



Statistical Machine Learning

Christian Walder

Machine Learning Research Group
CSIRO Data61

and

College of Engineering and Computer Science
The Australian National University

Canberra
Semester One, 2020.

Outlines

- Overview
- Introduction
- Linear Algebra
- Probability
- Linear Regression 1
- Linear Regression 2
- Linear Classification 1
- Linear Classification 2
- Kernel Methods
- Sparse Kernel Methods
- Mixture Models and EM 1
- Mixture Models and EM 2
- Neural Networks 1
- Neural Networks 2
- Principal Component Analysis
- Autoencoders
- Graphical Models 1
- Graphical Models 2
- Graphical Models 3
- Sampling
- Sequential Data 1
- Sequential Data 2

(Many figures from C. M. Bishop, "Pattern Recognition and Machine Learning")



Part XX

Approximate Inference

Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

*Variational Optimisation
applied to Inference*

Exponential Family

Expectation Propagation



- **Descriptive Statistics** describes the main features of a data set in quantitative terms.
- **Inference** : Drawing conclusions about a population from a random sample drawn from it, or, more generally, about a random process from its observed behavior during a finite period of time.
- Fitting a model to some observed data (finding the parameters of the model) is inference.

Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

*Variational Optimisation
applied to Inference*

Exponential Family

Expectation Propagation

Need for Approximate Inference



- Central task in the application of probabilistic models is the evaluation of the posterior distribution $p(\mathbf{Z} | \mathbf{X})$ of the latent variables \mathbf{Z} given the observed variables \mathbf{X} .
- This may be infeasible because
 - Posterior distribution has a too complex form for which expectations are not tractable.
 - Integration (for continuous variables) may not have a closed form solution.
 - Numerical integration (continuous variables) or sums over all possible configuration (discrete variables) is feasible for a small number of variables only.
- Need approximations.

Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

*Variational Optimisation
applied to Inference*

Exponential Family

Expectation Propagation



- **Stochastic Approximation**

- Produce **exact** results if given enough computational resources.
- However, sampling methods can be computational demanding.
- Example: Markov Chain Monte Carlo.

- **Deterministic Approximation:**

- Based on analytical approximation of the posterior distribution $p(\mathbf{Z} | \mathbf{X})$.
- Example: posterior is assumed to factorise in a particular way.
- However, these simplifying assumptions can never produce exact results for the original problem.

Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

*Variational Optimisation
applied to Inference*

Exponential Family

Expectation Propagation



- Originates from the **calculus of variations** (Euler, Lagrange, and others)
- Standard calculus uses a function to map one value to another value

$$y = f(x) \qquad x \xrightarrow{f} y$$

- A **functional** maps a function to a value.
Example: Entropy $H[p]$

$$H[p] = - \int p(x) \ln p(x) dx$$

- Functional derivative: How does the value change for infinitesimal changes of the input function.
- A large number of problems have the form: Find a function which optimises (maximises/minimises) a functional.
- Approximate solutions can be found by restricting the set of functions over which the functional will be optimised.



- Given a functional of the form

$$F[y] = \int_a^b G(y(x), y'(x), x) dx$$

with $f(a) = f_a$ and $f(b) = f_b$ fixed, find a function $y(\cdot)$ which optimises the functional under the given constraints.

- Consider an infinitesimal change of the function $y(\cdot)$ by $\eta(x)$ (use total derivative and integration by parts)

$$\begin{aligned} F[y + \epsilon\eta] &= F[y] + \epsilon \int_a^b \left\{ \frac{\partial G}{\partial y} \eta(x) + \frac{\partial G}{\partial y'} \eta'(x) \right\} dx + O(\epsilon^2) \\ &= F[y] + \epsilon \int_a^b \left\{ \frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \right\} \eta(x) dx + O(\epsilon^2) \end{aligned}$$

