

Debugging UB

How do we know when UB has converged? How do we know that our answer is right?

One check: UB lower bound

$$\mathcal{L} = \sum_{\mathbf{z}} \iiint q(\mathbf{z}, \pi, \mu, \lambda) \ln \left\{ \frac{p(x, \mathbf{z}, \pi, \mu, \lambda)}{q(\mathbf{z}, \pi, \mu, \lambda)} \right\} d\pi d\mu d\lambda$$

~~=~~ when this stops going up we can declare our inference algorithm to have converged. Dropping the \mathbb{E} 's and Expectation subscripts we have

$$= \mathbb{E}[\ln p(x, \mathbf{z}, \pi, \mu, \lambda)] - \mathbb{E}[\ln q(\mathbf{z}, \pi, \mu, \lambda)]$$

$$\begin{aligned} \textcircled{2} &= \mathbb{E}[\ln p(x|\mathbf{z}, \mu, \lambda)] + \mathbb{E}[\ln p(\mathbf{z}|\pi)] + \mathbb{E}[\ln p(\pi)] \\ &\quad + \mathbb{E}[\ln p(\mu, \lambda)] - \underbrace{\mathbb{E}[\ln q(\mathbf{z})] - \mathbb{E}[\ln q(\pi)] - \mathbb{E}[\ln q(\mu, \lambda)]} \end{aligned}$$

Note these are all entropy terms that can be looked up

We avoided doing these expectations before, we'll do one here just to show / highlight the required techniques.

Foremost is a conditional expectation trick.

When $p(a, b)$ naturally factorizes as $p(b|a)p(a) = p(a, b)$ - or is specified as such - we can perform expectations in a step wise manner

$$\mathbb{E}_{a,b}[f(x, a, b)] = \mathbb{E}_a[\mathbb{E}_{b|a}[f(x, a, b)]]$$

$$\begin{aligned} \text{pf.} &= \mathbb{E}_a\left[\sum_b f(x, a, b) p(b|a)\right] \\ &= \sum_a \left(\sum_b f(x, a, b) p(b|a)\right) p(a) \\ &= \sum_a \sum_b f(x, a, b) p(a, b) \quad \square \end{aligned}$$

all wrt. z^* dists

Let's work on first term

$$\mathbb{E}[\ln p(x|z, \mu, \Sigma)] = \mathbb{E}_z \left[\mathbb{E}_{\mu, \Sigma} \left[\mathbb{E}_z \left[\ln \prod_u \prod_k \mathcal{N}(x_u | \mu_k, \Sigma_k^{-1})^{z_{uk}} \right] \right] \right]$$

applying the log & taking the Expectation wrt. $z^*(z)$ gives

$$= \mathbb{E}_z \left[\mathbb{E}_{\mu, \Sigma} \left[\sum_u \sum_k \mathbb{E}[z_{uk}] \cdot \ln \mathcal{N}(x_u | \mu_k, \Sigma_k^{-1}) \right] \right]$$

we know this, $\mathbb{E}_{z^*(z)}[z_{uk}] = r_{uk}$

pushing the Expectations into the sums & taking the log of the above

$$= \sum_u \sum_k r_{uk} \mathbb{E}_z \left[\mathbb{E}_{\mu, \Sigma} \left[-\frac{1}{2} (x_u - \mu_k)^T \Sigma_k (x_u - \mu_k) - \frac{D}{2} \ln(2\pi) + \frac{1}{2} \ln |\Sigma_k| \right] \right]$$

$$\textcircled{1} = \sum_u \sum_k r_{uk} \left(\underbrace{\mathbb{E}_z \left[\mathbb{E}_{\mu, \Sigma} \left[-\frac{1}{2} (x_u - \mu_k)^T \Sigma_k (x_u - \mu_k) \right] \right]}_{\text{this term takes a little work}} - \frac{D}{2} \ln(2\pi) \right) + \frac{1}{2} \underbrace{\mathbb{E}_z \left[\ln |\Sigma_k| \right]}_{\text{this term can be looked up for Wishart distributions B.81 pg 693 PRML}}$$

this term takes a little work

this term can be looked up for Wishart distributions B.81 pg 693 PRML

Let's generalize this in the following way

$$\text{Let } \mu \sim \mathcal{N}(a, (\Sigma b)^{-1}) \text{ and } \Sigma \sim \mathcal{W}(\Psi, \nu)$$

$$\text{What is } \mathbb{E}_z \left[\mathbb{E}_{\mu, \Sigma} \left[(x - \mu)^T \Sigma (x - \mu) \right] \right] ?$$

this is a case where conditional expectation from previous page helps

Nice trick - recast problem slightly by adding and subtracting mean of μ , i.e.

