

Approximate Inference for Neural Networks

Mauricio A. Álvarez PhD,
H.F. Garcia C. Guarnizo (TA)



Universidad Tecnológica de Pereira, Pereira, Colombia

- 1 EM algorithm**
- 2 Variational EM**
- 3 Approximate inference for Neural Networks**



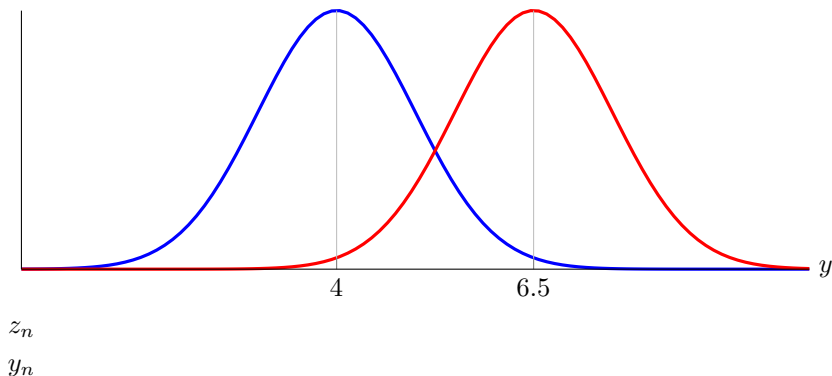
1 EM algorithm

2 Variational EM

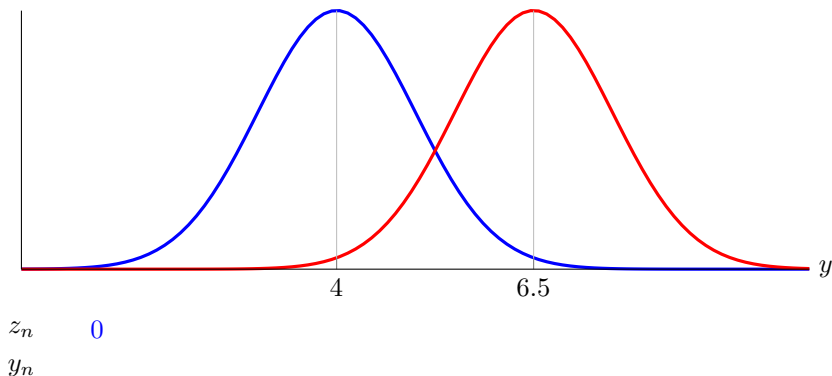
3 Approximate inference for Neural Networks



