

Approximate Bayesian Inference for Machine Learning^{1,2}

Louis Ellam

l.ellam@warwick.ac.uk

Warwick Centre of Predictive Modelling (WCPM)
The University of Warwick

13 October 2015



<http://www2.warwick.ac.uk/wcpm/>

¹Based on Louis Ellam, "Approximate Inference for Machine Learning", MSc Dissertation, Centre for Scientific Computing, University of Warwick

²Warwick Centre of Predictive Modelling (WCPM) Seminar Series, University of Warwick

- 1 Motivation
 - Machine Learning
 - Classical Approaches to Regression
 - Statistical Approaches to Regression
- 2 Monte Carlo Inference
 - Importance Sampling
 - Markov Chain Monte Carlo
- 3 Approximate Inference
 - Laplace Approximation
 - Variational Inference
 - Expectation Propagation
- 4 Summary of Methods

- 1 Motivation
 - Machine Learning
 - Classical Approaches to Regression
 - Statistical Approaches to Regression
- 2 Monte Carlo Inference
 - Importance Sampling
 - Markov Chain Monte Carlo
- 3 Approximate Inference
 - Laplace Approximation
 - Variational Inference
 - Expectation Propagation
- 4 Summary of Methods

Machine Learning

- Machine Learning is the study of algorithms that can learn from data

Regression

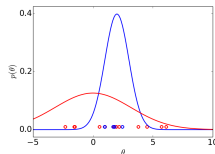
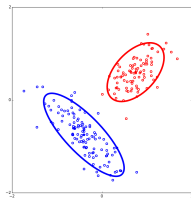
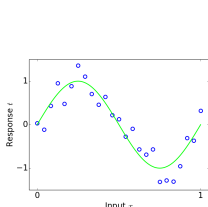
Given $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N, t_1, \dots, t_N\}$, find a model that allows us to accurately estimate the response of any input.

Clustering

Given data $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ and some fixed number of clusters K , assign each data point to a cluster so that points in each cluster are close together.

The Clutter Problem

Given data $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ comprising of a noisy signal that is immersed in clutter, recover the original signal θ i.e. determine $\mathbb{E}[\theta]$ and $p(\mathcal{D})$. Observations are drawn from $p(x|\theta) = (1 - w)\mathcal{N}(x|\theta, I) + w\mathcal{N}(x|\mathbf{0}, aI)$.



1 Motivation

- Machine Learning
- **Classical Approaches to Regression**
- Statistical Approaches to Regression

2 Monte Carlo Inference

- Importance Sampling
- Markov Chain Monte Carlo

3 Approximate Inference

- Laplace Approximation
- Variational Inference
- Expectation Propagation

4 Summary of Methods

Linear Model

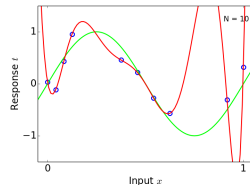
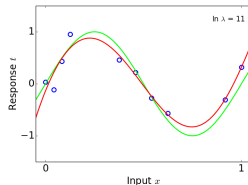
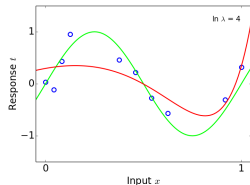
Use a linear model f defined over $M + 1$ linearly independent basis functions $\phi_0(x), \dots, \phi_M(x)$ with weights w_0, \dots, w_M

$$f(x, \mathbf{w}) = \sum_{m=0}^M w_m \phi_m(x).$$

Tikhonov Regularisation

$$E(\mathbf{w}) = \frac{1}{2} \|\Phi \mathbf{w} - \mathbf{t}\|^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2,$$

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} E(\mathbf{w}) = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{t}.$$



C. Bishop (2006), D. Calvetti et al. (2007)

- 1 Motivation
 - Machine Learning
 - Classical Approaches to Regression
 - **Statistical Approaches to Regression**
- 2 Monte Carlo Inference
 - Importance Sampling
 - Markov Chain Monte Carlo
- 3 Approximate Inference
 - Laplace Approximation
 - Variational Inference
 - Expectation Propagation
- 4 Summary of Methods

Bayes Theorem

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{D})}.$$

Likelihood and Prior

We expect that an observation to satisfy

$$t = f(\mathbf{x}, \mathbf{w}) + \eta,$$

where η denotes mean centered Gaussian noise with precision parameter β . Thus the likelihood function is defined as

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

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I}).$$

Posterior

$$\ln p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \alpha, \beta) = -\frac{\beta}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n))^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{const.}$$

R. Kass (1995)

Predictive Distribution

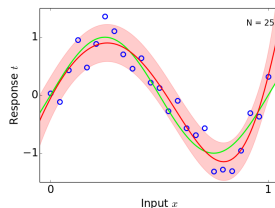
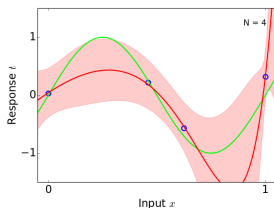
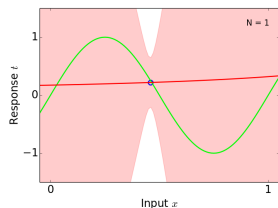
$$p(t|x, \mathbf{t}, \alpha, \beta) = \int p(t|\mathbf{w})p(\mathbf{w}|x, \mathbf{t}, \alpha, \beta)d\mathbf{w} = \mathcal{N}(t|\mathbf{m}_N^T\phi(x), \sigma_N^2(x)),$$

where

$$\sigma_N^2 = \frac{1}{\beta} + \phi(x)^T \mathbf{S}_N \phi(x),$$

$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi,$$

$$\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{y}.$$



- 'More Bayesian' if α is treated as a random variable, but analytically intractable.
- Classical integration methods are not suitable for high dimensional probability density functions.

