

---

# ***Variational Algorithms: Ising Model, Univariate Gaussian, Linear Regression, Model Selection***

*Prof. Nicholas Zabaras*

*Center for Informatics and Computational Science*

*<https://cics.nd.edu/>*

*University of Notre Dame*

*Notre Dame, IN, USA*

*Email: [nzabaras@gmail.com](mailto:nzabaras@gmail.com)*

*URL: <https://www.zabaras.com/>*

*March 29, 2018*



# Contents

---

- ❑ [Mean Field for the Ising Model](#), [Structured Mean Field](#)
- ❑ [Variational Inference for the univariate Gaussian](#), [Variational optimization and model selection](#).
- ❑ [Variational Linear Regression](#), [Predictive Distribution](#), [Lower Bound](#), [Selection of the order of the polynomial](#)
- ❑ [Variational Linear Regression with  \$\text{Gam}\(\beta|c\_0, d\_0\)\$](#)

Following:

- [Pattern Recognition and Machine Learning](#), Christopher M. Bishop, Chapter 10
- [Machine Learning: A Probabilistic Perspective](#), Kevin Murphy, Chapter 21.



# Mean Field for the Ising Model

Consider image denoising, where  $x_i \in \{-1, +1\}$  are the hidden pixel values of the clean image. We have a joint model of the form

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y}|\mathbf{x})$$

where the prior has the form

$$p(\mathbf{x}) = \frac{1}{Z_0} \exp(-E_0(\mathbf{x})), \quad E_0(\mathbf{x}) = - \sum_{i=1}^D \sum_{j \in \text{nbr}_i} W_{ij} x_i x_j$$

and the likelihood has the form

$$p(\mathbf{y}|\mathbf{x}) = \prod_i p(y_i|x_i) = \exp\left(\sum L_i(x_i)\right)$$

Therefore the posterior has the form

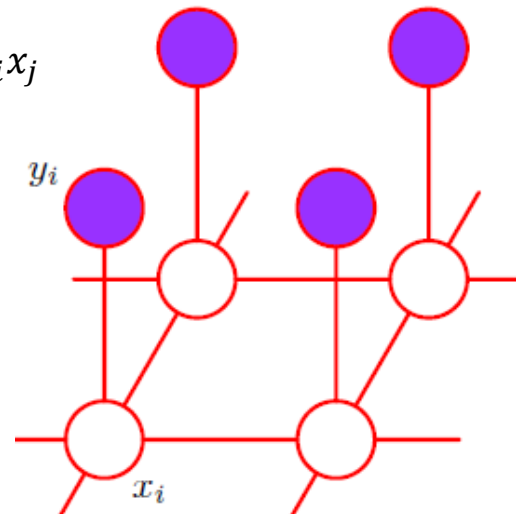
$$p(\mathbf{x}|\mathbf{y}) = 1/Z \exp(-E(\mathbf{x})), \quad E(\mathbf{x}) = E_0(\mathbf{x}) - \sum_i L_i(x_i)$$

We will now approximate this by a fully factored approximation

$$q(\mathbf{x}) = \prod_i q(x_i, \mu_i)$$

where  $\mu_i$  is the mean value of node  $i$ . To derive the update for the variational parameter  $\mu_i$ , we first write out  $\ln \tilde{p}(\mathbf{x}) = \ln p(\mathbf{x}, \mathbf{y}) = -E(\mathbf{x})$ , and dropping terms that do not involve  $x_i$ :

$$\ln \tilde{p}(\mathbf{x}) = x_i \sum_{j \in \text{nbr}_i} W_{ij} x_j + L_i(x_i) + \text{const}$$



# Mean Field for the Ising Model

$$\ln \tilde{p}(\mathbf{x}) = x_i \sum_{j \in \text{nbr}_i} W_{ij} x_j + L_i(x_i) + \text{const}$$

This only depends on the states of the neighboring nodes. Now we take expectations of this wrt  $\prod_{j \neq i} q_j(x_j)$  to get

$$q_i(x_i) \propto \exp \left( x_i \sum_{j \in \text{nbr}_i} W_{ij} \mu_j + L_i(x_i) \right)$$

Thus we replace the states of the neighbors by their average values. Let  $m_i = \sum_{j \in \text{nbr}_i} W_{ij} \mu_j$  be the mean field influence on node  $i$ . Also let  $L_i^+ \equiv L_i(+1)$ ,  $L_i^- \equiv L_i(-1)$ . The approximate marginal posterior is given by

$$q_i(x_i = 1) = \frac{e^{m_i + L_i^+}}{e^{m_i + L_i^+} + e^{-m_i + L_i^-}} = \frac{1}{1 + e^{-2m_i + L_i^- - L_i^+}} = \text{sigm}(2a_i), \quad a_i \equiv m_i + 0.5(L_i^+ - L_i^-)$$

Similarly  $q_i(x_i = -1) = \text{sigm}(-2a_i)$  and thus:  $\mu_i = \mathbb{E}_{q_i}[x_i] = q_i(x_i = 1)(+1) + q_i(x_i =$



# Mean Field for the Ising Model

---

$$\mu_i = \tanh\left(\sum_{j \in \text{nbr}_i} W_{ij} \mu_j + 0.5(L_i^+ - L_i^-)\right)$$

We can turn the above Eqs in to a fixed point algorithm by writing:

$$\mu_i^t = \tanh\left(\sum_{j \in \text{nbr}_i} W_{ij} \mu_j^{t-1} + 0.5(L_i^+ - L_i^-)\right)$$

To avoid checkerboarding effects, damped updates are often used as follows:

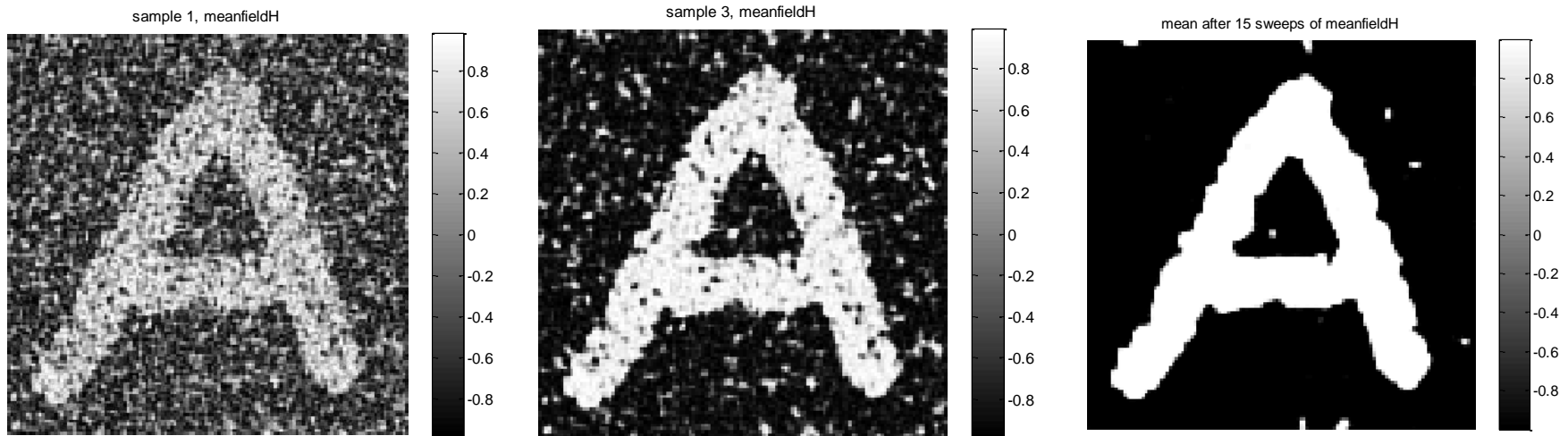
$$\mu_i^t = (1 - \lambda)\mu_i^{t-1} + \lambda \tanh\left(\sum_{j \in \text{nbr}_i} W_{ij} \mu_j^{t-1} + 0.5(L_i^+ - L_i^-)\right), 0 < \lambda < 1.$$

The nodes can be updated in parallel or asynchronously.



# Mean Field for the Ising Model

- Example of image denoising using mean field
- Parallel updates and  $\lambda = 0.5$ .
- Ising prior with  $W_{ij} = 1$  and a Gaussian noise model with  $\sigma = 2$ .
- Results are shown after 1, 3 and 15 iterations.



Run [isingImageDenoiseDemo](#) in the [PMTK3 toolbox](#)

# Structured Mean Field Approach

Assuming that all the variables are independent in the posterior is a very strong assumption that can lead to poor results.

Sometimes we can exploit **tractable substructure** in our problem to efficiently handle particular dependencies. This is called the **structured mean field approach**.

We group sets of variables together, and we update them simultaneously. We treat the variables in the  $i$ 'th group as a single “mega-variable” and we follow the earlier variational derivation.

If we can perform efficient inference in each  $q_i$ , the method is tractable overall.

A factorial HMM example is discussed next.

- Saul, L. and M. Jordan (1995). [Exploiting tractable substructures in intractable networks](#). In *NIPS*, Volume 8.
- Ghahramani, Z. and M. Jordan (1997). [Factorial hidden Markov models](#). *Machine Learning* 29, 245–273.
- Bouchard-Cote, A. and M. Jordan (2009). [Optimization of structured mean field objectives](#). In *UAI*.



