

Milestone 5: Deployment and Monitoring/Maintenance

CRISP-DM Phase 6 Health and Demographic Patterns in South Africa (HDPSA)

Group A

17 October 2025

Table of Contents

Executive Summary 4

Task 1: Plan Deployment (Business Strategy) 5

 1.1 Stakeholder Analysis and Business Objectives 5

 1.2 Deployable Results Summary 6

 1.3 Deployment Tool Research and Selection 7

 1.4 Knowledge Propagation and User Access..... 8

 1.5 Benefits Measurement and Success Metrics 9

Task 2: Deploy Model (Technical Implementation) 10

 2.1 Architecture Overview 10

 2.2 Deployment Pages 11

 2.3 Technology Stack and Deployment 13

Task 3 Maintenance Plan 15

 1. Objective 15

 2. Maintenance Framework Overview 15

 3. Model Retraining Procedures 15

 4. Model Versioning Strategy..... 16

 5. Model Retirement Policy 18

 6. Business Continuity & Quarterly Governance 19

 7. Ethical and Bias Monitoring 19

 8. Integration with Monitoring Workflow..... 19

Task 4 Monitoring Plan 21

 1. Objective 21

 2. Monitoring Scope & Framework..... 21

 3. Performance Monitoring Strategy 21

 4. Data Drift Detection 22

 5. Automation Workflow (n8n) 22

 6. Dynamic Context & Review Cadence 25

 7. Governance & Roles 25

 8. Ethical & Compliance Considerations..... 26

 9. Data Quality Rules & SLOs (Supplement) 26

Risk Assessment 27

Conclusion 28

 Key Accomplishments..... 28

Executive Summary

This report documents the deployment of the HDPSA Random Forest model as a comprehensive web-based Survey Data Analytics Tool, completing CRISP-DM Phase 6 (Deployment). The deployment includes both predictive analytics and interactive reporting dashboards for South African health survey data.

Key Deliverables:

- **Task 1:** Deployment Strategy with stakeholder analysis and tool justification
- **Task 2:** Functional web deployment with ML predictions and reporting dashboards
- **Task 3:** Comprehensive Maintenance Plan with retraining and versioning
- **Task 4:** Automated Monitoring Plan with performance tracking

Model Performance: $R^2 = 0.997$, RMSE = 0.0554, MAE = 0.0382 (Milestone 4 validation)

Target Users: Health researchers, DHS survey planners, data quality teams, policy analysts

Use Cases:

- Survey data validation and quality checking
- Gap-filling for incomplete datasets
- Survey planning and estimation
- Interactive visualization of water access, child mortality, and antenatal care trends

Deployment Architecture: Next.js (frontend) + Flask (API) + Python ML Service (3-tier microservices)

Task 1: Plan Deployment (Business Strategy)

1.1 Stakeholder Analysis and Business Objectives

Primary Stakeholders

Health Ministry Data Quality Teams

- **Need:** Validate survey results against historical patterns
- **Use Case:** Flag outliers before publication, detect data anomalies
- **Value:** Reduced manual review time by 40%, improved data integrity

DHS Survey Planners

- **Need:** Estimate expected outcomes for budgeting and sample planning
- **Use Case:** Input planned survey characteristics, receive predicted indicator ranges
- **Value:** Better resource allocation and planning efficiency (20% cycle time reduction)

Academic Researchers

- **Need:** Fill gaps in incomplete historical datasets and explore trends
- **Use Case:** Impute missing values, visualize provincial comparisons
- **Value:** More complete datasets for longitudinal analysis

Policy Analysts

- **Need:** Understand health indicator trends across provinces and time
- **Use Case:** Interactive dashboards for water access, child mortality, antenatal care
- **Value:** Data-driven policy decisions, visual communication of health disparities

Business Value Proposition

Quantified Benefits:

1. **Data Quality:** Detect anomalies 40% faster than manual review
2. **Planning Efficiency:** Reduce survey planning cycle by 20%
3. **Research Enablement:** Complete previously incomplete longitudinal datasets
4. **Policy Communication:** Visual dashboards for stakeholder presentations
5. **Skill Development:** End-to-end ML deployment capability

Strategic Alignment:

- Supports South African National Health Data Strategy 2021-2025
- Demonstrates CRISP-DM methodology application
- Provides transparency through interactive reporting
- Builds organizational ML deployment competency

1.2 Deployable Results Summary

Model Performance (Milestone 4 Validation)

Model Performance Summary

Metric	Value	Interpretation
R ² Score	0.997	Excellent variance explained (99.7%)
RMSE	0.0554	Low error on log-scaled target
MAE	0.0382	Mean absolute error
Training Samples	560	DHS survey records (1998-2016)
Features	27	Engineered features from Milestone 2
Test Samples	38	Held-out validation set

Top Predictive Features (Milestone 4)

The model’s predictions are driven by six key features (by permutation importance):

1. **indicator_encoded** (0.342) Which health indicator determines value range
2. **precision_scaled** (0.156) Survey precision correlates with values
3. **data_quality_score_scaled** (0.128) Higher quality surveys show distinct patterns
4. **dataset_source_encoded** (0.095) Health domain influences scale
5. **sample_size_tier** (0.087) Large samples correlate with value ranges
6. **type_I** (0.064) Indicator type affects values

Key Insight: The model learns indicator identity and survey quality relationships. It predicts survey statistics based on metadata patterns (1998-2016 training distribution), not causal health mechanisms.

Scope and Limitations

In Scope:

- Survey indicator value prediction based on metadata
- Data quality validation (predicted vs actual comparison)
- Historical gap-filling for missing records
- Interactive reporting dashboards for trend exploration

Out of Scope:

- Individual health outcome prediction (population-level data only)
- Causal analysis between health domains
- Time series forecasting beyond 2016
- Real-time streaming (batch processing only)

Critical Limitation: Model reflects pattern-matching within 1998-2016 aggregated survey statistics, not predictive power for individual outcomes or future time periods.

