

Homework 2

Due: November 1, 2017

Note: You may discuss these problems in groups. However, you must write up your own solutions and mention the names of the people in your group. Also, please do mention any books, papers or other sources you refer to. It is recommended that you typeset your solutions in \LaTeX .

1. **Biased coins strike back.** In class we considered the problem of distinguishing coins distributed according to the following two distributions:

$$P = \begin{cases} 1 & \text{w.p. } \frac{1}{2} - \varepsilon \\ 0 & \text{w.p. } \frac{1}{2} + \varepsilon \end{cases} \quad \text{and} \quad Q = \begin{cases} 1 & \text{w.p. } \frac{1}{2} \\ 0 & \text{w.p. } \frac{1}{2} \end{cases}$$

We derived matching upper and lower bounds (up to constants) of the form $\Theta(1/\varepsilon^2)$ on the number of coin tosses required to distinguish the two distributions. Consider now the problem of distinguishing two extremely biased coins with slightly differing biases:

$$P' = \begin{cases} 1 & \text{w.p. } \varepsilon \\ 0 & \text{w.p. } 1 - \varepsilon \end{cases} \quad \text{and} \quad Q' = \begin{cases} 1 & \text{w.p. } 2\varepsilon \\ 0 & \text{w.p. } 1 - 2\varepsilon \end{cases}$$

Find tight upper and lower bounds (up to constants) on the number of independent coin tosses required to distinguish coins distributed according to P' and Q' .

2. **Jensen-Shannon divergence.** While KL-divergence is sometimes used as a measure of the difference between two distributions, it is asymmetric and can be infinite. In some applications, one can instead consider the Jensen-Shannon divergence which addresses these issues.

- (a) For two distributions P and Q , we define the Jensen-Shannon divergence as

$$\text{JSD}(P, Q) := \frac{1}{2} \cdot D(P \| M) + \frac{1}{2} \cdot D(Q \| M) \quad \text{where} \quad M = \frac{P + Q}{2}.$$

Show that $0 \leq \text{JSD}(P, Q) \leq 1$.

- (b) Show that $\text{JSD}(P, Q) \geq \frac{1}{8 \ln 2} \cdot \|P - Q\|_1^2$.

- (c) The notion of Jensen-Shannon divergence can be generalized to an arbitrary number of distributions and an arbitrary convex combination. Let P_1, \dots, P_k be distributions on the same universe and let $\lambda = (\lambda_1, \dots, \lambda_k)$ be a tuple of non-negative weights such that $\sum_i \lambda_i = 1$. We define

$$\text{JSD}_\lambda(P_1, \dots, P_k) := \sum_i \lambda_i \cdot D(P_i \| M) \quad \text{where} \quad M = \sum_i \lambda_i P_i.$$

Show that $0 \leq \text{JSD}_\lambda(P_1, \dots, P_k) \leq H(\lambda)$, where $H(\lambda)$ denotes the entropy of λ , when viewed as a distribution over $[k]$.

3. **Counting using method of types (Problem 11.5 from the book).** Let U be a finite universe with $|U| = r$ and let $g : U \rightarrow \mathbb{R}$ be a real valued function. Let $S \subseteq U^n$ be the set of sequences x_1, \dots, x_n with each $x_i \in U$ defined as

$$S = \left\{ (x_1, \dots, x_n) \in U^n \mid \frac{1}{n} \sum_{i=1}^n g(x_i) \geq \alpha \right\}.$$

Let $\Pi = \{P \mid \sum_{a \in U} P(a)g(a) \geq \alpha\}$. Show that

$$|S| \leq (n+1)^r \cdot 2^{nH^*},$$

where $H^* = \sup_{P \in \Pi} H(P)$.

4. **Differential entropy of a Gaussian.** We saw in class that if the differential entropy $h(X)$ exists for a random variable X with a probability density, then $h(aX) = h(X) + \log |a|$, for $a \in \mathbb{R}$. For this problem, you may assume the n -dimensional generalization of this fact: if X is a random variable taking values in \mathbb{R}^n such that $h(X)$ exists, and $A \in \mathbb{R}^{n \times n}$ is a non-singular matrix, then

$$h(AX) = h(X) + \log |A|,$$

where $|A|$ denotes the absolute value of the determinant of A . We can use this to compute the entropy of a Gaussian random variable.

