

# Lecture 08: Finite Sample PCA I

Nikola Zlatanov

Innopolis University

Advanced Statistics

27-th of March to 3-th of April, 2023

# Motivation

- To do PCA of data vectors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \in \mathbb{R}^d$ , coming from some distribution  $f_{\mathbf{X}}(\mathbf{x})$ , we needed to have the covariance matrix  $\Sigma$ .
- But what happens in cases when we do not know the covariance matrix  $\Sigma$ . Can we still do PCA then as well?
- Specifically, assume we have data vectors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ , coming from some unknown distribution  $f_{\mathbf{X}}(\mathbf{x})$  and unknown covariance matrix  $\Sigma$ .
- How to do PCA then? And what will be the error we make compared to the case when we knew the covariance matrix?

# Finite Sample PCA

- Let us have a sample of  $N$  i.i.d. data vectors  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N \in \mathbb{R}^d$ , coming from some unknown distribution  $f_{\mathbf{X}}(\mathbf{x})$  and unknown covariance matrix  $\Sigma$ .
- Let us estimate the unknown covariance matrix  $\Sigma$  as  $\hat{\Sigma}_N$  using the  $N$  samples as

$$\hat{\Sigma}_N = \frac{1}{N} \sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i^T, \quad (1)$$

- Then, let's do PCA, using the estimated covariance matrix  $\hat{\Sigma}_N$ , on the data vectors  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N$ .
- Hopefully, the results using the estimated covariance matrix  $\hat{\Sigma}_N$  will be close to doing the PCA, using the actual covariance matrix  $\Sigma$ , on the data vectors  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N$ .

# Finite Sample PCA

- An immediate question in when using the estimated covariance matrix  $\hat{\Sigma}_N$  to do PCA is the following.
- How large should the sample size  $N$  be in order for the two PCA results to be close?
- Since each of the samples  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N$  is of dimension  $d$ ,  $\Sigma$  is a  $d \times d$  matrix.
- Then, for sure, the accuracy of the PCA using the estimated covariance matrix  $\hat{\Sigma}_N$  will depend on how  $N$  depends on  $d$ .
- How will  $N$  depend on  $d$ ? Maybe for a good accuracy we need  $N = O(d)$ , or maybe  $N = O(d^2)$  or maybe  $N = O(e^d)$

# Finite Sample PCA

- We know that the PCA result using the covariance matrix  $\Sigma$  is a function of its eigenvalues,  $\lambda_i(\Sigma)$ , and eigenvectors,  $\mathbf{v}_i(\Sigma)$ .
- Hence, if  $\lambda_i(\hat{\Sigma}_N)$  and  $\mathbf{v}_i(\hat{\Sigma}_N)$  are the eigenvalues and eigenvectors of the estimated covariance matrix,  $\hat{\Sigma}_N$ , then for the two PCA results to be close, the following has to hold

$$\begin{aligned}\lambda_i(\Sigma) &\approx \lambda_i(\hat{\Sigma}_N) \\ \text{AND} \\ \mathbf{v}_i(\Sigma) &\approx \mathbf{v}_i(\hat{\Sigma}_N),\end{aligned}\tag{2}$$

i.e., the eigenvalues and eigenvectors of  $\Sigma$  and  $\hat{\Sigma}_N$  are very close to each other.

- How large should the sample size  $N$  be in order for the eigenvalues and eigenvectors of  $\Sigma$  and  $\hat{\Sigma}_N$  to be very close to each other?

# Covariance Estimation Problem

Step 1:

- Let us measure how far  $\Sigma$  and  $\hat{\Sigma}_N$  are from each other using some measurement/metric.
- For matrices, so far, we learned that we can use the Frobenius norm or the operator norm.
- Let's try the operator norm:

$$\|\hat{\Sigma}_N - \Sigma\|_{\text{op}} = \sigma_1(\hat{\Sigma}_N - \Sigma), \quad (3)$$

where  $\sigma_1(\hat{\Sigma}_N - \Sigma)$  is the largest singular value of the matrix  $\hat{\Sigma}_N - \Sigma$

- Now, we have

$$\begin{aligned} \sigma_1(\hat{\Sigma}_N - \Sigma) &= \left| \lambda_1(\hat{\Sigma}_N - \Sigma) \right| = \max_{\mathbf{v}, \|\mathbf{v}\|_2=1} \left| \mathbf{v}^T (\hat{\Sigma}_N - \Sigma) \mathbf{v} \right| \\ &= \max_{\mathbf{v}, \|\mathbf{v}\|_2=1} \left| \mathbf{v}^T \hat{\Sigma}_N \mathbf{v} - \mathbf{v}^T \Sigma \mathbf{v} \right| \end{aligned} \quad (4)$$

