# Descriptives & Graphing

Image source: http://commons.wikimedia.org/wiki/File:3D_Bar_Graph_Meeting.jpg

## Lecture 3
Survey Research & Design in Psychology
James Neill, 2018
Creative Commons Attribution 4.0

---

# Getting to know data
(how to approach data)

4

---

## Overview:
## Descriptives & Graphing

1 Getting to know data
2 LOM & types of statistics
3 Descriptive statistics
4 Normal distribution
5 Non-normal distributions
6 Effect of skew on central tendency
7 Principles of graphing
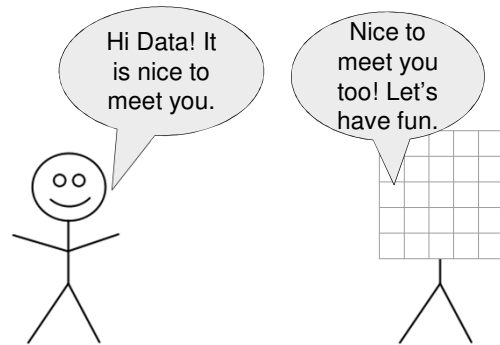8 Univariate graphical techniques

2

---

# Getting to know data



Hi Data! It is nice to meet you.

Nice to meet you too! Let's have fun.

Image source: https://commons.wikimedia.org/wiki/File:Stick_Figure.svg
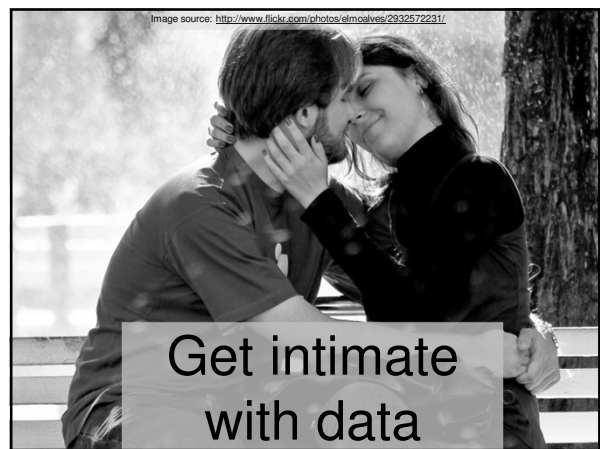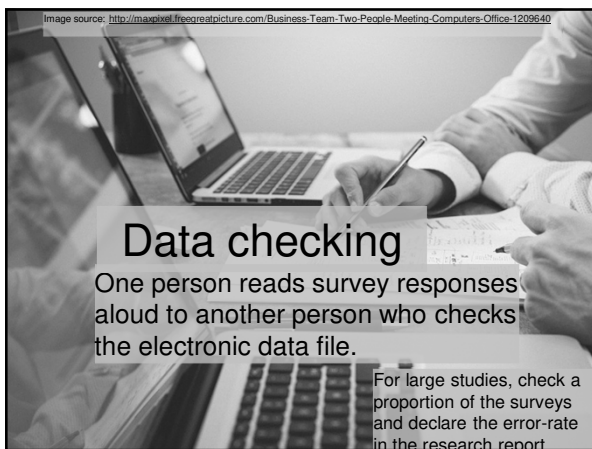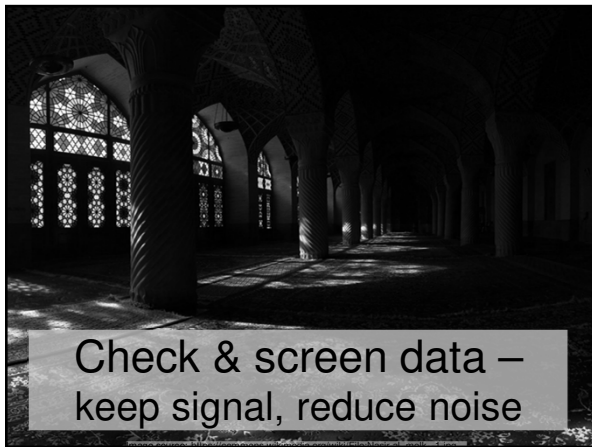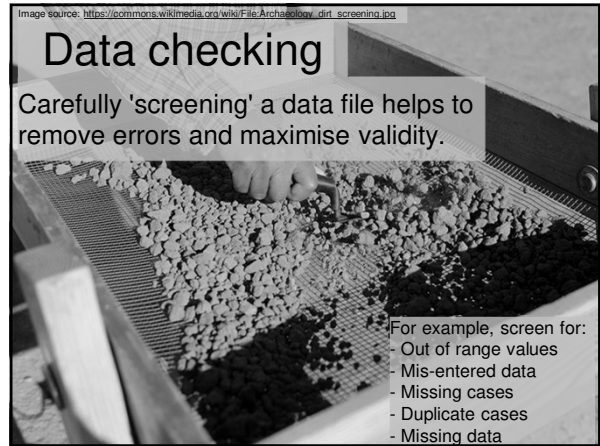
5

