

Bayesian Neural Networks

Sargur Srihari
srihari@cedar.buffalo.edu

Topics discussed here

1. Why Bayesian?
2. Difficulty of exact Bayesian treatment and need for approximation
3. Two approximate approaches
 - Variational
 - Laplace (one discussed here)
4. Bayesian neural network for regression
 - Posterior parameter distribution
 - Hyper-parameter optimization
5. Bayesian neural network for classification

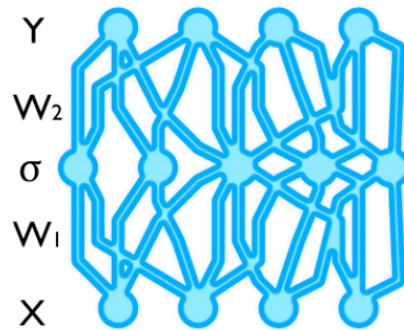
Why Bayesian?

- More complex models fit data better but generalize poorly
 - Linear with two free parameters, quadratic with three, cubic with four?
- Occam's razor says that unnecessarily complex models should not be preferred to simpler ones
- Neural networks are popular but notoriously lack objective grounding
- Bayesian approach allows different models to be compared (no of hidden units)

Classical Neural Network for Regression

- A neural network (deep learning too)
 - linearly transforms its input (bottom layer)
 - applies some non-linearity on each dimension (middle layer), and linearly transforms it again (top layer).
 - This model gives us point estimates with no uncertainty information.

Neural network for Regression.
Output activation f is identity function



$$y(x, w) = f\left(\sum_{j=1}^M w_j \phi_j(x)\right)$$

$$y_k(x, w) = f\left(\sum_{j=1}^M w_{kj}^{(2)} \sigma\left(\sum_{i=1}^D w_{ji}^{(1)} x_i\right)\right)$$

Sum of squared error/max likelihood:

$$E(w) = \frac{1}{2} \sum_{n=1}^N \|y(x_n, w) - t_n\|^2$$

Gradient Descent:

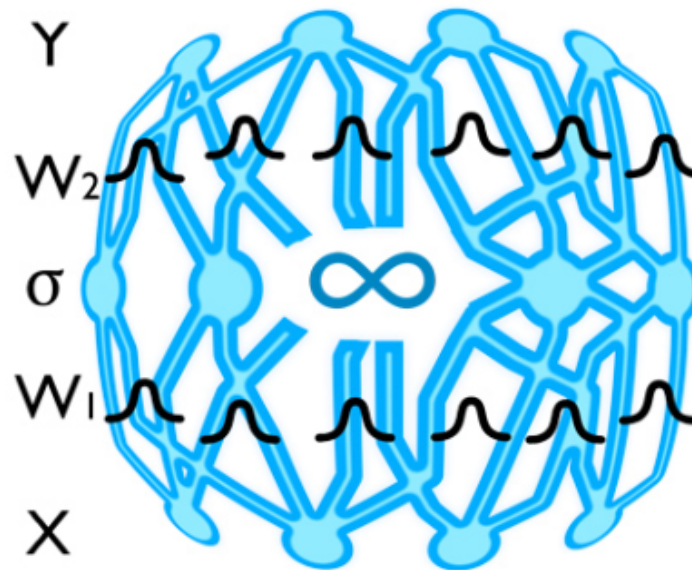
$$w^{(\tau+1)} = w^{(\tau)} - \eta \nabla E(w^{(\tau)})$$

$$\nabla E(w) = \begin{bmatrix} \frac{\partial E}{\partial w_0} \\ \frac{\partial E}{\partial w_1} \\ \frac{\partial E}{\partial w_T} \end{bmatrix}$$

Bayesian Neural Network

- A network with *infinitely* many weights with a distribution on each weight is a Gaussian process. The same network with *finitely* many weights is known as a *Bayesian neural network*

Distribution over
Weights induces a
Distribution over
outputs



Classical and Bayesian neural networks

- Classical neural networks use maximum likelihood

- To determine network parameters (weights and biases)
- Regularized maximum likelihood