- 1 Motivation
 - Machine Learning
 - Classical Approaches to Regression
 - Statistical Approaches to Regression
- 2 **Monte Carlo Inference**
 - **Importance Sampling**
 - Markov Chain Monte Carlo
- 3 Approximate Inference
 - Laplace Approximation
 - Variational Inference
 - Expectation Propagation
- 4 Summary of Methods

Monte Carlo Estimate

$$\mathbb{E}_f[q(\boldsymbol{\theta})] = \int q(\boldsymbol{\theta})f(\boldsymbol{\theta})d\boldsymbol{\theta} \approx \frac{1}{N} \sum_{n=1}^N q(\mathbf{X}_n) \text{ where } \mathbf{X}_n \stackrel{iid}{\sim} f(\boldsymbol{\theta}).$$

$$\mathbb{E}_f[q(\boldsymbol{\theta})] = \int q(\boldsymbol{\theta}) \underbrace{\frac{f(\boldsymbol{\theta})}{g(\boldsymbol{\theta})}}_{=:w(\boldsymbol{\theta})} g(\boldsymbol{\theta})d\boldsymbol{\theta} = \int q(\boldsymbol{\theta})w(\boldsymbol{\theta})g(\boldsymbol{\theta})d\boldsymbol{\theta}$$

Importance Sampling Estimate

$$\mathbb{E}_f[q(\boldsymbol{\theta})] \approx \frac{1}{N} \sum_{n=1}^N w(\mathbf{Z}_n)q(\mathbf{Z}_n) \text{ where } \mathbf{Z}_n \stackrel{iid}{\sim} g(\boldsymbol{\theta}).$$

Self-Normalised Importance Sampling Estimate

$$\mathbb{E}_f[q(\boldsymbol{\theta})] \approx \frac{1}{\sum_{m=1}^N \tilde{w}(\mathbf{Z}_m)} \sum_{n=1}^N \tilde{w}(\mathbf{Z}_n)q(\mathbf{Z}_n) \text{ where } \mathbf{Z}_n \stackrel{iid}{\sim} g(\boldsymbol{\theta}).$$

Monte Carlo Integration for the Clutter Problem

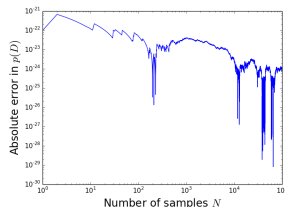
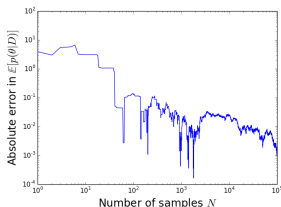
Define importance weights

$$w(\theta) = \frac{\tilde{f}(\theta)}{\tilde{g}(\theta)} = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\theta)} = p(\mathcal{D}|\theta).$$

Expectation (SNIS) and Evidence ('Vanilla' MC)

$$\mathbb{E}_{p(\theta|\mathcal{D})}[\theta] \approx \frac{1}{\sum_{m=1}^N p(\mathcal{D}|\mathbf{X}_m)} \sum_{n=1}^N \mathbf{X}_n p(\mathcal{D}|\mathbf{X}_n) \text{ where } \mathbf{X}_n \stackrel{iid}{\sim} p(\theta),$$

$$p(\mathcal{D}) = \mathbb{E}_{p(\theta)}[p(\mathcal{D}|\theta)] \approx \frac{1}{N} \sum_{n=1}^N p(\mathcal{D}|\mathbf{X}_n) \text{ where } \mathbf{X}_n \stackrel{iid}{\sim} p(\theta).$$



T. Minka (2001a), J. Liu (2008)

- 1 Motivation
 - Machine Learning
 - Classical Approaches to Regression
 - Statistical Approaches to Regression
- 2 **Monte Carlo Inference**
 - Importance Sampling
 - **Markov Chain Monte Carlo**
- 3 Approximate Inference
 - Laplace Approximation
 - Variational Inference
 - Expectation Propagation
- 4 Summary of Methods

Markov Chain Monte Carlo

- A discrete Markov chain is defined as a series of random variables $\mathbf{Z}^{(1)}, \dots, \mathbf{Z}^{(M)}$

$$p(\mathbf{Z}^{(M)} | \mathbf{Z}^{(1)}, \dots, \mathbf{Z}^{(M-1)}) = p(\mathbf{Z}^{(M)} | \mathbf{Z}^{(M-1)}).$$

- Draw new sample from proposal distribution and accept with probability

$$A(\mathbf{Z} | \mathbf{Z}^{(M-1)}) = \min \left\{ 1, \frac{f(\mathbf{Z})g(\mathbf{Z}^{(M-1)} | \mathbf{Z})}{f(\mathbf{Z}^{(M-1)})g(\mathbf{Z} | \mathbf{Z}^{(M-1)})} \right\}.$$

Algorithm 1 Metropolis-Hastings

Initialise $\mathbf{Z}^{(0)} = \{\mathbf{Z}_1^{(0)}, \dots, \mathbf{Z}_K^{(0)}\}$

for $M=1, \dots, T$ **do**

 Draw $\mathbf{Z} \sim g(\cdot | \mathbf{Z}^{(M-1)})$

 Compute $A(\mathbf{Z} | \mathbf{Z}^{(M-1)}) = \min \left\{ 1, \frac{f(\mathbf{Z})g(\mathbf{Z}^{(M-1)} | \mathbf{Z})}{f(\mathbf{Z}^{(M-1)})g(\mathbf{Z} | \mathbf{Z}^{(M-1)})} \right\}$

