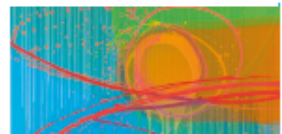# Lecture 01

## Concepts of Statistics
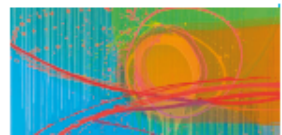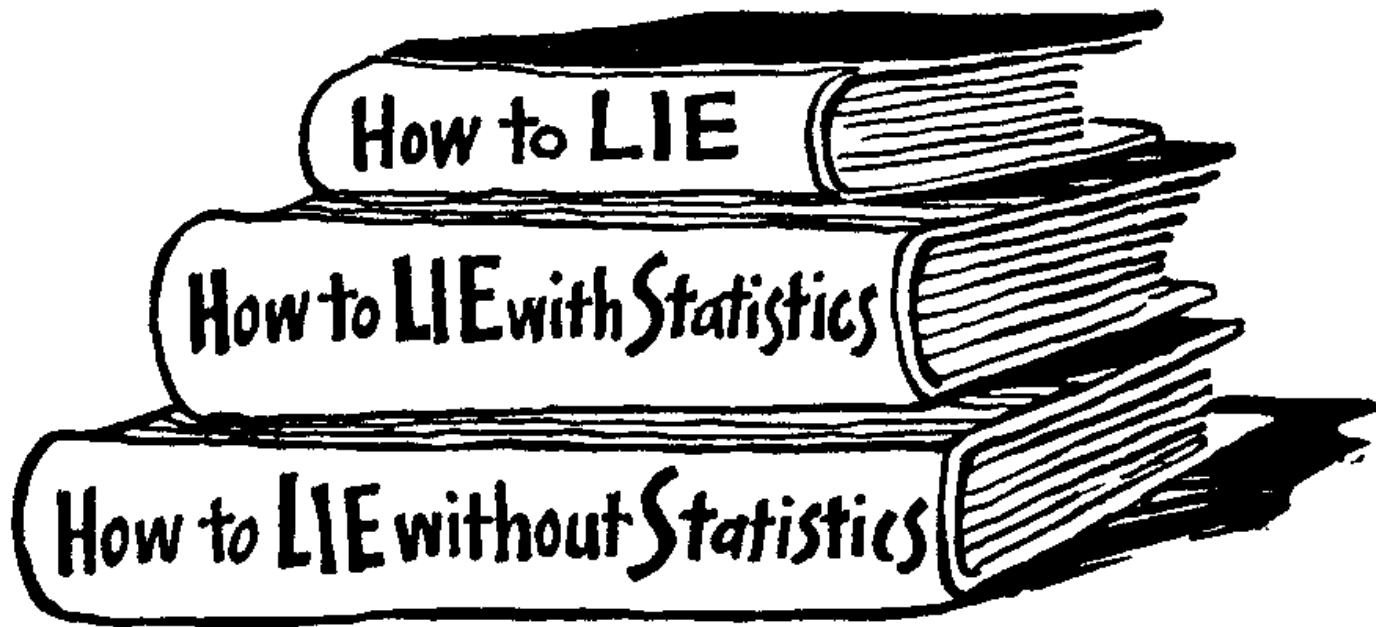
# Statistics=FORMULA?
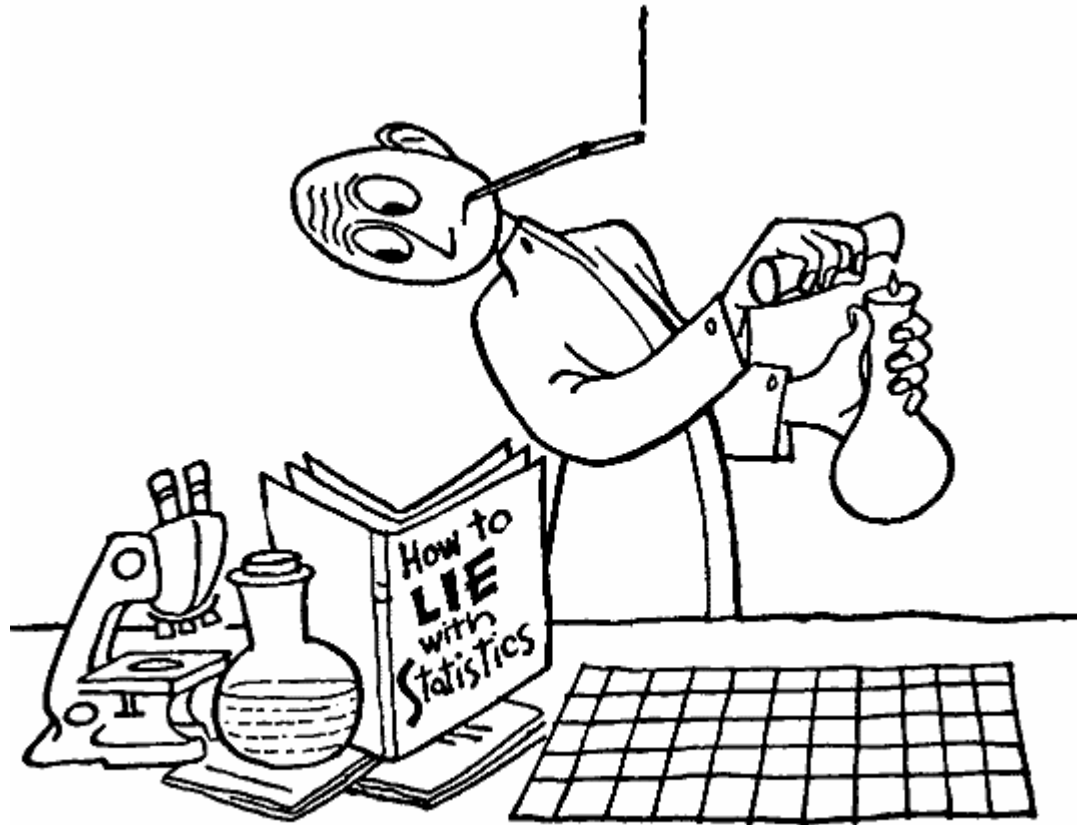
# First LESSON!

