# Predict readmission rate for diabetes patients

## Project introduction and background

- As part of the Affordable Health Care Act, hospitals are requested to improve the qualities of the health care services, one of which concerns with reducing readmission rate. As a result, Hospital Readmission Reduction Program (HRRP) is created to reduce reimbursement to hospitals with excessive readmissions. Hospitals need to account for this change and introduce some necessary measures to inspect their readmission situation. The first question the hospital want to find out is which patient and why they are readmitted. For further prevention, they want to predict whether their patients will be admitted again.

- Diabetes problem has affected a great population in the US. One in 10 people is diagnosed as having diabetes according to a report released by Ostling et al 2017. People who is diagnosed as diabetes are more likely to be readmitted to hospital stated in the report. Thus we try to predict the readmission rate for diabetes patients from a dataset containing the diabetes patient data from 130 US hospitals.

## Project Goal

We are aiming to build a classification model on predicting whether a diabetic patient will be readmitted to the hospital. The target varible *readmitted* is a binary object that shows whether the patient is readmitted within 30 days

## Dataset

The dataset is obtained from UCI repository with 101766 instances and 50 features. Relevant rearch use this dataset() The feature description, type, and data missingness are as followed:

| Feature name | Type | Description and values | % missing |
|---|---|---|---|
| Encounter ID | Numeric | Unique identifier of an encounter | 0% |
| Patient number | Numeric | Unique identifier of a patient | 0% |
| Race | Nominal | Values: Caucasian, Asian, African American, Hispanic, and other | 2% |
| Gender | Nominal | Values: male, female, and unknown/invalid | 0% |
| Age | Nominal | Grouped in 10-year intervals: [0, 10), [10, 20), . . ., [90, 100) | 0% |
| Weight | Numeric | Weight in pounds. | 97% |
| Admission type | Nominal | Integer identifier corresponding to 9 distinct values, for example, emergency, urgent, elective, newborn, and not available | 0% |
| Discharge disposition | Nominal | Integer identifier corresponding to 29 distinct values, for example, discharged to home, expired, and not available | 0% |
| Admission source | Nominal | Integer identifier corresponding to 21 distinct values, for example, physician referral, emergency room, and transfer from a hospital | 0% |
| Time in hospital | Numeric | Integer number of days between admission and discharge | 0% |
| Payer code | Nominal | Integer identifier corresponding to 23 distinct values, for example, Blue Cross\Blue Shield, Medicare, and self-pay | 52% |
| Medical specialty | Nominal | Integer identifier of a specialty of the admitting physician, corresponding to 84 distinct values, for example, cardiology, internal medicine, family\general practice, and surgeon | 53% |
| Number of lab procedures | Numeric | Number of lab tests performed during the encounter | 0% |
| Number of procedures | Numeric | Number of procedures (other than lab tests) performed during the encounter | 0% |
| Number of medications | Numeric | Number of distinct generic names administered during the encounter | 0% |
| Number of outpatient visits | Numeric | Number of outpatient visits of the patient in the year preceding the encounter | 0% |
| Number of emergency visits | Numeric | Number of emergency visits of the patient in the year preceding the encounter | 0% |
| Number of inpatient visits | Numeric | Number of inpatient visits of the patient in the year preceding the encounter | 0% |
| Diagnosis 1 | Nominal | The primary diagnosis (coded as first three digits of ICD9); 848 distinct values | 0% |
| Diagnosis 2 | Nominal | Secondary diagnosis (coded as first three digits of ICD9); 923 distinct values | 0% |
| Diagnosis 3 | Nominal | Additional secondary diagnosis (coded as first three digits of ICD9); 954 distinct values | 1% |
| Number of diagnoses | Numeric | Number of diagnoses entered to the system | 0% |
| Glucose serum test result | Nominal | Indicates the range of the result or if the test was not taken. Values: ">200," ">300," "normal," and "none" if not measured | 0% |
| A1c test result | Nominal | Indicates the range of the result or if the test was not taken. Values: ">8" if the result was greater than 8%, ">7" if the result was greater than 7% but less than 8%, "normal" if the result was less than 7%, and "none" if not measured. | 0% |
| Change of medications | Nominal | Indicates if there was a change in diabetic medications (either dosage or generic name). Values: "change" and "no change" | 0% |
| Diabetes medications | Nominal | Indicates if there was any diabetic medication prescribed. Values: "yes" and "no" | 0% |
| 24 features for medications | Nominal | For the generic names: metformin, repaglinide, nateglinide, chlorpropamide, glimepiride, acetohexamide, glipizide, glyburide, tolbutamide, pioglitazone, rosiglitazone, acarbose, miglitol, troglitazone, tolazamide, examide, sitagliptin, insulin, glyburide-metformin, glipizide-metformin, glimepiride-pioglitazone, metformin-rosiglitazone, and metformin-pioglitazone, the feature indicates whether the drug was prescribed or there was a change in the dosage. Values: "up" if the dosage was increased during the encounter, "down" if the dosage was decreased, "steady" if the dosage did not change, and "no" if the drug was not prescribed | 0% |
| Readmitted | Nominal | Days to inpatient readmission. Values: "<30" if the patient was readmitted in less than 30 days, ">30" if the patient was readmitted in more than 30 days, and "No" for no record of readmission. | 0% |

