# MODUL 11

# SAS® Viya® STEP-BY-STEP DATA SCIENCE PROJECT

## THEME DESCRIPTION

Students understand and are able to implement science projects step by step with an outline of the workflow starting from provide the data to identifying business opportunities, creating Train and Test data sets and analyzing target variables and predictor variables in order to prepare data models.

## WEEKLY LEARNING OUTCOMES (SUB-LESSONS)

CLO-3-Sub-CLO-11:

Able to implement the early stages of data science implementation including business opportunities as well as analysis and selection of effective data variables– C4.

**T**hrough the following learning steps**:**

1. Get the Data
2. Create Geographic Map
3. Exploring data files.
   a) Split data into TRAIN and TEST datasets at an 80/20 split
   b) Investigate the Target Variable
4. Explore the target variable:
   a) Eliminate outliers and create log transformed price variable
   b) Analyze the relationship between neighbourhood and zipcode variables
   c) Create global numeric variables
   d) Create global variables for character variables
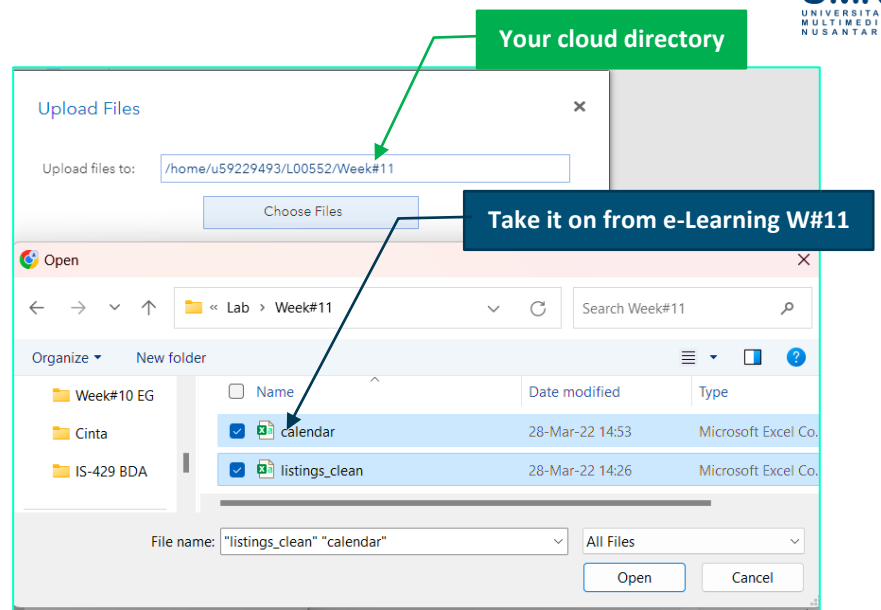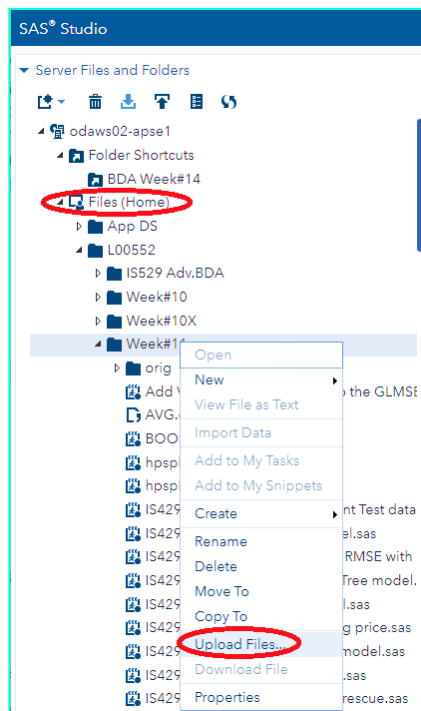5. Predictor Variable Analysis
6. Collinearity Analysis

## PRACTICUM SUPPORT

a. Windows Operating System
b. (any) Browser Application

## PRACTICUM STEPS

### 1. Get the Data

▸ The data provide of New York City Inside Airbnb data sets consisting of Listings and Calendar, have placed them in "/home/u59229493/L00552 shared-drive/u59229493/Week#11/" repository.
▸ These are the two main data sets that we will use for our analysis: listings_clean.csv and calendar.csv.
▸ If you fail to apply the L00552 shared-drive mentioned above then you can upload it directly from your local laptop/PC to your home directory using the following command:
a. Start, https://welcome.oda.sas.com

b. Open SAS®Studio, click options and click New SAS program (F4).

c. If you are not using the L00552 shared-drive, then change the folder according to the one you have!