---

## Readings

Howitt & Cramer (2014):
- Chapter 01 - Why statistics?
- Chapter 02 - Some basics: Variability and measurement
- Chapter 03 - Describing variables: Tables and diagrams
- Chapter 04 - Describing variables numerically: Averages, variation and spread
- Chapter 05 - Shapes of distributions of scores
- Chapter 06 - Standard deviation and $z$-scores: The standard unit of measurement in statistics

3

---



Play with data –
get to know it.

Image source: http://www.flickr.com/photos/analytik/1356366068/

Don't be afraid - you can't break data!

Image source: http://www.flickr.com/photos/rddave/5094020069

Image source: https://commons.wikimedia.org/wiki/File:Archaeology_dirt_screening.jpg

## Data checking

Carefully 'screening' a data file helps to remove errors and maximise validity.

For example, screen for:
- Out of range values
- Mis-entered data
- Missing cases
- Duplicate cases
- Missing data

Check & screen data –
keep signal, reduce noise

Explore data

Image source: http://maxpixel.freegreatpicture.com/Business-Team-Two-People-Meeting-Computers-Office-1209640

## Data checking

One person reads survey responses aloud to another person who checks the electronic data file.

For large studies, check a proportion of the surveys and declare the error-rate in the research report

Image source: http://www.flickr.com/photos/elmoalves/2932572231/

Get intimate
with data

## Describe data's main features

find a meaningful, accurate way to depict the true story' of the data

Image source: http://www.flickr.com/photos/lloydm/2429991235/

**16**

---

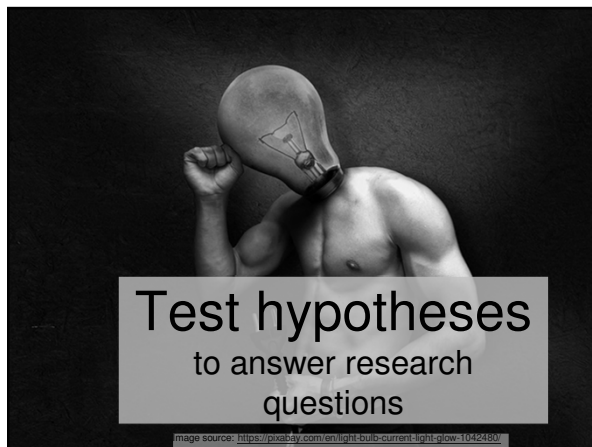### LOM → statistics

Level of measurement determines the type of statistics that can be used, including types of:
• descriptive statistics
• graphs
• inferential statistics

**17**

---

## Test hypotheses
to answer research questions

Image source: https://pixabay.com/en/light-bulb-current-light-glow-1042480/

---

### LOM - Parametric vs. non-parametric

Categorical & ordinal data DV
→ *non-parametric*
**(**Does not assume a normal distribution)

Interval & ratio data DV
→ *parametric*
(Assumes a normal distribution)
→ *non-parametric*
(If distribution is non-normal)

DVs = dependent variables

**15**

---

# Level of measurement & types of statistics

Image source: http://www.flickr.com/photos/peanutlen/2228077524/

**18**

---

### Parametric statistics

• Statistics which estimate **parameters** of a population, based on the **normal distribution**
  – **Univariate**:
    mean, standard deviation, skewness, kurtosis
  – **Bivariate:**
    correlation, linear regression, *t*-tests
  – **Multivariate**:
    multiple linear regression, ANOVAs

## Parametric statistics

- More powerful
  (more sensitive)
- More assumptions
  (population is normally distributed)
- Vulnerable to violations of
  assumptions
  (less robust)

**19**

## Summary: LOM & statistics

- If a normal distribution can be
  assumed, use parametric
  statistics (more powerful)
- If not, use non-parametric
  statistics (less power, but less
  sensitive to violations of
  assumptions)

**22**

## Non-parametric statistics

- Statistics which do not assume
  sampling from a population which
  is **normally distributed**
  – There are non-parametric alternatives for
    many parametric statistics
  – e.g., sign test, chi-squared, Mann-
    Whitney U test, Wilcoxon matched-pairs
    signed-ranks test.

