

Measuring Association Between Random Vectors

Oliver Grothe¹ Julius Schnieders¹ Johan Segers²

¹University of Cologne
Department of Economic and Social Statistics

²Université catholique de Louvain
Institut de statistique, biostatistique et sciences actuarielles

ESSEC La Défense, May 19, 2014

Association as a particular form of dependence

Positive/Negative association between random variables:

- ▶ Large values of one variable tend to go together with large/small values of the other one
- ▶ Particular form of *dependence*
 - ▶ If X is symmetrically distributed around 0 and if $Y = X^2$, then X and Y are perfectly dependent but not associated.

Of obvious practical interest:

- ▶ risk concentration/diversification

Quantifying association via linear correlation: the marginal distributions come into play

Example

If

$$X \sim N(0, 1), \quad Y = X^2$$

then

$$\begin{aligned}\text{cov}(X, Y) &= 0 \\ \text{cov}(e^X, Y) &= \sqrt{e} > 0 \\ \text{cov}(-e^{-X}, Y) &= -\sqrt{e} < 0\end{aligned}$$

Still, both functions $f_1(x) = e^x$ and $f_2(x) = -e^{-x}$ are increasing:
what is “large” and what is “small”?

Reasons to measure association in a margin-free way

We might want to measure association in a **margin-free** way if:

- ▶ The state spaces of X and Y are not directly comparable
 - ▶ Decathlon: high jump versus 100 metres
- ▶ Outliers have a disproportionately large influence
 - ▶ Very large positive or negative returns on asset prices

Margin-free: invariant w.r.t. increasing transformations

Examples:

- ▶ SPEARMAN's rho (1904)
- ▶ GINI's gamma (1914)
- ▶ KENDALL's tau (1938)
- ▶ BLOMQUIST's beta (1950)
- ▶ ...

Association within and between random vectors

Association between **scalar** outcomes:

- ▶ *within* the bivariate random vector $\mathbf{Z} = (X, Y)$
- ▶ *between* random variables X and Y

Possible generalizations to **vector**-valued outcomes:

- ▶ *within* a d -dimensional random vector $\mathbf{Z} = (Z_1, \dots, Z_d)$
- ▶ *between* random vectors $\mathbf{X} = (X_1, \dots, X_p)$ and $\mathbf{Y} = (Y_1, \dots, Y_q)$
→ *This talk*
- ▶ *between* random vectors $\mathbf{X}_1, \dots, \mathbf{X}_d$

Quantifying association between random vectors

- ▶ Classical measures of association:

- ▶ Canonical correlation [Hotelling 1936]
- ▶ RV coefficient [Escoufier 1973, Robert and Escoufier 1976]

However:

- ▶ positive values only: no distinction between positive/negative association
 - ▶ linear multivariate analysis
- ▶ Testing for independence between two random vectors
 - ▶ Matrices of Spearman and Kendall coefficients [El Maache and Lepage 2003]
 - ▶ Distances between densities [Székely, Rizzo and Bikorov 2007; Székely and Rizzo 2009]
 - ▶ Distances between copulas [Quesy 2010]

Again: no ‘direction’ for the association

Inference is to be rank-based and nonparametric

Invariance w.r.t. increasing transformations of the components

- ▶ No assumptions on existence of moments required
- ▶ Inference is based on vectors of **ranks**

Not even a parametric model for the ‘dependence’ (copula)

- ▶ Ideal for data exploration
- ▶ **Nonparametric** inference

Contributions

1. Measures of association between random vectors

- ▶ distinction between positive and negative associations
- ▶ margin free: invariance with respect to increasing transformations
- ▶ model free

2. Nonparametric, rank-based inference procedures

- ▶ point estimators taking the form of U -statistics
- ▶ consistent and asymptotically normal