d. Type it on these codes below:

```
1   LIBNAME MYDATA BASE "/home/u59229493/L00552 shared-drive/u59229493/Week#10";
2
3   FILENAME REFFILE '/home/u59229493/L00552 shared-drive/u59229493/Week#10/listings_clean.csv';
4
5   PROC IMPORT DATAFILE=REFFILE
6      DBMS=CSV
7      OUT= MYDATA.Listings;
8      GETNAMES=YES;
9   RUN;
10
11  PROC CONTENTS DATA=MYDATA.Listings; RUN;
12  /*end----2.1*/
13
14  /*Example work libraries*/
15  DATA WORK.TEST;
16    SET MYDATA.Listings;
17  RUN;
18
19  DATA TEST;
20    SET MYDATA.LISTINGS;
21  RUN;
22
23  FILENAME REFFILE2 '/home/u59229493/L00552 shared link/u59229493/Week#10/calendar.csv';
24
25  PROC IMPORT DATAFILE=REFFILE2
26     DBMS=CSV
27     OUT= MYDATA.Calendar;
28     GETNAMES=YES;
    RUN;
```

e.  After the above program is executed then you will get the data contents of Listings database as follows:

The CONTENTS Procedure

| Data Set Name | WORK.LISTINGS | Observations | 49056 |
|---|---|---|---|
| Member Type | DATA | Variables | 54 |

▸ The first two rows of the PROC CONTENTS output show that there are 49,056 observations with 54 variables in the Listings data set.
▸ The PROC CONTENTS procedure also produces a list of variables and attributes.
▸ The output of the above program produces several tables in the following libraries:
   1) MYDATA.Listings
   2) WORK.TEST
   3) MYDATA.LISTINGS

f.  Save it on your program into "IS429 Lab Week#11_Get the Data.sas".
g.  Print your results as shown on Ouput figure A and save it on  "IS-429 Lab Week#11A Get the Data NIMNAME".pdf

2.  **Create Geographic Map**

   ▸ SAS does an excellent job of creating picture-perfect maps with a few lines of code.
   a. Click New SAS program (F4).
   b. Type it on these codes below:

```
1  /* Create Geographic Map */
2  ODS GRAPHICS / RESET WIDTH=6.4in HEIGHT=4.8in;
3  PROC SGMAP plotdata=  MYDATA.LISTINGS;
4        openstreetmap;
5        scatter x=longitude y=latitude /
6  markerattrs=(size=3 symbol=circle);
7  RUN;
8  ODS GRAPHICS / RESET;
9
```

   ▸ They are highly customizable in Figure 10.1  with the ability to create choropleth maps and bubble maps along with lots of other customizable features.



Figure 10.1 choropleth maps and bubble maps

- For additional information, you can group the data points by a categorical feature.
- The program code below shows the colors of the Airbnb listing-coded by the Neighborhood Group feature."

c. Type it on these codes below:

```
10 /*Neighborhood Group Map View*/
11 ODS GRAPHICS / RESET WIDTH=6.4in HEIGHT=4.8in;
12 PROC SGMAP plotdata=MYDATA.LISTINGS;
13        openstreetmap;
14        scatter x=longitude y=latitude /
15        group=neighbourhood_group_cleansed
16        name='scatterPlot' markerattrs=(size=3 symbol=circle);
17        keylegend 'scatterPlot' /
18        title='neighbourhood_group_cleansed';
19 RUN;
20 ODS GRAPHICS / RESET;
```

- The map view that is grouped by neighborhood in Figure 10.2 provides an easy visual snapshot of the data. We can see that Manhattan and Brooklyn appear to have a majority of the listings while The Bronx and Staten Island have relatively few listings.
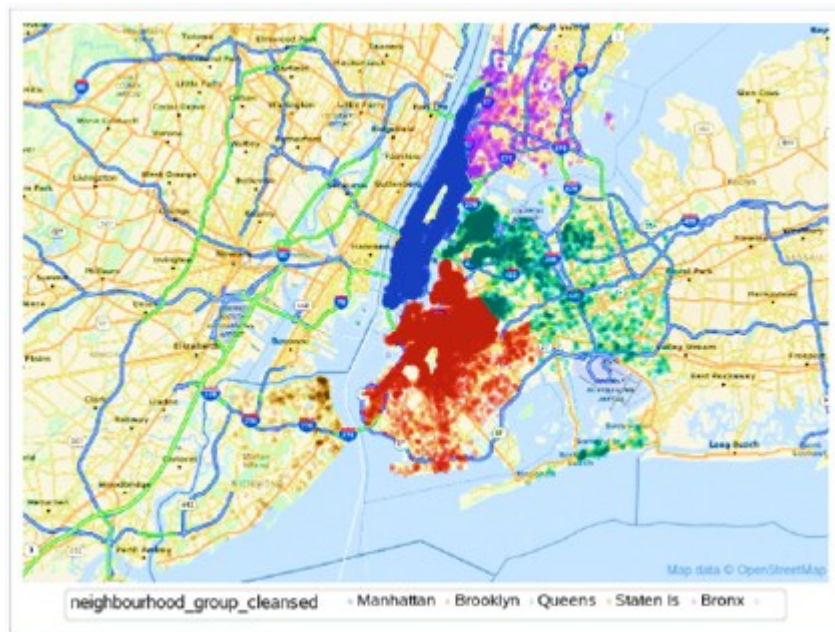


Figure 10.2 Map View by Neighborhood

d. Save it on your program into "IS429 Lab Week#11_Geographic Map.sas"
e. Screenshot your results as shown on Ouput figure B and save it on "IS-429 Lab Week#11B Geographic Map-NIMNAME.JPG/PNG/pdf".

