

MALTA Exercises

Giovanni Neglia and Samir M. Perlaza

October 9, 2024

Notation

Let $\alpha \in (0, 1)$ and $\beta \in [0, 1]$ be two fixed reals. The binary **relative entropy** (or Kullback-Leibler divergence) of the probability distribution represented by the vector $(\alpha, 1 - \alpha)$ with respect to the probability distribution represented by the vector $(\beta, 1 - \beta)$, denoted by $D(\alpha \parallel \beta)$, is

$$D(\alpha \parallel \beta) = \beta \log \left(\frac{\beta}{\alpha} \right) + (1 - \beta) \log \left(\frac{1 - \beta}{1 - \alpha} \right). \quad (1)$$

The binary **entropy** of the probability distribution represented by the vector $(\beta, 1 - \beta)$, denoted by $H(\beta)$, is

$$H(\beta) = -\beta \log(\beta) - (1 - \beta) \log(1 - \beta), \quad (2)$$

where $0 \log(0)$ is assumed to be zero in both (1) and (2).

Exercises Week 5

Ex. 1 — Consider a dataset

$$S = ((\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)) \in \left(\mathbb{R}^d \times \{-1, 1\} \right)^m,$$

and a function $h_{\mathbf{w}} : \mathbb{R}^d \rightarrow (0, 1)$ in the class of functions \mathcal{H}_{sig} of the form

$$h_{\mathbf{w}}(\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})}.$$

In class, it was shown that if for all $t \in \{1, 2, \dots, m\}$, the value $h_{\mathbf{w}}(\mathbf{x}_t)$ is interpreted as the posterior probability

$$\text{Prob}(Y = 1 | \mathbf{X} = \mathbf{x}_t),$$

then the likelihood of the dataset S , denoted by $P_{\mathbf{w}}(S)$, is

$$P_{\mathbf{w}}(S) = \left(\prod_{t \in \{i \in [m] : y_i = 1\}} h_{\mathbf{w}}(\mathbf{x}_t) \right) \left(\prod_{t \in \{i \in [m] : y_i = -1\}} 1 - h_{\mathbf{w}}(\mathbf{x}_t) \right), \quad (3)$$

which leads to the equality

$$-\frac{1}{m} \log(P_{\mathbf{w}}(S)) = \frac{1}{m} \sum_{t=1}^m \log(1 + \exp(-y_t \mathbf{w}^\top \mathbf{x}_t)). \quad (4)$$

1. Using the equality in (4), explain the connection between empirical risk minimization and maximum likelihood estimation in the context of logistic regressions.
2. The empirical risk minimization was claimed to be a convex optimization problem. Provide a proof.

3. The names of the labels do not influence the learning process. Assume that the dataset is transformed (by re-labeling) into

$$S = ((\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)) \in \left(\mathbb{R}^d \times \{0, 1\} \right)^m.$$

Show that in such a case, the following holds:

$$\log(P_{\mathbf{w}}(S)) = \sum_{t=1}^m y_t \log(h_{\mathbf{w}}(\mathbf{x}_t)) + (1 - y_t) \log(1 - h_{\mathbf{w}}(\mathbf{x}_t)). \quad (5)$$

Provide an interpretation on the right-hand side of (5) in terms of empirical risk minimization: What is the loss function?

4. The equality in (5) can be written in terms of information measures as follows:

$$-\log(P_{\mathbf{w}}(S)) = \sum_{t=1}^m D(y_t \| h_{\mathbf{w}}(\mathbf{x}_t)) + H(y_t). \quad (6)$$

Prove the equality in (6) and give an interpretation on the right-hand side of (6) in terms of empirical risk minimization: What is the loss function? **Hint:** Note that for all $t \in \{1, 2, \dots, m\}$, $H(y_t) = 0$.