

Part 1: Linear Algebra

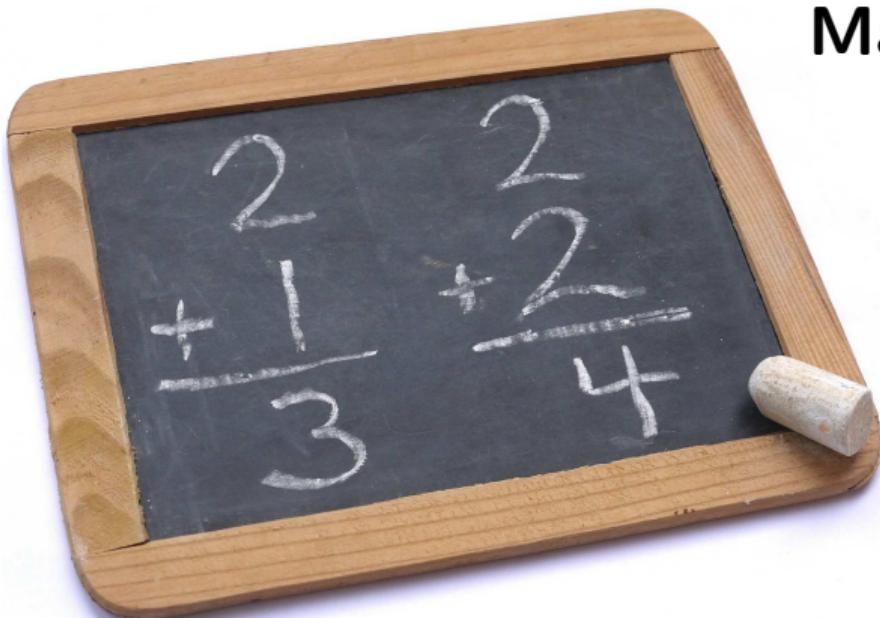
oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooo
ooooooo
ooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
ooooooo



A Refresher On
Machine Learning Maths

CAS Machine Learning
Module 1: Introduction

Michael Calonder
Hochschule Luzern HSLU

Oct 20/21, 2023

Image by Julia S from Pixabay

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
ooooooo

Outline

Part 1: Linear Algebra

Vectors and Matrices

Determinants, Definiteness, Trace, Inverse

Matrix Inversion and Decompositions

Part 2: Calculus

Derivatives: Polynomials, Power Functions, Product and Chain Rule, Multivariate Calculus

Mathematical Optimization and Newton's Method

Stochastic Gradient Descent

Integral Calculus

Part 3: Statistics and Probability

Random Variables

Density and Distribution Functions

Joint, Marginal, and Conditional Densities

Basic Distributions

Linear Algebra

Introduction

What is linear algebra and where is it used?

- Branch of mathematics concerned with linear equations such as

$$a_1x_1 + a_2x_2 + \cdots + a_nx_n = b$$

- Of central importance to most areas of science, mathematics, and engineering: ranging from differential equations, to modeling electrical circuits and control systems that fly aircraft, to econometry and genetics
- Not just “linear”: Non-linear systems may well be treated through a series of linear operations, e.g. first-order Taylor expansion
- At the heart of most machine learning algorithms
- Modern computers can process them very efficiently, still active research (DeepMind’s AlphaTensor)

Part 1: Linear Algebra

○●○○○○○○○○○○
○○○○
○○○○○○○○○○○○

Part 2: Calculus

○○○○○○○○○○○○○○○○
○○○○○○○○○○
○○○○○○○○
○○○○○○○○

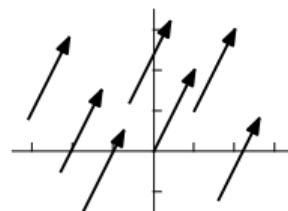
Part 3: Statistics and Probability

○○○○○○
○○○○○○
○○○○○○○○○○○○○○
○○○○○○○○

Vectors

Properties:

- A vector is an object in \mathbb{R}^n that is comprised of a magnitude and a direction
- Often visualized as an “arrow” for $n = 2$ or $n = 3$



Note that all above arrow represent the same vector

- Notation:

$$\mathbf{x} := \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

Part 1: Linear Algebra

○○●○○○○○○○○
○○○○
○○○○○○○○○○○○

Part 2: Calculus

○○○○○○○○○○○○○○○○
○○○○○○○○
○○○○○○
○○○○○○

Part 3: Statistics and Probability

○○○○○○
○○○○○○
○○○○○○○○○○○○○○
○○○○○○○○

Vectors

Linear Dependence

Given x_1, \dots, x_m be vectors in \mathbb{R}^n .

The set $\{x_1, \dots, x_m\}$ of these vectors *linearly independent* if

$$c_1 x_1 + c_2 x_2 + \dots + c_m x_m \quad \text{implies} \quad c_1 = c_2 = \dots = c_m = 0$$

Likewise, any two vectors x_i and x_j are *linearly dependent* if $x_1 = cx_2$ for some real $c \neq 0$

Part 1: Linear Algebra

○○○●○○○○○○○
○○○○
○○○○○○○○○○○

Part 2: Calculus

○○○○○○○○○○○○○○○
○○○○○○○○
○○○○○○
○○○○○○

Part 3: Statistics and Probability

○○○○○○
○○○○○○
○○○○○○○○○○○○○
○○○○○○○○



Determine whether or not the two vectors $(6, 2, 8, 4)^\top$ and $(21, 7, 28, 14)^\top$ are linearly dependent

Answer:

They are linearly dependent, $3.5 \times$ the first vector yields the second

Part 1: Linear Algebra

```
oooo●oooooooooooo
```

Part 2: Calculus

```
oooooooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo
```

Part 3: Statistics and Probability

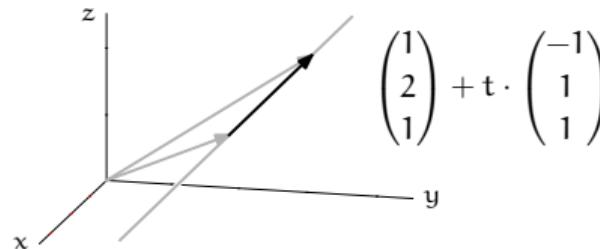
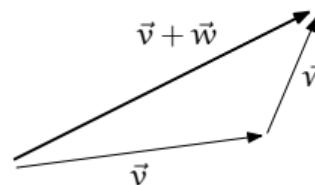
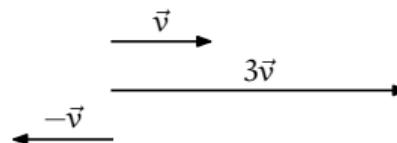
```
oooooooooooo
oooooooooooo
oooooooooooooooooooo
oooooooooooo
```

Vectors

Algebraic properties:

- Additive identity: $\forall x : \exists \mathbf{0}$ s.t. $x + \mathbf{0} = x$
- Multiplicative identity: $1x = x$
- Commutativity: $x + y = y + x$
- Associativity: $(x + y) + z = y + (x + z)$
- Additive inverse: $\forall v : \exists \bar{x} := -x$ s.t.
 $x + \bar{x} = x + (-x) = \mathbf{0}$
- Distributivity:
 - $\forall v, r: r(x + y) = rx + ry$
 - $\forall v, r, s: (r + s)x = rx + sx$

For $x, y, z \in \mathbb{R}^n$ and r, s scalars



Matrices

Basic properties:

- A matrix is an object in $\mathbb{R}^{m \times n}$, i.e. a rectangular array of numbers.
- Central object in linear algebra, enable actual operations, e.g. addition or multiplication
- Widely used, even outside linear algebra e.g. in graph theory (adjacency matrices)
- Notation:

$$\mathbf{A} := \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix}$$

- Note that for now we choose a_{ij} to be in \mathbb{R} , however, $a_{ij} \in \mathbb{C}$ in general
- Note that algebraic properties and identities differ from those above ones for vectors (see below)
- Important special case: *square matrix*, $m = n$

Part 1: Linear Algebra

```
oooooooo●oooooooo
    oooo
    oooooooooooooo
```

Part 2: Calculus

```
oooooooooooooooooooo
    ooooooo
    ooooooo
    ooooooo
```

Part 3: Statistics and Probability

```
ooooooo
    ooooooo
    ooooooooooooooooo
    oooooooo
```

Matrix Multiplication (MM)

Matrix product $\mathbf{C} = \mathbf{AB} \in \mathbb{R}^{m \times p}$ with $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{n \times p}$,

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix}, \quad \mathbf{B} = \begin{pmatrix} b_{11} & b_{12} & \cdots & b_{1p} \\ b_{21} & b_{22} & \cdots & b_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ b_{n1} & b_{n2} & \cdots & b_{np} \end{pmatrix}, \quad \mathbf{C} = \begin{pmatrix} c_{11} & c_{12} & \cdots & c_{1p} \\ c_{21} & c_{22} & \cdots & c_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ c_{m1} & c_{m2} & \cdots & c_{mp} \end{pmatrix}$$

is defined s.t.

$$c_{ij} = a_{i1}b_{1j} + a_{i2}b_{2j} + \cdots + a_{in}b_{nj} = \sum_{k=1}^n a_{ik}b_{kj}$$

For instance,

$$\begin{pmatrix} a_{11} & a_{12} \\ \vdots & \vdots \\ a_{31} & a_{32} \\ \vdots & \vdots \end{pmatrix} \begin{pmatrix} \cdot & b_{12} & b_{13} \\ \cdot & b_{22} & b_{23} \\ \cdot & \cdot & \cdot \end{pmatrix} = \begin{pmatrix} \cdot & c_{12} & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & c_{33} \\ \cdot & \cdot & \cdot \end{pmatrix}$$

MM Is Non-Commutative

MM is undefined if the inner dimensions differ.

But even when the MM is defined, changing the order of the factors will change the result.

For two square matrices $\mathbf{A} \in \mathbb{R}^n$, $\mathbf{B} \in \mathbb{R}^n$, and $n > 1$, in general

$$\mathbf{AB} \neq \mathbf{BA}$$

Example:

$$\begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \quad \neq \quad \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}$$

Part 1: Linear Algebra

oooooooo●ooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooooooo
oooooooooooo
oooooooo
oooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooooooo
ooooooo

MM Is Distributive

Distributivity: MM is distributive with respect to matrix addition. If A, B, C, D are matrices of size $m \times n, n \times p, n \times p, p \times q$, respectively, we have

- left distributivity: $A(B + C) = AB + AC$
- right distributivity: $(B + C)D = BD + CD$

Scalar multiplication:

- For any $c \in \mathbb{R}$ we have $cA = Ac$
- Also, $c(AB) = (cA)B$ and $(AB)c = A(BC)$

Part 1: Linear Algebra

oooooooooooo●○○
○○○○
ooooooooooooooo

Part 2: Calculus

ooooooooooooooooooo
oooooooooooo
oooooooo
oooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
ooooooo

Matrix Transpose

Transposition

For any matrix $A = [a_{ij}]$, its transpose is defined to be

$$A^T = [a_{ji}]$$

Example:

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{pmatrix}^T = \begin{pmatrix} 1 & 3 & 5 \\ 2 & 4 & 6 \end{pmatrix}$$

Matrix Transpose: Identities

Some useful matrix identities related to its transpose:

- $(\mathbf{A}^T)^T = \mathbf{A}$
- $(\mathbf{A} + \mathbf{B})^T = \mathbf{A}^T + \mathbf{B}^T$
- $(\mathbf{A}_1 \mathbf{A}_2 \dots \mathbf{A}_{k-1} \mathbf{A}_k)^T = \mathbf{A}_k^T \mathbf{A}_{k-1}^T \dots \mathbf{A}_2^T \mathbf{A}_1^T$ e.g., $(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$
- $(c\mathbf{A})^T = c\mathbf{A}^T$
- $\det \mathbf{A} = \det \mathbf{A}^T$
- If $\forall i, j : a_{ij} \in \mathbb{R} \Rightarrow \mathbf{A}^T \mathbf{A} \succeq 0$ i.e., \mathbf{A} PSD
- $(\mathbf{A}^{-1})^T = (\mathbf{A}^T)^{-1}$
- If $\mathbf{A} \in \mathbb{R}^{n \times n}$ then eigenvalues unchanged under transposition

Part 1: Linear Algebra

oooooooooooo●
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
ooooooo

Matrix Rank

Various definitions exists, all equivalent.

We define

rank A := number of linearly independent rows (or columns)

Example: The matrix

$$\begin{pmatrix} 1 & 0 & 1 \\ -2 & -3 & 1 \\ 3 & 3 & 0 \end{pmatrix}$$

is of rank 2 (do you see why?)

