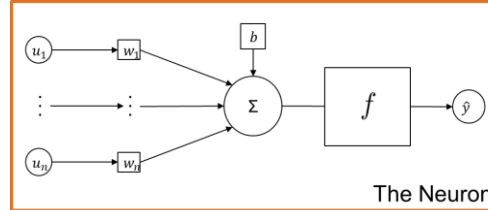


Introduction to Bayesian Neural Networks

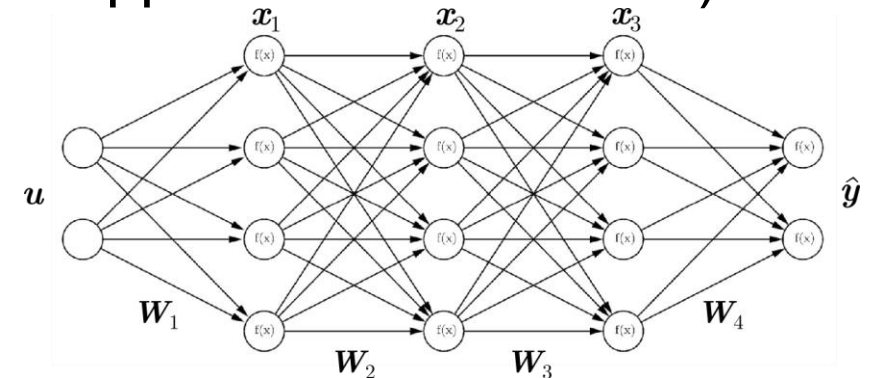
Anh Khoa Doan – n.a.k.doan@tudelft.nl

Recap from last week

- The neuron
 - “Base” function for regression problem
 - Single neuron with sigmoid: logistic regression (binary classification)



- Neural networks
 - Complex arrangement of neurons
 - Can “approximate” any function (universal approximator theorem)
 - Trained from data by MSE-minimization
 - With (form of) stochastic gradient descent
 - Efficient with “automatic differentiation”
 - And backpropagation



Recap from last week and why UQ with NN

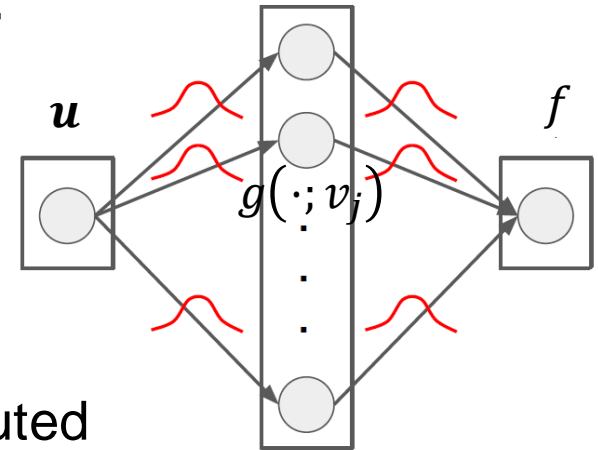
- Neural Networks
 - Can be used for regression (point-based estimate)
 - Can be used for probability distribution estimation over categories (classification with negative log-probability loss (cross entropy loss))
- What is the connection with UQ?
 - Smooth out the predictions by averaging over plausible explanations
 - Get confidence interval
 - Make robust decision

NN converges towards Gaussian process

- Let's start from a single hidden layer NN with N nodes:

$$f(\mathbf{u}) = b + \sum_j^N w_j g(\mathbf{u}; v_j)$$

- Now, let's assign some distributions. Assume:
 - $b \sim \mathcal{N}(0, \sigma_b^2)$ and $w_j \sim \mathcal{N}(0, \sigma_w^2)$
 - v_j (neuron weights) are independently and identically distributed
 - σ_w^2 scales as w^2/N then



$$\begin{aligned}\mathbb{E}(f(\mathbf{u})) &= 0 \\ \mathbb{E}(f(\mathbf{u})f(\mathbf{u}')) &= \sigma_b^2 + \sum_j^N \sigma_w^2 \mathbb{E}_x \left(g(\mathbf{u}; x_j) g(\mathbf{u}'; x_j) \right) \\ &= \sigma_b^2 + w^2 \mathbb{E}_x \left(g(\mathbf{u}; x_j) g(\mathbf{u}'; x_j) \right)\end{aligned}$$