# How to LIE with Statistics?

- http://www.physics.csbsju.edu/stats/display.html

ANDY FIELD

# Randomness



Unnatural Selection.

©T. McCracken mchumor.com
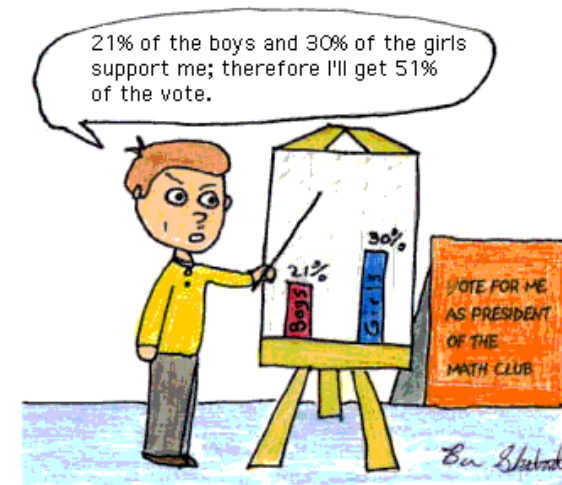
- Wolves with a roulette wheel spin to figure out which sheep they're going to hunt

# Types of Data Analysis

- **Quantitative Methods**
  - Testing theories using numbers
- **Qualitative Methods**
  - Testing theories using language
    - Magazine articles/Interviews
    - Conversations
    - Newspapers
    - Media broadcasts

# The Research Process



FIGURE 1.2
The research process

# Initial Observation

- Find something that needs explaining
  - Observe the real world
  - Read other research
- Test the concept: collect data
  - Collect data to see whether your hunch is correct
  - To do this you need to define variables
    - Anything that can be measured and can differ across entities or time.

