

# Bayesian Deep Learning and Restricted Boltzmann Machines

Narada Warakagoda

Forsvarets Forskningsinstitutt

*ndw@ffi.no*

November 1, 2018

- 1 Probability Review
- 2 Bayesian Deep Learning
- 3 Restricted Boltzmann Machines

# Probability Review

- Normal (Gaussian) Distribution

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- Categorical Distribution

$$P(x) = \prod_{i=1}^k p_i^{[x=i]}$$

- Sampling

$$\mathbf{x} \sim p(\mathbf{x})$$

- Independent variables

$$p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k) = \prod_{i=1}^k p(\mathbf{x}_i)$$

- Expectation

$$\mathbb{E}_{p(\mathbf{x})} f(\mathbf{x}) = \int f(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

or for discrete variables

$$\mathbb{E}_{p(\mathbf{x})} f(\mathbf{x}) = \sum_{i=1}^k f(\mathbf{x}_i) P(\mathbf{x}_i)$$

# Kullback Leibler Distance

$$\begin{aligned} KL(q(\mathbf{x}) || p(\mathbf{x})) &= \mathbb{E}_{q(\mathbf{x})} \log \left[ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right] \\ &= \int [q(\mathbf{x}) \log q(\mathbf{x}) - q(\mathbf{x}) \log p(\mathbf{x})] d\mathbf{x} \end{aligned}$$

For the discrete case

$$KL(Q(\mathbf{x}) || P(\mathbf{x})) = \sum_{i=1}^k [Q(\mathbf{x}_i) \log Q(\mathbf{x}_i) - Q(\mathbf{x}_i) \log P(\mathbf{x}_i)]$$

# Bayesian Deep Learning

- Joint distribution

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

- Marginalization

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

$$P(\mathbf{x}) = \sum_{\mathbf{y}} P(\mathbf{x}, \mathbf{y})$$

- Conditional distribution

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



- Prediction

$$p(\mathbf{y}|\mathbf{x}, \mathbf{w}) = \mathcal{N}(\mathbf{f}_{\mathbf{w}}(\mathbf{x}), \Sigma)$$

- Classification

$$P(y|\mathbf{x}, \mathbf{w}) = \prod_{i=1}^k f_{\mathbf{w}}^i(\mathbf{x})^{[y=i]}$$

# Training Criteria

- Maximum Likelihood(ML)

$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} p(\mathbf{Y}|\mathbf{X}, \mathbf{w})$$

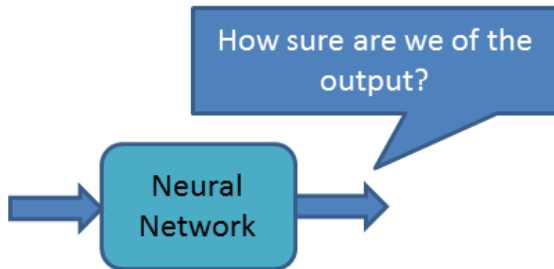
- Maximum A-Priori (MAP)

$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} p(\mathbf{Y}, \mathbf{w}|\mathbf{X}) = \arg \max_{\mathbf{w}} p(\mathbf{Y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w})$$

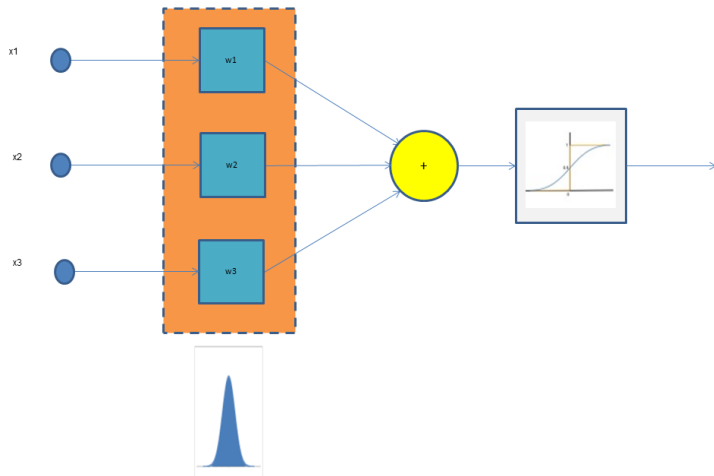
- Bayesian

$$p(\mathbf{w}|\mathbf{Y}, \mathbf{X}) = \frac{p(\mathbf{Y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w})}{P(\mathbf{Y}|\mathbf{X})} = \frac{p(\mathbf{Y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w})}{\int P(\mathbf{Y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}}$$

# Motivation for Bayesian Approach



# Motivation for Bayesian Approach



# Uncertainty with Bayesian Approach

- Not only prediction/classification, but their uncertainty can also be calculated
  - Since we have  $p(\mathbf{w}|\mathbf{Y}, \mathbf{X})$  we can sample  $\mathbf{w}$  and use each sample as network parameters in calculating the prediction/classification  $p(\hat{y}|\hat{x}, \mathbf{w})$  (i.e. network output for a given input ).
  - Prediction/classification is the mean of  $p(\hat{y}|\hat{x}, \mathbf{w})$

$$p_{out} = p(\hat{y}|\hat{x}, \mathbf{Y}, \mathbf{X}) = \int p(\hat{y}|\hat{x}, \mathbf{w}) p(\mathbf{w}|\mathbf{Y}, \mathbf{X}) d\mathbf{w}$$

