

# Combining data sources to be used in quantitative operational risk models

**M Gericke**

 [orcid.org/0000-0002-8604-624X](https://orcid.org/0000-0002-8604-624X)

Dissertation accepted in partial fulfilment of the requirements for  
the degree *Master of Science in Risk Analytics* at the North-  
West University

Supervisor: Prof PJ de Jongh

Co-supervisor: Prof H Raubenheimer

Graduation December 2022

12377023

## **Acknowledgements**

I would like to express my sincere gratitude towards my supervisor, Prof Riaan de Jongh. His patience and guidance in the completion of this dissertation was invaluable and it has been a privilege to learn from him.

I would also like to thank all my colleagues at the Centre for BMI, but especially Prof Helgard Raubenheimer for his valuable input and continuous support on this journey.

## **Abstract**

The management of financial losses is crucial for banks as they are required to set aside regulatory capital to absorb unexpected losses. Banks also need to calculate economic capital to ensure solvency according to their own risk profile. The main financial risks faced by banks are market, credit, and operational risk. Operational risk, the focus of the dissertation, includes fraud, improper business practices, regulatory risk, and others. Barings Bank's loss of over USD1 billion due to rogue trading activities is well-known, but an extreme example of such risk.

In order to calculate capital to withstand this risk, the aggregate distribution of operational losses for the next year is estimated. This distribution needs to be estimated in a forward-looking manner and for this, assessments by experts are often used. The extreme quantiles of this distribution are of specific interest. For instance, a bank should hold capital to survive a one-in-a-thousand-year aggregate operational loss (the 99.9% Value at Risk of the distribution). A methodology is described to calculate capital for different operational risk categories.

Banks often only have limited internal data available to accurately model the distribution and therefore use external sources and scenario assessments to supplement their data. Statistical methods are explored that could be used to combine limited historical data and scenario assessments provided by experts, to estimate the extreme quantiles of the aggregate distribution. This provides a way of constructing forward-looking distributions to calculate risk capital.

SAS® OpRisk Global Data is used to demonstrate how external data can be used and scaled for use in the risk modelling process. Some measures are suggested that could be used to challenge experts to adjust their scenario assessments based on available historical data.

The main contribution of the research is to provide a holistic view of how internal data, external data and scenario assessments can be used to create a consistent framework for modelling operational risk capital within a bank or other financial institution.

**Keywords: Operational risk management, operational risk quantification, operational risk measurement, capital models, loss distribution approach, external data, scenario assessments.**

## Table of contents

1.	Introduction .....	9
2.	Background .....	12
2.1	Introduction .....	12
2.2	An overview of risk management.....	12
2.3	A short history of operational risk.....	14
2.4	Operational risk categories.....	16
2.5	Regulation.....	18
2.5.1	The Basel Accords .....	18
2.5.2	Methods for calculating operational risk capital.....	19
2.6	Conclusion.....	21
3.	Quantitative risk modelling methodology .....	22
3.1	Introduction .....	22
3.2	Loss distribution approach (LDA).....	22
3.2.1	Frequency modelling .....	23
3.2.2	Severity modelling .....	24
3.2.3	Aggregate loss distribution .....	26
3.3	Capital estimation.....	28
3.4	Data sources .....	28
3.4.1	Internal data .....	29
3.4.2	External data.....	29
3.4.3	Scenario analysis.....	30
3.4.4	Business environment and internal control factors (BEICF's) .....	31
3.4.5	Combining sources of data .....	31
3.5	Conclusion.....	32
4.	Constructing forward-looking distributions using limited historical data and scenario assessments .....	33
4.1	Introduction .....	33
4.2	Scenario assessments .....	33
4.3	Estimating VaR.....	35

4.3.1	Naïve approach .....	35
4.3.2	The GPD approach .....	39
4.3.3	Venter’s approach .....	42
4.4	Conclusion.....	45
5.	Using external data sources to inform scenario assessments.....	46
5.1	Introduction .....	46
5.2	External databases.....	47
5.3	SAS® OpRisk Global Data .....	48
5.4	Preliminary data analysis and determination of explanatory variables.....	49
5.5	Allowing for reporting bias .....	54
5.6	Allowing for scaling bias .....	57
5.7	Model application.....	58
5.8	Model diagnostics and results .....	63
5.9	Modelling above a threshold.....	64
5.10	Results for an individual bank.....	66
5.11	Model application within a bank .....	71
5.12	Conclusion.....	72
6.	Concluding remarks and recommendations for future research .....	73
	Reference list .....	75
	Appendix A: Standardised measurement approach.....	78
	Appendix B: Guidelines for using the advanced measurement approach .....	80

## List of tables

Table 1 (a): Probability density functions .....	24
Table 1 (b): Probability distribution functions .....	24
Table 2: Breakdown of losses per business line .....	49
Table 3: Summary statistics per geographical region .....	51
Table 4: Summary statistics per event type.....	53
Table 5: Region dummy coding.....	60
Table 6: Business line dummy coding .....	60
Table 7: Event type dummy coding .....	61
Table 8: Estimated parameter values for the final model .....	62
Table 9: Adjusted probabilities for different values of $F1$ .....	66
Table 10: Bank of America Corporation loss data points per business line and event type ...	67
Table 11: Re-estimated parameter values for the model, excluding the Individual bank's data .....	68
Table 12: Estimated scenario points per business line for different models .....	70
Table A1: Marginal coefficients to calculate BIC .....	78

## List of figures

Figure 1:	Risk management process .....	13
Figure 2:	Illustration of the effects of VaR estimation using the naïve approach .....	38
Figure 3:	Illustration of VaR estimates obtained from a GPD fit on the oracle quantiles ...	41
Figure 4:	Comparison of VaR results for the GPD and Venter approaches .....	44
Figure 5:	Normalised quantile residual plot .....	63
Figure 6:	QQ plot of simulated losses vs observed losses .....	64



## List of abbreviations

AMA	: Advanced Measurement Approach
BEICF's	: Business environment and internal control factors
BIA	: Basic Indicator Approach
BCBS	: Basel Committee on Banking Supervision
ED	: External data
EVT	: Extreme value theory
GAMLSS	: Generalised Additive Models for Location Scale and Shape
GEV	: Generalised Extreme Value distribution
ILD	: Internal loss data
LDCE	: Loss Data Collection Exercise
MC	: Monte Carlo
ORC	: Operational risk category
ORX	: Operational Riskdata eXchange Association
SLA	: Single loss approximation
TSA	: Standardised Approach
UoM	: Unit of measure
VaR	: Value at Risk

## 1. Introduction

Financial risk is defined as an event or action that may adversely affect an organisation's ability to achieve its objectives and execute its strategies. The main financial risks faced by banks are market, credit, and operational risks. To protect against risk, or the probability of a shortfall in assets compared to liabilities, financial institutions need to set aside excess money, or risk capital. Operational risk capital is the focus of this dissertation and a short history and different categories of operational risk will be described in more detail in Chapter 2.

This dissertation is concerned with the modelling of operational risk capital. Specifically, the focus is on how various data sources available to banks or other financial institutions may be utilised to obtain estimates for the appropriate amount of operational risk capital to be set aside.

Two types of capital are important, namely regulatory and economic capital. Regulatory capital is a requirement by the regulator to guard against the collapse of the banking system and to ensure that banks remain solvent. In the face of unexpected losses, regulatory capital should absorb these losses so that the bank remains solvent and it is therefore important that banks do not underestimate the capital amount. Economic capital, on the other hand, is set aside to ensure a particular credit rating by the external rating agencies. The higher the rating that is required, the higher the capital amount to be held. Both regulatory and economic capital are held in liquid assets, and the return on investment is typically lower than what may be obtained in riskier investments. For this reason, banks can also not afford to overestimate the capital amount.

To ensure solvency, as well as return for shareholders, regulatory and economic capital should be determined as accurately as possible. Under Basel, the regulator has proposed various methods for calculating regulatory capital, namely a basic and standardised approach, as well as an advanced measurement approach (AMA). Although banks will be required to calculate operational risk capital using a new risk-

sensitive standardised measurement approach (SMA) from 1 January 2023, the AMA is still of interest. This approach is also used to estimate economic capital and is based on the so-called loss distribution approach (LDA) where the distribution of future losses is typically constructed using internal loss data, external data and scenario assessments by experts.

Under the LDA, an annual aggregate operational loss distribution is constructed. The tail of the distribution is of most importance, as the extreme quantiles of the distribution is required to calculate capital.

Basel II prescribes that a bank should hold sufficient capital to protect them against a one-in-a-thousand-year aggregate loss when explicitly dealing with operational losses, i.e., the 99.9% Value-at-Risk of the aggregate operational loss distribution. Ideally, should the bank have a thousand years of historical data, the bank can merely determine the maximum loss it had experienced during this time to determine the capital requirement. Another way to interpret this requirement is that given 1,000 banks with identical size and operational risk profiles, the regulatory capital requirement is the maximum aggregate operational risk loss amount experienced across all 1,000 banks in a given year. Given that banks generally only have 10 years' of operational risk loss data, one could potentially approach the problem using 100 banks with similar profiles, giving 1,000 annual aggregated operational risk loss amounts.

For this reason, the Basel Committee on Banking Supervision (BCBS) (2011b) suggests that banks should use loss data from external sources and scenario data in addition to their own internal loss data and controls to construct statistical models. However, there are significant challenges that banks need to address when combining data elements. BCBS (2011b) advises that the combination of data elements should be based on sound statistical methodologies. The focus of this dissertation is to recommend strategies for the optimal use of the various data sources in economic capital estimation and to suggest possible improvements to the loss distribution approach.

This dissertation is organised as follows: In the following chapter we will provide an overview of operational risk management, a short history of operational risk and a discussion on different operational risk categories. We also provide some background on the regulatory requirements for calculating operational risk capital. In Chapter 3, a methodology is described that could be used to calculate operational risk capital based on the loss distribution approach. Then, in Chapter 4, a recently proposed statistical modelling approach for the construction of forward-looking distributions, important in the loss distribution approach, is reviewed. The method combines limited internal loss data with the scenario assessment by experts and also incorporates a measure of agreement between the two data sources. This methodology has already successfully been implemented by two South African banks for modelling operational risk capital. Lastly, in Chapter 5 we use SAS® OpRisk Global Data to demonstrate how external data may be used to improve economic capital estimation. Banks often only have limited historical data of their own operational losses, and external data sources could be used and scaled to improve the banks statistical models. Finally, we make suggestions on how external data could be used to challenge experts to adjust their scenario assessments. Some concluding remarks and ideas for future research are made in the final chapter.

## **2. Background**

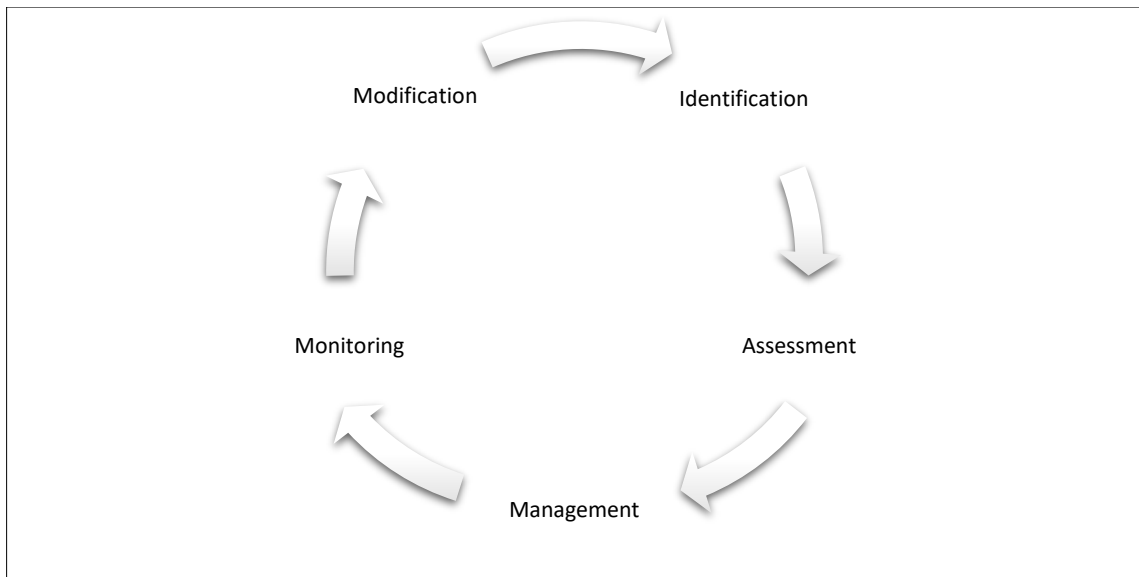
### **2.1 Introduction**

The focus of this dissertation is on quantitative operational risk models. A detailed review of the literature related to operational risk models was conducted and some of our findings are summarised in this chapter. The reader will be introduced to the risk management process and the increased importance of the management of operational risk will be highlighted. A short history of operational risk is given before the different risk categories are discussed. We also provide some background on the regulatory requirements for calculating operational risk capital. The reader that is familiar with operational risk may elect to glance through Chapter 2 and move to Chapter 3.

### **2.2 An overview of risk management**

McNeil *et al.* (2015) define financial risk as “any event or action that may adversely affect an organisation’s ability to achieve its objectives and execute its strategies”, or as “the quantifiable likelihood of loss”.

The process of managing risk is vital to understand the range of risks that an organisation may face at any point in time, and to realise that new risks may develop over time. Sweeting (2011) describes the risk management process as cyclical without a clear start and end, and Figure 1 gives a graphical presentation of the process before each stage of the process is briefly described.



**Figure 1: Risk management process**

The first stage in the risk management process is to identify all the risks faced by an organisation. This step also involves grouping risks in a coherent fashion and recording them consistently. This stage of the process is of utmost importance when considering a bank's internal data and grouping different risks into homogeneous risk categories as described later in Section 2.4. It is also critical if the organisation contributes its own data to external data sources or forms part of a data consortium. The use of internal and external data will be discussed in more detail in Chapter 3.4, and consortium data is specifically discussed in more detail in Section 5.2.

The assessment of each identified risk is the next step in the risk management process. Risk assessment includes deciding whether a risk can be quantified and how to sensibly aggregate risks. It also involves specifying risk measures to be used and acceptable values of those measures. The specific risk measure we will use in our risk capital models is the 99.9% Value at Risk.

The next step, risk management, involves responding to each risk, either by accepting, reducing, or removing the risk. The management stage is not the final step, and the treatment of each risk should be reviewed and adjusted where needed. The ongoing monitoring or review of the inputs and outputs of the process is important. Monitoring

does not only require that losses from risks should be carefully reported to an organisation's internal stakeholders, but also to external stakeholders such as regulators and shareholders. The final stage in the risk management process is modification, which involves a frequent review of the process and all its components.

The quantitative risk models referred to in this dissertation primarily forms part of the assessment stage of the risk management cycle described above. However, it will become clear in the remainder of this dissertation that there are touchpoints to almost every other stage of the process. We have already mentioned that the collection of internal data to be discussed in more detail later, forms part of the identification stage. In Section 2.5, the applicable regulation that forms part of the monitoring stage is discussed. Before we expand on these topics, we first give a short history of operational risk.

### **2.3 A short history of operational risk**

Banks and other financial institutions are faced with different risks that could be categorised under the main headings of market risk, credit risk, and operational risk. Market risk can be described as the risk inherent from exposure to capital markets, whereas credit risk refers to the risk associated with the default of a third party on a contract. Operational risks include a group of risks that impact on how a firm carry on business and include many different risks that often overlap each other to a significant extent (Sweeting, 2011). The Basel Committee on Banking Supervision (2006) defines operational risk as “the risk of loss resulting from inadequate or failed internal processes, people and systems or from external events. This definition includes legal risk, but excludes strategic and reputational risk.”

Operational risk management is the youngest of the three major risk branches, according to Peters *et al.* (2016). In the late 1990s, operational risk had a negative definition, namely “any risk that is not market or credit risk”, which they state was not very helpful to assess or manage operational risk. In the article “The invention of operational risk”, Power (2005) points out that although the generic term ‘operations risk’ had already been officially coined in 1991, it did not acquire widespread currency

until the mid to late 1990s when the Basel II proposals were developed and finally published in June 1999. A short history of the Basel Accords and the development of regulation around operational risk is discussed in more detail under Section 2.5 of this chapter.

