

Prediction on Diabetes Patient's Hospital Readmission¶



Problem Statement and Objective 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. For this reason, Centers for Medicare & Medicaid Services established the Hospital Readmissions Reduction Program which aims to improve quality of care for patients and reduce health care spending by applying payment penalties to hospitals that have more than expected readmission rates for certain conditions. Although diabetes is not yet included in the penalty measures, the program is regularly adding new disease conditions to the list, now totaling 6 for FY2018. In 2011, American hospitals spent over \$41 billion 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:

1. What factors are the strongest predictors of hospital readmission in diabetic patients?
2. How well can we predict hospital readmission in this dataset with limited features?

Data Set Description¶

VARIABLE NAMES: DESCRIPTION

- **Encounter ID** Unique identifier of an encounter
- **Patient number** Unique identifier of a patient
- **Race** Values: Caucasian, Asian, African American, Hispanic, and other
- **Gender** Values: male, female, and unknown/invalid
- **Age** Grouped in 10-year intervals: 0, 10), 10, 20), ..., 90, 100)
- **Weight** Weight in pounds
- **Admission type** Integer identifier corresponding to 9 distinct values, for example, emergency, urgent, elective, newborn, and not available
- **Discharge disposition** Integer identifier corresponding to 29 distinct values, for example, discharged to home, expired, and not available
- **Admission source** Integer identifier corresponding to 21 distinct values, for example, physician referral, emergency room, and transfer from a hospital
- **Time in hospital** Integer number of days between admission and discharge
- **Payer code** Integer identifier corresponding to 23 distinct values, for example, Blue Cross/Blue Shield, Medicare, and self-pay Medical

- **Medical specialty** Integer identifier of a specialty of the admitting physician, corresponding to 84 distinct values, for example, cardiology, internal medicine, family/general practice, and surgeon
- **Number of lab procedures** Number of lab tests performed during the encounter
- **Number of procedures** Numeric Number of procedures (other than lab tests) performed during the encounter
- **Number of medications** Number of distinct generic names administered during the encounter
- **Number of outpatient visits** Number of outpatient visits of the patient in the year preceding the encounter
- **Number of emergency visits** Number of emergency visits of the patient in the year preceding the encounter
- **Number of inpatient visits** Number of inpatient visits of the patient in the year preceding the encounter
- **Diagnosis 1** The primary diagnosis (coded as first three digits of ICD9); 848 distinct values
- **Diagnosis 2** Secondary diagnosis (coded as first three digits of ICD9); 923 distinct values
- **Diagnosis 3** Additional secondary diagnosis (coded as first three digits of ICD9); 954 distinct values
- **Number of diagnoses** Number of diagnoses entered to the system 0%
- **Glucose serum test result** Indicates the range of the result or if the test was not taken. Values: ">200," ">300," "normal," and "none" if not measured
- **A1c test result** 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.
- **Change of medications** Indicates if there was a change in diabetic medications (either dosage or generic name). Values: "change" and "no change"
- **Diabetes medications** Indicates if there was any diabetic medication prescribed. Values: "yes" and "no"
- 24 features for medications 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
- **Readmitted** 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

Data Preparation & Exploration

```
#Loading libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

```
#loading Dataset
df = pd.read_csv("../input/diabetic_data.csv")
```

```
#displaying first 10 rows of data
df.head(10).T
```

[illegible]

	0	1	2	3	4	5	6	
tolbutamide	No	No	No	No	No	No	No	No
pioglitazone	No	No	No	No	No	No	No	No
rosiglitazone	No	No	No	No	No	No	No	No
acarbose	No	No	No	No	No	No	No	No
miglitol	No	No	No	No	No	No	No	No
troglitazone	No	No	No	No	No	No	No	No
tolazamide	No	No	No	No	No	No	No	No
examide	No	No	No	No	No	No	No	No
citoglipton	No	No	No	No	No	No	No	No
insulin	No	Up	No	Up	Steady	Steady	Steady	No
glyburide-metformin	No	No	No	No	No	No	No	No
glipizide-metformin	No	No	No	No	No	No	No	No
glimepiride-pioglitazone	No	No	No	No	No	No	No	No
metformin-rosiglitazone	No	No	No	No	No	No	No	No
metformin-pioglitazone	No	No	No	No	No	No	No	No
change	No	Ch	No	Ch	Ch	No	Ch	No
diabetesMed	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes
readmitted	NO	>30	NO	NO	NO	>30	NO	>30

In [4]:

```
#checking shape of the dataset
df.shape
```

Out[4]:

```
(101766, 50)
```

In [5]:

```
#Checking data types of each variable
df.dtypes
```

Out[5]:

```
encounter_id          int64
patient_nbr           int64
race                  object
gender                object
age                  object
weight               object
admission_type_id     int64
discharge_disposition_id int64
admission_source_id   int64
time_in_hospital      int64
payer_code            object
medical_specialty     object
num_lab_procedures    int64
num_procedures        int64
num_medications       int64
number_outpatient     int64
number_emergency      int64
number_inpatient      int64
diag_1                object
diag_2                object
diag_3                object
number_diagnoses      int64
max_glu_serum         object
A1Cresult             object
metformin             object
repaglinide           object
nateglinide           object
```

chlorpropamide	object
glimepiride	object
acetoexamide	object
glipizide	object
glyburide	object
tolbutamide	object
pioglitazone	object
rosiglitazone	object
acarbose	object
miglitol	object
troglitazone	object
tolazamide	object
examide	object
citoglipton	object
insulin	object
glyburide-metformin	object
glipizide-metformin	object
glimepiride-pioglitazone	object
metformin-rosiglitazone	object
metformin-pioglitazone	object
change	object
diabetesMed	object
readmitted	object
dtype:	object

In [6]:

```
#Checking for missing values in dataset

#In the dataset missing values are represented as '?' sign

for col in df.columns:

    if df[col].dtype == object:

        print(col,df[col][df[col] == '?'].count())
```

```
race 2273
gender 0
age 0
weight 98569
payer_code 40256
medical_specialty 49949
diag_1 21
diag_2 358
diag_3 1423
max_glu_serum 0
A1Cresult 0
metformin 0
repaglinide 0
nateglinide 0
chlorpropamide 0
glimepiride 0
acetoexamide 0
glipizide 0
glyburide 0
tolbutamide 0
pioglitazone 0
rosiglitazone 0
acarbose 0
miglitol 0
troglitazone 0
tolazamide 0
examide 0
citoglipton 0
insulin 0
glyburide-metformin 0
glipizide-metformin 0
glimepiride-pioglitazone 0
```

```
metformin-rosiglitazone 0
metformin-pioglitazone 0
change 0
diabetesMed 0
readmitted 0
```

In [7]:

```
# gender was coded differently so we use a custom count for this one
print('gender', df['gender'][df['gender'] == 'Unknown/Invalid'].count())
```

```
gender 3
```

Dealing with Missing Values¶

Variable weight contains approximate 98% of the missing values so there is no significance in filling those missing values so we decided to drop these variables. Variable Payer code and medical specialty contains approximate 40% missing values so we also dropped these variables. Variables race, diag_1, diag_2, diag_3 and gender contains very less missing values as compared to other attributes which we dropped so for these attributes we also decided to drop those where missing values contains.

In [8]:

```
#dropping columns with large number of missing values
df = df.drop(['weight', 'payer_code', 'medical_specialty'], axis = 1)
```

In [9]:

```
drop_idx = set(df[(df['diag_1'] == '?') & (df['diag_2'] == '?') & (df['diag_3'] == '?')].index)

drop_idx = drop_idx.union(set(df['diag_1'][df['diag_1'] == '?'].index))
drop_idx = drop_idx.union(set(df['diag_2'][df['diag_2'] == '?'].index))
drop_idx = drop_idx.union(set(df['diag_3'][df['diag_3'] == '?'].index))
drop_idx = drop_idx.union(set(df['race'][df['race'] == '?'].index))
drop_idx = drop_idx.union(set(df[df['discharge_disposition_id'] == 11].index))
drop_idx = drop_idx.union(set(df['gender'][df['gender'] == 'Unknown/Invalid'].index))

new_idx = list(set(df.index) - set(drop_idx))
df = df.iloc[new_idx]
```

variables (drugs named citoglipton and examide), all records have the same value. So essentially these cannot provide any interpretive or discriminatory information for predicting readmission so we decided to drop these two variables

In [10]:

```
df = df.drop(['citoglipton', 'examide'], axis = 1)
```

In [11]:

```
#Checking for missing values in the data
for col in df.columns:
    if df[col].dtype == object:
        print(col,df[col][df[col] == '?'].count())

print('gender', df['gender'][df['gender'] == 'Unknown/Invalid'].count())
```

```
race 0
gender 0
age 0
diag_1 0
diag_2 0
diag_3 0
max_glu_serum 0
AlCresult 0
metformin 0
repaglinide 0
nateglinide 0
chlorpropamide 0
glimepiride 0
acetoexamide 0
glipizide 0
glyburide 0
tolbutamide 0
pioglitazone 0
rosiglitazone 0
acarbose 0
miglitol 0
troglitazone 0
tolazamide 0
insulin 0
glyburide-metformin 0
glipizide-metformin 0
glimepiride-pioglitazone 0
metformin-rosiglitazone 0
metformin-pioglitazone 0
change 0
diabetesMed 0
readmitted 0
gender 0
```

Feature Engineering

This is highly subjective, and partly depends on a knowledge of health care services, and making sense of the potential relationships between features. There are perhaps thousands of ways to try here. We tried some...

- 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.

In [12]:

```
df['service_utilization'] = df['number_outpatient'] + df['number_emergency'] + df['number_inpatient']
```

- 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.

In [13]:

```
keys = ['metformin', 'repaglinide', 'nateglinide', 'chlorpropamide', 'glimepiride', 'glipizide', 'glyburide', 'pioglitazone', 'rosiglitazone', 'acarbose', 'miglitol', 'insulin', 'glyburide-metformin', 'tolazamide', 'metformin-pioglitazone', 'metformin-rosiglitazone', 'glimepiride-pioglitazone', 'glipizide-metformin', 'troglitazone', 'tolbutamide', 'acetohehexamide']

for col in keys:
    colname = str(col) + 'temp'
    df[colname] = df[col].apply(lambda x: 0 if (x == 'No' or x == 'Steady') else 1)
df['numchange'] = 0

for col in keys:
    colname = str(col) + 'temp'
    df['numchange'] = df['numchange'] + df[colname]
    del df[colname]

df['numchange'].value_counts()
```

Out[13]:

```
0    70142
1    24922
2     1271
3      106
4         5
Name: numchange, dtype: int64
```

In [14]:

```
# re-encoding admission type, discharge type and admission source into fewer categories

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)

df['discharge_disposition_id'] = df['discharge_disposition_id'].replace(6,1)
df['discharge_disposition_id'] = df['discharge_disposition_id'].replace(8,1)
df['discharge_disposition_id'] = df['discharge_disposition_id'].replace(9,1)
df['discharge_disposition_id'] = df['discharge_disposition_id'].replace(13,1)
df['discharge_disposition_id'] = df['discharge_disposition_id'].replace(3,2)
df['discharge_disposition_id'] = df['discharge_disposition_id'].replace(4,2)
df['discharge_disposition_id'] = df['discharge_disposition_id'].replace(5,2)
df['discharge_disposition_id'] = df['discharge_disposition_id'].replace(14,2)
df['discharge_disposition_id'] = df['discharge_disposition_id'].replace(22,2)
df['discharge_disposition_id'] = df['discharge_disposition_id'].replace(23,2)
df['discharge_disposition_id'] = df['discharge_disposition_id'].replace(24,2)
df['discharge_disposition_id'] = df['discharge_disposition_id'].replace(12,10)
df['discharge_disposition_id'] = df['discharge_disposition_id'].replace(15,10)
df['discharge_disposition_id'] = df['discharge_disposition_id'].replace(16,10)
df['discharge_disposition_id'] = df['discharge_disposition_id'].replace(17,10)
df['discharge_disposition_id'] = df['discharge_disposition_id'].replace(25,18)
df['discharge_disposition_id'] = df['discharge_disposition_id'].replace(26,18)

df['admission_source_id'] = df['admission_source_id'].replace(2,1)
df['admission_source_id'] = df['admission_source_id'].replace(3,1)
df['admission_source_id'] = df['admission_source_id'].replace(5,4)
df['admission_source_id'] = df['admission_source_id'].replace(6,4)
df['admission_source_id'] = df['admission_source_id'].replace(10,4)
df['admission_source_id'] = df['admission_source_id'].replace(22,4)
df['admission_source_id'] = df['admission_source_id'].replace(25,4)
df['admission_source_id'] = df['admission_source_id'].replace(15,9)
df['admission_source_id'] = df['admission_source_id'].replace(17,9)
df['admission_source_id'] = df['admission_source_id'].replace(20,9)
df['admission_source_id'] = df['admission_source_id'].replace(21,9)
df['admission_source_id'] = df['admission_source_id'].replace(13,11)
df['admission_source_id'] = df['admission_source_id'].replace(14,11)

```

- Encoding 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 [15]:

```
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 [16]:

```
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('None', -99)
```

- 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:

In [17]:

```
# code age intervals [0-10) - [90-100) from 1-10
for i in range(0,10):
    df['age'] = df['age'].replace([''+str(10*i)+'-'+str(10*(i+1))+')', i+1)
df['age'].value_counts()
```

Out[17]:

```
8      24815
7      21521
6      16546
9      16223
5       9208
4       3538
10     2594
3      1471
2       466
1        64
Name: age, dtype: int64
```

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:

- Considering more than 2 readmissions across multiple encounters as readmission for collapsed record.
- Considering average stay at hospital across multiple encounters.
- Considering the percentage of the medication changes across multiple encounters
- Considering the total number of the encounters to replace the encounter unique ID
- Considering the combination of diagnoses across multiple encounters as a list However, taking the features such as “diagnosis”, for instance, we did not find it not meaningful to combine multiple categorical values into an array for building data model. We then considered first encounter and last encounter separately as possible representations of multiple encounters. However, last encounters gave extremely imbalanced data for readmissions (96/4 Readmissions vs No Readmissions) and thus, we decided to use first encounters of patients with multiple encounters. This resulted in dataset being reduced to about 70,000 encounters:

In [18]:

```
df2 = df.drop_duplicates(subset= ['patient_nbr'], keep = 'first')
df2.shape
(70442, 55)
```

Out[18]:

```
(70442, 55)
```

In [19]:

```
df.head().T
```

Out[19]:

	1	2	3	4	5
encounter_id	149190	64410	500364	16680	35754
patient_nbr	55629189	86047875	82442376	42519267	82637451
race	Caucasian	AfricanAmerican	Caucasian	Caucasian	Caucasian
gender	0	0	1	1	1
age	2	3	4	5	6
admission_type_id	1	1	1	1	1

	1	2	3	4	5
discharge_disposition_id	1	1	1	1	1
admission_source_id	7	7	7	7	1
time_in_hospital	3	2	2	1	3
num_lab_procedures	59	11	44	51	31
num_procedures	0	5	1	0	6
num_medications	18	13	16	8	16
number_outpatient	0	2	0	0	0
number_emergency	0	0	0	0	0
number_inpatient	0	1	0	0	0
diag_1	276	648	8	197	414
diag_2	250.01	250	250.43	157	411
diag_3	255	V27	403	250	250
number_diagnoses	9	6	7	5	9
max_glu_serum	-99	-99	-99	-99	-99
A1Cresult	-99	-99	-99	-99	-99
metformin	0	0	0	0	0
repaglinide	0	0	0	0	0
nateglinide	0	0	0	0	0
chlorpropamide	0	0	0	0	0
glimepiride	0	0	0	0	0
acetohexamide	0	0	0	0	0
glipizide	0	1	0	1	0
glyburide	0	0	0	0	0
tolbutamide	0	0	0	0	0
pioglitazone	0	0	0	0	0
rosiglitazone	0	0	0	0	0
acarbose	0	0	0	0	0
miglitol	0	0	0	0	0
troglitazone	0	0	0	0	0
tolazamide	0	0	0	0	0
insulin	1	0	1	1	1
glyburide-metformin	0	0	0	0	0
glipizide-metformin	0	0	0	0	0
glimepiride-pioglitazone	0	0	0	0	0
metformin-rosiglitazone	0	0	0	0	0
metformin-pioglitazone	0	0	0	0	0
change	1	0	1	1	0
diabetesMed	1	1	1	1	1
readmitted	>30	NO	NO	NO	>30
service_utilization	0	3	0	0	0
numchange	1	0	1	0	0

- Encoding 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 [20]:

```
df['readmitted'].value_counts()
```

Out[20]:

NO 50731

```
>30      34649
<30      11066
Name: readmitted, dtype: int64
```

In [21]:

```
df['readmitted'] = df['readmitted'].replace('>30', 0)
df['readmitted'] = df['readmitted'].replace('<30', 1)
df['readmitted'] = df['readmitted'].replace('NO', 0)
```

- **Categorization of diagnoses:** The dataset contained up to 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. Doing this in python was slightly cumbersome because, well, we are mapping the disease codes to certain category names. Below code should demonstrate this easily.