3. Exploring data files

- Before we start investigating the data and making any adjustments, we first need to split the data into TRAIN and TEST data sets.
- Create the TRAIN and TEST datasets using the **PROC SURVEY SELECT** procedure.
- This is a great tool to use which allows you **to set your sampling rate** (SAMPRATE). We use the 80/20 division for our TRAIN and TEST, with 80% of the data going to the TRAIN data set.

- For this Airbnb project, we applied a simple random sample without replacement. This sampling method is carried out by METHOD = SRS statement.
- Next is to build the seed. This determines the initial seed for random number generation. In the model development phase of the project, we want to make sure that we create the same data set every time we run our code. Seed formation is a great way to ensure that you can replicate your results.
- We use any number for the seed value. But for this Airbnb project we use 42 because **Douglas Adams** shows us that **42 is the "Answer to the Highest"** Questions about Life, the Universe, and Everything" in his book "**Hitchhiker's Guide Series for Galaxy**".

a. Still in the worksheet of SAS®Studio, click options and click New SAS program (F4).
b. Type it on the codes below and run.

```
1  /*Develop an 80/20 Split Indicator*/
2  /*Split data into TRAIN and TEST datasets at an 80/20 split*/
3  PROC SURVEYSELECT DATA=MYDATA.Listings SAMPRATE=0.20 SEED=42
4        OUT=Full OUTALL METHOD=SRS;
5  RUN;
6
7  DATA TRAIN TEST;
8     SET Full;
9        IF Selected=0 THEN OUTPUT TRAIN; ELSE OUTPUT TEST;
10       DROP Selected;
11 RUN;
```

- The PROC SURVEYSELECT procedure simply adds a field to the existing MYDATA.Listings data set.
- The Selected field is a binary variable that is coded as a 1 or 0. The 20% of observations that are randomly selected will have an indicator of 1, while the remaining 80% will have an indicator of 0.
- The output of the PROC SURVEYSELECT procedure provides you with a summary statement of the chosen parameters. The Sample Size column shows the number of observations that were randomly selected by the algorithm.

| The SURVEYSELECT Procedure | |
|---|---|
| Selection Method | Simple Random Sampling |
| Input Data Set | LISTINGS |
| Random Number Seed | 42 |
| Sampling Rate | 0.2 |
| Sample Size | 9812 |
| Selection Probability | 0.200016 |
| Sampling Weight | 4.999592 |
| Output Data Set | FULL |

- A quick confirmation shows above that 9812/4999592 = 20%.
- In addition to generating the sample size column above, the program generates three FULL TRAIN and TEST tables in your WORK repository.



```
NOTE: There were 49056 observations read from the data set WORK.FULL.
NOTE: The data set WORK.TRAIN has 39244 observations and 54 variables.
NOTE: The data set WORK.TEST has 9812 observations and 54 variables.
```

f. Screenshotted your log and WORK repository as shown in Output figure C and save it as **"IS-429 Lab Week#11C Exploring Data NIMNAME".JPG/PNG/pdf".**
g. Save it on your program into "IS429 Lab Week#11_Exploring Data Files.sas"
h. The output of the above program will produce the following analysis results:

1) The neighbourhood_group per zipcode using the FREQ procedure with the result that 3 missing frequencies were found

2) The neighbourhood_group per room_type using the FREQ procedure with the result that 3 missing frequencies were found

3) The output of the PROC SURVEYSELECT procedure provides you with a summary statement of the chosen parameters. The Sample Size column shows the number of observations that were randomly selected by the algorithm. A quick confirmation shows that 9812/49056 = 20%.

You can check the log to confirm that the TRAIN and TEST data sets are the expected sizes and that there are no errors in the processing.

```
NOTE: There were 49056 observations read from the data set WORK.FULL.
NOTE: The data set WORK.TRAIN has 39244 observations and 54 variables.
NOTE: The data set WORK.TEST has 9812 observations and 54 variables.
```

4. **Explore the target variable**

   ▸ Type it on below codes and run:

```
1
2  /* Explore the target variable */
3  PROC UNIVARIATE DATA=TRAIN; VAR Price; HISTOGRAM ; RUN;
4
5
6  /* Eliminate outliers and create log transformed price variable */
7  DATA Price;
8     SET TRAIN;
9     WHERE 30 le Price le 750;
10    Price_Log = LOG(Price);
11 RUN;
12
13
14 PROC UNIVARIATE DATA=WORK.Price;
15    VAR Price;
16    HISTOGRAM;
17 RUN;
18
19 PROC UNIVARIATE DATA=WORK.Price; VAR Price_Log; HISTOGRAM; RUN;
20
21
22 PROC MEANS DATA=Price (KEEP = _NUMERIC_) N NMISS MIN MAX MEAN MEDIAN STD; RUN;
```

