

# Mathematics for Machine Learning

**Marc Deisenroth**

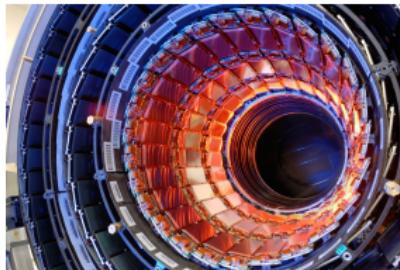
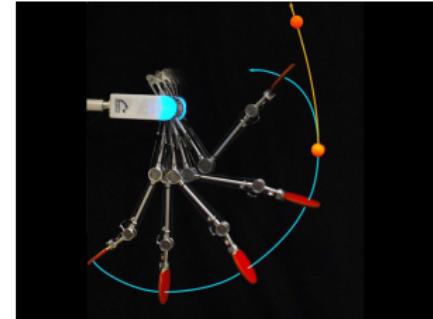
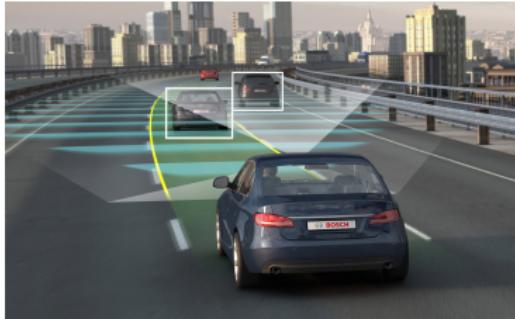
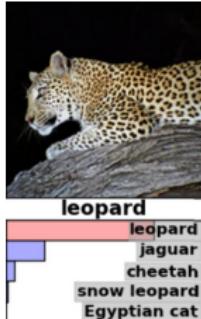
Statistical Machine Learning Group  
Department of Computing  
Imperial College London

 @mpd37

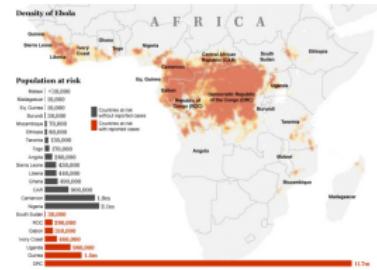
[m.deisenroth@imperial.ac.uk](mailto:m.deisenroth@imperial.ac.uk)  
[marc@prowler.io](mailto:marc@prowler.io)

Deep Learning 2017

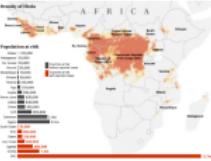
## Applications of Machine Learning



A screenshot of the Goodreads mobile application interface. At the top, a red banner reads "Today's Recommendations For You". Below it, a message says "Here's a daily sample of items recommended for you. Click here to see all recommendations". There are two main sections of recommendations. The first section, titled "LOOK INSIDE", features a book cover for "The Great Gatsby" by F. Scott Fitzgerald, with a star rating of 4.5 stars and 1,141 reviews. It includes a "Fix this recommendation" button. The second section, titled "LOOK INSIDE", features a book cover for "Simply JavaScript" by Kyle Simpson, with a star rating of 4.5 stars and 1,034 reviews. It also includes a "Fix this recommendation" button. At the bottom, there are tabs for "Any Category", "Algorithm", "Boxed Sets", "Business & Culture", "Java", "Graphic Design", "Horror", "Networking", "Networks, Protocols & APIs", and "SQL".



# Mathematical Concepts in Machine Learning



§ Linear algebra and matrix decomposition

§ **Differentiation**

§ Optimization

§ **Integration**

§ Probability theory and Bayesian inference

§ Functional analysis

# Outline

Introduction

Differentiation

Integration

# Overview

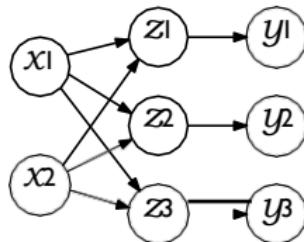
[Introduction](#)

[Differentiation](#)

[Integration](#)

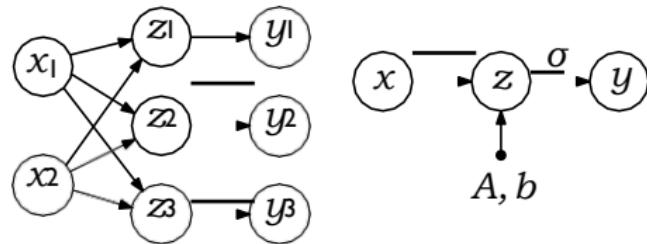
# Feedforward Neural Network

$$\begin{aligned}y &= \sigma p z q \\z &= Ax + b\end{aligned}$$



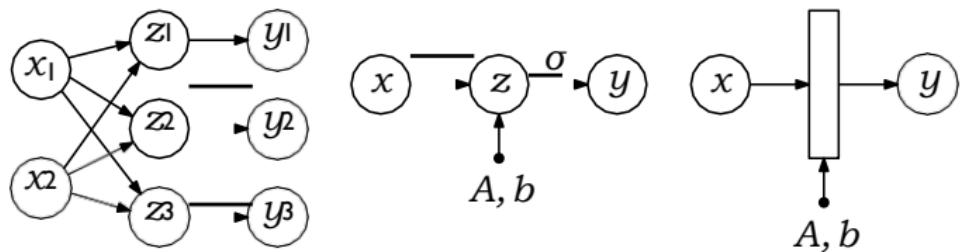
# Feedforward Neural Network

$$\begin{aligned}y &= \sigma p z q \\z &= Ax + b\end{aligned}$$



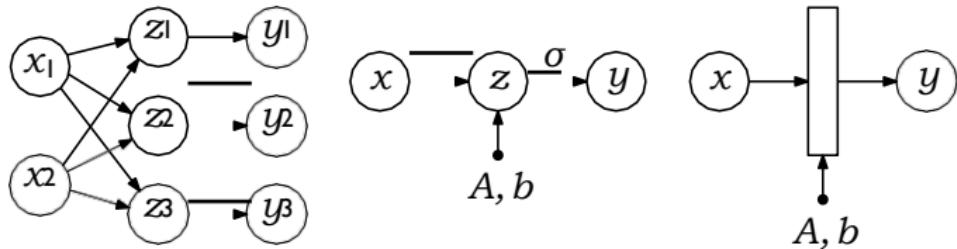
# Feedforward Neural Network

$$\begin{aligned}y &= \sigma p z q \\z &= Ax + b\end{aligned}$$



# Feedforward Neural Network

$$\begin{aligned}y &= \sigma p z q \\z &= Ax + b\end{aligned}$$



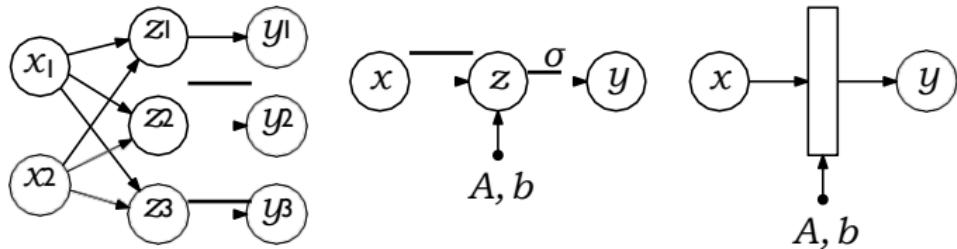
§ Training a neural network means parameter optimization:

Typically via some form of gradient descent

► **Challenge 1: Differentiation.** Compute gradients of a loss function with respect to neural network parameters  $A, b$

# Feedforward Neural Network

$$\begin{aligned}y &= \sigma(pz + q) \\z &= Ax + b\end{aligned}$$



- § Training a neural network means parameter optimization:  
Typically via some form of gradient descent
  - **Challenge 1: Differentiation.** Compute gradients of a loss function with respect to neural network parameters  $A, b$
- § Computing statistics (e.g., means, variances) of predictions
  - **Challenge 2: Integration.** Propagate uncertainty through a neural network

# Background: Matrix Multiplication

§ Matrix multiplication is not commutative, i.e.,  $AB \neq BA$

$$\begin{array}{ccc} \begin{matrix} \text{blue} \\ 3 \times 2 \end{matrix} & \begin{matrix} \text{yellow} \\ 2 \times 3 \end{matrix} & = & \begin{matrix} \text{green} \\ 3 \times 3 \end{matrix} \\ \begin{matrix} \text{yellow} \\ 2 \times 3 \end{matrix} & \begin{matrix} \text{blue} \\ 3 \times 2 \end{matrix} & = & \begin{matrix} \text{green} \\ 2 \times 2 \end{matrix} \end{array}$$

# Background: Matrix Multiplication

- § Matrix multiplication is not commutative, i.e.,  $AB \neq BA$

The diagram illustrates two examples of matrix multiplication. The first example shows a blue 3x3 matrix multiplied by an orange 3x3 matrix, resulting in a green 3x3 matrix. The second example shows an orange 2x3 matrix multiplied by a blue 3x2 matrix, resulting in a green 2x2 matrix. Both examples show that the order of multiplication matters.

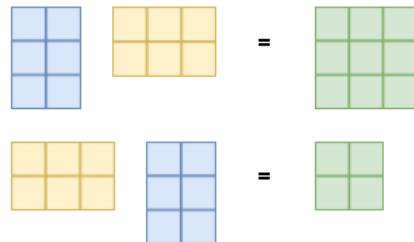
$$\begin{array}{ccc} \begin{matrix} \text{blue} \\ 3 \times 3 \end{matrix} & \begin{matrix} \text{orange} \\ 3 \times 3 \end{matrix} & = & \begin{matrix} \text{green} \\ 3 \times 3 \end{matrix} \\ \begin{matrix} \text{orange} \\ 2 \times 3 \end{matrix} & \begin{matrix} \text{blue} \\ 3 \times 2 \end{matrix} & = & \begin{matrix} \text{green} \\ 2 \times 2 \end{matrix} \end{array}$$

- § When multiplying matrices, the “neighboring” dimensions have to fit:

“**Left and Right**”  
Left      Right  
 $n^k$        $k^m$        $n^m$

# Background: Matrix Multiplication

§ Matrix multiplication is not commutative, i.e.,  $AB \neq BA$



§ When multiplying matrices, the “neighboring” dimensions have to fit:

“**Left and Right Dimension**”  
 $n^k$        $k^m$        $n^m$

$y = Ax$

$y = A \cdot \text{dot}(x)$

$y_i = \sum_j A_{ij} x_j$

$y = \text{np.einsum('ij, j', A, x)}$

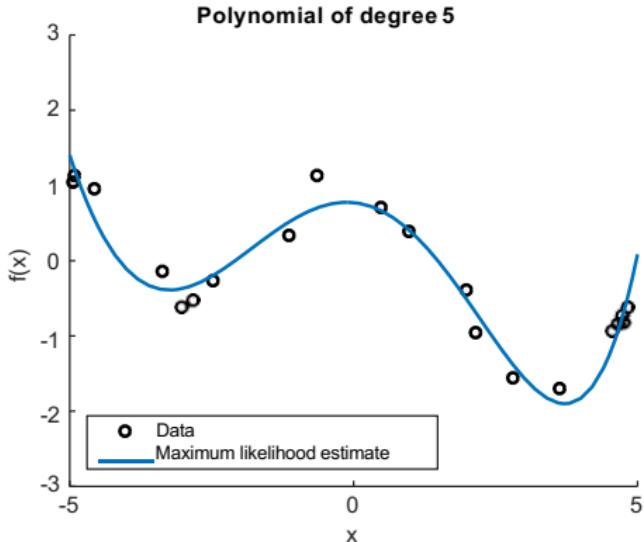
$C = AB$

$C = A \cdot \text{dot}(B)$

$C_{ij} = \sum_k A_{ik} B_{kj}$

$C = \text{np.einsum('ik, kj', A, B)}$

# Curve Fitting (Regression) in Machine Learning (1)

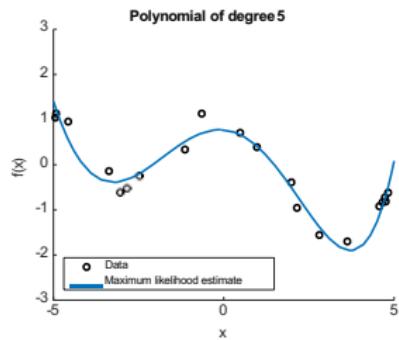


§ Setting: Given inputs  $x$ , predict outputs/targets  $y$   
§ Model  $f$  that depends on parameters  $\theta$ . Examples:

- § Linear model:  $f \propto x, \theta \in \mathbb{R}^D$
- § Neural network:  $f \propto x, \theta \in \text{NN}(x, \theta)$

