

*Timothy Chung*

COMP70015

Mathematics for Machine Learning

AUTUMN 2024



# Contents

<b>I</b>	<b>Fundamental Concepts in Mathematics for Machine Learning</b>	<b>5</b>
<b>0</b>	<b>Preliminaries</b>	<b>6</b>
0.1	Course Resources	6
0.2	Course Notation	6
0.3	Linear Algebra	6
0.4	Calculus	9
0.5	Multivariate Calculus	9
0.6	Probability and Statistics	10
0.7	Moments	11
0.8	Independent and Identically Distributed (i.i.d.) Random Variables	11
0.9	Einstein/Pythonic Index Notation	12
<b>1</b>	<b>Formalising ML Problem Settings</b>	<b>15</b>
1.1	Machine Learning from the Mathematical Perspective	15
1.2	Datasets	15
1.3	Types of Learning Tasks	17
1.4	The Linear Model	18
1.5	Data Preprocessing: Normalising the Data	19
<b>2</b>	<b>Linear Regression</b>	<b>21</b>
2.1	Learning as Optimisation	21
2.2	Ordinary Least Squares (OLS)	24
2.3	Basis Expansion	27
<b>3</b>	<b>Generalised Linear Models &amp; Gradient Descent</b>	<b>30</b>
3.1	Generalised Linear Models	30
3.2	Gradient Descent	33
<b>4</b>	<b>Deep Learning &amp; Automatic Differentiation</b>	<b>38</b>
4.1	From GLMs to Multi-Layer Perceptrons	38
4.2	Automatic Differentiation	43
<b>5</b>	<b>Convexity, Convergence and Optimisation</b>	<b>51</b>
5.1	Learning as Optimisation	51
5.2	Progress Bounds for Gradient Descent	52
5.3	Convex Optimisation	53
5.4	Convergence in Machine Learning	56
<b>6</b>	<b>Concentration Inequalities &amp; Generalisation Bounds</b>	<b>60</b>
6.1	Generalisation as a Concept	60
6.2	Test Sets	60

6.3	Concentration and the Weak Law of Large Numbers . . . . .	60
6.4	Universal Function Approximation . . . . .	62
6.5	Generalisation Error Bounds . . . . .	63
<b>7</b>	<b>Cross-Validation and Bias-Variance Trade-Off . . . . .</b>	<b>65</b>
7.1	Model Selection . . . . .	65
7.2	Regularisation . . . . .	65
7.3	Cross-Validation . . . . .	66
7.4	Bias and Variance of Estimators . . . . .	66
7.5	Bias-Variance Decomposition . . . . .	68
7.6	Bias-Variance Trade-Off in ML Models . . . . .	68
<b>8</b>	<b>Introduction to Probability . . . . .</b>	<b>73</b>
8.1	Sample and Event Space . . . . .	73
8.2	Sigma Algebra . . . . .	74
8.3	Measure Theory . . . . .	74
8.4	Joint Probability Distributions . . . . .	75
<b>9</b>	<b>Reference Boxes . . . . .</b>	<b>76</b>
9.1	Headings . . . . .	78
9.2	Environments . . . . .	78
<b>10</b>	<b>On the Use of the <code>tufte-book</code> Document Class . . . . .</b>	<b>80</b>
10.1	Page Layout . . . . .	80
10.2	Sidenotes . . . . .	80
10.3	References . . . . .	81
10.4	Figures and Tables . . . . .	81
10.5	Captions . . . . .	83
10.6	Full-width text blocks . . . . .	83
10.7	Typography . . . . .	84
10.8	Document Class Options . . . . .	84
<b>11</b>	<b>Customizing Tufte-LaTeX . . . . .</b>	<b>86</b>
11.1	File Hooks . . . . .	86
11.2	Numbered Section Headings . . . . .	86
11.3	Changing the Paper Size . . . . .	87
11.4	Customizing Marginal Material . . . . .	87
<b>12</b>	<b>Compatibility Issues . . . . .</b>	<b>89</b>
12.1	Converting from <code>article</code> to <code>tufte-handout</code> . . . . .	89
12.2	Converting from <code>book</code> to <code>tufte-book</code> . . . . .	89
<b>13</b>	<b>Troubleshooting and Support . . . . .</b>	<b>90</b>
13.1	Tufte-LaTeX Website . . . . .	90
13.2	Tufte-LaTeX Mailing Lists . . . . .	90
13.3	Getting Help . . . . .	90
13.4	Errors, Warnings, and Informational Messages . . . . .	90
13.5	Package Dependencies . . . . .	91

**Part I**  
**Fundamental Concepts in Mathematics for**  
**Machine Learning**

# 0

## Preliminaries

### 0.1 Course Resources

Course Resources	Reference 0.1.1
<ul style="list-style-type: none"><li>• The textbook, <i>Mathematics for Machine Learning</i>, to be referenced as Deisenroth et al. (2020), can be found <a href="#">here</a></li><li>• The textbook, <i>Machine Learning a Probabilistic Perspective</i>, to be referenced as Murphy (2012), can be found <a href="#">here</a></li><li>• The textbook, <i>Pattern Recognition and Machine Learning</i>, to be referenced as Bishop (2006), can be found <a href="#">here</a></li></ul>	
Additional Resources	Reference 0.1.2
<ul style="list-style-type: none"><li>• <b>Linear Algebra:</b> Chapter 2 (page 29) <i>Bengio et al. (2017)</i> <a href="#">link</a></li><li>• <b>Probability Theory Review:</b> Chapter 2 (page 27) <i>Murphy (2012)</i>. <i>Machine learning: a probabilistic perspective</i>. MIT Press.</li><li>• <b>Calculus Review:</b> Chapter 5 (page 139) <i>Deisenroth et al. (2020)</i>. <i>Mathematics for machine learning</i>. Cambridge University Press.</li></ul>	

### 0.2 Course Notation

Real values are denoted with lower-case Roman letters, e.g.  $x, y, z$ . As convention, we will use a bold Roman letter or a Greek symbol to denote a vector, e.g.  $\mathbf{x}, \boldsymbol{\theta}$ . Do not confuse this with random variables, which are denoted using bold, non-italic Roman letters, e.g.  $\mathbf{x}, \mathbf{y}, \mathbf{z}$ . Matrices are represented by upper-case Roman letters such as  $X$  or random matrices like  $\mathbf{X}$ . Greek letters will always have their type specified, except for  $\boldsymbol{\theta}$ , which will be used without loss of generality to represent model parameters.

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

### 0.3 Linear Algebra

A core prerequisite for this class is linear algebra. Please use this section and references herein to brush up on necessary concepts.

#### 0.3.1 Vectors

A vector is defined by a list of numbers. It is most useful to geometrically interpret vectors as points in space.

Notation Type	Example
Real Values	$x, y, z$
Vectors	$\mathbf{x}, \boldsymbol{\theta}$ (where $\boldsymbol{\theta}$ represents parameters)
Random Variables	$\mathbf{x}, \mathbf{y}, \mathbf{z}$ Example: $\mathbf{z} \sim \text{Unif}(0, 1)$
Matrices	$X = \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{pmatrix}$
Random Matrices	$\mathbf{X}$ Example: $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

Table 1: Summary of Course Notation with Examples

$$\mathbf{x} := \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad (1)$$

### 0.3.2 Vectors Operations

We have several vector operations:

- **Scalar Multiplication:**

$$\alpha \mathbf{x} = \begin{bmatrix} \alpha x_1 \\ \alpha x_2 \\ \vdots \\ \alpha x_n \end{bmatrix} \quad (2)$$

- **Vector Addition:**

$$\mathbf{x} + \mathbf{y} = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \vdots \\ x_n + y_n \end{bmatrix} \quad (3)$$

- **Dot Product:**

$$\mathbf{x}^\top \mathbf{y} = \sum_{i=1}^n x_i y_i \quad (4)$$

Symbol	Meaning
$\mathbf{v}^\top$	Transpose of $\mathbf{v}$
$\mathbf{u}^\top \mathbf{v}$	Inner (scalar) product
$\mathbf{u} \mathbf{v}^\top$	Outer product ( $n \times n$ matrix)
$\mathbf{u} \odot \mathbf{v}$	Hadamard/Elementwise product, $[u_1 v_1, \dots, u_n v_n]$
$\dim(\mathbf{v})$	Dimensionality of $\mathbf{v}$ (namely $n$ )
$\text{diag}(\mathbf{v})$	Diagonal $n \times n$ matrix made from vector $\mathbf{v}$
$\mathbf{1}$ or $\mathbf{1}_n$	Vector of ones (of length $n$ )
$\mathbf{0}$ or $\mathbf{0}_n$	Vector of zeros (of length $n$ )
$\ \mathbf{v}\ $ or $\ \mathbf{v}\ _2$	Euclidean or $\ell_2$ norm, $\sqrt{\sum_{i=1}^n v_i^2}$
$\ \mathbf{v}\ _1$	$\ell_1$ norm, $\sum_{i=1}^n  v_i $
$\mathbf{e}_i$	$i$ th-basis vector (1 in dimension $i$ and 0 elsewhere)

Table 2: Standard Notation for Vectors

### 0.3.3 Matrices

Matrices define linear transforms. Geometrically, it is best to interpret them as transformations whose columns are the new basis vectors.

$$A := \begin{bmatrix} A_{1,1} & A_{1,2} & \cdots & A_{1,n} \\ A_{2,1} & A_{2,2} & \cdots & A_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{m,1} & A_{m,2} & \cdots & A_{m,n} \end{bmatrix} \quad A \in \mathbb{R}^{m \times n} \quad (5)$$

Non-square matrices can be thought of as transformations from  $\mathbb{R}^m$  dimensional space to  $\mathbb{R}^n$  dimensional space.

### 0.3.4 Matrix Operations

We define several matrix operations as follows:

- **Scalar Multiplication:**

$$cA = cA_{i,j} \quad (6)$$

- **Hadamard Product (Element-wise Multiplication):**

$$A \circ B = A_{i,j} B_{i,j} \quad (7)$$

- **Matrix Subtraction:**

$$A - B = A_{i,j} - B_{i,j} \quad (8)$$

- **Matrix Product:**

$$C_{i,j} = \sum_k A_{i,k} B_{k,j} \quad (9)$$

- **Trace of a Matrix:**

$$\text{Tr}(A) = \sum_i A_{i,i} \quad (10)$$

- **Kronecker Product:**

$$A \otimes B = \begin{pmatrix} a_{11}B & a_{12}B & \cdots & a_{1n}B \\ a_{21}B & a_{22}B & \cdots & a_{2n}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}B & a_{m2}B & \cdots & a_{mn}B \end{pmatrix} \quad (11)$$

The Kronecker product  $A \otimes B$  results in a block matrix, where each element  $a_{ij}$  of matrix  $A$  is multiplied by the entire matrix  $B$ .

- **Positive Definite Matrix ( $S \succ 0$ ):**

$$\mathbf{x}^\top S \mathbf{x} > 0 \quad \forall \mathbf{x} \in \mathbb{R}^n, \mathbf{x} \neq \mathbf{0} \quad (12)$$

A matrix  $S$  is positive definite if for all non-zero vectors  $\mathbf{x}$ , the quadratic form  $\mathbf{x}^\top S \mathbf{x}$  is strictly greater than 0. Positive definite matrices are often used to indicate that a matrix represents a convex function.

**Dimensional Notation:**

$$A \in \mathbb{R}^{n \times m} : \mathbb{R}^n \rightarrow \mathbb{R}^m \quad (13)$$

$$B \in \mathbb{R}^{m \times k} : \mathbb{R}^m \rightarrow \mathbb{R}^k \quad (14)$$

$$AB \in \mathbb{R}^{n \times k} : \mathbb{R}^n \rightarrow \mathbb{R}^k \quad (15)$$



**Note 0.3.1** Matrix Multiplication as Dot Products of the Row and Column

This matrix multiplication is equivalent to the dot product between row  $i$  of matrix  $A$  and column  $j$  of matrix  $B$ :

$$\mathbf{x}^\top \mathbf{y} = \sum_{i=1}^n x_i y_i \quad (16)$$

Symbol	Meaning
$\mathbf{X}_{:,j}$	$j$ -th column of matrix
$\mathbf{X}_{i,:}$	$i$ -th row of matrix (treated as a column vector)
$\mathbf{X}_{ij}$	Element $(i, j)$ of matrix
$\mathbf{X}^\top$	Transpose of a matrix
$\text{diag}(\mathbf{S})$	Diagonal vector extracted from square matrix
$\mathbf{I}$ or $\mathbf{I}_n$	Identity matrix of size $n \times n$
$\mathbf{X} \odot \mathbf{Y}$	Hadamard/Elementwise product
$\mathbf{X} \otimes \mathbf{Z}$	Kronecker product
$\mathbf{S} \succ 0$	True iff $\mathbf{S}$ is a positive definite matrix
$\text{tr}(\mathbf{S})$	Trace of a square matrix
$\det(\mathbf{S})$	Determinant of a square matrix
$ \mathbf{S} $	Determinant of a square matrix

Table 3: Standard Notation for Matrices

## 0.4 Calculus

In this section, we cover basic multivariate calculus and key identities. These are considered prerequisite material, so students should spend time reviewing them if needed. For a more in-depth refresher, refer to Chapter 5 of Deisenroth et al. (2020).

Unqualified, a **function** is a mapping from the reals to the reals:  $f : \mathbb{R} \rightarrow \mathbb{R}$ . The mapping of a value  $x$  from the domain is denoted  $f(x)$ . The derivative is defined as:

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

assuming the limit exists. This is also known as the forward finite difference quotient.

**Leibniz** and **Lagrange** notations are often used for derivatives, as shown below:

Notation Type	Representation
Leibniz Notation	$\frac{dy}{dx}, \frac{d}{dx} f$
Lagrange Notation	$f'(x), f'(a)$ (for derivative at $a$ )

Table 4: Comparison of Leibniz and Lagrange Notations

The notation  $\left. \frac{dy}{dx} \right|_a$  is used for the derivative of  $f$  at  $a$ , while  $f'(a)$  is used in Lagrange notation.

## 0.5 Multivariate Calculus

For vector-valued functions  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , we define the partial derivative:

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

where  $e_i$  is the  $i$ -th basis vector. This extends the definition of the derivative to each variable.

The gradient  $\mathbf{g}$ , or  $\nabla f$ , is defined as the collection of all partial derivatives into a row vector:

$$\mathbf{g} = \nabla f = \left[ \frac{df}{dx_0}, \frac{df}{dx_1}, \dots, \frac{df}{dx_n} \right]$$

The gradient maps vectors to vectors:  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , which we call a **vector field**.

The derivative of a vector field  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is a matrix called the **Jacobian**, given by:

$$\nabla f = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

where  $f_i$  represents the  $i$ -th output of the function  $f$ . This matrix captures all partial derivatives and forms the foundation for calculating Jacobians of vector-valued functions.

## 0.6 Probability and Statistics

$x$  is a sample from a random variable  $X$ .

$P(X = x)$  refers to probability that the random variable  $X$  takes on the value  $x$ .

### 0.6.1 Probability Density Function (PDF)

The Probability Density Function (PDF) of a continuous random variable  $X$  is a function  $f_X(x)$  that describes the relative likelihood for this random variable to take on a given value. The PDF has the following properties:

- $f_X(x) \geq 0$  for all  $x$ .
- $\int_{-\infty}^{\infty} f_X(x) dx = 1$ .

Mathematically, the PDF is defined such that the probability that  $X$  lies within a particular interval  $[a, b]$  is given by:

$$P(a \leq X \leq b) = \int_a^b f_X(x) dx.$$

The most common example of a PDF is the normal distribution, which is characterised by two parameters: the mean  $\mu$  and the variance  $\sigma^2$ . The PDF of a normal distribution is given by:

$$\mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right).$$

### 0.6.2 Cumulative Distribution Function (CDF)

The Cumulative Distribution Function (CDF) of a continuous random variable  $X$  is a function  $F_X(x)$  that describes the probability that  $X$  will take a value less than or equal to  $x$ . It is defined as:

$$F_X(x) = P(X \leq x) = \int_{-\infty}^x f_X(t) dt.$$

The CDF has the following properties:

- $0 \leq F_X(x) \leq 1$  for all  $x$ .
- $F_X(x)$  is a non-decreasing function.
- $\lim_{x \rightarrow -\infty} F_X(x) = 0$ .
- $\lim_{x \rightarrow \infty} F_X(x) = 1$ .

Example PDF: Normal Distribution

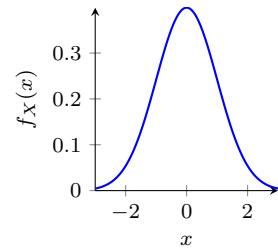


Figure 1: Probability Density Function of a standard normal distribution

Example CDF: Sigmoid Approximation

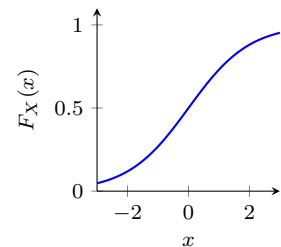


Figure 2: Cumulative Distribution Function

## 0.7 Moments

Moments help describe probability distributions. The **first moment** is the mean (or centre of mass), the **second moment** is the variance (or spread), and the **third moment** is the skewness. Higher moments are also useful but are less common in practice.

The **expectation** of a distribution, or its first moment, is given by:

$$\mu = \mathbb{E}[x] = \int_{-\infty}^{\infty} p(x=x)dx$$

where  $p(x)$  is the probability density or mass function.

The general form of the  $n$ -th moment is:

$$m_n = \frac{\mathbb{E}[(x-\mu)^n]}{\mathbb{E}[(x-\mu)^2]^{n/2}}$$

Moment	Name	Formula	Integration Formula	Description
1st	Mean	$\mu = \mathbb{E}[x]$	$\mu = \int_{-\infty}^{\infty} xp(x)dx$	Centre of mass, the average value of the random variable.
2nd	Variance	$\sigma^2 = \mathbb{E}[(x-\mu)^2]$	$\sigma^2 = \int_{-\infty}^{\infty} (x-\mu)^2 p(x)dx$	Spread of the distribution around the mean.
3rd	Skewness	$\frac{\mathbb{E}[(x-\mu)^3]}{\sigma^3}$	$\int_{-\infty}^{\infty} \frac{(x-\mu)^3}{\sigma^3} p(x)dx$	Asymmetry or tilt of the distribution.
4th	Kurtosis	$\frac{\mathbb{E}[(x-\mu)^4]}{\sigma^4}$	$\int_{-\infty}^{\infty} \frac{(x-\mu)^4}{\sigma^4} p(x)dx$	Measure of the tail thickness or sharpness of the peak.

## 0.8 Independent and Identically Distributed (i.i.d.) Random Variables

### 0.8.1 Independence of Random Variables

Independence refers to the idea that the values or outcomes of one observation in a dataset do not depend on or influence the values of any other observation.

For the set of random variables  $X_1, X_2, \dots, X_n$ , for the collection

$$\{x_1 \sim X_1, x_2 \sim X_2, \dots, x_N \sim X_N\}$$

we often assume independence, for any subset of observations  $\{X_{i_1}, X_{i_2}, \dots, X_{i_k}\}$  where  $i_1, i_2, \dots, i_k \in \{1, 2, \dots, N\}$ , the joint distribution factorises:

$$P(X_{i_1} = x_{i_1}, \dots, X_{i_k} = x_{i_k}) = P(X_{i_1} = x_{i_1}) \cdot \dots \cdot P(X_{i_k} = x_{i_k}) \quad (17)$$

The joint distribution of the subset is the product of the marginal distributions of the individual random variables.

**Why do we model samples as distinct random variables?**

1. By treating each sample as a distinct variable, we assume samples are i.i.d, allowing every sample to contribute independently to the likelihood of observing the data given the model parameters – so every sample provides unique information to estimate the parameters of the underlying distribution. If we treated all samples as a single random variable, we would lose the granularity of information, leading to a poorer estimate.
2. The assumption of i.i.d follows many results in probability and statistics, such as the Central Limit Theorem, which states that the sum of a large number of i.i.d. random variables is approximately normally distributed. It is also an assumed requirement for a model to generalise. The assumption simplifies the analysis and derivation process, allowing us to use techniques like maximum likelihood estimation (MLE) and empirical risk minimization (ERM).

#### Independence

#### Definition 0.8.1

The occurrence of any event does not affect the occurrence of others. This is a key assumption in many ML models that we will see the usefulness of in later chapters.

### 0.8.2 Identically Distributed Random Variables

The "identically distributed" part of i.i.d. refers to the idea that all random variables in the sample follow the same probability distribution. That is, they share the same probability density function (pdf) or probability mass function (pmf), depending on whether the data is continuous or discrete. If the set of random variables  $\{X_1, X_2, \dots, X_n\}$  are identically distributed, then:

$$f_{X_1}(x) = f_{X_2}(x) = \dots = f_{X_n}(x) = f_X(x) \quad (18)$$

$$F_{X_1}(x) = F_{X_2}(x) = \dots = F_{X_n}(x) = F_X(x) \quad (19)$$

This implies that the cumulative distribution function (CDF) is the same for all these random variables.

#### Identically Distributed

#### Definition 0.8.2

The random variables in a sample are drawn from the same probability distribution.

## 0.9 Einstein/Pythonic Index Notation

Most readers will have covered a lot of the linear algebra basics before, if so, this is the most important section to read!

Refer to this [video on Einstein summation convention](#) for more information.

As reasoning about matrix and vector products can sometimes be cumbersome, it is often useful to write out the operations we perform in index notation. The matrix product  $C = AB$  can be written as:  $C_{i,j} = \sum_k A_{i,k} B_{k,j}$ . It can also be useful to adopt a more "pythonic" index notation where we consider the system  $Ax = b$  and write the first entry of the vector  $b$  as:

$$A_{1,:}x = b_1 \quad (20)$$

which indicates that the first value of the result is simply the dot product of the first row of matrix  $A$  with the vector  $x$ .

Then, we can write the entire system as:

$$C_{i,j} = \sum A_{i,:} B_{:,j} \quad (21)$$

Generally in tensor calculus, a lot of expressions involve summing over particular indices.

$$\sum_{i=1}^3 a_i x_i = a_1 x_1 + a_2 x_2 + a_3 x_3 \quad (22)$$

**Einstein Notation consists of two different kinds of indices: the Dummy Index and the Free Index.** In **Einstein notation**, we can write this as:

$$a_i x_i \quad (i = 1, 2, 3) \quad (23)$$

This brings us to our first rule (and several others):

**Rule 1:** Any twice-repeated index in a single term is implicitly summed over.

Typically, this is from index 1 to 3 because most calculations are done in 3D space. For example, if we had:

$$a_{ij} b_j = a_{i1} b_1 + a_{i2} b_2 + a_{i3} b_3 \quad (24)$$

#### Fun fact

Yes, this was introduced by Albert Einstein in 1916!

*Non-Examinable 0.9.1*

We can then express this in the Einstein notation as:

$$a_{ij}b_j = a_{i\alpha}b_\alpha \quad \alpha \in \{1, 2, 3\} \quad (25)$$

This will begin to make more sense as we come up with a better way to describe indices that appear once, and twice-repeated indices.

**Rule 2:** The definitions of indices:

- We let  $j$  be the dummy index, because it is repeated only twice. One can thus replace  $j$  with any other index or letter, it is just a placeholder (thus called a dummy index). Although more rigorously:
  - One can replace any dummy index with a letter/index that is not already used in the expression.
  - This letter must be over the same range as the original dummy index, so in the case of replacing  $j$ , it must be over the range 1 to 3.

$$a_{ij}b_j = a_{i1}b_1 + a_{i2}b_2 + a_{i3}b_3 \quad (26)$$

$$= a_{i\alpha}b_\alpha \quad \alpha \in \{1, 2, 3\} \quad (27)$$

- $i$  is the free index, which can take on any value that  $j$  takes on, but it is not summed over and can only take on one value at a time  $i \in \{1, 2, 3\}$ .
- The free index occurs only **once** in the expression and **cannot be replaced by another free index**,

$$a_{ij}b_j \neq a_{kj}b_j \quad (28)$$

- To help avoid confusion, one tip is to use roman letters ( $i, j, k$ ) for free indices, and greek letters ( $\lambda, \mu, \rho$ ) for dummy indices.

**Rule 3:** No index may occur 3 or more times in a given term.

- $a_{ij}b_{ij}$  ✓
- $a_{ii}b_{ij}$  ✗
- $a_{ij}b_{ij}$  ✗
- $a_{ij}b_j + a_{ji}b_j$  ✓

In the last example, we are adding **multiple terms**, so the index occurrence rule only applies by term. So  $j$  is a dummy index for both terms since it occurs twice (per term).

**Rule 4:** In an equation involving Einstein notation, the free indices on the left-hand side must match the right-hand side.

- $x_i = a_{ij}b_j$  ✓
  - $i$  is a free index on both the LHS and RHS.
- $a_i = A_{xi}B_{xk}x_j + C_{ik}u_k$  ✓
  - $i$  is the free index on both the LHS and RHS.
- $x_i = A_{ij}$  ✗
  - $i$  is a free index on both the LHS and RHS, but  $j$  is a free index on the RHS that is not on the LHS.
- $x_j = A_{ik}u_k$  ✗
  - LHS free index:  $j$ .

- RHS free index:  $i$ .
- $x_i = A_{ik}u_k + c_j \quad \times$
- LHS free index:  $i$ .
- RHS free indices:  $i$  and  $j$ .

**Relating Einstein Notation to Pythonic Notation:** We had:

$$C_{i,j} = \sum_k A_{i,k} B_{k,j} \quad (29)$$

Since  $k$  is our dummy variable (it is summed over and doesn't appear in the final result), we can replace it with colon ( $:$ ) to indicate we are working with all elements along that dimension. So we have

$$C_{i,j} = A_{i,:} B_{:,j} \quad (30)$$

In Python, this is:

```
C[i, j] = np.dot(A[i, :], B[:, j])
```

# 1

## Formalising ML Problem Settings

### 1.1 *Machine Learning from the Mathematical Perspective*

What is machine learning?

Professor

Machine learning is the study of algorithms which automatically perform a task by processing an example dataset.

ChatGPT

Machine learning is a computational approach that allows systems to automatically learn and improve from experience without being explicitly programmed.

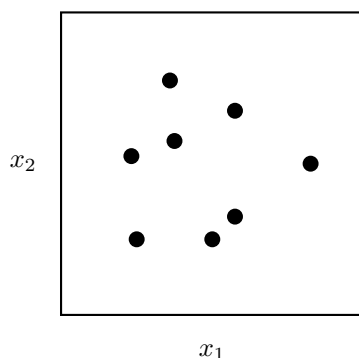
The professor himself much prefers his definition, which focuses on the word **algorithms**, **task**, and **dataset**.

The issue with the GPT definition is that the word ‘computational’ is too broad and vague, and the notion of ‘learn’ and ‘improving from experience’ is too colloquially used in the actual nature of machine learning.

### 1.2 *Datasets*

Data is the keystone of machine learning, used to infer and learn good models. We denote the input or feature space as  $x \in \mathbb{R}^n$  and outputs as  $y \in \mathbb{R}^m$ .

There are several ways we can look at data:



#### 1.2.1 *Set Perspective*

- Data is viewed as a collection of  $N$  2-tuples:  $\{(x^{(i)}, y^{(i)})\}_{i=1}^N$ , where each input-output pair is part of the dataset.
- This perspective is order/permutation invariant, focusing on the set of samples rather than the order in which they appear.
- The size of the set measures the sample complexity.

- Sets are often useful for reasoning about feature/label spaces in an abstract way.
- Example of dataset representation:

$$\{(\mathbf{x}^{(i)})_{i=1}^N\}$$

This represents a dataset consisting of a collection of  $N$  input vectors  $\mathbf{x}^{(i)}$ . Each  $\mathbf{x}^{(i)}$  is a feature vector or a sample, without any corresponding labels or outputs. This kind of dataset is typically used in unsupervised learning, where we only have inputs and are not concerned with outputs (e.g., clustering or dimensionality reduction).

$$\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})_{i=1}^N\}$$

This represents a set of  $N$  input-output pairs, where  $\mathbf{x}^{(i)}$  is the input feature vector and  $\mathbf{y}^{(i)}$  is the corresponding label or output. This is the standard representation for datasets in supervised learning, where the goal is to learn a function that maps inputs to outputs (e.g., regression or classification tasks).

$$\{(\mathbf{s}^{(i)}, \mathbf{a}^{(i)}, \mathbf{s}^{(i+1)})_{i=1}^N\}$$

This set represents a sequence of  $N$  tuples, where  $\mathbf{s}^{(i)}$  is the state at time step  $i$ ,  $\mathbf{a}^{(i)}$  is the action taken in that state, and  $\mathbf{s}^{(i+1)}$  is the resulting next state. This kind of dataset is commonly used in reinforcement learning (RL), where the goal is to learn a policy that dictates which actions to take in various states to maximize a reward signal over time.

