

# S&DS 602: High-Dimensional Probability and Applications

MDP Finds application to many fields:

- Statistics/ML: high-dim. data, parameter estimation
- CS: randomized algorithms, average-case complexity
- Engineering: multi-agent systems, randomized codes
- Math: combinatorial structures, Markov chains, high-dim. geometry

This course: common principles, methods, techniques.

## I. Concentration & measure

$X_1, \dots, X_n$  independent r.v.'s. If  $f(X_1, \dots, X_n)$  depends "not too much" on any individual  $X_i$ , then it is "typically close to" its expectation.

Ex: (CLT)  $X_1, \dots, X_n$  iid,  $\mathbb{E}X_i = 0$ ,  $\text{Var} X_i = \sigma^2$ .

$$P\left[\left|\frac{1}{n} \sum_{i=1}^n X_i\right| > t \cdot \frac{\sigma}{\sqrt{n}}\right] \rightarrow 2(1 - \underbrace{\Phi(t)}_{\text{standard } \mathcal{N}(0,1) \text{ CDF}}) \text{ as } n \rightarrow \infty.$$

Goals: (1) Non-asymptotic tail bounds

(2) Beyond sums of i.i.d. r.v.'s:

- sums of i.i.d. vectors and matrices
- U-statistics  $\frac{1}{\binom{n}{2}} \sum_{i < j} k(X_i, X_j)$
- polynomials  $p(X_1, \dots, X_n)$
- bounded differences

$$|f(x_1, \dots, x_i, \dots, x_n) - f(x_1, \dots, x'_i, \dots, x_n)| \leq C.$$

- Lipschitz functions, i.e.  $\|\nabla f(x_1, \dots, x_n)\|_2 \leq C$ .

II. Suprema of random processes,  $\sup_{t \in T} X_t$

If  $\{X_t\}_{t \in T}$  "sufficiently continuous," then size of  $\sup_{t \in T} X_t$  is controlled by "complexity" of  $T$ .

Examples:

(1) Norms of random vectors and matrices

$$\|v\|_2 = \sup_{u: \|u\|_2=1} u^T v, \quad \|A\|_{\text{op}} = \sup_{u, v: \|u\|_2=\|v\|_2=1} u^T A v$$

(2) Gaussian processes, e.g.

Gaussian complexity:  $\sup_{t \in T} g^T t, g_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0,1)$

Min-cut:  $\sup_{x \in \{-1,1\}^m} x^T A x, (A_{ij})_{ij} \stackrel{i.i.d.}{\sim} \mathcal{N}(0,1)$

(3) Empirical risk minimization:

$$\sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \ell(\theta, X_i) - \mathbb{E} \ell(\theta, X) \right|$$

Lecture 1: Chernoff bound, subgaussian r.v.'s, martingale method

Example:  $X_1, \dots, X_n \in \{0,1\}$ , i.i.d. Bernoulli( $p$ ),  $q > p$ .

$$\text{CLT: } P\left[\frac{1}{n} \sum_{i=1}^n X_i \geq q\right] \approx 1 - \Phi\left(\frac{q-p}{\sqrt{p(1-p)}}\right) \approx e^{-\frac{n(q-p)^2}{2p(1-p)}}$$

$$\text{For } q-p \sim \sqrt{\frac{p(1-p)}{n}}, \text{ large } n.$$

$$\text{Cramer's Thm: } P\left[\frac{1}{n} \sum_{i=1}^n X_i \geq q\right] \approx e^{-n D(q||p)}$$

$$D(q||p) = (1-q) \log \frac{1-q}{1-p} + q \log \frac{q}{p}, \text{ For } q-p \sim 1, \text{ large } n$$

$$\left[ \begin{array}{l} \text{(Taylor exp.)} \\ \approx \end{array} \frac{(q-p)^2}{2p(1-p)} \text{ For } q \approx p \right]$$

Non-asymptotic bound, valid for all  $q > p$ ?