- Euler-Lagrange equations

$$\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) = 0$$

Calculus of Variation - Shortest path

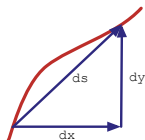


- Shortest path between two points in \mathbb{R}^N .
- In \mathbb{R}^N , the path length of a curve $y(x)$ is

$$ds = \sqrt{dx^2 + dy^2} = dx\sqrt{1 + (dy/dx)^2}$$

$$L = \int_a^b \sqrt{1 + y'(x)^2} dx$$

and therefore $G(y(x), y'(x), x) = \sqrt{1 + y'(x)^2}$



- Using the Euler-Lagrange equations

$$\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) = 0$$

- we get

$$\frac{d}{dx} y'(x) / \sqrt{1 + y'(x)^2} = 0$$

and therefore $y'(x) = \text{constant} \Rightarrow y(x) = \text{straight line}$.

Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

Variational Optimisation
applied to Inference

Exponential Family

Expectation Propagation

Variational Optimisation applied to Inference



- Variational optimisation is exact, there is no approximation involved.
- But it can be used to find approximate solutions to a given problem by **restricting** the class of functions.
- Examples: consider only quadratic functions, or linear combinations of fixed basis functions with variable coefficients.
- Probabilistic inference: assume the probability functions **factorise** in a specific way.

Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

*Variational Optimisation
applied to Inference*

Exponential Family

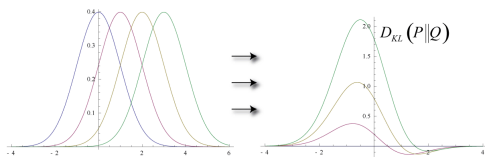
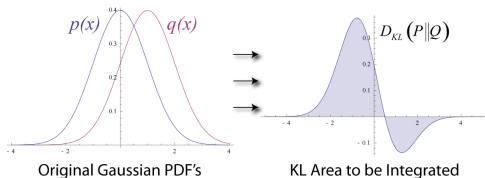
Expectation Propagation



- A widely used functional is the KL divergence between two distributions,

$$KL(p||q) \equiv D_{KL}(P||Q) \equiv \int p(x)(\log p(x) - \log q(x)) dx$$

- $KL(p||q) \geq 0$
- $KL(p||q) = 0 \Leftrightarrow p = q$
- **non-symmetric** (so not a distance metric)



Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

Variational Optimisation
applied to Inference

Exponential Family

Expectation Propagation



- From the lecture on EM

$$\ln p(\mathbf{X}) = \mathcal{L}(q) + \text{KL}(q||p)$$

$$\mathcal{L}(q) = \int q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} \right\} d\mathbf{Z}$$

$$\text{KL}(q||p) = - \int q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z} | \mathbf{X})}{q(\mathbf{Z})} \right\} d\mathbf{Z}.$$

- Here \mathbf{Z} contains
 - latent variables (as before), and
 - parameters θ which are now assumed to be stochastic (prior distribution over this variables provided).
- Given is the joint distribution $p(\mathbf{X}, \mathbf{Z})$.

Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

*Variational Optimisation
applied to Inference*

Exponential Family

Expectation Propagation



- From the lecture on EM

$$\ln p(\mathbf{X}) = \mathcal{L}(q) + \text{KL}(q||p)$$

- Maximise the lower bound of $\mathcal{L}(q)$ by optimisation.
- If all probability distributions $p(\mathbf{Z})$ are allowed this will be achieved for $\text{KL}(q||p) = 0$, or $q(\mathbf{Z}) = p(\mathbf{Z} | \mathbf{X})$.
- Assuming the true posterior is intractable, consider a restricted family $q(\mathbf{Z})$ and find a member of this family which minimises the KL divergence.

