

# Estimation of heavy tail dependence based on copulas for the precipitation analysis

N. D. Rassakhan\*, E. Yu. Shchetinin†

\* *Department of Applied Mathematics  
Moscow State Technology University "STANKIN"  
3a Vadkovsky Ln., Moscow, 127055, Russia*

† *FGU "All-Russian research institute on problems of civil defence and emergencies"  
of Emergency Control Ministry of Russia  
7 Davydkovskaya St., Moscow, 121352, Russia*

Email: rassahan@gmail.com, riviera-molto@mail.ru

Consideration of tail dependence is a very important part of risk analysis in many applied sciences that is measured in order to estimate the risk of simultaneous extreme events. Usually the tail dependence coefficient is the measurement in question. Pearson correlation coefficient unfortunately is not a suitable measure for estimating dependencies between two quantities in the context of simultaneous occurrence of extreme events when these events are of interest for the researcher because it takes extreme events into account with the same weight as it takes "normal" events although dependence of extreme values may slightly differ [2].

Present work emphasizes the importance of taking into consideration tail dependencies in bivariate statistical analysis using copulas. Due to increasing frequency of environmental cataclysms the issue of analyzing risks (e.g. economic losses) and their consequences comes to the fore. Moreover, researchers should take into account consequences of their joint occurrence. Three non-parametric estimators of tail dependence coefficients were compared in order to estimate correlation between daily cumulative rainfall totals recorded in central European part of Russia. Major part of existing estimators depends on threshold  $k$  and thus there is a trade-off between variance and bias during the calculation of the best value for  $k$ . For balancing an algorithm is presented that is based on using moving average and then searching the "stable" part of tail dependence coefficient [6]. Estimate of tail dependence coefficient is assumed to be equal to mean value on the "stable" part.

**Key words and phrases:** extreme value theory, spatial modelling, extreme precipitation, spatial structures of statistical dependence, tail dependence coefficient.

## 1. Introduction

One of the most important parts of multivariate extreme value analysis is exploration of extremal dependencies; basically tail dependence coefficient is used for this purpose. For bivariate vector  $(X_1, X_2)$  upper tail dependence coefficient has the following form [5]:

$$\lambda_U = \lim_{t \rightarrow 1^-} P(F_1(X_1) > t | F_2(X_2) > t), \quad (1)$$

where  $F_1, F_2$  are distribution functions of random variables  $X_1, X_2$  respectively,  $0 < t \leq 1$  is the threshold. We can define lower tail dependence coefficient in the same way:

$$\lambda_L = \lim_{t \rightarrow 0^+} P(F_1(X_1) < t | F_2(X_2) < t), \quad (2)$$

Using the copula function [8] equations (1),(2) can be written in alternative form:

$$\lambda_U = \lim_{t \rightarrow 1^-} \frac{1 - 2t + C(t, t)}{1 - t}; \quad \lambda_L = \lim_{t \rightarrow 0^+} \frac{C(t, t)}{t}.$$

## 2. Main section

Tail dependence coefficient estimation methods are essential analysis tools for extremal precipitation structures that are studied in this paper. Onward we will describe some of them. Foremost such estimators are non-parametric estimators based on empirical copula  $C^{(n)}(u, v)$  concept [4]:

$$C^{(n)}(u, v) = F^{(n)} \left( F_{(n)1}^{-1}(u), F_{(n)2}^{-1}(v) \right),$$

where  $F_{(n)}(\cdot)$  is empirical distribution function.

Let  $(X_1^{(1)}, X_2^{(1)}), \dots, (X_1^{(n)}, X_2^{(n)})$  be independent identically distributed copies of bivariate random vector  $(X_1, X_2)$ . Using the fact that their joint distribution function is defined as

$$F(x_1, x_2) = P(F_1(X_1) \leq F_1(x_1), F_2(X_2) \leq F_2(x_2)) = C(F_1(x_1), F_2(x_2))$$

and equation

$$\lambda = 2 - \lim_{t \rightarrow 1^-} \frac{1 - C(t, t)}{1 - t} = 2 - \lim_{t \rightarrow 0^+} \frac{1 - C(1 - t, 1 - t)}{t},$$

an estimator for upper tail dependence (1) coefficient can be derived [7]:

$$\hat{\lambda}^{SEC} \equiv \hat{\lambda}^{SEC}(k) = 2 - \frac{1 - \hat{C}\left(1 - \frac{k}{n}, 1 - \frac{k}{n}\right)}{\frac{k}{n}}, \quad 1 \leq k < n. \quad (3)$$

