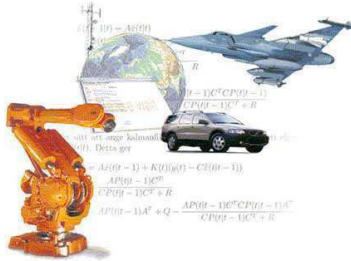


Machine Learning, Lecture 7 Approximate inference



Thomas Schön

Division of Automatic Control
Linköping University
Linköping, Sweden.

Email: schon@isy.liu.se,
Phone: 013 - 281373,
Office: House B, Entrance 27.



GM construction using latent variables

2(33)

Let z be a random variable having a 1-of- K coding scheme.

The marginal PDF of z is

$$p(z) = \prod_{k=1}^K \pi_k^{z_k},$$

where π_k are the mixture coefficients.

The conditional PDF of x given z is

$$p(x | z) = \prod_{k=1}^K \mathcal{N}(x | \mu_k, \Sigma_k)^{z_k}$$



E step (I/II)

3(33)

$$\begin{aligned} Q(\theta, \theta_i) &= \mathbb{E}_{\theta_i} [\ln p_{\theta}(Z, X) | X] \\ &= \mathbb{E}_{\theta_i} \left[\sum_{n=1}^N \sum_{k=1}^K z_{nk} (\ln \pi_k + \ln \mathcal{N}(x_n | \mu_k, \Sigma_k)) | X \right] \\ &= \sum_{n=1}^N \sum_{k=1}^K \underbrace{\mathbb{E}_{\theta_i} [z_{nk} | X]} (\ln \pi_k + \ln \mathcal{N}(x_n | \mu_k, \Sigma_k)) \end{aligned}$$

Hence, the E step amounts to finding $\mathbb{E}_{\theta_i} [z_{nk} | X]$, which is given by

$$\mathbb{E}_{\theta_i} [z_{nk} | X] = \sum_Z z_{nk} p_{\theta_i}(Z | X) = \sum_{z_{nk}} z_{nk} p_{\theta_i}(z_{nk} | X)$$



E step (I/II)

4(33)

$$\begin{aligned} \mathbb{E}_{\theta_i} [z_{nk} | X] &= \sum_{z_{nk}} z_{nk} \frac{p_{\theta_i}(x_n | z_{nk}) p_{\theta_i}(z_{nk})}{p_{\theta_i}(x_n)} \\ &= \frac{\sum_{z_{nk}} z_{nk} (\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k))^{z_{nk}}}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n | \mu_j, \Sigma_j)} \\ &= \frac{\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n | \mu_j, \Sigma_j)} \triangleq \gamma(z_{nk}), \end{aligned}$$

where the last equality follows from the fact that $z_{nk} \in \{0, 1\}$.



Algorithm 1 EM for Gaussian mixtures

- 1. Initialise:** Initialize $\mu_k^1, \Sigma_k^1, \pi_k^1$ and set $i = 1$.
- 2. While not converged do:**

- (a) Expectation (E) step:** Compute

$$\gamma(z_{nk}) = \frac{\pi_k^i \mathcal{N}(x_n | \mu_k^i, \Sigma_k^i)}{\sum_{j=1}^K \pi_j^i \mathcal{N}(x_n | \mu_j^i, \Sigma_j^i)}, \quad n = 1, \dots, N, k = 1, \dots, K.$$

- (b) Maximization (M) step:** Compute

$$\mu_k^{i+1} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) x_n, \quad \pi_k^{i+1} = \frac{N_k}{N}, \quad N_k = \sum_{n=1}^N \gamma(z_{nk})$$

$$\Sigma_k^{i+1} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (x_n - \mu_k^{i+1})(x_n - \mu_k^{i+1})^T$$