⇒ Optimising $\mathcal{L}(q)$ makes $q(\mathbf{Z}) \approx p(\mathbf{Z}|\mathbf{X})$ and bounds the model evidence $p(\mathbf{X})$

Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

*Variational Optimisation
applied to Inference*

Exponential Family

Expectation Propagation



- Partition \mathbf{Z} into disjoint sets \mathbf{Z}_i where $i = 1, \dots, M$.
- Assume $q(\mathbf{Z})$ factorises as

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

- No restriction on the functional form of the individual $q_i(\mathbf{Z}_i)$.
- This factorisation for variational inference corresponds to **Mean Field Theory** in physics, replacing the interaction of n particles with a model of one particle and a mean field (created by all the other $n - 1$ particles).

Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

*Variational Optimisation
applied to Inference*

Exponential Family

Expectation Propagation



- Insert $q(\mathbf{Z}) = \prod_{i=1}^M q_i(\mathbf{Z}_i)$ into

$$\mathcal{L}(q) = \int q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} \right\} d\mathbf{Z}$$

- Let $q_j = q_j(\mathbf{Z}_j)$ and keep $\{q_{i \neq j}\}$ fixed (only q_j can vary)

$$\begin{aligned} \mathcal{L}(q) &= \int \prod_{i=1}^M q_i \left\{ \ln p(\mathbf{X}, \mathbf{Z}) - \sum_i \ln q_i \right\} d\mathbf{Z} \\ &= \int q_j \ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) d\mathbf{Z}_j - \int q_j \ln q_j d\mathbf{Z}_j - \text{const} \end{aligned}$$

where

$$\ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) - \text{const} = \int \ln p(\mathbf{X}, \mathbf{Z}) \prod_{i \neq j} q_i d\mathbf{Z}_i = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})]$$

- $\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})]$ denotes an expectation with respect to the q distributions over all variables \mathbf{z}_i for $i \neq j$.

Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

*Variational Optimisation
applied to Inference*

Exponential Family

Expectation Propagation



- So, treating q_i for $i \neq j$ as constant,

$$\mathcal{L}(q) = \int q_j \ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) d\mathbf{Z}_j - \int q_j \ln q_j d\mathbf{Z}_j - \text{const}$$

- But $\mathcal{L}(q) + \text{const}$ can now be written as a negative Kullback-Leibler divergence

$$\mathcal{L}(q) + \text{const} = \int q_j \ln \left\{ \frac{\ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j)}{\ln q_j} \right\} = -\text{KL}(q_j \| \tilde{p}(\mathbf{X}, \mathbf{Z}_j))$$

- Maximising $\mathcal{L}(q)$ is equivalent to minimising $\text{KL}(q_j \| \tilde{p}(\mathbf{X}, \mathbf{Z}_j))$.

Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

*Variational Optimisation
applied to Inference*

Exponential Family

Expectation Propagation



Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

*Variational Optimisation
applied to Inference*

Exponential Family

Expectation Propagation

- Optimal solution for

$$\ln q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] - \text{const} \quad (3)$$

$$q_j^*(\mathbf{Z}_j) = \frac{\exp(\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})])}{\int \exp(\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})]) d\mathbf{Z}_j} \quad (4)$$

- The log of the optimal solution for factor q_j is obtained by considering the log of the joint distribution over all hidden and visible variables and taking expectations w.r.t. all other factors $\{q_i\}$ for $i \neq j$.
- Yields a set of self consistency equations.



$$\ln q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] - \text{const}$$

- 1 Initialise all factors.
- 2 Cycle through the factors and replace each with the revised estimate evaluated using the current estimate.
- 3 Convergence is guaranteed because the bound is convex w.r.t each of the factors.

Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

*Variational Optimisation
applied to Inference*

Exponential Family

Expectation Propagation



- The **exponential family** of distributions over \mathbf{x} , given parameters $\boldsymbol{\eta}$, is defined to be the set of distributions of the form

$$p(\mathbf{x} | \boldsymbol{\eta}) = h(\mathbf{x})g(\boldsymbol{\eta}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \}$$

where \mathbf{x} may be scalar or vector, and may be discrete or continuous.

