



# DATA EXPLORATION AND MINING

---

FEB 2021/ GROUP ASSESSMENT 1

## Exploration of a large medical dataset

## Question 1

Have identified 3 variables: Gender (Male/Female), Change (Ch  
they are binary in nature, these 3 variables can be mapped int

ange/No change) and DiabetesMed (Yes/No). Since  
o 0 and 1.

## Original data type

```
# read in diabetic patient dataset
df = pd.read_csv('D1.csv')
print(df.info())

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 51766 entries, 0 to 51765
Data columns (total 37 columns):
# binary
# binary
```

## Map gender to binary (int64)

```
1 # map gender into binary variable
2
3 gender_binary = {'Female':0, 'Male':1, 'Unknown/Invalid':np.Nan}
4 df['gender'] = df['gender'].map(gender_binary)
5
6
```

```
1 df['gender'].value_counts()
2
```

```
0.0    27718
1.0    24044
Name: gender, dtype: int64
```

## Map change to binary (int64)

```
# map change into binary variable
change_binary = {'No':0, 'Ch':1}
df['change'] = df['change'].map(change_binary)
```

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

```
1    26663
0    25103
Name: change, dtype: int64
```

## Map diabetesMed to binary (int64)

```
# map diabetesMed into binary variable
diabetesMed_binary = {'No':0, 'Yes':1}
df['diabetesMed'] = df['diabetesMed'].map(diabetesMed_binary)
```

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

```
1    40984
0    10782
Name: diabetesMed, dtype: int64
```

## Question 2

Using suitable statistical measures and functions:

1) Identify and report the skewness present in the variables.

As question 3 will use visualization, for this question we focus on descriptive statistics to view the numeric data.

in_hospital	num_lab_procedures	num_procedures	num_medications	number_outpatient	number_emergency	number_inpatient	number_diagnoses
766.000000	51766.000000	51766.000000	51766.000000	51766.000000	51766.000000	51766.000000	51766.000000
4.237337	43.853205	1.308214	16.735850	0.512846	0.269173	0.686879	7.928486
2.870999	19.945471	1.721694	8.063522	1.544866	1.146923	1.320968	1.693707
1.000000	1.000000	0.000000	1.000000	0.000000	0.000000	0.000000	1.000000
2.000000	33.000000	0.000000	11.000000	0.000000	0.000000	0.000000	7.000000
3.000000	45.000000	1.000000	16.000000	0.000000	0.000000	0.000000	9.000000
6.000000	58.000000	2.000000	21.000000	0.000000	0.000000	1.000000	9.000000
14.000000	132.000000	6.000000	75.000000	42.000000	76.000000	19.000000	16.000000

Below are the numeric variables:

- time\_in\_hospital: mean = 4.237 337, median = 3 ==> mean > median ==> positive skewed
- num\_lab\_procedures: mean = 43.853 205, median = 45 ==> mean < median ==> negative skewed
- num\_procedures: mean = 1.308214 median = 1 ==> mean > median ==> positive skewed
- num\_medications: mean = 16.7359, median = 16 ==> mean > median ==> positive skewed

- number\_outpatient: mean = 0.512846, median = 0 ==> mean > median => positive skewed
- number\_emergency: mean = 0.269173, median = 0 ==> mean > median => negative skewed
- number\_inpatient: mean = 0.686879, median = 0 ==> mean > median => positive skewed
- number\_diagnoses: mean = 7.928486, median = 9 ==> mean < median => negative skewed

Below are the nominal variables. Hence descriptive statistics are not applicable.

- encounter\_id
- patient\_nbr
- gender
- admission\_type\_id
- discharge\_disposition\_id
- admission\_source\_id
- change
- diabetesMed

2) There may be inconsistencies or errors in the data. List the errors identified and detail how you have identified them.

List out and review the unique values in each variable to identify if there are any unusual values.

```

1 # Looking at unique values in the dataset and see if there is any unusual values
2
3 for col in list(df):
4     print(col)
5     print(df[col].unique())
6
7

encounter_id
[150645834 150646782 150659562 ... 443854148 443857166 443867222]
patient_nbr
[ 66173058 104923647 30917367 ... 140199494 120975314 175429310]
race
['Caucasian' 'AfricanAmerican' 'Hispanic' '?' 'Other' 'Asian']
gender
[ 0.  1. nan]
age
[['[60-70)' '[80-90)' '[70-80)' '[40-50)' '[50-60)' '[90-100)' '[30-40)'
'[20-30)' '[10-20)' '[0-10)' '?']
weight
['?' '[100-125)' '[50-75)' '[75-100)' '[0-25)' '[125-150)' '[25-50)'
'[150-175)' '[175-200)' '>200']
admission_type_id
[2 3 1 5 6 8 7 4]
discharge_disposition_id
[ 1 2 6 3 22 11 7 23 13 4 15 9 28 5 19 8 25 14 24 27 18]
admission_source_id
[ 1 7 5 17 4 2 3 6 9 10 22 20 8 14 11 25 13]
time_in_hospital
[ 4 8 2 7 3 6 5 1 13 12 11 9 10 14]
payer_code
['?' 'MD' 'BC' 'MC' 'HM' 'CP' 'SP' 'OG' 'UN' 'DM' 'CM' 'PO' 'SI' 'WC' 'CH'
'OT' 'MP' 'FR' ]

```

Have identified the following variables with "?" as values:

- race
- age
- weight
- chlorpropamide
- payer\_code
- medical\_specialty
- diag\_2

Then check to see how many "?" in each variable:

```
1 # view the number of '?' data
2
3 num_missing = (df[['race', 'age', 'weight', 'chlorpropamide', 'payer_code', 'medical_specialty',
4                     'diag_2']] == "?").sum()
5 print(num_missing)
```

race	1016
age	10
weight	50431
chlorpropamide	9
payer_code	7601
medical_specialty	32203
diag_2	92
dtype: int64	

3a) What is the average time stay in hospital of a female patient who was readmitted in less than 30 days?

4.524752 days

```
1 print(df.groupby(['gender', 'readmitted'], as_index=False)[['time_in_hospital']].mean())
2 # Q1 mapped 0: Female
```

gender	readmitted	time_in_hospital
0.0	<30	4.524752
0.0	>30	4.465249
0.0	NO	4.207442
1.0	<30	4.482972
1.0	>30	4.192097
1.0	NO	4.012216

3b) Which age group has the highest risk of being readmitted within 30 days? 70-80 has the highest risk

```
4 df_readmitted_age_grouped = df.groupby(['age'])['readmitted'].value_counts()
5 df_readmitted_age_grouped
```

age	readmitted	value_counts
[40-50)	NO	2611
	>30	1537
	<30	471
[50-60)	NO	4789
	>30	2905
	<30	788
[60-70)	NO	6270
	>30	4074
	<30	1274
[70-80)	NO	6770
	>30	4767
	<30	1435
[80-90)	NO	4930
	>30	3611
	<30	1130
[90-100)	NO	940
	>30	483
	<30	193

Name: readmitted, dtype: int64

3c) How many age groups have more than 3000 cases of being readmitted?

4 age groups

```
1 filter_readmitted = df[(df['readmitted']==>30') | (df['readmitted']==<30')]
2 readmitted_by_age = filter_readmitted['age'].value_counts()
3 readmitted_by_age
```

[70-80)	6202
[60-70)	5348
[80-90)	4741
[50-60)	3693
[40-50)	2008
[30-40)	709
[90-100)	676
[20-30)	351
[10-20)	91
[0-10)	6
?	3

Name: age, dtype: int64

3d) Which are the top-three race categories according to the number of readmission cases?

