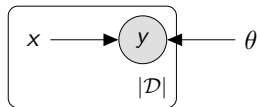# Bayesian Neural Networks
# for Text Classification and Regression

Wilker Aziz
ILLC @ UvA
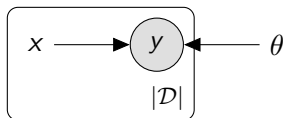
# Probabilistic models parameterised by NNs

The typical text classifier or regressor consists of



- some text input $x$ which we observe but do not model
  i.e. we do not attempt to estimate $p(x)$
- some random, though *observed*, target or response $y$
  e.g. a scalar in regression, a category in classification
- a deterministic mapping from input $x$ and parameters $\theta$
  to the *likelihood $p(y|x, \theta)$*
- mappings realised with the help of *NN architectures*

# Example: Regression



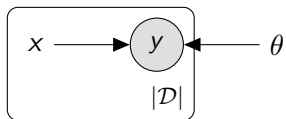$$Y|\theta, x \sim \mathcal{N}(\underbrace{\mu(x;\theta)}_{=u}, \underbrace{\sigma(x;\theta)^2}_{=s})$$

$$p(y|x,\theta) = \frac{1}{\sqrt{2\pi s^2}} \exp\left(-\frac{(x-u)^2}{2s^2}\right)$$

- $\mu(\cdot;\theta)$ maps $x$ to a location in $\mathbb{R}$
- $\sigma(\cdot;\theta)$ maps $x$ to a scale in $\mathbb{R}_{>0}$
- $\theta$ denotes parameters of our NN blocks;

---

affine$_d(\cdot)$ maps its input to $\mathbb{R}^d$ via an affine transformation

## Example: Regression



$$Y|\theta, x \sim \mathcal{N}(\underbrace{\mu(x;\theta)}_{=u}, \underbrace{\sigma(x;\theta)^2}_{=s})$$

$$p(y|x,\theta) = \frac{1}{\sqrt{2\pi s^2}} \exp\left(-\frac{(x-u)^2}{2s^2}\right)$$

- $\mu(\cdot;\theta)$ maps $x$ to a location in $\mathbb{R}$
- $\sigma(\cdot;\theta)$ maps $x$ to a scale in $\mathbb{R}_{>0}$
- $\theta$ denotes parameters of our NN blocks;
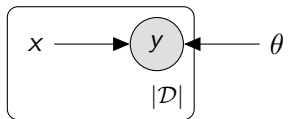
Example architecture:

$$\mathbf{e}_i = \text{emb}(x_i) \qquad \mu(x;\theta) = \text{affine}_1(\mathbf{h})$$
$$\mathbf{h} = \text{rnn}(\mathbf{e}_1^n) \qquad \sigma(x;\theta) = \text{softplus}(\text{affine}_1(\mathbf{h}))$$

---

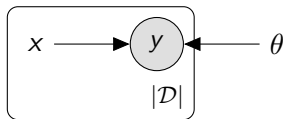affine$_d(\cdot)$ maps its input to $\mathbb{R}^d$ via an affine transformation

# Example: Classification



$$Y|\theta, x \sim \text{Cat}(\pi(x; \theta))$$
$$p(y|x, \theta) = \pi_y(x; \theta)$$

- $\pi(\cdot; \theta)$ maps $x$ to a $K$-dimensional probability vector
- $\theta$ denotes parameters of our NN blocks;

# Example: Classification

$$x \longrightarrow \boxed{y} \longleftarrow \theta$$
$$|\mathcal{D}|$$

$$Y|\theta, x \sim \text{Cat}(\pi(x; \theta))$$
$$p(y|x, \theta) = \pi_y(x; \theta)$$

- $\pi(\cdot; \theta)$ maps $x$ to a $K$-dimensional probability vector
- $\theta$ denotes parameters of our NN blocks;

Example architecture:

$$\mathbf{e}_i = \text{emb}(x_i)$$
$$\mathbf{h} = \text{rnn}(\mathbf{e}_1^n)$$
$$\pi(x; \theta) = \text{softmax}(\text{affine}_K(\mathbf{h}))$$

## Parameter Estimation

Choose $\theta$ that maximises the log-likelihood function

$$\mathcal{L}(\theta|\mathcal{D}) = \sum_{\langle x,y \rangle \in \mathcal{D}} \log p(y|x, \theta)$$

aka *maximum likelihood estimation* (MLE)

## Parameter Estimation

Choose $\theta$ that maximises the log-likelihood function

$$\mathcal{L}(\theta|\mathcal{D}) = \sum_{\langle x,y \rangle \in \mathcal{D}} \log p(y|x, \theta)$$

aka *maximum likelihood estimation* (MLE)

Search for $\theta$ using gradient-based methods
i.e. look for a solution to

$$\boldsymbol{\nabla}_\theta \mathcal{L}(\theta|\mathcal{D}) = \mathbf{0}$$

- backpropagation automates differentiation
- unbiased gradient estimates suffice
  i.e. $\theta^{(t+1)} = \theta^{(t)} + \gamma_t \mathbb{E}_{\mathcal{S} \sim \mathcal{D}} \left[ \boldsymbol{\nabla}_{\theta^{(t)}} \mathcal{L}(\theta^{(t)}|\mathcal{S}) \right]$

# Outline

# Bayes: what and why?

- What is a Bayesian neural net (BNN)?

# Bayes: what and why?

- What is a Bayesian neural net (BNN)?
- Why should we care about them?

# What's a BNN?

BNN



Regression

$$Y|\theta, x \sim \mathcal{N}(\mu(x; \theta), \sigma(x; \theta)^2)$$

Classification

$$Y|\theta, x \sim \text{Cat}(\pi(x; \theta))$$

with for example $\theta \sim \mathcal{N}(\underbrace{0, I}_{\alpha})$

- as before, NNs power the mapping from $x$ and $\theta$ to $p(y|x, \theta)$
- though now $\theta$ is a *random variable*
  distributed according to a prior $p(\theta|\alpha)$

# NNs and BNNs side by side
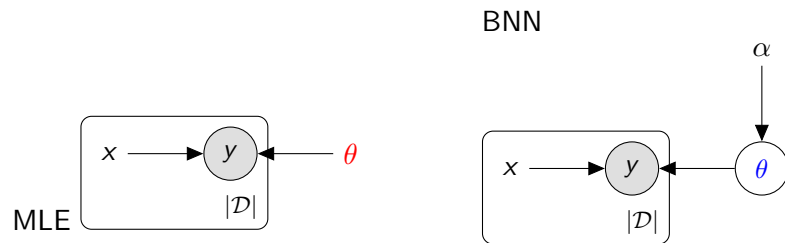


- MLE: assumes $\theta$ to be **given**
- BNN: all variables are treated alike, that is, they are random variables whether or not we call them *parameters*

BNNs also have deterministic parameters (e.g $\alpha$)
we call those *hyperparameters* and they are ideally *fixed*

---

Optimising $\alpha$, say using MLE, makes you an *Empirical Bayesian*, more on that later!

## Bayes

Being *Bayesian* seems to require specifying a prior distribution over parameters, though it's more than that

---

For the theory, including philosophical, mathematical, and statistical perspectives see (Bernardo and Smith, 2009)

## Bayes

Being *Bayesian* seems to require specifying a prior distribution over parameters, though it's more than that

- it's about acknowledging that most quantities are *unknown*

---

For the theory, including philosophical, mathematical, and statistical perspectives see (Bernardo and Smith, 2009)

# Bayes

Being *Bayesian* seems to require specifying a prior distribution over parameters, though it's more than that

- it's about acknowledging that most quantities are *unknown*
- and proceeding to reason probabilistically under uncertainty

---

For the theory, including philosophical, mathematical, and statistical perspectives see (Bernardo and Smith, 2009)

# Bayes

Being *Bayesian* seems to require specifying a prior distribution over parameters, though it's more than that

- it's about acknowledging that most quantities are *unknown*
- and proceeding to reason probabilistically under uncertainty
- priors are a means to this end, they specify what kinds of values are reasonable and with what expectation

---

For the theory, including philosophical, mathematical, and statistical perspectives see (Bernardo and Smith, 2009)

## Bayes

Being *Bayesian* seems to require specifying a prior distribution over parameters, though it's more than that

- it's about acknowledging that most quantities are *unknown*
- and proceeding to reason probabilistically under uncertainty
- priors are a means to this end, they specify what kinds of values are reasonable and with what expectation
- as we will see, acknowledging uncertainty and treating it seriously requires probabilistic inference

---

For the theory, including philosophical, mathematical, and statistical perspectives see (Bernardo and Smith, 2009)

## Bayes

Being *Bayesian* seems to require specifying a prior distribution over parameters, though it's more than that

- it's about acknowledging that most quantities are *unknown*
- and proceeding to reason probabilistically under uncertainty
- priors are a means to this end, they specify what kinds of values are reasonable and with what expectation
- as we will see, acknowledging uncertainty and treating it seriously requires probabilistic inference

Though there are philosophical reasons for adopting the Bayesian paradigm, I concentrate on practical ones.

---

For the theory, including philosophical, mathematical, and statistical perspectives see (Bernardo and Smith, 2009)

# Why Bayes?

How can we quantify the model's uncertainty about a *prediction*?

## Why Bayes?

How can we quantify the model's uncertainty about a *prediction*?

When I say *prediction*, does anyone think of the following?

$$y_* = \arg\max_y \; p(y|x_*, \theta)$$

# Why Bayes?

How can we quantify the model's uncertainty about a *prediction*?

When I say *prediction*, does anyone think of the following?

$$y_* = \arg\max_y \ p(y|x_*, \theta)$$

That's fair, but I meant something even more fundamental.

# Why Bayes?

How can we quantify the model's uncertainty about a *prediction*?

When I say *prediction*, does anyone think of the following?

$$y_* = \arg\max_{y} \ p(y|x_*, \theta)$$

That's fair, but I meant something even more fundamental.

This decision rule (and others) relies on the likelihood $p(y|x_*, \theta)$

# Why Bayes?

How can we quantify the model's uncertainty about a *prediction*?

When I say *prediction*, does anyone think of the following?

$$y_* = \arg\max_{y} \ p(y|x_*, \theta)$$

That's fair, but I meant something even more fundamental.

This decision rule (and others) relies on the likelihood $p(y|x_*, \theta)$

- think of the likelihood as a prediction on its own right
- an NN parameterised by $\theta$ has predicted this value from $x_*$
  (at least implicitly via the parameters of the chosen likelihood)

# Why Bayes?

How can we quantify the model's uncertainty about a *prediction*?

When I say *prediction*, does anyone think of the following?

$$y_* = \arg\max_y \ p(y|x_*, \theta)$$

That's fair, but I meant something even more fundamental.

This decision rule (and others) relies on the likelihood $p(y|x_*, \theta)$

- think of the likelihood as a prediction on its own right
- an NN parameterised by $\theta$ has predicted this value from $x_*$
  (at least implicitly via the parameters of the chosen likelihood)

Then I ask

- how do we know we can trust $p(y|x, \theta)$ for any pair $\langle x, y \rangle$?

# The importance of knowing what we don't know

How do we know we can trust a likelihood assessment $p(y|x, \theta)$?

- suppose we cannot, then neither can we rely on decisions based on likelihood

---

"The importance of knowing what we don't know" (Gal, 2016, Chapter 1).

# The importance of knowing what we don't know

How do we know we can trust a likelihood assessment $p(y|x, \theta)$?

- suppose we cannot, then neither can we rely on decisions based on likelihood

Intuition already suggests that
likelihood does not capture uncertainty

- the likelihood is itself a prediction
- one that's based on a particular instance of the model $(\theta)$
- it tells us that this model supports our data, but does our data support this model?

---

"The importance of knowing what we don't know" (Gal, 2016, Chapter 1).

# The importance of knowing what we don't know

How do we know we can trust a likelihood assessment $p(y|x, \theta)$?

- suppose we cannot, then neither can we rely on decisions based on likelihood

Intuition already suggests that
  likelihood does not capture uncertainty

- the likelihood is itself a prediction
- one that's based on a particular instance of the model $(\theta)$
- it tells us that this model supports our data, but does our data support this model?

If $p(y|x, \theta)$ **is not** about uncertainty, then what is?

---

"The importance of knowing what we don't know" (Gal, 2016, Chapter 1).

# The importance of knowing what we don't know

How do we know we can trust a likelihood assessment $p(y|x, \theta)$?

- suppose we cannot, then neither can we rely on decisions based on likelihood

Intuition already suggests that
likelihood does not capture uncertainty

- the likelihood is itself a prediction
- one that's based on a particular instance of the model $(\theta)$
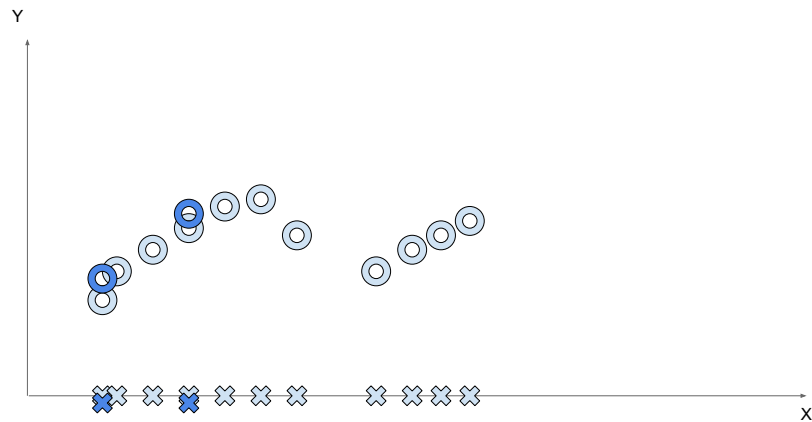- it tells us that this model supports our data, but does our data support this model?

