

1 Overview

In the last lecture, we showed the bound on approximate isometries. In this lecture, we focus on estimating the covariance matrix from finite sample, which has a well-known application in the principal components analysis (PCA). And the material follows from section 4.7, 5.2, 9.2 in [3].

2 Covariance estimation

Let $\mathbf{X} \in \mathbb{R}^d$ be a random vector with mean zero, the covariance matrix is $\Sigma := \mathbb{E}\mathbf{X}\mathbf{X}^\top$, i.e. $\Sigma_{ij} = \mathbb{E}\mathbf{X}_i\mathbf{X}_j$. We first show the following property of covariance matrix. For $\mathbf{u} \in \mathbb{S}^{d-1}$, $\mathbb{E}\langle \mathbf{X}, \mathbf{u} \rangle = 0$ as \mathbf{X} is mean zero, then

$$\text{Var}\langle \mathbf{X}, \mathbf{u} \rangle = \mathbb{E}(\langle \mathbf{X}, \mathbf{u} \rangle^2) = \mathbb{E}\mathbf{u}^\top \mathbf{X}\mathbf{X}^\top \mathbf{u} = \mathbf{u}^\top \mathbb{E}(\mathbf{X}\mathbf{X}^\top) \mathbf{u} = \mathbf{u}^\top \Sigma \mathbf{u} \quad (1)$$

Suppose there are n i.i.d samples $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(n)} \sim \mathbf{X}$, the sample covariance is defined as

$$\hat{\Sigma}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{X}^{(i)} \mathbf{X}^{(i)\top}$$

Put differently, let $\mathbf{A} \in \mathbb{R}^{n \times d}$ have row i is $\mathbf{X}^{(i)}$, then $\hat{\Sigma}_n = \frac{\mathbf{A}^\top \mathbf{A}}{n}$.

2.1 Covariance estimation for high dimensional distributions

Recall the following theorem we proved last time for sub-gaussian isotropic random vectors.

Theorem 1 (Theorem 4.6.1 in [3]). *Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ have rows that are independent, mean zero, sub-gaussian isotropic random vectors with $K = \max_{i=1, \dots, n} \|\mathbf{A}_i\|_{\psi_2}$, then with probability $1 - 2e^{-u}$,*

$$\left\| \frac{1}{n} \mathbf{A}^\top \mathbf{A} - \mathbf{I}_d \right\|_2 \leq CK^2 \left(\sqrt{\frac{d+u}{n}} + \frac{d+u}{n} \right)$$

Remark 2. *Since rows of \mathbf{A} are isotropic whose covariance are \mathbf{I}_d , $\left\| \frac{1}{n} \mathbf{A}^\top \mathbf{A} - \mathbf{I}_d \right\|_2 = \left\| \hat{\Sigma}_n - \Sigma \right\|_2$.*

Next, we extend to the general sub-gaussian setting.

Theorem 3 (Theorem 4.7.1 in [3]). Let \mathbf{X} be a sub-gaussian random vector in \mathbb{R}^d with invertible covariance matrix Σ s.t

$$\|\langle \mathbf{X}, \mathbf{u} \rangle\|_{\psi_2} \leq K \sqrt{\mathbb{E}(\langle \mathbf{X}, \mathbf{u} \rangle^2)} \quad \text{for all } \mathbf{u} \in \mathbb{S}^{d-1}$$

with probability $1 - 2e^{-u}$, then

$$\|\hat{\Sigma}_n - \Sigma\|_2 \leq CK^2 \|\Sigma\|_2 \left(\sqrt{\frac{d+u}{n}} + \frac{d+u}{n} \right).$$

Remark 4. $\mathbb{E}(\langle \mathbf{X}, \mathbf{u} \rangle^2) = \langle \mathbf{u}, \Sigma \mathbf{u} \rangle$ by (1).

