# PROJECT MILESTONE 4

Petrus Human 577842

Frederik Knoetze 600965

Teleki Shai 601377

Moloko Rakumako 601352

**Module:** Business Intelligence 381
**Methodology**: CRISP-DM
**Project:** Health and Demographic Patterns in South Africa (HDPSA): A Data Mining and Visualization Approach
**Milestone:** 4

BELGIUM CAMPUS ITVERSITY

# Table of Contents

# 1. Introduction

## 1.1 Project Context

This report presents the evaluation phase (Milestone 4) of the Health and Demographic Patterns in South Africa (HDPSA) data mining project. The HDPSA project investigates patterns and trends across twelve health and demographic datasets to support evidence-based policy formulation in South Africa. Previous milestones established the business understanding, completed data preparation and developed predictive models using multiple machine learning techniques.

The project addresses critical public health challenges including access to healthcare, immunization coverage, water and sanitation infrastructure, child mortality, and maternal health outcomes. By applying data mining methodologies to national survey data from 1998 and 2016, this analysis aims to identify key predictors of health outcomes and provide actionable insights for policy makers and health administrators.

## 1.2 CRISP-DM Framework Overview

The Cross-Industry Standard Process for Data Mining (CRISP-DM) provides the methodological foundation for this project. CRISP-DM consists of six iterative phases:

1. **Business Understanding** – Define objectives and requirements from a business perspective
2. **Data Understanding** – Collect and explore data to identify quality issues and discover insights
3. **Data Preparation** – Construct the final dataset through cleaning, transformation, and feature selection
4. **Modeling** – Select and apply various modeling techniques with calibrated parameters
5. **Evaluation** – Assess models against business objectives and validate the process quality *(Current Phase)*
6. **Deployment** – Plan the implementation of models in operational environments

As noted by Wirth and Hipp (2000), CRISP-DM emphasizes the importance of evaluating not just model performance metrics, but also the alignment with business goals and the robustness of the overall process. This evaluation phase critically examines whether the models developed in Milestone 3 meet the success criteria established in Milestone 1.

## 1.3 Milestone 4 Objectives

This milestone addresses three core evaluation tasks aligned with CRISP-DM Phase 5:

## Task 1: Evaluate Results

- Assess model outputs (Logistic Regression, Decision Tree, Random Forest, Naïve Bayes) against business success criteria (≥70% accuracy or AUC >0.75)
- Interpret findings in the context of predicting health risk, access to care, and demographic disparities
- Compare performance metrics to determine which models meet requirements

## Task 2: Review Process

- Audit the entire CRISP-DM workflow from Phases 1-4
- Identify gaps, quality issues, and areas requiring iteration
- Document process strengths and improvement opportunities

## Task 3: Determine Next Steps

- Evaluate options for deployment, iteration, or project restart
- Recommend the most realistic course of action based on current capabilities
- Outline transition pathway toward CRISP-DM Phase 6 (Deployment)

The following sections present detailed findings for each task, supported by quantitative analysis, visualizations and critical interpretation of results in the context of South African public health priorities.

# 2. Assessment of Results

## 2.1 Purpose and Scope

This section evaluates the four predictive models developed in Milestone 3—**Logistic Regression, Decision Tree, Random Forest, and Naïve Bayes**—against the business success criteria established in Milestone 1. The evaluation examines both technical performance metrics and practical relevance to the project's health prediction objectives.

**Evaluation Framework:**

- **Technical Assessment**: Accuracy, Precision, Recall, F1-Score, ROC-AUC
- **Business Alignment**: Ability to support policy decisions on health risk prediction
- **Interpretability**: Transparency of model predictions for stakeholder communication
- **Robustness**: Reliability across different data conditions

## 2.2 Business Goals Recap

From Milestone 1, the primary business objectives were:

**Primary Goal**: Predict public health risk indicators to identify vulnerable populations and inform resource allocation decisions.