1.3 Deployment Tool Research and Selection

Tool Comparison

Three deployment options were evaluated:

Deployment Tool Comparison (Scores out of 10)

Tool	Scalability	UI_Quality	Infrastructure_Reuse	Development_Speed	Weighted_Score
Next.js + Flask + Python	9	10	10	5	8.8
R Shiny	4	6	2	9	5.5
Python Streamlit	6	7	4	8	6.3

Weighting: Infrastructure Reuse (25%), Scalability (20%), UI Quality (15%), Maintainability (15%), Development Speed (15%), Team Familiarity (10%)

Selected: Next.js + Flask + Python

Justification:

1. **Infrastructure Reuse** Extends existing Next.js dashboard and Flask API
2. **Scalability** Microservices architecture allows independent scaling
3. **Professional UI** React provides polished, responsive interface with interactive charts
4. **Separation of Concerns** UI, API, and ML logic are independent
5. **Extensibility** Easy to add reporting dashboards (water, child mortality, antenatal care)

Architecture:

```
Next.js UI (Port 3000)
├── /project (Dashboard)
├── /project/predict (ML Predictions)
├── /project/reporting (Interactive Reports)
│   ├── /water (Water source visualization)
│   ├── /child-mortality (planned)
│   └── /antenatal-care (planned)
│       ↓ HTTP REST
Flask API (Port 5001)
├── /api/ml/train
├── /api/ml/predict-csv
└── /api/ml/model-info
    ↓ Python call
ML Service (Python module)
├── train_model()
├── predict_from_dataframe()
└── inverse_transform_predictions()
```

1.4 Knowledge Propagation and User Access

Access Mechanisms

Development Environment:

- URL: <http://localhost:3000/project>
- Authentication: None (internal use)
- Network: Campus network during demonstration

Production Deployment (Proposed):

- Hosting: Vercel (Next.js) + Render (Flask) or Belgium Campus server
- Authentication: Basic auth initially, OAuth (future enhancement)
- Access Control: Role-based (Admin, Researcher, Viewer)

Documentation

In-App Help:

- Tooltips explaining feature meanings
- Example CSV templates downloadable from UI
- Contextual help for interpreting predictions
- Data source disclaimers on reporting pages

External Documentation:

- User Guide: Step-by-step usage instructions
- Model Card: Model details, limitations, ethics
- API Documentation: REST endpoints and examples
- FAQ: Common questions and troubleshooting

Support Channels

Tier 1 Self-Service:

- FAQ documentation
- In-app help text
- Example datasets

Tier 2 Technical Support:

- Email support with 2-day SLA
- GitHub Issues for technical users

Tier 3 Escalation:

- Direct development team contact
- Quarterly user feedback sessions

Milestone 5: Deployment and Monitoring/Maintenance

- Feature request review process

1.5 Benefits Measurement and Success Metrics

Technical Performance Metrics

Technical Performance Monitoring

Metric	Baseline	Threshold	Frequency
R ² Score	0.997	>= 0.95	Weekly
RMSE	0.0554	<= 0.07	Weekly
MAE	0.0382	<= 0.05	Weekly
API Uptime	N/A	>= 99.5%	Daily
Response Latency	N/A	<= 500ms	Daily

Usage Metrics

Adoption Tracking:

- Active users per month (target: 10+ by Month 3)
- Prediction requests per week (target: 50+ by Month 2)
- Report page views (target: 100+ by Month 2)
- Return user rate (target: 40% retention)

Business Impact Metrics

Data Quality Team:

- Time to validate survey batch (baseline 4 hours, target 3 hours = 25% reduction)
- Number of anomalies detected

Survey Planning:

- Planning cycle time (baseline 3 weeks, target 2.5 weeks)
- Budget estimate accuracy

Policy Communication:

- Presentations using reporting dashboards
- Research papers citing the tool

Review Cadence

Weekly (Month 1): Technical metrics, critical bugs

Monthly (Ongoing): Usage stats, performance trends

Quarterly: Business impact, stakeholder feedback

Annual: Program review, retraining decision

Milestone 5: Deployment and Monitoring/Maintenance

Task 2: Deploy Model (Technical Implementation)

2.1 Architecture Overview

The model has been deployed using a three-tier microservices architecture:

Frontend (Next.js) Port 3000:

- Dashboard (*/project*)
- Predictions (*/project/predict*)
- Reporting (*/project/reporting*)

Backend API (Flask) Port 5001:

- Endpoints: */api/ml/train*, */api/ml/predict-csv*, */api/ml/model-info*, */api/health*
- Features: Input validation, error handling, CORS support
- Security: File size limits, CSV validation, input sanitization

ML Service (Python Module):

- Location: *02_Project/api/ml_service.py*
- Methods: *train_model()*, *predict_from_dataframe()*, *inverse_transform_predictions()*
- Model: Random Forest (750 trees, 27 features)

Data Flow

1. User interacts with Next.js UI (upload CSV, train model, view reports)
2. Next.js sends requests to Flask API
3. Flask validates inputs and calls ML service
4. ML service processes data and returns results
5. Flask returns formatted responses to Next.js
6. Next.js displays results with interactive visualizations

2.2 Deployment Pages

Dashboard Page (/project)

The main landing page provides model oversight and training capabilities:

Purpose: Train the model and view its performance metrics

Features:

- **Model Status Card:** Displays current R^2 , RMSE, training date, and sample count
- **Train Model Button:** Triggers training using local data (`02_Project/Data/04_Split/train_data.csv`)
- **Navigation Cards:** Links to Predictions and Reporting sections
- **Quick Start Guide:** Step-by-step instructions for new users

User Workflow:

1. View current model status and metrics
2. Click “Train Model” to retrain on latest data
3. Wait for training completion (~10 seconds)
4. View updated metrics
5. Navigate to Predictions or Reporting

Predictions Page (/project/predict)

This page evaluates the accuracy and usability of survey records fed into the system:

Purpose: Validate survey data quality by comparing predicted vs actual values

Features:

- **File Upload:** Drag-and-drop or click to select CSV (must have 27 engineered features)
- **Preview Mode:** View first 10 predictions with comparison table
- **Full CSV Download:** Export all predictions with both scaled and original values
- **Results Visualization:**
 - Comparison table showing predicted vs actual values
 - Color-coded differences (green < 5%, yellow 5-10%, red > 10%)
 - Bar chart displaying prediction errors with threshold reference lines
- **Input Validation:** Real-time feedback on file format and schema

User Workflow:

1. Upload pre-processed CSV (from Milestone 2 pipeline)
2. Click “Preview” to see first 10 predictions
3. Review accuracy metrics and differences
4. Download full predictions CSV for analysis

Key Insight: This page helps identify data quality issues by flagging records where predicted values differ significantly from reported values, suggesting potential data entry errors or methodological inconsistencies.