### 1.2.2 Probability Perspective

- Data is viewed as samples where each sample is a 2-tuple  $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$  drawn from a joint distribution  $P(\mathbf{X}, \mathbf{Y})$ , where the relationship between input  $\mathbf{X}$  and output  $\mathbf{Y}$  is of primary interest.
- We are interested in the probability of an output given we have observed a particular input:  $P(\mathbf{Y} = \mathbf{y} | \mathbf{X} = \mathbf{x})$ .
- This perspective is fundamental in probabilistic machine learning and provides tools for modelling uncertainty and reasoning about observations.
- Common assumptions like i.i.d. (independent and identically distributed) are often applied:

$$P(x_1, x_2, \dots, x_k) = P(x_1) \cdot P(x_2) \dots P(x_k)$$

- Example of dataset representation with probabilistic modeling:

$$\mathbf{x} \sim P(\mathbf{X})$$

This represents an input vector  $\mathbf{x}$  sampled from a probability distribution  $P(\mathbf{X})$ . In this case,  $\mathbf{X}$  refers to the input space, and each input  $\mathbf{x}$  is drawn according to its distribution.

$$(\mathbf{x}, \mathbf{y}) \sim P(\mathbf{X}, \mathbf{Y})$$

This represents an input-output pair  $(\mathbf{x}, \mathbf{y})$  sampled from the joint distribution  $P(\mathbf{X}, \mathbf{Y})$ , where  $\mathbf{X}$  refers to the input space, and  $\mathbf{Y}$  refers to the output space. The goal here is to model the relationship between inputs  $\mathbf{x}$  and corresponding outputs  $\mathbf{y}$ , as is typical in supervised learning.



$$\mathbf{s}^{(i+1)} \sim P(\mathbf{s}^{(i+1)} \mid \mathbf{s}^{(i)}, \mathbf{a}^{(i)})$$

This represents the next state  $\mathbf{s}^{(i+1)}$  sampled from the conditional probability distribution given the current state  $\mathbf{s}^{(i)}$  and action  $\mathbf{a}^{(i)}$ . This kind of representation is typical in reinforcement learning (RL), where we model transitions between states based on actions taken by the agent.

### 1.2.3 Linear Algebra/Statistics Perspective

- A bit misleading to call it a ‘linear algebra perspective’ since it’s used here in almost all of statistics.
- Data is treated as vectors and matrices. An input is a vector  $\mathbf{x} \in \mathbb{R}^{n \times 1}$ , collected into a matrix  $X \in \mathbb{R}^{m \times n}$ , with  $m$  samples and  $n$  features.

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad X = \begin{bmatrix} \mathbf{x}^{(1)\top} \\ \mathbf{x}^{(2)\top} \\ \vdots \\ \mathbf{x}^{(m)\top} \end{bmatrix} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix}$$

- Each row of the matrix represents a data sample (input vector), and columns correspond to features (called explanatory variables in statistics). The matrix is often referred to as the *design matrix*.
- This perspective allows the application of linear algebra tools to study properties like rank, invertibility, and dimensionality reduction.
- The labels (response variables)  $y \in \mathbb{R}^{m \times 1}$  are collected into a vector:

$$\mathbf{x} = \begin{pmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}$$

## 1.3 Types of Learning Tasks

There are several types of learning tasks that can be broadly classified into the following categories:

### 1.3.1 Supervised Learning

In supervised learning, we train models using input-output pairs, where the correct output (label) is known. The goal is to learn a mapping from inputs to outputs. For example, given a dataset of customer transaction data and their purchase categories, the task is to predict future purchases based on transaction history.

There are also two types of supervised learning tasks:

1. **Regression:** The output variable is continuous, such as predicting house prices based on square footage.
2. **Classification:** The output variable is discrete (categorical), such as classifying emails as spam or not spam.

### 1.3.2 Unsupervised Learning

Unsupervised learning involves learning from data without labeled outputs. The objective is to discover hidden patterns or groupings in the data. An example is clustering news articles into different topics based on the text content, without prior knowledge of the categories.

Unsupervised learning comes in convenient if we have huge amounts of data without labels, and we want to extract meaningful insights from it without having to manually label each data point (which is almost often done manually and can be expensive).

### 1.3.3 Generative Learning

Generative learning aims to model the underlying data distribution and generate new data samples that resemble the training data. For instance, a generative model could be trained on music tracks and then used to generate new compositions that sound like the originals. This can be in a supervised or unsupervised setting.

### 1.3.4 Algorithms

Back to the definition of machine learning:

Professor

Machine learning is the study of algorithms which automatically perform a task by processing an example dataset.

**Algorithms** are not the same as **models**. Instead, models are assumed/inferred and algorithms returned learned models.

Without getting into many technicalities, it is difficult to convey meaningful algorithmic difference, thus this section will be answered throughout the course.

## 1.4 The Linear Model

Choosing a model in machine learning is deciding how we think inputs and outputs can be related. The most basic linear model assumes that the output is a weighted linear combination of the input features. Weights are represented by vector  $\mathbf{w}$ , and the model is represented as:

$$\hat{y} = \mathbf{w}^\top \mathbf{x} := \sum_{i=1}^n w_i x_i \quad (1.1)$$

The above model would be learned by an **algorithm** from the dataset, where the algorithm decides what the weights should be. The model is then used to make predictions on new data points.

This is a dot product between the weight vector  $\mathbf{w}$  and the input feature vector  $\mathbf{x}$ . Keen readers will note out that this only model lines that cross through the origin, however there will be methods to extend this model to lines that do not cross through the origin.

One such way would be to normalize the data so it can always be represented by a line through the origin.

## 1.5 Data Preprocessing: Normalising the Data

Given the limitations of the linear model in the previous section, data normalisation allows us to transform the data so that it can be modelled by a line through the origin.

### 1.5.1 Naïve Method

Assuming no noise and outlier, normalisation can be naively done (s) by the following steps:

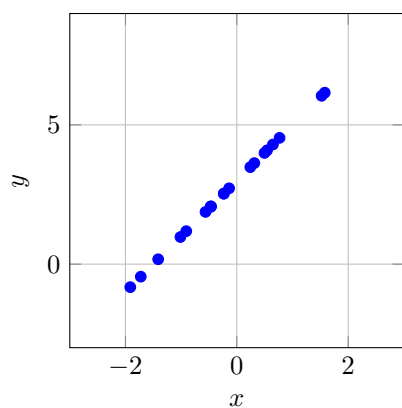
- Pick an arbitrary feature-label pair in the dataset,  $(\mathbf{x}^{(i)}, y^{(i)})$ .
- Subtract  $\mathbf{x}^{(i)}$  from all input features and  $y^{(i)}$  from all target labels.
- The data has been transformed.

#### Example for 1D $x$ and 1D $y$

#### Intuition 1.5.1

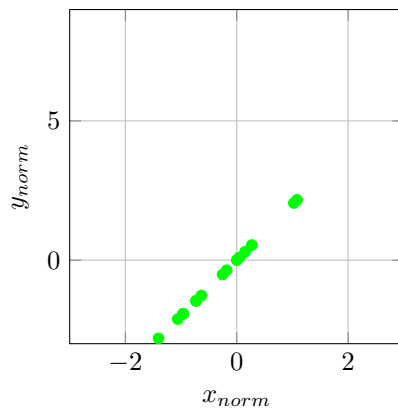
- The points can be described by the relationship  $y = mx + b$  where  $b$  is the bias.
- The bias is preventing us from modelling the data with the relationship  $y = mx$ .
- We choose an arbitrary point in the dataset to be subtracted from all points
- We can choose any point because all points contain the bias component.
- After subtracting, all bias components are removed, and the data can be modelled by a line through the origin.

Before Normalisation (Perfect Points)



(a) Before Normalisation (Perfect Points)

After Naive Normalisation (Perfect Points)



(b) After Naive Normalisation

Figure 1.1: Comparison of before and after naive normalisation.

We almost never use this method, as it is not robust to noise and outliers. We will discuss a more robust method in the next section.

### 1.5.2 Z-Score Normalisation

There is always be noise and outliers in the data. So instead of picking an arbitrary point to subtract from all points, we can use the mean of the input features and target values and subtract them from all points.

We also divide by the standard deviation to scale the data. This is called **standardisation** or **z-score normalisation**.

We begin by calculating the mean of the input features and the target values:

$$\mathbf{x}^\mu = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)} \quad y^\mu = \frac{1}{N} \sum_{i=1}^N y^{(i)} \quad (1.2)$$

Next, we compute the standard deviation for both the input features and target values:

$$\mathbf{x}^\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (\mathbf{x}^{(i)} - \mathbf{x}^\mu)^2} \quad y^\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (y^{(i)} - y^\mu)^2} \quad (1.3)$$

Once we have computed the mean and standard deviation, we can normalise the data by subtracting the mean and dividing by the standard deviation. This centres the data around the origin and scales it:

$$\tilde{\mathbf{x}}^{(i)} = \frac{\mathbf{x}^{(i)} - \mathbf{x}^\mu}{\mathbf{x}^\sigma} \quad \tilde{y}^{(i)} = \frac{y^{(i)} - y^\mu}{y^\sigma} \quad (1.4)$$

The normalised data,  $\tilde{\mathbf{x}}^{(i)}$  and  $\tilde{y}^{(i)}$ , can now be used in the linear model, which will allow it to handle cases where the data does not necessarily pass through the origin.

By centring and scaling the data, the linear model becomes more robust, especially when dealing with non-centred data. This is because the impact of outliers are reduced when subtracted by the mean of the points and scaled down by the standard deviation.

#### Note 1.5.1 Division by standard deviation

The curriculum notes only mention subtraction of the mean, and do not include the division by the standard deviation. Subtraction of the mean has a centring effect and an absolute reduction of outliers, however it is mostly the scaling by the standard deviation that reduces the (relative) impact of outliers.

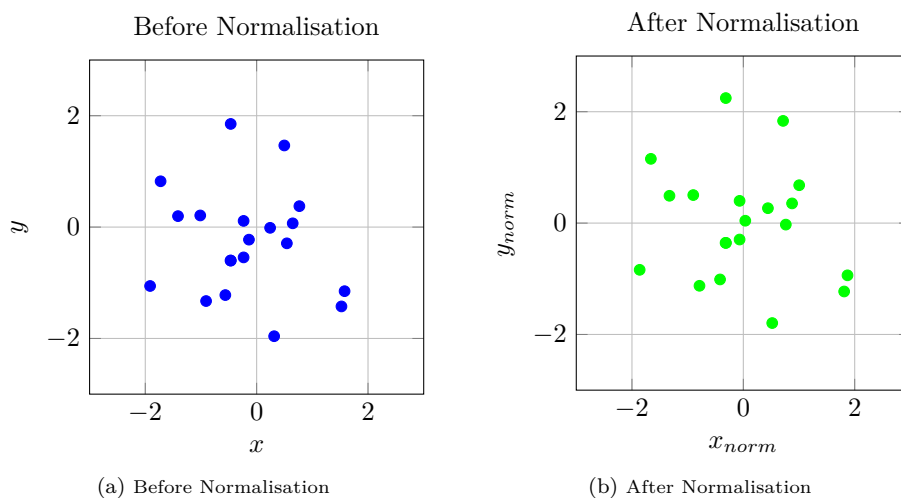


Figure 1.2: Example of a 1-dimension feature  $x$  and 1-dimensional output  $y$ . Comparison of points before and after normalisation.

## 2

# Linear Regression

## 2.1 Learning as Optimisation

Linear regression is exactly as its name suggests: performing regression with a linear model.

You may recall linear models written as affine transformations:  $y = mx + b$ , where  $b$  is the bias or constant term – this makes the model affine and not linear. Linear models refers to the relationship between model parameters and predictions via a linear transformation. This will be covered in later sections.

### Linear Transformation

### Definition 2.1.1

A linear transformation between two vector spaces  $V$  and  $W$  is a map

$$T : V \rightarrow W$$

such that:

- $T(v_1 + v_2) = T(v_1) + T(v_2) \quad \forall v_1, v_2 \in V$
- $T(\alpha v) = \alpha T(v) \quad \forall v \in V \text{ and scalar } \alpha$

Assuming a domain of  $\mathbb{R}^n$  and a one-dimensional co-domain, we can write our model as  $f(\mathbf{x}) = \mathbf{x}^\top \theta$ . Thus we have:

$$\hat{y}^{(i)} = \mathbf{x}^{(i)\top} \theta$$

The goal of learning is to find  $\theta$  such that  $\hat{y}^{(i)} \approx y^{(i)}$ . It is unclear what  $\approx$  means

in this context. This is because if we are predicting something like yearly salary, a prediction of £70,000 when the actual value is £70,300 is a decent prediction despite being £300 off. If you predicted age as 328 when the actual value is 28, that's way off. This is where standardisation comes in.

### Standardisation

### Definition 2.1.2

A form of data preprocessing to standardise response variables, allowing us to compare the performance of a model/algorithm without having to reason about with data-specific average error values.

An example of standardisation in our case would be to standardise response variables to be in the range  $[-1, 1]$ .

### 2.1.1 Norms as Error Models

After standardising, one way to measure a model's predictive error is to use **norms**. A **norm** is a function that assigns a non-negative length or size to a vector. Norms are widely used in mathematics to measure distances or lengths in vector spaces.

- **Euclidean Norm (2-norm)**: Also known as the straight-line distance, the Euclidean norm is given by:

$$\|\mathbf{x}\|_2 = \sqrt{\sum_i x_i^2} \quad (2.1)$$

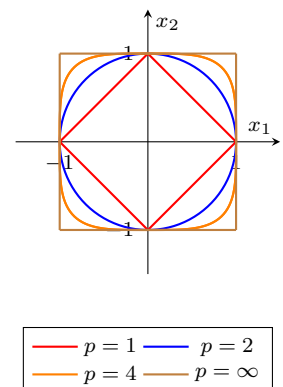


Figure 2.1: Different norms in 2D space.

- **p-norm:** The  $p$ -norm generalises the Euclidean norm to other values of  $p$ :

$$\|\mathbf{x}\|_p = \left( \sum_i |x_i|^p \right)^{1/p} \quad (2.2)$$

Special cases of the  $p$ -norm include:

- The 2-norm (Euclidean norm) when  $p = 2$ .
- The 1-norm, which represents the Manhattan distance, when  $p = 1$ .
- **Zero Norm (Hamming Distance):** The 0-norm counts the number of non-zero elements in a vector:

$$\|\mathbf{x}\|_0 = \sum_i \mathbb{I}(x_i \neq 0) \quad (2.3)$$

where  $\mathbb{I}(\cdot)$  is the indicator function, which returns 1 if the condition is true and 0 otherwise.

- **Infinity Norm:** Also called the *supremum norm*, the infinity norm is defined as the largest absolute value among the vector components:

$$\|\mathbf{x}\|_\infty = \sup_n |x_n| \quad (2.4)$$

- **Mahalanobis Distance:** The Mahalanobis distance takes into account the correlations between variables in a dataset:

$$d_S(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})^T S (\mathbf{x} - \mathbf{y}) \quad (2.5)$$

where  $S$  is a positive semi-definite matrix, typically the covariance matrix of the data.

Norms are typically used as a way to measure distances and magnitudes.

### Theorem 2.1.1 Geometric Interpretation of the Mahalanobis Distance

Given a positive semi-definite matrix  $S \in \mathbb{R}^{m \times m}$ , a feature vector  $\mathbf{x}' \in \mathbb{R}^m$ , and a similarity threshold  $\delta$ , all vectors  $\mathbf{x}'' \in \mathbb{R}^n$  satisfying  $d_S(\mathbf{x}', \mathbf{x}'') \leq \delta$  are contained within the axis-aligned orthotope:

$$[\mathbf{x}' - \delta\sqrt{\mathbf{d}}, \mathbf{x}' + \delta\sqrt{\mathbf{d}}]$$

where  $\mathbf{d} = \text{diag}(S)$ , the vector containing the elements along the diagonal of  $S$ .

## 2.1.2 Equivalence of Inner Product and Euclidean Norm

Taken from previous notes – notation is different!

We aim to prove the equivalence between the quadratic form and the squared Euclidean norm:

$$\text{Claim: } (\theta^T \mathbf{X} - \mathbf{y})^T (\theta^T \mathbf{X} - \mathbf{y}) = \|\theta^T \mathbf{X} - \mathbf{y}\|_2^2 \quad (2.6)$$

**Left-hand side (LHS):** The quadratic form on the LHS can be expanded as:

$$(\theta^T \mathbf{X} - \mathbf{y})^T (\theta^T \mathbf{X} - \mathbf{y}) = \left( \left( \sum_i \theta_i^T \mathbf{X}_{i,j} \right) - \mathbf{y}_j \right)^T \left( \left( \sum_i \theta_i^T \mathbf{X}_{i,j} \right) - \mathbf{y}_j \right) \quad (2.7)$$

$$= \sum_j \left( \left( \sum_i \theta_i^T \mathbf{X}_{i,j} \right) - \mathbf{y}_j \right)^T \left( \left( \sum_i \theta_i^T \mathbf{X}_{i,j} \right) - \mathbf{y}_j \right) \quad (2.8)$$

$$= \sum_j \left( \left( \sum_i \theta_i^T \mathbf{X}_{i,j} \right) - \mathbf{y}_j \right)^2 \quad (2.9)$$

**Right-hand side (RHS):** The squared Euclidean norm is defined as:

$$\|\theta^T \mathbf{X} - \mathbf{y}\|_2^2 \quad (2.10)$$

We know that the Euclidean norm is given by:

$$\|\mathbf{x}\|_2 = \sqrt{\sum_i x_i^2} \quad (2.11)$$

Therefore, we expand into:

$$\|\theta^T \mathbf{X} - \mathbf{y}\|_2^2 = \left( \sqrt{\sum_j \left( \sum_i \theta_i^T \mathbf{X}_{i,j} - \mathbf{y}_j \right)^2} \right)^2 \quad (2.12)$$

$$= \sum_j \left( \sum_i \theta_i^T \mathbf{X}_{i,j} - \mathbf{y}_j \right)^2 \quad (2.13)$$

**Conclusion:** Both the LHS and RHS are equivalent, as they both expand into the same form:

$$\sum_j \left( \sum_i \theta_i^T \mathbf{X}_{i,j} - \mathbf{y}_j \right)^2 \quad (2.14)$$

This shows that the inner product expression for  $(\theta^T \mathbf{X} - \mathbf{y})^T (\theta^T \mathbf{X} - \mathbf{y})$  is equivalent to the squared Euclidean norm  $\|\theta^T \mathbf{X} - \mathbf{y}\|_2^2$ , confirming the claim.

To finish it off in index notation, this is written as:

$$\theta_i \mathbf{X}_{i,j} \quad (2.15)$$

### 2.1.3 The best model is determined by the argmax parameter

We collect our model predictions into vector  $\hat{\mathbf{y}} := \theta^T \mathbf{x}^{(i)}$  and collect true values into corresponding vector  $\mathbf{y}$ . We then define the **error function** to compare the difference between predictions and true values.

We tend to use the Root Mean Squared Error (RMSE) as our error function, defined as:

$$\mathcal{L}_{RMSE}(\hat{\mathbf{y}}, \mathbf{y}) = \sqrt{\frac{1}{N} \sum_{i=1}^N (\hat{y}^{(i)} - y^{(i)})^2} \quad (2.16)$$

Assuming we've already standardised, the RMSE is defined as:

$$\mathcal{L}_{RMSE}(\hat{\mathbf{y}}, \mathbf{y}) = \|\hat{\mathbf{y}} - \mathbf{y}\|_2 \quad (2.17)$$

We also have the Mean Squared Error (MSE), defined as:

$$\mathcal{L}_{MSE}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{N} \sum_{i=1}^N (\hat{y}^{(i)} - y^{(i)})^2 = \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2 \quad (2.18)$$

We can rewrite  $\hat{\mathbf{y}}$  in terms of parameters as  $\mathbf{X}$  and  $\mathbf{y}$  are fixed, and express the above errors as functions of the parameter value, which we will call the **loss function**.

$$\mathcal{L}(\theta) = \|\mathbf{X}\theta - \mathbf{y}\|_2^2 \quad (2.19)$$

We use **loss functions** to determine if a parameter is better than another. Given two parameters,  $\theta'$  and  $\theta''$ , if  $\mathcal{L}(\theta') < \mathcal{L}(\theta'')$ , then  $\theta'$  is a better parameter than  $\theta''$ .

We then want to find the best parameter,  $\theta^*$ , which has the lowest loss from  $\mathcal{L}(\theta^*) = \min_{\theta} \mathcal{L}(\theta)$ . This is when the loss function is at its minimum. We don't care about *what* the minimum is, rather *where* it is. This is why we use the **argmin** to identify the argument where the function returns its minimum.

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^n} \mathcal{L}(\theta) \quad (2.20)$$

## 2.2 Ordinary Least Squares (OLS)

We want to find the parameter  $\theta^*$  that minimises the loss. For the linear model, the Mean Squared Error (MSE) can be written as:

$$L(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{k} \sum_i (\mathbf{y}_i - \mathbf{X}_{i,j} \theta_j)^2$$

This can also be expressed as:

$$L(\theta) = \frac{1}{k} (\mathbf{X}\theta - \mathbf{y})^\top (\mathbf{X}\theta - \mathbf{y})$$

We know that minimising functions occurs where the gradient is zero. To simplify differentiation, we expand:

$$\frac{1}{k} (\mathbf{X}\theta - \mathbf{y})^\top (\mathbf{X}\theta - \mathbf{y}) = \theta^\top (\mathbf{X}^\top \mathbf{X}) \theta - 2\mathbf{y}^\top \mathbf{X} \theta + \mathbf{y}^\top \mathbf{y}$$

Instead of differentiating manually, we can prove several differentiation rules to find out the derivative of the loss function with respect to the model parameters  $\theta$ .

### 2.2.1 Gradient Calculation

We calculate  $\nabla_{\theta} L = 0$ , using vector calculus identities:

- $\nabla_{\theta}(c^\top \theta) = c^\top$
- $\nabla_{\theta}(\theta^\top A \theta) = \theta^\top (A + A^\top)$

With these rules, we have:



$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{k} \sum_i (\mathbf{y}_i - \hat{\mathbf{y}}_i)^2 \quad (2.21)$$

$$= \frac{1}{k} \sum_i (\mathbf{y}_i - \mathbf{X}_{i,j} \theta_j)^2 \quad (2.22)$$

$$= \frac{1}{k} ((\mathbf{X}\theta - \mathbf{y})^\top (\mathbf{X}\theta - \mathbf{y})) \quad (2.23)$$

$$= \frac{1}{k} (\theta^\top \mathbf{X}^\top \mathbf{X} \theta - 2(\mathbf{y}^\top \mathbf{X} \theta) + \mathbf{y}^\top \mathbf{y}) \quad (2.24)$$

$$= \frac{1}{k} (\theta^\top [\mathbf{X}^\top \mathbf{X}] \theta - 2([\mathbf{y}^\top \mathbf{X}] \theta) + \mathbf{y}^\top \mathbf{y}) \quad (2.25)$$

$$\nabla_\theta \mathcal{L} = \frac{2}{k} (\mathbf{X}^\top \mathbf{X} \theta - \mathbf{y}^\top \mathbf{X}) \quad (2.26)$$

Setting this equal to zero and solving:

$$\begin{aligned} \theta^\top (\mathbf{X}^\top \mathbf{X}) &= \mathbf{y}^\top \mathbf{X} \\ \theta &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \end{aligned}$$

This is the **ordinary least squares estimate**.

#### General Steps for Checking Vector Calculus

#### Definition 2.2.1

Three general step to check vector calculus identities:

1. Break down the expression into index notation (Einstein summation notation).
2. Differentiate using standard calculus
3. Reassemble from index notation

#### Proving $\nabla_\theta \mathbf{c}^\top \theta = \mathbf{c}^\top$

#### Definition 2.2.2

Consider the expression  $\mathbf{c}^\top \theta$ , which in Einstein summation notation expands as:

$$\mathbf{c}^\top \theta = \sum_j c_j \theta_j$$

Now, taking the partial derivative with respect to  $\theta_j$  (this allows us to see what happens when not considering the dummy variable  $j$ ):

$$\frac{\partial \mathbf{c}^\top \theta}{\partial \theta_j} = c_j$$

This shows that the gradient of  $\mathbf{c}^\top \theta$  with respect to  $\theta$  is:

$$\nabla_\theta \mathbf{c}^\top \theta = \mathbf{c}^\top$$

**Proving  $\nabla_{\theta}(\theta^{\top} \mathbf{A} \theta) = \theta^{\top}(\mathbf{A} + \mathbf{A}^{\top})$** **Definition 2.2.3**

Now consider the quadratic form  $\theta^{\top} \mathbf{A} \theta$ . Expanding it in Einstein notation:

$$\theta^{\top} \mathbf{A} \theta = \sum_i \sum_j \theta_i \theta_j \mathbf{A}_{i,j}$$

Taking the derivative with respect to  $\theta_k$  (this allows us to see what happens when not considering the dummy variable  $k$ ):

$$\begin{aligned} \frac{\partial \theta^{\top} \mathbf{A} \theta}{\partial \theta_k} &= \sum_i \theta_i \mathbf{A}_{i,k} + \sum_j \mathbf{A}_{k,j} \theta_j \\ &= (\mathbf{A} \theta)_k + (\mathbf{A}^{\top} \theta)_k \end{aligned}$$

Thus, the gradient of  $\theta^{\top} \mathbf{A} \theta$  is:

$$\nabla_{\theta}(\theta^{\top} \mathbf{A} \theta) = \theta^{\top}(\mathbf{A} + \mathbf{A}^{\top})$$

**Note 2.2.1 The Derivative of Vectors are Row Vectors**

Some engineering and physics students may find this jarring, but the conventions of the notes will consider the derivative of vector-valued functions to always be row vectors.

$$\nabla_{\theta} \mathbf{c}^{\top} \theta = \mathbf{c}^{\top}$$

The reason why we use row vectors in the convention is because operations such as the chain rule and matrix-vector multiplications maintain a clean form.

Consider  $\mathbf{y} = \mathbf{A} \mathbf{x}$ , where  $\mathbf{A}$  is a matrix and  $\mathbf{x}$  is a vector. Let  $x = f(\mathbf{y})$  where  $f(\mathbf{y})$  is a scalar-valued function of  $\mathbf{y}$ .

If we want to compute  $\frac{\partial z}{\partial \mathbf{x}}$ :

**Column Vector Convention**

Since the gradient  $\nabla_{\mathbf{x}} z$  is a column vector, when applying the chain rule, we have to transpose the gradient to maintain compatibility with the matrix-vector multiplication:

$$\begin{aligned} \frac{\partial z}{\partial \mathbf{y}} &= \nabla_{\mathbf{y}} z \quad \text{this is a column vector} \\ \frac{\partial \mathbf{y}}{\partial \mathbf{x}} &= \mathbf{A} \quad \text{this is a matrix} \\ \frac{\partial z}{\partial \mathbf{x}} &= \frac{\partial z}{\partial \mathbf{y}} \cdot \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = (\nabla_{\mathbf{y}} z)^{\top} \cdot \mathbf{A} \end{aligned}$$

**Row Vector Convention**

Here, there's no need to transpose the result. The row vector  $\nabla_{\mathbf{y}} z$  directly multiplies the matrix  $\mathbf{A}$ :

$$\begin{aligned} \frac{\partial z}{\partial \mathbf{y}} &= \nabla_{\mathbf{y}} z \quad \text{this is a row vector} \\ \frac{\partial \mathbf{y}}{\partial \mathbf{x}} &= \mathbf{A} \quad \text{this is a matrix} \\ \frac{\partial z}{\partial \mathbf{x}} &= \frac{\partial z}{\partial \mathbf{y}} \cdot \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \nabla_{\mathbf{y}} z \cdot \mathbf{A} \end{aligned}$$

## 2.3 Basis Expansion

We recap on the linear model, which is described as a linear transformation:

Linear Transformation	Definition 2.3.1
A linear transformation between two vector spaces $V$ and $W$ is a map	
$T : V \rightarrow W$	
such that:	
<ul style="list-style-type: none"> <li>• <math>T(v_1 + v_2) = T(v_1) + T(v_2) \quad \forall v_1, v_2 \in V</math></li> <li>• <math>T(\alpha v) = \alpha T(v) \quad \forall v \in V \text{ and scalar } \alpha</math></li> </ul>	

The model of  $y = mx + b$  allows us to fit lines that don't pass through the origin without normalisation. However this does not fit the rules of a linear transformation, because the model becomes affine.

You may recall linear models written as affine transformations:  $y = mx + b$ , where  $b$  is the bias or constant term – this makes the model affine and not linear. Linear models refers to the relationship between model parameters and predictions via a linear transformation.

### 2.3.1 Affine Basis Expansion

The aforementioned limitation of linear models is easily fixed with a straightforward modification to capture the affine transformation  $f(x) = mx + b$ : we just add a feature to the input vector  $x$  that is always equal to 1, then the corresponding weight for this feature becomes the bias. Introduce:

$$\phi(x) : \mathbb{R} \Rightarrow \mathbb{R}^2 \quad \text{such that } \phi(x) = \begin{bmatrix} 1 \\ x \end{bmatrix} \quad (2.27)$$

We then need parameter vector  $\theta = \begin{bmatrix} b \\ m \end{bmatrix}$ .

We then have the model

$$\hat{y} = \phi(x)^\top \theta \Rightarrow \begin{bmatrix} 1 & x \end{bmatrix} \begin{bmatrix} b \\ m \end{bmatrix} = b + mx \quad (2.28)$$

Thus  $\theta_2$  is effectively the 'y-intercept'. It is the most simplistic form of basis expansion.

### 2.3.2 Polynomial Basis Expansion

Polynomial basis expansion is used to model non-linear functions. The key idea of basis expansion is to expand a one-dimensional feature into many dimensions, and use non-linear functions to increase the expressiveness of the model.

**Example 2.3.1** (Example Polynomial Basis Expansion)

A one dimensional domain  $x \in \mathbb{R}$  and a one dimensional co-domain  $y \in \mathbb{R}$  is assumed. Our model is  $\hat{y} = \phi(x)^\top \theta$ .

We choose the basis:

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix} \quad \phi(x) : \mathbb{R}^1 \rightarrow \mathbb{R}^3$$

and weights:

$$\theta \in \mathbb{R}^3$$

We finally have the fully expanded function (we use the dot product instead of transposing, for clarity of notation):

$$\hat{y} = \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix} \cdot \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} = \theta_0 + \theta_1 x + \theta_2 x^2$$

This example model is a quadratic polynomial.

