

# Modelling Uncertainty with Bayesian Neural Networks

Ellis Brown, Melanie Manko, Ethan Matlin

Columbia University

EECS6699

September 21, 2020

Neural Networks are really good at predicting stuff

But what about when they're wrong?

We want to know something about the certainty of a NN's prediction.

A natural framework for modelling uncertainty is the Bayesian paradigm.

# Experiment 1: Why Uncertainty Matters

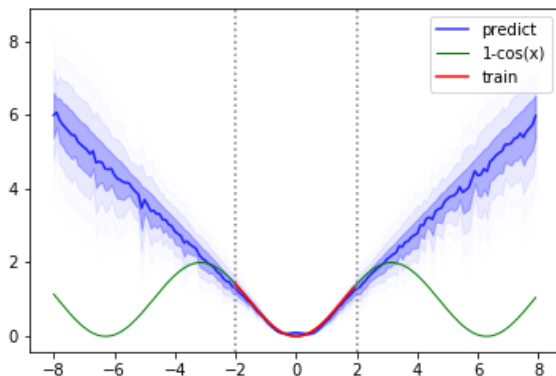


Figure: ReLU Network with 5 Hidden Layers and 1024 neurons per layer

But you just fabricated an example to fit your story...

Uncertainty following the Great Recession: Gross Domestic Product

Figure: ReLU Network with 5 Hidden Layers and 1024 neurons per layer

# Experiment 2: But what about Softmax? Doesn't that give uncertainty?

$f(x)$

$(f(x))$

# Experiment 2: But what about Softmax? Doesn't that give uncertainty?

$f(x)$

$(f(x))$

Neural Networks Need Better Notions of Uncertainty

Methods of Reasoning about Uncertainty

Infinite Bayesian NN  $\rightarrow$  Gaussian Process

Finite Bayesian NN: Use Numerical Techniques (MCMC, Variational Inference, Dropout)

How do priors on the weights translate to priors on functions?

Activation Functions and the Posterior Distribution

Shortfalls of the Approach

# Bayesian Neural Networks and Gaussian Processes

Infinite neural network  $\Leftrightarrow$  Gaussian Process (Neal [1995], Williams [1998])

1 hidden layer

Bounded nonlinearities

Weights ( $u_{ij}$  and  $w_{jk}$ ) are i.i.d with zero mean and finite variance

$$f_k(x) = \sum_{j=1}^H w_{jk} \left( \sum_{i=1}^I u_{ij} x_i \right) \quad (1)$$

Central Limit Theorem  $\Rightarrow f_k(x) \sim \mathcal{N} \left( 0; \frac{2}{b} + H^2 V(x) \right)$

$V(x)$  is the covariance function of the Gaussian Process corresponding to the network.



# Specific Covariance Functions (Williams [1998])

Gaussian Nonlinearity in NN()

$$V_G(x; x^0) = \frac{e^{-d}}{u} \exp\left(-\frac{x^T x}{2 \frac{2}{m}}\right)$$

Scaled Squared Exponential in GP

$$\exp\left(-\frac{(x - x^0)^T (x - x^0)}{2 \frac{2}{s}}\right) \exp\left(-\frac{x^{0T} x^0}{2 \frac{2}{m}}\right)$$

Sigmoid Nonlinearity in NN()

$$V_{\text{erf}}(x; x^0) = \frac{2}{\pi} \sin^{-1} \left[ \frac{2x^T x^0}{(1 + 2x^T x)(1 + 2x^{0T} x^0)} \right]$$

ArcSin Covariance Function in GP

# Experiment 3: In practice though we have only finite Networks

(a) NN (Gaussian Nonlinearity, 1 Hidden Layer);  
GP (Squared Exponential Covariance)

(b) NN (Gaussian Nonlinearity, 5 Hidden Layers);  
GP (Squared Exponential Covariance)

## How to evaluate a finite NN?

Monte Carlo techniques: HMC, RWMH, SMC, etc. Neal [1995]

Variational Inference: Barber & Bishop [1998], Graves [2011], Blundell et al. [2015]

Dropout is equivalent to Variational Inference in Gaussian Processes (Gal [2016], Gal & Ghahramani [2015b], Gal & Ghahramani [2016a], Gal & Ghahramani [2016b], Gal & Ghahramani [2015a])

Used as a regularization technique [Hinton et al. [2012], Srivastava et al. [2014]]

