

## Question

1 This work intend to answering the question “How well do DRG-predicted LOS values match actual values? Do any factors explain significant variability in the closeness of these values?”

### Import module

```
In [1]: 1 import pprint
2 import sqlite3
3 import pathlib
4
5 import pandas as pd
6 import numpy as np
7 import seaborn as sns
8 sns.set_style('whitegrid')
9 import matplotlib as plt
10 from scipy import stats
11
12 import matplotlib.pyplot as plt
13 plt.style.use('default')
14 plt.style.use ('seaborn-darkgrid')
15
16 import plotly.express as px
17 from statsmodels.graphics.gofplots import qqplot
18 from statsmodels.stats.weightstats import ztest
19
20
21 pd.set_option('display.max_rows', 500)
22 pd.set_option('display.max_columns', 500)
```

### Read from database

```
In [2]: 1 def query_sqlite_with_pandas(query):
2
3     #with sqlite3.Connection('LOS_data_cd41.sqlite') as conn:
4     with sqlite3.Connection('LOS_data_2023-02-01.sqlite') as conn:
5         query_results = pd.read_sql(query, conn)
6         conn.close()
7
8     return query_results
```

### Query the variables of interest

```
In [3]: 1 query = """
2 SELECT e.patient_encounter_tk, e.admission_datetime, e.discharge_datetime,           --Encounter information
3         pe.drg_code, d.drg_description, d.gmlos, d.clinical_type, d.mdc,             --DRG information
4         p.patient_tk, p.DOB, g.gender_abbreviation, p.ethnicity_tk, p.race_tk,       --Patient descriptives
5         e.last_inpatient_admission_datetime, e.patient_class_tk,                  --Encounter descriptives
6         f.facility_tk, f.hospital_system_tk, f.analytics_display_name              --Facility information
7
8 FROM cd_patient_encounters e
9     LEFT JOIN cd_patient_encounter_drugs pe ON e.patient_encounter_tk=pe.patient_encounter_tk
10    LEFT JOIN ref_drg_code_history d ON pe.drg_code=d.drg_code
11    LEFT JOIN cd_patients p ON e.patient_tk=p.patient_tk
12    LEFT JOIN ref_facilities f ON e.facility_tk=f.facility_tk
13    LEFT JOIN cd_genders g2 ON p.gender_tk=g2.gender_tk and f.facility_tk=g2.facility_tk
14    LEFT JOIN ref_genders g ON g2.master_gender_tk=g.master_gender_tk;
15 """
16
17 los_calc_df = query_sqlite_with_pandas(query)
```

```
In [4]: 1 ##### Make datetime datatype
2 los_calc_df['admission_datetime'] = pd.to_datetime(los_calc_df['admission_datetime'])
3 los_calc_df['discharge_datetime'] = pd.to_datetime(los_calc_df['discharge_datetime'],errors='coerce')
4 los_calc_df['last_inpatient_admission_datetime'] = pd.to_datetime(los_calc_df['last_inpatient_admission_datetime'],errors='coerce')
5 los_calc_df['dob'] = pd.to_datetime(los_calc_df['dob'])
```

## Initial Data Exploration and Understanding and Data Transformstion

### Create New Fields

Length of stay LOS is the actual duration of time used by a recorder patient encounter. This is calculated by subtracting the discharge time from admission time

```
In [5]: 1 los_calc_df['encounter_los'] = los_calc_df['discharge_datetime']-los_calc_df['admission_datetime']
```

```
In [6]: 1 los_calc_df['encounter_los'] = los_calc_df['encounter_los'] / pd.Timedelta(days=1)
```

```
In [293]: 1 #check for unique values in gender column
2 los_calc_df.gender_abbreviation.unique()
```

```
Out[293]: array(['M', 'F', 'U'], dtype=object)
```

#### Difference between DRG LOS and ENC LOS, Meet\_it criteria and return\_visit

LOS Difference is the calculated value as a difference between encounter\_los from the expected or given Geometric Mean Length of stay GMLOS. There is the MDC that group DRGs into different categories. The DRGs are mapped to MDC using a file called MDC mapper. Return flag is to code any patient that has been to the hospital system before i.e a returning patient. There are instances where admission date is less than date of birth this is specific for foetus who are not yet born as at admission. There is a meet\_it binary criteria that shows for each patient encounter if they exceed the GMLOS or not.

```
In [8]: 1 #calculate difference
2 los_calc_df['los_difference'] = los_calc_df['encounter_los']-los_calc_df['gmlos']
```

```
In [9]: 1 ### MDC Description
2 mdc_mapper = pd.read_csv("mdc_mapper.csv", dtype='object').set_index('mdc')
3 los_calc_df = los_calc_df.join(mdc_mapper, on='mdc')
4 los_calc_df['mdc_description'] = los_calc_df.mdc_description.fillna('')
```

```
In [10]: 1 # Return Visit Flag
2 los_calc_df['return_visit'] = 1*(~los_calc_df['last_inpatient_admission_datetime'].isna())
```

```
In [11]: 1 # patient : foetus i.e having DOB to be less than admission date
2 (los_calc_df['admission_datetime']<los_calc_df['dob']).any()
```

```
Out[11]: True
```

```
In [ ]: 1 los_calc_df[(los_calc_df['admission_datetime']<los_calc_df['dob'])]
```

```
In [13]: 1 los_calc_df.dtypes
```

```
Out[13]: patient_encounter_tk          int64
admission_datetime          datetime64[ns]
discharge_datetime          datetime64[ns]
drg_code                    object
drg_description              object
gmlos                       float64
clinical_type                object
mdc                         object
patient_tk                  int64
dob                         datetime64[ns]
gender_abbreviation          object
ethnicity_tk                int64
race_tk                     int64
last_inpatient_admission_datetime  datetime64[ns]
patient_class_tk            int64
facility_tk                  int64
hospital_system_tk          int64
analytics_display_name       object
encounter_los                float64
los_difference               float64
mdc_description              object
return_visit                 int32
dtype: object
```

```
In [14]: 1 los_calc_df.nunique()
```

```
Out[14]: patient_encounter_tk      160205
admission_datetime      133779
discharge_datetime      101921
drg_code                 759
drg_description         740
gmlos                   106
clinical_type            2
mdc                      27
patient_tk              130242
dob                     29487
gender_abbreviation      3
ethnicity_tk             23
race_tk                  240
last_inpatient_admission_datetime  3360
patient_class_tk         76
facility_tk               31
hospital_system_tk       4
analytics_display_name   31
encounter_los            29220
los_difference           52522
mdc_description          26
return_visit             2
dtype: int64
```

```
In [15]: 1 # Age at Admission
2 los_calc_df['age_at_admit'] = los_calc_df['admission_datetime']-los_calc_df['dob']
3 los_calc_df['age_at_admit'] = los_calc_df['age_at_admit'].astype('timedelta64[Y]').astype(int)
```

```
In [16]: 1 # meet_it (Is gmlos met?)
2 los_calc_df['meet_it'] = 1*(los_calc_df['encounter_los'] <= los_calc_df['gmlos'])
```

### Data Cleaning

```
In [17]: 1 los_calc_df[los_calc_df['encounter_los'] == 0].shape
2 #446 rows with exact same admit and discharge time
```

```
Out[17]: (446, 24)
```

```
In [18]: 1 encounter_records = los_calc_df[los_calc_df['encounter_los'] > 0]
2 #1279 rows have discharge at or before admit - these are dropped
3 #also drops 497 n/a discharges
```

```
In [19]: 1 encounter_records = encounter_records[encounter_records['encounter_los'] <= 180]
2 #drop 22 extreme outliers
```

```
In [20]: 1 encounter_records[encounter_records['gmlos'].isna()][['drg_code']].unique()
2 #8880 records in 18 drg codes do not map to a gmlos
```

```
Out[20]: array(['767', '', '765', '775', '133', '766', '774', '777', '238', '131',
              '484', '782', '129', '490', '778', '130', '781', '998', '132'],
              dtype=object)
```

```
In [21]: 1 encounter_records[encounter_records['drg_code']==''].shape
2 #8615 of those records have a blank drg
```

```
Out[21]: (8615, 24)
```

```
In [22]: 1 encounter_records_final = encounter_records[~encounter_records['encounter_los'].isna()]
```

### Create DRG Grouped Data

```
In [23]: 1 agg_funcs = {'encounter_los': stats.gmean,
2                 'age_at_admit': 'mean',
3                 'return_visit': 'mean',
4                 'meet_it': 'mean',
5                 'patient_encounter_tk': 'count'}
6 group_by_columns = ['drg_code', 'drg_description', 'mdc', 'mdc_description', 'clinical_type', 'gmlos'], #,
```

```
In [24]: 1 drg_records_prelim = encounter_records_final.groupby(group_by_columns).agg(agg_funcs).reset_index()
```

```
In [25]: 1 rename_cols = {"encounter_los": "actual_gmlos",
2                  "age_at_admit": "average_age",
3                  "return_visit": "pct_return_visits",
4                  "meet_it": "pct_meet_it",
5                  "patient_encounter_tk": "num_encounters"}
```

```
In [40]: 1 drg_records = drg_records_prelim.dropna().rename(columns=rename_cols)
2 drg_records['gmlos_difference'] = drg_records.eval("actual_gmlos-gmlos")
3 drg_records['drg_meet_it'] = 1*(drg_records['actual_gmlos'] <= drg_records['gmlos'])
4 drg_records.head()
5 drg_records_final = drg_records.copy()
```

```
In [ ]: 1
```

```
In [43]: 1 #select columns of interest
2 column = ['drg_code', 'mdc', 'clinical_type', 'average_age', 'pct_return_visits', 'num_encounters',
3          'gmlos', 'actual_gmlos', 'gmlos_difference', 'drg_meet_it']
4 df2=drg_records_final[column]
```

```
In [45]: 1 df2
```

```
Out[45]:
```

	drg_code	mdc	clinical_type	average_age	pct_return_visits	num_encounters	gmlos	actual_gmlos	gmlos_difference	drg_meet_it
0	001	PRE	SURG	60.571429	0.190476	21	30.1	20.212580	-9.887420	1
1	003	PRE	SURG	55.146341	0.012195	164	22.4	18.259437	-4.140563	1
2	004	PRE	SURG	60.181159	0.028986	138	20.0	22.124499	2.124499	0
3	011	PRE	SURG	72.421053	0.000000	19	10.9	0.771884	-10.128116	1
4	012	PRE	SURG	64.666667	0.000000	9	8.3	4.845043	-3.454957	1
...	...	...	...	...	...	...	...	...	...	...
734	982		SURG	60.717791	0.036810	163	4.6	4.288287	-0.311713	1
735	983		SURG	60.648649	0.000000	37	2.3	2.085249	-0.214751	1
736	987		SURG	65.827273	0.036364	110	7.7	6.624850	-1.075150	1
737	988		SURG	63.370690	0.017241	116	4.3	3.754208	-0.545792	1
738	989		SURG	56.647059	0.000000	17	2.3	2.481734	0.181734	0

739 rows × 10 columns

```
In [ ]: 1
```

```
In [ ]: 1
```

```
In [46]: 1 #select columns of interest
2 columns = ['drg_code', 'mdc', 'gender_abbreviation', 'clinical_type', 'age_at_admit', 'return_visit', 'gmlos', 'encounter_los']
```

```
In [47]: 1 encounter_records[columns].isna().sum()
```

```
Out[47]: drg_code      0
mdc      8880
gender_abbreviation  0
clinical_type  8880
age_at_admit    0
return_visit    0
gmlos      8880
encounter_los    0
los_difference  8880
meet_it        0
dtype: int64
```

```
In [312]: 1 # drop all rows that have no value, not a number
          2 df1=encounter_records[columns].dropna().reset_index(drop='Index')
          3 df1
```