 With probability $A(\mathbf{Z} | \mathbf{Z}^{(M-1)})$ set $\mathbf{Z}^{(M)} = \mathbf{Z}$, otherwise set $\mathbf{Z}^{(M)} = \mathbf{Z}^{(M-1)}$

end for

Algorithm 2 Gibbs Sampling

Initialise $\mathbf{Z}_1^{(0)}, \dots, \mathbf{Z}_K^{(0)}$

for $M=1, \dots, T$ **do**

 draw $\mathbf{Z}_1^{(M)} \sim \mathbf{Z}_1 | \mathbf{Z}_2^{(M-1)}, \dots, \mathbf{Z}_K^{(M-1)}$

 draw $\mathbf{Z}_2^{(M)} \sim \mathbf{Z}_2 | \mathbf{Z}_1^{(M)}, \mathbf{Z}_3^{(M-1)}, \dots, \mathbf{Z}_K^{(M-1)}$

 ⋮

 draw $\mathbf{Z}_K^{(M)} \sim \mathbf{Z}_K | \mathbf{Z}_1^{(M)}, \dots, \mathbf{Z}_{K-1}^{(M)}$

end for

Likelihood for GMM

$$p(\mathcal{D}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \prod_{n=1}^N \sum_{k=1}^K \pi_k \mathcal{N}(x_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k).$$

Introduce Latent Variables

$$p(z_{nk} = 1 | \pi_k) = \pi_k,$$

$$p(\mathbf{z} | \boldsymbol{\pi}) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}}.$$

Joint Distribution for GMM (Augmented)

$$p(\mathcal{D}, \mathbf{z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = p(\mathcal{D} | \mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) p(\mathbf{z} | \boldsymbol{\pi}) p(\boldsymbol{\pi}) p(\boldsymbol{\mu}) p(\boldsymbol{\Sigma}).$$

$$p(\mathcal{D} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{z}) = \prod_{n=1}^N \prod_{k=1}^K \mathcal{N}(x_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_{nk}}.$$

$$p(\boldsymbol{\mu}) p(\boldsymbol{\Sigma}) = \prod_{k=1}^K \mathcal{N}(\boldsymbol{\mu}_k | \mathbf{m}_0, \mathbf{V}_0) \mathcal{IW}(\boldsymbol{\Sigma}_k | \mathbf{W}_0, \nu_0),$$

$$p(\boldsymbol{\pi}) = \text{Dir}(\boldsymbol{\pi} | \boldsymbol{\alpha}_0).$$

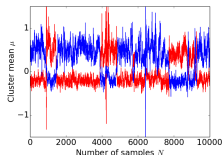
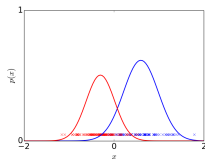
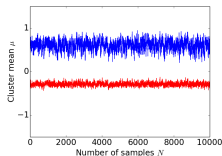
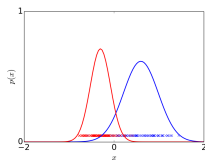
Full Conditionals

$$p(z_{nk} = 1 | x_n, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \frac{\pi_k \mathcal{N}(x_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)},$$

$$p(\boldsymbol{\pi} | \mathbf{z}) = \text{Dir}(\boldsymbol{\pi} | \boldsymbol{\alpha}_k),$$

$$p(\boldsymbol{\mu}_k | \boldsymbol{\Sigma}_k, \mathbf{z}, \mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_k | \mathbf{m}_k, \mathbf{V}_k),$$

$$p(\boldsymbol{\Sigma}_k | \boldsymbol{\mu}_k, \mathbf{z}, \mathbf{x}) = \mathcal{IW}(\boldsymbol{\Sigma}_k | \mathbf{W}_k, \nu_k),$$



- 1 Motivation
 - Machine Learning
 - Classical Approaches to Regression
 - Statistical Approaches to Regression
- 2 Monte Carlo Inference
 - Importance Sampling
 - Markov Chain Monte Carlo
- 3 **Approximate Inference**
 - **Laplace Approximation**
 - Variational Inference
 - Expectation Propagation
- 4 Summary of Methods

Laplace Approximation

Seek to approximate probability density function $p(\mathbf{z})$ of the form $p(\mathbf{z}) = \frac{1}{Z} f(\mathbf{z})$.

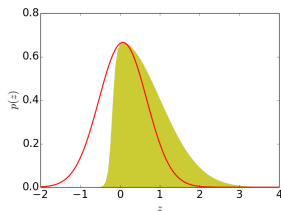
$$\ln f(\mathbf{z}) \approx \ln f(\mathbf{z}_0) + (\mathbf{z} - \mathbf{z}_0)^T \nabla \ln f(\mathbf{z})|_{\mathbf{z}=\mathbf{z}_0} + \frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T \nabla \nabla \ln f(\mathbf{z})|_{\mathbf{z}=\mathbf{z}_0} (\mathbf{z} - \mathbf{z}_0),$$
$$\nabla f(\mathbf{z})|_{\mathbf{z}=\mathbf{z}_0} = 0.$$

Laplace Approximation

$$q(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{z}_0, \mathbf{A}^{-1}), \text{ where } \mathbf{A} = -\nabla \nabla \ln f(\mathbf{z})|_{\mathbf{z}=\mathbf{z}_0},$$

