Prediction of hospital readmission in diabetic inpatients

Data analytics and Data driven decision, 2017/2018 Università degli studi di L'Aquila

Matricola ID № 256921, Stoil Yasenov Yanchev, stoilyanchev@gmail.com Matricola ID № 256915, Suleyman Erim Matricola ID № 256916, Sena Er

Prof. Fabrizio Rossi

Prof. Giovanni Felici

2017/2018 г.

Content:

| 1 | INT | FRODUCTION | 3 |
|---|------|--|----|
| 2 | PRE | PROCESSING | 4 |
| | 2.1 | Dealing with missing values: | 6 |
| | 2.2 | Creating and/or Recoding New Features | 7 |
| | 2.3 | Categorization of diagnoses | 8 |
| | 2.4 | Collapsing some other variables | 8 |
| | 2.5 | Recoding some variables | 9 |
| | 2.6 | Recoding the outcome variable | 10 |
| | 2.7 | Dealing with age | 10 |
| | 2.8 | Collapsing of Multiple Encounters for same patient | 10 |
| | 2.9 | Log Transformation | 11 |
| | 2.10 | Standardization | 12 |
| 3 | INT | ERACTION TERMS | 12 |

1 Introduction

A hospital readmission is when a patient who is discharged from the hospital, gets re-admitted again within a certain period of time. Hospital readmission rates for certain conditions are now considered an indicator of hospital quality, and also affect the cost of care adversely. Although diabetes is not yet included in the penalty measures, this is a big problem. In 2011, American hospitals spent over 41 billion dollars on diabetic patients who got readmitted within 30 days of discharge. Being able to determine factors that lead to higher readmission in such patients, and correspondingly being able to predict which patients will get readmitted can help hospitals save millions of dollars while improving quality of care. So, with that background in mind, we used a medical claims dataset (description below), to answer these questions:

What factors are the strongest predictors of hospital readmission in diabetic patients?

How well can we predict hospital readmission in this dataset with limited features?

Choosing a dataset

Finding a good dataset is one of the first challenges (besides defining a meaningful question), when trying out machine learning methods. The current state of the healthcare world is such that we can easily find datasets that rich (full

of useful information) but dirty (unstructured content or messy schemas) or datasets that are very clean but otherwise sterile in terms of information contained.

With this limitation, we picked a publicly available dataset from UCI repository (link) containing de-identified diabetes patient encounter data for 130 US hospitals containing 101,766 observations over 10 years. The dataset has over 50 features including patient characteristics, conditions, tests and 23 medications. Only diabetic encounters are included (i.e. at least one of three primary diagnosis was diabetes). This dataset has been used by Strack et al. in 2014 for an interesting analysis on the same topic. So we begin by loading the dataset (csv file downloaded from the link above) as a pandas dataframe

2 Preprocessing

The first thing to do when picking a dataset is to do some data profiling and look at key features, as we have done in the table below:

| Dataset | | Total observations | | Total features | |
|---|--|---|--------|----------------|-------|
| Diabetes encounter data US-130 hospitals (1999-2008) | | 101,766 | | 55 | |
| Continuous variables | Min | Mean | Median | Max | SD |
| Time in hospital | 1 | 4.40 | 4 | 14 | 2.98 |
| # of lab procedures | 1 | 43.10 | 44 | 132 | 19.67 |
| # of procedures | 0 | 1.34 | 1 | 6 | 1.71 |
| # of medications | 1 | 16.02 | 15 | 81 | 8.13 |
| # of outpatient visits | 0 | 0.37 | 0 | 42 | 1.27 |
| # of emergency visits | 0 | 0.19 | 0 | 76 | 0.93 |
| # of admissions | 0 | 0.64 | 0 | 21 | 1.26 |
| # of diagnoses | 1 | 7.42 | 8 | 16 | 1.93 |
| Categorical variables (| Details | | | | |
| Medication change (outco | me) | No change = 53.8%, Change = 46.2% | | | |
| HbA1c test | None = 83.3%, >8 =8.1%, Norm = 4.9%, >7 = 3.7% | | | | |
| Race | Caucasian = 74.8%, African American = 18.9%, Missing = 2.2%, | | | | |
| | | Hispanic = 2.0%, Other = 1.5%, Asian = 0.6% | | | |
| Gender | Female = 53.8%, Male = 46.2% | | | | |
| Age category | Most frequent = 70-80 years = 25.6% | | | | |
| Readmission | No = 53.9%, >30 days = 34.9%, <30 days = 11.2% | | | | |

What preprocessing and feature engineering techniques should be applied?

Before we can get to actual modeling, some wrangling with the data is almost always needed. We applied three types of methods here:

- 1. Cleaning tasks such as dropping bad data, dealing with missing values.
- 2. Modification of existing features e.g. standardization, log transforms etc.
- 3. Creation or derivation of new features, usually from existing ones.