- (c) $i \leftarrow i + 1$**



Consider the same Gaussian mixture as before,

$$p(x) = \underbrace{0.3}_{\pi_1} \mathcal{N}\left(x \mid \underbrace{\begin{pmatrix} 4 \\ 4.5 \end{pmatrix}}_{\mu_1}, \underbrace{\begin{pmatrix} 1.2 & 0.6 \\ 0.6 & 0.5 \end{pmatrix}}_{\Sigma_1}\right) + \underbrace{0.5}_{\pi_2} \mathcal{N}\left(x \mid \underbrace{\begin{pmatrix} 8 \\ 1 \end{pmatrix}}_{\mu_2}, \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}}_{\Sigma_2}\right) + \underbrace{0.2}_{\pi_3} \mathcal{N}\left(x \mid \underbrace{\begin{pmatrix} 9 \\ 8 \end{pmatrix}}_{\mu_3}, \underbrace{\begin{pmatrix} 0.6 & 0.5 \\ 0.5 & 1.5 \end{pmatrix}}_{\Sigma_3}\right)$$

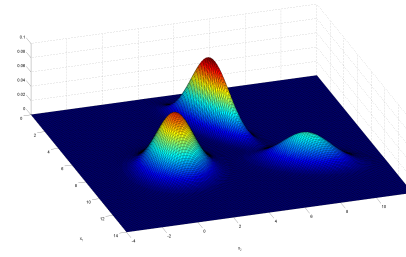
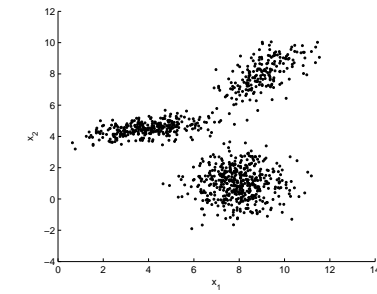


Figure: Probability density function.

Figure: $N = 1000$ samples from the Gaussian mixture $p(x)$.

- Apply the EM algorithm to estimate a Gaussian mixture with $K = 3$ Gaussians, i.e. use the 1000 samples to compute estimates of $\pi_1, \pi_2, \pi_3, \mu_1, \mu_2, \mu_3, \Sigma_1, \Sigma_2, \Sigma_3$.
- 200 iterations.

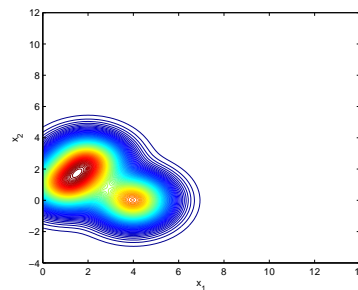


Figure: Initial guess.

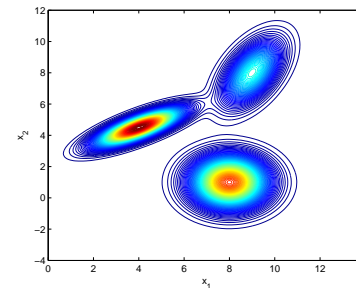


Figure: True PDF.

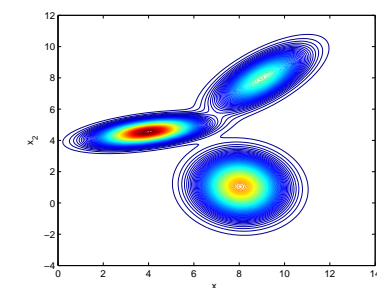


Figure: Estimate after 200 iterations of the EM algorithm.



Algorithm 2 K -means algorithm, a.k.a. Lloyd's algorithm

1. Initialize μ_k^1 and set $i = 1$.
2. Minimize J w.r.t. r_{nk} keeping $\mu_k = \mu_k^i$ fixed.

$$r_{nk}^{i+1} = \begin{cases} 1 & \text{if } k = \arg \min_j \|x_n - \mu_j^i\|^2 \\ 0 & \text{otherwise} \end{cases}$$

3. Minimize J w.r.t. μ_k keeping $r_{nk} = r_{nk}^{i+1}$ fixed.

$$\mu_k^{i+1} = \frac{\sum_{n=1}^N r_{nk}^{i+1} x_n}{\sum_{n=1}^N r_{nk}^{i+1}}.$$

4. If not converged, update $i := i + 1$ and return to step 2.

The name K -means stems from the fact that in step 3 of the algorithm, μ_k is given by the mean of all the data points assigned to cluster k .

Note the **similarities** between the K -means algorithm and the EM algorithm for Gaussian mixtures!

K -means is deterministic with “hard” assignment of data points to clusters (no uncertainty), whereas EM is a probabilistic method that provides a “soft” assignment.

