# Bayesian Deep Learning:
# Motivation and Model Definition

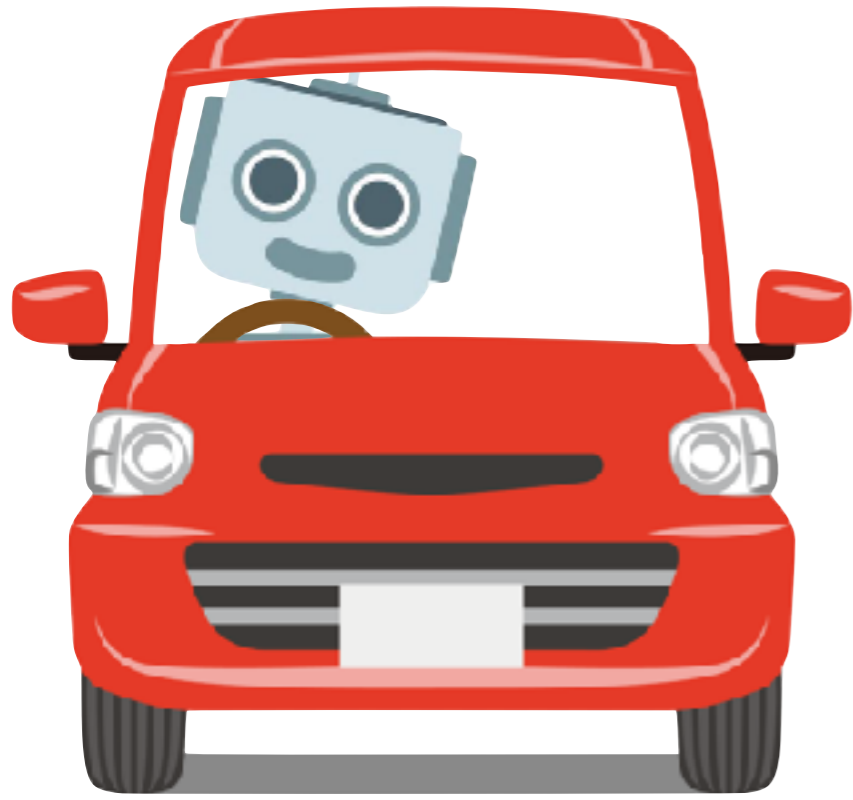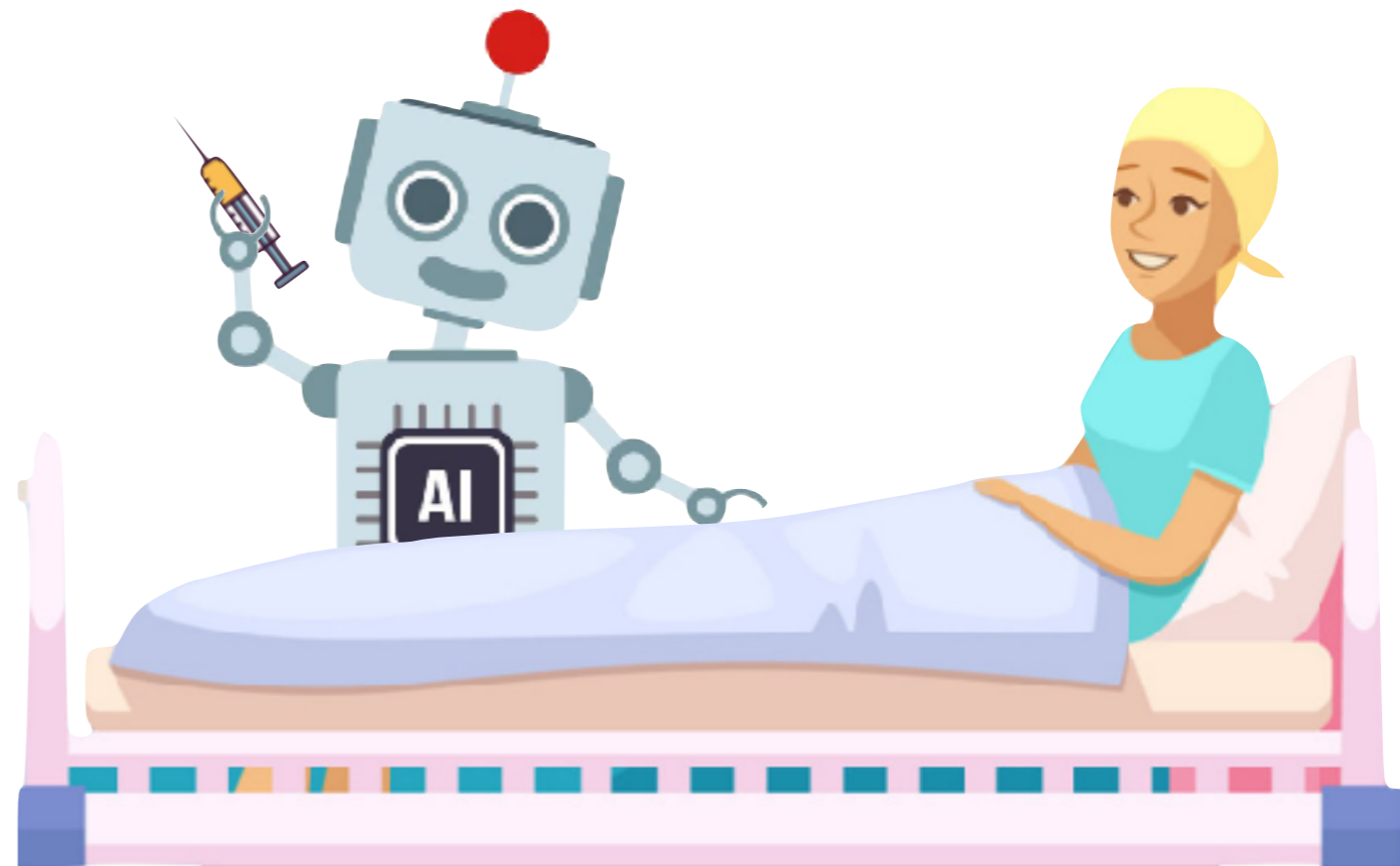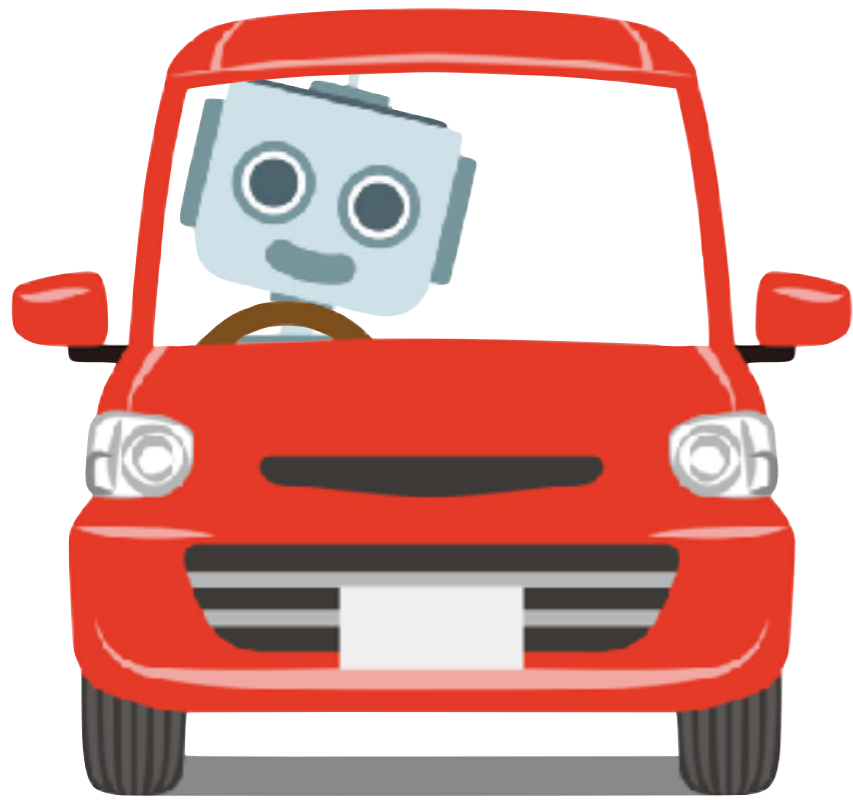## Eric Nalisnick
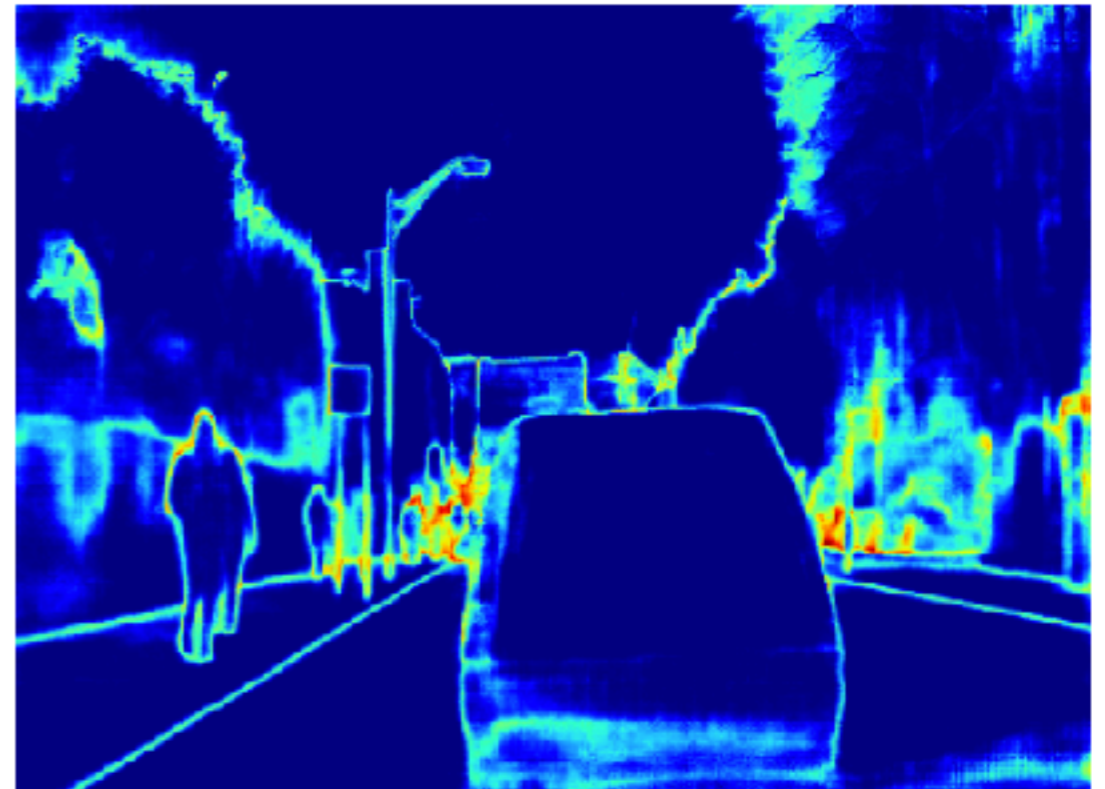
AMLAB
Amsterdam
Machine Learning Lab

