### TOPICS IN PROBABILITY. PART I: CONCENTRATION

#### EXERCISE SHEET 5: ENTROPY AND PROBABILISTIC METHOD

#### 1. Entropy and Variance

**Exercise 1** (Warm-up: Uniform distribution as a maximal entropy distribution). Consider a discrete probability space  $(S, \mathcal{P}(S), \mathbb{P})$  where S is finite. The entropy of the distribution is given by  $-\sum_{s\in S} \mathbb{P}[\{s\}] \log \mathbb{P}[\{s\}]$ . Prove that among all probability distributions on  $\{1,\ldots,n\}$  the uniform distribution has maximal entropy.

*Proof.* Observe that  $x \mapsto x \log x$  is convex. Thus, by discrete Jensen's inequality and since  $\sum_{s \in S} \mathbb{P}(\{s\}) = 1$ , we deduce that

$$\frac{1}{|S|} \sum_{s \in S} \mathbb{P}[\{s\}] \log \mathbb{P}[\{s\}] \ge \frac{1}{|S|} \sum_{s \in S} \mathbb{P}[\{s\}] \log \left(\frac{1}{|S|} \sum_{s \in S} \mathbb{P}[\{s\}]\right) = -\frac{\log |S|}{|S|}$$

such that the entropy is bounded, for whatever probability measure  $\mathbb{P}$  we choose, by

$$-\sum_{s\in S} \mathbb{P}[\{s\}] \log \mathbb{P}[\{s\}] \le \log(|S|).$$

It is easy to see that the uniform distribution actually achieves this upper bound, concluding the proof.  $\Box$ 

Exercise 2 (Variational characterization of variance). Let X be a square-integrable random variable on some probability space. Prove that

$$Var[X] = \sup_{Y} (2Cov(X, Y) - Var[Y]),$$

where supremum is taken over all square-integrable random variables Y on the same probability space.

*Proof.* By setting Y = X, we see that the r.h.s. is greater or equal than Var[X]. On the other hand, using the fact that  $2ab \le a^2 + b^2$ , we get that  $2Cov(X,Y) - Var[Y] \le Var[X]$ .  $\square$ 

**Exercise 3** (Entropy bounds variance). Show that for any non-negative random variable Z, it holds

$$Var[Z] \le Ent[Z^2].$$

Find a counterexample, which shows that the above inequality is not necessarily true if Z is not required to be non-negative.

Hint for the first part: for  $p \in [1,2)$ , consider  $\Psi_p(Z) = \mathbb{E}[Z^2] - (\mathbb{E}[Z^p])^{2/p}$ . Show that  $\lim_{p\uparrow 2} \frac{\Psi_p(Z)}{2-p} = \operatorname{Ent}[Z^2]/2$  and that  $p \mapsto \frac{\Psi_p(Z)}{1/p-1/2}$  is non-decreasing <sup>1</sup>.

<sup>&</sup>lt;sup>1</sup>Extra hint: show for example that  $\alpha(t) = t \log \mathbb{E}[Z^{1/t}]$  for  $t \in (1/2, 1]$  is convex and rewrite the above function in terms of  $\alpha$ .

*Proof.* Let  $p \in [1,2)$  and consider  $\Psi_p(Z) = \mathbb{E}[Z^2] - (\mathbb{E}[Z^p])^{2/p}$ . Note that  $\text{Var}[Z] = \Psi_1(Z)$ . By L'Hôpital's rule,

$$\lim_{p \uparrow 2} \frac{\Psi_p(Z)}{2 - p} = \lim_{p \uparrow 2} 2 \|Z\|_{L^p}^2 \left( \frac{1}{p} \frac{\mathbb{E}[Z^p \log Z]}{\mathbb{E}[Z^p]} - \frac{1}{p^2} \log \mathbb{E}[Z^p] \right)$$
$$= \frac{1}{2} \left( \mathbb{E}[Z^2 \log(Z^2)] - \mathbb{E}[Z^2] \log(\mathbb{E}[Z^2]) \right) = \frac{1}{2} \text{Ent}[Z^2].$$

Now if we show that  $\frac{\Psi_p(Z)}{1/p-1/2}$  as function in p is non-decreasing, then we will be able to conclude the desired result:

$$\operatorname{Var}[Z] = \Psi_1(Z) \le \frac{1}{2} \lim_{p \uparrow 2} \frac{\Psi_p(Z)}{1/p - 1/2} = 2 \lim_{p \uparrow 2} \frac{\Psi_p(Z)}{2 - p} = \operatorname{Ent}[Z^2].$$