If the Gaussian mixtures are modeled using covariance matrices

$$\Sigma_k = \epsilon I, \quad k = 1, \dots, K,$$

it can be shown that the EM algorithm for a mixture of K Gaussian's is **equivalent** to the K -means algorithm, when $\epsilon \rightarrow \infty$.

1. Summary of lecture 6
2. Bayesian reminder
3. Variational Bayesian inference
 - General derivation
 - Example – identification of an LGSS model
 - Example – Gaussian mixtures
4. Expectation propagation
 - General derivation
 - Example – state estimation

(Chapter 10)

This lecture builds on Umut Orguner's 2011 lecture.

The **Expectation Maximization (EM)** algorithm computes maximum likelihood estimates of unknown parameters in probabilistic models involving latent variables.

Expectation (E) step: Compute

$$\begin{aligned} Q(\theta, \theta_i) &= \mathbf{E}_{\theta_i} \{ \ln p_{\theta}(Z, X) \mid X \} \\ &= \int \ln p_{\theta}(Z, X) p_{\theta_i}(Z \mid X) dZ. \end{aligned}$$

Maximization (M) step: Compute

$$\theta_{i+1} = \arg \max_{\theta} Q(\theta, \theta_i).$$

We constructed a Gaussian mixture density using latent variables z (multinomial)

$$p(z) = \prod_{k=1}^K \pi_k^{z_k}, \quad p(x | z) = \prod_{k=1}^K \mathcal{N}(x | \mu_k, \Sigma_k)^{z_k}$$

This allowed us to derive an EM algorithm for estimating a Gaussian mixture.

Clustering is the problem of grouping N points $\{x_i\}_{i=1}^N$ into K clusters, where members of each cluster are “similar”.

The **K-means** algorithm tries to minimize $\sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$. We can show that the K-means algorithm is a (deterministic) special case of the EM algorithm.

In the Bayesian framework we are interested in the posterior density $p(Z|X)$ given by Bayes' rule as

$$p(Z|X) = \frac{p(X|Z)p(Z)}{p(X)},$$

where $X = x_1, \dots, x_N$ denotes the measurements and $Z = z_1, \dots, z_N$ denotes the latent variables.

Sometimes the posterior can be found exactly using the concept of **conjugate priors**.

- Gaussian case
- More generally the exponential family.

What happens when there is no exact solution?

Classic calculus involves functions and defines *derivatives* to optimize them.

The so-called **calculus of variations** investigates functions of functions which are called **functionals**.

$$\text{Example: Entropy } \mathcal{H}[p(\cdot)] = - \int p(x) \log(p(x)) dx.$$

The derivatives of functionals are called **variations**.

Calculus of variations has its origins in the 18th century and the most important result is probably the so-called Euler-lagrange equation

$$C(q) \triangleq \int \underbrace{L(t, q(t), q'(t))}_{\triangleq L(t, x, v)} dt, \quad L_x(t, q_*, q'_*) + \frac{d}{dt} L_v(t, q_*, q'_*) = 0,$$

which constitutes the core of optimal control theory.

In general variational methods, one generally assumes a predetermined form of the argument function, possibly parametric.

- Quadratic: $q(x) = x^T A x + b^T x + c$
- Basis functions: $q(x) = \sum_{i=1}^{N_\phi} w_i \phi_i(x)$

Variational inference: In the case of probabilistic inference, the variational approximation takes the form:

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

where $Z = \{Z_1, \dots, Z_M\}$ is a partitioning of the unknown variables.

Algorithm 3 Variational iteration

Solve the problem iteratively:

1. For $j = 1, \dots, M$

(a) Fix $\{q_i(Z_i)\}_{i=1}^M$ to their last estimated values $\{\hat{q}_i(Z_i)\}_{i=1}^M$, $i \neq j$.

(b) Find the solution of

$$\hat{q}_j(Z_j) = \arg \max_{q_j} \mathcal{L}(q)$$

2. Repeat 1 until convergence.

Consider the following Bayesian LGSS model

$$x_{k+1} = \theta x_k + v_k,$$

$$y_k = \frac{1}{2}x_k + e_k,$$

$$x_0 \sim \mathcal{N}(x_0; \bar{x}_0, \Sigma_0),$$

$$\theta \sim \mathcal{N}(\theta; 0, \sigma_\theta^2),$$

$$\begin{pmatrix} v_k \\ e_k \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_v^2 & 0 \\ 0 & \sigma_e^2 \end{pmatrix} \right).$$

Aim: Compute the posterior $p(\theta|y_{0:N})$ using the VB framework.