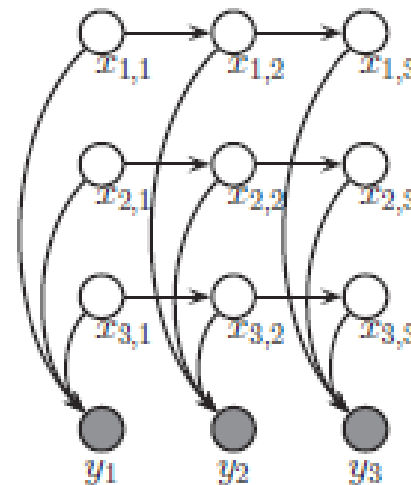
# Factorial HMM

Consider the factorial HMM model. Suppose there are  $M$  chains, each of length  $T$ , and suppose each hidden node has  $K$  states. The model is defined as follows

$$p(\mathbf{x}, \mathbf{y}) = \prod_m \prod_t p(x_{tm} | x_{t-1,m}) p(\mathbf{y}_t | \mathbf{x}_{tm})$$

where  $p(x_{tm} = k | x_{t-1,m} = j) = A_{mjk}$  is an entry in the **transition matrix for chain  $m$** ,  $p(x_{1m} = k | x_{0m}) = p(x_{1m} = k) = \pi_{mk}$ , **is the initial state distribution for chain  $m$** , and

$$p(\mathbf{y}_t | \mathbf{x}_t) = \mathcal{N}\left(\mathbf{y}_t \mid \sum_{m=1}^M \mathbf{W}_m \mathbf{x}_{tm}, \boldsymbol{\Sigma}\right)$$



is the observation model, where  $\mathbf{x}_{tm}$  is a 1-of- $K$  encoding of  $x_{tm}$  and  $\mathbf{W}_m$  is a  $D \times K$  matrix (assuming  $\mathbf{y}_t \in \mathbb{R}^D$ ).

Even though each chain is a priori independent, they become coupled in the posterior due to having an observed common child,  $\mathbf{y}_t$ . The junction tree algorithm applied to this graph takes  $\mathcal{O}(TMK^{M+1})$  time.

The structured mean field algorithm discussed next takes  $\mathcal{O}(TMK^2I)$  time, where  $I$  is the number of mean field iterations ( $I \sim 10$  for good performance).





# Factorial HMM

We can write the exact posterior as:

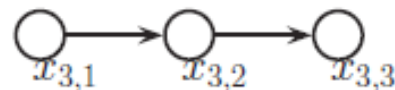
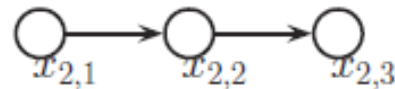
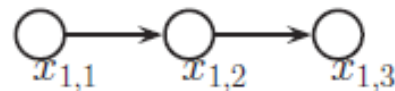
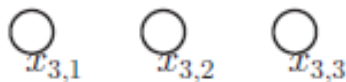
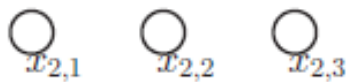
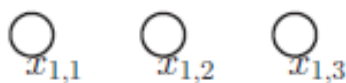
$$p(\mathbf{x}|\mathbf{y}) = \frac{1}{Z} \exp(-E(\mathbf{x}, \mathbf{y}))$$

$$E(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \sum_{t=1}^T \left( \mathbf{y}_t - \sum_m \mathbf{W}_m \mathbf{x}_{tm} \right)^T \boldsymbol{\Sigma}^{-1} \left( \mathbf{y}_t - \sum_m \mathbf{W}_m \mathbf{x}_{tm} \right) - \sum_m \mathbf{x}_{1m}^T \tilde{\boldsymbol{\pi}}_m - \sum_{t=2}^T \sum_m \mathbf{x}_{tm}^T \tilde{\mathbf{A}}_m \mathbf{x}_{t-1,m}$$

where pointwise  $\tilde{\mathbf{A}}_m = \ln \mathbf{A}_m$ , and  $\tilde{\boldsymbol{\pi}}_m = \ln \boldsymbol{\pi}_m$ .  $\mathbf{x}_{tm}$  is a vector with 1 in one location and zero elsewhere.

We can approximate the posterior as a product of marginals (see Fig on the left)

But a better approximation is to [use a product of chains](#) (Fig on the right) . Each chain can be tractably updated individually using the forwards-backwards algorithm.



# Factorial HMM

We assume:

$$q(\mathbf{x}|\mathbf{y}) = \frac{1}{Z_q} \prod_{m=1}^M \left( q(\mathbf{x}_{1m}|\xi_{1m}) \prod_{t=2}^T q(\mathbf{x}_{tm}|\mathbf{x}_{t-1,m}, \xi_{tm}) \right)$$

$$q(\mathbf{x}_{1m}|\xi_{1m}) = \prod_{k=1}^K (\xi_{1mk} \pi_{mk})^{x_{1mk}}$$

$$q(\mathbf{x}_{tm}|\mathbf{x}_{t-1,m}, \xi_{tm}) = \prod_{k=1}^K \left( \xi_{tmk} \prod_{j=1}^K (A_{mjk})^{x_{t-1,m,j}} \right)^{x_{tmk}}$$

We see that the parameters  $\xi_{1m}$  play the role of an approximate local evidence averaging out the effects of the other chains. Comparing with  $p(\mathbf{x}, \mathbf{y}) = \prod_m \prod_t p(\mathbf{x}_{tm}|\mathbf{x}_{t-1,m}) p(\mathbf{y}_t|\mathbf{x}_{tm})$ , the  $K \times 1$  vector  $\xi_{tm}$  plays the role of the probability of an observation  $p(\mathbf{y}_t|\mathbf{x}_{tm})$  for each of the  $K$  settings of  $\mathbf{x}_{tm}$ .

This is in contrast to the exact local evidence which coupled all the chains together. We can rewrite the approximate posterior as  $q(\mathbf{x}) = \frac{1}{Z_q} \exp(-E_q(\mathbf{x}))$  where

$$E_q(\mathbf{x}) = - \sum_{t=1}^T \sum_{m=1}^M \mathbf{x}_{tm}^T \tilde{\xi}_{tm} - \sum_{m=1}^M \mathbf{x}_{1m}^T \tilde{\pi}_m - \sum_{t=2}^T \sum_{m=1}^M \mathbf{x}_{tm}^T \tilde{A}_m \mathbf{x}_{t-1,m}, \quad \tilde{\xi}_{tm} = \ln \xi_{tm}$$



# Factorial HMM

$$q(\mathbf{x}|\mathbf{y}) = \frac{1}{Z_q} \prod_{m=1}^M q(\mathbf{x}_{1m}|\xi_{1m}) \prod_{t=2}^T q(\mathbf{x}_{tm}|\mathbf{x}_{t-1,m}, \xi_{tm}), \quad q(\mathbf{x}_{1m}|\xi_{1m}) = \prod_{k=1}^K (\xi_{1mk} \pi_{mk})^{x_{1mk}}$$

$$q(\mathbf{x}_{tm}|\mathbf{x}_{t-1,m}, \xi_{tm}) = \prod_{k=1}^K (\xi_{tmk} \prod_{j=1}^K (A_{mj k})^{x_{t-1,mj}})^{x_{tmk}}$$

$$E_q(\mathbf{x}) = - \sum_{t=1}^T \sum_{m=1}^M \mathbf{x}_{tm}^T \tilde{\xi}_{tm} - \sum_{m=1}^M \mathbf{x}_{1m}^T \tilde{\pi}_m - \sum_{t=2}^T \sum_{m=1}^M \mathbf{x}_{tm}^T \tilde{A}_m \mathbf{x}_{t-1,m}, \quad \tilde{\xi}_{tm} = \ln \xi_{tm}$$

This has the same temporal factors as the exact posterior but the local evidence is different. The objective function  $KL(q||p) = - \int q \ln \frac{p}{q} d\mathbf{x}$  is given by:

$$KL(q||p) = \int q \ln q d\mathbf{x} - \int q \ln p d\mathbf{x} = \int q \ln \frac{e^{-E_q}}{Z_q} d\mathbf{x} - \int q \ln \frac{e^{-E}}{Z} d\mathbf{x} = \mathbb{E}_q[E] - \mathbb{E}_q[E_q] - \ln Z_q + \ln Z$$

Let  $\bar{\mathbf{x}}_{tm} = \mathbb{E}_q[\mathbf{x}_{tm}]$ . Using,  $E(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \sum_{t=1}^T (\mathbf{y}_t - \sum_m \mathbf{W}_m \mathbf{x}_{tm})^T \Sigma^{-1} (\mathbf{y}_t - \sum_m \mathbf{W}_m \mathbf{x}_{tm}) - \sum_m \mathbf{x}_{1m}^T \tilde{\pi}_m - \sum_{t=2}^T \sum_m \mathbf{x}_{1m}^T \tilde{A}_m \mathbf{x}_{t-1,m}$ , we can write:

$$KL(q||p) = \sum_{m=1}^M \sum_{t=1}^T \bar{\mathbf{x}}_{tm}^T \tilde{\xi}_{tm} + \frac{1}{2} \sum_{t=1}^T \left[ \mathbf{y}_t^T \Sigma^{-1} \mathbf{y}_t - 2 \sum_{m=1}^M \mathbf{y}_t^T \Sigma^{-1} \mathbf{W}_m \bar{\mathbf{x}}_{tm} + \sum_{m=1}^M \sum_{n \neq m} \text{tr}(\mathbf{W}_m^T \Sigma^{-1} \mathbf{W}_n \bar{\mathbf{x}}_{tn} \bar{\mathbf{x}}_{tm}^T) + \sum_{m=1}^M \text{tr}(\mathbf{W}_m^T \Sigma^{-1} \mathbf{W}_m \text{diag}(\bar{\mathbf{x}}_{tm})) \right] - \ln Z_q + \ln Z$$

$\text{diag}(\cdot)$  is an operator that takes a vector and returns a square matrix with the elements of the vector along its diagonal.