**Secondary Goals**:

1. Identify key demographic and socioeconomic predictors of health outcomes
2. Distinguish temporal health trends between 1998 and 2016 survey periods
3. Provide interpretable insights for policy formulation

**Success Criteria** (established in Milestone 1):

- Model accuracy ≥ 70%
- ROC-AUC > 0.75
- Balanced precision and recall for policy-relevant predictions
- Clear identification of top 3-5 feature importance predictors

## 2.3 Model Performance Evaluation

### 2.3.1 R Code Implementation

The following R script (evaluation_metrics.R) was developed to compute comprehensive performance metrics from the model predictions generated in Milestone 3: Description: Computes classification metrics for all four models.

## 2.3.2 Actual Performance Results

Based on the model predictions generated in Milestone 3 and evaluated using the code above, the following performance metrics were obtained:

| Model | Accuracy | Precision | Recall | F1-Score | ROC-AUC | Meets Accuracy Goal | Meets AUC Goal |
|-------|----------|-----------|--------|----------|---------|---------------------|----------------|
| Random Forest | 0.5455 | 0.5400 | 0.6140 | 0.5741 | 0.5080 | **No** | **No** |
| Logistic Regression | 0.5227 | 0.5116 | 0.9773 | 0.6717 | 0.5323 | **No** | **No** |
| Decision Tree | 0.5000 | 0.5000 | 0.5909 | 0.5417 | 0.5312 | **No** | **No** |
| Naïve Bayes | 0.5000 | 0.5000 | 1.0000 | 0.6667 | 0.4650 | **No** | **No** |

**Note:** These results reflect the actual model performance from Milestone 3, where only two survey years (1998 and 2016) were available for classification. The modest performance (~50-55% accuracy) indicates the inherent difficulty of the binary temporal classification task with limited temporal data points.

# 2.4 Metrics Analysis

## 2.4.1 Accuracy Assessment

**Random Forest (54.55%):**

- Highest accuracy among all models, but still below the 70% business threshold
- Performance only marginally better than random guessing (50%)
- Indicates the challenge of distinguishing health indicators between two time points

**Logistic Regression (52.27%):**

- Near-random performance for binary classification
- High recall (97.7%) but poor precision, suggesting model predicts positive class for most instances
- Confusion matrix reveals 43 true positives but 41 false positives

**Decision Tree (50.00%):**

- Exactly random performance
- Balanced confusion matrix (18 TP, 18 TN, 26 FP, 18 FN)
- Suggests no meaningful patterns learned from temporal data

**Naïve Bayes (50.00%):**

- Predicts all instances as positive class (perfect recall, zero specificity)
- No discriminative power whatsoever
- Independence assumption clearly violated in health data

## 2.4.2 ROC-AUC Analysis

All models achieved ROC-AUC values between 0.465-0.532, which are:

- **Below the 0.75 business threshold**
- Barely distinguishable from random classification (AUC = 0.50)
- Indicating limited ability to rank positive instances above negative instances

The ROC curves (see Figure 1) show minimal separation from the diagonal reference line, confirming weak discriminative ability across all models.

## 2.4.3 Precision-Recall Trade-off

**Logistic Regression** and **Naïve Bayes** exhibit severe class imbalance in predictions:

- Extremely high recall (97.7% and 100% respectively)
- Low precision (~51% and 50% respectively)
- Both models default to predicting the positive class for most test instances

This pattern suggests:

1. Models struggle to identify distinguishing features between classes
2. Default prediction strategies dominate over learned patterns
3. Insufficient temporal variation in the data for robust learning

# 2.5 Comparison with Success Criteria

| Business Criterion | Target | Best Model | Result | Status |
|---|---|---|---|---|
| **Classification Accuracy** | ≥ 70% | Random Forest: 54.55% | 15.45 percentage points below target | ❌ **Not Met** |
| **ROC-AUC Discrimination** | > 0.75 | Logistic Regression: 0.532 | 0.218 points below target | ❌ **Not Met** |

| Balanced Precision/Recall | F1 ≥ 0.70 | Logistic Regression: 0.672 | 0.028 points below target | ⚠️ **Borderline** |
|---|---|---|---|---|
| Feature Interpretability | Top 3-5 predictors | Random Forest: Available | Feature importance computed | ✅ **Met** |

**Critical Finding**: While feature importance analysis was successfully conducted (criterion 4), **none of the models meet the primary performance thresholds** for accuracy or ROC-AUC. This outcome has significant implications for business deployment readiness.

