

LECTURE 20

Last Class: Covariance estimation problem for $X \in \mathbb{R}^d$, $\mathbb{E}X = 0$.

- (Population) covariance matrix: $\Sigma = \mathbb{E}XX^T$
- Sample covariance matrix: $\Sigma_n = \frac{1}{n} \sum_{i=1}^n X_i X_i^T$, where X_i are iid copies of X

$$\Sigma_n \stackrel{?}{\approx} \Sigma \quad \text{in the operator norm}$$

$$\|\Sigma_n - \Sigma\| = \max_{v \in S^{d-1}} |v^T (\Sigma_n - \Sigma) v| \quad (\text{HW})$$

$$\exists \varepsilon\text{-net } \mathcal{N} \subset S^{d-1}, \quad |\mathcal{N}| \leq \left(\frac{2}{\varepsilon} + 1\right)^d = 9^d \quad \text{if we choose } \boxed{\varepsilon = \frac{1}{4}}$$

$$\|\Sigma_n - \Sigma\| \leq \frac{1}{1-2\varepsilon} \max_{v \in \mathcal{N}} |v^T (\Sigma_n - \Sigma) v| \quad (\text{Last class \& HW})$$

$$= 2 \cdot \max_{v \in \mathcal{N}} \left| \underbrace{\frac{1}{n} \sum_{i=1}^n \langle X_i, v \rangle^2}_{S(v)} - \mathbb{E} \langle X, v \rangle^2 \right| \leq ? \quad (*)$$

• For each v , $S(v)$ = sum of indep. r.v.'s. Use:

TKM (Bernstein's inequality - lec. 11)

i.e. $P\{|Z_i| \geq t\} \leq 2 \exp(-\frac{t}{M})$ for smallest $M = \|Z_i\|_{\psi_1}$

If Z_i are independent subexponential r.v.'s, then

$$P\left\{ \sum_{i=1}^n (Z_i - \mathbb{E}Z_i) \geq t \right\} \leq 2 \exp \left[-c \cdot \min \left(\frac{t^2}{\sigma^2}, \frac{t}{K} \right) \right]$$

where $\sigma^2 = \sum_{i=1}^n \|Z_i\|_{\psi_1}^2$, $K = \max_i \|Z_i\|_{\psi_1}$

If all Z_i are independently distributed $\Rightarrow \sigma^2 = K^2 n$

$$P\left\{ \left| \frac{1}{n} \sum_{i=1}^n Z_i - \mathbb{E}Z \right| \geq \delta \right\} \stackrel{t = \delta n}{\leq} 2 \exp \left[-c \cdot \min \left(\frac{\delta^2 n^2}{K^2 n}, \frac{\delta n}{K} \right) \right] = 2 \exp \left[-c \cdot \min \left(\frac{\delta^2}{K^2}, \frac{\delta}{K} \right) n \right]$$

(**)

ALL TOOLS ARE READY TO PROVE:

Thm (Covariance Estimation) If $X_i \sim N(0, \Sigma)$ and $n \geq Cd$,

then $\|\Sigma_n - \Sigma\| \leq 0.1 \|\Sigma\|$
with probability at least $1 - 2e^{-cn}$.

Proof WLOG. $\|\Sigma\| = 1$.

① Fix any $v \in S^{d-1}$ ← unit sphere of \mathbb{R}^d .

$X = \text{mean } 0 \text{ normal r. vector} \Rightarrow \langle X, v \rangle \sim \text{mean } 0 \text{ normal r. vector}$

$\Rightarrow \langle X, v \rangle \sim N(0, \sigma^2)$.

↖ ? let's compute it.

v is unit

$$\sigma^2 = \mathbb{E} \langle X, v \rangle^2 = v^T (\mathbb{E} X X^T) v = v^T \Sigma v \leq \|\Sigma\| = 1.$$

$$\Rightarrow \mathbb{P} \{ |\langle X, v \rangle| \geq t \} \leq 2 \exp \left(-\frac{t^2}{2\sigma^2} \right) \leq 2 \exp(-t^2/2)$$

↑ gaussian tail (lec. 4)

$$\Rightarrow \mathbb{P} \{ \langle X, v \rangle^2 \geq s \} = \mathbb{P} \{ |\langle X, v \rangle| \geq \sqrt{s} \} \leq 2 \exp(-s/2)$$

$\Rightarrow \langle X, v \rangle^2$ is subexponential, $\|\langle X, v \rangle^2\|_{\psi_1} \leq 2$

