

Data Exploration using R

Statistics Refresher Workshop

Kai Xiong

k.xiong@auckland.ac.nz
Statistical Consulting Service
The Department of Statistics
The University of Auckland

July 1, 2011



- Checking the data for patterns, relationships, structures, and other features.
- Can be done through graphs and summary statistics.
- We recommend R for generating graphs, much more flexible than Excel, SPSS, SAS.



Roles of Exploratory Graphics

- 1 Are there any errors or outliers?
- 2 Are there any patterns in the data?
 - Symmetric, Skewed, Bimodal, Clusters?
- 3 Are there any relationships between the variables?
 - Linear (increasing, decreasing), polynomial, exponential...
- 4 What sort of model might be appropriate?



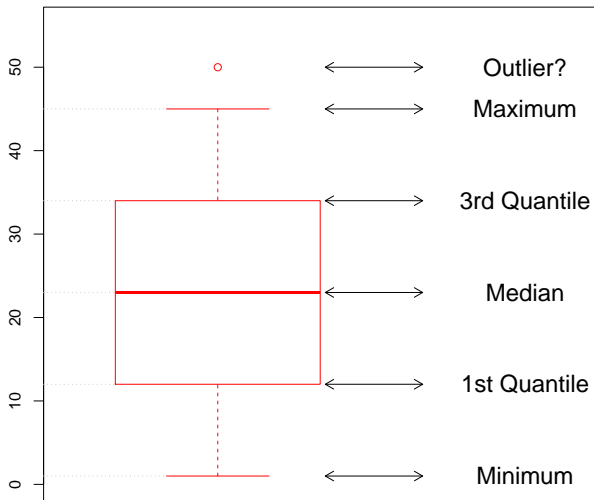
- 1 Quantitative
 - Continuous:
e.g. gene expression, Body length of mussels
 - Discrete:
e.g. number of photophores in lantern fish
- 2 Qualitative
 - Categorical
e.g. SNPs, Location of marine reserves
 - Ordinal
e.g. Position in a food chain

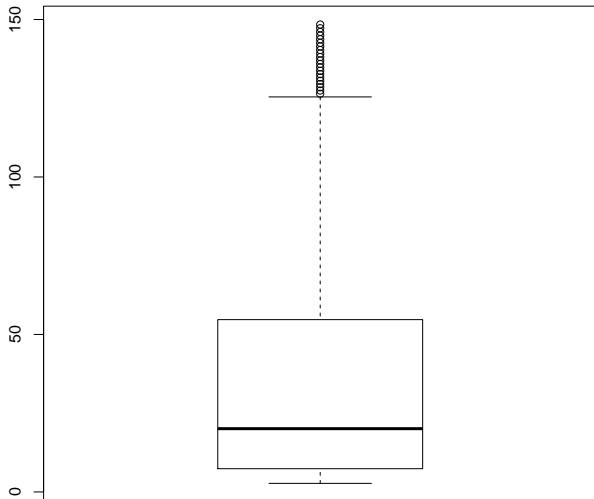


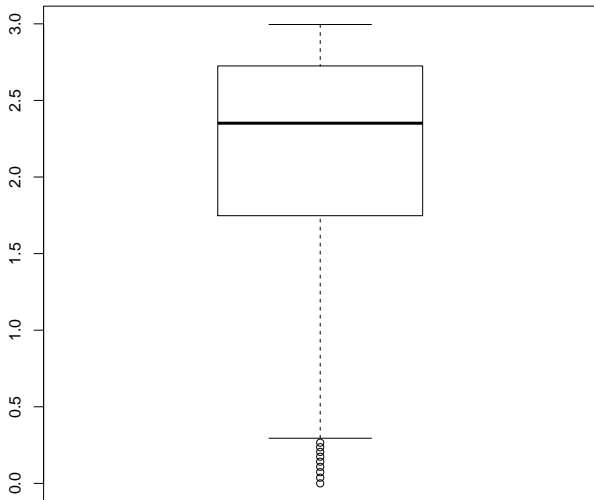
One Quantitative Variable

- Dotplots
- Five-number-summary statistics
 - 1 Minimum
 - 2 1st Quantile/Lower Quantile/25th percentile
 - 25% of the sorted variable is smaller than the 1st quantile
 - 75% of the sorted variable is greater than the 1st quantile
 - 3 2nd Quantile/Median/50th percentile
 - Divide the sorted variable in two equal halves
 - 4 3rd Quantile/Upper Quantile/75th percentile
 - 25% of the sorted variable is greater than the 3rd quantile
 - 75% of the sorted variable is smaller than the 3rd quantile
 - 5 Maximum
- Boxplot: visual display of 5-number-summary statistics.
- Normal QQ plot
- Histogram

