

ISyE 8803: Special Topics in Modern Mathematical Data Science
Homework 1

due on Friday, Feb 7th at 11:59 pm

Please submit electronically directly to Canvas in a PDF file.

1 MGF method vs. moment bounds

It is natural to compare the best bound on the tails obtained via MGF and by bounding the moments. As it turns out, the moment bounds are sharper, even if we only use the integer moments.

(a) Show that if $X > 0$ a.s., then for any $u > 0$,

$$\inf_{\lambda > 0} M_X(\lambda) e^{-\lambda u} \geq \inf_{k \in \mathbb{Z}_+} \mathbb{E} \left[X^k \right] u^{-k}.$$

(b) Show that if X is symmetric (i.e. X and $-X$ have the same distribution), then for any $u > 0$,

$$\inf_{\lambda > 0} M_X(\lambda) e^{-\lambda u} \geq \frac{1}{2} \inf_{k \in \mathbb{Z}_+} \mathbb{E} \left[X^{2k} \right] u^{-2k}.$$

2 Convexity of the cumulant-generating function

For any distribution X , the logarithm of the MGF

$$K_X(t) = \log \mathbb{E}[e^{tX}]$$

is called the cumulant-generating function, or the *log-partition function* of the distribution.

- (a) Show that K_X is convex. Use Young's inequality: for $a, b \in \mathbb{R}^d$ and $p, q \in [1, \infty]$ with $\frac{1}{p} + \frac{1}{q} = 1$,

$$|a^\top b| \leq \|a\|_p \|b\|_q.$$

You can assume that X has a discrete distribution.

3 Gaussian tails

3.1 Mills ratio

Let $\phi(\cdot)$ be the p.d.f. of $\mathcal{N}(0, 1)$, i.e. $\phi(t) = \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}}$. For any $u \geq 0$, let $\Phi(u) := \int_{t \geq u} \phi(t) dt$.

(a) Prove the following bounds (holding for all $u \geq 0$):

$$\left(\frac{1}{u} - \frac{1}{u^3} \right) \phi(u) \leq \Phi(u) \leq \frac{1}{u} \phi(u).$$

Hint 1: Try to prove the upper bound first.

Hint 2: Integrate by parts – first to prove the upper bound, then again for the lower bound.

(b) Capitalizing on the trick you have just figured out to get the lower bound from the upper bound, prove a new upper bound:

$$\Phi(u) \leq \left(\frac{1}{u} - \frac{1}{u^3} + \frac{3}{u^5} \right) \phi(u).$$

Note that this bound is sharper than the previous one for large enough u .

*(c) If we continue this approach, we obtain a power series in $1/u$ for the Mills ratio $\Phi(u)/\phi(u)$; see Theorem 2.1 from Lecture 2. Get yourself convinced in it (no need to prove).

3.2 Power series for c.d.f.

Show that

$$\frac{1}{2} - \Phi(u) = \frac{1}{\sqrt{2\pi}} \sum_{k=0}^{\infty} \frac{(-1)^k u^{2k+1}}{2^k k! (2k+1)}.$$

Hint: change variable to remove u from the integration limits; differentiate in u under the integral.

4 Paley-Zygmund and friends

(i) Prove the Paley-Zygmund inequality (it can be interpreted as a counterpart of Markov: a nonnegative random variable cannot be much *smaller* than its expectation):

If X is a non-negative random variable with $\mathbb{E}[X^2] < \infty$, then for any $t \in [0, 1]$ one has

$$\mathbb{P}(X \geq (1-t)\mathbb{E}X) \geq t^2 \frac{(\mathbb{E}X)^2}{\mathbb{E}[X^2]}. \quad (1)$$

(ii) Under the same assumptions, strengthen (1), to Cantelli's inequality:

$$\mathbb{P}(X \geq (1-t)\mathbb{E}X) \geq t^2 \frac{(\mathbb{E}X)^2}{t^2(\mathbb{E}X)^2 + \text{Var}[X]}.$$

This new inequality is sharp – give an example where it is attained.

(iii) Now: instead of $\mathbb{E}[X^2] < \infty$, assume that $\mathbb{E}[|X|^p] < \infty$ for some $p > 1$, and generalize (1) to

$$\mathbb{P}(X \geq (1-t)\mathbb{E}X) \geq \left(t^p \frac{(\mathbb{E}X)^p}{\mathbb{E}[|X|^p]} \right)^{\frac{1}{p-1}}.$$

Note that when $p > 2$, this gives an improvement over (1) for small t , which is important in applications where X is itself the sample average of i.i.d. Y_1, \dots, Y_n .

Hint: use Hölder's inequality: given $p, q \geq 1$ such that $1/p + 1/q = 1$, and random variables U, V on the same sample space, one has $\mathbb{E}[|UV|] \leq (\mathbb{E}|U|^p)^{\frac{1}{p}} \cdot (\mathbb{E}|V|^q)^{\frac{1}{q}}$.

