

Rank tests for PCA under weak identifiability

Davy Paindaveine, Laura Peralvo Maroto and Thomas
Verdebout

Université libre de Bruxelles (ULB)



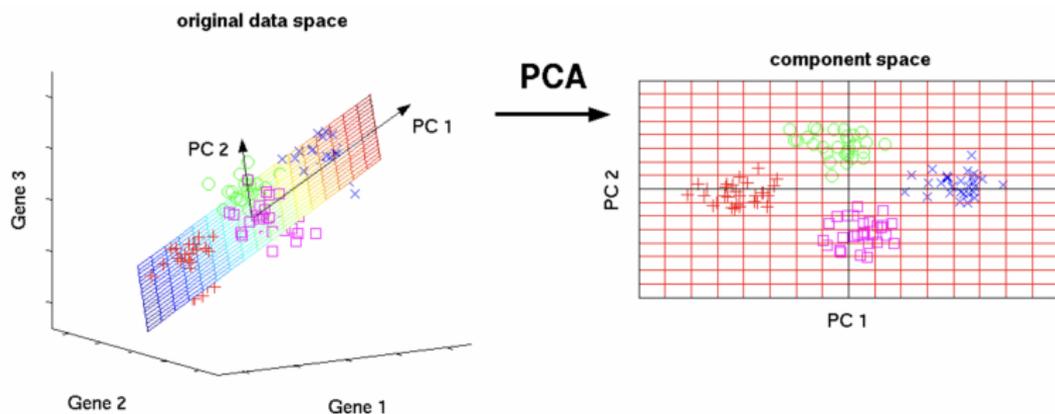
SPP : Mathematics PhD and Postdoc seminars, Bruxelles

- 1 PCA
- 2 Hypothesis testing and LAN
- 3 Weak identifiability
- 4 Weakly identifiable models
- 5 Asymptotic behavior of signed-rank procedures

- 1 PCA
- 2 Hypothesis testing and LAN
- 3 Weak identifiability
- 4 Weakly identifiable models
- 5 Asymptotic behavior of signed-rank procedures

Principal component analysis (PCA) is a classic multivariate statistical analysis technique.

↪ Objective: size reduction.



Let $\mathbf{X} = (X_1, \dots, X_p)$ be an observed p -dimensional random vector from the p variate distribution P with finite second-order moments and covariance matrix Σ .

We want to obtain ($S^{p-1} := \{\mathbf{u} \in \mathbb{R}^p, \mathbf{u}'\mathbf{u} = 1\}$)

$$\beta_1 := \operatorname{argmax}_{\beta \in S^{p-1}} \operatorname{Var}[\beta' \mathbf{X}] = \operatorname{argmax}_{\beta \in S^{p-1}} \beta' \Sigma \beta.$$

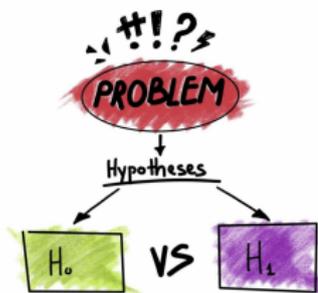
The first PC is then $Y_1 := \beta_1' \mathbf{X}$. Then,

$$\beta_2 := \operatorname{argmax}_{\beta \in S^{p-1}, \beta' \beta_1 = 0} \operatorname{Var}[\beta' \mathbf{X}].$$

The second PC is then $Y_2 := \beta_2' \mathbf{X}$.

It is well known that using the spectral decomposition of $\Sigma = \sum_{j=1}^p \lambda_j \theta_j \theta_j'$ ($\lambda_1 > \lambda_2 > \dots > \lambda_p$ are the eigenvalues and $\theta_1, \dots, \theta_p$ the associated eigenvectors of Σ), then $\beta_1 = \theta_1, \beta_2 = \theta_2, \dots$