1st - Caucasian, 2nd - AfricanAmerican, 3rd - Other (NOTE: excluding '?')

```
1 readmitted_by_race = filter_readmitted['race'].value_counts()
2 readmitted_by_race
```

Caucasian	19120
AfricanAmerican	3530
?	367
Other	347
Hispanic	344
Asian	120

Name: race, dtype: int64

**Question 3**

Using suitable visualisation plots:

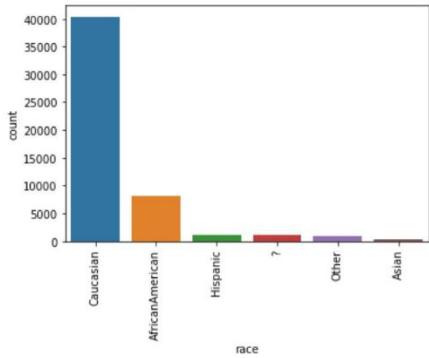
- 1) Understand the distribution of variables and identify data quality problems.

Findings are summarized in # lines on top of each plot.

```

1 # race
2 race_plot = sns.countplot(data=df, x='race')
3 race_plot.set_xticklabels(race_plot.get_xticklabels(), rotation=90)
4 plt.show()
5
6 # missing data with '?', Caucasian is the majority

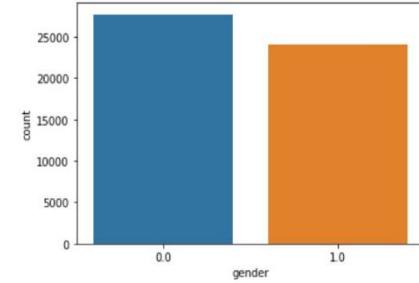
```



```

1 # gender (Female = 0, Male = 1)
2 gender_plot = sns.countplot(data=df, x='gender')
3 plt.show()
4
5 # Not include missing data with 'nan' in plot, about equal shares in gender

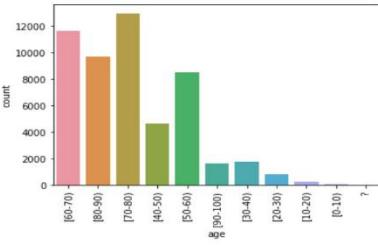
```



```

1 # age
2 age_plot = sns.countplot(data=df, x='age')
3 age_plot.set_xticklabels(age_plot.get_xticklabels(), rotation=90)
4 plt.show()
5
6 # missing data with '?', five skewed, 40 and above are the majority age group

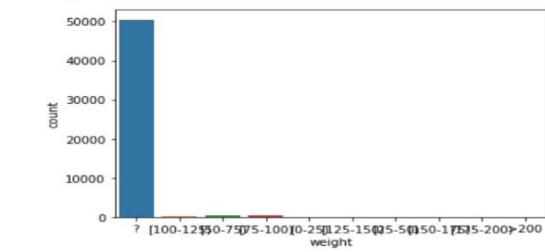
```



```

1 # weight
2 weight_plot = sns.countplot(data=df, x='weight')
3 plt.show()
4
5 # 50431 records are '?' (97.42% of data)

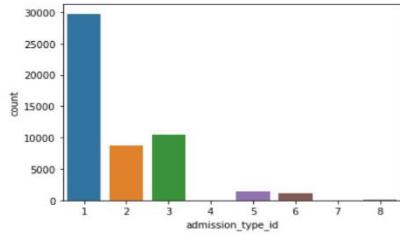
```



```

1 # admission type
2 admit_type_plot = sns.countplot(data=df, x='admission_type_id')
3 plt.show()
4
5 # Majority admitted under Emergency

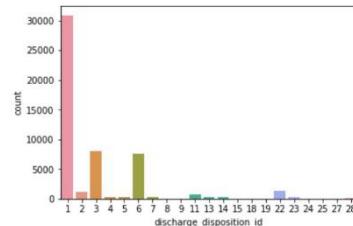
```



```

1 # discharge disposition
2 discharge_plot = sns.countplot(data=df, x='discharge_disposition_id')
3 plt.show()
4
5 # majority is discharged to home. followed by discharged/transferred to SNF and
6 # discharged/transferred to home with home health service

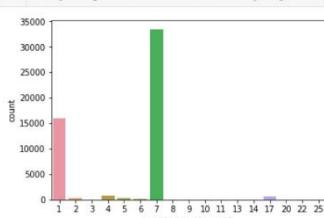
```



```

1 # admission source
2 admit_source_plot = sns.countplot(data=df, x='admission_source_id')
3 plt.show()
4
5 # 1: Physician Referral, 7: Emergency Room
6 # majority are admitted via emergency room, which co-incide with admission type

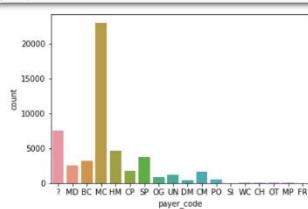
```



```

1 # payer_code (Unique identifier assigned to each insurance company)
2 payer_code_plot = sns.countplot(data=df, x='payer_code')
3 plt.show()
4
5 df['payer_code'].value_counts()
6
7 # missing data with "?" 7601 (about 1.5% of total data!)
8 # majority is "MC"
9

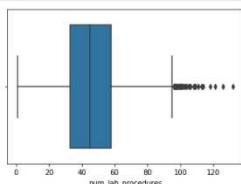
```



```

1 # number of lab procedures
2 lab_procedures_plot = sns.boxplot(df['num_lab_procedures'])
3 plt.show()
4
5 # per statistic median=mean, it is left skewed . With outliers.
6 # 20 records are >3sd
7 lab_procedures_outliers = (df['num_lab_procedures'] > (43.853205+19.945471*3)).value_counts()
8 lab_procedures_outliers

```



```

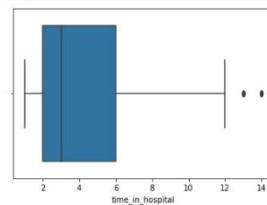
False 51746
True 20
Name: num_lab_procedures, dtype: int64

```

```

1 # time in hospital
2 time_plot = sns.boxplot(df['time_in_hospital'])
3 plt.show()
4
5 # positive skewed. 75% majority is below 6 days. median of around 3 days.
6 # with outliers of 13 & 14 days
7
8 # 377 records are >3sd (1.89% of data)
9 time_in_hospital_outliers = (df['time_in_hospital'] > (4.237+2.870999*3)).value_counts()
10 time_in_hospital_outliers

```



```

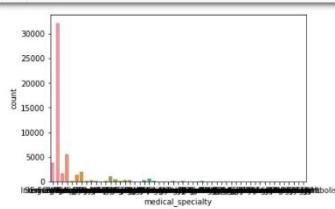
False 50789
True 977
Name: time_in_hospital, dtype: int64

```

```

1 # medical specialty (Indicates specialty of the admitting physician,for example, cardiology, surgeon, etc.)
2 medical_specialty_plot = sns.countplot(data=df, x='medical_specialty')
3 plt.show()
4
5 df['medical_specialty'].value_counts()
6
7 # missing data with "?" 32203 (about 62% of total data!)
8 # majority is 'internalMedicine'
9

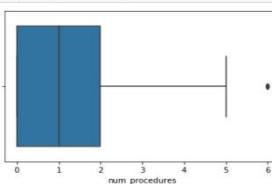
```



```

1 ??
2 Emergency/Trauma 32203
3 InternalMedicine 5616
4
5
6 df['medical_specialty'].value_counts()
7
8 # number of procedures (other than lab procedures)
9 no_procedures_plot = sns.boxplot(df['num_procedures'])
10 plt.show()
11
12 # positive skewed with outlier. Median is 1.
13 # 3873 records are >3sd (5.55%)
14 num_procedures_outliers = (df['num_procedures'] == 6).value_counts()
15 num_procedures_outliers

