

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

$$\mathcal{D} = (x_1, x_2, \dots, x_7) = (\text{head}, \text{head}, \text{tail}, \text{tail}, \text{head}, \text{head}, \text{head}).$$

We assume that all tosses x_1, x_2, \dots have been generated independently following the Bernoulli probability distribution

$$P(x | \theta) = \begin{cases} \theta & \text{if } x = \text{head} \\ 1 - \theta & \text{if } x = \text{tail}, \end{cases}$$

where $\theta \in [0, 1]$ is an unknown parameter.

(a) State the likelihood function $P(\mathcal{D}|\theta)$, that depends on the parameter θ .

$$P(\mathcal{D}|\theta) = \prod_{i=1}^7 p(x_i|\theta) = \theta^5 \cdot (1-\theta)^2$$

(b) Compute the maximum likelihood solution $\hat{\theta}$, and evaluate for this parameter the probability that the next two tosses are "head", that is, evaluate

$$P(x_8 = \text{head}, x_9 = \text{head} | \hat{\theta}).$$

$$P(\mathcal{D}|\theta) = \theta^5 (1-\theta)^2$$

$$\log P(\mathcal{D}|\theta) = 5 \cdot \log \theta + 2 \cdot \log(1-\theta) \quad \text{concave}$$

$$\frac{d}{d\theta} \log P(\mathcal{D}|\theta) = \frac{5}{\theta} - \frac{2}{1-\theta} \stackrel{!}{=} 0 \Rightarrow \hat{\theta} = \frac{5}{7}$$

$$P(h, h | \hat{\theta}) = \hat{\theta} \hat{\theta} = \frac{25}{49}$$

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

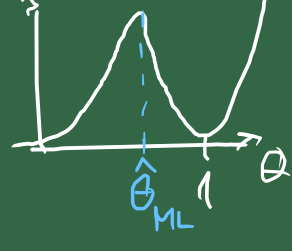
$$p(\theta) = \begin{cases} 1 & \text{if } 0 \leq \theta \leq 1 \\ 0 & \text{else.} \end{cases}$$

Compute the posterior distribution $p(\theta|\mathcal{D})$, and evaluate the probability that the next two tosses are head, that is,

$$\int P(x_8 = \text{head}, x_9 = \text{head} | \theta) p(\theta|\mathcal{D}) d\theta.$$

$$p(\theta|\mathcal{D}) = \frac{P(\mathcal{D}|\theta) p(\theta)}{\int_0^1 P(\mathcal{D}|\theta) p(\theta) d\theta} = \frac{\theta^5 (1-\theta)^2}{\int_0^1 \theta^5 (1-\theta)^2 d\theta} \quad (0 \leq \theta \leq 1)$$

$$= 168 \cdot \theta^5 (1-\theta)^2$$



$$\int_0^1 P(h, h | \theta) p(\theta|\mathcal{D}) d\theta$$

$$= \int_0^1 \theta^2 \cdot 168 \cdot \theta^5 (1-\theta)^2 d\theta$$

$$= \frac{7}{15} \left[\int_0^1 \theta^5 - 2\theta^6 + \theta^7 d\theta = \left[\frac{\theta^6}{6} - 2 \cdot \frac{\theta^7}{7} + \frac{\theta^8}{8} \right]_0^1 \right]$$

$$= \frac{7}{15} \left[\frac{1}{6} - \frac{2}{7} + \frac{1}{8} = \frac{56}{336} - \frac{96}{336} + \frac{42}{336} = \frac{2}{336} = \frac{1}{168} \right]$$

Exercise 2: Principal Component Analysis

We consider an unsupervised dataset $x_1, \dots, x_N \in \mathbb{R}^d$, where $\bar{x} = \frac{1}{N} \sum_{k=1}^N x_k$ is the empirical mean. The principal component analysis problem consists of finding the vector $e \in \mathbb{R}^d$ of norm 1 such that the data projected in this space has maximum variance, i.e. is a solution of the optimization problem

$$\max_{e \in \mathbb{R}^d} \frac{1}{N} \sum_{k=1}^N (e^T x_k - m)^2 \quad \text{subject to } \|e\|^2 = 1$$

where $m = \frac{1}{N} \sum_{k=1}^N e^T x_k$ is the mean of the projected data.