5 Tail bound for χ_d^2

Let $X \sim \chi_{2d}^2$ (chi-squared distribution with $2d$ degrees of freedom), that is $X = \|Z\|^2 = Z_1^2 + \dots + Z_{2d}^2$ where $Z \sim \mathcal{N}(0, \mathbf{I}_d)$ (equivalently, $Z_i \sim \mathcal{N}(0, 1)$ are i.i.d.). Define $M_{2d}(\cdot)$ as the MGF of $X \sim \chi_{2d}^2$,

$$M_{2d}(t) := \mathbb{E}[e^{tX}], \quad t \in \mathbb{R};$$

in particular, $M_2(t) = \mathbb{E}[e^{t(Z_1^2 + Z_2^2)}]$. Our ultimate goal here is to prove that, with probability $\geq 1 - \delta$,

$$X - 2d \leq \sqrt{Cd \log\left(\frac{1}{\delta}\right)} + c \log\left(\frac{1}{\delta}\right) \quad (2)$$

for some numerical constants $C, c > 0$. This bound is, in fact, optimal (see, e.g., [LM00, Lemma 1]).

(a) Derive the explicit form of $M_2(t)$:

$$M_2(t) = \frac{1}{1 - 2t}, \quad t < \frac{1}{2},$$

and $M_2 = +\infty$ for $t \geq \frac{1}{2}$. (To take the integral, pass to polar coordinates $(z_1, z_2) \mapsto (r, \theta)$ with $r = \sqrt{z_1^2 + z_2^2}$ —and don't forget the Jacobian, which equals r .) Claim that, as a corollary,

$$M_{2d}(t) = \frac{1}{(1 - 2t)^d}, \quad t < \frac{1}{2}.$$

(b) Using Chernoff's method, bound the tail function $\mathbb{P}(X > x)$, for any $x > 2d$, as follows:

$$\mathbb{P}(X > x) = \inf_{t < \frac{1}{2}} \frac{e^{-tx}}{(1 - 2t)^d} = \exp\left(d \log\left(\frac{x}{2d}\right) - \frac{x - 2d}{2}\right).$$

(Hint: it is convenient to take the logarithm, and use that $u \mapsto \log(u)$ on \mathbb{R}_+ is increasing.)

Note that, in terms of the deviation $z = x - 2d > 0$ above $2d$, this is equivalent to

$$\mathbb{P}(X - 2d > z) = \exp\left(d \log\left(\frac{2d + z}{2d}\right) - \frac{z}{2}\right).$$

**(c) Bonus. Bear with me, this part is a bit delicate – but we need it to reach the conclusion.*

(c.i) Show that

$$\mathbb{P}(X - 2d > z) \leq \begin{cases} \exp\left(-\frac{z^2}{16d}\right) & \text{for } 0 \leq z \leq 2d, \\ \exp\left(-\frac{z}{8}\right) & \text{for } z > 2d. \end{cases}$$

It is OK if you get some worse pair of constants $C > 16, c > 8$ leading to a weaker bound.

Hint: first show, using calculus, that

$$\log(1 + u) \leq u - \frac{1}{4} \min\{u, u^2\} \quad \forall u \geq 0$$

(c.ii) Reformulating the last bound as

$$\mathbb{P}(X - 2d > z) \leq \exp\left(-\min\left\{\frac{z^2}{16d}, \frac{z}{8}\right\}\right)$$

and letting $\mathbb{P}(X - 2d > z) = \delta$, “invert” the last inequality to get (2) with $C = 16$ and $c = 8$ (or with some worse constants). Hint: $\max\{a, b\} \leq a + b$ for $a, b \geq 0$.

6 Stein's paradox

Consider the problem of estimating the mean μ in the multivariate Gaussian location family

$$\mathbb{P}_\mu = \mathcal{N}(\mu, \mathbf{I}_d), \quad \mu \in \mathbb{R}^d, \quad (3)$$

where \mathbf{I}_d is the $d \times d$ identity matrix, from a single observation $X \sim \mathbb{P}_\mu$. Note that here, X itself is the maximum likelihood estimator (MLE) for μ . Defining for any estimator $\hat{\mu} = \hat{\mu}(X)$ of μ the variance

$$\text{Var}_\mu[\hat{\mu}] := \mathbb{E}_\mu[\|\hat{\mu} - \mathbb{E}[\hat{\mu}]\|^2]$$

and the quadratic risk

$$\text{Risk}_\mu[\hat{\mu}] := \mathbb{E}_\mu[\|\hat{\mu} - \mu\|^2],$$

where $\|x\| := (\sum_i x_i^2)^{1/2}$ is the Euclidean norm of $x = (x_1, \dots, x_d) \in \mathbb{R}^d$, we see that for any $\mu \in \mathbb{R}^d$,

$$\text{Risk}_\mu[X] = \text{Var}_\mu[X] = d.$$

Intuitively, one can suspect that no better estimator of X can be found: really, what can be done with only a single observation of the mean? Yet, this turns out to be false: one may improve over the MLE uniformly on the family (3) when $d > 2$. This celebrated result was established by James and Stein in 1976, and our goal is to reproduce it. But first, let us establish the terminology.