```



```

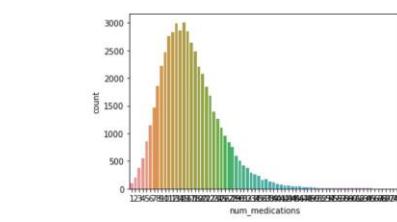
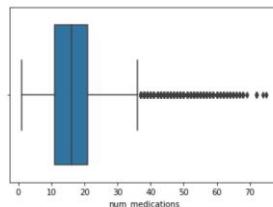
False 48893
True 2873
Name: num_procedures, dtype: int64

```

```

[49]: M 1 # number of medications
2 no_medications_plot = sns.boxplot(df['num_medications'])
3 plt.show()
4
5 no_medications_countplot = sns.countplot(data=df, x='num_medications')
6 plt.show()
7
8 # positive skewed. Median of 15
9 # 626 records are >3sd (1.21% of data)
10 no_medications_outliers = (df['num_medications'] > (16.735850+8.063522*3)).value_counts()
11 no_medications_outliers

```



```

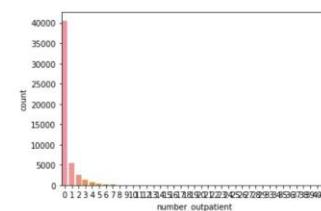
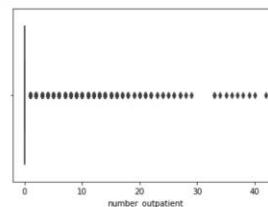
Out[49]: False 51140
True 626
Name: num_medications, dtype: int64

```

```

1 # Number of outpatient visits of the patient in the year preceding the encounter
2 no_outpatient_plot = sns.boxplot(df['number_outpatient'])
3 plt.show()
4
5 no_outpatient_countplot = sns.countplot(data=df, x='number_outpatient')
6 plt.show()
7
8 # positive skewed. Majority is 0. With outliers
9 # 737 records are >3sd (1.42% of data)
10 number_outpatient_outliers = (df['number_outpatient'] > (0.512846+1.544066*3)).value_counts()
11 number_outpatient_outliers

```



```

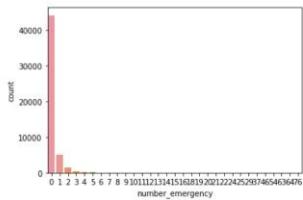
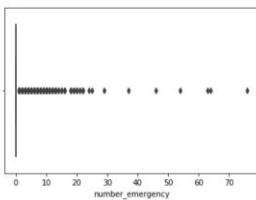
False 51029
True 737
Name: number_outpatient, dtype: int64

```

```

1 # Number of emergency visits of the patient in the year preceding the encounter
2 no_outpatient_plot = sns.boxplot(df['number_emergency'])
3 plt.show()
4
5 no_outpatient_count = sns.countplot(data=df, x='number_emergency')
6 plt.show()
7
8 # positive skewed. Majority is 0 with outliers
9 # 667 records are >3sd (1.29% of data)
10 number_emergency_outliers = (df['number_emergency'] > (0.269173+1.146923*3)).value_counts()
11 number_emergency_outliers

```

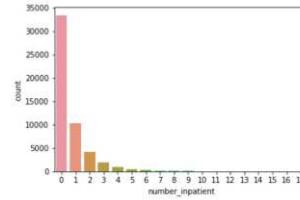
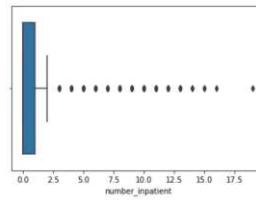


False 51099  
True 667  
Name: number\_emergency, dtype: int64

```

1 # Number of inpatient visits of the patient in the year preceding the encounter
2 no_outpatient_plot = sns.boxplot(df['number_inpatient'])
3 plt.show()
4
5 no_outpatient_count = sns.countplot(data=df, x='number_inpatient')
6 plt.show()
7
8 # positive skewed. Majority is 0 with outliers
9 # 1146 records are >3sd (2.21% of data)
10 number_inpatient_outliers = (df['number_inpatient'] > (0.686879+1.320968*3)).value_counts()
11 number_inpatient_outliers

```

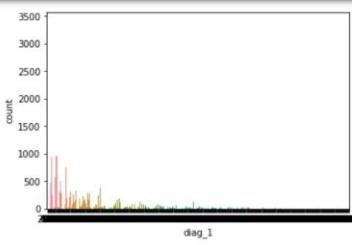


False 50620  
True 1146  
Name: number\_inpatient, dtype: int64

```

1 # diag_1 (primary diagnosis), in accordance with first three digits of ICD9
2 diag_1_count = sns.countplot(data=df, x='diag_1')
3 plt.show()
4
5 # majority is '428' (Congestive heart failure, unspecified) and
6 # '414' (Other forms of chronic ischemic heart disease)
7 # there are 649 types of diagnosis
8 df['diag_1'].value_counts()

```

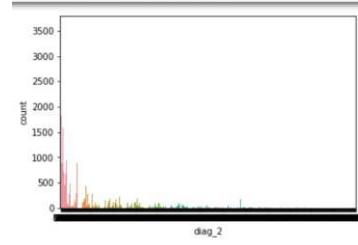


428 3395  
414 2746  
486 1808

```

1 # diag_2 (secondary diagnosis), in accordance with first three digits of ICD9
2 diag_2_count = sns.countplot(data=df, x='diag_2')
3 plt.show()
4
5 # majority is '276' (Hyperosmolality and/or hypernatremia) and
6 # '428' (Congestive heart failure, unspecified)
7 # there are 683 types of diagnosis
8 df['diag_2'].value_counts()

```

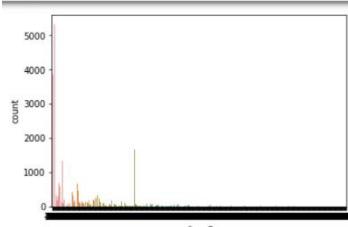


276 3604  
428 3177  
250 2685

```

1 # diag_3 (additional secondary), in accordance with first three digits of ICD9
2 diag_3_count = sns.countplot(data=df, x='diag_3')
3 plt.show()
4
5 # majority is '250' (Diabetes), '401' (Malignant essential hypertension)
6 # there are 726 types of diagnosis
7 df['diag_3'].value_counts()

```

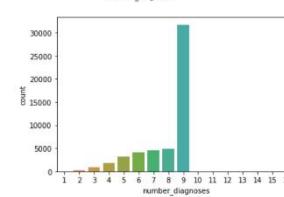
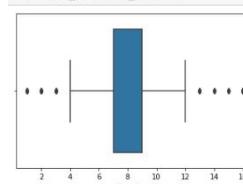


