



The Abdus Salam  
**International Centre  
for Theoretical Physics**



**2361-8**

**School on Large Scale Problems in Machine Learning and Workshop  
on Common Concepts in Machine Learning and Statistical Physics**

*20 - 31 August 2012*

**Large Scale Variational Bayesian Inference for Continuous Variable Models -  
Solutions to Exercises**

Matthias SEEGER

*Laboratory for Probabilistic Machine Learning, EPFL, CH-1015 Lausanne  
Switzerland*

ICTP School on Large Scale Problems in Machine Learning:  
Large Scale Variational Bayesian Inference for Continuous  
Variable Models — Solutions to Exercises

Matthias Seeger  
Probabilistic Machine Learning Laboratory  
Ecole Polytechnique Fédérale de Lausanne  
INR 112, Station 14, CH-1015 Lausanne  
*matthias.seeger@epfl.ch*

August 21, 2012

## 1 Super-Gaussian Bounding for Laplace Potentials

(b): If  $\pi_j = \gamma_j^{-1}$ , then

$$\frac{\partial -2 \log Z_Q}{\partial \pi_j} = \frac{-2}{Z_Q} \int P(\mathbf{y}|\mathbf{u}) \left( \prod_{i=1}^q e^{-\frac{1}{2}s_i^2/\gamma_i} \right) (-s_j^2/2) d\mathbf{u} = \mathbb{E}_Q[s_j^2].$$

(c):

$$\gamma_j \leftarrow \tau^{-1} \sqrt{\mathbb{E}_Q[s_j^2]}$$

(d): The equation in (c) is really coupled, since  $\mathbb{E}_Q[s_j^2]$  depends on  $\gamma_j$  as well, but appears on the right hand side only. The complete minimization can be done by iterating the fixed point equation. The marginal  $Q(s_j|\mathbf{y})$  does not have to be recomputed during this minimization, since we can always write

$$Q(s_j|\mathbf{y})' \propto Q(s_j|\mathbf{y}) e^{-\frac{1}{2}(\Delta\pi_j)s_j^2}, \quad \Delta\pi_j = \pi_j' - \pi_j.$$

## 2 Gaussian KL Minimization and Super-Gaussian Bounding

Obviously,

$$\phi_{\text{KL}}(\boldsymbol{\gamma}_*, \mathbf{0}) \leq \phi_{\text{SG}}(\boldsymbol{\gamma}_*)$$

implies the statement. By definition of super-Gaussianity,

$$-\log t_j(s_j) \leq s_j^2/(2\gamma_{*j}) + h_j(\gamma_{*j})/2,$$

so that

$$2\mathbb{E}_Q[-\log t_j(s_j) - s_j^2/(2\gamma_{*j})] \leq h_j(\gamma_{*j}).$$

### 3 Efficient Parameterization of Gaussian KL Minimization

(a) Immediate, given that

$$\nu_j = 2\mathbb{E}_Q[-\log t_j(s_j)]$$

and  $\mathbb{E}_Q[\mathbf{u}\mathbf{u}^T] = \Sigma + \boldsymbol{\mu}\boldsymbol{\mu}^T$ .

(b):

$$\nabla_{\Sigma}\phi = -\Sigma^{-1} + \mathbf{E} + \sum_{j=1}^q \pi_j \mathbf{b}_j \mathbf{b}_j^T.$$

Setting this equal to zero:

$$\Sigma_*^{-1} = \mathbf{E} + \mathbf{B}^T(\text{diag } \boldsymbol{\pi})\mathbf{B}.$$

### 4 Coordinate Update Algorithm for Gaussian KL Minimization

(a): Using the hint:

$$\Delta\pi_j = \frac{1}{\rho'_j} - \frac{1}{\rho_j}.$$

(b):

$$-p'_j + \mathbf{E}_{jj} + \frac{\partial \nu'_j}{\partial \rho'_j} = 0.$$

(c): This is just (a) in reverse. First, (b) implies that

$$\pi'_j = \frac{\partial \nu'_j}{\partial \rho'_j}.$$

Then,

$$\frac{1}{\rho'_j} = \Delta\pi_j + \frac{1}{\rho_j} \quad \Rightarrow \quad \rho'_j = \frac{\rho_j}{1 + (\Delta\pi_j)\rho_j}.$$

### 5 Spectral Analysis of Conjugate Gradients Algorithm

First,

$$P(\mathbf{A}) = P(\mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^T) = \sum_{j=0}^{k-1} \alpha_j (\mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^T)^j = \sum_{j=0}^{k-1} \alpha_j \mathbf{Q}\boldsymbol{\Lambda}^j \mathbf{Q}^T = \mathbf{Q}P(\boldsymbol{\Lambda})\mathbf{Q}^T,$$

because  $\mathbf{Q}^T\mathbf{Q} = \mathbf{I}$ . Then,

$$q(\mathbf{x}) = (1/2)\mathbf{x}^T \mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^T \mathbf{x} - \mathbf{b}^T \mathbf{Q}\mathbf{Q}^T \mathbf{x} = (1/2)\mathbf{y}^T \boldsymbol{\Lambda} \mathbf{y} - \bar{\mathbf{b}}^T \mathbf{y}.$$

