Principal Component Analysis for a Spiked Covariance Model with Largest Eigenvalues of the Same Asymptotic Order of Magnitude

Addy M. Bolívar Cimé

Centro de Investigación en Matemáticas A.C.

May 1, 2010

Contenido

Introduction

- Spiked Covariance Model
- Recent results

Contenido

Introduction

- Spiked Covariance Model
- Recent results



Result derived with RMT

Contenido

Introduction

- Spiked Covariance Model
- Recent results

Random Matrix Theory Result derived with RMT

3 Asymptotic results

- Consistency of eigenvalues
- Consistency of eigenvectors

An important problem in multivariate statistical analysis is the estimation of the population covariance matrix. In principal component analysis (PCA) one often fails to estimate the population eigenvalues and eigenvectors, since the sample covariance matrix is not a good approximation to the population covariance matrix when the data dimension is larger than the sample size. As pointed out in Johnstone (2001), one often observes one or a small number of large sample eigenvalues well separated from the rest. In this case, of special interest is the so-called *spiked covariance model*.

More specifically, suppose we have a $d \times n$ data matrix $X = [x_1, x_2, \ldots, x_n]$ with $d \gg n$, in the sense that $\frac{d}{n} \to \infty$, where $x_j = (x_{1j}, \ldots, x_{dj})^{\top}$, $j = 1, 2, \ldots, n$. Assume the columns of X are independent and identically distributed random vectors from a multivariate Gaussian distribution with mean zero and unknown covariance matrix Σ . The *spiked covariance model* considers a covariance matrix of the type

$$\Sigma = \operatorname{diag}(\tau_1, \tau_2, \dots, \tau_p, 1, \dots, 1), \tag{1}$$

with $\tau_1 \geq \tau_2 \geq \cdots \geq \tau_p > 1$, for some $1 \leq p < d$.

Ahn, Marron, Muller and Chi (2007) showed that for p = 1 and $\Sigma = \operatorname{diag}(d^{\alpha}, 1, \ldots, 1)$, with $\alpha > 1$, it is possible to estimate very well the first population eigenvalue and eigenvector using PCA when d and n are sufficiently large and $d \gg n$. For the case $p \ge 2$, Jung and Marron (2009) showed the consistency of the first p largest sample eigenvalues under the assumption that each of the p largest population eigenvalues, $\tau_1 \ge \tau_2 \ge \cdots \ge \tau_p$, has a different asymptotic order of magnitude when

dimension increases, that is

$$\frac{\tau_i}{d^{\alpha_i}} \to c_i, \qquad \text{as } d \to \infty,$$

where $\alpha_1 > \alpha_2 > \cdots > \alpha_p > 1$ and $c_i > 0$, $i = 1, 2, \dots, p$.

They showed that if $\hat{\tau}_1 \geq \hat{\tau}_2 \geq \cdots \geq \hat{\tau}_p$ are the *p* largest sample eigenvalues

$$\frac{\widehat{\tau}_i}{\tau_i} \stackrel{\text{w}}{\longrightarrow} \frac{\chi_n^2}{n}, \quad \text{as } d \to \infty, \tag{2}$$

for i = 1, 2, ..., p, where χ_n^2 is a Chi-square random variable with n degrees of freedom. Therefore, since $\chi_n^2/n \xrightarrow{w} 1$ as $n \to \infty$ the first p sample eigenvalues are consistent. They also show that in this case the corresponding sample eigenvectors are consistent.

To study the asymptotic behaviour of the sample eigenvalues and eigenvectors, when d, n are sufficiently large and $d \gg n$, in the case of the spiked covariance model where the p largest population eigenvalues $\tau_1, \ldots, \tau_p, p > 1$, have the same asymptotic order of magnitude when dimension increases, i.e. $\alpha_1 = \alpha_2 = \cdots = \alpha_p = \alpha > 1$ and $c_1 = c_2 = \cdots = c_p = c > 0$.

