# CP3403 DATA MINING PROJECT REPORT

## Adult Mortality Rate (2019 – 2021)

**19th April 2024**

**Members – Group 3**

Nang Kaung Shan Kham (14567218)

Hpu Hpu Thant Sin (14396720)

Nyan Lin Htut (14419795)

Kyaw Zaww Linn (14430237)

# Table of Contents

## Abstract

This project explores the application of data mining techniques in the global health insurance sector to enhance market segmentation and risk assessment strategies. Utilizing a comprehensive dataset on Adult Mortality Rates (2019-2021), this study employs the k-means clustering algorithm to identify distinct groups of countries based on health outcomes, economic status, and demographic factors. The primary objective is to discern patterns that can inform tailored insurance product offerings and strategic marketing efforts. By integrating these insights, the project aims to assist a global health insurance provider in optimizing their product alignment and market approach, catering specifically to regional health profiles and economic conditions. The findings demonstrate the value of advanced data mining in understanding complex global markets, ultimately aiding in the development of more effective and competitive insurance solutions tailored to meet diverse customer needs. This study not only highlights the practical implications of data mining for strategic decision-making in health insurance but also sets a framework for future research in this vital area of business analytics.

## Introduction

In the dynamic field of global health insurance, the ability to adapt and respond to the diverse needs of different regions is crucial. This project aims to leverage data mining techniques, particularly focusing on the Adult Mortality Rates from 2019 to 2021, to improve market segmentation and risk assessment in the insurance sector. By employing advanced clustering and classification methods, the study seeks to uncover patterns and correlations among demographic, economic, and health-related factors across various countries.

The project revolves around a comprehensive dataset that includes variables such as GDP per capita, health expenditure, and mortality rates. Through sophisticated algorithms such as k-means clustering, the team categorizes countries into distinct groups, thus providing a nuanced understanding of global health trends. This segmentation allows for the development of tailored insurance products that align with specific regional health profiles and economic conditions, thereby enhancing the strategic decision-making capabilities of health insurance providers.

The objective of this endeavor is twofold: to aid a global health insurance company in optimizing its product offerings and to set a foundation for further research in applying data mining for business analytics in the health insurance domain. As we navigate through the methodologies and findings, this report will highlight the potential of data-driven strategies in transforming the approach to global health insurance segmentation and risk management.

## Business Scenario

In the realm of global health insurance, understanding nuanced regional differences in health outcomes is critical for tailoring insurance products effectively. Our company is launching a sophisticated data mining initiative aimed at segmenting international markets based on a comprehensive dataset that includes "Adult Mortality Rate (2019-2021)." By employing advanced clustering techniques, we aim to uncover patterns in health metrics, economic status, and demographic factors across various countries.

The initiative involves dissecting the dataset to identify clusters of countries that exhibit similar characteristics in terms of economic indicators like GDP per capita, health expenditures, and demographic profiles such as population size. These clusters will help us recognize which markets share similar health risks and insurance needs, allowing us to customize our insurance offerings more precisely to meet local demands.

Moreover, our analysis will extend to health and risk profiling. Using key data points such as Adult Mortality Rates for males and females and the Average Crude Death Rate, we will map out regions based on prevailing health conditions and risks. This profiling is essential for adjusting our life and health insurance products to align better with regional medical realities and life expectancy trends.

Additionally, integrating data on the development level of each region with economic and health statistics will enable us to further refine our market segmentation. This will assist in determining the maturity of insurance markets in various regions, guiding strategic decisions regarding where to intensify marketing efforts and product deployment.

Our approach aims to provide a multi-dimensional view of potential insurance markets, guided by data-driven insights into health status, economic conditions, and demographic factors. This will not only help in identifying high-potential markets for expansion but also enhance our risk assessment capabilities, leading to more accurately priced and competitively positioned insurance products.

In conclusion, through this data mining project, we anticipate gaining a deeper understanding of global health trends and regional market potentials, which will empower us to offer more tailored insurance solutions. This strategic initiative is expected to bolster our market presence by aligning product offerings with the specific needs and challenges of diverse populations, ultimately enhancing customer satisfaction and business growth.