If $p(y|x, \theta)$ **is not** about uncertainty, then what is?
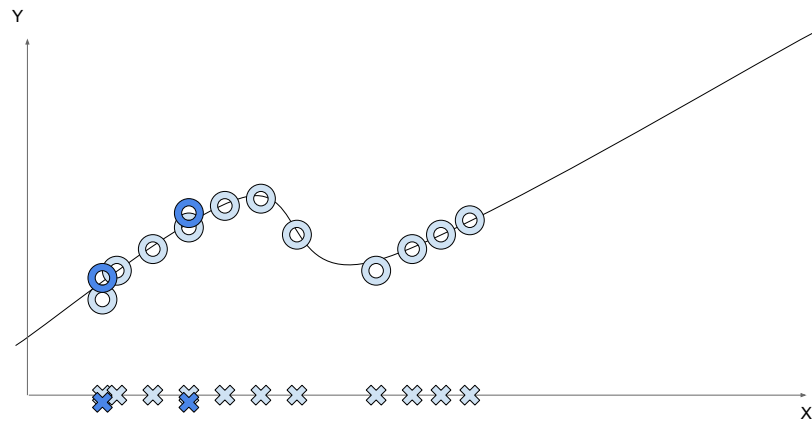The posterior $p(\theta|\mathcal{D})$.

---

"The importance of knowing what we don't know" (Gal, 2016, Chapter 1).

# Uncertainty illustrated



Suppose a regression problem for which we have observations

# Uncertainty illustrated



Let's approach with the help of NNs, i.e. $y = \text{NN}(x; \theta)$

# Uncertainty illustrated



By design, it extrapolates predictions to unseen inputs, e.g. $x_*$

# Uncertainty illustrated



Can we trust our model given $x_*$ is far from observations?

# Uncertainty illustrated



Let's fit Gaussians, i.e. $\mathcal{N}(\mu(x;\theta), \sigma(x;\theta)^2)$, around targets

## Uncertainty illustrated



Note that we never observe much variability for a given input $x$

# Uncertainty illustrated



$\mu(x; \theta)$ learns to be on average close to every response for $x$

# Uncertainty illustrated



$\sigma(x; \theta)$ instead learns to cover all responses for $x$, but no more

# Uncertainty illustrated



for MLE does not like covering more than observed responses

# Uncertainty illustrated



What is our expectation for $\sigma(x_*; \theta)$?

# Uncertainty illustrated



For all we know, $\sigma(\cdot; \theta)$ likes to predict small values

# Uncertainty illustrated



$\sigma(\cdot\,; \theta)$ seriously underestimates uncertainty for $x_*$

# Uncertainty illustrated



but what if, with probability $p(\theta^{(1)}|\mathcal{D})$, we consulted $\mu(x; \theta^{(1)})$?

# Uncertainty illustrated



and $\mu(x; \theta^{(2)})$, with probability $p(\theta^{(2)}|\mathcal{D})$

# Uncertainty illustrated



and $\mu(x; \theta^{(3)})$, with probability $p(\theta^{(3)}|\mathcal{D})$

# Uncertainty illustrated



and $\mu(x; \theta^{(4)})$, with probability $p(\theta^{(4)}|\mathcal{D})$

# Uncertainty illustrated



and $\mu(x; \theta^{(5)})$, with probability $p(\theta^{(5)}|\mathcal{D})$

# Uncertainty illustrated



and $\mu(x; \theta^{(6)})$, with probability $p(\theta^{(6)}|\mathcal{D})$?

# Uncertainty illustrated



Suddenly, we are a lot less certain about predictions for $x_*$

# Oh... This is a Fancy Ensemble?

# Oh... This is a Fancy Ensemble?

Not really, ensembles often *hide uncertainty* rather than expose it!

# Oh... This is a Fancy Ensemble?

Not really, ensembles often *hide uncertainty* rather than expose it!

Ensembles reduce sensitivity of a prediction to initial conditions
via average, i.e. $\hat{\mu}(x_*) = \frac{1}{S} \sum_{i=1}^{S} \mu(x_*; \theta^{(i)})$ for $\theta^{(i)} \sim \text{init}(\theta)$

# Oh... This is a Fancy Ensemble?

Not really, ensembles often *hide uncertainty* rather than expose it!

Ensembles reduce sensitivity of a prediction to initial conditions
via average, i.e. $\hat{\mu}(x_*) = \frac{1}{S} \sum_{i=1}^{S} \mu(x_*; \theta^{(i)})$ for $\theta^{(i)} \sim \text{init}(\theta)$

Note that eventually, it disregards the actual spread of $\mu(x_*; \theta)$



u(x; theta)

# Bayesian Reasoning

Reasoning with *parametric BNNs* involves averaging over parameters

- in a Bayesian sense a *model* is a set of assumptions
  e.g. conditional independences, choice of prior, choice of likelihood,
  architecture blocks, hyperparameters
- formally we should write $p(\mathcal{D}, \theta | \mathcal{M})$ where $\mathcal{M}$ is the model
  assumptions — we omit $\mathcal{M}$ for brevity

# Bayesian Reasoning

Reasoning with *parametric BNNs* involves averaging over parameters

- in a Bayesian sense a *model* is a set of assumptions
  e.g. conditional independences, choice of prior, choice of likelihood,
  architecture blocks, hyperparameters
- formally we should write $p(\mathcal{D}, \theta | \mathcal{M})$ where $\mathcal{M}$ is the model
  assumptions — we omit $\mathcal{M}$ for brevity
- $\theta$ is only one *hypothesis* under this model
- the "worth" of each $\theta$ is quantified by $p(\theta | \mathcal{D})$
- uncertainty estimates are based on this posterior probability

# Bayesian Reasoning

Reasoning with *parametric BNNs* involves averaging over parameters

- in a Bayesian sense a *model* is a set of assumptions
  e.g. conditional independences, choice of prior, choice of likelihood,
  architecture blocks, hyperparameters
- formally we should write $p(\mathcal{D}, \theta | \mathcal{M})$ where $\mathcal{M}$ is the model
  assumptions — we omit $\mathcal{M}$ for brevity
- $\theta$ is only one *hypothesis* under this model
- the "worth" of each $\theta$ is quantified by $p(\theta | \mathcal{D})$
- uncertainty estimates are based on this posterior probability

We will discuss *nonparametric Bayesian models* later and will see that
uncertainty estimates are based on yet more robust model assumptions.

# Bayesian Inference

*Posterior inference*

# Bayesian Inference

*Posterior inference*



$$p(\theta|\mathcal{D}) = \frac{\overbrace{p(\theta)}^{\text{prior}}\overbrace{p(\mathcal{D}|\theta)}^{\text{likelihood}}}{\underbrace{p(\mathcal{D})}_{\text{evidence}}}$$

# Bayesian Inference

*Posterior inference*



$$p(\theta|\mathcal{D}) = \frac{\overbrace{p(\theta)}^{\text{prior}}\overbrace{p(\mathcal{D}|\theta)}^{\text{likelihood}}}{\underbrace{p(\mathcal{D})}_{\text{evidence}}}$$

$$p(\mathcal{D}|\theta) = \prod_{\langle x,y \rangle \in \mathcal{D}} p(y|x,\theta)$$

# Bayesian Inference

*Posterior inference*



$$p(\theta|\mathcal{D}) = \frac{\overbrace{p(\theta)}^{\text{prior}}\overbrace{p(\mathcal{D}|\theta)}^{\text{likelihood}}}{\underbrace{p(\mathcal{D})}_{\text{evidence}}}$$

$$p(\mathcal{D}|\theta) = \prod_{\langle x,y \rangle \in \mathcal{D}} p(y|x,\theta)$$

$$p(\mathcal{D}) = \int p(\theta)p(\mathcal{D}|\theta)\mathrm{d}\theta$$

# Bayesian Inference

*Posterior inference*



$$p(\theta|\mathcal{D}) = \frac{\overbrace{p(\theta)}^{\text{prior}} \overbrace{p(\mathcal{D}|\theta)}^{\text{likelihood}}}{\underbrace{p(\mathcal{D})}_{\text{evidence}}}$$

$$p(\mathcal{D}|\theta) = \prod_{\langle x,y \rangle \in \mathcal{D}} p(y|x,\theta)$$

$$p(\mathcal{D}) = \int p(\theta)p(\mathcal{D}|\theta)\mathrm{d}\theta$$

*Posterior predictive distribution*



$$p(y_*|x_*, \mathcal{D}) =$$

# Bayesian Inference

*Posterior inference*

*Posterior predictive distribution*



$$p(\theta|\mathcal{D}) = \frac{\overbrace{p(\theta)}^{\text{prior}}\overbrace{p(\mathcal{D}|\theta)}^{\text{likelihood}}}{\underbrace{p(\mathcal{D})}_{\text{evidence}}}$$

$$p(\mathcal{D}|\theta) = \prod_{\langle x,y\rangle \in \mathcal{D}} p(y|x,\theta)$$

$$p(\mathcal{D}) = \int p(\theta)p(\mathcal{D}|\theta)\mathrm{d}\theta$$

$$p(y_*|x_*,\mathcal{D}) = \int p(y_*,\theta|\mathcal{D},x_*)\mathrm{d}\theta$$

# Bayesian Inference

*Posterior inference*



*Posterior predictive distribution*



$$p(\theta|\mathcal{D}) = \frac{\overbrace{p(\theta)}^{\text{prior}}\overbrace{p(\mathcal{D}|\theta)}^{\text{likelihood}}}{\underbrace{p(\mathcal{D})}_{\text{evidence}}}$$

$$p(\mathcal{D}|\theta) = \prod_{\langle x,y \rangle \in \mathcal{D}} p(y|x,\theta)$$

$$p(\mathcal{D}) = \int p(\theta)p(\mathcal{D}|\theta)\mathrm{d}\theta$$

$$p(y_*|x_*,\mathcal{D}) = \int p(y_*,\theta|\mathcal{D},x_*)\mathrm{d}\theta$$

$$= \int p(\theta|\mathcal{D})p(y_*|x_*,\theta)\mathrm{d}\theta$$

Conditional independence
$$Y_* \perp \mathcal{D} \mid \theta$$

# Terminology

Prior $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $p(\theta)$

Likelihood $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $p(\mathcal{D}|\theta)$

Posterior $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $p(\theta|\mathcal{D})$

Evidence (aka Marginal Likelihood) $\qquad\qquad\qquad\qquad\qquad$ $p(\mathcal{D})$

Posterior predictive distribution $\qquad\qquad\qquad\qquad$ $p(y_*|x_*, \mathcal{D})$

Probabilistic inference $\qquad\qquad\qquad$ marginalisation/expectation
(remark: in DL the word *inference* is used differently)

## Summary

BNNs are NNs with priors over parameters

The goal is to take uncertainty seriously

Uncertainty estimates help make decisions, e.g.
- model comparison and selection
- when a human should intervene

Other uses include

- reinforcement learning
- active learning
- meta-learning
- learn from streaming data

Bayesian reasoning requires probabilistic inference

# Literature

BDA3                                        Gelman et al. (2013, Chapter 1)

# Outline

# How does one choose a prior?

A prior is meant to capture our beliefs about the phenomenon we are modelling – in this case the relationship between $x$ and $y$

# How does one choose a prior?

A prior is meant to capture our beliefs about the phenomenon we are modelling – in this case the relationship between $x$ and $y$

Let's first consider a simple example: mixture model

$$Z|\pi \sim \text{Cat}(\pi)$$
$$X|\theta, z \sim \text{Cat}(\theta^{(z)})$$

*We first select a discrete mixture component z, this component then selects a Categorical distribution from which we generate a data point x*

# How does one choose a prior?

A prior is meant to capture our beliefs about the phenomenon we are modelling – in this case the relationship between $x$ and $y$

Let's first consider a simple example: mixture model

$$Z|\pi \sim \mathsf{Cat}(\pi)$$
$$X|\theta, z \sim \mathsf{Cat}(\theta^{(z)})$$

*We first select a discrete mixture component z, this component then selects a Categorical distribution from which we generate a data point x*

MLE

$$\pi = {}^1/\kappa \mathbf{1}_K$$
$$\theta = \langle \theta^{(1)}, \dots, \theta^{(K)} \rangle$$

# How does one choose a prior?

A prior is meant to capture our beliefs about the phenomenon we are modelling – in this case the relationship between $x$ and $y$

Let's first consider a simple example: mixture model

$$Z|\pi \sim \text{Cat}(\pi)$$
$$X|\theta, z \sim \text{Cat}(\theta^{(z)})$$

*We first select a discrete mixture component z, this component then selects a Categorical distribution from which we generate a data point x*

MLE

$$\pi = {}^1\!/\kappa \mathbf{1}_K$$
$$\theta = \langle \theta^{(1)}, \dots, \theta^{(K)} \rangle$$

Bayes

$$\pi|\alpha \sim \text{Dir}(\alpha \mathbf{1}_K)$$
$$\theta^{(k)}|\beta \sim \text{Dir}(\beta \mathbf{1}_V)$$

# What makes good mixing coefficients?

Say we have $K = 4$ components, I show a few samples for
$\pi \sim \text{Dir}(10 \times \mathbf{1}_K)$

# What makes good mixing coefficients?

Say we have $K = 4$ components, I show a few samples for
$\pi \sim \text{Dir}(10 \times \mathbf{1}_K)$



$\pi \sim \text{Dir}(1 \times \mathbf{1}_K)$

# What makes good mixing coefficients?

Say we have $K = 4$ components, I show a few samples for
$\pi \sim \text{Dir}(10 \times \mathbf{1}_K)$



$\pi \sim \text{Dir}(1 \times \mathbf{1}_K)$



$\pi \sim \text{Dir}(0.1 \times \mathbf{1}_K)$

# What makes good mixing coefficients?

Say we have $K = 4$ components, I show a few samples for
$\pi \sim \text{Dir}(10 \times \mathbf{1}_K)$



$\pi \sim \text{Dir}(1 \times \mathbf{1}_K)$



$\pi \sim \text{Dir}(0.1 \times \mathbf{1}_K)$



Can you make any assumptions before observing data?