Example

$$p(z) \propto \exp(-z^2/2)\sigma(20z + 4), \text{ where } \sigma(t) = (1 + e^{-t})^{-1}.$$



Want to find mode of GMM

$$\ln p(\mathcal{D}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}.$$

MLE of $\boldsymbol{\mu}_k$ is

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n,$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk}),$$
$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}.$$

MLE of $\boldsymbol{\Sigma}_k$ is

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T.$$

A. Dempster et al. (1977), C. Bishop (2006)

EM Algorithm (GMM)

The maximisation for $\boldsymbol{\pi}$ is constrained to $\sum_{k=1}^K \pi_k = 1$ so find MLE of

$$\ln p(\mathcal{D} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right),$$

which yields the following

$$\pi_k = \frac{N_k}{N}.$$

Algorithm 3 EM GMM

Initialise $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k$ for $k = 1, \dots, K$

repeat

E Step: Evaluate responsibilities

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(x_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}.$$

M Step: Define

$$N_k = \sum_{n=1}^N \gamma(z_{nk}),$$

then update parameters

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) x_n,$$

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (x_n - \boldsymbol{\mu}_k)(x_n - \boldsymbol{\mu}_k)^T,$$

$$\pi_k = \frac{N_k}{N}.$$

until convergence

Link to [code](#) / [.mp4](#)

- 1 Motivation
 - Machine Learning
 - Classical Approaches to Regression
 - Statistical Approaches to Regression
- 2 Monte Carlo Inference
 - Importance Sampling
 - Markov Chain Monte Carlo
- 3 **Approximate Inference**
 - Laplace Approximation
 - **Variational Inference**
 - Expectation Propagation
- 4 Summary of Methods

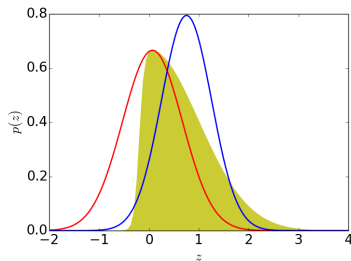
KL Divergence

Choose a proxy function $q(\theta)$ and optimise θ so that q is as close fit as possible to p . One such approach is to minimise

$$\text{KL}(q||p) := - \int q(\theta) \ln \frac{p(\theta|\mathcal{D})}{q(\theta)} d\theta.$$

Example

$$p(z) \propto \exp(-z^2/2)\sigma(20z + 4), \text{ where } \sigma(t) = (1 + e^{-t})^{-1}.$$



H. Attias (2000), Z. Ghahramani et al. (1999), C. Bishop (2006)

$$\begin{aligned}\ln p(\mathcal{D}) &= \ln \frac{p(\mathcal{D}, \theta)}{p(\theta|\mathcal{D})} = \int q(\theta) \ln \frac{p(\mathcal{D}, \theta)}{p(\theta|\mathcal{D})} d\theta = \int q(\theta) \ln \frac{p(\mathcal{D}, \theta)}{p(\theta|\mathcal{D})} \frac{q(\theta)}{q(\theta)} d\theta \\ &= \int q(\theta) \ln \frac{p(\mathcal{D}, \theta)}{q(\theta)} d\theta - \int q(\theta) \ln \frac{p(\theta|\mathcal{D})}{q(\theta)} d\theta \\ &= \mathcal{L}(q) + \text{KL}(q||p).\end{aligned}$$

Lowerbound

Minimising the KL divergence is equivalent to maximising the lower bound

$$\mathcal{L}(q) = \int q(\theta) \ln \frac{p(\mathcal{D}, \theta)}{q(\theta)} d\theta,$$

Lowerbound Evidence Approximation

It may be shown that

$$\text{KL}(q||p) \geq 0.$$

hence $p(\mathcal{D})$ may be approximated using

$$\begin{aligned}\ln p(\mathcal{D}) &\geq \mathcal{L}(q), \\ p(\mathcal{D}) &\rightarrow \exp \mathcal{L}(q) \quad \text{as} \quad \text{KL}(q||p) \rightarrow 0.\end{aligned}$$

Mean Field Approximation

Assume $q(\boldsymbol{\theta})$ can be expressed in factorised form

$$q(\boldsymbol{\theta}) = \prod_{m=1}^M q_m(\boldsymbol{\theta}_m).$$

Mean Field Lowerbound

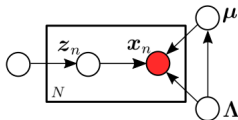
$$\begin{aligned}\mathcal{L}(q) &= \int \prod_{m=1}^M q(\boldsymbol{\theta}_m) \ln \frac{p(\mathcal{D}, \boldsymbol{\theta})}{\prod_{m=1}^M q(\boldsymbol{\theta}_m)} d\boldsymbol{\theta} = \int \prod_{m=1}^M q(\boldsymbol{\theta}_m) \left\{ \ln p(\mathcal{D}, \boldsymbol{\theta}) - \sum_{m=1}^M \ln q(\boldsymbol{\theta}_m) \right\} d\boldsymbol{\theta} \\ &= \int q(\boldsymbol{\theta}_j) \left\{ \int \ln p(\mathcal{D}, \boldsymbol{\theta}) \prod_{\substack{m=1 \\ m \neq j}}^M q(\boldsymbol{\theta}_m) d\boldsymbol{\theta}_m \right\} d\boldsymbol{\theta}_j - \int q(\boldsymbol{\theta}_j) \ln q(\boldsymbol{\theta}_j) d\boldsymbol{\theta}_j + \text{const} \\ &= \int q(\boldsymbol{\theta}_j) \mathbb{E}_{m \neq j} [\ln p(\mathcal{D}, \boldsymbol{\theta})] d\boldsymbol{\theta}_j - \int q(\boldsymbol{\theta}_j) \ln q(\boldsymbol{\theta}_j) d\boldsymbol{\theta}_j + \text{const} \\ &= -\text{KL}(q(\boldsymbol{\theta}_j) \parallel \exp \mathbb{E}_{m \neq j} [\ln p(\mathcal{D}, \boldsymbol{\theta})]) + \text{const}.\end{aligned}$$