$$\text{What is } \mathbb{E}_{z, \mu} \left[(x - a + a - \mu)^T \Sigma (x - a + a - \mu) \right] ?$$

This can be expanded like

$$= \mathbb{E}_{\Sigma} \left[\mathbb{E}_{\mu|\Sigma} \left[\begin{aligned} &(x-a)^T \Sigma (x-a) + (x-a)^T \Sigma (a-\mu) \\ &+ (a-\mu)^T \Sigma (x-a) + (a-\mu)^T \Sigma (a-\mu) \end{aligned} \right] \right]$$

which has some nice properties, namely

$$= \mathbb{E}_{\Sigma} \left[\begin{aligned} &(x-a)^T \Sigma (x-a) + \underbrace{(x-a)^T \Sigma (a - \mathbb{E}[\mu])}_{\rightarrow 0} \\ &+ \underbrace{(a - \mathbb{E}[\mu])^T \Sigma (x-a)}_0 + \underbrace{\mathbb{E}_{\mu|\Sigma} [(a-\mu)^T \Sigma (a-\mu)]}_{\text{almost a } \chi^2 \text{ RV.}} \end{aligned} \right]$$

$$= \mathbb{E}_{\Sigma} \left[(x-a)^T \Sigma (x-a) + \frac{1}{b} \mathbb{E}_{\mu|\Sigma} [(a-\mu)^T (b\Sigma) (a-\mu)] \right]$$

Now this is the expectation of a χ^2 RV.

$\mathbb{E}[Y] = D$ where $Y \sim \chi^2_D$

and D is dimension of μ

$$= (x-a)^T \mathbb{E}_{\Sigma} [\Sigma] (x-a) + \frac{D}{b}$$

this is the mean of a Wishart distribution, here Ψ_V

$$= v (x-a)^T \Psi (x-a) + \frac{D}{b}$$

In our case ~~the~~, plugging this in yields

$$\begin{aligned} &\mathbb{E}_{\mu_k} \left[\mathbb{E}_{\mu_k|\Sigma_k} \left[-\frac{1}{2} (x_u - \mu_k)^T \Sigma_k (x_u - \mu_k) \right] \right] \sim \mathcal{L}_k \sim W(W_k, V_k) \\ &= -\frac{1}{2} \left[(x_u - \mu_k)^T W_k (x_u - \mu_k) V_k + \frac{D}{\beta_k} \right] \end{aligned}$$

$\mu_k | \Sigma_k \sim \mathcal{N}(\mu_k, \beta_k \Sigma_k)$
 $\sim \mathcal{N}(\mu_k, (\beta_k \Sigma_k)^{-1})$

If we go back to ① and plug in everything we arrive at

$$\mathbb{E}[\ln p(x|z, \mu, \Lambda)] = \frac{1}{z} \sum_{k=1}^K N_k \left\{ \ln \tilde{\Lambda}_k - D \beta_k^{-1} - v_k \text{Tr}(S_k W_k) - v_k (\bar{x}_k - m_k)^T W_k (\bar{x}_k - m_k) - D \cdot \ln(2\pi) \right\}$$

where

$$\ln \tilde{\Lambda}_k \equiv \mathbb{E}[\ln |\Lambda_k|] = \sum_{i=1}^D \psi\left(\frac{v_k + 1 - i}{z}\right) + D \ln z + \ln |W_k|$$

and $\beta_k, v_k, S_k, W_k, \bar{x}_k, m_k, N_k$ are defined as before

If we go back to ② we see that we have only accounted for the first term in the sum of expectations. The rest are given on pg's 481-482 in PRML. For completeness the remaining terms are:

