

Krzysztof Basiaga, Tomasz Szkutnik

University of Economics in Katowice, Poland

tomasz.szkutnik@ue.katowice.pl

THE APPLICATION OF GENERALIZED PARETO DISTRIBUTION AND COPULA FUNCTIONS IN THE ISSUE OF OPERATIONAL RISK

Abstract: The article concerns the issue of modelling of operational risk in a bank. The area of analysis is related to two separate analytical areas composed of certain combinations of the Basel Matrix risk categories. The focus of interest is in the modelling of loss severity distributions in LDA models and in consideration of the power and character of dependences among the studied analytical areas. To model a single loss severity distribution, the authors used the approach based on extreme values theory EVT. GPD distribution was used to model the right tail. The *t*-Student copula function was used in the cases of consideration of power and character of dependences. The determined values describe the effects of the applied approach in relative scale.

Keywords: operational risk, copula functions, GPD distribution, VaR.

1. Introduction

The application of operational risk modelling started to be practiced in banking operations years ago. With the introduction of precautionary regulations based on the guidelines of the Basel Committee, the need for such a measure became particularly significant. The main goal of banks for the coming years is to lower capital load related to the lack of methods of the advanced measurement for operational risk.

In cases when a bank does not apply any advanced measurement method (further AMA, Advanced Measurement Approach), it is forced to use basic solutions that are proposed by supervisory institutions.

With regard to its specific character, operational risk is heterogeneous with respect to its internal structure and may embrace many aspects of bank operations, as well as many risk factors that influence a given bank. In order to unify the area of analysis related to the application of the AMA model, it is proposed to divide all risk factors with respect to business lines, thus creating the Basel Matrix, i.e. a system of 56 units, named further as risk categories. The value of risk in each unit is estimated individually on the basis of the relevant measure.

One of the most frequent AMA models in the estimation of operational risk is the LDA model (Loss Distribution Approach). This is an actuarial model that models separately distributions of loss frequency and loss severity. Aggregated distribution of loss and related risk measure are determined on the basis of the LDA model. In the considered case, risk measure is defined as α quantile for distribution F_L of a given random variable L , which can be written as $VaR_\alpha(L) = F_L^{-1}(\alpha) = \inf\{l \in \mathbb{R}: F_L(l) \geq \alpha\}$ [Embrecht et al. 2003, p. 147]. The method of determining the aggregated distribution of loss and its form are defined by the estimated distributions of loss frequency and severity. In particular, maladjustment of severity distribution in a distribution tail may significantly determine the aggregated value of VaR, thus causing an underestimation of the whole model and a definition of risk exposure degree at a far too low level.

The article uses the actual data related to losses with respect to operational risk for one of the Polish financial institutions. The data period includes 5 years in total. With regard to their structure and affinity to various risk categories, these data are divided into two sub-sets (each sub-set is further defined as an analytical area or section). Taking into consideration sensitivity of data and security requirements of the financial institution from which the data are derived, the authors do not present the parameters of the estimated distributions and other values that could identify the real extent of loss, its amount or affinity to particular risk categories of the Basel Matrix.

Section 2 presents the LDA model. There are defined distributions of loss severity and frequency, as well as the method of determining aggregated distribution of loss.

Section 3 presents the procedure of defining total risk for the bank on the basis of information related to the aggregated distributions of loss from particular analytical areas. It focuses on the application of the copula function in order to consider dependences among the aggregated distributions of loss and their comparison with the conservative approach proposed by the Basel Committee, where the way of defining the total value of VaR consists in summing up values for particular analysis areas.

Section 4 presents the application of the theory introduced in sections 2 and 3 to the real data related to operational loss and the summary of the achieved results.

2. LDA method and its outcomes

In the LDA method aggregated loss $L(i, j)$ is determined as a random sum of individual losses:

$$L(i, j) = \sum_{k=0}^{N(i, j)} X_k(i, j), \quad (1)$$

where $N(i, j)$ denotes the number of losses in the given time interval Δt , $X_k(i, j)$ means value of k -th loss in the given time interval Δt . Value i denotes business line

and j denotes risk factor. In the discussed case, indices i and j are limited to two analytical sections composed of separate risk categories defined by confidential elements (i, j) . In the LDA model the following values: loss frequency $N(i, j)$ and loss severity $X_k(i, j)$ must fulfil certain assumptions:

- Conditionally, with $N(i, j) = n$ random variables X_1, \dots, X_n are independent random variables with the same probability distributions.
- Conditionally, with $N(i, j) = n$ distribution of random variables is independent of n .
- Distribution of random variable $N(i, j)$ does not belong to distributions of variables X_1, \dots, X_n .

