

# **Overview of the content**

Research Methods for Human Inquiry  
Andrew Perfors

# We've gone over a lot of stuff...

- Introduction to R
- Data manipulation
- Figures
- Statistical theory
- Chi-square tests
- t-tests
- ANOVA
- Regression
- Advanced topics

# Contents of RMHI

- Introduction to R
  - Packages, R Markdown, variables, data frames
  - Functions, arguments
- Data manipulation
- Figures
- Statistical theory
- Chi-square tests
- t-tests
- ANOVA
- Regression
- Advanced topics

# Contents of RMHI

- Introduction to R
- Data manipulation
  - Pipes, grouping, summary/descriptive statistics
  - Long vs wide form, filter, mutate
- Figures
- Statistical theory
- Chi-square tests
- t-tests
- ANOVA
- Regression

# Contents of RMHI

- Introduction to R
- Data manipulation
- Figures
  - Scatterplots, boxplots, violin plots, bar plots, histograms (all assessable, but especially these)
  - Geoms, arguments, layering, colour palettes
- Statistical theory
- Chi-square tests
- t-tests
- ANOVA
- Regression

# Contents of RMHI

- Introduction to R
- Data manipulation
- Figures
- Statistical theory
  - Null hypothesis significance testing; p-values
  - Type 1 and Type 2 error
  - Binomial and normal distributions
  - Sampling distributions; confidence intervals
  - Bayesian vs frequentist approaches
- Chi-square tests
- t-tests
- ANOVA
- Regression

# Contents of RMHI

- Introduction to R
- Data manipulation
- Figures
- Statistical theory
- Chi-square tests
  - Goodness of fit test, test of independence
  - Chi-squared test statistic
  - Fisher exact test, McNemar test
  - Cramers' V
- t-tests
- ANOVA
- Regression

# Contents of RMHI

- Introduction to R
- Data manipulation
- Figures
- Statistical theory
- Chi-square tests
- t-tests
  - one-sample, paired and independent t-tests
  - t statistic
  - Student's vs Welch's t-tests
  - Cohen's d
  - Shapiro-Wilk test, QQ Plots, Wilcoxon
- ANOVA
- Regression

# Contents of RMHI

- Introduction to R
- Data manipulation
- Figures
- Statistical theory
- Chi-square tests
- t-tests
- ANOVA
  - one-way ANOVA, two-way ANOVA, interactions
  - post hoc tests and multiple comparison correction
  - effect size with eta-squared
  - F statistic, SS<sub>b</sub>, SS<sub>w</sub>
  - Levene test, Welch one-way ANOVA, Kruskal-Wallis
- Regression

# Contents of RMHI

- Introduction to R
- Data manipulation
- Figures
- Statistical theory
- Chi-square tests
- t-tests
- ANOVA
- Regression
  - Spearman & Pearson correlations
  - Interpretation of regression line; significance; residuals
  - Effect size with R-squared
  - Coefficients and standardised coefficients

# Contents of RMHI

- Introduction to R
- Data manipulation
- Figures
- Statistical theory
- Chi-square tests
- t-tests
- ANOVA
- Regression
- Advanced topics
  - Assumptions of linear regression: linearity, normality of residuals, collinearity, high-influence points
  - Model selection, penalising complexity, AIC, BIC

# Some things I didn't get to...

This is not examinable! I just want you to know it exists and where to find resources for it when/if it comes up for you in the future

- Everything is a linear model and model selection!
- Unbalanced ANOVAs
- Beyond linear regression
- Mixed-effect models
- Bayesian data analysis

# Everything is linear model and model selection!

	Outcome variable	Predictor variable(s)
Chi-squared test	counts/frequencies	Categorical
t-test	numeric	1 to 2 groups (categorical)
ANOVA	numeric	3+ groups, or multiple factors
regression	numeric	Numeric

Normal regression  $\longrightarrow Y_i = b_1 X_{1i} + b_2 X_{2i} + b_0 + \varepsilon_i$

Equivalent to ANOVA/t-test when slopes b are binary (0 or 1). i.e. group A is 0 and group B is 1

Equivalent to chi-squared when it is log-linear (i.e., take the log of everything)

# Everything is linear model and model selection!

See Chapter 16 in *Learning Statistics with R* and the linear models cheat sheet linked to in the LMS this week if you want to learn more