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

$$\max_{e \in \mathbb{R}^d} e^T C e \quad \text{subject to } \|e\|^2 = 1$$

where $C = \frac{1}{N} \sum_{k=1}^N (x_k - \bar{x}) \cdot (x_k - \bar{x})^T$ is the empirical covariance matrix.

$$\begin{aligned} \frac{1}{N} \sum_k \left[e^T x_k - \frac{1}{N} \sum_k e^T x_k \right]^2 &= \frac{1}{N} \sum_k e^T x_k x_k^T e - 2 \frac{1}{N} \sum_k e^T x_k \bar{x}^T e + e^T \bar{x} \bar{x}^T e \\ &= \frac{1}{N} \sum_k e^T (x_k x_k^T - 2 \bar{x} \bar{x}^T + \bar{x} \bar{x}^T) e \\ &= e^T \left(\frac{1}{N} \sum_k (x_k - \bar{x})(x_k - \bar{x})^T \right) e \\ &= e^T C e \end{aligned}$$

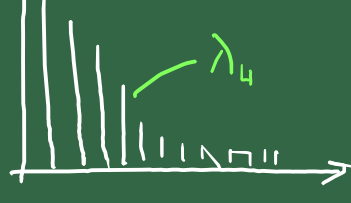
(b) Show using the method of Lagrange multipliers that the solution of the optimization problem above is an eigenvector of the matrix C .

$$\mathcal{L}(e, \lambda) = e^T C e - \lambda (\|e\|^2 - 1)$$

$$\frac{\partial \mathcal{L}}{\partial e} = 2Ce - 2\lambda e \stackrel{!}{=} 0 \Leftrightarrow \underline{Ce = \lambda e}$$

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

$$e^T C e = e^T (\lambda e) = \lambda \cdot \|e\|^2 = \lambda$$



Exercise 3: Neural Networks

We consider a neural network that takes two inputs x_1 and x_2 and produces an output y based on the following set of computations:

$$z_3 = x_1 \cdot w_{13} + x_2 \cdot w_{23}$$

$$a_3 = \tanh(z_3)$$

$$z_4 = x_1 \cdot w_{14} + x_2 \cdot w_{24}$$

$$a_4 = \tanh(z_4)$$

$$z_5 = a_3 \cdot w_{35} + a_4 \cdot w_{45}$$

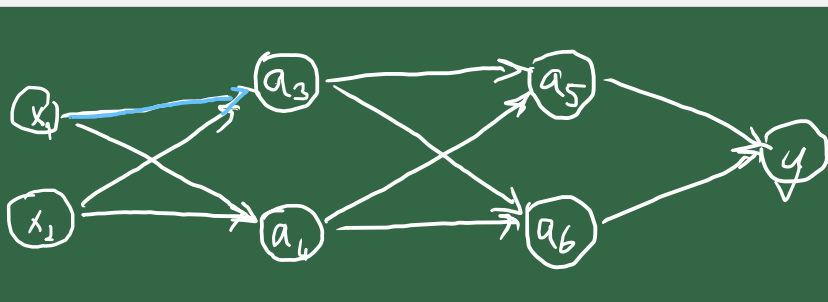
$$a_5 = \tanh(z_5)$$

$$z_6 = a_3 \cdot w_{36} + a_4 \cdot w_{46}$$

$$a_6 = \tanh(z_6)$$

$$y = a_5 + a_6$$

(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 $\partial y / \partial w_{13}$. Your answer should avoid redundant computations. Hint: $\tanh'(t) = 1 - (\tanh(t))^2$.