**20**

# Univariate
# descriptive
# statistics

**23**

## Non-parametric statistics

- Less powerful
  (less sensitive)
- Fewer assumptions
  (do not assume a normal distribution)
- Less vulnerable to assumption
  violation
  (more robust)

**21**

## Number of variables

**Univariate**
= one variable

mean, median, mode,
histogram, bar chart

**Bivariate**
= two variables

correlation, *t*-test,
scatterplot, clustered bar
chart

**Multivariate**
= more than two variables

reliability analysis, factor
analysis, multiple linear
regression

**24**

## What to describe?

- **Central tendency**(ies): e.g., frequencies, mode, median, mean
- **Distribution**:
  - **Spread** (dispersion): min., max., range, IQR, percentiles, variance, standard deviation
  - **Shape**: e.g., skewness, kurtosi

25

## Distribution

- Measures of shape, spread, dispersion, and deviation from the central tendency

**Non-parametric:**
- Min. and max.
- Range
- Percentiles

**Parametric:**
- *SD*
- Skewness
- Kurtosis

28

## Central tendency

Statistics which represent the "centre" of a frequency distribution:
- Mode (most frequent)
- Median (50$^{th}$ percentile)
- Mean (average)

Which ones to use depends on:
- Type of data (level of measurement)
- Shape of distribution (esp. skewness)

Reporting more than one may be appropriate.

26

## Distribution

| | Min / Max, Range | Percentile | *Var / SD* |
|---|---|---|---|
| Nominal | *x* | *x* | *x* |
| Ordinal | √ | If meaningful | *x* |
| Interval | √ | √ | √ |
| Ratio | √ | √ | √ |

29

## Central tendency

| | Mode / Freq. /%s | Median | Mean |
|---|---|---|---|
| Nominal | √ | *x* | *x* |
| Ordinal | √ | If meaningful | *x* |
| Interval | √ | √ | √ |
| Ratio | If meaningful | √ | √ |

27

## Descripives for nominal data

- **Nominal LOM** = Labelled categories
- Descriptive statistics:
  - Most frequent? (Mode – e.g., females)
  - Least frequent? (e.g., Males)
  - Frequencies (e.g., 20 females, 10 males)
  - Percentages (e.g. 67% females, 33% males)
  - Cumulative percentages
  - Ratios (e.g., twice as many females as males)

30

## Descripives for ordinal data

- **Ordinal LOM** = Conveys order but not distance (e.g., ranks)
- Descriptives approach is as for nominal (frequencies, mode etc.)
- Plus percentiles (including median) may be useful

31

## Mode (*Mo*)

- Most common score - highest point in a frequency distribution – a real score – the most common response
- Suitable for all levels of data, but may not be appropriate for ratio (continuous)
- Not affected by outliers
- Check frequencies and bar graph to see whether it is an accurate and useful statistic

34

## Descripives for interval data

- **Interval LOM** = order and distance, but no true 0 (0 is arbitrary).
- Central tendency (mode, median, mean)
- Shape/Spread (min., max., range, *SD*, skewness, kurtosis)

Interval data is discrete, but is often treated as ratio/continuous (especially for > 5 intervals)

## Frequencies (*f*) and percentages (%)

- # of responses in each category
- % of responses in each category
- Frequency table
- Visualise using a bar or pie chart

35

## Descriptives for ratio data

- **Ratio** = Numbers convey order and distance, meaningful 0 point
- As for interval, use median, mean, *SD*, skewness etc.
- Can also use ratios (e.g., Group A is twice as "large" as Group B)

33

## Median (*Mdn*)

- Mid-point of distribution (Quartile 2, 50th percentile)
- Not badly affected by outliers
- May not represent the central tendency in skewed data
- If Median is useful, other percentiles may also be worth reporting