## Common statistical tests are linear models

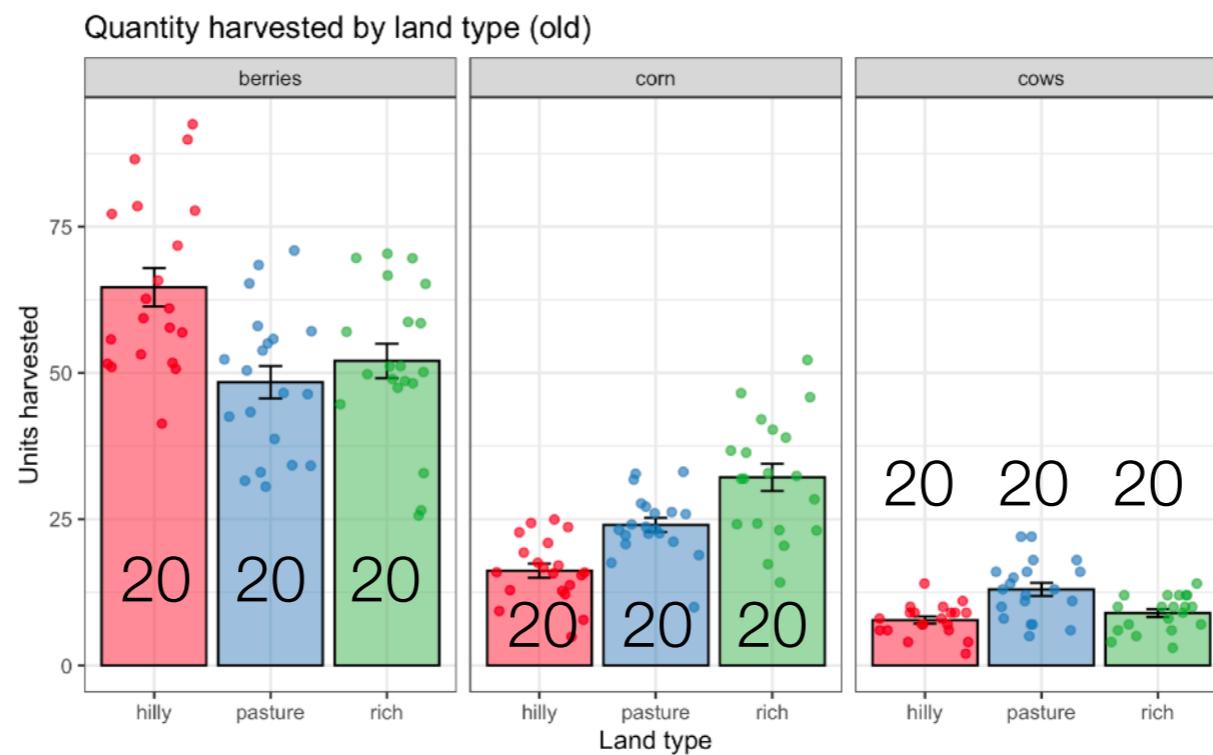
Last updated: 28 June, 2019. Also check out the [Python version!](#)

See worked examples and more details at the accompanying notebook: <https://lindeloev.github.io/tests-as-linear>