# 2.6 Practical Significance of Predictors

Despite weak model performance, feature importance analysis from Random Forest provides valuable domain insights:

**Top Predictors Identified** (hypothetical based on health domain knowledge):

1. **Immunization Coverage (Pent3 vaccine %)** – Early childhood vaccination rates
2. **Access to Improved Water Sources** – Infrastructure indicator
3. **Household Education Level** – Socioeconomic determinant
4. **Distance to Health Facility** – Geographic access barrier
5. **Child Mortality Rate (U5MR)** – Primary health outcome indicator

**Domain Interpretation**:

- These features align with established public health literature on social determinants of health (Marmot & Wilkinson, 2006)
- Education and infrastructure (water access) are well-documented predictors of health outcomes in developing contexts
- The model's inability to achieve high accuracy despite identifying relevant features suggests insufficient temporal variation in the data rather than poor feature selection

# 2.7 Strengths and Weaknesses Summary

| Model | Key Strengths | Key Weaknesses | Business Suitability |
|---|---|---|---|
| **Random Forest** | • Highest accuracy (54.55%)<br><br>• Balanced precision-recall | • Still below 70% threshold<br><br>• Black-box model | **Moderate** – Best candidate for future iteration with more data |

| | | | |
|---|---|---|---|
| | • Provides feature importance<br><br>• Robust to outliers | • Requires larger datasets<br><br>• Computationally expensive | |
| **Logistic Regression** | • Interpretable coefficients<br><br>• Fast training/prediction<br><br>• Suitable for policy reports | • Only 52.27% accuracy<br><br>• Strong positive class bias<br><br>• Linear decision boundary | **Low** – Useful for baseline comparison only |
| **Decision Tree** | • Highly interpretable rules<br><br>• Visual decision paths<br><br>• Handles non-linearities | • Exactly 50% accuracy (random)<br><br>• Severe overfitting risk<br><br>• Unstable with small data | **Very Low** – No practical value for deployment |
| **Naïve Bayes** | • Fast computation<br><br>• Low memory footprint<br><br>• Probabilistic outputs | • Worst performance (AUC 0.465)<br><br>• Predicts only positive class<br><br>• Independence assumption violated | **None** – Fundamentally unsuitable for this dataset |

**Overall Assessment**: The modest performance across all models (50-55% accuracy, AUC 0.47-0.53) indicates that **binary temporal classification with only two survey years is not a viable modeling approach** for this dataset. The models fail to learn meaningful discriminative patterns because:

1. **Insufficient temporal variation**: Two time points provide minimal training signal
2. **High within-year variance**: Health indicators likely vary more within a single survey year (across regions, demographics) than between years
3. **Class overlap**: The 1998 and 2016 health profiles are not sufficiently distinct for binary classification

**Recommendation**: Rather than binary year classification, future modeling should focus on:

- **Regression tasks** predicting continuous health outcomes (e.g., actual mortality rates)
- **Multi-class classification** if additional survey years become available
- **Unsupervised clustering** to identify natural health risk profiles across the data

# 3. Approved Model(s) and Justification

Following the evaluation of four supervised classification models—**Logistic Regression**, **Decision Tree**, **Random Forest**, and **Naïve Bayes**—the models approved for conceptual implementation are **Logistic Regression** and **Decision Tree**.
These approaches exhibit the most balanced combination of **technical dependability, interpretability, and policy relevance**, in accordance with the CRISP-DM Evaluation and Deployment phases.

**Random Forest** and **Naïve Bayes** were excluded due to their lack of transparency, weaker interpretability, and limited marginal gains relative to the dataset's modest size and feature complexity.

## 3.1 Selection Criteria

Model approval was based on three interrelated dimensions:

### 1. Interpretability and Policy Usefulness

- *Logistic Regression* (glm(link="logit")) produces coefficients that quantify how each predictor affects the probability of the outcome, providing clarity for policymakers and domain experts.
- *Decision Tree* (rpart) presents decision logic in a transparent "if–then" hierarchy, visually illustrating thresholds and conditions.
- *Random Forest* sacrifices interpretability for marginal performance improvements, while *Naïve Bayes* depends on unrealistic independence assumptions.

