

# The analysis of occupancy rates of beds in the Shelter Support and Housing Administration division's Shelter Management Information System\*

Fengyuan Tang

27 April 2022

## Abstract

In this paper, my goal is to analyze the variables that have significant effects on the proportion of actual bed capacity that is occupied for the reporting date. It was found that contributing factors include the programs' locations, gender, age, household size of the service user group, the type of overnight service provided, the program area. Furthermore, the number of beds showing as available for occupancy, rooms that a program has been approved to provide, rooms showing as occupied by a shelter user, rooms that are showing as available for occupancy that are not occupied, and rooms that are not currently available also have significant influences on the response variable. This matters because it provides updated information about the ways to effectively improve occupancy rates of beds in the shelter and overnight service programs administered by SSHA.

## 1 Introduction

The Daily Shelter & Overnight Service Occupancy & Capacity provides a list of active overnight shelters and allied services in the Shelter Support and Housing Administration division's Shelter Management Information System database (Toronto Open Data Portal 2021). It provides the daily updated information about shelters and overnight service programs administered by SSHA, including the program's operator, location, classification, occupancy and capacity (Toronto Open Data Portal 2021). This dataset was conducted in 2021 (Toronto Open Data Portal 2021). Based on the current events relating to the daily shelters and overnight service occupancy as well as capacity, we wanted to focus on the occupancy rates of beds in this system.

In the data section, I would examine all variables and possible data sets that are similar. I would also contain graphs which help people to understand what the variables look like. Furthermore, I would convey the features of this data, including summary statistics and relationships between the variables. In the model section, I conduct some explanatory data analysis as well as conduct the EDA. This includes the numerical Summaries and the graphical summaries. The first step is choosing a starting model, which is also called the full model. The second step is to ensure there is no multicollinearity in the model, so I build new models. Then, I check the 2 conditions to make sure that i can use residual plots to analyze the model. After that, I use residual plots to identify potential violations against model assumptions (linearity, normality, constant variance, and uncorrelatedness). The next step is to explore model transformations to correct assumption violations and fit a new model with transformed variables. When I begin to reduce the model, I conduct automated selection and manual selection. Model comparisons determine which model is better so I use ANOVA to compare these models. Then, I apply the diagnostic plots on both the manual reduced model and the auto reduced model. The results section displays the findings, including summary statistics, tables, graphs, images and statistical analysis.

The results interpret that contributing factors include the city of the location of the program (LOCATION\_CITY),the gender, age, household size of the service user group (SECTOR), the type of overnight

---

\*Code and data are available at: <https://github.com/FengyuanTang/Sta304finalpaper>

service provided (OVERNIGHT\_SERVICE\_TYPE), the PROGRAM\_AREA, which includes base shelter and overnight services system, or a temporary response program. Furthermore, the number of beds showing as available for occupancy (T\_CAPACITY\_ACTUAL\_BED), rooms that a program has been approved to provide (T\_CAPACITY\_FUNDING\_BED), rooms showing as occupied by a shelter user (T\_OCCUPIED\_BEDS), rooms that are showing as available for occupancy that are not occupied as of the occupancy date (T\_UNOCCUPIED\_BEDS), rooms that are not currently available in a program (T\_UNAVAILABLE\_BEDS) also have significant impacts on the proportion of actual bed capacity that is occupied for the reporting date. Moreover, I would include discussions of some interesting points and weaknesses of our paper.

The importance is that, as people have been impacted by COVID-19, we believe that there would be some important changes in the daily shelters and overnight service programs. Our aim is to analyze the factors that have significant effects on the proportion of actual bed capacity that is occupied for the reporting date. Therefore, I build models to find its important contributing factors. In this way, it promotes the efficiency of overnight shelters and allied services in the system. Furthermore, it could development government policies and make the society more satisfied. Moreover, there would be an improvement of living standards for those people who have a demand for shelters and overnight services.

## 2 Data

In this paper, we focused on analyzing the contributing factors of the proportion of actual bed capacity that is occupied for the reporting date. We used R programming language (R Core Team 2020) tidyverse (Wickham et al. 2019), janitor (Firke 2021), readxl (Wickham and Bryan 2019), knitr (Xie 2021), ggplot2 (Wickham 2016), dplyr (Wickham et al. 2021), patchwork (Pedersen 2020), car (Fox and Weisberg 2019) and readr (Wickham, Hester, and Bryan 2022).

The data-set is called “Daily Shelter & Overnight Service Occupancy & Capacity” (Toronto Open Data Portal 2021). The variables I used include LOCATION\_CITY, which is the city of the location of the program (Toronto Open Data Portal 2021). SETOR is defined as the means of categorizing homeless shelters based on the gender, age and household size of the service user group(s) served at the shelter location (Toronto Open Data Portal 2021). PROGRAM\_MODEL is a classification of shelter programs as either Emergency or Transitional (Toronto Open Data Portal 2021). OVERNIGHT\_SERVICE\_TYPE identifies the type of overnight service being provided (Toronto Open Data Portal 2021). PROGRAM\_AREA indicates whether the program is part of the base shelter and overnight services system, or is part of a temporary response program (Toronto Open Data Portal 2021). CAPACITY\_ACTUAL\_BED shows the number of beds showing as available for occupancy in the Shelter Management Information System (Toronto Open Data Portal 2021). CAPACITY\_FUNDING\_BED displays the number of beds that a program has been approved to provide (Toronto Open Data Portal 2021). OCCUPIED\_BEDS illustrates the number of beds showing as occupied by a shelter user in the Shelter Management Information System for this program for this date (Toronto Open Data Portal 2021). UNOCCUPIED\_BEDS is the number of beds that are showing as available for occupancy that are not occupied as of the occupancy date (Toronto Open Data Portal 2021). This is calculated as CAPACITY\_ACTUAL\_BED minus OCCUPIED\_BEDS (Toronto Open Data Portal 2021). UNAVAILABLE\_BEDS shows the number of beds that are not currently available in a program (Toronto Open Data Portal 2021). Specifically, this is calculated as CAPACITY\_FUNDING\_BED minus CAPACITY\_ACTUAL\_BED (Toronto Open Data Portal 2021). The response variable is called OCCUPANCY\_RATE\_BEDS, which displays the proportion of actual bed capacity that is occupied for the reporting date (Toronto Open Data Portal 2021).

Since the data-set is not tidy, I filter the missing values in some important variables and only keep the data that is “Bed Based Capacity”. Specifically, the variable CAPACITY\_TYPE is defined as whether the capacity for this program is measured in rooms or beds (Toronto Open Data Portal 2021). The project creates a histogram for the proportion of actual bed capacity that is occupied for the reporting date. This is calculated as OCCUPIED\_BEDS divided by CAPACITY\_ACTUAL\_BED (Toronto Open Data Portal 2021). I also make plots about some possible predictors, including CAPACITY\_ACTUAL\_BED, CAPACITY\_FUNDING\_BED, OCCUPIED\_BEDS, and UNOCCUPIED\_BEDS. Then, I choose my

starting model by using the results of EDA and my common sense. Drawing scatter-plots between  $y_i$  and  $\hat{y}_i$  and that between numerical predictors can be used to check Condition 1 and 2. The Residual vs. Fitted, Residual vs. Predictors and Residual QQ Plot can decide whether each regression modelling assumption is satisfied. Since power transform fails to work if variable contains 0, we add 0.0000001 instead. After applying the box-cox transformation, I use mutate to create transformed variables and fit a new model. This is called candidate model 1, and model comparisons can determine which model is better. Same methods are applied, including checking two conditions, creating residual plots, and so on.

## 2.1 Summary statistics

It shows that the proportion of actual bed capacity has a minimum of 7.14, and a maximum of 100. The average is 95.62. The number of rooms showing as available for occupancy has a minimum of 1, and a maximum of 235. The average is 34.73. The number of rooms that a program is has been approved to provide has a minimum of 2, and a maximum of 235. The average is 36.5. The number of rooms showing as occupied has a minimum of 1, and a maximum of 234. The average is 33.85. The number of beds that are showing as available for occupancy that are not occupied has a minimum of 0, and a maximum of 39. The average is 0.88. The number of beds that are not currently available in a program has a minimum of 0, and a maximum of 180.

```
##   Min. 1st Qu. Median Mean 3rd Qu. Max.
## 7.14  96.00 100.00 95.65 100.00 100.00
##   Min. 1st Qu. Median Mean 3rd Qu. Max.
## 1.00  17.00 27.00 34.72 45.00 235.00
##   Min. 1st Qu. Median Mean 3rd Qu. Max.
## 2.00  19.00 28.00 36.48 50.00 235.00
##   Min. 1st Qu. Median Mean 3rd Qu. Max.
## 1.00  16.00 25.00 33.85 44.00 234.00
##   Min. 1st Qu. Median Mean 3rd Qu. Max.
## 0.0000 0.0000 0.0000 0.8793 1.0000 39.0000
##   Min. 1st Qu. Median Mean 3rd Qu. Max.
## 0.000 0.000 0.000 1.755 1.000 180.000
```

## 2.2 Exploratory Data Analysis

I should conduct the exploratory data analysis in order to build the starting model. Thus, Figure 1 is a histogram of the response variable. The graph is extremely right skewed. It will be much better if it is normal distribution. This problem can be fixed later by using model transformation.

From Figure 2, Figure 3, Figure 4, and Figure 5, I infer that variables “CAPACITY\_ACTUAL\_BED”, “CAPACITY\_FUNDING\_BED”, “OCCUPIED\_BEDS” and “UNOCCUPIED\_BEDS” could have influence on the response variable. Therefore, I include them in the starting model.

Then, I create bar plots for some categorical variables

Figure 6 and Figure 7 show that the data volume is unbalanced. Thus, the credibility is weak and we do not need to include them.

