

EM 623 - Data Science and Knowledge Discovery

EXERCISE 03

Rushabh Barbhaya | Exercise Report | 2/23/2018

Contents

Missing Values 2

Metadata file 3

Outliers 3

Normalization of nbinge 6

Relation Matrix 8

List of Figures

Figure 1 Missing Values from Rattle 2

Figure 2 Missing Values using Excel 2

Figure 3 Normalizing option 6

Figure 4 Column selection 7

Figure 5 Normalized Values 7

Figure 6 Transformation options 8

Figure 7 Correlation 8

Figure 8 Relation Matrix 9

CLICKABLE

Missing Values

To calculate the total number of missing values from 'Hazardous Alcohol Consumption.csv', load the file in Rattle and press 'execute'. It will present the data with number of unique values in a particular column header along with missing values.

Data: **Explore** | Test | Transform | Cluster | Associate | Model | Evaluate | Log

Source: ☒ File ☐ ARFF ☐ ODBC ☐ R Dataset ☐ RData File ☐ Library ☐ Corpus ☐ Script

Filename: Separator: Decimal: ☒ Header

☐ Partition Seed:

☒ Input ☒ Ignore Weight Calculator:

Target Data Type: ☒ Auto ☐ Categorical ☐ Numeric ☐ Survival

No.	Variable	Data Type	Input	Target	Risk	Ident	Ignore	Weight	Comment
1	time	Constant	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	Unique: 1
2	sex	Categorical	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 3
3	age	Categorical	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 9
4	geo	Categorical	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 12
5	iscsed97	Categorical	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 4
6	day	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 50 Missing: 617
7	day_flag	Constant	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	Unique: 1 Missing: 1,113
8	month	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 176 Missing: 176
9	month_flag	Constant	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	Unique: 1 Missing: 1,053
10	week	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 135 Missing: 263
11	week_flag	Constant	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	Unique: 1 Missing: 1,071
12	lt1m	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 286 Missing: 68
13	lt1m_flag	Constant	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	Unique: 1 Missing: 1,017
14	nvr_occ	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 579
15	nvr_occ_flag	Constant	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	Unique: 1 Missing: 989
16	nbinge	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 432 Missing: 6
17	nbinge_flag	Constant	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	Unique: 1 Missing: 994
18	serialid	Ident	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 1,149

Figure 1 Missing Values from Rattle

You can add them all up coming to a total of 7367 missing values. This operation can also be performed in excel using the following function:

```
=COUNTBLANK(A1:R1150)
```

Figure 2 Missing Values using Excel

It amounts to the same value.

Metadata file

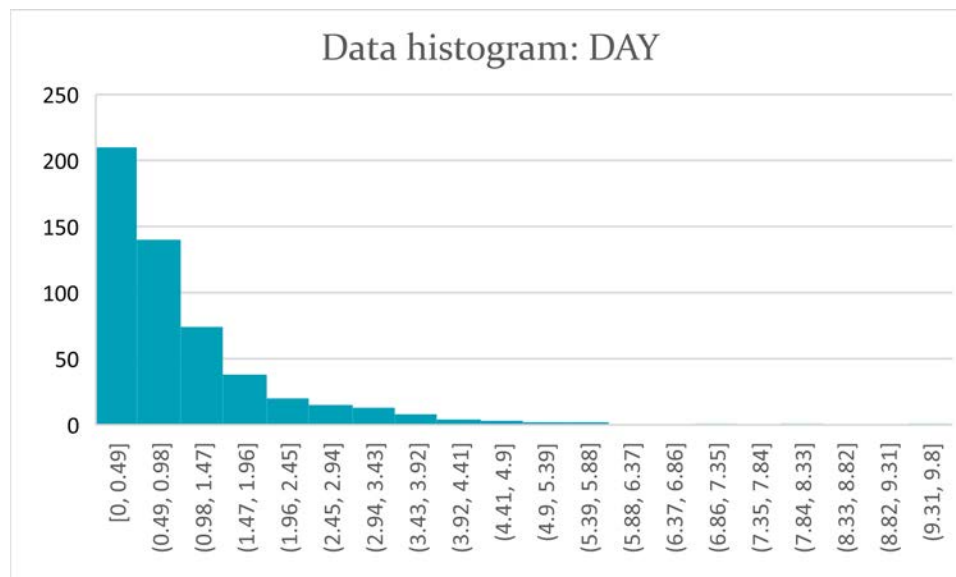
From the file 'Alcohol Consumption-Metadata.csv' we get a little description about the file 'Hazardous_Alcohol_Consumption.csv' we get the whole table name → **Hazardous Alcohol Consumption (Binge Drinking) by Sex Age and Educational Attainment Level**. We also get details about the source of the file, and its datapath.

