

## Lecture 5

①

Review: \* From tensorization of variance to tensorization of  $\Psi_X$ :

$$\text{Var}\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n \text{Var}(X_i) \quad \left| \quad \Psi_{\sum_{i=1}^n X_i}(\lambda) = \sum_{i=1}^n \Psi_{X_i}(\lambda)\right.$$

\* SubGaussian tails:

$$\Psi_X(\lambda) \leq \frac{\lambda^2 v}{2} \rightarrow P\left(\sum_{i=1}^n X_i > \sqrt{2n v \log \frac{1}{\delta}}\right) \leq \delta$$

\* less than subGaussian tails:

(a) Poisson-tail (right)

$$\Psi_X(\lambda) \leq v(e^{\lambda} - 1 - \lambda) \quad \left. \begin{array}{l} P\left(\sum_{i=1}^n X_i > \log \frac{1}{\delta} + \sqrt{cn v \log \frac{1}{\delta}}\right) \\ \leq \delta \end{array} \right\}$$

(b)  $\chi^2$ -tail (right)

(c) Exponential

$$\Psi_X(\lambda) \leq \frac{\lambda^2 v}{2}, \text{ for all } \lambda < \frac{c_1}{\sqrt{v}} \rightarrow P\left(\sum_{i=1}^n X_i > \sqrt{cn v \log \frac{1}{\delta}}\right) \leq \delta$$

\* Concentration bounds from tensorization of  $\Psi_X(\lambda)$

I. Hoeffding inequality:  $X_1, \dots, X_n$  indep,  $E[X_i] = 0$

Suppose  $X_i \in [a_i, b_i]$ . Then,

$$P\left(\sum_{i=1}^n X_i > t\right) \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$

## II. Bennett inequality

(2)

$$|X_i| < c, \quad \sigma^2 = \frac{1}{n} \cdot \sum_{i=1}^n \text{Var}(X_i)$$

$$P\left(\sum_{i=1}^n X_i > t\right) \leq \exp\left(\frac{-t^2}{2\sigma^2 n + \frac{2ct}{3}}\right).$$

\* Concentration using multiplicative property  
(martingale difference)

## I. Azuma inequality

$X_1, \dots, X_n$  are zero-mean and mutually uncorrelated

$$|X_i| < c_i, \quad 1 \leq i \leq n$$

$$P\left(\sum_{i=1}^n X_i > t\right) \leq \exp\left(\frac{-t^2}{2 \sum_{i=1}^n c_i^2}\right)$$

## II. McDiarmid inequality

$$f: \mathcal{X}_1 \times \dots \times \mathcal{X}_n \rightarrow \mathbb{R}$$

$X_1, \dots, X_n$  are indep, denote  $X = (X_1, \dots, X_n)$

Assume:  $f$  satisfies BDP with  $c = (c_1, \dots, c_n)$ , i.e.,

$$f(x) - f(x') \leq \sum_{i=1}^n c_i \mathbf{1}_{\{x_i \neq x'_i\}}$$

$$P(f(x) - \mathbb{E}f(x) > t) \leq \exp\left(\frac{-t^2}{2 \sum_{i=1}^n c_i^2}\right)$$

## Agenda: Applications

(3)

- [A] A bound for  $\mathbb{E}[\max_i X_i]$  for subGaussian  $X_i$
- [B]
  - longest increasing subsequence
  - longest common subsequence
  - chromatic number of Erdős-Rényi random graphs
- [C] Balls and bins

- [D] Concentration empirical measures
  - introduction to chaining
  - VC theory
  - Glivenko-Cantelli theorem

} next lecture

### A maximal inequality

Suppose  $X_1, \dots, X_n$  are subGaussian with variance

parameter  $\sigma^2$ . (no independence assumption!)

By Jensen's inequality, for  $\lambda > 0$ ,

$$\lambda \mathbb{E}[\max_i X_i] \leq \log \mathbb{E}[e^{\lambda \max_i X_i}].$$

Further,

$$\begin{aligned} \mathbb{E}[e^{\lambda \max_i X_i}] &= \mathbb{E}\left[\max_i e^{\lambda X_i}\right] \leq \mathbb{E}\left[\sum_{i=1}^n e^{\lambda X_i}\right] \\ &= \sum_{i=1}^n \mathbb{E}[e^{\lambda X_i}]. \end{aligned}$$

Thus,

$$\begin{aligned} \lambda \mathbb{E}[\max_i X_i] &\leq \log n \cdot e^{\lambda^2 \sigma^2 / 2} \\ &= \log n + \frac{\lambda^2 \sigma^2}{2}, \end{aligned}$$

i.e.,

$$\boxed{\mathbb{E}\left[\max_i X_i\right] \leq \frac{\log n}{\lambda} + \frac{\lambda \sigma^2}{2} \leq \sigma \sqrt{\frac{\log n}{2}}}$$

## B Concentration of some interesting quantities

(4)

### (a) Longest Common Sequence

Applications of  
McDiarmid

Let  $X_1, \dots, X_n, Y_1, \dots, Y_n$  be iid  $P_{XY}$ . ( $X_i, Y_i \in \mathcal{X}$ )

$$L = \max \{k \mid \exists i_1 < i_2 \dots < i_k \text{ and } j_1 < j_2 < \dots < j_k \}$$

s.t.  $X_{i_l} = Y_{j_l}, 1 \leq l \leq k$

Claim:  $P(|L - \mathbb{E}L| > t) \leq 2e^{-t^2/8n}$

If By changing any one pair  $(X_i, Y_i)$ ,  $L$  can change by at most 2. Use Hoeffding inequality.

