



Diabetes Readmission Prediction - EDA Documentation

Project: Healthcare AI Project: Diabetic Readmission Prediction

Phase: Exploratory Data Analysis (EDA) - Phase 1.2

Date: August 2025

Dataset: UCI Diabetes Dataset (101,766 records, 50+ features)

Author: Mohammad Javad Aghababaie Beni



Executive Summary

This document provides a comprehensive analysis of the Exploratory Data Analysis (EDA) phase completed for the Diabetes Readmission Prediction project. The EDA process has successfully transformed raw healthcare data into actionable insights, creating a solid foundation for feature engineering and machine learning modeling.



Key Achievements

- **Comprehensive Data Understanding:** Analyzed 101,766 patient records with 50+ features
- **Target Variable Creation:** Established binary classification for 30-day readmission prediction
- **Clinical Risk Stratification:** Created healthcare-specific risk categories
- **Advanced Feature Engineering:** Developed 8 new clinical features
- **Socioeconomic Analysis:** Integrated social determinants of health
- **Production-Ready Insights:** Clear roadmap for Week 2 modeling phase



Project Overview

Objective

Build an end-to-end ML system to predict 30-day readmission risk for diabetic patients using MLOps best practices and healthcare-specific insights.

Dataset Characteristics

- **Size:** 101,766 patient encounters
- **Features:** 50+ clinical and demographic variables
- **Time Period:** 1999-2008 (10 years of clinical care)
- **Source:** 130 US hospitals and integrated delivery networks
- **Target:** Binary classification (readmission within 30 days: Yes/No)



EDA Process & Findings

Phase 1: Data Overview & Quality Assessment

1.1 Dataset Dimensions & Structure

- **Total Records:** 101,766 patient encounters
- **Total Features:** 50 columns
- **Data Types:** 37 categorical (object), 13 numerical (int64)
- **Memory Usage:** Optimized for large-scale analysis



1.2 Missing Value Analysis

Critical Findings:

- **max_glu_serum:** 94.7% missing (96,420 records)
- **A1Cresult:** 83.3% missing (84,748 records)

1.3 Data Type Validation

Key Insights:

-  No mixed data types detected
-  High cardinality categorical variables identified:
 - diag_1: 717 unique values
 - diag_2: 749 unique values
 - diag_3: 790 unique values
 - medical_specialty: 73 unique values

Phase 2: Target Variable Analysis



2.1 Target Variable Creation

Binary Classification Established:

-  **0 (No Readmission): 90,409 patients (88.84%)**
- **1 (Readmission <30 days): 11,357 patients (11.16%)**

Phase 3: Clinical Risk Stratification

3.1 Risk Categories Created

Risk Category	Diagnosis Range	Patient Count	Percentage
 Low Risk	1-3 diagnoses	4,077	4.0%
 Medium Risk	4-6 diagnoses	27,091	26.6%
High Risk	7-10 diagnoses	70,500	69.3%
Critical Risk	11+ diagnoses	98	0.1%

3.2 Readmission Rates by Risk Category

Risk Category	Readmission Rate	Clinical Significance
Low Risk	<div><div></div>6.97%</div>	Minimal intervention needed
Medium Risk	<div><div></div>9.44%</div>	Standard monitoring
High Risk	12.06%	Enhanced monitoring required
Critical Risk	14.29%	Intensive intervention needed

Phase 4: Treatment Complexity Analysis

4.1 Complexity Score Creation

Composite Score Formula:

$$\text{Treatment Complexity} = (\text{Procedures} \times 0.3) + (\text{Medications} \times 0.4) + (\text{Diagnoses} \times 0.3)$$

4.2 Complexity Categories

Complexity Level	Patient Count	Readmission Rate	Clinical Action
Low	372 (0.4%)	<div><div></div>4.03%</div>	Standard care
Medium	4,690 (4.6%)	<div><div></div>7.55%</div>	Enhanced monitoring
High	14,705 (14.5%)	<div><div></div>9.17%</div>	Care coordination