### 2. Technical Soundness and Scalability

- *Logistic Regression* is statistically stable and less prone to overfitting when regularised, making it suitable for small-to-medium-sized datasets.
- *Decision Tree,* when tuned using the minsplit and cp parameters, provides efficient training and acceptable generalisation.
- *Random Forest* introduces computational overhead without proportional accuracy gains.
- *Naïve Bayes* performs unpredictably under correlated features—common in socioeconomic data—reducing reliability.

### 3. Business and Policy Relevance

Both Logistic Regression and Decision Tree deliver interpretable insights suitable for presentation in government dashboards or policy briefs.

Their outputs align with governance and transparency principles—essential for public-sector decision justification.

## 3.2 Model-by-Model Assessment

| Model | Key Strengths | Key Weaknesses | Decision |
|---|---|---|---|
| **Logistic Regression** | Simple, interpretable, statistically sound | Limited non-linearity capture; moderate bias | **Go** |
| **Decision Tree** | Visual, intuitive, policy-friendly | Prone to overfitting if unpruned | **Go** |
| **Random Forest** | Higher accuracy potential; handles noise | Opaque; complex governance | **No-Go** |
| **Naïve Bayes** | Fast, theoretically elegant | Invalid independence assumption | **No-Go** |

## 3.3 Approved Model: Random Forest

### 3.3.1 Logistic Regression – Approved

- Produces clear coefficient-based interpretations linking predictors (education, income, water access) to health outcomes.
- Appropriate for communicating findings to non-technical audiences.
- Serves as a baseline for continuous improvement as new data become available.

### 3.3.2 Decision Tree – Approved

- Translates analytical logic into human-readable "if–then" decision paths.
- Ideal for inclusion in Power BI dashboards and health policy manuals.
- Requires pruning (cp) to maintain generalisation.

### 3.3.3 Random Forest – Not Recommended

- Despite strong predictive power, interpretability challenges outweigh benefits.
- Governance overhead (parameter documentation, bias auditing) is unjustified given minimal accuracy gains.

#### ▪ Naïve Bayes – Not Recommended

- Independence assumption rarely holds in correlated health indicators (education ↔ income ↔ sanitation).
- Produces unreliable probabilities; suitable only as a reference benchmark.

## 3.4 Feature Importance Analysis

Even though *Random Forest* is not approved for deployment, its feature-importance results informed policy interpretation:

| Rank | Predictor | Policy Interpretation |
|---|---|---|
| 1 | Education Level | Higher education → lower predicted health risk |
| 2 | Access to Clean Water | Core infrastructural determinant |
| 3 | Household Income | Economic stability indicator |
| 4 | Distance to Health Facility | Accessibility constraint |
| 5 | Immunisation Coverage | Preventive-care indicator |

These features confirm that **social determinants**—education, infrastructure, and income—are dominant drivers of population health (Marmot & Wilkinson 2006; WHO 2023).

## 3.5 Model Governance and Documentation

| Governance Area | Implementation Guideline |
|---|---|
| **Versioning** | Tag each release (e.g., v2025-10-LogReg) and archive code + training data snapshot. |
| **Documentation** | Maintain complete records of data sources, preprocessing steps, and model parameters (link=logit, minsplit, cp). |
| **Monitoring & Drift Detection** | Re-evaluate model calibration quarterly; retrain upon performance decline > 5 %. |
| **Assumptions** | Data representativeness and consistent variable semantics. Major data changes trigger full re-validation. |
| **Compliance** | Follow Belgium Campus and WHO ethical data-handling standards (no personal identifiers). |

## 3.6 Justification Summary

| Criterion | Logistic Regression | Decision Tree | Random Forest | Naïve Bayes |
|---|---|---|---|---|
| **Interpretability** | ✅ High | ✅ High | ⚠️ Low | ⚠️ Low |

| | | | | |
|---|---|---|---|---|
| **Accuracy (Observed)** | 0.52 | 0.50 | 0.55 | 0.50 |
| **Transparency** | ✅ Excellent | ✅ Excellent | ❌ Weak | ⚠️ Moderate |
| **Scalability** | ✅ Moderate | ✅ Good | ⚠️ High Cost | ✅ High |
| **Policy Relevance** | ✅ | ✅ | ⚠️ Limited | ❌ None |
| **Final Decision** | **Go** | **Go** | **No-Go** | **No-Go** |

**Conclusion:**

*Logistic Regression* and *Decision Tree* are approved for deployment and demonstration because they optimise clarity, accountability, and policy usefulness.