From the website, PROGRAM\_AREA displays whether the program is part of the base shelter and overnight services system, or is part of a temporary response program (Toronto Open Data Portal 2021). Figure 8 illustrates that most programs are in the type of Base Shelter and Overnight Services System. These programs are intended to be regular, year-round, and permanent (Toronto Open Data Portal 2021). Also, the variable OVERNIGHT\_SERVICE\_TYPE shows the type of overnight service provided (Toronto Open Data Portal 2021). This includes Shelter, 24-Hour Respite, Motel/Hotel, Interim Housing, Warming Center, 24-Hour

the proportion of actual bed capacity that is occupied

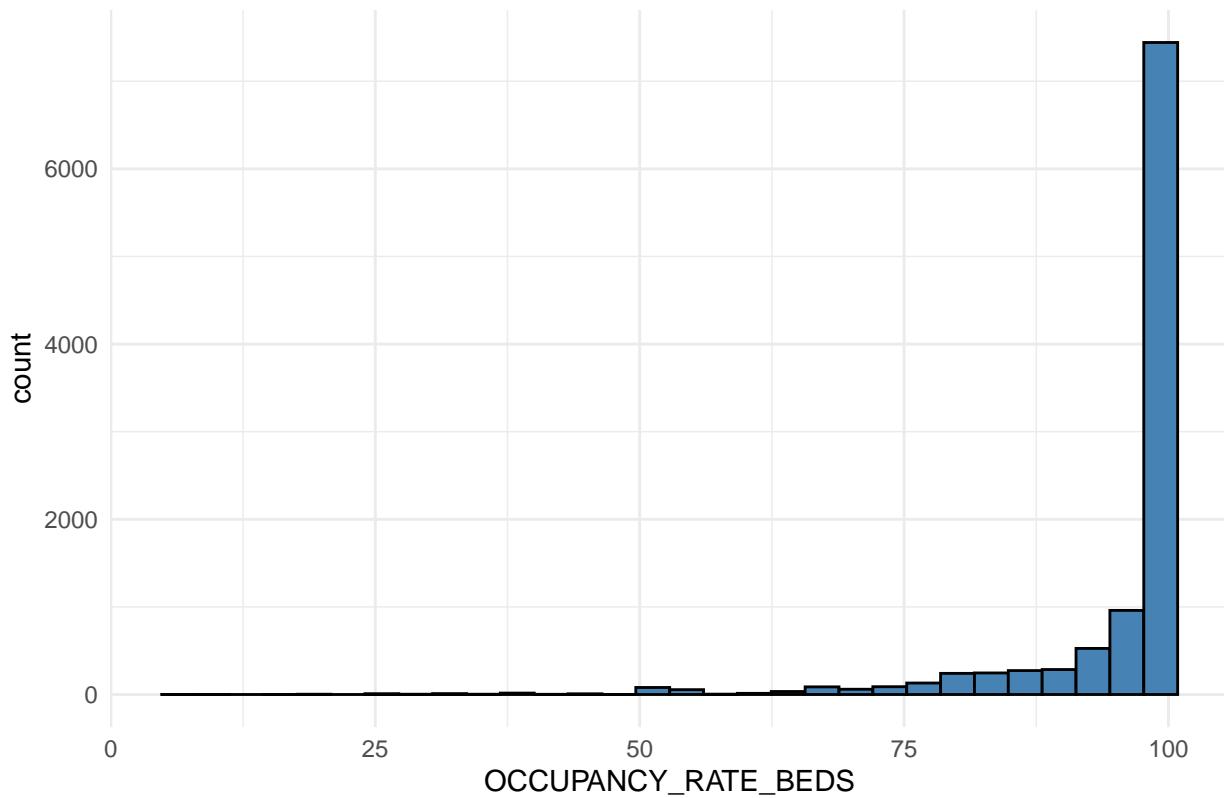


Figure 1: the proportion of actual bed capacity that is occupied for the reporting date

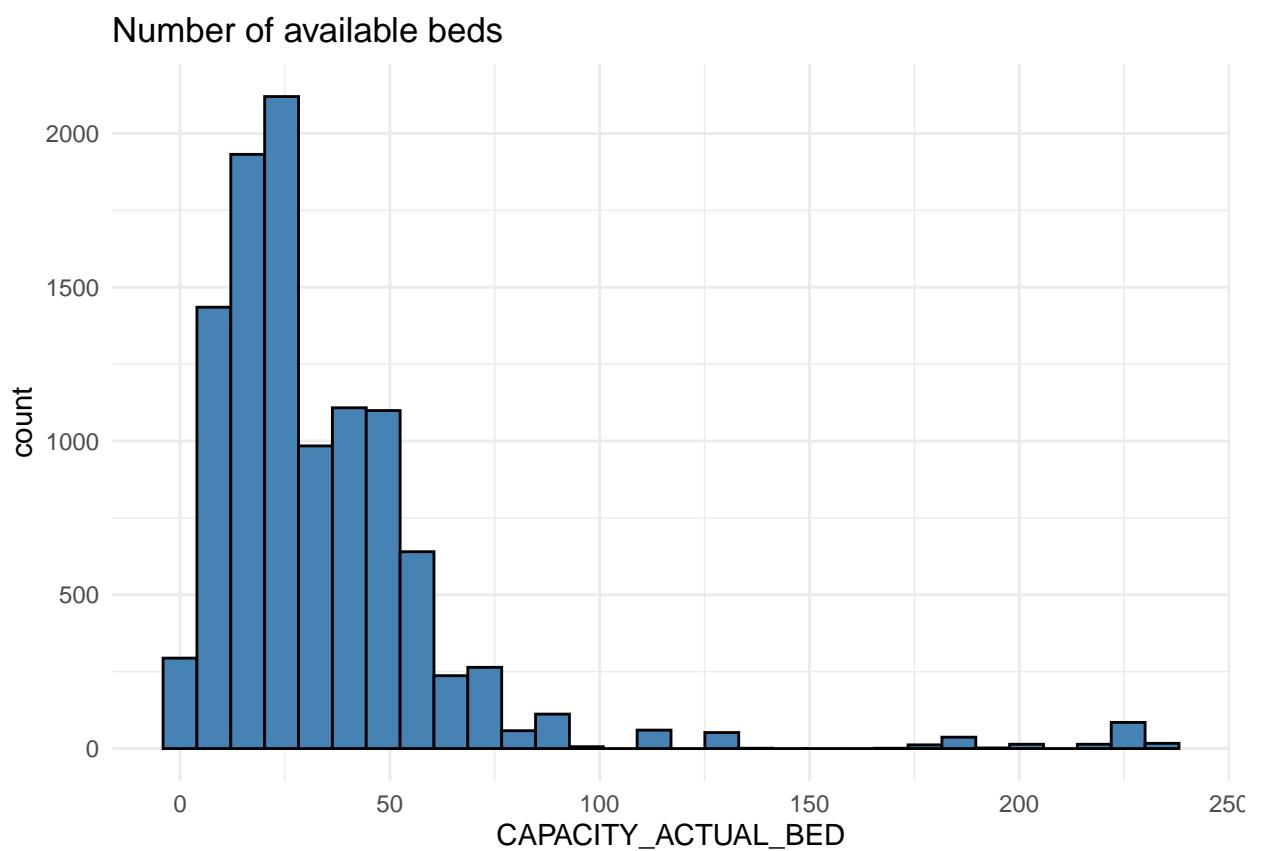


Figure 2: the number of beds showing as available for occupancy in the Shelter Management Information System

