
MATHEMATICS FOR MACHINE LEARNING

This version: November 28, 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
4.4	Unbiased estimators	35
4.5	Loss functions	37
4.6	Maximum likelihood estimator (MLE)	38
4.7	MLE in practice: three examples	39
4.7.1	Gaussian linear regression	39
4.7.2	Classification	40
4.7.3	Latent variables: <i>Expectation-Maximisation</i>	40
4.8	Enfoque bayesiano	41

4.9	Contexto y definiciones principales	41
4.10	Priors Conjugados	44
4.11	Estimadores bayesianos	48
4.12	Posterior predictiva	49
References		50

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.

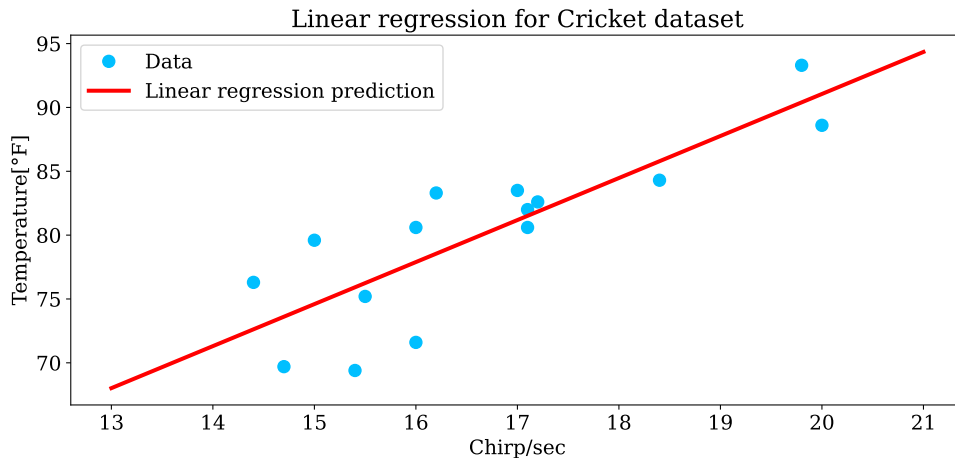


Fig. 1. Example of a linear regression model minimising the least squares loss.

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.

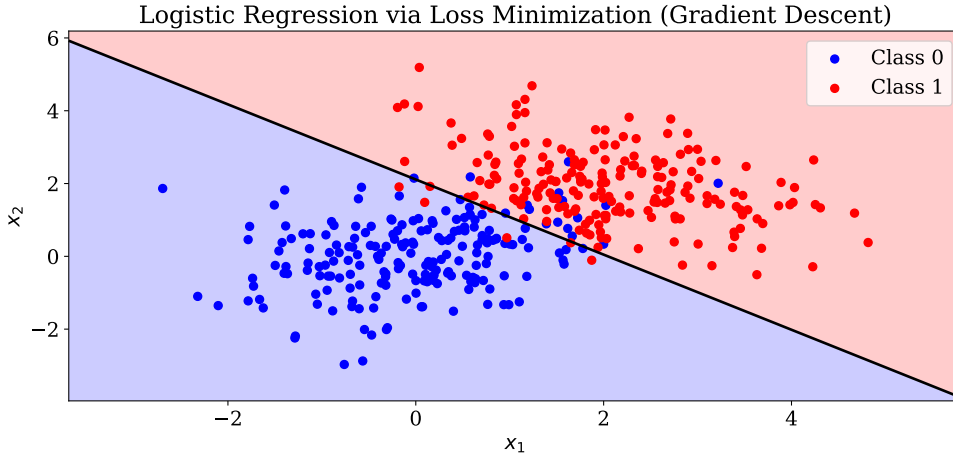


Fig. 2. Example of a logistic regression model to classify data from two Gaussians.

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:

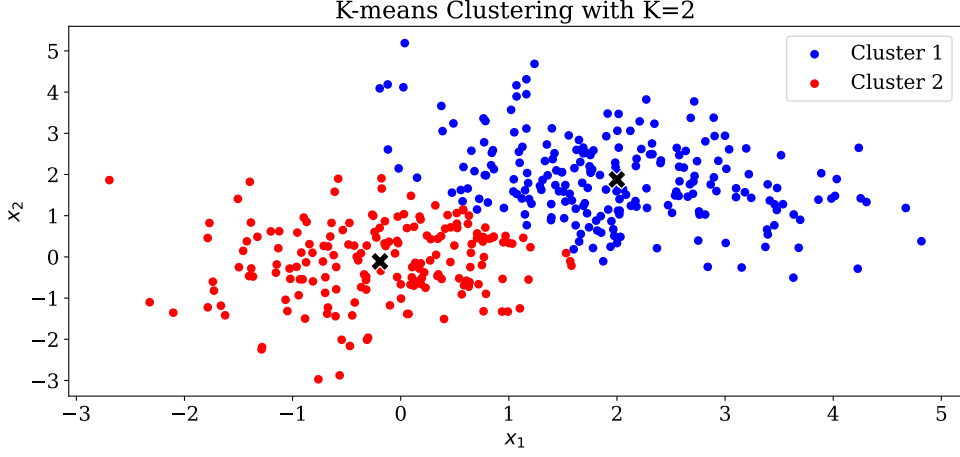


Fig. 3. Implementation of the K -means method to clustering data from two Gaussians. The figure shows the learnt centroids (black cross) and the colours correspond to the cluster assignments.

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

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 $\{(x, y) : x = \pm 1 \wedge y = \pm 1\}$.

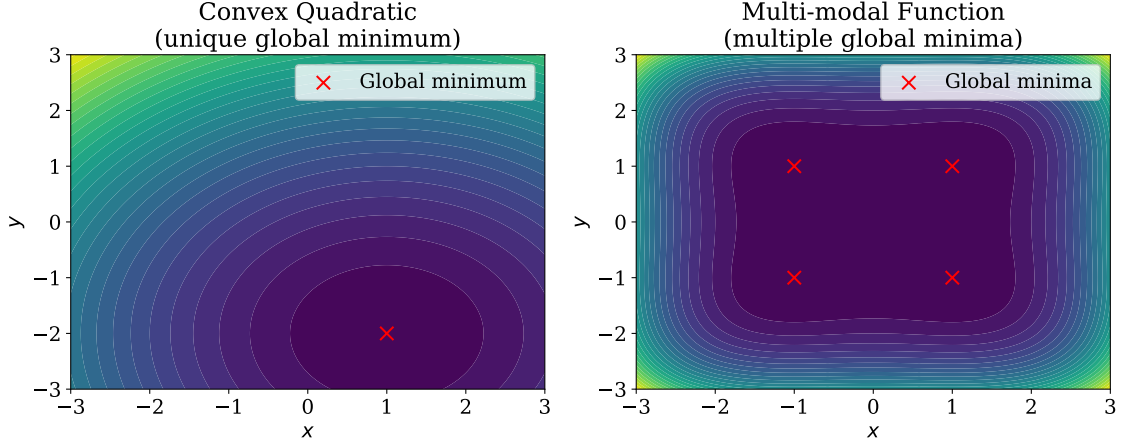


Fig. 4. Functions with a unique (left) and non-unique (right) global minima.

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.