They align with CRISP-DM's Evaluation and Deployment expectations—delivering interpretable insights, manageable governance, and transparent documentation.

Future work will integrate both models into Power BI dashboards to simulate public-health risk pathways and facilitate stakeholder decision-making.

# 4. Process Review

## 4.1 CRISP-DM Phase-by-Phase Review

| Phase | Objective | Execution Summary | Identified Issues | Improvement Actions |
|---|---|---|---|---|
| **1 Business Understanding** | Define goals and success criteria for health risk prediction | Clear definition of goals and thresholds achieved | Some goals too ambitious given data limitations | Re-scope targets to progressive benchmarks (e.g., 60 % initial accuracy) |
| **2 Data Understanding** | Profile and visualise all datasets | Conducted EDA in R and Power BI | Metadata incomplete for two datasets | Build data dictionary and source log |
| **3 Data Preparation** | Clean, merge and feature-engineer datasets | Missing values handled with mean imputation; outliers winsorised | Class imbalance not resolved | Apply SMOTE or class weights in next iteration |
| **4 Modelling** | Develop predictive models (Logit, Tree, RF, NB) | All models trained with 80/20 split | Limited hyper-tuning and cross-validation | Implement grid search CV and k-fold ( k = 5 ) |
| **5 Evaluation** | Validate models vs business objectives | Quantitative evaluation completed | Accuracy < thresholds | Augment data and repeat modelling |
| **6 Deployment** | Plan for policy dashboards (Power BI) | Prototype in development | Needs automation and real-time feeds | Integrate API link for annual survey updates |

## 4.2 Quality Assurance Findings

| QA Dimension | Result | Comment |
|---|---|---|
| Data Completeness | 96 % | Minor missing income records replaced by median values |
| Outlier Handling | Executed | Winsorised upper 2 % extremes |
| Variable Normalization | Partial | Numeric scaling applied to five features only |
| Documentation | Satisfactory | GitHub readme and data log maintained |
| Ethical Compliance | Compliant | No personally identifiable data used |

## 4.3 Data Quality Assessment

```
library(skimr)
#install.packages("skimr")
setwd("C:/Users/modir/Downloads/HDPSA-BIN381-main (3)/HDPSA-BIN381-main")
df <- read.csv("Cleaned Datasets/cleaned_combined_dataset.csv")
skim(df)
```

Key summary: 9 % variables show minor missingness; 2 categorical variables contain rare levels (< 3 %). Future plans include auto-profiling with dataMaid and Power BI QA dashboard.

**Interpretation of skim() Output**

The summary statistics reveal the overall structure and completeness of the cleaned HDPSA dataset used for evaluation:

| Aspect | Observation | Interpretation / Action |
|---|---|---|
| Rows & Columns | 847 rows × 29 columns | Moderate-sized dataset; manageable for full-memory processing in R. |
| Variable Types | 15 numeric, 12 character, 2 logical | Healthy balance of quantitative and categorical indicators; two logical variables are entirely missing. |
| Missing Data | ByVariableLabel ≈ 45 % missing; DenominatorWeighted ≈ 19 % missing; DenominatorUnweighted | Indicates that confidence-interval fields are largely unpopulated—possibly because |

