

Clustering, Gaussian mixture model and EM

MATH-412 - Statistical Machine Learning

K-means

Key assumption: Data composed of K “roundish” clusters of similar sizes with centroids (μ_1, \dots, μ_K) .

Problem can be formulated as:
$$\min_{\mu_1, \dots, \mu_K} \frac{1}{n} \sum_{i=1}^n \min_k \|\mathbf{x}_i - \mu_k\|^2.$$

Difficult (NP-hard) nonconvex problem.

K-means algorithm

- 1 Draw centroids at random (or use the “ k -means++” initialization)
- 2 Assign each point to the closest centroid

$$C_k \leftarrow \{i \mid \|\mathbf{x}_i - \mu_k\|^2 = \min_j \|\mathbf{x}_i - \mu_j\|^2\}$$

- 3 Recompute centroid as center of mass of the cluster

$$\mu_k \leftarrow \frac{1}{|C_k|} \sum_{i \in C_k} \mathbf{x}_i$$

- 4 Go to 2

K-means properties

Three remarks:

- K-means is a greedy algorithm
- It can be shown that K-means converges in a finite number of steps.
- The algorithm however typically gets stuck in local minima and in practice it is necessary to try several restarts of the algorithms with different initialization to have chances to obtain a better solution.
- Will typically fail if the clusters are not round or of too different sizes.

The EM algorithm for the Gaussian mixture model

Jensen's Inequality

Consider a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$

- ① if f is **convex** and if X is a random variable (with $\mathbb{E}[X] \in \mathbb{R}$), then

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$$

- ② When f is **strictly convex**, we have equality in the previous inequality if and only if X is constant almost surely.

The Kullback-Leibler divergence

Definition Let \mathcal{X} a finite (or countable) state space and p and q two distributions on \mathcal{X}

$$KL(p \parallel q) = \sum_x p(x) \log \frac{p(x)}{q(x)} = \mathbb{E}_{X \sim p} \left[\log \frac{p(X)}{q(X)} \right]$$

Entropy: $H(p) = -\sum_x p(x) \log p(x) \geq 0$

$$\text{So} \quad KL(p \parallel q) = \mathbb{E}_{X \sim p} [-\log q(X)] - H(p).$$

Property: $\forall p, q, KL(p \parallel q) \geq 0$ (could be infinite)

Proof:

$$KL(p \parallel q) = \sum_x p(x) \log \frac{p(x)}{q(x)} = \mathbb{E}_{X \sim p} \left[\log \frac{p(X)}{q(X)} \right] = \mathbb{E}_{X \sim p} \left[-\log \frac{q(X)}{p(X)} \right] =$$

The function $f(y) = -\log y$ is strictly convex so

$$\geq -\log \mathbb{E}_{X \sim p} \left[\frac{q(X)}{p(X)} \right] = -\log \sum_x p(x) \frac{q(X)}{p(X)} = -\log \sum_x q(x) = 0$$

with equality if and only if $p = q$ almost surely

Differential KL and entropies

Let P and Q two probability distributions with densities p and q with respect to a measure μ . Then, we can define

$$KL(p \parallel q) = \int_x \left(\log \frac{p(x)}{q(x)} \right) p(x) d\mu(x) = \mathbb{E}_{X \sim P} \left[\log \frac{p(X)}{q(X)} \right]$$

Differential entropy

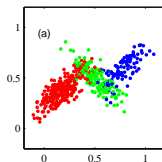
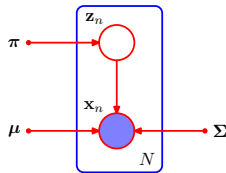
$$\mathcal{H}(p) = - \int_x p(x) \log p(x) d\mu(x)$$

Caveats: the differential entropy is dangerous

- $\mathcal{H}(p) \not\geq 0$
- \mathcal{H} depends on the choice of μ ...!

Gaussian mixture model

- K components
- $\mathbf{z} = (z_1, \dots, z_K)^\top \in \{0, 1\}^K$ indicator variable (one hot encoding)
- $\mathbf{z} \sim \mathcal{M}(1, (\pi_1, \dots, \pi_K))$
- $p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$, so $p(e_k) = \pi_k$
- $p(\mathbf{x}|\mathbf{z}; (\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)_k) = \sum_{k=1}^K z_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
- $p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
- Estimation: $\underset{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k}{\operatorname{argmax}} \log \left[\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right]$



Applying maximum likelihood to the Gaussian mixture

Let $\mathcal{Z} = \{z \in \{0, 1\}^K \mid \sum_{k=1}^K z_k = 1\}$

$$p(\mathbf{x}) = \sum_{z \in \mathcal{Z}} p(\mathbf{x}, z) = \sum_{z \in \mathcal{Z}} \prod_{k=1}^K \left[\pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right]^{z_k} = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Issue