$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- Is equivalent to MAP (maximum a posteriori) with Gaussian noise
- **Prior** $p(\mathbf{w}) = N(\mathbf{w} | \mathbf{m}_0, \alpha^{-1} I)$

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

- Bayesian treatment marginalizes over distribution of parameters in order to make prediction

Need for Approximation in Bayesian treatment

- In simple linear regression problem, under assumption of Gaussian noise
 - Posterior is Gaussian and evaluated exactly
 - Predictive distribution found in closed form
- In multilayered network
 - Highly nonlinear dependence of network function on parameter values
 - Exact Bayesian treatment not possible
 - Log of posterior distribution is non-convex implying multiple local minima of error function
 - Approximate methods are therefore necessary

Two Approaches To Bayesian Treatment

1. *Variational inference* using a factorized Gaussian approximation to the posterior distribution
 - Using a full covariance Gaussian
2. Most complete treatment is based on *Laplace approximation* and this is what is discussed here
 - Involves two approximations
 1. Replace posterior by a Gaussian centered at a mode of true posterior
 2. Covariance of Gaussian is small so that
 - Network function is approximately linear wrt parameters over region of parameter space for which the posterior probability is significantly nonzero

Regression: Posterior Parameter Distribution

- Simple case:
 - Predict single continuous target t from vector \mathbf{x} of inputs
 - Assume $p(t|\mathbf{x})$ is Gaussian with precision β
 - Mean is given by output of neural network $y(\mathbf{x}, \mathbf{w})$

$$p(t | \mathbf{x}, \mathbf{w}, \beta) = N(t | y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

- Prior over weights assumed Gaussian

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

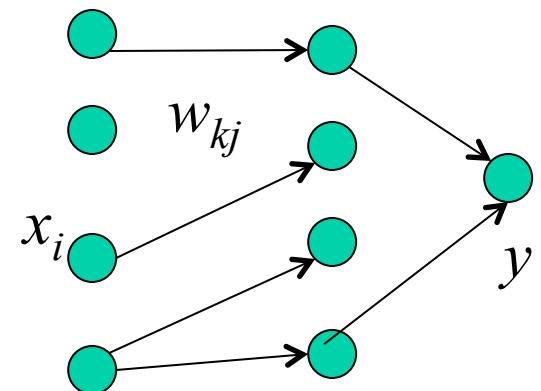
- Likelihood function

$$p(D | \mathbf{w}, \beta) = \prod_{n=1}^N N(t_n | y(\mathbf{x}_n, \mathbf{w}), \beta^{-1})$$

- Posterior is

$$p(\mathbf{w} | D, \alpha, \beta) \propto p(\mathbf{w} | \alpha) p(D | \mathbf{w}, \beta)$$

- Because of non-linear dependence of $y(\mathbf{x}, \mathbf{w})$ on \mathbf{w} it is non-Gaussian
- Therefore use Gaussian approximation using Laplace



Laplace approximation to posterior

- Find local maximum of posterior
 - Done using iterative numerical optimization
- Convenient to Maximize the logarithm of the posterior

$$\ln p(\mathbf{w} | D) = -\frac{\alpha}{2} \mathbf{w}^T \mathbf{w} - \frac{\beta}{2} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}) - t_n\}^2$$

- which corresponds to the regularized sum-of-squares error
- The maximum is denoted as \mathbf{w}_{MAP} which is found using nonlinear optimization methods
 - Which needs derivatives which are evaluated using backpropagation
- Having found the mode \mathbf{w}_{MAP} , we can build a local Gaussian approximation
- Matrix of second derivatives of negative log posterior is

$$A = -\nabla \nabla \ln p(\mathbf{w} | D, \alpha, \beta) = \alpha I + \beta H$$

Gaussian approximation to posterior

- Using algorithm for approximating the Hessian H

$$q(w | D) = N(w | w_{MAP}, A^{-1})$$

- Where A is in terms of the Hessian of the sum-of-squares error function
- Predictive distribution is obtained by marginalizing wrt this posterior distribution
 - Which leads to

$$p(t | x, D) = \int p(t | x, w) q(w | D) dw$$

$$p(t | x, D, \alpha, \beta) = N(t | y(x, w_{MAP}), \sigma^2(x))$$

- where $\sigma^2(x)$ is the input-dependent variance which is expressed in terms of β , A and g , the derivative of $y(x, w_{MAP})$ wrt w