250 5323  
401 3861  
276 2968

```

1 # Number of diagnoses entered to the system
2 no_diagnoses_plot = sns.boxplot(df['number_diagnoses'])
3 plt.show()
4
5 no_diagnoses_count = sns.countplot(data=df, x='number_diagnoses')
6 plt.show()
7
8 # mean = 7.928485878762122, median = 0, median > mean ==> negative skewed
9 # outliers on both sides
10
11 # 67 records are <=sd and 304 records are >3sd (0.71% of data)
12 number_diagnoses_outliers = (df['number_diagnoses'] > (7.928485878762122+1.693707*3)).value_counts()
13 number_diagnoses_outliers
14
15 number_diagnoses_outliers = (df['number_diagnoses'] < (7.928485878762122-1.693707*3)).value_counts()
16 number_diagnoses_outliers

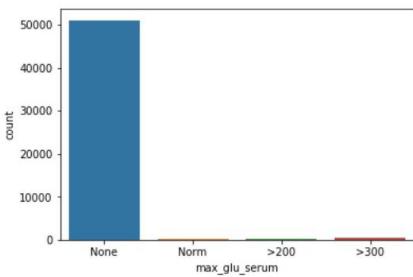
```



```

1 # Glucose serum test result.
2 # Indicates the range of the result or if the test was not taken("none")
3 max_glu_serum_count = sns.countplot(data=df, x='max_glu_serum')
4 plt.show()
5
6 # majority are None

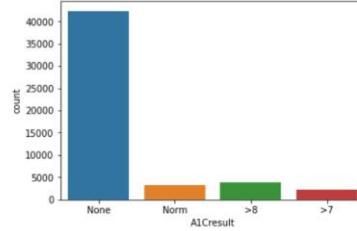
```



```

1 # Alc test result.
2
3 # Values include:>8" if the result was greater than 8%,
4 # ">7" if the result was greater than 7% but less than 8%,
5 # "normal" if the result was less than 7%, and
6 # "none" if not measured
7
8 Alcresult_count = sns.countplot(data=df, x='AlcResult')
9 plt.show()
10
11 # Majority are None.

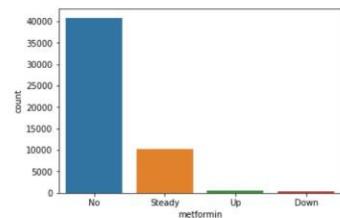
```



```

1 # metformin
2
3 # A diabetes medications
4 # The values of these variables indicate whether the drug was prescribed
5 # "up" if the dosage was increased during the encounter,
6 # "down" if the dosage was decreased,
7 # "steady" if the dosage did not change, and
8 # "no" if the drug was not prescribed.
9
10 metformin_count = sns.countplot(data=df, x='metformin')
11 plt.show()
12
13 # majority are not prescribed

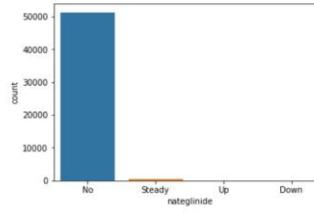
```



```

1 # nateglinide
2
3 # A diabetes medications
4 # The values of these variables indicate whether the drug was prescribed
5 # "up" if the dosage was increased during the encounter,
6 # "down" if the dosage was decreased,
7 # "steady" if the dosage did not change, and
8 # "no" if the drug was not prescribed.
9
10 nateglinide_count = sns.countplot(data=df, x='nateglinide')
11 plt.show()
12
13 # majority are not prescribed

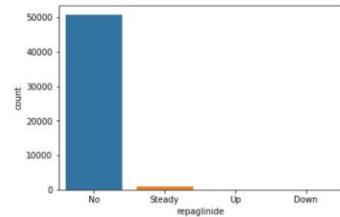
```



```

1 # repaglinide
2
3 # A diabetes medications
4 # The values of these variables indicate whether the drug was prescribed
5 # "up" if the dosage was increased during the encounter,
6 # "down" if the dosage was decreased,
7 # "steady" if the dosage did not change, and
8 # "no" if the drug was not prescribed.
9
10 repaglinide_count = sns.countplot(data=df, x='repaglinide')
11 plt.show()
12
13 # majority are not prescribed

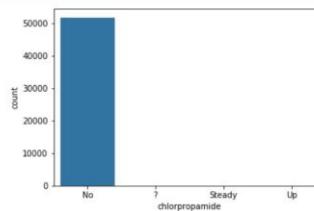
```



```

1 # chlorpropamide
2
3 # A diabetes medications
4 # The values of these variables indicate whether the drug was prescribed
5 # "up" if the dosage was increased during the encounter,
6 # "down" if the dosage was decreased,
7 # "steady" if the dosage did not change, and
8 # "no" if the drug was not prescribed.
9
10 chlorpropamide_count = sns.countplot(data=df, x='chlorpropamide')
11 plt.show()
12
13 # there is '?' values
14 # majority are not prescribed

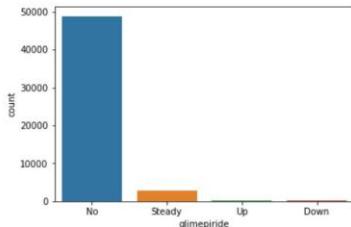
```



```

1 # glimepiride
2
3 # A diabetes medications
4 # The values of these variables indicate whether the drug was prescribed
5 # "up" if the dosage was increased during the encounter,
6 # "down" if the dosage was decreased,
7 # "steady" if the dosage did not change, and
8 # "no" if the drug was not prescribed.
9
10 glimepiride_count = sns.countplot(data=df, x='glimepiride')
11 plt.show()
12
13 # majority are not prescribed

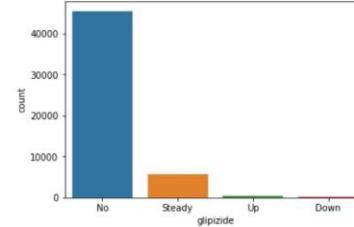
```



```

1 # glipizide
2
3 # A diabetes medications
4 # The values of these variables indicate whether the drug was prescribed
5 # "up" if the dosage was increased during the encounter,
6 # "down" if the dosage was decreased,
7 # "steady" if the dosage did not change, and
8 # "no" if the drug was not prescribed.
9
10 glipizide_count = sns.countplot(data=df, x='glipizide')
11 plt.show()
12
13 # majority are not prescribed

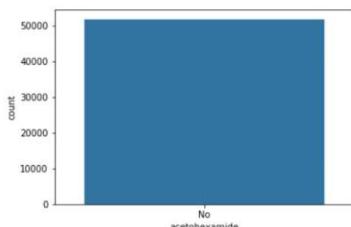
```



```

1 # acetohexamide
2
3 # A diabetes medications
4 # The values of these variables indicate whether the drug was prescribed
5 # "up" if the dosage was increased during the encounter,
6 # "down" if the dosage was decreased,
7 # "steady" if the dosage did not change, and
8 # "no" if the drug was not prescribed.
9
10 acetohexamide_count = sns.countplot(data=df, x='acetohexamide')
11 plt.show()
12
13 # NO patient is prescribed to it

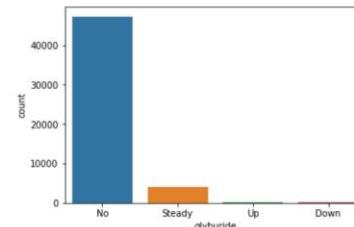
```



```

1 # glyburide
2
3 # A diabetes medications
4 # The values of these variables indicate whether the drug was prescribed
5 # "up" if the dosage was increased during the encounter,
6 # "down" if the dosage was decreased,
7 # "steady" if the dosage did not change, and
8 # "no" if the drug was not prescribed.
9
10 glyburide_count = sns.countplot(data=df, x='glyburide')
11 plt.show()
12
13 # majority are not prescribed

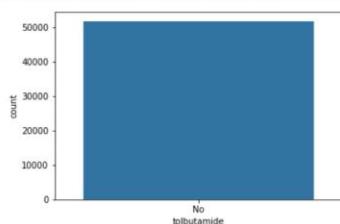
```



```

1 # tolbutamide
2
3 # A diabetes medications
4 # The values of these variables indicate whether the drug was prescribed
5 # "up" if the dosage was increased during the encounter,
6 # "down" if the dosage was decreased,
7 # "steady" if the dosage did not change, and
8 # "no" if the drug was not prescribed.
9
10 tolbutamide_count = sns.countplot(data=df, x='tolbutamide')
11 plt.show()
12
13 # NO patient is prescribed to it

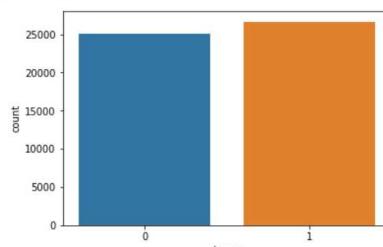
```



```

1 # Change of medications
2 # Indicates if there was a change in diabetic medications (either dosage or generic name)
3 # 1: "change" and 0: "no change"
4
5 change_count = sns.countplot(data=df, x='change')
6 plt.show()
7
8 # roughly similiar

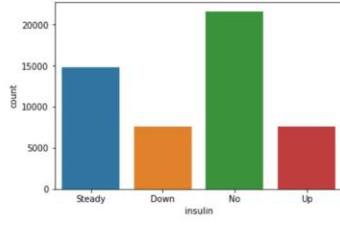
```



```

1 # insulin
2
3 # A diabetes medications
4 # The values of these variables indicate whether the drug was prescribed
5 # "up" if the dosage was increased during the encounter,
6 # "down" if the dosage was decreased,
7 # "steady" if the dosage did not change, and
8 # "no" if the drug was not prescribed.
9
10 insulin_count = sns.countplot(data=df, x='insulin')
11 plt.show()
12
13 # Good number of patients are on insulin

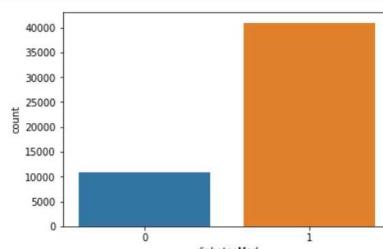
```



```

1 # Diabetes medications
2 # Indicates if there was any diabetic medication prescribed
3 # 1: "yes" and 0: "no"
4
5 diabetesMed_count = sns.countplot(data=df, x='diabetesMed')
6 plt.show()
7
8 # majority are not on diabetes medications

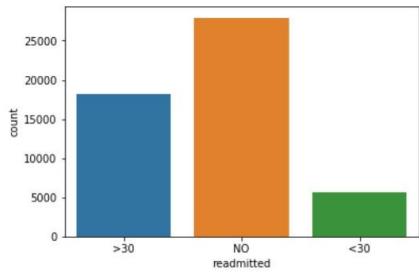
```



```

1 # Days to inpatient readmission.
2 # "<30" if the patient was readmitted in less than 30 days,
3 # ">30" if the patient was readmitted in more than 30 days, and
4 # "No" for no record of readmission.
5
6 readmitted_count = sns.countplot(data=df, x='readmitted')
7 plt.show()
8
9 # many are readmitted, majority of which are readmitted >30 days