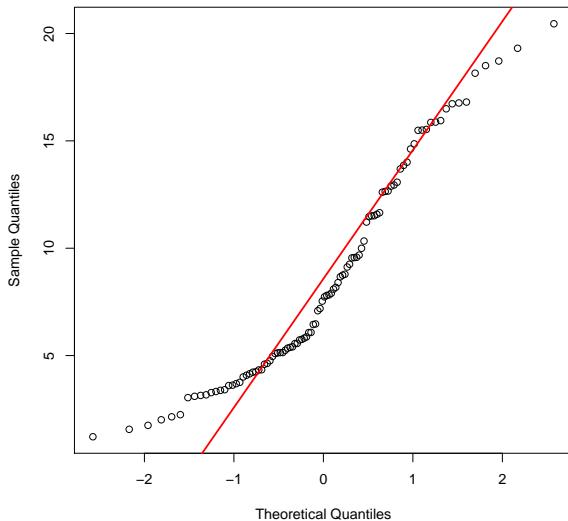




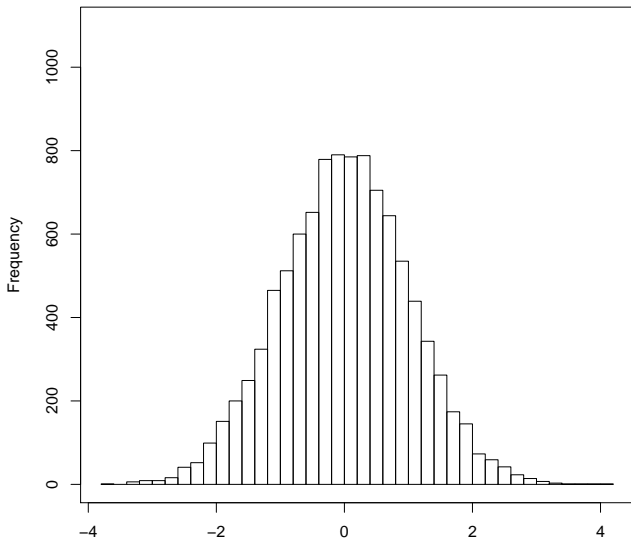




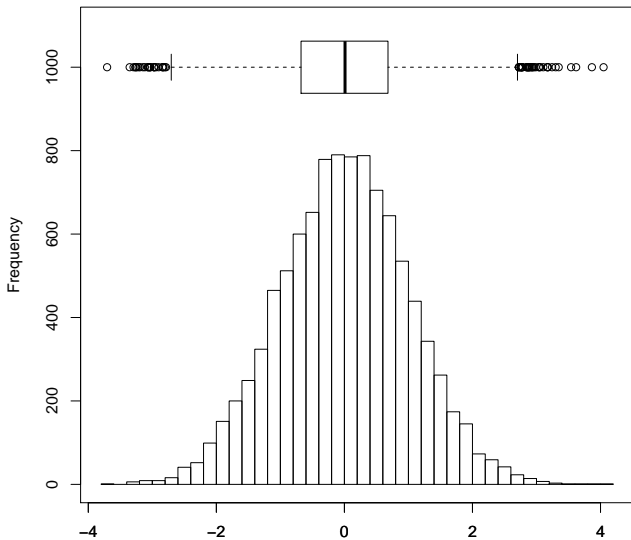
- Normal Quantile-Quantile (QQ) Plot



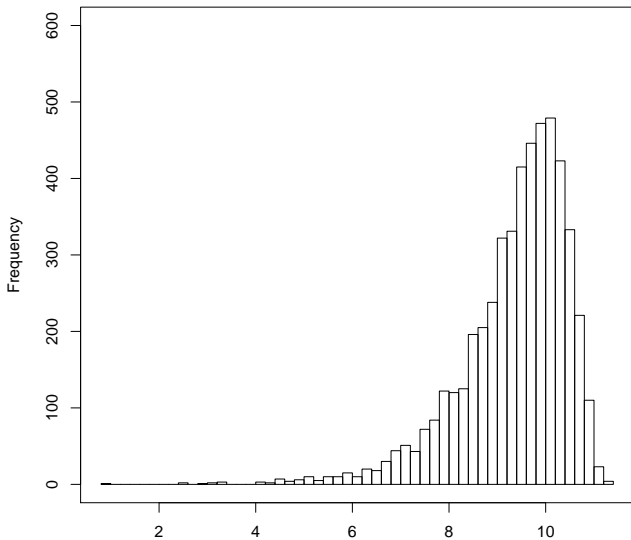
Histogram showing Normal Distribution



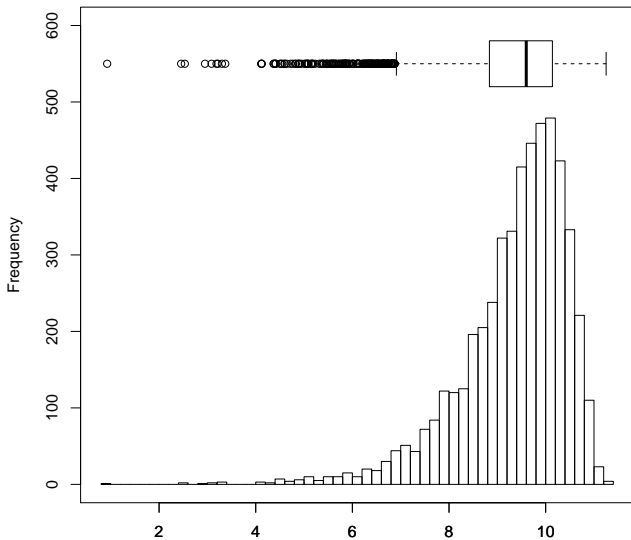
Histogram and boxplot showing Normal Distribution



