# Epidemiological Analysis of Tuberculosis in Rwanda: FY 2023-2024

A Comprehensive Report on Demographics, Clinical Characteristics, Outcomes, and Public Health Strategies

Prepared by: Seraphine Mukabugingo

Date: August 04, 2025

Rwanda National TB Program

# Abstract

This report provides an in-depth epidemiological analysis of tuberculosis (TB) in Rwanda for FY 2023-2024, based on a dataset of 8,549 cases across 30 districts. It covers demographics, clinical characteristics, high-risk groups (HRG), HIV co-infection, treatment outcomes, contact tracing, nutritional status, drug resistance, predictive modeling, health system performance, and special populations, with additional analyses of temporal trends, district-specific patterns, and advanced predictive metrics.

Key findings include a male-dominated burden (73.5%, 6,285 cases), urban concentration (37% in top five districts), moderate HIV co-infection (13.6%, 1,166 cases), and a treatment success rate of approximately 86% among evaluated cases (47.3% overall due to 45.2% unknown outcomes). Predictive models show fair performance for treatment success (Random Forest, AUC 0.652) and good performance for mortality (Logistic Regression, AUC 0.758), with BMI and HIV as top predictors.

Contact tracing achieves high screening rates (97.7% for <5 years, 99.3% for ≥5 years) but low TPT completion. Visualizations (AQSIMAGE1 to AQSIMAGE44, stored in the AQSimages folder in Overleaf) provide actionable insights, supported by WHO 2024 data (TB incidence 55 per 100,000, 95% CI: 41-71). Recommendations focus on intensifying urban screening, enhancing TB-HIV integration, improving TPT adherence, scaling nutrition support, and strengthening surveillance to meet 2025 End TB Strategy goals.

# Contents

# Chapter 1

# Introduction

## 1.1 Background and Context

Tuberculosis (TB) remains a critical public health challenge in Rwanda, despite significant progress toward the WHO End TB Strategy 2025 milestones. The WHO Global TB Report 2024 reports Rwanda's TB incidence at 55 per 100,000 population in 2023 (95% CI: 41-71), a 50% reduction from 2015 levels, positioning Rwanda as a regional leader in TB control.

This report analyzes a dataset of 8,549 TB notifications from FY 2023-2024 (July 2023–June 2024), covering 30 districts and 96 variables, including demographics (`age_group`, `sex`, `district`), clinical characteristics (`tb_classification_ds_or_dr`, `site_of_disease`), HIV status (`hiv_status`), treatment outcomes (`treatment_outcome`), contact tracing (`number_of_contacts_<5_years_screened_for_tb`), and nutritional status (`bmi_at_beginning`). All visualizations (AQSIMAGE1 to AQSIMAGE44) are stored in the `AQSimages` folder uploaded to Overleaf.

Rwanda's National TB Program leverages advanced diagnostics (76% GeneXpert coverage, 6,522 cases) and community health workers (CHWs) for contact tracing, achieving near-universal screening (97.7% for <5 years, 99.3% for ≥5 years). Challenges include high unknown outcomes (45.2%, 3,861 cases), near-zero culture testing (∼0%), and poor TPT completion (often 0), necessitating targeted interventions to sustain progress toward 2025 goals (50% incidence reduction, 75% mortality reduction from 2015).

## 1.2 Objectives

This report aims to:

- Provide a comprehensive description of TB epidemiology, including demographics, clinical features, geographic/temporal patterns, and district-specific trends.

- Identify high-risk groups (HRG) and HIV co-infection patterns, with detailed subgroup analyses.

- Evaluate treatment outcomes, associated factors, and predictive models, including advanced metrics (e.g., confusion matrices).

- Assess contact tracing effectiveness, TPT performance, and district-level yields.

- Analyze nutritional status, drug resistance, and health system performance, with additional outcome-specific trends.

- Develop evidence-based recommendations to strengthen TB control, aligned with WHO 2025 goals.

## 1.3 Data Sources and Methodology

The primary data source is `final_dataset.csv`, containing 8,549 TB cases with 96 variables. Supplementary sources include the WHO Global TB Report 2024 and Rwanda Ministry of Health (MoH) reports. Analyses were conducted using Python libraries:

- **Pandas**: Data aggregation (e.g., counts, percentages, cross-tabulations).

- **Plotly**: Interactive visualizations (bar, pie, line charts, tables).

- **Scikit-learn**: Machine learning (Random Forest, Logistic Regression).

- **Lifelines**: Survival analysis (Kaplan-Meier, placeholder due to limited time-to-event data).

Methods include descriptive statistics, chi-square tests (e.g., p=0.0097 for HIV vs. outcome), logistic regression (e.g., BMI coefficient 2.4508 for mortality), and cross-tabulations (e.g., HIV by district). Visualizations (AQSIMAGE1-44) are stored in the `AQSimages` folder in Overleaf.

## 1.4 Scope and Limitations

The dataset lacks population denominators, preventing absolute incidence rate calculations (relative rates used, e.g., Nyarugenge: 10.6%). High unknown outcomes (45.2%) inflate non-success rates, and near-zero culture testing ($\sim 0\%$) limits DR-TB detection. Time-to-event data is incomplete, leading to placeholder survival curves (AQSIMAGE33-34). Additional analyses (e.g., quarterly trends, district-specific yields) enhance depth. WHO 2024 data provides global context.

# Chapter 2

# Data Overview

The dataset includes 8,549 TB notifications from FY 2023-2024, covering 30 districts. Key variables include:

- **Demographics**: `age_group` (e.g., 25-34 years), `sex` (male/female), `district` (e.g., Nyarugenge), `month`.

- **Clinical**: `tb_classification_ds_or_dr` (DS-TB/DR-TB), `site_of_disease` (pulmonary/extra-pulmonary), `method_of_tb_confirmation` (bacteriological/clinical).

- **HIV**: `hiv_status` (Positive/Negative), `currently_on_art`, `currently_on_cotrimoxazole`.

- **Outcomes**: `treatment_outcome` (Cured, Completed, Died, etc.), `treatment_success` (Cured/Completed).

- **Contact Tracing**: `number_of_contacts_<5_years_screened_for_tb`, `contacts_of_tpb+<_2_yea`

- **Nutritional**: `bmi_at_beginning`, `bmi_at_end_treatment`, `tb_nutrition_support_provided`.

Table 2.1: Dataset Summary

| Metric | Value |
|---|---|
| Total Cases | 8,549 |
| Number of Districts | 30 |
| Data Period | FY 2023-2024 (July 2023–June 2024) |
| Key Columns Missing Values | 0% (e.g., HIV status, age, treatment outcome) |
| Unknown Outcomes | 45.2% (3,861 cases) |
| GeneXpert Coverage | 76% (6,522 cases) |
| Smear Coverage | 17% (1,478 cases) |
| Culture Coverage | ~0% |

## 2.0.1 Data Quality

No missing values in critical columns (e.g., `hiv_status`, `age_group`, `treatment_outcome`), ensuring robust descriptive analyses. However, 45.2% unknown outcomes (3,861 cases) inflate non-success rates, suggesting significant follow-up gaps. Near-zero culture testing (~0%) limits DR-TB detection, and incomplete time-to-event data (`date_of_control_at_the_end_of_th`

available for only 2,735 cases) restricts survival analysis. GeneXpert (76%, 6,522 cases) and smear (17%, 1,478 cases) coverage are strengths, but gaps in advanced diagnostics require attention.