① Chebyshev/Markov inequality:

$$\begin{aligned} P\left[\frac{1}{n}\sum_{i=1}^n X_i \geq q\right] &\leq P\left[\left(\frac{1}{n}\sum_{i=1}^n X_i - p\right)^2 \geq (q-p)^2\right] \\ &\leq \frac{E\left[\left(\frac{1}{n}\sum_{i=1}^n X_i - p\right)^2\right]}{(q-p)^2} \leq \frac{p(1-p)}{n(q-p)^2} \end{aligned}$$

②  $k^{\text{th}}$ -moment:

$$\begin{aligned} P\left[\frac{1}{n}\sum_{i=1}^n X_i \geq q\right] &\leq P\left[\left|\frac{1}{n}\sum_{i=1}^n X_i - p\right|^k \geq (q-p)^k\right] \\ &\leq \frac{E\left|\frac{1}{n}\sum_{i=1}^n X_i - p\right|^k}{(q-p)^k} \quad \text{Sharp (for } k(n) \gg 1) \\ &\quad \text{but harder to bound.} \end{aligned}$$

③ Moment/cumulant generating function:

$$\psi(\lambda) := \log E e^{\lambda(X_i - p)}$$

$$\begin{aligned} P\left[\frac{1}{n}\sum_{i=1}^n X_i \geq q\right] &= P\left[e^{\lambda\sum_{i=1}^n (X_i - p)} \geq e^{\lambda n(q-p)}\right] \quad \forall \lambda \geq 0 \\ &\leq e^{-\lambda n(q-p)} E e^{\lambda\sum_{i=1}^n (X_i - p)} = e^{-n(\lambda(q-p) - \psi(\lambda))} \end{aligned}$$

$$\text{Here: } \psi(\lambda) = \log E e^{\lambda X_i} - \lambda p = \log(1-p + p e^\lambda) - \lambda p.$$

$$\sup_{\lambda \geq 0} \lambda(q-p) - \psi(\lambda) = D(q||p).$$

$$\Rightarrow P\left[\frac{1}{n}\sum_{i=1}^n X_i \geq q\right] \leq e^{-n \cdot D(q||p)}$$

Abstracting this argument:

Then (Chernoff bound):  $X_1, \dots, X_n$  iid,  $\mathbb{E}X_i = 0$ . Let  $\psi(\lambda) = \log \mathbb{E} e^{\lambda X_i}$ .

Define the Fenchel-Legendre transform  $\psi^*(t) = \sup_{\lambda \geq 0} \lambda t - \psi(\lambda)$ .

Then  $\mathbb{P}\left[\frac{1}{n} \sum_{i=1}^n X_i \geq t\right] \leq e^{-n\psi^*(t)}$  for any  $t \geq 0$ .

Examples:

•  $X_i \sim \text{Bern}(p) - p$ .  $\psi(\lambda) = \log(1 - p + pe^\lambda) - \lambda p$ ,  $\psi^*(t) = D(t+p \| p)$ .

•  $X_i \sim \text{Pois}(\theta) - \theta$ .  $\psi(\lambda) = \theta(e^\lambda - 1)$ ,  $\psi^*(t) = (t+\theta) \log\left(1 + \frac{t}{\theta}\right) - t$ .

•  $X_i \sim \mathcal{N}(0, \sigma^2)$ .  $\psi(\lambda) = \frac{\lambda^2 \sigma^2}{2}$ ,  $\psi^*(t) = \frac{t^2}{2\sigma^2}$ .

Def: A mean-zero r.v.  $X$  is  $\sigma^2$ -subgaussian if  $\mathbb{E} e^{\lambda X} \leq e^{\frac{\lambda^2 \sigma^2}{2}}$

(i.e.  $\psi(\lambda) \leq \frac{\lambda^2 \sigma^2}{2}$ ) for all  $\lambda \in \mathbb{R}$ .

