

Bayesian Deep Learning

Narada Warakagoda

Forsvarets Forskningsinstitutt

ndw@ffi.no

October 19, 2020

- 1 Probability Review
- 2 Bayesian Deep Learning

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)

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

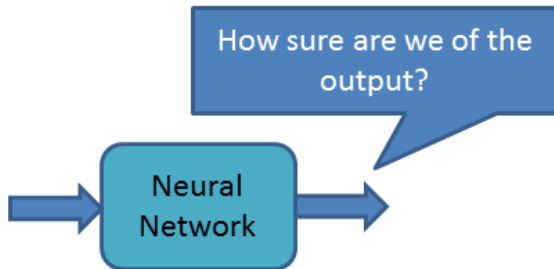
- Maximum A-Posteriori (MAP)

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

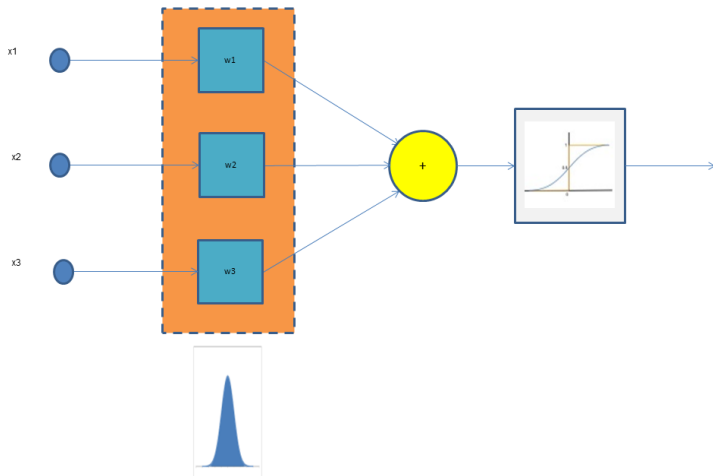
- Bayesian

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

Motivation for Bayesian Approach



Motivation for Bayesian Approach



Types of Uncertainty

- Epistemic Uncertainty
 - Due to lack of data and modelling error.
 - Model parameter distributions can tackle this.
 - We consider this.
- Aleatoric Uncertainty
 - Due to poor quality data.
 - Need to model observation quality (noise).
 - We do not consider this.

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)

Bayesian Approach vs ML and MAP

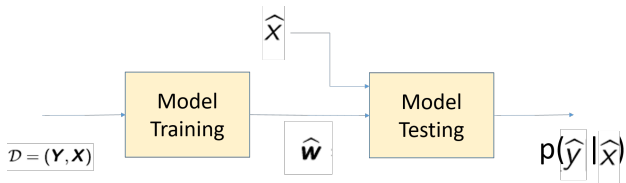


Figure: ML and MAP

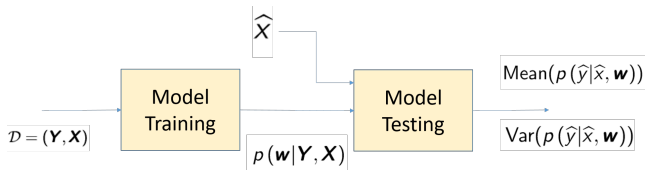


Figure: Bayesian approach

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
- Ensembles (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

Monte Carlo Techniques

- 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 ϵ 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
- Because $\frac{\tilde{p}(\mathbf{w}_t|\mathcal{D})}{\tilde{p}(\mathbf{w}_{t-1}|\mathcal{D})} = \frac{\frac{\tilde{p}(\mathbf{w}_t|\mathcal{D})}{p(\mathcal{D})}}{\frac{\tilde{p}(\mathbf{w}_{t-1}|\mathcal{D})}{p(\mathcal{D})}} = \frac{p(\mathbf{w}_t|\mathcal{D})}{p(\mathbf{w}_{t-1}|\mathcal{D})}$, sampling is valid for $p(\mathbf{w}|\mathcal{D})$ too.
- Since we sample \mathbf{w}_i from $p(\mathbf{w}|\mathcal{D})$, approximate the integral with

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

Other Monte Carlo and Related Techniques

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

Variational Inference

- 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

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

- We maximize ELBO with respect to $q(\mathbf{w})$
- First term $\mathbb{E}_{q(\mathbf{w})} \ln p(\mathcal{D}|\mathbf{w})$ 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 λ 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 λ

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

- Use the optimized q with respect to λ 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 λ .

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

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

Black Box Variational Inference

- Now we have

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

- 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

- Using log derivative trick, we get

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

- 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 λ .
- 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 λ

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

- 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 λ
- Use the re-parameterization trick

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

where ϵ^s is drawn from for example a standard Gaussian distribution.

Bayes by BackProp (BbB)

- The estimated ELBO now

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

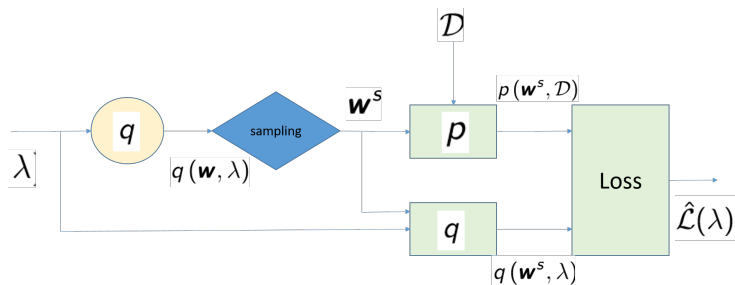
- Now this is a smooth function of λ and can differentiate

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

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

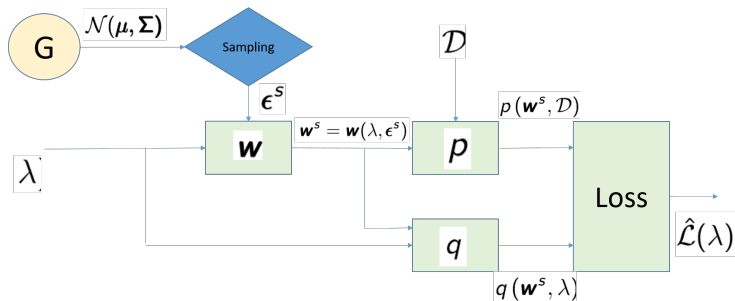
- Once the gradients are known, optimum λ^* and hence $q(\mathbf{w}, \lambda^*)$ can be found by gradient descent.

Bayes by Backprop - Schematic 1



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

Bayes by Backprop - Schematic 2



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

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 Ensembles

- Direct ensembles
- Indirect ensembles- Dropout

Bayesian Deep Learning with Direct Ensembles

- Train a set of models (say S models) with the same data set, but with different sets of initial values.
- Feed each network S with the test 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 $\bar{f}(s) = \frac{1}{S} \sum_s f(s)$

Bayesian Deep Learning with Dropout

- 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 test 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 $\bar{f}(s) = \frac{1}{S} \sum_s f(s)$