Intuitively, makes sense as a way to get a distribution with uncertainty

NN Variational Inference Objective Function:

$$\hat{\mathcal{L}} = \frac{N}{M} \sum_{i=1}^M \log p(y_i | f^{g(\cdot; \theta)}) + \text{KL}(q(\cdot) | p(\cdot)) \quad (2)$$

NN Loss Function:

$$\mathcal{L} = \log p(y | f^{g(\cdot; \theta)}) + \text{const.} + \sum_{1,j} W_{1,j} + \sum_{2,j} W_{2,j} + \sum_{3,j} b_{j,j} \quad (3)$$

KL condition:

$$\frac{\partial}{\partial} \text{KL}(Q(\cdot) | p(\cdot)) = \frac{\partial}{\partial} \sum_{1,j} W_{1,j} + \sum_{2,j} W_{2,j} + \sum_{3,j} b_{j,j} \quad (4)$$

=> the two are equivalent.

# Experiment 4: What prior should I choose? How do we priors correspond to function priors?

$$W \sim N(0; 0.05)$$

$$W \sim \text{Unif}[-0.05; 0.05]$$

$$W \sim N\left(0; \frac{2}{N_{\text{in}} + N_{\text{out}}}\right)$$

$$W \sim \text{Unif}\left[-\frac{6}{N_{\text{in}}}; \frac{6}{N_{\text{in}}}\right]$$

Figure: ReLU (Untrained), 5 Hidden Layers, 1024 neurons per layer

# Experiment 5: Activation Functions and Uncertainty Estimates

Figure: Linear; 5 Hidden Layers, 1024 Neurons per Layer

# Experiment 5: Activation Functions and Uncertainty Estimates

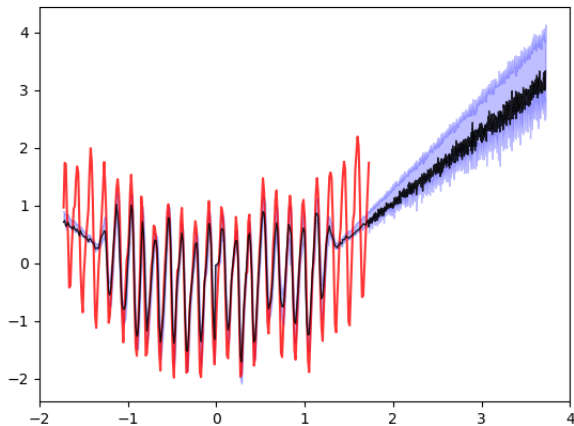


Figure: ReLu; 5 Hidden Layers; 1024 Neurons per Layer

# Experiment 5: What Activation Function?

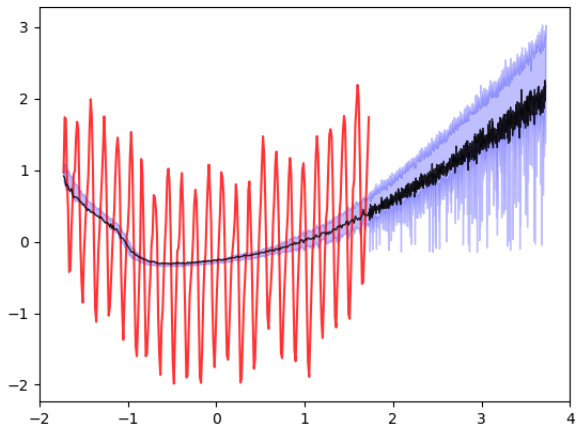


Figure: Softplus; 5 Hidden Layers; 1024 Neurons per Layer



# Experiment 5: What Activation Function?

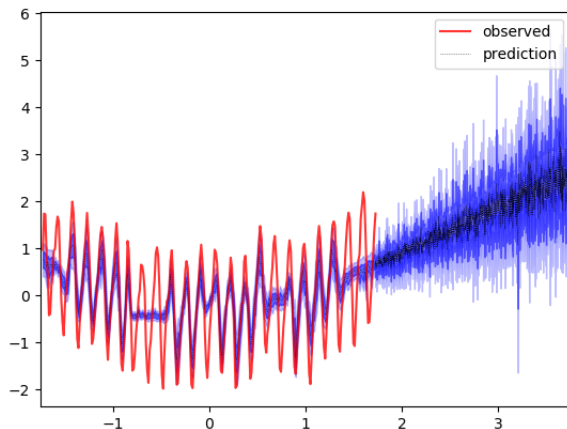


Figure: Softplus; 5 Hidden Layers; 1024 Neurons per Layer