Out[312]:

	drg_code	mdc	gender_abbreviation	clinical_type	age_at_admit	return_visit	gmlos	encounter_los	los_difference	meet_it
0	872	18	M	MED	55	0	3.5	4.162500	0.662500	0
1	441	07	F	MED	42	0	4.7	2.091667	-2.608333	1
2	871	18	F	MED	46	0	4.8	7.886806	3.086806	0
3	442	07	F	MED	47	0	3.2	1.006250	-2.193750	1
4	708	12	M	SURG	51	0	1.4	1.203472	-0.196528	1
...	...	...	...	...	...	...	...	...	...	...
149520	794	15	U	MED	0	0	3.4	1.873681	-1.526319	1
149521	026	01	F	SURG	64	0	3.8	1.264468	-2.535532	1
149522	787	14	F	SURG	37	0	3.5	2.150914	-1.349086	1
149523	982		F	SURG	35	0	4.6	1.933727	-2.666273	1
149524	795	15	M	MED	0	0	3.1	2.054375	-1.045625	1

149525 rows × 10 columns

```
In [313]: 1 df1.rename(columns={'gender_abbreviation':'gender','age_at_admit':'age','mdc':'mdc','encounter_los':'LOS','los_difference'
          2 df1
```

Out[313]:

	drg_code	mdc	gender	clinical_type	age	return_visit	gmlos	LOS	LOS_Diff	meet_it
0	872	18	M	MED	55	0	3.5	4.162500	0.662500	0
1	441	07	F	MED	42	0	4.7	2.091667	-2.608333	1
2	871	18	F	MED	46	0	4.8	7.886806	3.086806	0
3	442	07	F	MED	47	0	3.2	1.006250	-2.193750	1
4	708	12	M	SURG	51	0	1.4	1.203472	-0.196528	1
...	...	...	...	...	...	...	...	...	...	...
149520	794	15	U	MED	0	0	3.4	1.873681	-1.526319	1
149521	026	01	F	SURG	64	0	3.8	1.264468	-2.535532	1
149522	787	14	F	SURG	37	0	3.5	2.150914	-1.349086	1
149523	982		F	SURG	35	0	4.6	1.933727	-2.666273	1
149524	795	15	M	MED	0	0	3.1	2.054375	-1.045625	1

149525 rows × 10 columns

In [ ]:

1

## Business Intelligence

```
1 ### How well does GMLOS predict length of stay for DRG encounters?•
2 We like to understand how the given GMLON accurately explains the actual variation in hospital patients encounter lenght
  of stay (Actual_GMLOS). We plot a linear fit of the two using a scatter plot and a line of fit
```

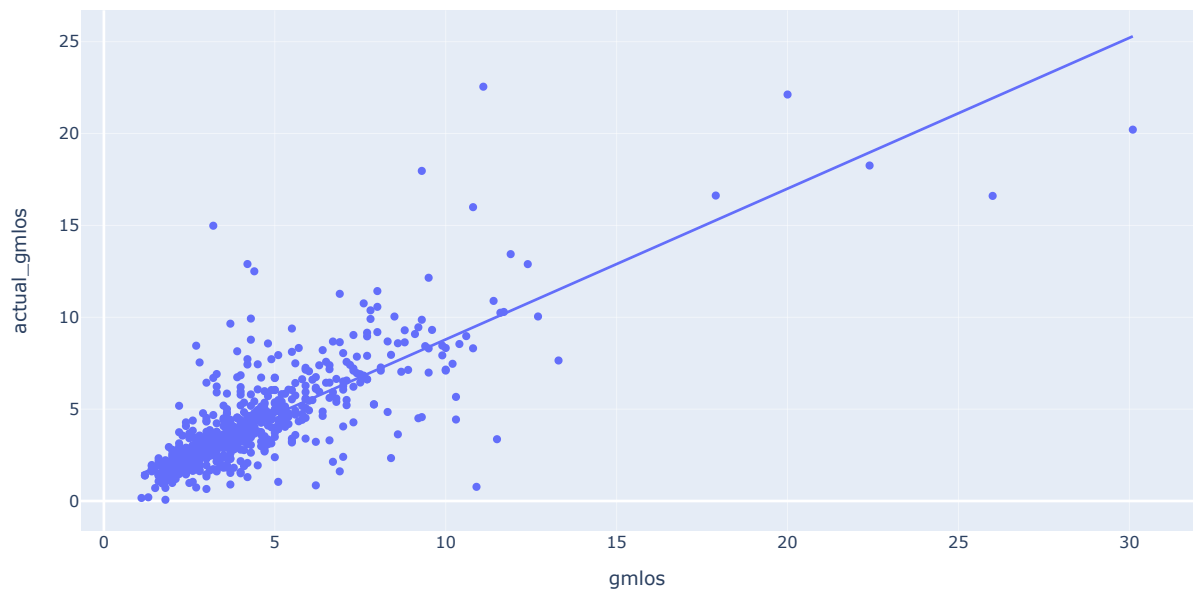
In [ ]:

1

In [177]:

```
1 def scatter(df,x,y,title):
2     '''scatter plot to plot a linear line of best fit between two continuos values x and y'''
3     fig = px.scatter(df, x=x, y=y, trendline='ols')
4     # Set plot title and center-align it
5     fig.update_layout(
6         title=title,
7         title_x=0.5 # Value of 0.5 centers the title
8     )
9     fig.show()
10 scatter(df=df2,x='gmlos',y='actual_gmlos',title='Scatter Plot of actual_GMLOS vs GMLOS')
```

Scatter Plot of actual\_GMLOS vs GMLOS



**Difference between Encounter LOS(Patient's Actual GMLOS) and given Geometric Mean LOS**

**At DRG level**

A histogram showing the distribution of Grouped DRGs and showing how many DRGs meet the government GMLOS. We expect data point to be at the mean of zero , slightly left skewed and a decent kurtois. Against our thought there are lots of DRGs that do not meet the GMLOS as the distribution is slightly more skewed to the right.

In [96]:

```
1 #aggregated dataframe used (grouped DRGs)
2 df2
```

Out[96]:

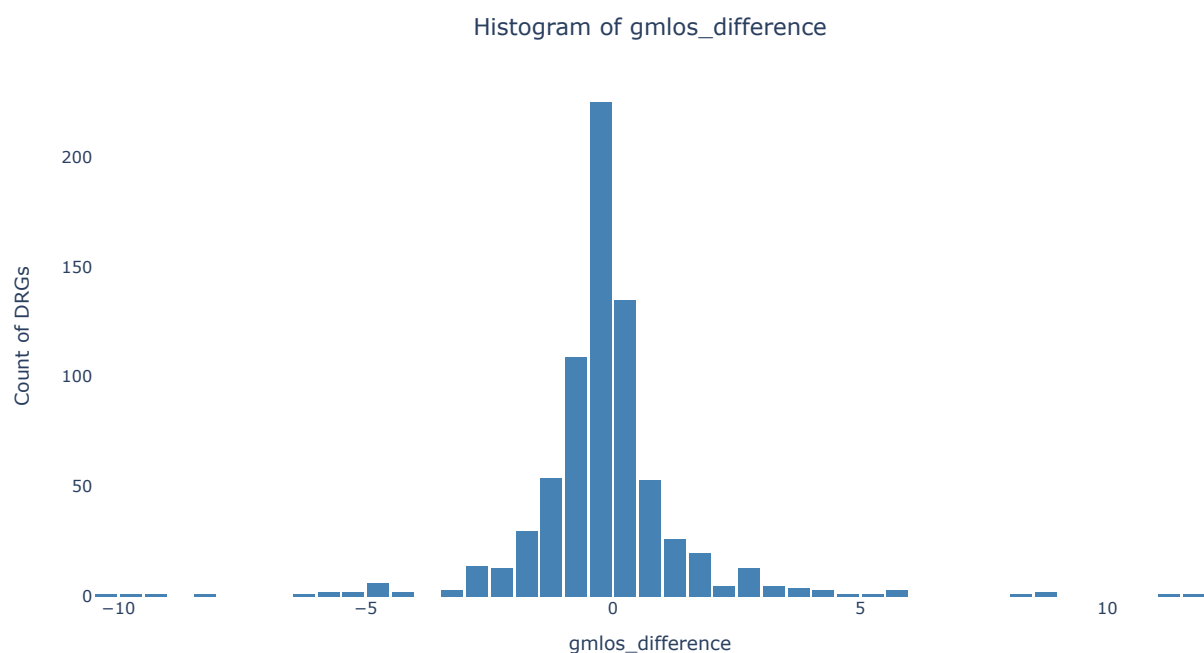
	drg_code	mdc	clinical_type	average_age	pct_return_visits	num_encounters	gmlos	actual_gmlos	gmlos_difference	drg_meet_it
0	001	PRE	SURG	60.571429	0.190476	21	30.1	20.212580	-9.887420	1
1	003	PRE	SURG	55.146341	0.012195	164	22.4	18.259437	-4.140563	1
2	004	PRE	SURG	60.181159	0.028986	138	20.0	22.124499	2.124499	0
3	011	PRE	SURG	72.421053	0.000000	19	10.9	0.771884	-10.128116	1
4	012	PRE	SURG	64.666667	0.000000	9	8.3	4.845043	-3.454957	1
...	...	...	...	...	...	...	...	...	...	...
734	982		SURG	60.717791	0.036810	163	4.6	4.288287	-0.311713	1
735	983		SURG	60.648649	0.000000	37	2.3	2.085249	-0.214751	1
736	987		SURG	65.827273	0.036364	110	7.7	6.624850	-1.075150	1
737	988		SURG	63.370690	0.017241	116	4.3	3.754208	-0.545792	1
738	989		SURG	56.647059	0.000000	17	2.3	2.481734	0.181734	0

739 rows × 10 columns

```

In [52]: 1 fig = px.histogram(df2, x='gmlos_difference')
2
3 # Set plot title and center-align it
4 fig.update_layout(
5     title='Histogram of gmlos_difference',
6     title_x=0.5 # Value of 0.5 centers the title
7 )
8
9 # Customize the bar colors
10 fig.update_traces(marker_color='steelblue')
11
12 # Set x-axis Label
13 fig.update_xaxes(title='gmlos_difference')
14
15 # Set y-axis Label
16 fig.update_yaxes(title='Count of DRGs')
17
18 # Customize the Layout
19 fig.update_layout(
20     plot_bgcolor='white', # Set plot background color
21     barmode='overlay', # Display bars in overlay mode
22     bargap=0.1, # Set gap between bars
23     bargroupgap=0.05, # Set gap between bar groups
24 )
25
26 fig.show()

```



In [ ]:

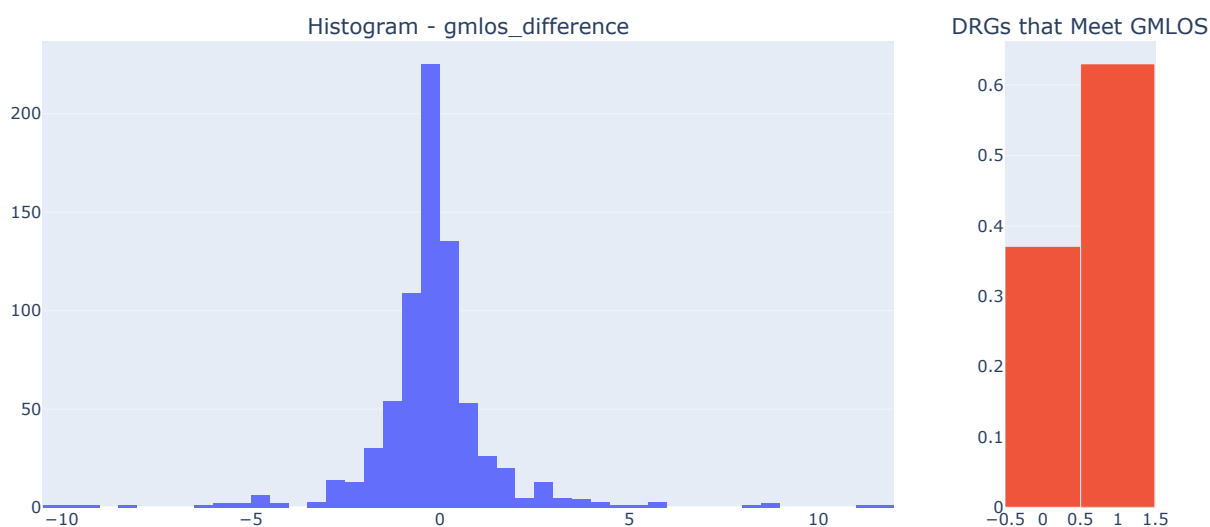
1

Plot distribution side by side to a bar chart showing the proportion

```

In [104]: 1
2 # Create a subplot with 1 row and 2 columns
3 fig = sp.make_subplots(rows=1, cols=2, subplot_titles=('Histogram - gmlos_difference', 'DRGs that Meet GMLOS'),
4                  column_widths=[0.85, 0.15])
5
6 # Histogram - gmlos_difference
7 fig.add_trace(go.Histogram(x=df2['gmlos_difference']), row=1, col=1)
8
9 # Bar Chart - drg_meet_it
10 meet_counts = df2['drg_meet_it'].value_counts(normalize=True)
11 fig.add_trace(go.Bar(x=meet_counts.index, y=meet_counts.values, #marker=dict(color=['skyblue', 'skyblue'])
12 #
13 # name='drg_meet_it',
14 # ), row=1, col=2)
15
16 # Update subplot layout
17 fig.update_layout(showlegend=False)
18 fig.update_layout(legend=dict(x=0, y=0))
19
20 fig.show()
21