• Apply Bernstein (**p.1) for $Z_i = \langle X_i, v \rangle^2$, $K=2$, $\delta=0.01 \Rightarrow$

$$\mathbb{P} \left\{ \left| \frac{1}{n} \sum_{i=1}^n \langle X_i, v \rangle^2 - \mathbb{E} \langle X, v \rangle^2 \right| \geq 0.01 \right\} \leq 2 \exp(-c'n).$$

$\underbrace{\hspace{10em}}_{S(v)}$

$$\textcircled{2} \quad \mathbb{P} \left\{ \max_{v \in \mathcal{N}} |S(v)| \geq 0.01 \right\} \leq \underbrace{|\mathcal{N}|}_{q^d \text{ (p.1)}} \cdot 2 \exp(-c'n) \leq 2 \exp \left(-\frac{c'n}{2} \right)$$

if $d < c'n$ ☺

③ The complement holds w/prob. $\geq 1 - 2 \exp(-c'n/2)$. When it holds,

$$\|\Sigma_n - \Sigma\| \stackrel{(**) \text{ p.1}}{\leq} 2 \cdot \max_{v \in \mathcal{N}} |S(v)| \leq 2 \cdot 0.01 \leq 0.1.$$

PERTURBATION THEORY

From $\Sigma_n \approx \Sigma$ to $\lambda_i(\Sigma_n) \approx \lambda_i(\Sigma)$, $v_i(\Sigma_n) \approx v_i(\Sigma)$:



THM (Weyl's inequality) $\forall d \times d$ symmetric matrices A, B :

$$\max_i |\lambda_i(A) - \lambda_i(B)| \leq \|A - B\|$$

Eigenvectors: $\begin{cases} \text{(a) up to sign} \\ \text{(b) spectral gap needed: if } \lambda_k(A) \approx \lambda_{k+1}(A), \text{ perturbation} \\ \text{can swap the order of eigenvectors.} \end{cases}$

THM (Davis-Kahan inequality)

Let A, B be $d \times d$ symmetric matrices, $1 \leq k \leq d$,

$P_A :=$ orthogonal projection onto $\text{span}\{v_1(A), \dots, v_k(A)\}$ and

$P_B :=$ orthogonal projection onto $\text{span}\{v_1(B), \dots, v_k(B)\}$. Then

$$\|P_A - P_B\| \leq \frac{\|A - B\|}{\lambda_k(A) - \lambda_{k+1}(A)}$$

• Combine Cov. Est. Thm (p. 2) with Weyl & Davis-Kahan



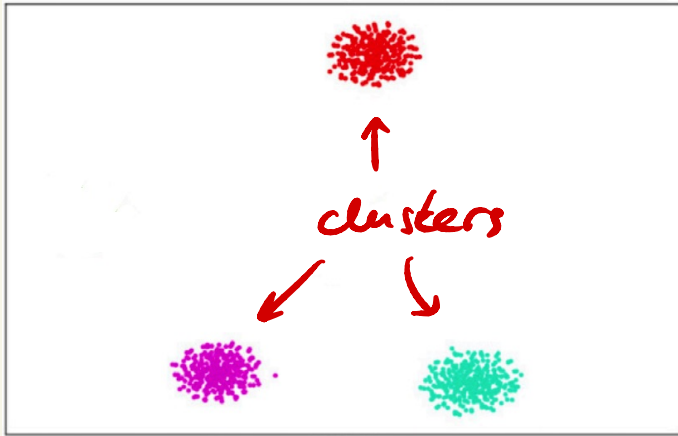
SUMMARY | $\text{PCA}_k(\text{sample}) \approx \text{PCA}_k(\text{population})$ as long as:

(a) sample size n is linear in $\dim d$ (or larger)

