

Data Mining in Actuarial Science: Reserves

Following the CRISP-DM methodology.

Juan Lara

December 5, 2023

Disclaimer

This document from Universidad Nacional de Colombia is for informational purposes only and is not intended as professional or financial advice. Efforts have been made to ensure the content’s accuracy, but no guarantees are made regarding its reliability or completeness. Universidad Nacional de Colombia is not liable for any decisions made based on this information, nor for any loss or damage resulting from its use. Links to external sites are for convenience; their inclusion does not imply endorsement. We are not responsible for any interruptions due to technical issues. Users should seek professional counsel before making business decisions, and by using this document, agree to the terms of this disclaimer.

Copyright

© [2023] [UNAL]

Contact

jlara@unal.edu.co

Changelog

v0.1	2023-08	Business Understanding: Introduced reserves and objectives, assessed risks, proposed a central machine learning model, detailed "ActuaSolutions" structure.
v0.3	2023-09	Data Understanding: Collected and integrated data from various sources, explored data characteristics and relationships, verified data quality, and conducted preliminary statistical analyses.
v0.5	2023-09	Data Preparation: Selected significant variables, structured data into matrices, performed data cleaning and transformation, applied normalization, and conducted dimensionality reduction with PCA.
v0.9	2023-10	Modeling: Developed Chain Ladder, Linear Regression, Generalized Linear, and Neural Network models; refined algorithms, and trained models with historical data.
v0.9.1	2023-11	Evaluation: Assessed model performance using MAPE, MSE, and R ² ; compared model predictions against actual data, and interpreted the results for each model.
v1.0	2023-12	Deployment: Finalized and documented the Jupyter Notebook for operational use, enabling stakeholders to run the models with new data sets; released the final version with complete user guidelines.

Table of Contents

1 Business Understanding	4
1.1 Business Objectives	4
1.2 Project Success Criteria	7
1.3 Assessing the Situation	7
1.4 Data Mining Objectives	10
1.5 Data Mining Project Plan	10
2 Data Understanding	11
2.1 Collecting Initial Data	11
2.2 Data Description	12
2.3 Exploring Data	13
3 Data Preparation	15
3.1 Data Cleaning and Structuring	15
3.2 Advanced Data Transformation	15
3.3 Data Normalization and Dimensionality Reduction	16
3.4 Correlation Analysis and Summary	16
4 Modelling	16
4.1 Chain Ladder Method	17
4.2 Linear Models	17
4.3 Generalized Linear Models (GLM)	18
4.4 Neural Networks	18
5 Evaluation	19
5.1 Analysis of Error Metrics	19
5.2 Interpretation of Graphical Results	19
5.3 Conclusion of Evaluation Phase	19
6 Deployment	21
Reference List	22

1 Business Understanding

Actuarial reserves, especially in the context of insurance, are estimates of future obligations that a company has due to already issued insurance policies. These estimates consider uncertain events, such as future claims, and aim to ensure that the company has the necessary funds to fulfill its responsibilities as they arise. Following the instructions in [2] we process describing the objectives.

1.1 Business Objectives

Proper management and calculation of reserves are essential to maintain the financial health of **ActuaSolutions** [1] and ensure its strength in the market. By refining this process, the company aims to:

- Ensure solvency and the ability to meet all future obligations.
- Strengthen the trust of its clients, shareholders, and regulators.
- Optimize the allocation of financial resources.
- Establish itself as a leader in the actuarial field by adopting innovative and accurate methodologies.

1.1.1 Organizational Structure

The company "ActuaSolutions" is a leading organization in the actuarial field that has firmly established itself in the market due to its commitment to accuracy, innovation, and adaptability. To better understand how the company operates, it is essential to grasp its organizational structure and how different departments collaborate to achieve business objectives.

At the heart of our actuarial company is a deep understanding of our structured organization which influences our approach to challenges and solutions:

- *Organizational Charts:* Our corporate landscape is detailed with divisions, departments, and project groups. This chart encompasses the hierarchy from top management to the operational level, detailing managers, their respective roles, and responsibilities.
- *Key Individuals:* We've recognized pivotal personnel across departments whose expertise and insights drive our operations forward. Their roles range from strategic planning, operations, to domain-specific expertise in actuarial science.
- *Internal Sponsorship:* Our internal sponsor bridges the financial and domain expertise, ensuring that the project receives both the necessary funds and knowledge backup.
- *Steering Committee:* Our committee provides direction and oversight. Comprising a diverse group of experts, this team ensures that the project aligns

with our corporate goals.

- *Affected Business Units:* Given the cross-functional nature of data mining, multiple business units including marketing, finance, and operations would be stakeholders in the project's outcome.

Interdepartmental collaboration is essential for the company to maintain its excellence and remain a leader in the actuarial field. We describe each department as:

1. Actuarial Reserves Department:

Functionality: This is the heart of the company. This department is responsible for calculating reserves, estimating future obligations, and ensuring that the company has the necessary resources to fulfill its responsibilities.

Key Teams:

- *Actuarial Analysts:* Specialized professionals who use mathematical and statistical techniques to assess risk and determine reserves.
- *Data Specialists:* Experts who work with large datasets, perform exploratory analyses, and prepare data for analysis.
- *IT Team:* Responsible for maintaining and optimizing the tools and platforms used for actuarial calculations.

2. Financial Department:

Functionality: Manages the company's financial resources, ensures obligations are met, and works closely with the actuarial department to adjust reserves.

