

**36-755, Fall 2017**  
**Homework 1**

Due Sep 20.

1. On the MLE in parametric models.

- (a) Recall that the Kullback-Leibler (KL) divergence between two probability measures  $P$  and  $Q$  on some measurable space  $(\mathcal{X}, \mathcal{B})$  with densities  $p$  and  $q$  with respect to a common dominating measure  $\mu$  is

$$K(P, Q) = \begin{cases} \int_{\mathcal{X}} p(x) \log \left( \frac{p(x)}{q(x)} \right) d\mu(x) & \text{if } P \ll Q \\ \infty & \text{otherwise.} \end{cases}$$

Use Jensen inequality to show that  $K(P, Q) \geq 0$  with equality if and only if  $P = Q$ .

- (b) Assume that  $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$  is a parametric model over the sample space  $(\mathcal{X}, \mathcal{B})$ , such that  $P_\theta \ll \mu$  for all  $\theta \in \Theta$ , for some  $\sigma$ -finite dominating measure  $\mu$ . Assume also that all the  $P_\theta$ 's have the same support and  $\theta \neq \theta'$  implies that  $P_\theta \neq P_{\theta'}$ . Let  $\mathbb{X}_n = (X_1, \dots, X_n) \stackrel{id}{\sim} P_{\theta_0}$  for some  $\theta_0 \in \Theta$  and write

$$L_n(\theta; \mathbb{X}_n) = \prod_{i=1}^n p_\theta(X_i),$$

for the likelihood function at  $\theta \in \Theta$ , where  $p_\theta$  is the density of  $P_\theta$  with respect to  $\mu$

Use the law of large numbers to show that, for any  $\theta \neq \theta_0$  in  $\Theta$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}(L_n(\mathbb{X}_n; \theta_0) > L_n(\mathbb{X}_n; \theta)) = 1$$

The previous result offers an asymptotic justification of why in this case the MLE is a sensible choice. *Hint: express the inequality in term of log-likelihood ratio and show that the ratio converges in probability to  $K(P_{\theta_0}, P_\theta)$ . You can use the law of large numbers.*

2. (Reading exercise. **Not to be graded for correctness, but only for effort**)

In this problem you are essentially required to reproduce a proof that can be found in the references given below. My intention is for you to read up and understand the proof rather than trying to solve this problem on your own, which would be challenging (though you are welcome to this challenge). Let  $X = (X_1, \dots, X_d) \in \mathbb{R}^d$  be a random vector with covariance matrix  $\Sigma$  such that  $\frac{X_i}{\sqrt{\Sigma_{i,i}}}$  is sub-Gaussian with parameter  $\nu^2$ , for all  $i = 1, \dots, d$ . Assume we observe  $n$  i.i.d. copies of  $X$  and compute the empirical covariance matrix  $\widehat{\Sigma}$ . Show that, for all  $i, j \in \{1, \dots, d\}$ ,

$$\mathbb{P}\left(\left|\widehat{\Sigma}_{i,j} - \Sigma_{i,j}\right| > \epsilon\right) \leq C_1 e^{-\epsilon^2 n C_2},$$

for some constants  $C_1$  and  $C_2$ . Conclude that

$$\max_{i,j} \left|\widehat{\Sigma}_{i,j} - \Sigma_{i,j}\right| = O_P\left(\sqrt{\frac{\log d}{n}}\right)$$

Thus, estimation of the covariance matrix in the  $L_\infty$  norm is possible even when  $d$  is much larger than  $n$ . Of course, you may ask yourself whether this is a good enough guarantee. In few weeks we will look at consistency rates for covariance estimation under more sensible norms and we will see that the requirements on  $d$  are much more stringent.

You may want to look these references:

- Lemma 12 in Yuan. M. (2010). High Dimensional Inverse Covariance Matrix Estimation via Linear Programming, JMLR, 11, 2261-2286.
- Lemma 1 in Ravikumar, P., Wainwright, M.J., Raskutti, G. and Yu, B. (2011). EJS, 5, 935-980.
- Lemma A.3 in Bickel, P.J. and Levina, E. (2008). Regularized estimation of large covariance matrices, the Annals of Statistics, 36(1), 199-227.

3. (Sampling with replacement). Let  $\mathcal{X}$  a finite set with  $N$  elements. Let  $X_1, \dots, X_n$  be a random sample without replacement from  $\mathcal{X}$  and  $Y_1, \dots, Y_n$  be a random sample with replacement from  $\mathcal{X}$ . Show that, for any convex function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ ,

$$\mathbb{E} \left[ f \left( \sum_{i=1}^n X_i \right) \right] \leq \mathbb{E} \left[ f \left( \sum_{i=1}^n Y_i \right) \right].$$

Use this result to show that all the inequalities derived for the sums of independent random variables  $\{Y_1, \dots, Y_n\}$  using Chernoff's bounding techniques remain true also for the sums of the  $X_i$ 's. (see Hoeffding, W. (1963). *Probability Inequalities for sums of Bounded Random Variables*, by W. Hoeffding, JASA, 58, 13-30., 1963).