Determinants

The determinant of a *square* matrix A is scalar that characterizes certain properties of A . Most importantly, the determinant is nonzero iff the matrix is invertible

Due to being expensive to compute, explicit use of determinants in machine learning applications is rare

Computation of the determinant of $A \in \mathbb{R}^n$:

- $n = 1$: $\det a = a$
- $n = 2$: $\det \begin{pmatrix} a & b \\ c & d \end{pmatrix} = ad - bc$
- $n = 3$: $\det \begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix} = aei + bfg + cdh - ceg - bdi - afh$
- $n > 3$: Usually computed via decomposition (see below)

Definiteness

For a *symmetric* matrix $A \in \mathbb{R}^{n \times n}$ and $x \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ consider

$$s := x^\top A x.$$

Then, A is called

- *positive-definite* (PSD) if $s > 0 \forall x$
- *positive-semidefinite* (PD) if $s \geq 0 \forall x$
- *negative-semidefinite* (NSD) if $s \leq 0 \forall x$
- *negative-definite* (ND) if $s < 0 \forall x$
- and *indefinite* otherwise

Example: The identity matrix $I = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ is positive-definite. To see this, let $x = \begin{pmatrix} a \\ b \end{pmatrix}$ and consider

$$x^\top I x = (a \quad b)^\top \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = a^2 + b^2$$

which is strictly positive for all $x \neq \mathbf{0}$.

Part 1: Linear Algebra

○○○○○○○○○○
○○●○
○○○○○○○○○○

Part 2: Calculus

○○○○○○○○○○○○○○
○○○○○○○○
○○○○○○
○○○○○○

Part 3: Statistics and Probability

○○○○○○
○○○○○○
○○○○○○○○○○○○○○
○○○○○○○○

Matrix Trace

The trace of a *square* matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is commonly defined as the sum of its diagonal elements:

$$\text{tr } \mathbf{A} = \sum_{i=1}^n a_{ii} = a_{11} + a_{22} + \cdots + a_{nn}.$$

Example:

$$\text{tr} \begin{pmatrix} 1 & 0 & 3 \\ 11 & 5 & 2 \\ 6 & 12 & -5 \end{pmatrix} = 1 + 5 + (-5) = 1$$

Properties:

- Linearity: $\text{tr}(\mathbf{A} + \mathbf{B}) = \text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B})$ and $\text{tr}(c\mathbf{A}) = c \text{tr}(\mathbf{A})$
- $\text{tr } \mathbf{A} = \text{tr } \mathbf{A}^\top$
- $\text{tr}(\mathbf{ABCD}) = \text{tr}(\mathbf{BCDA}) = \text{tr}(\mathbf{CDAB}) = \text{tr}(\mathbf{DABC})$ and as a special case $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$
- Trace as the sum of eigenvalues λ_i of \mathbf{A} : $\text{tr}(\mathbf{A}) = \sum_{i=1}^n \lambda_i$

Part 1: Linear Algebra

oooooooooooo
ooo●
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
ooooooo

Determine whether or not the real symmetric matrix



$$\mathbf{A} = \begin{pmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{pmatrix}$$

is positive-definite.

Answer:

\mathbf{A} is PD since for any non-zero column vector $\mathbf{z} = (a \ b \ c)^\top$:

$$\mathbf{z}^\top \mathbf{A} \mathbf{z} = (\mathbf{z}^\top \mathbf{M}) \mathbf{z} = ((2a - b) \ (-a + 2b - c) \ (-b + 2c)) \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \dots = a^2 + (a-b)^2 + (b-c)^2 + c^2 > 0$$

Matrix Inverse

The inverse matrix \mathbf{A}^{-1} to a regular matrix \mathbf{A} is defined s.t.

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$$

with \mathbf{I} then identity matrix

Example: The inverse of the real matrix

$$\mathbf{A} = \begin{pmatrix} 2 & 1 \\ 6 & 4 \end{pmatrix} \quad \text{is} \quad \mathbf{A}^{-1} = \begin{pmatrix} 2 & -\frac{1}{2} \\ -3 & 1 \end{pmatrix}$$

Verify this:

$$\mathbf{A}\mathbf{A}^{-1} = \begin{pmatrix} 2 & 1 \\ 6 & 4 \end{pmatrix} \begin{pmatrix} 2 & -\frac{1}{2} \\ -3 & 1 \end{pmatrix} = \begin{pmatrix} 4 - 3 & -1 + 1 \\ 12 - 12 & -3 + 4 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = \mathbf{I}$$

Part 1: Linear Algebra

```
oooooooooooo
ooooo
oooo
o●oooooooooooo
```

Part 2: Calculus

```
oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo
```

Part 3: Statistics and Probability

```
ooooooo
ooooooo
oooooooooooooooooooo
ooooooo
```

Matrix Inverse

When is a matrix \mathbf{A} invertible? Many equivalent conditions exist ("Invertible Matrix Theorem"):

- \mathbf{A} is row- or column-equivalent to the $n \times n$ identity matrix $I_n \Rightarrow$ Gaussian Elimination
- The columns of \mathbf{A} are linearly independent
- \mathbf{A} has full rank, $\text{rank } \mathbf{A} = n$
- 0 is *not* an eigenvalue of $\mathbf{A} \Rightarrow$ Definiteness
- $\det \mathbf{A} \neq 0$

Inversion plays an important part in the solution of linear systems of equations:

$$\mathbf{A}\mathbf{x} = \mathbf{b} \iff \mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$$

Closely related to rank: For instance, if $\text{rank } \mathbf{A} = n$ (full rank), the system

- $\mathbf{A}\mathbf{x} = \mathbf{0}$ has only the trivial solution $\mathbf{x} = \mathbf{0}$
- $\mathbf{A}\mathbf{x} = \mathbf{b}$ has exactly one solution for each $\mathbf{b} \in \mathbb{R}^n$

In practice, computing the inverse is costly and numerically delicate—avoid whenever possible

Part 1: Linear Algebra

○○○○○○○○○○
○○○○
○○●○○○○○○○○

Part 2: Calculus

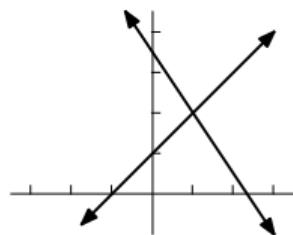
○○○○○○○○○○○○○○
○○○○○○○○
○○○○○○
○○○○○○

Part 3: Statistics and Probability

○○○○○
○○○○○○
○○○○○○○○○○○○
○○○○○○○

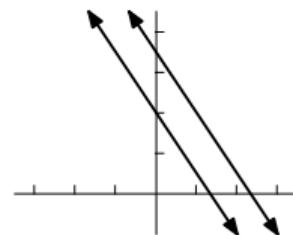
Linear Systems

Unique solution



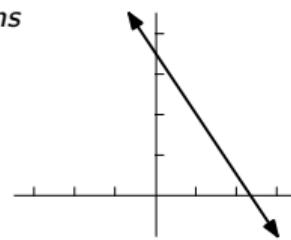
$$\begin{array}{rcl} 3x & + & 2y = 7 \\ x & - & y = -1 \end{array}$$

No solutions



$$\begin{array}{rcl} 3x & + & 2y = 7 \\ 3x & + & 2y = 4 \end{array}$$

Infinitely many solutions



$$\begin{array}{rcl} 3x & + & 2y = 7 \\ 6x & + & 4y = 14 \end{array}$$

Part 1: Linear Algebra

○○○○○○○○○○○○
○○○○
○○○●○○○○○○○○

Part 2: Calculus

○○○○○○○○○○○○○○○○
○○○○○○○○○○
○○○○○○○○
○○○○○○○○

Part 3: Statistics and Probability

○○○○○○
○○○○○○○
○○○○○○○○○○○○○○
○○○○○○○○



Find the coefficients a , b , and c so that the graph of $f(x) = ax^2 + bx + c$ passes through the points $(1, 2)$, $(-1, 6)$, and $(2, 3)$

Answer:

$$a = 1, b = -2, c = 3$$

Part 1: Linear Algebra

oooooooooooo
ooooo
oooo●oooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooooooo
oooooooooooo

Matrix Inverse: Identities

Some useful matrix identities to inversion:

- $(\mathbf{A}^{-1})^{-1} = \mathbf{A}$
- $(k\mathbf{A})^{-1} = k^{-1}\mathbf{A}^{-1}$ for any scalar $k \neq 0$
- $(\mathbf{A}^\top)^{-1} = (\mathbf{A}^{-1})^\top$
- For invertible matrices $\mathbf{A}_i \in \mathbb{R}^{n \times n}$, $(\mathbf{A}_1 \mathbf{A}_2 \cdots \mathbf{A}_{k-1} \mathbf{A}_k)^{-1} = \mathbf{A}_k^{-1} \mathbf{A}_{k-1}^{-1} \cdots \mathbf{A}_2^{-1} \mathbf{A}_1^{-1}$
- $\det \mathbf{A}^{-1} = (\det \mathbf{A})^{-1}$

Part 1: Linear Algebra

oooooooooooo
oooo
oooo●oooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
ooooooo

Matrix Decompositions

Eigendecomposition

The Eigendecomposition factorizes a diagonalizable matrix into its canonical form, represented by *eigenvalues* and *eigenvectors*. Particularly important special case in machine learning: real, symmetric matrices

A vector $x \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ is an eigenvector of a square $N \times N$ matrix A , if it satisfies

$$Ax = \lambda x$$

for some scalar λ . Then λ is called the eigenvalue corresponding to eigenvector x .

Computation

- Characteristic polynomial $p(\lambda) := \det(AI - \lambda I) = 0$ yields up to N distinct eigenvalues λ_i ;
- Each λ_i gives rise to a specific eigenvalue equation $(A - \lambda_i I)v_i = \mathbf{0}$, yielding the eigenvectors v_i

Matrix Decompositions

Lower-upper (LU) Decomposition

The LU decomposition factors a square matrix A into a lower triangular and an upper triangular matrix:

$$A = LU$$

Example: An invertible 3×3 matrix A can be factored as

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} = \begin{pmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{pmatrix} \begin{pmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{pmatrix}.$$

Remarks:

- Typical strategy used by computers to solve square systems of equations
- Often “LU” actually refers to “LUP” where P is a permutation matrix (“partial pivoting”)
- Any square matrix A admits LUP and PLU factorizations. If A is invertible, then it admits LU factorization



Matrix Decompositions

Lower-upper (LU) Decomposition

LU decompositions find wide-spread application from solving linear equations, to inverting a matrix, to computation of the determinant.

Example: Solve a system of linear equations with LUP

Given the system $\mathbf{Ax} = \mathbf{b}$ and the decomposition $\mathbf{PA} = \mathbf{LU}$, rewrite the problem

$$\mathbf{Ax} = \mathbf{b} \iff \mathbf{PAx} = \mathbf{Pb} \iff \mathbf{LUx} = \mathbf{Pb}$$

Then x can be found by

1. first solving $\mathbf{Ly} = \mathbf{Pb}$ for y (forward substitution)
2. second solving $\mathbf{Ux} = y$ for x (backward substitution)

Matrix Decompositions

Cholesky Decomposition

Why Cholesky? As a special case of LU for PD matrices, about twice as efficient to compute

The Cholesky decomposition of an invertible matrix \mathbf{A} is defined as

$$\mathbf{A} = \mathbf{L}\mathbf{L}^{\top}$$

where \mathbf{L} is a real lower triangular matrix with positive diagonal entries

Wide-spread applications:

- Linear least squares: efficient solutions of linear system arising from partial differential equations
- Non-linear optimization: numerical stability in Quasi-Newton by avoiding direct updates of the Hessian
- Monte Carlo simulations: produce correlated samples from uniform ones
- Kalman filters: forcing the system covariance matrix to remain positive semi-definite

Matrix Decompositions

Singular Value Decomposition (SVD)

SVD generalizes the eigendecomposition of a square matrix to any $m \times n$ matrix.

The SVD decomposes a real $m \times n$ matrix \mathbf{A} as follows,

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^\top$$

with \mathbf{U} and \mathbf{V} are real orthogonal matrices and Σ a $m \times n$ rectangular diagonal matrix w

If M is real, then U and V can be guaranteed to be real orthogonal matrices; in such contexts, the SVD is often denoted $\mathbf{U}\Sigma\mathbf{V}^\top$

Part 1: Linear Algebra

○○○○○○○○○○○○
○○○○
○○○○○○○○●○

Part 2: Calculus

○○○○○○○○○○○○○○○○
○○○○○○○○○
○○○○○○○
○○○○○○○

Part 3: Statistics and Probability

○○○○○○
○○○○○○
○○○○○○○○○○○○○○
○○○○○○○○

Paper exercises

- Does $(1, 0, 2, 1)$ lie on the line through $(-2, 1, 1, 0)$ and $(5, 10, -1, 4)$?
- For what value(s) of a does the system have nontrivial solutions?

$$x_1 + 2x_2 + x_3 = 0$$

$$-x_1 - x_2 + x_3 = 0$$

$$3x_1 + 4x_2 + ax_3 = 0$$

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo●

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

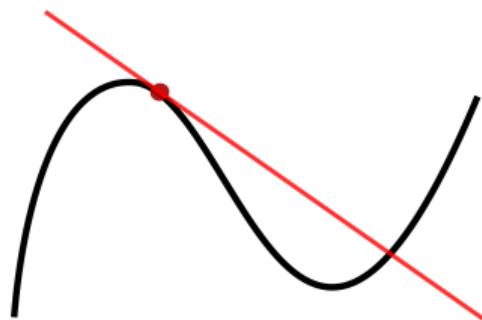
ooooooo
ooooooo
oooooooooooooooo
ooooooo

Exercises with Python

- Implement your own matrix multiplication in pure Python; measure speed for a 1000×1000 matrix; then use numpy for the multiplication and compare
- Compute the SVD of a random $m \times n$ matrix \mathbf{A} using your implementation of choice and numerically compare the reconstructed matrix $\mathbf{U}\Sigma\mathbf{V}^\top$ with \mathbf{A}
- Verify numerically the identity $(\mathbf{A}^{-1})^\top = (\mathbf{A}^\top)^{-1}$
- Implement your own solver for linear systems using LU decomposition; use `scipy.linalg.lu` for the decomposition and implement the foward- and backward substitution yourself
- Given a set of m vectors in \mathbb{R}^n , determine whether it is linearly independent

Intuition

Geometrically: The derivative is the tangent line to the curve at the point (x_0, y_0) is a line passing through (x_0, y_0) and 'flat against' the curve



Intuition: The derivative measures the sensitivity to change of the function value (output value) with respect to a change in its argument (input value)

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

●oooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
ooooooo

Definition

Roughly speaking, a function $f(x)$ is differentiable at a point a of its domain, if the limit

$$L = \lim_{h \rightarrow 0} \frac{f(a + h) - f(a)}{h}$$

exists.