# Curve Fitting (Regression) in Machine Learning (2)

- § Training data, e.g.,  $N$  pairs  $\{x_i, y_i\}$  of inputs  $x_i$  and observations  $y_i$
- § Training the model means finding parameters  $\theta^\circ$ , such that  $f(x_i, \theta^\circ) \approx y_i$



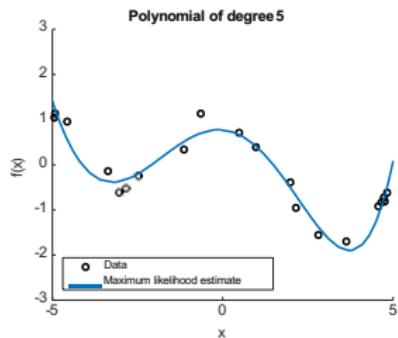
# Curve Fitting (Regression) in Machine Learning (2)

§ Training data, e.g.,  $N$  pairs  $\{x_i, y_i\}$  of inputs  $x_i$  and observations  $y_i$

§ Training the model means finding parameters  $\theta^*$ , such that  $f(x_i, \theta^*) \approx y_i$

§ Define a loss function, e.g.,  $\sum_{i=1}^N |y_i - f(x_i, \theta)|^2$ , which we want to optimize

§ Typically: Optimization based on some form of gradient descent  
► Differentiation required



# Overview

[Introduction](#)

[Differentiation](#)

[Integration](#)

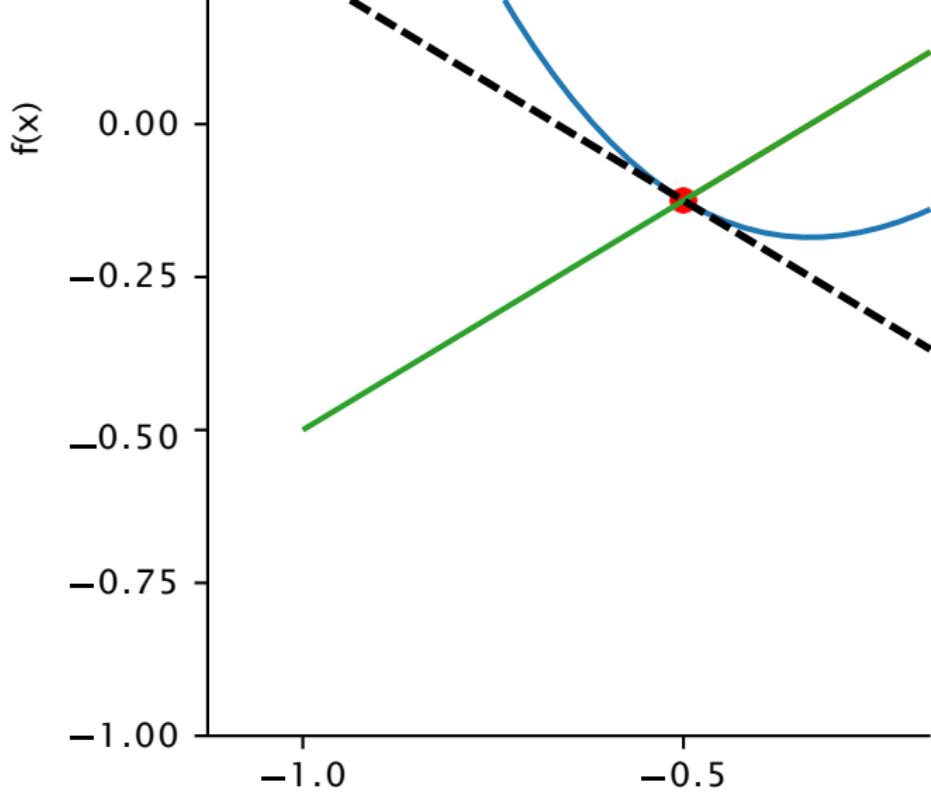
# Differentiation: Outline

1. Scalar differentiation:  $f: \mathbf{R} \rightarrow \mathbf{R}$
2. Multivariate case:  $f: \mathbf{R}^N \rightarrow \mathbf{R}$
3. Vector fields:  $f: \mathbf{R}^N \rightarrow \mathbf{R}^M$
4. General derivatives:  $f: \mathbf{R}^{M^N} \rightarrow \mathbf{R}^{P^Q}$

# Scalar Diff.

§ Derivative

► Slope



# Examples

$$f(p(x)) = x^n$$

$$f(p(x)) = \sin(p(x))$$

$$f(p(x)) = \tanh(p(x))$$

$$f(p(x)) = \exp(p(x))$$

$$f(p(x)) = \log(p(x))$$

$$f'(p(x)) = nx^{n-1}$$

$$f'(p(x)) = \cos(p(x))$$

$$f'(p(x)) = 1 - \tanh^2(p(x))$$

$$f'(p(x)) = \exp(p(x))$$

$$f'(p(x)) = \frac{1}{x}$$

# Rules

## § Sum Rule

$$f \text{pxq} + g \text{pxq} \stackrel{?}{=} f' \text{pxq} + g' \text{pxq} \quad \frac{df}{dx} + \frac{dg}{dx}$$

# Rules

## § Sum Rule

$$f \text{pxq} + g \text{pxq} = f' \text{pxq} + g' \text{pxq} = \frac{df}{dx} + \frac{dg}{dx}$$

## § Product Rule

$$(fg)' = f'g + fg'$$

# Rules

## § Sum Rule

$$f(p(x) + g(p(x)))' = f'(p(x)) \cdot g(p(x)) + \frac{df}{dx} + \frac{dg}{dx}$$

## § Product Rule

$$(f(p(x))g(p(x)))' = f'(p(x))g(p(x)) + f(p(x))g'(p(x)) \frac{df}{dx} + \frac{dg}{dx}$$

## § Chain Rule

$$(g(f(p(x))))' = g'(f(p(x))) \cdot f'(p(x)) \frac{dg}{df} \frac{df}{dx}$$

## Example: Chain Rule

$$pg'' f \frac{d}{dx} x^q = g p f x^{q-1} \cdot g \cancel{p} f x^{q-1} f \cancel{p} x^q = \frac{dg}{df} \frac{df}{dx}$$

$$\begin{aligned} g p z q &= \tanh p z q \\ z &= f p x q = x^n \\ pg'' f q^1 p x q &= \end{aligned}$$

## Example: Chain Rule

$$pg'' f \circ p x q = g p f p x q q^{-1} = g \circ f p x q f \circ p x q = \frac{dg}{df} \frac{df}{dx}$$

$$\begin{aligned} g p z q &= \tanh p z q \\ z &= f p x q = x^n \\ pg'' f \circ p x q &= \frac{1}{\cosh^2 p x} \tanh^{n-1} p x \end{aligned}$$

$\frac{dg}{df} \frac{df}{dx}$

$$f: \mathbf{R}^N \rightarrow \mathbf{R}$$

$$y = f(x_1, \dots, x_N)$$

§ Partial derivative (change one coordinate at a time):

$$\frac{\partial f}{\partial x_i} = \lim_{h \rightarrow 0} \frac{f(x_1, \dots, \underline{x_i + h}, \dots, x_N) - f(x_1, \dots, \underline{x_i}, \dots, x_N)}{h}$$

$$f: \mathbb{R}^N \rightarrow \mathbb{R}$$

$$y = f(x), \quad x = [x_1, \dots, x_N]^T$$

§ Partial derivative (change one coordinate at a time):

$$\frac{\partial f}{\partial x_i} = \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_{i-1}, x_i + h, x_{i+1}, \dots, x_N) - f(x)}{h}$$

§ Jacobian vector (gradient) collects all partial derivatives:

$$\frac{\partial f}{\partial x} = [\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_N}]^T \in \mathbb{R}^{1 \times N}$$

Note: This is a row vector.

# Example

$$f: \mathbf{R}^2 \rightarrow \mathbf{R}$$
$$f(x_1, x_2) = x_1^2 + x_2^2 \in \mathbf{R}$$

# Example

$$f: \mathbf{R}^2 \rightarrow \mathbf{R}$$
$$f(x_1, x_2) = x_1^2 + x_2^3 \in \mathbf{R}$$

§ Partial derivatives:

$$\frac{\partial f}{\partial x_1}(x) = 2x_1 x_2, \quad x_2^3$$

$$\frac{\partial f}{\partial x_2}(x) = x_1^2 + 3x_1 x_2^2$$

# Example

$$f: \mathbf{R}^2 \rightarrow \mathbf{R}$$

$$f(x_1, x_2) = x_1^2 + x_2^3 \in \mathbf{R}$$

§ Partial derivatives:

$$\frac{\partial f}{\partial x_1}(x) = 2x_1 x_2 + x_2^3$$

$$\frac{\partial f}{\partial x_2}(x) = x_1^2 + 3x_1 x_2^2$$

§ Gradient:

$$\frac{df}{dx} = \begin{pmatrix} \frac{\partial f}{\partial x_1}(x) \\ \frac{\partial f}{\partial x_2}(x) \end{pmatrix} = \begin{pmatrix} 2x_1 x_2 + x_2^3 \\ x_1^2 + 3x_1 x_2^2 \end{pmatrix} \in \mathbf{R}^{1 \times 2}.$$

# Rules

§ Sum Rule

$$\frac{B}{Bx} fpxq \quad gpxq \quad " \quad \frac{Bf}{Bx} \quad \frac{Bg}{Bx}$$

§ Product Rule

$$\frac{B}{Bx} fpxqgpxq \quad " \quad \frac{Bf}{Bx} gpxq \quad fpxq \quad \frac{Bg}{Bx}$$

§ Chain Rule

$$\frac{B}{Bx} pg \quad " \quad fpxq \quad " \quad \frac{B}{Bx} gp \quad fpxqq \quad " \quad \frac{Bg}{Bf} \frac{Bf}{Bx}$$

## Example: Chain Rule

§ Consider the function

$$L \text{eq} \left\{ \frac{1}{2} \|e\|^2 + e_2^T e_1 \right\}$$
$$e = y - Ax, \quad x \in \mathbb{R}^N, A \in \mathbb{R}^{M \times N}, e, y \in \mathbb{R}^M$$

§ Compute  $dL/dx$ . What is the dimension/size of  $dL/dx$ ?

# Example: Chain Rule

§ Consider the function

$$L \text{eq} \left\{ \frac{1}{2} e \right\}^2 - e_2^1 e^j \\ e = y' Ax, \quad x \in \mathbb{R}^N, A \in \mathbb{R}^{M \times N}, e, y \in \mathbb{R}^M$$

§ Compute  $dL/dx$ . What is the dimension/size of  $dL/dx$ ?

§  $dL/dx \in \mathbb{R}^{1 \times N}$

$$\frac{dL}{dx} = \frac{dL}{de} \frac{de}{dx} \\ \frac{dL}{de} \in \mathbb{R}^{1 \times M} \quad (1)$$

$$\frac{de}{dx} = A \in \mathbb{R}^{M \times N} \quad (2)$$

$$\tilde{n} \quad \frac{dL}{dx} = e^j p' A q = p'y' A x q' A \in \mathbb{R}^{1 \times N}$$

$$f: \mathbf{R}^N \tilde{\rightarrow} \mathbf{R}^M$$

$$\begin{array}{ccc} & y `` f \text{pxqP} \mathbf{R}^M, & x \mathbf{P} \mathbf{R}^N \\ \gg & y_1 \text{fi} & \gg f_1 \text{pxq} \text{fi} & \gg f_1 \text{px}_1 \dots, x_N \text{q} \text{fi} \\ \hline & \vdots \text{ffq} `` & \vdots \text{ffq} `` & \vdots \text{ffq} \\ y_M & f_M \text{pxq} & f_M \text{px}_1, \dots, x_N \text{q} \end{array}$$

$$f: \mathbf{R}^N \rightarrow \mathbf{R}^M$$

§ **Jacobian** matrix (collection of all partial derivatives)

$$\begin{array}{c} \Rightarrow \frac{dy_1}{d\bar{x}} \quad fi \quad \Rightarrow \frac{Bf_1}{Bx_1} \quad \dots \dots \quad \frac{Bf_1}{Bx_N} \quad fi \\ \hline \hline = : \quad f\ddot{f} \quad " \quad = : \quad \quad : \quad f\ddot{f} \quad PR^{M^* N} \\ \hline \hline \frac{dy_M}{d\bar{x}} \quad \quad \quad \frac{Bf_M}{Bx_1} \quad \dots \dots \quad \frac{Bf_M}{Bx_N} \end{array}$$