We also get information about column headers. Like the data was cumulated from January, ISCED97 stands for 'International Standard Classification of Education 1997'.

Outliers

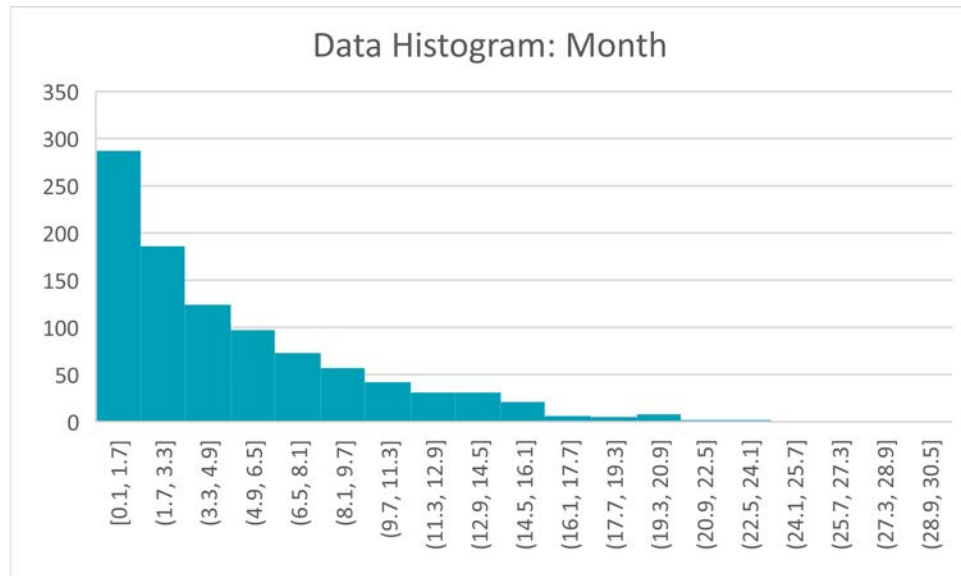
In order to determine the outliers I used a histogram for each data point.

A. Day

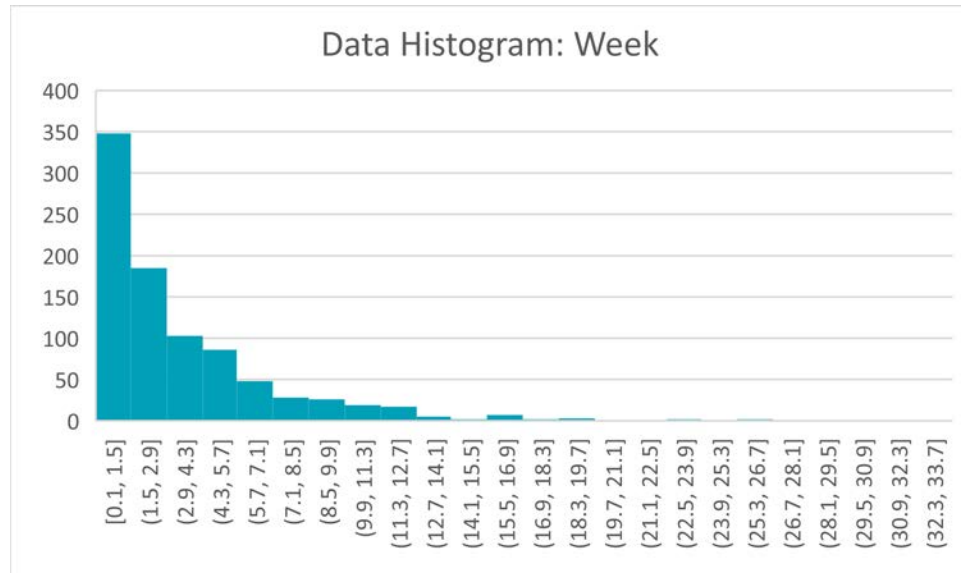


Graph 1 Data Histogram: DAY

From the data represented in this histogram the data above the range of 3.92 seem to be outliers. As they fall outside the normal curve of the graph.