Cummins *et al.* (2006) mention many very costly and highly publicised operational events that have heightened managerial and regulatory focus on operational risk. The infamous bankruptcy of Barings bank in 1995, which was triggered by a \$1.3 billion loss due to the actions of a rogue trader, is often cited in the literature as one of the best examples of operational risk. Cummins *et al.* (2006) also mention a loss of \$750 million at Allied Irish Bank in 2002 due to unauthorised trading and \$1.4 billion in fines levied in 2002 against several leading brokerage firms in the US for issuing misleading research reports to investors. Among insurance companies, the most significant operational losses were the Prudential Insurance Company of America which had to pay \$2 billion to settle claims for sales abuse during the late 1990s. Another example is State Farm Insurance which paid \$1.2 billion to motor insurance policyholders resulting from a breach of contract lawsuit in 1999. Cummins *et al.* (2006) argue that the increased focus on operational risk at the start of the 21st century likely emanated from two key developments: First, there was more emphasis on transparency in the financial reporting of organisations. Second, financial services organisations started using increasingly complex production technologies and concomitantly increased their exposure to operational risk.

The financial crisis of 2008 also brought new focus to the area of operational risk. De Jongh *et al.* (2013) conducted a detailed review of the impact of operational risks on the financial crisis. They describe the 2008 financial crisis as the worst crisis ever from an operational risk viewpoint and state that the amount associated with operational risk losses observed in 2008 was almost four times more than those observed in 2007. They conclude that the general public was more aware of the complexity of the global environment after the crisis. Therefore, it was no surprise that regulatory bodies also needed to strengthen the capital requirements for banks. The development of these regulatory requirements is discussed in more detail in the following sections.



Technology risk is another example of operational risk that has become increasingly important. Over the last two decades, banks and other financial organisations have benefited from an increased application of technology to financial services, with a concomitant rise in technology-related threats. The Covid-19 pandemic has exacerbated these operational risks and increased economic and business uncertainty. The pandemic has not only affected banks' information systems, personnel, and facilities but there has also been an increase in cyber threats like ransomware attacks and phishing. In addition, the potential for operational risk events caused by people, failed processes, and systems have increased due to greater reliance on working from home arrangements. The different operational risk event types are elaborated on in Section 2.4.

#### **2.4 Operational risk categories**

Operational risk is by definition heterogeneous and includes various risks associated with people, processes and systems. Kelliher *et al.* (2016) explain that it covers risks as diverse as systematic processing errors, cybercrime, health and safety breaches and product literature failings, to name a few. They, therefore, suggest that there is a need to break these risks down into more homogenous categories to enable exposures to be adequately understood and modelled. This process depends on a robust categorisation system with little scope for ambiguity regarding how losses and risks should be categorised.

Embrechts and Hofert (2011) explain that an operational risk category (ORC) is the level at which an organisation's model generates a separate distribution for estimating potential operational losses (e.g., the organisational unit, operational event type, or risk category). These models and statistical distributions will be expanded on in Chapter 3.

The ORC is also sometimes referred to as a unit of measure (UoM). Granularity is the term used to reflect the degree to which individual operational risk exposures are modelled. The lowest level of granularity implies using a single ORC to measure the organisation-wide exposure and allows all loss data to be pooled. On the other hand,

a high level of granularity introduces the challenge of adequately categorising sources of operational risk and may also pose a challenge when aggregating the different risk exposure estimates.

The BCBS (2019) does not prescribe the level of granularity required to model operational risk but merely states that the risk measurement system should be sufficiently 'granular' to capture the significant risk drivers affecting the shape of the tail distribution. Regulations require a bank to provide adequate analysis to show that its selection of ORCs is appropriate and a true reflection of its risk profile (e.g., analysis of the influence of variability in the ORC selection and correlation in and between ORCs). Embrechts and Hofert (2011) also suggest that the choice of granularity should be adequately supported by quantitative and qualitative analysis. The individual losses within a given ORC should be independent and identically distributed.

It is further stipulated that a bank must have well documented, objective criteria for allocating losses to specified business lines and event types. However, it is left to the bank to decide how it applies these categorisations in its internal operational risk measurement system.

Under Basel II (BCBS,2006), the proposed business lines are:

- Corporate finance,
- Trading and sales,
- Retail banking,
- Commercial banking,
- Payment and settlement,
- Agency services,
- Asset management, and
- Retail brokerage.

The following seven operational event types are considered:

- Internal fraud,
- External fraud,
- Employment practices and workplace safety,
- Clients, products and business practices,
- Damage to physical assets,
- Business disruption and system failures, and
- Execution, delivery and process management.

Using the eight business lines and seven event types specified by Basel, a bank should have 56 suggested ORCs. However, the actual number of ORCs will depend on the specific activities of the bank and their historical losses captured. The operational risk models discussed in later chapters are typically constructed per individual ORC, before the results are then aggregated.

Before we consider the methodology to construct risk models, we first give a brief history of the relevant regulation and how it applies to operational risk management.

## **2.5 Regulation**

The Basel Accords refer to a set of banking supervision regulations set by the Basel Committee on Banking Supervision (BCBS). In the following few sub-sections, the history and development of the capital requirements under the three accords are briefly discussed before the focus is shifted to the regulation related explicitly to operational risk (Basel Committee on Banking Supervision, n.d.).

### **2.5.1 The Basel Accords**

#### **Basel I**

The first Basel Capital Accord was released to banks in July 1988 and called for a minimum ratio of capital to risk-weighted assets (RWA) of 8%. Credit risk was the focus of the first Accord, and several amendments were made to the initial document in the early 1990s. In 1996, the Market Risk Amendment was issued, introducing a capital

requirement for the market risks arising from banks' exposures to foreign exchange, traded debt securities, equities, commodities, and options. For the first time, banks were allowed to use internal models (Value-at-Risk models) as a basis for measuring their market risk capital requirements, subject to strict quantitative and qualitative standards.

### **Basel II**

Basel II, an extension of Basel I, was introduced in 2004. Basel II created a more comprehensive risk management framework by creating standardised measures for credit and market risk, and for the first time also for operational risk. Under Basel II, the revised capital framework contained three pillars: minimum capital requirements, supervisory mechanisms and transparency, and market discipline.

### **Basel III**

The financial crisis of 2008 exposed the weaknesses of the international financial system and ultimately led to the creation of Basel III. The new standards under Basel III were first published in December 2010, and most of the reforms are still being phased in. The enhanced Basel framework revised and strengthened the three pillars established by Basel II and extended it in several areas. The BCBS completed its Basel III post-crisis reforms in 2017, with the publication of new standards for calculating capital requirements for credit risk, credit valuation adjustment risk and operational risk. The regulatory framework's revisions aimed to restore credibility in the calculation of RWA by enhancing the robustness and risk sensitivity of the standardised approaches for credit risk and operational risk and constraining internally modelled approaches and complementing the risk-based framework with a revised leverage ratio and output floor.

#### **2.5.2 Methods for calculating operational risk capital**

Basel II, still in force, provides three methods for calculating operational risk capital charges. These methods increase in sophistication and risk sensitivity: The Basic Indicator Approach (BIA), the Standardised Approach (TSA) and the Advanced Measurement Approach (AMA).

The BIA focuses on the bank's gross income to indicate its potential risk profile for operational risk losses. With the BIA, the bank must hold a fixed percentage of the average of three years' gross income as operational risk capital. For the TSA, a similar approach is followed to that of the BIA, except that the gross income element receives a more granular treatment in that the bank's respective business lines are isolated and specific percentage charges are multiplied by the gross income for those individual business lines.

In terms of Basel II, and still an option for banks to calculate their regulatory operational risk capital until 1 January 2023, is the use of the advanced measurement approach (AMA). Under the AMA, banks are allowed to use their own internal models to calculate risk capital. The Basel Committee on Banking Supervision (2006) allows greater flexibility for modelling practice when a bank uses the AMA. Subject to regulatory approval, the AMA permits the bank to directly analyse and model in detail its own operational risk profile.

In December 2017, as part of the post-crisis reforms, the BCBS published the rules for calculating operational risk capital effective 1 January 2023. The new risk-sensitive standardised measurement approach (SMA) will replace the existing standardised approaches (BIA and TSA) as well as the advanced measurement approach (AMA) that is based on banks' own internal models (BCBS, 2017). The interested reader is referred to Appendix A for a description of the methodology under the new SMA that should be used by banks to calculate their regulatory operational risk capital, as well as critique from industry against this approach. However, the remainder of this dissertation will focus on the methodology of developing risk models under the AMA.

Although the more sophisticated internal models developed under the AMA approach will no longer be allowed in determining minimum regulatory capital, these models will remain relevant for determining economic capital and influence decision making within banks and other financial institutions. According to the Bank of England's Prudential Regulation, regulators rely on these more advanced models for the

supervisory review process (Prudential Regulation Authority, 2020). It is also suggested that models based on the loss distribution approach (LDA) will continue to form an integral part of the supervisory review of a bank's internal operational risk management process. For this reason, we believe the LDA remains relevant and will continue to be studied and improved. Additionally, the SMA is calculated at an overall bank level which raises the question of how regulatory capital is allocated to lower levels of the organisation. The continued use of the LDA can then be used as an allocation tool for capital. The LDA is discussed in more detail in Chapter 3.

In addition, the principles applicable to the AMA are sound and remain applicable to internal risk models. For this reason, some of these principles from the guidelines to consider when developing operational risk models have been included in Appendix B.

## **2.6 Conclusion**

This chapter provided some background on operational risk capital and some of the concepts will be expanded on in the remainder of this dissertation. In Chapter 3, a review of the loss distribution approach is provided and a comprehensive methodology is discussed to calculate operational risk capital within a financial organisation.

### **3. Quantitative risk modelling methodology**

#### **3.1 Introduction**

In this chapter, we describe a methodology that can be used to calculate a bank's operational risk capital. This methodology is based on the loss distribution approach (LDA) which makes use of an annual aggregate loss distribution. The components of the aggregate loss distribution are discussed in more detail in this chapter and we also explain how various data sources can be used in the estimation of the severity distribution function.

#### **3.2 Loss distribution approach (LDA)**

The LDA is a popular method used by banks and other financial institutions to determine their operational risk capital. This approach is widely described in the literature (see, for example, Aue and Kalkbrener (2007), Benito and Lopez-Martin (2018), Lambrigger *et al.* (2007) and De Jongh *et al.* (2015)). The LDA is also sometimes referred to as the actuarial approach and apart from banks, general insurance companies often use this method to estimate claims against short-term insurance policies and determine the reserves needed to meet their obligations.

Aue and Kalkbrener (2007) argue that the LDA approach offers banks unparalleled flexibility in the way they determine their risk capital requirements. Under this approach, an organisation can estimate the probability distributions of both the severity and the one-year-event frequency using historical data. Having these two distributions, the organisation can then compute the probability distribution of the aggregate operational losses (Benito & Lopez-Martin, 2018).

Amin (2016) outlines many challenges associated with the LDA. For example, it does not differentiate between risks for which a large amount of historical data is available and those risks with minimal data. The methodologies discussed in Chapters 4 and 5 aim to address these challenges.

An aggregate loss distribution has to be determined for each of the ORCs within a bank. Based on the eight business lines and seven event types suggested by BCBS (2006) and discussed in Section 2.4, a bank would have a maximum of 56 ORCs. However, a bank may opt to increase granularity, for example, the “retail banking” business line can be broken down into different products such as vehicle finance, home loans, personal lending and credit card facilities if able to demonstrate that these products have different risk profiles, i.e. are heterogenous. The actual number of ORC’s will depend on the specific activities of the bank and the historical losses captured. The methodology to construct the aggregate loss distribution for each ORC is described in more detail below.

### 3.2.1 Frequency modelling

To predict the total loss amount that can be expected over one year in each ORC, we first need to estimate the annual frequency or number of operational loss events to occur in that specific ORC over the year. Let  $N$  be the random variable representing the annual number of loss events in an ORC. McNeil *et al.* (2015) suggest two possible probability mass functions to model  $N$ , namely the Poisson and the negative binomial distributions. They explain that using the Poisson model is very natural as a frequency distribution and is useful because of its aggregation properties discussed further in Section 3.2.3. The probability mass function of  $N \sim Poi(\lambda)$  is given by:

$$P(n) = \frac{e^{-\lambda} \lambda^n}{n!}, n = 0, 1, 2, \dots$$

Where sufficient historical data is available, the annual frequency in each ORC can be estimated by  $\hat{\lambda} = K/a$ , where  $K$  is the total number of loss events spread over  $a$  years. If no data is available for a specific ORC, the value for the frequency estimate could be determined as part of a scenario workshop, which will be discussed in more detail in Section 4.2.



### 3.2.2 Severity modelling

The next step is to determine a suitable severity distribution for losses within an ORC. Panman *et al.* (2019) showed that the severity distribution drives the shape of the aggregate distribution (see Section 3.2.3) and is therefore an important component in operational risk capital modelling.

Suppose the random variables  $Y_1, Y_2, \dots, Y_n$  denote the severities of the loss events in each ORC and assume that these loss events are independently and identically distributed. Suppose that the true severity distribution of  $Y_1, Y_2, \dots, Y_n$  is denoted by  $T$ . A suitable model for  $T$ , which can be a class of distributions  $F(y, \theta)$ , need to be determined, and the parameter(s)  $\theta$  would need to be estimated. Popular choices to model the severity of losses include the Burr, Gamma, generalised Pareto distributions (GPDs), Inverse Gaussian (Wald), Lognormal and Pareto distributions (De Jongh *et al.*, 2015). The probability density and distribution functions of each are given in Table 1 (a) and (b).

**Table 1 (a): Probability density functions**

Distribution	Par 1	Par 2	Par 3	Probability density function
Burr	$\mu > 0$	$\alpha > 0$	$\gamma > 0$	$f(y) = \alpha \gamma z^\gamma$
Gamma	$\mu > 0$	$\sigma > 0$		$f(y) = \frac{1}{\Gamma(\sigma)\mu^\sigma} y^{\sigma-1} e^{-\frac{y}{\mu}}$
Generalised Pareto	$\mu > 0$	$\xi > 0$		$f(y) = \frac{1}{\mu} (1 + \xi z)^{-1-\frac{1}{\xi}}$
Inverse Gaussian (Wald)	$\mu > 0$	$\alpha > 0$		$f(y) = \frac{1}{\mu} \sqrt{\frac{\alpha}{2\pi z^3}} \exp\left(-\frac{\alpha(z-1)^2}{2z}\right)$
Lognormal	$-\infty \leq \mu \leq \infty$	$\sigma > 0$		$f(y) = \frac{1}{y\sigma\sqrt{2\pi}} \exp\left(-\frac{(\ln(y) - \mu)^2}{2\sigma^2}\right)$
Pareto	$\mu > 0$	$\alpha > 0$		$f(y) = \frac{\alpha \mu^\alpha}{(y + \mu)^{\alpha+1}}$

**Table 1 (b): Probability distribution functions**

Distribution	Par 1	Par 2	Par 3	Probability distribution function
Burr	$\mu > 0$	$\alpha > 0$	$\gamma > 0$	$F(y) = 1 - \left(\frac{1}{1 + z^\gamma}\right)^\alpha$
Gamma	$\mu > 0$	$\sigma > 0$		$F(y) = \frac{\gamma(\sigma, z)}{\Gamma(\sigma)}$
Generalised Pareto	$\mu > 0$	$\xi > 0$		$F(x) = 1 - (1 + \xi z)^{-\frac{1}{\xi}}$
Inverse Gaussian (Wald)	$\mu > 0$	$\alpha > 0$		$F(y) = \Phi\left((z-1)\sqrt{\frac{\sigma}{z}}\right) + \Phi\left(-\left(z+1\right)\sqrt{\frac{\sigma}{z}}\right)\exp(2\sigma)$
Lognormal	$-\infty \leq \mu \leq \infty$	$\sigma > 0$		$F(y) = \Phi\left(\frac{\ln(y) - \mu}{\sigma}\right)$
Pareto	$\mu > 0$	$\alpha > 0$		$F(y) = 1 + \left(\frac{\theta}{y + \theta}\right)^\alpha$

Notes to Table 1:

- $z = \frac{y}{\mu}$ .
- $\mu$  denotes the scale parameter for all the distributions.
- $\gamma(a, b) = \int_0^b t^{a-1} \exp(-t) dt$ , the lower incomplete gamma function.
- $\Phi(y) = \frac{1}{2} \left(1 + \operatorname{erf}\left(\frac{y}{\sqrt{2}}\right)\right)$ , the standard normal cumulative density function.
- The function  $a(y, \phi)$  does not have an analytical expression and is evaluated using series expansion methods.

The body and tail of the severity distribution may also be modelled separately, i.e., a mixed distribution. The method to do this can be described as follows. Let  $q$  be a quantile of the severity distribution  $T$ . Here,  $q$  is the threshold that splice  $T$  in such a way that the interval below  $q$  is the expected part and the interval above  $q$  the unexpected part of the severity distribution.

Define two distribution functions:

$$\begin{aligned} T_e(y) &= T(y)/T(q) \text{ for } y \leq q \text{ and} \\ T_u(y) &= [T(y) - T(q)]/[1 - T(q)] \text{ for } y > q, \end{aligned} \quad (1)$$

i.e.,  $T_e(y)$  is the conditional distribution function of a random loss  $Y \sim T$  given that  $Y \leq q$  and  $T_u(y)$  is the conditional distribution function given that  $Y > q$ .