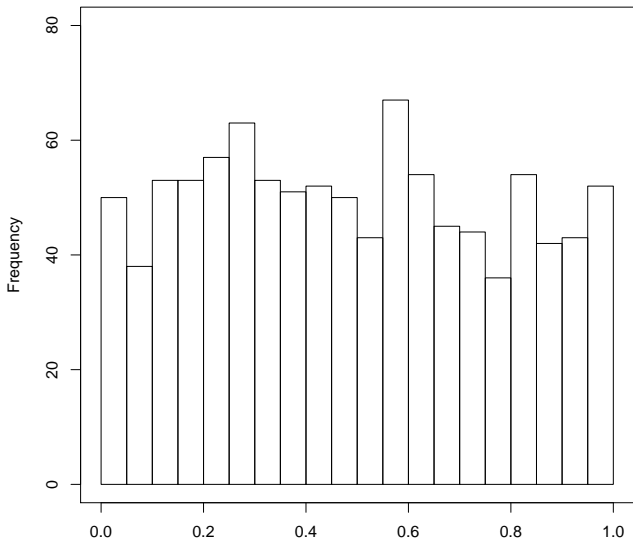
Histogram showing left skewed distribution



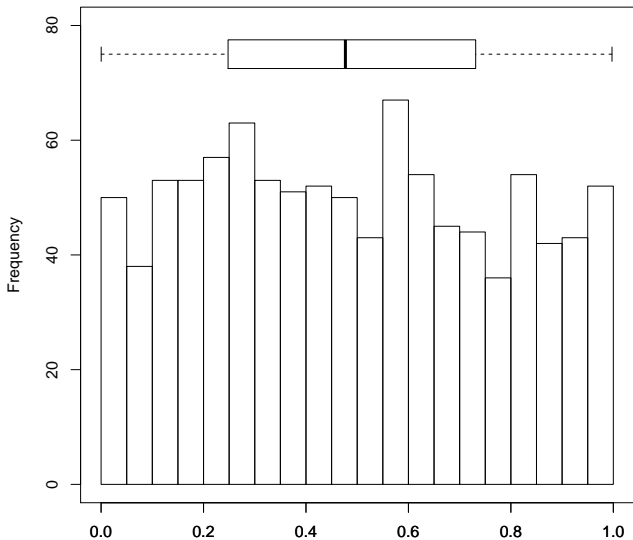
Histogram and boxplot showing left skewed istribution



Histogram showing uniform distribution



Histogram and boxplot showing uniform distribution



One Qualitative Variable

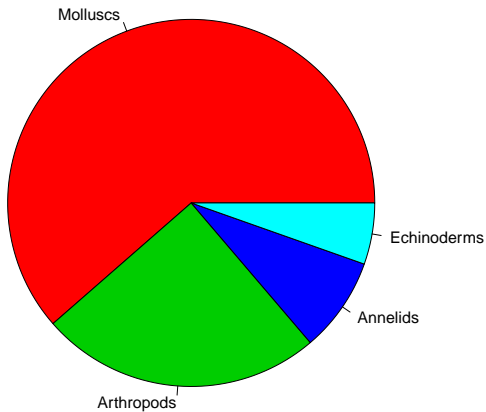
- Frequency table

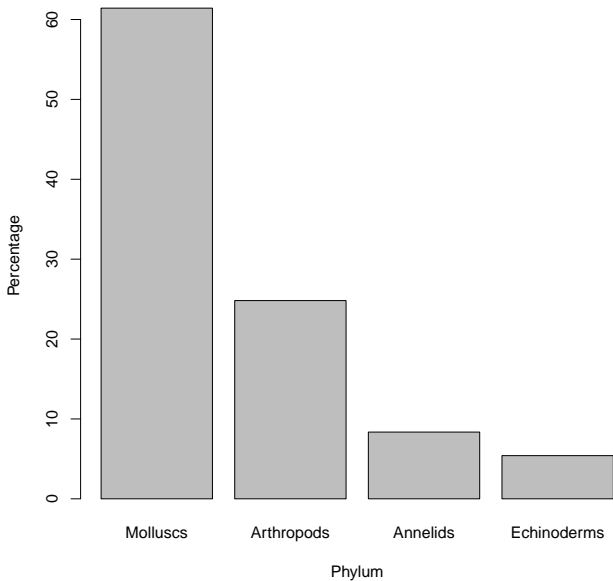
Table: One way frequency table

Phylum	Frequency	Percentage
Molluscs	250	61.4%
Annelids	34	8.4%
Arthropods	101	24.8%
Echinoderms	22	5.4%
Total	407	100%