- (a) Let $X \sim N(0, I_n)$ denote an n -dimensional Gaussian random variable (with mean 0) such that each coordinate is an *independent* one-dimensional Gaussian with mean 0 and variance 1 i.e., the covariance matrix is I_n . Calculate $h(X)$ (you may use the expression for the entropy of a one-dimensional Gaussian derived in class).
- (b) Let $Y \sim N(0, \Sigma)$ be an n -dimensional Gaussian with mean 0 and covariance matrix Σ i.e.,

$$\mathbb{E}[Y] = 0 \quad \text{and} \quad \mathbb{E}[YY^T] = \Sigma.$$

Assume that the covariance matrix Σ is *positive definite* and hence there exists a non-singular matrix R such that $\Sigma = R^2$. Use this to show that

$$H(Y) = \frac{n}{2} \log(2\pi e) + \frac{1}{2} \log |\Sigma| .$$

- (c) Use the above to show that for any two positive definite matrices Σ_1 and Σ_2 , and $\alpha \in [0, 1]$, we have

$$|\alpha \cdot \Sigma_1 + (1 - \alpha) \cdot \Sigma_2| \geq |\Sigma_1|^\alpha \cdot |\Sigma_2|^{1-\alpha} .$$

5. **Chernoff bound for read- k families.** We used Sanov's theorem to derive the Chernoff bound for independent random variables X_1, \dots, X_n taking values uniformly in $\{0, 1\}$. In particular, we showed that

$$\mathbb{P} \left[X_1 + \dots + X_n \geq \left(\frac{1}{2} + \varepsilon \right) n \right] \leq (n + 1) \cdot 2^{-n \cdot D(\frac{1}{2} + \varepsilon \| \frac{1}{2})} ,$$

where $D(\frac{1}{2} + \varepsilon \| \frac{1}{2})$ denotes the KL-divergence of two distributions on $\{0, 1\}$, with probabilities $(\frac{1}{2} + \varepsilon, \frac{1}{2} - \varepsilon)$ and $(\frac{1}{2}, \frac{1}{2})$. In this problem, we will consider functions f_1, \dots, f_r depending on the variables X_1, \dots, X_n and prove a concentration bound on the expression $f_1 + \dots + f_r$.

Let S_1, \dots, S_r be subsets of $[n]$ for each $i \in [r]$, let $f_i : \{0, 1\}^{S_i} \rightarrow \{0, 1\}$ be a function which depends only on the variables in S_i . We use the shorthand X_{S_i} to denote the variables $\{X_j\}_{j \in S_i}$. Moreover, we have the property that each variable is involved in only k functions i.e., $\forall j \in [n], |\{i \in [r] \mid j \in S_i\}| = k$. Such a family of functions is called a read- k family (it is not too hard to see that the lower bound extends to the case when each variable is in *at most* k functions).

- (a) Recall that for two random variables Z_1 and Z_2 distributed on *same universe* U , we also use $D(Z_1 \| Z_2)$ to mean $D(P_1 \| P_2)$. Let Y_1, \dots, Y_n be (not necessarily independent) random variables jointly distributed on $\{0, 1\}^n$ and let X_1, \dots, X_n be random variables as above, distributed uniformly and independently on $\{0, 1\}^n$. Let the sets $\{S_i\}_{i \in [r]}$ be as above. Use Shearer's lemma to show that

$$k \cdot D(Y_1, \dots, Y_n \| X_1, \dots, X_n) \geq \sum_{i \in [r]} D(Y_{S_i} \| X_{S_i}) .$$

- (b) Let $A = \left\{ (a_1, \dots, a_n) \in \{0, 1\}^n \mid \sum_{i \in [r]} f_i(\{a_j\}_{j \in S_i}) \geq t \right\}$. Let (Y_1, \dots, Y_n) be uniformly distributed over the set A (note that Y_1, \dots, Y_n are not necessarily independent). Prove that

$$\mathbb{P}_{X_1, \dots, X_n} \left[\sum_{i \in [r]} f_i(X_{S_i}) \geq t \right] = 2^{-D(Y_1, \dots, Y_n \| X_1, \dots, X_n)} ,$$

where the probability is over the uniform distribution for X_1, \dots, X_n .

- (c) For each $i \in [r]$, let $\mathbb{E}[f_i(X_{S_i})] = \mu_i$ and $\mathbb{E}[f_i(Y_{S_i})] = \nu_i$. Prove that

$$D(Y_{S_i} \| X_{S_i}) \geq D(\nu_i \| \mu_i),$$

where $D(\nu_i \| \mu_i)$ denotes the divergence of two distributions on $\{0, 1\}$ with probabilities $(\nu_i, 1 - \nu_i)$ and $(\mu_i, 1 - \mu_i)$.

- (d) Use the above bounds and the convexity of KL-divergence in both its arguments to show that for $\mu = \frac{1}{r} \cdot (\mu_1 + \dots + \mu_r)$,

$$\mathbb{P}_{X_1, \dots, X_n} [f_1(X_{S_1}) + \dots + f_r(X_{S_r}) \geq (\mu + \varepsilon) \cdot r] \leq 2^{-(r/k) \cdot D(\mu + \varepsilon \| \mu)}.$$