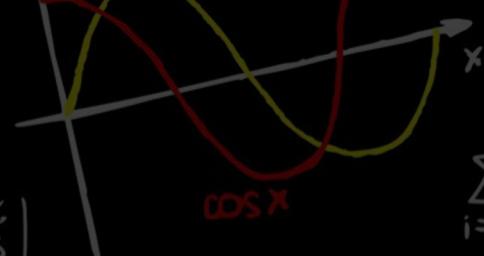


DATA ANALYSIS WITH R

BRUNO BELLISARIO, PHD

SESSION 4 STATISTICAL INFERENCE#2



$$Y_{i+1} = Y_i + b \cdot K_2$$

$$B = \begin{pmatrix} 2 & 1 & -1 & 0 \\ 3 & 0 & 1 & 2 \end{pmatrix}$$

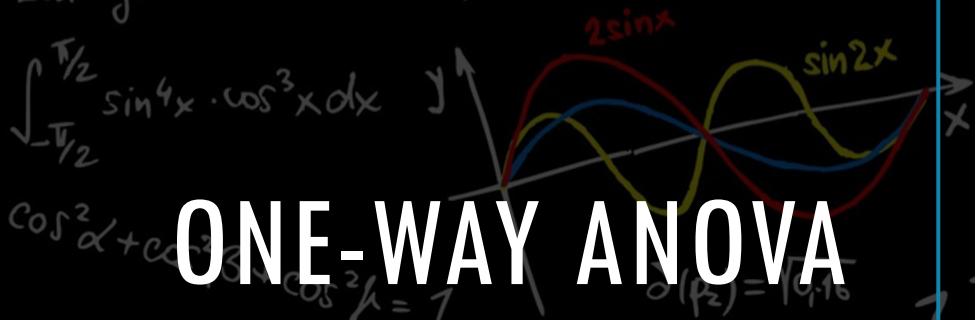
$$a^2 = b^2 + c^2 - 2bc \cos \alpha$$

$$\operatorname{tg} \frac{x}{2} = \frac{1 - \cos x}{\sin x} = \frac{\sin x}{1 + \cos x}$$

$$X_2 = \begin{pmatrix} -\lambda \\ -\beta \\ -\gamma \\ -\delta \end{pmatrix}$$

$$\iiint_M z dx dy dz = \int_0^{2\pi} \left(\int_0^2 \left(\int_{\frac{1}{2}\rho}^1 r^2 z d\sigma \right) dr \right) d\varphi$$

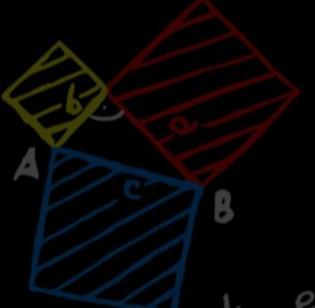
$$2 \arctg x - x = 0, I = (1, 10)$$



ONE-WAY ANOVA

$$\frac{\partial z}{\partial x} = 2, \frac{\partial z}{\partial y} = 0 \quad \vec{n} = (F_x, F_y, F_z)$$

$$\frac{x^2}{a^2} + \frac{y^2}{b^2} + \frac{z^2}{c^2} = 0$$



$$f(x) =$$

$$e^x - xy^2 = e; A[0; e; 1]$$

$$\frac{2x}{x^2 + 2y^2} = 2 \quad z = \frac{1}{x} \arctan \frac{y}{2}$$

$$\eta_1 = \lambda_1^2 - 3\lambda_1 + 1 + 0$$

$$\sin 2x = 2 \sin x \cdot \cos x$$

$$|Z| = \sqrt{a^2 + b^2}$$

$$\Delta(\frac{\partial F}{\partial z}) = 16 - x^2 + 16y^2 - 4z > 0$$

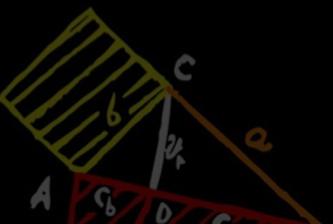
$$\sum_{i=0}^n (P_2(x_i) - y_i)^2 \quad \operatorname{tg} 2x = \frac{2 \operatorname{tg} x}{1 - \operatorname{tg}^2 x}$$

