# Regression Assumptions, Diagnostics, and Fraud Detection

# Regression Assumptions, Diagnostics, and Fraud Detection, Part I

There are several assumptions of linear regression.

1. The relationships are linear

2. The X variables (explanatory variables) are not correlated

3. Distribution of residuals

   a. The error terms have constant variance.

   b. The errors terms are not correlated.

   c. There are no outliers.

# Regression Assumptions, Diagnostics, and Fraud Detection, Part II

## Violations

- Nonlinearity
- Multicollinearity
- Heteroscedasticity
- Serial correlation
- Outliers

## Test

- RESET
- Variance inflation factor (VIF)
- Breusch-Pagan
- Durbin-Watson
- Bonferroni

# Regression Assumptions, Diagnostics, and Fraud Detection, Part III

## Violations

- Nonlinearity
- Multicollinearity
- Heteroscedasticity
- Serial correlation
- Outliers

## Correction

- Box-Cox and Box-Tidwell
- Factor analysis
- Huber regression
- Prais-Winsten
- Drop outlier

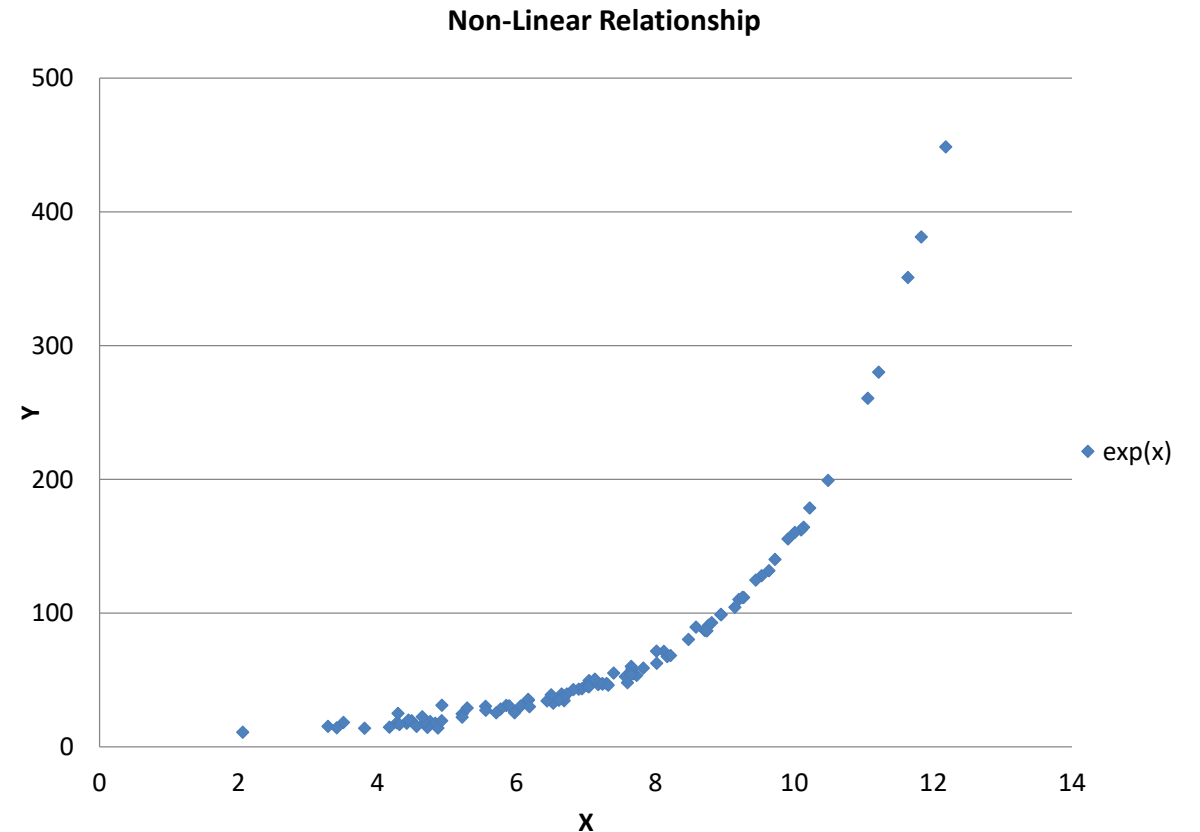Regression Assumptions, Diagnostics, and Fraud Detection