- Pie chart
- Barplot



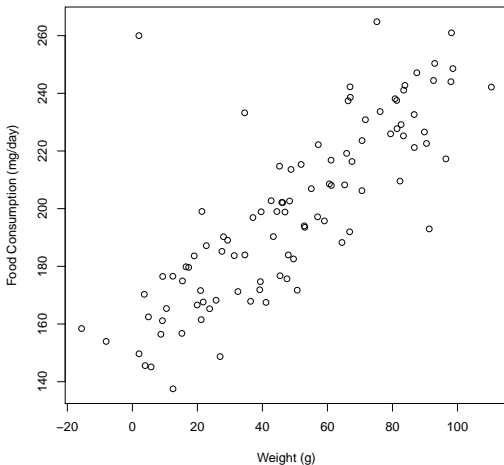




Two Variables

Two Quantitative

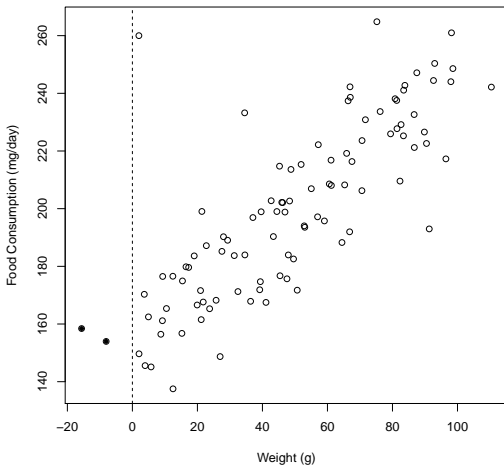
- Scatter Plot



Two Variables

Two Quantitative

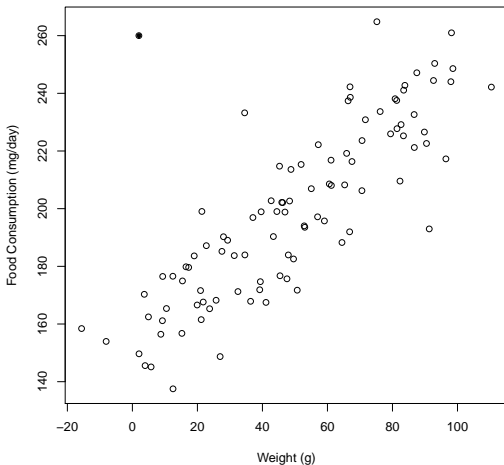
- Scatter Plot



Two Variables

Two Quantitative

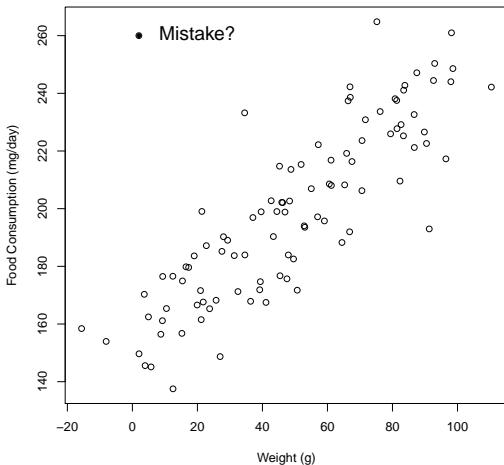
- Scatter Plot



Two Variables

Two Quantitative

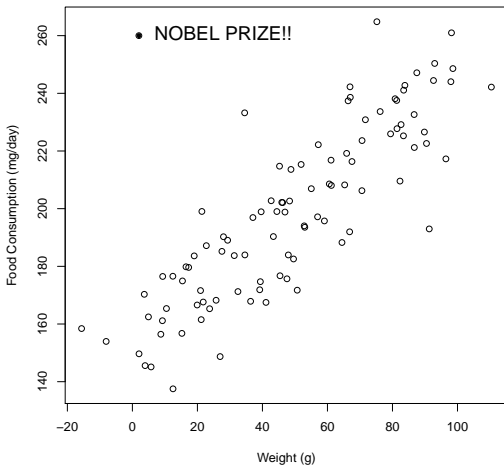
- Scatter Plot



Two Variables

Two Quantitative

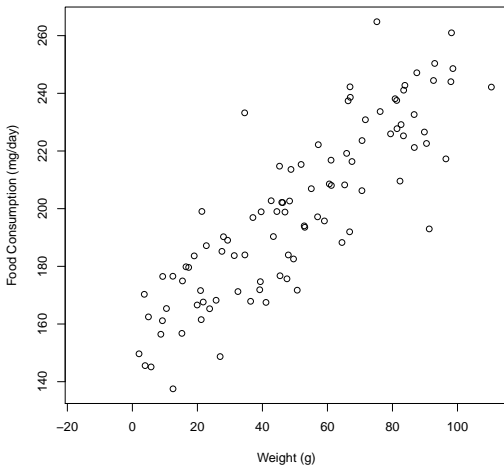
