Expectation Maximization (wrap-up) and Intro to Probabilistic PCA

Piyush Rai IIT Kanpur

Probabilistic Machine Learning (CS772A)

Feb 1, 2016

・ 同 ト ・ ヨ ト ・ ヨ ト

Parameter Estimation with Latent Variables

- Model $p(\mathbf{X}, \mathbf{Z}|\theta)$, observed data \mathbf{X} , latent variables \mathbf{Z} , model parameters θ
- Recall GMM, Z: cluster assignments, θ : GMM parameters $\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$
- Goal: Estimate the model parameters $\boldsymbol{\theta}$ via MLE

$$\hat{\theta} = \arg \max_{\theta} \log p(\mathbf{X}|\theta) = \arg \max_{\theta} \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta)$$

- Doing MLE in such models can be difficult because of the log-sum
- If we "knew" Z, sum over all possible Z not needed. Just define "complete data" {X, Z}, and do MLE on the complete data log-lik. log p(X, Z|θ)
- Assumption: MLE on log $p(\mathbf{X}, \mathbf{Z}|\theta)$ is easy
 - It often indeed is, especially when p(X, Z|θ) is exponential family distribution (or product of exponential family distributions)

イロン イロン イヨン イヨン 三日

Parameter Estimation with Latent Variables

- If MLE on log $p(\mathbf{X}, \mathbf{Z}|\theta)$ is easy then let's do it!
- ullet Problem: Well, we don't actually know f Z, so we are still stuck. igodot
- Solution: Use the posterior p(Z|X, θ) over latent variables Z to compute the expected complete data log-likelihood and do MLE on *that* objective.

$$\begin{split} \hat{\theta} &= \arg \max_{\theta} \ \mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\theta)] \\ &= \arg \max_{\theta} \ \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta) \log p(\mathbf{X}, \mathbf{Z}|\theta) \end{split}$$

But now we have a chicken-and-egg problem: the posterior p(Z|X, θ) over Z itself depends on the parameters θ

Solution: An Iterative Scheme (EM Algorithm)

Initialize the parameters: θ^{old} . Then alternate between these steps:

• E (Expectation) step:

- Compute the posterior $p(\mathbf{Z}|\mathbf{X}, \theta^{old})$ over latent variables \mathbf{Z} using θ^{old}
- Compute the expected complete data log-likelihood w.r.t. this posterior

$$\mathcal{Q}(\theta, \theta^{old}) = \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \theta^{old})}[\log p(\mathbf{X}, \mathbf{Z}|\theta)] = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{old}) \log p(\mathbf{X}, \mathbf{Z}|\theta)$$

- M (Maximization) step:
 - Maximize the expected complete data log-likelihood w.r.t. θ

$$\begin{array}{lll} \theta^{new} & = & \arg\max_{\theta} \mathcal{Q}(\theta, \theta^{old}) & (\text{if doing MLE}) \\ \theta^{new} & = & \arg\max_{\theta} \{\mathcal{Q}(\theta, \theta^{old}) + \log p(\theta)\} & (\text{if doing MAP}) \end{array}$$

• If the log-likelihood or the parameter values not converged then set $\theta^{old} = \theta^{new}$ and go to the E step.

Why is this doing the right thing?

ヘロト 不合 ト 不良 ト 不良 ト 一度

Illustration: EM for GMM

- Recall that the GMM parameters $\theta = \{\pi, \mu, \Sigma\} = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$
- The complete data likelihood

$$p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{n=1}^{N} \prod_{k=1}^{K} p(\boldsymbol{z}_n = k) p(\boldsymbol{x}_n | \boldsymbol{z}_n = k) = \prod_{n=1}^{N} \prod_{k=1}^{K} \pi_k^{z_{nk}} \mathcal{N}(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_{nk}}$$

• Taking the log, we get:

$$\log p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\pi}, \boldsymbol{\mu}, \mathbf{\Sigma}) = \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} \{ \log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \mu_k, \boldsymbol{\Sigma}_k) \}$$

• E-step computes the expected complete data log-likelihood:

$$\mathbb{E}_{p(\mathbf{Z}|\mathbf{X},\theta)}[\log p(\mathbf{X},\mathbf{Z}|\boldsymbol{\pi},\boldsymbol{\mu},\mathbf{\Sigma})] = \sum_{n=1}^{N} \sum_{k=1}^{K} \mathbb{E}[\boldsymbol{z}_{nk}] \{\log \pi_{k} + \log \mathcal{N}(\boldsymbol{x}_{n}|\boldsymbol{\mu}_{k},\boldsymbol{\Sigma}_{k})\}$$

where $\mathbb{E}[z_{nk}]$ is the expected value of z_{nk} under the posterior

Illustration: EM for GMM (Contd.)

