

Mathematical Foundation in Deep Learning

References:

deeplearningbook.org

Dive into Deep Learning, <https://d2l.ai/index.html>

introtodeeplearning.com

deeplearningindaba.com

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

- Maximum Likelihood Estimation , Cross Entropy, KL Divergence

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

Maximum Likelihood Estimation , Cross Entropy, KL Divergence

Differentiation

Where it is used in ML?

Types of differentiation

Gradients in Neural Networks

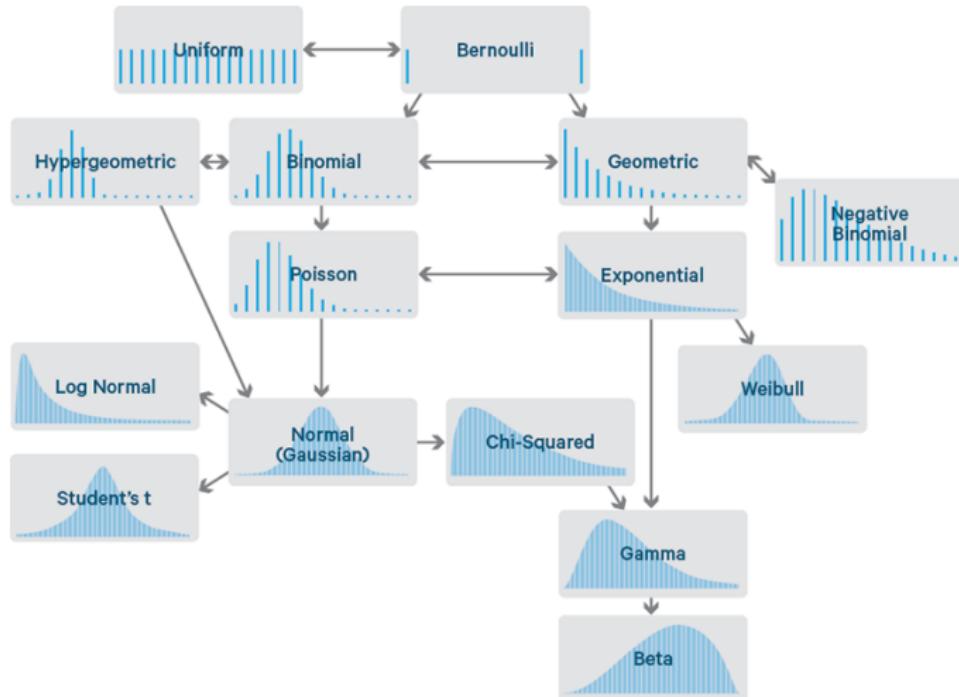
Probability density functions

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

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

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

Probability density functions



\blue{<http://www.math.wm.edu/~leemis/chart/UDR/UDR.html>}

Expectation and Variance

- Expectation of random variable X

Discrete: $E[X] = \sum_{x \in X} x p(x)$, Continuous: $E[X] = \int_X x p(x) dx$

- Let g be any function of X ,

Discrete: $E[g(X)] = \sum_{x \in X} g(x) p(x)$ Continuous: $E[g(X)] = \int_X g(x) p(x) dx$

- Variance:

$$Var[X] = E[(X - E[X])^2] = E[X^2] - (E[X])^2$$

- Covariance of variable X and Y:

$$Cov[X, Y] = E[(X - E[X])(Y - E[Y])]$$

- Some properties:

- ✓ $E[c] = c$, where c is a constant
- ✓ $E[cX] = cE(X)$
- ✓ $E[X + Y] = E[X] + E[Y]$
- ✓ $Var[c] = 0$
- ✓ $Var[c + X] = Var[X]$
- ✓ $Var[cX] = c^2 Var[X]$
- ✓ $Var[X + Y] = Var[X] + Var[Y]$ ONLY IF X and Y are independent!
- ✓ $Var[X + Y] = Var[X] + Var[Y] + 2Cov[X, Y]$ if X and Y are not independent!

Random vectors in high dimensions

- Random vector $X \in \mathbb{R}^n$, $X = \{X_1, X_2, \dots, X_n\}$
- Expectation vector: $\mu = E[X] = \{E[X_1], E[X_2], \dots, E[X_n]\} = \{\mu_1, \mu_2, \dots, \mu_n\}$
- Variance is replaced by the covariance matrix:
$$\Sigma = Cov[X] = E[(X - \mu)(X - \mu)^T] = E[XX^T] - \mu\mu^T$$
- Σ is $n \times n$, symmetric matrix, where
$$Cov[X]_{i,j} = E[(X_i - \mu_i)(X_j - \mu_j)] = E[X_i X_j] - \mu_i \mu_j$$
- For example: Bivariate Normal Distribution

$$\mathbf{x} = \begin{pmatrix} x \\ y \end{pmatrix} \quad \boldsymbol{\mu} = \begin{pmatrix} \mu_X \\ \mu_Y \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \sigma_X^2 & \rho\sigma_X\sigma_Y \\ \rho\sigma_X\sigma_Y & \sigma_Y^2 \end{pmatrix}$$

$$\text{Density } p(\mathbf{x}) = \frac{1}{2\pi} (\det \boldsymbol{\Sigma})^{-1/2} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