Key Teams:

- *Accounting:* Monitors and reports financial flows.
- *Investment Management:* Optimizes the company's investment returns.

3. Claims Department:

Functionality: Manages and processes claims submitted by policyholders.

Key Teams:

- *Claims Adjusters:* Analyze and assess the validity of submitted claims.
- *Customer Support Team:* Assist policyholders and facilitate the claims process.

4. Development and Strategy Department:

Functionality: Responsible for researching new markets, products, and strategies to expand and strengthen the company's position.

Key Teams:

- *Market Research:* Study market trends and identify opportunities.
- *Product Development:* Create and optimize the insurance products offered by the company.

5. Technology and Innovation Department:

Functionality: Explores and adopts emerging technologies, develops software solutions, and seeks ways to implement artificial intelligence and machine learning in the company's processes.

Key Teams:

- *Developers:* Build and maintain technological platforms and tools.
- *Machine Learning Specialists:* Design and train machine learning models for various purposes, including reserves.

Each department, with its unique teams and expertise, plays a vital role in the functioning of "ActuaSolutions".

1.1.2 Problem Area Identification

Our prime focus rests on the provision and actuarial domain, with the motivation to ensure accurate reserve estimations and enhanced financial health:

- **Domain:** Our challenge primarily lies within the actuarial realm, impacting areas such as risk assessment, financial obligations, and provisions.
- **Problem Description:** Our business faces the intricate challenge of accurately estimating future obligations based on past and current data.
- **Project Prerequisites:** This venture was conceived from the need for more precise provision estimations. Currently, while we utilize data-driven strategies, there is a gap in implementing advanced data mining techniques.
- **Status & Advocacy:** The project is in its nascent stage and requires both validation and advocacy. Recognizing this, we're ready to present the merits of data mining as a transformative solution to our board and stakeholders.

1.1.3 Current Solution Analysis

While we have legacy systems and processes in place, understanding their strengths and limitations is key:

- **Existing Solutions:** Our current strategy employs traditional actuarial methods combined with basic statistical tools to estimate provisions.
- **Advantages & Disadvantages:** While these methods have served us well over the years, providing consistency and a moderate degree of accuracy, they lack the adaptability and precision offered by advanced data mining. Furthermore, there's been a mixed reception within the organization about the current methods' effectiveness.

Understanding the above components provides a robust framework for our data mining project, ensuring that we're not just data-informed but also business-aligned.

1.2 Project Success Criteria

For a project as comprehensive and essential as ours, defining success is paramount. Recognizing this, we've broken down the criteria into two segments to capture the tangible and the intangible markers of accomplishment:

- **Accuracy Enhancement:** Achieve a specific percentage increase in the accuracy of our reserve estimations compared to our current methods.
- **Operational Efficiency:** Reduce the time taken for provision calculations by an agreed-upon percentage.
- **Financial Improvement:** Realize a quantifiable reduction in financial discrepancies or shortfalls due to more accurate provision calculations.

1.2.1 Subjective Success Criteria

These criteria, while not strictly quantifiable, are essential to ensure that our solutions align with the broader organizational vision and stakeholder expectations:

- **Enhanced Decision-making:** Uncover patterns and insights from our data that empower decision-makers to strategize more effectively.
- **Stakeholder Satisfaction:** Achieve a consensus among key stakeholders that the data mining model provides valuable and actionable insights. Here, feedback from department heads and the steering committee would be pivotal.
- **Discovery of Effective Actuarial Techniques:** While this is harder to pin down, our aim is to uncover and integrate innovative actuarial methods or clusters of techniques that resonate with our company's challenges and goals.

1.3 Assessing the Situation

To set the project on a solid foundation, it's crucial to understand our current standing, the resources we possess, and potential challenges we might face. The following assessment provides a comprehensive understanding of our present resources, needs, constraints, risks, and potential rewards.

1.3.1 Resource Inventory

Before plunging into the data mining process, it's imperative to evaluate our resources thoroughly:

- **Hardware Resources:** An assessment of the necessary hardware to support the data mining processes is vital. This encompasses server capabilities, storage, and other infrastructural necessities.
- **Data Sources and Knowledge Stores:**

- We have various data sources, each with its own type and format. It's essential to note these details as they dictate the preprocessing steps required.
- Currently, our data is stored in a combination of data warehouses and operational databases. We do have live access to these, ensuring real-time data analysis.
- External data acquisition, like demographic details, is under consideration. The costs and benefits will determine this decision.
- Security concerns, especially with sensitive data, have been flagged. We're ensuring that data access adheres to the company's privacy and security guidelines.
- **Personnel Resources:** Our team comprises business and data experts, database administrators, and additional support staff. They're the backbone of the project, ensuring seamless execution and valuable insights.

1.3.2 Requirements, Assumptions, and Constraints

Transparent understanding and communication regarding the project's boundaries and needs are pivotal:

- **Requirements:**
 - Adhering to legal and security restrictions on the data and results is a top priority.
 - Project scheduling has been communicated and agreed upon by all stakeholders.
 - Result deployment methods, including potential web publishing or database entries, have been outlined.
- **Assumptions:**
 - The project may incur additional costs, such as consulting fees, which have been factored into our projections.
 - Data quality has been assumed to a certain standard. Any deviation might necessitate additional preprocessing steps.
 - The management team primarily expects result-oriented insights, but a deeper dive into the models can be provided upon request.
- **Constraints:**
 - All necessary data access permissions have been obtained.
 - Legal constraints on data usage have been vetted.
 - The project budget encompasses all financial requirements.

