

# Evaluation Methods and Statistics revision#2

Chris Baber

# Aims of this Module

- This module provides an introduction to the use of empirical, scientific methods, including experimental design and statistics.
- This module is targeted at computer scientists with an interest in:
  - Developing systems that support human activity (Human-Computer Interaction)
  - Building computational models of human behaviour
  - Understanding human behaviour as an inspiration for Robotics, Machine Learning and Artificial Intelligence
  - **Designing and analysing experiments to evaluate system performance**

# Outcomes of this Module

- On successful completion of this module, you will be expected to be able to:
  - identify and apply research methodologies for investigating human behaviour;
  - recognise the appropriateness of statistical techniques in data analysis;
  - conduct and report statistical tests;
  - interpret and critique research findings that are supported by statistical tests;
  - demonstrate understanding of experimental design, including sampling, participant selection, task design and research ethics.

# New Exam Requirements

- You will full details from the School.
- What do I expect...
  - Not all students will have access to R (or Statistics packages), so any calculations should be possible to perform by hand
  - If you are able to use Statistics packages for calculations, I expect enough information in your answer to demonstrate how this answer was produced (don't just write a single number and expect full marks...)

# Doing Experiments

- **Ethics**

- What are the basic principles of the Declaration of Helsinki
- How can you ensure that participants' identity is protected and that they will not suffer from participating in the experiment?
- What are the basic principles of the (UK) Data Protection Act?

- **Ecological Validity**

- How 'true to life' is the activity that you are asking participants to do?
- How 'true to life' is the environment in which these tasks are performed?

- **Experimental Design**

- What is the Hypothesis to be tested?
- What are the Dependent and Independent variables?
- For the Independent variables, which is the 'control' and the 'experimental' condition?
- How can you manage confounding variables in the experiment?

# Experimental Design template

**Hypothesis:** Reaction time to congruent words will be faster than reaction time to incongruent words

**Independent Variable:** Congruent Words (colour of ink = name of word),  
Incongruent Words (colour of ink  $\neq$  name of word)

**Control Condition:**  
Congruent Words

**Experimental Condition:**  
Incongruent Words

**Dependent Variable(s):** Reaction Time

**Task:** participants will be asked to read, as quickly as possible, single words on a display. The words will be the names of colours and will be presented either in the same colour as the word's name or in a different colour

**Confounding Variables:** performance could be affected by ability to perceive colour ('colour-blindedness') and knowledge of the names of colour ('language skills')

# Hypothesis Testing

- Type I error
  - We could accept the Alternative hypothesis when it is false (false positive).
  - Many statistics tests are designed to minimise this error.
  - Type I errors define the significance level ( $\alpha$ ) that the experimenter will accept (conventionally 5%)
- Type II error
  - We could accept the Null hypothesis (fail to reject it) when it is false (false negative).
  - The probability of a Type II error is defined as  $\beta$
  - The probability of correctly rejecting a false null hypothesis is defined as  $1-\beta$ , which called Power.

# Statistics

- To apply a Parametric statistical test, we need to show that the data follow a Normal distribution
  - Applying shapiro-wilk tests (next slides)
- If the data are measured on at least an interval scale and are normally distributed, then you can use the t-statistic to compare means between two groups (for more groups, you need ANOVA; if the data are not normally distributed then you apply non-parametric tests)



# Shapiro-Wilk

- This is to test whether the set of data in an experiment are drawn from a normal distribution.
- The formula is:

$$W = \frac{\left(\sum_{i=1}^n a_i x_{(i)}\right)^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

- You will recognise the denominator as the formula for the Sum of Squares of the data; the numerator defines the expected mean, variance, and covariance of a sample size (n) from a normal distributed sample.
- While the sum of squares is easy to calculate, we use look-up tables for the expected values.
- I will provide a step-by-step tutorial on how to do this...

# Shapiro-Wilk tutorial

- Assume you collected these data:

Condition A	Condition B
65	74
61	35
63	72
86	68
70	45
55	58
M = 66.7	M = 59
sd = 10.7	sd = 15.8

# Shapiro-Wilk tutorial

- Step 1:
  - Combine ALL of the data for the experiment into a single table
  - Order the data in terms of size (smallest to largest)

All Experiment Data
35
45
55
58
61
63
65
68
70
72
74
86

# Shapiro-Wilk tutorial

- Step 2:
  - Calculate sum of squares
    - For each value, calculate difference between that value and the mean of the sample
    - Square this difference
    - Sum the squares

Data	X-M	(x-m) <sup>2</sup>
35	-27.67	765.4444
45	-17.67	312.1111
55	-7.67	58.77778
58	-4.67	21.77778
61	-1.67	2.777778
63	0.33	0.111111
65	2.33	5.444444
68	5.33	28.44444
70	7.33	53.77778
72	9.33	87.11111
74	11.33	128.4444
86	23.33	544.4444
M= 62.67		SS= 2008.667

# Shapiro-Wilk tutorial

- Step 3:
  - Estimate the expected values for the data *if* they were drawn from a normal distribution
  - Define sample size, N. In this case,  $N = 12$ .
  - Use look-up table to define coefficients for a when  $N = 12$ .

