



REPORT SERIES WITH DLOOKR

Exploratory Data Analysis Report

Author:
dlookr package

Version:
0.3.12

November 5, 2020

Contents

1	Introduction	3
1.1	Information of Dataset	3
1.2	Information of Variables	3
1.3	About EDA Report	4
2	Univariate Analysis	5
2.1	Descriptive Statistics	5
2.2	Normality Test of Numerical Variables	6
2.2.1	Statistics and Visualization of (Sample) Data	6
3	Relationship Between Variables	27
3.1	Correlation Coefficient	27
3.1.1	Correlation Coefficient by Variable Combination	27
3.1.2	Correlation Plot of Numerical Variables	27
4	Target based Analysis	29
4.1	Grouped Descriptive Statistics	29
4.1.1	Grouped Numerical Variables	29
4.1.2	Grouped Categorical Variables	29
4.2	Grouped Relationship Between Variables	29
4.2.1	Grouped Correlation Coefficient	29
4.2.2	Grouped Correlation Plot of Numerical Variables	29

Chapter 1

Introduction

The EDA Report provides exploratory data analysis information on objects that inherit `data.frame` and `data.frame`.

1.1 Information of Dataset

The dataset that generated the EDA Report is an 'data.frame' object. It consists of 28,534 observations and 21 variables.

1.2 Information of Variables

Table 1.1: Information of Variables

variables	types	missing_count	missing_percent	unique_count	unique_rate
idcode	numeric	0	0.0000000	4711	0.1651013
year	numeric	0	0.0000000	15	0.0005257
birth_yr	numeric	0	0.0000000	14	0.0004906
age	numeric	24	0.0841102	34	0.0011916
race	haven_labelled	0	0.0000000	3	0.0001051
msp	numeric	16	0.0560735	3	0.0001051
nev_mar	numeric	16	0.0560735	3	0.0001051
grade	numeric	2	0.0070092	20	0.0007009
collgrad	numeric	0	0.0000000	2	0.0000701
not_smsa	numeric	8	0.0280367	3	0.0001051
c_city	numeric	8	0.0280367	3	0.0001051
south	numeric	8	0.0280367	3	0.0001051
ind_code	numeric	341	1.1950655	13	0.0004556
occ_code	numeric	121	0.4240555	14	0.0004906
union	numeric	9296	32.5786781	3	0.0001051
wks_ue	numeric	5704	19.9901871	62	0.0021728
ttl_exp	numeric	0	0.0000000	4744	0.1662578
tenure	numeric	433	1.5174879	271	0.0094974
hours	numeric	67	0.2348076	86	0.0030139
wks_work	numeric	703	2.4637275	106	0.0037149
ln_wage	numeric	0	0.0000000	8173	0.2864302

The target variable of the data is 'NULL', and the data type of the variable is NULL(You did not specify a

target variable).

1.3 About EDA Report

EDA reports provide information and visualization results that support the EDA process. In particular, it provides a variety of information to understand the relationship between the target variable and the rest of the variables of interest.

Chapter 2

Univariate Analysis

2.1 Descriptive Statistics

```
Error in proxy[i, ..., drop = FALSE]: incorrect number of dimensions
Error in Hmisc::latex(x, file = ""): object 'x' not found
```

2.2 Normality Test of Numerical Variables

2.2.1 Statistics and Visualization of (Sample) Data

idcode

normality test : Shapiro-Wilk normality test
 statistic : 0.95281, p-value : 2.06327E-37

type	skewness	kurtosis
original	-0.0497	1.7943
log transformation	-2.1184	9.4828
sqrt transformation	-0.6177	2.4776

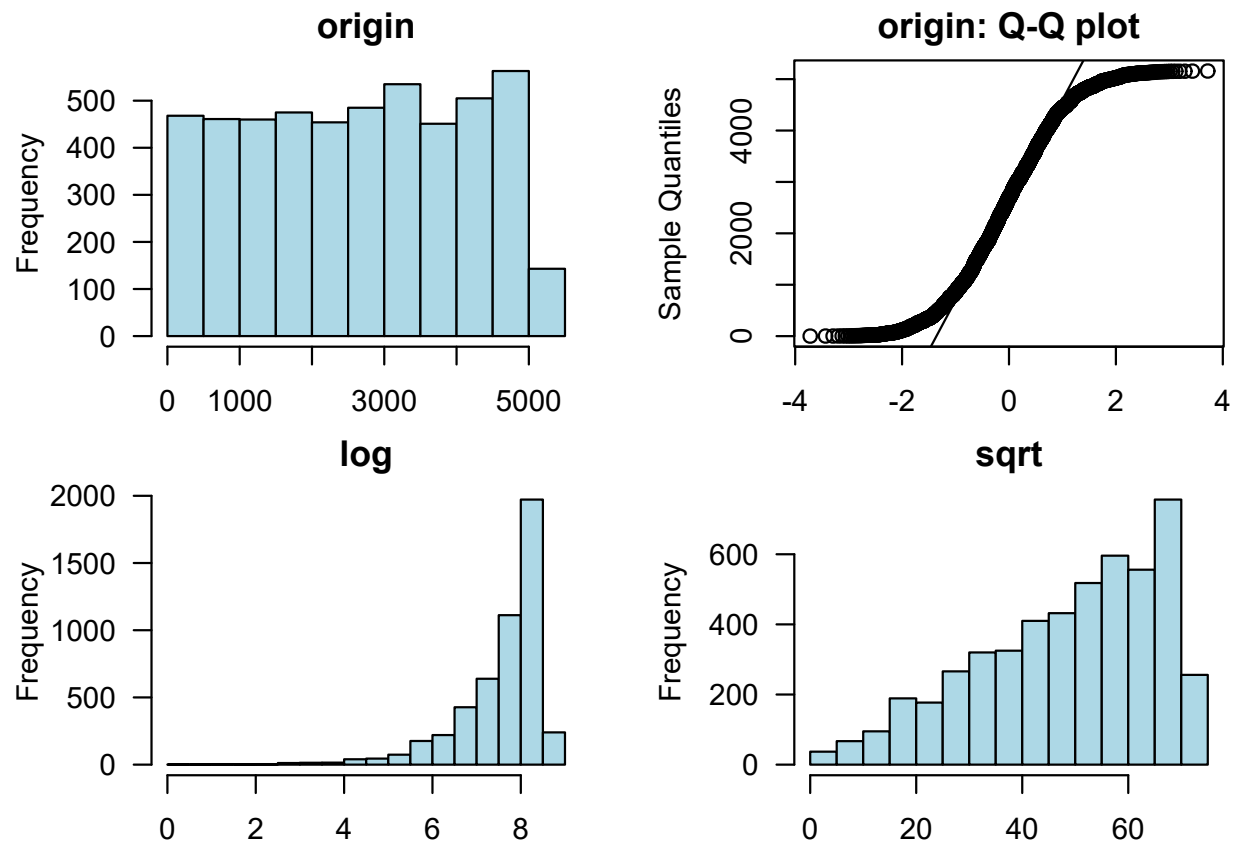


Figure 2.1: idcode

year

normality test : Shapiro-Wilk normality test
 statistic : 0.93035, p-value : 2.55668E-43

type	skewness	kurtosis
original	0.0956	1.6926
log transformation	0.0109	1.6887
sqrt transformation	0.0533	1.6887

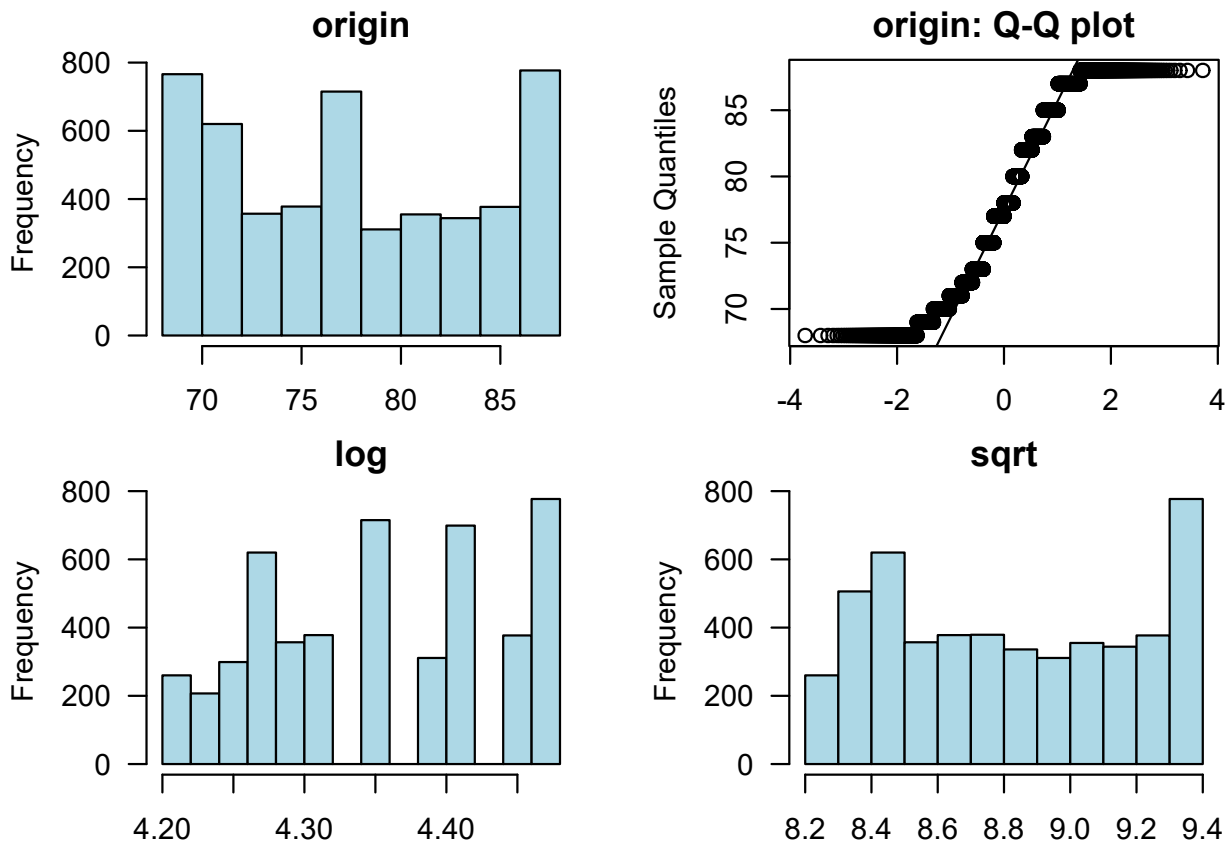


Figure 2.2: year

birth_yr

normality test : Shapiro-Wilk normality test
 statistic : 0.96021, p-value : 5.73494E-35

type	skewness	kurtosis
original	-0.1189	2.0054
log transformation	-0.2141	2.0580
sqrt transformation	-0.1663	2.0289