# What makes a good conditional?

Say we have $V = 10$ types of data points, I show samples $\theta^{(k)}|\beta \sim \text{Dir}(\beta)$

# Mixture Models are Simple to Understand

The unobservable random variables $\pi$ and $\theta^{(k)}$ are rather interpretable

- it's clear that we want assignments to be unambiguous
  sparse mixing weights
- it's clear that we want components to be rather selective
  sparse conditionals
- it's clear that we don't know the identify of clusters
  uniform marginals

All of that is essentially very clear a priori

- that is, before we collect observations
- by simply considering the nature of problem

# Mixture Models are Simple to Understand

The unobservable random variables $\pi$ and $\theta^{(k)}$ are rather interpretable

- it's clear that we want assignments to be unambiguous
  sparse mixing weights
- it's clear that we want components to be rather selective
  sparse conditionals
- it's clear that we don't know the identify of clusters
  uniform marginals

All of that is essentially very clear a priori

- that is, before we collect observations
- by simply considering the nature of problem

Meaning of weights in NNs are quite obscure!
  Who can tell what aspect of a classifier any of the LSTM parameters
controls?

## Learning Functions

We are essentially using NNs to learn some unknown function that maps from data to probabilities — which then support decisions

# Learning Functions

We are essentially using NNs to learn some unknown function that maps from data to probabilities — which then support decisions

It's only natural to ask, *What functions can NNs learn*?

# Learning Functions

We are essentially using NNs to learn some unknown function that maps from data to probabilities — which then support decisions

It's only natural to ask, *What functions can NNs learn*?

- usually one says "any continuous function" (Funahashi, 1989) given a wide enough hidden layer

# Learning Functions

We are essentially using NNs to learn some unknown function that maps from data to probabilities — which then support decisions

It's only natural to ask, *What functions can NNs learn*?

- usually one says "any continuous function" (Funahashi, 1989) given a wide enough hidden layer

Fair, but how about this,

- *What functions can we actually recover given finite observations?*
- *How about the fact that we employ convex optimisers?*

# Learning Functions

We are essentially using NNs to learn some unknown function that maps from data to probabilities — which then support decisions

It's only natural to ask, *What functions can NNs learn*?

- usually one says "any continuous function" (Funahashi, 1989) given a wide enough hidden layer

Fair, but how about this,

- *What functions can we actually recover given finite observations?*
- *How about the fact that we employ convex optimisers?*

It's hard to talk about what functions we can learn when the most important factors are *amount of data* and the *success of a local optimiser*

# Random functions

Let's consider what happens when our parameters are random following a given prior.

## Random functions

Let's consider what happens when our parameters are random following a given prior.

Sampling from these priors and performing forward passes with the network will expose a range of functions

- some might have specific properties
  e.g. smoothness, periodicity
- some will be preferred over others
- some may be impossible
  e.g. Brownian functions vs infinitely smooth functions

# Draws



Example from MacKay (1998)

# Priors

Priors are not about making things random for no reason
  they are about encoding assumptions, or inductive biases!

Let's turn to known priors over functions!

# Gaussian Processes

Consider the case of regression, where $y = f(x) + \epsilon$ for some
$\epsilon \sim \mathcal{N}(0, \tau^{-1})$

- this implies $Y|f(x) \sim \mathcal{N}(f(x), \tau^{-1})$
- let's design a prior for $f(x)$
  Note that a parametric way to do so is to say $f(x) = w^\top \phi(x)$ for
  some fixed feature function $\phi(x)$ and impose a prior on $w$, but then
  again, what are the properties of such a prior?

## Gaussian Processes

Consider the case of regression, where $y = f(x) + \epsilon$ for some
$\epsilon \sim \mathcal{N}(0, \tau^{-1})$

- this implies $Y|f(x) \sim \mathcal{N}(f(x), \tau^{-1})$
- let's design a prior for $f(x)$
  Note that a parametric way to do so is to say $f(x) = w^\top \phi(x)$ for
  some fixed feature function $\phi(x)$ and impose a prior on $w$, but then
  again, what are the properties of such a prior?

The probability distribution of a function $f(x)$ is a Gaussian process (GP)
if for any finite selection of points $x^{(1)}, \ldots, x^{(N)}$ the density
$p(f(x^{(1)}), \ldots, f(x^{(N)}))$ is a Gaussian.

A function represents an infinite object, but in ML we typically only reason
over finite datasets!

# GP prior

I'll employ boldfacing to denote a collection of $N$ datapoints, e.g. $\mathbf{x} = \{x^{(1)}, \ldots, x^{(N)}\}$ and $\mathbf{y} = \{y^{(1)}, \ldots, y^{(N)}\}$, indexing returns an element, e.g. $x_i \stackrel{\text{def}}{=} x^{(i)}$.

Check the *excellent* Kernel Cookbook by David Duvenaud

# GP prior

I'll employ boldfacing to denote a collection of $N$ datapoints, e.g. $\mathbf{x} = \{x^{(1)}, \ldots, x^{(N)}\}$ and $\mathbf{y} = \{y^{(1)}, \ldots, y^{(N)}\}$, indexing returns an element, e.g. $x_i \overset{\text{def}}{=} x^{(i)}$. Similarly, $\mathbf{f} = \{f(x^{(1)}), \ldots, f(x^{(N)})\}$ is a collection of latent function assessments.

---

Check the *excellent* Kernel Cookbook by David Duvenaud

# GP prior

I'll employ boldfacing to denote a collection of $N$ datapoints, e.g. $\mathbf{x} = \{x^{(1)}, \ldots, x^{(N)}\}$ and $\mathbf{y} = \{y^{(1)}, \ldots, y^{(N)}\}$, indexing returns an element, e.g. $x_i \stackrel{\text{def}}{=} x^{(i)}$. Similarly, $\mathbf{f} = \{f(x^{(1)}), \ldots, f(x^{(N)})\}$ is a collection of latent function assessments.

GP prior

$$\mathbf{F}|\mathbf{x} \sim \mathcal{N}(\mathbf{0}, k(\mathbf{x}, \mathbf{x}))$$

$$\mathbf{Y}|\mathbf{f} \sim \mathcal{N}(\mathbf{f}, \tau^{-1}\mathbf{I}_N)$$



The covariance matrix is defined by a kernel function $k(x, x')$

- I abuse notation and use $k(\mathbf{x}, \mathbf{x})$ to denote the $N \times N$ matrix $\mathbf{K}$ of kernel assessments, i.e. $K_{i,j} = k(x_i, x_j)$
- $k(x', \mathbf{x})$ denotes a row-vector of kernel assessments

---

Check the *excellent* Kernel Cookbook by David Duvenaud

# GP Inference

Given a collection $\mathbf{y}, \mathbf{f}|\mathbf{x}$ of jointly Gaussian variables, what's the family of the marginal?

$$p(\mathbf{y}|\mathbf{x}) = \int p(\mathbf{y}, \mathbf{f}|\mathbf{x})\mathrm{d}\mathbf{f} =$$

Bishop (2006, Chapter 2) for operations with Multivariate Gaussians

## GP Inference

Given a collection $\mathbf{y}, \mathbf{f}|\mathbf{x}$ of jointly Gaussian variables, what's the family of the marginal?

$$p(\mathbf{y}|\mathbf{x}) = \int p(\mathbf{y}, \mathbf{f}|\mathbf{x})\mathrm{d}\mathbf{f} = \int \underbrace{\mathcal{N}(\mathbf{f}|\mathbf{0}, k(\mathbf{x}, \mathbf{x}))\mathcal{N}(\mathbf{y}|\mathbf{f}, \tau^{-1}\mathbf{I}_N)}_{\text{jointly Gaussian}}\,\mathrm{d}\mathbf{f}$$

<span style="color:red">Recall: marginals of a multivariate Gaussian are Gaussians!</span>

Bishop (2006, Chapter 2) for operations with Multivariate Gaussians

## GP Inference

Given a collection $\mathbf{y}, \mathbf{f}|\mathbf{x}$ of jointly Gaussian variables, what's the family of the marginal?

$$p(\mathbf{y}|\mathbf{x}) = \int p(\mathbf{y}, \mathbf{f}|\mathbf{x})\mathrm{d}\mathbf{f} = \int \underbrace{\mathcal{N}(\mathbf{f}|\mathbf{0}, k(\mathbf{x}, \mathbf{x}))\mathcal{N}(\mathbf{y}|\mathbf{f}, \tau^{-1}\mathbf{I}_N)}_{\text{jointly Gaussian}}\mathrm{d}\mathbf{f}$$

Recall: marginals of a multivariate Gaussian are Gaussians!

Thus what's the family of the posterior?

$$p(\mathbf{f}|\mathbf{x}, \mathbf{y}) = \frac{p(\mathbf{y}, \mathbf{f}|\mathbf{x})}{p(\mathbf{y}|\mathbf{x})}$$

---

Bishop (2006, Chapter 2) for operations with Multivariate Gaussians

## GP Inference

Given a collection $\mathbf{y}, \mathbf{f}|\mathbf{x}$ of jointly Gaussian variables, what's the family of the marginal?

$$p(\mathbf{y}|\mathbf{x}) = \int p(\mathbf{y}, \mathbf{f}|\mathbf{x})\mathrm{d}\mathbf{f} = \int \underbrace{\mathcal{N}(\mathbf{f}|\mathbf{0}, k(\mathbf{x}, \mathbf{x}))\mathcal{N}(\mathbf{y}|\mathbf{f}, \tau^{-1}\mathbf{I}_N)}_{\text{jointly Gaussian}}\,\mathrm{d}\mathbf{f}$$

Recall: marginals of a multivariate Gaussian are Gaussians!

Thus what's the family of the posterior?

$$p(\mathbf{f}|\mathbf{x}, \mathbf{y}) = \frac{p(\mathbf{y}, \mathbf{f}|\mathbf{x})}{p(\mathbf{y}|\mathbf{x})}$$

Recall: conditioning on a subset of a multivariate Gaussian yields a multivariate Guassian

---

Bishop (2006, Chapter 2) for operations with Multivariate Gaussians

## Exact Inference with GPs

Posterior

$$\mathbf{F}|\mathbf{x}, \mathbf{y} \sim \mathcal{N}(\mathbf{m}_{\text{post}}, \mathbf{K}_{\text{post}})$$
$$\mathbf{m}_{\text{post}} = \mathbf{K}(\mathbf{K} + \tau^{-1}\mathbf{I}_N)^{-1}\mathbf{y}$$
$$\mathbf{K}_{\text{post}} = \mathbf{K} - \mathbf{K}(\mathbf{K} + \tau^{-1}\mathbf{I}_N)^{-1}\mathbf{K}^\top$$

Posterior predictive distribution:

$$Y_*|x_*, \mathbf{x}, \mathbf{y} \sim \mathcal{N}(k(x_*, \mathbf{x})(\mathbf{K} + \tau^{-1}\mathbf{I}_N)^{-1}\mathbf{y},$$
$$k(x_*, x_*) + \tau^{-1} - k(x_*, \mathbf{x})(\mathbf{K} + \tau^{-1}\mathbf{I}_N)^{-1}k(x_*, \mathbf{x})^\top)$$

# Uncertainty illustrated (revisited)



Let's get back to this

# Uncertainty illustrated (revisited)



What if uncertainty depended on the distance to observations?

# Uncertainty illustrated (revisited)



Kernels in GPs operationalise exactly this notion

## Terminology

Random functions: latent treatment to $f(x)$

Kernel: $k : \mathcal{X} \times X \to \mathbb{R}$ such that $k(x, x')$ is the covariance between $f(x)$ and $f(x')$

Gaussian process prior: $F \sim \mathcal{G}P(0, k)$

GP inference: marginals and conditionals are Gaussians

## Summary

Priors are as good as our understanding of what class of models they favour

NNs are meant to learn unknown functions

BNNs learn a distribution over functions by treating parameters as random variables

The effect of a parameter over the learned function is unclear

A prior over functions can be specified in a non-parametric way via specification of a covariance (kernel) function

A GP prior is a well-studied prior over functions

# Literature

David MacKay's pioneering work
  Bayesian interpolation                              MacKay (1992a)
  or go all the way through his PhD thesis       MacKay (1992b)

Priors for Infinite Networks                    Neal (1994, 1996)

Multivariate Gaussians               Bishop (2006, Chapter 2)

Introduction to GPs                      MacKay (1998)

GP summer school classes by Neil Laurence

Kernel Cookbook by David Duvenaud

---

Neal (1994): https://www.cs.toronto.edu/~radford/ftp/pin.pdf

# Outline

# GPs vs BNNs

GP

BNN



GP's a non-parametric models

- the complexity (or capacity) of the model grows with the data
- posterior predictive is known and tractable
- we know a lot about the random functions we get

BNNs are parametric models

- the complexity (or capacity) is pre-specified
- posterior predictive is unknown and intractable
- we know little about the random functions we get

# Why don't we always use GPs then?

Flexibility

- kernels for text are fewer, less convenient, and less well-understood
- $x$ can be very high-dimensional (and perhaps we have less intuitions to choose a kernel)

Computational complexity

- exact GP inference takes $O(N^3)$
- it's possible to scale them up, but that's an active research topic
- many solutions are specific to continuous inputs

# Let's then consider Bayesian inference for a BNN

This is essentially what we have to address

## Let's then consider Bayesian inference for a BNN

This is essentially what we have to address



That is,

$$p(y_*, \theta | \mathcal{D}, x_*) = \int p(y_*, \theta | \mathcal{D}) \mathrm{d}\theta = \int p(\theta | \mathcal{D}) p(y_* | \theta, x_*) \mathrm{d}\theta$$

# Let's then consider Bayesian inference for a BNN

This is essentially what we have to address



That is,

$$p(y_*, \theta | \mathcal{D}, x_*) = \int p(y_*, \theta | \mathcal{D}) \mathrm{d}\theta = \int p(\theta | \mathcal{D}) p(y_* | \theta, x_*) \mathrm{d}\theta$$

But $p(\theta | \mathcal{D}) = \int p(\theta | \alpha) p(\mathcal{D} | \theta) \mathrm{d}\theta$ is intractable

# Let's then consider Bayesian inference for a BNN

This is essentially what we have to address



That is,

$$p(y_*, \theta | \mathcal{D}, x_*) = \int p(y_*, \theta | \mathcal{D}) \mathrm{d}\theta = \int p(\theta | \mathcal{D}) p(y_* | \theta, x_*) \mathrm{d}\theta$$

But $p(\theta | \mathcal{D}) = \int p(\theta | \alpha) p(\mathcal{D} | \theta) \mathrm{d}\theta$ is intractable
Let's learn a proxy $q(\theta | \lambda)$!

