

EM for Naive Bayes and Gaussian Mixture Models, k -Means Clustering

Karl Stratos

June 27, 2018

EM Template

Input: model $P_{\Phi}(\mathbf{x}, \mathbf{z})$, unlabeled data $U = \{\mathbf{x}^{(i)}\}_{i=1}^n$, T

Output: local maximizer of $L_U(\Phi) := \sum_{i=1}^n \log P_{\Phi}(\mathbf{x}^{(i)})$

1. Initialize parameters $\Phi^{(0)}$.

2. For $t = 0 \dots T - 1$,

$$\Phi^{(t+1)} \leftarrow \arg \max_{\Phi} \sum_{i=1}^n \sum_{z=1}^m P_{\Phi^{(t)}}(z | \mathbf{x}^{(i)}) \times \log P_{\Phi}(\mathbf{x}^{(i)}, z)$$

3. Return $\Phi^{(T)}$.

EM Template

Input: model $P_{\Phi}(\mathbf{x}, \mathbf{z})$, unlabeled data $U = \{\mathbf{x}^{(i)}\}_{i=1}^n$, T

Output: local maximizer of $L_U(\Phi) := \sum_{i=1}^n \log P_{\Phi}(\mathbf{x}^{(i)})$

1. Initialize parameters $\Phi^{(0)}$.

2. For $t = 0 \dots T - 1$,

$$\Phi^{(t+1)} \leftarrow \arg \max_{\Phi} \sum_{i=1}^n \sum_{\mathbf{z}=1}^m P_{\Phi^{(t)}}(\mathbf{z}|\mathbf{x}^{(i)}) \times \log P_{\Phi}(\mathbf{x}^{(i)}, \mathbf{z})$$

3. Return $\Phi^{(T)}$.

See yesterday's lecture for how this is derived by alternating maximization of the ELBO $\mathcal{L}(\Phi, \Psi) \leq L_U(\Phi)$ where Ψ defines an auxiliary posterior $P_{\Psi}(y|\mathbf{x})$.

Overview

EM for Naive Bayes

Maximum Likelihood Estimation with Labeled Data

Maximum Likelihood Estimation with Unlabeled Data

EM for Gaussian Mixture Models

Maximum Likelihood Estimation with Labeled Data

Maximum Likelihood Estimation with Unlabeled Data

k -Means Clustering

Naive Bayes: Definition

A **naive Bayes (NB)** model with m labels and d binary-valued feature types has $m + 2dm$ parameters, denoted by Φ :

- ▶ $q(z) \geq 0$ for each $z \in \{1 \dots m\}$ such that

$$\sum_z q(z) = 1$$

- ▶ $q(0|z, j) \geq 0$ and $q(1|z, j) \geq 0$ such that

$$q(0|z, j) + q(1|z, j) = 1$$

for each $j \in \{1 \dots d\}$ and $z \in \{1 \dots m\}$

Naive Bayes: Definition

A **naive Bayes (NB)** model with m labels and d binary-valued feature types has $m + 2dm$ parameters, denoted by Φ :

- ▶ $q(z) \geq 0$ for each $z \in \{1 \dots m\}$ such that

$$\sum_z q(z) = 1$$

- ▶ $q(0|z, j) \geq 0$ and $q(1|z, j) \geq 0$ such that

$$q(0|z, j) + q(1|z, j) = 1$$

for each $j \in \{1 \dots d\}$ and $z \in \{1 \dots m\}$

Φ defines a joint distribution over $\mathbf{x} = (x_1 \dots x_d) \in \{0, 1\}^d$ and $z \in \{1 \dots m\}$ by

$$P_{\Phi}(\mathbf{x}, z) := q(z) \prod_{j=1}^d q(x_j|z, j)$$

Log Likelihood of Labeled Data

If $S = \{(\mathbf{x}^{(i)}, z^{(i)})\}_{i=1}^n$ is a set of n iid labeled samples, the log likelihood of S under Φ is

$$L_S(\Phi) = \sum_{i=1}^n \log P_{\Phi}(\mathbf{x}^{(i)}, z^{(i)})$$

Log Likelihood of Labeled Data

If $S = \{(\mathbf{x}^{(i)}, z^{(i)})\}_{i=1}^n$ is a set of n iid labeled samples, the log likelihood of S under Φ is

$$\begin{aligned} L_S(\Phi) &= \sum_{i=1}^n \log P_{\Phi}(\mathbf{x}^{(i)}, z^{(i)}) \\ &= \sum_{i=1}^n \log q(z^{(i)}) + \sum_{j=1}^m \log q(x_j^{(i)} | z, j) \end{aligned}$$