Cumulative distribution function for aggregated loss $L(i, j)$ in model LDA (1) may be expressed by means of:

$$F_{L(i,j)}(l(i,j)) = \begin{cases} \sum_{n=1}^{\infty} P(N(i,j) = n) F_X^{n*}(l(i,j)) & \text{for } l(i,j) > 0 \\ P(N(i,j) = 0) & \text{for } l(i,j) = 0 \end{cases}, \quad (2)$$

where F_X is a cumulative distribution function of random variable X , and F_X^{n*} denotes n -th convolution of distribution F_X with itself, $F_X^{n*} = P(\sum_k^n X_k \leq l(i,j))$.

With regard to the convolution of distribution F_X , the form of cumulative distribution function $F_{L(i,j)}$ is most frequently determined by means of the Monte Carlo simulation. Analytical solutions or other approximate solutions, e.g. by means of FFT (Fast Fourier Transform) or Panjer recursive algorithm are available only for the specified types of distributions of loss frequency and severity. In the article, the aggregated distribution of loss is determined by means of the Monte Carlo procedure, which is not described herein.

2.1. Distribution of loss frequency

Distribution of loss frequency is a discrete distribution, the aim of which is to model the number of losses in the given time interval, i.e. $N(i, j) \in \{0,1,2, \dots\}$. The most frequent discrete distributions that are used in the modelling of loss frequency are Poisson distributions and negative binomial distributions (a negative binomial distribution is a generalization of Poisson distribution).

The probability distribution function for the Poisson distribution is as follows:

$$P(N(i,j) = n) = \frac{e^{-\lambda(i,j)} \lambda(i,j)^n}{n!}, \quad (3)$$

where $P(N(i,j) = n)$ denotes the probability that the number of losses in the given time interval will equal n . For Poisson distribution mean value $E(N)$ equals variance $Var(N)$. This is a feature of the distribution and it is assumed, in practice, that the

relation $E(N) \approx Var(N)$ should be fulfilled. Poisson distribution has one parameter λ , thus it has certain limitations in the situations where the above mentioned feature is not fulfilled, i.e. when $E(N) > Var(N)$ or $E(N) < Var(N)$.

Negative binomial distribution, which has two parameters, is a more elastic distribution in the modelling of event frequency. The function of distribution is as follows:

$$P(N(i,j) = n) = \binom{\alpha(i,j) + n - 1}{n} \left(\frac{1}{1+\beta(i,j)}\right)^{\alpha(i,j)} \left(\frac{\beta(i,j)}{1+\beta(i,j)}\right)^n. \quad (4)$$

The relation between mean value $E(N)$ and variance $Var(N)$, in contrast to Poisson distribution, is as follows: $E(N) < Var(N)$. In the issues of operational risk, where variance is higher than the mean value for distributions of loss frequency, the negative binomial distribution may be a better candidate to describe the modelled process than Poisson distribution.

Mutual relations between characteristics $E(N)$ and $Var(N)$ of the considered distributions are certain qualitative comments and may be treated only as premise. These are not formal methods of selection of distribution because it is not known to what extent variance is to be higher than mean value to make negative binomial distribution more favourable [Klugman et al. 2008, p. 110].

A more objective criterion is the application of e.g. Chi-square goodness-of-fit test as a tool of selection of a relevant distribution [Klugman et al. 2008, p. 452].

2.2. Distribution of loss severity

Distribution of loss severity for operational risk is often modelled by means of log-normal, exponential, gamma or Weibull distribution. In cases when data concern operational risk, it may turn out that these distributions do not describe extreme losses in a relevant way and, thus, the aggregated distribution of losses may be underestimated in its right side. Extreme losses seldom occur, however, their influence on the total value of losses in the analysed time interval may be very high. It turns out in practice that the losses that constitute the greatest load for a bank are slightly modelled [Di Clemente, Romano 2004]. Therefore, it is necessary to model the right tail of distribution of loss severity separately. GPD distribution is used in order to model the right tail of the distribution.