1.3.3 Risks and Contingencies

Every venture has its potential pitfalls, and a key risk we must address is the possibility of not acquiring the necessary data for model formulation. Identifying and

mitigating these risks in advance is pivotal to ensuring a seamless and unhindered progression:

- **Risks:**

- Scheduling: Project delays.
- Financial: Budgetary concerns from the project sponsor.
- Data: Quality or coverage issues.
- Results: Not meeting initial expectations.
- Data Availability: Challenges in obtaining requisite data for model development.

- **Contingency Plans:**

- For **Scheduling Risk**: Regularly review project milestones, identify potential delays in advance, and adjust the project plan accordingly.
- For **Financial Risk**: Maintain open communication with the project sponsor, establish a buffer in the budget for unforeseen expenses, and seek alternative funding sources if needed.
- For **Data Risk**: Implement data quality assessment measures, invest in data cleansing and enrichment, and explore alternative data sources if necessary.
- For **Results Risk**: Continuously validate and recalibrate the model, involve stakeholders in setting realistic expectations, and prioritize transparency in reporting results.
- For **Data Availability Risk**: Establish multiple data procurement channels, collaborate closely with data providers, and consider using synthetic data as a temporary substitute.

1.3.4 Cost/Benefit Analysis

An in-depth analysis to ascertain the potential return on our investment:

- **Costs:**

- Data collection and any external data.
- Result deployment.
- Operational expenses.

- **Benefits:**

- Achieving the primary objective.
- Gaining deeper insights from data exploration.
- Advantages arising from a superior understanding of the data.

This comprehensive assessment not only lays the groundwork for our data mining project but also ensures that we are always cognizant of our strengths, challenges, and the value we aim to create.

1.4 Data Mining Objectives

The primary purpose is to develop a machine-learning model that can predict future obligations more accurately based on historical data and other relevant variables. This model seeks to identify complex patterns that traditional methodologies might overlook. Furthermore, the proposed model is expected to be adaptable to changes in the market environment and can be retrained and adjusted as needed.

1.5 Data Mining Project Plan

Given the pivotal role that a project plan plays in aligning team members, stakeholders, and resources to the project's goals, the plan must be comprehensive, accurate, and well-informed. Drawing inspiration from the guidelines shared, here is a meticulously crafted project plan for our data mining endeavor.

1.5.1 Project Plan Description

The project plan, which will serve as our central document, encapsulates our goals, resources, risks, and timelines for every phase of the data mining process. This ensures that everyone, from stakeholders to team members, is well-informed about the project's progress and objectives. To enhance accessibility and increase collaboration, we propose publishing this plan on the company's intranet.

The project plan creation involved discussions with all the stakeholders to ensure its feasibility. We have made sure to include time estimates for all tasks, highlight decision points, and emphasize stages that might require multiple iterations.

Phase	Time (Days)	Resources	Risks
Business Understanding	7	All analysts	Economic changes
Data Understanding	21	All analysts	Data-related issues, technological problems
Data Preparation	35	Data mining consultant, some database analyst time	Data-related issues, technological problems
Modeling	14	Data mining consultant, some database analyst time	Technological issues, finding a suitable model
Evaluation	7	All analysts	Economic changes, challenges in result implementation
Deployment	7	Data mining consultant, some database analyst time	Economic changes, challenges in result implementation

Table 1: Data Mining Project Plan

2 Data Understanding

2.1 Collecting Initial Data

Within our CRISP-DM methodology, our data journey begins with curating a pivotal dataset sourced from the Loss Reserving Data of the National Association of Insurance Commissioners (NAIC) Schedule P. This dataset, meticulously compiled for claims reserving studies, originates from a wealth of claims data encompassing major personal and commercial lines in the US property casualty insurance sector. Comprising six distinct lines of business, from private passenger auto liability to workers' compensation, this repository unveils the intricate tapestry of risk management and coverage within the diverse insurance landscape.

The data preparation process involved a meticulous three-step approach, ensur-

ing that the dataset was refined and robust for our analytical endeavors:

Step I: Pulling Triangle Data: Our journey begins with the extraction of triangle data from the Schedule P of the year 1997. Each of these triangles encapsulates claims spanning 10 accident years from 1988 to 1997, coupled with 10 development lags. This initial dataset forms the bedrock of our model development phase.

Step II: Triangle Squaring: In the pursuit of enhancing our model's efficacy, we square the triangles from the 1997 Schedule P dataset with outcomes from subsequent years. This cross-temporal integration allows us to validate and test our model's performance retrospectively.

Step III: Quality Assurance through Sampling: The dataset's integrity is of paramount importance. To ensure its quality, a meticulous sampling process was conducted. Insurers were retained in the final dataset based on stringent criteria. These include the completeness of observations, alignment of claims data, and non-zero net premiums written for all years.

The culmination of this thorough data preparation process is a refined dataset, encapsulating run-off triangles spanning six distinct lines of business. These triangles encapsulate claims from the accident year 1988 to 1997, paired with a 10-year development lag. Importantly, both upper and lower triangles are included, affording us the opportunity to model and test the performance of our strategies with a retrospective lens.

2.2 Data Description

The dataset for our claims reserving studies project is derived from the Loss Reserving Data provided by the National Association of Insurance Commissioners (NAIC) through their Schedule P. This meticulously curated dataset encompasses a total of 13,200 records, each containing 13 attributes, making it a rich and comprehensive source of insights.

The dataset's diverse attributes and comprehensive size make it an invaluable resource for our claims reserving studies, enabling us to analyze historical patterns and trends within the Workers' compensation line of business.