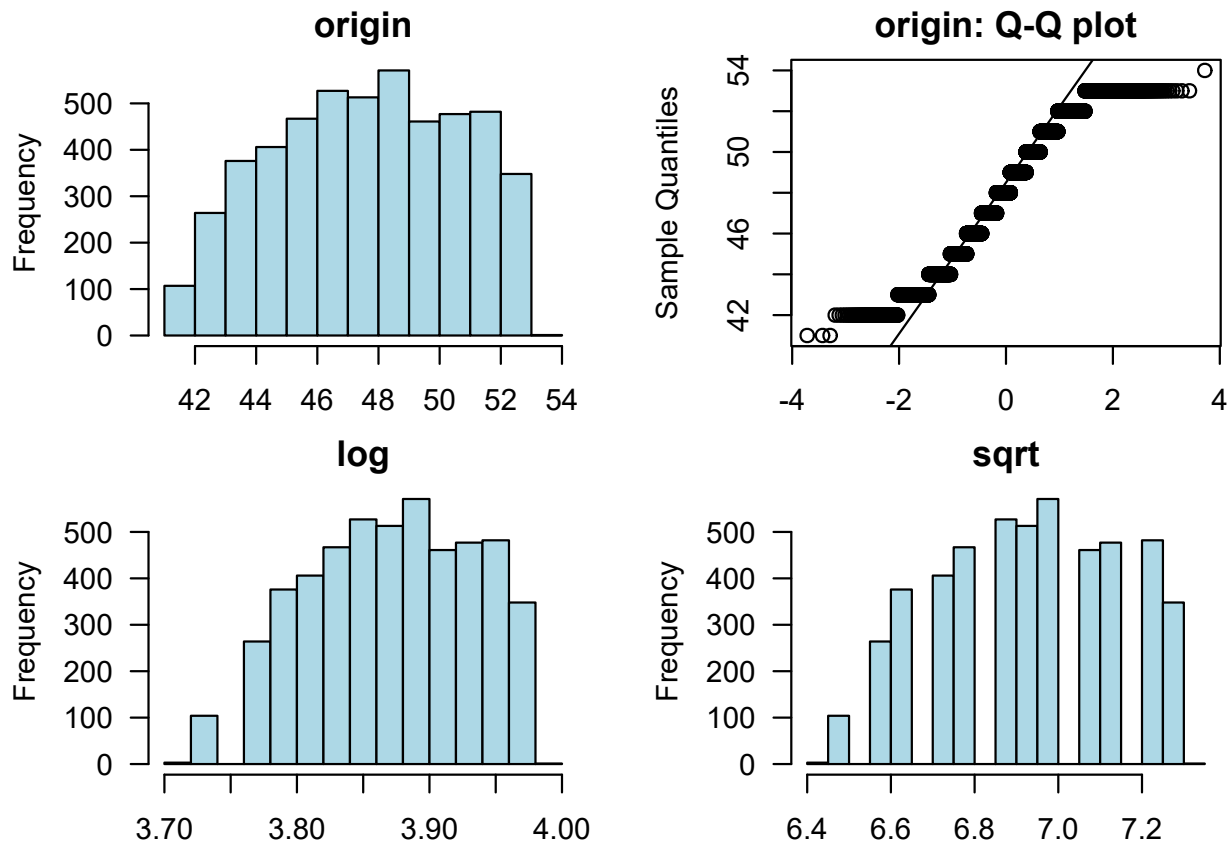


Figure 2.3: birth_yr

age

normality test : Shapiro-Wilk normality test
 statistic : 0.96826, p-value : 7.55026E-32

type	skewness	kurtosis
original	0.2598	2.0862
log transformation	-0.0872	2.0173
sqrt transformation	0.0872	2.0109

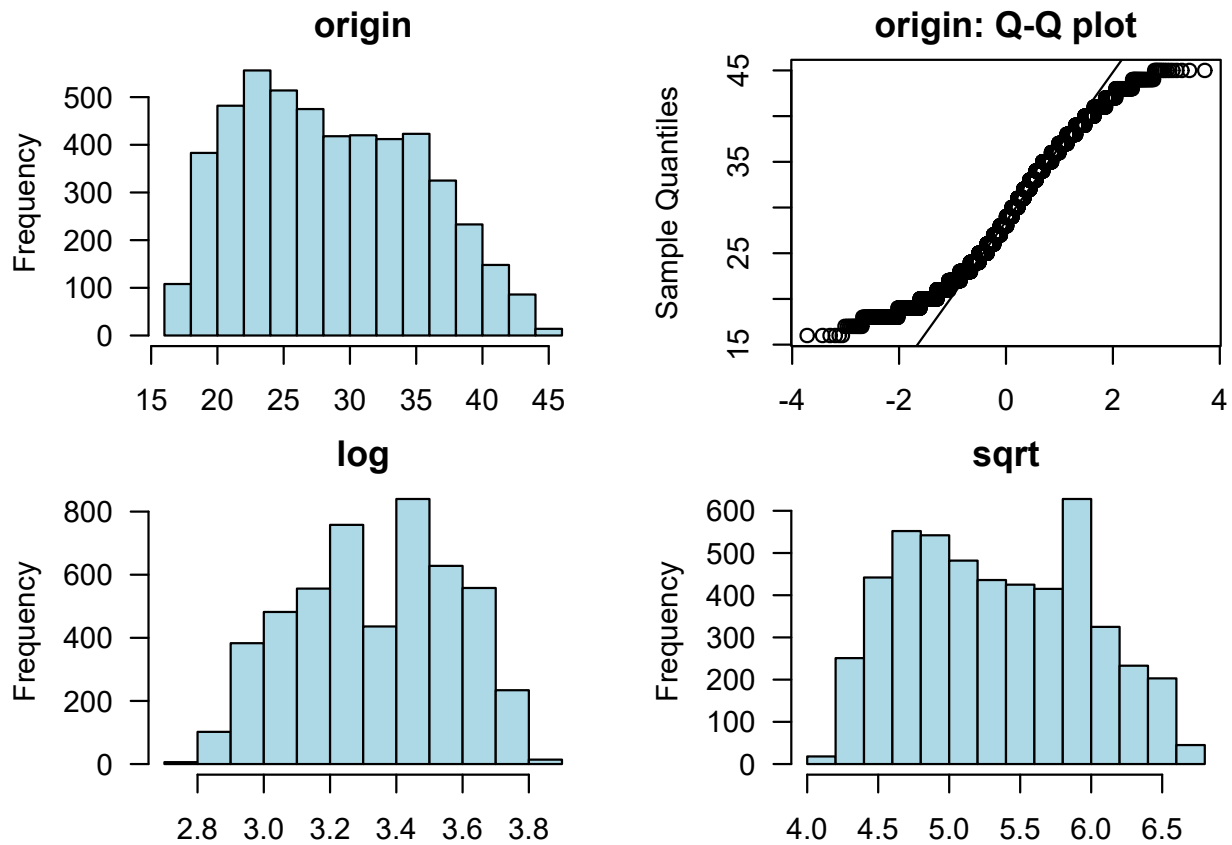


Figure 2.4: age

msp

normality test : Shapiro-Wilk normality test
 statistic : 0.62202, p-value : 3.3721E-74

type	skewness	kurtosis
original	-0.4069	1.1656
log transformation		
sqrt transformation	-0.4069	1.1656

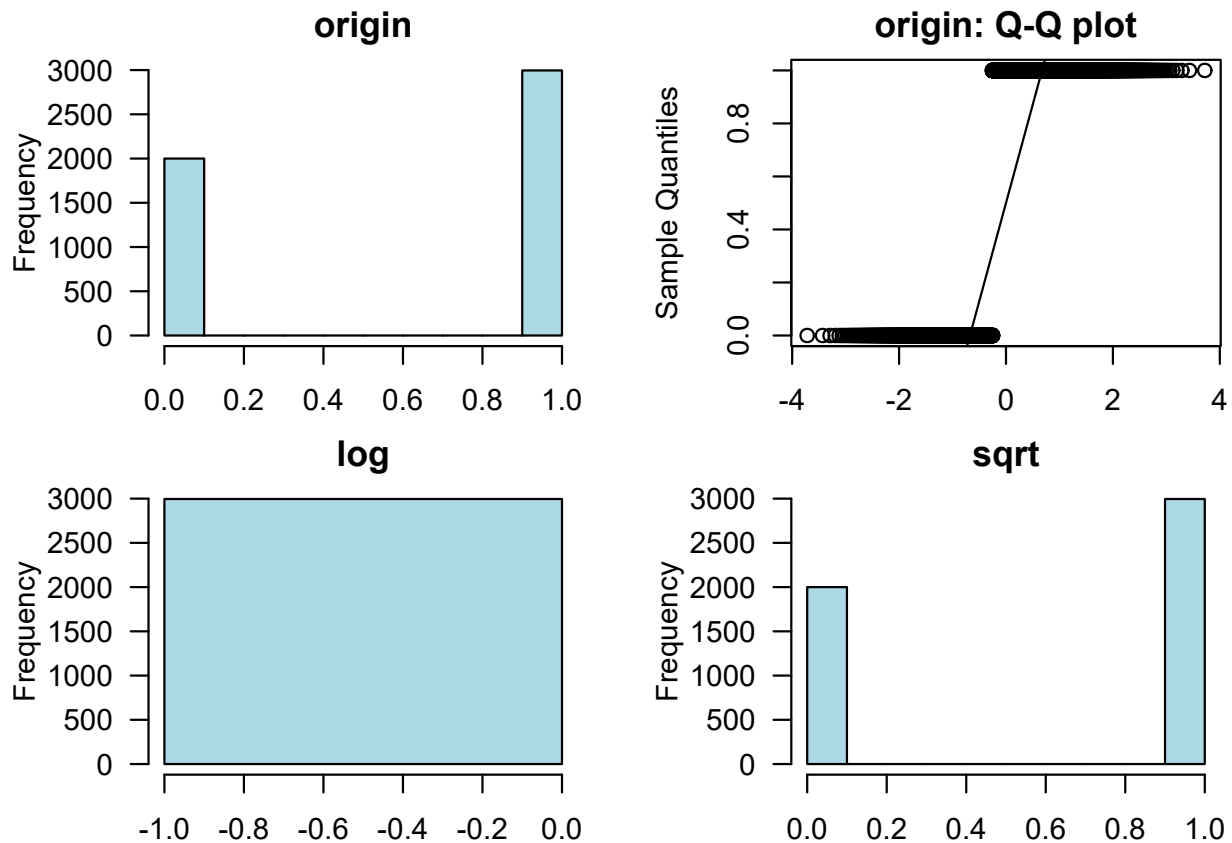


Figure 2.5: msp

nev_mar

normality test : Shapiro-Wilk normality test
 statistic : 0.51379, p-value : 1.56851E-79

type	skewness	kurtosis
original	1.3314	2.7725
log transformation		
sqrt transformation	1.3314	2.7725