$$\begin{aligned} \lambda x - y + z &= 1 \\ x + \lambda y + z &= \lambda \\ x + y + \lambda z &= \lambda^2 \end{aligned}$$



$$F_2 = 2 \times y^2 - 1 = 1$$

$$x_1 = \begin{pmatrix} 2\rho \\ 0 \end{pmatrix} \quad \begin{array}{l} y \\ y=x^2 \\ y=x^3 \end{array}$$



$$\lim_{n \rightarrow \infty}$$

$$n$$

$$\infty$$

$$\dots$$

$$\$$

$\cos x$ $\sin x$ $\tan \frac{x}{2} = \frac{1-\cos x}{\sin x} = \frac{\sin x}{1+\cos x}$
 $x_2 = \begin{pmatrix} -\alpha \\ \beta \\ -\gamma \\ -\delta \end{pmatrix}$ $\sum_{i=0}^n (P_2(x_i) - y_i)^2$ $\tan 2x = \frac{2\tan x}{1-\tan^2 x}$ $F_0 = 2 \times \sqrt{2} - 1 = 1$
 $\int \int \int_M z dx dy dz = \int_0^{2\pi} \left(\int_0^2 \left(\int_{\frac{1}{2}\pi}^{\pi} r^2 r d\sigma \right) dr \right) d\varphi$ $y = x^3$ $y = x^2$
 $2\arctan x - x = 0, I = (1, 10)$ $\lim_{n \rightarrow \infty} \sqrt[3]{3n^2 + 2n - 1}$ $y = x^4$
 $\int_{-\pi/2}^{\pi/2} \sin^4 x \cdot \cos^3 x dx$ $\sin 2x = 2\sin x \cos x$ $\cos^2 x + \sin^2 x = 1$ $A + B + C = 8$
 $\cos^2 \alpha + \cos^2 \beta + \cos^2 \gamma = 1$ $\partial(P_2) = \sqrt{16}$ $-3A - 7B + 2C = 10,3$
 $\partial_x = 2, \partial_y = 0$ $\vec{n} = (F_x, F_y, F_z)$ $-18A + 6B - 3C = 15$
 $\frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2} = 0$ $\alpha^2 + \beta^2 = c^2$ $f(x) = 2^{-x+1}, \epsilon = 0.005$
 $\sin^2 x = 2\sin x \cos x$ $\alpha, \beta, \gamma \in C$ $e^x - xy^2 = e, A[0; e; 1]$
 $|z| = \sqrt{a^2 + b^2}$ $\lim_{x \rightarrow 0} \frac{e^{2x} - 1}{5x} = \frac{2}{5}$ $\lambda_2 = i\sqrt{14}$
 $D(\frac{\partial F}{\partial z}) = 16 - x^2 + 16y^2 - 4z > 0$ $\lim_{x \rightarrow 0} \frac{e^{2x} - 1}{5x} = \frac{2}{5}$ $\frac{\sin x}{x} \leq \frac{x}{x} = 1$
 $\sin(x+y) = \sin x \cos y + \cos x \sin y$ $\lambda_1 = 1 + \sqrt{2}$
 $(x, 1+x^2, 1)$ $x=0, y=1, z=2$ $y' - \frac{\sqrt{y}}{x+2} = 0, y(0)=1$
 $\alpha^2 = b^2 + c^2 - 2bc \cos A$

ONE-WAY ANOVA

Appropriate data

- One-way data. That is, one measurement variable in two or more groups
- Dependent variable is interval/ratio, and is continuous
- Independent variable is a factor with two or more levels (two or more groups)
- Residuals are normally distributed (moderate deviation from normally-distributed residuals is permissible)
- Groups have the same variance (homoscedasticity)
- Observations among groups are independent.

