Bayesian Deep Learning with 10 % of the weights

Rob Romijnders

Introduction
Motivation

Method
The goal
Historical perspective
Bayesian inference
Parameter posterior
Uncertainty
Pruning

Experiments and results
Pruning
Uncertainties

Closing

# Bayesian Deep Learning with 10 % of the weights

## Practical approach to Bayesian deep learning

Rob Romijnders

robromijnders.github.io

PyData Amsterdam, 2018

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation

Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# Outline

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation

Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# Outline

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation

Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# Problems with neural networks

Neural networks have three problems:

1. Neural networks give no **uncertainty** in predictions
   $\rightarrow$ easily fooled by **adversarial examples**
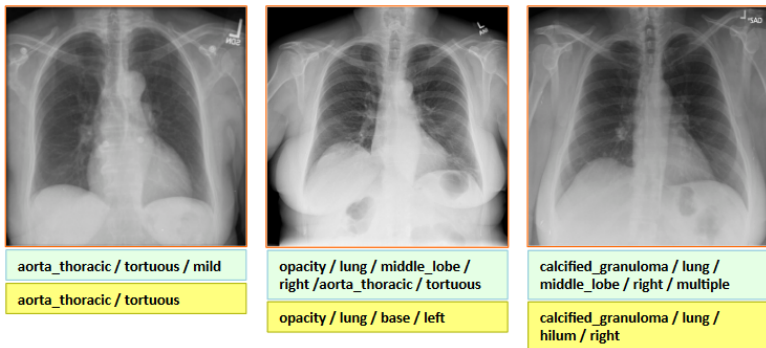2. Neural networks have **millions of parameters**

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

# Motivation



Figure: Uncertainty is important when making diagnoses

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation

Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# Motivation



Figure: Uncertainty is important when making a critical decision

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

# Motivation



Figure: Uncertainty is important when prediction bitcoin

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
**Motivation**

Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# Adversarial attack



Figure: Uncertainty is necessary to find adversarial examples

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation

Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# Embedded applications



Figure: Pruning reduces the memory and computation usage
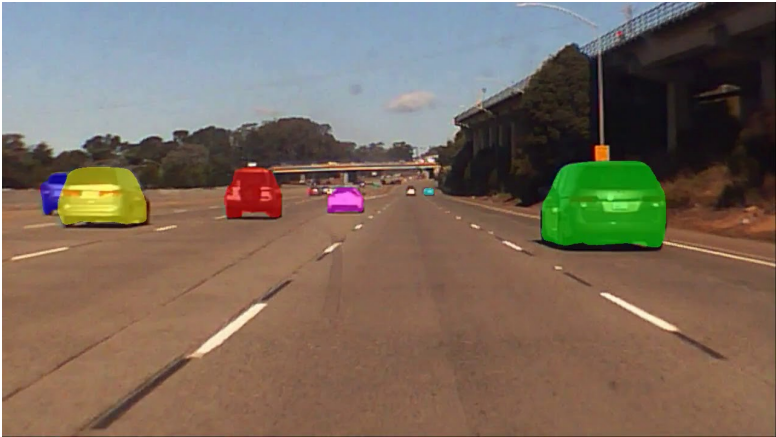(Pruning = dropping parameters)

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation

Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# Real time inference



Figure: Pruning reduces the computation requirements

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

# Outline

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation
Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning
Experiments
and results
Pruning
Uncertainties
Closing

# Pseudo code

In summary, this talk covers the following pseudo code

```
model = Model()

model.train(data)
model.prune()

# Actually, the next line is all we care about:
prediction, uncertainty = model.predict(input)
```

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

How to make a prediction?

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Outline

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation
Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning
Experiments
and results
Pruning
Uncertainties
Closing

# Historical perspective

This content is not new: it has been around for decennia/centuries

### Being Bayesian about neural networks

- Bayes lived in 18th century
- Variational inference for neural networks: Hinton and van Camp (1993)
- Bayesian Neural networks: Neal (1995)

### Uncertainties for a model

Shannon published information theory in 1948

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Outline

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation

Method
The goal
Historical
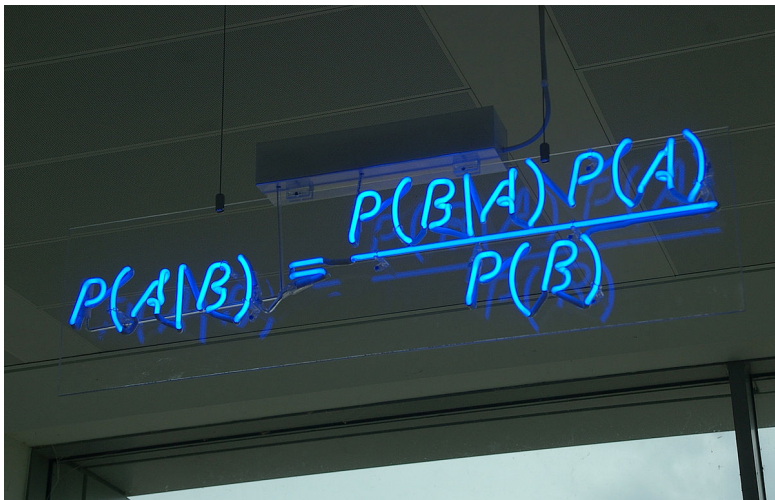perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# Bayes rule



