

This document provides a clear explanation of the Synthetic Hospital Readmission Dataset.

Dataset Overview

This synthetic dataset simulates hospital admission records¹. Its purpose is to train and test an AI model that predicts the likelihood of a patient being readmitted to the hospital within 30 days². Each row in the dataset represents a single hospital admission event³.

Data Dictionary

The table below explains each column in the dataset⁴.

Column Name	Description (Plain English)	Importance in a Hospital Setting	Typical Real-World Source
PatientID	A unique, random code for each patient ⁵⁵⁵⁵ .	Allows patient records to be tracked without using personal details ⁶⁶⁶⁶ .	Hospital EHR system ⁷ .
Age	The patient's age in years ⁸ .	A strong predictor of readmission risk, as older patients often have more complications ⁹ .	Patient demographics database ¹⁰ .
Gender	The patient's gender (Male, Female, Other) ¹¹¹¹¹¹¹¹ .	Can influence disease prevalence, treatment response, and	Patient registration system ¹³¹³¹³¹³ .

	the stay (e.g., Insulin, Antibiotics) ²⁹ .	associated with a higher likelihood of readmission ³⁰ .	
BMI	Body Mass Index, a measure of body weight relative to height ³² .	A very high or low BMI can signal health risks that impact readmission ³³ .	Vital signs / nursing assessments ³⁴ .
SmokingStatus	The patient's smoking history (never smoked, former smoker, or current smoker) ³⁵ .	Smoking is a significant risk factor for many chronic diseases and complications ³⁶ .	Patient lifestyle history ³⁷ .
AlcoholUse	The patient's alcohol consumption (none, moderate, or heavy) ³⁸ .	Alcohol use can impact recovery, medication adherence, and the risk of readmission ³⁹ .	Patient lifestyle history ⁴⁰ .
BloodPressure	The patient's blood pressure upon admission (e.g., 130/85) ⁴¹ .	High or low blood pressure is a key clinical indicator for many conditions ⁴² .	Vital signs monitoring ⁴³ .
CholesterolLevel	The patient's cholesterol level in mg/dL ⁴⁴ .	High cholesterol is linked to cardiovascular diseases and an increased risk of readmission ⁴⁵ .	Laboratory results ⁴⁶ .

HbA1c	A blood test result showing the average blood sugar level over 3 months (%) ⁴⁷ .	Monitors diabetes control, which is a strong factor in readmission risk ⁴⁸ .	Laboratory results ⁴⁹ .
FollowUpAppointmentScheduled	Indicates if a follow-up appointment was scheduled after discharge (1 = Yes, 0 = No) ⁵⁰ .	Scheduling follow-up care helps ensure continuity of treatment and reduces readmission risk ⁵¹ .	Discharge planning records ⁵² .
InsuranceType	The patient's payment method (Public, Private, Self-Pay) ⁵³ .	Can influence a patient's access to care and follow-up services ⁵⁴ .	Billing / administrative systems ⁵⁵ .
ReadmittedWithin30Days	The outcome variable, indicating if the patient was readmitted within 30 days (1 = Yes, 0 = No) ⁵⁶ .	This is the target variable that the AI model is designed to predict ⁵⁷ .	Hospital readmission tracking ⁵⁸ .
RecordGeneratedAt	The exact date and time the synthetic data record was created ⁵⁹ .	Useful for tracking when the data was generated during testing ⁶⁰ .	Synthetic data generator ⁶¹ .

Dataset Schema and Structure

The data fields are organized into logical groups to provide a complete view of a patient's

admission⁶².

- **Demographics:** Basic information about the patient, such as age, gender, and insurance type⁶³.
- **Clinical Information:** Details about the hospital stay, including admission/discharge dates, diagnosis, treatments, and follow-up plans⁶⁴.
- **Lifestyle Factors:** Patient habits and body metrics that influence health, like smoking status, alcohol use, and BMI⁶⁵.
- **Lab & Vital Signs:** Objective clinical measurements taken during the hospital stay, including blood pressure and cholesterol levels⁶⁶.
- **Outcomes:** The target variable for prediction (readmission) and the record's creation timestamp⁶⁷.

Why This Data Matters for AI

The inclusion of detailed patient information makes the dataset more realistic and useful for AI model training⁶⁸.

- **Lifestyle Factors:** Fields like smoking, alcohol use, and BMI are strong predictors of chronic disease outcomes⁶⁹.
- **Lab Results:** Clinical data such as cholesterol, HbA1c, and blood pressure add clinical depth to the dataset⁷⁰.
- **Care Planning:** Information on follow-up appointments and insurance reflects real-world hospital processes that affect readmission risk⁷¹.

Guide to Modifying Synthetic Data Generation

This table explains how to adjust the parameters used to generate the synthetic data⁷².

To Change This...	Edit This Code...	Example Modification	Effect on Data
BMI range	<code>bmi = round(random.uniform(18.0, 40.0), 1)</code>	Change the 18.0 and 40.0 values.	Alters the weight distribution of the synthetic patient population ⁷³ .
Smoking status options	<code>smoking_status = random.choice([...])</code>	Add or remove items from the list, e.g., "Never", "Former", "Current".	Reflects different population habits ⁷⁴ .
Alcohol use options	<code>alcohol_use = random.choice([...])</code>	Add or remove items from the list, e.g., "None", "Moderate", "Heavy".	Adjusts the lifestyle risk factors in the dataset ⁷⁵ .
Blood pressure range	<code>blood_pressure = f"{random.randint(100, 160)}/{random.randint(60, 100)}"</code>	Change the number ranges for systolic (100-160) and diastolic (60-100) values.	Simulates patient populations with different blood pressure profiles ⁷⁶ .
Cholesterol range	<code>cholesterol = random.randint(150, 300)</code>	Adjust the minimum (150) and maximum (300) values.	Generates patients with higher or lower cholesterol levels ⁷⁷ .
HbA1c range	<code>hba1c = round(random.uniform(4.5, 12.0), 1)</code>	Adjust the 4.5 and 12.0 values.	Simulates different levels of diabetes severity and control ⁷⁸ .

Follow-up appointment rate	follow_up = random.choice([0, 1])	Change to a weighted choice to generate more 'Yes' or 'No' values.	Controls the likelihood of a follow-up appointment being scheduled ⁷⁹ .
Insurance types	insurance_type = random.choice([...])	Add or remove items from the list, e.g., "Public", "Private", "Self-Pay".	Reflects different healthcare payment systems ⁸⁰ .