If the limit L exists, then this limit is called the *derivative* of f at a , and denoted

$$f'(a) \quad \text{or} \quad \frac{df}{dx}(a)$$

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oo●oooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
ooooooo

Derivatives of polynomials

There's just four simple facts which suffice to take the derivative of any polynomial plus somewhat more general things.

First, there is the rule for taking the derivative of a *power function* $f(x) = x^n$:

$$\frac{d}{dx} x^n = n x^{n-1}$$

Second, a special case of the rule above, the derivative of any constant c is zero:

$$\frac{d}{dx} c = 0$$

Third, for any function $f(x)$, multiplicative constants do not affect the result,

$$\frac{d}{dx} cf = c \frac{d}{dx} f$$

Fourth, for any two functions $f(x)$ and $g(x)$, the derivative of the sum is the sum of the derivatives:

$$\frac{d}{dx} (f + g) = \frac{d}{dx} f + \frac{d}{dx} g$$

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooo●oooooooooooo
oooooooooooo
oooooooo
oooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
ooooooo

Putting these four things together, we can write general formulas like

$$\frac{d}{dx}(ax^m + bx^n + cx^p) = amx^{m-1} + bn x^{n-1} + cp x^{p-1}$$

and so on, with any number of summands.

As a refresher, here are some examples:

$$\frac{d}{dx}5x^3 = 15x^2$$

$$\frac{d}{dx}(3x^7 + 5x^3 - 11) = 21x^6 + 15x^2$$

$$\frac{d}{dx}(2 - 3x^2 - 2x^3) = -6x - 6x^2$$

$$\frac{d}{dx}(-x^4 + 2x^5 + 1) = -4x^3 + 10x^4$$

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooo●oooooooooooo
oooooooooooo
oooooooo
oooooooo

Part 3: Statistics and Probability

ooooooo
oooooooo
oooooooooooooooooooo
oooooooooooo

More general power functions

Remember some of the other possibilities for the exponential notation x^n . For example

$$x^{\frac{1}{2}} = \sqrt{x}$$

$$x^{-1} = \frac{1}{x}$$

$$x^{-\frac{1}{2}} = \frac{1}{\sqrt{x}}$$

and so on.

Good news: Above rule for taking the derivative of powers of x is still correct here, even for exponents which are negative or fractions or even real numbers:

$$\frac{d}{dx} x^r = r x^{r-1}$$

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooo●oooooooooooo
oooooooooooo
oooooooo
oooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooooooo
ooooooo

Thus, in particular,

$$\frac{d}{dx} \sqrt{x} = \frac{d}{dx} x^{\frac{1}{2}} = \frac{1}{2} x^{-\frac{1}{2}}$$

$$\frac{d}{dx} \frac{1}{x} = \frac{d}{dx} x^{-1} = -1 \cdot x^{-2} = \frac{-1}{x^2}$$

When combined with the sum rule from above, we have the obvious possibilities:

$$\frac{d}{dx} \left(3x^2 - 7\sqrt{x} + \frac{5}{x^2} \right) = \frac{d}{dx} \left(3x^2 - 7x^{\frac{1}{2}} + 5x^{-2} \right) = 6x - \frac{7}{2}x^{-\frac{1}{2}} - 10x^{-3}$$

Expressing square roots, cube roots, inverses, etc., in terms of exponents often simplifies the rearranging of algebraic expressions

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooo●oooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooooooo
ooooooo

Quotient rule

Quotient rule: Given two differentiable functions $f(x)$ and $g(x)$, the derivative of the ratio of the two is given by:

$$\frac{d}{dx} \frac{f}{g} = \frac{f'g - g'f}{g^2}$$

Example: For $f(x) = x - 1$ and $g(x) = x - 2$ we get

$$\begin{aligned}\frac{d}{dx} \frac{x-1}{x-2} &= \frac{(x-1)'(x-2) - (x-1)(x-2)'}{(x-2)^2} = \frac{1 \cdot (x-2) - (x-1) \cdot 1}{(x-2)^2} \\ &= \frac{(x-2) - (x-1)}{(x-2)^2} = \frac{-1}{(x-2)^2}\end{aligned}$$

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooo●oooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
ooooooo

Special case of the quotient rule: Applying the quotient rule to the expression $1/g(x)$ yields

Reciprocal rule:

$$\frac{d}{dx} \frac{1}{g(x)} = -\frac{g'(x)}{g(x)^2}$$

Because

$$\frac{d}{dx} \frac{1}{g(x)} = \frac{0 \cdot g(x) - 1 \cdot g'(x)}{g(x)^2} = \frac{-g'(x)}{g(x)^2}$$

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooo●○oooo
oooooooo
oooooooo
oooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
ooooooo

Product rule

We have:

Product rule:

$$\frac{d}{dx}(fg) = f'g + fg'$$

Note the this rule behaves differently from the simple rule for sums above

Example: Differentiate $f(x) = x^2 \sin x$. We get $f'(x) = 2x \sin x + x^2 \cos x$

Note: “constant multiple rule” is a special case: $(cf(x))' = cf'(x)$ since $\frac{d}{dx}c = 0$

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooo●oooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooooooo
ooooooo

Find the derivative of the sigmoid function



$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

and express it using the sigmoid function itself. Hint: The identity $\frac{a}{1+a} = \frac{a+1-1}{1+a} = 1 - \frac{1}{1+a}$ should prove useful

Answer:

$$\frac{d}{d\sigma} \frac{1}{1 + e^{-x}} = \dots = \frac{1}{1 + e^{-x}} \frac{e^{-x}}{1 + e^{-x}} = \frac{1}{1 + e^{-x}} \frac{1 + e^{-x} - 1}{1 + e^{-x}} = \frac{1}{1 + e^{-x}} \left(1 - \frac{1}{1 + e^{-x}}\right) = \sigma(x)(1 - \sigma(x))$$

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooo●oooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
ooooooo

Chain rule

Function composition operator “ \circ ”:

Operator that takes two functions f and g and produces a function $h = g \circ f$ such that $h(x) = g(f(x))$. In other words, g is applied to the result of $f(x)$.

The chain rule states how the derivative of the composition of two differentiable functions f and g behaves:

Chain rule:

$$h' = (f \circ g)' = (f' \circ g) \cdot g'$$

or less formally

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

Note that the chain rule is a core utility in machine learning, e.g. forming the base of how neural networks are learning, among countless other uses.

Part 1: Linear Algebra

```
oooooooooooooo
oooo
oooooooooooo
```

Part 2: Calculus

```
oooooooooooooo●oooo
oooooooooooo
oooooooooooo
oooooooooooo
```

Part 3: Statistics and Probability

```
ooooooo
ooooooo
oooooooooooooooo
ooooooo
```

Use of chain rule is often straightforward. Just be clear what $f(x)$ and what $g(x)$ are.
Here are some examples.

Example 1:

$$\frac{d}{dx}(1 + x^2)^{100} = f'(g(x)) \cdot g'(x) = 100 g^{99}(x) \cdot 2x = 100(1 + x^2)^{99} \cdot 2x$$

Example 2:

$$\frac{d}{dx}\sqrt{3x+2} = \frac{d}{dx}(3x+2)^{1/2} = \frac{1}{2}(3x+2)^{-1/2} \cdot 3$$

Example 3:

$$\frac{d}{dx}(3x^5 - x + 14)^{11} = 11(3x^5 - x + 14)^{10} \cdot (15x^4 - 1)$$

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooo●ooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
ooooooo

Similarly the chain rule applies to composites of more than two functions:

$$(f \circ g \circ h)' = f' \circ (g \circ h) \cdot (g \circ h)' = f' \circ (g \circ h) \cdot g'(h) \cdot h' = (f' \circ g \circ h) \cdot (g' \circ h) \cdot h'$$

For example, compute $\frac{d}{dx} e^{\sin x^2}$:

- Define individual functions and compute their derivatives:

$$y := f(u) = e^u \implies \frac{dy}{du} = f'(u) = e^u = e^{\sin x^2}$$

$$u := g(v) = \sin v = \sin x^2 \implies \frac{du}{dv} = g'(v) = \cos v = \cos x^2$$

$$v := h(x) = x^2 \implies \frac{dv}{dx} = h'(x) = 2x$$

- Therefore:

$$\frac{d}{dx} e^{\sin x^2} = e^{\sin x^2} \cos(2x) 2x$$

Multivariate Chain Rule

Multivariate Chain Rule: Let

- $z = f(x_1, x_2, \dots, x_m)$ be differentiable in m independent variables, and let
- $\forall i \in 1, \dots, m : x_i = x_i(t_1, t_2, \dots, t_n)$ be differentiable in n independent variables

Then, $\forall j \in 1, 2, \dots, n$:

$$\frac{\partial z}{\partial t_j} = \sum_{i=1}^m \frac{\partial z}{\partial x_i} \frac{\partial x_i}{\partial t_j} = \frac{\partial z}{\partial x_1} \frac{\partial x_1}{\partial t_j} + \frac{\partial z}{\partial x_2} \frac{\partial x_2}{\partial t_j} + \cdots + \frac{\partial z}{\partial x_m} \frac{\partial x_m}{\partial t_j} = \nabla z \cdot \frac{\partial \mathbf{x}}{\partial t_j}$$

For example, in the case of $x(t)$, $y(t)$, and $z = f(x, y)$, we get

$$\frac{dz}{dt} = \frac{\partial z}{\partial x} \frac{dx}{dt} + \frac{\partial z}{\partial y} \frac{dy}{dt}$$

with $z = f(x(t), y(t))$ being differentiable at t .

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooo●oooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
oooooooooooo
oooooooooooooooo
oooooooooooo

Example: Multivariate Chain Rule

Given $u(x, y) = x^2 + 2y$ with $x(r, t) = r \sin(t)$ and $y(r, t) = \sin^2(t)$, we are interested in finding $\partial u / \partial r$ and $\partial u / \partial t$.

$$\frac{\partial u}{\partial r} = \frac{\partial u}{\partial x} \frac{\partial x}{\partial r} + \frac{\partial u}{\partial y} \frac{\partial y}{\partial r} = 2x \sin t + 2 \times 0 = 2r \sin^2 t$$

and

$$\begin{aligned}\frac{\partial u}{\partial t} &= \frac{\partial u}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial u}{\partial y} \frac{\partial y}{\partial t} \\ &= 2x \times r \cos t + 2 \times 2 \sin t \cos t \\ &= 2r \sin t \times r \cos t + 4 \sin t \cos t \\ &= 2(r^2 + 2) \sin t \cos t \\ &= (r^2 + 2) \sin(2t)\end{aligned}$$

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooo●
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
oooooooo
oooooooooooo
oooooooooooo



You are given 200 m of wire to make a fence around some rectangular area of side length a and b , respectively. How should you choose a and b s.t. the area which you are covering is maximized?

Answer:

$$\max_{a,b} ab \quad \text{s.t.} \quad 2a + 2b = 200 \iff \max_a a(100 - a) \iff a^* = 50 \quad (\text{square area})$$

Extreme Value Problems

Systematic procedure to find the minimum and maximum values of a function $f(x)$ on an interval $[a, b]$:

1. Solve $f'(x) = 0$ to find the list of critical points of f
2. Exclude any critical points not inside the interval $[a, b]$
3. Add to the list the *endpoints* a, b of the interval and any points of discontinuity or non-differentiability
4. At each point on the list, evaluate f : the biggest number that occurs is the maximum, and the littlest number that occurs is the minimum.

Example: Find the minima and maxima of the function $f(x) = x^4 - 8x^2 + 5$ on the interval $[-1, 3]$.

