

## Uncertainty analysis for neural network

What I need is uncertainty analysis of neural network. Because bayesian method doesn't seem to work. I instead do direct uncertainty analysis. Let's define probability of weights through error rate

$$p(\boldsymbol{\omega}|\text{data}) = N(\mathbf{y}_i - \mathbf{f}(\mathbf{x}_i, \boldsymbol{\omega})|0, \mathbf{x}_i, \mathbf{y}_i, \sigma^2)$$

And choose  $\sigma^2$  from precomputed optimal solution  $\boldsymbol{\omega}_0$  which is also a starting point. The related energy functions are

$$U(\boldsymbol{\omega}) = \frac{1}{2} \sum_i \|\mathbf{y}_i - \mathbf{f}(\mathbf{x}_i, \boldsymbol{\omega})\|^2 / \sigma^2$$
$$\nabla_{\boldsymbol{\omega}} U(\boldsymbol{\omega}) = \sum_i \nabla_{\boldsymbol{\omega}} \frac{1}{2} \|\mathbf{y}_i - \mathbf{f}(\mathbf{x}_i, \boldsymbol{\omega})\|^2 / \sigma^2$$

This method works for any number of dimensions but do not have prior information and overfits to the data.

*Implementation is in UHMC.cpp.*

## Hamiltonian Monte Carlo Bayesian Neural Network implementation notes (“Bayesian Gradient Descent”)

*This does only seem to work in 1-dimensional output for some reason. Implementation is in HMC.cpp.*

Tomas Ukkonen, 2015

Bayesian neural network extends basic feedforward network. This paper also describes general idea of “bayesian gradient descent” which arises when solving generic error minimization which maximizes likelihood of parameters given data instead of likelihood of data given parameters. This will then lead to qualitatively different and better solutions.

*NOTE: Feedforward networks are hopelessly outdated so after you get this one working somehow just move as quickly as possible to Stacked RBMs (Bayesian GB-RBM 80% done).*

The optimization method used is Hamiltonian Monte Carlo (HMC) which we assume we have a distribution  $\frac{1}{Z} p(\boldsymbol{\omega}) = e^{-U(\boldsymbol{\omega})}$  from which to sample from. Now according to bayesian inference rule we have

$$p(\boldsymbol{\omega}|\text{data}) = p(\text{data}|\boldsymbol{\omega}) p(\boldsymbol{\omega})$$

And we want to sample from  $p(\boldsymbol{\omega}|\text{data}) = \frac{1}{Z(\text{data})} e^{-U(\boldsymbol{\omega})}$  using HMC. We can rewrite data likelihood to be  $p(\mathbf{y}, \mathbf{x}|\boldsymbol{\omega}) = p(\mathbf{y}|\mathbf{x}, \boldsymbol{\omega}) p(\mathbf{x}|\boldsymbol{\omega})$ . We now define the first term to be squared error function with normally distributed noise  $\mathbf{y} \sim N(\mathbf{f}(\mathbf{x}, \boldsymbol{\omega}), \mathbf{I} \sigma^2)$  with fixed variance  $\sigma^2$ . For  $p(\mathbf{x}|\boldsymbol{\omega})$  we must use crude approximation and simply regularize/simplify the problem and model our data to have same distribution as the data by choosing  $\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$  and select similarly prior for weights to be gaussian ball with unit variance  $\frac{1}{2} \|\boldsymbol{\omega}\|^2$ . The known covariance matrix  $\boldsymbol{\Sigma}$  can be obtained before sampling  $\boldsymbol{\omega}$  by using some pre-optimization method and calculating the related covariance or by using gibbs sampling to sample from  $p(\sigma^2|\boldsymbol{\omega}, \text{data})$  before sampling from  $p(\boldsymbol{\omega}|\boldsymbol{\Sigma}, \text{data})$ . When taking logarithms from the both sides of the inference rule this will lead to equation:

$$U(\boldsymbol{\omega}) = \frac{1}{2} \sum_i \|\mathbf{y}_i - \mathbf{f}(\mathbf{x}_i, \boldsymbol{\omega})\|^2 / \sigma^2 + \log(Z(\boldsymbol{\omega})) + \frac{1}{2} \|\boldsymbol{\omega}\|^2 - \log(p_{\text{data}}(\mathbf{x})) + C(\text{data})$$