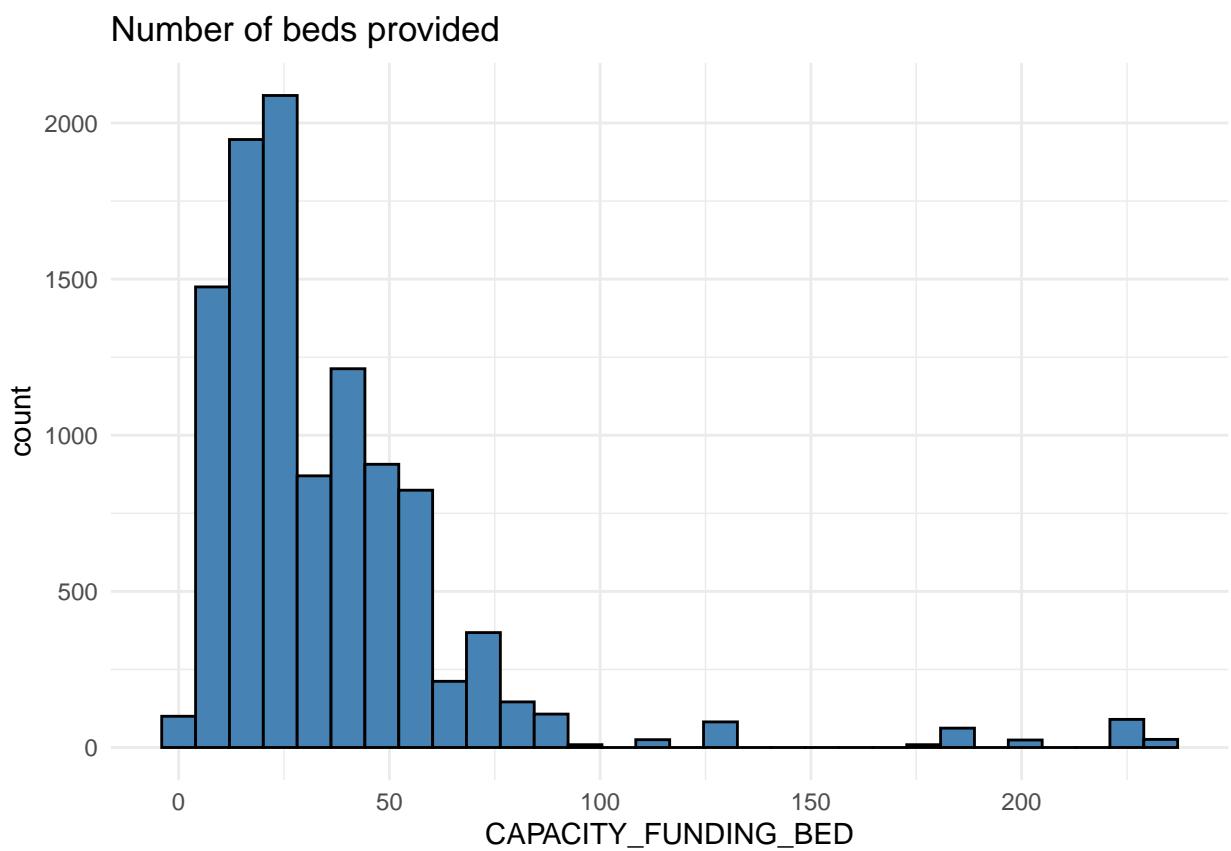


Figure 3: the number of beds that a program has been approved to provide

### Occupied beds and occupied rates

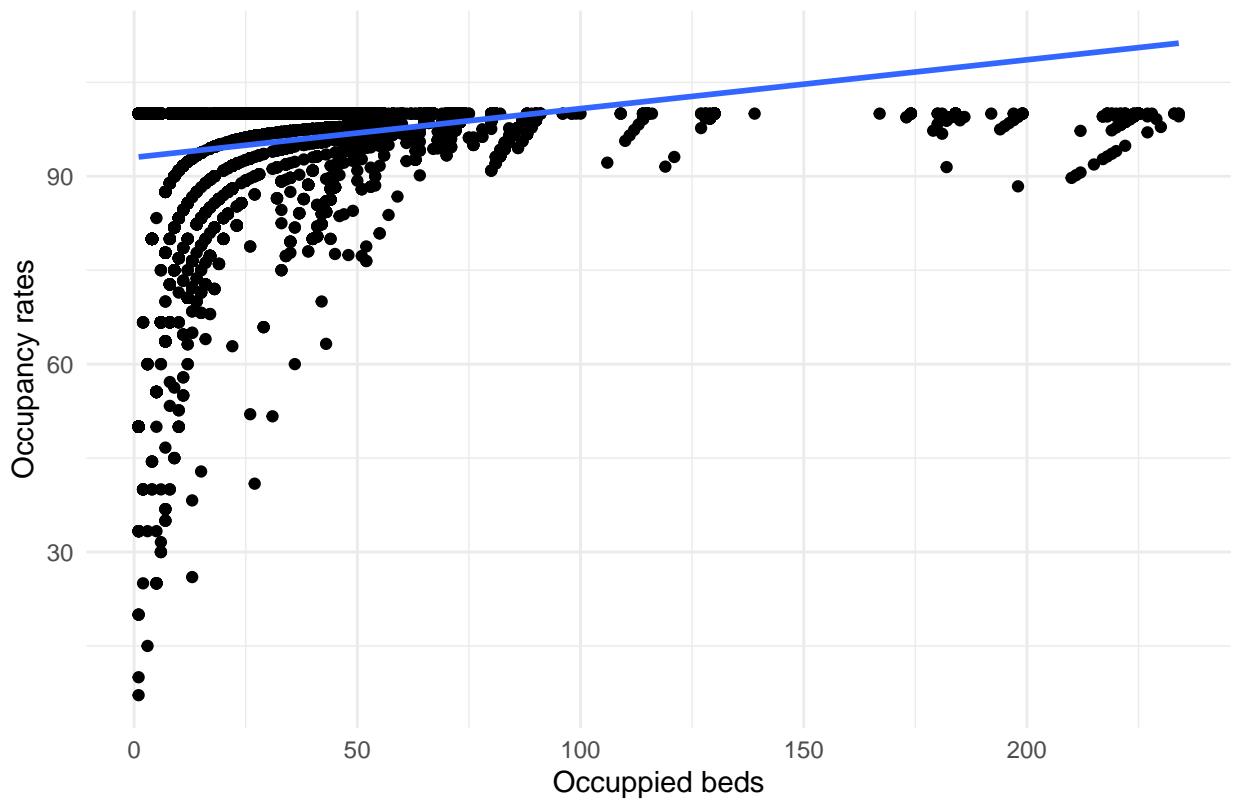


Figure 4: Occupied beds vs. occupied rates

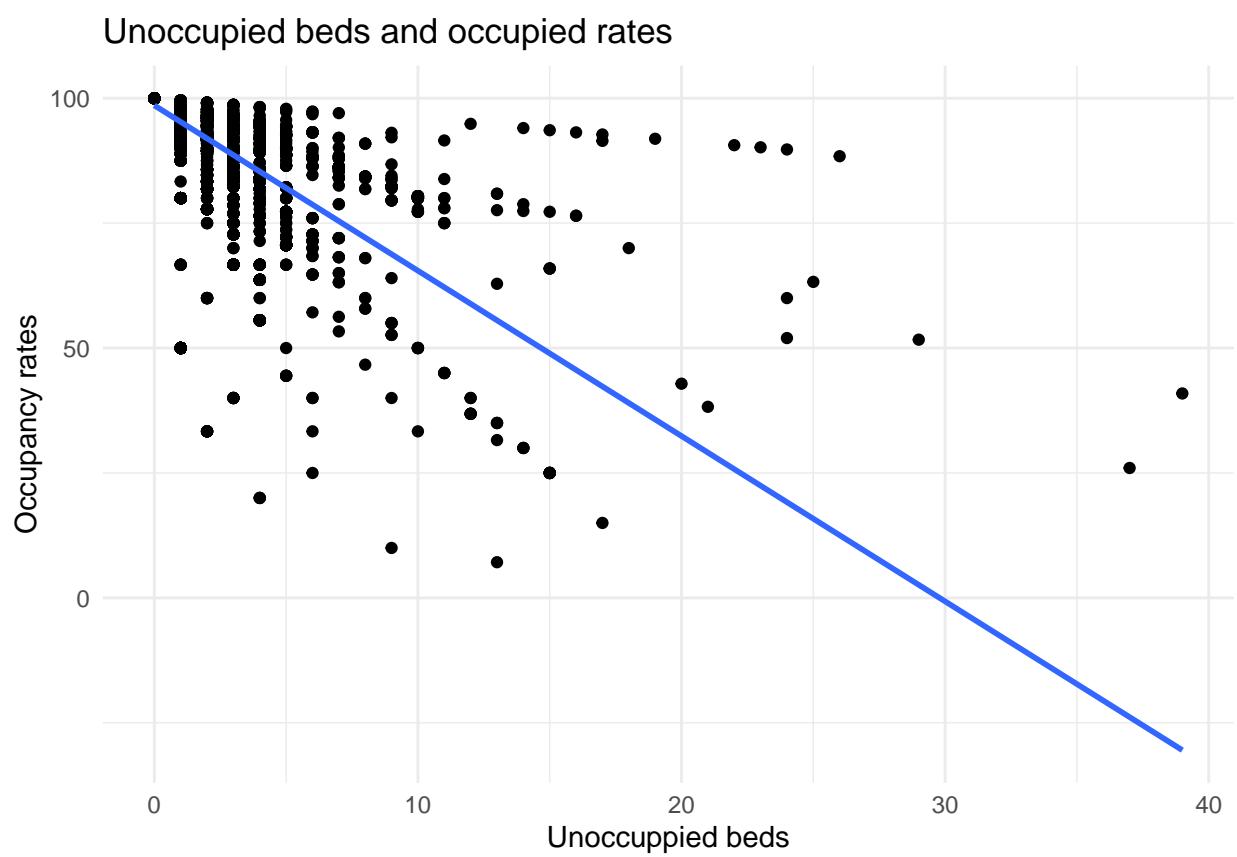


Figure 5: Unoccupied beds vs. occupied rates

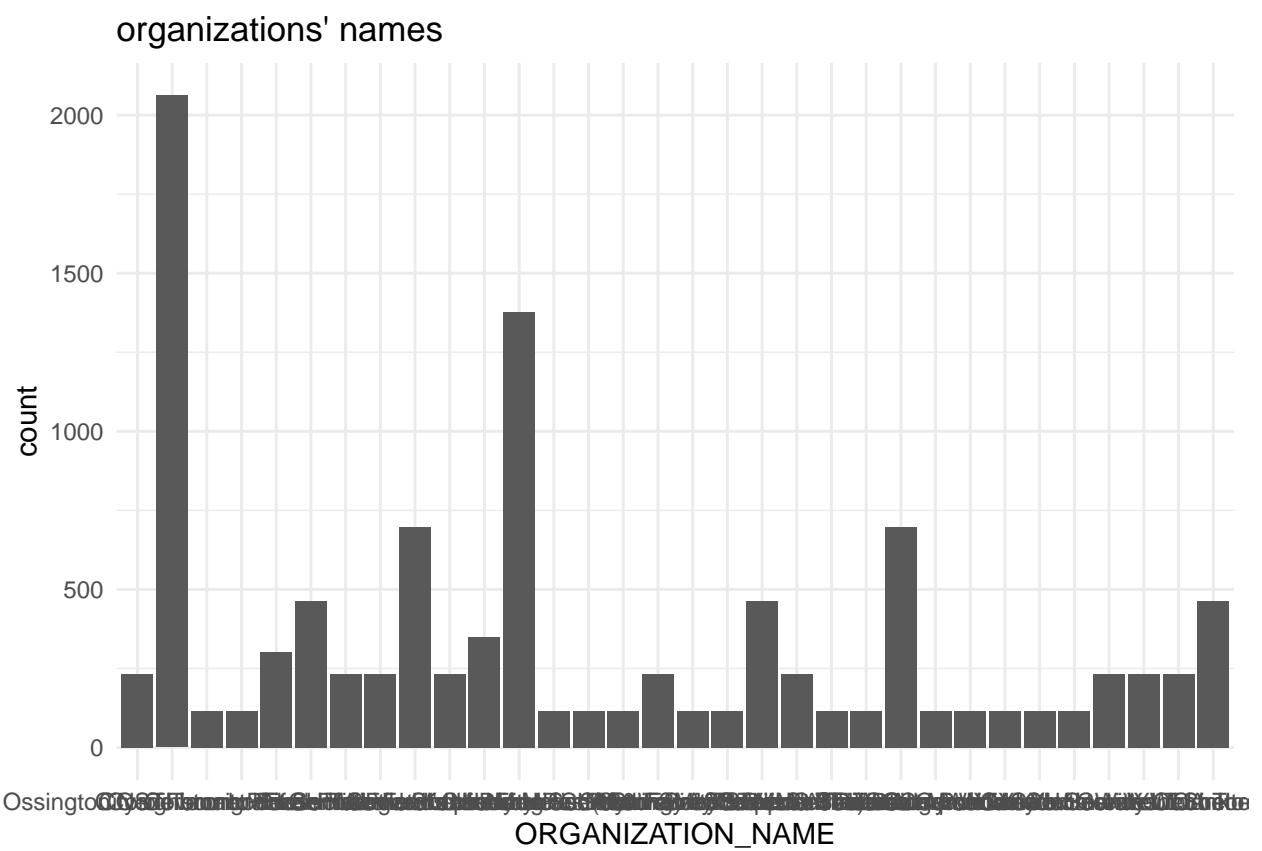


Figure 6: Name of the organization providing the overnight service

### the city of the location



Figure 7: The name of the location of the program

### plots of the program area

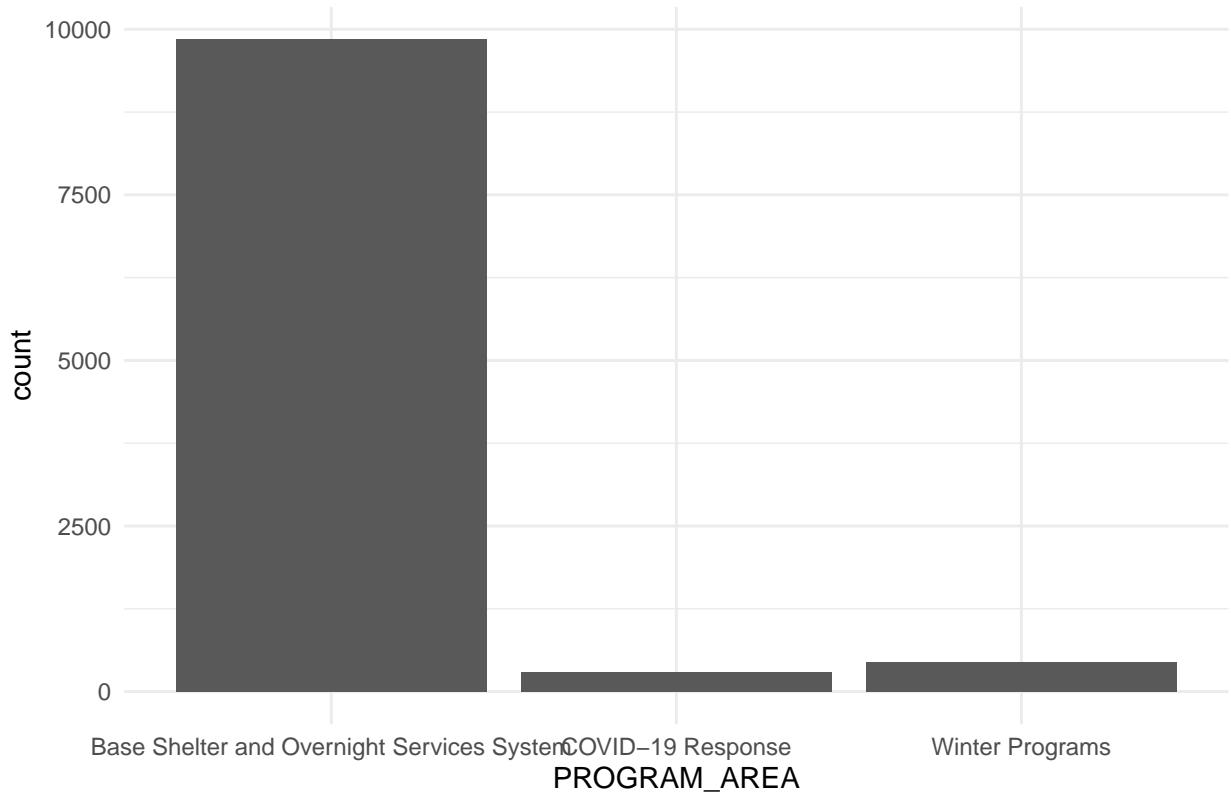


Figure 8: whether the program is part of the base shelter and overnight services system, or is part of a temporary response program

the type of overnight service being provided

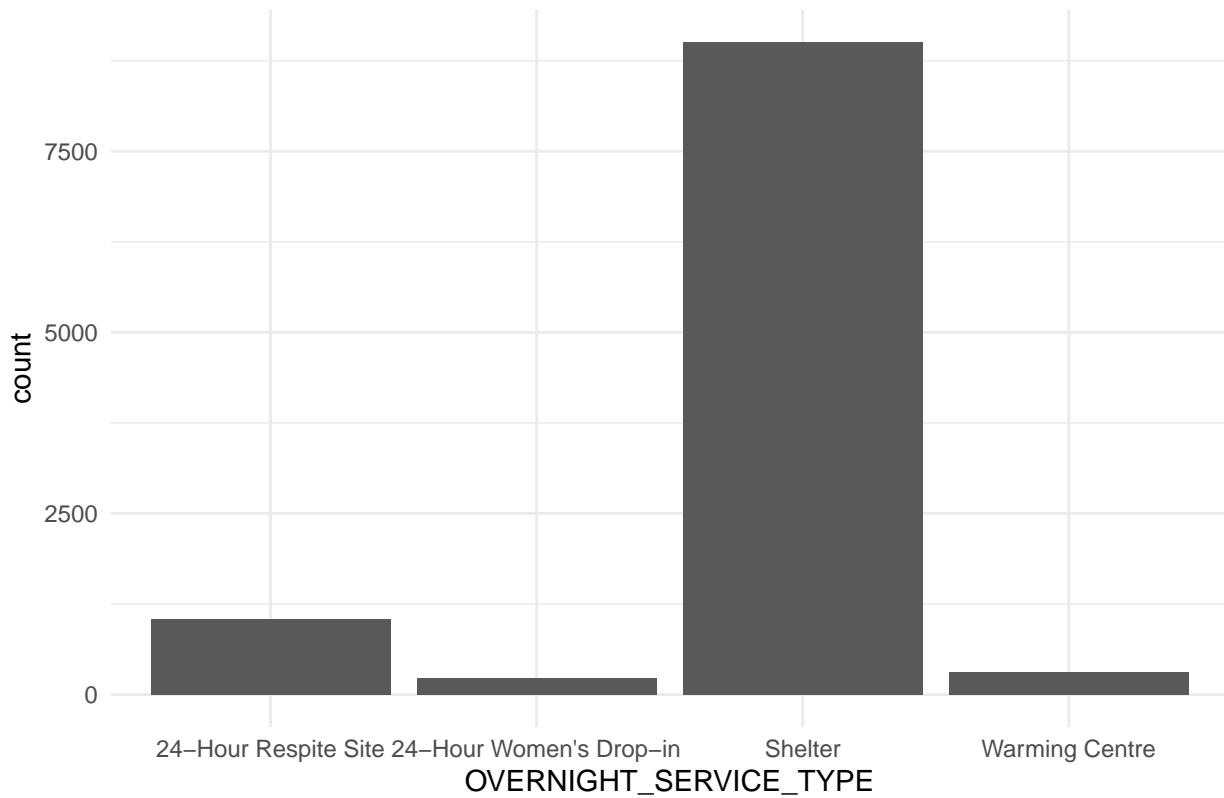


Figure 9: the type of overnight service being provided

Woman's Drop-in, Isolation/Recovery Site (Toronto Open Data Portal 2021). From Figure 9, it shows that shelters are the most type of overnight services provided in this system.

### 3 Model

The second step is to create my starting model by using common sense and results of EDA. Also, predictors can't have overlapping information.

```
##  
## Call:  
## lm(formula = OCCUPANCY_RATE_BEDS ~ LOCATION_CITY + SECTOR + PROGRAM_MODEL +  
##      OVERNIGHT_SERVICE_TYPE + PROGRAM_AREA + CAPACITY_ACTUAL_BED +  
##      CAPACITY_FUNDING_BED + OCCUPIED_BEDS + UNOCCUPIED_BEDS +  
##      UNAVAILABLE_BEDS, data = data_clean)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -64.094  -0.542    0.842   2.656  66.511  
##  
## Coefficients: (2 not defined because of singularities)  
##                                     Estimate Std. Error t value  
## (Intercept)                   97.179334  1.028510 94.486  
## LOCATION_CITYNorth York      -4.398259  0.603590 -7.287  
## LOCATION_CITYScarborough     -2.792785  0.557845 -5.006  
## LOCATION_CITYToronto          -3.532451  0.457137 -7.727  
## SECTORMen                     3.434248  0.874410  3.928  
## SECTORMixed Adult            3.261882  0.891226  3.660  
## SECTORWomen                   