1. $f'(x) = 4x^3 - 16x = 0 \iff x(x^2 - 4) = 0 \iff x(x + 2)(x - 2) = 0$
So critical points are $-2, 0, 2$
2. Remove -2 from list since outside interval
3. Add endpoints $-1, 3$ to list. Now list is $-1, 0, 2, 3$
4. Evaluate f at $-1, 0, 2, 3$ which yields $-2, 5, -11, 14$, respectively

Hence maximum is 14, which occurs at $x = 3$, and the minimum is -11 , which occurs at $x = 2$.

Global and Local Extrema

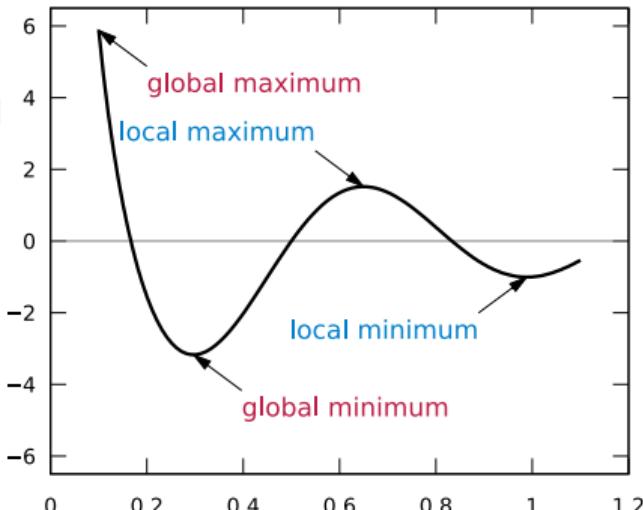
Use the *derivative test* to determine the type of each critical point.

In the scalar case:

- $f''(x) < 0$: local maximum
- $f''(x) > 0$: local minimum
- $f''(x) = 0$: test inconclusive \Rightarrow higher-order derivative test

Remarks:

- There are other tests to determine the nature of a critical point, e.g. first derivative test
- The above holds only for “reasonably well-behaved” functions
- The rules and intuition readily extend to the multi-dimensional case



Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oo●ooooooo
ooooooo
ooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
ooooooooooooooo
ooooooo

Mathematical Optimization

General optimization problem: Given: a function $f : A \rightarrow \mathbb{R}$ from some set A to the real numbers

Objective: find an element $x^* \in A$ such that

- $\forall x \in A : f(x^*) \leq f(x)$ “minimization”
- $\forall x \in A : f(x^*) \geq f(x)$ “maximization”

Optimization problems

- Two major classes: discrete and continuous optimization (focus on latter)
- We can always use *minimization* because

$$f(x^*) \geq f(x) \iff -f(x^*) \leq -f(x)$$

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
ooo●oooo
ooooooo
ooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
ooooooo

Optimization in Practice

In practice, algorithms that solve optimization problems fall into three broad classes:

- those that terminate in a finite number of steps, e.g. Simplex, combinatorial algos
- iterative methods that converge to a solution, e.g. Gradient descent, Newton's method, SQP
- heuristics that provide approx. solutions, e.g. genetic and evolutionary algos, PSO

Relevance for machine learning:

- Many machine learning problems are formulated as the minimization of some (so-called) loss function over a training set of examples
- Finding a solution may require to employ method from any of the three classes above

Part 1: Linear Algebra

```
oooooooooooooo
ooooo
oooo
oooooooooooooo
```

Part 2: Calculus

```
oooooooooooooooo
oooo●oooo
ooooooo
ooooooo
```

Part 3: Statistics and Probability

```
ooooooo
ooooooo
oooooooooooooooo
ooooooo
```

Nabla Operator: Gradients in \mathbb{R}^n are conveniently expressed with the *Nabla operator*:

$$\nabla \equiv \begin{pmatrix} \frac{\partial}{\partial x_1} \\ \frac{\partial}{\partial x_2} \\ \vdots \\ \frac{\partial}{\partial x_n} \end{pmatrix} \quad \text{or applied to some } f(\mathbf{x}) \quad \nabla f(\mathbf{x}) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}$$

Hessian Matrix: The Hessian is a square matrix of 2nd order partial derivatives of a scalar-valued function and a measure of local curvature:

$$H_f = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$$

Part 1: Linear Algebra
oooooooooooooo
ooooo
oooooooooooooo

Part 2: Calculus
oooooooooooooooooooo
ooooo●ooo
oooooooo
oooooooo

Part 3: Statistics and Probability
ooooooo
ooooooo
oooooooooooooooooooo
oooooooo

Gradient Descent

Gradient Descent (GD) is a first-order iterative optimization algorithm for finding a local minimum of a differentiable function:

- Assume $F(\mathbf{x})$ is defined and locally differentiable around \mathbf{a}
- Then $F(\mathbf{x})$ decreases fastest in the direction of the *negative* gradient of F at a , which is $-\nabla F(\mathbf{a})$
- For a sufficiently small step size (learning rate) $\gamma \in \mathbb{R}^+$, we have

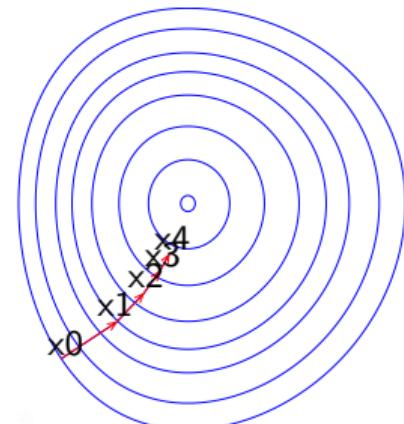
$$\mathbf{a}_{n+1} = \mathbf{a}_n - \gamma \nabla F(\mathbf{a}_n) \quad \text{and} \quad F(\mathbf{a}_n) \geq F(\mathbf{a}_{n+1})$$

Iterating the above we get the

Gradient Descent Update Rule:

$$\mathbf{x}_{n+1} = \mathbf{x}_n - \gamma_n \nabla F(\mathbf{x}_n), \quad n \geq 0$$

where $F(\mathbf{x}_0) \geq F(\mathbf{x}_1) \geq F(\mathbf{x}_2) \geq \dots$ is a monotonic sequence



Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooo●○o
ooooooo
ooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
ooooooooooooooo
ooooooo

Note: Convergence is *guaranteed* under certain reasonable assumptions and a proper choice of γ (e.g. via line search or the Barzilai Borwein method)

Wide range of applications: Linear and non-linear systems, least squares problems, anything “minimizable”

Part 1: Linear Algebra
oooooooooooooo
ooooo
oooooooooooooo

Part 2: Calculus
oooooooooooooooo
ooooooo●o
ooooooo
ooooooo

Part 3: Statistics and Probability
ooooooo
ooooooo
oooooooooooooooo
ooooooo

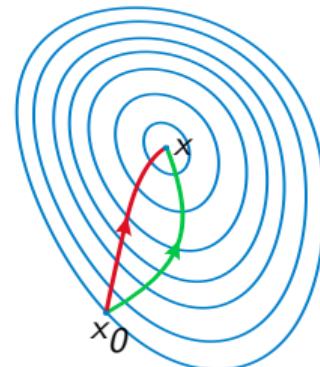
Newton's Method

Iterative procedure for finding roots x^* of a differentiable function F :

$$\{x_i^* : F(x_i^*) = 0\}$$

Intuitively, does so by repeatedly intersecting tangents with x -Axis:

$$y = F(x_k) + F'(x_k)(x - x_k) \stackrel{!}{=} 0 \implies x_{k+1} = x_k - \frac{F(x_k)}{F'(x_k)}$$



Applicable to minimization problem $\min_{x \in \mathbb{R}} f(x)$: Simply let $F(x) = f'(x)$ which yields:

Newton's Method (for optimization):

$$x_{k+1} = x_k - \frac{f'(x_k)}{f''(x_k)} \quad (\text{univariate case})$$

Corresponds to a second-order approximation of $f \Rightarrow$ fast convergence

Part 1: Linear Algebra

```
oooooooooooooo
ooooo
oooo
oooooooooooo
```

Part 2: Calculus

```
ooooooooooooooo
ooooooo●
ooooooo
ooooooo
```

Part 3: Statistics and Probability

```
ooooooo
ooooooo
ooooooooooooooo
ooooooo
```

Really Newton is just a second-order Taylor expansion of f around x_k :

$$f(x_k + t) \approx f(x_k) + f'(x_k)t + \frac{1}{2}f''(x_k)t^2$$

and minimizing its first derivative:

$$\frac{d}{dt}f(x_k + t) \approx f'(x_k) + f''(x_k)t \stackrel{!}{=} 0 \implies x_{k+1} = x_k - \frac{f'(x_k)}{f''(x_k)}$$

Similarly in the multi-variate case, with \mathbf{H}_f the Hessian of f :

$$f(\mathbf{x}_k + \mathbf{a}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top \mathbf{a} + \frac{1}{2}\mathbf{a}^\top \mathbf{H}_f \mathbf{a} \quad \text{and} \quad \frac{\partial}{\partial \mathbf{a}} f(\mathbf{x}_k + \mathbf{a}) \approx \nabla f(\mathbf{x}_k) + \mathbf{H}_f \stackrel{!}{=} \mathbf{0}$$

Newton's Method (for optimization):

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{H}_f^{-1} \nabla f(\mathbf{x}_k) \quad (\text{multivariate case})$$

Note that this yields the exact solution for any quadratic objective function in a *single* step.

Part 1: Linear Algebra
oooooooooooooo
ooooo
oooooooooooooo

Part 2: Calculus
oooooooooooooooooooo
oooooooooooo
●oooooooo
oooooooooooo

Part 3: Statistics and Probability
ooooooo
ooooooo
oooooooooooooooooooo
oooooooooooo

Stochastic Gradient Descent (SGD)

SGD is:

- an iterative optimization method, stochastic approximation of gradient descent
- the fundamental idea behind many machine learning algorithms such as linear support vector machines, logistic regression, graphical models
- the de facto standard training algorithm for neural networks

Given an objective function $Q_i(\mathbf{w})$ which computes some sort of error for sample i , e.g.
 $Q_i(\mathbf{w}) = (\hat{y}_i - y_i)^2$

Classical SGD:

```
Initialize parameter vector  $\mathbf{w}$  and set learning rate  $\eta$ 
while termination criterion not satisfied do
    randomly shuffle samples in training set
    for  $i \leftarrow 1, n$  do
         $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla Q_i(\mathbf{w})$ 
    end for
end while
```

Part 1: Linear Algebra

```
oooooooooooooo
ooooo
oooooooooooooo
```

Part 2: Calculus

```
ooooooooooooooo
oooooooooooo
●oooooo
oooooooo
```

Part 3: Statistics and Probability

```
ooooooo
ooooooo
oooooooooooooooo
oooooooo
```

Example: SGD to compute a least squares fit of a straight line

Given training data $(x_i, y_i)_{i=1}^N$, fit the model

$$\hat{y} = w_1 + w_2 x$$

Objective function (minimize):

$$Q(w) = \sum_{i=1}^n Q_i(w) = \sum_{i=1}^n (\hat{y}_i - y_i)^2 = \sum_{i=1}^n (w_1 + w_2 x_i - y_i)^2$$

The update step of the Basic SGD algorithm above becomes

$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \leftarrow \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} - \eta \begin{bmatrix} \frac{\partial}{\partial w_1} (w_1 + w_2 x_i - y_i)^2 \\ \frac{\partial}{\partial w_2} (w_1 + w_2 x_i - y_i)^2 \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} - \eta \begin{bmatrix} 2(w_1 + w_2 x_i - y_i) \\ 2x_i(w_1 + w_2 x_i - y_i) \end{bmatrix}.$$

Part 1: Linear Algebra
oooooooooooo
oooo
oooooooooooo

Part 2: Calculus
oooooooooooo
oooooooo
oo●oooo
oooooooo

Part 3: Statistics and Probability
ooooooo
ooooooo
oooooooooooo
ooooooo

SGD Extensions

Extension 1: Batch updates

Issue: Each iteration of classical SGD evaluates the gradient at a single point x_i only

- ⇒ often results in noisy estimates of the true gradient
- ⇒ computationally expensive (blocking vectorization)

Idea: Combine multiple gradient estimates of a so-called “mini-batch”, then do a single update

- ⇒ smoother convergence thanks to better estimates
- ⇒ speed-up owing to vectorization

Part 1: Linear Algebra

```
oooooooooooo
oooo
oooooooooooo
```

Part 2: Calculus

```
oooooooooooo
oooo
ooo●ooo
oooo
```

Part 3: Statistics and Probability

```
oooooo
oooooo
oooooooooooo
oooooo
```

Extension 2: Implicit updates (ISGD)**Issues:**

- Classical SGD sensitive to choice of η and may easily diverge
- Better convergence with some decay schedule of η over iterations, $\eta = \eta(\text{iter})$

Idea: Evaluate stochastic gradient at the next iteration rather than the current one:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla Q_i(\mathbf{w}_{k+1})$$

To see the effect, assuming an ordinary least squares objective

$$Q_i(\mathbf{w}) = \frac{1}{2}(\hat{y}_i - y_i)^2$$

in the usual samples $\{(x_i, y_i)\}_{i=1}^N$, compare the resulting update rules:

Classical SGD: $\mathbf{w}_{k+1} = \mathbf{w}_k + \eta(y_i - \mathbf{x}_i^\top \mathbf{w}_k) \mathbf{x}_i$

Implicit SGD (closed-form!): $\mathbf{w}_{k+1} = \mathbf{w}_k + \frac{\eta}{1 + \eta \|\mathbf{x}_i\|^2} (y_i - \mathbf{x}_i^\top \mathbf{w}_k) \mathbf{x}_i$

Note the normalizing effect on η , resulting in numerically stable procedure for virtually all η .