ANDY FIELD

# Survey VS Results

MCHUMOR.com by T. McCracken

SURVEY

RESULTS

©T. McCracken
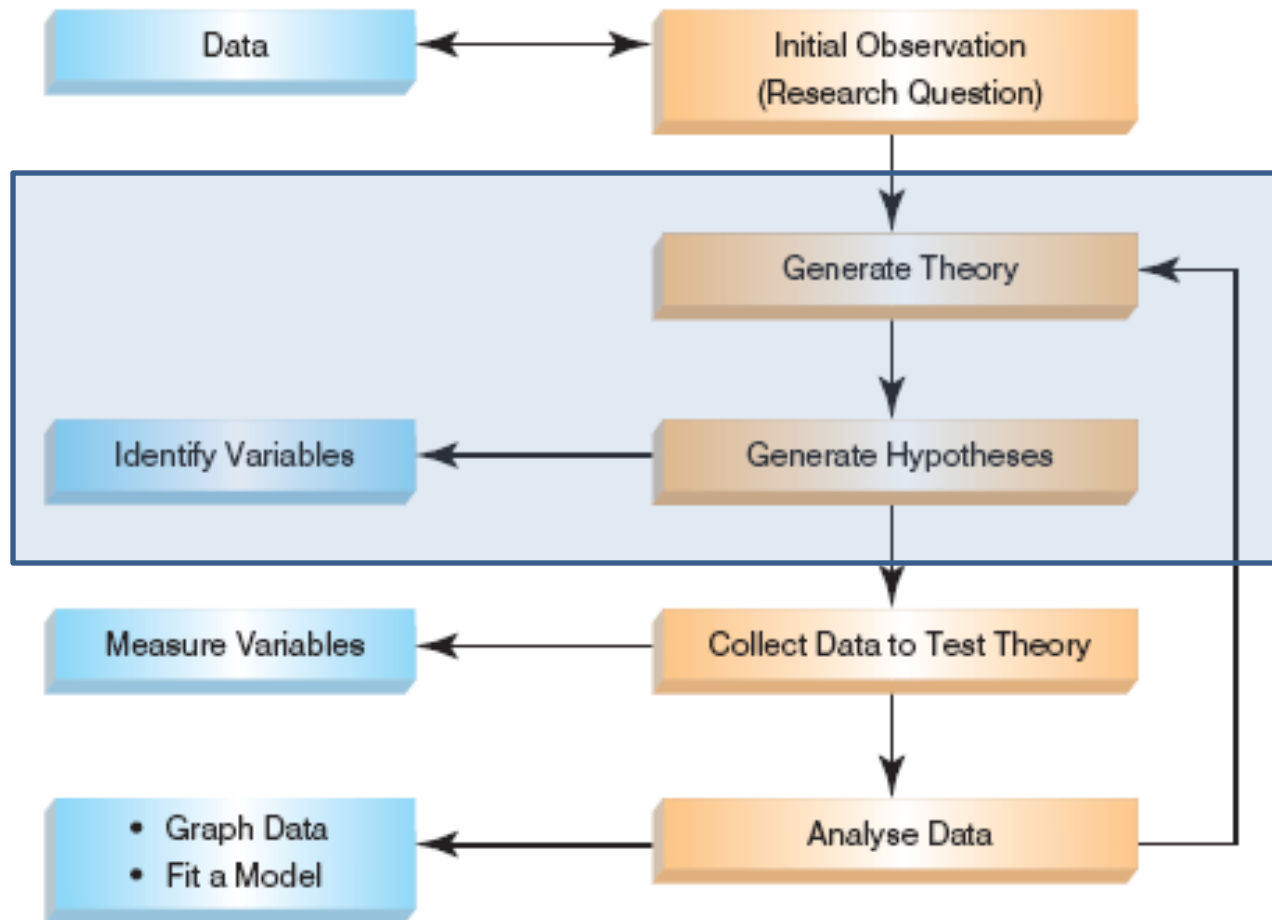www.mchumor.com

"If at first you don't get the survey results
you want on your first survey,
get mad and tear it up."
©T. McCracken mchumor.com

ANDY FIELD

# The Research Process



FIGURE 1.2
The research process

# Generating and Testing Theories

- **Theories**
  - An hypothesized general principle or set of principles that explain known findings about a topic and from which new hypotheses can be generated.

- **Hypothesis**
  - A prediction from a theory.
  - E.g. the number of people turning up for a Big Brother audition that have narcissistic personality disorder will be higher than the general level (1%) in the population.

- **Falsification**
  - The act of disproving a theory or hypothesis.

# Outlier

- An outlying observation, or outlier, is one that appears to deviate markedly from other members of the sample in which it occurs.

# Hypotheses?

# The Research Process



FIGURE 1.2
The research process

DISCOVERING STATISTICS USING SPSS
THIRD EDITION

ANDY FIELD

# Data Collection: What to Measure?



- Hypothesis:
  - ?.
- Independent Variable
  - The proposed cause
  - A predictor variable
  - A manipulated variable (in experiments)
- Dependent Variable
  - The proposed effect
  - An outcome variable
  - Measured not manipulated (in experiments)

# Levels of Measurement

- Categorical (entities are divided into distinct categories):
  - Binary variable: There are only two categories
    - e.g. dead or alive.
  - Nominal variable: There are more than two categories
    - e.g. whether someone is an omnivore, vegetarian, vegan, or fruitarian.
  - Ordinal variable: The same as a nominal variable but the categories have a logical order
    - e.g. whether people got a fail, a pass, a merit or a distinction in their exam.
- Continuous (entities get a distinct score):
  - Interval variable: Equal intervals on the variable represent equal differences in the property being measured
    - e.g. the difference between 6 and 8 is equivalent to the difference between 13 and 15.
  - Ratio variable: The same as an interval variable, but the ratios of scores on the scale must also make sense
    - e.g. a score of 16 on an anxiety scale means that the person is, in reality, twice as anxious as someone scoring 8.
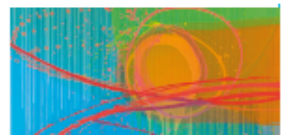
# Measurement Error

- ## Measurement error
  - The discrepancy between the actual value we're trying to measure, and the number we use to represent that value.

- ## Example:
  - You (in reality) weigh 80 kg.
  - You stand on your bathroom scales and they say 83 kg.
  - The measurement error is 3 kg.

ANDY FIELD

# Validity

- Whether an instrument measures what it set out to measure.
- Content validity
  - Evidence that the content of a test corresponds to the content of the construct it was designed to cover
- Ecological validity
  - Evidence that the results of a study, experiment or test can be applied, and allow inferences, to real-world conditions.

# Reliability
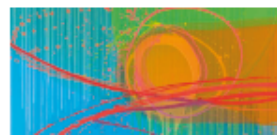
- Reliability
  - The ability of the measure to produce the same results under the same conditions.

- Test-Retest Reliability
  - The ability of a measure to produce consistent results when the same entities are tested at two different points in time.

# Data Collection: How to Measure

- ## Correlational research:
  - Observing what naturally goes on in the world without directly interfering with it.

- ## Cross-sectional research:
  - This term implies that data come from people at different age points with different people representing each age point.

- ## Experimental research:
  - One or more variable is systematically manipulated to see their effect (alone or in combination) on an outcome variable.
  - Statements can be made about cause and effect.