# Example

$$f \text{ is } "Ax, \quad f \in \mathbb{R}^M, \quad A \in \mathbb{R}^{M \times N}, \quad x \in \mathbb{R}^N"$$

§ Compute the gradient  $d f / dx$

§ Dimension of  $d f / dx$ :

## Example

$$f \text{ is } "Ax, \quad f \in \mathbb{R}^M, \quad A \in \mathbb{R}^{M \times N}, \quad x \in \mathbb{R}^N"$$

§ Compute the gradient  $df/dx$

§ Dimension of  $df/dx$ :

Since  $f: \mathbb{R}^N \rightarrow \mathbb{R}^M$ , it follows that  $df/dx \in \mathbb{R}^{M \times N}$

# Example

$$f \text{ is } "Ax, \quad f \in \mathbb{R}^M, \quad A \in \mathbb{R}^{M \times N}, \quad x \in \mathbb{R}^N"$$

§ Compute the gradient  $d f / dx$

§ Dimension of  $d f / dx$ :

Since  $f: \mathbb{R}^N \rightarrow \mathbb{R}^M$ , it follows that  $d f / dx \in \mathbb{R}^{M \times N}$

§ Gradient:

$$f_i = \sum_{j=1}^N A_{ij} x_j \quad \text{or} \quad \frac{\partial f_i}{\partial x_j} = A_{ij}$$

(3)

# Example

$$f \text{ is } "Ax, \quad f \in \mathbb{R}^M, \quad A \in \mathbb{R}^{M \times N}, \quad x \in \mathbb{R}^N"$$

§ Compute the gradient  $d f / dx$

§ Dimension of  $d f / dx$ :

Since  $f: \mathbb{R}^N \rightarrow \mathbb{R}^M$ , it follows that  $d f / dx \in \mathbb{R}^{M \times N}$

§ Gradient:

$$\begin{aligned} f_i &= \sum_{j=1}^N A_{ij} x_j \quad \text{if } \frac{\partial f_i}{\partial x_j} = A_{ij} \\ \frac{df}{dx} &= \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_N} & \text{if } \frac{\partial f_1}{\partial x_1} = A_{11} & \dots & \frac{\partial f_1}{\partial x_N} = A_{1N} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_M}{\partial x_1} & \dots & \frac{\partial f_M}{\partial x_N} & \frac{\partial f_M}{\partial x_1} = A_{M1} & \dots & \frac{\partial f_M}{\partial x_N} = A_{MN} \end{pmatrix} \quad (3) \end{aligned}$$

# Chain Rule

$$\frac{\partial}{\partial x} g(f(x)) = \frac{\partial}{\partial x} g(f) \cdot f'(x)$$

# Example

§ Consider  $f : \mathbf{R}^2 \rightarrow \mathbf{R}$ ,  $x : \mathbf{R} \rightarrow \mathbf{R}^2$

$$\begin{aligned} fpxq &= fpx_1 x_2 q \\ &\quad \ll \text{ff} \quad \ll \text{ff} \\ xptq &= x_1 ptq \quad " \quad \sin ptq \\ &\quad x_2 ptq \quad \cos ptq \end{aligned}$$

## Example

§ Consider  $f : \mathbf{R}^2 \rightarrow \mathbf{R}$ ,  $x : \mathbf{R} \rightarrow \mathbf{R}^2$

$$\begin{aligned} fpxq &= fpx_1 x_2 q = x_1^2 + 2x_1 x_2 \\ xptq &= x_1 ptq = \sin ptq \\ &\quad x_2 ptq = \cos ptq \end{aligned}$$

§ The dimensions  $df/dx$  and  $dx/dt$  are

## Example

§ Consider  $f : \mathbf{R}^2 \rightarrow \mathbf{R}$ ,  $x : \mathbf{R} \rightarrow \mathbf{R}^2$

$$\begin{aligned} fpxq &= fpx_1 x_2 q & x_1^2 & 2x_2 \\ &\quad \ll \text{ff} \quad \ll \text{ff} \\ xptq &= x_1 ptq & \sin ptq \\ &\quad x_2 ptq & \cos ptq \end{aligned}$$

§ The dimensions  $df/dx$  and  $dx/dt$  are  $1^2 2$  and  $2^1 1$ , respectively

§ Compute the gradient  $df/dt$  using the chain rule.

# Example

§ Consider  $f : \mathbf{R}^2 \rightarrow \mathbf{R}$ ,  $x : \mathbf{R} \rightarrow \mathbf{R}^2$

$$\begin{aligned} fpxq &= fpx_1 x_2 q = x_1^2 + 2x_2 \\ &\quad \text{ff} \quad \text{ff} \\ xptq &= x_1 ptq = \sin ptq \\ &\quad x_2 ptq \quad \cos ptq \end{aligned}$$

§ The dimensions  $df/dx$  and  $dx/dt$  are  $1 \times 2$  and  $2 \times 1$ , respectively

§ Compute the gradient  $df/dt$  using the chain rule.

$$\begin{aligned} \frac{df}{dt} &= \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t} \\ &= \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t} \\ &= \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t} \\ &= 2 \sin t \cdot 2 + \cos t \cdot \sin t \\ &= 2 \sin t \cos t + 2 \sin t = 2 \sin t \cos t + 1 \mathbf{q} \end{aligned}$$

**BREAK**

# Derivatives with Matrices

§ Re-cap: Gradient of a function  $f: \mathbf{R}^D \rightarrow \mathbf{R}^E$  is an  $E \times D$ -matrix: #

target dimensions  $\times$  #parameters

with

$$\frac{df}{dx} \in \mathbf{R}^{E \times D}, \quad df \text{ re, } ds = \frac{\partial f_e}{\partial x_d}$$

# Derivatives with Matrices

§ Re-cap: Gradient of a function  $f: \mathbf{R}^D \rightarrow \mathbf{R}^E$  is an  $E \times D$ -matrix: #  
target dimensions  $\times$  #parameters

with

$$\frac{df}{dx} \in \mathbf{R}^{E \times D}, \quad d\mathbf{f} = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_D} \end{pmatrix}$$

§ Generalization to cases, where the parameters ( $D$ ) or targets ( $E$ ) are matrices, apply immediately

# Derivatives with Matrices

§ Re-cap: Gradient of a function  $f: \mathbf{R}^D \rightarrow \mathbf{R}^E$  is an  $E \times D$ -matrix: #  
target dimensions  $\times$  #parameters

with

$$\frac{df}{dx} \in \mathbf{R}^{E \times D}, \quad df \in \mathbf{R}^{re, ds} \frac{\mathbf{B}f_e}{\mathbf{B}x_d}$$

§ Generalization to cases, where the parameters ( $D$ ) or targets ( $E$ ) are matrices, apply immediately

§ Assume  $f: \mathbf{R}^{M \times N} \rightarrow \mathbf{R}^{P \times Q}$ , then the gradient is a  $P \times Q \times M \times N$  object (tensor) where

$$df \in \mathbf{R}^{p, q, m, n} \frac{\mathbf{B}f_{pq}}{\mathbf{B}X_{mn}}$$

# Derivatives with Matrices: Example (1)

$$f = Ax, \quad f \in \mathbb{R}^M, A \in \mathbb{R}^{M \times N}, x \in \mathbb{R}^N$$

# Derivatives with Matrices: Example (1)

$$f = Ax, \quad f \in \mathbb{R}^M, A \in \mathbb{R}^{M \times N}, x \in \mathbb{R}^N$$

$$\frac{\partial f}{\partial A} = \frac{\partial f}{\partial A} \mathbb{R}^{M \times M \times N}$$
$$\frac{\partial f}{\partial A} = \frac{\partial f}{\partial A} \mathbb{R}^{M \times N} = \frac{\partial f}{\partial A} \mathbb{R}^{M \times N}$$
$$\frac{\partial f}{\partial A} = \frac{\partial f}{\partial A} \mathbb{R}^{M \times N}$$

## Derivatives with Matrices: Example (2)

$$f_i = \vec{y}^N A_{ij} x_j, \quad i = 1, \dots, M$$
$$\sum_j$$

(4)

## Derivatives with Matrices: Example (2)

$$f_i = \vec{y}^N A_{ij} x_j, \quad i = 1, \dots, M$$

$$\frac{\partial f_i}{\partial A_{iq}} = x_q$$

(4)

## Derivatives with Matrices: Example (2)

$$f_i = \hat{y}^N A_{ij} x_j, \quad i = 1, \dots, M$$
$$\sum_{j=1}^{N+1}$$

$$\frac{\partial f_i}{\partial A_{iq}} = x_q \quad \frac{\partial f_i}{\partial A_{i,:}} = x^J \in \mathbb{R}^{1 \times 1 \times N}$$

(4)

## Derivatives with Matrices: Example (2)

$$f_i = \hat{y}^N A_{ij} x_j, \quad i = 1, \dots, M$$

$$\frac{\partial f_i}{\partial A_{iq}} = x_q \quad \frac{\partial f_i}{\partial A_{i,:}} = x^J \mathbf{P} \mathbf{R}^{1^T 1^T N}$$

$$\frac{\partial f_i}{\partial A_{k \neq i,:}} = 0^J \mathbf{P} \mathbf{R}^{1^T 1^T N}$$

(4)

## Derivatives with Matrices: Example (2)

$$f_i = \vec{y}^N A_{ij} x_j, \quad i = 1, \dots, M$$
$$\sum_{j=1}^{N+1}$$

$$\frac{\partial f_i}{\partial A_{iq}} = x_q \quad \frac{\partial f_i}{\partial A_{i,:}} = x^J \mathbf{P} \mathbf{R}^{1 \times 1 \times N}$$

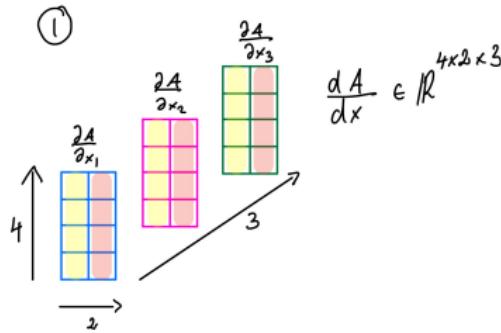
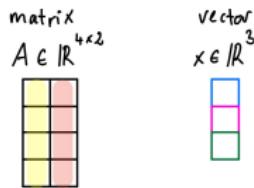
$$\frac{\partial f_i}{\partial A_{k \neq i,:}} = \mathbf{0}^J \mathbf{P} \mathbf{R}^{1 \times 1 \times N}$$
$$\Rightarrow \mathbf{0}^J f_i$$

$$\frac{\partial f_i}{\partial A} = \begin{matrix} \frac{\partial f_i}{\partial x^J} \\ \frac{\partial f_i}{\partial \mathbf{0}^J} \\ \vdots \\ \frac{\partial f_i}{\partial \mathbf{f}} \end{matrix} \mathbf{P} \mathbf{R}^{1 \times p \times M \times N \times q}$$
$$\mathbf{0}^J$$

(4)

# Example: Higher-Order Tensors

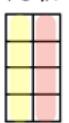
- § Consider a matrix  $A \in \mathbb{R}^{4^2}$  whose entries depend on a vector  $x \in \mathbb{R}^3$
- § We can compute  $\frac{dA}{dx}$  in two equivalent ways:



# Example: Higher-Order Tensors

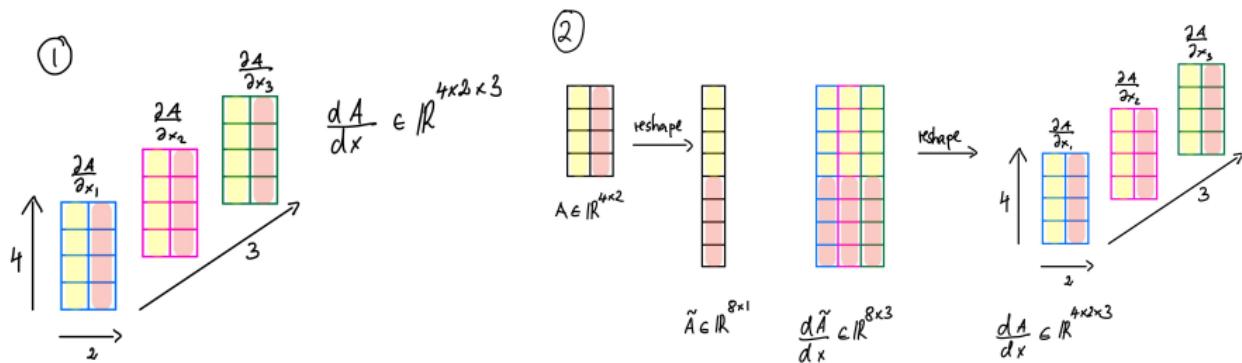
- § Consider a matrix  $A \in \mathbb{R}^{4^2}$  whose entries depend on a vector  $x \in \mathbb{R}^3$
- § We can compute  $\frac{dA}{dx}$  in two equivalent ways:

matrix  
 $A \in \mathbb{R}^{4 \times 2}$

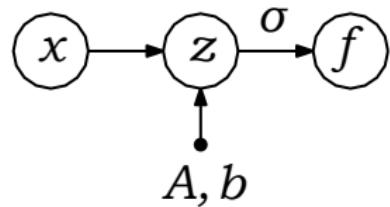


vector  
 $x \in \mathbb{R}^3$



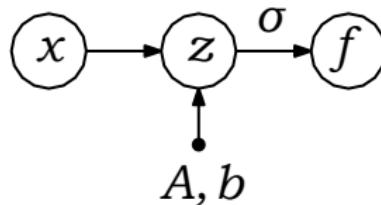


# Gradients of a Single-Layer Neural Network (1)



$f = \tanh(pA + bq)$   $\mathbf{P} \mathbf{R}^M$ ,     $x \mathbf{P} \mathbf{R}^N, A \mathbf{P} \mathbf{R}^{M \times N}, b \mathbf{P} \mathbf{R}^M$   
":  $z \mathbf{P} \mathbf{R}^M$