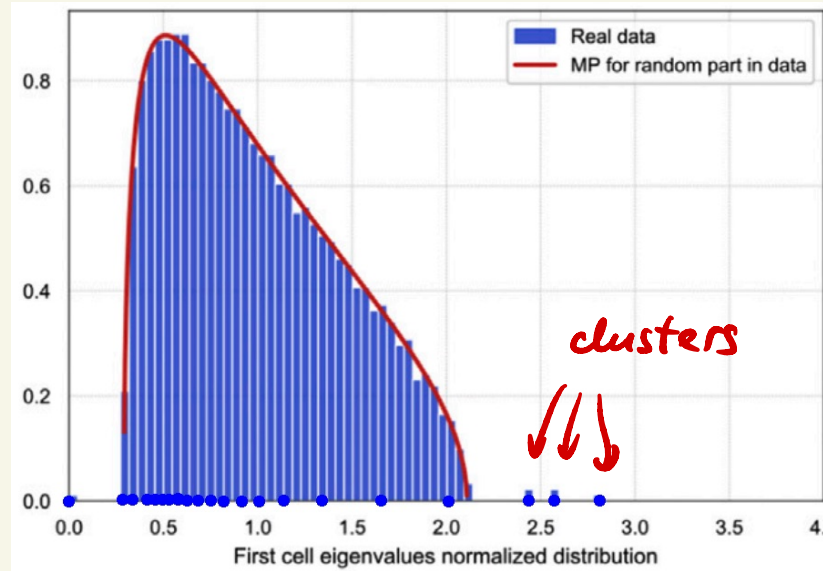
(b) there is an eigengap: the top k eigenvalues of Σ are separated from the rest of the spectrum.

HW
-state?

Example: single-cell data from [Aparicio et al. 2020]



→



noise
"bulk"
"signal"
outliers

$$X \sim N(0, I_n)$$

$$\Sigma = I_n$$

REMARKS : ① Linear sample size $n = O(d)$ is optimal (HW)

② Normal distribution $X \sim N(0, \Sigma)$ can be generalized to \forall subgaussian (simple - see book) HW?

③ Structured data \Rightarrow smaller sample size n ?

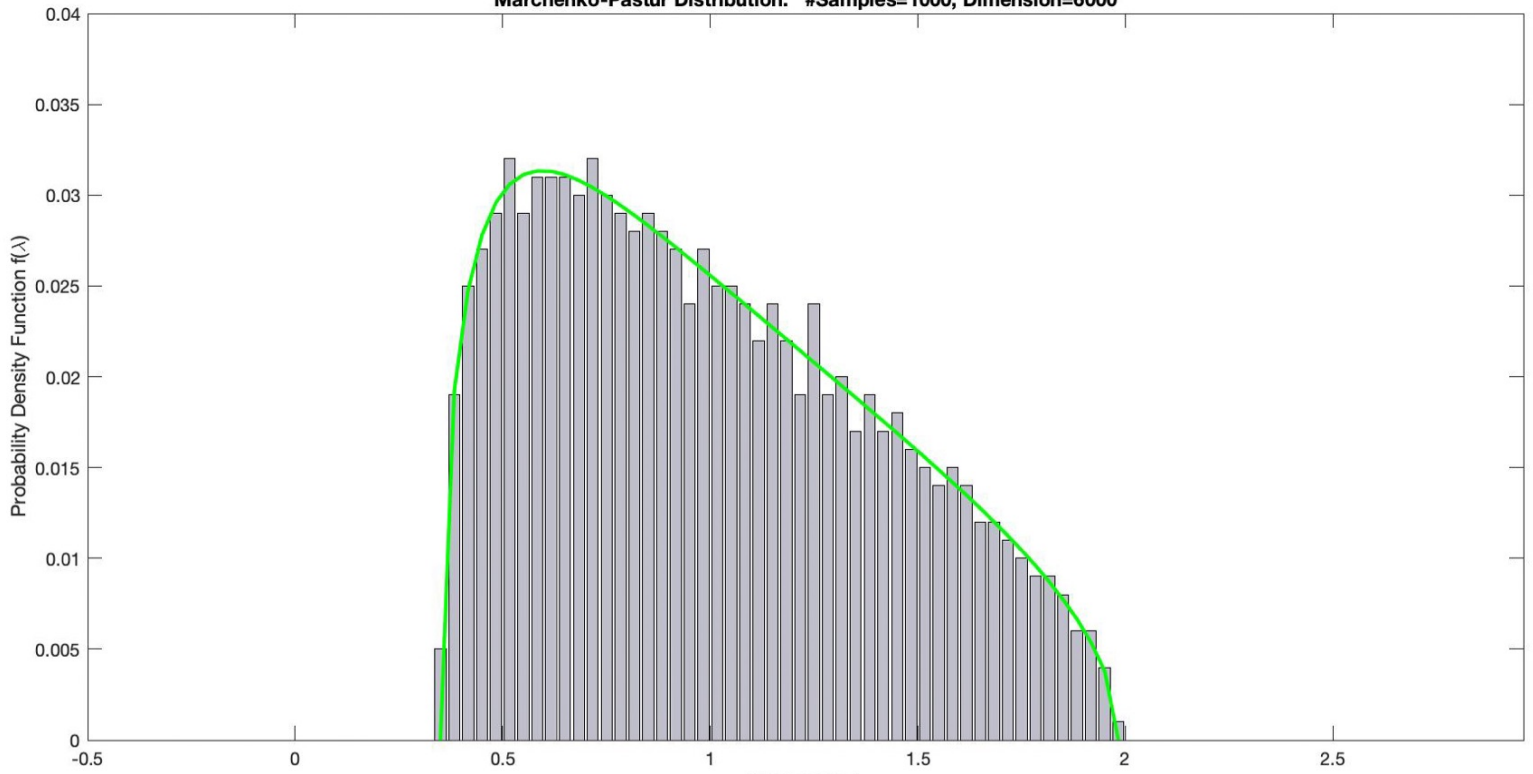
Yes : (a) trivially, if $X = (X_1, X_1, \dots, X_1)$ then represent $X \in \mathbb{R}^1 \Rightarrow n = O(1)$.