**Example 2.3.2** (Another Example Polynomial Basis Expansion)

Again, given our model is  $\hat{y} = \phi(\mathbf{x})^\top \theta$  where  $y \in \mathbb{R}$  and  $\mathbf{x} \in \mathbb{R}^2$ . We can use the basis

$$\phi(\mathbf{x}) \Rightarrow \phi \left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \end{bmatrix} \quad \phi(x) : \mathbb{R}^2 \rightarrow \mathbb{R}^6$$

and with a new set of corresponding weights, we have:

$$\hat{y} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \end{bmatrix} \cdot \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \\ \theta_5 \end{bmatrix} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 \quad \theta \in \mathbb{R}^6$$

**2.3.3 Radial Basis Function Kernel**

The polynomial basis expansion is just a single flavour of basis expansion. Another widely-used form of basis is the kernel basis expansion. One popular example is the radial basis function kernel (RBF kernel), which is a generalisation of the polynomial basis expansion.

It takes in a fixed parameter  $\gamma > 0$ , defined as

$$\kappa(x, x') = \exp(-\gamma \|x - x'\|^2)$$

where  $\|x - x'\|^2$  is the squared Euclidean distance (or more appropriately, the  $L_2$  Norm) between  $x$  and  $x'$ . In practice, one picks fixed centres  $x$  and the basis expansion computes the expanded feature set w.r.t the distance to these centres.<sup>1</sup>

<sup>1</sup> There is more nuance to this: you may have noticed that kernel functions take in two points, unlike the polynomial basis expansion  $\phi$  taking in a single point. This is because for each of the  $n$  points, we compute pairwise similarities with a point's other points, and then construct a feature vector for each point, where each point  $x$  is transformed into a  $n$ -dimensional feature vector where each dimension represents the similarity between  $x$  and one of the  $n$  points.

For example, given a dataset with three points  $x_1, x_2, x_3$ , and a radial basis function  $\kappa(x, x')$ , the RBF kernel matrix might look like this:

$$K = \begin{bmatrix} \kappa(x_1, x_1) & \kappa(x_1, x_2) & \kappa(x_1, x_3) \\ \kappa(x_2, x_1) & \kappa(x_2, x_2) & \kappa(x_2, x_3) \\ \kappa(x_3, x_1) & \kappa(x_3, x_2) & \kappa(x_3, x_3) \end{bmatrix}$$

where the feature vector for  $x_1$  as an example would be

$$\begin{bmatrix} \kappa(x_1, x_1) & \kappa(x_1, x_2) & \kappa(x_1, x_3) \end{bmatrix}$$

This is related to the use of the **kernel trick**.

- It is easy to see that the kernel basis expansion  $\kappa(x, x')$  has a minimum value of 0. It takes a maximum value of 1 when  $x = x'$ .
- When two points  $x$  and  $x'$  are far apart, the kernel value is closer to 0, and when they are closer together, the kernel value is closer to 1.
- It is quite akin to a similarity score, and that a smaller value of  $\gamma$  leads to larger similarity scores (see Figure 2.2).
- It is also symmetric,  $\kappa(x, x') = \kappa(x', x)$ , and is always positive.

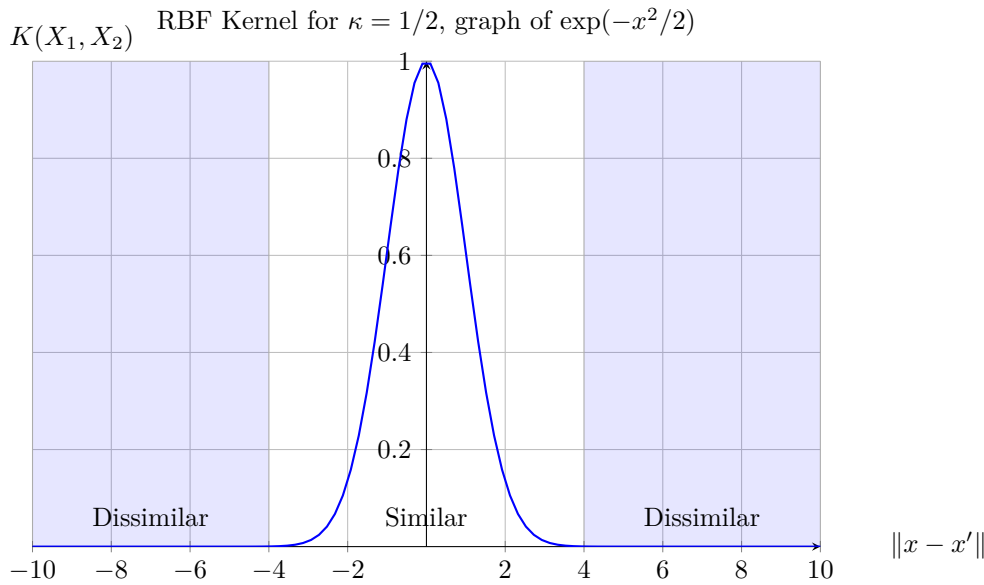


Figure 2.2: RBF Kernel for  $\kappa = 1/2$ , graph of  $\exp(-x^2/2)$ . Areas of dissimilarity are subjectively noted where the kernel values are negligible.

The RBF kernel is actually a special case of the polynomial basis expansion, where it is an infinite sum over a set of polynomial kernels. Proof [here](#).

#### Note 2.3.1 Syllabus Coverage

The full extent of kernel basis expansion is not covered in this section.

#### Example Radial Basis Function Kernel

#### Non-Examinable 2.3.1

Given a one-dimensional domain and a one-dimensional co-domain, we have the model  $\hat{y} = \phi(x)^\top \theta$ .

$$\phi(x) = \exp(-\gamma \|x - x'\|^2) \quad \phi(x) : \mathbb{R}^1 \rightarrow \mathbb{R}^1 \quad (2.29)$$

and with a new set of corresponding weights, we have:

$$\hat{y} = \exp(-\gamma \|x - x'\|^2) \cdot \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} = \theta_0 \exp(-\gamma \|x - x'\|^2) + \theta_1 \quad (2.30)$$

# 3

## Generalised Linear Models & Gradient Descent

### Note 3.0.1 Course Structure

This part's additional materials will only be helpful after Chapter 6. Divorcing GLMs from their probabilistic context is rarely done, and is only done here for pedagogical purposes.

Thus far we have utilised basis expansion to improve linear models. These models, especially with polynomial basis expansion, are great for modelling variables  $y$  that vary constantly with respect to either their features or basis-expanded features. Even though they are an arbitrarily powerful function approximator (formal statement will be discussed later in the course), they run the risk of overfitting (Chapter 5). We motivate that data assumptions and model assumptions are different.

### Gradient Notation

### Definition 3.0.1

In this course, vectors are assumed to be column vectors by default:

$$\mathbf{b} \in \mathbb{R}^n = \mathbf{b} \in \mathbb{R}^{n \times 1}$$

Gradients of vector-valued functions are, by convention, row vectors:

$$f(\mathbf{b}) : \mathbb{R}^n \rightarrow \mathbb{R} \implies \nabla_{\mathbf{b}} f(\mathbf{b}) \in \mathbb{R}^{1 \times n}$$

$$f(\mathbf{b}) : \mathbb{R}^n \rightarrow \mathbb{R}^m \implies \nabla_{\mathbf{b}} f(\mathbf{b}) \in \mathbb{R}^{m \times n}$$

This layout is referred to as the "numerator" layout.

### 3.1 Generalised Linear Models

Generalised Linear Models use a single non-linear function to encode our beliefs about the output of our model, with the general form:

$$\mathbf{y} = g^{-1}(\theta^\top \mathbf{x}) \tag{3.1}$$

$g$  is the "link" function, traditionally associated with a probability density function from the exponential family. However, we will later discuss when this interpretation will become very useful. A brief overview why GLMs are useful:

- They are flexible with different data distributions. Standard linear regression assumes a normal distribution of errors or outputs, but GLMs can model other distributions such as Poisson, Bernoulli/Binomial, or Poisson.
- Linear regression assumes that the output  $y$  is directly a linear function of the inputs (which we can transform from  $x$  as part of the basis expansion). GLM assumes that some transformation of  $y$  via the link function  $g$  is linear in the inputs. It transforms the output variable instead of the input variables.

### 3.1.1 Poisson Regression

In many applications, we want to predict count information: e.g. how many people in a house, or the number of accidents in a city yearly. However linear regression models do not have restricted outputs and can predict negative values. Poisson Regression is a GLM that addresses this.

Example: I want to predict the number of times a ship will get damaged in a particular period, given some information about the ship.

TYPE	YEAR	PERIOD	MONTHS	Y
b	1965-69	1975-79	20370	53
b	1970-74	1960-74	7064	12
b	1970-74	1975-79	13099	44
b	1975-79	1960-74	0	0
b	1975-79	1975-79	7117	18
b	1960-64	1960-74	44882	39
b	1960-64	1975-79	17176	29
b	1965-69	1960-74	28609	58
c	1960-64	1960-74	1179	1
c	1960-64	1975-79	552	1

Table 3.1: Types of damages and the number of times they occurred in a particular period.

Consider a design matrix  $\mathbf{X}' \in \mathbb{R}^{n \times m}$  and response variables are collected into response vector  $\mathbf{y} \in \mathbb{R}^n$ .

We augment the design matrix, adding a 1 to the end of each vector (affine basis expansion).

$$\mathbf{X}' \rightarrow \mathbf{X} \in \mathbb{R}^{(n+1) \times m} \quad (3.2)$$

To achieve the desired strictly positive co-domain, we simply exponentiate:

$$\hat{\mathbf{y}} = \exp(\mathbf{X}\theta) \quad (3.3)$$

The loss function is given by:

$$\mathcal{L}(\theta) = \sum_{i=1}^m \left( y^{(i)} \mathbf{x}^{(i)\top} \theta - \exp(\mathbf{x}^{(i)\top} \theta) \right) - \sum_{i=1}^m \log(y^{(i)}!) \quad (3.4)$$

The full context and derivation of this loss function will be covered in later chapters (early hint: it is the log-likelihood). This loss function does not have a closed-form solution.

Instead, we can try to use MSE loss instead to find its gradient:

$$\begin{aligned} \mathcal{L}(\theta) &= \frac{1}{m} (\hat{\mathbf{y}} - \mathbf{y})^2 \\ &= \frac{1}{m} (\exp(\mathbf{X}\theta) - \mathbf{y})^2 \\ &\propto (\exp(\mathbf{X}\theta) - \mathbf{y})^2 \\ &= (\exp(\mathbf{X}\theta) - \mathbf{y})^\top (\exp(\mathbf{X}\theta) - \mathbf{y}) \end{aligned}$$

To compute the gradient, the loss is expressed as a summation:

$$\mathcal{L}(\theta) = \sum_{i=1}^k (\exp(\mathbf{X}_{i,:}^\top \theta) - y_i)^2$$

where  $\mathbf{X}_{i,:}^\top$  is the  $i$ -th row of  $\mathbf{X}$ . Applying the chain rule:

$$\frac{\partial}{\partial \theta} (\exp(\mathbf{X}_{i,:}^\top \theta) - y_i)^2 = 2(\exp(\mathbf{X}_{i,:}^\top \theta) - y_i) \cdot \frac{\partial}{\partial \theta} \exp(\mathbf{X}_{i,:}^\top \theta)$$

The derivative of  $\exp(\mathbf{X}_{i,:}^\top \theta)$  with respect to  $\theta$  is:

$$\frac{\partial}{\partial \theta} \exp(\mathbf{X}_{i,:}^\top \theta) = \exp(\mathbf{X}_{i,:}^\top \theta) \cdot \mathbf{X}_{i,:}^\top$$

Thus, the gradient of each term is:

$$\frac{\partial}{\partial \theta} (\exp(\mathbf{X}_{i,:}^\top \theta) - y_i)^2 = 2(\exp(\mathbf{X}_{i,:}^\top \theta) - y_i) \cdot \exp(\mathbf{X}_{i,:}^\top \theta) \cdot \mathbf{X}_{i,:}^\top$$

Summing over all data points:

$$\nabla_{\theta} \mathcal{L}(\theta) = 2 \sum_{i=1}^k (\exp(\mathbf{X}_{i,:}^\top \theta) - y_i) \cdot \exp(\mathbf{X}_{i,:}^\top \theta) \cdot \mathbf{X}_{i,:}^\top$$

In matrix form:

$$\nabla_{\theta} \mathcal{L}(\theta) = 2 ((\exp(\mathbf{X}\theta) - \mathbf{y}) \odot \exp(\mathbf{X}\theta))^\top \mathbf{X}$$

Here,  $\odot$  denotes element-wise multiplication.

### 3.1.2 Logistic Regression

In many applications, we would like to predict the class of a particular input e.g. classifying if an email is spam or not. The simplest case is binary classification, where the output is either 0 or 1. We take  $\mathbf{y} \in [0, 1]$  in this case.

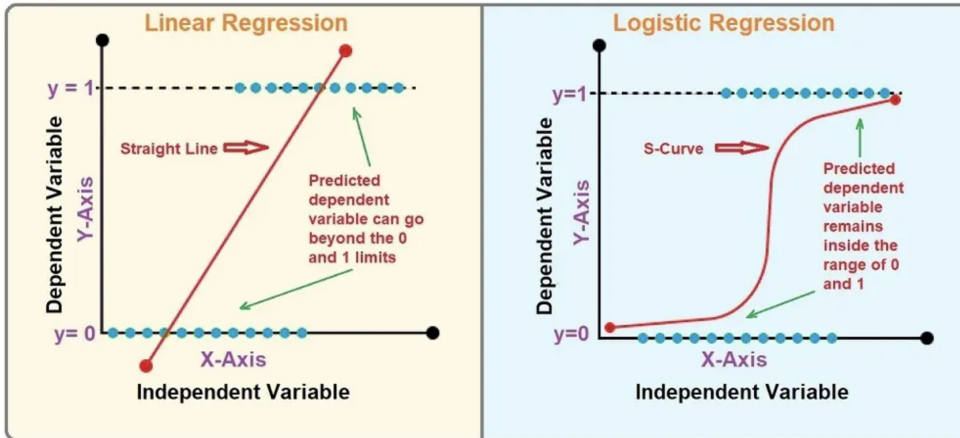


Figure 3.1: Using a contrived, unbalanced dataset, we visualise an intuitive case where the non-linearity in logistic regression helps us model the data more accurately compared to linear regression.

Consider a design matrix  $\mathbf{X}' \in \mathbb{R}^{n \times m}$  and response variables are collected into response vector  $\mathbf{y} \in \mathbb{R}^n$  where entries  $y_i \in \{0, 1\}$ .

We augment the design matrix as usual, adding a 1 to the end of each vector (affine basis expansion).

$$\mathbf{X}' \rightarrow \mathbf{X} \in \mathbb{R}^{(n+1) \times m} \quad (3.5)$$

We then apply the logistic function so that the range of the output is  $[0, 1]$ :

$$\sigma(x) = \frac{1}{1 + \exp(-x)} \quad (3.6)$$

The loss function we desire is the negative log-likelihood for the Bernoulli distribution which is:



$$\mathcal{L}(\theta) = -\frac{1}{m} \sum_{i=1}^m \left( y^{(i)} \log(\sigma(\theta^\top \mathbf{x}^{(i)})) + (1 - y^{(i)}) \log(1 - \sigma(\theta^\top \mathbf{x}^{(i)})) \right)$$

where  $\sigma(x) = \frac{1}{1+\exp(-x)}$  is the logistic function.

Since the minimiser of this equation cannot be solved in closed form, we use numerical methods such as gradient descent.

The gradient of the loss function with respect to  $\theta$  is:

$$\nabla_{\theta} \mathcal{L}(\theta) = -\frac{1}{m} \sum_{i=1}^m \left( y^{(i)} (1 - \sigma(\theta^\top \mathbf{x}^{(i)})) - (1 - y^{(i)}) \sigma(\theta^\top \mathbf{x}^{(i)}) \right) \mathbf{x}^{(i)}$$

This simplifies to:

$$\nabla_{\theta} \mathcal{L}(\theta) = -\frac{1}{m} \sum_{i=1}^m \left( y^{(i)} - \sigma(\theta^\top \mathbf{x}^{(i)}) \right) \mathbf{x}^{(i)}$$

### Linear Models for Classification

### Reference 3.1.1

**Linear models for classification, Chapter 4.1 - 4.1.3 (page 179) Bishop (2006).**

Note that we did not use the MSE approach here— for a more detailed discussion on why this is not ideal, refer to the above reference.

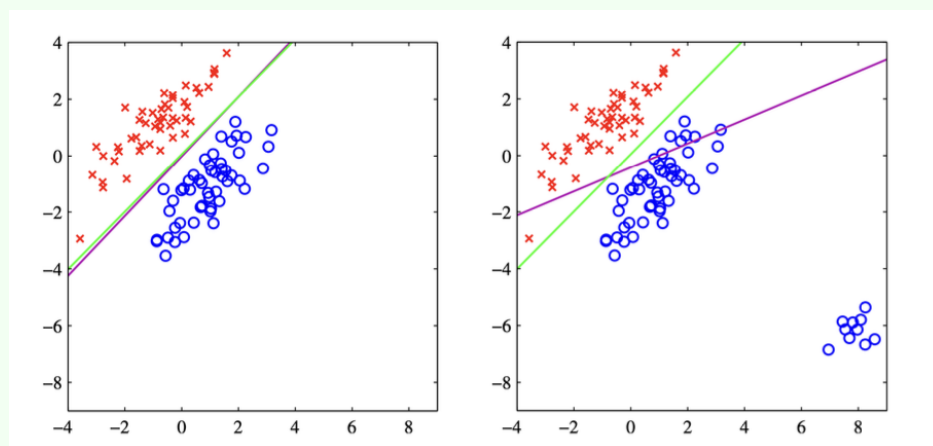


Figure and caption text from Bishop 2006: The left plot shows data from two classes, denoted by red crosses and blue circles, together with the decision boundary found by least squares (magenta curve) and also by the logistic regression model (green curve). The right-hand plot shows the corresponding results obtained when extra data points are added at the bottom left of the diagram, showing that least squares is highly sensitive to outliers, unlike logistic regression.

Why we do not use the MSE loss:

- The MSE itself is meant for continuous output values, so it does not make sense to use it for binary classification.
- The squaring of error in MSE amplifies the impact of large errors or outliers, thus a few outliers can disproportionately influence the whole model.

## 3.2 Gradient Descent

Gradient descent is the core algorithm behind the *learning* portion of many ML algorithms, by finding the minimum of a (loss) function. The idea is to iteratively

take steps in the direction of the steepest descent, which is determined by the gradient of the function. The gradient is denoted by  $\nabla$ , with a subscript indicating the variable we are differentiating with respect to.

The basic form of gradient descent is:

$$\theta^{(t+1)} = \theta^{(t)} + \alpha \nabla_{\theta} \mathcal{L}(f^{\theta}, \mathbf{X}, \mathbf{y}),$$

where  $\alpha$  is the learning rate, and  $\theta^{(0)} \sim \mathcal{N}(0, 1)$ .

We interpret (in the two-dimensional weight case)  $\mathcal{L}(f^{\theta}, \mathbf{X}, \mathbf{y})$  as a height we want to minimise, and visualise  $\theta$  as a puck being pushed in the direction that lowers the height.

Some factors that affect the convergence of gradient descent are:

- Learning rate: If the learning rate is too high, the puck may overshoot the minimum. If it is too low, the puck may take too long to reach the minimum.
- Iterations: A sufficient number of iterations is needed to reach the minimum – if not enough, the optimal loss is not achieved.

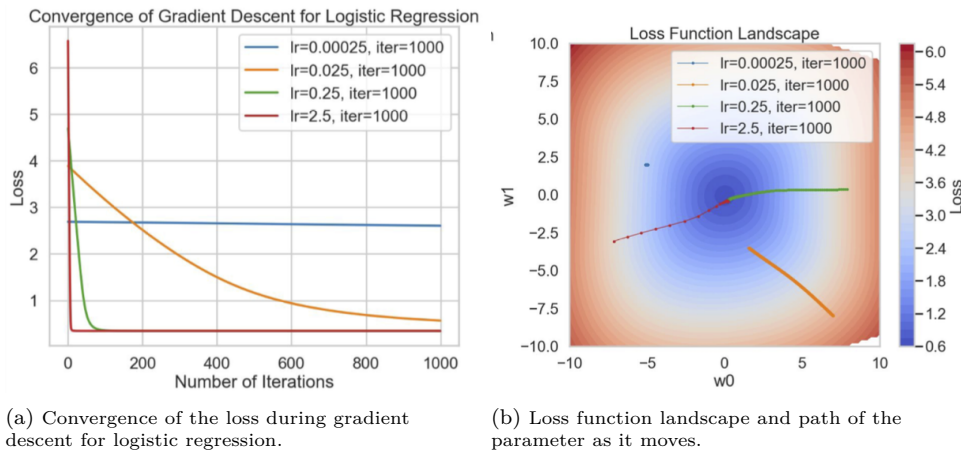


Figure 3.2: Effect of learning rate on logistic regression training. (a) The convergence of the loss over iterations, interpreted as the height of the puck in the bowl. (b) A visualisation of the parameter's path in the loss function landscape.

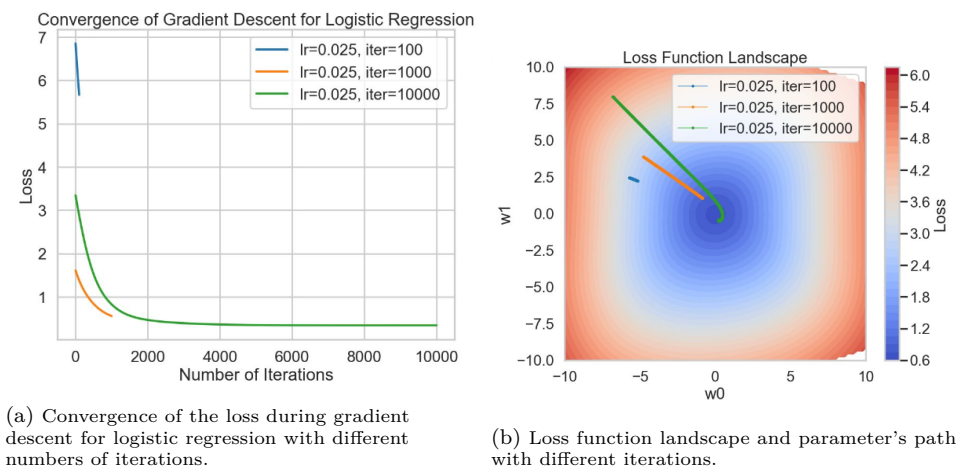


Figure 3.3: Effect of number of iterations on the learning of our logistic regression model. (a) The convergence of the loss over iterations, interpreted as the height of the puck in the bowl. (b) A visualization of the parameter's path in the loss function landscape for different iterations.

### 3.2.1 Complexity of OLS Closed-Form Solution

It is awkward to compare Ordinary Least Squares (OLS) to Gradient Descent (GD) in Generalised Linear Models (GLM), as we cannot directly use the OLS in GLMs.<sup>1</sup>

We can, however, consider the complexity of Ordinary Least Squares (OLS) with Gradient Descent (GD) for Generalised Linear Models (GLMs).

First, consider the OLS closed-form solution:

$$\theta = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

1. The matrix product  $\mathbf{X}^\top \mathbf{X}$  involves multiplying two matrices  $A \in \mathbb{R}^{n \times m}$  and  $B \in \mathbb{R}^{m \times n}$ , which has the general form:

$$C_{i,j} = \sum_{k=1}^l A_{i,k} B_{k,j}$$

- Since  $k$  is our dummy index, it gets summed over  $l$  times.
- It is repeated for each element in the resulting matrix  $C$ .
- Remember, the data matrix  $\mathbf{X}$  is of size  $n \times m$ , so that's  $n$  features and  $m$  samples.

This operation has a complexity of  $O(nlm)$ .

$$\begin{bmatrix} A_{1,1} & A_{1,2} & \dots & A_{1,m} \\ A_{2,1} & A_{2,2} & \dots & A_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ A_{n,1} & A_{n,2} & \dots & A_{n,m} \end{bmatrix} \times \begin{bmatrix} B_{1,1} & B_{1,2} & \dots & B_{1,n} \\ B_{2,1} & B_{2,2} & \dots & B_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ B_{m,1} & B_{m,2} & \dots & B_{m,n} \end{bmatrix} = \begin{bmatrix} C_{1,1} & C_{1,2} & \dots & C_{1,n} \\ C_{2,1} & C_{2,2} & \dots & C_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ C_{n,1} & C_{n,2} & \dots & C_{n,n} \end{bmatrix}$$

Figure 3.4: Matrix product  $\mathbf{X}^\top \mathbf{X}$  involves multiplying two matrices  $A \in \mathbb{R}^{n \times m}$  and  $B \in \mathbb{R}^{m \times n}$ , resulting in a matrix  $C \in \mathbb{R}^{n \times n}$  with a complexity of  $O(n^2m)$ .

Since  $\mathbf{X}^\top \mathbf{X}$  is of size  $n \times n$ , the total complexity of computing this matrix product is  $O(n^2m)$ .

$$\theta = \underbrace{(\mathbf{X}^\top \mathbf{X})^{-1}}_{n \times n} \underbrace{\mathbf{X}^\top}_{n \times m} \underbrace{\mathbf{y}}_{m \times 1}$$

2. Next, matrix inversion for  $(\mathbf{X}^\top \mathbf{X})^{-1}$  (which is square) has complexity<sup>2</sup>  $O(n^3)$ . Combining these operations, the current total complexity becomes:

$$O(n^2m + n^3)$$

3. We multiply the  $\mathbf{X}^\top \mathbf{y}$  term. This is a matrix vector product between  $A \in \mathbb{R}^{n \times m}$  and  $\mathbf{b} \in \mathbb{R}^{m \times 1}$  with complexity  $O(mn)$  (see Figure 3.5).

$$c_j = \sum_{i=1}^n A_{i,j} b_i$$

4. We then multiply  $(\mathbf{X}^\top \mathbf{X})^{-1}$  with  $\mathbf{X}^\top \mathbf{y}$ . This is a matrix vector product between  $A \in \mathbb{R}^{n \times n}$  and  $\mathbf{b} \in \mathbb{R}^{n \times 1}$  with complexity  $O(n^2)$  (see Figure 3.6).

$$c_j = \sum_{i=1}^n A_{i,j} b_i$$

<sup>1</sup> OLS assumes the errors in the model are normally distributed and have constant variance (called homoscedasticity), a common regression assumption. But Generalised Linear Models don't have this condition, and errors can be distributed via the Poisson, binomial, or exponential families. GLMs use Maximum Likelihood Estimation (MLE) instead of the OLS, which will be explained in later chapters.

<sup>2</sup> There are several algorithms for matrix inversion that are better than cubic, but we stick to  $O(n^3)$ .

$$\begin{bmatrix} X_{1,1} & X_{1,2} & \dots & X_{1,n} \\ X_{2,1} & X_{2,2} & \dots & X_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ X_{m,1} & X_{m,2} & \dots & X_{m,n} \end{bmatrix} \times \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_m \end{bmatrix}$$

Figure 3.5: Matrix-vector product for  $\mathbf{X}^\top \mathbf{y}$ 

$$\begin{bmatrix} X_{1,1} & X_{1,2} & \dots & X_{1,n} \\ X_{2,1} & X_{2,2} & \dots & X_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n,1} & X_{n,2} & \dots & X_{n,n} \end{bmatrix} \times \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{bmatrix}$$

Figure 3.6: Matrix-vector product for  $(\mathbf{X}^\top \mathbf{X})^{-1} \cdot (\mathbf{X}^\top \mathbf{y})$ 

We have

$$\theta = \underbrace{\left( \underbrace{\mathbf{X}^\top \mathbf{X}}_{O(n^2 m)} \right)^{-1}}_{O(n^3)} \underbrace{\mathbf{X}^\top \mathbf{y}}_{O(nm)} \quad (3.7)$$

and our final complexity is  $O(n^2 m + n^3)$ . This can be quite inefficient for large datasets where  $m$  is large! We alleviate this by using stochastic gradient descent.

This is because we had  $O(nm + n^2 m + n^2 + n^3) = O(n(m + nm + n) + n^3) \Rightarrow O(n^2 m + n^3)$

### 3.2.2 Complexity of Stochastic Gradient Descent

Instead of having to compute the full OLS solution by inverting the matrix, we can iterate towards an approximate solution by taking the gradient of the loss function and applying it to gradient descent.