```



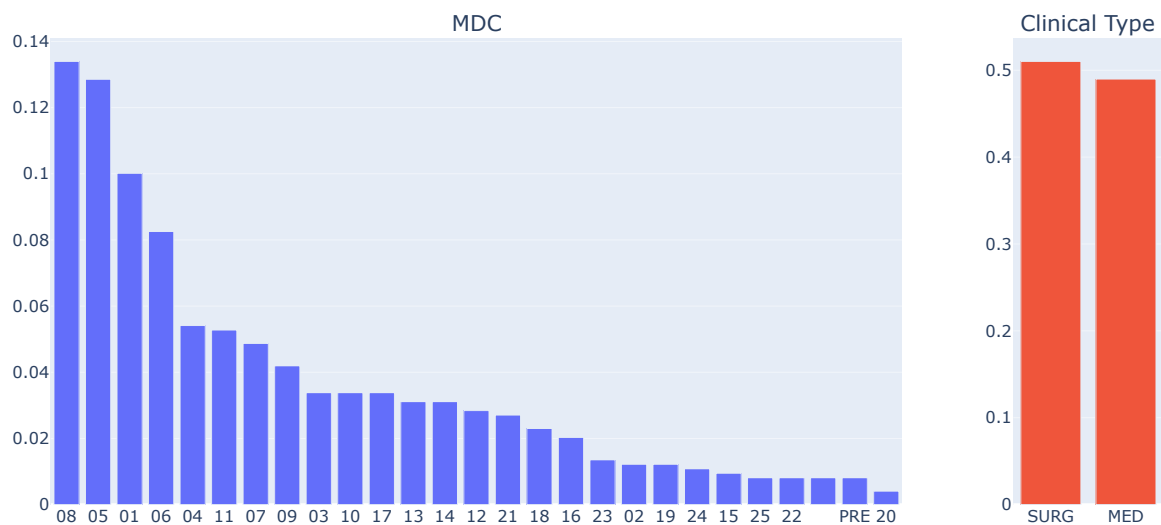
Plot bar chat of MDCs side by side to a bar chart showing the proportion of clinical type of patient



```

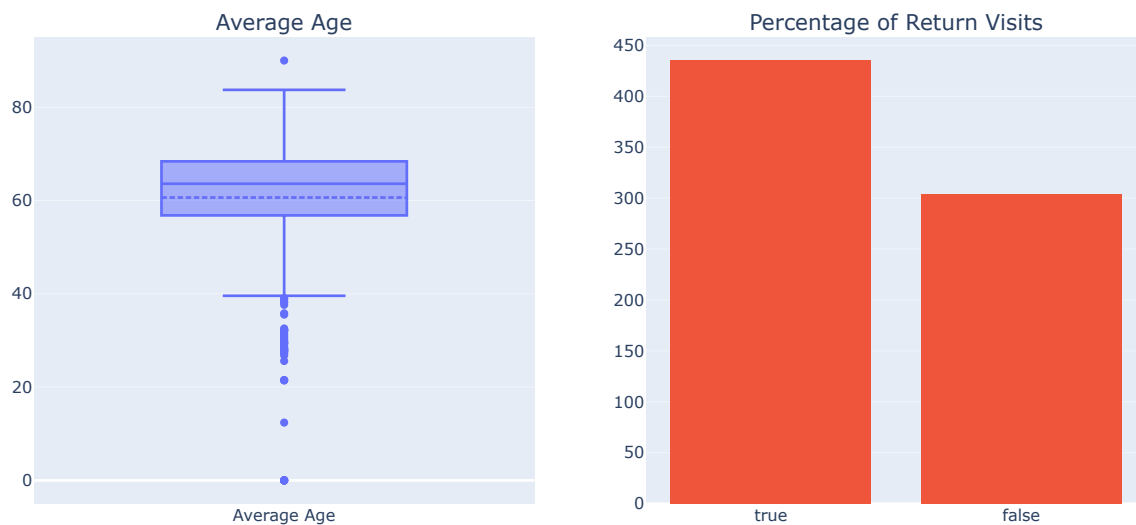
In [107]: 1 # Create a subplot with 1 row and 2 columns
2 fig = sp.make_subplots(rows=1, cols=2, subplot_titles=('Counts of MDC', 'Counts of Clinical Type'),
3               column_widths=[0.85, 0.15])
4
5
6 # Bar Chart - count of MDCs
7 meet_counts = df2['mdc'].value_counts(normalize=True)
8 fig.add_trace(go.Bar(x=meet_counts.index, y=meet_counts.values, #marker=dict(color=['skyblue', 'skyblue'])
9               name='drg_meet_it',
10              ), row=1, col=1)
11
12 # Bar Chart - count of Clinical type
13 meet_count = df2['clinical_type'].value_counts(normalize=True)
14 fig.add_trace(go.Bar(x=meet_count.index, y=meet_count.values, #marker=dict(color=['skyblue', 'skyblue'])
15               name='drg_meet_it',
16              ), row=1, col=2)
17
18 # Update subplot layout
19 fig.update_layout(showlegend=False)
20
21 fig.update_layout(legend=dict(x=0, y=0))
22
23 fig.show()
24

```



Plot Box plot of Age side by side to a bar chart showing the proportion of Return visit

```
In [121]: 1
2 # Create a subplot with 2 rows and 2 columns
3 fig = sp.make_subplots(rows=1, cols=2, subplot_titles=('Average Age', 'Percentage of Return Visits',
4                                                         ))
5
6 # Average Age
7 fig.add_trace(go.Box(y=df2['average_age'], boxmean=True, name='Average Age'), row=1, col=1)
8
9 # Percentage of Return Visits
10 fig.add_trace(go.Histogram(x=df2['pct_return_visits']>0, name=' Return Visits'), row=1, col=2)
11
12
13
14 # Update subplot layout
15 fig.update_layout(showlegend=False)
16
17 fig.show()
18
```



In [ ]:

1

Boxlot of average age and GMLOS grouped by meet\_it(0and 1)

```

In [ ]: 1 def plot(df1,title1,title2,y_text1,y_text2,target,col1,col2):
2
3     # Create a subplot with 1 row and 2 columns
4     fig = sp.make_subplots(rows=1, cols=2, subplot_titles=(title1, title2))
5
6     # Box plot for Average Age, grouped by drg_meet_it
7     for value in df1[target].unique():
8         fig.add_trace(go.Box(y=df1[df1[target] == value][col1],
9                             name='{} = {}'.format(target,value)), row=1, col=1)
10
11    # Box plot for Actual GMLOS, grouped by drg_meet_it
12    for value in df1[target].unique():
13        fig.add_trace(go.Box(y=df1[df1[target] == value][col2],
14                            name='{} = {}'.format(target,value)), row=1, col=2)
15
16    # Update subplot layout
17    fig.update_layout(showlegend=False)
18
19    # Set titles for the subplots
20    fig.update_annotations(
21        dict(text=title1, x=0.17, y=1.06, showarrow=False, font=dict(size=12)),
22        dict(text=title2, x=0.78, y=1.06, showarrow=False, font=dict(size=12))
23    )
24
25    # Update layout for each subplot
26    fig.update_xaxes(title_text=target, row=1, col=1)
27    fig.update_xaxes(title_text=target, row=1, col=2)
28    fig.update_yaxes(title_text=y_text1, row=1, col=1)
29    fig.update_yaxes(title_text=y_text2, row=1, col=2)
30
31    # Show the figure
32    fig.show()

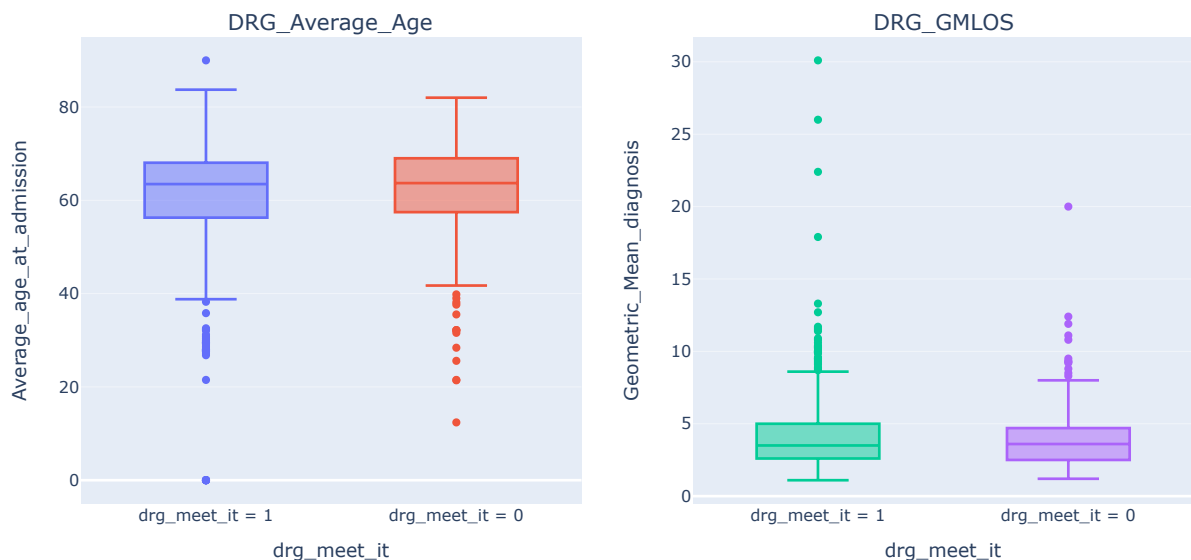
```

