



REPORT SERIES WITH DLOOKR

---

# Exploratory Data Analysis Report

---

*Author:*  
dlookr package

*Version:*  
0.3.12

November 3, 2020

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Information of Dataset	3
1.2	Information of Variables	3
1.3	About EDA Report	4
<b>2</b>	<b>Univariate Analysis</b>	<b>5</b>
2.1	Descriptive Statistics	5
2.2	Normality Test of Numerical Variables	6
2.2.1	Statistics and Visualization of (Sample) Data	6
<b>3</b>	<b>Relationship Between Variables</b>	<b>27</b>
3.1	Correlation Coefficient	27
3.1.1	Correlation Coefficient by Variable Combination	27
3.1.2	Correlation Plot of Numerical Variables	27
<b>4</b>	<b>Target based Analysis</b>	<b>29</b>
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 : 0.95578, p-value : 1.80748E-36

type	skewness	kurtosis
original	-0.0069	1.8157
log transformation	-2.0152	8.8369
sqrt transformation	-0.5824	2.4626

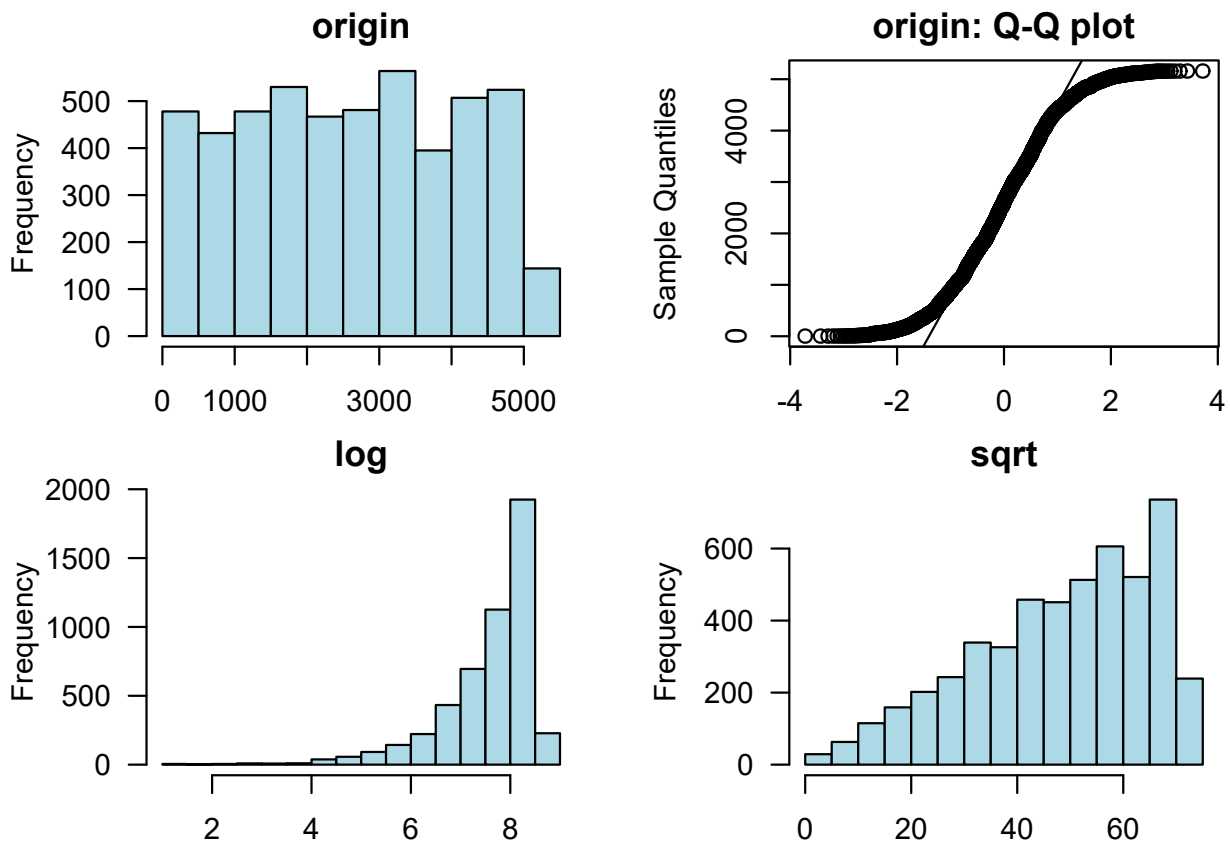


Figure 2.1: idcode

**year**

normality test : Shapiro-Wilk normality test  
 statistic : 0.93348, p-value : 1.34509E-42

type	skewness	kurtosis
original	0.0825	1.7129
log transformation	-0.0037	1.7084
sqrt transformation	0.0394	1.7086

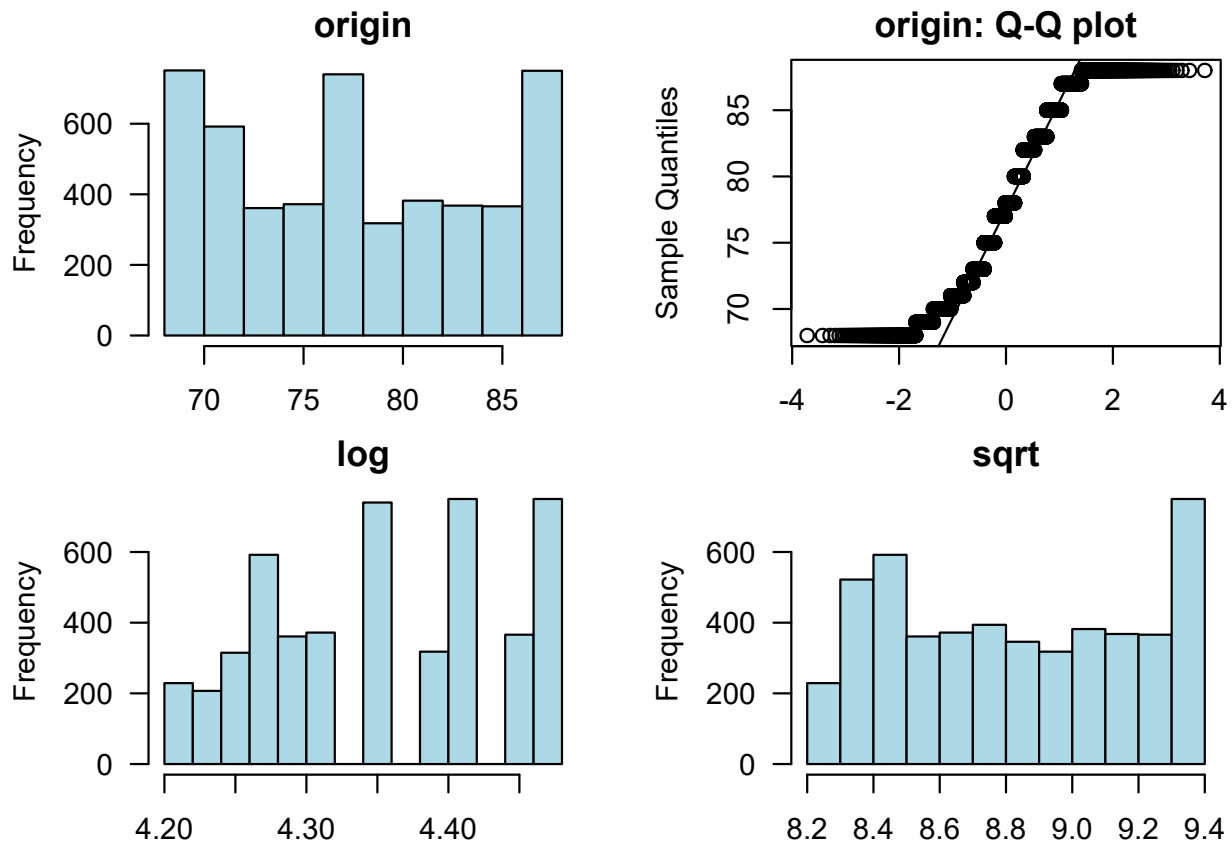


Figure 2.2: year

**birth\_yr**

normality test : Shapiro-Wilk normality test  
 statistic : 0.95932, p-value : 2.79664E-35

type	skewness	kurtosis
original	-0.0906	1.9831
log transformation	-0.1849	2.0320
sqrt transformation	-0.1375	2.0048

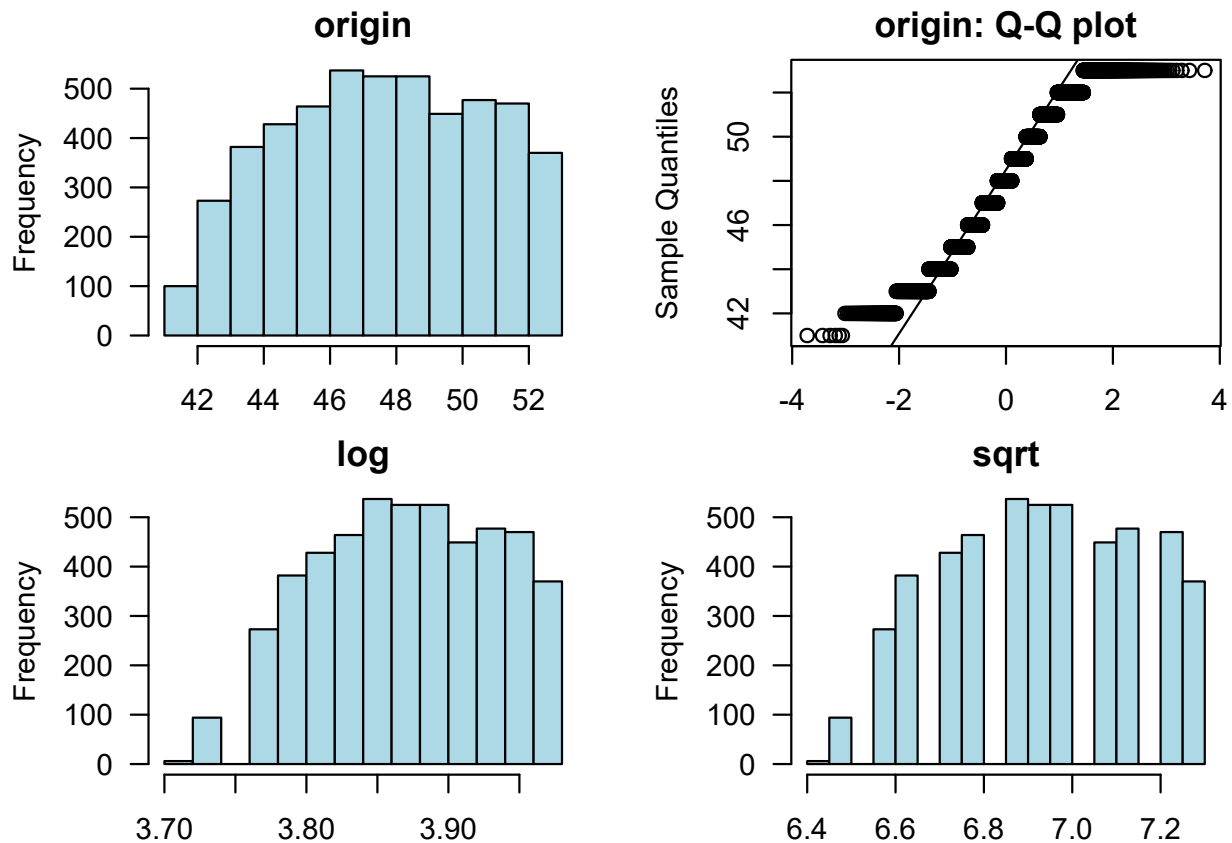


Figure 2.3: birth\_yr



age

normality test : Shapiro-Wilk normality test  
 statistic : 0.96678, p-value : 1.82097E-32

type	skewness	kurtosis
original	0.2642	2.0643
log transformation	-0.0795	2.0085
sqrt transformation	0.0938	1.9961

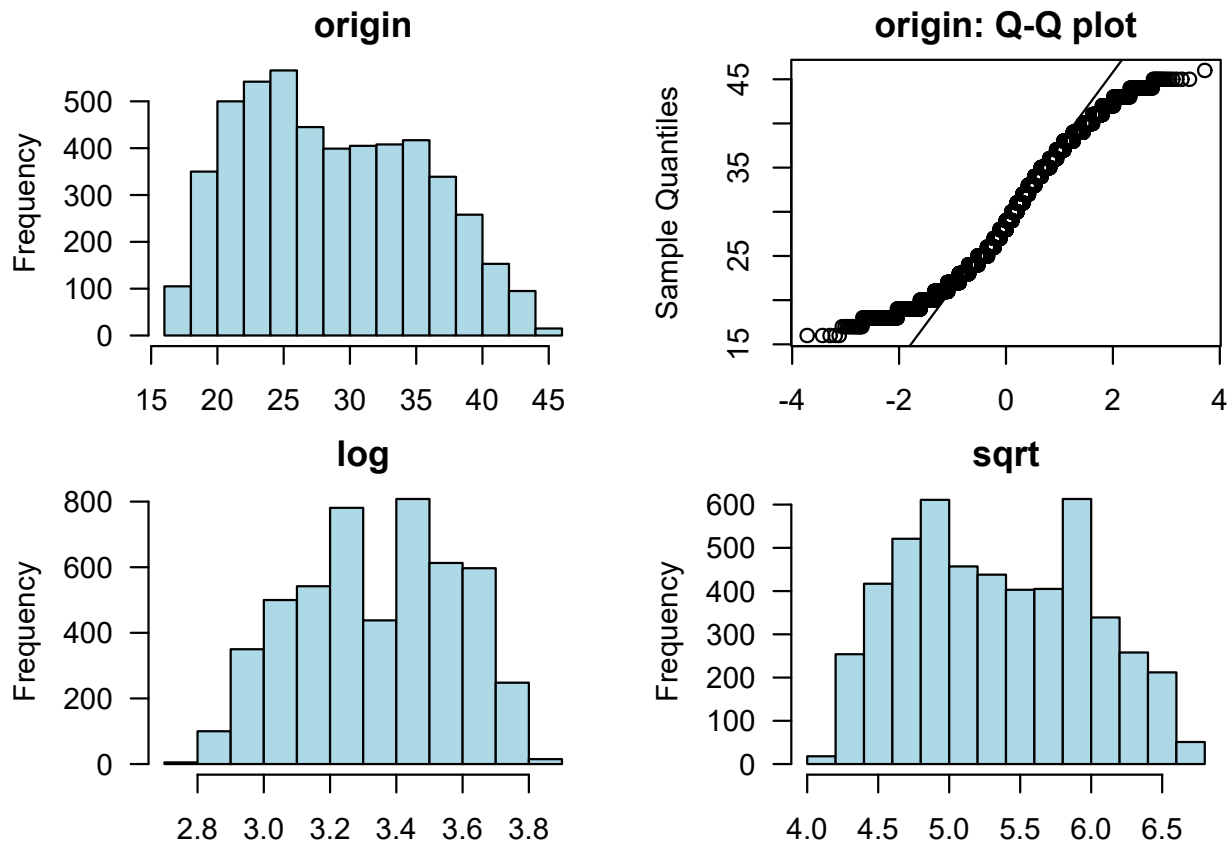


Figure 2.4: age

**msp**

normality test : Shapiro-Wilk normality test  
 statistic : 0.61989, p-value : 2.47088E-74

type	skewness	kurtosis
original	-0.4364	1.1905
log transformation		
sqrt transformation	-0.4364	1.1905

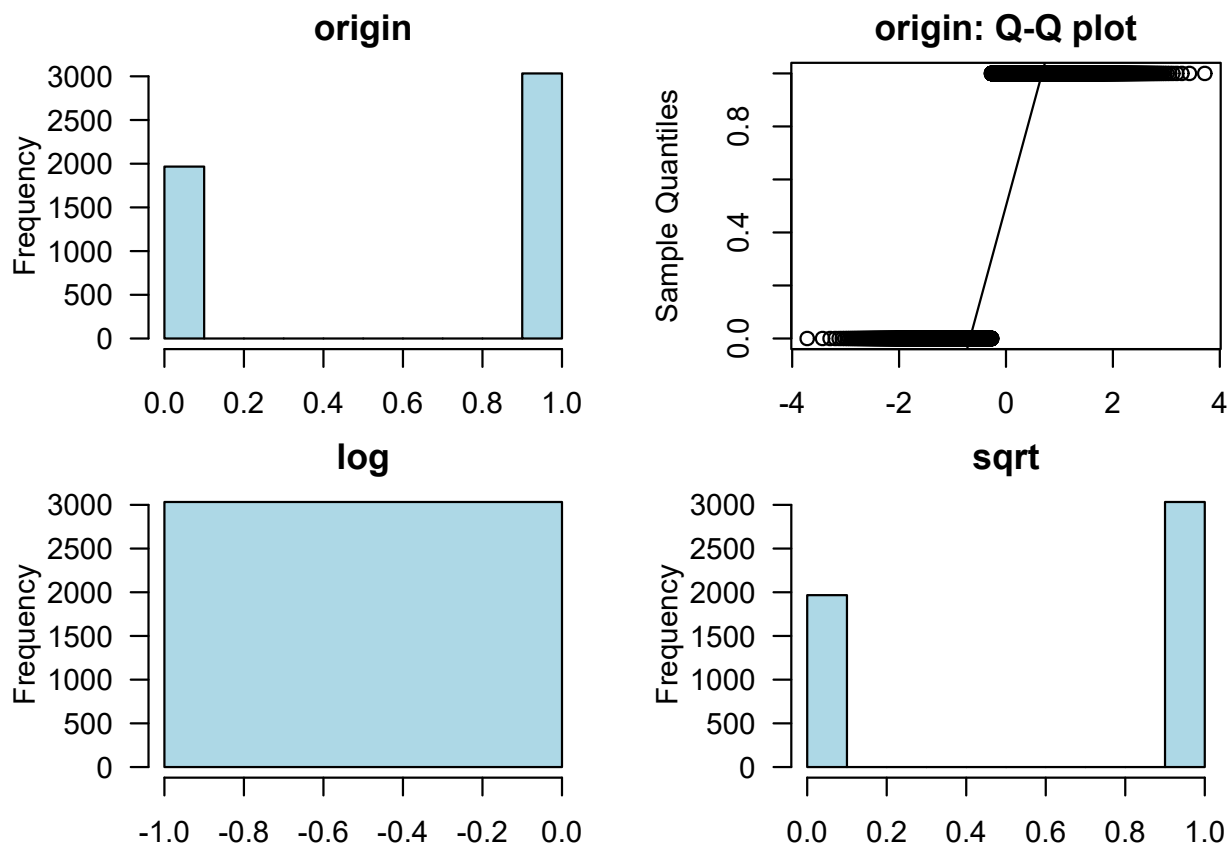


Figure 2.5: msp

**nev\_mar**

normality test : Shapiro-Wilk normality test  
 statistic : 0.51653, p-value : 2.07391E-79

type	skewness	kurtosis
original	1.3121	2.7215
log transformation		
sqrt transformation	1.3121	2.7215

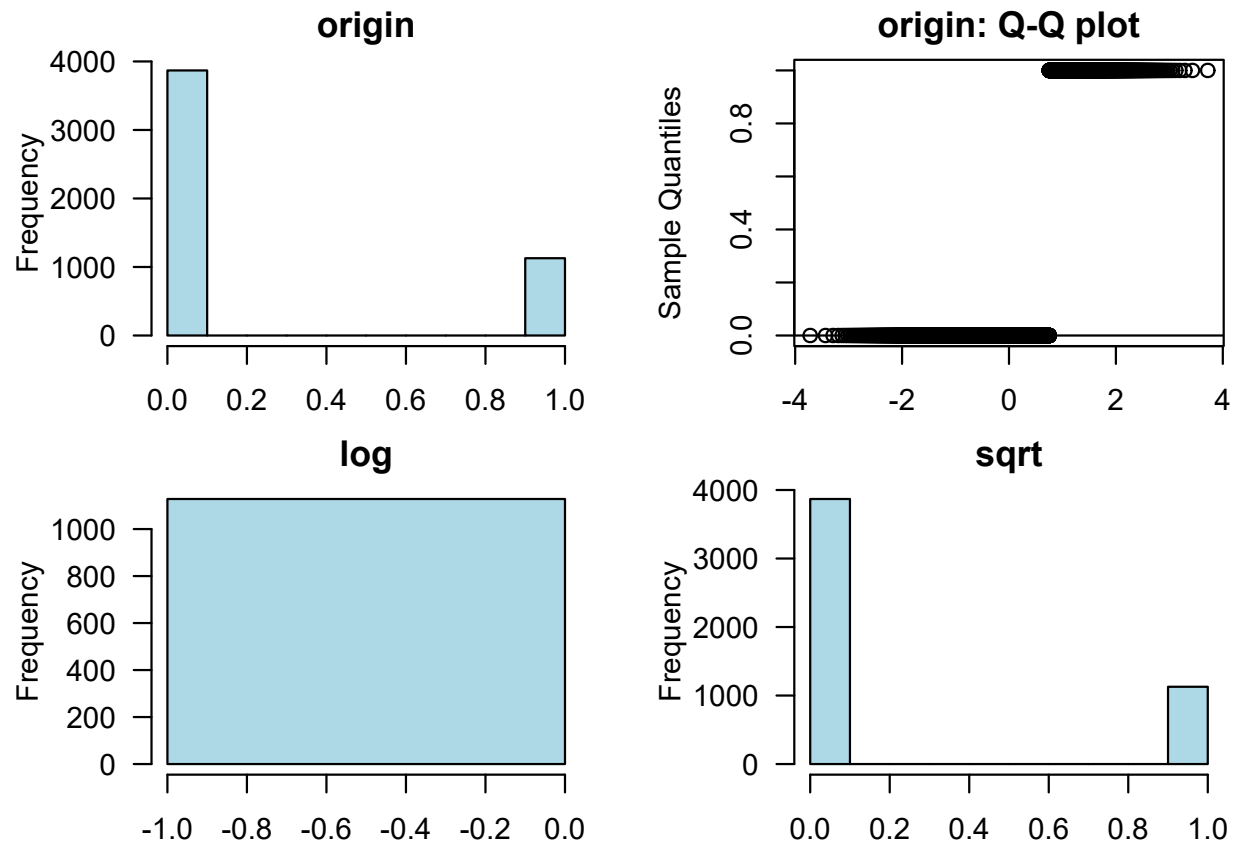


Figure 2.6: nev\_mar

**grade**

normality test : Shapiro-Wilk normality test  
 statistic : 0.87547, p-value : 5.76271E-53

type	skewness	kurtosis
original	0.0957	4.4801
log transformation		
sqrt transformation	-1.0308	12.2803

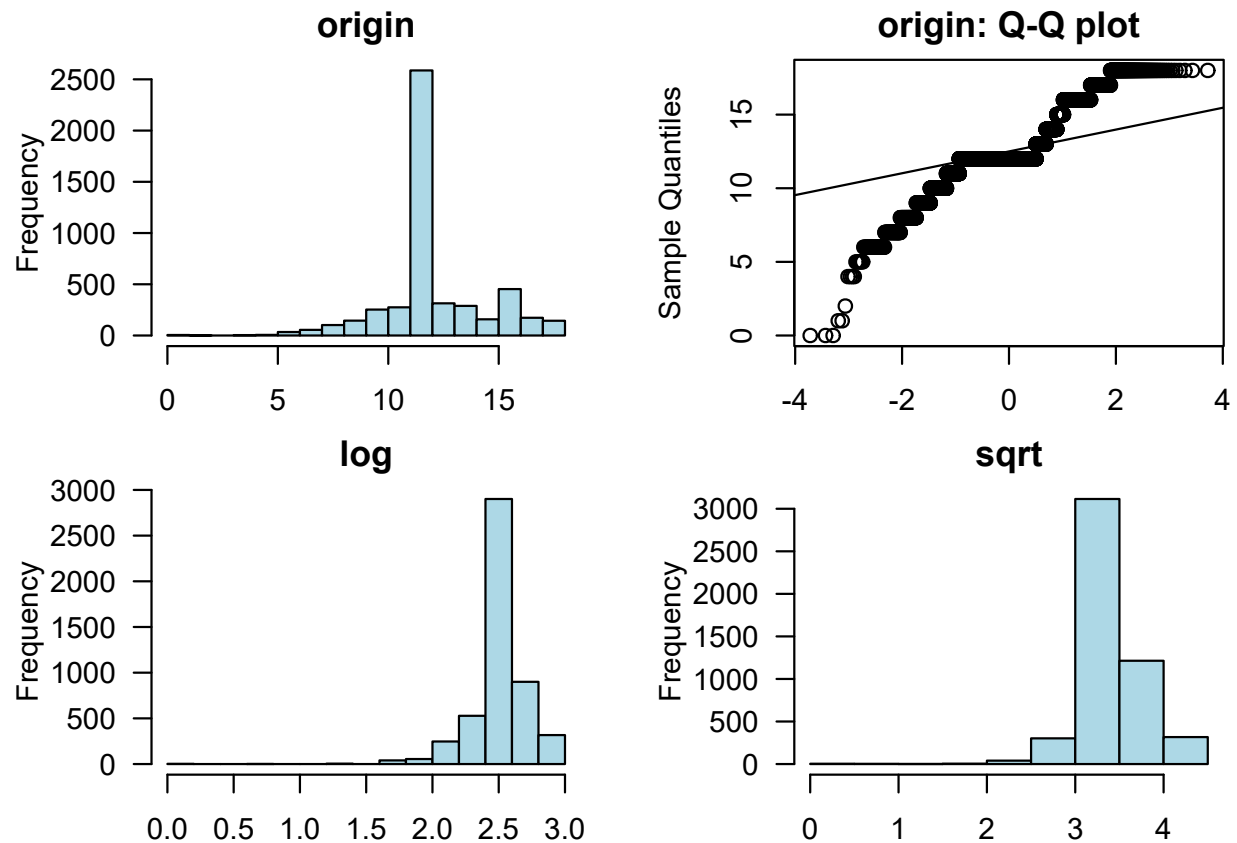


Figure 2.7: grade

**collgrad**

normality test : Shapiro-Wilk normality test  
 statistic : 0.44642, p-value : 2.29561E-82

type	skewness	kurtosis
original	1.8109	4.2795
log transformation		
sqrt transformation	1.8109	4.2795

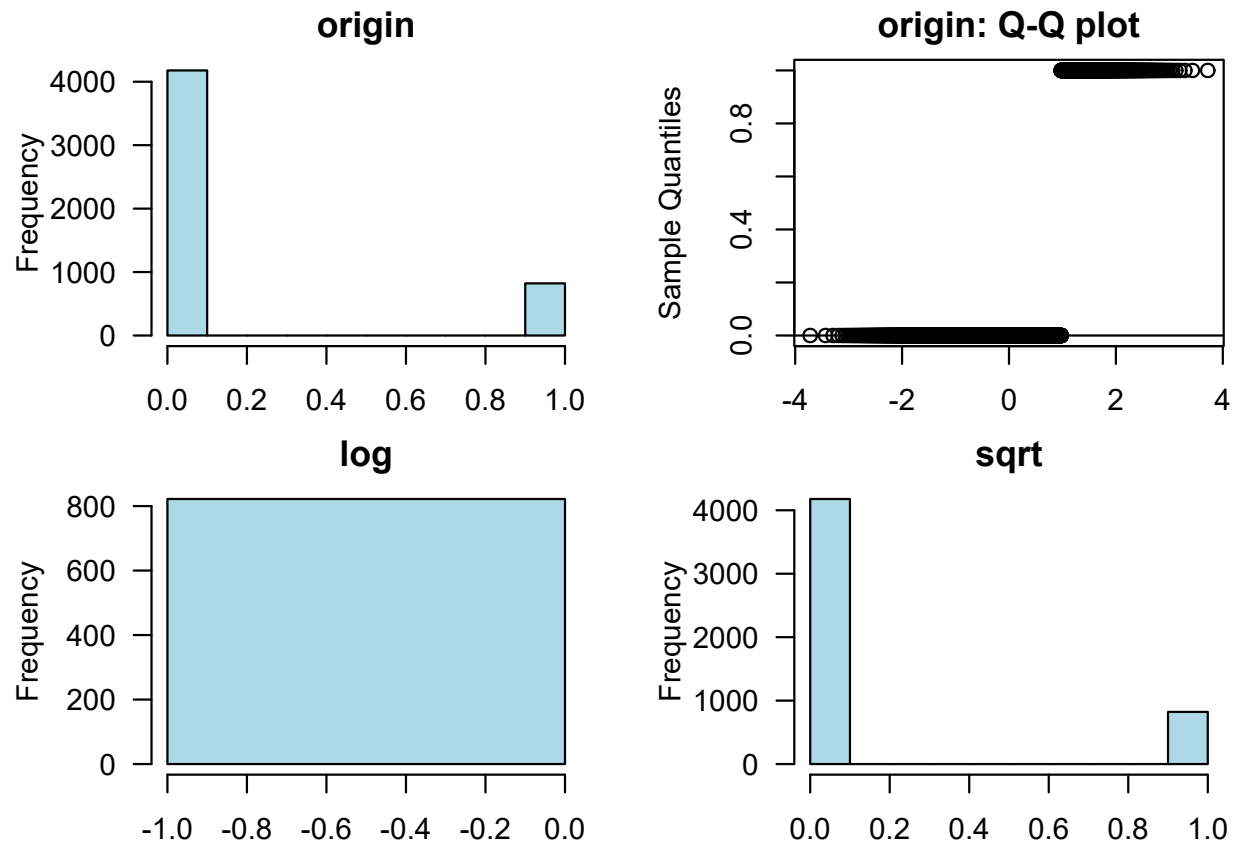


Figure 2.8: collgrad

**not\_smsa**

normality test : Shapiro-Wilk normality test  
 statistic : 0.55053, p-value : 7.39322E-78

type	skewness	kurtosis
original	1.0670	2.1384
log transformation		
sqrt transformation	1.0670	2.1384

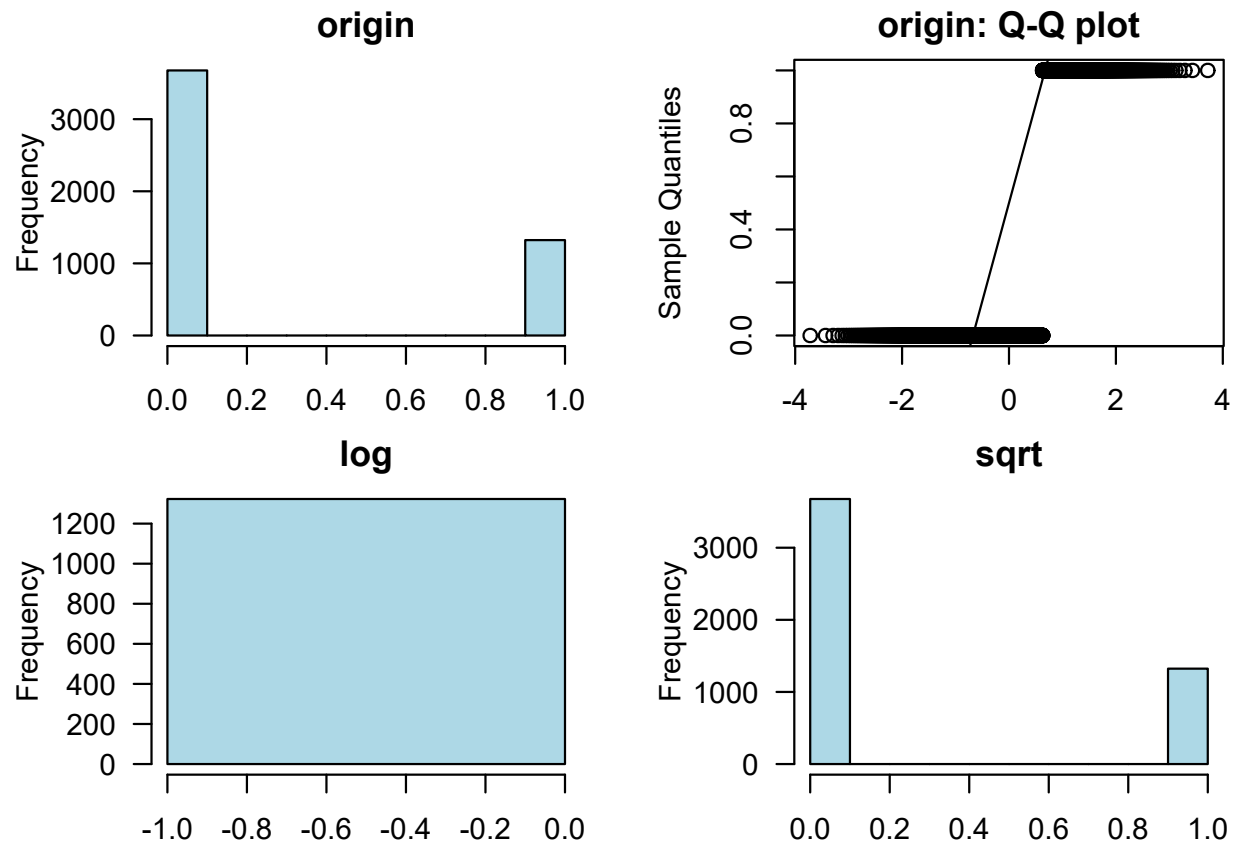


Figure 2.9: not\_smsa

**c\_city**

normality test : Shapiro-Wilk normality test  
 statistic : 0.61259, p-value : 9.91743E-75

type	skewness	kurtosis
original	0.5270	1.2777
log transformation		
sqrt transformation	0.5270	1.2777

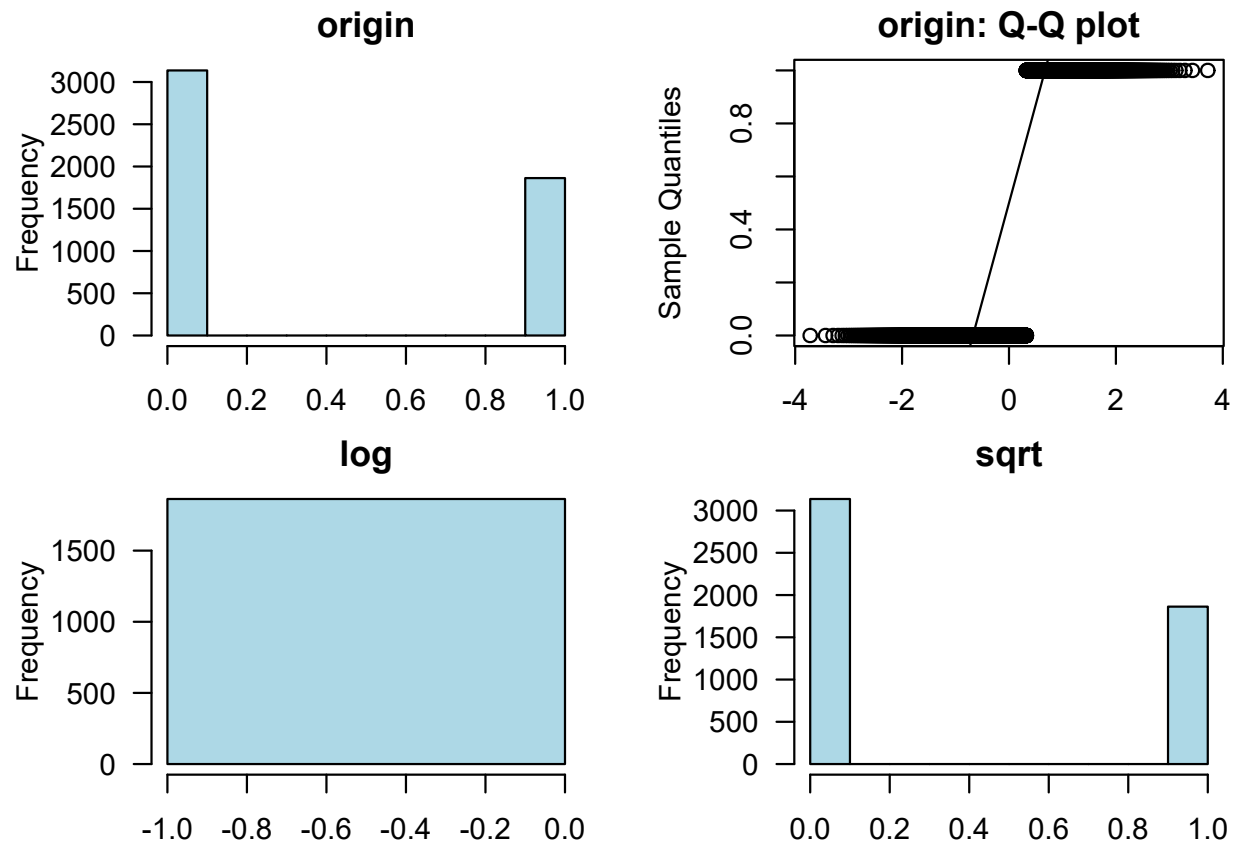


Figure 2.10: c\_city

**south**

normality test : Shapiro-Wilk normality test  
 statistic : 0.62522, p-value : 4.85386E-74

type	skewness	kurtosis
original	0.3584	1.1285
log transformation		
sqrt transformation	0.3584	1.1285

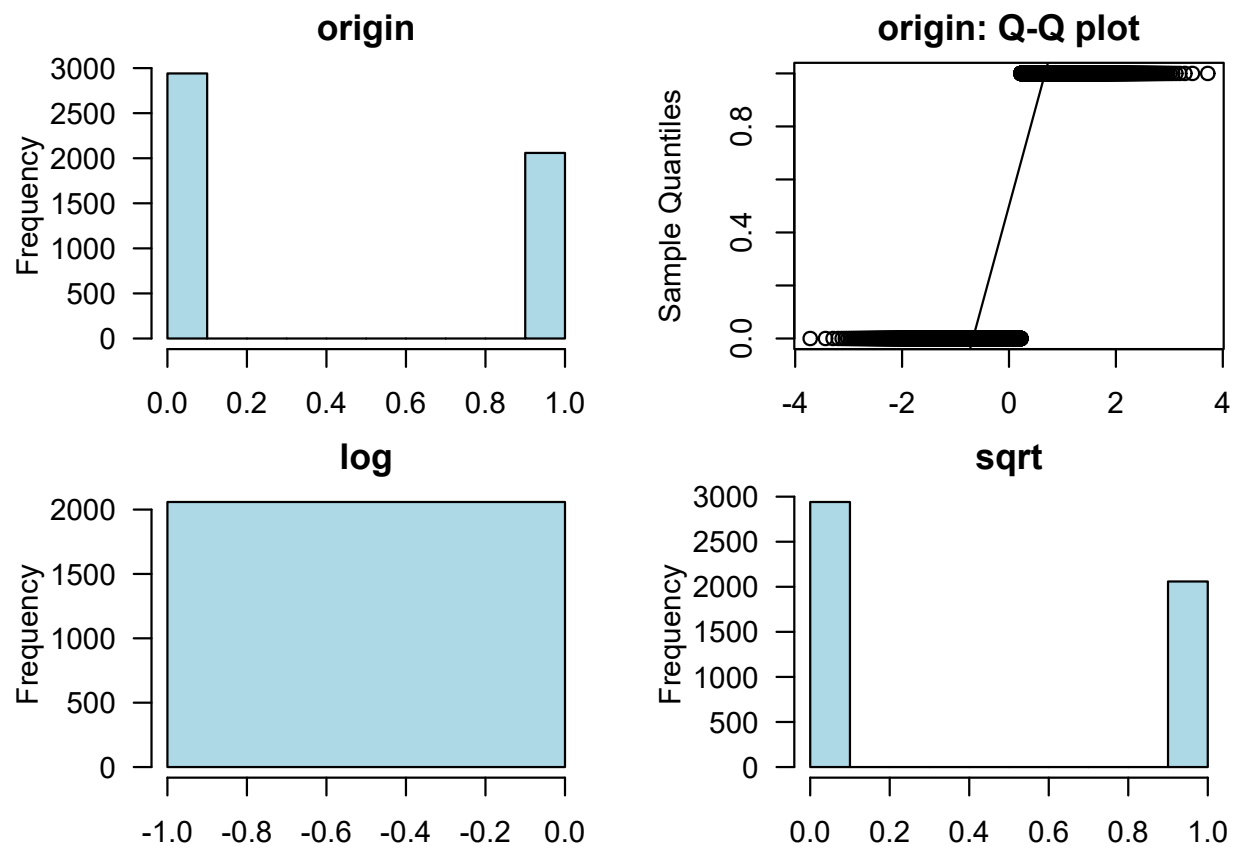


Figure 2.11: south



**ind\_code**

normality test : Shapiro-Wilk normality test  
 statistic : 0.86927, p-value : 1.29627E-53

type	skewness	kurtosis
original	-0.0002	1.5296
log transformation	-0.7853	4.0831
sqrt transformation	-0.2500	1.9867

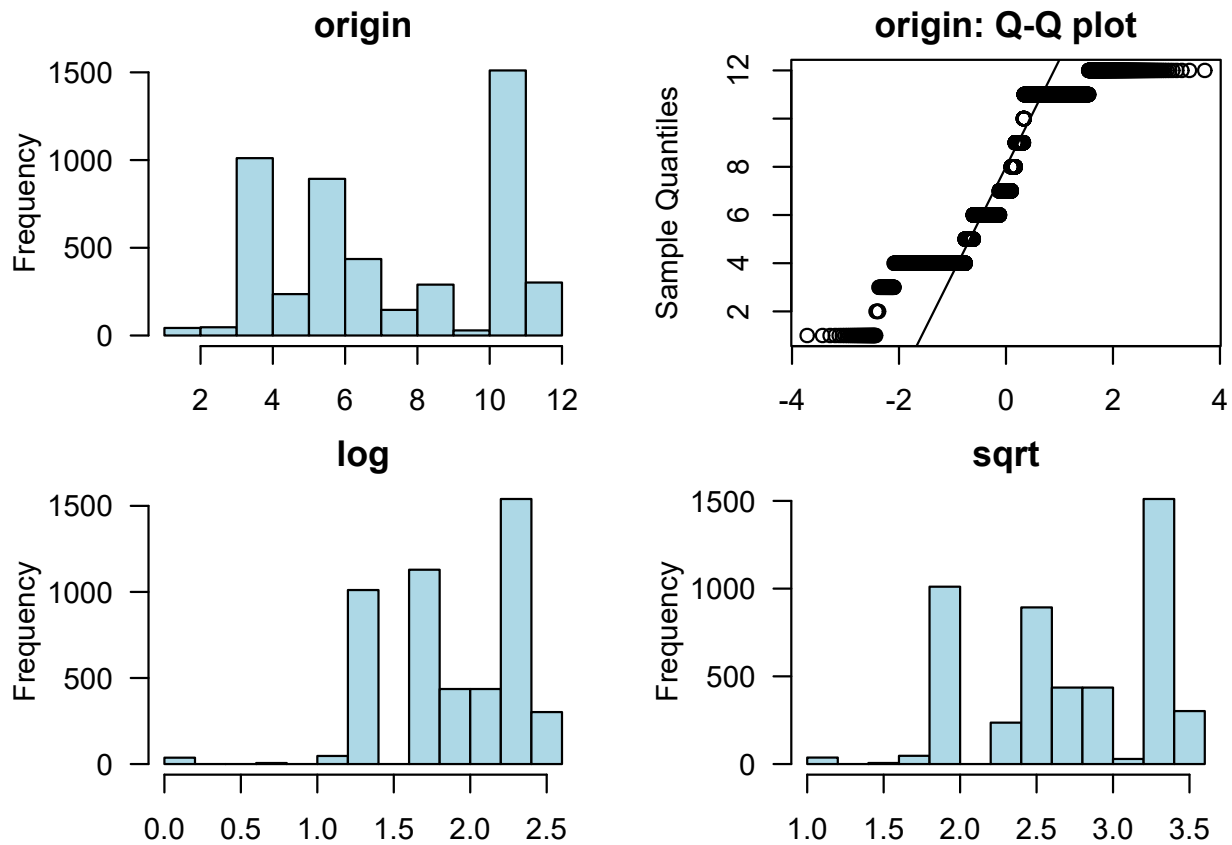


Figure 2.12: ind\_code

**occ\_code**

normality test : Shapiro-Wilk normality test  
 statistic : 0.84861, p-value : 2.27011E-56

type	skewness	kurtosis
original	1.1136	3.7482
log transformation	-0.2817	2.6717
sqrt transformation	0.4720	2.6659

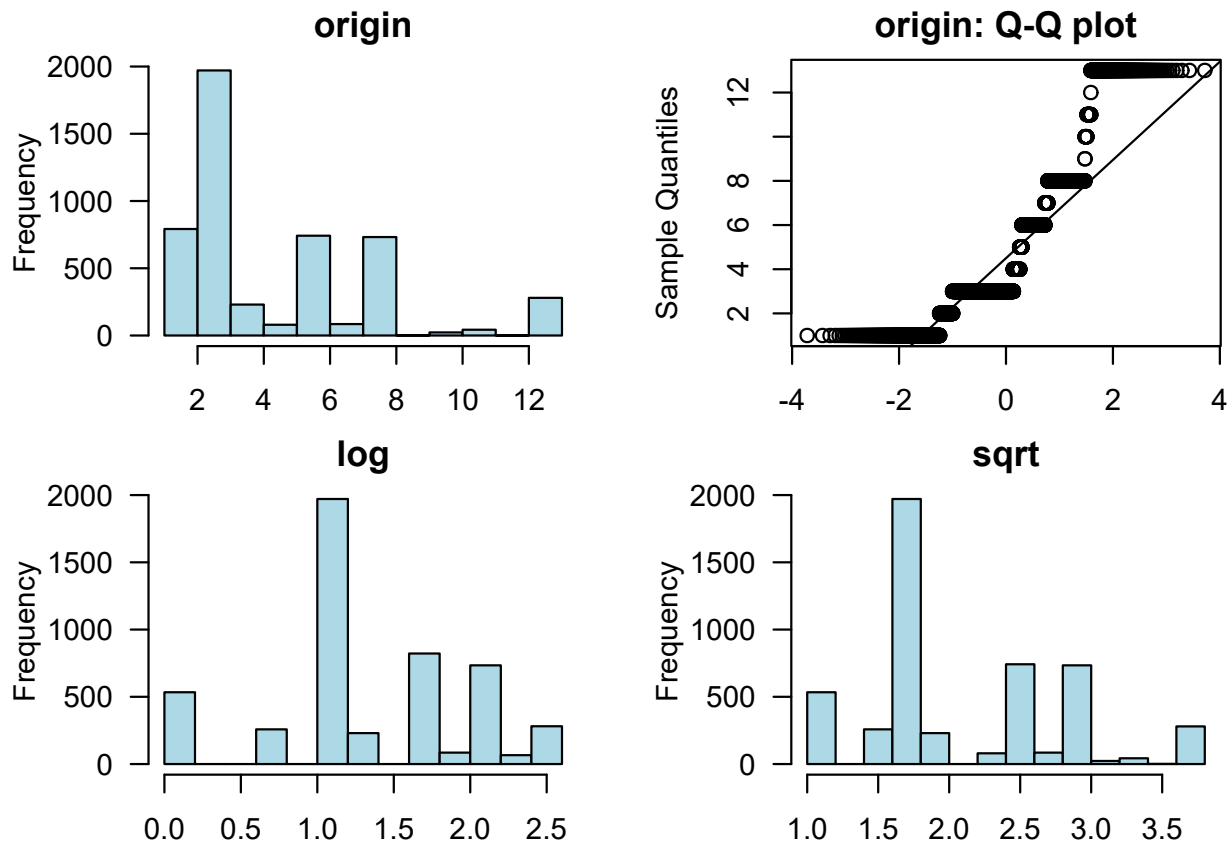


Figure 2.13: occ\_code

**union**

normality test : Shapiro-Wilk normality test  
 statistic : 0.52909, p-value : 6.48545E-70

type	skewness	kurtosis
original	1.2228	2.4952
log transformation		
sqrt transformation	1.2228	2.4952

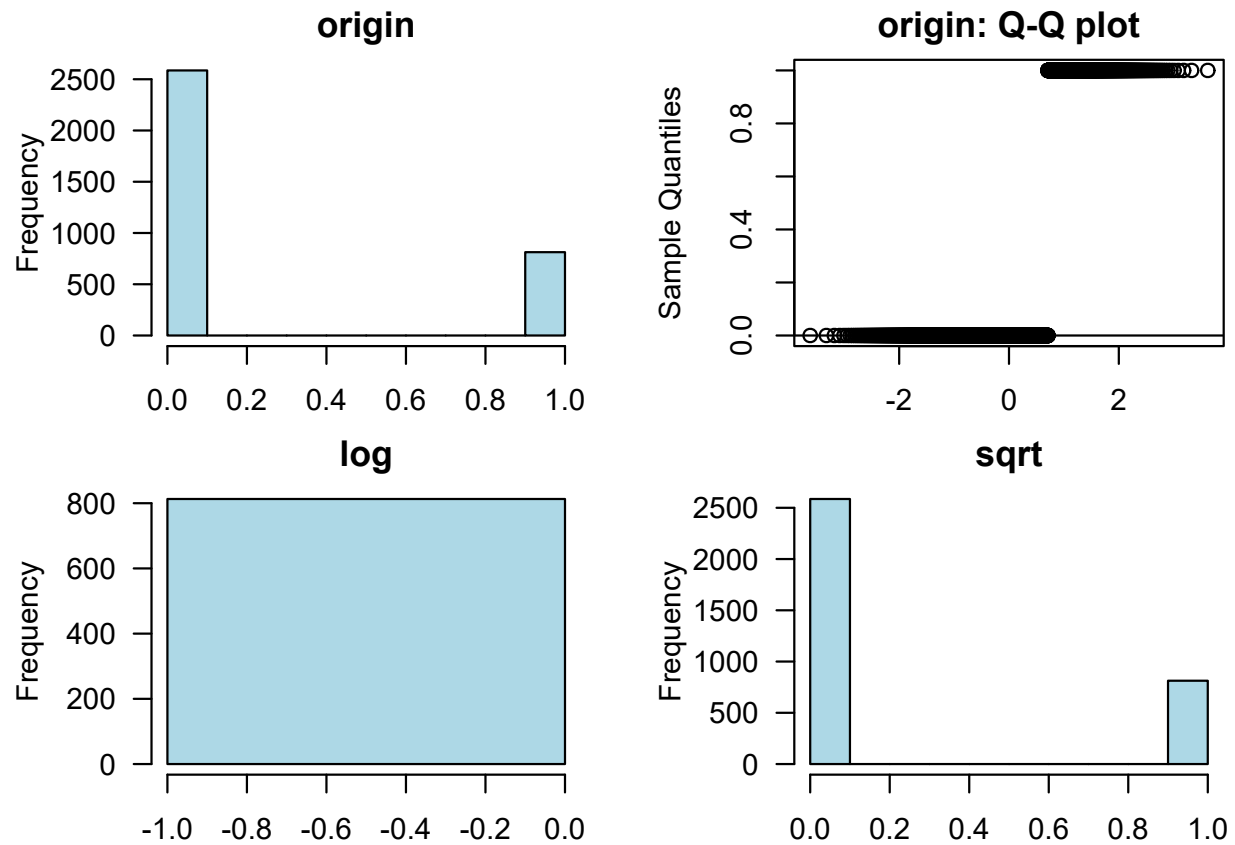


Figure 2.14: union

**wks\_ue**

normality test : Shapiro-Wilk normality test  
 statistic : 0.39506, p-value : 5.47147E-79

type	skewness	kurtosis
original	4.0307	20.8830
log transformation		
sqrt transformation	2.3900	8.2739

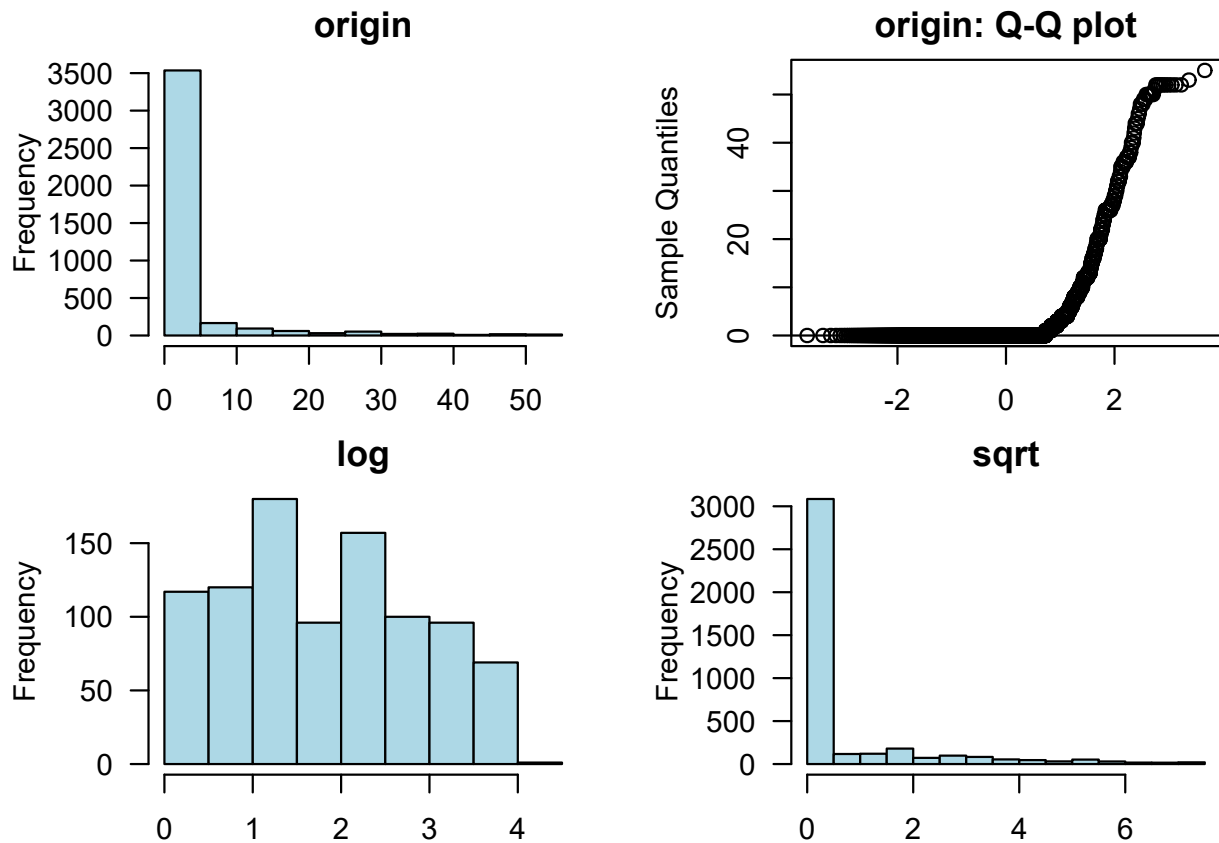


Figure 2.15: wks\_ue

**ttl\_exp**

normality test : Shapiro-Wilk normality test  
 statistic : 0.92512, p-value : 1.82139E-44

type	skewness	kurtosis
original	0.8395	2.9567
log transformation		
sqrt transformation	0.1223	2.2531

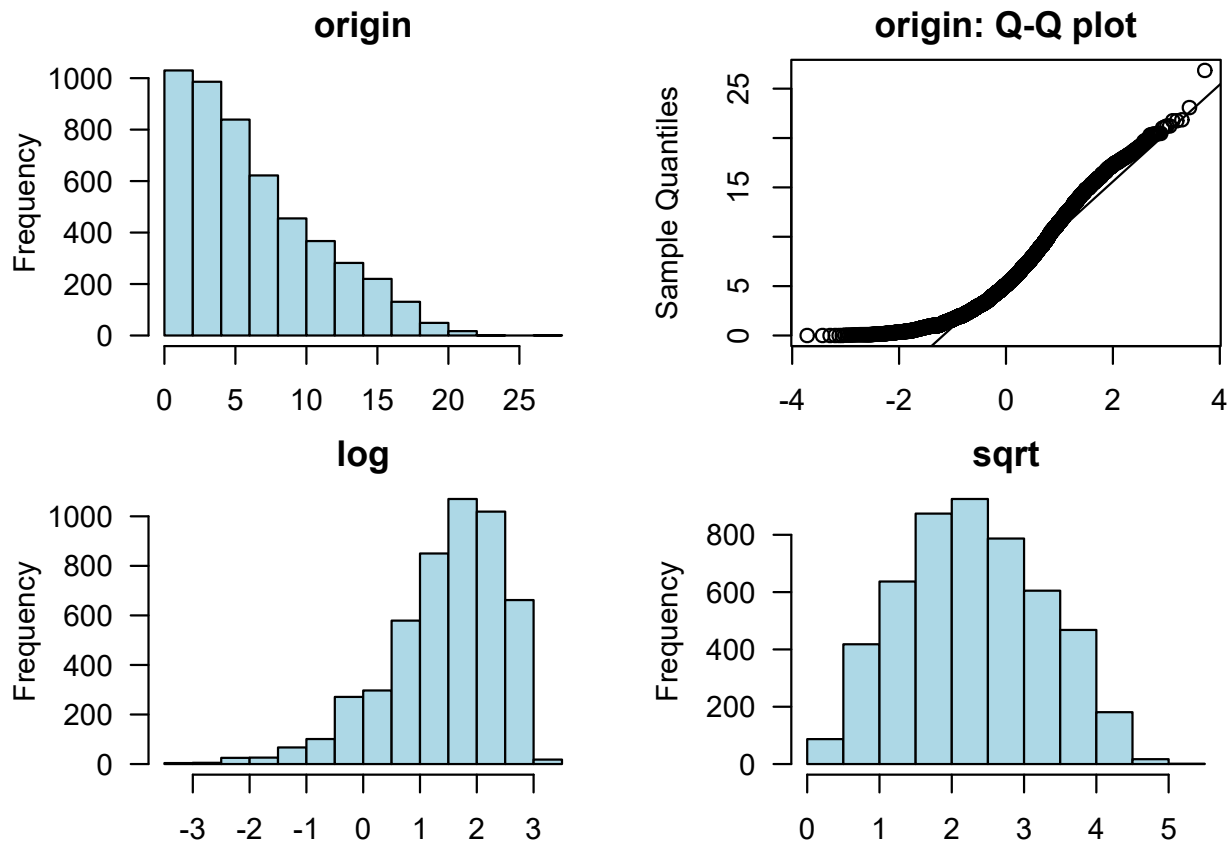


Figure 2.16: ttl\_exp

**tenure**

normality test : Shapiro-Wilk normality test  
 statistic : 0.76745, p-value : 3.78621E-64

type	skewness	kurtosis
original	1.9440	6.9703
log transformation		
sqrt transformation	0.7614	3.0796

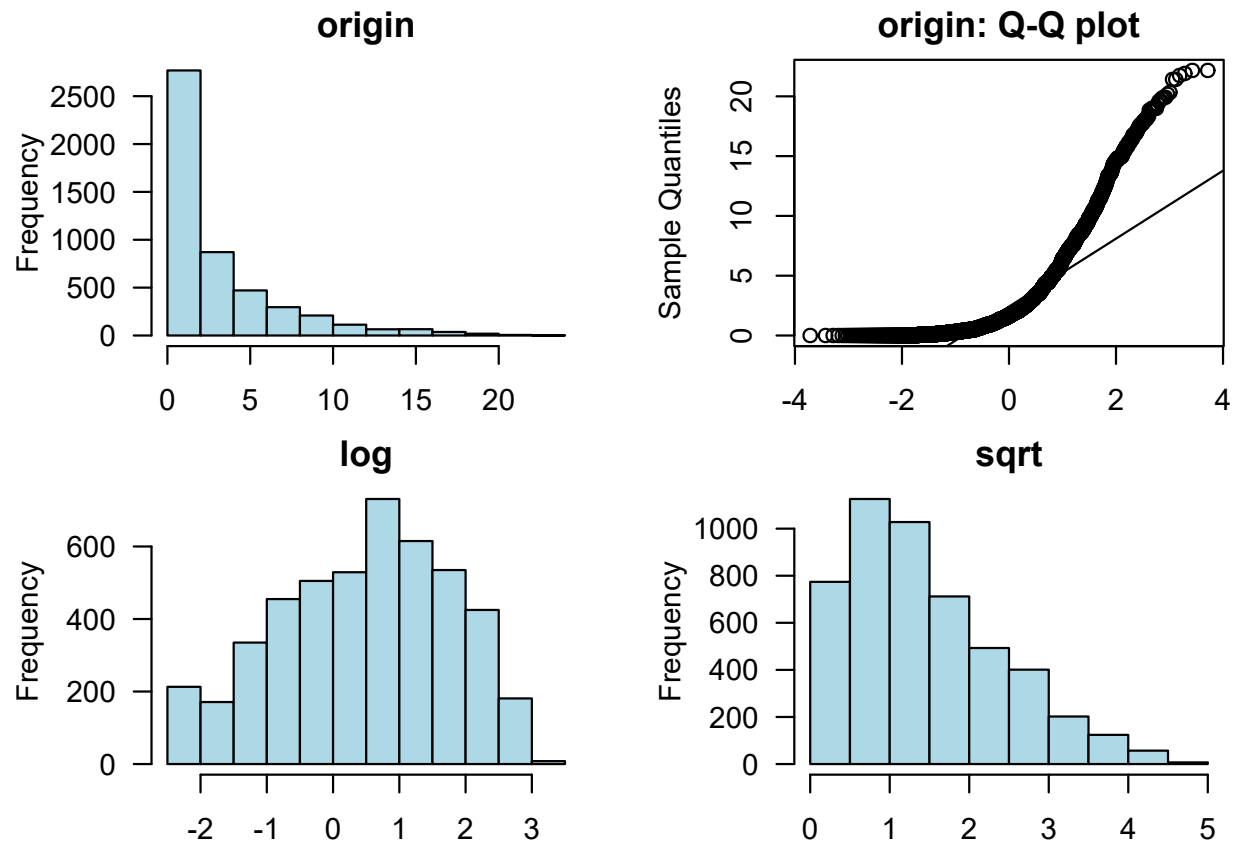


Figure 2.17: tenure

**hours**

normality test : Shapiro-Wilk normality test  
 statistic : 0.76988, p-value : 3.13373E-64

type	skewness	kurtosis
original	-0.6066	11.7546
log transformation	-3.2473	17.1790
sqrt transformation	-1.8946	8.4995

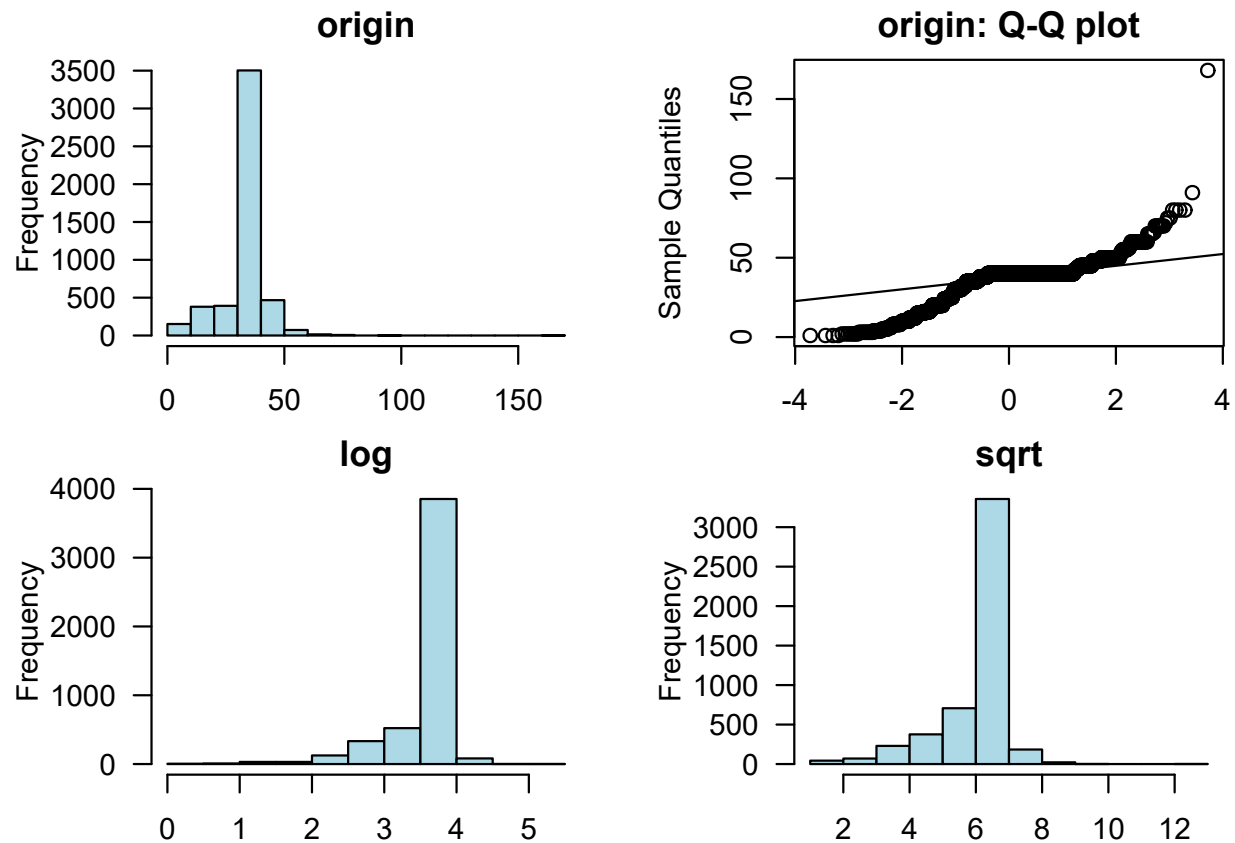


Figure 2.18: hours

**wks\_work**

normality test : Shapiro-Wilk normality test  
 statistic : 0.93814, p-value : 5.32086E-41

type	skewness	kurtosis
original	0.2038	2.3534
log transformation		
sqrt transformation	-0.7862	3.6660

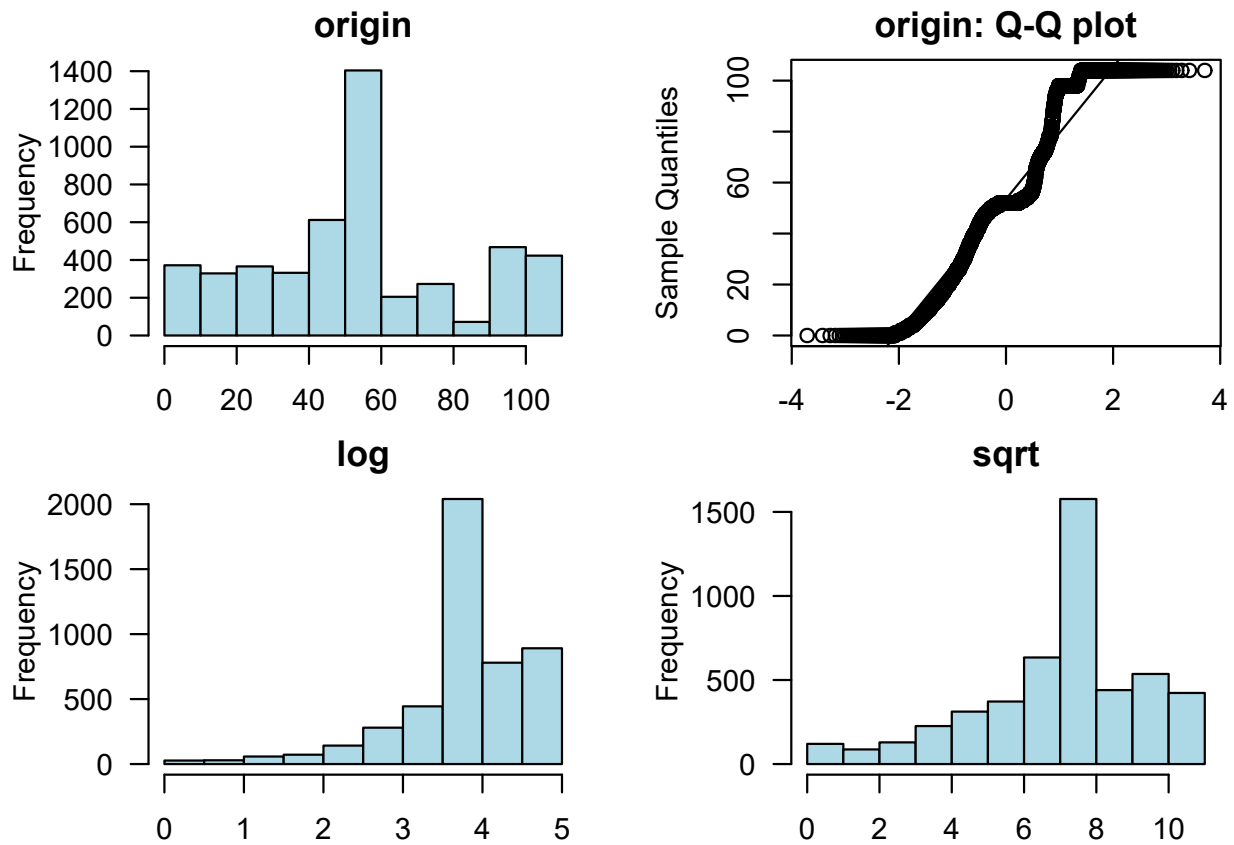


Figure 2.19: wks\_work



**ln\_wage**

normality test : Shapiro-Wilk normality test  
 statistic : 0.97722, p-value : 1.42972E-27

type	skewness	kurtosis
original	0.3756	5.1488
log transformation		
sqrt transformation	-0.7740	6.9490

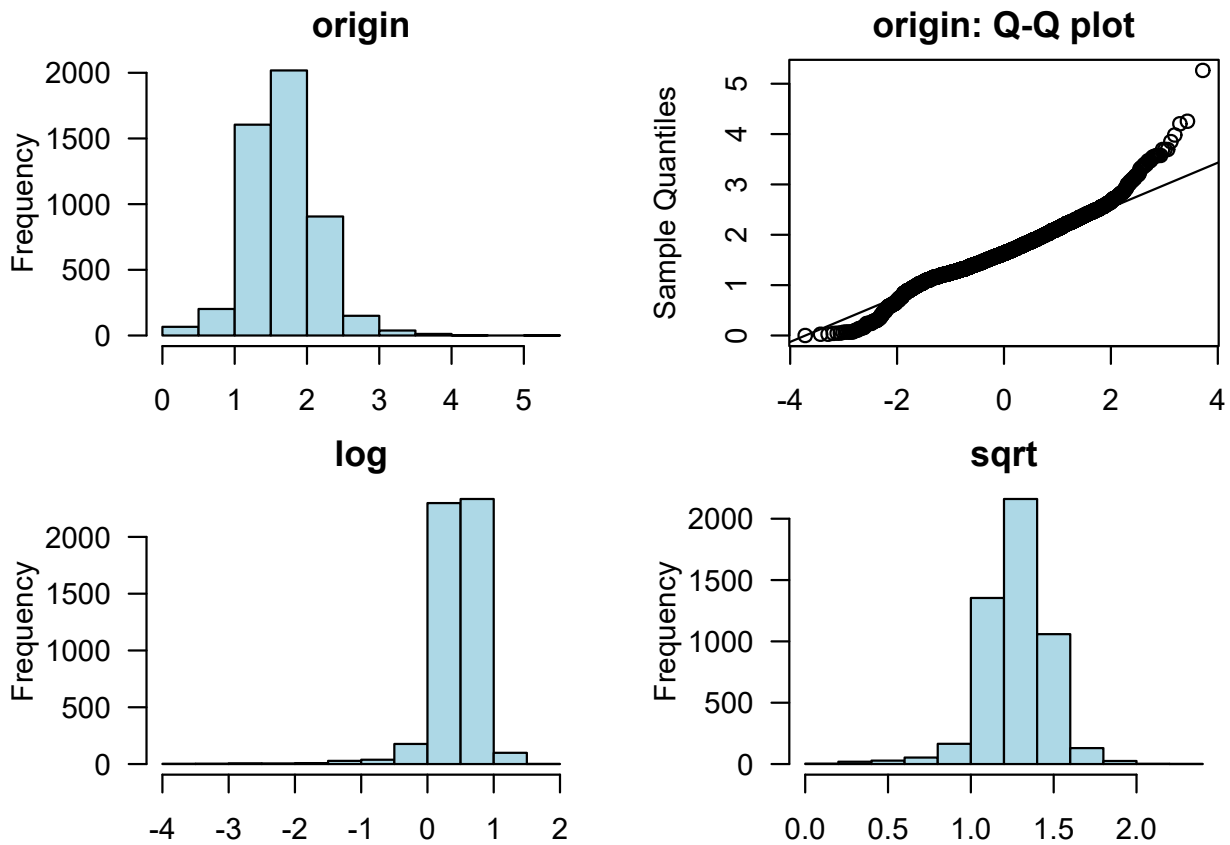


Figure 2.20: ln\_wage



## 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

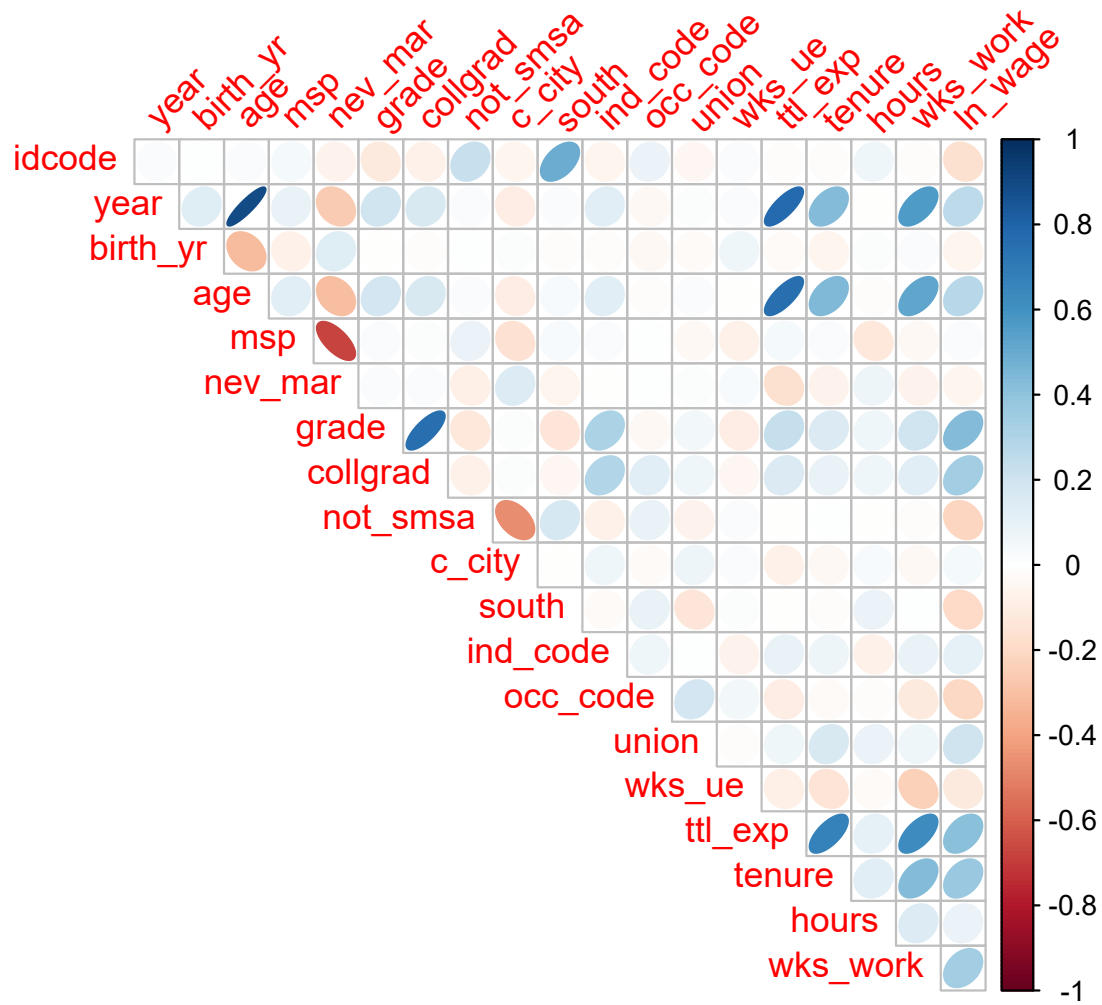


Figure 3.1: The correlation coefficient of numerical variables

# 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.