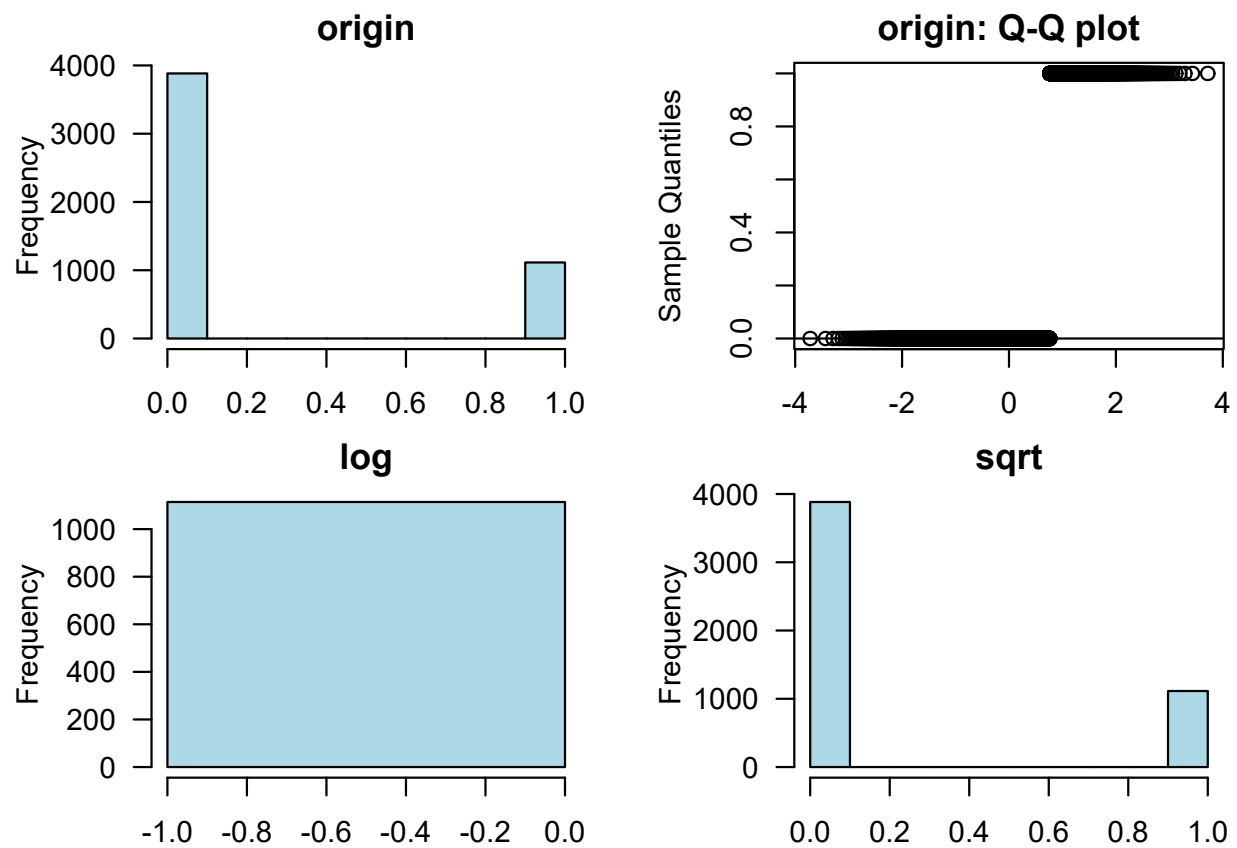


Figure 2.6: nev_mar

grade

normality test : Shapiro-Wilk normality test
 statistic : 0.88975, p-value : 7.25464E-51

type	skewness	kurtosis
original	0.2036	3.9287
log transformation		
sqrt transformation	-0.4752	7.1422

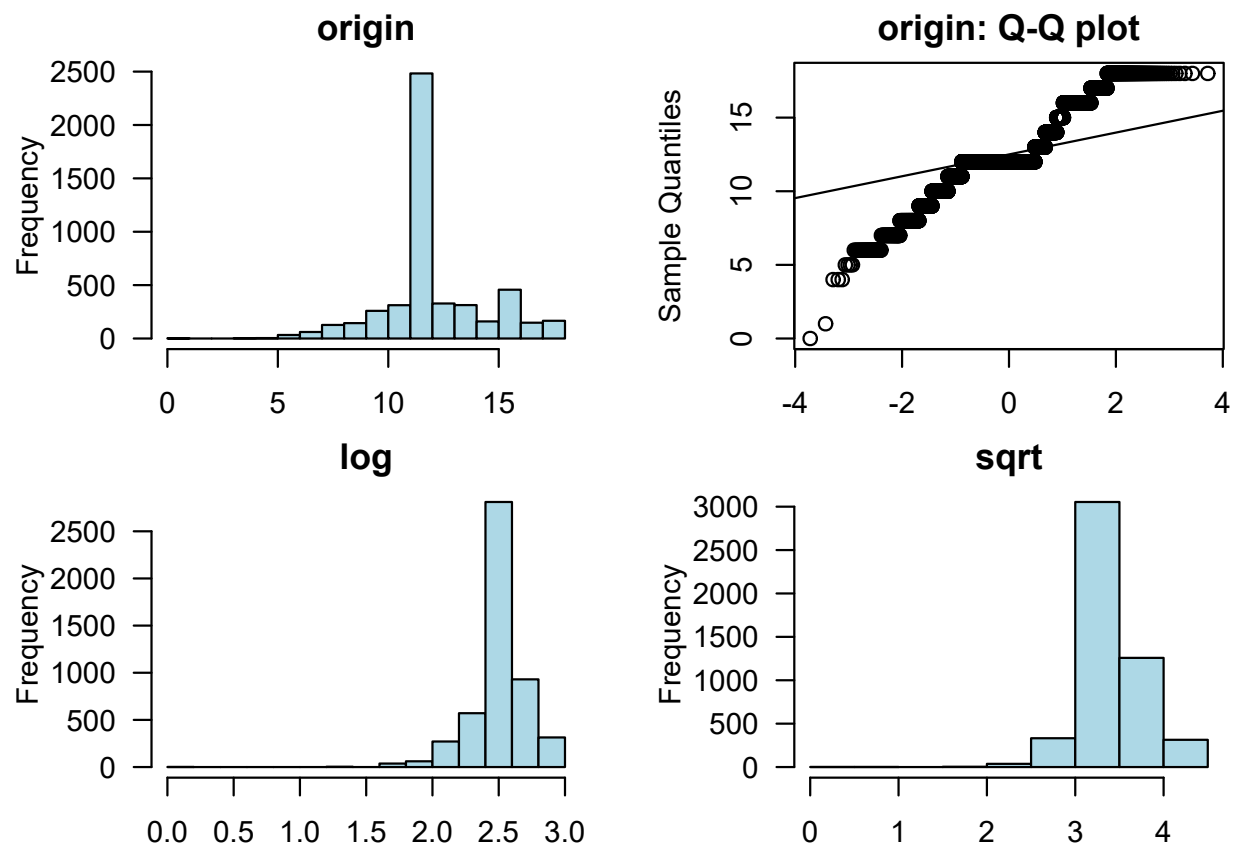


Figure 2.7: grade

collgrad

normality test : Shapiro-Wilk normality test
 statistic : 0.46199, p-value : 9.66261E-82

type	skewness	kurtosis
original	1.6980	3.8831
log transformation		
sqrt transformation	1.6980	3.8831