# Experimental Research Methods

- **Cause and Effect (Hume, 1748)**
  1. Cause and effect must occur close together in time (contiguity);
  2. The cause must occur before an effect does;
  3. The effect should never occur without the presence of the cause.
- **Confounding variables: the '*Tertium Quid*'**
  - A variable (that we may or may not have measured) other than the predictor variables that potentially affects an outcome variable.
  - E.g. The relationship between breast implants and suicide is confounded by self esteem.
- **Ruling out confounds (Mill, 1865)**
  - An effect should be present when the cause is present and that when the cause is absent the effect should be absent also.
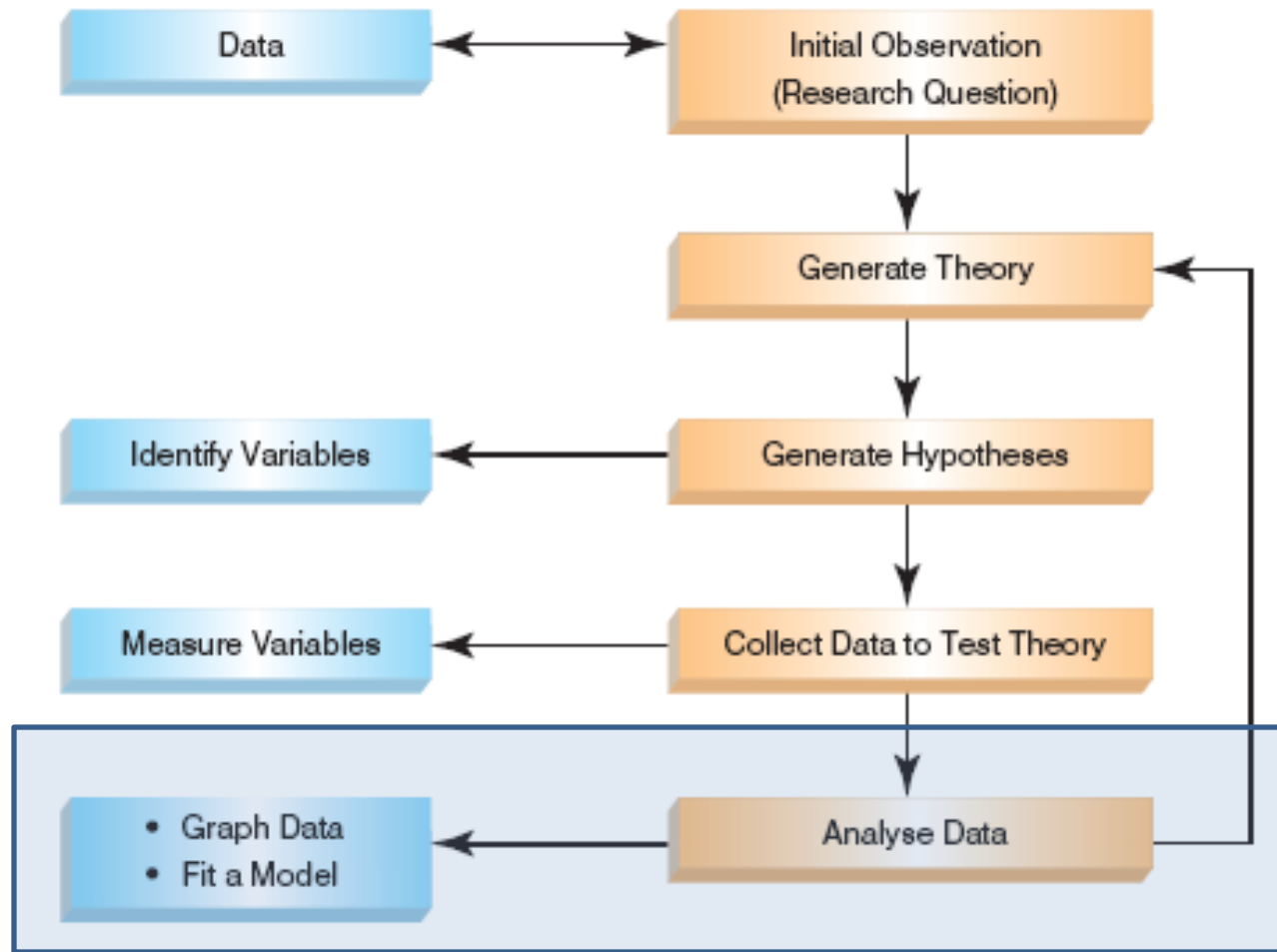  - Control conditions: the cause is absent.

# Methods of Data Collection

- **Between-group/Between-subject/independent**
  - Different entities in experimental conditions
- **Repeated measures (within-subject)**
  - The same entities take part in all experimental conditions.
  - Economical
  - Practice effects
  - Fatigue

# Types of Variation

- **Systematic Variation**
  - Differences in performance created by a specific experimental manipulation.

- **Unsystematic Variation**
  - Differences in performance created by unknown factors.
    - Age, Gender, IQ, Time of day, Measurement error etc.