- Scatter Plot



Two Variables

Two Quantitative

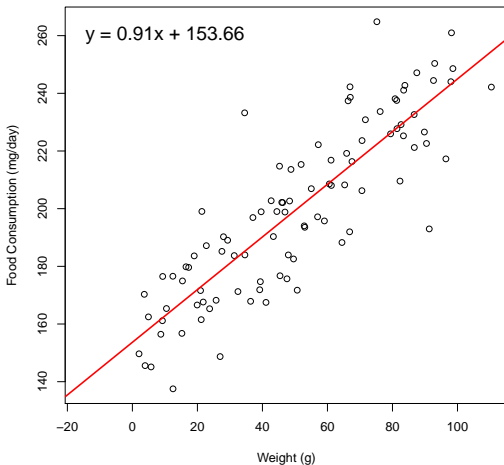
- Scatter Plot



Two Variables

Two Quantitative

- Scatter Plot



Two Variables

Two Qualitative

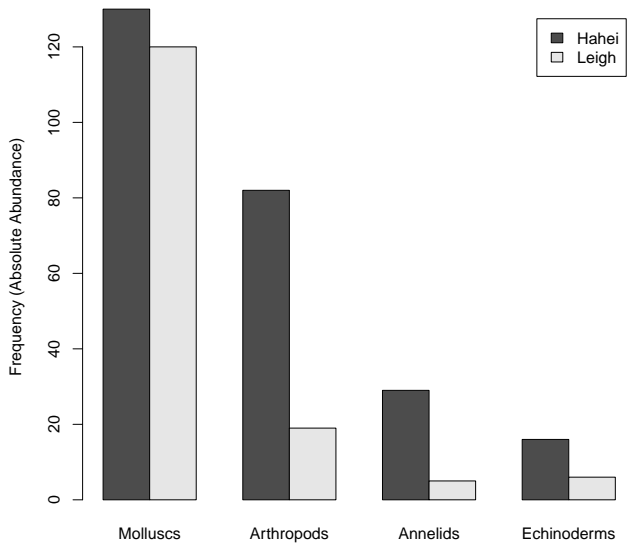
- Two way frequency table

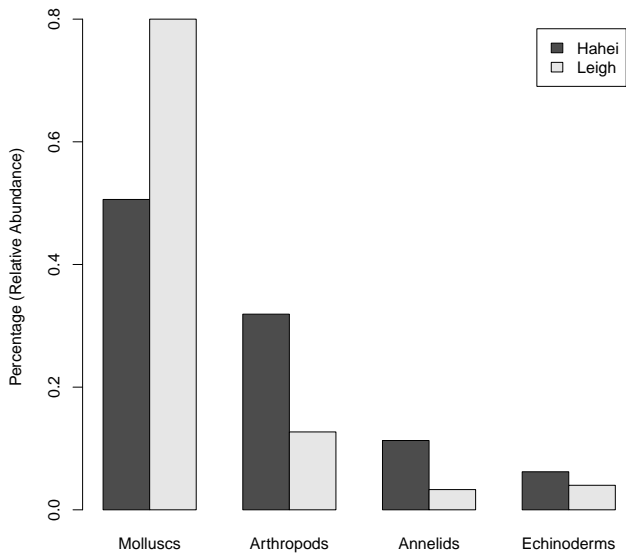
Table: Two way frequency table

Phylum	Hahei	Leigh	Total
Molluscs	130	120	250
Annelids	29	5	34
Arthropods	82	19	101
Echinoderms	16	6	22
Total	257	150	407

- Side by side barplot: frequencies or percentages?