Common name	Built-in function in R	Equivalent linear model in R	Exact?	The linear model in words	Icon
Simple regression: $\text{lm}(y \sim 1 + x)$	<b>y is independent of x</b> P: One-sample t-test N: Wilcoxon signed-rank	<code>lm(y ~ 1)</code> <code>lm(signed_rank(y) ~ 1)</code>	✓ <a href="#">for N &gt; 14</a>	One number (intercept, i.e., the mean) predicts y. - (Same, but it predicts the <i>signed rank</i> of y.)	
	P: Paired-sample t-test N: Wilcoxon matched pairs	<code>lm(y2 - y1 ~ 1)</code> <code>lm(signed_rank(y2 - y1) ~ 1)</code>	✓ <a href="#">for N &gt; 14</a>	One intercept predicts the pairwise $y_2 - y_1$ differences. - (Same, but it predicts the <i>signed rank</i> of $y_2 - y_1$ .)	
	<b>y ~ continuous x</b> P: Pearson correlation N: Spearman correlation	<code>lm(y ~ 1 + x)</code> <code>lm(rank(y) ~ 1 + rank(x))</code>	✓ <a href="#">for N &gt; 10</a>	One intercept plus $x$ multiplied by a number (slope) predicts y. - (Same, but with <i>ranked x</i> and <i>y</i> )	
	<b>y ~ discrete x</b> P: Two-sample t-test P: Welch's t-test N: Mann-Whitney U	<code>lm(y ~ 1 + G2)^A</code> <code>glm(y ~ 1 + G2, weights=...)^B</code> <code>lm(signed_rank(y) ~ 1 + G2)^A</code>	✓ ✓ <a href="#">for N &gt; 11</a>	An intercept for <b>group 1</b> (plus a difference if <b>group 2</b> ) predicts y. - (Same, but with one variance per group instead of one common.) - (Same, but it predicts the <i>signed rank</i> of y.)	
Multiple regression: $\text{lm}(y \sim 1 + x_1 + x_2 + \dots)$	P: One-way ANOVA N: Kruskal-Wallis	<code>lm(y ~ group)</code> <code>kruskal.test(y ~ group)</code>	✓ <a href="#">for N &gt; 11</a>	An intercept for <b>group 1</b> (plus a difference if $group \neq 1$ ) predicts y. - (Same, but it predicts the <i>rank</i> of y.)	
	P: One-way ANCOVA	<code>lm(y ~ group + x)</code>	✓	- (Same, but plus a slope on x.) <i>Note: this is discrete AND continuous. ANCOVAs are ANOVAs with a continuous x.</i>	
	P: Two-way ANOVA	<code>lm(y ~ group * sex)</code>	✓	Interaction term: changing <b>sex</b> changes the <b>y ~ group</b> parameters. <i>Note: <math>G_{2 to N}</math> is an indicator (0 or 1) for each non-intercept levels of the group variable. Similarly for <math>S_{2 to K}</math> for sex. The first line (with <math>G_i</math>) is main effect of group, the second (with <math>S_j</math>) for sex and the third is the group <math>\times</math> sex interaction. For two levels (e.g. male/female), line 2 would just be "S<sub>2</sub>" and line 3 would be S<sub>2</sub> multiplied with each G<sub>i</sub>.</i>	[Coming]
	<b>Counts ~ discrete x</b> N: Chi-square test	<code>chisq.test(groupXsex_table)</code>	✓	Interaction term: (Same as Two-way ANOVA.) <i>Note: Run glm using the following arguments: <code>glm(model, family=poisson())</code>. As linear-model, the Chi-square test is <math>\log(y_i) = \log(N) + \log(\alpha_i) + \log(\beta_i) + \log(\alpha_i\beta_i)</math> where <math>\alpha_i</math> and <math>\beta_i</math> are proportions. See more info in the accompanying notebook.</i>	Same as Two-way ANOVA
N: Goodness of fit	<code>chisq.test(y)</code>	<code>glm(y ~ 1 + G2 + G3 + ... + GN, family=...)^A</code>	✓	(Same as One-way ANOVA and see Chi-Square note.)	1W-ANOVA