- **GRCODE:** NAIC company code, including insurer groups and single insurers.
- **GRNAME:** NAIC company name, including insurer groups and single insurers.
- **AccidentYear:** Accident year (1988 to 1997).
- **DevelopmentYear:** Development year (1988 to 1997).
- **DevelopmentLag:** Development year ($AY-1987 + DY-1987 - 1$).
- **IncurLoss_D:** Incurred losses and allocated expenses reported at year end.

- **CumPaidLoss_D**: Cumulative paid losses and allocated expenses at year end.
- **BulkLoss_D**: Bulk and IBNR reserves on net losses and defense and cost containment expenses reported at year end.
- **EarnedPremDIR_D**: Premiums earned at incurral year - direct and assumed.
- **EarnedPremCeded_D**: Premiums earned at incurral year - ceded.
- **EarnedPremNet_D**: Premiums earned at incurral year - net.
- **Single**: 1 indicates a single entity, 0 indicates a group insurer.
- **PostedReserve97_D**: Posted reserves in year 1997 taken from the Underwriting and Investment Exhibit – Part 2A, including net losses unpaid and unpaid loss adjustment expenses.

Attributes such as 'GRCODE' and 'GRNAME' provide information about NAIC company codes and names. 'AccidentYear' and 'DevelopmentYear' span the years 1988 to 1997. Additional attributes like 'IncurLoss_D', 'CumPaidLoss_D', and others denote specific loss and premium data for the 'Workers' compensation' line of business.

2.3 Exploring Data

In this Data Exploration section, we present a comprehensive visualization of the key statistical characteristics of the run-off triangles, grouped by the indicators BulkLoss_D, CumPaidLoss_D, and IncurLoss_D. Each row in Figure 1 reflects one of these indicators, breaking down the data into measures of central tendency and variability. The first column displays the mean of the consolidated datasets, offering an overview of the average value across all triangles in our database. The second column exhibits the variance, providing quantification of data variability or dispersion. Finally, the third column presents the median, facilitating an understanding of the data distribution by indicating the central value. This systematic arrangement allows for direct and efficient comparison between different statistical metrics, essential for accurate interpretation of the underlying dynamics in actuarial reserve data.

The preceding Figure 1 visually synthesizes the essential statistical information of the run-off triangles, providing a solid foundation for further analysis in the process of Actuarial Data Mining. The orderly layout in a 3x3 format emphasizes the differences and similarities between the measures of mean, variance, and median, which are crucial for actuarial valuation and reserve prediction. The detailed graphical exploration of BulkLoss_D, CumPaidLoss_D, and IncurLoss_D paves the way for an in-depth discussion on the suitability of predictive models and the robustness of reserve estimates. The insights gained from these visualizations will be of paramount importance in the subsequent phases of the "Data Mining in Actuarial Science: Reserves" project, particularly when addressing the selection and application of advanced analytical techniques within the CRISP-DM methodol-

