# 1 Introduction and Motivation

## 1.1 Finding Words for Intuitions

Predictors、

training、

Type checking allows the reader to sanity 5check whether the equation that they are considering contains inputs and outputs of the correct type, and whether they are mixing different types of objects.

data ：we assume that the data has already been data appropriately converted into a numerical representation。

Models are simplified versions of reality, which capture model aspects of the real world that are relevant to the task。
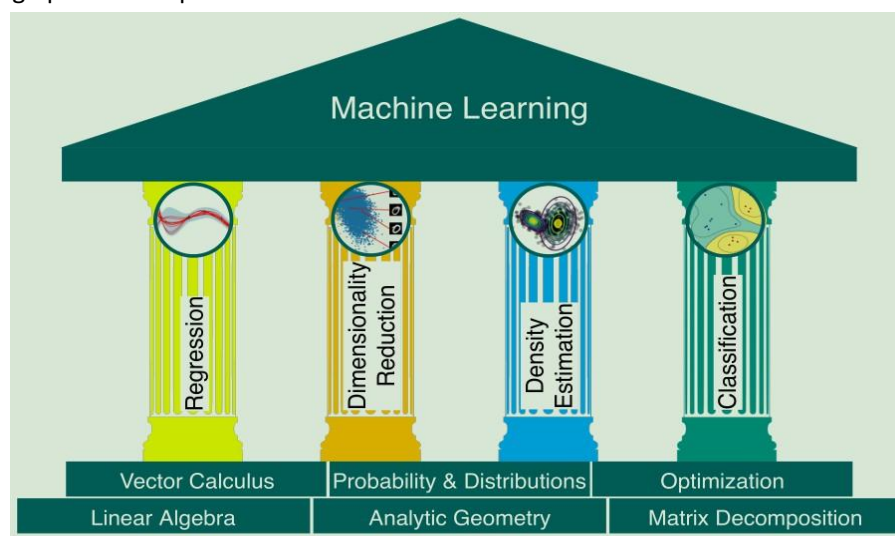
Let us summarize the main concepts of machine learning:
- We use domain knowledge to represent data as vectors.
- We choose an appropriate model, either using the probabilisitic or optimization view.
- We learn from past data by using numerical optimization methods with the aim that it performs well on unseen data.

## 1.2 Two Ways to Read this Book

Two strategies for understanding the mathematics for machine learning:
- Building up the concepts from foundational to more advanced.



- Drilling down from practical needs to more basic requirements.

Part I is about Mathematics

Chapter 2 linear algebra

Chapter 3 analytic geometry

Chapter 4 matrix decomposition

Chapter 5 calculus

Chapter 6 probability theory

Chapter 7 optimization

Part II is about Machine Learning

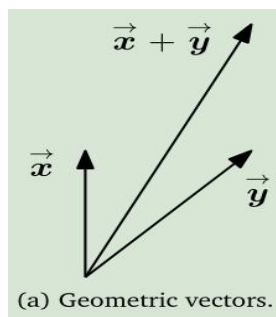|  | Supervised | Unsupervised |
|---|---|---|
| Continuous latent variables | Regression (Chapter 9) | Dimensionality reduction (Chapter 10) |
| Categorical latent variables | Classification (Chapter 12) | Density estimation (Chapter 11) |

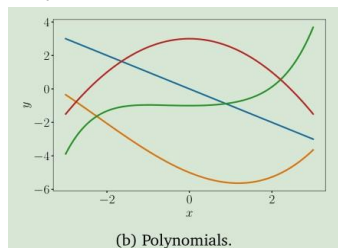1.3  Exercises and Feedback

# 2 Linear Algebra

In general, vectors are special objects that can be added together and multiplied by scalars to produce another object of the same kind. Any object that satisfies these two properties can be considered a vector.

some examples of such vector objects:

Geometric vectors



(a) Geometric vectors.

Polynomials are also vectors



(b) Polynomials.