- The marginal log-likelihood $\tilde{\ell}(\theta) = \sum_i \log(p(\mathbf{x}^{(i)}))$ with $\theta = (\boldsymbol{\pi}, (\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)_{1 \leq k \leq K})$ is now complicated
- $\log p(\mathbf{x}, z) = \sum_k \log p(\mathbf{x}, e_k) 1(z = e_k) = \sum_k z_k \log(\pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))$
- By contrast the complete log-likelihood has a rather simple form:

$$\tilde{\ell}(\theta) = \sum_{i=1}^M \log p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}) = \sum_{i, k} z_k^{(i)} \log \mathcal{N}(\mathbf{x}^{(i)}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) + \sum_{i, k} z_k^{(i)} \log(\pi_k),$$

Principle of the Expectation-Maximization Algorithm

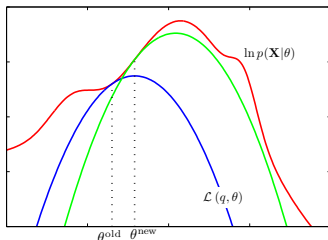
$$\begin{aligned}\log p(\mathbf{x}; \boldsymbol{\theta}) &= \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = \log \sum_{\mathbf{z}} q(\mathbf{z}) \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{q(\mathbf{z})} \\ &\geq \sum_{\mathbf{z}} q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{q(\mathbf{z})} = \mathbb{E}_q[\log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})] + H(q) =: \mathcal{L}(q, \boldsymbol{\theta})\end{aligned}$$

- This shows that $\mathcal{L}(q, \boldsymbol{\theta}) \leq \log p(\mathbf{x}; \boldsymbol{\theta})$
- $\boldsymbol{\theta} \mapsto \mathcal{L}(q, \boldsymbol{\theta})$ is *often*^a concave or easy to maximize
- It is possible to show that

$$\mathcal{L}(q, \boldsymbol{\theta}) = \log p(\mathbf{x}; \boldsymbol{\theta}) - KL(q \| p(\cdot | \mathbf{x}; \boldsymbol{\theta}))$$

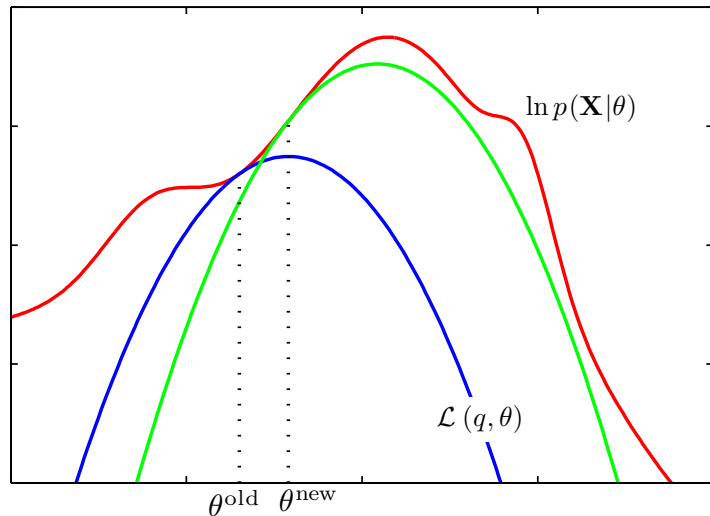
So that if we set $q(\mathbf{z}) = p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{(t)})$ then

$$\mathcal{L}(q, \boldsymbol{\theta}^{(t)}) = \log p(\mathbf{x}; \boldsymbol{\theta}^{(t)}).$$



^aIt is concave if $p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})_{\boldsymbol{\theta} \in \Theta}$ is a canonical exponential family, with $\boldsymbol{\theta}$ its natural parameter, i.e., $p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = h(\mathbf{x}, \mathbf{z}) \exp(\phi(\mathbf{x}, \mathbf{z})^\top \boldsymbol{\theta} - A(\boldsymbol{\theta}))$.

A graphical idea of the EM algorithm



Expectation Maximization algorithm

Initialize $\theta = \theta_0$

WHILE (Not converged)

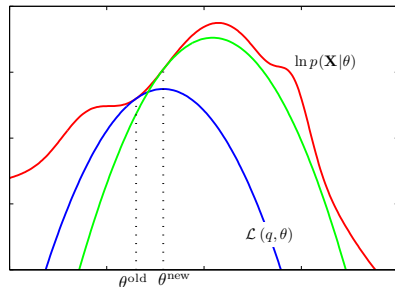
Expectation step

- 1 $q(z) = p(z \mid \mathbf{x}; \theta^{(t-1)})$
- 2 $\mathcal{L}(q, \theta) = \mathbb{E}_q[\log p(\mathbf{x}, \mathbf{Z}; \theta)] + H(q)$

Maximization step

- 1 $\theta^{(t)} = \operatorname{argmax}_{\theta} \mathbb{E}_q[\log p(\mathbf{x}, \mathbf{Z}; \theta)]$