Check:  $\mathcal{N}(0, \sigma^2)$ ,  $\text{Bern}(p) - p$  are subgaussian.  $\text{Pois}(\theta) - \theta$  is not.

Lemma (Hoeffding): If  $a \leq X \leq b$  a.s., then  $X - \mathbb{E}X$  is  $\frac{(b-a)^2}{4}$ -subgaussian.

Proof: Assume WLOG  $\mathbb{E}X = 0$ .

Define  $\tilde{X}$  w/ law  $\mathbb{P}[\tilde{X} \in A] = \mathbb{E}[\mathbb{1}_{X \in A} e^{\lambda X}] / \mathbb{E} e^{\lambda X}$ .

$$\Rightarrow \psi(\lambda) = \log \mathbb{E} e^{\lambda X}, \quad \psi(0) = 0$$

$$\psi'(\lambda) = \frac{\mathbb{E} X e^{\lambda X}}{\mathbb{E} e^{\lambda X}} = \mathbb{E} \tilde{X}, \quad \psi'(0) = \mathbb{E} X = 0.$$

$$\psi''(\lambda) = \frac{\mathbb{E} X^2 e^{\lambda X}}{\mathbb{E} e^{\lambda X}} - \left( \frac{\mathbb{E} X e^{\lambda X}}{\mathbb{E} e^{\lambda X}} \right)^2 = \mathbb{E} \tilde{X}^2 - (\mathbb{E} \tilde{X})^2 = \text{Var} \tilde{X} \leq \frac{(b-a)^2}{4}.$$

$$\text{So } \psi(\lambda) = \int_0^\lambda \psi'(t) dt = \int_0^\lambda \int_0^t \psi''(s) ds dt \leq \frac{(b-a)^2}{4} \int_0^\lambda \int_0^t ds dt = \frac{\lambda^2 (b-a)^2}{8}.$$

Prop: If  $X$  is mean-zero and  $\sigma^2$ -subgaussian, then

$$\mathbb{P}[X \geq t] \leq e^{-\frac{t^2}{2\sigma^2}} \quad \forall t \geq 0.$$

Proof:  $\mathbb{P}[X \geq t] \leq e^{-\lambda t} \mathbb{E} e^{\lambda X} \leq e^{-\lambda t + \frac{\lambda^2 \sigma^2}{2}}$ . Pick  $\lambda = \frac{t}{\sigma^2}$ .

Prop: The following are equivalent:

(a)  $\exists K_1 > 0$  s.t.  $\mathbb{P}[|X| \geq t] \leq 2e^{-t^2/K_1^2} \quad \forall t \geq 0$

(b)  $\exists K_2 > 0$  s.t.  $\|X\|_{L^p} := (\mathbb{E}|X|^p)^{1/p} \leq K_2 \sqrt{p} \quad \forall p \geq 1$ .

(c)  $\exists K_3 > 0$  s.t.  $\mathbb{E} e^{\frac{X^2}{K_3^2}} \leq 2$ .

If  $\mathbb{E} X = 0$ , then these are also equivalent to

(d)  $\exists K_4 > 0$  s.t.  $\psi(\lambda) \leq \lambda^2 K_4^2 \quad \forall \lambda \in \mathbb{R}$ .

Proof: (a)  $\Rightarrow$  (b): [Lemma: If  $Z \geq 0$  a.s., then  $E Z = \int_0^\infty P[Z \geq t] dt$ .

Proof:  $Z = \int_0^Z dt = \int_0^\infty \mathbb{1}\{t \leq Z\} dt$ . Take  $E$  on both sides.]