Deep Learning II,
University of Amsterdam

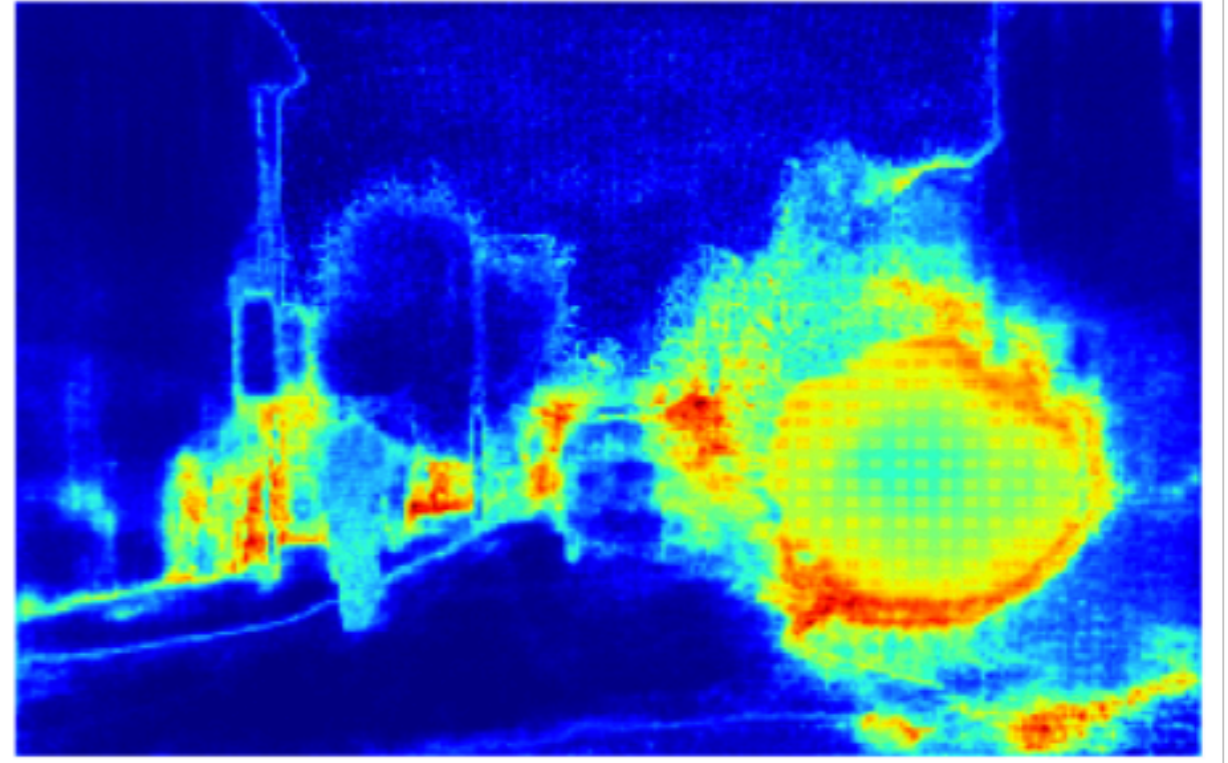Images from Kendall and Gal, "What uncertainties do we need in Bayesian deep learning for computer vision?", *NeurIPS 2017*.

