# NOVA IMS
Information Management School

# MAA

## Mestrado em Métodos Analíticos Avançados
Master Program in Advanced Analytics

## DEVELOPMENT AND IMPLEMENTATION OF THE PROFITABILITY RISK MODULE PROCESS

Mariana Coutinho de Lucena e Romão

Project Work presented as partial requirement for obtaining a Master's degree in Advanced Analytics

**NOVA Information Management School**
**Instituto Superior de Estatística e Gestão de Informação**
Universidade Nova de Lisboa

# DEVELOPMENT AND IMPLEMENTATION OF THE PROFITABILITY RISK MODULE PROCESS

by

Mariana Coutinho de Lucena e Romão

Project Work presented as partial requirement for obtaining a Master's degree in Advanced Analytics

**Advisor:** Vítor Duarte dos Santos

July 2023

# ACKNOWLEDGEMENTS

# ABSTRACT

The main objective of this report is to outline the project carried out at Neyond, where the main goal was developing an automated E.T.L. reporting process for an international bank, one of the clients.

This project played a crucial role in solidifying the knowledge gained and implementing the diverse techniques learned throughout the initial academic year. It also provided an opportunity to merge academic training with practical professional experience. This report provides an overview of the project's goals, methodologies, tools, technologies used, and the challenges encountered during its execution.

To accomplish this objective, this report describes the project and the main goals to achieve the desired result, the tools and technologies used, as well as some of the challenges and how they were surpassed.

# KEYWORDS

# INDEX

## LIST OF FIGURES

# LIST OF TABLES

# LIST OF ABBREVIATIONS AND ACRONYMS

**BD** Big Data

**BDA** Big Data Analytics

**ETL** Extract Transform Load

**FC** Financing Cost

**GCB** General terms and conditions of business, delivery and payment

**HDFS** Hadoop Distributed File System

**KPMs** Key Performance Measures

**MF** Mutual Funds

**OLAP** Online Analytical Processing

**ROF** Return on Investment Versus Failure

**SMV** Standard Minute Value

**SQL** Structured Query Language

**UDF** User-defined Function

# 1. INTRODUCTION

This paper will address the workflow of a project that has already been developed. However, due to constant changes and developments in the industry the implementation of a new process was necessary, the risk process. This has come as a new challenge and as an excuse to simplify the already implemented profitability process as the two are connected. Adding to this workflow the context and goals are described as well as the main technologies and methodologies, summed up in a conclusion.

The project report, which is equivalent to a Masters dissertation, is the final assignment to complete my Masters degree in Data Science and Advanced Analytics at the Nova Information Management School. This project at the technological team of Neyond will give me experience and key learnings as a Data Scientist.

## 1.1.  MOTIVATION

With more than 600 clients allied to the company my project focuses on developing one of them, a bank. This project resulted from the need for the management control department to have access to information regarding the profitability and risk associated with the banks and accounts associated with the bank, so my role was to create a process that facilitates what is already done manually and transform it into a more automated process, adding some criteria and metrics that the department needs to have, in order to manage it in the most efficient way.

Currently, information on risk management and performance metrics is calculated and reported manually by the area. This process being manual, like any manual processes, has a set of problems ranging from quality of information to human error. The goal of this project, as previously mentioned, is to improve and automate the reporting process and its generation in excel format, regarding a monthly or redeployed process of data, periodicity that will be discussed further in the paper.

ETL (Extract Transform Load) aims to collect, filter, process, and combine relevant data from various sources and databases to store in a data warehouse (Hendayun, M., Yulianto, E., Rusdi, J. F., Setiawan, A., & Ilman, B. 2021). This can be accomplished considering all the formulas for the metrics that are requested by the bank.

Overall, by using financial technology, commercial banks can improve their traditional business model by reducing bank operating costs, improving service efficiency, strengthening risk control capabilities, and creating enhanced customer-oriented business models for customers; thereby improving comprehensive competitiveness (Wang, Y., Xiuping, S., & Zhang, Q.  2021).

This creates a challenging project to develop and learn the skills required to be a data analyst, modeling all types of data and creating a final product that will consequentially facilitate the present processes implemented in the bank.

This process poses a challenge as even with the information provided by the client, a lot of the process needs to be reinvented because the data can be difficult to obtained, as it comes from different universes that contain different keys that need to be crossed with information from other data universes.

## 1.2. PROJECT GOALS

The main project goal is to implement a process to automate and calculate KPM's that the company needs to report in a determined period.

In order to achieve this goal, the following intermediate objectives were defined:

- Understand Data Fundamentals in Azure
- Get to know the process of a Developer
- Understand the Risk and Profitability Process
- Understand the formulas and KPM's (metrics that help to understand the business and its behavior in determined aspects) of the company
- Analyse the development and production dependencies with other teams
- Develop and Improve the Process, checking the tables and data that require retrieval to be the most efficient when creating the final product

## 2. LITERATURE REVIEW

The theoretical knowledge and framework conditions related to this project come from various fields such as finance, business process management and IT architecture. This chapter focusses on these areas as they represent the core aspects of the project's tasks. Both fields have been broadly scientifically covered and therefore the following sections reflect the essence of these topics relevant to the project.

### 2.1. RISK AND PROFITABILITY

When talking about risk we need to understand that the notion of risk is intimately linked to the return. Return includes ensuring remuneration of production factors and invested capital but also resource management in terms of efficiency and effectiveness. A full financial and economic diagnosis cannot be done without regard to the return-risk ratio.

In addition, risk analysis is useful in decision making concerning the use of economic-financial potential or investment decisions, in developing business plans, and also to inform partners about the enterprise's performance level (Solomon, D. C., & Muntean, M. 2012).

On one hand, financial risk analysis uses specific indicators such as: financial leverage, financial breakeven and leverage ratio (CLF) accompanying call to debt, presents a major interest to optimize the financial structure and viability of any company operating under a genuine market economy (Solomon, D. C., & Muntean, M. 2012). On the other hand, financial risk management avoids losses and maximizes profits, and hence is vital to most businesses as the task relies heavily on information-driven decision making (Mashhour, A., Luo, W., Zaidi, N. A., & Robles-Kelly, A. 2020). To sum up, it is possible to conclude that most business analytics tools are used to improve risk management therefore, risk management tools benefit from business analytics approaches (Steshenko, O., & Bondarenko, Y. 2021).

Furthermore, profitability is one form of expressing economic efficiency with summarizing the efforts made to obtain the expected results. Profitability rates measure the results obtained in relation to the activity of companies (commercial profitability) with economic means (economic profitability) or financial means (financial profitability) (Hada, I. D., & Mihalcea, M. M. 2020).

For this reason, financial performance is a major point of interest for both the internal and external environment of an economic entity. To be prosperous, attractive, efficient, and promise development, a company must obtain a profit. In the conditions of a dynamic economic environment, assailed by many changes, maximizing profitability or the ability to make a profit as a measure of performance is the main objective of the activity of an economic entity (Hada, I. D., & Mihalcea, M. M. 2020).

Ultimately, risk and profitability are closely linked in the financial industry, with an increase in risk often leading to an increase in potential profitability. Nonetheless, managing risk is crucial for long-term success. As stated in "The Relationship between Risk and Return" by Jorion (2007), "Risk and return are positively related, but the relationship is not linear: higher risk usually leads to higher expected return, but not always."

In order to effectively manage risk and increase profitability, organizations must have a strong risk management strategy in place. According to a study by McNeil, Frey, and Embrechts (2015) "Quantitative Risk Management: Concepts, Techniques and Tools," a good risk management strategy should include identifying potential risks, assessing their likelihood and impact, and taking action to mitigate or avoid them.

## 2.2. BUSINESS METRICS

The creation of prosperous business models has become essential in fluctuating business environments. However, the attention given to the role of metrics in crafting business models in literature is limited compared to the research on business modeling tools (Heikkilä, M., Bouwman, H., Heikkilä, J., Solaimani, S., & Janssen, W. 2016).

With the advancement of analytical capabilities, the traditional methods of selecting business metrics are changing, and it depends on various factors such as the scope of the business, the competitive environment, customer behavior, and channels of interaction. Analysts should take into account all sources of information while choosing the correct metrics, as well as consider who will be using this data to make business decisions (Panchyshyn, T., & Prokopovych-Pavlyuk, I. 2021).

Business Intelligence can aid in supporting operational to strategic business decisions, and it is necessary to optimize the metrics for the business to keep up with company development (Elveny, M., Nasution, M. K. M., Zarlis, M., & Efendi, S. 2021). Metrics are quantified measures that help identify what needs to be done and by whom. They should be linked to the organization's objectives and goals and indicate how well they are being met. Metrics should also motivate individual, group, or team action and continuous improvement (Sisco & Chorn, 2009).

Measuring the performance of business processes is vital for effective and efficient results. The choice of performance indicators is organization-dependent and must align with the business strategy (Van Looy, A., & Shafagatova, A. 2016). The timely and accurate measurement of appropriate business metrics is crucial in formulating new business strategies, and monitoring and measuring their effective use is essential for managers to

understand the relationship between their processes and the market space (Askar, M., Imam, S., & Prabhaker, P. R. 2009).

However, there are challenges in the organizational dimension of the decision-making process, defining strategic key metrics, data structuring, and data entry quality, which need to be considered for the optimal use of business analytics (Housbane, S., Khoubila, A., Ajbal, K., Serhier, Z., Agoub, M., Battas, O., & Othmani, M. B. 2020). Several studies have looked at the opportunities and challenges of using business metrics in organizations, highlighting the limited attention given to the role of metrics in designing business models and the challenges in the decision-making process, data structuration, and quality of data entry.

## 2.3. BIG DATA

Developing Big Data applications has become increasingly important in the last few years. In fact, several organizations from different sectors depend increasingly on knowledge extracted from huge volumes of data (Oussous, A., Benjelloun, F. Z., Ait Lahcen, A., & Belfkih, S. 2018).