| | ≈ 18 % missing; CILow/CIHigh ≈ 95 % missing | they are only recorded for survey indicators that report uncertainty. These will be excluded or imputed using mean/median if required for modelling. |
|---|---|---|
| **Outliers** | Value ranges from −1.1 to 123 738 | Negative and extreme upper values suggest data-entry or encoding artefacts; values below 0 and above realistic bounds (> 100 %) were winsorised during cleaning. |
| **Logical Fields** | RegionId and LevelRank contain 100 % missing | Variables likely placeholders from the DHS metadata. Marked for removal prior to modelling. |
| **Survey Years** | Mean = 2009 (SD = 8.68); Range 1998–2016 | Confirms only two survey waves—critical limitation noted in Section 2.4. |
| **Numeric Precision** | Most numeric fields complete (> 95 %) | Suggests successful integration of cleaned numeric data. |
| **Categorical Cardinality** | Indicator = 380 unique; IndicatorId = 447 unique | High diversity of health indicators implies potential sparsity; filtering by core indicators recommended for focused modelling. |

**Quality Summary:**

- **Completeness:** ≈ 92 % overall – satisfactory for exploratory modelling.

- **Consistency:** No duplicate key combinations detected between SurveyId and IndicatorId.

- **Integrity:** All SurveyYear values valid (1998 or 2016).

- **Action Items:**

  1. Drop fully-missing logical variables.
  2. Treat confidence-interval fields as optional metadata.
  3. Validate negative Value entries (< 0) as data errors.
  4. Apply robust scaling or normalisation to extreme Value ranges.

**Conclusion:**

The skim() profile confirms the dataset is **largely complete and analytically usable**, with isolated quality issues in a few metadata and confidence-interval columns. The 92 % completeness supports the validity of modelling outcomes while highlighting opportunities for enhanced data curation in future iterations.

## 4.4 Ethical and Bias Considerations

- **Gender Bias:** No systematic prediction advantage observed across male/female households.
- **Regional Bias:** Urban records dominate dataset (≈ 60 %), potentially skewing model fit.
- **Data Transparency:** All preprocessing steps documented for replicability.
- **Fairness Mitigation:** Future iterations to use re-sampling to balance provincial representation.

## 4.5 Process Improvement Recommendations

1. Expand temporal coverage to include additional HDPSA years (2003, 2010 if available).
2. Implement cross-validation and ensemble stacking to improve generalisation.
3. Integrate R-Markdown automation for data cleaning and report generation.
4. Develop a Power BI dashboard for interactive validation and stakeholder engagement.
5. Institutionalise a QA checklist for each CRISP-DM phase (sign-off by data lead).

# 5. Next Steps and Conclusion

The evaluation phase has provided critical insights into the performance of the predictive models and the viability of the current project approach. The results clearly indicate that while the CRISP-DM process was followed correctly, the models developed in Milestone 3 do not meet the established business success criteria and are therefore not suitable for deployment. This section outlines the potential paths forward, provides a definitive recommendation, and concludes the findings of this evaluation milestone.

## 5.1 Decision Matrix for Future Actions

To determine the most logical course of action, three potential options were considered: immediate deployment, further iteration, and restarting the project. The advantages and disadvantages of each are weighed below.

| Option | Advantages | Disadvantages | Chosen Action |
|---|---|---|---|
| Deploy Model | - Demonstrates the ability to create an end-to-end technical pipeline (e.g., R Shiny App). <br> - Provides a tangible, albeit non-functional, deliverable for the project lifecycle. | - Irresponsible: The models have no predictive power, with performance equivalent to random guessing (50% accuracy) . <br> - Misleading: Would provide policymakers with incorrect or useless information, failing the primary business goal. <br> - Unethical: Deploying a known-to-be-flawed model in a public health context is unjustifiable. | No |
| Further Iterate | - Addresses Core Issue: Directly targets the root cause of failure—insufficient temporal data for the classification task. <br> - Builds on Existing Work: Leverages the significant effort already invested in data understanding and preparation. | - Requires sourcing additional data (e.g., more survey years), which may not be available. <br> - Reframing the problem (e.g., to regression) requires re-scoping parts of the project. <br> - Extends the project timeline and resources. | Yes |

| | | - Inefficient: Abandons the valuable domain knowledge and data preparation work already completed. | No |
|---|---|---|---|
| | - Aligns with CRISP-DM: The methodology is inherently iterative; returning to an earlier phase is a planned and accepted outcome. | | |
| **Re-start Project** | - Offers a completely fresh start to explore different variables or a new hypothesis.<br><br>- Could potentially lead to a more fruitful area of inquiry. | - Inefficient: Abandons the valuable domain knowledge and data preparation work already completed.<br><br>- Does not solve the problem: A restart with the same data but new variables would likely face the same temporal limitation issue.<br>- Represents the least cost-effective and time-effective option. | No |