Mean Field Minimised KL Divergence

$$\ln q_j^*(\boldsymbol{\theta}_j) = \mathbb{E}_{m \neq j} [\ln p(\mathcal{D}, \boldsymbol{\theta})] + \text{const}.$$

Target Posterior

$$p(\mathcal{D}, \mathbf{z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = p(\mathcal{D} | \mathbf{z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) p(\mathbf{z} | \boldsymbol{\pi}) p(\boldsymbol{\pi}) p(\boldsymbol{\mu} | \boldsymbol{\Lambda}) p(\boldsymbol{\Lambda}).$$



Conjugate Prior

$$p(\boldsymbol{\mu} | \boldsymbol{\Lambda}) p(\boldsymbol{\Lambda}) = \prod_{k=1}^K p(\boldsymbol{\mu}_k | \boldsymbol{\Lambda}_k) p(\boldsymbol{\Lambda}_k) = \prod_{k=1}^K \mathcal{N}(\boldsymbol{\mu}_k | \mathbf{m}_0, (\beta_0 \boldsymbol{\Lambda}_k)^{-1}) \mathcal{W}(\boldsymbol{\Lambda}_k | \mathbf{W}_0, \nu_0).$$

Joint Distribution

$$p(\mathcal{D}, \mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}, \boldsymbol{\pi}) \propto \left(\prod_{n=1}^N \prod_{k=1}^K \left(\pi_k \mathcal{N}(x_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) \right)^{z_{nk}} \right) \cdot \text{Dir}(\boldsymbol{\pi} | \boldsymbol{\alpha}_0) \cdot \prod_{k=1}^K \mathcal{N}(\boldsymbol{\mu}_k | \mathbf{m}_0, \beta_0^{-1} \boldsymbol{\Lambda}_k) \mathcal{W}(\boldsymbol{\Lambda}_k | \mathbf{W}_0, \nu_0).$$

Mean Field Approximation

$$q(\mathbf{z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = q(\mathbf{z})q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}),$$

Mean Field Minimised KL Divergence

$$\ln q_j^*(\boldsymbol{\theta}_j) = \mathbb{E}_{m \neq j}[\ln p(\mathcal{D}, \boldsymbol{\theta})] + \text{const.}$$

Expression for Responsibilities

$$\ln q^*(\mathbf{z}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \ln \rho_{nk} + \text{const.},$$

where

$$\ln \rho_{nk} = \mathbb{E}[\ln \pi_k] + \frac{1}{2} \mathbb{E}[\ln |\boldsymbol{\Lambda}_k|] - \frac{D}{2} \ln 2\pi - \frac{1}{2} \mathbb{E}_{\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k} [(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Lambda}_k (\mathbf{x}_n - \boldsymbol{\mu}_k)],$$

which can be used to give an expression of the responsibilities

$$r_{nk} = \mathbb{E}[z_{nk}] = \frac{\rho_{nk}}{\sum_{j=1}^K \rho_{nj}}.$$

Define

$$N_k = \sum_{n=1}^N r_{nk},$$
$$\bar{\mathbf{x}}_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} \mathbf{x}_n,$$
$$\mathbf{S}_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \bar{\mathbf{x}}_k)(\mathbf{x}_n - \bar{\mathbf{x}}_k)^T.$$

Expression for Model Parameters

$$\ln q^*(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = \mathbb{E}_z[\ln p(\mathcal{D}, \mathbf{z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})] + \text{const.}$$

By inspection

$$q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = q(\boldsymbol{\pi})q(\boldsymbol{\mu}, \boldsymbol{\Lambda}) = q(\boldsymbol{\pi}) \prod_{k=1}^K q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k),$$

$$q(\boldsymbol{\pi}|\boldsymbol{\alpha}) = \text{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha}),$$

$$q^*(\boldsymbol{\mu}_k|\boldsymbol{\Lambda}_k) = \mathcal{N}(\boldsymbol{\mu}_k|\mathbf{m}_k, (\beta_k \boldsymbol{\Lambda}_k)^{-1}),$$

$$q^*(\boldsymbol{\Lambda}_k) = \mathcal{W}(\boldsymbol{\Lambda}_k|\mathbf{W}_k, \nu_k).$$

Expressions for Model Parameters

$$\begin{aligned}\alpha_k &= \alpha_0 + N_k, \quad \beta_k = \beta_0 + N_k \text{ and } \nu_k = \nu_0 + N_k, \\ \mathbf{m}_k &= \frac{1}{\beta_k} (\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k), \\ \mathbf{W}_k^{-1} &= \mathbf{W}_0^{-1} + N_k \mathbf{S}_k + \frac{\beta_0 N_k}{\beta_0 + N_k} (\bar{\mathbf{x}}_k - \mathbf{m}_0)(\bar{\mathbf{x}}_k - \mathbf{m}_0)^T.\end{aligned}$$