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

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

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

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

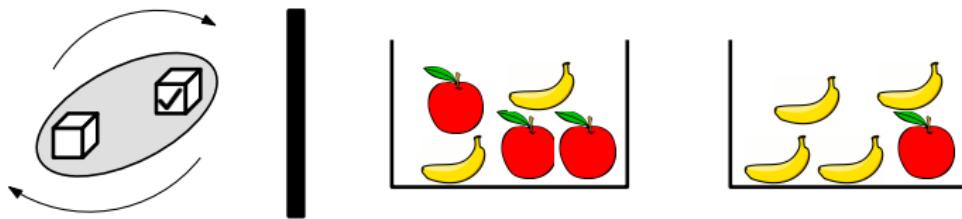
Sum rule/Marginalization property

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

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

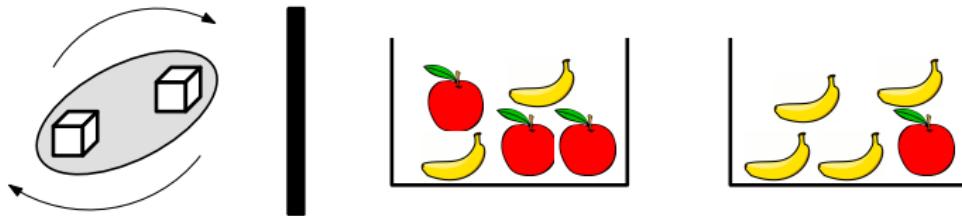


Two cartons $x = \{c_1, c_2\}$ with fruit in each $y = \{\text{banana, apple}\}$

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

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

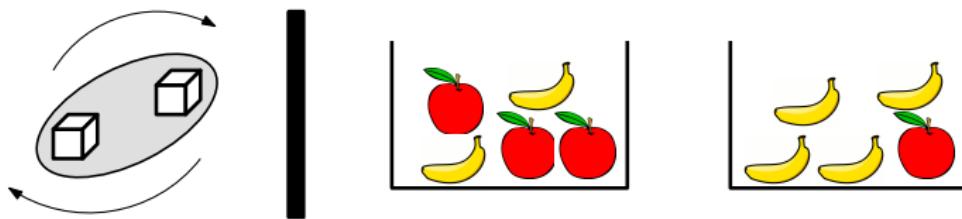
Illustration: Bayesian Inference



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

Prior

Illustration: Bayesian Inference



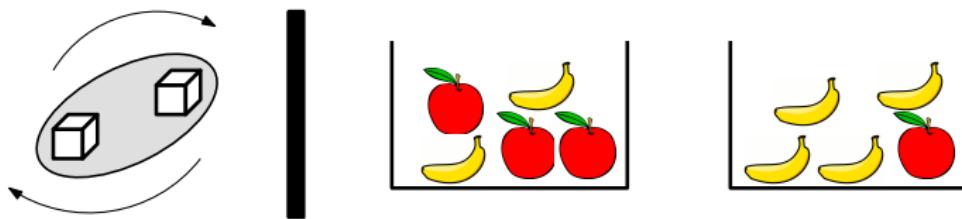
$$p(c_1) = 1/2 = p(c_2)$$

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

Prior

Likelihood

Illustration: Bayesian Inference



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

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

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

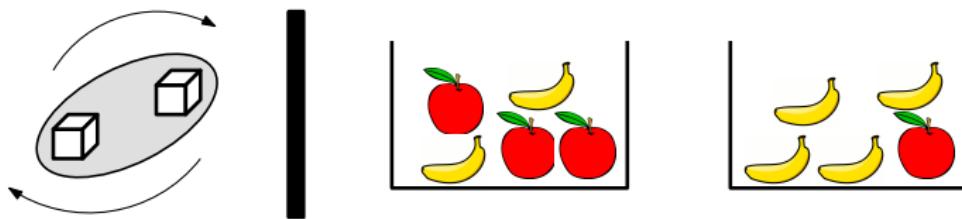
$$p(\text{banana}|c_1)p(c_1) + p(\text{banana}|c_2)p(c_2) = 6/10 = 3/5$$

Prior

Likelihood

Evidence

Illustration: Bayesian Inference



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

Prior

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

Likelihood

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

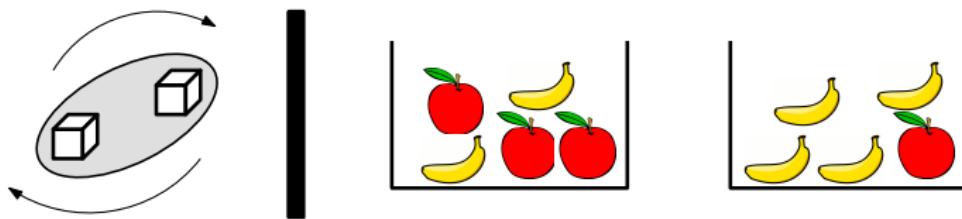
Evidence

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

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

Posterior

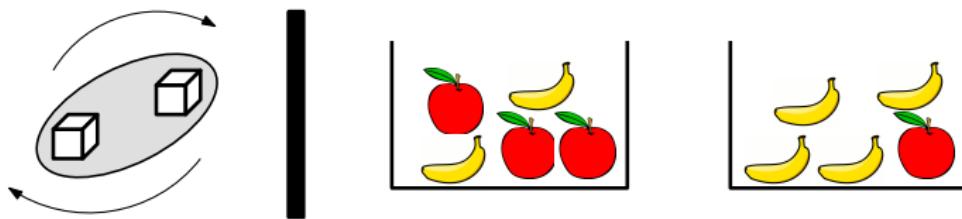
Illustration: Bayesian Inference



$$\begin{aligned} p(c_1) &= 1/2 = p(c_2) && \text{Prior} \\ p(\text{banana}|c_1) &= 2/5, \quad p(\text{banana}|c_2) = 4/5 && \text{Likelihood} \\ p(\text{banana}) &= \sum_x p(x, y) = \sum_x p(y|x)p(x) = && \text{Evidence} \\ &\quad p(\text{banana}|c_1)p(c_1) + p(\text{banana}|c_2)p(c_2) = 6/10 = 3/5 && \\ \text{What is } p(x = c_1 | y = \text{banana})? &&& \text{Posterior} \end{aligned}$$

Bayes' Theorem: Posterior = $\frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$

Illustration: Bayesian Inference



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

Prior

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

Likelihood

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

Evidence

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

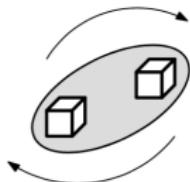
What is $p(x = c_1|y = \text{banana})$?

Posterior

Bayes' Theorem: Posterior = $\frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$

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

This example is analogous to sentiment detection



C: Positive or negative



d1: {good, fair, wonderful, ...}



d2: {poor, bad, fair, ...}

- For a document d and a class c

Naïve Bayes Classifier (I)

$$P(c|d) = \frac{P(d|c)P(c)}{P(d)}$$

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(c|d)$$

$$= \operatorname{argmax}_{c \in C} \frac{P(d|c)P(c)}{P(d)}$$

$$= \operatorname{argmax}_{c \in C} P(d|c)P(c)$$

MAP is "maximum a posteriori" = most likely class

Bayes Rule

Dropping the denominator

Where $p(d|c) = p(w_1, w_2, \dots, w_n|c)$
 $= \prod_{i=1}^n p(w_i|c)$ (independence assumption)

