

# GML Autumn 23 - HW 1: Concentration bounds

Instructions:

- Please upload to gradescope **by Thursday, 12.10.23 by 23:59**
- Typeset (Latex, Markdown etc.), do not add your name on homeworks and midterms to avoid biases when grading, start the answer to each question on a new page
- See [website](#) for details regarding collaboration and honor code
- MW refers to Martin Wainwright's book
- You will receive partial points if you attempt a question and zero if you don't
- Optional means that you will not receive points for the grade for solving this assignment.
- Please de-register if you do not want to solve the homework
- if you solve all questions that are not marked as bonus or optional, you can get full points. However, if you solve a question marked as bonus you get additional points if you lose some in the other questions

## 1 Optional Warm-up: Optimality of polynomial Markov

Chernoff's bound is obtained via Markov's inequality. In this question we show that Markov's inequality is actually tight. Furthermore, the  $k$ -th moment Markov bounds are in fact never worse than the Chernoff bound based on the moment generating function.

- a) **Find a** non-negative random variable  $X$  for which Markov's inequality is met with equality at a point  $a > 0$ .
- b) Suppose that  $X \geq 0$  and that  $\mathbb{E}e^{\lambda X}$  exists in an interval around zero. Given some  $\delta > 0$  and integer  $k = 1, 2, \dots$  **show that**

$$\inf_{k=0,1,\dots} \frac{\mathbb{E}|X|^k}{\delta^k} \leq \inf_{\lambda>0} \frac{\mathbb{E}e^{\lambda X}}{e^{\lambda\delta}}$$

## 2 Concentration and kernel density estimation

Let  $\{X_i\}_{i=1}^n$  be an i.i.d. sequence of random variables drawn from a density  $f$  on the real line. A standard estimate of  $f$  is the kernel density estimate:

$$f_n(x) := \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right),$$

where  $K : \mathbb{R} \rightarrow [0, \infty)$  is a kernel function satisfying  $\int_{-\infty}^{\infty} K(t) dt = 1$ , and  $h > 0$  is a bandwidth parameter.

Suppose that we assess the quality of  $f_n$  using the  $L_1$ -norm:

$$\|f_n - f\|_1 := \int_{-\infty}^{\infty} |f_n(t) - f(t)| dt.$$

**Prove that:**

$$P[\|f_n - f\|_1 \geq \mathbb{E}[\|f_n - f\|_1] + \delta] \leq e^{-\frac{n\delta^2}{18}}.$$

### 3 Sub-Gaussian maxima

In this exercise we prove an inequality used repeatedly in later lectures.

Let  $\{X_i\}_{i=1}^n$  be a sequence of zero-mean random variables, each subgaussian with parameter  $\sigma$ . The random variables  $X_i$  are *not* assumed to be independent.

a) **Prove that** for all  $n \geq 1$  we have

$$\mathbb{E} \max_{i=1, \dots, n} X_i \leq \sqrt{2\sigma^2 \log n}.$$

*Hint:* the exponential is a convex function.

b) **Prove that** for all  $n \geq 2$  we have

$$\mathbb{E} \max_{i=1, \dots, n} |X_i| \leq \sqrt{2\sigma^2 \log(2n)} \leq 2\sqrt{\sigma^2 \log n}.$$

### 4 Bonus: Sharper tail bounds for bounded variables: Bennett's inequality

**Read MW Chapter 2** and learn about subexponential tail bounds and Bernstein's inequality, yielding some more tail bounds for empirical means of random variables satisfying conditions other than the subgaussian one. Bernstein's inequality is sometimes tighter for bounded variables than when applying the subgaussian bound. In this problem we prove an even tighter bound for bounded variables, known as Bennett's inequality

a) Consider a zero-mean random variable such that  $|X_i| \leq b$  for some  $b > 0$ . **Prove that**

$$\log \mathbb{E} e^{\lambda X_i} \leq \sigma_i^2 \lambda^2 \frac{e^{\lambda b} - 1 - \lambda b}{(\lambda b)^2}$$

for all  $\lambda \geq 0$ , where  $\sigma_i^2 = \text{Var}(X_i)$ .