$$\delta_j = \frac{\partial y}{\partial z_j}$$

$$\delta_6 = 1 - a_6^2$$

$$\delta_5 = 1 - a_5^2$$

$$\delta_3 = (\delta_5 \cdot w_{35} + \delta_6 \cdot w_{36}) \cdot (1 - a_3^2)$$

$$\frac{\partial y}{\partial w_{13}} = x_1 \cdot \delta_3$$

Regarding Slater's condition, since the constraints are linear and the two hyperplanes typically do not coincide, there exists a strictly feasible point satisfying $w_1^T x > 1, w_2^T x < -1$. Thus, Slater's condition holds, ensuring strong duality, and we can solve the problem using Lagrangian duality.

then strong duality holds

Slater: linear separability of two hyperplanes

Exercise 4: Support Vector Machines

The primal problem for the linear hard margin SVM is

$$\min_{w, \theta} \|w\|^2 \quad \text{subject to } y_i (w^T x_i + \theta) \geq 1, \quad \text{for } 1 \leq i \leq N,$$

where $\|\cdot\|$ denotes the Euclidean norm, and the minimization is performed in $w \in \mathbb{R}^d, \theta \in \mathbb{R}$, while the data $x_i \in \mathbb{R}^d, y_i \in \{-1, 1\}$ are regarded as fixed constants.

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

$$\mathcal{L}(w, \theta, \alpha) = \frac{1}{2} \|w\|^2 - \sum_i \alpha_i [y_i (w^T x_i + \theta) - 1]$$

$$\max_{\alpha} \min_{w, \theta} \mathcal{L}(w, \theta, \alpha) \quad \text{s.t. } \forall_i \alpha_i \geq 0$$

Slater: linear separability of the two classes

(b) Show that the Lagrange dual takes the form of a quadratic optimization problem w.r.t. the dual variables $\alpha_1, \dots, \alpha_N$.

$$\frac{\partial \mathcal{L}}{\partial w} = w - \sum_i \alpha_i y_i x_i \stackrel{!}{=} 0 \Rightarrow w = \sum_i \alpha_i y_i x_i$$

$$\frac{\partial \mathcal{L}}{\partial \theta} = -\sum_i \alpha_i y_i \stackrel{!}{=} 0 \Rightarrow \sum_i \alpha_i y_i = 0$$

$$\mathcal{L}(\tilde{\alpha}) = -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum_i \alpha_i$$

$$\max_{\alpha} -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum_i \alpha_i$$

$$\text{s.t. } \forall_i \alpha_i \geq 0$$

$$\sum_i \alpha_i y_i = 0$$

The dual formulation is preferred when the number of constraints is much smaller than the number of primal variables (i.e., when d is large). This is because solving the dual problem involves optimizing over a lower-dimensional space, making it computationally more efficient. However, if d is small or the number of constraints is large, solving the primal problem directly may be more straightforward.

当约束的数量远小于原始变量的数量时 (即 when is large), 则首选对偶公式。这是因为解决对偶问题涉及在较低维空间上进行优化, 从而提高计算效率。但是, 如果很小或者 constraints 的数量很多, 直接解决 the primal problem 可能更直接。

Exercise 5: Kernels

A kernel function $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ generalizes the linear scalar product between two vectors. The kernel must satisfy positive semi-definiteness, that is, for any sequence of data points $x_1, \dots, x_n \in \mathbb{R}^d$ and coefficients $c_1, \dots, c_n \in \mathbb{R}$ the following inequality should hold:

$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) \geq 0$$

We consider the kernel function $k(x, x') = \langle x, x' \rangle$.

$$K_{ij} = k(x_i, x_j)$$

$$C^T K C \geq 0$$

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

$$\sum_i \sum_j c_i c_j \left(\sum_k x_{ik} \cdot x_{jk} \right)^2$$

$$= \sum_i \sum_j c_i c_j \left(\sum_k x_{ik} x_{jk} \right) \left(\sum_l x_{il} x_{jl} \right)$$

$$= \sum_i \sum_j c_i c_j \sum_k \sum_l x_{ik} x_{jk} x_{il} x_{jl}$$

$$= \sum_k \sum_l \sum_i \sum_j c_i x_{ik} x_{il} \cdot c_j x_{jk} x_{jl}$$

$$= \sum_k \sum_l \left(\sum_i c_i x_{ik} x_{il} \right)^2 \geq 0$$