   ▸ Let's start by focusing on the cross-sectional data in the TRAIN data set. We are attempting to predict the per diem rate of Airbnb listings in New York City as of 12/06/2018.
   ▸ The data set contains a few different price variables for us to choose from:
      ● Price
      ● Weekly price
      ● Monthly price
      ● Security deposit
      ● Cleaning fee
      ● Extra guest fee

   ▸ The Price variable looks like the right one to use for our model. The other variables relating to cost are important pieces of information that we might use later, but for now, we will investigate our selected target variable.
   ▸ We have specified that we want to look at the Price variable. We have also requested that the output produce a histogram of the Price variable.

a. The distribution of Price using UNIVARIATE procedure

▸ The output of the PROC UNIVARIATE procedure contains a lot of information, so let's look at it one section at a time.

▸ The first section of Output shows that the data contained 49,053 observations and that the average value for price is $152.30.

<table>
<tr><th colspan="4">The UNIVARIATE Procedure<br>Variable: price</th></tr>
<tr><th colspan="4">Moments</th></tr>
<tr><td>N</td><td>49053</td><td>Sum Weights</td><td>49053</td></tr>
<tr><td>Mean</td><td>152.301653</td><td>Sum Observations</td><td>7470853</td></tr>
<tr><td>Std Deviation</td><td>227.800206</td><td>Variance</td><td>51892.934</td></tr>
<tr><td>Skewness</td><td>20.4539323</td><td>Kurtosis</td><td>697.82581</td></tr>
<tr><td>Uncorrected SS</td><td>3683275463</td><td>Corrected SS</td><td>2545452199</td></tr>
<tr><td>Coeff Variation</td><td>149.571723</td><td>Std Error Mean</td><td>1.02854033</td></tr>
</table>

▸ The next section of the PROC UNIVARIATE output shows the basic statistical measures.

▸ This section contains the average value of the Price variable (mean, median, and mode) along with measures of variability (std dev, variance, range, and interquartile range).

<table>
<tr><th colspan="4">Basic Statistical Measures</th></tr>
<tr><th colspan="2">Location</th><th colspan="2">Variability</th></tr>
<tr><td>Mean</td><td>152.3017</td><td>Std Deviation</td><td>227.80021</td></tr>
<tr><td>Median</td><td>110.0000</td><td>Variance</td><td>51893</td></tr>
<tr><td>Mode</td><td>150.0000</td><td>Range</td><td>10000</td></tr>
<tr><td></td><td></td><td>Interquartile Range</td><td>109.00000</td></tr>
</table>

▸ There is a significant difference between the mean and median values. This difference is further evidence that outlier values are affecting the mean value.

▸ The range metric shows us the difference between the highest and lowest value. There is a $10,000 per diem range for price. Wow, that high end must be an awesome place to live.

▸ The final section of the PROC UNIVARIATE output below contains more information about the details of the variable distribution, along with outliers and missing values.

<table>
<tr><th colspan="2">Quantiles (Definition 5)</th></tr>
<tr><th>Level</th><th>Quantile</th></tr>
<tr><td>100% Max</td><td>10000</td></tr>
<tr><td>99%</td><td>750</td></tr>
<tr><td>95%</td><td>355</td></tr>
<tr><td>90%</td><td>272</td></tr>
<tr><td>75% Q3</td><td>179</td></tr>
<tr><td>50% Median</td><td>110</td></tr>
<tr><td>25% Q1</td><td>70</td></tr>
<tr><td>10%</td><td>50</td></tr>
<tr><td>5%</td><td>40</td></tr>
<tr><td>1%</td><td>30</td></tr>
<tr><td>0% Min</td><td>0</td></tr>
</table>

<table>
<tr><th colspan="4">Extreme Observations</th></tr>
<tr><th colspan="2">Lowest</th><th colspan="2">Highest</th></tr>
<tr><th>Value</th><th>Obs</th><th>Value</th><th>Obs</th></tr>
<tr><td>0</td><td>42926</td><td>9999</td><td>13960</td></tr>
<tr><td>0</td><td>38781</td><td>10000</td><td>10346</td></tr>
<tr><td>0</td><td>31405</td><td>10000</td><td>12395</td></tr>
<tr><td>0</td><td>31397</td><td>10000</td><td>17985</td></tr>
<tr><td>0</td><td>31339</td><td>10000</td><td>20017</td></tr>
</table>

<table>
<tr><th colspan="4">Missing Values</th></tr>
<tr><th rowspan="2">Missing Value</th><th rowspan="2">Count</th><th colspan="2">Percent Of</th></tr>
<tr><th>All Obs</th><th>Missing Obs</th></tr>
<tr><td>.</td><td>3</td><td>0.01</td><td>100.00</td></tr>
</table>