Log Likelihood of Labeled Data

If $S = \{(\mathbf{x}^{(i)}, z^{(i)})\}_{i=1}^n$ is a set of n iid labeled samples, the log likelihood of S under Φ is

$$\begin{aligned}L_S(\Phi) &= \sum_{i=1}^n \log P_{\Phi}(\mathbf{x}^{(i)}, z^{(i)}) \\&= \sum_{i=1}^n \log q(z^{(i)}) + \sum_{j=1}^m \log q(x_j^{(i)} | z, j) \\&= \sum_{z=1}^m \mathbf{count}(z) \log q(z) + \sum_{z=1}^m \sum_{j=1}^m \sum_{x \in \{0,1\}} \mathbf{count}(z, j, x) \log q(x | z, j)\end{aligned}$$

where

$$\mathbf{count}(z) := \sum_{\substack{i=1: \\ z^{(i)}=z}}^n 1 \qquad \mathbf{count}(z, j, x) := \sum_{\substack{i=1: \\ z^{(i)}=z \\ x_j^{(i)}=x}}^n 1$$

MLE with Labeled Data

What are the parameter values $q(z)$ and $q(x|z, j)$ that maximize

$$\sum_{z=1}^m \mathbf{count}(z) \log q(z) + \sum_{z=1}^m \sum_{j=1}^m \sum_{x \in \{0,1\}} \mathbf{count}(z, j, x) \log q(x|z, j)$$

under the constraints that they are nonnegative, $\sum_z q(z) = 1$, and $\sum_x q(x|z, j) = 1$?

MLE with Labeled Data

What are the parameter values $q(z)$ and $q(x|z, j)$ that maximize

$$\sum_{z=1}^m \mathbf{count}(z) \log q(z) + \sum_{z=1}^m \sum_{j=1}^m \sum_{x \in \{0,1\}} \mathbf{count}(z, j, x) \log q(x|z, j)$$

under the constraints that they are nonnegative, $\sum_z q(z) = 1$, and $\sum_x q(x|z, j) = 1$?

Answer: See the lemma in yesterday's lecture for why.

$$q(z) = \frac{\mathbf{count}(z)}{n}$$
$$q(x|z, j) = \frac{\mathbf{count}(z, j, x)}{\mathbf{count}(z, j, 0) + \mathbf{count}(z, j, 1)}$$

Overview

EM for Naive Bayes

Maximum Likelihood Estimation with Labeled Data

Maximum Likelihood Estimation with Unlabeled Data

EM for Gaussian Mixture Models

Maximum Likelihood Estimation with Labeled Data

Maximum Likelihood Estimation with Unlabeled Data

k -Means Clustering

Log Likelihood of Unlabeled Data

If $U = \{\mathbf{x}^{(i)}\}_{i=1}^n$ is a set of n iid unlabeled samples, the log likelihood of U under Φ is

$$L_U(\Phi) = \sum_{i=1}^n \log P_{\Phi}(\mathbf{x}^{(i)})$$

Log Likelihood of Unlabeled Data

If $U = \{\mathbf{x}^{(i)}\}_{i=1}^n$ is a set of n iid unlabeled samples, the log likelihood of U under Φ is

$$\begin{aligned}L_U(\Phi) &= \sum_{i=1}^n \log P_{\Phi}(\mathbf{x}^{(i)}) \\ &= \sum_{i=1}^n \log \left(\sum_{z=1}^m P_{\Phi}(\mathbf{x}^{(i)}, z) \right)\end{aligned}$$

Log Likelihood of Unlabeled Data

If $U = \{\mathbf{x}^{(i)}\}_{i=1}^n$ is a set of n iid unlabeled samples, the log likelihood of U under Φ is

$$\begin{aligned}L_U(\Phi) &= \sum_{i=1}^n \log P_{\Phi}(\mathbf{x}^{(i)}) \\&= \sum_{i=1}^n \log \left(\sum_{z=1}^m P_{\Phi}(\mathbf{x}^{(i)}, z) \right) \\&= \sum_{i=1}^n \log \left(\sum_{z=1}^m \log q(z) + \sum_{j=1}^m \log q(x_j^{(i)}|z, j) \right)\end{aligned}$$

Unfortunately, finding valid parameter values $q(z)$ and $q(x|z, j)$ that maximize this **marginalized** log likelihood is not as trivial (e.g., there is no closed-form solution).