(b) more generally, if data is approximately k -dimensional, $k = O(n)$ is enough.

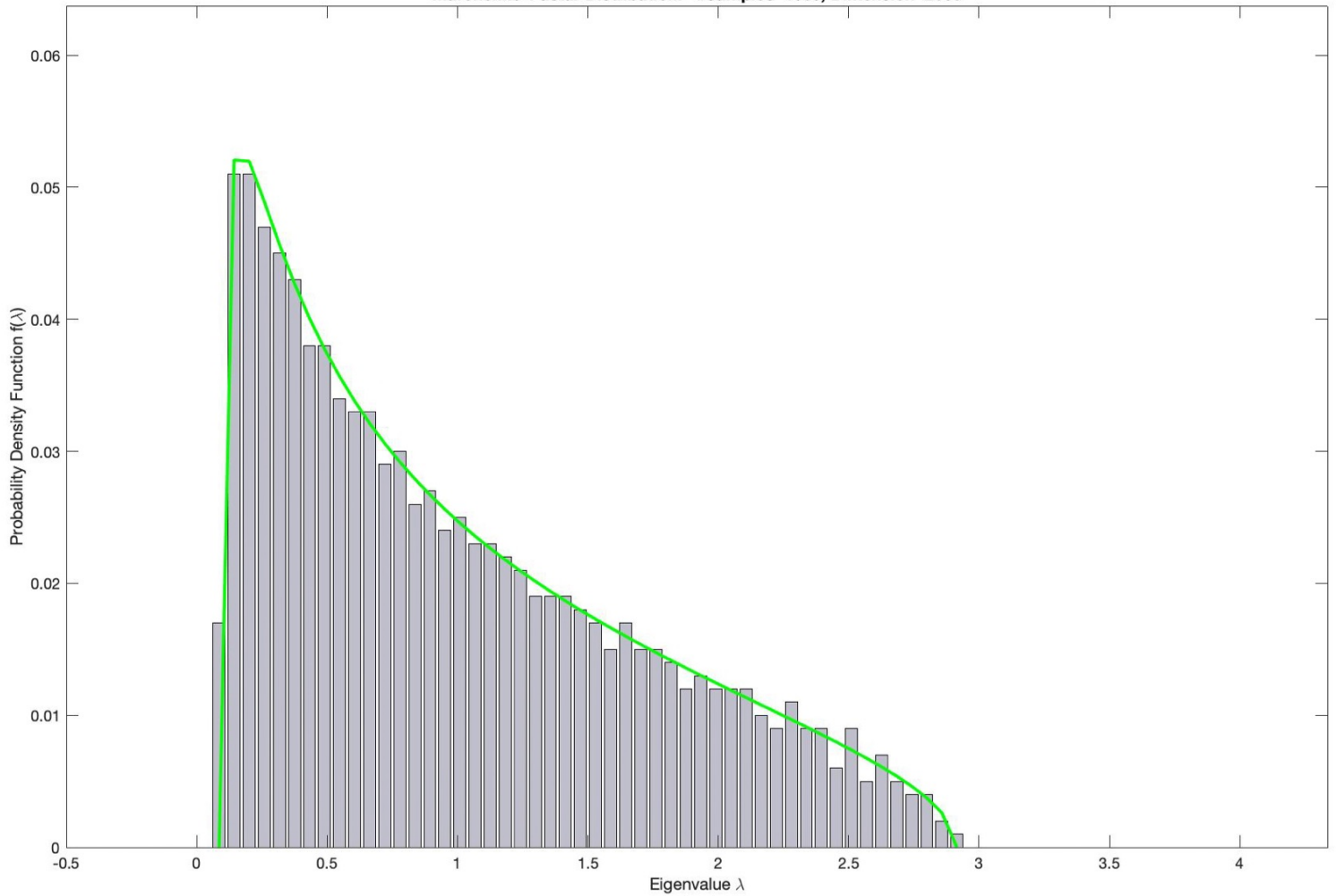
[Koltchinskii-Couillard 2017] = Thm. 9.2.4 in the book.

