

Slides: Marc Deisenroth (Imperial College London/Prowler.io), 2017

Mathematics for Machine Learning

Avishkar Bhoopchand, *DeepMind*

avishkar@google.com

Cynthia Mulenga, *Mwabu*

cynthiammulenga@gmail.com

Daniela Massiceti, *Univ Oxford*

daniela@robots.ox.ac.uk

Kendi Muchungi, *Africa Nazarenne Univ*

kmuchungi@anu.ac.ke



DEEP
LEARNING
INDABA

Stellenbosch University

9th September 2018

Before we begin...

Introduce yourselves to your neighbours!

(you'll be working with them throughout this lecture)

Before we begin...

Colab notebook with all examples in this lecture:

[https://colab.research.google.com/github/
deep-learning-indaba/indaba-2018/blob/master/
Mathematics_for_Machine_Learning_Examples.
ipynb](https://colab.research.google.com/github/deep-learning-indaba/indaba-2018/blob/master/Mathematics_for_Machine_Learning_Examples.ipynb)

What is Machine Learning?

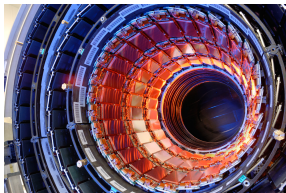
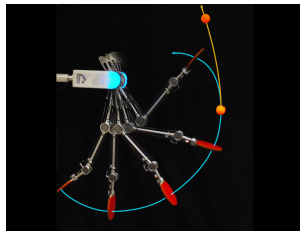
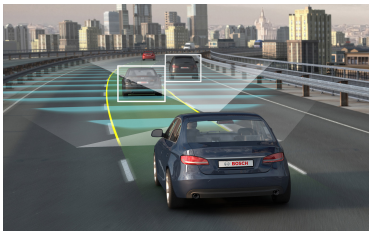
Improving a machine's performance at some task with experience

- ▶ Define task **T** - classification, regression
- ▶ Establish performance measure **P** - top-5/top-1 classification error rate, Mean Square Error (MSE)
- ▶ Provide some experience **E** - training examples

Applications of Machine Learning



leopard



Today's Recommendations For You

Here's a daily sample of items recommended for you. Click here to [see all recommendations](#)

Data Science in Python
by Steve Doulos
\$29.99

Fix this recommendation

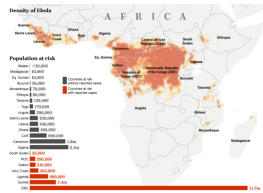
Simply JavaScript (Paperback)
by Kevin Yank
\$26.97

Fix this recommendation

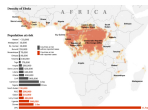
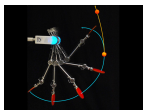
The Art & Science of JavaScript (Paperback)
\$29.99

Fix this rec

Any Category | Algorithms | Boxed Sets | Business & Culture | Java | Graphic Design | Microsoft | Networking | Networks, Protocols & APIs | New | SQL

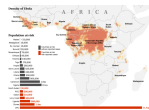


Mathematical Concepts in Machine Learning



- ▶ Probability theory and Bayesian inference
- ▶ Linear algebra and matrix decomposition
- ▶ Differentiation
- ▶ Optimisation
- ▶ Integration
- ▶ Functional analysis

Mathematical Concepts in Machine Learning



- ▶ **Probability theory** and Bayesian inference
- ▶ Linear algebra and matrix decomposition
- ▶ **Differentiation**
- ▶ Optimisation
- ▶ Integration
- ▶ Functional analysis

Outline

Key Concepts in Probability

- Probability density function

- Sum & Product Rules

- Bayes' Theorem

Differentiation

- Where it is used in ML?

- Types of differentiation

- Gradients in Neural Networks

Overview

Key Concepts in Probability

- Probability density function

- Sum & Product Rules

- Bayes' Theorem

Differentiation

- Where it is used in ML?

- Types of differentiation

- Gradients in Neural Networks

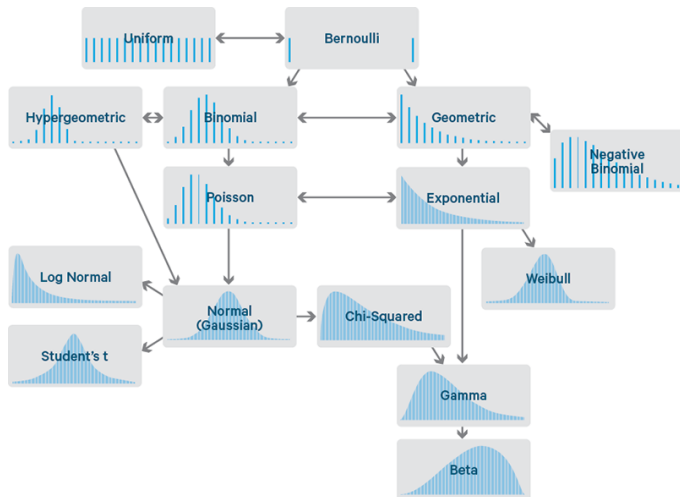
Probability density functions

A function $p : \mathbb{R}^D \rightarrow \mathbb{R}$ is called a probability density function if

1. Its integral exists,
2. $\forall \mathbf{x} \in \mathbb{R}^D : p(\mathbf{x}) \geq 0$ and
3. $\int_{\mathbb{R}^D} p(\mathbf{x}) d\mathbf{x} = 1$ *

* Here, $\mathbf{x} \in \mathbb{R}^D$ is a (continuous) random variable. For discrete random variables, the integral is replaced with a sum.

Probability density functions



<http://www.math.wm.edu/~leemis/chart/UDR/UDR.html>

Sum & Product Rules

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

Sum rule/Marginalization property

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

Product rule

- ▶ $p(\mathbf{x}, \mathbf{y})$ is the joint probability distribution of \mathbf{x} and \mathbf{y}
- ▶ $p(\mathbf{x})$ and $p(\mathbf{y})$ are the marginal probability distributions of \mathbf{x} and \mathbf{y}
- ▶ $p(\mathbf{y}|\mathbf{x})$ is the conditional probability distribution of \mathbf{y} given \mathbf{x}

Sum & Product Rules

$$p(x) = \int p(x, y) dy$$

Sum rule/Marginalization property

$$p(x, y) = p(y|x)p(x)$$

Product rule

- ▶ $p(x, y)$ is the joint probability distribution of x and y
- ▶ $p(x)$ and $p(y)$ are the marginal probability distributions of x and y
- ▶ $p(y|x)$ is the conditional probability distribution of y given x

Sum & Product Rules

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

Sum rule/Marginalization property

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

Product rule

- ▶ $p(\mathbf{x}, \mathbf{y})$ is the joint probability distribution of \mathbf{x} and \mathbf{y}
- ▶ $p(\mathbf{x})$ and $p(\mathbf{y})$ are the marginal probability distributions of \mathbf{x} and \mathbf{y}
- ▶ $p(\mathbf{y}|\mathbf{x})$ is the conditional probability distribution of \mathbf{y} given \mathbf{x}

Sum & Product Rules

$$p(x) = \int p(x, y) dy$$

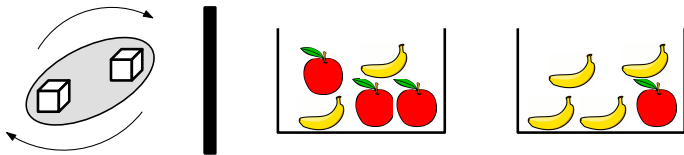
Sum rule/Marginalization property

$$p(x, y) = p(y|x)p(x)$$

Product rule

- ▶ $p(x, y)$ is the joint probability distribution of x and y
- ▶ $p(x)$ and $p(y)$ are the marginal probability distributions of x and y
- ▶ $p(y|x)$ is the conditional probability distribution of y given x

Illustration: Bayesian Inference

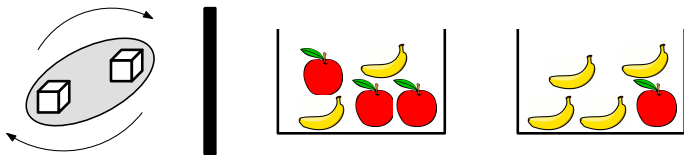


Two cartons $\mathbf{x} = \{c_1, c_2\}$ with fruit in each $\mathbf{y} = \{\text{banana, apple}\}$

Say we pick a banana ($\mathbf{y} = \text{banana}$)