Boxplot of AGE and GMLOS grouped by drg\_meet\_it at DRG level

```

In [146]: 1 plot(df2, 'DRG_Average_Age', 'DRG_GMLOS', 'Average_age_at_admission', 'Geometric_Mean_diagnosis', 'drg_meet_it', 'average_age', '

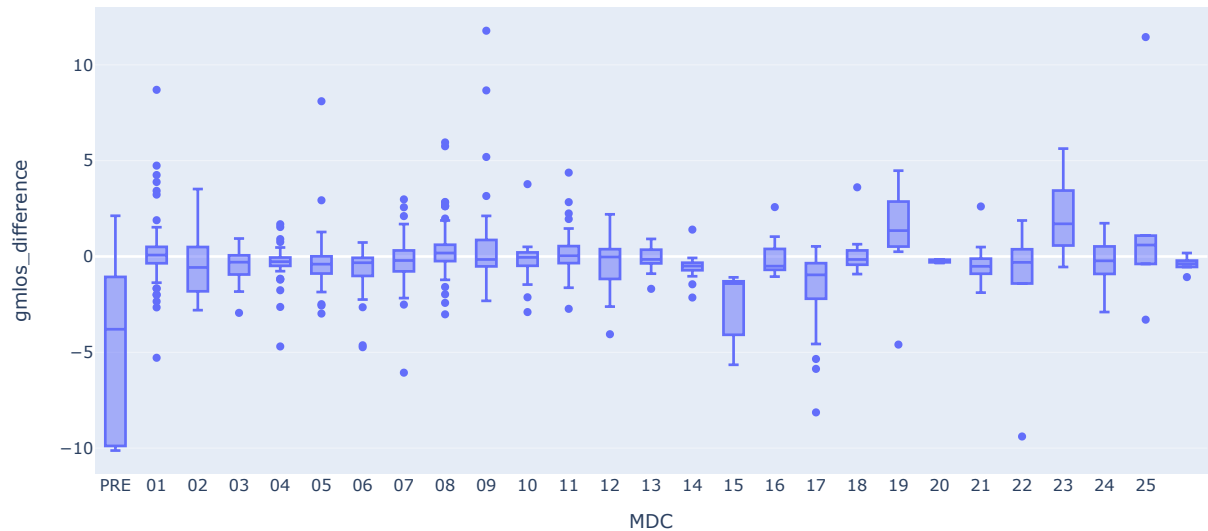
```



Boxlot of GMLOS difference grouped by MDC

```
In [156]: 1 def cat_box(df,x,y,title,xaxes_title,yaxes_title):
2         # Create the box plot
3         fig = px.box(df, x=x , y=y , title=title ,)
4
5         # Update Layout
6         fig.update_xaxes(title=xaxes_title)
7         fig.update_yaxes(title=yaxes_title)
8
9         # Show the figure
10        fig.show()
11        cat_box(df2,'mdc','gmlos_difference','Box Plot - gmlos_difference Grouped by MDC','MDC','gmlos_difference')
```

Box Plot - gmlos\_difference Grouped by MDC



```
In [ ]: 1
```

**Difference between Encounter LOS(Patient's Actual GMLOS) and given Geometric Mean LOS**

**- ENCOUNTER LEVEL**

We checked for the distribution for each encounter in the as recorded system in the most granular form

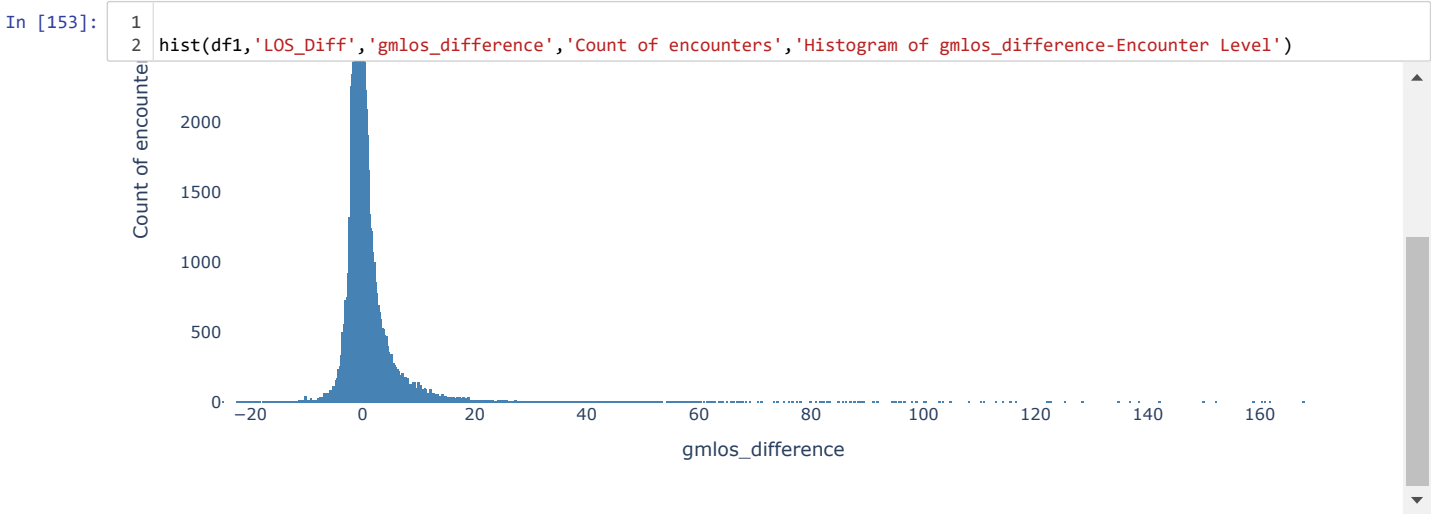
```
In [54]: 1 df1
```

	drg_code	mdc	gender	clinical_type	age	return_visit	gmlos	LOS	LOS_Diff	meet_it
0	872	18	M	MED	55	0	3.5	4.162500	0.662500	0
1	441	07	F	MED	42	0	4.7	2.091667	-2.608333	1
2	871	18	F	MED	46	0	4.8	7.886806	3.086806	0
3	442	07	F	MED	47	0	3.2	1.006250	-2.193750	1
4	708	12	M	SURG	51	0	1.4	1.203472	-0.196528	1
...	...	...	...	...	...	...	...	...	...	...
149520	794	15	U	MED	0	0	3.4	1.873681	-1.526319	1
149521	026	01	F	SURG	64	0	3.8	1.264468	-2.535532	1
149522	787	14	F	SURG	37	0	3.5	2.150914	-1.349086	1
149523	982		F	SURG	35	0	4.6	1.933727	-2.666273	1
149524	795	15	M	MED	0	0	3.1	2.054375	-1.045625	1

149525 rows x 10 columns

The distribution is very right skewed as a lot of patients are not meeting the specific GMLOS for their ailment as suggested by the government

```
In [ ]: 1 def hist(df,x,xaxes_title,yaxes_title,layout_title):
2         fig = px.histogram(df, x)
3
4
5         # Customize the bar colors
6         fig.update_traces(marker_color='steelblue')
7
8         # Set x-axis Label
9         fig.update_xaxes(title=xaxes_title)
10
11        # Set y-axis Label
12        fig.update_yaxes(title=yaxes_title)
13
14        # Set plot title and center-align it
15        fig.update_layout(
16            title=layout_title,
17            title_x=0.5 # Value of 0.5 centers the title
18        )
19
20        # Customize the layout
21        fig.update_layout(
22            plot_bgcolor='white', # Set plot background color
23        )
24
25        fig.show()
26
```



1

To understand what is happening we reduced the skewness by eliminating those who exceed the limit over 25 days differnt

```
In [150]: 1 #slice the frame limiting to those patient who do not exceed their GMLOS more than 25 days
2 data = df1[df1.LOS_Diff<25]
3 data
```

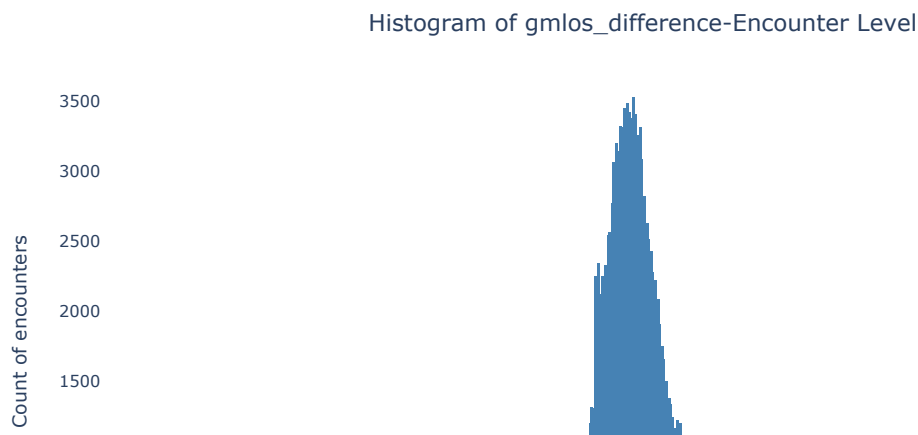
Out[150]:

	drp_code	mdc	gender	clinical_type	age	return_visit	gmlos	LOS	LOS_Diff	meet_it
0	872	18	M	MED	55	0	3.5	4.162500	0.662500	0
1	441	07	F	MED	42	0	4.7	2.091667	-2.608333	1
2	871	18	F	MED	46	0	4.8	7.886806	3.086806	0
3	442	07	F	MED	47	0	3.2	1.006250	-2.193750	1
4	708	12	M	SURG	51	0	1.4	1.203472	-0.196528	1
...	...	...	...	...	...	...	...	...	...	...
149520	794	15	U	MED	0	0	3.4	1.873681	-1.526319	1
149521	026	01	F	SURG	64	0	3.8	1.264468	-2.535532	1
149522	787	14	F	SURG	37	0	3.5	2.150914	-1.349086	1
149523	982		F	SURG	35	0	4.6	1.933727	-2.666273	1
149524	795	15	M	MED	0	0	3.1	2.054375	-1.045625	1

148603 rows × 10 columns

```
In [ ]: 1
```

```
In [152]: 1 hist(data, 'LOS_Diff', 'gmlos_difference', 'Count of encounters', 'Histogram of gmlos_difference-Encounter Level')
```



```
In [ ]: 1
```

### Scatter plot of age vs los\_difference

This does not give a particular correlation the data point is clustered around zero level as expected

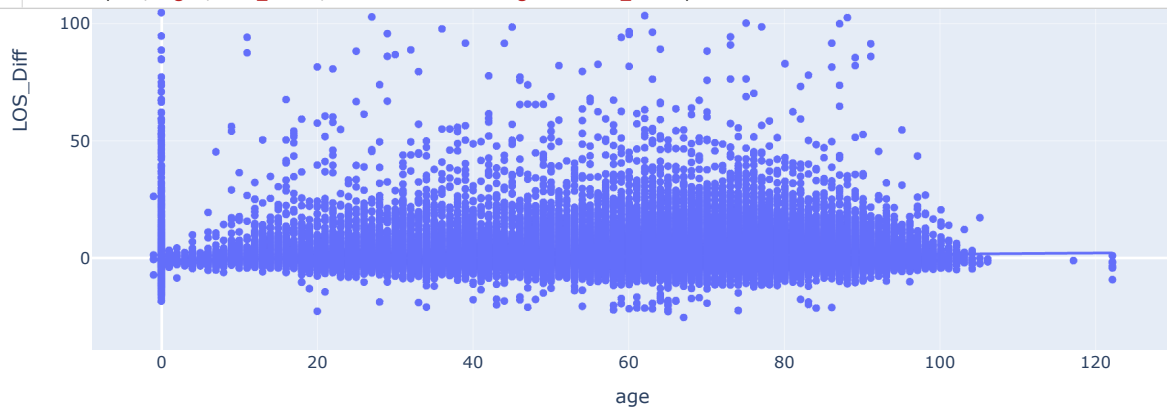
```
In [315]: 1 df1
```

Out[315]:

	drg_code	mdc	gender	clinical_type	age	return_visit	gmlos	LOS	LOS_Diff	meet_it
0	872	18	M	MED	55	0	3.5	4.162500	0.662500	0
1	441	07	F	MED	42	0	4.7	2.091667	-2.608333	1
2	871	18	F	MED	46	0	4.8	7.886806	3.086806	0
3	442	07	F	MED	47	0	3.2	1.006250	-2.193750	1
4	708	12	M	SURG	51	0	1.4	1.203472	-0.196528	1
...	...	...	...	...	...	...	...	...	...	...
149520	794	15	U	MED	0	0	3.4	1.873681	-1.526319	1
149521	026	01	F	SURG	64	0	3.8	1.264468	-2.535532	1
149522	787	14	F	SURG	37	0	3.5	2.150914	-1.349086	1
149523	982		F	SURG	35	0	4.6	1.933727	-2.666273	1
149524	795	15	M	MED	0	0	3.1	2.054375	-1.045625	1