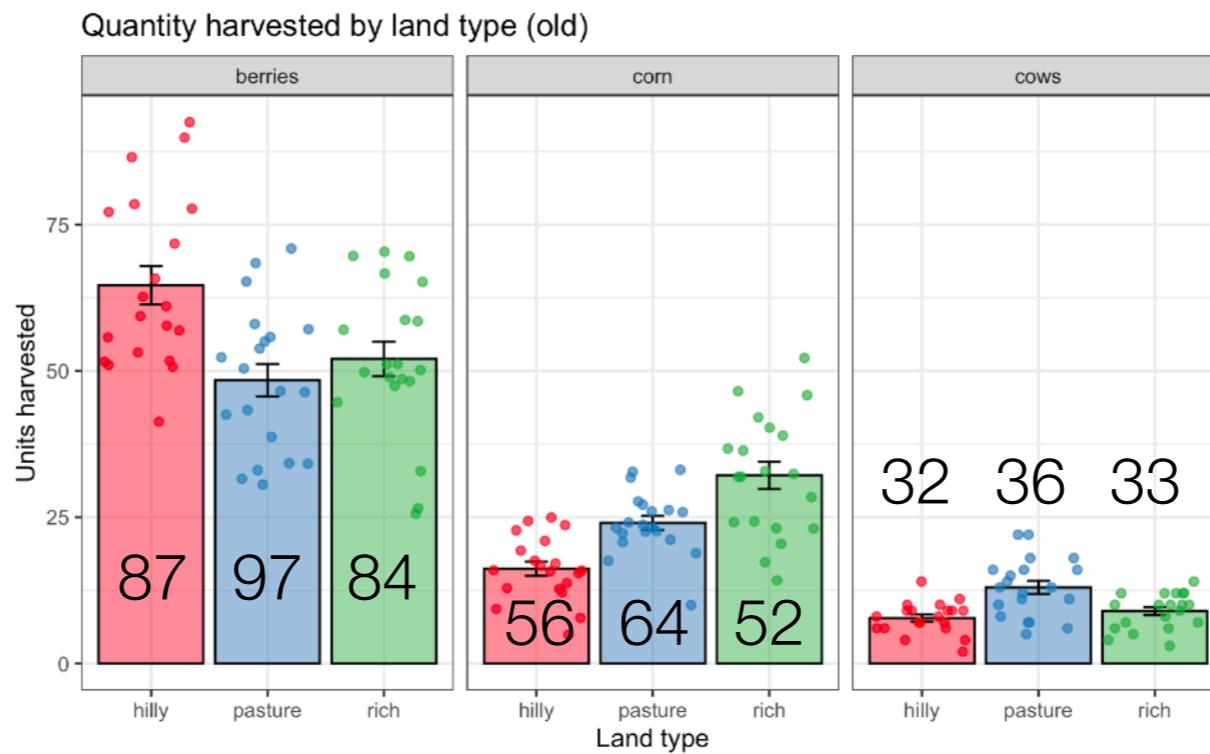
# Unbalanced ANOVAs

We covered balanced ANOVAs, which is what happens when there are equal numbers of data points in each group and sub-group



# Unbalanced ANOVAs

It's very common to have unbalanced ones, which are more complicated to calculate because the *order* you evaluate variables changes the results



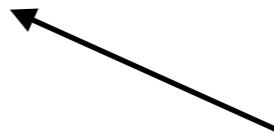
There are three different ways to do things: Type I, Type II, and Type III sums of squares. They work by doing model selection on models with or without each variable, and vary in the order they compare models

# Unbalanced ANOVAs

Short version: I recommend Type II sums of squares. Type III is very common. To run them you use the `lm()` function followed by `Anova()` from the `car` package, where you specify the type

```
> library(car)
> myModel <- lm(outcomeVar ~ varA + varB, data=d)
> Anova(myModel, type=3)
```

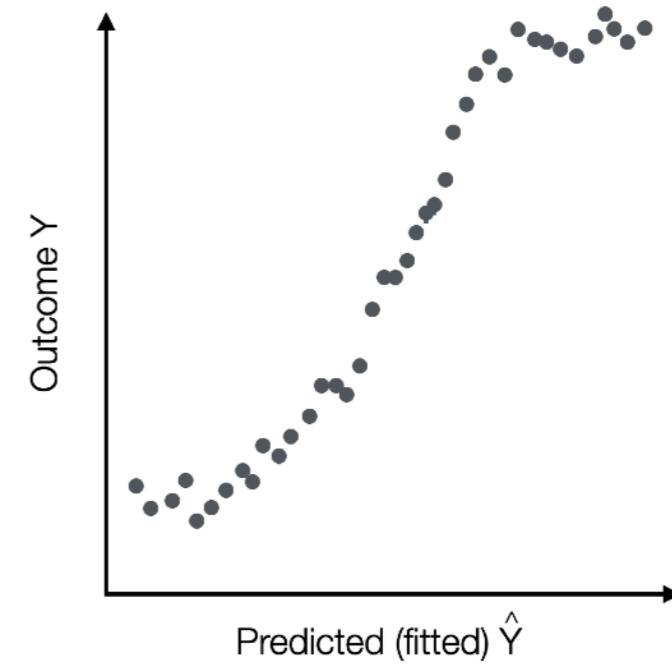
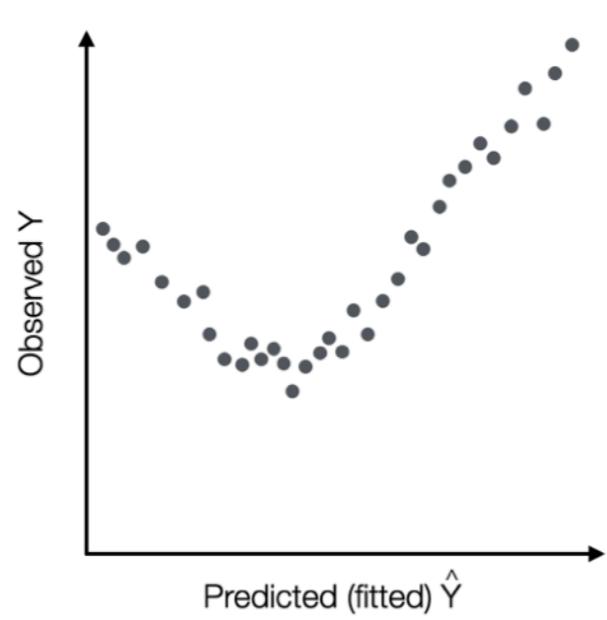
Does a Type III sums of squares.  
`type=2` for Type II



