
Problem Set 7 (Graded Homework - To be Submitted on Dec 23)
For the Exercise Sessions on Dec 02, Dec 09 and Dec 16

Last name	First name	SCIPER Nr	Points

Problem 1: Exponential Families and Maximum Entropy 1

Let $Y = X_1 + X_2$. Find the maximum entropy of Y under the constraint $\mathbb{E}[X_1^2] = P_1$, $\mathbb{E}[X_2^2] = P_2$:

- (a) If X_1 and X_2 are independent.
- (b) If X_1 and X_2 are allowed to be dependent.

Problem 2: Exponential Families and Maximum Entropy 2

Find the maximum entropy density f , defined for $x \geq 0$, satisfying $\mathbb{E}[X] = \alpha_1$, $\mathbb{E}[\ln X] = \alpha_2$. That is, maximize $-\int f \ln f$ subject to $\int x f(x) dx = \alpha_1$, $\int (\ln x) f(x) dx = \alpha_2$, where the integral is over $0 \leq x < \infty$. What family of densities is this?

Problem 3: Exponential Families and Maximum Entropy 3

For $t > 0$, consider a family of distributions supported on $[t, +\infty]$ such that $\mathbb{E}[\ln X] = \frac{1}{\alpha} + \ln t$, $\alpha > 0$.

1. What is the parametric form of a maximum entropy distribution satisfying the constraint on the support and the mean?
2. Find the exact form of the distribution.

Problem 4: Exponential Families and Maximum Entropy 4: I -projections

Let P denote the zero-mean and unit-variance Gaussian distribution. Assume that you are given N iid samples distributed according to P and let \hat{P}_N be the empirical distribution.

Let Π denote the set of distributions with second moment $\mathbb{E}[X^2] = 2$. We are interested in

$$\lim_{N \rightarrow \infty} \frac{1}{N} \log \Pr\{\hat{P}_N \in \Pi\} = - \inf_{Q \in \Pi} D(Q \| P).$$

- (a) Determine $-\operatorname{arginf}_{Q \in \Pi} D(Q \| P)$, i.e., determine the element Q for which the infimum is taken on.
- (b) Determine $-\inf_{Q \in \Pi} D(Q \| P)$.

Problem 5: Choose the Shortest Description

Suppose $\mathcal{C}_0 : \mathcal{U} \rightarrow \{0, 1\}^*$ and $\mathcal{C}_1 : \mathcal{U} \rightarrow \{0, 1\}^*$ are two prefix-free codes for the alphabet \mathcal{U} . Consider the code $\mathcal{C} : \mathcal{U} \rightarrow \{0, 1\}^*$ defined by

$$\mathcal{C}(u) = \begin{cases} [0, \mathcal{C}_0(u)] & \text{if } \text{length}\mathcal{C}_0(u) \leq \text{length}\mathcal{C}_1(u) \\ [1, \mathcal{C}_1(u)] & \text{else.} \end{cases}$$

Observe that $\text{length}(\mathcal{C}(u)) = 1 + \min\{\text{length}(\mathcal{C}_0(u)), \text{length}(\mathcal{C}_1(u))\}$.

- (a) Is \mathcal{C} a prefix-free code? Explain.
 (b) Suppose $\mathcal{C}_0, \dots, \mathcal{C}_{K-1}$ are K prefix-free codes for the alphabet \mathcal{U} . Show that there is a prefix-free code \mathcal{C} with

$$\text{length}(\mathcal{C}(u)) = \lceil \log_2 K \rceil + \min_{0 \leq k < K-1} \text{length}(\mathcal{C}_k(u)).$$

- (c) Suppose we are told that U is a random variable taking values in \mathcal{U} , and we are also told that the distribution p of U is one of K distributions p_0, \dots, p_{K-1} , but we do not know which. Using (b) describe how to construct a prefix-free code \mathcal{C} such that

$$\mathbb{E}[\text{length}(\mathcal{C}(U))] \leq \lceil \log_2 K \rceil + 1 + H(U).$$

[Hint: From class we know that for each k there is a prefix-free code \mathcal{C}_k that describes each letter u with at most $\lceil -\log_2 p_k(u) \rceil$ bits.]

Problem 6: Prediction and coding

After observing a binary sequence u_1, \dots, u_i , that contains $n_0(u^i)$ zeros and $n_1(u^i)$ ones, we are asked to estimate the probability that the next observation, u_{i+1} will be 0. One class of estimators are of the form

$$\hat{P}_{U_{i+1}|U^i}(0|u^i) = \frac{n_0(u^i) + \alpha}{n_0(u^i) + n_1(u^i) + 2\alpha} \quad \hat{P}_{U_{i+1}|U^i}(1|u^i) = \frac{n_1(u^i) + \alpha}{n_0(u^i) + n_1(u^i) + 2\alpha}.$$

We will consider the case $\alpha = 1/2$, this is known as the Krichevsky–Trofimov estimator. Note that for $i = 0$ we get $\hat{P}_{U_1}(0) = \hat{P}_{U_1}(1) = 1/2$.