Images from Besnier et al., "Learning Uncertainty For Safety-Oriented Semantic Segmentation In Autonomous Driving", *ICIP 2021*.

# Neural Network



input layer

hidden layer 1  hidden layer 2

output layer

Data set of
hand-written digits

8

[Louizos & Welling, ICML 2017]

# Traditional NN

Images from Blundell et al., "Weight Uncertainty in Neural Networks", *ICML 2015*.

# Traditional NN      Bayesian NN

Images from Blundell et al., "Weight Uncertainty in Neural Networks", *ICML 2015*.

$$p(y^* \mid x^*, D) = \int_w p(y^* \mid x^*, w)\ p(w \mid D)\ dw$$

$$p(y* \mid x*, D) = \int_{w} p(y* \mid x*, w) \, p(w \mid D) \, dw$$

# Model Definition

Data model

$$y \sim p\left(y \mid x, W_1, \ldots, W_L\right)$$

Assume NNs are fully-connected, feedforward, unless stated otherwise.

Data model

$$y \sim p\left(y \mid x, W_1, \ldots, W_L\right)$$

Assume NNs are fully-connected, feedforward, unless stated otherwise.

For real-valued regression...

$$y \sim N\left(\mu = f\left(x; W_1, W_2, W_3\right), \sigma_0^2\right)$$

Data model

$$y \sim p\left(y \mid x, W_1, \ldots, W_L\right)$$

Assume NNs are fully-connected, feedforward, unless stated otherwise.

For classification…

$$y \sim \text{Categorical}\left(\pi = f\left(x; W_1, W_2, W_3\right)\right)$$

# Data model

$$y \sim p\left(y \mid x, W_1, \ldots, W_L\right)$$

# Prior per weight

$$\mathbf{w} \sim p(\mathbf{w})$$

WEIGHT MATRIX

# Prior per weight

$$w \sim N(0, \sigma^2)$$

WEIGHT MATRIX

# Prior per layer

$$\mathsf{W}_l \sim p(\mathsf{W}_l)$$

WEIGHT MATRIX

# Prior per layer

$$W_l \sim N(\mathbf{0}, \Sigma)$$



WEIGHT MATRIX

$D_{l-1}$

$D_l$

# Prior per layer

$$W_l \sim N(\mathbf{0}, \Sigma)$$

WEIGHT MATRIX

Size: $(D_{l-1} + D_l)^2$

$D_{l-1}$

$D_l$

# Joint prior

$$\mathsf{W}_1, \ldots, \mathsf{W}_L \sim p\left(\mathsf{W}_1, \ldots, \mathsf{W}_L\right)$$

# Joint prior

$$W_1, \ldots, W_L \sim p\left(W_1, \ldots, W_L\right)$$

WEIGHT MATRIX

WEIGHT MATRIX

WEIGHT MATRIX

NOT USUALLY DONE

Joint prior

$$W_1, \ldots, W_L \sim p\left(W_1, \ldots, W_L\right)$$

$$W_1, \ldots, W_L \sim N(\mathbf{0}, \Sigma)$$

Size: (# total weights)$^2$

# Posterior

$$p\left(\mathsf{W}_1, \ldots, \mathsf{W}_L \,|\, \mathsf{y}, \mathsf{x}\right) =$$

# Posterior

$$p\left(\mathsf{W}_1, \ldots, \mathsf{W}_L \mid \mathsf{y}, \mathsf{x}\right) =$$

$$\frac{p\left(\mathsf{y} \mid \mathsf{x}, \mathsf{W}_1, \ldots, \mathsf{W}_L\right) \prod_{l=1}^{L} p(\mathsf{W}_l)}{p\left(\mathsf{y} \mid \mathsf{x}\right)}$$

## Posterior

$$p\left(\mathsf{W}_1, \ldots, \mathsf{W}_L \mid \mathsf{y}, \mathsf{x}\right) =$$

$$\frac{p\left(\mathsf{y} \mid \mathsf{x}, \mathsf{W}_1, \ldots, \mathsf{W}_L\right) \prod_{l=1}^{L} p(\mathsf{W}_l)}{p\left(\mathsf{y} \mid \mathsf{x}\right)}$$

# Posterior

$$p\left(\mathsf{W}_1, \ldots, \mathsf{W}_L \mid \mathsf{y}, \mathsf{x}\right) =$$

$$\frac{p\left(\mathsf{y} \mid \mathsf{x}, \mathsf{W}_1, \ldots, \mathsf{W}_L\right) \prod_{l=1}^{L} p(\mathsf{W}_l)}{\int_{\mathsf{W}_1, \ldots, \mathsf{W}_L} p(\mathsf{y} \mid \mathsf{x}, \mathsf{W}_1, \ldots, \mathsf{W}_L) \prod_{l} p(\mathsf{W}_l) \, d\mathsf{W}_1, \ldots, \mathsf{W}_L}$$

# Posterior Predictive

$$p(y^* \mid x^*, y, x) =$$

$$\int_{W_1,\ldots,W_L} p(y^* \mid x^*, W_1, \ldots, W_L) \, p(W_1, \ldots, W_L \mid y, x) \, dW_1, \ldots, W_L$$