149525 rows × 10 columns

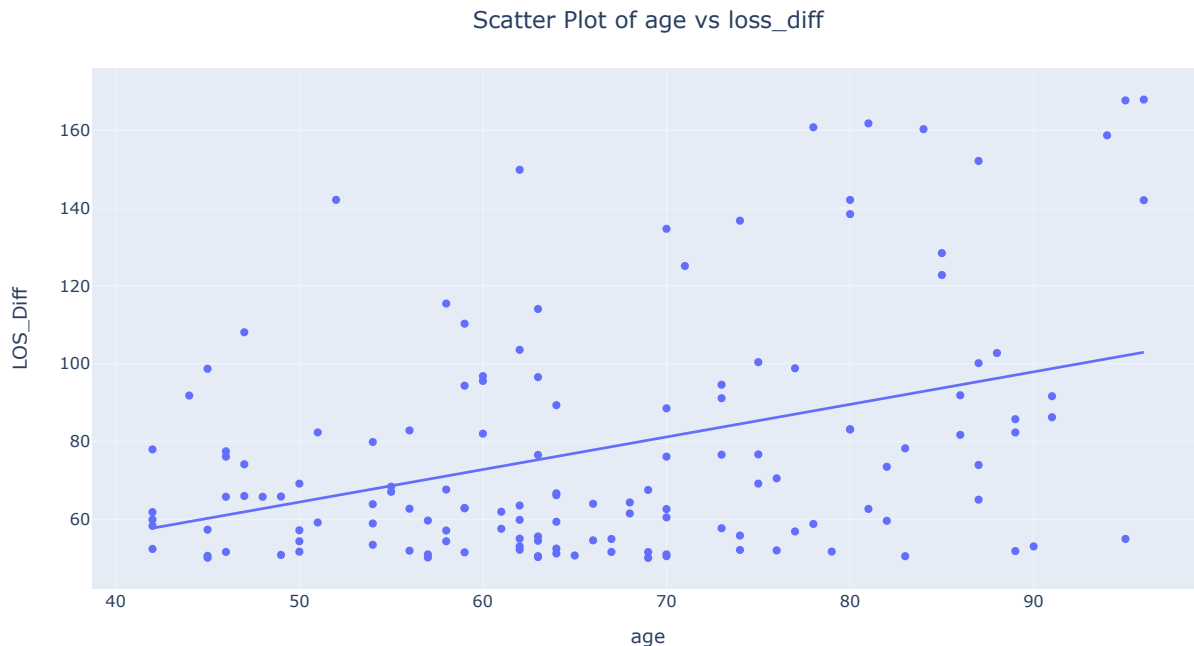
```
In [316]: 1
2 scatter(df1, 'age', 'LOS_Diff', 'Scatter Plot of age vs loss_diff')
```



### Scatter plot for slized Age vs Difference in Lenght of stay (Actual GMLOS -GMLOS)

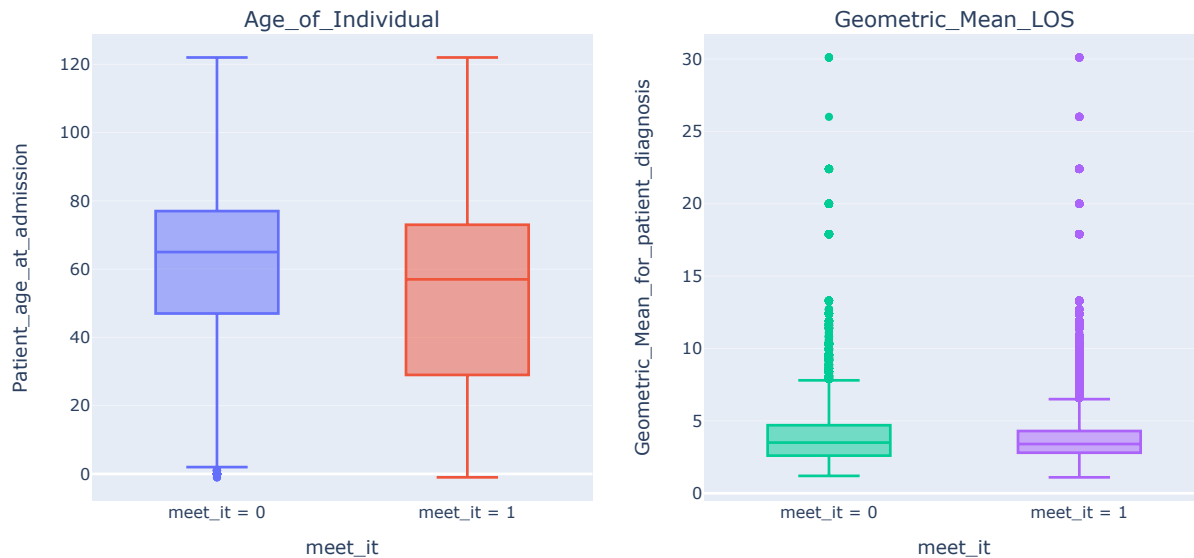
When we removed those who constantly meet it e.g foetus we see a clearer linear correlation between LOS\_Difference and age of patient. It shows that as patient get older there is a higher tendency not to meet the GMLOS

```
In [319]: 1 ##sliced dataframe for those LOS
2 df_sliced = df1[(df1['LOS_Diff'] > 50) & (df1['age'] > 40)]
3 scatter(df_sliced, 'age', 'LOS_Diff', 'Scatter Plot of age vs loss_diff')
```



Boxplot of AGE and GMLOS grouped by drg\_meet\_it at Encounter level

```
In [147]: 1 plot(df1, 'Age_of_Individual', 'Geometric_Mean_LOS', 'Patient_age_at_admission', 'Geometric_Mean_for_patient_diagnosis', 'meet_
```



```
In [ ]: 1
```

Boxplot of AGE and GMLOS grouped by clinical\_type at Encounter level

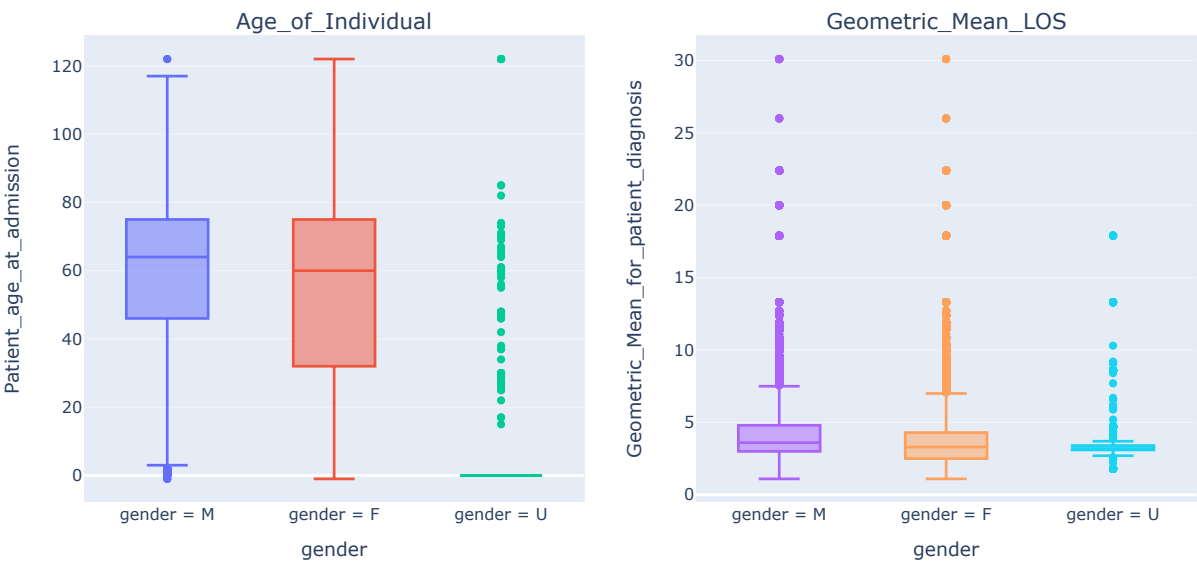
```
In [148]: 1 plot(df1, 'Age_of_Individual', 'Geometric_Mean_LOS', 'Patient_age_at_admission', 'Geometric_Mean_for_patient_diagnosis', 'clinical_type')
```



```
In [ ]: 1
```

Boxplot of AGE and GMLOS grouped by Gender(Female , Male and Unknown) at Encounter level

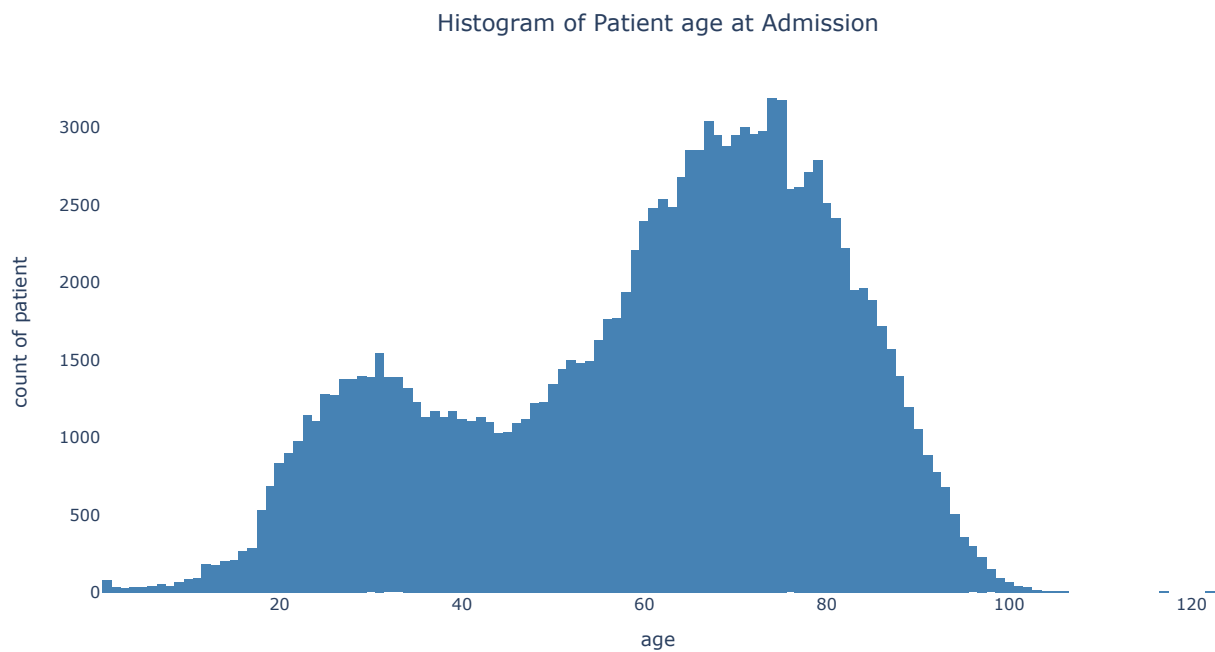
```
In [149]: 1 plot(df1, 'Age_of_Individual', 'Geometric_Mean_LOS', 'Patient_age_at_admission', 'Geometric_Mean_for_patient_diagnosis', 'gender')
```



Histogram of Age for patients Greater than 0 yrs



```
In [216]: 1 hist(df1[df1.age>0], 'age', 'age', 'count of patient', 'Histogram of Patient age at Admission')
```



```
In [ ]: 1
```

```
In [ ]: 1 df=pd.read_csv('encounter_records.csv')
2 df.drop(columns=['Unnamed: 0'],inplace=True)
3
4 ###replace the MDC without name with 00
5 df.mdc.replace(to_replace=' ', value='00',inplace=True)
6
```