4.417486  0.882677  5.005  
## SECTORYouth                   0.856632  0.881479  0.972  
## PROGRAM_MODELTransitional   -1.591021  0.160337 -9.923  
## OVERNIGHT_SERVICE_TYPE24-Hour Women's Drop-in  0.635646  0.530108  1.199  
## OVERNIGHT_SERVICE_TYPEShelter  -0.465552  0.275149 -1.692  
## OVERNIGHT_SERVICE_TYPEWarming Centre  -0.885442  0.715161 -1.238  
## PROGRAM_AREACOVID-19 Response  -2.776146  0.508910 -5.455  
## PROGRAM_AREAWinter Programs   0.287212  0.622278  0.462  
## CAPACITY_ACTUAL_BED          -3.187474  0.034661 -91.961  
## CAPACITY_FUNDING_BED          0.011108  0.008811  1.261  
## OCCUPIED_BEDS                 3.249190  0.033637  96.596  
## UNOCCUPIED_BEDS                NA          NA          NA  
## UNAVAILABLE_BEDS                NA          NA          NA  
##  
## (Intercept)                  < 2e-16 ***  
## LOCATION_CITYNorth York      3.40e-13 ***  
## LOCATION_CITYScarborough     5.64e-07 ***  
## LOCATION_CITYToronto          1.20e-14 ***  
## SECTORMen                     8.64e-05 ***  
## SECTORMixed Adult            0.000253 ***  
## SECTORWomen                   5.69e-07 ***  
## SECTORYouth                   0.331167  
## PROGRAM_MODELTransitional   < 2e-16 ***  
## OVERNIGHT_SERVICE_TYPE24-Hour Women's Drop-in  0.230521  
## OVERNIGHT_SERVICE_TYPEShelter  0.090677 .  
## OVERNIGHT_SERVICE_TYPEWarming Centre  0.215707  
## PROGRAM_AREACOVID-19 Response  5.01e-08 ***
```

```

## PROGRAM_AREAWinter Programs          0.644415
## CAPACITY_ACTUAL_BED                < 2e-16 ***
## CAPACITY_FUNDING_BED               0.207428
## OCCUPIED_BEDS                     < 2e-16 ***
## UNOCCUPIED_BEDS                   NA
## UNAVAILABLE_BEDS                  NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.5 on 10394 degrees of freedom
##   (173 observations deleted due to missingness)
## Multiple R-squared:  0.5499, Adjusted R-squared:  0.5492
## F-statistic: 793.8 on 16 and 10394 DF, p-value: < 2.2e-16

```

Here are the two conditions we need to check before assessing the model assumptions. I want to make sure that we can use residual plots to analyze the model.