# p(y* | x*, $D$)

$$w \sim N(0, \sigma^2 = 5)$$

$$\sigma_0 \sim Gamma\,(1/2,\ 1)$$

$$y \sim N\left(\mu = f\left(x; W_1, W_2, W_3\right), \sigma_0^2\right)$$

# Bayesian Deep Learning: Priors

## Eric Nalisnick

Deep Learning II,
University of Amsterdam

$$\frac{p\left(\text{y} \mid \text{x}, \text{W}_1, \ldots, \text{W}_L\right) \prod_{l=1}^{L} p(\text{W}_l)}{p\left(\text{y} \mid \text{x}\right)}$$

*Garbage in:* arbitrary priors
*Garbage out:* uncontrollable error bars

Michael I. Jordan, *MLSS* (2017)

# Normal Prior



WEIGHT MATRIX

$p(\mathbf{w})$

0

$\mathbf{w}$

# Normal Prior

As NN becomes infinitely wide, it converges to a *Gaussian process*

$$w \sim N(0, \sigma^2/H)$$

# Normal Prior

As NN becomes infinitely wide, it converges to a *Gaussian process*

$$\mathbf{w} \sim \mathbf{N}(0, \, \sigma^2/H)$$

*"With Gaussian priors the contributions of individual units are all negligible, and consequently, these units do not represent 'hidden features' that capture important aspects of the data"* [Neal, 1995]

# Normal Prior

As NN becomes infinitely wide, it converges to a *Gaussian process*



[Matthews et al., 2018]

# Multivariate Normal

# Multivariate Normal

# Multivariate Normal

$$\begin{bmatrix} \sigma_1^2 & & & & & \\ & \sigma_2^2 & & & & \\ & & \sigma_3^2 & & & \\ & & & \sigma_4^2 & & \\ & & & & \sigma_5^2 & \\ & & & & & \sigma_6^2 \end{bmatrix}$$

$$\Sigma$$

# Multivariate Normal

$$\begin{bmatrix} \sigma_1^2 & & & & & \\ & \sigma_2^2 & & & & \\ & & \sigma_3^2 & & & \\ & & & \sigma_4^2 & & \\ & & & & \sigma_5^2 & \\ & & & & & \sigma_6^2 \end{bmatrix}$$

$\Sigma$



$\mathsf{h}_{l-1}\mathsf{W}_l$

# Hierarchical Priors

$$\tau \sim p(\tau)$$

$$\mathtt{w} \sim p(\mathtt{w} \,|\, \tau)$$

# Hierarchical Priors

$$\tau \sim p(\tau)$$

$$\mathbf{w} \sim N(0, \tau^2)$$

# Hierarchical Priors: Structure

$$\tau_i \sim p(\tau)$$

$$\mathsf{w}_{i,j} \sim \mathsf{N}(0, \ \tau_i^2)$$



$D_l$

$D_{l-1}$

*i* indexes rows

# Hierarchical Priors: Structure

$$\tau_i \sim p(\tau)$$

$$\mathbf{w}_{i,j} \sim \mathbf{N}(0, \tau_i^2)$$

$D_l$

MacKay, 1994

"Automatic Relevance Determination"

$D_{l-1}$

# Hierarchical Priors: Heavy-Tails

$$\tau^2 \sim \Gamma^{-1}(\alpha, \beta)$$

$$w \sim N(0, \tau^2)$$

$$t(w) = \int_\tau N(w; 0, \tau^2) \, \Gamma^{-1}(\tau^2; \alpha, \beta) \, d\tau$$

# Hierarchical Priors: Heavy-Tails

$$\tau^2 \sim \text{Cauchy}^+(\sigma)$$

$$w \sim N(0, \tau^2)$$



horseshoe(3)

N(0, 1)

# Hierarchical Priors: Heavy-Tails

# Hierarchical Priors: Heavy-Tails



Forget regularization: "bounded Influence"

# Hierarchical Priors: Heavy-Tails

Infinitely wide NN no longer converges to a *Gaussian process;* instead a *jump* process.

# Discrete Priors

$$w \sim \text{Bernoulli}(\pi)$$

Interesting due to their computational efficiency [Soudry et al., 2014] and biological plausibility [Baldassi et al., 2007].

But no access to gradients.

# Other Architectures: ResNets

# Other Architectures: ResNets



$$h \times W$$

$d_{l-1}$

$d_{l-1}$

$d_l$

SKIP CONNECTION

Allows information to bypass interaction with the weights

# Other Architectures: ResNets

Scale shared across matrix

$d_{l-1}$

**h**

X

**W**

$d_{l-1}$

$d_l$

SKIP CONNECTION

Allows information to bypass interaction with the weights

# Other Architectures: ResNets



Bayesian shrinkage can control the effective depth of the network

# Other Architectures: ResNets



Bayesian shrinkage can control the effective depth of the network

$$w_{l,i,j} \sim \mathsf{N}(0, \gamma_l^2) \qquad \gamma_l \sim p(\gamma)$$

# Other Architectures: ResNets



Bayesian shrinkage can control the effective depth of the network

$$w_{l,i,j} \sim N(0, \tau_i^2 \gamma_l^2) \quad \gamma_l \sim p(\gamma) \quad \tau_i \sim p(\tau)$$

# Other Architectures: ConvNet

# Other Architectures: LSTM

# Tuning the Prior: Type II MLE

$$p(W; \psi)$$

# Tuning the Prior: Type II MLE

$$p(W; \psi)$$

$$p\left(y \mid x; \psi\right)$$

# Tuning the Prior: Type II MLE

$$p(W; \psi)$$

$$p\left(y \mid x; \psi\right) = \int_W p(y \mid x, W) \; p(W; \psi) \; dW$$

# Tuning the Prior: Type II MLE

$$p(W; \psi)$$

$$p(y \mid x; \psi) = \int_W p(y \mid x, W) \, p(W; \psi) \, dW$$

# Summary

⊗ **Normal priors**: easy to implement, correspond to Gaussian process in the infinite limit.

⊗ **Hierarchical priors**: good for inducing structure and heavy-tails.

⊗ **Discrete priors**: efficient but hard to implement.

# Priors for Neural Networks

Herbert K. H. Lee

Department of Applied Mathematics and Statistics

University of California, Santa Cruz

herbie@ams.ucsc.edu

### Abstract

Neural networks are commonly used for classification and regression. The Bayesian approach may be employed, but choosing a prior for the parameters presents challenges. This paper reviews several priors in the literature and introduces Jeffreys priors for neural network models. The effect on the posterior is demonstrated through an example.