## Data Description

The data set used is the Adult Mortality Rate (2019 - 2021). This data set includes 156 instances. The following attributes, which are of the nominal and numeric data type, have been kept for further analysis :

- Countries

- Continent

- Average_Pop(thousands people)

- Average_GDP(M$)

- Average_GDP_per_capita

- Average_HEXP($)

- Development_level

- AMR_female(per_1000_female_adults)

- AMR_male(per_1000_male_adults)

- Average_CDR

## Data Preprocessing

### Data Cleaning

This foundational step involves several key actions to ensure the dataset is free of errors and inconsistencies:

- Removing or correcting erroneous data

- Filling or removing missing values in important variables

- Identifying and handling outliers in data distributions

As the data set used has already been cleaned initially, additional preprocessing and cleaning to remove error data and missing values were unnecessary to perform.

Specific countries were identified as outliers and removed based on unique economic and demographic characteristics that could distort the overall analysis:

- Luxembourg (ID 147): Removed due to its extremely high GDP per capita and small population, which could skew economic indicators.

- Ethiopia (ID 139): Excluded due to its status as one of the least economically developed countries, presenting unique challenges that are not representative of the broader global market.

- Qatar (ID 152): Removed due to its high GDP coupled with low life expectancy and a significant expatriate population, creating anomalies in health expenditure and mortality rate analysis.

**Rationale for Removing Outliers**

The decision to remove these outliers was based on their potential to disproportionately influence statistical analyses and the relevance of their unique conditions to the global insurance strategy. This ensures that our models and strategies are built on data that reflect the typical markets the business aims to serve, enhancing the accuracy and applicability of our predictive models and strategic decisions.

**Discretization, Normalization and Standardization**

In the revised approach to data pre-processing for our data mining project, we have streamlined the treatment of the dataset by focusing our efforts on discretization for only two attributes: GDP per capita and health expenditure. This targeted discretization allows us to

categorize these economically significant variables into 'low', 'medium', and 'high' groups. Such categorization aids in the analysis by simplifying the economic conditions relative to insurance affordability and health service utilization across different regions, making it easier to interpret the impact of economic status on health outcomes.

Furthermore, we have opted not to apply normalization and standardization processes across our dataset. The decision to exclude these transformations is based on our specific analytical needs and the nature of our data, which does not require scaling for the algorithms we intend to use. This approach reduces complexity in our data processing and maintains the original scale of the data, which might be crucial for interpreting results in a business context where actual values have intrinsic meaning.

Additionally, all numeric attributes in the dataset have been converted to normal attributes. This conversion involves treating numeric data as categorical, which can be particularly beneficial for certain types of analysis where binning or direct comparisons between grouped categories are more insightful than dealing with a wide range of numbers. This transformation simplifies the data and can enhance the performance and interpretability of certain classification algorithms by reducing the noise and variability that often accompany raw numerical data.

## Methods and Algorithms

In our project on the 'Adult Mortality Rate (2019-2021)' dataset, the initial data preparation will ensure all data is formatted correctly, handling any missing values, and normalizing significant economic indicators such as GDP to maintain consistency across different scales (Han, Pei, & Kamber, 2011). We will segment continuous variables, such as age and GDP, into categorized bins to facilitate analysis (James, Witten, Hastie, & Tibshirani, 2013).

For the analytical part, we start with clustering techniques, specifically implementing k-means clustering due to its efficiency with large datasets, helping us identify patterns across demographic and economic factors effectively (Han, Pei, & Kamber, 2011). Features like health expenditure per capita and GDP per capita will be central to our analysis, aiming to segment countries into clusters with similar mortality characteristics (World Health Organization, 2021).

Following clustering, we will develop classification models using J48, IBK, Naive Bayes and OneR to predict mortality risk categories (Kuhn & Johnson, 2013). The data will be split into training and testing sets to validate the accuracy of our models, employing cross-validation to avoid overfitting (James, Witten, Hastie, & Tibshirani, 2013).

Additionally, we aim to integrate insights from both clustering and classification in a comprehensive risk assessment model, crucial for international health organizations and insurance companies in their planning and strategy development (Wang et al., 2020).

The final step in our methodology involves validating these models through statistical methods to ensure reliability and refining them based on real-world feedback and application outcomes (Rajulton, Ravanera, & Beaujot, 2007). This iterative process will help us refine our predictions and enhance the overall effectiveness of our models, enabling stakeholders to make informed decisions based on robust data-driven insights (International Monetary Fund, 2021).