Chain Rule (or Product Rule) of Conditional Probabilities

- Any joint probability distribution over many random variables may be decomposed into conditional distributions over only one variable:

$$P(x^{(1)}, \dots, x^{(n)}) = P(x^{(1)}) \prod_{i=2}^n P(x^{(i)} \mid x^{(1)}, \dots, x^{(i-1)})$$

$$P(a, b, c) = P(a \mid b, c)P(b, c)$$

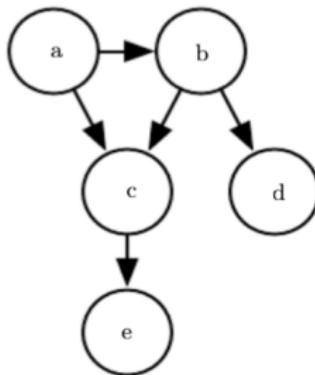
$$P(b, c) = P(b \mid c)P(c)$$

$$P(a, b, c) = P(a \mid b, c)P(b \mid c)P(c)$$

- Sometimes, a variable may only depend on a few other variables, e.g. variable a only depends on b . Thus, $p(a|b, c) = p(a|b)$
- Structured probabilistic model or graph model:** factorizations into conditional probability distributions

Structured Probabilistic Models

- Machine learning algorithms often involve probability distributions over a very large number of random variables, which may interact with each other.
- For example, for a document d , $p(d) = p(w_1, w_2, \dots, w_n)$
- Structured probabilistic model or graphical model:
 - ✓ Node: random variable
 - ✓ Edge: conditional probability of target node given source node



$$p(a, b, c, d, e) = p(a)p(b | a)p(c | a, b)p(d | b)p(e | c).$$

Probability Chain Rule and Language Model

- Try the following tongue twister!

Peter Piper picked a peck of pickled pepper.

Where's the pickled pepper that Peter Piper picked?

Bigram model: $p(w_{n+1}|w_1, w_2, \dots, w_n) = p(w_{n+1}|w_n)$

$P(\text{"Where's the pickled pepper that Peter Piper picked"}) = ?$

$$\begin{aligned} & P\left(\begin{array}{l} \text{Peter Piper picked a peck of pickled pepper} \\ \text{that Peter Piper picked} \end{array}\right) \\ = & P(\text{Peter}|\langle s \rangle) \times P(\text{Piper}|\text{Peter}) \times P(\text{picked}|\text{Piper}) \times \\ & P(\text{a}|\text{picked}) \times P(\text{peck}|\text{a}) \times P(\text{of}|\text{peck}) \times \\ & P(\text{pickled}|\text{of}) \times P(\text{pepper}|\text{pickled}) \times P(\text{that}|\text{pepper}) \times \\ & P(\text{Peter}|\text{that}) \times P(\text{Piper}|\text{Peter}) \times P(\text{picked}|\text{Piper}) \times \\ & P(\langle /s \rangle|\text{picked}) \\ = & .5 \times 1 \times 1 \times .5 \times 1 \times 1 \times 1 \times .5 \times 1 \times 1 \times .5 \\ = & .0625 \end{aligned}$$

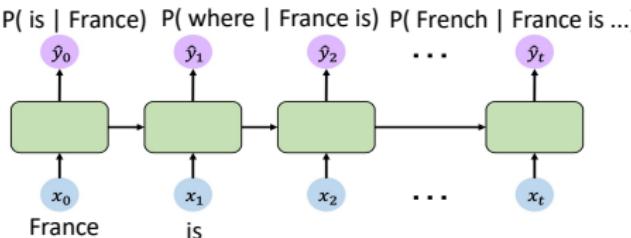
$$P(\text{Piper}|\text{Peter}) = \frac{|\text{Peter Piper}|}{|\text{Peter}|} = \frac{2}{2} = 1$$

Reference:

<https://faculty.sbs.arizona.edu/hammond/archive/ling696f-sp03/snlp4.pdf>

- Neural Network Language Model:

France is where I grew up, but I now live in Boston. I speak fluent ____.



Maximum Likelihood Estimation for supervised learning

- Assume supervised learning with dataset $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$
- A model with parameters θ
- $q_\theta(y_i|x_i)$: the predicted likelihood of y_i from the model for sample x_i .
- Find θ such that the product of the likelihood for all the samples in the dataset is maximized:

$$\operatorname{argmax}_\theta \prod_{i=1}^n q_\theta(y_i|x_i)$$

Equivalent to:

$$\operatorname{argmax}_\theta \sum_{i=1}^n \log q_\theta(y_i|x_i) \quad \text{or}$$

$$\operatorname{argmin}_\theta - \sum_{i=1}^n \log q_\theta(y_i|x_i)$$

Shannon Entropy

- Given a distribution p over a given variable X ,

$$\text{Shannon Entropy: } H(p) = E_p[-\log p]$$

- If X is continuous: $H(p) = -\int_X p(x)\log p(x)dx$
- If X is discrete: $H(p) = -\sum_{x \in X} p(x)\log p(x)$

- Intuition:

- ✓ $H(p)$ quantifies the amount of uncertainty in p
- ✓ It gives a lower bound on the number of bits needed on average to encode symbols drawn from p

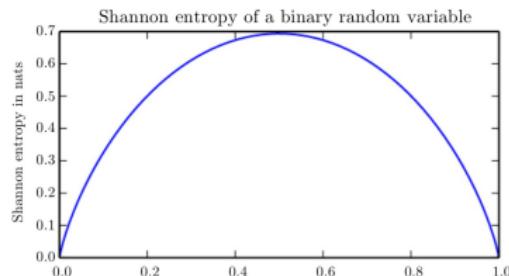


Figure 3.5: This plot shows how distributions that are closer to deterministic have low Shannon entropy while distributions that are close to uniform have high Shannon entropy. On the horizontal axis, we plot p , the probability of a binary random variable being equal to 1. The entropy is given by $(p-1)\log(1-p) - p\log p$. When p is near 0, the distribution is nearly deterministic, because the random variable is nearly always 0. When p is near 1, the distribution is nearly deterministic, because the random variable is nearly always 1. When $p = 0.5$, the entropy is maximal, because the distribution is uniform over the two outcomes.

Quick Quiz!

Cross Entropy

- Given two distributions p and q over a given variable X , cross entropy of q relative to p is:

$$\text{Cross Entropy } H(p, q) = E_p[-\log q]$$

- If X is continuous:

$$H(p, q) = - \int_X p(x) \log q(x) dx$$

- If X is discrete:

$$H(p, q) = - \sum_{x \in X} p(x) \log q(x)$$

- $H(p, q)$ measures how many bits needed on average to encode data if a coding scheme is optimized for an estimated probability distribution q rather than the true distribution p
- $H(p, q)$ is used as a cost function. The more $p(x)$ approximates $q(x)$, the smaller the cost. Assuming x is binary, try the following:
 - If $p(x = 1) = 0.99$ while $q(x = 1) = 0.01$, $H(p, q) = ?$
 - If $p(x = 1) = 0.99$ while $q(x = 1) = 0.98$, $H(p, q) = ?$

Question: how to calculate cross entropy cost
in single label or multi-label classification?