Figure: Every presentation on Bayesian machine learning has this image. So this presentation too

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

# Bayes rule

$$posterior \propto likelihood \times prior$$

$$p(w|data) \propto p(data|w)p(w)$$

$$logp(w|data) = logp(data|w) + logp(w) + constant$$

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

# Bayes rule

We have been using Bayes' rule all the time

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

# Weight decay .. L2 regularisation

$$-log\ posterior = -log\ likelihood \quad -log\ prior \quad +constant$$
$$loss = classification\ loss \quad +\lambda \sum_i w_i^2 \quad +constant$$

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation

Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# Stochastic gradient descent



Figure: But we inferred only one parameter vector

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation
Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# From this...



Figure: Used with kind permission of Eric Ma

# ...to this

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation
Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

Figure: Used with kind permission of Eric Ma

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation
Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning
Experiments
and results
Pruning
Uncertainties
Closing

# Parameters of a Gaussian



Figure: For a Gaussian, we need parameters $\mu$ and $\sigma$

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Outline

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation
Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning
Experiments
and results
Pruning
Uncertainties
Closing

# Posterior probability

The parameter posterior will:

- Enable more samples for prediction $\rightarrow$ uncertainty over prediction
- Tell us which parameters have high zero-probability $\rightarrow$ pruning

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation
Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning
Experiments
and results
Pruning
Uncertainties
Closing

# Loss functions

## Loss functions

*old loss*

$$= \textit{classification loss} + \sum_i \underbrace{\lambda w_i^2}_{\text{L2 penalty}} + \textit{constant}$$

*new loss*

$$= \textit{classification loss} + \sum_i \underbrace{\frac{1}{2}\lambda \mu_i^2}_{\text{L2 penalty}} \underbrace{- \log \sigma_i + \frac{1}{2}\lambda \sigma_i^2}_{\text{penalty on } \sigma} + \textit{constant}$$

## Interpretation

- L2 penalty on the parameter remains

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation

Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# Intuition

$$loss = \underbrace{classification\ loss + \sum_i \frac{1}{2}\lambda\mu_i^2}_{\text{loss on location of weights}} \underbrace{- \log \sigma_i + \frac{1}{2}\lambda\sigma_i^2}_{\text{loss on }\sigma} + constant$$

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation

Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# Summary

- **What do we care about?**
  Uncertainties and pruning

- **How we do it?**
  Find many parameter vectors and average

- **How we do that?**
  Bayesian inference

- **How we do that?**
  Approximate the parameter posterior

- **What do we do in the end?**
  Minimize the loss function on the previous slide

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation
Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
**Uncertainty**
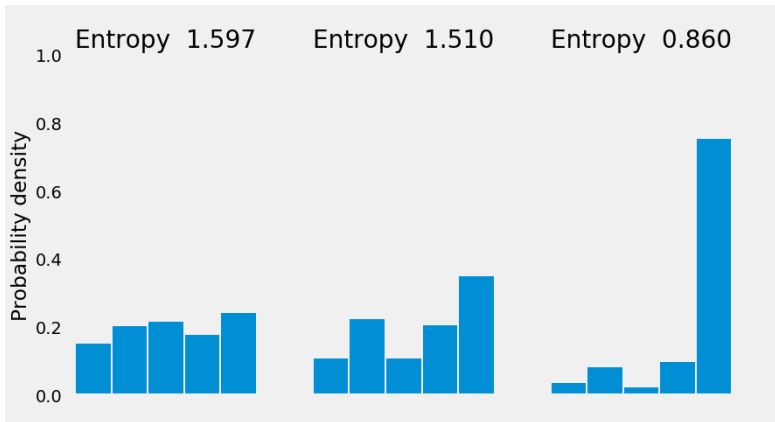Pruning
Experiments
and results
Pruning
Uncertainties

Closing

# Outline

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

# Use entropy as uncertainty metric



Figure: Which prediction has least uncertainty?

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

# Use entropy as uncertainty metric



Figure: Which prediction has least uncertainty?

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation
Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning
Experiments
and results
Pruning
Uncertainties

Closing

# Outline

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation

Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# From this...