- By theorem central limit, f converges to a Gaussian process as $N \rightarrow \infty$

Structure of the lecture

- **Reminder on Bayesian Inference**
- Bayesian Neural Network

Small reminder on Bayesian Inference

- Bayesian inference: predictions by averaging over all likely explanations under the posterior distribution
- Posterior estimated with Bayes' rule

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

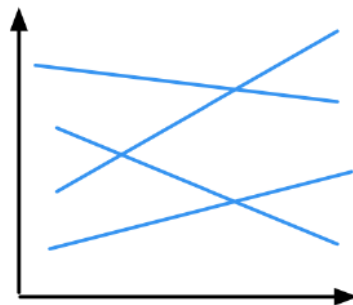
\mathcal{D} : distribution of observed values

- Prediction using the posterior predictive distribution:

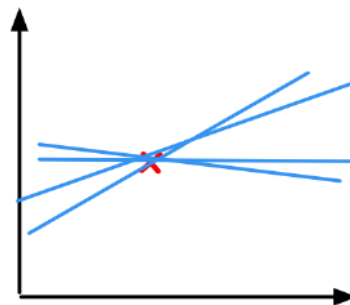
$$p(y|\mathbf{u}, \mathcal{D}) = \int p(\mathbf{w}|\mathcal{D})p(y|\mathbf{u}, \mathbf{w})d\mathbf{w}$$

Bayesian linear regression

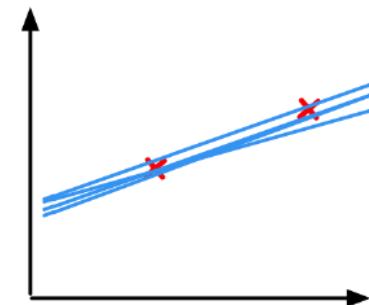
- Bayesian linear regression: considers various “plausible” explanation for data generation
- Prediction obtained using all possible regression weights, weighted by their posterior probability



