

# Detecting independence via Lévy metrics

(+)	B. Böttcher	2018	KSTA
	M. Kalter-Ressel	2019	AoS
(+)	B. Böttcher	2020	Open Stats.

Key Words

distance covariance  
distance multivariance  
Székely, Rizzo, Bakirov  
 $\approx 2007/2009$   
2 r.v.

$n$  r.v.  
RS et al.

## § Problem

- (+) measure (in-)dependence of data sets  $X, Y$
- (-) high dim  $\Rightarrow$  computation!
- (-)  $\dim X \neq \dim Y$

classical correlation

Pearson corr.

$\text{COV} = 0 \nRightarrow \perp$   
2 r.v only  
same dim

Szebelj:

$$1 \stackrel{=}{\circ} f_Z(u) = \mathbb{E} e^{iuZ} \text{ dkw fn.}$$

$$f_{(X,Y)}(s,t) - f_X(s)f_Y(t) \stackrel{!}{=} 0$$

$$2 \stackrel{=}{\circ} V^2(X, Y) = \iint |f_{XY}(s,t) - f_X(s)f_Y(t)|^2 \frac{\omega(s,t) ds dt}{\text{weight}}$$

$$\omega(s,t) = C_{\alpha,m} |s|^{-\alpha-m} C_{\alpha,n} |t|^{-\alpha-n}$$

Lévy measure,  $\alpha$ -stable

$$3 \stackrel{=}{\circ} V^2(X, Y) = \mathbb{E} |X - X'|^\alpha |Y - Y'|^\alpha + \mathbb{E} |X - X'|^\alpha \mathbb{E} |Y - Y'|^\alpha - 2 \mathbb{E} |X - \tilde{X}|^\alpha |Y - \tilde{Y}|^\alpha$$

$(X', Y')$ ,  $(X, Y)$  iid copies

$$\tilde{X} \sim X, \tilde{Y} \sim Y \oplus \tilde{X} \perp \tilde{Y}$$

$$4 \stackrel{=}{\circ} V^2(X, Y) < \infty \text{ if } \mathbb{E} |X|^\alpha, \mathbb{E} |Y|^\alpha < \infty$$

$$V^2(X, Y) = 0 \iff X \perp \tilde{Y}$$

$\exists$  estimator for " $\perp \tilde{Y}$ "

$$V_N^2(\tilde{X}, \tilde{Y}) = \bar{T}_1 + \bar{T}_2 - \bar{T}_3$$

$$\bar{T}_1 = \frac{1}{N^2} \sum_{k,e} |x_k - x_e|^\alpha |y_k - y_e|^\alpha$$

$$\bar{T}_2 = \frac{1}{N^2} \frac{1}{N^2} \sum_{k,e} |x_k - x_e|^\alpha \sum_{e,e'} |y_e - y_{e'}|^\alpha$$

$$T_3 = \frac{1}{N^3} \sum_{k,l,m} |x_k - x_l|^\alpha |y_k - y_l|^\alpha$$

$\exists$  good properties: consistent  
+ approx. properties

§ First generalization  $| \cdot |^\alpha \rightsquigarrow \psi(\cdot)$

$C_{\alpha,m}$  is  $|\cdot|^{-\alpha-m}$  Le'vy measure

$$|x|^\alpha = C_{\alpha,m} \int (1 - \cos xs) \frac{ds}{|s|^{\alpha+m}}, \quad x \in \mathbb{R}^m$$

Le'vy-Khintchine

$$\psi(x) = \int (1 - \cos xs) \mu(ds) + \frac{1}{2} x^T Q x \geq 0$$

$\uparrow$   $\int s^2 \lambda(s) ds < \infty$

$m \times m$   
matrix  
pos. semidef.

chow. exponent of a Le'vy process (symm.)