We then can optimize the computation by sub-sampling a *batch* of the data, leading to *stochastic gradient descent* (SGD). The gradient for linear regression is given by:

$$\mathbf{X}^\top (\mathbf{X}\theta - \mathbf{y})$$

The operations comprise:

- Matrix-vector multiplication  $\mathbf{X}\theta$  has a complexity of  $O(bn)$  if  $\mathbf{X} \in \mathbb{R}^{b \times n}$ , where  $b$  is the batch size,  $\theta \in \mathbb{R}^{n \times 1}$ ,  $n$  is the number of features.

$$\mathbf{c}_j = \sum_{i=1}^n \mathbf{X}_{i,j} \theta_i$$

$$\begin{bmatrix} X_{1,1} & X_{1,2} & \dots & X_{1,n} \\ X_{2,1} & X_{2,2} & \dots & X_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ X_{b,1} & X_{b,2} & \dots & X_{b,n} \end{bmatrix} \times \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_b \end{bmatrix}$$

- $(\mathbf{X}\theta - \mathbf{y}) \in \mathbb{R}^{b \times 1}$  is a vector subtraction, which has a complexity of  $O(b)$ .
- Gradient computation  $\mathbf{X}^\top (\mathbf{X}\theta - \mathbf{y})$  also has a complexity of  $O(bn)$ .  
 $\mathbf{X}^\top \in \mathbb{R}^{n \times b}$ ,  $(\mathbf{X}\theta - \mathbf{y}) \in \mathbb{R}^{b \times 1}$ .

$$\mathbf{c}_j = \sum_{i=1}^b \mathbf{X}_{i,j} \theta_i$$

$$\begin{bmatrix} X_{1,1} & X_{1,2} & \dots & X_{1,b} \\ X_{2,1} & X_{2,2} & \dots & X_{2,b} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n,1} & X_{n,2} & \dots & X_{n,b} \end{bmatrix} \times \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_b \end{bmatrix} = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{bmatrix}$$

- We have

$$\underbrace{\mathbf{X}^\top \underbrace{(\underbrace{\mathbf{X}\boldsymbol{\theta}}_{O(b)} - \mathbf{y})}_{O(bn)}}_{O(bn)}$$

Performing *stochastic gradient descent* in batches of size  $b$ , iterating  $k$  times results in a total complexity of:

$$O(kbn)$$

For the whole dataset where we just set the batch size to the size of the entire dataset  $n$ , the complexity is:

$$O(kmn)$$

### 3.2.3 Complexity Comparison: OLS vs Gradient Descent

Comparing the complexities:

- OLS:  $O(n^2m + n^3)$
- Gradient Descent:  $O(kmn)$ , with SGD reducing this to  $O(kbn)$

If  $k < n$ , gradient descent becomes more efficient than OLS. Thus, gradient descent (and especially stochastic gradient descent) can be computationally preferable for large datasets.

### 3.2.4 End Notes

#### SGD + Momentum

#### Non-Examinable 3.2.1

We will provide lots of analysis of SGD in this course, but one situation we will not formally analyse is **momentum**:

$$\begin{aligned} \mathbf{z}^{(t+1)} &= \beta \mathbf{z}^{(t)} + \nabla_{\boldsymbol{\theta}} \mathcal{L} \\ \boldsymbol{\theta}^{(t+1)} &= \boldsymbol{\theta}^{(t)} - \alpha \mathbf{z}^{(t+1)} \end{aligned}$$

# 4

## Deep Learning & Automatic Differentiation

Additional Resources	Reference 4.0.1
<ul style="list-style-type: none"> <li>Automatic Differentiation, Section 5.6 (page 159) Deisenroth et al. (2020).</li> <li>Online resource with implementation in Rust.</li> <li>A polished video walk through of automatic differentiation.</li> </ul>	

### 4.1 From GLMs to Multi-Layer Perceptrons

Previously we saw that for generalised linear models, which are only slightly more complex than neural networks, do not have a closed-form solution for finding the optimal parameters (unlike something like linear regression). Recall the generalised linear model (GLM) is given by:

$$\mathbf{y} = \sigma(\theta^T \mathbf{x}) \quad (4.1)$$

where  $\sigma$  is the link function, previously denoted as  $g$ . We perform basis expansion and use summation notation to represent the model as:

$$\mathbf{y} = \sigma\left(\sum_{i=1}^n \theta_i x_i\right) \quad (4.2)$$

This is the model of a single perceptron. A 2-D example is visualised in Figure 4.1.

Now to expand this to a network of multiple perceptrons, we remove the summation notation to add another perceptron:

$$[y_0, y_1] = [\sigma(b_0 + \theta_0^T \mathbf{x}), \sigma(b_1 + \theta_1^T \mathbf{x})]$$

Notice how we now have a second output related to the same input. It then makes more sense to combine the parameter vectors  $\theta_0$  and  $\theta_1$  into a matrix  $\theta := [\theta_0, \theta_1]$ . This allows us to write the model as:

$$\mathbf{y} = \sigma(\theta^T \mathbf{x})$$

Where  $\theta$  is a matrix of parameters,  $\mathbf{y} \in \mathbb{R}^2$  is the output vector, and  $\mathbf{x} \in \mathbb{R}^n$  is the input vector.

#### 4.1.1 The Single-Layer Perceptron

We have previously discussed the  $n$ -dimensional generalisation of the generalised linear model (GLM). This model can now be viewed as a network of neurons. A matrix of size  $m \times n$  represents the weighted connections in a bipartite graph on  $m + n$  nodes, where each entry  $i, j$  of the matrix represents the "strength" of the connection between node  $i$  (the  $i$ th input dimension) and node  $j$  (the  $j$ th output dimension).

By denoting the parameter matrix  $\mathbf{W} \in \mathbb{R}^{m \times n}$ , we can define the multi-layer perceptron as an extension of the GLM:

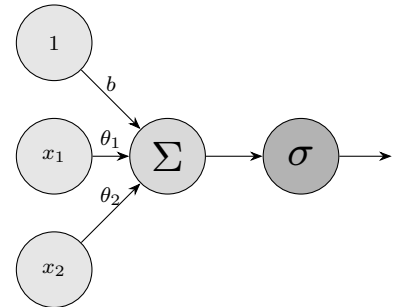


Figure 4.1: A single perceptron.

$$y = g^{-1}(\theta^\top \sigma(\mathbf{W}\mathbf{x} + \mathbf{b}))$$

To clarify, we break this into two stages, highlighting that  $\theta \in \mathbb{R}^{m \times 1}$ :

$$z = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b}) \quad (4.3)$$

$$y = g^{-1}(\theta^\top z) \quad (4.4)$$

- The first stage, Equation 4.3, represents the hidden neurons in the bipartite graph with weights  $\mathbf{W}$ .
- The second stage, Equation 4.4, maps the hidden layer of  $m$  dimensions to the output space, forming another bipartite graph with weights determined by  $\theta$ .

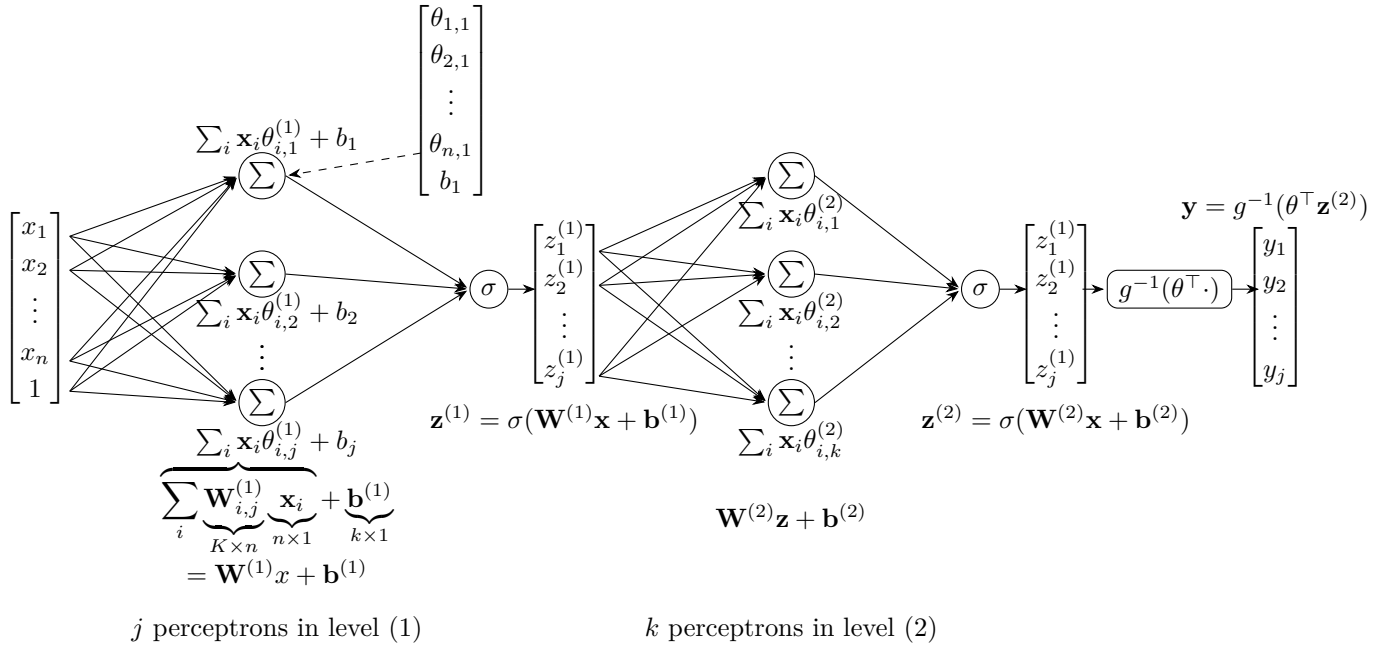


Figure 4.2: A multi-layer perceptron with notation for the input vector  $\mathbf{x}$ , the parameter matrices  $\mathbf{W}^{(1)}$  and  $\mathbf{W}^{(2)}$ , and the output vector  $\mathbf{y}$ . Note that every node marked with a  $\Sigma$  represents a single perceptron. So a Multi-layer perceptron is a network that can have multiple perceptrons per single layer, and have multiple layers.

#### 4.1.2 The Multi-Layer Perceptron

The next step in constructing a perceptron is to allow for multiple layers. By adding additional hidden layers, we can generalise the single-layer perceptron to a multi-layer perceptron. The construction of a two-layer perceptron is as follows, visualised by Figure 4.2:

$$\begin{aligned} \mathbf{z}^{(1)} &= \sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) \\ \mathbf{z}^{(2)} &= \sigma(\mathbf{W}^{(2)}\mathbf{z}^{(1)} + \mathbf{b}^{(2)}) \\ y &= g^{-1}(\theta^\top \mathbf{z}^{(2)}) \end{aligned}$$

- The extension to an arbitrary number of hidden layers follows naturally from this structure, as visualised in Figure 4.3.
- Further questions on this construction are addressed in exercises, and in future lectures there will be more rigorous theoretical analysis.

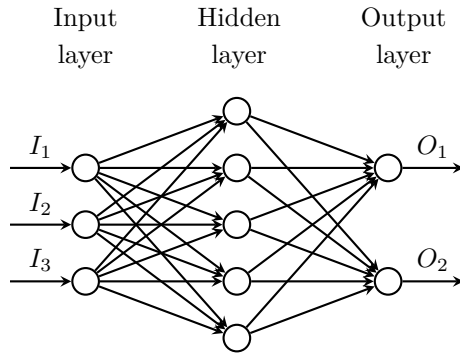


Figure 4.3: The hidden layer is an abstraction that contains the perceptrons in the middle layers in between the input and output layers.

### 4.1.3 Neural Network Architectures

In this section, we explore various neural network layers, focusing on their key structures without constructing them from first principles. We provide some links for further reading.

#### Convolutional Layers

Convolutional layers are designed for processing data with grid-like topology, such as images. The key operation is the convolution, where a filter (kernel) slides over the input data to produce feature maps. Each filter detects features such as edges, textures, or more complex patterns.

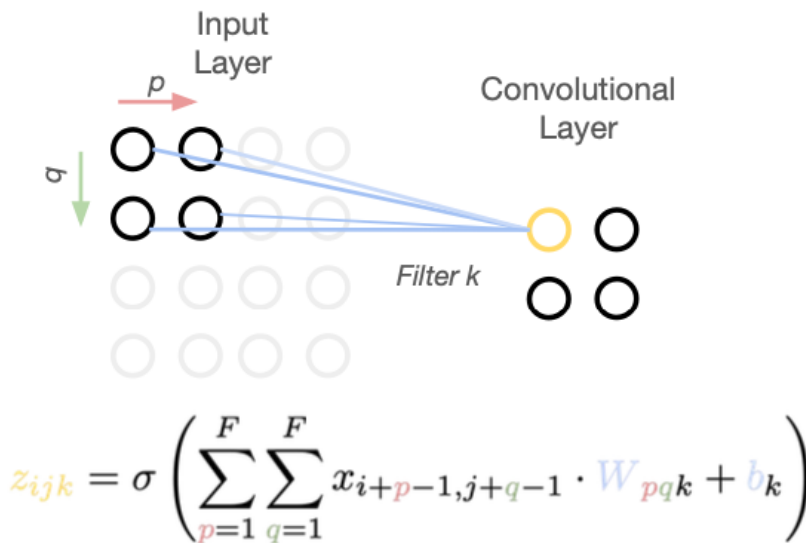


Figure 4.4: A diagram of a single convolutions filter acting on a single channel input with color-coded equation for demonstration purposes. The filter we display is two by two (or in general,  $F \times F$  and the weight contributions (in blue) are multiplied by the input values (in black) and are summed over (red and green).

The convolutional layer is mathematically defined as:

$$z_{ijk} = \sigma \left( \sum_{p=1}^F \sum_{q=1}^F x_{i+p-1, j+q-1} \cdot W_{pqk} + b_k \right)$$

Given an input  $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$  (height  $H$ , width  $W$ , and  $C$  channels) and a set of  $K$  filters  $W_k \in \mathbb{R}^{F \times F \times C}$ , the output feature map  $y \in \mathbb{R}^{H' \times W' \times K}$  is computed as:

$$y_{ijk} = \sigma \left( \sum_{c=1}^C \sum_{p=1}^F \sum_{q=1}^F x_{i+p-1, j+q-1, c} \cdot W_{pqck} + b_k \right)$$

Here,  $\sigma$  is a non-linear activation function, and  $b_k$  is a bias term. The output dimensions  $H'$  and  $W'$  depend on the stride and padding used in the convolution operation.

**Stride:** Stride refers to the step size with which the filter moves across the input. A larger stride reduces the spatial dimensions of the output.

**Padding:** Padding refers to the addition of extra pixels (typically zeros) around the input to control the output size. Padding ensures that the spatial dimensions of the input and output can be maintained or adjusted based on the chosen convolutional operation.



Convolutional layers are foundational in Convolutional Neural Networks (CNNs), which are widely used in tasks such as image classification, object detection, and segmentation.

### Recurrent Layers

Recurrent layers are designed to handle sequential data (order of data points is important). These layers are the core components of Recurrent Neural Networks (RNNs), which are used in tasks such as time series prediction, natural language processing, and speech recognition.

In a recurrent layer, the output at each time step depends on both the current input and the hidden state from the previous time step. This allows the network to maintain a memory of past inputs, capturing temporal dependencies.

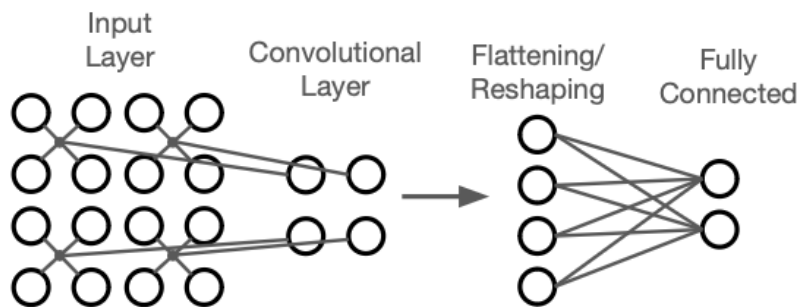


Figure 4.5: A diagram of a simplistic neural network. The first layer is a convolution, followed by flattening and a fully connected layer representing the final output.

Formally, given an input sequence  $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T]$ , the hidden state  $\mathbf{h}_t$  at time step  $t$  is computed as:

$$\mathbf{h}_t = \sigma(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_x \mathbf{x}_t + \mathbf{b}_h)$$

where  $\mathbf{W}_h$  and  $\mathbf{W}_x$  are weight matrices,  $\mathbf{b}_h$  is a bias term, and  $\sigma$  is a non-linear activation function.

Variants such as Long Short-Term Memory (LSTM) networks and Gated Recurrent Units (GRUs) have been developed to address issues like vanishing gradients and to better capture long-term dependencies.

### Attention Layers

Attention layers have revolutionised many areas of machine learning, particularly in natural language processing (NLP). The key idea behind attention mechanisms is to allow the model to focus on different parts of the input sequence when making predictions, rather than treating all parts of the sequence equally. This approach is a significant improvement to the model's ability to handle long sequences and complex dependencies.

A common type of attention mechanism is **self-attention**, where the model computes attention scores within a single sequence. Given an input sequence  $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T]$ , the attention mechanism computes a set of output vectors  $\mathbf{y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_T]$  as:

$$\mathbf{y}_i = \sum_{j=1}^T \alpha_{ij} \mathbf{x}_j$$

In this equation,  $\alpha_{ij}$  represents the attention weight that tells the model how much focus should be given to each input  $\mathbf{x}_j$  when generating the output  $\mathbf{y}_i$ . The attention weights are computed using a compatibility function denoted by  $a_\phi$ .

In the self-attention mechanism, the function  $a_\phi$  is used to compute the attention weights  $\alpha_{ij}$ , which determine how much focus each input  $\mathbf{x}_j$  should receive when constructing the output  $\mathbf{y}_i$ .

The function  $a_\phi$  is called a compatibility function, as it measures how relevant or aligned two inputs  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are. This compatibility score is then converted into a probability-like weight  $\alpha_{ij}$  through a softmax function.

For example, in the commonly used **dot-product attention**, the function  $a_\phi$  is simply the dot product between two input vectors:

$$a_\phi(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$$

This gives a raw score of how similar or aligned  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are. To compute the final attention weights  $\alpha_{ij}$ , these scores are normalised across all inputs using a softmax function:

$$\begin{aligned} \alpha_{ij} &= \frac{\exp(a_\phi(\mathbf{x}_i, \mathbf{x}_j))}{\sum_{k=1}^T \exp(a_\phi(\mathbf{x}_i, \mathbf{x}_k))} \\ &= \frac{\exp(\mathbf{x}_i^\top \mathbf{x}_j)}{\sum_{k=1}^T \exp(\mathbf{x}_i^\top \mathbf{x}_k)} \end{aligned}$$

The process of computing attention scores and generating weighted sums for each output is visually depicted in Figure 4.6. This figure illustrates how each input vector  $\mathbf{e}_j$  is transformed through a function  $f_\psi$  before being scaled by attention weights  $\alpha_{ij}$ . The weighted values are then summed to produce the output vector  $\mathbf{e}'_j$ , demonstrating the flow of information in a self-attention layer.

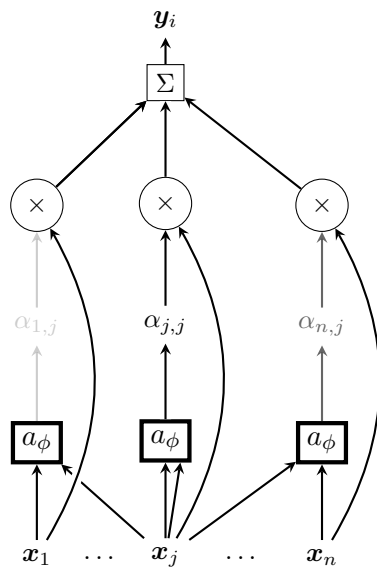


Figure 4.6: Self-attention mechanism: Each input vector  $\mathbf{x}_j$  is weighted by attention weights  $\alpha_{ij}$  and then combined to form the output vector  $\mathbf{y}_i$ .

Self-attention is a key component of modern architectures like **Transformers** which rely heavily on this mechanism. Transformers have become state-of-the-art in many NLP tasks, such as machine translation, text generation, and question answering.

For more detailed readings, consider the seminal papers:

- “Convolutional Neural Networks for Visual Recognition” by Yann LeCun
- “Long Short-Term Memory” by Hochreiter and Schmidhuber
- “Attention is All You Need” by Vaswani et al.

## 4.2 Automatic Differentiation

In the context of Generalised Linear Models (GLMs), gradients are used for learning algorithms such as gradient descent.

$$\theta^{(t+1)} = \theta^{(t)} + \alpha \nabla_{\theta} \mathcal{L}(f^{\theta}, \mathbf{X}, \mathbf{y})$$

We will rewrite this as:

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta} \mathcal{L}(\theta^{(t)})$$

It can sometimes be challenging to compute gradients manually. For neural networks, even small models can make manual gradient computation infeasible. To address this, we use a systematic approach called *automatic differentiation*, which automates the process of computing derivatives.

### 4.2.1 A Running Example

Consider a one-hidden-layer neural network for regression:

$$\begin{aligned} \mathbf{z}^{(0)} &= \mathbf{x} \\ \zeta^{(1)} &= \tanh(\mathbf{W}^{(1)} \mathbf{z}^{(0)} + \mathbf{b}^{(1)}) \\ \mathbf{z}^{(2)} &= \mathbf{W}^{(2)} \zeta^{(1)} + \mathbf{b}^{(2)} \end{aligned}$$

Assuming we have access to i.i.d. data samples, we aim to find parameters that minimise the loss function:

$$\mathcal{L}(\mathbf{y}, \mathbf{z}^{(2)}) = (\mathbf{y} - \mathbf{z}^{(2)})^2$$

To compute gradients such as  $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(1)}}$ , we would need to compute several intermediate derivatives.

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(1)}} &= \left( \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{(2)}} \cdot \frac{\partial \mathbf{z}^{(2)}}{\partial \zeta^{(1)}} \cdot \frac{\partial \zeta^{(1)}}{\partial (\mathbf{W}^{(1)} \mathbf{z}^{(0)} + \mathbf{b}^{(1)})} \right) \frac{\partial (\mathbf{W}^{(1)} \mathbf{z}^{(0)} + \mathbf{b}^{(1)})}{\partial \mathbf{W}^{(1)}} \\ &= \left( (\mathbf{z}^{(2)} - \mathbf{y})^{\top} \mathbf{W}^{(2)} \circ \left( 1 - \tanh^2(\mathbf{W}^{(1)} \mathbf{z}^{(0)} + \mathbf{b}^{(1)}) \right) \right) \mathbf{z}^{(0)\top} \end{aligned}$$

For small networks, manual computation may be feasible, but for deeper models, this process becomes impractical. Plus, everytime we make a minor modification we would need to re-derive our gradients.

An alternative method to compute derivatives is the **finite difference approach**, which approximates the derivative of a function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ . For each input  $x_i$ , the partial derivative can be approximated by:

$$\begin{aligned} \frac{\partial f}{\partial x_1} &= \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(x_1, x_2, \dots, x_n)}{h} \\ \frac{\partial f}{\partial x_2} &= \lim_{h \rightarrow 0} \frac{f(x_1, x_2 + h, \dots, x_n) - f(x_1, x_2, \dots, x_n)}{h} \\ &\vdots \\ \frac{\partial f}{\partial x_n} &= \lim_{h \rightarrow 0} \frac{f(x_1, x_2, \dots, x_n + h) - f(x_1, x_2, \dots, x_n)}{h} \end{aligned}$$

The complexity of this method is  $O(2nT)$ , where  $n$  is the number of parameters and  $T$  is the complexity of evaluating the function or model once. We evaluate the function twice at  $\mathbf{x}$  and  $\mathbf{x} + (\mathbf{1} \cdot h)$  for  $n$  parameters.

However, applying this method to each parameter can be cumbersome and prone to compounded approximation errors, particularly in high-dimensional settings or complex functions.

### 4.2.2 Introducing Automatic Differentiation

The main concept of automatic differentiation is to break down complex equations into simpler operations, each of which can be easily differentiated. We can then apply the chain rule systematically. By understanding the sequence of operations, we can express the overall model as a composition of functions:

$$f(\mathbf{x}) = h(g(\mathbf{x}))$$

Previously, we sought an analytical expression for the Jacobian matrix,  $\nabla_{\mathbf{x}} f$ , to solve for optimal parameters. Consider two functions  $g : \mathbb{R}^n \rightarrow \mathbb{R}^k$  and  $h : \mathbb{R}^k \rightarrow \mathbb{R}^m$ . By applying the chain rule, the Jacobian of the composed function is given by:

$$\mathbf{J}_f = \mathbf{J}_{h \circ g} = \mathbf{J}_h(g(\mathbf{x})) \mathbf{J}_g(\mathbf{x})$$

This general principle allows us to compute **directional derivatives**, which are used to assemble gradients. Given a series of  $L$  compositions of functions, the directional derivative computed via automatic differentiation is:

$$\mathbf{J}\mathbf{x} = \mathbf{J}^{(L)} \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(1)} \mathbf{x}$$

Here, we are not computing the full Jacobian matrix, but rather the **Jacobian-vector product** with respect to the vector  $\mathbf{x}$ , which is referred to as the **directional** derivative.

### 4.2.3 Forward-Mode Automatic Differentiation

The first paradigm of automatic differentiation we will cover is forward-mode. The key idea is that we compute the necessary derivatives of our function concurrently with the values during the forward pass. Forward mode can be viewed as breaking down the Jacobian-vector product in a sequential manner, as shown below:

$$\begin{aligned} \mathbf{J}\mathbf{x} &= \mathbf{J}^{(L)} \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(1)} \mathbf{x} \\ &= \mathbf{J}^{(L)} \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(2)} (\mathbf{J}^{(1)} \mathbf{x}^{(1)}) \\ &= \mathbf{J}^{(L)} \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(3)} (\mathbf{J}^{(2)} \mathbf{x}^{(2)}) \\ &\dots \\ \mathbf{J}\mathbf{x} &= \mathbf{J}^{(L)} \mathbf{x}^{(L-1)} \end{aligned} \tag{4.5}$$

Here,  $\mathbf{x}^{(L-1)}$  denotes the intermediate result after the last Jacobian has been applied.

Now, let us apply this process to a simple neural network model. In practice, when training neural networks, we aim to compute the gradient with respect to the parameters. However, to demonstrate the forward-mode process, we will compute the input gradient step by step. First, recall our example:

$$\begin{aligned} \mathbf{z}^{(0)} &= \mathbf{x} \\ \zeta^{(1)} &= \tanh(\mathbf{W}^{(1)} \mathbf{z}^{(0)} + \mathbf{b}^{(1)}) \\ \mathbf{z}^{(2)} &= \mathbf{W}^{(2)} \zeta^{(1)} + \mathbf{b}^{(2)} \end{aligned}$$

Then, we express the forward pass and the loss function as a sequence of functions that we can program:

$$\mathbf{z}^{(2)} = \text{Matmul}(\mathbf{W}^{(1)}, \tanh(\text{Matmul}(\mathbf{W}^{(0)} \mathbf{z}^{(0)}))) \tag{4.6}$$

#### Note 4.2.1 Side Note on AutoDiff

We are trying to compute  $\alpha \nabla_{\theta} \mathcal{L}(\theta^{(t)})$  in

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta} \mathcal{L}(\theta^{(t)})$$

$\mathbf{J} := \nabla_{\theta} \mathcal{L}(\theta) \in \mathbb{R}^{m \times n}$ ,  
 $\theta^{(t)} \in \mathbb{R}^{n \times 1}$ ,  
 $\mathcal{L}(\theta) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  
 $\alpha > 0$ ,

$$\mathcal{L} = \left( \mathbf{y} - \mathbf{z}^{(2)} \right)^2 \quad (4.7)$$

The purpose of this step is to identify the elementary operations involved in computing the loss, which will help us build the computational graph.

In automatic differentiation literature, these operations are referred to as **primitives**. Each primitive represents a basic operation that is easy to differentiate. For example, the primitives we commonly deal with include functions like addition, multiplication, or applying activation functions (like tanh or raising to a power). For each primitive, the differentiation rules from calculus apply.

Operation	Value Update	Derivative Update
Addition of a constant $c$	$g(w) + c$	$\frac{d}{dw}(g(w) + c) = \frac{d}{dw}g(w)$
Multiplication by a constant $c$	$cg(w)$	$\frac{d}{dw}(cg(w)) = c \frac{d}{dw}g(w)$
Raising to a power $n$	$g(w)^n$	$\frac{d}{dw}(g(w)^n) = n(g(w)^{n-1}) \frac{d}{dw}g(w)$
Applying Tanh	$\tanh(g(w))$	$\frac{d}{dw}(\tanh(g(w))) = 1 - \tanh^2(g(w)) \frac{d}{dw}g(w)$

Figure 4.7: Value and Derivative Updates for Various Operations

Returning to our equations— we ignore the bias terms because they are relegated to affine basis expansion for brevity. We can now compute the forward pass and the derivative of the loss function, step by step:

$$\begin{aligned}
 \mathbf{z}^{(0)} &= \mathbf{x} \\
 \zeta^{(1)} &= \tanh \left( \mathbf{W}^{(1)} \mathbf{z}^{(0)} + \mathbf{b}^{(1)} \right) \\
 \mathbf{z}^{(2)} &= \mathbf{W}^{(2)} \zeta^{(1)} + \mathbf{b}^{(2)} \\
 &= \text{Matmul} \left( \mathbf{W}^{(1)}, \tanh \left( \text{Matmul} \left( \mathbf{W}^{(0)} \mathbf{z}^{(0)} \right) \right) \right) \\
 \mathcal{L} &= \left( \mathbf{y} - \mathbf{z}^{(2)} \right)^2
 \end{aligned}$$

Value	Derivative
$x_0 = x$	$d_0 = \mathbf{1}$
$x_1 = W^{(0)} x_0$	$d_1 = W^{(0)\top} d_0$
$x_2 = \tanh(x_1)$	$d_2 = 1 - \tanh^2(x_1) d_1$
$x_3 = W^{(1)} x_2$	$d_3 = W^{(1)\top} d_2$
$x_4 = (y - x_3)$	$d_4 = -d_3$
$x_5 = (x_4)^2$	$d_4 = 2(x_4) d_4$

Figure 4.8: Forward and Backward Pass for Value and Derivative Calculation.