## Relevant projects

- In the research on the impact of HbA1c Measurement on Hospital Readmission Rate, Beata Strack et al(2014) tried to provide an assessment and to find future directions which might lead to improvements in diabetic patient safety. They used multivariable

logistic regression to fit the relationship between the measurement of HbA1c and early readmission while controlling for covariates such as demographics, severity and type of the disease, and type of admission. Results show that the measurement of HbA1c was performed infrequently (18.4%) in the inpatient setting. The statistical model suggests that the relationship between the probability of readmission and the HbA1c measurement depends on the primary diagnosis. The data suggest further that the greater attention to diabetes reflected in HbA1c determination may improve patient outcomes and lower cost of inpatient care.

## Data Preprocesing Procedures

- Inspect and drop features that have great missingness. For example, the feature *Weight* has more than 97% of missing values, we decided to drop this feature. Payer code is also irrelevant to our target variable and it is dropped in this dataset. For 24 features for medications, some features might have uniform answers and thus might not bring more information to this data set. We first check the unique values of these columns.

- Transform the diagnosis data to string according to diagnosis table

| Group name | icd9 codes | Number of encounters | % of encounter | Description |
|---|---|---|---|---|
| Circulatory | 390–459, 785 | 21,411 | 30.6% | Diseases of the circulatory system |
| Respiratory | 460–519, 786 | 9,490 | 13.6% | Diseases of the respiratory system |
| Digestive | 520–579, 787 | 6,485 | 9.3% | Diseases of the digestive system |
| Diabetes | 250.xx | 5,747 | 8.2% | Diabetes mellitus |
| Injury | 800–999 | 4,697 | 6.7% | Injury and poisoning |
| Musculoskeletal | 710–739 | 4,076 | 5.8% | Diseases of the musculoskeletal system and connective tissue |
| Genitourinary | 580–629, 788 | 3,435 | 4.9% | Diseases of the genitourinary system |
| Neoplasms | 140–239 | 2,536 | 3.6% | Neoplasms |
| | 780, 781, 784, 790–799 | 2,136 | 3.1% | Other symptoms, signs, and ill-defined conditions |
| | 240–279, without 250 | 1,851 | 2.6% | Endocrine, nutritional, and metabolic diseases and immunity disorders, without diabetes |
| | 680–709, 782 | 1,846 | 2.6% | Diseases of the skin and subcutaneous tissue |
| | 001–139 | 1,683 | 2.4% | Infectious and parasitic diseases |
| Other (17.3%) | 290–319 | 1,544 | 2.2% | Mental disorders |
| | E–V | 918 | 1.3% | External causes of injury and supplemental classification |
| | 280–289 | 652 | 0.9% | Diseases of the blood and blood-forming organs |
| | 320–359 | 634 | 0.9% | Diseases of the nervous system |
| | 630–679 | 586 | 0.8% | Complications of pregnancy, childbirth, and the puerperium |
| | 360–389 | 216 | 0.3% | Diseases of the sense organs |
| | 740–759 | 41 | 0.1% | Congenital anomalies |

- Uses Ordinal Encoder for categorical features like age. Medications are stated with values like 'NO','Steady, 'Up','Down', which shows the level of the severity of the

current situation of patients. With the use of Ordinal Encoder we could quantify patient's situation.
med_columns='metformin', 'repaglinide', 'nateglinide', 'chlorpropamide', 'glimepiride', 'acetohexamide', 'glipizide', 'glyburide', 'tolbutamide', 'pioglitazone', 'rosiglitazone', 'acarbose', 'miglitol', 'troglitazone', 'tolazamide', 'insulin', 'glyburide-metformin', 'glipizide-metformin', 'glimepiride-pioglitazone', 'metformin-rosiglitazone', 'metformin-pioglitazone'

- Uses OneHot Encoder to dummify some categorical features such as 'race', 'gender','medical_specialty','diag_1_group','diag_2_group','diag_3_group','max_glu_serum','A1Cresult'

- Inspect continous variables and choose to use MinMax and standard Scalar to transform data.
  MinMax_columns='age','num_procedures','number_diagnoses'
  std_columns='time_in_hospital','num_lab_procedures','num_medications','number_outpatient','number_emergency','number_inpatient'

- Uses Label Encoder to transform target variable 'Readmitted'

[ ]