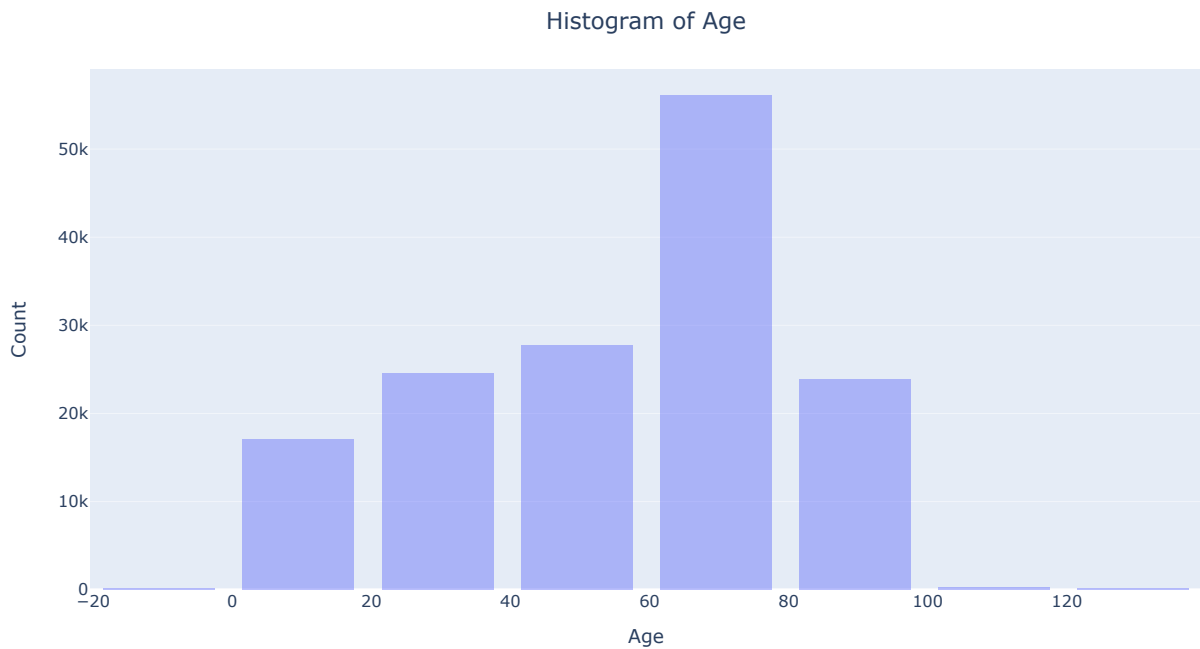
```
In [220]: 1 #Create new column for flagging patients with complication
2 df['with_cc_mcc']=((df['drg_description'].str.contains('WITH MCC')== True) | (df['drg_description'].str.contains('WITH CC/
```

```
In [221]: 1 #Create new column for flagging patient using ventilator
2 df['with_vent']=((df['drg_description'].str.contains('WITH MV')== True) | (df['drg_description'].str.contains('WITH VENTIL
```

```
In [222]: 1 df1=df[['gender_abbreviation','clinical_type','drg_code',
2 'mdc','age_at_admit','with_cc_mcc','with_vent','return_visit','meet_it']]
3
4 df=df1.copy()
5 df.rename(columns={'gender_abbreviation':'gender','age_at_admit':'age','mdc':'mdc','with_cc_mcc':'cc_mc','with_vent':'vent',
6 ,inplace=True)
7
8 df['gender']=df.gender.map({'M':'M','U':'M','F':'F'})
9
10 bins = [-np.inf, 20, 40, 60, 70, 80, np.inf]
11 names = ['<20', '20_40', '40_60', '60_70', '70_80', '80more']
12
13 df['age_bin'] = pd.cut(df['age'], bins, labels=names)
```

**Histogram of Binned Age**

```
In [262]: 1 def hist_bin(df,x,nbins,title,xaxis_title):
2         # Create a histogram using Plotly with adjusted bar space
3         fig = px.histogram(df1, x=x, nbins= nbins, barmode='overlay')
4
5         # Update Layout and axis Labels
6         fig.update_layout(title=title,
7                           title_x=0.5, # Set the title position to the center
8                           xaxis_title=xaxis_title,
9                           yaxis_title='Count',
10                          bargap=0.2) # Adjust the value to increase/decrease the bar space
11
12         # Update aesthetic styles
13         fig.update_traces(opacity=0.90) # Change the bar color and opacity
14         fig.update_layout( showlegend=False) # Set background color and hide Legend
15
16         # Show the histogram
17         fig.show()
18 hist_bin(df,x='age_at_admit',nbins=10,title='Histogram of Age',xaxis_title='Age')
```



In [ ]:

1

Bar chart of proportion of encounters that meet it or not (1 or 0)

```

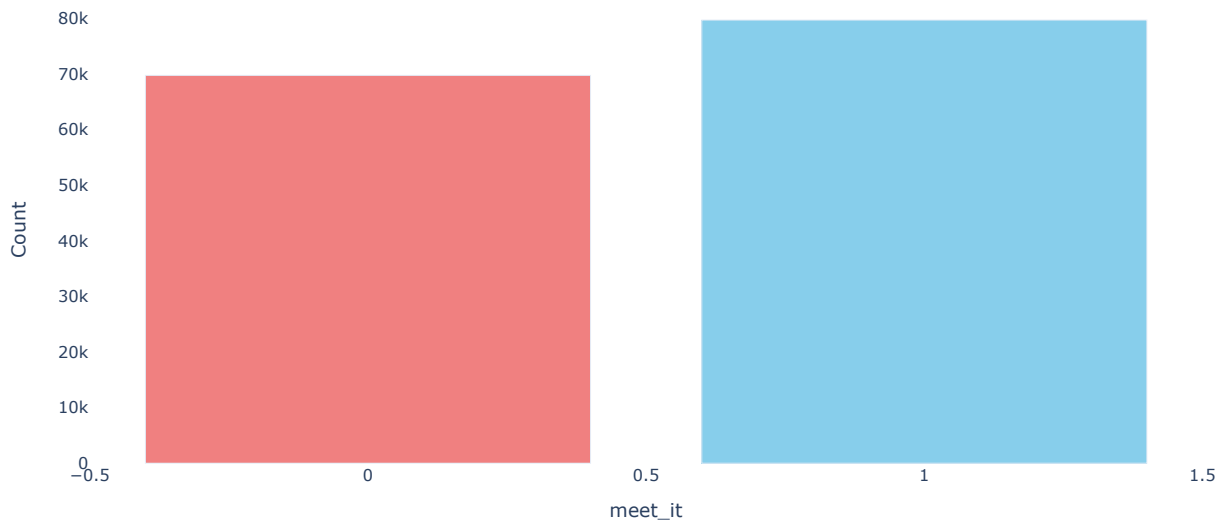
In [277]: 1 def barchart(df,target,subject):
2           # Count the number of 1s and 0s in the 'meet_it' column
3           meet_it_counts = df1[target].value_counts()
4
5           # Create a bar chart using Plotly
6           fig = px.bar(x=meet_it_counts.index, y=meet_it_counts.values,
7                        color=meet_it_counts.index, labels={'x': target, 'y': 'Count'},
8                        title='Bar Chart of Encounter that {}'.format(subject))
9
10          # Update aesthetic styles
11          fig.update_traces(marker_color=['skyblue', 'lightcoral']) # Set custom colors for the bars
12          fig.update_layout(plot_bgcolor='white', title_x=0.5, showlegend=False) # Set background color and hide Legend
13          fig.update_xaxes(title_text=target) # Set x-axis label
14
15          # Show the bar chart
16          fig.show()
17          barchart(df1,'meet_it','meet GMLOS')

```

C:\ProgramData\Anaconda3\lib\site-packages\numpy\core\numeric.py:2327: FutureWarning:

elementwise comparison failed; returning scalar instead, but in the future will perform elementwise comparison

Bar Chart of Encounter that meet GMLOS

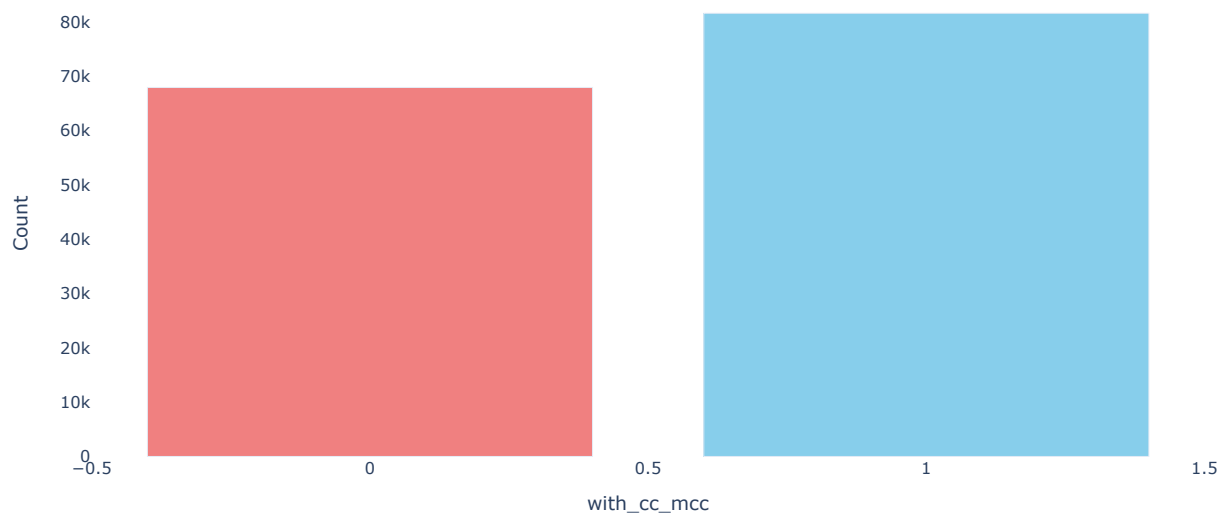


In [ ]: 1

Bar chart of proportion of encounters that has complications or not (1 or 0)

```
In [279]: 1 barchart(df1,'with_cc_mcc', " are with Complication during Admission")
```

Bar Chart of Encounter that are with Complication during Admission



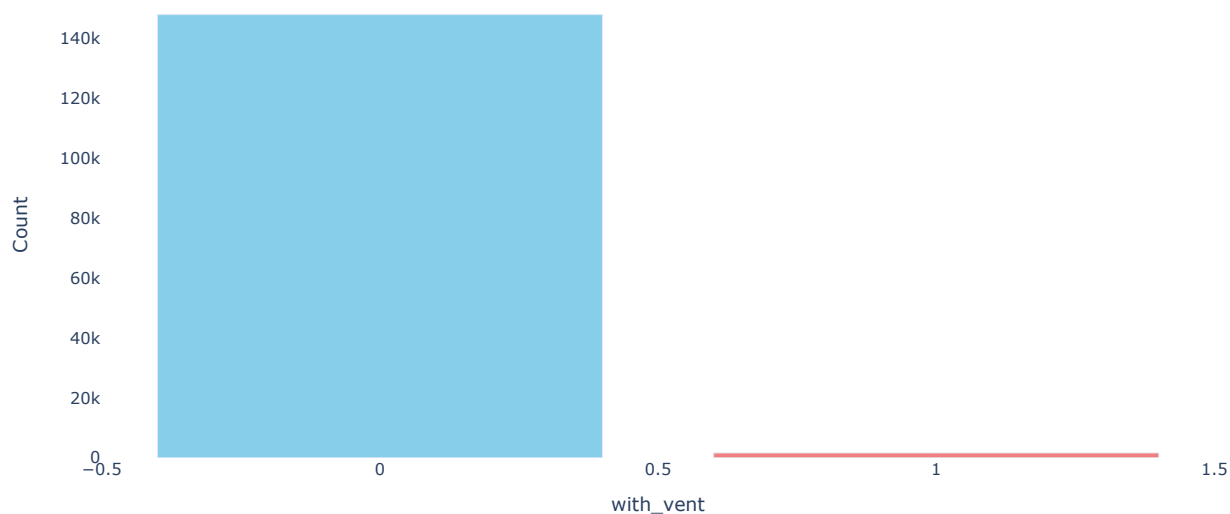
Bar chart of proportion of encounters that use ventilator or not (1 or 0)

```
In [281]: 1 barchart(df1,'with_vent', " are using Ventilation")
```