continuous neg. definite function

Schoenberg 1938

- $d_\psi(x, x') := \sqrt{\psi(x-x')}$

metric in  $\mathbb{R}^m$

- $\|\psi\|_\infty < \infty \iff \psi(\mathbb{R}^m) \text{ bounded}$

$$Q \equiv 0$$

BB

- $E \psi(X-z) = E \psi(\tilde{X}-z) \quad \forall z$   
 $\Rightarrow X \sim \tilde{X}$  fpr  $\psi$  is  
"characterizing"  
 $\downarrow$   
 $\text{supp } \mu = \mathbb{R}^n \setminus \{0\}$

Def  $X \in \mathbb{R}^m$ ,  $Y \in \mathbb{R}^n$ ,  $\mu, \nu$  Le'vy meas.  
cntrf  $\psi, \varphi$ . Set

$$V^2(X, Y) = \| f_{(X,Y)} - f_X \otimes f_Y \|_{L^2(\mu \otimes \nu)}^2 \in [0, \infty]$$

Then a)  $V^2(X, Y) < \infty$  if  $E \psi(X), E \varphi(Y) < \infty$   
 $\leq 16 E \psi(X) E \varphi(Y)$

b)  $\text{supp } \mu = \mathbb{R}^n \setminus \{0\}$   
 $\text{supp } \nu = \mathbb{R}^n \setminus \{0\} \Rightarrow \begin{bmatrix} V^2(X, Y) = 0 \\ \iff X \perp\!\!\!\perp Y \end{bmatrix}$

$\hookrightarrow V^2(X, Y) = E \psi(X_1 - X_4) \varphi(Y_1 - Y_4)$

$$+ E \psi(X_1 - X_4) E \varphi(Y_1 - Y_4)$$

$\circlearrowleft -2 E \psi(X_1 - X_2) \varphi(Y_1 - Y_3)$

$(X_i, Y_i)$  iid copies of  $(X, Y)$

here we need finiteness,  $\infty - \infty$

### § Estimators

$\vec{x}, \vec{y}$  empirical r.v's

samples  $(x_i, y_i)_{i=1, \dots, N}$

$$V_N^2(\vec{x}, \vec{y}) = \frac{1}{N^2} \sum_{i, \ell=1}^N \psi(x_i - x_\ell) \varphi(y_i - y_\ell)$$

$$+ \frac{1}{N^2} \sum_{i, j} \underbrace{\psi(x_i - x_j)}_{a_N} \underbrace{\varphi(y_i - y_j)}_{b_N} \frac{1}{N^2} \sum_{\ell, c} \varphi(y_\ell - y_c)$$

$$- \frac{2}{N^3} \sum_{i, j, m} \psi(x_i - x_j) \varphi(y_i - y_m)$$

④ alpha

$$\Rightarrow c_N \sum_{ij} + c'_N \sum_{ije} + c''_N \sum_{ije,k}$$

$\underbrace{\phantom{c'_N \sum_{ije}}}_{\text{distinct}} \quad \uparrow$

TJ-statistics!

Lemma  $V_N^2(\vec{x}, \vec{y}) = \frac{1}{N^2} \text{trace } (\mathcal{B}^T \mathcal{A})$

$$\mathcal{A} = C^T a C, \quad \mathcal{B} = C^T b C$$

$$a = (-\psi(x_e - x_\ell))_{\ell, e}$$

$$b = (-\varphi(y_e - y_\ell))_{\ell, e}$$

distance matrices

$$C = \begin{pmatrix} 1 - \frac{1}{N} & -\frac{1}{N} \\ -\frac{1}{N} & 1 - \frac{1}{N} \end{pmatrix}$$

Theorem 2  $\mathbb{E} \varphi(X) < \infty, \mathbb{E} \varphi(Y) < \infty$

$$NV_N^2(\vec{x}, \vec{y}) \xrightarrow[\text{a.s.}]{} V^2(x, y)$$

consistency.