**Clustering**

In our analysis of the 'Adult Mortality Rate (2019-2021)' dataset, we employ the k-means clustering algorithm, a robust method for grouping data based on similar characteristics. To determine the most appropriate number of clusters for our analysis, we utilize the Elbow Method. This method involves plotting the within-cluster sum of squares (WCSS) against a range of potential cluster counts (k-values). The WCSS represents the total squared distance between each point in a cluster and the centroid of that cluster. We look for the 'elbow' point in this plot where the decrease in WCSS begins to level off, indicating that adding more clusters does not significantly improve the tightness of the clusters, thus suggesting an optimal k-value for our analysis (Han, Pei, & Kamber, 2011; Jain, 2010). This technique is commonly applied in strategic management research to determine the number of clusters that provide the most meaningful interpretation (Ketchen & Shook, 1996).

To perform clustering using the SimpleKMeans method, start by selecting the 'Cluster' option in Weka software. Choose 'SimpleKMeans' from the available algorithms. Then, click on 'Ignore Attributes'. This step is crucial as it excludes irrelevant attributes from the analysis,

focusing the clustering process on the most significant features that influence the formation of distinct groups. Ignoring these attributes can enhance the algorithm's effectiveness by reducing noise and unnecessary complexity in the dataset.



Next, configure the SimpleKMeans algorithm by setting 'numClusters' to 4, which dictates the algorithm to partition the data into four distinct clusters. Keeping the 'seed' at its default value ensures consistency in the results across multiple runs, provided all other parameters remain unchanged. This step is essential to maintain reproducibility in the analysis while the chosen number of clusters allows for a manageable yet insightful breakdown of the dataset into meaningful groups.

Continuing with the clustering process, once all settings are in place, use all attributes available in the dataset for a comprehensive analysis. Execute the clustering algorithm, specifically focusing on two key areas: Health and Risk Profile, and Development-Based Clustering. By segmenting the data according to these distinct categories, we gain insights into different facets of the data set, identifying unique patterns and relationships within each cluster that pertain to these

specific areas. This approach allows for a nuanced understanding of the dataset, facilitating targeted analysis and decision-making based on clustered group characteristics.

**Classification**

In our comprehensive study, we have strategically harnessed the capabilities of several classification algorithms to scrutinize the "Adult Mortality Rate (2019-2021)" dataset. The overarching aim of this analysis is to forecast potential mortality risks and to accurately delineate the spectrum of insurance needs prevalent across a multitude of countries. We meticulously selected algorithms that harmonize with the distinct attributes of our data corpus:

- **Decision Trees**: Our preference for Decision Trees is rooted in their exceptional ability to render complex decision-making processes into an easily interpretable visual structure. This attribute is of paramount importance as it allows us to efficaciously segment our dataset. By capitalizing on key mortality rate predictors, such as health expenditure and GDP per capita, Decision Trees provide a clear roadmap for understanding the factors that underpin mortality risks in various regions.

- **Naive Bayes**: We champion the Naive Bayes algorithm for its unparalleled efficiency, particularly when working with datasets wherein predictors are presumed to function independently. When confronted with our preprocessed categorical data—now ingeniously categorized into discrete age brackets and economic classifications—the Naive Bayes algorithm performs admirably, offering a fluid and coherent analysis pipeline.

- **IBk (Instance-Based Learning by k-nearest neighbors)**: The inherent flexibility of the IBk approach in capturing non-linear patterns makes it an indispensable tool in our analytical arsenal. Its ability to model the intricate web of dependencies that exist among variables such as health expenditures and mortality rates is crucial. IBk's method of discerning patterns through an examination of instance similarities lends itself to a finely tuned interpretation of data, wherein contextual nuances and the relational proximity of data points are given due weight.