- We have some latent variables $x_{0:N} \triangleq \{x_0, \dots, x_N\}$.
- Different notation compared to Bishop! The observations are denoted y and the latent variables are denoted x .

With latent variables

$$p(\theta|y_{0:N}) = \int p(\theta, x_{0:N}|y_{0:N}) dx_{0:N}$$

There is still no exact form for the joint density $p(\theta, x_{0:N}|y_{0:N})$.

Variational Approximation

- Approximate the posterior $p(\theta, x_{0:N}|y_{0:N})$ as

$$p(\theta, x_{0:N}|y_{0:N}) \approx q_\theta(\theta)q_x(x_{0:N})$$

- Find $q_\theta(\theta)$ and $q_x(x_{0:N})$ using

$$\log q_\theta(\theta) = E_{q_x} [\log p(y_{0:N}, x_{0:N}, \theta)] + \text{const.}$$

$$\log q_x(x_{0:N}) = E_{q_\theta} [\log p(y_{0:N}, x_{0:N}, \theta)] + \text{const.}$$

Variational Bayes formulas are

$$\log q_\theta(\theta) = E_{q_x} [\log p(y_{0:N}, x_{0:N}, \theta)] + \text{const.}$$

$$\log q_x(x_{0:N}) = E_{q_\theta} [\log p(y_{0:N}, x_{0:N}, \theta)] + \text{const.}$$

We have the joint density $p(y_{0:N}, x_{0:N}, \theta)$ as

$$\begin{aligned} p(y_{0:N}, x_{0:N}, \theta) &= p(y_{0:N}|x_{0:N})p(x_{1:N}|x_{0:N-1}, \theta)p(x_0)p(\theta) \\ &= \prod_{i=0}^N p(y_i|x_i) \prod_{i=1}^N p(x_i|x_{i-1}, \theta)p(x_0)p(\theta) \end{aligned}$$

Taking the logarithm and separating the constant terms

$$\begin{aligned} \log p(y_{0:N}, x_{0:N}, \theta) &= - \sum_{i=0}^N \frac{0.5}{\sigma_e^2} (y_i - 0.5x_i)^2 - \sum_{i=1}^N \frac{0.5}{\sigma_v^2} (x_i - \theta x_{i-1})^2 \\ &\quad - 0.5/\sigma_0^2 (x_0 - \bar{x}_0)^2 - 0.5/\sigma_\theta^2 \theta^2 + \text{const.} \end{aligned}$$

Back to the Bishop's notation: x now denotes a measurement.

- Suppose we have $x_{1:N}$ i.i.d. and distributed as

$$x_i \sim p(x|\pi_{1:K}, \mu_{1:K}, \Lambda_{1:K}) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Lambda_k^{-1})$$

- In the Bayesian framework, all the unknowns $\{\pi_{1:K}, \mu_{1:K}, \Lambda_{1:K}\}$ are random.

$$\pi_{1:K} \sim \text{Dir}(\pi_{1:K}|\alpha_0) \propto \prod_{k=1}^K \pi_k^{\alpha_0-1}$$

$$\mu_{1:K}, \Lambda_{1:K} \sim p(\mu_{1:K}, \Lambda_{1:K}) \triangleq \prod_{k=1}^K \mathcal{N}(\mu_k; m_0, (\beta_0 \Lambda_k)^{-1}) \mathcal{W}(\Lambda_k | W_0, \nu_0)$$

- Define the latent variables $z_i \triangleq [z_{i1}, \dots, z_{iK}]^T$ as in EM. Then

$$p(x_{1:N}, z_{1:N}) = \prod_{i=1}^N \prod_{k=1}^K \pi_k^{z_{ik}} \mathcal{N}(x; \mu_k, \Lambda_k^{-1})^{z_{ik}}$$

- The Bayesian framework then asks for the posterior density $p(z_{1:N}, \pi_{1:K}, \mu_{1:K}, \Lambda_{1:K} | x_{1:N})$.