The most troublesome term in the equation is the  $Z(\boldsymbol{\omega})$  partition function that integrates *over data*, is function of  $\boldsymbol{\omega}$  and cannot be typically ignored (but is ignored in some contexts/books):

$$Z(\boldsymbol{\omega}) = \int e^{-\frac{1}{2} \sum_i \|\mathbf{y}_i - \mathbf{f}(\mathbf{x}_i, \boldsymbol{\omega})\|^2 / \sigma^2} p_{\text{data}}(\mathbf{x}_1 \dots \mathbf{x}_N) d\mathbf{x}_1 \dots d\mathbf{x}_N d\mathbf{y}_1 \dots d\mathbf{y}_N$$

Luckily we don't have to calculate  $Z(\boldsymbol{\omega})$  directly because in HMC we don't need value of  $U(\boldsymbol{\omega})$  but its gradient and difference (or in general MCMC we just need the difference):

$$\nabla_{\omega} U(\omega) = \nabla_{\omega} e(\omega) + \frac{1}{Z(\omega)} \nabla_{\omega} Z(\omega) + \omega$$

$$U(\omega_{n+1}) - U(\omega_n) = e(\omega_{n+1}) - e(\omega_n) + \log \left( \frac{Z(\omega_{n+1})}{Z(\omega_n)} \right) + \frac{1}{2} (\|\omega_{n+1}\|^2 - \|\omega_n\|^2)$$

$$e(\omega) = \frac{1}{2} \sum_i \|\mathbf{y}_i - f(\mathbf{x}_i, \omega)\|^2 / \sigma^2$$

It is rather straightforward to simplify the gradient of the partition function further:

$$\frac{1}{Z(\omega)} \nabla_{\omega} Z(\omega) = - \int \nabla_{\omega} e(\omega) p(\text{data}|\omega) d \text{data}$$

This last equation is interesting one because it defines sum of expected errors when  $(\mathbf{x}, \mathbf{y})$  is distributed according to our neural network model  $\mathbf{y} = f(\mathbf{x}, \omega)$ . If we make assumption that (in our model) each pair  $(\mathbf{x}, \mathbf{y})$  is independently distributed and have the same distribution we can write

$$\int \nabla_{\omega} e(\omega) p(\text{data}|\omega) d \text{data} = N \int \frac{1}{2} \nabla_{\omega} \|\mathbf{y} - f(\mathbf{x}, \omega)\|^2 / \sigma^2 p_{\text{model}}(\mathbf{x}, \mathbf{y}|\omega) d \mathbf{x} d \mathbf{y}$$

By inspecting the calculation methods we can see that gradient descent matches distribution of model  $p_{\text{model}}(\mathbf{y}|\mathbf{x}, \omega)$  as close as possible to distribution of data  $p_{\text{data}}(\mathbf{y}|\mathbf{x})$ .

### Learning variance parameter

In practice, variance term must be learnt in order to match model into data as well as possible. We can extend our model by writing

$$p(\omega, \sigma^2 | \mathbf{x}, \mathbf{y}) = p(\omega | \sigma^2, \mathbf{x}, \mathbf{y}) p(\sigma^2 | \mathbf{x}, \mathbf{y})$$

Weight parameters  $\omega$  can be sampled from the first distribution and we can rewrite the equation another way around

$$p(\sigma^2, \omega | \mathbf{x}, \mathbf{y}) = p(\sigma^2 | \omega, \mathbf{x}, \mathbf{y}) p(\omega | \mathbf{x}, \mathbf{y})$$

And use again the first probability term to sample/estimate  $\sigma^2$ . We can then alternate these sampling steps to generate samples from joint probability. In practice we need to estimate

$p(\sigma^2 | \omega, \mathbf{x}, \mathbf{y})$  which can be done simply by calculating average squared error without sampling (kind of maximum likelihood estimate).