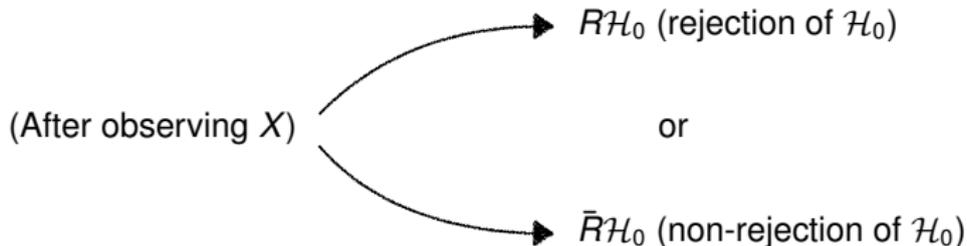
- 1 PCA
- 2 Hypothesis testing and LAN**
- 3 Weak identifiability
- 4 Weakly identifiable models
- 5 Asymptotic behavior of signed-rank procedures



Let X be an observation described by a statistical model $\{P_\theta : \theta \in \Theta\}$.

Consider a partition of Θ into $\Theta = \mathcal{H}_0 \oplus \mathcal{H}_1$ (\oplus denotes a disjoint union) where \mathcal{H}_0 is called the “null hypothesis” and \mathcal{H}_1 is the “alternative”.

A test problem is a decision problem in which only two decisions are possible.



Two mistakes can be made.

	$\theta \in \mathcal{H}_0$	$\theta \in \mathcal{H}_1$
$R\mathcal{H}_0$	Type I error (probability = α)	correct decision (probability = $1 - \beta$)
$\bar{R}\mathcal{H}_0$	correct decision probability = $1 - \alpha$	Type II error probability = β

→ The ideal test would be one that minimizes both first- and second-class risks.

Power of a test = $P[R\mathcal{H}_0] = 1 - \beta$ when $\theta \in \mathcal{H}_1$.

↔ we want to maximize the power.

Let f_{θ} be a density of P_{θ} with respect to some measure μ .

Definition

The model $(P_{\theta} : \theta \in \Theta)$ is called *differentiable in quadratic mean at θ* if there exists measurable functions $\dot{\ell}_{\theta}$ such that, as $\mathbf{h} \rightarrow 0$,

$$\frac{1}{\|\mathbf{h}\|^2} \int \left\{ f_{\theta+\mathbf{h}}^{1/2} - f_{\theta}^{1/2} - \frac{1}{2} \mathbf{h}' \dot{\ell}_{\theta} f_{\theta}^{1/2} \right\}^2 d\mu = o(1).$$

Definition

The sequence of statistical models $(P_{\theta}^{(n)} : \theta \in \Theta)$ is *locally asymptotically normal (LAN)* at θ if there exist matrices \mathbf{r}_n and Γ_{θ} and random vectors $\Delta_{\theta}^{(n)}$ such that, under $P_{\theta}^{(n)}$, for every converging sequence $\mathbf{h}_n \rightarrow \mathbf{h}$,

$$\log \frac{dP_{\theta+\mathbf{r}_n^{-1}\mathbf{h}_n}^{(n)}}{dP_{\theta}^{(n)}} = \mathbf{h}' \Delta_{\theta}^{(n)} - \frac{1}{2} \mathbf{h}' \Gamma_{\theta} \mathbf{h} + o_{\mathbb{P}}(1) \quad \text{and} \quad \Delta_{\theta}^{(n)} \xrightarrow{\mathcal{L}} \mathcal{N}(\mathbf{0}, \Gamma_{\theta}).$$

Note : if the experiment $(P_{\theta} : \theta \in \Theta)$ is differentiable in quadratic mean, then the sequence of model $(P_{\theta}^{(n)} : \theta \in \Theta)$ is LAN (with norming matrices $r_n = \sqrt{n}I$).

- 1 PCA
- 2 Hypothesis testing and LAN
- 3 Weak identifiability**
- 4 Weakly identifiable models
- 5 Asymptotic behavior of signed-rank procedures

Testing problem: throughout the presentation, we consider the following problem test :

$$\begin{cases} \mathcal{H}_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0 \\ \mathcal{H}_1 : \boldsymbol{\theta} \neq \boldsymbol{\theta}_0 \end{cases}$$

where $\boldsymbol{\theta}$ is the eigenvector associated with the largest eigenvalue λ_1 of the underlying covariance matrix and $\boldsymbol{\theta}_0$ is a given unit vector of \mathbb{R}^p .

We will consider situations where $\lambda_1 - \lambda_2$ is small, that is a situation of **weak identifiability** of $\boldsymbol{\theta}$.