36

## Summary: Descriptive statistics

- **Level of measurement** and **normality** determines whether data can be treated as **parametric**
- Describe the **central tendency**
  - Frequencies, Percentages
  - Mode, Median, Mean
- Describe the **distribution**:
  - Min., Max., Range, Quartiles
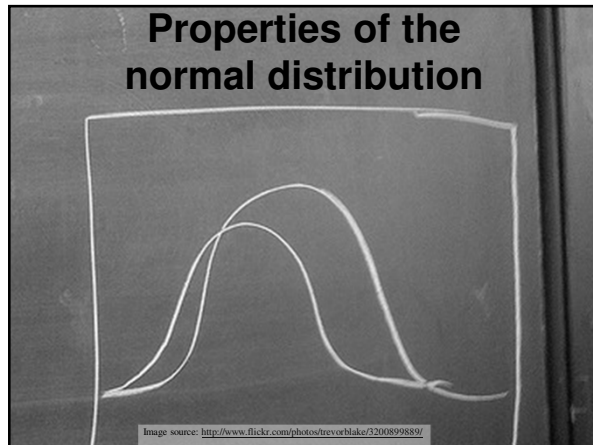  - Standard Deviation, Variance

37

## Four moments of a normal distribution

Four mathematical qualities (parameters) can describe a continuous distribution which at least roughly follows a bell curve shape:

- $1^{st}$ = mean (central tendency)
- $2^{nd}$ = *SD* (dispersion)
- $3^{rd}$ = skewness (lean / tail)
- $4^{th}$ = kurtosis (peakedness / flattness)

40

## Properties of the normal distribution

Image source: http://www.flickr.com/photos/trevorblake/3200899889/

## Mean (1st moment)

- Average score
$$Mean = \Sigma X / N$$
- Use for normally distributed ratio or interval (if treating as continuous) data.
- Influenced by extreme scores (outliers)

41

## Four moments of a normal distribution

Mean

←SD→

kew         +ve Skew

Kurtosis→

## Beware inappropriate averaging

With your head in an oven and your feet in ice you would feel, **on average**, just fine

The majority of people have more than the average number of legs ($M = 1.999$).

42

## Standard deviation (2nd moment)

- $SD$ = square root of the variance

$$= \frac{\Sigma \, (X - \overline{X})^2}{N - 1}$$

- Use for normally distributed interval or ratio data
- Affected by outliers
- Can also derive Standard Error (SE) = $SD$ / square root of $N$

43

## Kurtosis (4th moment)

- Flatness vs. peakedness of distribution:

  +ve = peaked
  -ve = flattened
- Altering the X &/or Y axis can artificially make a distribution look more peaked or flat – add a normal curve to help judge kurtosis visually.

46

## Skewness (3rd moment)

- Lean of distribution
  - +ve = tail to right
  - -ve = tail to left
- Skew be caused by an outlier, or ceiling or floor effects
- Skew be accurate

  (e.g., cars owned per person would have a skewed distribution)

44

## Kurtosis (4th moment)



Image source: https://classconnection.s3.amazonaws.com/65/flashcards/2185065/jpg/kurtosis-142C1127AF2178FB244.jpg

## Skewness (3rd moment)
### (with ceiling and floor effects)
Image source http://www.visualstatistics.net/Visual%20Statistics%20Multimedia/normalization.htm



- Negative skew
- Ceiling effect

- Positive skew
- Floor effect

## Severity of skewness and kurtosis

- View histogram with normal curve
- Deal with outliers
- Rule of thumb:

  Skewness and kurtosis > -1 or < 1 is generally considered to sufficiently normal for meeting the assumptions of parametric inferential statistics
- Significance tests of skewness:

  Tend to be overly sensitive
  (therefore avoid using)