- **OneR**: The OneR algorithm is lauded for its straightforward yet profound modus operandi. It stands as a beacon of simplicity, rapidly sifting through demographic and economic predictors to stratify countries into distinct mortality risk categories. This early-phase analysis affords us a lucid grasp of the variables with the most gravitas, setting a solid foundation for subsequent, more elaborate data modeling.
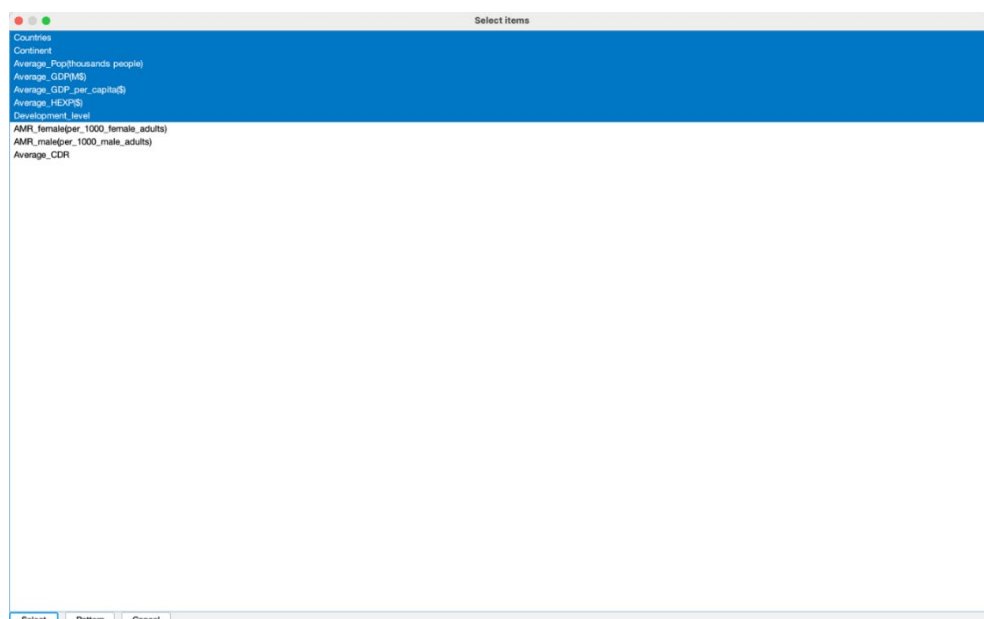
In conclusion, these algorithms form the backbone of our classification strategy, each contributing its unique strengths to the cause. Through their concerted application, we strive to shed light on the complicated patterns that dictate mortality risks, thus empowering insurance entities to tailor their offerings with a high degree of precision and foresight. The judicious application of these classification techniques paves the way for a nuanced understanding of global health trends and ushers in a new era of data-informed decision-making in the domain of health insurance.

## Findings and Discussions

**Clustering**

**Health and Risk Profile Clustering**

The data analysis begins by discarding irrelevant attributes to streamline the dataset.

We then apply the elbow method, testing k-values from 2 to 10, to identify the optimal cluster count by looking for the point at which further clustering yields minimal benefit, an essential step in k-means cluster analysis.

Elbow method for Optimal k





We drew the graph out in microsoftExcel manually so that we can clearly find out the number of clusters we should apply. According to the figure, the significant number of clusters which is "k" tends to be 4 (k=4).

The Elbow Plot and WEKA clustering output demonstrate the application of the k-means algorithm to stratify data into k=4 distinct clusters, focused on health and risk profiles. The plot illustrates the WCSS versus the number of clusters, which decreases as more clusters are added. The leveling off of the WCSS curve beyond k=4 suggests that this is a suitable choice for the number of clusters, as further increases in k result in minimal gains in data modeling efficiency.

The analysis, using a k=4 cluster model, divides the dataset into groups that likely represent different health and risk characteristics. These may encapsulate variables such as health expenditures, morbidity and mortality rates, and other health-related indicators across different populations or geographic regions. Each cluster encapsulates a unique combination of these attributes, which could have significant implications for public health planning, insurance risk assessment, and policy development. By segmenting countries or regions into these four clusters, stakeholders can tailor health interventions, resource distribution, and risk mitigation strategies to the specific needs and profiles of each group.

This granular approach to clustering provides a clearer understanding of the underlying patterns within the health and risk data. It allows for a more nuanced understanding of the disparities or similarities that exist across the dataset, which is crucial for addressing public health concerns and managing risks effectively. The k=4 clusters serve as a means to categorize and compare various demographics, enabling focused analysis and potential identification of areas for improvement or investment in healthcare systems and policies.

## Development-Based Clustering

Starting with the elimination of superfluous attributes, the data analysis process is streamlined.

Following this, the elbow method is utilized, exploring k-values ranging from 2 to 10 to pinpoint the most efficient number of clusters.



Elbow method for Optimal k

Based on the figure shown above, it appears that the optimal number of clusters, which is represented by "k," is likely to be 3 (k=3).