Then taking into consideration  $\log(1 - t) \sim -t$ ,  $t \approx 0$  next estimator can be obtained:

$$\hat{\lambda}^{LOG} \equiv \hat{\lambda}^{LOG}(k) = 2 - \frac{\log \hat{C}\left(1 - \frac{k}{n}, 1 - \frac{k}{n}\right)}{\log\left(1 - \frac{k}{n}\right)}, \quad 1 \leq k < n. \quad (4)$$

Here  $\hat{C}$  denotes empirical copula defined as

$$\hat{C}\left(1 - \frac{k}{n}, 1 - \frac{k}{n}\right) = \frac{1}{n} \sum_{i=1}^n 1_{\{\hat{F}_1(X_1^{(i)}) \leq 1 - \frac{k}{n}, \hat{F}_2(X_2^{(i)}) \leq 1 - \frac{k}{n}\}}, \quad 1 \leq k < n,$$

where 1 is indicator function and  $\hat{F}_j, j = 1, 2$  are marginal empirical distribution functions of  $X_1$  and  $X_2$  respectively. To increase estimation's accuracy it is often considered that

$$\hat{F}_j(u) = \frac{1}{n+1} \sum_{i=1}^n 1_{\{X_j^{(i)} \leq u\}}, \quad j = 1, 2.$$

Note that both estimators depend on choice of threshold  $k$  and thereafter  $k^{th}$  order statistic. It is very important to choose right value for  $k$  that is not an easy task due to trade-off between variance and bias.

Another estimator for upper tail dependence coefficient is suggested in works [9, 10]:

$$\hat{\lambda}^{CFG} = 2 - 2 \exp \left[ \frac{1}{n} \sum_{i=1}^n \log \left\{ \frac{\sqrt{\log \frac{1}{\hat{F}_1(X_1^{(i)})} \log \frac{1}{\hat{F}_2(X_2^{(i)})}}}{\log \frac{1}{\max(\hat{F}_1(X_1^{(i)}), \hat{F}_2(X_2^{(i)}))^2}} \right\} \right]. \quad (5)$$

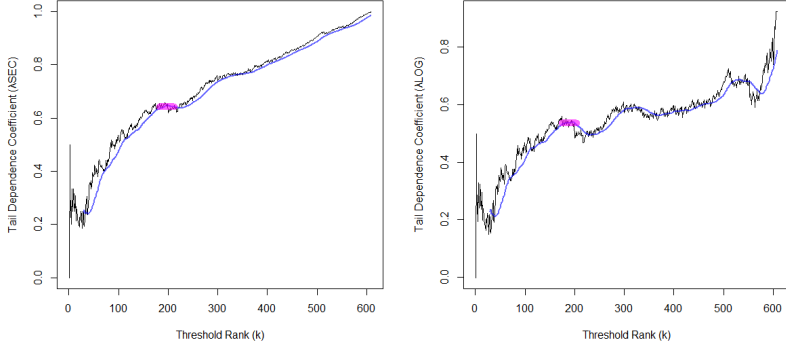


Figure 1. Implementation of the algorithm for finding "stable" part of TDC for  $\hat{\lambda}^{SEC}$  (left) and  $\hat{\lambda}^{LOG}$  (right). Moscow and Ryazan were used as a pair of cities from the area under study

Main advantage of this equation is that  $\hat{\lambda}$  doesn't depend on  $k$ . However, copula  $C(X_1, X_2)$  must be well approximated with extreme-value copulas for correctness of the estimator.