The loss severity model for operational risk that takes into account a separate modelling of its right tail by means of GPD distribution may be presented as follows [Di Clemente, Romano 2004]:

$$F_{i,j}(x) = \begin{cases} \Phi\left(\frac{\ln x - \mu(i,j)}{\sigma(i,j)}\right) & \text{for } 0 < x < u(i,j) \\ 1 - \frac{N_{u(i,j)}}{N_{i,j}} \left(1 + \xi(i,j) \frac{x - u(i,j)}{\beta(i,j)}\right)^{-1/\xi(i,j)} & \text{for } u(i,j) \leq x \end{cases}, \quad (5)$$

where $N_{u(i,j)}$ is the number of losses that exceed the value of $u(i,j)$. Values $u(i,j), \beta(i,j), \xi(i,j)$ are parameters of GPD distribution, respectively for the position, scale and shape.

Threshold value $u(i,j)$ is the starting point for modelling of data by means of GPD distribution. As pointed out in the subject literature [Di Clemente, Romano 2004, p. 196], in order to maintain the monotone of cumulative distribution function (5), the selection of threshold value $u(i,j)$ is to define the highest value of x , such that:

$$\Phi\left(\frac{\ln x - \mu(i,j)}{\sigma(i,j)}\right) < 1 - \frac{N_x}{N}, \tag{6}$$

where N_x is the number of historical loss that exceed value x .

Distribution of loss above threshold $u(i,j)$ or, in other words, distribution of excess above threshold $u(i,j)$ is defined as follows:

$$F_u(x) = P(X - u \leq x | X > u), x \geq 0. \tag{7}$$

For high values of $u(i,j)$, the right tail of distribution may be approximated by GPD distribution. Of course the selection of threshold is a key value and it influences the quality of achieved estimators [Embrecht et al. 1997]. The application of GPD distribution in the modelling of tail of another distribution is motivated by extreme value theory, as well as Pickands theorem and Balkem-de-Haan theorem. The precise introduction into the issue of extreme value theory may be found in [Klugman et al. 2008; Embrecht et al. 1997].

3. Total loss distribution for the whole institution

The division of bank operations into separate analytical sections defines the aggregated distributions of losses for each of them and, at the same time, allows for the determination of **VaR** value for each process individually.

In order to specify the level of operational risk for the whole bank, the already defined individual values of **VaR** should be aggregated.

The easiest way of aggregation consists in the summing up of value **VaR** for the already defined areas of operation. Such a procedure assumes a perfect dependence among particular analytical sections, which, in practice, means the assumption that loss will occur simultaneously in all areas of the operation. This is consistent with the initial requirements of the Basel Committee and is regarded as a conservative approach towards the modelling of operational risk.

In cases when particular analytical sections do not show a perfect dependence among themselves, there is a possibility to use the effect of diversification of risk in particular areas of operation. To that end, copula functions are applied with this view. It is known that copula functions allow for the expression of character of dependences separately from marginal distributions and they may model non-linear dependences.

3.1. Copula functions

Copula function \mathcal{C} is a n -dimensional cumulative distribution function with marginal distributions from uniform distribution in section [0,1].

One of the fundamental theorems in this case is the Sklar theorem.

Theorem 1 (Sklar). Let X_1, \dots, X_n be random variables with a total n -dimensional cumulative distribution function F and marginal cumulative distribution functions F_1, \dots, F_n . Then, there is such a copula function \mathcal{C} that [Copula methods... 2004]:

$$F(x_1, \dots, x_n) = \mathcal{C}(F_1(x_1), \dots, F_n(x_n)), \quad (8)$$

when cumulative distribution functions F_1, \dots, F_n are continuous, there exists a unique copula \mathcal{C} .

In order to specify the total value of operational risk for the whole bank, which takes into account the degree of dependence among the aggregated distributions of losses, copula function \mathcal{C} is applied. It is given as follows:

$$\mathcal{C}_t(u_1, \dots, u_n) = t_{v,P}(t_v^{-1}(u_1), \dots, t_v^{-1}(u_n)), \quad (9)$$

where $t_v(x)$ is a cumulative distribution function of t -Student distribution with v -degrees of freedom, and $t_{v,P}(x_1, \dots, x_n)$ is a n -dimensional cumulative distribution function of t -Student distribution with v degrees of freedom and correlation matrix \mathbf{P} .

The article uses the procedure of estimation of copula function \mathcal{C} parameters that is presented in [Di Clemente, Romano 2004, p. 11].