Data Mining in Actuarial Science: Reserves

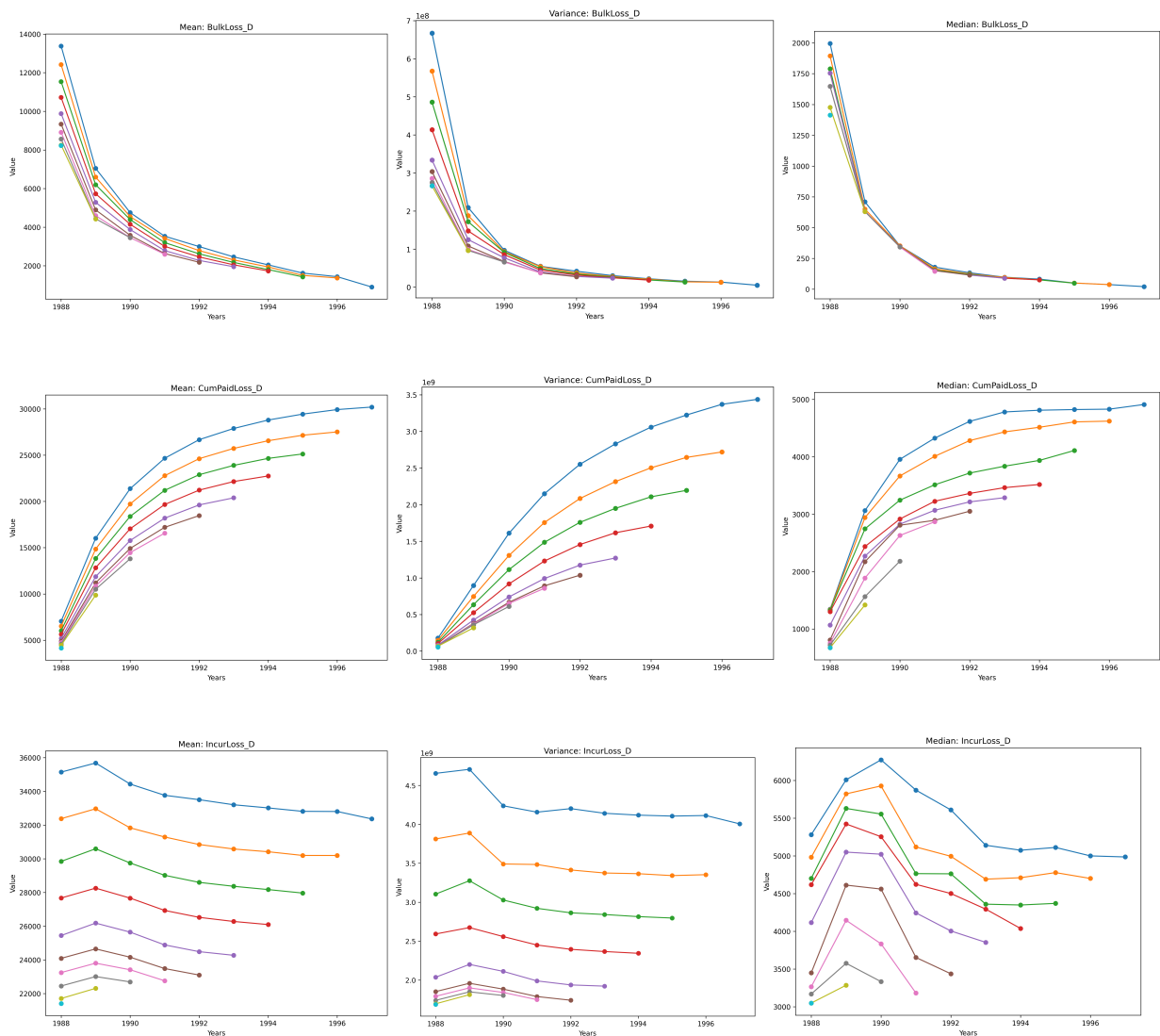


Figure 1: Statistical Overview of Run-Off Triangle Metrics

ogy. The interpretation of these figures will contribute to building a comprehensive framework for strategic decision-making in actuarial risk management.

3 Data Preparation

The Data Preparation phase of our actuarial claims reserving project is a crucial stage where we transform and refine raw data to create a robust analytical foundation.

3.1 Data Cleaning and Structuring

Our initial focus was on data cleaning and structuring. The process began with the meticulous selection of essential variables that form the core of our analysis. We identified and extracted key attributes such as incurred losses, paid losses, bulk reserves, and earned premiums. These variables were chosen due to their significant role in understanding the financial behavior of claims over time.

With these critical variables at hand, we proceeded to structure the data into a decade-long development format. The 10x10 matrix structure we employed serves as a fundamental tool for visualizing the progression of claim settlements and reserve allocations. This step was not only about organizing data into a readable format but also about preparing it for the application of actuarial techniques.

3.2 Advanced Data Transformation

Following the initial structuring, we engaged in advanced data transformation techniques. Triangulation, a critical process in actuarial data preparation, was executed to align our data with industry-standard reserving methods like the chain-ladder and Bornhuetter-Ferguson techniques. This transformation involved converting our structured matrices into triangular forms, facilitating the modeling and estimation of claims payment development.

In parallel, we conducted a thorough initial data check to ensure data quality. This involved verifying the absence of missing values and ensuring the appropriate data type assignments for each variable. Special attention was given to the financial columns, where we identified the presence of negative values. Recognizing their potential impact as adjustments or returns in financial records, we decided to preserve these values to maintain the integrity and completeness of our actuarial analysis.

3.3 Data Normalization and Dimensionality Reduction

A key step in our data preparation was normalization, which we performed on the numerical columns. The purpose of this step was to equalize the contribution of each variable to the analysis, ensuring that the variance in variable scales did not skew our subsequent analyses.

Building upon the normalized data, we applied Principal Component Analysis (PCA) for dimensionality reduction. The PCA process allowed us to distill the dataset into principal components that encapsulated a significant proportion of the dataset's variance. Interestingly, our PCA results showed that the first two components accounted for approximately 89% of the total variance, highlighting substantial underlying patterns within the data.

3.4 Correlation Analysis and Summary

To further enhance our understanding of the dataset, we conducted a correlation analysis. This involved examining a correlation matrix of the financial figures to ascertain significant relationships between variables. The insights gained from this analysis supported the application of PCA, confirming its ability to effectively consolidate correlated variables.

Overall, the Data Preparation phase has been instrumental in setting the stage for sophisticated actuarial modeling and analysis. We have diligently cleaned, normalized, and reduced the complexity of the data, ensuring that it is primed for extracting valuable insights in the subsequent phases of our project.

4 Modelling

In the modeling phase of our actuarial reserving analysis, we implemented a comprehensive approach using various predictive models. Our objective was to analyze and optimize reserve data, identifying patterns and risks that could inform effective financial decisions. We employed the following models: Chain Ladder Method, Linear Models, Generalized Linear Models (GLM), and Neural Networks. Each of these models was chosen for its relevance and capability to handle the complexities of actuarial data.