Measuring Association Between Random Vectors

Measures of Association

Association Between Random Variables

Association Between Random Vectors

Inference

A Crash Course on U -Statistics

Estimating the Association Measures

Numerical Examples

Implementation

Simulation Study

Data Examples

Conclusion

Measuring Association Between Random Vectors

Measures of Association

Association Between Random Variables

Association Between Random Vectors

Inference

A Crash Course on U -Statistics

Estimating the Association Measures

Numerical Examples

Implementation

Simulation Study

Data Examples

Conclusion

Measuring Association Between Random Vectors

Measures of Association

Association Between Random Variables

Association Between Random Vectors

Inference

A Crash Course on U -Statistics

Estimating the Association Measures

Numerical Examples

Implementation

Simulation Study

Data Examples

Conclusion

Covariance and positive quadrant dependence

Random variables X and Y are *positive quadrant dependent* if

$$P(X \leq x, Y \leq y) \geq P(X \leq x) P(Y \leq y), \quad \forall x, y \in \mathbb{R}$$

[Lehmann 1966]

If X and Y have finite second moments, then

$$\text{cov}(X, Y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \{P(X \leq x, Y \leq y) - P(X \leq x) P(Y \leq y)\} dx dy$$

[Lehmann 1966, acknowledging Hoeffding 1940]

Lemma

X and Y are positive quadrant dependent iff

$$\text{cov}\{f(X), g(Y)\} \geq 0 \quad \forall \text{non-decreasing } f, g$$

[Esary, Proschan and Walkup 1967]

Measuring association in a margin-free way: Weighing the margins

Choose measures μ and ν on \mathbb{R}^2 and consider

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} P(X \leq x, Y \leq y) \, d\mu(x, y) \\ - \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} P(X \leq x) P(Y \leq y) \, d\nu(x, y)$$

μ	ν	association measure
Lebesgue	Lebesgue	covariance
$P_{X,Y}$	$P_X \otimes P_Y$	proportional to Kendall's tau
$P_X \otimes P_Y$	$P_X \otimes P_Y$	proportional to Spearman's rho

Special case: Kendall's tau

If (X_1, Y_1) and (X_2, Y_2) are iid (X, Y) , then

$$\begin{aligned} & \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} P(X \leq x, Y \leq y) dP(X \leq x, Y \leq y) \\ & - \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} P(X \leq x) P(Y \leq y) dP(X \leq x) dP(Y \leq y) \\ & = \text{cov}(\mathbf{1}\{X_1 \leq X_2\}, \mathbf{1}\{Y_1 \leq Y_2\}) \end{aligned}$$

Standardizing to obtain a result in $[-1, 1]$, we find **Kendall's tau**:

$$\begin{aligned} \tau(X, Y) &= \text{cor}(\mathbf{1}\{X_1 \leq X_2\}, \mathbf{1}\{Y_1 \leq Y_2\}) \\ &= 4 \text{cov}(\mathbf{1}\{X_1 \leq X_2\}, \mathbf{1}\{Y_1 \leq Y_2\}) \end{aligned}$$

provided the distributions of X and Y are continuous.

Properties of Kendall's tau

Special values:

$\tau = 1$ iff $Y = f(X)$ with f increasing

$\tau = 0$ if (but not only if) X and Y are independent

$\tau = -1$ iff $Y = f(X)$ with f decreasing

Compare: $\text{cor}(X, Y) = \pm 1$ iff $Y = aX + b$ with $\text{sign}(a) = \pm 1$

Margin free: If f and g are increasing, then

$$\tau(X, Y) = \tau(f(X), g(Y))$$

Special case: Spearman's rho

By the Hoeffding/Lehmann covariance formula,
provided $F(x) = P(X \leq x)$ and $G(y) = P(Y \leq y)$ are continuous:

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \{P(X \leq x, Y \leq y) - P(X \leq x)P(Y \leq y)\} dF(x) dG(y) \\ = \text{cov}\{F(X), G(Y)\}$$