# Factorial HMM

$$KL(q||p) = \sum_{m=1}^M \sum_{t=1}^T \bar{\mathbf{x}}_{tm}^T \tilde{\boldsymbol{\xi}}_{tm} + \frac{1}{2} \sum_{t=1}^T \left[ \mathbf{y}_t^T \boldsymbol{\Sigma}^{-1} \mathbf{y}_t - 2 \sum_{m=1}^M \mathbf{y}_t^T \boldsymbol{\Sigma}^{-1} \mathbf{W}_m \bar{\mathbf{x}}_{tm} + \sum_{m=1}^M \sum_{n \neq m} \text{tr}(\mathbf{W}_m^T \boldsymbol{\Sigma}^{-1} \mathbf{W}_n \bar{\mathbf{x}}_{tn} \bar{\mathbf{x}}_{tm}^T) + \sum_{m=1}^M \text{tr}(\mathbf{W}_m^T \boldsymbol{\Sigma}^{-1} \mathbf{W}_m \text{diag}(\bar{\mathbf{x}}_{tm})) \right] - \ln Z_q + \ln Z$$

Using  $E_q(\mathbf{x}) = \sum_{t=1}^T \sum_{m=1}^M \mathbf{x}_{tm}^T \tilde{\boldsymbol{\xi}}_{tm} - \sum_{m=1}^M \mathbf{x}_{1m}^T \tilde{\boldsymbol{\pi}}_m - \sum_{t=2}^T \sum_{m=1}^M \mathbf{x}_{tm}^T \tilde{\mathbf{A}}_m \mathbf{x}_{t-1,m}$ , and since

$$\frac{\partial Z_q}{\partial \tilde{\boldsymbol{\xi}}_{\tau n}} = \int -e^{-E_q} \frac{\partial E_q}{\partial \tilde{\boldsymbol{\xi}}_{\tau n}} d\mathbf{x} = \int e^{-E_q} \mathbf{x}_{\tau n} d\mathbf{x} = Z_q \bar{\mathbf{x}}_{\tau n}, \quad \frac{\partial \ln Z_q}{\partial \tilde{\boldsymbol{\xi}}_{\tau n}} = \bar{\mathbf{x}}_{\tau n}$$

we can write:

$$\frac{\partial KL}{\partial \tilde{\boldsymbol{\xi}}_{\tau n}} = \bar{\mathbf{x}}_{\tau n} + \left( \frac{\partial \bar{\mathbf{x}}_{tm}}{\partial \tilde{\boldsymbol{\xi}}_{\tau n}} \right)^T \sum_{t=1}^T \sum_{m=1}^M \left[ \tilde{\boldsymbol{\xi}}_{tm} - \mathbf{W}_m^T \boldsymbol{\Sigma}^{-1} \mathbf{y}_t + \sum_{\ell \neq m}^M \mathbf{W}_m^T \boldsymbol{\Sigma}^{-1} \mathbf{W}_\ell \bar{\mathbf{x}}_{t\ell} + \frac{1}{2} \delta_m \right] - \bar{\mathbf{x}}_{\tau n} = 0$$

Thus the term in parenthesis is zero and the update Eqs are as follows:

$$\tilde{\boldsymbol{\xi}}_{tm} = \exp \left( \mathbf{W}_m^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{y}}_{tm} - \frac{1}{2} \delta_m \right), \quad \delta_m = \text{diag}(\mathbf{W}_m^T \boldsymbol{\Sigma}^{-1} \mathbf{W}_m), \quad \tilde{\mathbf{y}}_{tm} = \mathbf{y}_t - \sum_{\ell \neq m}^M \mathbf{W}_\ell \bar{\mathbf{x}}_{t\ell}$$

$\delta_m$  is the vector of diagonal elements of  $\mathbf{W}_m^T \boldsymbol{\Sigma}^{-1} \mathbf{W}_m$ .



# Factorial HMM

---

The  $\xi_{tm}$  parameter plays the role of the local evidence, averaging over the neighboring chains.

Having computed this for each chain, we can perform forwards-backwards in parallel, using these approximate local evidence terms to compute  $q(\mathbf{x}_{t,m} | \mathbf{y}_{1:T})$  for each  $m$  and  $t$ .

The update cost is  $\mathcal{O}(TMK^2)$  for a full “sweep” over all the variational parameters, since we have to run forwards-backwards  $M$  times, for each chain independently.

This is the same cost as a fully factorized approximation, but is much more accurate.



# Variational Bayes (VB) for the Univariate Gaussian

Suppose our goal is to infer the posterior distribution for the mean  $\mu$  and precision  $\tau$  given data  $\mathcal{D} = \{x_1, \dots, x_n\}$  assumed to be drawn independently from the Gaussian.

The likelihood function is given by:

$$p(\mathcal{D}|\mu, \tau) = \left(\frac{\tau}{2\pi}\right)^{N/2} \exp\left\{-\frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2\right\}$$

We introduce conjugate prior distributions for  $\mu$  and  $\tau$  given by

$$\begin{aligned} p(\mu|\tau) &= \mathcal{N}(\mu|\mu_0, (\lambda_0\tau)^{-1}) \\ p(\tau) &= \text{Gam}(\tau|a_0, b_0) \end{aligned}$$

Since we chose conjugate priors, [this can be solved analytically](#). However, we consider a factorized approximation of the form:

$$q(\mu, \tau) = q_\mu(\mu)q_\tau(\tau)$$

- MacKay, D. J. C. (2003). [Information Theory, Inference and Learning Algorithms](#). Cambridge University Press.
- MacKay, D. (1995a). [Developments in probabilistic modeling with neural networks — ensemble learning](#). In *Proc. 3rd Ann. Symp. Neural Networks*.
- Attias, H. (2000). [A variational Bayesian framework for graphical models](#). In *NIPS-12*.
- Beal, M. and Z. Ghahramani (2006). [Variational Bayesian Learning of Directed Graphical Models with Hidden Variables](#). *Bayesian Analysis* 1(4).
- Smidl, V. and A. Quinn (2005). [The Variational Bayes Method in Signal Processing](#). Springer.



# Computing $q_\mu^*(\mu)$

Computing  $\ln q_j^*(Z_j) = \mathbb{E}_{i \neq j} [\ln p(X, Z)] + \text{const}$ , we get:

$$\begin{aligned}\ln q_\mu^*(\mu) &= \mathbb{E}_\tau [\ln p(\mathcal{D}|\mu, \tau) + \ln p(\mu|\tau)] + \text{const} \\ &= -\frac{\mathbb{E}[\tau]}{2} \{ \lambda_0 (\mu - \mu_0)^2 + \sum_{n=1}^N (x_n - \mu)^2 \} + \text{const}\end{aligned}$$

By completing the square over  $\mu$  we see that  $q_\mu(\mu) = \mathcal{N}(\mu|\mu_N, \lambda_N^{-1})$  with

$$\mu_N = \frac{\lambda_0 \mu_0 + N \bar{x}}{\lambda_0 + N} \text{ (no need updating)}$$

$$\lambda_N = (\lambda_0 + N) \mathbb{E}[\tau]$$

Note that as  $N \rightarrow \infty$  the MLE estimate is recovered and the precision is infinite:

$$\begin{aligned}\mu_N &= \frac{\lambda_0 \mu_0 + N \bar{x}}{\lambda_0 + N} \xrightarrow{N \rightarrow \infty} \bar{x} \\ \lambda_N &= (\lambda_0 + N) \mathbb{E}[\tau] \xrightarrow{N \rightarrow \infty} \infty\end{aligned}$$



# Computing $q_{\tau}^*(\tau)$

Similarly:

$$\begin{aligned}\ln q_{\tau}^*(\tau) &= \mathbb{E}_{\mu}[\ln p(\mathcal{D}|\mu, \tau) + \ln p(\mu|\tau)] + \ln p(\tau) + \text{const} \\ &= (a_0 - 1) \ln \tau - b_0 \tau + \frac{1}{2} \ln \tau + \frac{N}{2} \ln \tau - \frac{\tau}{2} \mathbb{E}_{\mu} \left[ \sum_{n=1}^N (x_n - \mu)^2 + \lambda_0 (\mu - \mu_0)^2 \right] + \text{const}\end{aligned}$$