- ▸ The first section of the above Output shows the quantiles of the variable.
- ▸ This contains the minimum and maximum values as well as values at certain percentage cutoff points. The maximum value for the Price variable is $10,000 while the minimum is $0.
- ▸ At this point, we need to ask ourselves a question. Does it make any sense that someone would charge $0 per night to stay at their house?
- ▸ Maybe they are very lonely and just want some company. It is probably just the result of bad data entry.
- ▸ We also see that in the missing values section of the output, there are two observations with missing price data. This is actually good news. Usually, user-entered data is messier than this.
- ▸ These statistics show that less than 1% of the data would need to be either adjusted or deleted. We will make those decisions in the near future.
- ▸ On the high end of the distributions, we see that the 99% cutoff value is $750. There is a massive difference between the 99% and 100% values. That top 1% of observations are extreme outlier data points.
- ▸ The final piece of output that we can examine is the histogram of the Price variable. The distribution chart gives a visual demonstration of the statistics that we just reviewed.
- ▸ The histogram Output above shows a highly skewed distribution with a very long right-sided tail. The story that this chart tells us matches perfectly with the skewness value of 20.25 that we see in the statistics table above.
- ▸ All of this information tells us that the Price variable has significant outliers on the high end that we need to fix before we can do any modeling.
- ▸ So, we have investigated the target variable and determined that values with missing or $0 entries do not make sense. We have also investigated the high-end values, and we have determined that a cutoff point of $750 would eliminate extreme outliers. Let's code it up...

b. Save it on your program into "IS429 Lab Week#11_Explore the target variable"

**Investigate the target variable:**

- ▸ The data provide of New York City Inside Airbnb data sets consisting of Listings and Calendar, have placed them in "/home/u59229493/L00552 shared-drive/u59229493/Week#11/" repository or you can use as previous instructions to get in from your local laptop/ PC.
- ▸ These are the two main data sets that we will use for our analysis.
- a. Eliminate outliers and create log transformed price variable
  - ▸ The log shows us that the adjusted data set contains 38,527 observations. The original data set contained 39,244; therefore, we have lost only 717 observations.

```
NOTE: There were 38527 observations read from the data set WORK.TRAIN.
      WHERE (Price>=30 and Price<=750);
NOTE: The data set WORK.PRICE has 38527 observations and 55 variables
```

Table: WORK.PRICE ▾ | View: OUTPUT DATA ▾ | ▸ WORK

Columns    ⊙    Total rows: 38527 Total columns: 55    ▸ TRAIN Total rows: 39244 Total columns: 54

2. The majority of these excluded observations were for properties that have a $0 per diem rate.
- ▸ Although our target variable is much cleaner now that we have established our upper and lower boundaries, we expect that the distribution of the data will still be skewed to the right.

▸ A PROC UNIVARIATE procedure developed on the Price variable shows a right tail skewness with a value of 2.16

| The UNIVARIATE Procedure Variable: price | | | |
|---|---|---|---|
| Moments | | | |
| N | 38527 | Sum Weights | 38527 |
| Mean | 139.258702 | Sum Observations | 5365220 |
| Std Deviation | 102.106469 | Variance | 10425.731 |
| Skewness | 2.1610293 | Kurtosis | 6.43560029 |
| Uncorrected SS | 1148815284 | Corrected SS | 401661713 |
| Coeff Variation | 73.3214282 | Std Error Mean | 0.52020038 |


Distribution of price

▸ In anticipation of this skewness, We created a log-transformed version of the Price variable in the previous DATA step:
```
PROC MEANS DATA=Price (KEEP = _NUMERIC_) N NMISS MIN MAX MEAN MEDIAN STD;
RUN;
```
▸ This approach will transform the Price variable into a normally distributed variable. A PROC UNIVARIATE developed on the log-transformed Price variable shows a skewness value of 0.26

| The UNIVARIATE Procedure Variable: Price_Log | | | |
|---|---|---|---|
| Moments | | | |
| N | 38527 | Sum Weights | 38527 |
| Mean | 4.72538955 | Sum Observations | 182055.083 |
| Std Deviation | 0.63568501 | Variance | 0.40409543 |
| Skewness | 0.27656183 | Kurtosis | -0.3750729 |
| Uncorrected SS | 875849.368 | Corrected SS | 15568.1804 |
| Coeff Variation | 13.4525418 | Std Error Mean | 0.00323862 |


Distribution of Price_Log

▸ Great! We now have a data set where we have eliminated extreme outliers and nonsensical values.
▸ We have also transformed our target variable into a normally distributed variable.
▸ Save it on your program into "IS429 Lab Week#11_Explore the target variable.sas"
▸ **Screenshotted your Price table and Price_log distribution chart as show on Output figure D.**
▸ **Save it as "IS-429 Lab Week#11D Target Variable NIM Name.JPG/PNG/pdf"**