**Key Words:** nonparametric classification; nonparametric regression; Bayesian statistics; prior sensitivity

## 1 Introduction

Neural networks are a popular tool for nonparametric classification and regression. They offer a computationally tractable model that is fully flexible, in the sense of being able to approximate a wide range of functions (such as all continuous functions). Many references on neural networks are available (Bishop, 1995; Fine, 1999; Ripley, 1996). The Bayesian approach is appealing as it allows full accounting for uncertainty in the model and the choice of model (Lee, 2001; Neal, 1996). An important decision in any Bayesian analysis is the choice of prior. The idea is that your prior should reflect your current beliefs (either from previous data

# Chapter 3

## Survey of Neural Network Priors

[Lee, 2004]

*We demand rigidly defined areas of doubt and uncertainty!*

Douglas Adams
The Hitchhiker's Guide to the Galaxy

Having covered the basics of Bayesian NNs and strategies for inferring their posterior, I now turn to the focal point of the dissertation: prior distributions for both conditional NNs and density networks. Surprisingly, a broad review of Bayesian NN priors has been performed by only Robinson [2001], which is now considerably out of date. Thus, in this chapter I survey the existing work on NN priors, some of which was performed in the early days of Bayesian NNs and therefore also discussed by Robinson [2001]. However, most of the work is recent, some having been conducted concurrently with my own work to be presented in the coming chapters.

NNs have been applied to a myriad of different problems over the past thirty years, and this of course makes it impossible to discuss every prior ever used for a NN. Instead, I attempt to summarize broad themes from the literature that pertain to core NN methodology. For instance,

[Nalisnick, 2018]

eural Networks

K. H. Lee

Mathematics and Statistics

lifornia, Santa Cruz

us.ucsc.edu

stract

cation and regression. The Bayesian approach may be presents challenges. This paper reviews several priors neural network models. The effect on the posterior is

rametric regression; Bayesian statistics; prior sensitiv-

Neural networks are a popular tool for nonparametric classification and regression. They offer a computationally tractable model that is fully flexible, in the sense of being able to approximate a wide range of functions (such as all continuous functions). Many references on neural networks are available (Bishop, 1995; Fine, 1999; Ripley, 1996). The Bayesian approach is appealing as it allows full accounting for uncertainty in the model and the choice of model (Lee, 2001; Neal, 1996). An important decision in any Bayesian analysis is the choice of prior. The idea is that your prior should reflect your current beliefs (either from previous data

# Chapter 3

## Survey of Neural Network Priors

*We demand rigidly defined areas of doubt and uncertainty!*

Douglas A[...]

The Hitchhiker's Guide to the G[...]

Having covered the basics of Bayesian NNs and strategies for inferring their posterior, I [...] turn to the focal point of the dissertation: prior distributions for both conditional NNs and [...] sity networks. Surprisingly, a broad review of Bayesian NN priors has been performed by [...] Robinson [2001], which is now considerably out of date. Thus, in this chapter I survey th[...] isting work on NN priors, some of which was performed in the early days of Bayesian NNs [...] therefore also discussed by Robinson [2001]. However, most of the work is recent, some ha[...] been conducted concurrently with my own work to be presented in the coming chapters.

NNs have been applied to a myriad of different problems over the past thirty years, and [...] of course makes it impossible to discuss every prior ever used for a NN. Instead, I attemp[...] summarize broad themes from the literature that pertain to core NN methodology. For insta[...]

[Nalisnick, 2018]

eural Networks

K. H. Lee

[Lee, 2004]

Neural networks are a popular tool for [...] tationally tractable model that is fully fl[...] functions (such as all continuous functions [...] Fine, 1998; Ripley, 1996). The Bayesian a[...] the model and the choice of model (Lee, 20[...] the choice of prior. The idea is that your [...]

## PRIORS IN BAYESIAN DEEP LEARNING: A REVIEW

**Vincent Fortuin**
Department of Computer Science
ETH Zürich
Zürich, Switzerland
fortuin@inf.ethz.ch

### ABSTRACT

While the choice of prior is one of the most critical parts of the Bayesian inference workflow, recent Bayesian deep learning models have often fallen back on vague priors, such as standard Gaussians. In this review, we highlight the importance of prior choices for Bayesian deep learning and present an overview of different priors that have been proposed for (deep) Gaussian processes, variational autoencoders, and Bayesian neural networks. We also outline different methods of learning priors for these models from data. We hope to motivate practitioners in Bayesian deep learning to think more carefully about the prior specification for their models and to provide them with some inspiration in this regard.

### 1 Introduction

Bayesian models have gained a stable popularity in data analysis [1] and machine learning [2]. Especially in recent years, the interest in combining these models with deep learning has surged[1]. The main idea of Bayesian modeling is to infer a *posterior* distribution over the parameters $\theta$ of the model given some observed data $\mathcal{D}$ using Bayes' theorem [3, 4] as

$$p(\theta \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \theta)\, p(\theta)}{p(\mathcal{D})} = \frac{p(\mathcal{D} \mid \theta)\, p(\theta)}{\int p(\mathcal{D} \mid \theta)\, p(\theta)\, d\theta} \tag{1}$$

[Fortuin, 2021]

# Bayesian Deep Learning:
## Posterior Inference

Eric Nalisnick



Deep Learning II,
University of Amsterdam

$$p\left(\mathsf{W}_1, \dots, \mathsf{W}_L \mid \mathsf{y}, \mathsf{x}\right)$$

# Non-identifiability

$$h_1 = f(w_{1,1}x)$$

$$h_2 = f(w_{1,2}x)$$

$$\hat{y} = w_{2,1}h_1 + w_{2,2}h_2$$

# Non-identifiability

$$h_1 = f(w_{1,1}x)$$

$$h_2 = f(w_{1,2}x)$$

$$\hat{y} = w_{2,1}h_1 + w_{2,2}h_2$$

# Non-identifiability

$$h_1 = f(w_{1,1}x)$$

$$h_2 = f(w_{1,2}x)$$

$$\hat{y} = w_{2,1}h_1 + w_{2,2}h_2$$

# Non-identifiability

⊗ Permutation invariance.

⊗ Scale invariance for ReLUs:

$$\text{ReLU}(x) = (1/\alpha) \cdot \text{ReLU}(\alpha \cdot x), \quad \forall \alpha > 0$$

# Conjugacy?

Not in general…

# Conjugacy?

Not in general…

But sometimes for the last layer:

$$\log N\left(y \mid x, W_1, \dots, W_L\right) =$$
$$\frac{-1}{2\sigma_0^2}\left(y - h_{L-1}W_L\right)^2 + \dots$$

# Conjugacy?

Not in general…

But sometimes for the last layer.