两个重要网站

https://www.youtube.com/playlist?list=PLlXfTHzgMRUKXD88IdzS14F4NxAZudSmv

https://ocw.mit.edu/courses/mathematics/18-06-linear-algebra-spring-2010/video-lectures/

**2.1 Systems of Linear Equations**

**2.2 Matrices**

**2.2.1 Matrix Addition and Multiplication**

**2.2.2 Inverse and Transpose**

**2.2.3 Multiplication by a Scalar**

- Distributivity:
$$(\lambda + \psi)C = \lambda C + \psi C, \quad C \in \mathbb{R}^{m \times n}$$
$$\lambda(B + C) = \lambda B + \lambda C, \quad B, C \in \mathbb{R}^{m \times n}$$
- Associativity:
$$(\lambda \psi)C = \lambda(\psi C), \quad C \in \mathbb{R}^{m \times n}$$
$$\lambda(BC) = (\lambda B)C = B(\lambda C) = (BC)\lambda, \quad B \in \mathbb{R}^{m \times n}, C \in \mathbb{R}^{n \times k}.$$
Note that this allows us to move scalar values around.
- $(\lambda C)^\top = C^\top \lambda^\top = C^\top \lambda = \lambda C^\top$ since $\lambda = \lambda^\top$ for all $\lambda \in \mathbb{R}$.

2.2.4 Compact Representations of Systems of Linear Equations

2.3 Solving Systems of Linear Equations

2.3.1 Particular and General Solution

2.3.2 Elementary Transformations

2.3.3 The Minus-1 Trick

2.3.4 Algorithms for Solving a System of Linear Equations

$$Ax = b \iff A^\top A x = A^\top b \iff x = (A^\top A)^{-1} A^\top b$$

2.4 Vector Spaces

2.4.1 Groups

Groups play an important role in computer science. Besides providing a fundamental framework for operations on sets, they are heavily used in cryptography, coding theory and graphics.

**Definition 2.6** (Group). Consider a set $\mathcal{G}$ and an operation $\otimes : \mathcal{G} \times \mathcal{G} \to \mathcal{G}$ defined on $\mathcal{G}$.
  Then $G := (\mathcal{G}, \otimes)$ is called a *group* if the following hold:

1. *Closure* of $\mathcal{G}$ under $\otimes$: $\forall x, y \in \mathcal{G} : x \otimes y \in \mathcal{G}$
2. *Associativity*: $\forall x, y, z \in \mathcal{G} : (x \otimes y) \otimes z = x \otimes (y \otimes z)$
3. *Neutral element*: $\exists e \in \mathcal{G} \, \forall x \in \mathcal{G} : x \otimes e = x$ and $e \otimes x = x$
4. *Inverse element*: $\forall x \in \mathcal{G} \, \exists y \in \mathcal{G} : x \otimes y = e$ and $y \otimes x = e$. We often write $x^{-1}$ to denote the inverse element of $x$.

If additionally $\forall x, y \in \mathcal{G} : x \otimes y = y \otimes x$ then $G = (\mathcal{G}, \otimes)$ is an *Abelian group* (commutative).

2.4.2 Vector Spaces

The mapping

$$\Phi : \mathbb{R}^4 \to \mathbb{R}^2, \quad \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} \mapsto \begin{bmatrix} 1 & 2 & -1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} x_1 + 2x_2 - x_3 \\ x_1 + x_4 \end{bmatrix}$$

(2.124)

$$= x_1 \begin{bmatrix} 1 \\ 1 \end{bmatrix} + x_2 \begin{bmatrix} 2 \\ 0 \end{bmatrix} + x_3 \begin{bmatrix} -1 \\ 0 \end{bmatrix} + x_4 \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad (2.125)$$

is linear. To determine $\mathrm{Im}(\Phi)$ we can take the span of the columns of the transformation matrix and obtain

$$\mathrm{Im}(\Phi) = \mathrm{span}[\begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 2 \\ 0 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}]. \quad (2.126)$$

To compute the kernel (null space) of $\Phi$, we need to solve $\boldsymbol{Ax} = \boldsymbol{0}$, i.e., we need to solve a homogeneous equation system. To do this, we use Gaussian elimination to transform $\boldsymbol{A}$ into reduced row echelon form:

$$\begin{bmatrix} 1 & 2 & -1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix} \rightsquigarrow \cdots \rightsquigarrow \begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & -\frac{1}{2} & -\frac{1}{2} \end{bmatrix}. \quad (2.127)$$