5. **Predictor Variable Analysis**

▸ Due to the goal of the Airbnb project, We want to develop a tool that property owners could use to give them the optimal per diem price for their property on Airbnb.
▸ The predictor variables will be the values that property owners enter to describe the features of their property. You have the option of including additional features based on regional factors (tourist attractions, crime rate, events, and so on), but for now, we will focus on the data that we have already gathered.
▸ Under this scenario, we could make some assumptions:
   ➲ We can expect that the traditional home features will be significant in the model:
      ◦ Location
      ◦ Property type
      ◦ Number of people the property will accommodate
      ◦ Number of beds, baths, rooms, and so on
   ➲ Some data will not be available at the point of application:
      ◦ User review scores
      ◦ Review volume

- Further analysis of this variable shows that there are differences between property owners who have a single property compared to those that have between two and ten properties and those who are major property owners who have more than ten properties.
- We have created by Clean table, a categorical variable for these levels and discarded the original numeric variable.

**The FREQ Procedure**

| host_count_cat | Frequency | Percent | Cumulative Frequency | Cumulative Percent |
|---|---|---|---|---|
| Level 1 | 23534 | 61.08 | 23534 | 61.08 |
| Level 2 | 12350 | 32.06 | 35884 | 93.14 |
| Level 3 | 2643 | 6.86 | 38527 | 100.00 |

a. Import Clean.csv from "/home/u59229493/L00552 shared-drive/u59229493/Week#11/clean.csv" or you can use as previous instructions to get in from your local laptop/ PC to your WORK space repository.



b. Type it on these codes below and run them all:

```
1  /*Import Clean.csv from /home/u59229493/L00552 shared-drive/u59229493/Week#10/clean.csv to your WORK repository*/
2  /* proc univariate data=clean; var bedrooms bathrooms beds; run; */
3  PROC FREQ DATA=Clean; TABLES host_count_cat; RUN;
4
5  PROC MEANS DATA=Clean;
6     CLASS host_count_cat;
7     VAR Price Price_Log;
8  RUN;
9
10 /* Create global numeric variables */
11 PROC CONTENTS NOPRINT DATA=Clean (KEEP=_NUMERIC_ DROP=id host_id
12 latitude longitude Price Price_Log) OUT=var1 (KEEP=name);
13 RUN;
14
15 PROC SQL NOPRINT;
16    SELECT name INTO:varx separated by " " FROM var1;
17 QUIT;
18
19 %PUT &varx;
20
21 /* Create correlation analysis */
22 PROC CORR DATA=Clean;
23    VAR &varx.;
24 RUN;
```

## 6. Collinearity Analysis

- One of the founding assumptions of regression models is that the predictor variables need to be independent (not correlated with one another).
- This assumption is important because if two or more variables are closely related, it is difficult to separate the individual effects of those variables on the response variable. In the end, collinearity reduces the accuracy of the estimates of the regression coefficients because it inflates the standard error of those coefficients.

▶ The table below provides a good rule of thumb for interpreting correlation coefficients. Any correlation coefficient above 0.7 or below -0.7 is considered highly correlated.

| Size of Correlation | Interpretation |
|---|---|
| 0.9 to 1.0 (-0.9 to -1.0) | Very high correlation |
| 0.7 to 0.9 (-0.7 to -0.9) | High correlation |
| 0.5 to 0.7 (-0.5 to -0.7) | Moderate correlation |
| 0.3 to 0.5 (-0.3 to -0.5) | Low correlation |
| 0.0 to 0.3 (0.0 to -0.3) | Negligible correlation |

▶ SAS provides two great ways to identify multicollinearity through the CORR and REG procedures.
▶ PROC CORR creates a correlation matrix that contains the correlation coefficient for each variable combination.

