

STA 437/2005:  
Methods for Multivariate Data  
Weeks 8: Covariance matrix estimation

Piotr Zwiernik

University of Toronto

# Covariance matrix estimation

We started our discussion of PCA on the population level.

- maximizing  $\mathbf{u}^\top \Sigma \mathbf{u}$  gives a direction of the highest variance of  $\mathbf{X} \in \mathbb{R}^m$ .

In practice we have no access to  $\Sigma \in \mathbb{S}_+^m$ .

The main approach is to estimate  $\Sigma$  using the sample covariance matrix  $S_n$ .

Recall that  $S_n$  is almost unbiased.

It can be shown that it is a consistent estimator of  $\Sigma$ .

In other words, if  $n$  is “very large”,  $S_n$  should be a good estimator of  $\Sigma$ .

# High-dimensional problems

How large  $n$  has to be generally depends on  $m$ .

This is intuitively clear because  $\Sigma$  has  $\binom{m}{2} = \frac{m(m+1)}{2}$  parameters to estimate.

**Classical asymptotics** lets  $n \rightarrow \infty$  keeping  $m$  fixed.

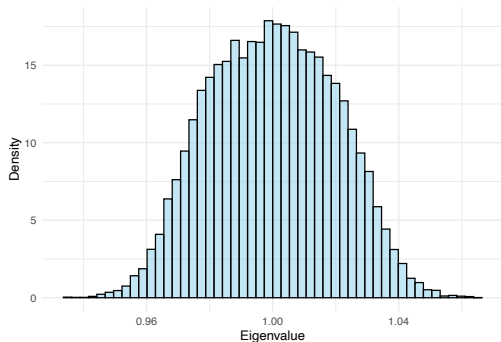
**High-dimensional asymptotics** studies estimation when both  $m, n \rightarrow \infty$ .

- We assume  $m/n \rightarrow \gamma \in [0, 1)$ .

# Why does it matter?

Suppose that  $\Sigma = I_m$ . If  $S_n$  is close to  $\Sigma$  all its eigenvalues should be close to 1.

Consider a simple example:  $m = 3$ ,  $n = 1000$ . Sample  $S_n$  several times and look at the histogram of eigenvalues.



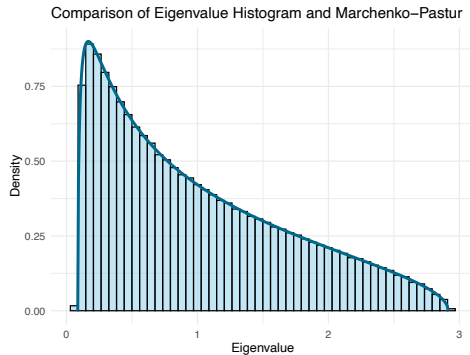
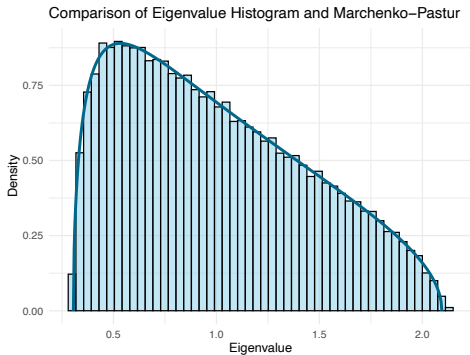
Indeed! A **sharp** concentration around 1.

Arguably, this is a very extreme situation.

In typical applications the ration  $n/m$  is much smaller.

Consider now the eigenvalue distribution in the same setting but with much higher  $m$ .

Take  $n = 1000$ ,  $m = 200$  and  $m = 500$ .



The eigenvalues deviate from 1, following the **Marchenko–Pastur law**.

# Marchenko-Pastur Law

Marchenko-Pastur Law gives the limiting distribution of the eigenvalues of  $S_n$  ( $\Sigma = I_m$ ) in the limiting case when  $m/n \rightarrow \gamma$ .