### Analysis of $U(\omega)$ and effect of $Z(\omega)$

At optimum points we have  $\nabla_{\omega} U(\omega) = 0$ . Therefore we have an equation:

$$\int \nabla_{\omega} e(\omega) [p_{\text{data}}(\mathbf{y}|\mathbf{x}) - p_{\text{model}}(\mathbf{y}|\mathbf{x}, \omega)] p_{\text{data}}(\mathbf{x}) d \mathbf{x} d \mathbf{y} \approx 0$$

So in optimum points distributions of data and model are matched, or if that is not possible,  $U(\omega)$  increases or decreases endlessly when generating sampling points using sampler. In practice  $U(\omega)$  tends to increase to lower probability meaning that  $p_{\text{data}}(\mathbf{y}|\mathbf{x}) < p_{\text{model}}(\mathbf{y}|\mathbf{x}, \omega)$ . Additionally, we can see from the equations that  $Z(\omega)$  terms cause model to maximize generic model error while data terms minimize error of fitting function to data. Integrating gradient equation we can get formula for  $U(\omega)$  without  $Z$  function.

$$U(\omega) = N (E_{\text{data}}[e(\omega)] - E_{\text{model}}[e(\omega)]) + \frac{1}{2} \|\omega\|^2$$

This equation seem to have a problem that in practice the  $E_{\text{model}}[e(\omega)]$  tends to dominate in our approximated models meaning that  $N E_{\text{data}}[e(\omega)] + \frac{1}{2} \|\omega\|^2$  can increase too much. In a simplified analysis there is too much variance in our model but in bayesian probabilistic analysis data and model has equal amount of variance (it might be possible to get beyond bayesian probability analysis somehow by altering  $E_{\text{model}}[e(\omega)]$  this might include negative probabilities, complex number valued probabilities or other exotic calculation methods - altering  $E_{\text{model}} / Z$  means that we don't anymore force our probability mass to have fixed value).

We diverge when data has smaller variance (error) than modelling variance. Therefore we must match modelling variance to data variance but our modelling errors tends to increase variance/error of model. When our model too simple variance can grow too much.

Assuming we have initial solution  $\omega_0$  which is known to be a good function/model parameters for the data, we can write altered form (our model now integrates over  $\omega$  to give a single good model  $\omega_0$ ):

$$U^*(\omega) = N (E_{\text{data}}[e(\omega)] - E_{\text{model}, \omega_0}[e(\omega)])$$

And the negative term becomes constant because we force our model to has  $\omega_0$  parameters while data term is really sum of squared error values (only approximating  $E_{\text{data}}$ ) and not given parameter  $\omega_0$ . This mean we will ignore  $Z(\omega)$  term which maximizes function complexity (variance/error). After this we will only look for solutions  $\omega$  which allows small probabilistic variations in errors. This ignorance of  $Z(\omega)$  will also allow for probability mass of the function to be infinite?.

### Solving inverse values of function (not needed anymore)

If input space  $\mathbf{x}$  is small we can do parallel search (fast) by picking  $\mathbf{x}$  points randomly (or more preferably from training data) and then use gradient descent to minimize error function  $e(\mathbf{x}) = \frac{1}{2}(\mathbf{y} - \mathbf{f}(\mathbf{x}))^2$  and keep the solutions which are “close enough zero” (or at least  $N$  smallest error cases). This requires calculating  $\frac{d\mathbf{f}}{d\mathbf{x}}$  which we can compute using chain rule. For example, assuming two layer neural network layers we have:

$$\begin{aligned} \mathbf{f}(\mathbf{x}) &= \mathbf{f}_2(\mathbf{A}_2 \mathbf{f}_1(\mathbf{A}_1 \mathbf{x} + \mathbf{a}_1) + \mathbf{a}_2) \\ \frac{d\mathbf{f}}{d\mathbf{x}} &= \mathbf{f}_2'(\mathbf{A}_2 \mathbf{f}_1(\mathbf{A}_1 \mathbf{x} + \mathbf{a}_1) + \mathbf{a}_2) \mathbf{A}_2 \mathbf{f}_1'(\mathbf{A}_1 \mathbf{x} + \mathbf{a}_1) \mathbf{A}_1 \end{aligned}$$

And it is straightforward to extend this to multiple layers. However, the theoretical problem is that there are probably very large or infinite many possible input values  $\mathbf{x}$  for many cases of  $\mathbf{y}$  leading to very difficult to estimate  $p(\mathbf{x}|\mathbf{y})$ . We can limit the number of solutions if we restrict ourselves to the convex space (+ little extra around corners to do extrapolation..) and set up by training points and calculate our inverse function and probability function only in that space.

