
MATHEMATICS FOR MACHINE LEARNING

This version: November 29, 2025

Latest version: github.com/felipe-tobar/Maths-for-ML

Felipe Tobar
Department of Mathematics
Imperial College London

f.tobar@imperial.ac.uk
www.ma.ic.ac.uk/~ft410

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 4 |
| 2 | Optimisation | 5 |
| 2.1 | Terminology | 6 |
| 2.2 | Continuous unconstrained optimisation | 8 |
| 2.3 | Convex optimisation | 9 |
| 2.4 | First order methods | 12 |
| 2.4.1 | Role of the step size | 12 |
| 2.4.2 | Momentum | 13 |
| 2.4.3 | Newton method | 13 |
| 2.5 | Stochastic gradient descent | 14 |
| 2.5.1 | Adagrad | 15 |
| 2.5.2 | RMSProp | 15 |
| 2.5.3 | Adam | 15 |
| 2.6 | Training and Regularisation of Neural Networks | 16 |
| 2.6.1 | Network Architecture | 16 |
| 2.6.2 | Cost Function and Output Units | 17 |
| 2.6.3 | Forward Propagation | 17 |
| 2.6.4 | Backward Propagation – Preliminaries | 17 |
| 2.6.5 | Backward Propagation – Hidden Layers | 18 |
| 2.6.6 | Backward Propagation – Output Layer | 18 |
| 2.6.7 | Regularisation | 19 |
| 3 | Probability | 20 |
| 3.1 | Introduction | 20 |
| 3.1.1 | Definitions | 20 |
| 3.1.2 | Basic properties | 21 |
| 3.2 | Random variables | 21 |
| 3.2.1 | Discrete RVs | 21 |
| 3.2.2 | Continuous RVs | 22 |
| 3.2.3 | Properties and identities | 23 |
| 3.2.4 | Moments | 25 |
| 3.3 | Some particular RVs | 26 |
| 3.3.1 | Bernoulli and binomial | 26 |
| 3.3.2 | Uniform distribution | 27 |
| 3.3.3 | Gaussian distribution | 28 |
| 3.3.4 | Other distributions | 29 |
| 3.4 | Transformations of RVs distributions | 29 |
| 3.5 | Sampling | 30 |
| 4 | Statistics | 33 |
| 4.1 | The statistical model | 33 |
| 4.2 | Statistics | 34 |
| 4.3 | Estimators | 35 |
| 4.3.1 | Unbiased estimators | 35 |
| 4.3.2 | Loss functions | 37 |
| 4.4 | Maximum likelihood estimator (MLE) | 37 |
| 4.5 | MLE in practice: three examples | 39 |
| 4.5.1 | Gaussian linear regression | 39 |
| 4.5.2 | Classification | 40 |
| 4.5.3 | Latent variables: Expectation-Maximisation | 40 |
| 4.6 | Bayesian inference | 41 |

| | | |
|-------------------|--------------------------------|-----------|
| 4.7 | Conjugate priors | 43 |
| 4.8 | Bayesian estimators | 46 |
| 4.8.1 | Posterior predictive | 47 |
| References | | 49 |

1 Introduction

MISSING

2 Optimisation

NB: in this chapter, we follow (Murphy, 2022) and (Deisenroth, Faisal, & Ong, 2020). For a deeper presentation of the topics covered here, also see (Nesterov, 2018) and (Boyd & Vandenberghe, 2004).

Optimisation is central to ML, since models are *trained* by minimising a loss function (or optimising a reward function). In general, model design involves the definition of a training objective or **loss**, that is, a function that denotes **how well a model fits the data**. This training objective is a function of the training data and a chosen model, the latter usually represented by its parameters. The best model is thus chosen by optimising the loss function.

Example: Linear regression (LR)

In the LR setting, we aim to determine the function

$$\begin{aligned} f: \mathbb{R}^M &\rightarrow \mathbb{R} \\ x &\mapsto f(x) = a^\top x + b, \quad a \in \mathbb{R}^M, b \in \mathbb{R} \end{aligned} \quad (2.1)$$

conditional to a set of observations

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N \subset \mathbb{R}^M \times \mathbb{R}. \quad (2.2)$$

Using least squares, the function f is chosen via minimisation of the sum of the square differences between observations $\{y_i\}_{i=1}^N$ and predictions $\{f(x_i)\}_{i=1}^N$. That is, we aim to minimise the loss:

$$J(\mathcal{D}, f) = \sum_{i=1}^N (y_i - f(x_i))^2 = \sum_{i=1}^N (y_i - a^\top x_i - b)^2. \quad (2.3)$$

We show an example of a linear model learnt from data in Figure 1.

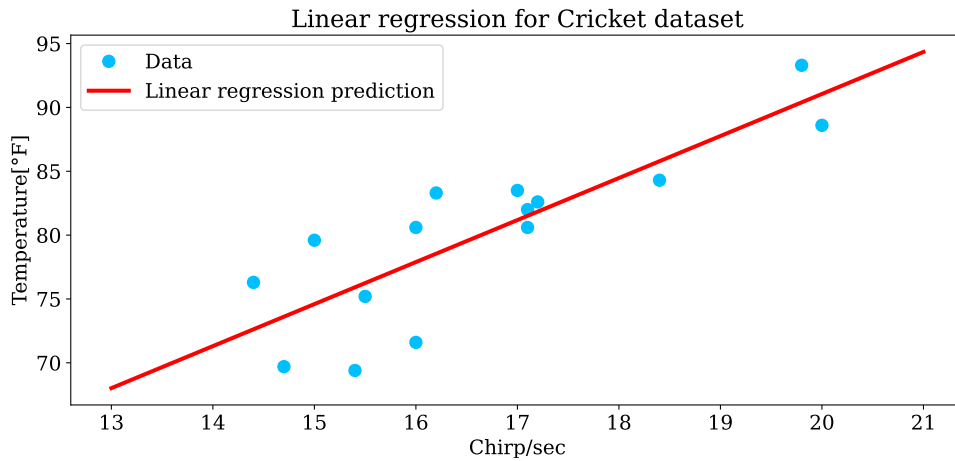


Fig. 1. Example of a linear regression model minimising the least squares loss.

Example: Logistic regression

Here, we aim to determine the function

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

$$x \mapsto f(x) = \frac{1}{1 + e^{-\theta^\top x + b}}, \quad \theta \in \mathbb{R}^M, b \in \mathbb{R}, \quad (2.4)$$

conditional to the observations

$$\mathcal{D} = \{(x_i, c_i)\}_{i=1}^N \subset \mathbb{R}^M \times \{0, 1\}. \quad (2.5)$$

The standard loss function for the classification problem is the cross entropy, given by:

$$J(\mathcal{D}, f) = -\frac{1}{N} \sum_{i=1}^N (c_i \log f(x_i) + (1 - c_i) \log(1 - f(x_i))) \quad (2.6)$$

$$= \frac{1}{N} \sum_{i=1}^N \left(\log(1 + e^{-\theta^\top x + b}) - y_i(-\theta^\top x + b) \right). \quad (2.7)$$

Figure 2 shows an example of a binary classification task minimising the cross entropy to separate data generated from two Gaussian distributions.

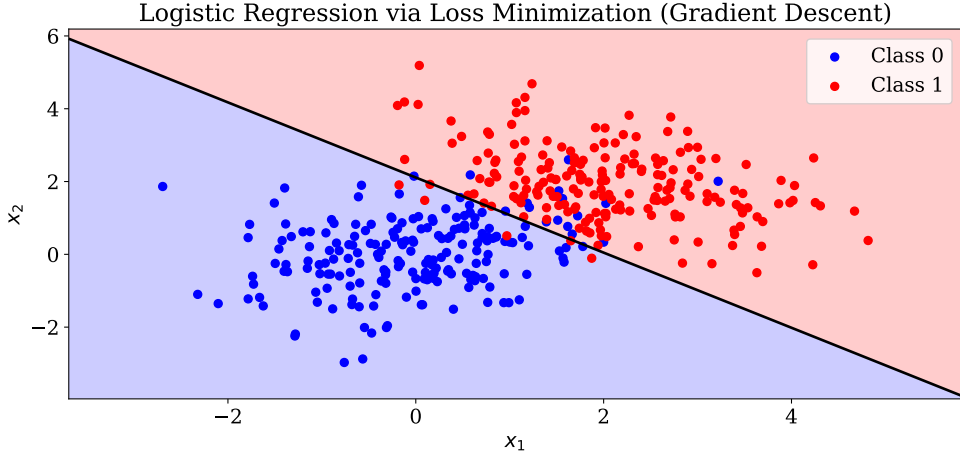


Fig. 2. Example of a logistic regression model to classify data from two Gaussians.

Example: Clustering (K -means)

Given a set of observations

$$\mathcal{D} = \{x_i\}_{i=1}^N \subset \mathbb{R}^M, \quad (2.8)$$

we aim to find cluster centres (or prototypes) $\mu_1, \mu_2, \dots, \mu_K$ and *assignment variables* $\{r_{ik}\}_{i,k=1}^{N,K}$, to minimise the following loss

$$J(\mathcal{D}, f) = \sum_{i=1}^N \sum_{k=1}^K r_{ik} \|x_i - \mu_k\|^2. \quad (2.9)$$

An example of a K -means model after the loss minimisation is shown on Figure 3

2.1 Terminology

We denote an optimisation problem as follows:

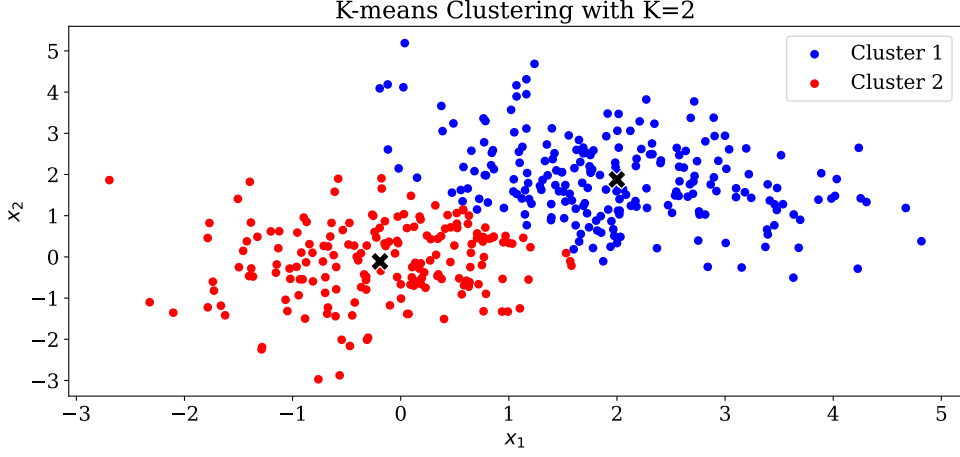


Fig. 3. Implementation of the K -means method to clustering data from two Gaussians. The figure shows the learnt centroids (black cross) and the colours correspond to the cluster assignments.

$$\min_{x \in \mathcal{X}} f(x), \quad \text{s.t.}, \quad g_i(x) \leq 0, \quad h_j(x) = 0, \quad i = 1, \dots, I, \quad j = 1, \dots, J. \quad (2.10)$$

We describe the components of this statement in detail:

- **Objective function:** The function $f : \mathcal{X} \rightarrow \mathbb{R}$ is the quantity to be minimised, with respect to x .
- **Optimisation variable:** Minimising f requires finding the value of x , such that $f(x)$ is minimum. This is also written as

$$x_\star = \arg \min_{x \in \mathcal{X}} f(x), \quad \text{s.t.}, \quad g_i(x) \leq 0, \quad h_j(x) = 0. \quad (2.11)$$

- **Restrictions:** These are denoted by the functions g_i and h_i above, which describe the requirements for the optimiser in the form of equalities and inequalities, respectively.
- **Feasible region:** This is the subset of the domain that complies with the restrictions, that is

$$C = \{x \in \mathcal{X}, \quad \text{s.t.}, \quad g_i(x) \leq 0, \quad h_j(x) = 0, \quad i = 1, \dots, I, \quad j = 1, \dots, J\}. \quad (2.12)$$

- **Local / global optima.** Values for the optimisation variable that solve the optimisation problem either locally or globally. More formally:

$$x_\star \text{ is a local optima} \iff \exists \lambda > 0 \quad \text{s.t.} \quad x_\star = \arg \min_{x \in \mathcal{X} \quad \text{s.t.} \quad \|x - x_\star\| \leq \lambda} f(x). \quad (2.13)$$

$$x_\star \text{ is a global optima} \iff x_\star = \arg \min_{x \in \mathcal{X}} f(x). \quad (2.14)$$

Example: Unique and non-unique closed form minima

In unconstrained optimisation, we have functions that have a unique global minimum and others that have more than one. For instance, Figure 4 shows $(x-1)^2 + (y+2)^2$ on the left, which has a unique minimiser on $(x, y) = (1, -2)$ and $(x^2-1)^2 + (y^2-1)^2$ on the right, where every global minimiser belongs to the set $\{(x, y) : x = \pm 1 \wedge y = \pm 1\}$.

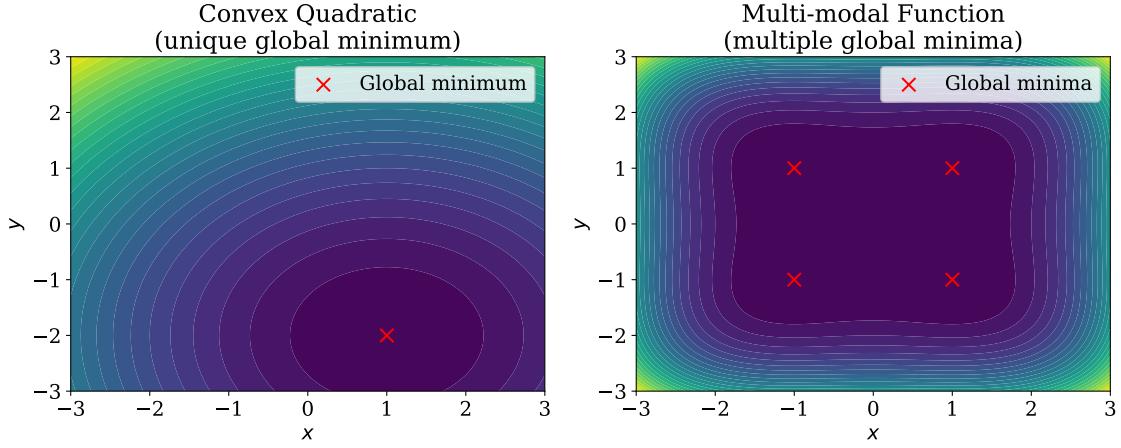


Fig. 4. Functions with a unique (left) and non-unique (right) global minima.

Interplay between constraints and local/global optima

On the other hand, we may be presented with functions such as $f(x, y) = \sin(x) * \sin(y) + 0.1(x^2 + y^2)$, which has a global minimum but also local minima. Figure 5 shows how different restrictions change the number and type of optima. In this case, when restricting the minimisation problem to the circle centred in $(1.5, 1)$ with radius 2, we observe a minimiser that is not the global minimum of the unconstrained problem (an unfeasible point) nor one of the local minima.

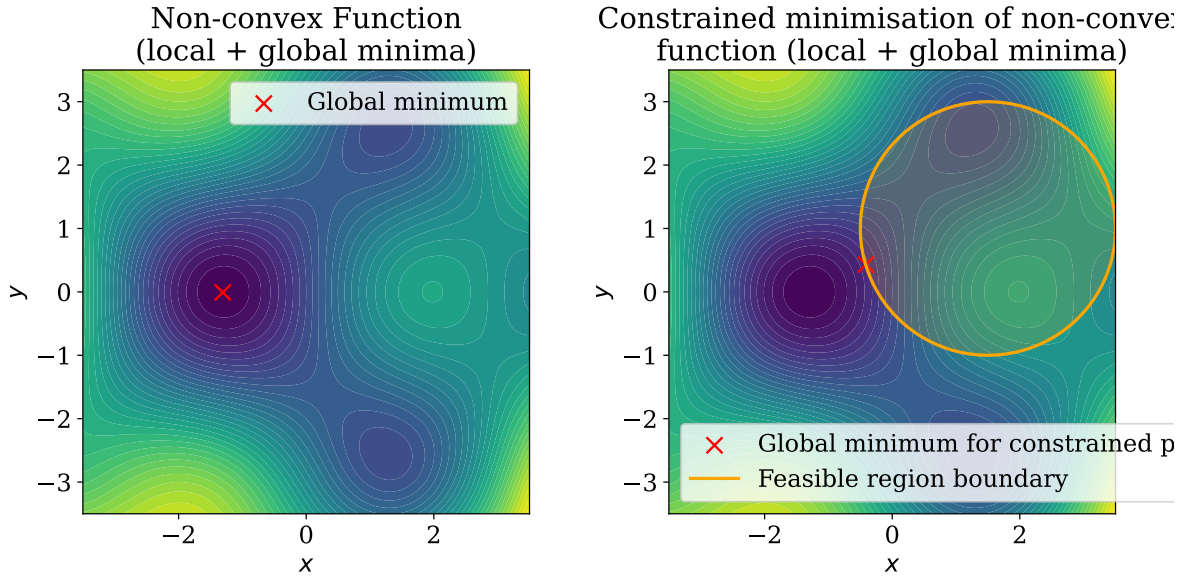


Fig. 5. Minima for constrained (left) and unconstrained (right) problems.