- **Randomization**
  - Minimizes unsystematic variation.
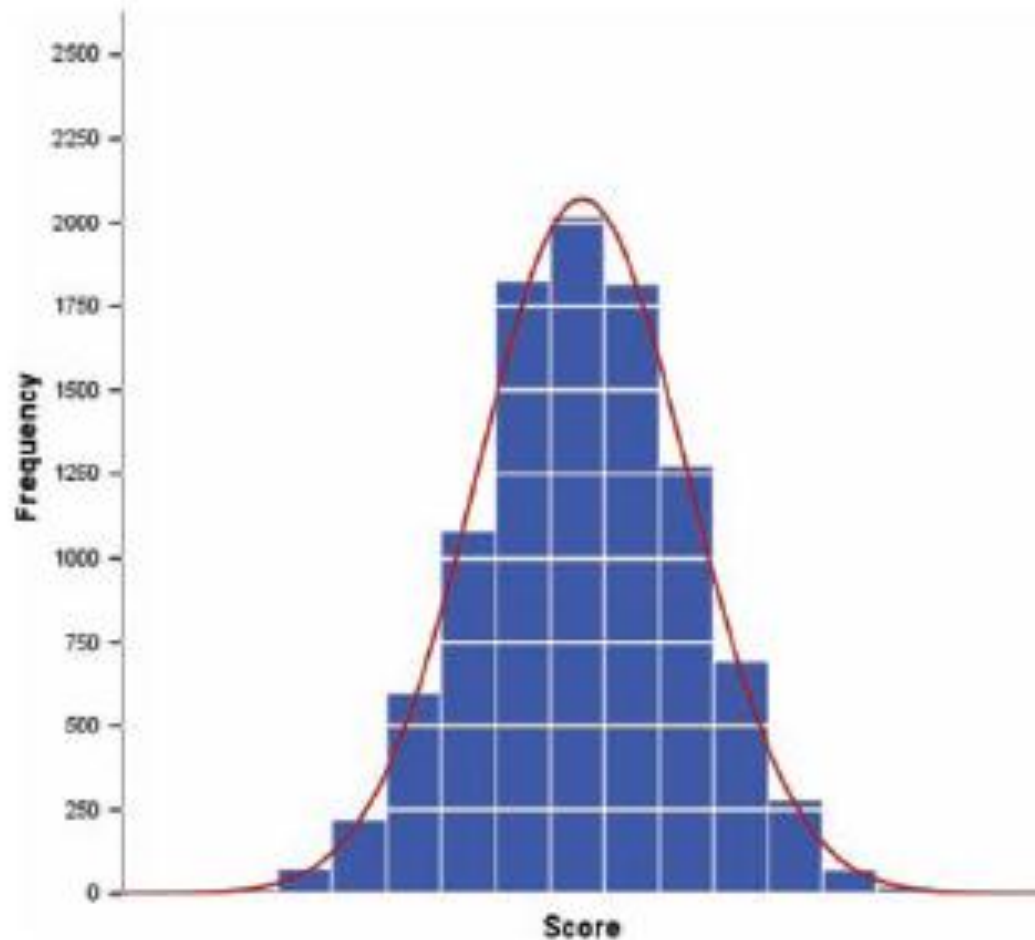
# The Research Process



FIGURE 1.2
The research process

# Analysing Data: Histograms

- ## Frequency Distributions (aka Histograms)
  - A graph plotting values of observations on the horizontal axis, with a bar showing how many times each value occurred in the data set.