Cross Entropy Cost and Maximum Likelihood Estimation

- Again, supervised learning with dataset $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$
- $p(y|x_i)$: the distribution for the ground truth label (y) for feature x_i :

$$p(y|x_i) = \begin{cases} 1, & y = y_i \\ 0, & \text{otherwise} \end{cases}$$

- $q_\theta(y|x_i)$: the predicted likelihood of y from the model for sample x_i
- Cross entropy of $q_\theta(y|x_i)$ relative to $p(y|x_i)$:

$$H_i(p, q_\theta) = - \sum_{y \in Y} p(y|x_i) \log q_\theta(y|x_i) = -\log q_\theta(y_i|x_i)$$

- Total cross entropy cost (loss function) :

$$L = \sum_{i=1}^n H_i(p, q_\theta) = \sum_{i=1}^n -\log q_\theta(y_i|x_i)$$

- Thus the optimization goal is:

$$\operatorname{argmin}_\theta L = \sum_{i=1}^n -\log q_\theta(y_i|x_i)$$

- This is exactly Maximum Likelihood Estimation

Cross Entropy Cost and Maximum Likelihood Estimation

- Again, supervised learning with dataset $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$
- $p(y|x_i)$: the distribution for the ground truth label (y) for feature x_i :
$$p(y|x_i) = \begin{cases} 1, & y = y_i \\ 0, & \text{otherwise} \end{cases}$$
- $q_\theta(y|x_i)$: the predicted likelihood of y from the model for sample x_i
- Cross entropy of $q_\theta(y|x_i)$ relative to $p(y|x_i)$:

$$H_i(p, q_\theta) = - \sum_{y \in Y} p(y|x_i) \log q_\theta(y|x_i) = - \log q_\theta(y_i|x_i)$$

- Total **cross entropy cost** (loss function) :

$$L = \sum_{i=1}^n H_i(p, q_\theta) = \sum_{i=1}^n - \log q_\theta(y_i|x_i)$$

- Thus the optimization goal is:

$$\operatorname{argmin}_\theta L = \sum_{i=1}^n - \log q_\theta(y_i|x_i)$$

- This is exactly **Maximum Likelihood Estimation**

Kullback–Leibler (KL) Divergence

- Given two distributions p and q over a given variable X , KL divergence measures how different these two distributions are:

$$D_{KL}(p||q) = E_p[-\log \frac{q}{p}]$$

- If X is continuous: $D_{KL}(p||q) = -\int_X p(x)\log \frac{q(x)}{p(x)} dx$
- If X is discrete: $D_{KL}(p||q) = -\sum_{x \in X} p(x)\log \frac{q(x)}{p(x)}$
- $D_{KL}(p||q)$ properties:
 - Non-negative (let $f(x) = -\log \frac{q(x)}{p(x)}$, show it by using Jensen inequality: $E(f(x)) \geq f(E(x))$ when $f(x)$ is convex)
 - $D_{KL}(p||q) = 0$ if and only if p and q are equal
 - Not symmetric, $D_{KL}(p||q) \neq D_{KL}(q||p)$ (not a true distance measure)
- It is also called relative entropy:

$$D_{KL}(p||q) = E_p[-\log \frac{q}{p}] = E_p[-\log q + \log p] = H(p, q) - H(p)$$

- Thus, $H(p, q) = D_{KL}(p||q) + H(p)$

Cross Entropy and KL Divergence

- Optimizing cross entropy :

$$\operatorname{argmin}_{\theta} L = \operatorname{argmin}_{\theta} \sum_{i=1}^n H_i(p, q_{\theta})$$

$$= \operatorname{argmin}_{\theta} \sum_{i=1}^n [D_{KL,i}(p||q_{\theta}) + H_i(p)]$$

$$= \operatorname{argmin}_{\theta} \sum_{i=1}^n [D_{KL,i}(p||q_{\theta})]$$

because $H_i(p)$ is independent of θ

- Thus, in classification problems, optimizing the sum of **cross entropy** over all the training samples is equivalent to
 - Optimizing the sum of **KL divergence** over all the training samples
 - Maximum likelihood estimation**

Overview

Key Concepts in Probability

Probability density function

Sum & Product Rules

Bayes' Theorem

Maximum Likelihood Estimation , Cross Entropy, KL Divergence

Differentiation

Linear algebra and notations

Where it is used in ML?

Types of differentiation

Gradients in Neural Networks

Linear Algebra Notation

- **Scalar:** A single number
- **Vector:** An array of numbers, i.e. $x \in \mathbb{R}^n$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \text{ or } x = [x_1, x_2, \dots, x_n]^T$$

- **Matrix:** A 2-D array of numbers, e.g. $A \in \mathbb{R}^{m \times n}$
 - Each element is identified by two indices, e.g. $A_{i,j}$
 - $A_{i,:}$ is the i th row and $A_{:,j}$ is the j th column
 - Transpose: $A^T \in \mathbb{R}^{n \times m}$, $(A^T)_{i,j} = A_{j,i}$
 - A **vector** can be thought of as a **matrix with only one column**
 - The **transpose of a vector** is therefore a **matrix with only one row**
 - A scalar a can be thought of as a matrix with only one value. $a^T = a$.
- **Tensor:** An array of numbers with more than two axes is known as a tensor. The element of a tensor $\mathbf{A} \in \mathbb{R}^{m \times n \times p}$ at coordinates (i, j, k) is denoted as $\mathbf{A}_{i,j,k}$.