So now it is left to show that  $\frac{\Psi_p(Z)}{1/p-1/2}$  as function in p is non-decreasing. To this end, define  $\alpha(t)=t\log\mathbb{E}[Z^{1/t}]$  for  $t\in(1/2,1]$ . Let  $s,t\in(1/2,1], a\in(0,1)$ , then by the generalized Hölder's inequality (take  $\theta_1=a,\theta_2=1-a,f_1=Z,f_2=Z,q=\frac{1}{at+(1-a)s},q_1=1/t,q_2=1/s)$  we get that  $\alpha$  is convex on (1/2,1]. Recall that generalized Hölder's inequality is the following: let  $\theta_1,\theta_2\in(0,1)$  such that  $\theta_1+\theta_2=1, q_1,q_2\in[1,\infty)$  and  $\frac{1}{q}=\frac{\theta_1}{q_1}+\frac{\theta_2}{q_2}$ , then  $\||f_1|^{\theta_1}|f_2|^{\theta_2}\|_q\leq\|f_1\|_{q_1}^{\theta_1}\|f_2\|_{q_2}^{\theta_2}$ . Since  $\alpha$  is convex, then so is  $\beta(t):=e^{2\alpha(t)}=\mathbb{E}[Z^{1/t}]^{2t}$ . Now by the properties of convex functions we know that  $\frac{\beta(t)-\beta(1/2)}{t-1/2}$  is non-decreasing on (1/2,1]. The claim follows since

$$\Psi_p(Z) = \frac{\beta(1/2) - \beta(1/p)}{1/p - 1/2}.$$

## 2. Probabilistic method: Random objects used to prove deterministic statements

**Exercise 4** (Warm-up: Orthogonal projection in expectation). Recall the  $m \times N$ -matrix P from the proof of Johnson-Lindenstrauss lemma, which has i.i.d. standard Gaussian entries. Extend this matrix to an  $N \times N$ -matrix  $\tilde{P}$  by filling the new rows with zeroes. Show that  $\mathbb{E}[\tilde{P}^2] = I_m$ , where  $I_m$  is a diagonal matrix with first m values on the diagonal being 1's and 0 otherwise. Note that  $I_m$  is an orthogonal projection from  $\mathbb{R}^N$  onto m-dimensional subspace  $\mathbb{R}^m$  of  $\mathbb{R}^N$ .

Remark: for a different insight on P being close to a projection read discussion in Section 1.1 in the following notes.

*Proof.* Let  $\tilde{P}$  be as in the exercise, then  $(\tilde{P}^2)_{ij} = (\sum_{k=1}^m X_{ik} X_{kj})$  if  $i \leq m$  and 0 otherwise. Note that  $\mathbb{E}[X_{ik} X_{kj}] = \mathbf{1}_{i=k=j}$ . Therefore, for i, j > m or  $i \neq j$ ,  $\mathbb{E}[(\tilde{P}^2)_{ij}] = 0$ , and for  $i \leq m$ ,  $\mathbb{E}[(\tilde{P}^2)_{ii}] = 1$ .

**Exercise 5** (Erdös theorem, 1959). Given  $k \geq 3$ ,  $g \geq 3$ , there exists a graph with girth at least g and chromatic number  $\chi(G)$  at least k. Recall that the girth of a graph is the length of its shortest circle; and the chromatic number of graph G is the least amount of colors necessary to color the vertices such that no two adjacent vertices share the same color. To prove this theorem proceed as follows:

- Instead of working with  $\chi(G)$  directly (as it is difficult), consider independence number  $\alpha(G)$ , which is the size of the largest set of mutually non-adjacent vertices. Show that  $\chi(G) \geq \frac{|V(G)|}{\alpha(G)}$ .
- Let n be large and consider an Erdös-Renyi graph on these vertices with probability p of edge being present. Choose p small enough such that the expected number of short cycles (of length less than g) is small, but large enough such that the existence of independent sets is unlikely. More precisely, adjust p such that with high probability there are at most n/2 cycles of length less than g and G contains no independent sets of size greater or equal to n/(2k).<sup>2</sup>
- Remove a vertex in each cycle of the graph from Step 2. Argue that the resulting graph with high probability is the desired one.

*Proof.* Suppose that  $\chi(G) = k$ , then a k-coloring of vertices of G gives a disjoint partition  $V(G) = V_1 \cup \ldots \cup V_k$  such that each  $V_i$  is an independent set. Hence,  $|V_i| \leq \alpha(G)$  and  $|V(G)| = \sum_{i \leq k} |V_i| \leq k\alpha(G)$ . We get,  $k = \chi(G) \geq \frac{|V(G)|}{\alpha(G)}$ .