$$E |X|^p = \int_0^\infty P[|X|^p \geq t] dt = \int_0^\infty P[|X|^p \geq u^p] \cdot p u^{p-1} du$$

$$\stackrel{(a)}{\leq} \int_0^\infty 2e^{-u^2/k^2} \cdot p u^{p-1} du \leq 2p \left(\frac{k_1}{\sqrt{2}}\right)^p \int_0^\infty e^{-v^2/2} v^{p-1} dv \leq (CK_1/p)^p$$

$$(b) \Rightarrow (c): E e^{\lambda^2 X^2} = \sum_{k \geq 0} \frac{\lambda^{2k}}{k!} E X^{2k} \stackrel{(b)}{\leq} \sum_{k \geq 0} \frac{\lambda^{2k}}{k!} (K_2 \sqrt{2k})^{2k}$$

$$\text{Apply } k! \geq (k/2)^k: E e^{\lambda^2 X^2} \leq \sum_{k \geq 0} (2e K_2^2 \lambda^2)^k \leq 2 \text{ for } \lambda = \frac{1}{CK_2}$$

$$(c) \Rightarrow (a): P[|X| \geq t] \leq P[e^{X^2/k^2} \geq e^{t^2/k^2}] \\ \leq e^{-t^2/k^2} \cdot E e^{X^2/k^2} \stackrel{(c)}{\leq} 2e^{-t^2/k^2}$$

If  $EX = 0$ :

$$(b) \Rightarrow (d): E e^{\lambda X} = 1 + \sum_{k \geq 2} \frac{\lambda^k}{k!} E X^k$$

$$\text{For } k \geq 1: E |X|^{2k+1} \leq \frac{1}{2} (E |X|^{2k} + E |X|^{2k+2})$$

$$\Rightarrow E e^{\lambda X} \leq \sum_{k \geq 0} \frac{2^k \lambda^{2k} E X^{2k}}{(2k)!} \stackrel{(b)}{\leq} \sum_{k \geq 0} \frac{2^k \lambda^{2k} (K_2 \sqrt{2k})^{2k}}{(2k)!}$$

$$\text{Apply } k^k \leq e^k \cdot k!, (2k)! \geq (2^k \cdot k!)^2: E e^{\lambda X} \leq \sum_{k \geq 0} \frac{(CK_2 \lambda^2)^k}{k!} = e^{CK_2 \lambda^2}$$

(d)  $\Rightarrow$  (a): Apply previous prop. to  $X$  and  $-X$ .

Def:  $\|X\|_{\psi_2} = \inf \{K > 0: \mathbb{E} e^{X^2/K^2} \leq 2\}$  is the subgaussian norm.

[Vershynin takes  $\|X\|_{\psi_2} < \infty$  to be definition of subgaussian if  $\mathbb{E}X \neq 0$ ]

Then (Hoeffding's inequality): If  $X_1, \dots, X_n$  are independent,  $\mathbb{E}X_i = 0$ ,  $X_i$  is  $\sigma_i^2$ -subgaussian, then  $\sum_{i=1}^n X_i$  is  $\sum_{i=1}^n \sigma_i^2$ -subgaussian.

Thus,  $\mathbb{P}\left[\sum_{i=1}^n X_i \geq t\right] \leq e^{-\frac{t^2}{2\sum_{i=1}^n \sigma_i^2}}$  for all  $t \geq 0$ .

Proof:  $\mathbb{E} e^{\lambda \sum_{i=1}^n X_i} = \prod_{i=1}^n \mathbb{E} e^{\lambda X_i} \leq \prod_{i=1}^n e^{\frac{\lambda^2 \sigma_i^2}{2}} = e^{\frac{\lambda^2}{2} \sum_{i=1}^n \sigma_i^2}$  □

Cor: If  $a_i \leq X_i \leq b_i$  a.s., then

$\mathbb{P}\left[\sum_{i=1}^n X_i - \mathbb{E}X_i \geq t\right] \leq e^{-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}}$  for all  $t \geq 0$ .

Proof: Follows from Hoeffding's Lemma. □

Example:  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$ . Then  $\mathbb{P}\left[\frac{1}{n} \sum_{i=1}^n X_i - p \geq t\right] \leq e^{-2nt^2}$ .

(Cruder than the above Chernoff bound, esp. for  $p$  far from  $1/2$ )

Martingale method

How to show concentration of non-linear functions  $f(X_1, \dots, X_n)$ ?