Background: Matrix Multiplication

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

The diagram illustrates the non-commutativity of matrix multiplication through two examples. In the first example, a 2x3 blue matrix is multiplied by a 3x2 yellow matrix, resulting in a 2x2 green matrix. In the second example, a 3x2 yellow matrix is multiplied by a 2x3 blue matrix, resulting in a 3x3 green matrix.

$$\begin{matrix} \text{blue} & \times & \text{yellow} & = & \text{green} \\ \begin{matrix} 1 & 2 \\ 3 & 4 \end{matrix} & \times & \begin{matrix} 5 & 6 \\ 7 & 8 \\ 9 & 0 \end{matrix} & = & \begin{matrix} 15 & 16 & 17 \\ 27 & 28 & 29 \end{matrix} \\ \text{yellow} & \times & \text{blue} & = & \text{green} \\ \begin{matrix} 5 & 6 \\ 7 & 8 \\ 9 & 0 \end{matrix} & \times & \begin{matrix} 1 & 2 \\ 3 & 4 \end{matrix} & = & \begin{matrix} 17 & 18 & 19 \\ 39 & 40 & 41 \\ 61 & 62 & 63 \end{matrix} \end{matrix}$$

Background: Matrix Multiplication

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

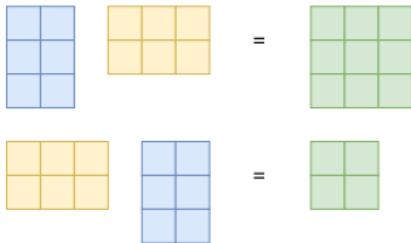
$$\begin{array}{ccc} \begin{matrix} \text{blue} \\ 2 \times 3 \end{matrix} & \times & \begin{matrix} \text{yellow} \\ 3 \times 2 \end{matrix} \\ = & & \begin{matrix} \text{green} \\ 2 \times 2 \end{matrix} \end{array}$$
$$\begin{array}{ccc} \begin{matrix} \text{yellow} \\ 3 \times 2 \end{matrix} & \times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

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

$$\underbrace{\begin{matrix} A \\ n \times k \end{matrix}}_{\text{dimensions}} \quad \underbrace{\begin{matrix} B \\ k \times m \end{matrix}}_{\text{dimensions}} = \underbrace{\begin{matrix} C \\ n \times m \end{matrix}}_{\text{dimensions}}$$

Background: Matrix Multiplication

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



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

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

| | |
|--------------------------------|--|
| | <code>import numpy as np</code> |
| $y = Ax$ | <code>y = A.dot(x)</code> |
| $y_i = \sum_j A_{ij}x_j$ | <code>y = np.einsum('ij, j', A, x)</code> |
| $C = AB$ | <code>C = A.dot(B)</code> |
| $C_{ij} = \sum_k A_{ik}B_{kj}$ | <code>C = np.einsum('ik, kj', A, B)</code> |

Multiplying Matrices and Vectors

- **Matrix Product:**

$$\mathbf{C} = \mathbf{AB}, \text{ where } C_{i,j} = \sum_k A_{i,k}B_{k,j}$$

- **Element wise product** (Hadamard product):

$$\mathbf{C} = \mathbf{A} \odot \mathbf{B}, \text{ where } C_{i,j} = A_{i,j}B_{i,j}$$

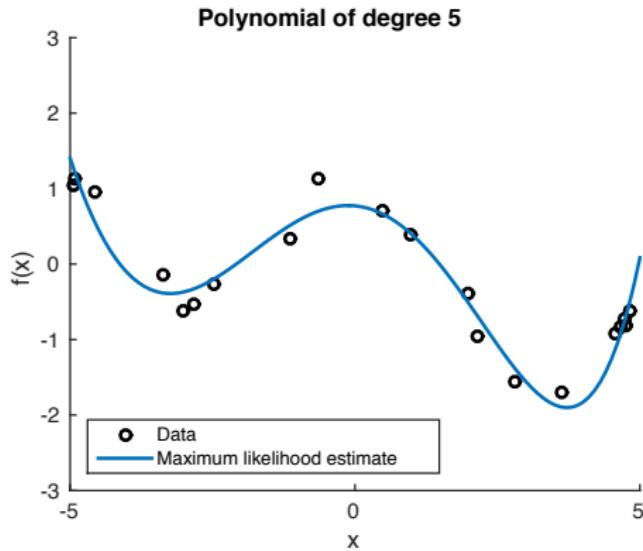
- **Dot product** between two vectors \mathbf{x} and \mathbf{y} ($\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$) is $\mathbf{x}^T\mathbf{y}$

- What is the dimension of $\mathbf{x}^T\mathbf{y}$?
- We can think of the matrix product $\mathbf{C} = \mathbf{AB}$ as computing $C_{i,j}$ as the dot product between row i of \mathbf{A} and column j of \mathbf{B} .

- Useful properties:

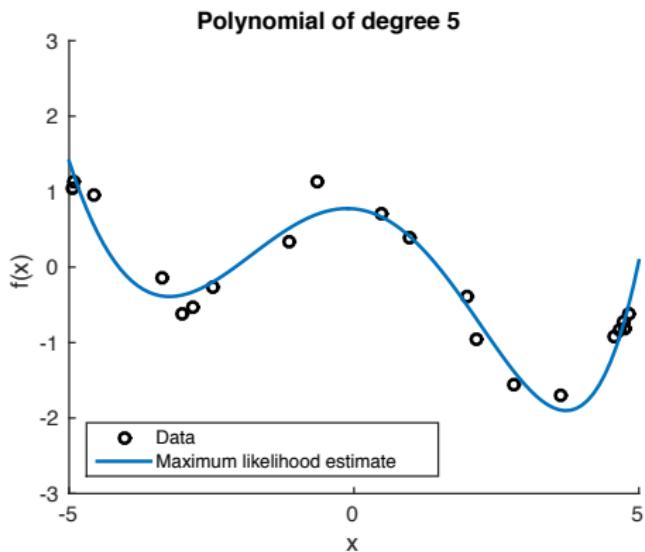
- $\mathbf{AB} \neq \mathbf{BA}, \quad \mathbf{A(BC)} = (\mathbf{AB})\mathbf{C}, \quad \mathbf{A(B+C) = AB+AC}$
- $(\mathbf{AB})^T = \mathbf{B}^T\mathbf{A}^T$
- $\mathbf{x}^T\mathbf{y} = \mathbf{y}^T\mathbf{x}$ (how to demonstrate it?)