$$\mathbb{E}_{\pi, z}[\ln p(z|\pi)] = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \ln \hat{\pi}_k$$

where $\ln \hat{\pi}_k = \mathbb{E}[\ln \pi_k] = \psi(\alpha_k) - \psi(\hat{\alpha})$
 where, as before $\alpha_k = \alpha_0 + N_k$ and $\hat{\alpha} = K\alpha_0 + N$

$$\mathbb{E}_{\pi}[\ln p(\pi)] = \ln C(\alpha_0) + (\alpha_0 - 1) \sum_{k=1}^K \ln \hat{\pi}_k$$

$$\mathbb{E}_{\mu, \Lambda}[\ln p(\mu, \Lambda)] = \frac{1}{z} \sum_{k=1}^K \left\{ D \ln\left(\frac{\beta_0}{z\pi}\right) + \ln \tilde{\Lambda}_k - \frac{D\beta_0}{\beta_0} - \beta_0 v_k (m_k - m_0)^T W_k (m_k - m_0) \right\}$$

$$+ K \ln B(W_0, v_0)$$

$$+ \frac{(v_0 + D - 1)}{z} \sum_{k=1}^K \ln \tilde{\Lambda}_k - \frac{1}{z} \sum_{k=1}^K v_k \text{Tr}(W_0^{-1} W_k)$$

$$\mathbb{E}[\ln q(\mathbf{z})] = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \ln r_{nk}$$

$$\mathbb{E}[\ln q(\boldsymbol{\pi})] = \sum_{k=1}^K (\alpha_k - 1) \ln \hat{\pi}_k + \ln C(\boldsymbol{\alpha})$$

$$\mathbb{E}[\ln q(\boldsymbol{\mu}, \boldsymbol{\Lambda})] = \sum_{k=1}^K \left\{ \frac{1}{2} \ln \hat{\Sigma}_k + \frac{D}{2} \ln \left(\frac{\beta_k}{2\pi} \right) - \frac{D}{2} - H[q(\boldsymbol{\mu}_k)] \right\}$$

The $H \approx B$ terms can be looked up in the PRML appendix.

Summary: the variational lower bound is a lower bound on the evidence of the data under the model.

- a) This can be useful for model comparison
- b) Is useful for debugging

Overall summary:

VB re-estimation equations can be cycled through, factor by factor. In optimizing each factor we arrive at parameters of distributions necessary to optimize other factors.

Other Uses of Variational approximation to the posterior distribution:

Prediction

Let \hat{x} be a new data point, what is the predictive distribution (posterior) for \hat{x}

$$p(\hat{x}|X) = \sum_{\mathbf{z}} \int \int p(\hat{x}|\mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) p(\mathbf{z}|\boldsymbol{\pi}) \underbrace{p(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}|X)}_{\text{posterior}} d\boldsymbol{\pi} d\boldsymbol{\mu} d\boldsymbol{\Lambda}$$

Remember, $q(\boldsymbol{\theta}) = p(\boldsymbol{\theta}|X)$ is the approximate posterior distribution so we can plug it into expressions like this and derive an approximate posterior predictive

For instance

$$p(\hat{x}|X) = \sum_{k=1}^K \iiint \pi_k \mathcal{N}(\hat{x} | \mu_k, \Sigma_k^{-1}) q(\pi) q(\mu_k, \Sigma_k) d\pi d\mu_k d\Sigma_k$$

where several integrals have been performed implicitly to get to this step.

The remaining integrals can be performed analytically to arrive at a mixture of Student-t dist's.

$$p(\hat{x}|X) = \frac{1}{Z} \sum_{k=1}^K \alpha_k \text{St}(\hat{x} | \mu_k, L_k, \nu_k + 1 - D)$$

where $L_k = \frac{(\nu_k + 1 - D) \beta_k}{(1 + \beta_k)} W_k$

Last note: learning the # of components. The prior $\alpha \sim \pi$ allows effective control over the "sparsity" of the model.

Induced factorizations

Factorization of approximating dist. factorized further as a consequence of the conditional independencies implied by embodied by the graphical model.

Computationally important to utilize these additional factorizations (space and time)