Idea: Write  $M_i = \mathbb{E}[f(X_1, \dots, X_n) | X_1, \dots, X_i]$ ,

$$f(X_1, \dots, X_n) - \mathbb{E} f(X_1, \dots, X_n) = M_n - M_0 = \sum_{i=1}^n \underbrace{M_i - M_{i-1}}_{\Delta_i}$$

$\{M_i\}_{i=1}^n$  is a martingale (w.r.t. filtration  $\mathcal{F}_i = \sigma(X_1, \dots, X_i)$ )

i.e.  $\mathbb{E}[M_i | \mathcal{F}_{i-1}] = M_{i-1}$ .

Then (Azuma-Hoeffding): Suppose  $\{M_i\}_{i=1}^n$  is a martingale adapted to  $\{\mathcal{F}_i\}_{i=1}^n$  and  $\Delta_i := M_i - M_{i-1}$  satisfy

$$\mathbb{E}[e^{\lambda \Delta_i} | \mathcal{F}_{i-1}] \leq e^{\frac{\lambda^2 \sigma_i^2}{2}} \text{ a.s. } (\Delta_i \text{ is conditionally } \sigma_i^2 \text{-subgaussian)}$$

Then  $\sum_{i=1}^n \Delta_i = M_n - M_0$  is  $\sum_{i=1}^n \sigma_i^2$ -subgaussian, and

$$\mathbb{P}\left[\sum_{i=1}^n \Delta_i \geq t\right] \leq e^{-\frac{t^2}{2 \sum_{i=1}^n \sigma_i^2}}$$

Proof:  $\mathbb{E} e^{\lambda \sum_{i=1}^k \Delta_i} = \mathbb{E}\left[\mathbb{E}\left[e^{\lambda \sum_{i=1}^k \Delta_i} \mid \mathcal{F}_{k-1}\right]\right]$

$$= \mathbb{E}\left[e^{\lambda \sum_{i=1}^{k-1} \Delta_i} \cdot \mathbb{E}\left[e^{\lambda \Delta_k} \mid \mathcal{F}_{k-1}\right]\right] \leq e^{\frac{\lambda^2 \sigma_k^2}{2}} \mathbb{E} e^{\lambda \sum_{i=1}^{k-1} \Delta_i}$$

Then by induction,  $\mathbb{E} e^{\lambda \sum_{i=1}^n \Delta_i} \leq e^{\frac{\lambda^2}{2} \sum_{i=1}^n \sigma_i^2}$ , and tail bound follows as before. □

Corollary (McDiarmid's bounded differences): Let

$$\|D_i f\|_\infty = \sup_{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n} \left[ \sup_z f(x_1, \dots, x_{i-1}, z, x_{i+1}, \dots, x_n) - \inf_z f(x_1, \dots, x_{i-1}, z, x_{i+1}, \dots, x_n) \right]$$

If  $X_1, \dots, X_n$  are independent, then

$$P \left[ f(X_1, \dots, X_n) - \mathbb{E} f(X_1, \dots, X_n) \geq t \right] \leq e^{-\frac{2t^2}{\sum_{i=1}^n \|D_i f\|_\infty^2}}$$

Proof: Take

$$\Delta_i = \mathbb{E}[f(X_1, \dots, X_n) | X_1, \dots, X_i] - \mathbb{E}[f(X_1, \dots, X_n) | X_1, \dots, X_{i-1}]$$

$$\text{Hence } \Delta_i \geq \underbrace{\mathbb{E} \left[ \inf_z f(x_1, \dots, x_{i-1}, z, x_{i+1}, \dots, x_n) - f(x_1, \dots, x_n) | X_1, \dots, X_{i-1} \right]}_{:= A_i}$$

$$\Delta_i \leq \underbrace{\mathbb{E} \left[ \sup_z f(x_1, \dots, x_{i-1}, z, x_{i+1}, \dots, x_n) - f(x_1, \dots, x_n) | X_1, \dots, X_{i-1} \right]}_{:= B_i}$$

where  $A_i, B_i$  are  $\mathcal{F}_{i-1}$ -measurable,  $|B_i - A_i| \leq \|D_i f\|_\infty$  a.s.

