

# Roman Vershynin - covariance estimation

-1-

Bud'  $p(x)$  hustota na  $\mathbb{R}^m$  a necht  $X$  je náhodná veličina s hustotou  $p$ . Tedy  $X$  je náhodný vektor se složkami

$$X = (X_1, \dots, X_m) \in \mathbb{R}^m.$$

Typická otázka: jsou  $X_1, \dots, X_m$  nezávislé, resp. jaká je covarianční matice  $\text{Cov}(X_1, \dots, X_m)$ ?

$$\begin{aligned} \text{Označení: } \text{Cov}(X) &= \text{Cov}(X_1, \dots, X_m) = \begin{pmatrix} E X_1 X_1 & \dots & E X_1 X_m \\ \vdots & \ddots & \vdots \\ E X_m X_1 & \dots & E X_m X_m \end{pmatrix} \\ &= E \begin{pmatrix} X_1 X_1 & \dots & X_1 X_m \\ \vdots & \ddots & \vdots \\ X_m X_1 & \dots & X_m X_m \end{pmatrix} = E \begin{bmatrix} X_1 \\ \vdots \\ X_m \end{bmatrix} (X_1, \dots, X_m) = E[XX^T] \end{aligned}$$

Mějme dámo  $m$  pozorování  $X_1, \dots, X_m$  ... tedy  $X_1, \dots, X_m$  jsou kopie  $X$ .

Sample covariance matrix:

$$\text{Cov}_m(X) = \frac{1}{m} \sum_{i=1}^m X_i X_i^T$$

$X_i \in \mathbb{R}^m$   
NEZÁVISLE!

"Population expectation" nahrazeno "sample expectation".

Zákon velkých čísel:  $\text{Cov}_m(X) \rightarrow \text{Cov}(X)$  almost surely.

... jak velká je třeba mít  $m$ , aby chyba byla malá?!

# 1, Bez náhodných matic

-2-

Budeme predpokladať, že  $X$  je sub-gaussovský

- momenty & tail bounds nižšie alebo stejné jako pro gaussovské  $X$

- def. např. pomocí  $\mathbb{P}(|X| \geq t) \leq 2 \exp(-t^2/K_1^2)$ ,  $t \geq 0$

nebo  $\|X\|_p = (E|X|^p)^{1/p} \leq K_2 \sqrt{p}$ ,  $p \geq 1$

nebo  $E \exp(X^2/K_3^2) \leq 2$

- $\|X\|_{\psi_2} = \inf \{t > 0: E \exp(X^2/t^2) \leq 2\}$  ... sub-gaussian norm

- Subgaussian vectors:  $\|X\|_{\psi_2} = \sup_{\|x\|_2 \leq 1} \|\langle X, x \rangle\|_{\psi_2} = \sup_{\|x\|_2 \leq 1} \left\| \sum_{j=1}^m x_j \cdot X_j \right\|_{\psi_2}$

Věta (4.7.1. u R. Vershynin): Necht  $X$  je sub-gaussovský náhodný vektor v  $\mathbb{R}^m$ . Pak ex.  $K$  (závislá na  $\|X\|_{\psi_2}$ ) tak, že pro každé  $m \geq 1$  je

$$E \|\text{Cov}_m(X) - \text{Cov}(X)\| \leq CK^2 \left( \sqrt{\frac{m}{m}} + \frac{m}{m} \right) \|\text{Cov}(X)\|.$$

"Dikar": Zároveň na odhadu vlastních (resp. singulárních) čísel a "dobry před Bernsteinem":

- Necht  $A$  je  $m \times m$  matice sloupců  $A_i$ , které jsou nezávislé, centrovány, sub-gaussovské a izotropické vektory v  $\mathbb{R}^m$ .

Pak  $\sqrt{m} - CK^2(\sqrt{m} + t) \leq \lambda_m(A) \leq \lambda_1(A) \leq \sqrt{m} + CK^2(\sqrt{m} + t)$

spravd.  $\geq 1 - 2 \exp(-t^2)$ . Zde je  $K := \max_i \|A_i\|_{\psi_2}$ .

... i.  $\| \frac{1}{m} A^T A - I \| \leq K^2 \max(\sigma, \sigma^2)$ ,  $\sigma = C \left( \sqrt{\frac{m}{m}} + \frac{t}{\sqrt{m}} \right)$  [Isotropic  $E[XX^T] = \text{Cov}(X) = I$ ]

• Je-li  $Z$  isotopický a centrováný, <sup>me'</sup>  $X = \mu + \Sigma^{1/2} Z$  ( $\Sigma \succeq 0$  fixed)  
~~me'~~ střední hodnota  $\mu$  a  $\text{Cov}(X) = \Sigma$

• Vyjádříme  $X_i = \Sigma^{1/2} Z_i = [\text{Cov}(X)]^{1/2} Z_i$ ,  $Z_i$  w isotop.,  $\|Z_i\|_2 \leq K$

$$\& \|\text{Cov}_m(X) - \text{Cov}(X)\| = \left\| \frac{1}{m} \sum_{i=1}^m X_i X_i^T - \text{Cov}(X) \right\|$$

$$= \left\| \frac{1}{m} \sum_{i=1}^m [\text{Cov}(X)]^{1/2} Z_i Z_i^T [\text{Cov}(X)]^{1/2} - \text{Cov}(X) \right\|$$

$$= \left\| \text{Cov}(X)^{1/2} \cdot \left( \frac{1}{m} \sum_{i=1}^m Z_i Z_i^T \right) \cdot \text{Cov}(X)^{1/2} - \text{Cov}(X)^{1/2} \cdot \text{Id} \cdot \text{Cov}(X)^{1/2} \right\|$$

$$= \left\| \text{Cov}(X)^{1/2} \underbrace{\left[ \frac{1}{m} \sum_{i=1}^m Z_i Z_i^T - \text{Id} \right]}_{R_m} \text{Cov}(X)^{1/2} \right\| \leq \|\text{Cov}(X)\| \cdot \|R_m\|$$