Proof. We first bring the random vectors $\mathbf{X}_1, \dots, \mathbf{X}_n$ to the isotropic position. By slides(<https://people.math.wisc.edu/~roch/hdps/roch-hdps-slides12.pdf>) or Exercise 3.2.2 in [3], there exist isotropic mean zero vectors $\mathbf{Z}^{(1)}, \dots, \mathbf{Z}^{(n)}$ s.t

$$\mathbf{X}^{(i)} = \Sigma^{1/2} \mathbf{Z}^{(i)}, \text{ for all } i = 1, \dots, n.$$

We proved previously, $\|\mathbf{Z}^{(i)}\|_{\psi_2} \leq K$. Let $\hat{\mathbf{R}}_n := \frac{1}{n} \sum_{i=1}^n \mathbf{Z}^{(i)} \mathbf{Z}^{(i)\top} - \mathbf{I}_d$, then

$$\|\hat{\Sigma}_n - \Sigma\|_2 = \|\Sigma^{1/2} \hat{\mathbf{R}}_n \Sigma^{1/2}\|_2 \leq \|\hat{\mathbf{R}}_n\|_2 \|\Sigma\|_2$$

where we apply the property $\|\mathbf{A}\mathbf{B}\|_2 \leq \|\mathbf{A}\|_2 \|\mathbf{B}\|_2$ and $\|\mathbf{A}^{1/2}\|_2^2 = \left(\sqrt{\sigma_{\max}(\mathbf{A})}\right)^2 = \|\mathbf{A}\|_2$.

Consider the $n \times d$ random matrix \mathbf{A} whose rows are $\mathbf{Z}^{(i)}$, then

$$\frac{1}{n} \mathbf{A}^\top \mathbf{A} - \mathbf{I}_d = \frac{1}{n} \sum_{i=1}^n \mathbf{Z}^{(i)} \mathbf{Z}^{(i)\top} - \mathbf{I}_d = \hat{\mathbf{R}}_n$$

Then we conclude by applying Theorem 1 for \mathbf{A} . □

2.2 Covariance estimation for lower-dimensional distributions

We found that the covariance matrix Σ of an n -dimensional distribution can be estimated from $m = O(n)$ sample points for sub-gaussian distributions. For approximately lower-dimensional distributions, smaller sample can be sufficient for covariance estimation, which means that the distribution tends to concentrate near a small subspace.

Theorem 5 (Theorem 9.2.4 in [3]). Let \mathbf{X} be a sub-gaussian random vector in \mathbb{R}^d with invertible covariance matrix Σ s.t

$$\|\langle \mathbf{X}, \mathbf{u} \rangle\|_{\psi_2} \leq K \sqrt{\mathbb{E}(\langle \mathbf{X}, \mathbf{u} \rangle^2)} \quad \text{for all } \mathbf{u} \in \mathbb{S}^{d-1}$$

with probability $1 - 2e^{-u}$, then

$$\|\hat{\Sigma}_n - \Sigma\|_2 \leq CK^4 \|\Sigma\|_2 \left(\sqrt{\frac{r+u}{n}} + \frac{r+u}{n} \right).$$

where $r = \text{tr}(\Sigma) / \|\Sigma\|_2$.

2.3 General covariance estimation

Next, we state a more general version of covariance estimation by applying matrix Bernstein inequality, see the slides(<https://people.math.wisc.edu/~roch/hdps/roch-hdps-slides20.pdf>) for more details.

Theorem 6 (Theorem 5.6.1 in [3]). *Let \mathbf{X} be a random vector in \mathbb{R}^d , $d \geq 2$. Assume that for some $K \geq 1$,*

$$\|\mathbf{X}\|_2 \leq K \left(\mathbb{E} \|\mathbf{X}\|_2^2 \right)^{1/2} \quad \text{almost surely.}$$

Then, with probability $1 - 2e^{-u}$, we have

$$\left\| \hat{\Sigma}_n - \Sigma \right\|_2 \leq C \|\Sigma\|_2 \left(\sqrt{\frac{K^2 r (\log d + u)}{n}} + \frac{K^2 r (\log d + u)}{n} \right)$$

where $r = \text{tr}(\Sigma) / \|\Sigma\|_2 \leq n$.

References

- [1] Rick Durrett, *Probability—theory and examples (fifth edition)*, Cambridge University Press, 2019.
- [2] Ramon van Handel, *APC 550: Probability in High Dimension*, Lecture Notes, 2016. <https://web.math.princeton.edu/~rvan/APC550.pdf>
- [3] Roman Vershynin, *High-dimensional probability: An introduction with applications in data science*, Cambridge University Press, 2018.