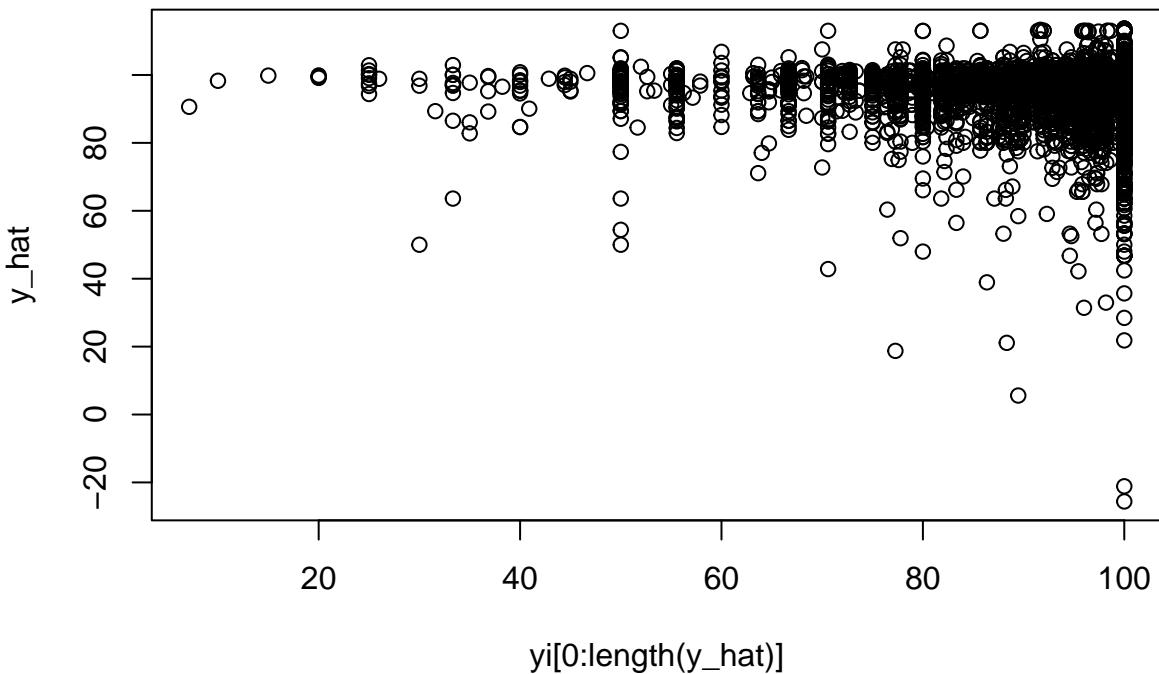


Figure 10: Condition 1: draw a scatter plot between  $y_i$  and  $y_{\text{hat}}$

Now, I check the two conditions, which are Figure 10 and Figure 11. I draw a scatter-plot between  $Y_i$  and  $Y_{\text{hat}}$  to check condition 1. It shows that there is a strong pattern between them. Thus, Condition 1 is satisfied. Next, we draw scatter plots between numerical predictors. This graph displays that there is no or linear relationship between these predictors. Therefore, Condition 2 is satisfied.

After that, I use plots to identify potential violations against model assumptions, which are linearity, normality, constant variance, and uncorrelatedness. Figure 12 is Residual vs. Fitted.

In this graph, the linearity holds, but the independence and constant variance can be improved.

Figure 13 is Residual vs. Predictors.

Figure 14 is Normal Quantile-Quantile (QQ) plots.

To summarize, the linearity should hold. However, the constant variance, independence and normality may be violated and can be improved by using model transformations. There is a severe deviation in the Normal

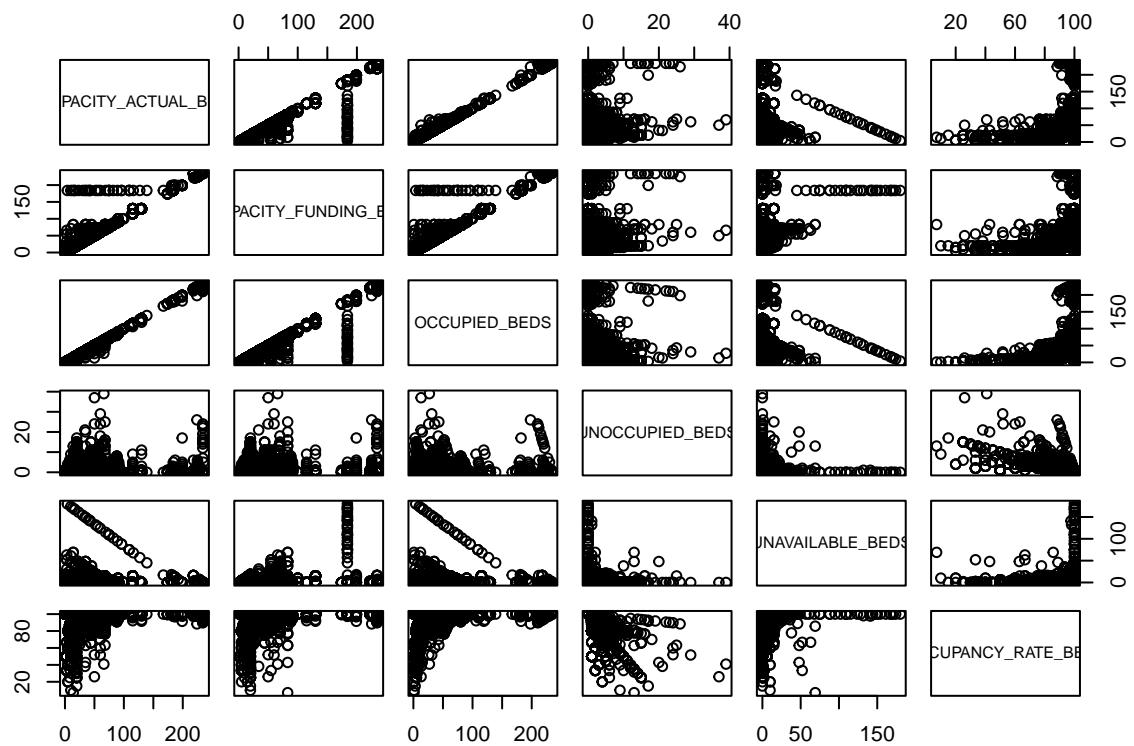


Figure 11: Condition 2: draw scatter plots between numerical predictors

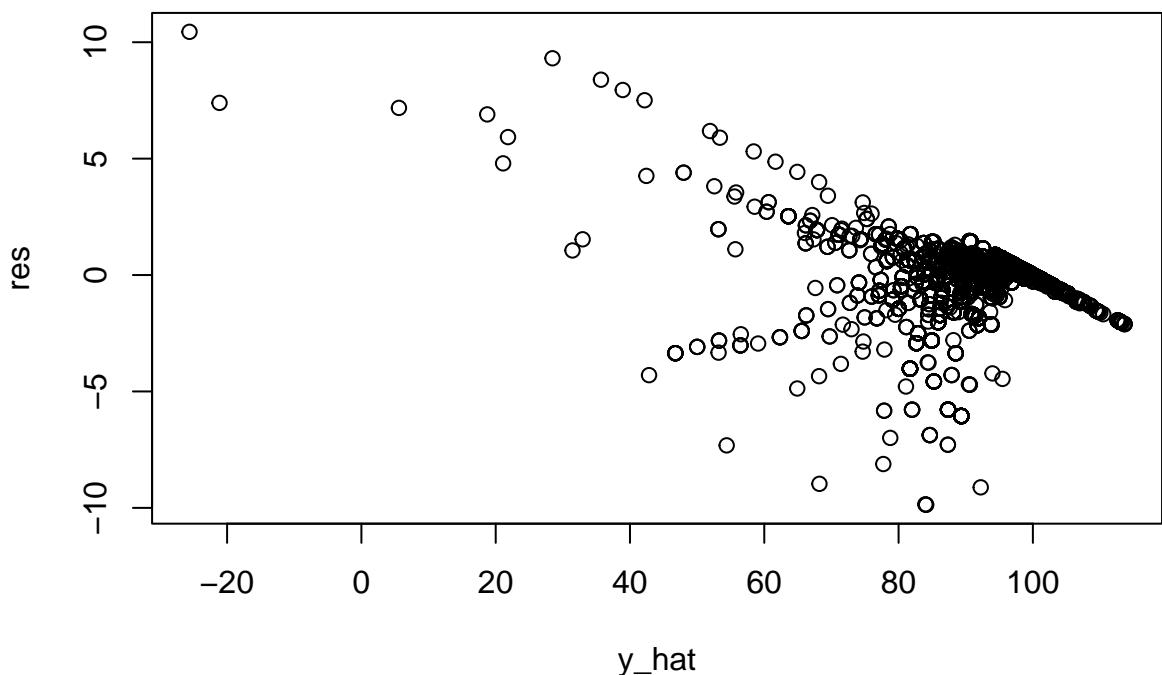


Figure 12: Residual vs. Fitted.