For the given correlation matrix \mathbf{P} for copula function \mathcal{C}_t , one can determine the total distribution of loss on the basis of the procedure presented below. Such a total distribution takes into account the character of dependences described by fitted function \mathcal{C} .

3.2. Determination of total distribution of dependent random variables

The procedure described below concerns only the two-dimensional case because such a case will be analysed in part 3.

For given distributions of aggregated values of loss for two categories of risk L_1 and L_2 , with cumulative distribution functions F_1 and F_2 and estimated parameters for copula function \mathcal{C}_t , i.e. correlation matrix \mathbf{P} as well as value v , one should:

1) generate a random vector for function \mathcal{C}_t , with marginal values u_1 and $u_2 \sim U(0; 1)$;

2) with the use of $F_i^{-1}(u_i) = \inf\{x: F_i(x) \geq u_i\}$ dla $i \in \{1,2\}$, determine a vector of random variables from step 1 with distributions defined by cumulative distribution functions F_1 and F_2 ; defined in such a way, losses l_1 and l_2 for two

categories of risk will constitute a two-dimension random variable of dependent loss losses;

3) add up elements of the achieved loss vector from step 2 to achieve one hypothetical realisation of the total loss that may occur for the considered categories of risk;

4) repeat steps 1–3 many times to achieve distribution of total losses. On the basis of such distribution, one determines the value of VaR measure at the accepted significance level.

4. Empirical study

As mentioned in the introduction, the data concern operational losses for a period of 5 years. With regard to the mentioned aspect related to data confidence, information that may identify the real and approximate extent of loss in terms of frequency and severity cannot be presented. Estimated parameters of the analyzed distributions of frequency and severity, basic descriptive statistics of data sets, goodness-of-fit test χ^2 for the distribution as well as results of Kolmogorov-Smirnov and Anderson-Darling goodness-of-fit tests will not be presented.

In order to fit a relevant function of distribution of loss frequency, the 5-year-period was divided into months. As a result, the authors achieved 60 monthly observations of loss frequency for each of the analytical areas.

Two types of distributions were achieved for the data from these 60 months. Poisson distribution and negative binomial distribution were fitted in the case of both analytical areas. In each of the considered cases, the initial premise related to the selection of distribution form were fulfilled, i.e. value of variance $Var(N)$ was several times higher than the expected value $E(N)$. This was confirmed in the case of goodness-of-fit test at the non-specified values of distribution parameters, where null hypothesis informs about consistency of distribution in the sample with distribution in the general population. In both analytical areas, the test indicated a lack of grounds for the rejection of null hypothesis for the negative binomial; in the case of Poisson distribution, this test indicated the rejection of null hypothesis for the benefit of alternative hypothesis at significance level 0.05. Therefore, negative binomial distributions were selected for analyses of both analytical areas.

In cases of study of loss severity with respect to operational risk, the authors applied the approach based on the separate modelling of distribution tail by means of GPD distribution and the remaining lower part by means of log-normal distribution. In the case of the first analytical area, log-normal distribution was a referential distribution. The determined value from the above formula (6), the observation of which is modelled on the basis of GPD distribution, constituted 86 percentile (real value of threshold u in the formula of loss excess distribution cannot be directly given with regard to the confidential character of data).

A similar procedure was applied for the second analytical area. In this case, threshold value u in the formula of loss excess distribution was defined at the level of 89 percentile.

In order to model only a tail of an unknown distribution, the selection of threshold u should be based on the features of GPD distribution. In these cases the selection of threshold u may be made, e.g. on the basis of mean excess function or the Gertebgarbe-Werner graph [Embrecht et al. 1997].

Figure 1 and 2 present an empirical cumulative distribution function together with the fitted values of distributions: log-normal (marked as LN) for the whole scope of data modelled by LN distribution for the observation below value u and GPD distribution for the observation above threshold u (marked in the graph as the right tail of distribution GPD). The lack of scale for the values of random variables is to prevent the estimation of the financial extent of loss. The limited scope of value of the probability accumulated only for a value below 0.95 and 0.9, respectively for Figure 1 and Figure 2, is to allow for the estimation of extent of fit of the theoretical distributions in the right tail of the distribution.