[illegible]

# Shapiro-Wilk tutorial

- Step 4:
  - Apply these coefficients to your data.
    - Divide the data into pairs to correspond with the number of coefficients.
    - In this case, there are 6 values for  $a$  (from the table).
    - We divide our data into 6 pairs, subtracting the largest from the smallest values.

	$x_l - x_s$	$a$	$a * (x_l - x_s)$
$x_{12} - x_1$	51	0.5475	27.9225
$x_{11} - x_2$	29	0.3325	9.6425
$x_{10} - x_3$	17	0.2347	3.9899
$x_9 - x_4$	12	0.1586	1.9032
$x_8 - x_5$	7	0.0922	0.6454
$x_7 - x_6$	2	0.0303	0.0606
			$b = 44.1641$
			$b^2 = 1950.468$

# Shapiro-Wilk tutorial

- Step 5:

- Calculate W

$$W = 1950.468 / 2008.667$$
$$= 0.971026$$

$$W = \frac{\left(\sum_{i=1}^n a_i x_{(i)}\right)^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

# Shapiro-Wilk tutorial

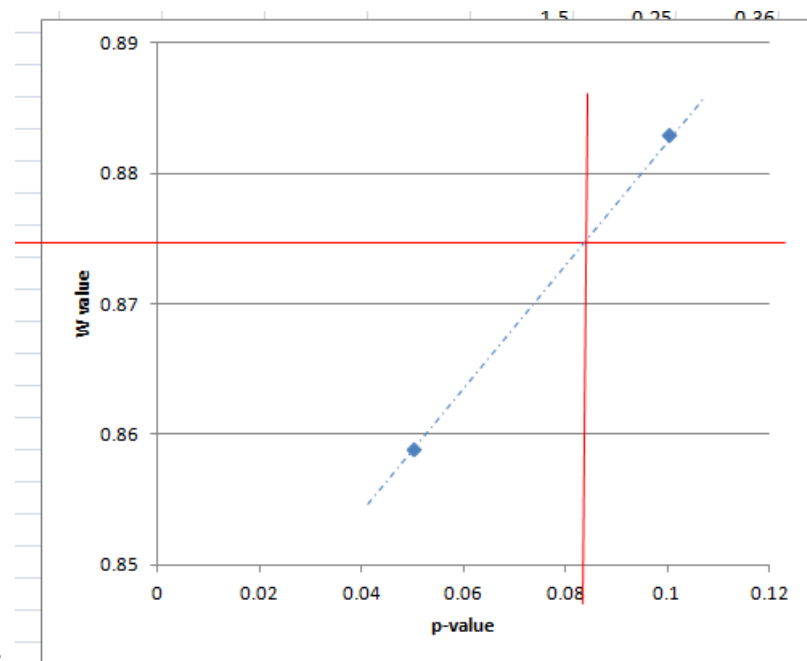
- Step 7:
  - Determine significance of this result by using look-up tables
  - Our calculated value (0.971026) lies between 0.974 and 0.943, i.e., between  $p=0.9$  and  $p=0.5$ .
  - As the smallest p-value is  $> 0.05$ , we accept the null hypothesis and conclude that the data are normally distributed

$n \backslash p$	0.01	0.02	0.05	0.1	0.5	0.9	0.95	0.98	0.99
3	0.753	0.756	0.767	0.789	0.959	0.998	0.999	1.000	1.000
4	0.687	0.707	0.748	0.792	0.935	0.987	0.992	0.996	0.997
5	0.686	0.715	0.762	0.806	0.927	0.979	0.986	0.991	0.993
6	0.713	0.743	0.788	0.826	0.927	0.974	0.981	0.986	0.989
7	0.730	0.760	0.803	0.838	0.928	0.972	0.979	0.985	0.988
8	0.749	0.778	0.818	0.851	0.932	0.972	0.978	0.984	0.987
9	0.764	0.791	0.829	0.859	0.935	0.972	0.978	0.984	0.986
10	0.781	0.806	0.842	0.869	0.938	0.972	0.978	0.983	0.986
11	0.792	0.817	0.850	0.876	0.940	0.973	0.979	0.984	0.986
12	0.805	0.828	0.859	0.883	0.943	0.973	0.979	0.984	0.986
13	0.814	0.837	0.866	0.889	0.945	0.974	0.979	0.984	0.986
14	0.825	0.846	0.874	0.895	0.947	0.975	0.980	0.984	0.986
15	0.835	0.855	0.881	0.901	0.950	0.975	0.980	0.984	0.987
16	0.844	0.863	0.887	0.906	0.952	0.976	0.981	0.985	0.987
17	0.851	0.869	0.892	0.910	0.954	0.977	0.981	0.985	0.987



# Shapiro-Wilk tutorial

- What if the value of  $W$  was 0.8752?
  - The corresponding p-values (for  $N = 12$ ) are 0.05 and 0.1.
  - The table values are 0.859 and 0.883
  - By linear interpolation, we can estimate a p-value that is around 0.08. As this is  $>0.05$ , the data are normally distributed



# T-statistic

- Should you use Independent or Repeated measures test?
- How do you define the significance level of the result?
- How do you use Cohen's  $d$  to calculate effect size?