Part 1: Linear Algebra

```
oooooooooooo
oooo
oooooooooooo
```

Part 2: Calculus

```
oooooooooooo
oooo
oooo●ooo
oooo
```

Part 3: Statistics and Probability

```
oooooo
oooooo
oooooooooooooooo
oooooo
```

You are in the middle of a minimization with ISGD:

- Iteration $k = 348$, number of samples $n = 50'000$
- Current weight vector: $\mathbf{w}_{348} = \begin{pmatrix} 5 & -3 \end{pmatrix}^\top$
- Current sample: $\mathbf{x}_{1283} = \begin{pmatrix} 3 & -4 \end{pmatrix}^\top$, $y_{1283} = 57$
- Learning rate: $\eta = 0.04$



Continue the iteration for one more step to find the next weight vector \mathbf{w}_{349} .

Answer:

$$\begin{aligned}
 \mathbf{w}_{349} &= \mathbf{w}_{348} + \frac{\eta}{1 + \eta \|\mathbf{x}_i\|^2} \left(y_i - \mathbf{x}_i^\top \mathbf{w}_k \right) \mathbf{x}_i \\
 &= \begin{pmatrix} 5 \\ -3 \end{pmatrix} + \frac{0.04}{1 + 0.04 \times 25} (57 - 27) \begin{pmatrix} 3 \\ -4 \end{pmatrix} \\
 &= \begin{pmatrix} 6.8 \\ -5.4 \end{pmatrix}
 \end{aligned}$$

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooo●●
oooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
oooooooo

Extension 3: Momentum

Issue: Classical SGD frequently shows oscillations, hampering convergence

Idea: Introduce some “memory” $\Delta \mathbf{w}_k$ of the current direction, i.e., do not take “hard turns” at every step:

$$\Delta \mathbf{w}_{k+1} = -\eta \nabla Q_i(\mathbf{w}_k) + \alpha \Delta \mathbf{w}_k$$

resulting in the update rule

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \Delta \mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla Q_i(\mathbf{w}_k) + \alpha \Delta \mathbf{w}_k$$

with $\alpha \in [0, 1)$ an exponential decay factor (hyperparameter)

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooo●
oooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
oooooooo

More Extensions

More sophisticated extensions that adapt the learning rate exists:

- AdaGrad (adaptive gradient): offers per-parameter learning rate
- Adam (adaptive momentum): also per-parameter learning rate; basis for most modern algorithms (AdaMax, AMSGrad)
- Second-order methods: stochastic analogue of the standard (deterministic) Newton algorithm

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooo
oooooooo
ooooooo
●oooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooo
ooooooo

Integration

Definite integrals of real-valued functions $f(x)$ w.r.t. a real variable x on an interval $[a, b]$ is written as

$$\int_a^b f(x) dx$$

with limits a and b . A function f is *integrable* if its integral over its domain is finite

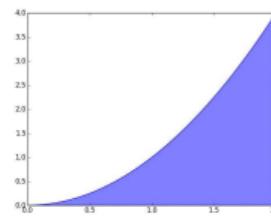
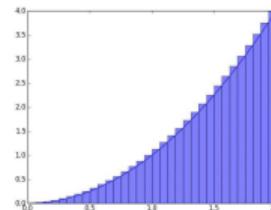
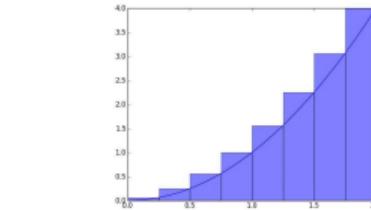
Indefinite integrals (antiderivative), written

$$F(x) := \int f(x) dx$$

represent a class of functions whose derivative is the integrand, hence

$$F' = f$$

Please note that all statements around *integrals* and *integrability* are meant in the *Riemann* sense



Part 1: Linear Algebra

```
oooooooooooo
oooo
oooooooooooo
```

Part 2: Calculus

```
oooooooooooo
oooo
oooo
o●oooo
```

Part 3: Statistics and Probability

```
oooooo
oooooo
oooooooooooo
oooooo
```

Property 1 (linearity): The integral of a linear combination is the linear combination of the integrals:

$$\int_a^b (\alpha f + \beta g)(x) dx = \alpha \int_a^b f(x) dx + \beta \int_a^b g(x) dx$$

Property 2: Flipping integration bounds flips the sign of the result,

$$\int_a^b f(x) dx = - \int_b^a f(x) dx$$

Property 3: Given $\int_a^b f(x) dx$ exists with $a \leq b$, then for any $c \in [a, b]$, we have

$$\int_a^b f(x) dx = \int_a^c f(x) dx + \int_c^b f(x) dx$$

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oo●oooo

Part 3: Statistics and Probability

ooooooo
ooooooo
ooooooooooooooo
ooooooo

Closed-form Integration

Refresher by example:

- $\int a \, dx = ax + C$

- $\int x^n \, dx = \frac{x^{n+1}}{n+1} + C, \quad n \neq -1$

- $\int (ax + b)^n \, dx = \frac{(ax + b)^{n+1}}{a(n+1)} + C, \quad n \neq -1$

- $\int \frac{1}{x} \, dx = \ln|x| + C$

- $\int e^{ax} \, dx = \frac{1}{a}e^{ax} + C$

- $\int f'(x)e^{f(x)} \, dx = e^{f(x)} + C$

- $\int a^x \, dx = \frac{a^x}{\ln a} + C$

- $\int \ln x \, dx = x \ln x - x + C$

- $\int \log_a x \, dx = \frac{x \ln x - x}{\ln a} + C$

- $\int \sin x \, dx = -\cos x + C$

- $\int \cos x \, dx = \sin x + C$

- $\int \tan x \, dx = \ln|\sec x| + C = -\ln|\cos x| + C$

- $\int_0^\infty \sqrt{x} e^{-x} \, dx = \frac{1}{2}\sqrt{\pi}$ Gamma function

- $\int_0^\infty e^{-ax^2} \, dx = \frac{1}{2}\sqrt{\frac{\pi}{a}}$ Gaussian integral, $a > 0$

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooo
ooooooo
oooo●ooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
ooooooo

Integration By Parts

If $y(x) = u(x)v(x)$, then

$$y' = (uv)' = u'v + uv' \iff uv' = (uv)' - u'v$$

Therefore we get:

Integration by parts (helpful to integrate products of functions):

$$\int uv' dx = \int (uv)' dx - \int u'v dx = uv - \int vu' dx$$

Example: To find the antiderivative $\int x \sin(x)dx$, we choose $u := x$ and $v := \sin x$,

$$\int x \sin x dx = -x \cos x - \int (-\cos x) dx = \sin x - x \cos x + C$$

Note that the “right” choice of u and v is crucial

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooo●oo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
ooooooo

Integration By Substitution

Transforms one integral into another integral that is easier to compute

Integration by substitution (helpful to integrate compositions of functions):

To transform $\int f(g(x)) dx$

- choose a *suitable* substitution $u = g(x)$
- compute $du = g'(x) dx$ (not obvious)
- substitute integrand and differential, then solve the (simpler) integral
- back-substitute

Example: To find the antiderivative of $x \cos(x^2 + 1)$, choose $u := x^2 + 1$, $du = 2 dx$. Then

$$\int x \cos(x^2 + 1) dx = \frac{1}{2} \int 2x \cos(x^2 + 1) dx = \frac{1}{2} \int \cos u du = \frac{1}{2} \sin u + C = \frac{1}{2} \sin(x^2 + 1) + C$$

Numerical Integration

Issue: In practice, many integrals possess no closed-form solution

Resort: Quadrature (numerical integration) for 1-dimensional integrals works well, e.g.

$$\int_a^b f(x) dx \approx \sum_{i=1}^n f(x_i^*) \Delta x$$

Multi-dimensional integrals (arising e.g. as normalizers in Bayesian stats) require more sophisticated techniques such as Monte Carlo integration. Popular algorithms in this class include Metropolis-Hastings and Gibbs sampling (both instances of Markov Chain Monte Carlo)

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooo
oooooooo
ooooooo
oooooo●

Part 3: Statistics and Probability

ooooooo
ooooooo
ooooooooooooooo
ooooooo

Practice!

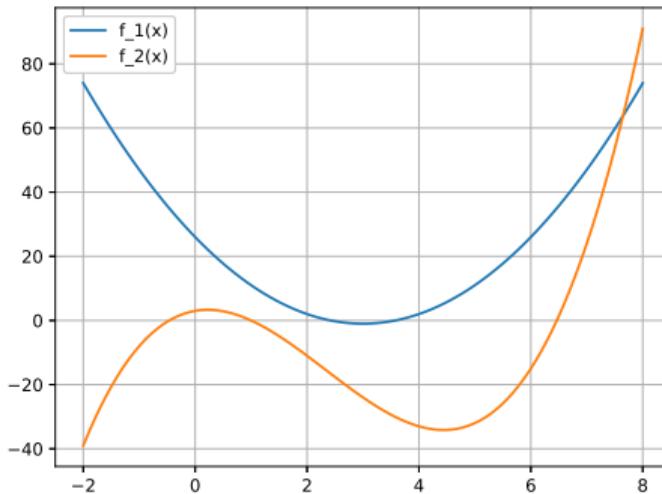
Implement the 1-dimensional version of

- Gradient Descent
- Basic SGD

Then compare convergence of the two algorithms
on one or both test functions

$$f_1(x) = 3(x-3)^2 - 1 \quad \text{and} \quad f_2(x) = (x-1)^3 - 4x^2 + 4$$

Also vary the initial guess x_0 and the learning rate η . What happens e.g. for f_2 with $x_0 = 0$?



Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooooooo
oooooooooooo
oooooooo
oooooooo

Part 3: Statistics and Probability

●oooooo
oooooooo
oooooooooooooooooooo
oooooooo

Introduction

Background

- Probability is the mathematical field concerned with reasoning under uncertainty
- Probabilistic models enable us to reason about the *likelihood* of various events
- Use of probabilities to describe the frequencies of repeatable events (like coin tosses) is fairly uncontroversial

Schools of thoughts

- Frequentists scholars adhere to an interpretation of probability that applies only to such repeatable events
- Bayesian scholars use the language of probability more broadly to formalize our reasoning under uncertainty

Bayesian probability is characterized by two unique features:

- assigning prior belief e.g., $P[\text{"moon is made out of cheese"}]$, opinions may differ
- evidence, that is hard facts that we measure

Bayesian probability provides unambiguous rules for how we should *update* our belief in light of new evidence, while allows for different individuals to start off with different prior beliefs

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooooooo
oooooooooooo
oooooooo
oooooooo

Part 3: Statistics and Probability

●○○○○
○○○○○○
oooooooooooooooooooo
○○○○○○

Terminology

Imagine the experiment:

- Toss a coin n times. The probability for heads is $p_h = 1/2$.
- We observe n_h heads and $n_t = n - n_h$ tails
- Fraction of heads that we *expect* to see should exactly match the expected fraction of tails
- The larger n the less likely $n_h = n_t$ becomes

Terminology:

- p_h is called a *probability*, captures the certainty with which any given toss will come up heads
- Probabilities assign scores between 0 and 1 to *events* (i.e. the outcomes of interest)
- Here the event of interest is “heads” and we denote the corresponding probability $P(\text{heads})$
- $P(\cdot) = 1$ indicates absolute certainty, and $P(\cdot) = 0$ impossibility
- n_h/n and n_t/n are called *statistics* (not probabilities)—empirical quantities computed based on the observed data
- *Estimators* are special statistics that use the data to produce estimates of model parameters (like probabilities). *Consistent estimators* converge to the corresponding probability (cf. next slide)