```



- 2) Determine if there is any relationship between the variables num\\_medications and readmitted? How would you handle these two variables in the data modelling if a relationship exists?

Low correlation of 0.067557 existed between num\\_medications and readmitted. If high correlation do exist between features, it could cause multicollinearity and affect the model. Therefore, removing correlated features is recommended.

```

1 # Q3(2) Relationship between num_medications (numeric) and readmitted (categorical)
2 # Very low correlation of 0.067557
3
4 cat_map={">30":2, "<30":1, "NO":0}
5 df2=df
6 df2['readmitted']=df_new['readmitted'].map(cat_map)
7 df1=pd.DataFrame(df2, columns=['readmitted', 'num_medications'])
8 df1.corr(method='spearman')

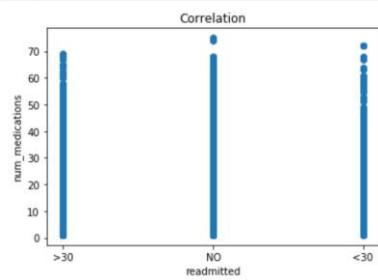
```

	readmitted	num_medications
readmitted	1.000000	0.067557
num_medications	0.067557	1.000000

```

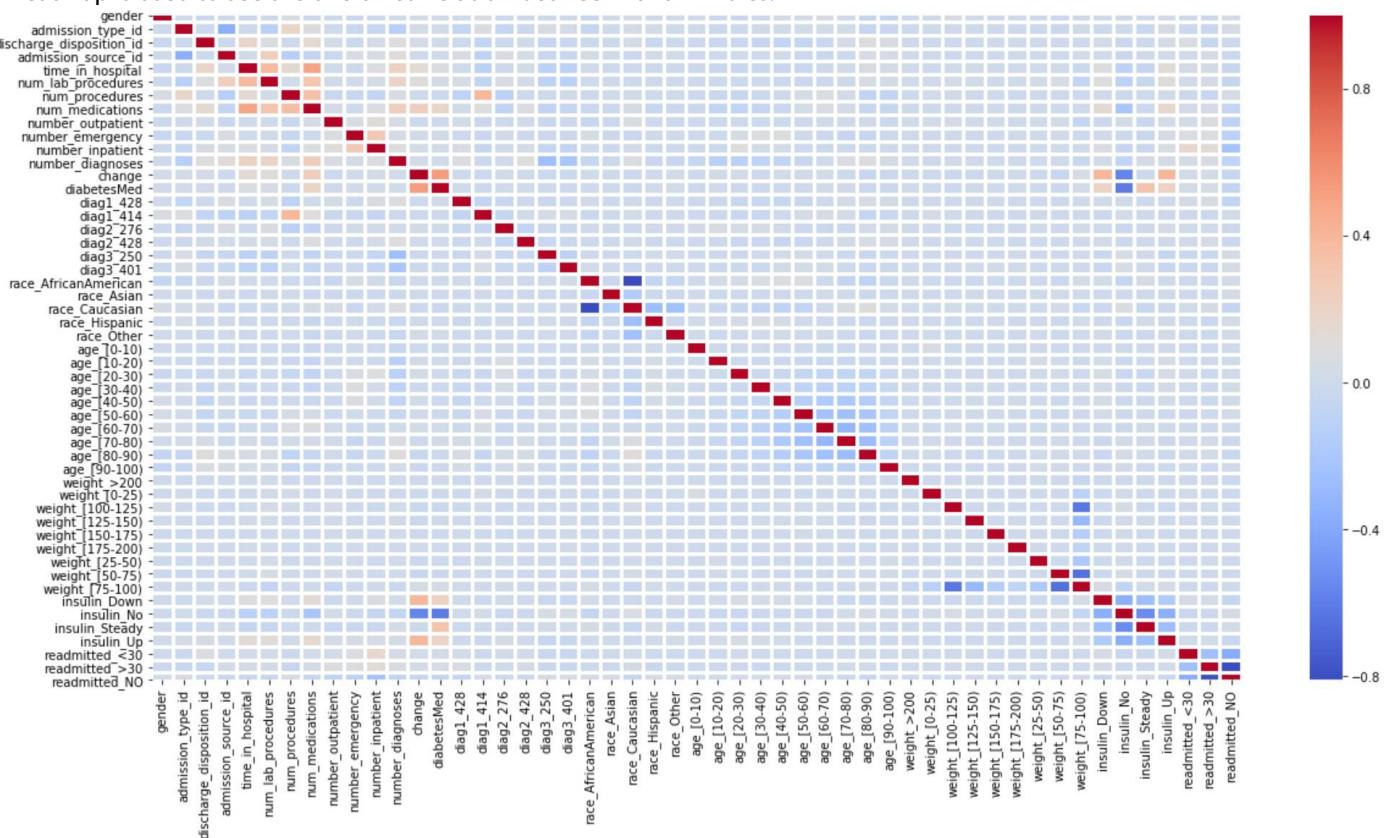
1 # Q3(2) Relationship between num_medications and readmitted
2 plt.scatter(df['readmitted'], df['num_medications'])
3 plt.title('Correlation')
4 plt.xlabel('readmitted')
5 plt.ylabel('num_medications')
6 plt.show()

```



- 3) Identify the highly correlated variable pairs and elaborate such a case.