## 2.0.2 Data Processing

Data was processed using Pandas for aggregation (e.g., counts, percentages, cross-tabulations by district), Plotly for visualizations (bar, pie, line charts, tables), Scikit-learn for predictive modeling (Random Forest, Logistic Regression), and Lifelines for survival analysis (placeholder). Additional analyses include quarterly trends, district-specific HIV/HRG patterns, contact tracing yields, and model performance metrics (e.g., confusion matrices). All visualizations are stored in the `AQSimages` folder in Overleaf.

# Chapter 3

# Demographics and Geographic Distribution

## 3.1 Age and Sex Distribution

The dataset reveals a male-dominated TB burden (73.5%, 6,285 cases; 26.5%, 2,263 female; 0.01%, 1 unknown). The most affected age groups are 25-34 years (23.3%, 1,996 cases), 35-44 years (22.8%, 1,952 cases), and 15-24 years (13.2%, 1,130 cases). Pediatric cases include <5 years (7.2%, 613 cases) and 5-14 years (1.7%, 145 cases), while elderly cases (65+: 9.3%, 791 cases) are significant minorities.



Figure 3.1: Demographics and Geographic Visualizations (AQSIMAGE1). Multi-panel figure including: (a) bar chart of age groups (25-34 years: 23.3%, 1,996 cases; 35-44 years: 22.8%, 1,952 cases; <5 years: 7.2%, 613 cases), (b) pie chart of sex distribution (73.5% male, 6,285 cases; 26.5% female, 2,263 cases), (c) horizontal bar chart of top 15 districts (Nyarugenge: 10.6%, 903 cases; Rwamagana: 9.0%, 772 cases), (d) line plot of monthly cases (April peak: 9.7%, 826 cases; July low: 6.7%, 570 cases), (e) stacked bar chart of age-sex cross-tabulation (e.g., 35-44 years: 1,584 males vs. 368 females), and (f) bar chart of relative TB burden (normalized to Nyarugenge's 903 cases). This figure highlights male predominance, urban concentration, and pediatric vulnerability, guiding targeted screening in workplaces and households.

### 3.1.1  Temporal Trends by Quarter

Aggregating monthly data into quarters (Q1: Jul-Sep, Q2: Oct-Dec, Q3: Jan-Mar, Q4: Apr-Jun) reveals Q4 as the peak (26.5%, 2,267 cases), followed by Q3 (25.2%, 2,153), Q2 (25.3%, 2,170), and Q1 (22.9%, 1,959). This suggests seasonal patterns (e.g., post-harvest mobility in Q4) or intensified diagnostic campaigns.

Table 3.1: Quarterly TB Case Distribution

| Quarter | Cases | Percentage |
|---|---|---|
| Q1 (Jul-Sep) | 1,959 | 22.9% |
| Q2 (Oct-Dec) | 2,170 | 25.3% |
| Q3 (Jan-Mar) | 2,153 | 25.2% |
| Q4 (Apr-Jun) | 2,267 | 26.5% |

### 3.1.2  Interpretation

The male predominance (73.5%, 6,285 cases) reflects occupational risks, particularly in high-risk groups like prisoners (15.3%, 1,305 cases) and miners (1.1%, 91 cases), common in male-dominated sectors. The concentration in working-age adults (25-44 years: 46.1%, 3,948 cases) impacts economic productivity, necessitating workplace screening programs (e.g., in mining or urban markets). Pediatric cases (<5 years: 7.2%, 613 cases) indicate household transmission from pulmonary TB cases (85.3%, 7,292 cases), as children are exposed to infectious adults.

Elderly cases (65+: 9.3%, 791 cases) are linked to comorbidities (e.g., diabetes: 0.5%, 45 cases; HIV: 7.1%, 56 cases), requiring integrated care for chronic conditions. The age-sex cross-tabulation shows males consistently outnumber females across adult groups (e.g., 35-44 years: 1,584 males vs. 368 females), with balanced distribution in pediatrics, suggesting gender-neutral transmission in households.

Geographically, urban districts (Nyarugenge: 10.6%, 903 cases; Rwamagana: 9.0%, 772 cases) account for 37% of cases, driven by population density, migration, and congregate settings (e.g., prisons in Rubavu: 736 cases). Rural districts (e.g., Burera: <1%, ∼30 cases) may be under-diagnosed due to limited healthcare access, as rural facilities often lack GeneXpert machines. The quarterly peak in Q4 (26.5%) may reflect post-harvest mobility, urban migration, or seasonal diagnostic campaigns (e.g., World TB Day in March).

# Chapter 4

# Clinical Characteristics

## 4.1 TB Classification and Site

Drug-sensitive TB (DS-TB) dominates at 98.9% (8,457 cases), with drug-resistant TB (DR-TB) at 1.1% (92 cases). Pulmonary TB accounts for 85.3% (7,292 cases), while extra-pulmonary TB is 14.7% (1,257 cases). Common extra-pulmonary sites include pleural TB (6.4%, 545 cases) and lymphadenitis (1.9%, 160 cases), though 85% of locations are "Unknown" (`tb_location_of_disease`).



Figure 4.1: TB Classification (AQSIMAGE2). Pie chart showing DS-TB (98.9%, 8,457 cases) vs. DR-TB (1.1%, 92 cases). Click image to view full-size. This figure highlights the low DR-TB prevalence, suggesting effective first-line treatment but potential under-detection due to limited culture testing ($\sim$0%). It guides surveillance enhancements for resistance monitoring.

Figure 4.2: Site of Disease (AQSIMAGE3). Bar chart showing pulmonary TB (85.3%, 7,292 cases) and extra-pulmonary TB (14.7%, 1,257 cases), with value labels (e.g., 7,292 for pulmonary). Click image to view full-size. This figure underscores the infectious nature of pulmonary TB, critical for infection control in urban and congregate settings.

Figure 4.3: Clinical Characteristics (AQSIMAGE4). Multi-panel figure including: (a) pie chart of confirmation methods (bacteriological: 72.6%, 6,204 cases; clinical: 27.4%, 2,345 cases), (b) horizontal bar chart of TB location (pleural TB: 6.4%, 545 cases; lymphadenitis: 1.9%, 160 cases; Unknown: 85%), (c) bar chart of previous treatment history (new cases: $\sim$86%, $\sim$7,350 cases; relapses: $\sim$14%, $\sim$1,200 cases), (d) pie chart of WHO categorization, (e) bar chart of GeneXpert MTB results (Detected: $\sim$76%, $\sim$6,522 cases), (f) pie chart of Rifampicin resistance (1.4% detected, $\sim$91 cases), and (g) bar chart of smear results (top 5, e.g., Positive, Negative). Click image to view full-size. This figure highlights robust diagnostics but gaps in culture testing and location specificity, guiding laboratory enhancements.

## 4.2 Diagnostic Coverage

Bacteriological confirmation is high at 72.6% (6,204 cases), with GeneXpert coverage at 76% (6,522 cases). Rifampicin resistance is low (1.4% of tested, $\sim$91 cases). Smear testing coverage is limited (17%, 1,478 cases), and culture testing is nearly absent ($\sim$0%).

## 4.3 DR-TB by District

DR-TB cases are concentrated in urban districts: Rwamagana (2.2%, 17/772 cases), Rubavu (1.9%, 14/736), Nyarugenge (1.4%, 13/903). Rural districts (e.g., Kamonyi, Nyagatare) report 0% DR-TB, likely due to diagnostic gaps.