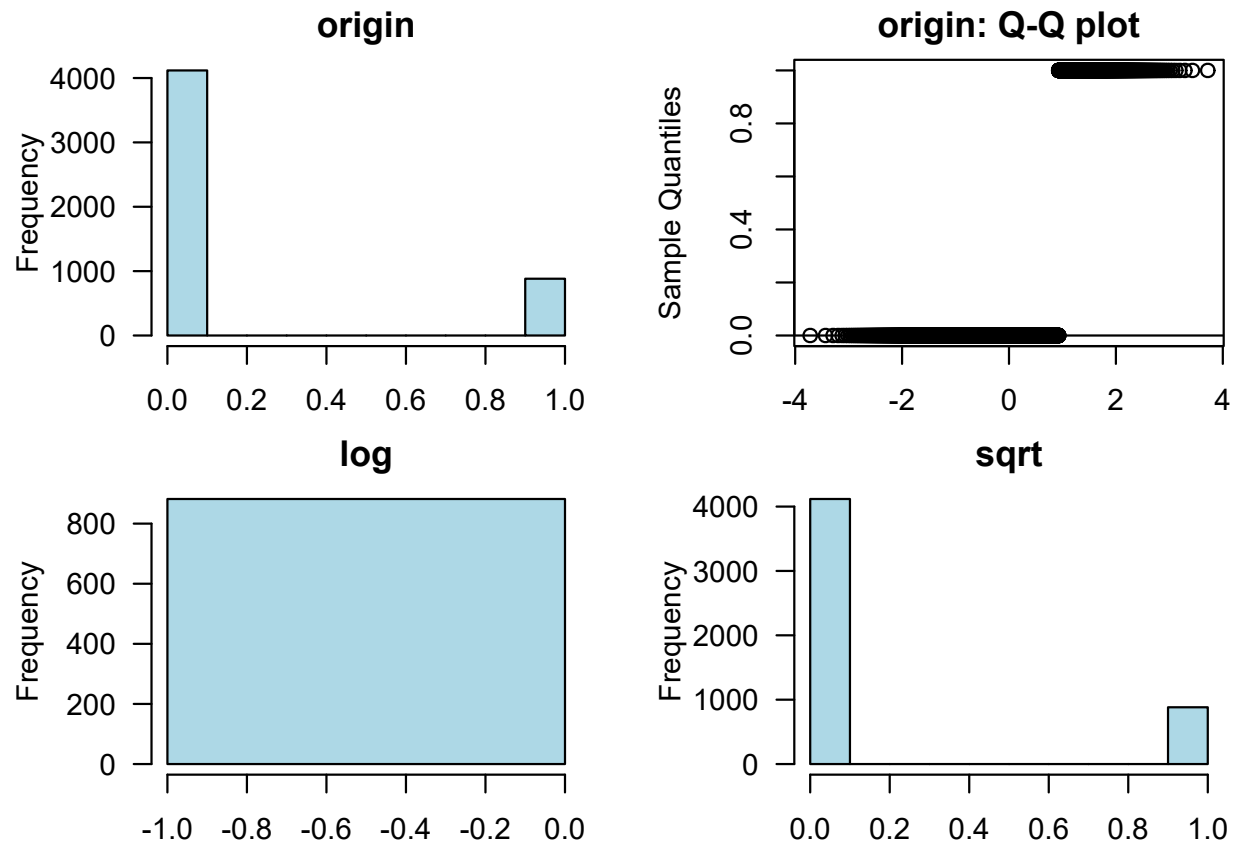


Figure 2.8: collgrad

not_smsa

normality test : Shapiro-Wilk normality test
 statistic : 0.56096, p-value : 2.33271E-77

type	skewness	kurtosis
original	0.9885	1.9772
log transformation		
sqrt transformation	0.9885	1.9772

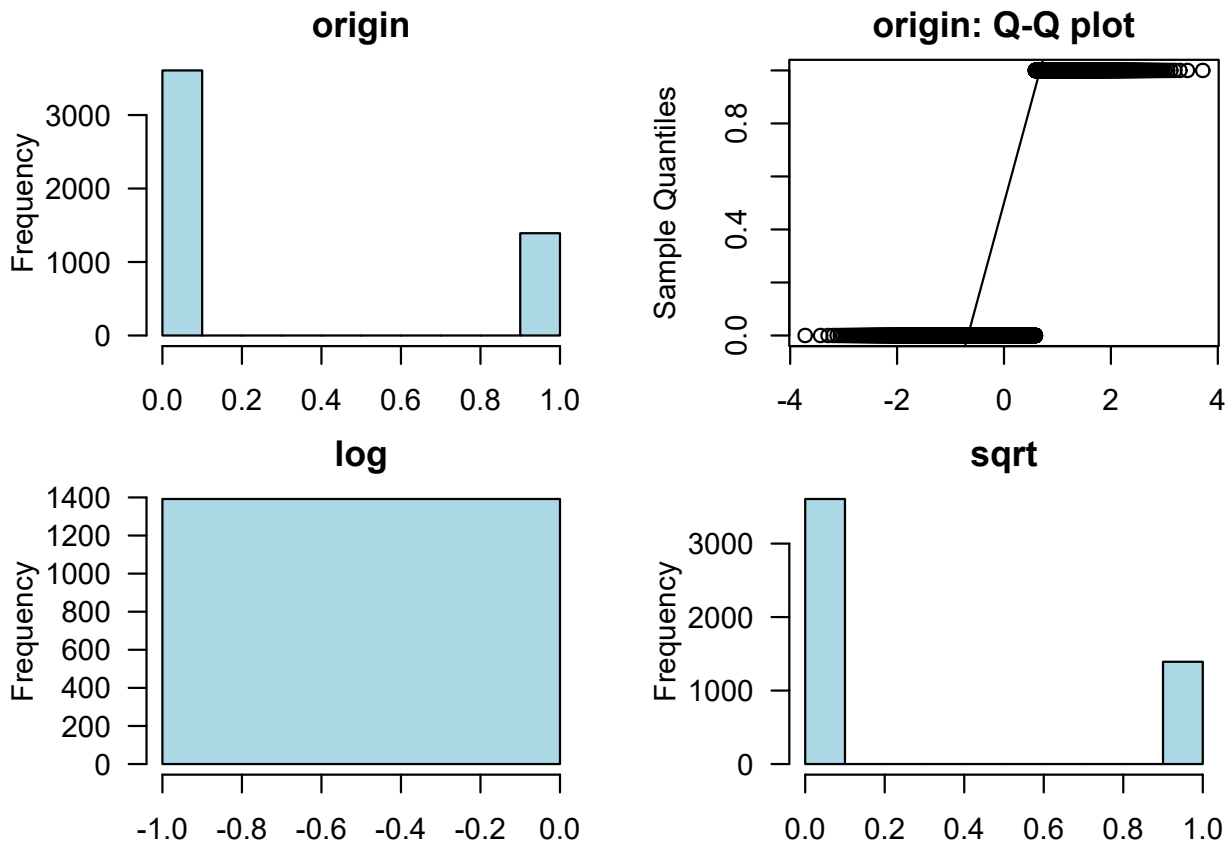


Figure 2.9: not_smsa

c_city

normality test : Shapiro-Wilk normality test
 statistic : 0.60508, p-value : 3.98038E-75

type	skewness	kurtosis
original	0.6085	1.3702
log transformation		
sqrt transformation	0.6085	1.3702

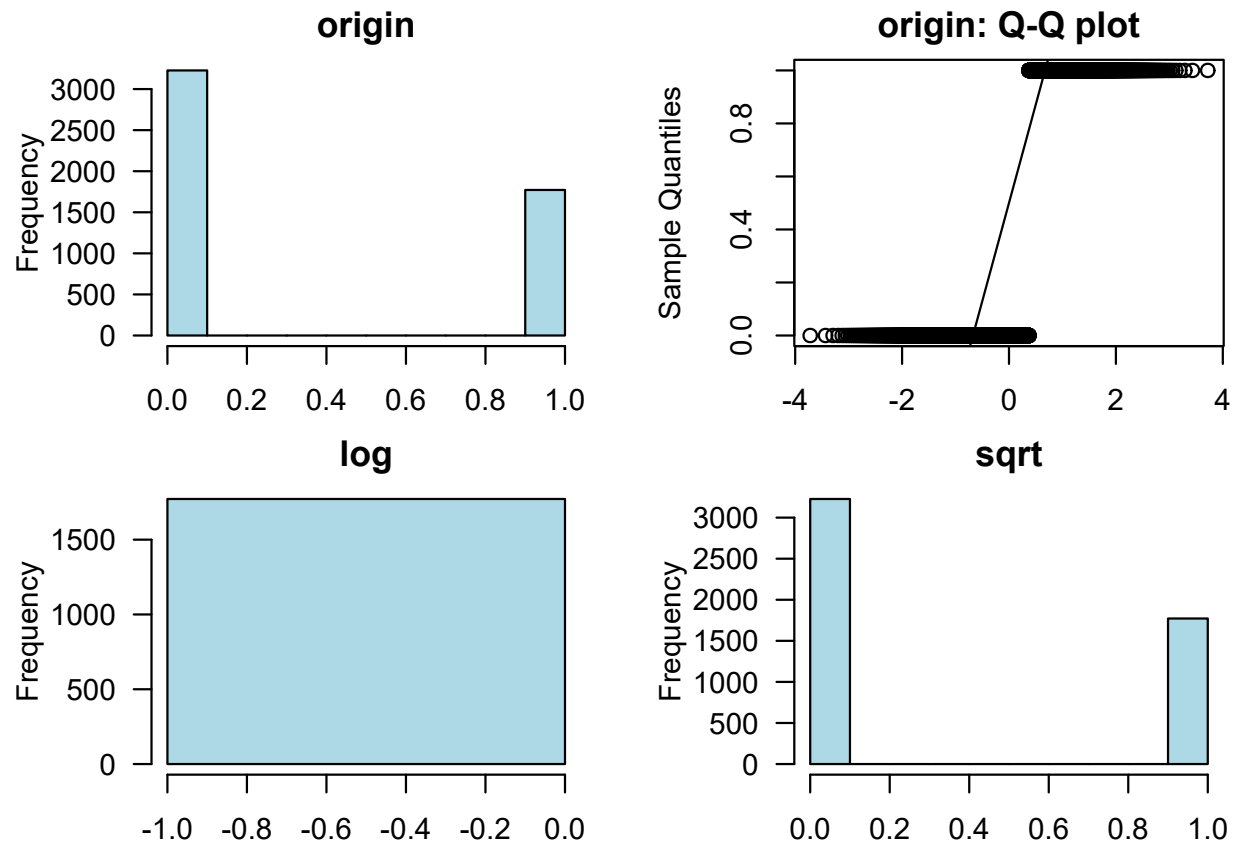


Figure 2.10: c_city

south

normality test : Shapiro-Wilk normality test
 statistic : 0.62252, p-value : 3.55727E-74

type	skewness	kurtosis
original	0.3996	1.1597
log transformation		
sqrt transformation	0.3996	1.1597