C:\ProgramData\Anaconda3\lib\site-packages\numpy\core\numeric.py:2327: FutureWarning:

elementwise comparison failed; returning scalar instead, but in the future will perform elementwise comparison

Bar Chart of Encounter that are using Ventilation

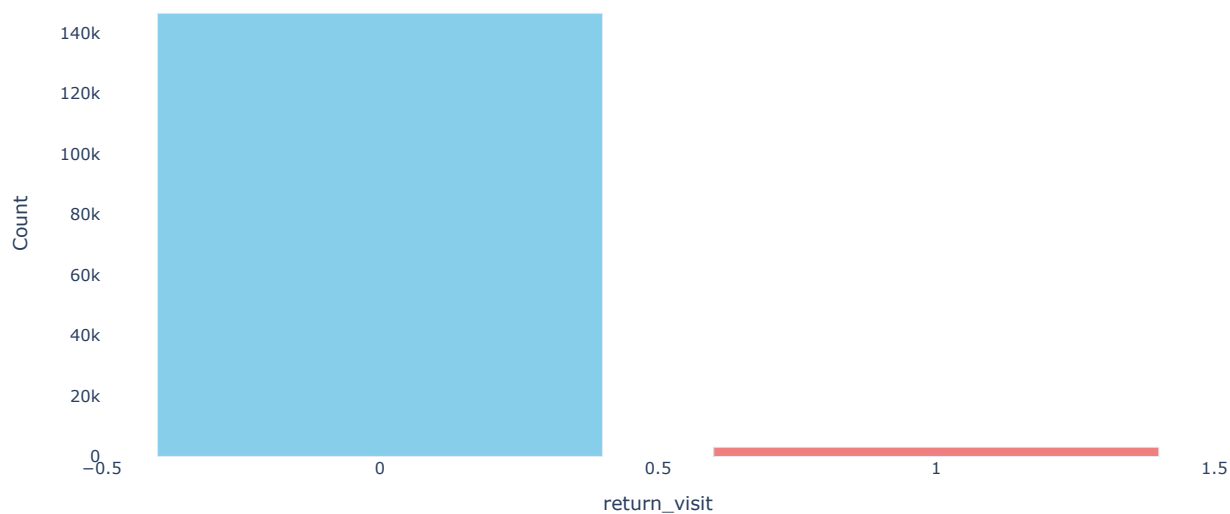


```
In [ ]: 1
```

Bar chart of proportion of encounters that return to the hospital or not (1 or 0)

```
In [282]: 1 barchart(df1,'return_visit', 'have been in the system before')
```

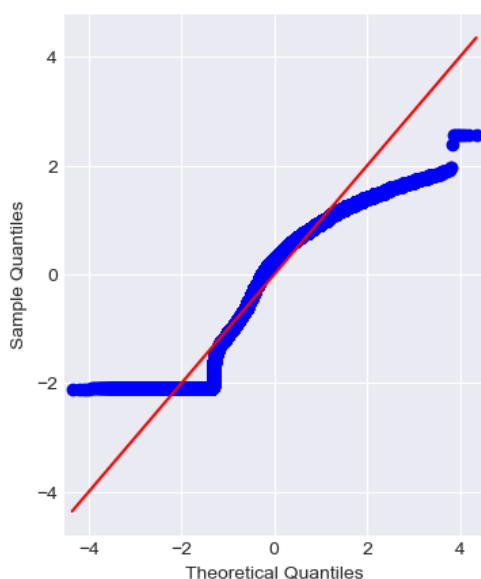
Bar Chart of Encounter that have been in the system before



### Check for normality in age distribution

As we have known by now , there is a violation of normality for the distribution of age . A lot of outliers present

```
In [244]: 1 # check normality
2 fig, ax = plt.subplots(figsize=(4,5))
3 qqplot(df1['age_at_admit'], fit=True, line='s',ax=ax)
4 plt.show()
```



## Modelling

- Used a variety of statistical and machine learning models to assess the data, with the goal of building an accurate model that could predict whether a patient will meet or not meet the GMLOS. Some of these models gave a likelihood to meet the GMLOS. In all of the models, we were attempting to predict for each encounter a binary variable that represented whether the encounter met the GMLOS or not.
-

```

3 We first tried logistic regression. This model was 59.1% accurate. We tried another logistic regression but removed any
  patients with an admit age of 0 or below, just to avoid the extreme violation of normality, to remove newborn patients,
  which exhibited different behavior than the rest of our patients. This model was less accurate than the first. Then tried
  a standard decision tree model on the full data, which was 59.26% accurate, and a gradient boosted classifier, which was
  59.48% accurate. Finally tried a random forest model which returned our most accurate model by a thin margin at 59.5%. For
  reference, 53.23% should be used as a benchmark as this is what percentage of our data met the GMLOS. It can be seen that
  these models performed marginally better than a blanket guess of "Yes", meaning that their interpretation could provide
  some value, but may not be extremely accurate as a predictor at admission for a patient as there is significant lack of
  fit .

```

### Grouping the Categories

```

In [ ]: 1 df['mdc']=df.mdc.astype('category')
        2 df['gender']=df.gender.astype('category')
        3 df['clinical_type']=df.clinical_type.astype('category')
        4 df['drg_code']=df.drg_code.astype('category')
        5
        6 df1=df.copy()
        7 df.drop(columns=['drg_code'],inplace=True)

```

### Creating Dummies

```

In [227]: 1 dummies_age_bin = pd.get_dummies(df['age_bin'])
          2 dummies_age_bin.drop(columns='<20',inplace=True)
          3
          4 dummies_mdc = pd.get_dummies(df['mdc'],prefix='mdc')
          5 dummies_mdc.drop(columns='mdc_00',inplace=True)
          6
          7 dummies_clinical_type = pd.get_dummies(df['clinical_type'])
          8 dummies_clinical_type.drop(columns='SURG',inplace=True)
          9
          10 dummies_gender = pd.get_dummies(df['gender'])
          11 dummies_gender.drop(columns='M',inplace=True)
          12
          13 categorical_features = ['clinical_type', 'gender', 'mdc', 'age_bin']
          14 continuous_features = ['age']
          15 binary_features = ['return_visit', 'cc_mc', 'vent']
          16 target = ['meet_it']
          17 target=df[target]

```

### Final dataframe

```

In [228]: 1 data = pd.concat([df[binary_features], dummies_gender], axis=1)
          2 data_1 = pd.concat([data, dummies_clinical_type], axis=1)
          3 data_2 = pd.concat([data_1, dummies_mdc], axis=1)
          4 data_3 = pd.concat([data_2, dummies_age_bin], axis=1)
          5 df = pd.concat([data_3, target], axis=1)
          6

```

```

In [258]: 1 df

```

Out[258]:

	return_visit	cc_mc	vent	F	MED	mdc_01	mdc_02	mdc_03	mdc_04	mdc_05	mdc_06	mdc_07	mdc_08	mdc_09	mdc_10	mdc_11	mdc_12	mi
0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	1	0	1	1	0	0	0	0	0	0	1	0	0	0	0	0	0
2	0	1	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0
3	0	1	0	1	1	0	0	0	0	0	0	1	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
149520	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
149521	0	1	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0
149522	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
149523	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
149524	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0

149525 rows × 37 columns

```

In [ ]: 1

```

```
In [283]: 1 # rename columns
          2 df.rename(columns={'20-40': '20_40', '40-60': '40_60', '60-70': '60_70', '70-80': '70_80', '80+': '80more'})
```

Out[283]:

	return_visit	cc_mc	vent	F	MED	mdc_01	mdc_02	mdc_03	mdc_04	mdc_05	mdc_06	mdc_07	mdc_08	mdc_09	mdc_10	mdc_11	mdc_12	mi
0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	1	0	1	1	0	0	0	0	0	0	1	0	0	0	0	0	0
2	0	1	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0
3	0	1	0	1	1	0	0	0	0	0	0	1	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
149520	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
149521	0	1	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0
149522	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
149523	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
149524	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0

149525 rows × 37 columns

```
In [ ]: 1
```

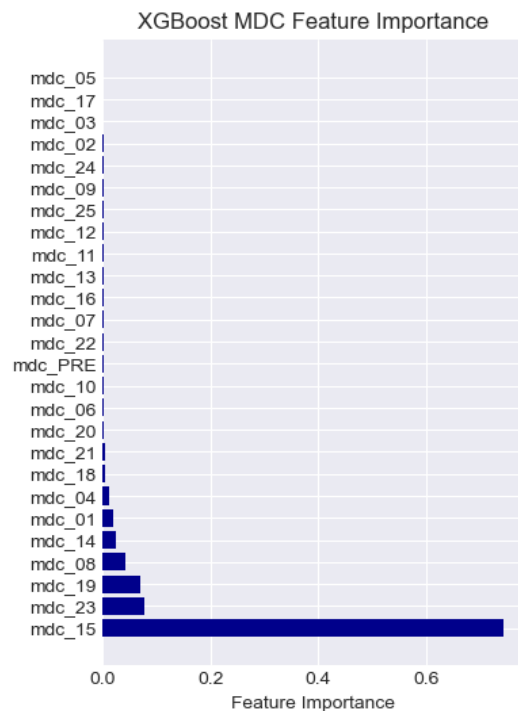
Check for feature Importance

The features were grouped into different sects and then checked their importance using XGbosst feature importance capability. The Features was latter combined and checked for the top 10 most important feature in the data

```

In [287]: 1 # plot feature importance manually
2 from numpy import loadtxt
3 from xgboost import XGBClassifier
4 import matplotlib.pyplot as plt
5
6
7
8
9 ### MDC features
10
11 features_mdc=['mdc_01', 'mdc_02',
12             'mdc_03', 'mdc_04', 'mdc_05', 'mdc_06', 'mdc_07', 'mdc_08', 'mdc_09',
13             'mdc_10', 'mdc_11', 'mdc_12', 'mdc_13', 'mdc_14', 'mdc_15', 'mdc_16',
14             'mdc_17', 'mdc_18', 'mdc_19', 'mdc_20', 'mdc_21', 'mdc_22', 'mdc_23',
15             'mdc_24', 'mdc_25', 'mdc_PRE',
16             ]
17 target =['meet_it']
18
19 X = df[features_mdc]
20 y = df[target]
21
22 # fit model on training data
23 model = XGBClassifier()
24 model.fit(X, y)
25
26 # feature importance
27 feature_importance = model.feature_importances_
28
29 # sort feature importance in descending order
30 sorted_idx = feature_importance.argsort()[::-1]
31 sorted_features = [features_mdc[i] for i in sorted_idx]
32 sorted_importance = feature_importance[sorted_idx]
33
34 # create a horizontal bar plot
35 fig, ax = plt.subplots(figsize=(4, 6))
36 ax.barh(sorted_features, sorted_importance, color='darkblue')
37 ax.set_xlabel('Feature Importance')
38 ax.set_title('XGBoost MDC Feature Importance')
39
40 plt.show()
41 #plt.savefig('MDC_feat_import.png',dpi=360)
42