2.2 Continuous unconstrained optimisation

We will ignore constraints in this section, and focus on problems of the form

$$\theta \in \arg \min_{\theta \in \Theta} L(\theta). \quad (2.15)$$

We emphasise that if θ_* satisfies the above, then

$$\forall \theta \in \Theta, L(\theta_*) \leq L(\theta), \quad (2.16)$$

meaning that it is a **global** optimum. However, as this might be very hard to find, we are also interested in local optima, that is, θ_* such that

$$\exists \delta > 0 \text{ s.t. } \forall \theta \in \Theta \quad \|\theta - \theta_*\| < \delta \Rightarrow L(\theta_*) \leq L(\theta). \quad (2.17)$$

We now review the **optimality conditions**.

Assumption 2.1. The loss function L is twice differentiable.

Denoting $g(\theta) = \nabla_\theta L(\theta)$ and $H(\theta) = \nabla_\theta^2 L(\theta)$, we can state the following optimality conditions.

- **First order necessary condition:** If θ_* is a local minimum, then
 - $\nabla_\theta L(\theta_*) = 0$.
- **Second order necessary condition:** If θ_* is a local minimum, then
 - $\nabla_\theta L(\theta_*) = 0$
 - $\nabla_\theta^2 L(\theta_*)$ is positive semidefinite
- **Second order sufficient condition:** If θ_* is a local minimum if and only if
 - $\nabla_\theta L(\theta_*) = 0$
 - $\nabla_\theta^2 L(\theta_*)$ is positive definite

Example: different stationary points

Let us consider the function

$$\begin{aligned} f: \mathbb{R}^2 &\rightarrow \mathbb{R} \\ x &\mapsto f(x) = (p-1)x^2 + (p+1)y^2, \quad p \in \mathbb{R} \end{aligned} \quad (2.18)$$

Observe that

$$\nabla f = \begin{bmatrix} 2(p-1)x \\ 2(p+1)y \end{bmatrix}, \quad (2.19)$$

meaning that the only stationary points is $(x, y) = (0, 0)$. Furthermore,

$$\nabla^2 f = \begin{bmatrix} 2(p-1) & 0 \\ 0 & 2(p+1) \end{bmatrix}, \quad (2.20)$$

where we have 3 possible cases:

- $p > 1$: The stationary point is a minimum
- $-1 < p < 1$: The stationary point is a *saddle point*
- $p < -1$: The stationary point is a maximum

Figure 6 shows the function behaviour for different p . What happens when $|p| = 1$?

2.3 Convex optimisation

This setting is defined by having a convex objective function and a convex feasible region. Critically, in the setting of convex optimisation a local minimum (according to the first/second order conditions presented above) is a global minimum. We next formally provide the relevant definitions.

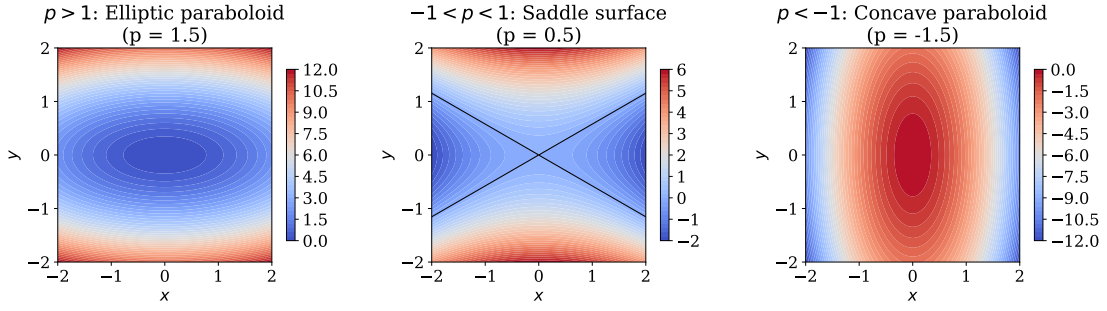


Fig. 6. Different types of critical points for $f(x, y) = (p - 1)x^2 + (p + 1)y^2$.

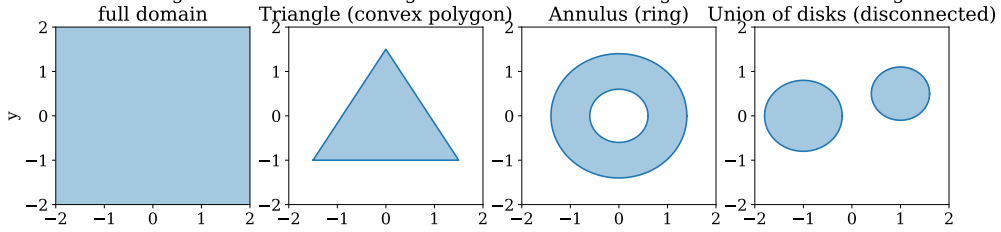


Fig. 7. Examples of convex and non-convex sets, highlighted in blue.

Definition 2.1 (Convex set). \mathcal{S} is a convex set if $\forall x, x' \in \mathcal{S}$, we have:

$$\lambda x + (1 - \lambda)x' \in \mathcal{S}, \quad \forall \lambda \in [0, 1]. \quad (2.21)$$

For illustration, Figure 7 shows two convex and two non-convex sets.

Definition 2.2 (Epigraph of a function). The epigraph of a function $f : \mathcal{X} \rightarrow \mathbb{R}$ is the set defined by the region above the graph of the function, that is,

$$\text{epi}(f) = \{ (x, t) \in \mathcal{X} \times \mathbb{R} \mid f(x) \leq t \}. \quad (2.22)$$

Definition 2.3 (Convex function). f is a convex function if its epigraph is convex. Equivalently, f is convex if it is supported on a convex set and $\forall x, x' \in \mathcal{X}$

$$f(\lambda x + (1 - \lambda)x') \leq \lambda f(x) + (1 - \lambda)f(x'), \quad \forall \lambda \in [0, 1]. \quad (2.23)$$

Furthermore, if the inequality is strict, we say that the function is **strictly convex**.

Example: Convex functions (in 1D)

The following are convex functions from \mathbb{R} to \mathbb{R} :

- $f(x) = x^2$
- $f(x) = e^{ax}$, $a \in \mathbb{R}$
- $f(x) = -\log x$
- $f(x) = x^a$, $a > 1$, $x > 0$
- $f(x) = |x|^a$, $a \geq 1$
- $f(x) = x \log x$, $x > 0$

Figure 8 shows the epigraphs for these functions.

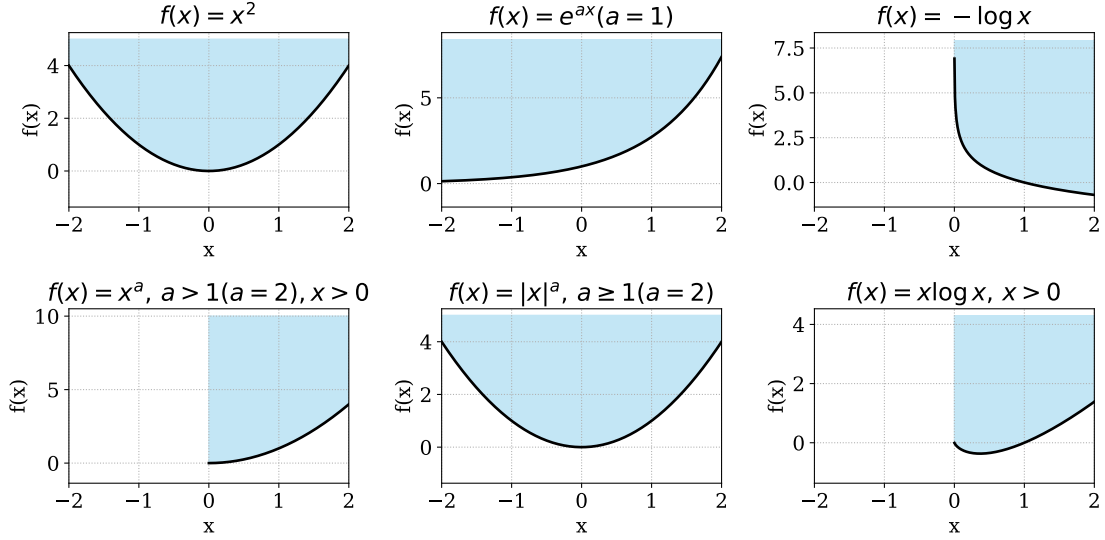


Fig. 8. Convex 1D functions with their epigraph in light blue.

We now review some important results in convex optimisation

Proposition 2.1. Consider $f : \mathcal{X} \subset \mathbb{R} \rightarrow \mathbb{R}$ differentiable. We have that if $f'(x) \geq 0 \forall x \in \mathbb{R}$, f is non-decreasing

Proof. By the fundamental theorem of calculus, we have that for $a, b \in \mathbb{R}, a < b$,

$$f(b) - f(a) = \int_a^b f'(x) dx, \quad (2.24)$$

since $f'(x) \geq 0, \forall x \in [a, b]$, we have $\int_a^b f'(x) dx \geq 0$, therefore $f(b) \geq f(a)$, which means that f is non-decreasing. ■

Proposition 2.2. Consider $f : \mathcal{X} \subset \mathbb{R}^d \rightarrow \mathbb{R}$ differentiable. The direction of maximum growth of f at x_0 is along its gradient $\nabla f(x_0)$

Proof. Let us consider $x' = x_0 + \rho u$, where $u \in \mathcal{X}, \|u\| = 1$, and $\rho > 0$ is a small constant. We find the maximum growth direction by maximising $f(x') - f(x_0)$ with respect to u . We consider the Taylor expansion

$$f(x') = f(x_0) + \nabla f(x_0) \rho u + \mathcal{O}(\rho^2), \quad (2.25)$$

and thus conclude that $f(x') - f(x_0) \simeq \nabla f(x_0) \rho u$, meaning that the maximum growth can be achieved by choosing u parallel to $\nabla f(x_0)$. That is, $\nabla f(x_0)$ is the direction of maximum growth for f at x_0 . ■

Teorema 2.1. Suppose $f : \mathcal{X} \subset \mathbb{R}^d \rightarrow \mathbb{R}$ twice differentiable, then f is convex if and only if ∇^2 is positive semi definite.

Proof. We consider $d = 1$. Using the FTC,

$$f'(b) - f'(a) = \int_a^b f''(x) dx \geq 0, \quad (2.26)$$

which implies that f' is non-decreasing. Therefore (using FTC again),

$$f(b) - f(a) = \int_a^b f'(x) dx \geq (b - a) f'(a), \quad (2.27)$$

equivalently,

$$f(b) \geq f(a) + (b - a) f'(a), \quad (2.28)$$

meaning that the function f is *always above its tangent*. Evaluating (2.28) for (a, z) and (b, z) , where $z = (1 - t)a + tb$, we have

$$f(z) \geq f(a) + (z - a)f'(a) \quad (2.29)$$

$$f(z) \geq f(b) + (z - b)f'(b). \quad (2.30)$$

Then, multiplying the above equations by $(1 - t)$ and t respectively and summing them, we obtain:

$$f(z) \geq (1 - t)f(a) + tf(b) + (1 - t)(tb - ta)f'(a) + t[(1 - t)a - (1 - t)b]f'(b) \quad (2.31)$$

$$= (1 - t)f(a) + tf(b) + (1 - t)t(b - a)[f'(a) - f'(b)] \quad (2.32)$$

$$\geq (1 - t)f(a) + tf(b). \quad (2.33)$$

This concludes the proof. Discuss in class how to extend this to arbitrary dimension $d > 1$. ■

Example: Explore some functions

[TODO: Choose some functions, compute the derivative and Hessian, analyse them]

2.4 First order methods

In general, finding a minimum by setting $\nabla f(x) = 0$ and solving for x is not possible. For that reason, we will consider iterative methods based on gradients. The idea here is to go *downhill* following the gradient towards the minimum (ignoring the curvature information for now).

We will specify a starting point x_0 and calculate

$$x_{t+1} = x_t + \eta_t d_t, \quad (2.34)$$

where η_t is a *step size* and d_t is a *descent direction*, such as $-\nabla f$. Here, the subindex \cdot_t represents the iteration number (starting from iteration $t = 0$). We iterate until convergence, that is, until the elements in the sequence $x_t, x_{t+1}, x_{t+2}, \dots$ become constant (or very similar). If convergence is achieved, we will assume the minimum has been found.

Note that there are several *descent directions*, that is, directions d_t such that

$$L(x_t + \eta_t d_t) \leq L(x_t). \quad (2.35)$$

In fact, as long as $d_t^\top \nabla f \leq 0$, d_t is a descent direction. Clearly, choosing $d_t = -\nabla f(x_t)$ is the *steepest descent direction*.

2.4.1 Role of the step size

The step size η_t is also known as *learning rate*. Furthermore, we refer to the set $\{\eta_1, \eta_2, \dots\}$ as the learning rate schedule. We will usually consider a constant learning rate, that is, $\eta_t = \eta, \forall t \in \mathbb{N}$. Though this is the simplest choice, there are some concerns to it: if η is too large, the iteration may fail to converge; whereas if it is too small, it may not converge at all.

Example: convergence for a parabola

Let us consider the function

$$J = (\theta - 3)^2. \quad (2.36)$$

In Figure 9 we show how the steepest descent converges/diverges for different learning rates.

The learning rate is usually tuned based on heuristics. When implementing the rule

$$x_{t+1} = x_t + \eta \nabla f(x_t), \quad (2.37)$$

we will usually set $\eta < \|\nabla f\|^{-1}$, as this will result in a stable autoregressive system for the sequence x_t (discuss this in class).

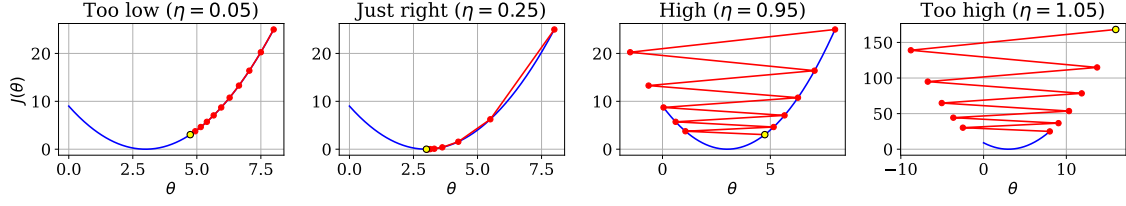


Fig. 9. Gradient based optimisation with different learning rates.

2.4.2 Momentum

In higher dimensions, we want to move faster in some directions and slower in others, depending on the value of the gradient in each coordinate. This can be achieved by:

$$m_t = \beta m_{t-1} + \nabla f(x_{t-1}) \quad (2.38)$$

$$x_t = x_{t-1} - \nabla_t m_t, \quad (2.39)$$

where m_t is a smoothed version of the gradient, and $\beta \in [0, 1]$ is a design (memory) parameter. This way, previous values of the gradient have effect on future updates: if a particular coordinate of the gradient is consistently large, then that coordinate receives updates of a higher magnitude. This is particularly useful when the evaluation of the gradient is noisy. Figure 17 shows how learning changes using momentum, with the same learning rates as the previous example.

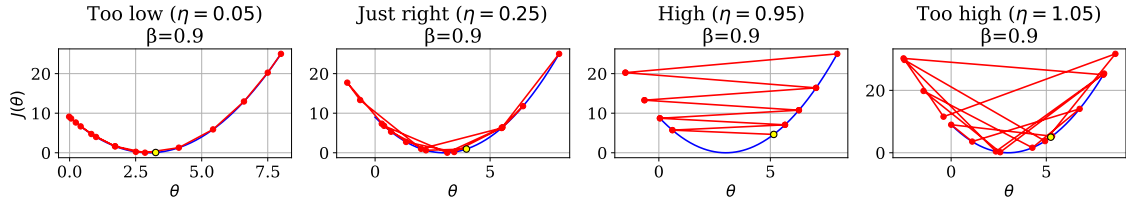


Fig. 10. Gradient based optimisation with momentum.

2.4.3 Newton method

Newton's method is

$$x_{t+1} = x_t - \eta_t H_t^{-1} \nabla f(x_t), \quad (2.40)$$

where recall that $H_t = \nabla^2 f(x_t)$ denotes the Hessian of f at x_t . This update follows from considering the second order approximation of the loss function around the current point, that is:

$$L(x) \simeq L(x_t) + (\nabla f(x_t))^\top (x - x_t) + \frac{1}{2} (x - x_t)^\top H_t (x - x_t), \quad (2.41)$$

the minimum of which is given by

$$x_\star = x_t - H_t^{-1} \nabla L(x_t), \quad (2.42)$$

where the learning rate can also be used to accelerate (or de-accelerate) convergence.

Example: convergence for a parabola (2)

Let us consider the same parabolic function as before. This time we will use the Newton method presented above.

As shown in Figure 11, this method converges in one step (for a quadratic function). That is, the update corresponds to the closed form solution of the minimisation problem.

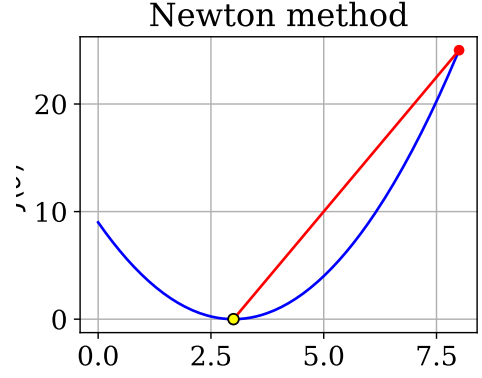


Fig. 11. One-step parabola minimisation with the Newton method.

2.5 Stochastic gradient descent

We now consider stochastic optimisation, where

$$L(x) = \mathbb{E}_{q(z)} L(x, z), \quad (2.43)$$

that is, when the loss function is random and we aim to minimise its expected value. For instance, in linear regression we have $L(\theta) = \mathbb{E}_{q(y|x)}(y - \theta^\top x)$.

In practice, we are unable to compute this expectation since the law $q(z)$ is unknown. However, since we usually have samples of q , we can do a sample approximation of the expectation. In fact, we will consider

$$L(\theta) = \frac{1}{N} \sum_{i=1}^n L(x, z_i). \quad (2.44)$$

In the linear regression example, this would be $L(\theta) = \frac{1}{N} \sum_{i=1}^n (y_i - \theta^\top x_i)^2$.

The gradient is then also approximated using a batch of, say, B samples. That is,

$$\nabla L(\theta) \simeq \frac{1}{N} \sum_{i=1}^B \nabla L(x, z_i). \quad (2.45)$$

Example: random loss function for linear regression

Consider a toy example of linear regression where $q(y|x) = \mathcal{N}(y; x, \sigma^2)$ for a given $\sigma > 0$. Consequently, we can compare the true expectation (equal to the variance in this case) with the empirical approximation for increasing values of N .

Figure 12 shows how the empirical loss changes as we consider more and more samples.

Recall that, in general, we will consider parameter updates of the form

$$\theta_{t+1} = \theta_t - M_t^{-1} g_t, \quad (2.46)$$

where M_t is an estimate of the magnitudes of the coordinates of the gradient g_t , as this ensures convergence (or prevents divergence). We next explore some different choices of this estimate.

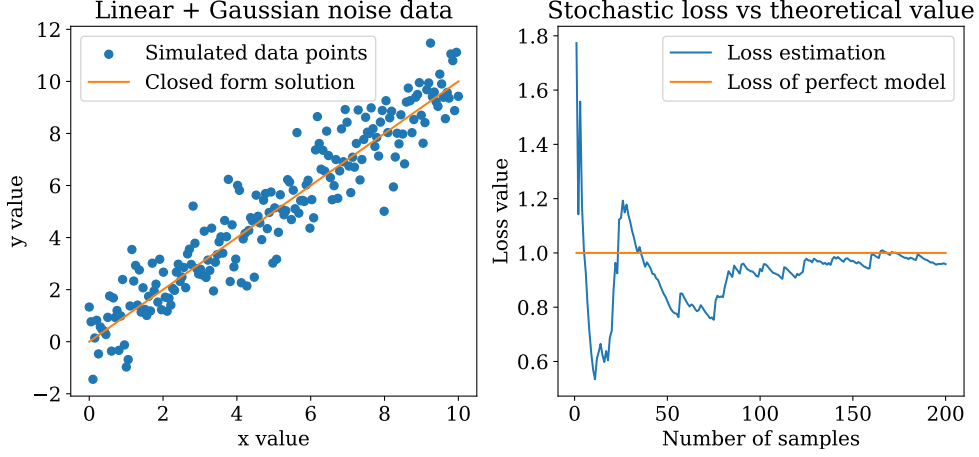


Fig. 12. Estimation of the true loss via the empirical value from data.

2.5.1 Adagrad

This considers the following iterative rule (d denotes the coordinate):

$$\theta_{t+1,d} = \theta_{t,d} - \eta_t \frac{1}{\sqrt{s_{t,d} + \epsilon}} g_{t,d}, \quad (2.47)$$

where $s_{t,d} = \sum_{i=1}^t g_{i,d}^2$. This can also be expressed in vector format as

$$\Delta\theta_t = -\eta_t \frac{1}{s_t} g_t. \quad (2.48)$$

The aim of Adagrad is to control the magnitude of the step size for each coordinate of the parameter independently, based on a rolling estimate using previous gradient evaluations. Notice that the computation of $s_{t,d}$ assigns the same importance to all gradient evaluations in time.

2.5.2 RMSProp

Adagrad might fail to represent the current gradient magnitude by not emphasising current evaluations of the gradient. This can be solved by replacing the sum in the computation of $s_{t,d}$ by a moving average. That is,

$$s_{t+1,d} = \beta s_{t,d} + (1 - \beta) g_{t,d}^2, \quad (2.49)$$

where the *forgetting factor* β is usually chosen close to 0.9. Then, the update rule also follows

$$\Delta\theta_t = -\eta_t \frac{1}{s_t} g_t. \quad (2.50)$$

2.5.3 Adam

The de facto optimiser in deep learning is Adam (Adaptive Moment Estimation) proposed by (Kingma & Ba, 2015), which combines RMSProp and momentum. Adam's update follows:

$$m_t = \beta_m m_{t-1} + (1 - \beta_m) g_t \quad (2.51)$$

$$s_t = \beta_s s_{t-1} + (1 - \beta_s) g_t^2. \quad (2.52)$$

Then,

$$\theta_t = \theta_{t-1} - \eta_t \frac{m_t}{\sqrt{s_t + \epsilon}}. \quad (2.53)$$

Example: different optimisers

We will consider the function

$$f(x, y) = 0.1x^2 + y^2. \quad (2.54)$$

Notice that this function is convex, hence we would expect a reasonable gradient based method to converge. However, the speed and way in which they converge might be different in every case. This text function will illustrate this, since its curvature is much more pronounced along the y axis.

We visualise the effect of the different optimisers presented above in Figure 13. Notice that SGD converges more slowly than RMSprop and Adam. Adagrad, on the other hand, slows down too early.

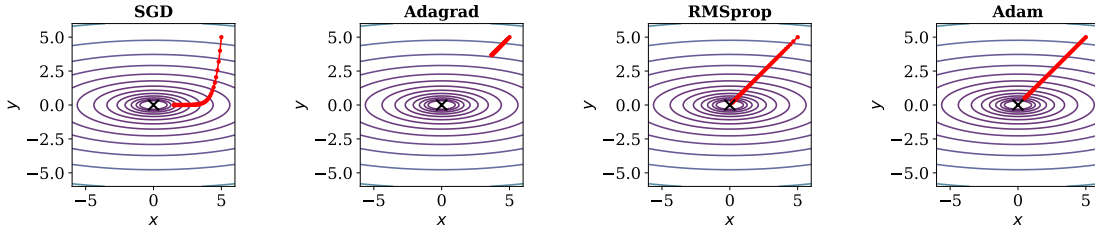


Fig. 13. Convergence comparison for different optimisers.

2.6 Training and Regularisation of Neural Networks

2.6.1 Network Architecture

The **architecture** of a neural network refers to its overall structure: number of layers, neurons per layer, connectivity, and activation functions.

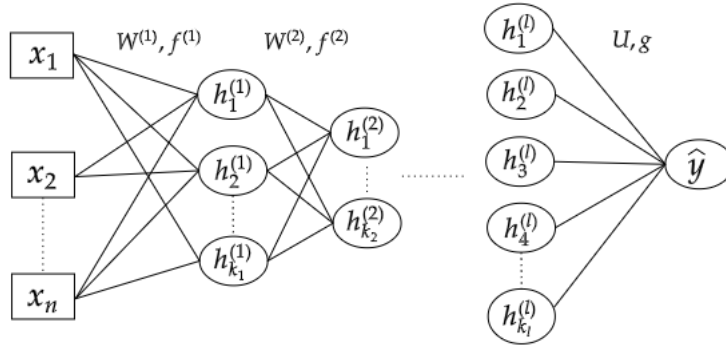


Fig. 14. A neural network of depth l .

We adopt the following notation (as illustrated in Fig. 14):

$$h^{(k)} = f^{(k)}(h^{(k-1)}W^{(k)} + b^{(k)}), \quad k \in \{1, \dots, l\}, \quad h^{(0)} = x, \quad (2.55)$$

$$\hat{y} = g(h^{(l)}U + c), \quad (2.56)$$

where f is the activation function for all (hidden) neurons, and g is the activation function of the **output neurons**. Furthermore, $h^{(k)}$ denotes the output of the k -th layer.

2.6.2 Cost Function and Output Units

Unlike linear models, the use of nonlinear activations renders the cost function generally non-convex, so gradient descent does not guarantee convergence to a global optimum. Although least squares may be used, in practice one minimises the negative log-likelihood of a probabilistic model $p(y|\mathbf{x}; \boldsymbol{\theta})$:

$$J(\boldsymbol{\theta}) = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{y}|\mathbf{x})]. \quad (2.57)$$

The choice of output unit determines the appropriate likelihood and loss:

1. **Linear:** $\hat{\mathbf{y}} = h^{(l)}U + c$

Models $\mathcal{N}(\mathbf{y}|\hat{\mathbf{y}}, \mathbf{I})$ and is equivalent to minimising MSE; suited for regression.

2. **Sigmoid:** $\hat{y} = \text{sig}(h^{(l)}U + c)$

Used for binary classification with $p(y|\mathbf{x}) = \text{Bernoulli}(y|p)$ and $\hat{y} = p(y = 1|\mathbf{x})$.

3. **Softmax:** $\hat{\mathbf{y}} = \text{softmax}(h^{(l)}U + c)$

For multiclass problems, with

$$\text{softmax}(z)_i = \frac{e^{z_i}}{\sum_j e^{z_j}}.$$

Recall that, in practice, since p_{data} in eq. (2.57) is unknown, we optimise an empirical estimate of that loss given by

$$J(\boldsymbol{\theta}) = -\sum_i \log p_{\text{model}}(\mathbf{y}_i|\mathbf{x}_i). \quad (2.58)$$

The rest of this section is devoted to the calculation of this loss and its gradient.

2.6.3 Forward Propagation

In a feedforward network, information flows from input \mathbf{x} to output $\hat{\mathbf{y}}$. During training, this process computes the cost $J(\mathbf{X}, \boldsymbol{\theta})$.

Algorithm 1 Forward Propagation

Require: Depth l , weights $\mathbf{W}^{(k)}$, biases $\mathbf{b}^{(k)}$, input \mathbf{x} , output \mathbf{y} , and output parameters U, c, g .

- 1: $\mathbf{h}^{(0)} \leftarrow \mathbf{x}$
 - 2: **for** $k = 1, \dots, l$ **do**
 - 3: $\mathbf{u}^{(k)} \leftarrow \mathbf{h}^{(k-1)}\mathbf{W}^{(k)} + \mathbf{b}^{(k)}$
 - 4: $\mathbf{h}^{(k)} \leftarrow f^{(k)}(\mathbf{u}^{(k)})$
 - 5: $\hat{\mathbf{y}} \leftarrow g(\mathbf{h}^{(l)}U + c)$
 - 6: $J \leftarrow L(\hat{\mathbf{y}}, \mathbf{y})$
-

2.6.4 Backward Propagation – Preliminaries

The **backpropagation** algorithm propagates cost information backwards to compute gradients efficiently. In practice, this is achieved by means of **automatic differentiation**, which provides fast and scalable computation, thus enabling modern deep learning.

Assume training is performed in **mini-batches** $(\mathbf{x}^d)_{d=1}^N$, enabling efficient matrix operations on GPUs. For an MSE loss in regression:

$$J(\mathbf{X}, \boldsymbol{\theta}) = \frac{1}{2N} \sum_{d=1}^N (\hat{y}_d - y_d)^2,$$

where $J_d = (\hat{y}_d - y_d)^2$, our goal is to update all weights $w_{ij}^{(k)}$. Using the chain rule:

$$\frac{\partial J_d}{\partial w_{ij}^{(k)}} = \frac{\partial J_d}{\partial u_{dj}^{(k)}} \frac{\partial u_{dj}^{(k)}}{\partial w_{ij}^{(k)}},$$

where we define the **error term**

$$\delta_{dj}^{(k)} \equiv \frac{\partial J_d}{\partial u_{dj}^{(k)}},$$

and note that

$$\frac{\partial u_{dj}^{(k)}}{\partial w_{ij}^{(k)}} = h_{di}^{(k-1)}.$$

Hence,

$$\frac{\partial J_d}{\partial w_{ij}^{(k)}} = \delta_{dj}^{(k)} h_{di}^{(k-1)},$$

and the total gradient (ignoring the constant $1/N$ for now) is

$$\frac{\partial J}{\partial W^{(k)}} = (h^{(k-1)})^T @ \delta^{(k)}. \quad (2.59)$$

2.6.5 Backward Propagation – Hidden Layers

From the chain rule:

$$\delta_{dj}^{(k)} = \frac{\partial J_d}{\partial u_{dj}^{(k)}} = \sum_{a=1}^{k_{k+1}} \frac{\partial J_d}{\partial u_{da}^{(k+1)}} \frac{\partial u_{da}^{(k+1)}}{\partial u_{dj}^{(k)}} = \sum_{a=1}^{k_{k+1}} \delta_{da}^{(k+1)} \frac{\partial u_{da}^{(k+1)}}{\partial u_{dj}^{(k)}},$$

where, since $\frac{\partial u_{da}^{(k+1)}}{\partial u_{dj}^{(k)}} = w_{ja}^{(k+1)} f'(u_{dj}^{(k)})$, we have

$$\delta_{dj}^{(k)} = f'(u_{dj}^{(k)}) \sum_{a=1}^{k_{k+1}} w_{ja}^{(k+1)} \delta_{da}^{(k+1)}.$$

In matrix form:

$$\delta^{(k)} = f'(u^{(k)}) * (\delta^{(k+1)} @ (W^{(k+1)})^T). \quad (2.60)$$

2.6.6 Backward Propagation – Output Layer

For a regression problem with linear output and MSE loss:

$$\delta_{d1}^{(l)} = (\hat{y}_d - y_d)(\hat{y}_d)' = (\hat{y}_d - y_d),$$

and with the normalisation term:

$$\delta^{(l)} = \frac{1}{N}(\hat{y} - y). \quad (2.61)$$

Furthermore,

$$\frac{\partial J}{\partial b^{(k)}} = \text{Sum}_1(\delta^{(k)}), \quad \frac{\partial J}{\partial U} = (h^{(l)})^T @ \delta^{(l)}, \quad \frac{\partial J}{\partial c} = \text{Sum}_1(\delta^{(l)}), \quad (2.62)$$

where $\text{Sum}_1(A)$ denotes summation over the first axis.

Algorithm 2 Backward Propagation

Require: Learning rate λ .

- 1: Perform forward propagation, store $(\hat{y}), (u^{(k)}), (h^{(k)})$.
 - 2: Compute output error $\delta^{(l)}$ according to the output unit.
 - 3: Update $U \leftarrow U - \lambda \frac{\partial J}{\partial U}, \quad c \leftarrow c - \lambda \frac{\partial J}{\partial c}$.
 - 4: **for** $k = l, \dots, 1$ **do**
 - 5: Compute $\delta^{(k)}$ using the hidden-layer recursion.
 - 6: Evaluate $\frac{\partial J}{\partial W^{(k)}}, \frac{\partial J}{\partial b^{(k)}}$.
 - 7: Update $W^{(k)} \leftarrow W^{(k)} - \lambda \frac{\partial J}{\partial W^{(k)}}$.
 - 8: Update $b^{(k)} \leftarrow b^{(k)} - \lambda \frac{\partial J}{\partial b^{(k)}}$.
-

Remark 2.1. The algorithm is applied to each mini-batch for several **epochs**. Increasing epochs generally decreases training error, but excessive training may degrade generalisation due to **overfitting**.

2.6.7 Regularisation

Neural networks are often large models; controlling their complexity is crucial. **Regularisation** techniques aim to reduce the generalisation error while retaining sufficient model capacity.

L2 Regularisation. Adds a quadratic penalty on weights:

$$\tilde{J}(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \frac{\alpha}{2} \|\boldsymbol{\theta}\|_2^2.$$

The hyperparameter α controls the regularisation strength. Usually, only weights are penalised, not biases. In practice this corresponds to *weight decay*:

$$W^{(k)} \leftarrow (1 - \beta)W^{(k)} - \lambda \frac{\partial J}{\partial W^{(k)}}, \quad \beta = \lambda\alpha.$$

3 Probability

3.1 Introduction

NB: in this chapter, we follow (Murphy, 2022; Deisenroth et al., 2020). For additional material see (Durrett, 2019; Bertsekas & Tsitsiklis, 2008; Grimmett & Stirzaker, 2020)

The field of probability studies, from a quantitative perspective, how *likely* an event is. Conceptually, and perhaps historically, there are two main interpretations of probability. The first one is **frequentist probability**, which relates to frequency of occurrence, and then applies only to events that can be repeated an infinite number of times, such as throwing a dice or flopping a coin. As a consequence, this approach to probability fails to assign a probability to events that are impossible to repeat, such as the average temperature of the Earth’s surface reaching an all-time maximum in the year 2025. A second interpretation is that of **Bayesian probability**, which represents uncertainty about the occurrence of an event. This uncertainty might come from different sources, such unknown features in the experiments (epistemological uncertainty) or random components (aleatoric uncertainty). In this case, events need not be repeatable be assigned with a probability.

A basic yet rigorous understanding of probability theory, definitions and main results is fundamental is central to the construction and applications of ML methods. This is because in ML we design, train and deploy mathematical models that aim to i) capture/quantify uncertainty, and ii) deal with noise-corrupted training data.

3.1.1 Definitions

To start studying probability, we will consider the outcome ω of a hypothetical experiment. For instance, this can be throwing a dice, where ω takes values $\{1, 2, 3, 4, 5, 6\}$. In this context, we define:

Definition 3.1 (Sample space). The set containing all the possible outcomes ω of an experiment is called sample space and is denoted by Ω . For the dice example, $\Omega = \{1, 2, 3, 4, 5, 6\}$.

Definition 3.2 (Event space). The set \mathcal{A} , referred to as event space, contains all possible subset of the sample space Ω . Therefore, each element $A \in \mathcal{A}$ represents a possible results of the experiment.

Definition 3.3 (Probability). The function

$$\mathbb{P} : \mathcal{A} \rightarrow [0, 1] \tag{3.1}$$

$$x \mapsto \mathbb{P}(x) \tag{3.2}$$

denotes the probability of the result of the experiment falling inside the elements of \mathcal{A} , that is, $\mathbb{P}(A) = \mathbb{P}(\omega \in A)$. The function \mathbb{P} needs to fulfil some standard properties such as $\mathbb{P}(\Omega) = 1$, $\mathbb{P}(\emptyset) = 0$, and $\mathbb{P}(A^c) = 1 - \mathbb{P}(A)$.

Discussion: \mathcal{A} vs Ω

Why is the probability defined over \mathcal{A} and not over Ω ? Elaborate over intuitive examples

Examples

[TODO: Consider basic examples (dice, coin, uniform, rain), and present the sample space, event space, and probability]

Definition 3.4. We refer to the triplet $(\Omega, \mathcal{A}, \mathbb{P})$ as **probability space**.

3.1.2 Basic properties

The joint probability of events A and B is denoted by

$$\mathbb{P}(A \wedge B) = \mathbb{P}(A \cap B) = \mathbb{P}(A, B). \quad (3.3)$$

This follows from the fact that if $\omega \in A$ and $\omega \in B$, then, $\omega \in A \cap B$.

The **conditional probability** of A occurring, given that the event B occurred, is denoted by

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A, B)}{\mathbb{P}(B)}, \quad (3.4)$$

which is only valid when $\mathbb{P}(B) > 0$.

Additionally, we say that events A and B are **independent** iff $\mathbb{P}(A, B) = \mathbb{P}(A)\mathbb{P}(B)$. Observe that this implies that

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A, B)}{\mathbb{P}(B)} = \frac{\mathbb{P}(A)\mathbb{P}(B)}{\mathbb{P}(B)} = \mathbb{P}(A), \quad (3.5)$$

meaning that when A and B are independent, the latter provides no **information** for A .

Lastly, the probability of intersection, i.e., the probability of the events $\omega \in A$ or $\omega \in B$, is

$$\mathbb{P}(A \vee B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \wedge B). \quad (3.6)$$

3.2 Random variables

In general, other than basic toy examples, we do not refer to the underlying experiment and its sample/event spaces. We instead consider *quantities of interest* that result from the outcome of the experiment. Therefore, let us consider a map from the sample space to a target space \mathcal{T} , e.g., $\mathcal{T} = \mathbb{R}$, $\mathcal{T} = \mathbb{N}$ or $\mathcal{T} = \{1, 2, 3, \dots, N\}$.

Definition 3.5 (Random variable). A function

$$X : \Omega \rightarrow \mathcal{T} \quad (3.7)$$

$$\omega \mapsto X(\omega) \quad (3.8)$$

is referred to as random variable. In general, we will denote the function as X and its value as $X(\omega) = x$, ignoring the explicit dependence on ω .

Remark 3.1. From now on, we will consider the outcome of the experiment (living in the sample space) as the value of the RV; this aims towards a streamlined setup, avoiding the need of a map from Ω to \mathcal{T} . Accordingly, the event space \mathcal{A} and the probability \mathbb{P} correspond to the outcomes (values) of X . We will denote the sample space by \mathcal{X} .

3.2.1 Discrete RVs

If the sample space is finite or countably infinite, we will say that the RV is discrete. In this case, we denote the probability of the event $X = x$ by $\mathbb{P}(X = x)$.

Definition 3.6 (Probability mass function (pmf)). For a discrete RV $X \in \mathcal{X}$, the function

$$p_X : \mathcal{X} \rightarrow [0, 1] \quad (3.9)$$

$$x \mapsto p_X(x) = \mathbb{P}(X = x) \quad (3.10)$$

denotes the probability of the event where the RV X takes the value $x \in \mathcal{X}$. Evidently,

$$\sum_{x \in \mathcal{X}} p_X(x) = 1. \quad (3.11)$$

Discrete RV: time to first heads

Consider an infinite sequence of independent coin flips, each taking values

$$X_i = \begin{cases} 1, & \text{(Heads)} \\ 0, & \text{(Tails)}, \end{cases} \quad \text{with } \mathbb{P}(X_i = 1) = p, \mathbb{P}(X_i = 0) = 1 - p.$$

Define the discrete random variable

$$T = \min\{k \geq 1 : X_k = 1\},$$

the *time of the first heads*. Then T takes values in $\{1, 2, 3, \dots\}$ and

$$\mathbb{P}(T = k) = (1 - p)^{k-1}p, \quad k = 1, 2, \dots$$

This is the *geometric distribution*. Its cdf is

$$F_T(k) = \mathbb{P}(T \leq k) = 1 - (1 - p)^k, \quad k = 1, 2, \dots$$

and

$$\mathbb{E}[T] = \frac{1}{p}, \quad \mathbb{V}(T) = \frac{1 - p}{p^2}.$$

For the particular case of $p = 1/2$, we have

$$\mathbb{P}(T = k) = \frac{1}{2^k}, \quad k = 1, 2, \dots$$

Where it holds that $\sum_{k=1}^{\infty} \mathbb{P}(T = k) = \sum_{k=1}^{\infty} \frac{1}{2^k} = 1$.

3.2.2 Continuous RVs

If $\mathcal{X} = \mathbb{R}^d$, with $d \geq 1$, or any other infinitely-uncountable space, we say that the RV is continuous. In this case, defining a pmf as in Def. 3.6 is not possible: since there is an uncountable number of symbols in the sample space \mathcal{X} , it is not possible to assign a probability strictly greater than zero to each of them, while having a total probability mass equal to one. Therefore, we use the event space and assign probabilities **to intervals** rather than particular values.

Definition 3.7 (Cumulative distribution function (cdf)). For a continuous RV $X \in \mathcal{X}$, we define the probability of X to be less or equal than a given value $x \in \mathcal{X}$ by

$$P_X : \mathcal{X} \rightarrow [0, 1] \tag{3.12}$$

$$x \mapsto P_X(x) = \mathbb{P}(X \leq x). \tag{3.13}$$

Observe that this also allows to denote the probability of X lying in a bounded interval, that is,

$$\mathbb{P}(a \leq X \leq b) = P_X(b) - P_X(a). \tag{3.14}$$

The cdf is monotonically non-decreasing by construction, and, when $\mathcal{X} = \mathbb{R}$, we have

$$\lim_{x \rightarrow -\infty} P_X(x) = 0 \tag{3.15}$$

$$\lim_{x \rightarrow \infty} P_X(x) = 1. \tag{3.16}$$

The idea to assign probabilities to intervals, suggest that there exists a **density of probability**, that is, an amount of probability mass per unit of length (or *measure*) of the sample space. This is formalised via the following definition.

Definition 3.8 (Probability density function (pdf)). For a continuous RV $X \in \mathcal{X}$, the probability density function of X is given by the function

$$p_X : \mathcal{X} \rightarrow \mathbb{R} \quad (3.17)$$

$$x \mapsto p_X(x) = \frac{d}{dx} P_X(x). \quad (3.18)$$

This is possible when P_X is differentiable.

Knowing the pdf, we can express:

$$\mathbb{P}(a \leq X \leq b) = \int_a^b p_X(x) dx = P_X(b) - P_X(a), \quad (3.19)$$

which is a direct particular case of the fundamental theorem of Calculus.

Also, for Δx small, we have

$$\mathbb{P}(x \leq X \leq x + \Delta x) = \int_x^{x+\Delta x} p_X(x) dx \simeq p(x) \Delta x. \quad (3.20)$$

Lastly, if the cdf is **strictly monotonically increasing**, its inverse exists and it is known as the **quantile function**.

Discussion

Why is the quantile function useful?

Gaussian distribution

In Figure 15 we can see the probability density function and the cumulative density function of the one dimensional Gaussian distribution. We will define this distribution in Section 3.3.3.

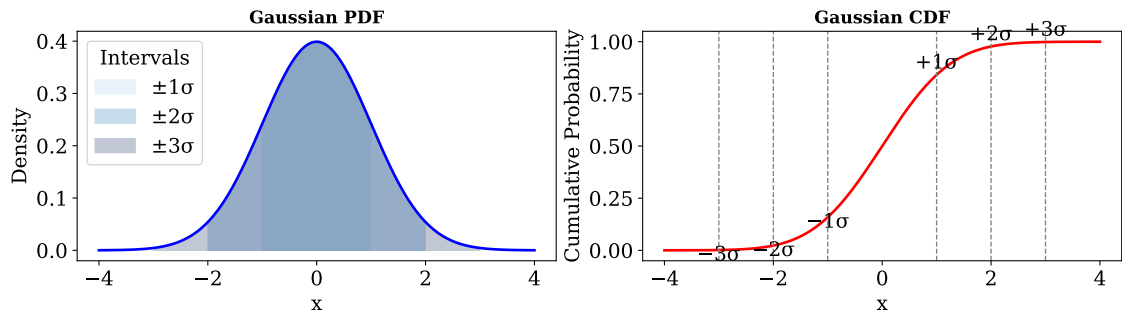


Fig. 15. Probability and cumulative density function.

3.2.3 Properties and identities

Probability can be thought of as an extension of logical reasoning, and the following properties are consistent with such extension. Let us first consider, akin to the definition of event probabilities, the following results.

Both for discrete (pmf) and discrete (pdf) RVs, we can denote $p(x, y)$ the joint pdf/pmf for RVs $X \in \mathcal{X}$ and $Y \in \mathcal{Y}$. This can be seen as the probability of the *intersection* event, where $X = x$ and $Y = y$.

Furthermore, we can write

$$p_X(x) = \sum_{y \in \mathcal{Y}} p(x, y) \quad (\text{discrete}) \quad (3.21)$$

$$p_X(x) = \int_{\mathcal{Y}} p(x, y) dy \quad (\text{continuous}). \quad (3.22)$$

In the ML community, this is known as **sum rule**, since it allows us to express the density $p(y)$ from a joint density by summing over all possible values.

The so called **product rule** is the decomposition of the joint law as follows

$$p(x, y) = p(y|x)p(x), \quad (3.23)$$

which allows us to decompose the joint pmf/pdf using the law of the conditional event $X = x$ given that $Y = y$, and the marginal law of Y . A key consequence of the above expression is that, since the factorisation works in both ways, we have

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}, \quad (3.24)$$

which is known as the **Bayes theorem**. This is a fundamental result in probabilistic ML, since we usually observe measurements y and want to infer parameters or latent variables x .

An expression combining the sum and the product rules is the so called **law of total probability**, where the sum rule is expressed via conditional probabilities. That is,

$$p_X(x) = \sum_{y \in \mathcal{Y}} p(x|y)p(y) \quad (\text{discrete}) \quad (3.25)$$

$$p_X(x) = \int_{\mathcal{Y}} p(x|y)p(y)dy \quad (\text{continuous}). \quad (3.26)$$

This operation is usually called **marginalisation**, or **disintegration** (of y). Here, we say that y has been **integrated out**.

Example: The Monty Hall problem

You are on a game show with three doors: behind one is a car (C), behind the others goats (G). You choose door 1. The host, who knows what's behind each door, opens door 3 revealing a goat. Should you switch to door 2?

Let A_i be the event that the car is behind door i , and B the event that the host opens door 3.

$$P(A_1) = P(A_2) = P(A_3) = \frac{1}{3}.$$

Given that you chose door 1, we have:

$$P(B | A_1) = \frac{1}{2}, \quad P(B | A_2) = 1, \quad P(B | A_3) = 0.$$

Bayes' theorem states:

$$P(A_1 | B) = \frac{P(B | A_1)P(A_1)}{P(B)} \quad (3.27)$$

$$= \frac{P(B | A_1)P(A_1)}{P(B | A_1)P(A_1) + P(B | A_2)P(A_2) + P(B | A_3)P(A_3)} \quad (3.28)$$

$$= \frac{(1/2)(1/3)}{(1/2)(1/3) + (1)(1/3)} \quad (3.29)$$

$$= \frac{1}{3}. \quad (3.30)$$

Hence

$$P(A_2 | B) = \frac{2}{3}.$$

Conclusion: Switching doubles your chance of winning.

3.2.4 Moments

Definition 3.9 (Mean). The mean, or expected value, of an RV $X \in \mathcal{X}$, often denoted by $\mu = \mu_X$, is defined as

$$\mathbb{E}[X] = \int_{\mathcal{X}} xp_X(x)dx \quad (\text{continuous}) \quad (3.31)$$

$$\mathbb{E}[X] = \sum_{x \in \mathcal{X}} xp_X(x) \quad (\text{discrete}), \quad (3.32)$$

$$(3.33)$$

where in the continuous case we have assumed the existence of a pdf.

Observe that $\mathbb{E}[\cdot]$ is a linear operator, meaning that

$$\mathbb{E}[aX + b] = a\mathbb{E}[X] + b \quad a, b \in \mathbb{R} \quad (3.34)$$

$$\mathbb{E}\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n \mathbb{E}[X_i], \quad (3.35)$$

$$(3.36)$$

where X_1, X_2, \dots, X_n is a collection of arbitrary RVs with finite mean.

The definition above also allows for calculating the mean of a transformation of the RV, this is because for a function $f : \mathcal{X} \rightarrow \mathbb{R}$, the quantity $f(X)$ is also an RV, and thus its mean is $\mathbb{E}[f(X)] = \int_{\mathcal{X}} f(x)p_X(x)dx$. In particular, this enables the following definition.

Definition 3.10 (Variance). The variance measures the **spread** of a distribution (with respect to its mean) and it is given by

$$\mathbb{V}[X] = \mathbb{E}[(X - \mu_X)^2] \quad (3.37)$$

$$= \int_{\mathcal{X}} (x - \mu_X)^2 p_X(x)dx \quad (3.38)$$

$$= \int_{\mathcal{X}} (x^2 - 2x\mu_X + \mu_X^2) p_X(x)dx \quad (3.39)$$

$$= \mathbb{E}[X^2] - \mu_X^2, \quad (3.40)$$

where the discrete case is obtained similarly.

Definition 3.11 (Covariance). The covariance measures the **joint variability** of two random variables X and Y , and it is given by

$$\text{Cov}[X, Y] = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)] \quad (3.41)$$

$$= \int_{\mathcal{X} \times \mathcal{Y}} (x - \mu_X)(y - \mu_Y) p_{X,Y}(x, y) dx dy \quad (3.42)$$

$$= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]. \quad (3.43)$$

A positive covariance indicates that X and Y tend to increase together, while a negative value indicates that when one increases, the other tends to decrease.

Additionally, a related quantity is the **standard deviation**, given by

$$\text{std}[X] = \sqrt{\mathbb{V}[X]}. \quad (3.44)$$

A list of relevant properties of the variance are:

- **Definition:** $\mathbb{V}(X) = \mathbb{E}[(X - \mathbb{E}[X])^2]$.
- **Equivalent form:** $\mathbb{V}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$.
- **Shift invariance:** $\mathbb{V}(X + c) = \mathbb{V}(X)$.
- **Scaling:** $\mathbb{V}(aX) = a^2 \mathbb{V}(X)$.
- **Additivity (independent X, Y):** $\mathbb{V}(X + Y) = \mathbb{V}(X) + \mathbb{V}(Y)$.
- **General case:** $\mathbb{V}(X + Y) = \mathbb{V}(X) + \mathbb{V}(Y) + 2\mathbb{V}(X, Y)$.
- **Nonnegativity:** $\mathbb{V}(X) \geq 0$, with equality iff X constant.

As a consequence, let us consider an arbitrary sequence of **independent** RVs X_1, X_2, \dots, X_n and $a, b \in \mathbb{R}$, to highlight two additional key properties of the variance:

$$\mathbb{V}[aX_1 + b] = a^2 \mathbb{V}[X_1] \quad (3.45)$$

$$\mathbb{V}\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n \mathbb{V}[X_i] \quad (3.46)$$

3.3 Some particular RVs

3.3.1 Bernoulli and binomial

Consider tossing a coin, where $\mathbb{P}(\text{heads}) = \theta$ with $0 \leq \theta \leq 1$. Note that, as a consequence, $\mathbb{P}(\text{tails}) = 1 - \theta$. This is called the Bernoulli distribution, and it is written as

$$Y \sim \text{Ber}(\theta). \quad (3.47)$$

Bernoulli's pmf is given by

$$p_{\text{Ber}}(y) = \begin{cases} \theta, & y = 1, \\ 1 - \theta, & y = 0, \\ 0, & \text{otherwise.} \end{cases}$$

More concisely, we can write the pmf as

$$p_{\text{Ber}}(y) = \theta^y (1 - \theta)^{1-y}. \quad (3.48)$$

Bernoulli likelihood

The expression $p_{\text{Ber}}(y) = \theta^y (1 - \theta)^{1-y}$, will be fundamental in ML, since this close form expression will allow for optimising the parameters of the Bernoulli distribution in the light of data.

Let us now consider N Bernoulli trials, denoted $Y_i \sim \text{Ber}(\cdot|\theta), i = 1, \dots, N$ – this can be thought of as tossing a coin N times. Define the RV S as the total number of heads, that is

$$S = \sum_{i=1}^N I(y_i = 1). \quad (3.49)$$

The distribution of S is

$$\text{Bin}(s|N, \theta) = \binom{N}{s} \theta^s (1 - \theta)^{N-s}, \quad (3.50)$$

which is known as the Binomial distribution, where $\binom{N}{s} = \frac{N!}{(N-s)!s!}$.

Logistic regression

Within ML, the Bernoulli distribution is largely used for (probabilistic) classification. In such case we want to predict a binary variable $y \in \{0, 1\}$, representing a sample being class 1 (and not class 0), conditional to some observed features $x \in \mathbb{R}^d$. We can then model $p(y = 1)$ as

$$p(y|x, \theta) = \text{Ber}(y|f(x, \theta)). \quad (3.51)$$

Since $f(x, \theta)$ represents the parameter of $\text{Ber}(y)$, we need $0 \leq f(x, \theta) \leq 1$. However, to avoid that requirement, we can leave $f(x, \theta)$ unconstrained and write

$$p(y|\theta, x) = \text{Ber}(y|\sigma(f(x, \theta))), \quad (3.52)$$

where

$$\sigma(a) = \frac{1}{1 + e^{-a}} \quad (3.53)$$

is the logistic function. When we choose $f(x, \theta) = \theta_1^\top x + \theta_2$, we obtain the logistic regression model.

We show the sigmoid function in Figure 16. A visualisation of an adjusted logistic regression model was shown in Figure 1.

The extension of the Bernoulli (resp. Binomial) distribution to the multivariate case, i.e., when the output of the experiments take values on categories $\{1, 2, \dots, K\}$ with $K > 2$ is known as the categorical (resp. Multinomial). These distributions will be left for personal study.

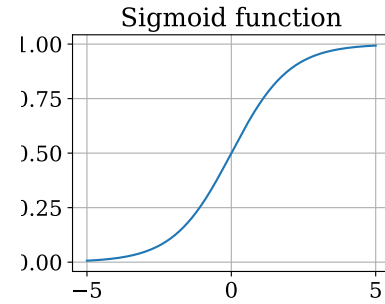


Fig. 16. Visualisation of a sigmoid function.

3.3.2 Uniform distribution

The uniform distribution models situations in which all outcomes within a set are **equally likely**.

Discrete case. Consider a random variable Y that can take K equally probable values $\{1, 2, \dots, K\}$. Then

$$Y \sim \text{Uniform}(\{1, \dots, K\}), \quad (3.54)$$

and its pmf is

$$p_{\text{Uniform}}(y) = \begin{cases} \frac{1}{K}, & y \in \{1, \dots, K\}, \\ 0, & \text{otherwise.} \end{cases} \quad (3.55)$$

The mean and variance are

$$\mathbb{E}[Y] = \frac{K+1}{2}, \quad \mathbb{V}[Y] = \frac{K^2-1}{12}.$$

Continuous case. Let X be a continuous random variable uniformly distributed over the interval $[a, b]$:

$$X \sim \text{Uniform}(a, b), \quad (3.56)$$

with pdf

$$p_{\text{Uniform}}(x) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b, \\ 0, & \text{otherwise.} \end{cases} \quad (3.57)$$

The expectation and variance are given by

$$\mathbb{E}[X] = \frac{a+b}{2}, \quad \mathbb{V}[X] = \frac{(b-a)^2}{12}.$$

Uniform likelihood

The uniform distribution represents **maximum uncertainty** within a range — every value is equally probable. In machine learning, uniform priors are often used to express a lack of prior knowledge about a parameter, or as initial distributions in random sampling (e.g., parameter initialisation or Monte Carlo methods).

3.3.3 Gaussian distribution

The pdf the (scalar) Gaussian RV X is given by

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}, \quad (3.58)$$

when $\mu = 0$ and $\sigma = 1$, we say that X is a standard normal.

Why Gaussian

The Gaussian distribution is widely used for modelling continuous RVs, this is, in part, because:

- Central limit theorem: The scaled sum of multiple independent RVs, with different distributions, converges to a Gaussian RVs
- The Gaussian distribution is the distribution with the maximum entropy (fewer structural assumptions) for a fixed variance
- Its parameters are easy to interpret: the pdf is in fact parametrised by its mean and variance.
- The maths are simple: this is key in applications

The Gaussian distribution is central to model **observation noise**, and therefore is key in regression models. That is, one can consider a linear regression models where

$$Y|x \sim \mathcal{N}(y|a^\top x + b, \sigma^2), \quad (3.59)$$

which is equivalent to

$$y = a^\top x + b + \epsilon, \quad (3.60)$$

where $\epsilon \sim \mathcal{N}(0, \sigma^2)$.

Multivariate Gaussian. Let $\mathbf{X} \in \mathbb{R}^d$ be a d -dimensional random vector. The multivariate Gaussian distribution is defined as

$$\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma), \quad (3.61)$$

where $\boldsymbol{\mu} \in \mathbb{R}^d$ is the mean vector and $\Sigma \in \mathbb{R}^{d \times d}$ is the covariance matrix (symmetric positive definite). Its pdf is

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}. \quad (3.62)$$

Some observations about the MVN listed as follows:

- The mean vector $\boldsymbol{\mu}$ gives the expected value of each component.
- The covariance matrix Σ encodes the variance of each component and the pairwise correlations.
- Linear transformations of Gaussian vectors are still Gaussian: if $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ and A is a matrix, $A\mathbf{X} + b \sim \mathcal{N}(A\boldsymbol{\mu} + b, A\Sigma A^\top)$. This enables sampling of multivariate Gaussians with arbitrary mean and covariance.
- Central in probabilistic ML: used in multivariate regression, Gaussian processes, PCA, and Bayesian inference.

Recall that Figure 1 shows a mixture of two two-dimensional Gaussian distributions with different covariance matrices.

3.3.4 Other distributions

There are several other distributions that are widely used in statistics and machine learning, but we will not cover them in detail due to time constraints. These include, among others:

- **Heavy-tailed:** Student's t , Cauchy
- **Positive-only:** Chi-square, Gamma, inv-Gamma, Exponential, half Gaussian
- **Bounded / shape-flexible:** Beta, Laplace

These distributions are frequently used as priors, for modelling non-Gaussian noise, or in robust statistical methods.

3.4 Transformations of RVs distributions

In ML, we are interested in the construction of expressive **generative models**, i.e., probability distributions. Despite having a large collection of known distributions, sometimes we need to design purpose-specific models; this is achieved by transforming a given RV.

Consider $X \sim p_X(\cdot)$ and $Y = f(X)$ a deterministic transformation. We are interested in computing $p_Y(y)$.

The discrete case is straightforward, as we just need to sum over the possible (discrete) values. That is,

$$p_Y(y) = \sum_{x: y=f(x)} p_X(x). \quad (3.63)$$

Uniform into binary

Consider

$$p_X = \text{Uniform}(\{0, 1, \dots, 9\}), \quad (3.64)$$

and

$$y = f(x) = \begin{cases} 1, & \text{if } x \text{ is odd,} \\ 0, & \text{if } x \text{ is even.} \end{cases}$$

The continuous case is a bit more involved, since we cannot directly sum over all possible values of the RV. However, we can work with the cdf to show that

$$P_Y(y) = \mathbb{P}(Y \leq y) = \mathbb{P}(f(x) \leq y) = \mathbb{P}(x \in \{x | f(x) \leq y\}). \quad (3.65)$$

Let us notice that if f is invertible, we can rely on the **change of variable theorem**, which states that

$$p_Y(y) = p_X(x) \left| \frac{dx}{dy} \right|. \quad (3.66)$$

This identity can be proven using the cdf. When f is invertible, we can denote $g = f^{-1}$ and note that (assuming that f is non-decreasing)

$$P_Y(y) = \mathbb{P}(Y \leq y) = \mathbb{P}(f(x) \leq y) = \mathbb{P}(x \leq g(y)) = P_X(g(y)). \quad (3.67)$$

We can now take the derivative, to give:

$$p_Y(y) = \frac{d}{dy} P_X(g(y)) = \frac{dP_X(g(y))}{dx} \frac{dx}{dy} = p_X(g(y)) \frac{dx}{dy}. \quad (3.68)$$

If f had been non-increasing, we would have obtained the same expression with the reversed sign. Therefore, we conclude eq. (3.66).

In the multivariate case, the CVT reads

$$p_Y(y) = p_X(x) |\det J_q(y)|, \quad (3.69)$$

where $J_q(y) = \frac{dg(y)}{dy}$ is the Jacobian of $g = f^{-1}$.

Example: χ^2 family of distributions.

A random variable Z is said to distribute according to a χ^2 (denoted $Z \sim \chi^2$) when $Z = \sum_{i=1}^k X_i^2$ for $X_i \sim \mathcal{N}(0, 1)$. We provide a visualisation of the probability density function of χ^2 distributions for various k in Figure 17.

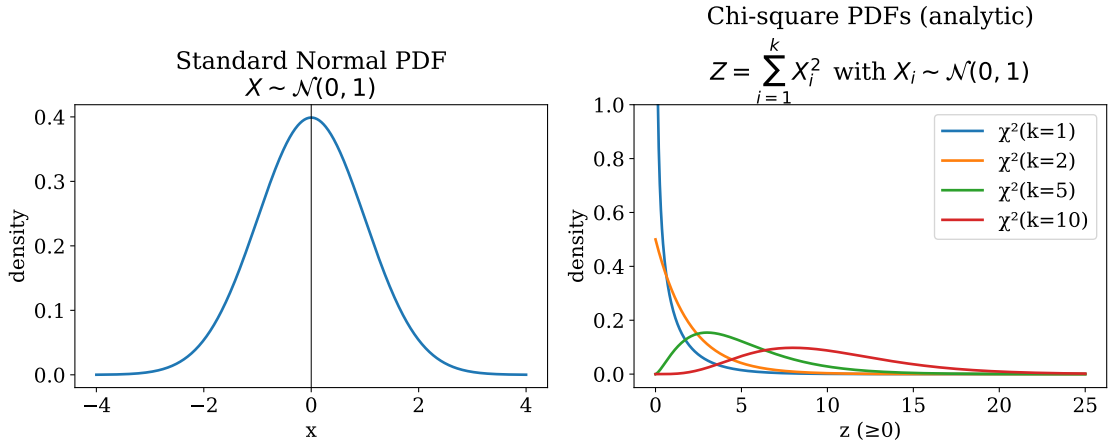


Fig. 17. Gaussian distribution and χ^2 distribution for different k parameters.

3.5 Sampling

Consider a set $\{x_1, x_2, \dots, x_N\}$ of independent and identically distributed (i.i.d.) samples of an RV $X \sim p_X$. We are interested in approximating the law p_X , and also computing expectations wrt it. An empirical estimate of the law of X is the so called **histogram**

$$P_N(x) = \frac{1}{N} \sum_{i=1}^N \delta_{x_i}(x), \quad (3.70)$$

where $\delta_{x_i}(\cdot)$ denotes the delta-Dirac mass in x_i .

Real approximation?

Does $P_N(x)$ become more and more similar to the true law of X as N grows?

Notice that P_N allows to compute (approximate) expectations. Let us consider a function $f : x \mapsto f(x)$, and

$$I_f \stackrel{\text{def}}{=} \int f(x)p(x)dx \simeq f(x)P_N(x) = \frac{1}{N} \sum_{i=1}^N f(x_i) \stackrel{\text{def}}{=} I_N. \quad (3.71)$$

Due to the Strong Law of Large Numbers (SLLN), we know that I_N gets **arbitrarily close** to I_f as $N \rightarrow \infty$. This approach to computing expectation is referred to as **Monte Carlo approximation**.

A relevant question is how to generate the samples needed to approximate P_N , this is known as **sampling**. When sampling directly from a distribution is not possible, there are three main approaches to sampling as follows:

- **Rejection sampling:** Draw candidate samples $x' \sim q(x)$ from an easy-to-sample *proposal distribution* $q(x)$, and accept x' with probability

$$\alpha = \frac{p(x')}{Mq(x')}, \quad \text{for } M \geq \sup_x \frac{p(x)}{q(x)}.$$

Accepted samples are distributed according to the target $p(x)$. Simple but inefficient if $q(x)$ poorly matches $p(x)$.

- **Importance sampling:** Draw samples $x_i \sim q(x)$ and compute weighted estimates of expectations:

$$\mathbb{E}_p[f(X)] = \int f(x)p(x)dx \simeq \frac{1}{N} \sum_{i=1}^N f(x_i)w_i, \quad w_i = \frac{p(x_i)}{q(x_i)}.$$

Useful when direct sampling from p is difficult. Forms the basis for many Bayesian inference methods.

- **Markov Chain Monte Carlo (MCMC):** Construct a Markov chain (X_1, X_2, \dots) such that its stationary distribution is $p(x)$. For example, the Metropolis-Hastings algorithm updates $X_t \rightarrow X_{t+1}$ via:

$$X_{t+1} = \begin{cases} x', & \text{with probability } \min\left(1, \frac{p(x')q(x_t|x')}{p(x_t)q(x'|x_t)}\right), \\ X_t, & \text{otherwise.} \end{cases}$$

Powerful for high-dimensional or complex distributions, but requires monitoring convergence.

These methods are central to probabilistic machine learning, Bayesian inference, and situations where exact integration is intractable.

Hierarchical sampling. In some settings involving more than one random variable, sampling is naturally performed by **sampling in stages**. In this case, a random variable X is generated by first sampling an auxiliary variable Z , and then sampling X conditional on Z . This is called **hierarchical sampling** (or sampling from a hierarchical model).

Formally, if $Z \sim p(z)$ and $X | Z = z \sim p(x | z)$ is possible to sample, then sampling from the **marginal** distribution $p(x)$ is obtained by:

$$z \sim p(z), \quad x \sim p(x | z).$$

Example: Gaussian mixture.

The Gaussian mixture model (GMM) is a random variable with a pdf given by a mixture of Gaussian components. For mixture weights $\pi_k \geq 0$, with $\sum_{k=1}^K \pi_k = 1$, and component densities $X | Z = k \sim \mathcal{N}(\mu_k, \Sigma_k)$, the marginal density of X is

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x | \mu_k, \Sigma_k).$$

When $K = 2$, we can define the component index as the RV $Z \in \{1, 2\}$. Therefore, a component can be chosen with

$$\mathbb{P}(Z = 1) = \pi, \quad \mathbb{P}(Z = 2) = 1 - \pi,$$

and then a sample of X can be obtained via:

$$X | (Z = 1) \sim \mathcal{N}(\mu_1, \sigma_1^2), \quad X | (Z = 2) \sim \mathcal{N}(\mu_2, \sigma_2^2).$$

The extension to arbitrary $K > 2$ is left as an exercise. Notice that in this case, the component index Z is given by a Categorical distribution with parameters $\pi_k \geq 0$, with $\sum_{k=1}^K \pi_k = 1$.

4 Statistics

4.1 The statistical model

This part of the course focuses on *mathematical statistics*, a methodological approach to inference based on probability, algebra, geometry, optimisation and measure theory. In this setup, we assume the existence of a dataset generated from an unknown *statistical model*, aka probabilistic or generative model. Our objective is to use the data to estimate such model (and in particular its parameters) to ultimately learn the underlying properties of the data generating process and make predictions.

Definition 4.1 (Statistical model). A statistical model is a set of probability distribution that can be considered as *candidates* for the data-generating mechanism.

In some cases, the statistical model can be *parametric*, that is, represented by a **finite** set of parameters. This includes the Normal distribution, which is parametrised by its mean and variance. In such cases, inference boils down to identifying the parameters. Furthermore, we will consider a dataset with elements x belonging to an abstract space \mathfrak{X} , where typically $\mathfrak{X} = \mathbb{R}^n$, and where x is the realisation of an RV $X \in \mathfrak{X}$. The statistical model can also be interpreted as the space of the possible hypotheses that explain the observed data.

The main aim of this part is statistical inference, that is, given an observation x , what is the distribution of X , or at least some of its features? To answer this question we will restrict ourselves to parametric statistical models; this requires the rigorous definition of the parameters and their space.

Definition 4.2 (Parameters and parameter space). A parameter is a fixed but unknown quantity that specifies a feature of a random variable's distribution (e.g., its mean, variance, or proportion). Parameters will be denoted with the symbol $\theta \in \Omega$, where Ω is the set of all possible values for the parameter called *parameter space*.

We denote the parametric family \mathcal{P} as

$$\mathcal{P} = \{\mathcal{P}_\theta | \theta \in \Omega\},$$

where \mathcal{P}_θ is a probability distribution *indexed* by a parameter $\theta \in \Omega$. We will consider finite-dimensional Ω , that is, $\Omega \subseteq \mathbb{R}^n$. Therefore, we denote:

$$\theta = [\theta_1, \dots, \theta_n]^\top.$$

In summary:

- θ is the parameter to be estimated from data
- Ω is the parameter space, where $\Omega \subseteq \mathbb{R}^n$
- \mathcal{P}_θ is probability over \mathfrak{X} (as a function of θ)
- X is a RV in \mathfrak{X}
- x is the data, a realisation of X and a generic element of \mathfrak{X} .

Example 4.1 (Computer manufacturer). A computer manufacturer wishes to estimate the lifetime of a particular component in its computers. To do this, data is first collected from computers that have been used under normal conditions. After consulting with experts, they decide to use a normal distribution to model the time it will take for the component of interest to fail. The useful life of the components is then modelled with an average lifetime θ and variance σ^2 , with θ and σ^2 unknown parameters. If there are N components, the random variables that model the useful life of each component will be identified as X_1, \dots, X_N , with $X_i \sim \mathcal{N}(\theta, \sigma^2), \forall i = 1, \dots, N$. What do you think of this model?

Statistical inference is a tool that allows us to solve a number of data-driven challenges. The most relevant ones involve identification, that is, to discover the model that generated the data, and prediction, where we estimate a quantity that has not yet been observed. Needless to say, a statistically-oriented approach to these challenges requires appropriately modelling the associated uncertainty.

4.2 Statistics

The initial setup in statistical inference features, in addition to our own assumptions, an observation dataset. Therefore, the first course of action is to apply transformations to the data; this underpins the construction of the so called *statistics*.

Definition 4.3 (Statistic). Let $(\mathcal{T}, \mathcal{A}, \mu)$ be a probability space and $X \in \mathcal{X}$ a RV with parametric distribution $\mathcal{P} = \{P_\theta \mid \theta \in \Theta\}$. A statistic is a function of the realisation $X = x$ that does not depend from the parameter θ (and the distribution P_θ).

$$T : \mathcal{X} \rightarrow \mathcal{T}$$

$$x \mapsto T(x).$$

Remark 4.1. It is critical to understand the difference between the value of the statistic $T(x)$ and the application of the function $T(\cdot)$ to the RV X . The former is a “fixed” value, while the latter is a random variable with its own probability distribution induced by P_θ and $T(\cdot)$.

Some example statistics are:

$$T(x) = \frac{1}{n} \sum_{i=1}^n x_i, \quad T'(x) = x, \quad T''(x) = \min(x), \quad T'''(x) = c \in \mathbb{R}.$$

Sufficiency The objective of a statistic is to encapsulate or summarize the information contained in a sample $x = (x_1, x_2, \dots, x_n)$ that is useful for determining parameters of the distribution of X . Therefore, the identity function —or even the mean in some cases— seems to fulfil this mission, at least intuitively. This is because, in practice, we want to extract as much information as possible from the data; this is achieved by the statistic T (which summarizes all the data) and the statistic T' (which contains all the data). On the contrary, the statistic T'' loses information, since it only retains the minimum value of all the data, thus losing the representation of, e.g., the dispersion of the sample. The same analysis can be made for the constant statistic T''' , which contains no information from the data whatsoever.

Informally, the idea of *sufficiency* of a statistic wrt a parameter can be expressed as

“...no other statistic that can be calculated from the same sample provides any additional information as to the value of the parameter.” (Ronald Fisher, On the mathematical foundations of theoretical statistics (1922))

A formal definition is

Definition 4.4 (Sufficient statistic). Let (S, \mathcal{A}, μ) be a probability space and $X \in \mathcal{X}$ a RV with parametric distribution $\mathcal{P} = \{P_\theta \mid \theta \in \Theta\}$. We say that the function $T : \mathcal{X} \rightarrow \mathcal{T}$ is a sufficient statistic for θ (or for X , or for \mathcal{P}) if the conditional distribution $X|T(X)$ does not depend on the parameter θ , that is,

$$P_\theta(X \in A | T(X)), \quad A \in \mathcal{B}(X), \text{ does not depend on } \theta.$$

Example 4.2 (Trivial sufficient statistic). For a given parametric family \mathcal{P} , the statistic defined by

$$T(x) = x,$$

is a sufficient statistic. Indeed, $P_\theta(X \in A | X = x) = \mathbf{1}_A(x)$ does not depend on the parameter θ .

Example 4.3 (Sufficient statistic for a Bernoulli RV). Let $x = (x_1, \dots, x_n) \sim \text{Ber}(\theta)$, $\theta \in \Theta = [0, 1]$, that is

$$P_\theta(X = x) = \theta^{\sum x_i} (1 - \theta)^{n - \sum x_i}.$$

Observe that $T(x) = \sum_{i=1}^n x_i$ is a sufficient statistic for θ . In fact (by definition):

$$\begin{aligned} P(X = x | T(X) = t) &= \frac{P(T(X) = t | X = x)P(X = x)}{P(T(X) = t)} && \text{(Bayes thm)} \\ &= \frac{\mathbf{1}_{T(x)=t} \theta^{\sum x_i} (1 - \theta)^{n - \sum x_i}}{\binom{n}{t} \theta^t (1 - \theta)^{n-t}} && \text{(model, and sum of Bernoulli is Binomial)} \\ &= \mathbf{1}_{T(x)=t} \binom{n}{t}^{-1} && \text{(since } T(x) = t) \end{aligned}$$

Therefore, $T(x) = \sum_{i=1}^n x_i$ is a sufficient statistic.

In practice, sufficiency is no assessed by definition, but via a result known as Fisher-Neyman theorem. This will be left for personal study.

4.3 Estimators

We now focus on a particular class of statistics that enables parameter estimation.

Definition 4.5 (Estimator). Let $g : \Omega \rightarrow \mathbb{R}^n$, a function of the unknown parameter. An estimator of g is a statistic $\hat{g} : \mathfrak{X} \rightarrow g(\Omega)$. We will say that $\hat{g}(x)$ is the estimation of $g(\theta)$.

Remark 4.2. Estimators are particular instances of statistics: they are functions of the data with an image space that is equal to the image space of Ω through $g(\cdot)$.

Remark 4.3. Estimators can be used to identify parameters, e.g., when $g(\theta) = \theta$, or other related relevant quantities. For instance, for the Gaussian model, where the parameter is given by $\theta = [\mu, \sigma^2]$, we might be interested in estimating the 95% confidence interval, given by

$$g(\theta) = [\mu - 2\sigma, \mu + 2\sigma]. \quad (4.1)$$

Example 4.4 (Estimator for the mean of a Gaussian RV). Let us consider $X = (X_1, \dots, X_n) \sim \mathcal{N}(\mu, \sigma^2)$. An estimator of $g(\theta) = g(\mu, \sigma) = \mu$ is the statistic

$$\hat{g}(X) = \frac{1}{n} \sum_{i=1}^n X_i.$$

4.3.1 Unbiased estimators

Estimators, being a deterministic function of the RV X , are themselves RVs. Therefore, we can assess an estimator in terms of its mean: we would ideally want that an estimator for $g(\theta)$ has a mean equal to that quantity.

Definition 4.6 (Unbiased estimator). Let $\hat{g}(X)$ be an estimator of $g(\theta)$. This is an unbiased estimator if

$$\mathbb{E} \hat{g}(X) = g(\theta),$$

where the *bias* of \hat{g} is defined as

$$b_{\hat{g}}(\theta) = \mathbb{E} \hat{g}(X) - g(\theta).$$

Furthermore, we say that the estimator is **asymptotically unbiased** if:

$$\lim_n \mathbb{E} \hat{g}(X_1, \dots, X_n) = g(\theta),$$

meaning that, the estimator only becomes unbiased when an *infinite amount of data* are available.

Remark 4.4. The relevance of unbiased estimators in statistical inference stems from the fact that they recover the true parameter on average. However, one should not focus on unbiasedness exclusively: performing well on average does not guarantee anything about the estimator's variability (variance), or about how large a sample is needed for the estimator to be reliable.

The following examples illustrate the role of unbiased estimators in two parametric families.

Example 4.5 (Unbiased estimator for the Gaussian mean). The estimator of $g(\theta) = \mu$ described in Example 4.4 is unbiased. Indeed:

$$\mathbb{E} \hat{g}(X) = \mathbb{E} \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} \sum_{i=1}^n \mathbb{E} X_i = \frac{1}{n} \sum_{i=1}^n \mu = \mu.$$

The following example shows a **biased** estimator of the variance and how an **unbiased** estimator can be constructed based on it.

Example 4.6 (Pythagoras). Let us consider a parametric family $\mathcal{P} = \{P_\theta \mid \theta \in \Theta\}$, and denote its mean and variance by μ and σ^2 respectively. Using observations x_1, x_2, \dots, x_n , we can calculate the variance of the estimator of the mean, given by $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ according to

$$\mathbb{V}_\theta(\bar{x}) = \mathbb{V}_\theta\left(\frac{1}{n} \sum_{i=1}^n x_i\right) \underset{\text{i.i.d.}}{=} \frac{1}{n^2} \sum_{i=1}^n \mathbb{V}_\theta(x_i) = \frac{\sigma^2}{n}. \quad (4.2)$$

This means that the estimator of the mean using n observations has a variance σ^2/n .

Let us now consider the following estimator for the variance:

$$S_2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2, \quad (4.3)$$

and observe that the expected value of such estimator is

$$\begin{aligned} \mathbb{E}_\theta(S_2) &= \mathbb{E}_\theta\left(\frac{1}{n} \sum_{i=1}^n (x_i - \mu + \mu - \bar{x})^2\right) \\ &= \mathbb{E}_\theta\left(\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 + 2\frac{1}{n} \sum_{i=1}^n (x_i - \mu)(\mu - \bar{x}) + \frac{1}{n} \sum_{i=1}^n (\mu - \bar{x})^2\right) \\ &= \mathbb{E}_\theta\left(\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 - 2(\mu - \bar{x})^2 + (\mu - \bar{x})^2\right) \\ &= \mathbb{E}_\theta\left(\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 - (\mu - \bar{x})^2\right) \\ &= \mathbb{V}_\theta(x_i) - \mathbb{V}_\theta(\bar{x}) \quad \text{see eq. (4.2)} \\ &= \sigma^2 - \sigma^2/n = \left(\frac{n-1}{n}\right) \sigma^2. \end{aligned} \quad (4.4)$$

Therefore, the bias of the estimator in eq. (4.3) is asymptotically unbiased. Additionally, we can remove this bias by multiplying the estimator S_2 in eq. (4.3) by $n/(n-1)$, giving

$$S'_2 = \frac{n}{n-1} S_2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2, \quad (4.5)$$

for which it holds

$$\mathbb{E}_\theta(S'_2) = \left(\frac{n}{n-1}\right) \mathbb{E}_\theta(S_2) \underset{\text{eq. (4.4)}}{=} \left(\frac{n}{n-1}\right) \left(\frac{n-1}{n}\right) \sigma^2 = \sigma^2. \quad (4.6)$$

As a consequence, the estimator S'_2 in eq. (4.5) is unbiased.

4.3.2 Loss functions

We now consider loss functions to determine the cost of estimating a parameter by a given estimator. From now on, we will consider estimators of $g(\theta) = \theta$ for simplicity of notation.

Definition 4.7 (Loss function). Let $\theta \in \Omega$ be a parameter and $a \in \Omega$ an estimator. The cost of estimating θ by a is given by:

$$L : (\Omega \times \Omega) \rightarrow \mathbb{R} \quad (4.7)$$

$$(\theta \times a) \mapsto L(\theta, a). \quad (4.8)$$

Example 4.7 (Quadratic loss function). An extensively-used loss for comparing estimators is the quadratic loss, given by

$$L_2(\theta, a) = \|\theta - a\|^2.$$

Remark 4.5. Why do we use the exponent 2 used and not other such as 4 or 1?

Example 4.8 (0 – 1 loss). When estimating parameters in a space with no order relation, the 0 – 1 loss is usually considered. This is given by

$$L_{01}(\theta, a) = \mathbf{1}_{\theta \neq a}.$$

The loss is computed as a function of the estimator, which is in turn computed as a function of X , this means that the loss is itself random. We are thus interested in the expected loss, which is known as *risk*.

Definition 4.8 (Risk). Given a parameter $\theta \in \Omega$, an estimator $a \in \Omega$, and a loss function $L(\theta, a)$, the *risk* of estimating θ by a is defined as the expected loss:

$$R : (\Omega \times \Omega) \rightarrow \mathbb{R}, \quad (4.9)$$

$$(\theta \times a) \mapsto \mathbb{E}[L(\theta, a)]. \quad (4.10)$$

In particular, the risk associated to the quadratic loss in ex. 4.7 for an estimator ϕ of the parameter θ , is given by:

$$\begin{aligned} R(\theta, \phi) &= \mathbb{E}(\theta - \phi)^2 \\ &= \mathbb{E}(\theta - \bar{\phi} + \bar{\phi} - \phi)^2; \quad \text{denoting } \bar{\phi} = \mathbb{E}\phi \\ &= \mathbb{E}(\theta - \bar{\phi})^2 + 2(\theta - \bar{\phi})(\bar{\phi} - \phi) + (\bar{\phi} - \phi)^2 \\ &= \underbrace{(\theta - \bar{\phi})^2}_{=b_\phi^2 \text{ (bias}^2)} + \underbrace{\mathbb{E}(\bar{\phi} - \phi)^2}_{=V_\phi \text{ (variance)}}. \end{aligned} \quad (4.11)$$

From this result, we can see a justification for the use of the quadratic loss: its risk splits automatically in two terms expressing the accuracy (how unbiased it is) and precision (how disperse) of the estimator.

4.4 Maximum likelihood estimator (MLE)

In general, loss functions are not a pathway to construct estimators. This is because optimising a defined loss is unfeasible, as it depends explicitly on the true parameter. Luckily, we can construct an estimator based directly on: i) the pdf of $X \in \mathcal{X}$, where the parameters θ appears explicitly, and ii) an observation dataset. To this end, the following definition is crucial.

Definition 4.9 (Likelihood function). Let X be a random variable with probability density (or mass) function $p_\theta(x)$ indexed by a parameter $\theta \in \Omega$. For an observed value x , the *likelihood function* is defined as the mapping

$$\mathcal{L} : \Omega \rightarrow \mathbb{R}_{\geq 0}, \quad (4.12)$$

$$\theta \mapsto p_\theta(x), \quad (4.13)$$

which measures how plausible each parameter value θ is given the observation x .

Remark 4.6. The likelihood function is not a probability density function.

Remark 4.7 (Log-likelihood). It is often convenient to work with the *log-likelihood*,

$$\ell(\theta) = \log \mathcal{L}(\theta),$$

since the logarithm is a strictly increasing function and therefore preserves maximisers of the likelihood. Moreover, $\ell(\theta)$ is typically easier to manipulate and optimise, both analytically and computationally.

Remark 4.8. In general (yet not always) we will assume i.i.d. observations $\mathcal{D} = X_1, \dots, X_n$, $X_i \sim p(x|\theta)$. In such case, the likelihood factorises as $L_{\mathcal{D}}(\theta) = \prod_{i=1}^n L_{X_i}(\theta)$, and the log-likelihood takes the form

$$l_{\mathcal{D}}(\theta) = \sum_{i=1}^n l_{X_i}(\theta). \quad (4.14)$$

Example 4.9 (Likelihood of the mean of a Gaussian). Let us consider the observations $\mathcal{D} = \{x_1, \dots, x_n\}$, where x_i is the realisation of a RV $X_i \sim \mathcal{N}(\mu, \sigma^2)$ with known σ^2 . The likelihood function of μ is given by:

$$\begin{aligned} \mathcal{L}(\mu) &= p(\mathcal{D}|\mu) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-1}{2\sigma^2}(x_i - \mu)^2\right) \\ &= \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(\frac{-1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right). \end{aligned} \quad (4.15)$$

Thus, the log-likelihood is given by:

$$l(\mu) = \log \mathcal{L}(\mu) = -\frac{n}{2} \log(2\pi\sigma^2) + \frac{-1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2. \quad (4.16)$$

Definition 4.10 (Maximum likelihood estimator, MLE). Let x be an observation (or realisation) of the random variable X , and let $\mathcal{L}(\theta | x)$ be the corresponding likelihood function. The *maximum likelihood estimator* is defined by

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} \mathcal{L}(\theta | x). \quad (4.17)$$

Exercise 4.1. Find the MLE for $\theta = (\mu, \Sigma)$ for the RV $X \sim \mathcal{N}(\mu, \Sigma)$ in ex. 4.9.

Remark 4.9. The MLE can be defined either in terms of the likelihood itself or any non-decreasing function of it, and it may also fail to exist or to be unique. Furthermore, we will often even ignore constants in the (log) likelihood, as these do not affect the maximiser of $L(\theta)$.

Example 4.10 (MLE: Bernoulli). Let $X_1, \dots, X_n \sim \text{Ber}\theta$, the likelihood of θ is given by

$$L(\theta) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i}, \quad (4.18)$$

and its log-likelihood by $l(\theta) = (\sum_{i=1}^n x_i) \log \theta + (n - \sum_{i=1}^n x_i) \log(1 - \theta)$. The MLE can be found by $\frac{\partial l(\theta)}{\partial \theta} = 0$:

$$\begin{aligned} \frac{\partial l(\theta)}{\partial \theta} = 0 &\Rightarrow \left(\sum_{i=1}^n x_i\right) \theta^{-1} = (n - \sum_{i=1}^n x_i) (1 - \theta)^{-1} \\ &\Rightarrow \sum_{i=1}^n x_i (1 - \theta) = (n - \sum_{i=1}^n x_i) \theta \\ &\Rightarrow \theta = \sum_{i=1}^n x_i / n. \end{aligned}$$

Example 4.11. Let $X \sim \text{Uniform}\theta$, that is, $p(x) = \theta^{-1}\mathbf{1}_{0 \leq x \leq \theta}$. To calculate the likelihood function, observe that

$$L(\theta) = \prod_{i=1}^n p_{\theta}(x_i), \quad (4.19)$$

and note that $p_{\theta}(x_i) = 0$ if $x_i > \theta$. Therefore, $L(\theta) > 0$ only if θ is greater than each observation, in particular, if $\theta \geq \max_{i=1, \dots, n} \{x_i\}$. Additionally, if $\theta \geq \max_{i=1, \dots, n} \{x_i\}$, then notice that $p_{\theta}(x_i) = 1/\theta$. As a consequence, the likelihood is given by

$$L(\theta) = \theta^{-n}, \quad \theta \geq \max_{i=1, \dots, n} \{x_i\} \quad (4.20)$$

and the MLE is $\theta_{\text{MLE}} = \max_{i=1, \dots, n} \{x_i\}$.

Plot $l(\theta)$
and MLE

The maximum likelihood estimator (MLE) satisfies several important theoretical properties:

- **Consistency:** Under the assumption that the statistical model is *identifiable*—i.e., different parameter values correspond to different probability distributions—the MLE converges to the true parameter as the number of observations grows. Intuitively, maximising the likelihood asymptotically minimises the Kullback–Leibler divergence between the true distribution and the distribution induced by a candidate parameter.
- **Equivariance:** If $\hat{\theta}_{\text{MLE}}$ is the MLE of θ , then for any transformation g , the MLE of $g(\theta)$ is $g(\hat{\theta}_{\text{MLE}})$. This property allows us to compute MLEs under reparametrisations directly.
- **Asymptotic normality:** For large sample sizes, the MLE is approximately normally distributed around the true parameter with covariance matrix given by the inverse Fisher information. Formally,

$$\sqrt{n}(\hat{\theta}_{\text{MLE}} - \theta) \xrightarrow{d} \mathcal{N}(0, I(\theta)^{-1}),$$

where $I(\theta)$ is the Fisher information matrix.

- **Asymptotic efficiency:** As a consequence of asymptotic normality, the MLE achieves the Cramér–Rao lower bound for the variance in the limit of large n , making it asymptotically optimal among unbiased estimators.

In practice, these properties justify the widespread use of the MLE: it not only converges to the true parameter under mild assumptions, but also allows for straightforward reparametrisations and provides an estimator with minimal asymptotic variance.

4.5 MLE in practice: three examples

4.5.1 Gaussian linear regression

Let us consider the task of modelling the number of passengers travelling each month in an airline. Experts have established that this variable as a quadratic growing trend and an oscillatory component of annual frequency. These phenomena can be explained by population growth, diminishing costs of operation, and the yearly seasonality of economic activity.

Assuming that the number of passenger is a RV, we can model its (conditional) distribution wrt time t by a normal distribution parametrised by

$$X \sim \mathcal{N}(\theta_0 + \theta_1 t^2 + \theta_2 \cos(2\pi t/12), \theta_3^2), \quad (4.21)$$

where $\theta_0, \theta_1, \theta_2$ control the mean and θ_3 the variance.

Consequently, if the observations are given by $\{(t_i, x_i)\}_{i=1}^n$, the log-likelihood of θ is

$$\begin{aligned} l(\theta) &= \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\theta_3^2}} \exp \left(-\frac{(x_i - \theta_0 - \theta_1 t_i^2 - \theta_2 \cos(2\pi t_i/12))^2}{2\theta_3^2} \right) \right) \\ &= \frac{n}{2} \log(2\pi\theta_3^2) - \frac{1}{2\theta_3^2} \sum_{i=1}^n (x_i - \theta_0 - \theta_1 t_i^2 - \theta_2 \cos(2\pi t_i/12))^2 \end{aligned} \quad (4.22)$$

where we can see that θ_{MLE} can be calculated explicitly and it is a function of $\{(t_i, x_i)\}_{i=1}^n$, since eq. (4.22) is quadratic in $[\theta_0, \theta_1, \theta_2]$.

plot and
show the
ML soln

4.5.2 Classification

$\hat{\theta}_{\text{MLE}}$ was explicitly calculated above because the Gaussian model with a linearly parametrised mean leads to a quadratic log-likelihood, where the maximum is unique due to convexity and can be written in closed form. However, in many situations, a linear Gaussian model is not appropriate.

For instance, in credit scoring, a bank officer must decide whether or not to grant the requested credit to a prospective client based on a set of descriptive *features*. To make this decision, the officer reviews the bank's database and identifies clients who, in the past, either repaid or defaulted on their loans, in order to determine the profile of a *payer* and a *non-payer*. A new client can then be *classified* as a payer or non-payer based on their similarity to these groups.

Formally, let us denote the client's features by $t \in \mathbb{R}^N$, and assume that the client repays the loan with probability $\sigma(t)$ and defaults with probability $1 - \sigma(t)$, where the function $\sigma(t)$ is to be defined. This is equivalent to constructing the random variable X .

$$X|t \sim \text{Ber}(\sigma(t)), \quad (4.23)$$

where $X = 1$ represent a paying client, and $X = 0$ the opposite. A usual choice for the function $\sigma(\cdot)$ is the logistic function applied to a linear transformation of t , that is,

$$\Pr(X = 1|t) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 t)}}. \quad (4.24)$$

Observe that this is a linear classifier, where $\theta = [\theta_0, \theta_1]$ defines a hyperplane in \mathbb{R}^N where clients such that $t \in \{t | 0 \leq \theta_0 + \theta_1 t\}$ pay with probability greater or equal to $1/2$, and the rest with probability less or equal to $1/2$. This is known as **logistic regression**.

Therefore, using the bank data $\{(x_i, t_i)\}_{i=1}^n$, what is the MLE for $\theta = [\theta_0, \theta_1]$? To answer this, notice that log-likelihood ℓ can be written as

$$\begin{aligned} \ell(\theta) &= \log \prod_{i=1}^n p(x_i|t) \\ &= \sum_{i=1}^n x_i \log \sigma(t) + \left(n - \sum_{i=1}^n x_i\right) \log(1 - \sigma(t)) \\ &= \sum_{i=1}^n x_i \log \left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 t)}} \right) + \left(n - \sum_{i=1}^n x_i\right) \log \left(1 - \frac{1}{1 + e^{-(\theta_0 + \theta_1 t)}} \right). \end{aligned}$$

This expression does not have a global minima and, even though its gradient can be computed, we cannot solve $\partial \ell(\theta) / \partial \theta = 0$ analytically. Finding the MLE thus requires gradient descent.

4.5.3 Latent variables: Expectation-Maximisation

In certain scenarios, the process of interest can be generated by a mixture of sources. For example, consider the distribution of heights in a population: we can naturally model this as a mixture of marginal distributions for male and female heights separately, i.e.,

$$\begin{aligned} \mathcal{L}(\theta) &= \prod_{i=1}^n p \mathcal{N}(X|\mu_H, \Sigma_H) + (1 - p) \mathcal{N}(X|\mu_M, \Sigma_M) \\ &= \prod_{i=1}^n \left(p \frac{1}{\sqrt{2\pi\Sigma_H^{-1}}} \exp \left(\frac{-(x_i - \mu_H)^2}{2\Sigma_H^2} \right) + (1 - p) \frac{1}{\sqrt{2\pi\Sigma_M^{-1}}} \exp \left(\frac{-(x_i - \mu_M)^2}{2\Sigma_M^2} \right) \right), \end{aligned}$$

where p is referred to as the *class-conditional probability*.

Optimising this expression with respect to all five components of θ is difficult, particularly because of the sum inside the product, which prevents simplification via the logarithm.

A key difference of this model compared to previous ones is the implicit introduction of a *latent variable*, that indicates which Gaussian generated each observation. If we knew this latent variable, the problem would become much simpler. Indeed, assume for a moment that we have access to observations $\{z_i\}_{i=1}^n$ of the random variables $\{Z_i\}_{i=1}^n$, where each Z_i indicates which model generated x_i (e.g., $Z_i = 0$ for male, $Z_i = 1$ for female).

In this case, let us consider the **complete data** $\{(x_i, z_i)\}_{i=1}^n$ and write the complete-data log-likelihood as

$$\begin{aligned}\ell(\theta|z_i, x_i) &= \log \prod_{i=1}^n \mathcal{N}(X|\mu_H, \Sigma_H)^{z_i} \mathcal{N}(X|\mu_M, \Sigma_M)^{(1-z_i)} \\ &= \sum_{i=1}^n \left(z_i \log \frac{1}{\sqrt{2\pi\Sigma_H^{-1}}} \exp\left(\frac{-(x_i - \mu_H)^2}{2\Sigma_H^2}\right) + (1 - z_i) \log \frac{1}{\sqrt{2\pi\Sigma_M^{-1}}} \exp\left(\frac{-(x_i - \mu_M)^2}{2\Sigma_M^2}\right) \right).\end{aligned}$$

This objective function is much easier to optimise. However, the observations on which this likelihood builds are not truly *observed*, because the RV Z is latent. To address this, we can take the conditional expectation of the above expression with respect to Z , and then maximise this expression (now independent of Z) with respect to θ , and iterate. Specifically, since the expression is linear in z_i , it suffices to take its expectation and plug it in the above expression.

$$\begin{aligned}\mathbb{E}_\theta(Z_i|\theta_t, x_i) &= 1 \cdot \mathbb{P}(Z_i = 1|\theta_t, x_i) + 0 \cdot \mathbb{P}(Z_i = 0|\theta_t, x_i) \\ &= \frac{\mathbb{P}(x_i|\theta_t, z_i = 1)p(z_i = 1)}{p(x_i|\theta)} \\ &= \frac{\mathbb{P}(x_i|\theta_t, z_i = 1)p(z_i = 1)}{p(x_i|z = 1, \theta)p(z = 1) + p(x_i|z = 0, \theta)p(z = 0)}\end{aligned}$$

4.6 Bayesian inference

We will study a complementary perspective to finding the parameter θ which incorporates a model for the parameter's uncertainty. In the Bayesian paradigm, we will assume that the parameter is a RV Θ that takes a value θ . This allows us to incorporate *a priori* information into the estimation process. A conceptual difference between the two inferential approaches is that the frequentist approach aims to avoid subjectivity, while the Bayesian one builds on the expert knowledge of the scientist to propose and assess possible hypotheses.

Definition 4.11 (A priori distribution). The information, inductive biases, and any other known property of Θ can be encoded into the parameter's probability distribution, referred to as *a priori distribution*, or simply *prior*. We assume this distribution has a density denoted $p(\theta)$.

With the above definition, we can see that the joint density of the RVs X, Θ can be expressed as

$$p(x, \theta) = p(x|\theta)p(\theta), \quad (4.25)$$

where we have written $p(x|\theta)$ instead of $p_\theta(x)$ to emphasise that we now consider the parameter to be a RV.

Definition 4.12 (Marginal distribution). The distribution of X can be obtained through the disintegration of Θ from (X, Θ) , that is

$$p(x) = \int_{\Omega} p(x|\theta)p(\theta)d\theta. \quad (4.26)$$

We refer to this distribution as the marginal distribution of X .

We can now define the central object of the Bayesian inference paradigm.

Definition 4.13 (Posterior distribution). Given a set of observations \mathcal{D} , the *posterior* distribution of the parameter, which considers the information provided by the data \mathcal{D} and the prior knowledge, is given through Bayes theorem by

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})} \propto p(\mathcal{D}|\theta)p(\theta), \quad (4.27)$$

where:

- $p(\theta)$ is the parameter's prior law.
- $p(\theta|\mathcal{D})$ is the parameter's posterior law.
- $p(\mathcal{D}|\theta)$ is the likelihood.
- $p(\mathcal{D}) = \int_{\Omega} p(\mathcal{D}|\theta)p(\theta)d\theta$ is the data marginal density.

The *transition* from prior to posterior can be interpreted as the process of incorporating the evidence provided by the data (through the likelihood function) in order to reduce uncertainty about the value of the parameter Θ . From eq. (4.27) we can see that this process, sometimes referred to as *Bayesian update*, amounts to multiplying by the likelihood and then normalising, ensuring that $p(\theta|\mathcal{D})$ is indeed a probability density.

Remark 4.10. The symbol \propto in eq. (4.27) is used to indicate that the left-hand side is equal to the right-hand side up to a proportionality constant that depends on the \mathcal{D} and not on θ . Therefore, the posterior can be calculated up to a *proportional version*, which can then be refined via normalisation.

Example 4.12 (Posterior distribution for Bernoulli). Let θ be the probability to obtain heads when throwing a coin, and let X_1, \dots, X_n be n i.i.d. realisations of this experiment. Without prior knowledge of θ , it is reasonable to assume a uniform prior $\theta \sim \text{Uniform}(0, 1)$. Recall that this prior encodes our beliefs for the parameter before the experiment. Modelling $X_1, \dots, X_n \sim \text{Ber}(\theta)$, we have:

$$p(X_1, \dots, X_n|\theta) = \prod_{i=1}^n \theta^{X_i} (1 - \theta)^{1-X_i} = \theta^{\sum_{i=1}^n X_i} (1 - \theta)^{n - \sum_{i=1}^n X_i}.$$

Also, $p(X_1, \dots, X_n)$ is given by

$$p(X_1, \dots, X_n) = \int_0^1 \theta^{\sum_{i=1}^n X_i} (1 - \theta)^{n - \sum_{i=1}^n X_i} d\theta = B\left(\sum_{i=1}^n X_i + 1, n - \sum_{i=1}^n X_i + 1\right),$$

where $B(x, y)$ is the beta function:

$$B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}.$$

Let $s = \sum_{i=1}^n X_i$. Then, the posterior is

$$p(\theta|X_1, \dots, X_n) = \frac{p(X_1, \dots, X_n|\theta)}{p(X_1, \dots, X_n)} = \frac{1}{B(s+1, n-s+1)} \theta^s (1 - \theta)^{n-s}.$$

In some cases, the observations x_1, \dots, x_n are available sequentially, or *on line*. Therefore, assuming a prior $p(\theta)$, it is possible to perform on online, or continual, Bayesian update as new data are observed. After seeing x_1 , the posterior $p(\theta|x_1)$ can be computed as:

$$p(\theta|x_1) \propto p(x_1|\theta)p(\theta).$$

Then, after observing x_2 , and given that X_1 and X_2 are conditionally independent given θ , we have:

$$p(\theta|x_1, x_2) \propto p(x_2|\theta)p(\theta|x_1) \propto p(x_1|\theta)p(x_2|\theta)p(\theta).$$

In the general case we have

$$p(\theta|x_1, \dots, x_n) \propto p(x_n|\theta)p(\theta|x_1, \dots, x_{n-1}) \propto p(\theta) \prod_{i=1}^n p(\theta|x_i).$$

Remark 4.11. In the online Bayesian update, the n -stage posterior is the $n+1$ -stage prior.

include
a plot:
prior,
likeli-
hood,
posterior.

4.7 Conjugate priors

Bayesian updating may yield a posterior that is known only up to a proportionality constant (when the marginal distribution $p(x)$ cannot be computed), or a posterior that does not belong to any familiar distributional family. A key tool that guarantees the computability of the posterior distribution, including its normalising constant, and ensures that it takes a known form is the use of **conjugate priors**.

Definition 4.14. Consider a model with likelihood $p(x|\theta)$ and a prior $p(\theta)$. We say that $p(\theta)$ is conjugate to the likelihood $p(x|\theta)$ if the posterior

$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)} \quad (4.28)$$

is in the *same family* as the prior $p(\theta)$. This means that both prior and posterior have the same functional form, e.g., $f_\lambda(\theta)$, but with different parameter λ .

Example 4.13 (continuation of Ex. 4.12). Verify if the prior in Ex. 4.12 is indeed a conjugate prior.

Example 4.14 (Multinomial distribution). Consider a (k -dimensional) multinomial random variable $X \sim \text{Mult}(n, \theta)$. The i -th component of this variable models the number of times that event i occurs among k possible events over n realisations. The parameter of this distribution is known as the *probability vector*, denoted θ , which belongs to the simplex

$$\{\theta \in [0, 1]^k : \theta_1 + \dots + \theta_k = 1\}.$$

For example, if we roll a fair die 100 times, the vector containing the count of occurrences of each face can be modelled as

$$X_{\text{die}} \sim \text{Mult}\left(100, \left[\frac{1}{6}, \dots, \frac{1}{6}\right]\right).$$

Let $X = [x_1, \dots, x_k]$. Any multinomial sample satisfies

$$x_i \in \{0, 1, \dots, n\}, \quad \sum_{i=1}^k x_i = n.$$

The multinomial distribution is given by

$$\text{Mult}(X; n, \theta) = \frac{n!}{x_1! \dots x_k!} \theta_1^{x_1} \dots \theta_k^{x_k},$$

and it generalises the following distributions:

- the *Bernoulli* distribution when $k = 2$ and $n = 1$;
- the *categorical* (or *multinoulli*) distribution when $n = 1$;
- the *binomial* distribution when $k = 2$.

Observe that the parameter θ for the multinomial (and the other related three distributions) is a discrete probability function. This means that constructing the prior $p(\theta)$ implies the definition of a probability distribution over probability distributions.

Definition 4.15 (Dirichlet distribution). Let us consider the distribution

$$\theta \sim \text{Dir}(\theta|\alpha) = \frac{1}{B(\alpha)} \prod_{i=1}^k \theta_i^{\alpha_i-1}, \quad (4.29)$$

where $\alpha = (\alpha_1, \dots, \alpha_k)$ is a concentration parameter and the normalising constant is given by $B(\alpha) = \prod_{i=1}^k \Gamma(\alpha_i) / \Gamma(\sum_{i=1}^k \alpha_i)$. The support of this distribution is the simplex.

For $k = 3$, Dirichlet's distribution can be plotted in 2 dimensions. Fig. 18 shows three plots for different values of the concentration parameter.

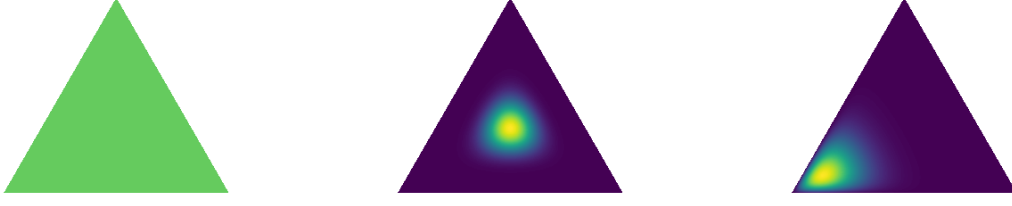


Fig. 18. Dirichlet distribution for $k = 3$ with concentration parameters α (left to right) $[1, 1, 1]$, $[10, 10, 10]$ y $[10, 2, 2]$.

We will see now that the Dirichlet distribution is conjugate to the Multinomial, and, as a consequence, to Bernoulli, Categorical and Binomial. In fact, if $\theta \sim \text{Dir}(\theta; \alpha)$ y $X \sim \text{Mult}(X; n, \theta)$, then

$$\begin{aligned} p(\theta|x) &= \frac{\text{Mult}(x; n, \theta) \text{Dir}(\theta; \alpha)}{p(x)} \\ &= \frac{n!}{x_1! \cdots x_k! p(x) B(\alpha)} \prod_{i=1}^k \theta_i^{x_i + \alpha_i - 1} \\ &= \frac{1}{B(\alpha')} \prod_{i=1}^k \theta_i^{\alpha'_i - 1} \end{aligned} \quad (4.30)$$

where $\alpha' = (\alpha'_1, \dots, \alpha'_k) = (\alpha'_1 + x_1, \dots, \alpha'_k + x_k)$ is the updated concentration parameter.

Remark 4.12. How difficult was it to compute the normalising constant in the above example?

Example 4.15. Let us consider $\alpha = [1, 2, 3, 4, 5]$ and generate a sample from $\theta \sim \text{Dir}(\theta|\alpha)$. The following snippet generates, plots and prints this sample.

```
1 import numpy as np
2 alpha = np.array([1,2,3,4,5])
3 theta = np.random.dirichlet(alpha)
4 plt.bar(np.arange(5)+1, theta);
5 print(f'theta = {theta}')
```

The obtained parameters are $\theta = [0.034, 0.171, 0.286, 0.185, 0.324]$.

Now, using a Dirichlet prior over θ con $\alpha_p = [1, 1, 1, 1, 1]$, let us calculate the posterior according to eq. (4.30). Fig. 19 shows 50 samples from the posterior for different amounts of observations between 0 and 10^5 .

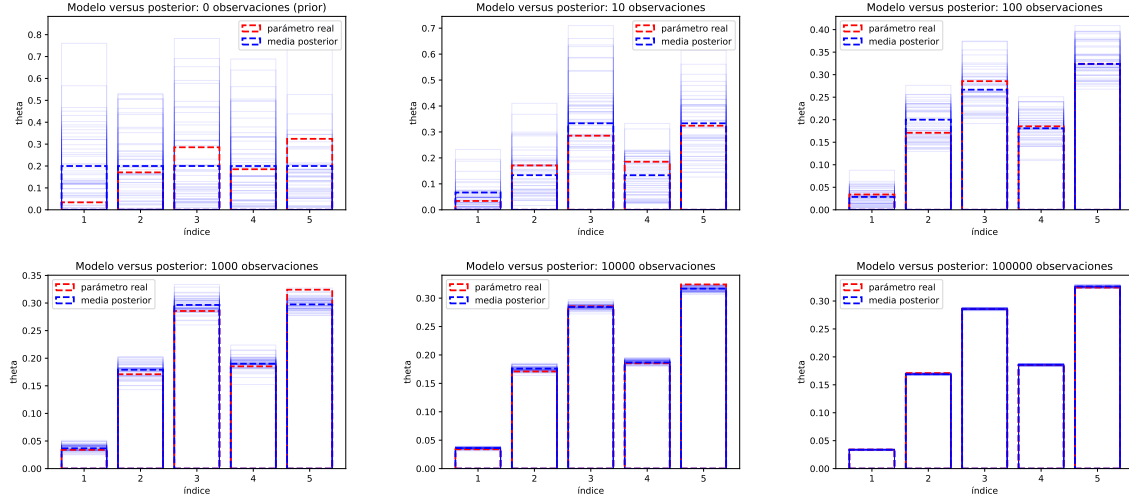


Fig. 19. Posterior concentration around the real parameter for a model $X \sim \text{Mult}(\theta)$ and the prior Dirichlet distribution $\theta \sim \text{Dir}(\alpha)$. Observations considered ranged between 0 and 10^5 , and each plot whows the the true parameter in (broken red line), the posterior mean (broken blue line) and 50 samples from thr posterior (light blue). Observe how the effect of the prior (broken blue line in the first plot) vanishes as the number of observations grows.

Example 4.16. Gaussian model (known σ^2). Let us consider observations x_1, \dots, x_n from a Gaussian pdf $\mathcal{N}(\mu, \sigma^2)$, and consider the prior on the mean $p(\mu) = \mathcal{N}(\mu_0, \sigma_0^2)$, in which case the posterior is given by

$$p(\mu|\mathcal{D}) \propto \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu)^2\right) \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp\left(-\frac{1}{2\sigma_0^2}(\mu - \mu_0)^2\right) \quad (4.31)$$

$$\propto \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 - \frac{1}{2\sigma_0^2}(\mu - \mu_0)^2\right), \quad (4.32)$$

where the proportionality arises from ignoring the constant $p(\mathcal{D})$ in the first line, and from disregarding all constants that do not depend on μ in the second line. Recall that these constants for μ include the variance of x , σ^2 , so omitting this quantity is only possible because we are considering the case in which σ^2 is known. Completing the square in μ inside the exponential in Eq. (4.32), we obtain (we shall define μ_n and σ_n^2 shortly)

$$p(\mu|\mathcal{D}) \propto \exp\left(-\frac{1}{2\sigma_n^2}(\mu - \mu_n)^2\right). \quad (4.33)$$

Since $p(\mu|\mathcal{D})$ must integrate to one, the only probability density proportional to the right-hand side of the previous equation is the Gaussian with mean μ_n and variance σ_n^2 . That is, the proportionality constant required for equality in the expression above is

$$\int_{\mathbb{R}} \exp\left(-\frac{1}{2\sigma_n^2}(\mu - \mu_n)^2\right) d\mu = (2\pi\sigma_n^2)^{n/2}. \quad (4.34)$$

Consequently, we confirm that the chosen prior is indeed conjugate to the Gaussian likelihood, so the posterior is given by the following (Gaussian) density:

$$p(\mu|\mathcal{D}) = \mathcal{N}(\mu; \mu_n, \sigma_n^2) = \frac{1}{(2\pi\sigma_n^2)^{N/2}} \exp\left(-\frac{1}{2\sigma_n^2}(\mu - \mu_n)^2\right), \quad (4.35)$$

re-run
this plot
with En-
glish la-
bels

where the mean and variance are given, respectively, by

$$\mu_n = \frac{1}{\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}} \left(\frac{1}{\sigma_0^2} \mu_0 + \frac{n}{\sigma^2} \bar{x} \right), \quad \text{donde } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (4.36)$$

$$\sigma_n^2 = \left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right)^{-1}. \quad (4.37)$$

Remark 4.13. Bayesian updating transforms the parameters of the prior on μ from μ_0 and σ_0^2 to μ_n and σ_n^2 in Eqs. (4.36) and (4.37), respectively. Notice that the posterior parameters are combinations (which can be given a natural interpretation) of the prior parameters and the observed data. In particular, μ_n is a weighted average of μ_0 (our prior candidate for μ) with weight σ_0^{-2} and the sample mean \bar{x} with weight $(\sigma^2/n)^{-1}$, which is the maximum likelihood estimator.

It is also worth noticing that these weights correspond to the inverse variances, i.e., the precisions, of μ_0 and \bar{x} . Finally, observe that σ_n^2 represents a *parallel sum* of the variances: if we rewrite Eq. (4.37) in terms of precisions, we see that the initial precision σ_0^{-2} is incremented by a contribution from each observed data point, which makes sense because with more information the precision should increase rather than the uncertainty (here represented by the variance).

4.8 Bayesian estimators

So far, we have studied the role of the prior in Bayesian inference, however, we have not leveraged this setup for building estimators. In particular, recall that the MLE does not incorporate a priori knowledge of the parameter. This can be addressed from a Bayesian standpoint.

Definition 4.16 (Bayesian risk). For a loss function $L(\theta, \hat{\theta})$ and observations \mathcal{D} , the Bayesian risk is the posterior expectation of the loss function, that is

$$R(\hat{\theta}) = \int_{\Omega} L(\theta, \hat{\theta}) p(\theta | \mathcal{D}) d\theta. \quad (4.38)$$

Definition 4.17 (Bayesian estimator). Given observations \mathcal{D} and the Bayesian risk $R(\theta)$, an estimator that minimises the Bayesian risk

$$\theta_{\text{Bayes}}(\mathcal{D}) = \arg \min_{\Omega} R(\theta) \quad (4.39)$$

is known as Bayesian estimator. We have ignored the explicit dependence of $R(\cdot)$ in \mathcal{D} .

We will now define Bayesian estimators corresponding to different loss or risk functions.

Definition 4.18 (Bayesian model average). The standard case is the quadratic loss function $L_2(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$, which leads to the estimator given by the posterior mean $\theta_{\text{Bayes}}(\mathcal{D}) = \mathbb{E}[\theta | \mathcal{D}]$.

Definition 4.19 (Minimum absolute-error (MAE)). Similarly, the cost function $L_1(\theta, \hat{\theta}) = |\theta - \hat{\theta}|_1$ results in the estimator that gives the posterior median.

Finding a loss function for the maximum a posteriori (MAP) estimator is less straightforward. Let us first consider the case where $\theta \in \Omega$ is discrete and the “0-1” loss.

$$L_{0-1}(\theta, \hat{\theta}) = \begin{cases} 0 & \text{if } \theta = \hat{\theta}, \\ 1 & \text{if not.} \end{cases}$$

The Bayes risk associated to $L_{0-1}(\theta, \hat{\theta})$ (in the discrete case) takes the form

$$R(\hat{\theta}) = \mathbb{P}(\theta \neq \hat{\theta} | \mathcal{D}) = 1 - \mathbb{P}(\theta = \hat{\theta} | \mathcal{D}), \quad (4.40)$$

which is minimised by choosing $\hat{\theta}$ such that $\mathbb{P}(\theta = \hat{\theta} | \mathcal{D})$ is maximised, i.e., the MAP.

Definition 4.20 (Maximum a posteriori estimator). Let $\theta \in \Theta$ be a parameter with posterior $p(\theta|D)$ defined over Θ . We refer to its point estimate given by :

$$\theta_{MAP} = \arg \max_{\theta \in \Theta} p(\theta|D),$$

as the *maximum a posteriori* (MAP) estimator.

Remark 4.14. Notice that it is possible to find the MAP estimate using only a *proportional* version of the posterior, as is common in Bayesian inference, or equivalently by maximising its logarithm. Additional, notice the relationship between the MLE and MAP given by

$$\theta_{MAP} = \arg \max_{\theta \in \Theta} p(\theta|D) = \arg \max_{\theta \in \Theta} p(D|\theta)p(\theta) = \arg \max_{\theta \in \Theta} \left(\underbrace{\log p(D|\theta)}_{l(\theta)} + \log p(\theta) \right).$$

Remark 4.15. Observe that the MAP estimator is a *modification* of the MLE, since both share part of the same objective function (the likelihood), with the difference that the MAP additionally includes the *log-prior* term. This can be understood as a form of regularisation of the MLE solution, where the additional term encodes properties of the estimator that may not be revealed by the data alone.

Example 4.17 (MAP for the mean of the Gaussian model). For the linear Gaussian model considered so far, we can compute θ_{MAP} for a Gaussian prior with zero mean and variance σ_θ^2 . This is given by (assuming the noise variance σ_ϵ^2 is known):

$$\begin{aligned} \theta_{MAP}^* &= \operatorname{argmax} p(Y|\theta, X)p(\theta) \\ [\text{ind., def.}] &= \operatorname{argmax} \prod_{i=1}^N \mathcal{N}(y_i; \theta^\top x_i, \sigma_\epsilon^2) \mathcal{N}(\theta; 0, \sigma_\theta^2) \\ &= \operatorname{argmax} \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma_\epsilon} \exp\left(\frac{-1}{2\sigma_\epsilon^2}(y_i - \theta^\top x_i)^2\right) \frac{1}{(\sqrt{2\pi}\sigma_\theta)^{M+1}} \exp\left(\frac{-\|\theta\|^2}{2\sigma_\theta^2}\right) \\ &= \operatorname{argmax} \frac{1}{\sqrt{2\pi}\sigma_\epsilon} \frac{1}{(\sqrt{2\pi}\sigma_\theta)^{M+1}} \exp\left(\sum_{i=1}^N \frac{-1}{2\sigma_\epsilon^2}(y_i - \theta^\top x_i)^2 - \frac{\|\theta\|^2}{2\sigma_\theta^2}\right) \\ [\log.] &= \operatorname{argmin} \sum_{i=1}^N (y_i - \theta^\top x_i)^2 + \frac{\sigma_\epsilon^2}{\sigma_\theta^2} \|\theta\|^2. \end{aligned}$$

We can see that by choosing a uniform prior or a Gaussian prior with very large variance, the MAP is equivalent to the MLE. What does this mean? What different behaviour does the MAP promote compared to the MLE in this case?

4.8.1 Posterior predictive

In Bayesian inference, predictions play an important role, since after performing inference on a statistical model, we are generally interested in studying how future data will be generated by the model. To this end, we define the Bayesian prediction as

Definition 4.21 (Posterior predictive). For a dataset \mathcal{D} and a parameter θ , the posterior predictive density is given by

$$p(x|\mathcal{D}) = \int_{\Omega} p(x|\theta)p(\theta|\mathcal{D})d\theta = \mathbb{E}[p(x|\theta)|\mathcal{D}], \quad (4.41)$$

that is, the expected value of the statistical model with respect to the posterior distribution of the parameter (model).

We can now consider the posterior predictive as our *learned* model and generate data from it, where we face the same dilemma as with a point estimator in the previous case: it is possible to consider random samples, the mean, the median, or some interval.

illustration
if possible

Remark 4.16. The posterior predictive is generally different from the *plug-in* prediction, in which we consider the statistical model $p_{\hat{\theta}}$ based on some (point) estimator $\hat{\theta}$. From this perspective, the posterior predictive is equivalent to considering point estimators and models, but integrating over all of them with respect to the posterior distribution.

an example of posterior predictions vs plug-in predictions - there might be one in the ML lecture notes

References

- Bertsekas, D. P., & Tsitsiklis, J. N. (2008). *Introduction to probability* (2nd ed.). Belmont, MA: Athena Scientific. Retrieved from <https://athenasc.com/probbook.html>
- Boyd, S., & Vandenberghe, L. (2004). *Convex optimization*. Cambridge, UK: Cambridge University Press. Retrieved from <https://web.stanford.edu/~boyd/cvxbook/>
- Deisenroth, M. P., Faisal, A. A., & Ong, C. S. (2020). *Mathematics for machine learning*. Cambridge University Press. Retrieved from <https://mml-book.com/>
- Durrett, R. (2019). *Probability: Theory and examples* (5th ed.). Cambridge, UK: Cambridge University Press. Retrieved from <https://doi.org/10.1017/9781108591034>
- Grimmett, G. R., & Stirzaker, D. R. (2020). *Probability and random processes* (4th ed.). Oxford, UK: Oxford University Press.
- Kingma, D. P., & Ba, J. (2015). Adam: A method for stochastic optimization. *International Conference on Learning Representations (ICLR)*. Retrieved from <https://arxiv.org/abs/1412.6980>
- Murphy, K. P. (2022). *Probabilistic machine learning: An introduction*. MIT Press. Retrieved from probml.ai
- Nesterov, Y. (2018). *Lectures on convex optimization* (2nd ed.). Cham, Switzerland: Springer International Publishing. doi: 10.1007/978-3-319-91578-4