Actually, this is a special case of a more general property known as negative dependence. The punch-line is that negatively dependent variables obeys the same Chernoff's bound as independent variables. Perhaps the most prominent example of negatively dependent variables is that of Multinomial variables. For more information see <http://www.brics.dk/RS/96/25/BRICS-RS-96-25.pdf>

4. From tail bounds to moment bounds and high probability bounds.

- (a) Suppose that the random variable  $X$  has mean zero and satisfies the inequality

$$\mathbb{P}(|X| \geq t) \leq c_1 e^{-c_2 n t^a}, \quad \forall t > 0$$

where  $a \in \{1, 2\}$ ,  $n$  is a positive integer and  $c_1$  and  $c_2$  are positive numbers.

- Show that, when  $a = 2$ ,  $\mathbb{V}[X] \leq \frac{c_1}{nc_2}$ .
- Show that

$$\mathbb{E}[|X|] \leq c_3 n^{-1/a}$$

and express  $c_3$  as a function of  $c_1$  and  $c_2$ .

- (b) (From Hoeffding/Bernstein exponential inequality to high probability bounds). Suppose that, for all  $t > 0$ , and some positive constants  $a, b, c$  and a non-negative constant  $d$ ,

$$\mathbb{P}(|X| \geq t) \leq a \exp \left\{ -\frac{nbt^2}{c + dt} \right\}.$$

Then show that, for any  $\delta \in (0, 1)$ ,

$$|X| \leq \sqrt{\frac{c}{nb} \ln \frac{a}{\delta}} + \frac{d}{nb} \ln \frac{a}{\delta},$$

with probability at least  $1 - \delta$ .

5. Let  $X$  be distributed like a  $N_d(0, I_d)$ , where  $I_d$  is the  $d$ -dimensional identity matrix. Then,  $\|X\|^2 = \sum_{i=1}^d X_i^2 \sim \chi_d^2$ .

- (a) Show that, for any  $\epsilon \in (0, 1)$

$$\mathbb{P}(|\|X\|^2 - d| \geq d\epsilon) \leq 2e^{-d\epsilon^2/8}.$$

You can use the following fact: the moment generating function of a  $\chi_d^2$  is  $(1 - 2\lambda)^{-d/2}$  for all  $\lambda < 1/2$ . Alternatively, use the version of Bernstein inequality for sum of sub-exponential variables given in class. This results says that, in high dimensions,  $X$  is concentrated around a sphere of radius  $\sqrt{d}$ .

See, e.g., Lemma 2 in *A Probabilistic Analysis of EM for Mixtures of Separated, Spherical Gaussians*, by S. Dasgupta and L. Schulman, *JMLR*, 8, 203–26, 2007.

You should convince yourself that the same result holds for any vector  $X$  whose entries are i.i.d. sub-Gaussians.

- (b) Now assume that  $X$  and  $Y$  are both  $\in N_d(0, I_d)$  and are independent. Argue **very informally (it is OK to use heuristics)** that

$$\frac{|X^\top Y|}{\|X\|\|Y\|} \sim \frac{1}{\sqrt{d}},$$

with high probability. Thus conclude that in high-dimensions, independent isotropic Gaussian vectors are orthogonal with high probability, the more so the higher the dimension.

You may use the fact that if  $X \sim N_d(0, I_n)$ , then  $\frac{X}{\|X\|}$  and  $\|X\|$  are independent.

Again, the assumption of Gaussianity can be replaced by that of sub-Gaussianity.

6. Suppose that  $X_1, \dots, X_n$  are such that  $X_i \in SG(\sigma_i^2)$ , not necessarily independent. Show that  $\sum_{i=1}^n X_i \in SG(\tau^2)$  and find  $\tau$ . What if  $X_i \in SE(\tau_i^2, \alpha_i)$  for all  $i$ ?
7. (Random Projection and the Johnson-Lindenstrauss Lemma).

See D. Achlioptas, *Database friendly random projections: Johnson-Lindenstrauss with binary coins*, *Journal of Computer and System Sciences* 66 (2003) 671687.