no observations



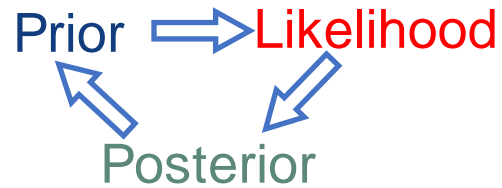
one observation



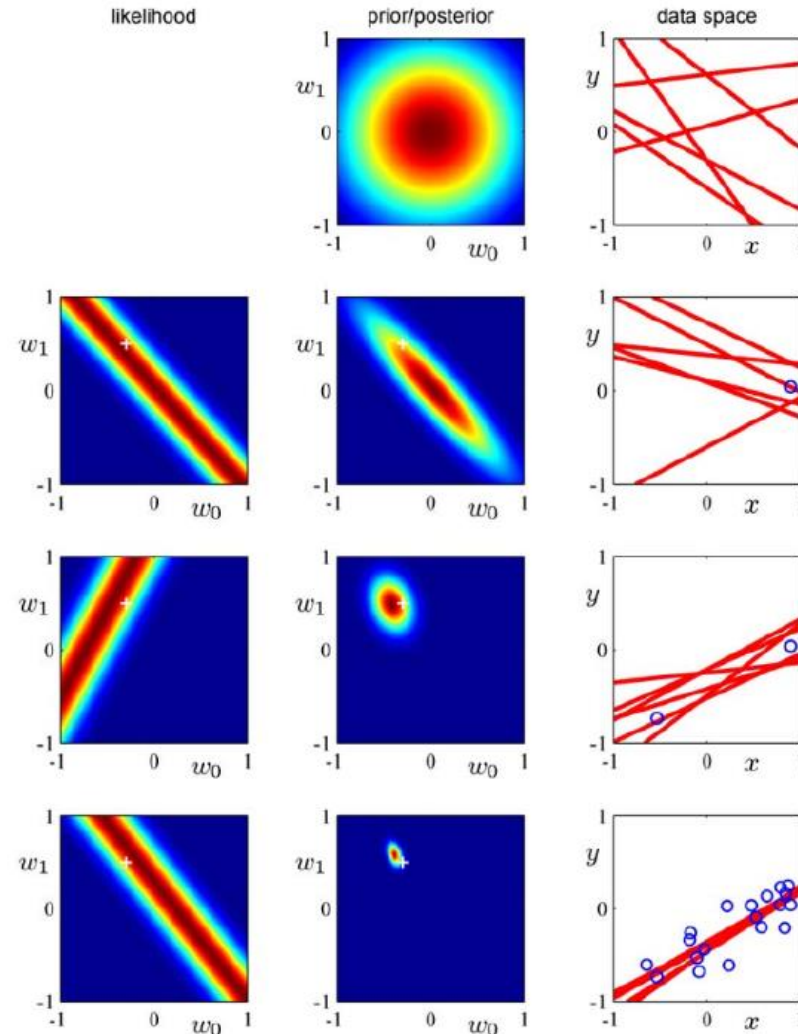
two observations

- Prior distribution: $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{S})$
- Likelihood: $p(y|\mathbf{u}, \mathbf{w}) \sim \mathcal{N}(\mathbf{w}^T \psi(\mathbf{u}), \sigma^2)$
- Assumed fixed/known \mathbf{S} and σ^2 is a big assumption

Bayesian linear regression



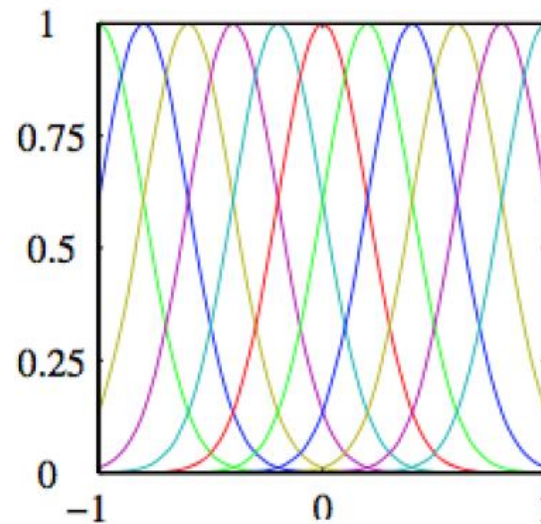
1. Prior distribution
2. 3. 4. Model distribution is adjusted to data