Part 1: Linear Algebra

○○○○○○○○○○○○
○○○○
○○○○○○○○○○○○

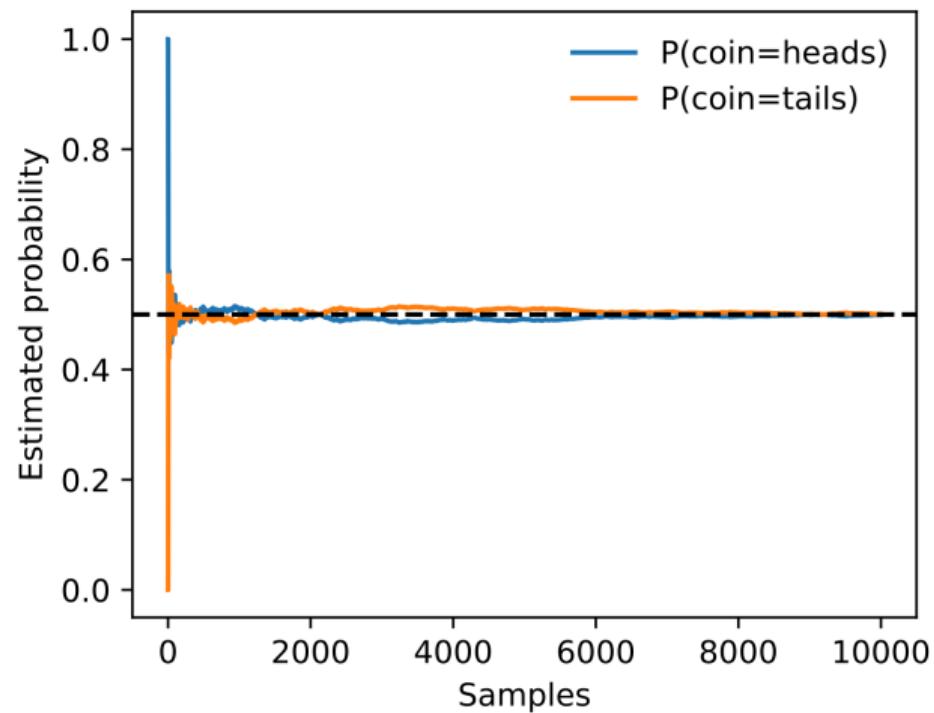
Part 2: Calculus

○○○○○○○○○○○○○○○○
○○○○○○○○○○
○○○○○○○○
○○○○○○○○

Part 3: Statistics and Probability

○○●○○○
○○○○○○
○○○○○○○○○○○○○○
○○○○○○○○

Estimators and Consistency



Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooooooo
oooooooooooo
oooooooo
oooooooo

Part 3: Statistics and Probability

ooo●ooo
oooooooo
oooooooooooooooooooo
oooooooo

Formal Definition

Denote the set of possible outcomes \mathcal{S} , and call it the *sample space* or *outcome space*, e.g.:

- Rolling a single coin: $\mathcal{S} = \{\text{heads, tails}\}$
- Rolling a single die, $\mathcal{S} = \{1, 2, 3, 4, 5, 6\}$
- flipping two coins, possible outcomes are $\{(\text{heads, heads}), (\text{heads, tails}), (\text{tails, heads}), (\text{tails, tails})\}$

Events are subsets of the sample space.

Axioms of Probability: A *probability function* maps events onto real values $P : \mathcal{A} \subseteq \mathcal{S} \rightarrow [0, 1]$. The probability of an event \mathcal{A} in the given sample space \mathcal{S} , denoted $P(\mathcal{A})$, satisfies the following properties (Kolmogorov, 1933):

- The probability of any event \mathcal{A} is a non-negative real number, i.e., $P(\mathcal{A}) \geq 0$;
- The probability of the entire sample space is 1, i.e., $P(\mathcal{S}) = 1$;
- For any countable sequence of events $\mathcal{A}_1, \mathcal{A}_2, \dots$ that are *mutually exclusive* ($\mathcal{A}_i \cap \mathcal{A}_j = \emptyset$ for all $i \neq j$), the probability that any of them happens is equal to the sum of their individual probabilities, i.e., $P(\bigcup_{i=1}^{\infty} \mathcal{A}_i) = \sum_{i=1}^{\infty} P(\mathcal{A}_i)$.

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

oooo●o
oooooooooooo
oooooooooooooooooooo
oooooooooooo

Random Variables

A *random variable* (RV) is a mapping from an underlying sample space to a set of values

- Random variables can be much coarser than the raw sample space. E.g. define RV “greater than 0.5” is binary but has an infinite sample space
- Multiple RVs may share the same underlying sample space. E.g. for RVs “home alarm goes off” and “my house was burglarized”, knowing the value taken by one random variable can tell us something about the likely value of another random variable

Every value taken by a random variable corresponds to a subset of the underlying sample space. Thus the occurrence where the random variable X takes value v , denoted by $X = v$, is an *event* and $P(X = v)$ denotes its probability.

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

oooo●
ooooooo
oooooooooooo
ooooooo

Notation in practice

- $P(X)$ may refer broadly to the *distribution* of X
- $P(X, Y) = P(X)P(Y)$, may be shorthand to express a statement that is true for all of the values that the random variables X and Y , i.e., $\forall i, j : P(X = i \text{ and } Y = j) = P(X = i)P(Y = j)$
- $P(v)$ may be used when the random variable is clear from the context
- $P(1 \leq X \leq 3)$ may denote the probability of the event $\{1 \leq X \leq 3\}$

In general, there's two types of random variables:

- Discrete, e.g. flips of a coin or tosses of a die
- Continuous, e.g. height H of a person. Note that in this case, asking e.g. for $P(H = 1.8796549851915233234)$ does not make sense. Instead, we work with probability *densities* p

Part 1: Linear Algebra
oooooooooooo

Part 2: Calculus
oooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability
oooooooo
●oooooooo
oooooooooooooooo
oooooooooooo

Probability Density Functions (PDFs)

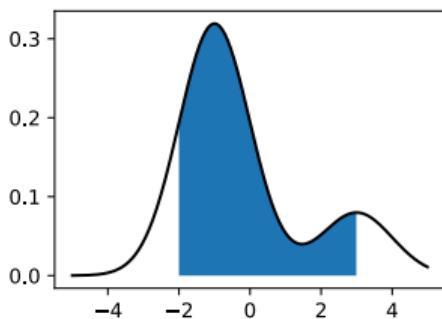
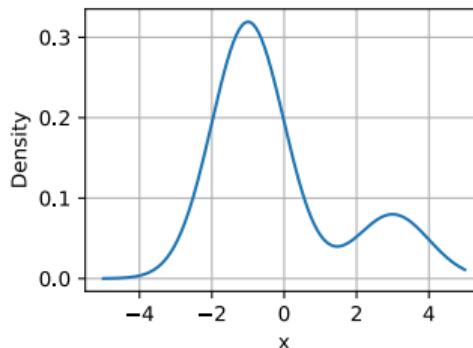
PDF: Informally, a PDF $p(x)$ is characterized by the following properties:

- $p(x) \geq 0$
- $\int_{-\infty}^{\infty} p(x) dx = 1$

From this it follows that the probability of any single value x is 0, $P(X = x) = 0$ and that $p(x) > 1$ is admissible

To compute actual probabilities we integrate $p(x)$ over an interval:

$$P(X \in (a, b]) = \int_a^b p(x) dx$$



Cumulative Distribution Functions (CDFs)

The CDF $F(x)$ of a random variable X with density $p(x)$ is given by

$$F(x) = \int_{-\infty}^x p(x) dx = P(X \leq x)$$

Note that $F(x)$ represents an actual probability

Properties:

- $F(x) \rightarrow 0$ as $x \rightarrow -\infty$
- $F(x) \rightarrow 1$ as $x \rightarrow \infty$
- $F(x)$ is non-decreasing, i.e. $y > x \implies F(y) \geq F(x)$
- $F(x)$ is continuous (has no jumps) if X is a continuous random variable

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
oo●ooo
oooooooooooooooo
oooooooooooo

Characterizing Distributions

Expected value

Expected value of a RV (mean):

- $\mu_X = E[X] = \sum_i x_i p_i$ (for discrete RV X)
- $\mu_X = E[X] = \int_{-\infty}^{\infty} x p(x) dx$ (for continuous RV X)

Properties:

- For any RV X and numbers a and b , we have that $\mu_{aX+b} = a\mu_X + b$
- For any two random variables X and Y , we have $\mu_{X+Y} = \mu_X + \mu_Y$

Means characterize average behavior of a RV, which is not sufficient however: A profit of CHF10 \pm CHF1 per sale is very different from making CHF10 \pm CHF15 per sale.

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooo●ooo
oooooooooooooooooooo
oooooooooooo

One possibility to quantify risk r as expected loss. That is, by weighting the loss l (e.g., in CHF) which some event causes with the probability p_e that this event actually occurs:

$$r := p_e l$$



Now assume you want to park your car in the city of Zurich. The parking garage will charge you CHF 24. If you park in some quiet corner, you may get fined and pay CHF 40. The chance that you'll get caught is $2/3$. What's the strictly economic choice here, parking garage or taking your chances?

Answer:

The parking garage as it has the lower expected loss: $24 < \frac{2}{3} \times 40 \approx 26.67$

Part 1: Linear Algebra
oooooooooooooo
ooooo
oooooooooooooo

Part 2: Calculus
oooooooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability
ooooooo
oooo●ooo
oooooooooooooooooooo
oooooooooooo

Characterizing Distributions

Variance and Standard Deviation

Variance: The variance measures the variation of a RV about its mean:

$$\sigma_X^2 = \text{Var}(X) = E[(X - \mu_X)^2]$$

Properties:

- For any RV X , $\text{Var}(X) \geq 0$, with $\text{Var}(X) = 0$ iff X constant
- For any RV X and numbers a and b , we have that $\text{Var}(aX + b) = a^2\text{Var}(X)$
- For any two independent RV X and Y , we have $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$
- $E[(X - \mu_X)^2] = E[X^2] - \mu_X^2$

Standard Deviation (SD): The SD also measures the variation of a RV about its mean, but has the same units as the mean:

$$\sigma_X = \sqrt{\text{Var}(X)}$$

Analogous properties as in the case of variance apply

Continuous Distributions

A Well Behaved Example

Task: Compute the first two moments of the uniform distribution

$$p(x) = \begin{cases} 1 & x \in [0, 1], \\ 0 & \text{otherwise} \end{cases}$$

Find expected value:

$$\mu_X = \int_{-\infty}^{\infty} xp(x) dx = \int_0^1 x dx = \frac{1}{2}$$

Find variance:

$$\sigma_X^2 = \int_{-\infty}^{\infty} x^2 p(x) dx - \left(\frac{1}{2}\right)^2 = \frac{1}{3} - \frac{1}{4} = \frac{1}{12}$$

where we used above identity $E[(X - \mu_X)^2] = E[X^2] - \mu_X^2$

Continuous Distributions

An Less Well-Behaved Example

Task: Compute the first two moments of a Cauchy-like distribution with the PDF given by

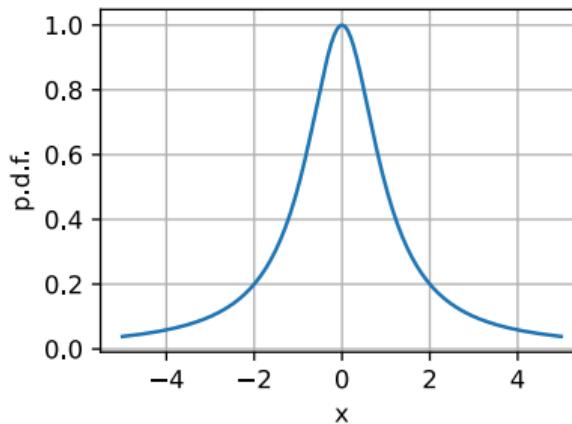
$$p(x) = \frac{1}{1+x^2}$$

Trying to compute the variance yields

$$\sigma_X^2 = \int_{-\infty}^{\infty} x^2 p(x) dx = \int_{-\infty}^{\infty} \frac{x^2}{1+x^2} dx = \infty$$

Hence the variance of $p(x)$ is undefined.

A similar problem arises for the mean, which is undefined too.



Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
●oooooooooooo
ooooooo

Joint Density

Similar to the univariate case, a PDF can be defined for two or more variables.

The two RVs X and Y can be defined to have a *joint density* $p(x, y)$ or $p_{X,Y}$. Then properties similar to the univariable case hold:

- $p(x, y) \geq 0$
- $\int_{\mathbb{R}^2} p(x, y) dx dy = 1$
- $P((X, Y) \in \mathcal{D}) = \int_{\mathcal{D}} p(x, y) dx dy$

This generalizes to any number of RVs.

Part 1: Linear Algebra

○○○○○○○○○○○○
○○○○
○○○○○○○○○○○○

Part 2: Calculus

○○○○○○○○○○○○○○○○
○○○○○○○○○○
○○○○○○○○
○○○○○○○○

Part 3: Statistics and Probability

○○○○○○
○○○○○○○○○○○○○○○○
○●○○○○○○○○○○○○○○○○
○○○○○○○○

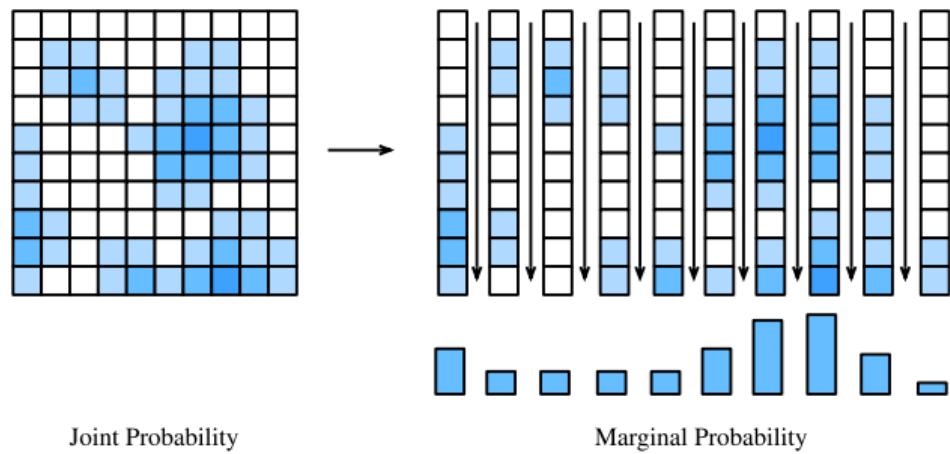
Marginal Distributions

When dealing with multiple variables, we often ask “how is this one variable distributed?”. The distribution that answers this question is the *marginal distribution*.