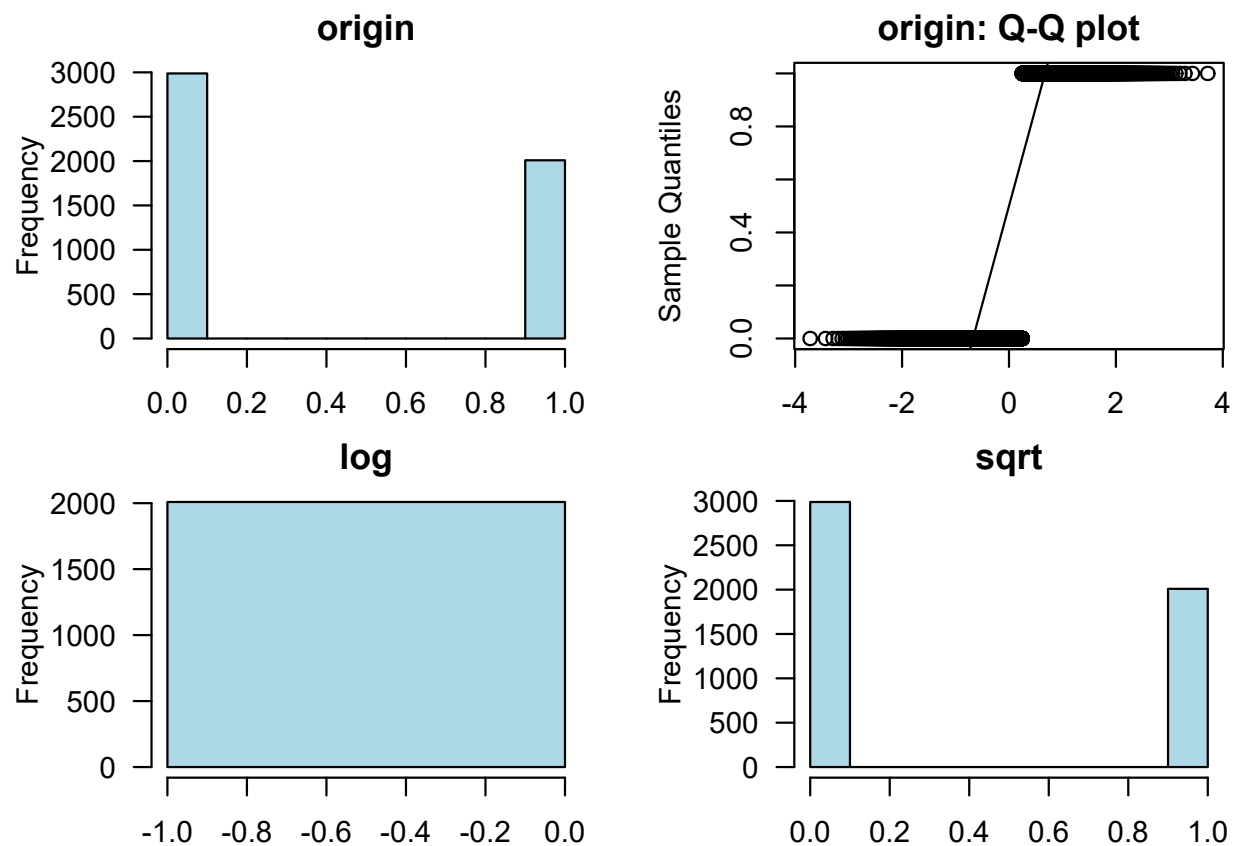


Figure 2.11: south

ind_code

normality test : Shapiro-Wilk normality test
 statistic : 0.87176, p-value : 3.25176E-53

type	skewness	kurtosis
original	-0.0153	1.5413
log transformation	-0.8503	4.2946
sqrt transformation	-0.2787	2.0417

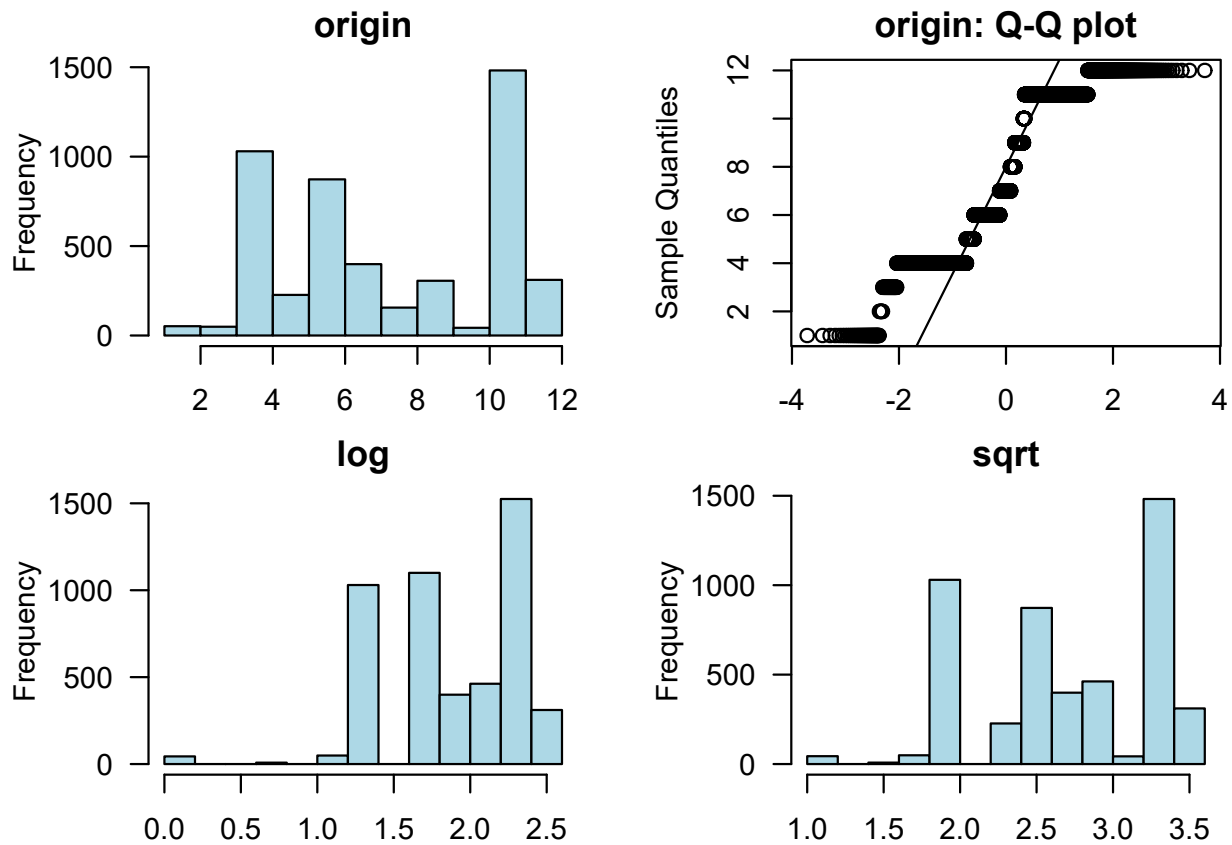


Figure 2.12: ind_code

occ_code

normality test : Shapiro-Wilk normality test
 statistic : 0.85196, p-value : 5.91043E-56

type	skewness	kurtosis
original	1.0908	3.7081
log transformation	-0.3001	2.6692
sqrt transformation	0.4486	2.6404

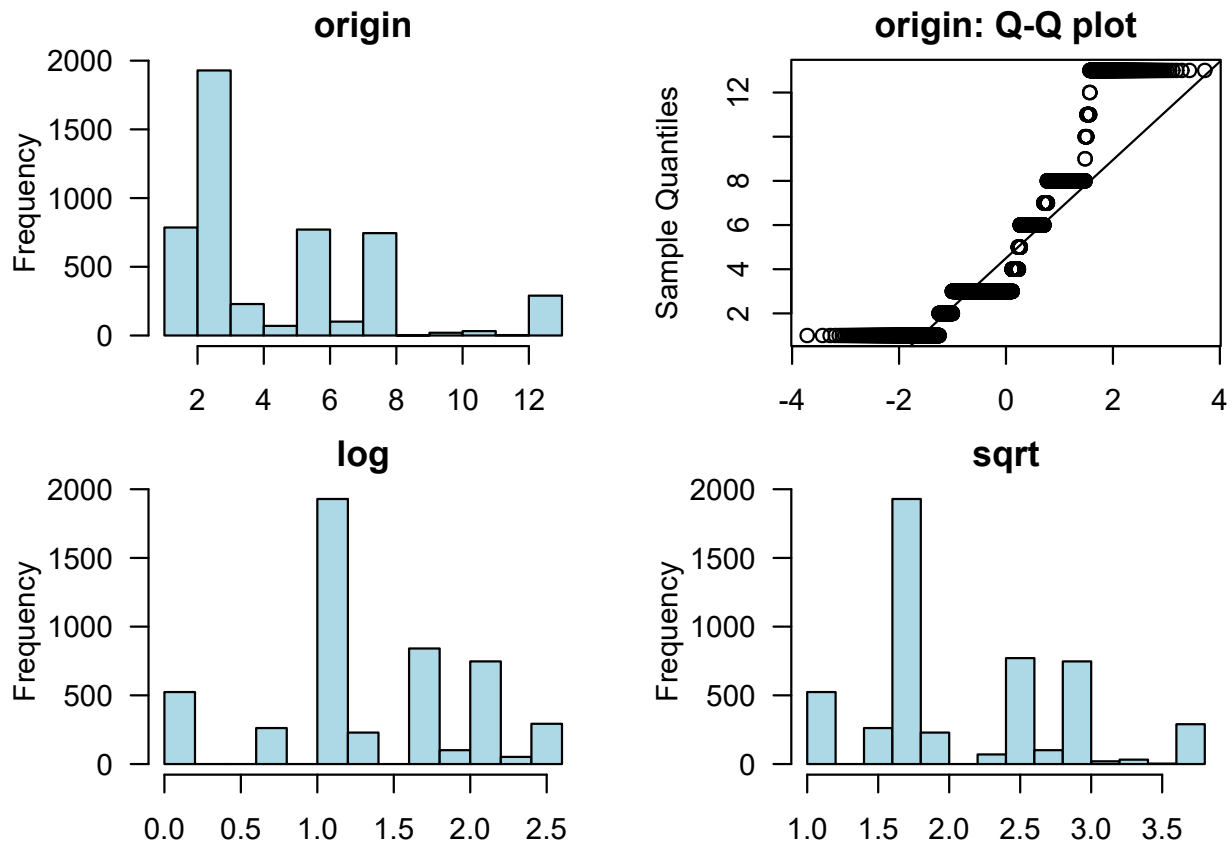


Figure 2.13: occ_code

union

normality test : Shapiro-Wilk normality test
 statistic : 0.51861, p-value : 4.23067E-70

type	skewness	kurtosis
original	1.2973	2.6829
log transformation		
sqrt transformation	1.2973	2.6829