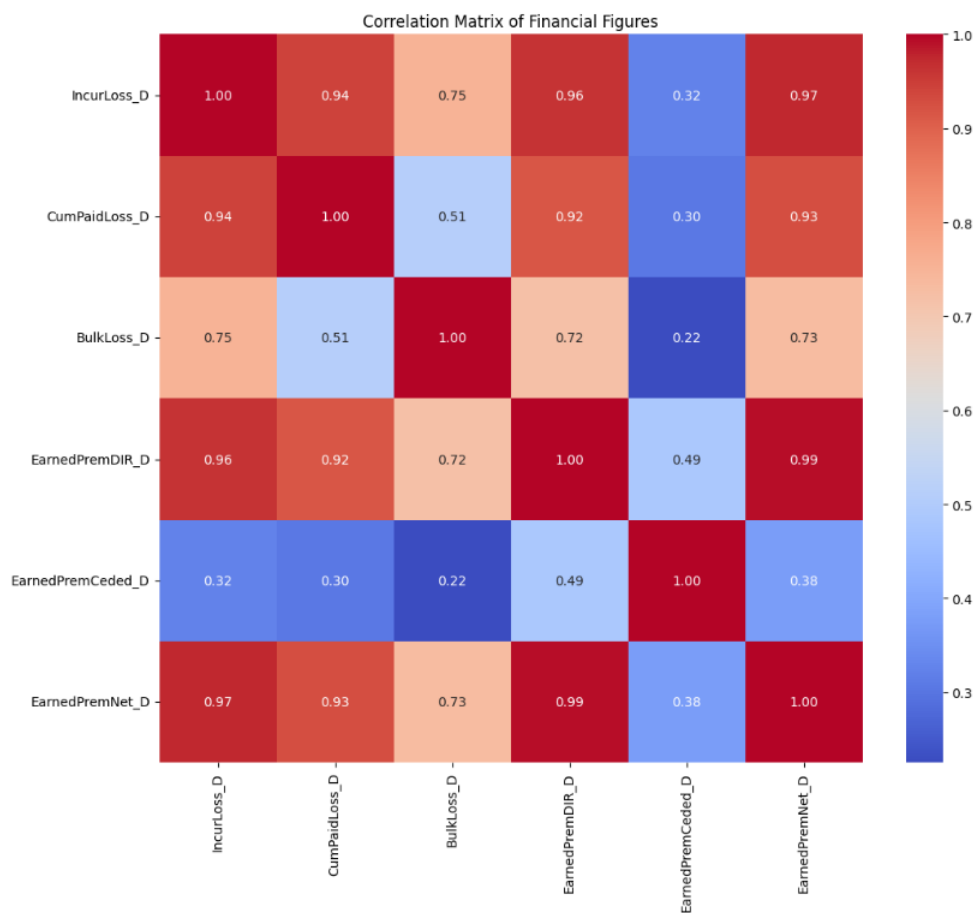


Figure 2: Correlation Matrix

4.1 Chain Ladder Method

The Chain Ladder Method (CLM) stands as one of the quintessential techniques in actuarial science for projecting future claims. Renowned for its straightforward application and robust historical performance, the CLM predicates on the assumption that past claims development is indicative of future trends. Our application of CLM entailed an intricate process of calculating age-to-age factors, which represent the ratios of claims development from one period to the subsequent one. These factors were then methodically applied across the development triangle, enabling us to extrapolate and estimate the ultimate claims settlements. Despite its traditional roots, the CLM’s enduring relevance stems from its adaptability to diverse claims data scenarios, providing a foundational benchmark against which we gauged more complex models.

4.2 Linear Models

Our exploration of predictive capabilities further led us to the realm of Linear Models (LMs), which are revered for their analytical clarity in revealing relationships between predictors and claim amounts. By adopting both simple and multiple re-

gression techniques, we could distill the essence of linear associations within our data. Simple linear regression offered insights into bivariate relationships, while multiple regression encompassed a broader spectrum of factors simultaneously. This dual approach facilitated a granular understanding of how various predictors, such as policyholder demographics and economic indicators, influenced the magnitude of claims. The application of LMs served as a crucial stepping stone in our analysis, setting a precedent for comparison with more intricate models and enriching our comprehension of the data's linear structure.

4.3 Generalized Linear Models (GLM)

The inherent complexity of actuarial data, often characterized by non-normal distribution patterns, necessitated the deployment of Generalized Linear Models (GLMs). As an extension of LMs, GLMs transcend the conventional normality constraint by embracing diverse distributions for the response variable—most notably, Poisson for claim counts and Binomial for binary outcomes. This versatility was pivotal in our analysis, allowing us to adeptly model the frequency and severity of claims with a precision tailored to the actual data behavior. By leveraging GLMs, we could accommodate the skewness and kurtosis inherent in our actuarial datasets, thereby achieving a more nuanced fit and enhancing the fidelity of our predictions.

4.4 Neural Networks

In pursuit of capturing the intricate and non-linear intricacies embedded within our data, we ventured into the domain of Neural Networks (NNs). Specifically, we architected a multi-layer perceptron (MLP), a class of NNs renowned for its proficiency in pattern recognition. Our MLP consisted of multiple layers of neurons, each interconnected to form a web of nodes capable of complex representations. Through an iterative process of learning and adjustment, known as backpropagation, our MLP was trained on historical data, discerning subtle patterns and interactions that more traditional models could potentially overlook. The NN's prowess in handling vast datasets and its ability to internalize non-linearity rendered it an invaluable asset in our analytical arsenal, significantly contributing to the depth and breadth of our analysis.

5 Evaluation

In the evaluation phase of our study, each predictive model was scrutinized using a variety of statistical measures to assess accuracy and predictive power. Our analysis utilized Mean Absolute Percentage Error (MAPE), Mean Squared Error (MSE), and the Coefficient of Determination (R^2). The results of this comprehensive evaluation revealed distinct performances across the models, shedding light on their respective strengths and limitations in the context of actuarial reserving.

5.1 Analysis of Error Metrics

For IncurLoss_D, the Linear Model Method displayed superior precision with the lowest MAPE, indicating its predictions were the most aligned with actual values. It also boasted the smallest MSE, suggesting fewer deviations from observed data, and achieved the highest R^2 , reflecting a robust model fit.

Conversely, for CumPaidLoss_D, the same Linear Model Method maintained its superior performance, leading in all three metrics and demonstrating high accuracy and an excellent fit.