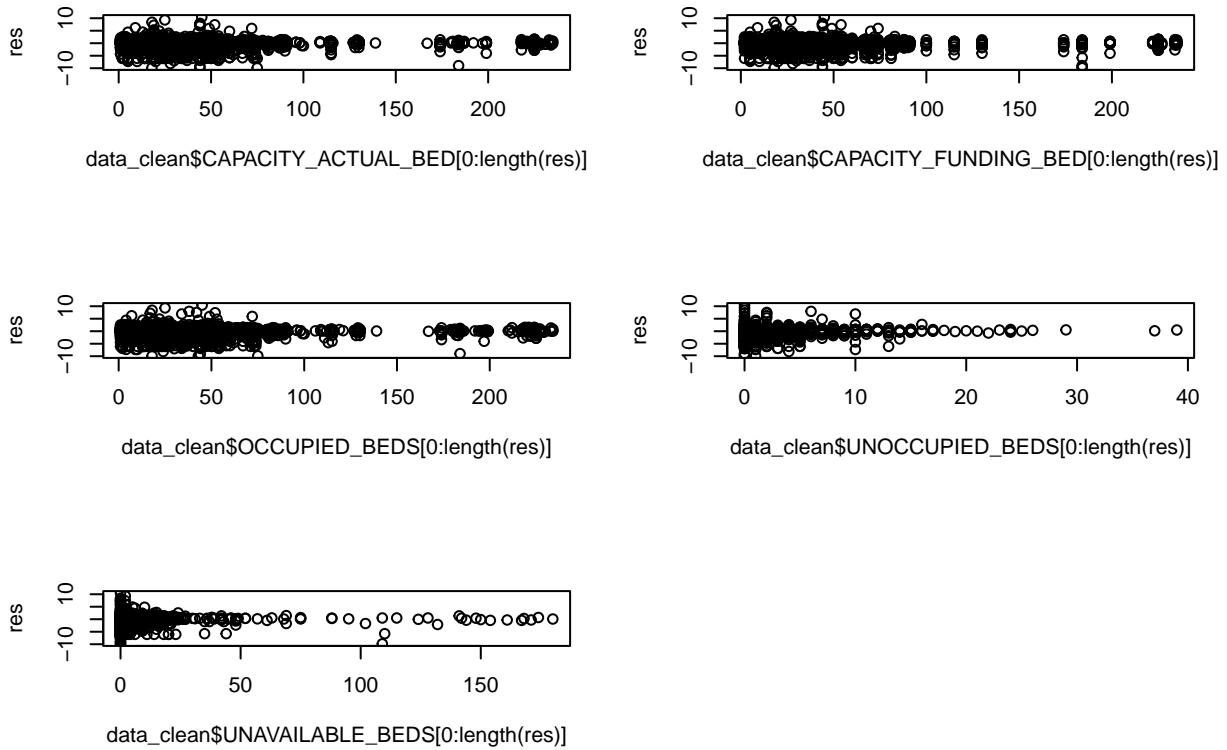


Figure 13: Residual vs. Predictors.

### Normal Q-Q Plot

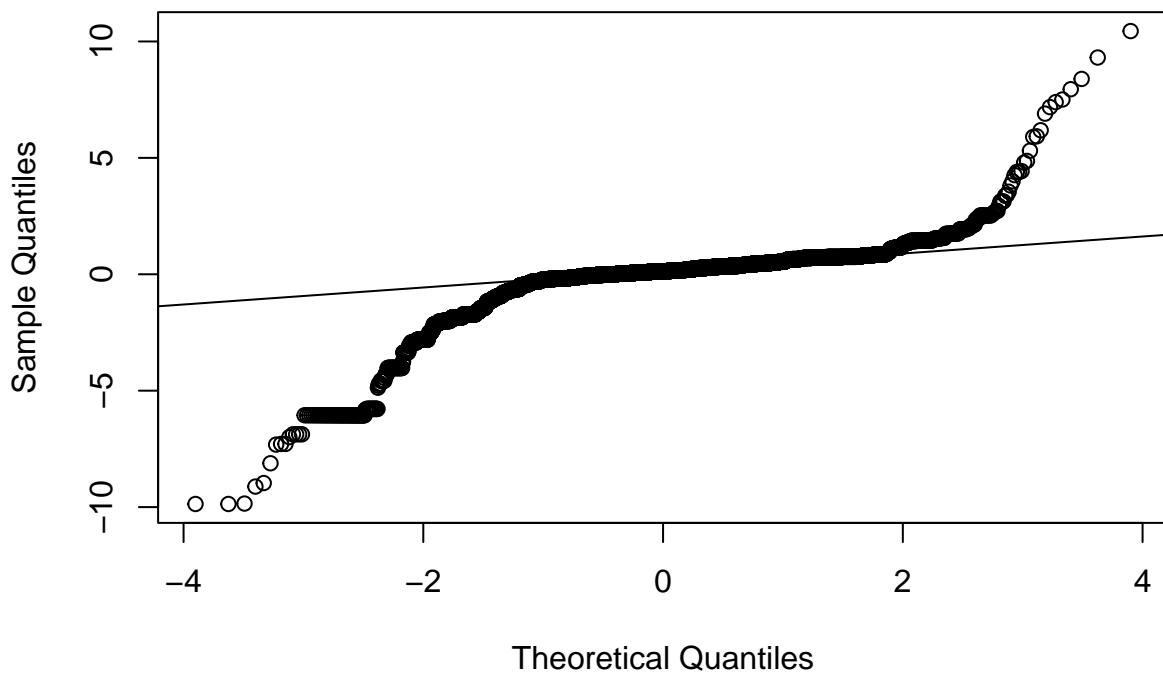


Figure 14: Normal Quantile-Quantile (QQ) plots.

QQ plot.

The next step is to explore model transformations to correct assumption violations. Since the power transform fails to work if any variable contains 0, one way to fix this problem is to add 0.0000001 to this variable. After that, I apply box-cox transformation for numerical variables.

```
## bcPower Transformations to Multinormality
##   Est Power Rounded Pwr Wald Lwr Bnd Wald Upr Bnd
## Y1    0.4592      0.46     0.4512     0.4671
## Y2    0.3758      0.38     0.3667     0.3850
## Y3    0.4711      0.47     0.4632     0.4789
## Y4   -0.1145     -0.11    -0.1189    -0.1101
## Y5   -0.2001     -0.20    -0.2051    -0.1951
##
## Likelihood ratio test that transformation parameters are equal to 0
## (all log transformations)
##                               LRT df      pval
## LR test, lambda = (0 0 0 0 0) 25645.78 5 < 2.22e-16
##
## Likelihood ratio test that no transformations are needed
##                               LRT df pval
## LR test, lambda = (1 1 1 1 1) NaN 5 NA
```

Then, I create the transformed variables.

The next step is to fit a new model with transformed variables.

```
##
## Call:
## lm(formula = OCCUPANCY_RATE_BEDS ~ LOCATION_CITY + SECTOR + PROGRAM_MODEL +
##     OVERNIGHT_SERVICE_TYPE + PROGRAM_AREA + T_CAPACITY_ACTUAL_BED +
##     T_CAPACITY_FUNDING_BED + T_OCCUPIED_BEDS + T_UNOCCUPIED_BEDS +
##     T_UNAVAILABLE_BEDS, data = T_data_clean)
##
## Residuals:
##   Min     1Q Median     3Q    Max
## -40.895 -0.689  0.302  1.619  55.779
##
## Coefficients:
##                               Estimate Std. Error t value
## (Intercept)                94.243701  0.643430 146.471
## LOCATION_CITYNorth York   -2.920369  0.341893 -8.542
## LOCATION_CITYScarborough -1.207561  0.319442 -3.780
## LOCATION_CITYToronto       -1.578499  0.259331 -6.087
## SECTORMen                  3.163350  0.501824  6.304
## SECTORMixed Adult          2.740196  0.508180  5.392
## SECTORWomen                 3.757262  0.505755  7.429
## SECTORYouth                 2.470550  0.505417  4.888
## PROGRAM_MODELTransitional  0.037263  0.094898  0.393
## OVERNIGHT_SERVICE_TYPE24-Hour Women's Drop-in  1.152000  0.300417  3.835
## OVERNIGHT_SERVICE_TYPEShelter -0.162620  0.156615 -1.038
## OVERNIGHT_SERVICE_TYPEWarming Centre  0.688733  0.407093  1.692
## PROGRAM_AREACOVID-19 Response -3.273190  0.276092 -11.855
## PROGRAM_AREAWinter Programs   0.265448  0.354548  0.749
## T_CAPACITY_ACTUAL_BED        -46.579215 0.305245 -152.596
## T_CAPACITY_FUNDING_BED       2.630519  0.189755  13.863
```

```

## T_OCCUPIED_BEDS          43.428823   0.262910  165.185
## T_UNOCCUPIED_BEDS        0.399050   0.021101   18.911
## T_UNAVAILABLE_BEDS       0.014135   0.003576    3.953
##
## Pr(>|t|)
## (Intercept)                < 2e-16 ***
## LOCATION_CITYNorth York      < 2e-16 ***
## LOCATION_CITYScarborough     0.000158 ***
## LOCATION_CITYToronto         1.19e-09 ***
## SECTORMen                   3.02e-10 ***
## SECTORMixed Adult           7.11e-08 ***
## SECTORWomen                  1.18e-13 ***
## SECTORYouth                  1.03e-06 ***
## PROGRAM_MODELTransitional    0.694581
## OVERNIGHT_SERVICE_TYPE24-Hour Women's Drop-in 0.000126 ***
## OVERNIGHT_SERVICE_TYPEShelter 0.299134
## OVERNIGHT_SERVICE_TYPEWarming Centre 0.090707 .
## PROGRAM_AREACOVID-19 Response < 2e-16 ***
## PROGRAM_AREAWinter Programs 0.454059
## T_CAPACITY_ACTUAL_BED       < 2e-16 ***
## T_CAPACITY_FUNDING_BED       < 2e-16 ***
## T_OCCUPIED_BEDS              < 2e-16 ***
## T_UNOCCUPIED_BEDS             < 2e-16 ***
## T_UNAVAILABLE_BEDS            7.77e-05 ***
##
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.679 on 10392 degrees of freedom
##   (173 observations deleted due to missingness)
## Multiple R-squared:  0.8559, Adjusted R-squared:  0.8556
## F-statistic:  3428 on 18 and 10392 DF,  p-value: < 2.2e-16

