

More on EM and variational inference

Julia Hockenmaier

juliahmr@illinois.edu
3324 Siebel Center

<http://www.cs.uiuc.edu/class/sp10/cs598jhm>

Exponential Families

Today

- The connection between EM and variational inference
- Exponential families

Bishop (2006) sections 2.4, 9.3, 9.4, 10.1, 10.4

CS598JHM: **Advanced NLP**
2

Exponential families

- Most parametric distributions that we've seen so far belong to the exponential family of distributions

Distributions in the exponential family are nice because:

- They have **conjugate priors** (other distributions generally don't)
- The likelihood and posterior can be expressed in terms of **sufficient statistics**

CS598JHM: **Advanced NLP**
4

Definition

The **exponential family of distributions** over x

(x can be scalar or vector; discrete or continuous)

given the “natural” parameters η

is defined as the set of distributions

$$p(x|\eta) = h(x)g(\eta)\exp(\eta^T u(x))$$

- $g(\eta)$: normalization coefficient:

$$g(\eta) = (\int \exp(\eta^T u(x)))^{-1}$$

- $u(x)$: some function of x

Likelihood

Given a sequence of i.i.d observations $Y=(y_1, \dots, y_n)$, the likelihood $P(Y|\eta)$ is:

$$P(Y|\eta) = \left[\prod_{i=1}^n h(y_i) \right] g(\eta)^n \exp \left(\eta^T \sum_{i=1}^n u(y_i) \right)$$

Define a function $t(Y)$, called **sufficient statistics**:

$$t(Y) = \sum_{i=1}^n u(y_i)$$

$$\begin{aligned} \text{Thus: } P(Y|\eta) &= \left[\prod_{i=1}^n h(y_i) \right] g(\eta)^n \exp \left(\eta^T t(Y) \right) \\ &\propto g(\eta)^n \exp \left(\eta^T t(Y) \right) \end{aligned}$$

Conjugate priors

It is straightforward to define a conjugate prior for members of the exponential family:

$$\begin{aligned} \text{Likelihood } P(Y|\eta) &\propto g(\eta)^n e^{\eta^T t(Y)} \\ \text{Prior } P(\eta) &\propto g(\eta)^\mu e^{\eta^T \nu} \\ \text{Posterior } P(\eta|Y) &\propto P(\eta)P(Y|\eta) \\ &= g(\eta)^{\mu+n} e^{\eta^T (\nu+t(Y))} \end{aligned}$$

Expectation Maximization revisited

The EM algorithm

The goal of EM: Find the maximum likelihood solution for a model consisting of parameters θ , given observed (incomplete) data \mathbf{X} and latent variables \mathbf{Z}

$$\ln p(\mathbf{X}|\theta) = \ln \left\{ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) \right\}$$

Note: Even if complete likelihood $p(\mathbf{X}, \mathbf{Z} | \theta)$ is in exponential family, incomplete likelihood $p(\mathbf{X} | \theta)$ may not be.

We just have **incomplete data \mathbf{X}** , so don't know $p(\mathbf{X}, \mathbf{Z} | \theta)$.

We can only infer \mathbf{Z} from posterior $p(\mathbf{Z} | \mathbf{X}, \theta)$.

We will compute the **expectation** of $p(\mathbf{X}, \mathbf{Z} | \theta)$ wrt. $p(\mathbf{Z} | \mathbf{X}, \theta)$

The EM algorithm

1. **Initialization:** Choose initial θ^{old}

2. **Expectation step:**

Compute **posterior of the latent variables** $p(\mathbf{Z} | \mathbf{X}, \theta^{old})$

3. **Maximization step:**

Find θ^{new} which **maximize the expected log-likelihood of the joint** $p(\mathbf{Z}, \mathbf{X} | \theta^{new})$ under $p(\mathbf{Z} | \mathbf{X}, \theta^{old})$:

$$\theta^{new} = \arg \max_{\theta} \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \theta^{old}) \ln p(\mathbf{X}, \mathbf{Z} | \theta)$$

4. **Check for convergence.**

Stop, or set $\theta^{old} := \theta^{new}$ and go to 2.

Another view of EM

We want to maximize $\ln p(\mathbf{X}|\theta) = \ln \left\{ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) \right\}$

By the product rule: $\ln p(\mathbf{X}, \mathbf{Z}|\theta) = \ln p(\mathbf{Z}|\mathbf{X}, \theta) + \ln p(\mathbf{X}|\theta)$

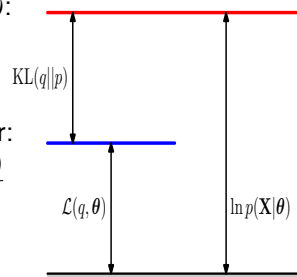
Define a functional of distribution $q(\mathbf{Z})$:

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

KL-divergence btw. $q(\mathbf{Z})$ and posterior:

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

Thus $\ln p(\mathbf{X}|\theta) = \mathcal{L}(q, \theta) + KL(q||p)$



EM again...

$\mathcal{L}(q, \theta)$ is a lower bound on log-likelihood $\ln p(\mathbf{X} | \theta)$

E-step:

Maximize $\mathcal{L}(q, \theta^{old})$ wrt. $q(\mathbf{Z})$, keep θ^{old} fixed.

This happens when $KL(q||p) = 0$.

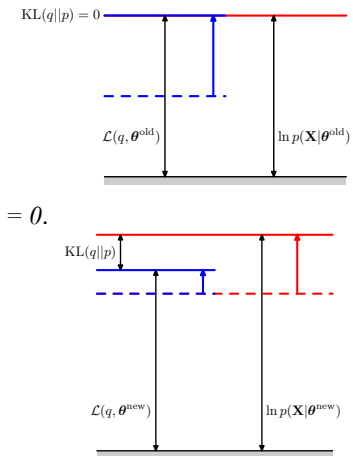
M-step:

Maximize $\mathcal{L}(q, \theta^{old})$ wrt. θ , keep $q(\mathbf{Z})$ fixed.

$\mathcal{L}(q, \theta)$ will increase.

Thus $\ln p(\mathbf{X} | \theta)$ will increase.

Hence, now: $KL(q||p) > 0$



Bayesian model

- In a fully Bayesian model, all parameters θ are stochastic variables with priors.
- Now \mathbf{Z} consists of latent variables and priors.
- We still want to maximize (incomplete) log-likelihood:

$$\ln p(\mathbf{X}) = \mathcal{L}(q) + KL(q||p)$$

$$\mathcal{L}(q) = \int q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z}$$

$$KL(q||p) = - \int q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}|\mathbf{X})}{q(\mathbf{Z})} d\mathbf{Z}$$

CS598JHM: Advanced NLP

14

Variational inference for Bayesian models

Factorized distributions

Assume q factorizes:

$$q(\mathbf{Z}) = \prod_{i=1}^M q_i(\mathbf{Z}_i)$$

We still want to maximize $\mathcal{L}(q)$.

We can do this by optimizing with respect to each factor q_i in turn

$$\begin{aligned} \mathcal{L}(q) &= \int \prod_i q_i \left\{ \ln p(\mathbf{X}, \mathbf{Z}) - \sum_i \ln q_i \right\} d\mathbf{Z} \\ &= \int q_j \langle \ln p(\mathbf{X}, \mathbf{Z}) \rangle_{i \neq j} + c' d\mathbf{Z}_j - \int q_i \ln q_j d\mathbf{Z}_j + c \end{aligned}$$

CS598JHM: Advanced NLP

15