The financial field is one of them, as it is deeply involved in the calculation of big data events. As a result, hundreds of millions of financial transactions occur in the financial world each day. Therefore, financial practitioners and analysts consider it an emerging issue of data management and analytics of different financial products and services. Also, big data has significant impacts on financial products and services. Therefore, identifying the financial issues where big data has a significant influence is also an important issue to explore (Hasan, M. M., Popp, J., & Oláh, J. 2020).

According to Davenport and Harris (2007), organizations that utilize data-driven decision making enabled by the collection and analysis of big data have a significant competitive advantage. They also suggest that effective use of big data leads to improved decision-making processes and better performance outcomes.

 Lazer, Kennedy, King, and Vespignani (2014) argue that big data has the potential to revolutionize decision making by providing access to new types of data and analytical methods. They further highlight its benefits in various industries such as finance, healthcare, and manufacturing.

In the realm of marketing, Chiang and Storey (2016) find that big data enhances the accuracy of customer segmentation, a crucial metric.

Chen and Chen (2020) emphasize the role of big data analytics in financial management, stating that it provides insights into financial performance and risk management. They

suggest that big data analytics helps organizations identify patterns, trends, and hidden information to make better decisions and improve financial outcomes.

In addition, big data can also be used to improve the efficiency of business processes Wang, Q., Liu, L., & Wu, J. (2018) found that big data can be used to improve supply chain management by providing access to real-time data on inventory levels, production schedules, and shipping schedules.

It is worth highlighting that although big data holds great potential for enhancing business performance, it presents certain challenges that demand attention. These challenges encompass data privacy and security, data governance and management, as well as data integration. Therefore, businesses need to acknowledge these obstacles and establish strategies to mitigate their impact.

To summarize, by leveraging big data, businesses can enhance the precision and effectiveness of their business metrics, thereby facilitating better decision making and overall performance improvement. Nonetheless, it is crucial for businesses to proactively address and manage the potential challenges associated with big data.

## 2.3.1. Analytics

Big data refers to the large and complex datasets that are generated by various sources such as social media, sensors, and transactional systems. The sheer volume, velocity, and variety of big data make it difficult to be processed and analyzed using traditional data processing techniques (Mayer-Schönberger, V., & Cukier, K. 2013).

Furthermore, Big Data Analytics (BDA) is increasingly becoming a trending practice that many organizations are adopting with the purpose of constructing valuable information from BD (Sivarajah, U., Kamal, M. M., Irani, Z., & Weerakkody, V. 2017). The world has become excited about big data and advanced analytics not just because the data sets are big but also because the potential for impact is big (Court, D. 2015). And so, an organization that quickly adopts new tools and adapts itself to capture their potential is more likely to achieve large-scale benefits from its data-analytics efforts (Court, D. 2015).

The analytics process, including the deployment and use of BDA tools, is seen by organizations as a tool to improve operational efficiency though it has strategic potential, drive new revenue streams and gain competitive advantages over business rivals (Sivarajah, U., Kamal, M. M., Irani, Z., & Weerakkody, V. 2017).

Furthermore, the practice of big data analytics entails the thorough investigation, refinement, conversion, and shaping of substantial datasets in order to reveal valuable

insights and understanding. This intricate process incorporates the utilization of sophisticated analytical tools and methodologies, including machine learning, natural language processing, and statistical analysis, to extract meaningful observations from extensive data (Feldman, R., & Sanger, J. 2013).

The use of big data analytics can provide organizations with a competitive advantage by providing them with insights into customer behavior, market trends, and operational efficiencies. For example, by analyzing customer data, organizations can identify patterns of customer behavior that can be used to improve marketing strategies and personalize customer experiences (Zikopoulos, P., & Eaton, C. 2012).

Big data analytics, despite its benefits, presents certain obstacles such as safeguarding data privacy, ensuring data security, and addressing ethical considerations. Enterprises must prioritize adherence to data protection laws and regulations while implementing big data analytics. They should also establish robust technical and organizational measures to safeguard their customers' personal information (Wang, Y., & Li, Y. 2017).

In addition, it is crucial for organizations to take into account the reliability and precision of the data they are examining. Large volumes of information frequently lack organization and may include inaccuracies, prejudices, and gaps that could result in misleading findings. Hence, organizations must guarantee that the data undergoes meticulous cleansing, preprocessing, and validation prior to its utilization in the analytics procedure (Gandomi, A., & Haider, M. 2015).

In conclusion, big data analytics has the potential to provide organizations with valuable insights, but it also poses some challenges such as data privacy, security, and data quality. Organizations need to be aware of these challenges and address them accordingly to ensure the successful implementation of big data analytics.

## 2.3.2. Challenges and Opportunities in Big Data

Big data analytics has emerged as a powerful tool for businesses to gain insights and make data-driven decisions. However, the vast amount of data and the complexity of analyzing it presents several opportunities and challenges.

One notable opportunity for businesses lies in the potential to gain a competitive edge through effective employment of big data analytics. Furthermore, big data analytics can aid businesses in identifying fresh revenue streams and uncovering market opportunities (KPMG, 2013).

However, big data analytics also presents several challenges. Among these, data quality and accuracy emerge as key concerns (Liu et al., 2015). Without proper data cleaning and preparation, the insights derived from big data analytics may lack accuracy. Additionally, there exists a demand for skilled personnel capable of effectively analyzing and interpreting the data (Fayyad et al., 2016). Consequently, businesses must allocate resources towards training and hiring data scientists and analysts to fully leverage the potential of their big data.

In conclusion, big data analytics provides an extent of opportunities for businesses to gain a competitive advantage and identify new avenues for revenue generation. However, challenges related to data quality, the need for skilled personnel, and privacy and security concerns must be effectively addressed in order to harness the full potential of big data analytics.

## 2.4. CHALLENGES AND OPPORTUNITIES IN RISK MANAGEMENT

Identifying, assessing, and mitigating risks poses a significant challenge in risk management (Hausberg & Wied, 2015). This complexity is often attributed to the extensive data that organizations must process in order to detect risks.

To address this challenge, big data analytics offers valuable assistance by equipping organizations with tools to process and analyze vast amounts of data, enabling more effective risk identification and assessment (Xie, Zhang, & Wang, 2018).

Another hurdle in risk management lies in predicting and mitigating emerging risks (KPMG, 2015). Big data analytics comes to the rescue once again, empowering organizations with tools to discern patterns and trends within large datasets, aiding in the prediction and mitigation of emerging risks (Hsu, 2017).

Nevertheless, incorporating big data analytics into risk management introduces its own set of challenges, including data quality, data privacy, and regulatory compliance (KPMG, 2016). Hence, organizations must carefully consider and address these challenges during the implementation of big data analytics in risk management.

On a positive note, big data analytics in risk management presents numerous opportunities. It enhances the effectiveness of risk identification and assessment, leading to more efficient risk management practices (Xie, Zhang, & Wang, 2018). Additionally, it enables the identification of patterns and trends within extensive data, contributing to the anticipation and mitigation of emerging risks (Hsu, 2017). Furthermore, big data analytics empowers organizations to make informed decisions, thereby enhancing risk management outcomes.

In summary, big data analytics equips organizations with the capability to process and analyze extensive datasets, facilitating more effective risk identification and assessment. However, it is crucial for organizations to consider challenges such as data quality, data privacy, and regulatory compliance when integrating big data analytics into their risk management frameworks.

## 3. WORK PLAN

This chapter will discuss the methodologies applied in the project work on ETL and reporting. The project aimed to efficiently extract, transform, and load data from various sources into a centralized database, and generate meaningful reports to support data-driven decision making.

### 3.1. PROJECT METHODOLOGY

To achieve the goal of this project it was followed a methodological path divided in three different phases: exploration phase, analytical phase, and conclusive phase. Each phase is divided in specific steps to attain the proposed goal.

**Exploration Phase**
- Theoretical framework
- Methodology review and improvment

**Analythical Phase**
- Development of the project
- Presentation of the project

**Conclusive Phase**
- Final product contributions
- Final product presentation
- Final revisions

Figure 3.1 Methodology Diagram

In the exploration phase the theoretical framework will be developed, identifying the theoretical basis for the project, it will describe in detail the main topics necessary to understand the development done in the next two phases. Topics such as Risk and Profitability, Business Metrics, Big Data, Analytics, Challenges and Opportunities in Big Data, Challenges and Opportunities in Risk Management. During this stage previous works that may contribute to the global understanding of the project will be identified. In the end of this phase the methodology can be refined.

In the analytical phase, a project plan provided by the manager of the team and discussed with me will be developed and studied to be applied in order to obtain the results that the

client's wants. To accomplish the final product, it will be necessary to study some aspects of the client needs and what the final product should look like, as it will help in the development of the process to know the global image and the business.

| n | Tasks |
|---|---|
| **1** | **Phase 1** |
| 1.1 | Dump The process of the Book of Business input tables in cloud |
| 1.2 | Migrate process from on prem to DataBricks and run it |
| 1.3 | Unit Tests - Phase 1 (Execution of the Book of Business Process) |
| **2** | **Phase 2** |
| 2.1 | Addition of a column to the parameterization table that defines whether it is the Book of Business Module or Risk Module |
| 2.2 | Adaptation to the Book of Business process to add the possibility of executing the Risk Module |
| 2.3 | Unit Tests - Phase 2 (Execution of the Risks Module process that will only get the intended rules) |
| ***3*** | **Phase 3 – Incorporation of the aggregators script in the Cloud** |
| *3.1* | Analyze aggregator scripts and outline how to adapt to the Book of Business process |
| *3.2* | Integration of the Redeployment module in the Book of Business scripts in the cloud |
| *3.3* | Loading the process input tables of the aggregators in the cloud |
| *3.4* | Addition of the construction of aggregator tables to the process |
| *3.5* | Logic integration in the process if it is necessary to run the Aggregators Module or not |
| *3.6* | Alteration of the old aggregator generation processes in the Profitability module, to only integrate the creation of non-aggregated kpms tables. |
| *3.7* | Added stage field in aggregator tables |
| *3.8* | Integration in the aggregator generation process for obtaining the Stage field, coming from the universe table. |
| *3.9* | Analysis on how to obtain the stage, since the universe table does not cover all records |
| *3.10* | Unit Tests - Phase 3 (New execution of the process, ensuring the construction of the aggregator tables and then the execution of the Risks Module) |
| ***4*** | **Phase 4 – Creation of the Macro Excell** |
| *4.1* | Macro excel creation for Risk Module |
| *4.2* | Testing - Phase 4 (Running the Risks Mod completely, applying all rules, and building Excel) |
| ***5*** | **Phase 5 – Specify the Fields** |
| *5.1* | APRs - Continuation of the work already carried out in defining the universe of the various APRs |
| *5.2* | Line ID of Risk Module - Definition of the formulas to be applied to each line of the Risk Module |

Table 3.1 Tasks of the Project

This plan offers a comprehensive overview of the main steps involved in achieving the desired goal, which is to create a report containing all the necessary information. The work plan consists of five phases, each summarizing the key developments within them.