**"neural linear model"**

$$\frac{-1}{2\sigma_0^2}\left(y - h_{L-1}W_L\right)^2 + \dots$$

# MAP Estimation

$$p\left(\mathsf{W}_1, \ldots, \mathsf{W}_L \,|\, \mathsf{y}, \mathsf{x}\right) \propto$$

$$\log p\left(\mathsf{y} \,|\, \mathsf{x}, \mathsf{W}_1, \ldots, \mathsf{W}_L\right) + \sum_{l=1}^{L} \log p(\mathsf{W}_l)$$

# MAP Estimation

$$p\left(\mathsf{W}_1, \ldots, \mathsf{W}_L \mid \mathsf{y}, \mathsf{x}\right) \propto$$

$$\log p\left(\mathsf{y} \mid \mathsf{x}, \mathsf{W}_1, \ldots, \mathsf{W}_L\right) + \sum_{l=1}^{L} \log p(\mathsf{W}_l)$$

For normal priors…

$$-\sum_{l=1}^{L} \frac{1}{2\sigma_l^2} ||\mathsf{W}_l||_2^2 + \text{const}.$$

# MAP Estimation

$$p\left(\mathsf{W}_1, ..., \mathsf{W}_L \mid \mathsf{y}, \mathsf{x}\right) \propto$$

**Equivalent to weight decay**

For normal priors…

$$-\sum_{l=1}^{L} \frac{1}{2\sigma_l^2} ||\mathsf{W}_l||_2^2 + \text{const}.$$

# MAP Estimation

Caution: MAP estimates have very different characteristics than the true posterior (e.g. sparsity)

# Markov Chain Monte Carlo (MCMC)

# Markov Chain Monte Carlo (MCMC)

$$p\left(W_1, \ldots, W_L \mid y, x\right) \approx \frac{1}{S}\sum_{s=1}^{S} \delta\left[\hat{W}_{1,s}, \ldots, \hat{W}_{L,s}\right]$$

Initialize $\mathbf{w}^0$

For t=1 to T:

.
.

Initialize $\mathbf{w}^0$

For t=1 to T:

 Sample $u \sim \text{Uniform}(0,1)$

 :

Initialize $\mathbf{w}^0$

For t=1 to T:

    Sample   $u \sim \text{Uniform}(0,1)$

    Sample   $\mathbf{w}^* \sim q(\mathbf{w}^* | \mathbf{w}^{t-1})$

    :

Initialize $\mathbf{w}^0$

For t=1 to T:

  Sample $\quad u \sim \mathsf{Uniform}(0,1)$

  Sample $\quad \mathbf{w}^* \sim \mathsf{q}(\mathbf{w}^* | \mathbf{w}^{t-1})$

  If $\quad u < \min \left\{ 1, \dfrac{\mathsf{p}(y, \mathbf{w}^* | x) \; \mathsf{q}(\mathbf{w}^{t-1} | \mathbf{w}^*)}{\mathsf{p}(y, \mathbf{w}^{t-1} | x) \; \mathsf{q}(\mathbf{w}^* | \mathbf{w}^{t-1})} \right\} :$

    $\mathbf{w}^t = \mathbf{w}^*$

Initialize $\mathbf{w}^0$

For t=1 to T:

    Sample   $u \sim \mathsf{Uniform}(0,1)$

    Sample   $\mathbf{w}^* \sim \mathsf{q}(\mathbf{w}^* | \mathbf{w}^{t-1})$

    If  $u < \min \left\{ 1, \dfrac{\mathsf{p}(y, \mathbf{w}^* | x)\ \mathsf{q}(\mathbf{w}^{t-1} | \mathbf{w}^*)}{\mathsf{p}(y, \mathbf{w}^{t-1} | x)\ \mathsf{q}(\mathbf{w}^* | \mathbf{w}^{t-1})} \right\}$ :

        $\mathbf{w}^t = \mathbf{w}^*$

    Else:

        $\mathbf{w}^t = \mathbf{w}^{t-1}$

# Hamiltonian Monte Carlo (HMC)

Generate proposal by iterating:

(assuming the identity matrix for the mass)

$$v^{m+1} = v^m + \alpha \nabla_w \log p(w^m | y, x)$$

$$w^{m+1} = w^m + \alpha' \cdot v^m$$

where $\alpha$ and $\alpha'$ are step sizes and $v^0 \sim N(0,1)$

# Hamiltonian Monte Carlo (HMC)

Generate proposal by iterating:

(assuming the identity matrix for the mass)

$$v^{m+1} = v^m + \alpha \nabla_w \log p(w^m \,|\, y, x)$$

$$w^{m+1} = w^m + \alpha' \cdot v^m$$

where $\alpha$ and $\alpha'$ are step sizes and $v^0 \sim N(0,1)$

Propose: $w^* = w^M$

# Hamiltonian Monte Carlo (HMC)

**What Are Bayesian Neural Network Posteriors Really Like?**

| Pavel Izmailov | Sharad Vikram | Matthew D. Hoffman | Andrew Gordon Wilson |
| New York University | Google Research | Google Research | New York University |

Computation done
on 512 TPUs

# Hamiltonian Monte Carlo (HMC)

**What Are Bayesian Neural Network Po**

Pavel Izmailov    Sharad Vikram    Matthew D. Hof
New York University    Google Research    Google Resea

Computation done
on 512 TPUs

# HMC to Langevin Dynamics

One step iteration of HMC:

$$w^1 = w^0 + \alpha' \cdot v^1$$
$$= w^0 + \alpha' \cdot \left( \alpha \, \nabla_w \log p(w^0 | y, x) + v^0 \right)$$

# HMC to Langevin Dynamics

One step iteration of HMC:

$$w^1 = w^0 + \alpha' \cdot v^1$$

$$= w^0 + \alpha' \cdot \left( \alpha \, \nabla_w \log p(w^0 | y, x) + v^0 \right)$$

Langevin Dynamics:

$$w^{m+1} = w^m + \alpha'' \cdot \nabla_w \log p(w^m | y, x) + \hat{v}$$

$$\hat{v} \sim N(0, \epsilon)$$

# Langevin Dynamics

$$w^{m+1} = w^m + \alpha'' \cdot \nabla_w \log p(w^m | y, x) + \hat{v}$$

$$\hat{v} \sim N(0, \epsilon)$$

⊗ "Adjusted": Run accept-reject step

⊗ "Unadjusted": Always accept proposal

⊗ Can also use stochastic gradients

# MCMC for ResNet-20 on CIFAR-10

| METRIC | HMC (REFERENCE) | SGMCMC | | | |
|---|---|---|---|---|---|
| | | SGLD | SGHMC | SGHMC CLR | SGHMC CLR-PREC |
| ACCURACY | 89.64 ±0.25 | 89.32 ±0.23 | 89.38 ±0.32 | **89.63 ±0.37** | 87.46 ±0.21 |
| AGREEMENT | 94.01 ±0.25 | 91.54 ±0.15 | 91.98 ±0.35 | **92.67 ±0.52** | 90.96 ±0.24 |
| TOTAL VAR | 0.074 ±0.003 | 0.110 ±0.001 | 0.109 ±0.001 | **0.099 ±0.006** | 0.111 ±0.002 |

# Variational Inference

$$p\left(W_1, \ldots, W_L \mid y, x\right) \approx q\left(W_1, \ldots, W_L; \phi\right)$$



Image from Blei et al., "Variational Inference: A Review for Statisticians," JASA 2017

We usually need to assume some degree of factorization.

