

## Homework 2

Due: Friday, October 28, 2016

**Notes:** For positive integers  $k$ ,  $[k] := \{1, \dots, k\}$  denotes the set of the first  $k$  positive integers. When  $X \sim p$  and  $Y \sim q$  are random variables over the same sample space,  $D(X||Y)$ ,  $D(X||q)$ , and  $D(p||Y)$  should all be read as  $D(p||q)$ . The homework is out of 75 points – 5 points per part.

## 1. Maximum Entropy of Independent Bernoulli Sums

In this problem, we will show that the binomial and (optionally) Poisson distributions are maximum entropy (MaxEnt) distributions over an appropriate class  $\mathcal{P}$  of distributions, and derive several useful properties of KL divergence along the way.

For any positive integer  $n$  and  $p \in [0, 1]$ , let  $\text{Binomial}(n, p)$  denote the binomial distribution (the sum of  $n$  IID Bernoulli events of probability  $p$ ), which has density function

$$\text{Binomial}_{n,p}(k) = \binom{n}{k} p^k (1-p)^{1-k}.$$

For  $\lambda \geq 0$ , let  $\Pi(\lambda)$  denote the mean- $\lambda$  Poisson distribution, which has density function

$$\text{Poisson}_\lambda(k) = \frac{\lambda^k}{k!} e^{-\lambda}, \quad \forall k \in \mathbb{N} \cup \{0\}.$$

The class  $\mathcal{P}_\lambda$  of distributions is that of sums  $S_n := \sum_{i=1}^n X_i$  of  $n$  independent (but not necessarily identically distributed) binary variables  $\{X_i\}_{i=1}^n$  constrained such that  $\mathbb{E}[S_n] = \lambda$ , for some  $\lambda \in [0, n]$ . Note that any  $p \in \mathcal{P}_\lambda$  can be parametrized by  $(p_1, \dots, p_n) \in [0, 1]^n$ , with  $\sum_{i=1}^n p_i = \lambda$ . We will show that the Binomial case  $p_1 = \dots = p_n = \frac{\lambda}{n}$  is the MaxEnt distribution over  $\mathcal{P}_\lambda$ , and that the Poisson distribution is the limit as  $n \rightarrow \infty$ .

- Derive the maximum likelihood estimate of  $\lambda$  under the assumption that you observe  $n$  IID samples  $X_1, \dots, X_n$  from a Poisson distribution.
- Define  $D(X) := \min_{\lambda \geq 0} D(X||\Pi(\lambda))$ . Derive a closed form for  $D(X)$  in terms of  $X$ .<sup>1</sup>
- Show that the KL divergence  $D(p||q)$  is convex in  $p$ .
- Let

$$\begin{aligned} \mathcal{P}_\lambda(p_3, \dots, p_n) &= \{q \in \mathcal{P}_\lambda : q_3 = p_3, \dots, q_n = p_n, \} \\ &= \left\{ (x_1, x_2, p_3, \dots, p_n) : x_1 + x_2 = \lambda - \sum_{i=3}^n p_i \right\} \end{aligned}$$

denote the subspace of  $\mathcal{P}_\lambda$  with all but two coordinates fixed. Show that  $H(S_n)$  is strictly concave on  $\mathcal{P}_\lambda(p_3, \dots, p_n)$ . (*Hint: Use parts (b) and (c) to reduce this to showing  $\mathbb{E}[\log(S_n!)]$  is strictly concave on  $\mathcal{P}_\lambda(p_3, \dots, p_n)$ . Then, since*

$$\mathbb{E}[\log(S_n!)] = \mathbb{E}[\mathbb{E}[\log(S_n!) | X_3, \dots, X_n]],$$

<sup>1</sup> $X$  may have any distribution over  $\{0, 1, 2, \dots\}$ , but you may assume any necessary functionals of  $X$  are finite.

which is a linear functional of  $\mathbb{E}[\log(S_n!)|X_3, \dots, X_n]$ , show that  $\mathbb{E}[\log(S_n!)|X_3, \dots, X_n]$  is strictly concave on  $\mathcal{P}_\lambda(p_3, \dots, p_n)$ , for any values of  $X_3, \dots, X_n$ .)

- (e) Use part (d) to show that  $\text{Binomial}(n, \lambda/n)$  is the unique MaxEnt distribution over  $\mathcal{P}$ .
- (f) Given independent random variables  $X$  and  $Y$  taking values on  $\mathbb{N}$ , show that

$$D(X + Y) \leq D(X) + D(Y). \quad (1)$$

(Hint: Use the General Data Processing Inequality from Homework 1 and the fact that the sum of two Poisson-distributed variables with means  $\lambda_1$  and  $\lambda_2$  is itself Poisson-distributed with mean  $\lambda_1 + \lambda_2$ .)