We want to estimate $p(\mathbf{x} = c_1 | \mathbf{y} = \text{banana})$

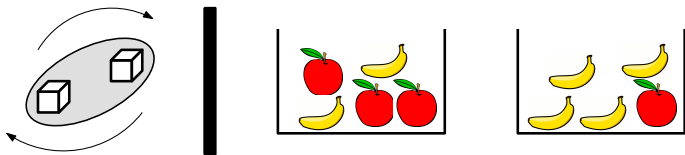
Illustration: Bayesian Inference



$$p(c_1) = 1/2 = p(c_2)$$

Prior

Illustration: Bayesian Inference

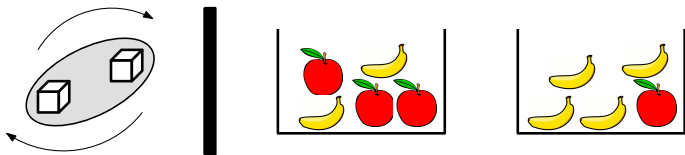


$$p(c_1) = 1/2 = p(c_2)$$

$$p(\text{banana}|c_1) = 2/5, \quad p(\text{banana}|c_2) = 4/5$$

Prior
Likelihood

Illustration: Bayesian Inference



$$p(c_1) = 1/2 = p(c_2)$$

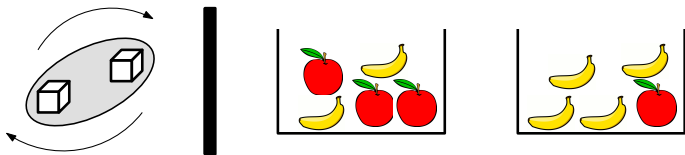
$$p(\text{banana}|c_1) = 2/5, \quad p(\text{banana}|c_2) = 4/5$$

$$p(\text{banana}) = \sum_x p(\mathbf{x}, \mathbf{y}) = \sum_x p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) =$$

$$p(\text{banana}|c_1)p(c_1) + p(\text{banana}|c_2)p(c_2) = 6/10 = 3/5$$

Prior
Likelihood
Evidence

Illustration: Bayesian Inference



$$p(c_1) = 1/2 = p(c_2)$$

$$p(\text{banana}|c_1) = 2/5, \quad p(\text{banana}|c_2) = 4/5$$

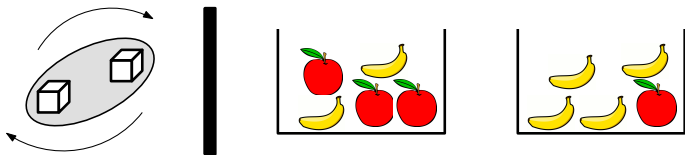
$$p(\text{banana}) = \sum_x p(\mathbf{x}, \mathbf{y}) = \sum_x p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) =$$

$$p(\text{banana}|c_1)p(c_1) + p(\text{banana}|c_2)p(c_2) = 6/10 = 3/5$$

What is $p(\mathbf{x} = c_1|\mathbf{y} = \text{banana})$?

Prior
Likelihood
Evidence
Posterior

Illustration: Bayesian Inference



$$p(c_1) = 1/2 = p(c_2)$$

$$p(\text{banana}|c_1) = 2/5, \quad p(\text{banana}|c_2) = 4/5$$

$$p(\text{banana}) = \sum_x p(\mathbf{x}, \mathbf{y}) = \sum_x p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) =$$

$$p(\text{banana}|c_1)p(c_1) + p(\text{banana}|c_2)p(c_2) = 6/10 = 3/5$$

What is $p(x = c_1 | y = \text{banana})$?

Prior

Likelihood

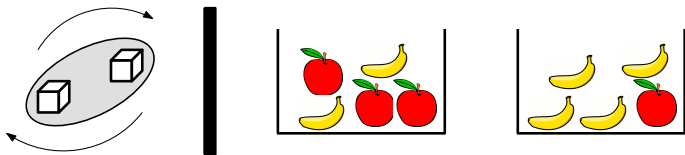
Evidence

Posterior

Bayes' Theorem:

$$\text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$$

Illustration: Bayesian Inference



$$p(c_1) = 1/2 = p(c_2)$$

$$p(\text{banana}|c_1) = 2/5, \quad p(\text{banana}|c_2) = 4/5$$

$$p(\text{banana}) = \sum_x p(\mathbf{x}, \mathbf{y}) = \sum_x p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) =$$

$$p(\text{banana}|c_1)p(c_1) + p(\text{banana}|c_2)p(c_2) = 6/10 = 3/5$$

What is $p(\mathbf{x} = c_1 | \mathbf{y} = \text{banana})$?

Prior

Likelihood

Evidence

Posterior

Bayes' Theorem:

$$\text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$$

$$p(\mathbf{x} = c_1 | \mathbf{y} = \text{banana}) = \frac{p(\text{banana}|c_1)p(c_1)}{p(\text{banana})} = \frac{2/5 * 1/2}{3/5} = 1/3$$

Bayes' Theorem

Bayes' Theorem (Probabilistic Inverse)

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x}) p(\mathbf{x})}{p(\mathbf{y})}, \quad \mathbf{x} : \text{hypothesis}, \quad \mathbf{y} : \text{measurement}$$

- ▶ Posterior belief
- ▶ Prior belief
- ▶ Likelihood (measurement model)
- ▶ Marginal likelihood (normalisation constant)



Overview

Key Concepts in Probability

Probability density function

Sum & Product Rules

Bayes' Theorem

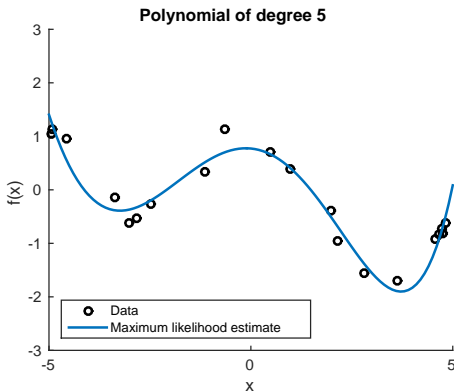
Differentiation

Where it is used in ML?

Types of differentiation

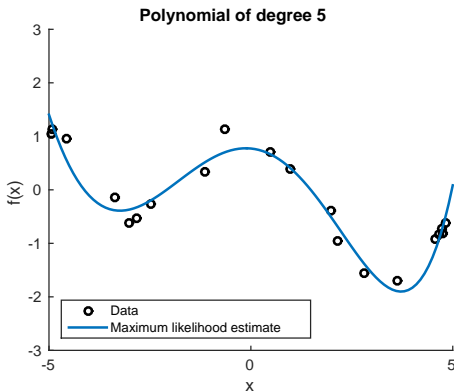
Gradients in Neural Networks

Scenario 1: Find optimal parameters to fit a 1D curve



- **Given:** N data pairs (x_i, y_i) (black dots) where x_i is input & y_i is output/target

Scenario 1: Find optimal parameters to fit a 1D curve



- ▶ **Given:** N data pairs (x_i, y_i) (black dots) where x_i is input & y_i is output/target
- ▶ **Goal:** Learn a model f (blue curve) that best fits the data

$$f(x_i) \approx y_i$$

Scenario 1: Find optimal parameters to fit a 1D curve

- ▶ Model f is D-degree polynomial parametrised by θ , or f_{θ}

$$f_{\theta}(x) = \theta_0x^0 + \theta_1x^1 + \theta_2x^2 + \dots + \theta_{D-1}x^{D-1} + \theta_Dx^D$$

Scenario 1: Find optimal parameters to fit a 1D curve

- ▶ Model f is D-degree polynomial parametrised by θ , or f_{θ}

$$f_{\theta}(x) = \theta_0 x^0 + \theta_1 x^1 + \theta_2 x^2 + \dots + \theta_{D-1} x^{D-1} + \theta_D x^D$$

$$= \begin{bmatrix} \theta_0 & \theta_1 & \theta_2 & \dots & \theta_{D-1} & \theta_D \end{bmatrix} \begin{bmatrix} 1 \\ x \\ x^2 \\ \dots \\ x^{D-1} \\ x^D \end{bmatrix}$$

Scenario 1: Find optimal parameters to fit a 1D curve

- Model f is D-degree polynomial parametrised by θ , or f_{θ}

$$f_{\theta}(x) = \theta_0 x^0 + \theta_1 x^1 + \theta_2 x^2 + \dots + \theta_{D-1} x^{D-1} + \theta_D x^D$$

$$= \begin{bmatrix} \theta_0 & \theta_1 & \theta_2 & \dots & \theta_{D-1} & \theta_D \end{bmatrix} \begin{bmatrix} 1 \\ x \\ x^2 \\ \dots \\ x^{D-1} \\ x^D \end{bmatrix}$$

$$= \theta^{\top} \mathbf{x} \quad \mathbf{x}, \theta \in \mathbb{R}^{D+1}$$

Scenario 1: Find optimal parameters to fit a 1D curve

- ▶ Model f is D-degree polynomial parametrised by θ , or f_{θ}

$$f_{\theta}(x) = \theta_0 x^0 + \theta_1 x^1 + \theta_2 x^2 + \dots + \theta_{D-1} x^{D-1} + \theta_D x^D$$

$$= \begin{bmatrix} \theta_0 & \theta_1 & \theta_2 & \dots & \theta_{D-1} & \theta_D \end{bmatrix} \begin{bmatrix} 1 \\ x \\ x^2 \\ \dots \\ x^{D-1} \\ x^D \end{bmatrix}$$

$$= \theta^{\top} \mathbf{x} \quad \mathbf{x}, \theta \in \mathbb{R}^{D+1}$$

- ▶ Training data: N pairs (x_i, y_i)
- ▶ Training the model means finding parameters θ^* , such that

$$f_{\theta^*}(x_i) \approx y_i$$

Scenario 1: Find optimal parameters to fit a 1D curve

- ▶ Model f is D -dimensional polynomial parametrised by θ , or f_θ

$$f_\theta(x) = \theta_0 x^0 + \theta_1 x^1 + \theta_2 x^2 + \dots + \theta_{D-1} x^{D-1} + \theta_D x^D$$

$$= \begin{bmatrix} \theta_0 & \theta_1 & \theta_2 & \dots & \theta_{D-1} & \theta_D \end{bmatrix} \begin{bmatrix} 1 \\ x \\ x^2 \\ \dots \\ x^{D-1} \\ x^D \end{bmatrix}$$

$$= \theta^\top x \quad x, \theta \in \mathbb{R}^{D+1}$$

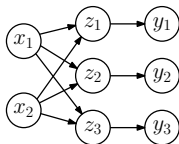
- ▶ Define a **loss function**, $L_\theta(x, y) = \sum_{i=1}^N (y_i - f_\theta(x_i))^2$
- ▶ **Minimisation** of loss function is typically based on **gradient descent**
 - ▶ **Differentiation** required to compute gradients of L_θ w.r.t θ