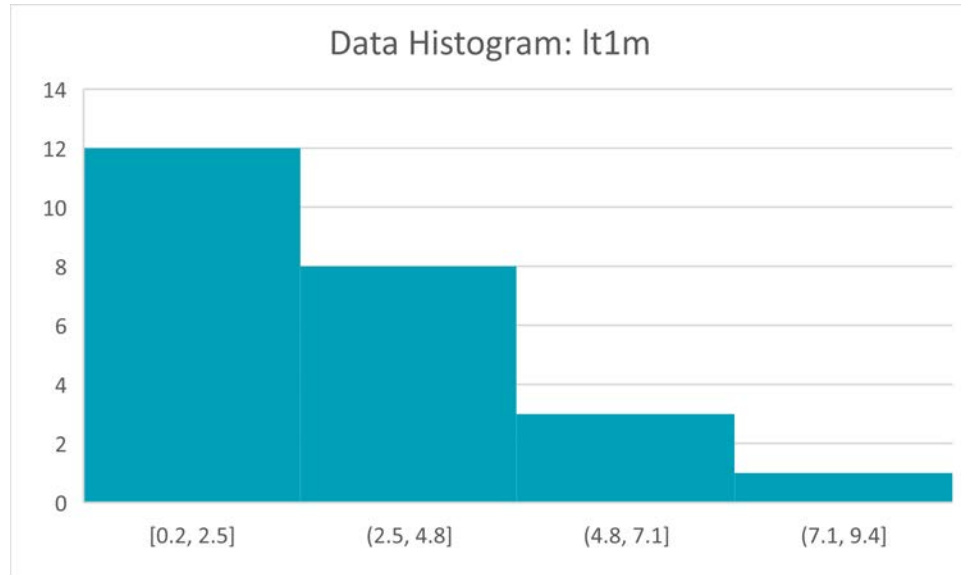
B. Month*Graph 2 Data Histogram: Month*

From the data of month represented in the histogram above it can be assumed that the data above 19.3 seem to be outliers. As they fall outside the normal flow of graph.

C. Week*Graph 3 Data Histogram: Week*

From the data shown above it can be assumed that the data above 15.5 can be considered as outliers as they fall outside the normal curve of the graph.

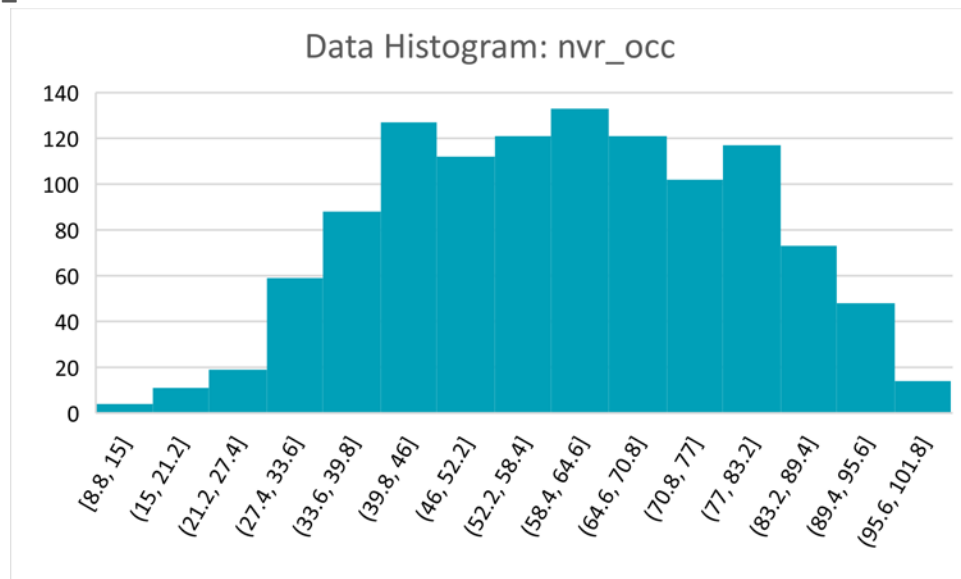
D. lt1m



Graph 4 Data Histogram: lt1m

From the data in the graph above it can be assumed that there are no outliers as all the data fall under the normal curve of the graph.