Under this last assumption, we first prove a multivariate extension of (2), namely

$$(\frac{\widehat{\tau}_1}{\tau_1}, \frac{\widehat{\tau}_2}{\tau_2}, \dots, \frac{\widehat{\tau}_p}{\tau_p})^\top \xrightarrow{w} (cn)^{-1} (\ell_1, \ell_2, \dots, \ell_p)^\top \text{ as } d \to \infty$$

where $\ell_1 \ge \ell_2 \ge \cdots \ge \ell_p$ are the nonzero eigenvalues of a Wishart random matrix with p degrees of freedom and $n \times n$ covariance matrix cI_n .

Wishart Distribution

Addy Bolívar (CIMAT)

Lemma 2.1

If $\ell_1 \geq \ell_2 \geq \cdots \geq \ell_p$ are the nonzero eigenvalues of $U \sim W_n(p, cI_n)$ with n > p and c > 0, then the joint density of $(\ell_1, \ell_2, \dots, \ell_p)^\top$ is given by

$$f_{\ell_{1},\ell_{2},...,\ell_{p}}(\ell_{1},\ell_{2},...,\ell_{p}) = \Delta \sum_{\alpha \in S_{p}} \operatorname{sign}(\alpha) \exp(-\frac{1}{2c} \sum_{k=1}^{p} \ell_{k}) \\ * \prod_{k=1}^{p} \ell_{p+1-k}^{\alpha_{k}+(n-p-1)/2-1},$$
(3)

with $\ell_1 > \ell_2 > \cdots > \ell_p$, where

$$\Delta = \frac{\pi^{p^2/2}}{(2c)^{pn/2}\Gamma_p(\frac{p}{2})\Gamma_p(\frac{n}{2})}.$$
(4)

Proposition 2.1

If $\ell_1 \geq \ell_2 \geq \cdots \geq \ell_p$ are the nonzero eigenvalues of $U \sim \mathcal{W}_n(p, cI_n)$ with n > p and c > 0, then the characteristic function of $(\ell_1, \ell_2, \dots, \ell_p)^{\top}$ is given by

$$\begin{aligned} \mathcal{L}_{\ell_{1},\ell_{2},...,\ell_{p}}(t_{1},t_{2},...,t_{p}) &= \\ \Delta\left(\frac{2c}{p}\right)^{pn/2} \widehat{G}\left(\sum_{j=1}^{p} t_{j};\frac{pn}{2},\frac{2c}{p}\right) \sum_{k_{1}=0}^{\infty} \cdots \sum_{k_{p-1}=0}^{\infty} \Gamma\left(\frac{pn}{2} + \sum_{j=1}^{p-1} k_{j}\right) \\ &* C_{n,p}(k_{1},...,k_{p-1}) \prod_{r=1}^{p-1} \left(\frac{r}{p}\right)^{k_{r}} \frac{\widehat{G}(\sum_{j=1}^{p} t_{j};k_{r},\frac{2c}{p})}{\widehat{G}\left(\sum_{j=p+1-r}^{p} t_{j};k_{r},\frac{2c}{r}\right)}, \end{aligned}$$
(5)

where $\hat{G}(t; a, b)$ is the characteristic function of the gamma distribution Gamma(a, b) and

$$C_{n,p}(k_1,\ldots,k_{p-1}) = \sum_{\alpha \in S_p} \operatorname{sign}(\alpha) \prod_{r=1}^{p-1} \frac{\Gamma(\sum_{j=1}^r \alpha_j + \frac{r(n-p-1)}{2} + \sum_{j=1}^{r-1} k_j)}{\Gamma(\sum_{j=1}^r \alpha_j + \frac{r(n-p-1)}{2} + \sum_{j=1}^r k_j + 1)}.$$
(6)