## Variational Inference

Let's learn a proxy $q(\theta|\lambda)$ to $p(\theta|\mathcal{D})$ and solve

$$p(y_*, \theta|\mathcal{D}, x_*) = \int p(\theta|\mathcal{D})p(y_*|\theta, x_*)\mathrm{d}\theta \approx \int q(\theta|\lambda)p(y_*|\theta, x_*)\mathrm{d}\theta$$

---

An alternative with guarantees is MCMC – as discussed in ML2. Example: MCMC for a mixture of Gaussians by David Blei.

# Variational Inference

Let's learn a proxy $q(\theta|\lambda)$ to $p(\theta|\mathcal{D})$ and solve

$$p(y_*, \theta|\mathcal{D}, x_*) = \int p(\theta|\mathcal{D})p(y_*|\theta, x_*)\mathrm{d}\theta \approx \int q(\theta|\lambda)p(y_*|\theta, x_*)\mathrm{d}\theta$$

Principle: choose an approximation that minimises KL-divergence

$$\underset{q(\theta)}{\arg\min} \ \ \mathrm{KL}(q(\theta)||p(\theta|\mathcal{D}))$$

$$= \underset{q(\theta)}{\arg\min} \ \ \mathbb{E}_{q(\theta)}\left[\log \frac{q(\theta|\lambda)}{p(\theta|\mathcal{D})}\right] \quad \text{definition of KL}$$

---

An alternative with guarantees is MCMC – as discussed in ML2. Example: MCMC for a mixture of Gaussians by David Blei.

# Evidence Lowerbound (ELBO)

Principle: choose an approximation that minimises KL-divergence

$$\arg\min_{q(\theta)} \ \mathrm{KL}(q(\theta) \| p(\theta|\mathcal{D}))$$

$$= \arg\min_{q(\theta)} \ \mathbb{E}_{q(\theta)} \left[ \log \frac{q(\theta|\lambda)}{p(\theta|\mathcal{D})} \right] \qquad \text{definition of KL}$$

# Evidence Lowerbound (ELBO)

Principle: choose an approximation that minimises KL-divergence

$$\underset{q(\theta)}{\arg\min} \quad \mathrm{KL}(q(\theta)||p(\theta|\mathcal{D}))$$

$$= \underset{q(\theta)}{\arg\min} \quad \mathbb{E}_{q(\theta)}\left[\log \frac{q(\theta|\lambda)}{p(\theta|\mathcal{D})}\right] \qquad \text{definition of KL}$$

$$= \underset{q(\theta)}{\arg\min} \quad \mathbb{E}_{q(\theta)}\left[\log \frac{q(\theta|\lambda)p(\mathcal{D})}{p(\theta, \mathcal{D})}\right] \qquad \text{definition of posterior}$$

# Evidence Lowerbound (ELBO)

Principle: choose an approximation that minimises KL-divergence

$$\arg\min_{q(\theta)} \; \text{KL}(q(\theta)||p(\theta|\mathcal{D}))$$

$$= \arg\min_{q(\theta)} \; \mathbb{E}_{q(\theta)}\left[\log \frac{q(\theta|\lambda)}{p(\theta|\mathcal{D})}\right] \qquad\qquad \text{definition of KL}$$

$$= \arg\min_{q(\theta)} \; \mathbb{E}_{q(\theta)}\left[\log \frac{q(\theta|\lambda)p(\mathcal{D})}{p(\theta, \mathcal{D})}\right] \qquad\qquad \text{definition of posterior}$$

$$= \arg\min_{q(\theta)} \; \mathbb{E}_{q(\theta)}\left[\log \frac{q(\theta|\lambda)}{p(\theta, \mathcal{D})}\right] + \log p(\mathcal{D}) \qquad\qquad \text{constant}$$

# Evidence Lowerbound (ELBO)

Principle: choose an approximation that minimises KL-divergence

$$\underset{q(\theta)}{\arg\min} \quad \mathrm{KL}(q(\theta)||p(\theta|\mathcal{D}))$$

$$= \underset{q(\theta)}{\arg\min} \quad \mathbb{E}_{q(\theta)}\left[\log \frac{q(\theta|\lambda)}{p(\theta|\mathcal{D})}\right] \qquad\qquad \text{definition of KL}$$

$$= \underset{q(\theta)}{\arg\min} \quad \mathbb{E}_{q(\theta)}\left[\log \frac{q(\theta|\lambda)p(\mathcal{D})}{p(\theta, \mathcal{D})}\right] \qquad \text{definition of posterior}$$

$$= \underset{q(\theta)}{\arg\min} \quad \mathbb{E}_{q(\theta)}\left[\log \frac{q(\theta|\lambda)}{p(\theta, \mathcal{D})}\right] + \log p(\mathcal{D}) \qquad\qquad \text{constant}$$

$$= \underset{q(\theta)}{\arg\min} \quad - \mathbb{E}_{q(\theta)}\left[\log \frac{p(\theta, \mathcal{D})}{q(\theta|\lambda)}\right] \qquad\qquad \text{property of log}$$

# Evidence Lowerbound (ELBO)

Principle: choose an approximation that minimises KL-divergence

$$\arg\min_{q(\theta)} \; \text{KL}(q(\theta)||p(\theta|\mathcal{D}))$$

$$= \arg\min_{q(\theta)} \; \mathbb{E}_{q(\theta)}\left[\log \frac{q(\theta|\lambda)}{p(\theta|\mathcal{D})}\right] \qquad\qquad \text{definition of KL}$$

$$= \arg\min_{q(\theta)} \; \mathbb{E}_{q(\theta)}\left[\log \frac{q(\theta|\lambda)p(\mathcal{D})}{p(\theta, \mathcal{D})}\right] \qquad\qquad \text{definition of posterior}$$

$$= \arg\min_{q(\theta)} \; \mathbb{E}_{q(\theta)}\left[\log \frac{q(\theta|\lambda)}{p(\theta, \mathcal{D})}\right] + \log p(\mathcal{D}) \qquad\qquad \text{constant}$$

$$= \arg\min_{q(\theta)} \; - \mathbb{E}_{q(\theta)}\left[\log \frac{p(\theta, \mathcal{D})}{q(\theta|\lambda)}\right] \qquad\qquad \text{property of log}$$

$$= \arg\max_{q(\theta)} \; \mathbb{E}_{q(\theta)}\left[\log p(\theta, \mathcal{D})\right] + \mathbb{H}(q(\theta)) \qquad\qquad \text{ELBO}$$

# Evidence Lowerbound (ELBO)

Principle: choose an approximation that minimises KL-divergence

$$\underset{q(\theta)}{\arg\min} \ \ \mathrm{KL}(q(\theta)||p(\theta|\mathcal{D}))$$

$$= \underset{q(\theta)}{\arg\min} \ \mathbb{E}_{q(\theta)}\left[\log \frac{q(\theta|\lambda)}{p(\theta|\mathcal{D})}\right] \qquad \text{definition of KL}$$

$$= \underset{q(\theta)}{\arg\min} \ \mathbb{E}_{q(\theta)}\left[\log \frac{q(\theta|\lambda)p(\mathcal{D})}{p(\theta,\mathcal{D})}\right] \qquad \text{definition of posterior}$$

$$= \underset{q(\theta)}{\arg\min} \ \mathbb{E}_{q(\theta)}\left[\log \frac{q(\theta|\lambda)}{p(\theta,\mathcal{D})}\right] + \log p(\mathcal{D}) \qquad \text{constant}$$

$$= \underset{q(\theta)}{\arg\min} \ -\mathbb{E}_{q(\theta)}\left[\log \frac{p(\theta,\mathcal{D})}{q(\theta|\lambda)}\right] \qquad \text{property of log}$$

$$= \underset{q(\theta)}{\arg\max} \ \mathbb{E}_{q(\theta)}[\log p(\theta,\mathcal{D})] + \mathbb{H}(q(\theta)) \qquad \text{ELBO}$$

All quantities are either tractable or easy to estimate by sampling!

# ELBO continued

Parametric assumption

$$\underset{q(\theta)}{\arg\max}\ \mathbb{E}_{q(\theta)}\left[\log p(\theta, \mathcal{D})\right] + \mathbb{H}(q(\theta))$$

$$= \underset{\lambda}{\arg\max}\ \mathbb{E}_{q(\theta|\lambda)}\left[\log p(\theta, \mathcal{D})\right] + \mathbb{H}(q(\theta|\lambda))$$

Recall

$$p(\theta, \mathcal{D}) = p(\theta) \prod_{i=1}^{N} p(y^{(i)}|\theta, x^{(i)})$$

# ELBO continued

Parametric assumption

$$\arg\max_{q(\theta)} \ \mathbb{E}_{q(\theta)}\left[\log p(\theta, \mathcal{D})\right] + \mathbb{H}(q(\theta))$$

$$= \arg\max_{\lambda} \ \mathbb{E}_{q(\theta|\lambda)}\left[\log p(\theta, \mathcal{D})\right] + \mathbb{H}(q(\theta|\lambda))$$

Recall

$$p(\theta, \mathcal{D}) = p(\theta) \prod_{i=1}^{N} p(y^{(i)}|\theta, x^{(i)})$$

And thus the ELBO evaluates to

$$\mathbb{E}_{q(\theta|\lambda)}\left[\log p(\theta) + \sum_{i=1}^{N} \log p(y^{(i)}|\theta, x^{(i)})\right] + \mathbb{H}(q(\theta|\lambda))$$

$$= \mathbb{E}_{q(\theta|\lambda)}\left[\sum_{i=1}^{N} \log p(y^{(i)}|\theta, x^{(i)})\right] - \text{KL}(q(\theta|\lambda)||p(\theta))$$

# Mean Field Assumption

Let $\theta \in \mathbb{R}^D$. The simplest approximate posterior is

---

We can group parameters and assume independence of groups (e.g. layers).

# Mean Field Assumption

Let $\theta \in \mathbb{R}^D$. The simplest approximate posterior is

$$q(\theta|\lambda) = \prod_{d=1}^{D} q(\theta_d|\lambda)$$

where we assume independence amongst $\theta_d$.

---

We can group parameters and assume independence of groups (e.g. layers).

## Mean Field Assumption

Let $\theta \in \mathbb{R}^D$. The simplest approximate posterior is

$$q(\theta|\lambda) = \prod_{d=1}^{D} q(\theta_d|\lambda)$$

where we assume independence amongst $\theta_d$.

If we have the same exponential family for $p(\theta)$ and $q(\theta|\lambda)$,

---

We can group parameters and assume independence of groups (e.g. layers).

# Mean Field Assumption

Let $\theta \in \mathbb{R}^D$. The simplest approximate posterior is

$$q(\theta|\lambda) = \prod_{d=1}^{D} q(\theta_d|\lambda)$$

where we assume independence amongst $\theta_d$.

If we have the same exponential family for $p(\theta)$ and $q(\theta|\lambda)$, then

$$\text{KL}(q(\theta|\lambda)||p(\theta)) = \sum_{d=1}^{D} \underbrace{\text{KL}(q(\theta_d|\lambda)||p(\theta_d))}_{\text{closed form}}$$

is known in closed form.

---

We can group parameters and assume independence of groups (e.g. layers).

# Choosing $\lambda$

How should we choose $\lambda$?

# Choosing $\lambda$

How should we choose $\lambda$?

How can we approach the following problem?

$$\arg\max_{\lambda} \; \mathbb{E}_{q(\theta|\lambda)}\left[\log p(\mathcal{D}|\theta)\right] - \underbrace{\mathsf{KL}(q(\theta|\lambda)||p(\theta))}_{\text{closed form}}$$

# Choosing $\lambda$

How should we choose $\lambda$?

How can we approach the following problem?

$$\arg\max_{\lambda} \; \mathbb{E}_{q(\theta|\lambda)}\left[\log p(\mathcal{D}|\theta)\right] - \underbrace{\mathsf{KL}(q(\theta|\lambda)\|p(\theta))}_{\text{closed form}}$$

What's the workhorse of optimisation in deep learning?

# Gradient-based optimisation for $\lambda$

We take steps in the direction that maximises the ELBO

$$\boldsymbol{\nabla}_\lambda \mathsf{ELBO} = \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \sum_{i=1}^{N} \log p(y^{(i)}|\theta, x^{(i)}) \right] - \boldsymbol{\nabla}_\lambda \, \mathsf{KL}(q(\theta|\lambda)||p(\theta))$$

# Gradient-based optimisation for $\lambda$

We take steps in the direction that maximises the ELBO

$$\boldsymbol{\nabla}_\lambda \text{ELBO} = \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \sum_{i=1}^{N} \log p(y^{(i)}|\theta, x^{(i)}) \right] - \boldsymbol{\nabla}_\lambda \text{KL}(q(\theta|\lambda)||p(\theta))$$

By assumption (mean field and exponential families), KL is tractable (we pack it in a node and autodiff does the job)!