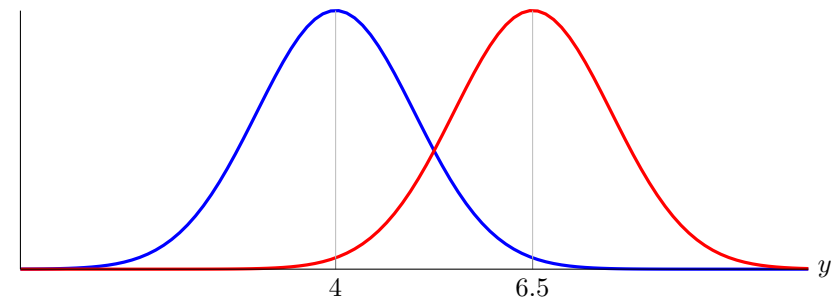
Problem definition - Example



Problem definition - Example



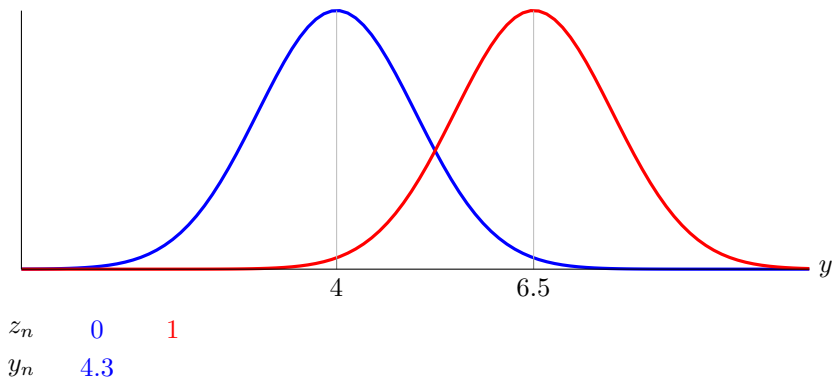
Problem definition - Example



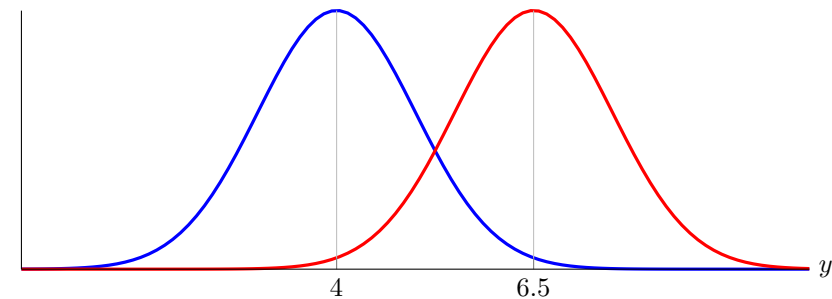
z_n 0
 y_n 4.3



Problem definition - Example



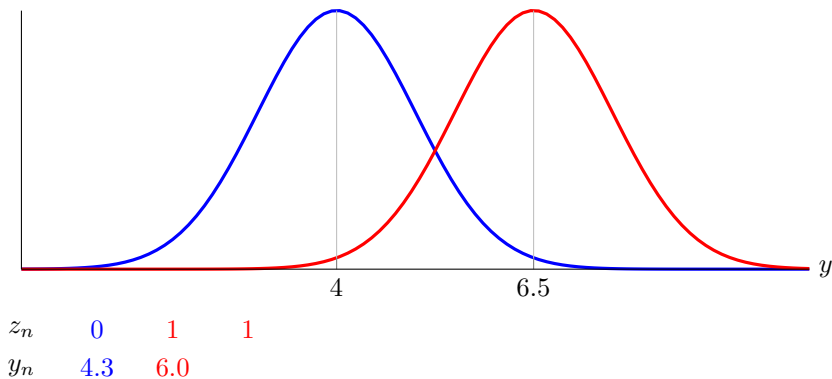
Problem definition - Example



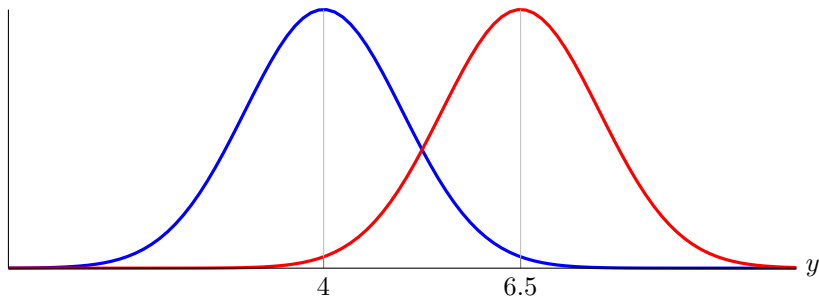
z_n	0	1
y_n	4.3	6.0



Problem definition - Example



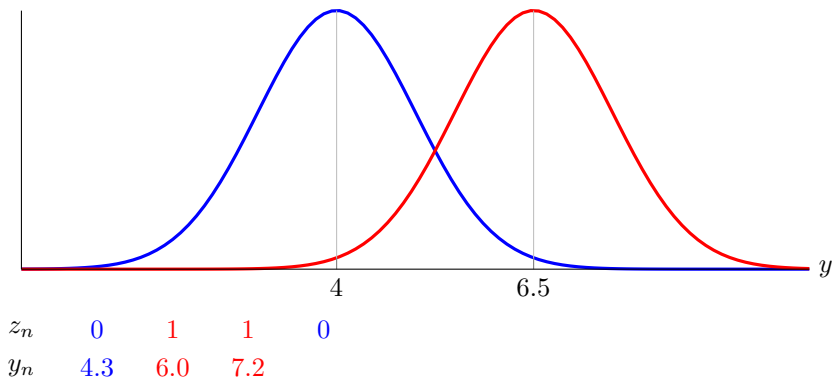
Problem definition - Example



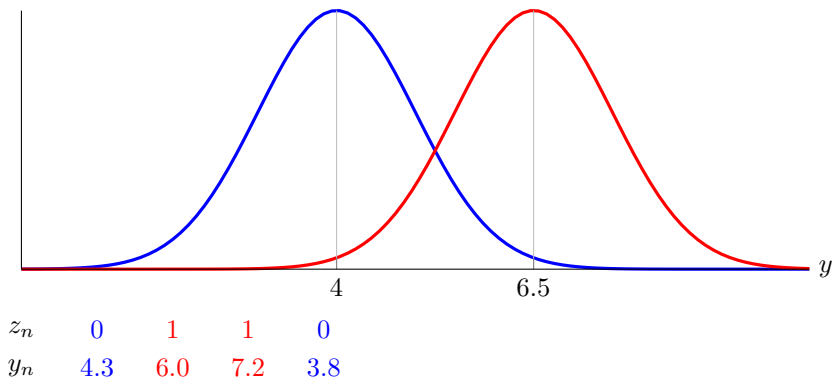
z_n	0	1	1
y_n	4.3	6.0	7.2



Problem definition - Example



Problem definition - Example



Problem definition - Example

For z :

$$p(z) = \pi^z (1 - \pi)^{1-z}, \quad p(z = 1) = \pi, \quad p(z = 0) = 1 - \pi.$$

For y :

$$p(y|z) = (\mathcal{N}(y|\mu_1, \sigma_1^2))^z (\mathcal{N}(y|\mu_2, \sigma_2^2))^{1-z}.$$

The the joint probability:

$$p(y, z) = (\pi \mathcal{N}(y|\mu_1, \sigma_1^2))^z ((1 - \pi) \mathcal{N}(y|\mu_2, \sigma_2^2))^{1-z}.$$

Additionally:

$$\begin{aligned} p(z = 1|y) &= \frac{p(y|z = 1)p(z = 1)}{p(y)} \\ &= \frac{\mathcal{N}(y|\mu_1, \sigma_1^2)\pi}{\pi \mathcal{N}(y|\mu_1, \sigma_1^2) + (1 - \pi) \mathcal{N}(y|\mu_2, \sigma_2^2)} \end{aligned}$$



Problem definition - Example

Let's assume that we have $(\mathbf{y}, \mathbf{z}) = \{y_n, z_n\}_{n=1}^N$.