```



```

In [321]: 1 ### Other feature group

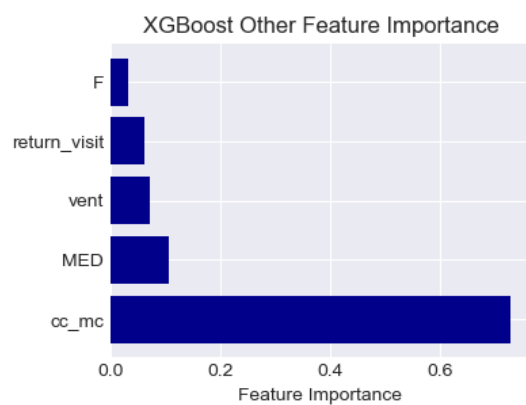
```



```

In [288]: 1 features_others=['return_visit', 'cc_mc', 'vent', 'F', 'MED']
          2
          3
          4 X = df[features_others]
          5 y = df[target]
          6
          7 # fit model on training data
          8 model = XGBClassifier()
          9 model.fit(X, y)
         10
         11 # feature importance
         12 feature_importance = model.feature_importances_
         13
         14 # sort feature importance in descending order
         15 sorted_idx = feature_importance.argsort()[::-1]
         16 sorted_features = [features_others[i] for i in sorted_idx]
         17 sorted_importance = feature_importance[sorted_idx]
         18
         19 # create a horizontal bar plot
         20 fig, ax = plt.subplots(figsize=(4, 3))
         21 ax.barh(sorted_features, sorted_importance, color='darkblue')
         22 ax.set_xlabel('Feature Importance')
         23 ax.set_title('XGBoost Other Feature Importance')
         24
         25 plt.show()
         26 #plt.savefig('other_feat_import.png',dpi=360)

```



```

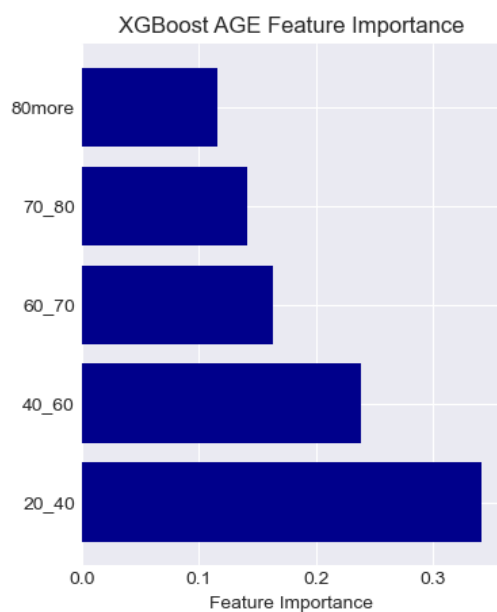
In [322]: 1 ##### Age feature group

```

```

In [290]: 1 features_age=['20_40', '40_60', '60_70', '70_80', '80more']
          2
          3
          4 X = df[features_age]
          5 y = df[target]
          6
          7 # fit model on training data
          8 model = XGBClassifier()
          9 model.fit(X, y)
         10
         11 # feature importance
         12 feature_importance = model.feature_importances_
         13
         14 # sort feature importance in descending order
         15 sorted_idx = feature_importance.argsort()[::-1]
         16 sorted_features = [features_age[i] for i in sorted_idx]
         17 sorted_importance = feature_importance[sorted_idx]
         18
         19 # create a horizontal bar plot
         20 fig, ax = plt.subplots(figsize=(4, 5))
         21 ax.barh(sorted_features, sorted_importance, color='darkblue')
         22 ax.set_xlabel('Feature Importance')
         23 ax.set_title('XGBoost AGE Feature Importance')
         24
         25 plt.show()
         26 #plt.savefig('age_feat_import.png',dpi=360)

```



In [ ]:

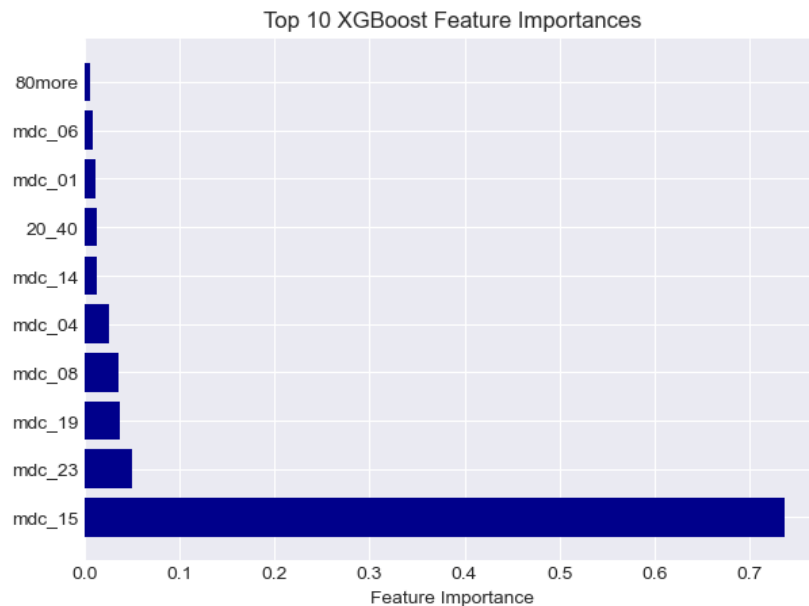
1

**All features combined**

```

In [291]: 1 features=['return_visit', 'cc_mc', 'vent', 'F', 'MED', 'mdc_01', 'mdc_02',
2           'mdc_03', 'mdc_04', 'mdc_05', 'mdc_06', 'mdc_07', 'mdc_08', 'mdc_09',
3           'mdc_10', 'mdc_11', 'mdc_12', 'mdc_13', 'mdc_14', 'mdc_15', 'mdc_16',
4           'mdc_17', 'mdc_18', 'mdc_19', 'mdc_20', 'mdc_21', 'mdc_22', 'mdc_23',
5           'mdc_24', 'mdc_25', 'mdc_PRE', '20_40', '40_60', '60_70', '70_80',
6           '80more']
7 target = ['meet_it']
8
9 df[features]
10
11 X = df[features]
12 y = df[target]
13
14 # fit model on training data
15 model = XGBClassifier()
16 model.fit(X, y)
17
18 # feature importance
19 feature_importance = model.feature_importances_
20
21 # sort feature importance in descending order
22 sorted_idx = feature_importance.argsort()[::-1]
23 sorted_features = [features[i] for i in sorted_idx]
24 sorted_importance = feature_importance[sorted_idx]
25
26 # create a horizontal bar plot of the top 10 feature importances
27 top_k = 10
28 fig, ax = plt.subplots(figsize=(7, 5))
29 ax.barh(sorted_features[:top_k], sorted_importance[:top_k], color='darkblue')
30 ax.set_xlabel('Feature Importance')
31 ax.set_title(f'Top {top_k} XGBoost Feature Importances')
32
33 plt.show()
34 #plt.savefig('Top_10_allfeat_import.png', dpi=360)

```



## Model Building

```

In [72]: 1 from sklearn.model_selection import train_test_split
2         from sklearn import tree
3
4 X = df[features].copy()
5 y = df[target].copy()
6
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30, random_state=0)

```

```

In [106]: 1 ## Decision Tree
2
3 clf = tree.DecisionTreeClassifier(random_state=0,)
4 clf = clf.fit(X_train, y_train)
5 clf.score(X_train, y_train)

```

Out[106]: 0.6077830008407857

```
In [107]: 1 clf.score(X_test, y_test)
```

```
Out[107]: 0.5926074549670055
```

## Random Forest

```
In [123]: 1 from sklearn.ensemble import RandomForestClassifier
2
3 #y_train=y_train.meet_it.ravel()
4 clf = RandomForestClassifier(random_state=70,max_depth=13)
5 clf = clf.fit(X_train, y_train)
6 clf.score(X_train, y_train)
```

```
Out[123]: 0.6037701597492929
```

```
In [124]: 1 clf.score(X_test, y_test)
```

```
Out[124]: 0.5954164437310505
```

## Logistic Regression

```
In [131]: 1 from sklearn.linear_model import LogisticRegression
2 from sklearn.preprocessing import StandardScaler
3
4 X = df[features].copy()
5 y = df[target].copy()
6
7 scaler = StandardScaler()
8 X_train = scaler.fit_transform(X_train)
9 X_test = scaler.transform(X_test)
10
11 y_train = y_train.meet_it.ravel()
12 y_test = y_test.meet_it.ravel()
13
14 logisticRegr = LogisticRegression(max_iter=1000)
15 logisticRegr.fit(X_train, y_train)
16
17 predictions = logisticRegr.predict(X_test)
18
```

```
In [132]: 1 # Use score method to get accuracy of model
2 score = logisticRegr.score(X_test, y_test)
3 print(score)
```

```
0.5911137863385054
```

```
In [133]: 1 from sklearn.metrics import confusion_matrix
2 y_true = y_test
3 y_pred = predictions
4 confusion_matrix(y_true, y_pred)
```

```
Out[133]: array([[11080, 9906],
 [ 8435, 15435]], dtype=int64)
```

```
In [135]: 1 # lst=[[11107, 9879],
2 #       [ 8468, 15402]]
3 lst=confusion_matrix(y_true, y_pred).tolist()
4 lst
```

```
Out[135]: [[11080, 9906], [8435, 15435]]
```

```
In [140]: 1 tot_P=(y_test==1).sum()
2 tot_N=(y_test==0).sum()
3 true_Neg = lst[0][0]/tot_N
4 false_Neg = lst[0][1]/tot_N
5 false_Pos = lst[1][0]/tot_P
6 true_Pos = lst[1][1]/tot_P
7 print(f'true positive rate is: {true_Pos},false positive rate is: {false_Pos}')
8 print(f'true negative rate is: {true_Neg}, false negative rate is: {false_Neg}')
9
```

```
true positive rate is: 0.6466275659824047,false positive rate is: 0.3533724340175953
true negative rate is: 0.527971028304584, false negative rate is: 0.47202897169541597
```

```
In [151]: 1 precision=(true_Pos/(true_Pos+false_Pos))
```

```
In [155]: 1 recall=(true_Pos/(true_Pos+false_Neg))
```

```
In [156]: 1 f1=2*(precision*recall/(precision+recall))
```

```
In [157]: 1 f1
```

```
Out[157]: 0.6104128295293665
```

## XGBoost

```
In [143]: 1 from numpy import loadtxt
2 from xgboost import XGBClassifier
3 from sklearn.model_selection import train_test_split
4 from sklearn.metrics import accuracy_score
```

```
In [144]: 1 X = df[features]
2 Y = df[target]
3
```

```
In [145]: 1 # split data into train and test sets
2 seed = 7
3 test_size = 0.33
4 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=test_size, random_state=seed)
5
```

## Train Model

```
In [146]: 1
2 # fit model no training data
3 model = XGBClassifier()
4 model.fit(X_train, y_train)
```

```
Out[146]: XGBClassifier(base_score=None, booster=None, callbacks=None,
                        colsample_bylevel=None, colsample_bynode=None,
                        colsample_bytree=None, early_stopping_rounds=None,
                        enable_categorical=False, eval_metric=None, feature_types=None,
                        gamma=None, gpu_id=None, grow_policy=None, importance_type=None,
                        interaction_constraints=None, learning_rate=None, max_bin=None,
                        max_cat_threshold=None, max_cat_to_onehot=None,
                        max_delta_step=None, max_depth=None, max_leaves=None,
                        min_child_weight=None, missing=nan, monotone_constraints=None,
                        n_estimators=100, n_jobs=None, num_parallel_tree=None,
                        predictor=None, random_state=None, ...)
```

```
In [147]: 1
2 # make predictions for test data
3 y_pred = model.predict(X_test)
4 predictions = [round(value) for value in y_pred]
```

```
In [148]: 1
2 # evaluate predictions
3 accuracy = accuracy_score(y_test, predictions)
4 print("Accuracy: %.2f%%" % (accuracy * 100.0))
```

Accuracy: 59.48%

```
In [ ]: 1
```

```
In [ ]: 1
```

## Conclusion

In this work, it was found that in overall, some DRGs do better at meeting these expected GMLOS than others. We found that DRGs with a high expected length of stay are less likely to be met than those with short stays.

insights were given into what the data is made up, which will help healthcare professionals better understand which patients overstay their GMLOS. The visualizations allow for detailed investigation of how different aspects of a patient and their condition contribute to their ability to meet the DRG-defined GMLOS. The machine learning models, can be used to essentially triage new patients that may require more attention to meet the GMLOS. These insights will provide value to both the patient and hospital

```
In [ ]: 1
```