The Chain Ladder Method, while a traditional approach in actuarial science, showed significant shortcomings. It registered an exceedingly high MAPE, particularly for BulkLoss_D, which signaled considerable predictive errors. Its MSE was also notably high, and the R^2 was profoundly negative, indicative of a poor fit.

In comparison, the Generalized Linear Model (GLM) and Neural Network (NN) exhibited similar levels of performance. Both models encountered challenges, reflected by elevated MAPE and MSE values and negative R^2 values for IncurLoss_D and CumPaidLoss_D, suggesting a less than satisfactory model fit.

5.2 Interpretation of Graphical Results

The graphical representations of these metrics further illustrate the disparities between the models. As seen in Figure 3, the MAPE values for the Linear Model are significantly lower than those for the other methods, reinforcing its superior accuracy. Figure 4 displays the MSE results, where again, the Linear Model Method demonstrates its enhanced predictive capability with notably lower values than its counterparts. Finally, in Figure 5, the R^2 values for the Linear Model Method are closer to 1, especially for IncurLoss_D, highlighting its strong fit.

5.3 Conclusion of Evaluation Phase

The evaluation phase has provided us with critical insights into the predictive models' performance. The Linear Model Method consistently emerged as the superior approach, excelling across all metrics and variables. In contrast, the Chain

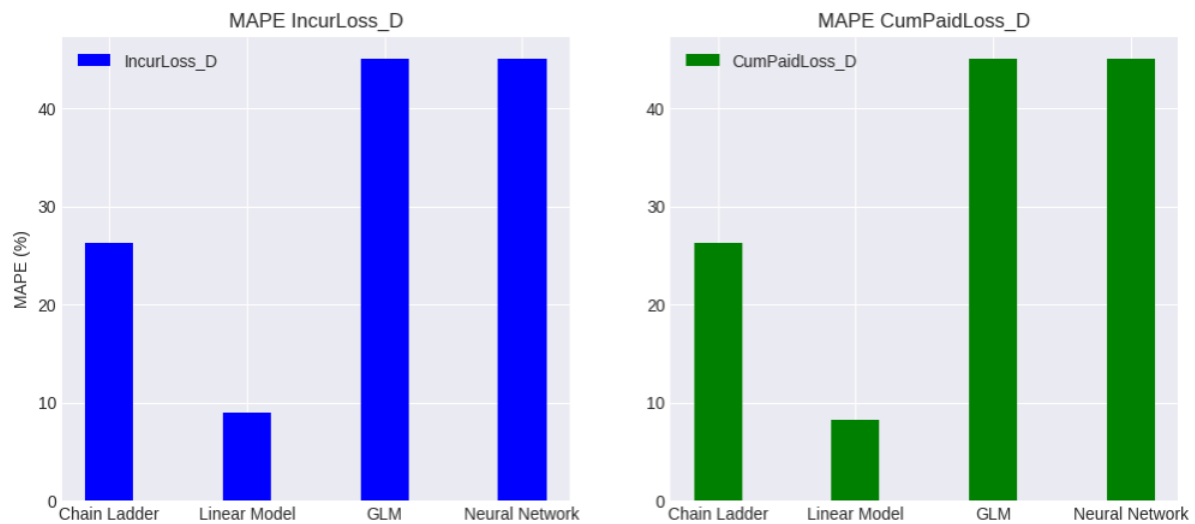


Figure 3: Comparison of MAPE across different models

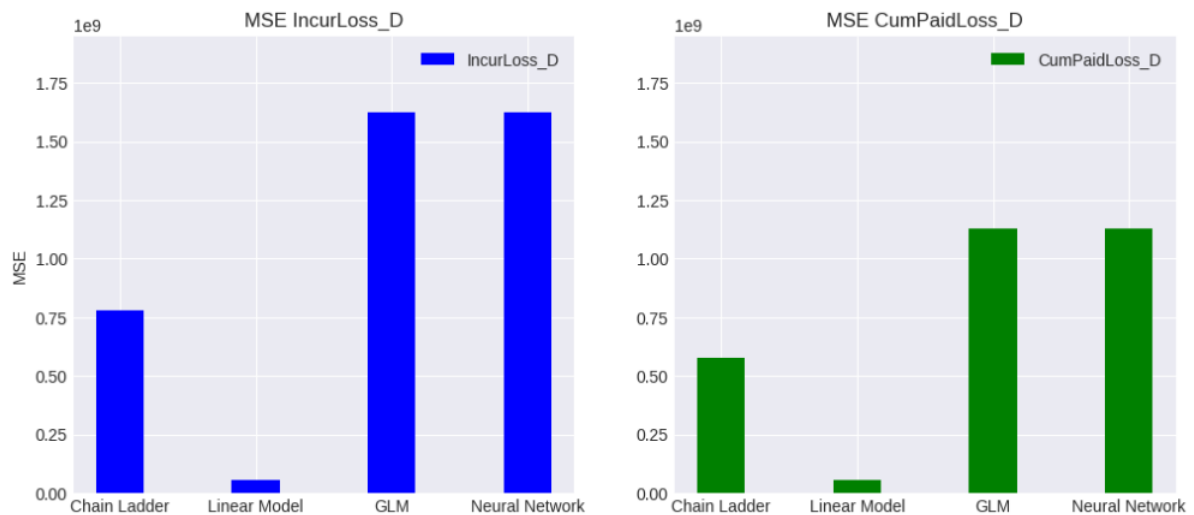


Figure 4: Comparison of MSE across different models

Ladder Method's performance was less satisfactory, struggling particularly with the BulkLoss_D variable. Both GLM and NN showed limitations, with suboptimal fit and accuracy for certain variables. These findings are integral to our understanding of the models' capabilities and will guide future actuarial reserving efforts.