The first phase focuses on migrating from an on-premises system to the cloud while ensuring the preservation of the existing process structure. Moving on to the second phase, the process is adjusted to integrate the new risk module process rules. Once this adaptation is complete, the third phase involves developing the aggregators new process and adding the crucial field stage, which is essential to the new module process.

In the fourth phase, a new macro excel is created, providing the required structure for reporting. Lastly, in the fifth phase, the necessary formulas are defined for each line representing the modules universe, thereby concluding the analytical phase.

During the conclusive phase, the final product is presented, highlighting the main changes within the company and the contributions made from both a scientific perspective and to the client. The limitations and challenges encountered throughout the project are also addressed. Furthermore, this phase allows for an analysis of potential improvements and additional features that can enhance the final product, catering to the client's needs.

## 4. APPLICATIONS AND TECHNOLOGIES

In today's rapidly evolving digital landscape, organizations are constantly challenged with the task of efficiently managing massive volumes of data to drive insights and informed decision-making.

This chapter explores applications and technologies that have emerged as indispensable allies in tackling this data management problem. We explore an array of cutting-edge tools and techniques, spanning from robust cloud computing platforms to sophisticated data warehousing frameworks and versatile programming languages, all geared towards empowering organizations to harness the full potential of their data assets for enhanced business outcomes.

### 4.1. AZURE DATABRICKS

Databricks is a cloud-based platform designed for big data processing and analytics, enabling organizations to manage, process, and analyze extensive data volumes efficiently and cost-effectively. Databricks Workspace, Databricks SQL, and Databricks MLflow are key technologies that provide organizations with a comprehensive suite of services for data management and analysis.

Extensive research has explored the advantages of using Databricks in data management and analysis. For instance, a case study conducted by Smith (2020) demonstrated how Databricks Workspace enhanced data collaboration and project management capabilities for a large financial services company, leading to more informed business decisions. Additionally, Johnson (2019) highlighted the reduced time and effort required for data integration using Databricks SQL, which effectively processed and queried data from multiple sources.

Moreover, Brown (2018) investigated the use of Databricks flow for building and scaling machine learning models. The study emphasized the platform's ability to track, version, and share machine learning models, making it invaluable for organizations seeking to deploy complex machine learning workloads.

In conclusion, Databricks offers a comprehensive suite of services for effective data management and analysis. Its robustness, scalability, and cost-effectiveness make it an appealing choice for organizations looking to drive data-driven initiatives. By leveraging Databricks, organizations can efficiently manage, process, and analyze large data volumes, enabling them to make informed decisions and achieve their business objectives.

## 4.2. HIVE

Hive is an open-source data warehousing framework that facilitates data summarization, querying, and analysis. Originally developed by Facebook and now part of the Apache Software Foundation, Hive provides a SQL-like interface for performing operations on large datasets stored in the Hadoop Distributed File System (HDFS) and other storage systems such as Apache HBase. Hive supports custom user-defined functions (UDFs) and enables the analysis of structured and semi-structured data.

Wang et al. (2012) conducted a study comparing the performance of Hive with traditional relational databases for online analytical processing (OLAP) queries. Their findings revealed that Hive delivered competitive query performance for large-scale datasets while offering scalability and fault tolerance.

Zheng et al. (2013) focused on optimizing Hive performance for large-scale data analysis. Their proposed cost-based query optimization framework for Hive, which considered both data and query characteristics, resulted in improved query performance. This framework rendered Hive more suitable for large-scale data processing.

Furthermore, Hive has been integrated with other big data technologies such as Spark and Presto. Borthakur et al. (2015) detailed the integration of Hive with Spark, enabling faster query processing by offloading computations to Spark. The authors demonstrated that the combination of Hive and Spark offers a powerful and scalable platform for big data analysis.

To summarize, Hive is a well-established and extensively studied data warehousing framework, providing a SQL-like interface for data analysis and integration with other big data technologies. Its performance and scalability make it an attractive option for organizations seeking to process and analyze large datasets.

## 4.3. PYTHON

Python is a widely used high-level programming language renowned for its data manipulation and analysis capabilities. Its simplicity, ease of use, and availability of numerous data analysis libraries have made it a popular choice among data scientists and engineers. One prominent library for data analysis in Python is Pandas, which offers fast and flexible data structures and analysis tools.

McKinney (2010) conducted a study focusing on the use of Pandas for data analysis and manipulation. The study showcased how Pandas facilitates common data manipulation tasks

such as filtering, aggregating, grouping, and merging datasets. The findings emphasized that Pandas provides a convenient and efficient approach to data manipulation in Python.

Apache Spark is another crucial technology for data processing and analysis. It is an open-source distributed computing framework that offers a fast and flexible platform for big data processing. Spark supports Python, making it an appealing choice for data scientists and engineers already familiar with the language.

Zaharia et al. (2010) delved into the use of Spark for large-scale data processing. The authors discussed Spark's architecture and its application in distributed computing tasks like data processing and machine learning on extensive datasets. The study showcased Spark's high-level API for data processing and its ability to process data significantly faster than traditional MapReduce-based systems.

To execute Spark applications in a cluster, users can utilize the "spark-submit" command bundled with the Spark distribution. This command enables users to submit Spark applications for execution in a cluster, facilitating parallel processing of large datasets.

In summary, Python and Spark are powerful technologies for data manipulation and analysis. Pandas, a Python library, offers a convenient and efficient means of manipulating data. Spark provides a fast and flexible platform for large-scale data processing. The combination of Python and Spark presents a robust and user-friendly platform for data analysis and processing.

## 5. PROJECT

This project focuses on automating a reporting process of risk-related Key Performance Measures (KPMs) in a bank. Currently, the reporting is done manually in Excel, which presents various challenges such as data collection, transformation, and field calculations.

The project involves restructuring the profitability module and implementing a streamlined Extract, Transform, Load (ETL) process. This restructuring is driven by significant changes within the module and its associated tables, to simplify the process and increase efficiency. By validating tables and minimizing redundant reconstruction, the process can be executed more effectively and with reduced resource consumption.

The workflow for the risk module process is presented through visual representations, showcasing the essential steps required to obtain the necessary information. The process varies depending on whether it is a monthly process or a redeployed process, and it incorporates validations, user-defined parameters, and table constructions specific to each type.

### 5.1. ORGANIZATION CONTEXT

Neyond, founded in Portugal 2004, is a company that provides various services, one of them being IT Consulting that focuses its activity on the provision of strategic and operational consulting services in the area of technologies. As a service provider, Neyond works on projects that can be for the company itself or more often for other companies.

Neyond was established in 2004 with the aim of providing business and technology consulting services and strategic outsourcing in the areas of accounting, finance, tax, human resources and business support processes. It has more than 470 highly qualified professionals in its staff with offices in Lisbon, Porto and Madrid. It integrates a network of international partners "Process Solutions Network", having a strong presence at European level in a total of 34 countries (Neyond. 2015).

Neyond offers four different types of services:

- Business Process Outsourcing (Expert Outsourcing, Functional Outsourcing, Full Outsourcing): Activity focused on the provision of services of a continuous nature and strategic outsourcing of financial and administrative functions (Neyond. 2015);
- Business Consulting (Strategy & Market Services, Operations & Efficiency Services, Finance & Analytics Services, Risk & Regulation Services): Activity focused on the provision of strategic, operational and financial consulting services (Neyond. 2015);
- IT Consulting (IT Strategy & Governance Services, Digital & Mobility Strategy Services, Solution & Mobile Apps Development Services, Analytics & Business Insights

Services): Activity focused on the provision of strategic and operational consulting services in the area of technologies (Neyond. 2015);

- IT Outsourcing (Expert Services, Application Support & Management Services, Quality Assurance Services, Content & Information Management): Activity focused on the provision of services of a continuous nature and outsourcing in information technologies (Neyond. 2015).

What distinguishes Neyond from other companies is the global offer in terms of services like full stack, that is, what we conceptualize and define, implement, maintain, and operate. Deep knowledge of the sectors in which it operates resulting from their track record of clients and projects. The presence of a lot of senior consultants and specialists in the area working day by day near the client brings an effective presence into the market.

So, the choice to do a project in this corporation became obvious, due to the fact that this company could offer me a valuable experience as well as a high education in the domain, that is necessary to obtain the master's degree. Neyond offers a start in the job market and a preparation to the challenges that it brings, as it has a position of interest in the market, it enables the work for renowned companies.

This chapter will describe in more depth the process that has been implemented for one of their clients an international bank, as well as the changes that will occur by applying the risk module to the profitability process workflow, complementing with some background on Big Data.

## 5.2.    PROFITABILITY PROCESS

In regards of creating a context, the bank already had a profitability process implemented that resulted from the need to aggregate all the various types of information and applications used by the client. This process was created due to the necessity of a repository containing all the information as well as the necessity of querying to the operation level.

In order to understand in more depth this project background we need to understand its main goal, that is the creation of an organized and non-redundant tabular base with Management Control information at the most granular level (with sufficient historical depth for the purposes of future Pre-Provision Net Revenue or other models) as well as for all pre-defined dimensions.