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

- 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

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

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

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

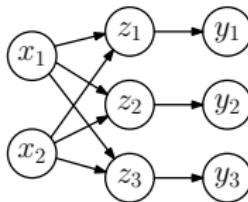
- Define a **loss function**, $L_\theta(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^N (y_i - f_\theta(x_i))^2$
 - Minimisation of loss function is typically based on **gradient descent**
- **Differentiation** required to compute gradients of L_θ w.r.t $\boldsymbol{\theta}$

Scenario 2: Find optimal parameters of a Neural Network

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

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

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



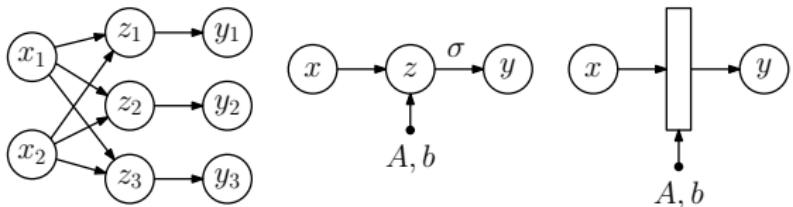
Scenario 2: Find optimal parameters of a Neural Network

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

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

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

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



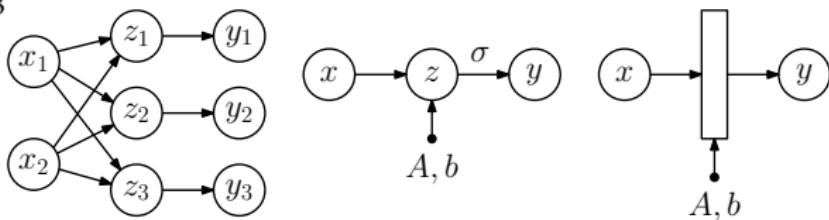
Scenario 2: Find optimal parameters of a Neural Network

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

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

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

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



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

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

Similarly, parameters \mathbf{A} , \mathbf{b} are tuned to minimize the prediction error through gradient descent

Again, differentiation is used to compute gradients of the loss function with regard to matrix \mathbf{A} and \mathbf{b}

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

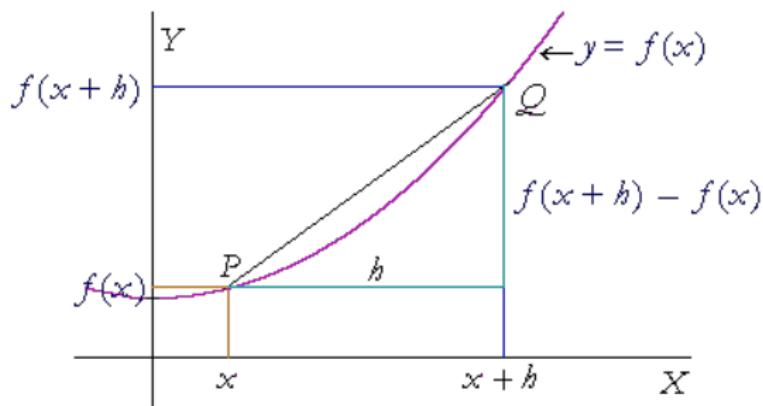
$$\left(\frac{f(x)}{g(x)}\right)' = \frac{f(x)'g(x) - f(x)g(x)'}{(g(x))^2} = \frac{\frac{df}{dx}g(x) - f(x)\frac{dg}{dx}}{(g(x))^2}$$

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 $f(x+h)$



Example: Scalar Chain Rule

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

Beginner

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

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

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

$$= -12$$

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:

$$\frac{df}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \cdots & \frac{\partial f}{\partial x_N} \end{bmatrix} \in \mathbb{R}^{1 \times N}$$

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

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

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

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

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

Example: Multivariate 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 $dL/d\mathbf{x}$. What is the dimension/size of $dL/d\mathbf{x}$?

$$\frac{dL}{d\mathbf{x}} = \frac{dL}{d\mathbf{e}} \frac{d\mathbf{e}}{d\mathbf{x}}$$

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

$$\frac{d\mathbf{e}}{d\mathbf{x}} = -\mathbf{A} \in \mathbb{R}^{M \times N}$$

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

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)

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

Example: Vector field differentiation

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

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

- ▶ Compute the gradient $df/d\mathbf{x}$
 - ▶ Dimension of $df/d\mathbf{x}$:

Example: Vector field differentiation

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

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

- ▶ Compute the gradient $df/d\mathbf{x}$

- ▶ Dimension of $df/d\mathbf{x}$: Since $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$, $df/d\mathbf{x} \in \mathbb{R}^{M \times N}$

- ▶ Gradient:

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

What about
 df/dA ?
Dimension?

$$\Rightarrow \frac{df}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_N} \\ \vdots & \vdots & \vdots \\ \frac{\partial f_M}{\partial x_1} & \cdots & \frac{\partial f_M}{\partial x_N} \end{bmatrix} = \begin{bmatrix} A_{11} & \cdots & A_{1N} \\ \vdots & & \vdots \\ A_{M1} & \cdots & A_{MN} \end{bmatrix} = A$$

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)

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

Example 1: Derivatives w.r.t Matrices

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

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

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

Example 1: Derivatives w.r.t Matrices

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

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

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

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

Example 2: Derivatives w.r.t Matrices

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

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

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

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

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

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

Example 2: Derivatives w.r.t Matrices

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

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

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

Example 2: Derivatives w.r.t Matrices

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

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

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

Example 2: Derivatives w.r.t Matrices

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

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

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

Example 2: Derivatives w.r.t Matrices

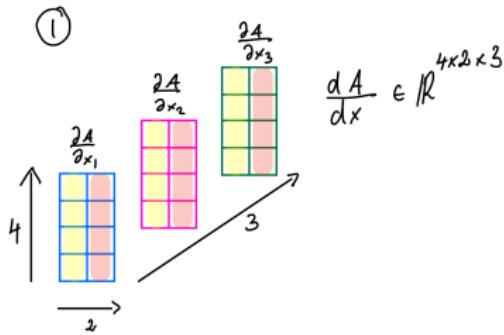
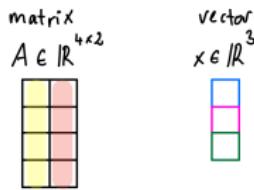
$$f_i = \sum_{j=1}^N A_{ij}x_j, \quad i = 1, \dots, M$$