- Uncertainty of prediction/classification is the variance of  $p(\hat{y}|\hat{x}, \mathbf{w})$

$$\text{Var}(p(\hat{y}|\hat{x}, \mathbf{w})) = \int [p(\hat{y}|\hat{x}, \mathbf{w}) - p_{out}]^2 p(\mathbf{w}|\mathbf{Y}, \mathbf{X}) d\mathbf{w}$$

- Uncertainty is important in safety critical applications (eg: self-driving cars, medical diagnosis, military applications)

# Other Advantages of Bayesian Approach

- Natural interpretation for regularization
- Model selection
- Input data selection (active learning)

# Main Challenge of Bayesian Approach

- We calculate
  - For continuous case:

$$p(\mathbf{w}|\mathbf{Y}, \mathbf{X}) = \frac{p(\mathbf{Y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w})}{\int P(\mathbf{Y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}}$$

- For discrete case:

$$P(\mathbf{w}|\mathbf{Y}, \mathbf{X}) = \frac{p(\mathbf{Y}|\mathbf{X}, \mathbf{w}) P(\mathbf{w})}{\sum_{\mathbf{w}} p(\mathbf{Y}|\mathbf{X}, \mathbf{w}) P(\mathbf{w})}$$

- Calculating denominator is often intractable
  - Eg: Consider a weight vector  $\mathbf{w}$  of 100 elements, each can have two values. Then there are  $2^{100} = 1.2 \times 10^{30}$  different weight vectors. Compare this with universe's age 13.7 billion years.
- We need approximations

# Different Approaches

- Monte Carlo techniques (Eg: Markov Chain Monte Carlo -MCMC)
- Variational Inference
- Introducing random elements in training (eg: Dropout)



# Advantages and Disadvantages of Different Approaches

- Markov Chain Monte Carlo - MCMC
  - Asymptotically exact
  - Computationally expensive
- Variational Inference
  - No guarantee of exactness
  - Possibility for faster computation

- We are interested in

$$p_{out} = \text{Mean}(p(\hat{y}|\hat{x}, \mathbf{w})) = p(\hat{y}|\hat{x}, \mathbf{Y}, \mathbf{X}) = \int p(\hat{y}|\hat{x}, \mathbf{w}) p(\mathbf{w}|\mathbf{Y}, \mathbf{X}) d\mathbf{w}$$

$$\text{Var}(p(\hat{y}|\hat{x}, \mathbf{w})) = \int [p(\hat{y}|\hat{x}, \mathbf{w}) - p_{out}]^2 p(\mathbf{w}|\mathbf{Y}, \mathbf{X}) d\mathbf{w}$$

- Both are integrals of the type

$$I = \int F(\mathbf{w}) p(\mathbf{w}|\mathcal{D}) d\mathbf{w}$$

where  $\mathcal{D} = (\mathbf{Y}, \mathbf{X})$  is training data.

- Approximate the integral by sampling  $\mathbf{w}_i$  from  $p(\mathbf{w}|\mathcal{D})$

$$I \approx \frac{1}{L} \sum_{i=1}^L F(\mathbf{w}_i).$$

- Challenge: We don't have the posterior

$$p(\mathbf{w}|\mathcal{D}) = p(\mathbf{w}|\mathbf{Y}, \mathbf{X}) = \frac{p(\mathbf{Y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w})}{\int P(\mathbf{Y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}}$$

- "Solution": Use importance sampling by sampling from a proposal distribution  $q(\mathbf{w})$

$$I = \int F(\mathbf{w}) \frac{p(\mathbf{w}|\mathcal{D})}{q(\mathbf{w})} q(\mathbf{w}) d\mathbf{w} \approx \frac{1}{L} \sum_{i=1}^L F(\mathbf{w}_i) \frac{p(\mathbf{w}_i|\mathcal{D})}{q(\mathbf{w}_i)}$$

- Problem: We still do not have  $p(\mathbf{w}|\mathcal{D})$

- Problem: We still do not have  $p(\mathbf{w}|\mathcal{D})$
- Solution: use unnormalized posterior  $\tilde{p}(\mathbf{w}|\mathcal{D}) = p(\mathbf{Y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w})$  where normalization factor  $Z = \int P(\mathbf{Y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}$  such that

$$p(\mathbf{w}|\mathcal{D}) = \frac{\tilde{p}(\mathbf{w}|\mathcal{D})}{Z}$$

- Integral can be calculated with:

$$I \approx \frac{\sum_{i=1}^L F(\mathbf{w}_i) \tilde{p}(\mathbf{w}_i|\mathcal{D}) / q(\mathbf{w}_i)}{\sum_{i=1}^L \tilde{p}(\mathbf{w}_i|\mathcal{D}) / q(\mathbf{w}_i)}$$

# Weakness of Importance Sampling

- Proposal distribution must be close to the non-zero areas of original distribution  $p(\mathbf{w}|\mathcal{D})$ .
- In neural networks,  $p(\mathbf{w}|\mathcal{D})$  is typically small except for few narrow areas.
- Blind sampling from  $q(\mathbf{w})$  has a high chance that they fall outside non-zero areas of  $p(\mathbf{w}|\mathcal{D})$
- We must actively try to get samples that lie close to  $p(\mathbf{w}|\mathcal{D})$
- Markov Chain Monte Carlo (MCMC) is such technique.