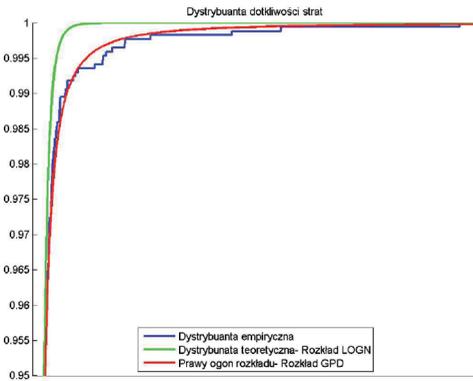


Figure 1. Distribution fit. First analytical area

Source: own calculations.

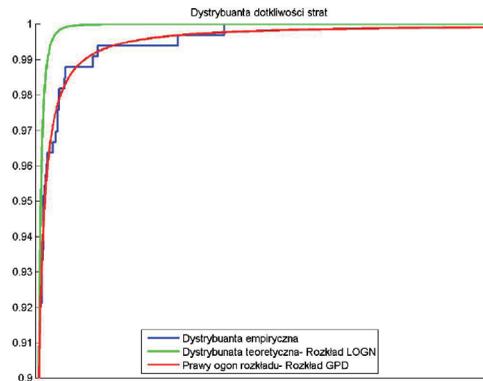


Figure 2. Distribution fit. Second analytical area

Source: own calculations.

Modelling of the distribution of loss severity in total by means of log-normal distribution does not give the desired effects in the situation when the main point of interest is a good fit in the right tail of the distribution. Extreme loss will significantly determine the value of VaR, defined on the basis of aggregated distribution of loss. Separate modelling of the tail itself by means of GPD distribution provides much better results. Extreme losses that seldom occur may exceed the determined VaR value, even for high significance levels, if the distributing of loss severity is not sensitive to extreme loss, as it is in the case of log-normal distribution.

Quantification of total Value-at-Risk

An additional aspect in the article concerns the application for copula functions in studying dependences among various analytical sections and their influence on the total value of risk, i.e. value of VaR.

The correlation coefficient estimated on the basis of the method mentioned in part 2 defines the level of dependence between the analytical sections. In this case it equals $r_{12} = 0.32$. The following tables present the percentage values that describe the relative increase of changes for VaR value.

With the use of value $VaR_{\alpha}(L_i^j)$ we defined the value of VaR for confidence level α of aggregated distribution L_i^j , where i -denotes the studied area of bank operation and j -denotes the accepted type of loss severity distribution (i.e. $j = LN$ for log-normal distribution and $j = GPD$ for GPD distribution that models the tail of loss severity distribution).

Table 1. Relative changes of VaR for the same values of i and changes with respect to α

	$\frac{VaR_{0.999}(L_1^{LN})}{VaR_{0.99}(L_1^{LN})}$	$\frac{VaR_{0.999}(L_1^{GPD})}{VaR_{0.99}(L_1^{GPD})}$	$\frac{VaR_{0.999}(L_2^{LN})}{VaR_{0.99}(L_2^{LN})}$	$\frac{VaR_{0.999}(L_2^{GPD})}{VaR_{0.99}(L_2^{GPD})}$
Relative change	138%	259%	362%	651%

Source: own calculations.

Table 1 presents percentage changes among VAR values for various values of percentiles within the same analytical areas and accepted severity distribution.

Table 2. Relative changes of VaR for the same values of i and α change with respect to j

	$\frac{VaR_{0.99}(L_1^{GPD})}{VaR_{0.99}(L_1^{LN})}$	$\frac{VaR_{0.999}(L_1^{GPD})}{VaR_{0.999}(L_1^{LN})}$	$\frac{VaR_{0.99}(L_2^{GPD})}{VaR_{0.99}(L_2^{LN})}$	$\frac{VaR_{0.999}(L_2^{GPD})}{VaR_{0.999}(L_2^{LN})}$
Relative change	127%	239%	291%	524%

Source: own calculations.

Table 2 presents percentage changes among VaR values for various loss severity distributions within the same area and level.

Similar combinations are presented in Table 3. They concern the application of copula functions in the determination of the total VaR value. The total value, with the given severity distributions j and confidence level α , is denoted as $VaR_{\alpha}^{C_r}(L_{1+2}^j)$, where the character of dependence is introduced by the considered copula function C , and as $VaR_{\alpha}^{covr=1}(L_{1+2}^j)$, where the character of dependence is not considered and there is assumption related to the perfect dependence among the analytical areas. The influence of diversification is noticeable and its power depends on the adopted distributions, analytical area and confidence level α . The summary is presented in Table 3.