Explanation of EM for This Problem

- ▶ EM is a **local search** algorithm to iteratively optimize

$$L_U(\Phi) = \sum_{i=1}^n \log \left(\sum_{z=1}^m P_{\Phi}(\mathbf{x}^{(i)}, z) \right)$$

That is, it calculates $\Phi^{(1)} \dots \Phi^{(T)}$ such that

$$L_U(\Phi^{(1)}) \leq L_U(\Phi^{(2)}) \leq \dots \leq L_U(\Phi^{(T-1)}) \leq L_U(\Phi^{(T)})$$

Explanation of EM for This Problem

- ▶ EM is a **local search** algorithm to iteratively optimize

$$L_U(\Phi) = \sum_{i=1}^n \log \left(\sum_{z=1}^m P_{\Phi}(\mathbf{x}^{(i)}, z) \right)$$

That is, it calculates $\Phi^{(1)} \dots \Phi^{(T)}$ such that

$$L_U(\Phi^{(1)}) \leq L_U(\Phi^{(2)}) \leq \dots \leq L_U(\Phi^{(T-1)}) \leq L_U(\Phi^{(T)})$$

- ▶ Importantly, each EM update *is trivial*: it has a closed-form solution.

Explanation of EM for This Problem

- ▶ EM is a **local search** algorithm to iteratively optimize

$$L_U(\Phi) = \sum_{i=1}^n \log \left(\sum_{z=1}^m P_{\Phi}(\mathbf{x}^{(i)}, z) \right)$$

That is, it calculates $\Phi^{(1)} \dots \Phi^{(T)}$ such that

$$L_U(\Phi^{(1)}) \leq L_U(\Phi^{(2)}) \leq \dots \leq L_U(\Phi^{(T-1)}) \leq L_U(\Phi^{(T)})$$

- ▶ Importantly, each EM update *is trivial*: it has a closed-form solution.
- ▶ As usual with local search algorithms, it only finds a local optimum and is not guaranteed to find a global optimum.

Posterior Probabilities

At each iteration t , we use the current parameter estimates

$$\Phi^{(t)} = \left\{ q^{(t)}(z), q^{(t)}(x|z, j) \right\}$$

to calculate the **posterior probabilities** on *individual* samples $\mathbf{x}^{(i)}$. This can be easily precomputed by Bayes rule: for every $i \in \{1 \dots n\}$ and $z \in \{1 \dots m\}$, calculate

$$P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) = \frac{P_{\Phi^{(t)}}(\mathbf{x}^{(i)}, z)}{P_{\Phi^{(t)}}(\mathbf{x}^{(i)})} = \frac{q^{(t)}(z) \prod_{j=1}^d q^{(t)}(x_j^{(i)}|z, j)}{\sum_{z=1}^m q^{(t)}(z) \prod_{j=1}^d q^{(t)}(x_j^{(i)}|z, j)}$$

Expected Log Likelihood of Labeled Data Under $\Phi^{(t)}$

$$\sum_{i=1}^n \sum_{z=1}^m P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \times \log P_{\Phi}(\mathbf{x}^{(i)}, z)$$

Expected Log Likelihood of Labeled Data Under $\Phi^{(t)}$

$$\begin{aligned} & \sum_{i=1}^n \sum_{z=1}^m P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \times \log P_{\Phi}(\mathbf{x}^{(i)}, z) \\ &= \sum_{i=1}^n \sum_{z=1}^m P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \times \left(\log q(z) + \sum_{j=1}^m \log q(x_j^{(i)}|z, j) \right) \end{aligned}$$

Expected Log Likelihood of Labeled Data Under $\Phi^{(t)}$

$$\begin{aligned} & \sum_{i=1}^n \sum_{z=1}^m P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \times \log P_{\Phi}(\mathbf{x}^{(i)}, z) \\ &= \sum_{i=1}^n \sum_{z=1}^m P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \times \left(\log q(z) + \sum_{j=1}^m \log q(x_j^{(i)}|z, j) \right) \\ &= \sum_{z=1}^m \sum_{i=1}^n P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \log q(z) + \sum_{z=1}^m \sum_{j=1}^m \sum_{i=1}^n P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \log q(x_j^{(i)}|z, j) \end{aligned}$$

Expected Log Likelihood of Labeled Data Under $\Phi^{(t)}$

$$\begin{aligned} & \sum_{i=1}^n \sum_{z=1}^m P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \times \log P_{\Phi}(\mathbf{x}^{(i)}, z) \\ &= \sum_{i=1}^n \sum_{z=1}^m P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \times \left(\log q(z) + \sum_{j=1}^m \log q(x_j^{(i)}|z, j) \right) \\ &= \sum_{z=1}^m \sum_{i=1}^n P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \log q(z) + \sum_{z=1}^m \sum_{j=1}^m \sum_{i=1}^n P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \log q(x_j^{(i)}|z, j) \\ &= \sum_{z=1}^m \widehat{\mathbf{count}}_t(z) \log q(z) + \sum_{z=1}^m \sum_{j=1}^m \sum_{x \in \{0,1\}} \widehat{\mathbf{count}}_t(z, j, x) \log q(x|z, j) \end{aligned}$$

where

$$\widehat{\mathbf{count}}_t(z) := \sum_{i=1}^n P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \quad \widehat{\mathbf{count}}_t(z, j, x) := \sum_{i=1: x_j^{(i)}=x}^n P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)})$$