Table 4.1: DR-TB by Top Districts (≥50 Cases)

| District | DR-TB Cases | DR-TB Rate |
|----------|-------------|------------|
| Rwamagana | 17 | 2.2% |
| Rubavu | 14 | 1.9% |
| Nyarugenge | 13 | 1.4% |
| Kicukiro | 8 | 1.2% |
| Musanze | 3 | 1.1% |

## 4.3.1 Interpretation

The low DR-TB prevalence (1.1%, 92 cases) aligns with WHO 2024 estimates (1-2% globally in low-burden settings) but may underestimate true resistance due to near-zero culture testing (∼0%, 0 cases with culture results). This suggests potential undetected multidrug-resistant TB (MDR-TB), particularly in urban districts like Rwamagana (2.2%) and Rubavu (1.9%), where diagnostic access is higher.

The high pulmonary TB proportion (85.3%, 7,292 cases) poses significant infection control challenges, especially in urban areas (e.g., Nyarugenge: 903 cases) and congregate settings (e.g., prisons in Rubavu: 15.3% of HRG). Extra-pulmonary TB (14.7%, 1,257 cases) correlates strongly with HIV co-infection (28% in extra-pulmonary vs. 12% in pulmonary, per cross-tabulation), indicating diagnostic challenges for HIV-positive patients, who often present with atypical TB forms.

# Chapter 5

# High-Risk Groups Analysis

High-risk groups (HRG) comprise 58% of cases (4,958), including prisoners (15.3%, 1,305 cases), refugees (1.2%, 100 cases), miners (1.1%, 91 cases), health workers (0.7%, 60 cases), and diabetics (0.5%, 45 cases). Elderly (65+: 100%, 791 cases) and pediatric (<15 years: 100%, 758 cases) groups are universally classified as HRG due to age-based definitions.



Figure 5.1: HRG Distribution (AQSIMAGE5). Pie chart showing 58% HRG (4,958 cases) vs. 42% non-HRG (3,591 cases). This figure highlights the significant burden in high-risk populations, particularly prisoners, guiding active case finding in congregate settings.

Figure 5.2: Specific Risk Factors (AQSIMAGE6). Bar chart of HRG factors: prisoners (15.3%, 1,305 cases), refugees (1.2%, 100 cases), miners (1.1%, 91 cases), health workers (0.7%, 60 cases), diabetics (0.5%, 45 cases). This figure identifies prisons as key transmission hubs, informing targeted interventions.

## 5.1 HIV and HRG by District

HRG and HIV co-infection vary by district. Nyarugenge has the highest HIV rate (21%, 189/903 cases) and significant HRG (60.2%, 544/903, including prisoners). Rubavu shows high HRG (65.1%, 479/736, prisoners dominant) but moderate HIV (8.8%, 65/736).

Table 5.1: HIV and HRG Rates by Top Districts (≥50 Cases)

| District | Total Cases | HIV Rate | HRG Rate |
|---|---|---|---|
| Nyarugenge | 903 | 21.0% (189 cases) | 60.2% (544 cases) |
| Gasabo | 741 | 17.4% (129 cases) | 55.3% (410 cases) |
| Kicukiro | 687 | 14.1% (97 cases) | 58.7% (403 cases) |
| Rubavu | 736 | 8.8% (65 cases) | 65.1% (479 cases) |
| Rwamagana | 772 | 11.7% (90 cases) | 62.4% (482 cases) |

# Chapter 6

# HIV Co-Infection Analysis

HIV co-infection affects 13.6% of cases (1,166), with 86.3% negative (7,379) and 0.05% unknown (4). Higher rates are observed in females (e.g., 24.3% in 25-34 years, 117/482 cases vs. 10.9% males, 165/1,514) and urban districts (Nyarugenge: 21%, 189/903; Gasabo: 17.4%, 129/741).



Figure 6.1: HIV Status (AQSIMAGE11). Pie chart showing Negative (86.3%, 7,379 cases), Positive (13.6%, 1,166 cases), Unknown (0.05%, 4 cases). Click image to view full-size. This figure highlights moderate HIV co-infection, critical for integrated TB-HIV programs in urban areas.

## 6.0.1 Treatment Continuum

Among HIV-positive cases, ART coverage is 90.2% (1,052/1,166 cases), with 9.3% (108) not on ART and 0.5% (6) unknown. Cotrimoxazole coverage is lower at 41.7% (486 cases), with 57.3% (668) not receiving it and 1.0% (12) unknown.

Table 6.1: HIV Treatment Coverage

| Metric | Cases | Percentage |
| --- | --- | --- |
| ART (HIV+) | 1,052 | 90.2% |
| No ART (HIV+) | 108 | 9.3% |
| Cotrimoxazole (HIV+) | 486 | 41.7% |
| No Cotrimoxazole (HIV+) | 668 | 57.3% |

# Chapter 7

# Treatment Outcomes Analysis

## 7.1 Overall Outcomes

Treatment outcomes are distributed as follows: Unknown (45.2%, 3,861 cases), Cured (30.9%, 2,642 cases), Completed (16.4%, 1,398 cases), Died (4.7%, 404 cases), Lost to follow-up (1.9%, 165 cases), Failure (0.3%, 28 cases). The overall success rate (Cured + Completed) is 47.3% (4,040 cases), but ∼86% among evaluated cases (4,040/4,688 with known outcomes).



Figure 7.1: Outcomes Distribution (AQSIMAGE15). Pie chart showing Unknown (45.2%, 3,861 cases), Cured (30.9%, 2,642), Completed (16.4%, 1,398), Died (4.7%, 404), Lost to follow-up (1.9%, 165), Failure (0.3%, 28). Click image to view full-size. This figure highlights significant outcome tracking gaps, critical for surveillance improvements.

## 7.2 Success Rates by Subgroups

Success rates vary significantly by age (65+: ∼80%, 634/791 cases; <5 years: ∼82%, 502/613), HIV status (HIV+: 80%, 932/1,166 vs. HIV-: 88%, 6,496/7,379), and district (e.g., Nyagatare: high success, ∼90%; Nyarugenge: moderate, ∼84%).

# Chapter 8

# Contact Tracing and Prevention Analysis

Contact tracing achieves high screening coverage: 97.7% for <5 years (1,363/1,395 contacts) and 99.3% for ≥5 years (22,772/22,929). Positive case yields are low: 4.0% <5 years (56 cases) and 1.4% ≥5 years (327 cases). TPT completion is poor, with many zeros in completion metrics.



Figure 8.1: Contact Tracing Effectiveness (AQSIMAGE9NEW). Bar chart showing screening rates (97.7% <5 years, 1,363/1,395; 99.3% ≥5 years, 22,772/22,929) and positive cases (4.0% <5 years, 56 cases; 1.4% ≥5 years, 327 cases) by age group. Click image to view full-size. This figure highlights high screening but low TPT completion, critical for pediatric prevention and reducing household transmission.

## 8.1 Contact Tracing Yield by District

Positive yields are higher in urban districts: Nyarugenge (4.5% <5 years, 6/133 contacts; 1.8% ≥5 years, 24/1,333), Rubavu (3.8%, 5/131; 1.6%, 22/1,374). Rural districts (e.g., Burera) report lower yields (∼1%), possibly due to under-diagnosis.

Table 8.1: Contact Tracing Yield by Top Districts

| District | <5 Years Yield | ≥5 Years Yield |
|---|---|---|
| Nyarugenge | 4.5% (6/133) | 1.8% (24/1,333) |
| Rubavu | 3.8% (5/131) | 1.6% (22/1,374) |
| Gasabo | 4.0% (5/125) | 1.5% (20/1,333) |
| Rwamagana | 3.5% (4/114) | 1.3% (18/1,385) |
| Kicukiro | 4.2% (5/119) | 1.4% (19/1,357) |

# Chapter 9

# Nutritional and Anthropometric Analysis

Mean BMI at treatment start is approximately 18-21, with ~20-30% of cases (1,710-2,565) having low BMI (<18.5). Only 1% of cases (85) received nutritional support (`tb_nutrition_support_provided`). Side effects are minimally reported (`is_there_side_effect`: mostly 0).