48

## Areas under the normal curve
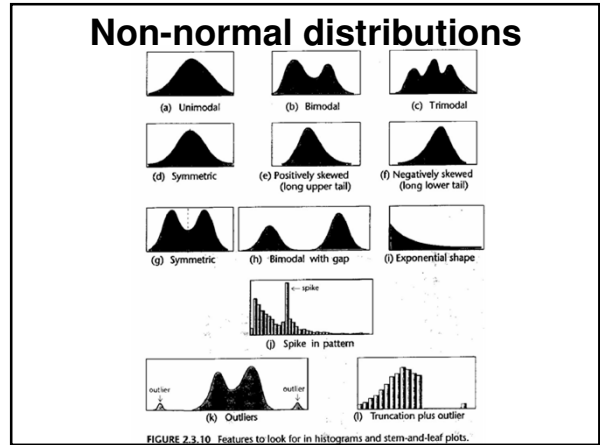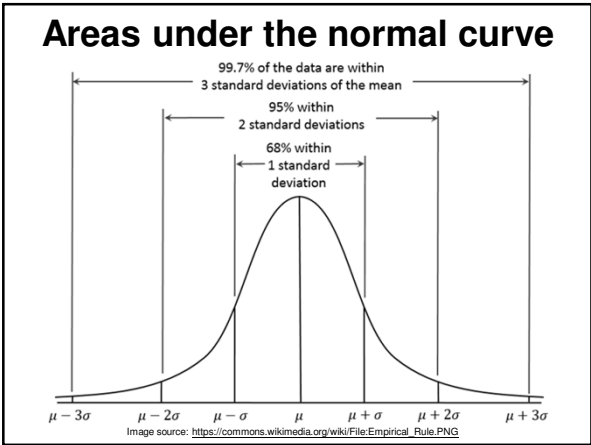
If distribution is normal
(bell-shaped):
~68% of scores within +/- 1 *SD* of *M*
~95% of scores within +/- 2 *SD* of *M*
~99.7% of scores within +/- 3 *SD* of *M*
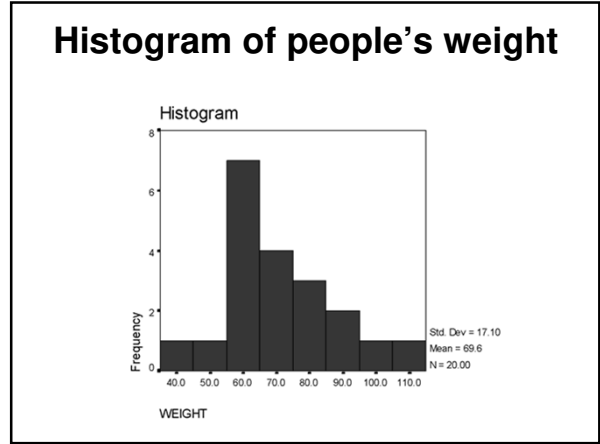
49

## Non-normal distributions

• Modality
  – Uni-modal (one peak)
  – Bi-modal (two peaks)
  – Multi-modal (more than two peaks)
• Skewness
  – Positive (tail to right)
  – Negative (tail to left)
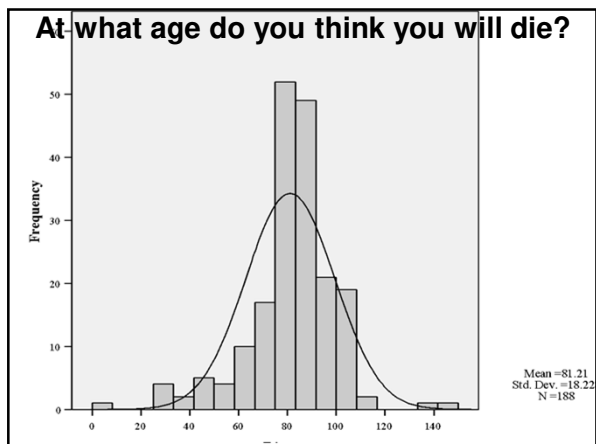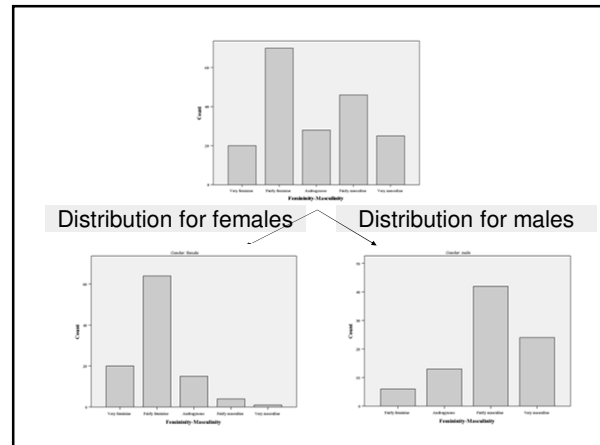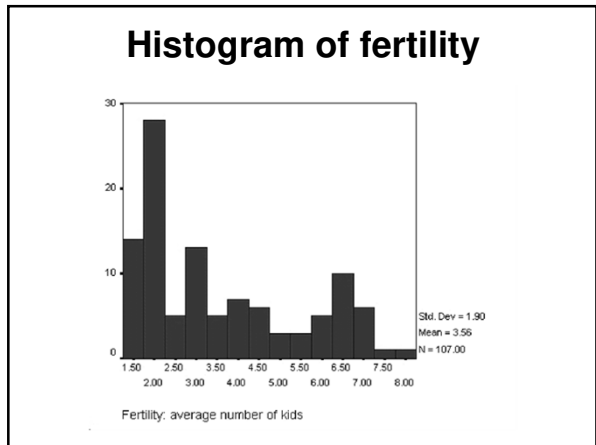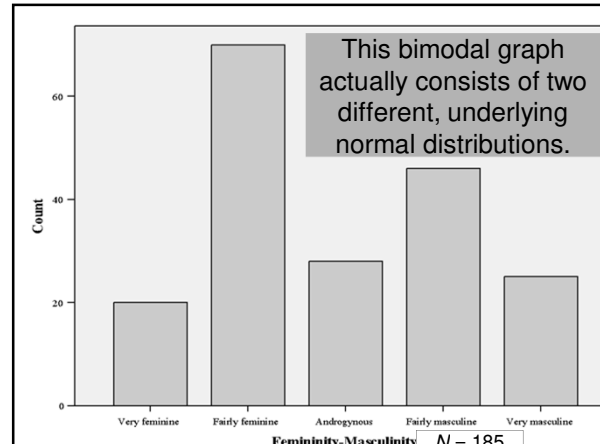• Kurtosis
  – Platykurtic (Flat)
  – Leptokurtic (Peaked)