Mathematically (often referred to as *integrating out* or *marginalizing out* unneeded variables):

$$p_X(x) = \int_{-\infty}^{\infty} p_{X,Y}(x,y) dy$$

Visually:



Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooooooo
oooooooooooo
oooooooo
oooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oo●oooooooooooo
ooooooo

Conditional Probability and Statistical Independence

The conditional distribution quantifies the probability of an event occurring, given that another event (e.g. by assertion or by evidence) has already occurred.

Conditional Probability: We define the conditional probability of " \mathcal{A} given \mathcal{B} ", written $P(\mathcal{A} | \mathcal{B})$ to be:

$$P(\mathcal{A} | \mathcal{B}) = \frac{P(\mathcal{A}, \mathcal{B})}{P(\mathcal{B})} \quad \text{or equivalently} \quad P(\mathcal{A}, \mathcal{B}) = P(\mathcal{A} | \mathcal{B}) P(\mathcal{B})$$

Based on the above we define independence:

Statistical Independence: Events \mathcal{A} and \mathcal{B} are defined to be *statistically independent* if

$$P(\mathcal{A}, \mathcal{B}) = P(\mathcal{A}) P(\mathcal{B})$$

Therefore independence implies both $P(\mathcal{A} | \mathcal{B}) = P(\mathcal{A})$ and $P(\mathcal{B} | \mathcal{A}) = P(\mathcal{B})$

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
ooo●oooooooooooo
ooooooo

Conditional Independence

Conditional probabilities are proper probabilities \Rightarrow concepts of (in)dependence also apply to them.

Conditional Independence: Two RVs A and B are *conditionally independent* given a third RV C iff

$$P(A, B | C) = P(A | C) P(B | C)$$

Two independent RV can become dependent when conditioning on a third RV C . Occurs when RVs A and B correspond to causes of RV C .

Example:

- A broken bone and lung cancer might be independent in the general population
- Now condition on being in the hospital
- Conditionally, the broken bone is negatively correlated with lung cancer
- We say the broken bone *explained away* why some person is in the hospital, in turn lowering the probability that the person has lung cancer

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
oooooooooooo
oooo●oooooooooooo
oooooooooooo

Bayes' Theorem

By definition of conditional probabilities, $P(A, B) = P(B | A) P(A)$ and $P(A, B) = P(A | B) P(B)$, which is equivalent to

Bayes' Theorem:

$$P(A | B) = \frac{P(B | A) P(A)}{P(B)}$$

Please note the profound implications: It allows us to reverse the order of conditioning. If we know how to estimate $P(B | A)$, $P(A)$, and $P(B)$, then we can estimate $P(A | B)$.

Example:

Given the prevalence of symptoms for a given disease, and the overall prevalences of the disease and symptoms, respectively, we can determine how likely someone is to have the disease based on their symptoms.

Special case: Without direct access to $P(B)$, we still know that

$$P(A | B) \propto P(B | A) P(A)$$

and use the normalization constraint $\sum_a P(A = a | B) = 1$.

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
ooooo●ooooooo
ooooooo

Marginalization

Marginalization (“sum rule”) is key ingredient to make Bayes theorem work.

It simply sums or integrates out, respectively, other variables from the joint distribution:

$$P(B) = \sum_a P(A, B)$$

$$p_B = \int_{-\infty}^{\infty} p(a, b) da$$

The result is called marginal probability.

In practice, where joint probabilities consist of many variables, this “normalizer” term in Bayes’ theorem often becomes a computational bottleneck.

It's 1963 and you participate in "Let's Make a Deal" (TV Show). You play this game:

- There's three doors: behind two of them there's a goat, behind the third a Ferrari
- First you pick a door
- Then the host opens one of the two other doors with a goat behind
- Then you can optionally switch to the other closed door or stay with your initial one
- Then your door will be opened and if it's the Ferrari, the car is yours



If you want the car, what's best to do? a) stay, b) switch, c) makes no difference
Hint: Use Python or Bayes' theorem

Answer:

With loss of generality, you pick door 1 and the host opens door 2 or 3. Let

$D := \text{"number of the door with the car behind"} \in \{1, 2, 3\}$, and likewise

$H := \text{number of door which host opens}$. At the end of the game there's two options, $d = 1$ and $d = 3$:

$$P(D = d | H = 2) = \frac{P(H = 2 | D = d) P(D = d)}{P(H = 2)}$$

$$P(D = 1 | H = 2) = \frac{1/2 \times 1/3}{1/3 + 1/6} = \frac{1}{3} \quad P(D = 3 | H = 2) = \frac{1 \times 1/3}{1/3 + 1/6} = \frac{2}{3}$$

Part 1: Linear Algebra

oooooooooooooo
ooooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

oooooooooooo
oooooooooooo
oooooooooooo●oooooooo
oooooooooooo

Covariance

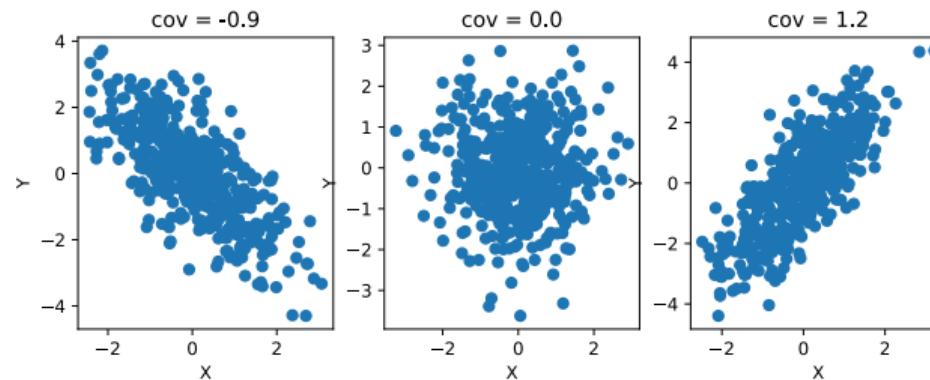
Covariance: Covariance measures the degree that two random variables fluctuate together:

$$\sigma_{XY} = \text{Cov}(X, Y) = E[XY] - E[X]E[Y]$$

Spelled out this means for the discrete and continuous case, respectively,

$$\sigma_{XY} = \sum_{i,j} (x_i - \mu_X)(y_j - \mu_Y)p_{ij} \quad \text{and} \quad \sigma_{XY} = \int_{\mathbb{R}^2} (x - \mu_X)(y - \mu_Y)p(x, y) dx dy$$

Visually:



Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooo●oooo
ooooooo

Summarizing the main properties of covariance:

- For any RV X , $\text{Cov}(X, X) = \text{Var}(X)$
- For any RVs X and Y and numbers a and b , $\text{Cov}(aX + b, Y) = \text{Cov}(X, aY + b) = a\text{Cov}(X, Y)$
- If X and Y are independent then $\text{Cov}(X, Y) = 0$

Correlation

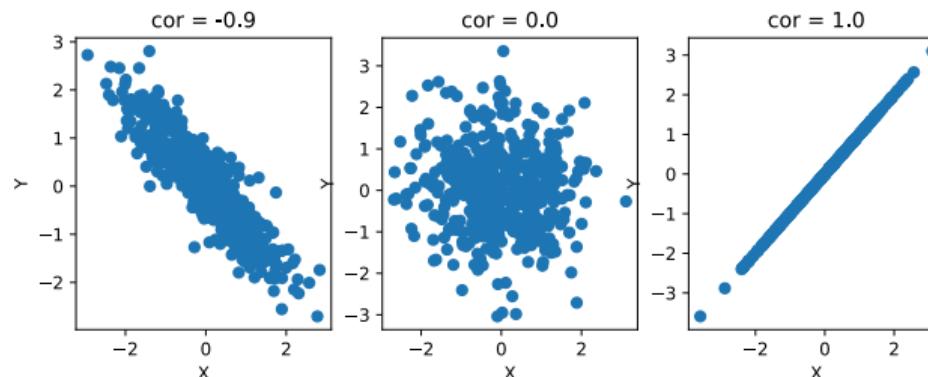
Correlation coefficient: Define the *correlation coefficient* to be

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

Properties:

- For any RV X , $\rho(X, X) = 1$
- For any RVs X and Y and numbers a and b , $\rho(aX + b, Y) = \rho(X, aY + b) = \rho(X, Y)$
- If X and Y are independent with non-zero variance then $\rho(X, Y) = 0$

Visually:



Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooo●oooo
ooooooo

Example

Let X be any random variable, and $Y = aX + b$ as linear deterministic function of X . Then, one can compute that

$$\sigma_Y = \sigma_{aX+b} = |a| \sigma_X$$

$$\text{Cov}(X, Y) = \text{Cov}(X, aX + b) = a\text{Cov}(X, X) = a\text{Var}(X),$$

and by definition

$$\rho(X, Y) = \frac{a\text{Var}(X)}{|a|\sigma_X^2} = \frac{a}{|a|} = \text{sign}(a)$$

Thus the correlation is $+1$ for any $a > 0$, and -1 for any $a < 0$.

This shows that correlation measures the *degree and directionality* of variation of the two RV, not the scale that the variation takes.

Part 1: Linear Algebra

```
oooooooooooo
oooo
oooo
oooooooooooo
```

Part 2: Calculus

```
oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo
```

Part 3: Statistics and Probability

```
ooooooo
ooooooo
oooooooooooo●ooo
ooooooo
```

Maximum Likelihood

Maximum Likelihood:

- A very common way of thinking in machine learning. Really it is a “point of view”
- Given a probabilistic model with unknown parameters, the most likely parameters are those under which the chance that the model generated the data are maximized

Technically, given our data X , we want

$$\hat{\theta} = \arg \max_{\theta} P(\theta | X)$$

where $\theta = (\theta_1, \dots, \theta_k)^\top$ the parameter vector of interest.

Using Bayes' theorem the Maximum Likelihood Estimator (MLE) for θ is

$$\hat{\theta} = \arg \max_{\theta} \frac{P(X | \theta)P(\theta)}{P(X)} = \arg \max_{\theta} P(X | \theta)P(\theta)$$

since $P(X)$ is just scaling the objective w.r.t. θ .

Exact same arguments apply to the densities of the continuous case.

Part 1: Linear Algebra
oooooooooooooo
oooo
oooooooooooo

Part 2: Calculus
oooooooooooooooooooo
oooooooooooo
oooooooo
oooooooo

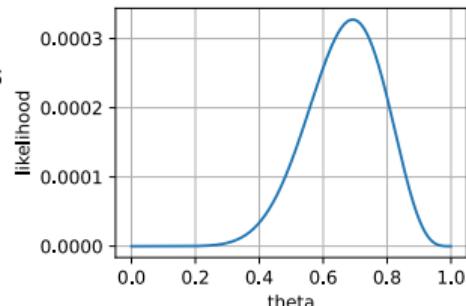
Part 3: Statistics and Probability
oooooooo
oooooooooooo
oooooooooooo●o
oooooooo

Maximum Likelihood: An Example

Setting:

- Given a single, unknown parameter θ : Probability that a coin flip is heads
- Assume “uninformative prior” \Rightarrow drop $P(\theta)$ from Bayes’ theorem
- Data: “HHHTHTTTHHHHHT”, i.e. $n_H = 9$ heads, $n_T = 4$ tails

Question: What is the probability that the coin comes up heads?



Answer: Assume samples independent \Rightarrow probabilities multiply. Likelihood becomes

$$P(X | \theta) = \theta^{n_H} (1 - \theta)^{n_T} = P(X | \theta) = \theta^9 (1 - \theta)^4$$

Then maximize the likelihood: The first-order condition

$$\frac{d}{d\theta} P(X | \theta) = \frac{d}{d\theta} \theta^9 (1 - \theta)^4 = 9\theta^8(1 - \theta)^4 - 4\theta^9(1 - \theta)^3 = \theta^8(1 - \theta)^3(9 - 13\theta) = 0$$

has three solutions: 0, 1 and $9/13$. First two are minima since $P(X | \theta = 0) = P(X | \theta = 1) = 0$. Last solution assigns non-zero probability to our data, hence the maximum likelihood estimate is $\hat{\theta} = 9/13$.