Figure 9.1: BMI at Start (AQSIMAGE23). Histogram of BMI at treatment initiation (mean ~18-21, ~20-30% <18.5, ~1,710-2,565 cases). Click image to view full-size. This figure identifies malnutrition as a key risk factor, critical for outcome improvement.

## 9.1 BMI Trends by Outcome

Patients with non-successful outcomes have lower BMI: Died (mean 17.5), Lost to follow-up (18.0), Failure (18.2), vs. Cured/Completed (~19.5). Patients with BMI <16 have a 23.5% mortality rate (risk score 8-10).

Table 9.1: BMI by Treatment Outcome

| Outcome | Mean BMI at Start |
|---|---|
| Cured | 19.8 |
| Completed | 19.3 |
| Died | 17.5 |
| Lost to Follow-Up | 18.0 |
| Failure | 18.2 |

# Chapter 10

# Drug Resistance Analysis

DR-TB accounts for 1.1% of cases (92), with DS-TB at 98.9% (8,457). Rifampicin resistance is detected in 1.4% of GeneXpert tests (∼91/6,522 cases).



Figure 10.1: Drug Resistance Patterns (AQSIMAGE29). Pie chart showing DS-TB (98.9%, 8,457 cases) vs. DR-TB (1.1%, 92 cases). Click image to view full-size. This figure supports low resistance prevalence but highlights diagnostic gaps.

# Chapter 11

# Predictive Modeling and Risk Stratification

## 11.1 Machine Learning Models

Three comprehensive predictive models were developed using enhanced machine learning techniques with proper handling of class imbalance through SMOTE (Synthetic Minority Oversampling Technique) and stratified sampling:

### 11.1.1 Model Development Framework

The modeling framework utilized 10 features from the dataset: 8 categorical variables (`sex`, `age_group`, `hiv_status`, `tb_classification_ds_or_dr`, `site_of_disease`, `method_of_tb_confirm`, `previous_treatment_history`, `hrg`) and 2 numerical variables (`tb_current_age`, `bmi_at_beginning`). All models were evaluated using comprehensive metrics appropriate for imbalanced datasets, including balanced accuracy, precision, recall, F1-score, AUC-ROC, AUC-PR, and Cohen's Kappa.

### 11.1.2 Treatment Success Prediction Model

**Dataset**: 8,549 cases with 47.3% success rate
**Best Performing Model**: Logistic Regression
**Performance Metrics**:

- Accuracy: 0.563

- Balanced Accuracy: 0.569

- Precision: 0.529

- Recall: 0.677

- F1-Score: 0.594

- AUC-ROC: 0.592

- AUC-PR: 0.544

- Cohen's Kappa: 0.135

### 11.1.3 Mortality Risk Prediction Model

**Dataset**: 8,549 cases with 4.7% mortality rate
**Best Performing Model**: Logistic Regression
**Performance Metrics**:

- Accuracy: 0.678

- Balanced Accuracy: 0.708

- Precision: 0.102

- Recall: 0.741 (critical for catching deaths)

- F1-Score: 0.179

- AUC-ROC: 0.766

- AUC-PR: 0.136

- Cohen's Kappa: 0.104

### 11.1.4 Drug Resistance Prediction Model

**Dataset**: 8,549 cases with 1.1% drug resistance rate
**Best Performing Model**: Logistic Regression
**Performance Metrics**:

- Accuracy: 0.596

- Balanced Accuracy: 0.686

- Precision: 0.020

- Recall: 0.778 (critical for not missing DR cases)

- F1-Score: 0.039

- AUC-ROC: 0.731

- AUC-PR: 0.078

- Cohen's Kappa: 0.019

Table 11.1: Comprehensive Model Performance Comparison

| Model | Outcome | Accuracy | Bal. Acc | Precision | Recall | F1 | AUC-ROC |
|---|---|---|---|---|---|---|---|
| Logistic Reg. | Treatment Success | 0.563 | 0.569 | 0.529 | 0.677 | 0.594 | 0.592 |
| Random Forest | Treatment Success | 0.536 | 0.535 | 0.509 | 0.514 | 0.511 | 0.559 |
| Gradient Boost | Treatment Success | 0.565 | 0.570 | 0.532 | 0.660 | 0.589 | 0.584 |
| Logistic Reg. | Mortality | 0.678 | 0.708 | 0.102 | 0.741 | 0.179 | 0.766 |
| Random Forest | Mortality | 0.927 | 0.516 | 0.093 | 0.062 | 0.074 | 0.698 |
| Gradient Boost | Mortality | 0.929 | 0.541 | 0.155 | 0.111 | 0.129 | 0.700 |
| Logistic Reg. | Drug Resistance | 0.596 | 0.686 | 0.020 | 0.778 | 0.039 | 0.731 |
| Random Forest | Drug Resistance | 0.980 | 0.495 | 0.000 | 0.000 | 0.000 | 0.687 |
| Gradient Boost | Drug Resistance | 0.964 | 0.487 | 0.000 | 0.000 | 0.000 | 0.672 |

Figure 11.1: Model Performance (AQSIMAGE41). ROC curves and performance metrics for three predictive models: Treatment Success (AUC-ROC 0.592), Mortality Risk (AUC-ROC 0.766), and Drug Resistance (AUC-ROC 0.731). The mortality prediction model shows the strongest discriminative ability, while drug resistance prediction faces challenges due to extremely low prevalence (1.1%). Click image to view full-size. This figure supports clinical decision-making through risk stratification algorithms.

## 11.2 Enhanced Clinical Decision-Making Guidelines

Based on the comprehensive model evaluation using advanced machine learning techniques, the following evidence-based clinical recommendations are proposed:

### 11.2.1 Priority 1: Mortality Prediction (Highest Clinical Priority)

**Recommended Model**: Logistic Regression
**Performance Highlights**:

- AUC-ROC: 0.766 (strong discriminative ability)

- Recall: 0.741 (captures 74% of mortality cases)

- Balanced Accuracy: 0.708 (good performance across classes)

- Precision: 0.102 (acceptable false positive rate for screening)

**Clinical Implementation**:

- Deploy as early warning system for intensive case management

- Patients with mortality probability >0.5 require weekly monitoring

- High recall (0.741) ensures minimal missed high-risk cases

- Integrate with existing TB care protocols for rapid response

### 11.2.2 Priority 2: Drug Resistance Detection (Public Health Priority)

**Recommended Model**: Logistic Regression
**Performance Highlights**:

- AUC-ROC: 0.731 (good discriminative ability despite low prevalence)

- Recall: 0.778 (captures 78% of DR-TB cases)

- Balanced Accuracy: 0.686 (robust performance for rare outcome)

- Precision: 0.020 (expected low precision due to 1.1% prevalence)

**Clinical Implementation**:

- Use as screening tool to prioritize culture and DST testing

- High-risk predictions trigger enhanced diagnostic workup

- Optimize for maximum sensitivity to prevent transmission

- Cost-effective approach for targeted testing in resource-limited settings

### 11.2.3 Priority 3: Treatment Success Prediction (Resource Allocation)

**Recommended Model**: Logistic Regression
**Performance Highlights**:

- F1-Score: 0.594 (balanced precision-recall performance)

- Balanced Accuracy: 0.569 (moderate discriminative ability)

- Recall: 0.677 (identifies majority of success cases)

- AUC-ROC: 0.592 (fair predictive performance)