```

Now, i begin to reduce our model. Specifically, I apply automated selection because I want to know which one can be removed.

```

## Start:  AIC=27140.85
## OCCUPANCY_RATE_BEDS ~ LOCATION_CITY + SECTOR + PROGRAM_MODEL +
##   OVERNIGHT_SERVICE_TYPE + PROGRAM_AREA + T_CAPACITY_ACTUAL_BED +
##   T_CAPACITY_FUNDING_BED + T_OCCUPIED_BEDS + T_UNOCCUPIED_BEDS +
##   T_UNAVAILABLE_BEDS
##
##                                     Df Sum of Sq   RSS   AIC
## - PROGRAM_MODEL                 1      2 140635 27139
## <none>                           140633 27141
## - T_UNAVAILABLE_BEDS            1     211 140844 27154
## - OVERNIGHT_SERVICE_TYPE        3     402 141035 27165
## - LOCATION_CITY                  3     1049 141682 27212
## - PROGRAM_AREA                   2     1940 142573 27279
## - SECTOR                          4     2140 142773 27290
## - T_CAPACITY_FUNDING_BED         1     2601 143234 27330
## - T_UNOCCUPIED_BEDS             1     4840 145473 27491
## - T_CAPACITY_ACTUAL_BED          1     315120 455753 39380
## - T_OCCUPIED_BEDS                1     369257 509890 40549
##
## Step:  AIC=27139
## OCCUPANCY_RATE_BEDS ~ LOCATION_CITY + SECTOR + OVERNIGHT_SERVICE_TYPE +

```

```

##      PROGRAM_AREA + T_CAPACITY_ACTUAL_BED + T_CAPACITY_FUNDING_BED +
##      T_OCCUPIED_BEDS + T_UNOCCUPIED_BEDS + T_UNAVAILABLE_BEDS
##
##                                Df Sum of Sq     RSS     AIC
## <none>                            140635 27139
## + PROGRAM_MODEL                  1      2 140633 27141
## - T_UNAVAILABLE_BEDS            1      213 140848 27153
## - OVERNIGHT_SERVICE_TYPE        3      401 141036 27163
## - LOCATION_CITY                 3      1063 141698 27211
## - PROGRAM_AREA                  2      1941 142576 27278
## - SECTOR                        4      2146 142781 27289
## - T_CAPACITY_FUNDING_BED       1      2604 143239 27328
## - T_UNOCCUPIED_BEDS            1      5215 145850 27516
## - T_CAPACITY_ACTUAL_BED        1      315557 456192 39388
## - T_OCCUPIED_BEDS              1      369397 510032 40550

```

It shows the predictors that are removed, which is “PROGRAM\_MODEL”. The AIC in the last model is 26000.22.

```

## Analysis of Variance Table
##
## Model 1: OCCUPANCY_RATE_BEDS ~ LOCATION_CITY + SECTOR + OVERNIGHT_SERVICE_TYPE +
##           PROGRAM_AREA + T_CAPACITY_ACTUAL_BED + T_CAPACITY_FUNDING_BED +
##           T_OCCUPIED_BEDS + T_UNOCCUPIED_BEDS + T_UNAVAILABLE_BEDS
## Model 2: OCCUPANCY_RATE_BEDS ~ LOCATION_CITY + SECTOR + PROGRAM_MODEL +
##           OVERNIGHT_SERVICE_TYPE + PROGRAM_AREA + T_CAPACITY_ACTUAL_BED +
##           T_CAPACITY_FUNDING_BED + T_OCCUPIED_BEDS + T_UNOCCUPIED_BEDS +
##           T_UNAVAILABLE_BEDS
##   Res.Df     RSS Df Sum of Sq    F Pr(>F)
## 1 10393 140635
## 2 10392 140633  1    2.0865 0.1542 0.6946

```

The P value here is 0.5843, which is large. This means that the auto reduced model is better than the full model. Therefore, I can create diagnostic plots for Auto\_reduced\_model.

## 4 Results

Same as before, Figure 15 and Figure 16 display that both Condition 1 and 2 are held in the reduced model.

From Figure 17 and Figure 18, the linearity and the constant variance should hold. However, the independence and normality may be violated.

The distribution of my response variable is right skewed, which may cause problems. Moreover, it is observed that there is a positive relationship between the number of rooms showing as occupied and the proportion of actual bed capacity that is occupied for the reporting date. It also displays that the number of rooms that are showing as available for occupancy that are not occupied as of the occupancy date has a negative relationship with the proportion of actual bed capacity that is occupied. There is a strong pattern between  $y_i$  and  $\hat{y}_i$ , and there is no/linear relationship between predictors, so the two conditions hold. According to the following graphs, the linearity holds. Regarding to independence, there appear to be some evidence of grouping in residual plots. Constant variance is difficult to tell since some residuals vs. predictors plots contain points that are far away. For normality, there is lifting in the tails. After fitting a new model with transformed variables, we apply automated selection to determine which variables can be removed. To be specific, predictor removed is PROGRAM\_MODEL”.

The results means LOCATION\_CITY, SECTOR, OVERNIGHT\_SERVICE\_TYPE, PROGRAM\_AREA, T\_CAPACITY\_ACTUAL\_BED, T\_CAPACITY\_FUNDING\_BED, T\_OCCUPIED\_BEDS, T\_UNOCCUPIED\_BEDS,

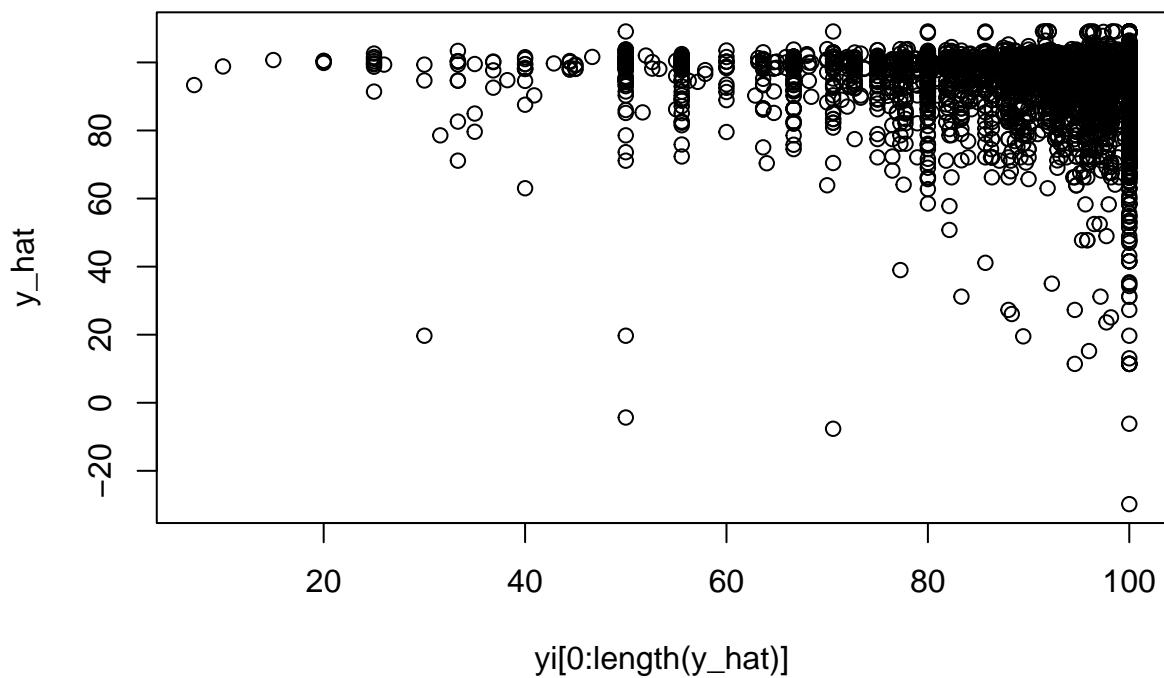


Figure 15: Condition 1: draw a scatter plot between  $\text{yi}$  and  $\text{y\_hat}$ .