In this process it was necessary to take into consideration the origin of the data, due to the fact that the records can be automatic, manual (examples are types of additions or required manual adjustments), from commercial areas or even from corporate activities. Regarding this, it is relevant to have the total universe of records ensuring its uniform treatment. When

treating a redeployment process it was necessary to consider the current terms and keep monthly photographs of each redeployment.

Allowing the total and exclusive production of the financial and risk modules of the Book of Business producing them exclusively from this data source.

## 5.2.1. Concepts

This solution will make it possible to respond to a series of initiatives related to the Bank's capital area, designated as Capital Enhancing.

All KPM used by Management Control were identified:

- Commercial Banking

- Majority Bank

- Corporate Activities

The information was collected with the level of granularity existing in one of the applications. It is possible to explore at the highest level of detail and aggregated information. This information was also generated through aggregations and calculation formulas for each KPM.

The loading of adjustments by users is done via a standard format that was defined together with the different Areas as well as the rules that were used in the incorporation of adjustments in existing data.

This solution was developed taking into consideration some Fundamental Principles:

- The presented solution aims to define the necessary structures to store the Profitability information, such as the incorporation of manual information. It also aims to define the Aggregation rules so that the existing information is sufficient to respond to the needs of existing Reports or to exist in the future and the respective periodicity of generation of that information

- Profitability information (KPM) must be stored in a single structure with standardization of information layouts regardless of the Area or type of Adjustment

- There is no redundancy of information

- Aggregation does not create or calculate new KPM

- From the previous point, it is corollary that the entire KPM will always be represented at a granular level.

## 5.2.2. Main Structure

For a better understanding of this process, it is relevant to know the sources of information for this procedure to work. This information is processed, and is regulated by the filters and rules implemented, and described further in this report.

The system will receive information from, various sources such as:

- Application
- Files uploaded by users
- Parameterizations
    - Filters
    - Rules

### 5.2.2.1.    Automatic Application information

The information coming from the application will be distributed, initially, by 5 tables:

- Table of Contracts
- Table of Contracts Data Redeployed
- Table of Override Contracts
- KPM Table
- KPM Commissions Table

The information via file will be stored in the Settings partitions. The data in this table will be in the format of Override, KPM or Commission KPM. The contracts contained in the application will be reflected in the key translator, thus existing the applications vs master relationship. There are several applications that send contract information to other tables in the database, so in order to be able to relate all the information and obtain the necessary information for reporting, there are tables for translating contract keys from the various existing applications.

The remaining tables in the schema will be populated by automatic processes, considering information that is loaded manually. Aggregate KPM and CORP Aggregate KPM Table – these

tables will be automatically generated every week (last day of each week) and monthly (last day of each month). The User will have the possibility to generate new aggregators after uploading manual Adjustments/KPM/Overrides files. The aggregated KPM table will have information aggregated by local concepts (KPM AGGREGATE) and the KPM table aggregated corporate concepts (KPM AGGREGATE CORP) will have aggregated information by Spain concepts contained in the Book of Businesses. The automatic information from the application will not be manipulated by the system, as there will only be a calculation of the KPM.

The granular information of Contracts and KPM will be saved with the following deadlines:

- Information of the last 7 working days – when generating the information of the day, the information of the previous 8th day is eliminated (except Friday and the last day of each month)
- Closing information for the last 4 weeks – when processing on Friday, the previous 5th Friday is deleted (except if it is the end of the month)
- End of month information is not deleted – granular information for each month end will not be deleted for a minimum of 15 years

| DATA | HISTORICAL DATA |
|---|---|
| CONTRACTS | HISTORY CONTRACTS |
| CONTRACTS GIVEN REDEPLOYMENT | CONTRACTS GIVEN RED. HISTORY |
| OVERRIDE CONTRACT | OVERRIDE CONTRACT HISTORY |
| KPM TABLE | HISTORY KPM TABLE |
| KPM COMMISSION TABLE | KPM COMMISSIONS TABLE HISTORY |
| KPM AGGREGATE | KPM AGGREGATE |
| KPM AGGREGATE CORP | KPM AGGREGATE CORP |

Table 5.1 Data Structure

**CONTRACTS TABLE:** In this table the 'static' Contract data will be stored.

**CONTRACTS TABLE (REDEPLOYMENT DATA):** This table will store data that can be changed and that result in the redeployment of results. This table will always keep the data provided by the inputs and there will be differentiation between Automatic registers and Manual registers.

**OVERRIDE CONTRACTS TABLE (REDEPLOYMENT DATA):** This table will consider changes to some fields that are made via file by Users, that is, this table will have the Contract data to

be used in the aggregation process. This table will have History management, since in the redeployment process the aggregators will be recalculated for the last 36 months at most (two previous years plus the months of the current year that have already passed). Therefore, at each redeployed process, the last version of this table will be transferred to the history so that granular information can be consulted in a view prior to the redeployed process. In the redeployment process, this table will be generated for the last 2 years and for the past months of the current year.

**KPM TABLE:** In this table will be the KPM that are not referring to Commissions. The KPM will be in columns having rows for each date.

**KPM COMMISSION TABLE:** This table will have the KPM referring to commissions. They are placed in a separate structure since the METACOMMISSION, one of the variables in this table is liable to be redeployed, and therefore the values of these KPM may vary.

**HISTORY TABLE:** These tables have the history of the corresponding tables before the redeployed is processed. It will be the same as the table from which it originates, with two additional fields in order to identify which redeployment the history refers to.

### 5.2.2.2. Manual Information

The Manual Information - Manual Adjustments/KPM and Overrides comes from a layout where it is possible to enter manual information in the Profitability module. The Adjustments/KPM Manuals and Overrides files are loaded in the Front to be developed for this purpose.

All information entered via file will always be translated as a new contract entry, even when referring to an existing contract, Manual Adjustments/KPM will always be translated into granular information.

Each Adjustment has the Adjustment Reference Date given as front input (by default this will be a date before the due date – typically this will be the end of the previous month). The adjustments will take effect on that date and for cases of Overrides file will persist until the Expiration Date indicated on the front. If an Override validity date is not indicated, then the validity will only be for the informed reference date and will only be applied to the aggregation of the information of the introduced reference date.

Manual Adjustments/KPM will not have persistence.

The User will be able to disable previously loaded Settings files when dealing with files of the incremental type. If the input files are of the total typology, then loading a new file will overwrite the previously loaded information for that date.

The files that have this information are:

- Overrides File: The User can load Overrides for previous reporting periods. The Reference Date will be input at the front, and it will be for that date that the registrations will take effect.
- Adjustment File/KPM Manuals: The User will be able to load Manual Adjustments/KPM and Overrides for previous reporting periods. The Reference Date is input at the front
- Files with GCB Contracts, Corporate Activities and Costs: The Management Control Areas that have portfolios that are not represented in the application provide their information in the previously defined layout referring to Adjustments and Manual KPM.

### 5.2.2.3. Filters

The information to be considered is filtered. The filter rules can be defined in the rules parameterization screen that was detailed in the Universe. Access to the parameterization of this filter is restricted to just a few Users/Areas since they need a user that has access to the filters in order to define the rules and so the application filters are managed by different areas.

### 5.2.2.4. Rules

The system allows the creation of rules so that the User can define in terms of Dimensions, Business Area, Product and Commission to be assigned to a set of Records. These rules can be set / Changed on the Front (application) screen. On this screen, the User selects the type of rules he wants to handle. In this case they will be application Rules – Profitability.

Afterwards, they must choose the dimension that they want to define the value and insert that value. It then adds the dimensions to check for the rule and the values each of those dimensions can take. Note here that each rule will have a limit of five dimensions to be checked and that when a rule validates more than one dimension, the rule will be made if all dimensions fulfill the validation.

## 5.2.3. KPM

The KPM´s originated in order to develop the final report on one hand some of them already come from the application in the final formula, on the other hand some of the KPM´s need to be calculated in the process, these KPM's use the already calculated KPM's with formulas.

The output has a total of two hundred and fifteen KPM's.

As an example, here is a short description of two of the KPM´s reported by the process.

- Equity: Equity can be defined as the amount of money the owner of an asset would be paid after selling it and any debts associated with the asset were paid off.
- Taxes: Taxes are mandatory contributions levied on individuals or corporations by a government entity.

| KPM | DETAILED CALCULATION |
|---|---|
| **Equity** | Monthly SMV (Loan - Mortgage Loan - FC Loan) * 8% + SMV (Mortgage Loan + FC Loan) * 4% |
| **Taxes** | (Monthly MF (Business Volume + Previous Months MF (Business Volume) + Interest Arrears in the Month (products) + Interest Carrying Portfolio Overdue Credit (Products) + Other MF + ROF + Direct Personnel + Direct General Expenses) + Charged Personnel Costs + Charged General Expenses + Amortizations + Return on Capital + Guarantee Fund Deposits + Net Commissions + Flow Provisions Impairment) *30%/31% (depending on the Year) *-1 |

Table 5.2 Examples of KPM's

### 5.2.4. Aggregator

The process will aggregate, the information processed each day, after loading the Adjustments and processing them, the aggregator is processed again and now considers the loaded adjustments and Overrides, generating new aggregated information.

The aggregation is done considering the granular information of each period, including the Adjustments/Overrides that have been loaded. Adjustments loaded for a given reference date do not have persistence, that is, they will only be considered in the reference date aggregator and will not be considered in subsequent periods.

The retention periods of the structures will be the same as those referred to in the Contracts and KPM structures of the automatic information. The latest version of the aggregator will be used to generate automatic reports. Granular information that justifies a given aggregation will be given according to whether a given Adjustment was considered in that aggregation.

### 5.2.5. Redeployment

When the redeployment process occurs, Data Lake will recalculate the Month Close of the Aggregator of the last 2 years and past months of the current year. This process will take into account the latest version of the Contracts table (redeployment fields) and respect the frequency of redeployment, which is currently quarterly.

The tables that are passed to history will be registered with the redeployment code that identifies the period to which it refers. This code will be the same in all tables so that they can be cross-referenced to previous redeployment periods.

The Override Contracts table will be stored in history (with redeployment period code) for the target months of redeployment and those months will be reconstructed considering the last version of each contract. This means that if a contract currently has a different segment, the segment will be updated to the current value in the previous 2 years and in the elapsed months of the current year.

Afterwards, the Override Contracts table are replaced for the months affected by the redeployment, that is, the new values of the fields that can be reallocated from the automatic information will be considered. If there are manual overrides in the months in which you are reassigning, these Overrides will be considered in the construction of the Override Contracts table. This means that the validations described in the adjustment's files will be done again and the construction of the Override Contracts table will follow the same

rule – that is, manual information has priority over automatic information even in redeployments.

The values of the fields that can be redeployed will be updated taking into consideration the current value in the Customer Data table. This data will be updated in the Contracts table (redeployment data) automatically and manually. In the case of manuals, it will only update those that entered without filling in the values.

The redeployment will not alter information loaded via Adjustments/Override, that is, if there is an Override where the Segment of a contract is changed, that Override will prevail on the reference date for which it was inserted, regardless of the Contract in question having a change of Segment via redeployment. The same applies to Adjustments at Contract level where the segment or other field that can be redeployed has been filled in.

Note that the Overrides validations are repeated, which means that if there is an Override by the Management Center without informing the Segment, then the Segment will be updated to the new value (situation in which the Segment changed) and it will remain the Override of the Management Center (with the possibility that now the Management Center will already be the one that had been changed manually at the time).

Analogously, the validations of the Area will be carried out again. That is, the automatic redeployment will be effective in adjustments at contract level that affect KPM (for redeployment variables that are not adjusted), and in adjustments at customer level that affect KPM (for redeployment variables that are not adjusted).

Variables changed in adjustments will not be effective for redeployment, whether at contract, customer, or segment/center level.

The calculation of the Aggregator will be made for all the target months of redeployment. In this calculation, the existing Adjustments in force for each of the closing months will be considered. The new versions generated will be the ones that will be in effect for future reports, keeping the previous version in history.

In the situation of having to make an Adjustment to a previous month that has already been closed and if you want that adjustment to be reflected in the accruals of the following month and months or if you make adjustments in several previous months and you want everything to be reprocessed, the User should be able to request redeployment of every month from a certain date. This reprocessing will be done at the end of the day and will run from the earliest date of reprocessing requests from a given date forward.

## 5.3.    Risk Module Process Project

The inception of this project stems from the need to report risk-related Key Performance Measures (KPMs) to the head office on a monthly basis, with a reporting deadline of D+9. Currently, this report is manually prepared in Excel, which entails various challenges such as data collection, transformation, and calculations.

The primary objective of the bank is to establish an automated process that enhances the quality of the reported information. Additionally, the aim is to incorporate additional data that is currently excluded due to the limitations of the manual process, thus improving the comprehensiveness of the report.

In essence, the goal of this project is to automate a mandatory reporting process for the main corporation, eliminating manual efforts and ensuring consistent, accurate, and comprehensive reporting of risk-related KPMs. By transitioning to an automated approach, the bank can streamline operations, enhance data quality, and provide timely and comprehensive insights to the head office.

### 5.3.1.1.    Restructure of the Profitability Process

In order to implement the new process, it was determined that certain existing processes within the profitability module required restructuring. This decision stemmed from the significant changes that occurred within the module and its associated tables. The goal was to simplify the process and make it more efficient.

Considering that the construction of aggregator tables did not need to be repeated each time the profitability process was run, it became essential to establish a flow that showcased the reconstruction and implementation of the process. To avoid unnecessary repetition and reconstruction of well-constructed tables, a simple validation step was introduced. This step ensures that the tables are valid before proceeding, preventing redundant re-running and reconstruction.

Overall, the restructuring of the profitability module and the introduction of a streamlined E.T.L. process have resulted in a more simplified and efficient implementation. By validating the tables and minimizing unnecessary reconstruction, the process can be executed with enhanced effectiveness and reduced resource consumption.
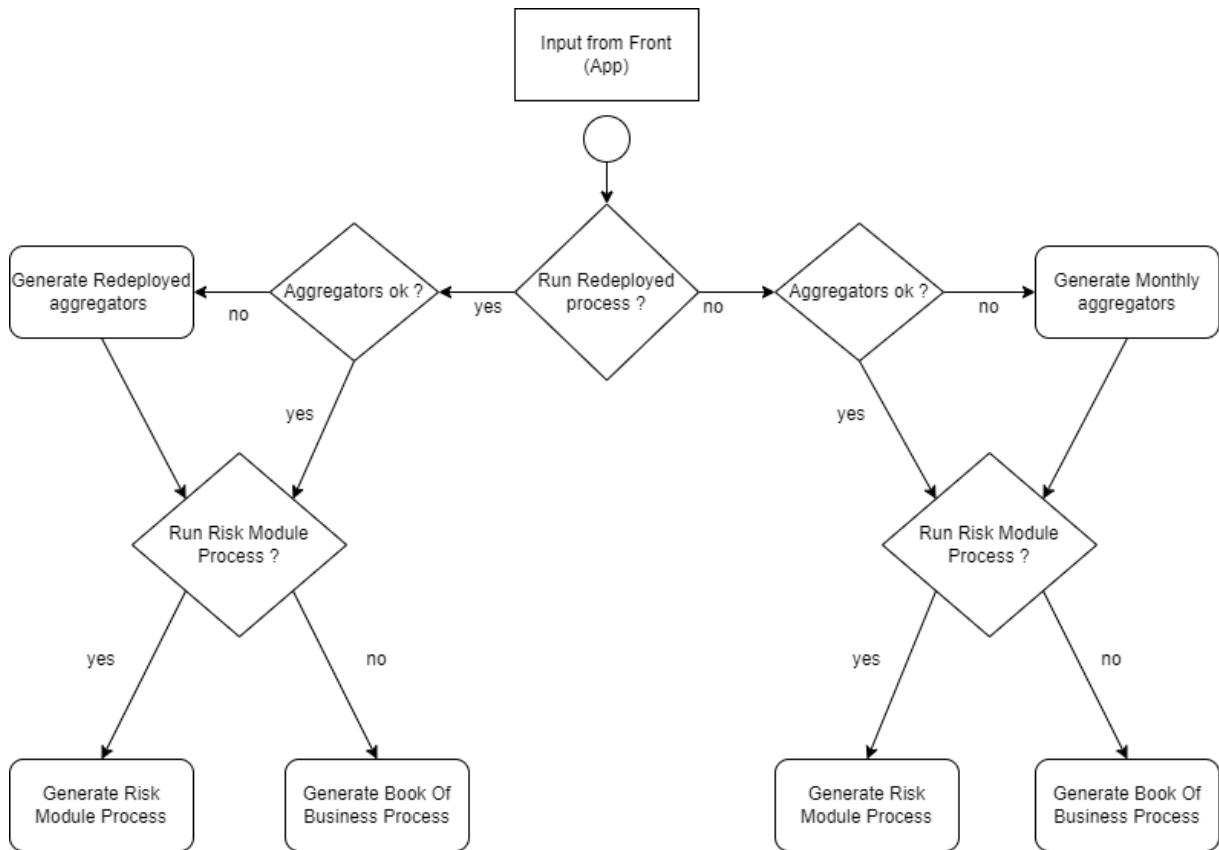
Figure 5.1 Workflow Risk Module Process

The implemented alteration has transformed the process into a more efficient one, featuring a scalable structure that can be applied to other modules with dependencies on these tables. Within the flow (Figure 5.1), another module called Book of Business Process can be observed, which also relies on these tables as dependencies. While the structure of the Book of Business Process is similar to the Risk module process, the applied rules differ due to the variation in the reported information.

Through this modification, efficiency is assured. The inclusion of a table validation step prevents unnecessary execution, minimizing server space utilization and allowing resources to be allocated to other processes. Additionally, the modified process runs more swiftly, ensuring a rapid response and timely delivery of the output file to the client.

Input from Front (App) :
- procID
- codProc
- rent_type
- propParam

Generate Redeployed aggregators — no — Aggregators ok ? — Redeployed Process (id_report,procID, propParams,rent_type) — yes — Run Redeployed process ? — no — Monthly Process (id_report,procID) — Aggregators ok ? — no — Generate Monthly aggregators

Aggregators ok ? — yes — Run Risk Module Process ?

Run Risk Module Process ? — yes — Generate Risk Module Process

Run Risk Module Process ? — no — Generate Book Of Business Process

Aggregators ok ? — yes — Run Risk Module Process ?

Run Risk Module Process ? — yes — Generate Risk Module Process

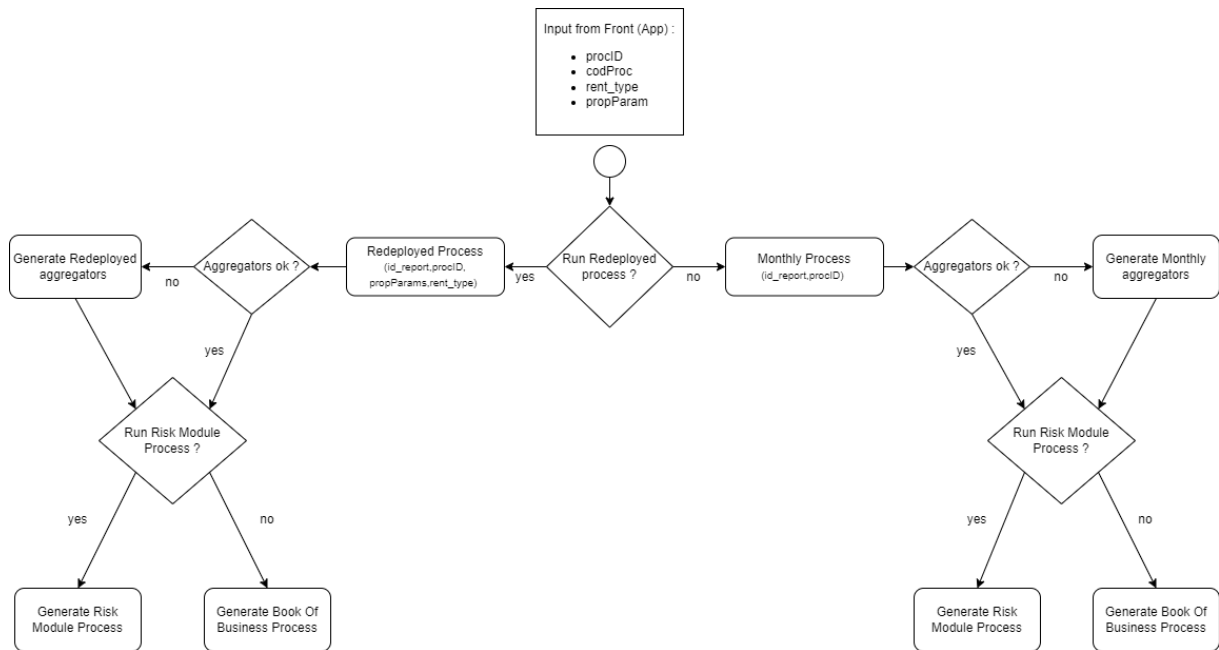Run Risk Module Process ? — no — Generate Book Of Business Process

