

Exercise 1: Maximum Likelihood vs. Bayes

An unfair coin is tossed seven times and the event (head or tail) is recorded at each iteration. The observed sequence of events is

D = (x1, x2, ..., x7) = (head, head, tail, tail, head, head, head).

We assume that all tosses x1, x2, ... have been generated independently following the Bernoulli probability distribution

P(x | θ) = { θ if x = head, 1 - θ if x = tail,

where θ ∈ [0, 1] is an unknown parameter.

(a) State the likelihood function P(D|θ), that depends on the parameter θ.

P(D|θ) = Π_{i=1}^7 P(x_i | θ) = θ · θ · (1-θ) · (1-θ) · θ · θ · θ = θ^{#heads} · (1-θ)^{#tails} = θ^5 · (1-θ)^2

(b) Compute the maximum likelihood solution θ̂, and evaluate for this parameter the probability that the next two tosses are "head", that is, evaluate

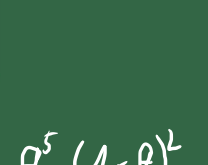
P(x8 = head, x9 = head | θ̂).

log P(D|θ) = 5 · log θ + 2 · log(1-θ) (concave) ∂/∂θ log P(D|θ) = 5/θ - 2/(1-θ) = 0 => θ̂ = 5/7

P(x8=head, x9=head) = θ̂ · θ̂ = 25/49

(c) We now adopt a Bayesian view on this problem, where we assume a prior distribution for the parameter θ defined as:

p(θ) = { 1 if 0 ≤ θ ≤ 1, 0 else. ∫ p(θ) dθ = 1, p(θ) ≥ 0



Compute the posterior distribution p(θ|D), and evaluate the probability that the next two tosses are head, that is,

∫ P(x8 = head, x9 = head | θ) p(θ|D) dθ.

p(θ|D) = P(D|θ)p(θ) / ∫ P(D|θ)p(θ)dθ = (θ^5 · (1-θ)^2) / ∫_0^1 θ^5 (1-θ)^2 dθ = 168 · θ^5 · (1-θ)^2 / 168 = θ^5 · (1-θ)^2 ∫ P(x8=head, x9=head | θ) · 168 · θ^5 · (1-θ)^2 dθ = ∫_0^1 θ^2 · 168 · θ^5 (1-θ)^2 dθ = 7/15

Exercise 2: Principal Component Analysis

We consider an unsupervised dataset x1, ..., xN ∈ ℝ^d, where x̄ = 1/N ∑_{k=1}^N x_k is the empirical mean. The principal component analysis problem consists of finding the vector e ∈ ℝ^d of norm 1 such that the data projected in this space has maximum variance, i.e. is a solution of the optimization problem

max_{e ∈ ℝ^d} 1/N ∑_{k=1}^N (e^T x_k - m)^2 subject to ||e||^2 = 1

where m = 1/N ∑_{k=1}^N e^T x_k is the mean of the projected data.

(a) Show that the problem can be rewritten as the quadratic program

max_{e ∈ ℝ^d} e^T C e subject to ||e||^2 = 1

where C = 1/N ∑_{k=1}^N (x_k - x̄) · (x_k - x̄)^T is the empirical covariance matrix.

1/N ∑ e^T x_k = e^T (1/N ∑ x_k) = e^T x̄

1/N ∑ (e^T x_k - 1/N ∑ e^T x_k)^2 = 1/N ∑ e^T x_k x_k^T e - 2 e^T x̄ x̄^T e + e^T x̄ x̄^T e = e^T (1/N ∑ (x_k - x̄)(x_k - x̄)^T) e = e^T C e

(b) Show using the method of Lagrange multipliers that the solution of the optimization problem above is an eigenvector of the matrix C.

L(e, λ) = e^T C e - λ (||e||^2 - 1) ∂L/∂e = 2Ce - 2λe = 0 => Ce = λe

(c) Show that, among all possible eigenvectors of C, the solution of the optimization problem above is the one with highest associated eigenvalue.

