# ÉCOLE POLYTECHNIQUE FÉDÉRALE DE LAUSANNE
## School of Computer and Communication Sciences

Foundations of Data Science      Assignment date: Friday, February 3rd, 2023, 9:15 am

Fall 2022      Due date: Friday, Feburary 3rd, 2023, 12:15 noon

# Final Exam – SG0211

This exam is open book. No electronic devices of any kind are allowed. There are 4 problems. Choose the ones you find easiest and collect as many points as possible. We do not necessarily expect you to finish all of them. Good luck!

Name: _____

| | |
|---|---|
| Problem 1 | / 15 |
| Problem 2 | / 15 |
| Problem 3 | / 12 |
| Problem 4 | / 20 |
| **Total** | /62 |

**Problem 1** (*Fisher Goes Exponential*). [15 pts]

Let $p_\theta(x)$ denote a family of distributions parameterized by $\theta$. Define the Fisher information as

$$I_\theta = \mathbb{E}_\theta[\nabla_\theta \log p_\theta(X)(\nabla_\theta \log p_\theta(X))^T].$$

(1) [5pts] Let $p_\theta(x) = h(x)e^{\langle\theta,\phi(x)\rangle - A(\theta)}$ be an exponential family. What is the Fisher information in terms of the parameters of the family?

(2) [5pts - 1pt per question] Consider distributions of the form $p_\lambda(x) = \lambda e^{-\lambda x}$, where $\lambda \in \mathbb{R}^+$.

    1. Write it in the form of an exponential family.

    2. What is $\Theta = \{\theta \in \mathbb{R} : A(\theta) < \infty\}$.

    3. Is the family regular?

    4. Is it minimal?

    5. What is the Fisher information?

(3) [5pts - 1pt per question] Consider distributions of the form $p_p(k) = (1 - p)^k p$, where $p \in (0, 1)$ and $k \in \mathbb{N}$.

    1. Write it in the form of an exponential family.

    2. What is $\Theta = \{\theta \in \mathbb{R} : A(\theta) < \infty\}$.

    3. Is the family regular?

    4. Is it minimal?

    5. What is the Fisher information?

(space for problem 1)

(space for problem 1)

(space for problem 1)

**Problem 2** (*Compression*). [15 pts]

Suppose $\mathcal{P} \in \Pi(\mathcal{X}, \mathcal{Y})$ be a probability distribution on $\mathcal{X} \times \mathcal{Y}$ and $(X, Y)$ be a joint random variable with distribution $P_{XY}$ with marginals $P_X$ and $P_Y$.

In what follows, assume that **all codes are optimal, prefix-free, and binary**. Optimal here means having smallest possible average length. All logs are to the base 2.

(1) [1 pt] Let $c_X : \mathcal{X} \to \{0,1\}^*$ and $c_Y : \mathcal{Y} \to \{0,1\}^*$ be optimal prefix free codes. What are lower and upper bounds for the expected length of these codes $c_X$ and $c_Y$?

(2) [1 pt] Let $c_{XY} : \mathcal{X} \times \mathcal{Y} \to \{0,1\}^*$ be an optimal prefix free code. What are lower and upper bounds for the expected length of this code?

(3) [10 pts total] In this sub problem, assume that $X, Y$ have a joint distribution according to the following table:

|        | Y=0   | Y=1   |
|--------|-------|-------|
| X=0    | $1/4$ | $0$   |
| X=1    | $1/8$ | $1/8$ |
| X=2    | $1/8$ | $1/8$ |
| X=3    | $0$   | $1/4$ |

(a) [4 pts] What are lower and upper bounds for the expected lengths of $c_X$ and $c_Y$? Are the lower bounds tight?

(b) [3 pts] What are lower and upper bounds for the expected lengths of $c_{XY}$? Is the lower bound tight?

(c) [3 pts] For the above joint distribution, is it more efficient to compress separately and concatenate the individual code words (which, as we saw in the lecture, is guaranteed to yield a prefix free code), or to compress $(X, Y)$ jointly (again, in a prefix free manner)?