Definition 1. An estimator $\hat{\mu}$ is *dominated* by some other estimator $\hat{\mu}'$ if $\text{Risk}_\mu[\hat{\mu}'] \leq \text{Risk}_\mu[\hat{\mu}]$ for any μ , and there exists a parameter value $\bar{\mu}$ such that $\text{Risk}_{\bar{\mu}}[\hat{\mu}'] < \text{Risk}_{\bar{\mu}}[\hat{\mu}]$.

Definition 2. An estimator $\hat{\mu}$ is called *admissible* if it is not dominated by any other estimator. Otherwise, it is called *inadmissible*.

As statisticians, ideally we would like to compare two estimators over the whole family at once, without specifying a value of μ . Two admissible estimators cannot be compared this way, but at the very least we can rule out any *inadmissible* estimator, as for it there exists a uniformly better one.

You will show that the MLE is inadmissible when $d \geq 3$, by constructing a dominating estimator.

- (a) Consider *shrinkage estimators* $\hat{\mu} = sX$ with $s \in \mathbb{R}$, and compute their risks for any s . Show that one can restrict attention to $s \in [0, 1]$ (hence “shrinkage”) by finding a dominating estimator for $\hat{\mu}$ with $s < 0$ or $s > 1$.
- (b) Show that, for given μ , the best value of s —i.e., the one minimizing the risk—is given by

$$s^* = \frac{\|\mu\|^2}{d + \|\mu\|^2} = 1 - \frac{d}{d + \|\mu\|^2}.$$

- (c) Unfortunately, $\hat{\mu}^* = s^*X$ is not a proper estimator. (*Why?*) Instead of it, one may consider

$$\left(1 - \frac{d}{\|X\|^2}\right) X,$$

which is an actual estimator. Can you explain the heuristic motivation behind this estimator?

*(d) Assuming that $d \geq 2$, derive the *James-Stein estimator*

$$\hat{\mu}^{JS} = \left(1 - \frac{d-2}{\|X\|^2}\right) X \quad (4)$$

by minimizing over $\delta \in \mathbb{R}$ the risk of the estimator

$$\hat{\mu}^\delta = \left(1 - \frac{\delta}{\|X\|^2}\right) X$$

for a fixed μ . In order to show that $R(\delta) = \text{Risk}_\mu[\hat{\mu}^\delta]$ is minimized at $d-2$, use Stein's lemma:

Lemma 1. *Let $X \sim \mathcal{N}(\mu, I)$ and $g(x)$ be a function on \mathbb{R}^d differentiable almost everywhere, and such that $\mathbb{E}_\mu \left[\left| \frac{\partial}{\partial x_i} g(X) \right| \right] < \infty$ and $\mathbb{E}_\mu [(X_i - \mu_i)g(X)] < \infty$ for any $i \in [d] := \{1, 2, \dots, d\}$. Then*

$$\mathbb{E}_\mu [(X_i - \mu_i)g(X)] = \mathbb{E}_\mu \left[\frac{\partial}{\partial x_i} g(X) \right], \quad i \in [d].$$

When applying Stein's lemma to the right function $g(X)$, please do check the absolute integrability conditions in its premise, and explain why the argument does not work for $d = 1$. Finally, verify that $R(\delta)$ is strictly convex when $d \geq 3$ (thus $\hat{\mu}^{JS}$ indeed dominates the MLE).

7 Planar Venn diagrams

A (congruent) *Venn diagram* in \mathbb{R}^d for n sets is the following object: you choose a “base” set $A \subset \mathbb{R}^d$ and n locations $a_1, \dots, a_n \in \mathbb{R}^d$ such that the shifted sets A_1, A_2, \dots, A_n , where $A_j := \{a + a_j, a \in A\}$, intersect in all possible combinations: for any subset of indices $I \subseteq \{1, 2, \dots, n\}$, the set $A_I := \cap_{i \in I} A_i$ must be nonempty. Prove the following result:

One cannot draw a planar ($d = 2$) Venn diagram for $n \geq 5$ sets by shifting a circle.

Use **Euler’s formula**: any planar graph with V vertices, E edges, and F faces (subsets in which \mathbb{R}^2 is partitioned by the graph) satisfies

$$V - E + F = 2.$$

For example, in the case of a triangle $V = E = 3$ and $F = 2$.

*Hint: estimate V_n, E_n, F_n in a Venn diagram for n sets in terms of $V_{n-1}, E_{n-1}, F_{n-1}$ respectively.*¹

¹In fact, $n = 4$ is also impossible, but I am not aware of a purely combinatorial (and elegant) proof.

References

- [LM00] *B. Laurent and P. Massart. Adaptive estimation of a quadratic functional by model selection.*
The Annals of Statistics, 28(5):1302–1338, 2000.