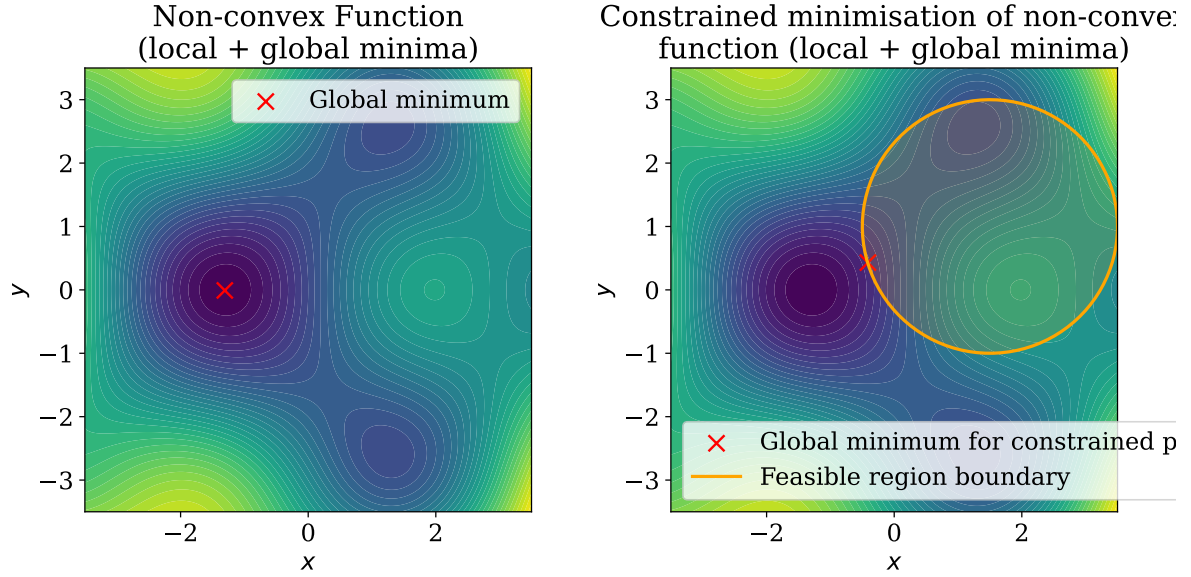


Fig. 5. Minima for constrained (left) and unconstrained (right) problems.

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.

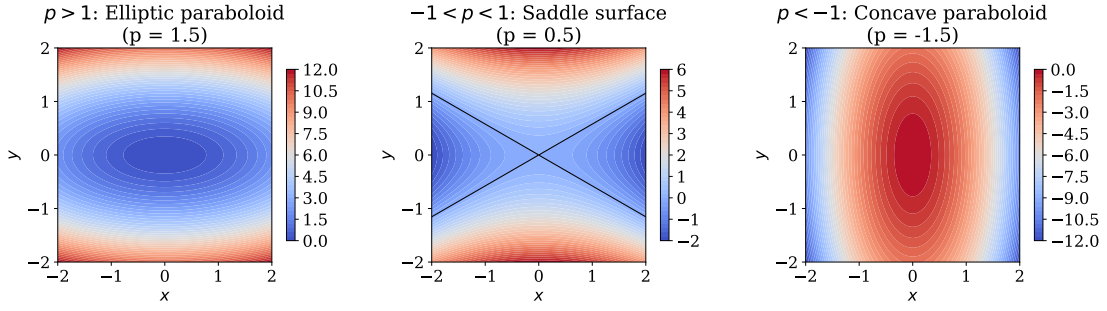


Fig. 6. Different types of critical points for $f(x, y) = (p - 1)x^2 + (p + 1)y^2$.

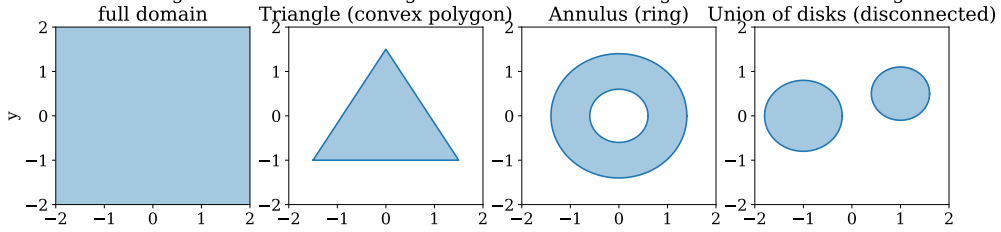


Fig. 7. Examples of convex and non-convex sets, highlighted in blue.

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.

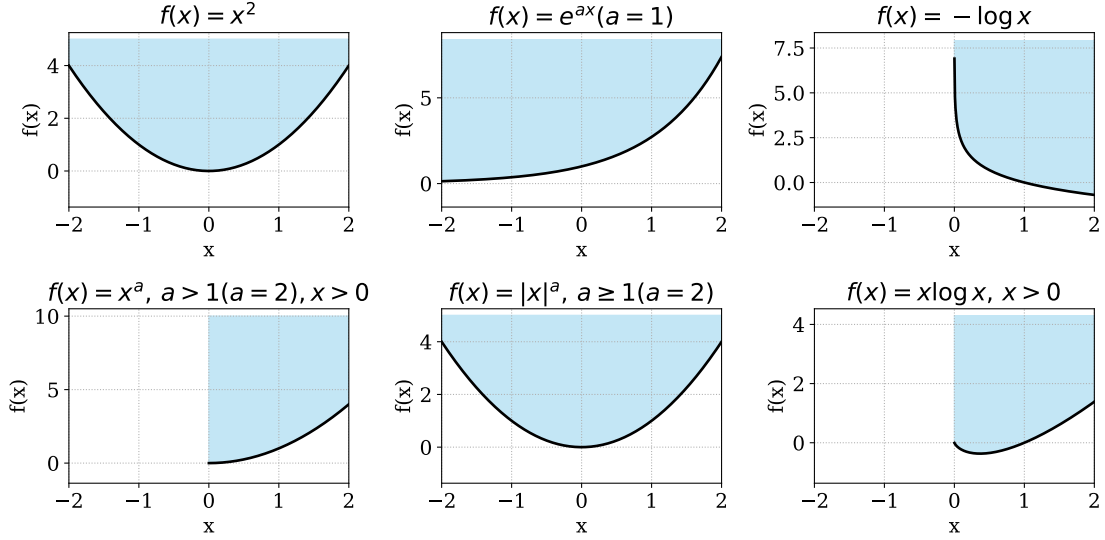


Fig. 8. Convex 1D functions with their epigraph in light blue.

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).

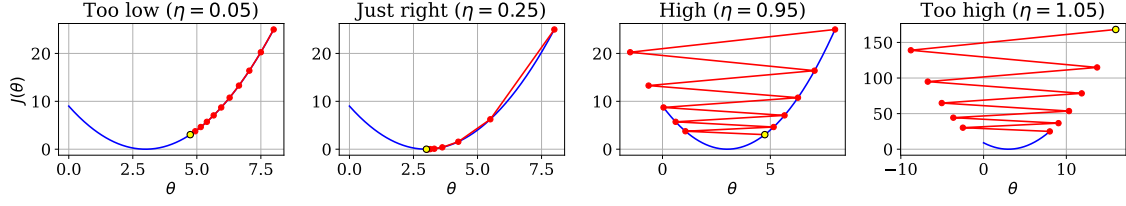


Fig. 9. Gradient based optimisation with different learning rates.

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.

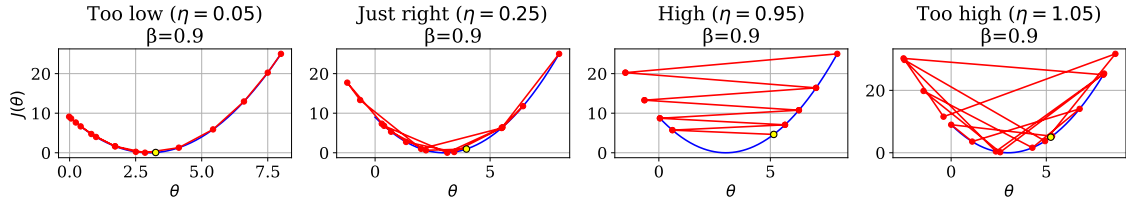


Fig. 10. Gradient based optimisation with momentum.