MLE in the t -th Iteration of EM

What are the parameter values $q(z)$ and $q(x|z, j)$ that maximize

$$\sum_{z=1}^m \widehat{\mathbf{count}}_t(z) \log q(z) + \sum_{z=1}^m \sum_{j=1}^m \sum_{x \in \{0,1\}} \widehat{\mathbf{count}}_t(z, j, x) \log q(x|z, j)$$

under the constraints that they are nonnegative, $\sum_z q(z) = 1$, and $\sum_x q(x|z, j) = 1$?

MLE in the t -th Iteration of EM

What are the parameter values $q(z)$ and $q(x|z, j)$ that maximize

$$\sum_{z=1}^m \widehat{\mathbf{count}}_t(z) \log q(z) + \sum_{z=1}^m \sum_{j=1}^m \sum_{x \in \{0,1\}} \widehat{\mathbf{count}}_t(z, j, x) \log q(x|z, j)$$

under the constraints that they are nonnegative, $\sum_z q(z) = 1$, and $\sum_x q(x|z, j) = 1$?

Answer:

$$q(z) = \frac{\widehat{\mathbf{count}}_t(z)}{n} \qquad q(x|z, j) = \frac{\widehat{\mathbf{count}}_t(z, j, x)}{\sum_{x \in \{0,1\}} \widehat{\mathbf{count}}_t(z, j, x)}$$

EM for NB

1. Initialize NB parameters $\Phi^{(0)}$.
2. For $t = 0 \dots T - 1$,
 - 2.1 For $i = 1 \dots n$ and $z = 1 \dots m$, calculate current posterior

$$P_{\Phi^{(t)}}(z | \mathbf{x}^{(i)}) \leftarrow \frac{q^{(t)}(z) \prod_{j=1}^d q^{(t)}(x_j^{(i)} | z, j)}{\sum_{z=1}^m q^{(t)}(z) \prod_{j=1}^d q^{(t)}(x_j^{(i)} | z, j)}$$

- 2.2 “Count” $\widehat{\text{count}}_t(z) \leftarrow \sum_{i=1}^n P_{\Phi^{(t)}}(z | \mathbf{x}^{(i)})$ and
 $\widehat{\text{count}}_t(z, j, x) \leftarrow \sum_{i=1: x_j^{(i)}=x}^n P_{\Phi^{(t)}}(z | \mathbf{x}^{(i)})$ and set
 $\Phi^{(t+1)} = \{q^{(t+1)}(z), q^{(t+1)}(x|z, j)\}$ by

$$q^{(t+1)}(z) \leftarrow \frac{\widehat{\text{count}}_t(z)}{n} \quad q^{(t+1)}(x|z, j) \leftarrow \frac{\widehat{\text{count}}_t(z, j, x)}{\sum_{x \in \{0,1\}} \widehat{\text{count}}_t(z, j, x)}$$

3. Return $\Phi^{(T)}$.

Overview

EM for Naive Bayes

- Maximum Likelihood Estimation with Labeled Data

- Maximum Likelihood Estimation with Unlabeled Data

EM for Gaussian Mixture Models

- Maximum Likelihood Estimation with Labeled Data

- Maximum Likelihood Estimation with Unlabeled Data

k -Means Clustering

Gaussian Mixture Model: Definition

A **Gaussian mixture model (GMM)** with m clusters with identity covariance matrix in \mathbb{R}^d has $m + dm$ parameters, denoted by Φ :

- ▶ $\pi(z) \geq 0$ for each $z \in \{1 \dots m\}$ such that

$$\sum_z \pi(z) = 1$$

- ▶ $\mu_z \in \mathbb{R}^d$ for each $z \in \{1 \dots m\}$.

Gaussian Mixture Model: Definition

A **Gaussian mixture model (GMM)** with m clusters with identity covariance matrix in \mathbb{R}^d has $m + dm$ parameters, denoted by Φ :

- ▶ $\pi(z) \geq 0$ for each $z \in \{1 \dots m\}$ such that

$$\sum_z \pi(z) = 1$$

- ▶ $\mu_z \in \mathbb{R}^d$ for each $z \in \{1 \dots m\}$.