(4) [3 pts] Assume that $(X, Y)$ has some generic joint distribution. Assume further that $I(X; Y) > 1$. Show that in this case optimal joint prefix free compression is more efficient than compressing individually and concatenating.

(space for problem 2)

(space for problem 2)

(space for problem 2)

**Problem 3** (*Stability implies Generalization*). [12 pts]

Let $S = \{(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)\}$ be a training dataset composed of $n$ i.i.d. samples drawn from $\mathcal{D}$. As usual, we denote $L_{\mathcal{D}}(h) = E_{(x,y)\sim\mathcal{D}}[l(h(x), y)]$ and $L_{\mathcal{S}}(h) = \frac{1}{n}\sum_{i=1}^{n} l(h(x_i), y_i)$ the true and empirical risks of a hypothesis $h$, respectively. For simplicity, let us denote by $h_S$ the output of a learning algorithm when trained with dataset $S$.

An important property of learning algorithms is their ability to generalize, i.e., the true and empirical risks of the output hypothesis should be close in expectation. Formally, we say that a learning algorithm $\mathcal{A}$ $\epsilon$-generalizes in expectation if

$$|E_S[L_S(h_S) - L_{\mathcal{D}}(h_S)]| < \epsilon . \tag{1}$$

An interesting connection arises when we investigate the *stability* of a learning algorithm. Formally, we call a learning algorithm $\epsilon$-*uniformly stable* if $\forall S, S'$ datasets of size $n$ that differ in at most one sample we have

$$\sup_{(x,y)} l(h_S(x), y) - l(h_{S'}(x), y) < \epsilon . \tag{2}$$

<u>Notations:</u> $(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n), (\widetilde{x}_1, \widetilde{y}_1), \ldots, (\widetilde{x}_n, \widetilde{y}_n)$ are $2n$ independently sampled training examples. We define $S = \{(x_1, y_1), \ldots, (x_n, y_n)\}$, $\widetilde{S} = \{(\widetilde{x}_1, \widetilde{y}_1), \ldots, (\widetilde{x}_n, \widetilde{y}_n)\}$ and $S^{(i)} = \{(x_1, y_1), \ldots, (x_{i-1}, y_{i-1}), (\widetilde{x}_i, \widetilde{y}_i), (x_{i+1}, y_{i+1}), \ldots, (x_n, y_n)\}$.

(1) [2 pts] Prove that $L_{\mathcal{D}}(h_S) = E_{\widetilde{S}}[\frac{1}{n}\sum_{i=1}^{n} l(h_S(\widetilde{x}_i), \widetilde{y}_i)]$.

(2) [3 pts] Prove that $E_{S,\widetilde{S}}[l(h_S(\widetilde{x}_i), \widetilde{y}_i)] = E_{S,S^{(i)}}[l(h_{S^{(i)}}(x_i), y_i)]$.

(3) [7 pts] Prove that an $\epsilon$-uniformly stable learning algorithm $\epsilon$-generalizes in expectation, by justifying each step in the following sequence.

$$|E_S[L_S(h_S) - L_{\mathcal{D}}(h_S)]| \overset{(a)}{=} \left|E_S\left[L_S(h_S) - E_{\tilde{S}}\left[\frac{1}{n}\sum_{i=1}^{n} l(h_S(\tilde{x}_i), \tilde{y}_i)\right]\right]\right|$$

$$\overset{(b)}{=} \left|E_S\left[L_S(h_S)\right] - E_{S,\tilde{S}}\left[\frac{1}{n}\sum_{i=1}^{n} l(h_S(\tilde{x}_i), \tilde{y}_i)\right]\right|$$

$$\overset{(c)}{=} \left|E_S\left[L_S(h_S)\right] - \frac{1}{n}\sum_{i=1}^{n} E_{S,\tilde{S}}\left[l(h_S(\tilde{x}_i), \tilde{y}_i)\right]\right|$$

