

MATHEMATICS OF DATA SCIENCE

December 12, 2023

Today: Last lecture (and last concentration inequalities)

RANDOM MATRICES AND COVARIANCE ESTIMATION

① Random matrices

→ Random matrices can be generated from random variables, random vectors or directly from a random matrix distribution

Example of matrix distribution: Wishart distribution

$A \in \mathbb{R}^{m \times m}$ follows a Wishart distribution $W_n(\Sigma, m)$ where $m \geq 1$
 $n \geq 1$
and $\Sigma = \Sigma^T > 0$ if its probability density is

$$p(A=B) = \frac{1}{2^{\frac{mn}{2}} \det(\Sigma)^{\frac{m}{2}} \Gamma_n\left(\frac{m}{2}\right)} \det(B)^{\frac{m-1}{2}} e^{-\frac{1}{2} \text{trace}(\Sigma^{-1}B)}$$

↑
Gamma function

→ Generalization of the χ^2 distribution ("chi-squared")

↳ Concentration inequalities for random matrices are obtained thanks to assumptions on the coefficients of these matrices or on the matrices themselves.

Theorem Let $A \in \mathbb{R}^{m \times m}$ be a rectangular random matrix with independent, mean zero, subgaussian entries.

(Every A_{ij} is mean zero and subgaussian $\forall i=1..m, j=1..m$,
and the A_{ij} s are independent)

Then $\exists C > 0$ such that $\forall t \geq 0$,

$$\mathbb{P}(\|A\| \leq CK(\sqrt{m} + \sqrt{m} + t)) \geq 1 - 2\exp(-t^2)$$

where $\|A\| = \max_{\substack{x \in \mathbb{R}^m \\ x \neq 0}} \frac{\|Ax\|}{\|x\|}$ and $K = \max_{\substack{1 \leq i \leq m \\ 1 \leq j \leq m}} \|A_{ij}\|_{\psi_2}$

subgaussian norm

Corollary

Let $A \in \mathbb{R}^{m \times m}$ is a random symmetric matrix such that $\{A_{ij}\}_{1 \leq i \leq j \leq m}$ are independent, mean zero, subgaussian random variables.

Then, $\forall t \geq 0$,

$$\mathbb{P}(\|A\| \leq CK(\sqrt{m} + t)) \geq 1 - 4\exp(-t^2)$$

where $C > 0$, $K = \max_{1 \leq i \leq j \leq m} \|A_{ij}\|_{\psi_2}$.

Proof: Split A into its upper and lower triangular part

$$A = \begin{bmatrix} A_{11} & & & \\ & \ddots & & \\ & & \ddots & \\ A_{1m} & \dots & & A_{mm} \end{bmatrix} = \underbrace{\begin{bmatrix} A_{11} & & & \\ & \ddots & & \\ & & \ddots & \\ 0 & & & A_{mm} \end{bmatrix}}_{A^u} + \underbrace{\begin{bmatrix} 0 & & & \\ A_{12} & \dots & & \\ & \ddots & \ddots & \\ A_{1m} & & & A_{mm} \end{bmatrix}}_{A^l}$$

$$\|A\| \leq \|A^u\| + \|A^l\|$$

$$\mathbb{P}(\|A\| \leq CK(\sqrt{m} + t)) \geq \mathbb{P}(\|A^u\| \leq CK(\sqrt{m} + t) \text{ and } \|A^l\| \leq CK(\sqrt{m} + t))$$

Remark: Matrix Bernstein's inequality does not require the coefficients of A_1, \dots, A_N to be independent

(2) Covariance matrix estimation

↳ Consider y_1, \dots, y_N iid random vectors in \mathbb{R}^m

Let $\mu = \mathbb{E}[y_i] \in \mathbb{R}^m$ (mean vector)

and $\Sigma = \mathbb{E}[(y_i - \mathbb{E}[y_i])(y_i - \mathbb{E}[y_i])^T] \in \mathbb{R}^{m \times m}$
(covariance matrix)

→ Estimating μ : "easy" with $\frac{1}{N} \sum_{i=1}^N y_i \Rightarrow$ this concentrates around μ

→ Estimating Σ :
Use $\Sigma_N = \frac{1}{N} \sum_{i=1}^N \left(y_i - \frac{1}{N} \sum_{j=1}^N y_j \right) \left(y_i - \frac{1}{N} \sum_{j=1}^N y_j \right)^T$

(if mean known, use $\Sigma_N = \frac{1}{N} \sum_{i=1}^N (y_i - \mu)(y_i - \mu)^T$)

Theorem

Suppose that the y_i 's satisfy

$$\|y_i\|^2 \leq C \mathbb{E}[\|y_i\|^2]$$
$$\|y_i - \mu\|^2 \lesssim \|y_i - \mu\|^2$$

almost surely
for some $C > 0$

Then,

$$\mathbb{E}_{y_1, \dots, y_N} [\|\Sigma_N - \Sigma\|] \leq \hat{C} \|\Sigma\| \left(\sqrt{\frac{m \ln(m)}{N}} + \frac{m \ln(m)}{N} \right)$$

Corollary $E[\|\Sigma_N - \Sigma\|] \leq \varepsilon \|\Sigma\|$ (for small $\varepsilon > 0$) holds

when $N \geq \bar{C} \varepsilon^{-2} \ln(m) m$ for $\bar{C} > 0$

Proof idea: Apply matrix Bernstein to the random matrices

$$\left\{ \underbrace{(y_i - \mu)(y_i - \mu)^T}_{\substack{m \times 1 \quad 1 \times m \\ m \times m}} - \Sigma \right\}_{i=1, \dots, N}$$

$\in \mathbb{R}^{m \times m}$

→ independent
→ symmetric
→ zero mean

$$E[(y_i - \mu)(y_i - \mu)^T - \Sigma] \\ = E[(y_i - \mu)(y_i - \mu)^T] - \Sigma = \Sigma - \Sigma = 0$$

To apply matrix Bernstein, we show that $\|(y_i - \mu)(y_i - \mu)^T - \Sigma\| \leq K$ almost surely

$$\begin{aligned} \|(y_i - \mu)(y_i - \mu)^T - \Sigma\| &\leq \|(y_i - \mu)(y_i - \mu)^T\| + \|\Sigma\| \\ &= \|y_i - E[y_i]\|^2 + \|\Sigma\| \\ &\leq \frac{1}{2} \|y_i\|^2 + \frac{1}{2} \|E[y_i]\|^2 + \|\Sigma\| \\ &\leq \frac{1}{2} C E[\|y_i\|^2] + \frac{1}{2} \|\mu\|^2 + \|\Sigma\| \\ &\leq \frac{1}{2} \max(C, 1) \left[\underbrace{E[\|y_i\|^2] + \|\mu\|^2}_{= \text{trace}(\Sigma)} \right] + \|\Sigma\| \end{aligned}$$

Concluding remarks

→ Concentration inequalities: apply to random variables, random vectors and random matrices

⇒ Guarantee $P(Y \geq t) \rightarrow 0$
Exponentially fast as $t \rightarrow \infty$

→ The point is NOT to remember all the inequalities!

→ Two cases of study:

- Derive concentration inequalities using convex optimization
- Apply them to provide probabilistic guarantees