# 3 Analytic Geometry

**Definition 3.1** (Norm). A *norm* on a vector space $V$ is a function

$$\|\cdot\| : V \to \mathbb{R}, \qquad\qquad (3.1)$$
$$x \mapsto \|x\|, \qquad\qquad (3.2)$$

which assigns each vector $x$ its *length* $\|x\| \in \mathbb{R}$, such that for all $\lambda \in \mathbb{R}$ and $x, y \in V$ the following hold:

- Absolutely homogeneous: $\|\lambda x\| = |\lambda| \|x\|$
- *Triangle inequality*: $\|x + y\| \leqslant \|x\| + \|y\|$
- *Positive definite*: $\|x\| \geqslant 0$ and $\|x\| = 0 \iff x = 0$.

Manhattan Norm

**Example 3.1 (Manhattan Norm)**
The *Manhattan norm* on $\mathbb{R}^n$ is defined for $x \in \mathbb{R}^n$ as

$$\|x\|_1 := \sum_{i=1}^{n} |x_i|, \qquad\qquad (3.3)$$

where $|\cdot|$ is the absolute value. The left panel of Figure 3.3 indicates all vectors $x \in \mathbb{R}^2$ with $\|x\|_1 = 1$. The Manhattan norm is also called $\ell_1$ norm.

Euclidean Norm

**Example 3.2 (Euclidean Norm)**
The length of a vector $x \in \mathbb{R}^n$ is given by

$$\|x\|_2 := \sqrt{\sum_{i=1}^{n} x_i^2} = \sqrt{x^\top x}, \qquad\qquad (3.4)$$

which computes the *Euclidean distance* of $x$ from the origin. This norm is called the *Euclidean norm*. The right panel of Figure 3.3 shows all vectors $x \in \mathbb{R}^2$ with $\|x\|_2 = 1$. The Euclidean norm is also called $\ell_2$ norm.

3.2 Inner Products

3.2.1 Dot Product

**3.2.2 General Inner Products**

**3.2.3 Symmetric, Positive Definite Matrices**

3.3 Lengths and Distances

3.4 Angles and Orthogonality

$$-1 \leqslant \frac{\langle x, y \rangle}{\|x\| \|y\|} \leqslant 1.$$

nique $\omega \in [0, \pi]$ with

$$\cos \omega = \frac{\langle x, y \rangle}{\|x\| \|y\|},$$

# 4 Matrix Decompositions

**Theorem 4.17.** *Cholesky Decomposition: A symmetric positive definite matrix $\boldsymbol{A}$ can be factorized into a product $\boldsymbol{A} = \boldsymbol{L}\boldsymbol{L}^\top$, where $\boldsymbol{L}$ is a lower triangular matrix with positive diagonal elements:*

$$
\begin{bmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{bmatrix} = \begin{bmatrix} l_{11} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ l_{n1} & \cdots & l_{nn} \end{bmatrix} \begin{bmatrix} l_{11} & \cdots & l_{n1} \\ \vdots & \ddots & \vdots \\ 0 & \cdots & l_{nn} \end{bmatrix}. \qquad (4.43)
$$

$\boldsymbol{L}$ *is called the Cholesky factor of $\boldsymbol{A}$.*

4.4 Eigendecomposition and Diagonalization

4.5 Singular Value Decomposition

4.6 Matrix Approximation

**4.7 Matrix Phylogeny**

# Vector Calculus

5.1 Differentiation of Univariate Functions

5.1.1 Taylor Series

5.1.2 Differentiation Rules

Product Rule: $\quad (f(x)g(x))' = f'(x)g(x) + f(x)g'(x)$ $\qquad (5.37)$

Quotient Rule: $\quad \left( \dfrac{f(x)}{g(x)} \right)' = \dfrac{f'(x)g(x) - f(x)g'(x)}{(g(x))^2}$ $\qquad (5.38)$

Sum Rule: $\quad (f(x) + g(x))' = f'(x) + g'(x)$ $\qquad (5.39)$

Chain Rule: $\quad \big(g(f(x))\big)' = (g \circ f)'(x) = g'(f(x))f'(x)$ $\qquad (5.40)$

Here, $g \circ f$ is a function composition $x \mapsto f(x) \mapsto g(f(x))$.