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.

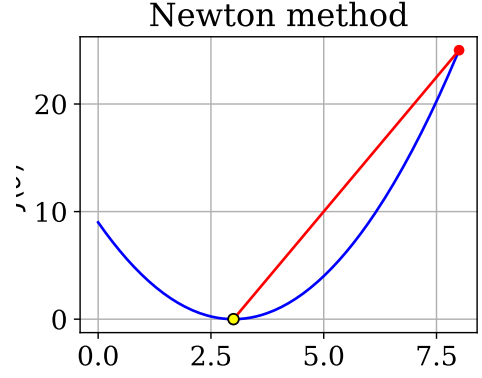


Fig. 11. One-step parabola minimisation with the Newton method.

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.

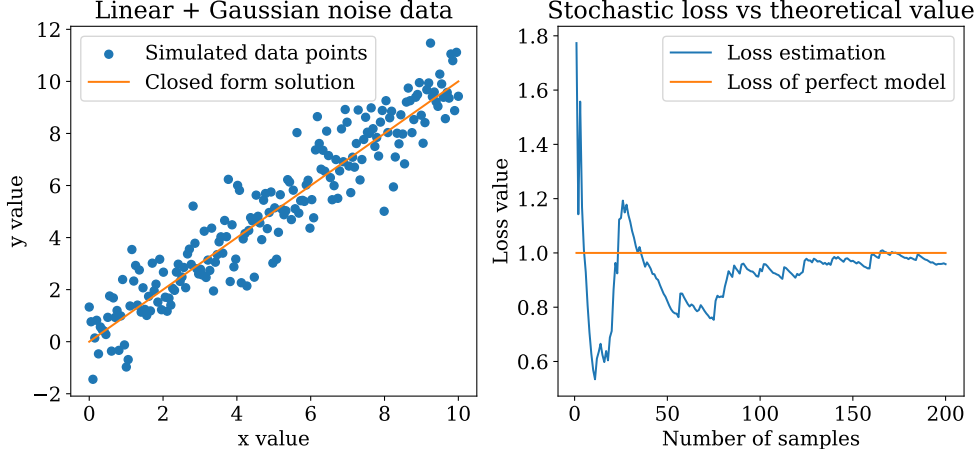


Fig. 12. Estimation of the true loss via the empirical value from data.

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.

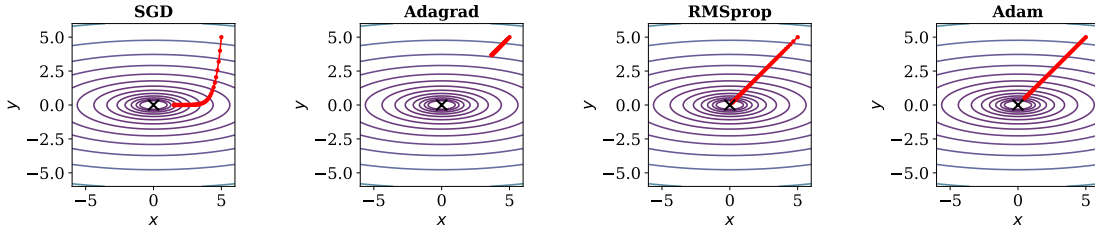


Fig. 13. Convergence comparison for different optimisers.

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.

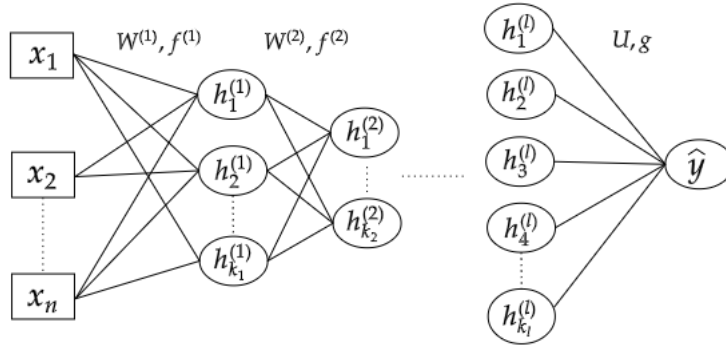


Fig. 14. A neural network of depth l .

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 $b^{(k)} \leftarrow b^{(k)} - \lambda \frac{\partial J}{\partial b^{(k)}}$.
-

Observación 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

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

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

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 .

Observación 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.

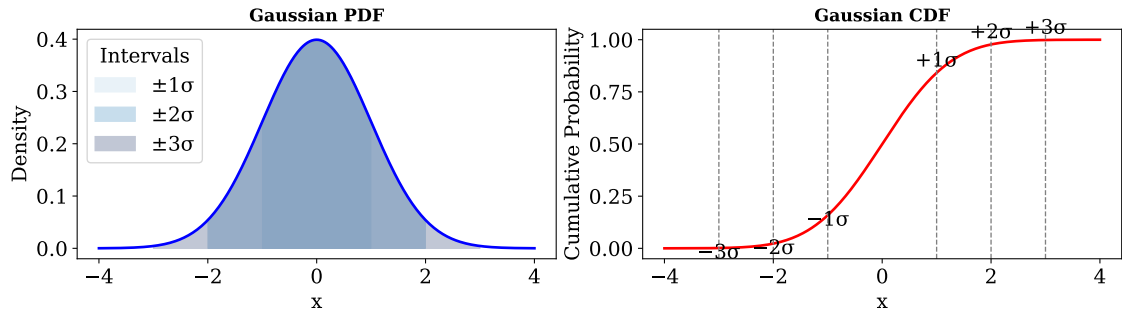


Fig. 15. Probability and cumulative density function.

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.

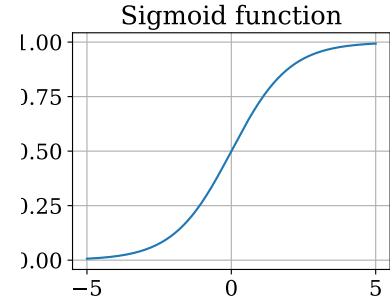


Fig. 16. Visualisation of a sigmoid function.

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.

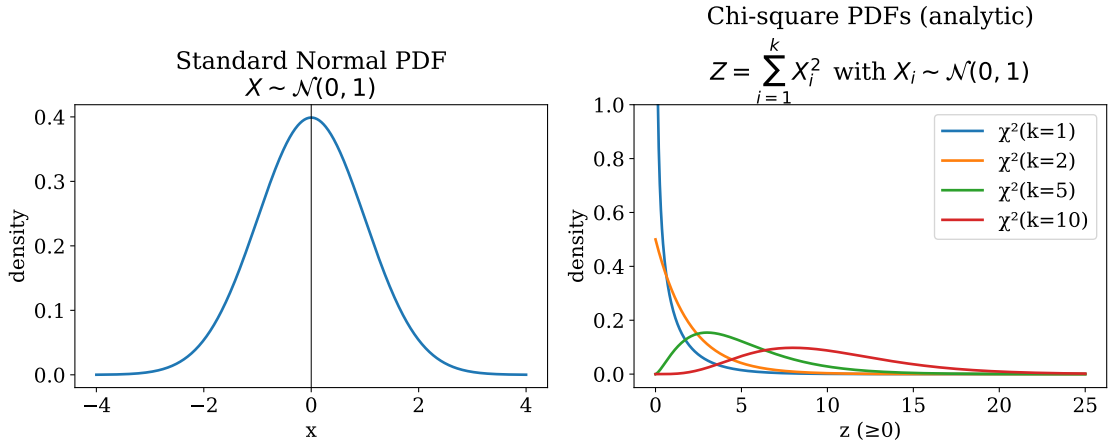


Fig. 17. Gaussian distribution and χ^2 distribution for different k parameters.