- (g) Show that  $D(\text{Binomial}(n, \frac{\lambda}{n})) \rightarrow 0$  as  $n \rightarrow \infty$ . This is (a fairly strong form of) the “Law of Rare Events” (a.k.a. the “Poisson Limit Theorem”), which states that the frequency of a large number of unlikely events is approximately Poisson-distributed and justifies many applications of the Poisson distribution. (Hint: Show  $D(X_i) \leq p_i^2$  and apply (1).)
- (h) **(This part is optional.)** Show that  $H(\Pi(\lambda)) = \lim_{n \rightarrow \infty} H(B(n, \lambda/n))$ . (Hint: Use the equivalence

$$H(p) + D(p||q) = \mathbb{E}_{X \sim p} [\log q(x)],$$

discussed in Lecture 1. Note that one step of this proof requires switching a limit and an infinite summation. If you are not familiar with the dominated convergence theorem, you may wish to take this step for granted.)

## 2. Wavelet Denoising with CRM

In this problem, we will analyze the convergence rate of a wavelet-based denoising estimator.

*Haar wavelets and quantization:* Recall that Haar wavelets over  $\mathcal{X} := [0, 1)$  are piecewise constant functions  $\psi_{j,k} : \mathcal{X} \rightarrow \{-2^{j/2}, 0, 2^{j/2}\}$  such that

$$\psi_{j,k}(x) = 2^{j/2} (1_{[k2^{-j}, (k+1/2)2^{-j})} - 1_{[(k+1/2)2^{-j}, (k+1)2^{-j})}),$$

for all  $j \in \mathbb{N} \cup \{0\}$ ,  $k \in \{0, \dots, 2^j - 1\}$ ,  $x \in \mathcal{X}$ . Since Haar wavelets form a basis for  $L^2(\mathcal{X})$ , for any  $\ell \in \mathbb{N} \cup \{0\}$ , if we define the projection

$$f_\ell := \sum_{j=0}^{\ell} \sum_{k=0}^{2^j-1} \langle \psi_{j,k}, f \rangle,$$

of  $f$  onto the first  $\ell + 1$  scales of the Haar basis, then  $f_\ell \rightarrow f$  as  $\ell \rightarrow \infty$ . To encode the projection  $f_\ell$ , we also need to quantize the coefficients. Quantized projections lie in the set

$$Q_{\ell, \varepsilon} := \left\{ \sum_{j=0}^{\ell} \sum_{k=0}^{2^j-1} a_{j,k} \psi_{j,k} \in L^2(\mathcal{X}) : a_{j,k} = 2b_{j,k}\varepsilon, \text{ for some integer } b_{j,k} \right\},$$

so that their wavelet coefficients are multiples of  $\varepsilon$ . Our quantized projection of  $f$  is then

$$f_{\ell, \varepsilon} := \operatorname{argmin}_{g \in Q_{\ell, \varepsilon}} \|f - g\|_2.$$

Thus,  $f_{\ell,\varepsilon}$  is the best (in  $L^2$  distance) representation of  $f$  in terms of Haar wavelets of scale at most  $\ell$  and coefficient precision  $\varepsilon$ .

*CRM Denoising:* We will assume the true function  $f$  lies in the class  $\mathcal{F}_{s,M} \subseteq L^2(\mathcal{X})$  of piecewise constant functions with at most  $s$  discontinuities and bounded  $L^\infty$  norm  $\|f\|_\infty = \sup_{x \in \mathcal{X}} |f(x)| \leq M$ . We observe  $n$  noisy IID pairs  $\{(X_i, Y_i)\}_{i=1}^n$ , where each  $X_1, \dots, X_n \sim U(\mathcal{X})$  is uniformly distributed and, for  $\varepsilon_1, \dots, \varepsilon_n \sim \mathcal{N}(0, \sigma^2)$ ,  $Y_i = f(X_i) + \varepsilon_i$ .

For  $\delta \in (0, 1)$ , the complexity-penalized empirical risk minimizing (CRM) estimator <sup>2</sup> is

$$\hat{f}_{\ell,\varepsilon,\delta} := \operatorname{argmin}_{g_{\ell,\varepsilon} \in Q_{\ell,\varepsilon}} \left[ \|g_{\ell,\varepsilon} - f\|_2^2 + \frac{c(g_{\ell,\varepsilon}) - \ln \delta}{n} \right],$$

where  $c(g_{\ell,\varepsilon})$  denotes the number of bits required to encode  $g_{\ell,\varepsilon}$ . In class, we derived the following excess risk bound for CRM estimators:

$$R(\hat{f}_{\ell,\varepsilon,\delta}) - R^* = \|\hat{f}_{\ell,\varepsilon,\delta} - f\|_2^2 \leq \inf_g \left[ \|g_{\ell,\varepsilon} - f\|_2^2 + \frac{c(g_{\ell,\varepsilon}) - \ln \delta}{n} \right] + \delta. \quad (2)$$

In this problem, we will analyze the terms of (2) to derive a convergence rate bound in terms of the complexity  $s$  of  $f$  and the sample size  $n$ .

