

T.61.5140 Machine Learning: Advanced Probabilistic Methods

Hollmén, Raiko (Spring 2008)

Problem session, 18th of April, 2008

<http://www.cis.hut.fi/Opinnot/T-61.5140/>

On notation:

Kullback-Leibler divergence can be written in two ways:

$$\text{KL}(q(x) \parallel p(x)) = E_{q(x)} \left\{ \ln \frac{q(x)}{p(x)} \right\} \quad (1)$$

$$= -E_{q(x)} \left\{ \ln \frac{p(x)}{q(x)} \right\} \quad (2)$$

Tapani has used the former and Bishop's book uses the latter. The lower bound $\mathcal{L}(q)$ on page 463 is the same as $-C_{\text{vb}}$. Maximization of $\mathcal{L}(q)$ is thus the same thing as minimization of C_{vb} . Also, Bishop absorbs the parameters θ into \mathbf{Z} when describing variational inference. This makes some of the formulas nicer, but the VB-EM algorithm cannot be described easily.

Let us use the shorthand $\langle \cdot \rangle = E_q\{\cdot\}$ to denote the expectation over the posterior approximation q .

1. Consider the probabilistic principal component analysis (PPCA) model

$$\mathbf{x}(t) = \mathbf{A}s(t) + \mathbf{n}(t) \quad (3)$$

with 3-dimensional data $\mathbf{x}(t)$, 1-dimensional source $s(t)$, and Gaussian noise $p(n_i(t)) = N(n_i(t) | 0, \sigma_i^2)$. We would like to estimate the noise level σ_i^2 for each data dimension $i = 1, 2, 3$ using the maximum likelihood estimator. What happens when the weight matrix \mathbf{A} goes to $\mathbf{A} = (1 \ 0 \ 0)^T$, and the source $s(t)$ copies the first dimension of the data: $s(t) = x_1(t)$?

Solution:

First, let us compute the reconstruction error, or the noise term

$$\mathbf{n}(t) = \mathbf{x}(t) - \mathbf{A}s(t) = \begin{pmatrix} 0 \\ x_2(t) \\ x_3(t) \end{pmatrix}, \quad (4)$$

which hints that there might be a problem with the variance parameters σ_1^2 . The likelihood of the data is

$$p(\mathbf{x}(t) | \mathbf{A}, s(t)) = \prod_{t=1}^T \prod_{i=1}^3 \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp \left[-\frac{(x_i(t) - \mathbf{A}_i s(t))^2}{2\sigma_i^2} \right] \quad (5)$$

$$= \prod_{t=1}^T \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp \left[-\frac{0}{2\sigma_1^2} \right] \prod_{i=2}^3 \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp \left[-\frac{(x_i(t) - \mathbf{A}_i s(t))^2}{2\sigma_i^2} \right] \quad (6)$$

$$= \text{const} \frac{1}{\sigma_1} \quad (7)$$

so we notice that when σ_1 goes to zero, the likelihood goes to infinity.

Maximum likelihood criterion considers this degenerate solution infinitely good, while it is in fact useless. This is an extreme case of overfitting, which is related to the fact that the posterior distribution has a high but narrow peak at the mentioned solution. Using any methodology that is sensitive to probability mass rather than density avoids this problem.

2. Variational Bayesian cost function C_{vb} is often a sum of simple terms, one for each variable. Consider the variable s whose model is $p(s | m, v) = N(s | m, \exp(-v))$. Assume that s, m, v are independent a posteriori, that is, $q(s, m, v) = q(s)q(m)q(v)$ and that $q(s) = N(s | \bar{s}, \tilde{s})$. Show that the term of the cost for s , $C_{\text{vb}}(s) = E_q \left\{ \ln \frac{q(s)}{p(s|m,v)} \right\}$, is

$$C_{\text{vb}}(s) = \frac{1}{2} \left(E_q \{ \exp v \} \left[(\bar{s} - E_q \{ m \})^2 + \text{Var}_q \{ m \} + \tilde{s} \right] - E_q \{ v \} + \ln(2\pi) \right) - \frac{1}{2} \ln(2\pi e \tilde{s}). \quad (8)$$

Solution:

Let us first show that $\langle mv \rangle = \langle m \rangle \langle v \rangle$ when $q(m, v) = q(m)q(v)$.

$$\langle mv \rangle = \int \int q(m, v) m v \, dm \, dv \quad (9)$$

$$= \int q(m) m \, dm \int q(v) v \, dv \quad (10)$$

$$= \langle m \rangle \langle v \rangle \quad (11)$$

Let us address the cost function $C_{vb}(s) = \langle \ln q(s) \rangle - \langle \ln p(s | m, v) \rangle$ in two parts.

$$- \langle \ln p(s | m, v) \rangle = - \langle \ln N(s | m, \exp(-v)) \rangle \quad (12)$$

$$= - \left\langle \ln \left[(2\pi \exp(-v))^{-1/2} \exp \frac{-(s-m)^2}{2 \exp(-v)} \right] \right\rangle \quad (13)$$

$$= - \left\langle \ln [2\pi \exp(-v)]^{-1/2} \right\rangle - \left\langle \ln \exp \left[-\frac{1}{2}(s-m)^2 \exp v \right] \right\rangle \quad (14)$$

$$= \frac{1}{2} \ln 2\pi - \frac{1}{2} \langle v \rangle + \frac{1}{2} \langle (s-m)^2 \rangle \langle \exp v \rangle. \quad (15)$$