Remark.  $\mathbb{E}L = \frac{n}{|\mathcal{X}|^2}$  (check: Chvatal-Sankoff)

### (b) Longest increasing sequence

$X_1, \dots, X_n$  are iid  $P_X$

$$f(X^n) = \max \{k \mid \exists i_1 < i_2 \dots < i_k, X_{i_1} < X_{i_2} \dots < X_{i_k}\}$$

$f$  satisfies BDP with  $(1, \dots, 1)$ .

### (c) Erdős-Rényi random graphs

$G(n, p)$ : Graph on  $n$  vertices,  $(i, j) \in E$  with prob.  $p$   
- indep. for all pairs  $(i, j)$

$$X_{ij} = \mathbf{1}(\exists \text{ edge b/w } i \text{ and } j), \quad 1 \leq i < j \leq n$$

A coloring of  $G$  is an assignment of colors to the vertices of  $G$  s.t. no two adjacent vertices have the same color.

The chromatic number of  $G$  is the min. # of colors required for coloring  $G$ .

$X(G_1(n, p)) \equiv$  chromatic # of  $G_1(n, p)$

(5)

Claim:  $P(|X(G_1(n, p)) - \mathbb{E} X(G_1(n, p))| > t) \leq 2 e^{-\frac{t^2}{2n}}$ .

Remark:  $\mathbb{E} X(G_1(n, p)) = \begin{cases} \frac{n}{\log n} & \text{when } p > \frac{\log n}{n} \\ c & \text{when } p < \frac{\log n}{n} \end{cases}$

### C Balls and bins

\* Throw m balls independently and uniformly at random into n bins

$Z = \# \text{ of empty bins}$

$$= \sum_{i=1}^n \underbrace{\mathbf{1}(\text{bin } i \text{ is empty})}_{X_i}$$

$$\mathbb{E} Z = n P(X_i = 1) = n \left(1 - \frac{1}{n}\right)^m \leq e^{-\frac{m}{n} + \log n}$$

Thus,  $\mathbb{E} Z \geq 1$  roughly only when  $m \ll n \log n$ .

Note that  $X_1, \dots, X_n$  are not indep.

Let  $B_j = \text{index of the bin where ball } j \text{ falls}$

$$\in \{1, \dots, n\}, \text{ for } 1 \leq j \leq m.$$

Then,  $Z = f(B_1, \dots, B_m)$ . Clearly,  $f$  satisfies BDP with  $(1, \dots, 1)$ . Thus, by McDiarmid's inequality,

$$P(Z - \mathbb{E} Z > t) \leq \exp\left(-\frac{t^2}{2m}\right).$$

that is

(6)

$$P\left(Z > EZ + \sqrt{2m \log \frac{1}{\delta}}\right) \leq \delta$$

$$\Leftrightarrow P\left(Z \leq EZ + \sqrt{2m \log \frac{1}{\delta}}\right) \geq 1 - \delta.$$

To study the max.  $m$  s.t. no bin is empty, it is better to consider

$Z' = \# \text{ of non-empty bins}$

and bound  $P(Z' \geq n)$ . As before

$$P\left(Z' - EZ' > \sqrt{2m \log \frac{1}{\delta}}\right) \leq \delta$$

Also

$$EZ' = n \left(1 - \left(1 - \frac{1}{n}\right)^m\right) \approx m$$

Thus,

$P(Z' \geq n)$  is small if (roughly)

$$m + \sqrt{m} \leq n. \quad \begin{aligned} & \left( m + \sqrt{m} \leq (\sqrt{m} + 0.5)^2 \leq n \right) \\ & \Leftrightarrow \sqrt{m} \leq \sqrt{n} - 0.5 \\ & \Leftrightarrow \sqrt{m} \leq \frac{\sqrt{n}}{2} \end{aligned}$$

$\uparrow$

$$m \leq \frac{n}{4}$$

That is, up to  $m = \frac{n}{c}$  there is at least one nonempty bin.

Also, similarly,  $P\left(Z' < EZ - \sqrt{2m \log \frac{1}{\delta}}\right) \leq \delta$ .

Thus, if  $m - \sqrt{m} > n$ ,  $P(Z' < n) \leq \delta$ . This again holds once  $m \geq c'n$ .

\* What happens for  $m=n$ ?

We seek to study the max. load in this regime.

(7)

$X_i^{(m)}$  # of balls in bin i

$$X_i^{(m)} \sim \text{Bin}\left(m, \frac{1}{n}\right)$$

Thus,

$$P\left(\max_i X_i^{(m)} > t\right) \leq n P\left(X_1^{(m)} > t\right)$$

(by union bound)

By Bennett inequality,

$$P\left(X_1^{(m)} - \frac{m}{n} > t\right) \leq \exp\left(-\frac{t^2}{2\frac{m}{n} + 2t}\right)$$

In particular, for  $m=n$ ,

$$P\left(X_1^{(m)} - 1 > c \log \frac{n}{\delta}\right) \leq \frac{\delta}{n}$$

$$\Rightarrow P\left(\max_{1 \leq i \leq n} X_i^{(m)} > 1 + c \log \frac{n}{\delta}\right) \leq \delta$$

$$\text{Max. load} = O(\log n)$$

Reading assignment: Poisson approximation  
 (Mitzenmacher and Upfal, "Probability and computing")