Note that despite everyone would correctly guess $9/13$, MLE provides a structured approach which equally applies to vastly more complex situations

Part 1: Linear Algebra
oooooooooooooo
ooooo
oooooooooooooo

Part 2: Calculus
oooooooooooooooooooo
oooooooooooo
oooooooo
oooooooooooo

Part 3: Statistics and Probability
oooooooo
oooooooooooo
oooooooooooo●
oooooooooooo

Numerical Optimization and the Negative Log-Likelihood

In practice, we may have billions of parameters and examples.

- Independence assumption results in long product of small values, say, e.g. $\prod_{i=1}^{10^9} \frac{1}{2}$, something far below machine precision
- However, the log of this product is a sum, $-\sum_{i=0}^{10^9} \log_{10} 2 = -10^9 \times \log_{10} 2 \approx -301029995.664$, which fits into a 32-bit floating point value

Hence should work with log-likelihood,

$$\log P(X | \theta)$$

which will not affect the result since log is monotonically increasing.

Often working with loss functions, to be *minimized*, hence use negative log-likelihood:

$$-\log P(X | \theta)$$

Example: For the coin flipping problem from before, this means minimizing

$$-\log(P(X | \theta)) = -\log(\theta^{n_H}(1 - \theta)^{n_T}) = -n_H \log(\theta) - n_T \log(1 - \theta)$$

which again leads to $\hat{\theta} = n_H / (n_H + n_T)$

Distributions

Bernoulli

Simple RV encoding e.g. a coin flip which comes up 1 (heads) with probability p and 0 (tails) with probability $1 - p$.

Bernoulli Distribution:

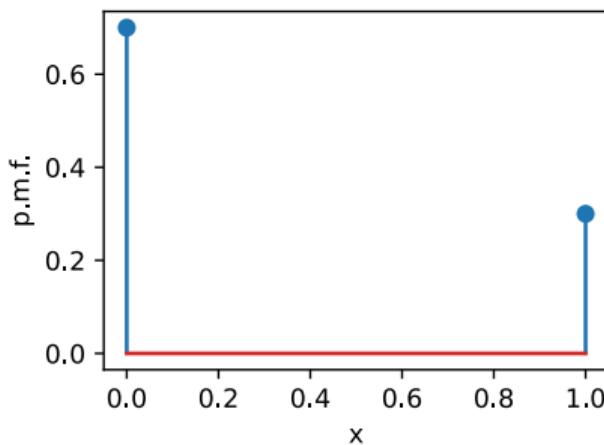
$$X \sim \text{Bernoulli}(p)$$

Cumulative distribution function:

$$F(x) = \begin{cases} 0 & x < 0, \\ 1 - p & 0 \leq x < 1, \\ 1 & x \geq 1. \end{cases}$$

Moments:

- $\mu_X = p$
- $\sigma_X^2 = p(1 - p)$



Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooo
o●ooooooo

Distributions

Discrete Uniform

Every value out of given set is equally likely has a probability of $p_i = \frac{1}{n}$ to occur

Discrete Uniform Distribution: Assume the set of possible outcomes is $\{1, 2, \dots, n\}$.

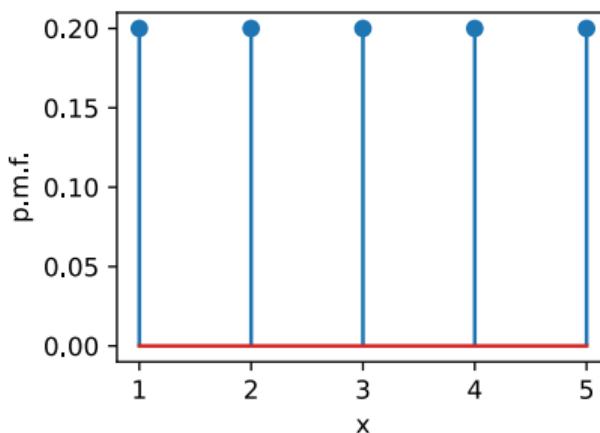
$$X \sim U(n)$$

Cumulative distribution function:

$$F(x) = \begin{cases} 0 & x < 1, \\ \frac{k}{n} & k \leq x < k + 1 \text{ with } 1 \leq k < n, \\ 1 & x \geq n. \end{cases}$$

Moments:

- $\mu_X = \frac{1+n}{2}$
- $\sigma_X^2 = \frac{n^2-1}{12}$



Distributions

Continuous Uniform

Similar to discrete uniform distribution: Starting from $U(n)$ with $n \rightarrow \infty$ and linearly scaling it to the interval $[a, b]$ will approach a continuous random variable that just picks an arbitrary value in $[a, b]$ with equal probability

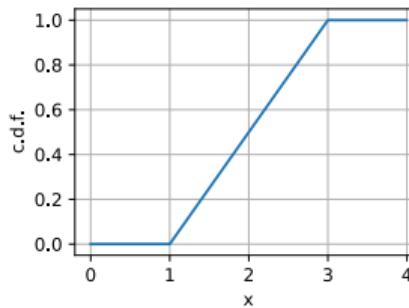
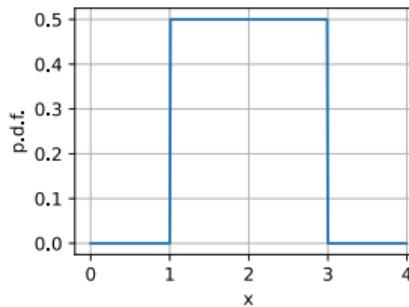
Continuous Uniform Distribution:

$$X \sim U(a, b)$$

PDF: $p(x) = \begin{cases} \frac{1}{b-a} & x \in [a, b] \\ 0 & x \notin [a, b] \end{cases}$

CDF: $F(x) = \begin{cases} 0 & x < a \\ \frac{x-a}{b-a} & x \in [a, b] \\ 1 & x \geq b \end{cases}$

Moments: $\mu_X = \frac{a+b}{2}$ $\sigma_X^2 = \frac{(b-a)^2}{12}$



Part 1: Linear Algebra

```
oooooooooooooo
ooooo
oooo
oooooooooooo
```

Part 2: Calculus

```
ooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo
```

Part 3: Statistics and Probability

```
ooooooo
ooooooo
ooooooooooooooo
ooo●ooo
```

Distributions

Binomial

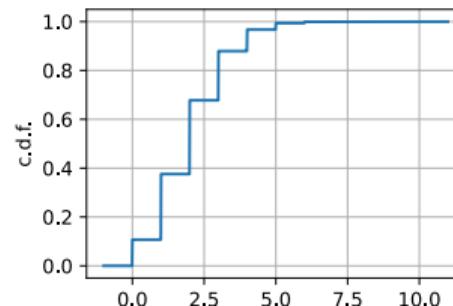
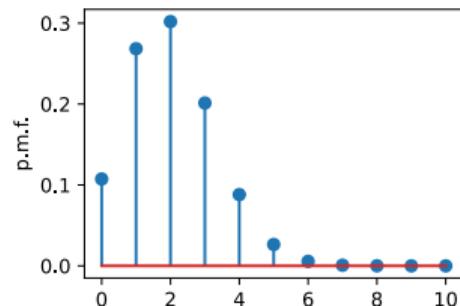
Given n independent experiments, each with success probability p . Interested in number of successes X :

$$X = \sum_{i=1}^n X_i \quad \text{with} \quad X_i \sim \text{Bernoulli}(p) \quad \iff \quad X \sim \text{Binomial}(n, p)$$

Binomial Distribution:

$$\text{CDF: } F(x) = \begin{cases} 0 & x < 0 \\ \sum_{m \leq k} \binom{n}{m} p^m (1-p)^{n-m} & k \leq x < k+1 \text{ with } 0 \leq k < n \\ 1 & x \geq n \end{cases}$$

$$\text{Moments: } \mu_X = np \quad \sigma_X^2 = np(1-p)$$



Part 1: Linear Algebra

```
oooooooooooo
oooo
oooo
oooooooooooo
```

Part 2: Calculus

```
oooooooooooo
oooo
oooo
oooo
oooo
```

Part 3: Statistics and Probability

```
oooo
oooo
oooo
oooooooooooo
oooo●ooo
```

Distributions

Poisson

Model for arrival rates, e.g. number of buses at a given stop per time interval. Call this rate λ (with units [1/s])

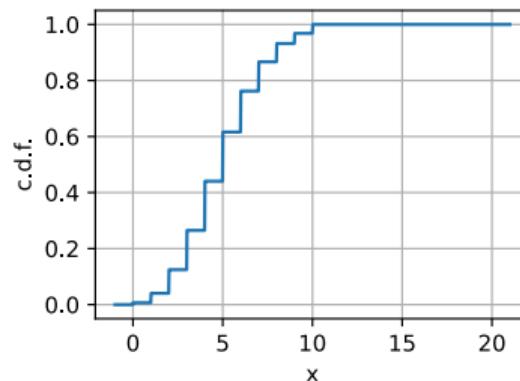
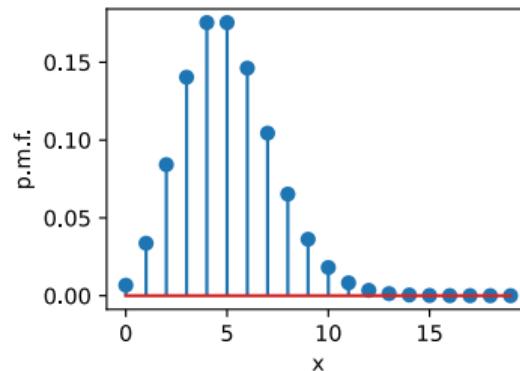
Poisson Distribution:

$$X \sim U(a, b)$$

PMF: $p_k = \frac{\lambda^k e^{-\lambda}}{k!}$

CDF: $F(x) = \begin{cases} 0 & x < 0, \\ e^{-\lambda} \sum_{m=0}^k \frac{\lambda^m}{m!} & k \leq x < k + 1 \text{ with } 0 \leq k \end{cases}$

Moments: $\mu_X = \lambda \quad \sigma_X^2 = \lambda$



Distributions

Gaussian

The central distribution in probability theory with many unique properties because of the *central limit theorem*.

Informally, given N arbitrarily distributed, i.i.d. RVs X_i and

$$X^{(N)} = \sum_{i=1}^N X_i$$

their sum. Then, under mild conditions,

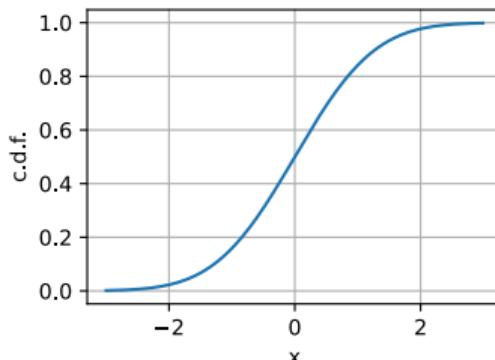
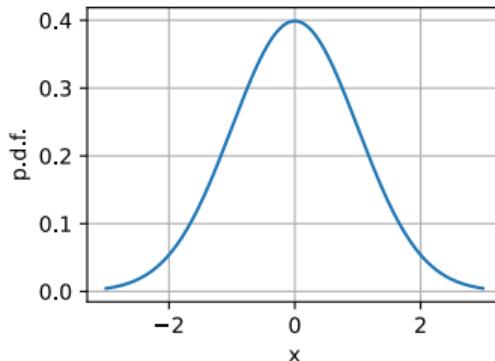
$$\frac{X^{(N)} - \mu_{X^{(N)}}}{\sigma_{X^{(N)}}} \sim \text{approximately Gaussian}$$

Gaussian Distribution:

$$X \sim \mathcal{N}(\mu, \sigma^2)$$

PDF: $p_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

Moments: $\mu_X = \mu$ $\sigma_X^2 = \sigma^2$



Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooooooo
ooooooo●o

Distributions: Summary

Things to remember about distributions:

- Bernoulli RVs can be used to model events with a yes/no outcome
- Discrete uniform distributions model selects from a finite set of possibilities
- Continuous uniform distributions select from an interval
- Binomial distributions model a series of Bernoulli RVs, and count the number of successes
- Poisson RVs model the arrival of rare events
- Gaussian RVs model the result of adding a large number of independent RVs together
- All the above distributions belong to the exponential family

Part 1: Linear Algebra

oooooooooooo
oooo
oooooooooooo

Part 2: Calculus

oooooooooooooooooooo
oooooooooooo
oooooooooooo
oooooooooooo

Part 3: Statistics and Probability

ooooooo
ooooooo
oooooooooooooooooooo
oooooooo●

License

This work is based on the following works:

- Jim Hefferson, *Linear Algebra*, 4th Edition, freely available at <https://hefferon.net/linearalgebra/>
- Wikipedia, The Free Encyclopedia, <https://www.wikipedia.org/>
- D2L.ai: Interactive Deep Learning Book with Multi-Framework Code, Math, and Discussions, freely available at <https://github.com/d2l-ai/d2l-en>

The above three works are based on Creative Commons Attribution-ShareAlike (CC BY-SA)-compatible licenses. Hence this work is again released under the same permissive Creative Commons Zero v1.0 Universal License. Sources are available at <https://github.com/openteaching/cas-m11>