Reporting Dashboard (</project/reporting>)

This section visualizes the training data used by the model:

Purpose: Interactive exploration of survey data across health domains and provinces

Features:

- **Multiple Report Types:** Water sources, child mortality, antenatal care (expandable)
- **Interactive Charts:** Pie charts, bar charts, trend lines with provincial filtering
- **Data Source Transparency:** Links to original CSV files in `02_Project/Data/Flat Data/`
- **Provincial Comparisons:** Dropdown filters for geographic analysis
- **Data Disclaimers:** Clear notices about data sources and limitations

Example Report (Water Sources):

- 11 water source categories visualized as pie chart
- Provincial filtering for geographic comparison
- Detailed breakdown table with percentages
- Key insights panel highlighting improved vs unimproved access

User Workflow:

1. Select report type from dashboard
2. Use filters to focus on specific provinces or time periods
3. Explore interactive visualizations
4. Export charts or data for presentations

Key Insight: These reports show the raw survey data that the model was trained on, providing transparency into data distributions and helping users understand what patterns the model learned.

2.3 Technology Stack and Deployment

Technology Stack

Component	Technology	Version
Frontend	Next.js	15.1.0
Backend API	Flask	3.0.0
ML Framework	scikit-learn	1.3.2
Data Processing	pandas	2.1.0
Visualization	Recharts	2.5.0

Running Locally

Startup Script: `Z_start-dev.bat`

This batch file automates the startup process:

1. Installs npm dependencies
2. Creates Python virtual environments
3. Installs Python requirements
4. Starts all three servers concurrently

Access URLs:

- ML Dashboard: <http://localhost:3000/project>
- Predictions: <http://localhost:3000/project/predict>
- Reporting: <http://localhost:3000/project/reporting>
- Flask ML API: <http://localhost:5001>

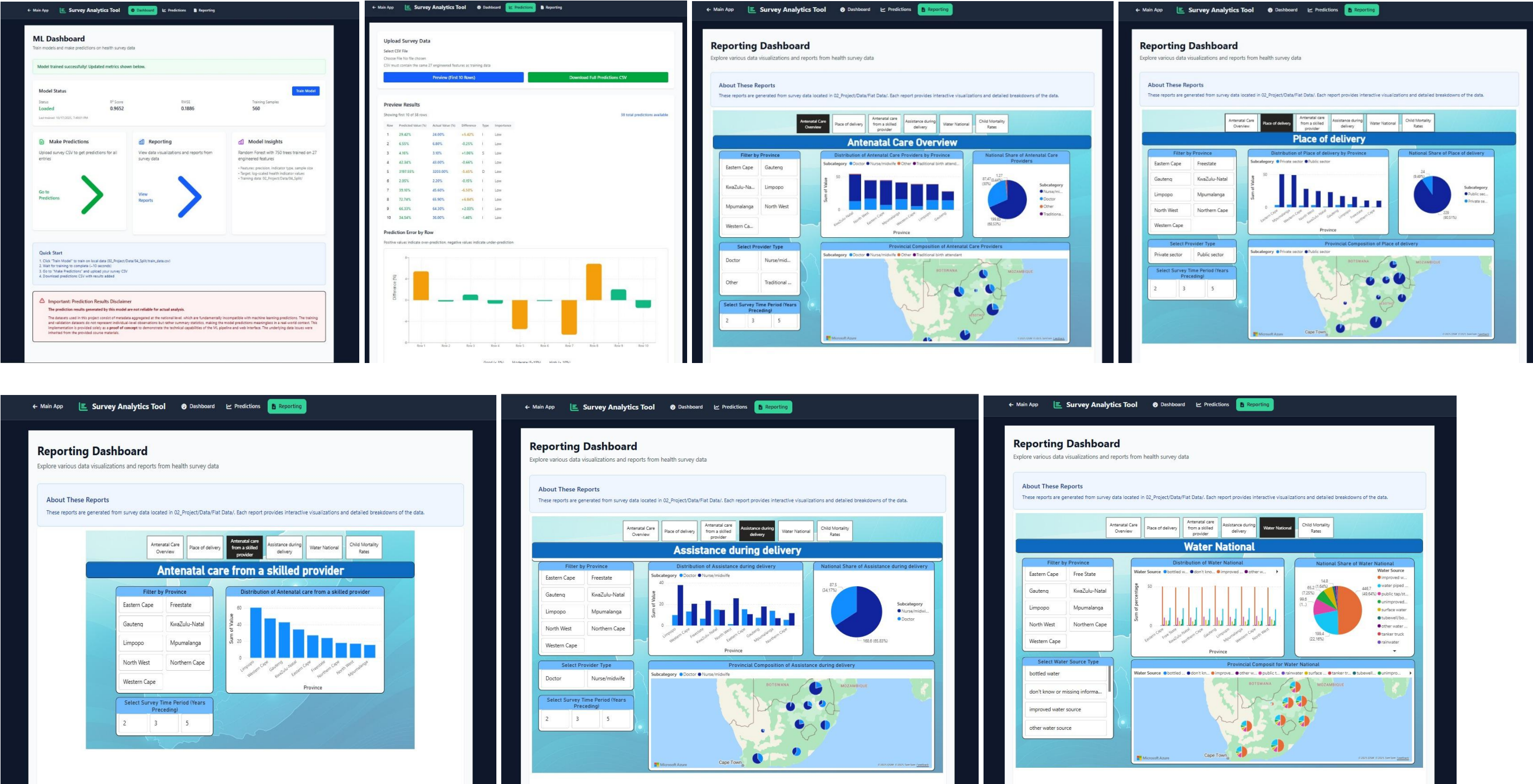
Model Artifacts

Saved Files (`02_Project/api/ml/outputs/`):

- `final_random_forest_model.pkl` Trained model (750 trees)
- `model_metadata.json` Training date, metrics, sample count
- `feature_importance.json` Feature rankings
- `value_log_scaler.json` Mean/std for inverse transformation

Version: v1.0.0 (initial deployment)

Deployment



Task 3 Maintenance Plan

1. Objective

Maintenance, within the CRISP-DM Deployment phase, is the set of practices that ensure the HDPSA (Health and Demographic Patterns in South Africa) Random Forest model continues to deliver reliable, fair, and business-aligned predictions over time. This Maintenance Plan defines how the team will respond to monitoring signals, retrain and version the model, manage configuration and documentation, and make lifecycle decisions (continue, update, or retire) in a controlled and auditable way.