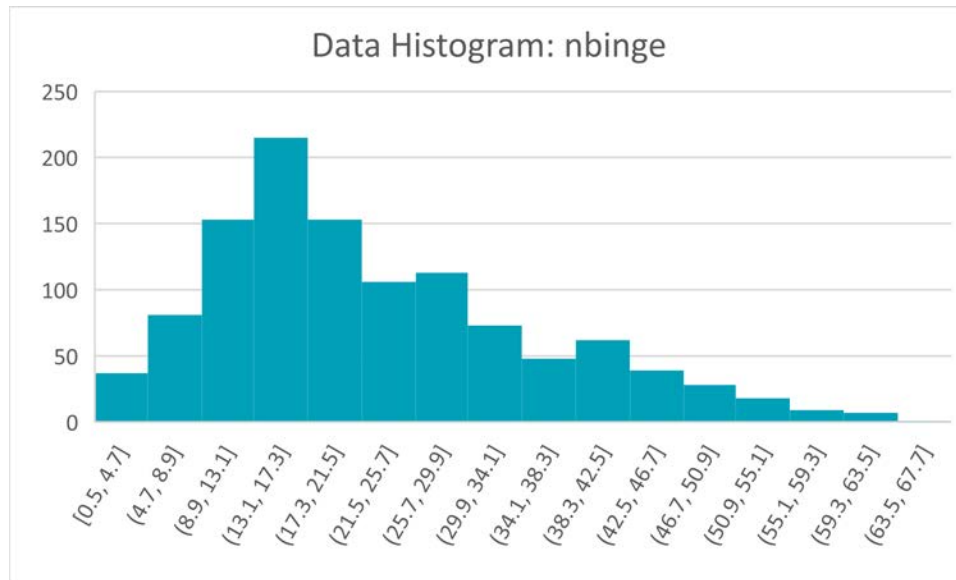
E. nvr_occ



Graph 5 Data Histogram: nvr_occ

From the data in the graph above it can be assumed that there are no outliers as all the data fall under the normal curve of the graph.

F. nbinge



Graph 6 Data Histogram: nbinge

From the data in the graph above it can be assumed that there are no outliers as all the data fall under the normal curve of the graph.

Normalization of nbinge

First load the file in data tab and execute it. Then select the Scale [0-1] option in 'Rescale' type from 'Transform' tab [shown in figure 3]

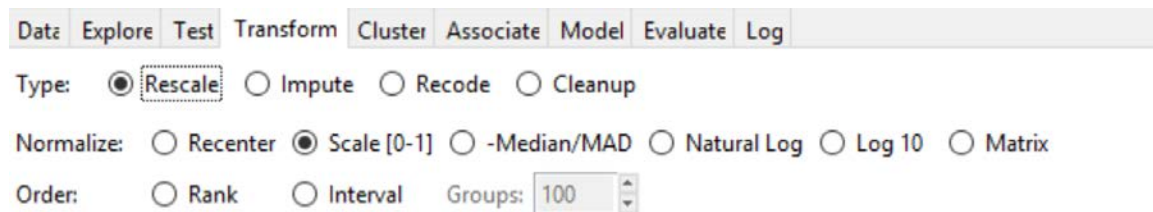


Figure 3 Normalizing option

Then select the column to be rescaled.

14	nvr_occ	Numeric [8.80 to 100.00; unique=579; mean=59.91; median=60.40].
15	nvr_occ_flag	Categorical [1 levels; miss=989; ignored].
16	nbinge	Numeric [0.50 to 63.80; unique=432; mean=22.68; median=19.50; miss=6; ignored].
17	nbinge_flag	Categorical [1 levels; miss=994; ignored].
18	serialid	Numeric [1 to 1149; unique=1149; mean=575; median=575].

Figure 4 Column selection

After pressing execute a new column will be created with a suffix of 'R01_' followed by the name of the column selected for rescaling. After that select view option from data tab and we'll be able to see the normalized table. [shown in figure 5]

R01_nbinge
0.3917852
0.4581359
0.5766193
0.4739336
0.1327014
0.2496051
0.3112164
0.2227488
0.1769352
0.3048973

Figure 5 Normalized Values

Relation Matrix

Rattle cannot display sex and age parameters as they are not numeric values. To convert those to numeric values. To do that we must select 'as numeric' under 'recode' type/option from 'transform' tab and select the rows we need to transform.