The following identity then exists:

$$T(y) = T(q)T_e(y) + [1 - T(q)]T_u(y), \text{ for all } y. \quad (2)$$

This identity represents  $T(y)$  as a mixture of the two conditional distributions. Instead of modelling  $T(y)$  with a class of distributions  $F(y, \theta)$ ,  $T_e(y)$  is modelled with  $F_e(y, \theta)$  and  $T_u(y)$ , with  $F_u(y, \theta)$ . For  $F_e(y, \theta)$ , the empirical distribution can be used, or De Jongh *et al.* (2015) has shown that the Burr distribution may also be a good choice (see Table 1). Borrowing from extreme value theory (EVT), a popular choice for  $F_u(y, \theta)$  could be the generalised Pareto distribution (GPD). According to the Pickands-Balkema-de Haan limit theorem (McNeil *et al.*, 2015), the conditional tail of all distributions in the domain of attraction of the Generalised Extreme Value distribution (GEV) tends to a GPD distribution. The distributions in the domain of attraction of the GEV are a broad class of distributions, which includes most distributions of interest. Again, De Jongh *et al.* (2015) suggested that the GPD distribution is a good choice. However, one could also consider alternative distributions for modelling the tail of a severity distribution, as is discussed in more detail in Chapter 5.

### 3.2.3 Aggregate loss distribution

Given the frequency  $N \sim Poi(\lambda)$  and loss severity function  $Y \sim T$ , the annual aggregate loss is obtained by  $A = \sum_{n=1}^N Y_n$  and the distribution of  $A$  is known as a compound Poisson distribution or the so-called aggregate loss distribution. In order to determine the aggregate loss distribution, estimates for  $\lambda$ , the frequency, and  $T$ , the severity distribution is needed.

The 99.9% Value-at-Risk of this aggregate loss distribution is of interest for operational risk capital estimation. However, it is not easy to do this analytically in most cases, and for this reason, Monte Carlo (MC) simulation is often used. The following algorithm can be used to calculate the Value-at-Risk of the compound Poisson distribution (see De Jongh *et al.*, 2015):

- Step 1: Generate  $N$  distributed according to the assumed frequency distribution.
- Step 2: Generate  $Y_1, Y_2, \dots, Y_N$  independent and identically distributed (i.i.d.), according to the true severity distribution  $T$  and calculate  $A = \sum_{n=1}^N Y_n$ .
- Step 3: Repeat Step 1 and Step 2  $I$  times independently to obtain  $A_i, i = 1, 2, \dots, I$  and then approximate the 99.9% VaR by  $A_{([0.999 \cdot I] + 1)}$  where  $A_{(i)}$  denotes the  $i$ -th order statistic and  $[k]$  the largest integer contained in  $k$ .

Where a mixed distribution is used as described in the previous section, Steps 1 and 2 of the above algorithm would look slightly different, but this is revisited in Chapter 4.

The frequency and severity distributions assumed in Step 1 and Step 2 are not known in practice and they need to be estimated as explained in Sections 3.2.1 and 0 above, and would be based on actual loss data.

Other numerical methods can also be used, like the single-loss approximation (SLA) method suggested by Böcker and Klüppelberg (2005). The SLA method is summarised as follows: if  $T$  is the true underlying severity distribution function of the individual losses, and  $\lambda$  is the true annual frequency, then the  $100(1 - \gamma)\%$  VaR of the compound loss distribution may be approximated by  $T^{-1}(1 - \gamma/\lambda)$ .

### **3.3 Capital estimation**

The VaR of the aggregate distribution calculated in Section 3.2.3 above is used to estimate the standalone capital for each ORC. The sum of the standalone capital for all ORCs is known as the undiversified capital.

Correlations or copulas may be used to capture dependencies of operational risk losses across business lines or event types. Such dependencies may result from business cycles, bank-specific factors, or cross-dependence of significant events. Banks employing more granular modelling approaches may incorporate a dependence structure for operational risk losses incurred across those business lines and event types for which separate operational risk models are used. Note that when using correlations to measure dependence, it is generally not true that higher correlations imply a higher risk capital outcome (e.g. for extremely heavy-tailed distributions) (see Embrechts and Hofert, 2011).

A Student-t copula is often used to model the dependence structure between the aggregate loss distributions for the different ORCs, as it includes a degree of tail dependence. In order to simulate the dependence structure between ORCs, the correlation between the ORCs first need to be estimated from historical data. The detail of this process is not expanded on in this dissertation, but the interested reader may refer to McNeil et al., 2015.

### **3.4 Data sources**

As explained under Section 3.2.3, the frequency and severity distributions need to be estimated from actual loss data. We therefore move our focus to the various data sources available to banks to estimate these distributions and their parameters.

It is standard practice in operational risk management to use different data sources for modelling future losses. Banks have typically been collecting their own loss data for some time, referred to as internal data. In addition, various external loss databases exist, including publicly available data, insurance data and consortium data. The Basel Accord (2011) also suggests the use of scenario assessments to improve severity

distribution estimation and business environment and internal control factors. Each of these data sources will be discussed in more detail below.

#### **3.4.1 Internal data**

When constructing statistical models, a bank's own historical losses is probably the most appropriate data to use. However, operational risk is a relatively new risk, and most banks have only been collecting operational losses for the past ten or fifteen years. Once this data is divided into homogeneous groups, to ensure independent and identically distributed losses (for the individual ORCs previously mentioned), the number of data points per homogeneous group, decreases significantly.

Aue and Kalkbrener (2007) explain that although an organisation's internal loss data is the most objective risk indicator available to them, it suffers from two main drawbacks. The first is that it is backwards-looking and therefore does not allow for changes in the organisation's control environment. It is also not available in sufficient quantities to build reliable statistical models, specifically for extreme losses used to inform the capital estimates.

It is for this reason that external data sources may be considered to enhance the banks' data. In addition, history is not always the best predictor of the future, and there is a requirement for capital models to be forward-looking. Scenario assessments by business experts are valuable in giving this forward-looking view. The role of external data and scenario assessments are expanded on below.

#### **3.4.2 External data**

External data is expected to complement a bank's internal data when modelling the loss severity. It includes information on significant actual losses that the individual bank may not have experienced (Basel Committee on Banking Supervision, 2011b). Regulatory supervisors expect banks to use external data to estimate the loss severity as it may contain valuable information to inform the tail of the loss distribution(s). In addition, they argue that it may be a necessary input into scenario analysis and

potentially other uses beyond providing information on large losses for modelling purposes. The potential use of external data in informing the scenario assessments of experts is central to this research study.

In Chapter 5, it is shown how the estimation of an appropriate severity distribution  $F(y, \theta)$  can be done using data from an external database. This is useful where a bank does not have sufficient internal data.

### **3.4.3 Scenario analysis**

Scenario analysis is the third source of loss data to be considered. It is a crucial tool in the identification and management of operational risk. Scenario analysis is described as a process whereby the opinions of risk managers or experts within a specific business line can be obtained to identify possible risk events and establish their potential outcomes (see Basel Committee on Banking Supervision 2011a).

According to the Risk Management Association (2011), the financial industry has reached a consensus view on the importance of scenario analysis and how it can support the risk management process. This is mainly due to the numerous tail events that have resulted in enormous financial loss and reputational damage since the turn of the century. However, scenario analysis is one of the most challenging aspects of the AMA despite the emerging recognition of its value. For example, different practices exist for the use of scenario analysis when measuring risk. As part of the 2008 Loss Data Collection Exercise (LDCE) conducted by the Operational Risk Subgroup of the Standards Implementation Group and published by the Basel Committee on Banking Supervision, banking supervisors collected scenario analysis data on an international basis. They reference three types of scenario approaches: the individual approach, the interval approach, and the percentile approach (Basel Committee on Banking Supervision, 2009). Unfortunately, they do not describe the different scenario approaches or how it is used in practice.

The Risk Management Association (2011) observed that no single accepted conceptual or technical practice has emerged for incorporating scenario analysis into operational

risk measurement. De Jongh *et al.* (2015) investigated various approaches to incorporate scenario assessments into modelling the severity distribution. They proposed a new method for integrating limited historical data (whether internal loss data or external data) with scenario assessments.

In Chapter 4, we expand on the significance of qualitative scenario assessments as a potential data source for quantitative operational risk models and how the method described by De Jongh *et al.* (2015) can incorporate these data points.

#### **3.4.4 Business environment and internal control factors (BEICF's)**

The BCBS (2006) explains that in addition to using loss data, whether actual or scenario-based, a bank's firm-wide risk assessment methodology must capture the critical business environment and internal control factors that can change its operational risk profile. These factors will make a bank's risk assessments more forward-looking and directly reflect the quality of the bank's control and operating environments. It will also help align capital assessments with risk management objectives and recognise improvements and deterioration in operational risk profiles more immediately.

The use of business environment and internal control factors are not expanded on in the remainder of this dissertation.

#### **3.4.5 Combining sources of data**

The Basel Committee on Banking Supervision (2011b) suggests that there may be a need for different combinations of the different data sources, and the onus is on the bank to show that their process of combining the data is sufficient for the purpose it is intended, i.e. to estimate capital.

Various methods have been proposed to combine historical or internal data with the scenario assessments of experts. For example, Dutta and Babbel (2014) suggested using a "Change of Measure" approach to evaluate each scenario's impact on the total estimate for operational risk capital. They conclude that each scenario will change the



historical probability associated with a given severity range. Their paper claims that the successful use of scenario data in capital models has often failed due to the incorrect interpretation or implementation of the data.

Lambrigger *et al.* (2007), on the other hand, combine three sources, namely internal data, external data, and expert opinions, to estimate the parameters of the risk frequency and severity distributions. Their approach uses Bayesian inference as the statistical technique to combine the three sources of data.

De Jongh *et al.* (2015) introduced a simple method whereby the severity distribution of the aggregated losses can be estimated using both historical data and scenario assessments. Their method incorporates a measure of agreement between the two data sources, assessing the quality of both.

Our research shows how both internal and external data can be combined to estimate the severity distribution of losses and how this could be used to inform or challenge the scenario assessments of business experts.

### **3.5 Conclusion**

The proposed methodology for constructing quantitative operational risk models is based on the loss distribution approach. Constructing an aggregate loss distribution from which the risk capital required for each operational risk category within a bank can be estimated, was explained in detail in this chapter.

In the next chapter, this methodology is expanded, and it is shown how scenario analysis is combined with limited historical data to estimate the severity distribution. In Chapter 5, external data sources are considered where a bank does not have sufficient internal data. It is shown how an appropriate severity distribution  $F(y, \theta)$  can be estimated using an external database and utilising certain explanatory variables to scale the severity distribution for an individual bank.

## 4. Constructing forward-looking distributions using limited historical data and scenario assessments

### 4.1 Introduction

In this chapter, statistical methods are explored that could be used to combine limited historical data and scenario assessments to estimate extreme quantiles of a loss distribution. The methodology discussed in this chapter provides a way of constructing forward-looking distributions to ultimately determine risk capital. Extracts from this chapter were published in the book “Linear and non-linear financial econometrics: Theory and Practice” (De Jongh *et al.*, 2021).

### 4.2 Scenario assessments

It was previously mentioned that scenario assessments by experts are sometimes used to augment the limited internal data available by banks. According to Wei *et al.* (2018), the two subjective data elements of scenario analysis and BEICFs are less used than the two objective data elements of internal and external loss data in operational risk measurement. However, scenario analysis is valuable to inform the scale of extreme losses or the tail of the loss distribution.

Three different types of scenario approaches were mentioned in Chapter 3.4.3. However, in this chapter, the focus is on the percentile scenario approach explained in detail by De Jongh *et al.* (2015). They also refer to it as the 1-in- $c$  years' scenario approach. In the 1-in- $c$  years' scenario approach, scenario makers are asked the following question: “What loss level  $q_c$  is expected to be exceeded only once every  $c$  years?” Popular choices for  $c$  range between 5 and 100 years, and often three values for  $c$  are used. Typically, the first choice would be motivated by the number of years of historical data available to a bank, for example, seven or ten years. In this case, the maximum loss experienced by the bank over ten years may serve as a guide for choosing  $q_{10}$ , because the loss level has actually been reached once over the ten years. Scenario makers may, however, want to provide a lower (or higher) assessment of  $q_{10}$  based on whether they believe that the future will be better (or worse) than the past. De Jongh *et al.* (2015) explain that the other choices of  $c$  are selected to obtain a

scenario spread within the range that we can expect a reasonable improvement in accuracy from the experts' inputs, such as 20, 50 or even 100 years. The choice of  $c = 100$  is questionable, given that a judgement on a 1-in-100 years loss level likely fall outside the scope of any expert's experience. For this reason, they may also take additional guidance from external data of similar banks, which in effect amplifies the number of years for which historical data are available. It is argued that this is an essential input into scenario analysis according to BCBS (2011b) and also expanded on in Chapter 5.

De Jongh *et al.* (2015) showed that if the annual loss frequency is  $Poi(\lambda)$  distributed and the true underlying severity distribution is  $T$ , and if the experts are of oracle quality in the sense of actually knowing  $\lambda$  and  $T$ , then the assessments provided should be

$$q_c = T^{-1}\left(1 - \frac{1}{c\lambda}\right). \quad (3)$$

To see this, let  $N_c$  denote the number of loss events experienced in  $c$  years and let  $M_c$  denote the number of these that are greater than  $q_c$ . Then  $N_c \sim Poi(c\lambda)$  and the conditional distribution of  $M_c$  given  $N_c$  is binomial with parameters  $N_c$  and  $1 - p_c = P(Y \geq q_c) = 1 - T(q_c)$  with  $Y \sim T$  and  $p_c = T(q_c) = 1 - \frac{1}{c\lambda}$ . Therefore  $E[M_c] = E[E(M_c|N_c)] = E[N_c(1 - p_c)] = c\lambda(1 - T(q_c))$ . Requiring that  $E[M_c] = 1$ , yields (3).

As an illustration of the complexity of the experts' task, take  $\lambda = 50$  then  $q_{10} = F^{-1}(0.998)$ ,  $q_{20} = F^{-1}(0.999)$  and  $q_{100} = F^{-1}(0.9998)$  which implies that the quantiles that have to be estimated are very extreme.

Considering the equation derived for the scenario assessments in (3), the expert must know both the true severity distribution and the annual frequency when assessment is provided. In order to simplify the task of the expert, De Jongh *et al.* (2015) made use of the mixed model in (1) first introduced in Chapter 3 to show that  $T_u(q_c) = 1 - \frac{b}{c}$ .

Using the equation for  $T_u$  and taking  $q = q_b$ ,  $b < c$ , the new equation for  $T_u(q_c)$  does not depend on the annual frequency  $\lambda$ . They argue that when  $b = 1$ ,  $q_1$  would be the experts' answer to the question 'What loss level is expected to be exceeded once annually?' and a reasonably accurate assessment of  $q_1$  should be possible. As a result,  $T_u(q_c) = 1 - 1/c$  or  $1 - T_u(q_c) = 1/c$ . Keeping in mind the conditional probability meaning of  $T_u$ ,  $q_c$  would be the answer to the question: 'Among those losses that are larger than  $q_1$ , what level is expected to be exceeded only once in  $c$  years?'. Conditioning on losses larger than  $q_1$  has the effect that the annual frequency of all losses drops out of consideration when an answer is sought. The remainder of the chapter assumes that this question is posed to the experts when making their assessments.

### 4.3 Estimating VaR

Suppose historical loss data  $Y_1, Y_2, \dots, Y_n$  for  $a$  years are available, as well as three scenario assessments  $\tilde{q}_{10}$ ,  $\tilde{q}_{20}$  and  $\tilde{q}_{100}$  are provided by experts. Because real scenario makers are not oracles, we denote their assessments by  $\tilde{q}_c$ . The estimation of the 99.9% VaR of the aggregate loss distribution is of interest. In Section 3.2.3 it was explained how the VaR can be approximated using Monte Carlo simulation. Below three approaches are considered to estimate the VaR, namely the naïve approach, the GPD approach and Venter's approach. The naïve approach uses only historical data. The GPD approach (based on the mixed model formulation introduced in Section 0) and Venter's approach use historical data and scenario assessments. As far as estimating VaR is concerned, it is shown that Venter's approach is preferred to the GPD and naïve approaches.

#### 4.3.1 Naïve approach

Assume only historical data are available with a total of  $K$  loss events spread over  $a$  years and denote these observed or historical losses by  $y_1, \dots, y_K$ . Then the annual frequency is estimated by  $\hat{\lambda} = K/a$ . Let  $F(y; \theta)$  denote a suitable family of distributions to model the true loss severity distribution  $T$ . The fitted distribution is denoted by  $F(y; \hat{\theta})$ , with  $\hat{\theta}$  denoting the (maximum likelihood) estimate of the

parameter(s)  $\theta$ . The source of the historical data in this process would be a bank's own operational losses, although in Chapter 5 it is shown how data from external sources can be used and scaled to estimate  $F(y; \hat{\theta})$  where sufficient internal data is not available. In order to estimate VaR, a slight adjustment of the Monte Carlo approximation approach discussed earlier is necessary.