Will probably cover later.

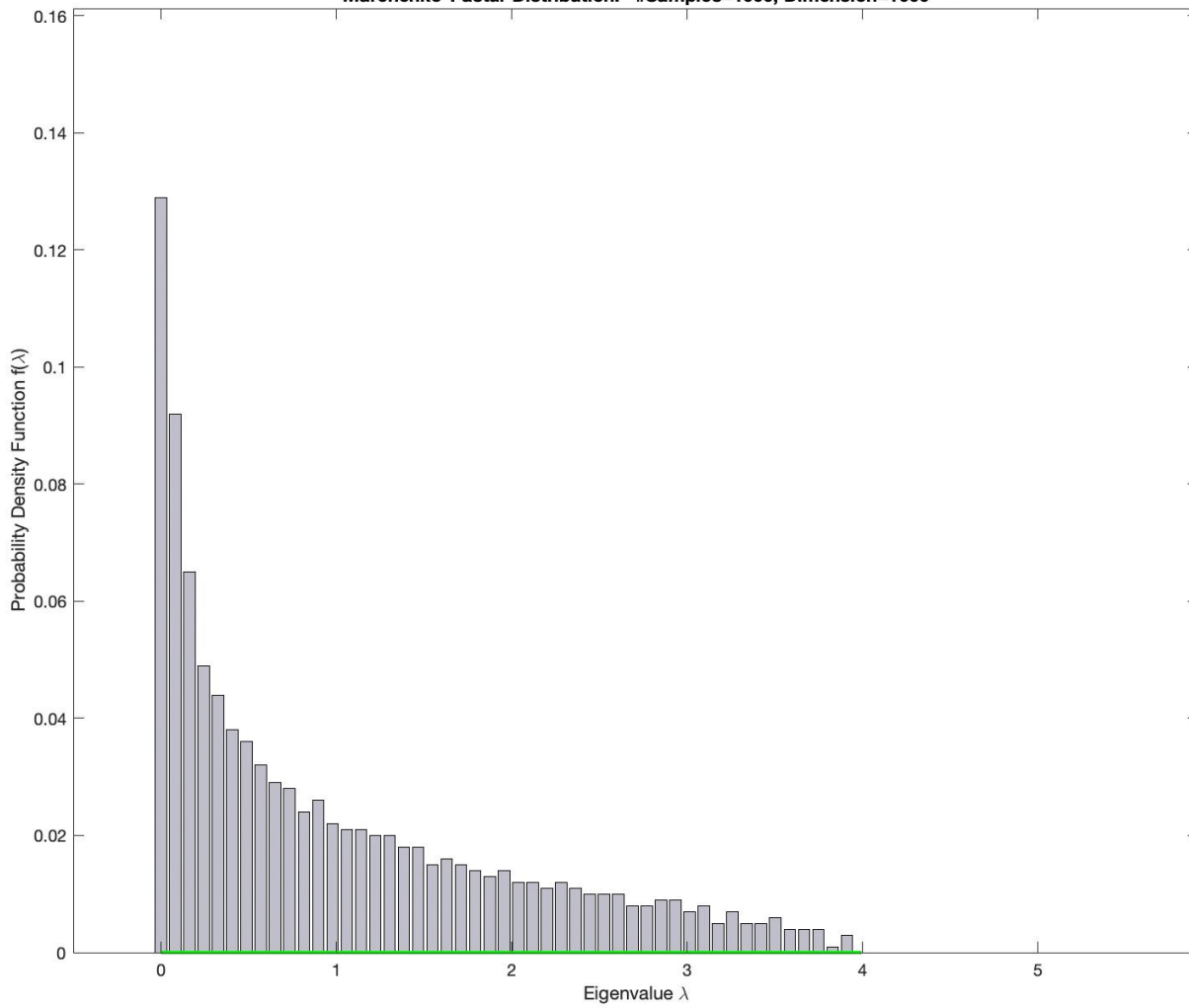
Marchenko-Pastur Distribution. #Samples=1000, Dimension=6000