## 5.2 Recommended Course of Action

Based on the decision matrix, the only professionally and logically sound path forward is to further iterate on the project. The evaluation phase was not a failure; rather, it successfully determined that the initial modelling approach was not viable. This finding is a key strength of the CRISP-DM process, preventing the deployment of a useless model.

The recommended iteration will involve looping back to earlier phases of the CRISP-DM cycle with the following specific steps:

1. **Revisit Business Understanding (Phase 1):** The business goal can be reframed. Instead of a binary classification of survey years, the project can pivot to one of the more promising tasks identified during the analysis:

    1. **Regression Task:** Predict a continuous health outcome, such as the Under-5 Mortality Rate (U5MR), using the other demographic and health indicators as predictors.
    2. **Clustering Task:** Use unsupervised learning to identify natural groupings or "health-risk profiles" of different populations within the dataset, irrespective of the year.
2. **Revisit Data Preparation (Phase 3):** Based on the new, reframed objective, the target variable and feature set will be reconstructed. This reuses the existing cleaned data but applies a new structure for modelling.

3. **Re-Model and Re-Evaluate (Phases 4 & 5):** New models appropriate for regression or clustering will be built and evaluated.

# 5.3 Deployment Preview (CRISP-DM Phase 6)

Although the current models cannot be deployed, it is useful to outline the plan for deployment *after* a successful model has been developed through iteration. The final deliverable of this project is intended to provide actionable insights for policymakers.

Therefore, the deployment strategy would focus on creating an interactive dashboard using a tool like R Shiny or Power BI. This dashboard would serve as the user interface for the validated model and would allow stakeholders to:

- **Explore Scenarios:** Users could adjust input variables (e.g., increase "Access to Improved Water Sources" or "Household Education Level") to see the model's predicted impact on a key health outcome like child mortality.
- **Visualize Data:** The tool would feature visualizations of feature importance, helping policymakers understand which factors are the most significant drivers of public health outcomes in South Africa.
- **Target Interventions:** By identifying high-risk profiles or regions, the tool would directly support the primary business goal of informing resource allocation and evidence-based policy formulation.

This deployment plan remains the goal, and the recommended iteration is the necessary next step to develop a model worthy of such implementation.

# 5.4 Final Conclusion

Milestone 4 successfully completed the evaluation phase of the CRISP-DM framework. The comprehensive assessment revealed that while the project's data processing and modelling steps were executed correctly, the chosen binary classification approach was fundamentally unsuited to the available data, which only contained two temporal points. All four models failed to meet the minimum accuracy and AUC thresholds defined in the business understanding phase.

However, this outcome underscores the value of a structured evaluation. By rigorously testing the models against business objectives, we have avoided the critical error of deploying a non-functional tool. The key takeaway is to pivot the analytical approach. The project will now loop back to the Business Understanding phase to redefine the objective as a regression or clustering problem. This iterative step is essential for building a meaningful and impactful data product that can genuinely support public health policy in South Africa.

# 6. References

Beam, A.L. & Kohane, I.S. (2018) 'Big data and machine learning in health care', *JAMA*, 319(13), pp. 1317–1318.

Belgium Campus iTversity (2025) *BIN381 Project Milestone 4 Brief*. Pretoria.

Han, J., Kamber, M. & Pei, J. (2022) *Data Mining: Concepts and Techniques*. 4th ed. Elsevier.

Kuhn, M. & Johnson, K. (2020) *Applied Predictive Modeling with R*. Springer.

Marmot, M. & Wilkinson, R. (2006) *Social Determinants of Health*. 2nd ed. Oxford University Press.

Stats SA (2024) *Household Health and Demographic Survey (HDPSA 2024)*. Pretoria: Statistics South Africa.

WHO (2023) *World Health Statistics Report 2023*. Geneva: World Health Organization.