Hypotheses

- Null hypothesis: The means of the measurement variable for each group are equal.
- Alternative hypothesis (two-sided): The means of the measurement variable for among groups are not equal.

ONE-WAY ANOVA W/BLOCKS

Blocks are used in an analysis of variance or similar models in order to account for suspected variation from factors other than the treatments or main independent variables being investigated.

Traditionally, in agricultural experiments, plots would be arranged into blocks according to factors in the field that could not be controlled.

The experiment might be designed in a randomized complete block design in which each block had a plot with each treatment.

$$\tan \frac{x}{2} = \frac{1 - \cos x}{\sin x} = \frac{\sin x}{1 + \cos x}$$

$$F_2 = 2 \times y^2 - 1 = 1$$

$$y = x^3$$

$$y = x^2$$

$$y = x^4$$

$$y = x^5$$

$$y = x^6$$

$$A + B + C = 8$$

$$-3A - 7B + 2C = 10,3$$

$$-18A + 6B - 3C = 15$$

ONE-WAY ANOVA W/BLOCKS

$$B = \begin{pmatrix} 2 & 1 & -1 & 0 \\ 3 & 0 & 1 & 2 \end{pmatrix}$$

$$a^2 = b^2 + c^2 - 2bc \cos \alpha$$

$$Y_{i+1} = Y_i + b \cdot k_2$$

$$\sum_{i=0}^n (P_2(x_i) - y_i)^2$$

$$\tan 2x = \frac{2 \tan x}{1 - \tan^2 x}$$

Appropriate data

→ One-way data, with blocks. That is, one measurement variable in two or more groups, where each group is also distributed among at least two blocks

→ Dependent variable is interval/ratio, and is continuous

→ Independent variable is a factor with two or more levels.

→ A second independent variable is a blocking factor variable with two or more levels

→ Residuals are normally distributed (moderate deviation from normally-distributed residuals is permissible)

→ Groups have the same variance (homoscedasticity)

→ Observations among groups are independent.

Hypotheses

→ Null hypothesis: The means of the measurement variable for each group are equal

→ Alternative hypothesis (two-sided): The means of the measurement variable among groups are not equal

→ An additional null hypothesis is tested for the effect of blocks: The means of the measurement variable for each block are equal

$$X_2 = \begin{pmatrix} \alpha \\ \beta \\ -\gamma \\ -\delta \end{pmatrix}$$

$$\iiint_M z dx dy dz = \int_0^{2\pi} \left(\int_0^{\pi} \left(\int_{1/2}^1 r^2 r dr \right) dr \right) d\varphi$$

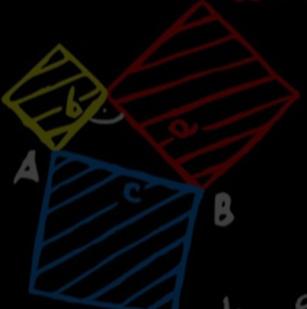
$$2 \arctan x - x = 0, I = (1, 10)$$



$$\cos^2 \alpha + \cos^2 \beta + \cos^2 \mu = 1$$

$$\frac{\partial z}{\partial x} = 2, \frac{\partial z}{\partial y} = 0 \quad \vec{n} = (F_x, F_y, F_z)$$