# Metropolis Algorithm

- Metropolis algorithm is an example of MCMC
- Draw samples repeatedly from random walk  $\mathbf{w}_{t+1} = \mathbf{w}_t + \epsilon$  where  $\epsilon$  is a small random vector,  $\epsilon \sim q(\epsilon)$  (eg: Gaussian noise)
- Drawn sample at  $t = t$  is either accepted based on the ratio  $\frac{\tilde{p}(\mathbf{w}_t|\mathcal{D})}{\tilde{p}(\mathbf{w}_{t-1}|\mathcal{D})}$ 
  - If  $\tilde{p}(\mathbf{w}_t|\mathcal{D}) > \tilde{p}(\mathbf{w}_{t-1}|\mathcal{D})$  accept sample
  - If  $\tilde{p}(\mathbf{w}_t|\mathcal{D}) < \tilde{p}(\mathbf{w}_{t-1}|\mathcal{D})$  accept sample with probability  $\frac{\tilde{p}(\mathbf{w}_t|\mathcal{D})}{\tilde{p}(\mathbf{w}_{t-1}|\mathcal{D})}$
  - If sample accepted use it for calculating  $I$
- Can use the same formula for calculating  $I$

$$I \approx \frac{\sum_{i=1}^L F(\mathbf{w}_i) \tilde{p}(\mathbf{w}_i|\mathcal{D}) / q(\mathbf{w}_i)}{\sum_{i=1}^L \tilde{p}(\mathbf{w}_i|\mathcal{D}) / q(\mathbf{w}_i)}$$

- Hybrid Monte Carlo (Hamiltonian Monte Carlo)
  - Similar to Metropolis algorithm
  - But uses gradient information rather than a random walk.
- Simulated Annealing

- Goal: computation of posterior  $p(\mathbf{w}|\mathcal{D})$ , i.e. the parameters of the neural network  $\mathbf{w}$  given data  $\mathcal{D} = (\mathbf{Y}, \mathbf{X})$
- But this computation is often intractable
- Idea: find a distribution  $q(\mathbf{w})$  from a family of distributions  $Q$  such that  $q(\mathbf{w})$  can closely approximate  $p(\mathbf{w}|\mathcal{D})$
- How to measure the distance between  $q(\mathbf{w})$  and  $p(\mathbf{w}|\mathcal{D})$  ?
  - Kullback-Leibler Distance  $\text{KL}(q(\mathbf{w})||p(\mathbf{w}|\mathcal{D}))$
- The problem can be formulated as

$$\hat{p}(\mathbf{w}|\mathcal{D}) = \arg \min_{q(\mathbf{w})} \text{KL}(q(\mathbf{w})||p(\mathbf{w}|\mathcal{D}))$$



# Minimizing KL Distance

- Using the definition of KL distance

$$\text{KL}(q(\mathbf{w})||p(\mathbf{w}|\mathcal{D})) = \int q(\mathbf{w}) \ln \frac{q(\mathbf{w})}{p(\mathbf{w}|\mathcal{D})} d\mathbf{w}$$

- Cannot minimize this directly, because we do not know  $p(\mathbf{w}|\mathcal{D})$
- But we can manipulate it further, and transform it to another equivalent optimization problem involving a quantity known as Evidence Lower Bound (ELBO)

# Evidence Lower Bound (ELBO)

$$\begin{aligned}\text{KL}(q(\mathbf{w})||p(\mathbf{w}|\mathcal{D})) &= \int q(\mathbf{w}) \ln \frac{q(\mathbf{w})}{p(\mathbf{w}|\mathcal{D})} d\mathbf{w} \\ &= \int q(\mathbf{w}) \ln \frac{q(\mathbf{w}) p(\mathcal{D})}{p(\mathbf{w}, \mathcal{D})} d\mathbf{w} \\ &= \int q(\mathbf{w}) \ln \frac{q(\mathbf{w})}{p(\mathbf{w}, \mathcal{D})} d\mathbf{w} + \int q(\mathbf{w}) \ln p(\mathcal{D}) d\mathbf{w} \\ &= \mathbb{E}_{q(\mathbf{w})} \ln \frac{q(\mathbf{w})}{p(\mathbf{w}, \mathcal{D})} + \ln p(\mathcal{D}) \int q(\mathbf{w}) d\mathbf{w}\end{aligned}$$

$$\ln p(\mathcal{D}) = \boxed{\mathbb{E}_{q(\mathbf{w})} \ln \frac{p(\mathbf{w}, \mathcal{D})}{q(\mathbf{w})}} + \text{KL}(q(\mathbf{w})||p(\mathbf{w}|\mathcal{D}))$$

- Since  $\ln p(\mathcal{D})$  is constant, minimizing  $\text{KL}(q(\mathbf{w})||p(\mathbf{w}|\mathcal{D}))$  is equivalent to maximizing ELBO