Naïve VaR estimation algorithm:

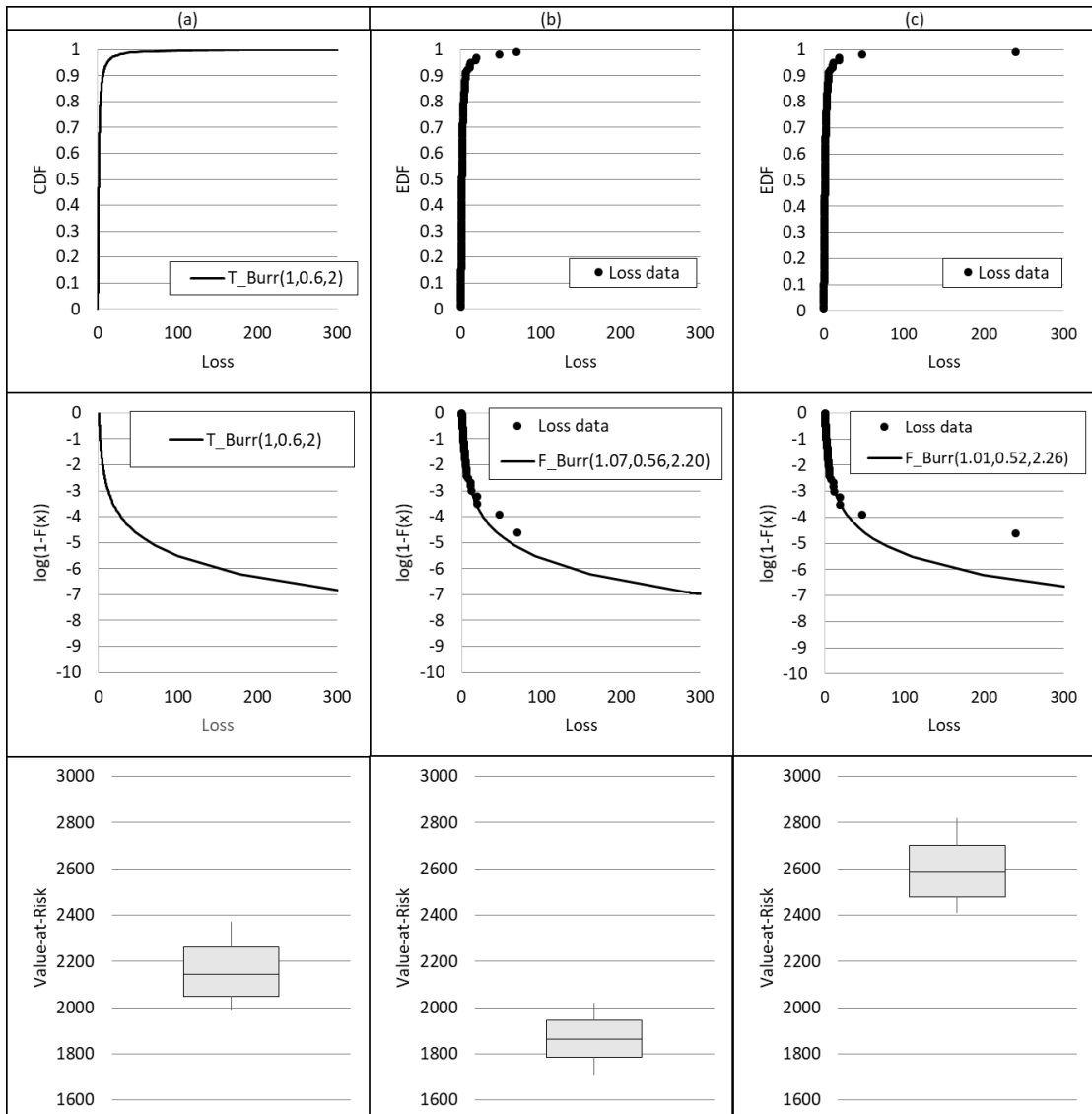
- Step 1: Generate  $N$  from the Poisson distribution with parameter  $\hat{\lambda}$ ;
- Step 2: Generate  $Y_1, \dots, Y_N \sim iid F(y; \hat{\theta})$  and calculate  $A = \sum_{n=1}^N Y_n$ ;
- Step 3: Repeat steps 1 and 2  $I$  times independently to obtain  $A_i, i = 1, 2, \dots, I$ .  
Then the 99.9% VaR is estimated by  $A_{([0.999 * I] + 1)}$  where  $A_{(i)}$  denotes the  $i$ -th order statistic and  $[k]$  the largest integer contained in  $k$ .

Other percentiles can also be used in this algorithm, but as previously mentioned this dissertation focusses on the 99.9<sup>th</sup> percentile that is required for regulatory capital.

The naïve approach should be used cautiously as a single extreme loss can cause drastic changes in estimating the means and variances of severity distributions. This could be true even if a large amount of loss data is available. Annual aggregate losses will typically be driven by the value of the most extreme losses and the high quantiles of the aggregate annual loss distribution are primarily determined by the high quantiles of the severity distributions containing the extreme losses. Two different severity distributions for modelling the individual losses may both fit the data well in terms of goodness-of-fit statistics, yet may provide capital estimates which may differ by billions. Certain deficiencies of the naïve estimation approach, in particular, the estimation of the severity distribution and the subsequent estimation of an extreme VaR of the aggregate loss distribution, are highlighted in Embrechts and Hofert (2011).

In Figure 2 we used the naïve approach to illustrate the effect of some of the above-mentioned claims. In Figure 2(a) we assumed a Burr distribution, i.e.,  $T\_Burr(1, 0.6, 2)$ , as our true underlying severity distribution. In the top panel we show the distribution function and in the middle the log of 1 minus the distribution

function. This gives us more accentuated view of the tail of the distribution. Then in the bottom panel the Monte Carlo results of the VaR approximations are given by means of a box plot using the 5<sup>th</sup> and 95<sup>th</sup> percentiles for the box. One million simulations were used to approximate VaR and the VaR calculations were repeated 1000 times. In Figure 2(b) we assume  $\lambda = 10, a = 10$  and generated  $\lambda a = 100$  observations from the  $T\_Burr(1, 0.6, 2)$  distribution. The observations generated are plotted in the top panel and in the middle panel the fitted distribution and the maximum likelihood estimates of the parameters are depicted as  $F\_Burr(1.07, 0.56, 2.2)$ . In the bottom panel the results of the VaR estimates using the naïve approach are provided. Note how the distribution of the VaR estimates differ from those obtained using the true underlying severity distribution. Of course, sampling error is present, and the generation of another sample will result in a different box plot. To illustrate this, we study the effect of extreme observations by moving the maximum value further into the tail of the distribution, and fitting another distribution. The data set is depicted in the top panel of Figure 2(c) and the fitted distribution in the middle as  $F\_Burr(1.01, 0.52, 2.26)$ . Again, the resulting VaR estimates are shown in the bottom panel. In this case the introduction of the extreme loss has a profound boosting effect on the resulting VaR estimates.



**Figure 2: Illustration of the effects of VaR estimation using the naïve approach**

In practice, and due to imprecise loss definitions, risk managers may incorrectly group two losses into one extreme loss that has a profound boosting effect on VaR estimates. In light of this, the manager must be aware of the process of generating the data and the importance of clear definitions of loss events.

### 4.3.2 The GPD approach

This modelling approach is based on the mixed model formulation in Equation (2) as explained in Section 0. As before,  $a$  years of historical loss data is available,  $y_1, y_2, \dots, y_K$  and scenario assessments  $\tilde{q}_{10}$ ,  $\tilde{q}_{20}$  and  $\tilde{q}_{100}$ . The annual frequency  $\lambda$  can again be estimated as  $\hat{\lambda} = K/a$ . Next,  $b$  and the threshold  $q = q_b$  must be specified. One possibility is to take  $b$  as the smallest of the scenario  $c$ -year multiples and to estimate  $q_b$  as the corresponding smallest of the scenario assessments  $\tilde{q}_b$  provided by the experts, in this case  $\tilde{q}_{10}$ .  $T_e(y)$  can be estimated by fitting a parametric family  $F_e(y, \theta)$  (such as the Burr distribution) to the data  $y_1, y_2, \dots, y_K$  or by calculating the empirical distribution and then conditioning it to the interval  $(0, \tilde{q}_b]$ . Either of these estimates is a reasonable choice, especially if  $K$  is large and the parametric family is well chosen. Whichever estimate is used, denote it by  $\tilde{F}_e(y)$ . For future notational consistency, tildes are used on all estimates of distribution functions which involve the use of the scenario assessments.

Next,  $F_u(y)$  is modelled by the  $GPD(y; \sigma, \xi, q_b)$  distribution. See Table 1 in Chapter 3 for the density function of the GPD distribution. Suppose there are actual scenario assessments  $\tilde{q}_{10}$ ,  $\tilde{q}_{20}$  and  $\tilde{q}_{100}$  and thus take  $b = 10$  and estimate  $q_b$  by  $\tilde{q}_{10}$ . Substituting these scenario assessments into  $F_u(q_c) = 1 - \frac{b}{c}$ ; with  $b = 10$ ,  $c = 20$ , 100 yields two equations

$$\begin{aligned} F_u(\tilde{q}_{20}) &= GPD(\tilde{q}_{20}; \sigma, \xi, \tilde{q}_{10}) = 0.5, \text{ and} \\ F_u(\tilde{q}_{100}) &= GPD(\tilde{q}_{100}; \sigma, \xi, \tilde{q}_{10}) = 0.9, \end{aligned} \tag{4}$$

that can be solved to obtain estimates  $\tilde{\sigma}$  and  $\tilde{\xi}$  of the parameters  $\sigma$  and  $\xi$  in the GPD that is based on the scenario assessments. De Jongh *et al.* (2015) showed that a solution exists only if  $\frac{\tilde{q}_{100} - \tilde{q}_{10}}{\tilde{q}_{20} - \tilde{q}_{10}} > \frac{\ln(100) - \ln(10)}{\ln(20) - \ln(10)} = 3.32$ . This fact should be borne in mind when the experts do their assessments.

With more than three scenario assessments, fitting techniques can be based on (4), which links the quantiles of the GPD to the scenario assessments. An example would



be to minimise  $\sum_c |GPD(\tilde{q}_c; \sigma, \xi, \tilde{q}_b) - (1 - b/c)|$ , although this is not expanded on here. The final estimate of  $F_u(y)$  is denoted by  $\tilde{F}_u(y)$  and all these components are now substituted into Equation 2 to yield the estimate  $\tilde{F}(y)$  of  $F(y)$ , namely

$$\hat{\lambda}\tilde{F}(y) = \left(\hat{\lambda} - \frac{1}{10}\right)\tilde{F}_e(y) + \frac{1}{10}\tilde{F}_u(y). \quad (5)$$

Showing the practical use of equation (5), the algorithm below summarises the integration of the historical data with the 1-in- $c$  years scenarios following the MC approach.

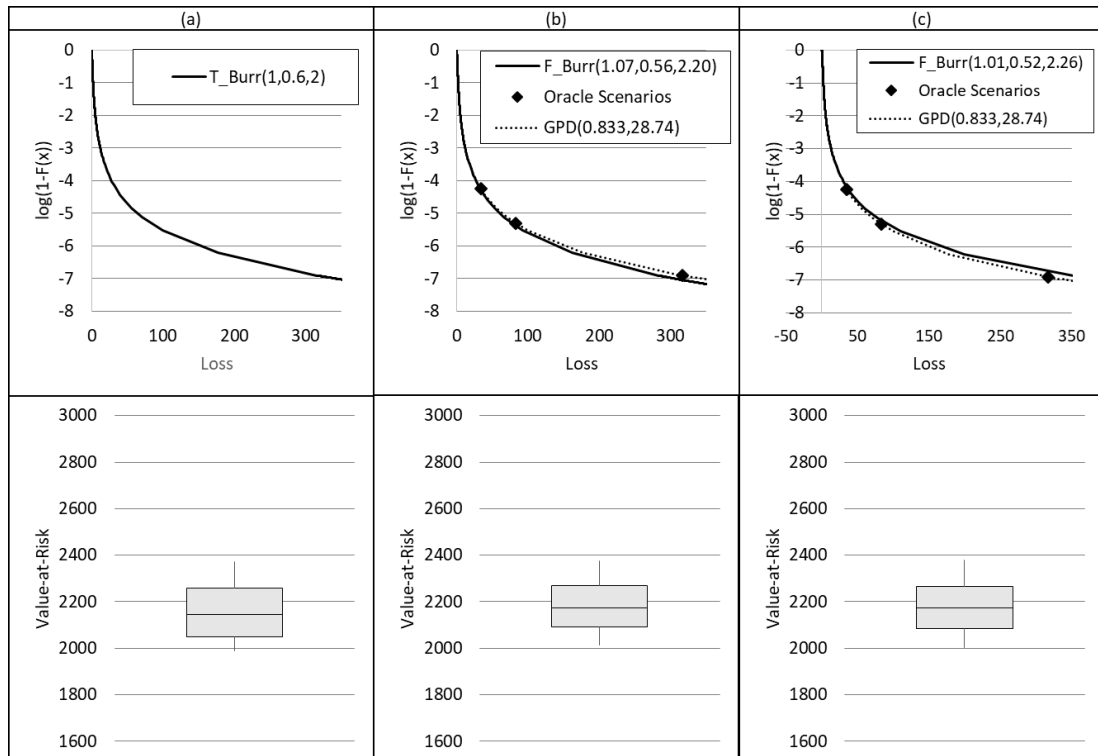
GPD VaR estimation algorithm:

- Step 1: Generate  $N_e \sim Poi\left(\hat{\lambda} - \frac{1}{10}\right)$  and  $N_u \sim Poi\left(\frac{1}{10}\right)$ ;
- Step 2: Generate  $Y_1, \dots, Y_{N_e} \sim iid \tilde{F}_e$  and  $Y_{N_e+1}, \dots, Y_{N_e+N_u} \sim iid \tilde{F}_u$  and calculate  $A = \sum_{n=1}^N Y_n$  where  $N = N_u + N_e$ . It follows that  $A$  is distributed as a random sum of  $N$  i.i.d. losses from  $\tilde{F}$ .
- Step 3: Repeat steps 1 and 2  $I$  times independently to obtain  $A_i, i = 1, 2, \dots, I$  and estimate the 99.9% VaR by the corresponding empirical quantile of these  $A_i$ 's as before.

Using the GPD 1-in- $c$  years integration approach to model the severity distribution, the 99.9% VaR of the aggregate distribution is almost exclusively determined by the scenario assessments and their reliability fundamentally affects the reliability of the VaR estimate. The SLA supports this conclusion. As noted before, the SLA implies that we need to estimate  $q_{1000} = T^{-1}\left(1 - \frac{1}{1000\lambda}\right)$  and its estimate would be  $\hat{q}_{1000} = GPD^{-1}\left(\frac{\left(1 - \frac{1}{1000\hat{\lambda}}\right)}{1 - \left(1 - \frac{1}{10\hat{\lambda}}\right)}, \hat{\sigma}, \hat{\xi}, \hat{q}_b\right)$ .

Therefore, the 99.9% VaR largely depends on the GPD fitted with the scenario assessments. In Figure 3 we depict the VaR estimation results by fitting  $\tilde{F}_e$  assuming a Burr distribution and  $\tilde{F}_u$  assuming a GPD. The top panel in Figure 3 (a) depicts the tail behaviour of the true severity distribution which is assumed as a Burr distribution and

denoted as  $T\_Burr(1, 0.6, 2)$ . Using the VaR approximation technique discussed in Section 3.2.3 and assuming  $\lambda = 10$ ,  $I = 1\,000\,000$  and 1 000 repetitions, the VaR approximations are depicted in the bottom panel in the form of a box plot as before. Assuming that we were supplied with quantile assessments by the oracle we use the two samples discussed in Figure 2 and apply the GDP approach. The results are displayed in Figure 3 (b) and (c) below.



**Figure 3: Illustration of VaR estimates obtained from a GPD fit on the oracle quantiles**

The GPD fitted to the oracle quantiles produce similar box plots, which in turn is very similar to the box plot of the VaR approximations. Clearly the fitted Burr has little effect on the VaR estimates. The VaR estimates obtained through the GPD approach is clearly dominated by the oracle quantiles. Of course, if the assessments are supplied by experts and not oracles the results would differ significantly. This is illustrated when we compare the GPD with Venter's approach.