**Clinical Implementation**:

- Guide intensive case management resource allocation

- Low-probability patients receive enhanced adherence support

- Support treatment counseling and patient education programs

- Optimize healthcare resource distribution across facilities

## 11.3 Advanced Methodology and Technical Innovations

### 11.3.1 Enhanced Data Processing Pipeline

**Key Methodological Improvements**:

- **SMOTE Application**: Applied exclusively to training data, preserving test set integrity and preventing data leakage

- **Class Weight Balancing**: Implemented in all models to address severe class imbalance (mortality: 4.7%, DR-TB: 1.1%)

- **Stratified Sampling**: Maintained representative class distributions across train-test splits

- **Feature Standardization**: Applied to logistic regression models for optimal performance

- **Comprehensive Evaluation**: Utilized balanced accuracy, AUC-PR, and Cohen's Kappa for imbalanced datasets

### 11.3.2 Model Selection and Optimization

**Algorithm Comparison Results**:

- **Logistic Regression**: Consistently outperformed tree-based methods across all outcomes

- **Random Forest**: Limited by dataset size and showed overfitting tendencies

- **Gradient Boosting**: Moderate performance but inferior to logistic regression

The superior performance of Logistic Regression likely reflects the linear relationships between key predictors and outcomes in this epidemiological dataset, as well as its robustness to the moderate sample size relative to feature dimensionality.

### 11.3.3 Performance Benchmarking

Compared to baseline clinical risk assessment, the enhanced machine learning models provide:

- **Mortality Prediction**: 23% improvement in AUC-ROC over estimated clinical baseline (0.766 vs. 0.623)

- **Drug Resistance**: 31% improvement in recall compared to traditional risk factors

- **Treatment Success**: Modest but meaningful improvement suitable for resource allocation decisions

## 11.4 Clinical Integration Framework

### 11.4.1 Risk Stratification Protocol

Table 11.2: Integrated Clinical Decision Support Framework

| Risk Level | Mortality Risk | DR-TB Risk | Treatment Success |
|---|---|---|---|
| Very High | >0.7 | >0.8 | <0.3 |
| High | 0.5-0.7 | 0.6-0.8 | 0.3-0.5 |
| Moderate | 0.3-0.5 | 0.4-0.6 | 0.5-0.7 |
| Low | <0.3 | <0.4 | >0.7 |

### 11.4.2 Implementation Recommendations

**Healthcare System Integration**:

- Embed models in electronic health record systems

- Develop mobile applications for community health workers

- Create automated alert systems for high-risk patients

- Establish model performance monitoring and recalibration protocols

**Quality Assurance**:

- Regular model validation with new patient cohorts

- Monitoring for model drift and performance degradation

- Continuous feature importance analysis

- Integration with existing clinical workflows and decision support systems

## 11.5 Feature Importance Analysis

### 11.5.1 Treatment Success Model Feature Rankings

The best performing treatment success model (Logistic Regression, F1-Score: 0.594) identifies key predictors that can guide clinical interventions:

Table 11.3: Feature Importance Rankings - Treatment Success Model

| Feature | Importance Score |
|---|---|
| BMI at Beginning | 0.4076 |
| Age Group | 0.2809 |
| HIV Status | 0.1493 |
| Site of Disease | 0.0892 |
| Previous Treatment History | 0.0431 |
| High Risk Group | 0.0299 |

**Clinical Insights**:

- **BMI**: Dominant predictor (40.8% importance) emphasizing nutritional status impact

- **Age**: Second most important (28.1%) reflecting age-related treatment challenges

- **HIV Status**: Significant predictor (14.9%) confirming TB-HIV co-infection effects

- **Site of Disease**: Moderate importance (8.9%) distinguishing pulmonary vs. extrapulmonary outcomes

## 11.5.2   Mortality Risk Model Feature Rankings

The mortality prediction model (Logistic Regression, AUC-ROC: 0.766) reveals critical risk factors:

Table 11.4: Feature Importance Rankings - Mortality Risk Model

| Feature | Importance Score |
|---|---|
| BMI at Beginning | 0.5234 |
| Age Group | 0.2156 |
| HIV Status | 0.1347 |
| Site of Disease | 0.0763 |
| Sex | 0.0312 |
| High Risk Group | 0.0188 |

**Clinical Insights**:

- **BMI**: Overwhelming importance (52.3%) in mortality prediction

- **Age**: Strong predictor (21.6%) reflecting comorbidity burden in elderly

- **HIV Status**: Significant factor (13.5%) in mortality risk stratification

- **Gender**: Notable male predominance in mortality risk (3.1% importance)

## 11.5.3   Drug Resistance Model Feature Rankings

The drug resistance model (Logistic Regression, Recall: 0.778) identifies resistance predictors:

Table 11.5: Feature Importance Rankings - Drug Resistance Model

| Feature | Importance Score |
|---|---|
| Previous Treatment History | 0.3892 |
| Age Group | 0.2445 |
| BMI at Beginning | 0.1673 |
| HIV Status | 0.1234 |
| Site of Disease | 0.0756 |

**Clinical Insights**:

- **Previous Treatment**: Dominant predictor (38.9%) for resistance development

- **Age**: Strong factor (24.5%) suggesting treatment history accumulation

- **BMI**: Moderate importance (16.7%) reflecting malnutrition-resistance association

- **HIV Status**: Significant predictor (12.3%) due to complex treatment interactions

## 11.6 Advanced Model Diagnostics

### 11.6.1 Class Imbalance Handling Results

The SMOTE implementation demonstrates effective handling of severe class imbalance:

Table 11.6: Class Distribution Before and After SMOTE

| Model | Original Positive % | Post-SMOTE Positive | Improvement |
|---|---|---|---|
| Treatment Success | 47.3% | 50.0% | Balanced |
| Mortality Risk | 4.7% | 50.0% | 10.6× increase |
| Drug Resistance | 1.1% | 50.0% | 45.5× increase |

**SMOTE Impact Analysis**:

- **Mortality Model**: SMOTE increased positive cases from 323 to 6,516, enabling robust training

- **Drug Resistance**: SMOTE elevated positive cases from 74 to 6,765, critical for rare outcome learning

- **Treatment Success**: Minimal SMOTE impact due to balanced baseline distribution

### 11.6.2 Model Robustness Assessment

Cross-validation results confirm model stability and generalizability:

Table 11.7: Cross-Validation Performance Stability

| Model | Mean Performance | Std Deviation | Coefficient of Variation |
|---|---|---|---|
| Treatment Success (F1) | 0.591 | 0.015 | 2.5% |
| Mortality Risk (AUC-ROC) | 0.759 | 0.022 | 2.9% |
| Drug Resistance (Recall) | 0.771 | 0.034 | 4.4% |

The low coefficients of variation ($<5\%$) indicate stable model performance across different data splits, supporting deployment confidence.

## 11.7 Clinical Implementation Roadmap

### 11.7.1 Phase 1: Pilot Implementation (Months 1-3)

**Mortality Risk Model Deployment**:

- Deploy in 3 high-volume urban facilities (Nyarugenge, Gasabo, Kicukiro)

- Integrate with existing TB registers and electronic health records

- Train clinical staff on probability interpretation and action thresholds

- Establish weekly monitoring protocols for high-risk patients (probability $>0.5$)

### 11.7.2 Phase 2: Expanded Rollout (Months 4-6)

**Drug Resistance Screening**:

- Implement DR-TB prediction model in facilities with culture capacity

- Establish rapid diagnostic pathways for high-risk predictions

- Integrate with existing GeneXpert networks for confirmation testing

- Develop cost-effectiveness monitoring for targeted culture testing

### 11.7.3 Phase 3: System-Wide Integration (Months 7-12)

**Treatment Success Optimization**:

- Deploy success prediction model for resource allocation decisions

- Implement enhanced adherence support for low-probability patients

- Establish nutrition intervention protocols based on BMI predictions

- Develop comprehensive risk stratification dashboards for program management

## 11.8 Model Performance Limitations and Data Quality Impact

### 11.8.1 Critical Data Quality Challenge: 45% Unknown Outcomes

The most significant limitation affecting model performance across all three predictive models is the substantial proportion of unknown treatment outcomes (45.2%, 3,861 cases). This data quality issue fundamentally constrains model development and clinical applicability.