5.2 Partial Differentiation and Gradients

5.2.1 Basic Rules of Partial Differentiation

Product Rule: $\quad \dfrac{\partial}{\partial \boldsymbol{x}} (f(\boldsymbol{x})g(\boldsymbol{x})) = \dfrac{\partial f}{\partial \boldsymbol{x}} g(\boldsymbol{x}) + f(\boldsymbol{x}) \dfrac{\partial g}{\partial \boldsymbol{x}}$ $\qquad (5.54)$

Sum Rule: $\quad \dfrac{\partial}{\partial \boldsymbol{x}} (f(\boldsymbol{x}) + g(\boldsymbol{x})) = \dfrac{\partial f}{\partial \boldsymbol{x}} + \dfrac{\partial g}{\partial \boldsymbol{x}}$ $\qquad (5.55)$

Chain Rule: $\quad \dfrac{\partial}{\partial \boldsymbol{x}} (g \circ f)(\boldsymbol{x}) = \dfrac{\partial}{\partial \boldsymbol{x}} (g(f(\boldsymbol{x}))) = \dfrac{\partial g}{\partial f} \dfrac{\partial f}{\partial \boldsymbol{x}}$ $\qquad (5.56)$

5.2.2 Chain Rule

$$
\frac{\mathrm{d}f}{\mathrm{d}t} = \begin{bmatrix} \dfrac{\partial f}{\partial x_1} & \dfrac{\partial f}{\partial x_2} \end{bmatrix} \begin{bmatrix} \dfrac{\partial x_1(t)}{\partial t} \\ \dfrac{\partial x_2(t)}{\partial t} \end{bmatrix} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t}
$$

## 5.3 Gradients of Vector-Valued Functions

For a function $f : \mathbb{R}^n \to \mathbb{R}^m$ and a vector $\boldsymbol{x} = [x_1, \ldots, x_n]^\top \in \mathbb{R}^n$, the corresponding vector of function values is given as

$$f(\boldsymbol{x}) = \begin{bmatrix} f_1(\boldsymbol{x}) \\ \vdots \\ f_m(\boldsymbol{x}) \end{bmatrix} \in \mathbb{R}^m . \tag{5.64}$$

**Definition 5.6** (Jacobian). The collection of all first-order partial derivatives of a vector-valued function $f : \mathbb{R}^n \to \mathbb{R}^m$ is called the *Jacobian*. The Jacobian $\boldsymbol{J}$ is an $m \times n$ matrix, which we define and arrange as follows:

$$\boldsymbol{J} = \nabla_{\boldsymbol{x}} f = \frac{d f(\boldsymbol{x})}{d \boldsymbol{x}} = \begin{bmatrix} \frac{\partial f(\boldsymbol{x})}{\partial x_1} & \cdots & \frac{\partial f(\boldsymbol{x})}{\partial x_n} \end{bmatrix} \tag{5.67}$$

$$= \begin{bmatrix} \frac{\partial f_1(\boldsymbol{x})}{\partial x_1} & \cdots & \frac{\partial f_1(\boldsymbol{x})}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m(\boldsymbol{x})}{\partial x_1} & \cdots & \frac{\partial f_m(\boldsymbol{x})}{\partial x_n} \end{bmatrix} , \tag{5.68}$$

$$\boldsymbol{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} , \quad J(i, j) = \frac{\partial f_i}{\partial x_j} . \tag{5.69}$$

## 5.4 Gradients of Matrices

**Example 5.11 (Gradient of Vectors with Respect to Matrices)**
Let us consider the following example, where

$$f = Ax, \quad f \in \mathbb{R}^M, A \in \mathbb{R}^{M \times N}, x \in \mathbb{R}^N \tag{5.97}$$

and where we seek the gradient $df/dA$. Let us start again by determining the dimension of the gradient as

$$\frac{df}{dA} \in \mathbb{R}^{M \times (M \times N)}. \tag{5.98}$$

By definition, the gradient is the collection of the partial derivatives:

$$\frac{df}{dA} = \begin{bmatrix} \frac{\partial f_1}{\partial A} \\ \vdots \\ \frac{\partial f_M}{\partial A} \end{bmatrix}, \quad \frac{\partial f_i}{\partial A} \in \mathbb{R}^{1 \times (M \times N)}. \tag{5.99}$$