$$\overset{(d)}{=} \left|E_S\left[L_S(h_S)\right] - \frac{1}{n}\sum_{i=1}^{n} E_{S^{(i)},(x_i,y_i)}\left[l(h_{S^{(i)}}(x_i), y_i)\right]\right|$$

$$\overset{(e)}{=} \left|E_S\left[\frac{1}{n}\sum_{i=1}^{n} l(h_S(x_i), y_i)\right] - \frac{1}{n}\sum_{i=1}^{n} E_{S,S^{(i)}}\left[l(h_{S^{(i)}}(x_i), y_i)\right]\right|$$

$$\overset{(f)}{=} \left|\frac{1}{n}\sum_{i=1}^{n} E_{S,S^{(i)}}\left[l(h_S(x_i), y_i)) - l(h_{S^{(i)}}(x_i), y_i)\right]\right|$$

$$\overset{(g)}{\leq} \frac{1}{n}\sum_{i=1}^{n} \epsilon = \epsilon$$

(space for problem 3)

(space for problem 3)

(space for problem 3)

**Problem 4** (*Multi-arm Bandits* ). [20 pts]

We consider the following game where in each round $t$ we can choose between $[N] = \{1, 2, \ldots, N\}$ different actions. After we choose an action $a_t \in [N]$ an adversary reveals the loss of each action in this round, call it $l_i^t \in [0, 1]$, $i \in [N]$. Note that this is an adversarial setting, where the losses do not come from a probability distribution. This setting differs from what we had discussed in class where only the loss for the chosen action was revealed.

Our goal is to design a randomized algorithm $\mathcal{A}$ which maintains a probability distribution $p^t$ over actions, and achieves a sub-linear regret, i.e., $\mathcal{R}(T) = \max_i\{\sum_{t=1}^{T} E_{A_t \sim p^t}[l_{A_t}^t - l_i^t]\} \leq o(T)$. We also note that the adversary may know the probability distribution $p^t$, but does not know the realizations $A_t$. We will analyze the following algorithm:

---
**Algorithm 1:** Multiplicative Weights Update
---
    **Input:** learning parameter $\epsilon$
    **Initialization:** $p_i^1 = 1/N, w_i^1 = 1, \forall i \in [N], \Phi^1 = N$
    **for** $t = 1$ to $T$ **do**
        $A_t \sim p^t$
        Adversary reveals the loss vector $l^t$ and we suffer $l_{A_t}^t$
        Update weights $w_i^{t+1} = w_i^t \cdot \exp(-\epsilon \cdot l_i^t), \forall i \in [N]$ and let $\Phi^{t+1} = \sum_i w_i^{t+1}$
        Update the probability distribution: $p_i^{t+1} = w_i^{t+1}/\Phi^{t+1}, \forall i$
    **end for**
---

(1) [2 pts] Prove that $w_i^{T+1} = \exp(-\epsilon \cdot \sum_{t=1}^{T} l_i^t), \forall i \in [N]$

(2) [8 pts] Prove that $\Phi^{t+1} \leq \Phi^t \cdot \exp(\epsilon^2 - \epsilon \langle p^t, \ l^t \rangle)$

    *Hint:* Note that $w_i^{t+1} = p_i^{t+1} \cdot \Phi^{t+1}$ and use the inequalities: (a) $e^x \leq 1 + x + x^2, \forall x \in [0, 1]$ and (b) $e^x \geq x + 1, \forall x$.

(3) [2 pts] Prove that $\Phi^{T+1} \leq \Phi^1 \cdot \exp(\epsilon^2 \cdot T - \epsilon \sum_{t=1}^{T} \langle p^t, \ l^t \rangle)$

(4) [8 pts] By noting that $\Phi^1 \cdot \exp(\epsilon^2 \cdot T - \epsilon \sum_{t=1}^{T} \langle p^t, \ l^t \rangle) \geq \Phi^{T+1} \geq w_i^{T+1}, \forall i \in [N]$ set the learning parameter $\epsilon$ so that $\mathcal{R}(T) \leq 2\sqrt{\log(N) \cdot T}$.

(space for problem 4)

(space for problem 4)

(space for problem 4)