Standardizing to obtain a result in $[-1, 1]$, we find **Spearman's rho**:

$$\rho_S(X, Y) = \text{cor}\{F(X), G(Y)\} \\ = 12 \text{cov}\{F(X), G(Y)\}$$

Enjoys similar properties as Kendall's tau

Measuring Association Between Random Vectors

Measures of Association

Association Between Random Variables

Association Between Random Vectors

Inference

A Crash Course on U -Statistics

Estimating the Association Measures

Numerical Examples

Implementation

Simulation Study

Data Examples

Conclusion

For random vectors, quantify association in a similar way

Random vectors \mathbf{X} in \mathbb{R}^p and \mathbf{Y} in \mathbb{R}^q . Distribution functions

$$H(\mathbf{x}, \mathbf{y}) = P(\mathbf{X} \leq \mathbf{x}, \mathbf{Y} \leq \mathbf{y}),$$

$$F(\mathbf{x}) = P(\mathbf{X} \leq \mathbf{x}),$$

$$G(\mathbf{y}) = P(\mathbf{Y} \leq \mathbf{y})$$

Measuring association: compare $H(\mathbf{x}, \mathbf{y})$ and $F(\mathbf{x}) G(\mathbf{y})$:

$$\int_{\mathbb{R}^p} \int_{\mathbb{R}^q} H(\mathbf{x}, \mathbf{y}) \, d\mu(\mathbf{x}, \mathbf{y}) \\ - \int_{\mathbb{R}^p} \int_{\mathbb{R}^q} F(\mathbf{x}) G(\mathbf{y}) \, d\nu(\mathbf{x}, \mathbf{y})$$

for appropriate choices of μ and ν .

A generalization of Kendall's tau: just as in the bivariate case

If $(\mathbf{X}_1, \mathbf{Y}_1)$ and $(\mathbf{X}_2, \mathbf{Y}_2)$ are iid (\mathbf{X}, \mathbf{Y}) , then

$$\begin{aligned} & \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} H(\mathbf{x}, \mathbf{y}) \, dH(\mathbf{x}, \mathbf{y}) \\ & - \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} F(\mathbf{x}) G(\mathbf{y}) \, dF(\mathbf{x}) \, dG(\mathbf{y}) \\ & = \text{cov}(\mathbf{1}\{\mathbf{X}_1 \leq \mathbf{X}_2\}, \mathbf{1}\{\mathbf{Y}_1 \leq \mathbf{Y}_2\}) \end{aligned}$$

Standardizing, we obtain a multivariate analogue of **Kendall's tau**:

$$\begin{aligned} \tau(\mathbf{X}, \mathbf{Y}) &= \text{cor}(\mathbf{1}\{\mathbf{X}_1 \leq \mathbf{X}_2\}, \mathbf{1}\{\mathbf{Y}_1 \leq \mathbf{Y}_2\}) \\ &= \frac{\text{cov}(\mathbf{1}\{\mathbf{X}_1 \leq \mathbf{X}_2\}, \mathbf{1}\{\mathbf{Y}_1 \leq \mathbf{Y}_2\})}{\sqrt{\text{var}(\mathbf{1}\{\mathbf{X}_1 \leq \mathbf{X}_2\}) \text{var}(\mathbf{1}\{\mathbf{Y}_1 \leq \mathbf{Y}_2\})}} \end{aligned}$$

If $p + q \geq 3$, the standardization depends on F and G .