Objective : Paindaveine, Remy and Verdebout (2020) studied purely Gaussian procedures for this problem in the Gaussian case.

↪ the objective of this work is to extend the results of Paindaveine, Remy and Verdebout (2020) to

- (i) the elliptical case;
- (ii) signed-rank procedures.

- 1 PCA
- 2 Hypothesis testing and LAN
- 3 Weak identifiability
- 4 Weakly identifiable models**
- 5 Asymptotic behavior of signed-rank procedures

We consider triangular arrays of elliptically symmetric observations \mathbf{X}_{ni} , $i = 1, \dots, n$, $n = 1, 2, \dots$ where $\mathbf{X}_{n1}, \dots, \mathbf{X}_{nn}$ form an observed n -tuple of mutually independent p -dimensional random vectors with probability density function of the form

$$f_{\sigma_n^2, \mathbf{V}_n, f_1}(\mathbf{x}) := c_{p, f_1} \frac{1}{\sigma_n^p |\mathbf{V}_n|^{1/2}} f_1 \left(\frac{1}{\sigma_n} (\mathbf{x}' \mathbf{V}_n^{-1} \mathbf{x})^{1/2} \right), \quad \mathbf{x} \in \mathbb{R}^p, \quad (1)$$

where, $\boldsymbol{\Sigma}_n := \mathbf{I}_p + r_n \mathbf{v} \boldsymbol{\theta} \boldsymbol{\theta}'$, $\sigma_n := |\boldsymbol{\Sigma}_n|^{1/(2p)} = (1 + r_n \nu)^{1/2p}$ is a *scale parameter* in $\mathbb{R}_0^+ := (0, \infty)$,

$$\mathbf{V}_n := \boldsymbol{\Sigma}_n / \sigma_n^2 := (1 + r_n \nu)^{-1/p} (\mathbf{I}_p + r_n \mathbf{v} \boldsymbol{\theta} \boldsymbol{\theta}')$$

is a *shape parameter* with eigenvalues

$$\lambda_{n1, \mathbf{V}_n} = (1 + r_n \nu)^{(\rho-1)/p} \quad \text{et} \quad \lambda_{n2, \mathbf{V}_n} = \dots = \lambda_{np, \mathbf{V}_n} = (1 + r_n \nu)^{-1/p},$$

où r_n et ν are positive real numbers and $f_1 : \mathbb{R}_0^+ \rightarrow \mathbb{R}^+ := [0, \infty)$ is a *standardized radial density* ($f_1 \in \mathcal{F}_1$).

The resulting hypothesis will be denoted as $P_{\theta, r_n, \nu, f_1}^{(n)}$.

Examples :

- (i) the p -variate multinormal distribution, with radial density

$$f_1(r) = \phi_1(r) := \exp(-a_p r^2/2);$$

- (ii) the p -variate Student distributions, with radial densities (for $\nu \in \mathbb{R}_0^+$ degrees of freedom) $f_1(r) = f_{1,\nu}^t(r) := (1 + a_{p,\nu} r^2/\nu)^{-(p+\nu)/2}$;

where the positive constants a_p and $a_{p,\nu}$ are such that $f_1 \in \mathcal{F}_1$.

We denote by

$$d_{ni} := d_{ni}(\mathbf{V}_n) := \|\mathbf{V}_n^{-1/2} \mathbf{X}_{ni}\| \quad \text{and} \quad \mathbf{U}_{ni} := \mathbf{U}_{ni}(\mathbf{V}_n) := \mathbf{V}_n^{-1/2} \mathbf{X}_{ni}/d_{ni}$$

respectively the *standardized elliptical distances* and the *multivariate signs*, $i = 1, \dots, n$.

We studied likelihood ratio of the form

$$\Lambda_n := \log \frac{dP_{\boldsymbol{\theta}_0 + \nu_n \boldsymbol{\tau}_n, r_n, \nu, f_1}^{(n)}}{dP_{\boldsymbol{\theta}_0, r_n, \nu, f_1}^{(n)}}$$

for some admissible perturbations $\boldsymbol{\tau}_n$, where (ν_n) is a positive real sequence.

Some sequences (r_n) does not provide LAN (locally and asymptotically normal) experiments!