A heatmap is used to see the overall correlation between variables:



Taking a deeper look at the variable with high correlation (the dark blue):

```
1 # What is the correlation b/w race_AfricanAmerican and race_Caucasian
2 race_corr=pd.DataFrame(df_dummy, columns=['race_AfricanAmerican', 'race_Caucasian'])
3 race_corr.corr(method='spearman')
4
```

	race_AfricanAmerican	race_Caucasian
race_AfricanAmerican	1.000000	-0.810038
race_Caucasian	-0.810038	1.000000

```
1 # What is the correlation b/w readmitted_>30 and readmitted_NO
2 readmitted_corr=pd.DataFrame(df_dummy, columns=['readmitted_>30', 'readmitted_NO'])
3 readmitted_corr.corr(method='spearman')
4
```

	readmitted_>30	readmitted_NO
readmitted_>30	1.000000	-0.797729
readmitted_NO	-0.797729	1.000000

1) race\_AfricanAmerican is highly negatively correlated with race\_Caucasian (-81%) ==> came from same variable "race"

2) readmitted\_>30 is highly negatively correlated with readmitted\_NO (-79.77%) ==> came from same variable "readmitted"

Will use race and readmitted as variables as a whole and not

on how to treat these variables in the mining process in

bles:

on how to treat these variables in the mining process in

**Question 4**

Summarization and data preparation:

- 1) Summarize your findings based on data exploration.

Error in data (values with '?')

- race - 1016 records (1% of data)
- age - 10 records
- weight - 50431 records (about 97% of total data)
- chlorpropamide - 9 records
- payer\_code - 7601 records (about 15% of total data)
- medical\_specialty - 32203 records (about 62% of total data)
- diag\_2 - 92 records (0.18% of total data)

NaN values

- gender 4 records - mapped from 'Unknown/Invalid':np.NaN

Outliers

They are not noise. They are anomalies. Do not drop.

- time\_in\_hospital - 977 records are  $>3\text{sd}$  - (1.89% of data) (positive skewed. 75% majority is below 6 days) (mean = 4.237, sd = 2.87, median = 3, min = 1, 1stQ = 2, 3rdQ = 6, max = 14)

```

1 # time in hospital
2 time_plot = sns.boxplot(df['time_in_hospital'])
3 plt.show()
4
5 # positive skewed. 75% majority is below 6 days. median of around 3 days.
6 # with outliers of 13 & 14 days
7
8 # 977 records are >3sd (1.89% of data)
9 time_in_hospital_outliers = (df['time_in_hospital'] > (4.237+2.870999*3)).value_counts()
10 time_in_hospital_outliers

```

A box plot for the 'time\_in\_hospital' column. The x-axis ranges from 2 to 14. The box represents the interquartile range (IQR) from approximately 2 to 6. The median is at 3. Whiskers extend to 1 and 12. Two outliers are marked with asterisks at 13 and 14.

```

False    50789
True     977
Name: time_in_hospital, dtype: int64

```

- num\_lab\_procedures - 20 records are  $>3\text{sd}$  - (median>mean, it is -ive skewed.) (mean = 43.853205, sd = 19.945471, median = 45, min = 1, 1stQ = 33, 3rdQ = 58, max = 132)

```

1 # number of lab procedures
2 lab_procedures_plot = sns.boxplot(df['num_lab_procedures'])
3 plt.show()
4
5 # per statistic median>mean, it is -ive skewed . With outliers.
6 # 20 records are >3sd
7 lab_procedures_outliers = (df['num_lab_procedures'] > (43.853205+19.945471*3)).value_counts()
8 lab_procedures_outliers

```

A box plot for the 'num\_lab\_procedures' column. The x-axis ranges from 0 to 120. The box represents the IQR from approximately 20 to 60. The median is at 45. Whiskers extend to 0 and 100. Numerous outliers are marked with asterisks between 100 and 130.

```

False    51746
True      20
Name: num_lab_procedures, dtype: int64

```

- num\_procedures - 2873 records are  $>3\text{sd}$  (5.55%) - = 1.308214, sd = 1.721694, median = 1, min = 0, 1stQ =

(positive skewed with outliers Median is 1 procedure) (mean 0, 3rdQ = 2, max = 6)

```

1 # number of procedures (other than lab procedures)
2 no_procedures_plot = sns.boxplot(df['num_procedures'])
3 plt.show()
4
5 # positive skewed with outlier. Median is 1.
6 # 2873 records are >3sd (5.55%)
7 num_procedures_outliers = (df['num_procedures'] > 6).value_counts()
8 num_procedures_outliers

```

Box plot of num\_procedures. The x-axis ranges from 0 to 6. The median is at 1. The box spans from approximately 0.5 to 2. The whiskers extend to 3 and 5. There is one outlier at 6.

```

False    48893
True     2873
Name: num_procedures, dtype: int64

```

- num\_medications - 626 records are  $>3\text{sd}$  (1.21% of data) - (positive skewed. with outliers Median of 15) (mean = 16.735850, sd = 8.063522, median = 16, min = 1, 1stQ = 11, 3rdQ = 21, max = 75)

```

[49]: M
1 # number of medications
2 no_medications_plot = sns.boxplot(df['num_medications'])
3 plt.show()
4
5 no_medications_countplot = sns.countplot(data=df, x='num_medications')
6 plt.show()
7
8 # positive skewed. Median of 15
9 # 626 records are >3sd (1.21% of data)
10 no_medications_outliers = (df['num_medications'] > (16.735850+8.063522*3)).value_counts()
11 no_medications_outliers

```

Box plot of num\_medications. The x-axis ranges from 0 to 70. The median is at 15. The box spans from approximately 10 to 20. The whiskers extend to 0 and 40. There are many outliers between 40 and 70. Below is a countplot of the same data.

Countplot of num\_medications. The x-axis ranges from 0 to 75. The y-axis is labeled 'count' and ranges from 0 to 3000. The distribution is highly right-skewed, peaking at 15 with a count of approximately 3000. Other counts are significantly lower, with a long tail extending towards 75.

```

Out[49]: False    51140
True      626
Name: num_medications, dtype: int64

```

- number\_outpatient - 737 records are  $>3\text{sd}$  (1.42% of data) - (positive skewed. Majority is 0) (mean = 0.512846, sd = 1.544866, median = 0, min = 0, 1stQ = 0, 3rdQ = 0, max = 42)

```

1 # Number of outpatient visits of the patient in the year preceding the encounter
2 no_outpatient_plot = sns.boxplot(df['number_outpatient'])
3 plt.show()
4
5 no_outpatient_count = sns.countplot(data=df, x='number_outpatient')
6 plt.show()
7
8 # positive skewed. Majority is 0. With outliers
9 # 737 records are >3sd (1.42% of data)
10 number_outpatient_outliers = (df['number_outpatient'] > (0.512846+1.544866*3)).value_counts()
11 number_outpatient_outliers

```

Box plot of number\_outpatient. The x-axis ranges from 0 to 40. The median is at 0. The box spans from approximately 0 to 10. The whiskers extend to 0 and 30. There are many outliers between 30 and 40. Below is a countplot of the same data.