Generalizing Spearman's rho:

the joint survival function comes into play

Let $(X_1, Y_1), (X_2, Y_2), (X_3, Y_3)$ be iid (X, Y) . Then

$$\begin{aligned}\iint \{H(\mathbf{x}, \mathbf{y}) - F(\mathbf{x}) G(\mathbf{y})\} dF(\mathbf{x}) dG(\mathbf{y}) &= \text{cov}(\mathbf{1}\{X_1 \leq X_2\}, \mathbf{1}\{Y_1 \leq Y_3\}) \\ &= \text{cov}(\bar{F}(X), \bar{G}(Y))\end{aligned}$$

$$\begin{aligned}\iint \{\bar{H}(\mathbf{x}, \mathbf{y}) - \bar{F}(\mathbf{x}) \bar{G}(\mathbf{y})\} dF(\mathbf{x}) dG(\mathbf{y}) &= \text{cov}(\mathbf{1}\{X_1 \geq X_2\}, \mathbf{1}\{Y_1 \geq Y_3\}) \\ &= \text{cov}(F(X), G(Y))\end{aligned}$$

in terms of the *joint survival functions*

$$\bar{H}(\mathbf{x}, \mathbf{y}) = P(X \geq \mathbf{x}, Y \geq \mathbf{y}),$$

$$\bar{F}(\mathbf{x}) = P(X \geq \mathbf{x}),$$

$$\bar{G}(\mathbf{y}) = P(Y \geq \mathbf{y})$$

If $p \geq 2$, it is no longer true that $\bar{F}(\mathbf{x})$ equals $1 - F(\mathbf{x})$ etc.

Generalizations of Spearman's rho: pick your choice

Several alternative generalizations of Spearman's rho:

$$\begin{aligned}\rho(\mathbf{X}, \mathbf{Y}) &= \text{cor}(F(\mathbf{X}), G(\mathbf{Y})) \\ &= \frac{\text{cov}(F(\mathbf{X}), G(\mathbf{Y}))}{\sqrt{\text{var } F(\mathbf{X}) \text{ var } G(\mathbf{Y})}}\end{aligned}$$

$$\begin{aligned}\bar{\rho}(\mathbf{X}, \mathbf{Y}) &= \text{cor}(\bar{F}(\mathbf{X}), \bar{G}(\mathbf{Y})) \\ &= \frac{\text{cov}(\bar{F}(\mathbf{X}), \bar{G}(\mathbf{Y}))}{\sqrt{\text{var } \bar{F}(\mathbf{X}) \text{ var } \bar{G}(\mathbf{Y})}}\end{aligned}$$

$$\rho^*(\mathbf{X}, \mathbf{Y}) = \frac{1}{2}(\rho(\mathbf{X}, \mathbf{Y}) + \bar{\rho}(\mathbf{X}, \mathbf{Y}))$$

Again, if $p + q \geq 3$, the standardizations depend on F and G .

Measuring Association Between Random Vectors

Measures of Association

Association Between Random Variables

Association Between Random Vectors

Inference

A Crash Course on U -Statistics

Estimating the Association Measures

Numerical Examples

Implementation

Simulation Study

Data Examples

Conclusion

Measuring Association Between Random Vectors

Measures of Association

Association Between Random Variables

Association Between Random Vectors

Inference

A Crash Course on U -Statistics

Estimating the Association Measures

Numerical Examples

Implementation

Simulation Study

Data Examples

Conclusion

Kendall's tau can be represented in terms of estimable parameters of degree two

Let $(X_1, Y_1), (X_2, Y_2)$ be iid (X, Y) . We have

$$\begin{aligned}\tau(\mathbf{X}, \mathbf{Y}) &= \text{cor}(\mathbf{1}\{X_1 \leq X_2\}, \mathbf{1}\{Y_1 \leq Y_2\}) \\ &= \frac{p_{X,Y} - p_X p_Y}{\sqrt{p_X (1 - p_X) p_Y (1 - p_Y)}}\end{aligned}$$

where

$$\begin{aligned}p_{X,Y} &= P(X_1 \leq X_2, Y_1 \leq Y_2), \\ p_X &= P(X_1 \leq X_2), \\ p_Y &= P(Y_1 \leq Y_2)\end{aligned}$$

The latter three probabilities can be estimated using ***U-statistics*** of degree $m = 2$.