# Another Look at ELBO

$$\begin{aligned}\text{ELBO} &= \mathbb{E}_{q(\mathbf{w})} \ln \frac{p(\mathbf{w}, \mathcal{D})}{q(\mathbf{w})} \\ &= \int q(\mathbf{w}) \ln p(\mathbf{w}, \mathcal{D}) d\mathbf{w} - \int q(\mathbf{w}) \ln q(\mathbf{w}) d\mathbf{w} \\ &= \int q(\mathbf{w}) \ln [p(\mathcal{D}|\mathbf{w})p(\mathbf{w})] d\mathbf{w} - \int q(\mathbf{w}) \ln q(\mathbf{w}) d\mathbf{w} \\ &= \int q(\mathbf{w}) \ln p(\mathcal{D}|\mathbf{w}) d\mathbf{w} - \int q(\mathbf{w}) \ln \frac{q(\mathbf{w})}{p(\mathbf{w})} d\mathbf{w} \\ &= \mathbb{E}_{q(\mathbf{w})} p(\mathcal{D}|\mathbf{w}) - \text{KL}(q(\mathbf{w})||p(\mathbf{w}))\end{aligned}$$

- We maximize ELBO with respect to  $q(\mathbf{w})$
- First term  $\mathbb{E}_{q(\mathbf{w})} p(\mathcal{D}|\mathbf{w})$  is equivalent to maximizing  $q(\mathbf{w})$ 's ability explain training data
- Second term  $\text{KL}(q(\mathbf{w})||p(\mathbf{w}))$  is equivalent to minimizing  $q(\mathbf{w})$ 's distance to  $p(\mathbf{w})$

# Outline of Procedure with ELBO

- Start with ELBO

$$\text{ELBO} = \mathcal{L} = \mathbb{E}_{q(\mathbf{w})} \ln \frac{p(\mathbf{w}, \mathcal{D})}{q(\mathbf{w})} = \mathbb{E}_{q(\mathbf{w})} [\ln p(\mathbf{w}, \mathcal{D}) - \ln q(\mathbf{w})]$$

- Rewrite with parameter  $\lambda$  of  $q(\mathbf{w})$  and expand expectation

$$\mathcal{L}(\lambda) = \int \ln[p(\mathbf{w}, \mathcal{D})]q(\mathbf{w}, \lambda) d\mathbf{w} - \int \ln[q(\mathbf{w}, \lambda)]q(\mathbf{w}, \lambda) d\mathbf{w}$$

- Maximize  $\mathcal{L}(\lambda)$  with respect to  $\lambda$

$$\lambda^* = \arg \max_{\lambda} \mathcal{L}(\lambda)$$

- Use the optimized  $q$  with respect to  $\lambda$  as posterior

$$q(\mathbf{w}, \lambda^*) = p(\mathbf{w}, \mathcal{D})$$

# How to Maximize ELBO

- Analytical methods are not practical for deep neural networks
- We resort to gradient methods with Monte Carlo sampling
- We discuss two methods:
  - Black box variational inference: Based on log derivative trick
  - Bayes by Backprop: Based on re-parameterization trick

# Black Box Variational Inference

- Start with ELBO:

$$\mathcal{L}(\lambda) = \int \ln[p(\mathbf{w}, \mathcal{D})]q(\mathbf{w}, \lambda) d\mathbf{w} - \int \ln[q(\mathbf{w}, \lambda)]q(\mathbf{w}, \lambda) d\mathbf{w}$$

- Differentiate with respect to  $\lambda$ .

$$\begin{aligned}\nabla_{\lambda}\mathcal{L}(\lambda) &= \int \ln[p(\mathbf{w}, \mathcal{D})]\nabla_{\lambda}[q(\mathbf{w}, \lambda)]d\mathbf{w} \\ &\quad - \int \ln[q(\mathbf{w}, \lambda)]\nabla_{\lambda}[q(\mathbf{w}, \lambda)]d\mathbf{w} \\ &\quad - \int \nabla_{\lambda}[\ln[q(\mathbf{w}, \lambda)]]q(\mathbf{w}, \lambda) d\mathbf{w}\end{aligned}$$

- The last term is zero (Can you prove it?)

# Black Box Variational Inference

- Now we have

$$\begin{aligned}\nabla_{\lambda}\mathcal{L}(\lambda) &= \int \ln[p(\mathbf{w}, \mathcal{D})]\nabla_{\lambda}[q(\mathbf{w}, \lambda)]d\mathbf{w} \\ &\quad - \int \ln[q(\mathbf{w}, \lambda)]\nabla_{\lambda}[q(\mathbf{w}, \lambda)]d\mathbf{w} \\ &= \int \left[ [p(\mathbf{w}, \mathcal{D})] - \ln[q(\mathbf{w}, \lambda)] \right] \nabla_{\lambda}[q(\mathbf{w}, \lambda)]d\mathbf{w}\end{aligned}$$

- We want to write this as an expectation with respect to  $q$
- Use the log derivative trick

$$\nabla_{\lambda}[q(\mathbf{w}, \lambda)] = \nabla_{\lambda}[\ln q(\mathbf{w}, \lambda)]q(\mathbf{w}, \lambda)$$