It follows from the equations (3), (4) that estimators depend on choice of threshold  $k$  choice of which is determined by balancing variance and bias for estimator according to stability theorem for  $\lambda_U$  [8]. Increasing the value of  $k$  leads to reduction of bias and increase in variance; it goes the same the other way around. For big enough data sample size  $n$  balance between bias and variance is described by the "stable" part of  $\lambda_U$  plot. An algorithm for finding this "stable" part is presented in paper [6]:

1. Empirical estimation is smoothed with moving average filter which window size is equal to  $b = \text{int}(0.05n)$ . Sequence  $\hat{\lambda}_1, \dots, \hat{\lambda}_{n-2b}$  is obtained as a result.
2. Now we can find vector  $(\hat{\lambda}_k, \dots, \hat{\lambda}_{k+m-1})$ ,  $k = 1, \dots, n - 2b + m - 1$ ,  $m = \text{int}(\sqrt{n - 2b})$  from the sequence  $\hat{\lambda}_1, \dots, \hat{\lambda}_{n-2b}$  by a sequential search.
3. If the current vector satisfies

$$\sum_{i=k+1}^{k+m-1} |\bar{\lambda}_i - \bar{\lambda}_k| \leq 2\sigma,$$

where  $\sigma$  is standard deviation of  $\hat{\lambda}_1, \dots, \hat{\lambda}_{n-2b}$  then final evaluation takes the form of

$$\lambda_U = \frac{1}{m} \sum_{i=1}^m \bar{\lambda}_{k+i-1}.$$

If the condition is not satisfied after sequential searching, then

$$\lambda_U = 0.$$

In this study precipitation data of the All-Russian Research Institute of Hydrometeorological Information - the World Data Center of the Russian Federation is used, which is daily precipitation in 11 cities of the European part of Russia. The data is freely available on the website <http://aisori.meteo.ru/ClimaterR> and is represented

by a set of tables (a separate table for each city); each table contains daily rainfall value for the period 1966-2016 years.

Implementation of this algorithm is presented in Fig.1. Both plots are using monthly maximum of precipitation in Moscow and Kostroma to evaluate upper tail dependence coefficient using estimators  $\hat{\lambda}^{SEC}$  (left) and  $\hat{\lambda}^{LOG}$  (right). Black line corresponds to  $\hat{\lambda}(k)$ ; blue smooth line is  $\hat{\lambda}(k)$  after applying moving average filter to it. Pink transparent plateau is the resulting value for  $\hat{\lambda}_U$  where placement of plateau corresponds to indexes of vector  $\hat{\lambda}_k, \dots, \hat{\lambda}_{k+m-1}$  from the algorithm above.

All three estimators (3), (4), (5) for upper tail dependence coefficient were calculated for 55 pairs of 11 cities under study. Furthermore, Pearson's correlation coefficient (PCC) was also calculated with the view to compare it with estimators. Results for some pairs are represented by Table 1. As we can see, PCC is quite different from all other estimators sometimes (Sp. Petersburg - N. Novgorod pair as an example) but it actually is close enough to at least one of estimators for the most of pairs.

Finally we want to find (or see) some correlation between estimators for upper tail dependence and the distance between cities under study so we plot a scatterplot for four values compared in Table 1. Resulting scctterplots are presented in Fig. 2.

It is obvious that there must be an inverse relation between the distance  $h$  and dependence estimators. Therefore  $\hat{\lambda}^{SEC}$  is bad estimator for the problem under study. Three other estimators ( $\hat{\lambda}^{LOG}$ ,  $\hat{\lambda}^{CFG}$  and PCC) show roughly the same with some correction, that's why they are considered to be more trustworthy. So it is proposed to take the average of  $\hat{\lambda}^{LOG}$  and  $\hat{\lambda}^{CFG}$  or just  $\hat{\lambda}^{CFG}$  as the resulting estimator for  $\lambda_U$ .

### 3. Conclusions

This paper highlights the importance of taking into account the tail dependence coefficient in the context of multivariate frequency analysis using copulas. The three following nonparametric estimators ( $\hat{\lambda}^{SEC}$ ,  $\hat{\lambda}^{LOG}$ ,  $\hat{\lambda}^{CFG}$ ) have been compared. The aim of this comparison was to choose the best estimator in the context of our application. No estimator works in every case yet some of them show poor performance thus they need to be excluded. It is therefore important to pursue research in this field to get the right estimation for  $\lambda_U$  based on values of  $\hat{\lambda}^{SEC}$ ,  $\hat{\lambda}^{LOG}$  and  $\hat{\lambda}^{CFG}$ .

