CS6462
*Probabilistic and Explainable AI*

# Lesson 21
# *Bayesian Neural Networks*
# *
# Posterior Variational Inference*

by Emil Vassev

March 10, 2025

# BNN Probabilistic Models

*Bayesian Neural Network:*
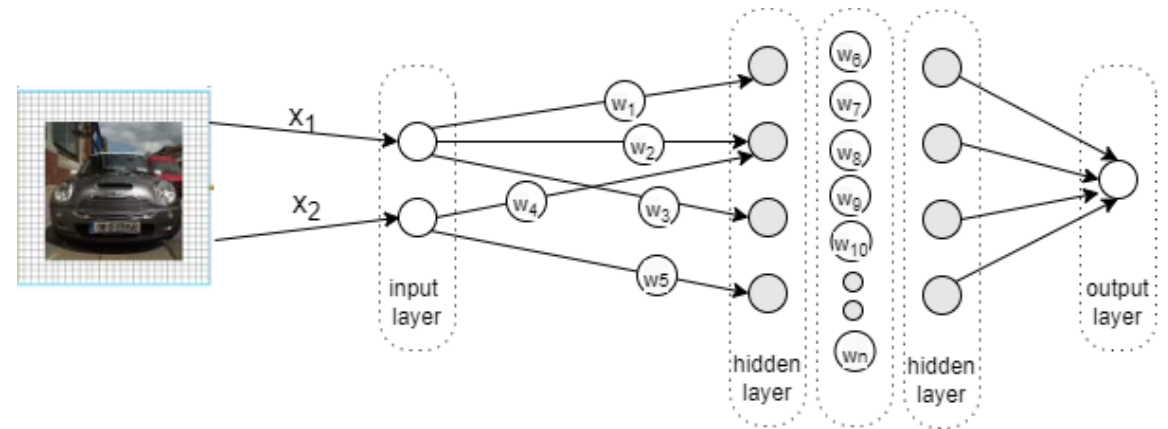
- based on Bayes' Theorem

- probabilistic model of BNN

  $p(y|x, w)$

  - inputs $x = \{x_1, x_2, ..., x_n\}$
  - weight factors $w = \{w_1, w_2, ..., w_n\}$
  - for classification:
    - $y$ - set of classes
    - $p(y|x, w)$ - Categorical distribution
  - for regression:
  - $y$ - continuous random variable
    - $p(y|x, w)$ - Gaussian distribution

$$Posterior = \frac{Likelyhood * Prior}{Evidence} \propto Likelyhood * Prior$$

$$p(w|D) = \frac{p(D|w) * p(w)}{p(D)} \propto p(D|w) * p(w)$$

# BNN Probabilistic Models (cont.)

*Bayesian Neural Network Training:*

BNN Probabilistic model: $p(y|x,w)$

- training dataset
$D = \{x^{(i)}, y^{(i)}\}$

- $x^{(i)}$ - the input vector of the i-th training example

- $y^{(i)}$ - class label for the i-th training example

*Maximum Likelihood Estimate (MLE)*:

- likelihood distribution:

   $p(D|w) = \prod_i p(y^{(i)}|x^{(i)}, w)$ - function of the weight factors $w$

   maximizing $p(D|w)$ - Maximum Likelihood Estimate (MLE) of $w$

- optimization objective – negative log likelihood:
   - Categorical distribution: Cross Entropy Error function
   - Gaussian distribution: proportional to the sum of Square Error function
   - MLE can lead to heavy overfitting

# BNN Probabilistic Models (cont.)

*Maximum a Posteriori (MAP) Estimate:*

BNN Probabilistic model: $p(y|x,w)$

- $p(w|D)$ - posterior distribution

- product $p(D|w) * p(w)$ - proportional ($\propto$) to $p(w|D)$

$$p(w|D) \propto p(D|w) * p(w)$$

- maximizing $p(D|w) * p(w)$ - Maximum a Posteriori (MAP) estimate of $w$

- computing MAP - can prevent overfitting

- optimization objective:
  - negative log likelihood
  - regularization term with log prior

*Posterior Predictive Distribution:*

$$p(y|x,D) = \int p(y|x,w) * p(w|D) * dw$$

- full posterior distribution over parameters - predictions with weight uncertainty into account

