

The Randomly Stopped Sums (RSS) Model

1.1 Overall methodology of RSS Regression Model

Let $\{X_1, X_2, \dots\}$ be a sequence of independent and identically distributed continuous random variables representing the severity of operational losses with support $(0, +\infty)$, and let N be a counting random variable, i.e. an integer-valued, non-negative and non-degenerate at zero random variable, representing operational loss frequency in a certain period of time. Moreover, assume X and N are independent. The total amount of operational loss is a randomly stopped sum and can be expressed as¹:

Equation 1. Definition of the randomly stopped sums

$$S = \begin{cases} 0, & N = 0 \\ \sum_{i=1}^N X_i, & N > 0 \end{cases}$$

The distribution of this quantity is named a Randomly Stopped Sum (RSS) distribution, also known as a compound distribution in other financial literatures. The distribution function of the randomly stopped sum is shown in Equation 2.

Equation 2. The distribution function of the randomly stopped sum

$$F_S(s) = \Pr(S \leq s) = \sum_{n=0}^{\infty} \Pr(S \leq s | N = n) \Pr(N = n)$$

Clearly, we can also view the distribution as a bivariate joint distribution $f_{S,N}(s, n) = f_{S|N}(s) f_N(n)$. Since the distribution function involves an infinite sum, the functional form cannot be generally expressed as a closed form. An important observation from this distribution is that there could be a single probability mass at 0, at the lower bound of the support, since losses may not occur in a year in a given risk segment. This is in contrast with most of the loss distributions such as Lognormal whose support doesn't include 0. The probability mass at 0 is simply the probability mass of the count distribution at 0. For instance, if the count distribution is chosen to be Poisson, the mass at 0 will be:

Equation 3. The probability of the randomly stopped sum at 0

$$\Pr(S = 0) = e^{-\lambda},$$

where λ is the parameter of the Poisson distribution. Thus, while the RSS model is applicable to all risk segments, it is a model of choice especially when excess 0 counts are present in the data.

¹ S. Klugman, H. Panjer, G. Willmot. Loss Models: From Data To Decisions. 3rd ed, 1998. §9.3

The central tendency of the frequency and severity components in the model can be described by a set of covariates, i.e. the Federal Reserve Released macroeconomic and the bank specific predictor variables in stress testing models. Therefore, the expected total operational losses in a defined time period can be forecasted by the model.

Let μ_N denote the mean of the frequency distribution, then $h(\mu_N) = x'\beta$. The mean is determined by a vector of predictor variables x' , β is a vector of the coefficients, and $h(\cdot)$ is a twice differentiable monotonic link function. If the distribution of the loss severity is from the exponential family, the sum of n i.i.d losses will still be the same distribution with mean $n\mu_X$ due to the independence of frequency and severity, with μ_X the mean of the severity distribution. Thus, given n losses occur in a quarter, the expected total loss would be $n\mu_X$. As such, the mean of the total loss can be described as $g(\mu_S) = z'\gamma + \ln(n)$, where z' is a vector of variables that predicts mean loss, γ is the vector of coefficients and $g(\cdot)$ is a twice differentiable monotonic link function. The frequency and total loss models may use a same or different set of predictors. The link functions could be arbitrary, but some are commonly used. Natural log is one of the frequently used link functions. If natural log is used for both frequency and severity models, then the expected total loss can be conveniently expressed as $e^{x'\beta + z'\gamma}$. In our application, the log link function will be chosen modeling. So, the offset in aggregate loss estimation can be expressed as $\ln(n)$ as shown above.

Let the pdf of loss severity X be defined as $f_X(x; \mu_X)$, then $S|(N = n) \sim f_X(s; n\mu_X)$. Now we have

Equation 4. Conditional distribution of the sums of severity

$$f(s|N) = \begin{cases} 1, & s = 0 \text{ and } n = 0 \\ f_X(s; n\mu_X), & s > 0. \end{cases}$$

The log likelihood of observing the total operational loss is $\ln[f_{S,N}(y, n)] = \ln[f_{S|N}(s)f_N(n)] = \ln[f_X(s; n\mu_X)] + \ln[f_N(n; \mu_N)]$. As a result, the maximum can be achieved in two steps: first by maximizing $\ln[f_N(n; \mu_N)]$ with respect to the parameters involved in the determination of μ_N , i.e. β ; and then by maximizing $\ln[f_X(s; n\mu_X)]$ with respect to the parameters involved in the determination of μ_X , i.e. γ .