Since  $\mathbf{A}^{-1} = \mathbf{Q}\boldsymbol{\Lambda}^{-1}\mathbf{Q}^T$ , we have that  $q_* = -(1/2)\mathbf{b}^T \mathbf{Q}\boldsymbol{\Lambda}^{-1}\mathbf{Q}^T \mathbf{b} = -(1/2)\bar{\mathbf{b}}^T \boldsymbol{\Lambda}^{-1} \bar{\mathbf{b}}$ .

Next,  $\mathbf{x}_k \in \mathcal{K}_k$ , which is spanned by  $\mathbf{A}^j \mathbf{b}$  for  $j < k$ . This means that  $\mathbf{x}_k = P_k(\mathbf{A})\mathbf{b}$  for some polynomial with  $\deg(P_k) < k$ . Therefore,

$$\mathbf{y}_k = \mathbf{Q}^T P_k(\mathbf{A})\mathbf{b} = P_k(\mathbf{\Lambda})\mathbf{Q}^T \mathbf{b} = P_k(\mathbf{\Lambda})\bar{\mathbf{b}}.$$

Using the solutions from above,

$$\begin{aligned} q(\mathbf{x}_k) - q_* &= (1/2) \min_{P_k} \sum_{i=1}^n (\lambda_i y_{k,i}^2 - \bar{b}_i y_{k,i} + \bar{b}_i^2 / \lambda_i) = (1/2) \min_{P_k} \sum_{i=1}^n \bar{b}_i^2 (\lambda_i P_k(\lambda_i)^2 - P_k(\lambda_i) + 1/\lambda_i) \\ &= (1/2) \min_{P_k} \sum_{i=1}^n (\bar{b}_i^2 / \lambda_i) (\lambda_i P_k(\lambda_i) - 1)^2. \end{aligned}$$

Let  $Q_k(t) := tP_k(t) - 1$ . As  $P_k$  runs over polynomials of degree  $< k$ ,  $Q_k$  runs over polynomials of degree  $\leq k$  with  $Q_k(0) = -1$ . The bound is nondecreasing in the  $Q_k(\lambda_i)^2$  (which is why we can also use  $-Q_k$  for the argument). Without assumptions on  $\mathbf{b}$ , we have to strive for small  $|Q_k(\lambda_i)|$ , especially for the smaller  $\lambda_i$ .

If  $\{\lambda_1, \dots, \lambda_n\} = \{\kappa_1, \dots, \kappa_k\}$ , pick the polynomial  $Q_k(t) = [\prod_{j=1}^k (t - \kappa_j)] / [\prod_j \kappa_j]$ . Here,  $\prod_j \kappa_j > 0$  because all  $\kappa_j > 0$  ( $\mathbf{A}$  is positive definite). Then,  $|Q_k(0)| = 1$  and  $Q_k(\lambda_i) = 0$  for all  $i$ , so that  $q(\mathbf{x}_k) = q_*$ . Since  $q$  is strictly convex, it has a unique minimum, so  $\mathbf{x}_k = \mathbf{x}_*$ .

## 6 Super-Gaussian Bounding for Bernoulli Potentials

(a):

$$\frac{1}{1 + e^{-ys}} = \frac{e^{ys/2}}{e^{ys/2} + e^{-ys/2}},$$

so that

$$b = y/2, \quad \tilde{t}(s) = \frac{1}{2 \cosh(bs)}.$$

$\tilde{t}(s)$  is even. Using the hint,  $\log \tilde{t}(s) = -\log \cosh(bs) - \log 2$  is a convex function of  $x = s^2$ . This means that  $\tilde{t}(s)$ , and therefore  $t(s)$ , are super-Gaussian.

(b): It suffices to show that  $f(x) = \log(1 + e^x)$  is convex.

$$f'(x) = \frac{e^x}{1 + e^x} = \frac{1}{1 + e^{-x}}, \quad f''(x) = \frac{e^{-x}}{(1 + e^{-x})^2} = f'(x)f'(-x) > 0,$$

so that  $f(x)$  is strictly convex. This means that  $t(s)$  is log-concave, and so is  $\tilde{t}(s)$ , since  $\log \tilde{t}(s) = \log t(s) - bs$ .