- ## The 'Normal' Distribution
  - Bell shaped
  - Symmetrical around the centre

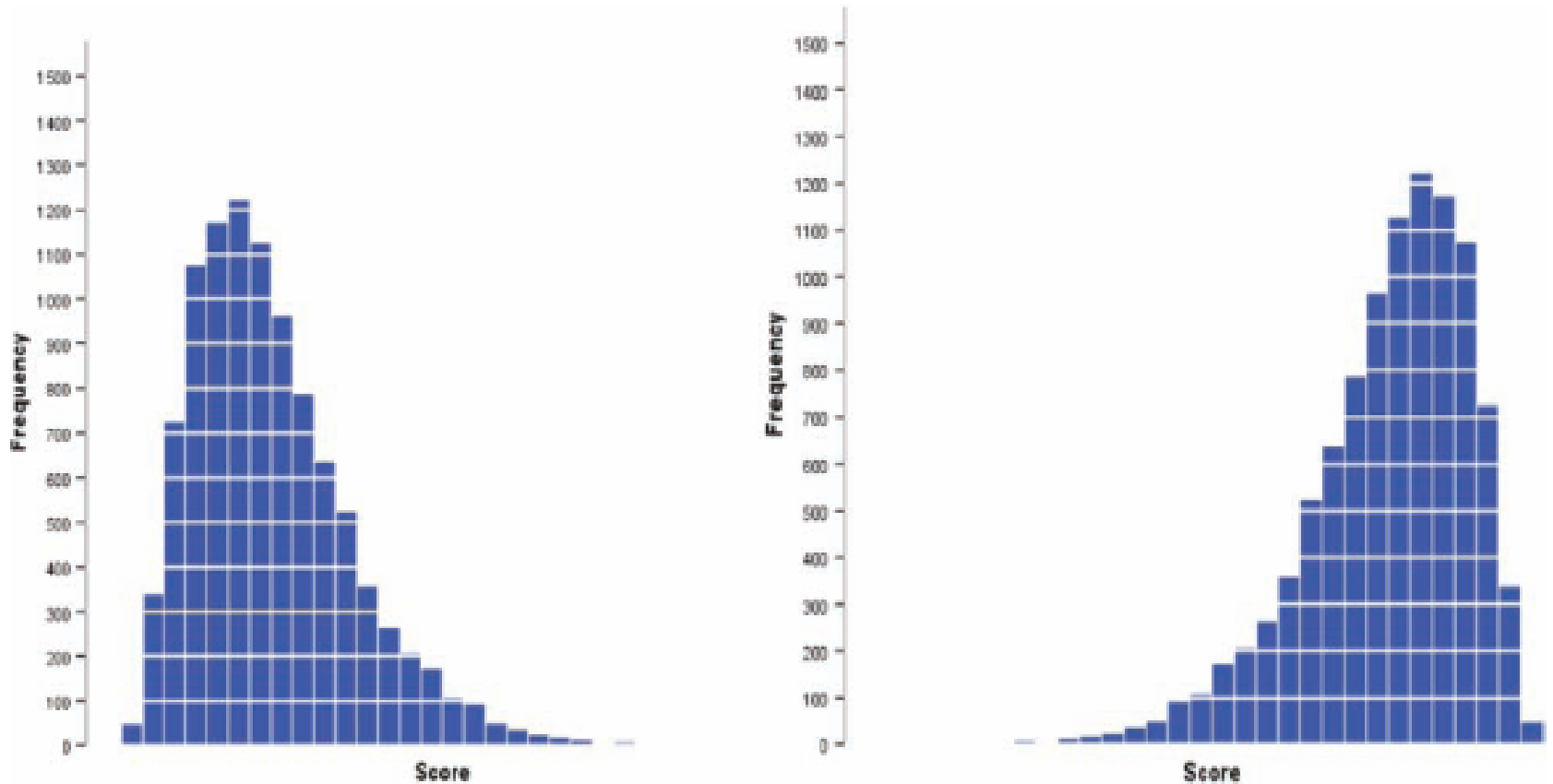# The Normal Distribution



FIGURE 1.3
A 'normal' distribution (the curve shows the idealized shape)
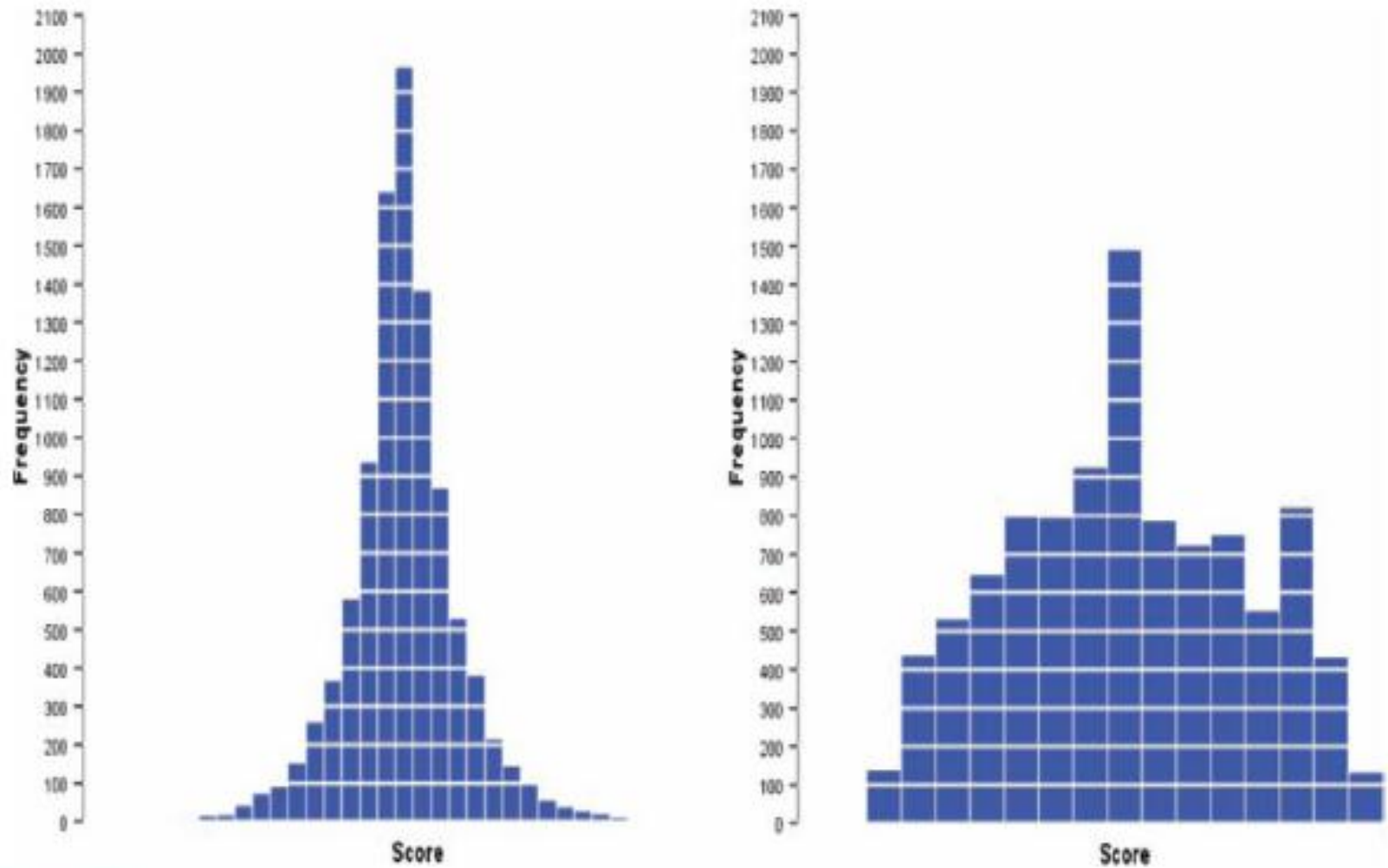
# Properties of Frequency Distributions