We first obtain a general result to derive the asymptotic behavior of Λ_n . In the following result, $f_{\boldsymbol{\vartheta}_n}$ is a density associated with a certain distribution $P_{\boldsymbol{\vartheta}_n, f}$ with respect to a dominated measure μ and $\mathbf{Z}_{n1}, \dots, \mathbf{Z}_{nn}$ form a random sample from the distribution $P_{\boldsymbol{\vartheta}_n, f}$.

Proposition

Let (ν_n) be a positive real sequence and suppose that for all sequence $(\vartheta_n) \in \Theta \subset \mathbb{R}^p$, there exists real valued functions $\dot{\ell}_{\vartheta_n, \tau_n}^{(n)}$, where τ_n is a bounded sequence in \mathbb{R}^p , such that

$$\int \left(f_{\vartheta_n + \nu_n \tau_n}^{1/2}(\mathbf{z}) - f_{\vartheta_n}^{1/2}(\mathbf{z}) - \frac{1}{2} \dot{\ell}_{\vartheta_n, \tau_n}^{(n)}(\mathbf{z}) f_{\vartheta_n}^{1/2}(\mathbf{z}) \right)^2 d\mu(\mathbf{z}) = o(n^{-1})$$

as $n \rightarrow \infty$. Assume furthermore that

$$\mathbb{E}_{\mathbb{P}_{\vartheta_n, f}^{(n)}} [n \dot{\ell}_{\vartheta_n, \tau_n}^{(n)}(\mathbf{Z}_{n1})] = o(1), \quad \Gamma_{\vartheta_n, \tau_n}^{(n)} := \mathbb{E}_{\mathbb{P}_{\vartheta_n, f}^{(n)}} [n (\dot{\ell}_{\vartheta_n, \tau_n}^{(n)}(\mathbf{Z}_{n1}))^2] = O(1),$$

$$\left(\sum_{i=1}^n (\dot{\ell}_{\vartheta_n, \tau_n}^{(n)}(\mathbf{Z}_{ni}))^2 \right) - \Gamma_{\vartheta_n, \tau_n}^{(n)} = o_{\mathbb{P}_{\vartheta_n, f}^{(n)}}(1)$$

and

$$\mathbb{E}_{\mathbb{P}_{\vartheta_n, f}^{(n)}} [n (\dot{\ell}_{\vartheta_n, \tau_n}^{(n)}(\mathbf{Z}_{n1}))^2 \mathbb{I}[n (\dot{\ell}_{\vartheta_n, \tau_n}^{(n)}(\mathbf{Z}_{n1}))^2 \geq n\epsilon^2]] = o(1)$$

as $n \rightarrow \infty$. Then, as $n \rightarrow \infty$ under $\mathbb{P}_{\vartheta_n, f}^{(n)}$,

$$\log \frac{d\mathbb{P}_{\vartheta_n + \nu_n \tau_n, f}^{(n)}}{d\mathbb{P}_{\vartheta_n, f}^{(n)}} = \sum_{i=1}^n \dot{\ell}_{\vartheta_n, \tau_n}^{(n)}(\mathbf{Z}_{ni}) - \frac{1}{2} \Gamma_{\vartheta_n, \tau_n}^{(n)} + o_{\mathbb{P}}(1).$$

We need some further mild regularity conditions on f_1 in the sequel.

Define $\varphi_{f_1} := -\dot{f}_1/f_1$ where \dot{f}_1 is the a.e. derivative of f_1 .

In the next result, f_1 is supposed to be in the collection of all absolutely continuous radial standardized densities \mathcal{F}_1^a for which

$$\mathcal{J}_\rho(f_1) := E[\varphi_{f_1}^2(d_{ni}/\sigma)(d_{ni}/\sigma)^2] < \infty.$$

The following theorem describes the asymptotic behavior of Λ_n in the following four regimes:

- (i) $r_n \equiv 1$: *away from contiguity*,
- (ii) $r_n = o(1)$ with $\sqrt{nr_n} \rightarrow \infty$: *above contiguity*,
- (iii) $r_n = 1/\sqrt{n}$: *under contiguity*,
- (iv) $r_n = o(1/\sqrt{n})$: *under strict contiguity*.