In practice, the inverse solutions can be extremely large (outside the convex space setup by data). If we assume our data has zero mean and unit variance, then we can limit the solutions space by regularizing/using prior for  $\mathbf{x}$  by forcing it to be taken from gaussian ball and define error function to be:

$$\begin{aligned} E(\mathbf{x}) &= \frac{1}{2}(\mathbf{y} - \mathbf{f}(\mathbf{x}))^2 + \frac{1}{2}\|\mathbf{x}\|^2/\alpha^2 \\ \nabla E(\mathbf{x}) &= (\mathbf{y} - \mathbf{f}(\mathbf{x}))^T \nabla \mathbf{f}(\mathbf{x}) + \mathbf{x}/\alpha^2 \end{aligned}$$

Where alpha is accuracy/inverse variance parameter which we can set  $\alpha = 3$ .

After solving inverse values we can then assign probability to each  $\mathbf{x}$  according to probability density function, normalize to unity and select one according to probability distribution.

### Solving difference equation

After we have computed gradient of the network we must solve ratio  $\log \left( \frac{Z(\omega_{n+1})}{Z(\omega_n)} \right)$  in order to make accept/reject decisions. We can estimate partition function ratios by using (reference):

$$\frac{Z(\omega_{n+1})}{Z(\omega_n)} = \int \frac{1}{Z(\omega_n)} p^*(\mathbf{x}, \mathbf{y}|\omega_{n+1}) d\mathbf{x} d\mathbf{y} = \int \frac{p^*(\mathbf{x}, \mathbf{y}|\omega_{n+1})}{p^*(\mathbf{x}, \mathbf{y}|\omega_n)} p(\mathbf{x}, \mathbf{y}|\omega_n) d\mathbf{x} d\mathbf{y}$$

So we must calculate estimate  $E_{\mathbf{x}, \mathbf{y}|\omega_n} \left[ \frac{p^*(\mathbf{x}, \mathbf{y}|\omega_{n+1})}{p^*(\mathbf{x}, \mathbf{y}|\omega_n)} \right]$  by calculating ratio of biased probability distributions of data. In practice we have multiple data points in our data likelihood. This leads to equation:

$$E_{\{\mathbf{x}_i, \mathbf{y}_i\}|\omega_n} \left[ \prod_i \frac{p^*(\mathbf{x}_i, \mathbf{y}_i|\omega_{n+1})}{p^*(\mathbf{x}_i, \mathbf{y}_i|\omega_n)} \right]$$

Now our terms are identically distributed and independent from each other  $E[AB] = E[A] E[B]$ , this means

$$E_{\{x_i, y_i\}|\omega_n} \left[ \prod_i \frac{p^*(x_i, y_i|\omega_{n+1})}{p^*(x_i, y_i|\omega_n)} \right] = \prod_i E_{x, y|\omega_n} \left[ \frac{p^*(x, y|\omega_{n+1})}{p^*(x, y|\omega_n)} \right] = E_{x, y|\omega_n} \left[ \frac{p^*(x, y|\omega_{n+1})}{p^*(x, y|\omega_n)} \right]^N$$

$$\frac{p^*(x, y|\omega_{n+1})}{p^*(x, y|\omega_n)} = e^{\frac{1}{2}\|y - f(x, \omega_n)\|^2/\sigma^2 - \frac{1}{2}\|y - f(x, \omega_{n+1})\|^2/\sigma^2}$$

And we sample from  $p_{\text{model}}(\mathbf{y}, \mathbf{x}|\omega_n) = p_{\text{model}}(\mathbf{y}|\mathbf{x}, \omega_n) p_{\text{data}}(\mathbf{x})$  in feedforward neural network because we are only interested in conditional distribution  $p(\mathbf{y}|\mathbf{x})$  which is in our model is normal distribution.