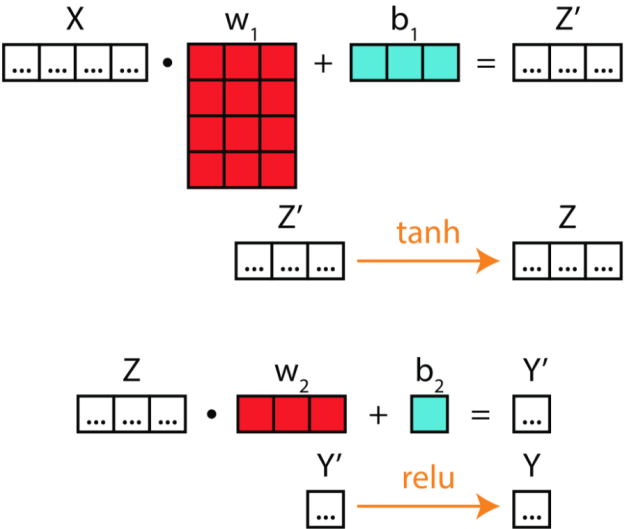Figure: Used with kind permission of Eric Ma

Bayesian Deep
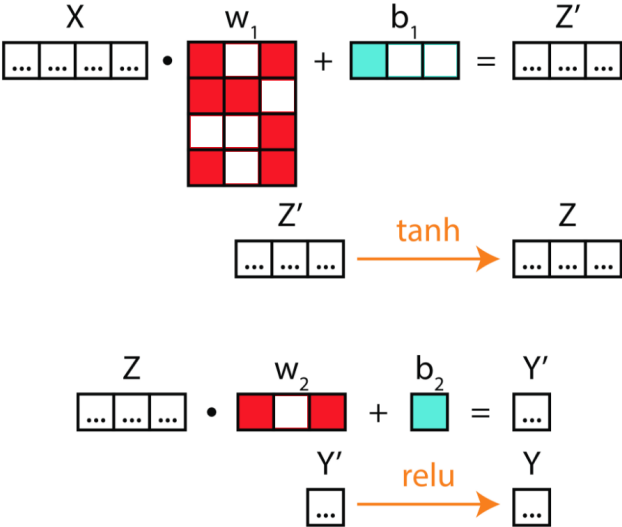Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation
Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# From this...



Figure: Used with kind permission of Eric Ma

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation

Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# Pruning according to posterior



Figure: Which parameter would you rather prune?

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation
Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# Data sets

### Fun
No deep learning project is complete without **MNIST**

### Serious
Two most common applications of deep learning:

- Image recognition: **CIFAR10** data set
- Time series classification: **UCR** - **ECG's**
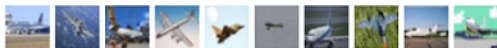  - Train set only 500 time series $\rightarrow$ Bayesian's don't overfit

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation

Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# MNIST examples



Figure: Examples of MNIST. Train set: 50k samples. Test set: 10k samples

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

# CIFAR examples



Figure: Examples of CIFAR. Train set: 50k samples. Test set: 10k samples

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

# ECG examples

??? MAKE ECG EXAMPLES HERE

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation
Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning
Experiments
and results
Pruning
Uncertainties
Closing

# Outline

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation

Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning
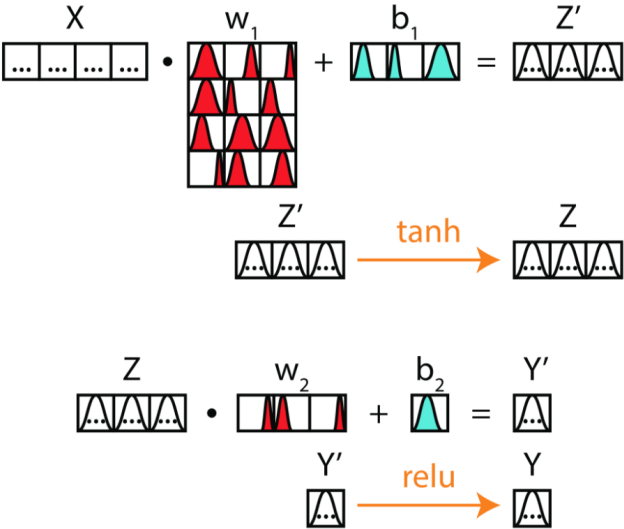
Experiments
and results
Pruning
Uncertainties