Consider now the joint distribution $\hat{P}(u^n)$ on $\{0, 1\}^n$ induced by this estimator,

$$\hat{P}(u^n) = \prod_{i=1}^n \hat{P}_{U_i|U^{i-1}}(u_i|u^{i-1}).$$

- (a) Show, by induction on n that, for any n and any $u^n \in \{0, 1\}^n$,

$$\hat{P}(u_1, \dots, u_n) \geq \frac{1}{2\sqrt{n}} \left(\frac{n_0}{n}\right)^{n_0} \left(\frac{n_1}{n}\right)^{n_1},$$

where $n_0 = n_0(u^n)$ and $n_1 = n_1(u^n)$.

[Hint: if $0 \leq m \leq n$, then $(1 + 1/n)^{n+1/2} \geq \frac{m+1}{m+1/2} (1 + 1/m)^m$]

- (b) Conclude that there is a prefix-free code $\mathcal{C} : \mathcal{U} \rightarrow \{0, 1\}^*$ such that

$$\text{length}\mathcal{C}(u_1, \dots, u_n) \leq nh_2\left(\frac{n_0(u^n)}{n}\right) + \frac{1}{2} \log n + 2,$$

with $h_2(x) = -x \log x - (1 - x) \log(1 - x)$.

(c) Show that if U_1, \dots, U_n are i.i.d. Bernoulli, then

$$\frac{1}{n} \mathbb{E}[\text{length } \mathcal{C}(U_1, \dots, U_n)] \leq H(U_1) + \frac{1}{2n} \log n + \frac{2}{n}$$

Problem 7: Universal codes

Suppose we have an alphabet \mathcal{U} , and let Π denote the set of distributions on \mathcal{U} . Suppose we are given a family of S of distributions on \mathcal{U} , i.e., $S \subset \Pi$. For now, assume that S is finite.

Define the distribution $Q_S \in \Pi$

$$Q_S(u) = Z^{-1} \max_{P \in S} P(u)$$

where the normalizing constant $Z = Z(S) = \sum_u \max_{P \in S} P(u)$ ensures that Q_S is a distribution.

- (a) Show that $D(P||Q) \leq \log Z \leq \log |S|$ for every $P \in S$.
- (b) For any S , show that there is a prefix-free code $\mathcal{C} : \mathcal{U} \rightarrow \{0, 1\}^*$ such that for any random variable U with distribution $P \in S$,

$$E[\text{length } \mathcal{C}(U)] \leq H(U) + \log Z + 1.$$

(Note that \mathcal{C} is designed on the knowledge of S alone, it cannot change on the basis of the choice of P .) [Hint: consider $L(u) = -\log_2 Q_S(u)$ as an ‘almost’ length function.]

- (c) Now suppose that S is not necessarily finite, but there is a finite $S_0 \subset \Pi$ such that for each $u \in \mathcal{U}$, $\sup_{P \in S} P(u) \leq \max_{P \in S_0} P(u)$. Show that $Z(S) \leq |S_0|$.

Now suppose $\mathcal{U} = \{0, 1\}^m$. For $\theta \in [0, 1]$ and $(x_1, \dots, x_m) \in \mathcal{U}$, let

$$P_\theta(x_1, \dots, x_m) = \prod_i \theta^{x_i} (1 - \theta)^{1-x_i}.$$

(This is a fancy way to say that the random variable $U = (X_1, \dots, X_m)$ has i.i.d. Bernoulli θ components). Let $S = \{P_\theta : \theta \in [0, 1]\}$.

- (d) Show that for $u = (x_1, \dots, x_m) \in \{0, 1\}^m$

$$\max_{\theta} P_\theta(x_1, \dots, x_m) = P_{k/m}(x_1, \dots, x_m)$$

where $k = \sum_i x_i$.

- (e) Show that there is a prefix-free code $\mathcal{C} : \{0, 1\}^m \rightarrow \{0, 1\}^*$ such that whenever X_1, \dots, X_m are i.i.d. Bernoulli,

$$\frac{1}{m} \mathbb{E}[\text{length } \mathcal{C}(X_1, \dots, X_m)] \leq H(X_1) + \frac{1 + \log_2(1 + m)}{m}.$$