• The only expectation we need to compute $\mathbb{E}_{p(Z|X,\theta)}[\log p(X, Z|\pi, \mu, \Sigma)]$ is

$$\mathbb{E}[z_{nk}] = \sum_{z_{nk} = \{0,1\}} z_{nk} p(z_{nk} | \mathbf{x}_n, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = p(z_{nk} = 1 | \mathbf{x}_n, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \mu_j, \boldsymbol{\Sigma}_j)} = \gamma_{nk}$$

• Thus the expected complete data log-likelihood

$$\mathbb{E}_{p(\mathbf{Z}|\mathbf{X},\theta)}[\log p(\mathbf{X},\mathbf{Z}|\boldsymbol{\pi},\boldsymbol{\mu},\mathbf{\Sigma})] = \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma_{nk} \{\log \pi_{k} + \log \mathcal{N}(\mathbf{x}_{n}|\mu_{k},\boldsymbol{\Sigma}_{k})\}$$

- M-step maximizes the the exp. complete data log-likelihood w.r.t. π_k, μ_k, Σ_k
- The update equations for these will be (shown on the board)

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} \boldsymbol{x}_n, \quad \boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} (\boldsymbol{x}_n - \mu_k) (\boldsymbol{x}_n - \mu_k)^\top, \quad \boldsymbol{\pi}_k = \frac{N_k}{N}$$

where $N_k = \sum_{n=1}^{N} \gamma_{nk}$ is "effective" num. of examples assigned to k^{th} Gaussian

Why does EM work?

э

< ロ > < 同 > < 回 > < 回 > < 回 > <

Justification 1

• Consider the log likelihood on "incomplete" data X

$$\log p(\mathbf{X}|\theta) = \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) = \log \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{q(\mathbf{Z})} \quad \text{(where } q(\mathbf{Z}) \text{ is some distribution)}$$

$$\geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{q(\mathbf{Z})} \quad \text{(using Jensen's inequality for concave log)}$$

$$\log p(\mathbf{X}|\theta) \geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\theta) - \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \log q(\mathbf{Z})}_{\text{doesn't depend on } \theta} = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}|\theta) + \text{const.}$$

• If we set $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \theta)$ then the above inequality becomes equality

$$\sum_{\mathbf{Z}} q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{q(\mathbf{Z})} = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta) \log \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{p(\mathbf{Z}|\mathbf{X}, \theta)} = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta) \log \frac{p(\mathbf{Z}|\mathbf{X}, \theta)p(\mathbf{X}|\theta)}{p(\mathbf{Z}|\mathbf{X}, \theta)}$$
$$= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta) \log p(\mathbf{X}|\theta)$$
$$= \log p(\mathbf{X}|\theta) \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta) = \log p(\mathbf{X}|\theta)$$

• Thus for $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \theta)$, we have

$$\log p(\mathbf{X}|\theta) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X},\theta) \log p(\mathbf{X},\mathbf{Z}|\theta) + \text{const.} = \mathbb{E}[\log p(\mathbf{X},\mathbf{Z}|\theta)] + \text{const.}$$

• EM maximizes $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z}|\theta)]$, a tight lower-bound on $\log p(\mathbf{X}|\theta)$

Justification 2

• We can also write the incomplete log likelihood

 $\log p(\mathbf{X}| heta) = \mathcal{L}(q, heta) + \mathsf{KL}(q||p_z)$

where q is some distr. on \mathbf{Z} , $p_{Z} = p(\mathbf{Z}|\mathbf{X}, \theta)$ is the posterior over \mathbf{Z} , and $\mathcal{L}(q, \theta) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{q(\mathbf{Z})} \right\}$ $\operatorname{KL}(q||p_{z}) = -\sum_{\mathbf{Z}} q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \theta)}{q(\mathbf{Z})} \right\}$ $\operatorname{KL}(q||p)$ $\operatorname{KL}(q||p)$ $\operatorname{KL}(q||p)$

(to verify, use $\log p(\mathbf{X}, \mathbf{Z}|\theta) = \log p(\mathbf{Z}|\mathbf{X}, \theta) + \log p(\mathbf{X}|\theta)$ in the expression of $\mathcal{L}(q, \theta)$)

• Since $\mathsf{KL}(q||p_z) \ge 0$, $\mathcal{L}(q, \theta)$ is a lower-bound on $\log p(\mathbf{X}|\theta)$ for any q

Picture courtesy: PRML (Bishop, 2006)

Probabilistic Machine Learning (CS772A)

・ 同 ト ・ ヨ ト ・ ヨ ト

Justification 2 (contd.)