The individual steps are described in detail below. However, note that what looks like a nice sequence of steps now, was a result of many trial and error attempts to see what works well for getting our data into best shape.

2.1 Dealing with missing values:

This gives us a long list but the following variables had missing values:

Now the important part—deciding what to do:

Weight is missing in over 98% records. Owing to the poor interpretability of missing values and little predictive generalizability to other patients, best thing is to just drop it.

Medical Specialty of treating physician also have 40–50% missing values. We decided to drop this, but there are other ways too to deal with such missing values.

Primary (diag_1), Secondary (diag_2) and Additional (diag_3) diagnoses were have very few missing values. Technically, if all three are missing, that's bad data. So we only drop those records where all three diagnoses are missing.

Gender has only 3 missing or invalid values so we decided to drop these records.

Also, one more cleaning step that depends on understanding the data and some common sense: since we are trying to predict readmissions, those patients who died during this hospital admission, have zero probability of readmission. So we should remove those records (discharge_disposition = 11).

2.2 Creating and/or Recoding New Features

This is highly subjective, and partly depends on a knowledge of healthcare services, and making sense of the potential relationships between features. There are perhaps thousands of ways to try here. We tried some (none are perfect) and here's why.

Service utilization: The data contains variables for number of inpatient (admissions), emergency room visits and outpatient visits for a given patient in the previous one year. These are (crude) measures of how much hospital/clinic services a person has used in the past year. We added these three to create a new variable called service utilization (see figure below). The idea was to see which version gives us better results. Granted, we did not apply any special weighting to the three ingredients of service utilization but we wanted to try something simple at this stage.



Number of medication changes: The dataset contains 23 features for 23 drugs (or combos) which indicate for each of these, whether a change in that medication was made or not during the current hospital stay of patient. Medication change for diabetics upon admission has been shown

by previous research to be associated with lower readmission rates. We decided to count how many changes were made in total for each patient, and declared that a new feature. The reasoning here was to both simplify the model and possibly discover a relationship with number of changes regardless of which drug was changed.

Number of medication used: Another possibly related factor could be the total number of medications used by the patient (which may indicate severity of their condition and/or the intensity of care). So we created another feature by counting the medications used during the encounter.

2.3 Categorization of diagnoses

The dataset contained upto three diagnoses for a given patient (primary, secondary and additional). However, each of these had 700–900 unique ICD codes and it is extremely difficult to include them in the model and interpret meaningfully. Therefore, we collapsed these diagnosis codes into 9 disease categories in an almost similar fashion to that done in the original publication using this dataset. These 9 categories include Circulatory, Respiratory, Digestive, Diabetes, Injury, Musculoskeletal, Genitourinary, Neoplasms, and Others. Although we did this for primary, secondary and additional diagnoses, we eventually decided to use only the primary diagnosis in our model.

2.4 Collapsing some other variables

Just like diagnoses, there were quite a few categories for admission source, admission type and discharge disposition. We collapsed these variables into fewer categories where it made sense. For example, admission types 1, 2 and 7 correspond to Emergency, Urgent Care and Trauma, and thus were combined into a single category as these are all non-elective situations.

```
In [24]: df['admission_type_id'] = df['admission_type_id'].replace(2,1)
    df['admission_type_id'] = df['admission_type_id'].replace(7,1)
    df['admission_type_id'] = df['admission_type_id'].replace(6,5)
    df['admission_type_id'] = df['admission_type_id'].replace(8,5)
```

2.5 Recoding some variables

The original dataset used string values for gender, race, medication change, and each of the 23 drugs used. To better fit those variables into our model, we interpret the variables to numeric binary variables to reflect their nature. For example, we encoded the "medication change" feature from "No" (no change) and "Ch" (changed) into 0 and 1.

```
In [25]: df['change'] = df['change'].replace('Ch', 1)
    df['change'] = df['change'].replace('No', 0)
    df['gender'] = df['gender'].replace('Male', 1)
    df['gender'] = df['gender'].replace('Female', 0)
    df['diabetesMed'] = df['diabetesMed'].replace('Yes', 1)
    df['diabetesMed'] = df['diabetesMed'].replace('No', 0)

# keys is the same as before
for col in keys:
    df[col] = df[col].replace('No', 0)
    df[col] = df[col].replace('Steady', 1)
    df[col] = df[col].replace('Up', 1)
    df[col] = df[col].replace('Down', 1)
```