U -statistics: generalizations of sample means

Let X_1, \dots, X_m be iid X , in some space \mathcal{X} .

Aim: estimate a 'parameter' θ of the form

$$\theta = \mathbb{E}[g(X_1, \dots, X_m)]$$

for some known function $g : \mathcal{X}^m \rightarrow \mathbb{R}$.

- ▶ Example: $p_X = P(X_1 \leq X_2) = \mathbb{E}[\mathbf{1}(X_1 \leq X_2)]$

The U -statistic estimator for θ based on a sample X_1, \dots, X_n is

$$\hat{\theta}_m = \frac{1}{n(n-1)\cdots(n-m+1)} \sum g(X_{i_1}, \dots, X_{i_m})$$

Sum over all permutations (i_1, \dots, i_m) of $\{1, \dots, n\}$ of length m

Hoeffding's decomposition theorem: Linear expansion of a U -statistic

If $E[g^2(X_1, \dots, X_m)] < \infty$, then

$$\sqrt{n}(\hat{\theta}_n - \theta) = \frac{m}{\sqrt{n}} \sum_{i=1}^n h_1(X_i) + o_p(1)$$

where

$$h_1(x_1) = E[g_{\text{sym}}(x_1, X_2, \dots, X_m)] - \theta$$
$$g_{\text{sym}}(x_1, \dots, x_m) = \frac{1}{m!} \sum g(x_{i_1}, \dots, x_{i_m})$$

[Hoeffding 1948]

Asymptotic normality

Hoeffding's decomposition yields **joint asymptotic normality** of a vector of U -statistics, e.g.

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N(0, m^2 \sigma_1^2) \quad (n \rightarrow \infty)$$

$$\sigma_1^2 = \text{var } h_1(X_1)$$

The asymptotic (co)variance(s) can be estimated consistently by

- ▶ U -statistics
- ▶ jackknife [Lee 1990]
- ▶ the sample (co)variance of

$$\hat{h}_{1,n}(X_i), \quad i = 1, \dots, n$$
$$\hat{h}_{1,n}(x) = U\text{-statistic}$$

Measuring Association Between Random Vectors

Measures of Association

Association Between Random Variables

Association Between Random Vectors

Inference

A Crash Course on U -Statistics

Estimating the Association Measures

Numerical Examples

Implementation

Simulation Study

Data Examples

Conclusion

Kendall's tau can be represented in terms of estimable parameters of degree two

Let $(X_1, Y_1), (X_2, Y_2)$ be iid (X, Y) . We have

$$\begin{aligned}\tau(\mathbf{X}, \mathbf{Y}) &= \text{cor}(\mathbf{1}\{X_1 \leq X_2\}, \mathbf{1}\{Y_1 \leq Y_2\}) \\ &= \frac{p_{X,Y} - p_X p_Y}{\sqrt{p_X (1 - p_X) p_Y (1 - p_Y)}}\end{aligned}$$

where

$$\begin{aligned}p_{X,Y} &= P(X_1 \leq X_2, Y_1 \leq Y_2), \\ p_X &= P(X_1 \leq X_2), \\ p_Y &= P(Y_1 \leq Y_2)\end{aligned}$$

The latter three probabilities can be estimated using U -statistics of degree $m = 2$.

U -statistic estimator for Kendall's tau

Represent $p_{X,Y}$ as an estimable parameter of degree $m = 2$:

$$\begin{aligned} p_{X,Y} &= P(\mathbf{X}_1 \leq \mathbf{X}_2, \mathbf{Y}_1 \leq \mathbf{Y}_2) \\ &= E[\underbrace{\mathbf{1}(\mathbf{X}_1 \leq \mathbf{X}_2) \mathbf{1}(\mathbf{Y}_1 \leq \mathbf{Y}_2)}_{=g((\mathbf{X}_1, \mathbf{Y}_1), (\mathbf{X}_2, \mathbf{Y}_2))}] \end{aligned}$$

The corresponding U -statistic estimator is

$$\hat{p}_{X,Y;n} = \frac{1}{n(n-1)} \sum_{i \neq j} \mathbf{1}(\mathbf{X}_i \leq \mathbf{X}_j) \mathbf{1}(\mathbf{Y}_i \leq \mathbf{Y}_j)$$

Similarly for $p_X = p_{X,X}$ and $p_Y = p_{Y,Y}$. Finally, put

$$\hat{\tau}_n(\mathbf{X}, \mathbf{Y}) = \frac{\hat{p}_{X,Y;n} - \hat{p}_{X,n} \hat{p}_{Y,n}}{\sqrt{\hat{p}_{X,n} (1 - \hat{p}_{X,n}) \hat{p}_{Y,n} (1 - \hat{p}_{Y,n})}}$$

The estimator is asymptotically normal

Hoeffding's decomposition yields joint asymptotic normality of

$$\sqrt{n} \begin{pmatrix} \hat{p}_{\mathbf{X},\mathbf{Y},n} - p_{\mathbf{X},\mathbf{Y}} \\ \hat{p}_{\mathbf{X},n} - p_{\mathbf{X}} \\ \hat{p}_{\mathbf{Y},n} - p_{\mathbf{Y}} \end{pmatrix}$$

with explicit 3×3 covariance matrix Σ .

From the delta method, we get asymptotic normality of

$$\sqrt{n}(\hat{\tau}_n(\mathbf{X}, \mathbf{Y}) - \tau(\mathbf{X}, \mathbf{Y}))$$

The expression for the asymptotic variance is longish, but explicit, and can be estimated consistently.

Spearman's rho can be estimated via a function of a vector of U -statistics of degree three

Recall

$$\begin{aligned}\rho(\mathbf{X}, \mathbf{Y}) &= \text{cor}(\bar{F}(\mathbf{X}), \bar{G}(\mathbf{Y})) \\ &= \frac{\text{cov}(\bar{F}(\mathbf{X}), \bar{G}(\mathbf{Y}))}{\sqrt{\text{var} \bar{F}(\mathbf{X}) \text{var} \bar{G}(\mathbf{Y})}}\end{aligned}$$

If $(\mathbf{X}_1, \mathbf{Y}_1), (\mathbf{X}_2, \mathbf{Y}_2), (\mathbf{X}_3, \mathbf{Y}_3)$ are iid (\mathbf{X}, \mathbf{Y}) , then

$$\begin{aligned}\text{cov}(\bar{F}(\mathbf{X}), \bar{G}(\mathbf{Y})) &= \text{cov}(\mathbf{1}(\mathbf{X}_1 \leq \mathbf{X}_2), \mathbf{1}(\mathbf{Y}_1 \leq \mathbf{Y}_3)) \\ &= P(\mathbf{X}_1 \leq \mathbf{X}_2, \mathbf{Y}_1 \leq \mathbf{Y}_3) - P(\mathbf{X}_1 \leq \mathbf{X}_2)P(\mathbf{Y}_1 \leq \mathbf{Y}_3)\end{aligned}$$

Similarly for $\text{var} \bar{F}(\mathbf{X})$ and $\text{var} \bar{G}(\mathbf{Y})$.

All probabilities occurring in these expressions can be estimated by U -statistics of degree $m \in \{2, 3\}$.

Measuring Association Between Random Vectors

Measures of Association

Association Between Random Variables

Association Between Random Vectors

Inference

A Crash Course on U -Statistics

Estimating the Association Measures

Numerical Examples

Implementation

Simulation Study

Data Examples

Conclusion

Measuring Association Between Random Vectors

Measures of Association

Association Between Random Variables

Association Between Random Vectors

Inference

A Crash Course on U -Statistics

Estimating the Association Measures

Numerical Examples

Implementation

Simulation Study

Data Examples

Conclusion

Write the formulas in such a way as to reduce the number of nested loops

- ▶ A naive implementation of the U -statistics formulas requires multiple nested loops over the sample:
 - ▶ all permutations (i_1, \dots, i_m) of $\{1, \dots, n\}$ of length mPatience needed if $n = 1\,000$ and $m = 3 \dots$
- ▶ In case of τ and ρ , a clever rewriting of the formulas reduces the calculations to at most double loops
 - ▶ Compare with the formulas for the sample variance:

$$\hat{\sigma}_n^2 = \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j=1}^n \frac{1}{2} (X_i - X_j)^2 \quad \text{double loop}$$

$$= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \quad \text{single loop}$$

Perform the comparisons in a smart order

Estimators involve statistics of the form

$$\sum_{i=1}^n \sum_{j:j \neq i} \mathbf{1}_{\underbrace{\{X_{i,1} \leq X_{j,1}, \dots, X_{i,p} \leq X_{j,p}\}}_{\mathbf{X}_i \leq \mathbf{X}_j}}$$

For a given $i = 1, \dots, n$:

1. Let $R_{i,k}$ be the rank of $X_{i,k}$ within $X_{1,k}, \dots, X_{n,k}$
2. Let $k(i)$ be the component of \mathbf{X}_i with the highest rank
3. Restrict the inner sum to those j such that $R_{j,k(i)} \geq R_{i,k(i)}$

Requires preliminary sorting of the data: $O(n \log n)$ algorithm

Implement the core computations in a low-level programming language

Even after all these tricks, nested loops are required: $O(n^2)$ steps

- ▶ Painful in a simulation study when n is large

Solution: do the computations in a low-level programming language

- ▶ for instance C or Fortran
- ▶ to be called from Matlab or R for convenience
- ▶ faster than computing directly in Matlab or R, even if vectorized

Measuring Association Between Random Vectors

Measures of Association

Association Between Random Variables

Association Between Random Vectors

Inference

A Crash Course on U -Statistics

Estimating the Association Measures

Numerical Examples

Implementation

Simulation Study

Data Examples

Conclusion

Aims of the simulation study

To answer the following questions:

- ▶ Finite-sample bias and variance of the estimators of the association measures τ and ρ^* ?
- ▶ Accuracy of the estimators of the standard errors?
- ▶ Influence of a number of factors:
 - ▶ Model
 - ▶ Dimension
 - ▶ Sample size
- ▶ Computational feasibility

Set-up

- Models:
- ▶ multivariate Gaussian
 - ▶ D-vine copula with bivariate Clayton margins

[Clayton 1978; Joe 1997; Bedford and Cooke 2001, 2002]

Sample sizes: $n = 50, 100, 1\,000$

Dimensions: $p = q = 3$ and $p = q = 4$

Number of repetitions: 10 000

‘True’ values: based on 30 repetitions of samples of size 50 000

Results: Everything works as it should

Imagine here some
huge multi-way tables
with tons of digits

- ▶ $\hat{\tau}_n$ has lower standard error and bias than $\hat{\rho}_n^*$
- ▶ Finite-sample bias, decreasing with sample size
 - ▶ Bias is highest for $\hat{\rho}_n^*$ at weak dependence and low sample size
- ▶ Standard deviations are estimated accurately for sample sizes above 100

Results: Everything works as it should

Imagine here some
huge multi-way tables
with tons of digits

- ▶ $\hat{\tau}_n$ has lower standard error and bias than $\hat{\rho}_n^*$
- ▶ Finite-sample bias, decreasing with sample size
 - ▶ Bias is highest for $\hat{\rho}_n^*$ at weak dependence and low sample size
- ▶ Standard deviations are estimated accurately for sample sizes above 100

Measuring Association Between Random Vectors

Measures of Association

Association Between Random Variables

Association Between Random Vectors

Inference

A Crash Course on U -Statistics

Estimating the Association Measures

Numerical Examples

Implementation

Simulation Study

Data Examples

Conclusion

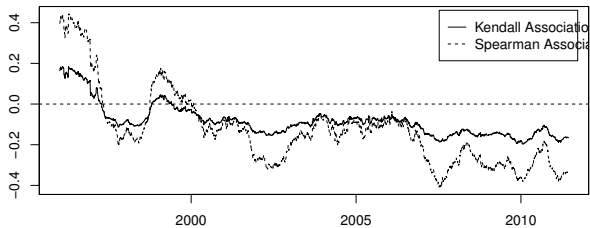
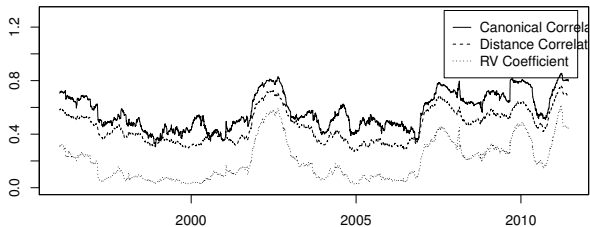
Stocks versus Bonds

- ▶ In terms of crises, association between stocks and bonds can become negative
 - ▶ Investors searching for safe havens like bonds of strong countries

[Gulko 2002, Ilmanen 2003]

- ▶ Data: daily returns
 - ▶ stock market indices of 5 major countries
 - ▶ All Ordinaries, CAC 40, DAX, Nikkei 225, S&P 500
 - ▶ government bond indices for the same countries
 - ▶ indices from The Bank of America Merrill Lynch
 - ▶ From January 3, 1996, to November 15, 2012
 - ▶ Forward looking moving window of 150 days

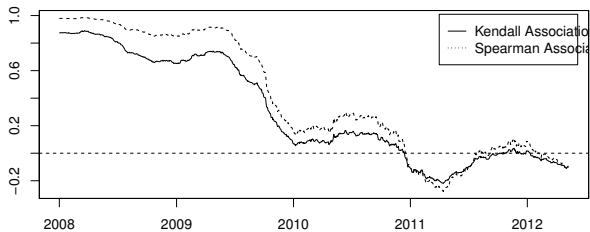
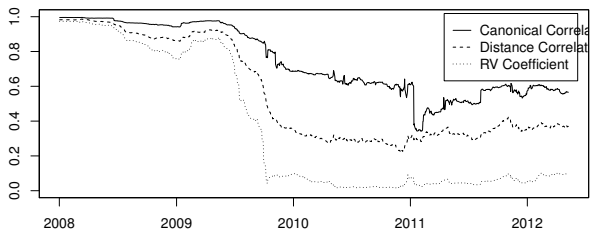
From positive to negative association



North versus South

- ▶ European sovereign debt crisis
- ▶ Association between north and south European bond markets changes as credit-worthiness of countries evolve
- ▶ Data: daily returns of Merrill Lynch government bond indices
 - ▶ North: France, Germany, the Netherlands
 - ▶ South: Italy, Portugal, Spain
 - ▶ From January 1, 2007 to November 15, 2012
 - ▶ Forward looking moving window of 150 days

With the crisis, association becomes negative



Measuring Association Between Random Vectors

Measures of Association

Association Between Random Variables

Association Between Random Vectors

Inference

A Crash Course on U -Statistics

Estimating the Association Measures

Numerical Examples

Implementation

Simulation Study

Data Examples

Conclusion

Contributions

1. Measures of association between random vectors
 - ▶ positive/negative associations
 - ▶ margin free (copula-based)
 - ▶ model free
2. Nonparametric, rank-based inference procedures
 - ▶ U -statistics > consistent, asymptotically normal
3. Numerical examples
 - ▶ Finite-sample performance is satisfactory
 - ▶ Case studies demonstrate potential applicability