Bayesian Regression

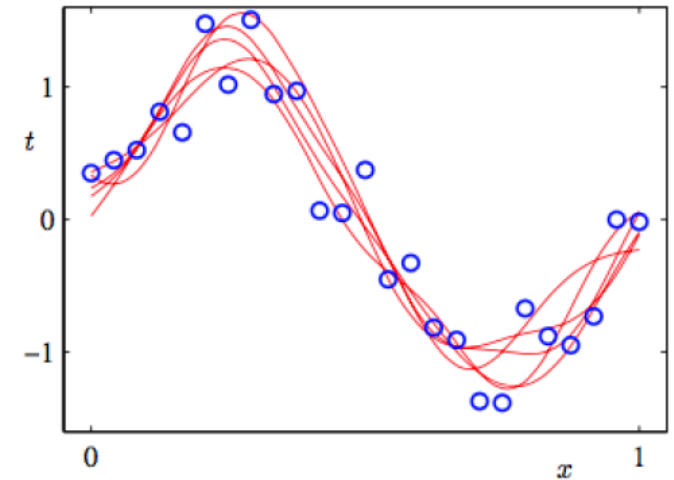
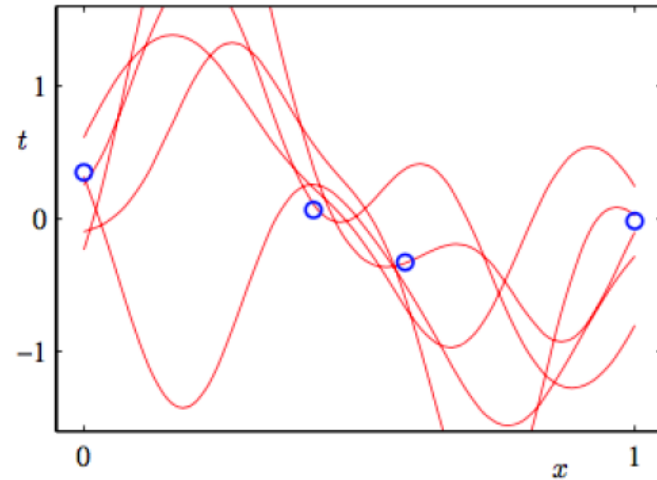
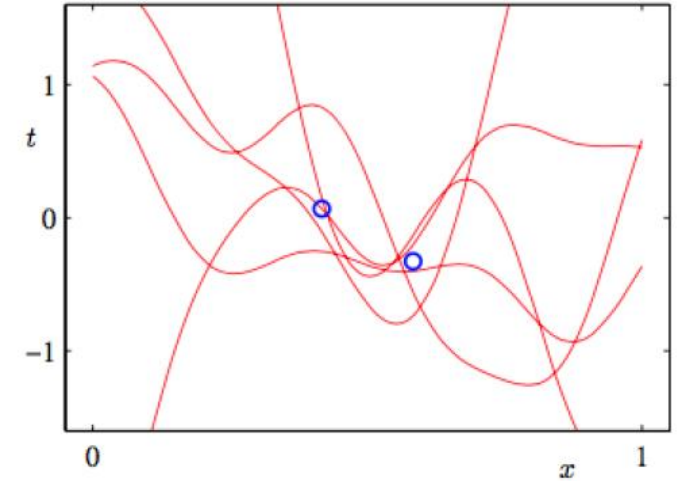
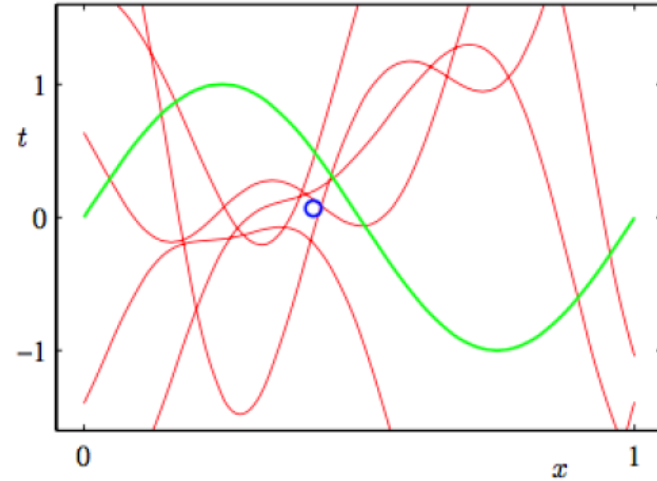
- Example with radial basis function features:

$$\phi_j(u) = \exp\left(-\frac{(x - \mu_j)^2}{2s^2}\right)$$



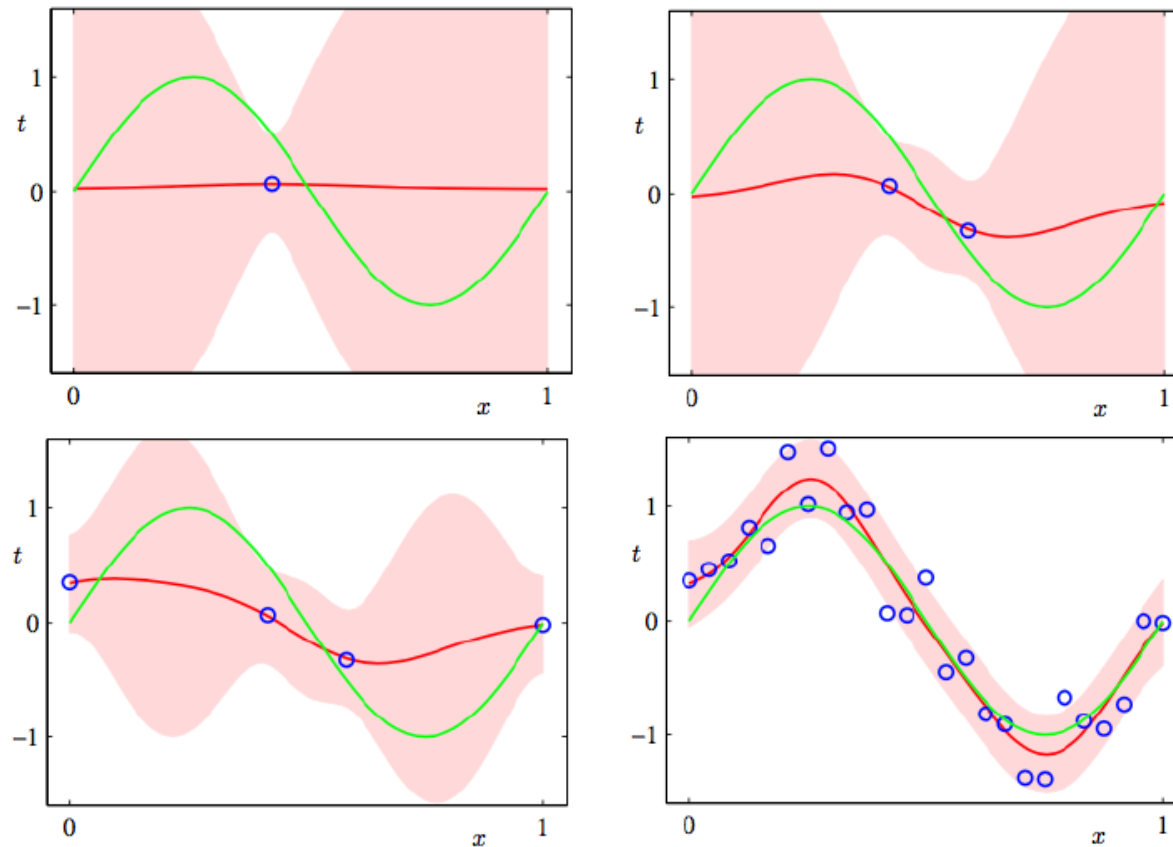
Bayesian Regression

Functions sampled from the posterior:



Bayesian Regression

- Visualization of confidence intervals based on posterior predictive mean and variance



Structure of the lecture

- Reminder on Bayesian Inference
- **Bayesian Neural Network**

Bayesian neural network

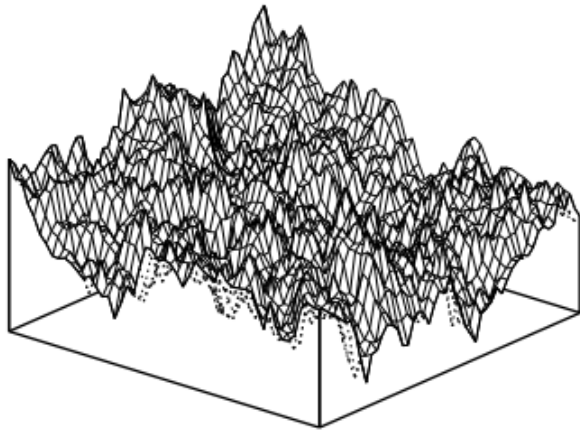
- So far: fixed basis functions...
- Can we combine the advantage of neural networks with Bayesian models?
 - Place a prior on the weights of the network, e.g. $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \mathbf{0}, \eta \mathbf{I})$
 - Define an observation model, e.g. $p(y|\mathbf{u}, \mathbf{w}) = \mathcal{N}(y; f_{\mathbf{w}}(\mathbf{u}), \sigma^2)$
 - Apply Bayes' Rule:

$$p(\mathbf{w}|\mathcal{D}) \propto p(\mathbf{w}) \prod_{i=1}^N p(y^{(i)}|\mathbf{u}^{(i)}, \mathbf{w})$$

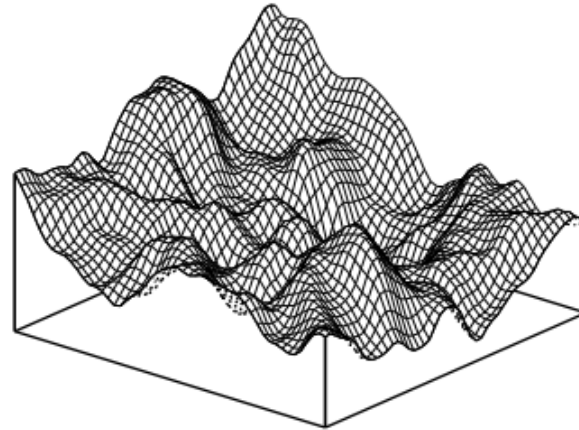
\mathcal{D} : dataset distribution

Samples from the prior

- How to interpret prior $p(\mathbf{w})$ with network kernel?
- Prior samples for a BNN with one hidden layer and 10,000 units