The parameters $\theta = \{\pi, \mu_1, \mu_2, \sigma_1^2, \sigma_2^2\}$, can be estimated as

$$\begin{aligned}\theta_{\text{ML}} &= \arg \max_{\theta} \ln (p(\mathbf{y}, \mathbf{z}|\theta)), \\ &= \arg \max_{\theta} \ln \left(\prod_{n=1}^N p(y_n, z_n|\theta) \right), \\ &= \arg \max_{\theta} \sum_{n=1}^N \ln (p(y_n, z_n|\theta)).\end{aligned}$$

Where:

$$\begin{aligned}\ln p(y_n, z_n|\theta) &= z_n \left(\ln(\pi) + \ln \mathcal{N}(y_n|\mu_1, \sigma_1^2) \right) \\ &\quad + (1 - z_n) \left(\ln(1 - \pi) + \ln \mathcal{N}(y_n|\mu_2, \sigma_2^2) \right).\end{aligned}$$



Problem definition - Example

Let's assume that $\mathbf{y} = \{y_n\}_{n=1}^N$, \mathbf{y} \mathbf{z} is unknown (latent).
The parameters $\theta = \{\pi, \mu_1, \mu_2, \sigma_1^2, \sigma_2^2\}$, can be estimated as

$$\begin{aligned}\theta_{\text{ML}} &= \arg \max_{\theta} \sum_{n=1}^N \ln (p(y_n|\theta)) , \\ &= \arg \max_{\theta} \sum_{n=1}^N \ln \left(\sum_z p(y_n, z_n|\theta) \right) .\end{aligned}$$

Where:

$$\ln \left(\sum_z p(y_n, z_n|\theta) \right) = \ln \left(\pi \mathcal{N}(y_n|\mu_1, \sigma_1^2) + (1 - \pi) \mathcal{N}(y_n|\mu_2, \sigma_2^2) \right) .$$



Problem definition - General

The model consists of observations \mathbf{y} and a latent random variable \mathbf{z} . Then

$$\begin{aligned} p(\mathbf{y}|\theta) &= \prod_{n=1}^N p(y_n|\theta) \\ &= \prod_{n=1}^N \sum_{\mathbf{z}} p(y_n, \mathbf{z}|\theta) = \prod_{n=1}^N \sum_{\mathbf{z}} p(y_n|\mathbf{z}, \theta) p(\mathbf{z}|\theta) \end{aligned}$$

The estimation of θ_{ML} is given by

$$\begin{aligned} \theta_{\text{ML}} &= \arg \max_{\theta} \ln(p(\mathbf{y}|\theta)) \\ &= \arg \max_{\theta} \sum_{n=1}^N \ln \left(\sum_{\mathbf{z}} p(y_n|\mathbf{z}, \theta) p(\mathbf{z}|\theta) \right) \end{aligned}$$



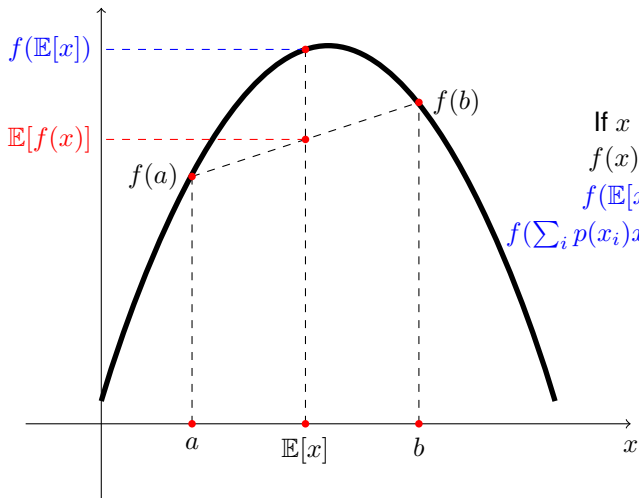
Problem definition

$$\theta_{\text{ML}} = \arg \max_{\theta} \sum_{n=1}^N \ln \left(\sum_{\mathbf{z}} p(y_n, |\mathbf{z}, \theta) p(\mathbf{z}|\theta) \right)$$

- The sum over \mathbf{z} couples the parameters θ .
- Gradients have no closed form.
- \mathbf{z} and θ are also coupled.



Jensen's Inequality



If x is a r.v. and
 $f(x)$ is concave:

$$f(\mathbb{E}[x]) \geq \mathbb{E}[f(x)]$$

$$f(\sum_i p(x_i)x_i) \geq \sum_i p(x_i)f(x_i)$$

EM algorithm - Jensen's Inequality

$$\log p(\mathbf{y}|\theta) = \sum_{n=1}^N \ln \left(\sum_{\mathbf{z}} p(y_n, \mathbf{z}|\theta) \right)$$



EM algorithm - Jensen's Inequality

$$\begin{aligned}\log p(\mathbf{y}|\theta) &= \sum_{n=1}^N \ln \left(\sum_{\mathbf{z}} p(y_n, \mathbf{z}|\theta) \right) \\ &= \sum_{n=1}^N \ln \left(\sum_{\mathbf{z}} q(\mathbf{z}) \frac{p(y_n, \mathbf{z}|\theta)}{q(\mathbf{z})} \right)\end{aligned}$$