- **Skew**
  - The symmetry of the distribution.
  - Positive skew (scores bunched at low values with the tail pointing to high values).
  - Negative skew (scores bunched at high values with the tail pointing to low values).

- **Kurtosis**
  - The 'heaviness' of the tails.
  - Leptokurtic = heavy tails.
  - Platykurtic = light tails.

# Skew



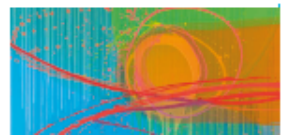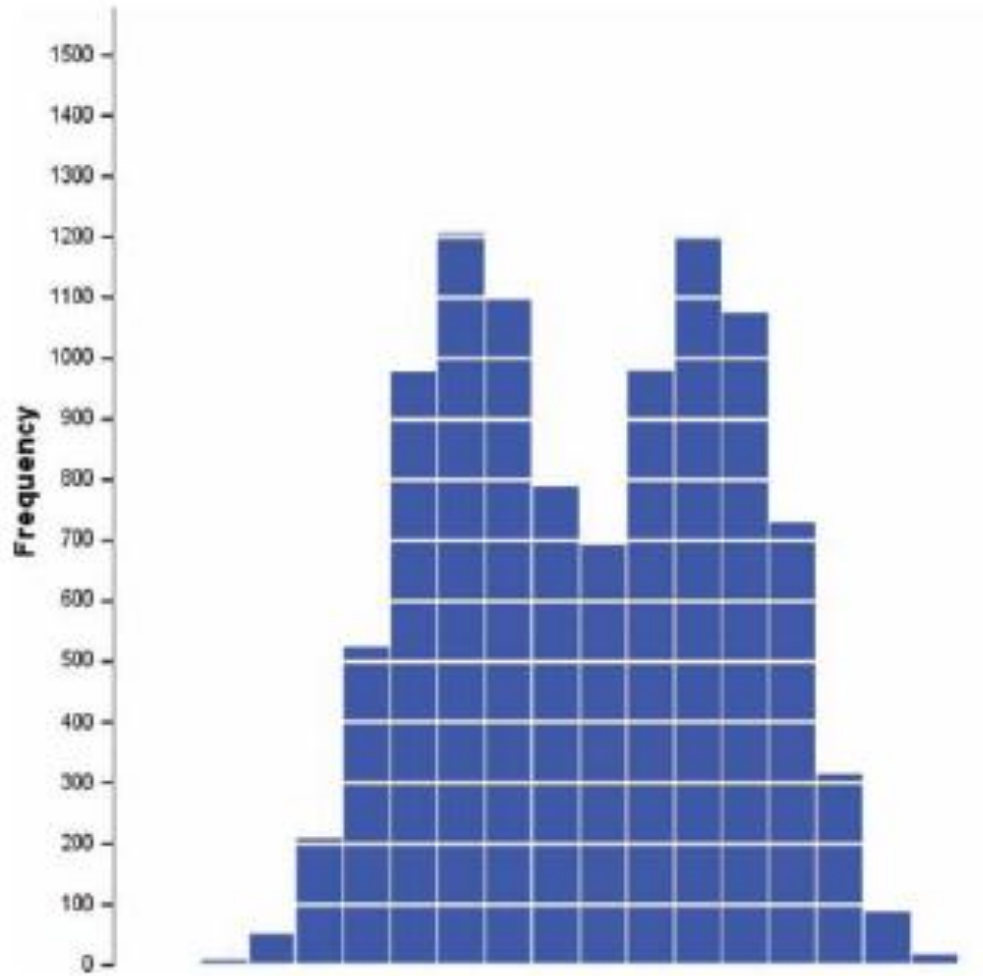FIGURE 1.4 A positively (left figure) and negatively (right figure) skewed distribution

# Kurtosis



FIGURE 1.5 Distributions with positive kurtosis (leptokurtic, left figure) and negative kurtosis (platykurtic, right figure)