Scenario 2: Find optimal parameters of a Neural Network

$$\mathbf{x} \in \mathbb{R}^2, \mathbf{A} \in \mathbb{R}^{2 \times 3}, \mathbf{b} \in \mathbb{R}^3$$

$$\mathbf{z} = \mathbf{Ax} + \mathbf{b} \in \mathbb{R}^3$$

$$\mathbf{y} = \sigma(\mathbf{z}) \in \mathbb{R}^3$$

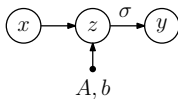
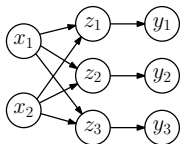


Scenario 2: Find optimal parameters of a Neural Network

$$\mathbf{x} \in \mathbb{R}^2, \mathbf{A} \in \mathbb{R}^{2 \times 3}, \mathbf{b} \in \mathbb{R}^3$$

$$\mathbf{z} = \mathbf{Ax} + \mathbf{b} \in \mathbb{R}^3$$

$$\mathbf{y} = \sigma(\mathbf{z}) \in \mathbb{R}^3$$



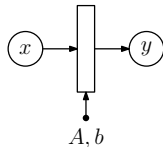
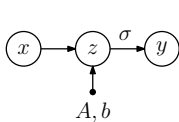
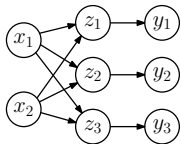
Scenario 2: Find optimal parameters of a Neural Network

$$\mathbf{x} \in \mathbb{R}^2, \mathbf{A} \in \mathbb{R}^{2 \times 3}, \mathbf{b} \in \mathbb{R}^3$$

$$\mathbf{z} = \mathbf{Ax} + \mathbf{b} \in \mathbb{R}^3$$

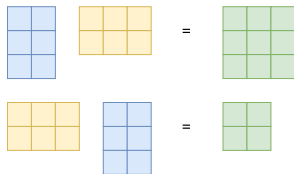
$$\mathbf{y} = \sigma(\mathbf{z}) \in \mathbb{R}^3$$

$$\mathbf{y} = f_{\mathbf{A}, \mathbf{b}}(\mathbf{x})$$



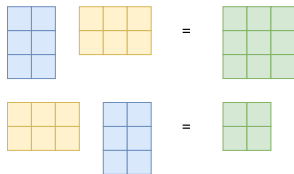
Background: Matrix Multiplication

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



Background: Matrix Multiplication

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

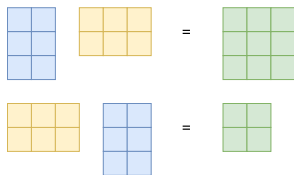


- ▶ When multiplying matrices, the “neighbouring” dimensions have to fit:

$$\underbrace{A}_{n \times k} \underbrace{B}_{k \times m} = \underbrace{C}_{n \times m}$$

Background: Matrix Multiplication

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



- ▶ When multiplying matrices, the “neighbouring” dimensions have to fit:

$$\underbrace{A}_{n \times k} \underbrace{B}_{k \times m} = \underbrace{C}_{n \times m}$$

$$y = Ax$$

$$y_i = \sum_j A_{ij} x_j$$

$$C = AB$$

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

```
import numpy as np
```

```
y = A.dot(x)
```

```
y = np.einsum('ij, j', A, x)
```

```
C = A.dot(B)
```

```
C = np.einsum('ik, kj', A, B)
```

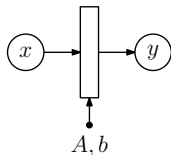
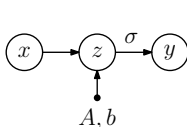
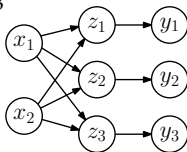
Scenario 2: Find optimal parameters of a Neural Network

$$\mathbf{x} \in \mathbb{R}^2, \mathbf{A} \in \mathbb{R}^{2 \times 3}, \mathbf{b} \in \mathbb{R}^3$$

$$\mathbf{z} = \mathbf{A}\mathbf{x} + \mathbf{b} \in \mathbb{R}^3$$

$$\mathbf{y} = \sigma(\mathbf{z}) \in \mathbb{R}^3$$

$$\mathbf{y} = f_{\mathbf{A}, \mathbf{b}}(\mathbf{x})$$



- ▶ Training data: N pairs $(\mathbf{x}_i, \mathbf{y}_i)$
- ▶ **Training the NN** means finding parameters $\mathbf{A}^*, \mathbf{b}^*$, such that

$$f_{\mathbf{A}^*, \mathbf{b}^*}(\mathbf{x}_i) \approx \mathbf{y}_i$$

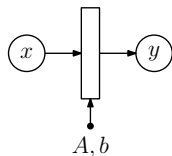
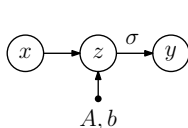
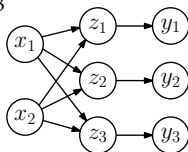
Scenario 2: Find optimal parameters of a Neural Network

$$\mathbf{x} \in \mathbb{R}^2, \mathbf{A} \in \mathbb{R}^{2 \times 3}, \mathbf{b} \in \mathbb{R}^3$$

$$\mathbf{z} = \mathbf{Ax} + \mathbf{b} \in \mathbb{R}^3$$

$$\mathbf{y} = \sigma(\mathbf{z}) \in \mathbb{R}^3$$

$$\mathbf{y} = f_{\mathbf{A}, \mathbf{b}}(\mathbf{x})$$



- ▶ Training data: N pairs $(\mathbf{x}_i, \mathbf{y}_i)$
- ▶ **Training the NN** means finding parameters $\mathbf{A}^*, \mathbf{b}^*$, such that

$$f_{\mathbf{A}^*, \mathbf{b}^*}(\mathbf{x}_i) \approx \mathbf{y}_i$$

- ▶ Define a **loss function**, $L_{\mathbf{A}, \mathbf{b}}(\mathbf{x}, \mathbf{y})$
- ▶ **Minimisation** of loss function is typically based on **gradient descent**
- ▶ **Differentiation** required to compute gradients of $L_{\mathbf{A}, \mathbf{b}}$ w.r.t neural network parameters \mathbf{A}, \mathbf{b}

Types of differentiation

1. Scalar differentiation: $f : \mathbb{R} \rightarrow \mathbb{R}$

$y \in \mathbb{R}$ w.r.t. $x \in \mathbb{R}$

2. Multivariate case: $f : \mathbb{R}^N \rightarrow \mathbb{R}$

$y \in \mathbb{R}$ w.r.t. vector $\mathbf{x} \in \mathbb{R}^N$

3. Vector fields: $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$

vector $\mathbf{y} \in \mathbb{R}^M$ w.r.t. vector $\mathbf{x} \in \mathbb{R}^N$

4. General derivatives: $f : \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^{P \times Q}$

matrix $\mathbf{y} \in \mathbb{R}^{P \times Q}$ w.r.t. matrix $\mathbf{x} \in \mathbb{R}^{M \times N}$

Rules

- ▶ Sum Rule

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

Rules

- ▶ Sum Rule

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

- ▶ Product Rule

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

Rules

- ▶ Sum Rule

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

- ▶ Product Rule

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

- ▶ Chain Rule

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

Rules

- ▶ Sum Rule

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

- ▶ Product Rule

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

- ▶ Chain Rule

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

- ▶ Quotient Rule

$$\left(\frac{f(x)}{g(x)}\right)' = \frac{f'(x)g(x) - f(x)g'(x)}{(g(x))^2} = \frac{\frac{df}{dx}g(x) - f(x)\frac{dg}{dx}}{(g(x))^2}$$