Theorem 3  $\mathbb{E} \varphi(X) < \infty, \mathbb{E} \varphi(Y) < \infty$

$$\mathbb{E} \log(1 + \|X\|^2)^{1+\varepsilon} < \infty$$

$$\mathbb{E} \log(1 + \|Y\|^2)^{1+\varepsilon} < \infty$$

$$X \perp \!\!\! \perp Y \Rightarrow NV_N^2(\vec{x}, \vec{y}) \xrightarrow{d} \iint |G_7(s, t)|^2 \mu(ds) \nu(dt)$$

of centered Gaussian field.

Corollary  $\mathbb{E} \varphi(X) < \infty, \mathbb{E} \varphi(Y) < \infty$

a)  $X \perp \!\!\! \perp Y + \log \text{constant}$

$$\Rightarrow \frac{NV_N^2}{a_N b_N} \xrightarrow{d} Q$$

late by  
BB

b)  $X \not\perp \!\!\! \perp Y : \frac{NV_N^2}{a_N b_N} \xrightarrow{P} \infty$

### § More Than 2 r.v. multivariate

$f = f_1 \otimes \dots \otimes f_n$  in  $L^2$  by meas.

in  $\mathbb{R}^{d_1}, \dots, \mathbb{R}^{d_n}$ ,  $\psi_1, \dots, \psi_n$

$$M_p^2(x_1, \dots, x_n) = \left\| \mathbb{E} \prod_{k=1}^n (e^{i x_k} - f_{X_k}(\cdot)) \right\|_{L^2(p)}^2$$

Issue 1  $M_p^2 = 0 \not\Rightarrow x_1, \dots, x_n$  indep.

$\exists x_1 \equiv c, x_2, \dots, x_n$   
any choice

Issue 2  $x_1, \dots, x_n \Leftarrow$  from  $H \{i_1, \dots, i_m\}, 1 \leq i_1 < \dots < i_m$   
indep.  $x_{i_1}, \dots, x_{i_m}$  indep.  
 $m$ -independence

$$\binom{n}{m} \approx 2^n \text{ choices!}$$

total distance covariance

$$\widetilde{M}_p^2(x_1, \dots, x_n) = \sum_{\substack{1 \leq i_1 < \dots < i_m \leq n \\ 2 \leq m \leq n}} M_m^2 \circ \otimes_{j=1}^m (x_{i_1}, \dots, x_{i_m})$$

$$\left\| \prod_{k=1}^n (1 + x_k) \right\|$$

$$= 1 + x_1 + \dots + x_n + \sum_{i_1, \dots, i_m} x_{i_1} \dots x_{i_m}$$

## Theorem 4 TFAE

(1)  $X_1, \dots, X_n$  iid.

$$(2) \quad \overline{M}_\phi(X_1, \dots, X_n) = 0$$

BB 2020  
 $n=1$

Theorem 5  $f_i \Leftrightarrow \psi_i$ ,  $E[\psi_i(X_i)] < \infty$

$$\overline{M}_\phi^2(X_1, \dots, X_n) = E \prod_{i=1}^n (1 + \psi_i(X_i, X'_i))$$

$$\psi_i(X_i, X'_i) =$$

$$= -\psi(X_i - X'_i) + E(\psi(X_i - X'_i) | X_i)$$

$$\underbrace{\psi_i(X_i - X'_i)}_{X_i \text{ not id, any } n} + E(\psi(X_i - X'_i) | X'_i) - E\psi(X_i - X'_i)$$

$(X_1, \dots, X_n), (X'_1, \dots, X'_n)$  iid

$\Rightarrow \exists$  consistent estimator

$\exists$  computationally fast

$\exists$  same failure as in  $n=2$

$\exists$  detect dependence clusters

i.e.

$(X_1, \dots, X_n)$  dependent

$\rightsquigarrow 33$  2020.

$$|x_i - x_j|^{\frac{2}{\alpha}}$$

$$\psi(x_i - x_j)$$