# Tables for a (Shapiro-Wilk)

[illegible][illegible]

# More Tables for a (Shapiro-Wilk)

[illegible]

# More Tables for a (Shapiro-Wilk)

[illegible]

# P-values (Shapiro-Wilk)

n \ P	0.01	0.02	0.05	0.1	0.5	0.9	0.95	0.98	0.99		0.01	0.02	0.05	0.1	0.5	0.9	0.95	0.98	0.99
3	0.753	0.756	0.767	0.789	0.959	0.998	0.999	1.000	1.000	26	0.891	0.904	0.920	0.933	0.965	0.982	0.985	0.988	0.989
4	0.687	0.707	0.748	0.792	0.935	0.987	0.992	0.996	0.997	27	0.894	0.906	0.923	0.935	0.965	0.982	0.985	0.988	0.990
5	0.686	0.715	0.762	0.806	0.927	0.979	0.986	0.991	0.993	28	0.896	0.908	0.924	0.936	0.966	0.982	0.985	0.988	0.990
6	0.713	0.743	0.788	0.826	0.927	0.974	0.981	0.986	0.989	29	0.898	0.910	0.926	0.937	0.966	0.982	0.985	0.988	0.990
7	0.730	0.760	0.803	0.838	0.928	0.972	0.979	0.985	0.988	30	0.900	0.912	0.927	0.939	0.967	0.983	0.985	0.988	0.990
8	0.749	0.778	0.818	0.851	0.932	0.972	0.978	0.984	0.987	31	0.902	0.914	0.929	0.940	0.967	0.983	0.986	0.988	0.990
9	0.764	0.791	0.829	0.859	0.935	0.972	0.978	0.984	0.986	32	0.904	0.915	0.930	0.941	0.968	0.983	0.986	0.988	0.990
10	0.781	0.806	0.842	0.869	0.938	0.972	0.978	0.983	0.986	33	0.906	0.917	0.931	0.942	0.968	0.983	0.986	0.989	0.990
11	0.792	0.817	0.850	0.876	0.940	0.973	0.979	0.984	0.986	34	0.908	0.919	0.933	0.943	0.969	0.983	0.986	0.989	0.990
12	0.805	0.828	0.859	0.883	0.943	0.973	0.979	0.984	0.986	35	0.910	0.920	0.934	0.944	0.969	0.984	0.986	0.989	0.990
13	0.814	0.837	0.866	0.889	0.945	0.974	0.979	0.984	0.986	36	0.912	0.922	0.935	0.945	0.970	0.984	0.986	0.989	0.990
14	0.825	0.846	0.874	0.895	0.947	0.975	0.980	0.984	0.986	37	0.914	0.924	0.936	0.946	0.970	0.984	0.987	0.989	0.990
15	0.835	0.855	0.881	0.901	0.950	0.975	0.980	0.984	0.987	38	0.916	0.925	0.938	0.947	0.971	0.984	0.987	0.989	0.990
16	0.844	0.863	0.887	0.906	0.952	0.976	0.981	0.985	0.987	39	0.917	0.927	0.939	0.948	0.971	0.984	0.987	0.989	0.991
17	0.851	0.869	0.892	0.910	0.954	0.977	0.981	0.985	0.987	40	0.919	0.928	0.940	0.949	0.972	0.985	0.987	0.989	0.991
18	0.858	0.874	0.897	0.914	0.956	0.978	0.982	0.986	0.988	41	0.920	0.929	0.941	0.950	0.972	0.985	0.987	0.989	0.991
19	0.863	0.879	0.901	0.917	0.957	0.978	0.982	0.986	0.988	42	0.922	0.930	0.942	0.951	0.972	0.985	0.987	0.989	0.991
20	0.868	0.884	0.905	0.920	0.959	0.979	0.983	0.986	0.988	43	0.923	0.932	0.943	0.951	0.973	0.985	0.987	0.990	0.991
21	0.873	0.888	0.908	0.923	0.960	0.980	0.983	0.987	0.989	44	0.924	0.933	0.944	0.952	0.973	0.985	0.987	0.990	0.991
22	0.878	0.892	0.911	0.926	0.961	0.980	0.984	0.987	0.989	45	0.926	0.934	0.945	0.953	0.973	0.985	0.988	0.990	0.991
23	0.881	0.895	0.914	0.928	0.962	0.981	0.984	0.987	0.989	46	0.927	0.935	0.945	0.953	0.974	0.985	0.988	0.990	0.991
24	0.884	0.898	0.916	0.930	0.963	0.981	0.984	0.987	0.989	47	0.928	0.936	0.946	0.954	0.974	0.985	0.988	0.990	0.991
25	0.888	0.901	0.918	0.931	0.964	0.981	0.985	0.988	0.989	48	0.929	0.937	0.947	0.954	0.974	0.985	0.988	0.990	0.991
										49	0.929	0.939	0.947	0.955	0.974	0.985	0.988	0.990	0.991
										50	0.930	0.938	0.947	0.955	0.974	0.985	0.988	0.990	0.991