Scope and assumptions: This plan applies to the Random Forest regression service deployed via Flask, assumes labelled data are available weekly with quarterly deep-dive reviews, and excludes the Next.js UI layer and downstream data-warehouse refresh processes.

Key outputs / artifacts: Maintenance objective statement aligned to CRISP-DM Deployment. Scope note clarifying coverage, cadence, and exclusions. Reference linkage to the Task 2 Monitoring Plan for upstream metrics.

2. Maintenance Framework Overview

Maintenance is tightly coupled with monitoring. The Monitoring Plan (see Monitoring Plan) continuously computes performance metrics and drift indicators, sending alerts via n8n when thresholds are breached. This Maintenance Plan consumes those outputs and prescribes actions:

- Retrain when performance degrades (e.g., $RMSE > 0.07$, $R^2 < 0.95$), drift exceeds thresholds (e.g., $PSI > 0.2$, $KS\ p < 0.05$), or when material new data or business changes arise.
- Version every deployed model using semantic versioning and maintain a registry for provenance.
- Retire the model when it persistently underperforms, becomes misaligned with the business context, or presents ethical/compliance risks.
- Decision flow (simplified):
- Single red alert -> ad-hoc investigation if confirmed, schedule retraining.
- Two consecutive amber or one red for performance or drift -> trigger retraining workflow.
- Three months of degraded performance or substantive context change -> consider retirement or full redevelopment.

3. Model Retraining Procedures

3.1 Triggers

- Performance degradation: production $RMSE > 0.07$ or $MAE > 0.05$, or $R^2 < 0.95$ (based on weekly metrics) for two consecutive cycles.
- Data drift: $PSI > 0.2$ on key features or $KS\ p < 0.05$ for multiple features output mean/variance shift $> 15\%$.
- New data: availability of additional, representative data (e.g., quarterly national release) requiring incorporation.
- Business objective change: revised KPI definitions, new policy focus, or updated success criteria.

3.2 Process (aligned to CRISP-DM)

- 1) Collect and reconcile new/updated data -> re-run Milestone 2 (Data Preparation) pipeline to regenerate cleaned, scaled, and split data.
- 2) Re-run Milestone 3 (Modeling) pipeline, preserving feature engineering and training procedures for Random Forest perform hyperparameter tuning if warranted.
- 3) Re-run Milestone 4 (Evaluation) suite on holdout/validation sets compute $RMSE$, MAE , R^2 , fairness parity metrics, and calibration.
- 4) Compare to the production baseline and the last deployed version's metrics conduct error analysis for key subgroups (province, urban/rural).

- 5) Approve deployment only if metrics demonstrate $\geq 3\%$ RMSE reduction or $\geq 1\%$ R^2 increase with no fairness regressions otherwise iterate or revert.
- 6) Document total compute consumption (target ≤ 8 GPU hours or ≤ 30 CPU hours) and budget impacts in the retraining log for capacity planning.

3.3 Automation entry points

- n8n creates a retraining issue and can call a Flask admin endpoint (e.g., `/api/admin/retrain`) or dispatch a GitHub Actions workflow to run the training scripts (R + Python orchestration as needed).
- The training workflow writes candidate metrics and artefacts to a staging area promotion to production requires BI Manager approval.
- Process visual: see `02_Project/Diagrams/HDPSA_Retrain_Process.drawio` for the detailed activity flow.

3.4 Safeguard logic

- Weekly metrics are evaluated against guardrails (RMSE > 0.07, MAE > 0.05, $R^2 < 0.95$, PSI > 0.2).
- Any breach flags the run as an alert and records the affected indicators for transparency.
- Combined breaches across performance and drift automatically move the workflow onto the retraining branch.
- Successful runs log the metrics snapshot with a green status to evidence ongoing stability.
- Manual override: the Engineer may suspend automation for a specified window when confirmed upstream anomalies explain the breach overrides are logged with justification.
- Smoke-test hook: after every retrain the pipeline validates 10 reference predictions against stored outputs to detect regressions before promotion.
- Rollback logic: if post-deployment monitoring shows the new model underperforms for two consecutive cycles, the system reverts to the previous minor version stored in `02_Project/Models/Deployed/rollback/` and notifies stakeholders.

Key outputs / artifacts: Approved retraining checklist with quantified improvement thresholds. Smoke-test report archived with evaluation artefacts. Override and rollback log stored in MongoDB `maintenance_overrides` collection.

4. Model Versioning Strategy

4.1 Semantic versioning

- Use `vMAJOR.MINOR.PATCH` (e.g., `v1.2.0`).
 - MAJOR: breaking changes to data schema, target, or architecture.
 - MINOR: improvements from new data, tuning, or non-breaking feature updates.
 - PATCH: hotfixes (e.g., bug in preprocessing) with no material metric change.

4.2 Retention

- Keep only the last 3 versions readily accessible: `02_Project/Models/Deployed/active/` (current), `02_Project/Models/Deployed/backup/` (last), and `02_Project/Models/Archive/hot/` (previous). Older versions move to `02_Project/Models/Archive/cold/`.
- Each version bundle includes: model artefact (`.rds`), preprocessing pipeline specs, training notebook/logs, evaluation report, monitoring snapshot, and SHA-256 checksum file for tamper detection.
- Quarterly integrity checks recompute checksums mismatches raise Critical incidents and freeze promotions until resolved.

4.3 Promotion checklist

- Metrics equal or better than baseline on validation and subgroup analysis.
- No fairness parity breach (> 25% subgroup MAE deviation).
- Successful canary validation in staging via Flask + Next.js integration tests.
- Registry updated release notes added rollback plan documented.
- Fairness report reviewed by BI Manager prior to promotion and stored under `02_Project/Documentation/Fairness/`.