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 availability of a dataset generated from a *statistical model*, (aka, probabilistic or generative model) which is unknown. Here, our objective is to use the data to determine such models (and in particular its parameters) to ultimately learn the underlying properties of the data generating process and make predictions. A first step in this regard, is to define our *statistical model*.

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, parametrised by its mean and variance. In such cases, inference boils down to identifying the parameters.

In this part of the course, we will consider a dataset x belonging to an abstract space \mathfrak{X} , where typically $\mathfrak{X} = \mathbb{R}^n$. We will assume that x is the realisation of an RV $X \in \mathfrak{X}$. We can understand the statistical model as the space of the possible hypotheses explaining the observed data. In this sense, a relevant question is: what is the law of X ?

In this course we will focus on 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 now denote the parametric family \mathcal{P} as follows:

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

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} .

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 problems. 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. Of course, we seek to achieve both goals statistically, that is, by appropriately modelling the associated uncertainty.

4.2 Statistics

The initial setup in statistical inference features an observation dataset and our own assumptions. Therefore, the starting point in this regard 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 $\mathcal{P} = \{P_\theta \text{ t.q. } \theta \in \Theta\}$. A statistic is a function of the realisation $X = x$ that is independent from the parameter θ (and the distribution P_θ).

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

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

Observación 4.1. It is critical to understand the difference between the value of the statistic $T(x)$ as a function of the data x (which we consider to be the realization $X = x$ of the random variable) and the application of the function $T(\cdot)$ to the RV X , i.e., $T(X)$ —which is also a RV. The former is a “fixed” value, while the latter is a random variable with its own probability distribution induced by P_θ and by the function T (called the pushforward distribution $T_{\#}P_\theta$).

Some example statistics are:

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

Suficiencia 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 or some other property. Therefore, the identity function or the mean seem to fulfil this mission, at least intuitively. This is because, in practice, we want to extract as much information as possible from the data, which is achieved by the statistic T (which summarizes all the data) and the statistic T' (which contains all the data). On the contrary, note that the statistic T'' loses information, since only the minimum value of all the data obtained is extracted, thus losing the representation of, e.g., the dispersion of the sample. The same analysis can be made for the constant statistic, which contains no information from the data.

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)

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 $\mathcal{P} = \{P_\theta \text{ t.q. } \theta \in \Theta\}$. 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,

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

Let us consider the following examples.

Trivial sufficient statistic

For any 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 .

Ejemplo 4.1 (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 (by definition), since

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

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

In practice, sufficiency is not 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 $g(\theta)$.

Observación 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 $g(\cdot)$.

Observación 4.3. Estimators can be used to identify parameters, where $g(\theta) = \theta$, 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

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

Ejemplo 4.2 (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

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

4.4 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

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

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.

Observación 4.4. Unbiased estimators are important in statistical theory and practice because 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.

Ejemplo 4.3 (Unbiased estimator for the Gaussian mean). The estimator of $g(\theta) = \mu$ described in Example 4.2 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.

Ejemplo 4.4 (Pythagoras). Let us consider a parametric family $\mathcal{P} = \{P_\theta \text{ t.q. } \theta \in \Theta\}$, 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

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

Therefore, the bias of the estimator in eq. (4.3) is asymptotically unbiased. However, we can fix this bias by multiplying the estimator S_2 in eq. (4.3), by $n/(n+1)$. This results in a corrected estimator given by

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

for which it holds

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

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

4.5 Loss functions

We will now consider loss functions for parameters, that is, a function that determines 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:

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

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

Ejemplo 4.5 (Quadratic loss function). A widely used loss function for comparing estimators is the **quadratic loss**, given by

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

Discussion: why do we use the exponent 2 and not other?

Ejemplo 4.6 (0 – 1 loss). When estimating parameters in a space with no order relation, the 0 – 1 loss is usually considered. This is given by

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

Since the estimator is a RV, so is the loss and we can calculate its expectation. This 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. 4.5 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.

4.6 Maximum likelihood estimator (MLE)

Since the parameter θ is unknown, calculating the loss related to an estimator $\hat{\theta} = \hat{\theta}(X)$ is not possible. To bypass the definition of a loss, 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 likelihood function defined as follows will be crucial

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

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

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

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

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

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} \mathcal{L}(\theta | x). \quad (4.14)$$

Observación 4.5 (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.

Observación 4.6. Clearly, 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. In particular, we will focus on finding $\hat{\theta}_{\text{MLE}}$ by maximising the log-likelihood $l(\theta) = \log L(\theta)$, which is usually easier to optimise computationally or analytically. In fact, we will often even ignore constants in the (log) likelihood, as these do not affect the maximiser of $L(\theta)$.

Ejemplo 4.7 (MLE: Bernoulli). Let $X_1, \dots, X_n \sim \text{Ber}\theta$, the likelihood of θ is given by

$$L(\theta) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i}, \quad (4.15)$$

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 $\frac{\partial l(\theta)}{\partial \theta} = 0$:

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

Ejercicio 4.1. Plot $l(\theta)$ in ex. 4.7.

Ejercicio 4.2. Find the MLE for $\theta = (\mu, \Sigma)$ for the RV $X \sim \mathcal{N}(\mu, \Sigma)$.

Ejemplo 4.8. Let $X \sim \text{Uniforme}(\theta)$, that is, $p(x) = \theta^{-1} \mathbf{1}_{0 \leq x \leq \theta}$. To calculate the likelihood function, observe that

$$L(\theta) = \prod_{i=1}^n p_{\theta}(x_i), \quad (4.16)$$

and note that $p_{\theta}(x_i) = 0$ si $x_i > \theta$. Therefore, $L(\theta) > 0$ only if θ is greater than each observation, in particular, if $\theta \geq \max\{x_i\}_1^n$.

Additionally, if we have $\theta \geq \max\{x_i\}_1^n$, then notice that $p_{\theta}(x_i) = 1/\theta$. As a consequence, the likelihood is given by

$$L(\theta) = \theta^{-n}, \quad \theta \geq \max\{x_i\}_1^n \quad (4.17)$$

and the MLE is $\theta_{\text{MV}} = \max\{x_i\}_1^n$.

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.

4.7 MLE in practice: three examples

4.7.1 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

$$X \sim \mathcal{N}(\theta_0 + \theta_1 t^2 + \theta_2 \cos(2\pi t/12), \theta_3^2), \quad (4.18)$$

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

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

$$\begin{aligned} l(\theta) &= \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\theta_3^2}} \exp \left(-\frac{(x_i - \theta_0 - \theta_1 t_i^2 - \theta_2 \cos(2\pi t_i/12))^2}{2\theta_3^2} \right) \right) \\ &= \frac{n}{2} \log(2\pi\theta_3^2) - \frac{1}{2\theta_3^2} \sum_{i=1}^n (x_i - \theta_0 - \theta_1 t_i^2 - \theta_2 \cos(2\pi t_i/12))^2 \end{aligned} \quad (4.19)$$

where we can see that θ_{MV} can be calculated explicitly and it is a function of $\{(t_i, x_i)\}_{i=1}^n$, since eq. (4.19) is quadratic in $[\theta_0, \theta_1, \theta_2]$.

4.7.2 Classification

The reason why $\hat{\theta}_{MLE}$ could be calculated explicitly is that the Gaussian model with a linearly parametrised mean leads to a quadratic log-likelihood, where the maximum is unique and can be written in closed form. However, in many situations, a linear Gaussian model is not appropriate.