Table 1  
Values for estimators of  $\lambda_U$  and Pearson's correlation coefficient calculated for some pairs of cities in the European part of Russia

Pair of cities	$\hat{\lambda}^{SEC}$	$\hat{\lambda}^{LOG}$	$\hat{\lambda}^{CFG}$	PCC
Moscow - Kolomna	0.6417793	0.5358919	0.5210492	0.6021248
Kolomna - Ryazan	0.5561837	0.4623016	0.5022938	0.5509277
Pskov - Smolensk	0.54967	0.437089	0.3722651	0.4314
Kostroma - N.Novgorod	0.6346736	0.415878	0.3856968	0.4766541
Bryansk - Mozhaik	0.6846937	0.1764084	0.433269	0.6021248
St. Petersburg - N.Novgorod	0.1470922	0.0770384	0.2608262	0.3353886
Smolensk - Moscow	0.488586	0.3740484	0.3968146	0.3952423
St. Petersburg - Pskov	0.5338348	0.4111191	0.4124836	0.4867774
N.Novgorod - Tambov	0.2271567	0.1740808	0.3475754	0.3898402
Pskov - Kostroma	0.5107715	0.4133031	0.3863028	0.5022300

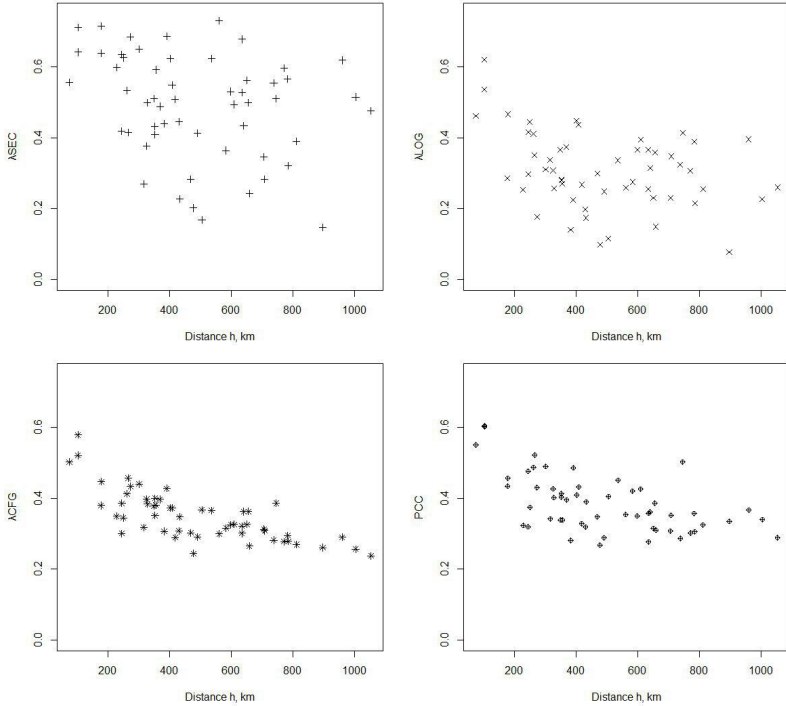


Figure 2. Comparison of 4 estimators and their dependence on the distance  $h$  between observation points:  $\hat{\lambda}^{SEC}$  - upper left,  $\hat{\lambda}^{LOG}$  - upper right,  $\hat{\lambda}^{LOG}$  - lower left, PCC - lower right.

Most of the nonparametric estimators have to deal with the choice of the number  $k$  of order statistics to be considered in the production of an estimate. This is not an easy task since it requires a trade-off between variance and bias (small values of  $k$  cause large variance and large values of  $k$  increase the bias).

Frahm et al. [6] introduced a simple plateau finding algorithm after smoothing the latter plot by some box kernel in order to find the optimal threshold  $k$ . Since this very simple algorithm revealed some potential, we intend to develop this idea further. But it is important to not overcomplicate things because PCC is not that bad compared to estimators presented in this paper.