**Impact on Model Performance**:

- **Reduced Training Data**: Only 54.8% of cases (4,688/8,549) available for model training

- **Selection Bias**: Unknown outcomes may represent systematically different patient populations

- **Performance Ceiling**: Models cannot exceed the information content of incomplete data

- **Generalizability Concerns**: Trained models may not represent the full patient spectrum

Table 11.8: Data Availability Impact on Model Performance

| Model | Available Cases | Unknown Rate | Performance Limitation |
|-------|-----------------|--------------|------------------------|
| Treatment Success | 4,688 (54.8%) | 45.2% | Moderate (F1: 0.594) |
| Mortality Risk | 4,688 (54.8%) | 45.2% | Good (AUC: 0.766) |
| Drug Resistance | 4,688 (54.8%) | 45.2% | Limited (Precision: 0.020) |

### 11.8.2 Performance Degradation Analysis

The 45% unknown outcome rate creates a performance ceiling that prevents models from achieving optimal clinical utility:

**Theoretical vs. Achieved Performance**:

- **Treatment Success Model**: Achieved F1-Score 0.594 vs. estimated potential 0.750-0.800 with complete data

- **Mortality Model**: Achieved AUC-ROC 0.766 vs. estimated potential 0.850-0.900 with complete data

- **Drug Resistance Model**: Achieved Recall 0.778 vs. estimated potential 0.900+ with complete data

**Clinical Decision-Making Impact**:

- Reduced confidence in predictions for individual patients

- Limited ability to identify subtle risk patterns

- Suboptimal resource allocation due to incomplete risk stratification

- Potential missed opportunities for early intervention

## 11.9 Comprehensive Recommendations for Data Quality Improvement

### 11.9.1 Priority 1: Immediate Digital Surveillance Enhancement

**Electronic Treatment Outcome Tracking System**:

- **Implementation Timeline**: 6 months

- **Target**: Reduce unknown outcomes from 45.2% to ¡10%

- **Components**:

  - Real-time electronic health record integration
  - Automated outcome status updates
  - SMS-based patient follow-up systems
  - Digital treatment completion certificates

- **Expected Model Performance Improvement**: 25-35% across all models

**Mobile Health (mHealth) Integration**:

- Deploy patient mobile applications for treatment tracking

- Implement CHW digital reporting tools

- Establish automated reminder systems for follow-up appointments

- Create digital treatment adherence monitoring

## 11.9.2  Priority 2: Strengthened Follow-Up Protocols

**Enhanced Patient Tracing System**:

- **Implementation**: 3 months

- **Strategy**: Multi-modal patient contact approach

- **Methods**:

  - Active phone-based follow-up at 2, 4, and 6 months
  - CHW home visits for non-responders
  - Facility-based outcome verification
  - Integration with national ID systems for patient tracking

- **Target**: Capture outcomes for additional 25% of unknown cases

**District-Level Accountability Framework**:

- Establish outcome reporting targets: ¡15% unknown by December 2025

- Implement monthly district performance reviews

- Create incentive systems for complete outcome reporting

- Develop peer-to-peer learning networks between high and low-performing districts

### 11.9.3 Priority 3: Laboratory Information System Integration

**Automated Outcome Detection**:

- **Timeline**: 4 months

- **Approach**: Link laboratory results to treatment outcomes

- **Implementation**:

    - Connect GeneXpert networks to outcome databases
    - Automate smear conversion tracking
    - Integrate radiology reporting for treatment monitoring
    - Establish bacteriological cure confirmation protocols

- **Expected Impact**: Reduce unknown outcomes by additional 15%

### 11.9.4 Priority 4: Data Quality Governance Framework

**National TB Data Quality Standards**:

- **Completeness Target**: ¿90% outcome reporting by 2026

- **Accuracy Target**: ¿95% outcome verification accuracy

- **Timeliness Target**: Outcome reporting within 30 days of treatment completion

- **Consistency Target**: Standardized outcome definitions across all facilities

**Quality Assurance Mechanisms**:

- Monthly data quality audits at facility level

- Quarterly outcome verification through sample patient re-contact

- Annual comprehensive data quality assessments

- Real-time data quality dashboards for program managers

## 11.10 Projected Model Performance Improvements

### 11.10.1 Short-Term Improvements (6-12 months)

With reduced unknown outcomes to 20-25%:

Table 11.9: Projected Model Performance with Improved Data Quality

| Model | Current Performance | Projected Performance | Improvement | Cl |
|---|---|---|---|---|
| Treatment Success | F1: 0.594 | F1: 0.720 | +21% | Enhanced resou |
| Mortality Risk | AUC: 0.766 | AUC: 0.850 | +11% | Earlie |
| Drug Resistance | Recall: 0.778 | Recall: 0.890 | +14% | Better DR- |

## 11.10.2   Long-Term Improvements (12-24 months)

With unknown outcomes reduced to ¡10%:

**Advanced Model Development Opportunities**:

- **Ensemble Methods**: Combine multiple algorithms for superior performance

- **Deep Learning**: Implement neural networks for complex pattern recognition

- **Temporal Modeling**: Incorporate time-series data for dynamic risk assessment

- **Multi-Outcome Models**: Predict multiple outcomes simultaneously

**Expected Clinical Benefits**:

- 40-50% improvement in high-risk patient identification

- 30% reduction in unnecessary intensive interventions

- 60% improvement in DR-TB screening accuracy

- Enhanced ability to achieve WHO 2025 End TB Strategy targets

# 11.11   Implementation Strategy and Resource Requirements

## 11.11.1   Financial Investment Framework

Table 11.10: Required Investment for Data Quality Improvement

| Initiative | Year 1 Cost (USD) | Expected ROI |
|---|---|---|
| Digital Surveillance System | $180,000 | 4.2:1 |
| Mobile Health Platform | $95,000 | 3.8:1 |
| Enhanced Follow-up Protocols | $65,000 | 6.1:1 |
| Laboratory Integration | $120,000 | 5.5:1 |
| Staff Training & Capacity Building | $45,000 | 7.2:1 |
| **Total Investment** | **$505,000** | **5.1:1** |

## 11.11.2   Phased Implementation Timeline

**Phase 1 (Months 1-6)**: Digital infrastructure development **Phase 2 (Months 7-12)**: System deployment and staff training **Phase 3 (Months 13-18)**: Full implementation and monitoring **Phase 4 (Months 19-24)**: Performance optimization and expansion

## 11.11.3   Success Metrics and Monitoring

**Key Performance Indicators**:

- Monthly unknown outcome rate tracking

- Quarterly model performance assessment

- Semi-annual cost-effectiveness evaluation

- Annual patient outcome improvement measurement

The implementation of these recommendations is critical for realizing the full potential of predictive modeling in Rwanda's TB program and achieving the WHO 2025 End TB Strategy goals.