From the Elbow Plot and WEKA clustering output, it is evident that a k-means clustering algorithm has been applied to the dataset, with the aim of discovering inherent groupings based on health and risk profiles. The Elbow Plot shows a decline in the within-cluster sum of squares (WCSS) as more clusters are introduced, but this decrease diminishes, flattening out past the k=3 mark, indicating that additionalclusters do not significantly improve the model's fit to the data.

The WEKA output complements this by detailing a three-cluster solution. These clusters, each comprising a subset of the dataset, are distinguished by their respective health expenditure, GDP per capita, and AMR rates, among other attributes. The clusters presumably encapsulate variations in healthcare infrastructure, economic capacity, and health risk factors of different groups or regions.

For instance, one cluster may represent regions with lower economic and health indices, suggesting populations with less access to healthcare and lower health-related expenditures.

Another cluster mightencompass regions with higher economic indices yet still display moderate development levels, perhaps indicating an uneven distribution of wealth and access to health services.

These cluster insights are pivotal for designing targeted health policies or medical resource allocation. Regions in the cluster with lower economic and health indices might benefit from interventions aimed at improving healthcare access and affordability, whereas regions in the more economically robust clustermay require policies that address inequality and encourage sustainable development.

By segmenting the data into these clusters, policymakers and healthcare providers can tailor their strategies to the specific needs and circumstances of each group, leading to more efficient and effectivehealthcare interventions. The use of three clusters offers a balance between detail and manageability, enabling stakeholders to address the complex landscape of health and risk profiles with a data-informed approach.

**Classification**

In the analysis of the 'Adult Mortality Rate (2019-2021)' dataset, four different classification models were employed to assess their predictive performances.

The **OneR** model, known for its simplicity, uses one rule for prediction. In our dataset, this model achieved an accuracy of 67.3077%, reflecting a basic level of predictive capability. While not highly accurate, OneR models are easy to understand and interpret, making them useful for preliminary analysis.

The **Naive Bayes** classifier, which assumes that the presence (or absence) of a particular feature of a class is unrelated to the presence (or absence) of any other feature, showed better performance. It attained an accuracy of 87.8205%. This model's strength lies in its foundation on probability, making it particularly suited for datasets where the predictors can be assumed to be independent.

For the **IBk** model, which is an instance-based learning algorithm or k-nearest neighbor classifier, the results were quite robust with an accuracy of 86.5385%. This model classifies new instances by analyzing their closest historical examples in the dataset, making it versatile and adaptive to complex data patterns.

Lastly, the **Decision Tree** algorithm outshone all others with a remarkable accuracy of 98.0769%. Decision Trees are valuable for their interpretability and depth, as they use an if-then-else decision rule approach to predict outcomes. Their high accuracy in our analysis demonstrates their potential as a powerful tool for detailed data classification and decision-making processes.

## Results

Based on the classification models evaluated, the results are as follows:

The **OneR** model demonstrated a fundamental level of predictive ability with a modest accuracy rate of approximately 67.31%. This result positions OneR as a baseline model, useful for its simplicity and speed, but perhaps less reliable for complex decision-making tasks where higher precision is required.

The **Naive Bayes** classifier achieved a higher accuracy rate of about 87.82%, indicating a strong performance, especially considering its underlying assumption of feature independence. This suggests that for the given dataset, the Naive Bayes classifier was able to capture the essential patterns efficiently, even if the features may not be entirely independent in reality.



The **IBk** model, employing the k-nearest neighbors algorithm, showed a commendable performance with an accuracy rate of approximately 86.54%. This reflects its capability to effectively classify instances based on similarity measures, confirming its utility for datasets where patterns can be recognized in terms of proximity or resemblance among data points.

Finally, the **Decision Tree** model demonstrated superior accuracy at approximately 98.08%. This exceptional level of accuracy underscores the model's capability to dissect and analyze the dataset's complex structures, delivering precise and actionable classifications. It reaffirms the Decision Tree's strength in dealing with intricate datasets and highlights its potential as a reliable tool for in-depth data analysis in various practical applications.

## The Decision Tree



The decision tree analysis, centered around the 'Average_GDP_per_capita' metric, provides a compelling visual stratification of data that is instrumental for economic segmentation. From its root, the tree branches out into various decision nodes, each representing a level of GDP.