EM algorithm - Jensen's Inequality

$$\begin{aligned}\log p(\mathbf{y}|\theta) &= \sum_{n=1}^N \ln \left(\sum_{\mathbf{z}} p(y_n, \mathbf{z}|\theta) \right) \\ &= \sum_{n=1}^N \ln \left(\sum_{\mathbf{z}} q(\mathbf{z}) \frac{p(y_n, \mathbf{z}|\theta)}{q(\mathbf{z})} \right) \\ &\geq \sum_{n=1}^N \sum_{\mathbf{z}} q(\mathbf{z}) \ln \left(\frac{p(y_n, \mathbf{z}|\theta)}{q(\mathbf{z})} \right)\end{aligned}$$



Assuming a value of $\theta^{(t)}$, what form $q(\mathbf{z})$ should have to maximize

$$\ln p(\mathbf{y}|\theta^{(t)}) \geq \sum_{n=1}^N \sum_{\mathbf{z}} q(\mathbf{z}) \ln \left(\frac{p(\mathbf{z}|y_n, \theta^{(t)})p(y_n|\theta^{(t)})}{q(\mathbf{z})} \right)$$



Assuming a value of $\theta^{(t)}$, what form $q(\mathbf{z})$ should have to maximize

$$\begin{aligned}\ln p(\mathbf{y}|\theta^{(t)}) &\geq \sum_{n=1}^N \sum_{\mathbf{z}} q(\mathbf{z}) \ln \left(\frac{p(\mathbf{z}|y_n, \theta^{(t)})p(y_n|\theta^{(t)})}{q(\mathbf{z})} \right) \\ &\geq \sum_{n=1}^N \sum_{\mathbf{z}} q(\mathbf{z}) \left[\ln \left(\frac{p(\mathbf{z}|y_n, \theta^{(t)})}{q(\mathbf{z})} \right) + \ln p(y_n|\theta^{(t)}) \right]\end{aligned}$$



Distribution $q(\mathbf{z})$

Assuming a value of $\theta^{(t)}$, what form $q(\mathbf{z})$ should have to maximize

$$\begin{aligned}\ln p(\mathbf{y}|\theta^{(t)}) &\geq \sum_{n=1}^N \sum_{\mathbf{z}} q(\mathbf{z}) \ln \left(\frac{p(\mathbf{z}|y_n, \theta^{(t)})p(y_n|\theta^{(t)})}{q(\mathbf{z})} \right) \\ &\geq \sum_{n=1}^N \sum_{\mathbf{z}} q(\mathbf{z}) \left[\ln \left(\frac{p(\mathbf{z}|y_n, \theta^{(t)})}{q(\mathbf{z})} \right) + \ln p(y_n|\theta^{(t)}) \right]\end{aligned}$$

organizing,

$$\sum_{n=1}^N \ln p(y_n|\theta^{(t)}) \geq \sum_{n=1}^N \sum_{\mathbf{z}} q(\mathbf{z}) \ln \left(\frac{p(\mathbf{z}|y_n, \theta^{(t)})}{q(\mathbf{z})} \right) + \sum_{n=1}^N \ln p(y_n|\theta^{(t)}),$$



Distribution $q(\mathbf{z})$

Assuming a value of $\theta^{(t)}$, what form $q(\mathbf{z})$ should have to maximize

$$\begin{aligned}\ln p(\mathbf{y}|\theta^{(t)}) &\geq \sum_{n=1}^N \sum_{\mathbf{z}} q(\mathbf{z}) \ln \left(\frac{p(\mathbf{z}|y_n, \theta^{(t)})p(y_n|\theta^{(t)})}{q(\mathbf{z})} \right) \\ &\geq \sum_{n=1}^N \sum_{\mathbf{z}} q(\mathbf{z}) \left[\ln \left(\frac{p(\mathbf{z}|y_n, \theta^{(t)})}{q(\mathbf{z})} \right) + \ln p(y_n|\theta^{(t)}) \right]\end{aligned}$$

organizing,

$$\sum_{n=1}^N \ln p(y_n|\theta^{(t)}) \geq \sum_{n=1}^N \sum_{\mathbf{z}} q(\mathbf{z}) \ln \left(\frac{p(\mathbf{z}|y_n, \theta^{(t)})}{q(\mathbf{z})} \right) + \sum_{n=1}^N \ln p(y_n|\theta^{(t)}),$$

then, $q(\mathbf{z}) = p(\mathbf{z}|y_n, \theta^{(t)})$.