# Gradients of a Single-Layer Neural Network (1)



$f = \text{tanh}(Ax + b) \in \mathbb{R}^M$ ,  $x \in \mathbb{R}^N, A \in \mathbb{R}^{M \times N}, b \in \mathbb{R}^M$

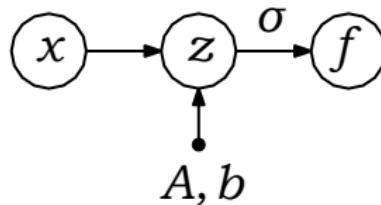
$\frac{\partial f}{\partial x}$

$$\underline{Bf} = \begin{matrix} \underline{Bf} \\ \vdots \\ \underline{Bf} \end{matrix} \in \mathbb{R}^{M \times M}$$

$$\underline{Bb} = \begin{matrix} \underline{Bb} \\ \vdots \\ \underline{Bb} \end{matrix} \in \mathbb{R}^{M \times M}$$

$$M^M \quad M^M$$

# Gradients of a Single-Layer Neural Network (1)



$f = \text{tanh}(Ax + b) \in \mathbb{R}^M$ ,  $x \in \mathbb{R}^N$ ,  $A \in \mathbb{R}^{M \times N}$ ,  $b \in \mathbb{R}^M$

$\frac{\partial f}{\partial z} \in \mathbb{R}^M$

$$\underline{Bf} = \begin{matrix} \underline{Bf} & \underline{Bz} \end{matrix} \in \mathbb{R}^{M \times M}$$

$Bb = \begin{matrix} Bz \\ Bb \end{matrix}$  ~~dot product~~

$M \times M \quad M \times M$

$$\underline{Bf} = \begin{matrix} \underline{Bf} & \underline{Bz} \end{matrix} \in \mathbb{R}^{M \times pM \times Nq}$$

$BA = \begin{matrix} Bz \\ Bn \\ Bn \\ Bn \end{matrix}$  ~~dot product~~

$M \times M \quad M \times pM \times Nq$

# Gradients of a Single-Layer Neural Network (2)

$$f = \tanh(\mathbf{A}x + \mathbf{b}) \in \mathbb{R}^M, \quad x \in \mathbb{R}^N, \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{b} \in \mathbb{R}^M$$

$$\nabla f = \begin{matrix} \underline{\mathbf{B}_f} & \underline{\mathbf{B}_z} \\ \mathbf{I} & \mathbf{B}_x \end{matrix} \in \mathbb{R}^{M \times M}$$

$\mathbf{B}_f \in \mathbb{R}^{M \times M}$     $\mathbf{B}_z \in \mathbb{R}^{M \times M}$

$\mathbf{I} \in \mathbb{R}^{M \times M}$     $\mathbf{B}_x \in \mathbb{R}^{M \times M}$

# Gradients of a Single-Layer Neural Network (2)

$$f = \tanh(pA + b) \in \mathbb{R}^M, \quad x \in \mathbb{R}^N, A \in \mathbb{R}^{M \times N}, b \in \mathbb{R}^M$$
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial z} \frac{\partial z}{\partial x} \in \mathbb{R}^M$$

$$\begin{aligned} \frac{\partial f}{\partial b} &= \frac{\partial f}{\partial z} \frac{\partial z}{\partial b} \in \mathbb{R}^{M \times M} \\ \frac{\partial b}{\partial x} &\stackrel{\text{defn}}{=} \frac{\partial z}{\partial x} \end{aligned}$$
$$\frac{\partial z}{\partial x} \in \mathbb{R}^{M \times M}$$

$$\frac{\partial f}{\partial z} = \text{diag}(1 - \tanh^2 p z q) \in \mathbb{R}^{M \times M}$$

## Gradients of a Single-Layer Neural Network (2)

$$f = \tanh(pA + b) \in \mathbb{R}^M, \quad x \in \mathbb{R}^N, A \in \mathbb{R}^{M \times N}, b \in \mathbb{R}^M$$
$$\frac{\partial f}{\partial x} = z \in \mathbb{R}^M$$

$$\begin{aligned} \frac{\partial f}{\partial b} &= \frac{\partial f}{\partial p} \frac{\partial p}{\partial b} \in \mathbb{R}^{M \times M} \\ \frac{\partial b}{\partial x} &= I \in \mathbb{R}^{M \times M} \end{aligned}$$

$$\frac{\partial f}{\partial z} = \text{diag}(1 - \tanh^2 p) q \in \mathbb{R}^{M \times M}$$

$$\frac{\partial z}{\partial b} = I \in \mathbb{R}^{M \times M}$$

$$\frac{\partial f}{\partial b}_{ri, js} = \sum_{l=1}^M \frac{\partial f}{\partial z}_{ri, ls} \frac{\partial z}{\partial b}_{rl, js}$$

(5)

## Gradients of a Single-Layer Neural Network (2)

$$f = \tanh(p_A x + b) \in \mathbb{R}^M, \quad x \in \mathbb{R}^N, A \in \mathbb{R}^{M \times N}, b \in \mathbb{R}^M$$

$$\begin{aligned} \underline{\mathbf{B}}\underline{\mathbf{f}} &= \frac{\partial f}{\partial x} = \frac{\partial z}{\partial x} \in \mathbb{R}^{M \times M} \\ \underline{\mathbf{B}}\underline{\mathbf{b}} &= \frac{\partial f}{\partial b} = \frac{\partial z}{\partial b} \in \mathbb{R}^{M \times M} \end{aligned}$$

$$\frac{\partial f}{\partial z} = \text{diag}(1 - \tanh^2 p_z q) \in \mathbb{R}^{M \times M}$$

$$\frac{\partial z}{\partial b} = I \in \mathbb{R}^{M \times M}$$

$$\frac{\partial f}{\partial b}_{ri, js} = \sum_{l=1}^M \frac{\partial f}{\partial z}_{ri, ls} \frac{\partial z}{\partial b}_{rl, js}$$

$$\text{dfdb} = \text{np.einsum('il, lj', dfdz, dzdb)}$$

# Gradients of a Single-Layer Neural Network (3)

$$f = \tanh(pA + bq) \in \mathbb{R}^M, \quad x \in \mathbb{R}^N, A \in \mathbb{R}^{M \times N}, b \in \mathbb{R}^M$$
$$\hat{z} \in \mathbb{R}^M$$

$$\frac{\partial f}{\partial A} = \frac{\partial f}{\partial p} \frac{\partial p}{\partial A} + \frac{\partial f}{\partial q} \frac{\partial q}{\partial A}$$
$$\in \mathbb{R}^{M \times M}$$
$$\frac{\partial f}{\partial p} = \frac{\partial f}{\partial z} \frac{\partial z}{\partial p} \in \mathbb{R}^{M \times N}$$
$$\frac{\partial f}{\partial q} = \frac{\partial f}{\partial z} \frac{\partial z}{\partial q} \in \mathbb{R}^{M \times N}$$

## Gradients of a Single-Layer Neural Network (3)

$$f = \tanh(pA + bq) \in \mathbb{R}^M, \quad x \in \mathbb{R}^N, A \in \mathbb{R}^{M \times N}, b \in \mathbb{R}^M$$
$$\hat{z} \in \mathbb{R}^M$$

$$\frac{\partial f}{\partial A} = \frac{\partial f}{\partial p} \frac{\partial p}{\partial A} + \frac{\partial f}{\partial q} \frac{\partial q}{\partial A}$$
$$\in \mathbb{R}^{M \times M} \quad \in \mathbb{R}^{M \times N}$$

$$\frac{\partial f}{\partial z} = \text{diag}(1 - \tanh^2(pz)) \in \mathbb{R}^{M \times M} \quad (6)$$

# Gradients of a Single-Layer Neural Network (3)

$$f = \tanh(pA + bq) \in \mathbb{R}^M, \quad x \in \mathbb{R}^N, A \in \mathbb{R}^{M \times N}, b \in \mathbb{R}^M$$
$$\hat{z} \in \mathbb{R}^M$$

$$\frac{\partial f}{\partial A} = \frac{\partial f}{\partial p} \frac{\partial p}{\partial A} + \frac{\partial f}{\partial q} \frac{\partial q}{\partial A}$$
$$= \mathbb{R}^{M \times M} \mathbb{R}^{M \times N}$$
$$= \mathbb{R}^{M \times M}$$

$$\frac{\partial f}{\partial z} = \text{diag}(1 - \tanh^2(pz)) \in \mathbb{R}^{M \times M} \quad (6)$$

$\frac{\partial z}{\partial A} \rightarrow \text{See (4)}$

$$\frac{\partial f}{\partial A}_{ri, j, ks} = \sum_{l=1}^M \frac{\partial f}{\partial z}_{ri, ls} \frac{\partial z}{\partial A}_{rl, j, ks}$$

# Gradients of a Single-Layer Neural Network (3)

$$f = \tanh(pA + bq) \in \mathbb{R}^M, \quad x \in \mathbb{R}^N, A \in \mathbb{R}^{M \times N}, b \in \mathbb{R}^M$$
$$\nabla f = \begin{matrix} \text{Bf} \\ \text{Bz} \end{matrix} \in \mathbb{R}^{M \times M}$$

$$\frac{\partial f}{\partial A} = \frac{\partial f}{\partial p} \frac{\partial p}{\partial A} + \frac{\partial f}{\partial q} \frac{\partial q}{\partial A}$$
$$\frac{\partial f}{\partial p} = \frac{\partial f}{\partial z} \frac{\partial z}{\partial p} = \frac{\partial f}{\partial z} \frac{\partial z}{\partial q} \frac{\partial q}{\partial p}$$
$$\frac{\partial f}{\partial q} = \frac{\partial f}{\partial z} \frac{\partial z}{\partial q} = \frac{\partial f}{\partial z} \frac{\partial z}{\partial p} \frac{\partial p}{\partial q}$$
$$\frac{\partial f}{\partial z} = \text{diagp1}' \tanh^2 p z q q \in \mathbb{R}^{M \times M}$$

$$\frac{\partial f}{\partial z} = \text{diagp1}' \tanh^2 p z q q \in \mathbb{R}^{M \times M} \quad (6)$$

$\frac{\partial f}{\partial A} \rightarrow \text{See (4)}$

$$\frac{\partial f}{\partial A}_{ri, j, ks} = \sum_{l=1}^M \frac{\partial f}{\partial z}_{ri, ls} \frac{\partial z}{\partial A}_{rl, j, ks}$$

$\text{dfdA} = \text{np.einsum('il, ljk', dfdz, dzdA)}$

# Putting Things Together

§ Inputs  $x$ , observed outputs  $y \sim f(pz, \theta)$   $\mid g, g_0 \sim N(0, \Sigma)$

# Putting Things Together

- § Inputs  $x$ , observed outputs  $y \sim f(pz, \theta)$   $\mid g, g_0 \sim N(0, \Sigma)$
- § Train single-layer neural network with  
 $f(pz, \theta) = \tanh(pzq), \quad z = Ax + b, \quad \theta = tA, bu$

# Putting Things Together

- § Inputs  $x$ , observed outputs  $y \sim f(pz, \theta)$   $\sim g, g \sim N(0, \Sigma)$
- § Train single-layer neural network with

$$f(pz, \theta) = \tanh(pzq), \quad z = Ax + b, \quad \theta = tA, bu$$

- § Find  $A, b$ , such that the squared loss

$$L(\theta) = \frac{1}{2} \|e\|^2, \quad e = y - f(pz, \theta)$$

is minimized

# Putting Things Together

§ Inputs  $x$ , observed outputs  $y \sim f(pz, \theta)$ ,  $\mathbf{g}, \mathbf{g} \sim N(0, \Sigma)$