# The End

# R

Regression Linearity Test

# R: Regression Linearity Test

- Problem: data is nonlinear
- Test: RESET
- Correction: Box-Cox and Box-Tidwell



**Non-Linear Relationship**
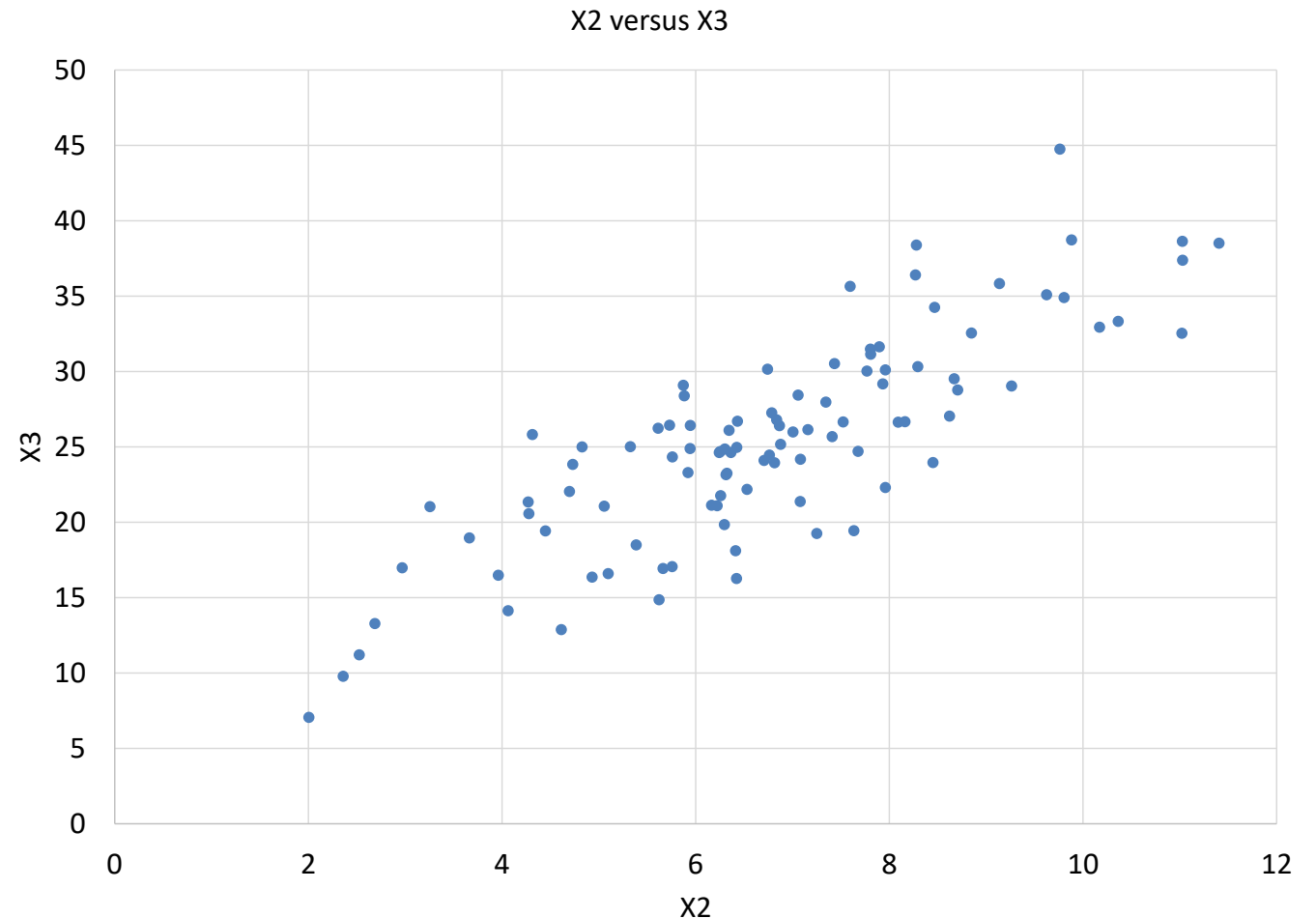
R: Regression Linearity Test

# The End

# R

Collinearity Test

# R: Collinearity Test

- Problem: X-variables are correlated
- Test: variance inflation factor (VIF)
- Correction: factor analysis



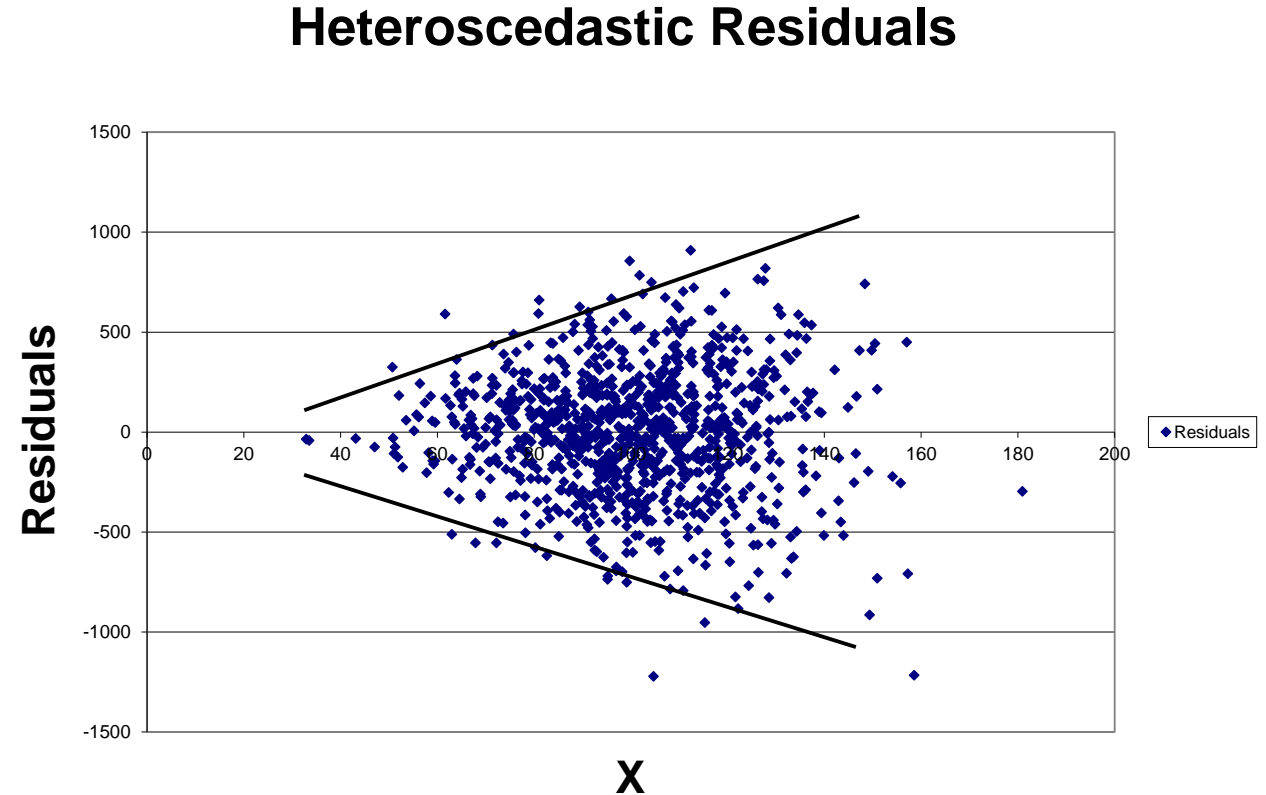X2 versus X3

R: Collinearity Test

# The End

# R

Heteroscedasticity Test

# R: Heteroscedasticity

- Problem: magnitude of residuals changes as X changes (nonconstant variance)
- Test: Breusch-Pagan
- Correction: Huber regression

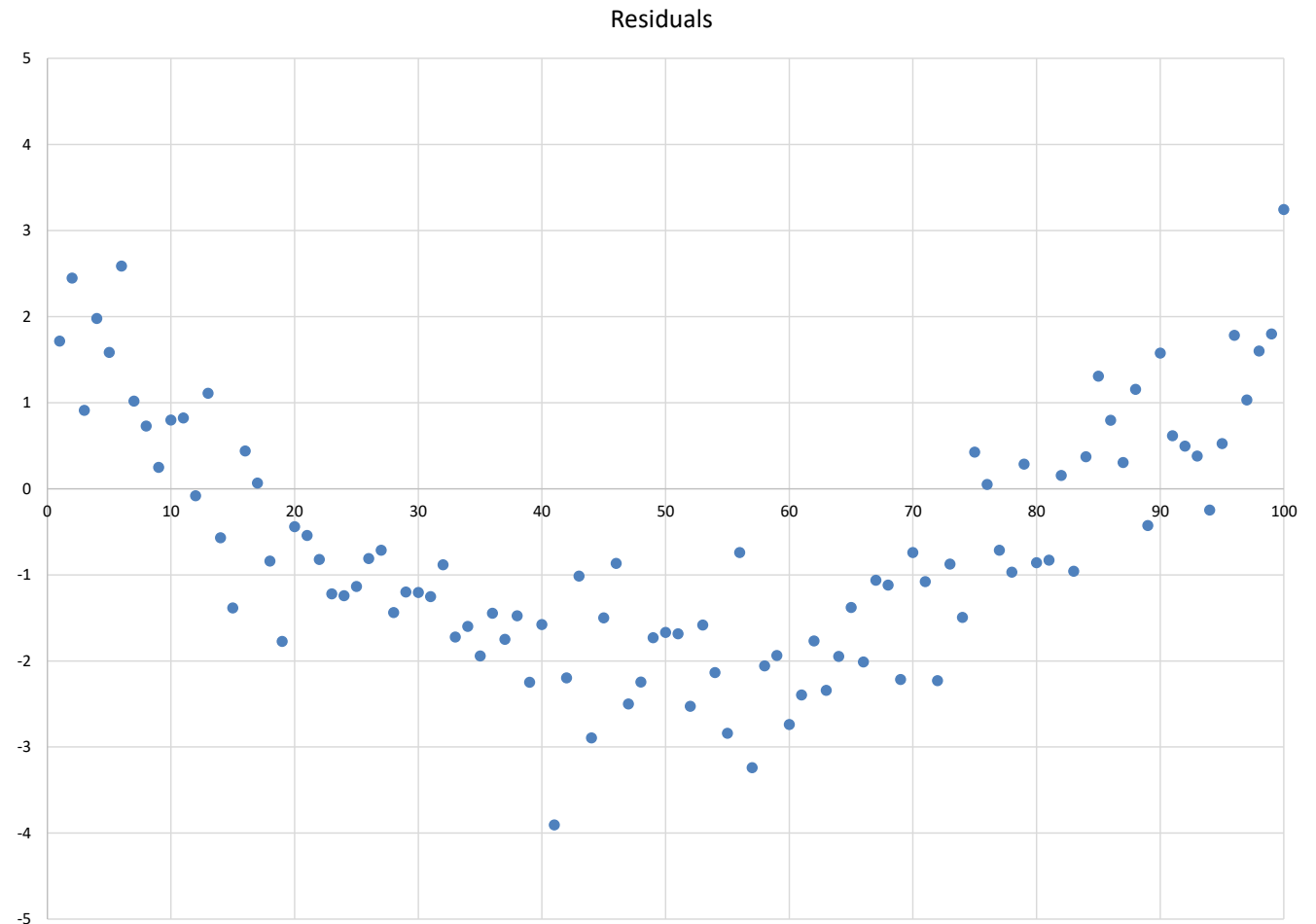**Heteroscedastic Residuals**

R: Heteroscedasticity

# The End

# R

Serial Correlation Test (Durbin-Watson)

# R: Serial Correlation

- Problem: error terms are correlated
- Test: Durbin-Watson
- Correction: Prais-Winsten, Cochrane-Orcutt, ARCH, rho-differencing
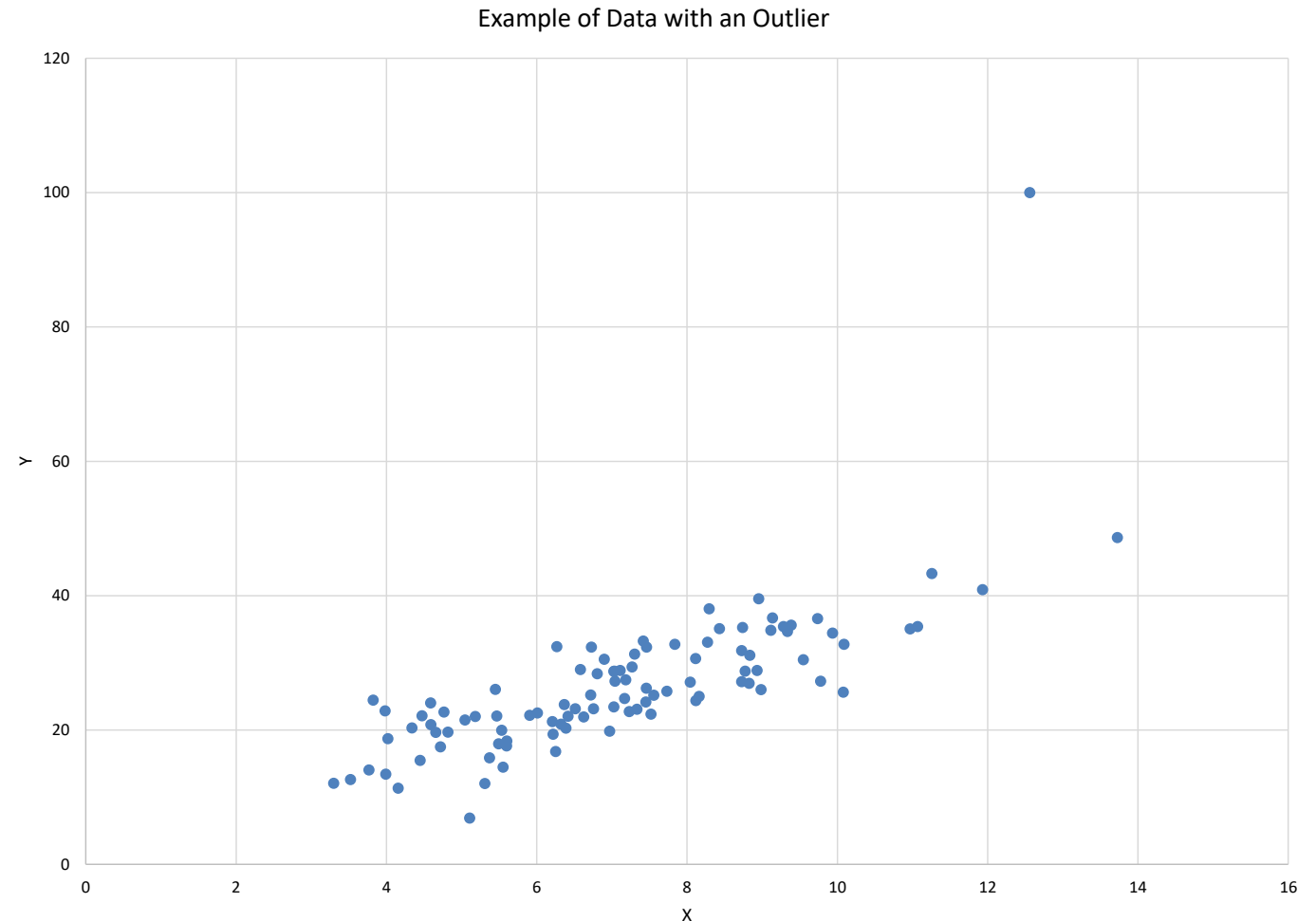


Residuals

R: Serial Correlation

# The End

# R

Outlier Test

# R: Outliers

- Problem: outliers are data points far from the line, artificially influencing the slope
- Test: Bonferroni
- Correction: drop outliers

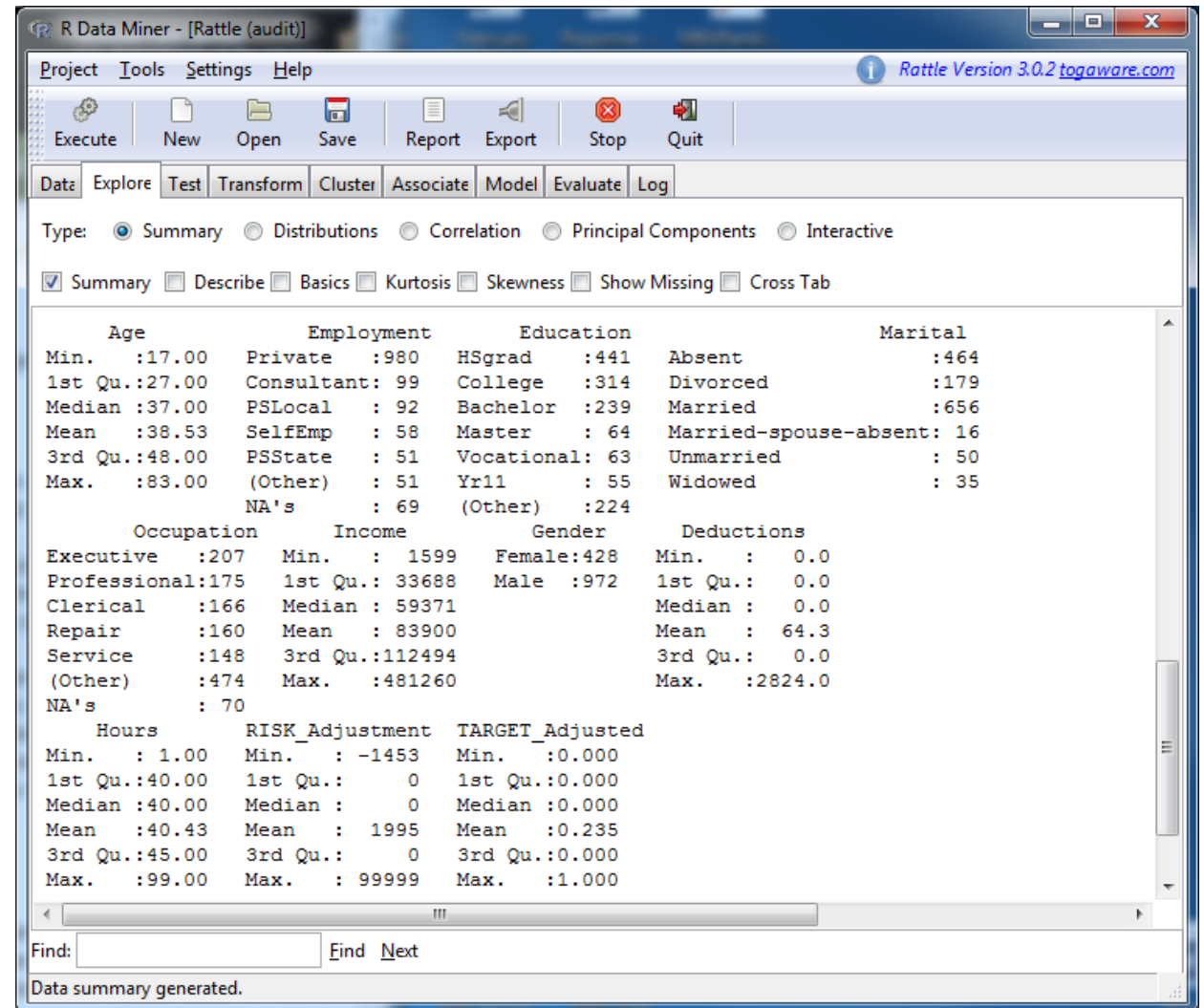Example of Data with an Outlier

# The End

# R

Data Mining Using Rattle

# R: Data Mining Using Rattle

Rattle is a data mining package in R which can create summary statistics including:

• Quartiles across multiple variables simultaneously

R: Data Mining Using Rattle
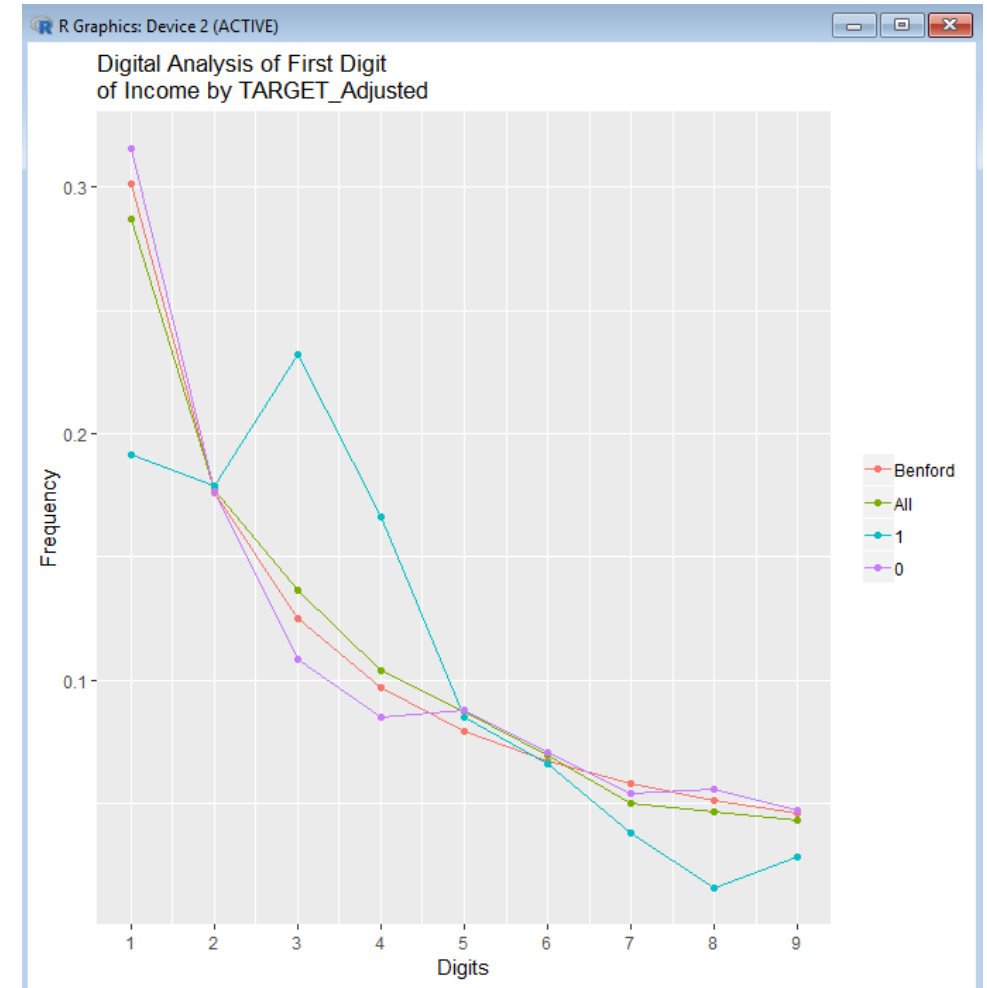
# The End

# R

Benford's Law

# R: Benford's Law

- Benford's law states that financial numbers tend to start with smaller digits.
- Deviations from this distribution are often a signal of fraud.
- Benford's law can detect fraud in expense reports, accounts receivable, accounts payable, financial statements, and income tax returns.

| First digit | Frequency |
|:---:|:---:|
| 1 | 30.1% |
| 2 | 17.6% |
| 3 | 12.5% |
| 4 | 9.7% |
| 5 | 7.9% |
| 6 | 6.7% |
| 7 | 5.8% |
| 8 | 5.1% |
| 9 | 4.6% |

# R: Benford's Law

- Deviations from Benford's distribution can signal fraud.
- The graph on the right shows the first digit of income on tax returns
- Group 1 returns were found to be fraudulent; group 0 were not fraudulent

R: Benford's Law

# The End