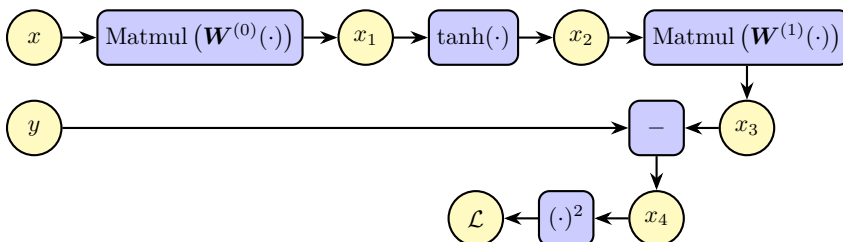


Figure 4.9: Computational Graph for the example neural network.

For an example  $x = 1$ ,  $W^{(0)} = 2$ ,  $W^{(1)} = 1.2$ ,  $y = 2$ , the calculations would look like the following:

Value	Derivative
$x_0 = x \Rightarrow 1$	$d_0 = \mathbf{1}$
$x_1 = W^{(0)}x_0 \Rightarrow 1 \times 2 = 2$	$d_1 = W^{(0)\top}d_0 \Rightarrow 2 \times d_0 = 2$
$x_2 = \tanh(x_1) \Rightarrow \tanh(2) \approx 0.964$	$d_2 = (1 - \tanh^2(2))d_1 \Rightarrow 0.14 \times d_1 \approx 0.14$
$x_3 = W^{(1)}x_2 \Rightarrow 1.2 \times 0.964 \approx 1.15$	$d_3 = W^{(1)\top}d_2 \Rightarrow 1.2 \times 0.14 \approx 0.169$
$x_4 = y - x_3 \Rightarrow 2 - 1.15 = 0.85$	$d_4 = -d_3 \Rightarrow -0.169$
$x_5 = (x_4)^2 \Rightarrow (0.85)^2 = 0.722$	$d_5 = 2(x_4)d_4 \Rightarrow 2(0.85) \times (-0.169) \approx -0.283$

Figure 4.10: Forward and Backward Pass for Value and Derivative Calculation with extensions.

To clarify,  $\mathbf{x}_0$  is the input,  $\mathbf{x}_1$  is the result of the first matrix multiplication,  $\mathbf{x}_2$  is after applying the activation function,  $\mathbf{x}_3$  is the second matrix multiplication,  $\mathbf{x}_4$  represents the absolute error with respect to the label, and  $\mathbf{x}_5$  is the squared error. We computed  $d_5 = \frac{d\mathcal{L}}{dx} = -0.283$ , but this does not provide  $\frac{d\mathcal{L}}{d\mathbf{W}^{(1)}}$ . To compute this, another forward differentiation would be needed, which is beyond this example.

#### 4.2.4 Complexity of Forward Mode Differentiation

Let the forward pass complexity be  $T$ , the number of operations required to compute the function. The computational complexity of finding the Jacobian with forward mode automatic differentiation is  $O(2nT)$ . This is identical to the complexity of finite-difference approximation, but a much less approximation error since approximations are not compounded.

#### 4.2.5 Reverse-Mode Automatic Differentiation

The idea of reverse mode is to begin with an output and propagate backwards, rather than beginning with an input and propagating forwards. This is particularly useful when the number of inputs is much larger than the number of outputs, as is the case with neural networks.

In forward-mode automatic differentiation (AD), we track both the nominal values and the tangents (directional derivatives). We could do this by using the associativity of matrix products that give us our Jacobian matrix, as in Equation 4.5.

Reverse-mode AD, however, focuses on computing adjoints, which can be defined as:

$$\bar{x} = \frac{\partial z}{\partial x}$$

##### ‘Adjoint’ in Automatic Differentiation (AD)

##### Definition 4.2.1

‘Adjoint’ comes from linear algebra, and in the context of automatic differentiation (AD), particularly reverse-mode AD, it refers to the propagation of derivatives.

The adjoint of a variable  $x$  is the derivative of the target function (often the loss function in ML) with respect to  $x$ . It is framed towards the reverse accumulation of derivatives in a system.

$$\bar{x} = \frac{\partial z}{\partial x}$$

This adjoint must have the correct shape to be multiplied by the Jacobian. To illustrate, consider the output of a model,  $\mathbf{y} = f(\mathbf{x})$ , and compute the derivative in reverse order with respect to the output:

$$\begin{aligned}
\mathbf{J}\mathbf{x} &= \bar{\mathbf{y}}\mathbf{J} = (\bar{\mathbf{y}}\mathbf{J}^{(L)})\mathbf{J}^{(L-1)} \dots \mathbf{J}^{(1)} \\
&= (\bar{\mathbf{y}}^{(L)}\mathbf{J}^{(L-1)})\mathbf{J}^{(L-2)} \dots \mathbf{J}^{(1)} \\
&= (\bar{\mathbf{y}}^{(L-1)}\mathbf{J}^{(L-2)})\mathbf{J}^{(L-3)} \dots \mathbf{J}^{(1)} \\
&\dots \\
&= \bar{\mathbf{y}}^{(2)}\mathbf{J}^{(1)}
\end{aligned}$$

In reverse mode, we compute the necessary Jacobians in reverse order, meaning we need intermediate results from the forward pass. For a fully connected neural network with  $L$  layers, the forward pass is given by:

$$\begin{aligned}
\mathbf{z}^{(1)} &= \mathbf{x} \\
\mathbf{z}^{(i+1)} &= \sigma\left(\mathbf{W}^{(i)}\mathbf{z}^{(i)} + \mathbf{b}^{(i)}\right)
\end{aligned}$$

Now, let us consider the weights and biases to be joined into parameters denoted  $\theta^{(i)} = \{\mathbf{W}^{(i)}, \mathbf{b}^{(i)}\}$ . We aim to compute the gradient of our function starting from the last layer using the chain rule:

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \theta^{(L-1)}} &= \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{(L)}} \frac{\partial \mathbf{z}^{(L)}}{\partial \theta^{(L-1)}} \\
\frac{\partial \mathcal{L}}{\partial \theta^{(L-2)}} &= \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{(L)}} \frac{\partial \mathbf{z}^{(L)}}{\partial \mathbf{z}^{(L-1)}} \frac{\partial \mathbf{z}^{(L-1)}}{\partial \theta^{(L-2)}} \\
\frac{\partial \mathcal{L}}{\partial \theta^{(L-3)}} &= \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{(L)}} \frac{\partial \mathbf{z}^{(L)}}{\partial \mathbf{z}^{(L-1)}} \frac{\partial \mathbf{z}^{(L-1)}}{\partial \mathbf{z}^{(L-2)}} \frac{\partial \mathbf{z}^{(L-2)}}{\partial \theta^{(L-3)}} \\
\frac{\partial \mathcal{L}}{\partial \theta^{(i)}} &= \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{(L)}} \frac{\partial \mathbf{z}^{(L)}}{\partial \mathbf{z}^{(L-1)}} \frac{\partial \mathbf{z}^{(L-1)}}{\partial \mathbf{z}^{(L-2)}} \dots \frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{z}^{(1)}} \frac{\partial \mathbf{z}^{(1)}}{\partial \theta^{(i)}}
\end{aligned}$$

We can observe that by chain rule, we can compute the partial derivative of the parameters in layer  $i$  by computing the **partial derivative of the layer w.r.t its input**, and then multiply it with the **partial derivative of the layer output w.r.t its parameters**. Using this recursive application of the chain rule, we propagate the error back through the layers, computing the gradient of the loss  $\mathcal{L}$  with respect to each parameter  $\theta^{(i)}$ .

#### Note 4.2.2 Side Note on AutoDiff

We are trying to compute  $\alpha \nabla_{\theta} \mathcal{L}(\theta^{(t)})$  in

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta} \mathcal{L}(\theta^{(t)})$$

$\mathbf{J} := \nabla_{\theta} \mathcal{L}(\theta) \in \mathbb{R}^{m \times n}$ ,  
 $\theta^{(t)} \in \mathbb{R}^{n \times 1}$ ,  
 $\mathcal{L}(\theta) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  
 $\alpha > 0$ ,

In contrast to Forward Mode:

$$\begin{aligned}
\mathbf{J}\mathbf{x} &= \mathbf{J}^{(L)}\mathbf{J}^{(L-1)} \dots \mathbf{J}^{(1)}\mathbf{x} \\
&= \mathbf{J}^{(L)}\mathbf{J}^{(L-1)} \dots \mathbf{J}^{(2)}(\mathbf{J}^{(1)}\mathbf{x}^{(1)}) \\
&= \mathbf{J}^{(L)}\mathbf{J}^{(L-1)} \dots \mathbf{J}^{(3)}(\mathbf{J}^{(2)}\mathbf{x}^{(2)}) \\
&\dots \\
\mathbf{J}\mathbf{x} &= \mathbf{J}^{(L)}\mathbf{x}^{(L-1)}
\end{aligned}$$

### 4.2.6 General Adjoint/Reverse Mode Algorithm

We now describe the general reverse-mode automatic differentiation algorithm. Let  $x_1, x_2, \dots, x_D$  be the input variables, and  $x_{d+1}, \dots, x_{D-1}$  be the intermediate variables, with  $x_D$  as the output variable. The flow of the computational graph is given by:

$$x_i = g_i(x_{\text{Pa}(x_i)}) \quad \forall i \in [d+1, D] \quad (4.8)$$

Here,  $g_i$  is the elementary operation associated with the node in the graph, and  $\text{Pa}(x_i)$  is the function that returns the set of parent (incoming) nodes to variable  $x_i$ .

We are trying to avoid graph theory notation as much as possible in the notes.

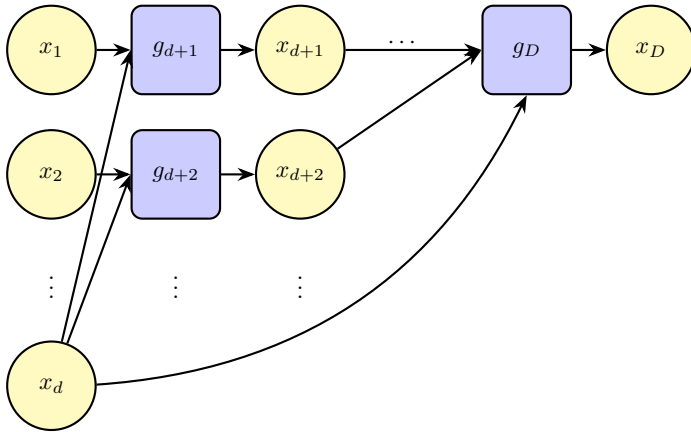


Figure 4.11: Computational Graph notation

The chain rule can then be applied to compute the derivative:

$$\frac{\partial f}{\partial x_i} = \sum_{x_j: i \in \text{Pa}(x_j)} \frac{\partial x_j}{\partial x_i} \frac{\partial f}{\partial x_j} = \sum_{x_j: i \in \text{Pa}(x_j)} \frac{\partial g_j}{\partial x_i} \frac{\partial f}{\partial x_j} \quad (4.9)$$



**Example 4.2.1** (Example of finding the computational graph)

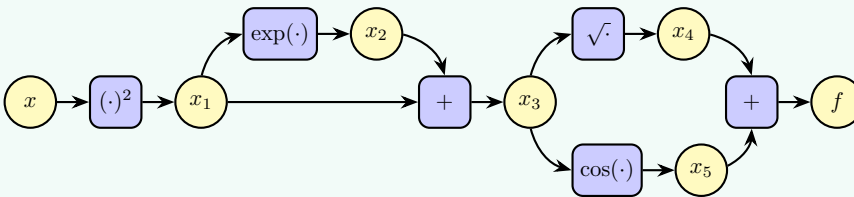
Consider the function:

$$f(x) = \sqrt{x^2 + \exp(x^2)} + \cos(x^2 + \exp(x^2))$$

We break this equation down into its elementary operations: square root, exponential, power, cosine, addition :

$x_1 = x^2$	$\frac{\partial x_1}{\partial x} = 2x$
$x_2 = \exp(x_1)$	$\frac{\partial x_2}{\partial x_1} = \exp(x_1)$
$x_3 = x_1 + x_2$	$\frac{\partial x_3}{\partial x_1} = \frac{\partial x_3}{\partial x_2} = 1$
$x_4 = \sqrt{x_3}$	$\frac{\partial x_4}{\partial x_3} = \frac{1}{2\sqrt{x_3}}$
$x_5 = \cos(x_3)$	$\frac{\partial x_5}{\partial x_3} = -\sin(x_3)$
$x_6 = x_4 + x_5$	$\frac{\partial x_6}{\partial x_4} = \frac{\partial x_6}{\partial x_5} = 1$

Here,  $x_6$  represents the complete function  $f(z)$ . Before performing the backward pass, it is useful to write out the derivative of each function at each step.



Now that we have defined the derivatives, we can construct a computational graph by working backwards from the output. The chain rule is applied step by step:

$$\begin{aligned}
 \frac{\partial x_6}{\partial x_3} &= \frac{\partial x_6}{\partial x_4} \frac{\partial x_4}{\partial x_3} + \frac{\partial x_6}{\partial x_5} \frac{\partial x_5}{\partial x_3} \\
 \frac{\partial x_6}{\partial x_2} &= \frac{\partial x_6}{\partial x_3} \frac{\partial x_3}{\partial x_2} \\
 \frac{\partial x_6}{\partial x_1} &= \frac{\partial x_6}{\partial x_2} \frac{\partial x_2}{\partial x_1} + \frac{\partial x_6}{\partial x_3} \frac{\partial x_3}{\partial x_1} \\
 \frac{\partial x_6}{\partial x} &= \frac{\partial x_6}{\partial x_1} \frac{\partial x_1}{\partial x}
 \end{aligned}$$

While working backwards, ensure all parent nodes are accounted for (summed) to avoid incorrect partial derivatives. Though we only derive the recipe for the derivative and not the actual values, understanding the process outlined in Equations 4.5 is important.

#### 4.2.7 Complexity of Reverse Mode Differentiation

- The function must be evaluated once (forward pass) to get the output, in order to compute all the intermediate values. This has complexity  $T$ .
- The function must be evaluated (with complexity  $T$ ) for each of the  $m$  outputs.
- We get  $O(mT + T)$ . In many cases,  $m$  is much smaller than the number of

inputs, so reverse mode (with complexity  $O(mT + T)$ ) is more efficient than forward mode (with complexity  $O(2nT)$ ).

# 5

## Convexity, Convergence and Optimisation

### Additional Resources

### Reference 5.0.1

- Covering basic ideas: *Numerical Computation*, Chapter 4 (Page 78) Bengio et al. (2017)
- Practical but introductory discussion of gradient descent in deep models, Chapter 8 (Page 271) Bengio et al. (2017)
- More comprehensive overall coverage: *Continuous Optimization*, Chapter 7 (Page 225) up to 7.4 Deisenroth et al. (2020)

### 5.1 Learning as Optimisation

In previous lectures, we explored gradient descent as a method for finding the optimal parameter for a given model and dataset. In this lecture, we will delve deeper into this process by revisiting the optimisation problem we aim to solve:

$$\theta^* = \arg \min_{\theta} \frac{1}{N} \sum_i^N \mathcal{L}(f^{\theta}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})$$

Here,  $\theta^*$  is referred to as the optimal parameter. The argmin operator returns the parameter value that minimises the following loss function:

$$\min_{\theta} \frac{1}{N} \sum_i^N \mathcal{L}(f^{\theta}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})$$

The core algorithm used to solve this minimisation problem is gradient descent, which we previously introduced. The algorithm, stated both in practical and mathematical terms, is as follows:

---

**Algorithm 1:** Gradient Descent

**Input:**  $X$  - Inputs,  $Y$  - Labels,  $\alpha$  - Learning rate,  $K$  - Number of iterations

---

```
1:  $\theta^{(1)} \leftarrow$  Random Initialisation
2: for  $i \in [K]$  do
3:    $l \leftarrow \mathcal{L}(Y, f^{\theta}(X))$ 
4:    $\theta^{(i+1)} \leftarrow \theta^{(i)} - \alpha \nabla_{\theta} l$ 
5: end for
6: return  $\theta_K$ 
```

---

Gradient descent iteratively updates the parameter  $\theta$  using the gradient of the loss function, scaled by the learning rate  $\alpha$ , until convergence or after  $K$  iterations.

## 5.2 Progress Bounds for Gradient Descent

Given our understanding of gradient descent, we have:

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla \mathcal{L}(\theta^{(t)})$$

- $\mathcal{L}(\theta) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,
- $\nabla_{\theta} \mathcal{L}(\theta) \in \mathbb{R}^{m \times n}$ ,
- $\alpha > 0$ ,
- $\mathbf{J} := \nabla_{\theta} \mathcal{L}(\theta)$  (Jacobian matrix).

One critical question about gradient descent is: *does it work?* While practical experiments have established it as a de facto workhorse in machine learning, we will use mathematical tools to understand why and when it works. To do this, we assume that the function is **Lipschitz Continuous**.

### Lipschitz Continuous Gradient

#### Definition 5.2.1

A function  $\mathcal{L}(\theta) : \mathbb{R}^n \rightarrow \mathbb{R}$  has a Lipschitz continuous gradient if there exists a constant  $L > 0$  such that for all  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ , the following inequality holds:

$$\|\nabla \mathcal{L}(\mathbf{u}) - \nabla \mathcal{L}(\mathbf{v})\| \leq L \|\mathbf{u} - \mathbf{v}\|.$$

The constant  $L$  is called the Lipschitz constant, which bounds the rate of change of the gradient.

### Showing $L \geq \|H(\mathbf{v})\|$

#### Definition 5.2.2

For any twice-differentiable function, by Lipschitz continuity, it can be shown that

$$L \geq \|H(\mathbf{v})\|$$

To show this, rewrite the second line assumption of Lipschitz continuity as

$$L \geq \frac{\|\nabla \mathcal{L}(\mathbf{w}) - \nabla \mathcal{L}(\mathbf{v})\|}{\|\mathbf{w} - \mathbf{v}\|}$$

We return to the Taylor Expansion and differentiate again:

$$\begin{aligned} \mathcal{L}(\mathbf{w}) &\approx \mathcal{L}(\mathbf{v}) + \nabla \mathcal{L}(\mathbf{v})^\top (\mathbf{w} - \mathbf{v}) + \frac{1}{2} (\mathbf{w} - \mathbf{v})^\top H(\mathbf{v}) (\mathbf{w} - \mathbf{v}) \\ \nabla \mathcal{L}(\mathbf{w}) &= \nabla \mathcal{L}(\mathbf{v}) + H(\mathbf{v})(\mathbf{w} - \mathbf{v}) + o(\|\mathbf{w} - \mathbf{v}\|) \end{aligned}$$

$$\begin{aligned} \nabla \mathcal{L}(\mathbf{w}) - \nabla \mathcal{L}(\mathbf{v}) &\geq H(\mathbf{v})(\mathbf{w} - \mathbf{v}) \\ \frac{\nabla \mathcal{L}(\mathbf{w}) - \nabla \mathcal{L}(\mathbf{v})}{\|\mathbf{w} - \mathbf{v}\|} &\geq \frac{H(\mathbf{v})(\mathbf{w} - \mathbf{v})}{\|\mathbf{w} - \mathbf{v}\|} \\ \frac{\|\nabla \mathcal{L}(\mathbf{w}) - \nabla \mathcal{L}(\mathbf{v})\|}{\|\mathbf{w} - \mathbf{v}\|} &\geq \|H(\mathbf{v})\| \frac{\|\mathbf{w} - \mathbf{v}\|}{\|\mathbf{w} - \mathbf{v}\|} \end{aligned}$$

Rearranging the rewritten Lipschitz continuity assumption into the equation, we have

$$L \geq \frac{\|\nabla \mathcal{L}(\mathbf{w}) - \nabla \mathcal{L}(\mathbf{v})\|}{\|\mathbf{w} - \mathbf{v}\|} \geq \|H(\mathbf{v})\|$$

Where  $\mathbf{u}$  is the unit norm with  $\|\mathbf{u}\| = 1$ , giving us

$$L \geq \|H(\mathbf{v})\|$$

### 5.2.1 Proving Gradient Descent Minimises $\mathcal{L}$ per step

For a twice-differentiable function  $\mathcal{L}(\theta)$ , we can expand it around a point  $\theta_0$  using a multivariate Taylor expansion. Given two points  $\mathbf{v}$  and  $\mathbf{w} \in \mathbb{R}^n$ , we have:

$$\mathcal{L}(\mathbf{w}) \approx \mathcal{L}(\mathbf{v}) + \nabla \mathcal{L}(\mathbf{v})^\top (\mathbf{w} - \mathbf{v}) + \frac{1}{2}(\mathbf{w} - \mathbf{v})^\top H(\mathbf{v})(\mathbf{w} - \mathbf{v}),$$

where  $H(\mathbf{v})$  is the Hessian matrix of  $\mathcal{L}$  at  $\mathbf{v}$ .

By assuming that the gradient of  $\mathcal{L}$  is Lipschitz continuous with constant  $L$ , it can be shown  $L \geq \|H(\mathbf{v})\|$ , so we can bound the second-order term, giving:

$$\mathcal{L}(\mathbf{w}) \leq \mathcal{L}(\mathbf{v}) + \nabla \mathcal{L}(\mathbf{v})^\top (\mathbf{w} - \mathbf{v}) + \frac{L}{2}\|\mathbf{w} - \mathbf{v}\|^2.$$

This inequality is known as the *descent lemma*, which provides an upper bound on the function value in terms of the gradient and the Lipschitz constant.

Recalling the gradient descent update rule, where at iteration  $t$  we update the parameter  $\theta$  as follows:

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla \mathcal{L}(\theta^{(t)}),$$

we subtract  $\theta^{(t)}$  from both sides to get:

$$\theta^{(t+1)} - \theta^{(t)} = -\alpha \nabla \mathcal{L}(\theta^{(t)}).$$

Recall the gradient descent lemma— we then substitute  $\mathcal{L}(\mathbf{w})$  with  $\mathcal{L}(\theta^{(t+1)})$  and  $\mathcal{L}(\mathbf{v})$  with  $\mathcal{L}(\theta^{(t)})$ :

$$\mathcal{L}(\mathbf{w}) \leq \mathcal{L}(\mathbf{v}) + \nabla \mathcal{L}(\mathbf{v})^\top (\mathbf{w} - \mathbf{v}) + \frac{L}{2}\|\mathbf{w} - \mathbf{v}\|^2.$$

We have:

$$\mathcal{L}(\theta^{(t+1)}) \leq \mathcal{L}(\theta^{(t)}) + \nabla \mathcal{L}(\theta^{(t)})^\top (\theta^{(t+1)} - \theta^{(t)}) + \frac{L}{2}\|\theta^{(t+1)} - \theta^{(t)}\|^2$$

We then substitute the first occurrence of  $\theta^{(t+1)} - \theta^{(t)}$  with  $-\alpha \nabla \mathcal{L}(\theta^{(t)})$ :

$$\mathcal{L}(\theta^{(t+1)}) \leq \mathcal{L}(\theta^{(t)}) - \nabla \mathcal{L}(\theta^{(t)})^\top (\alpha \nabla \mathcal{L}(\theta^{(t)})) + \frac{L}{2}\|\theta^{(t+1)} - \theta^{(t)}\|^2$$

Let  $\alpha = \frac{1}{L}$ , then we have:

$$\mathcal{L}(\theta^{(t+1)}) \leq \mathcal{L}(\theta^{(t)}) - \frac{1}{L}\|\nabla \mathcal{L}(\theta^{(t)})\|^2 + \frac{L}{2}\|\theta^{(t+1)} - \theta^{(t)}\|^2$$

where  $\|\theta^{(t+1)} - \theta^{(t)}\| = \frac{1}{L}\|\nabla \mathcal{L}(\theta^{(t)})\|$ , we obtain:

$$\mathcal{L}(\theta^{(t+1)}) \leq \mathcal{L}(\theta^{(t)}) - \frac{1}{2L}\|\nabla \mathcal{L}(\theta^{(t)})\|^2.$$

This shows that the function value decreases by at least  $\frac{1}{2L}\|\nabla \mathcal{L}(\theta^{(t)})\|^2$  at each step, proving that gradient descent makes progress towards minimising  $\mathcal{L}$ .  $\square$

## 5.3 Convex Optimisation

We have shown of gradient descent moves us in a direction that improves the loss of our model. But as we continue in that direction, where will we end up? We begin answering this question by considering **convex optimisation**.

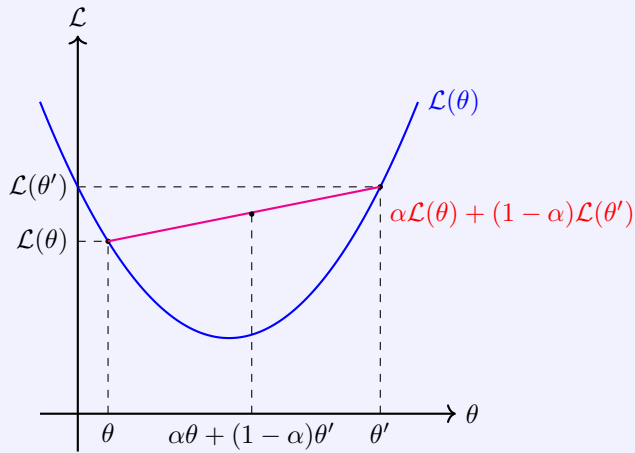
An optimisation problem is convex if and only if the function we are optimising is convex and the domain of the problem is convex. For a function to be convex, it must satisfy:

**Definition 5.3.1. Convex Function****Definition 5.3.1**

We say that a function  $\mathcal{L}(\theta) : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex if and only if:

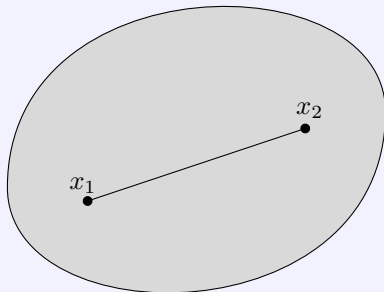
$$\mathcal{L}(\alpha\theta + (1 - \alpha)\theta') \leq \alpha\mathcal{L}(\theta) + (1 - \alpha)\mathcal{L}(\theta')$$

This definition is formulated in terms of the secant of a function and essentially says that every function value between  $\mathcal{L}(\theta)$  and  $\mathcal{L}(\theta')$  must lie below or on the secant. We also require that the domain of the function is a convex set.

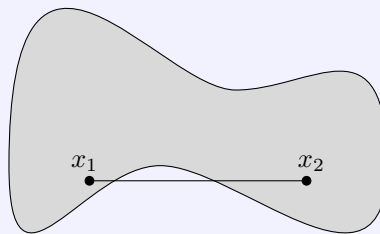
**Definition 5.3.2. Convex Set****Definition 5.3.2**

We say that a set  $C \subseteq \mathbb{R}^n$  is convex if  $\forall x, y \in C$  and  $\forall \alpha \in [0, 1]$ :

$$\alpha x + (1 - \alpha)y \in C$$



Convex set



Non-convex set

An important practical tool for determining if a function is convex is by examining its second-order derivative, known as the Hessian. For twice-differentiable functions, convexity can be characterised by the positive semi-definiteness of the Hessian matrix.

**Theorem 5.3.1** Theorem 5.3.3. Hessian Test for Convexity

Let  $\mathcal{L}(\theta) : \mathbb{R}^n \rightarrow \mathbb{R}$  be a twice-differentiable function. The function  $\mathcal{L}$  is convex if and only if its Hessian matrix  $H(\theta) = \nabla^2 \mathcal{L}(\theta)$  is positive semidefinite (PSD) for all  $\theta \in \mathbb{R}^n$ . That is, for all  $v$  and for all vectors  $v \in \mathbb{R}^n$ ,

$$v^\top H(\theta) v \geq 0$$

*Proof.* The proof follows from the fact that for any twice-differentiable function  $\mathcal{L}(\theta)$ , the second-order Taylor expansion of  $\mathcal{L}$  around a point  $\theta_0$  is given by:

$$\mathcal{L}(\theta) \approx \mathcal{L}(\theta_0) + \nabla \mathcal{L}(\theta_0)^\top (\theta - \theta_0) + \frac{1}{2} (\theta - \theta_0)^\top H(\theta_0) (\theta - \theta_0)$$

If the Hessian  $H(\theta_0)$  is positive semidefinite, the quadratic term is non-negative, meaning that the function lies above its tangent plane at  $\theta_0$ . This is a defining property of a convex function (that you are asked to show as an exercise), and thus, the function  $\mathcal{L}$  is convex.  $\square$

**5.3.1 Global Optimum in Convex Optimisation**

A convex function has a unique global minimum over a convex set, provided it attains a minimum. Thus, if the optimisation problem posed by a machine learning model is convex, there is only one possible value that  $\theta^*$  can take.

**Theorem 5.3.2** Theorem 5.3.4. Global Maximum of a Convex Function

Let  $\mathcal{L}(\theta) : \mathbb{R}^n \rightarrow \mathbb{R}$  be a convex function, and let  $C \subseteq \mathbb{R}^n$  be a convex set. If  $\mathcal{L}$  attains a local maximum at  $\theta^* \in C$ , then  $\theta^*$  is a global maximum. Moreover, the maximum is unique if  $\mathcal{L}$  is strictly convex.

*Proof.* Suppose  $\mathcal{L}$  is convex, and let  $\theta^*$  be a local minimum of  $\mathcal{L}$  in  $C$ . For some neighbourhood  $N \subseteq C$  around  $\theta^*$ , we have  $\mathcal{L}(\theta) \geq \mathcal{L}(\theta^*)$  for all  $\theta \in N$ .