Scalar Differentiation $f : \mathbb{R} \rightarrow \mathbb{R}$

- ▶ Derivative defined as the limit of the difference quotient

$$f'(x) = \frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

- ▶▶ Slope of the secant line through $f(x)$ and $f(x+h)$

Some examples

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

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

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

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

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

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

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

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

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

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

Example: Scalar Chain Rule

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

Beginner

$$g(z) = 6z + 3$$

$$z = f(x) = -2x + 5$$

$$(g \circ f)'(x) =$$

Advanced

$$g(z) = \tanh(z)$$

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

$$(g \circ f)'(x) =$$

Work it out with your neighbours

Example: Scalar Chain Rule

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

Beginner

$$g(z) = 6z + 3$$

$$z = f(x) = -2x + 5$$

$$\begin{aligned}(g \circ f)'(x) &= \underbrace{(6)}_{dg/df} \underbrace{(-2)}_{df/dx} \\ &= -12\end{aligned}$$

Advanced

$$g(z) = \tanh(z)$$

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

$$(g \circ f)'(x) = \underbrace{(1 - \tanh^2(x^n))}_{dg/df} \underbrace{nx^{n-1}}_{df/dx}$$

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

$$y = f(\mathbf{x}), \quad \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} \in \mathbb{R}^N$$

- ▶ **Partial derivative** (change one coordinate at a time):

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

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

$$y = f(\mathbf{x}), \quad \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} \in \mathbb{R}^N$$

- ▶ **Partial derivative** (change one coordinate at a time):

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

- ▶ **Jacobian** vector (**gradient**) collects all partial derivatives:

$$\frac{df}{d\mathbf{x}} = \left[\frac{\partial f}{\partial x_1} \quad \dots \quad \frac{\partial f}{\partial x_N} \right] \in \mathbb{R}^{1 \times N}$$

Note: This is a row vector.

Example: Multivariate differentiation

Beginner

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

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

Advanced

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

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

Partial derivatives?

Work it out with your neighbours

Example: Multivariate differentiation

Beginner

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

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

Advanced

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

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

Partial derivatives

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

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

$$\frac{\partial f(x_1, x_2)}{\partial x_1} = 2(x_1 + 2x_2^3) \underbrace{\frac{d}{dx_1}(x_1 + 2x_2^3)}_{(1)}$$

$$\frac{\partial f(x_1, x_2)}{\partial x_2} = 2(x_1 + 2x_2^3) \underbrace{(6x_2^2)}_{\frac{d}{dx_2}(x_1 + 2x_2^3)}$$

Example: Multivariate differentiation

Beginner

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

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

Advanced

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

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

Partial derivatives

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

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

$$\frac{\partial f(x_1, x_2)}{\partial x_1} = 2(x_1 + 2x_2^3) \underbrace{\frac{d}{dx_1}(x_1 + 2x_2^3)}_{(1)}$$

$$\frac{\partial f(x_1, x_2)}{\partial x_2} = 2(x_1 + 2x_2^3) \underbrace{(6x_2^2)}_{\frac{d}{dx_2}(x_1 + 2x_2^3)}$$

Gradient $\frac{df}{dx} = \left[\frac{\partial f(x_1, x_2)}{\partial x_1} \quad \frac{\partial f(x_1, x_2)}{\partial x_2} \right] \in \mathbb{R}^{1 \times 2}$

$$\frac{df}{dx} = [2x_1 x_2 + x_2^3 \quad x_1^2 + 3x_1 x_2^2]$$

$$\frac{df}{dx} = [2(x_1 + 2x_2^3) \quad 12(x_1 + 2x_2^3)x_2^2]$$

Multivariate Differentiation Rules

- ▶ Sum Rule

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

- ▶ Product Rule

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

- ▶ Chain Rule

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

- ▶ Quotient Rule

$$\left(\frac{f(x)}{g(x)}\right)' = \frac{f'(x)g(x) - f(x)g'(x)}{(g(x))^2} = \frac{\frac{df}{dx}g(x) - f(x)\frac{dg}{dx}}{(g(x))^2}$$

Example: Multivariate Chain Rule

- ▶ Consider the function

$$L(\mathbf{e}) = \frac{1}{2} \|\mathbf{e}\|^2 = \frac{1}{2} \mathbf{e}^\top \mathbf{e}$$

$$\mathbf{e} = \mathbf{y} - \mathbf{A}\mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^N, \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{e}, \mathbf{y} \in \mathbb{R}^M$$

- ▶ Compute $dL/d\mathbf{x}$. What is the dimension/size of $dL/d\mathbf{x}$?

Example: Multivariate Chain Rule

- ▶ Consider the function

$$L(\mathbf{e}) = \frac{1}{2} \|\mathbf{e}\|^2 = \frac{1}{2} \mathbf{e}^\top \mathbf{e}$$

$$\mathbf{e} = \mathbf{y} - \mathbf{A}\mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^N, \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{e}, \mathbf{y} \in \mathbb{R}^M$$

- ▶ Compute dL/dx . What is the dimension/size of dL/dx ?

$$\frac{dL}{dx} = \frac{dL}{de} \frac{de}{dx}$$

$$\frac{dL}{de} = \mathbf{e}^\top \in \mathbb{R}^{1 \times M}$$

$$\frac{de}{dx} = -\mathbf{A} \in \mathbb{R}^{M \times N}$$

$$\Rightarrow \frac{dL}{dx} = \mathbf{e}^\top (-\mathbf{A}) = -(\mathbf{y} - \mathbf{A}\mathbf{x})^\top \mathbf{A} \in \mathbb{R}^{1 \times N}$$

Vector field differentiation $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$

$$\mathbf{y} = f(\mathbf{x}) \in \mathbb{R}^M, \quad \mathbf{x} \in \mathbb{R}^N$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} f_1(x_1, \dots, x_N) \\ \vdots \\ f_M(x_1, \dots, x_N) \end{bmatrix}$$

Vector field differentiation $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$

$$\mathbf{y} = f(\mathbf{x}) \in \mathbb{R}^M, \quad \mathbf{x} \in \mathbb{R}^N$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} f_1(x_1, \dots, x_N) \\ \vdots \\ f_M(x_1, \dots, x_N) \end{bmatrix}$$

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

$$\begin{bmatrix} \frac{dy_1}{dx} \\ \vdots \\ \frac{dy_M}{dx} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_i} & \cdots & \frac{\partial f_1}{\partial x_N} \\ \vdots & & \vdots & & \vdots \\ \frac{\partial f_M}{\partial x_1} & \cdots & \frac{\partial f_M}{\partial x_i} & \cdots & \frac{\partial f_M}{\partial x_N} \end{bmatrix} \in \mathbb{R}^{M \times N}$$

Example: Vector field differentiation

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(x) \\ \vdots \\ f_M(x) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1(N-1)}x_{N-1} + A_{1N}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{M(N-1)}x_{N-1} + A_{MN}x_N \end{bmatrix}$$

- ▶ Compute the gradient df/dx
 - ▶ Dimension of df/dx :

Example: Vector field differentiation

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(x) \\ \vdots \\ f_M(x) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1(N-1)}x_{N-1} + A_{1N}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{M(N-1)}x_{N-1} + A_{MN}x_N \end{bmatrix}$$

► Compute the gradient df/dx

► Dimension of df/dx : Since $f: \mathbb{R}^N \rightarrow \mathbb{R}^M$, $df/dx \in \mathbb{R}^{M \times N}$

Example: Vector field differentiation

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(x) \\ \vdots \\ f_M(x) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1(N-1)}x_{N-1} + A_{1N}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{M(N-1)}x_{N-1} + A_{MN}x_N \end{bmatrix}$$

► Compute the gradient df/dx

► Dimension of df/dx : Since $f: \mathbb{R}^N \rightarrow \mathbb{R}^M$, $df/dx \in \mathbb{R}^{M \times N}$

► Gradient:

$$f_i(x) = \sum_{j=1}^N A_{ij}x_j \quad \Rightarrow \quad \frac{\partial f_i}{\partial x_j} = A_{ij}$$

Example: Vector field differentiation

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(x) \\ \vdots \\ f_M(x) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1(N-1)}x_{N-1} + A_{1N}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{M(N-1)}x_{N-1} + A_{MN}x_N \end{bmatrix}$$

► Compute the gradient df/dx

► Dimension of df/dx : Since $f: \mathbb{R}^N \rightarrow \mathbb{R}^M$, $df/dx \in \mathbb{R}^{M \times N}$

► Gradient:

$$f_i(x) = \sum_{j=1}^N A_{ij}x_j \quad \Rightarrow \quad \frac{\partial f_i}{\partial x_j} = A_{ij}$$
$$\Rightarrow \frac{df}{dx} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_N} \\ \vdots & & \vdots \\ \frac{\partial f_M}{\partial x_1} & \cdots & \frac{\partial f_M}{\partial x_N} \end{bmatrix} = \begin{bmatrix} A_{11} & \cdots & A_{1N} \\ \vdots & & \vdots \\ A_{M1} & \cdots & A_{MN} \end{bmatrix} = A$$

Chain Rule

$$\frac{\partial}{\partial \mathbf{x}}(g \circ f)(\mathbf{x}) = \frac{\partial}{\partial \mathbf{x}}(g(f(\mathbf{x}))) = \frac{\partial g}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial \mathbf{x}}$$

Example: Combined Chain Rule

- ▶ Consider $f : \mathbb{R}^2 \rightarrow \mathbb{R}$, $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^2$

$$f(\mathbf{x}) = f(x_1, x_2) = x_1^2 + 2x_2,$$

$$\mathbf{x}(t) = \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix} = \begin{bmatrix} \sin(t) \\ \cos(t) \end{bmatrix}$$

Example: Combined Chain Rule

- ▶ Consider $f : \mathbb{R}^2 \rightarrow \mathbb{R}$, $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^2$

$$f(\mathbf{x}) = f(x_1, x_2) = x_1^2 + 2x_2,$$

$$\mathbf{x}(t) = \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix} = \begin{bmatrix} \sin(t) \\ \cos(t) \end{bmatrix}$$

- ▶ What are the dimensions of $df/d\mathbf{x}$ and $d\mathbf{x}/dt$?