§ Train single-layer neural network with

$$f(pz, \theta) = \tanh(pzq), \quad z = Ax + b, \quad \theta = tA, bu$$

§ Find  $A, b$ , such that the squared loss

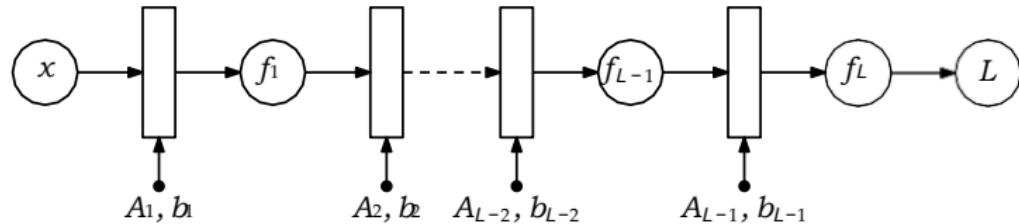
$$L(\theta) = \frac{1}{2} \|e\|^2, \quad e = y - f(pz, \theta)$$

is minimized

§ Partial derivatives:

$$\begin{array}{lll} \cancel{\frac{\partial L}{\partial \theta}} = \cancel{\frac{\partial L}{\partial p} \frac{\partial p}{\partial \theta} + \frac{\partial L}{\partial z} \frac{\partial z}{\partial \theta}} & ; & \frac{\partial L}{\partial p} \rightarrow (1) \\ \frac{\partial L}{\partial A} = \cancel{\frac{\partial L}{\partial p} \frac{\partial p}{\partial A} + \frac{\partial L}{\partial z} \frac{\partial z}{\partial A}} & ; & \frac{\partial L}{\partial p} \rightarrow (2), (3) \\ \cancel{\frac{\partial L}{\partial \theta}} = \cancel{\frac{\partial L}{\partial p} \frac{\partial p}{\partial \theta} + \frac{\partial L}{\partial z} \frac{\partial z}{\partial \theta}} & ; & \cancel{\frac{\partial L}{\partial p}} \rightarrow (6) \\ \frac{\partial L}{\partial b} = \cancel{\frac{\partial L}{\partial p} \frac{\partial p}{\partial b} + \frac{\partial L}{\partial z} \frac{\partial z}{\partial b}} & ; & \cancel{\frac{\partial L}{\partial p}} \rightarrow (5) \end{array}$$

# Gradients of a Multi-Layer Neural Network

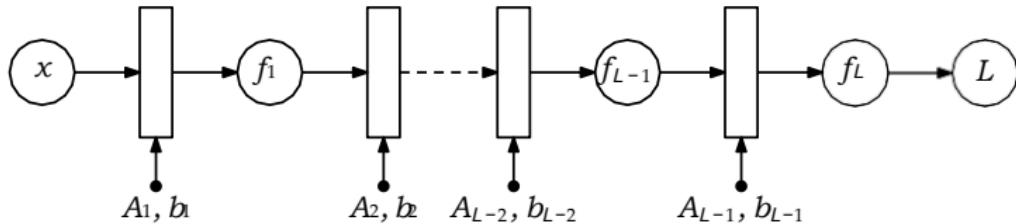


- § Inputs  $x$ , observed outputs  $y$
- § Train multi-layer neural network with

$$f_0 \leftarrow x$$

$$f_i \leftarrow \sigma_i \circ A_i \circ f_{i-1} + b_i, \quad i = 1, \dots, L$$

# Gradients of a Multi-Layer Neural Network



§ Inputs  $x$ , observed outputs  $y$

§ Train multi-layer neural network with

$$f_0 \leftarrow x$$

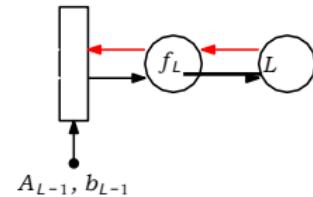
$$f_i \leftarrow \sigma_i \circ A_i \circ f_{i-1} + b_i, \quad i = 1, \dots, L$$

§ Find  $A_j, b_j$  for  $j = 0, \dots, L - 1$ , such that the squared loss

$$L(\theta) = \|y - f_L(\theta, x)\|^2$$

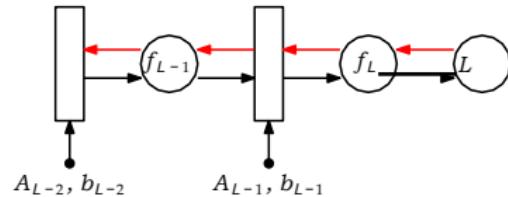
is minimized, where  $\theta = [A_j, b_j]^\top$ ,  $j = 0, \dots, L - 1$

# Gradients of a Multi-Layer Neural Network



$$\frac{\underline{\mathbf{B}}\underline{\mathbf{L}}}{\mathbf{B}\boldsymbol{\theta}_{L'1}} \quad \frac{\underline{\mathbf{B}}\underline{\mathbf{L}} \mathbf{B}\mathbf{f}_L}{\mathbf{B}\mathbf{f}_L \underline{\mathbf{B}\boldsymbol{\theta}_{L'1}}}$$

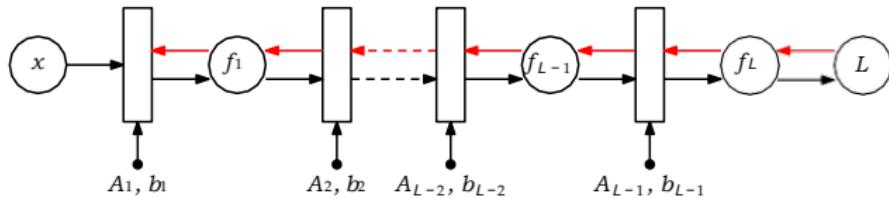
# Gradients of a Multi-Layer Neural Network



$$\frac{\underline{BL}}{\underline{B\theta_{L'1}}} = \frac{\underline{BL} \underline{Bf_L}}{\underline{Bf_L} \underline{B\theta_{L'1}}}$$

$$\frac{\underline{BL}}{\underline{B\theta_{L'2}}} = \frac{\underline{BL}}{\underline{Bf_L}} \boxed{\frac{\underline{Bf_L}}{\underline{Bf_{L'1}}} \frac{\underline{Bf_{L'1}}}{\underline{B\theta_{L'2}}}}$$

# Gradients of a Multi-Layer Neural Network

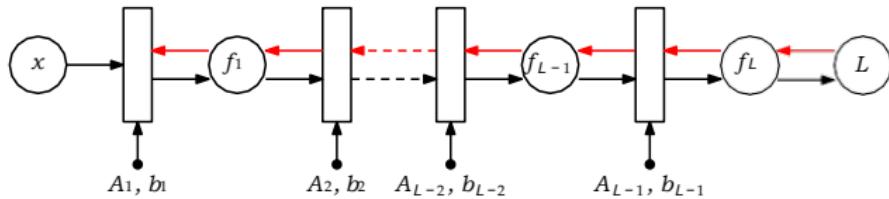


$$\frac{\underline{BL}}{\underline{B\theta_{L'1}}} = \frac{\underline{BL} \underline{Bf_L}}{\underline{Bf_L} \underline{B\theta_{L'1}}}$$

$$\frac{\underline{BL}}{\underline{B\theta_{L'2}}} = \frac{\underline{BL}}{\underline{Bf_L}} \boxed{\frac{\underline{Bf_L}}{\underline{Bf_{L'1}}} \frac{\underline{Bf_{L'1}}}{\underline{B\theta_{L'2}}}}$$

$$\frac{\underline{BL}}{\underline{B\theta_{L'3}}} = \frac{\underline{BL} \underline{Bf_L}}{\underline{Bf_L} \underline{Bf_{L'1}}} \boxed{\frac{\underline{Bf_{L'1}}}{\underline{Bf_{L'2}}} \frac{\underline{Bf_{L'2}}}{\underline{B\theta_{L'3}}}}$$

# Gradients of a Multi-Layer Neural Network



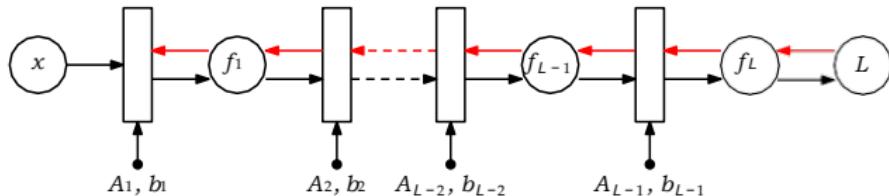
$$\frac{\partial L}{\partial \theta_{L'1}} = \frac{\partial L}{\partial f_L} \frac{\partial f_L}{\partial \theta_{L'1}}$$

$$\frac{\partial L}{\partial \theta_{L'2}} = \frac{\partial L}{\partial f_L} \frac{\partial f_L}{\partial f_{L'1}} \frac{\partial f_{L'1}}{\partial \theta_{L'2}}$$

$$\frac{\partial L}{\partial \theta_{L'3}} = \frac{\partial L}{\partial f_L} \frac{\partial f_L}{\partial f_{L'1}} \frac{\partial f_{L'1}}{\partial f_{L'2}} \frac{\partial f_{L'2}}{\partial \theta_{L'3}}$$

$$\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial f_L} \frac{\partial f_L}{\partial f_{L'1}} \dots \frac{\partial f_{i'2}}{\partial f_{i'1}} \frac{\partial f_{i'1}}{\partial \theta_i}$$

# Gradients of a Multi-Layer Neural Network



$$\frac{\underline{BL}}{\underline{B\theta_{L'1}}} = \frac{\underline{BL} Bf_L}{\underline{Bf_L} \underline{B\theta_{L'1}}}$$

$$\frac{\underline{BL}}{\underline{B\theta_{L'2}}} = \frac{\underline{BL}}{\underline{Bf_L}} \boxed{\frac{Bf_L}{Bf_{L'1}} \frac{Bf_{L'1}}{B\theta_{L'2}}}$$

$$\frac{\underline{BL}}{\underline{B\theta_{L'3}}} = \frac{\underline{BL} Bf_L}{\underline{Bf_L} \underline{Bf_{L'1}}} \boxed{\frac{Bf_{L'1}}{Bf_{L'2}} \frac{Bf_{L'2}}{B\theta_{L'3}}}$$

$$\frac{\underline{BL}}{\underline{B\theta}} = \frac{\underline{BL} Bf_L}{\underline{Bf_L} \underline{Bf_{L'1}}} \dots \boxed{\frac{Bf_{i'2}}{Bf_{i'1}} \frac{Bf_{i'1}}{B\theta}}$$

► More details (including efficient implementation) later this week

# Training Neural Networks as Maximum Likelihood Estimation

- § Training a neural network in the above way corresponds to maximum likelihood estimation:
  - § If  $y \sim \text{NN}(x, \theta)$ ,  $\theta \in \mathbb{R}^n$ ,  $\mathbb{R}^m$ ,  $N(0, I)$  then the log-likelihood is
$$\log p(y|X, \theta) = -\frac{1}{2}\|y - \text{NN}(x, \theta)\|^2$$

# Training Neural Networks as Maximum Likelihood Estimation

§ Training a neural network in the above way corresponds to **maximum likelihood estimation**:

§ If  $y \sim \text{NN}(x, \theta)$ ,  $\theta \in \mathcal{G}$ ,  $\mathcal{G} \subseteq \mathbb{R}^n$ , then the **log-likelihood** is

$$\log p(y|X, \theta) = -\frac{1}{2}\|y - \text{NN}(x, \theta)\|^2$$

§ Find  $\hat{\theta}$  by **minimizing the negative log-likelihood**:

$$\hat{\theta} = \arg \min_{\theta} -\log p(y|x, \theta)$$

$$= \arg \min_{\theta} \frac{1}{2}\|y - \text{NN}(x, \theta)\|^2$$

$$= \arg \min_{\theta} L(\theta)$$

# Training Neural Networks as Maximum Likelihood Estimation

- § Training a neural network in the above way corresponds to maximum likelihood estimation:
  - § If  $y \sim \text{NN}(x, \theta)$ ,  $\theta \in \mathbb{R}^n$ ,  $\mathbb{R}^m \rightarrow \mathbb{R}^n$ , then the log-likelihood is
$$\log p(y|X, \theta) = -\frac{1}{2}\|y - \text{NN}(x, \theta)\|^2$$