It can then be shown that  $q_{\tau}(\tau) = \text{Gam}(\tau | a_N, b_N)$  with

$$a_N = a_0 + \frac{N + 1}{2}, \quad b_N = b_0 + \frac{1}{2} \mathbb{E}_{\mu} \left[ \sum_{n=1}^N (x_n - \mu)^2 + \lambda_0 (\mu - \mu_0)^2 \right]$$

Using the expressions for the [mean and variance of the Gamma](#) distribution, we note:

$$\begin{aligned}\mathbb{E}[\tau] &= \frac{a_N}{b_N} = \frac{2a_0 + N + 1}{2b_0 + \mathbb{E} \left[ \lambda_0 (\mu - \mu_0)^2 + \sum_{n=1}^N (x_n - \mu)^2 \right]} \xrightarrow{N \rightarrow \infty} \frac{N}{\mathbb{E} \left[ \sum_{n=1}^N (x_n - \mu)^2 \right]} \\ \text{var}[\tau] &= \frac{a_N}{b_N^2} = \frac{\mathbb{E}[\tau]}{b_0 + \frac{1}{2} \mathbb{E} \left[ \lambda_0 (\mu - \mu_0)^2 + \sum_{n=1}^N (x_n - \mu)^2 \right]} \xrightarrow{N \rightarrow \infty} 0\end{aligned}$$





# Self-consistent iterative process

$$q_\mu(\mu) = \mathcal{N}(\mu | \mu_N, \lambda_N^{-1}), \quad \mu_N = \frac{\lambda_0 \mu_0 + N \bar{x}}{\lambda_0 + N}, \quad \lambda_N = (\lambda_0 + N) \mathbb{E}[\tau]$$

$$q_\tau(\tau) = \text{Gam}(\tau | a_N, b_N),$$

$$a_N = a_0 + \frac{N + 1}{2}, \quad b_N = b_0 + \frac{1}{2} \mathbb{E}_\mu \left[ \sum_{n=1}^N (x_n - \mu)^2 + \lambda_0 (\mu - \mu_0)^2 \right]$$

The eqs for the two optimal distributions  $q_\mu(\mu)$ ,  $q_\tau(\tau)$  can be iterated until convergence.

For example after computing  $q_\mu(\mu)$ , we use the moments  $\mathbb{E}_\mu[\mu]$ ,  $\mathbb{E}_\mu[\mu^2]$  for updating  $q_\tau(\tau)$ , etc.

# Using non-informative priors

Assume non-informative priors:  $\mu_0 = a_0 = b_0 = \lambda_0 = 0$ . The mean and the precision of the optimal  $q_\mu(\mu)$  are simplified leading to:

$$\mathbb{E}[\mu] = \bar{x}, \quad \mathbb{E}[\mu^2] = \bar{x}^2 + \frac{1}{N\mathbb{E}[\tau]}$$

Substituting these equations in the expression for  $\mathbb{E}[\tau]$  gives

$$\frac{1}{\mathbb{E}[\tau]} = \frac{1}{N+1} \mathbb{E} \left[ \sum_{n=1}^N (x_n - \mu)^2 \right] = \frac{N}{N+1} \left( \overline{x^2} - 2\bar{x}^2 + \underbrace{\mathbb{E}[\mu^2]}_{\bar{x}^2 + \frac{1}{N\mathbb{E}[\tau]}} \right) \Rightarrow \frac{1}{\mathbb{E}[\tau]} = \overline{x^2} - \bar{x}^2$$

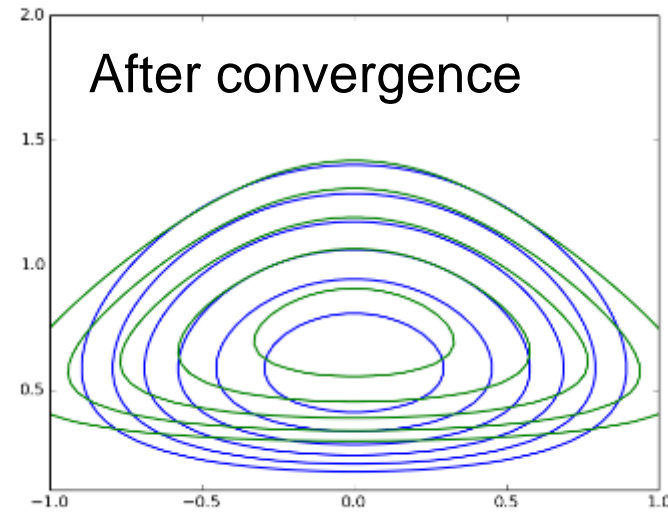
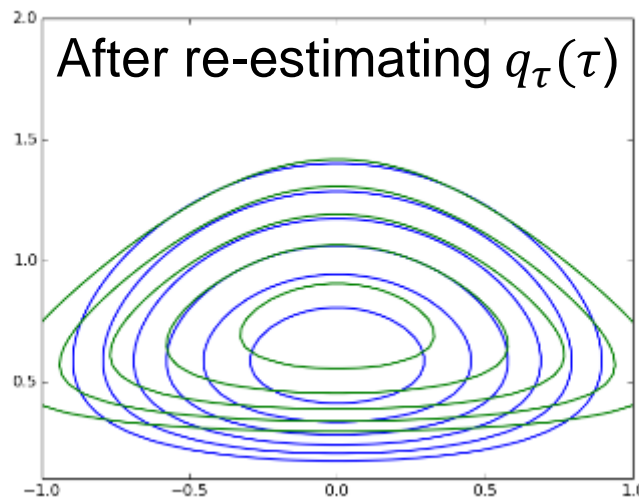
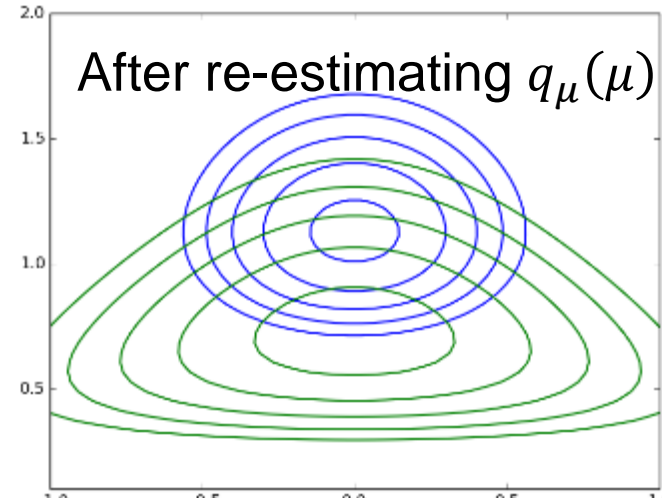
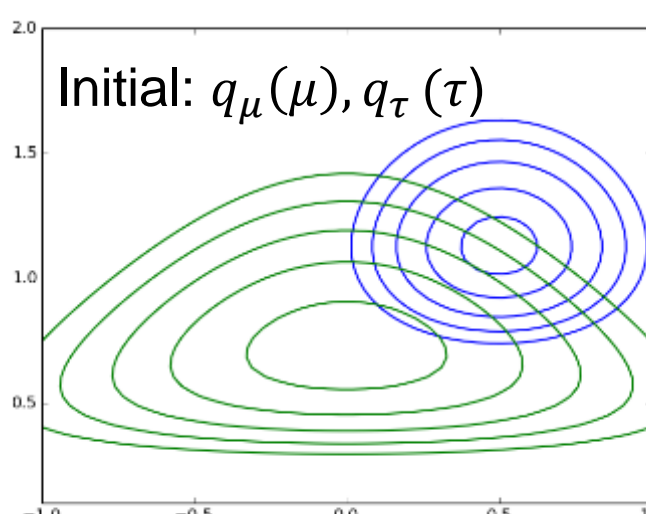
$$\frac{1}{\mathbb{E}[\tau]} = \overline{x^2} - \bar{x}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})^2 \text{ (sample variance)}$$

For a Bayesian treatment for the Gaussian including a comparison with MLE see the reference below.