Example: Combined Chain Rule

- ▶ Consider $f : \mathbb{R}^2 \rightarrow \mathbb{R}$, $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^2$

$$f(\mathbf{x}) = f(x_1, x_2) = x_1^2 + 2x_2,$$

$$\mathbf{x}(t) = \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix} = \begin{bmatrix} \sin(t) \\ \cos(t) \end{bmatrix}$$

- ▶ What are the dimensions of $df/d\mathbf{x}$ and $d\mathbf{x}/dt$?

$$1 \times 2 \text{ and } 2 \times 1$$

- ▶ Compute the gradient df/dt using the chain rule:

Work it out with your neighbours

Example: Combined Chain Rule

- ▶ Consider $f : \mathbb{R}^2 \rightarrow \mathbb{R}$, $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^2$

$$f(\mathbf{x}) = f(x_1, x_2) = x_1^2 + 2x_2,$$

$$\mathbf{x}(t) = \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix} = \begin{bmatrix} \sin(t) \\ \cos(t) \end{bmatrix}$$

- ▶ What are the dimensions of $df/d\mathbf{x}$ and $d\mathbf{x}/dt$?

$$1 \times 2 \text{ and } 2 \times 1$$

- ▶ Compute the gradient df/dt using the chain rule:

$$\begin{aligned} \frac{df}{dt} &= \frac{df}{d\mathbf{x}} \frac{d\mathbf{x}}{dt} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix} \begin{bmatrix} \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial t} \end{bmatrix} = \begin{bmatrix} 2 \sin t & 2 \end{bmatrix} \begin{bmatrix} \cos t \\ -\sin t \end{bmatrix} \\ &= 2 \sin t \cos t - 2 \sin t = 2 \sin t (\cos t - 1) \end{aligned}$$

5 MINUTE BREAK

Derivatives with respect to Matrices

- ▶ Recall: function $f : \mathbb{R}^D \rightarrow \mathbb{R}^E$ has a gradient that is $E \times D$ -matrix

target dimensions \times # input dimensions

with

$$\frac{df}{dx} \in \mathbb{R}^{E \times D}, \quad df[e, d] = \frac{\partial f_e}{\partial x_d}$$

Derivatives with respect to Matrices

- ▶ Recall: function $f : \mathbb{R}^D \rightarrow \mathbb{R}^E$ has a gradient that is $E \times D$ -matrix

target dimensions \times # input dimensions

with

$$\frac{df}{dx} \in \mathbb{R}^{E \times D}, \quad df[e, d] = \frac{\partial f_e}{\partial x_d}$$

- ▶ This generalises to when the inputs (D) or targets (E) are **matrices**

Derivatives with respect to Matrices

- ▶ Recall: function $f : \mathbb{R}^D \rightarrow \mathbb{R}^E$ has a gradient that is $E \times D$ -matrix

target dimensions \times # input dimensions

with

$$\frac{df}{dx} \in \mathbb{R}^{E \times D}, \quad df[e, d] = \frac{\partial f_e}{\partial x_d}$$

- ▶ This generalises to when the inputs (D) or targets (E) are **matrices**
- ▶ Function $f : \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^{P \times Q}$, has a gradient that is a $(P \times Q) \times (M \times N)$ object (tensor)

$$\frac{df}{d\mathbf{X}} \in \mathbb{R}^{(P \times Q) \times (M \times N)}, \quad df[p, q, m, n] = \frac{\partial f_{pq}}{\partial X_{mn}}$$

Example 1: Derivatives w.r.t Matrices

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1(N-1)}x_{N-1} + A_{1N}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{M(N-1)}x_{N-1} + A_{MN}x_N \end{bmatrix}$$

$$\frac{d\mathbf{f}}{d\mathbf{A}} \in \mathbb{R}^?$$

Example 1: Derivatives w.r.t Matrices

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(x) \\ \vdots \\ f_M(x) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1(N-1)}x_{N-1} + A_{1N}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{M(N-1)}x_{N-1} + A_{MN}x_N \end{bmatrix}$$

$$\frac{df}{dA} \in \mathbb{R}^{\# \text{ target dim} \times \# \text{ input dim}} = M \times (M \times N)$$

$$\frac{df}{dA} = \begin{bmatrix} \frac{\partial f_1}{\partial A} \\ \vdots \\ \frac{\partial f_M}{\partial A} \end{bmatrix}, \quad \frac{\partial f_i}{\partial A} \in \mathbb{R}^{1 \times (M \times N)}$$

Example 2: Derivatives w.r.t Matrices

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_i(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1(N-1)}x_{N-1} + A_{1N}x_N \\ \vdots \\ A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{i(N-1)}x_{N-1} + A_{iN}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{M(N-1)}x_{N-1} + A_{MN}x_N \end{bmatrix}$$

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

$$\frac{\partial f_i}{\partial A_{ij}} = ?$$

$$\frac{\partial f_i}{\partial A_{k \neq i, :}} = ?$$

$$\frac{\partial f_i}{\partial \mathbf{A}} = ?$$

Example 2: Derivatives w.r.t Matrices

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_i(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1(N-1)}x_{N-1} + A_{1N}x_N \\ \vdots \\ A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{i(N-1)}x_{N-1} + A_{iN}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{M(N-1)}x_{N-1} + A_{MN}x_N \end{bmatrix}$$

$$\frac{\partial f_i}{\partial A_{iq}} = \underbrace{x_q}_{\in \mathbb{R}} \quad \frac{\partial f_i}{\partial A_{i,:}} = ? \quad \frac{\partial f_i}{\partial A_{k \neq i,:}} = ? \quad \frac{\partial f_i}{\partial \mathbf{A}} = ?$$

Example 2: Derivatives w.r.t Matrices

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_i(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1(N-1)}x_{N-1} + A_{1N}x_N \\ \vdots \\ A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{i(N-1)}x_{N-1} + A_{iN}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{M(N-1)}x_{N-1} + A_{MN}x_N \end{bmatrix}$$

$$\frac{\partial f_i}{\partial A_{iq}} = \underbrace{x_q}_{\in \mathbb{R}} \quad \frac{\partial f_i}{\partial A_{i,:}} = \underbrace{\mathbf{x}^\top}_{\in \mathbb{R}^{1 \times 1 \times N}} \quad \frac{\partial f_i}{\partial A_{k \neq i,:}} = ? \quad \frac{\partial f_i}{\partial \mathbf{A}} = ?$$

Example 2: Derivatives w.r.t Matrices

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_i(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1(N-1)}x_{N-1} + A_{1N}x_N \\ \vdots \\ A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{i(N-1)}x_{N-1} + A_{iN}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{M(N-1)}x_{N-1} + A_{MN}x_N \end{bmatrix}$$

$$\frac{\partial f_i}{\partial A_{iq}} = \underbrace{x_q}_{\in \mathbb{R}} \quad \frac{\partial f_i}{\partial A_{i,:}} = \underbrace{\mathbf{x}^\top}_{\in \mathbb{R}^{1 \times 1 \times N}} \quad \frac{\partial f_i}{\partial A_{k \neq i,:}} = \underbrace{\mathbf{0}^\top}_{\in \mathbb{R}^{1 \times 1 \times N}} \quad \frac{\partial f_i}{\partial \mathbf{A}} = ?$$

Example 2: Derivatives w.r.t Matrices

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

$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_i(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1(N-1)}x_{N-1} + A_{1N}x_N \\ \vdots \\ A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{i(N-1)}x_{N-1} + A_{iN}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{M(N-1)}x_{N-1} + A_{MN}x_N \end{bmatrix}$$

$$\frac{\partial f_i}{\partial A_{iq}} = \underbrace{x_q}_{\in \mathbb{R}} \quad \frac{\partial f_i}{\partial A_{i,:}} = \underbrace{\mathbf{x}^\top}_{\in \mathbb{R}^{1 \times 1 \times N}} \quad \frac{\partial f_i}{\partial A_{k \neq i,:}} = \underbrace{\mathbf{0}^\top}_{\in \mathbb{R}^{1 \times 1 \times N}} \quad \frac{\partial f_i}{\partial \mathbf{A}} = \underbrace{\begin{bmatrix} \mathbf{0}^\top \\ \vdots \\ \mathbf{x}^\top \\ \vdots \\ \mathbf{0}^\top \end{bmatrix}}_{\in \mathbb{R}^{1 \times (M \times N)}}$$