# Central tendency: The Mode

- Mode
  - The most frequent score

- Bimodal
  - Having two modes

- Multimodal
  - Having several modes

# A Bimodal Distribution



FIGURE 1.6
A bimodal distribution

# Central Tendency: The Median

- Median
  - The middle score when scores are ordered.
- Example
  - Number of friends of 11 Facebook users.
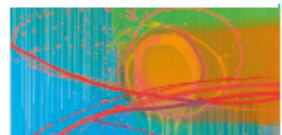
22, 40, 53, 57, 93, 98, 103, 108, 116, 121, 252
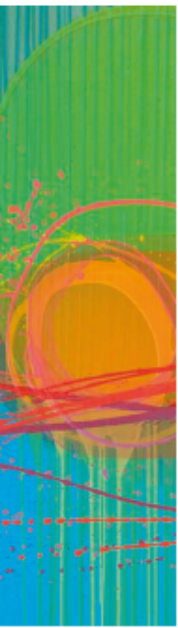
Median

ANDY FIELD

# Central Tendency: The Mean

- Mean
  - The sum of scores divided by the number of scores.
  - Number of friends of 11 Facebook users.

$$\overline{X} = \frac{\sum_{i=1}^{n} x_i}{n}$$

$$\sum_{i=1}^{n} x_i = 22 + 40 + 53 + 57 + 93 + 98 + 103 + 108 + 116 + 121 + 252$$

$$= 1063$$

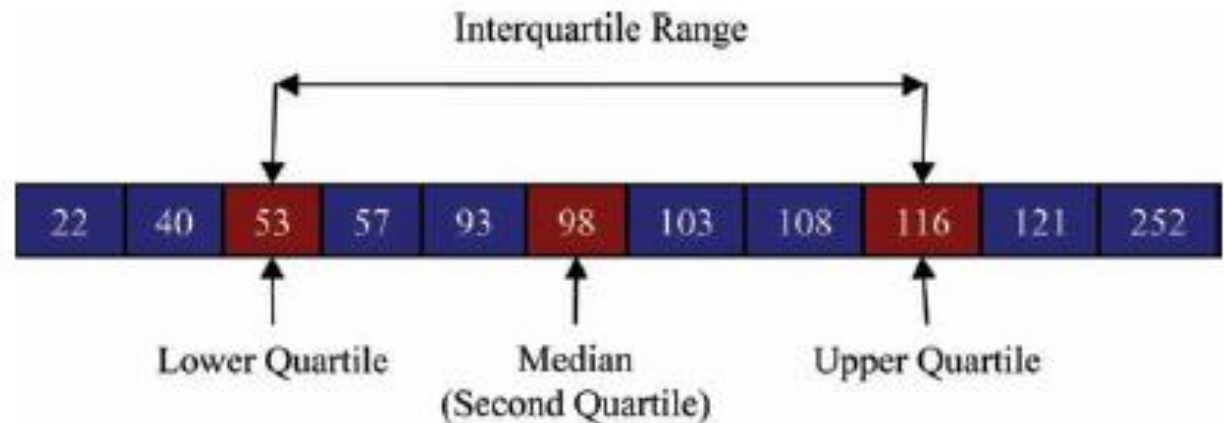$$\overline{X} = \frac{\sum_{i=1}^{n} x_i}{n} = \frac{1063}{11} = 96.64$$

# The Dispersion: Range

- ## The Range
  - The smallest score subtracted from the largest
- ## Example
  - Number of friends of 11 Facebook users.
  - 22, 40, 53, 57, 93, 98, 103, 108, 116, 121, 252
  - Range = 252 − 22 = 230
  - Very biased by outliers

ANDY FIELD

# The Dispersion: The Interquartile range

- Quartiles
  - The three values that split the sorted data into four equal parts.
  - Second Quartile = median.
  - Lower quartile = median of lower half of the data
  - Upper quartile = median of upper half of the data

FIGURE 1.7
Calculating quartiles and the interquartile range



Interquartile Range

| 22 | 40 | 53 | 57 | 93 | 98 | 103 | 108 | 116 | 121 | 252 |

Lower Quartile   Median (Second Quartile)   Upper Quartile

# Going beyond the data: Z-scores

- Z-scores
  - Standardising a score with respect to the other scores in the group.
  - Expresses a score in terms of how many standard deviations it is away from the mean.
  - The distribution of *z*-scores has a mean of 0 and *SD* = 1.

$$z = \frac{X - \overline{X}}{s}$$

ANDY FIELD

# Properties of *z*-scores

- 1.96 cuts off the top 2.5% of the distribution.
- −1.96 cuts off the bottom 2.5% of the distribution.
- As such, 95% of z-scores lie between −1.96 and 1.96.
- 99% of z-scores lie between −2.58 and 2.58,
- 99.9% of them lie between −3.29 and 3.29.

# Types of Hypotheses

- ## Null hypothesis, $H_0$
  - There is no effect.
  - E.g. Big Brother contestants and members of the public will not differ in their scores on personality disorder questionnaires

- ## The alternative hypothesis, $H_1$
  - AKA the experimental hypothesis
  - E.g. Big Brother contestants will score higher on personality disorder questionnaires than members of the public