The definition of pdf in Equation 4 implies that when the sum of loss is zero, in other words when the loss count is 0, the observation does not contribute to the log likelihood function in severity model estimation, as should be expected. This is because non-zero loss severities are only observed contingent on the occurrence of losses. The sum of loss modeling thus will be performed only on non-zero losses.

Given a set of predictor variables, the severity could be distributed as other continuous distributions such as Lognormal. Lognormal distribution is not closed under convolution. One of

the solutions is to logarithm transform the aggregate loss to make the loss normally distributed. Then, the normally distributed variable can be modeled as usual.

1.2 The Rationale for Developing the RSS Model

Aggregate operational loss can be decomposed as a frequency component and a severity component. The frequency component takes non-negative integer values, and the loss severity must be positive values. These two components can be assumed to be independent and to be modeled separately. The loss frequency is usually described by discrete distributions such as the Poisson distribution, and the severity of losses can be delineated by various continuous distributions such as Gamma and Inverse Gaussian distributions. The definition of randomly stopped sums model suits such a setting well.

1.3 The Forecast from RSS

The RSS model produces forecasts not only for the loss frequency but also for the aggregate loss, in contrast to the loss frequency models. The forecasted expected losses from RSS are always positive.

1.4 Literature Review

The application of RSS can be found in research areas such as insurance and psychological studies. In the paper of G. Heller², et al., the randomly stopped sums model was used to model the aggregate loss or the pure premium in general insurance (e.g. property, casualty, and liability) portfolios. In such portfolios, the aggregate losses can be further decomposed into and modeled as a frequency component and a severity component, respectively.

M. Smithson³ and Y. Shou described the application of the model in psychology. The Smithson paper proposed the randomly stopped sums model and demonstrated the application of the model to three psychological experiments including the total duration of eye fixation on print advertisements, the eye movement patterns in viewing portal web page, and total time spent in making charitable donation decisions. Evaluation and diagnosis of the sums models of these examples were also provided. The authors recommended the model be used in many other psychological applications.

Operational loss closely resembles the applications introduced above. Thus, RSS can also be applied.

1.5 Advantage of RSS

With the RSS model, the total operational loss severity in a quarter can be modeled. The information that comes from frequency and severity is used in its entirety. The fitted regression

²G. Heller, M. Stasinopoulos and R. Rigby. Randomly Stopped Models. Paper attached to the appendix.

³M. Smithson and Y. Shou (2014). Randomly stopped sums: models and psychological applications. *Front Psychol.* 2014; 5: 1279.

model to loss severity reflects the dynamics of the loss severity over the model development and the forecast period. If the regression model can be established, the variance of the severity prediction can also be reduced compared with the unconditional estimates, improving the accuracy of aggregate loss prediction. Furthermore, the RSS model incorporates the interconnectedness between loss severity and possible predictor variables as shown in last section.

Operational loss distribution is usually skewed. Linear regression can suffer from problems such as the violation of normality. Due to the high skewness of data and the relatively small sample size, asymptotic normality might not hold. Also, negative losses could be produced in loss forecasting with linear regression or robust regression.

The severity modeling in RSS provides a good solution. The severity distributions in RSS, such as Inverse Gaussian or Lognormal, can be highly skewed. With carefully selected link functions, the forecasted losses will be always positive, in line with intuition. Again, in the application of the RSS model in this document, log link will be used.

General GLM features are also applicable to RSS models, so the evaluation of the model can be done following regular GLM procedures.

The model is particularly valuable when the total loss exhibits excess zero losses across the loss history as discussed in section1.1.