52

## Areas under the normal curve



Image source: https://commons.wikimedia.org/wiki/File:Empirical_Rule.PNG

## Non-normal distributions



FIGURE 2.3.10   Features to look for in histograms and stem-and-leaf plots.

# Non-normal distributions

51

## Histogram of people's weight



9

## Histogram of daily calorie intake



N = 75



This bimodal graph actually consists of two different, underlying normal distributions.

Femininity-Masculinity    N = 185

## Histogram of fertility



Std. Dev = 1.90
Mean = 3.56
N = 107.00

Fertility: average number of kids



Distribution for females    Distribution for males

## At what age do you think you will die?



Mean =81.21
Std. Dev.=18.22
N =188

## Non-normal distribution:
### Use non-parametric descriptive statistics

- Min. & Max.
- Range = Max. - Min.
- Percentiles
- Quartiles
  - Q1
  - Median (Q2)
  - Q3
  - IQR (Q3-Q1)

60

## Effects of skew on measures of central tendency

**+vely skewed distributions**
　　mode < median < mean
**symmetrical (normal) distributions**
　　mean = median = mode
**-vely skewed distributions**
　　mean < median < mode

61

## Review questions

1. If a survey question produces a "floor effect", where will the mean, median and mode lie in relation to one another?

64

## Effects of skew on measures of central tendency



positive skew

## Review questions

2. Would the mean # of cars owned in Australia exceed the median?

65

## Transformations

- Converts data using various formulae to achieve normality and allow more powerful tests
- Loses original metric
- Complicates interpretation

63

## Review questions

3. Would the mean score on an easy test exceed the median performance?

66

**Graphical techniques**

Image source: http://www.flickr.com/photos/pagedooley/2121472112/

---

# Is Pivot a turning point for web exploration?
**(Gary Flake)**

Image source:http://commons.wikimedia.org/wiki/File:Parodyfilm.png

## (TED talk - 6 min.)

70

---

**Visualisation**

Image source: http://en.wikipedia.org/wiki/File:FAE_visualization.jpg

"Visualization is any technique
for creating images, diagrams, or
animations to communicate a message."
- Wikipedia

---

**Principles of graphing**

• Clear purpose
• Maximise clarity
• Minimise clutter
• Allow visual comparison

71

---

# Science is beautiful
**(Nature Video)**

Image source::http://commons.wikimedia.org/wiki/File:Parodyfilm.png

## (Youtube – 5:30 mins)

69

---

**Graphs (Tufte)**

• Visualise data
• Reveal data
  – Describe
  – Explore
  – Tabulate
  – Decorate
• Communicate complex ideas
  with clarity, precision, and
  efficiency

72

## Graphing steps

1 Identify purpose of the graph (make large amounts of data coherent; present many #s in small space; encourage the eye to make comparisons)
2 Select type of graph to use
3 Draw and modify graph to be clear, non-distorting, and well-labelled (maximise clarity, minimise clarity; show the data; avoid distortion; reveal data at several levels/layers)

**73**

## Cleveland's hierarchy

Image source: https://priceonomics.com/how-william-cleveland-turned-data-visualization/



## Graphing software

**1 Statistical packages**
  - e.g., SPSS Graphs or via Analyses
**2 Spreadsheet packages**
  - e.g., MS Excel
**3 Word-processors**
  - e.g., MS Word – Insert – Object – Micrograph Graph Chart