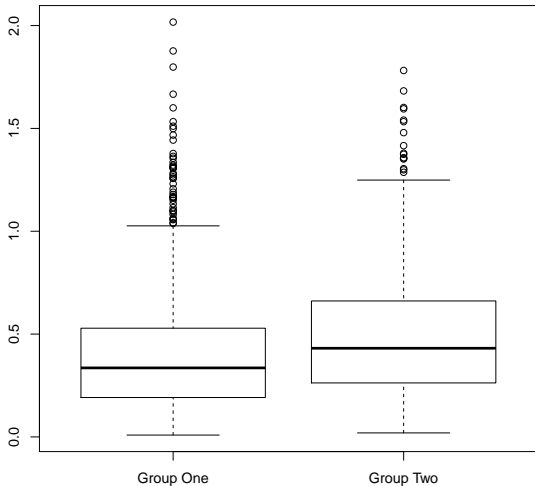




Two Variables

One Qualitative and One Quantitative

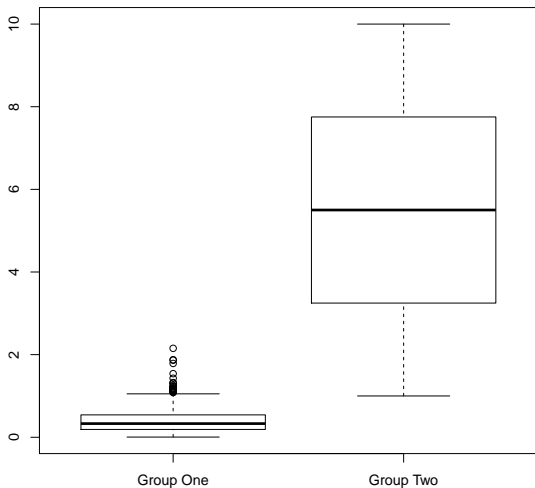
- Side-by-side boxplot



Two Variables

One Qualitative and One Quantitative

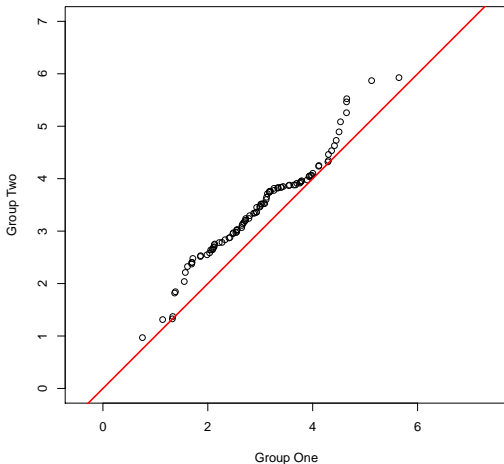
- Side-by-side boxplot



Two Variables

One Qualitative and One Quantitative

- QQ Plot



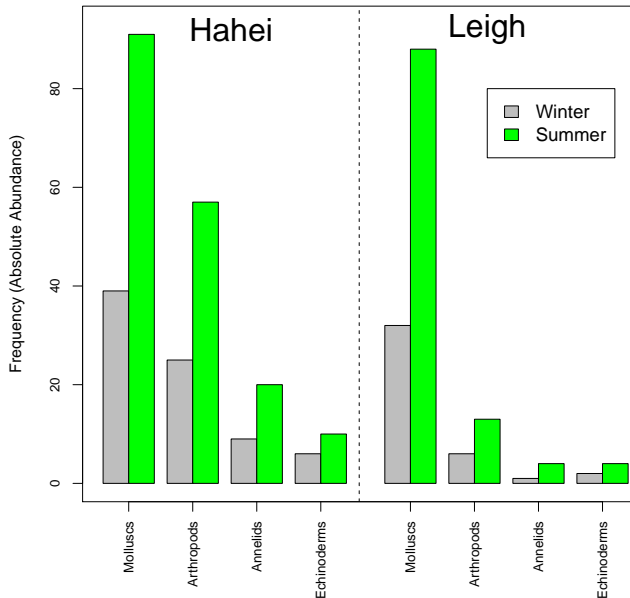
Three Variables

Three Qualitative Variables

Table: Three way frequency table

	Location			
	Hahei		Leigh	
Phylum	Winter	Summer	Winter	Summer
Molluscs	39	91	32	88
Annelids	9	20	1	4
Arthropods	25	57	6	13
Echinoderms	6	10	2	4



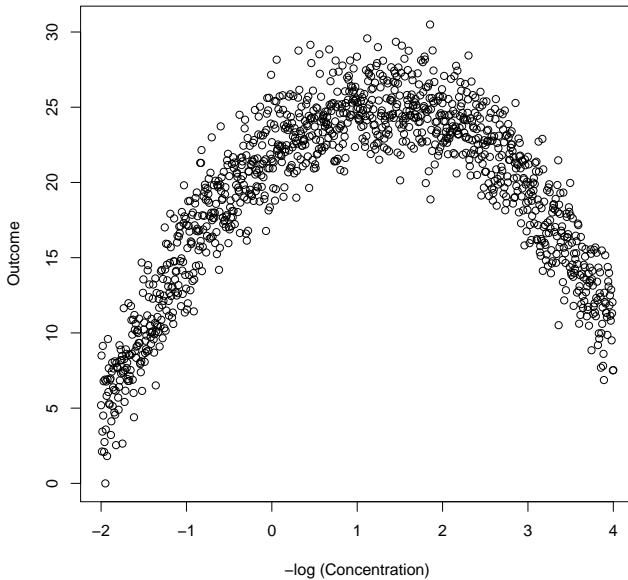


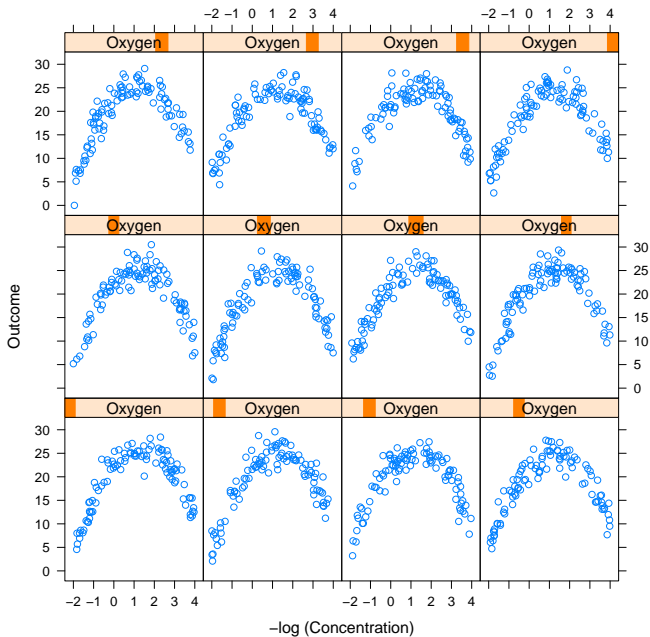
Three Variables

Three Quantitative Variables

- Trellis Graphs: Display a variable or the relationship between variables, conditioned on one or more other variables.
 - Three continuous variables $-\log(\text{Concentration})$, *Outcome* and *Oxygen level*.
 - What is the relationship between $-\log(\text{Concentration})$ and *Outcome*, conditioned on *Oxygen level*? In other words, how does the relationship between $-\log(\text{Concentration})$ and *Outcome* change over different *Oxygen level*?
 - 1 Plot $-\log(\text{Concentration})$ against *Outcome* first.
 - 2 Break *Oxygen level* into ordinal groups. e.g. Divide *Oxygen level* into 12 equally spaced interval.