Then  $\Delta_i$  is conditionally  $\frac{\|D_i f\|_\infty^2}{4}$ -subgaussian by

Hoeffding Lemma, and result follows from Azuma-Hoeffding.  $\square$

Example (Rademacher complexity): Let  $\varepsilon_1, \dots, \varepsilon_n$  be iid Rademacher variables, i.e.  $P[\varepsilon_i = \pm 1] = \frac{1}{2}$ . Let  $T \subseteq \mathbb{R}^n$ ,

$$f(\varepsilon_1, \dots, \varepsilon_n) = \sup_{t \in T} \varepsilon^\top t.$$

Then  $\|D_n f\|_\infty = 2 \sup_{t \in T} |t_i|$ , so

$$P[|f(\varepsilon_1, \dots, \varepsilon_n) - \mathbb{E}f(\varepsilon_1, \dots, \varepsilon_n)| \geq u] \leq 2e^{-\frac{u^2}{2\sigma^2}}, \quad \sigma^2 = \sum_{i=1}^n \sup_{t \in T} t_i^2.$$

Later in the course: Improve this to  $\sigma^2 = \sup_{t \in T} \sum_{i=1}^n t_i^2 = \sup_{t \in T} \|t\|_2^2$ .

Example (U-statistic): Let  $X_1, \dots, X_n$  be iid,  $h: \mathbb{R}^2 \rightarrow \mathbb{R}$  with  $\|h\|_\infty \leq B$ ,  $h(x, y) = h(y, x)$ . Consider

$$f(X_1, \dots, X_n) = \frac{1}{\binom{n}{2}} \sum_{i < j} h(X_i, X_j).$$

$$\begin{aligned} \text{Then } \|D_n f\|_\infty &\leq \sup_{x, x'} \frac{1}{\binom{n}{2}} \sum_{j \neq i} |h(x, x_j) - h(x', x_j)| \\ &\leq \frac{(n-1) \cdot 2B}{\binom{n}{2}} = \frac{4B}{n}. \end{aligned}$$

$$\text{So } P[|f(X_1, \dots, X_n) - \mathbb{E}f(X_1, \dots, X_n)| \geq \epsilon] \leq 2e^{-\frac{n\epsilon^2}{8B^2}}.$$

More on this example later.

Example (Shamir, Spencer '86): Let  $G \sim \mathcal{G}(n, p)$  be an Erdős-Rényi graph, i.e.  $P[i \sim j] = p$  independently for all  $i, j \in \{1, \dots, n\}$ .

Let  $\chi(G)$  be the "chromatic number" — minimal # colors needed to color all vertices s.t. no edge has 2 vertices of same color.

Let  $X_1 =$  all edges  $1 \rightarrow \{2, \dots, n\}$ .

$X_2 =$  all edges  $2 \rightarrow \{3, \dots, n\}$

$\vdots$

$X_{n-1} =$  edge  $n-1 \rightarrow n$ .

Fixing all but  $X_i$ ,  $\chi(G)$  smallest when  $X_i = (0, 0, \dots, 0)$

largest when  $X_i = (1, 1, \dots, 1)$ , and  $\|D_i \chi\|_\infty \leq 1$ .

$$\Rightarrow P[|\chi(G) - \mathbb{E}\chi(G)| \geq t] \leq 2e^{-\frac{2t^2}{n-1}}$$

For any  $p \in (0, 1)$ , as  $n \rightarrow \infty$ ,  $\chi(G) = \mathbb{E}\chi(G) + O_p(\sqrt{n})$ .

[It is known  $\mathbb{E}\chi(G) \sim \frac{n}{\log n}$ .]