Hyperparameter Optimization

- Practical procedure for determining α and β
- Point estimates are obtained for α and β by maximizing $\ln p(D|\alpha, \beta)$ given by

$$p(D|\alpha, \beta) = \int p(D|w, \beta) p(w|\alpha) dw$$

- Gaussian approximation using the Laplacian approach is used

$$\ln p(D|\alpha, \beta) \simeq -E(w_{MAP}) - \frac{1}{2} \ln |A| + \frac{W}{2} \ln \alpha + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

- where

$$E(w_{MAP}) = \frac{\beta}{2} \sum_{n=1}^N \{y(x_n, w_{MAP}) - t_n\}^2 + \frac{\alpha}{2} w_{MAP}^T w_{MAP}$$

Comparison of Models

- Neural networks with different numbers of hidden units
- Need to evaluate model evidence $p(D)$
- Approximated by using

$$\ln p(D | \alpha, \beta) \simeq -E(w_{MAP}) - \frac{1}{2} \ln |A| + \frac{W}{2} \ln \alpha + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

- And substituting values of α and β obtained by iterative optimization of these hyper-parameters

Bayesian Neural Network for Classification

- We have used Laplace approximation to develop a Bayesian treatment of neural network regression
- Next, modifications needed to go from regression to classification
- Consider single logistic sigmoid output for 2-classes

Evidence Framework

Two class data set

Green: Optimal Decision boundary

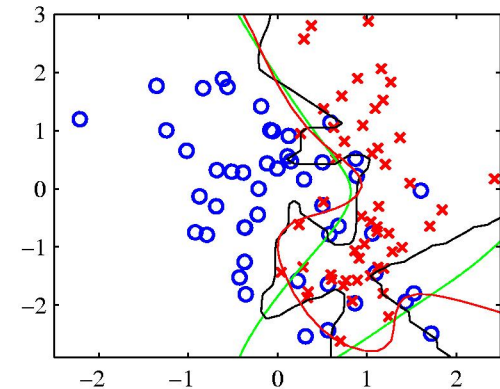
Black: two-layer neural network

with 8 hidden units by maximum likelihood

Red: including a regularizer

in which α is optimized using the evidence procedure

Starting from initial value $\alpha = 0$

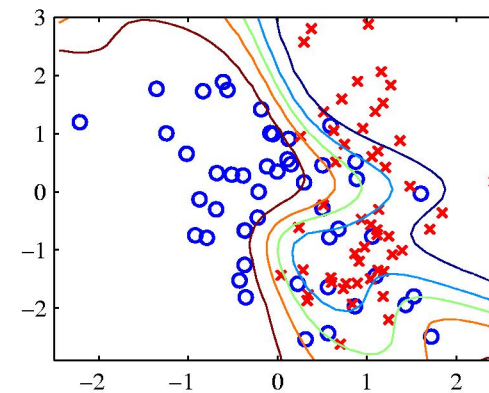


Evidence procedure greatly reduces over-fitting

Example of Laplace Approximation of a Bayesian Neural Network

- Eight hidden units
 - with tanh activation function
 - Single logistic sigmoid output
- Green: $y = 0.5$ decision boundary
- Others: probabilities of $y=0.1$, 0.3, 0.7 and 0.9
- Effect of marginalization is to leave the $y=0.5$ contour unaffected but spread the others out
 - Predictions are less confident

Simple Approximation based on point Estimate w_{MAP} of parameters



Based on predictive distribution