Theorem

(i) (away from contiguity) if $r_n \equiv 1$, then, with $\nu_n = 1/\sqrt{n}$,

$$\Delta_{f_1}^{(n)} := \frac{\nu}{\sqrt{1+\nu}} \sqrt{n} (\mathbf{I}_p - \boldsymbol{\theta}_0 \boldsymbol{\theta}_0') \left(\frac{1}{n} \sum_{i=1}^n \varphi_{f_1} \left(\frac{d_{ni}}{\sigma_n} \right) \frac{d_{ni}}{\sigma_n} \mathbf{U}_{ni} \mathbf{U}_{ni}' - \mathbf{I}_p \right) \boldsymbol{\theta}_0$$

and

$$\Gamma_{f_1} = \frac{\mathcal{J}_p(f_1) \nu^2}{p(p+2)(1+\nu)} (\mathbf{I}_p - \boldsymbol{\theta}_0 \boldsymbol{\theta}_0'),$$

we have that, under $\mathbb{P}_{\boldsymbol{\theta}_0, r_n, \nu, f_1}^{(n)}$,

$$\Lambda_n = \boldsymbol{\tau}_n' \Delta_{f_1}^{(n)} - \frac{1}{2} \boldsymbol{\tau}_n' \Gamma_{f_1} \boldsymbol{\tau}_n + o_{\mathbb{P}}(1)$$

and that $\Delta_{f_1}^{(n)}$ is asymptotically normal with mean zero and covariance matrix Γ_{f_1} ;

Theorem

(ii) (above contiguity) if r_n is $o(1)$ with $\sqrt{nr_n} \rightarrow \infty$, then, with $\nu_n = 1/(\sqrt{nr_n})$,

$$\Delta_{f_1}^{(n)} := \nu\sqrt{n}(\mathbf{I}_p - \boldsymbol{\theta}_0\boldsymbol{\theta}_0') \left(\frac{1}{n} \sum_{i=1}^n \varphi_{f_1} \left(\frac{d_{ni}}{\sigma_n} \right) \frac{d_{ni}}{\sigma_n} \mathbf{U}_{ni} \mathbf{U}_{ni}' - \mathbf{I}_p \right) \boldsymbol{\theta}_0$$

and

$$\Gamma_{f_1} = \frac{\mathcal{J}_p(f_1)\nu^2}{\rho(\rho+2)} (\mathbf{I}_p - \boldsymbol{\theta}_0\boldsymbol{\theta}_0'),$$

we have that, under $\mathbb{P}_{\boldsymbol{\theta}_0, r_n, \nu, f_1}^{(n)}$,

$$\Lambda_n = \boldsymbol{\tau}_n' \Delta_{f_1}^{(n)} - \frac{1}{2} \boldsymbol{\tau}_n' \Gamma_{f_1} \boldsymbol{\tau}_n + o_{\mathbb{P}}(1)$$

and that $\Delta_{f_1}^{(n)}$ is asymptotically normal with mean zero and covariance matrix Γ_{f_1} ;

Theorem

(iii) (under contiguity) if $r_n = 1/\sqrt{n}$, then, letting $\nu_n \equiv 1$,

$$\Lambda_n = \boldsymbol{\tau}'_n \left[\nu \sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n \varphi_{f_1} \left(\frac{d_{ni}}{\sigma_n} \right) \frac{d_{ni}}{\sigma_n} \mathbf{U}_{ni} \mathbf{U}'_{ni} - \mathbf{I}_p \right) \left(\boldsymbol{\theta}_0 + \frac{1}{2} \boldsymbol{\tau}_n \right) \right] \\ - \frac{\mathcal{J}_p(f_1) \nu^2}{\rho(\rho+2)} \left(\frac{\|\boldsymbol{\tau}_n\|^2}{2} - \frac{\|\boldsymbol{\tau}_n\|^4}{8} \right) + o_P(1),$$

under $P_{\boldsymbol{\theta}_0, r_n, \nu, f_1}^{(n)}$, where, if $\boldsymbol{\tau}_n \rightarrow \boldsymbol{\tau}$, then

$$\boldsymbol{\tau}'_n \nu \sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n \varphi_{f_1} \left(\frac{d_{ni}}{\sigma_n} \right) \frac{d_{ni}}{\sigma_n} \mathbf{U}_{ni} \mathbf{U}'_{ni} - \mathbf{I}_p \right) \left(\boldsymbol{\theta}_0 + \frac{1}{2} \boldsymbol{\tau}_n \right)$$

is asymptotically normal with mean zero and covariance matrix

$$\frac{\mathcal{J}_p(f_1)}{\rho(\rho+2)} \left(\|\boldsymbol{\tau}\|^2 - \frac{\|\boldsymbol{\tau}\|^4}{4} \right);$$