We usually need to assume some degree of factorization.

Over layers:

$$q\left(W_1, \ldots, W_L; \phi\right) = \prod_{l=1}^{L} q\left(W_l; \phi_l\right)$$

We usually need to assume some degree of factorization.

Over layers:

$$q\left(W_1, \ldots, W_L; \phi\right) = \prod_{l=1}^{L} q\left(W_l; \phi_l\right)$$

Over weights ("mean-field"):

$$= \prod_{l=1}^{L} \prod_{d=1}^{D_l} q\left(w_{l,d}; \phi_{l,d}\right)$$

Normals are common, for instance.

Over layers:

$$q\left(W_1, \ldots, W_L; \phi\right) = \prod_{l=1}^{L} N\left(\boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l\right)$$

Over weights ("mean-field"):

$$= \prod_{l=1}^{L} \prod_{d=1}^{D_l} N\left(\mu_{l,d}, \sigma_{l,d}^2\right)$$

# Optimization Objective

$$\phi^* = \text{argmin}_\phi \; \mathbb{D}\left[ q\left(w; \phi\right) \mid\mid p\left(w \mid y, x\right) \right]$$

# Optimization Objective

$$\phi* = \text{argmin}_\phi \text{ KLD}\left[q\left(w; \phi\right) || p\left(w | y, x\right)\right]$$

# Optimization Objective

$$\phi* = \text{argmin}_\phi \ \text{KLD}\left[q\left(w;\phi\right)||p\left(w|y,x\right)\right]$$

$$= \text{argmin}_\phi \ \int_w q\left(w;\phi\right) \log \frac{q\left(w;\phi\right)}{p\left(w|y,x\right)} \, dw$$

# Optimization Objective

$$\text{KLD}\left[q\left(w;\phi\right)||p\left(w|y,x\right)\right] =$$

# Optimization Objective

$$\mathrm{KLD}\left[q\left(w;\phi\right)||p\left(w|y,x\right)\right] =$$

$$\mathbb{E}_{q_\phi}\left[-\log p\left(y|x,w\right)\right]+$$

$$\mathrm{KLD}\left[q(w;\phi)||p(w)\right] + \mathrm{const}.$$

# Optimization Objective

$$\text{KLD}\left[q\left(w;\phi\right)||p\left(w|y,x\right)\right]=$$

$$\mathbb{E}_{q_\phi}\left[-\log p\left(y|x,w\right)\right]+$$

$$\text{KLD}\left[q(w;\phi)||p(w)\right]+\text{const}.$$

# Reparameterization Trick

$$\mathbb{E}_{q_\phi} \left[ -\log p \left( y \mid x, w \right) \right]$$

# Reparameterization Trick

$$\mathbb{E}_{q_\phi} \left[ -\log p\left(y \mid x, w\right) \right]$$

$$= \mathbb{E}_\eta \left[ -\log p\left(y \mid x, w = g(\eta; \phi)\right) \right]$$

# Reparameterization Trick

$$\boxed{\mathbb{E}_{\mathsf{q}_\phi} \left[ -\log \mathsf{p} \left( \mathsf{y} \mid \mathsf{x}, \mathsf{w} \right) \right]}$$

$$= \mathbb{E}_\eta \left[ -\log \mathsf{p} \left( \mathsf{y} \mid \mathsf{x}, \mathsf{w} = g(\eta; \phi) \right) \right]$$

$$\approx \frac{1}{S} \sum_s -\log \mathsf{p} \left( \mathsf{y} \mid \mathsf{x}, \mathsf{w} = g(\hat{\eta}_s; \phi) \right)$$

$$\frac{\partial}{\partial \phi} \mathbb{E}_{q_\phi} \left[ -\log p \left( y \mid x, w \right) \right]$$

$$\approx -\frac{1}{S} \sum_s \frac{\partial}{\partial \phi} \log p \left( y \mid x, w = g(\hat{\eta}_s; \phi) \right)$$

$$\frac{\partial}{\partial \phi} \mathbb{E}_{q_\phi} \left[ -\log \mathsf{p} \left( \mathsf{y} \mid \mathsf{x}, \mathsf{w} \right) \right]$$

$$\approx -\frac{1}{S} \sum_s \frac{\partial}{\partial \phi} \log \mathsf{p} \left( \mathsf{y} \mid \mathsf{x}, \mathsf{w} = g(\hat{\eta}_s; \phi) \right)$$

Blundell et al., 2015

"Bayes by Backprop"

$$\text{If } q(w; \phi) = N(\mu, \sigma):$$

$$\hat{w} = g(\hat{\eta}; \phi) = \mu + \sigma \cdot \hat{\eta}, \quad \eta \sim N(0,1)$$

If $q(\mathrm{w}; \phi) = \mathrm{N}(\mu, \sigma)$:

$$\hat{\mathrm{w}} = g(\hat{\eta}; \phi) = \mu + \sigma \cdot \hat{\eta}, \quad \eta \sim \mathrm{N}(0,1)$$

$$h \cdot \hat{\mathrm{w}} = h \cdot (\mu + \sigma \cdot \hat{\eta})$$

If $q(w; \phi) = N(\mu, \sigma)$:

$$\hat{w} = g(\hat{\eta}; \phi) = \mu + \sigma \cdot \hat{\eta}, \quad \eta \sim N(0,1)$$

$$h \cdot \hat{w} = h \cdot (\mu + \sigma \cdot \hat{\eta})$$

Or for a general q:

$$\hat{w} = CDF_q^{-1}(\hat{\eta}; \phi), \quad \eta \sim Uniform(0,1)$$

# MCMC for ResNet-20 on CIFAR-10

| | | | | SGMCMC | | |
| Metric | HMC (Reference) | MFVI | SGLD | SGHMC | SGHMC CLR | SGHMC CLR-Prec |
|---|---|---|---|---|---|---|
| Accuracy | 89.64 ±0.25 | 86.45 ±0.27 | 89.32 ±0.23 | 89.38 ±0.32 | **89.63** ±**0.37** | 87.46 ±0.21 |
| Agreement | 94.01 ±0.25 | 88.75 ±0.24 | 91.54 ±0.15 | 91.98 ±0.35 | **92.67** ±**0.52** | 90.96 ±0.24 |
| Total Var | 0.074 ±0.003 | 0.136 ±0.000 | 0.110 ±0.001 | 0.109 ±0.001 | **0.099** ±**0.006** | 0.111 ±0.002 |