The decision tree analysis, centering on 'Average_GDP_per_capita', serves as a strategic tool for classifying populations based on economic status. With GDP per capita as the decision node, the tree segments countries into distinct categories such as 'Short', 'Average', and 'High'. These classifications likely correspond to quantitative thresholds of GDP values, and the numeric figures in parentheses suggest the count of countries or data points falling into each category. This data-driven segmentation aids in identifying the economic diversity across different regions, which is crucial for developing targeted economic or health interventions. The clarity of this tree allows for the extraction of actionable insights that inform policy decisions, resource allocation, and strategic market targeting.

## Discussion

The discussion surrounding the classification models for the 'Adult Mortality Rate (2019-2021)' dataset reveals a spectrum of predictive capabilities. The simplicity of the OneR model, yielding 67.31% accuracy, offers a valuable starting point for analysis despite its limited precision. The Naive Bayes and IBk models show stronger performance, achieving 87.82% and 86.54%

accuracy respectively, indicating their ability to capture essential patterns and similarities among instances for more nuanced classifications.

Dominating the results, the Decision Tree model, with its 98.08% accuracy, stands out for its ability to interpret and predict with high precision, highlighting the significant potential of using decision trees in complex data environments. Its success illustrates the importance of choosing the right model based on the specific requirements and nature of the dataset at hand. The clarity and depth provided by the Decision Tree model, in particular, underscore its value in extracting detailed and strategic insights for decision-making processes. These results collectively demonstrate the effectiveness of different data mining techniques and emphasize the importance of model selection in the field of predictive analytics.

**Classification Performance Table:**

| *Models* | Cross-Validation (10) |
|---|---|
| *One R* | 67.3077% |
| *Naive Bayes* | 87.8205% |
| *IBK* | 86.5385% |
| *Decision Tree* | 98.0769% |

The Classification Performance Table presents the comparative effectiveness of four distinct models as per their cross-validation accuracy. The One R model, with a predictive accuracy of 67.3077%, appears to be the least effective among the tested models. Naive Bayes and IBK models show improved performance with accuracies of 87.8205% and 86.5385%, respectively, indicating a substantial increase in predictive reliability. Notably, the Decision Tree model outperforms the others with a remarkable accuracy of 98.0769%. This high level of accuracy suggests that the Decision Tree model is significantly more adept at handling the dataset for the given task. Its success may be attributed to its ability to capture intricate patterns within the data,

leading to more accurate predictions during the cross-validation process. The substantial margin by which the Decision Tree model leads implicates its potential as a robust tool for predictive analytics in the context of this study.

# Conclusion

In concluding the evaluation of classification models for the 'Adult Mortality Rate (2019-2021)' dataset, our analysis spans from the OneR model's basic rule-based approach to the intricate decision-making framework of the Decision Tree algorithm. The OneR model's simplicity, paired with an accuracy of 67.3077%, provides an accessible entry point for data analysis, albeit with limited detail. Naive Bayes and IBk models offer more robust predictive power with accuracies of 87.8205% and 86.5385%, harnessing probabilistic foundations and instance-based learning to reveal key data patterns. However, the Decision Tree model, with an outstanding 98.0769% accuracy, clearly excels, demonstrating a high degree of predictive accuracy and interpretability. It validates the Decision Tree's utility as a tool for complex data analysis, cementing its role in extracting strategic insights. This breadth of methodologies underscores the versatility of classification algorithms in data mining and the critical importance of aligning model selection with the nuanced characteristics of the dataset.

**Problems and Justifications**

1. **Data Complexity and Volume**

One of the significant challenges faced during this project was managing the complexity and volume of the data. The dataset included detailed demographic, economic, and health indicators from multiple countries over a span of three years, making it inherently complex. To address this issue, we used advanced clustering algorithms capable of handling large datasets efficiently, ensuring robustness in our analysis without compromising on the processing speed.

2. **Algorithm Selection**

Selecting the appropriate data mining algorithms was a pivotal challenge due to the strategic implications of the project. The need to balance accuracy with computational efficiency led us to opt for k-means clustering, known for its effectiveness in large datasets. The decision was justified as k-means provided the needed scalability and was able to segment the data into meaningful clusters that were instrumental for strategic decision-making.

3. **No Missing Data or Outliers**

A unique aspect of our dataset was the absence of missing data and outliers, which is unusual in large-scale data collections. This anomaly raised concerns about the potential pre-processing done by data providers, which might have included over-smoothing or error corrections that could bias the results. We addressed this by reviewing the data collection and cleaning processes with the providers to ensure that the integrity of the data was maintained. Confirming the data quality upfront allowed us to proceed with confidence in our analytical processes.