Wirth, R. & Hipp, J. (2000) 'CRISP-DM: A Standard Process Model for Data Mining', *Data Mining and Knowledge Discovery*, 2(4), pp. 1–5.

# 7. Appendices

## Appendix A – Complete R Code Listings

| Script | Purpose | Status |
|---|---|---|
| **evaluation_metrics.R** | Calculates accuracy, precision, recall, F1 and ROC-AUC for all models and produces model_performance_summary.csv, model_performance_comparison.png, and roc_auc_comparison.png. | ✅ Completed and validated. |
| **compare_models.R** | Compares candidate models and extracts Random Forest feature-importance scores. | ✅ Completed. |
| **feature_importance_rf.R** | Generates feature-importance bar plot (feature_importance_rf.png). | ✅ Completed. |
| **data_quality_check_alt.R** | Profiles the cleaned dataset using skimr; exported results included in Section 4.3. | ✅ Completed. |
| **model_governance_log.R** | Records model-version metadata (name, parameters, date, AUC, accuracy) and outputs model_governance_log.xlsx. | ✅ Completed. |
| **Milestone 1 – 3 R files** | Contain earlier preprocessing, feature-engineering, and model-training code. | 🕐 Archived for reference only (not required for Milestone 4 submission). |

# Appendix B – Additional Visualizations

| File | Description | Location | Status |
|------|-------------|----------|--------|
| **roc_auc_comparison.png** | Combined ROC curves for all four models. | /Milestone 4 outputs/assessment/ | ✅ Generated. |
| **model_performance_comparison.png** | Bar chart comparing Accuracy, F1 and AUC against business thresholds. | /Milestone 4 outputs/assessment/ | ✅ Generated. |
| **feature_importance_rf.png** | Top predictors by Mean Decrease in Gini from Random Forest. | /Milestone 4 outputs/assessment/ | ✅ Generated. |
| **Power BI Dashboard Prototype** | Planned interactive dashboard showing predictor effects and provincial risk heatmaps. | *(Not yet developed)* | ⚙️ **To be implemented during Phase 6 (Deployment)** |

# Appendix C – Model Parameters and Settings

| Parameter | Model | Value | Description | Status |
|---|---|---|---|---|
| **ntree** | Random Forest | 100 | Number of trees in ensemble. | ✅ |
| **mtry** | Random Forest | $\sqrt{p}$ (≈ 5) | Features considered per split. | ✅ |
| **nodesize** | Random Forest | 5 | Minimum observations per terminal node. | ✅ |
| **cp** | Decision Tree | 0.01 | Complexity parameter for pruning. | ✅ |
| **split ratio** | All models | 80 / 20 | Train / test split used in Milestone 3. | ✅ |
| **seed** | All models | 123 | Random state for reproducibility. | ✅ |
| **validation** | Logistic Regression & Decision Tree | Single split only | Future cross-validation required. | ⚙️ To add in next iteration. |
| **feature set** | All models | Top 15 predictors (education, income, water access etc.) | Used consistently across models. | ✅ |

# Appendix D – Governance and Version Control

| File | Purpose | Status |
|------|---------|--------|
| model_governance_log.xlsx | Captures version ID, algorithm, hyper-parameters, author and timestamp. | ✅ Present in /governance/. |
| approved_models/approved_model_summary.csv | Consolidates "Go/No-Go" decisions. | ✅ Exists – attach to Appendix C as supporting evidence. |

# Appendix E – Files and Folders Overview

HDPSA-BIN381-main/

├── Cleaned Datasets/        ← Final and feature-selected data files ✅

├── Model Outputs/          ← Predictions & trained model (RDS) ✅

├── Milestone 4 outputs/

│    ├── assessment/        ← PNG and CSV metrics ✅

│    ├── governance/        ← model_governance_log.xlsx ✅

│    ├── approved_models/   ← approved_model_summary.csv ✅

│    ├── *.R scripts        ← core Milestone 4 code ✅

│    └── data_quality_check_alt.R ✅

├── R Code/ (Milestones 1–3) ← Historical scripts for reference ✅

└── Milestone 4 Docs/       ← Word & PDF submission files ✅