Pearson Correlation Coefficients, N = 38527
Prob > |r| under H0: Rho=0

| | accommodates | availability_30 | availability_60 | availability_90 | availability_365 | bath_per_accom | bathrooms | bedrooms | beds | beds_per_accom | guests_included | maximum_nights | minimum_nights |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| accommodates | 1.00000 | 0.04635 <.0001 | 0.09961 <.0001 | 0.11312 <.0001 | 0.14772 <.0001 | 0.84018 <.0001 | 0.30334 <.0001 | 0.63635 <.0001 | 0.76936 <.0001 | 0.39991 <.0001 | 0.57702 <.0001 | 0.04360 <.0001 | -0.01487 0.0035 |
| availability_30 | 0.04635 <.0001 | 1.00000 | 0.88138 <.0001 | 0.80279 <.0001 | 0.56338 <.0001 | 0.03323 <.0001 | 0.04367 <.0001 | 0.03176 <.0001 | 0.04436 <.0001 | -0.02820 <.0001 | 0.03163 <.0001 | 0.00976 0.0553 | 0.21671 <.0001 |
| availability_60 | 0.09961 <.0001 | 0.88138 <.0001 | 1.00000 | 0.97204 <.0001 | 0.70076 <.0001 | 0.08516 <.0001 | 0.05175 <.0001 | 0.06033 <.0001 | 0.09276 <.0001 | -0.01705 0.0008 | 0.08975 <.0001 | -0.00442 0.3856 | 0.18545 <.0001 |
| availability_90 | 0.11312 <.0001 | 0.80279 <.0001 | 0.97204 <.0001 | 1.00000 | 0.74642 <.0001 | 0.09653 <.0001 | 0.05632 <.0001 | 0.06843 <.0001 | 0.10651 <.0001 | -0.01498 0.0033 | 0.10349 <.0001 | -0.00735 0.1491 | 0.18173 <.0001 |
| availability_365 | 0.14772 <.0001 | 0.56338 <.0001 | 0.70076 <.0001 | 0.74642 <.0001 | 1.00000 | 0.11822 <.0001 | 0.08195 <.0001 | 0.08579 <.0001 | 0.14514 <.0001 | -0.01108 0.0297 | 0.11724 <.0001 | 0.07323 <.0001 | 0.27086 <.0001 |
| bath_per_accom | 0.84018 <.0001 | 0.03323 <.0001 | 0.08516 <.0001 | 0.09653 <.0001 | 0.11822 <.0001 | 1.00000 | -0.15773 <.0001 | 0.42211 <.0001 | 0.60927 <.0001 | 0.43778 <.0001 | 0.46848 <.0001 | 0.03376 <.0001 | -0.02654 <.0001 |
| bathrooms | 0.30334 <.0001 | 0.04367 <.0001 | 0.05175 <.0001 | 0.05632 <.0001 | 0.08195 <.0001 | -0.15773 <.0001 | 1.00000 | 0.37430 <.0001 | 0.31388 <.0001 | -0.01151 0.0239 | 0.18395 <.0001 | 0.03093 <.0001 | 0.04715 <.0001 |
| bedrooms | 0.63635 <.0001 | 0.03176 <.0001 | 0.06033 <.0001 | 0.06843 <.0001 | 0.08579 <.0001 | 0.42211 <.0001 | 0.37430 <.0001 | 1.00000 | 0.65316 <.0001 | 0.01766 0.0005 | 0.41964 <.0001 | 0.01157 0.0232 | 0.01101 0.0307 |
| beds | 0.76936 <.0001 | 0.04436 <.0001 | 0.09276 <.0001 | 0.10651 <.0001 | 0.14514 <.0001 | 0.60927 <.0001 | 0.31388 <.0001 | 0.65316 <.0001 | 1.00000 | -0.08219 <.0001 | 0.48548 <.0001 | 0.03024 <.0001 | 0.00570 0.2629 |
| beds_per_accom | 0.39991 <.0001 | -0.02820 <.0001 | -0.01705 0.0008 | -0.01498 0.0033 | -0.01108 0.0297 | 0.43778 <.0001 | -0.01151 0.0239 | 0.01766 0.0005 | -0.08219 <.0001 | 1.00000 | 0.15311 <.0001 | 0.02783 <.0001 | -0.03786 <.0001 |
| guests_included | 0.57702 <.0001 | 0.03163 <.0001 | 0.08975 <.0001 | 0.10349 <.0001 | 0.11724 <.0001 | 0.46848 <.0001 | 0.18395 <.0001 | 0.41964 <.0001 | 0.48548 <.0001 | 0.15311 <.0001 | 1.00000 | -0.01932 0.0001 | -0.07096 <.0001 |
| maximum_nights | 0.04360 <.0001 | 0.00976 0.0553 | -0.00442 0.3856 | -0.00735 0.1491 | 0.07323 <.0001 | 0.03376 <.0001 | 0.03093 <.0001 | 0.01157 0.0232 | 0.03024 <.0001 | 0.02783 <.0001 | -0.01932 0.0001 | 1.00000 | 0.06014 <.0001 |
| minimum_nights | -0.01487 0.0035 | 0.21671 <.0001 | 0.18545 <.0001 | 0.18173 <.0001 | 0.27086 <.0001 | -0.02654 <.0001 | 0.04715 <.0001 | 0.01101 0.0307 | 0.00570 0.2629 | -0.03786 <.0001 | -0.07096 <.0001 | 0.06014 <.0001 | 1.00000 |

▶ The default correlation statistic is the Pearson's r. The results of the procedure are below. We have color-coded the top 10% of high and low correlation values.