Variational Approximation

- Approximate the posterior as

$$p(z_{1:N}, \pi_{1:K}, \mu_{1:K}, \Lambda_{1:K} | x_{1:N}) \approx q_z(z_{1:N}) q_{\pi, \mu, \Lambda}(\pi_{1:K}, \mu_{1:K}, \Lambda_{1:K})$$

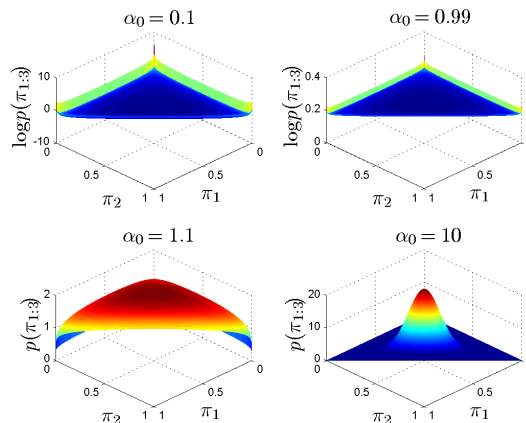
- Find $q_z(z_{1:N})$ and $q_{\pi, \mu, \Lambda}(\pi_{1:K}, \mu_{1:K}, \Lambda_{1:K})$ iteratively.

Symmetric Dirichlet distribution for $K = 3$.

$$\pi_{1:3} \sim \text{Dir}(\pi_{1:3}|\alpha_0)$$

$$\propto \prod_{k=1}^3 \pi_k^{\alpha_0-1}$$

$$= (\pi_1 \pi_2 (1 - \pi_1 - \pi_2))^{\alpha_0-1}$$

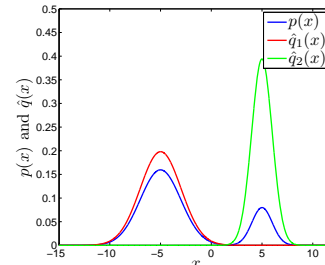


Suppose we have

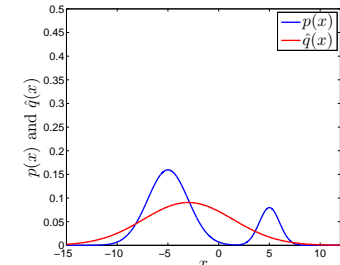
$$p(x) = 0.2\mathcal{N}(x; 5, 1) + 0.8\mathcal{N}(x, -5, 2^2)$$

Let $q_{\mu, \sigma}(x) \triangleq \mathcal{N}(x; \mu, \sigma^2)$

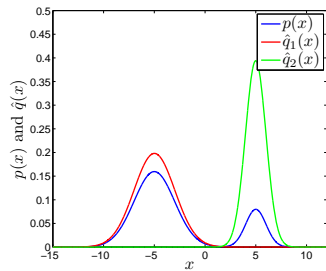
Find $\min_{\mu, \sigma} \text{KL}(q_{\mu, \sigma} || p)$



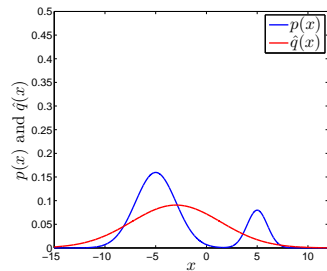
Find $\min_{\mu, \sigma} \text{KL}(p || q_{\mu, \sigma})$



Find $\min_{\mu, \sigma} \text{KL}(q_{\mu, \sigma} || p)$



Find $\min_{\mu, \sigma} \text{KL}(p || q_{\mu, \sigma})$



$$\text{KL}(q_{\mu, \sigma} || p) \triangleq \int q_{\mu, \sigma}(x) \log \frac{q_{\mu, \sigma}(x)}{p(x)} dx$$

zero-forcing

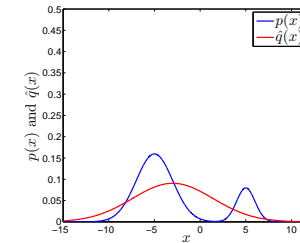
$$\text{KL}(p || q_{\mu, \sigma}) \triangleq \int p(x) \log \frac{p(x)}{q_{\mu, \sigma}(x)} dx$$

non-zero-forcing

This second form of optimization

$$\text{KL}(p || q_{\mu, \sigma}) \triangleq \int p(x) \log \frac{p(x)}{q_{\mu, \sigma}(x)} dx$$

has the following attractive property.