- § Find  $\hat{\theta}$  by minimizing the negative log-likelihood:

$$\hat{\theta} = \arg \min_{\theta} -\log p(y|x, \theta)$$

$$= \arg \min_{\theta} \frac{1}{2}\|y - \text{NN}(x, \theta)\|^2$$

$$= \arg \min_{\theta} L(\theta)$$

- § Maximum likelihood estimation can lead to overfitting (interpret noise as signal)

## Example: Linear Regression (1)

§ Linear regression with a polynomial of order  $M$ :

$$y \sim fpx, \theta | g, \quad g \sim N(0, \sigma_g^2)$$

$$fpx, \theta = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_M x^M = \sum_{i=0}^M \theta_i x^i$$

## Example: Linear Regression (1)

§ Linear regression with a polynomial of order  $M$ :

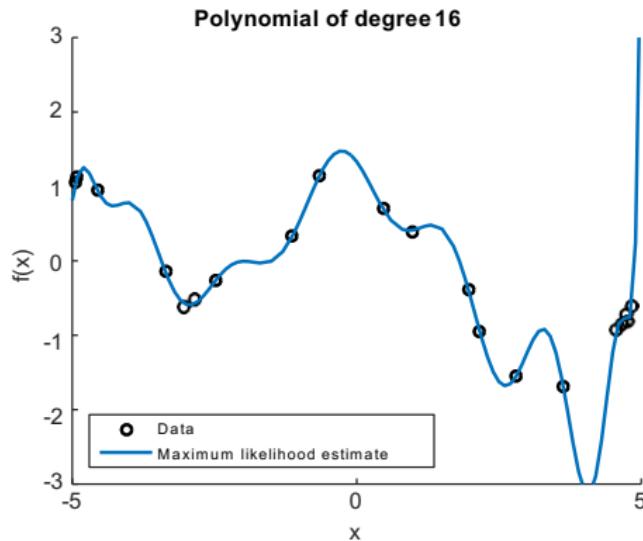
$$y \sim fpx, \theta | g, \quad g \sim N(0, \sigma_g^2)$$

$$fpx, \theta = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_M x^M = \sum_{i=0}^M \theta_i x^i$$

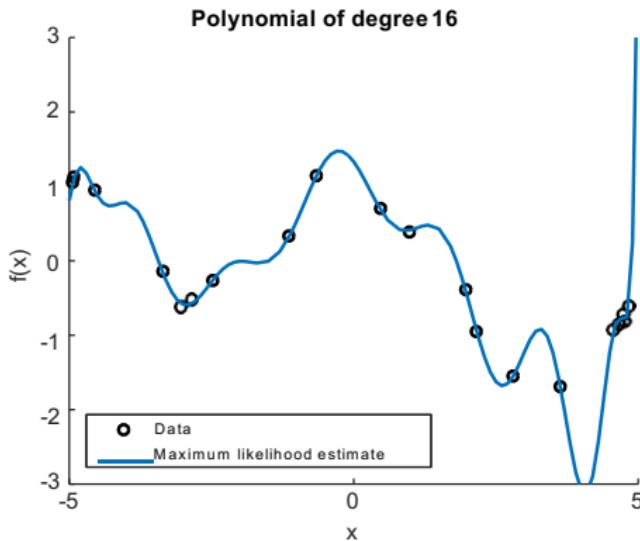
Given inputs  $x_i$  and corresponding (noisy) observations  $y_i$ ,  
 $i = 1, \dots, N$ , find parameters  $\theta = [\theta_0, \dots, \theta_M]$ , that minimize the  
squared loss (equivalently: maximize the likelihood)

$$L(\theta) = \sum_{i=1}^N (y_i - f(x_i, \theta))^2$$

## Example: Linear Regression (2)



## Example: Linear Regression (2)



- § Regularization, model selection etc. can address overfitting
  - ▶ Tutorials later this week
- § Alternative approach based on integration

# Overview

[Introduction](#)

[Differentiation](#)

[Integration](#)

# Integration: Outline

1. Motivation
2. Monte-Carlo estimation
3. Basic sampling algorithms

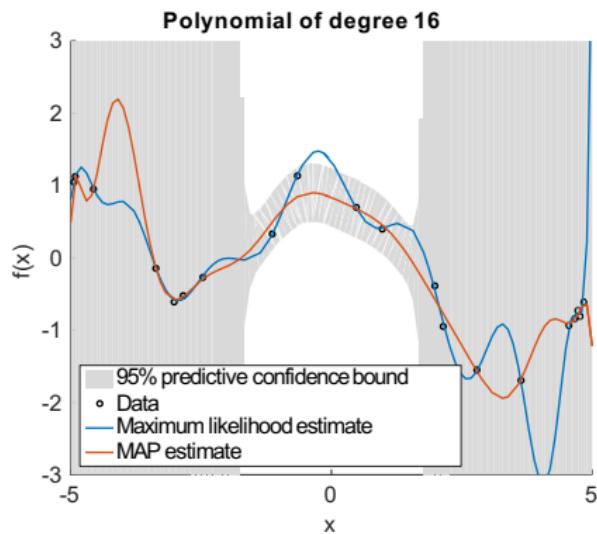
# Bayesian Integration to Avoid Overfitting

§ Instead of fitting a single set of parameters  $\theta^*$ , we can average over all plausible parameters

► Bayesian integration:

ż

$$p(y|x) \propto p(y|x, \theta) p(\theta)$$



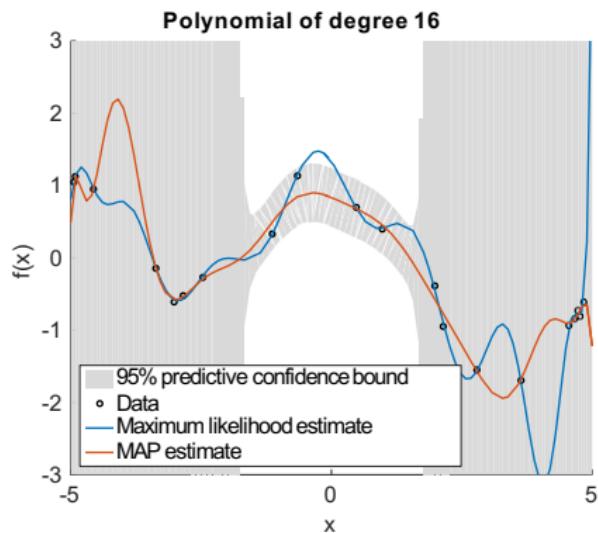
# Bayesian Integration to Avoid Overfitting

§ Instead of fitting a single set of parameters  $\theta^*$ , we can average over all plausible parameters

► Bayesian integration:

ż

$$p(\theta|x) \propto p(\theta)p(x|\theta)$$



§ More details on what  $p(\theta|x)$  is   ► Tutorials later this week

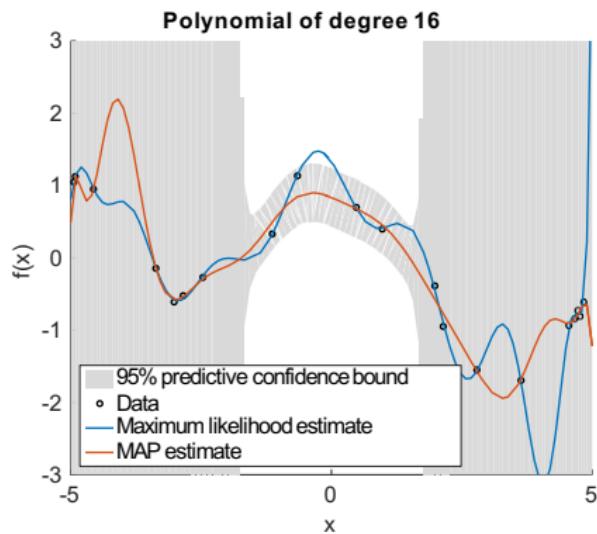
# Bayesian Integration to Avoid Overfitting

§ Instead of fitting a single set of parameters  $\theta^*$ , we can average over all plausible parameters

► Bayesian integration:

ż

$$p(\theta|x) \propto p(x|\theta)p(\theta)$$



- § More details on what  $p(\theta|x)$  is ► Tutorials later this week
- § For neural networks this integration is intractable
- Approximations

# Computing Statistics of Random Variables

§ Computing means/(co)variances also requires solving integrals:

$$\mathbf{E}_{x \sim p(x)} = \int_{-\infty}^{\infty} x p(x) dx : \mu_x$$

$$\mathbf{V}_{x \sim p(x)} = \int_{-\infty}^{\infty} (x - \mu_x)^2 p(x) dx$$

$$\text{Cov}_{x,y} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - \mu_x)(y - \mu_y) p(x, y) dx dy$$

# Computing Statistics of Random Variables

§ Computing means/(co)variances also requires solving integrals:

$$\mathbf{E}_{x \sim p(x)} = \int_{-\infty}^{\infty} x p(x) dx : \mu_x$$

$$\mathbf{V}_{x \sim p(x)} = \int_{-\infty}^{\infty} (x - \mu_x)^2 p(x) dx$$

$$\text{Cov}_{x,y} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - \mu_x)(y - \mu_y) p(x, y) dx dy$$

§ These integrals can often not be computed in closed form

► Approximations

# Approximate Integration

- § Numerical integration (low-dimensional problems)
- § Bayesian quadrature, e.g., O'Hagan (1987, 1991); Rasmussen & Ghahramani (2003)
- § Variational Bayes, e.g., Jordan et al. (1999)
- § Expectation Propagation, Opper & Winther (2001); Minka (2001)
- § Monte-Carlo Methods, e.g., Gilks et al. (1996), Robert & Casella (2004), Bishop (2006)

# Monte Carlo Methods – Motivation

- § Monte Carlo methods are computational techniques that make use of **random numbers**
- § Two typical problems:
  1. **Problem 1:** Generate samples  $x^{psq} u$  from a given probability distribution  $p(x|q)$ , e.g., for simulation or representations of data distributions

# Monte Carlo Methods – Motivation

- § Monte Carlo methods are computational techniques that make use of **random numbers**
- § Two typical problems:
  1. **Problem 1:** Generate samples  $x^{psq} u$  from a given probability distribution  $p(x)$ , e.g., for simulation or representations of data distributions
  2. **Problem 2:** Compute expectations of functions under that distribution:

$$\hat{z} = \mathbf{E}[f(x)] = \int f(x) p(x) dx$$

# Monte Carlo Methods – Motivation

- § Monte Carlo methods are computational techniques that make use of **random numbers**
- § Two typical problems:
  1. **Problem 1:** Generate samples  $x^{psq}$  from a given probability distribution  $p(x)$ , e.g., for simulation or representations of data distributions
  2. **Problem 2:** Compute expectations of functions under that distribution:

$$\mathbf{E}[f(x)] = \int f(x)p(x)dx$$

► Example: Means/variances of distributions, predictions  
**Complication:** Integral cannot be evaluated analytically

## Problem 2: Monte Carlo Estimation

§ Computing expectations via statistical sampling:

$$\mathbf{E}_{\mathbf{f} \sim p_{\text{true}}} f(\mathbf{x}) = \int p_{\text{true}}(\mathbf{x}) f(\mathbf{x}) d\mathbf{x}$$
$$\approx \frac{1}{S} \sum_{s=1}^S f(\mathbf{x}^{psq}_s), \quad \mathbf{x}^{psq}_s \sim p_{\text{sample}}$$

## Problem 2: Monte Carlo Estimation

§ Computing expectations via statistical sampling:

$$\mathbf{E}_{\theta} f(p_{\theta}(x|q)) = \int p_{\theta}(x|q) f(x) dx$$

$$\approx \frac{1}{S} \sum_{s=1}^S f(p_{\theta^{sq}}(x|q), x^{sq}) p_{\theta}(x|q)$$

§ Making predictions (e.g., Bayesian regression with inputs  $x$  and targets  $y$ )

$$p(y|x) = \frac{p(y|\theta, x)p(\theta)}{\text{norm}}$$

$$\approx \frac{1}{S} \sum_{s=1}^S p(y|\theta^{sq}, x, \theta^{sq}) p(\theta|q)$$

Parameter distribution

## Problem 2: Monte Carlo Estimation

§ Computing expectations via statistical sampling:

ż

$$\mathbf{E} f p_{\mathbf{xq}} = \int p_{\mathbf{xq}} f p_{\mathbf{xq}} d\mathbf{x}$$

$$\approx \frac{1}{S} \sum_{s=1}^S f p_{\mathbf{x}^{psq}, \mathbf{q}}, \quad \mathbf{x}^{psq}, \quad p_{\mathbf{xq}}$$