One example of this is the problem of credit scoring, where, based on a set of *features* describing a client, a bank officer must decide whether or not to grant the requested credit. To make this decision, the officer can review the bank's database and identify 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 .

$$X|t \sim \text{Ber}\sigma(t) \quad (4.20)$$

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,

$$\Pr(X = 1|t) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 t)}}. \quad (4.21)$$

Observe that this is a linear classifier, where $\theta = [\theta_0, \theta_1]$ defines a hyperplane in \mathbb{R}^N where clients $t \in \{t | 0 \leq \theta_0 + \theta_1 t\}$ 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 ℓ can be written as

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

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

4.7.3 Latent variables: *Expectation-Maximisation*

In certain scenarios, it is natural to assume that our data come from a mixture of models. 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.,

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

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

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

This objective function is much easier to optimise; however, it is not *observed* because the random variable Z is not available. One way to resolve this is to take the conditional expectation of the above expression with respect to Z , given the observed data and the current parameter estimates, then maximise this expectation with respect to θ , and iterate. Specifically, since the expression is linear in z_i , it suffices to take its expectation:

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

4.8 Enfoque bayesiano

En esta sección complementaremos el enfoque visto hasta ahora en cuanto a la incorporación de un modelo para la incertidumbre asociada al parámetro θ . En el paradigma bayesiano, consideraremos que el parámetro es una variable aleatoria, es decir, Θ , la cual para una realización particular tomar el valor $\Theta = \theta$. Esto nos permite información *a priori* sobre la estimación a realizar, lo que permite, en muchos casos, ayudar a la inferencia. Una diferencia conceptual entre ambos enfoques, es que la estadística frecuentista evita la subjetividad, mientras que la estadística bayesiana se basa en la convicción del(a) investigador(a), para emitir juicios sobre una hipótesis.

En este capítulo, se estudiará la estadística bayesiana y se introducirán los mismos conceptos vistos anteriormente desde el punto de vista bayesiano.

4.9 Contexto y definiciones principales

Definition 4.11 (Distribución a priori). La información, sesgos y cualquier otra característica conocida de Θ codificadas mediante la propia ley de probabilidad de esta VA, la cual tiene densidad $p(\theta)$, nos referimos a esta como la *densidad a priori* o simplemente *prior*.

Con esta definición, podemos ver que la densidad conjunta de las VAs X, Θ pueden ser expresadas combinando la densidad a priori con el modelo visto en las secciones anteriores, es decir,

$$p(x, \theta) = p(x|\theta)p(\theta) \quad (4.22)$$

donde hemos escrito $p(x|\theta)$ en vez de $p_\theta(x)$ para hacer explícito que ahora consideramos el parámetro como una variable aleatoria.

Adicionalmente, con la distribución conjunta en la ecuación (4.23), podemos definir:

Definition 4.12 (Distribución marginal). La distribución de X , obtenida mediante la desintegración de parámetro Θ del par (X, Θ) , es decir

$$p(x) = \int_{\Omega} p(x|\theta)p(\theta)d\theta \quad (4.23)$$

es conocida como distribución marginal de X .

Consideremos ahora que tenemos un conjunto de observaciones denotado por \mathcal{D} , de un modelo estadístico con parámetro Θ , entonces podemos definir

Definition 4.13 (Función de verosimilitud). La densidad de probabilidad evaluada en un conjunto de observaciones \mathcal{D} como función del valor del parámetro Θ , es decir

$$L : \Omega \rightarrow \mathbb{R} \quad (4.24)$$

$$\theta \mapsto l(\theta) = L_{\mathcal{D}}(\theta) = p(\mathcal{D}|\theta), \quad (4.25)$$

recibe el nombre de función de verosimilitud, o en inglés, *likelihood*.

Observación 4.7. La función de verosimilitud no es una densidad de probabilidad, es decir, no es cierto que

$$\int_{\Omega} L(\theta) d\theta = 1 \quad (4.26)$$

Observación 4.8. Dado que la función función de verosimilitud usualmente adquiere una forma exponencial (como por ejemplo en el caso de la familia exponencial), hay ocasiones en donde es conveniente usar la *log-verosimilitud*, esto es,

$$l(\theta) = \log L(\theta) = \log p(\mathcal{D}|\theta). \quad (4.27)$$

Esta formulación será particularmente útil cuando queramos optimizar la verosimilitud.

Observación 4.9. En general (pero no siempre) asumimos observaciones $\mathcal{D} = X_1, \dots, X_n$, $X_i \sim p(x|\theta)$, que son i.i.d. En cuyo caso, la verosimilitud factoriza de la forma $L_{\mathcal{D}}(\theta) = \prod_{i=1}^n L_{X_i}(\theta)$, con lo cual la log-verosimilitud toma la forma:

$$l_{\mathcal{D}}(\theta) = \sum_{i=1}^n l_{X_i}(\theta) \quad (4.28)$$

Ejemplo 4.9. Considere los datos $\mathcal{D} = \{x_1, \dots, x_n\}$, donde x_i es la observación de una VA $X_i \sim \mathcal{N}(\mu, \sigma^2)$ iid con σ^2 conocido. La función de verosimilitud de μ está dada por:

$$\begin{aligned} L(\mu) = p(\mathcal{D}|\mu) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-1}{2\sigma^2}(x_i - \mu)^2\right) \\ &= \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(\frac{-1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right). \end{aligned} \quad (4.29)$$

Luego, la log-verosimilitud está dada por:

$$l(\mu) = \log L(\mu) = -\frac{n}{2} \log(2\pi\sigma^2) + \frac{-1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2. \quad (4.30)$$

Ahora estamos en condiciones de definir el elemento central de la inferencia bayesiana, sobre el cual todo el proceso de inferencia toma lugar.

Definition 4.14 (Distribución posterior). Dado el conjunto de observación \mathcal{D} la distribución *posterior* del parámetro, es decir, considerando la información reportada por los datos \mathcal{D} , está dada por el teorema de Bayes mediante

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})} \propto p(\mathcal{D}|\theta)p(\theta) \quad (4.31)$$

donde:

- $p(\theta)$ es el prior del parámetro.
- $p(\theta|\mathcal{D})$ es la posterior del parámetro.
- $p(\mathcal{D}|\theta)$ es la verosimilitud
- $p(\mathcal{D}) = \int \Omega p(\mathcal{D}|\theta)p(\theta)d\theta$ es la densidad marginal de los datos

La *transición* de prior a posterior puede ser interpretada como el proceso de incorporar la evidencia de los datos (a través de la función de verosimilitud) para reducir la incertidumbre con respecto del valor del parámetro Θ . De la ecuación (4.32) podemos ver que este proceso, a veces referido como *actualización bayesiana*, equivale a multiplicar por la verosimilitud, para luego normalizar, garantizando que $p(\theta|\mathcal{D})$ es en efecto una densidad de probabilidad.

Observación 4.10. El símbolo \propto en la ecuación (4.32) es usado para indicar que el lado izquierdo es igual al lado derecho salvo una constante de proporcionalidad que depende de \mathcal{D} y no de θ . Con esto, cuando estemos calculando la posterior, solo nos enfocaremos en *una versión proporcional*, pues luego la densidad posterior se puede encontrar mediante la normalización de esta última.

Ejemplo 4.10 (Posterior modelo Bernoulli). Sea θ la probabilidad de obtener cara al lanzar una moneda, y sean X_1, \dots, X_n n resultados obtenidos al lanzar la moneda. Si no sabemos nada de θ antes del experimento, hace sentido tomar su prior como una distribución que de igual probabilidad a todo espacio de parámetros, es decir: $\theta \sim \text{Unif}(0, 1)$. Notemos que el prior encapsula la información que tenemos antes del experimento. Modelamos $X_1, \dots, X_n \sim \text{Bernoulli}(\theta)$. Entonces:

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

Notemos que en este caso, podemos calcular la distribución $p(X_1, \dots, X_n)$:

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

donde $B(x, y)$ es la función beta:

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

Sea $s = \sum_{i=1}^n X_i$. Entonces la distribución a posteriori será:

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

Usualmente, en experimentos reales, los datos x_1, \dots, x_n son recibidos de forma secuencial, es decir, *en línea*. De esta forma, es relevante notar que en primer lugar se observa x_1 primero, luego x_2 , y así sucesivamente.