Now let G be the Erdös-Renyi graph with probability parameter p (to be determined) on n vertices (n is large). Let us first adjust p to the requirement on the cycles. Let X be the number of cycles of length less than g. Then  $\mathbb{E}[X] = \sum_{i=2}^{g-1} \sum_{u} \mathbb{E}[\mathbf{1}_{\{u \text{ is a cycle of length } i\}}] = \sum_{i=2}^{g-1} \binom{n}{i} p^i \leq \sum_{i=2}^{g-1} (np)^i$ . We want the probability parameter p to be sufficiently large (so that existence of independent sets is unlikely), but so that  $\sum_{i=2}^{g-1} (np)^i = o(n)$ . A good choice seems to be  $p = n^{-1+\varepsilon}$  with  $\varepsilon \in (0, 1/g)$ . Let us fix  $\varepsilon = 1/(2g)$  and check that this choice satisfies the desired requirements of the second bulletpoint:

- (1) By Markov's inequality,  $\mathbb{P}[X \ge n/2] \le \frac{g\sqrt{n}}{n} \to 0$  as  $n \to \infty$ .
- (2) Let  $r = \lfloor n/(2k) \rfloor$  and  $\overline{G}$  be the complement graph of G, i.e.,  $\overline{G}$  has the same set of vertices and exactly the edges which are not present in G.

 $\mathbb{P}[\alpha(G) \geq r] \leq \mathbb{P}[\exists \text{ a complete subgraph of } \overline{G} \text{ on } r+1 \text{ vertices}]$ 

$$= \binom{n}{r+1} (1-p)^{\binom{r+1}{2}} < (ne^{-pr/2})^{r+1}.$$

As n tends to infinity,  $ne^{-pr/2} = ne^{-n^{1/(2g)}/(4k)} \rightarrow 0$ .

We have shown that for n sufficiently large there exists a graph G with n vertices that with high probability has fewer than n/2 cycles of length smaller than g and with  $\alpha(G) \leq n/(2k)$ . Since removal of vertices from a graph does not increase the independence number (i.e., if G' is obtained from G by removing vertices and edges incident to them,  $\alpha(G') \leq \alpha(G)$ ), deleting a vertex from each short (of length less than g) cycle delivers a graph with girth at least g and (by step 1) chromatic number at least k.

# $\star$ For those who know or would like to learn something about packings and coverings.

**Exercise 6** (Tightness of Johnson-Lindenstrauss). Johnson-Lindenstrauss lemma proven in the lecture shows that for any k points in a Hilbert space H (or  $\mathbb{R}^N$  for simplicity) can be mapped into  $\mathbb{R}^m$  with  $m \geq \log k$  while distorting the distances between them by at most a constant factor. Show that  $m \geq \log k$  is a necessary condition. Assume that the mapping

<sup>&</sup>lt;sup>2</sup>For example,  $p = n^{-1+1/(2g)}$ , Markov inequality might be useful.

considered in the statement of the lemma are exclusively linear maps.

Hint: show that the image of k orthonormal vectors  $x_1, \ldots, x_n$  in H (or  $\mathbb{R}^N$  as in the lecture) under a map  $T: H \to \mathbb{R}^m$  that nearly preserves distances (corresponds to  $\frac{P}{\sqrt{m}}$  in the lecture notes) is a packing of a ball in  $\mathbb{R}^m$ .

Literature on packings and coverings: see Section 5.2 in "Probability in High Dimension", Ramon van Handel or Section 4.2 in "High-Dimensional Probability", Roman Vershynin.

Proof. Let  $e_1, \ldots e_k$  be first k standard vectors in  $\mathbb{R}^N$ . Suppose that JL lemma holds for  $A := \{e_1, \ldots, e_k, 0\}$ . Let T be the linear map which is almost isometric on A - A. We, in particular, get that all  $Te_i$  belong to the ball centered at T0 = 0 (as T linear) of radius at most  $1 + \varepsilon$ , and that  $||Te_i - Te_j|| \ge (1 - \varepsilon) ||e_i - e_j|| = 2(1 - \varepsilon)$ . By definition,  $(Te_i)_i$  is a  $2(1 - \varepsilon)$ -packing of  $B^m(0, 1 + \varepsilon)$  (m-dimensional ball). By proposition 4.2.12 in Vershynin's book, the  $2(1 - \varepsilon)$ -packing number  $P(B^m(0, 1 + \varepsilon), 2(1 - \varepsilon))$  satisfies the following bound:

$$P(B^m(0,1+\varepsilon),2(1-\varepsilon)) \le \frac{|B^m(0,1+\varepsilon) + (1-\varepsilon)B^m(0,1)|}{(1-\varepsilon)B^m(0,1)} \le \left(1 + \frac{1+\varepsilon}{1-\varepsilon}\right)^m.$$

Thus, we should have  $k \leq \left(\frac{2}{1-\varepsilon}\right)^m$ . By taking logarithm we can see that  $m \gtrsim \log k$  is indeed a necessary condition for linear almost isometric T to exist.