Closing

# Remember the model
## ...to this

Bayesian Deep
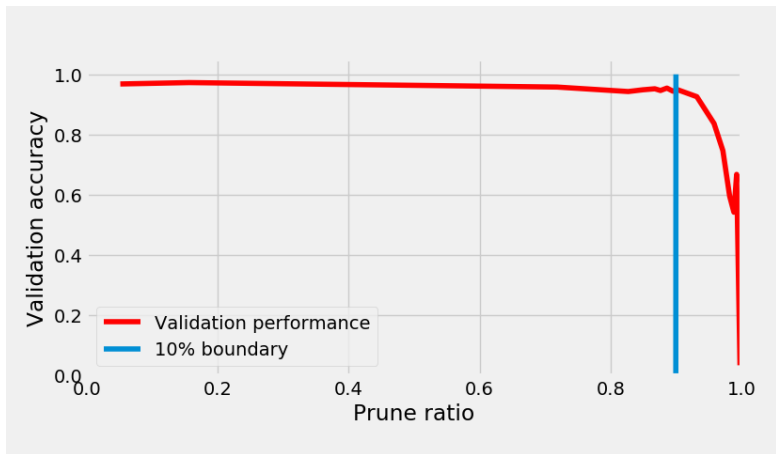Learning with
10 % of the
weights

Rob
Romijnders

# Pruning MNIST



Figure: Pruning curve for MNIST

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

# Pruning CIFAR



Figure: Pruning curve for CIFAR

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation

Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
Uncertainties

Closing

# Pruning ECG



Figure: Pruning curve for ECG

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation
Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning
Experiments
and results
Pruning
**Uncertainties**
Closing

Outline

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation

Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results
Pruning
**Uncertainties**

Closing

# Experiment uncertainty

How to mutilate images to raise uncertainty?

- Add noise
- Warping

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Here will be slides with experiments to show uncertainty on
CIFAR10 when we add noise or rotate or warp

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

### Take aways

- Get uncertainty for critical predictions
- Robust against adversarial attacks
- Prune networks for small memory and small compute

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation
Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning
Experiments
and results
Pruning
Uncertainties
Closing

# Questions?

robromijnders.github.io

## Material
github.com/RobRomijnders/weight_uncertainty

- All code
- Further reading
- More explanation

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

# Additional slides

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

# Learning the sigma's



Figure: The VI objective increases the sigma's by itself!!

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

# Loss on $\sigma$

What does the loss for $\sigma$ look like?



Plot of loss on $\sigma$: $log(\sigma_i) + 1/(2\sigma_i^2)$

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation
Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning
Experiments
and results
Pruning
Uncertainties

Closing

# Make predictions

### Sampling

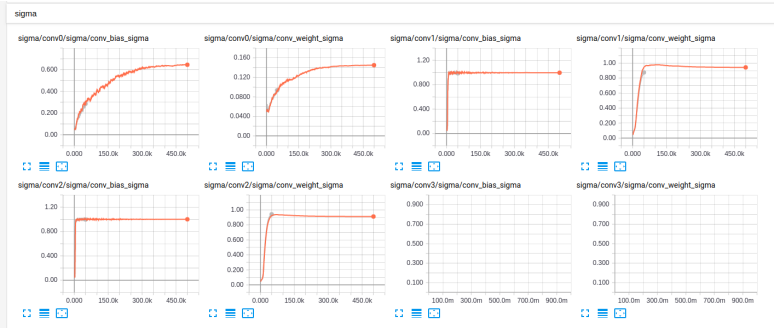Make multiple predictions with sampled parameters. One can think of this sampling as an ensemble method

```
def make_prediction(input):
    for param_vec in param_vecs:
        yield model.get_output(input, param_vec)
prediction = np.mean(make_prediction(input))
```

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation

Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning

Experiments
and results

Pruning
Uncertainties

Closing

# Pseudo code

Pseudo code for training our neural network

```
# OLD CODE
while not converged:
  # Get the loss
  x, y = sample_batch()
  loss = loss_function(x, y, w)

  #Update the parameters
  w_grad = gradient(loss, w)
  w = update(w, w_grad)

##################################################
# NEW CODE
while not converged:
  # Get the loss
  x, y = sample_batch()
  w = approximation.sample()
  loss = loss_function(x, y, w)

  # Update the approximation
  w_grad = gradient(loss, w)
  approximation = update(approximation, w_grad)
```

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

```
while not converged :
    # Get the loss
    x, y = sample_batch()
    w = approximation.sample()
    loss = loss_function(x, y, w)

    # Update the approximation
    w_grad = gradient(loss, w)
    approximation = update(approximation, w_grad)
```

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

# Research

### Pruning: speed

Bayesian compression for deep learning, Louizos @ NIPS2017

### Uncertainty: adversarial attack

Adversarial phenomenon in Bayesian deep learning, Rawat, 2017

Bayesian Deep
Learning with
10 % of the
weights

Rob
Romijnders

Introduction
Motivation
Method
The goal
Historical
perspective
Bayesian
inference
Parameter
posterior
Uncertainty
Pruning
Experiments
and results
Pruning
Uncertainties

Closing

## Gaussian approximation

Approximate with a normal distribution

- Captures local structure of the posterior, which indicates
  the uncertainty
- Simple for parameter pruning

# Anything is better than point estimation !!!