## Marchenko-Pastur Law

Let  $\lambda_{\min} := (1 - \sqrt{\gamma})^2$ ,  $\lambda_{\max} := (1 + \sqrt{\gamma})^2$ . Then MP Law has density

$$f_{\text{MP}}(\lambda) = \frac{1}{2\pi\gamma\lambda} \sqrt{(\lambda_{\max} - \lambda)(\lambda - \lambda_{\min})} \quad \text{for } \lambda \in [\lambda_{\min}, \lambda_{\max}].$$

# Alternative Estimators

This example shows that  $S_n$  is not a good estimator of  $\Sigma$  when  $m/n$  is too large.

General approach: if there is some additional structure in  $\Sigma$ , exploit it.

- This stabilizes the estimators.

This approach may be problematic if you exploited structure that is not there.

We now review some common approaches that work well in a wide-range of scenarios.

# Alternative Estimators Overview

**Linear Shrinkage:**  $\hat{\Sigma}_{ls} = (1 - \lambda)S_n + \lambda I_m$  for some  $\lambda \in (0, 1)$ .

- Reduces variance by shrinking towards  $I_m$ .

**Graphical Lasso:** We consider penalized Gaussian log-likelihood. Define

$$\hat{K} := \arg \min_{K \in \mathbb{S}_+^m} \{ \text{tr}(S_n K) - \log \det(K) + \lambda \|K\|_1 \},$$

where  $\|K\|_1 = \sum_{i \neq j} |K_{ij}|$  ( $\ell_1$ -penalty). Finally,  $\hat{\Sigma}_{\text{glasso}} = \hat{K}^{-1}$ .

- Promotes sparsity in the precision matrix.



# Alternative Estimators Continued

**Factor Models:** Suppose  $\Sigma$  has the form  $\Sigma = LL^\top + \Psi$  where  $L \in \mathbb{R}^{m \times r}$  for  $r < m$  and  $\Psi$  is diagonal.

- ▶ We will show how to exploit this in estimation.
- ▶ Probabilistic PCA gives one example with  $\Psi = \sigma^2 I_m$ .

**Thresholding-Based Methods:** If  $\Sigma$  has zeros, it is natural to estimate

$$\hat{\Sigma}_{\text{thresh}} = \{(S_n)_{ij} \cdot \mathbb{I}(|(S_n)_{ij}| > \tau)\}_{i,j}.$$

Sets small covariance entries to zero.

# Tyler's Scatter Estimator

This is a popular estimator in robust statistics.

If  $X \sim E(\mathbf{0}, \Sigma)$  then  $Z = \Sigma^{-1/2}X$  is spherical;  $Z/\|Z\|$  is uniform on the unit sphere.

Then  $\frac{1}{m}I_m = \text{var}(Z/\|Z\|) = \mathbb{E}(\frac{1}{\|Z\|^2}ZZ^\top) = \mathbb{E}(\frac{1}{X^\top \Sigma^{-1}X} \Sigma^{-1/2}XX^\top \Sigma^{-1/2})$ .

Equivalently  $\mathbb{E}(\frac{1}{X^\top \Sigma^{-1}X} XX^\top) = \frac{1}{m}\Sigma$ . Consider a sample version of this equation:

$$\sum_{i=1}^n \frac{1}{x^{(i)\top} \Sigma^{-1} x^{(i)}} x^{(i)} x^{(i)\top} = \frac{n}{m} \Sigma.$$

Under mild conditions, there is a unique solution; computed using fixed-point iterations:

$$\hat{\Sigma}^{(k+1)} = \frac{m}{n} \sum_i \frac{x^{(i)} x^{(i)\top}}{x^{(i)\top} (\hat{\Sigma}^{(k)})^{-1} x^{(i)}}.$$

# Summary

Estimating the covariance matrix in modern applications raises many challenges.

If  $\Sigma$  satisfies some structure, we could exploit it to stabilize estimation.

We study some structures that can appear in practice.

- ▶ Diagonal plus low rank.
- ▶  $\Sigma$  or  $\Sigma^{-1}$  sparse.

This is an active area of research. Links to random matrix theory.