- parameters $w$ are marginalized

# Variational Inference

*Posterior with a variational distribution:*

$$p(y|x, D) = \int p(y|x, w) * p(w|D) * dw$$

- $p(w|D)$ - posterior distribution

- difficult analytical solution to $p(w|D)$ in BNN

- solution: approximate the true posterior with a variational distribution

- $q(w|\theta)$ - variational distribution;

- $\theta$ – set of parameters we want to estimate

- new posterior with variational distribution:
  - cost function $\boldsymbol{F}(\boldsymbol{D}, \boldsymbol{\theta})$
  - minimizes the *Kullback–Leibler divergence\** between $q(w|\theta)$ and $p(w|D)$

*Kullback–Leibler divergence on Wikipedia

# Variational Inference(cont)

*Cost function (variational free energy):*

- corresponding optimization objective or cost function (*variational free energy*):

$$F(D, \theta) = KL\big(q(w|\theta) \parallel p(w)\big) - \mathbb{E}_{q(w|\theta)}[\log(p(D|w))]$$

$$F(D, \theta) = complexity\ cost - likelyhood\ cost$$

$KL\big(q(w|\theta) \parallel p(w)\big)$ - *measures the statistical distance how $q(w|\theta)$ is different from $p(w)$*

$\mathbb{E}_{q(w|\theta)}[\log(p(D|w))]$ - likelihood cost: the expected value of $p(D|w)$ with respect to $q(w|\theta)$

$\mathbb{E}_{q(w|\theta)}$ - energy expectation function

# Variational Inference (cont.)

*Cost function:*

$$F(D, \theta) = KL\big(q(w|\theta) \,\|\, p(w)\big) - \mathbb{E}_{q(w|\theta)}[\log(p(D|w))]$$

- rearranging the complexity cost component:

$$\boldsymbol{F(D, \theta)} = \boldsymbol{KL}\big(\boldsymbol{q(w|\theta)} \,\|\, \boldsymbol{p(w)}\big) - \mathbb{E}_{\boldsymbol{q(w|\theta)}}[\boldsymbol{\log(p(D|w))}] =$$

$$\boldsymbol{F(D, \theta)} = \mathbb{E}_{\boldsymbol{q(w|\theta)}}[\boldsymbol{\log(q(w|\theta))}] - \mathbb{E}_{\boldsymbol{q(w|\theta)}}[\boldsymbol{\log(p(w))}] - \mathbb{E}_{\boldsymbol{q(w|\theta)}}[\boldsymbol{\log(p(D|w))}]$$

- all three terms are energy expectations with respect to variational distribution $q(w|\theta)$

# Variational Inference (cont.)

*Cost function (cont.):*

$$F(D, \theta) = \mathbb{E}_{q(w|\theta)}[\log(q(w|\theta))] - \mathbb{E}_{q(w|\theta)}[\log(p(w))] - \mathbb{E}_{q(w|\theta)}[\log(p(D|w))]$$

- cost function $F(D, \theta)$ can be approximated by drawing samples $w^{(i)}$ from $q(w|\theta)$

$$F(D, \theta) \approx \frac{1}{N} \sum_{i=1}^{N} [\log(q(w^{(i)}|\theta)) - \log(p(w^{(i)})) - \log(p(D|w^{(i)}))]$$

*Example:*

- Gaussian distribution for the variational posterior $q(w|\theta)$, parameterized by $\theta = \{\mu, \sigma\}$
- $\mu$ - mean vector of the distribution
- $\sigma$ - standard deviation vector

We do not parameterize BNN on weights directly but on $\mu$ and $\sigma$ – we double the number of parameters compared to a plain BNN.

# Summary

Bayesian Neural Networks – *Posterior Variational Inference*

*BNN Probabilistic Model:* $p(y|x, w)$

*Bayesian Neural Network Training*

- *Training Dataset:* $D = \{x^{(i)}, y^{(i)}\}$

- *Maximum Likelihood Estimate (MLE)*

- *Maximum a Posteriori Estimate (MAP)*

*Posterior Variational Inference*

- *Variational Distribution:* $q(w|\theta)$

- *Cost Function:* $F(D, \theta) = KL\big(q(w|\theta) \parallel p(w)\big) - \mathbb{E}_{q(w|\theta)}[log(p(D|w))]$

*Next Lesson:*

- Causal Inference

# Thank You!

Questions?