# Black Box Variational Inference

- Now we get

$$\begin{aligned}\nabla_{\lambda}\mathcal{L}(\lambda) &= \int \ln[p(\mathbf{w}, \mathcal{D})]\nabla_{\lambda}[\ln q(\mathbf{w}, \lambda)]q(\mathbf{w}, \lambda) d\mathbf{w} \\ &\quad - \int \ln[q(\mathbf{w}, \lambda)]\nabla_{\lambda}[\ln q(\mathbf{w}, \lambda)]q(\mathbf{w}, \lambda) d\mathbf{w}\end{aligned}$$

- Rearranging terms

$$\nabla_{\lambda}\mathcal{L}(\lambda) = \int \left[ \ln[p(\mathbf{w}, \mathcal{D})] - \ln q(\mathbf{w}, \lambda) \right] \nabla_{\lambda}[\ln q(\mathbf{w}, \lambda)]q(\mathbf{w}, \lambda) d\mathbf{w}$$

- This is the same as Expectation with respect to  $q$

$$\nabla_{\lambda}\mathcal{L}(\lambda) = \mathbb{E}_{q(\mathbf{w}, \lambda)} \left[ \ln[p(\mathbf{w}, \mathcal{D})] - \ln q(\mathbf{w}, \lambda) \right] \nabla_{\lambda}[\ln q(\mathbf{w}, \lambda)]$$



# BBVI optimization procedure

- Assume a distribution  $q(\mathbf{w}, \lambda)$  parameterized by  $\lambda$ .
- Draw  $S$  samples of  $\mathbf{w}$  from the distribution using the current value of  $\lambda = \lambda_t$
- Estimate the gradient of ELBO using the sample values:

$$\nabla_{\lambda} \hat{\mathcal{L}}(\lambda) = \frac{1}{S} \sum_{s=1}^S \left[ \ln[p(\mathbf{w}^s, \mathcal{D})] - \ln q(\mathbf{w}^s, \lambda) \right] \nabla_{\lambda} [\ln q(\mathbf{w}^s, \lambda)]$$

- Update  $\lambda$

$$\lambda_{t+1} = \lambda_t + \rho \nabla_{\lambda} \hat{\mathcal{L}}(\lambda)$$

- repeat from step 2

- Try to approximate ELBO directly by sampling from the  $q(\mathbf{w}, \lambda)$

$$\text{ELBO} = \mathcal{L}(\lambda) = \mathbb{E}_{q(\mathbf{w}, \lambda)} [\ln p(\mathbf{w}, \mathcal{D}) - \ln q(\mathbf{w}, \lambda)]$$

with

$$\hat{\mathcal{L}}(\lambda) = \frac{1}{S} \sum_{s=1}^S [\ln p(\mathbf{w}^s, \mathcal{D}) - \ln q(\mathbf{w}^s, \lambda)]$$

- But we need  $\nabla_{\lambda} \hat{\mathcal{L}}(\lambda)$  and we can not differentiate  $\hat{\mathcal{L}}(\lambda)$  because it is not a smooth function of  $\lambda$
- Use the re-parameterization trick

$$\mathbf{w}^s = \mathbf{w}(\lambda, \epsilon^s)$$

where  $\epsilon^s$  is drawn from for example a standard Gaussian distribution.

# Bayes by BackProp (BbB)

- The estimated ELBO now

$$\hat{\mathcal{L}}(\lambda) = \frac{1}{S} \sum_{s=1}^S [\ln p(\mathbf{w}(\lambda, \epsilon^s), \mathcal{D}) - \ln q(\mathbf{w}(\lambda, \epsilon^s), \lambda)]$$

- Now this is a smooth function of  $\lambda$  and can differentiate

$$\nabla_{\lambda} \hat{\mathcal{L}}(\lambda) = \frac{1}{S} \sum_{s=1}^S \left[ \frac{\partial \hat{\mathcal{L}}_s}{\partial \mathbf{w}} \frac{\partial \mathbf{w}}{\partial \lambda} + \frac{\partial \hat{\mathcal{L}}_s}{\partial \lambda} \right]$$

where  $\hat{\mathcal{L}}_s = \ln p(\mathbf{w}(\lambda, \epsilon^s), \mathcal{D}) - \ln q(\mathbf{w}(\lambda, \epsilon^s), \lambda)$

- Once the gradients are known, optimum  $\lambda^*$  and hence  $q(\mathbf{w}, \lambda^*)$  can be found by gradient descent.