# Distance Between True and Sampled Covariance Matrices

- The first element in (4) is

$$\begin{aligned}
 \mathbf{v}^T \hat{\Sigma}_N \mathbf{v} &= \frac{1}{N} \sum_{i=1}^N \mathbf{v}^T \mathbf{X}_i \mathbf{X}_i^T \mathbf{v} = \frac{1}{N} \sum_{i=1}^N Y_i Y_i = \frac{1}{N} \sum_{i=1}^N Y_i^2 \\
 &= \frac{1}{N} \sum_{i=1}^N \langle \mathbf{v}, \mathbf{X}_i \rangle^2,
 \end{aligned} \tag{5}$$

where  $Y_i = \mathbf{v}^T \mathbf{X}_i = \mathbf{X}_i^T \mathbf{v} = \langle \mathbf{v}, \mathbf{X}_i \rangle$ .

- The second element in (4) is

$$\begin{aligned}
 \mathbf{v}^T \Sigma \mathbf{v} &= \mathbf{v}^T E[\mathbf{X} \mathbf{X}^T] \mathbf{v} = E[\mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v}] = E[Y Y] = E[Y^2] \\
 &= E[\langle \mathbf{v}, \mathbf{X} \rangle^2]
 \end{aligned} \tag{6}$$

where  $Y = \mathbf{v}^T \mathbf{X} = \mathbf{X}^T \mathbf{v} = \langle \mathbf{v}, \mathbf{X} \rangle$ .

# Distance Between True and Sampled Covariance Matrices

- Inserting (5) and (6) into (4), and then (4) into (3), we obtain

$$\|\hat{\Sigma}_N - \Sigma\|_{\text{op}} = \max_{\mathbf{v}, \|\mathbf{v}\|_2=1} \left| \frac{1}{N} \sum_{i=1}^N \langle \mathbf{v}, \mathbf{X}_i \rangle^2 - E[\langle \mathbf{v}, \mathbf{X} \rangle^2] \right|, \quad (7)$$

which is a random variable (RV).



# Covariance Estimation Theorem

- Thm (Covariance Estimation Theorem): Let  $\mathbf{X}_i \in \mathbb{R}^d$  be  $N$  sample vectors independently and identically drawn from the Gaussian distribution  $\mathcal{N}(0, \mathbf{\Sigma})$ . Let each element of  $\mathbf{X}_i$  be with variance  $\sigma^2$ . Let  $\delta \geq 0$ . Then, the following holds

$$\Pr \left\{ \|\hat{\mathbf{\Sigma}}_N - \mathbf{\Sigma}\|_{\text{op}} \leq \delta \right\} \geq 1 - 2 \exp \left( -d \sqrt{\frac{2}{\epsilon}} \right) \exp \left( -N \frac{(1-\epsilon)\delta}{2\sigma^2} \min \left\{ \frac{(1-\epsilon)\delta}{2\sigma^2}, 1 \right\} \right) \quad (8)$$

- What this Thm says is that the distance between  $\hat{\mathbf{\Sigma}}_N$  and  $\mathbf{\Sigma}$  can be made smaller than  $\delta$ , with probability that goes to one as  $N$  increases, as long as  $N$  satisfies

$$N \frac{(1-\epsilon)\delta}{2\sigma^2} \min \left\{ \frac{(1-\epsilon)\delta}{2\sigma^2}, 1 \right\} > d \sqrt{\frac{2}{\epsilon}} \quad (9)$$

# Covariance Estimation Theorem

- What the Covariance Estimation Theorem tells us is that  $N = o(d)$ .
- In other words, the number of vector samples  $N$  we need to accurately estimate the covariance matrix  $\Sigma$  using

$$\hat{\Sigma}_N = \frac{1}{N} \sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i^T \quad (10)$$

depends linearly on the vectors  $\mathbf{X}_i$  length  $d$ .

- This is very good news since there is no curse-of-dimensionality in this problem.
- Imagine if it was  $N = o(e^d)$ . Then, it would have been impossible to estimate  $\hat{\Sigma}_N$  using (10) since for  $d = 100$  we would need at least  $e^{100}$  samples, which is more than the number of atoms in the observable universe.