In [22]:

```
# Creating additional columns for diagnosis# Creati
df['level1_diag1'] = df['diag_1']
df['level2_diag1'] = df['diag_1']
df['level1_diag2'] = df['diag_2']
df['level2_diag2'] = df['diag_2']
df['level1_diag3'] = df['diag_3']
df['level2_diag3'] = df['diag_3']
```

In [23]:

```
df.loc[df['diag_1'].str.contains('V'), ['level1_diag1', 'level2_diag1']] = 0
df.loc[df['diag_1'].str.contains('E'), ['level1_diag1', 'level2_diag1']] = 0
df.loc[df['diag_2'].str.contains('V'), ['level1_diag2', 'level2_diag2']] = 0
df.loc[df['diag_2'].str.contains('E'), ['level1_diag2', 'level2_diag2']] = 0
df.loc[df['diag_3'].str.contains('V'), ['level1_diag3', 'level2_diag3']] = 0
df.loc[df['diag_3'].str.contains('E'), ['level1_diag3', 'level2_diag3']] = 0
df['level1_diag1'] = df['level1_diag1'].replace('?', -1)
df['level2_diag1'] = df['level2_diag1'].replace('?', -1)
df['level1_diag2'] = df['level1_diag2'].replace('?', -1)
df['level2_diag2'] = df['level2_diag2'].replace('?', -1)
df['level1_diag3'] = df['level1_diag3'].replace('?', -1)
```

```
df['level2_diag3'] = df['level2_diag3'].replace('?', -1)
```

In [24]:

```
df['level1_diag1'] = df['level1_diag1'].astype(float)
df['level2_diag1'] = df['level2_diag1'].astype(float)
df['level1_diag2'] = df['level1_diag2'].astype(float)
df['level2_diag2'] = df['level2_diag2'].astype(float)
df['level1_diag3'] = df['level1_diag3'].astype(float)
df['level2_diag3'] = df['level2_diag3'].astype(float)
```

In [25]:

```
for index, row in df.iterrows():
    if (row['level1_diag1'] >= 390 and row['level1_diag1'] < 460) or (np.floor(row[
'level1_diag1']) == 785):
        df.loc[index, 'level1_diag1'] = 1
    elif (row['level1_diag1'] >= 460 and row['level1_diag1'] < 520) or (np.floor(ro
w['level1_diag1']) == 786):
        df.loc[index, 'level1_diag1'] = 2
    elif (row['level1_diag1'] >= 520 and row['level1_diag1'] < 580) or (np.floor(ro
w['level1_diag1']) == 787):
        df.loc[index, 'level1_diag1'] = 3
    elif (np.floor(row['level1_diag1']) == 250):
        df.loc[index, 'level1_diag1'] = 4
    elif (row['level1_diag1'] >= 800 and row['level1_diag1'] < 1000):
        df.loc[index, 'level1_diag1'] = 5
    elif (row['level1_diag1'] >= 710 and row['level1_diag1'] < 740):
        df.loc[index, 'level1_diag1'] = 6
    elif (row['level1_diag1'] >= 580 and row['level1_diag1'] < 630) or (np.floor(ro
w['level1_diag1']) == 788):
        df.loc[index, 'level1_diag1'] = 7
    elif (row['level1_diag1'] >= 140 and row['level1_diag1'] < 240):
        df.loc[index, 'level1_diag1'] = 8
    else:
        df.loc[index, 'level1_diag1'] = 0

    if (row['level1_diag2'] >= 390 and row['level1_diag2'] < 460) or (np.floor(row[
'level1_diag2']) == 785):
        df.loc[index, 'level1_diag2'] = 1
    elif (row['level1_diag2'] >= 460 and row['level1_diag2'] < 520) or (np.floor(ro
w['level1_diag2']) == 786):
        df.loc[index, 'level1_diag2'] = 2
    elif (row['level1_diag2'] >= 520 and row['level1_diag2'] < 580) or (np.floor(ro
w['level1_diag2']) == 787):
```

```

df.loc[index, 'level1_diag2'] = 3
elif (np.floor(row['level1_diag2']) == 250):
    df.loc[index, 'level1_diag2'] = 4
elif (row['level1_diag2'] >= 800 and row['level1_diag2'] < 1000):
    df.loc[index, 'level1_diag2'] = 5
elif (row['level1_diag2'] >= 710 and row['level1_diag2'] < 740):
    df.loc[index, 'level1_diag2'] = 6
elif (row['level1_diag2'] >= 580 and row['level1_diag2'] < 630) or (np.floor(row['level1_diag2']) == 788):
    df.loc[index, 'level1_diag2'] = 7
elif (row['level1_diag2'] >= 140 and row['level1_diag2'] < 240):
    df.loc[index, 'level1_diag2'] = 8
else:
    df.loc[index, 'level1_diag2'] = 0

if (row['level1_diag3'] >= 390 and row['level1_diag3'] < 460) or (np.floor(row['level1_diag3']) == 785):
    df.loc[index, 'level1_diag3'] = 1
elif (row['level1_diag3'] >= 460 and row['level1_diag3'] < 520) or (np.floor(row['level1_diag3']) == 786):
    df.loc[index, 'level1_diag3'] = 2
elif (row['level1_diag3'] >= 520 and row['level1_diag3'] < 580) or (np.floor(row['level1_diag3']) == 787):
    df.loc[index, 'level1_diag3'] = 3
elif (np.floor(row['level1_diag3']) == 250):
    df.loc[index, 'level1_diag3'] = 4
elif (row['level1_diag3'] >= 800 and row['level1_diag3'] < 1000):
    df.loc[index, 'level1_diag3'] = 5
elif (row['level1_diag3'] >= 710 and row['level1_diag3'] < 740):
    df.loc[index, 'level1_diag3'] = 6
elif (row['level1_diag3'] >= 580 and row['level1_diag3'] < 630) or (np.floor(row['level1_diag3']) == 788):
    df.loc[index, 'level1_diag3'] = 7
elif (row['level1_diag3'] >= 140 and row['level1_diag3'] < 240):
    df.loc[index, 'level1_diag3'] = 8
else:
    df.loc[index, 'level1_diag3'] = 0

```

In [26]:

```

for index, row in df.iterrows():
    if (row['level2_diag1'] >= 390 and row['level2_diag1'] < 399):
        df.loc[index, 'level2_diag1'] = 1
    elif (row['level2_diag1'] >= 401 and row['level2_diag1'] < 415):

```

```

df.loc[index, 'level2_diag1'] = 2
elif (row['level2_diag1'] >= 415 and row['level2_diag1'] < 460):
    df.loc[index, 'level2_diag1'] = 3
elif (np.floor(row['level2_diag1']) == 785):
    df.loc[index, 'level2_diag1'] = 4
elif (row['level2_diag1'] >= 460 and row['level2_diag1'] < 489):
    df.loc[index, 'level2_diag1'] = 5
elif (row['level2_diag1'] >= 490 and row['level2_diag1'] < 497):
    df.loc[index, 'level2_diag1'] = 6
elif (row['level2_diag1'] >= 500 and row['level2_diag1'] < 520):
    df.loc[index, 'level2_diag1'] = 7
elif (np.floor(row['level2_diag1']) == 786):
    df.loc[index, 'level2_diag1'] = 8
elif (row['level2_diag1'] >= 520 and row['level2_diag1'] < 530):
    df.loc[index, 'level2_diag1'] = 9
elif (row['level2_diag1'] >= 530 and row['level2_diag1'] < 544):
    df.loc[index, 'level2_diag1'] = 10
elif (row['level2_diag1'] >= 550 and row['level2_diag1'] < 554):
    df.loc[index, 'level2_diag1'] = 11
elif (row['level2_diag1'] >= 555 and row['level2_diag1'] < 580):
    df.loc[index, 'level2_diag1'] = 12
elif (np.floor(row['level2_diag1']) == 787):
    df.loc[index, 'level2_diag1'] = 13
elif (np.floor(row['level2_diag1']) == 250):
    df.loc[index, 'level2_diag1'] = 14
elif (row['level2_diag1'] >= 800 and row['level2_diag1'] < 1000):
    df.loc[index, 'level2_diag1'] = 15
elif (row['level2_diag1'] >= 710 and row['level2_diag1'] < 740):
    df.loc[index, 'level2_diag1'] = 16
elif (row['level2_diag1'] >= 580 and row['level2_diag1'] < 630):
    df.loc[index, 'level2_diag1'] = 17
elif (np.floor(row['level2_diag1']) == 788):
    df.loc[index, 'level2_diag1'] = 18
elif (row['level2_diag1'] >= 140 and row['level2_diag1'] < 240):
    df.loc[index, 'level2_diag1'] = 19
    elif row['level2_diag1'] >= 240 and row['level2_diag1'] < 280 and (np.floor(row
['level2_diag1']) != 250):
        df.loc[index, 'level2_diag1'] = 20
    elif (row['level2_diag1'] >= 680 and row['level2_diag1'] < 710) or (np.floor(ro
w['level2_diag1']) == 782):
        df.loc[index, 'level2_diag1'] = 21
    elif (row['level2_diag1'] >= 290 and row['level2_diag1'] < 320):

```



```
df.loc[index, 'level2_diag1'] = 22
else:
    df.loc[index, 'level2_diag1'] = 0

if (row['level2_diag2'] >= 390 and row['level2_diag2'] < 399):
    df.loc[index, 'level2_diag2'] = 1
elif (row['level2_diag2'] >= 401 and row['level2_diag2'] < 415):
    df.loc[index, 'level2_diag2'] = 2
elif (row['level2_diag2'] >= 415 and row['level2_diag2'] < 460):
    df.loc[index, 'level2_diag2'] = 3
elif (np.floor(row['level2_diag2']) == 785):
    df.loc[index, 'level2_diag2'] = 4
elif (row['level2_diag2'] >= 460 and row['level2_diag2'] < 489):
    df.loc[index, 'level2_diag2'] = 5
elif (row['level2_diag2'] >= 490 and row['level2_diag2'] < 497):
    df.loc[index, 'level2_diag2'] = 6
elif (row['level2_diag2'] >= 500 and row['level2_diag2'] < 520):
    df.loc[index, 'level2_diag2'] = 7
elif (np.floor(row['level2_diag2']) == 786):
    df.loc[index, 'level2_diag2'] = 8
elif (row['level2_diag2'] >= 520 and row['level2_diag2'] < 530):
    df.loc[index, 'level2_diag2'] = 9
elif (row['level2_diag2'] >= 530 and row['level2_diag2'] < 544):
    df.loc[index, 'level2_diag2'] = 10
elif (row['level2_diag2'] >= 550 and row['level2_diag2'] < 554):
    df.loc[index, 'level2_diag2'] = 11
elif (row['level2_diag2'] >= 555 and row['level2_diag2'] < 580):
    df.loc[index, 'level2_diag2'] = 12
elif (np.floor(row['level2_diag2']) == 787):
    df.loc[index, 'level2_diag2'] = 13
elif (np.floor(row['level2_diag2']) == 250):
    df.loc[index, 'level2_diag2'] = 14
elif (row['level2_diag2'] >= 800 and row['level2_diag2'] < 1000):
    df.loc[index, 'level2_diag2'] = 15
elif (row['level2_diag2'] >= 710 and row['level2_diag2'] < 740):
    df.loc[index, 'level2_diag2'] = 16
elif (row['level2_diag2'] >= 580 and row['level2_diag2'] < 630):
    df.loc[index, 'level2_diag2'] = 17
elif (np.floor(row['level2_diag2']) == 788):
    df.loc[index, 'level2_diag2'] = 18
elif (row['level2_diag2'] >= 140 and row['level2_diag2'] < 240):
```

```

        df.loc[index, 'level2_diag2'] = 19
    elif row['level2_diag2'] >= 240 and row['level2_diag2'] < 280 and (np.floor(row
['level2_diag2']) != 250):
        df.loc[index, 'level2_diag2'] = 20
    elif (row['level2_diag2'] >= 680 and row['level2_diag2'] < 710) or (np.floor(ro
w['level2_diag2']) == 782):
        df.loc[index, 'level2_diag2'] = 21
    elif (row['level2_diag2'] >= 290 and row['level2_diag2'] < 320):
        df.loc[index, 'level2_diag2'] = 22
    else:
        df.loc[index, 'level2_diag2'] = 0

if (row['level2_diag3'] >= 390 and row['level2_diag3'] < 399):
    df.loc[index, 'level2_diag3'] = 1
elif (row['level2_diag3'] >= 401 and row['level2_diag3'] < 415):
    df.loc[index, 'level2_diag3'] = 2
elif (row['level2_diag3'] >= 415 and row['level2_diag3'] < 460):
    df.loc[index, 'level2_diag3'] = 3
elif (np.floor(row['level2_diag3']) == 785):
    df.loc[index, 'level2_diag3'] = 4
elif (row['level2_diag3'] >= 460 and row['level2_diag3'] < 489):
    df.loc[index, 'level2_diag3'] = 5
elif (row['level2_diag3'] >= 490 and row['level2_diag3'] < 497):
    df.loc[index, 'level2_diag3'] = 6
elif (row['level2_diag3'] >= 500 and row['level2_diag3'] < 520):
    df.loc[index, 'level2_diag3'] = 7
elif (np.floor(row['level2_diag3']) == 786):
    df.loc[index, 'level2_diag3'] = 8
elif (row['level2_diag3'] >= 520 and row['level2_diag3'] < 530):
    df.loc[index, 'level2_diag3'] = 9
elif (row['level2_diag3'] >= 530 and row['level2_diag3'] < 544):
    df.loc[index, 'level2_diag3'] = 10
elif (row['level2_diag3'] >= 550 and row['level2_diag3'] < 554):
    df.loc[index, 'level2_diag3'] = 11
elif (row['level2_diag3'] >= 555 and row['level2_diag3'] < 580):
    df.loc[index, 'level2_diag3'] = 12
elif (np.floor(row['level2_diag3']) == 787):
    df.loc[index, 'level2_diag3'] = 13
elif (np.floor(row['level2_diag3']) == 250):
    df.loc[index, 'level2_diag3'] = 14
elif (row['level2_diag3'] >= 800 and row['level2_diag3'] < 1000):

```

```

df.loc[index, 'level2_diag3'] = 15
elif (row['level2_diag3'] >= 710 and row['level2_diag3'] < 740):
    df.loc[index, 'level2_diag3'] = 16
elif (row['level2_diag3'] >= 580 and row['level2_diag3'] < 630):
    df.loc[index, 'level2_diag3'] = 17
elif (np.floor(row['level2_diag3']) == 788):
    df.loc[index, 'level2_diag3'] = 18
elif (row['level2_diag3'] >= 140 and row['level2_diag3'] < 240):
    df.loc[index, 'level2_diag3'] = 19
elif row['level2_diag3'] >= 240 and row['level2_diag3'] < 280 and (np.floor(row
['level2_diag3']) != 250):
    df.loc[index, 'level2_diag3'] = 20
elif (row['level2_diag3'] >= 680 and row['level2_diag3'] < 710) or (np.floor(ro
w['level2_diag3']) == 782):
    df.loc[index, 'level2_diag3'] = 21
elif (row['level2_diag3'] >= 290 and row['level2_diag3'] < 320):
    df.loc[index, 'level2_diag3'] = 22
else:
    df.loc[index, 'level2_diag3'] = 0

```

Data Visualization¶

Distribution of Readmission¶

Our target variable is imbalance. Number of readmitted patient are quite less as compared to Not readmitted

In [27]:

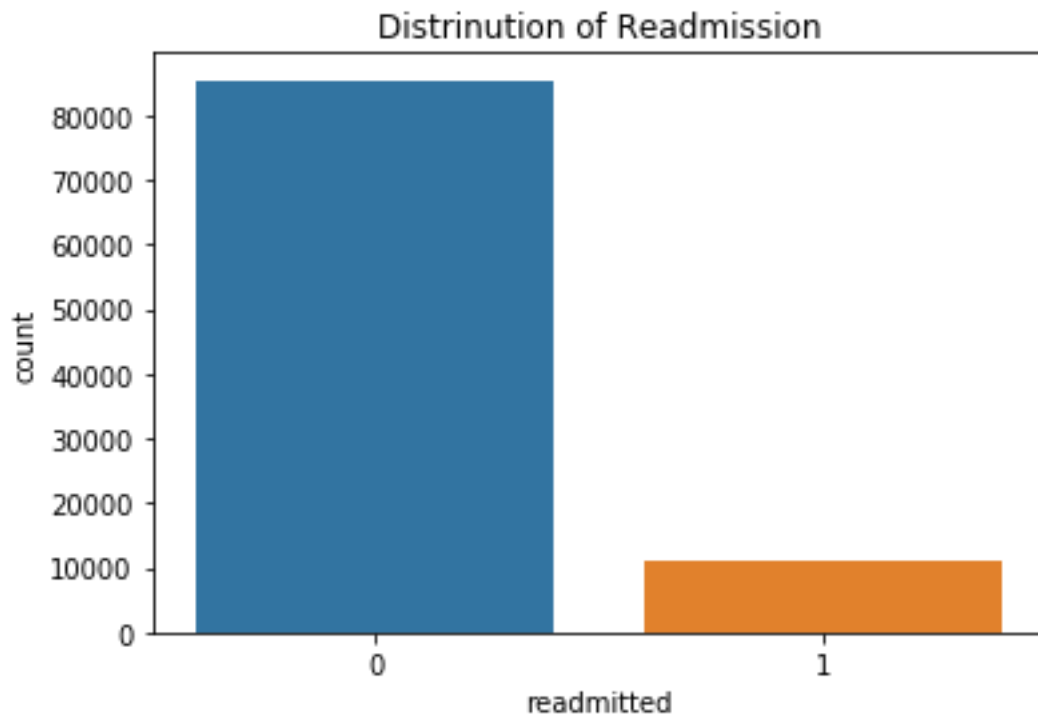
```

# Distribution of Readmission
sns.countplot(df['readmitted']).set_title('Distrinution of Readmission')

```

Out[27]:

```
Text(0.5, 1.0, 'Distrinution of Readmission')
```



Time in Hospital and Readmission ¶

In [28]:

```
fig = plt.figure(figsize=(13,7),)

ax=sns.kdeplot(df.loc[(df['readmitted'] == 0),'time_in_hospital'] , color='b',shade
=True,label='Not Readmitted')

ax=sns.kdeplot(df.loc[(df['readmitted'] == 1),'time_in_hospital'] , color='r',shade
=True, label='Readmitted')

ax.set(xlabel='Time in Hospital', ylabel='Frequency')

plt.title('Time in Hospital VS. Readmission')
```