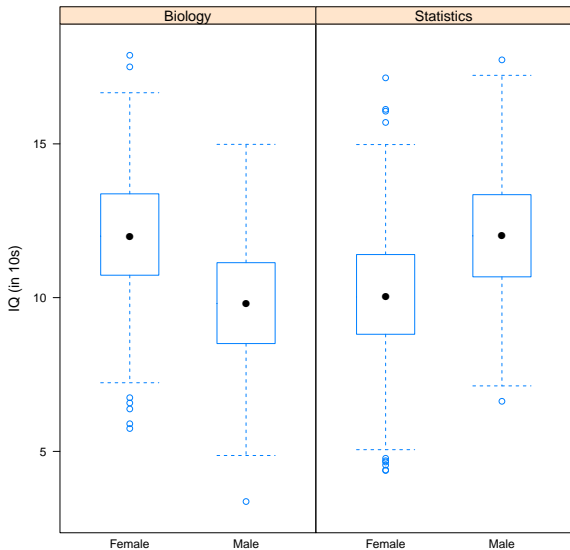






Three Variables

Two Qualitative Variables and One Quantitative Variable



Three Variables

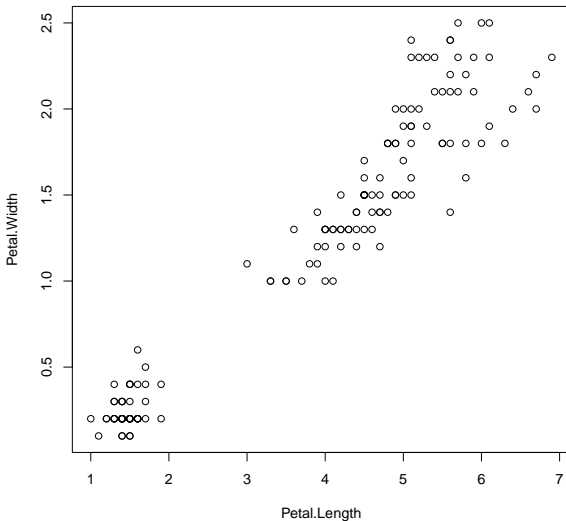
Two Quantitative Variables and One Qualitative Variable

- *Iris* flower data set
- Four continuous measure on Petal Length, Petal Width, Sepal Length and Sepal Width.
- One nominal variable, Species.
- Suppose we are interested in the relationship between Petal Length and Petal Width, and how such relationship changes for different species.
 - 1 Plot Petal Length against Petal Width.
 - 2 Use different plotting character/color to represent different species.