Expressions for Responsibilities

$$\begin{aligned}\mathbb{E}_{\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k} [(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Lambda}_k (\mathbf{x}_n - \boldsymbol{\mu}_k)] &= \nu_k (\mathbf{x}_n - \mathbf{m}_k)^T \mathbf{W}_k (\mathbf{x}_n - \mathbf{m}_k) + D \beta_k^{-1}, \\ \ln \tilde{\pi}_k := \mathbb{E}[\ln |\pi_k|] &= \psi(\alpha_k) - \psi\left(\sum_{k=1}^K \alpha_k\right), \\ \ln \tilde{\Lambda}_k := \mathbb{E}[\ln |\boldsymbol{\Lambda}_k|] &= \sum_{i=1}^D \psi\left(\frac{\nu_k + 1 - i}{2}\right) + D \ln 2 + \ln |\mathbf{W}_k|,\end{aligned}$$

Mixing Coefficients

$$\mathbb{E}[\pi_k] = \frac{\alpha_0 + N_k}{K \alpha_0 + N},$$

Algorithm 4 Variational Bayes GMM

Initialise $\pi, \mu, \Lambda, \mathbf{m}, \mathbf{W}, \beta_0, \nu_0, \alpha_0$

repeat

E Step: Compute the responsibilities as follows

$$r_{nk} \propto \tilde{\pi}_k \tilde{\Lambda}_k^{1/2} \exp \left\{ -\frac{D}{2\beta_k} - \frac{\nu_k}{2} (\mathbf{x}_n - \mathbf{m}_k)^T \mathbf{W}_k (\mathbf{x}_n - \mathbf{m}_k) \right\},$$

$$\ln \tilde{\pi}_k = \psi(\alpha_k) - \psi \left(\sum_{k=1}^K \alpha_k \right),$$

$$\ln \tilde{\Lambda}_k = \sum_{i=1}^D \left(\psi \left(\frac{\nu_k + 1 - i}{2} \right) \right) + D \ln 2 + \ln |\mathbf{W}_k|,$$

M Step: Update model parameters as follows

$$N_k = \sum_{n=1}^N r_{nk},$$

$$\bar{\mathbf{x}}_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} \mathbf{x}_n,$$

$$\mathbf{S}_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \bar{\mathbf{x}}_k)(\mathbf{x}_n - \bar{\mathbf{x}}_k)^T,$$

$$\alpha_k = \alpha_0 + N_k, \beta_k = \beta_0 + N_k \text{ and } \nu_k = \nu_0 + N_k,$$

$$\mathbf{m}_k = \beta_k^{-1} (\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k),$$

$$\mathbf{W}_k^{-1} = \mathbf{W}_0^{-1} + N_k \mathbf{S}_k + \frac{\beta_0 N_k}{\beta_0 + N_k} (\bar{\mathbf{x}}_k - \mathbf{m}_0)(\bar{\mathbf{x}}_k - \mathbf{m}_0)^T.$$

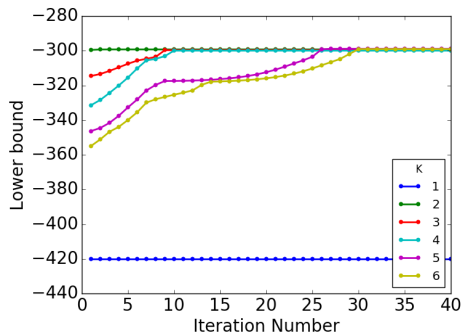
until convergence

Link to [code](#) / [.mp4](#)

Link to [code](#) / [.mp4](#)

Lowerbound

$$\begin{aligned}
 \mathcal{L}(q) &= \sum_{\mathbf{z}} \iiint q(\mathbf{z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \ln \left\{ \frac{p(\mathcal{D}, \mathbf{z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})}{q(\mathbf{z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})} \right\} d\boldsymbol{\pi} d\boldsymbol{\mu} d\boldsymbol{\Lambda} \\
 &= \mathbb{E}[\ln p(\mathcal{D}, \mathbf{z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})] - \mathbb{E}[\ln q(\mathbf{z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})] \\
 &= \mathbb{E}[\ln p(\mathcal{D} | \mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\Lambda})] + \mathbb{E}[\ln p(\mathbf{z} | \boldsymbol{\pi})] + \mathbb{E}[\ln p(\boldsymbol{\pi})] + \mathbb{E}[\ln p(\boldsymbol{\mu}, \boldsymbol{\Lambda})] \\
 &\quad - \mathbb{E}[\ln q(\mathbf{z})] - \mathbb{E}[\ln q(\boldsymbol{\pi})] - \mathbb{E}[\ln q(\boldsymbol{\mu}, \boldsymbol{\Lambda})].
 \end{aligned}$$



Model Selection with Lowerbound

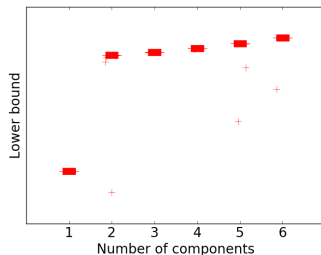
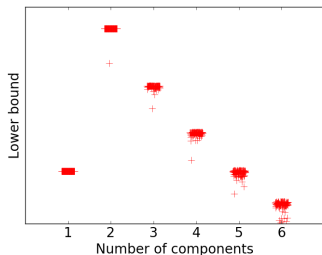
$$p(K|\mathcal{D}) = \frac{K!p(\mathcal{D}|K)}{\sum_{K'} K'!p(\mathcal{D}|K')},$$