The challenge is therefore to find a way of integrating the historical data and scenario assessments such that both sets of information are adequately utilised in the process. It would be beneficial to have measures indicating whether the experts' scenario assessments are in line with the observed historical data and, if not, to require them to produce reasons why their assessments are so different. In Chapter 5 we suggest a way of using external data to challenge experts' scenario assessments.

### 4.3.3 Venter's approach

De Jongh *et al.* (2015) explain that their proposed approach was first suggested by their colleague, Hennie Venter. Given the quantiles  $q_{10}$ ,  $q_{20}$  and  $q_{100}$ , the distribution function  $T$  can be written as follows:

$$T(y) = \begin{cases} \frac{p_{10}}{T(q_{10})} T(y) & \text{for } y \leq q_{10} \\ p_{10} + \frac{p_{20} - p_{10}}{T(q_{20}) - T(q_{10})} [T(y) - T(q_{10})] & \text{for } q_{10} < y \leq q_{20} \\ p_{20} + \frac{p_{100} - p_{20}}{T(q_{100}) - T(q_{20})} [T(y) - T(q_{20})] & \text{for } q_{20} < y \leq q_{100} \\ p_{100} + \frac{1 - p_{20}}{1 - T(q_{100})} [T(y) - T(q_{100})] & \text{for } q_{100} < y < \infty. \end{cases} \quad (6)$$

Again  $T(q_c) = p_c = 1 - \frac{1}{c\lambda}$  and therefore, the expressions on the right reduces to  $T(y)$ . The definition of  $T(y)$  could also be extended for more quantiles. They model  $T(y)$  by  $F(y, \theta)$  and estimate it by  $F(y, \hat{\theta})$  using historical data and maximum likelihood estimates and also estimate the annual frequency by  $\hat{\lambda} = K/a$ . Given scenario assessments  $\tilde{q}_{10}$ ,  $\tilde{q}_{20}$  and  $\tilde{q}_{100}$ ,  $T(q_c)$  can be estimated by  $F(\tilde{q}_c, \hat{\theta})$  and  $p_c$  by  $\hat{p}_c = 1 - \frac{1}{c\hat{\lambda}}$ . The estimated ratios are then defined by

$$\begin{aligned} R(10) &= \frac{\hat{p}_{10}}{F(\tilde{q}_{10}; \hat{\theta})}, \\ R(10,20) &= \frac{\hat{p}_{20} - \hat{p}_{10}}{F(\tilde{q}_{20}; \hat{\theta}) - F(\tilde{q}_{10}; \hat{\theta})}, \\ R(20,100) &= \frac{\hat{p}_{100} - \hat{p}_{20}}{F(\tilde{q}_{100}; \hat{\theta}) - F(\tilde{q}_{20}; \hat{\theta})}, \text{ and} \\ R(100) &= \frac{1 - \hat{p}_{100}}{1 - F(\tilde{q}_{100}; \hat{\theta})}. \end{aligned} \quad (7)$$

Their new method is to estimate the true severity distribution function  $T$  by an adjusted form of  $F(x, \hat{\theta})$ , then Hennie's distribution  $\tilde{H}$  is defined as follows (see De Jongh *et al.* 2015):

$$\tilde{H}(y) = \begin{cases} R(10)F(y; \hat{\theta}) & \text{for } y \leq \tilde{q}_{10} \\ \hat{p}_{10} + R(10,20)[F(y; \hat{\theta}) - F(\tilde{q}_{10}; \hat{\theta})] & \text{for } \tilde{q}_{10} < y \leq \tilde{q}_{20} \\ \hat{p}_{20} + R(20,100)[F(y; \hat{\theta}) - F(\tilde{q}_{20}; \hat{\theta})] & \text{for } \tilde{q}_{20} < y \leq \tilde{q}_{100} \\ \hat{p}_{100} + R(100)[F(y; \hat{\theta}) - F(\tilde{q}_{100}; \hat{\theta})] & \text{for } \tilde{q}_{100} < y < \infty. \end{cases} \quad (8)$$

If all estimators are precisely equal to what they are estimating, this estimate is consistent because it reduces to  $T$ . Their new severity distribution estimate  $\tilde{H}$  'believes' the scenario quantile information but follows the distribution fitted on the historical data to the left of, within, and right of the scenario intervals. The ratios  $R(10)$ ,  $R(10,20)$ ,  $R(20,100)$  and  $R(100)$  in (7) can be viewed as measures of agreement between the historical data and the scenario assessments and could help assess their validities. The steps required to estimate VaR using the Venter method are set out below.

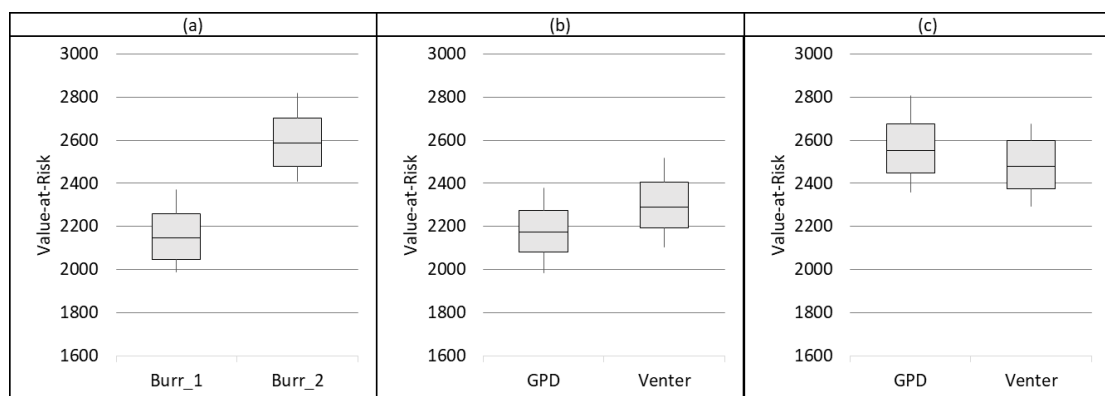
Venter method VaR estimation algorithm:

- Step 1: Generate  $N \sim Poi(\hat{\lambda})$ ;
- Step 2: Generate  $Y_1, \dots, Y_N \sim iid \tilde{H}$  and calculate  $A = \sum_{n=1}^N Y_n$ ;
- Step 3: Repeat steps 1 and 2  $I$  times independently to obtain  $A_i, i = 1, 2, \dots, I$  and estimate the 99.9% VaR by the corresponding empirical quantile of these  $A_i$ 's as before.

The single loss approximation introduced earlier is considered and how it applies to this method. The SLA implies that they need to estimate  $q_{1000} = T^{-1}\left(1 - \frac{1}{1000\lambda}\right)$  and its estimate would be  $\hat{q}_{1000} = \tilde{H}^{-1}\left(1 - \frac{1}{1000\lambda}\right) = \tilde{H}^{-1}(\hat{p}_{1000})$ . The equation  $F(\hat{q}_{1000}; \hat{\theta}) = F(\tilde{q}_{100}; \hat{\theta}) + (\hat{p}_{1000} - \hat{p}_{100})/R(100)$  needs to be solved for  $\hat{q}_{1000}$ . It is relatively easy when using the Burr family as there is an explicit expression for its

quantile function. It shows that a combination of the historical data and scenario assessments are involved, and not exclusively the latter. In as much as the SLA provides an approximation to the actual VaR of the aggregate loss distribution, it is expected that the same will hold for Venter’s approach.

In order to illustrate the properties of this approach we assume that the true underlying severity distribution is the  $Burr(1, 0.6, 2)$  as before. We then construct a ‘false’ severity distribution as the fitted distribution to the distorted sample depicted in Figure 2 (c), i.e. the  $Burr(1.01, 0.52, 2.26)$ . We refer to the true severity distribution as  $Burr_1$  and the false one as  $Burr_2$ . In Figure 4 (a) the box plots of the VaR approximations of the two distributions are given (using the same input for the MC simulations in Figure 2 (a) and Figure 2 (c) respectively). We then illustrate the performance of the GPD and Venter approach in two cases. The first case assumes that the correct (oracle) quantiles of  $Burr_1$  are supplied, but that the loss data are distributed according to the false distribution  $Burr_2$ . In the second case, the quantiles of the false severity distribution are supplied, but the loss data follows the true severity distribution. The box plots of the VaR estimates are given in Figure 4 (b) for case 1 and Figure 4 (c) for case 2.



**Figure 4: Comparison of VaR results for the GPD and Venter approaches**

The Venter approach is affected by both the loss data and quantiles supplied. In the example studied here it seems as if the method is more affected by the quantiles than by the data. The role of the data relative to the quantiles changes positively with the volume of loss data being supplied.

#### **4.4 Conclusion**

Historical data was used to estimate frequency and severity distributions and it was shown how these can be integrated with scenario assessments. Where a bank has sufficient data of its own historical losses, this would be the most important source to estimate these frequency and severity distributions. However, it was previously mentioned that banks often have limited data in a specific ORC, and it may thus be difficult to estimate the severity distribution. For this reason, banks may also make use of external data sources for one or more ORC where their own data is scarce.

This chapter considered three approaches to estimate the VaR, namely the naïve approach, the GPD approach, and Venter's approach. The Venter approach estimates the severity distribution using historical data and experts' scenario assessments jointly. The way in which historical data and scenario assessments are integrated incorporates measures of agreement between these data sources, which can be used to evaluate the quality of both. Major international banks have already implemented this method with great success.

In the next chapter, SAS® OpRisk Global Data is used to demonstrate how external data can be scaled for use in the modelling process.

## **5. Using external data sources to inform scenario assessments**

### **5.1 Introduction**

In the previous chapter we explored statistical methods to combine historical data and scenario assessments to estimate the extreme quantiles of the aggregate loss distribution. Although it was not specified whether this historical data had to be internal or external data, it was noted in Section 3.4.1 that a bank's own historical losses would be most appropriate for this modelling exercise. However, most banks have only been collecting operational loss data for a relatively short period of time, and once the data is divided into homogeneous groups, the number of data points in each ORC decreases significantly. This could make the process of fitting a parametric distribution difficult.

To address this shortfall, banks often subscribe to external data consortiums or use other external data sources to supplement their own data. External databases are extremely valuable because they pool together the industry's experiences and provide an estimate of exposure for an average bank in the industry.

First, external data can be used for ORC's where no or limited data is available, to estimate both frequency and severity distributions. Note that the objective is not to combine internal and external data where sufficient internal data is available, but rather to obtain estimates for frequency and severity distributions where no or too few internal data points are available. The proposed process of scaling external data to be appropriate for use to an individual bank will be discussed in this chapter.

Second, for ORC's where sufficient internal data is available, the use of external data may still be beneficial. For example, the quantiles of the aggregate distribution fitted to external data can potentially be used to inform or challenge scenario assessments.

In this chapter, we consider the use of external data in operational risk models and how it could be utilised in the ultimate decision making of a bank. Specifically, we suggest a scaling methodology to estimate the severity distribution where no or

limited data is available, but also how external data can inform or challenge subjective scenario assessments.

First, we explain the difference between the two types of external databases, namely a public database and a consortium.

## 5.2 External databases

Two types of external databases exist in practice (Baud *et al.*, 2002):

- The first type of database records publicly reported losses. Wei *et al.* (2018) explain that commercial vendors and researchers construct these public databases by collecting information on operational loss events from public data sources such as newspapers, websites, and others. SAS® OpRisk Global Data, OpBase, and the Willis Towers Watson's (WTW) database are examples of these public databases constructed by commercial vendors. The losses included in these databases are typically considered far too important to be concealed from the public eye and are therefore openly reported. In Section 5.3 the public database used in this study, namely SAS® OpRisk Global Data, is discussed in more detail.
- The second type of external database is based on a consortium of organisations. It works on an agreement among a group of member organisations who commit to supply the database with their own internal losses, provided that some confidentiality principles are respected. In return, the member organisations are allowed to use the data to supplement their internal data. An example of a consortium database is Operational Riskdata eXchange Association (ORX). A group of banks initially set up ORX to provide a global platform for the secure and anonymous exchange of data. They now consist of more than 80 member banks, including some of the largest financial institutions around the globe.

Baud *et al.* (2002) explain that the two types of databases will differ by how losses are truncated, i.e., only losses above a certain threshold will be captured, and that threshold amount is not necessarily the same. With publicly reported losses, the



truncation point is expected to be higher than a consortium's database, given that it is only losses in the public's interest that are openly reported and subsequently recorded. The consortium can specify to their member organisations the threshold at which a loss should be included, and for example, ORX includes all losses above EUR 20 million in their global database (Wei *et al.*, 2018). Making specific allowance for these thresholds as part of the modelling process is discussed in more detail later in this chapter.

### **5.3 SAS® OpRisk Global Data**

SAS® OpRisk Global Data was used for this research study. The SAS® OpRisk Global Data is a comprehensive and accurate repository of publicly reported operational losses above USD 100 000, containing more than 32 000 events across all industries worldwide (SAS, 2021). For each publicly available operational loss, the SAS dataset provides the loss amount and additional information about the company and industry where the loss occurred. This includes, among others, a description of the loss event, the region, the size of assets of the organisation and other information associated with the loss.

The SAS® OpRisk Global Data is updated monthly, and the database published in August 2020 was used for this research study. Although there were almost 37 000 losses captured at that date, only losses in the Financial Services industry and exceeding USD 1 million were included. Because the focus is on risk modelling within a bank or financial institution, losses in non-financial lines of business would not be considered representative. Losses below USD 1 million were omitted, as for capital estimation, the focus is on modelling the tail of the distribution.

Table 2 shows a breakdown of the number of losses for each of the nine business lines. Almost a third of the losses (3 630 of 11 190) occurred in Retail Banking, although the median log-loss of USD1.6 million in this business line is relatively low compared to the other business lines.

It should be noted that the data includes a business line for “Insurance”, which was not part of the business lines specified by the BCBS and set out in Section 2.4. Given that banking groups often also provide insurance business, these datapoints were not removed. A more detailed discussion of the distribution of the data is provided in Section 5.4.

**Table 2: Breakdown of losses per business line**

Business line	No of losses	% of losses	Log-losses (USD Million)		
			50 <sup>th</sup> percentile	90 <sup>th</sup> percentile	99 <sup>th</sup> percentile
Agency Services	174	1.6%	3.18	5.62	7.78
Asset Management	513	4.6%	2.61	4.99	7.52
Commercial Banking	2 093	18.7%	2.26	4.93	7.25
Corporate Finance	581	5.2%	2.92	5.70	8.26
Insurance	1 899	17.0%	2.20	4.81	7.02
Payment & Settlement	227	2.0%	2.40	5.84	7.68
Retail Banking	3 630	32.4%	1.60	4.61	7.68
Retail Brokerage	808	7.2%	1.49	3.77	6.22
Trading & Sales	1 265	11.3%	3.17	6.09	8.41
	<b>11 190</b>	<b>100.0%</b>			

#### 5.4 Preliminary data analysis and determination of explanatory variables

Some authors have previously suggested that the severity of operational losses experienced by organisations may be related to specific exposure indicators, for example, the size of the organisation or the region in which the organisation operates. This suggests that if the same loss event had to occur at two differently sized organisations, the larger organisation is expected to experience a larger loss, all else being equal. Further, one may expect that the business, legal and regulatory environments in which organisations operate will differ from one region to another and impact the severity of operational losses. We include a summary from the

literature on those variables that could potentially impact operational losses, for ease of reference.

Shih *et al.* (2000) showed that three variables associated with the size of an organisation, namely revenue, assets, and number of employees, were correlated with the size of operational losses. They found that revenue had the most substantial relationship with loss size and that the logarithm of the scale variables showed a stronger relationship than the raw variables. Cope and Labbi (2008) investigated whether an organisation or a business line's size is correlated with operational risk loss and if the organisation's geographical region impacts the loss size. They used quantile regression techniques to characterise differences in loss distributions for banks of different sizes and conducting business in different geographies. Their study was based on ORX data, and an interesting finding of their study was that there are certain business lines and event types where the loss severity decreases when the size of the bank increases. A possible explanation is that bigger banks may have invested relatively more in their internal controls. Dahlen and Dionne (2010) constructed models to scale both the severity and frequency of external losses for integration with internal data. Their ordinary least squares estimation results show that size, business line, and risk type variables can be used to explain external loss amounts. In their paper, they state that total revenue, total assets, or the number of employees can be used as a proxy to estimate the size of the bank where the loss occurred. However, they chose total assets for their study, as it was the variable most correlated with losses in the Algo OpData they used.

For our investigation, the relevant literature discussed above and the available data were considered. For each reported loss, SAS® OpRisk Global Data provides five possible data fields that could indicate the organisation's size where the loss occurred. These include revenue, assets, net income, number of employees and shareholder equity. It is reasonable to assume that there is a positive correlation between these variables. Therefore, only a single variable was selected to represent the organisation's size, namely the logarithm of the organisation's assets.