Theorem

(iv) (under strict contiguity) if $r_n = o(1/\sqrt{n})$, then, even with $\nu_n \equiv 1$, we have that $\Lambda_n = o_P(1)$ under $P_{\theta_0, r_n, \nu, f_1}^{(n)}$.

It follows that, for any fixed $\nu > 0$ and for any fixed sequence (r_n) associated with regime (i) or regime (ii), the sequence of models is LAN with central sequence

$$\Delta_{\delta, f_1}^{(n)} := \frac{\sqrt{n\nu}}{\sqrt{1 + \delta\nu}} (\mathbf{I}_p - \boldsymbol{\theta}_0 \boldsymbol{\theta}'_0) \left(\frac{1}{n} \sum_{i=1}^n \varphi_{f_1} \left(\frac{d_{ni}}{\sigma_n} \right) \frac{d_{ni}}{\sigma_n} \mathbf{U}_{ni} \mathbf{U}'_{ni} - \mathbf{I}_p \right) \boldsymbol{\theta}_0$$

and Fisher information matrix

$$\Gamma_{\delta, f_1} = \frac{\mathcal{J}_p(f_1) \nu^2}{p(p+2)(1 + \delta\nu)} (\mathbf{I}_p - \boldsymbol{\theta}_0 \boldsymbol{\theta}'_0)$$

where $\delta := 1$ if regime (i) is considered and $\delta := 0$ otherwise.

- 1 PCA
- 2 Hypothesis testing and LAN
- 3 Weak identifiability
- 4 Weakly identifiable models
- 5 Asymptotic behavior of signed-rank procedures**

Below, we denote by $R_{ni}(\mathbf{V}_n)$ the rank of $d_{ni}(\mathbf{V}_n)$ among $d_{n1}(\mathbf{V}_n), \dots, d_{nn}(\mathbf{V}_n)$.

The rank-based test $\phi_K = \phi_K^{(n)}$ of Hallin, Paindaveine and Verdebout (2010) rejects the null hypothesis (at asymptotic level α) when

$$Q_K := \frac{np(p+2)}{\mathcal{J}_p(K)} \sum_{j=2}^p (\tilde{\boldsymbol{\theta}}_j' \mathbf{S}_K^{(n)} \boldsymbol{\theta}_0)^2 > \chi_{p-1; 1-\alpha}^2$$

where $\mathcal{J}_p(K)$ is a constant, $\tilde{\boldsymbol{\theta}}_j$ stands for a constrained estimator of $\hat{\mathbf{V}}_{\text{Tyler}}$'s j th eigenvector for the shape estimator $\hat{\mathbf{V}}_{\text{Tyler}}$ of Tyler (1987) and the signed-rank covariance matrix is of the form

$$\mathbf{S}_K^{(n)} := \frac{1}{n} \sum_{i=1}^n K\left(\frac{R_{ni}}{n+1}\right) \mathbf{U}_{ni} \mathbf{U}_{ni}',$$

with $K : (0, 1) \rightarrow \mathbb{R}$ stands for some *score function*, $\mathbf{U}_{ni} = \mathbf{U}_{ni}(\tilde{\boldsymbol{\theta}}_0 \hat{\Lambda}_{\text{Tyler}} \tilde{\boldsymbol{\theta}}_0')$ et $R_{ni} = R_{ni}(\tilde{\boldsymbol{\theta}}_0 \hat{\Lambda}_{\text{Tyler}} \tilde{\boldsymbol{\theta}}_0')$ with $\tilde{\boldsymbol{\theta}}_0 := (\boldsymbol{\theta}_0, \tilde{\boldsymbol{\theta}}_2, \dots, \tilde{\boldsymbol{\theta}}_p)$.