where

$$p(\mathcal{D}|K) \approx \exp \mathcal{L}(K).$$

Alternative View

$$\mathcal{L}_M(K) \approx \mathcal{L}(K) + \ln K!.$$



A. Corduneanu (2001), C. Bishop (2006), K. Murphy (2012)

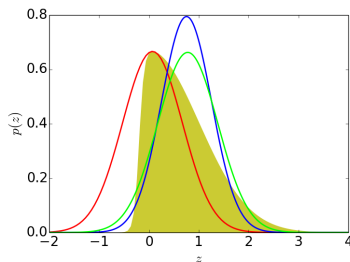
- 1 Motivation
 - Machine Learning
 - Classical Approaches to Regression
 - Statistical Approaches to Regression
- 2 Monte Carlo Inference
 - Importance Sampling
 - Markov Chain Monte Carlo
- 3 **Approximate Inference**
 - Laplace Approximation
 - Variational Inference
 - **Expectation Propagation**
- 4 Summary of Methods

Minimise Reverse KL Divergence

$$\text{KL}(p||q) = - \int p(\theta|\mathcal{D}) \ln \frac{q(\theta)}{p(\theta|\mathcal{D})} d\theta.$$

Example

$$p(z) \propto \exp(-z^2/2)\sigma(20z + 4), \text{ where } \sigma(t) = (1 + e^{-t})^{-1}.$$



T. Minka (2001a), T. Minka (2001b)

Exponential Family

$$q(\boldsymbol{\theta}) = h(\boldsymbol{\theta})g(\boldsymbol{\eta}) \exp\{\boldsymbol{\eta}^T \mathbf{u}(\boldsymbol{\theta})\},$$

Reverse KL Divergence for Exponential Family

$$\begin{aligned} \text{KL}(p||q) &= \int p(\boldsymbol{\theta}) \ln \left\{ \frac{p(\boldsymbol{\theta})}{q(\boldsymbol{\theta})} \right\} d\boldsymbol{\theta} \\ &= \int p(\boldsymbol{\theta}) \{ \ln p(\boldsymbol{\theta}) - \ln (h(\boldsymbol{\theta})g(\boldsymbol{\eta})) - \boldsymbol{\eta}^T \mathbf{u}(\boldsymbol{\theta}) \} d\boldsymbol{\theta} \\ &= -\ln g(\boldsymbol{\eta}) - \boldsymbol{\eta}^T \mathbb{E}_{p(\boldsymbol{\theta})}[\mathbf{u}(\boldsymbol{\theta})] + \text{const.} \end{aligned}$$

Minimised Reverse KL Divergence for Exponential Family

$$-\nabla \ln g(\boldsymbol{\eta}) = \mathbb{E}_{p(\boldsymbol{\theta})}[\mathbf{u}(\boldsymbol{\theta})],$$

i.e. when

$$\mathbb{E}_{q(\boldsymbol{\theta})}[\mathbf{u}(\boldsymbol{\theta})] = \mathbb{E}_{p(\boldsymbol{\theta})}[\mathbf{u}(\boldsymbol{\theta})].$$

Difficult, as trying to taking expectation w.r.t. unknown posterior.

Assumed Form of Joint Distribution

$$p(\boldsymbol{\theta}, \mathcal{D}) = \prod_n f_n(\boldsymbol{\theta}),$$

Approximating Distribution

An approximating posterior belonging to the **exponential family** is sought of the following form where Z is a normalisation constant

$$q(\boldsymbol{\theta}) = \frac{1}{Z} \prod_n \tilde{f}_n(\boldsymbol{\theta}).$$

- The idea is to refine each factor \tilde{f}_n by making the approximating distribution

$$q^{\text{new}}(\boldsymbol{\theta}) \propto \tilde{f}_n(\boldsymbol{\theta}) \prod_{j \neq n} \tilde{f}_j(\boldsymbol{\theta}),$$

a closer approximation to

$$f_n(\boldsymbol{\theta}) \prod_{j \neq n} \tilde{f}_j(\boldsymbol{\theta}).$$

- To proceed, factor n is first removed from the approximating distribution by

$$q^{\setminus n}(\theta) = \frac{q(\theta)}{\tilde{f}_n(\theta)}.$$

- A new distribution is defined by

$$p^n(\theta) = \frac{1}{Z_n} f_n(\theta) q^{\setminus n}(\theta),$$

with normalisation constant

$$Z_n = \int f_n(\theta) q^{\setminus n}(\theta) d\theta.$$

- The factor \tilde{f}_n is updated such that the sufficient statistics of q are matched to those of p^n .

Algorithm 5 Expectation Propagation

Initialise approximating factors $\tilde{f}_n(\theta)$

Initialise approximating function $q(\theta) \propto \prod_n \tilde{f}_n(\theta)$

repeat

 Choose a factor $\tilde{f}_n(\theta)$ for refinement

 Set $q^{\setminus n}(\theta) = q(\theta) / \tilde{f}_n(\theta)$

 Set sufficient statistics of $q^{\text{new}}(\theta)$ to those of $q^{\setminus n}(\theta) f_n(\theta)$ including $Z = \int q^{\setminus n}(\theta) f_n(\theta) d\theta$

 Set $\tilde{f}_n(\theta) = Z_n q^{\text{new}}(\theta) / q^{\setminus n}(\theta)$

until convergence

Factorised Posterior Distribution

$$f_n(\boldsymbol{\theta}) = (1 - w)\mathcal{N}(x_n|\boldsymbol{\theta}, \mathbf{I}) + w\mathcal{N}(x_n|\mathbf{0}, a\mathbf{I}) \text{ for } n = 1, \dots, N$$
$$f_0(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{0}, a\mathbf{I})$$