- [Thomas P. Minka, Bayesian inference of a uniform distribution](#), 2001 (Msft Research)



# ***The univariate Gaussian: Example***



[Python Code](#)



# The Univariate Gaussian: Lower Bound

It is a good practice to compute the lower bound and monitor its values during the iterations (it should always monotonically increase).

$$\mathcal{L}(q) = \iint q(\mu, \tau) \ln \frac{p(\mathcal{D}, \mu, \tau)}{q(\mu, \tau)} d\mu d\tau = \mathbb{E}[\ln p(\mathcal{D}|\mu, \tau)] + \mathbb{E}[\ln p(\mu|\tau)] + \mathbb{E}[\ln p(\tau)] - \mathbb{E}[\ln q(\mu)] - \mathbb{E}[\ln q(\tau)]$$

We compute the various terms below using the posteriors we arrived at with the variational approach.

$$-\mathbb{E}[\ln q(\mu)] = -\frac{1}{2} \ln \lambda_N + \frac{1}{2} (1 + \ln(2\pi)) \text{ (entropy of } \mathcal{N}(\mu|\mu_N, \lambda_N^{-1}) \text{)}$$

$$-\mathbb{E}[\ln q(\tau)] = \ln \Gamma(a_N) - (a_N - 1) \psi(a_N) - \ln b_N + a_N \text{ (entropy of } \text{Gam}(\tau|a_N, b_N) \text{)}$$

$$\mathbb{E}_{q(\tau)}[\ln p(\tau)] = \mathbb{E}_{q(\tau)}[a_0 \ln b_0 - \ln \Gamma(a_0) + (a_0 - 1) \ln \tau - b_0 \tau] = a_0 \ln b_0 - \ln \Gamma(a_0) + (a_0 - 1)(\psi(a_N) - \ln b_N) - b_0 \frac{a_N}{b_N} \text{ (recall for } \tau \sim \text{Gamma}(a, b): \mathbb{E}_{q(\tau)}[\ln \tau] = \psi(a) - \ln b, \mathbb{E}_{q(\tau)}[\tau] = \frac{a}{b} \text{)}.$$

$$\begin{aligned} \mathbb{E}_{q(\mu, \tau)}[\ln p(\mathcal{D}|\mu, \tau)] &= \mathbb{E}_{q(\mu, \tau)} \left[ -\frac{N}{2} \ln(2\pi) + \frac{N}{2} \ln \tau - \frac{1}{2} \tau \sum_{n=1}^N (x_n - \mu)^2 \right] = -\frac{N}{2} \ln(2\pi) + \frac{N}{2} \mathbb{E}[\ln \tau] - \\ &\frac{1}{2} \mathbb{E}[\tau] \mathbb{E}[\sum_{n=1}^N (x_n - \mu)^2] = -\frac{N}{2} \ln(2\pi) + \frac{N}{2} (\psi(a_N) - \ln b_N) - \frac{N}{2} \frac{a_N}{b_N} (\overline{x^2} - 2\bar{x}\mathbb{E}[\mu] + \mathbb{E}[\mu^2]) = \\ &-\frac{N}{2} \ln(2\pi) + \frac{N}{2} (\psi(a_N) - \ln b_N) - \frac{N}{2} \frac{a_N}{b_N} \left( \overline{x^2} - 2\bar{x}\mu_N + \mu_N^2 + \frac{1}{\lambda_N} \right) \end{aligned}$$

$$\begin{aligned} \mathbb{E}_{q(\mu, \tau)}[\ln p(\mu|\tau)] &= \mathbb{E}_{q(\mu, \tau)} \left[ \frac{1}{2} \ln \left( \frac{\lambda_0}{2\pi} \right) + \frac{1}{2} \ln \tau - \frac{1}{2} \lambda_0 \tau (\mu - \mu_0)^2 \right] = \frac{1}{2} \ln \left( \frac{\lambda_0}{2\pi} \right) + \frac{1}{2} (\psi(a_N) - \ln b_N) - \\ &\frac{\lambda_0}{2} \frac{a_N}{b_N} \left( (\mu_N - \mu_0)^2 + \frac{1}{\lambda_N} \right) \end{aligned}$$



# The Univariate Gaussian: Lower Bound

Finally:

$$\begin{aligned}
 \mathcal{L}(q) = & -\frac{N}{2} \ln(2\pi) + \frac{N}{2} (\psi(a_N) - \ln b_N) - \frac{N}{2} \frac{a_N}{b_N} \left( \overline{x^2} - 2\bar{x}\mu_N + \mu_N^2 + \frac{1}{\lambda_N} \right) \\
 & + \frac{1}{2} \ln \left( \frac{\lambda_0}{2\pi} \right) + \frac{1}{2} (\psi(a_N) - \ln b_N) - \frac{\lambda_0}{2} \frac{a_N}{b_N} \left( (\mu_N - \mu_0)^2 + \frac{1}{\lambda_N} \right) \\
 & + a_0 \ln b_0 - \ln \Gamma(a_0) + (a_0 - 1) (\psi(a_N) - \ln b_N) - b_0 \frac{a_N}{b_N} \\
 & + \ln \Gamma(a_N) - (a_N - 1) \psi(a_N) - \ln b_N + a_N \\
 & - \frac{1}{2} \ln \lambda_N + \frac{1}{2} (1 + \ln(2\pi))
 \end{aligned}$$

Light blue terms are obviously constants.

The terms in black when combined and using ,  $b_N = b_0 + \frac{1}{2} \mathbb{E}_\mu [\sum_{n=1}^N (x_n - \mu)^2 + \lambda_0 (\mu - \mu_0)^2] = b_0 + \frac{1}{2} \left( \overline{x^2} N - 2\bar{x} N \mu_N + N \mu_N^2 + \frac{N}{\lambda_N} + \lambda_0 (\mu_N - \mu_0)^2 + \frac{\lambda_0}{\lambda_N} \right)$  lead to a constant.

Red terms contribute to the final expression below using  $a_N = a_0 + \frac{N+1}{2}$ . We conclude:

$$\mathcal{L}(q) = -\frac{1}{2} \ln \lambda_N + \ln \Gamma(a_N) - a_N \ln b_N + \text{const}$$



# Variational Optimization and Model Selection

Consider comparing a set of candidate models, labelled by  $m$ , and having prior probabilities  $p(m)$ .

Our goal is then to approximate the posterior probabilities  $p(m|\mathbf{X})$ , where  $\mathbf{X}$  is the observed data.

Different models may have different structure and different dimensionality for the hidden variables  $\mathbf{Z}$ . We factorize the joint posterior as follows:

$$q(\mathbf{Z}, m) = q(\mathbf{Z}|m)q(m)$$

Starting with the definition of the lower bound (see below)  $\mathcal{L}$ , we can show:

$$\begin{aligned}\mathcal{L} &= \sum_m \sum_{\mathbf{Z}} q(\mathbf{Z} | m) q(m) \ln \frac{p(\mathbf{Z}, \mathbf{X}, m)}{q(\mathbf{Z} | m) q(m)} = \sum_m \sum_{\mathbf{Z}} q(\mathbf{Z} | m) q(m) \ln \frac{p(\mathbf{Z}, m | \mathbf{X}) p(\mathbf{X})}{q(\mathbf{Z} | m) q(m)} \\ &= \sum_m \sum_{\mathbf{Z}} q(\mathbf{Z} | m) q(m) \ln \frac{p(\mathbf{Z}, m | \mathbf{X})}{q(\mathbf{Z} | m) q(m)} + \ln p(\mathbf{X})\end{aligned}$$

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



# Variational Optimization and Model Selection

We first optimize each  $p(\mathbf{Z}|m)$  by optimization of the lower bound of  $\mathcal{L}_m =$

$\sum_{\mathbf{Z}} q(\mathbf{Z}|m) \ln \frac{p(\mathbf{Z}, \mathbf{X}|m)}{q(\mathbf{Z}|m)}$  for each model  $m$ . Then note that the overall lower bound  $\mathcal{L}$  is:

$$\begin{aligned} \mathcal{L} &= \sum_m \sum_{\mathbf{Z}} q(\mathbf{Z} | m) q(m) \ln \frac{p(\mathbf{Z}, \mathbf{X}, m)}{q(\mathbf{Z} | m) q(m)} = \sum_m \sum_{\mathbf{Z}} q(\mathbf{Z} | m) q(m) (\ln p(\mathbf{Z}, \mathbf{X} | m) + \ln p(m) - \ln q(\mathbf{Z} | m) - \ln q(m)) \\ &= \sum_m q(m) \left( \ln p(m) - \ln q(m) + \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z} | m) (\ln p(\mathbf{Z}, \mathbf{X} | m) - \ln q(\mathbf{Z} | m))}_{\mathcal{L}_m} \right) = \\ &= \sum_m q(m) (\ln p(m) - \ln q(m) + \ln e^{\mathcal{L}_m}) = \sum_m q(m) \left( \ln \frac{p(m) e^{\mathcal{L}_m}}{q(m)} \right) \end{aligned}$$

This is as the distance  $-KL(q(m), p(m)e^{\mathcal{L}_m})$  which is maximized when  
 $q(m) \propto p(m) \exp(\mathcal{L}_m)$

$$\text{where: } \mathcal{L}_m = \sum_{\mathbf{Z}} q(\mathbf{Z} | m) \ln \frac{p(\mathbf{Z}, \mathbf{X} | m)}{q(\mathbf{Z} | m)} \quad (\text{lower bound for model } m)$$

After normalization, we can use  $q(m)$  for model selection or model averaging.



# Variational Linear Regression

Consider the Bayesian linear regression model.

Recall that the likelihood function is given by

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \boldsymbol{\phi}_n, \beta^{-1})$$

where  $\boldsymbol{\phi}_n = \boldsymbol{\phi}(\mathbf{x}_n)$ .

We use the following conjugate prior distributions (i.e. Gaussian-Gamma):

$$\begin{aligned} p(\mathbf{w}|\alpha) &= \mathcal{N}(\mathbf{w} | 0, \alpha^{-1} \mathbf{I}) \\ p(\alpha) &= \text{Gam}(\alpha | a_0, b_0) \end{aligned}$$

To simplify discussion the noise precision parameter  $\beta$  is fixed and assumed to be known. The framework can be extended with this assumption being relaxed.