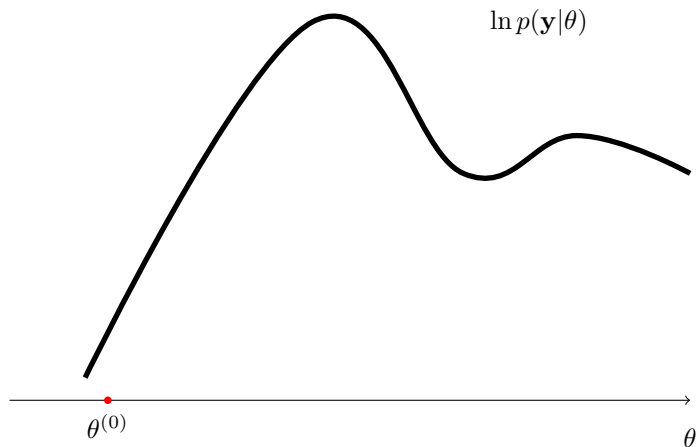
Parameter estimation

If $q(\mathbf{z}) = p(\mathbf{z}|y_n, \theta^{(t)})$,

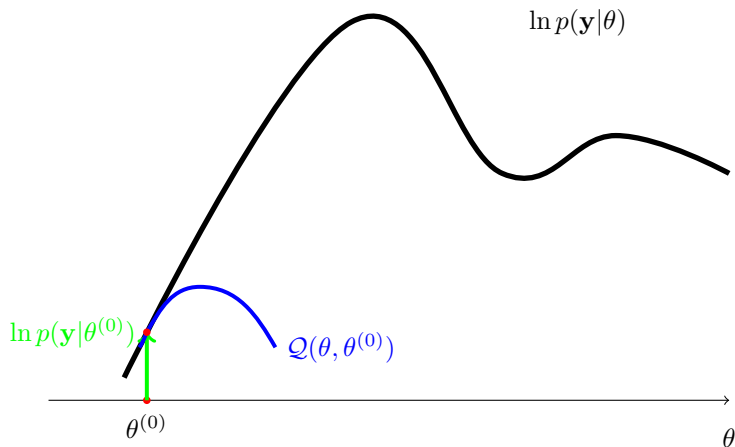
$$\begin{aligned}\theta^{(t+1)} &= \arg \max_{\theta} \sum_{n=1}^N \sum_{\mathbf{z}} p(\mathbf{z}|y_n, \theta^{(t)}) \ln \left(\frac{p(y_n, \mathbf{z}|\theta)}{p(\mathbf{z}|y_n, \theta^{(t)})} \right) \\ &= \arg \max_{\theta} \sum_{n=1}^N \mathbb{E}_{q(\mathbf{z})} [\ln (p(y_n, \mathbf{z}|\theta))] \\ &= \arg \max_{\theta} \mathcal{Q}(\theta, \theta^{(t)})\end{aligned}$$