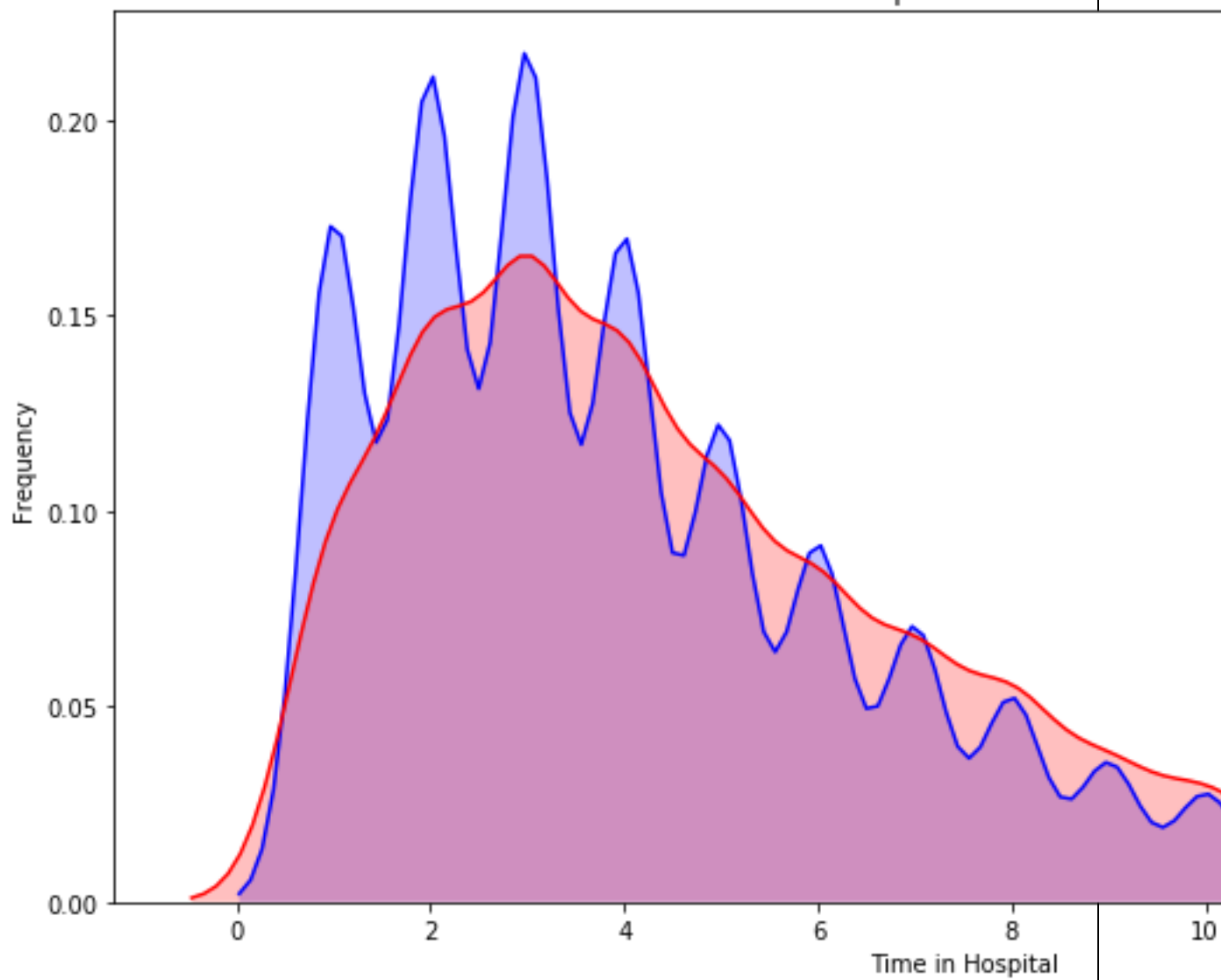
/opt/conda/lib/python3.6/site-packages/scipy/stats/stats.py:1713: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

```
return np.add.reduce(sorted[indexer] * weights, axis=axis) / sumval
```

Out[28]:

```
Text(0.5, 1.0, 'Time in Hospital VS. Readmission')
```

Time in Hospital VS. Readmission



Age and Readmission

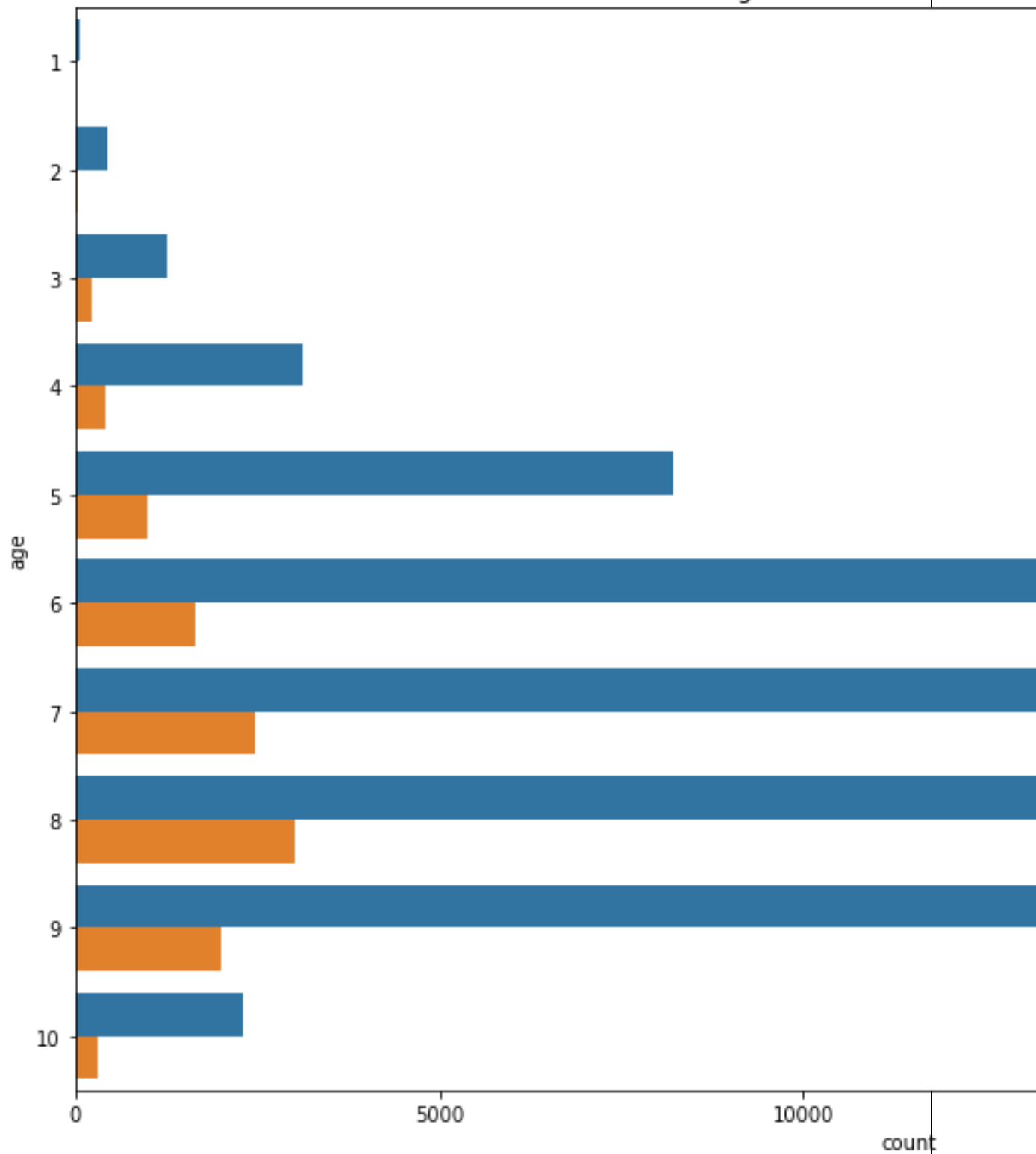
In [29]:

```
fig = plt.figure(figsize=(15,10))
sns.countplot(y= df['age'], hue = df['readmitted']).set_title('Age of Patient VS. R
eadmission')
```

Out[29]:

Text(0.5, 1.0, 'Age of Patient VS. Readmission')

Age of Patient VS. Readmission



Ethnicity of patient and Readmission

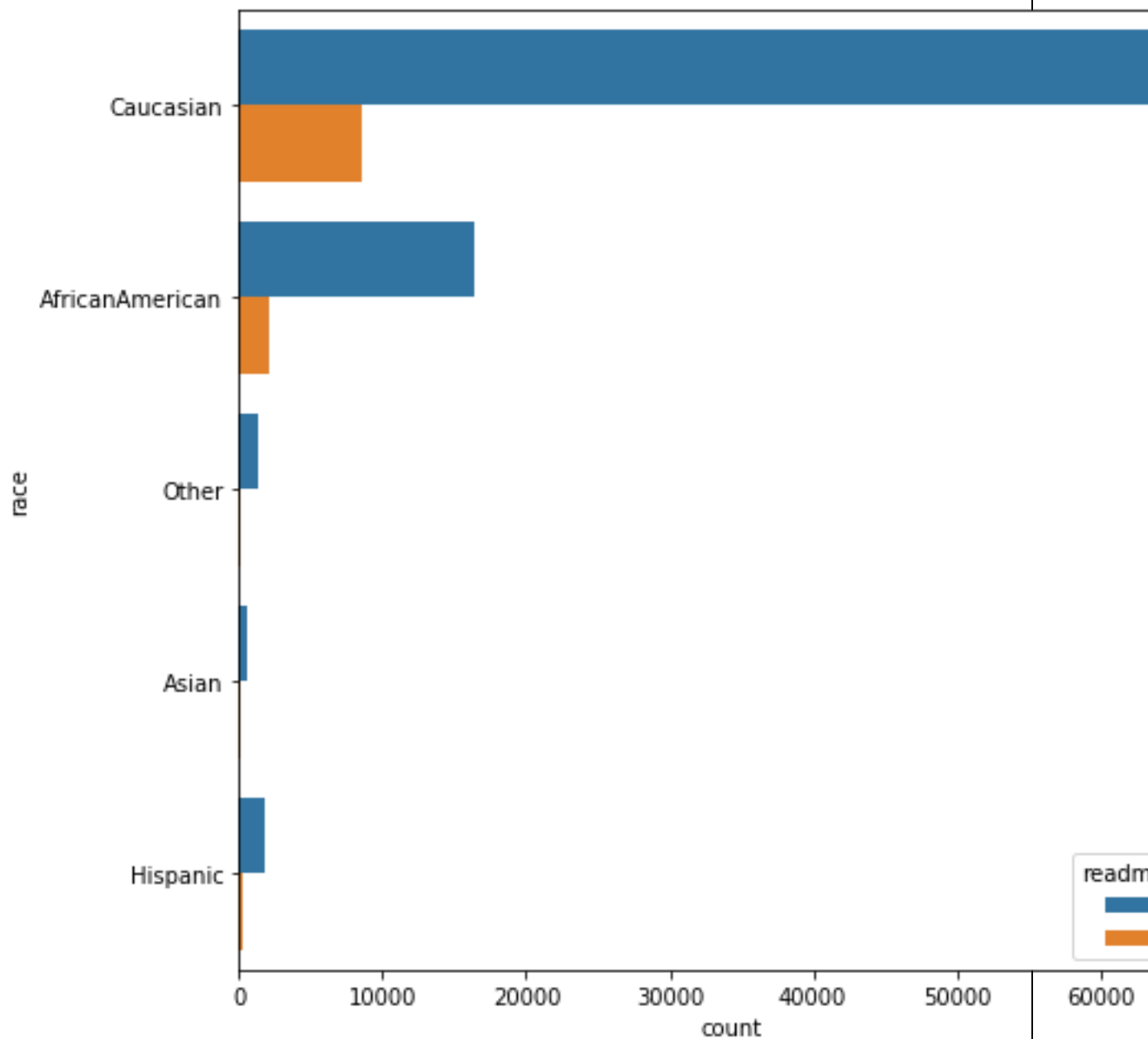
In [30]:

```
fig = plt.figure(figsize=(8,8))
```

```
sns.countplot(y = df['race'], hue = df['readmitted'])
```

Out[30]:

<matplotlib.axes._subplots.AxesSubplot at 0x7f0fba49d1d0>



Number of medication used and Readmission¶

In [31]:

```
fig = plt.figure(figsize=(8,8))
sns.barplot(x = df['readmitted'], y = df['num_medications']).set_title("Number of medication used VS. Readmission")
```

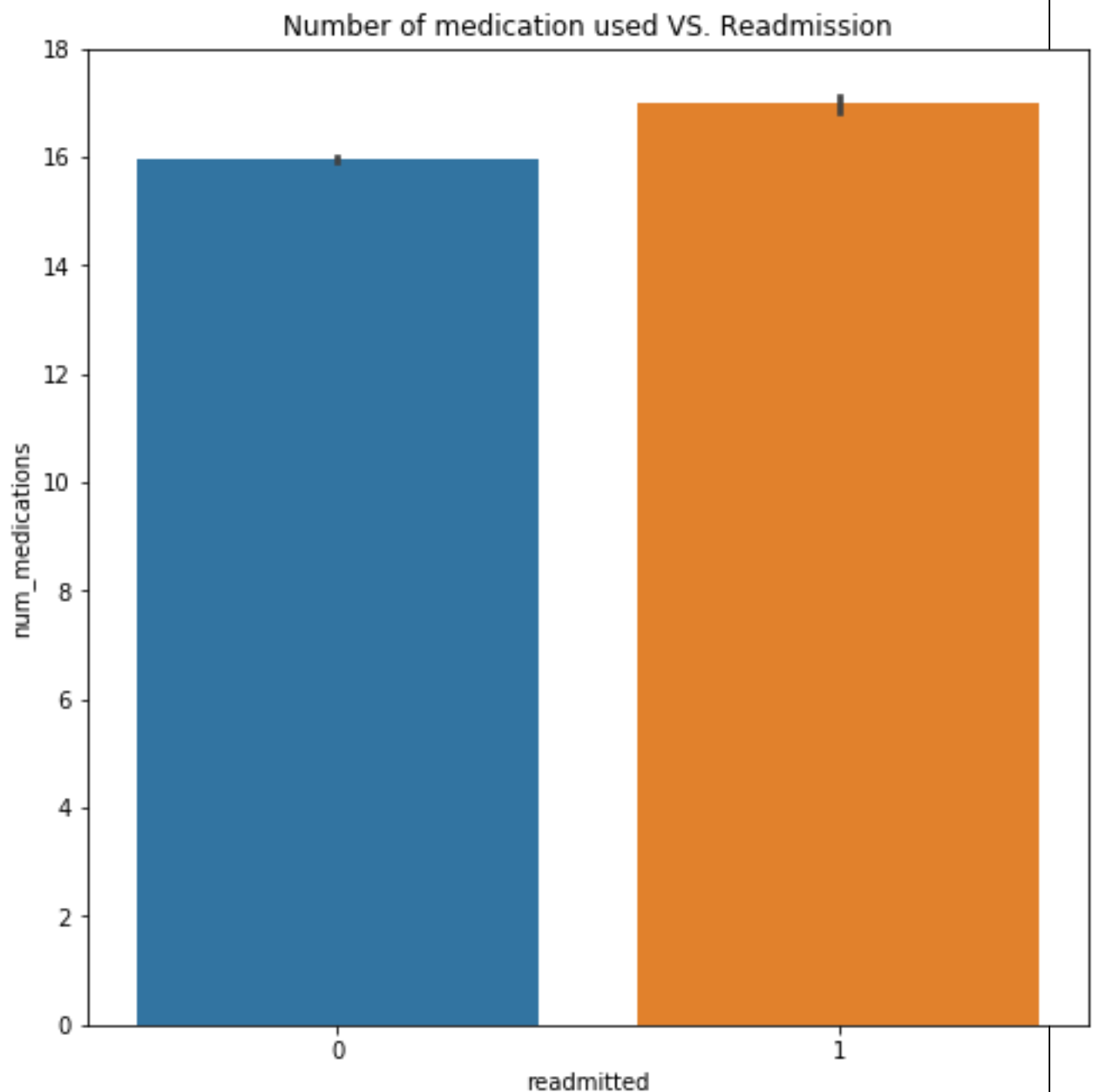
/opt/conda/lib/python3.6/site-packages/scipy/stats/stats.py:1713: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array

```
index, `arr[np.array(seq)]`, which will result either in an error or a different result.
```

```
return np.add.reduce(sorted[indexer] * weights, axis=axis) / sumval
```

Out[31]:

```
Text(0.5, 1.0, 'Number of medication used VS. Readmission')
```



Gender and Readmission

- Male = 1
- Female = 0

In [32]:

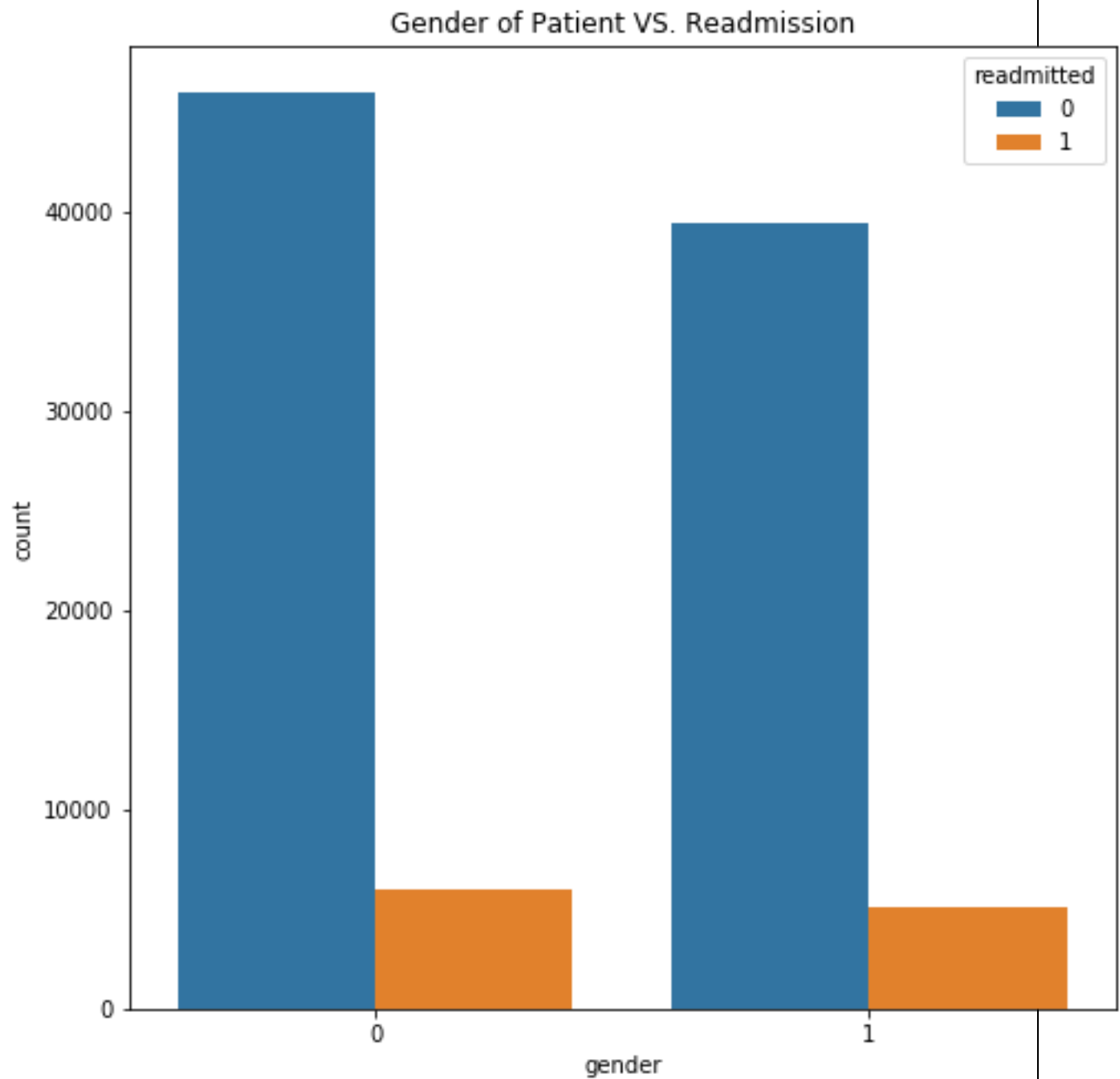
```
fig = plt.figure(figsize=(8,8))
```



```
sns.countplot(df['gender'], hue = df['readmitted']).set_title("Gender of Patient VS  
. Readmission")
```

Out[32]:

```
Text(0.5, 1.0, 'Gender of Patient VS. Readmission')
```



Change of Medication and Readmission ¶

- Change = 1
- No Change = 0

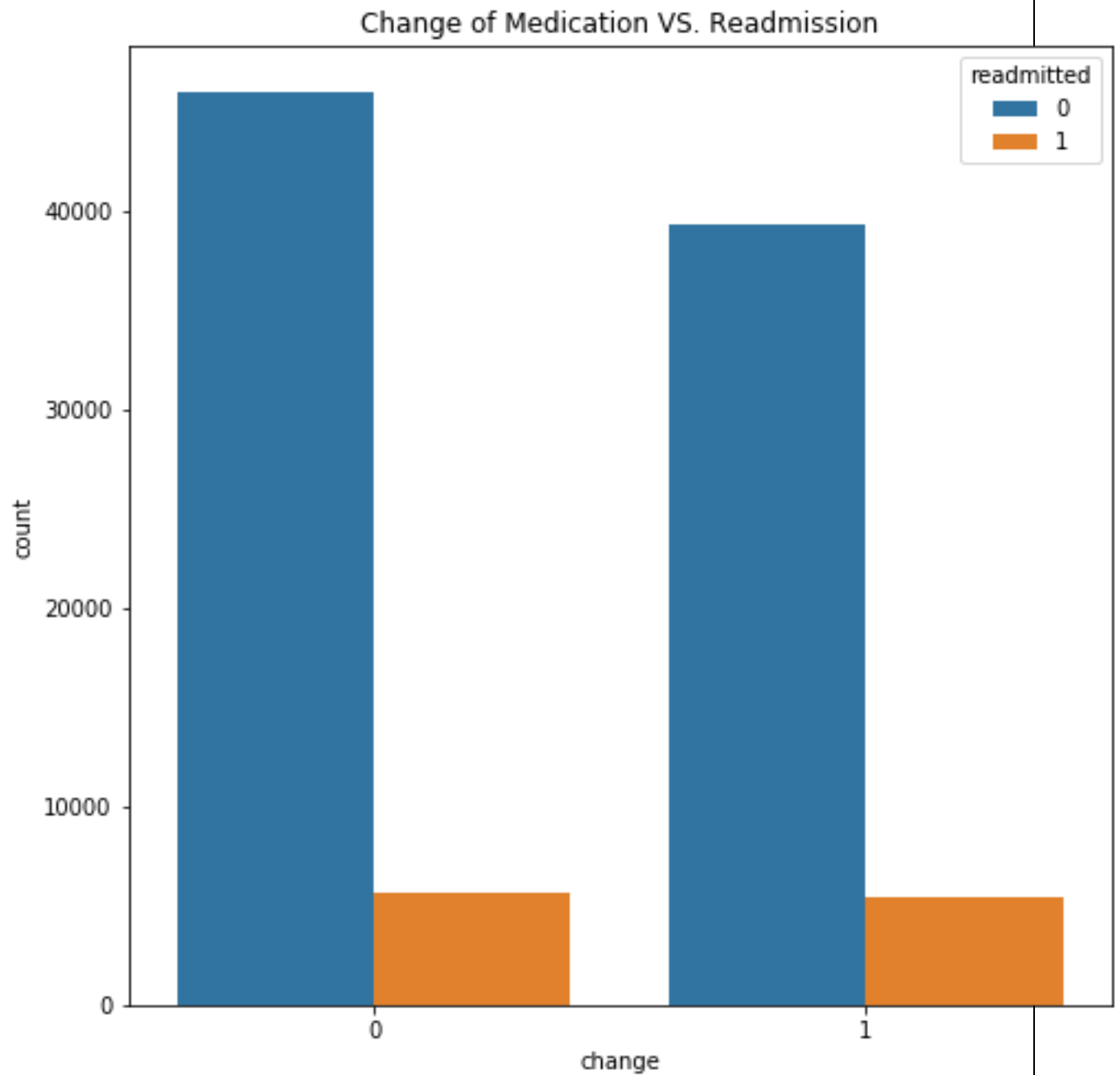
In [33]:

```
fig = plt.figure(figsize=(8,8))
```

```
sns.countplot(df['change'], hue = df['readmitted']).set_title('Change of Medication  
VS. Readmission')
```

Out[33]:

```
Text(0.5, 1.0, 'Change of Medication VS. Readmission')
```



Diabetes Medication prescribed and Readmission¶

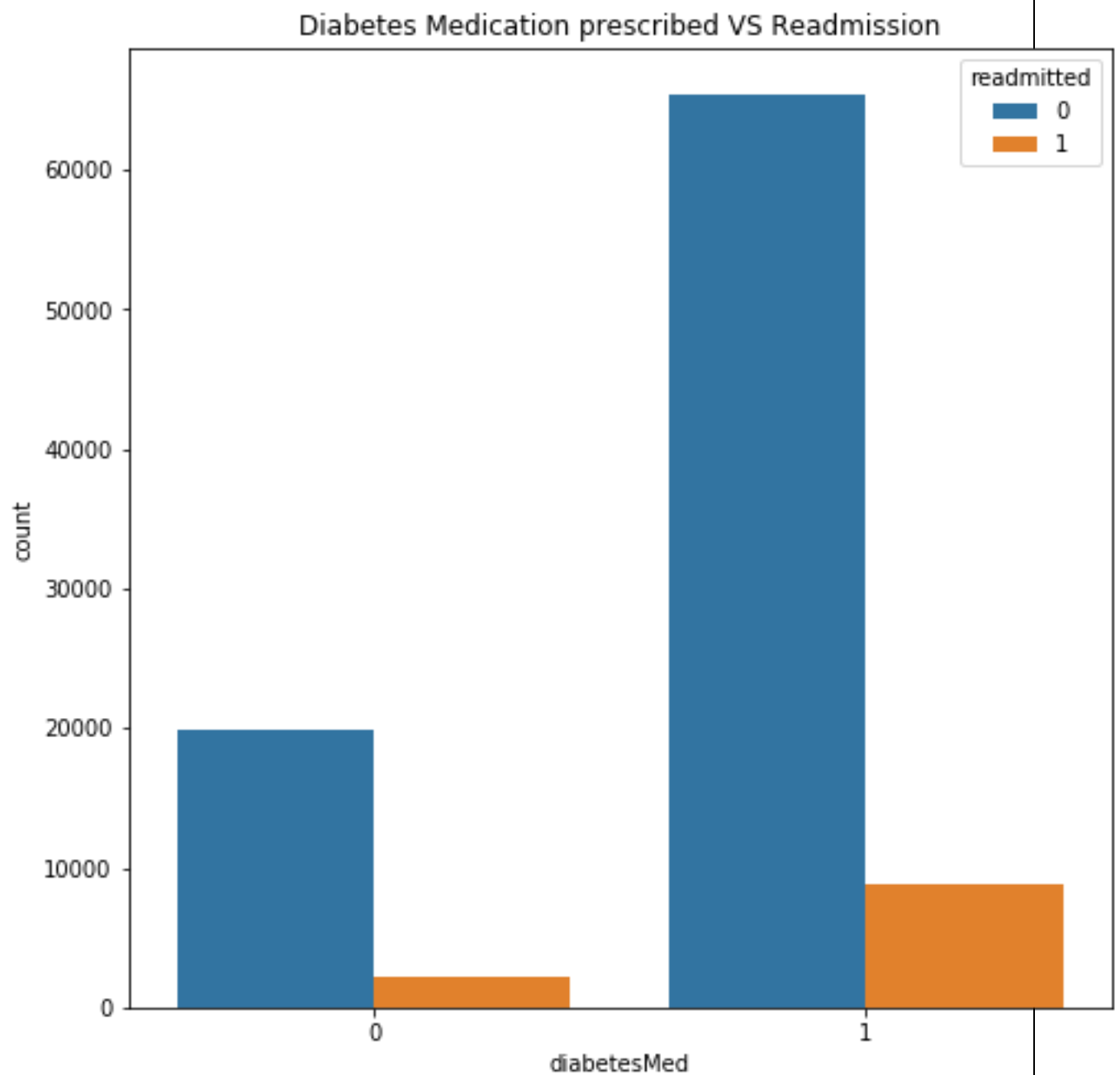
- Diabetes Medication - medications Nominal Indicates if there was any diabetic medication prescribed.
- Values: “yes” : 1 “no” : 0

In [34]:

```
fig = plt.figure(figsize=(8,8))
sns.countplot(df['diabetesMed'], hue = df['readmitted']).set_title('Diabetes Medication prescribed VS Readmission')
```

Out[34]:

Text(0.5, 1.0, 'Diabetes Medication prescribed VS Readmission')



Service Utilization and Readmission ¶

In [35]:

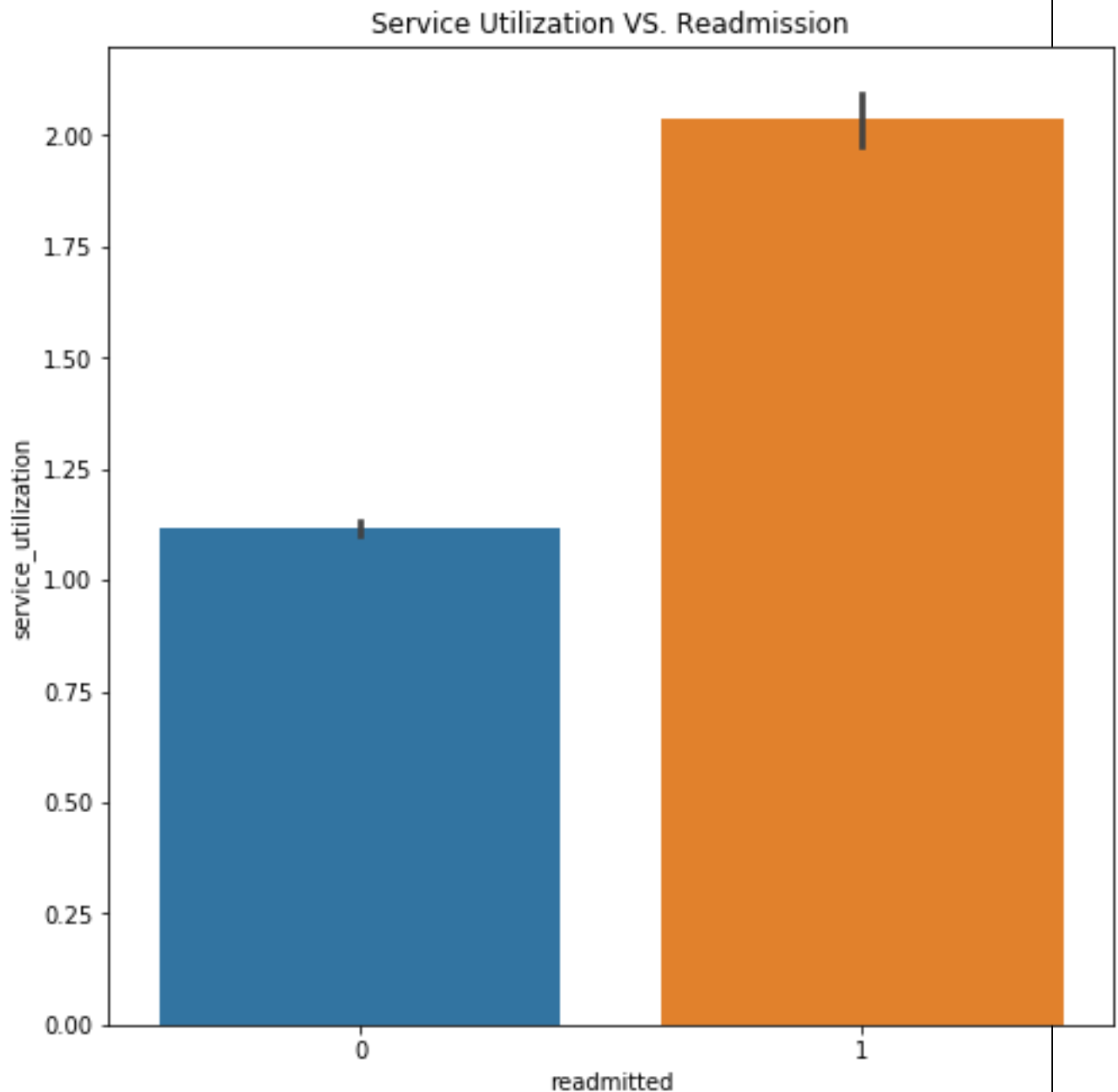
```
fig = plt.figure(figsize=(8,8))
sns.barplot(y = df['service_utilization'], x = df['readmitted']).set_title('Service Utilization VS. Readmission')
```

```
/opt/conda/lib/python3.6/site-packages/scipy/stats/stats.py:1713: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.
```

```
return np.add.reduce(sorted[indexer] * weights, axis=axis) / sumval
```

Out[35]:

```
Text(0.5, 1.0, 'Service Utilization VS. Readmission')
```



Glucose serum test result and Readmission¶

Glucose Serum test - A blood glucose test is used to find out if your blood sugar levels are in the healthy range. It is often used to help diagnose and monitor diabetes.

- '>200' : 1 = indicates diabetes
- '>300' : 1 = Indicates diabetes

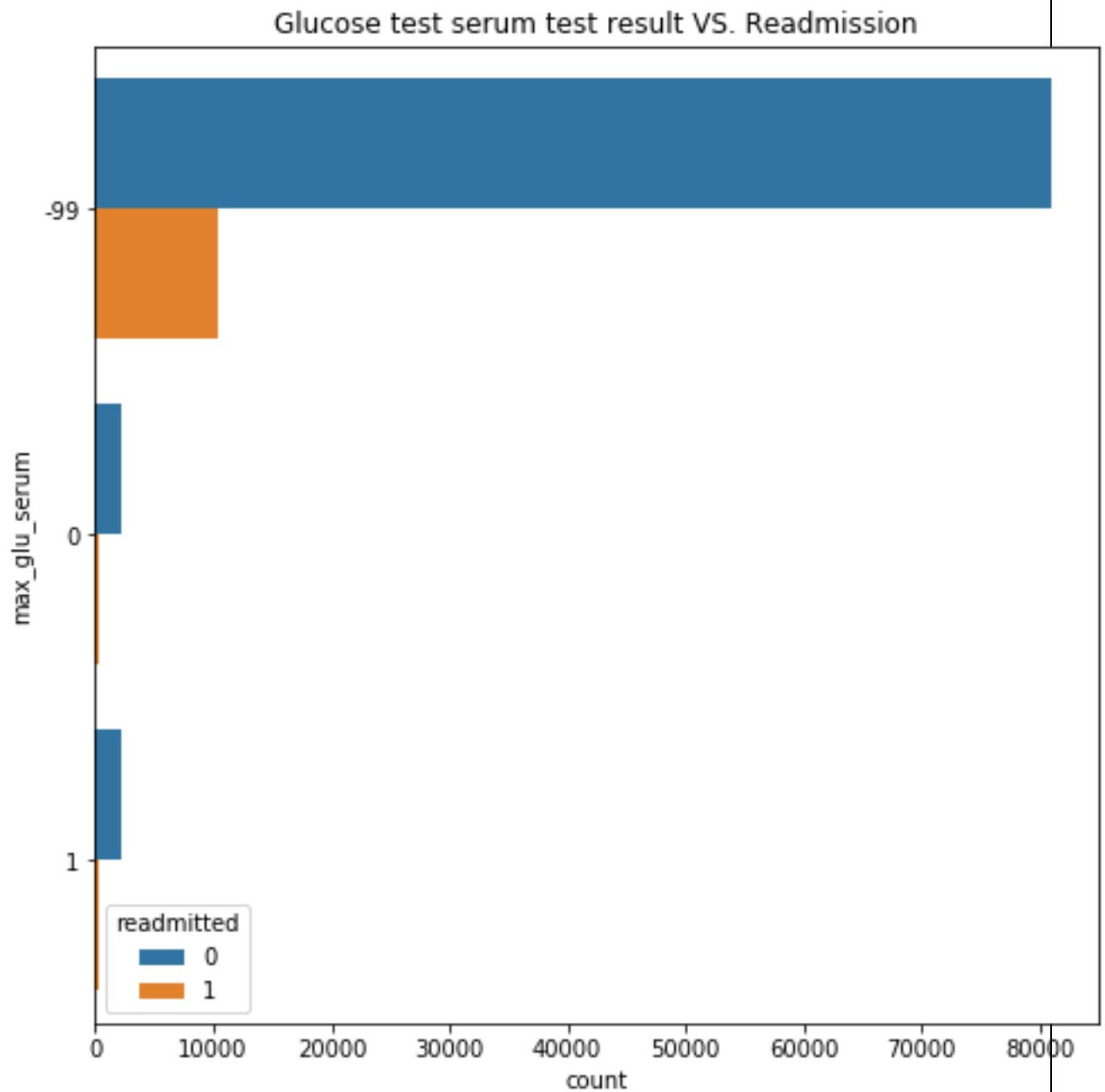
- 'Norm' : 0 = Normal
- 'None' : -99 = test was not taken

In [36]:

```
fig = plt.figure(figsize=(8,8))  
  
sns.countplot(y = df['max_glu_serum'], hue = df['readmitted']).set_title('Glucose test serum test result VS. Readmission')
```

Out[36]:

Text(0.5, 1.0, 'Glucose test serum test result VS. Readmission')



A1C result and Readmission A1C test - The A1C test is a blood test that provides information about your average levels of blood glucose, also called blood sugar, over the past 3 months.