Approximating Distribution

Taken to be an isotropic Gaussian of the form

$$q(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}, \nu\mathbf{I}),$$

- q is treated as a product of $N + 1$ approximating factors

$$q(\boldsymbol{\theta}) \propto \prod_{n=0}^N \tilde{f}_n(\boldsymbol{\theta}).$$

- The factors have been defined as follows where factor \tilde{f}_n has scaling constant s_n , mean \mathbf{m}_n and variance $\nu_n\mathbf{I}$

$$\tilde{f}_n(\boldsymbol{\theta}) = s_n\mathcal{N}(\boldsymbol{\theta}|\mathbf{m}_n, \nu_n\mathbf{I}),$$

in terms of a scaled Gaussian distribution.

- The approximating distribution is therefore Gaussian

$$q(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}, \nu\mathbf{I}).$$

- $p(\mathcal{D}) \approx \int q(\boldsymbol{\theta})d\boldsymbol{\theta}$ is obtained using results for Gaussian products.

Expectation Propagation for Clutter Problem

Algorithm 6 Expectation Propagation for the Clutter Problem

Initialise the prior factor $\tilde{f}_0(\theta) = \mathcal{N}(\theta|\mathbf{0}, bI)$

Initialise the remaining factors $\tilde{f}_n = 1, v_n = \infty, m_n = \mathbf{0}$

Initialise the approximating function $m = \mathbf{0}, v = b$

repeat

for $n=1, \dots, N$ **do**

 Remove the current factor from $q(\theta)$

$$(v \setminus^n)^{-1} = v^{-1} - v_n^{-1},$$

$$m \setminus^n = m + v_n^{-1} v \setminus^n (m - m_n).$$

 Evaluate the new posterior $q(\theta)$ by matching sufficient statistics

$$Z_n = (1 - w) \mathcal{N}(x_n | m \setminus^n, (v \setminus^n + 1)I) + w \mathcal{N}(x_n | \mathbf{0}, aI),$$

$$\rho_n = 1 - \frac{w}{Z_n} \mathcal{N}(x_n | \mathbf{0}, aI),$$

$$m = m \setminus^n + \rho_n \frac{v \setminus^n}{v \setminus^n + 1} (x_n - m \setminus^n),$$

$$v = v \setminus^n - \rho_n \frac{(v \setminus^n)^2}{v \setminus^n + 1} + \rho_n (1 + \rho_n) \frac{(v \setminus^n)^2 \|x_n - m \setminus^n\|^2}{D(v \setminus^n + 1)^2}.$$

 Evaluate and store the new factor $\tilde{f}_n(\theta)$

$$v_n^{-1} = v^{-1} - (v \setminus^n)^{-1},$$

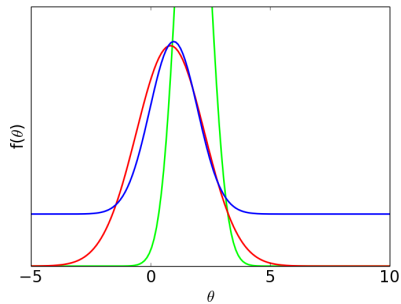
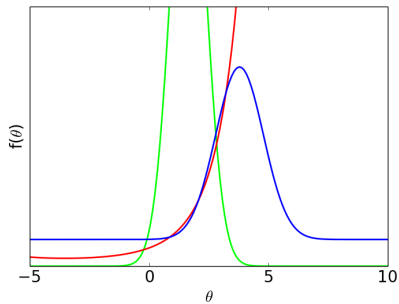
$$m_n = m \setminus^n + (v_n + v \setminus^n)(v \setminus^n)^{-1} (m - m \setminus^n),$$

$$s_n = \frac{Z_n}{\mathcal{N}(m_n | m \setminus^n, (v_n + v \setminus^n)I)}.$$

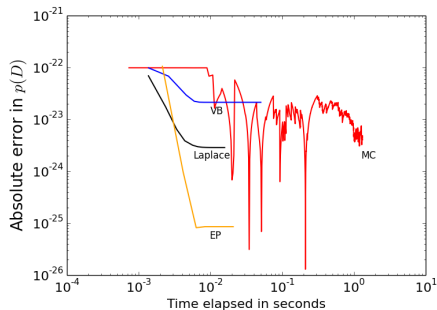
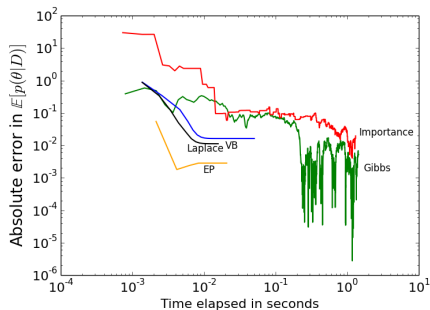
end for

until convergence

Expectation Propagation for the Clutter Problem



Summary of Methods for the Clutter Problem



T. Minka (2001a), T. Minka (2001b)

MSc Dissertation



L. Ellam.

Approximate Bayesian Inference for Machine Learning.

The MSc dissertation may be downloaded from [here](#).

Software

The Python code for implementing the various algorithms reported in this presentation may be obtained by clicking on the respective figures. All software and datasets may be obtained in zip form from [here](#).

Acknowledgments

- Professor N. Zabaras (Thesis Advisor) and Dr P. Brommer (Second Marker)
- EPSRC Strategic Package Project EP/L027682/1 for research at WCPM