### Additional implementation notes

In practice, we need to estimate negative phase gradient and zratio by sampling from distribution  $p_{\text{model}}(\mathbf{x}, \mathbf{y}|\omega)$ , which we can only crudely approximate. We are estimating mean value of z-ratio and mean value of gradients (vector) and we much estimate error levels in both of them in order to gather enough samples.

TODO

### Divergence of HMC sampler

In practice HMC sampler diverge to larger and larger error levels after successfully visiting small levels of error. The reason for this is currently unclear but only way for sampler to keep accepting the new samples is that Z-ratio becomes very small to always dominate the results which then also means that gradient and is based on Z-value and other terms are ignored.

More detailed analysis seems to imply that the problem is variance (or covariance) term in bayesian neural network. If variance is set to unity, it sets the accuracy (level of detail) of our predicted outcome. Too high variance when compared to data variance means that model “underfits” to the data. Moreover, gradient descent rule now keeps increasing  $U(\omega)$  leading to flat posterior distribution  $p(\omega)$  [because any parameters can fit to data with too much variance in model] and high model error levels. This is because tails of the model distribution (extreme values) are more likely in  $p_{\text{model}}(\mathbf{x}, \mathbf{y}|\omega)$  which means that following the negative gradient causes  $U(\omega)$  to increase

$$\nabla_{\omega} U(\omega) = \int \nabla_{\omega} \frac{1}{2} \|\mathbf{y} - \mathbf{f}(\mathbf{y}|\omega)\|^2 [p_{\text{data}}(\mathbf{y}|\mathbf{x}) p_{\text{data}}(\mathbf{x}) - p_{\text{model}}(\mathbf{y}|\mathbf{x}, \omega) p_{\text{data}}(\mathbf{x})] + \omega$$

Additionally, we can also see that gradient descent rule will converge when model distribution has converged to data distribution. And because in practice we are only interested in  $p(\mathbf{y}|\mathbf{x})$  distribution we don’t have to calculate inverse of  $\mathbf{f}(\mathbf{x})$  in order to try to (hopelessly) match the whole distribution but sample from  $p_{\text{data}}(\mathbf{x})$  and then from  $p_{\text{model}}(\mathbf{y}|\mathbf{x}, \omega)$  in order to generate negative particles. So in order for unit variance method to work, our data must have **larger** variance than unit variance and “the level of detail” in our data must have unit st.dev.

In order to fix this, we must update variance parameter  $\sigma^2$  at every sampling step using Gibbs sampling. First we have  $\omega$  and we sample from  $p(\sigma^2|\omega, \text{data}) \propto p(\text{data}|\omega, \sigma^2) p(\sigma^2)$ . After we have sampled from  $\sigma^2$  we use our previous results to sample using fixed  $\sigma^2$  from  $p(\omega|\text{data}, \sigma^2)$  and store our sample  $(\omega_i, \sigma_i^2)$ .

We can directly sample from distribution

$$p(\sigma^2|\omega, \text{data}) \sim X^2(n d - 1, s^2), s^2 = \sum_i \|\mathbf{y}_i - \mathbf{f}(\mathbf{x}_i, \omega)\|^2 / (n d)$$

We can sample from this by generating  $n d - 1$ , normally distributed zero mean unit variance variables and calculating their squared sum.  $Z = \sum_i X_i^2$ , after this we scale by  $s^2$ .

### Sampling covariance matrix $\Sigma$

Similarly, we can do calculate posterior distribution  $p(\Sigma|\text{data}, \omega)$  for covariance matrix (InvWishart distribution) and sample from it by using standard multivariate normal distribution inference.