**74**

## Univariate graphs

- Bar graph } Non-parametric i.e., nominal or ordinal
- Pie chart
- Histogram
- Stem & leaf plot
- Data plot / Error bar } Parametric i.e., normally distributed interval or ratio
- Box plot

**77**

## Cleveland's hierarchy

Image source: http://www.processtrends.com/TOC_data_visualization.htm



Position along a common scale
Position along nonaligned scales
Length
Angle/ slope
Area
Volume
Color

**Worst**

*Based on graphic (Figure 2) in Presentation Graphics (white paper) by Leland Wilkinson, SPSS, Inc and Northwestern Uiv.*

## Bar chart

- Allows comparison of heights of bars
- X-axis: Collapse if too many categories
- Y-axis: Count/Frequency or % - truncation exaggerates differences
- Can add data labels (data values for each bar)



Note truncated Y-axis

## Pie chart

- Use a bar chart instead
- Hard to read
  - Difficult to show
    - Small values
    - Small differences
  - Rotation of chart and position of slices influences perception
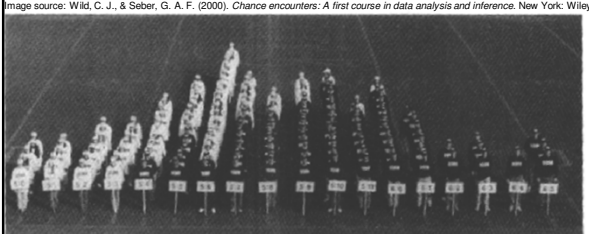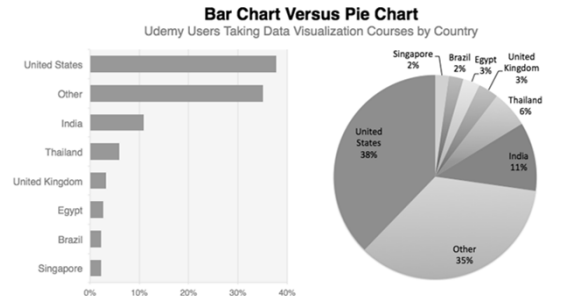
**79**

## Histogram of male and female heights

FIGURE 2.3.11 Histogram of heights constructed using the people. Photograph by Peter Morenus in conjunction with Prof. Linda Strausberg, University of Connecticut. Subjects are University of Connecticut genetics students, females in white tops, males in dark tops.

Wild & Seber (2000)

## Pie chart → Use bar chart instead



Image source: https://priceonomics.com/how-william-cleveland-turned-data-visualization/

**80**

## Stem and leaf plots

- Use for ordinal, interval and ratio data (if rounded)
- May look confusing to unfamiliar reader