Complexity Level	Patient Count	Readmission Rate	Clinical Action
Critical	80,671 (79.3%)	11.77%	Intensive management

Phase 5: Socioeconomic Analysis

5.1 Risk Score Creation

Scoring System:

- **Medicaid (MC):** +2 points (higher risk)
- **Medicare (MD):** +1 point
- **African American Race:** +1 point
- **Unknown Age/Weight:** +1 point each

Phase 6: Advanced Feature Engineering

6.1 New Clinical Features Created

8 Advanced Features Developed:

1. **Medication Adherence Score** - Mean: 1.94, Range: 0-5
2. **Hospital Utilization Score** - Mean: 0.38
3. **Lab Efficiency Score** - Mean: 9.77
4. **Age Group Categorization** - Young, Middle, Senior, Elderly
5. **Length of Stay Risk** - Low, Medium, High, Critical
6. **Diagnosis Complexity** - Mean: 0.08
7. **Insurance-Age Interaction** - Combined risk factor analysis
8. **Clinical Severity Index** - Mean: 7.02



Why These EDA Steps Were Essential

1. Healthcare Domain Expertise

- **Risk Stratification:** Standard practice in healthcare
- **Treatment Complexity:** Directly impacts outcomes
- **Socioeconomic Factors:** Social determinants of health
- **Age Group Analysis:** Age-related risk patterns

2. Machine Learning Optimization

- **Feature Engineering Foundation:** Created clinically relevant predictors
- **Data Quality:** Ensured modeling data integrity
- **Clinical Metrics:** Domain-specific performance measures

3. Production Readiness

- **MLOps Integration:** Feature pipeline and monitoring requirements defined
- **Business Impact:** Cost analysis and ROI calculation framework
- **Clinical Integration:** Workflow and decision support requirements



How EDA Supports Next Steps

Week 2: Feature Engineering & Modeling

Immediate Actions:

1. **Feature Selection:** Choose top 20-30 features from 58+ available
2. **Data Preprocessing:** Handle missing values and encode categorical variables
3. **Baseline Models:** Train Logistic Regression, Random Forest, XGBoost
4. **Performance Evaluation:** Healthcare-specific metrics and validation



Key Insights for Stakeholders

Clinical Teams

- **Risk Stratification:** Clear patient categorization for care planning
- **Treatment Complexity:** Resource allocation guidance
- **Readmission Prevention:** High-risk patient identification
- **Clinical Decision Support:** Data-driven intervention strategies

Hospital Administration

- **Resource Planning:** Treatment complexity-based staffing
- **Cost Analysis:** Readmission cost vs. prevention investment
- **Quality Metrics:** Performance benchmarking by risk category
- **Strategic Planning:** Population health management insights

Success Metrics & Validation

EDA Quality Metrics

- **Data Understanding:** 100% features analyzed
- **Target Definition:** Clear binary classification established
- **Feature Engineering:** 8 clinically relevant features created
- **Risk Stratification:** Healthcare-standard categories defined
- **Documentation:** Comprehensive analysis and insights

Conclusion

The EDA phase has successfully transformed raw healthcare data into actionable clinical insights, creating a robust foundation for machine learning modeling. The comprehensive analysis demonstrates exceptional healthcare domain expertise and positions the project for successful development and production deployment.

Key Achievements:

1. **Established Clinical Understanding:** Deep insights into diabetes readmission patterns
2. **Created Actionable Features:** 8 new clinical features for modeling
3. **Defined Risk Categories:** Healthcare-standard risk stratification
4. **Identified Socioeconomic Factors:** Social determinants of health integration
5. **Prepared for Modeling:** Clear roadmap for Week 2 development



Next Steps

Immediate Actions (Week 2):

1. Feature selection from 8 engineered features
2. Data preprocessing pipeline implementation
3. Baseline model training and evaluation
4. Hyperparameter optimization with Optuna

Document Version: 1.0

Last Updated: August 2025

Next Review: Week 2 completion

Status: EDA Phase Complete 