# Performance of BBVI and BbB

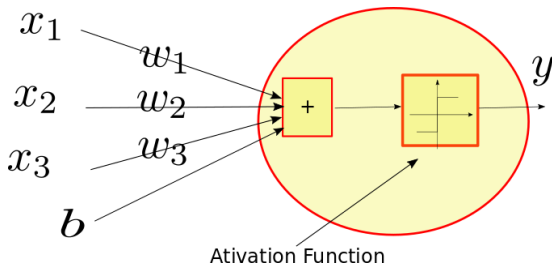
- Both methods estimate approximate gradients by sampling
- High variance of the estimated gradients is a problem
- In practice, these algorithms need modifications to tackle high variance
- BbB tends to have a lower variance estimates than BBVI

# Bayesian Deep Learning through Randomization in Training

- Stochastic gradient descent and Dropout can be given Bayesian interpretations
- Dropout procedure in testing can be used for estimating the uncertainty of model outputs (Monte Carlo Dropout).
  - Enable dropout and feed the network  $S$  times with data and collect the outputs  $f(s)$ ,  $s = 1, 2, \dots, S$
  - Output variance =  $\frac{1}{S} \sum_s (f(s) - \bar{f}(s))^2$  where  $\bar{f}(s) = \frac{1}{S} \sum_s f(s)$

# Restricted Boltzmann Machines

# Stochastic Neurons



Deterministic Neuron  $y = \sigma(b + \sum_i w_i x_i)$

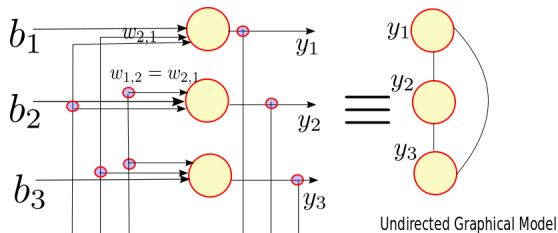
Stochastic Neuron  $p(y = 1) = \sigma(b + \sum_i w_i x_i)$

- We consider stochastic binary neurons, i.e.  $y$  can be either 1 or 0

$$p(y = 1) = \sigma(b + \sum_i w_i x_i)$$

$$p(y = 0) = 1 - p(y = 1)$$

# Boltzmann Machine

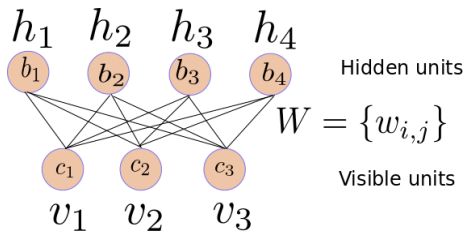


Stochastic Recurrent Neural Network

- A Boltzmann machine is a recurrent network with stochastic neurons
- Weights are symmetrical
- At the equilibrium, the relationships of the neuron outputs can be represented using an undirected graphical model



# Restricted Boltzmann Machine (RBM)



- Neurons are divided into two groups: Visible and Hidden
- Restricted architecture: No connections within visible group or hidden group
- Network parameters:
  - Bias vector hidden units,  $\mathbf{b} = [b_1, b_2, \dots, b_H]$
  - Bias vector visible units,  $\mathbf{c} = [c_1, c_2, \dots, c_V]$
  - Connection weights,  $W = \{w_{i,j}\}$
- Network values are binary random vectors:  $\mathbf{v} = [v_1, v_2, \dots, v_V]$  and  $\mathbf{h} = [h_1, h_2, \dots, h_H]$

# How the network parameters and values are related?

- Through the definition of an Energy function
- In RBM, the energy function is defined as

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{h}^T \mathbf{W} \mathbf{v} - \mathbf{c}^T \mathbf{v} - \mathbf{b}^T \mathbf{h}$$

- We assign probabilities to  $(\mathbf{v}, \mathbf{h})$  based on Boltzmann distribution

$$p(\mathbf{v}, \mathbf{h}) = \frac{\exp(-E(\mathbf{v}, \mathbf{h}))}{Z}$$

where

$$Z = \sum_{\mathbf{v}', \mathbf{h}'} \exp(-E(\mathbf{v}', \mathbf{h}'))$$

# What can we do with RBM?

Assume that the network parameters  $\mathbf{W}, \mathbf{b}, \mathbf{c}$  are known.

- Can we calculate the probability of a given pair of vectors  $(\hat{\mathbf{v}}, \hat{\mathbf{h}})$ ?
  - This is generally not tractable, because calculating  $Z$  requires to sum all combinations  $\mathbf{v}$  and  $\mathbf{h}$  values.
- Can we calculate the probability of  $\mathbf{h}$  given  $\mathbf{v}$  or vice-versa?
  - Yes, this is "inference" and possible.

Assume that a data set of  $\mathbf{v}$  vectors given.

- Can we estimate the network parameters  $\mathbf{W}, \mathbf{b}, \mathbf{c}$  ?
  - Yes, this is training and possible

We want to find  $p(\mathbf{h}|\mathbf{v})$  assuming  $\mathbf{W}, \mathbf{b}, \mathbf{c}$  are known.

- We start with the Bayes rule

$$\begin{aligned} p(\mathbf{h}|\mathbf{v}) &= \frac{p(\mathbf{h}|\mathbf{v})}{\sum_{\mathbf{h}'} p(\mathbf{h}', \mathbf{v}')} \\ &= \frac{\exp(\mathbf{h}^T \mathbf{W} \mathbf{v} + \mathbf{c}^T \mathbf{v} + \mathbf{b}^T \mathbf{h}) / Z}{\sum_{\mathbf{h}' \in \{0,1\}^H} \exp(\mathbf{h}'^T \mathbf{W} \mathbf{v}' + \mathbf{c}^T \mathbf{v}' + \mathbf{b}^T \mathbf{h}') / Z} \end{aligned}$$

- Canceling common factors and expanding vector-matrix multiplication as a summation

$$p(\mathbf{h}|\mathbf{v}) = \frac{\exp\left(\sum_j (h_j \mathbf{W}_j \mathbf{v} + b_j h_j)\right)}{\sum_{h'_1 \in \{0,1\}} \sum_{h'_2 \in \{0,1\}} \cdots \sum_{h'_H \in \{0,1\}} \exp(\sum_j (h'_j \mathbf{W}_j \mathbf{v} + b_j h'_j))}$$

We want to find  $p(\mathbf{h}|\mathbf{v})$  assuming  $\mathbf{W}$ ,  $\mathbf{b}$ ,  $\mathbf{c}$  are known.