Assume, for contradiction, there exists  $\hat{\theta} \in C$  such that  $\mathcal{L}(\hat{\theta}) < \mathcal{L}(\theta^*)$ . Consider the line segment  $S(\alpha) = \alpha\theta^* + (1 - \alpha)\hat{\theta}$  for  $\alpha \in (0, 1)$ , which lies in  $C$  since  $C$  is convex. By convexity of  $\mathcal{L}$ , we have:

$$\mathcal{L}(S(\alpha)) \leq \alpha\mathcal{L}(\theta^*) + (1 - \alpha)\mathcal{L}(\hat{\theta}) < \alpha\mathcal{L}(\theta^*) + (1 - \alpha)\mathcal{L}(\theta^*) = \mathcal{L}(\theta^*).$$

Since  $\alpha$  can be chosen close enough to 1 such that  $S(\alpha) \in N$ , this contradicts the assumption that  $\theta^*$  is a local minimum.

Thus,  $\mathcal{L}(\theta^*) \leq \mathcal{L}(\theta)$  for all  $\theta \in C$ , meaning  $\theta^*$  is a global minimum of  $\mathcal{L}$  on  $C$ .

To show uniqueness, assume  $\theta_1^*$  and  $\theta_2^*$  are two distinct global minima. For any  $\alpha \in (0, 1)$ , by strict convexity:

$$\mathcal{L}(\alpha\theta_1^* + (1 - \alpha)\theta_2^*) < \alpha\mathcal{L}(\theta_1^*) + (1 - \alpha)\mathcal{L}(\theta_2^*) = \mathcal{L}(\theta_1^*) = \mathcal{L}(\theta_2^*),$$

which contradicts the assumption that both  $\theta_1^*$  and  $\theta_2^*$  are global minima.

Thus, the global minimum must be unique.  $\square$

Here are improved notes for Section 5.4, written concisely while maintaining all the mathematical content:

## 5.4 Convergence in Machine Learning

In machine learning, we aim to determine if gradient descent (GD) will converge to a solution and, if so, under what conditions and how quickly. Convergence is the process by which a sequence approaches a limit. Formally:

Convergence	Definition 5.4.1
A series $x_1, x_2, \dots, x_n$ converges to a limit $T$ if for any $\epsilon > 0$ , there exists an integer $K$ such that $\forall M > K,  x_M - T  < \epsilon$ .	

In GD, our parameters  $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(t)}$  form a sequence. The goal is for this sequence to converge to the optimal parameter:

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(f^\theta(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})),$$

where  $\mathcal{L}$  is the loss function,  $f^\theta$  is the model, and  $N$  is the number of data points. To assess convergence, we ask:

- Under what conditions does the sequence converge?
- What does it converge to?
- How quickly does it converge?

### 5.4.1 Convergence in Linear Regression

Recall the linear model with input features  $\mathbf{x} \in \mathbb{R}^{D \times 1}$  and labels  $\mathbf{y} \in \mathbb{R}$ . Recall the model is:

$$f^\theta(\mathbf{x}, \theta) = \mathbf{x}^\top \theta \quad \hat{y} = f^\theta(\mathbf{x}, \theta).$$

In the context of linear regression, the exact solution for the parameter vector  $\theta^*$  is given by:

$$\theta^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

**Definition:**  $\theta^*$  represents the optimal parameter vector that minimizes the mean squared error (MSE) between the predicted values and the actual target values. Considering the mean squared error (MSE) loss function:

$$\mathcal{L}(\theta) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\theta\|_2^2,$$

the gradient descent (GD) update rule is expressed as:

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta^{(t)}} \mathcal{L}(\theta),$$

where the gradient of the loss function with respect to  $\theta$  is:

$$\nabla_{\theta^{(t)}} \mathcal{L}(\theta) = \mathbf{X}^\top (\mathbf{X}\theta^{(t)} - \mathbf{y}).$$

Substituting this gradient into the GD update rule, we obtain:

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \mathbf{X}^\top (\mathbf{X}\theta^{(t)} - \mathbf{y}).$$

Simplifying the above equation, we get:

$$\theta^{(t+1)} = (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X}) \theta^{(t)} + \alpha \mathbf{X}^\top \mathbf{y}.$$



Thus, we derive the recurrence relation for the gradient descent update rule in linear regression:

$$\theta^{(t+1)} = (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})\theta^{(t)} + \alpha \mathbf{X}^\top \mathbf{y}.$$

## Recurrence Relation and Convergence Analysis

To analyse the convergence of the gradient descent algorithm, we will develop the recurrence relation step by step.

Applying the GD update rule:<sup>1</sup>

<sup>1</sup> Note: The notes for Lecture 5 derive this rather differently.

$$\begin{aligned}\theta^{(1)} &= (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})\theta^{(0)} + \alpha \mathbf{X}^\top \mathbf{y} \\ \theta^{(2)} &= (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})\theta^{(1)} + \alpha \mathbf{X}^\top \mathbf{y} \\ &= (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^2\theta^{(0)} + (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})\alpha \mathbf{X}^\top \mathbf{y} + \alpha \mathbf{X}^\top \mathbf{y} \\ \theta^{(3)} &= (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})\theta^{(2)} + \alpha \mathbf{X}^\top \mathbf{y}, \\ &= (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^3\theta^{(0)} + (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^2\alpha \mathbf{X}^\top \mathbf{y} + (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})\alpha \mathbf{X}^\top \mathbf{y} + \alpha \mathbf{X}^\top \mathbf{y}.\end{aligned}$$

Observing the pattern from the first three iterations, we can generalise the parameter vector at the  $t^{\text{th}}$  iteration as:

$$\theta^{(t)} = (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^t \theta^{(0)} + \sum_{k=0}^{t-1} (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^k \alpha \mathbf{X}^\top \mathbf{y}.$$

This pattern resembles a geometric sequence where each term involves increasing powers of  $(\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})$ . To generalise the pattern, we define the matrix  $\mathbf{A}$  as:

$$\mathbf{A} = \alpha \mathbf{X}^\top \mathbf{X}.$$

Using this definition, the recurrence relation can be expressed as:

$$\begin{aligned}\theta^{(t)} &= (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^t \theta^{(0)} + \left( \sum_{k=0}^{t-1} (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^k \right) \alpha \mathbf{X}^\top \mathbf{y} \\ &= (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^t \theta^{(0)} + \left( \sum_{k=0}^{t-1} (\mathbf{I} - \mathbf{A})^k \right) \alpha \mathbf{X}^\top \mathbf{y}\end{aligned}$$

This summation represents a geometric series. The closed-form solution for the sum of a geometric series involving matrices is given by:

$$\sum_{k=0}^{t-1} (\mathbf{I} - \mathbf{A})^k = (\mathbf{I} - (\mathbf{I} - \mathbf{A})^t) \mathbf{A}^{-1}.$$

provided that  $\mathbf{A}$  is invertible, so  $\mathbf{A}^{-1} = \frac{1}{\alpha} (\mathbf{X}^\top \mathbf{X})^{-1}$ . Substituting this into the expression for  $\theta^{(t)}$ , we obtain:

$$\begin{aligned}\theta^{(t)} &= (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^t \theta^{(0)} + \left( \sum_{k=0}^{t-1} (\mathbf{I} - \mathbf{A})^k \right) \alpha \mathbf{X}^\top \mathbf{y} \\ &= (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^t \theta^{(0)} + (\mathbf{I} - (\mathbf{I} - \mathbf{A})^t) \mathbf{A}^{-1} \alpha \mathbf{X}^\top \mathbf{y} \\ &= (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^t \theta^{(0)} + (\mathbf{I} - (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^t) \left( \frac{1}{\alpha} \mathbf{X}^\top \mathbf{X} \right)^{-1} \alpha \mathbf{X}^\top \mathbf{y} \\ &= (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^t \theta^{(0)} + (\mathbf{I} - (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^t) \theta^*\end{aligned}$$

### Theorem 5.4.1 Geometric Series for Matrices

The summation formula for a geometric series of matrices is given by:

$$\sum_{k=0}^{t-1} (\mathbf{I} - \mathbf{A})^k = (\mathbf{I} - (\mathbf{I} - \mathbf{A})^t) \mathbf{A}^{-1}.$$

*Proof.* Let

$$\mathbf{S}_t = \sum_{k=0}^{t-1} (\mathbf{I} - \mathbf{A})^k.$$

To simplify the sum, multiply  $\mathbf{S}_t$  by  $(\mathbf{I} - \mathbf{A})$ :

$$\begin{aligned}(\mathbf{I} - \mathbf{A})\mathbf{S}_t &= \sum_{k=0}^{t-1} (\mathbf{I} - \mathbf{A})^{k+1} \\ &= \sum_{k=1}^t (\mathbf{I} - \mathbf{A})^k\end{aligned}$$

Notice that this sum can be rewritten as:

$$(\mathbf{I} - \mathbf{A})\mathbf{S}_t = \mathbf{S}_t - \mathbf{I} + (\mathbf{I} - \mathbf{A})^t.$$

Rearranging terms, we get:

$$\mathbf{S}_t - (\mathbf{I} - \mathbf{A})\mathbf{S}_t = \mathbf{I} - (\mathbf{I} - \mathbf{A})^t.$$

Factoring out  $\mathbf{S}_t$  on the left-hand side:

$$(\mathbf{I} - (\mathbf{I} - \mathbf{A}))\mathbf{S}_t = \mathbf{I} - (\mathbf{I} - \mathbf{A})^t.$$

Simplify  $\mathbf{I} - (\mathbf{I} - \mathbf{A}) = \mathbf{A}$ :

$$\mathbf{A}\mathbf{S}_t = \mathbf{I} - (\mathbf{I} - \mathbf{A})^t.$$

Finally, solve for  $\mathbf{S}_t$  by multiplying both sides by  $\mathbf{A}^{-1}$ :

$$\mathbf{S}_t = (\mathbf{I} - (\mathbf{I} - \mathbf{A})^t) \mathbf{A}^{-1}.$$

This completes the proof.

Rearranging the above equation, we have:

$$\theta^{(t)} = (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^t (\theta^{(0)} - \theta^*) + \theta^*.$$

This equation explicitly shows the dependence of  $\theta^{(t)}$  on the optimal parameter vector  $\theta^*$  and the matrix  $(\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^t$ .

### 5.4.2 Conditions for Convergence

For the gradient descent algorithm to converge to the optimal parameter vector  $\theta^*$ , the term  $(\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^t (\theta^{(0)} - \theta^*)$  must approach the zero vector as  $t \rightarrow \infty$ . This occurs if and only if the eigenvalues of  $\mathbf{A} = \mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X}$  have magnitudes less than one.

#### Eigenvalue Analysis:

Let  $\lambda_i$  be the eigenvalues of  $\mathbf{X}^\top \mathbf{X}$ . Then, the eigenvalues of  $\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X}$  are:

$$\mu_i = 1 - \alpha \lambda_i.$$

#### Convergence Criteria:

For  $\mathbf{A}^t$  to approach the zero matrix as  $t \rightarrow \infty$ , all eigenvalues  $\mu_i$  must satisfy:

$$|\mu_i| = |1 - \alpha \lambda_i| < 1 \quad \forall i.$$

Solving the inequality:

$$-1 < 1 - \alpha \lambda_i < 1.$$

This leads to:

$$0 < \alpha \lambda_i < 2 \quad \forall i.$$

$$0 < \alpha < \frac{2}{\lambda_i} \quad \forall i.$$

#### Deriving the Learning Rate Condition:

To satisfy the above condition for all eigenvalues simultaneously, the learning rate  $\alpha$  must satisfy:

$$\alpha < \frac{2}{\lambda_{\max}(\mathbf{X}^\top \mathbf{X})},$$

where  $\lambda_{\max}(\mathbf{X}^\top \mathbf{X})$  is the largest eigenvalue of  $\mathbf{X}^\top \mathbf{X}$ .

#### Implications:

Setting the learning rate  $\alpha$  according to the above condition ensures that all eigenvalues of  $\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X}$  lie within the interval  $(-1, 1)$ . This guarantees that:

$$\lim_{t \rightarrow \infty} (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^t (\theta^{(0)} - \theta^*) = \mathbf{0},$$

thereby ensuring that:

$$\lim_{t \rightarrow \infty} \theta^{(t)} = \theta^*.$$

#### Summary of Convergence:

To ensure convergence of the gradient descent algorithm in linear regression, the learning rate  $\alpha$  must be chosen such that:

$$\alpha < \frac{2}{\lambda_{\max}(\mathbf{X}^\top \mathbf{X})}.$$

This condition guarantees that the parameter vector  $\theta^{(t)}$  converges to the optimal solution  $\theta^*$  as the number of iterations  $t$  approaches infinity. We will discuss why in the next part.

### Distance to Optimal Parameter Vector

The distance from  $\theta^{(t)}$  to  $\theta^*$  in the  $\ell_2$ -norm is given by:

$$\begin{aligned} \|\theta^{(t)} - \theta^*\|_2^2 &= \|(\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^t (\theta^{(0)} - \theta^*)\|_2^2 \\ &= |(\theta^{(0)} - \theta^*)^\top (\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^{2t} (\theta^{(0)} - \theta^*)| \end{aligned}$$

Using the eigenvalue decomposition of  $\mathbf{X}^\top \mathbf{X}$ , we can bound this distance as follows: recall from our review of linear algebra that we can use eigen decomposition to bound the quadratic form involving a symmetric matrix. The relevant fact is:

$$\lambda_{\min}(\mathbf{A}) \|\mathbf{x}\|_2^2 \leq \mathbf{x}^\top \mathbf{A} \mathbf{x} \leq \lambda_{\max}(\mathbf{A}) \|\mathbf{x}\|_2^2.$$

When substituting for the relevant terms, we obtain:

$$\begin{aligned} \|\theta^{(t)} - \theta^*\|_2^2 &\geq \lambda_{\min}((\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^{2t}) \|\theta^{(0)} - \theta^*\|_2^2, \\ \|\theta^{(t)} - \theta^*\|_2^2 &\leq \lambda_{\max}((\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^{2t}) \|\theta^{(0)} - \theta^*\|_2^2. \end{aligned}$$

$$\lambda_{\min}((\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^{2t}) \|\theta^{(0)} - \theta^*\|_2^2 \leq \|\theta^{(t)} - \theta^*\|_2^2 \leq \lambda_{\max}((\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^{2t}) \|\theta^{(0)} - \theta^*\|_2^2.$$

### Convergence Properties in Different Cases:

1. If  $\lambda_{\max} < 1$ : The system always converges.
2. If  $\lambda_{\min} \geq 1$ : The system always diverges.
3. If  $\lambda_{\min} < 1$  but  $\lambda_{\max} \geq 1$ : Convergence depends on the initial condition  $\theta^{(1)}$ .

Here,  $\lambda_{\min}$  and  $\lambda_{\max}$  denote the minimum and maximum eigenvalues of the matrix  $(\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^{2t}$ , respectively. As  $t$  increases, provided that  $|1 - \alpha \lambda_i| < 1$  for all eigenvalues  $\lambda_i$ , both bounds approach zero, indicating that  $\theta^{(t)}$  converges to  $\theta^*$ .

To ensure convergence of our algorithm, we must control the maximum eigenvalue of the matrix. Using eigenvalue properties, let  $\lambda$  be an eigenvalue of  $\mathbf{X}^\top \mathbf{X}$ . Recall the following:

- If  $\lambda$  is an eigenvalue of  $\mathbf{A}$ , then  $c\lambda$  is an eigenvalue of  $c\mathbf{A}$ .
- If  $\lambda$  is an eigenvalue of  $\mathbf{A}$ , then  $1 - \lambda$  is an eigenvalue of  $\mathbf{I} - \mathbf{A}$ .
- If  $\lambda$  is an eigenvalue of  $\mathbf{A}$ , then  $\lambda^2$  is an eigenvalue of  $\mathbf{A}^2$ .

Combining these facts, if  $\lambda$  is an eigenvalue of  $\mathbf{X}^\top \mathbf{X}$ , then  $(1 - \alpha \lambda)^2$  is an eigenvalue of  $(\mathbf{I} - \alpha \mathbf{X}^\top \mathbf{X})^2$ . To guarantee convergence, we require  $(1 - \alpha \lambda_{\max})^2 < 1$ , which simplifies to:

$$\alpha < \frac{2}{\lambda_{\max}(\mathbf{X}^\top \mathbf{X})}.$$

# 6

## Concentration Inequalities & Generalisation Bounds

### 6.1 Generalisation as a Concept

In the previous lecture, we analysed the process of learning. This lecture focuses on evaluating the result of that process. Specifically, after training a model, how do we determine if it is *good*? What does "good" mean? We will examine the concept of *trial deployment* to measure the relative risk of classifiers and decide which one to use in practice. This discussion ties into the idea of generalisation. We will explore *concentration of measure* to theoretically validate trial deployment, particularly from a statistical perspective.

### 6.2 Test Sets

To compare two models, a practical approach is to subject them to a trial deployment and measure performance. We use a *risk function*  $R : \mathbb{R}^{m \times m} \rightarrow \mathbb{R}$ , which takes as input the desired and predicted outputs, and returns a score where larger values indicate worse errors. For  $k$  predictions, the risk of a model  $f^\theta$  can be estimated as:

$$\hat{R}_k(f^\theta) := \frac{1}{k} \sum_{i=1}^k R(f^\theta(\mathbf{x}^{(i)}), \mathbf{y}).$$

Given another model  $g^\omega$ , if  $\hat{R}_k(f^\theta) < \hat{R}_k(g^\omega)$ , we should prefer  $f$  over  $g$ . In practice, instead of actual deployment, we hold out a set of  $k$  input-output pairs from training. These points are not used during training (e.g., gradient descent). This *holdout set* forms the basis of our evaluation.

How large should  $k$  be to conclusively determine whether  $f$  is better than  $g$ ? To answer, we need *concentration inequalities*.

### 6.3 Concentration and the Weak Law of Large Numbers

The length of a trial impacts its outcome. To formalise this, we adopt a probabilistic view of the dataset, treating each test sample as a random variable. For a given model  $f^\theta$ , let  $r = R(f^\theta(x), y)$  be the risk for a sample  $(x, y)$ . The randomness arises from the sample distribution. Ideally, we want our estimated risk  $\hat{r}_k$  to converge to the true risk  $\mathbb{E}[\mathbf{r}]$ , computed as:

$$\mathbb{E}[\mathbf{r}] = \sum_r r P(\mathbf{r} = r).$$

However, we lack access to all possible risk values and their probabilities. Instead, we desire convergence of our empirical estimate  $\hat{r}_k$  to the true risk.

### 6.3.1 Convergence for a Sequence of Random Variables

#### Convergence in Probability

#### Definition 6.3.1

A sequence of random variables  $\{\hat{r}_k\}_{k=1}^{\infty}$  converges **in probability** to a limit  $T$  if, for every  $\epsilon > 0$ :

$$\lim_{k \rightarrow \infty} P(|\hat{r}_k - T| > \epsilon) = 0.$$

This definition implies that the probability of  $\hat{r}_k$  deviating from  $T$  by more than  $\epsilon$  approaches zero as  $k \rightarrow \infty$ . In our context, we aim for  $\mathbb{E}[\mathbf{r}] \approx \hat{r}_k$ , i.e., the sample mean should converge to the true mean.

### 6.3.2 Concentration Inequalities

Recalling the expectation of a discrete random variable  $\mathbf{r}$ :

$$\mathbb{E}[\mathbf{r}] = \sum_r rP(\mathbf{r} = r),$$

we derive the following useful inequality:

$$\mathbb{E}[\mathbf{r}] \geq aP(\mathbf{r} > a).$$

This is **Markov's Inequality**, relating expectation to the probability of exceeding a threshold. By considering the random variable  $(\mathbf{r} - \mu)^2$ , where  $\mu = \mathbb{E}[\mathbf{r}]$ , and setting  $a = \epsilon^2$ , we obtain:

$$P(|\mathbf{r} - \mu| \geq \epsilon) \leq \frac{\sigma^2}{\epsilon^2},$$

known as **Chebyshev's Inequality**. This states that for a small variance, the probability of significant deviation from the mean is small.

### 6.3.3 Weak Law of Large Numbers

The *Weak Law of Large Numbers (WLLN)* ensures that as  $N \rightarrow \infty$ , the empirical risk  $\hat{r}_N$  converges to the true risk  $\mu$ . For  $\hat{r}_N = \frac{1}{N} \sum_{i=1}^N \mathbf{r}_i$ , where  $\mathbf{r}_i$  are i.i.d. samples:

$$\mathbb{E}[\hat{r}_N] = \mu, \quad \mathbb{V}[\hat{r}_N] = \frac{\sigma^2}{N}.$$

Using Chebyshev's inequality, we bound the probability:

$$P(|\hat{r}_N - \mu| \geq \epsilon) \leq \frac{\sigma^2}{N\epsilon^2}.$$

#### Markov's Inequality

To prove Markov's Inequality, consider a non-negative random variable  $\mathbf{r}$  and a threshold  $a > 0$ . By the definition of expectation:

$$\mathbb{E}[\mathbf{r}] = \int_0^{\infty} rf(r) dr,$$

where  $f(r)$  is the probability density function of  $\mathbf{r}$ . Split this integral as:

$$\mathbb{E}[\mathbf{r}] = \int_0^a rf(r) dr + \int_a^{\infty} rf(r) dr.$$

For the second term, observe that  $r \geq a$  when  $r > a$ , so:

$$\begin{aligned} \int_a^{\infty} rf(r) dr &\geq \int_a^{\infty} a f(r) dr \\ &= aP(\mathbf{r} > a) \end{aligned}$$

Thus,

$$\mathbb{E}[\mathbf{r}] \geq aP(\mathbf{r} > a),$$

which is Markov's Inequality.

#### Definition 6.3.2

#### Chebyshev's Inequality

To prove Chebyshev's Inequality, start with the definition of variance for a random variable  $\mathbf{r}$  with mean  $\mu = \mathbb{E}[\mathbf{r}]$ :

$$\sigma^2 = \mathbb{E}[(\mathbf{r} - \mu)^2].$$

For any  $\epsilon > 0$ , let  $A = \{|\mathbf{r} - \mu| \geq \epsilon\}$ . By the definition of probability and expectation:

$$\sigma^2 = \mathbb{E}[(\mathbf{r} - \mu)^2] \geq \int_A (\mathbf{r} - \mu)^2 f(r) dr,$$

since  $(\mathbf{r} - \mu)^2 \geq \epsilon^2$  when  $|\mathbf{r} - \mu| \geq \epsilon$ . Therefore:

$$\sigma^2 \geq \epsilon^2 \int_A f(r) dr = \epsilon^2 P(|\mathbf{r} - \mu| \geq \epsilon).$$

Rearranging gives:

$$P(|\mathbf{r} - \mu| \geq \epsilon) \leq \frac{\sigma^2}{\epsilon^2}.$$

#### Definition 6.3.3

**Weak Law of Large Numbers (WLLN)**

To sketch the proof of the Weak Law of Large Numbers, consider the empirical mean:

$$\hat{r}_N = \frac{1}{N} \sum_{i=1}^N \mathbf{r}_i,$$

where  $\mathbf{r}_i$  are i.i.d. random variables with mean  $\mu$  and variance  $\sigma^2$ . The expectation and variance of  $\hat{r}_N$  are:

$$\mathbb{E}[\hat{r}_N] = \mu, \quad \mathbb{V}[\hat{r}_N] = \frac{\sigma^2}{N}.$$

Using Chebyshev's inequality, we bound the probability of deviation from the mean:

$$P(|\hat{r}_N - \mu| \geq \epsilon) \leq \frac{\mathbb{V}[\hat{r}_N]}{\epsilon^2} = \frac{\sigma^2}{N\epsilon^2}.$$

As  $N \rightarrow \infty$ , the right-hand side converges to 0, implying:

$$P(|\hat{r}_N - \mu| \geq \epsilon) \rightarrow 0.$$

This shows that  $\hat{r}_N$  converges in probability to  $\mu$ , which is the statement of the Weak Law of Large Numbers.

**Definition 6.3.4****6.3.4 Hoeffding's Inequality****Theorem 6.3.1 Hoeffding's Inequality**

For independent random variables  $X_1, X_2, \dots, X_n$  satisfying  $a_i \leq X_i \leq b_i$  almost surely, let  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ . Then, for any  $\epsilon > 0$ :

$$P(|\bar{X}_n - \mathbb{E}[\bar{X}_n]| \geq \epsilon) \leq 2 \exp\left(-\frac{2n\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$

When  $a_i = 0$  and  $b_i = 1$ , this simplifies to:

$$P(|\bar{X}_n - \mathbb{E}[\bar{X}_n]| \geq \epsilon) \leq 2 \exp(-2n\epsilon^2).$$

A proof of Hoeffding's inequality typically requires proof of Chernoff's bound which is a related and useful concentration inequality that requires moment generating functions which is beyond the scope of this class.

**6.4 Universal Function Approximation**

To generalise well, models should align with the data-generating process.

Polynomial basis functions, such as  $\phi(x) = [x, 1]^T$  for affine models or  $\phi'(x) = [x^2, x, 1]^T$  for quadratic models, can approximate functions of arbitrary degree. However, blindly increasing model complexity risks overfitting.

**Weierstrass Theorem****Definition 6.4.1**

For any continuous function  $f$  on  $[a, b]$  and any  $\epsilon > 0$ , there exists a polynomial  $P(x)$  such that:

$$\|f - P\|_\infty = \sup_{x \in [a, b]} |f(x) - P(x)| < \epsilon.$$

This shows that polynomials can uniformly approximate continuous functions. Similarly, neural networks exhibit universal approximation properties:

**Theorem 6.4.1 Universal Approximation (Cybenko, 1989)**

Let  $\sigma(x)$  be a sigmoidal activation function. For any continuous  $f$  on a compact subset of  $\mathbb{R}^n$  and  $\epsilon > 0$ , there exist weights  $w_i$ , biases  $b_i$ , and  $N$  such that:

$$F(x) = \sum_{i=1}^N w_i \sigma(w_i^T x + b_i),$$

satisfies  $|F(x) - f(x)| < \epsilon$ .

### 6.4.1 No-Free Lunch Theorems

- Universal function approximators, such as polynomial basis expansion and neural networks, can both represent a wide range of functions. However, in practice, neural networks are often favoured due to practical considerations.
- Key reasons for choosing neural networks include:
  - **Low sample complexity:** Neural networks often require fewer training examples for effective generalisation.
  - **Low relative computational complexity:** They are typically easier to optimise for many real-world problems.
- The exact reasons behind these advantages remain an open area of research.
- The **No-Free Lunch Theorems** formalise the idea that no single model is universally superior:
  - Given two models that perfectly fit the training data (e.g., a polynomial model and a neural network), there is no inherent reason to prefer one over the other.
  - These theorems state that it is impossible to guarantee which model will perform better on unseen data.
- A major focus in model comparison is studying **inductive biases**:
  - Inductive biases refer to the types of functions that a model is predisposed to learn.
  - Understanding these biases helps in selecting models suited for specific types of data or tasks.
- Future discussions and research often focus on comparing models trained on the same dataset to better understand their performance differences.

## 6.5 Generalisation Error Bounds

- Concentration inequalities allow us to derive bounds on the generalisation error for a hypothesis space  $\Theta$ .
- We study the probability that the true risk  $R(\theta)$  and empirical risk  $R_{\text{emp}}(\theta)$  differ by more than  $\epsilon$  for any  $\theta \in \Theta$ :

$$\mathbb{P} \left[ \sup_{\theta \in \Theta} |R(\theta) - R_{\text{emp}}(\theta)| > \epsilon \right].$$

- The left-hand side is converted into a union bound to manage probabilities across the hypothesis space.

### 6.5.1 Union Bound

- For events  $A_1, A_2, \dots, A_n$ , the **union bound** states:

$$\mathbb{P} \left[ \bigcup_{i=1}^n A_i \right] \leq \sum_{i=1}^n \mathbb{P}(A_i).$$

- Applied to the generalisation problem:

$$\mathbb{P} \left[ \sup_{\theta \in \Theta} |R(\theta) - R_{\text{emp}}(\theta)| > \epsilon \right] \leq \sum_{\theta \in \Theta} \mathbb{P}[|R(\theta) - R_{\text{emp}}(\theta)| > \epsilon].$$

### 6.5.2 Assembling the Bound

- Using Hoeffding's inequality:

$$\mathbb{P}[|R(\theta) - R_{\text{emp}}(\theta)| > \epsilon] \leq 2 \exp(-2n\epsilon^2).$$

- Substituting into the union bound:

$$\mathbb{P} \left[ \sup_{\theta \in \Theta} |R(\theta) - R_{\text{emp}}(\theta)| > \epsilon \right] \leq |\Theta| \cdot 2 \exp(-2n\epsilon^2).$$

- Denoting  $\delta = \mathbb{P}[\sup_{\theta \in \Theta} |R(\theta) - R_{\text{emp}}(\theta)| > \epsilon]$  and solving for  $\epsilon$ , we obtain:

$$|R(\theta) - R_{\text{emp}}(\theta)| \leq \sqrt{\frac{\log(|\Theta|) + \log(2/\delta)}{2n}}.$$

### 6.5.3 Final Generalisation Bound

- The generalisation error is bounded as:

$$R(\theta) \leq R_{\text{emp}}(\theta) + \sqrt{\frac{\log(|\Theta|) + \log(2/\delta)}{2n}}.$$

- This bound shows the trade-off between:
  - **Model complexity:** Larger hypothesis spaces  $|\Theta|$  increase the generalisation error.
  - **Dataset size:** Larger datasets  $n$  decrease the generalisation error.
- The bound highlights the need for regularisation to manage model complexity.