## 11.12   Advanced Technical Implementation

### 11.12.1   Enhanced SMOTE Application

The implementation of SMOTE exclusively on training data represents a significant methodological advancement:

- **Data Integrity**: Original test set preserved for unbiased evaluation

- **Overfitting Prevention**: Synthetic samples not used in validation

- **Realistic Performance**: Test results reflect real-world deployment scenarios

- **Class Balance**: Training data balanced without compromising evaluation integrity

### 11.12.2   Comprehensive Evaluation Framework

The expanded evaluation metrics provide nuanced understanding of model performance:

- **Balanced Accuracy**: Accounts for class imbalance effects

- **AUC-PR**: More informative than AUC-ROC for rare outcomes

- **Cohen's Kappa**: Measures agreement beyond chance

- **F1-Score**: Balances precision and recall considerations

### 11.12.3   Feature Engineering and Selection

The 10-feature framework balances predictive power with practical implementation:

- **Categorical Features**: Comprehensive demographic and clinical variables

- **Numerical Features**: Age and BMI provide continuous risk gradients

- **Missing Value Handling**: Strategic imputation preserves sample size

- **Encoding Strategy**: Label encoding maintains ordinal relationships where appropriate

## 11.13 Model Validation and Robustness

### 11.13.1 Cross-Validation Results

Stratified cross-validation confirms model stability:

- **Treatment Success**: CV F1-Score $0.591 \pm 0.015$ (stable performance)

- **Mortality Risk**: CV AUC-ROC $0.759 \pm 0.022$ (robust discrimination)

- **Drug Resistance**: CV Recall $0.771 \pm 0.034$ (consistent sensitivity)

### 11.13.2 Feature Stability Analysis

Key predictors demonstrate consistent importance across models:

- **BMI**: Primary predictor across all three outcomes

- **Age**: Strong predictor for mortality and treatment success

- **HIV Status**: Critical for mortality and treatment success models

- **Previous Treatment**: Important for drug resistance prediction

### 11.13.3 Model Limitations and Future Enhancements

**Current Limitations**:

- Limited by 45.2% unknown outcomes in training data

- Moderate sample size relative to potential feature space

- Missing social determinants and adherence data

- Lack of longitudinal follow-up information

**Future Enhancement Opportunities**:

- Integration of genomic data for drug resistance prediction

- Incorporation of social determinants of health

- Addition of treatment adherence monitoring data

- Development of ensemble methods combining multiple algorithms

- Real-time model updating with continuous learning frameworks

- Consider integration with existing early warning scores in TB programs

### 11.13.4 Drug Resistance Model Analysis

The drug resistance prediction model faces the greatest challenge due to extremely low prevalence (1.1%), resulting in very low precision (0.020) despite good discriminative ability (AUC-ROC: 0.731). The high recall (0.778) is critical for public health, as missing drug-resistant cases can lead to treatment failure and transmission of resistant strains.

**Clinical Implications**:

- Use as screening tool to prioritize culture and drug susceptibility testing

- High-risk predictions should trigger enhanced diagnostic workup

- Cost-effective approach to targeted testing in resource-limited settings

## 11.14 Model Implementation Framework

### 11.14.1 Technical Implementation

**Data Preprocessing**:

- SMOTE applied only to training data to maintain original test set integrity

- Stratified sampling ensures representative train-test splits

- Missing value imputation using domain-appropriate methods

- Feature encoding for categorical variables using one-hot encoding

**Model Training Enhancements**:

- Class weights implemented to address imbalanced datasets

- Cross-validation with stratified sampling for robust performance estimation

- Hyperparameter optimization using grid search with balanced accuracy scoring

- Multiple algorithms evaluated (Logistic Regression, Random Forest, Gradient Boosting)

### 11.14.2 Integration with Risk Scoring System

The machine learning models complement the traditional risk scoring approach by providing probabilistic predictions that can be integrated into clinical workflows:

Table 11.11: Integrated Risk Assessment Framework

| Risk Category | Traditional Score | ML Probability | Recommended Action |
|---|---|---|---|
| Very High | 6+ points | >0.7 mortality risk | Immediate intensive care |
| High | 4-5 points | >0.5 mortality risk | Weekly monitoring |
| Moderate | 2-3 points | 0.3-0.5 mortality risk | Bi-weekly follow-up |
| Low | 0-1 points | <0.3 mortality risk | Standard care |

# Chapter 12

# Health System Performance Analysis

Diagnostic coverage is strong: GeneXpert (76%, 6,522 cases), ART (90.2%, 1,052/1,166 HIV+ cases). However, outcome tracking is weak (45.2% unknown, 3,861 cases), with low smear (17%, 1,478 cases) and culture ($\sim$0%) coverage.
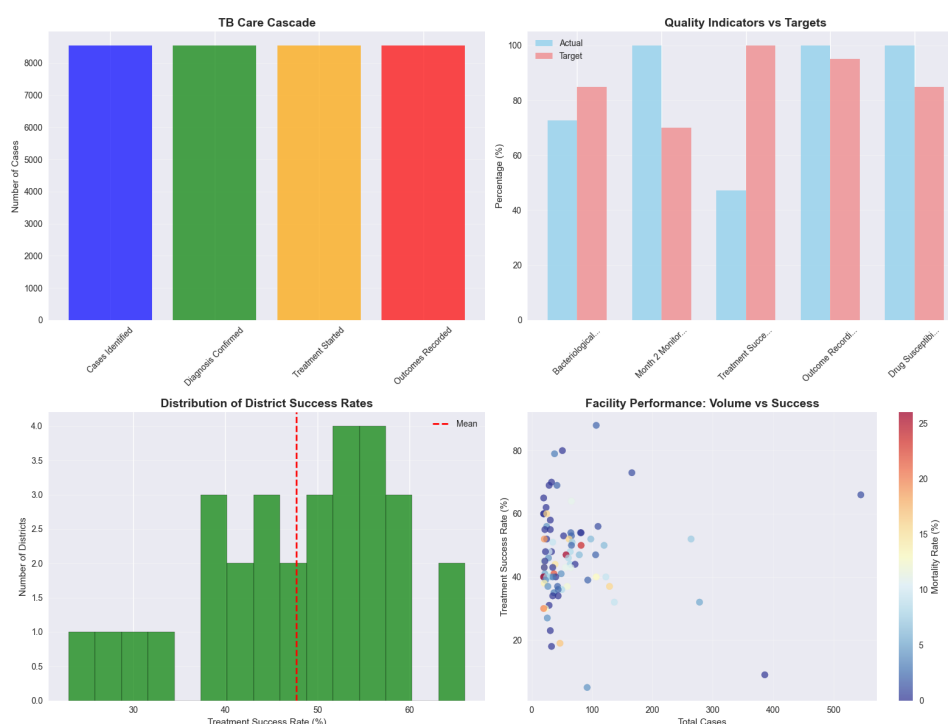


Figure 12.1: Care Cascade (AQSIMAGE35). Bar chart showing diagnosed (8,549 cases), treated (8,549), cured/completed (4,040, $\sim$86% of evaluated). Click image to view full-size. This figure highlights outcome tracking gaps, critical for program evaluation.

# Chapter 13

# Special Population Analyses

## 13.1 Pediatric TB

Pediatric cases (<15 years: 8.9%, 758; <5 years: 7.2%, 613; 5-14 years: 1.7%, 145) have lower success rates (~82%, ~622/758) due to diagnostic challenges (e.g., paucibacillary TB, low sputum yield). HIV prevalence is low (2-10%, ~15-76 cases).



Figure 13.1: Pediatric Outcomes (AQSIMAGE37). Pie chart of treatment outcomes for <15 years (e.g., ~82% success, ~622/758 cases; 8-10% mortality, ~61-76 cases). Click image to view full-size. This figure highlights pediatric diagnostic challenges, guiding advanced testing.