Proposition

Fix a unit p -vector θ_0 , $v > 0$ and $g_1 \in \mathcal{F}_1$. Then, for any sequence r_n , Q_K converges weakly to a chi-square random variable with $p - 1$ degrees of freedom under $P_{\theta_0, r_n, v, g_1}^{(n)}$.

↪ ϕ_K is robust to weak identifiability.

↪ Impact on asymptotic efficiency properties?

Proposition

Fix a unit p -vector θ_0 , $v > 0$, $g_1 \in \mathcal{F}_1$. Then,

- (a) when (i) $r_n \equiv 1$ ($\delta = 1$) or (ii) $r_n = o(1)$ with $\sqrt{nr_n} \rightarrow \infty$ ($\delta = 0$), we have that under $P_{\theta_0 + \tau_n / (\sqrt{nr_n}), r_n, v, g_1}^{(n)}$, the test statistic Q_K converges weakly to a chi-square random variable with $p - 1$ degrees of freedom and non-centrality parameter

$$\frac{\mathcal{J}_p^2(K, g_1) v^2}{\mathcal{J}_p(K) p(p+2)(1+\delta v)} \|\tau\|^2,$$

where $\tau := \lim_{n \rightarrow \infty} \tau_n$;

Proposition

(b) under $P_{\theta_0 + \tau_n, 1/\sqrt{n}, v, g_1}^{(n)}$, Q_K converges weakly to a chi-square random variable with $p - 1$ degrees of freedom and with non-centrality parameter

$$\frac{v^2 \mathcal{J}_p^2(K, g_1)}{16 \mathcal{J}_p(K) p(p+2)} \|\tau\|^2 (4 - \|\tau\|^2) (2 - \|\tau\|^2)^2; \quad (2)$$

(c) when $r_n = o(1)$ with $\sqrt{n}r_n \rightarrow 0$, we have that under $P_{\theta_0 + \tau_n, r_n, v, g_1}^{(n)}$, Q_K converges weakly to a chi-square random variable with $p - 1$ degrees of freedom.

Therefore :

$$\text{ARE}_{\theta_0, r_n, v, g_1}(\phi_K / \phi_G) := \frac{(1 + \kappa_p(g_1)) \mathcal{J}_p^2(K, g_1)}{p(p+2) \mathcal{J}_p(K)}$$

↔ not affected by weak identifiability.

Simulations

First simulation exercise : for any $b = 0, 1, \dots, 5$, we generate $M = 2500$ mutually independent random samples $\mathbf{X}_i^{(b,r)}$, $i = 1, \dots, n = 200$, from the $p = 3$ -variate t_1 ($r = 1$), t_5 ($r = 2$) and normal ($r = 3$) distributions, with mean zero and scatter matrix

$$\Sigma_n^{(b)} := \mathbf{I}_p + n^{-b/6} \boldsymbol{\theta}_0 \boldsymbol{\theta}_0'$$

where $\boldsymbol{\theta}_0 = (1, 0, 0)'$. This covers regimes (i) ($b = 0$), (ii) ($b = 1, 2$), (iii) ($b = 3$) and (iv) ($b = 4, 5$).

We perform, at nominal level 5%, the Wilcoxon test ϕ_{K_1} and the van der Waerden test ϕ_{K_Φ} for $\mathcal{H}_0^{(n)} : \boldsymbol{\theta} = \boldsymbol{\theta}_0$.