b) Given independent random variables  $X_1, \dots, X_n$  satisfying the condition of part (a), let  $\sigma^2 := \frac{1}{n} \sum_{i=1}^n \sigma_i^2$  be the average variance. **Prove Bennett's inequality**

$$\mathbb{P} \left( \frac{1}{n} \sum_{i=1}^n X_i \geq \delta \right) \leq e^{-\frac{n\sigma^2}{b^2} h\left(\frac{b\delta}{\sigma^2}\right)}$$

where  $h(t) := (1+t) \log(1+t) - t$  for  $t \geq 0$ .

c) **Show that** Bennett's inequality is at least as good as Bernstein's inequality.

### 5 Sharp upper bounds on binomial tails

Let  $\{X_i\}_{i=1}^n$  be an i.i.d. sequence of Bernoulli variables with parameter  $\alpha \in (0, \frac{1}{2}]$ , and consider the binomial random variable  $Z_n = \sum_{i=1}^n X_i$ . The goal of this exercise is to prove, for any  $\delta \in (0, \alpha)$ , a sharp upper bound on the tail probability  $P[Z_n \leq \delta n]$ .

a) **Show that**

$$P[Z_n \leq \delta n] \leq e^{-nD(\delta \parallel \alpha)},$$

where the quantity

$$D(\delta \parallel \alpha) := \delta \log \frac{\delta}{\alpha} + (1 - \delta) \log \frac{(1 - \delta)}{(1 - \alpha)}$$

is the Kullback–Leibler divergence between the Bernoulli distributions with parameters  $\delta$  and  $\alpha$ , respectively.

b) **Show that** the bound from part (a) is strictly better than the Hoeffding bound for all  $\delta \in (0, \alpha)$ .

## 6 Robust estimation of the mean

Suppose we want to estimate the mean  $\mu$  of a 1-dimensional random variable  $X$  with variance  $\sigma^2$  from a sample  $X_1, \dots, X_n$ , drawn independently from the distribution of  $X$ . We want an  $\epsilon$ -accurate estimate of the mean, i.e., one that falls with probability  $\geq 1 - \delta$  in the interval  $[\mu - \epsilon, \mu + \epsilon]$ .

**Show that** a sample size of  $N = O\left(\log(\delta^{-1}) \frac{\sigma^2}{\epsilon^2}\right)$  suffices to compute an  $\epsilon$ -accurate estimate of the mean with probability at least  $1 - \delta$ . *Hint: Compute the median of  $\log(\delta^{-1})$  weak estimates.*

## 7 Best-arm identification

We now look at an interesting application of concentration bounds. Assume that we have  $K$  newly developed drugs to cure a disease and denote with  $\mu_k \in [0, 1]$  the probability of getting cured by the  $k$ -th drug, which is assumed to be unknown. In order to determine the best drug  $k^*$  with the highest chance of a successful treatment  $\mu^* = \mu_{k^*} = \max_k \mu_k$ , we treat different volunteers in a clinical trial with one drug each and record the outcome. We model the observation of the outcome on one patient as sampling from a Bernoulli distribution with parameter  $\mu_k$ . We denote with  $X_{k,i} \in \{0, 1\}$  the random variable indicating whether the  $i$ -th volunteer treated with the  $k$ -th drug was successful.

In a randomized control trial, all drugs would have the same probability of getting assigned to any patient throughout the trial. In this exercise, we want to study an adaptive algorithm that assigns treatment depending on the outcome of previous treatments. The goal is to assign the drugs in a way such that for some  $\delta \in (0, 1)$ , with probability  $\geq 1 - \delta$ , the algorithm finds the best drug  $k^*$  in *as few volunteers as possible*. This is ethically more reasonable than assigning a “bad” drug to patients even when their results are clearly inferior to others in the trial.