Φ defines a joint distribution over $\mathbf{x} \in \mathbb{R}^d$ and $z \in \{1 \dots m\}$ by

$$P_{\Phi}(\mathbf{x}, z) := \pi(z) \times \underbrace{\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \|\mathbf{x} - \mu_z\|_2^2\right)}_{\mathcal{N}(\mathbf{x}|\mu_z, I_d)}$$

Log Likelihood of Labeled Data and MLE

If $S = \{(\mathbf{x}^{(i)}, z^{(i)})\}_{i=1}^n$ is a set of n iid labeled samples, the log likelihood of S under Φ is

$$L_S(\Phi) = \sum_{i=1}^n \log P_{\Phi}(\mathbf{x}^{(i)}, y^{(i)})$$

Log Likelihood of Labeled Data and MLE

If $S = \{(\mathbf{x}^{(i)}, z^{(i)})\}_{i=1}^n$ is a set of n iid labeled samples, the log likelihood of S under Φ is

$$\begin{aligned}L_S(\Phi) &= \sum_{i=1}^n \log P_{\Phi}(\mathbf{x}^{(i)}, y^{(i)}) \\ &= \sum_{i=1}^n \log \pi(z^{(i)}) - \frac{1}{2} \left\| \mathbf{x}^{(i)} - \mu_{z^{(i)}} \right\|_2^2 - \log \sqrt{2\pi}\end{aligned}$$

Log Likelihood of Labeled Data and MLE

If $S = \{(\mathbf{x}^{(i)}, z^{(i)})\}_{i=1}^n$ is a set of n iid labeled samples, the log likelihood of S under Φ is

$$\begin{aligned}L_S(\Phi) &= \sum_{i=1}^n \log P_{\Phi}(\mathbf{x}^{(i)}, y^{(i)}) \\&= \sum_{i=1}^n \log \pi(z^{(i)}) - \frac{1}{2} \left\| \mathbf{x}^{(i)} - \mu_{z^{(i)}} \right\|_2^2 - \log \sqrt{2\pi} \\&= \left(\sum_{z=1}^m \mathbf{count}(z) \log \pi(z) \right) + \left(-\frac{1}{2} \sum_{i=1}^n \left\| \mathbf{x}^{(i)} - \mu_{z^{(i)}} \right\|_2^2 \right) + C\end{aligned}$$

Log Likelihood of Labeled Data and MLE

If $S = \{(\mathbf{x}^{(i)}, z^{(i)})\}_{i=1}^n$ is a set of n iid labeled samples, the log likelihood of S under Φ is

$$\begin{aligned}L_S(\Phi) &= \sum_{i=1}^n \log P_{\Phi}(\mathbf{x}^{(i)}, y^{(i)}) \\&= \sum_{i=1}^n \log \pi(z^{(i)}) - \frac{1}{2} \left\| \mathbf{x}^{(i)} - \mu_{z^{(i)}} \right\|_2^2 - \log \sqrt{2\pi} \\&= \left(\sum_{z=1}^m \mathbf{count}(z) \log \pi(z) \right) + \left(-\frac{1}{2} \sum_{i=1}^n \left\| \mathbf{x}^{(i)} - \mu_{z^{(i)}} \right\|_2^2 \right) + C\end{aligned}$$

Parameter values $\pi(z)$ (with probability constraints) and μ_z (with no constraints) that maximize $L_S(\Phi)$ are thus

$$\pi(z) = \frac{\mathbf{count}(z)}{n} \quad \mu_z = \frac{1}{\mathbf{count}(z)} \sum_{i=1: z^{(i)}=z}^n \mathbf{x}^{(i)}$$

Overview

EM for Naive Bayes

- Maximum Likelihood Estimation with Labeled Data

- Maximum Likelihood Estimation with Unlabeled Data

EM for Gaussian Mixture Models

- Maximum Likelihood Estimation with Labeled Data

- Maximum Likelihood Estimation with Unlabeled Data

k -Means Clustering

Log Likelihood of Unlabeled Data

If $U = \{\mathbf{x}^{(i)}\}_{i=1}^n$ is a set of n iid unlabeled samples, the log likelihood of U under Φ is

$$L_U(\Phi) = \sum_{i=1}^n \log P_{\Phi}(\mathbf{x}^{(i)})$$

Log Likelihood of Unlabeled Data

If $U = \{\mathbf{x}^{(i)}\}_{i=1}^n$ is a set of n iid unlabeled samples, the log likelihood of U under Φ is

$$\begin{aligned}L_U(\Phi) &= \sum_{i=1}^n \log P_{\Phi}(\mathbf{x}^{(i)}) \\ &= \sum_{i=1}^n \log \left(\sum_{z=1}^m P_{\Phi}(\mathbf{x}^{(i)}, z) \right)\end{aligned}$$