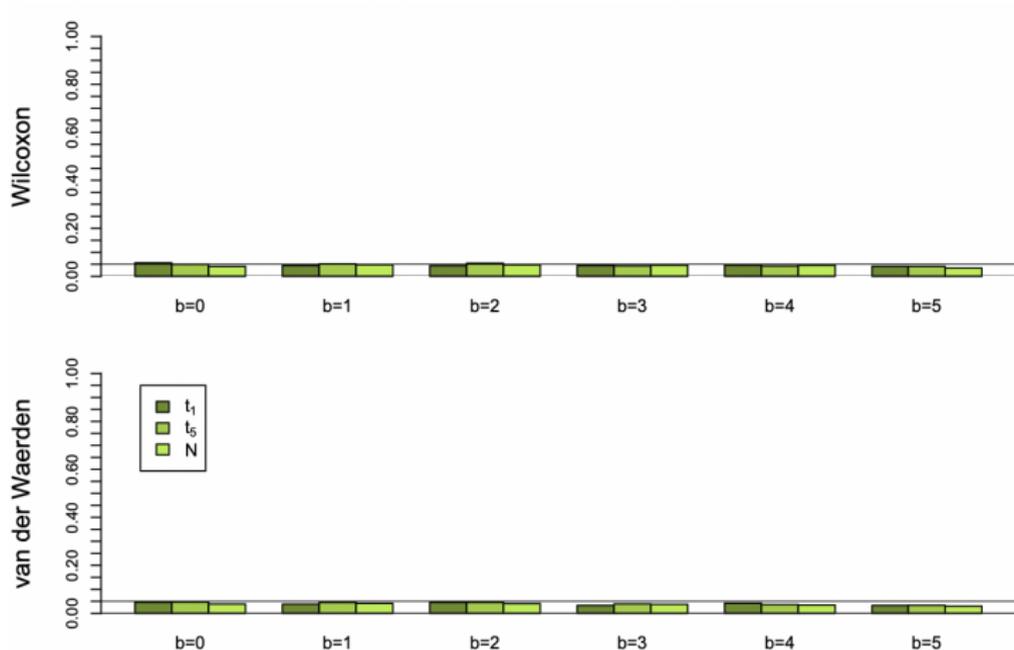


Figure: Empirical rejection frequencies, under the null hypothesis, of the Wilcoxon test ϕ_{K_1} and the van der Waerden test ϕ_{K_ϕ} , performed at nominal level 5%.

Simulations

Second simulation exercise : we generate $M = 100\,000$ mutually independent random samples $\mathbf{X}_i^{(\ell)}$, $i = 1, \dots, n = 10\,000$, $\ell = 0, 1, \dots, 20$, from the $p = 2$ -variate t_1 distribution with mean zero and scatter matrix

$$\Sigma_n^{(\ell)} := \mathbf{I}_p + n^{-1/2}(\boldsymbol{\theta}_0 + \boldsymbol{\tau}_\ell)(\boldsymbol{\theta}_0 + \boldsymbol{\tau}_\ell)',$$

where $\boldsymbol{\theta}_0 = (1, 0)'$ et $\boldsymbol{\theta}_0 + \boldsymbol{\tau}_\ell = (\cos(\ell\pi/40), \sin(\ell\pi/40))'$.

The value $\ell = 0$ is associated to the null hypothesis, while values $\ell = 1, \dots, 20$ provide increasingly severe alternatives.

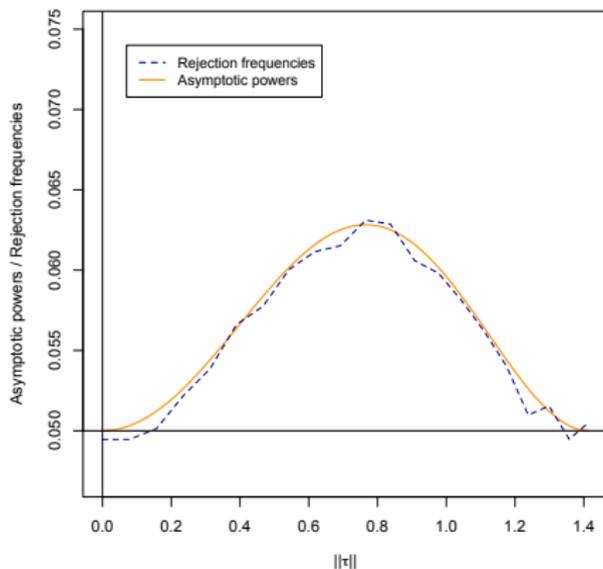


Figure: Empirical rejection frequencies, under the null hypothesis and local alternatives, of the van der Waerden test (ϕ_{K_ϕ}), performed at nominal level 5%.

References

- Hallin, M., Paindaveine, D. and Verdebout, Th. (2010). Optimal rank-based testing for principal components. *Annals of Statistics*, 38, 3245–3299.
- Paindaveine, D., Remy, J. and Verdebout, Th. (2020). Testing for principal component directions under weak identifiability. *Annals of Statistics*, 48, 324-345.
- Tyler, D. E. (1987). A distribution-free M-estimator of multivariate component vectors. *Annals of Statistics*, 11, 1243-1250.