# Distance Between Finite Sample PCA and True PCA

- So far we only proved that we can make  $\|\hat{\Sigma}_N - \Sigma\|_{\text{op}}$  vary small by increasing  $N$ , given that condition (9) holds.
- But PCA does not depend on  $\|\hat{\Sigma}_N - \Sigma\|_{\text{op}}$  explicitly, and instead depends on the distances between the eigenvalues and eigenvectors of  $\Sigma$  and  $\hat{\Sigma}_N$ .
- In fact, for the PCA results using  $\hat{\Sigma}_N$  to be very close to the PCA results using  $\Sigma$ , we needed the following to hold for all dimensions  $i$  over which we do the PCA (what happens in the other dimensions is irrelevant)

$$\lambda_i(\Sigma) \approx \lambda_i(\hat{\Sigma}_N) \quad \text{AND} \quad \mathbf{v}_i(\Sigma) \approx \mathbf{v}_i(\hat{\Sigma}_N), \quad (11)$$

i.e., the eigenvalues and eigenvectors of  $\Sigma$  and  $\hat{\Sigma}_N$  are very close to each other in all dimensions  $i$  over which we do the PCA.

- In the following, we will show that by making  $\|\hat{\Sigma}_N - \Sigma\|_{\text{op}}$  vary small, we will implicitly achieve (11).

## Distance Between Eigenvalues

- Thm (Weyl's Inequality):  $\forall d \times d$  symmetric matrices  $\mathbf{A}$  and  $\mathbf{B}$ , the following holds

$$\max_i |\lambda_i(\mathbf{A}) - \lambda_i(\mathbf{B})| \leq \|\mathbf{A} - \mathbf{B}\|_{\text{op}} \quad (12)$$

- We will not prove it!
- What this theorem says is that the largest distance between the eigenvalues of any two symmetric matrices  $\mathbf{A}$  and  $\mathbf{B}$  is upper bounded by the operator norm of  $\mathbf{A} - \mathbf{B}$
- Hence, if we make  $\|\hat{\Sigma}_N - \Sigma\|_{\text{op}}$  vary small, we implicitly make the distance between the eigenvalues of  $\Sigma$  and  $\hat{\Sigma}_N$  also very small.  
i.e., the following holds

$$\lambda_i(\Sigma) \approx \lambda_i(\hat{\Sigma}_N) \quad (13)$$

- How about the eigenvectors. Now, they are more tricky, but we can still do it as shown in the following.

## Distance Between Eigenvectors

- Thm (Davis-Kahan Inequality): For any symmetric matrices  $\Sigma$  and  $\hat{\Sigma}_N$  with eigenvectors  $\mathbf{v}_1(\Sigma), \mathbf{v}_2(\Sigma), \dots, \mathbf{v}_k(\Sigma)$  and  $\mathbf{v}_1(\hat{\Sigma}_N), \mathbf{v}_2(\hat{\Sigma}_N), \dots, \mathbf{v}_k(\hat{\Sigma}_N)$ , respectively, let  $P_\Sigma$  and  $P_{\hat{\Sigma}_N}$  be defined as

$$P_\Sigma = \sum_{i=1}^k \mathbf{v}_i(\Sigma) \mathbf{v}_i^T(\Sigma) \quad (14)$$

$$P_{\hat{\Sigma}_N} = \sum_{i=1}^k \mathbf{v}_i(\hat{\Sigma}_N) \mathbf{v}_i^T(\hat{\Sigma}_N). \quad (15)$$

Then, the following holds

$$\|P_{\hat{\Sigma}_N} - P_\Sigma\|_{\text{op}} \leq \frac{\|\hat{\Sigma}_N - \Sigma\|_{\text{op}}}{\lambda_k(\hat{\Sigma}_N) - \lambda_{k+1}(\hat{\Sigma}_N)} \quad (16)$$

- We will not prove it!

## Distance Between Eigenvectors

- What the Davis-Kahan Inequality Thm says is that by making  $\|\hat{\Sigma}_N - \Sigma\|_{\text{op}}$  vary small, we implicitly make the first  $k$  eigenvectors of  $\hat{\Sigma}_N$  to be very close to the first  $k$  eigenvectors of  $\Sigma_N$ , as long as there is a ‘spectral gap’ between the  $k$ -th and the  $(k + 1)$ -th eigenvalue of  $\hat{\Sigma}_N$ .
- ‘Spectral gap’ between the  $k$ -th and the  $(k + 1)$ -th eigenvalue of  $\hat{\Sigma}_N$ , means that they are as far away as possible. See this Figure:
- This Thm also gives us a rule of thumb used in the current PCA algos: Do PCA dimension reduction up to the point when we observe a spectral gap in the eigenvalues of  $\hat{\Sigma}_N$ . Anything beyond the spectral gap consider as noise and throw it away.
- In the next lecture, we will arrive at the same concept, noise that needs to be thrown away, via a different method which will reveal us even more insights into random matrices.