...  $A$  ...  $m \times m$  mat. matice středky  $Z_i^T \dots \frac{1}{m} A^T A - I = R_m$

&  $\|R_m\| \leq K^2 \max(\delta, \delta^2)$ ,  $\delta = C \left( \sqrt{\frac{1}{m}} + \frac{t}{\sqrt{m}} \right)$  -- & asymptot. great.

• Organizim Bernstajnozmerovosti

Opet  $X \in \mathbb{R}^m$  je náh. vektor,  $\text{Cov}(X) = \mathbb{E}XX^T$

Věta (5.6.1 u R. Veršynin)

Nechť  $X$  je náhodný vektor v  $\mathbb{R}^m$ ,  $m \geq 2$ . Předpokládejme,

že  $\|X\|_2 \leq K(\mathbb{E}\|X\|_2^2)^{1/2}$ . Pak pro  $m \geq 1$  platí

$$\mathbb{E}\|\text{Cov}_m(X) - \text{Cov}(X)\| \leq C \left( \sqrt{\frac{K^2 m \log m}{m}} + \frac{K^2 m \log m}{m} \right) \|\text{Cov}(X)\|$$

"Důkaz":  $\bullet \mathbb{E}\|X\|_2^2 = \mathbb{E} \sum_{i=1}^m X_i^2 = \sum_{i=1}^m \mathbb{E}[X_i^2] = \text{tr} \text{Cov}(X)$

$$\Rightarrow \|X\|_2^2 \leq K^2 \cdot \text{tr}(\text{Cov}(X))$$

$$\mathbb{E}\|\text{Cov}_m(X) - \text{Cov}(X)\| = \mathbb{E}\left\| \frac{1}{m} \sum_{i=1}^m X_i X_i^T - \text{Cov}(X) \right\|$$

$$= \mathbb{E}\left\| \frac{1}{m} \sum_{i=1}^m (X_i X_i^T - \text{Cov}(X)) \right\| = \frac{1}{m} \mathbb{E}\left\| \sum_{i=1}^m Z_i \right\|$$

$\underbrace{\sum_{i=1}^m Z_i}_{\substack{Z_i \\ i=1 \\ m}}$

$$\lesssim \frac{1}{m} (\sigma \sqrt{\log m} + M \log m)$$

kde  $\bullet \sigma^2 = \left\| \sum_{i=1}^m \mathbb{E}(X_i X_i^T - \text{Cov}(X)) \right\|^2 = m \left\| \mathbb{E}(XX^T - \text{Cov}(X))^2 \right\|$

$\bullet M \geq \|XX^T - \text{Cov}(X)\|$  a.s.

"Stac" vyjádřit  $\sigma^2$  a  $M$ .

$$\begin{aligned} \bullet M: \|XX^T - \text{Cov}(X)\| &\leq \underbrace{\|XX^T\|}_{\text{rank-1 matrix}} + \|\text{Cov}(X)\| \leq \|X\|_2^2 + \underbrace{\|\text{Cov}(X)\|}_{\leq t(\text{Cov}(X))} \\ &\leq K^2 t(\text{Cov}(X)) + t(\text{Cov}(X)) \leq 2K^2 t(\text{Cov}(X)) \quad (K \geq 1) \end{aligned}$$

$$\begin{aligned} \bullet G^2: \mathbb{E}(XX^T - \text{Cov}(X))^2 &= \mathbb{E}\{XX^T XX^T - \underbrace{XX^T \text{Cov}(X)}_{\text{Cov}(X)} - \underbrace{\text{Cov}(X) XX^T}_{\text{Cov}(X)} + \text{Cov}(X)^2\} \\ &= \mathbb{E}(XX^T)^2 - \text{Cov}(X)^2 \leq \mathbb{E}(XX^T)^2 \end{aligned}$$

$$\bullet (XX^T)^2 = \underbrace{XX^T}_{\|X\|_2^2} XX^T \leq \|X\|_2^2 \cdot XX^T \leq K^2 t(\text{Cov}(X)) \cdot XX^T$$

$$\bullet \mathbb{E}(XX^T)^2 \leq K^2 t(\text{Cov}(X)) \cdot \text{Cov}(X); \|\cdot\| \leq K^2 \cdot n \cdot \|\text{Cov}(X)\|^2$$

$$\begin{aligned} \Rightarrow \mathbb{E}\|\text{Cov}_m(X) - \text{Cov}(X)\| &\leq C \cdot \frac{1}{m} \left( \sqrt{m K^2 \|\text{Cov}(X)\|^2 \cdot m \log m} \right. \\ &\quad \left. + 2K^2 \cdot n \|\text{Cov}(X)\| \cdot \log m \right) \\ &= C \cdot \|\text{Cov}(X)\| \cdot \left\{ \sqrt{\frac{m}{m} \cdot K^2 \cdot \log m} + \frac{m}{m} \cdot K^2 \cdot \log m \right\} \quad \square \end{aligned}$$