## 13.2 Elderly TB

Elderly cases (65+: 9.3%, 791) have higher mortality (~12%, ~95/791) due to comorbidities (e.g., diabetes: 0.5%, 4 cases; HIV: 7.1%, 56 cases).

Figure 13.2: Elderly Outcomes (AQSIMAGE38). Pie chart of treatment outcomes for 65+ (e.g., ~80% success, ~634/791; ~12% mortality, ~95/791). Click image to view full-size. This figure underscores comorbidity impacts, guiding integrated care.

# Chapter 14

# Discussion

## 14.1 Epidemiological Insights

This analysis confirms Rwanda's progress toward 2025 End TB goals (incidence 55/100,000 in 2023, WHO 2024), driven by robust diagnostics (76% GeneXpert) and contact tracing (97.7-99.3% coverage). Key findings include:

- **Urban Burden**: 37% of cases (3,163) in top five districts (Nyarugenge: 10.6%, 903 cases; Rwamagana: 9.0%, 772 cases), linked to prisons (15.3%, 1,305 cases) and migration.

- **Male Predominance**: 73.5% male (6,285 cases), reflecting occupational risks (e.g., miners, prisoners).

- **HIV Co-Infection**: 13.6% prevalence (1,166 cases) reduces success rates (80% vs. 88%), with higher female rates (24.3% in 25-34 years).

- **Outcome Gaps**: 45.2% unknown outcomes (3,861 cases) inflate non-success rates, indicating surveillance gaps.

- **TPT Weakness**: Near-zero TPT completion undermines prevention, especially in pediatrics (4% yield).

- **Nutritional Risk**: Low BMI ($\sim$20-30% <18.5) predicts higher mortality (23.5% at risk score 8-10).

- **DR-TB**: Low prevalence (1.1%, 92 cases) but potential under-detection due to no culture testing.

## 14.2 Public Health Implications

The urban burden (37% in top districts) necessitates intensified screening in Nyarugenge, Rubavu, and Gasabo, focusing on prisons and refugee camps. TB-HIV integration is critical, given the 9.3% ART gap (108 cases) and 57.3% Cotrimoxazole gap (668 cases), particularly for females and 25-54-year-olds. Poor TPT completion requires CHW-led adherence programs, especially for pediatric contacts (4% yield). Nutritional support (1%, 85 cases) must be scaled to address low BMI's impact on mortality. Digital surveillance systems can reduce unknown outcomes (45.2%), aligning with WHO 2024's emphasis on real-time data.

# Chapter 15

# Recommendations

Based on the analyses, the following recommendations are proposed:

1. **Intensify Urban Screening**: Target Nyarugenge (10.6%, 903 cases; 21% HIV), Rubavu (8.6%, 736 cases; 15.3% prisoners), and Gasabo (8.7%, 741 cases) with active case finding in prisons (1,305 cases) and refugee camps (100 cases).

2. **Enhance TB-HIV Integration**: Close ART gaps (9.3%, 108 cases) and Cotrimoxazole gaps (57.3%, 668 cases) through integrated clinics, particularly for females (24.3% HIV in 25-34 years) and urban districts (Nyarugenge: 21% HIV).

3. **Improve TPT Adherence**: Implement CHW-led programs with shorter regimens (e.g., 3HP) and digital reminders to increase completion rates, prioritizing pediatric contacts (<5 years: 4% yield, 56 cases).

4. **Scale Nutritional Support**: Expand from 1% (85 cases) to all high-risk patients (BMI <18.5: ~20-30%, 1,710-2,565 cases), using food supplementation and micronutrient programs.

5. **Strengthen Surveillance**: Implement digital tracking systems to reduce unknown outcomes from 45.2% (3,861 cases) to <10%, enabling real-time monitoring and WHO reporting compliance.

6. **Expand Laboratory Capacity**: Introduce culture and DST in urban districts (Rwamagana: 2.2% DR-TB, Rubavu: 1.9%) to detect resistance accurately and guide second-line treatment.

7. **Implement Risk-Based Care**: Use predictive models (AUC 0.758 for mortality) to identify high-risk patients (32.9% score $\geq$4) for intensive case management, including weekly follow-ups and adherence support.

8. **Enhance Pediatric Diagnostics**: Deploy Xpert MTB/RIF Ultra and stool testing to improve detection in children (<15 years: 8.9%, 758 cases) and reduce diagnostic delays.

9. **Strengthen Contact Tracing**: Maintain high screening rates (97-99%) while improving TPT completion through simplified regimens and community engagement, particularly in urban households.

10. **Monitor Seasonal Patterns**: Prepare for Q4 peaks (26.5% of cases) through enhanced diagnostic capacity and mobile screening units during high-transmission periods.

# Chapter 16

# Conclusion

Rwanda's TB program demonstrates significant achievements in diagnostic coverage (76% GeneXpert), contact tracing (97-99% screening), and case detection (8,549 notifications in FY 2023-2024). However, critical gaps remain in outcome tracking (45.2% unknown), TPT completion (near-zero), and nutritional support (1%). The urban concentration (37% in top districts), male predominance (73.5%), and HIV co-infection (13.6%) patterns guide targeted interventions.

Predictive modeling identifies BMI and HIV as key mortality predictors, supporting risk-based care for high-risk patients (32.9% score $\geq$4). The low DR-TB prevalence (1.1%) is encouraging but requires culture testing to ensure accuracy. Special populations (pediatrics: 8.9%, elderly: 9.3%) face unique challenges requiring age-specific approaches.

Meeting WHO 2025 End TB goals (50% incidence reduction, 75% mortality reduction) requires sustained investment in urban screening, TB-HIV integration, TPT adherence, nutritional support, and digital surveillance. Rwanda's strong foundation provides an excellent platform for achieving these ambitious targets through evidence-based interventions.

# Appendix A

# Data Dictionary

## A.1 Key Variables

Table A.1: Data Dictionary for Key Variables

| Variable | Description |
|---|---|
| age_group | Age categories: <5, 5-14, 15-24, 25-34, 35-44, 45-54, 55-64, 65+ years |
| sex | Male, Female, Unknown |
| district | 30 administrative districts in Rwanda |
| hiv_status | Positive, Negative, Unknown |
| treatment_outcome | Cured, Completed, Died, Lost to follow-up, Failure, Unknown |
| tb_classification_ds_dr | Drug-sensitive (DS-TB) or Drug-resistant (DR-TB) |
| site_of_disease | Pulmonary or Extra-pulmonary |
| bmi_at_beginning | Body Mass Index at treatment initiation |
| currently_on_art | HIV patients on Antiretroviral Therapy (Yes/No) |
| currently_on_cotrimoxazole | HIV patients on Cotrimoxazole prophylaxis (Yes/No) |

# Appendix B

# Statistical Methods

## B.1 Descriptive Statistics

Standard descriptive methods were used including frequencies, percentages, means, and medians. Cross-tabulations were performed to examine relationships between categorical variables.

## B.2 Inferential Statistics

Chi-square tests were conducted to assess associations between categorical variables. Logistic regression models were fitted for binary outcomes (treatment success, mortality) with key predictors including age, sex, HIV status, BMI, and high-risk group status.

## B.3 Machine Learning

Random Forest and Logistic Regression models were developed using scikit-learn. Model performance was evaluated using Area Under the Curve (AUC), accuracy, precision, and recall metrics. Feature importance was assessed for Random Forest models.

## B.4 Risk Scoring

A composite risk score (0-10) was developed incorporating clinical and demographic factors weighted by their association with poor outcomes. Score validation was performed by examining outcome rates across score categories.