Countplot of number\_outpatient. The x-axis ranges from 0 to 40. The y-axis is labeled 'count' and ranges from 0 to 40000. The distribution is highly right-skewed, with the vast majority of patients having 0 outpatient visits (count ~50000). A few individuals have between 10 and 40 visits, with counts decreasing as the number of visits increases.

```

False    51029
True      737
Name: number_outpatient, dtype: int64

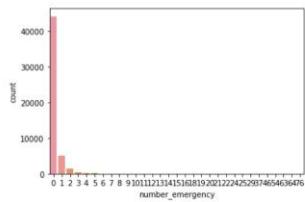
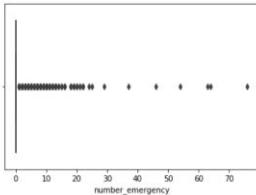
```

- number\_emergency - 667 records are  $>3\text{sd}$  (1.29% of data) - (positive skewed. Majority is 0) (mean = 0.269173,  $\text{sd} = 1.146923$ , median = 0, min = 0, 1stQ = 0, 3rdQ = 0, max = 76)

```

1 # Number of emergency visits of the patient in the year preceding the encounter
2 no_outpatient_plot = sns.boxplot(df['number_emergency'])
3 plt.show()
4
5 no_outpatient_count = sns.countplot(data=df, x='number_emergency')
6 plt.show()
7
8 # positive skewed. Majority is 0 with outliers
9 # 667 records are >3sd (1.29% of data)
10 number_emergency_outliers = (df['number_emergency'] > (0.269173+1.146923*3)).value_counts()
11 number_emergency_outliers

```



```

False      51099
True       667
Name: number_emergency, dtype: int64

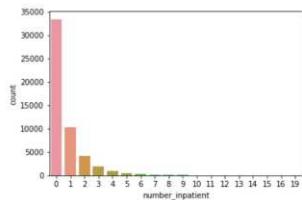
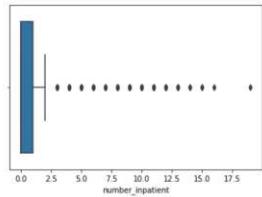
```

- number\_inpatient - 1146 records are  $>3\text{sd}$  (2.21% of data) - (positive skewed. Majority is 0) (mean = 0.686879,  $\text{sd} = 1.320968$ , median = 0, min = 0, 1stQ = 0, 3rdQ = 1, max = 19)

```

1 # Number of inpatient visits of the patient in the year preceding the encounter
2 no_outpatient_plot = sns.boxplot(df['number_inpatient'])
3 plt.show()
4
5 no_outpatient_count = sns.countplot(data=df, x='number_inpatient')
6 plt.show()
7
8 # positive skewed. Majority is 0 with outliers
9 # 1146 records are >3sd (2.21% of data)
10 number_inpatient_outliers = (df['number_inpatient'] > (0.686879+1.320968*3)).value_counts()
11 number_inpatient_outliers

```

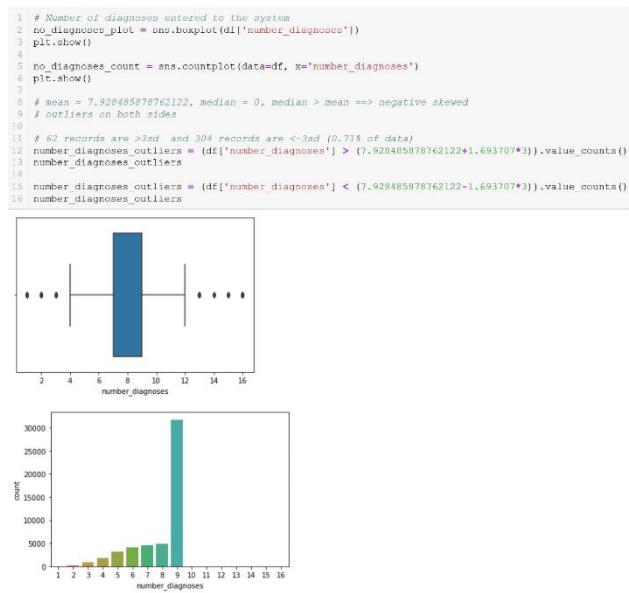


```

False      50620
True       1146
Name: number_inpatient, dtype: int64

```

- number\_diagnoses - 62 records are  $>3\text{sd}$  and 304 records are  $<-3\text{sd}$  (0.71% of data) a) (median > mean ==> negative skewed) (mean = 7.928486, sd = 1.693707, median = 9, min = 1, 1stQ = 7, 3rdQ = , max = 16)



- 2) Elaborate on the data preparation steps required (by correcting and transformation) to address the data quality problems that you encountered during data exploration.

In general, the following can be done when errors are found:

- ignore
- treat errors as separate value
- impute with mean/median (numeric) or mode (categorical)
- other techniques - eg. regression, interpolation
- flagging

In addition, there are other considerations should be taken into account on whether variables should be dropped. They include:

- Unary variable - A variable that takes only one unique value. Since all the values are same, it will not change the outcome of any data mining model.
- Redundant variables - Same information represented in two or more different ways. All of them except one should be removed as no new information is provided by such variables.
- Derived variable - A variable that is derived from another variable. In most cases, it is not necessary to use the derived variable in the analysis.

Below the analysis on data preparations steps:

#### Error in data (values with '?')

- race - 1016 records - difficult to approximate race. 1016 records is about 1% of data. Will keep them.  
Map '?' to np.NaN

#### Weight

- weight – although there are 50431 records with '?' representing about 97% of total data. Nevertheless, it is an important indicator in diabetes patients. Hence, will impute the missing values using mode.

Drop column

- encounter\_id – nominal - not useful for analysis
- patient\_nbr – nominal - not useful for analysis
- payer\_code - 7601 records - (about 15% of total data) – not useful for the analysis
- medical\_specialty - 32203 records - (about 62% of total data) – not useful for analysis
- max\_glu\_serum (most value is "none" i.e. test not taken)
- A1Cresult (most value is "none" i.e. not measured)

Drop row

- age with '?' - 10 records – few records
- chlorpropamide with '?' - 9 records – few records

Diabetes Medicine

As majority of the below diabetes medicines are not prescribed plus it already has a variable on diabetesMed indicating if taking any diabetes medication. Hence, the following are drop due to redundancies and to reduce dimensionality.

- metformin (diabetes medicine)
- repaglinide (diabetes medicine)
- nateglinide (diabetes medicine)
- chlorpropamide (diabetes medicine)
- glimepiride (diabetes medicine)
- acetoexamide (diabetes medicine)
- glipizide (diabetes medicine)
- glyburide (diabetes medicine)
- tolbutamide (diabetes medicine)

diag\_1, diag\_2, and diag\_3

```
1 # check to see how many unique values in diag_1
2 diag1_unique = len(pd.unique(df_process['diag_1']))
3 diag1_unique
```

649

```
1 # check to see how many unique values in diag_2
2 diag2_unique = len(pd.unique(df_process['diag_2']))
3 diag2_unique
```

683

```
1 # check to see how many unique values in diag_3
2 diag3_unique = len(pd.unique(df_process['diag_3']))
3 diag3_unique
```

726

- diag\_1 – there are 649 different diagnoses in diag\_1. Need to reduce dimension in the analysis. As a start, will create derived variables on the top 2 diagnosis. '428' diagnosis (Congestive heart failure, unspecified) and '414' diagnosis (Other forms of chronic ischemic heart disease) to see if these add predictive power of the model. Other diagnoses are drop.
- diag\_2 - there are 683 different diagnoses in diag\_2. Need to reduce dimension in the analysis. As a start, will create derived variables on the top 2 diagnosis. '276' diagnosis (Hyperosmolality and/or hypernatremia) and '428' diagnosis (Congestive heart failure, unspecified) to see if these add predictive power of the model. Other diagnoses are drop.
- diag\_3 - there are 726 different diagnoses in diag\_2. Need to reduce dimension in the analysis. As a start, will create derived variables on the top 2 diagnosis. '250' diagnosis (Diabetes) and '401' diagnosis (Malignant essential hypertension) to see if these add predictive power of the model. Other diagnoses are drop.