Table 3. The effect of risk diversification depending on VaR confidence level and loss severity distributions

	$\frac{VaR_{0,999}^{Cr}(L_{1+2}^{LN})}{VaR_{0,99}^{corr=1}(L_{1+2}^{LN})}$	$\frac{VaR_{0,999}^{Cr}(L_{1+2}^{LN})}{VaR_{0,999}^{corr=1}(L_{1+2}^{LN})}$	$\frac{VaR_{0,99}^{Cr}(L_{1+2}^{GPD})}{VaR_{0,99}^{corr=1}(L_{1+2}^{GPD})}$	$\frac{VaR_{0,999}^{Cr}(L_{1+2}^{GPD})}{VaR_{0,999}^{corr=1}(L_{1+2}^{GPD})}$
Relative change	83.71%	90.13%	89.75%	95.03%

Source: own calculations.

5. Conclusions

In spite of the limited insight into the study due to security reasons and bank secrecy related to the publication of data and study results, it can be observed that the initial hypotheses are confirmed in the empirical study. First of all, in Tables 1 and 2 one can observe the increase of risk value in the situation when the modelling process includes the approach based on the modelling of GPD tail separately from its remaining part, where LN distribution is a referential distribution. Moreover, the second part of the study uses the effect of diversification in order to lower the total VaR value. Owing to the application of the copula function theory, it was easy to introduce dependences among the aggregated loss distributions. It is perfectly visible that the effect of diversification is suppressed by the adopted GPD distributions in the modelling of loss severity distribution tail, however, one can notice, in a maximum case, the decrease of total VaR value at the level of 10%.

It should be realised that the modelling based on a separate fit of GPD distribution to the distribution tail and further determination of aggregated distribution of loss shows pacifying trends with respect to risk measure. This means that the model will be very sensitive to extreme loss and VaR value may be at a much higher level than in the case of models based on one parametric distribution for the whole scope of loss severity. One should also take into consideration such a high confidence level in the quantification of operational risk (currently, value of $\alpha = 0.999$, defined by supervisory institutions) in the situation when the model focuses on the modelling of extreme loss, and all the changes in the distribution tail that result from even single events will be immediately reflected in the estimated risk values.

Literature

Chernobai A., Rachev S., Fabozzi F. (2007), *Operational Risk. A Guide to Basel II Capital Requirements, Models and Analysis*, Wiley & Sons, Hoboken, NJ.

Cherubini U., Luciano E., Vecchiato W. (2004), *Copula Methods in Finance*, Wiley & Sons, Chichester.

Di Clemente A., Romano C. (2004), A copula-extreme value theory approach for modelling operational risk, [in:] M. Cruz (ed.), *Operational Risk Modelling and Analysis. Theory and Practice*, Risk Books, London.

- Embrecht P., Kluppelberg C., Mikosch T. (1997), *Modeling Extremal Events*, Springer-Verlag, Berlin.
- Embrecht P., Hoing A., Juri A. (2003), Using copulae to bound the Value-at-Risk for functions of dependent risks, *Finance and Stochastics* 7.
- Klugman S., Panjer H., Willmot G. (2008), *Loss Models from Data to Decisions*, Wiley & Sons, Hoboken, NJ.

ZASTOSOWANIE UOGÓLNIONEGO ROZKŁADU PARETA I FUNKCJI ŁĄCZĄCYCH W ZAGADNIENIU RYZYKA OPERACYJNEGO

Streszczenie: Temat niniejszego artykułu dotyczy problemu modelowania ryzyka operacyjnego w banku. Obszar analizy dotyczy dwóch rozdzielnych obszarów analitycznych złożonych z pewnych kombinacji kategorii ryzyka macierzy bazylejskiej. Uwagę skupiono na modelowaniu rozkładów dotkliwości strat w modelach LDA oraz uwzględnieniu siły i charakteru zależności pomiędzy badanymi obszarami analitycznymi. Do modelowania pojedynczego rozkładu dotkliwości strat wykorzystano podejście oparte na teorii wartości ekstremalnych EVT i rozkładzie GPD do modelowania jego prawego ogona. W przypadku uwzględnienia siły i charakteru zależności wykorzystano funkcję łączącą t -Studenta. Wyznaczone wielkości opisują w wartościach względnych efekty zastosowanego podejścia.

Słowa kluczowe: ryzyko operacyjne, funkcje łączące, rozkład GPD, VaR.