hard threshold activations



tanh activations

BNN: what to maximize

- More generally, what we consider is:
 - Prior weight distribution: $p(\mathbf{w})$
 - Given a training dataset sampled distribution \mathcal{D}
 - Likelihood function: $p(\mathcal{D}|\mathbf{w}) = p(\mathbf{w}) \prod_{i=1}^N p(y^{(i)}|\mathbf{u}^{(i)}, \mathbf{w})$
- We could try to:
 - maximize $p(\mathcal{D}|\mathbf{w}) \rightarrow$ Maximum likelihood estimation (MLE)
 - Can be biased by the data \rightarrow to be avoided when few data
- Alternate: Bayes rule gives:

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

- We could try to maximize:
 - $p(\mathcal{D}|\mathbf{w})p(\mathbf{w})$: Maximum a posteriori (MAP) estimation

Posterior Inference

- Possible use of posterior distribution: sample set of values $\mathbf{w}_1, \dots, \mathbf{w}_K$ from the posterior distribution $p(\mathbf{w}|\mathcal{D})$ and average their predictive distributions:

$$p(y|\mathbf{u}, \mathcal{D}) \approx \frac{1}{K} \sum_{k=1}^K p(y|\mathbf{u}, \mathbf{w}_k)$$

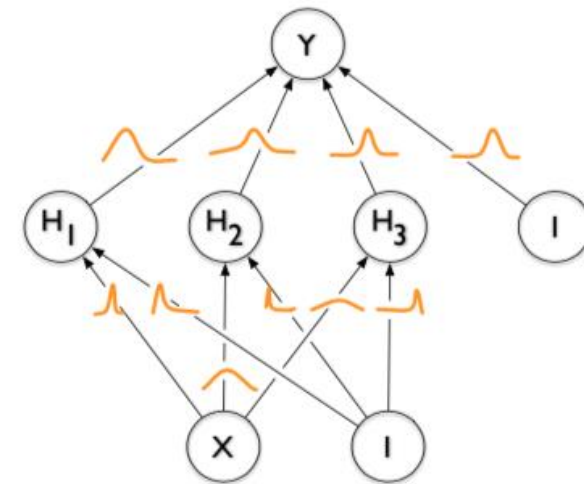
- Sample from posterior: can be obtained approximately with Markov Chain Monte Carlo
 - But can be expensive with very large dataset...
- How can we obtain $p(\mathbf{w}|\mathcal{D})$?
- Instead: Variational inference to estimate $p(\mathbf{w}|\mathcal{D})$

Posterior Inference: Variational Bayes

- Idea: Approximate complex posterior distribution with simpler (analytical) variational approximation q with parameters θ
 - Eg: Assume Gaussian posterior with diagonal covariance

$$q(\mathbf{w}; \theta) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$
$$= \prod_j^D \mathcal{N}(w_j; \mu_j, \sigma_j)$$

“each weight of the network has its own mean and variance”



Posterior Inference: Variational Bayes

- From $q(\mathbf{w}; \boldsymbol{\theta})$, we now need to estimate the “best” $\boldsymbol{\theta}$ that gives us the best approximation of the posterior $p(\mathbf{w}|\mathcal{D})$
- We can use the Kullback-Leibler divergence (“measure of distance between two distributions”), D_{KL} , between the two:

$$\begin{aligned} D_{KL}(q(\mathbf{w}; \boldsymbol{\theta}) || p(\mathbf{w}|\mathcal{D})) &\equiv \int q(\mathbf{w}; \boldsymbol{\theta}) \log \frac{q(\mathbf{w}; \boldsymbol{\theta})}{p(\mathbf{w}|\mathcal{D})} d\mathbf{w} \\ &= \mathbb{E}_{q(\mathbf{w}; \boldsymbol{\theta})} \log \frac{q(\mathbf{w}; \boldsymbol{\theta})}{p(\mathbf{w}|\mathcal{D})} \end{aligned}$$