$$\hat{\mu} = E_{\hat{q}}(x) = E_p(x)$$

$$\hat{\sigma}^2 = E_{\hat{q}}[(x - E_{\hat{q}}(x))^2] = E_p[(x - E_p(xx^T))^2]$$

- Similar properties hold for the entire exponential family.
- A variational method using this type of KL-divergence minimization and hence the expectation equations above is **Expectation Propagation (EP)**.

- Suppose we have a posterior distribution in the form of

$$p(X|Y) \propto \prod_{i=1}^I f_i(X)$$

which is intractable or too computationally costly to compute.

- Then EP approximates the posterior as

$$p(X|Y) \approx q(X) \triangleq \prod_{i=1}^I q_i(X) = \prod_{i=1}^I \mathcal{N}(X; \mu_i, \Sigma_i)$$

- **Ideally** we want to minimize the KL divergence between the true posterior and the approximation,

$$\hat{q}(X) = \arg \min_q \text{KL} \left(\frac{1}{Z} \prod_{i=1}^I f_i(X) || \prod_{i=1}^I q_i(X) \right)$$

Solving this is intractable, make the approximation that we minimize the KL divergence between pairs of factors $f_i(X)$ and $q_i(X)$.

- The terms $q_j(x_j)$ are estimated iteratively as in VB by keeping the last estimates of $\{\hat{q}_i\}_{i \neq j}^I$.

$$\hat{q}_j(X) = \arg \min_{q_j} \text{KL} \left(f_j(X) \prod_{i \neq j} \hat{q}_i(X) || q_j(X) \prod_{i \neq j} \hat{q}_i(X) \right)$$

- This is in the Gaussian case obtained by solving the equations

$$E_{q_j \prod_{i \neq j} \hat{q}_i}(X) = E_{f_j \prod_{i \neq j} \hat{q}_i}(X)$$

$$E_{q_j \prod_{i \neq j} \hat{q}_i}(XX^T) = E_{f_j \prod_{i \neq j} \hat{q}_i}(XX^T)$$

for the mean μ_i and the covariance Σ_i of $\hat{q}_i(\cdot)$.

Consider the following LGSS model

$$\begin{aligned} x_{k+1} &= x_k + v_k, & x_0 &= 0 \text{ is known} \\ y_k &= x_k + e_k, & v_k &\sim \mathcal{N}(v_k; 0, \sigma_v^2) \end{aligned}$$

$$e_k \sim p_e(e_k) \triangleq 0.9\mathcal{N}(e_k; 0, \sigma_e^2) + 0.1\mathcal{N}(e_k; 0, (10\sigma_e)^2)$$

Aim: Compute the posterior density $p(x_{1:N}|y_{1:N})$.

- Recall that the true posterior factorizes as

$$p(x_{1:N}|y_{1:N}) \propto \prod_{i=1}^N p(y_i|x_i)p(x_i|x_{i-1})$$

- The true posterior in this case is a Gaussian mixture with 2^N components which is not feasible to compute.

- Make the variational approximation

$$p(x_{1:N}|y_{1:N}) \approx q(x_{1:N}) \triangleq \prod_{i=1}^N \mathcal{N}(x_i; \mu_i, \sigma_i^2)$$

- Consider the density for x_j given as

$$\begin{aligned} \bar{p}(x_j) &\propto \int \int p(y_j|x_j)p(x_{j+1}|x_j)p(x_j|x_{j-1}) \\ &\quad \times \mathcal{N}(x_{j+1}; \mu_{j+1}, \sigma_{j+1}^2)\mathcal{N}(x_{j-1}; \mu_{j-1}, \sigma_{j-1}^2) dx_{j+1}dx_{j-1} \end{aligned}$$