Example 3: Derivatives w.r.t Higher-Order Tensors

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

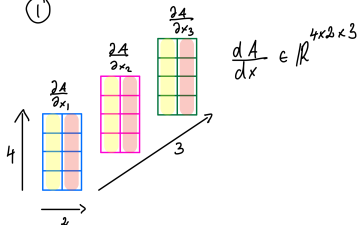
matrix
 $A \in \mathbb{R}^{4 \times 2}$



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



①



Example 3: Derivatives w.r.t Higher-Order Tensors

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

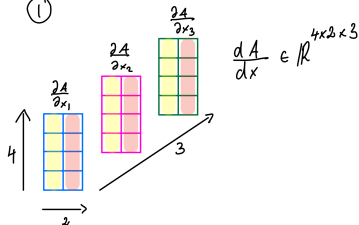
matrix
 $A \in \mathbb{R}^{4 \times 2}$



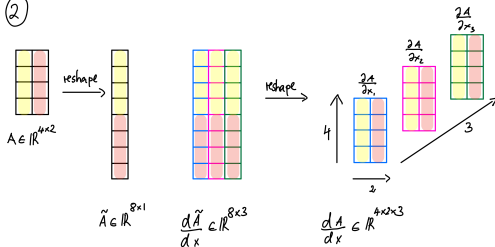
vector
 $x \in \mathbb{R}^3$



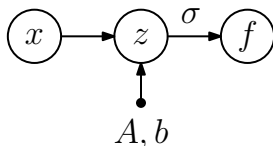
①



②



Gradients of a Single-Layer Neural Network



$$f = \tanh(\underbrace{Ax + b}_{=:z \in \mathbb{R}^M}) \in \mathbb{R}^M, \quad x \in \mathbb{R}^N, A \in \mathbb{R}^{M \times N}, b \in \mathbb{R}^M$$

Gradients of a Single-Layer Neural Network

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

$$\frac{\partial f}{\partial \mathbf{b}} =$$

$$\frac{\partial f}{\partial \mathbf{A}} =$$

Gradients of a Single-Layer Neural Network

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

$$\frac{\partial \mathbf{f}}{\partial \mathbf{b}} = \underbrace{\frac{\partial \mathbf{f}}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{b}}}_{M \times M} \in \mathbb{R}^{M \times M}$$

$$\frac{\partial \mathbf{f}}{\partial \mathbf{A}} =$$

Gradients of a Single-Layer Neural Network

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

$$\frac{\partial \mathbf{f}}{\partial \mathbf{b}} = \underbrace{\frac{\partial \mathbf{f}}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{b}}}_{M \times M} \in \mathbb{R}^{M \times M}$$

$$\frac{\partial \mathbf{f}}{\partial \mathbf{A}} = \underbrace{\frac{\partial \mathbf{f}}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{A}}}_{M \times (M \times N)} \in \mathbb{R}^{M \times (M \times N)}$$

Gradients of a Single-Layer Neural Network

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

$$\frac{\partial f}{\partial \mathbf{b}} = \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{b}}}_{M \times M} \in \mathbb{R}^{M \times M}$$

$$\frac{\partial f}{\partial \mathbf{A}} = \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{A}}}_{M \times (M \times N)} \in \mathbb{R}^{M \times (M \times N)}$$

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

Gradients of a Single-Layer Neural Network

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

$$\frac{\partial \mathbf{f}}{\partial \mathbf{b}} = \underbrace{\frac{\partial \mathbf{f}}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{b}}}_{M \times M} \in \mathbb{R}^{M \times M}$$

$$\frac{\partial \mathbf{f}}{\partial \mathbf{A}} = \underbrace{\frac{\partial \mathbf{f}}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{A}}}_{M \times (M \times N)} \in \mathbb{R}^{M \times (M \times N)}$$

$$\frac{\partial \mathbf{f}}{\partial \mathbf{z}} = \underbrace{\text{diag}(1 - \tanh^2(\mathbf{z}))}_{\in \mathbb{R}^{M \times M}} \quad \frac{\partial \mathbf{z}}{\partial \mathbf{b}} = \underbrace{\mathbf{I}}_{\in \mathbb{R}^{M \times M}}$$

Gradients of a Single-Layer Neural Network

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

$$\frac{\partial f}{\partial \mathbf{b}} = \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{b}}}_{M \times M} \in \mathbb{R}^{M \times M}$$

$$\frac{\partial f}{\partial \mathbf{A}} = \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{M \times M} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{A}}}_{M \times (M \times N)} \in \mathbb{R}^{M \times (M \times N)}$$

$$\frac{\partial f}{\partial \mathbf{z}} = \underbrace{\text{diag}(1 - \tanh^2(\mathbf{z}))}_{\in \mathbb{R}^{M \times M}} \quad \frac{\partial \mathbf{z}}{\partial \mathbf{b}} = \underbrace{\mathbf{I}}_{\in \mathbb{R}^{M \times M}} \quad \frac{\partial \mathbf{z}}{\partial \mathbf{A}} = \underbrace{\begin{bmatrix} \mathbf{x}^\top & \cdot & \mathbf{0}^\top & \cdot & \mathbf{0}^\top \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \mathbf{0}^\top & \cdot & \mathbf{x}^\top & \cdot & \mathbf{0}^\top \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \mathbf{0}^\top & \cdot & \mathbf{0}^\top & \cdot & \mathbf{x}^\top \end{bmatrix}}_{\in \mathbb{R}^{M \times (M \times N)}}$$

```

import numpy as np
def NN(x, A, b):
    M, N = A.shape
    z = A.dot(x) + b
    f = np.tanh(z)

    # partial derivatives
    dfdz = 1-f**2
    dzdx = A
    dzdb = np.eye(M)
    dzdA = np.zeros((M, M, N))
    for i in range(M):
        dzdA[i,i,:] = x.T

    # gradients
    dfdx = np.einsum('il, lj', dfdz, dzdx)
    dfdb = np.einsum('il, lj', dfdz, dzdb)
    dfdA = np.einsum('il, ljk', dfdz, dzdA)

    return f, dfdA, dfdb, dfdx

```

$$\frac{\partial f}{\partial \mathbf{x}}[i, j] = \sum_{l=1}^M \frac{\partial f}{\partial z}[i, l] \frac{\partial z}{\partial \mathbf{x}}[l, j]$$

$$\frac{\partial f}{\partial \mathbf{b}}[i, j] = \sum_{l=1}^M \frac{\partial f}{\partial z}[i, l] \frac{\partial z}{\partial \mathbf{b}}[l, j]$$

$$\frac{\partial f}{\partial \mathbf{A}}[i, j, k] = \sum_{l=1}^M \frac{\partial f}{\partial z}[i, l] \frac{\partial z}{\partial \mathbf{A}}[l, j, k]$$

Putting Things Together

- ▶ Inputs $x \in \mathbb{R}^N$

Putting Things Together

- ▶ Inputs $\mathbf{x} \in \mathbb{R}^N$
- ▶ Observed outputs $\mathbf{y} = f_{\theta}(\mathbf{z}) + \boldsymbol{\epsilon} \in \mathbb{R}^M, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$

Putting Things Together

- ▶ Inputs $\mathbf{x} \in \mathbb{R}^N$
- ▶ Observed outputs $\mathbf{y} = f_{\boldsymbol{\theta}}(\mathbf{z}) + \boldsymbol{\epsilon} \in \mathbb{R}^M, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$
- ▶ Train single-layer neural network with

$$f_{\boldsymbol{\theta}}(\mathbf{z}) = \tanh(\mathbf{z}) \in \mathbb{R}^M, \quad \mathbf{z} = \mathbf{A}\mathbf{x} + \mathbf{b} \in \mathbb{R}^M, \quad \boldsymbol{\theta} = \{\mathbf{A}, \mathbf{b}\}$$

Putting Things Together

- ▶ Inputs $\mathbf{x} \in \mathbb{R}^N$
- ▶ Observed outputs $\mathbf{y} = f_{\boldsymbol{\theta}}(\mathbf{z}) + \boldsymbol{\epsilon} \in \mathbb{R}^M, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$
- ▶ Train single-layer neural network with

$$f_{\boldsymbol{\theta}}(\mathbf{z}) = \tanh(\mathbf{z}) \in \mathbb{R}^M, \quad \mathbf{z} = \mathbf{A}\mathbf{x} + \mathbf{b} \in \mathbb{R}^M, \quad \boldsymbol{\theta} = \{\mathbf{A}, \mathbf{b}\}$$

- ▶ Find $\mathbf{A}^*, \mathbf{b}^*$, such that the squared loss

$$L(\boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{e}\|^2 \in \mathbb{R}, \quad \mathbf{e} = \mathbf{y} - f_{\boldsymbol{\theta}}(\mathbf{z}) \in \mathbb{R}^M$$

is minimised

Putting Things Together

Partial derivatives:

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{A}} &= \frac{\partial L}{\partial \mathbf{e}} \frac{\partial \mathbf{e}}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{A}} \\ \frac{\partial L}{\partial \mathbf{b}} &= \frac{\partial L}{\partial \mathbf{e}} \frac{\partial \mathbf{e}}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{b}}\end{aligned}$$

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{e}} &= \underbrace{\mathbf{e}^\top}_{\in \mathbb{R}^{1 \times M}} & \frac{\partial \mathbf{e}}{\partial \mathbf{f}} &= \underbrace{-\mathbf{I}}_{\in \mathbb{R}^{M \times M}} & \frac{\partial \mathbf{f}}{\partial \mathbf{z}} &= \underbrace{\text{diag}(1 - \tanh^2(\mathbf{z}))}_{\in \mathbb{R}^{M \times M}} \\ \frac{\partial \mathbf{z}}{\partial \mathbf{A}} &= \underbrace{\begin{bmatrix} \mathbf{x}^\top & \cdot & \mathbf{0}^\top & \cdot & \mathbf{0}^\top \\ \cdot & & \cdot & & \cdot \\ \mathbf{0}^\top & \cdot & \mathbf{x}^\top & \cdot & \mathbf{0}^\top \\ \cdot & & \cdot & & \cdot \\ \mathbf{0}^\top & \cdot & \mathbf{0}^\top & \cdot & \mathbf{x}^\top \end{bmatrix}}_{\in \mathbb{R}^{M \times (M \times N)}} & \frac{\partial \mathbf{z}}{\partial \mathbf{b}} &= \underbrace{\mathbf{I}}_{\in \mathbb{R}^{M \times M}}\end{aligned}$$

Gradient descent for neural network optimisation

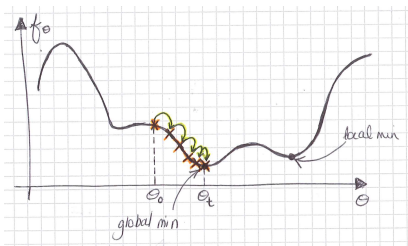
Goal: Find the local optimum of loss $L(\theta^*) \in \mathbb{R}$ where $\theta^* = \{A^*, b^*\}$

Gradient descent algorithm:

1. Initialise neural network parameters, $\theta_0 = \{A_0, b_0\}$
2. Compute gradient of L w.r.t. parameters, $\frac{\partial L}{\partial \theta_0}$
3. Update initial guess θ_0 using chain rule:

$$\theta_{i+1} = \theta_i + \Delta\theta_i = \theta_i - \gamma \left(\frac{\partial L}{\partial \theta_i}(\theta_i) \right)^\top$$

► Note: parameter update is the direction of steepest **descent**



Gradient descent for neural network optimisation

Goal: Find the local optimum of loss $L(\boldsymbol{\theta}^*) \in \mathbb{R}$ where $\boldsymbol{\theta}^* = \{\mathbf{A}^*, \mathbf{b}^*\}$

Gradient descent algorithm:

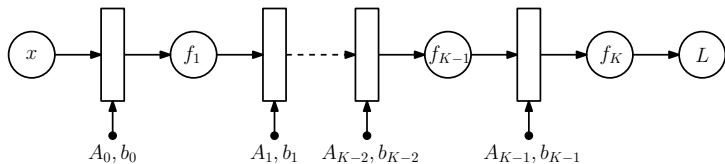
1. Initialise neural network parameters, $\boldsymbol{\theta}_0 = \{\mathbf{A}_0, \mathbf{b}_0\}$
2. Compute gradient of L w.r.t. parameters, $\frac{\partial L}{\partial \boldsymbol{\theta}_0}$
3. Update initial guess $\boldsymbol{\theta}_0$ using chain rule:

$$\boldsymbol{\theta}_{i+1} = \boldsymbol{\theta}_i + \Delta \boldsymbol{\theta}_i = \boldsymbol{\theta}_i - \gamma \left(\frac{\partial L}{\partial \boldsymbol{\theta}_i}(\boldsymbol{\theta}_i) \right)^\top$$

▶▶ Note: parameter update is the direction of steepest **descent**

For suitable step-size γ , the sequence $L(\boldsymbol{\theta}_0) \geq L(\boldsymbol{\theta}_1) \geq \dots$ converges to a local minimum $L(\boldsymbol{\theta}^*)$.

Gradients of a Multi-Layer Neural Network

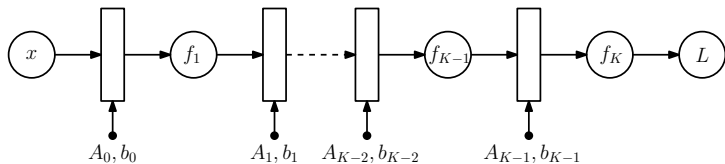


- ▶ Inputs x , observed outputs y
- ▶ Train K-layer neural network with

$$f_0 = x$$

$$f_i = \sigma_i(A_{i-1}f_{i-1} + b_{i-1}), \quad i = 1, \dots, K$$

Gradients of a Multi-Layer Neural Network



- ▶ Inputs x , observed outputs y
- ▶ Train K -layer neural network with

$$f_0 = x$$

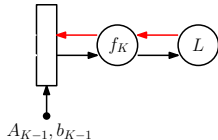
$$f_i = \sigma_i(A_{i-1}f_{i-1} + \mathbf{b}_{i-1}), \quad i = 1, \dots, K$$

- ▶ Find A_k^*, \mathbf{b}_k^* for $k = 0, \dots, K-1$, such that the squared loss

$$L(\theta) = \|\mathbf{y} - \mathbf{f}_{K,\theta}(x)\|^2$$

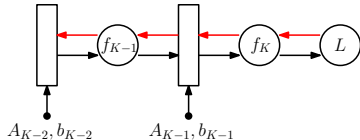
is minimised, where $\theta = \{A_k, \mathbf{b}_k\}$, $k = 0, \dots, K-1$

Gradients of a Multi-Layer Neural Network



$$\frac{\partial L}{\partial \boldsymbol{\theta}_{K-1}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial \boldsymbol{\theta}_{K-1}}$$

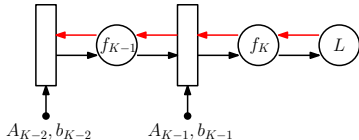
Gradients of a Multi-Layer Neural Network



$$\frac{\partial L}{\partial \theta_{K-1}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial \theta_{K-1}}$$

$$\frac{\partial L}{\partial \theta_{K-2}} = \frac{\partial L}{\partial f_K} \boxed{\frac{\partial f_K}{\partial f_{K-1}} \frac{\partial f_{K-1}}{\partial \theta_{K-2}}}$$

Gradients of a Multi-Layer Neural Network

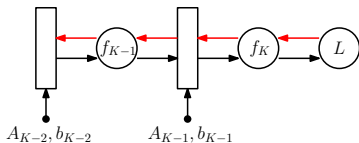


$$\frac{\partial L}{\partial \theta_{K-1}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial \theta_{K-1}}$$

$$\frac{\partial L}{\partial \theta_{K-2}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial f_{K-1}} \frac{\partial f_{K-1}}{\partial \theta_{K-2}}$$

$$\frac{\partial L}{\partial \theta_{K-3}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial f_{K-1}} \frac{\partial f_{K-1}}{\partial f_{K-2}} \frac{\partial f_{K-2}}{\partial \theta_{K-3}}$$

Gradients of a Multi-Layer Neural Network



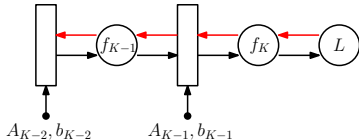
$$\frac{\partial L}{\partial \theta_{K-1}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial \theta_{K-1}}$$

$$\frac{\partial L}{\partial \theta_{K-2}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial f_{K-1}} \frac{\partial f_{K-1}}{\partial \theta_{K-2}}$$

$$\frac{\partial L}{\partial \theta_{K-3}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial f_{K-1}} \frac{\partial f_{K-1}}{\partial f_{K-2}} \frac{\partial f_{K-2}}{\partial \theta_{K-3}}$$

$$\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial f_{K-1}} \dots \frac{\partial f_{i+2}}{\partial f_{i+1}} \frac{\partial f_{i+1}}{\partial \theta_i}$$

Gradients of a Multi-Layer Neural Network



$$\frac{\partial L}{\partial \theta_{K-1}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial \theta_{K-1}}$$

$$\frac{\partial L}{\partial \theta_{K-2}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial f_{K-1}} \frac{\partial f_{K-1}}{\partial \theta_{K-2}}$$

$$\frac{\partial L}{\partial \theta_{K-3}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial f_{K-1}} \frac{\partial f_{K-1}}{\partial f_{K-2}} \frac{\partial f_{K-2}}{\partial \theta_{K-3}}$$

$$\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial f_{K-1}} \dots \frac{\partial f_{i+2}}{\partial f_{i+1}} \frac{\partial f_{i+1}}{\partial \theta_i}$$