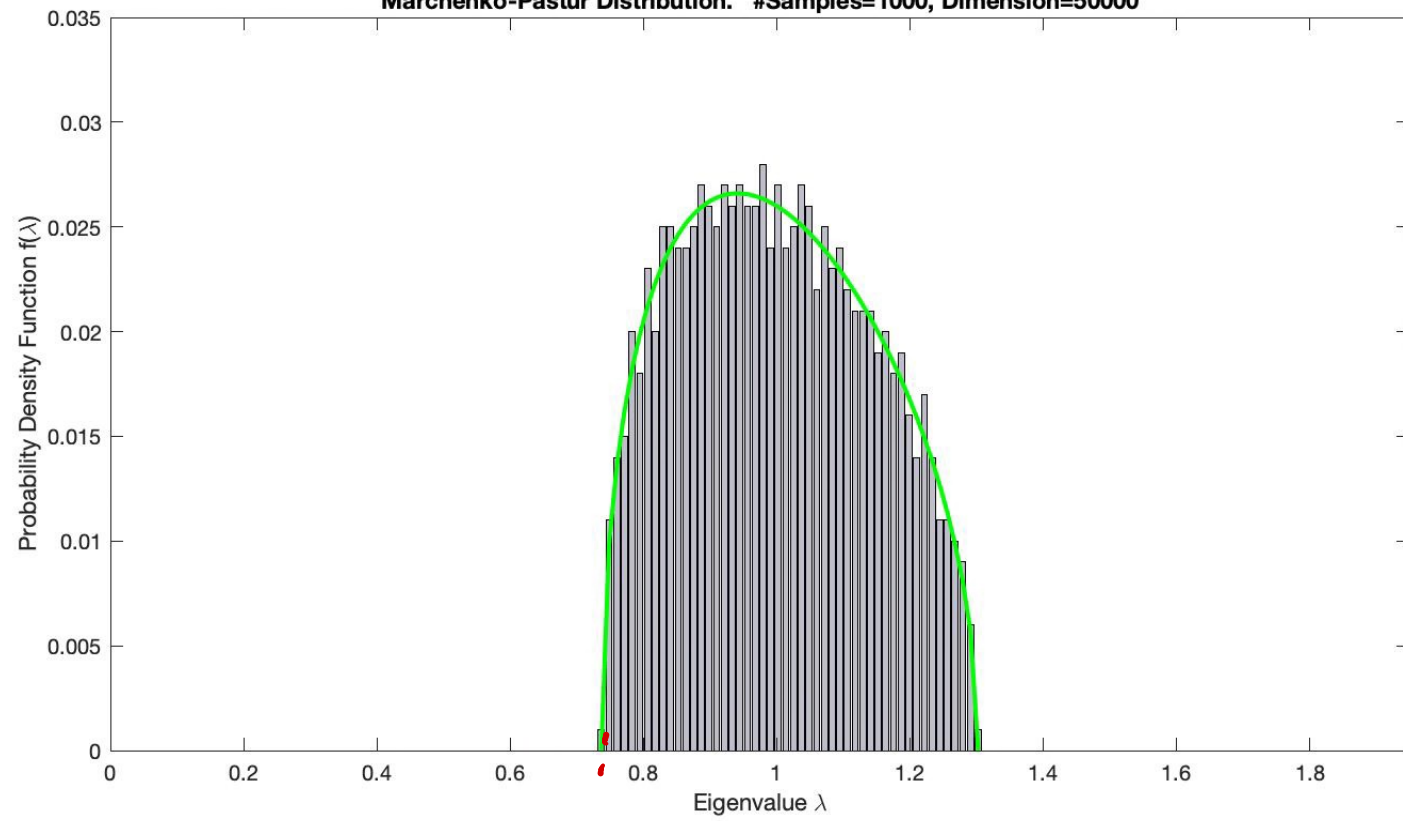
Marchenko-Pastur Distribution. #Samples=1000, Dimension=2000



Marchenko-Pastur Distribution. #Samples=1000, Dimension=1000



Marchenko-Pastur Distribution. #Samples=1000, Dimension=50000



- histogram $e_i z_i^s(\Sigma_n)$ seems \rightarrow density as $d, n \rightarrow \infty$, $\frac{d}{n} \rightarrow \lambda$
 compactly supported! No outliers (good for PCA)