# Gradient-based optimisation for $\lambda$

We take steps in the direction that maximises the ELBO

$$\boldsymbol{\nabla}_\lambda \mathsf{ELBO} = \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \sum_{i=1}^{N} \log p(y^{(i)}|\theta, x^{(i)}) \right] - \boldsymbol{\nabla}_\lambda \mathsf{KL}(q(\theta|\lambda)||p(\theta))$$

By assumption (mean field and exponential families), KL is tractable (we pack it in a node and autodiff does the job)!

How about the first term? What if $N$ is prohibitively large?

# Gradient-based optimisation for $\lambda$

We take steps in the direction that maximises the ELBO

$$\boldsymbol{\nabla}_\lambda \text{ELBO} = \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \sum_{i=1}^{N} \log p(y^{(i)}|\theta, x^{(i)}) \right] - \boldsymbol{\nabla}_\lambda \text{KL}(q(\theta|\lambda)||p(\theta))$$

By assumption (mean field and exponential families), KL is tractable (we pack it in a node and autodiff does the job)!

How about the first term? What if $N$ is prohibitively large?

$$\sum_{i=1}^{N} \log p(y^{(i)}|\theta, x^{(i)}) \quad \text{is certainly prohibitive!}$$

# Stochastic gradients are allowed

Noisy, but unbiased, gradients:

$$\nabla_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \sum_{i=1}^{N} \log p(y^{(i)}|\theta, x^{(i)}) \right]$$

# Stochastic gradients are allowed

Noisy, but unbiased, gradients:

$$\boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \sum_{i=1}^{N} \log p(y^{(i)}|\theta, x^{(i)}) \right]$$

$$= \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ N \sum_{i=1}^{N} \frac{1}{N} \log p(y^{(i)}|\theta, x^{(i)}) \right] \qquad \text{multiply by } N/N$$

## Stochastic gradients are allowed

Noisy, but unbiased, gradients:

$$
\boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \sum_{i=1}^{N} \log p(y^{(i)}|\theta, x^{(i)}) \right]
$$

$$
= \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ N \sum_{i=1}^{N} \frac{1}{N} \log p(y^{(i)}|\theta, x^{(i)}) \right] \qquad \text{multiply by } N/N
$$

$$
= \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ N \mathbb{E}_{I \sim \mathcal{U}(1/N)} \left[ \log p(y^{(I)}|\theta, x^{(I)}) \right] \right]
$$

## Stochastic gradients are allowed

Noisy, but unbiased, gradients:

$$\boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \sum_{i=1}^{N} \log p(y^{(i)}|\theta, x^{(i)}) \right]$$

$$= \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ N \sum_{i=1}^{N} \frac{1}{N} \log p(y^{(i)}|\theta, x^{(i)}) \right] \qquad \text{multiply by } N/N$$

$$= \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ N \mathbb{E}_{I \sim \mathcal{U}(1/N)} \left[ \log p(y^{(I)}|\theta, x^{(I)}) \right] \right]$$

$$= N \boldsymbol{\nabla}_\lambda \mathbb{E}_{I \sim \mathcal{U}(1/N)} \left[ \mathbb{E}_{q(\theta|\lambda)} \left[ \log p(y^{(I)}|\theta, x^{(I)}) \right] \right] \qquad \text{swap expectations}$$

## Stochastic gradients are allowed

Noisy, but unbiased, gradients:

$$\boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \sum_{i=1}^N \log p(y^{(i)}|\theta, x^{(i)}) \right]$$

$$= \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ N \sum_{i=1}^N \frac{1}{N} \log p(y^{(i)}|\theta, x^{(i)}) \right] \qquad \text{multiply by } N/N$$

$$= \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ N \mathbb{E}_{I \sim \mathcal{U}(1/N)} \left[ \log p(y^{(I)}|\theta, x^{(I)}) \right] \right]$$

$$= N \boldsymbol{\nabla}_\lambda \mathbb{E}_{I \sim \mathcal{U}(1/N)} \left[ \mathbb{E}_{q(\theta|\lambda)} \left[ \log p(y^{(I)}|\theta, x^{(I)}) \right] \right] \qquad \text{swap expectations}$$

$$= N \mathbb{E}_{I \sim \mathcal{U}(1/N)} \left[ \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \log p(y^{(I)}|\theta, x^{(I)}) \right] \right] \qquad \text{linearity}$$

## Stochastic gradients are allowed

Noisy, but unbiased, gradients:

$$\boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \sum_{i=1}^N \log p(y^{(i)}|\theta, x^{(i)}) \right]$$

$$= \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ N \sum_{i=1}^N \frac{1}{N} \log p(y^{(i)}|\theta, x^{(i)}) \right] \qquad \text{multiply by } N/N$$

$$= \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ N \mathbb{E}_{I \sim \mathcal{U}(1/N)} \left[ \log p(y^{(I)}|\theta, x^{(I)}) \right] \right]$$

$$= N \boldsymbol{\nabla}_\lambda \mathbb{E}_{I \sim \mathcal{U}(1/N)} \left[ \mathbb{E}_{q(\theta|\lambda)} \left[ \log p(y^{(I)}|\theta, x^{(I)}) \right] \right] \qquad \text{swap expectations}$$

$$= N \mathbb{E}_{I \sim \mathcal{U}(1/N)} \left[ \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \log p(y^{(I)}|\theta, x^{(I)}) \right] \right] \qquad \text{linearity}$$

$$\overset{\text{MC}}{\approx} \frac{N}{M} \sum_{i=1}^M \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \log p(y^{(i)}|\theta, x^{(i)}) \right] \qquad I \sim \mathcal{U}(1/N)$$

Sample a batch, solve expected value under $q$, then take gradient.

# Challenge

$$\boldsymbol{\nabla}_\lambda \mathsf{ELBO} =$$

$$\boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \sum_{i=1}^{N} \log p(y^{(i)}|\theta, x^{(i)}) \right] - \boldsymbol{\nabla}_\lambda \, \mathsf{KL}(q(\theta|\lambda)||p(\theta))$$

# Challenge

$\boldsymbol{\nabla}_\lambda \mathsf{ELBO} =$

$\boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \sum_{i=1}^{N} \log p(y^{(i)}|\theta, x^{(i)}) \right] - \boldsymbol{\nabla}_\lambda \mathsf{KL}(q(\theta|\lambda)||p(\theta))$

$= N \mathbb{E}_{\mathcal{U}(1/N)} \left[ \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \log p(y^{(I)}|\theta, x^{(I)}) \right] \right] - \boldsymbol{\nabla}_\lambda \mathsf{KL}(q(\theta|\lambda)||p(\theta))$

## Challenge

$\nabla_\lambda \, \mathsf{ELBO} =$

$$\nabla_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \sum_{i=1}^{N} \log p(y^{(i)}|\theta, x^{(i)}) \right] - \nabla_\lambda \, \mathsf{KL}(q(\theta|\lambda)||p(\theta))$$

$$= N \mathbb{E}_{\mathcal{U}(1/N)} \left[ \nabla_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \log p(y^{(I)}|\theta, x^{(I)}) \right] \right] - \nabla_\lambda \, \mathsf{KL}(q(\theta|\lambda)||p(\theta))$$

- we can compute KL and thus differentiate it
- mini-batching is allowed, so we can compute the first term for a few datapoints at a time
- but can we really solve $\mathbb{E}_{q(\theta|\lambda)} \left[ \log p(y^{(i)}|\theta, x^{(i)}) \right]$ for even a single instance?

## Reparameterisation

Let there be a bijective transformation $t(\epsilon, \lambda)$, such that

$$\epsilon \sim \phi(\epsilon) \quad \Leftrightarrow \quad t(\epsilon, \lambda) \sim q(\theta|\lambda)$$

## Reparameterisation

Let there be a bijective transformation $t(\epsilon, \lambda)$, such that

$$\epsilon \sim \phi(\epsilon) \quad \Leftrightarrow \quad t(\epsilon, \lambda) \sim q(\theta | \lambda)$$

Then it's clear that

- $\theta = t(\epsilon, \lambda)$ for $\epsilon \sim \phi(\cdot)$
- and $\epsilon = t^{-1}(\theta, \lambda)$ for $\theta \sim q(\cdot | \lambda)$

and note that $t$ **absorbs** $q$'s **dependence** on $\lambda$

## Reparameterisation

Let there be a bijective transformation $t(\epsilon, \lambda)$, such that

$$\epsilon \sim \phi(\epsilon) \quad \Leftrightarrow \quad t(\epsilon, \lambda) \sim q(\theta | \lambda)$$

Then it's clear that

- $\theta = t(\epsilon, \lambda)$ for $\epsilon \sim \phi(\cdot)$
- and $\epsilon = t^{-1}(\theta, \lambda)$ for $\theta \sim q(\cdot | \lambda)$

and note that $t$ **absorbs** $q$'s **dependence** on $\lambda$

Example: location-scale families such as the Normal distribution

$$\theta_d = \underbrace{\mu_d + \epsilon_d \sigma_d}_{=t(\epsilon_d, \lambda_d)} \qquad\qquad \epsilon_d = \frac{\theta_d - \mu_d}{\sigma_d} \sim \mathcal{N}(0, 1)$$

$$\lambda_d = (\mu_d, \sigma_d) \text{ and } \sigma_d > 0$$

# Reparameterisation (cont.)

Let there be a bijective transformation $t(\epsilon, \lambda)$, such that

$$\epsilon \sim \phi(\epsilon) \quad \Leftrightarrow \quad t(\epsilon, \lambda) \sim q(\theta|\lambda)$$

Then it's clear that

- $\theta = t(\epsilon, \lambda)$ for $\epsilon \sim \phi(\cdot)$
- and $\epsilon = t^{-1}(\theta, \lambda)$ for $\theta \sim q(\cdot|\lambda)$

and note that $t$ **absorbs** $q$'s **dependence** on $\lambda$

---

Can you prove this law? Tip: use integration by substitution and change of density theorem.

## Reparameterisation (cont.)

Let there be a bijective transformation $t(\epsilon, \lambda)$, such that

$$\epsilon \sim \phi(\epsilon) \quad \Leftrightarrow \quad t(\epsilon, \lambda) \sim q(\theta | \lambda)$$

Then it's clear that

- $\theta = t(\epsilon, \lambda)$ for $\epsilon \sim \phi(\cdot)$
- and $\epsilon = t^{-1}(\theta, \lambda)$ for $\theta \sim q(\cdot | \lambda)$

and note that $t$ **absorbs** $q$'s **dependence** on $\lambda$

The *law of the unconscious statistician* says

$$\mathbb{E}_{q(\theta|\lambda)}[f(\theta)] = \mathbb{E}_{\phi(\epsilon)}[f(\theta = t(\epsilon, \lambda))]$$

---

Can you prove this law? Tip: use integration by substitution and change of density theorem.

# Reparameterised Gradients

Assume we pick $q(\theta|\lambda)$ from a reparameterisable family
e.g. location-scale distributions

## Reparameterised Gradients

Assume we pick $q(\theta|\lambda)$ from a reparameterisable family
e.g. location-scale distributions

$\boldsymbol{\nabla}_\lambda \text{ELBO} =$

$N \mathbb{E}_I \left[ \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \log p(y^{(I)}|\theta, x^{(I)}) \right] \right] - \boldsymbol{\nabla}_\lambda \, \text{KL}(q(\theta|\lambda)||p(\theta))$

# Reparameterised Gradients

Assume we pick $q(\theta|\lambda)$ from a reparameterisable family
e.g. location-scale distributions

$$\boldsymbol{\nabla}_\lambda \, \text{ELBO} =$$

$$N\mathbb{E}_I \left[ \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \log p(y^{(I)}|\theta, x^{(I)}) \right] \right] - \boldsymbol{\nabla}_\lambda \, \text{KL}(q(\theta|\lambda)||p(\theta))$$

$$= N\mathbb{E}_I \left[ \boldsymbol{\nabla}_\lambda \mathbb{E}_{\underline{\phi(\epsilon)}} \left[ \log p(y^{(I)}|\underline{\theta = t(\epsilon, \lambda)}, x^{(I)}) \right] \right] - \boldsymbol{\nabla}_\lambda \, \text{KL}(q(\theta|\lambda)||p(\theta))$$

# Reparameterised Gradients

Assume we pick $q(\theta|\lambda)$ from a reparameterisable family
e.g. location-scale distributions

$\boldsymbol{\nabla}_\lambda \, \text{ELBO} =$

$N \mathbb{E}_I \left[ \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \log p(y^{(I)}|\theta, x^{(I)}) \right] \right] - \boldsymbol{\nabla}_\lambda \, \text{KL}(q(\theta|\lambda)||p(\theta))$

$= N \mathbb{E}_I \left[ \boldsymbol{\nabla}_\lambda \mathbb{E}_{\underline{\phi(\epsilon)}} \left[ \log p(y^{(I)}|\underline{\theta = t(\epsilon, \lambda)}, x^{(I)}) \right] \right] - \boldsymbol{\nabla}_\lambda \, \text{KL}(q(\theta|\lambda)||p(\theta))$

$= N \mathbb{E}_I \left[ \mathbb{E}_{\phi(\epsilon)} \left[ \boldsymbol{\nabla}_\lambda \log p(y^{(I)}|t(\epsilon, \lambda), x^{(I)}) \right] \right] - \boldsymbol{\nabla}_\lambda \, \text{KL}(q(\theta|\lambda)||p(\theta))$

# Reparameterised Gradients

Assume we pick $q(\theta|\lambda)$ from a reparameterisable family
e.g. location-scale distributions

$\boldsymbol{\nabla}_\lambda \text{ELBO} =$

$N\mathbb{E}_I \left[ \boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)} \left[ \log p(y^{(I)}|\theta, x^{(I)}) \right] \right] - \boldsymbol{\nabla}_\lambda \text{KL}(q(\theta|\lambda)||p(\theta))$