We still need to compute the expectation of $(s-m)^2$

$$\langle (s-m)^2 \rangle = \langle s^2 - 2sm + m^2 \rangle \quad (16)$$

$$= \langle s^2 \rangle - 2 \langle sm \rangle + \langle m^2 \rangle \quad (17)$$

$$= \bar{s}^2 + \tilde{s} - 2\bar{s} \langle m \rangle + \langle m \rangle^2 + \text{Var} \{m\} \quad (18)$$

$$= (\bar{s} - \langle m \rangle)^2 + \tilde{s} + \text{Var} \{m\}. \quad (19)$$

Substituting $\langle (s-m)^2 \rangle$ back in, we get:

$$- \langle \ln p(s | m, v) \rangle = \frac{1}{2} \left\{ \langle \exp v \rangle \left[(\bar{s} - \langle m \rangle)^2 + \text{Var} \{m\} + \tilde{s} \right] - \langle v \rangle + \ln 2\pi \right\}. \quad (20)$$

The second part of the cost function is

$$\langle \ln q(s) \rangle = \langle \ln N(s | \bar{s}, \tilde{s}) \rangle \quad (21)$$

$$= \left\langle \ln \left[(2\pi \tilde{s})^{-1/2} \exp \frac{-(s-\bar{s})^2}{2\tilde{s}} \right] \right\rangle \quad (22)$$

$$= -\frac{1}{2} \ln 2\pi \tilde{s} + \frac{-\langle (s-\bar{s})^2 \rangle}{2\tilde{s}} \quad (23)$$

$$= -\frac{1}{2} \ln 2\pi \tilde{s} + \frac{-\tilde{s}}{2\tilde{s}} \quad (24)$$

$$= -\frac{1}{2} \ln 2\pi e \tilde{s}. \quad (25)$$

3. Consider the following model:

$$p(x(t) | m, v) = N(x(t) | m, \exp(-v)) \quad (26)$$

$$p(m) = N(m | 0, \exp(5)) \quad (27)$$

$$p(v) = N(v | 0, \exp(5)), \quad (28)$$

where $x(t), t = 1, \dots, T$ are the observed data and m and v are latent variables. Use the posterior approximation $q(m, v) = q(m)q(v) = N(m | \bar{m}, \tilde{m})N(v | \bar{v}, \tilde{v})$. Assuming that $q(v)$ is fixed, find $q(m)$ that minimizes $C_{vb} = E_q \left\{ \ln \frac{q(m, v)}{p(\{x(t)\}_{t=1}^T, m, v)} \right\}$.

Solution:

Let us first write the whole cost function.

$$C_{vb} = \left\langle \ln \frac{q(m, v)}{p(\{x(t)\}_{t=1}^T, m, v)} \right\rangle \quad (29)$$

$$= \langle \ln q(m) \rangle + \langle \ln q(v) \rangle + \sum_{t=1}^T$$

$$- \langle \ln p(x(t) | m, v) \rangle - \langle \ln p(m) \rangle - \langle \ln p(v) \rangle \quad (30)$$

$$= -\frac{1}{2} \ln(2\pi e \tilde{m}) - \frac{1}{2} \ln(2\pi e \tilde{v})$$

$$+ \sum_{t=1}^T \frac{1}{2} \left\{ \langle \exp v \rangle \left[(x(t) - \bar{m})^2 + \tilde{m} \right] - \langle v \rangle + \ln 2\pi \right\}$$

$$+ \frac{1}{2} \left[\langle \exp 5 \rangle (\bar{m}^2 + \tilde{m}) - \langle -5 \rangle + \ln 2\pi \right]$$

$$+ \frac{1}{2} \left[\langle \exp 5 \rangle (\bar{v}^2 + \tilde{v}) - \langle -5 \rangle + \ln 2\pi \right] \quad (31)$$

Now if we only leave the terms that depend on $q(m)$ we have

$$C_{vb} = \text{const} - \frac{1}{2} \ln(2\pi e \tilde{m}) + \sum_{t=1}^T \frac{1}{2} \langle \exp v \rangle \left[(x(t) - \bar{m})^2 + \tilde{m} \right] \\ + \frac{1}{2} \langle \exp 5 \rangle (\bar{m}^2 + \tilde{m}). \quad (32)$$

The minimum of this can be found at the zero of the gradient with respect to mean \bar{m} and the variance \tilde{m} :

$$\frac{\partial C_{vb}}{\partial \bar{m}} = \sum_{t=1}^T \langle \exp v \rangle (\bar{m} - x(t)) + \exp(-5)\bar{m} = 0 \quad (33)$$

$$\bar{m} = \frac{\sum_{t=1}^T x(t)}{\frac{\exp(-5)}{\exp v} + T} \quad (34)$$

$$\frac{\partial C_{vb}}{\partial \tilde{m}} = -\frac{1}{2\tilde{m}} + \sum_{t=1}^T \frac{1}{2} \langle \exp v \rangle + \frac{1}{2} \exp(-5) = 0 \quad (35)$$

$$\tilde{m} = \frac{1}{\exp(-5) + T \exp v}. \quad (36)$$

We can note that the mean \bar{m} is close to the average of the data, only the prior makes it go slightly towards zero. Also we note that the posterior variance \tilde{m} is the same as prior variance in case that there are no observations, or $T = 0$, and with each observation, the uncertainty (or variance) decreases. It is interesting that the actual data values $x(t)$ do not affect \tilde{m} .

The fourth problem will be solved next week.