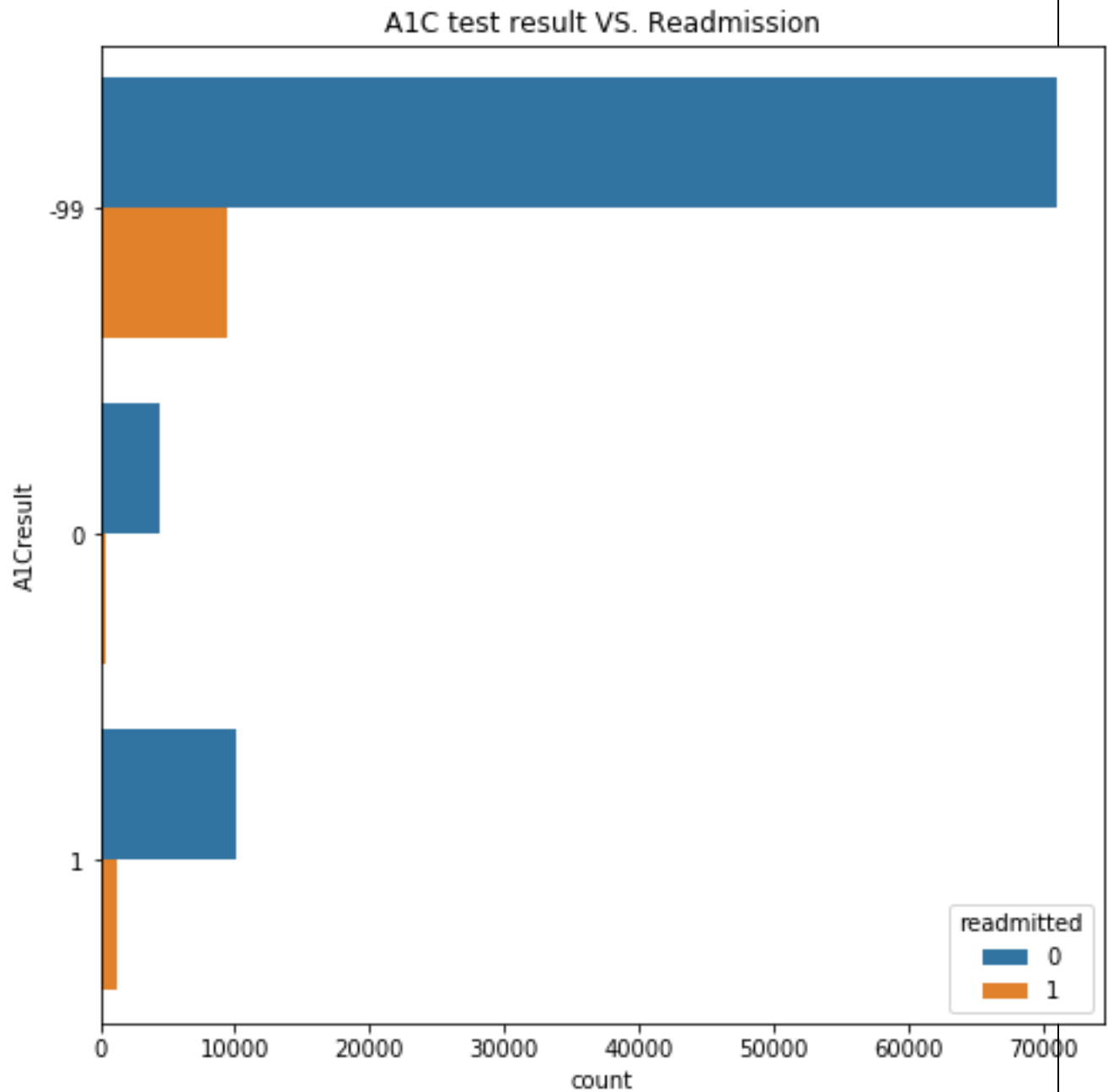
- '>7': 1
- '>8': 1
- Norm : 0 = Normal
- None : -99 = Test was not taken

In [37]:

```
fig = plt.figure(figsize=(8,8))  
sns.countplot(y= df['A1Cresult'], hue = df['readmitted']).set_title('A1C test result  
t VS. Readmission')
```

Out[37]:

Text(0.5, 1.0, 'A1C test result VS. Readmission')



Number of lab procedure and Readmission

In [38]:

```
fig = plt.figure(figsize=(15,6),)

ax=sns.kdeplot(df.loc[(df['readmitted'] == 0),'num_lab_procedures'] , color='b',shade=True,label='Not readmitted')

ax=sns.kdeplot(df.loc[(df['readmitted'] == 1),'num_lab_procedures'] , color='r',shade=True, label='readmitted')

ax.set(xlabel='Number of lab procedure', ylabel='Frequency')

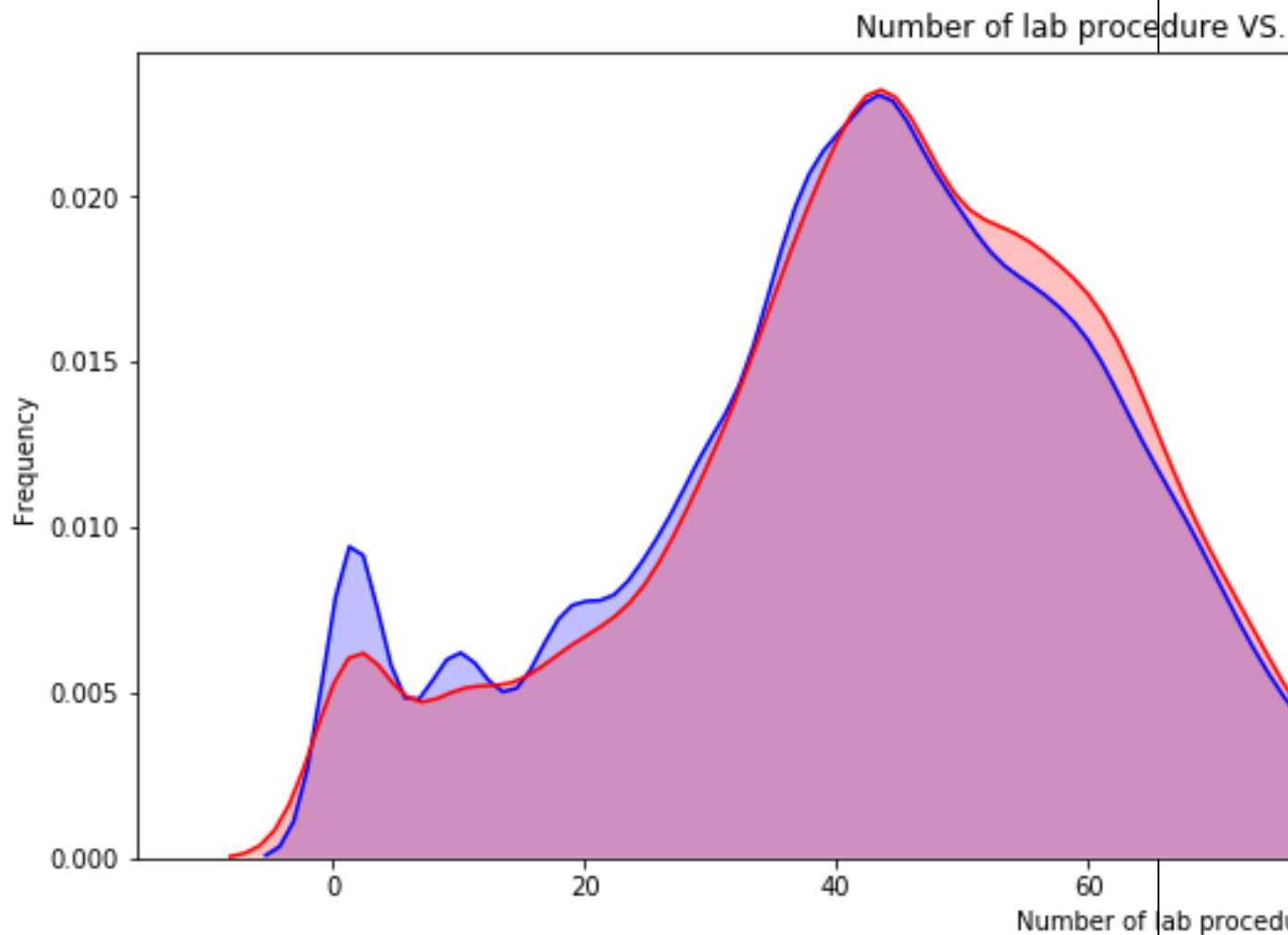
plt.title('Number of lab procedure VS. Readmission')
```

/opt/conda/lib/python3.6/site-packages/scipy/stats/stats.py:1713: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

```
return np.add.reduce(sorted[indexer] * weights, axis=axis) / sumval
```

Out[38]:

Text(0.5, 1.0, 'Number of lab procedure VS. Readmission')



Pre-Modeling Data Preprocessing

This code converts age as categorical variable to a continuous approximation by assuming mid-point of each age-category as the actual age value. This is done to avoid having to deal with age as a dummy variable in the models which makes interpretation very cumbersome. Also, since age category is not purely nominal but ordinal, we do not want to lose that information by treating it as a simple categorical variable

In [39]:

```
df['age'] = df['age'].astype('int64')
print(df.age.value_counts())

# convert age categories to mid-point values
age_dict = {1:5, 2:15, 3:25, 4:35, 5:45, 6:55, 7:65, 8:75, 9:85, 10:95}
df['age'] = df.age.map(age_dict)
print(df.age.value_counts())
```

```
8      24815
7      21521
6      16546
9      16223
5       9208
4       3538
10      2594
3       1471
2        466
1         64
Name: age, dtype: int64
75      24815
65      21521
55      16546
85      16223
45       9208
35       3538
95       2594
25       1471
15        466
5         64
Name: age, dtype: int64
```

In [40]:

```
# convert data type of nominal features in dataframe to 'object' type
i = ['encounter_id', 'patient_nbr', 'gender', 'admission_type_id', 'discharge_dispo-
sition_id', 'admission_source_id', \
      'A1Cresult', 'metformin', 'repaglinide', 'nateglinide', 'chlorpropamide',
'glimepiride', 'acetoexamide', \
      'glipizide', 'glyburide', 'tolbutamide', 'pioglitazone', 'rosiglitazone',
'acarbose', 'miglitol', \
      'troglitazone', 'tolazamide', 'insulin', 'glyburide-metformin', 'glipizid-
e-metformin', \
      'glimepiride-pioglitazone', 'metformin-rosiglitazone', 'metformin-pioglit-
azone', 'change', 'diabetesMed', \
      'age', 'A1Cresult', 'max_glu_serum', 'level1_diag1', 'level1_diag2', 'lev-
el1_diag3', 'level2_diag1', 'level2_diag2', 'level2_diag3']

df[i] = df[i].astype('object')
```


In [41]:

```
df.dtypes
```

Out[41]:

```
encounter_id          object
patient_nbr           object
race                  object
gender                object
age                   object
admission_type_id     object
discharge_disposition_id object
admission_source_id   object
time_in_hospital      int64
num_lab_procedures    int64
num_procedures         int64
num_medications        int64
number_outpatient      int64
number_emergency       int64
number_inpatient       int64
diag_1                 object
diag_2                 object
diag_3                 object
number_diagnoses       int64
max_glu_serum          object
A1Cresult              object
metformin              object
repaglinide            object
nateglinide            object
chlorpropamide         object
glimepiride            object
acetohexamide          object
glipizide              object
glyburide              object
tolbutamide            object
pioglitazone           object
rosiglitazone          object
acarbose               object
miglitol               object
troglitazone           object
tolazamide             object
insulin                object
glyburide-metformin    object
glipizide-metformin    object
glimepiride-pioglitazone object
metformin-rosiglitazone object
metformin-pioglitazone object
change                 object
diabetesMed            object
readmitted             int64
service_utilization    int64
numchange              int64
level1_diag1           object
level2_diag1           object
level1_diag2           object
level2_diag2           object
level1_diag3           object
level2_diag3           object
dtype: object
```

- 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 (keys variable in code below is continued from above):

In [42]:

```
df['nummed'] = 0

for col in keys:
    df['nummed'] = df['nummed'] + df[col]
df['nummed'].value_counts()
```

Out[42]:

```
1    44589
0    22156
2    20901
3     7448
4     1290
5        57
6         5
Name: nummed, dtype: int64
```

In [43]:

```
# get list of only numeric features
num_col = list(set(list(df._get_numeric_data().columns)) - {'readmitted'})
num_col
```

Out[43]:

```
['service_utilization',
 'number_outpatient',
 'num_procedures',
 'num_medications',
 'number_emergency',
 'time_in_hospital',
 'number_inpatient',
 'num_lab_procedures',
 'number_diagnoses',
 'numchange']
```

In [44]:

```
# Removing skewness and kurtosis using log transformation if it is above a threshold value - 2

statdataframe = pd.DataFrame()
statdataframe['numeric_column'] = num_col

skew_before = []
skew_after = []

kurt_before = []
kurt_after = []

standard_deviation_before = []
```

```

standard_deviation_after = []

log_transform_needed = []

log_type = []

for i in num_col:
    skewval = df[i].skew()
    skew_before.append(skewval)

    kurtval = df[i].kurtosis()
    kurt_before.append(kurtval)

    sdval = df[i].std()
    standard_deviation_before.append(sdval)

    if (abs(skewval) >2) & (abs(kurtval) >2):
        log_transform_needed.append('Yes')

        if len(df[df[i] == 0])/len(df) <=0.02:
            log_type.append('log')
            skewvalnew = np.log(pd.DataFrame(df[train_data[i] > 0])[i]).skew()
            skew_after.append(skewvalnew)

            kurtvalnew = np.log(pd.DataFrame(df[train_data[i] > 0])[i]).kurtosis()
            kurt_after.append(kurtvalnew)

            sdvalnew = np.log(pd.DataFrame(df[train_data[i] > 0])[i]).std()
            standard_deviation_after.append(sdvalnew)

        else:
            log_type.append('log1p')
            skewvalnew = np.log1p(pd.DataFrame(df[df[i] >= 0])[i]).skew()
            skew_after.append(skewvalnew)

            kurtvalnew = np.log1p(pd.DataFrame(df[df[i] >= 0])[i]).kurtosis()
            kurt_after.append(kurtvalnew)

            sdvalnew = np.log1p(pd.DataFrame(df[df[i] >= 0])[i]).std()
            standard_deviation_after.append(sdvalnew)

```

```

else:
    log_type.append('NA')
    log_transform_needed.append('No')

    skew_after.append(skewval)
    kurt_after.append(kurtval)
    standard_deviation_after.append(sdval)

statdataframe['skew_before'] = skew_before
statdataframe['kurtosis_before'] = kurt_before
statdataframe['standard_deviation_before'] = standard_deviation_before
statdataframe['log_transform_needed'] = log_transform_needed
statdataframe['log_type'] = log_type
statdataframe['skew_after'] = skew_after
statdataframe['kurtosis_after'] = kurt_after
statdataframe['standard_deviation_after'] = standard_deviation_after

```

In [45]:

```
statdataframe
```

Out[45]:

	numeric_column	skew_before	kurtosis_before	standard_deviation_before	log_transform_needed	log_type	
0	service_utilization	5.312374	67.194018	2.315789	Yes	log1p	1
1	number_outpatient	8.767489	146.244961	1.280061	Yes	log1p	2
2	num_procedures	1.313236	0.856100	1.703183	No	NA	1
3	num_medications	1.339187	3.549325	8.072516	No	NA	1
4	number_emergency	22.695921	1165.140400	0.948089	Yes	log1p	3
5	time_in_hospital	1.127510	0.839050	2.982330	No	NA	1
6	number_inpatient	3.566269	20.044813	1.269975	Yes	log1p	1
7	num_lab_procedures	-0.240626	-0.253275	19.656782	No	NA	-
8	number_diagnoses	-0.807741	-0.372558	1.836659	No	NA	-
9	numchange	1.426548	1.451898	0.488614	No	NA	1