Ganegoda and Evans (2012) also suggest that the equity ratio, being the proportion of equity used to finance the company's assets, can indicate the risk-taking tendency of management. It provides a measure of leverage used and given that both the assets and shareholder equity are provided in the SAS data, this ratio could easily be computed. It was used as the second explanatory variable in the scaling model.

The third explanatory variable included in our model was the geographic region in which the organisation operates: Africa, Asia, Europe, North America, Other Americas, or Other. Wilson (2007) explains that all operational losses arise due to a specific set of circumstances and a lack of, or failure in, controls. The reason for including region as an explanatory variable is based on the assumption that circumstances should be similar in a specific geographic region and therefore impact on operational losses in different regions.

Table 3 provides summary statistics about the losses in the different regions.

**Table 3: Summary statistics per geographical region**

Region	No of losses	% of losses	Log-losses (USD Million)		
			50 <sup>th</sup> percentile	90 <sup>th</sup> percentile	99 <sup>th</sup> percentile
Africa	124	1.1%	1.47	5.05	6.14
Asia	1 514	13.5%	2.12	5.02	7.16
Europe	2 820	25.2%	2.44	5.72	7.91
North America	6 292	56.2%	1.99	4.77	7.32
Other	317	2.9%	1.84	4.61	6.86
Other Americas	123	1.1%	2.86	5.60	7.46
	<b>11 190</b>	<b>100.0%</b>			

Note: Given the relatively small number of losses reported in *Africa*, *Other* and *Other Americas*, we have grouped these losses.

The Basel Committee on Banking Supervision (2005) specifies that a bank's activities should be categorised into business lines. Some business lines are considered riskier than others and may potentially suffer higher losses. Hence, the severity distribution will be impacted by the business line, being the fourth explanatory variable. The observed log-losses per business line were already provided in Table 2. Given the relatively small number of losses reported under *Agency Services*, *Asset Management* and *Payment and Settlement*, the losses in these categories were grouped. This category is referred to as *AS*, *AM* and *PS* later in the chapter.

A comprehensive set of non-overlapping operational event types should be defined and applied across the various business lines. The final explanatory variable was event type, as it is found that different types of loss events are associated with different sized losses. A list of the event types used as part of the analysis is provided in Table 4.

**Table 4: Summary statistics per event type**

Event type	No of losses	% of losses	Log-losses (USD Million)		
			50 <sup>th</sup> percentile	90 <sup>th</sup> percentile	99 <sup>th</sup> percentile
Business Disruption and System Failures	62	0.6%	3.01	5.20	6.25
Clients, Products & Business Practices	5 747	51.4%	2.63	5.55	7.92
Damage to Physical Assets	102	0.9%	3.14	6.91	9.04
Employment Practices and Workplace Safety	359	3.2%	1.86	4.24	6.56
Execution, Delivery & Process Management	508	4.5%	1.47	4.10	6.74
External Fraud	2 297	20.5%	1.49	3.96	6.39
Internal Fraud	2 115	18.9%	1.62	4.59	7.09
	<b>11 190</b>	<b>100.0%</b>			

Note: Given the relatively small number of losses reported under *Business Disruption and System Failures*, *Damage to Physical Assets* and *Employment Practices and Workplace Safety*, the losses in these categories were grouped. This category is referred to as *SF, D and EP* later in the chapter.

To summarise, the explanatory variables that would be investigated as part of the model application elaborated on in Section 5.7 are the company's assets (using the logarithm of this value), the equity ratio, geographical region, business line, and event type. Wilson (2007) outlines three types of biases inherent in external data: reporting bias, control bias, and scale bias. Allowing for reporting and scale bias are discussed in Sections 5.5 and 5.6 respectively.

## 5.5 Allowing for reporting bias

Before explaining the model to be used, it is essential first to consider potential biases in the data. As previously explained, the SAS® OpRisk Global Data contains information obtained from several online information providers and other publications. A team of seasoned SAS operational risk research analysts maintain the database according to strict data quality standards and review it periodically to update it and ensure accuracy and completeness. Ganegoda and Evans (2012) argue that most external databases, especially those maintained by vendors collecting publicly reported losses, suffer from reporting bias. Wilson (2007) explains that larger losses (and those associated with larger firms) are more likely reported in the media due to factors such as size and nature of loss. This is because not all operational losses are reported on public platforms, especially smaller losses. As a result, public databases may contain a disproportionately high number of large or significant losses. One should make allowance for this bias when fitting a statistical model, or else the tail of the distribution will be overestimated.

Ganegoda and Evans (2012) draws on a method first introduced by De Fontnouvelle *et al.* (2006) to assign a weight to each loss in the external database. This means that smaller losses will carry greater weight, and greater losses will carry a smaller weight. We briefly explain their methodology below.

They firstly assume that a (log) loss  $y_i$  is only reported in the public domain if it exceeds a certain truncation or observation point,  $t_i$ . This truncation point  $t_i$  is a stochastic variable and should not be confused with the threshold at which losses are captured in the database, being USD100 000 in the case of the SAS database. To explain this, if

a loss is greater than the unobserved truncation point but lower than the USD100 000 threshold, the research analyst responsible for compiling the database will observe the loss but not include it in the database. On the other hand, if the unobserved truncation point is higher than the USD100 000 threshold and the observed loss, the analyst will not take note of the loss, and it will, for this reason, also not be included in the database. Therefore, only losses greater than both  $t_i$  and USD100 000 will be included in the database.

Because they assume that the loss amount,  $y_i$  and truncation point,  $t_i$  are independent, the distribution of losses in the database is given by:

$$f(y_i|y_i > t_i) = \frac{f(y_i)G(y_i)}{\int_{\mathbb{R}} f(y)G(y) dy},$$

where  $f(y_i)$  is the marginal densities of  $y_i$  and  $G(\cdot)$  is the cumulative distribution function of the random truncation point  $t_i$ .

They recommend using a logistic distribution for  $G(\cdot)$ , which is given by

$$G(t_i; \tau, a) = \frac{1}{1 + \exp\left[\frac{-(t_i - \tau)}{a}\right]},$$

where  $\tau$  is the location parameter which indicates the log loss, with a 50% probability of being reported in the database, and  $a$  is the scale parameter which dictates the rate at which the probability of being reported increases with the size of the loss.

If  $z_i = y_i - u$  is defined as the excess log loss over a high enough threshold  $u$ , it is shown that  $z_i$  can be approximated using an exponential distribution. They obtain the following likelihood equation

$$L(b, \tau, a) = \prod_{i=1}^n \frac{h(z_i; b)G(z_i; \tau, a)}{\int_{\mathbb{R}} h(z; b)G(z; \tau, a) dz},$$



where  $h(z_i; b) = \frac{1}{b} \exp(-\frac{z_i}{b})$  and  $\tau^* = (\tau - u)$ . The parameters  $b$ ,  $\tau^*$  and  $a$  are estimated by maximising the likelihood function, and finally, the normalised weights to be assigned to each loss is calculated as

$$w'_i = \frac{nw_i}{\sum_{i=1}^n w'_i}$$

where:

$$w_i = \frac{1}{G(y_i|\tau, a)}$$

In order to confirm the existence of reporting bias in the SAS dataset, likelihood ratio tests were carried out for the restriction that the reporting probabilities are constant across all losses for each threshold level (i.e., there is no reporting bias in the data). The  $p$ -values of the likelihood ratio tests for all the threshold values were less than 0.01, confirming the existence of reporting bias.

The parameters  $b$ ,  $\tau^*$  and  $a$  for different choices of the threshold  $u$  were estimated, and it was found that the parameter values for  $b$  and  $a$  stabilised after the USD9 million threshold. Therefore, the associated estimates  $\hat{a} = 1.11809$  and  $\hat{t} = (3.50179 + \log(9))$  was used to calculate corresponding weights for all the losses reported in the SAS database.

Providing a practical example of the weights applied to the SAS data, a current log loss of USD0.78 million would be allocated a weight of 1.9654, whereas a log loss of USD7.78 million would only carry a weight of 0.0276. This illustrates the earlier explanation that larger losses are more likely to be included in databases with publicly reported losses and should therefore carry a smaller weight not to skew the modelling results.

## 5.6 Allowing for scaling bias

Scale bias occurs when data is collected from institutions of a different size. The scaling methodology applied below to the external SAS data to model the severity distribution of operational losses will correct the scale bias.

The method used in this chapter was first introduced by Ganegoda and Evans (2012) and uses regression analysis based on the generalised additive models for location scale and shape (GAMLSS) framework to model the scaling properties of operational losses. They explain that the GAMLSS framework can model all the distributional parameters and therefore offers flexibility in estimating the scaling properties of a model.

In their paper, Ganegoda and Evans (2012) argue that a good scaling model should allow for variations in model parameters for different business lines and event types. The discussion below provides the technical background to their approach.

Consider log losses denoted by  $\mathbf{y} = (y_1, \dots, y_n)^T$ , a random sample of independent observations. Assume that these log losses follow some parametric distribution  $f(\mathbf{y}; \theta)$  with parameter vector  $\theta$ . For the sake of simplicity and in line with Ganegoda and Evans' (2012) notation, assume that  $\theta = (\mu, \sigma)^T$  is a vector of only two distributional parameters.

A set of link functions are defined that specifies the relationship between the linear predictor and the distributional parameters of each distribution component distribution as:

$$\begin{aligned}\log \mu &= \beta_{11} + \beta_{12}X_{i12} + \dots + \beta_{1p}X_{i1p}, \\ \log \sigma &= \beta_{21} + \beta_{22}X_{i22} + \dots + \beta_{2p}X_{i2p},\end{aligned}\tag{9}$$

for  $i = 1, \dots, n$ , where  $X_{ijp}$  is the value of the  $p$ th explanatory variable relating to the observation  $y_i$  in the  $j$ th distributional parameter, and  $\beta_{jp}$  is the parameter corresponding to  $X_{ijp}$ .

The set of equations can also be written in matrix notation as follows:

$$\log \mu = \mathbf{X}_1 \boldsymbol{\beta}_1,$$

$$\log \sigma = \mathbf{X}_2 \boldsymbol{\beta}_2,$$

where,  $\mathbf{X}_j$  are the matrix of the  $j$ th distributional parameter, and  $\boldsymbol{\beta}_j$  are the corresponding parameter vectors.

Ganegoda and Evans (2012) suggested using the log link function, which is also found to be appropriate given the choice of distribution as discussed as part of the model application in Section 5.7.

The maximum likelihood estimates of  $\boldsymbol{\beta}_1$  and  $\boldsymbol{\beta}_2$  are then obtained by solving:

$$\max_{\boldsymbol{\beta}_1, \boldsymbol{\beta}_2} \sum_{i=1}^n w_i' \log f(y_i; \boldsymbol{\beta}_1, \boldsymbol{\beta}_2).$$

To solve the above equation, the PROC NLP function in SAS Enterprise Guide was used.

## 5.7 Model application

The first step in the model selection process was to find a base model that closely follows the data without considering any of the explanatory variables set out above. In other words, an appropriate probability distribution assumption was selected to be used in the subsequent model fitting. For this purpose, the SEVERITY procedure in SAS was used. Six different parametric models were considered to model the severity of log losses, namely the Burr, Gamma, Generalised Pareto, Inverse Gaussian (Wald), Lognormal and Pareto. The density and cumulative distribution functions for all these distributions were provided in Section 0.

In order to select the best base model, three goodness of fit tests are considered. These are twice the negative log-likelihood (-2LogLikelihood), the Akaike's Information

Criterion (AIC) and the Bayesian Information Criterion (BIC). The AIC and BIC are based on the  $-2\text{LogLikelihood}$  and smaller values of all these criteria indicate a better model. Both the AIC and BIC penalise models with more parameters, but the BIC even more so, and for this reason, the BIC was selected as the main determining factor in selecting the best-fit model.

The BIC is defined as:

$$BIC = -2\text{LogLikelihood} + k \cdot \ln(n),$$

where  $k$  is the number of estimated parameters in the model and  $n$  is the number of observations used in the model.

The results of the model fitting process showed that the Burr distribution had the lowest BIC value, but the fact that it has three parameters introduced potential complications for the scaling model. For this reason, a decision was taken to use the Gamma distribution that ranked second among the potential models and this was also the severity distribution used by Ganegoda and Evans (2012).

The probability density function of the Gamma distribution is given by:

$$f(y) = \frac{1}{\Gamma(\sigma)\mu^\sigma} y^{\sigma-1} e^{-\frac{y}{\mu}},$$

where  $\sigma > 0$  is the shape parameter and  $\mu > 0$  is the scale parameter.

Since the Gamma distribution was selected as the appropriate distribution function fitted to all the data, it was assumed that the Gamma distribution would also be appropriate as the base to continue the modelling process. The parameters were estimated to be  $\hat{\sigma} = 1.0993$  and  $\hat{\mu} = 0.8857$  without accounting for the explanatory variables.

The Gamma distribution was fitted to the data again, this time allowing for the explanatory variables introduced in the previous section. The log-assets and equity ratio for each loss was calculated using the SAS data provided. For the categorical explanatory variables, namely region, business line and event type, the dummy variables summarised in Table 5, 6 and 7 were used as part of the coding.

**Table 5: Region dummy coding**

<b>Region_domicile</b>	<b>d_RE1</b>	<b>d_RE2</b>	<b>d_RE3</b>
North America	0	0	0
Africa, Other Americas, Other	1	0	0
Asia	0	1	0
Europe	0	0	1

**Table 6: Business line dummy coding**

<b>Business line 1</b>	<b>d_BL1</b>	<b>d_BL2</b>	<b>d_BL3</b>	<b>d_BL4</b>	<b>d_BL5</b>	<b>d_BL6</b>
Trading and sales	0	0	0	0	0	0
Agency Services, Asset Management, and Payment and Settlement	1	0	0	0	0	0
Commercial Banking	0	1	0	0	0	0
Corporate Finance	0	0	1	0	0	0
Insurance	0	0	0	1	0	0
Retail Banking	0	0	0	0	1	0
Retail Brokerage	0	0	0	0	0	1

**Table 7: Event type dummy coding**

Event type	d_ET1	d_ET2	d_ET3	d_ET4
Business Disruption and System Failures, Damage to Physical Assets, and Employment Practices and Workplace Safety	0	0	0	0
Clients, Products & Business Practices	1	0	0	0
Execution, Delivery & Process Management	0	1	0	0
External Fraud	0	0	1	0
Internal Fraud	0	0	0	1

We carried out a stepwise selection of these variables to determine the parameters  $\mu$  and  $\sigma$  of the Gamma distribution using the link functions introduced in Equation 9. The first step involved only determining an intercept for both  $\mu$  and  $\sigma$ . The next step involved a forward selection of variables only for  $\mu$ , testing the addition of each variable using the model fit criterion that the  $p$ -value should be lower than 0.05. This process was followed by the forward selection of variables for  $\sigma$  given the model already obtained for  $\mu$ . Afterwards, a backward elimination of variables for  $\mu$ , given the selected models for both  $\mu$  and  $\sigma$  was performed and finally a backward elimination of variables for  $\sigma$ .

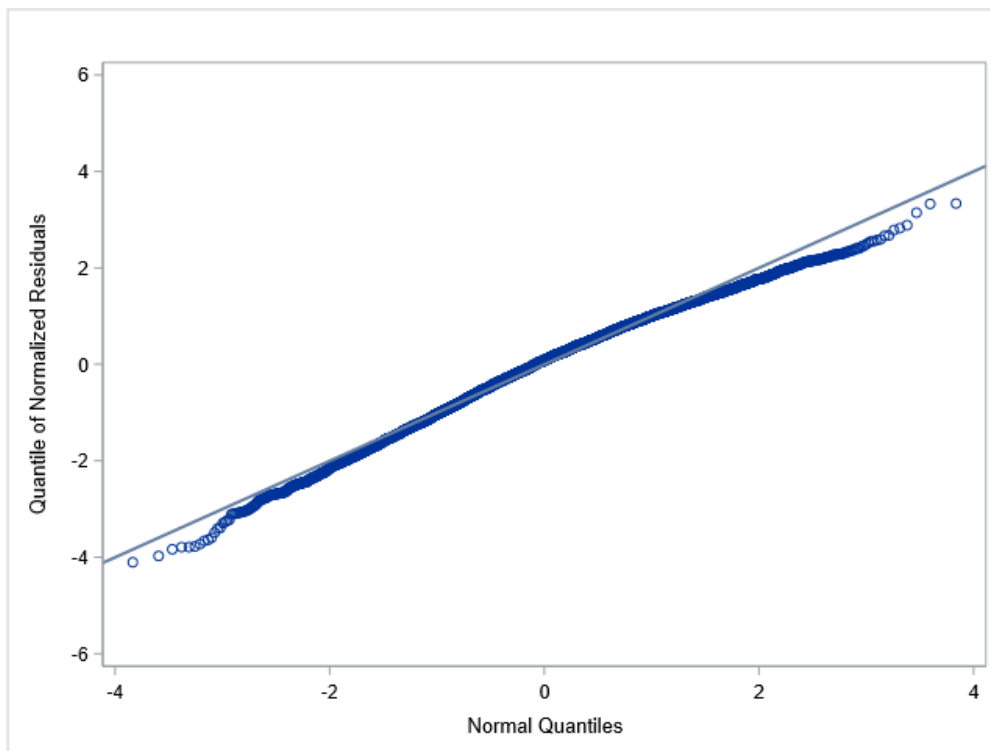
Based on the stepwise selection method described above, it was found that log-assets and seven other business line and event type explanatory variables were significant to the scale parameter  $\mu$ . None of the region variables was found to be significant for either  $\mu$  or  $\sigma$ . For the shape parameter  $\sigma$ , log-assets, the business line variables, Commercial Banking and Retail Brokerage were significant explanatory variables, as well as the event type variable Execution, Delivery & Process Management. The parameter estimates of the final model given by the stepwise selection method are shown in Table 8.