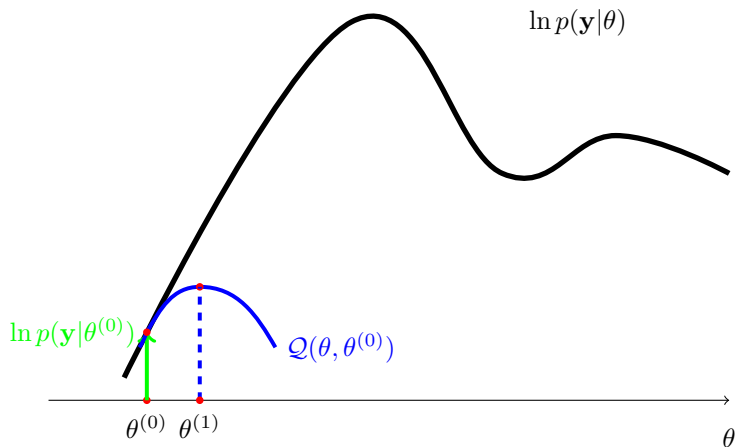
EM algorithm - Visual explanation



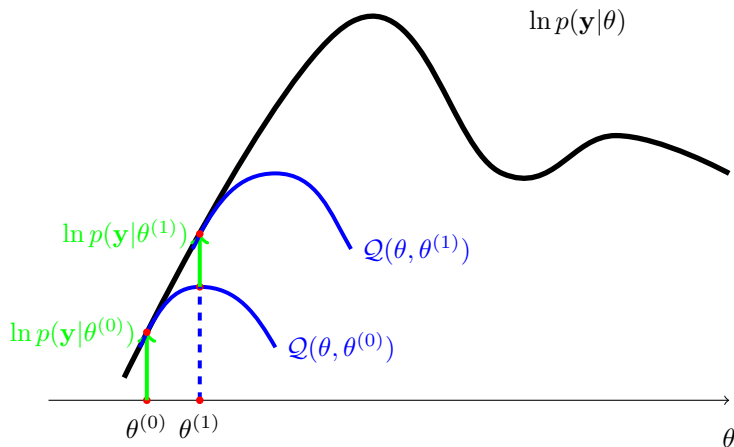
EM algorithm - Visual explanation



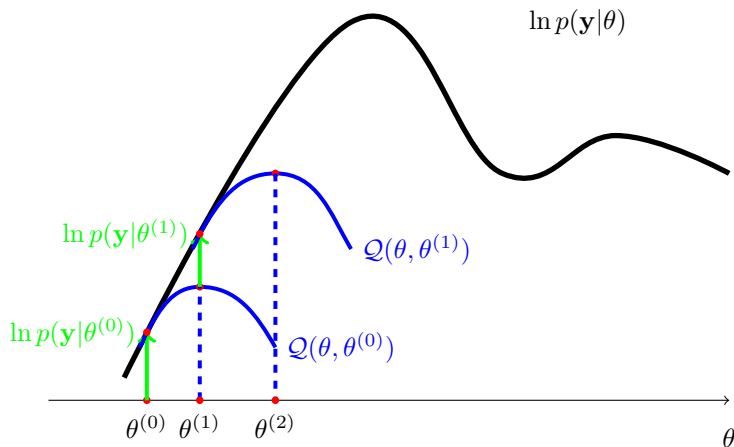
EM algorithm - Visual explanation



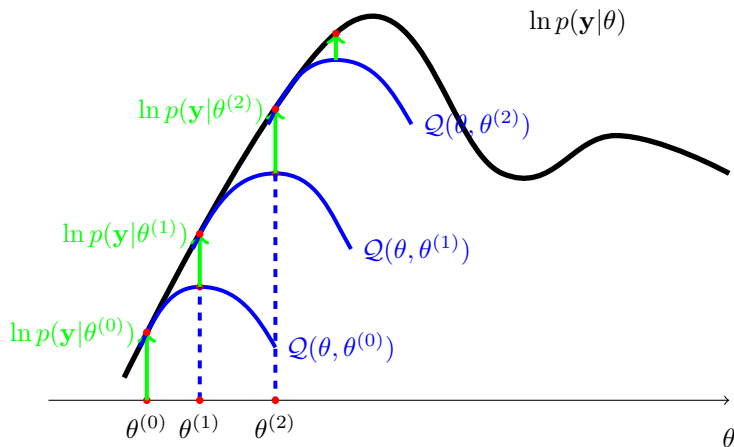
EM algorithm - Visual explanation



EM algorithm - Visual explanation



EM algorithm - Visual explanation



EM Algorithm - Pseudo-code

Given the joint distribution $p(\mathbf{y}, \mathbf{z}|\theta)$, the objective is to maximize $p(\mathbf{y}|\theta)$ w.r.t. θ .

1 Select the initial parameters $\theta^{(t)} \leftarrow \theta^{(0)}$.

2 Evaluate E-step, $q(\mathbf{z}) = p(\mathbf{z}|\mathbf{y}, \theta^{(t)})$.

3 Evaluate M-step,

$$\theta^{(t+1)} = \arg \max_{\theta} \mathcal{Q}(\theta, \theta^{(t)})$$

4 Verify convergence. If it doesn't satisfy, then

$$\theta^{(t)} \leftarrow \theta^{(t+1)},$$

return to step 2.



EM algorithm - Remark

$$\log p(\mathbf{y}|\theta) \geq \sum_{n=1}^N \sum_{\mathbf{z}} q(\mathbf{z}) \ln \left(\frac{p(y_n, \mathbf{z}|\theta)}{q(\mathbf{z})} \right)$$

$$\sum_{n=1}^N \ln p(y_n|\theta) \geq \sum_{n=1}^N \sum_{\mathbf{z}} q(\mathbf{z}) \ln \left(\frac{p(\mathbf{z}|y_n, \theta)}{q(\mathbf{z})} \right) + \sum_{n=1}^N \ln p(y_n|\theta),$$

$$\sum_{n=1}^N \ln p(y_n|\theta) \geq \sum_{n=1}^N \sum_{\mathbf{z}} q(\mathbf{z}) \ln (p(y_n|\mathbf{z}, \theta)) + \sum_{\mathbf{z}} q(\mathbf{z}) \ln \left(\frac{p(\mathbf{z}|\theta)}{q(\mathbf{z})} \right).$$

$$\sum_{n=1}^N \ln p(y_n|\theta) \geq \sum_{n=1}^N \sum_{\mathbf{z}} q(\mathbf{z}) \ln (p(y_n, \mathbf{z}|\theta)) - \sum_{\mathbf{z}} q(\mathbf{z}) \ln (q(\mathbf{z})).$$



EM algorithm - Summary

- 1 We just need to calculate the expected value of the joint distribution w.r.t. the posterior of the latent variables.
- 2 We are able to estimate the parameters θ by maximizing an easier objective function.



1 EM algorithm

2 Variational EM

3 Approximate inference for Neural Networks



Problem formulation

In this case,

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

is intractable. Most of the time because

$$\begin{aligned} p(\mathbf{y}) &= \int_{\mathbf{z}} p(\mathbf{y}, \mathbf{z}) \, d\mathbf{z} \\ &= \int_{\mathbf{z}} p(\mathbf{y}|\mathbf{z})p(\mathbf{z}) \, d\mathbf{z} \\ &= \mathbb{E}_{p(\mathbf{z})} [p(\mathbf{y}|\mathbf{z})] \end{aligned}$$

is difficult to calculate.



Variational EM - Objective

We are interested in estimating $p(\mathbf{z}|\mathbf{y})$. This can be achieved by solving the following optimization problem,

$$q^*(\mathbf{z}) = \arg \min \text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{y}))$$

where

$$\text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{y})) = \int q(\mathbf{z}) \ln \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{y})} d\mathbf{z}$$



Variational EM - Objective

Analysing the KL divergence between $q(\mathbf{z})$ and $p(\mathbf{z}|\mathbf{y})$,

$$\begin{aligned}\text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{y})) &= \int q(\mathbf{z}) \ln \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{y})} d\mathbf{z} \\ &= \int q(\mathbf{z}) \ln q(\mathbf{z}) d\mathbf{z} - \int q(\mathbf{z}) \ln p(\mathbf{z}|\mathbf{y}) d\mathbf{z} \\ &= \int q(\mathbf{z}) \ln q(\mathbf{z}) d\mathbf{z} - \int q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \mathbf{y})}{p(\mathbf{y})} d\mathbf{z} \\ &= \int q(\mathbf{z}) \ln q(\mathbf{z}) d\mathbf{z} - \int q(\mathbf{z}) \ln p(\mathbf{z}, \mathbf{y}) d\mathbf{z} + \ln p(\mathbf{y})\end{aligned}$$



Variational EM - The evidence lower bound

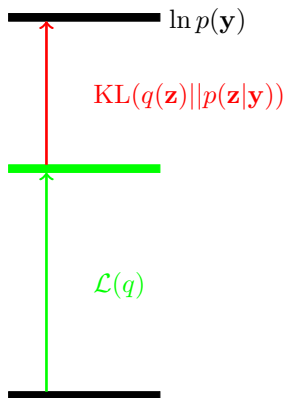
From previous equation we have,

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

Given that $\text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{y})) \geq 0$, then the lower bound of $\ln p(\mathbf{y})$ is $\mathcal{L}(q)$.