| | accommodates | availability 30 | availability 60 | availability 90 | availability 365 | bathrooms | bedrooms | beds | guests included | maximum nights | minimum nights |
|---|---|---|---|---|---|---|---|---|---|---|---|
| accommodates | | 0.05 | 0.10 | 0.11 | 0.15 | 0.31 | 0.63 | 0.77 | 0.57 | 0.04 | (0.02) |
| availability 30 | 0.05 | | 0.88 | 0.80 | 0.56 | 0.04 | 0.03 | 0.04 | 0.03 | 0.01 | 0.22 |
| availability 60 | 0.10 | 0.88 | | 0.97 | 0.70 | 0.05 | 0.06 | 0.09 | 0.09 | (0.01) | 0.18 |
| availability 90 | 0.11 | 0.80 | 0.97 | | 0.75 | 0.05 | 0.06 | 0.10 | 0.10 | (0.01) | 0.18 |
| availability 365 | 0.15 | 0.56 | 0.70 | 0.75 | | 0.07 | 0.08 | 0.14 | 0.12 | 0.07 | 0.27 |
| bathrooms | 0.31 | 0.04 | 0.05 | 0.05 | 0.07 | | 0.37 | 0.31 | 0.18 | 0.03 | 0.04 |
| bedrooms | 0.63 | 0.03 | 0.06 | 0.06 | 0.08 | 0.37 | | 0.65 | 0.42 | 0.01 | 0.01 |
| beds | 0.77 | 0.04 | 0.09 | 0.10 | 0.14 | 0.31 | 0.65 | | 0.49 | 0.03 | 0.01 |
| guests included | 0.57 | 0.03 | 0.09 | 0.10 | 0.12 | 0.18 | 0.42 | 0.49 | | (0.02) | (0.07) |
| maximum nights | 0.04 | 0.01 | (0.01) | (0.01) | 0.07 | 0.03 | 0.01 | 0.03 | (0.02) | | 0.06 |
| minimum nights | (0.02) | 0.22 | 0.18 | 0.18 | 0.27 | 0.04 | 0.01 | 0.01 | (0.07) | 0.06 | |

▶ Variables are considered highly correlated if the correlation coefficient is above 0.7 or below -0.7.
▶ The bottom 10% of correlation values are highlighted in green. We can see that negative correlation is not a problem since the strongest negative correlation is -0.07.

- For the positively correlated variables, we see that there are several cases where the correlation coefficient is above 0.7. It makes sense that the availability variables are highly correlated. If a property is available within 30 days, it would also be available within the next 60, 90, and 365 days. Also, the correlation between accommodates and beds makes sense. If a property has a high number of beds, it can accommodate more people.
- Save it on your program into "IS429 Lab Week#11_Predictor Variable Analysis.sas"

c. <mark>Print your results as shown in Output figure E into "IS-429 Lab Week#11E Predictor Variable AnalysisNIMNAME".pdf</mark>

b. **Scatter Matrix**

- An additional method of analyzing our numeric variables is to produce scatter plots of each variable's relationship with the target variable. SAS provides a great visual demonstration of these relationships with the scatter matrix graph.
- This graph contains a series of scatter plots and allows the researcher to easily see the relationships that each of the numeric predictor variables has with the target variable. We can also see the relationship that the numeric predictors have with one another.
- In the Output Figure E, we can already see the strong positive relationship between the price and accommodates variables, as well as between price and guests included.

a. Still on SAS®Studio, click options and click New SAS program (F4), type in below codes and run it.

```
1  /* ------------------------------------------------------------------
2     Run the standardize procedure and Create Scatter Matrix
3     ------------------------------------------------------------------ */
4  PROC STANDARD DATA=Clean OUT=Stnd_Clean
5     MEAN=0 STD=1 REPLACE;
6     VAR accommodates bathrooms bedrooms beds guests_included
7         minimum_nights maximum_nights availability_30 beds_per_accom bath_per_accom
8       poly_accom poly_bath poly_guests poly_min poly_max poly_avail ;
9  RUN;
10
11 PROC SGSCATTER DATA=Stnd_Clean;
12    TITLE 'Scatter Plot Matrix';
13    MATRIX Price_Log accommodates guests_included minimum_nights maximum_nights/
14    START=TOPLEFT ELLIPSE = (ALPHA=0.05 TYPE=PREDICTED) NOLEGEND;
15 RUN;
```

b. <mark>Save your Scatter Matrix as "IS-429 Lab Week#11F ScatterMatrix NIMNAME.pdf".</mark>
c. Save your program into "IS429 Lab Week#11_Scatter Matrix.sas".

⮑ Finally, today's practicum is over..
⮑ <mark>Zipped your screenshots and pdf files  and submit immediately today to e-Learning IS-429 Practicum Week#11 with the naming format IS-429 BDA Lab Week#11 NIM yourName.zip.</mark>

# RESULTS/ OUTPUT

A.  Get the Data



B.  Create Geographic Map

## C. Exploring data files.



## D. Explore the target variable

## E. Predictor Variable Analysis



## F. Scatter Matirx

The
End

## REFERENCE

1. Gearhart, James. 2020. End-to-End Data Science with SAS®: A Hands-On Programming Guide. Cary, NC: SAS Institute Inc.
2. SAS Institute Inc. 2020. SAS® Viya® Programming: Getting Started. SAS Institute Inc. Cary, NC, USA.
3. SAS® Support | Documentation
4. Other additional references are excerpts from various Online Learning/websites.