

Unsupervised Detection of Anomalous Wellness Profiles: Identifying Statistical Outliers in Self-Reported Mental Health and Lifestyle Data

Collaborative Final Project
CSST101 – Machine Learning
CSST102 – Knowledge Representation and Reasoning

Submitted by: Group Name: Group 2

Group Members:

- Dela Paz, Zyril
- Delos Reyes, Axcel Andrei
- Hornilla, John Benedict

PROJECT OVERVIEW

Combine an Isolation Forest anomaly detector with rule-based reasoning to identify individuals with unusual wellness indicators and produce risk-level recommendations and next-step actions. Objectives

OBJECTIVES

General Objective: Detect and prioritize potential mental-health risk cases using a hybrid ML + knowledge-based system.

Specific Objectives:

- Build an unsupervised anomaly detection model to score wellness-related records.
- Define knowledge-based rules to translate anomaly signals into human-understandable risk levels.
- Produce an interpretable report listing top anomalies and recommended actions.
- Package model and pipeline for simple deployment and repeatable evaluation.

SYSTEM ARCHITECTURE

User Input → Preprocessing → Isolation Forest (ML) → KRR Rules (CSST102) → Final Risk Level → Recommendations

Report Artifacts: Model saved as output/isolation_forest.joblib, anomaly outputs in output/anomaly_scores.csv and output/top_anomalies.csv Machine Learning Component (CSST101)

MACHINE LEARNING COMPONENT (CSST101)

- **Algorithm Used:** Isolation Forest (unsupervised anomaly detection)Isolation Forest (unsupervised anomaly detection)
- **Dataset:** See archive (1)/Mental Health Dataset.csv/Mental%20Health%20Dataset.csv) for raw records
- **Model Evaluation:** Use AUC/precision@k where labeled data available, plus inspection of top anomalies and cluster summaries in output/

MACHINE LEARNING PIPELINE

- **Data Collection:** Survey/CSV dataset containing wellness indicators and demographic fields (source: archive (1)/Mental Health Dataset.csv).
- **Data Preprocessing:** Handle missing values, normalize numeric features, encode categorical features, remove low-variance columns, engineer aggregated wellness scores.
- **Model Training:** Train IsolationForest on cleaned features; tune n_estimators, contamination, and max_samples via cross-validation or domain-guided selection.
- **Model Evaluation:** Validate using known anomalies (if any), inspect distribution of anomaly scores, check stability across folds, review output/cluster_summary.md.
- **Model Deployment:** Save model to output/isolation_forest.joblib; provide simple inference API or notebook cell in anomaly_analysis.ipynb.

DATASET DESCRIPTION

- **Dataset Type:** Tabular CSV (survey/clinical indicators)
- **Number of Records:** See archive (1)/Mental Health Dataset.csv/Mental%20Health%20Dataset.csv) for the exact count
- **Target Variable:** No explicit label (unsupervised); target concept = anomalous/unusual wellness profile

KNOWLEDGE REPRESENTATION & REASONING (CSST102)

- **Rule 1:** IF anomaly_score ≥ 0.85 THEN Risk = High AND Recommend = "Clinical follow-up within 48 hours"
- **Rule 2:** IF 0.6 ≤ anomaly_score < 0.85 AND severe_symptom_flag = True THEN Risk = High AND Recommend = "Contact support hotline"
- **Rule 3:** IF 0.4 ≤ anomaly_score < 0.6 THEN Risk = Medium AND Recommend = "Schedule clinician review"
- **Rule 4:** IF anomaly_score < 0.4 AND no_concerning_flags THEN Risk = Low AND Recommend = "Routine monitoring"
- **Rule 5:** IF demographic_vulnerability = True AND anomaly_score ≥ 0.6 THEN Elevate Risk by one level AND Recommend = "Prioritized outreach"

HYBRID DECISION LOGIC

Logic: Compute anomaly score from Isolation Forest (0–1), apply deterministic KRR thresholds above, and combine with domain flags (e.g., `severe_symptom_flag`, `demographic_vulnerability`) using simple weighted rules to produce final Risk Level (Low/Medium/High) and corresponding recommendation text.

Aggregation: If multiple rules trigger, pick the highest risk; include all matching rule explanations in the report.

SYSTEM FEATURES

- Wellness risk prediction: Detects anomalous profiles.
- Rule-based recommendations: Maps scores to actionable steps.
- Reports & CSV outputs: `output/anomaly_report.md`, `output/top_anomalies.csv`.
- Model persistence: `output/isolation_forest.joblib`.
- Notebook & API-ready: `anomaly_analysis.ipynb` for interactive use.

TESTING AND EVALUATION

Test Case 1 | Input Summary: Synthetic record with extreme scores | Expected Output: $\text{anomaly_score} \geq 0.85 \rightarrow \text{Risk} = \text{High}$

Test Case 2 | Input Summary: Mild deviations, no flags | Expected Output: $\text{anomaly_score} \sim 0.5 \rightarrow \text{Risk} = \text{Medium}$

Test Case 3 | Input Summary: Normal profile | Expected Output: $\text{anomaly_score} < 0.4 \rightarrow \text{Risk} = \text{Low}$ Validation: Compare top anomalies with domain expert review; run tests/test_pipeline.py to verify end-to-end pipeline.

CONCLUSION

Outcome: A lightweight hybrid system that flags unusual wellness records, provides interpretable risk levels, and outputs recommendations for follow-up.

GROUP CONTRIBUTION

Member Name | Contribution

Dela Paz, Zyril - Main coder, Proof Reader

Delos Reyes, Axcel - Debugger, Lead Documenter

Hornilla, John Benedict - Dataset provider, Lead Presenter

REFERENCES

- Paganelli, A. I., Mondéjar, A. G., da Silva, A. C., Silva-Calpa, G., Teixeira, M. F., Carvalho, F., Raposo, A., & Endler, M. (2022). Real-time data analysis in health monitoring systems: A comprehensive systematic literature review. *Journal of Biomedical Informatics*, 127, 104009. <https://www.sciencedirect.com/science/article/pii/S1532046422000259>
- Sharma, S. K., Alutaibi, A. I., Khan, A. R., Tejani, G. G., Ahmad, F., & Mousavirad, S. J. (2025). Early detection of mental health disorders using machine learning models using behavioral and voice data analysis. *Scientific Reports*, 15(1). <https://doi.org/10.1038/s41598-025-00386-8>
- Bijlani, N., Nilforooshan, R., & Kouchaki, S. (2022). An Unsupervised Data-driven Anomaly Detection Approach for Detection of Adverse Health Conditions in People Living with Dementia: Cohort Study (Preprint). *JMIR Aging*. <https://doi.org/10.2196/38211>