We also reduced both A1C test result and Glucose serum test result into categories of Normal, Abnormal and Not tested.

```
In [26]: df['A1Cresult'] = df['A1Cresult'].replace('>7', 1)
    df['A1Cresult'] = df['A1Cresult'].replace('>8', 1)
    df['A1Cresult'] = df['A1Cresult'].replace('Norm', 0)
    df['A1Cresult'] = df['A1Cresult'].replace('None', -99)

df['max_glu_serum'] = df['max_glu_serum'].replace('>200', 1)
    df['max_glu_serum'] = df['max_glu_serum'].replace('>300', 1)
    df['max_glu_serum'] = df['max_glu_serum'].replace('Norm', 0)
    df['max_glu_serum'] = df['max_glu_serum'].replace('Norm', -99)
```

2.6 Recoding the outcome variable

The outcome we are looking at is whether the patient gets readmitted to the hospital within 30 days or not. The variable actually has < 30, > 30 and No Readmission categories. To reduce our problem to a binary classification, we combined the readmission after 30 days and no readmission into a single category:

```
In [27]: df['readmitted'] = df['readmitted'].replace('>30', 0)
    df['readmitted'] = df['readmitted'].replace('<30', 1)
    df['readmitted'] = df['readmitted'].replace('NO', 0)</pre>
```

2.7 Dealing with age

There are different ways to deal with this. The dataset only gives us age as 10 year categories, so we don't know the exact age of each patient. The previous study on this dataset used age categories as nominal variables, but we wanted to be able to see the effect of increasing age on readmission, even if in a crude way. To do that, we assume that age of the patient on average lies at the midpoint of the age category. For example, if the patient's age category is 20–30 years, then we assume the age = 25 years. So we converted age categories to midpoints, resulting in a numeric variable:

2.8 Collapsing of Multiple Encounters for same patient

Some patients in the dataset had more than one encounter. We could not count them as independent encounters because that bias the results towards those patients who had multiple encounters. Thus we tried multiple techniques to collapse and consolidate multiple encounters for same patient such as:

- 1. Considering more than 2 readmissions across multiple encounters as readmission for collapsed record.
- 2. Considering average stay at hospital across multiple encounters.
- 3. Considering the percentage of the medication changes across multiple encounters
- 4. Considering the total number of the encounters to replace the encounter unique ID
- 5. Considering the combination of diagnoses across multiple encounters as a list

2.9 Log Transformation

A preliminary analysis of our numerical features revealed that many of these were highly skewed and had high kurtosis. As a reference, the skew of a normal distribution is 0 and the excess kurtosis (difference of actual kurtosis from ideal normal distribution value of 3), as returned by the kurtosis() function for a normal distribution is 0, which would impact standardization. Features such as number of emergency visits, service utilization, number of inpatient admissions and number of outpatient visits had high skew and kurtosis. Thus, we performed log transformation where

a skew or kurtosis beyond the limits of $-2 \le$ skew and kurtosis ≤ 2 . Also, since log (0) is not defined, we decided to use the following rule:

- 1. Compute log(x) for any feature x if percentage of 0s in $x \le 2\%$, after removing the zeros. This ensured that we didn't bulk-remove records that hold predictive power for other columns.
- 2. Compute log1p(x) otherwise (log1p(x) means log(x+1), while retaining the zeros.

2.10 Standardization

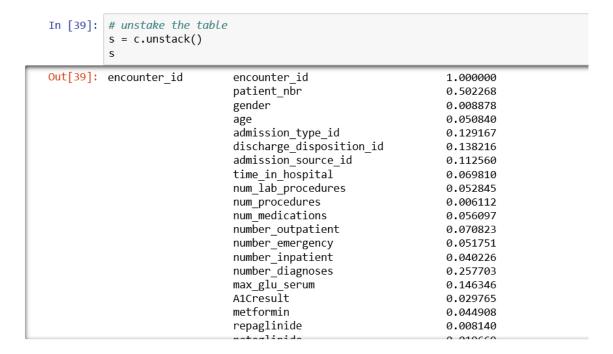
Since we had used log transformation to ensure that the numeric variables had a Gaussian-like or normal distribution (before the log transformation) or were log transformed to ensure a normal distribution, we decided to standardize our numerical features using the formula:

$$New value = \frac{Value - Mean(Values)}{Standard Deviation (Values)}$$

3 Interaction Terms

Variables can have interdependent effects on readmission, called interactions. We can identify possible candidates for interaction terms by looking at what makes theoretical sense and by observing a correlation matrix of the predictor variables to see which ones seem highly correlated.





In this list, there are two kinds of situations:

1. One variable is contained in/derivative of another: In the above examples, number of outpatient visits is part of service utilization. However, we created this feature ourselves, so in this case our decision is to not put these in same feature set (we used two feature sets described later). Similarly, diabetesMed (any diabetic medication prescribed) is contained in the number of medications used. We decided to drop

diabetesMed from analysis in this case because, well, it is understood that all of these patients are getting some diabetic medication.

2. Possible actual co-variance: The other situation is an actual co-variance between two variables. This seems to be the case for number of medication and whether a change was made or not, and it also makes some intuitive sense. So we created interaction terms for such cases.