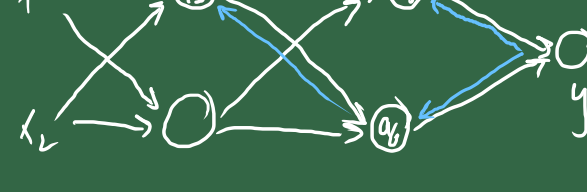
e^T C e = λ · e^T e = λ ||e||^2 = λ

Exercise 3: Neural Networks

We consider a neural network that takes two inputs x1 and x2 and produces an output y based on the following set of computations:

z3 = x1 · w13 + x2 · w23, z5 = z3 · w35 + a4 · w45, y = z5 + z6, a3 = tanh(z3), a5 = tanh(z5), z4 = x1 · w14 + x2 · w24, z6 = z3 · w36 + a4 · w46, a4 = tanh(z4), a6 = tanh(z6)

(a) Draw the neural network graph associated to this set of computations.



(b) Write the set of backward computations that leads to the evaluation of the partial derivative ∂y/∂w13. Your answer should avoid redundant computations. Hint: tanh'(t) = 1 - (tanh(t))^2.

δ5 = 1 - a5^2 = ∂y/∂z5, δ6 = 1 - a6^2, δ3 = (w35 · δ5 + w36 · δ6) · (1 - a3^2) = ∂y/∂z3, ∂y/∂w13 = x1 · δ3

Exercise 4: Support Vector Machines

The primal program for the linear hard margin SVM is

min_{w, θ} ||w||^2, subject to y_i (w^T x_i + θ) ≥ 1, for 1 ≤ i ≤ N,

where ||·|| denotes the Euclidean norm, and the minimization is performed in w ∈ ℝ^d, θ ∈ ℝ, while the data xi ∈ ℝ^d, yi ∈ {-1, 1} are regarded as fixed constants.

(a) State the Lagrangian dual of the constrained optimization problem above and determine when the Slater's conditions for strong duality are satisfied.

L(w, θ, α1, α2, ..., αN) = 1/2 ||w||^2 + ∑ α_i (1 - y_i (w^T x_i + θ)) max_{α, w, θ} min_{w, θ} L(w, θ, α) s.t. α_i ≥ 0, ∃ (w, θ) s.t. y_i (w^T x_i + θ) > 1 => strong duality



(b) Show that the Lagrange dual takes the form of a quadratic optimization problem w.r.t. the dual variables α1, ..., αN.

∂L/∂w = w - ∑ α_i y_i x_i = 0 => w = ∑ α_i y_i x_i, ∂L/∂θ = -∑ α_i y_i = 0 => ∑ α_i y_i = 0 (1) into (2): L(θ, α) = 1/2 [∑ ∑ α_i α_j y_i y_j x_i^T x_j] + ∑ α_i - [∑ α_i y_i ∑ y_j x_j^T x_i] => L(α) = ∑ α_i - 1/2 ∑ ∑ α_i α_j y_i y_j x_i^T x_j s.t. ∑ α_i y_i = 0

Exercise 5: Kernels

A kernel function k: ℝ^d × ℝ^d → ℝ generalizes the linear scalar product between two vectors. The kernel must satisfy positive semi-definiteness, that is, for any sequence of data points x1, ..., xn ∈ ℝ^d and coefficients c1, ..., cn ∈ ℝ the following inequality should hold:

∑_{i=1}^n ∑_{j=1}^n c_i c_j k(x_i, x_j) ≥ 0

We consider the kernel function k(x, x') = <x, x'>^2.



(a) Show that this kernel is positive semi-definite.

∑ ∑ c_i c_j <x_i, x_j>^2 = ∑ ∑ c_i c_j (∑_{k=1}^d x_{ik} · x_{jk})^2 = ∑ ∑ c_i c_j (∑_k x_{ik} x_{jk}) (∑_l x_{il} x_{jl}) = ∑ ∑ c_i c_j x_{ik} x_{il} · c_j x_{jk} x_{jl} = ∑_k ∑_l (∑_i c_i x_{ik} x_{il}) (∑_j c_j x_{jk} x_{jl}) = ∑_k ∑_l (∑_i c_i x_{ik} x_{il})^2 ≥ 0

(b) Show that this kernel can be rewritten as a dot product k(x, x') = <φ(x), φ(x')>.