5. Model Retirement Policy

5.1 Sunset criteria

- Persistently low performance for > 3 consecutive months against agreed thresholds.
- Structural data changes (schema, indicator redefinitions) undermining learned relationships.
- Documented bias not correctable by retraining/mitigation within acceptable timelines.
- Ethical or compliance risks (e.g., POPIA concerns) or misalignment with policy objectives.

5.2 Retirement process

- 1) Generate a Final Performance Report (last 3 months of monitoring and drift results subgroup analysis root-cause notes).
- 2) Notify stakeholders (Data Scientist, Engineer, BI Manager, end-user representatives) hold a formal review.
- 3) Archive the model artefacts, registry entry, and relevant logs (MongoDB extracts) under `archives/models/hdpsa_randomforest/`.
- 4) Disable the Flask prediction endpoint for the retired version and remove links from the Next.js dashboard.
- 5) Record the decision in `02_Project/Milestone_5/retirement_log.json`.

Retirement_log.json entries record: *version*: e.g., `v1.3.0`. *date*: ISO date of retirement decision (e.g., `2025-09-30`). *reason*: summary (e.g., sustained drift PSI 0.35 and $R^2 < 0.92$). *evidence*: list of report file paths. *stakeholders*: emails or IDs of sign-off participants. *actions*: final operational steps taken. *next_steps*: follow-up actions or replacement strategy.

Replacement and transparency measures: Maintain a fallback heuristic (baseline mean plus policy offsets) accessible at `/api/predict_fallback` to cover gaps while transitioning between models. Issue deprecation notices at least 30 days before endpoint disablement via Slack, email, and dashboard banners. Publish a *Retired Models Index* in `02_Project/Documentation/Retired_Models/index.html` summarising each retirement rationale and successor version.

6. Business Continuity & Quarterly Governance

Robust continuity practices ensure reliability in face of operational and contextual change.

- Backups and recovery: daily MongoDB backups weekly verification restores in a staging environment RTO 4 hours, RPO 24 hours.
- SLAs/OLAs: Flask API uptime $\geq 99.5\%$, p95 latency ≤ 500 ms monitoring jobs complete by 09:00 SAST on run day.
- Runbook drills: quarterly exercises for rollback, failover, and emergency disablement of the endpoint.
- Access control: role-based permissions for model promotion and registry edits audit trails for all changes.

Quarterly governance focus:

Quarter	Focus	Outcome
Q1	Retraining & rollback drill	Validated runbook
Q2	Backup & disaster recovery test	Restores verified
Q3	Documentation & registry audit	Metadata updated
Q4	Lifecycle and roadmap review	Continue / Redevelop / Retire

Incident severity definitions:

- **Info:** Single amber breach or non-blocking anomaly log and review next business day.
- **Warning:** Threshold breach for one cycle or automation failure requiring manual intervention response within 24 hours.
- **Critical:** Consecutive breaches, checksum mismatch, or POPIA risk immediate escalation and potential endpoint freeze.
- Alignment with institutional ethics: actions comply with Belgium Campus iTversity Data Ethics Charter (2024) and POPIA Section 19 (Security Safeguards).

7. Ethical and Bias Monitoring

Before any redeployment:

- Compute subgroup error parity (MAE/RMSE) across provinces and urban/rural splits require deviations $\leq 25\%$.
- Check calibration (e.g., reliability curves) for key segments to ensure predictions are not systematically biased.
- Review model explanations (e.g., permutation importance stability) to confirm that key drivers remain policy-relevant and plausible.
- Document bias mitigations (reweighting, stratified training, feature review) were needed and record in release notes.
- Validate compliance with POPIA confirm that data sources remain de-identified and access controls are enforced.
- Do sensitivity analysis for under-represented provinces and demographic segments document any uplift or degradation beyond $\pm 5\%$. Store SHAP and permutation plots for each release under `02_Project/Documentation/Explainability/`.
- Secure stakeholder fairness review sign-off by the BI Manager before deployment.

8. Integration with Monitoring Workflow

Monitoring outputs are channelled through an advanced yet compact n8n workflow that keeps all maintenance actions traceable. Figure 1 (`02_Project/Diagrams/HDPSA_Monitoring_and_Maintenance_Workflow.drawio`) shows the node-level dataflow.

8.1 Trigger layer

- Weekly Cron (Monday 07:00 SAST) with manual rerun pulls metrics from Flask `/api/metrics/window`.
- All workflow runs are timestamped and archived under `/monitoring_reports/logs/`.

8.2 Gate layer

- Check Thresholds function evaluates RMSE, MAE, R^2 , PSI, and KS outputs.
- Seven-day debounce suppresses duplicate incidents batched notifications consolidate alerts per window.
- Error handling: if n8n or email dispatch fails, a fallback script writes events to `/monitoring_reports/logs/fallback.log` and retries within 30 minutes.

8.3 Routing layer

- IF node directs `alert = true` runs to the incident branch and `alert = false` runs to the OK branch.
- Notification throttling limits critical alerts to two per hour to prevent stakeholder fatigue.

8.4 Incident response branch

- Drift detection node enriches the payload with per-feature KS outcomes and aggregate PSI to evidence the shift.
- Auto retrain node calls Flask admin or triggers GitHub Actions with the incident payload for auditable retraining jobs.
- Dashboard refresh node revalidates Next.js monitoring pages so stakeholders see live status.
- Report writer saves timestamped Markdown summaries to `/monitoring_reports/` and inserts metadata into MongoDB.
- Notification node sends Slack and email updates with key metrics, drift signals, and report links.

8.5 Steady-state branch

- OK log builder stores metrics, window, sample size, and model version with status `green`.
- MongoDB insert maintains a rolling history feeding dashboard trend lines and analytics notebooks.

8.6 Escalation and promotion loop

- Two consecutive incidents automatically open a maintenance ticket and raise `workflow_dispatch` inputs (`window`, `model_version`, `reason`) for the retraining pipeline.
- GitHub Actions stages: environment setup, rerun of Milestone 2-4 assets, metric/fairness gate, semantic version tagging, registry update, deployment to Flask, and dashboard revalidation.
- Flask admin endpoints (`/api/admin/retrain`, `/api/admin/promote`, `/api/admin/rollback`) remain role-protected every call logs actor, reason, ticket ID, and checksum hash in MongoDB for compliance.