$= N\mathbb{E}_I \left[ \boldsymbol{\nabla}_\lambda \mathbb{E}_{\underline{\phi(\epsilon)}} \left[ \log p(y^{(I)}|\underline{\theta = t(\epsilon, \lambda)}, x^{(I)}) \right] \right] - \boldsymbol{\nabla}_\lambda \text{KL}(q(\theta|\lambda)||p(\theta))$

$= N\mathbb{E}_I \left[ \mathbb{E}_{\phi(\epsilon)} \left[ \boldsymbol{\nabla}_\lambda \log p(y^{(I)}|t(\epsilon, \lambda), x^{(I)}) \right] \right] - \boldsymbol{\nabla}_\lambda \text{KL}(q(\theta|\lambda)||p(\theta))$

$= N\mathbb{E}_{\phi(\epsilon)} \left[ \mathbb{E}_I \left[ \boldsymbol{\nabla}_\lambda \log p(y^{(I)}|t(\epsilon, \lambda), x^{(I)}) \right] \right] - \boldsymbol{\nabla}_\lambda \text{KL}(q(\theta|\lambda)||p(\theta))$

# Reparameterised Gradient Estimate

$$\boldsymbol{\nabla}_\lambda \, \mathsf{ELBO} =$$

$$N\mathbb{E}_{\phi(\epsilon)} \left[ \mathbb{E}_I \left[ \boldsymbol{\nabla}_\lambda \log p(y^{(I)}|t(\epsilon, \lambda), x^{(I)}) \right] \right] - \boldsymbol{\nabla}_\lambda \, \mathsf{KL}(q(\theta|\lambda)||p(\theta))$$

# Reparameterised Gradient Estimate

$\boldsymbol{\nabla}_\lambda \, \text{ELBO} =$

$$N\mathbb{E}_{\phi(\epsilon)} \left[ \mathbb{E}_I \left[ \boldsymbol{\nabla}_\lambda \log p(y^{(I)}|t(\epsilon, \lambda), x^{(I)}) \right] \right] - \boldsymbol{\nabla}_\lambda \, \text{KL}(q(\theta|\lambda)||p(\theta))$$

$$\stackrel{\text{MC}}{\approx} \left( \frac{M}{NK} \sum_{k=1}^{K} \sum_{i=1}^{M} \boldsymbol{\nabla}_\lambda \log p(y^{(i)}| \underbrace{t(\epsilon^{(k)}, \lambda)}_{=\theta^{(k)}}, x^{(i)}) \right)$$

$$- \boldsymbol{\nabla}_\lambda \, \text{KL}(q(\theta|\lambda)||p(\theta))$$

where $\epsilon^{(k)} \sim \phi(\epsilon)$, $\theta^{(k)} = t(\epsilon^{(k)}, \lambda)$, and $I \sim \mathcal{U}(1/N)$

## Reparameterised Gradient Estimate

$\boldsymbol{\nabla}_\lambda \, \text{ELBO} =$

$$N\mathbb{E}_{\phi(\epsilon)} \left[ \mathbb{E}_I \left[ \boldsymbol{\nabla}_\lambda \log p(y^{(I)}|t(\epsilon, \lambda), x^{(I)}) \right] \right] - \boldsymbol{\nabla}_\lambda \, \text{KL}(q(\theta|\lambda)||p(\theta))$$

$$\stackrel{\text{MC}}{\approx} \left( \frac{M}{NK} \sum_{k=1}^{K} \sum_{i=1}^{M} \boldsymbol{\nabla}_\lambda \log p(y^{(i)}|\underbrace{t(\epsilon^{(k)}, \lambda)}_{=\theta^{(k)}}, x^{(i)}) \right)$$

$$- \boldsymbol{\nabla}_\lambda \, \text{KL}(q(\theta|\lambda)||p(\theta))$$

where $\epsilon^{(k)} \sim \phi(\epsilon)$, $\theta^{(k)} = t(\epsilon^{(k)}, \lambda)$, and $I \sim \mathcal{U}(1/N)$

- Sample parameters via deterministic reparameterisation
- Sample batch
- Compute likelihood and KL: forward
- Sampling parameters first allows for efficient parallel implementation

# After training?

Training now gives you a point estimate for $\lambda$
  so we are not training $p$, we are training $q$!

# After training?

Training now gives you a point estimate for $\lambda$
so we are not training $p$, we are training $q$!

After training, we don't have 1 model, we have a distribution $q(\theta|\lambda)$ over
"all possible models"

- $q(\theta|\mathcal{D})$ approximates the true posterior $p(\theta|\mathcal{D})$
- it should prefer models that are likely after observing data $\mathcal{D}$
  in light of whatever prior assumptions we made
- there are no convergence guarantees and most approximating families
  are too simple (underestimate variance)

# After training?

Training now gives you a point estimate for $\lambda$
so we are not training $p$, we are training $q$!

After training, we don't have 1 model, we have a distribution $q(\theta|\lambda)$ over "all possible models"

- $q(\theta|\mathcal{D})$ approximates the true posterior $p(\theta|\mathcal{D})$
- it should prefer models that are likely after observing data $\mathcal{D}$
  in light of whatever prior assumptions we made
- there are no convergence guarantees and most approximating families
  are too simple (underestimate variance)

After training we make inferences using $q(\theta|\lambda)$

- $p(y_*|\mathcal{D}, x_*) \approx \int q(\theta|\lambda)p(y_*|\theta, x_*)\mathrm{d}\theta$
  which we typically further approximate via sampling

# After training?

Training now gives you a point estimate for $\lambda$
 so we are not training $p$, we are training $q$!

After training, we don't have 1 model, we have a distribution $q(\theta|\lambda)$ over
"all possible models"

- $q(\theta|\mathcal{D})$ approximates the true posterior $p(\theta|\mathcal{D})$
- it should prefer models that are likely after observing data $\mathcal{D}$
  in light of whatever prior assumptions we made
- there are no convergence guarantees and most approximating families
  are too simple (underestimate variance)

After training we make inferences using $q(\theta|\lambda)$

- $p(y_*|\mathcal{D}, x_*) \approx \int q(\theta|\lambda)p(y_*|\theta, x_*)\mathrm{d}\theta$
  which we typically further approximate via sampling

We can also estimate $p(\mathcal{D})$ using $q$ and the importance sampling

## Terminology

Approximate posterior $\qquad\qquad\qquad\qquad\qquad\qquad q(\theta|\lambda)$

Variational inference $\qquad\qquad\qquad\qquad \arg\min_{q(\theta)} \; \mathrm{KL}(q(\theta)||p(\theta|\mathcal{D}))$

ELBO $\qquad\qquad\qquad\qquad \mathbb{E}_{q(\theta)}[\log p(\mathcal{D}|\theta)] - \mathrm{KL}(q(\theta)||p(\theta))$

Mean field assumption $\qquad\qquad\qquad\qquad q(\theta|\lambda) = \prod_{d=1}^{D} q(\theta_d|\lambda_d)$

Reparameterised gradients

$$\boldsymbol{\nabla}_\lambda \mathbb{E}_{q(\theta|\lambda)}[f(\theta)] = \mathbb{E}_{\phi(\epsilon)}[\boldsymbol{\nabla}_\theta f(\theta)\boldsymbol{\nabla}_\lambda t(\epsilon, \lambda)]$$

Posterior predictive distribution

$$p(y_*|\mathcal{D}, x_*) \approx \int q(\theta|\lambda)p(y_*|\theta, x_*)\mathrm{d}\theta$$

## Summary

VI tuns inference into optimisation and gives you a proxy to $p(\theta|\mathcal{D})$

Estimates of posterior predictive mean and variance
- help you decide whether or not to make a decision
  - in classification: consider plotting precision and recall against predictive variance
  - in regression: interval in which you expect a response to be

Estimates of marginal likelihood
- help you compare models under different hyperparameters

Caveat: limited understanding about the impact of our priors

# Literature

Variational inference                                          Blei et al. (2017)

Stochastic VI                                            Hoffman et al. (2013)
 for nonconjugate inference          Titsias and Lázaro-Gredilla (2014)

Bayes by backprop                                    Blundell et al. (2015)

Model comparison                                          MacKay (1992a)

# Outline

# Dropout

A **very** simple technique to make MLE more robust

- stochastic training: with probability $1 - p$, "drop" inputs to a fully connected layer
- possibly use $L_2$ regularisation (because why not?)
- deterministic test: disable "dropout" and scale weights by $p$

Srivastava et al. (2014)

# Relate dropout to BNNs

BNNs come with a somewhat disappointing fact, that we have no clue what classes of random functions a given prior leads to

Many BNNs however can be seen as an approximation to a GP

- and the nonlinearities we employ correspond to a certain known kernel

Then let's see that variational inference for this model, using a pretty specific approximation $q(\theta|\lambda)$, turns dropout into approximate inference for an approximate GP.

The consequence is that we gain access to estimates of marginal likelihood and posterior predictive distribution

## Notation

$\mathbf{w}$ is a $C$-dimensional column vector of parameters

$\mathbf{W} = [\mathbf{w}_r^\top]_{r=1}^R$ stacks row vectors into a $R \times C$ matrix

$\mathbf{x}$ is an $I$-dimensional input

$\mathbf{y}$ is an $O$-dimensional output

$\mathbf{x}_{1:N}, \mathbf{y}_{1:N}$ is a collection of input-output pairs

$\mathbf{y}_{:,d}$ gather the $d$th output of each observation

$\mathcal{N}(\mathbf{y}_{1:N}|\mathbf{m}_{1:N}, \mathbf{I}_N)$ denotes $O$ independent multivariate Gaussians
i.e. $\prod_{d=1}^O \mathcal{N}(\mathbf{y}_{:,d}|\mathbf{m}_{:,d}, \mathbf{I}_N)$

# A GP approximation

We specify a GP prior by specifying a kernel. Valid kernels can be composed into other valid kernels.

---

Formal properties of Kernels (Shawe-Taylor and Cristianini, 2004, Chapter 3)

# A GP approximation

We specify a GP prior by specifying a kernel. Valid kernels can be composed into other valid kernels. For example:

$$k(\mathbf{x}, \mathbf{x}') = \int p(\mathbf{w})p(b)\sigma(\mathbf{w}^\top\mathbf{x} + b)\sigma(\mathbf{w}^\top\mathbf{x}' + b)\mathrm{d}\mathbf{w}\mathrm{d}b$$

This particular kernel is not tractable to assess: we cannot solve the integral for general $p(\mathbf{w}, b)$ and nonlinearity $\sigma(\cdot)$.

---

Formal properties of Kernels (Shawe-Taylor and Cristianini, 2004, Chapter 3)

# A GP approximation

We specify a GP prior by specifying a kernel. Valid kernels can be composed into other valid kernels. For example:

$$k(\mathbf{x}, \mathbf{x}') = \int p(\mathbf{w})p(b)\sigma(\mathbf{w}^\top \mathbf{x} + b)\sigma(\mathbf{w}^\top \mathbf{x}' + b)\mathrm{d}\mathbf{w}\mathrm{d}b$$

This particular kernel is not tractable to assess: we cannot solve the integral for general $p(\mathbf{w}, b)$ and nonlinearity $\sigma(\cdot)$.

But we know MC estimation! For a fixed number of samples $H$

$$\hat{k}(\mathbf{x}, \mathbf{x}') = \frac{1}{H} \sum_{s=1}^{H} \sigma(\mathbf{w}_s^\top \mathbf{x} + b_s)\sigma(\mathbf{w}_s^\top \mathbf{x}' + b_s)$$

where $\mathbf{w}_s \sim p(\mathbf{w})$ and $b_s \sim p(b)$

Formal properties of Kernels (Shawe-Taylor and Cristianini, 2004, Chapter 3)

# GP based on $\hat{k}$

We can immediately define a GP using this new kernel

$$
\begin{aligned}
b_s &\sim \mathcal{N}(0, l_0^{-2}) \\
\mathbf{b} &= [b_1, \ldots, b_H]^\top \\
\mathbf{w}_s &\sim \mathcal{N}(0, l^{-2}\mathbf{I}_I) \\
\mathbf{W}_1 &= [\mathbf{w}_s^\top]_{s=1}^H
\end{aligned}
$$

$$
\begin{aligned}
\hat{\mathbf{K}} &= \hat{k}(\mathbf{x}_{1:N}, \mathbf{x}_{1:N}) \\
\mathbf{F}|\mathbf{x}, \mathbf{W}_1, \mathbf{b} &\sim \mathcal{N}(\mathbf{0}, \hat{\mathbf{K}}) \\
\mathbf{Y}|\mathbf{f} &\sim \mathcal{N}(\mathbf{f}, \tau^{-1}\mathbf{I}_N)
\end{aligned}
$$

Recall that $\hat{K}_{ij} = \frac{1}{H}\sum_{s=1}^H \sigma(\mathbf{w}_s^\top \mathbf{x}_i + b)\sigma(\mathbf{w}_s^\top \mathbf{x}_j + b)$

# GP based on $\hat{k}$

We can immediately define a GP using this new kernel

$$b_s \sim \mathcal{N}(0, l_0^{-2})$$
$$\mathbf{b} = [b_1, \ldots, b_H]^\top$$
$$\mathbf{w}_s \sim \mathcal{N}(0, l^{-2}\mathbf{I}_I)$$
$$\mathbf{W}_1 = [\mathbf{w}_s^\top]_{s=1}^H$$

$$\hat{\mathbf{K}} = \hat{k}(\mathbf{x}_{1:N}, \mathbf{x}_{1:N})$$
$$\mathbf{F}|\mathbf{x}, \mathbf{W}_1, \mathbf{b} \sim \mathcal{N}(\mathbf{0}, \hat{\mathbf{K}})$$
$$\mathbf{Y}|\mathbf{f} \sim \mathcal{N}(\mathbf{f}, \tau^{-1}\mathbf{I}_N)$$

Note that we have essentially parameterised our kernel and imposed a prior on the kernel parameters.