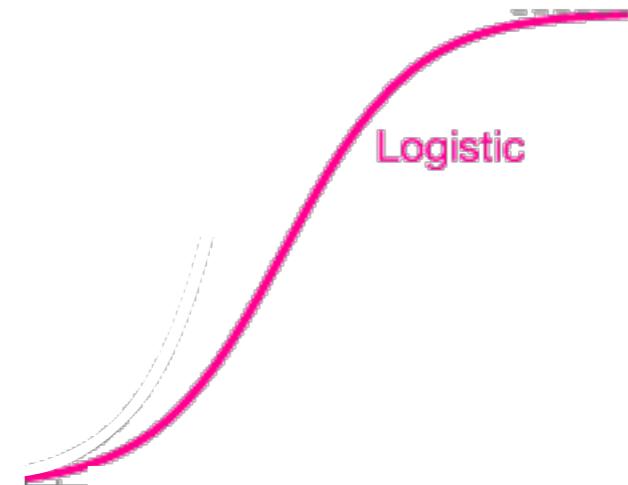
See Chapter 16 in *Learning Statistics with R* for lots more details!

# Beyond linear regression

What if you have numeric outcomes and predictors, but the relationship is not linear?

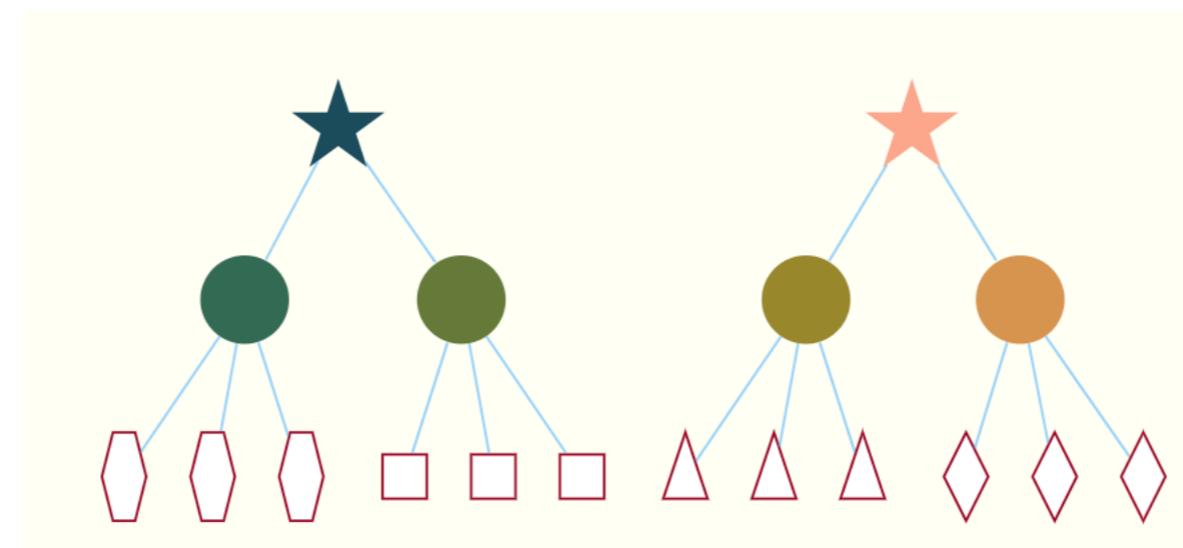


Answer: use other kinds of equation!



# Mixed models

Occur when some predictors are nested in others (e.g. each participant contributes multiple trials; you're studying students, each of which occurs in some classroom)



Very useful for making the most out of non-independent data!

The LMS contains some useful resources on how to do these

# Bayesian data analysis

Instead of testing a null hypothesis, you compare hypotheses directly by calculating the **likelihood** of the data under each and combining it with its **prior** probability

$$\frac{P(H_2 | D)}{P(H_1 | D)} = \frac{P(D | H_2) P(H_2)}{P(D | H_1) P(H_1)}$$

Posterior probability that each hypothesis is true, given the data

Likelihood that the data is true, given the hypothesis

Prior probability that the hypothesis is true

See Week 9 Seminar, Chapter 17 of *Learning Statistics with R*, or  
*Statistical Rethinking* by McElreath

# Wow. That's a lot of stuff

In the last video I will try to give a bit more of a “forest” view of all of this, rather than the trees we’ve been mired in lately