Variational EM - The evidence lower bound



Variational EM - Summary

- 1 If the posterior $p(\mathbf{Z}|\mathbf{y})$ is intractable or complex (main issue is to calculate $p(\mathbf{y})$).
- 2 We can approximate the posterior by minimizing $\text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{y}))$, which results in maximizing $\mathcal{L}(q)$.
- 3 The E-step is an iterative process.



$$\text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{y})) = \int q(\mathbf{z}) \ln \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{y})} \mathrm{d}\mathbf{z}$$

$$\text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{y})) = \int q(\mathbf{z}) \ln q(\mathbf{z}) \mathrm{d}\mathbf{z} - \int q(\mathbf{z}) \ln p(\mathbf{z}, \mathbf{y}) \mathrm{d}\mathbf{z} + \ln p(\mathbf{y})$$

$$\text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{y})) = - \int q(\mathbf{z}) \ln p(\mathbf{y}|\mathbf{z}) \mathrm{d}\mathbf{z} + \int q(\mathbf{z}) \ln \frac{q(\mathbf{z})}{p(\mathbf{z})} \mathrm{d}\mathbf{z} + \ln p(\mathbf{y})$$



1 EM algorithm

2 Variational EM

3 Approximate inference for Neural Networks



Approx. NNs - Problem definition

Let's assume we have a dataset $\mathcal{D} = \{\mathbf{X}, \mathbf{Y}\}$, with $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ and $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$

$$\begin{aligned}\text{KL}(q_{\theta}(\boldsymbol{\omega}) \| p(\boldsymbol{\omega} | \mathcal{D})) &\propto - \int q_{\theta}(\boldsymbol{\omega}) \log p(\mathbf{Y} | \mathbf{X}, \boldsymbol{\omega}) d\boldsymbol{\omega} + \text{KL}(q_{\theta}(\boldsymbol{\omega}) \| p(\boldsymbol{\omega})), \\ &= - \sum_{i=1}^N \int q_{\theta}(\boldsymbol{\omega}) \log p(\mathbf{y}_i | \mathbf{f}^{\boldsymbol{\omega}}(\mathbf{x}_i)) d\boldsymbol{\omega} + \text{KL}(q_{\theta}(\boldsymbol{\omega}) \| p(\boldsymbol{\omega})),\end{aligned}$$

where $\boldsymbol{\omega}$ is a vector composed by the stacking of all weights.



Approx. NNs - Problem definition

- 1 The summed-over terms log-likelihood are not tractable for BNNs with more than a single hidden layer.
- 2 This objective requires us to perform computations over the entire dataset, which can be too costly for large N .



We can reduce the data size problem by data sub-sampling (also referred to as mini-batch optimization).

$$\hat{\mathcal{L}}_{\text{VI}}(\theta) := -\frac{N}{M} \sum_{i \in S} \int q_{\theta}(\omega) \log p(\mathbf{y}_i | \mathbf{f}^{\omega}(\mathbf{x}_i)) d\omega + \text{KL}(q_{\theta}(\omega) \| p(\omega))$$

with a random index set S of size M .



Approx. NNs - Data size solution

We can reduce the data size problem by data sub-sampling (also referred to as mini-batch optimization).

$$\hat{\mathcal{L}}_{\text{VI}}(\theta) := -\frac{N}{M} \sum_{i \in S} \int q_{\theta}(\omega) \log p(\mathbf{y}_i | \mathbf{f}^{\omega}(\mathbf{x}_i)) d\omega + \text{KL}(q_{\theta}(\omega) \| p(\omega))$$

with a random index set S of size M .



Approx. NNs - Expected value solution

We can reduce the data size problem by data sub-sampling (also referred to as mini-batch optimization).

$$\mathbb{E}_{q_{\theta}(\boldsymbol{\omega})}[\log p(\mathbf{y}_i | \mathbf{f}^{\boldsymbol{\omega}}(\mathbf{x}_i))] = \int q_{\theta}(\boldsymbol{\omega}) \log p(\mathbf{y}_i | \mathbf{f}^{\boldsymbol{\omega}}(\mathbf{x}_i)) d\boldsymbol{\omega}$$

$$q_{\theta}(\boldsymbol{\omega}) = \prod_{i=1}^L q_{\theta}(\mathbf{W}_i) = \prod_{i=1}^L \prod_{j=1}^{K_i} \prod_{k=1}^{K_{i+1}} q_{m_{ijk}, \sigma_{ijk}}(w_{ijk}) = \prod_{i,j,k} \mathcal{N}(w_{ijk}; m_{ijk}, \sigma_{ijk}^2)$$

$$\mathbb{E}_{q_{\theta}(\boldsymbol{\omega})}[\log p(\mathbf{y}_i | \mathbf{f}^{\boldsymbol{\omega}}(\mathbf{x}_i))] \approx \frac{1}{M} \sum_{m=1}^M q_{\theta}(\boldsymbol{\omega}_m) \log p(\mathbf{y}_i | \mathbf{f}^{\boldsymbol{\omega}_m}(\mathbf{x}_i))$$