For any given parameter configuration $\mathbf{W}_1$ and $\mathbf{b}$, we get a GP:

$$p(\mathbf{y}_{1:N}|\mathbf{x}_{1:N}) = \int p(\mathbf{b})p(\mathbf{W}_1)p(\mathbf{y}_{1:N}|\mathbf{f}_{1:N})p(\mathbf{f}_{1:N}|\mathbf{W}_1, \mathbf{b}, \mathbf{x}_{1:N})\mathrm{d}\mathbf{f}_{1:N}\mathrm{d}\mathbf{W}_1\mathrm{d}\mathbf{b}$$

---

Recall that $\hat{K}_{ij} = \frac{1}{H}\sum_{s=1}^H \sigma(\mathbf{w}_s^\top \mathbf{x}_i + b)\sigma(\mathbf{w}_s^\top \mathbf{x}_j + b)$

# Regression vs Classification

Regression

$$\mathbf{Y}|\mathbf{f} \sim \mathcal{N}(\mathbf{f}, \tau^{-1}\mathbf{I}_N)$$

Classification

$$\mathbf{Y}|\mathbf{f} \sim \mathcal{N}(\mathbf{f}, 0\mathbf{I}_N)$$
$$C|\mathbf{y} \sim \text{Cat}(\text{softmax}(\mathbf{y}))$$

- the GP is essentially inducing a distribution over logits
- note the change of notation (to be closer to Gal's), here $\mathbf{y}$ is not an observation, $c$ is

# Marginalise $\mathbf{f}$

We can marginalise $\mathbf{f}$ for an assignment of $\mathbf{W}_1$ and $\mathbf{b}$

Can you show that $\hat{\mathbf{K}} = \Phi\Phi^\top$?

# Marginalise $\mathbf{f}$

We can marginalise $\mathbf{f}$ for an assignment of $\mathbf{W}_1$ and $\mathbf{b}$

Note that $\hat{\mathbf{K}} = \Phi\Phi^\top$ for $\Phi = [\phi(\mathbf{x}_n, \mathbf{W}_1, \mathbf{b})^\top]_{n=1}^N$
  with feature vectors $\phi(\mathbf{x}, \mathbf{W}_1, \mathbf{b}) = \sqrt{1/H}\sigma(\mathbf{W}_1\mathbf{x} + \mathbf{b})$

---

Can you show that $\hat{\mathbf{K}} = \Phi\Phi^\top$?

## Marginalise $\mathbf{f}$

We can marginalise $\mathbf{f}$ for an assignment of $\mathbf{W}_1$ and $\mathbf{b}$

Note that $\hat{\mathbf{K}} = \Phi\Phi^\top$ for $\Phi = [\phi(\mathbf{x}_n, \mathbf{W}_1, \mathbf{b})^\top]_{n=1}^N$
  with feature vectors $\phi(\mathbf{x}, \mathbf{W}_1, \mathbf{b}) = \sqrt{1/H}\sigma(\mathbf{W}_1\mathbf{x} + \mathbf{b})$

From marginalisation of a subset of jointly Gaussian variables

$$p(\mathbf{y}_{1:N}|\mathbf{x}_{1:N}) = \int p(\mathbf{W}_1)p(\mathbf{b})p(\mathbf{y}_{1:N}|\mathbf{f}_{1:N})p(\mathbf{f}_{1:N}|\mathbf{W}_1, \mathbf{b}, \mathbf{x}_{1:N})\mathrm{d}\mathbf{f}_{1:N}\mathrm{d}\mathbf{W}_1\mathrm{d}\mathbf{b}$$

$$= \int \mathcal{N}(\mathbf{y}_{1:N}|\mathbf{0}, \Phi\Phi^\top + \tau^{-1}\mathbf{I}_N)p(\mathbf{W}_1)p(\mathbf{b})\mathrm{d}\mathbf{W}_1\mathrm{d}\mathbf{b}$$

We have "parameters", but note the model remains non-parametric
  complexity (capacity) adjusts with data size

---

Can you show that $\hat{\mathbf{K}} = \Phi\Phi^\top$?

## Parametric assumption

Recall that every Gaussian is the marginal of some other Gaussian.

## Parametric assumption

Recall that every Gaussian is the marginal of some other Gaussian.

In particular, it's true that

$$\mathcal{N}(\mathbf{y}_{:,d}|\mathbf{0}, \Phi\Phi^\top + \tau^{-1}\mathbf{I}_N) = \int \mathcal{N}(\mathbf{y}_{:,d}|\Phi\mathbf{w}_d, \tau^{-1}\mathbf{I}_N)\mathcal{N}(\mathbf{w}_d|\mathbf{0}, \mathbf{I}_H)\mathrm{d}\mathbf{w}_d$$

where $\mathbf{w}_d \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_H)$ is an $H$-dimensional auxiliary random vector

## Parametric assumption

Recall that every Gaussian is the marginal of some other Gaussian.

In particular, it's true that

$$\mathcal{N}(\mathbf{y}_{:,d}|\mathbf{0}, \Phi\Phi^\top + \tau^{-1}\mathbf{I}_N) = \int \mathcal{N}(\mathbf{y}_{:,d}|\Phi\mathbf{w}_d, \tau^{-1}\mathbf{I}_N)\mathcal{N}(\mathbf{w}_d|\mathbf{0}, \mathbf{I}_H)\mathrm{d}\mathbf{w}_d$$

where $\mathbf{w}_d \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_H)$ is an $H$-dimensional auxiliary random vector

We introduce $O$ such vectors, $\mathbf{W}_2 = [\mathbf{w}_d^\top]_{d=1}^O$ each $\mathbf{w}_d \sim p(\mathbf{w})$

$$\begin{aligned}
p(\mathbf{y}_{1:N}|\mathbf{x}_{1:N}) &= \int \mathcal{N}(\mathbf{y}_{1:N}|\mathbf{0}, \Phi\Phi^\top + \tau^{-1}\mathbf{I}_N)p(\mathbf{W}_1)p(\mathbf{b})\mathrm{d}\mathbf{W}_1\mathrm{d}\mathbf{b} \\
&= \int p(\mathbf{y}_{1:N}|\mathbf{x}_{1:N}, \mathbf{W}_1, \mathbf{W}_2, \mathbf{b})p(\mathbf{W}_1)p(\mathbf{b})p(\mathbf{W}_2)\mathrm{d}\mathbf{W}_2\mathrm{d}\mathbf{W}_1\mathrm{d}\mathbf{b}
\end{aligned}$$

## Digest

In the "first layer" we project inputs to $H$-dimensional *feature vectors* and apply a nonlinearity $\sigma(\cdot)$. This is a "hidden layer".

The parameters $\mathbf{W}_1$ and $\mathbf{b}$ of this projection are stochastic.

## Digest

In the "first layer" we project inputs to $H$-dimensional *feature vectors* and apply a nonlinearity $\sigma(\cdot)$. This is a "hidden layer".

The parameters $\mathbf{W}_1$ and $\mathbf{b}$ of this projection are stochastic.

In the "second layer" we decompose the covariance matrix using an $H$-dimensional vector (regardless of data size $N$). We do so independently for each of the $O$ outputs of the model. This is an "output layer".

Again the parameters $\mathbf{W}_2$ of this decomposition are stochastic.

## Digest

In the "first layer" we project inputs to $H$-dimensional *feature vectors* and apply a nonlinearity $\sigma(\cdot)$. This is a "hidden layer".

The parameters $\mathbf{W}_1$ and $\mathbf{b}$ of this projection are stochastic.

In the "second layer" we decompose the covariance matrix using an $H$-dimensional vector (regardless of data size $N$). We do so independently for each of the $O$ outputs of the model. This is an "output layer".

Again the parameters $\mathbf{W}_2$ of this decomposition are stochastic.

In the limit of an infinite hidden layer, marginalising over the stochastic parameters gives us a GP likelihood.

The induced kernel depends on the non-linearity $\sigma(\cdot)$

# A parametric approximation to a GP

A BNN with Gaussian priors over parameters is a parametric approximation to a GP

$$p(\mathbf{y}_{1:N}|\mathbf{x}_{1:N}) = \int p(\mathbf{y}_{1:N}|\mathbf{x}_{1:N}, \mathbf{W}_1, \mathbf{W}_2, \mathbf{b})p(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b})\mathrm{d}\mathbf{W}_2\mathrm{d}\mathbf{W}_1\mathrm{d}\mathbf{b}$$

Obviously the marginal is intractable, after all this is a BNN!

---

This is true for tanh, relu, step functions, sigmoid, for example.

# VI

Recipe

- propose a parametric proxy $q(\theta|\lambda)$
  $\theta = \{\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}\}$
- which you can reparameterise
- choose $\lambda$ to maximise the ELBO via stochastic gradient-based optimisation
- now because we count on autodiff, make sure you use differentiable nonlinearities (step function is no longer an option)

# Mean field

Independence across parameter groups

$$q(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}|\lambda) = q(\mathbf{W}_1|\lambda_1)q(\mathbf{W}_2|\lambda_2)q(\mathbf{b}|\lambda_b)$$

# Mean field

Independence across parameter groups

$$q(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}|\lambda) = q(\mathbf{W}_1|\lambda_1)q(\mathbf{W}_2|\lambda_2)q(\mathbf{b}|\lambda_b)$$

Diagonal Gaussian for biases

$$q(\mathbf{b}|\lambda) = \mathcal{N}(\mathbf{b}|\mathbf{m}, \text{diag}(\boldsymbol{\sigma}^2)) \quad \text{nicely reparameterisable!}$$

# Mean field

Independence across parameter groups

$$q(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}|\lambda) = q(\mathbf{W}_1|\lambda_1)q(\mathbf{W}_2|\lambda_2)q(\mathbf{b}|\lambda_b)$$

Diagonal Gaussian for biases

$$q(\mathbf{b}|\lambda) = \mathcal{N}(\mathbf{b}|\mathbf{m}, \mathrm{diag}(\boldsymbol{\sigma}^2)) \quad \text{nicely reparameterisable!}$$

Mixture of Gaussians posterior with $p \in (0, 1)$

$$q(\mathbf{w}_r|\lambda) = p\mathcal{N}(\mathbf{w}_r|\mathbf{m}_r, \mathrm{diag}(\boldsymbol{\sigma}^2)) + (1 - p)\mathcal{N}(\mathbf{w}_r|\mathbf{0}, \mathrm{diag}(\boldsymbol{\sigma}^2))$$

## Mean field

Independence across parameter groups

$$q(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}|\lambda) = q(\mathbf{W}_1|\lambda_1)q(\mathbf{W}_2|\lambda_2)q(\mathbf{b}|\lambda_b)$$

Diagonal Gaussian for biases

$$q(\mathbf{b}|\lambda) = \mathcal{N}(\mathbf{b}|\mathbf{m}, \text{diag}(\boldsymbol{\sigma}^2)) \quad \text{nicely reparameterisable!}$$

Mixture of Gaussians posterior with $p \in (0, 1)$

$$q(\mathbf{w}_r|\lambda) = p\mathcal{N}(\mathbf{w}_r|\mathbf{m}_r, \text{diag}(\boldsymbol{\sigma}^2)) + (1 - p)\mathcal{N}(\mathbf{w}_r|\mathbf{0}, \text{diag}(\boldsymbol{\sigma}^2))$$

Can we reparameterise samples from this posterior?

## Mean field

Independence across parameter groups

$$q(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}|\lambda) = q(\mathbf{W}_1|\lambda_1)q(\mathbf{W}_2|\lambda_2)q(\mathbf{b}|\lambda_b)$$

Diagonal Gaussian for biases

$$q(\mathbf{b}|\lambda) = \mathcal{N}(\mathbf{b}|\mathbf{m}, \text{diag}(\boldsymbol{\sigma}^2)) \quad \text{nicely reparameterisable!}$$

Mixture of Gaussians posterior with $p \in (0, 1)$

$$q(\mathbf{w}_r|\lambda) = p\mathcal{N}(\mathbf{w}_r|\mathbf{m}_r, \text{diag}(\boldsymbol{\sigma}^2)) + (1-p)\mathcal{N}(\mathbf{w}_r|\mathbf{0}, \text{diag}(\boldsymbol{\sigma}^2))$$

Can we reparameterise samples from this posterior?