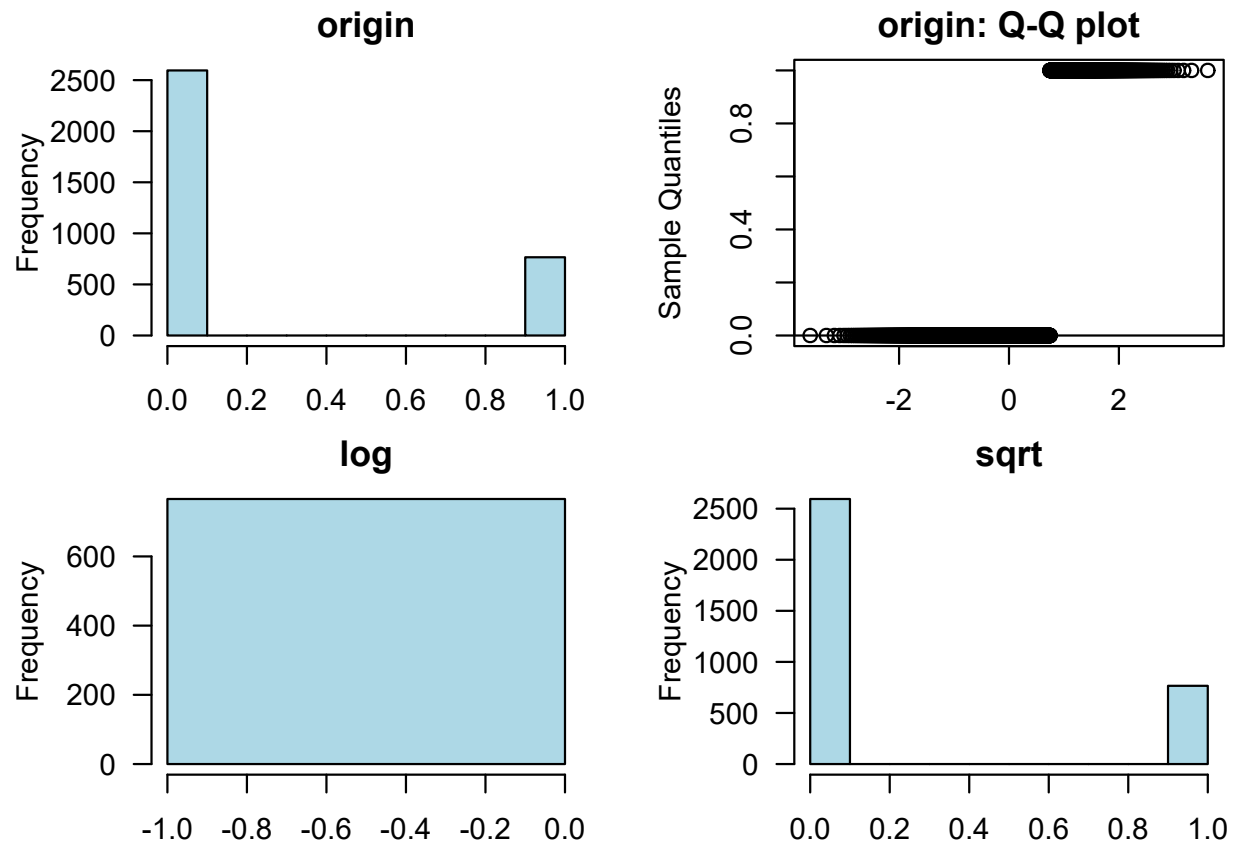


Figure 2.14: union

wks_ue

normality test : Shapiro-Wilk normality test
 statistic : 0.40683, p-value : 2.49059E-78

type	skewness	kurtosis
original	4.0592	22.6931
log transformation		
sqrt transformation	2.3165	7.8738

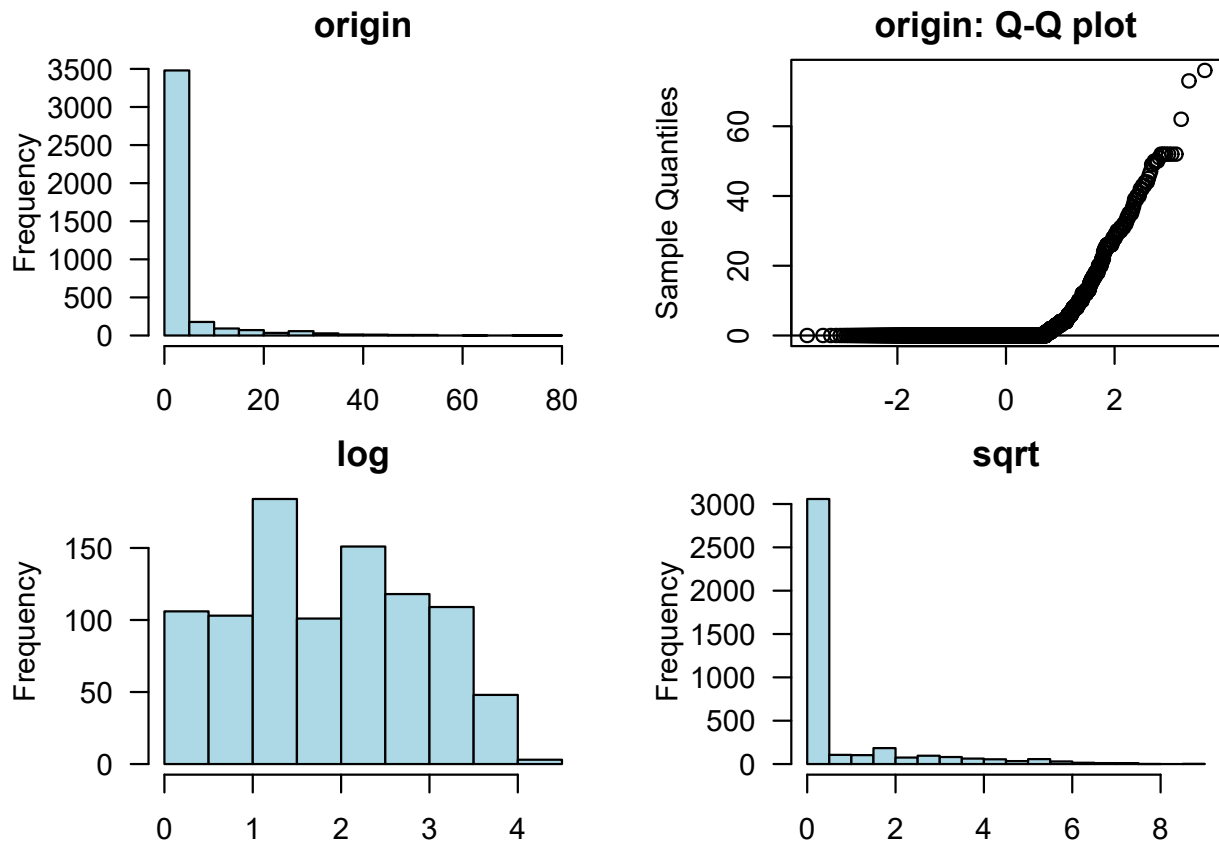


Figure 2.15: wks_ue

ttl_exp

normality test : Shapiro-Wilk normality test
 statistic : 0.9246, p-value : 1.41021E-44

type	skewness	kurtosis
original	0.8471	3.0397
log transformation		
sqrt transformation	0.1424	2.2231

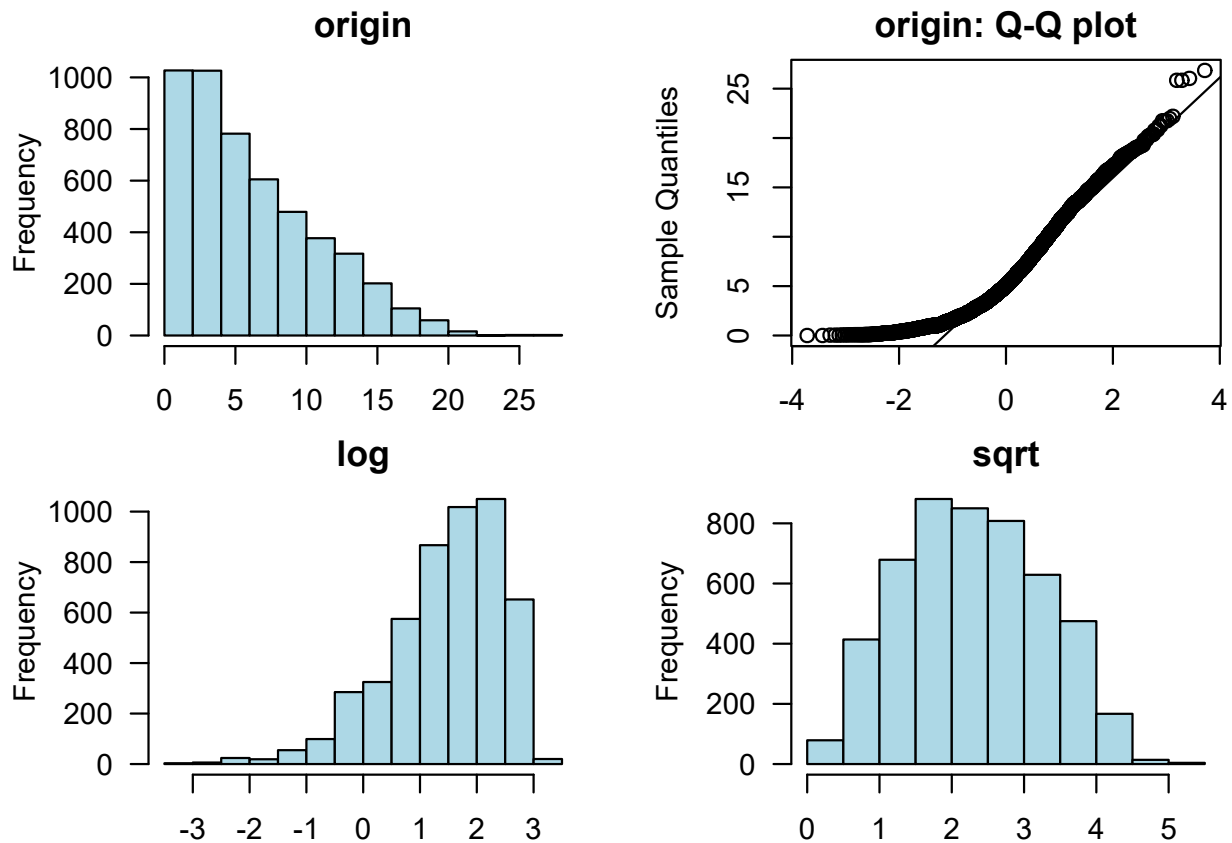


Figure 2.16: ttl_exp

tenure

normality test : Shapiro-Wilk normality test
 statistic : 0.76529, p-value : 3.11793E-64

type	skewness	kurtosis
original	1.9250	6.7410
log transformation		
sqrt transformation	0.7737	3.0973

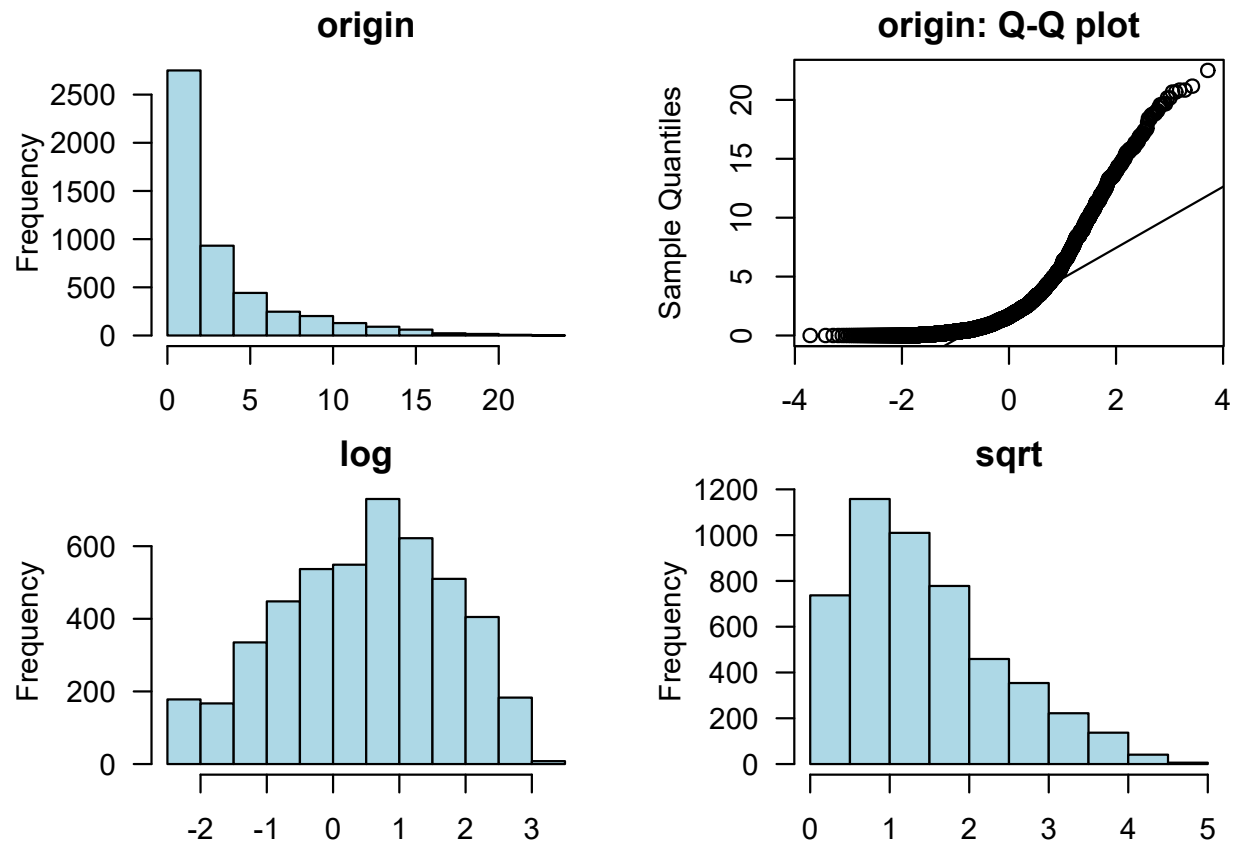


Figure 2.17: tenure

hours

normality test : Shapiro-Wilk normality test
 statistic : 0.76794, p-value : 2.20078E-64

type	skewness	kurtosis
original	-0.4788	11.8507
log transformation	-3.2248	16.6093
sqrt transformation	-1.8604	8.6082

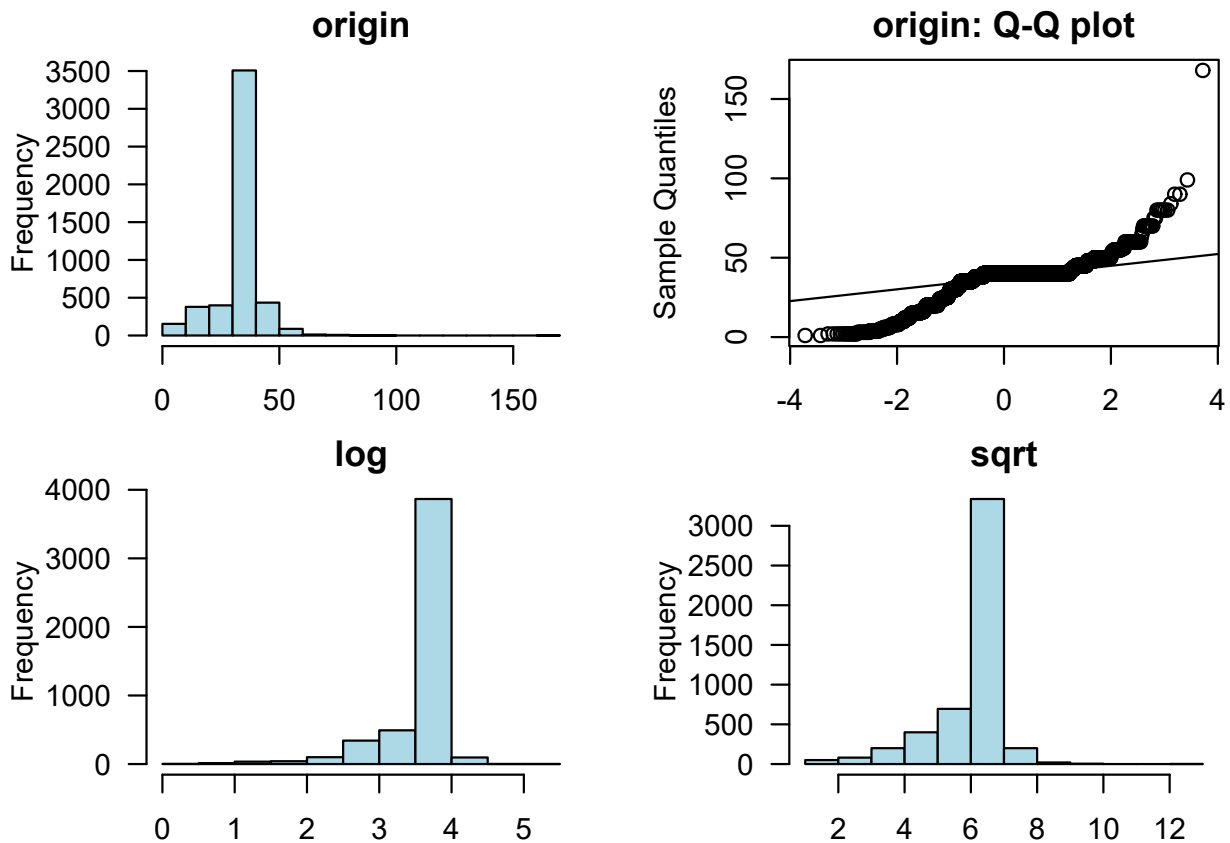


Figure 2.18: hours

wks_work

normality test : Shapiro-Wilk normality test
 statistic : 0.93791, p-value : 3.83498E-41

type	skewness	kurtosis
original	0.1763	2.3197
log transformation		
sqrt transformation	-0.8122	3.6663

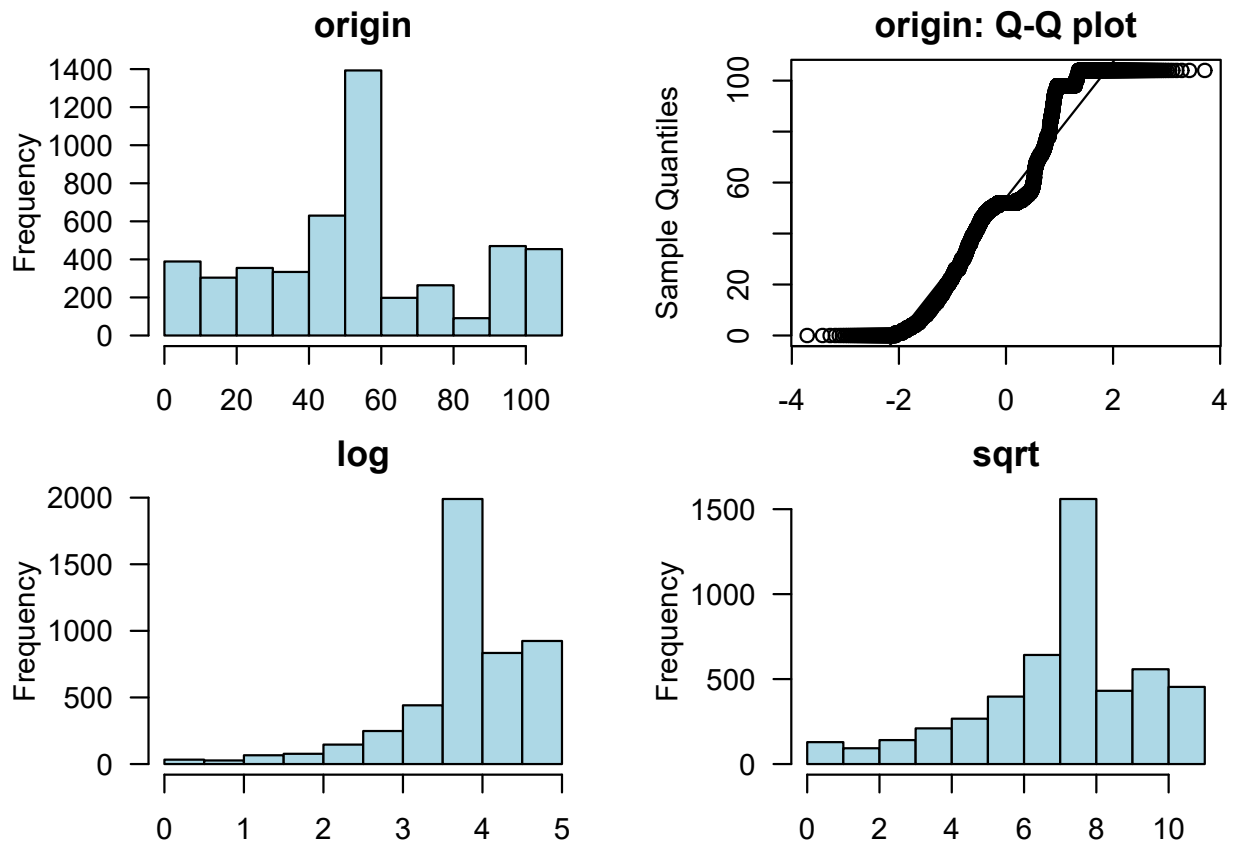


Figure 2.19: wks_work

ln_wage

normality test : Shapiro-Wilk normality test
statistic : 0.98572, p-value : 4.30583E-22

type	skewness	kurtosis
original	0.2035	4.1959
log transformation	-3.5222	31.7807
sqrt transformation	-0.8083	6.2651

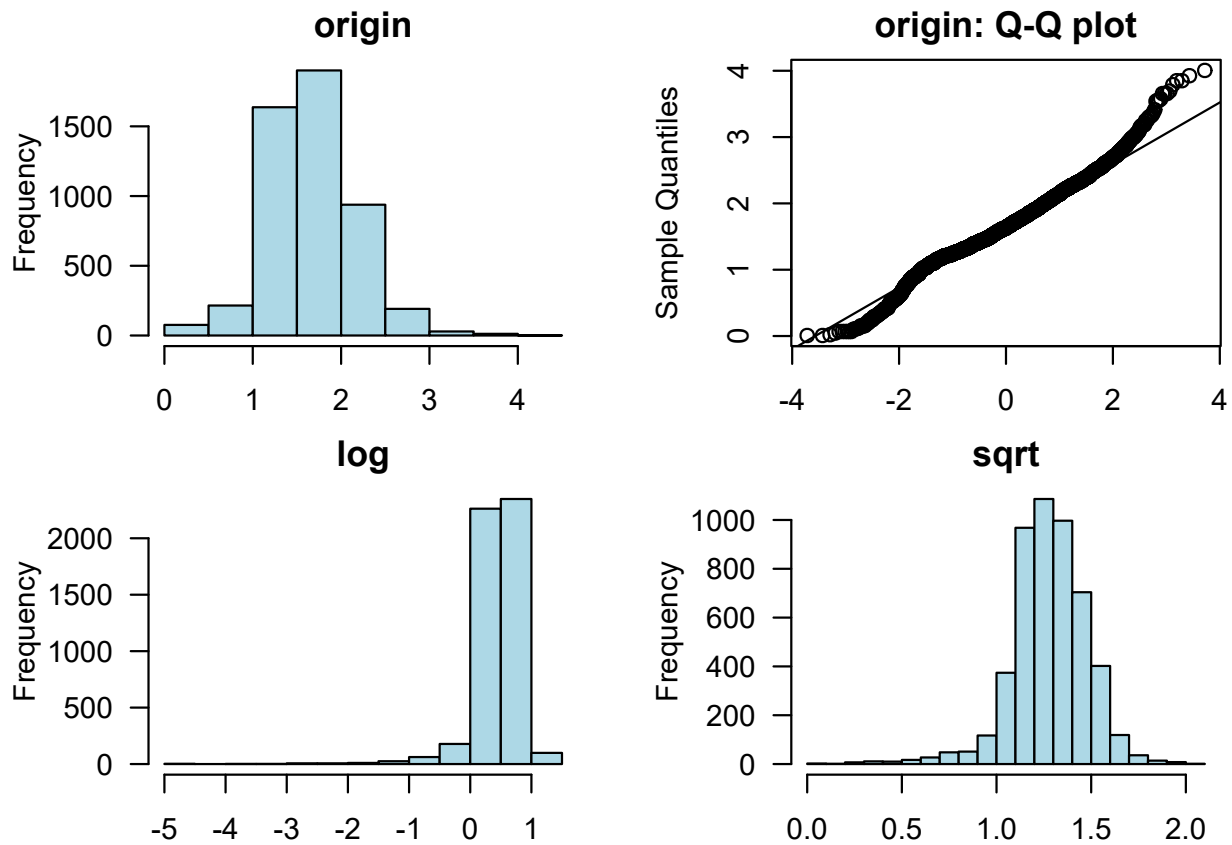


Figure 2.20: ln_wage

Chapter 3

Relationship Between Variables

3.1 Correlation Coefficient

3.1.1 Correlation Coefficient by Variable Combination

Table 3.1: The correlation coefficients (0.5 or more)

Variable1	Variable2	Correlation Coefficient
age	year	0.895
ttl_exp	year	0.777
collgrad	grade	0.757
ttl_exp	age	0.756
tenure	ttl_exp	0.674
nev_mar	msp	-0.673
wks_work	ttl_exp	0.630
wks_work	year	0.565
wks_work	age	0.525

3.1.2 Correlation Plot of Numerical Variables

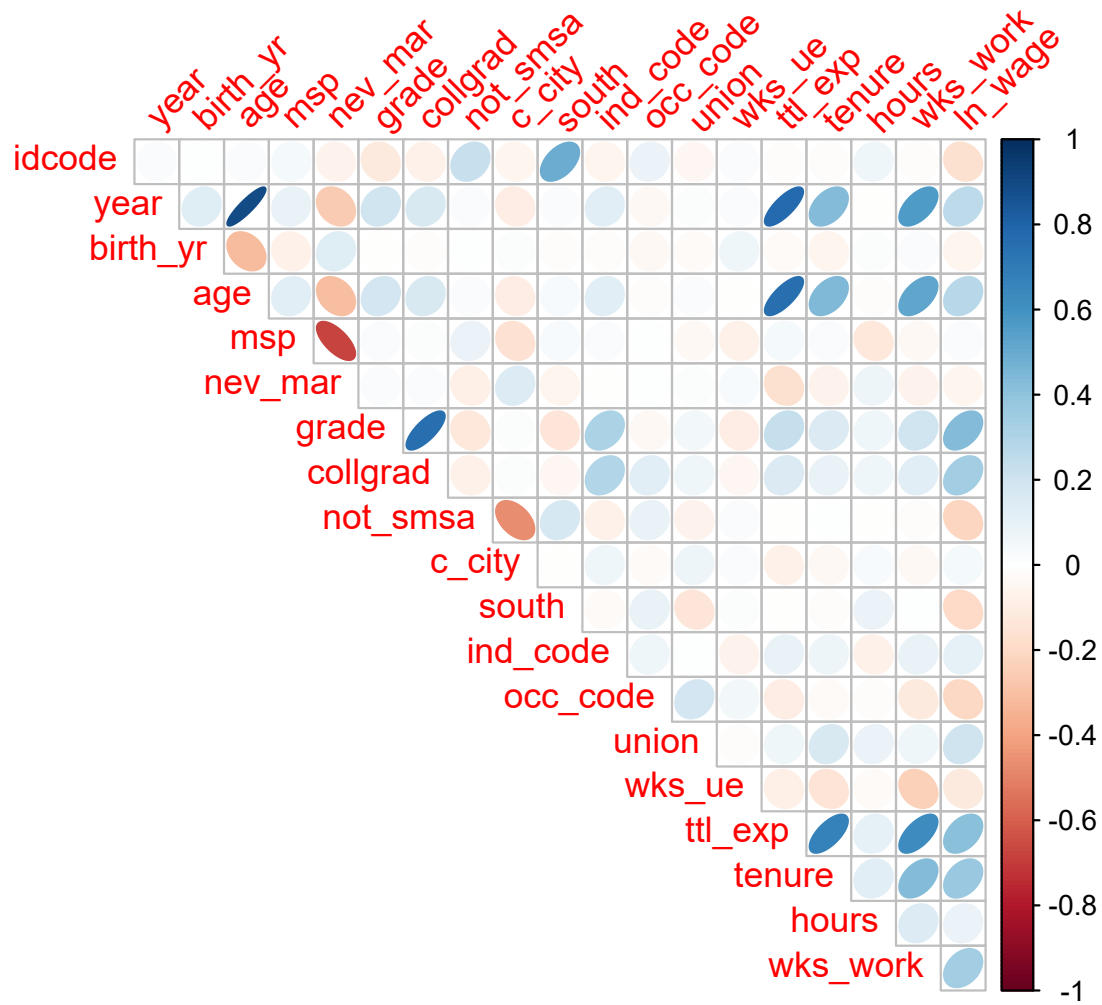


Figure 3.1: The correlation coefficient of numerical variables

Chapter 4

Target based Analysis

4.1 Grouped Descriptive Statistics

4.1.1 Grouped Numerical Variables

There is no target variable.

4.1.2 Grouped Categorical Variables

There is no target variable.

4.2 Grouped Relationship Between Variables

4.2.1 Grouped Correlation Coefficient

There is no target variable.

4.2.2 Grouped Correlation Plot of Numerical Variables

There is no target variable.