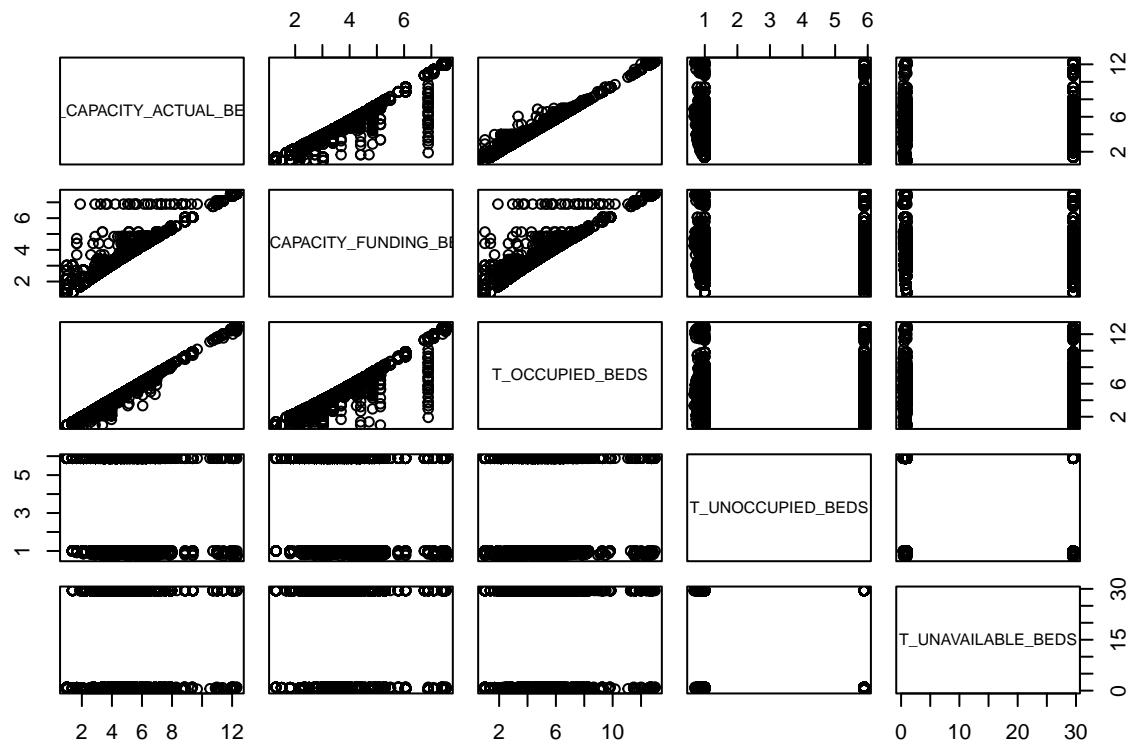


Figure 16: Condition 2: draw scatter plots between numerical predictors

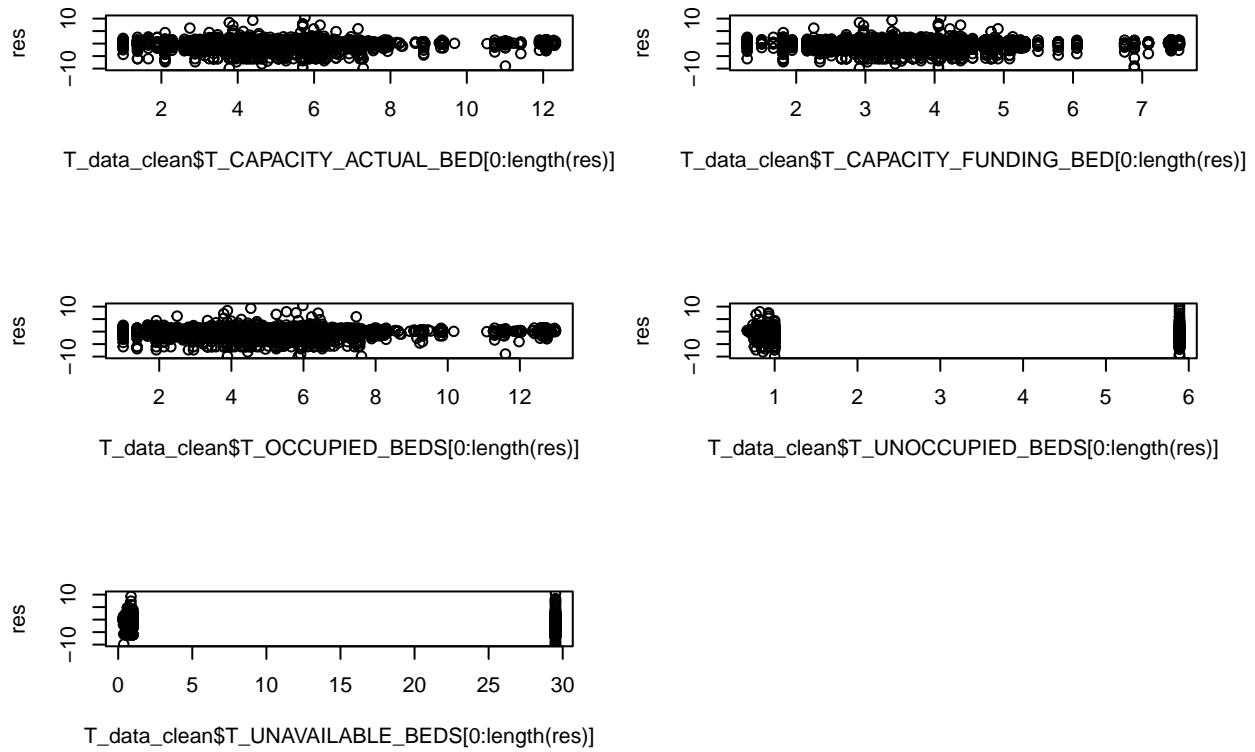


Figure 17: Residual vs. Predictors

### Normal Q-Q Plot

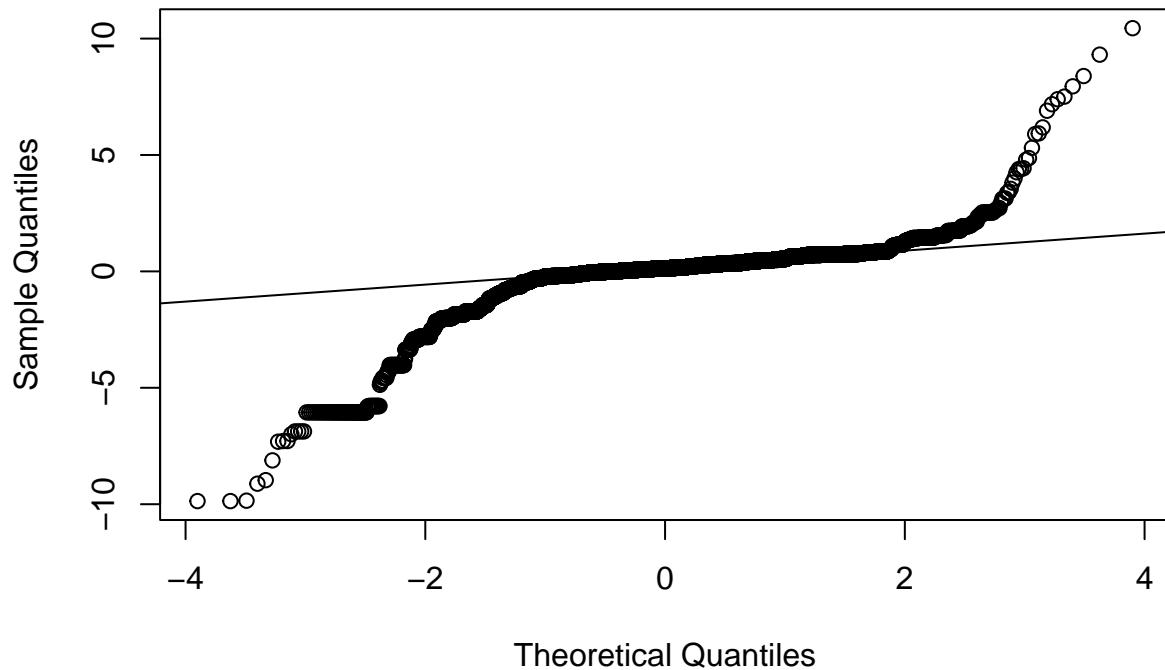


Figure 18: Residual QQ Plot

T\_UNAVAILABLE\_BEDS are important predictors. The Auto reduced model has ANOVA p-value of 0.5843. This means the auto reduced model is better. For both Auto\_reduced\_model, the linearity holds. There appear to be some evidence of grouping. The constant variance can be improved. There is lifting in the tails.

In conclusion, I will state some of the possible ways for the program to increase its beds' occupied rates. This includes increasing the number of beds showing as available for occupancy and the beds showing as occupied by a shelter user. Moreover, it is also useful to decrease the number of beds that are showing as available for occupancy that are not occupied as of the occupancy date, and beds that are not currently available.

## 5 Discussion

### 5.1 First discussion point

In this paper, I aim to analyze the variables that have significant effects on the proportion of actual bed capacity that is occupied for the reporting date. The variables I used include LOCATION\_CITY, SETOR, PROGRAM\_MODEL, OVERNIGHT\_SERVICE\_TYPE, PROGRAM\_AREA, CAPACITY\_ACTUAL\_BED, CAPACITY\_FUNDING\_BED, OCCUPIED\_BEDS, UNOCCUPIED\_BEDS, UNAVAILABLE\_BEDS (Toronto Open Data Portal 2021). The response variable is OCCUPANCY\_RATE\_BEDS (Toronto Open Data Portal 2021). To begin with, I conduct the exploratory data analysis in order to build the starting model. After creating my starting model by common sense and results of EDA, I check the two conditions before assessing the model assumptions. After that, I use plots to identify potential violations against model assumptions, which are linearity, normality, constant variance, and uncorrelatedness. The next step is to explore model transformations to correct assumption violations. After that, I apply box-cox transformation for numerical variables and create transformed variables. The next step is to fit a new model with transformed variables by applying automated selection. Then, I apply same analysis on my reduced model, which is better than the full model. Overall, It was found that contributing factors include the programs' locations, gender, age, household size of the service user group, the type of overnight service provided, the program area. Furthermore, the number of beds showing as available for occupancy, rooms that a program has been approved to provide, rooms showing as occupied by a shelter user, rooms that are showing as available for occupancy that are not occupied, and rooms that are not currently available also have significant influences on the response variable.

### 5.2 Second discussion point

From this paper, we learn that increasing the number of beds showing as available for occupancy and the beds showing as occupied by a shelter user is an effective way to increase its beds' occupied rates. Moreover, it is also useful to decrease the number of beds that are showing as available for occupancy that are not occupied as of the occupancy date, and beds that are not currently available.

### 5.3 Third discussion point

Furthermore, it's important to find the patterns and trends of the proportion of actual room capacity that is occupied for the reporting date. This report shows its contributing factors, so program managers can utilize this to make disciplines that can solve this issue more effectively. In this way, people can have a better living standards, and the economic development will be improved.

### 5.4 Weaknesses and next steps

The limitations include that the EDA part shows the response variable and predictors don't follow normal distributions. Specifically, there is unbalanced data volume. This may cause the results to be incredible and may be the reasons why the QQ plots contain violations. Also, I didn't checkthe leverage points, so it's possible there do contain pointsthat aren't credible. Moreover, violations still exist in residual plots of my reduced model. For the next step, we should try to avoid these problems and focus on other possible factors that can influence the proportion of actual bed capacity that is occupied for the reporting date.

## References

- Firke, Sam. 2021. *Janitor: Simple Tools for Examining and Cleaning Dirty Data*. <https://CRAN.R-project.org/package=janitor>.
- Fox, John, and Sanford Weisberg. 2019. *An R Companion to Applied Regression*. Third. Thousand Oaks CA: Sage. <https://socialsciences.mcmaster.ca/jfox/Books/Companion/>.
- Pedersen, Thomas Lin. 2020. *Patchwork: The Composer of Plots*.
- R Core Team. 2020. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Toronto Open Data Portal, City of. 2021. *Open Data Dataset*. <https://open.toronto.ca/dataset/daily-shelter-overnight-service-occupancy-capacity/>.
- Wickham, Hadley. 2016. *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. <https://ggplot2.tidyverse.org>.
- Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D'Agostino McGowan, Romain François, Garrett Grolemund, et al. 2019. "Welcome to the tidyverse." *Journal of Open Source Software* 4 (43): 1686. <https://doi.org/10.21105/joss.01686>.
- Wickham, Hadley, and Jennifer Bryan. 2019. *Readxl: Read Excel Files*. <https://CRAN.R-project.org/package=readxl>.
- Wickham, Hadley, Romain François, Lionel Henry, and Kirill Müller. 2021. *Dplyr: A Grammar of Data Manipulation*. <https://CRAN.R-project.org/package=dplyr>.
- Wickham, Hadley, Jim Hester, and Jennifer Bryan. 2022. *Readr: Read Rectangular Text Data*.
- Xie, Yihui. 2021. *Knitr: A General-Purpose Package for Dynamic Report Generation in r*. <https://yihui.org/knitr/>.