To compute the partial derivatives, it will be helpful to explicitly write out the matrix vector multiplication:

$$f_i = \sum_{j=1}^{N} A_{ij} x_j, \quad i = 1, \ldots, M, \tag{5.100}$$

and the partial derivatives are then given as

$$\frac{\partial f_i}{\partial A_{iq}} = x_q. \tag{5.101}$$

This allows us to compute the partial derivatives of $f_i$ with respect to a row of $A$, which is given as

$$\frac{\partial f_i}{\partial A_{i,:}} = x^\top \in \mathbb{R}^{1 \times 1 \times N}, \tag{5.102}$$

$$\frac{\partial f_i}{\partial A_{k \neq i,:}} = 0^\top \in \mathbb{R}^{1 \times 1 \times N} \tag{5.103}$$

5.5 Useful Identities for Computing Gradients

In the following, we list some useful gradients that are frequently required in a machine learning context (Petersen and Pedersen, 2012):

$$\frac{\partial}{\partial X} f(X)^\top = \left(\frac{\partial f(X)}{\partial X}\right)^\top \tag{5.112}$$

$$\frac{\partial}{\partial X} \mathrm{tr}(f(X)) = \mathrm{tr}\left(\frac{\partial f(X)}{\partial X}\right) \tag{5.113}$$

$$\frac{\partial}{\partial X} \det(f(X)) = \det(f(X))\mathrm{tr}\left(f^{-1}(X)\frac{\partial f(X)}{\partial X}\right) \tag{5.114}$$

$$\frac{\partial}{\partial X} f^{-1}(X) = -f^{-1}(X)\frac{\partial f(X)}{\partial X} f^{-1}(X) \tag{5.115}$$

$$\frac{\partial a^\top X^{-1} b}{\partial X} = -(X^{-1})^\top a b^\top (X^{-1})^\top \tag{5.116}$$

$$\frac{\partial x^\top a}{\partial x} = a^\top \tag{5.117}$$

$$\frac{\partial a^\top x}{\partial x} = a^\top \tag{5.118}$$

$$\frac{\partial a^\top X b}{\partial X} = ab^\top \tag{5.119}$$

$$\frac{\partial x^\top B x}{\partial x} = x^\top(B + B^\top) \tag{5.120}$$

$$\frac{\partial}{\partial s}(x - As)^\top W(x - As) = -2(x - As)^\top W A \quad \text{for symmetric } W \tag{5.121}$$

Here, we use tr as the trace operator (see Definition 4.3) and det is the determinant (see Section 4.1).

5.6 Backpropagation and Automatic Differentiation

5.6.1 Gradients in a Deep Network

5.6.2 Automatic Differentiation

**Example 5.13**
Consider the function

$$f(x) = \sqrt{x^2 + \exp(x^2)} + \cos\left(x^2 + \exp(x^2)\right) \qquad (5.135)$$

from (5.122). If we were to implement a function $f$ on a computer, we would be able to save some computation by using *intermediate variables*:
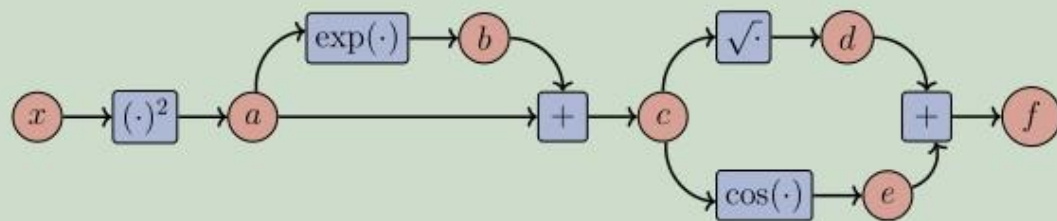
$$a = x^2, \qquad (5.136)$$
$$b = \exp(a), \qquad (5.137)$$
$$c = a + b, \qquad (5.138)$$
$$d = \sqrt{c}, \qquad (5.139)$$
$$e = \cos(c), \qquad (5.140)$$
$$f = d + e. \qquad (5.141)$$

# 6 Probability and Distributions