Figure 5.2 Risk Module Process

This detailed workflow (Figure 5.2) outlines the essential steps required to obtain the necessary information for reporting to headquarters. The process varies depending on whether it is a monthly process or a redeployed process, as the construction of tables differs accordingly. Therefore, the following validations are crucial: the construction of aggregator tables and their validity, as well as determining whether the client desires a monthly or redeployed process for the module.

This workflow (Figure 5.2) incorporates the previously mentioned validation, which verifies the construction and validity of the aggregator tables. Additionally, it introduces the next validation, which involves a parameter that enables the process to determine whether it should run as a monthly or redeployed process. In the front application, the user selects the desired process type and specifies the date for processing or reprocessing the data.

By incorporating these validations and user-defined parameters, the workflow ensures a customized and efficient process tailored to the client's needs. This approach allows for flexibility in generating the required reports and facilitates the timely processing of accurate and relevant information.
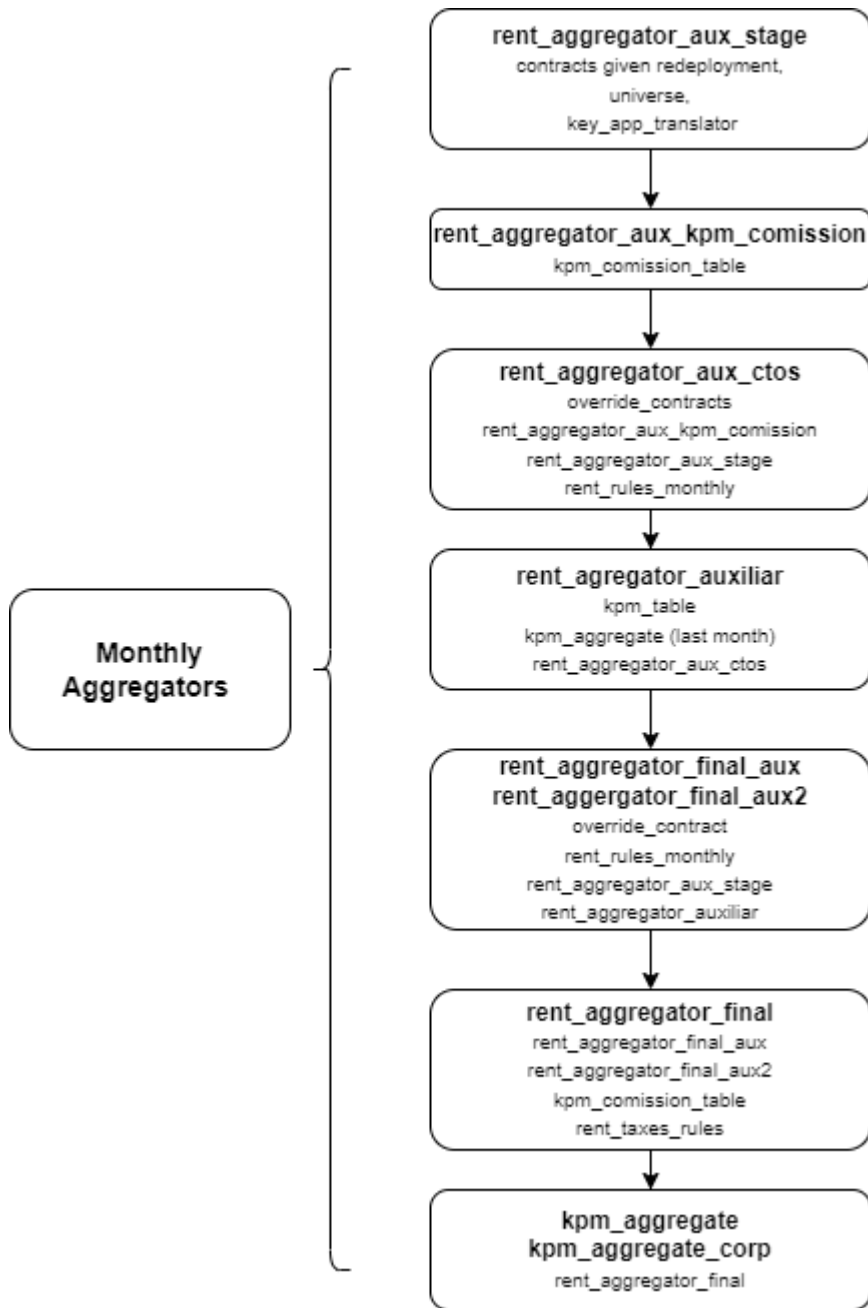
Figure 5.3 Monthly Aggregators

The monthly process as demonstrated in Figure 5.3, involves a complex construction of the main tables necessary for the reported metrics. This process, known as Extract, Transform, Load (E.T.L.), utilizes three primary tables: the Contracts Table (Redeployment Data), Override Contracts Table (Redeployment Data), and KPM Table. Due to the intricate formulas and calculations involved in the Key Performance Measures (KPMs), several steps are required.

The primary focus during the construction of this process is to ensure the accuracy of the data used and the proper execution of joins between the tables, considering all the necessary keys. The first step involves a table that calculates the stage, which will be discussed in detail in subsequent sections of this paper.

Next, the process aggregates the KPM Commission Table and performs calculations on the KPMs. Subsequently, the Override Contracts Table, the already aggregated auxiliary table based on the Commission Table, the stage table, and the table containing the rules for monthly process execution are all integrated.

To obtain all the KPMs, the process incorporates the KPM Table and the KPM Aggregate for the previous month's information. Using the aforementioned auxiliary table, further aggregations and calculations are performed. The KPM Commission Table is then added to include the final KPMs required for the process.

Finally, the output tables, namely KPM Aggregate and KPM Aggregate Corp, are populated with information from the auxiliary table that underwent the final aggregation.

Overall, this process ensures the proper construction of the main tables and the accurate calculation of KPMs, providing reliable metrics for reporting purposes. The step-by-step approach guarantees the integrity and validity of the data used in the process, resulting in meaningful and valuable insights.
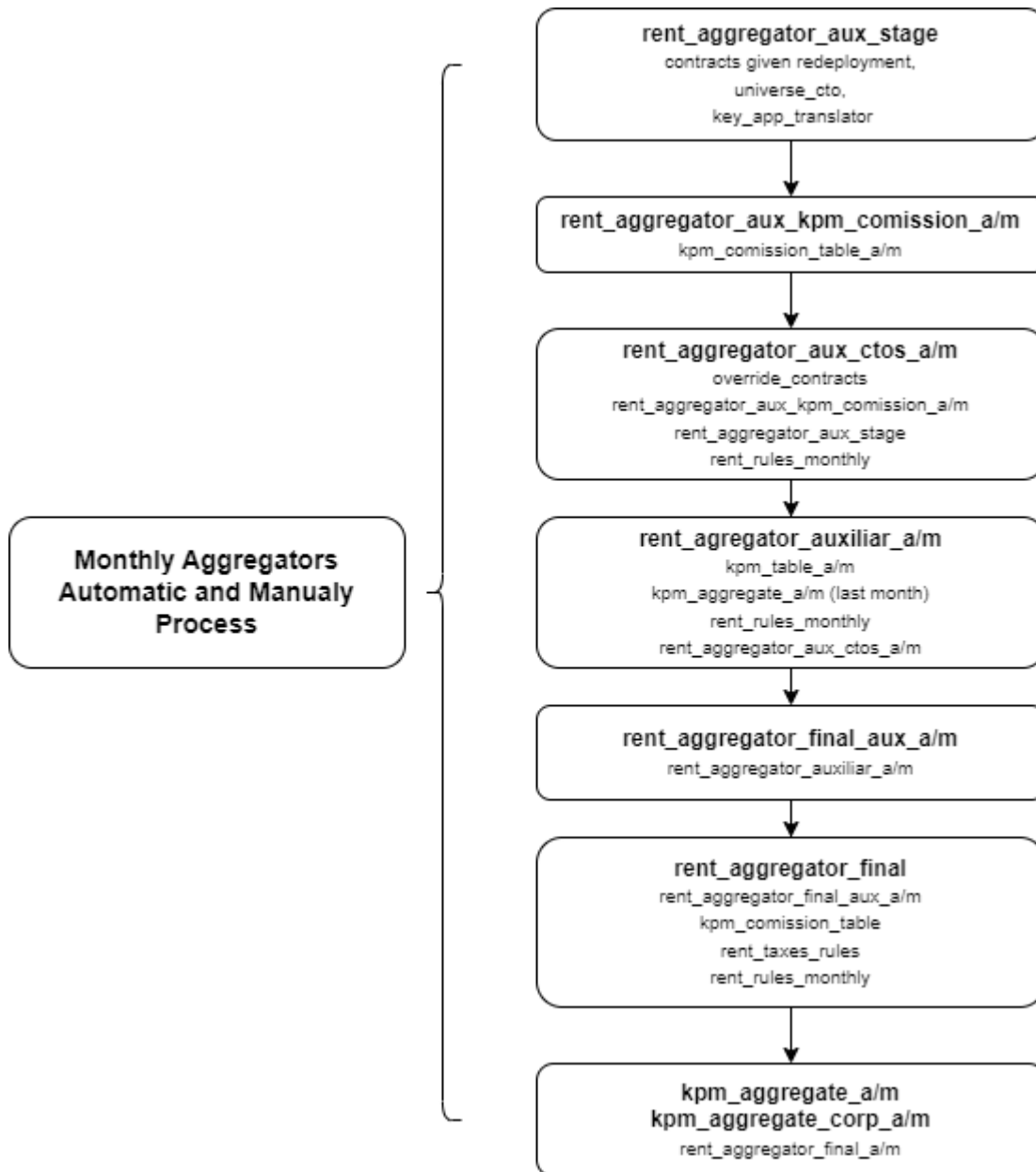
Figure 5.4 Automatic and Manual Monthly Aggregators

The Mensal process is divided into total aggregators, mensal aggregators, and automatic aggregators, which simplifies the visualization of information. The total aggregators consume all available information, while the automatic or manual aggregators consume automatic or manual information, respectively. Automatic information is derived from the process itself, while manual information includes adjustments and manually added data.