ENDWHILE



$$\theta^{\text{old}} = \theta^{(t-1)}$$

$$\theta^{\text{new}} = \theta^{(t)}$$

Expected complete log-likelihood

With the notation: $q_{ik}^{(t)} = \mathbb{P}_{q_i^{(t)}}(z_k^{(i)} = 1) = \mathbb{E}_{q_i^{(t)}}[z_k^{(i)}]$, we have

$$\begin{aligned}\mathbb{E}_{q^{(t)}}[\tilde{\ell}(\theta)] &= \mathbb{E}_{q^{(t)}}[\log p(\mathbf{X}, \mathbf{Z}; \theta)] \\&= \mathbb{E}_{q^{(t)}}\left[\sum_{i=1}^M \log p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}; \theta)\right] \\&= \mathbb{E}_{q^{(t)}}\left[\sum_{i,k} z_k^{(i)} \log \mathcal{N}(\mathbf{x}^{(i)}, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) + \sum_{i,k} z_k^{(i)} \log(\pi_k)\right] \\&= \sum_{i,k} \mathbb{E}_{q_i^{(t)}}[z_k^{(i)}] \log \mathcal{N}(\mathbf{x}^{(i)}, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) + \sum_{i,k} \mathbb{E}_{q_i^{(t)}}[z_k^{(i)}] \log(\pi_k) \\&= \sum_{i,k} q_{ik}^{(t)} \log \mathcal{N}(\mathbf{x}^{(i)}, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) + \sum_{i,k} q_{ik}^{(t)} \log(\pi_k)\end{aligned}$$

Expectation step for the Gaussian mixture

We computed previously $q_i^{(t)}(\mathbf{z}^{(i)})$, which is a multinomial distribution defined by

$$q_i^{(t)}(\mathbf{z}^{(i)}) = p(\mathbf{z}^{(i)} | \mathbf{x}^{(i)}; \theta^{(t-1)})$$

Abusing notation we will denote $(q_{i1}^{(t)}, \dots, q_{iK}^{(t)})$ the corresponding vector of probabilities defined by

$$q_{ik}^{(t)} = \mathbb{P}_{q_i^{(t)}}(z_k^{(i)} = 1) = \mathbb{E}_{q_i^{(t)}}[z_k^{(i)}]$$

$$q_{ik}^{(t)} = p(z_k^{(i)} = 1 | \mathbf{x}^{(i)}; \theta^{(t-1)}) = \frac{\pi_k^{(t-1)} \mathcal{N}(\mathbf{x}^{(i)}, \boldsymbol{\mu}_k^{(t-1)}, \boldsymbol{\Sigma}_k^{(t-1)})}{\sum_{j=1}^K \pi_j^{(t-1)} \mathcal{N}(\mathbf{x}^{(i)}, \boldsymbol{\mu}_j^{(t-1)}, \boldsymbol{\Sigma}_j^{(t-1)})}$$

Maximization step for the Gaussian mixture

$$(\boldsymbol{\pi}^t, (\boldsymbol{\mu}_k^{(t)}, \boldsymbol{\Sigma}_k^{(t)})_{1 \leq k \leq K}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \mathbb{E}_{q^{(t)}} [\tilde{\ell}(\boldsymbol{\theta})]$$

This yields the updates:

$$\boldsymbol{\mu}_k^{(t)} = \frac{\sum_i \mathbf{x}^{(i)} q_{ik}^{(t)}}{\sum_i q_{ik}^{(t)}}, \quad \boldsymbol{\Sigma}_k^{(t)} = \frac{\sum_i (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k^{(t)}) (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k^{(t)})^\top q_{ik}^{(t)}}{\sum_i q_{ik}^{(t)}}$$

and

$$\pi_k^{(t)} = \frac{\sum_i q_{ik}^{(t)}}{\sum_{i,k'} q_{ik'}^{(t)}}$$

Final EM algorithm for the Gaussian mixture model

Initialize $\theta = \theta_0$

WHILE (Not converged)

Expectation step

$$q_{ik}^{(t)} \leftarrow \frac{\pi_k^{(t-1)} \mathcal{N}(\mathbf{x}^{(i)}, \boldsymbol{\mu}_k^{(t-1)}, \boldsymbol{\Sigma}_k^{(t-1)})}{\sum_{j=1}^K \pi_j^{(t-1)} \mathcal{N}(\mathbf{x}^{(i)}, \boldsymbol{\mu}_j^{(t-1)}, \boldsymbol{\Sigma}_j^{(t-1)})}$$

Maximization step

$$\boldsymbol{\mu}_k^{(t)} = \frac{\sum_i \mathbf{x}^{(i)} q_{ik}^{(t)}}{\sum_i q_{ik}^{(t)}}, \quad \boldsymbol{\Sigma}_k^{(t)} = \frac{\sum_i (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k^{(t)}) (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k^{(t)})^\top q_{ik}^{(t)}}{\sum_i q_{ik}^{(t)}}$$

and

$$\pi_k^{(t)} = \frac{\sum_i q_{ik}^{(t)}}{\sum_{i,k'} q_{ik'}^{(t)}}$$

ENDWHILE