**Table 8: Estimated parameter values for the final model**

Explanatory variable	$\mu$		$\sigma$	
	Estimate	Std. Error	Estimate	Std. Error
Intercept	-0.376027	0.082815	0.335437	0.080145
Log-assets	0.024578	0.006180	-0.016364	0.006857
Equity ratio	-	-	-	-
Africa, Other Americas, Other	-	-	-	-
Asia	-	-	-	-
Europe	-	-	-	-
Corporate finance	0.197507	0.048580	-	-
AS, AM & PS	0.126969	0.039538	-	-
Commercial Banking	0.134076	0.051698	0.111317	0.055369
Insurance	-	-	-	-
Retail Banking	-0.097547	0.023836	-	-
Retail Brokerage	-	-	-0.214888	0.044202
Clients, Products & Business Practices	0.109277	0.038579	-	-
Execution, Delivery & Process Management	-	-	-0.217156	0.057255
External Fraud	-0.298827	0.041982	-	-
Internal Fraud	-0.189108	0.041876	-	-

## 5.8 Model diagnostics and results

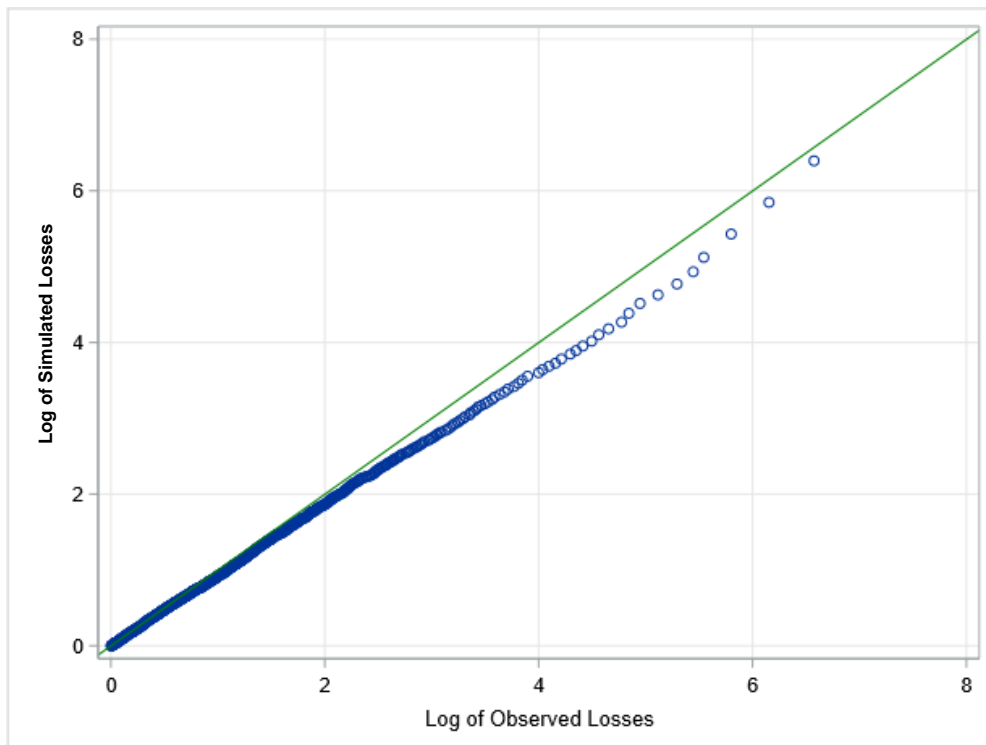
Ganegoda and Evans (2012) use normalised quantile residuals,  $\hat{r}_i$ , to verify the adequacy of the fitted GAMLSS models. For a response variable  $Y$  with a continuous cumulative distribution function  $F(y; \theta)$ , the normalised quantile residuals are defined as  $\hat{r}_i = \Phi^{-1}[F(y; \hat{\theta})]$ , where  $\Phi^{-1}$  is the inverse cumulative distribution function of the standard Normal distribution. According to Rigby and Stasinopoulos (2005), the error  $\hat{r}_i$  should be standard Normally distributed if the model is adequate. We show the QQ plot of the estimated residuals against the theoretical quantiles of the standard Normal distribution in Figure 5, graphically indicating the normality of the estimated residuals.



**Figure 5: Normalised quantile residual plot**

As further validation of the proposed model, 1 000 000 losses were simulated from the model with parameters presented in Table 8. The goodness-of-fit was tested by comparing the quantiles of the simulated losses with the observed losses using a QQ plot. The QQ plot is shown in Figure 6. The QQ plot follows a fairly straight line confirming the model.





**Figure 6: QQ plot of simulated losses vs observed losses**

### 5.9 Modelling above a threshold

It was previously suggested that the purpose of this model could be to assist banks with their scenario analysis process, which may be specifically helpful for banks with limited internal data. Equation 3 derived in Chapter 4 is used to determine the quantiles of the aggregate loss distribution and compare it with the assessments provided by experts as part of the scenario analysis process. However, to show how this proposed model can be used to determine these quantiles, consideration should be given to the fact that the model is based on amounts above a certain threshold.

For the discussion below, recall the percentile or 1-in- $c$  years approach referred to in Section 4.2. It was explained that scenario makers are asked the following question: “What loss level  $q_c$  is expected to be exceeded only once every  $c$  years”, with popular choices for  $c$  being 10, 20 and 100 years.

Previously, a spliced distribution function was constructed, using backwards-looking historical information for the “expected” (or “body”) part of the distribution and

forward-looking scenario information for the “unexpected” (or “tail”) part. A number  $b$  was selected with the corresponding quantile  $q_b$ , and  $T_e(y)$  was the conditional distribution function of a random loss where  $Y \leq q_b$ .  $T_u(y)$  was the conditional distribution function given that  $Y > q_b$ . From Equation 1, the distribution function for  $F_u(q_c)$  was given by:

$$F_u(q_c) = \frac{[F(q_c) - F(q_b)]}{[1 - F(q_b)]} \text{ for } q_c > q_b. \quad (10)$$

Because the model described under Section 5.7 only model losses greater than USD1 million, only the unexpected part of the severity distribution explained above is effectively modelled. An appropriate allowance needs to be made for the expected part of the distribution if the model is to be used to determine capital estimates in the tail of the distribution. If the necessary allowance for losses below USD1 million is not made, the model will under-estimate the required risk capital. To explain this further using the notation set out in Equation 7,  $q_b$  is equal to USD1 million, although this is a pre-determined amount and not explicitly related to the loss amount only exceeded every  $b$  years. The probability that losses would exceed USD1 million is not known, i.e.,  $P(X > 1) = 1 - F(1)$ , and for comparative purposes it is assumed that  $F(1)$  is between 0.95 and 0.98. Although not exact, this assumption is based on data from the Loss Data Collection Exercise done by Basel in 2008.

This means that the quantiles need to be adjusted because we are conditionally modelling above USD1 million. Table 9 shows the adjusted probabilities for different values of  $F(1)$ , i.e., the cumulative probability that losses would be less than USD1 million. The probabilities are calculated using Equations 3 and 10 and assuming an annual frequency of 6.58627. The reason for selecting this value for the annual frequency is explained in the following section.

**Table 9: Adjusted probabilities for different values of  $F(1)$** 

Scenario point	Cumulative prob. on $F(\cdot)$	Cumulative prob. on $F_u(\cdot)$ for values of $F(1)$			
		0.95	0.96	0.97	0.98
1-in-10 year	0.984817	0.696338	0.620423	0.493897	0.240845
1-in-20 year	0.992408	0.848169	0.810211	0.746948	0.620423
1-in-100 year	0.998482	0.969634	0.962042	0.949390	0.924085
99.9% VaR	0.999848	0.996963	0.996204	0.994939	0.992408

**5.10 Results for an individual bank**

In this section, it is shown how an individual bank can utilise the model that was built on SAS data. It is assumed that the internal loss data for an individual bank is available. For this purpose, the loss data for the Bank of America Corporation were extracted from the SAS database. In the remainder of this chapter, Bank of America Corporation is referred to as the individual bank under question. Table 10 summarises the number of operational losses above USD1 million for the individual bank and reported in the SAS database. The bank itself is expected to have a significantly higher number of data points, given that the internal data would not suffer from reporting bias. In addition, it is expected that the internal data would include information on losses below USD1 million. This information could be used to model the body of the severity distribution, although in using the proposed model, the quantiles are adjusted to allow for this fact.

**Table 10: Bank of America Corporation loss data points per business line and event type**

Event type	Business line					Total
	Clients, Products & Business Practices	Employment Practices & Workplace Safety	Execution, Delivery & Process Management	External Fraud	Internal Fraud	
Agency Services	2	-	-	-	-	<b>2</b>
Asset Management	10	-	1	-	-	<b>11</b>
Commercial Banking	1	-	-	3	1	<b>5</b>
Corporate Finance	10	1	-	-	-	<b>11</b>
Insurance	1	-	-	-	-	<b>1</b>
Payment and Settlement	-	-	1	-	-	<b>1</b>
Retail Banking	26	2	1	21	7	<b>57</b>
Retail Brokerage	29	9	6	-	7	<b>51</b>
Trading & Sales	34	1	5	1	3	<b>44</b>

The model was rerun, but this time excluding the loss data of the individual bank. The results of the new model are shown in Table 11.

**Table 11: Re-estimated parameter values for the model, excluding the Individual bank's data**

Explanatory variable	$\mu$		$\sigma$	
	Estimate	Std. Error	Estimate	Std. Error
Intercept	-0.412680	0.085442	0.402985	0.081920
Log-assets	0.027787	0.006581	-0.02033	0.007212
Equity ratio	-	-	-	-
Africa. Other				
Americas. Other	-	-	-	-
Asia			-0.089890	0.037436
Europe	-	-	-	-
AS. AM & PS	0.137837	0.040684	-	-
Commercial				
Banking	0.250548	0.031946	-	-
Corporate Finance	0.193883	0.049640		
Insurance	-	-	-	-
Retail Banking	-0.093410	0.024465	-	-
Retail Brokerage	-	-	-0.212340	0.046812
Clients. Products & Business Practices	0.082871	0.043878	-	-
Execution. Delivery & Process Management	-	-	-0.283320	0.072153
External Fraud	-0.333850	0.047060	-	-
Internal Fraud	-0.203750	0.047088	-	-

1 000 000 losses were simulated for two business lines, namely retail banking and retail brokerage, using the new model results. It would have been ideal to simulate losses only for one event type within a business line (for example, external fraud in retail banking). However, given the limited number of data points per individual bank, all losses were grouped over event type within a single business line, i.e., assuming that event types are independent.

In order to obtain estimates for a 1-in-10 year, 1-in-20 year and 1-in-100 year loss per business line, an assumption had to be made for  $\lambda$ , the annual loss frequency of losses. For this, refer to Ganegoda and Evans (2012), where they approximated that a bank with USD1 billion assets would experience 0.00823 losses per year, based on data from a Loss Data Collection Exercise done by Basel in 2008. They further show that the total number of losses per year can be weighted to obtain a frequency for each business line and event type within the bank. Using a similar approach and assuming that the individual bank has assets of USD2 trillion (based on the SAS data), the estimated annual frequencies for retail banking was 6.586 and 1.485 for retail brokerage.

In Table 12, the model estimates for scenario points for a 1-in-10 year, 1-in-20 year and 1-in-100 year loss are shown for the two business lines. These correspond to the quantiles for the adjusted probabilities of the fitted distribution, as shown in Table 9. Note that only the scenario point estimates for the assumption that  $F(1) = 0.98$  is shown, i.e. there is a 0.02 probability that losses are above USD1 million. The 1-in-1000 year estimate is also provided. This would be the amount corresponding to the 99.9% Value-at-Risk and the regulatory capital required for the business line. In addition to the point estimates, we show the distribution free 90% confidence intervals for these quantiles.

Given that loss data specific to the individual bank or “internal loss data” is also available, this data could be used in isolation to fit a model specific to the individual bank. The concern with this approach is that the data, especially when working within a specific business line and event type, is limited, as shown in Table 10.

Only publicly available data for the individual bank is considered. However, even if one had access to all the bank’s collected data, it tends to be limited and even more so for higher losses. This point also illustrates the need for banks to augment their own internal data with data from external sources. For the same two business lines under consideration, a Gamma distribution is fitted only to the internal data points. The same quantiles estimated from these models are compared to the estimates from the GAMLSS model described above.

Table 12 provides a summary of the results obtained from the two models for the two business lines.

**Table 12: Estimated scenario points per business line for different models**

	Retail banking		Retail brokerage	
	Individual bank’s data	Model	Individual bank’s data	Model
1-in-10 year	0.139162	0.265663 (0.265; 0.266)	-	- (-)
1-in-20 year	0.684096	0.877871 (0.876; 0.879)	-	- (-)
1-in-100 year	2.179431	2.296219 (2.292; 2.300)	0.901160	0.845764 (0.844; 0.847)
1-in-1000 year	4.469461	4.404938 (4.392; 4.419)	2.518588	2.917611 (2.911; 2.924)

Table 12 shows that the estimated scenario points for the retail banking business line are similar for both models. The first model is based on internal data, and the second model is on external data but tailored for the unique explanatory variables specific to the individual bank (and include 90% confidence intervals). For the retail brokerage business line, where the internal data is even more scarce, the difference between the estimates of the two models is more significant.

The estimated scenario points for 1-in-10 years and 1-in-20 years are zero for the retail brokerage business line. This is because the estimated annual frequency of losses in this business line is only 1.485. As a result, the individual bank is not expected to observe losses higher than USD1 million in this business line in 10 or even 20 years.

### **5.11 Model application within a bank**

This chapter showed how the SAS® OpRisk Global Data could be used to estimate the severity distribution of losses. It is assumed that experts or scenario makers are asked to answer the following question: ‘What loss level is expected to be exceeded once in  $c$  years?’. Given the explanatory variables for a specific bank, the distribution  $F(y; \hat{\theta})$  may be used to determine quantiles of the aggregate loss distribution, which can be compared to the scenario assessments of the experts. Once an appropriate distribution function has been selected, the quantiles can be determined that relate to the scenario assessments provided by the experts. For example, the 1-in-100-year loss predicted by the expert should be in line with the 99% quantile of the aggregate loss distribution. Therefore, if the loss scenario points provided by the experts deviate too far from the quantiles of the loss distribution that was estimated by the data, one can revert to the expert and request them to justify the difference. Using internal and external data, specifically for units of measure where adequate historical data is available, one should model future expected losses reasonably well. However, the more significant benefit of the scaling model is for banks where very limited or no internal within a business line is available. In such a case, the bank may use the model based on external data and use its own characteristics to infer values for expected future losses.



## 5.12 Conclusion

It was shown how SAS® OpRisk Global Data could be used by a bank when they do not have their own internal loss data, to build statistical capital models. Suggestions were made on how a bank can use a model only based on external data to inform or challenge the scenario assessments provided by experts. Scenario assessments are often used as a significant component of operational risk management. However, given the subjective nature of these assessments, it is vital to have an objective measure to check whether the expert's opinion is not biased or completely unrealistic. Although experts may not change their views based on the results of statistical models, they may be required to justify why their assessments deviate from the data. The suggested model considers the reporting bias included in any external database and shows that operational losses depend on certain factors specific to a bank, such as size, region, business line and event type associated with operational losses.

## **6. Concluding remarks and recommendations for future research**

The focus of this dissertation was on quantitative operational risk models. We showed how the various data sources available to financial institutions can be utilised to obtain accurate estimates of operational risk capital.

A review of the literature was done to gain a better understanding of the risk management process and the importance of operational risk management within financial organisations. Despite the standardisation of regulation relating to operational risk capital calculations, we motivated the continued need for advanced statistical models in determining economic capital.