which can be calculated as

$$\begin{aligned} \bar{p}(x_j) &= w_1(\mu_{j\pm 1}, \sigma_{j\pm 1})\mathcal{N}(x_j; \eta_1(\mu_{j\pm 1}, \sigma_{j\pm 1}), \rho_1^2(\mu_{j\pm 1}, \sigma_{j\pm 1})) \\ &\quad + w_2(\mu_{j\pm 1}, \sigma_{j\pm 1})\mathcal{N}(x_j; \eta_2(\mu_{j\pm 1}, \sigma_{j\pm 1}), \rho_2^2(\mu_{j\pm 1}, \sigma_{j\pm 1})) \end{aligned}$$

$$\begin{aligned} \bar{p}(x_j) &= w_1(\mu_{j\pm 1}, \sigma_{j\pm 1})\mathcal{N}(x_j; \eta_1(\mu_{j\pm 1}, \sigma_{j\pm 1}), \rho_1^2(\mu_{j\pm 1}, \sigma_{j\pm 1})) \\ &\quad + w_2(\mu_{j\pm 1}, \sigma_{j\pm 1})\mathcal{N}(x_j; \eta_2(\mu_{j\pm 1}, \sigma_{j\pm 1}), \rho_2^2(\mu_{j\pm 1}, \sigma_{j\pm 1})) \end{aligned}$$

where the parameters $w_{1,2}, \eta_{1,2}$ and $\rho_{1,2}$ are

$$\begin{aligned} \eta_1 &= \rho_1^2 \left(\frac{\bar{\eta}}{\bar{\rho}^2} + \frac{y_j}{\sigma_e^2} \right) & \eta_2 &= \rho_2^2 \left(\frac{\bar{\eta}}{\bar{\rho}^2} + \frac{y_j}{(10\sigma_e)^2} \right) \\ \rho_1^2 &= \left(\frac{1}{\bar{\rho}^2} + \frac{1}{\sigma_e^2} \right)^{-1} & \rho_2^2 &= \left(\frac{1}{\bar{\rho}^2} + \frac{1}{(10\sigma_e)^2} \right)^{-1} \\ w_1 &\propto 0.9\mathcal{N}(y_j; \bar{\eta}, \bar{\rho}^2 + \sigma_e^2) & w_2 &\propto 0.1\mathcal{N}(y_j; \bar{\eta}, \bar{\rho}^2 + (10\sigma_e)^2) \\ \bar{\eta} &= \bar{\rho}^2 \left(\frac{\mu_{j-1}}{\sigma_{j-1}^2 + \sigma_v^2} + \frac{\mu_{j+1}}{\sigma_{j+1}^2 + \sigma_v^2} \right) & \bar{\rho}^2 &= \left(\frac{1}{\sigma_{j-1}^2 + \sigma_v^2} + \frac{1}{\sigma_{j+1}^2 + \sigma_v^2} \right)^{-1} \end{aligned}$$

$$\begin{aligned} \bar{p}(x_j) &= w_1(\mu_{j\pm 1}, \sigma_{j\pm 1})\mathcal{N}(x_j; \eta_1(\mu_{j\pm 1}, \sigma_{j\pm 1}), \rho_1^2(\mu_{j\pm 1}, \sigma_{j\pm 1})) \\ &\quad + w_2(\mu_{j\pm 1}, \sigma_{j\pm 1})\mathcal{N}(x_j; \eta_2(\mu_{j\pm 1}, \sigma_{j\pm 1}), \rho_2^2(\mu_{j\pm 1}, \sigma_{j\pm 1})) \end{aligned}$$

The EP solution for $q_j(x_j) = \mathcal{N}(x_j; \mu_j, \sigma_j^2)$ is obtained by matching (propagating) expectations between $q_j(\cdot)$ and $\bar{p}(x_j)$.

$$\begin{aligned} \mu_j &= w_1\eta_1 + w_2\eta_2 \\ \sigma_j^2 &= w_1(\rho_1^2 + (\eta_1 - \mu_j)^2) + w_2(\rho_2^2 + (\eta_2 - \mu_j)^2) \end{aligned}$$

Variational Inference: Approximate Bayesian inference where factorial approximations are made on the form of the posteriors.

Kullback-Leibler (KL) Divergence: A cost function to find optimal approximations for the posteriors in two different forms.

Variational Bayes: A form of variational inference where $\text{KL}(q||p)$ is used for the optimization.

Expectation Propagation: A form of variational inference where $\text{KL}(p||q)$ is used for the optimization.