- Writing exponential of sums as product of exponentials

$$\begin{aligned}
 p(\mathbf{h}|\mathbf{v}) &= \frac{\prod_j (\exp(h_j \mathbf{W}_j \mathbf{v} + b_j h_j))}{\sum_{h'_1 \in \{0,1\}} \sum_{h'_2 \in \{0,1\}} \cdots \sum_{h'_H \in \{0,1\}} \prod_j (\exp(h'_j \mathbf{W}_j \mathbf{v} + b_j h'_j))} \\
 &= \frac{\prod_j (\exp(h_j \mathbf{W}_j \mathbf{v} + b_j h_j))}{(\sum_{h'_1 \in \{0,1\}} \exp(h'_1 \mathbf{W}_1 \mathbf{v} + b_1 h'_1)) \cdots (\sum_{h'_H \in \{0,1\}} \exp(h'_H \mathbf{W}_H \mathbf{v} + b_H h'_H))} \\
 &= \frac{\prod_j (\exp(h_j \mathbf{W}_j \mathbf{v} + b_j h_j))}{\prod_j (\sum_{h'_j \in \{0,1\}} \exp(h'_j \mathbf{W}_j \mathbf{v} + b_j h'_j))} \\
 &= \frac{\prod_j (\exp(h_j \mathbf{W}_j \mathbf{v} + b_j h_j))}{\prod_j \exp(\mathbf{W}_j \mathbf{v} + b_j)} \\
 &= \prod_j \frac{(\exp(h_j \mathbf{W}_j \mathbf{v} + b_j h_j))}{\exp(\mathbf{W}_j \mathbf{v} + b_j)}
 \end{aligned}$$

- This implies that calculation of  $p(\mathbf{h}|\mathbf{v})$  is tractable

- Let's try to interpret

$$p(\mathbf{h}|\mathbf{v}) = \prod_j \frac{\exp(h_j \mathbf{W}_j \mathbf{v} + b_j h_j)}{\exp(\mathbf{W}_j \mathbf{v} + b_j)}$$

- Consider the quantity  $q(h_j) = \frac{\exp(h_j \mathbf{W}_j \mathbf{v} + b_j h_j)}{\exp(\mathbf{W}_j \mathbf{v} + b_j)}$
- $q(h_j)$  takes two values,  $q(0)$  and  $q(= 1)$ . And sum of these values are 1. Therefore it is a probability measure of  $h_j$ .
- Since we assumed  $\mathbf{v}$  is given,  $q(h_j)$  is actually  $p(h_j|\mathbf{v})$
- A simple manipulation shows that  $p(h_j = 1) = \sigma(\mathbf{W}_j \mathbf{v} + b_j)$  i.e. The activation function of a stochastic neuron.

- We consider maximum likelihood training with a given dataset  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N\}$  with respect to the log likelihood  $L = \log \prod_i^N p(\mathbf{v}_i) = \sum_i^N \log p(\mathbf{v}_i)$
- We use gradient descent and therefore calculate  $\frac{\partial L}{\partial \theta}$  the gradient of  $L$  with respect to a model parameter  $\theta$
- Derive the gradient for a single sample  $\frac{\partial \log(p(\mathbf{v}))}{\partial \theta}$

- By definition we know that

$$p(\mathbf{v}, \mathbf{h}) = \frac{\exp(-E(\mathbf{v}, \mathbf{h}))}{Z} \quad (1)$$

where

$$Z = \sum_{\mathbf{v}', \mathbf{h}'} \exp(-E(\mathbf{v}', \mathbf{h}')) \quad (2)$$

- Therefore

$$p(\mathbf{v}) = \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}) = \sum_{\mathbf{h}} \frac{\exp(-E(\mathbf{v}, \mathbf{h}))}{Z} \quad (3)$$

- Take log and differentiate wrt  $\theta$

$$\frac{\partial \log p(\mathbf{v})}{\partial \theta} = \frac{\partial \log \sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h}))}{\partial \theta} - \frac{\partial \log Z}{\partial \theta} \quad (4)$$



- Consider the first term

$$\frac{\partial \log \sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h}))}{\partial \theta} = - \frac{\sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h})) \frac{\partial(E(\mathbf{v}, \mathbf{h}))}{\partial \theta}}{\sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h}))} \quad (5)$$

$$= - \sum_{\mathbf{h}} \frac{\exp(-E(\mathbf{v}, \mathbf{h}))}{\sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h}))} \frac{\partial(E(\mathbf{v}, \mathbf{h}))}{\partial \theta} \quad (6)$$