In [46]:

```

# performing the log transformation for the columns determined to be needing it above.

for i in range(len(statdataframe)):
    if statdataframe['log_transform_needed'][i] == 'Yes':
        colname = str(statdataframe['numeric_column'][i])

        if statdataframe['log_type'][i] == 'log':
            df = df[df[colname] > 0]
            df[colname + "_log"] = np.log(df[colname])

```

```

elif statdataframe['log_type'][i] == 'loglp':
    df = df[df[colname] >= 0]
    df[colname + "_loglp"] = np.loglp(df[colname])

```

In [47]:

```

df = df.drop(['number_outpatient', 'number_inpatient', 'number_emergency', 'service_
utilization'], axis = 1)

```

In [48]:

```
df.shape
```

Out[48]:

```
(96446, 54)
```

In [49]:

```

# get list of only numeric features
numerics = list(set(list(df._get_numeric_data().columns))- {'readmitted'})
numerics

```

Out[49]:

```

['num_procedures',
 'num_medications',
 'time_in_hospital',
 'service_utilization_loglp',
 'number_inpatient_loglp',
 'num_lab_procedures',
 'number_diagnoses',
 'number_emergency_loglp',
 'numchange',
 'number_outpatient_loglp']

```

In [50]:

```

# show list of features that are categorical
df.encounter_id = df.encounter_id.astype('int64')
df.patient_nbr = df.patient_nbr.astype('int64')
df.diabetesMed = df.diabetesMed.astype('int64')
df.change = df.change.astype('int64')

# convert data type of nominal features in dataframe to 'object' type for aggregati
ng
i = ['metformin', 'repaglinide', 'nateglinide', 'chlorpropamide', 'glimepiride', 'a
cetoexamide', \
      'glipizide', 'glyburide', 'tolbutamide', 'pioglitazone', 'rosiglitazone',
'acarbose', 'miglitol', \
      'troglitazone', 'tolazamide', 'insulin', 'glyburide-metformin', 'glipizid
e-metformin', \
      'glimepiride-pioglitazone', 'metformin-rosiglitazone', 'metformin-pioglit
azone', 'AlCresult']
df[i] = df[i].astype('int64')

```

```
df.dtypes
```

Out[50]:

```
encounter_id          int64
patient_nbr           int64
race                  object
gender                object
age                   object
admission_type_id     object
discharge_disposition_id object
admission_source_id   object
time_in_hospital      int64
num_lab_procedures    int64
num_procedures         int64
num_medications        int64
diag_1                object
diag_2                object
diag_3                object
number_diagnoses       int64
max_glu_serum          object
A1Cresult              int64
metformin              int64
repaglinide            int64
nateglinide            int64
chlorpropamide         int64
glimepiride            int64
acetoexamide           int64
glipizide              int64
glyburide              int64
tolbutamide            int64
pioglitazone           int64
rosiglitazone          int64
acarbose               int64
miglitol               int64
troglitazone           int64
tolazamide             int64
insulin                int64
glyburide-metformin    int64
glipizide-metformin    int64
glimepiride-pioglitazone int64
metformin-rosiglitazone int64
metformin-pioglitazone int64
change                 int64
diabetesMed            int64
readmitted             int64
numchange              int64
level1_diag1           object
level2_diag1           object
level1_diag2           object
level2_diag2           object
level1_diag3           object
level2_diag3           object
nummed                 object
service_utilization_loglp float64
number_outpatient_loglp float64
number_emergency_loglp float64
number_inpatient_loglp float64
dtype: object
```

In [51]:

```
dfcopy = df.copy(deep = True)
```

In [52]:

```
df['readmitted'] = df['readmitted'].apply(lambda x: 0 if x == 2 else x)
```

In [53]:

```
# drop individual diagnosis columns that have too granular disease information
# also drop level 2 categorization (which was not comparable with any reference)
# also drop level 1 secondary and tertiary diagnoses
df.drop(['diag_1', 'diag_2', 'diag_3', 'level2_diag1', 'level1_diag2', 'level2_diag2', 'level1_diag3', 'level2_diag3'], axis=1, inplace=True)
```

In [54]:

```
interactionterms = [('num_medications', 'time_in_hospital'),
('num_medications', 'num_procedures'),
('time_in_hospital', 'num_lab_procedures'),
('num_medications', 'num_lab_procedures'),
('num_medications', 'number_diagnoses'),
('age', 'number_diagnoses'),
('change', 'num_medications'),
('number_diagnoses', 'time_in_hospital'),
('num_medications', 'numchange')]
```

In [55]:

```
for inter in interactionterms:
    name = inter[0] + '|' + inter[1]
    df[name] = df[inter[0]] * df[inter[1]]
```

In [56]:

```
df[['num_medications', 'time_in_hospital', 'num_medications|time_in_hospital']].head()
```

Out[56]:

	num_medications	time_in_hospital	num_medications time_in_hospital
1	18	3	54
2	13	2	26
3	16	2	32
4	8	1	8
5	16	3	48

In [57]:

```
# Feature Scaling
datf = pd.DataFrame()
datf['features'] = numerics
datf['std_dev'] = datf['features'].apply(lambda x: df[x].std())
```

```
datf['mean'] = datf['features'].apply(lambda x: df[x].mean())
```

In [58]:

```
# dropping multiple encounters while keeping either first or last encounter of these patients
df2 = df.drop_duplicates(subset= ['patient_nbr'], keep = 'first')
df2.shape
```

Out[58]:

```
(67580, 55)
```

In [59]:

```
# standardize function
def standardize(raw_data):
    return ((raw_data - np.mean(raw_data, axis = 0)) / np.std(raw_data, axis = 0))
```

In [60]:

```
df2[numerics] = standardize(df2[numerics])
import scipy as sp
df2 = df2[(np.abs(sp.stats.zscore(df2[numerics])) < 3).all(axis=1)]
```

/opt/conda/lib/python3.6/site-packages/pandas/core/frame.py:3140: SettingWithCopyWarning: A value is trying to be set on a copy of a slice from a DataFrame. Try using `.loc[row_indexer,col_indexer] = value` instead

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

```
self[k1] = value[k2]
```

In [61]:

```
from matplotlib.colors import ListedColormap
my_cmap = ListedColormap(sns.light_palette((250, 100, 50), input="husl", n_colors=50).as_hex())
table = df2.drop(['patient_nbr', 'encounter_id'], axis=1).corr(method='pearson')
table.style.background_gradient(cmap=my_cmap, axis = 0)
```

/opt/conda/lib/python3.6/site-packages/matplotlib/colors.py:512: RuntimeWarning: invalid value encountered in less

```
xa[xa < 0] = -1
```

Out[61]:

	time_in_hospital	num_lab_dures	num_procedures	num_medications	number_of_diagnoses	A1C_last	met_for_surgery	revascularization	nateglinide	chlorzoxipone	glimepiride	acetaminophen	glipizide	glyburide	tolbutamide	pioglitazone	rosiglitazone
time_in_hospital	1	0.320135	0.144392	0.44167	0.224011	0.0741271	-0.012421	0.0311302	0.00665259	0.00116397	0.00849595	nan	0.0136978	0.0246103	0.0050625	-0.00437611	-0.000031

	time_i n_hospital	num_la b_proce dures	num_ proce dures	num_ medic ations	nume r_diagn oses	A1C resu lt	met for min	repa glini de	nate glini de	chlor propa mide	glim epiri de	aceto hexa mide	glipi zide	glyb uride	tolb uta mide	piog litaz one	ro lita one
num_lab_proce dures	0.3201 35	1	0.0139 234	0.2308 88	0.138470 15953	0.27 05	0.05 4736 05	0.00 4736 61	0.00 7754 78	0.002 3442	0.00 9527 73	nan	0.01 655 25	0.00 2659 38	0.00 0594 595	0.01 7769 6	0.01 0.01 43
num_procedure s	0.14430 92	0.01392 34	1	0.3349 18	0.060380 049517	0.02 732 4	0.05 3595 61	0.00 3595 87	0.00 5481 43	0.005 28647 26	0.00 5021 26	nan	0.00 463 863	0.00 2827 74	5.66 207e -05	0.01 2495 32	0.01 0.01 80
num_medicatio ns	0.44110 67	0.23088 8	0.3349 18	1	0.246870 97807	0.01 7807	0.07 921 32	0.02 5270 2	0.02 09484 9	0.002 09745 6	0.04 3242 6	nan	0.05 536 21	0.04 3031 1	0.00 15079 38	0.07 4775 4	0.01 0.01 4
number_diagno ses	0.22400 11	0.13847 1	0.0603 804	0.2468 79	1	0.00 2964 7	0.07 213 64	0.03 1346 1	0.01 4550 5982	0.015 1961 5	0.01 1961 5	nan	0.01 260 79	0.02 4172 1	1.18 069e -05	0.00 4104 12	0.01 0.01 67
A1Cresult	0.07410 271	0.27595 3	- 0.0295 174	0.0178 07	0.00296 47	1	0.03 775 76	0.01 7017 6	0.00 1161 39	- 0.001 98419 4	0.01 8481 4	nan	0.00 805 466	0.00 3427 09	- 0.00 1884 14	0.00 3550 19	0.01 0.01 66
metformin	- 0.01240 21	- 0.05078 05	- 0.0573 261	0.0792 132	- 0.07213 64	0.03 77571 6	0.00 48075 84	0.01 56990 9	- 0.011 7282 2	0.04 0138 2	nan	0.07 435 03	0.13 013 77	- 0.00 8117 69	0.00 3039 9	0.05 0.05 42	0.01 0.01 9
repaglinide	0.03110 302	0.00473 661	0.0035 9587	0.0252 702	0.03134 6	0.01 7017480 6784	0.00 10980 23	0.00 10980 23	- 0.003 69721 12	0.00 8599 12	nan	0.02 027 74	0.02 3441 2	0.00 1789 2	0.01 503975 2	0.01 0.01 6	
nateglinide	0.0066 5259	- 0.00775 478	- 0.0054 8143	0.0294 849	0.01455 51	0.00 1161569 39	0.01 10981 99	0.00 10981 23	- 0.002 78349 37	0.00 8967 37	nan	0.01 582 24	0.02 2630 5	- 0.00 13473	0.02 873484 8	0.01 0.01 8	
chlorpropamid e	0.0011 6397	- 0.00234 42	0.0052 8647	0.0020 9745	- 0.01559 82	0.00 1984172 19	0.01 36972783 21	0.00 0.00 49	1	0.00 7629 32	nan	0.01 080 74	0.00 64940 59	0.00 519723703 284	0.00 0.00 68	0.01 0.01 39	
glimepiride	0.0084 9595	- 0.00952 773	0.0050 2126	0.0432 426	0.01196 15	0.01 8481013 4	0.04 8599 12	0.00 89670 37	- 0.007 62932	1	nan	0.07 477 24	0.07 27403 7	0.00 3692 01	0.05 073465 8	0.01 0.01 5	
acetoexamide	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan
glipizide	0.01360 978	0.01655 25	0.0046 3863	0.0553 621	- 0.01260 79	0.00 8054435 66	0.07 02775822 4	0.01 8074 4	- 0.010 4772 4	0.07 4772 4	nan	1	0.11 011 2234	0.00 2234598 598	0.03 0.03 62	0.01 0.01 6	
glyburide	0.02460 103	0.00265 938	0.0028 2774	0.0430 311	- 0.02417 21	0.00 024173427 09	0.13 34412630 2	0.02 0.02 5	- 0.006 49459 7	0.07 2740 7	nan	0.11 223 4	1	0.00 2134 31	0.02 126899 3	0.01 0.01 2	
tolbutamide	0.00500 625	0.00059 4595	5.6620 7e-05	0.0050 7938	1.180690 e-05	0.00 1884811 14	0.00 0.00 769	0.00 0.00 17	- 0.00 1347 4	0.00 0.00 01	nan	0.00 598	0.00 2134 31	1	0.00 445341 7	0.01 0.01 39	

	time_i	num_la	num_	num_	nume	A1C	met	repa	nate	chlor	glim	aceto	glibi	glyb	tolb	piog	ro
	n_hos	b_proce	proce	medic	r_diagn	resu	for	glini	glini	propa	epiri	hexa	glipi	glyb	tolb	piog	ro
	pital	dures	dures	ations	oses	lt	min	de	de	mide	de	mide	zide	uride	uta	litaz	litaz
																one	one
pioglitazone	-	-	0.0124	0.0747	0.0041	100.00	0.05	0.01	0.02	-	0.05	nan	0.03	0.02	-	-	-
	0.0043	0.0177	0.0124	0.0747	0.0041	100.00	0.05	0.01	0.02	-	0.05	nan	0.03	0.02	-	-	-
	7611	96	95	754	412	3550	19	9	2	3	2376	8	892	1268	4453	1	0.0
rosiglitazone	-	-	0.0080	0.0582	-	0.00	0.09	0.01	0.01	-	0.03	nan	0.03	0.02	-	-	-
	0.0008	0.0088	0.0080	0.0582	0.0067	7748	2542	7575	8490	0.000	0.03	nan	0.03	0.02	-	-	-
	1331	943	1632	544	696	66	39	6	8	4	38039	5	627	9915	4192	0.06	1
acarbose	0.0008	-	-	0.0139	0.0068	208	2179	4238	3626	0.001	0.01	nan	0.02	0.01	-	0.01	0.0
	16241	0.0019	0.0018	0.0139	0.0068	208	2179	4238	3626	0.001	0.01	nan	0.02	0.01	-	0.01	0.0
	16241	971	3317	269	689	783	509	8	92	6494	79	5442	442	7804	0798	5307	09
miglitol	-	-	0.0010	0.0030	-	0.00	0.01	0.02	0.01	-	0.01	nan	2.09	-	-	-	-
	0.0022	0.0036	0.0010	0.0030	0.0012	0.00	0.01	0.02	0.01	-	0.01	nan	2.09	-	-	-	-
	0634	5	088	5175	699	2302	179	569	7073	0.000	0.01	nan	664	0.00	0.00	0.00	0.0
troglitazone	0.0056	0.0051	-	0.0042	0.0052	20.00	0.00	0.00	0.00	0.000	0.00	nan	-	-	-	-	-
	9745	974	0.0057	0.0042	0.0052	20.00	0.00	0.00	0.00	0.000	0.00	nan	-	-	-	-	-
	9745	974	3509	5403	184	3313	362	0800	0602	23220	9044	nan	267	2473	0112	1991	76
tolazamide	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	0.0066	0.0010	0.0002	0.0015	0.0104	1724	118	2400	1807	69676	4953	nan	0.00	0.00	0.00	0.00	0.0
	5848	684	8984	25763	94	78	311	66	37	3	85	nan	802	7422	0337	2950	39
insulin	0.0973	0.1001	0.0024	0.1943	0.0903	7772	0.09	0.00	0.00	-	0.00	nan	-	-	-	-	-
	714	7	5433	39	0.0903	7772	0.09	0.00	0.00	-	0.00	nan	-	-	-	-	-
	714	7	5433	39	0.0903	7772	0.09	0.00	0.00	-	0.00	nan	-	-	-	-	-
glyburide-metformin	-	-	-	0.0087	-	0.00	0.02	0.00	0.00	-	0.00	nan	0.02	0.01	-	0.03	0.0
	0.0003	0.0142	0.0091	0.0087	0.0024	3250	5602	5950	2443	0.002	0.00	nan	0.02	0.01	-	0.03	0.0
	01704	74	1033	4587	533	54	82	07	44	76706	9730	nan	763	0859	1339	0716	91
glipizide-metformin	0.0023	-	-	0.0042	0.0017	0.00	0.00	0.00	0.00	0.000	0.00	nan	0.00	0.00	-	0.00	-
	3415	0.0037	0.0033	0.0042	0.0017	0.00	0.00	0.00	0.00	0.000	0.00	nan	0.00	0.00	-	0.00	-
	3415	084	7336	6288	275	3027	173	1222	0920	35471	2521	nan	408	1194	0171	2898	28
glimepiride-pioglitazone	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan
metformin-pioglitazone	0.0027	-	-	0.0028	-	0.00	0.00	0.00	0.00	0.000	0.00	nan	-	-	-	0.01	-
	9377	0.0039	0.0009	0.0028	0.0070	1913	799	0461	0347	13406	0953	nan	154	1428	763e	4566	10
	9377	099	35588	3758	432	17	196	907	752	3	162	nan	385	1	-05	8	34
change	0.0944	0.0645	-	0.2332	0.0423	0.09	0.35	0.07	0.05	-	0.14	nan	0.20	0.19	-	0.21	0.2
	059	98	0.0161	0.2332	0.0423	0.09	0.35	0.07	0.05	-	0.14	nan	0.20	0.19	-	0.21	0.2
	059	98	973	56	37	1	2	6	9	87703	7047	nan	8	9645	3177	4111	35
diabetesMed	0.0598	0.0374	-	0.1900	0.0225	0.07	0.29	0.06	0.04	0.019	0.13	nan	0.21	0.20	-	0.16	0.1
	106	58	0.0173	0.1900	0.0225	0.07	0.29	0.06	0.04	0.019	0.13	nan	0.21	0.20	-	0.16	0.1
	106	58	448	55	25	9	7	6	2	0128	5177	nan	894	2533	9200	3065	34
readmitted	0.0551	0.0256	0.0020	0.0387	0.0396	-	0.01	-	-	-	-	nan	0.00	0.00	-	-	-
	478	38	8551	684	02	0.00	0.01	5210	0.00	0.004	0.00	nan	859	7831	-	-	-
	478	38	8551	684	02	0.00	0.01	5	0.00	47192	0.00	nan	946	6	-	-	-

