

Lecture 3

Last time

- ▷ Chernoff Method
- ▷ Sub-Gaussians
- ▷ Hoeffding's inequality

Today

- ▷ Robust mean estimation
- ▷ Chernoff's inequality
- ▷ Degree of random graphs.

Robust Mean Estimation

Using Hoeffding's inequality we can obtain confidence intervals for mean estimation of sub-Gaussians

If X_1, \dots, X_n iid σ^2 -sub-Gaussians, then

$\hat{x} = \frac{1}{n} \sum_{i=1}^n x_i$ satisfies that $\forall p > 0$

$$P(|\hat{x} - \mathbb{E}x_1| \leq \sqrt{2\sigma^2 \log(1/p)}) \geq 1-p.$$

A natural question is whether one can derive a similar guarantee for some estimator \hat{x} if the x_i are not sub-Gaussian?

The answer is yes! Try to think how before reading the following.

The intuition is simple: one can bootstrap an estimator by ignoring outliers.

Theorem: Consider X_1, \dots, X_n iid with $\mathbb{E}X_i = \mu$ and $\text{Var}(X) = \sigma^2$. Fix $p \in [0, 1]$. Split the samples into $K = 18 \log(1/p)$ (equally sized) bins and let the empirical means be \hat{x}_j for $j=1, \dots, K$. Then, $\hat{X} = \text{median}(\hat{x}_1, \dots, \hat{x}_K)$ satisfies

$$P(|\hat{X} - \mu| \leq \sqrt{\frac{3\sigma^2 \log(1/p)}{n}}) \geq 1-p.$$

Recall that $\#\{i \mid \hat{x}_i \leq \text{median}\} \geq \frac{1}{2}$,
 $\#\{i \mid \hat{x}_i \geq \text{median}\} \geq \frac{1}{2}$.

Remarks:

1. We only require second moments and so the MGF might not exist.
2. The estimator depends on p . This is unavoidable (See Exercise 2.17 of Vershynin's).

Proof: First the main idea: suppose that we have a $t > 0$ s.t.



More than $\frac{1}{2}$ of the points \hat{x}_i here

Then, $|\text{median}(\hat{x}_1, \dots, \hat{x}_k) - \mu| \leq t$. (why?)
 Consider the events

$$E_i = \left\{ |\hat{x}_i - \mu| \leq \sqrt{\frac{3\sigma^2 k}{n}} \right\}.$$

If more than half of these events hold, then, we get $|\bar{x} - \mu| \leq \sqrt{\frac{3\sigma^2 k}{n}}$. This equivalent to

$$\frac{1}{k} \sum \mathbb{1}_{E_i} > \frac{1}{2} \quad (\star)$$

By Chebychev

$$P(E_i^c) \leq \frac{\text{Var}(x_i)}{\frac{3\sigma^2 k}{n}} = \frac{\frac{\sigma^2}{(n/k)}}{\frac{3\sigma^2 k}{n}} = \frac{1}{3}.$$

So $E \mathbb{1}_{E_i} > \frac{2}{3}$, and if (\star) does not hold

$$\begin{aligned} \frac{1}{k} \sum \mathbb{1}_{E_i} &\leq \frac{1}{2} \Rightarrow \frac{1}{k} \sum \mathbb{1}_{E_i} \leq \frac{2}{3} - \frac{1}{6} \\ &\leq E \mathbb{1}_{E_i} - \frac{1}{6} \\ \Rightarrow \left| \frac{1}{k} \sum \mathbb{1}_{E_i} - E \mathbb{1}_{E_i} \right| &\geq \frac{1}{6} \end{aligned}$$

Thus,

$$P\left[\frac{1}{k} \sum_{i=1}^k \mathbb{1}_{E_i} \leq \frac{1}{2}\right] \leq P\left[\left| \frac{1}{k} \sum \mathbb{1}_{E_i} - E \mathbb{1}_{E_i} \right| \geq \frac{1}{6}\right]$$

$$\text{Hoeffding} \rightarrow \leq \exp\left(-\frac{k}{18}\right). \quad \square$$

One of the potential final projects has to do with how to do this in higher dimensions than one.

Chernoff's inequality.

A weakness of Hoeffding's for bounded r.v. is that it only takes into account the worst possible spread ($b-a$) instead of the average one.

Example: If $X_i \sim \text{Ber}(p) = \begin{cases} 1 & \text{w.p. } p \\ 0 & \text{w.p. } 1-p. \end{cases}$

Then, as $p \rightarrow 0$, concentration is better, yet Hoeffding's does not improve. →

A much better bound holds for Bernoullis.

Theorem (Chernoff's inequality)

Let $X_i \sim \text{Ber}(p_i)$ be iid. Then,

$S_n = \sum_{i=1}^n X_i$ has mean $E S_n = \sum p_i = \mu$

and

$$\mathbb{P}(S_n \geq t) \leq e^{-\mu} \left(\frac{e\mu}{t}\right)^t \quad \forall t \geq \mu.$$

Proof: We apply the same idea as in the proof of Hoeffding's:

$$\mathbb{P}(S_n \geq t) = \mathbb{P}(e^{\lambda S_n} \geq e^{\lambda t})$$

$$\text{Markov's} \rightarrow \leq e^{-\lambda t} \mathbb{E} e^{\lambda S_n}$$

$$\text{Independence} \rightarrow = e^{-\lambda t} \prod_{i=1}^n \mathbb{E} e^{\lambda X_i}$$

$$= e^{-\lambda t} \prod (p_i e^\lambda + (1-p_i))$$

To combine the terms in the product we use $1+x \leq e^x$, which gives

$$\begin{aligned} \mathbb{P}(S_n \leq t) &\leq e^{-\lambda t} \prod \exp((e^\lambda - 1)p_i) \\ &= e^{-\lambda t} \exp((e^\lambda - 1) \sum_{i=1}^n p_i) \\ &= \exp(-\lambda t + (e^\lambda - 1)\mu) \end{aligned}$$

Minimizing for λ yields $\lambda = \ln(t/\mu)$
 (check!) and so

$$\begin{aligned} \mathbb{P}(S_n \leq t) &\leq \exp(-\ln(t/\mu)t + t - \mu) \\ &= e^{-\mu} \left(\frac{e\mu}{t}\right)^t \end{aligned}$$

By assumption $t \geq \mu$

Remarks:

- ① One can extend the result

$$P(S_n \leq t) \leq e^{-\mu} \left(\frac{e\mu}{t}\right)^t \quad \forall 0 \leq t \leq \mu.$$

- ② Chernoff is optimal in that

$$P(S_N \geq t) \geq \left(\frac{\mu}{t}\right)^t \quad \forall 1 \leq t \leq n.$$

- ③ When t is large (Large deviations)

$$e^{-\mu} \left(\frac{e\mu}{t}\right)^t \sim t^{-t} \leq e^{-t \log t}$$

which is much heavier than $e^{-t^2/2}$.

- ④ When $t \approx \mu$ (small deviations)

say $t = (1+\delta)\mu$

$$e^{-\mu} \left(\frac{e\mu}{t}\right)^t = e^{-\mu} \left(\frac{e}{1+\delta}\right)^{(1+\delta)\mu}$$

$$= e^{\delta\mu} \left(\frac{1}{1+\delta}\right)^{(1+\delta)\mu}$$

$$= \exp[\mu(\delta - (1+\delta) \log(1+\delta))]$$

$$\leq \exp\left(-\frac{\delta^2}{6}\mu\right)$$

$$\delta + \frac{\delta^2}{2} + O(\delta^3)$$

(Taylor expansion)