- **Natural parameter** $\boldsymbol{\eta}$
- And \mathbf{u} is some function of \mathbf{x} .
- The function $g(\boldsymbol{\eta})$ can be interpreted as the coefficient ensuring normalisation

$$g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \} \, d\mathbf{x} = 1$$

- Other form with $g(\boldsymbol{\eta}) = \exp\{-G(\boldsymbol{\eta})\}$ we get

$$p(\mathbf{x} | \boldsymbol{\eta}) = h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) - G(\boldsymbol{\eta}) \}$$

Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

*Variational Optimisation
applied to Inference*

Exponential Family

Expectation Propagation

Example - Normal Distribution



- Normal distribution with mean μ and standard deviation σ

$$\mathcal{N}(x | \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/(2\sigma^2)}$$

$$\boldsymbol{\eta} = \left(\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2} \right)^T$$

$$h(x) = \frac{1}{\sqrt{2\pi}}$$

$$\mathbf{u}(x) = (x, x^2)^T$$

$$g(\boldsymbol{\eta}) = \sqrt{-2\eta_2} \exp\left(\frac{\eta_1^2}{4\eta_2}\right)$$

$$G(\boldsymbol{\eta}) = -\frac{1}{2} \ln(-2\eta_2) - \frac{\eta_1^2}{4\eta_2}$$

Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

*Variational Optimisation
applied to Inference*

Exponential Family

Expectation Propagation

Exponential Family moments from partition (1)



Differentiation of $-\ln g(\boldsymbol{\eta})$ provides the moments

$$\frac{d}{d\boldsymbol{\eta}} - \ln g(\boldsymbol{\eta}) = \mathbb{E} [\mathbf{u}(\mathbf{x})]$$

Prove using $g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \} d\mathbf{x} = 1$:

$$\begin{aligned} \frac{d}{d\boldsymbol{\eta}} - \ln g(\boldsymbol{\eta}) &= \frac{d}{d\boldsymbol{\eta}} - \ln \left[\left(\int h(\mathbf{x}) e^{\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x})} d\mathbf{x} \right)^{-1} \right] \\ &= \left(\int h(\mathbf{x}) e^{\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x})} d\mathbf{x} \right)^{-1} \frac{d}{d\boldsymbol{\eta}} \int h(\mathbf{x}) e^{\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x})} d\mathbf{x} \\ &= g(\boldsymbol{\eta}) \int \mathbf{u}(\mathbf{x}) h(\mathbf{x}) e^{\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x})} d\mathbf{x} \\ &= \int \mathbf{u}(\mathbf{x}) p(\mathbf{x}|\boldsymbol{\eta}) d\mathbf{x}. \end{aligned}$$

Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

Variational Optimisation
applied to Inference

Exponential Family

Expectation Propagation

Exponential Family moments from partition (2)



Generally differentiation of $-\ln g(\boldsymbol{\eta})$ provides the moments

$$\frac{d}{d\boldsymbol{\eta}} -\ln g(\boldsymbol{\eta}) = \mathbb{E} [\mathbf{u}(\mathbf{x})]$$

$$\frac{d^2}{d\boldsymbol{\eta}^2} -\ln g(\boldsymbol{\eta}) = \text{cov}[\mathbf{u}(\mathbf{x})]$$

etc.

or using $G(\boldsymbol{\eta})$ defined by $g(\boldsymbol{\eta}) = \exp\{-G(\boldsymbol{\eta})\}$, e.g.

$$\frac{d}{d\boldsymbol{\eta}} G(\boldsymbol{\eta}) = \mathbb{E} [\mathbf{u}(\mathbf{x})]$$

in terms of the **log partition function** G , where

$$p(\mathbf{x} | \boldsymbol{\eta}) = h(\mathbf{x})g(\boldsymbol{\eta}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \} = h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) - G(\boldsymbol{\eta}) \}$$