$$\frac{x^2}{a^2} + \frac{y^2}{b^2} + \frac{z^2}{c^2} = 0$$



$$a^2 + b^2 = c^2$$

$$f(x) = 2^{-x} \cdot 1, f = 0.005$$

$$\lim_{x \rightarrow 0} \frac{e^{2x} - 1}{5x} = \frac{2}{5}$$

$$\sin 2x = 2 \sin x \cdot \cos x$$

$$|Z| = \sqrt{a^2 + b^2}$$

$$\Delta(\frac{\partial F}{\partial z}) = 16 - x^2 + 16y^2 - 4z > 0$$

$$\begin{pmatrix} x, 1+x^2, 1 \\ 1, 1-x^2, 1 \end{pmatrix}, x=0, y=1, z=2$$

$$y' - \frac{\sqrt{y}}{x+2} = 0, y(0) = 1$$

$$\sin(x+y) = \sin x \cos y + \cos x \sin y$$



ONE-WAY ANOVA W/RANDOM BLOCKS

Here, the analysis is done with a **mixed effects model**, with the treatments treated as a fixed effect and the blocks treated as a random effect.

In analysis of variance, blocking variables are often treated as random variables.

This is because the blocking variable represents a random selection of levels of that variable.

The analyst wants to take the effects of the blocking variable into account, but the identity of the specific levels of the blocks are not of interest.

TWO-WAY ANOVA

A two-way ANOVA can investigate the main effects of each of two independent factor variables, as well as the effect of the interaction of these variables.

Appropriate data

- Two-way data
- Dependent variable is interval/ratio, and is continuous
- Independent variables are a factor with two or more levels
- Residuals are normally distributed (moderate deviation from normally-distributed residuals is permissible)
- Groups have the same variance (homoscedasticity)
- Observations among groups are independent

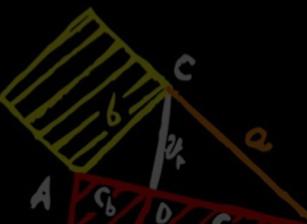
Hypotheses

- Null hypothesis 1: The means of the dependent variable for each group in the first independent variable are equal
- Alternative hypothesis 1 (two-sided): The means of the dependent variable for each group in the first independent variable are not equal

REPEATED MEASURE ANOVA

When an experimental design takes measurements on the same experimental unit over time, the analysis of the data must take into account the probability that measurements for a given experimental unit will be correlated in some way.

For example, if we were measuring calorie intake for students, we would expect that if one student had a higher intake at Time 1, that that student would have a higher intake others at Time 2, and so on.