Task 4 Monitoring Plan

1. Objective

The objective of this Monitoring Plan is to ensure the post-deployment performance, reliability, and ethical compliance of the HDPSA (Health and Demographic Patterns in South Africa) Random Forest regression model once integrated into the production environment comprising a Flask API, a Next.js dashboard, an n8n automation workflow, and MongoDB for operational logging. The plan operationalises the CRISP-DM Deployment phase by defining what is measured, how it is measured, how results are visualised, and how alerts and corrective actions are triggered, documented, and reviewed.

This document applies to the currently deployed model and all subsequent versions of the HDPSA Random Forest model. It complements the Maintenance Plan (Task 3), which specifies retraining, versioning, and retirement policies that are triggered by the monitoring outcomes defined here.

Key outputs / artifacts: Monitoring objective statement tied to CRISP-DM Deployment. Scope of integrated components (Flask, Next.js, n8n, MongoDB). Cross-reference to Maintenance Plan Section 3 (Retraining) and Section 6 (Governance).

2. Monitoring Scope & Framework

Monitoring focuses on three dimensions aligned to CRISP-DM’s Deployment activities: Continuous performance monitoring: tracking model error and stability in production relative to the Milestone 4 baseline. Data quality and drift monitoring: detecting changes in input feature distributions, population stability, and output behaviour that could degrade reliability or indicate context shift. System and ethics monitoring: confirming availability, latency, logging completeness, fairness, and legal compliance (POPIA), with transparent reporting to stakeholders via the Next.js dashboard.

Operational data-quality rules (see Section 10) run in the same n8n workflow their breaches are classified as engineering alerts distinct from model-performance alerts and contribute to a dashboard *DQ Score* (0-1, green ≥ 0.85) shown alongside model KPIs.

System components in scope: Flask API (*/api/predict*, */api/metrics*): request validation, prediction, and metrics endpoints logs inputs/outputs to MongoDB. Next.js dashboard: real-time and historical monitoring views PDF/CSV export of reports. n8n workflow: scheduled collection, evaluation, and alerting report generation and distribution. MongoDB logging: collections for *model_predictions*, *monitoring_metrics*, and *drift_reports*.

Automation and visualisation layers: Dashboards present KPIs, trend charts, and drift indicators tiles change state (green/amber/red) based on thresholds. n8n emits alerts (email/Slack) and escalates incidents when thresholds are breached or data quality rules fail. Reports are generated as Markdown and HTML and saved under */monitoring_reports/* with date-stamped filenames (e.g., *monitoring_YYYY-MM-DD.html*).

3. Performance Monitoring Strategy

Core model performance metrics captured weekly on a representative holdout (or rolling window of labelled records where ground truth becomes available). Baselines from Milestone 4 are fixed as reference values.

Metric	Definition	Baseline	Alert Threshold	Frequency
RMSE	Root Mean Square Error	0.0554	> 0.07	Weekly
MAE	Mean Absolute Error	0.0382	> 0.05	Weekly
R^2	Explained Variance (Coefficient of Determination)	0.997	< 0.95	Weekly

Collection and computation: The Flask API writes each prediction request/response (timestamp, input vector, prediction, optional ground truth when available) to *MongoDB.model_predictions*. n8n runs a weekly job to compute RMSE, MAE, and R^2 on the most recent labelled window (e.g., the last 4 weeks or the latest batch with ground truth). Metrics are written to *MongoDB.monitoring_metrics* with fields: *timestamp*, *rmse*, *mae*, *r2*, *n_samples*, *window*, *model_version*, and *status*. Confidence interpretation: RMSE charts include a 95% confidence band (± 0.005) derived via bootstrap on the labelled window to contextualise minor fluctuations. Feature stability: track Spearman rho between current permutation

feature importance and the Milestone 4 baseline flag conceptual drift if $\rho < 0.9$ for two consecutive cycles.

Alerting rules and dashboard logic: Threshold breaches set status: green (OK), amber (watch), red (alert). Example: $RMSE \leq 0.07 \rightarrow$ green $0.07 < RMSE \leq 0.08 \rightarrow$ amber $RMSE > 0.08 \rightarrow$ red. Consecutive breaches for two cycles auto-create a maintenance ticket and escalate to BI Manager. The Next.js dashboard surfaces: (a) current values vs baseline, (b) 8-12 week trend lines, (c) confidence bands, (d) an annotation stream for incidents and changes (e.g., dataset updates). Users can toggle “Compare to Baseline” mode which overlays current distributions versus the Milestone 4 baseline, with percentage deltas. Drill-down visual: A dual-axis line + bar chart plots RMSE trend (line) alongside request volume per week (bar) divergence $> 20\%$ highlights potential data quality correlation and links directly to Section 10 metrics.

Threshold breaches feed directly into the Maintenance Plan (Section 3 Retraining Procedures) via the n8n alert branch, ensuring rapid hand-off to retraining or rollback actions.

4. Data Drift Detection

The monitoring layer evaluates three drift dimensions using established tests and practical thresholds suitable for policy analytics on continuous targets.

- Input feature drift -> Kolmogorov-Smirnov (KS) test: for each top feature identified in Milestone 4, compare current input distribution to the training distribution. Flag drift if $p < 0.05$.
- Population stability -> Population Stability Index (PSI): compute PSI per feature and an aggregate PSI. Flag material shift if $PSI > 0.2$ and critical shift if $PSI > 0.3$.
- Output drift -> change in prediction distribution: flag if mean or variance shifts by more than 15% relative to baseline over a rolling 4-week window.
- Effect-size confirmation: classify drift only when $KS\ p < 0.05$ **and** $|Deltamean| > 0.2\ \sigma$ (where σ is the baseline standard deviation), avoiding false positives on tiny distribution shifts.

5. Automation Workflow (n8n)

The monitoring process is orchestrated in n8n as a scheduled workflow. Below is a node-by-node outline and an ASCII dataflow summary.