Log Likelihood of Unlabeled Data

If $U = \{\mathbf{x}^{(i)}\}_{i=1}^n$ is a set of n iid unlabeled samples, the log likelihood of U under Φ is

$$\begin{aligned}L_U(\Phi) &= \sum_{i=1}^n \log P_{\Phi}(\mathbf{x}^{(i)}) \\&= \sum_{i=1}^n \log \left(\sum_{z=1}^m P_{\Phi}(\mathbf{x}^{(i)}, z) \right) \\&= \sum_{i=1}^n \log \left(\sum_{z=1}^m \pi(z) \times \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{1}{2} \left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_z \right\|_2^2 \right) \right)\end{aligned}$$

Log Likelihood of Unlabeled Data

If $U = \{\mathbf{x}^{(i)}\}_{i=1}^n$ is a set of n iid unlabeled samples, the log likelihood of U under Φ is

$$\begin{aligned}L_U(\Phi) &= \sum_{i=1}^n \log P_{\Phi}(\mathbf{x}^{(i)}) \\&= \sum_{i=1}^n \log \left(\sum_{z=1}^m P_{\Phi}(\mathbf{x}^{(i)}, z) \right) \\&= \sum_{i=1}^n \log \left(\sum_{z=1}^m \pi(z) \times \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{1}{2} \left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_z \right\|_2^2 \right) \right)\end{aligned}$$

Again, finding valid parameter values $\pi(z)$ and $\boldsymbol{\mu}_z$ that maximize this **marginalized** log likelihood is not as trivial (there is no closed-form solution).

Log Likelihood of Unlabeled Data

If $U = \{\mathbf{x}^{(i)}\}_{i=1}^n$ is a set of n iid unlabeled samples, the log likelihood of U under Φ is

$$\begin{aligned}L_U(\Phi) &= \sum_{i=1}^n \log P_{\Phi}(\mathbf{x}^{(i)}) \\&= \sum_{i=1}^n \log \left(\sum_{z=1}^m P_{\Phi}(\mathbf{x}^{(i)}, z) \right) \\&= \sum_{i=1}^n \log \left(\sum_{z=1}^m \pi(z) \times \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{1}{2} \left\| \mathbf{x}^{(i)} - \mu_z \right\|_2^2 \right) \right)\end{aligned}$$

Again, finding valid parameter values $\pi(z)$ and μ_z that maximize this **marginalized** log likelihood is not as trivial (there is no closed-form solution).

EM is useful here again because each iteration *does* have a trivial solution.

Posterior Probabilities

At each iteration t , we use the current parameter estimates

$$\Phi^{(t)} = \left\{ \pi^{(t)}(z), \mu_z^{(t)} \right\}$$

to calculate the **posterior probabilities** on *individual* samples $\mathbf{x}^{(i)}$. This can again be easily precomputed by Bayes rule: for every $i \in \{1 \dots n\}$ and $z \in \{1 \dots m\}$, calculate

$$P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) = \frac{P_{\Phi^{(t)}}(\mathbf{x}^{(i)}, z)}{P_{\Phi^{(t)}}(\mathbf{x}^{(i)})} = \frac{\pi^{(t)}(z) \times \mathcal{N}(\mathbf{x}|\mu_z^{(t)}, I_d)}{\sum_{z=1}^m \pi^{(t)}(z) \times \mathcal{N}(\mathbf{x}|\mu_z^{(t)}, I_d)}$$

Expected Log Likelihood of Labeled Data Under $\Phi^{(t)}$

$$\sum_{i=1}^n \sum_{z=1}^m P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \times \log P_{\Phi}(\mathbf{x}^{(i)}, z)$$

Expected Log Likelihood of Labeled Data Under $\Phi^{(t)}$

$$\begin{aligned} & \sum_{i=1}^n \sum_{z=1}^m P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \times \log P_{\Phi}(\mathbf{x}^{(i)}, z) \\ &= \sum_{i=1}^n \sum_{z=1}^m P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \times \left(\log \pi(z) - \frac{1}{2} \left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_z \right\|_2^2 - \log \sqrt{2\pi} \right) \end{aligned}$$

Expected Log Likelihood of Labeled Data Under $\Phi^{(t)}$

$$\begin{aligned} & \sum_{i=1}^n \sum_{z=1}^m P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \times \log P_{\Phi}(\mathbf{x}^{(i)}, z) \\ &= \sum_{i=1}^n \sum_{z=1}^m P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \times \left(\log \pi(z) - \frac{1}{2} \left\| \mathbf{x}^{(i)} - \mu_z \right\|_2^2 - \log \sqrt{2\pi} \right) \\ &= \underbrace{\sum_{z=1}^m \sum_{i=1}^n P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \log \pi(z)}_{\widehat{\text{count}}_t(z)} - \sum_{i=1}^n \sum_{z=1}^m P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \left\| \mathbf{x}^{(i)} - \mu_z \right\|_2^2 \end{aligned}$$