4. **Justification for Not Using More Complex Models**

We considered the application of more complex machine learning models like neural networks for segmentation. However, given the clarity in clustering patterns identified by the k-means algorithm and the project's time constraints, we decided against employing more computationally intensive models. This decision was justified as the k-means algorithm sufficiently met the project's objectives, providing clear and actionable segmentation with considerable computational efficiency.

These problem-solving steps were crucial in overcoming the challenges encountered during the project, ensuring the delivery of reliable and actionable insights for strategic decision-making in global health insurance.

**Conclusion**

The insights gained from this project underscore the transformative potential of data mining in the global health insurance industry. The application of k-means clustering and various classification algorithms has not only enabled the identification of distinct health and economic

profiles but also facilitated a deeper understanding of the intricate relationship between demographic factors and mortality rates. These findings have practical implications for insurance providers, allowing for more precise risk assessments and the development of insurance products that are better tailored to meet the specific needs of different regions.

Moreover, the project demonstrates the critical role of data preprocessing and the strategic selection of analytical methods in ensuring the accuracy and relevance of the results. The successful segmentation of countries based on health outcomes and economic status has paved the way for targeted marketing strategies and product customization, enhancing customer satisfaction and competitive positioning in the market.

In conclusion, this study represents a significant step forward in the use of data mining for strategic decision-making in health insurance. It not only provides valuable insights into global health trends but also sets a precedent for future research in this crucial area of business analytics. The methodologies and findings detailed in this report could serve as a blueprint for other companies seeking to harness the power of data to refine their operations and strategies in the ever-evolving landscape of global health insurance.

**References**

Burges, C. J. C. (1998). A tutorial on support vector machines for pattern recognition. *Data Mining and Knowledge Discovery, 2*(2), 121-167. https://doi.org/10.1023/A:1009715923555

Han, J., Pei, J., & Kamber, M. (2011). Data Mining: Concepts and Techniques (3rd ed.). Morgan Kaufmann.

Han, J., Pei, J., & Kamber, M. (2011). *Data mining: Concepts and techniques* (3rd ed.). Morgan Kaufmann Publishers.

Holte, R. C. (1993). Very simple classification rules perform well on most commonly used datasets. *Machine Learning, 11*(1), 63-90. https://doi.org/10.1023/A:1022631118932

International Monetary Fund. (2021). World Economic Outlook Reports. Retrieved from https://www.imf.org/en/Publications/WEO

Jain, A. K. (2010). Data clustering: 50 years beyond K-means. *Pattern Recognition Letters, 31*(8), 651-666.

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). An Introduction to Statistical Learning: with Applications in R. Springer.

Ketchen, D. J., & Shook, C. L. (1996). The application of cluster analysis in strategic management research: An analysis and critique. *Strategic Management Journal, 17*(6), 441-458.

Kuhn, M., & Johnson, K. (2013). Applied Predictive Modeling. Springer.

Loh, W. Y. (2011). Classification and regression trees. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 1*(1), 14-23. https://doi.org/10.1002/widm.8

Mikhail, (n.d.). Adult mortality rate (2019-2021) [Data set]. Kaggle. Retrieved April 15, 2024, from https://www.kaggle.com/datasets/mikhail1681/adult-mortality-rate-2019-2021/data?select=Adult+mortality+rate+%282019-2021%29.csv

Rajulton, F., Ravanera, Z. R., & Beaujot, R. (2007). Measuring social cohesion: An experiment using the Canadian National Survey of Giving, Volunteering, and Participating. Social Indicators Research, 80(2), 461-492.

Wang, H., Abbas, K. M., Abbasifard, M., Abbasi-Kangevari, M., Abbastabar, H., Abd-Allah, F., ... & Murray, C. J. L. (2020). Global age-sex-specific fertility, mortality, healthy life expectancy (HALE), and population estimates in 204 countries and territories, 1950–2019: a comprehensive demographic analysis for the Global Burden of Disease Study 2019. The Lancet, 396(10258), 1160-1203.

World Health Organization. (2021). World Health Statistics 2021: Monitoring health for the SDGs. WHO.

Zhang, H. (2004). The optimality of Naive Bayes. *Proceedings of the Seventeenth International Florida Artificial Intelligence Research Society Conference*, 562-567. Retrieved from http://www.aaai.org/Papers/FLAIRS/2004/Flairs04-097.pdf