$$\begin{aligned} A+B+C &= 8 \\ -3A-7B+2C &= 10,3 \\ BC &= 15 \end{aligned}$$

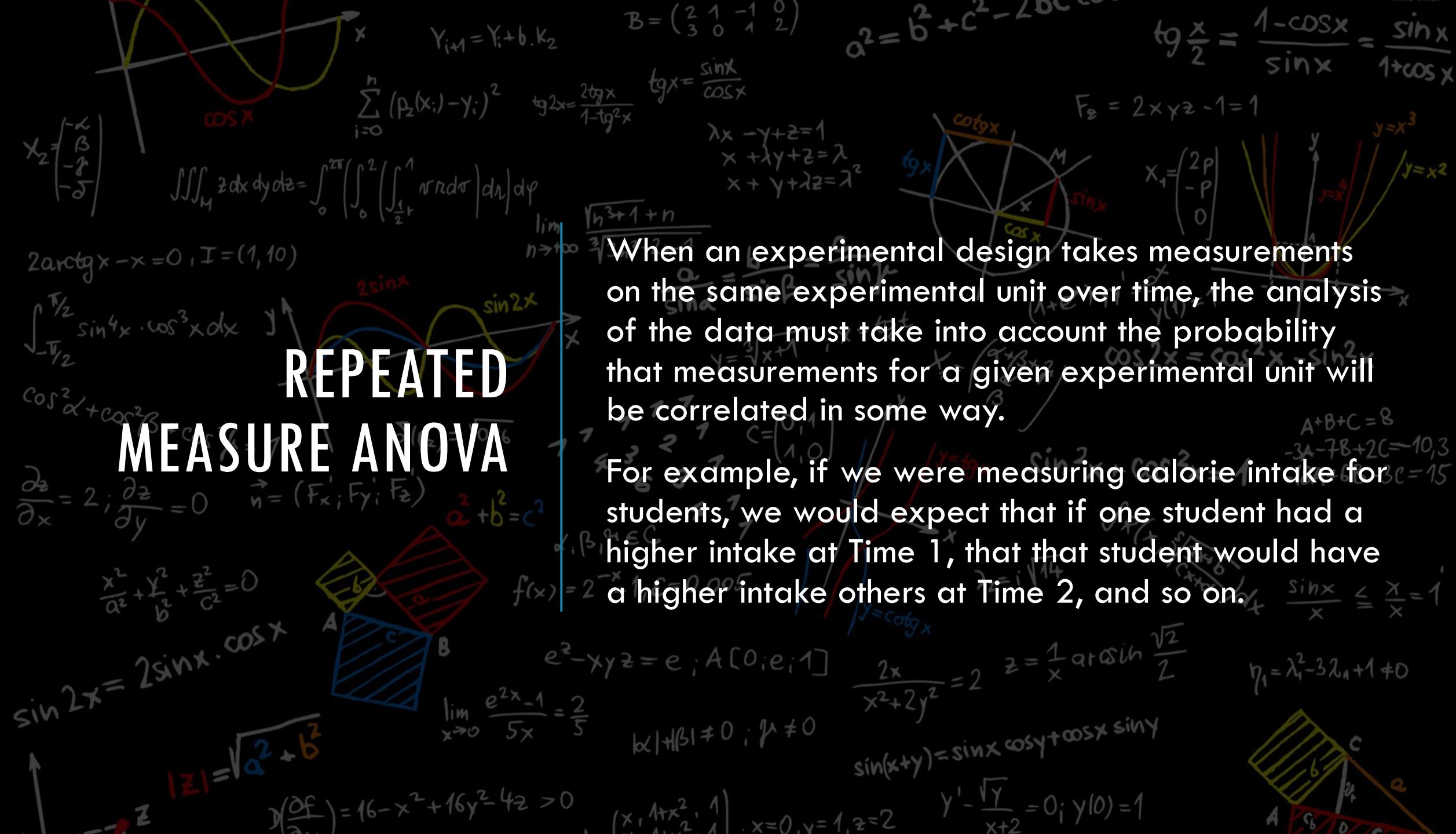
$$\frac{\sin x}{x} \leq \frac{x}{x} = 1$$

$$\eta_1 = \lambda_1^2 - 3\lambda_1 + 1 + 0$$

$$\frac{2x}{x^2+2y^2} = 2 \quad z = \frac{1}{x} \arctan \frac{\sqrt{2}}{2}$$

$$\sin(x+y) = \sin x \cos y + \cos x \sin y$$

$$y' - \frac{\sqrt{y}}{x+2} = 0 \quad y(0) = 1$$



CORRELATION & LINEAR REGRESSION

CORRELATION & LINEAR REGRESSION

- Linear regression specifies one variable as the independent variable and another as the dependent variable.
- The resultant model relates the variables with a linear relationship.
- The tests associated with linear regression are parametric and assume normality, homoscedasticity, and independence of residuals, as well as a linear relationship between the two variables.

CORRELATION & LINEAR REGRESSION

CORRELATION & LINEAR REGRESSION

The Kendall correlation method measures the correspondence between the ranking of x and y variables in a different way.

The total number of possible pairings of x with y observations is $n(n-1)/2$, where n is the size of x and y .