Data Explore Test Transform Cluster Associate Model Evaluate Log

Type: ☐ Rescale ☐ Impute ☒ Recode ☐ Cleanup

Binning: ☐ Quantiles ☐ KMeans ☐ Equal Width Number:

☐ Indicator Variable ☐ Join Categories ☐ As Categorical ☒ As Numeric

No.	Variable	Data Type and Number Missing
1	time	Numeric [2008 to 2008; unique=1; mean=2008; median=2008; ignored].
2	sex	Categorical [3 levels; ignored].
3	age	Categorical [9 levels; ignored].
4	geo	Categorical [12 levels].
5	isced97	Categorical [4 levels].

Figure 6 Transformation options

After that we can plot the 'correlation' plot from 'explore' tab. It will show something like this.

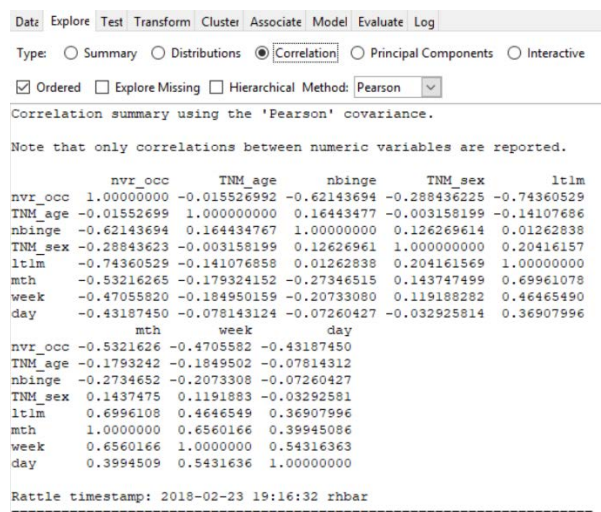


Figure 7 Correlation

To save the plot just open to Rstudio's plot workgroup. The matrix plot can be saved from that option.

Correlation Hazardous_Alcohol_Consumption.csv using Pearson

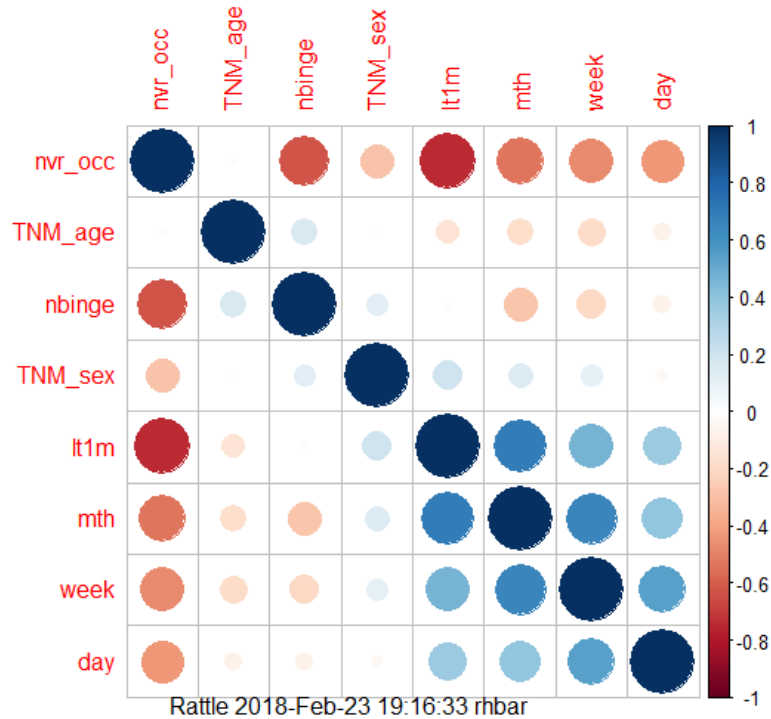


Figure 8 Relation Matrix

It is clear that the 'age' and 'sex' have no correlation at all. But 'age' is so slightly inversely correlated to 'lt1m', 'month', 'week' and 'day' and slightly directly correlated with 'nbinge'.

Similarly 'sex' is directly correlated with 'nbinge', 'lt1m', 'month' and 'week' but inversely correlated with 'nvr_occ' ever so slightly.