Posterior Inference: Variational Bayes

$$D_{KL}(q(\mathbf{w}; \boldsymbol{\theta}) || p(\mathbf{w}|\mathcal{D})) = \mathbb{E}_{q(\mathbf{w}; \boldsymbol{\theta})} \log \frac{q(\mathbf{w}; \boldsymbol{\theta})}{p(\mathbf{w}|\mathcal{D})}$$

- With Bayes' rule for $p(\mathbf{w}|\mathcal{D})$

$$D_{KL}(q(\mathbf{w}; \boldsymbol{\theta}) || p(\mathbf{w}|\mathcal{D})) = \mathbb{E}_{q(\mathbf{w}; \boldsymbol{\theta})} \log \frac{q(\mathbf{w}; \boldsymbol{\theta})}{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})} p(\mathcal{D})$$

$$= \mathbb{E}_{q(\mathbf{w}; \boldsymbol{\theta})} \log \frac{q(\mathbf{w}; \boldsymbol{\theta})}{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})} p(\mathcal{D})$$

$$= \mathbb{E}_{q(\mathbf{w}; \boldsymbol{\theta})} [\log q(\mathbf{w}; \boldsymbol{\theta}) - \log p(\mathcal{D}|\mathbf{w}) - \log p(\mathbf{w}) + \log p(\mathcal{D})]$$

$$D_{KL}(q(\mathbf{w}; \boldsymbol{\theta}) || p(\mathbf{w}|\mathcal{D})) = D_{KL}(q(\mathbf{w}; \boldsymbol{\theta}) || p(\mathbf{w})) - \mathbb{E}_{q(\mathbf{w}; \boldsymbol{\theta})} [\log p(\mathcal{D}|\mathbf{w})] + \log p(\mathcal{D})$$

Posterior Inference: Variational Bayes

$$D_{KL}(q(\mathbf{w}; \boldsymbol{\theta}) || p(\mathbf{w}|\mathcal{D})) = D_{KL}(q(\mathbf{w}; \boldsymbol{\theta}) || p(\mathbf{w})) - \mathbb{E}_{q(\mathbf{w}; \boldsymbol{\theta})}[\log p(\mathcal{D}|\mathbf{w})] + \log p(\mathcal{D})$$

- $\mathcal{F}(\mathcal{D}, \boldsymbol{\theta}) \equiv D_{KL}(q(\mathbf{w}; \boldsymbol{\theta}) || p(\mathbf{w})) - \mathbb{E}_{q(\mathbf{w}; \boldsymbol{\theta})}[\log p(\mathcal{D}|\mathbf{w})]$
 - Called “variational free energy”
- We can minimize \mathcal{F} with respect to $\boldsymbol{\theta}$ and that will minimize $D_{KL}(q(\mathbf{w}; \boldsymbol{\theta}) || p(\mathbf{w}|\mathcal{D}))$
 - And $q(\mathbf{w}; \boldsymbol{\theta})$ will approximate the posterior
- Alternate name: $\mathcal{L}(\mathcal{D}, \boldsymbol{\theta}) \equiv -\mathcal{F}(\mathcal{D}, \boldsymbol{\theta})$ “evidence lower bound”
 - As $\mathcal{L} \leq \log p(\mathcal{D})$ (as $D_{KL} \geq 0$)

Posterior Inference: Variational Bayes

- Now we need to minimize:

$$\mathcal{F}(\mathcal{D}, \boldsymbol{\theta}) \equiv D_{KL}(q(\mathbf{w}; \boldsymbol{\theta}) || p(\mathbf{w})) - \mathbb{E}_{q(\mathbf{w}; \boldsymbol{\theta})} [\log p(\mathcal{D} | \mathbf{w})]$$