Expected Log Likelihood of Labeled Data Under $\Phi^{(t)}$

$$\begin{aligned} & \sum_{i=1}^n \sum_{z=1}^m P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \times \log P_{\Phi}(\mathbf{x}^{(i)}, z) \\ &= \sum_{i=1}^n \sum_{z=1}^m P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \times \left(\log \pi(z) - \frac{1}{2} \left\| \mathbf{x}^{(i)} - \mu_z \right\|_2^2 - \log \sqrt{2\pi} \right) \\ &= \underbrace{\sum_{z=1}^m \sum_{i=1}^n P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \log \pi(z)}_{\widehat{\text{count}}_t(z)} - \sum_{i=1}^n \sum_{z=1}^m P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \left\| \mathbf{x}^{(i)} - \mu_z \right\|_2^2 \end{aligned}$$

MLE in the t -th iteration of EM

$$\pi(z) = \frac{\widehat{\text{count}}_t(z)}{n} \quad \mu_z = \frac{\sum_{i=1}^n P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) \mathbf{x}^{(i)}}{\widehat{\text{count}}_t(z)}$$

EM for GMMs

1. Initialize GMM parameters $\Phi^{(0)}$.
2. For $t = 0 \dots T - 1$,
 - 2.1 For $i = 1 \dots n$ and $y = 1 \dots m$, calculate current posterior

$$P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)}) = \frac{\pi^{(t)}(z) \times \mathcal{N}(\mathbf{x}|\mu_z^{(t)}, I_d)}{\sum_{z=1}^m \pi^{(t)}(z) \times \mathcal{N}(\mathbf{x}|\mu_z^{(t)}, I_d)}$$

- 2.2 Set $\Phi^{(t+1)} = \left\{ \pi^{(t+1)}(z), \mu_z^{(t+1)} \right\}$ by

$$\pi(z) = \frac{\widehat{\text{count}}_t(z)}{n} \quad \mu_z = \frac{\sum_{i=1}^n P_{\Phi^{(t)}}(z|\mathbf{x}^{(i)})\mathbf{x}^{(i)}}{\widehat{\text{count}}_t(z)}$$

3. Return $\Phi^{(T)}$.

Overview

EM for Naive Bayes

- Maximum Likelihood Estimation with Labeled Data

- Maximum Likelihood Estimation with Unlabeled Data

EM for Gaussian Mixture Models

- Maximum Likelihood Estimation with Labeled Data

- Maximum Likelihood Estimation with Unlabeled Data

k-Means Clustering

Non-Probabilistic Clustering

- ▶ You can train a GMM Φ with k clusters with EM and obtain a “soft” k -clustering given by the posterior

$$P_{\Phi}(z|\mathbf{x}^{(i)}) = \frac{\pi(z) \times \mathcal{N}(\mathbf{x}|\mu_z, I_d)}{\sum_{z=1}^k \pi(z) \times \mathcal{N}(\mathbf{x}|\mu_z, I_d)}$$

- ▶ If all you want is to cluster n points $\mathbf{x}^{(1)} \dots \mathbf{x}^{(n)} \in \mathbb{R}^d$ into k clusters, you can do **k -means clustering**.

k -Means Clustering

Input: points $U = \{\mathbf{x}^{(i)}\}_{i=1}^n$ in \mathbb{R}^d , number of clusters k , T

Output: cluster assignments $a_1 \dots a_n \in \{1 \dots k\}$

1. Initialize centroids: $\boldsymbol{\nu}_1^{(0)} \dots \boldsymbol{\nu}_k^{(0)} \in \mathbb{R}^d$.
2. For $t = 0 \dots T - 1$,
 - 2.1 Assign each point to its closest centroid:

$$a_i^{(t)} \leftarrow \arg \min_{j=1}^k \left\| \mathbf{x}^{(i)} - \boldsymbol{\nu}_j^{(t)} \right\|_2^2$$

- 2.2 Update centroids: denoting $C_j^{(t)} := \{\mathbf{x}^{(i)} : a_i^{(t)} = j\}$,

$$\boldsymbol{\nu}_j^{(t+1)} \leftarrow \frac{1}{|C_j^{(t)}|} \sum_{\mathbf{x} \in C_j^{(t)}} \mathbf{x}$$