(c): We follow the argument given in the course. The MAP problem for the setup here would be

$$\min_{\mathbf{u}_*} \sigma^{-2} \|\mathbf{u}_*\|^2 - 2 \sum_{j=1}^q \log t(s_{*j}), \quad \mathbf{s}_* = \mathbf{B}\mathbf{u}_*.$$

Now,  $t(s_{*j}) = e^{bs_{*j}} \tilde{t}(s_{*j})$ . In the course, we showed that for an even potential, the IL problem differs from MAP in that  $\tilde{t}(s_{*j})$  is replaced by  $\tilde{t}((z_j + s_{*j})^{1/2})$ . Here,

$$\log t(s_{*j}) = bs_{*j} + \log \tilde{t}(s_{*j}) \rightarrow bs_{*j} + \log \tilde{t}((z_j + s_{*j}^2)^{1/2}),$$

so the IL problem is

$$\min_{\mathbf{u}_*} \sigma^{-2} \|\mathbf{u}_*\|^2 - 2 \sum_{j=1}^q \left( b s_{*j} + \log \tilde{t}((z_j + s_{*j}^2)^{1/2}) \right).$$

From (b),  $\log \tilde{t}(s_{*j})$  is log-concave, which implies the convexity of  $h(\gamma_j)$  and therefore of  $\tilde{t}((z_j + s_{*j}^2)^{1/2})$  (this result is quoted in the course and proved in [1]), and  $b s_{*j}$  is linear, therefore convex.

## 7 Proximal Map for Inner Loop Optimization Problem

(a): If

$$s' = \text{prox}(r) = \underset{s}{\text{argmin}} f(s; r),$$

then  $f(s; -r) = f(-s; r)$ , so that  $\text{prox}(-r) = -\text{prox}(r)$ . If  $r > 0$  and  $s < 0$ , then  $f(-s; r) < f(s; r)$ , since  $(-s - r)^2 < (s - r)^2$ . Therefore,  $\text{prox}(r) \geq 0$ .

(b): Recall that  $s \geq 0$ . Define  $y = (1 + s^2)^{1/2}$ , so that

$$f(s; r) = \kappa y + \frac{1}{2}(s - r)^2.$$

The stationary equation is  $df/ds = 0$ :

$$\frac{\kappa s'}{y'} = r - s'.$$

$s' \leq r$  follows from  $s'/y' \geq 0$ . Also,  $r > 0$  implies  $s' > 0$ . Moreover,  $r = s'(1 + \kappa/y') < s'(1 + \kappa)$ , since  $y' > 1$ , so that  $s' > r/(1 + \kappa)$ . Finally,  $y' > |s'|$ , so that  $s'/y' < 1$  and  $r - s' < \kappa$ . These inequalities can be used to bracket a solution for  $s'$ .

(c): Squaring both sides of the stationary equation gives

$$\kappa^2 s^2 = (1 + s^2)(r - s)^2 \Leftrightarrow s^4 - 2rs^3 + (r^2 + 1 - \kappa^2)s^2 - 2rs + r^2 = 0.$$

## 8 Bound on Marginal Variances

First,

$$\mathbf{A} = \sigma^{-2} \mathbf{X}^T \mathbf{X} + \mathbf{B}^T \mathbf{\Gamma}^{-1} \mathbf{B}.$$

We have that  $\text{Var}_Q[s_j | \mathbf{y}] = \mathbf{b}_j^T \mathbf{A}^{-1} \mathbf{b}_j$ , where  $\mathbf{b}_j = \mathbf{B}^T \boldsymbol{\delta}_j$  is the  $j$ -th row of  $\mathbf{B}$ . Then,

$$\begin{aligned} \mathbf{b}_j^T \mathbf{A}^{-1} \mathbf{b}_j &= \max_{\mathbf{x}} 2\mathbf{b}_j^T \mathbf{x} - \mathbf{x}^T \mathbf{A} \mathbf{x} = \max_{\mathbf{x}} 2\boldsymbol{\delta}_j^T \mathbf{B} \mathbf{x} - \sigma^{-2} \|\mathbf{X} \mathbf{x}\|^2 - (\mathbf{B} \mathbf{x})^T \mathbf{\Gamma}^{-1} (\mathbf{B} \mathbf{x}) \\ &\leq \max_{\mathbf{x}} 2\boldsymbol{\delta}_j^T (\mathbf{B} \mathbf{x}) - (\mathbf{B} \mathbf{x})^T \mathbf{\Gamma}^{-1} (\mathbf{B} \mathbf{x}) \leq \max_{\mathbf{w}} 2\boldsymbol{\delta}_j^T \mathbf{w} - \mathbf{w}^T \mathbf{\Gamma}^{-1} \mathbf{w} = \boldsymbol{\delta}_j^T \mathbf{\Gamma} \boldsymbol{\delta}_j = \gamma_j. \end{aligned}$$

The first  $\leq$  is due to  $\|\mathbf{X} \mathbf{x}\|^2 \geq 0$ , the second due to the fact that  $\mathbf{B} \mathbf{x}$  runs over a subspace of  $\mathbf{w} \in \mathbb{R}^q$ . The first and last  $=$  are applications of the identity provided in the hint.

## References

- [1] M. Seeger and H. Nickisch. Large scale Bayesian inference and experimental design for sparse linear models. *SIAM Journal of Imaging Sciences*, 4(1):166–199, 2011.