Suppose we have a (deterministic) vector  $x$  in  $\mathbb{R}^D$  and, for  $\epsilon \in (0, 1/2)$  we would like to find a random mapping  $f: \mathbb{R}^D \rightarrow \mathbb{R}^d$ , where  $d$  is smaller than  $D$ , such that

$$(1 - \epsilon)\|f(x)\|^2 \leq \|x\|^2 \leq (1 + \epsilon)\|f(x)\|^2$$

with high probability. One way is to construct a  $d \times D$  matrix  $A$  with iid entries from the  $N(0, 1)$  distribution and then take

$$f(x) = \frac{1}{\sqrt{d}}Ax, \quad x \in \mathbb{R}^D.$$

You can think of  $f$  as being a random projection from a high-dimensional space  $\mathbb{R}^D$  into the smaller space  $\mathbb{R}^d$ .

Show that

$$(a) \quad \|x\|^2 = \mathbb{E}[\|f(x)\|^2].$$

- (b) For each  $\epsilon \in (0, 1/2)$

$$\mathbb{P}(|\|f(x)\|^2 - \|x\|^2| > \epsilon\|x\|^2) < 2 \exp\{-d/4(\epsilon^2 - \epsilon^3)\}.$$

- (c) Using the above result, show that, if we are given  $n$  deterministic vectors  $(x_1, \dots, x_n)$  in  $\mathbb{R}^D$  and we compute their projections  $f(x_1), \dots, f(x_n)$  in  $\mathbb{R}^d$ , we are guaranteed that the all the pairwise squared distances between the projected points are distorted by at most a factor of  $\epsilon \in (0, 1/2)$  with probability at least  $1 - \delta$  if  $d \geq \frac{4(\log(1/\delta) + 2 \log(n))}{\epsilon^2 - \epsilon^3}$ . That is,

$$\|x_i - x_j\|^2(1 - \epsilon) \leq \|f(x_i) - f(x_j)\|^2 \leq \|x_i - x_j\|^2(1 + \epsilon), \quad \forall i \neq j,$$

with probability at least  $1 - \delta$ .

For parts (a) and (b) proceed as follows: show that the squared norm of  $\frac{\sqrt{d}f(x)}{\|x\|}$  is equal in distribution to the sum of  $d$  squared standard normals, and therefore has a  $\chi_d^2$  distribution. In your subsequent derivation, you may use the following facts:

- (a) The mfg of a  $\chi_1^2$  at any  $\lambda < 1/2$  is  $(1 - 2\lambda)^{-1/2}$ .  
(b) For any  $\epsilon \in (0, 1/2)$ , setting  $\lambda = \frac{\epsilon}{2(1+\epsilon)} < 1/2$ , we get

$$\frac{e^{-2(1+\epsilon)\lambda}}{1 - 2\lambda} = (1 + \epsilon)e^{-\epsilon} < e^{-1/2(\epsilon^2 - \epsilon^3)}$$

and setting  $\lambda = \frac{\epsilon}{2(1-\epsilon)} < 1/2$  we get

$$\frac{e^{2(1-\epsilon)\lambda}}{1 + 2\lambda} = (1 - \epsilon)e^{\epsilon} < e^{-1/2(\epsilon^2 - \epsilon^3)}$$

What is striking about this result is that the dimension  $D$  of the original space does not appear anywhere in these bounds!

This is an instance of what is also known as the Johnson-Lindenstrauss Lemma, which loosely speaking, states that a random projection of  $n$  points from a high-dimensional space into a  $d$  dimensional space preserves the pairwise squared distances up to a multiplicative factor of  $\epsilon$  with high probability if  $d$  is of order  $\frac{\log n}{\epsilon^2}$ , independently of the dimension of the original space.

Notice that instead of using independent  $N(0, 1)$  variables to populate  $A$ , we could have used any sub-Gaussian distribution.

8. Show that if  $X \in SG(\sigma^2)$  than  $X^2 \in SE(\nu^2, \alpha)$  where

$$\nu = \alpha = 16\sigma^2.$$

*Hint: For this problem you may find it helpful to use the following facts:*

- (a) **The  $C_r$  inequality:** If  $X$  and  $Y$  are random variables such that  $\mathbb{E}|X|^r < \infty$  and  $\mathbb{E}|Y|^r < \infty$ , where  $r \geq 1$ , then

$$\mathbb{E}|X + Y|^r \leq 2^{r-1} (\mathbb{E}|X|^r + \mathbb{E}|Y|^r)$$

- (b) The following bound, proved in class

$$\mathbb{E}|X|^r \leq (2\sigma^2)^{r/2} r \Gamma(r/2) \quad r \geq 1.$$

Feel free to prove the claim by other methods and/or by obtaining sharper (smaller) bounds on  $\nu^2$  and/or  $1/\alpha$ .