- $\mathbf{w}_r = z(\mathbf{m}_r + \boldsymbol{\sigma} \odot \epsilon) + (1-z)\boldsymbol{\sigma} \odot \epsilon$
  for $Z \sim \text{Bern}(p)$ and $\epsilon \sim \phi(\epsilon)$

# Mixture of Deltas

Let $\boldsymbol{\sigma} \to \mathbf{0}$

- from mixture of Gaussians to mixture of Deltas

$$Z \sim \text{Bern}(p)$$

$$\mathbf{w}_r = \begin{cases} \mathbf{m}_r + \mathbf{0} \odot \epsilon & \text{if } z = 1 \\ (1-z)\mathbf{0} \odot \epsilon & \text{if } z = 0 \end{cases}$$

$$= z\mathbf{m}_r$$

- biases are deterministic because $\boldsymbol{\sigma} \to \mathbf{0}$ the Gaussian tends to $\delta(\mathbf{m} - \mathbf{b})$

## Mixture of Deltas

Let $\sigma \to \mathbf{0}$

- from mixture of Gaussians to mixture of Deltas

$$Z \sim \text{Bern}(p)$$

$$\mathbf{w}_r = \begin{cases} \mathbf{m}_r + \mathbf{0} \odot \epsilon & \text{if } z = 1 \\ (1-z)\mathbf{0} \odot \epsilon & \text{if } z = 0 \end{cases}$$

$$= z\mathbf{m}_r$$

- biases are deterministic because $\sigma \to \mathbf{0}$ the Gaussian tends to $\delta(\mathbf{m} - \mathbf{b})$

The affine transform in fully connected layers becomes

$$([(z_s \mathbf{m}_s)^\top]_{s=1}^H)\mathbf{x} + \mathbf{m} = ([\mathbf{m}_s^\top]_{s=1}^H)(\mathbf{z} \odot x) + \mathbf{m}$$

with probability $p$, we essentially drop inputs

Same happens with $\mathbf{W}_2$ (weights of the second layer)

# ELBO

Recall the ELBO

$$N \, \mathbb{E}_I \left[ \mathbb{E}_Z \left[ \log p(y^{(I)}|x^{(I)}, \theta = t(z, \lambda)) \right] \right] - \mathsf{KL}(q(\theta|\lambda)||p(\theta))$$

where

- $\theta = \{\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}\}$ and $\lambda = \{\mathbf{M}_1, \mathbf{M}_2, \mathbf{m}\}$
- $Z_{1,i} \sim \mathsf{Bern}(p)$ and $Z_{2,s} \sim \mathsf{Bern}(p)$

# ELBO

Recall the ELBO

$$N \, \mathbb{E}_I \left[ \mathbb{E}_Z \left[ \log p(y^{(I)}|x^{(I)}, \theta = t(z, \lambda)) \right] \right] - \mathsf{KL}(q(\theta|\lambda)||p(\theta))$$

where

- $\theta = \{\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}\}$ and $\lambda = \{\mathbf{M}_1, \mathbf{M}_2, \mathbf{m}\}$
- $Z_{1,i} \sim \mathsf{Bern}(p)$ and $Z_{2,s} \sim \mathsf{Bern}(p)$

We can sample with a reparameterisation

- draws from $\mathsf{Bern}(p)$ are used to mask the inputs to layers

# ELBO

Recall the ELBO

$$N \, \mathbb{E}_I \left[ \mathbb{E}_Z \left[ \log p(y^{(I)}|x^{(I)}, \theta = t(z, \lambda)) \right] \right] - \mathsf{KL}(q(\theta|\lambda)||p(\theta))$$

where

- $\theta = \{\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}\}$ and $\lambda = \{\mathbf{M}_1, \mathbf{M}_2, \mathbf{m}\}$
- $Z_{1,i} \sim \mathsf{Bern}(p)$ and $Z_{2,s} \sim \mathsf{Bern}(p)$

We can sample with a reparameterisation

- draws from $\mathsf{Bern}(p)$ are used to mask the inputs to layers

We can assess the likelihood

- because we were the ones to choose it

# ELBO

Recall the ELBO

$$N \, \mathbb{E}_I \left[ \mathbb{E}_Z \left[ \log p(y^{(I)} | x^{(I)}, \theta = t(z, \lambda)) \right] \right] - \text{KL}(q(\theta|\lambda) || p(\theta))$$

where

- $\theta = \{\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}\}$ and $\lambda = \{\mathbf{M}_1, \mathbf{M}_2, \mathbf{m}\}$
- $Z_{1,i} \sim \text{Bern}(p)$ and $Z_{2,s} \sim \text{Bern}(p)$

We can sample with a reparameterisation

- draws from $\text{Bern}(p)$ are used to mask the inputs to layers

We can assess the likelihood

- because we were the ones to choose it

We are only missing a KL term: which can be nicely approximated by $L_2$ on $\lambda$

# KL approximation

In regression (recall $\tau$ is the prior precision for the likelihood $Y|f$)

$$\frac{p}{2\tau N}||\mathbf{M}_1||_2^2 + \frac{p}{2\tau N}||\mathbf{M}_2||_2^2 + \frac{1}{2\tau N}||\mathbf{m}||_2^2 \tag{1}$$

In classification (we assume a degenerate Gaussian for $Y|f$)

$$\frac{p}{2N}||\mathbf{M}_1||_2^2 + \frac{p}{2N}||\mathbf{M}_2||_2^2 + \frac{1}{2N}||\mathbf{m}||_2^2 \tag{2}$$

See Gal and Ghahramani (2016b)

## Prior parameters

The prior length-scale is the inverse of the standard deviation of the
distribution over the scaling weights in affine layers, it controls the rate of
change of the sampled functions.

The prior precision (in regression) controls the observation noise, smaller
precision leads to bigger error bars

In classification we let $\tau^{-1} \to 0$.

These are hyperparameters you have to search for. You can also relate
them to the weight decay if your NN library offers weight decay out of the
box.

# Approximate posterior parameters

$\mathbf{M}_1$ variational mean for input-to-hidden, $\mathbf{m}$ variational mean of bias vector, $\mathbf{M}_2$ variational mean for hidden-to-output

- these are the only trainable parameters

In principle, we can have one Bernoulli parameter per layer, generally we don't because Bernoulli sampling is nondifferentiable and thus we have to tune the parameter by hand.

## Inferences

Marginal likelihood: get estimates via importance sampling to compare
different hyperparameters

## Inferences

Marginal likelihood: get estimates via importance sampling to compare different hyperparameters

Posterior predictive distribution
- do not disable dropout at test time

## Inferences

Marginal likelihood: get estimates via importance sampling to compare different hyperparameters

Posterior predictive distribution

- do not disable dropout at test time
- compute empirical estimates of predictive mean and variance from stochastic outputs

## Inferences

Marginal likelihood: get estimates via importance sampling to compare different hyperparameters

Posterior predictive distribution
- do not disable dropout at test time
- compute empirical estimates of predictive mean and variance from stochastic outputs
- for classification, consider uncertainty over class probabilities, for example, boxplot samples of the probability $p(y|x_*, \theta)$ of each outcome $y$ for $\theta \sim q(\theta|\lambda)$

## Inferences

Marginal likelihood: get estimates via importance sampling to compare different hyperparameters

Posterior predictive distribution

- do not disable dropout at test time
- compute empirical estimates of predictive mean and variance from stochastic outputs
- for classification, consider uncertainty over class probabilities, for example, boxplot samples of the probability $p(y|x_*, \theta)$ of each outcome $y$ for $\theta \sim q(\theta|\lambda)$
- multiple forward passes, but a single trained model

## Inferences

Marginal likelihood: get estimates via importance sampling to compare different hyperparameters

Posterior predictive distribution

- do not disable dropout at test time
- compute empirical estimates of predictive mean and variance from stochastic outputs
- for classification, consider uncertainty over class probabilities, for example, boxplot samples of the probability $p(y|x_*, \theta)$ of each outcome $y$ for $\theta \sim q(\theta|\lambda)$
- multiple forward passes, but a single trained model
- unlike in an ensemble, we are sampling from $q(\theta|\lambda)$, an approximation to $p(\theta|\mathcal{D})$

## Extensions

Note a few things

- it was crucial to use a mixture of *deltas* as variational approximation
- this allowed us to sample parameters by having a mask over inputs (rather than over parameters)
- this trick seems general, but it does depend on the type of layer we deal with

## Extensions

Note a few things

- it was crucial to use a mixture of *deltas* as variational approximation
- this allowed us to sample parameters by having a mask over inputs (rather than over parameters)
- this trick seems general, but it does depend on the type of layer we deal with

An RNN is just a FFNN dynamically unfolded through time

- all we need is to sample the mask once per data point
- and reuse the same mask for all steps in the sequence

## Extensions

Note a few things

- it was crucial to use a mixture of *deltas* as variational approximation
- this allowed us to sample parameters by having a mask over inputs (rather than over parameters)
- this trick seems general, but it does depend on the type of layer we deal with

An RNN is just a FFNN dynamically unfolded through time

- all we need is to sample the mask once per data point
- and reuse the same mask for all steps in the sequence

CNNs can be reformulated as a linear operation followed by a pooling non-linearity, in this view we need to drop outputs of the linear operation (a parameterised inner product) before pooling

## Summary

Bayesian posterior inference for NNs with negligible training effort

Bayesian posterior predictive at linear cost (one forward pass per sample)

Future research
- better posterior approximations (more correlations)
- better handle on properties of kernels
- BNNs typically underestimate variance

# Literature

Yarin Gal's thesis and blogpost                                    Gal (2016)

Dropout as Bayesian approximation     Gal and Ghahramani (2016b, esp appendix)

CNN and RNN variants                 Gal and Ghahramani (2016a,c)

# References I

José M Bernardo and Adrian FM Smith. *Bayesian theory*, volume 405. John Wiley & Sons, 2009.

Christopher M. Bishop. *Pattern recognition and machine learning (information science and statistics)*. Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2006. ISBN 0-387-31073-8. tex.date-added: 2019-09-19 08:15:14 +0000 tex.date-modified: 2019-09-19 08:15:14 +0000.

David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017. Publisher: Taylor & Francis.

# References II

Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In Francis Bach and David Blei, editors, *Proceedings of the 32nd international conference on machine learning*, volume 37 of *Proceedings of machine learning research*, pages 1613–1622, Lille, France, July 2015. PMLR. URL http://proceedings.mlr.press/v37/blundell15.html. tex.date-added: 2019-09-24 16:11:39 +0000 tex.date-modified: 2019-09-24 16:11:46 +0000 tex.pdf: http://proceedings.mlr.press/v37/blundell15.pdf.

K. Funahashi. On the approximate realization of continuous mappings by neural networks. *Neural Netw.*, 2(3):183–192, May 1989. ISSN 0893-6080. doi: 10.1016/0893-6080(89)90003-8. URL http://dx.doi.org/10.1016/0893-6080(89)90003-8. Number of pages: 10 Publisher: Elsevier Science Ltd. tex.acmid: 71290 tex.address:

## References III

Oxford, UK, UK tex.date-added: 2019-09-19 05:15:10 +0000
tex.date-modified: 2019-09-19 05:15:17 +0000 tex.issue_date: 1989.

Yarin Gal. *Uncertainty in deep learning*. PhD thesis, PhD thesis, University
of Cambridge, 2016. tex.date-added: 2019-09-19 11:45:00 +0000
tex.date-modified: 2019-09-24 16:02:09 +0000.

Yarin Gal and Zoubin Ghahramani. Bayesian convolutional neural networks
with Bernoulli approximate variational inference. In *ICLR - workshop
track*, 2016a.

Yarin Gal and Zoubin Ghahramani. Dropout as a Bayesian Approximation:
Representing Model Uncertainty in Deep Learning. In Maria Florina
Balcan and Kilian Q. Weinberger, editors, *Proceedings of The 33rd
International Conference on Machine Learning*, volume 48 of
*Proceedings of Machine Learning Research*, pages 1050–1059, New
York, New York, USA, June 2016b. PMLR. URL
http://proceedings.mlr.press/v48/gal16.html.

# References IV

Yarin Gal and Zoubin Ghahramani. A theoretically grounded application of dropout in recurrent neural networks. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in neural information processing systems 29*, pages 1019–1027. Curran Associates, Inc., 2016c. URL `http://papers.nips.cc/paper/ 6241-a-theoretically-grounded-application-of-dropout-in-rec pdf`.

Andrew Gelman, John Carlin, Hal Stern, David Dunson, Aki Vehtari, and Donald Rubin. *Bayesian Data Analysis*. Chapman and Hall/CRC, third edition, 2013.

## References V

Matthew D. Hoffman, David M. Blei, Chong Wang, and John Paisley. Stochastic variational inference. *J. Mach. Learn. Res.*, 14(1):1303–1347, May 2013. ISSN 1532-4435. URL http://dl.acm.org/citation.cfm?id=2502581.2502622. Number of pages: 45 Publisher: JMLR.org tex.acmid: 2502622 tex.date-added: 2019-09-24 13:46:13 +0000 tex.date-modified: 2019-09-24 13:46:13 +0000 tex.issue_date: January 2013.

David J. C. MacKay. Bayesian interpolation. *Neural Comput.*, 4(3): 415–447, May 1992a. ISSN 0899-7667. doi: $10.1162/neco.1992.4.3.415$. URL http://dx.doi.org/10.1162/neco.1992.4.3.415. Number of pages: 33 Publisher: MIT Press tex.acmid: 148163 tex.address: Cambridge, MA, USA tex.date-added: 2019-09-18 21:10:45 +0000 tex.date-modified: 2019-09-18 21:10:45 +0000 tex.issue_date: May 1992.

# References VI

David JC MacKay. *Bayesian methods for adaptive models*. PhD thesis, California Institute of Technology, 1992b. tex.date-added: 2019-09-18 21:12:59 +0000 tex.date-modified: 2019-09-18 21:12:59 +0000.

David JC MacKay. Introduction to gaussian processes. *NATO ASI Series F Computer and Systems Sciences*, 168:133–166, 1998. Publisher: Springer Verlag tex.date-added: 2019-09-18 21:12:30 +0000 tex.date-modified: 2019-09-18 21:12:30 +0000.

Radford M. Neal. Priors for infinite networks. Technical report, Dept. of Computer Science, University of Toronto, 1994. tex.date-added: 2019-09-18 21:06:36 +0000 tex.date-modified: 2019-09-18 21:07:22 +0000.

Radford M. Neal. Priors for infinite networks. In *Bayesian learning for neural networks*, pages 29–53. Springer New York, New York, NY, 1996. ISBN 978-1-4612-0745-0.

# References VII

John Shawe-Taylor and Nello Cristianini. *Kernel methods for pattern analysis*. Cambridge University Press, New York, NY, USA, 2004. ISBN 0-521-81397-2. tex.date-added: 2019-09-25 04:24:57 +0000 tex.date-modified: 2019-09-25 04:25:07 +0000.

Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15 (56):1929–1958, 2014. URL http://jmlr.org/papers/v15/srivastava14a.html.

Michalis Titsias and Miguel Lázaro-Gredilla. Doubly Stochastic Variational Bayes for non-Conjugate Inference. In *International Conference on Machine Learning*, pages 1971–1979. PMLR, June 2014. URL http://proceedings.mlr.press/v32/titsias14.html. ISSN: 1938-7228.