## Histogram

- For continuous data (Likert?, Ratio)
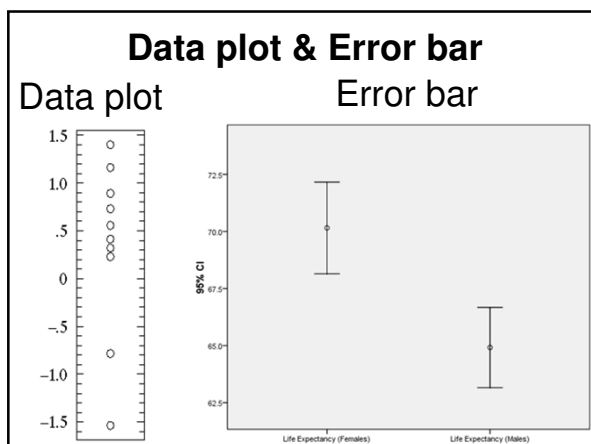- X-axis needs a happy medium for # of categories
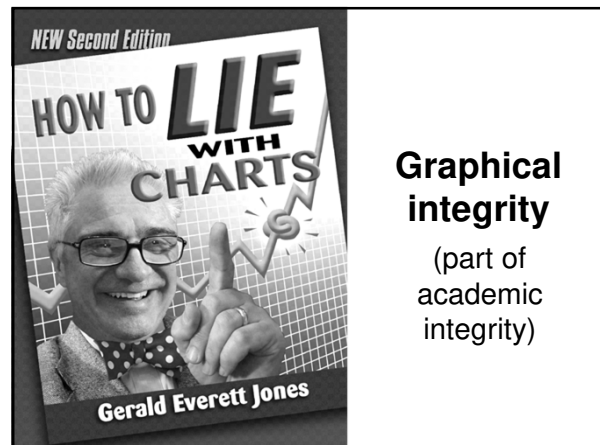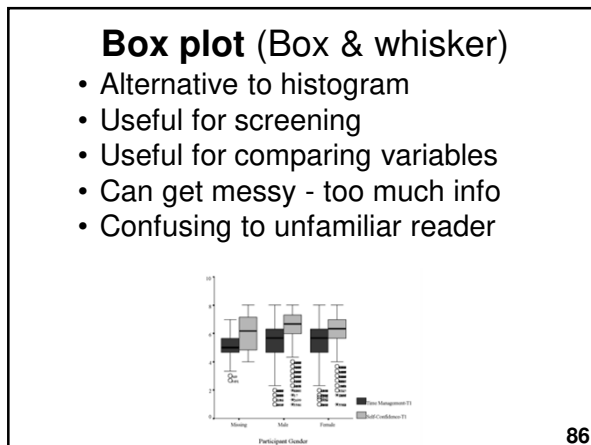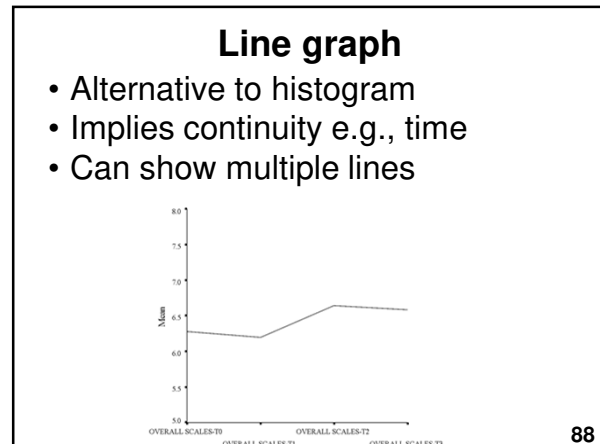- Y-axis matters (can exaggerate)



**81**

## Stem and leaf plot

- Contains actual data
- Collapses tails
- Underused alternative to histogram



**84**

## Box plot
### (Box & whisker)

- Useful for interval and ratio data
- Represents min., max, median, quartiles, & outliers



## Line graph

- Alternative to histogram
- Implies continuity e.g., time
- Can show multiple lines



88

## Box plot (Box & whisker)

- Alternative to histogram
- Useful for screening
- Useful for comparing variables
- Can get messy - too much info
- Confusing to unfamiliar reader



86



## Graphical integrity
(part of academic integrity)

## Data plot & Error bar

Data plot          Error bar



"Like good writing, good graphical displays of data communicate ideas with clarity, precision, and efficiency. Like poor writing, bad graphical displays distort or obscure the data, make it harder to understand or compare, or otherwise thwart the communicative effect which the graph should convey."

Michael Friendly – Gallery of Data          90

## Tufte's graphical integrity

- Some lapses intentional, some not
- Lie Factor = $\frac{\text{size of effect in graph}}{\text{size of effect in data}}$
- Misleading uses of area
- Misleading uses of perspective
- Leaving out important context
- Lack of taste and aesthetics

91

## Next lecture

**Correlation**
- Covariation
- Purpose of correlation
- Linear correlation
- Types of correlation
- Interpreting correlation
- Assumptions / limitations

94

## Review exercise:
### Fill in the cells in this table

| Level | Properties | Examples | Descriptive Statistics | Graphs |
|---|---|---|---|---|
| Nominal /Categorical | | | | |
| Ordinal / Rank | | | | |
| Interval | | | | |
| Ratio | | | | |

Answers: http://goo.gl/Ln9e1

92

## References

1 Chambers, J., Cleveland, B., Kleiner, B., & Tukey, P. (1983). *Graphical methods for data analysis*. Boston, MA: Duxbury Press.
2 Cleveland, W. S. (1985). *The elements of graphing data*. Monterey, CA: Wadsworth.
3 Jones, G. E. (2006). *How to lie with charts*. Santa Monica, CA: LaPuerta.
4 Tufte, E. R. (1983). *The visual display of quantitative information*. Cheshire, CT: Graphics Press.
5 Tufte. E. R. (2001). *Visualizing quantitative data*. Cheshire, CT: Graphics Press.
6 Tukey J. (1977). *Exploratory data analysis*. Addison-Wesley.
7 Wild, C. J., & Seber, G. A. F. (2000). *Chance encounters: A first course in data analysis and inference*. New York: Wiley.

93