*Context:* This problem is often referred to as a best-arm identification problem. In adaptive or online learning scenarios, where at each time step we sample from one of  $k$  distributions  $\{\mathbb{P}_1, \dots, \mathbb{P}_K\}$  is often called a multi-armed bandit. Pulling an arm  $k$  then corresponds to sampling from  $\mathbb{P}_k$ . In our case they are Bernoulli distributions with means  $\{\mu_1, \dots, \mu_K\}$ .

In this exercise, we analyze a specific type of algorithm to solve the problem called the Successive Elimination algorithm.

---

**Algorithm 1:** Successive Elimination

---

$S_0 = \{1, \dots, K\}$  ;

**for**  $1 \leq t \leq \infty$  **do**

    Pull all arms in  $S_{t-1}$  to obtain samples  $X_{k,t} \sim \mathcal{D}_k$  with  $k \in S_{t-1}$ ;

    Update  $S_t = S_{t-1} - \{i \in S_{t-1} : \exists k \in S_{t-1} : \hat{\mu}_{k,t} - U(t, \delta/K) > \hat{\mu}_{i,t} + U(t, \delta/K)\}$ ;

    Stop when  $|S_t| = 1$ ;

**end**

---

Notation:

- $S_t$ : The active set of arms.
- $\hat{\mu}_{k,t} := \frac{1}{t} \sum_{i=1}^t X_{k,i}$ : Estimated mean of the reward  $\mu_k$  for arm  $k$  after  $t$  pulls.
- $U(t, \delta)$ : An **any-time confidence interval**, such that for any arm  $k$ ,

$$\mathbb{P} \left( \bigcup_{t=1}^{\infty} \{|\hat{\mu}_{k,t} - \mu_k| \geq U(t, \delta)\} \right) \leq \delta.$$

The goal of this exercise is to prove Theorem 1 where we show that the Successive Elimination algorithm is correct and derive an upper bound on the maximum amount of steps needed to for the algorithm to terminate.

**Theorem 1.** *With probability  $\geq 1 - \delta$ :*

1. *For any  $t \geq 1$ , the best arm  $k^*$  is contained in the set  $S_t$ .*
2. *There exists an any-time confidence interval  $U$  such that the Successive Elimination algorithm terminates after  $O(\sum_{k \neq k^*}^K \Delta_k^{-2} \log(K \Delta_k^{-1}))$  samples with  $\Delta_k := \mu^* - \mu_k$  and the  $O$  notation is with respect to  $K$  and  $\Delta_k$  for a constant  $\delta$ .*

We first prove that with high probability the best arm stays in the active set  $S_t$  for all  $t$  until termination.

- a) Define  $\mathcal{E}$  as the event that for any  $t \geq 1$ , the estimated reward  $\hat{\mu}_{k,t}$  of any arm  $k$  is not contained in the confidence interval  $U(t, \delta/K)$  around the true mean  $\mu_k$ , i.e.

$$\mathcal{E} := \bigcup_{k=1}^K \bigcup_{t=1}^{\infty} \{|\hat{\mu}_{k,t} - \mu_k| > U(t, \delta/K)\}.$$

**Show that  $\mathbb{P}(\mathcal{E}) \leq \delta$ .**

- b) **Prove** statement 1 in Theorem 1.

It is not yet shown whether and after how many steps the algorithm terminates. To do so, we derive a sufficiently tight any-time confidence interval  $U$  based on the concentration inequalities discussed in the lecture.

- c) Let  $\{Z_t\}_{t=1}^{\infty}$  be i.i.d bounded random variables with  $Z_t \in [a, b]$  with  $a \leq b$ . **Show that**

$$U = \sqrt{\frac{(b-a)^2 \log(4t^2/\delta)}{2t}}$$

is a valid any-time confidence interval for the random variable  $Z_t$ . *Hint:* Use Hoeffding's bound and union bound.

- d) *Bonus:* **Prove** statement 2 in Theorem 1.