	time_in_hospital	num_lab_procedures	num_procedures	num_medications	number_diagnoses	A1C_last_visit	metformin	repaglinide	nateglinide	chlorpropamide	glimepiride	acetohexamide	glipizide	glyburide	tolbutamide	pioglitazone	rosiglitazone
numchange	0.13967	0.111876	-0.0109966	0.193801	0.065893	20.111597	0.0196183	0.0316958	0.0043850	-0.0213918	0.0090934	nan	0.00320806	-0.000152969	-0.00604123	0.0237586	0.000006
service_utilization_log1p	0.00686184	-0.00178267	-0.0305084	0.0450339	0.086091	30.0507292	0.00837825	0.0096226	0.0120621	-0.00329107	0.00465963	nan	0.00432188	0.00284534	-0.00335433	0.0183022	0.000003
number_outpatient_log1p	-0.0276142	-0.0412072	-0.0172952	0.0315759	0.0579024	-0.0326165	0.0263800	0.0044324	0.0124780	-0.00466847	0.000338977	nan	0.00511965	0.0127531	-0.00207377	0.0196698	0.000017
number_emergency_log1p	-0.00182889	0.00678354	-0.026767	0.0133994	0.050317	-0.011517	0.00202658	0.00906256	0.002664056	-0.00772182	0.00484999	nan	0.00214713	-0.0137833	-0.00373678	0.0172993	0.000005
number_inpatient_log1p	0.05140963	0.0491532	-0.0126986	0.0348021	0.0515976	0.0439998	0.018785	0.0107275	0.000004	0.00410726	0.00555767	nan	0.00100386	0.00170352	0.00105754	0.00290046	0.000019
num_medications time_in_hospital	0.859364	0.328828	0.2637	0.752726	0.244135	0.0627198	0.0180484	0.0314325	0.015039	0.0075431	0.020308	nan	0.027586	0.0325266	0.002701	0.024894	0.000005
num_medications num_procedures	0.270542	0.114754	0.877783	0.599321	0.111045	0.00169956	-0.0122985	0.00860907	0.00128817	0.00956255	0.0132681	nan	0.0176849	0.0178769	0.0040167	0.0276992	0.000008
time_in_hospital num_lab_procedures	0.853015	0.674595	0.138598	0.439898	0.219112	0.179762	0.0356974	0.0250283	-0.00112905	0.00088621	0.0017135	nan	0.016832	0.0161035	0.00469782	-0.0124606	0.000056
num_medications num_lab_procedures	0.502133	0.732844	0.24444	0.772237	0.25072	0.178755	0.00803661	0.0215162	0.0126760	-0.0177266	0.021408	nan	0.0414828	0.0235942	0.00260141	0.0312066	0.000006
num_medications number_diagnoses	0.456259	0.258944	0.292403	0.898603	0.599241	0.0282219	0.0219442	0.0351259	0.0230200	-0.0214166	0.038733	nan	0.0365288	0.0217995	0.0051018	0.0586323	0.000003
change num_medications	0.231169	0.135242	0.109633	0.530894	0.123807	0.0782336	0.300672	0.0693626	0.0539000	-0.0511084	0.130815	nan	0.181777	0.1605521	0.000765066	0.197647	0.000081
number_diagnoses time_in_hospital	0.938944	0.322335	0.153126	0.459725	0.484368	0.0615688	-0.0341624	0.0397095	0.0112960	-0.005613	0.00994339	nan	0.00529891	0.0130817	0.0038351	-0.00292311	0.000081
num_medications numchange	0.229845	0.155882	0.0750897	0.392592	0.120178	0.0995713	0.0216115	0.0318127	0.0048440	-0.0004966	0.0260481	nan	0.0105357	0.00170784	-0.00486772	0.0362543	0.000046

In [62]:

```
df2['level1_diag1'] = df2['level1_diag1'].astype('object')

df_pd = pd.get_dummies(df2, columns=['gender', 'admission_type_id', 'discharge_disposition_id',
                                     'admission_source_id', 'max_glu_serum', 'A1Cresult', 'level1_diag1'], drop_first = True)

just_dummies = pd.get_dummies(df_pd['race'])

df_pd = pd.concat([df_pd, just_dummies], axis=1)

df_pd.drop(['race'], inplace=True, axis=1)
```

In [63]:

```
non_num_cols = ['race', 'gender', 'admission_type_id', 'discharge_disposition_id',
                'admission_source_id',
                'max_glu_serum', 'A1Cresult', 'level1_diag1' ]
```

In [64]:

```
num_cols = list(set(list(df._get_numeric_data().columns))- {'readmitted', 'change'})

num_cols
```

Out[64]:

```
['num_medications|time_in_hospital',
 'num_medications|num_procedures',
 'time_in_hospital',
 'num_medications|numchange',
 'metformin',
 'encounter_id',
 'insulin',
 'service_utilization_loglp',
 'repaglinide',
 'acetoexamide',
 'pioglitazone',
 'time_in_hospital|num_lab_procedures',
 'numchange',
 'nateglinide',
 'glipizide-metformin',
 'num_procedures',
 'num_medications',
 'A1Cresult',
 'number_inpatient_loglp',
 'num_lab_procedures',
 'number_diagnoses',
 'glyburide-metformin',
 'rosiglitazone',
 'number_diagnoses|time_in_hospital',
 'glimepiride-pioglitazone',
 'num_medications|number_diagnoses',
 'num_medications|num_lab_procedures',
 'number_emergency_loglp',
 'patient_nbr',
 'acarbose',
 'miglitol',
 'metformin-pioglitazone',
 'tolbutamide',
 'glipizide',
 'glyburide',
```

```
'number_outpatient_loglp',
'glimepiride',
'metformin-rosiglitazone',
'chlorpropamide',
'troglitazone',
'diabetesMed',
'tolazamide',
'change|num_medications']
```

In [65]:

```
new_non_num_cols = []
for i in non_num_cols:
    for j in df_pd.columns:
        if i in j:
            new_non_num_cols.append(j)
```

In [66]:

```
new_non_num_cols
```

Out[66]:

```
['gender_1',
'admission_type_id_3',
'admission_type_id_4',
'admission_type_id_5',
'discharge_disposition_id_2',
'discharge_disposition_id_7',
'discharge_disposition_id_10',
'discharge_disposition_id_18',
'discharge_disposition_id_19',
'discharge_disposition_id_20',
'discharge_disposition_id_27',
'discharge_disposition_id_28',
'admission_source_id_4',
'admission_source_id_7',
'admission_source_id_8',
'admission_source_id_9',
'admission_source_id_11',
'max_glu_serum_0',
'max_glu_serum_1',
'A1Cresult_0',
'A1Cresult_1',
'level1_diag1_1.0',
'level1_diag1_2.0',
'level1_diag1_3.0',
'level1_diag1_4.0',
'level1_diag1_5.0',
'level1_diag1_6.0',
'level1_diag1_7.0',
'level1_diag1_8.0']
```

In [67]:

```
l = []
for feature in list(df_pd.columns):
    if '|' in feature:
        l.append(feature)

l
```

Out[67]:

```
['num_medications|time_in_hospital',
 'num_medications|num_procedures',
 'time_in_hospital|num_lab_procedures',
 'num_medications|num_lab_procedures',
 'num_medications|number_diagnoses',
 'age|number_diagnoses',
 'change|num_medications',
 'number_diagnoses|time_in_hospital',
 'num_medications|numchange']
```

In [68]:

```
df_pd.head().T
```

Out[68]:

	1	2	3	4	5
encounter_id	149190	64410	500364	16680	35754
patient_nbr	55629189	86047875	82442376	42519267	82637451
age	15	25	35	45	55
time_in_hospital	-0.444872	-0.784109	-0.784109	-1.12335	-0.444872
num_lab_procedures	0.803605	-1.60206	0.0518349	0.402661	-0.5997
num_procedures	-0.816784	2.02679	-0.24807	-0.816784	2.5955
num_medications	0.270081	-0.335198	0.0279693	-0.940477	0.0279693
number_diagnoses	0.881754	-0.704844	-0.175978	-1.23371	0.881754
metformin	0	0	0	0	0
repaglinide	0	0	0	0	0
nateglinide	0	0	0	0	0
chlorpropamide	0	0	0	0	0
glimepiride	0	0	0	0	0
acetohexamide	0	0	0	0	0
glipizide	0	1	0	1	0
glyburide	0	0	0	0	0
tolbutamide	0	0	0	0	0
pioglitazone	0	0	0	0	0
rosiglitazone	0	0	0	0	0
acarbose	0	0	0	0	0
miglitol	0	0	0	0	0
troglitazone	0	0	0	0	0
tolazamide	0	0	0	0	0
insulin	1	0	1	1	1
glyburide-metformin	0	0	0	0	0
glipizide-metformin	0	0	0	0	0
glimepiride-pioglitazone	0	0	0	0	0
metformin-rosiglitazone	0	0	0	0	0
metformin-pioglitazone	0	0	0	0	0
change	1	0	1	1	0
...
discharge_disposition_id_2	0	0	0	0	0
discharge_disposition_id_7	0	0	0	0	0
discharge_disposition_id_10	0	0	0	0	0
discharge_disposition_id_18	0	0	0	0	0
discharge_disposition_id_19	0	0	0	0	0
discharge_disposition_id_20	0	0	0	0	0
discharge_disposition_id_27	0	0	0	0	0
discharge_disposition_id_28	0	0	0	0	0
admission_source_id_4	0	0	0	0	0
admission_source_id_7	1	1	1	1	0

	1	2	3	4	5
admission_source_id_8	0	0	0	0	0
admission_source_id_9	0	0	0	0	0
admission_source_id_11	0	0	0	0	0
max_glu_serum_0	0	0	0	0	0
max_glu_serum_1	0	0	0	0	0
A1Cresult_0	0	0	0	0	0
A1Cresult_1	0	0	0	0	0
level1_diag1_1.0	0	0	0	0	1
level1_diag1_2.0	0	0	0	0	0
level1_diag1_3.0	0	0	0	0	0
level1_diag1_4.0	0	0	0	0	0
level1_diag1_5.0	0	0	0	0	0
level1_diag1_6.0	0	0	0	0	0
level1_diag1_7.0	0	0	0	0	0
level1_diag1_8.0	0	0	0	1	0
AfricanAmerican	0	1	0	0	0
Asian	0	0	0	0	0
Caucasian	1	0	1	1	1
Hispanic	0	0	0	0	0
Other	0	0	0	0	0

81 rows x 5 columns

Modeling

In [69]:

```
feature_set = ['age', 'time_in_hospital', 'num_procedures', 'num_medications', 'num
ber_outpatient_loglp',

               'number_emergency_loglp', 'number_inpatient_loglp', 'number_diagno
ses', 'metformin',

               'repaglinide', 'nateglinide', 'chlorpropamide', 'glimepiride', 'gl
ipizide', 'glyburide',

               'pioglitazone', 'rosiglitazone', 'acarbose', 'tolazamide', 'insuli
n', 'glyburide-metformin',

               'AfricanAmerican', 'Asian', 'Caucasian', 'Hispanic', 'Other', 'gen
der_1',

               'admission_type_id_3', 'admission_type_id_5', 'discharge_dispositi
on_id_2', 'discharge_disposition_id_7',

               'discharge_disposition_id_10', 'discharge_disposition_id_18', 'adm
ission_source_id_4',

               'admission_source_id_7', 'admission_source_id_9', 'max_glu_serum_0
', 'max_glu_serum_1', 'A1Cresult_0',

               'A1Cresult_1', 'num_medications|time_in_hospital', 'num_medication
s|num_procedures',

               'time_in_hospital|num_lab_procedures', 'num_medications|num_lab_pr
ocedures', 'num_medications|number_diagnoses',

               'age|number_diagnoses', 'change|num_medications', 'number_diagnose
s|time_in_hospital',

               'num_medications|numchange', 'level1_diag1_1.0', 'level1_diag1_2.0
', 'level1_diag1_3.0', 'level1_diag1_4.0',
```

```
'level1_diag1_5.0','level1_diag1_6.0', 'level1_diag1_7.0', 'level1_diag1_8.0']
```

In [70]:

```
X = df_pd[feature_set]
y = df_pd['readmitted']
```

Logistic Regression ¶

In [71]:

```
df_pd['readmitted'].value_counts()
```

Out[71]:

```
0    54635
1     5071
Name: readmitted, dtype: int64
```

In [72]:

```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import cross_val_score

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.20, random_state=0)

logit = LogisticRegression(fit_intercept=True, penalty='l1')
logit.fit(X_train, y_train)
```

/opt/conda/lib/python3.6/site-packages/sklearn/linear_model/logistic.py:433: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
FutureWarning)

Out[72]:

```
LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
                    intercept_scaling=1, max_iter=100, multi_class='warn',
                    n_jobs=None, penalty='l1', random_state=None, solver='warn',
                    tol=0.0001, verbose=0, warm_start=False)
```

In [73]:

```
logit_pred = logit.predict(X_test)

pd.crosstab(pd.Series(y_test, name = 'Actual'), pd.Series(logit_pred, name = 'Predict'), margins = True)
```

Out[73]:

Predict	0	All
Actual		
0	1466	1466
1	167	167
All	1633	1633

In [74]:


```

from sklearn.metrics import accuracy_score, precision_score, recall_score
print("Accuracy is {0:.2f}".format(accuracy_score(y_test, logit_pred)))
print("Precision is {0:.2f}".format(precision_score(y_test, logit_pred)))
print("Recall is {0:.2f}".format(recall_score(y_test, logit_pred)))

```

```

Accuracy is 0.91
Precision is 0.00
Recall is 0.00
/opt/conda/lib/python3.6/site-packages/sklearn/metrics/classification.py:1143: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no predicted samples.
  'precision', 'predicted', average, warn_for)

```

Since our target variable is having class imbalance problem, So will use SMOTE technique to resolve it

In [75]:

```

from imblearn.over_sampling import SMOTE
from collections import Counter
print('Original dataset shape {}'.format(Counter(y_train)))
sm = SMOTE(random_state=20)
train_input_new, train_output_new = sm.fit_sample(X_train, y_train)
print('New dataset shape {}'.format(Counter(train_output_new)))

```

```

Using TensorFlow backend.
Original dataset shape Counter({0: 43711, 1: 4053})
New dataset shape Counter({0: 43711, 1: 43711})

```

In [76]:

```

train_input_new = pd.DataFrame(train_input_new, columns = list(X.columns))
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import cross_val_score
X_train, X_test, y_train, y_test = train_test_split(train_input_new, train_output_new, test_size=0.20, random_state=0)
logit = LogisticRegression(fit_intercept=True, penalty='l1')
logit.fit(X_train, y_train)

```

```

/opt/conda/lib/python3.6/site-packages/sklearn/linear_model/logistic.py:433: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
  FutureWarning)

```

Out[76]:

```

LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
  intercept_scaling=1, max_iter=100, multi_class='warn',
  n_jobs=None, penalty='l1', random_state=None, solver='warn',
  tol=0.0001, verbose=0, warm_start=False)

```

In [77]:

```

logit_pred = logit.predict(X_test)

```

```
pd.crosstab(pd.Series(y_test, name = 'Actual'), pd.Series(logit_pred, name = 'Predict'), margins = True)
```

Out[77]:

Predict	0	1	All
Actual			
0	5805	2901	8706
1	3895	4884	8779
All	9700	7785	17485

In [78]:

```
print("Accuracy is {0:.2f}".format(accuracy_score(y_test, logit_pred)))
print("Precision is {0:.2f}".format(precision_score(y_test, logit_pred)))
print("Recall is {0:.2f}".format(recall_score(y_test, logit_pred)))

accuracy_logit = accuracy_score(y_test, logit_pred)
precision_logit = precision_score(y_test, logit_pred)
recall_logit = recall_score(y_test, logit_pred)
```

```
Accuracy is 0.61
Precision is 0.63
Recall is 0.56
```

Decision Tree

In [79]:

```
feature_set_no_int = ['age', 'time_in_hospital', 'num_procedures', 'num_medications',
                      'number_outpatient_loglp',
                      'number_emergency_loglp', 'number_inpatient_loglp', 'number_diagnoses', 'metformin',
                      'repaglinide', 'nateglinide', 'chlorpropamide', 'glimepiride', 'glipizide',
                      'glyburide', 'pioglitazone', 'rosiglitazone', 'acarbose',
                      'tolazamide', 'insulin', 'glyburide-metformin',
                      'AfricanAmerican', 'Asian', 'Caucasian',
                      'Hispanic', 'Other', 'gender_1',
                      'admission_type_id_3', 'admission_type_id_5',
                      'discharge_disposition_id_2', 'discharge_disposition_id_7',
                      'discharge_disposition_id_10', 'discharge_disposition_id_18',
                      'admission_source_id_4', 'admission_source_id_7',
                      'admission_source_id_9', 'max_glu_serum_0',
                      'max_glu_serum_1', 'A1Cresult_0', 'A1Cresult_1',
                      'level1_diag1_1.0',
                      'level1_diag1_2.0',
                      'level1_diag1_3.0',
                      'level1_diag1_4.0',
```

```
'level1_diag1_5.0',  
'level1_diag1_6.0',  
'level1_diag1_7.0',  
'level1_diag1_8.0']
```

In [80]:

```
X = df_pd[feature_set_no_int]  
y = df_pd['readmitted']  
df_pd['readmitted'].value_counts()
```

Out[80]:

```
0    54635  
1     5071  
Name: readmitted, dtype: int64
```

In [81]:

```
print('Original dataset shape {}'.format(Counter(y)))  
smt = SMOTE(random_state=20)  
train_input_new, train_output_new = smt.fit_sample(X, y)  
print('New dataset shape {}'.format(Counter(train_output_new)))  
train_input_new = pd.DataFrame(train_input_new, columns = list(X.columns))  
X_train, X_test, y_train, y_test = train_test_split(train_input_new, train_output_n  
ew, test_size=0.20, random_state=0)
```

```
Original dataset shape Counter({0: 54635, 1: 5071})  
New dataset shape Counter({0: 54635, 1: 54635})
```

In [82]:

```
from sklearn.tree import DecisionTreeClassifier  
  
dtree = DecisionTreeClassifier(max_depth=28, criterion = "entropy", min_samples_spl  
it=10)  
  
dtree.fit(X_train, y_train)
```

