

1 Information and learning

Here we will try to get a closer understanding of the relationship between information theoretic measures we have discussed in class and the objectives in machine learning.

Problem 1 [17 points]

Let $X \in [m]$ be a random variable, taking one of m according to some distribution $P \in \mathbb{S}^m$. I.e., P is a point on the m -simplex, in other words, P is an m -dimensional vectors whose m elements are all between 0 and 1, and sum to 1.

Show that the P maximizing the entropy $H(X)$ is the uniform distribution over $[m]$, i.e., $P(i) = 1/m$ for all $i \in [m]$.

Advice: There are at least two ways to do this problem. One is to set up an optimization problem, use Lagrange multipliers, and solve. Another is to let U be the uniform distribution over $[m]$, calculate $D_{KL}(U \| P)$ for any arbitrary distribution P , and use this to conclude the desired result. The second method is cleaner.

End of problem 1

• Let P be any arbitrary distribution $[p_1, \dots, p_m]$.

Let U be the uniform distribution $[1/m, \dots, 1/m]$.

$$\begin{aligned} D_{KL}(P \| U) &= -\sum_{i=1}^m p_i \log \frac{1}{m} + \sum_{i=1}^m p_i \log p_i \\ &= -\log \frac{1}{m} + \sum_{i=1}^m p_i \log p_i \end{aligned}$$

• As $H(x) = -\sum_{i=1}^m p_i \log p_i$ is the entropy of an arbitrary distribution, we want

$$\begin{aligned} \max \{H(x)\} &= \max \{ -D_{KL}(P \| U) - \log \frac{1}{m} \} \\ &= \min \{ D_{KL}(P \| U) \} - \log \frac{1}{m} \end{aligned}$$

• As $D_{KL}(P \| U)$ achieves a minimum of zero at $P = U$, we obtain that the entropy maximizing distribution is U .

Problem 2 [17 points]

Consider a parameter estimation problem where we wish to maximize conditional log-likelihood of a predictive model on a labeled (classification) data set. That is, we're given the population $\{(x_i, y_i)\}$, and a parametric likelihood function $p(y|x; \theta)$. show that the objective of *maximizing* the conditional log-likelihood of the model on the training data is equivalent to the objective of *minimizing* the KL divergence between the empirical distribution of the data $P_{\text{data}}(x, y)$ and the posterior distribution predicted by the model:

$$\operatorname{argmax}_{\theta} \sum_i \log p(y_i | x_i; \theta) = \operatorname{argmin}_{\theta} D_{KL}(P_{\text{data}}(y|x) || p(y|x; \theta)). \quad (1)$$

Advice: Recall that the empirical data distribution is defined by the specific data set:

$$P_{\text{data}}(x, y) = P_{\text{data}}^x(x) P_{\text{data}}^y(y|x)$$

where P_{data}^x assigns probability of $1/N$ on every x_i , and conditional probability P_{data}^y assigns probability 1 to y_i given x_i (and 0 to other values of y).

Also, note that the KL divergence between the conditional distributions, like the one that appears in (1), mean the expected KL divergence between the conditional distributions, with the expectation taken over the conditioning variable. That is, suppose we have joint distributions $p(a, b)$ and $q(a, b)$ over random variables a and b , from which we can derive the conditional $p(b|a)$, $q(b|a)$ and marginal $p(a)$, $q(a)$. Then

$$D_{KL}(p(b|a) || q(b|a)) = E_{a \sim p(a)} [D_{KL}(p(b|a) || q(b|a))]$$

End of problem 2

$$\text{RHS } D_{KL}(P_{\text{data}}(y|x) || p(y|x; \theta))$$

$$= E_{\bar{x}_i \sim P_{\text{data}}(x_i)} [D_{KL}(P_{\text{data}}(y|\bar{x}_i) || p(y|\bar{x}_i; \theta))]$$

$$= E_{\bar{x}_i \sim P_{\text{data}}(x_i)} \left[\sum_{j=1}^N P_{\text{data}}(y_j | \bar{x}_i) \cdot \log \frac{P_{\text{data}}(y_j | \bar{x}_i)}{p(y_j | \bar{x}_i; \theta)} \right]$$

$$= \sum_{i=1}^N \frac{1}{N} \cdot \left[\sum_{j=1}^N P_{\text{data}}(y_j | \bar{x}_i) \cdot \log \frac{P_{\text{data}}(y_j | \bar{x}_i)}{p(y_j | \bar{x}_i; \theta)} \right]$$

$$P_{\text{data}}(y_j | \bar{x}_i) = 0 \quad \text{if } j \neq i$$

$$P_{\text{data}}(y_j | \bar{x}_i) = 1 \quad \text{if } j = i$$

$$= \sum_{i=1}^N \frac{1}{N} \cdot (-\log p(y_i | x_i; \theta))$$

$$\text{Therefore, } \operatorname{argmin}_{\theta} D_{KL}(P_{\text{data}}(y|x) || p(y|x; \theta))$$

$$= \operatorname{argmin}_{\theta} \left\{ \frac{1}{N} \sum_{i=1}^N (-\log p(y_i | x_i; \theta)) \right\}$$

$$= \operatorname{argmin}_{\theta} \left\{ \sum_{i=1}^N \log p(y_i | x_i; \theta) \right\}.$$