§ Making predictions (e.g., Bayesian regression with inputs  $x$  and targets  $y$ )

ż

$$p_{\mathbf{y}|\mathbf{xq}} = p_{\mathbf{y}|\theta, \mathbf{xq}}$$

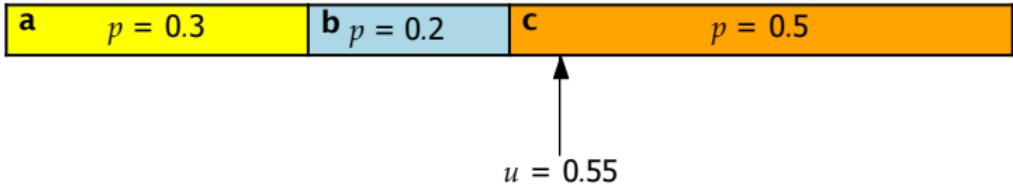
Parameter distribution

$$\approx \frac{1}{S} \sum_{s=1}^S p_{\mathbf{y}|\theta^{psq}, \mathbf{xq}}, \quad \theta^{psq}, \quad p_{\theta}$$

§ Key problem: Generating samples from  $p_{\mathbf{xq}}$  or  $p_{\theta}$

► Need to solve Problem 1

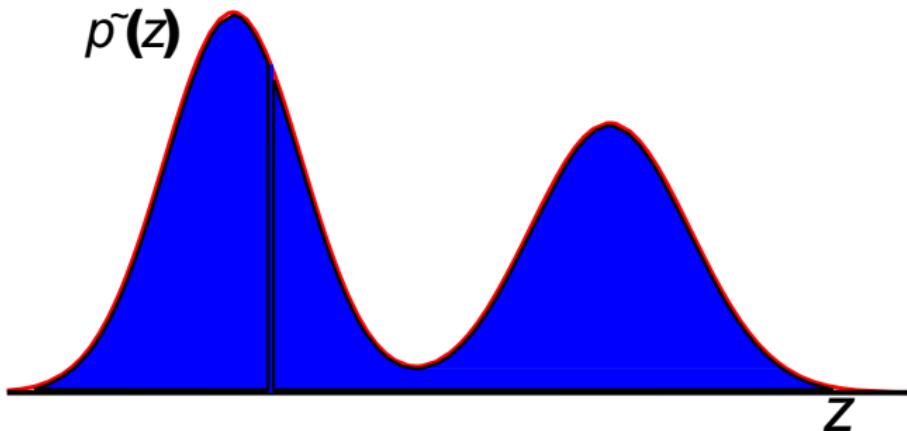
# Sampling Discrete Values



$\$ u \sim U[0, 1]$ , where  $U$  is the uniform distribution

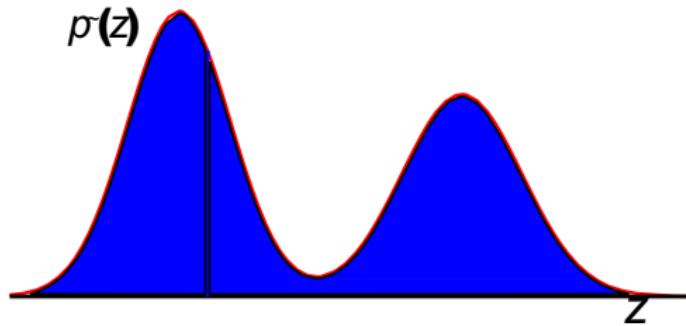
$\$ u = 0.55 \rightarrow x = c$

# Continuous Variables



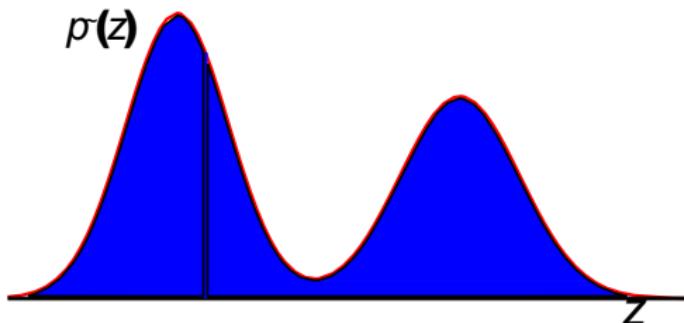
- § More complicated
- § Geometrically, we wish to sample uniformly from the area under the curve
- § Two algorithms here:
  - § Rejection sampling
  - § Importance sampling

# Rejection Sampling: Setting



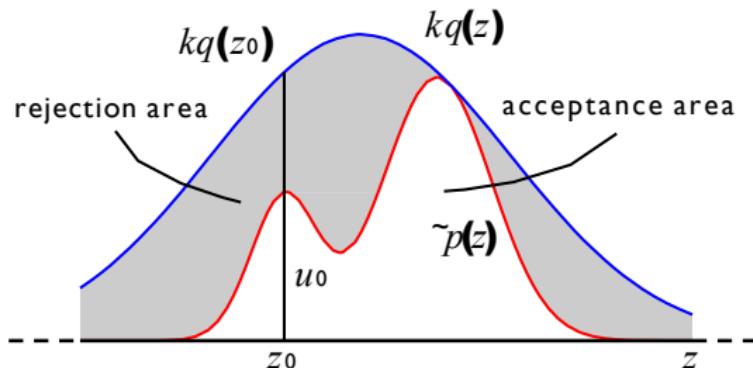
- § Assume:
  - § Sampling from  $p_{\text{pzq}}$  is difficult
  - § Evaluating  $p_{\text{pzq}}^{\sim}$  “ $Z$ ”  $p_{\text{pzq}}$  is easy (and  $Z$  may be unknown)

# Rejection Sampling: Setting



- § Assume:
  - § Sampling from  $p(z)$  is difficult
  - § Evaluating  $p(z)$  “ $Z$ ” is easy (and  $Z$  may be unknown)
- § Find a simpler distribution (**proposal distribution**)  $q(z)$  from which we can easily draw samples (e.g., Gaussian, Laplace)
- § Find an **upper bound**  $kq(z) \leq p(z)$

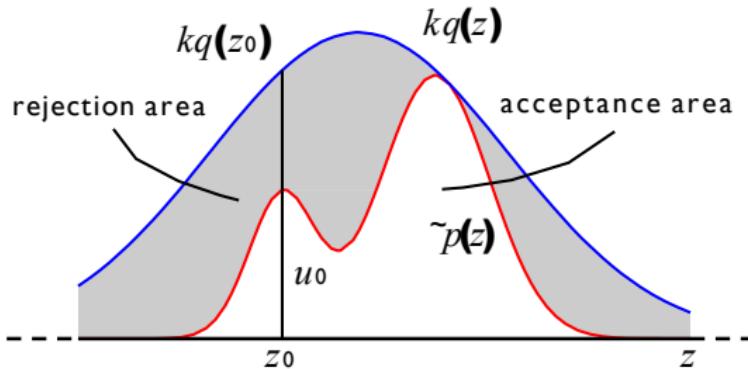
# Rejection Sampling: Algorithm



Adapted from PRML (Bishop, 2006)

1. Generate  $z_0 \sim q(z)$
2. Generate  $u_0 \sim U(0, kq(z_0))$
3. If  $u_0 \leq p(z_0)$ , reject the sample. Otherwise, retain  $z_0$

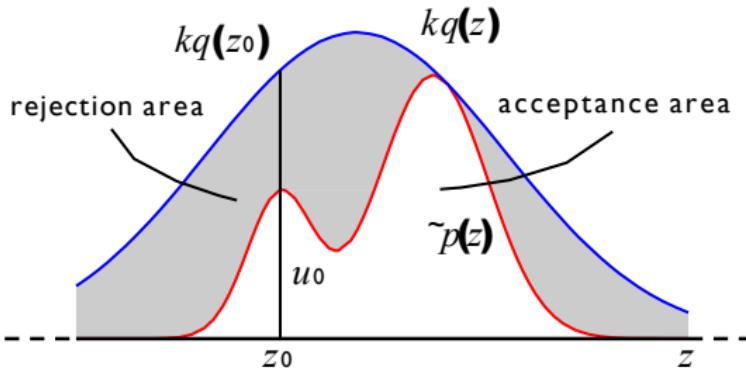
# Properties



Adapted from PRML (Bishop, 2006)

§ Accepted pairs  $p_z, u_q$  are uniformly distributed under  $\tilde{p} \tilde{p}_z q$

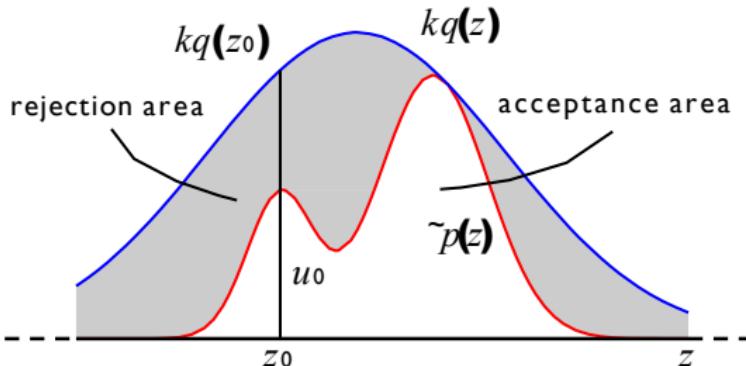
# Properties



Adapted from PRML (Bishop, 2006)

- § Accepted pairs  $p_z, u_q$  are uniformly distributed under  $\tilde{p} \tilde{p}_z q$
- § Probability density of the  $z$ -coordinates of accepted points must be proportional to  $p p_z q$

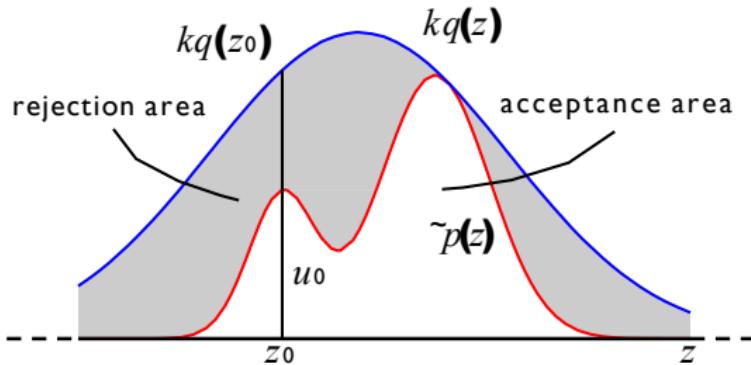
# Properties



Adapted from PRML (Bishop, 2006)

- § Accepted pairs  $p_z, u_q$  are uniformly distributed under  $\tilde{p}(z|q)$
- § Probability density of the  $z$ -coordinates of accepted points must be proportional to  $p(p_z|q)$
- § Samples are independent samples from  $p(p_z|q)$

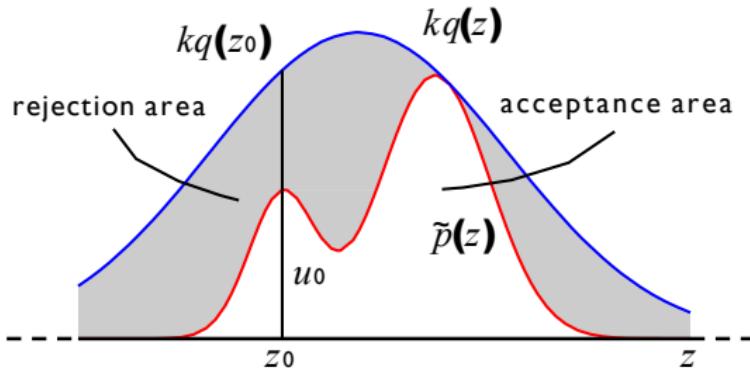
# Shortcomings



Adapted from PRML (Bishop, 2006)

§ Finding the upper bound  $k$  is tricky

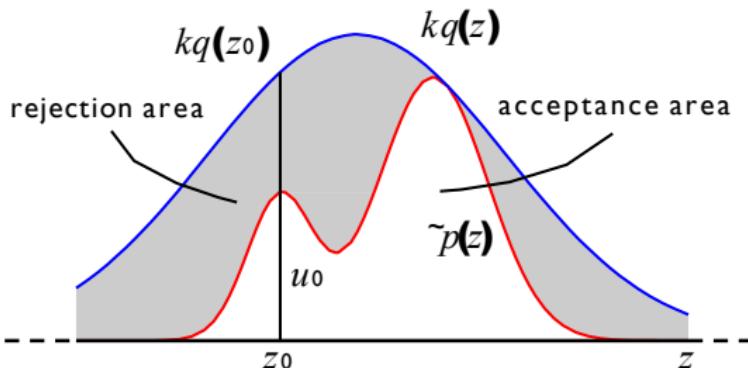
# Shortcomings



Adapted from PRML (Bishop, 2006)