## References

1. *Ferreira M.* Nonparametric estimation of the Tail-dependence coefficient. *Revstat*, Vol. 11, no. 1. — P. 1–16, 2013.
2. *Poulin A., Huard D., Favre A.-C., Pugin S.* Importance of Tail Dependence in Bivariate Frequency Analysis. *Journal of hydrologic engineering* — 2007.
3. *Ferreira M., Silva S.* An Analysis of a Heuristic Procedure to Evaluate Tail (in)dependence. *Journal of Probability and Statistics*, vol. 2014 — 2014

4. *Sibuya M.* Bivariate extreme statistics, I. *Annals of the Institute of Statistical Mathematics*, Vol. 11, no. 2. — P. 195–210, 1959.
5. *Draisma G., Drees H., Ferreira A., De Haan L.* Bivariate tail estimation: dependence in asymptotic independence. *Bernoulli*, Vol. 10, no. 2. — P. 251–280, 2004.
6. *Frahm G., Junker M., Schmidt R.* Estimating the tail-dependence coefficient: Properties and pitfalls. *Insurance: Mathematics and Economics*, Vol. 37, no. 1. — P. 80–100, 2005.
7. *Joe H., Smith R.L., Weissman I.* Bivariate Threshold Methods for Extremes. *Journal of the Royal Statistical Society. Series B*, Vol. 54, no. 1. — P. 171–183, 1992.
8. *Coles S., Heffernan J., Tawn J.* Dependence Measures for Extreme Value Analyses. *Extremes*, Vol. 2, no. 4. — P. 339–365, 1999.
9. *Caperua P., Fougères A.-L., Genest C.* A nonparametric estimation procedure for bivariate extreme value copulas. *Biometrika*, Vol. 84, no. 5. — P. 567–577, 1997.
10. *Schmidt R., Stadtmüller U.* Non-parametric Estimation of Tail Dependence. *Scandinavian Journal of Statistics*, Vol. 33, no. 2. — P. 307–335, 2006.

УДК 519.246.5

## Оценка зависимостей тяжелых хвостов при помощи аппарата копул для анализа осадков

Н. Д. Рассахан\*, Е. Ю. Щетинин†

\* Кафедра прикладной математики  
Московский государственный технологический университет "Станкин"  
Вадковский пер., д.3а, Москва, Россия, 127055

† Всероссийский научно-исследовательский институт  
по проблемам гражданской обороны и чрезвычайных ситуаций МЧС России  
ул. Давыдовская, д. 7, Москва, Россия, 121352

Email: rassahan@gmail.com, riviera-molto@mail.ru

Измерение хвостовой зависимости является важной задачей во многих прикладных науках для оценивания риска совместного наступления экстремальных событий. Обычно мерой зависимости является коэффициент хвостовой зависимости. Корреляция Пирсона не является подходящей мерой для оценивания зависимости двух величин в контексте совместного наступления экстремальных событий, когда они представляют интерес для исследователя, так как она учитывает экстремальные события так же (с тем же весом), что и "рядовые" события, хотя зависимость между экстремальными событиями может сильно отличаться от общей картины.

Данная работа подчеркивает важность учитывания хвостовой зависимости в контексте двумерного анализа при помощи копул. В связи с учащающимися природными катаклизмами резко встает вопрос об оценивании разнотипных рисков (в т.ч. экономических) и последствий их совместного наступления с учетом пространственных связей между наблюдениями. Сравниваются 3 непараметрические оценки коэффициента хвостой зависимости для оценки зависимости между ежедневными наблюдениями осадков в городах Европейской части России. Большинство существующих оценок зависит от порога  $k$  и, следовательно, при выборе используемого значения  $k$  происходит трейд-офф между смещением и вариацией. Для установления баланса в работе представлен алгоритм, основанный на использовании скользящего среднего и поиска "стабильного участка" коэффициента хвостовой зависимости. Именно среднее значение на "стабильном" участке и принимается за значения оценки коэффициента хвостовой зависимости.

**Ключевые слова:** теория экстремальных величин, пространственное моделирование, экстремальные осадки, пространственные структуры статистической зависимости, коэффициент хвостовой зависимости.