# 7

## Cross-Validation and Bias-Variance Trade-Off

### 7.1 Model Selection

- Model selection aims to find the best combination of hyper-parameters (e.g., learning rate, number of layers, regularisation strength, etc.) to minimise risk.
- Regularisation strength  $\lambda$  is key to controlling model complexity:
  - Higher  $\lambda$ : Reduces flexibility, preventing overfitting.
  - Lower  $\lambda$ : Allows for a more complex model, increasing the risk of overfitting.
- Optimisation problem for hyper-parameters  $\lambda^*$  and  $t^*$ :

$$\lambda^*, t^* = \underset{\lambda, t}{\operatorname{argmin}} \mathbb{E} \left[ R \left( f^{M(\theta', \lambda, t)}(\mathbf{x}), \mathbf{y} \right) \right]$$

Where  $M$  is a function  $M : \mathbb{R}^{n+2} \rightarrow \mathbb{R}$  that given an initial parameter, a regularisation parameter, and amount of time, returns a trained parameter  $M(\theta', \lambda, t) = \theta^{(t)}$ .

- Using the training set for hyper-parameter selection may lead to models that minimise training loss perfectly, making comparisons difficult.
- Using the test set for hyper-parameter selection is flawed, as it leads to overfitting to the test data, making it an unreliable measure of generalisation.

#### Hyper-parameter Selection on the Test Set

#### Intuition 7.1.1

The other perspective is to use the test set. Unfortunately, this is a seriously flawed approach, as it leads to overfitting by effectively "training" on the test set. A simple way to see that this is flawed is to imagine "selecting" the value of  $\theta_1$ . Thus, the optimisation becomes:

$$\lambda^*, t^*, \theta_1^* = \underset{\lambda, t, \theta_1}{\operatorname{argmin}} \min_{\theta/\theta_1} \mathbb{E} \left[ R \left( f^{M(\theta', \lambda, t)}(\mathbf{x}), \mathbf{y} \right) \right]$$

Taking this reasoning to the extreme, we could "select" a value for every parameter and perform comparisons. While this is not a realistic approach to learning, it highlights the issues with using the test set to both select a model and measure performance. Once a model is selected on the basis of specific samples, those samples can no longer serve as an objective measure of model performance.

### 7.2 Regularisation

- Ridge regression (or  $\ell_2$  regularisation) modifies the loss function to add the L2 norm of the parameters to be minimised, with  $\lambda$  as the regularisation strength:

$$\mathcal{L}_{\text{ridge}}(\theta) = \mathcal{L}(\theta) + \lambda \|\theta\|_2^2.$$

#### Regularisation

Regularisation introduces a parameter (e.g., weight decay) to limit model complexity and reduce overfitting.

#### Definition 7.2.1

- The ridge regression loss with  $\lambda$  penalises large parameter magnitudes, enforcing simpler models:

$$\begin{aligned}\mathcal{L}_{\text{ridge}}(\theta) &= \frac{1}{k} \|\mathbf{y} - \mathbf{X}\theta\|^2 + \lambda \|\theta\|_2^2 \\ &= \frac{1}{k} (\mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\theta + \theta^\top \mathbf{X}^\top \mathbf{X}\theta) + \lambda \theta^\top \theta\end{aligned}$$

Taking the derivative with respect to  $\theta$ , we have:

$$\nabla_{\theta} \mathcal{L}_{\text{ridge}} = \frac{2}{k} (\mathbf{X}^\top \mathbf{X}\theta - \mathbf{X}^\top \mathbf{y}) + 2\lambda \theta$$

Setting  $\nabla_{\theta} \mathcal{L}_{\text{ridge}} = 0$  and solving for  $\theta$ :

$$\begin{aligned}0 &= \frac{2}{k} (\mathbf{X}^\top \mathbf{X}\theta - \mathbf{X}^\top \mathbf{y}) + 2\lambda \theta \\ &= \mathbf{X}^\top \mathbf{X}\theta - \mathbf{X}^\top \mathbf{y} + k\lambda \theta \\ (\mathbf{X}^\top \mathbf{X} + k\lambda \mathbf{I})\theta &= \mathbf{X}^\top \mathbf{y} \\ \theta &= (\mathbf{X}^\top \mathbf{X} + k\lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}\end{aligned}$$

### 7.3 Cross-Validation

- We have shown in the previous section why we cannot use the same test-set to measure generalisation and make our hyper-parameter choices.
- A secret “third option” for selecting hyper-parameters is splitting the training set further into a validation set, keeping the test set purely for evaluation.
- $k$ -fold cross-validation:
  - Divide data into  $k$  folds.
  - Train  $k$  times, using  $k - 1$  folds for training and 1 fold for validation.
  - Average the performance across  $k$  iterations:

$$\text{CV}(\mathcal{D}, f) = \frac{1}{k} \sum_{j=1}^k \frac{1}{|\mathcal{D}_j|} \sum_{i=1}^{|\mathcal{D}_j|} R\left(f^{M(\theta', \mathcal{D}_{k \neq j})}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)}\right).$$

- Where  $\mathcal{D}_{k \neq j}$  denotes the set of all folds not equal to the  $j$ th fold.
- Leave-One-Out Cross-Validation (LOOCV) is a special case where  $k$  equals the total number of samples.
  - In each iteration, a single data point is used for validation, and the model is trained on the remaining  $n - 1$  points.
  - This process is repeated  $n$  times, with the final performance averaged across all iterations.

### 7.4 Bias and Variance of Estimators

#### Statistic

#### Definition 7.4.1

A statistic  $S$  is a random variable that is a function of data  $\mathcal{D}$ :  $S = g(\mathcal{D})$ .  $S$  is a random variable not because  $g$  is random, but because the inputs to  $g$  are random.

**Estimator****Definition 7.4.2**

An estimator  $\hat{S}_n$  is a statistic intended to estimate a parameter of the data distribution, where  $\mathcal{D}$  is the data distribution with  $|\mathcal{D}| = n$  datapoints.

An example of an estimator is

$$\hat{r}_k = \frac{1}{n} \sum_{i=1}^k r_i$$

where  $r_i$  is the  $i$ th sample of the random variable  $r$ . This estimator estimates the mean of the random variable  $r$ .

**7.4.1 Linearity of Expectation**

- For any random variables  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$ :
- The expected value of their sum is the sum of the expectations.<sup>1</sup>

$$\mathbb{E} \left[ \sum_{i=1}^k \mathbf{x}_i \right] = \sum_{i=1}^k \mathbb{E}[\mathbf{x}_i].$$

<sup>1</sup> This property holds regardless of whether the random variables are independent, making it widely useful for analysing expectations in various scenarios.

- Constants can be factored out:

$$\mathbb{E}[a \cdot \mathbf{x}] = a \cdot \mathbb{E}[\mathbf{x}].$$

- Sums of random variables and composed can be decomposed:

$$\mathbb{E}[a \cdot \mathbf{x} + b \cdot \mathbf{y}] = a \cdot \mathbb{E}[\mathbf{x}] + b \cdot \mathbb{E}[\mathbf{y}]$$

**Example 7.4.1 (Using Linearity of Expectation)**

As an example, if we need the expectation of

$$\frac{1}{k} \sum_{i=1}^k (r_i - \mu_r)^2,$$

we can use the linearity of expectation to distribute the expectation across the sum:

$$\mathbb{E} \left[ \frac{1}{k} \sum_{i=1}^k (r_i - \mu_r)^2 \right] = \frac{1}{k} \sum_{i=1}^k \mathbb{E} [(r_i - \mu_r)^2].$$

This simplification is crucial when calculating expectations of quantities involving multiple terms, as it allows us to analyse each component individually. These properties will be applied in subsequent steps to study the bias of our variance estimate.

**Bias of an Estimator****Definition 7.4.3**

The bias of an estimator quantifies the systematic deviation of the estimator's expected value from the true parameter it estimates. Let  $\hat{S}$  be an estimator for a parameter  $S$ , and let  $\mathbb{E}[\hat{S}]$  denote the expected value of the estimator. The bias is defined as:

$$\text{Bias}(\hat{S}) = \mathbb{E}[\hat{S}] - S.$$

An estimator is **unbiased** if  $\text{Bias}(\hat{S}) = 0$ , meaning its expected value equals the true parameter  $S$ . Otherwise, a non-zero bias implies the estimator systematically overestimates or underestimates  $S$ .

**Variance of an Estimator****Definition 7.4.4**

The variance of an estimator measures the average squared deviation of the estimator from its expected value, quantifying its variability. For an estimator  $\hat{S}$ , the variance is given by:

$$\text{Var}(\hat{S}) = \mathbb{E}[(\hat{S} - \mathbb{E}[\hat{S}])^2].$$

Smaller variance indicates more consistent estimates across samples, while larger variance suggests higher dispersion around the expected value.

## 7.5 Bias-Variance Decomposition

- The bias-variance decomposition simply states that error of an estimate can be totally accounted for with these two quantities.
- Total error of an estimator decomposes into bias and variance:

$$\text{Err}(\hat{S}) = \mathbb{E}[(\hat{S} - S)^2] = \text{Bias}^2 + \text{Var}.$$

- We expand:

$$\mathbb{E}[(\hat{S} - S)^2] = \mathbb{E}[(\hat{S} - \mu)^2] + (\mu - S)^2.$$

- We have:

$$\mathbb{E}[(\hat{S} - S)^2] = \mathbb{E}\left[(\hat{S} - \mu + \mu - S)^2\right] \quad (7.1)$$

$$= \mathbb{E}\left[(\hat{S} - \mu) + (\mu - S)\right]^2 \quad (7.2)$$

$$= \mathbb{E}\left[(\hat{S} - \mu)^2 + 2(\hat{S} - \mu)(\mu - S) + (\mu - S)^2\right] \quad (7.3)$$

$$= \mathbb{E}[(\hat{S} - \mu)^2] + \mathbb{E}\left[2(\hat{S} - \mu)(\mu - S)\right] + \mathbb{E}[(\mu - S)^2] \quad (7.4)$$

- Next, observe that the middle term

$$\mathbb{E}\left[2(\hat{S} - \mu)(\mu - S)\right]$$

drops out. This is because  $\mu = \mathbb{E}(\hat{S})$ , so  $\mathbb{E}[\hat{S} - \mu] = 0$ , making the product  $(\hat{S} - \mu)(\mu - S)$  have an expectation of zero:

$$\begin{aligned} \mathbb{E}[(\hat{S} - S)^2] &= \mathbb{E}[(\hat{S} - \mu)^2] + (\mu - S)^2 \\ &= \text{Var}(\hat{S}) + \text{Bias}(\hat{S})^2 \end{aligned}$$

## 7.6 Bias-Variance Trade-Off in ML Models

- In Chapter 6, the estimator was for the risk of our model at deployment time.
- Bias-variance decomposition allows us to reason about the bias-variance trade-off of a ML model.
- Are its errors more due to the bias term or variance term?
- Total prediction error decomposes as:

$$\mathbb{E}_{\mathcal{D}}[(f_{\mathcal{D}}^{\theta}(x) - y)^2] = \underbrace{\mathbb{E}_{\mathcal{D}}[(f_{\mathcal{D}}^{\theta}(x) - \bar{f}(x))^2]}_{\text{Variance}} + \underbrace{(\bar{f}(x) - \mu_y)^2}_{\text{Bias}^2} + \underbrace{\mathbb{E}[(\mu_y - y)^2]}_{\text{Label Noise}}.$$

A statistic  $S$  is just some function of our data  $S = g(\mathcal{D})$ . If we let it replace  $\hat{S}$ , then we have:

$$\text{Err}(\hat{S}) = \mathbb{E}[(g(\mathcal{D}) - S)^2]$$

Which is just the MSE. So what is left to do is prove the first equation in Equation 7.1.

### 7.6.1 General Prediction

The **bias-variance trade-off** arises because simple models tend to have high bias but low variance, while complex models exhibit low bias but high variance.

To analyse the sources of error in a machine learning model, we decompose the model's prediction error into three components: variance, squared bias, and label noise. For a given input  $x$  and target  $y$ , the prediction error of a model  $f_{\mathcal{D}}^{\theta}(x)$ , trained on dataset  $\mathcal{D}$ , is expressed as:

$$\mathbb{E}_{\mathcal{D}}[(f_{\mathcal{D}}^{\theta}(x) - y)^2],$$

where the expectation is over all possible datasets  $\mathcal{D}$ . Our goal is to break down this expression to quantify the distinct sources of error. To better understand this, we define the *mean model prediction*:

$$\bar{f}(x) := \mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}^{\theta}(x)],$$

which represents the average output of the model across all datasets.

By adding and subtracting  $\bar{f}(x)$ , we can rewrite the error as:

$$\mathbb{E}_{\mathcal{D}}[(f_{\mathcal{D}}^{\theta}(x) - y)^2] = \mathbb{E}_{\mathcal{D}}[(f_{\mathcal{D}}^{\theta}(x) - \bar{f}(x) + \bar{f}(x) - y)^2].$$

Expanding the square, we obtain:

$$\mathbb{E}_{\mathcal{D}}[(f_{\mathcal{D}}^{\theta}(x) - y)^2] = \mathbb{E}_{\mathcal{D}}[(f_{\mathcal{D}}^{\theta}(x) - \bar{f}(x))^2] + 2\mathbb{E}_{\mathcal{D}}[(f_{\mathcal{D}}^{\theta}(x) - \bar{f}(x))(\bar{f}(x) - y)] + (\bar{f}(x) - y)^2.$$

Since  $\mathbb{E}_{\mathcal{D}}[f_{\mathcal{D}}^{\theta}(x) - \bar{f}(x)] = 0$ , the cross (middle) term vanishes:

$$\mathbb{E}_{\mathcal{D}}[(f_{\mathcal{D}}^{\theta}(x) - y)^2] = \mathbb{E}_{\mathcal{D}}[(f_{\mathcal{D}}^{\theta}(x) - \bar{f}(x))^2] + (\bar{f}(x) - y)^2.$$

Next, we decompose  $(\bar{f}(x) - y)^2$  into its components. Let  $\mu_y := \mathbb{E}[y]$  represent the true mean of the target  $y$ . Adding and subtracting  $\mu_y$ , we get:

$$(\bar{f}(x) - y)^2 = (\bar{f}(x) - \mu_y + \mu_y - y)^2.$$

Expanding this, we have:

$$(\bar{f}(x) - y)^2 = (\bar{f}(x) - \mu_y)^2 + 2(\bar{f}(x) - \mu_y)(\mu_y - y) + (\mu_y - y)^2.$$

Taking the expectation over  $y$ , the cross term  $2\mathbb{E}[(\bar{f}(x) - \mu_y)(\mu_y - y)]$  vanishes because  $\mathbb{E}[\mu_y - y] = 0$ . Thus:

$$\mathbb{E}[(\bar{f}(x) - y)^2] = (\bar{f}(x) - \mu_y)^2 + \mathbb{E}[(\mu_y - y)^2].$$

Finally, substituting this result back, we obtain the full decomposition of the prediction error:

$$\mathbb{E}_{\mathcal{D}}[(f_{\mathcal{D}}^{\theta}(x) - y)^2] = \underbrace{\mathbb{E}_{\mathcal{D}}[(f_{\mathcal{D}}^{\theta}(x) - \bar{f}(x))^2]}_{\text{Variance}} + \underbrace{(\bar{f}(x) - \mu_y)^2}_{\text{Squared Bias}} + \underbrace{\mathbb{E}[(\mu_y - y)^2]}_{\text{Label Noise}}.$$

#### Prediction Error Decomposition

#### Definition 7.6.1

For a model  $f_{\mathcal{D}}^{\theta}(x)$  trained on dataset  $\mathcal{D}$ , the prediction error can be expressed as:

$$\mathbb{E}_{\mathcal{D}}[(f_{\mathcal{D}}^{\theta}(x) - y)^2] = \underbrace{\mathbb{E}_{\mathcal{D}}[(f_{\mathcal{D}}^{\theta}(x) - \bar{f}(x))^2]}_{\text{Variance}} + \underbrace{(\bar{f}(x) - \mu_y)^2}_{\text{Squared Bias}} + \underbrace{\mathbb{E}[(\mu_y - y)^2]}_{\text{Label Noise}}.$$

Here:

- **Variance:** Variability of predictions due to training data.
- **Squared Bias:** Systematic error between the mean prediction  $\bar{f}(x)$  and the true mean  $\mu_y$ .
- **Label Noise:** Irreducible error in the target  $y$ .

### 7.6.2 Bias-Variance in Linear Regression

For a linear regression model, assuming there is no modelling error:

$$\text{P.Err}(\theta) = \mathbb{E}[\|\mathbf{y} - \mathbf{X}\theta\|_2^2],$$

assuming  $\mathbf{y} = \mathbf{X}\theta^*$ . Using ridge regression:

$$\theta^{\text{Ridge}} = (\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y},$$

the bias depends on  $\lambda$ , where increasing  $\lambda$  induces bias but reduces variance. For OLS:

$$\theta^{\text{OLS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y},$$

we observe zero bias but potentially high variance in ill-conditioned settings.

### 7.6.3 Bias-Variance in Linear Regression

To decompose the prediction error of a linear regression model, we analyse its bias and variance. Let the model be parameterised by  $\theta$  and trained on a dataset  $\mathcal{D} := (\mathbf{X}, \mathbf{y})$ . Assuming no modelling error,<sup>2</sup> the predictive error is given by:

$$\text{P.Err}(\theta) = \mathbb{E}[\|\mathbf{y} - \mathbf{X}\theta\|_2^2].$$

Expanding  $\mathbf{y} = \mathbf{X}\theta^*$  and denoting  $\text{Err}(\theta) := \theta - \theta^*$ , we have:

$$\text{P.Err}(\theta) = \mathbb{E}[(\mathbf{X}\text{Err}(\theta))^\top (\mathbf{X}\text{Err}(\theta))].$$

To further understand this, we decompose the error of the learned parameter  $\theta$  as:

$$\begin{aligned} \text{Err}(\theta) &= \mathbb{E}_{\mathcal{D}}[(\theta(\mathcal{D}) - \theta^*)^\top (\theta(\mathcal{D}) - \theta^*)] \\ &= \mathbf{b}(\theta)^\top \mathbf{b}(\theta) + \mathbb{V}(\theta), \end{aligned}$$

where  $\mathbf{b}(\theta) = \mathbb{E}_{\mathcal{D}}[\theta(\mathcal{D})] - \theta^*$  is the *bias*, and  $\mathbb{V}(\theta)$  is the *variance* of  $\theta$ .

#### Ridge and OLS Estimators

The bias and variance decomposition applies to specific estimators:

- **Ridge Regression:**

$$\theta^{\text{Ridge}} = (\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

- **Ordinary Least Squares (OLS):**

$$\theta^{\text{OLS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

Under the no-modelling-error assumption ( $\mathbf{y} = \mathbf{X}\theta^*$ ), we can express the OLS estimator as<sup>3</sup>:

$$\theta^{\text{OLS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\theta^*).$$

Simplifying, we see that:

$$\mathbb{E}[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\theta^*)] = \theta^*,$$

which proves that OLS is an *unbiased* estimator.

For Ridge regression, substituting the true model similarly gives:

$$\mathbb{E}[(\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\theta^*)] = (\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\theta^*.$$

Since the matrix term does not completely cancel, Ridge regression introduces bias:

$$\mathbf{b}(\lambda) := \mathbb{E}_{\mathcal{D}}[\theta^{\text{Ridge}}] - \theta^* = -\mathbb{E}[(\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1}] \lambda \theta^*.$$

<sup>2</sup> (i.e., there exists  $\theta^*$  such that  $\mathbf{y} = \mathbf{X}\theta^*$ )

<sup>3</sup> All we have done is substituted in the true model for  $\mathbf{y}$  since there is no mismatch. Proving that the OLS is an unbiased estimator of  $\theta^*$  is shown by verifying:

When  $\lambda = 0$ , Ridge regression reduces to OLS, recovering the unbiased case.

### Conclusion of Ridge vs OLS

### Intuition 7.6.1

Ridge regression trades increased bias for reduced variance, stabilising the solution in ill-conditioned or high-dimensional settings. OLS, while unbiased, can exhibit high variance when  $\mathbf{X}^\top \mathbf{X}$  is nearly singular.

## 7.6.4 Variance of Linear Models

The variance of an estimator quantifies how much the learned parameter  $\theta$  fluctuates around its expected value due to changes in the training dataset  $\mathcal{D}$ . Formally, the variance is defined as:

### Variance

### Definition 7.6.2

$$\mathbb{V}(\theta) = \mathbb{E}_{\mathcal{D}}[(\theta(\mathcal{D}) - \mathbb{E}_{\mathcal{D}}[\theta(\mathcal{D})])^\top (\theta(\mathcal{D}) - \mathbb{E}_{\mathcal{D}}[\theta(\mathcal{D})])].$$

### Variance of OLS Estimator

### Definition 7.6.3

For the ordinary least squares (OLS) estimator, assuming  $\mathbf{y} = \mathbf{X}\theta^*$ , we have:

$$\theta^{\text{OLS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \theta^*.$$

Simplifying gives  $\theta^{\text{OLS}} = \theta^*$ . Thus, in this idealised case, OLS is unbiased, and its variance depends entirely on the variance of the dataset  $\mathbf{X}^\top \mathbf{X}$ . If this variance is high (e.g., due to ill-conditioned  $\mathbf{X}^\top \mathbf{X}$ ), the estimate for  $\theta^*$  becomes unstable.

### Variance of Ridge Estimator

### Definition 7.6.4

For ridge regression, the estimator is:

$$\theta^{\text{Ridge}} = (\lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \theta^*.$$

Simplifying gives:

$$\mathbb{V}(\theta^{\text{Ridge}}) = \mathbb{E}_{\mathcal{D}}[(\theta^{\text{Ridge}} - \mathbb{E}_{\mathcal{D}}[\theta^{\text{Ridge}}])^\top (\theta^{\text{Ridge}} - \mathbb{E}_{\mathcal{D}}[\theta^{\text{Ridge}}])].$$

The variance reflects the stabilising effect of the regularisation term  $\lambda$ . By adding  $\lambda \mathbf{I}$ , ridge regression reduces variance, particularly when  $\mathbf{X}^\top \mathbf{X}$  is nearly singular or ill-conditioned.

- Ill-conditioned (eigenvalues that are very small or close to zero)  $\mathbf{X}^\top \mathbf{X}$  can lead to high variance in OLS, because the inverse becomes unstable.
- When  $\mathbf{X}$  has multicollinearity (highly correlated predictors), the eigenvalues of  $\mathbf{X}^\top \mathbf{X}$  tend to be small in such cases.
- Adding  $\lambda \mathbf{I}$  ensures that the eigenvalues of  $\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}$  are shifted by at least  $\lambda$ , which increases the smallest eigenvalues, and improves the conditioning of the matrix.
- Better conditioning reduces numerical instability, making the solution more robust.

**Summary****Intuition 7.6.2**

OLS achieves unbiased estimates but can exhibit high variance when  $\mathbf{X}^\top \mathbf{X}$  is poorly conditioned. Ridge regression trades a small bias for reduced variance, stabilising the parameter estimates by adding  $\lambda \mathbf{I}$  to  $\mathbf{X}^\top \mathbf{X}$ .



# 8

## Introduction to Probability

In previous lectures, we approached

Sigma Algebra and Borel Measures	<i>Non-Examinable 8.0.1</i>
<ul style="list-style-type: none"> <li>The definition of Sigma Algebra</li> <li>Borel Measurability</li> <li>Proof requiring you think about what is measurable and what isn't</li> </ul>	

The ideas of CDFs, PMFs, PDFs, and their corresponding probability tools will only be examined.

### 8.0.1 Notation Preliminaries

Notation	Description
$\mathbf{x} \sim p(\mathbf{x})$	$\mathbf{x}$ distributed according to the pdf $p(\mathbf{x})$
$x \sim p(\mathbf{x})$	$x$ sampled from the r.v. $\mathbf{x}$
$\mathbb{E}[\mathbf{x}]$	The expectation of the random variable $\mathbf{x}$
$\mathbb{E}_{\mathbf{x}}[x]$	The expectation of the random variable $\mathbf{x}$

Table 8.1: Notation for distributions and expectations

Not random	Random
$x, \mathbf{x}$	$\mathbf{X}, \mathbf{X} x, \mathbf{X}, X$

Table 8.2: Notation distinguishing between deterministic and random variables. The  $X$  is the only random variable that is not bolded– the exception is so it's more easily handwritten during an exam!

### 8.1 Sample and Event Space

Sample Space	Definition 8.1.1
<p>The sample space is the collection of all possible outcomes. For example, a coin flip has the sample space <math>\{H, T\}</math>, where <math>H</math> is heads and <math>T</math> is tails. A dice roll has the sample space <math>\{1, 2, 3, 4, 5, 6\}</math>.</p> <p>The sample space of two consecutive coin flips is <math>\{HH, HT, TH, TT\}</math>.</p>	

Event Space	Definition 8.1.2
<p>This is the power set of the sample space and comprises every outcome of an experiment we can measure. The sample space of two consecutive coin flips has the event space:</p> $\{\emptyset, \{HH\}, \{HT\}, \{TH\}, \{TT\}, \{HH, HT\}, \{HH, TH\}, \{HH, TT\}, \{HT, TH\}, \{HT, TT\}, \{TH, TT\}, \{HH, HT, TH\}, \{HH, HT, TT\}, \{HH, TH, TT\}, \{HT, TH, TT\}, \{HH, HT, TH, TT\}\}$	

### 8.1.1 Probability of Events

Given our event space, we want to build our idea of probability. There are a few distinct ways to go about reasoning about probabilities:

1. Long-run relative frequency of events (Frequentist)
2. Express our beliefs about the relative likelihood of events (Bayesian)

## 8.2 Sigma Algebra

There are two rules of a sigma algebra:

1. Closed under complements
2. Closed under countable unions

**Examples of Discrete Sigma Algebras:**

$$A = \{\emptyset, \sigma\}$$

$$A = \{\emptyset, E, E^c, \sigma\}$$

**Examples of Continuous Sigma Algebras:**

$$\sigma = \mathbb{R}$$

$$\tau = \{(a, b) \mid \forall a < b \in \mathbb{R}\}$$

$A = B(\tau)$   $B$  is the operation that returns the smallest sigma algebra on  $\tau$  or the Borel set

## 8.3 Measure Theory

### Definition 8.2.1: $\sigma$ -Algebra

A  $\sigma$ -algebra  $\mathcal{A}$  on a non-empty set  $\Omega$  is a collection  $\mathcal{A} \subseteq \mathcal{P}(\Omega)$ , the power set of  $\Omega$ , which is closed under complements and countable unions. Formally:

1. **Closed under complements:** If  $B \in \mathcal{A}$ , then  $B^c \in \mathcal{A}$ .
2. **Closed under countable unions:** If  $B_1, B_2, \dots \in \mathcal{A}$ , then  $\bigcup_{i=1}^{\infty} B_i \in \mathcal{A}$ .

### Definition 8.3.1

#### The Rules

We can generate probability distributions that work for our models so long as they satisfy these three rules of a probability measure.

#### Intuition 8.3.1

### Definition 8.2.2: Probability Measure

Some abuse of notation here!

$$\Omega = \{i\}_{i=1}^n$$

$$\mathcal{A} = \mathcal{P}(\Omega)$$

$$P(A \in \mathcal{A}) = \frac{|A|}{n}$$

A probability measure  $P$  on a measurable space  $(\Omega, \mathcal{A})$  is a function  $P : \mathcal{A} \rightarrow [0, 1]$  that satisfies:

1.  $P(\emptyset) = 0$  and  $P(\Omega) = 1$ .
2. **Non-negativity:** For all  $B \in \mathcal{A}$ ,  $P(B) \geq 0$ .
3. **Countable additivity:** For any pairwise disjoint collection  $\{B_i\}_{i=1}^{\infty} \subseteq \mathcal{A}$ , we have

$$P\left(\bigcup_{i=1}^{\infty} B_i\right) = \sum_{i=1}^{\infty} P(B_i).$$

### Definition 8.3.2

**Definition 8.2.3: Cumulative Distribution Function** **Definition 8.3.3**

A cumulative distribution function  $F : \mathbb{R} \rightarrow [0, 1]$  for a probability measure satisfies:

1. **Monotonicity:** If  $x \leq y$ , then  $F(x) \leq F(y)$ .
2. **Right-continuity:**  $\lim_{x \rightarrow y^+} F(x) = F(y)$ .
3.  $\lim_{x \rightarrow -\infty} F(x) = 0$  and  $\lim_{x \rightarrow \infty} F(x) = 1$ .

If the CDF is differentiable, it has an associated probability density function (PDF),  $f(x) = \frac{dF(x)}{dx}$ , which represents probabilities over small intervals.

**Rules for the CDF****Intuition 8.3.2**

If we have a CDF that satisfies the above rules, we can then write its derivative,  $p(x)$ , as the probability density function (PDF).

Using this we can define the measure corresponding to the CDF.

$$P((a, b]) = \int_a^b p(x) dx$$

## 8.4 Joint Probability Distributions

We might have probability distributions over two variables:

$$P(\mathbf{x}, \mathbf{y}) = P(\mathbf{x} = x, \mathbf{y} = y)$$

We can view this as a surface over the two variables, where the height of the surface is the probability of that point.