- $D_{KL}(q(\mathbf{w}; \boldsymbol{\theta}) || p(\mathbf{w}))$: “complexity cost”
 - $\mathbb{E}_{q(\mathbf{w}; \boldsymbol{\theta})} [\log p(\mathcal{D} | \mathbf{w})]$: “likelihood cost”
 - Some re-arrangement
- $$\mathcal{F}(\mathcal{D}, \boldsymbol{\theta}) = \mathbb{E}_{q(\mathbf{w}; \boldsymbol{\theta})} [\log q(\mathbf{w}; \boldsymbol{\theta})] - \mathbb{E}_{q(\mathbf{w}; \boldsymbol{\theta})} [\log p(\mathbf{w})] - \mathbb{E}_{q(\mathbf{w}; \boldsymbol{\theta})} [\log p(\mathcal{D} | \mathbf{w})]$$
- All terms are expectations with respect to $q(\mathbf{w}; \boldsymbol{\theta})$, so we can approximate them by drawing samples for $\mathbf{w}^{(i)}$ from $q(\mathbf{w}; \boldsymbol{\theta})$

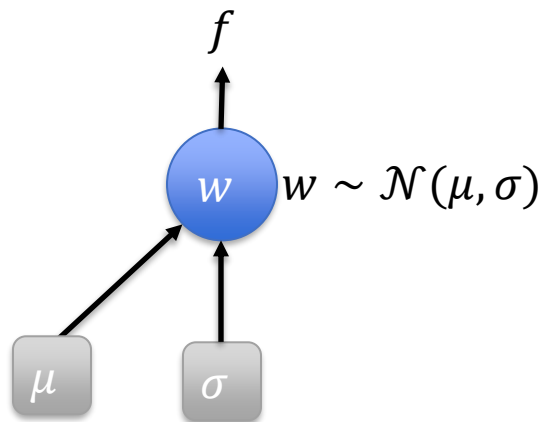
$$\mathcal{F}(\mathcal{D}, \boldsymbol{\theta}) \approx \frac{1}{N} \sum_i^N \log q(\mathbf{w}^{(i)}; \boldsymbol{\theta}) - \log p(\mathbf{w}^{(i)}) - \log p(\mathcal{D} | \mathbf{w}^{(i)})$$

→ We can minimize that with respect to $\boldsymbol{\theta}$


Note on backpropagation with random nodes

- How do we backpropagate with random variables involved?
- “Re-parametrization” trick

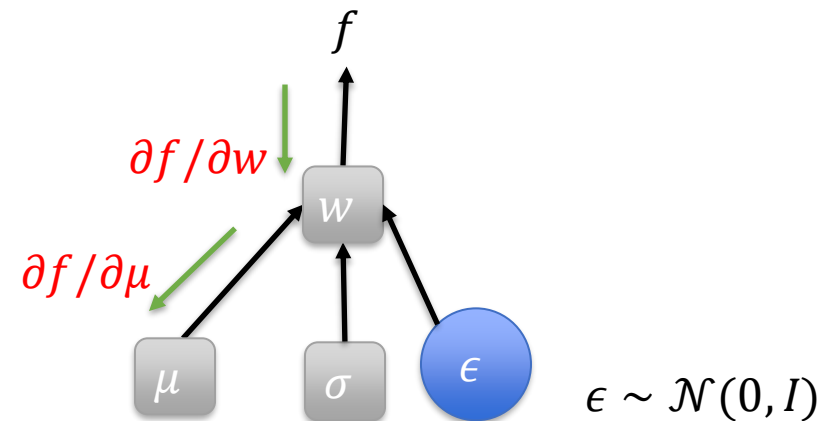
Original form:



 : deterministic node

 : random node

Reparametrized form:



1. Sample ϵ
2. Let $w = \mu + \sigma\epsilon$ (or $w = \mu + (1 + \log \sigma)\epsilon$)
3. Backpropagate

Training with Bayesian NN

- Train a Bayesian NN with:

$$w_j = \mu_j + \sigma_j \epsilon_j$$

With $\epsilon_j \sim \mathcal{N}(0,1)$

- ϵ_j sampled at the beginning of training, independent of μ_j, σ_j
 - \rightarrow Deterministic graph, backpropagation algorithm can be used
- If all $\sigma_j = 0$, then $\theta_j = \mu_j$, and ordinary backpropagation with deterministic neural network can be used

Structure of the lecture

- Reminder on Bayesian Inference
- Bayesian Neural Network

Bayesian Neural Network - Summary

- Approach that combines Bayesian principle with NN
- Can propagate uncertainty/account for uncertainty