Recall log $p(\mathbf{X}|\theta) = \mathcal{L}(q, \theta) + KL(q||p_z)$. EM can also be seen as:

• With θ fixed to θ^{old} , maximize $\mathcal{L}(q, \theta^{old})$ w.r.t. q

$$\hat{q} = rg\max_{q} \mathcal{L}(q, heta^{\textit{old}})$$

which is equivalent to making $KL(q||p_z) = 0$ or setting $\hat{q} = p(\mathbf{Z}|\mathbf{X}, \theta^{old})$ (This step makes $\mathcal{L}(\hat{q}, \theta^{old}) = \log p(\mathbf{X}|\theta^{old})$; see next slide)

• With \hat{q} fixed at $p(\mathbf{Z}|\mathbf{X}, \theta^{old})$, maximize $\mathcal{L}(\hat{q}, \theta)$ w.r.t. θ , where

$$\mathcal{L}(\hat{q},\theta) = \sum_{\mathsf{Z}} p(\mathsf{Z}|\mathsf{X},\theta^{old}) \log p(\mathsf{X},\mathsf{Z}|\theta) - \underbrace{\sum_{\mathsf{Z}} p(\mathsf{Z}|\mathsf{X},\theta^{old}) \log p(\mathsf{Z}|\mathsf{X},\theta^{old})}_{\text{constant w.r.t. } \theta}$$

$$\mathcal{Q}(\theta, \theta^{old}) + \text{const}$$

$$\theta^{\textit{new}} = \arg\max_{\theta} \mathcal{Q}(\theta, \theta^{\textit{old}})$$

(This step ensures that $\log p(\mathbf{X}|\theta^{new}) \ge \log p(\mathbf{X}|\theta^{old})$; see next slide)

Justification 2 (contd.)

E-step: $\mathcal{L}(q, \theta^{old})$ increases and becomes equal to $\log p(\mathbf{X}|\theta^{old})$, $\mathsf{KL}(q||p_z)$ becomes 0 because we set $q = p(\mathbf{Z}|\mathbf{X}, \theta)$



M-step: θ^{new} makes $\mathcal{L}(q, \theta^{new})$ go further up, makes $\mathsf{KL}(q||p_z) > 0$ again because $q \neq p(\mathbf{Z}|\mathbf{X}, \theta^{new})$ and thus ensures that $\log p(\mathbf{X}|\theta^{new}) \geq \log p(\mathbf{X}|\theta^{old})$



Thus the E and M steps never decrease the log-likelihood $p(\mathbf{X}|\theta)$

Picture courtesy: PRML (Bishop, 2006)

▶ < 글 > < 글 > ...

A View in the Parameter Space

- E-step: Update of q makes the $\mathcal{L}(q, \theta)$ curve touch the log $p(\mathbf{X}|\theta)$ curve
- M-step gives the maxima θ^{new} of $\mathcal{L}(q, \theta)$
- Next E-step readjusts $\mathcal{L}(q, \theta)$ curve (green) to meet log $p(\mathbf{X}|\theta)$ curve again
- This continues until a local maxima of log $p(\mathbf{X}|\theta)$ is reached



< 回 ト < 三 ト < 三 ト

Some EM Variants

- Generalized EM: M step doesn't require maximization w.r.t. θ; even if the new θ just increases the lower bound, we will still converge to a local optima
- Variational EM and MCMC EM: If the E step of computing the posterior p(Z|X, θ) is intractable, we can use variational Bayes (VB) or MCMC to approximate the posterior
- Expectation Conditional Maximization: Parameters are partitioned in groups. M step consists of multiple steps (each optimizing one group of parameters, treating all other groups as fixed)
- Online/incremental EM: E step only processes one (or a small number of) observation, computing posteriors/expectations only w.r.t. that minibatch of data. For exponential famility distributions, the sufficient statistics needed in the M step can be easily updated incrementally, leading to simple form of incremental parameter updates. Very useful for scalable inference. See:

 Online EM Algorithm for Latent Data Models (Cappé & Moulines, 2009)
 - (2) Online EM for Unsupervised Models (Liang & Klein, 2009)

Next up: Probabilistic PCA and Factor Analysis

周下 イモト イモト