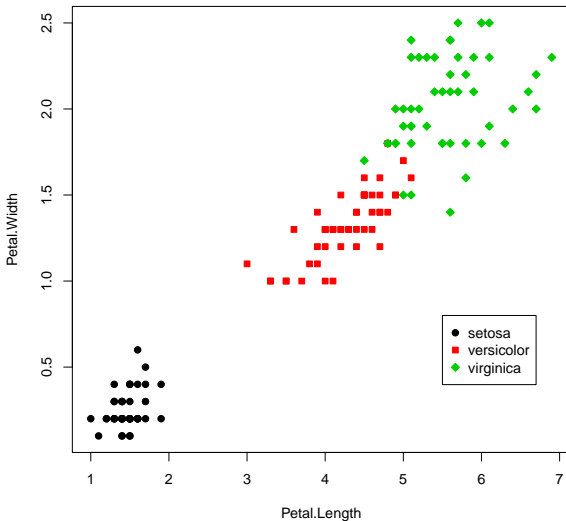
Three Variables

Two Quantitative Variables and One Qualitative Variable



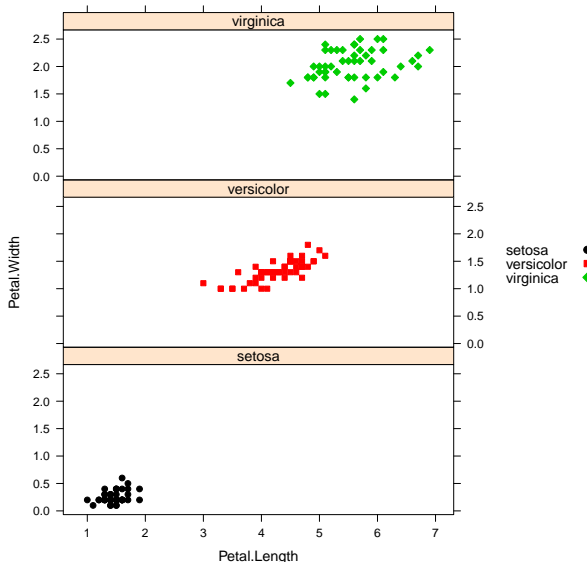
Three Variables

Two Quantitative Variables and One Qualitative Variable



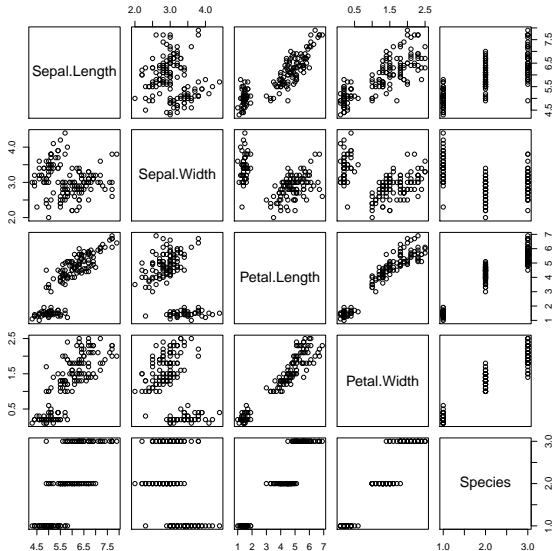
Three Variables

Two Quantitative Variables and One Qualitative Variable

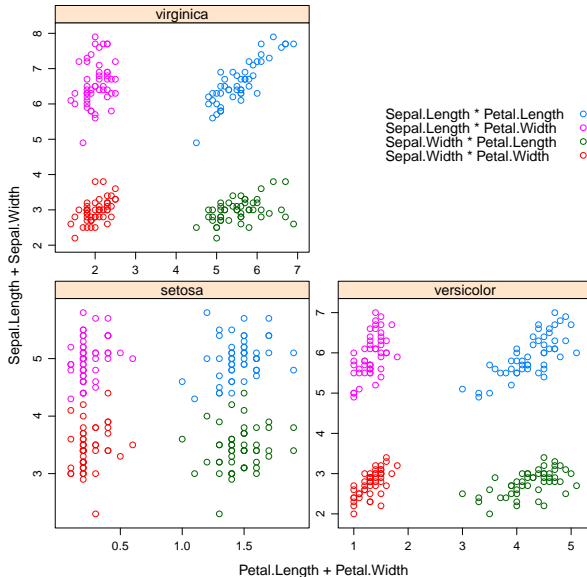


More Than Three Variables

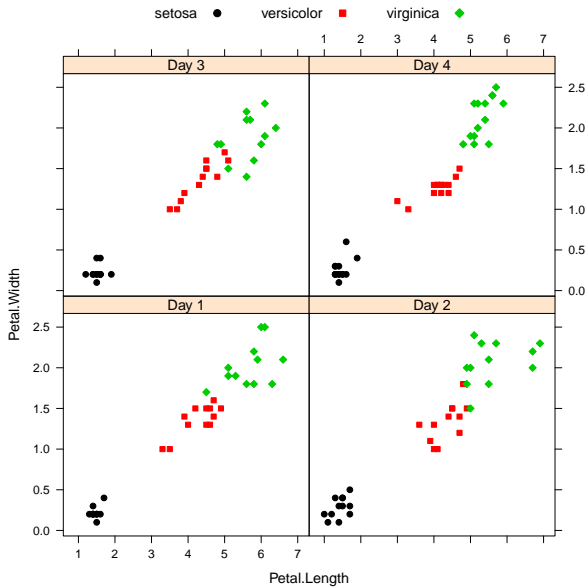
Pairs Plot



More Than Three Variables



More Than Three Variables



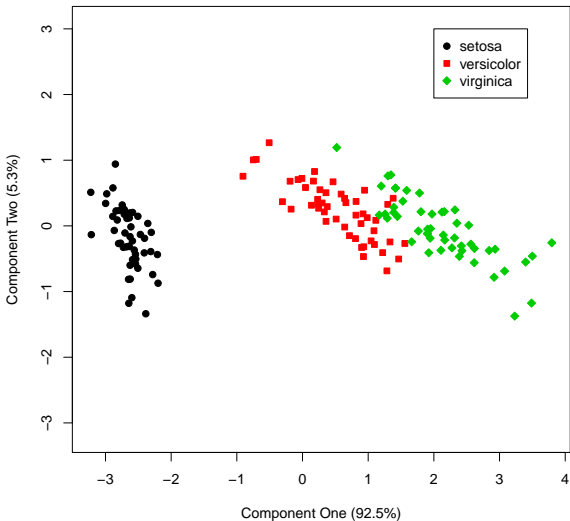
- Suppose a dataset has n observations (rows) and p variables (columns).
- The n observations lie in a p dimensional space.
- To visualise the data, we apply some mathematical procedures to reduce the dimensionality of the data cloud, optimally into two dimensions.
 - 1 Throw away some/lots of variances/information, TANSTAAFL.
 - 2 Plot the reduced data into a two/three dimensional graphs, and visualise the MAIN components of variance.



Multivariate

Principal Component Analysis (PCA)

97.8% of the total variance is explained
in this two dimensional graph





Kai.

<http://www.stat.auckland.ac.nz/~kxio001>.



Paul Murrell.

R Graphics.

Chapman & Hall/CRC, Boca Raton, FL, 2005.