The prior for the normally distributed zero mean data  $N(\mathbf{y}_i - \mathbf{f}(\mathbf{x}_i, \mathbf{w}) | \boldsymbol{\mu}, \boldsymbol{\Sigma})$  is zero mean normal distribution  $\boldsymbol{\mu} \sim N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}/\kappa_0)$  and  $\boldsymbol{\Sigma}^{-1} \sim \text{Wishart}_{v_0}(\boldsymbol{\Lambda}_0^{-1})$ . For mean we decide  $\boldsymbol{\mu}_0 = 0$  and  $\kappa_0 = \infty$  and we get a posterior distribution for  $\boldsymbol{\Sigma} | \text{data}, \boldsymbol{\omega} \sim \text{InvWishart}_{v_n}(\boldsymbol{\Lambda}_n^{-1})$ .

$$v_n = v_0 + n$$

$$\boldsymbol{\Lambda}_n = \boldsymbol{\Lambda}_0^{-1} + \sum_i \mathbf{z}_i \mathbf{z}_i^T, \quad \mathbf{z}_i = \mathbf{y}_i - \mathbf{f}(\mathbf{x}_i, \boldsymbol{\omega})$$

And we can get non-informative prior by taking parameter values  $\boldsymbol{\Lambda}_0 = 0$  and  $v_0 = -1$ . Now we need to sample from InvWishart distribution. This means  $\boldsymbol{\Sigma}^{-1} \sim \text{Wishart}_{v_n}(\boldsymbol{\Lambda}_n^{-1})$  and we sample from Wishart distribution. This can be done by sampling  $\boldsymbol{\alpha}_i \sim N(\mathbf{0}, \boldsymbol{\Lambda}_n^{-1})$  and calculating  $\sum_i^{v_n} \boldsymbol{\alpha}_i \boldsymbol{\alpha}_i^T$ . (*Bayesian Data Analysis. Andrew Gelman. Appendix A.*)

When number of observations is less than  $\dim(\mathbf{z})$ :  $n \leq \dim(\mathbf{z})$  we set identity matrix prior  $v_0 = \dim(\mathbf{z}) - n$ ,  $\boldsymbol{\Lambda}_0 = \mathbf{I}$  and otherwise we use non-informative prior  $\boldsymbol{\Lambda}_0 = 0$  and  $v_0 = -1$ .

In practice we choose  $\sigma^2 = \min(\text{eig}(\boldsymbol{\Sigma}))$ .

### Calculating backpropagation gradient of error term

Calculating gradient of error term  $\nabla_{\boldsymbol{\omega}} e(\boldsymbol{\omega})$  is complicated because the basic backpropagation algorithm calculates gradient minimizing non whitened error term. The error function:

$$e(\boldsymbol{\omega}) = \frac{1}{2} \sum_i (\mathbf{y}_i - \mathbf{f}(\mathbf{x}_i, \mathbf{w}))^T \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i - \mathbf{f}(\mathbf{x}_i, \mathbf{w}))$$

can be rewritten using  $\mathbf{W} = \mathbf{D}^{1/2} \mathbf{X}^T$  where  $\boldsymbol{\Sigma}^{-1} = \mathbf{X} \mathbf{D} \mathbf{X}^T$ ,

$$e(\boldsymbol{\omega}) = \frac{1}{2} \sum_i (\mathbf{W} \mathbf{y}_i - \mathbf{W} \mathbf{f}(\mathbf{x}_i, \mathbf{w}))^T (\mathbf{W} \mathbf{y}_i - \mathbf{W} \mathbf{f}(\mathbf{x}_i, \mathbf{w}))$$

$$\tilde{e}(\boldsymbol{\omega}) = \frac{1}{2} \sum_i (\tilde{\mathbf{y}}_i - \tilde{\mathbf{f}}(\mathbf{x}_i, \mathbf{w}))^T (\tilde{\mathbf{y}}_i - \tilde{\mathbf{f}}(\mathbf{x}_i, \mathbf{w}))$$

Which have preferred squared error form. However, our neural network now has one extra layer  $\mathbf{W}$  which scales and rotates actual output of the function. Therefore we must look into backpropagation algorithm and make slight modifications to calculate gradient when the output layer is changed with additional extra linear layer.

Update rule for local gradient and the update rule for the weights is:

$$\boldsymbol{\delta}_{l-1} = \text{diag}(\boldsymbol{\varphi}'(\mathbf{v}_{l-1})) \mathbf{W}_l^T \boldsymbol{\delta}_l, \quad \mathbf{W}_{ji}^{l-1} = \boldsymbol{\delta} \mathbf{y}^T$$