► Intermediate derivatives are stored during the forward pass

Training NNs as Maximum Likelihood Estimation

- ▶ $\mathbf{y} = f_{\theta}(\mathbf{x}) + \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- ▶ **Goal:** Maximise log-likelihood of training data (\mathbf{X}, \mathbf{Y}) :
 $\log p_{\theta}(\mathbf{Y}|\mathbf{X})$

Training NNs as Maximum Likelihood Estimation

▶ $\mathbf{y} = f_{\theta}(\mathbf{x}) + \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

▶ **Goal:** Maximise log-likelihood of training data (\mathbf{X}, \mathbf{Y}) :

$$\log p_{\theta}(\mathbf{Y}|\mathbf{X}) = \sum_{n=1}^N \log p_{\theta}(\mathbf{y}_n|\mathbf{x}_n) = \sum_{n=1}^N \log \mathcal{N}(\mathbf{y}_n = f_{\theta}(\mathbf{x}_n), \mathbf{I})$$

Training NNs as Maximum Likelihood Estimation

▶ $\mathbf{y} = f_{\theta}(\mathbf{x}) + \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

▶ **Goal:** Maximise log-likelihood of training data (\mathbf{X}, \mathbf{Y}) :

$$\log p_{\theta}(\mathbf{Y}|\mathbf{X}) = \sum_{n=1}^N \log p_{\theta}(\mathbf{y}_n|\mathbf{x}_n) = \sum_{n=1}^N \log \mathcal{N}(\mathbf{y}_n = f_{\theta}(\mathbf{x}_n), \mathbf{I})$$

▶ Then:

$$\log p_{\theta}(\mathbf{y}_n|\mathbf{x}_n) = \frac{N}{2} \log(2\pi) - \frac{1}{2} \|\mathbf{y} - f_{\theta}(\mathbf{x})\|^2$$

Training NNs as Maximum Likelihood Estimation

▶ $\mathbf{y} = f_{\theta}(\mathbf{x}) + \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

▶ **Goal:** Maximise log-likelihood of training data (\mathbf{X}, \mathbf{Y}) :

$$\log p_{\theta}(\mathbf{Y}|\mathbf{X}) = \sum_{n=1}^N \log p_{\theta}(\mathbf{y}_n|\mathbf{x}_n) = \sum_{n=1}^N \log \mathcal{N}(\mathbf{y}_n = f_{\theta}(\mathbf{x}_n), \mathbf{I})$$

▶ Then:

$$\log p_{\theta}(\mathbf{y}_n|\mathbf{x}_n) = \frac{N}{2} \log(2\pi) - \frac{1}{2} \|\mathbf{y} - f_{\theta}(\mathbf{x})\|^2$$

▶ Find θ^* by minimising the negative log-likelihood:

$$\begin{aligned} \theta^* &= \arg \min_{\theta} [-\log p_{\theta}(\mathbf{Y}|\mathbf{X})] = \arg \min_{\theta} \frac{1}{2} \sum_{n=1}^N \|\mathbf{y} - f_{\theta}(\mathbf{x})\|^2 \\ &= \arg \min_{\theta} L(\theta) \end{aligned}$$

▶▶ **Equivalent to minimising loss of neural network**

Training NNs as Maximum Likelihood Estimation

▶ $\mathbf{y} = f_{\theta}(\mathbf{x}) + \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

▶ **Goal:** Maximise log-likelihood of training data (\mathbf{X}, \mathbf{Y}) :

$$\log p_{\theta}(\mathbf{Y}|\mathbf{X}) = \sum_{n=1}^N \log p_{\theta}(\mathbf{y}_n|\mathbf{x}_n) = \sum_{n=1}^N \log \mathcal{N}(\mathbf{y}_n = f_{\theta}(\mathbf{x}_n), \mathbf{I})$$

▶ Then:

$$\log p_{\theta}(\mathbf{y}_n|\mathbf{x}_n) = \frac{N}{2} \log(2\pi) - \frac{1}{2} \|\mathbf{y} - f_{\theta}(\mathbf{x})\|^2$$

▶ Find θ^* by **minimising the negative log-likelihood**:

$$\begin{aligned} \theta^* &= \arg \min_{\theta} [-\log p_{\theta}(\mathbf{Y}|\mathbf{X})] = \arg \min_{\theta} \frac{1}{2} \sum_{n=1}^N \|\mathbf{y} - f_{\theta}(\mathbf{x})\|^2 \\ &= \arg \min_{\theta} L(\theta) \end{aligned}$$

▶▶ **Equivalent to minimising loss of neural network**

▶ ML estimation can lead to **overfitting** (interprets noise as signal)

Example: Linear Regression with Polynomials

- ▶ Linear regression with a polynomial of order M :

$$y = f_{\theta}(x) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_{\epsilon}^2)$$

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

Example: Linear Regression with Polynomials

- ▶ Linear regression with a polynomial of order M :

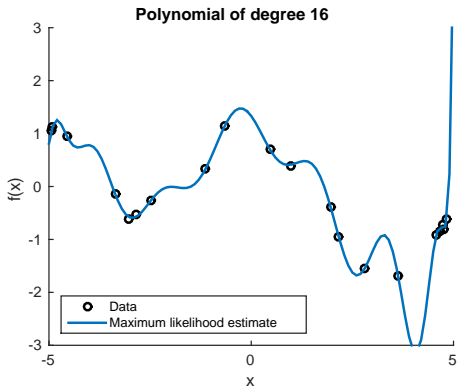
$$y = f_{\boldsymbol{\theta}}(x) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_{\epsilon}^2)$$

$$f_{\boldsymbol{\theta}}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_M x^M = \sum_{i=0}^M \theta_i x^i$$

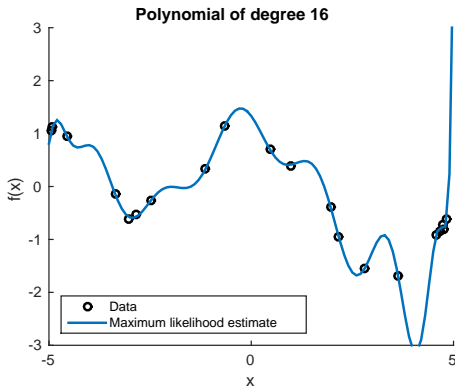
- ▶ Given inputs x_i and corresponding (noisy) observations y_i , $i = 1, \dots, N$, find parameters $\boldsymbol{\theta} = [\theta_0, \dots, \theta_M]^{\top}$ that minimise the squared loss (equivalently: maximise the likelihood)

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

Example: Linear Regression with Polynomials



Example: Linear Regression with Polynomials



- ▶ Overfitting!
- ▶ Regularization, model selection etc. necessary
- ▶▶ Tutorials later this week

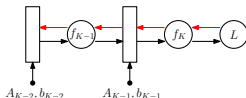
Summary



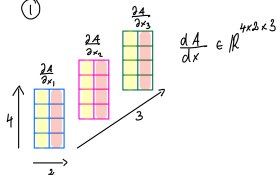
matrix
 $A \in \mathbb{R}^{4 \times 2}$



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



①



- ▶ Key concepts in probability - PDFs, Sum & Product Rules, Bayes
- ▶ **Differentiation** for optimising parameters of ML models
 - ▶▶ Vector calculus
 - ▶▶ Chain rule
 - ▶▶ Single- & Multi-layer NNs
 - ▶▶ Maximum Likelihood Estimation Equivalence

Useful Resources

- ▶ [Indaba 2017 slides/videos](#)
- ▶ Textbooks:
 - ▶ [Mathematics for ML](#) - Marc Deisenroth, Aldo Faisal, Cheng Soon Ong
 - ▶ [Deep Learning](#) - Yoshua Bengio, Ian Goodfellow, Aaron Courville
 - ▶ [Matrix Cookbook](#) - Kaare Petersen, Michael Pedersen
- ▶ Online courses:
 - ▶ [fast.ai Deep Learning by Jeremy Howard](#)
 - ▶ [deeplearning.ai/coursera Deep Learning Specialisation by Andrew Ng](#)
 - ▶ [Udacity Machine Learning Nanodegree Programs](#)
 - ▶ [Coursera Machine Learning by Stanford University](#)
- ▶ Blogs, podcasts:
 - ▶ [Reddit ML](#), [KDnuggets](#), [TWiML&AI](#), [Talking Machines](#)

More textbooks, online courses, papers, datasets and other resources at github.com/ChristosChristofidis/awesome-deep-learning

Rest of the week...

Monday

Deep Learning Fundamentals - [Moustapha Cisse](#)

Convolutional Models - [Naila Murray](#)

Tuesday

Probabilistic Thinking - [Yabebal Fantaye](#)

Recurrent Neural Networks - [Kyunghyun Cho](#)

Wednesday

Generative Models - [Asja Fischer](#)

Reinforcement Learning - [Katja Hofmann](#)

Thursday - Parallel Sessions

Computer Vision, Non-Recurrent Sequence Models, AI for Africa, Generative Modelling and Healthcare, Natural Language Processing, Life of a ML Start-up

Friday - More Parallel Sessions

Reinforcement Learning, ML in Production, AI Ethics & Policy

Learn. Connect. Strengthen.
Enjoy.



DEEP
LEARNING
INDABA

Stellenbosch University
9th September 2018