This generalises to many variables:

$$p(x_1, x_2, \dots, x_n) = P(\mathbf{x}_1 = x_1, \mathbf{x}_2 = x_2, \dots, \mathbf{x}_n = x_n)$$

Y	X			
	0	1	2	3
0	0.15	0.1	0.0875	0.0375
1	0.1	0.175	0.1125	0
2	0.0875	0.1125	0	0
3	0.0375	0	0	0

Table 8.3: Joint probability distribution of  $X$  and  $Y$

We can consider the probability of one variable given the other is equal to some known value:

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

Given a joint distribution  $P(\mathbf{x}, \mathbf{y})$ , we may be interested in the marginal distributions  $P(\mathbf{x})$  or  $P(\mathbf{y})$ . For instance, after assessing the likelihood of the friend's claim, we may want to examine  $P(\mathbf{x} = x)$ , representing the distribution of signal bars in our area. Marginalization allows us to isolate  $P(\mathbf{x})$  from  $P(\mathbf{x}, \mathbf{y})$  by summing over values of other variables. For continuous cases, this becomes:

$$P(\mathbf{x} = x) = \int_{\Omega_y} P(\mathbf{x} = x, \mathbf{y} = y) dy$$

# 9

## Reference Boxes

Theorem 9.0.1 Theorem Name

This is the statement of the theorem.

Corollary 9.0.1 Corollary Name

This is the statement of the corollary.

Lemma 9.0.1 Lemma Name

This is the statement of the lemma.

Claim 9.0.1 Claim Name

This is the statement of the claim.

Example 9.0.1 (Example Name)

This is the explanation of the example.

Note 9.0.1 Side Note Box

This is a side note.

This is a block of highlighted text

Definition Title

Definition 9.0.2

This is an example definition.

Extra Title

Non-Examinable 9.0.2

This is an example box with extra information.

Example Title

Example Q 9.0.2

This is an example question.

Answer here

Q1c - 2018

Exam Q 9.0.2

This is an example exam question.

Reference Title

Reference 9.0.2

This is an example reference to source material.

Note 9.0.1 Side Note Box

This is a smaller side note.

This is a block of highlighted text that's smaller.

Small Def Title

Example Text

Definition 9.0.1

Small Title

Example Text

Non-Examinable 9.0.1

Small Title

Example Text

Answer

Example Q 9.0.1

Q1c - 2018

This is an example exam question, smaller.

Exam Q 9.0.1

Reference Title

This is an example reference to source material.

Reference 9.0.1

## Intuition Title

This is an example of an intuitive explanation, smaller

Intuition 9.0.1

## Intuition Title

## Intuition 9.0.2

This is an example of an intuitive explanation.

THE FRONT MATTER of a book refers to all of the material that comes before the main text. The following table from shows a list of material that appears in the front matter of *The Visual Display of Quantitative Information*, *Envisioning Information*, *Visual Explanations*, and *Beautiful Evidence* along with its page number. Page numbers that appear in parentheses refer to folios that do not have a printed page number (but they are still counted in the page number sequence).

Page content	Books			
	<i>VDQI</i>	<i>EI</i>	<i>VE</i>	<i>BE</i>
Blank half title page	(1)	(1)	(1)	(1)
Frontispiece <sup>1</sup>	(2)	(2)	(2)	(2)
Full title page	(3)	(3)	(3)	(3)
Copyright page	(4)	(4)	(4)	(4)
Contents	(5)	(5)	(5)	(5)
Dedication	(6)	(7)	(7)	7
Epigraph	–	–	(8)	–
Introduction	(7)	(9)	(9)	9

<sup>1</sup> The contents of this page vary from book to book. In *VDQI* this page is blank; in *EI* and *VE* this page holds a frontispiece; and in *BE* this page contains three epigraphs.

The design of the front matter in Tufte's books varies slightly from the traditional design of front matter. First, the pages in front matter are traditionally numbered with lowercase roman numerals (*e.g.*, i, ii, iii, iv, ...). Second, the front matter page numbering sequence is usually separate from the main matter page numbering. That is, the page numbers restart at 1 when the main matter begins. In contrast, Tufte has enumerated his pages with arabic numerals that share the same page counting sequence as the main matter.

There are also some variations in design across Tufte's four books. The page opposite the full title page (labeled "frontispiece" in the above table) has different content in each of the books. In *The Visual Display of Quantitative Information*, this page is blank; in *Envisioning Information* and *Visual Explanations*, this page holds a frontispiece; and in *Beautiful Evidence*, this page contains three epigraphs. The dedication appears on page 6 in *VDQI* (opposite the introduction), and is placed on its own spread in the other books. In *VE*, an epigraph shares the spread with the opening page of the introduction.

None of the page numbers (folios) of the front matter are expressed except in *BE*, where the folios start to appear on the dedication page.

THE FULL TITLE PAGE of each of the books varies slightly in design. In all the books, the author's name appears at the top of the page, the title is set just above the center line, and the publisher is printed along the bottom margin. Some of the differences are outlined in the following table.

## On the side note of...

yeah

Non-Examinable 9.0.3

Feature	<i>VDQI</i>	<i>EI</i>	<i>VE</i>	<i>BE</i>
Author				
Typeface	serif	serif	serif	sans serif
Style	italics	italics	italics	upright, caps
Size	24 pt	20 pt	20 pt	20 pt
Title				
Typeface	serif	serif	serif	sans serif
Style	upright	italics	upright	upright, caps
Size	36 pt	48 pt	48 pt	36 pt
Subtitle				
Typeface	–	–	serif	–
Style	–	–	upright	–
Size	–	–	20 pt	–
Edition				
Typeface	sans serif	–	–	–
Style	upright, caps	–	–	–
Size	14 pt	–	–	–
Publisher				
Typeface	serif	serif	serif	sans serif
Style	italics	italics	italics	upright, caps
Size	14 pt	14 pt	14 pt	14 pt

THE TABLES OF CONTENTS in Tufte’s books give us our first glimpse of the structure of the main matter. *The Visual Display of Quantitative Information* is split into two parts, each containing some number of chapters. His other three books only contain chapters—they’re not broken into parts.

## 9.1 Headings

Tufte’s books include the following heading levels: parts, chapters,<sup>2</sup> sections, subsections, and paragraphs. Not defined by default are: sub-subsections and subparagraphs.

<sup>2</sup> Parts and chapters are defined for the `tufte-book` class only.

Heading	Style	Size
Part	roman	24/36×40 pc
Chapter	italic	20/30×40 pc
Section	italic	12/16×26 pc
Subsection	italic	11/15×26 pc
Paragraph	italic	10/14

Table 9.1: Heading styles used in *Beautiful Evidence*.

*Paragraph* Paragraph headings (as shown here) are introduced by italicized text and separated from the main paragraph by a bit of space.

## 9.2 Environments

The following characteristics define the various environments:

Environment	Font size	Notes
Body text	10/14×26 pc	
Block quote	9/12×24 pc	Block indent (left and right) by 1 pc
Sidenotes	8/10×12 pc	Sidenote number is set inline, followed by word space
Captions	8/10×12 pc	

Table 9.2: Environment styles used in *Beautiful Evidence*.

Column 1	Column 2	Column 3
Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris.	Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris.	Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris.
Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris.	Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris.	Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris.
Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris.	Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris.	Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris.

Table 9.3: Example table with limited column widths

# 10

## On the Use of the `tufte-book` Document Class

The Tufte- $\text{\LaTeX}$  document classes define a style similar to the style Edward Tufte uses in his books and handouts. Tufte’s style is known for its extensive use of sidenotes, tight integration of graphics with text, and well-set typography. This document aims to be at once a demonstration of the features of the Tufte- $\text{\LaTeX}$  document classes and a style guide to their use.

### 10.1 *Page Layout*

#### 10.1.1 *Headings*

This style provides A- and B-heads (that is, `\section` and `\subsection`), demonstrated above.

If you need more than two levels of section headings, you’ll have to define them yourself at the moment; there are no pre-defined styles for anything below a `\subsection`. As Bringhurst points out in *The Elements of Typographic Style*,<sup>1</sup> you should “use as many levels of headings as you need: no more, and no fewer.” The Tufte- $\text{\LaTeX}$  classes will emit an error if you try to use `\subsubsection` and smaller headings.

<sup>1</sup> Bringhurst2005

IN HIS LATER BOOKS,<sup>2</sup> Tufte starts each section with a bit of vertical space, a non-indented paragraph, and sets the first few words of the sentence in SMALL CAPS. To accomplish this using this style, use the `\newthought` command:

<sup>2</sup> Tufte2006

```
\newthought{In his later books}, Tufte starts...
```

### 10.2 *Sidenotes*

One of the most prominent and distinctive features of this style is the extensive use of sidenotes. There is a wide margin to provide ample room for sidenotes and small figures. Any `\footnotes` will automatically be converted to sidenotes.<sup>3</sup> If you’d like to place ancillary information in the margin without the sidenote mark (the superscript number), you can use the `\marginnote` command. The specification of the `\sidenote` command is:

<sup>3</sup> This is a sidenote that was entered using the `\footnote` command.

This is a margin note. Notice that there isn’t a number preceding the note, and there is no number in the main text where this note was written.

```
\sidenote[⟨number⟩][⟨offset⟩]{Sidenote text.}
```

Both the `⟨number⟩` and `⟨offset⟩` arguments are optional. If you provide a `⟨number⟩` argument, then that number will be used as the sidenote number. It will change of the number of the current sidenote only and will not affect the numbering sequence of subsequent sidenotes.

Sometimes a sidenote may run over the top of other text or graphics in the margin space. If this happens, you can adjust the vertical position of the sidenote by providing a dimension in the `⟨offset⟩` argument. Some examples of valid dimensions are:

```
1.0in    2.54cm    254mm    6\baselineskip
```

If the dimension is positive it will push the sidenote down the page; if the dimension is negative, it will move the sidenote up the page.



While both the  $\langle number \rangle$  and  $\langle offset \rangle$  arguments are optional, they must be provided in order. To adjust the vertical position of the sidenote while leaving the sidenote number alone, use the following syntax:

```
\sidenote[][\langle offset \rangle]{Sidenote text.}
```

The empty brackets tell the `\sidenote` command to use the default sidenote number.

If you *only* want to change the sidenote number, however, you may completely omit the  $\langle offset \rangle$  argument:

```
\sidenote[\langle number \rangle]{Sidenote text.}
```

The `\marginnote` command has a similar *offset* argument:

```
\marginnote[\langle offset \rangle]{Margin note text.}
```

## 10.3 References

References are placed alongside their citations as sidenotes, as well. This can be accomplished using the normal `\cite` command.<sup>4</sup>

The complete list of references may also be printed automatically by using the `\bibliography` command. (See the end of this document for an example.) If you do not want to print a bibliography at the end of your document, use the `\nobibliography` command in its place.

To enter multiple citations at one location,<sup>5</sup> you can provide a list of keys separated by commas and the same optional vertical offset argument:

```
\cite{Tufte2006,Tufte1990}.
```

```
\cite[\langle offset \rangle]{bibkey1,bibkey2,...}
```

<sup>4</sup> The first paragraph of this document includes a citation.

<sup>5</sup> **Tufte2006, Tufte1990**

## 10.4 Figures and Tables

Images and graphics play an integral role in Tufte’s work. In addition to the standard `figure` and `tabular` environments, this style provides special figure and table environments for full-width floats.

Full page-width figures and tables may be placed in `figure*` or `table*` environments. To place figures or tables in the margin, use the `marginfigure` or `marginfigure` environments as follows (see figure 10.1):

```
\begin{marginfigure}
  \includegraphics{helix}
  \caption{This is a margin figure.}
  \label{fig:marginfig}
\end{marginfigure}
```

The `marginfigure` and `marginfigure` environments accept an optional parameter  $\langle offset \rangle$  that adjusts the vertical position of the figure or table. See the “Sidenotes” section above for examples. The specifications are:

```
\begin{marginfigure}[\langle offset \rangle]
  ...
\end{marginfigure}

\begin{marginfigure}[\langle offset \rangle]
  ...
\end{marginfigure}
```

Figure ?? is an example of the `figure*` environment and figure ?? is an example of the normal `figure` environment.

Figure 10.1: This is a margin figure. The helix is defined by  $x = \cos(2\pi z)$ ,  $y = \sin(2\pi z)$ , and  $z = [0, 2.7]$ . The figure was drawn using Asymptote (<http://asymptote.sf.net/>).

As with sidenotes and marginnotes, a caption may sometimes require vertical adjustment. The `\caption` command now takes a second optional argument that enables you to do this by providing a dimension  $\langle offset \rangle$ . You may specify the caption in any one of the following forms:

```
\caption{long caption}
\caption[short caption]{long caption}
\caption[][\langle offset \rangle]{long caption}
\caption[short caption][\langle offset \rangle]{long caption}
```

A positive  $\langle offset \rangle$  will push the caption down the page. The short caption, if provided, is what appears in the list of figures/tables, otherwise the “long” caption appears there. Note that although the arguments  $\langle short\ caption \rangle$  and  $\langle offset \rangle$  are both optional, they must be provided in order. Thus, to specify an  $\langle offset \rangle$  without specifying a  $\langle short\ caption \rangle$ , you must include the first set of empty brackets `[]`, which tell `\caption` to use the default “long” caption. As an example, the caption to figure ?? above was given in the form

```
\caption[Hilbert curves...][6pt]{Hilbert curves...}
```

Table 10.1 shows table created with the `booktabs` package. Notice the lack of vertical rules—they serve only to clutter the table’s data.

Margin	Length
Paper width	8 <sup>1</sup> / <sub>2</sub> inches
Paper height	11 inches
Textblock width	6 <sup>1</sup> / <sub>2</sub> inches
Textblock/sidenote gutter	3/ <sub>8</sub> inches
Sidenote width	2 inches

Table 10.1: Here are the dimensions of the various margins used in the Tufte-handout class.

OCCASIONALLY L<sup>A</sup>T<sub>E</sub>X will generate an error message:

```
Error: Too many unprocessed floats
```

L<sup>A</sup>T<sub>E</sub>X tries to place floats in the best position on the page. Until it’s finished composing the page, however, it won’t know where those positions are. If you have a lot of floats on a page (including sidenotes, margin notes, figures, tables, etc.), L<sup>A</sup>T<sub>E</sub>X may run out of “slots” to keep track of them and will generate the above error.

L<sup>A</sup>T<sub>E</sub>X initially allocates 18 slots for storing floats. To work around this limitation, the Tufte-L<sup>A</sup>T<sub>E</sub>X document classes provide a `\morefloats` command that will reserve more slots.

The first time `\morefloats` is called, it allocates an additional 34 slots. The second time `\morefloats` is called, it allocates another 26 slots.

The `\morefloats` command may only be used two times. Calling it a third time will generate an error message. (This is because we can’t safely allocate many more floats or L<sup>A</sup>T<sub>E</sub>X will run out of memory.)

If, after using the `\morefloats` command twice, you continue to get the `Too many unprocessed floats` error, there are a couple things you can do.

The `\FloatBarrier` command will immediately process all the floats before typesetting more material. Since `\FloatBarrier` will start a new paragraph, you should place this command at the beginning or end of a paragraph.

The `\clearpage` command will also process the floats before continuing, but instead of starting a new paragraph, it will start a new page.

You can also try moving your floats around a bit: move a figure or table to the next page or reduce the number of sidenotes. (Each sidenote actually uses *two* slots.)

After the floats have placed, L<sup>A</sup>T<sub>E</sub>X will mark those slots as unused so they are available for the next page to be composed.

## 10.5 Captions

You may notice that the captions are sometimes misaligned. Due to the way L<sup>A</sup>T<sub>E</sub>X’s float mechanism works, we can’t know for sure where it decided to put a float. Therefore, the Tufte-L<sup>A</sup>T<sub>E</sub>X document classes provide commands to override the caption position.

*Vertical alignment* To override the vertical alignment, use the `\setfloatalignment` command inside the float environment. For example:

```
\begin{figure}[btp]
  \includegraphics{sinewave}
  \caption{This is an example of a sine wave.}
  \label{fig:sinewave}
  \setfloatalignment{b}% forces caption to be bottom-aligned
\end{figure}
```

The syntax of the `\setfloatalignment` command is:

```
\setfloatalignment{<pos>}
```

where *<pos>* can be either `b` for bottom-aligned captions, or `t` for top-aligned captions.

*Horizontal alignment* To override the horizontal alignment, use either the `\forceversofloat` or the `\forcerectofloat` command inside of the float environment. For example:

```
\begin{figure}[btp]
  \includegraphics{sinewave}
  \caption{This is an example of a sine wave.}
  \label{fig:sinewave}
  \forceversofloat% forces caption to be set to the left of the float
\end{figure}
```

The `\forceversofloat` command causes the algorithm to assume the float has been placed on a verso page—that is, a page on the left side of a two-page spread. Conversely, the `\forcerectofloat` command causes the algorithm to assume the float has been placed on a recto page—that is, a page on the right side of a two-page spread.

## 10.6 Full-width text blocks

In addition to the new float types, there is a `fullwidth` environment that stretches across the main text block and the sidenotes area.

```
\begin{fullwidth}
  Lorem ipsum dolor sit amet...
\end{fullwidth}
```

*Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.*

## 10.7 *Typography*

### 10.7.1 *Typefaces*

If the Palatino, Helvetica, and Bera Mono typefaces are installed, this style will use them automatically. Otherwise, we'll fall back on the Computer Modern typefaces.

### 10.7.2 *Letterspacing*

This document class includes two new commands and some improvements on existing commands for letterspacing.

When setting strings of ALL CAPS or SMALL CAPS, the letterspacing—that is, the spacing between the letters—should be increased slightly.<sup>6</sup> The `\allcaps` command has proper letterspacing for strings of FULL CAPITAL LETTERS, and the `\smallcaps` command has letterspacing for SMALL CAPITAL LETTERS. These commands will also automatically convert the case of the text to upper- or lowercase, respectively.

<sup>6</sup>Bringhurst2005

The `\textsc` command has also been redefined to include letterspacing. The case of the `\textsc` argument is left as is, however. This allows one to use both uppercase and lowercase letters: THE INITIAL LETTERS OF THE WORDS IN THIS SENTENCE ARE CAPITALIZED.

## 10.8 *Document Class Options*

The `tufte-book` class is based on the L<sup>A</sup>T<sub>E</sub>X `book` document class. Therefore, you can pass any of the typical book options. There are a few options that are specific to the `tufte-book` document class, however.

The `a4paper` option will set the paper size to A4 instead of the default US letter size.

The `sfsidenotes` option will set the sidenotes and title block in a sans serif typeface instead of the default roman.

The `twoside` option will modify the running heads so that the page number is printed on the outside edge (as opposed to always printing the page number on the right-side edge in `oneside` mode).

The `symmetric` option typesets the sidenotes on the outside edge of the page. This is how books are traditionally printed, but is contrary to Tufte's book design which sets the sidenotes on the right side of the page. This option implicitly sets the `twoside` option.

The `justified` option sets `alldocclsoptdef` and `right`). The default is to set the text ragged right. The body text of Tufte's books are set ragged right. This prevents needless hyphenation and makes it easier to read the text in the slightly narrower column.

The `bidirectional` option loads the  `bidi` package which is used with X<sub>Y</sub>L<sup>A</sup>T<sub>E</sub>X to typeset bi-directional text. Since the  `bidi` package needs to be loaded before the sidenotes and cite commands are defined, it can't be loaded in the document preamble.

The `debug` option causes the Tufte-L<sup>A</sup>T<sub>E</sub>X classes to output debug information to the log file which is useful in troubleshooting bugs. It will also cause the graphics to be replaced by outlines.

The `nofonts` option prevents the Tufte-L<sup>A</sup>T<sub>E</sub>X classes from automatically loading the Palatino and Helvetica typefaces. You should use this option if you wish to load your own fonts. If you're using X<sub>Y</sub>L<sup>A</sup>T<sub>E</sub>X, this option is implied (*i.e.*, the Palatino and Helvetica fonts aren't loaded if you use X<sub>Y</sub>L<sup>A</sup>T<sub>E</sub>X).

The `nols` option inhibits the letterspacing code. The Tufte-L<sup>A</sup>T<sub>E</sub>X classes try to load the appropriate letterspacing package (either pdf<sub>T</sub>E<sub>X</sub>'s `letterspace` package

or the `soul` package). If you're using  $\text{\LaTeX}$  with `fontenc`, however, you should configure your own letterspacing.

The `notitlepage` option causes `\maketitle` to generate a title block instead of a title page. The `book` class defaults to a title page and the `handout` class defaults to the title block. There is an analogous `titlepage` option that forces `\maketitle` to generate a full title page instead of the title block.

The `notoc` option suppresses Tufte- $\text{\LaTeX}$ 's custom table of contents (TOC) design.

The current TOC design only shows unnumbered chapter titles; it doesn't show sections or subsections. The `notoc` option will revert to  $\text{\LaTeX}$ 's TOC design.

The `nohyper` option prevents the `hyperref` package from being loaded. The default is to load the `hyperref` package and use the `\title` and `\author` contents as metadata for the generated PDF.

# 11

## Customizing Tufte- $\text{\LaTeX}$

The Tufte- $\text{\LaTeX}$  document classes are designed to closely emulate Tufte’s book design by default. However, each document is different and you may encounter situations where the default settings are insufficient. This chapter explores many of the ways you can adjust the Tufte- $\text{\LaTeX}$  document classes to better fit your needs.

### 11.1 *File Hooks*

If you create many documents using the Tufte- $\text{\LaTeX}$  classes, it’s easier to store your customizations in a separate file instead of copying them into the preamble of each document. The Tufte- $\text{\LaTeX}$  classes provide three file hooks:

`tufte-common-local.tex`, `tufte-book-local.tex`, and  
`tufte-handout-local.tex`.

*tufte-common-local.tex* If this file exists, it will be loaded by all of the Tufte- $\text{\LaTeX}$  document classes just prior to any document-class-specific code. If your customizations or code should be included in both the book and handout classes, use this file hook.

*tufte-book-local.tex* If this file exists, it will be loaded after all of the common and book-specific code has been read. If your customizations apply only to the book class, use this file hook.

*tufte-common-handout.tex* If this file exists, it will be loaded after all of the common and handout-specific code has been read. If your customizations apply only to the handout class, use this file hook.

### 11.2 *Numbered Section Headings*

While Tufte dispenses with numbered headings in his books, if you require them, they can be enabled by changing the value of the `secnumdepth` counter. From the table below, select the heading level at which numbering should stop and set the `secnumdepth` counter to that value. For example, if you want parts and chapters numbered, but don’t want numbering for sections or subsections, use the command:

```
\setcounter{secnumdepth}{0}
```

The default `secnumdepth` for the Tufte- $\text{\LaTeX}$  document classes is `-1`.

Heading level	Value
Part (in <code>tufte-book</code> )	<code>-1</code>
Part (in <code>tufte-handout</code> )	<code>0</code>
Chapter (only in <code>tufte-book</code> )	<code>0</code>
Section	<code>1</code>
Subsection	<code>2</code>
Subsubsection	<code>3</code>
Paragraph	<code>4</code>
Subparagraph	<code>5</code>

Table 11.1: Heading levels used with the `secnumdepth` counter.

### 11.3 Changing the Paper Size

The Tufte-L<sup>A</sup>T<sub>E</sub>X classes currently only provide three paper sizes: A4, B5, and US letter. To specify a different paper size (and/or margins), use the `\geometrysetup` command in the preamble of your document (or one of the file hooks). The full documentation of the `\geometrysetup` command may be found in the `geometry` package documentation.<sup>1</sup>

<sup>1</sup> `pkg-geometry`

### 11.4 Customizing Marginal Material

Marginal material includes sidenotes, citations, margin notes, and captions.

Normally, the justification of the marginal material follows the justification of the body text. If you specify the `justified` document class option, all of the margin material will be fully justified as well. If you don't specify the `justified` option, then the marginal material will be set ragged right.

You can set the justification of the marginal material separately from the body text using the following document class options: `sidenote`, `marginnote`, `caption`, `citation`, and `marginals`. Each option refers to its obviously corresponding marginal material type. The `marginals` option simultaneously sets the justification on all four marginal material types.

Each of the document class options takes one of five justification types:

*justified* Fully justifies the text (sets it flush left and right).

*raggedleft* Sets the text ragged left, regardless of which page it falls on.

*raggedright* Sets the text ragged right, regardless of which page it falls on.

*raggedouter* Sets the text ragged left if it falls on the left-hand (verso) page of the spread and otherwise sets it ragged right. This is useful in conjunction with the `symmetric` document class option.

*auto* If the `justified` document class option was specified, then set the text fully justified; otherwise the text is set ragged right. This is the default justification option if one is not explicitly specified.

For example,

```
\documentclass[symmetric,justified,marginals=raggedouter]{tufte-book}
```

will set the body text of the document to be fully justified and all of the margin material (sidenotes, margin notes, captions, and citations) to be flush against the body text with ragged outer edges.

THE FONT AND STYLE of the marginal material may also be modified using the following commands:

```
\setsidenotefont{\font commands}
\setcaptionfont{\font commands}
\setmarginnotefont{\font commands}
\setcitationfont{\font commands}
```

The `\setsidenotefont` sets the font and style for sidenotes, the `\setcaptionfont` for captions, the `\setmarginnotefont` for margin notes, and the `\setcitationfont` for citations. The `\font commands` can contain font size changes (e.g., `\footnotesize`, `\Huge`, etc.), font style changes (e.g., `\sffamily`, `\ttfamily`, `\itshape`, etc.), color changes (e.g., `\color{blue}`), and many other adjustments.

If, for example, you wanted the captions to be set in italic sans serif, you could use:

```
\setcaptionfont{\itshape\sffamily}
```



# 12

## Compatibility Issues

When switching an existing document from one document class to a Tufte- $\text{\LaTeX}$  document class, a few changes to the document may have to be made.

### *12.1 Converting from article to tufte-handout*

The following `article` class options are unsupported: `10pt`, `11pt`, `12pt`, `a5paper`, `b5paper`, `executivepaper`, `legalpaper`, `landscape`, `onecolumn`, and `twocolumn`.  
The following headings are not supported: `\subsubsection` and `\subparagraph`.

### *12.2 Converting from book to tufte-book*

The following `report` class options are unsupported: `10pt`, `11pt`, `12pt`, `a5paper`, `b5paper`, `executivepaper`, `legalpaper`, `landscape`, `onecolumn`, and `twocolumn`.  
The following headings are not supported: `\subsubsection` and `\subparagraph`.

# 13

## Troubleshooting and Support

### 13.1 *Tufte- $\LaTeX$ Website*

The website for the Tufte- $\LaTeX$  packages is located at <http://code.google.com/p/tufte-latex/>. On our website, you'll find links to our SVN repository, mailing lists, bug tracker, and documentation.

### 13.2 *Tufte- $\LaTeX$ Mailing Lists*

There are two mailing lists for the Tufte- $\LaTeX$  project:

*Discussion list* The `tufte-latex` discussion list is for asking questions, getting assistance with problems, and help with troubleshooting. Release announcements are also posted to this list. You can subscribe to the `tufte-latex` discussion list at <http://groups.google.com/group/tufte-latex>.

*Commits list* The `tufte-latex-commits` list is a read-only mailing list. A message is sent to the list any time the Tufte- $\LaTeX$  code has been updated. If you'd like to keep up with the latest code developments, you may subscribe to this list. You can subscribe to the `tufte-latex-commits` mailing list at <http://groups.google.com/group/tufte-latex-commits>.

### 13.3 *Getting Help*

If you've encountered a problem with one of the Tufte- $\LaTeX$  document classes, have a question, or would like to report a bug, please send an email to our mailing list or visit our website.

To help us troubleshoot the problem more quickly, please try to compile your document using the `debug` class option and send the generated `.log` file to the mailing list with a brief description of the problem.

### 13.4 *Errors, Warnings, and Informational Messages*

The following is a list of all of the errors, warnings, and other messages generated by the Tufte- $\LaTeX$  classes and a brief description of their meanings.

Error: `\subparagraph` is undefined by this class.

The `\subparagraph` command is not defined in the Tufte- $\LaTeX$  document classes. If you'd like to use the `\subparagraph` command, you'll need to redefine it yourself. See the "Headings" section on page 80 for a description of the heading styles available in the Tufte- $\LaTeX$  document classes.

Error: `\subsubsection` is undefined by this class.

The `\subsubsection` command is not defined in the Tufte- $\LaTeX$  document classes. If you'd like to use the `\subsubsection` command, you'll need to redefine

it yourself. See the “Headings” section on page 80 for a description of the heading styles available in the Tufte-L<sup>A</sup>T<sub>E</sub>X document classes.

**Error:** You may only call `\morefloats` twice. See the Tufte-LaTeX documentation for other workarounds.

L<sup>A</sup>T<sub>E</sub>X allocates 18 slots for storing floats. The first time `\morefloats` is called, it allocates an additional 34 slots. The second time `\morefloats` is called, it allocates another 26 slots.

The `\morefloats` command may only be called two times. Calling it a third time will generate this error message. See page 82 for more information.

**Warning:** Option ‘*⟨class option⟩*’ is not supported -- ignoring option.

This warning appears when you’ve tried to use *⟨class option⟩* with a Tufte-L<sup>A</sup>T<sub>E</sub>X document class, but *⟨class option⟩* isn’t supported by the Tufte-L<sup>A</sup>T<sub>E</sub>X document class. In this situation, *⟨class option⟩* is ignored.

**Info:** The ‘*symmetric*’ option implies ‘*twoside*’

You specified the `symmetric` document class option. This option automatically forces the `twoside` option as well. See page 84 for more information on the `symmetric` class option.

## 13.5 Package Dependencies

The following is a list of packages that the Tufte-L<sup>A</sup>T<sub>E</sub>X document classes rely upon. Packages marked with an asterisk are optional.

- xifthen
- natbib *and* bibentry
- ifpdf\*
- optparams
- ifxetex\*
- placeins
- hyperref
- mathpazo\*
- geometry
- helvet\*
- ragged2e
- fontenc
- chngpage *or* changepage
- beramono\*
- paralist
- fancyhdr
- textcase
- xcolor
- soul\*
- textcomp
- letterspace\*
- titlesec
- setspace
- titletoc