Assume X as above with unknown covariance matrix given by the spiked covariance model (1) with p < n < d and where $\tau_1 \ge \tau_2 \ge \cdots \ge \tau_p > 1$ have the same asymptotic order of magnitude, that is $\tau_i = \tau_i(d)$, $d^{-\alpha}\tau_i(d) \rightarrow c$ as $d \rightarrow \infty$ for some c > 0 and $\alpha > 1$, $i = 1, 2, \ldots, p$. Let $\hat{\tau}_1 \ge \hat{\tau}_2 \ge \cdots \ge \hat{\tau}_p$ be the p largest sample eigenvalues of the sample covariance matrix $S = n^{-1}XX^{\top}$. Then

$$(\frac{\widehat{\tau}_1}{\tau_1}, \frac{\widehat{\tau}_2}{\tau_2}, \dots, \frac{\widehat{\tau}_p}{\tau_p})^\top \xrightarrow{w} (cn)^{-1} (\ell_1, \ell_2, \dots, \ell_p)^\top,$$

when $d \to \infty$, where $\ell_1 \ge \ell_2 \ge \cdots \ge \ell_p$ are the nonzero eigenvalues of a Wishart random matrix with distribution $W_n(p, cI_n)$.

Proposition 3.1

Let c > 0 and $\ell_1 \ge \ell_2 \ge \cdots \ge \ell_p$ be the nonzero eigenvalues of $U \sim W_n(p, cI_n)$ with n > p. Then

$$(cn)^{-1}(\ell_1,\ell_2,\ldots,\ell_p)^{ op} \xrightarrow{w} (1,1,\ldots,1)^{ op}, \qquad \text{as } n \to \infty.$$

Under the assumptions of Theorem 3.1, we have that for all $\epsilon > 0$

$$\lim_{n \to \infty} \lim_{d \to \infty} P\left(\| \left(\frac{\widehat{\tau}_1}{\tau_1}, \frac{\widehat{\tau}_2}{\tau_2}, \dots, \frac{\widehat{\tau}_p}{\tau_p} \right)^\top - (1, 1, \dots, 1)^\top \| > \epsilon \right) = 0.$$
 (7)

Under the same conditions of Theorem 3.1 we have

$$d^{-\alpha} n(\widehat{\tau}_1, \widehat{\tau}_2, \dots, \widehat{\tau}_n)^\top \xrightarrow{\mathsf{w}} (\ell_1, \ell_2, \dots, \ell_p, 0, \dots, 0)^\top$$

as $d \to \infty$, where $\ell_1 \ge \ell_2 \ge \cdots \ge \ell_p$ are the nonzero eigenvalues of a Wishart random matrix with distribution $W_n(p, cI_n)$.

We define

$$\begin{aligned} \mathsf{Angle}(v_j, \mathsf{span}\{u_i : i \in J\}) &= \arccos(\frac{v_j^\top [\mathsf{Proj}_{\mathsf{span}\{u_i : i \in J\}} v_j]}{\parallel v_j \parallel \parallel \mathsf{Proj}_{\mathsf{span}\{u_i : i \in J\}} v_j \parallel}) \\ &= \arccos(\frac{v_j^\top (\sum_{i \in J} (u_i^\top v_j) u_i)}{\parallel v_j \parallel \parallel \sum_{i \in J} (u_i^\top v_j) u_i \parallel}). \end{aligned}$$

Following the definition given in Jung and Marron (2009), we say that the sample eigenvector v_j , with $j \in J$, is *subspace consistent* if

$$\mathsf{Angle}(v_j,\mathsf{span}\{u_i:i\in J\})\overset{P}{\to}\mathsf{0}.$$

Under the same conditions of Theorem 3.1, let v_1, v_2, \ldots, v_p be the sample eigenvectors corresponding to the first p sample eigenvalues $\hat{\tau}_1 \geq \hat{\tau}_2 \geq \cdots \geq \hat{\tau}_p$. Then if $\epsilon > 0$

 $\lim_{n \to \infty} \lim_{d \to \infty} P(\operatorname{Angle}(v_i, \operatorname{span}\{e_j : j = 1, 2, \dots, p\}) > \epsilon) = 0, \quad (8)$

for i = 1, 2, ..., p.