# Laplace Approximation

$$p\left(w \mid y, x\right)$$

$$\approx N\left(\hat{w}_{\text{MAP}}, \bar{H}^{-1}(\hat{w}_{\text{MAP}})\right)$$

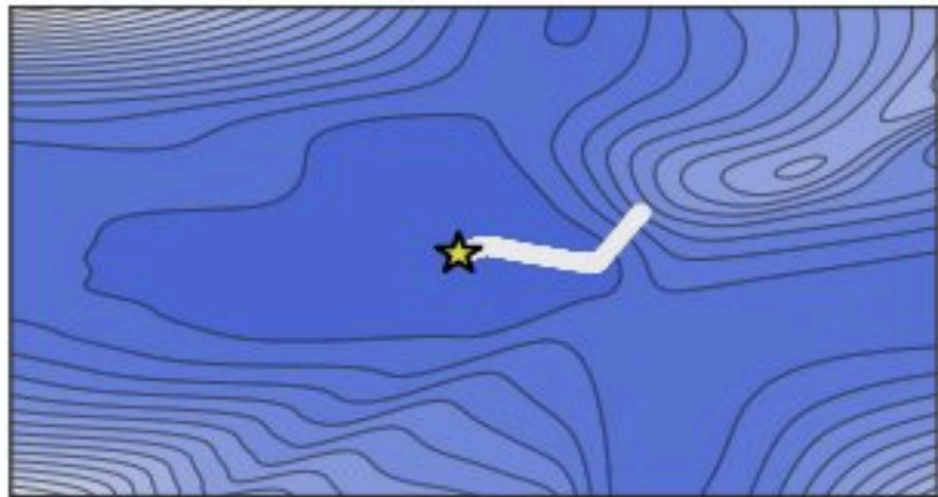## Laplace Approximation

$$p\left(w \mid y, x\right)$$

$$\approx N\left(\hat{w}_{\text{MAP}}, \bar{H}^{-1}(\hat{w}_{\text{MAP}})\right)$$

$$\bar{H}(w) = -\sum_{n=1}^{N} \frac{\partial^2 \log p(y_n, w \mid x_n)}{\partial w^2}$$
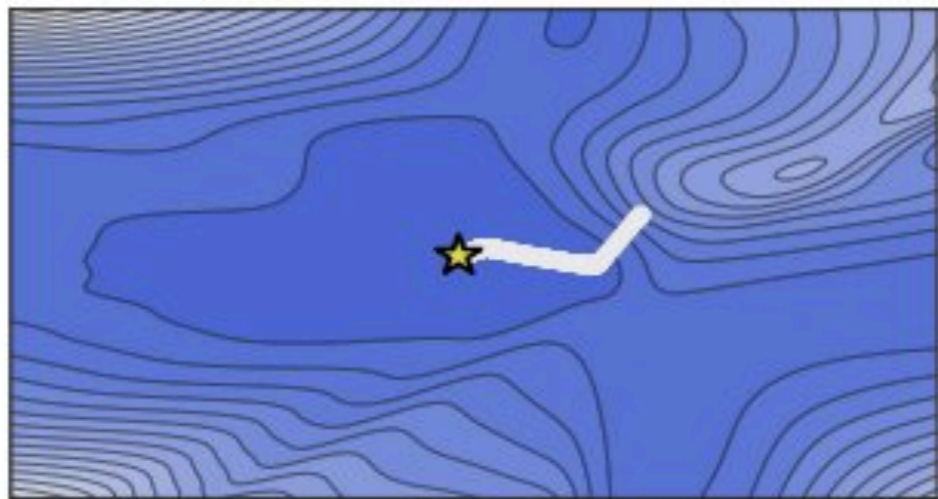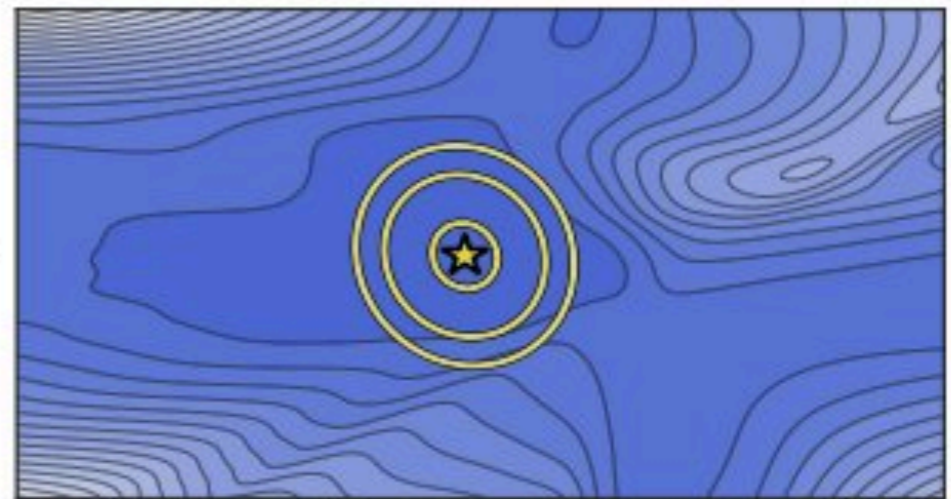
# Laplace Approximation



find $\hat{w}_{MAP}$

# Laplace Approximation



find $\hat{w}_{MAP}$

$N\left(\hat{w}_{MAP},\ H^{-1}(\hat{w}_{MAP})\right)$

# Laplace Approximation

⊗ Pro: can apply to a pre-trained model by assuming parameters are at the 'MAP'

⊗ Con: Hessian matrix can be numerically unstable, need to assume structure (e.g. low-rank, diagonal).

# Summary

⊗ **Conjugacy** for last layer (sometimes)

⊗ **MCMC** is possible but will require approximations

⊗ **Variational inference** is practical but usually has inferior performance (compared to MCMC).