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

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

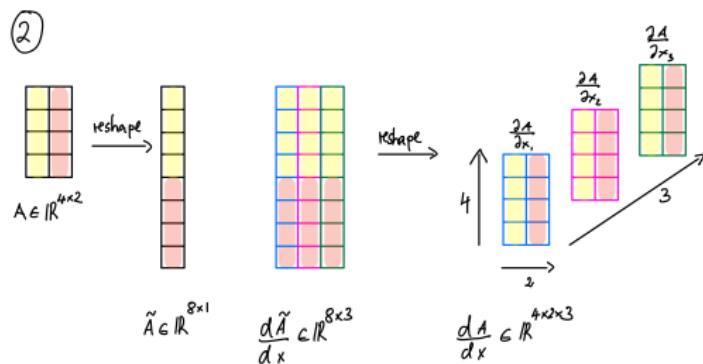
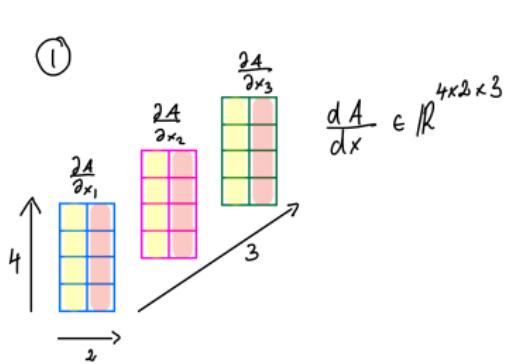
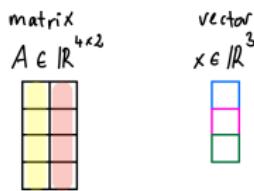
Example 3: Derivatives w.r.t Higher-Order Tensors

- 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:

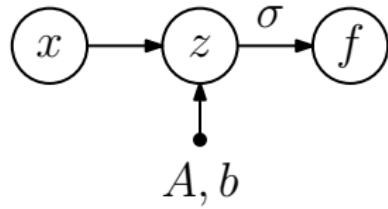


Example 3: Derivatives w.r.t Higher-Order Tensors

- 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:



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

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

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

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

Gradients of a Single-Layer Neural Network

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

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

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

Gradients of a Single-Layer Neural Network

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

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

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

Gradients of a Single-Layer Neural Network

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

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

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

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

Gradients of a Single-Layer Neural Network

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

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

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

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

Gradients of a Single-Layer Neural Network

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

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

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

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

Putting Things Together

- ▶ Inputs $x \in \mathbb{R}^N$

Putting Things Together

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

Putting Things Together

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

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

Putting Things Together

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

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

- ▶ Find A^*, b^* , such that the squared loss

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

is minimised

Putting Things Together

Partial derivatives:

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{A}} &= \underbrace{\frac{\partial L}{\partial \mathbf{e}}}_{\in \mathbb{R}^{1 \times M}} \underbrace{\frac{\partial \mathbf{e}}{\partial f}}_{\in \mathbb{R}^{M \times M}} \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{\in \mathbb{R}^{M \times N}} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{A}}}_{\in \mathbb{R}^{N \times M}} \\ \frac{\partial L}{\partial \mathbf{b}} &= \underbrace{\frac{\partial L}{\partial \mathbf{e}}}_{\in \mathbb{R}^{1 \times M}} \underbrace{\frac{\partial \mathbf{e}}{\partial f}}_{\in \mathbb{R}^{M \times M}} \underbrace{\frac{\partial f}{\partial \mathbf{z}}}_{\in \mathbb{R}^{M \times N}} \underbrace{\frac{\partial \mathbf{z}}{\partial \mathbf{b}}}_{\in \mathbb{R}^{N \times M}}\end{aligned}$$

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

Gradient descent for neural network optimisation

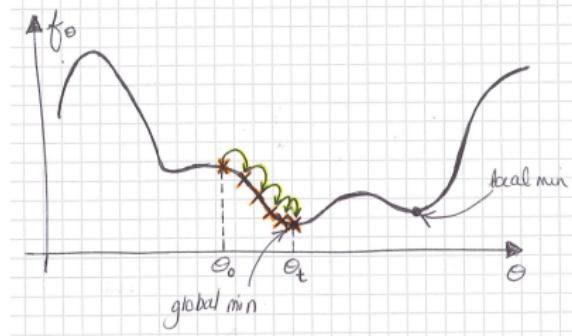
Goal: Find the local optimum of loss $L(\theta^*) \in \mathbb{R}$ where $\theta^* = \{A^*, b^*\}$

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 $\boldsymbol{\theta}^* = \{A^*, \mathbf{b}^*\}$

Gradient descent algorithm:

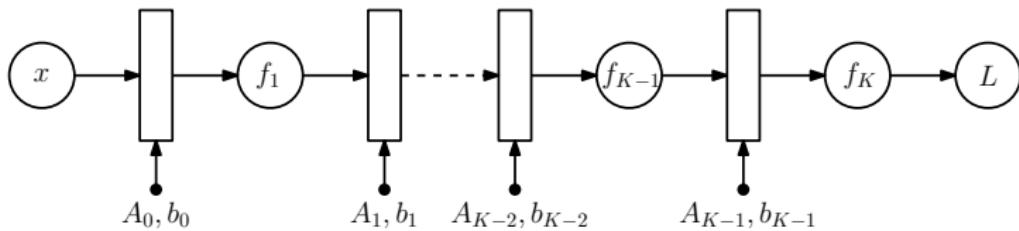
1. Initialise neural network parameters, $\boldsymbol{\theta}_0 = \{A_0, \mathbf{b}_0\}$
2. Compute gradient of L w.r.t. parameters, $\frac{\partial L}{\partial \boldsymbol{\theta}_0}$
3. Update initial guess $\boldsymbol{\theta}_0$ using chain rule:

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

► 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

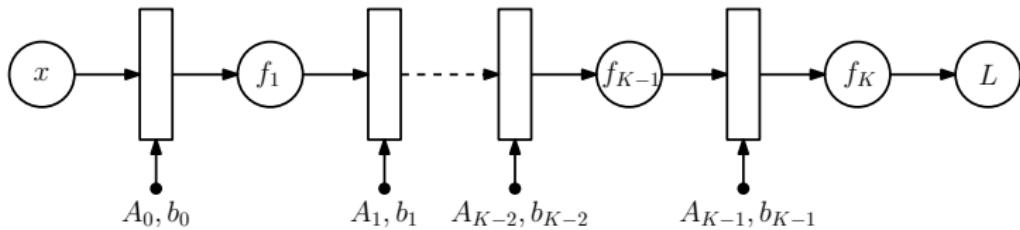


- ▶ Inputs x , observed outputs y
- ▶ Train K-layer neural network with

$$f_0 = x$$

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

Gradients of a Multi-Layer Neural Network



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

$$f_0 = x$$

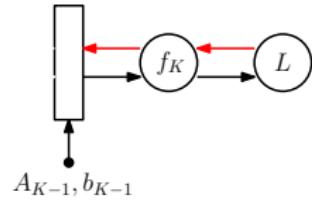
$$f_i = \sigma_i(A_{i-1}f_{i-1} + b_{i-1}), \quad i = 1, \dots, K$$