- 3) Demonstrate the data preparation by including a screenshot(s) of the Python code and its outputs that show the steps on how you had corrected all the identified data quality problems in this dataset.

```

1 # create derived variables for diag_1, diag_2, diag_3 for their top two diagnoses
2
3 # For diag_1
4 # create new column for diag_1 with '428' diagnosis (Congestive heart failure, unspecified)
5 df_process['diag1_428'] = df_process['diag_1'].apply(lambda x: 1 if x == '428' else 0)
6
7 # create new column for diag_1 with '414' diagnosis(Other forms of chronic ischemic heart disease)
8 df_process['diag1_414'] = df_process['diag_1'].apply(lambda x: 1 if x == '414' else 0)
9
10
11 # For diag_2
12 # create new column for diag_2 with '276' diagnosis(Hyperosmolality and/or hypernatremia)
13 df_process['diag2_276'] = df_process['diag_2'].apply(lambda x: 1 if x == '276' else 0)
14
15 # create new column for diag_2 with '428' diagnosis(Congestive heart failure, unspecified)
16 df_process['diag2_428'] = df_process['diag_2'].apply(lambda x: 1 if x == '428' else 0)
17
18
19
20 # For diag_3
21 # create new column for diag_3 with '250' diagnosis(Diabetes)
22 df_process['diag3_250'] = df_process['diag_3'].apply(lambda x: 1 if x == '250' else 0)
23
24 # create new column for diag_3 with '401' diagnosis(Malignant essential hypertension)
25 df_process['diag3_401'] = df_process['diag_3'].apply(lambda x: 1 if x == '401' else 0)
26
27

```

```

1 ### Pre-process data
2
3
4 # Drop columns : encounter_id, patient_nbr, payer_code, medical_specialty
5 df_process = df_process.drop(['encounter_id', 'patient_nbr', 'payer_code',
6 'medical_specialty'], axis=1)
7
8 # Drop rows : age with ? (10 records)
9 df_process = df_process[df_process.age != '?']
10
11 # Drop rows : chlorpropamide (9 records)
12 df_process = df_process[df_process.chlorpropamide != '?']
13
14 # Drop rows : gender (4 records) (NaN)
15 df_process = df_process.dropna(how='any', subset=['gender'])
16
17 # Race - replace "?" to nan
18 df_process['race'] = df_process['race'].replace(['?'], np.nan)
19
20 # Impute missing weight values using mode
21 df_process['weight'] = df_process['weight'].replace(['?'], np.nan)
22 stats.mode(df_process['weight'])
23 # replace NaN weight with mode '[75-100]'
24 df_process['weight'] = df_process['weight'].replace(np.nan, '[75-100]')
25
26 # Dropping other variables in df_process
27 df_process = df_process.drop(['diag_1', 'diag_2', 'diag_3',
28 'max_glu_serum', 'A1Cresult',
29 'metformin', 'repaglinide', 'nateglinide',
30 'chlorpropamide', 'glimepiride', 'acetohexamide', 'glipizide',
31 'glyburide', 'tolbutamide'], axis=1)
32

```

	race	gender	age	weight	admission_type_id	discharge_disposition_id	admission_source_id	time_in_hospital	num_lab_procedures	num_procedures	...	insulin	change	diabetesMed	readmitted	diag1_428	diag1_414	diag2_276	diag2_428	diag3_250	diag3_401
0	Caucasian	0.0	[60-70)	[75-100)	2	1	1	4	43	0	...	Steady	0	1	>30	1	0	0	0	0	0
1	Caucasian	0.0	[80-90)	[75-100)	3	2	1	8	48	6	...	Down	1	1	>30	0	1	0	0	0	0
2	Caucasian	0.0	[80-90)	[75-100)	2	1	1	2	39	0	...	No	0	0	NO	0	0	0	0	0	1
3	Caucasian	0.0	[60-70)	[75-100)	3	1	1	2	54	0	...	No	0	0	NO	0	0	0	0	0	0
4	Caucasian	0.0	[60-70)	[75-100)	3	2	1	7	70	1	...	Down	1	1	>30	0	0	0	0	0	0
5	Caucasian	1.0	[70-80)	[75-100)	2	1	7	3	1	4	...	Steady	0	1	NO	0	0	0	0	0	0
6	Caucasian	1.0	[60-70)	[75-100)	2	1	1	6	28	1	...	Up	1	1	NO	0	0	0	0	0	1
7	AfricanAmerican	0.0	[40-50)	[75-100)	1	1	7	3	48	0	...	Down	1	1	NO	0	0	0	0	1	0
8	AfricanAmerican	1.0	[70-80)	[75-100)	1	1	7	5	42	2	...	No	0	1	<30	0	0	0	0	0	0
9	AfricanAmerican	0.0	[70-80)	[75-100)	1	1	7	1	43	1	...	Up	1	1	>30	0	0	0	0	0	0
10	Caucasian	0.0	[70-80)	[75-100)	3	1	1	3	27	6	...	Steady	1	1	NO	0	1	0	0	0	0

**Question 5**

Selection of data mining task and feature selection:

- 1) Identify the most suitable data mining task (i.e. classification, clustering or association mining) that can be performed on this dataset. Justify your choice.

Assuming we would like to predict patients who will be readmitted to the hospital in less than 30 days. The most suitable technique for this data mining task is classification, as it is based on historical labelled data and can model to make prediction on the target variable (ie. readmitted<30 days). As target class is already known hence training data is already labelled with answers.

Classification data mining requires sufficient size of records to be met with about 51700 records. The dataset required to be able to train the model, in this case, the criteria can divided into training, validation and testing sub-sets in order to create the classification model.

- 2) What variables will you include in this data mining task and why? Describe here if you will create any derived variables. Identify the roles (i.e., input or target) of each variable.

With the assumption that we would like to predict patients who will be readmitted to the hospital within 30 days.

Target : readmission

Input :

1. gender
2. admission\_type\_id
3. discharge\_disposition\_id
4. admission\_source\_id
5. time\_in\_hospital
6. num\_lab\_procedures
7. num\_medications
8. number\_outpatient
9. number\_diagnoses
10. change
11. diabetesMed
12. diag1\_428
13. diag1\_414
14. diag2\_276
15. diag2\_428
16. diag3\_250
17. diag3\_401
18. race
19. age
20. weight
21. insulin

Will use the above variables to build an initial model, then by process of elimination using the features selection techniques (eg. Univariate filter method, multivariate filter method, wrapper methods etc..) to select features to reduce overfitting.

Normally, derived variables are not necessary as it is highly correlated with the 'parent' variables. In this case, we will use the derived variables of the top two diagnoses in diag\_1, diag\_2 and diag\_3 for the model to see if it helps with prediction. At the same time, due to dimensionality reduction, diag\_1, diag\_2 and diag\_3 variables are dropped.