Consecuentemente, si se asume el prior para el parámetro θ dado por $p(\theta)$, es posible hacer la actualización bayesiana *en línea* (o de forma adaptativa o continual), lo cual implica una corrección del modelo cada vez que se observan más datos.

Luego de observar x_1 , la posterior $p(\theta|x_1)$ puede ser calculada como:

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

Luego, al observar x_2 , usamos el hecho que X_1 y X_2 son condicionalmente independientes dado θ y obtenemos:

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

Con lo que para el caso general tenemos que

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

Observación 4.11. Cuando las observaciones \mathcal{D} son condicionalmente independientes dado el parámetro θ , entonces, la posterior $p(\theta|\mathcal{D})$ factoriza en las verosimilitudes de cada uno de los datos.

Observación 4.12. En la actualización bayesiana en línea, la posterior de la etapa n sirve de prior de la etapa $n + 1$.

4.10 Priors Conjugados

La actualización bayesiana puede resultar en una posterior solo conocida de forma proporcional (cuando no es posible calcular la distribución marginal $p(x)$) o bien en una distribución que no pertenece a una familia conocida. Una herramienta que asegure el cálculo de las distribuciones posteriores (incluyendo la constante de normalización) y que esta adopta una forma conocida es a través del uso de **priors conjugados**.

Definition 4.15. Sea un modelo con verosimilitud $p(x|\theta)$ y un prior sobre θ con densidad $p(\theta)$. Decimos que $p(\theta)$ es conjugado con la verosimilitud $p(x|\theta)$ si la posterior

$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)} \quad (4.32)$$

pertenece a *la misma familia* que el prior $p(\theta)$. Donde pertenecer a la misma familia quiere decir que ambas tienen una densidad de probabilidad definida por la misma forma funcional, e.g., $f_\lambda(\theta)$ pero con distintos valores para el *parámetro* λ , el cual es un *hiperparámetro* del modelo.

Ejemplo 4.11 (continuación de Ejemplo 4.10). Tarea: Verifique si el Ejemplo 4.10 es en efecto uno de prior conjugado.

Ejemplo 4.12 (Distribución Multinomial). Consideremos una variable aleatoria multinomial $X \sim \text{Mult}(n, \theta)$ donde θ pertenece al simplex

$$\{\theta \in [0, 1]^k : \theta_1 + \dots + \theta_k = 1\}. \quad (4.33)$$

La distribución multinomial genera vectores $X \in \mathbb{N}^k$ cuya i -ésima componente modela la cantidad de veces que ocurre el evento i dentro de k eventos en n intentos. Por ejemplo, si lanzamos un dado balanceado 100 veces, el vector que contiene el conteo de veces que obtenemos cada cara puede modelarse como

$$\theta_{\text{dado}} \sim \text{Mult}\left(100, \left[\frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}\right]\right). \quad (4.34)$$

Denotando $X = [x_1, \dots, x_n]$, observemos que una muestra multinomial $X \sim \text{Mult}(n, \theta)$ cumple con

$$\{x_i\}_{i=1}^k \subset \{0, 1, \dots, n\}, \quad \sum_{i=1}^k x_i = n. \quad (4.35)$$

Finalmente, la distribución Multinomial está dada por

$$\text{Mult}(X; n, \theta) = \frac{n!}{x_1! \dots x_k!} \theta_1^{x_1} \dots \theta_k^{x_k}, \quad (4.36)$$

y es la generalización de las distribuciones:

- Bernoulli cuando $k = 2$ y $n = 1$; pues $\text{Ber}X; \theta = \theta^x (1 - \theta)^{1-x}$
- Categórica (o *multinoulli*): cuando $n = 1$; pues $\text{Cat}(X; \theta) = \theta_1^{x_1} \dots \theta_k^{x_k}$
- Binomial: cuando $k = 2$; pues $\text{Bin}X; n, \theta = \binom{n}{x} \theta^x (1 - \theta)^{n-x}$

Observemos que el parámetro θ en la distribución multinomial (y las otras tres) es precisamente una distribución de probabilidad (discreta). Es decir, el construir un prior $p(\theta)$ implica definir una distribución sobre distribuciones discretas.

Definition 4.16 (Distribución de Dirichlet). Consideremos la distribución de Dirichlet

$$\theta \sim \text{Dir}(\theta|\alpha) = \frac{1}{B(\alpha)} \prod_{i=1}^k \theta_i^{\alpha_i-1}, \quad (4.37)$$

donde $\alpha = (\alpha_1, \dots, \alpha_k)$ es el parámetro de concentración y la constante de normalización está dada por $B(\alpha) = \prod_{i=1}^k \Gamma(\alpha_i) / \Gamma(\sum_{i=1}^k \alpha_i)$. El soporte de esta distribución es el simplex presentado en la ecuación (4.34).

En el caso $k = 3$, la distribución de Dirichlet puede ser graficada en el simplex de 2 dimensiones. La Figura 18 presenta tres gráficos para distintos valores del parámetro de concentración.

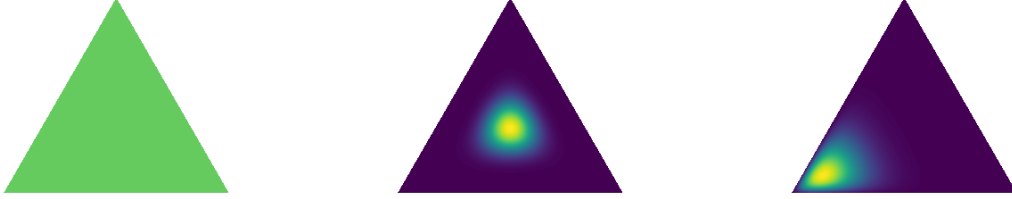


Fig. 18. Distribuciones Dirichlet para $k = 3$ con parámetros de concentración α (desde izquierda a derecha) dado por $[1, 1, 1]$, $[10, 10, 10]$ y $[10, 2, 2]$.

Veamos a continuación que la distribución de Dirichlet es conjugada al modelo Multinomial, y consecuentemente para Bernoulli, Categórica y Binomial. En efecto, si $\theta \sim \text{Dir}(\theta; \alpha)$ y $X \sim \text{Mult}(X; n, \theta)$, entonces

$$\begin{aligned} p(\theta|x) &= \frac{\text{Mult}(x; n, \theta) \text{Dir}(\theta; \alpha)}{p(x)} \\ &= \frac{n!}{x_1! \cdots x_k! p(x) B(\alpha)} \prod_{i=1}^k \theta_i^{x_i + \alpha_i - 1} \\ &= \frac{1}{B(\alpha')} \prod_{i=1}^k \theta_i^{\alpha'_i - 1} \end{aligned} \quad (4.38)$$

donde $\alpha' = (\alpha'_1, \dots, \alpha'_k) = (\alpha'_1 + x_1, \dots, \alpha'_k + x_k)$ es el nuevo parámetro de concentración.

Ejemplo 4.13. Consideremos $\alpha = [1, 2, 3, 4, 5]$ y generemos una muestra de $\theta \sim \text{Dir}(\theta|\alpha)$. El siguiente código genera, grafica e imprime esta muestra.

```
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}')
```

En nuestro caso, obtuvimos los parámetros $\theta = [0.034, 0.171, 0.286, 0.185, 0.324]$.

Ahora, usaremos un prior Dirichlet sobre θ con $\alpha_p = [1, 1, 1, 1, 1]$ para calcular la posterior de acuerdo a la ecuación (4.39). La Figura 19 muestra 50 muestras de la distribución posterior para distintas cantidades de observaciones entre 0 y 10^5 .

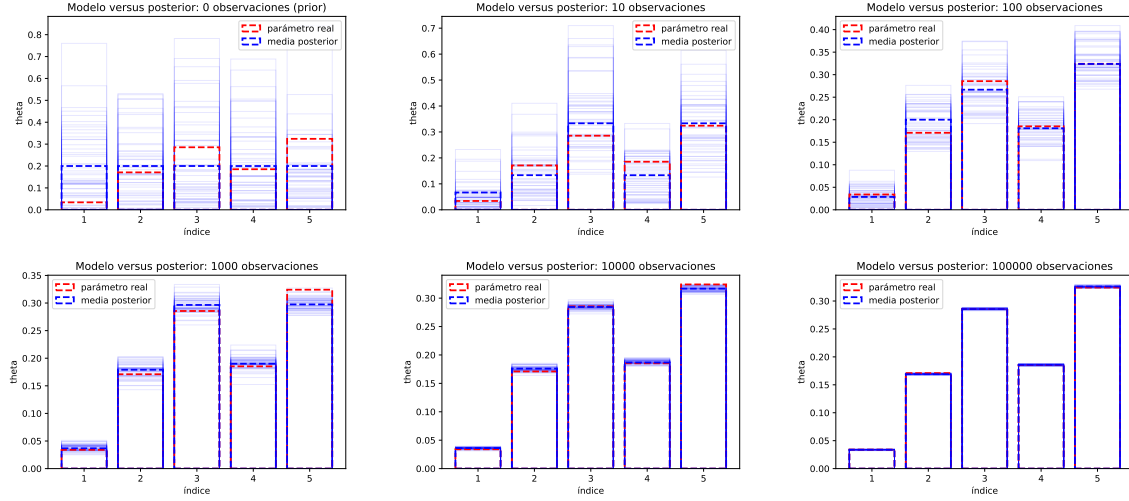


Fig. 19. Concentración de la distribución posterior en torno al parámetro real para un modelo $X \sim \text{Mult}(\theta)$ y una distribución a priori Dirichlet $\theta \sim \text{Dir}(\alpha)$. Se considera desde 0 hasta 10^5 observaciones y cada gráfico (desde izquierda-arriba hasta derecha-abajo) muestra el parámetro real (línea roja quebrada), la media posterior (línea azul quebrada) y 50 muestras de la posterior (azul claro). Observe cómo la distribución a priori (línea azul quebrada en la primera figura) pierde importancia a medida que el número de observaciones aumenta.

Ejemplo 4.14. Modelo gaussiano (σ^2 conocido). Consideremos el prior sobre la media $p(\mu) = \mathcal{N}(\mu_0, \sigma_0^2)$, con lo que la posterior está dada por

$$p(\mu|\mathcal{D}) \propto \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu)^2\right) \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp\left(-\frac{1}{2\sigma_0^2}(\mu - \mu_0)^2\right) \quad (4.39)$$

$$\propto \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 - \frac{1}{2\sigma_0^2}(\mu - \mu_0)^2\right), \quad (4.40)$$

donde la proporcionalidad viene de ignorar la constante $p(\mathcal{D})$ en la primera línea e ignorar todas las constantes que no dependen de μ en la segunda línea. Recordemos que estas constantes para μ incluyen a la varianza de x , σ^2 , por lo que ignorar esta cantidad es solo posible debido a que estamos considerando el caso en que σ^2 es conocido. Completando la forma cuadrática para μ dentro de la exponencial en la ec. (4.41), obtenemos

$$p(\mu|\mathcal{D}) \propto \exp\left(-\frac{1}{2\sigma_n^2}(\mu - \mu_n)^2\right), \quad (4.41)$$

donde (ya definiremos μ_n y σ_n^2 en breve) como $p(\mu|\mathcal{D})$ debe integrar uno, la única densidad de probabilidad proporcional al lado derecho de la ecuación anterior es la Gaussiana de media μ_n y varianza σ_n^2 . Es decir, la constante de proporcionalidad necesaria para la igualdad en la expresión anterior es

$$\int_{\mathbb{R}} \exp\left(-\frac{1}{2\sigma_n^2}(\mu - \mu_n)^2\right) d\mu = (2\pi\sigma_n^2)^{n/2}. \quad (4.42)$$

Consecuentemente, confirmamos que el prior elegido era efectivamente conjugado con la verosimilitud gaussiana, con lo que la posterior está dada por la siguiente densidad (gaussiana):

$$p(\mu|\mathcal{D}) = \mathcal{N}(\mu; \mu_n, \sigma_n^2) = \frac{1}{(2\pi\sigma_n^2)^{N/2}} \exp\left(-\frac{1}{2\sigma_n^2}(\mu - \mu_n)^2\right), \quad (4.43)$$

donde la media y la varianza están dadas respectivamente por

$$\mu_n = \frac{1}{\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}} \left(\frac{1}{\sigma_0^2} \mu_0 + \frac{n}{\sigma^2} \bar{x} \right), \quad \text{donde } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (4.44)$$

$$\sigma_n = \left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right)^{-1}. \quad (4.45)$$

Observación 4.13. La actualización bayesiana transforma los parámetros del prior de μ desde μ_0 y σ_0^2 hacia μ_n y σ_n^2 en las ecs. (4.45) y (4.46) respectivamente. Notemos que los parámetros de la posterior son combinaciones (interpretables por lo demás) entre los parámetros del prior y los datos, en efecto, la μ_n es el promedio ponderado entre μ_0 (que es nuestro candidato para μ antes de ver datos) con factor σ_0^{-2} y el promedio de los datos \bar{x} con factor $(\sigma^2/n)^{-1}$, que a su vez es el estimador de máxima verosimilitud. Es importante también notar que estos factores son las varianzas inversas—i.e., precisión—de μ_0 y de \bar{x} . Finalmente, observemos que σ_n es la *suma paralela* de las varianzas, pues si expresamos la ec. (4.46) en términos de *precisiones*, vemos que la precisión inicial σ_0^2 aumenta un término σ^2 con cada dato que vemos; lo cual tiene sentido pues con más información es la precisión la que debe aumentar y no la incertidumbre (en este caso representada por la varianza).

Ejemplo 4.15. Modelo gaussiano (μ conocido). Ahora procedemos con el siguiente prior para la varianza, llamado Gamma-inverso:

$$p(\sigma^2) = \text{inv-}\Gamma(\sigma^2; \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)(\sigma^2)^{\alpha+1}} \exp(-\beta/\sigma^2) \quad (4.46)$$

esta densidad recibe dicho nombre pues es equivalente a modelar la precisión, definida como el recíproco de la varianza $1/\sigma^2$, mediante la distribución Gamma. Los hiperparámetros α y β son conocidos como parámetros de forma y de tasa (o precisión) respectivamente.

Con este prior, la posterior de la varianza toma la forma:

$$\begin{aligned} p(\sigma^2 | \mathcal{D}) &\propto \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu)^2\right) \frac{\beta^\alpha}{\Gamma(\alpha)(\sigma^2)^{\alpha+1}} \exp(-\beta/\sigma^2) \\ &\propto \frac{1}{(\sigma^2)^{N/2+\alpha+1}} \exp\left(-\frac{1}{\sigma^2} \left(\frac{1}{2} \sum_{i=1}^N (x_i - \mu)^2 + \beta \right)\right) \end{aligned} \quad (4.47)$$

donde nuevamente la proporcionalidad ha sido mantenida debido a la remoción de las constantes. Esta última expresión es proporcional a una distribución Gamma inversa con hiperparámetros α y β ajustados en base a los datos observados.

Hay ocasiones en las que el conocimiento a priori sobre el parámetro no puede ser convenientemente expresado mediante una densidad de probabilidad pero sí una densidad que no necesariamente integra uno o incluso es (Lebesgue) integrable. Para reflejar esta idea, se usan priors impropios.

Definition 4.17 (Prior impropia). Una distribución a priori impropia es una distribución que no es necesariamente de probabilidad (i.e., no integra 1), pero que de todas formas puede ser utilizada como distribución a priori en el contexto de inferencia bayesiana, pues la distribución posterior correspondiente si es una distribución de probabilidad apropiada.

Observación 4.14. No es necesario usar la constante de normalización en las densidades a priori Gaussianas (o ninguna otra en realidad).

Observación 4.15. Veamos que un prior impropio puede incluso tener integral infinita, en el caso de la distribución normal $X \sim \mathcal{N}(X; \mu, 1)$, $\mu \in \mathbb{R}$, podemos elegir $p(\mu) \propto 1$ y escribir

$$p(\mu|x) \propto p(x|\mu) \cdot 1 = \mathcal{N}(x; \mu, 1) = \mathcal{N}(\mu; x, 1). \quad (4.48)$$

Considerar distribuciones uniformes impropias como priors no informativas parece tener sentido, pues intuitivamente no estamos dando preferencia (mayor probabilidad a priori) a ningún valor del parámetro por sobre otro. Sin embargo, este procedimiento sufre de una desventaja conceptual.

4.11 Estimadores bayesianos

Si bien ya hemos estudiado el rol del prior en la inferencia bayesiana, hasta ahora no lo hemos considerado en la construcción de estimadores. En particular, el EMV no incorpora conocimiento a priori del parámetro. Con el objetivo de incorporar este conocimiento a priori en el cálculo de estimadores puntuales, consideramos que en el caso general, podemos considerar otros estimadores puntuales a través de una función de pérdida asociada a estimar el parámetro θ mediante el estimador $\hat{\theta}$ dada por $L(\theta, \hat{\theta})$. Con esto podemos definir los conceptos de riesgo y estimador bayesiano.

Definition 4.18 (Riesgo bayesiano). Para una función de pérdida $L(\theta, \hat{\theta})$ y un conjunto de observaciones \mathcal{D} , el riesgo bayesiano es la esperanza posterior de dicha función de pérdida, es decir

$$R(\hat{\theta}) = \int_{\Omega} L(\theta, \hat{\theta}) p(\theta | \mathcal{D}) d\theta. \quad (4.49)$$

Definition 4.19 (Estimador bayesiano). Dado un conjunto de datos \mathcal{D} y un riesgo bayesiano $R(\theta)$, un estimador bayesiano es uno que minimiza el riesgo bayesiano:

$$\theta_{\text{Bayes}}(\mathcal{D}) = \arg \min_{\Omega} R(\theta). \quad (4.50)$$

donde es implícito que $R(\cdot)$ se define con \mathcal{D} .

A continuación, se definirán estimadores bayesianos con distintas funciones de costo o de riesgo.

Definition 4.20 (Bayes' Least-Squares (BLS)). El caso estándar es la función de pérdida cuadrática $L_2(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$ la cual resulta en el estimador dado por la media posterior $\theta_{\text{Bayes}}(\mathcal{D}) = \mathbb{E} \theta | \mathcal{D}$

Definition 4.21 (Minimum absolute-error (MAE)). De forma similar, la función de costo $L_1(\theta, \hat{\theta}) = |\theta - \hat{\theta}|_1$ resulta en el estimador dado por la mediana posterior.

Encontrar una función de pérdida para el máximo a posteriori es menos directo. Consideremos en primer lugar el caso $\theta \in \Omega$ discreto y la pérdida “0-1”

$$L_{0-1}(\theta, \hat{\theta}) = \begin{cases} 0 & \text{si } \theta = \hat{\theta}, \\ 1 & \text{si no.} \end{cases}$$

El riesgo de Bayes asociado a $L_{0-1}(\theta, \hat{\theta})$ (en el caso discreto) toma la forma

$$R(\hat{\theta}) = \mathbb{P}(\theta \neq \hat{\theta} | \mathcal{D}) = 1 - \mathbb{P}(\theta = \hat{\theta} | \mathcal{D}), \quad (4.51)$$

lo cual es minimizado eligiendo $\hat{\theta}$ tal que $\mathbb{P}(\theta = \hat{\theta} | \mathcal{D})$ es máximo, es decir, el MAP. ¿por qué no es posible proceder de esta forma para el caso continuo? ¿cuál es la función de costo asociada al MAP en el caso continuo?

Definition 4.22 (Estimador máximo a posteriori). Sea $\theta \in \Theta$ un parámetro con distribución a posteriori $p(\theta | D)$ definida en todo Θ . Entonces nos referiremos a su estimación puntual dada por:

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

como el estimador *máximo a posteriori* (MAP). Se utiliza la siguiente función de costo:

$$C(a, b) = \begin{cases} 1, & |a - b| > 0 \\ 0, & \sim \end{cases}$$

Observación 4.16. Es posible encontrar el MAP solo teniendo acceso a una versión *proporcional* a la distribución posterior, un escenario usual en inferencia bayesiana, o también mediante la maximización del logaritmo de ésta última. En efecto,

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

donde hemos encontrado la maximización de la función de log-verosimilitud, pero ahora junto al log-prior.

Observación 4.17. Es relevante notar que el estimador MAP es una *modificación* del EMV, pues ambos comparten una parte de la misma función objetivo (verosimilitud) con la diferencia que el MAP además incluye el término *log-prior*. Esto puede entenderse como una regularización de la solución del problema de MV, en donde el término adicional puede representar las propiedades del estimador más allá de que las pueden ser exclusivamente revelada por los datos.

Ejemplo 4.16 (Máximo a posterior para el modelo gaussiano). En particular, para el modelo lineal y gaussiano que hemos considerado hasta ahora, podemos calcular θ_{MAP} para un prior Gaussiano de media cero y varianza σ_θ^2 . Éste está dado por (asumimos la varianza del ruido σ_ϵ^2 conocida):

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

Podemos ver que eligiendo un prior uniforme o de normal de varianza muy amplia, el MAP es equivalente al EMV. ¿qué significa esto? ¿qué comportamiento diferente de EMV promueve el MAP en este caso?

4.12 Posterior predictiva

En la inferencia bayesiana las predicciones ocupan un rol relevante, pues luego de realizar inferencia sobre un modelo estadístico, en general estamos interesados estudiar cómo serán los siguientes datos generados por el modelo. Para esto definiremos la predicción bayesiana de la forma

Definition 4.23 (Posterior predictiva). Para un conjunto de datos \mathcal{D} y un parámetro θ , la densidad posterior predictiva está dada por

$$p(x|\mathcal{D}) = \int_{\Omega} p(x|\theta)p(\theta|\mathcal{D})d\theta = \mathbb{E} p(x|\theta)|\mathcal{D}, \quad (4.52)$$

es decir, el valor esperado del modelo estadístico con respecto a la ley posterior del parámetro (modelo).

Podemos ahora considerar la posterior predictiva como nuestro modelo *aprendido* y generar datos de él, donde nos encontramos frente al mismo dilema de un estimador puntual como en el caso anterior: es posible considerar muestras aleatorias, la media, la mediana o algún intervalo.

Observación 4.18. La posterior predictiva es distinta (en general) a la predicción *plug-in*, en donde consideramos en modelo estadístico $p_{\hat{\theta}}$ en base a un estimador (puntual) cualquiera $\hat{\theta}$. Desde esa perspectiva, la posterior predictiva equivale a considerar estimadores y modelos puntuales pero integrar todos ellos con respecto a la ley posterior.

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