Approx. NNs - Expected value solution

A better way is to use Pathwise Gradient Estimator

$$\mathbb{E}_{q_{\theta}(\omega)} [f_{\theta}(\omega)] \longrightarrow \mathbb{E}_{q_0(\epsilon)} [f_{\theta}(g(\theta; \epsilon))]$$

$$q_{ijk}(w_{ijk}) = g(\theta_{ijk}, \epsilon_{ijk}) = m_{ijk} + \sigma_{ijk}\epsilon_{ijk}$$

then

$$\hat{\mathcal{L}}_{\text{VI}}(\theta) := -\frac{N}{M} \sum_{i \in S} \int q_{\theta}(\omega) \log p(\mathbf{y}_i | \mathbf{f}^{\omega}(\mathbf{x}_i)) d\omega + \text{KL}(q_{\theta}(\omega) \| p(\omega))$$

becomes

$$\hat{\mathcal{L}}_{\text{MC}}(\theta) = -\frac{N}{M} \sum_{i \in S} \log p(\mathbf{y}_i | \mathbf{f}^{g(\theta, \epsilon)}(\mathbf{x}_i)) + \text{KL}(q_{\theta}(\omega) \| p(\omega))$$

where $\mathbb{E}_{S, \epsilon} [\hat{\mathcal{L}}_{\text{MC}}(\theta)] = \mathcal{L}_{\text{VI}}(\theta)$.



Algorithm 1 Minimise divergence between $q_{\theta}(\omega)$ and $p(\omega|X, Y)$

- 1: Given dataset \mathbf{X}, \mathbf{Y} ,
- 2: Define learning rate schedule η ,
- 3: Initialise parameters θ randomly.
- 4: **repeat**
- 5: Sample M random variables $\hat{\epsilon}_i \sim p(\epsilon)$, S a random subset of $\{1, \dots, N\}$ of size M .
- 6: Calculate stochastic derivative estimator w.r.t. θ :

$$\widehat{\Delta\theta} \leftarrow -\frac{N}{M} \sum_{i \in S} \frac{\partial}{\partial \theta} \log p(\mathbf{y}_i | \mathbf{f}^{q(\theta, \hat{\epsilon}_i)}(\mathbf{x}_i)) + \frac{\partial}{\partial \theta} \text{KL}(q_{\theta}(\omega) || p(\omega)).$$

- 7: Update θ :
 $\theta \leftarrow \theta + \eta \widehat{\Delta\theta}$.
 - 8: **until** θ has converged.
-



$$\begin{aligned}\hat{\mathbf{y}} &= \hat{\mathbf{h}}\mathbf{M}_2 \\ &= (\mathbf{h} \odot \hat{\boldsymbol{\epsilon}}_2)\mathbf{M}_2 \\ &= (\mathbf{h} \cdot \text{diag}(\hat{\boldsymbol{\epsilon}}_2))\mathbf{M}_2 \\ &= \mathbf{h}(\text{diag}(\hat{\boldsymbol{\epsilon}}_2)\mathbf{M}_2) \\ &= \sigma(\hat{\mathbf{x}}\mathbf{M}_1 + \mathbf{b})(\text{diag}(\hat{\boldsymbol{\epsilon}}_2)\mathbf{M}_2) \\ &= \sigma((\mathbf{x} \odot \hat{\boldsymbol{\epsilon}}_1)\mathbf{M}_1 + \mathbf{b})(\text{diag}(\hat{\boldsymbol{\epsilon}}_2)\mathbf{M}_2) \\ &= \sigma(\mathbf{x}(\text{diag}(\hat{\boldsymbol{\epsilon}}_1)\mathbf{M}_1) + \mathbf{b})(\text{diag}(\hat{\boldsymbol{\epsilon}}_2)\mathbf{M}_2)\end{aligned}$$



- 1 EM algorithm
- 2 Variational EM
- 3 Approximate inference for Neural Networks



Approx. NNs - Laplace approximation

The standard Laplace approximation is obtained by taking the second-order Taylor expansion around a mode of a distribution.

If we approximate the log posterior over the weights of a network given some data \mathcal{D} around a MAP estimate ω^* , we obtain

$$\log p(\omega|\mathcal{D}) \approx \log p(\omega^*|\mathcal{D}) - \frac{1}{2} (\omega - \omega^*)^\top \overline{H} (\omega - \omega^*),$$

where $\overline{H} = \mathbb{E}[H]$ is the average Hessian of the negative log posterior. The posterior over the weights is then approximated as Gaussian:

$$\omega \sim \mathcal{N}(\omega^*, \overline{H}^{-1})$$





Bishop, C. M. (2006).

Pattern Recognition and Machine Learning (Information Science and Statistics).

Springer-Verlag, Berlin, Heidelberg.



Gal, Y. (2016).

Uncertainty in Deep Learning.

PhD thesis, University of Cambridge.



Ritter, H., Botev, A., and Barber, D. (2018).

A Scalable Laplace Approximation for Neural Networks.

In International Conference on Learning Representations.