- ▶ Find A_k^*, b_k^* for $k = 0, \dots, K - 1$, such that the squared loss

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

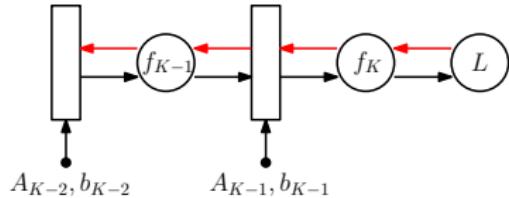
is minimised, where $\theta = \{A_k, b_k\}$, $k = 0, \dots, K - 1$

Gradients of a Multi-Layer Neural Network



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

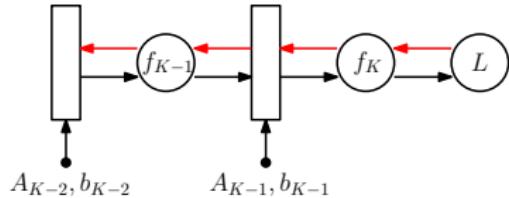
Gradients of a Multi-Layer Neural Network



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

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

Gradients of a Multi-Layer Neural Network

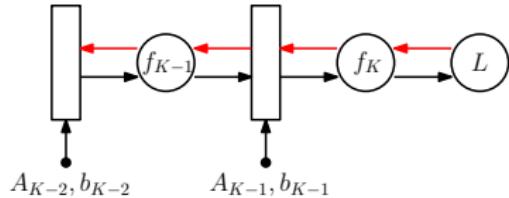


$$\frac{\partial L}{\partial \theta_{K-1}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial \theta_{K-1}}$$

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

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

Gradients of a Multi-Layer Neural Network



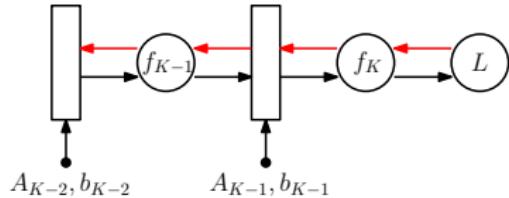
$$\frac{\partial L}{\partial \theta_{K-1}} = \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial \theta_{K-1}}$$

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

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

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

Gradients of a Multi-Layer Neural Network



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

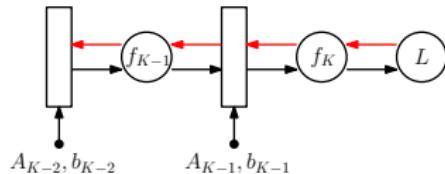
$$\frac{\partial L}{\partial \theta_{K-2}} = \frac{\partial L}{\partial f_K} \boxed{\frac{\partial f_K}{\partial f_{K-1}} \frac{\partial f_{K-1}}{\partial \theta_{K-2}}}$$

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

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

► Intermediate derivatives are stored during the forward pass

Gradients of a Multi-Layer Neural Network


 δ_K

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

 \implies

$$\delta_K = \frac{\partial L}{\partial f_K}$$

Backpropagation

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

 \implies

$$\delta_{K-1} = \delta_K \frac{\partial f_K}{\partial f_{K-1}}$$

Backpropagate δ from K to $K-1, K-2, \dots, i, \dots$

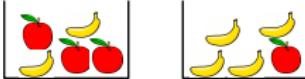
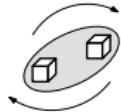
$$\frac{\partial L}{\partial \theta_{K-3}} = \boxed{\frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial f_{K-1}} \frac{\partial f_{K-1}}{\partial f_{K-2}} \frac{\partial f_{K-2}}{\partial \theta_{K-3}}}$$

 \vdots
 \vdots

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

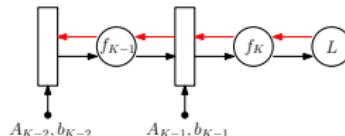
How to calculate gradients:
 1. Given x and θ , first do a forward pass to calculate L and f at each layer
 2. Then follow these equations to get partial derivatives of each θ

Summary



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

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



(1)

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

- ▶ 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

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

$$f(x) = f(x_1, x_2) = x_1^2 + 2x_2,$$

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

Example: Combined Chain Rule

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

$$f(x) = f(x_1, x_2) = x_1^2 + 2x_2,$$

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

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

1×2 and 2×1

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

Work it out with your neighbours

Example: Combined Chain Rule

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

$$f(x) = f(x_1, x_2) = x_1^2 + 2x_2,$$

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

- What are the dimensions of df/dx and dx/dt ?

1×2 and 2×1

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

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

```

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

```

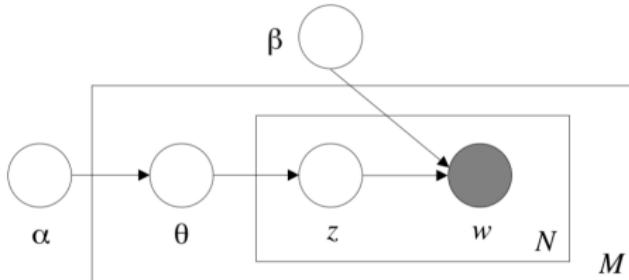
$$\frac{\partial f}{\partial \mathbf{x}}[i,j] = \sum_{l=1}^M \frac{\partial f}{\partial z}[i,l] \frac{\partial z}{\partial \mathbf{x}}[l,j]$$

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

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

Structured Probabilistic Model of LDA

- Remember this?



Graphical model representation of LDA. The boxes are “plates” representing replicates. The outer plate represents documents, while the inner plate represents the repeated choice of topics and words within a document.

$$p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta) = p(\theta | \alpha) \prod_{n=1}^N p(z_n | \theta) p(w_n | z_n, \beta)$$

LDA assumes the following generative process for each document \mathbf{w} in a corpus D :

1. Choose $N \sim \text{Poisson}(\xi)$.
2. Choose $\theta \sim \text{Dir}(\alpha)$.
3. For each of the N words w_n :

- (a) Choose a topic $z_n \sim \text{Multinomial}(\theta)$.
- (b) Choose a word w_n from $p(w_n | z_n, \beta)$, a multinomial probability conditioned on the topic z_n .

Note: β is the word probabilities per topic, a $k \times V$ matrix, $\beta_{ij} = p(w_j = 1 | z_i = 1)$. β is a parameter to be estimated.