- (a) Show that the projections  $f_\ell$  and  $f_{\ell,\varepsilon}$  can each have at most  $C_0 s \ell + 1$  nonzero coefficients, for some constant  $C_0$ .
- (b) Bound the approximation errors  $\|f - f_\ell\|_2^2$  and  $\|f - f_{\ell,\varepsilon}\|_2^2$ .
- (c) How many bits  $c(f)$  are required to encode  $f_{\ell,\varepsilon}$  (for known  $s$ ,  $M$ ,  $\ell$ , and  $\varepsilon$ )?
- (d) By choosing  $\varepsilon > 0$ ,  $\ell \in \mathbb{N}$ , and  $\delta > 0$  appropriately, use parts (b) and (c) with the bound(2) show <sup>3</sup>

$$\|\hat{f} - f\|_2^2 \in O\left(\frac{s \log^2 n}{n}\right).$$

Note that, up to log factors, this is a parametric rate with  $s$  parameters.

---

<sup>2</sup>Recall that  $\hat{f}_{\ell,\varepsilon,\delta}$  can be easily computed by hard-thresholding.

<sup>3</sup>Here, treat  $M$  as a constant.

### 3. Universal Prediction with Exponential Weights

Fix a (potentially infinite) countable class of predictors  $\mathcal{F}$ . Recall that, in the universal prediction setting, at each time point  $t \in \{1, \dots, T\}$  up to a predetermined time horizon  $T$ , we see some data  $x_t$  and choose a predictor  $\hat{f}_t \in \mathcal{F}$ , before then seeing a true label  $y_t$  and suffering loss  $\ell(\hat{f}_t(x_t), y_t) \in [0, 1]$ . Since we are allowing, for example, adversarial sequences  $\{(x_t, y_t)\}_{t=1}^T$ , a randomized algorithm is needed to provide any guarantees. Given a learning rate  $\eta > 0$  and prior  $\pi$  over  $\mathcal{F}$ , the exponential weights algorithm proposes to draw  $\hat{f}_t$  according to a distribution  $q_t$  defined such that  $q_1 = \pi$  and each

$$q_{t+1}(f) \propto q_t(f) \exp(-\eta \ell(f(x_t), y_t)).$$

For each  $f \in \mathcal{F}$  and  $t \in [T]$ , let

$$L_t(f) := \sum_{\tau=1}^t \ell(f(x_\tau), y_\tau) \quad \text{and} \quad L_t(\hat{f}) := \sum_{\tau=1}^t \ell(\hat{f}_\tau(x_\tau), y_\tau)$$

denote the cumulative losses of  $f$  and our predictions, respectively, at time  $t$ . Define

$$W_t = \mathbb{E}_{f \sim \pi} [\exp(-\eta L_t(f))], \quad \forall t \in \{1, \dots, T\}.$$

(a) Show that  $\ln W_T \geq -\inf_{f \in \mathcal{F}} [\eta L_T(f) - \log \pi(f)]$ .

(b) Show that

$$\frac{W_{t+1}}{W_t} = \mathbb{E}_{f \sim q_{t+1}} [\exp(-\eta \ell(f(x_{t+1}), y_{t+1}))].$$

(c) Use part (b) to show that

$$\ln W_T \leq -\eta \sum_{t=1}^T \mathbb{E}_{f \sim q_t} [\ell(f(x_t), y_t)] + \frac{\eta^2 T}{8}.$$

*Hint: Recall Hoeffding's Lemma: for a random variable  $X$  with  $X \in [a, b]$  a.s.,*

$$\ln \mathbb{E}[e^{sX}] \leq s \mathbb{E}[X] + \frac{s^2(b-a)^2}{8}.$$

(d) Use parts (a) and (c) and a convenient choice of  $\eta$  to bound the expected loss of the exponential weights algorithm by

$$\mathbb{E}[L_T(\hat{f})] \leq \inf_{f \in \mathcal{F}} \left[ L_T(f) + (1 - \log \pi(f)) \sqrt{\frac{T}{8}} \right].$$

If  $\mathcal{F}$  is finite, give a simple sufficient condition on the prior  $\pi$  such that the regret

$$\mathbb{E}[L_T(\hat{f})] - \inf_{f \in \mathcal{F}} L_T(f) \in O(T^{1/2}).$$