- Drugowitsch, J. (2008). [Bayesian linear regression](#). Technical report, U. Rochester.



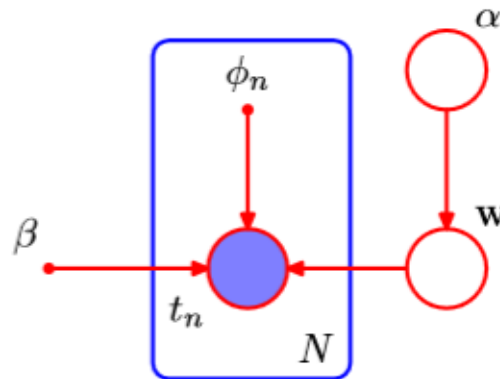


# Variational Linear Regression

Collectively we have the joint distribution

$$p(\mathbf{t}, \mathbf{w}, \alpha) = p(\mathbf{t}|\mathbf{w})p(\mathbf{w}|\alpha)p(\alpha)$$

which can be represented by the following graphical model



We assume a variational posterior of the following form to obtain the unknown hyperparameters

$$q(\mathbf{w}, \alpha) = q(\mathbf{w})q(\alpha)$$

# Variational Linear Regression

Using the same framework as previously discussed, we obtain for  $q(\alpha)$ :

$$\begin{aligned}\ln q^*(\alpha) &= \ln p(\alpha) + \mathbb{E}_{\mathbf{w}}[\ln p(\mathbf{w}|\alpha)] + \text{const} \\ &= (a_0 - 1) \ln \alpha - b_0 \alpha + \frac{M}{2} \ln \alpha - \frac{\alpha}{2} \mathbb{E}[\mathbf{w}^T \mathbf{w}] + \text{const}\end{aligned}$$

As expected this gives a [Gamma distribution](#)

$$q^*(\alpha) = \text{Gam}(\alpha|a_N, b_N)$$

with

$$\begin{aligned}a_N &= a_0 + \frac{M}{2} \\ b_N &= b_0 + \frac{1}{2} \mathbb{E}[\mathbf{w}^T \mathbf{w}]\end{aligned}$$



# Variational Linear Regression

Similarly we obtain for  $q(\mathbf{w})$ :

$$\begin{aligned}\ln q^*(\mathbf{w}) &= \ln p(\mathbf{t}|\mathbf{w}) + \mathbb{E}_\alpha[\ln p(\mathbf{w}|\alpha)] + \text{const} \\ &= -\frac{\beta}{2} \sum_{n=1}^N \{\mathbf{w}^T \boldsymbol{\phi}_n - t_n\}^2 - \frac{1}{2} \mathbb{E}[\alpha] \mathbf{w}^T \mathbf{w} + \text{const} \\ &= -\frac{1}{2} \mathbf{w}^T (\mathbb{E}[\alpha] \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}) \mathbf{w} + \beta \mathbf{w}^T \boldsymbol{\Phi}^T \mathbf{t} + \text{const}\end{aligned}$$

As expected this gives a normal distribution

$$q^*(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$$

with

$$\begin{aligned}\mathbf{m}_N &= \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t} \\ \mathbf{S}_N &= (\mathbb{E}[\alpha] \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1}\end{aligned}$$

Recall that  $\boldsymbol{\Phi}^T \boldsymbol{\Phi} = (\boldsymbol{\phi}_1 \quad \dots \quad \boldsymbol{\phi}_N) \begin{pmatrix} \boldsymbol{\phi}_1^T \\ \vdots \\ \boldsymbol{\phi}_N^T \end{pmatrix} = \sum_{n=1}^N \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T$

These are the same formulas obtained when  $\alpha$  is treated as constant. In the VI approach,  $\alpha$  is replaced by  $\mathbb{E}[\alpha]$ .



# Variational Linear Regression

By standard properties for Gaussian and [Gamma distributions](#) we finally have

$$\mathbb{E}[\alpha] = \frac{a_N}{b_N}$$
$$\mathbb{E}[\mathbf{w}\mathbf{w}^T] = \mathbf{m}_N \mathbf{m}_N^T + \mathbf{S}_N$$

which collectively give

$$\mathbb{E}[\alpha] = \frac{a_N}{b_N} = \frac{a_0 + M/2}{b_0 + \mathbb{E}[\mathbf{w}^T \mathbf{w}]/2}$$

For the case of  $a_0=b_0=0$  (infinitely broad prior over  $\alpha$ ):  $\mathbb{E}[\alpha] = \frac{M}{\mathbf{m}_N^T \mathbf{m}_N + \text{tr}(\mathbf{S}_N)}$ . This is similar to the result obtained using [the model evidence approximation](#).

It then remains a task to cycle between computing  $a_N$ ,  $b_N$ ,  $\mathbf{m}_N$  and  $\mathbf{S}_N$  until some convergence criterion is met.

These results are consistent with those obtained by [maximizing the evidence using EM](#) (except that the point estimate of  $\alpha$  is now replaced by  $\mathbb{E}[\alpha]$ ) and give identical results in the case of an infinitely broad prior for which  $\mathbb{E}[\alpha] =$

$$\frac{M}{\mathbf{m}_N^T \mathbf{m}_N + \text{tr}(\mathbf{S}_N)}.$$



# Predictive Distribution

The predictive distribution of  $t$  for a new input  $\mathbf{x}$  is given as:

$$\begin{aligned} p(t|\mathbf{x}, t) &= \int p(t|\mathbf{x}, \mathbf{w})p(\mathbf{w}|t)d\mathbf{w} \approx \int p(t|\mathbf{x}, \mathbf{w})q(\mathbf{w})d\mathbf{w} \\ &= \int \mathcal{N}(t|\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}), \beta^{-1})\mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N) d\mathbf{w} = \mathcal{N}(t|\mathbf{m}_N^T \boldsymbol{\phi}(\mathbf{x}), \sigma^2(\mathbf{x})) \end{aligned}$$

Here we used an earlier result for linear Gaussian models. Given a marginal Gaussian distribution for  $\mathbf{x}$  and a conditional Gaussian distribution for  $\mathbf{y}$  given  $\mathbf{x}$  in the form

$$\begin{aligned} p(\mathbf{x}) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{\Lambda}^{-1}) \\ p(\mathbf{y}|\mathbf{x}) &= \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}) \end{aligned}$$

the marginal distribution of  $\mathbf{y}$  is given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\mathbf{\Lambda}^{-1}\mathbf{A}^T)$$

Using the above result, the input dependent variance is given as:

$$\sigma^2(\mathbf{x}) = \frac{1}{\beta} + \boldsymbol{\phi}(\mathbf{x})^T \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x}), \text{ with } \mathbf{S}_N = (\mathbb{E}[\alpha]\mathbf{I} + \beta\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1}$$

This is the same form as that obtained when  $\alpha$  is treated as constant but now  $\alpha$  is replaced by  $\mathbb{E}[\alpha]$ .



# Lower Bound

But how can we select M (degree of polynomial)? We can compute the lower bound

$$\mathcal{L}(q) = \mathbb{E}[\ln p(\mathbf{w}, \alpha, \mathbf{t})] - \mathbb{E}[\ln q(\mathbf{w}, \alpha)]$$

$$= \mathbb{E}_{\mathbf{w}}[\ln p(\mathbf{t}|\mathbf{w})] + \mathbb{E}_{\mathbf{w}, \alpha}[\ln p(\mathbf{w}|\alpha)] + \mathbb{E}_{\alpha}[\ln p(\alpha)] - \mathbb{E}_{\mathbf{w}}[\ln q(\mathbf{w})] - \mathbb{E}_{\alpha}[\ln q(\alpha)]$$

where

$$\begin{aligned} \mathbb{E}_{\mathbf{w}}[\ln p(\mathbf{t}|\mathbf{w})] &= -\frac{N}{2} \ln(2\pi) + \frac{N}{2} \ln(\beta) - \frac{\beta}{2} \mathbb{E} \left[ \sum_{n=1}^N \{\mathbf{w}^T \boldsymbol{\phi}_n - t_n\}^2 \right] \\ &= -\frac{N}{2} \ln(2\pi) + \frac{N}{2} \ln(\beta) - \frac{\beta}{2} \{ \mathbf{t}^T \mathbf{t} - 2 \mathbb{E}[\mathbf{w}^T] \boldsymbol{\Phi}^T \mathbf{t} + \text{Tr}(\mathbb{E}[\mathbf{w} \mathbf{w}^T] \boldsymbol{\Phi}^T \boldsymbol{\Phi}) \} \\ &= \frac{N}{2} \ln \frac{\beta}{2\pi} - \frac{\beta}{2} \mathbf{t}^T \mathbf{t} + \beta \mathbf{m}_N^T \boldsymbol{\Phi}^T \mathbf{t} - \frac{\beta}{2} \text{Tr} [\boldsymbol{\Phi}^T \boldsymbol{\Phi} (\mathbf{m}_N \mathbf{m}_N^T + \mathbf{S}_N)] \end{aligned}$$

$$\begin{aligned} \mathbb{E}_{\mathbf{w}, \alpha}[\ln p(\mathbf{w}|\alpha)] &= -\frac{M}{2} \ln 2\pi + \frac{M}{2} \mathbb{E}[\ln \alpha] - \frac{\mathbb{E}[\alpha]}{2} \mathbb{E}[\mathbf{w}^T \mathbf{w}] \\ &= -\frac{M}{2} \ln 2\pi + \frac{M}{2} (\psi(a_N) - \ln b_N) - \frac{a_N}{2b_N} [\mathbf{m}_N^T \mathbf{m}_N + \text{Tr}(\mathbf{S}_N)] \end{aligned}$$