The underlying methodology used in our operational risk capital model is the loss distribution approach (LDA) which makes use of an annual aggregate loss distribution. The components of the aggregate loss distribution, namely the frequency and severity distributions, were discussed in some detail and the subsequent chapters focused on the improvement of the estimation of the severity distribution by using various data sources.

It was noted that there are four main sources of data available to financial institutions for use in their capital risk models, namely internal data, external loss data, scenario assessment and business environment and internal control factors. The latter were not considered in any detail in this dissertation. We have however investigated the other three data sources in some detail and also specifically how these data sources can be combined or used to complement each other.

Statistical methods were explored that could be used to combine limited historical data and scenario assessments to estimate extreme quantiles. We also showed how external data, namely SAS® OpRisk Global Data, could be used by a bank when they do not have their own internal loss data, to build statistical capital models. Suggestions were made on how a bank can use a model only based on external data to inform or challenge the scenario assessments provided by experts.

As far as future research is concerned, we aim to investigate the effectiveness of using the ratios suggested in Chapter 4 in assisting scenario experts with their assessments. We further aim to investigate the use of other external data sources that could be more appropriate for banks operating in South Africa specifically.

## Reference list

- (BCBS), B. C. o. B. S., 2011b. *Operational risk: Supervisory guidelines for the advanced measurement approaches. Report 196.*, s.l.: s.n.
- Amin, Z., 2016. *Quantification of operational risk: A scenario-based approach*, s.l.: s.n.
- Aue, F. & Kalkbrener, M., 2007. LDA at work: Deutsche Bank's approach to quantifying operational risk. *Journal of Operational Risk*, 1(4), pp. 49-93.
- Basel Committee on Banking Supervision, 2006. *International convergence of capital measurement and capital standards. A revised framework.*, Basel, Switzerland: s.n.
- Basel Committee on Banking Supervision, 2009. *Results from the 2008 loss data collection Exercise for operational risk*, s.l.: s.n.
- Basel Committee on Banking Supervision, 2011a. *Principles for the sound management of operational risk. Report 195.*, s.l.: s.n.
- Basel Committee on Banking Supervision, 2011b. *Operational risk: Supervisory guidelines for the advanced measurement approaches. Report 196.*, s.l.: s.n.
- Basel Committee on Banking Supervision, 2017. *High-level summary of Basel III reforms*, s.l.: s.n.
- Basel Committee on Banking Supervision, 2019. *OPE 30 - Advanced Measurement Approaches*, s.l.: s.n.
- Baud, N., Frachot, A. & Roncalli, T., 2002. *Internal data, external data and consortium data for operational risk measurement: How to pool data properly?*, Paris: Groupe de Recherche Opérationnelle, Crédit Lyonnais, France.
- Benito, S. & Lopez-Martin, C., 2018. A review of the state of the art in quantifying operational risk. *Journal of Operational Risk*, 13(4), pp. 89-129.
- Böcker, K. & Klüppelberg, C., 2005. Operational VaR: a Closed-Form Approximation. *Risk*, pp. 90-93.
- Cope, E. & Labbi, A., 2008. Operational loss scaling by exposure indicators: evidence from the ORX database. *The Journal of Operational Risk*, 3(4), pp. 25-45.
- Cummins, J., Lewis, C. M., & Wei, R. (2006). The market value impact of operational loss events for US bank and insurers. *Journal of Banking & Finance*, 30, 2605-2634.
- Dahen, H. & Dionne, G., 2010. Scaling models for the severity and frequency of external operational loss data. *Journal of Banking & Finance*, 34(7), pp. 1484-1496.
- De Fontnouvelle, P., DeJesus-Rueff, V., Jordan, J. & Rosengren, E., 2006. Capital and risk: New evidence on implications of large operational losses. *Journal of Money, Credit, and Banking*, p. 1819–1846.
- De Jongh, E., De Jongh, D. & De Jongh, R., 2013. A review of operational risk in banks and its role in the financial crisis. *The South African Journal of Economic and Management Sciences*, Issue 4, pp. 364-382.

- De Jongh, P., de Wet, T., Raubenheimer, H. & Venter, J., 2015. Combining scenario and historical data in the loss distribution approach: a new procedure that incorporates measures of agreement between scenarios and historical data. *Journal of Operational Risk*, 10(1), pp. 45-76.
- De Jongh, P., De Wet, T., Raubenheimer, H. & Venter, J., 2015. Combining scenario and historical data in the loss distribution approach: a new procedure that incorporates measures of agreement between scenarios and historical data. *Journal of Operational Risk*, 10(1), pp. 45-76.
- De Jongh, R., Raubenheimer, H. & Gericke, M., 2021. Construction of forward-looking distributions using limited historical data and scenario assessments. In: M. K. T. a. G. Djurovic, ed. *Linear and non-linear financial econometrics: Theory and Practice*. London, UK: IntechOpen, pp. 13-30.
- Dutta, K. K. & Babbel, D. F., 2014. Scenario analysis in the measurement of operational risk capital: A change of measure approach. *The journal of risk and insurance*, 81(2), pp. 303-334.
- Embrechts, P. & Hofert, M., 2011. Practices and issues in operational risk modeling under Basel II. *Lithuanian Mathematical Journal*, 51(2), pp. 180-193.
- Embrechts, P. & Hofert, M., 2011. Practices and issues in operational risk modelling under Basel II.. *Lithuanian Mathematical Journal*, 51(2), pp. 180-193.
- Ganegoda, A. & Evans, J., 2012. A scaling model for severity of operational losses using generalized additive models for location scale and shape (GAMLSS). *Annals of Actuarial Science*, 7(1), pp. 61-100.
- Kelliher, P. et al., 2016. *Good practice guide to setting inputs for operational risk models*, Edinburgh: Institute and Faculty of Actuaries.
- Lambrigger, D. D., Shevchenko, P. V. & Wüthrich, M. V., 2007. The quantification of operational risk using internal data, relevant external data and expert opinions. *The journal of operational risk*, 2(3), pp. 3-27.
- McNeil, A., Frey, R. & Embrechts, P. 2015. *Quantitative risk management: Concepts, techniques and tools. Revised Edition*.. Princeton and Oxford: Princeton University Press.
- Panman, K., Van Biljon, L., Haasbroek, L. J., Schutte, W.D. and Verster, T. 2019. Quantification of the estimation risk inherent in loss distribution approach models. *The Journal of Risk Model Validation*, 13(4), 1-25.
- Peters, G. W., Shevchenko, P. V., Hassani, B. & Chapelle, A., 2016. *Standardized Measurement Approach for operational risk: Pros and Cons*. [Online] Available at: <https://ssrn.com/abstract=2789006> [Accessed 26 April 2021].
- Power, M., 2005. The invention of operational risk. *Review of International Political Economy*, October.
- Prudential Regulation Authority, 2020. *The PRA's methodologies for setting Pillar 2 capital*, London: Bank of England.
- Rigby, R. & Stasinopoulos, D., 2005. Generalized additive models for location, scale and shape. *Journal of the Royal Statistical Society*, pp. 507-554.
- SAS, 2021. [Online] Available at: [https://www.sas.com/content/dam/SAS/en\\_us/doc/productbrief/sas-oprisk-global-data-101187.pdf](https://www.sas.com/content/dam/SAS/en_us/doc/productbrief/sas-oprisk-global-data-101187.pdf)

Shih, J., Samad-Khan, A. & Medapa, P., 2000. Is the size of an operational loss related to firm size?. *Operational risk magazine.*, pp. 1-2.

Supervision, B. C. o. B., n.d. *History of the Basel Committee*. [Online]

Available at: <https://www.bis.org/bcbs/history.htm>

[Accessed 27 July 2021].

Sweeting, P., 2011. *Financial Enterprise Risk Management*. First edition ed. Cambridge: Cambridge University Press.

The Actuarial Education Company, 2018. *Subject ST9 Combined Materials Pack*. s.l.:s.n.

The Risk Management Association, 2011. *Scenario analysis: Perspectives and principles*, USA: s.n.

Wei, L., Li, J. & Zhu, X., 2018. Operational Loss Data Collection: A Literature Review. *Annals of Data Science*.

Wilson, S., 2007. *A review of correction techniques for inherent biases in external operational risk loss data*, s.l.: s.n.

## Appendix A: The standardised measurement approach

A brief overview of the standardised measurement approach (SMA) that will apply from 1 January 2023 is provided in this section. We also summarise some of the critique against this approach.

### Standardised measurement approach

The new standardised measurement approach calculates operational risk capital requirements based on measures for both the bank's income and using the bank's historical losses. The operational risk capital requirement can be summarised as follows:

$$\text{Operational risk capital} = \text{Business Indicator Component (BIC)} * \text{Internal loss multiplier (ILM)}$$

The BIC is calculated by multiplying the Business Indicator (BI) by marginal coefficients. The BI consists of the following three elements, all calculated as an average over three years:

- The interest, leases and dividend component;
- The services component; and
- The financial component.

In order to calculate the BIC, the BI is multiplied with the marginal coefficients according to a progressive sliding scale (similar to standard tax formulas) as set out in Table A1.

**Table A1: Marginal coefficients to calculate BIC**

Bucket	BI range (in EUR billions)	BI marginal coefficient
1	$\leq 1$	12%
2	$1 < BI \leq 30$	15%
3	$> 30$	18%

The Internal Loss Multiplier (ILM) is calculated from the Loss Component (LC), being 15 times a bank's average historical losses over the previous ten years, and the BIC, using the following formula:

$$ILM = \ln \left( \exp(1) - 1 + \left( \frac{LC}{BIC} \right)^{0.8} \right).$$

The ILM increases as the  $(LC/BIC)$  ratio increases, although at a decreasing rate. National supervisors are awarded some discretion in setting the ILM for all banks in their jurisdiction to 1, i.e., the capital requirement would only depend on the value of the BIC. However, all banks would still be required to disclose their historical losses.

Below, we summarise some of the critique against the new standardised approach and provide a motivation for continued research into quantitative operational risk models.

### **Critique against the new standardised approach**

A group of academics (Peters *et al.*, 2016) compiled an initial response to the proposed standardised measurement approach (SMA), independently from corporate or individual interests. In their view, the SMA: introduces capital instability and extreme sensitivity to the dominant loss process; reduces risk responsiveness and interpretability; incentivises enhanced risk-taking; fails to utilise a range of data sources and fails to provide risk management insight, and introduce the possibility of superadditive capital calculation.

Some of the other critique against the standardised approach include that it allows national supervisors to decide whether historical loss data will be included in calculating operational risk capital. Capital levels will entirely be based on the Business Indicator when historical loss data is excluded, i.e., it would be purely based on a bank's income level, rather than factoring in any historical operational risk-related losses. Even where historical loss data is included, the fact that the natural log is used in calculating the loss component means variations in losses produce only a small effect on capital.

It is widely opined that the new approach is not as valuable for managing operational risk as the internal models developed under the Advanced Measurement Approach. Banks have already spent considerable resources developing these models. For this reason, these models will continue to be used for regulatory capital purposes until the introduction of the standardised approach in 2023 at the earliest. It is therefore prudent for banks to continue



refining their own internal operational risk models. There is also strong evidence to suggest that internal models will continue to be used for Pillar II capital and as part of the supervisory review process. The investment in these models is, therefore, not only for short-term use.

## **Appendix B: Guidelines for using the advanced measurement approach**

The BCBS published OPE30 in December 2019 (Basel Committee on Banking Supervision, 2019), setting out the criteria that banks have to meet to calculate operational risk capital requirements based on their internal risk measurement systems. It specifies that the regulatory capital requirement should equal the risk measure generated by the bank's internal operational risk measurement system and be based on quantitative and qualitative criteria for the advanced measurement approach (AMA) as set out in that document. The regulatory guidelines provided under Basel II's advanced measurement approach remain relevant and provide valuable guidelines for developing operational risk models appropriate for the future.

There are several qualitative standards worth noting when a bank intends to use the AMA, and only the most relevant are highlighted below:

- The bank must have an independent operational risk management function responsible for designing and implementing the operational risk management framework.
- The internal operational risk measurement system must be integrated into the day-to-day risk management processes. The bank's measurement system must support an allocation of economic capital for operational risk across business lines in a manner that creates incentives to improve business line operational risk management.
- Operational risk exposures must be regularly reported to business unit management, senior management, and directors.
- The internal operational risk measurement system must be well documented.
- The operational risk management processes must be subject to validation and regular independent review.
- External auditors or supervisors must review the operational risk assessment system regularly.

The quantitative standards that should apply when using the AMA, include:

- The bank must be able to demonstrate that its approach captures potentially severe tail-loss events. Further, it must demonstrate that its chosen operational risk measure meets a soundness standard comparable to the internal ratings-based approach for credit risk (one year holding period, 0.999-quantile).
- Banks must have and maintain rigorous procedures for operational risk model development and validation.
- Any internal measurement system must be consistent with the scope of operational risk and the following seven events: Internal fraud; external fraud; employment practices and workplace safety; clients, products, and business practices; damage to physical assets; business disruption and system failures; execution, delivery, and process management.
- The bank is required to calculate its regulatory capital charge as the sum of expected loss (i.e., the mean of the loss distribution) and unexpected loss (i.e., the difference between the Value-at-Risk (VaR) and expected loss).
- The risk measurement system must be sufficiently granular to capture the significant drivers of operational risk.
- Risk measures for different operational risk estimates must be added for calculating the minimum capital requirement. However, the bank may use internally determined correlations of losses, provided that they are determined with sound methods, implemented with integrity, consider the uncertainty of the correlation estimates, and those correlation assumptions are validated using quantitative and qualitative techniques.
- Any operational risk measurement system must use the four data elements, namely internal data, relevant external data, scenario analysis, and business environment and internal control factors.
- The bank needs to have a credible, transparent, well-documented, and verifiable approach for weighting these elements.
- The bank's internal measurement system must reasonably estimate unexpected losses based on the combined use of internal and relevant external loss data, scenario analysis and bank-specific business environment and internal control factors. The guidelines

that should be followed in using the four data sources, are summarised under the headings below.

OPE30 (Basel Committee on Banking Supervision, 2019) also provides the following guidelines with respect to the four main data sources referred to in this dissertation, namely internal data, external loss data, scenario assessments and business environment and internal control factors.

### **Internal data**

- The tracking of internal loss event data is an essential prerequisite to the development and functioning of a credible operational risk measurement system.
- Risk measure estimates based on internal-loss data must be calculated with a minimum of five years of historical observations.
- Aside from gross loss amounts, the bank should collect information about the date of the event, recoveries of the gross loss amount, and descriptive information about the drivers of the event.
- The internal loss data must be comprehensive in that it captures all material activities and exposures. There must be an appropriate minimum gross loss threshold for internal loss data collection, e.g., 10 000 Euros.
- A bank must have documented procedures for assessing the ongoing relevance of historical loss data, e.g., for scaling historical data, and who is authorised to make the corresponding decisions.
- A bank must map its historical internal loss data into the eight business lines, and seven event types specified in the Basel accords and provide these data to supervisors. They must develop criteria for assigning loss data from an event in a centralised function (e.g., information technology department) or events in time to this eight-by-seven matrix of business lines and event types. The allocation criteria must be documented and objective.
- Operational risk losses that were historically included in the bank's credit risk database should be treated as credit risk to calculate the minimum capital requirement. However, banks must identify these losses for internal operational risk management.
- Operational risk losses related to market risk are treated as operational risk losses to calculate the minimum capital requirement.

**External loss data:**

- External loss data, comprising the operational risk losses experienced by third parties, can offset the scarcity of internal operational risk loss data in areas where a bank has a potential risk but has not experienced significant losses (i.e., from a forward-looking perspective).
- The bank's operational risk management system must use relevant external data in the form of loss amounts, information on the scale of business operations, and information on the causes and circumstances of the loss events. Further, a bank must have a systematic process for determining the situations which require external data and the methodologies used to incorporate the data.

**Scenario assessments:**

- The bank must use scenario analysis of expert opinion in conjunction with external data to evaluate its exposure to high-severity events (e.g., the expert assessments could be expressed as parameters of a loss distribution).
- Scenario analysis should also be used to assess the impact of deviations from the correlation assumptions embedded in the bank's operational risk measurement framework, in particular, to evaluate potential losses arising from simultaneous operational risk loss events.
- Scenario assessments need to be validated and re-assessed through comparison to actual loss experience over time to ensure their reasonableness.

**Business environment and internal control factors:**

- Each factor should be a meaningful driver of risk, based on experience and expert judgment. The factors should be translatable to quantitative measures.
- The sensitivity of the estimated risk to changes in the factors and the relative weighting of the factors must be well reasoned.
- The use of the factors in a bank's risk measurement framework must be documented and be subject to independent review.
- Over time, the processes and outcomes need to be compared to actual internal loss experience and relevant external data, and appropriate adjustments must be made.