3. Return $a_i^{(T)} \leftarrow \arg \min_{j=1}^k \left\| \mathbf{x}^{(i)} - \boldsymbol{\nu}_j^{(T)} \right\|_2^2$.

Loss of k -Means Clustering

Using indicator $[[A]]$ which is 1 if A is true and 0 otherwise,

$$L(\boldsymbol{\nu}_1 \dots \boldsymbol{\nu}_k, a_1 \dots a_n) := \sum_{i=1}^n \sum_{j=1}^k [[a_i = j]] \left\| \mathbf{x}^{(i)} - \boldsymbol{\nu}_j \right\|_2^2$$

Loss of k -Means Clustering

Using indicator $[[A]]$ which is 1 if A is true and 0 otherwise,

$$L(\boldsymbol{\nu}_1 \dots \boldsymbol{\nu}_k, a_1 \dots a_n) := \sum_{i=1}^n \sum_{j=1}^k [[a_i = j]] \left\| \mathbf{x}^{(i)} - \boldsymbol{\nu}_j \right\|_2^2$$

k -means is an **alternating minimization** algorithm for this loss.

1. Fix centroids $\boldsymbol{\nu}_1 \dots \boldsymbol{\nu}_k$, optimize over assignments $a_1 \dots a_n$:

$$a_i \leftarrow \arg \min_{j=1}^k \left\| \mathbf{x}^{(i)} - \boldsymbol{\nu}_j \right\|_2^2$$

Loss of k -Means Clustering

Using indicator $[[A]]$ which is 1 if A is true and 0 otherwise,

$$L(\boldsymbol{\nu}_1 \dots \boldsymbol{\nu}_k, a_1 \dots a_n) := \sum_{i=1}^n \sum_{j=1}^k [[a_i = j]] \left\| \mathbf{x}^{(i)} - \boldsymbol{\nu}_j \right\|_2^2$$

k -means is an **alternating minimization** algorithm for this loss.

1. Fix centroids $\boldsymbol{\nu}_1 \dots \boldsymbol{\nu}_k$, optimize over assignments $a_1 \dots a_n$:

$$a_i \leftarrow \arg \min_{j=1}^k \left\| \mathbf{x}^{(i)} - \boldsymbol{\nu}_j \right\|_2^2$$

2. Fix assignments $a_1 \dots a_n$, optimize over centroids $\boldsymbol{\nu}_1 \dots \boldsymbol{\nu}_k$:

$$\boldsymbol{\nu}_j \leftarrow \arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1: a_i=j}^n \left\| \mathbf{x}^{(i)} - \mathbf{w} \right\|_2^2 = \frac{1}{|C_j|} \sum_{\mathbf{x} \in C_j} \mathbf{x}$$

Thus k -means can only decrease the loss in each step.

Generalization of k -Means

Choose a “distortion” function $D(\mathbf{x}, \mathbf{y}) \geq 0$ and do alternating minimization of

$$L(\boldsymbol{\nu}_1 \dots \boldsymbol{\nu}_k, a_1 \dots a_n) := \sum_{i=1}^n \sum_{j=1}^k [[a_i = j]] D(\mathbf{x}^{(i)}, \boldsymbol{\nu}_j)$$

1. Fix centroids $\boldsymbol{\nu}_1 \dots \boldsymbol{\nu}_k$, optimize over assignments $a_1 \dots a_n$:

$$a_i \leftarrow \arg \min_{j=1}^k D(\mathbf{x}^{(i)}, \boldsymbol{\nu}_j)$$

2. Fix assignments $a_1 \dots a_n$, optimize over centroids $\boldsymbol{\nu}_1 \dots \boldsymbol{\nu}_k$:

$$\boldsymbol{\nu}_j \leftarrow \arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1: a_i=j}^n D(\mathbf{x}^{(i)}, \mathbf{w})$$

Choice of Distortion Function

- ▶ The standard k -means clustering uses squared Euclidean distance $D(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2$ and

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1: a_i=j}^n D(\mathbf{x}^{(i)}, \mathbf{w})$$

is given by the **mean** of C_j .

- ▶ It turns out that for a wide class of distortion functions called the **Bregman divergence**, this optimization is always given by the mean of C_j .
- ▶ Examples of Bregman divergence: squared Euclidean norm, KL divergence (this only makes sense if data points are probability distributions).
- ▶ So we can swap in any Bregman divergence and perform exactly the same updates.

k -Medians Clustering

- ▶ Use the Manhattan distance in the algorithm:

$$D(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_1 := \sum_{l=1}^d |x_l - y_l|$$

- ▶ The solution of

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1: a_i=j}^n \|\mathbf{x}^{(i)} - \mathbf{w}\|_1$$

is given by the **element-wise median** of C_j .