Exponential Family: $\min KL \Leftrightarrow$ match moments



KL minimisation \Leftrightarrow moment matching:

$$q_{\eta}(\mathbf{x}) = h(\mathbf{x}) \exp \{ \eta^T \mathbf{u}(\mathbf{x}) - G(\eta) \}$$

$$\begin{aligned} \text{KL}(p \| q_{\eta}) &\propto - \int p(\mathbf{x}) \log q_{\eta}(\mathbf{x}) \, d\mathbf{x} \\ &= - \int p(\mathbf{x}) \log [h(\mathbf{x}) \exp \{ \eta^T \mathbf{u}(\mathbf{x}) - G(\eta) \}] \, d\mathbf{x} \end{aligned}$$

$$\text{KL}(p \| q_{\eta}) = - \int p(\mathbf{x}) [\log h(\mathbf{x}) - \{ \eta^T \mathbf{u}(\mathbf{x}) - G(\eta) \}] \, d\mathbf{x}$$

$$\nabla_{\eta} \text{KL}(p \| q_{\eta}) = 0 = \int p(\mathbf{x}) \{ \mathbf{u}(\mathbf{x}) - \nabla_{\eta} G(\eta) \} \, d\mathbf{x}$$

The optimality condition is

$$\Rightarrow \int p(\mathbf{x}) \mathbf{u}(\mathbf{x}) \, d\mathbf{x} = \nabla_{\eta} G(\eta)$$

and using the property of G established earlier

$$\Rightarrow \underbrace{\mathbb{E}_{p(\mathbf{x})} [\mathbf{u}(\mathbf{x})]}_{\text{expectation of } \mathbf{u} \text{ under } p} = \underbrace{\mathbb{E}_{q_{\eta}(\mathbf{x})} [\mathbf{u}(\mathbf{x})]}_{\text{expectation of } \mathbf{u} \text{ under } q_{\eta}}$$



- Given a joint distribution over observed data \mathcal{D} and variables θ in the form of a product of (**exact**) factors

$$p(\mathcal{D}, \theta) = \prod_i f_i(\theta)$$

- Goal: Approximate the posterior distribution $p(\theta | \mathcal{D})$ by a distribution of the form of a product of (**approximate**) factors

$$q(\theta) = \frac{1}{Z} \prod_i \tilde{f}_i(\theta).$$

- Goal: Estimate the model evidence $p(\mathcal{D})$.

By letting q be the **exact** posterior for the **approximate** terms \tilde{f}_i , we can delete any single \tilde{f}_i and replace it with a refined approximation given all the other approximate terms.

Expectation Propagation with Exponential Family



- 1 Initialise all approximating factors $\tilde{f}_i(\theta)$.
- 2 Initialise the posterior approximation by setting

$$q(\theta) \propto \prod_i \tilde{f}_i(\theta)$$

- 3 Choose a factor $\tilde{f}_j(\theta)$ to refine.
- 4 Delete $\tilde{f}_j(\theta)$ from the posterior by division

$$q^{\setminus j}(\theta) = \frac{q(\theta)}{\tilde{f}_j(\theta)}$$

- 5 Evaluate the new posterior by setting the sufficient statistics (moments) of q^{new} equal to those of $q^{\setminus j} \tilde{f}_j(\theta)$ including evaluation of the normalisation constant

$$Z_j = \int q^{\setminus j}(\theta) \tilde{f}_j(\theta) d\theta$$

- 6 Evaluate and store the new factor, and goto 3.

$$\tilde{f}_j(\theta) = Z_j \frac{q^{\text{new}}(\theta)}{q^{\setminus j}(\theta)}$$



- Finally, approximate the model evidence

$$p(\mathcal{D}) \approx \int \prod_i \tilde{f}_i(\boldsymbol{\theta}) \, d\boldsymbol{\theta}$$

- Expectation propagation is not guaranteed to converge.
- A nice description of EP:

Matthias Seeger,
Expectation Propagation for Exponential Families,
Technical Report, Berkeley, 2007.

Approximate Inference

Approximation Schemes

Variational Optimisation

Calculus of Variation

*Variational Optimisation
applied to Inference*

Exponential Family

Expectation Propagation