Out[82]:

```
DecisionTreeClassifier(class_weight=None, criterion='entropy', max_depth=28,  
    max_features=None, max_leaf_nodes=None,  
    min_impurity_decrease=0.0, min_impurity_split=None,  
    min_samples_leaf=1, min_samples_split=10,  
    min_weight_fraction_leaf=0.0, presort=False, random_state=None,  
    splitter='best')
```

In [83]:

```
dtree_pred = dtree.predict(X_test)  
  
pd.crosstab(pd.Series(y_test, name = 'Actual'), pd.Series(dtree_pred, name = 'Predi  
ct'), margins = True)
```

Out[83]:

Predict	0	1	All
Actual			
0	10244	692	10936

Predict	0	1	All
Actual			
1	1130	9788	10918
All	11374	10480	21854

In [84]:

```
print("Accuracy is {0:.2f}".format(accuracy_score(y_test, dtree_pred)))
print("Precision is {0:.2f}".format(precision_score(y_test, dtree_pred)))
print("Recall is {0:.2f}".format(recall_score(y_test, dtree_pred)))

accuracy_dtree = accuracy_score(y_test, dtree_pred)
precision_dtree = precision_score(y_test, dtree_pred)
recall_dtree = recall_score(y_test, dtree_pred)
```

```
Accuracy is 0.92
Precision is 0.93
Recall is 0.90
```

In [85]:

```
# Create list of top most features based on importance
feature_names = X_train.columns
feature_imports = dtree.feature_importances_

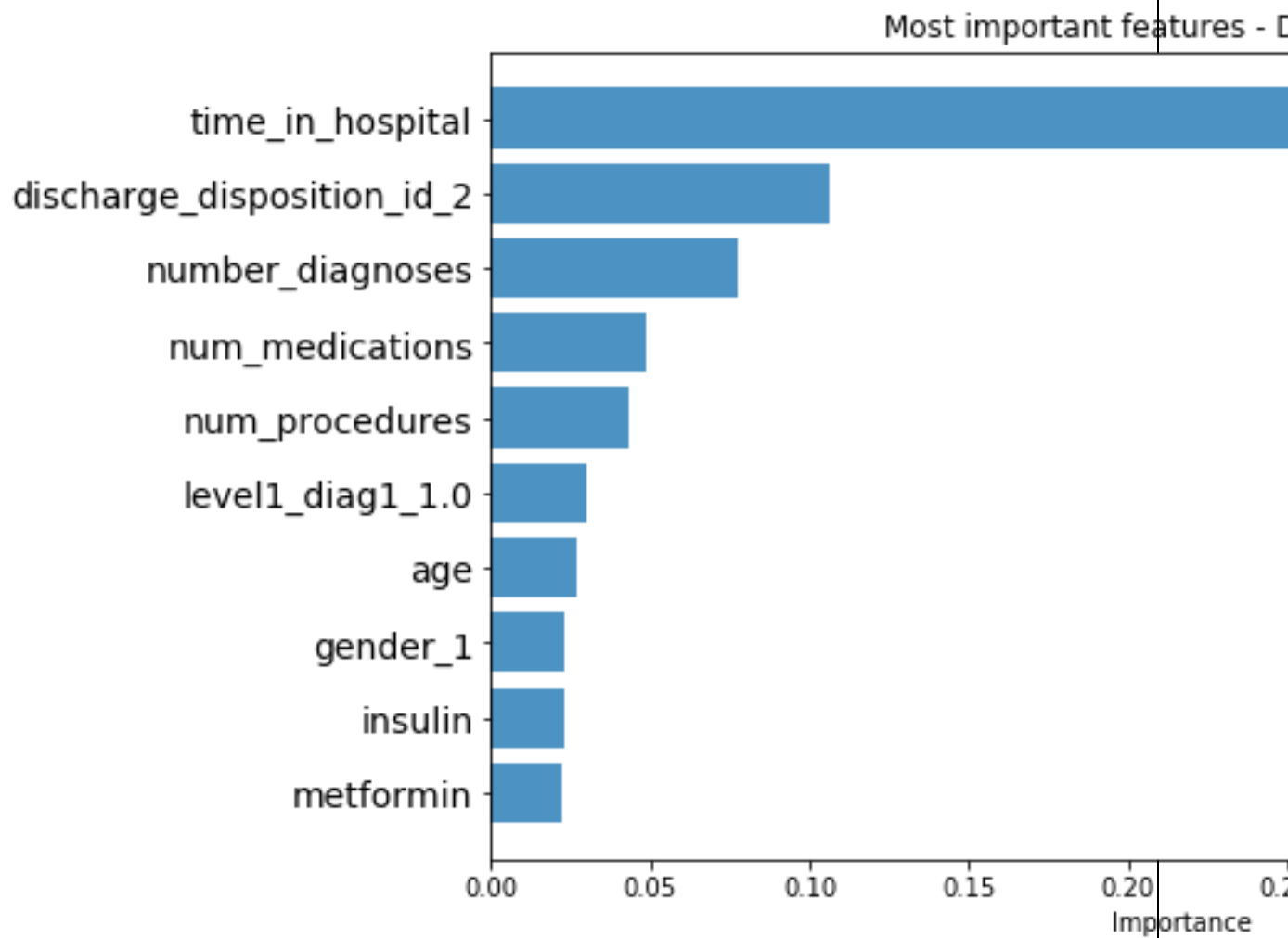
most_imp_features = pd.DataFrame([f for f in zip(feature_names, feature_imports)], columns=["Feature", "Importance"]).nlargest(10, "Importance")
most_imp_features.sort_values(by="Importance", inplace=True)
print(most_imp_features)

plt.figure(figsize=(10,6))

plt.barh(range(len(most_imp_features)), most_imp_features.Importance, align='center', alpha=0.8)

plt.yticks(range(len(most_imp_features)), most_imp_features.Feature, fontsize=14)
plt.xlabel('Importance')
plt.title('Most important features - Decision Tree')
plt.show()
```

	Feature	Importance
8	metformin	0.022040
19	insulin	0.022958
26	gender_1	0.023195
0	age	0.026882
40	levell_diag1_1.0	0.029614
2	num_procedures	0.043294
3	num_medications	0.048912
7	number_diagnoses	0.077183
29	discharge_disposition_id 2	0.106147
1	time_in_hospital	0.413369



Random Forest

In [86]:

```
X = df_pd[feature_set_no_int]
y = df_pd['readmitted']

print('Original dataset shape {}'.format(Counter(y)))
smt = SMOTE(random_state=20)
train_input_new, train_output_new = smt.fit_sample(X, y)
print('New dataset shape {}'.format(Counter(train_output_new)))
train_input_new = pd.DataFrame(train_input_new, columns = list(X.columns))
X_train, X_test, y_train, y_test = train_test_split(train_input_new, train_output_new, test_size=0.20, random_state=0)
```

```
Original dataset shape Counter({0: 54635, 1: 5071})
New dataset shape Counter({0: 54635, 1: 54635})
```

In [87]:

```
from sklearn.ensemble import RandomForestClassifier

rm = RandomForestClassifier(n_estimators = 10, max_depth=25, criterion = "gini", min_samples_split=10)

rm.fit(X_train, y_train)
```

Out[87]:

```
RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                        max_depth=25, max_features='auto', max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=10,
                        min_weight_fraction_leaf=0.0, n_estimators=10, n_jobs=None,
                        oob_score=False, random_state=None, verbose=0,
                        warm_start=False)
```

In [88]:

```
rm_prd = rm.predict(X_test)

pd.crosstab(pd.Series(y_test, name = 'Actual'), pd.Series(rm_prd, name = 'Predict')
, margins = True)
```

Out[88]:

Predict	0	1	All
Actual			
0	10795	141	10936
1	1141	9777	10918
All	11936	9918	21854

In [89]:

```
print("Accuracy is {0:.2f}".format(accuracy_score(y_test, rm_prd)))
print("Precision is {0:.2f}".format(precision_score(y_test, rm_prd)))
print("Recall is {0:.2f}".format(recall_score(y_test, rm_prd)))

accuracy_rm = accuracy_score(y_test, rm_prd)
precision_rm = precision_score(y_test, rm_prd)
recall_rm = recall_score(y_test, rm_prd)
```

```
Accuracy is 0.94
Precision is 0.99
Recall is 0.90
```

In [90]:

```
# Create list of top most features based on importance
feature_names = X_train.columns

feature_imports = rm.feature_importances_

most_imp_features = pd.DataFrame([f for f in zip(feature_names,feature_imports)], columns=["Feature", "Importance"]).nlargest(10, "Importance")

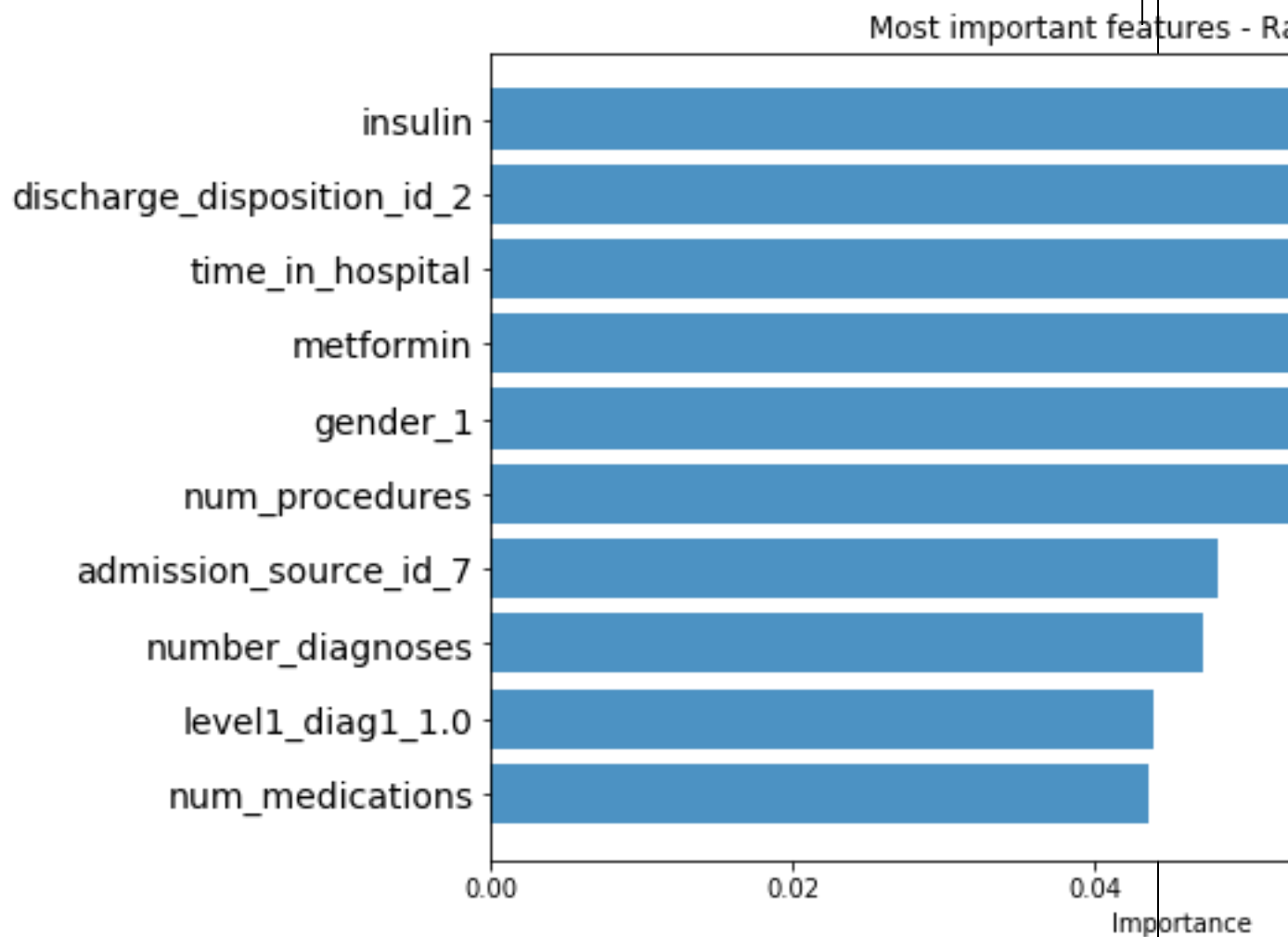
most_imp_features.sort_values(by="Importance", inplace=True)

plt.figure(figsize=(10,6))

plt.barh(range(len(most_imp_features)), most_imp_features.Importance, align='center', alpha=0.8)

plt.yticks(range(len(most_imp_features)), most_imp_features.Feature, fontsize=14)
```

```
plt.xlabel('Importance')
plt.title('Most important features - Random Forest ')
plt.show()
```



Model Comparision

In [91]:

```
plt.figure(figsize=(14, 7))
ax = plt.subplot(111)

models = ['Logistic Regression', 'Decision Tree', 'Random Forests']
values = [accuracy_logit, accuracy_dtree, accuracy_rm]
model = np.arange(len(models))

plt.bar(model, values, align='center', width = 0.15, alpha=0.7, color = 'red', label= 'accuracy')
plt.xticks(model, models)
```

```

ax = plt.subplot(111)

models = ['Logistic Regression', 'Decision Tree', 'Random Forests']
values = [precision_logit, precision_dtree, precision_rm]
model = np.arange(len(models))

plt.bar(model+0.15, values, align='center', width = 0.15, alpha=0.7, color = 'blue'
, label = 'precision')
plt.xticks(model, models)


ax = plt.subplot(111)

models = ['Logistic Regression', 'Decision Tree', 'Random Forests' ]
values = [recall_logit, recall_dtree, recall_rm, ]
model = np.arange(len(models))

plt.bar(model+0.3, values, align='center', width = 0.15, alpha=0.7, color = 'green'
, label = 'recall')
plt.xticks(model, models)


plt.ylabel('Performance Metrics for Different models')
plt.title('Model')

# removing the axis on the top and right of the plot window
ax.spines['right'].set_visible(False)
ax.spines['top'].set_visible(False)
ax.legend()

plt.show()

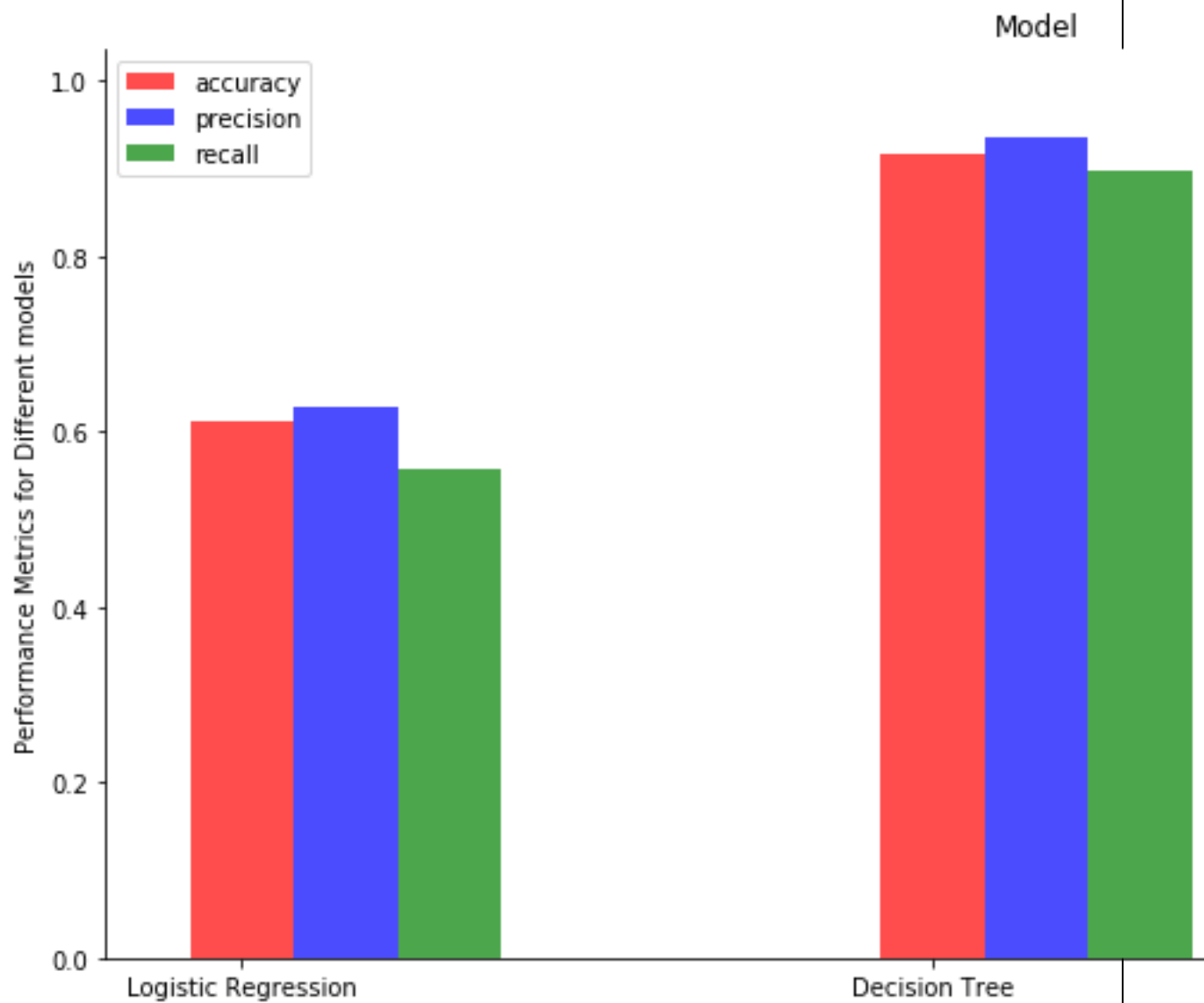
```

/opt/conda/lib/python3.6/site-packages/matplotlib/figure.py:98: MatplotlibDeprecationWarning:
Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

"Adding an axes using the same arguments as a previous axes "
/opt/conda/lib/python3.6/site-packages/matplotlib/figure.py:98: MatplotlibDeprecationWarning:

Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

"Adding an axes using the same arguments as a previous axes "



In [92]: