



Basic Data Analytics Using R

**Chapter 3 from “Data Science and Big Data Analytics:
Discovering, Analyzing, Visualizing and Presenting Data”**

1st Edition by [EMC Education Services](#)

Ka-Chun Wong, Department of Computer Science City University of Hong Kong
Charles Tappert Seidenberg, School of CSIS, Pace University



Data Analytics Using R

- A. An overview of R
- B. Using R to perform exploratory data analysis tasks using visualization



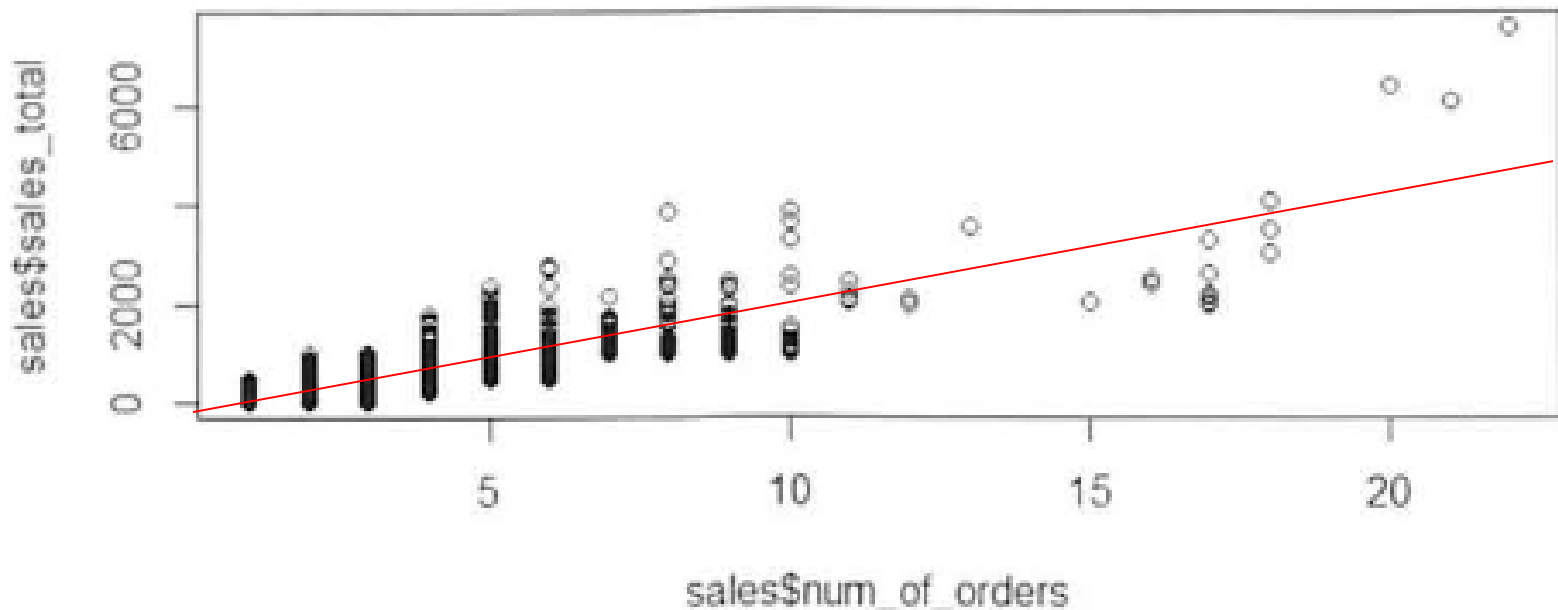
3.1 Introduction to R

- Generic R functions are functions that share the same name but behave differently depending on the type of arguments they receive (polymorphism)
- Some important functions used in chapter (most are generic)
 - `head()` displays first six records of a file
 - `summary()` generates descriptive statistics
 - `plot()` can generate a scatter plot of one variable against another
 - `lm()` applies a linear regression model between two variables
 - `hist()` generates a histogram
 - `help()` provides details of a function

3.1 Introduction to R

Example: number of orders vs sales

```
plot(sales$num_of_orders, sales$sales_total,  
     main="Number of Orders vs. Sales")
```



```
abline(lm(formula = (sales$sales_total ~ sales$num_of_orders)))
```

intercept = -154.1 slope = 166.2



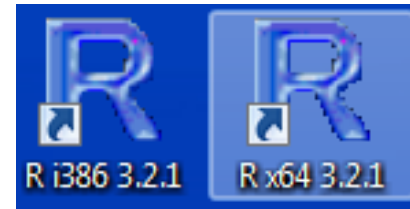
3.1 Introduction to R

- 3.1.1 R Graphical User Interfaces
 - Getting R and RStudio
- 3.1.2 Data Import and Export
 - Necessary for project work
- 3.1.3 Attributes and Data Types
 - Vectors, matrices, data frames
- 3.1.4 Descriptive Statistics
 - `summary()`, `mean()`, `median()`, `sd()`

3.1.1 Getting R and RStudio

- Download R and install (32-bit and 64-bit)

- <https://www.r-project.org/>

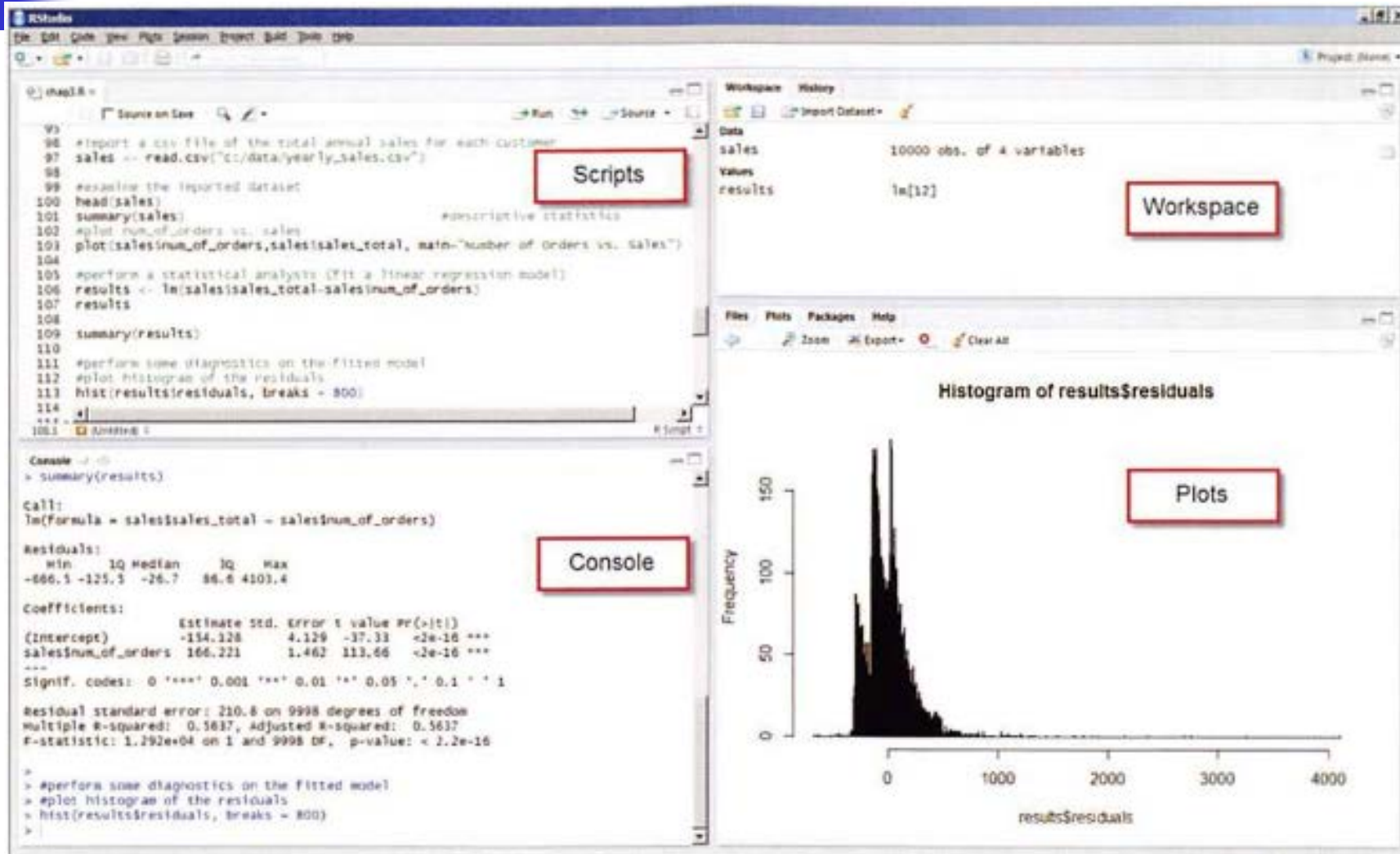


- Download RStudio and install

- <https://www.rstudio.com/>

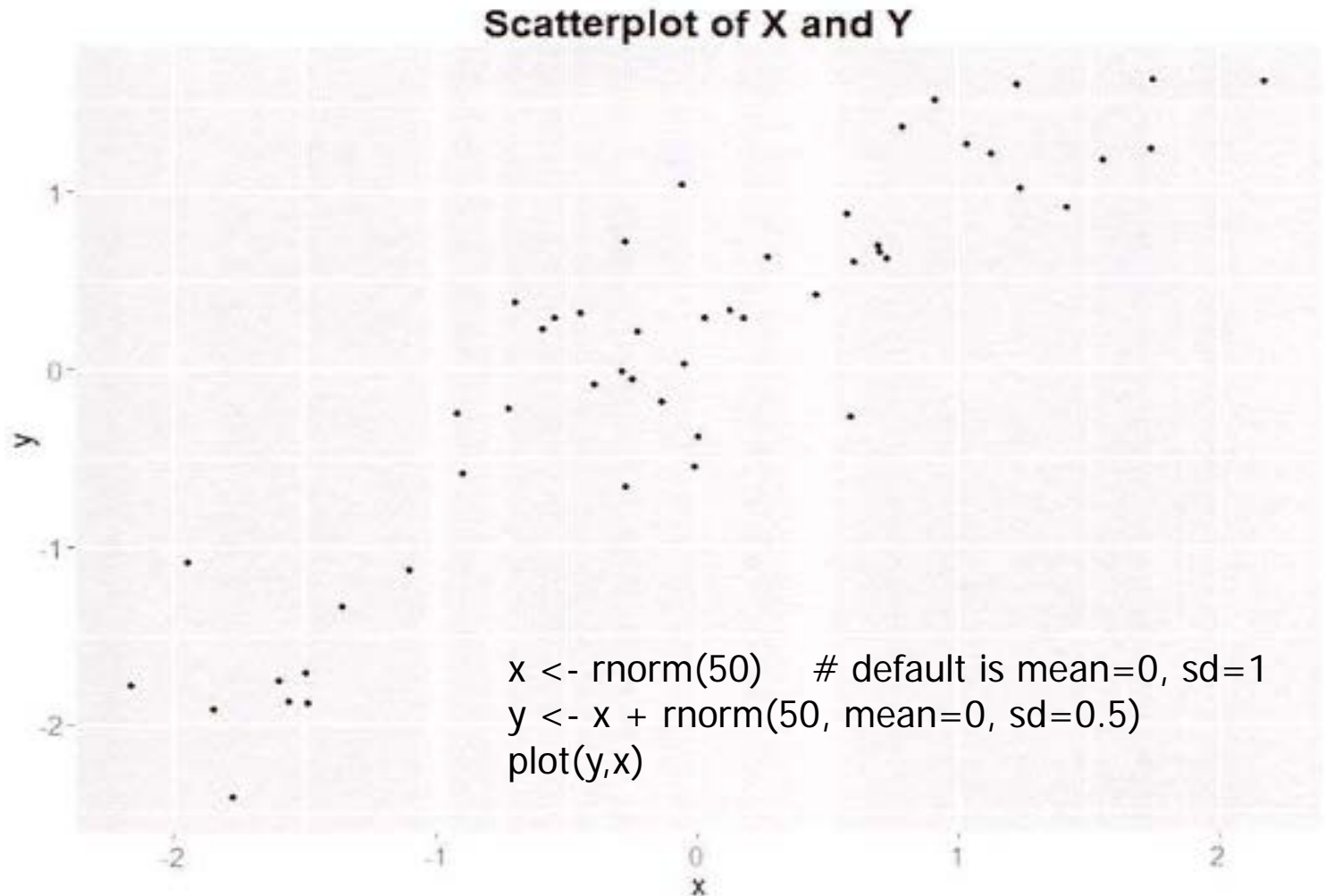


3.1.1 RStudio GUI



3.2 Exploratory Data Analysis

Scatterplots show possible relationships





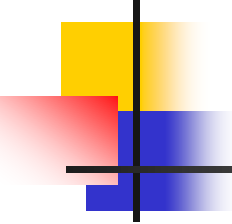
3.2 Exploratory Data Analysis

- 3.2.1 Visualization before Analysis
- 3.2.2 Dirty Data
- 3.2.3 Visualizing a Single Variable
- 3.2.4 Examining Multiple Variables
- 3.2.5 Data Exploration versus Presentation

3.2.1 Visualization before Analysis

Anscombe's quartet

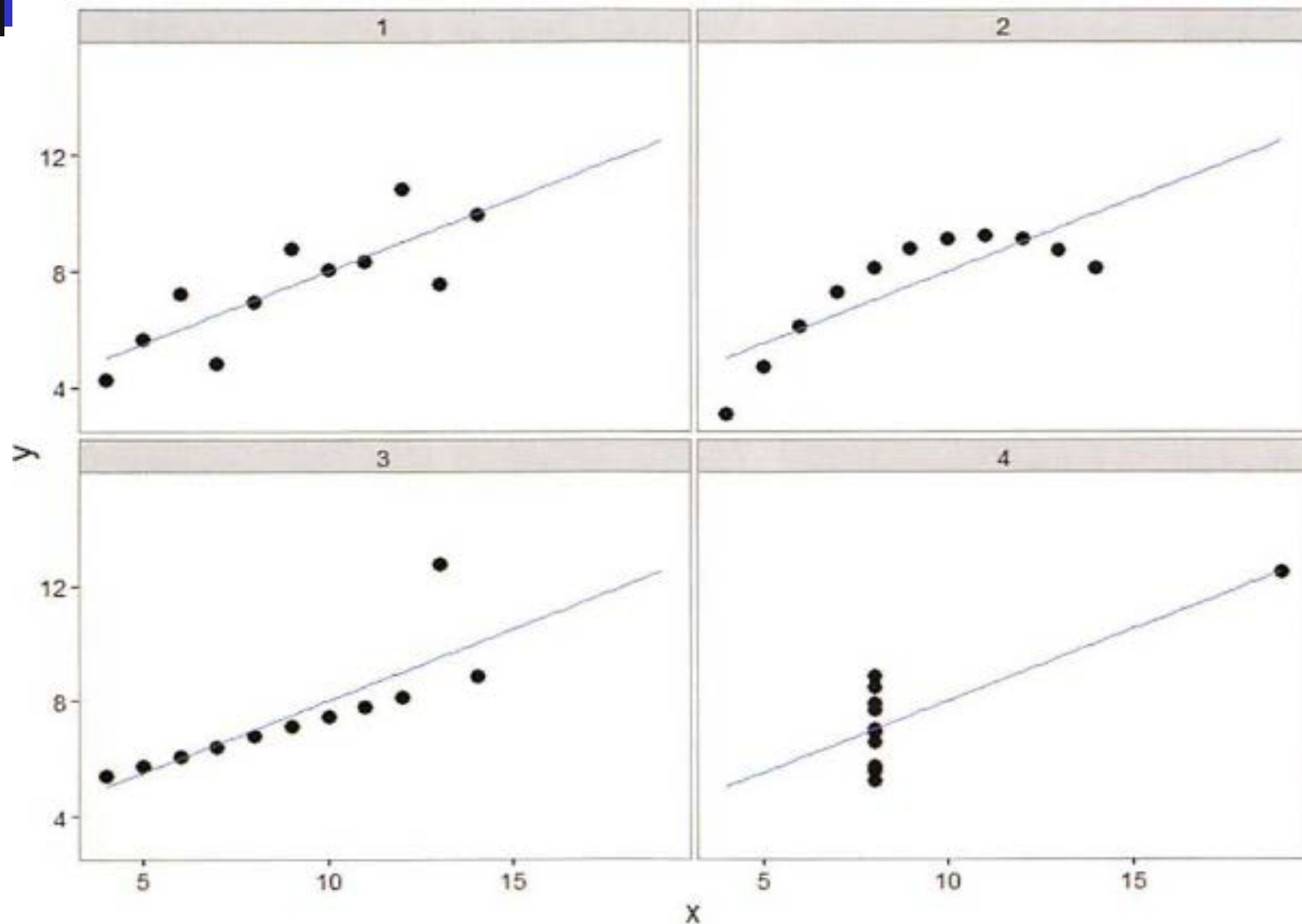
4 datasets with the same statistics



Statistical Property	Value
Mean of x	9
Variance of x	11
Mean of y	7.50 (to 2 decimal points)
Variance of y	4.12 or 4.13 (to 2 decimal points)
Correlations between x and y	0.816
Linear regression line	$y = 3.00 + 0.50x$ (to 2 decimal points)

3.2.1 Visualization before Analysis

Anscombe's quartet – visualized



3.2.1 Visualization before Analysis

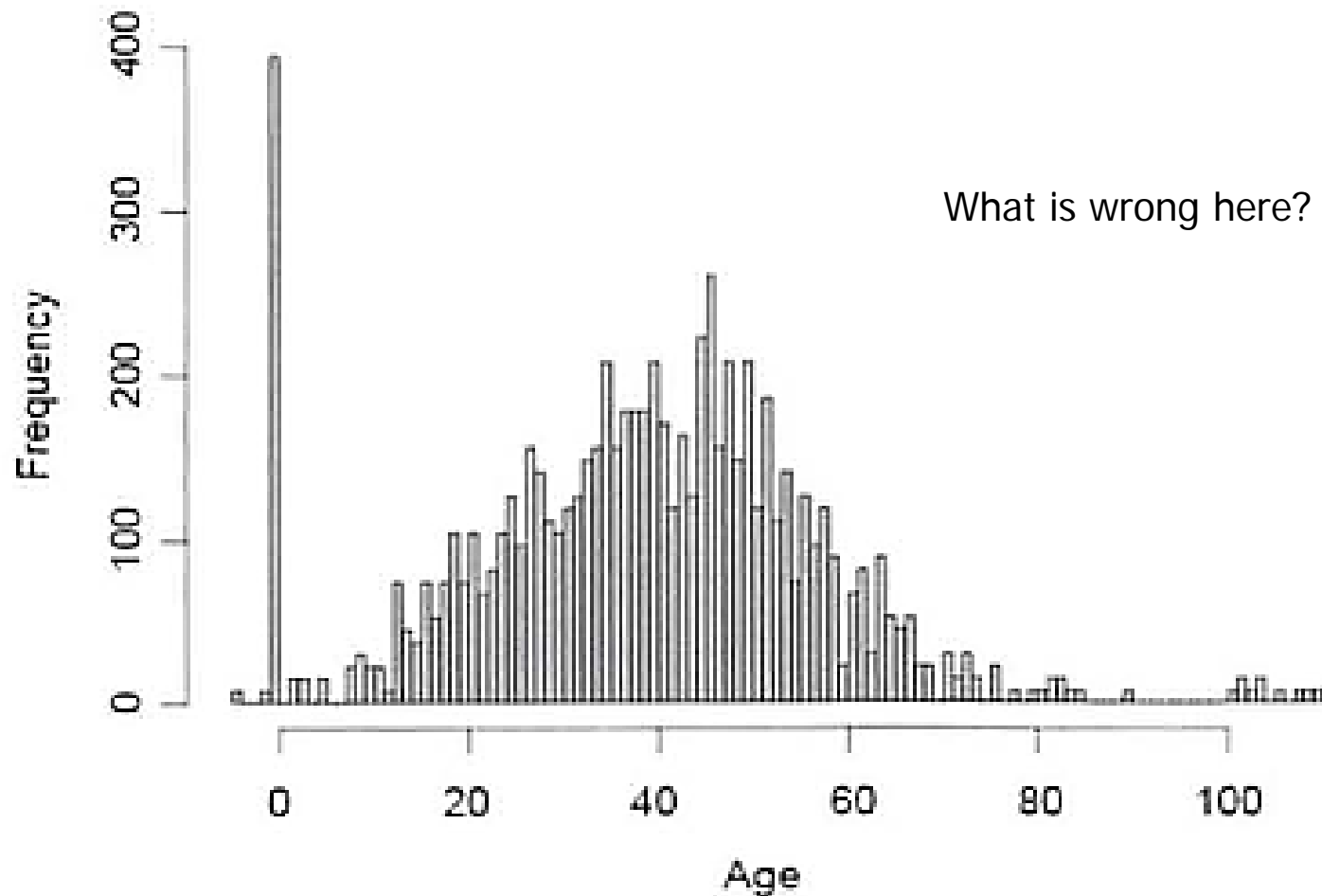
Anscombe's quartet – Rstudio exercise

- Enter and plot Anscombe's dataset #3
- and obtain the linear regression line

```
x <- 4:14
x
y <- c(5.39,5.73,6.08,6.42,6.77,7.11,7.46,7.81,8.15,12.74,8.84)
y
summary(x)
var(x)
summary(y)
var(y)
plot(y~x)
lm(y~x)
```

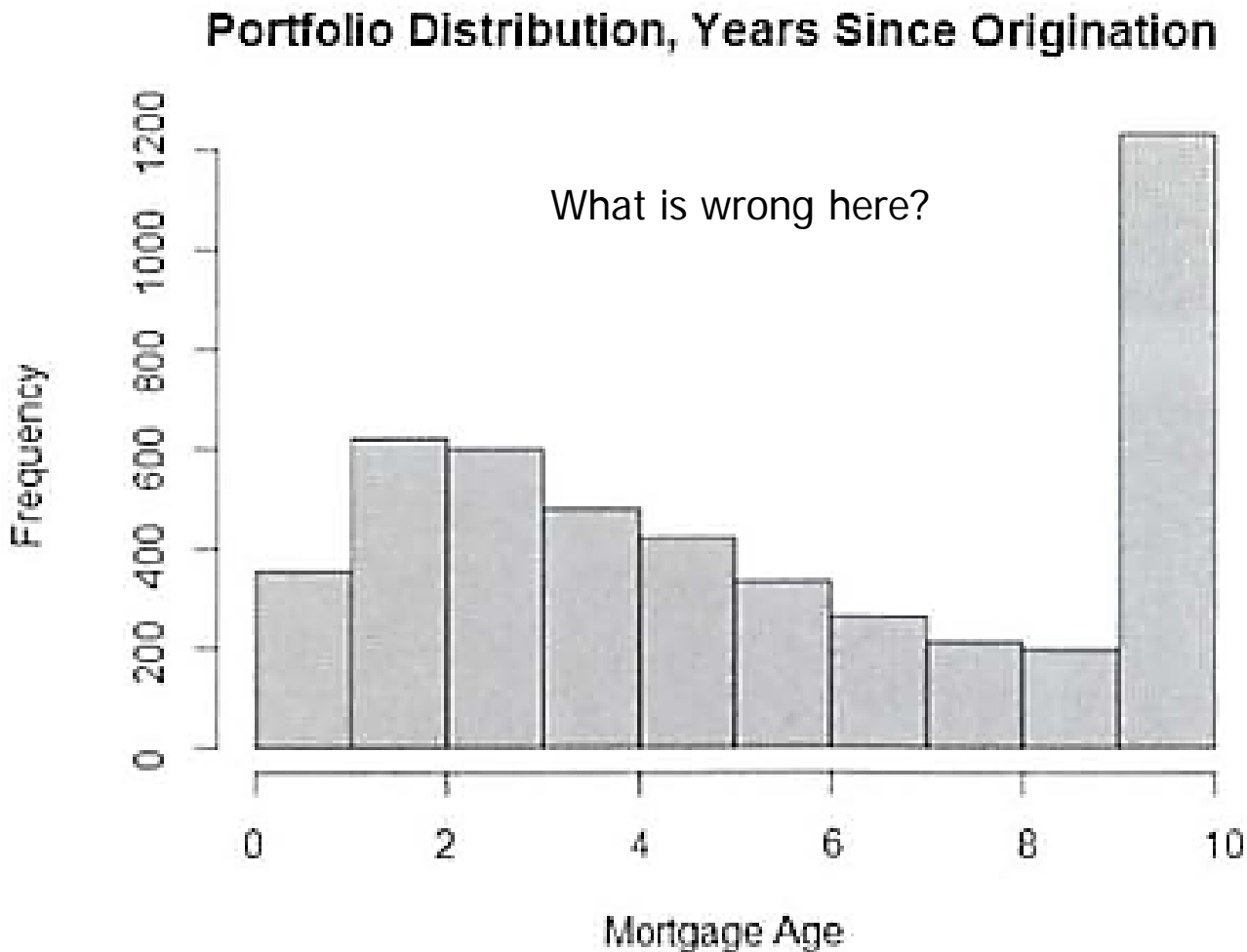
3.2.2 Dirty Data

Age Distribution of bank account holders



3.2.2 Dirty Data

Age of Mortgage



3.2.3 Visualizing a Single Variable

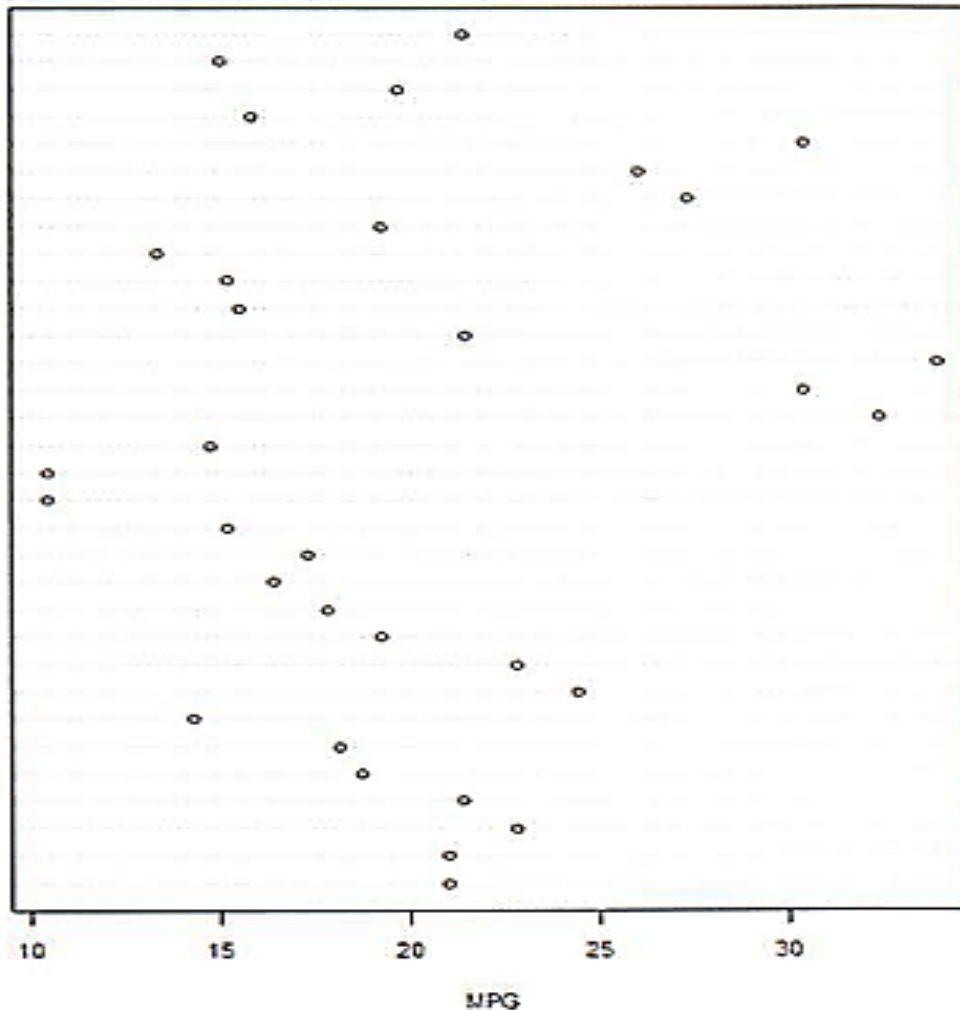
Example Visualization Functions

Function	Purpose
<code>plot (data)</code>	Scatterplot where x is the index and y is the value; suitable for low-volume data
<code>barplot (data)</code>	Barplot with vertical or horizontal bars
<code>dotchart (data)</code>	Cleveland dot plot [12]
<code>hist (data)</code>	Histogram
<code>plot (density (data))</code>	Density plot (a continuous histogram)
<code>stem (data)</code>	Stem-and-leaf plot
<code>rug (data)</code>	Add a rug representation (1-d plot) of the data to an existing plot

3.2.3 Visualizing a Single Variable

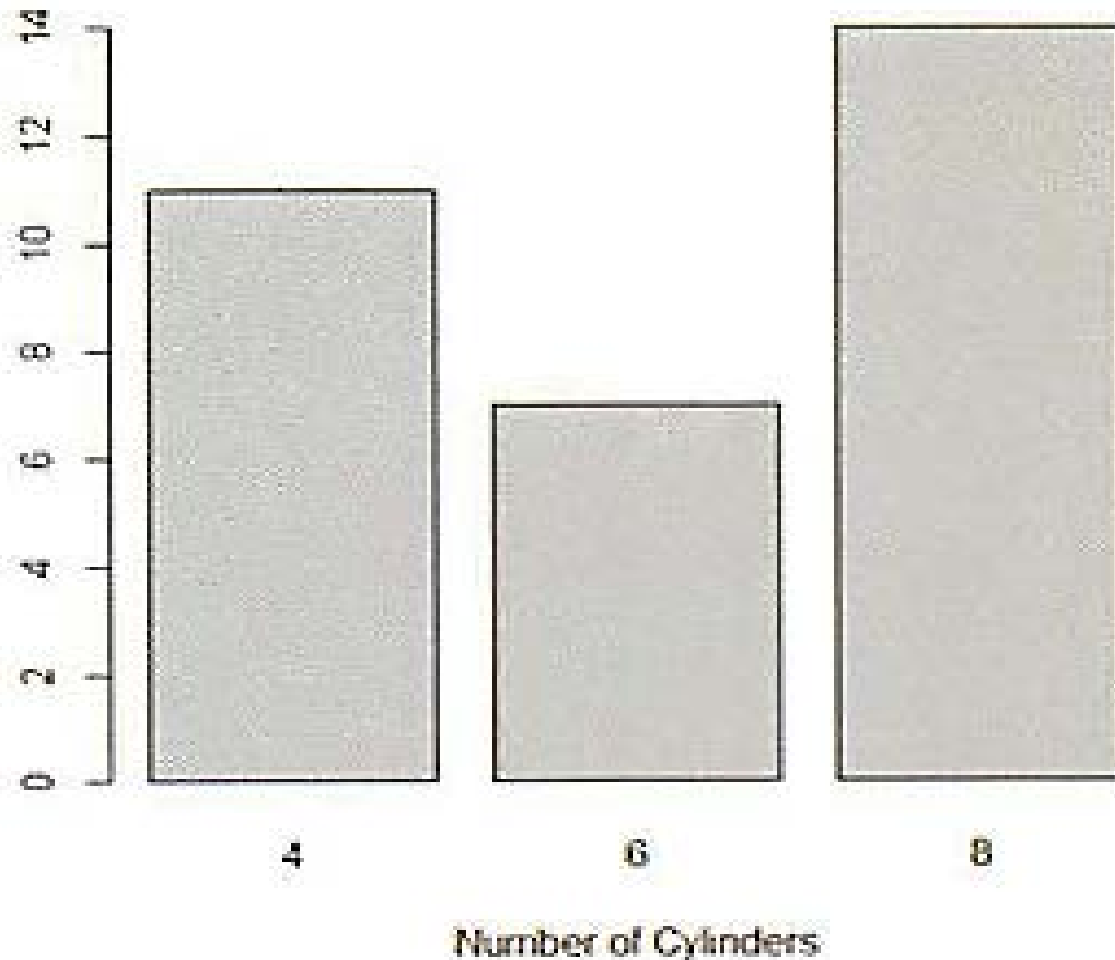
Dotchart – MPG of Car Models

Volvo 142E
Maserati Bora
Ferrari Dino
Ford Pantera L
Lotus Europa
Porsche 914-2
Fiat X1-9
Pontiac Firebird
Camaro Z28
AMC Javelin
Dodge Challenger
Toyota Corona
Toyota Corolla
Honda Civic
Fiat 128
Chrysler Imperial
Lincoln Continental
Cadillac Fleetwood
Merc 450SLC
Merc 450SL
Merc 450SE
Merc 280C
Merc 280
Merc 230
Merc 240D
Duster 360
Valant
Hornet Sportabout
Hornet 4 Drive
Datsun 710
Mazda RX4 Wag
Mazda RX4



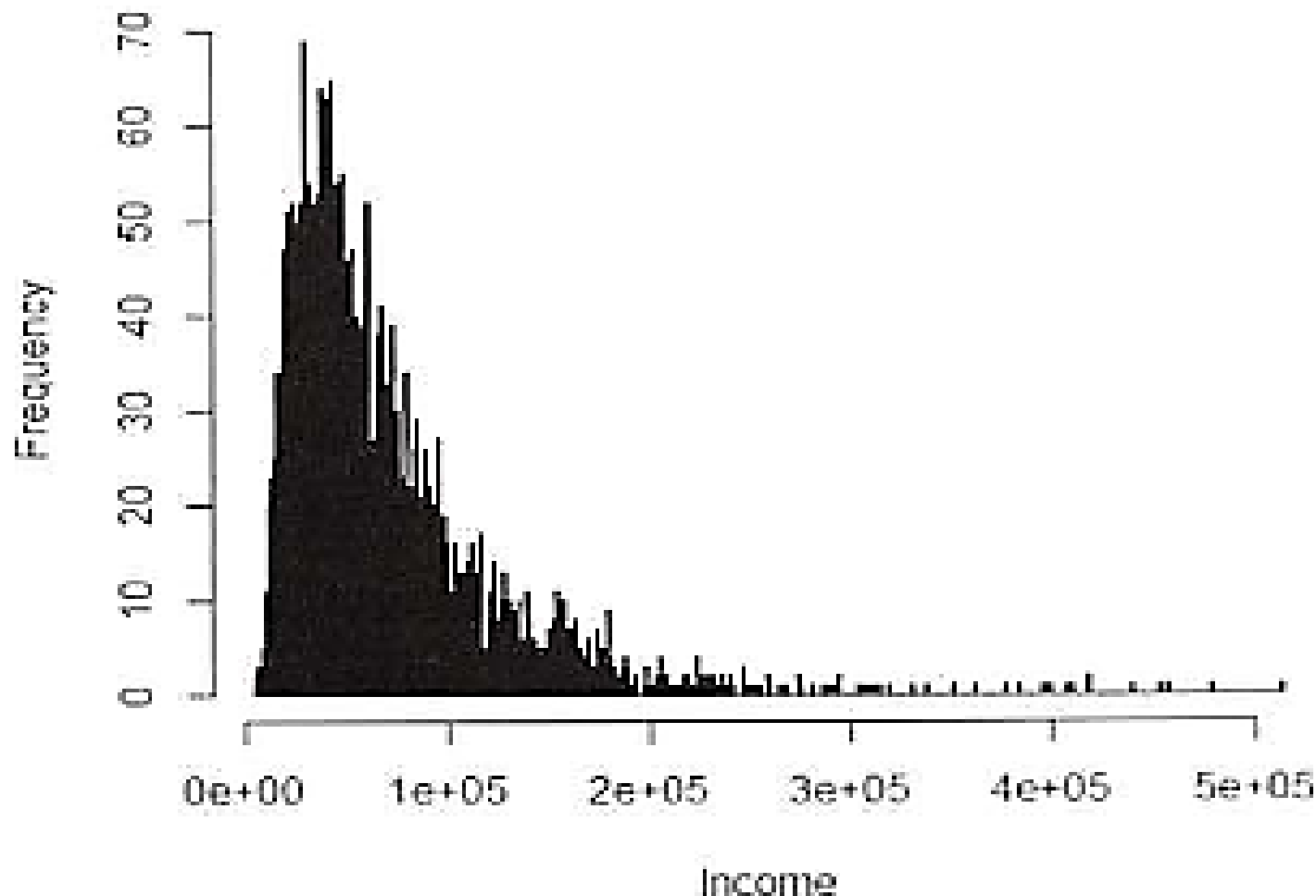
3.2.3 Visualizing a Single Variable

Barplot – Distribution of Car Cylinder Counts



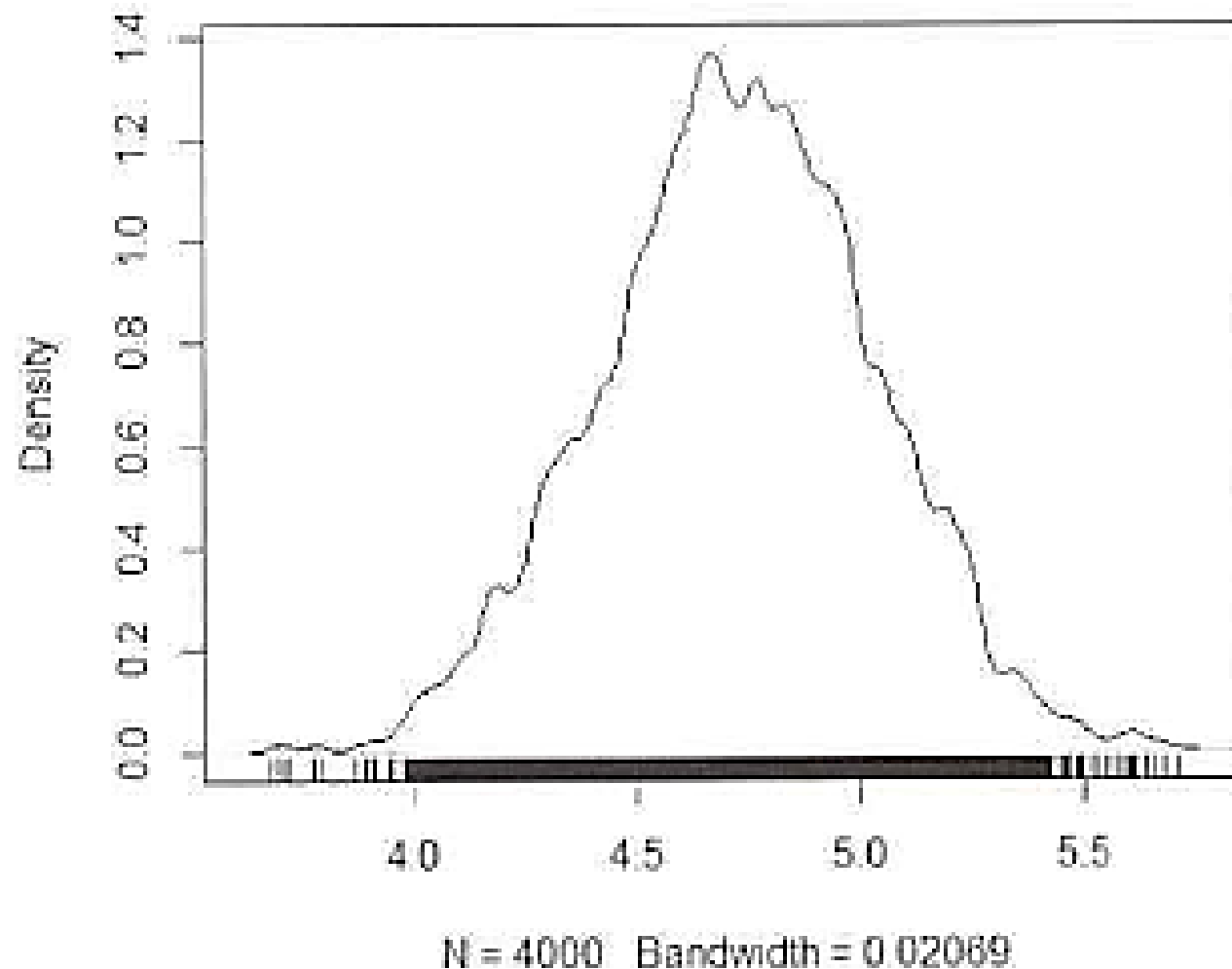
3.2.3 Visualizing a Single Variable

Histogram – Income



3.2.3 Visualizing a Single Variable

Density – Income (log10 scale)





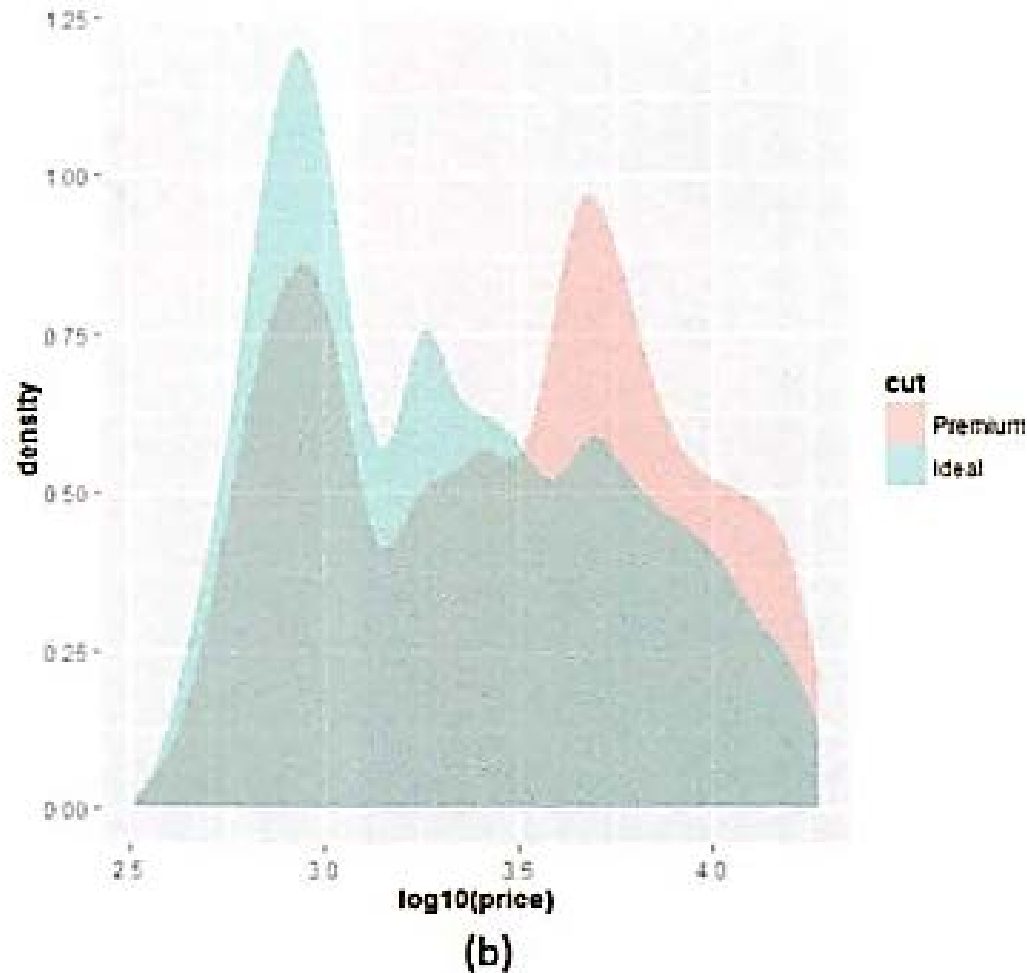
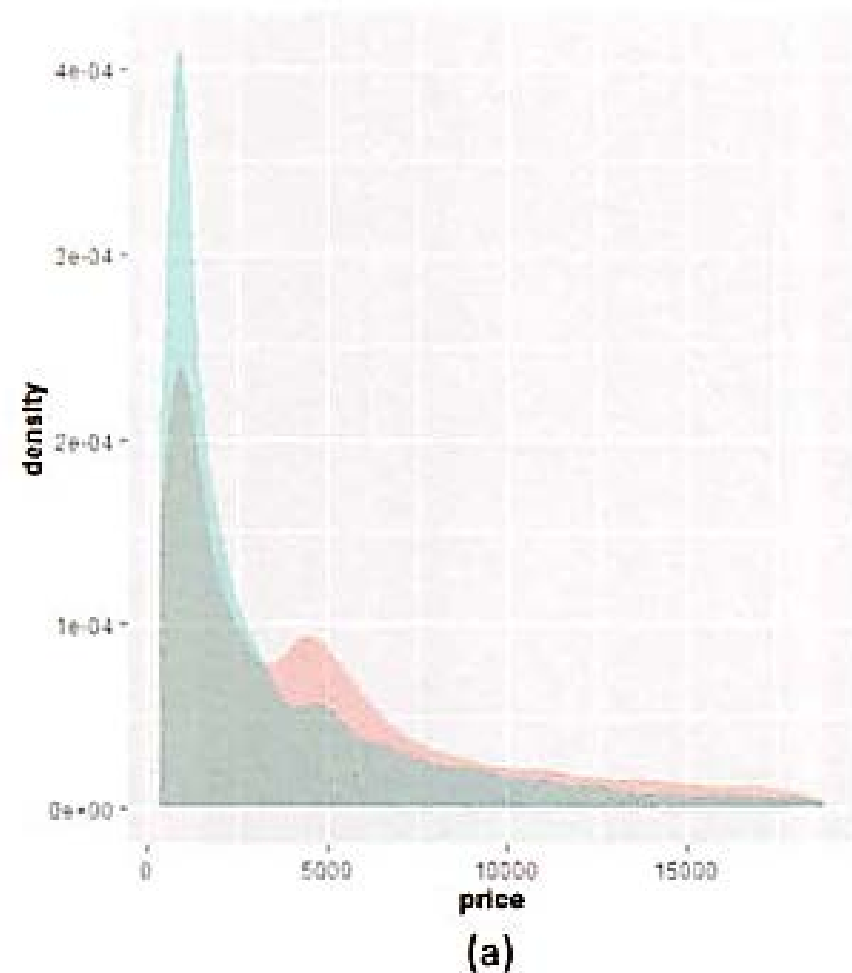
3.2.3 Visualizing a Single Variable

Density – Income (log10 scale)

- In this case, the log density plot emphasizes the log nature of the distribution
- The rug() function at the bottom creates a one-dimensional density plot to emphasize the distribution

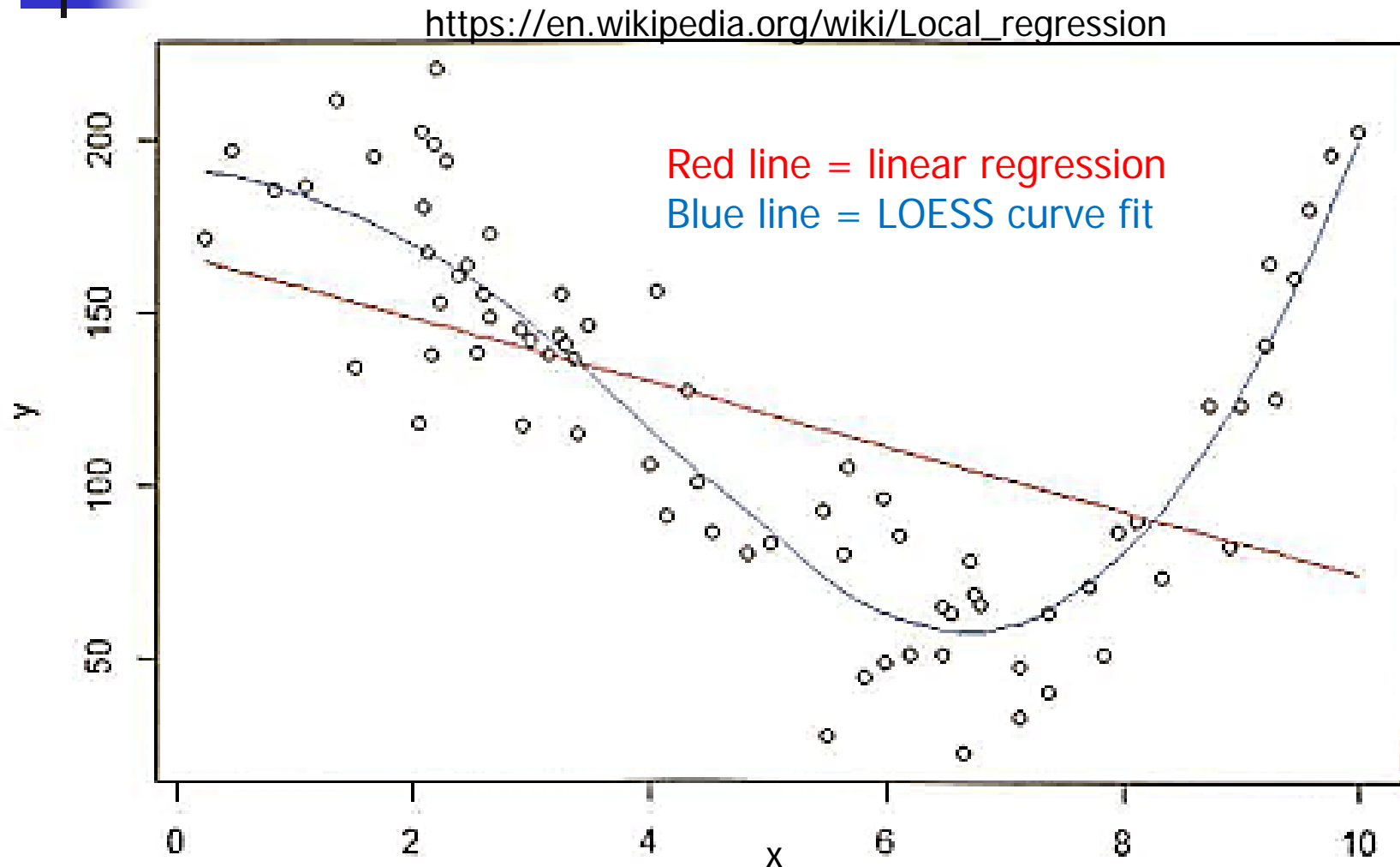
3.2.3 Visualizing a Single Variable

Density plots – Diamond prices, log of same



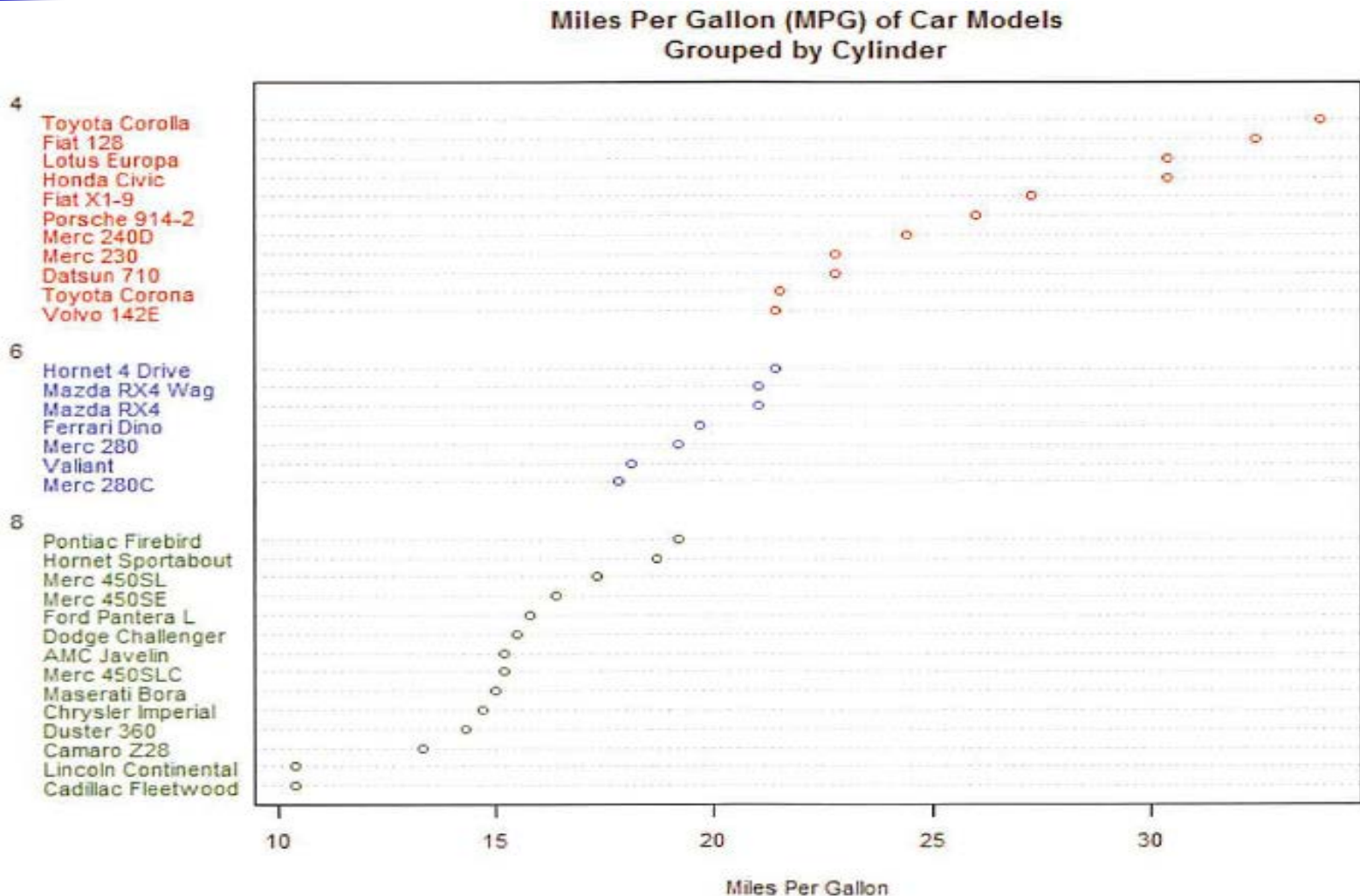
3.2.4 Examining Multiple Variables

Examining two variables with regression



3.2.4 Examining Multiple Variables

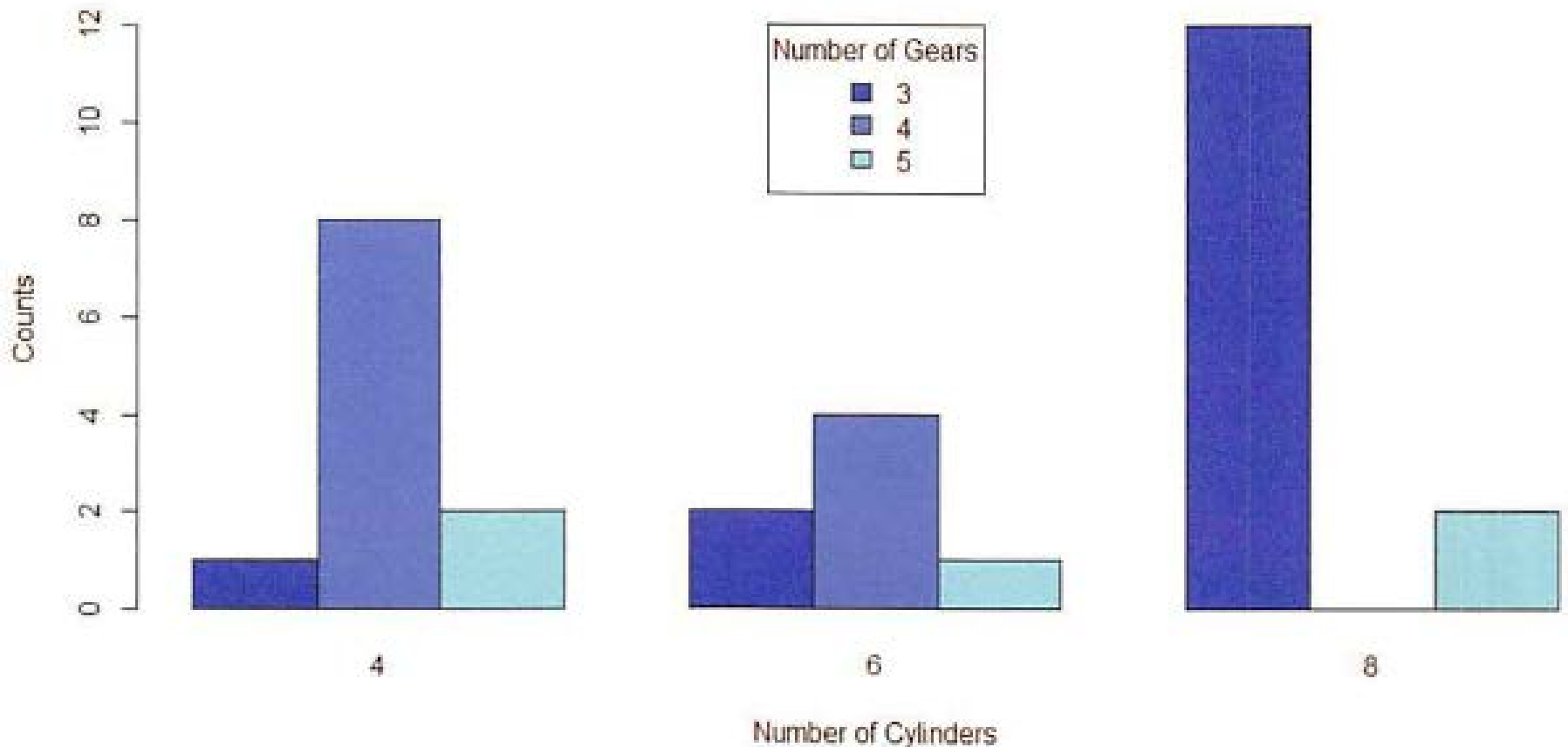
Dotchart: MPG of car models grouped by cylinder



3.2.4 Examining Multiple Variables

Barplot: visualize multiple variables

Distribution of Car Cylinder Counts and Gears



3.2.4 Examining Multiple Variables

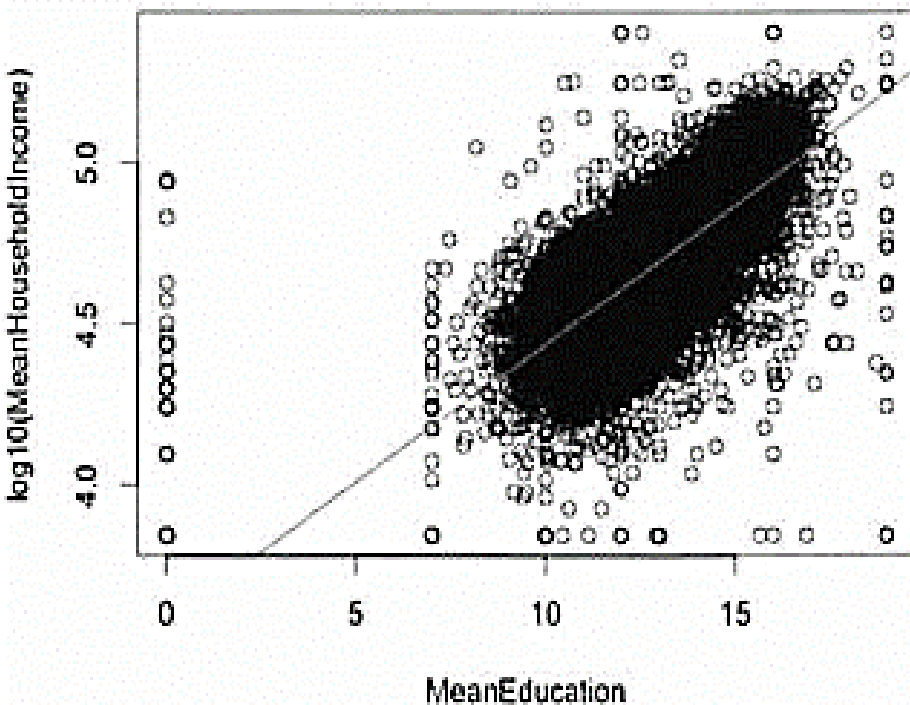
Box-and-whisker plot: income versus region



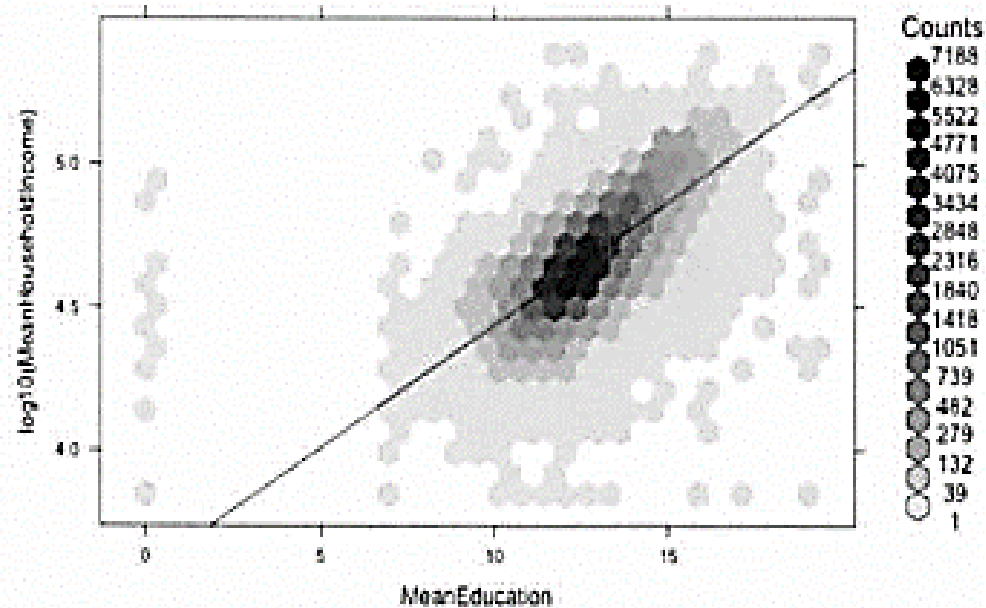
3.2.4 Examining Multiple Variables

Scatterplot (a) & Hexbinplot – income vs education

The hexbinplot combines the ideas of scatterplot and histogram
For high-volume big data hexbinplot may be better than scatterplot



(a)

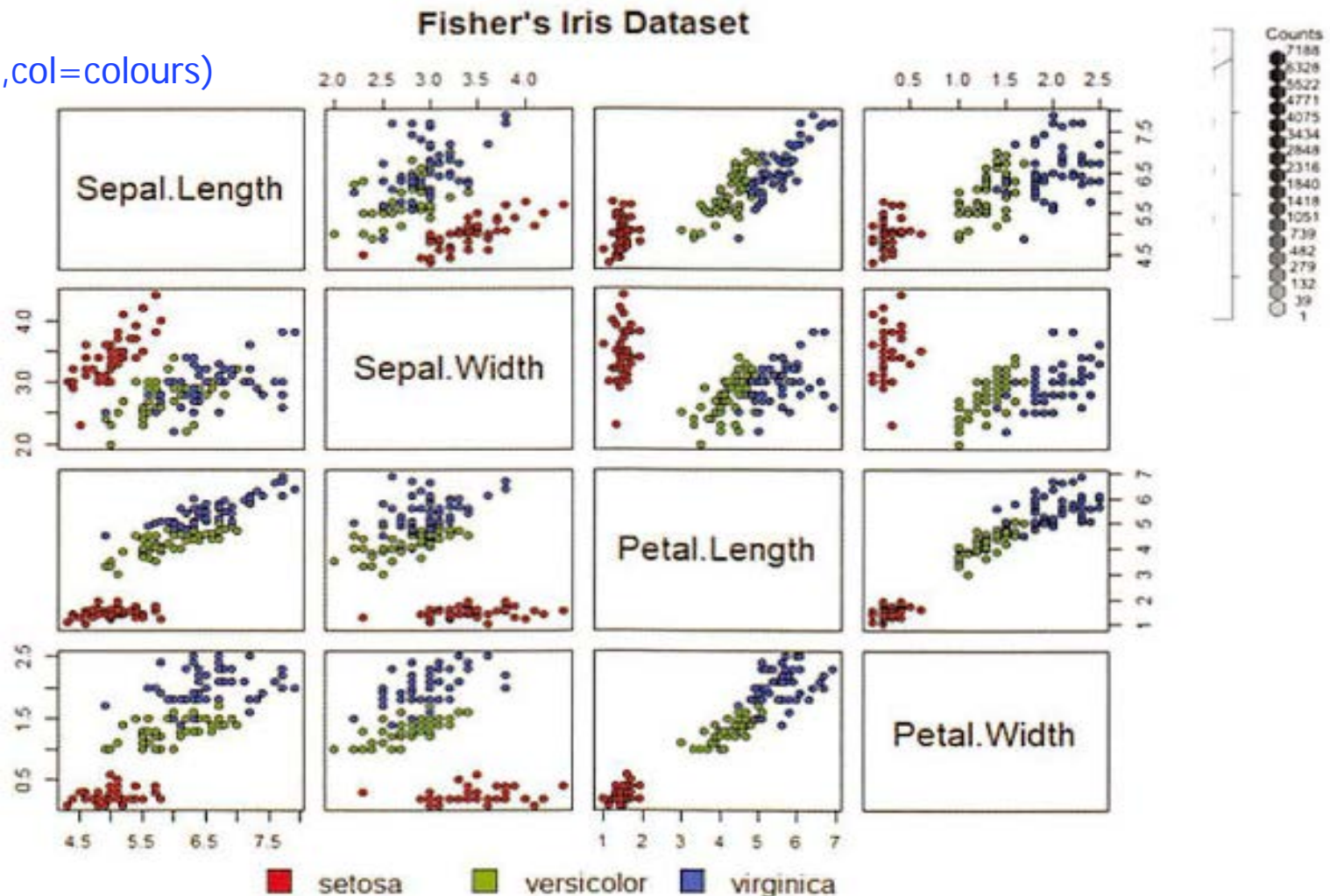


(b)

3.2.4 Examining Multiple Variables

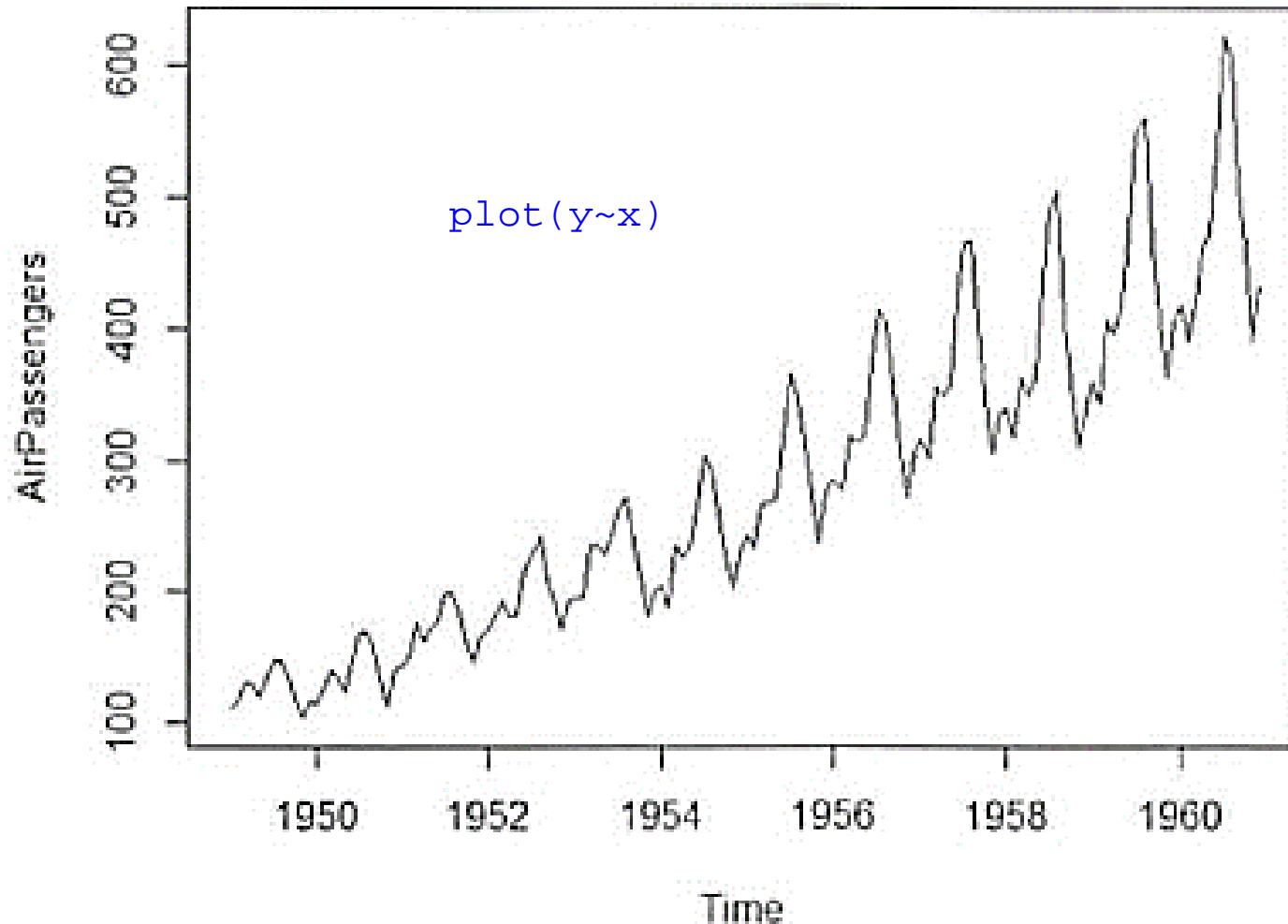
Matrix of Scatterplots

`pairs(data,col=colours)`



3.2.4 Examining Multiple Variables

Variable over time – airline passenger counts





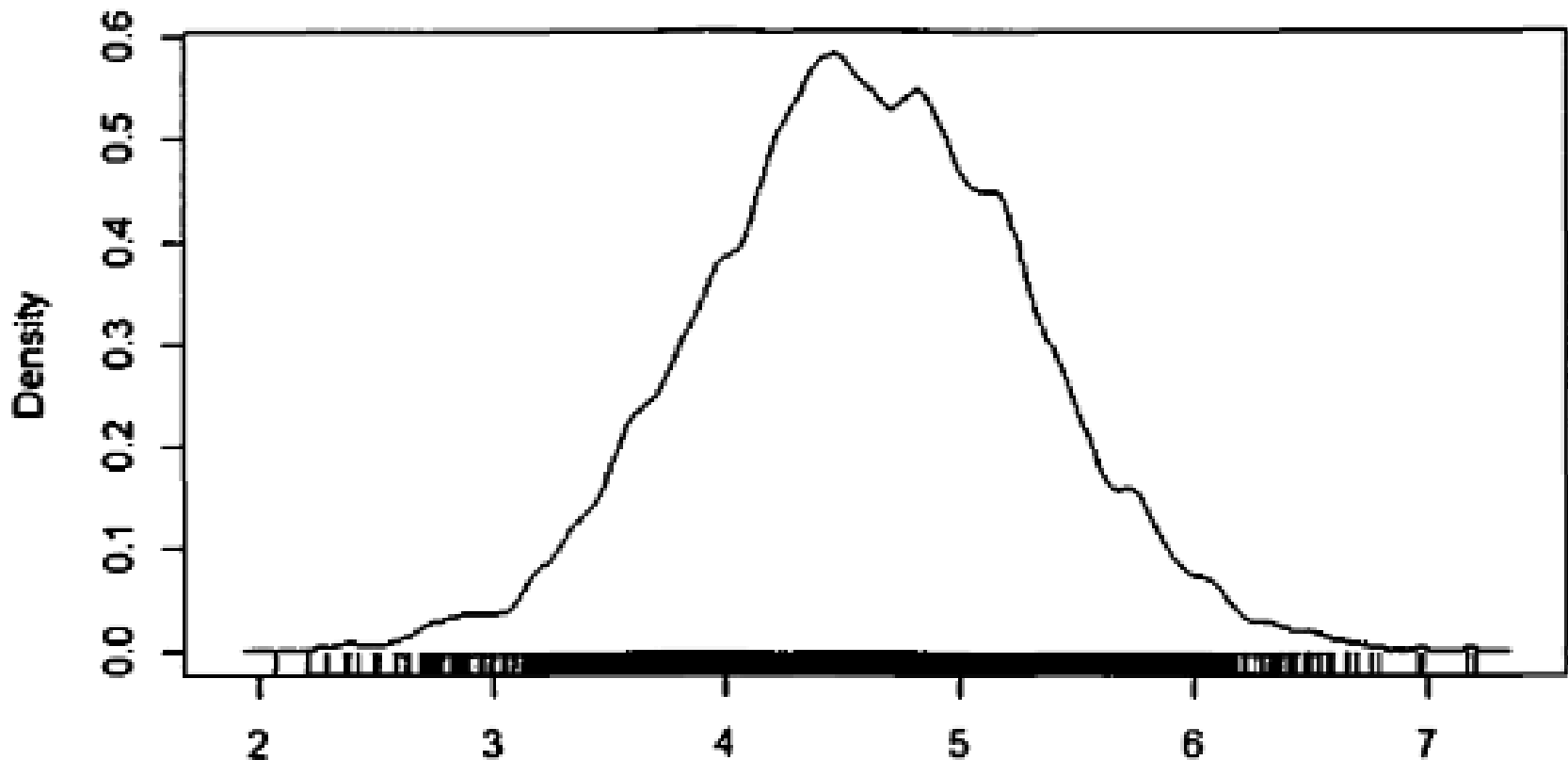
3.2.5 Exploration vs Presentation

- Data visualization for data exploration is different from presenting results to stakeholders
 - Data scientists prefer graphs that are technical in nature
 - Nontechnical stakeholders prefer simple and clear graphics that focus on the message rather than the data

3.2.5 Exploration vs Presentation

Density plots better for data scientists

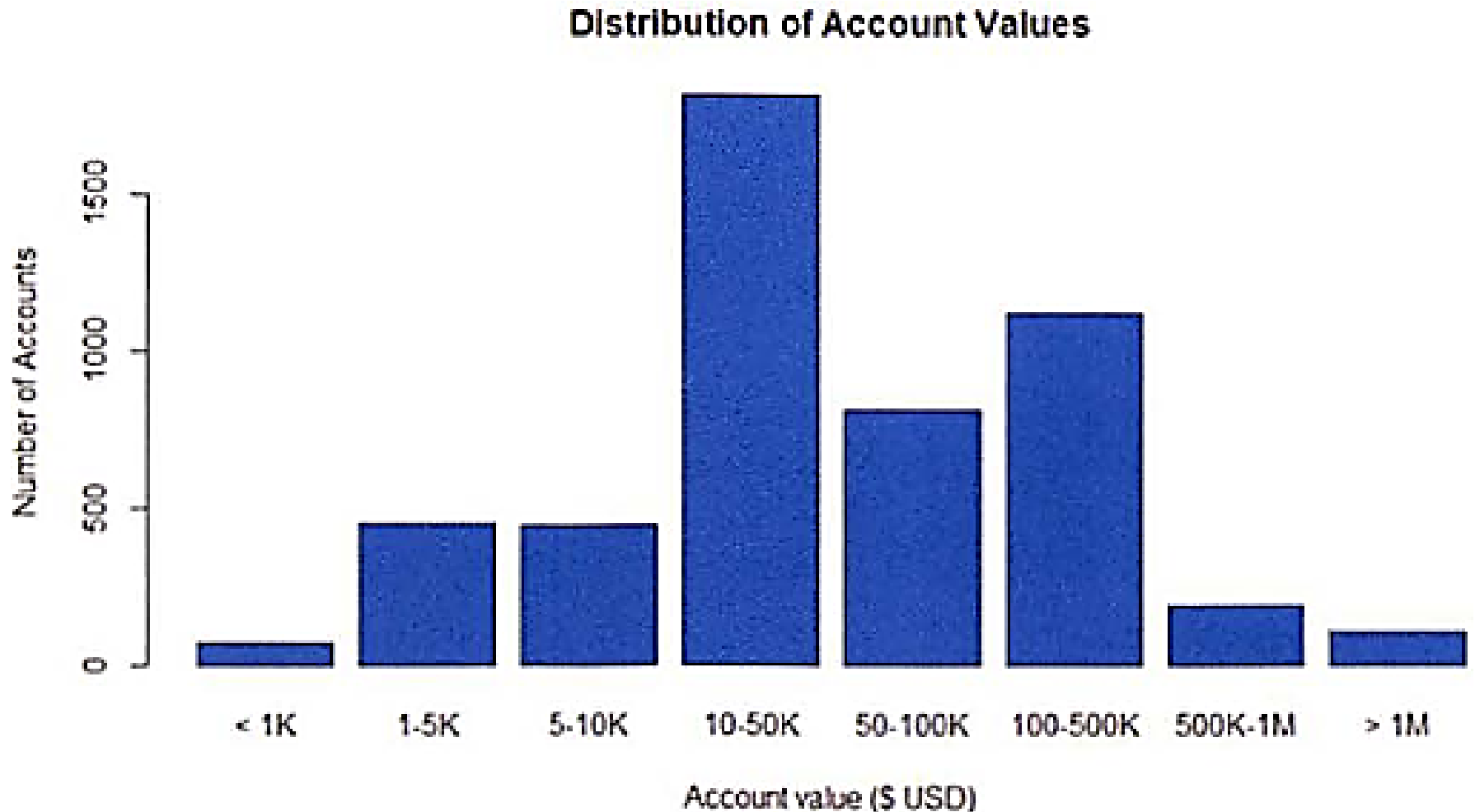
Distribution of Account Values (log10 scale)

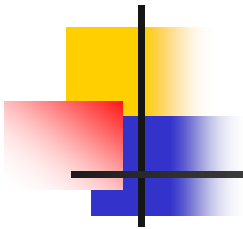


N = 5000 Bandwidth = 0.05759

3.2.5 Exploration vs Presentation

Histograms better to show stakeholders





Package Usage in R

Many useful R function come in packages, free libraries of code written by R's active user community. To install an R package, open an R session and type at the command line

```
install.packages("<the package's name>")
```

R will download the package from CRAN, so you'll need to be connected to the internet. Once you have a package installed, you can make its contents available to use in your current R session by running

```
library("<the package's name>")
```

There are thousands of helpful R packages for you to use, but navigating them all can be a challenge. To help you out, we've compiled this guide to some of the best. We've used each of these, and found them to be outstanding – we've even written some of them. But you don't have to take our word for it, these packages are also some of the top most downloaded R packages.



To load data in R

[DBI](#) - The standard for communication between R and relational database management systems. Packages that connect R to databases depend on the DBI package.

[odbc](#) - Use any ODBC driver with the odbc package to connect R to your database.

[RMySQL](#), [RPostgresSQL](#), [RSQLite](#) - If you'd like to read in data from a database, these packages are a good place to start. Choose the package that fits your type of database.

[XLConnect](#), [xlsx](#) - These packages help you read and write Microsoft Excel files from R. You can also just export your spreadsheets from Excel as .csv's.

[foreign](#) - Want to read a SAS data set into R? Or an SPSS data set? Foreign provides functions that help you load data files from other programs into R.

[haven](#) - Enables R to read and write data from SAS, SPSS, and Stata.

R can handle plain text files – no package required. Just use the functions `read.csv`, `read.table`, and `read.fwf`. If you have even more exotic data, consult the CRAN [guide](#) to data import and export.

and more.....



To manipulate data in R

[dplyr](#) - Essential shortcuts for subsetting, summarizing, rearranging, and joining together data sets. dplyr is our go to package for fast data manipulation.

[tidyr](#) - Tools for changing the layout of your data sets. Use the gather and spread functions to convert your data into the [tidy format](#), the layout R likes best.

[stringr](#) - Easy to learn tools for regular expressions and character strings.

[lubridate](#) - Tools that make working with dates and times easier.

and more.....



To visualize data in R

[ggplot2](#) - R's famous package for making beautiful graphics. ggplot2 lets you use the [grammar of graphics](#) to build layered, customizable plots.

[ggvis](#) - Interactive, web based graphics built with the grammar of graphics.

[rgl](#) - Interactive 3D visualizations with R

[htmlwidgets](#) - A fast way to build interactive (javascript based) visualizations with R.

Packages that implement htmlwidgets include:

- [leaflet](#) (maps)
- [dygraphs](#) (time series)
- [DT](#) (tables)
- [diagrammeR](#) (diagrams)
- [network3D](#) (network graphs)
- [threeJS](#) (3D scatterplots and globes)

[googleVis](#) - Let's you use Google Chart tools to visualize data in R.

and more.....



To model data in R

[car](#) - car's [Anova](#) function is popular for making type II and type III Anova tables.

[mgcv](#) - Generalized Additive Models

[lme4](#)/[nlme](#) - Linear and Non-linear mixed effects models

[randomForest](#) - Random forest methods from machine learning

[multcomp](#) - Tools for multiple comparison testing

[vcd](#) - Visualization tools and tests for categorical data

[glmnet](#) - Lasso and elastic-net regression methods with cross validation

[survival](#) - Tools for survival analysis

[caret](#) - Tools for training regression and classification models

and more.....

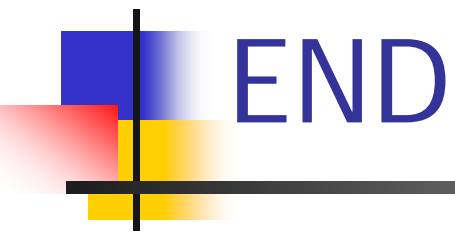
The figure displays a 6x6 grid of 36 diverse plots illustrating data visualization techniques for temperature and stress level data. The plots include:

- Heatmaps:** Visualizing the relationship between temperature and stress levels, showing patterns of high and low values.
- Line Graphs:** Tracking temperature and stress levels over time, highlighting seasonal and daily cycles.
- Box Plots:** Summarizing the distribution of stress levels across different temperature categories.
- Scatter Plots:** Exploring the correlation between temperature and stress levels, with some plots showing regression lines.
- Histograms:** Displaying the frequency distribution of stress levels for different temperature ranges.
- Violin Plots:** Comparing the distribution of stress levels across various temperature categories.
- Density Plots:** Visualizing the probability density of stress levels for different temperature ranges.
- Other Visualizations:** Including a pie chart showing the proportion of days exceeding a stress level, a bar chart showing the proportion of days exceeding a stress level, and a series of small plots showing the distribution of stress levels for different temperature ranges.



Summary

- 3.1 Introduction to R and RStudio
- 3.2 Exploratory Data Analysis using R
 - 3.2.1 Visualization before Analysis
 - 3.2.2 Dirty Data
 - 3.2.3 Visualizing a Single Variable
 - 3.2.4 Examining Multiple Variables
 - 3.2.5 Data Exploration versus Presentation
- Exercises available in the textbook*



END

3.3 Statistical Methods for Evaluation

Statistics helps answer data analytics questions

- Model Building

- What are the best input variables for the model?
- Can the model predict the outcome given the input?

- Model Evaluation

- Is the model accurate?
- Does the model perform better than an obvious guess?
- Does the model perform better than other models?

- Model Deployment

- Is the prediction sound?
- Does model have the desired effect (e.g., reducing cost)?



3.3 Statistical Methods for Evaluation Subsections

- 3.3.1 Hypothesis Testing
- 3.3.2 Difference of Means
- 3.3.3 Wilcoxon Rank-Sum Test
- 3.3.4 Type I and Type II Errors
- 3.3.5 Power and Sample Size
- 3.3.6 ANOVA (Analysis of Variance)



3.3.1 Hypothesis Testing

- Basic concept is to form an assertion and test it with data
- Common assumption is that there is no difference between samples (default assumption)
- Statisticians refer to this as the *null hypothesis* (H_0)
- The *alternative hypothesis* (H_A) is that there is a difference between samples



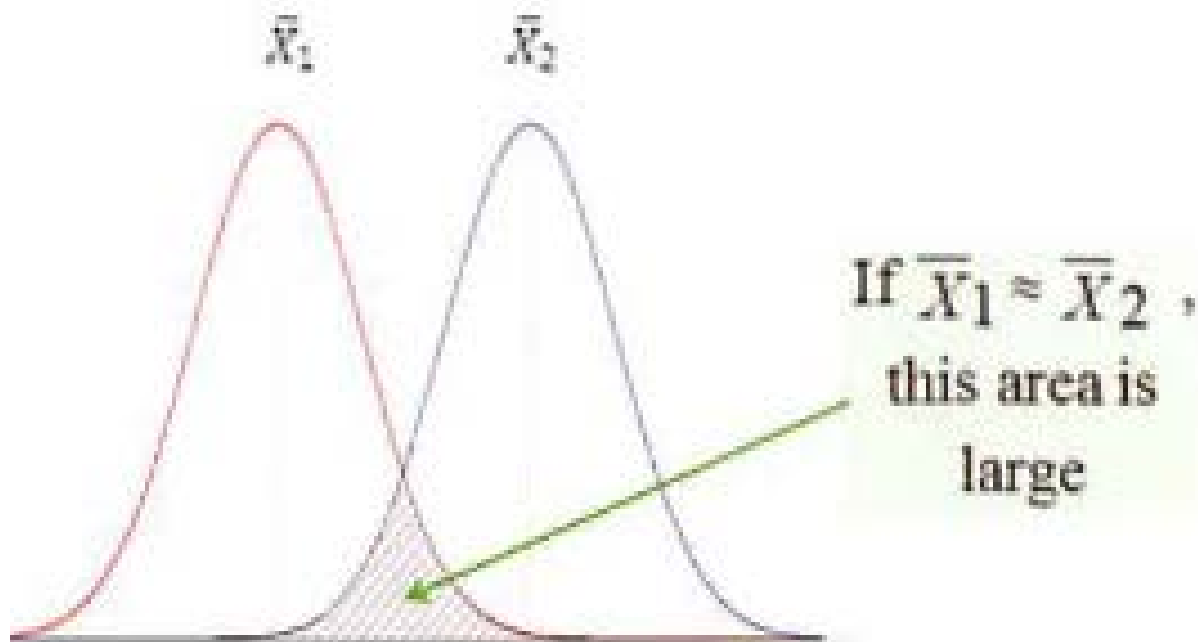
3.3.1 Hypothesis Testing

Example Null and Alternative Hypotheses

Application	Null Hypothesis	Alternative Hypothesis
Accuracy Forecast	Model X <i>does not predict</i> better than the existing model.	Model X <i>predicts</i> better than the existing model.
Recommendation Engine	Algorithm Y <i>does not produce</i> better recommendations than the current algorithm being used.	Algorithm Y <i>produces</i> better recommendations than the current algorithm being used.
Regression Modeling	This variable <i>does not affect</i> the outcome because its coefficient is zero.	This variable <i>affects</i> outcome because its coefficient is not zero.

3.3.2 Difference of Means

Two populations – same or different?





3.3.2 Difference of Means

Two Parametric Methods

- Student's t-test
 - Assumes two normally distributed populations, and that they have equal variance
- Welch's t-test
 - Assumes two normally distributed populations, and they don't necessarily have equal variance



3.3.3 Wilcoxon Rank-Sum Test

A Nonparametric Method

- Makes no assumptions about the underlying probability distributions



3.3.4 Type I and Type II Errors

- An hypothesis test may result in two types of errors
 - Type I error – rejection of the null hypothesis when the null hypothesis is TRUE
 - Type II error – acceptance of the null hypothesis when the null hypothesis is FALSE



3.3.4 Type I and Type II Errors

	H_0 is true	H_0 is false
H_0 is accepted	Correct outcome	<i>Type II Error</i>
H_0 is rejected	<i>Type I error</i>	Correct outcome

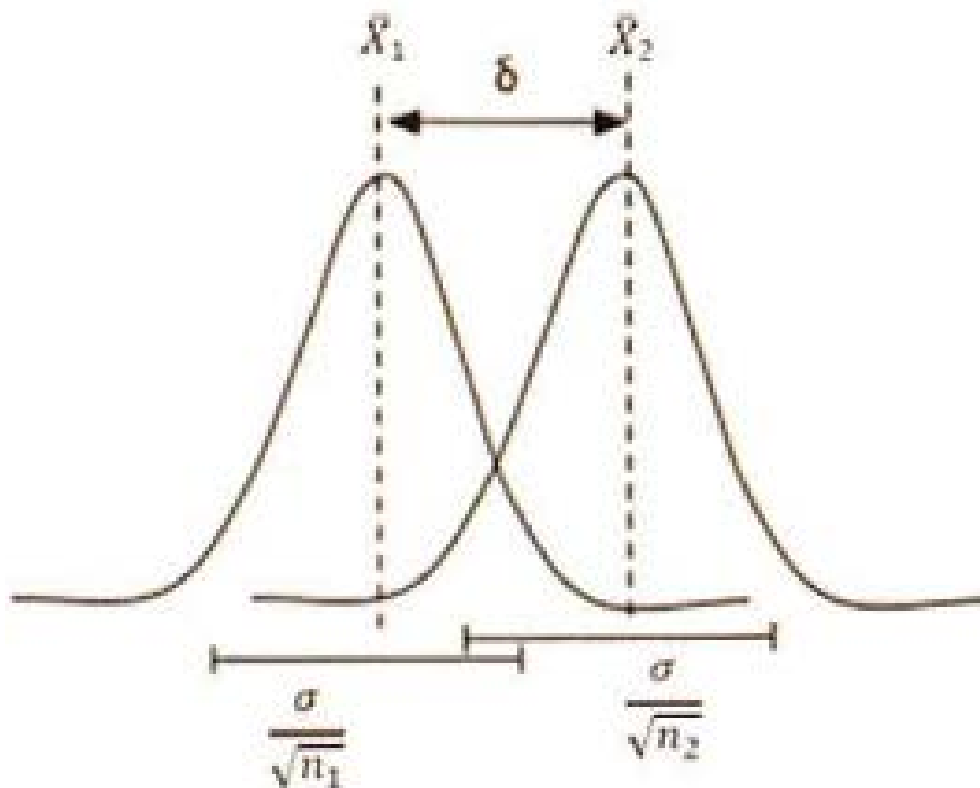


3.3.5 Power and Sample Size

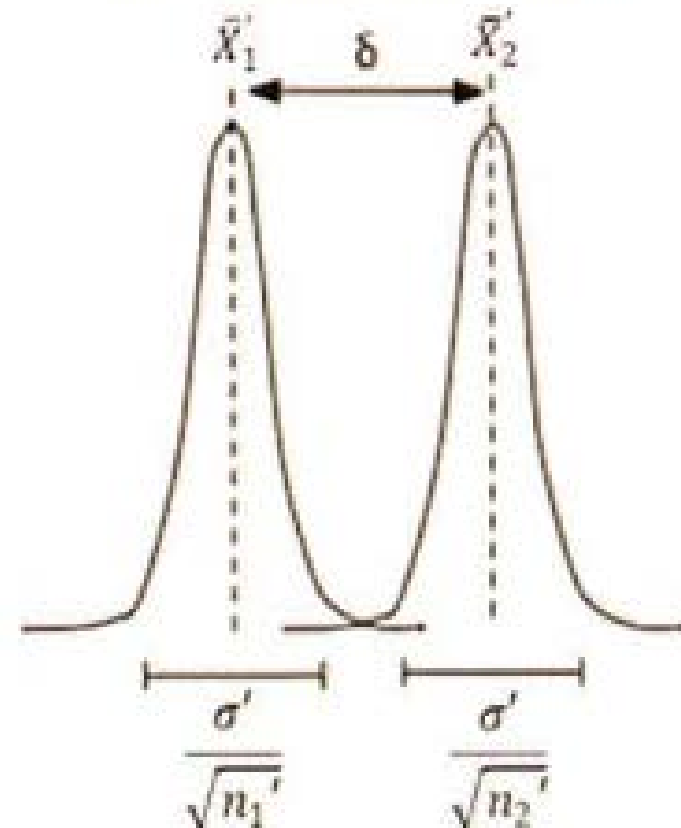
- The *power of a test* is the probability of correctly rejecting the null hypothesis
- The power of a test increases as the sample size increases
- *Effect size δ* = difference between the means
- It is important to consider an appropriate effect size for the problem at hand

3.3.5 Power and Sample Size

Moderate Sample Size



Larger Sample Size





3.3.6 ANOVA (Analysis of Variance)

- A generalization of the hypothesis testing of the difference of two population means
- Good for analyzing more than two populations
- ANOVA tests if any of the population means differ from the other population means