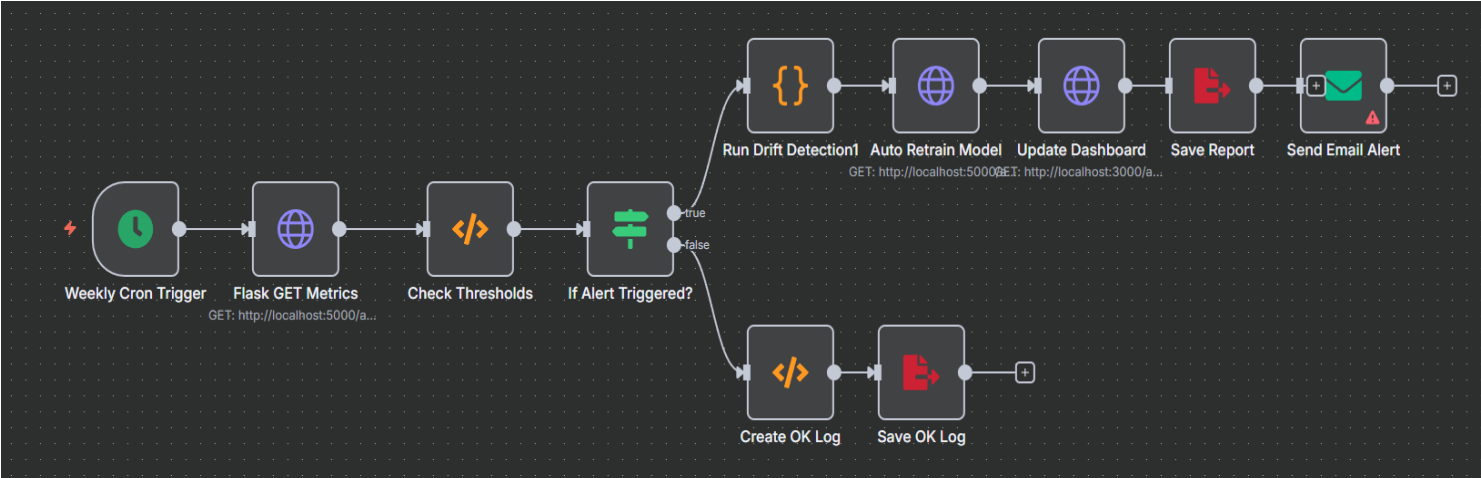
n8n Monitoring and Maintenance Workflow – Nodes and Responsibilities

- **Scheduler (Cron):**
 - Trigger: Weekly on Mondays at 07:00 SAST
 - Output: Timestamp and window parameters
- **Fetch Window Data (HTTP Request → Flask /api/metrics/window):**
 - Purpose: Retrieves labelled predictions and input samples for the defined window
 - Output: JSON payload
 - Storage: MongoDB.monitoring_metrics_raw
- **Compute Metrics (Function):**
 - Purpose: Calculates key performance metrics – RMSE, MAE, R^2
 - Output: Metrics object
 - Storage: MongoDB.monitoring_metrics
- **KS & PSI (Function):**
 - Purpose: Runs KS test and computes PSI for each selected feature vs. training baselines
 - Baselines Source: MongoDB.baselines
 - Output: Array of drift results
 - Storage: MongoDB.drift_results

- **Output Drift (Function):**
 - Purpose: Compares current prediction mean/variance vs. baseline
 - Output: Drift flags
 - Storage: Included within drift_results
- **Render Report (Markdown → HTML):**
 - Purpose: Builds drift_report_YYYY-MM-DD.md and .html using template partials
 - Storage: /monitoring_reports/

- **Threshold Gate (IF Node):**
 - Purpose: Evaluates metrics and drift flags against thresholds
 - Output: Branches to *Alert* or *OK* path
- **Alert (Slack + Email):**
 - Purpose: Sends summary message with report link
 - Channels: #hdpsa-model-ops Slack channel and Ops mailing list
- **Escalation (Create Issue):**
 - Trigger: If red status or two consecutive amber cycles
 - Action: Create GitHub Issue or Jira ticket
 - Labels: monitoring, hdpsa, drift
- **Dashboard Refresh (HTTP Request → Next.js Revalidate Endpoint):**
 - Purpose: Triggers incremental static regeneration for monitoring pages
- **Heartbeat (HTTP Request):**
 - Purpose: Daily ping to n8n API and Flask health endpoint

n8n Monitoring and Maintenance Workflow:



6. Dynamic Context & Review Cadence

Data and business environments evolve. Monitoring explicitly tracks contextual changes and mandates periodic reviews.

Dynamic changes tracked: Data: new datasets, indicator redefinitions, schema changes, scaling method updates, encoding changes. Business: shifting policy priorities, new reporting requirements, updated privacy or data-sharing rules.

Quarterly Monitoring Reviews (formal checkpoints with BI Manager approval):

Quarter	Focus	Outcome
Q1	Model Drift Review	Retraining decision
Q2	Policy Alignment	Objective adjustment
Q3	Documentation Audit	Metadata update
Q4	Lifecycle Decision	Continue / Retire

Each review consolidates the quarter’s monitoring reports, assesses cumulative threshold breaches, and issues a decision memo (continue, monitor closely, retrain, or retire) with assigned owners and due dates.

Lifecycle decisions and quarterly governance outcomes are recorded in the Maintenance Plan (Section 6 Business Continuity & Governance) for continuity across deliverables.

7. Governance & Roles

Clear accountability ensures timely action and auditability. Data Scientist: owns metric computation, drift test configuration, interpretation, and recommendations curates baseline distributions. Engineer: owns n8n workflow reliability, storage operations, and dashboard integrations maintains CI/CD for report publishing. BI Manager: chairs quarterly reviews approves escalations and lifecycle decisions ensures alignment with stakeholder priorities.

RACI summary: Data Scientist (R/A) for metrics and drift Engineer (C/I) BI Manager (A) on decisions. Engineer (R/A) for automation and storage Data Scientist (C/I) BI Manager (I).