Use [expectation of the log of a Gamma distributed variable](#)

$$\begin{aligned} \mathbb{E}_{\alpha}[\ln p(\alpha)] &= a_0 \ln b_0 + (a_0 - 1) \mathbb{E}[\ln \alpha] - b_0 \mathbb{E}[\alpha] - \ln \Gamma(a_0) \\ &= a_0 \ln b_0 + (a_0 - 1) [\psi(a_N) - \ln b_N] - \frac{b_0 a_N}{b_N} - \ln \Gamma(a_0) \end{aligned}$$



# Lower Bound

The final two terms in  $\mathcal{L}(q)$  represent the entropies of the Gaussian and Gamma distributions:

$$-\mathbb{E}_{\mathbf{w}}[\ln q(\mathbf{w})] = \frac{1}{2} \ln |\mathbf{S}_N| + \frac{M}{2} [1 + \ln 2\pi]$$

$$-\mathbb{E}_{\alpha}[\ln q(\alpha)] = \ln \Gamma(a_N) - (a_N - 1)\psi(a_N) - \ln b_N + a_N$$

We substitute in the Eqs. above the following expressions for the moments:

$$\mathbb{E}[\mathbf{w}] = \mathbf{m}_N$$

$$\mathbb{E}[\mathbf{w}\mathbf{w}^T] = \mathbf{m}_N \mathbf{m}_N^T + \mathbf{S}_N$$

$$\mathbb{E}[\alpha] = \frac{a_N}{b_N}$$

$$\mathbb{E}[\ln \alpha] = \psi(a_N) - \ln b_N$$

to obtain the final expression for  $\mathcal{L}(q)$  :

$$\mathcal{L}(q) = \mathbb{E}_{\mathbf{w}}[\ln p(\mathbf{t}|\mathbf{w})] + \mathbb{E}_{\mathbf{w},\alpha}[\ln p(\mathbf{w}|\alpha)] + \mathbb{E}_{\alpha}[\ln p(\alpha)] - \mathbb{E}_{\mathbf{w}}[\ln q(\mathbf{w})] - \mathbb{E}[\ln q(\alpha)]$$



# Lower Bound Vs the Order of the Polynomial

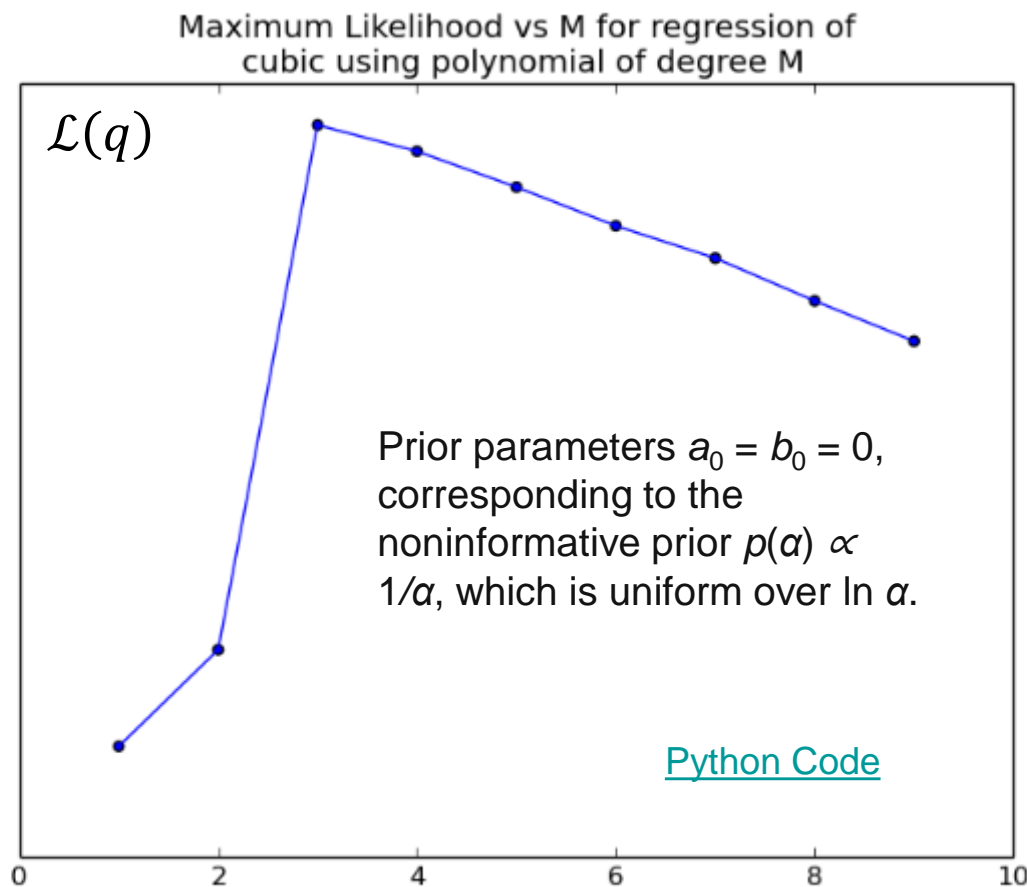
10 data points have been generated from a polynomial of degree 3 over  $[-5, 5]$  with additive Gaussian noise. The lower bound is maximized for  $M=3$  corresponding to the true model form which the data was generated.

$\mathcal{L}$  represents lower bound on the log marginal likelihood  $\ln p(\mathbf{t}/M)$  for the model.

If we assign equal prior probabilities  $p(M)$  to the different values of  $M$ , then we can interpret  $\mathcal{L}$  as an approximation to  $p(M/\mathbf{t})$ .

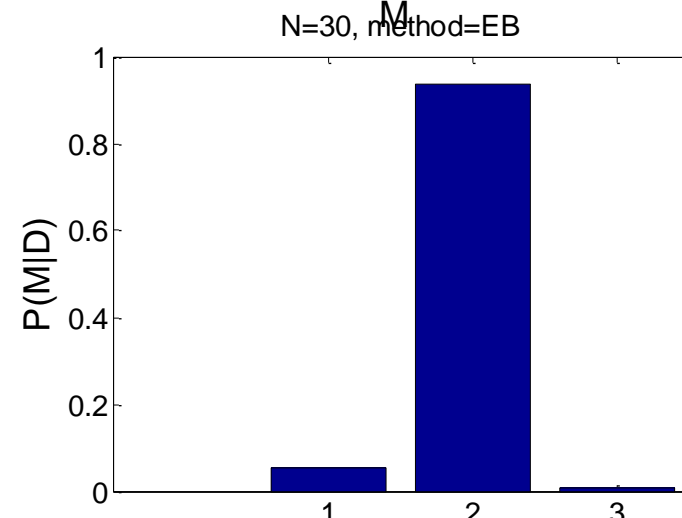
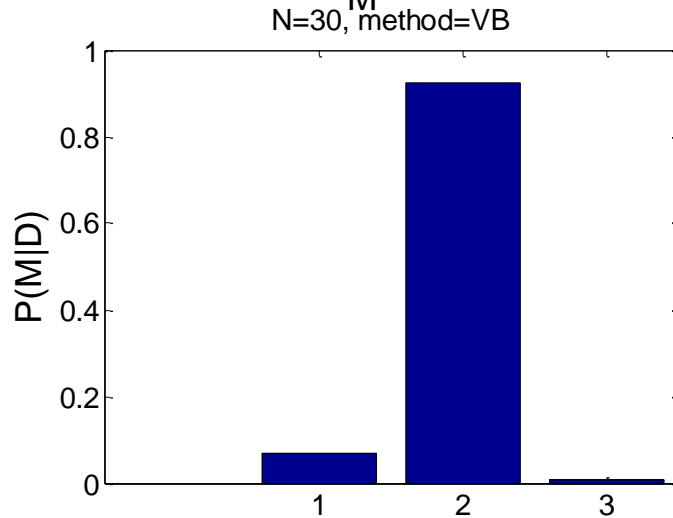
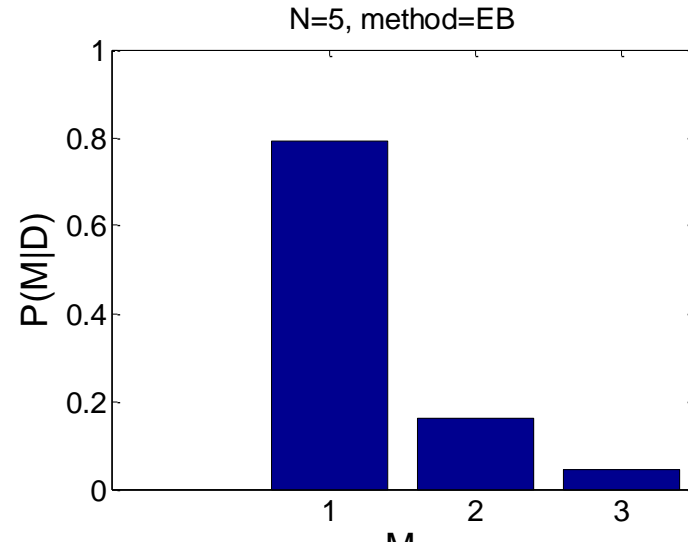
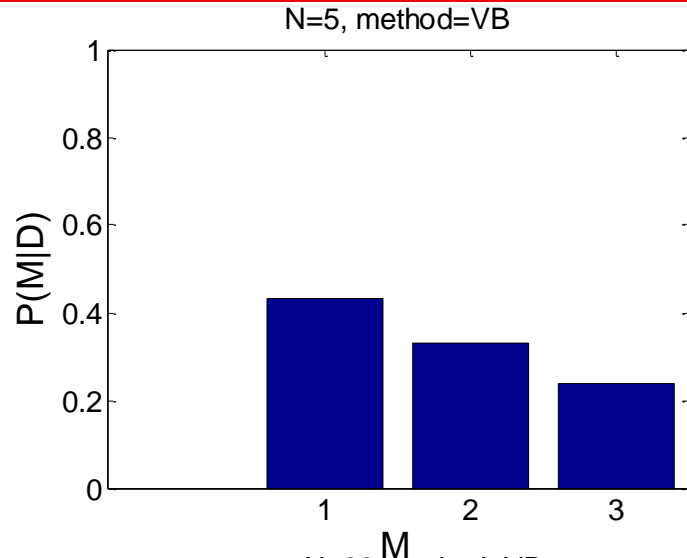
Thus the variational framework assigns the highest probability to the model with  $M = 3$ .

