbrace{p(\mathbf{X} | \mathbf{Z}, \theta)}_{L(\theta; \mathbf{X}, \mathbf{Z})} p(\mathbf{Z} | \theta) d\mathbf{Z} = \mathbb{E}_{\mathbf{Z} \sim p(\mathbf{Z} | \theta)} [L(\theta; \mathbf{X}, \mathbf{Z})] \end{aligned} \tag{16}$$

A Generalized Formulation of The EM Algorithm

The objective function of EM:

$$\min_{\theta} -\mathbb{E}_{\mathbf{Z} \sim p(\mathbf{Z}|\theta)} [L(\theta; \mathbf{X}, \mathbf{Z})]$$

A Generalized Formulation of The EM Algorithm

The objective function of EM:

$$\min_{\theta} -\mathbb{E}_{\mathbf{Z} \sim p(\mathbf{Z}|\theta)} [L(\theta; \mathbf{X}, \mathbf{Z})] \quad (17)$$

E-step: Construct the expected value of the log-likelihood function of θ given current conditional distribution of \mathbf{Z} given \mathbf{X} and current $\theta^{(t)}$:

$$Q(\theta; \theta^{(t)}) := -\mathbb{E}_{\mathbf{Z} \sim p(\mathbf{Z}|\theta^{(t)})} [L(\theta; \mathbf{X}, \mathbf{Z})].$$

A Generalized Formulation of The EM Algorithm

The objective function of EM:

$$\min_{\theta} -\mathbb{E}_{\mathbf{Z} \sim p(\mathbf{Z}|\theta)} [L(\theta; \mathbf{X}, \mathbf{Z})] \quad (17)$$

E-step: Construct the expected value of the log-likelihood function of θ given current conditional distribution of \mathbf{Z} given \mathbf{X} and current $\theta^{(t)}$:

$$Q(\theta; \theta^{(t)}) := -\mathbb{E}_{\mathbf{Z} \sim p(\mathbf{Z}|\theta^{(t)})} [L(\theta; \mathbf{X}, \mathbf{Z})]. \quad (18)$$

M-step: Maximize the expectation:

$$\theta^{(t+1)} = \arg \max_{\theta} Q(\theta; \theta^{(t)}).$$

A Generalized Formulation of The EM Algorithm

The objective function of EM:

$$\min_{\theta} -\mathbb{E}_{\mathbf{Z} \sim p(\mathbf{Z}|\theta)} [L(\theta; \mathbf{X}, \mathbf{Z})] \quad (17)$$

E-step: Construct the expected value of the log-likelihood function of θ given current conditional distribution of \mathbf{Z} given \mathbf{X} and current $\theta^{(t)}$:

$$Q(\theta; \theta^{(t)}) := -\mathbb{E}_{\mathbf{Z} \sim p(\mathbf{Z}|\theta^{(t)})} [L(\theta; \mathbf{X}, \mathbf{Z})]. \quad (18)$$

M-step: Maximize the expectation:

$$\theta^{(t+1)} = \arg \max_{\theta} Q(\theta; \theta^{(t)}). \quad (19)$$

Question: In GMM, what is \mathbf{Z} ? and what is $p(\mathbf{Z}|\theta^{(t)})$?

In Summary

- ▶ Generative modeling and Gaussian mixture models
- ▶ EM algorithm
- ▶ Revisit K-means from a statistical viewpoint

Next...

- ▶ A Bayesian viewpoint of GMMs
- ▶ Kernel density estimation
- ▶ Mean-shift algorithm