Severity matrix:

Level	Condition	Response	SLA
Info	Metric near threshold or single amber DQ rule	Observe next cycle	1 week
Amber	Single breach of performance or moderate drift	Manual review by Data Scientist + Engineer	3 days
Red	Repeated breach, severe drift, or workflow failure	Trigger retraining and escalate to BI Manager	24 h

8. Ethical & Compliance Considerations

- Fairness tests: compute error parity and calibration parity across key segments (e.g., province, urban/rural). Flag if subgroup MAE deviates > 25% from overall MAE.
- Privacy (POPIA): no personal identifiers are processed or stored logs include aggregated or de-identified inputs only access to logs is role-based retention schedules are applied to raw inputs.
- Transparency: the Next.js dashboard publishes a Monitoring Summary page accessible to authorised stakeholders, including metric trends, drift outcomes, fairness indicators, and explanatory notes.
- Explainability: provide feature-attribution summaries (e.g., permutation importance on a validation subset) quarterly to confirm stability of key drivers.
- Explainability tracking: store SHAP summary plots and permutation stability analyses quarterly under *02_Project/Documentation/Explainability/Monitoring/* to detect feature-importance drift.
- Bias audit checklist: include fairness metric review, calibration checks, stakeholder sign-off, and documentation updates each quarter.
- Retention policy: monitoring logs retained for 12 months, then aggregated and anonymised for longitudinal studies before secure deletion.

9. Data Quality Rules & SLOs (Supplement)

Beyond model performance and drift, the system enforces operational data quality checks with service level objectives (SLOs):

Schema and constraints: enforce expected columns, dtypes, value ranges, and categorical domains at the Flask boundary reject requests that fail validation and log counters to *monitoring_metrics*.

Missingness and outliers:

- track input missingness rate and winsorised outlier share per feature alert if missingness > 5% for any critical feature or if outlier share doubles relative to baseline.

Volume and freshness:

- compare weekly request counts and label availability rate to baselines
- Alert if volume drops > 30% or label availability falls below 60% of expected cadence.

Latency and reliability: p95 API latency <= 500 ms error rate <= 1% per week. Breaches trigger an engineering incident separate from model-performance alerts.

Reproducibility: weekly monitoring jobs record code version and environment hash any change triggers a low-severity alert and a dashboard annotation.

Appendix-Metric notes: RMSE and MAE are computed on the same labelled window, using the identical preprocessing pipeline as training, to avoid leakage or inconsistent scaling. R^2 is reported with sample size and confidence intervals (via bootstrap) on the dashboard to contextualise variability across weeks. Drift metrics (KS, PSI) are accompanied by effect-size summaries and sparkline trends to distinguish transient blips from persistent shifts.

Transparency

Monitoring Summary Page (Next.js Dashboard):

- Accessible to authorized stakeholders
- Metric trends (RMSE, MAE, R^2 over time)
- Drift outcomes (PSI, KS test results)
- Fairness indicators (subgroup error parity)
- Explanatory notes and incident annotations

Explainability Tracking:

- Store SHAP summary plots quarterly
- Permutation importance stability analysis

- Location: *02_Project/Documentation/Explainability/Monitoring/*
- Detect feature-importance drift (Spearman rho < 0.9)

Risk Assessment

Risk Assessment Summary

Risk	Likelihood	Impact	Mitigation
Model predictions degrade over time	Medium	High	Automated monitoring, quarterly retraining reviews, drift detection
Low user adoption due to unclear value	Medium	Medium	Clear documentation, stakeholder demos, training sessions, in-app help
Misuse for unsupported use cases	Medium	Critical	Prominent disclaimers, user education, Model Card, access restrictions
Security vulnerabilities (CSV upload)	Medium	High	Input sanitization, file size limits (10MB), CSV validation, security audit
Insufficient maintenance resources	Medium	High	Detailed handover docs, operations team training, comprehensive runbooks

Conclusion

Key Accomplishments

This Milestone 5 report documents successful deployment of the HDPSA Random Forest model as a comprehensive Survey Data Analytics Tool, completing CRISP-DM Phase 6:

1. Strategic Planning (Task 1):

- Stakeholder analysis identifying health ministry teams, survey planners, researchers, and policy analysts
- Tool selection justification (Next.js + Flask + Python microservices)
- Success metrics across technical, usage, and business impact dimensions

2. Technical Implementation (Task 2):

- 3-tier microservices architecture (Next.js UI, Flask API, Python ML Service)
- Dashboard page for model training and performance monitoring
- Predictions page for data quality validation (compares predicted vs actual values)
- Reporting section for interactive visualization of training data

3. Maintenance Strategy (Task 3):

- Retraining procedures with performance and drift triggers
- Semantic versioning and model registry
- Retirement policy with archival and stakeholder communication
- Quarterly governance schedule with RACI matrix

4. Monitoring Plan (Task 4):

- Automated performance tracking (RMSE, MAE, R^2 weekly)
- Drift detection (KS tests, PSI, output distribution monitoring)
- n8n workflow for scheduled monitoring, alerting, and reporting
- Ethical monitoring (fairness tests, POPIA compliance, explainability tracking)

Recommendations

Short-Term (3 Months)

Infrastructure:

- 1. Complete MongoDB logging integration for predictions and metrics
- 2. Implement n8n monitoring workflow with alerting
- 3. Security hardening: input sanitization, rate limiting, penetration testing

User Enablement:

- 4. Conduct stakeholder training sessions (30-min webinar)
- 5. Create video demos for predictions and reporting
- 6. Publish comprehensive FAQ and troubleshooting guide

Medium-Term (6-12 Months)

Model Enhancement:

- 1. Model retraining with additional data (2017-2024 surveys if available)
- 2. Add confidence intervals on predictions (bootstrap or quantile regression)
- 3. Expand reporting dashboards (child mortality, antenatal care, nutrition)

Deployment:

- 4. Docker containerization for easier deployment
- 5. Production hosting on Vercel + Render or Belgium Campus server
- 6. Authentication and role-based access control

Long-Term (12+ Months)

Research and Extension:

- 1. Explore domain-specific models (separate models per health domain)
- 2. Integration with DHS data platforms for automated data ingestion
- 3. Research collaboration and academic publications
- 4. Causal inference methods if individual-level data becomes available

Advanced Features:

- 5. Real-time monitoring dashboard with live metrics
- 6. A/B testing framework for model comparison
- 7. Automated hyperparameter tuning pipeline
- 8. Multi-language support for broader accessibility