While similar to the monthly process, this particular process differs in terms of the tables used, which depend on the type of data involved as shown in Figure 5.4. For automated

processes, tables containing automated information are utilized, whereas manual processes require tables with manual information. However, it is important to note that in order to maintain data consistency, unlike the monthly process for the total aggregators process which incorporates the OVERRIDE CONTRACT table, this process cannot utilize that table due to inconsistent information. Therefore, it becomes necessary to incorporate the CONTRACTS GIVEN REDEPLOYMENT table to ensure the accuracy and coherence of the data.

Additionally, the process requires tables that contain the rules for calculating key factors in the output table. These tables play a vital role in ensuring accurate calculations and producing reliable results. By incorporating the appropriate tables and considering the specific data types, the Mensal process maintains coherence and accuracy throughout its execution.
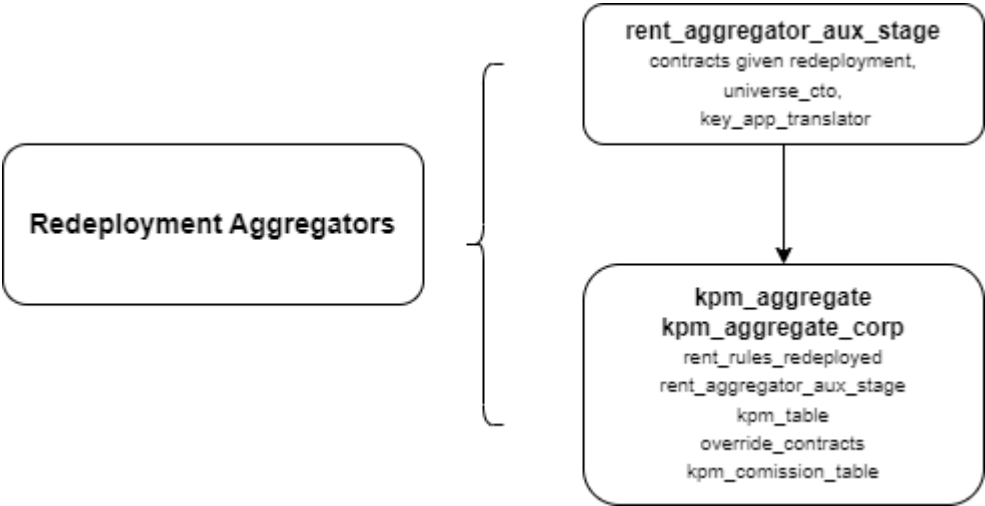


Figure 5.5 Redeployment Process

In a redeployment process (Figure 5.5) within the risk module, the metrics are already calculated, making the construction of final tables a simpler task. In this process, the primary objective is to obtain the stage information from the universal table and combine it with other relevant data to recalculate the tables for specific dates.

Since the redeployment process involves pre-calculated KPMs, the key steps include aggregating the CONTRACTS GIVEN REDEPLOYMENT table, the universe table, and the translator to derive the stage information. In the final step, all the pertinent tables, including rules tables, the OVERRIDE CONTRACT table, the KPM TABLE, the KPM COMMISSION TABLE, and the previously generated stage table, are aggregated.

This straightforward process ensures that any overrides made by the client, which may have affected the data, are taken into account. By redeploying the data, the quality and coherence of the information are safeguarded, thereby ensuring accurate and reliable results.

### 5.3.1.2.    Inclusion of the field Stage to the Process

The incorporation of the stage field is vital to the construction of the risk module process as it provides crucial insights into the types of risks and the corresponding KPMs associated with them. This field is derived from a comprehensive table that represents the entire universe of contracts within the company.

To utilize this information effectively, a translator table is required. This translator table serves as a bridge between the universal table and the process's reporting and information retrieval needs. In the Data Lake, where all the relevant information is stored, a dedicated table is established to facilitate the conversion and extraction of the stage information.

By leveraging the translator table and the corresponding Data Lake table, the risk module process can accurately obtain the stage field, enabling comprehensive risk analysis and reporting. The inclusion of this field enhances the overall effectiveness and reliability of the process, ensuring that the right KPMs are associated with the appropriate risk types.
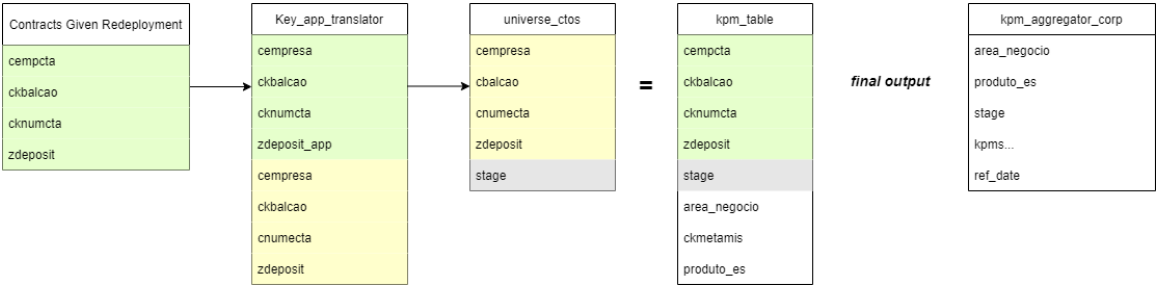


Figure 5.6 Obtaining Stage

The visual representation (Figure 5.6) simplifies the understanding of a translator and its role in obtaining the stage information for the risk module process. The picture highlights the

primary table, CONTRACTS GIVEN REDEPLOYMENT, which comprises four essential components. In order to extract the stage, this information needs to be cross-referenced with the universe table, and further cross-referenced with its own set of four key components.

To achieve this, the translator table plays a crucial role. It enables the seamless integration of the three tables, allowing the extraction of the desired stage information. By leveraging the translator table, the process successfully combines the relevant data from CONTRACTS GIVEN REDEPLOYMENT, the universe table, and its corresponding key components. Ultimately, the goal of obtaining the stage information is reached through the utilization of this translator table and the subsequent joining of the tables.

### 5.3.1.3.    Future work

As described throughout this paper the project development had five phases with many activities that were purposed, the last two weren't finished, phases four and five, that contemplated the creation of the Macro Excel and the Specification the Fields.

Due to unexpected external factors, the project was unable to reach its intended conclusion within the initially projected timeline. As a result, the conclusion of the project will be deferred until the client determines the appropriate time to resume the development process. The creation of the Excel file, which relies on the output table generated in the risk module process, will be addressed in subsequent stages.

Despite the delays encountered, the project remains adaptable and responsive to the client's evolving requirements, ensuring that the final outcome aligns with their objectives.

# 6. CONCLUSIONS

The outcome of this project would have an impact on how the management control department reports to the client's headquarters, with a more specific spectrum of metrics that can then be analyzed and taken into consideration when making future decisions. This automation of the process can then be replicated to deliver other relevant data as well as placed into other departments that need their processes rethought.

One of the main focuses of the project besides delivery a good product to the client is to take into consideration its scalability since it can facilitate future client purposes and alterations to the product. Making a final product scalable will enable adjustments a lot faster and efficient, creating a more promising worker-client relationship since it will demonstrate how a good developer works, and how the product can be manipulated easily.

This project will in the future help the company build knowledge and facilitate learning in ways to report information regarding essential risk and profitability, key performance indicators and metrics that will help in a more successful business. As some companies are exploring digital transformation initiatives such as big data analytics and artificial intelligence to find the "magic bullet" which will unlock here to fore unrealized value (Williams, M. 2022).

From a scientific standpoint, this thesis holds significant value as a reference for others in their pursuit of scientific research, thereby fostering new opportunities for exploration. The remarkable advancement in computing power, data storage, and sensing technology has ushered in an era where vast amounts of data can be captured and analyzed. This has sparked increasing interest and led to the development of numerous studies focused on this subject matter (Wall, J., & Krummel, T., 2020).

## 6.1. CHALLENGES AND LESSONS LEARNED

In the process of undertaking this project, I have acquired invaluable knowledge and experienced profound growth that has greatly influenced my problem-solving strategies. The power of collaboration became evident as a formidable asset, demonstrating the advantages of collective thinking and harnessing collective knowledge.

Moreover, recognizing the inherent significance of comprehending the essence of data emerged as a focal factor in attaining substantial outcomes. Additionally, we confronted distinct obstacles concerning data performance and time allocation, which presented distinctive challenges throughout the project's development.

Throughout the course of this project, the most valuable lesson I have learned is the power of collaboration. I firmly believe that two minds working together yield better results than one. When faced with challenges and obstacles that seem overwhelming, it is often more beneficial to participate in open discussions and seek input from others. Even if someone may not possess a complete understanding of the project's intricacies, they can still ask crucial questions that shed light on potential issues and help pinpoint the root of the problem.