- But dividing equation 1 by equation 3 we get

$$\frac{p(\mathbf{v}, \mathbf{h})}{p(\mathbf{v})} = p(\mathbf{h}|\mathbf{v}) = \frac{\exp(-E(\mathbf{v}, \mathbf{h}))}{\sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h}))} \quad (7)$$

- Substitute equation 7 in equation 6

$$\frac{\partial \log \sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h}))}{\partial \theta} = - \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}) \frac{\partial(E(\mathbf{v}, \mathbf{h}))}{\partial \theta} \quad (8)$$

# Gradients

- Consider the second term in equation 4 and substitute for  $Z$  from equation 2

$$\frac{\partial \log Z}{\partial \theta} = \frac{\partial \log \sum_{\mathbf{v}', \mathbf{h}'} \exp(-E(\mathbf{v}', \mathbf{h}'))}{\partial \theta} \quad (9)$$

$$= - \frac{\sum_{\mathbf{v}', \mathbf{h}'} \exp(-E(\mathbf{v}', \mathbf{h}')) \frac{\partial(E(\mathbf{v}', \mathbf{h}'))}{\partial \theta}}{\sum_{\mathbf{v}', \mathbf{h}'} \exp(-E(\mathbf{v}', \mathbf{h}'))} \quad (10)$$

$$= - \sum_{\mathbf{v}, \mathbf{h}} \frac{\exp(-E(\mathbf{v}, \mathbf{h}))}{\sum_{\mathbf{v}', \mathbf{h}'} \exp(-E(\mathbf{v}', \mathbf{h}'))} \frac{\partial(E(\mathbf{v}, \mathbf{h}))}{\partial \theta} \quad (11)$$

- From equations 1 and 2 it is clear that  $\frac{\exp(-E(\mathbf{v}, \mathbf{h}))}{\sum_{\mathbf{v}', \mathbf{h}'} \exp(-E(\mathbf{v}', \mathbf{h}'))}$  is

$p(\mathbf{v}, \mathbf{h})$

- Therefore

$$\frac{\partial \log Z}{\partial \theta} = - \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial(E(\mathbf{v}, \mathbf{h}))}{\partial \theta} \quad (12)$$

- From equations 4, 8 and 12

$$\frac{\partial \log p(\mathbf{v})}{\partial \theta} = - \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}) \frac{\partial(E(\mathbf{v}, \mathbf{h}))}{\partial \theta} + \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial(E(\mathbf{v}, \mathbf{h}))}{\partial \theta} \quad (13)$$

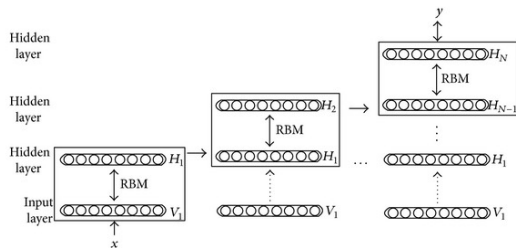
$$\boxed{\frac{\partial \log p(\mathbf{v})}{\partial \theta} = -\mathbb{E}_{p(\mathbf{h}|\mathbf{v})} \left[ \frac{\partial(E(\mathbf{v}, \mathbf{h}))}{\partial \theta} \right] + \mathbb{E}_{p(\mathbf{v}, \mathbf{h})} \left[ \frac{\partial(E(\mathbf{v}, \mathbf{h}))}{\partial \theta} \right]} \quad (14)$$

- The first term of equation 14
  - Known as positive phase
  - Depends on training data
  - Can be computed exactly
- The second term of equation 14
  - Known as negative phase
  - independent of training data, completely model dependent
  - Must be estimated through Gibb's sampling and a procedure known as Contrastive Divergence

# Applications of RBMs

- Deep belief networks
- Collaborative filtering

# Deep Belief Networks



## Method for initializing a multilayer network

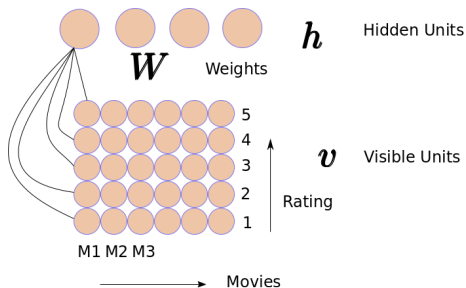
- 1 Train an RBM with training data
- 2 Initialize the current layer with the trained parameters
- 3 Present training data to the RBM and sample the hidden layer values
- 4 Use the hidden layer values as training data and repeat from step 1.

# Collaborative Filtering

	M1	M2	M3	M4	M5	M6
U1				3		
U2	5		1			
U3		3	5			
U4	4		?			5
U5			4			
U6					2	

- Application in recommendation systems. Eg: Movie rating/recommendation
- Different users rate different items (eg: movies) using a rating scale such as 1 to 5
- Problem is to estimate the rating for an unrated item by a given user

# Collaborative filtering with RBM



- Train a different RBM for each user. But share weights across users
- Visible units correspond to the ratings given to each movie
- In training movies with missing ratings are omitted
- For prediction of a missing rating, find  $p(\mathbf{h}|\mathbf{v})$  and back to  $p(\mathbf{v}|\mathbf{h})$

# The End