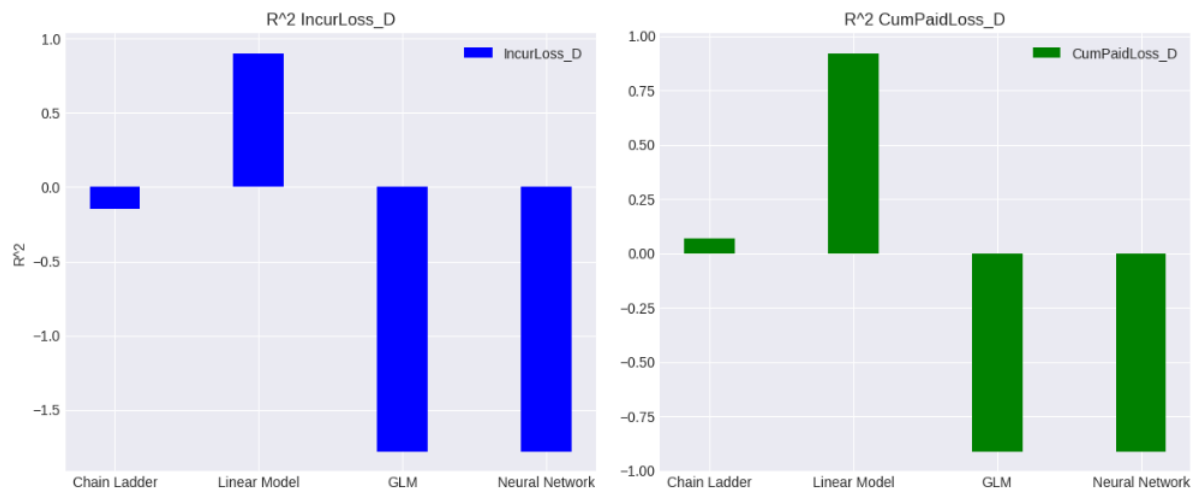


Figure 5: Comparison of R² across different models

6 Deployment

Deployment marks the transition of our actuarial models from a developmental stage to practical application. To execute the models, stakeholders are guided to access the Jupyter Notebook located in the ‘notebooks’ directory of our GitHub repository. This interactive platform is not only a walkthrough of the analytical journey but also a tool for applying the models to new data sets.

Users begin by cloning the repository and navigating to the notebook, which presents the data mining process through executable code cells. Each predictive model, from the Chain Ladder Method to Neural Networks, is encapsulated as a callable function within the notebook. These functions are pre-programmed to accept new runoff triangle data, perform the necessary computations, and output predictions along with dynamic visualizations reflecting the results.

The notebook is structured with clear instructions on data formats, model execution, and result interpretation, ensuring that even those with limited coding expertise can deploy the models efficiently. It is maintained under version control for ongoing enhancements and encourages collaborative improvement through user feedback and contributions.

Through this deployment strategy, our project delivers a comprehensive toolset for actuaries to perform reserve analysis with confidence, ensuring our models’ continued relevance and utility in the dynamic landscape of actuarial science.

Reference List

- [1] *Actua Solutions*. es-ES. URL: <https://actuasolutions.com> (visited on 08/25/2023).
- [2] Salem Al Gharbi et al. "Using Data-Mining CRISP-DM Methodology to Predict Drilling Troubles in Real-Time". In: *Day 1 Tue, November 17, 2020*. Virtual: SPE, Nov. 2020, D013S103R013. DOI: 10.2118/202326-MS. URL: <https://onepetro.org/SPEAPOG/proceedings/20APOG/1-20APOG/Virtual/451699> (visited on 08/25/2023).
- [3] Ilham Battas, Hicham Behja, and Laurent Deshayes. "DMAICS 2 CRISP DM' approach for improving and optimising the performance of an industrial mining production process". en. In: *International Journal of Six Sigma and Competitive Advantage* 14.4 (2023), pp. 408–436. ISSN: 1479-2494, 1479-2753. DOI: 10.1504/IJSSCA.2023.134444. URL: <http://www.inderscience.com/link.php?id=134444> (visited on 12/05/2023).
- [4] Andriy Burkov. *The hundred-page machine learning book*. eng. Polen: Andriy Burkov, 2019. ISBN: 9781999579500 9781999579517.
- [5] Michael Franke. *An Introduction to Data Analysis*. URL: <https://michael-franke.github.io/intro-data-analysis/index.html> (visited on 12/05/2023).
- [6] Trevor Hastie, Robert Tibshirani, and J. H. Friedman. *The elements of statistical learning: data mining, inference, and prediction*. 2nd ed. Springer series in statistics. New York, NY: Springer, 2009. ISBN: 9780387848570 9780387848587.
- [7] Paweł Mielcarz, Dmytro Osiichuk, and Paweł Wnuczak. "Actuarial Reserves, Provisions and Contingent Liabilities in DCF Valuation". In: *Zeszyty Naukowe Uniwersytetu Szczecińskiego Finanse Rynki Finansowe Ubezpieczenia* 79 (2016), pp. 289–298. ISSN: 2450-7741, 2300-4460. DOI: 10.18276/frfu.2016.79-22. URL: <https://wnus.edu.pl/frfu/pl/issue/655/article/10387/> (visited on 12/05/2023).