An additional significant lesson learned was the importance of comprehending the meaning of the data. It is crucial to understand how the data relates to other tables and universes, ensuring that it maintains coherence and contributes to the overall business knowledge. By grasping the contextual significance of the data, we can effectively establish connections and derive valuable insights for informed decision-making.

During the development of this project, we encountered various challenges, including data performance issues. Specifically, the creation of tables resulted in significant memory consumption. Therefore, it became crucial to develop efficient queries that minimized both memory usage and processing time. Striving for optimized data performance was essential to ensure smooth execution and maintain efficient resource utilization.

Another significant challenge I encountered during the project was time management. While there was no strict deadline for project completion, it proved difficult to estimate delivery dates to my manager. This challenge stemmed from my limited experience and unfamiliarity with potential problems that could arise.

In summary, the project journey has offered me invaluable insights and knowledge that will greatly influence my future accomplishments. The true power of collaboration has been firmly established, emphasizing the immense value of teamwork and diverse perspectives when addressing problems.

Furthermore, the experience of effectively managing time has highlighted the significance of experience and adaptability in accurately estimating and meeting project timelines.

These valuable lessons and challenges have not only contributed to the triumph of this project but have also encouraged personal growth and professional progress that will resonate in future projects.


## 6.2. FUTURE PERSPECTIVES


This significant experience has played a vital role in uncovering my true passions and providing clarity regarding my professional aspirations. Moving forward, my intention is to

continue to work in this field, actively seeking opportunities to expand my value, deepen my knowledge, and refine my skills.

As I embark on this journey toward my future occupation, I recognize that challenges and uncertainties will inevitably arise. However, armed with the lessons and resilience cultivated through this project, I approach these obstacles with optimism.

Beyond personal growth and professional ambitions, my vision for the future extends to making a positive impact on society. I aspire to collaborate with like-minded individuals, leveraging innovation and a deep sense of purpose to drive meaningful change.

In summary, this experience has not only provided clarity in terms of my interests and professional goals, but it has also ignited an unwavering determination to continuously evolve and excel.

## BIBLIOGRAPHY

Askar, M., Imam, S., & Prabhaker, P. R. (2009). Business metrics: A key to competitive advantage. Advances in Competitiveness Research, 17(1), 91–110.

Borthakur, D., Saha, B., Sen Sarma, J., Elkin, M., & Evans, R. (2015). Apache Hadoop goes real-time at Facebook. In Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data (pp. 1461-1466).

Brown, M. (2018). Building and scaling machine learning models with Databricks MLflow. Data Science Journal, 17(4), 201-212.

Chen, Y., & Chen, H. (2020). Big data analytics in financial management. Journal of Business Research, 117, 365-373.

Chiang, R. H., & Storey, V. C. (2016). Big data and customer segmentation. Journal of Business Research, 69(12), 5497-5504.

Court, D. (2015). Getting big impact from big data. McKinsey & Company. https://www.mckinsey.com/capabilities/mckinsey-digital/our-insights/getting-big-impact-from-big-data

Davenport, T. H., & Harris, J. G. (2007). Competing on Analytics: The New Science of Winning.

Elveny, M., Nasution, M. K. M., Zarlis, M., & Efendi, S. (2021). An advantage optimization for profiling business metrics competitive with robust nonparametric regression. Journal of Theoretical and Applied Information Technology, 99(1), 114–124.

Fayyad, U., Piatetsky-Shapiro, G., & Smyth, P. (2016). From data mining to knowledge discovery in databases. AI magazine, 17(3), 37.

Feldman, R., & Sanger, J. (2013). The Text Mining Handbook: Advanced Approaches in Analyzing Unstructured Data. Cambridge University Press.

Gandomi, A., & Haider, M. (2015). Beyond the hype: Big data concepts, methods, and analytics. International Journal of Information Management, 35(2), 137-144.

Hada, I. D., & Mihalcea, M. M. (2020). The importance of profitability indicators in assessing the financial performance of economic entities. The Annals of the University of Oradea, Economic Sciences, 219–228.

Hasan, M. M., Popp, J., & Oláh, J. (2020). Current landscape and influence of big data on finance. Journal of Big Data, 7(1). https://doi.org/10.1186/s40537-020-00291-z

Heikkilä, M., Bouwman, H., Heikkilä, J., Solaimani, S., & Janssen, W. (2016). Business model metrics: An open repository. Information Systems and E-Business Management, 14(2), 337–366. https://doi.org/10.1007

Housbane, S., Khoubila, A., Ajbal, K., Serhier, Z., Agoub, M., Battas, O., & Othmani, M. B. (2020). Monitoring mental healthcare services using business analytics. Healthcare Informatics Research, 26(2), 146–152. https://doi.org/10.4258/hir.2020.26.2.146

Hausberg, J. P., & Wied, D. (2015). Risk management in practice: a guide for decision makers. Springer.

Hsu, W. (2017). Big data analytics for emerging risks. Journal of Emerging Risks, 1(1), 1-14.

Johnson, K. (2019). "Data integration using Databricks SQL: A case study." In Proceedings of the International Conference on Big Data and Analytics, pp. 87-93.

Jorion, P. (2007). The relationship between risk and return.

KPMG. (2013). Unleashing the power of big data. KPMG International.

KPMG. (2015). Emerging risks in the 21st century: a KPMG report. KPMG.

KPMG. (2016). Big data analytics in risk management: opportunities and challenges. KPMG.

Lazer, D., Kennedy, R., King, G., & Vespignani, A. (2014). The parable of Google Flu: Traps in big data analysis. Science, 343, 1203-1205. https://doi.org/10.1126/science.1248506

Liu, B., Bu, D., Chen, L., & Zhang, Y. (2015). Big data: A survey. Mobile Networks and Applications, 20(2), 171-209.

Mashrur, A., Luo, W., Zaidi, N. A., & Robles-Kelly, A. (2020). Machine learning for financial risk management: A survey. IEEE Access. Institute of Electrical and Electronics Engineers Inc. https://doi.org/10.1109/ACCESS.2020.3036322

Mayer-Schönberger, V., & Cukier, K. (2013). Big data: A revolution that will transform how we live, work, and think. New York: Houghton Mifflin Harcourt.

McKinney, W. (2010). Data Structures for Statistical Computing in Python. In Proceedings of the 9th Python in Science Conference, pp. 51-56.

McNeil, A.J., Frey, R., Embrechts, P. (2015). Quantitative risk management: Concepts, techniques, and tools. Princeton University Press.

Neyond. (2015). Neyond - New Ways To Go Beyond. Neyond.pt. https://www.neyond.pt/en/

Oussous, A., Benjelloun, F. Z., Ait Lahcen, A., & Belfkih, S. (2018). Big Data technologies: A survey. Journal of King Saud University - Computer and Information Sciences, 29(2), 132-140. https://doi.org/10.1016/j.jksuci.2017.06.001

Panchyshyn, T., & Prokopovych-Pavlyuk, I. (2021). Modern business metrics for evaluation of efficiency of marketing measures. Odessa National University Herald. Economy, 26(1(86)), 22-25. https://doi.org/10.32782/2304-0920/1-86-22

Smith, J. (2020). The impact of Databricks Workspace on data collaboration and project management. Journal of Cloud Computing, 9(2), 123-135. https://doi.org/10.1186/s13677-020-00143-9

Sisco, C., & Chorn, B. (2009). Key Performance Indicators for Responsible Sourcing. BSR. https://www.bsr.org/.../BSR_Responsible_Sourcing_

Sivarajah, U., Kamal, M. M., Irani, Z., & Weerakkody, V. (2017). Critical analysis of Big Data challenges and analytical methods. Journal of Business Research, 70, 263-286. https://doi.org/10.1016/j.jbusres.2016.08.001

Solomon, D. C., & Muntean, M. (2012). Assessment of Financial Risk in Firm's Profitability Analysis. Economy Transdisciplinarity Cognition, 15(2), 58-67.

Steshenko, O., & Bondarenko, Y. (2021). Business analytics in financial risk management. Market Infrastructure, 55, 28-31. https://doi.org/10.32843/infrastruct55-28

Van Looy, A., & Shafagatova, A. (2016). Business process performance measurement: A structured literature review of indicators, measures, and metrics. SpringerPlus, 5(1), 1-14. https://doi.org/10.1186/s40064-016-3498-1

Wall, J., & Krummel, T. (2020). The digital surgeon: How big data, automation, and artificial intelligence will change surgical practice. Journal of Pediatric Surgery, 55, 47-50. https://doi.org/10.1016/j.jpedsurg.2019.09.008

Wang, Y., & Li, Y. (2017). Big data analytics: A literature review and a framework. Journal of Big Data, 4(1), 1-17.

Wang, Y., Xiuping, S., & Zhang, Q. (2021). Can fintech improve the efficiency of commercial banks? —An analysis based on big data. Research in International Business and Finance, 55, 101338. https://doi.org/10.1016/j.ribaf.2020.101338

Wang, Q., Liu, L., & Wu, J. (2018). Big data and supply chain management: A review and research agenda. Journal of Business Research, 85, 208-218.

Wang, X., Zhou, Y., Yin, H., & Yu, J. X. (2012). Hive-based online analytical processing on large data sets. Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data, 1-12.

Williams, M. (2022). Process automation platforms. In Integration and Optimization of Unit Operations (pp. 239–247). https://doi.org/10.1016/b978-0-12-823502-7.00024-4

Xie, M., Zhang, J., & Wang, H. (2018). Big data analytics for risk assessment in financial institutions. Journal of Big Data, 5(1), 1-16.

Zaharia, M., Chowdhury, M., Franklin, M. J., Shenker, S., & Stoica, I. (2010). Spark: Cluster computing with working sets. In Proceedings of the 2nd USENIX Conference on Hot Topics in Cloud Computing, 10-10.

Zheng, Y., Liu, Y., & Zhang, Y. (2013). Cost-based query optimization for Hive. Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data, 1117-1128.

Zikopoulos, P., & Eaton, C. (2012). Understanding big data: Analytics for enterprise class Hadoop and streaming data. McGraw-Hill Osborne Media.