- § Finding the upper bound  $k$  is tricky
- § In high dimensions the factor  $k$  is probably huge

# Shortcomings



Adapted from PRML (Bishop, 2006)

- § Finding the upper bound  $k$  is tricky
- § In high dimensions the factor  $k$  is probably huge
- § **Low acceptance rate/high rejection rate** of samples

# Importance Sampling

**Key idea:** Do not throw away all rejected samples, but give them lower weight by rewriting the integral as an expectation under a simpler distribution  $q$  (**proposal distribution**):

$$\mathbf{E}_{p \mid f} f(x) q(x) = \int p(x) q(p(x)) f(x) dx$$

# Importance Sampling

**Key idea:** Do not throw away all rejected samples, but give them lower weight by rewriting the integral as an expectation under a simpler distribution  $q$  (**proposal distribution**):

$$\begin{aligned} \mathbf{E}_p r f p | q s &= \int_{\mathbb{R}} f p | q p p | q d x \\ &= \int_{\mathbb{R}} f p | q p p | q \frac{q}{q} d x \end{aligned}$$

# Importance Sampling

**Key idea:** Do not throw away all rejected samples, but give them lower weight by rewriting the integral as an expectation under a simpler distribution  $q$  (**proposal distribution**):

$$\mathbf{E}_p r f p x q s = \int_{\mathbb{R}} f p x q p p x q d x$$
$$= \int_{\mathbb{R}} f p x q p p x q \frac{q p x q}{q p x q} d x = \int_{\mathbb{R}} f p x q \frac{p p x q}{q p x q} q p x q d x$$

# Importance Sampling

**Key idea:** Do not throw away all rejected samples, but give them lower weight by rewriting the integral as an expectation under a simpler distribution  $q$  (**proposal distribution**):

$$\begin{aligned} & \mathbf{E}_p \left[ f(p|x) q(x) \right] = \int p(x) f(p|x) q(x) dx \\ &= \int p(x) \frac{f(p|x)}{\frac{f(p|x)}{q(x)}} q(x) dx = \int p(x) \frac{p(x)}{q(x)} q(x) dx \\ &= \mathbf{E}_q \left[ f(p|x) \frac{p(x)}{q(x)} \right] \end{aligned}$$

# Importance Sampling

**Key idea:** Do not throw away all rejected samples, but give them lower weight by rewriting the integral as an expectation under a simpler distribution  $q$  (**proposal distribution**):

$$\begin{aligned} & \mathbf{E}_p r f p(x) q(x) dx \\ &= \mathbf{E}_p \frac{f p(x) q(x)}{f p(x) q(x)} dx = \mathbf{E}_q \frac{f p(x) \cancel{q(x)}}{\cancel{f p(x) q(x)}} dx \\ &= \mathbf{E}_q \frac{f p(x)}{q(x)} dx \end{aligned}$$

If we choose  $q$  in a way that we can easily sample from it, we can approximate this last expectation by Monte Carlo:

$$E_q \frac{f p(x)}{q(x)} \stackrel{\square}{\ll} \frac{1}{S} \sum_{s=1}^S \frac{f p(x^{sq})}{q(x^{sq})}, \quad x^{sq} \sim q(x)$$

# Importance Sampling

**Key idea:** Do not throw away all rejected samples, but give them lower weight by rewriting the integral as an expectation under a simpler distribution  $q$  (**proposal distribution**):

$$\begin{aligned} \mathbf{E}_p [f p(x)] &= \int f p(x) q(p(x)) dx \\ &= \int f p(x) q(p(x)) \frac{q(p(x))}{q(p(x))} dx = \int f p(x) \frac{p(x)}{q(p(x))} q(p(x)) dx \\ &= \mathbf{E}_q \left[ f p(x) \frac{p(x)}{q(p(x))} \right] \end{aligned}$$

If we choose  $q$  in a way that we can easily sample from it, we can approximate this last expectation by Monte Carlo:

$$E_q \left[ f p(x) \frac{p(x)}{q(p(x))} \right] \approx \frac{1}{S} \sum_{s=1}^S f p(x^{psq}) \frac{p(x^{psq})}{q(x^{psq})}$$

“ $w_s f p(x^{psq})$ ,  $x^{psq}$ ”,  $q(x^{psq})$

# Properties

§ Unbiased if  $q \neq 0$  where  $p \neq 0$  and if we can evaluate  $p$

# Properties

- § Unbiased if  $q \neq 0$  where  $p \neq 0$  and if we can evaluate  $p$
- § Breaks down if we do not have enough samples (puts nearly all weight on a single sample)
  - Degeneracy, see also Particle Filtering and SMC  
(Thrun et al., 2005; Doucet et al., 2000)

# Properties

- § Unbiased if  $q \neq 0$  where  $p \neq 0$  and if we can evaluate  $p$
- § Breaks down if we do not have enough samples (puts nearly all weight on a single sample)
  - Degeneracy, see also Particle Filtering and SMC  
(Thrun et al., 2005; Doucet et al., 2000)
- § Many draws from proposal density  $q$  required, especially in high dimensions

# Properties

- § Unbiased if  $q \neq 0$  where  $p \neq 0$  and if we can evaluate  $p$
- § Breaks down if we do not have enough samples (puts nearly all weight on a single sample)
  - Degeneracy, see also Particle Filtering and SMC  
(Thrun et al., 2005; Doucet et al., 2000)
- § Many draws from proposal density  $q$  required, especially in high dimensions
- § Requires to be able to evaluate true  $p$ . Generalization exists for  $\tilde{p}$ .  
This generalization is biased (but consistent).

# Properties

- § Unbiased if  $q \neq 0$  where  $p \neq 0$  and if we can evaluate  $p$
- § Breaks down if we do not have enough samples (puts nearly all weight on a single sample)
  - Degeneracy, see also Particle Filtering and SMC  
(Thrun et al., 2005; Doucet et al., 2000)
- § Many draws from proposal density  $q$  required, especially in high dimensions
- § Requires to be able to evaluate true  $p$ . Generalization exists for  $\tilde{p}$ .  
This generalization is biased (but consistent).
- § Does not scale to interesting (high-dimensional) problems

# Properties

- § Unbiased if  $q \neq 0$  where  $p \neq 0$  and if we can evaluate  $p$
  - § Breaks down if we do not have enough samples (puts nearly all weight on a single sample)
    - Degeneracy, see also Particle Filtering and SMC  
(Thrun et al., 2005; Doucet et al., 2000)
  - § Many draws from proposal density  $q$  required, especially in high dimensions
  - § Requires to be able to evaluate true  $p$ . Generalization exists for  $\tilde{p}$ .  
This generalization is biased (but consistent).
  - § Does not scale to interesting (high-dimensional) problems
- Different approach to sample from complicated (high-dimensional) distributions: Markov Chain Monte Carlo (e.g., Gilks et al., 1996)

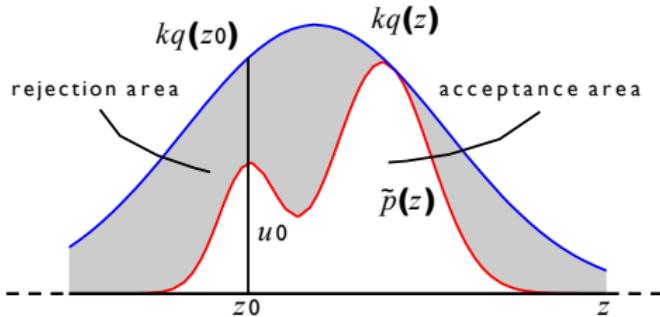
# Summary

matrix  
 $A \in \mathbb{R}^{4 \times 2}$

vector  
 $x \in \mathbb{R}^3$

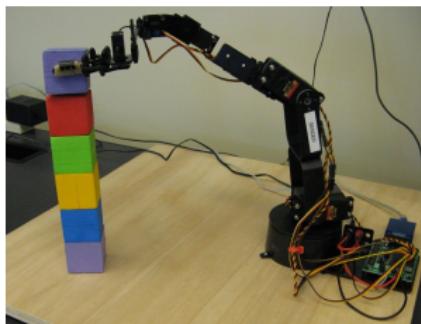
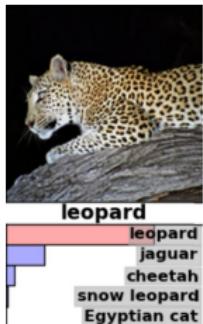
①

$$\frac{dA}{dx} \in \mathbb{R}^{4 \times 2 \times 3}$$



- § Two mathematical challenges in machine learning
  - § **Differentiation** for optimizing parameters of machine learning models
    - Vector calculus and chain rule
  - § **Integration** for computing statistics (e.g., means, variances) and as a principled way to address overfitting issue
    - Monte-Carlo integration to solve intractable integrals

# Some Application Areas



- § **Image/speech/text/language processing** using deep neural networks (e.g., Krizhevsky et al., 2012 or overview in Goodfellow et al., 2016)
- § **Data-efficient reinforcement learning and robot learning** using Gaussian processes (e.g., Deisenroth & Rasmussen, 2011)
- § **High-energy physics** using deep neural networks or Gaussian processes (e.g., Sadowski et al. 2014; Bertone et al., 2016)

# References I

- 1 G. Bertone, M. P. Deisenroth, J. S. Kim, S. Liem, R. R. de Austri, and M. Welling. Accelerating the BSM Interpretation of LHC Data with Machine Learning. arXiv preprint arXiv:1611.02704, 2016.
- 2 C. M. Bishop. *Pattern Recognition and Machine Learning*. Information Science and Statistics. Springer-Verlag, 2006.
- 3 M. P. Deisenroth, D. Fox, and C. E. Rasmussen. Gaussian Processes for Data-Efficient Learning in Robotics and Control. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(2):408–423, Feb. 2015.
- 4 M. P. Deisenroth and C. E. Rasmussen. PILCO: A Model-Based and Data-Efficient Approach to Policy Search. In *Proceedings of the International Conference on Machine Learning*, pages 465–472. ACM, June 2011.
- 5 A. Doucet, S. J. Godsill, and C. Andrieu. On Sequential Monte Carlo Sampling Methods for Bayesian Filtering. *Statistics and Computing*, 10:197–208, 2000.
- 6 W. R. Gilks, S. Richardson, and D. J. Spiegelhalter, editors. *Markov Chain Monte Carlo in Practice: Interdisciplinary Statistics*. Chapman & Hall, 1996.
- 7 I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT press, 2016.
- 8 M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul. An Introduction to Variational Methods for Graphical Models. *Machine Learning*, 37:183–233, 1999.
- 9 S. Kamthe and M. P. Deisenroth. Data-Efficient Reinforcement Learning with Probabilistic Model Predictive Control. *arXiv:1706.06491*, abs/1706.06491, 2017.
- 10 A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet Classification with Deep Convolutional Neural Networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- 11 T. P. Minka. *A Family of Algorithms for Approximate Bayesian Inference*. PhD thesis, Massachusetts Institute of Technology, Cambridge, MA, USA, Jan. 2001.
- 12 R. M. Neal. *Bayesian Learning for Neural Networks*. PhD thesis, Department of Computer Science, University of Toronto, 1996.
- 13 A. O'Hagan. Monte Carlo is Fundamentally Unsound. *The Statistician*, 36(2/3):247–249, 1987.

## References II

- 14 A. O'Hagan. Bayes-Hermite Quadrature. *Journal of Statistical Planning and Inference*, 29:245–260, 1991.
- 15 M. Opper and O. Winther. Adaptive and Self-averaging Thouless-Anderson-Palmer Mean-field Theory for Probabilistic Modeling. *Physical Review E: Statistical, Nonlinear, and Soft Matter Physics*, 64:056131, 2001.
- 16 C. E. Rasmussen and Z. Ghahramani. Bayesian Monte Carlo. In S. Becker, S. Thrun, and K. Obermayer, editors, *Advances in Neural Information Processing Systems 15*, pages 489–496. The MIT Press, Cambridge, MA, USA, 2003.
- 17 C. P. Robert and G. Casella. *Monte Carlo Methods*. Wiley Online Library, 2004.
- 18 P. Sadowski, J. Collado, D. Whiteson, and P. Baldi. Deep Learning, Dark Knowledge, and Dark Matter. In *NIPS Workshop on High-energy Physics and Machine Learning*, pages 81–87, 2014.
- 19 S. Thrun, W. Burgard, and D. Fox. *Probabilistic Robotics*. The MIT Press, Cambridge, MA, USA, 2005.