This should be contrasted with the MLE result, which assigns ever smaller residual error to models of increasing complexity until the residual error is driven to zero, causing MLE to over-fitted models.





# Lower Bound Vs the Order of the Polynomial



Run [linregEbModelSelVsN](#) in the [PMTK3 toolbox](#)



# Variational Linear Regression with $\text{Gam}(\beta|c_0, d_0)$

We now extend the variational treatment of Bayesian linear regression to include a gamma hyperprior  $\text{Gam}(\beta|c_0, d_0)$  over  $\beta$  and solve variationally, by assuming a factorized variational distribution of the form  $q(\mathbf{w})q(\alpha)q(\beta)$ .

We modify the joint distribution of all variables as:

$$p(\mathbf{t}, \mathbf{w}, \alpha, \beta) = p(\mathbf{t}|\mathbf{w}, \beta)p(\mathbf{w}|\alpha)p(\alpha)p(\beta)$$

The formulae for  $p(\alpha)$  remain the same:

$$q^*(\alpha) = \text{Gam}(\alpha|a_N, b_N), a_N = a_0 + \frac{M}{2}, b_N = b_0 + \frac{1}{2} \mathbb{E}[\mathbf{w}^T \mathbf{w}]$$

For  $q^*(\mathbf{w})$  we have:

$$\begin{aligned} \ln q^*(\mathbf{w}) &= \ln p(\mathbf{t}|\mathbf{w}, \beta) + \mathbb{E}_\alpha[\ln p(\mathbf{w}|\alpha)] + \text{const} \\ &= -\frac{\mathbb{E}[\beta]}{2} \sum_{n=1}^N \{\mathbf{w}^T \boldsymbol{\phi}_n - t_n\}^2 - \frac{1}{2} \mathbb{E}[\alpha] \mathbf{w}^T \mathbf{w} + \text{const} \\ &= -\frac{1}{2} \mathbf{w}^T (\mathbb{E}[\alpha] \mathbf{I} + \mathbb{E}[\beta] \boldsymbol{\Phi}^T \boldsymbol{\Phi}) \mathbf{w} + \mathbb{E}[\beta] \mathbf{w}^T \boldsymbol{\Phi}^T \mathbf{t} + \text{const} \end{aligned}$$

Thus  $q^*(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N)$  with

$$\begin{aligned} \mathbf{m}_N &= \mathbb{E}[\beta] \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t} \\ \mathbf{S}_N &= (\mathbb{E}[\alpha] \mathbf{I} + \mathbb{E}[\beta] \boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \end{aligned}$$



# Variational Linear Regression with $\text{Gam}(\beta|c_0, d_0)$

For  $q^*(\beta)$

$$\begin{aligned}\ln q^*(\beta) &= \mathbb{E}[\ln p(\mathbf{t}|\mathbf{w}, \beta)] + \ln p(\beta) + \text{const} \\ &= \frac{N}{2} \ln \beta - \frac{\beta}{2} \mathbb{E}_{\mathbf{w}} \left[ \sum_{n=1}^N \{\mathbf{w}^T \boldsymbol{\phi}_n - t_n\}^2 \right] + (c_0 - 1) \ln \beta - d_0 \beta + \text{const}\end{aligned}$$

We recognize the log of a Gamma distribution with:

$$\begin{aligned}q^*(\beta) &= \text{Gam}(\beta|c_N, d_N), c_N = c_0 + \frac{N}{2}, \\ d_N &= d_0 + \frac{1}{2} \mathbb{E}[\sum_{n=1}^N \{\mathbf{w}^T \boldsymbol{\phi}_n - t_n\}^2] = d_0 + \frac{1}{2} (\text{Tr}(\boldsymbol{\Phi}^T \boldsymbol{\Phi}) \mathbb{E}[\mathbf{w} \mathbf{w}^T] + \mathbf{t}^T \mathbf{t} - \mathbf{t}^T \boldsymbol{\Phi} \mathbb{E}[\mathbf{w}]) \\ &= d_0 + \frac{1}{2} (\|\mathbf{t} - \boldsymbol{\Phi} \mathbf{m}_N\|^2 + \text{Tr}(\boldsymbol{\Phi}^T \boldsymbol{\Phi} \mathbf{S}_N))\end{aligned}$$

Where we used:

$$\mathbb{E}[\mathbf{w} \mathbf{w}^T] = \mathbf{m}_N \mathbf{m}_N^T + \mathbf{S}_N \text{ and}$$

$$\mathbb{E}[\mathbf{w}] = \mathbf{m}_N, \text{ with } \mathbf{m}_N = \mathbb{E}[\beta] \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t}, \mathbf{S}_N = (\mathbb{E}[\alpha] \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1}$$

$$\text{Thus } \mathbb{E}[\beta] = \frac{c_N}{d_N}$$



# Variational Linear Regression with $\text{Gam}(\beta|c_0, d_0)$

The lower bound also needs to be modified. Starting with the modified log-likelihood and using  $\mathbb{E}[\mathbf{w}] = \mathbf{m}_N$ ,  $\mathbb{E}[\beta] = \frac{c_N}{d_N}$ ,  $\mathbb{E}[\mathbf{w}\mathbf{w}^T] = \mathbf{m}_N\mathbf{m}_N^T + \mathbf{S}_N$  and  $\mathbb{E}[\ln\beta] = \psi(c_N) - \ln d_N$ :

$$\begin{aligned}\mathbb{E}_\beta[\mathbb{E}_{\mathbf{w}}[\ln p(\mathbf{t}|\mathbf{w})]] &= \frac{N}{2}(\mathbb{E}[\beta] - \ln(2\pi)) - \frac{\mathbb{E}[\beta]}{2} - \mathbb{E}[\|\mathbf{t} - \Phi\mathbf{w}\|^2] \\ &= \frac{N}{2}(\psi(c_N) - \ln d_N - \ln(2\pi)) - \frac{c_N}{2d_N}(\|\mathbf{t} - \Phi\mathbf{w}\|^2 + \text{Tr}(\Phi^T\Phi\mathbf{S}_N))\end{aligned}$$

Next using  $\mathbb{E}[\beta] = \frac{c_N}{d_N}$  and  $\mathbb{E}[\ln\beta] = \psi(c_N) - \ln d_N$ , we consider the term corresponding to log prior over  $\beta$ :

$$\begin{aligned}\mathbb{E}[\ln p(\beta)] &= (c_0 - 1)\mathbb{E}[\ln\beta] - d_0\mathbb{E}[\beta] + c_0 \ln d_0 - \ln\Gamma(c_0) \\ &= (c_0 - 1)(\psi(c_N) - \ln d_N) - \frac{d_0 c_N}{d_N} + c_0 \ln d_0 - \ln\Gamma(c_0)\end{aligned}$$

Finally we compute the negative entropy of the posterior over  $\beta$ :

$$- \mathbb{E}[\ln q^*(\beta)] = (c_N - 1)\psi(c_N) + \ln d_N - c_N - \ln\Gamma(c_N)$$

Finally the predictive distribution is given as:

$$p(\mathbf{t}|\mathbf{x}, \mathbf{t}) \approx \mathcal{N}(\mathbf{t}|\mathbf{m}_N^T\Phi(\mathbf{x}), \sigma^2(\mathbf{x})), \quad \sigma^2(\mathbf{x}) = \frac{1}{\mathbb{E}[\beta]} + \Phi(\mathbf{x})^T\mathbf{S}_N\Phi(\mathbf{x})$$