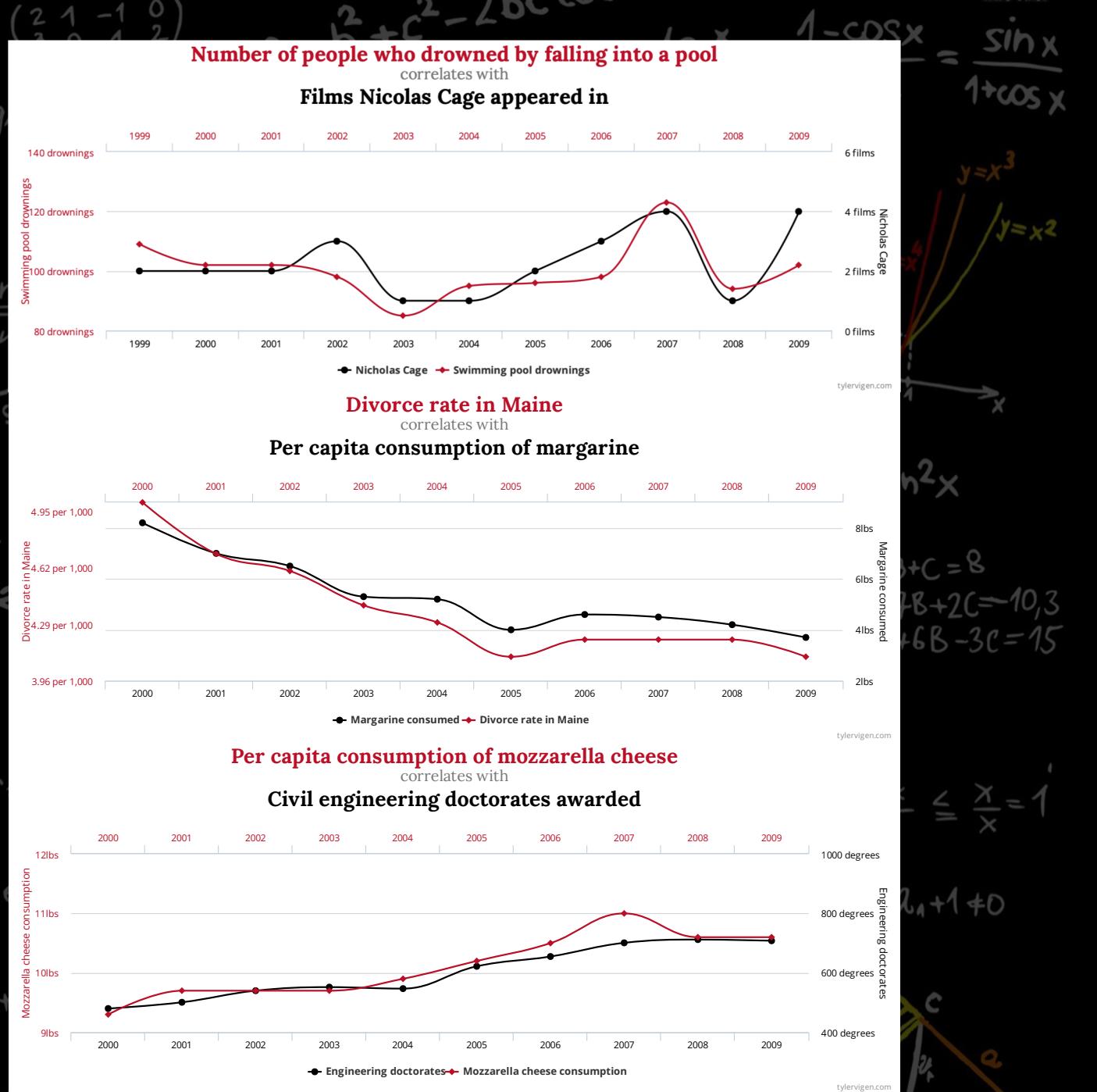
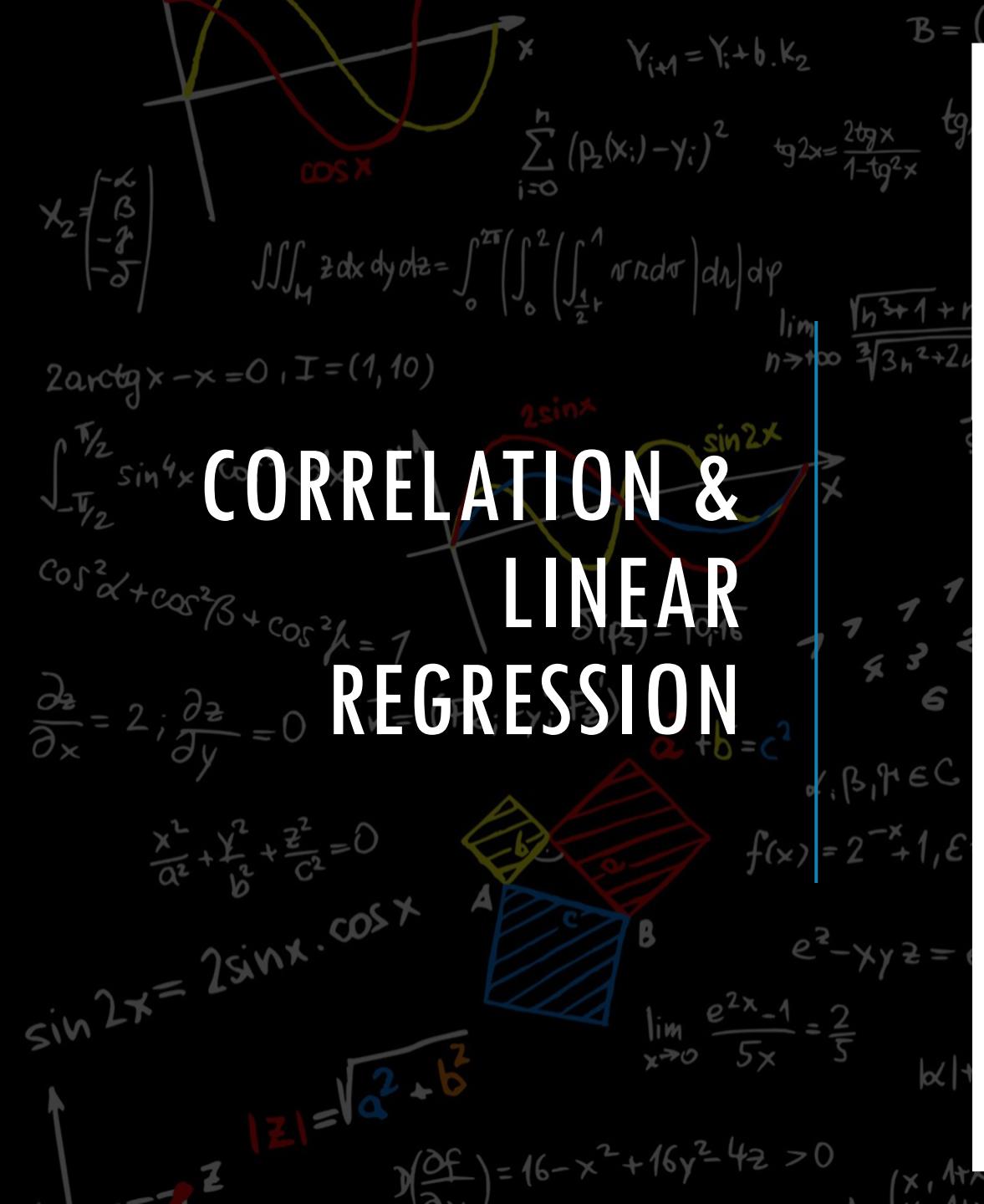
In the formula below, n_c is total number of concordant pairs n_d is total number of discordant pairs and n is the size of x and y

CORRELATION & LINEAR REGRESSION

Beware of spurious correlations!

In statistics, a spurious relationship (or spurious correlation) is a mathematical relationship in which two or more events or variables are associated but not causally related, due to either coincidence or the presence of a certain third, unseen factor (referred to as a "common response variable", "confounding factor", or "lurking variable").

CORRELATION & LINEAR REGRESSION



CORRELATION & LINEAR REGRESSION

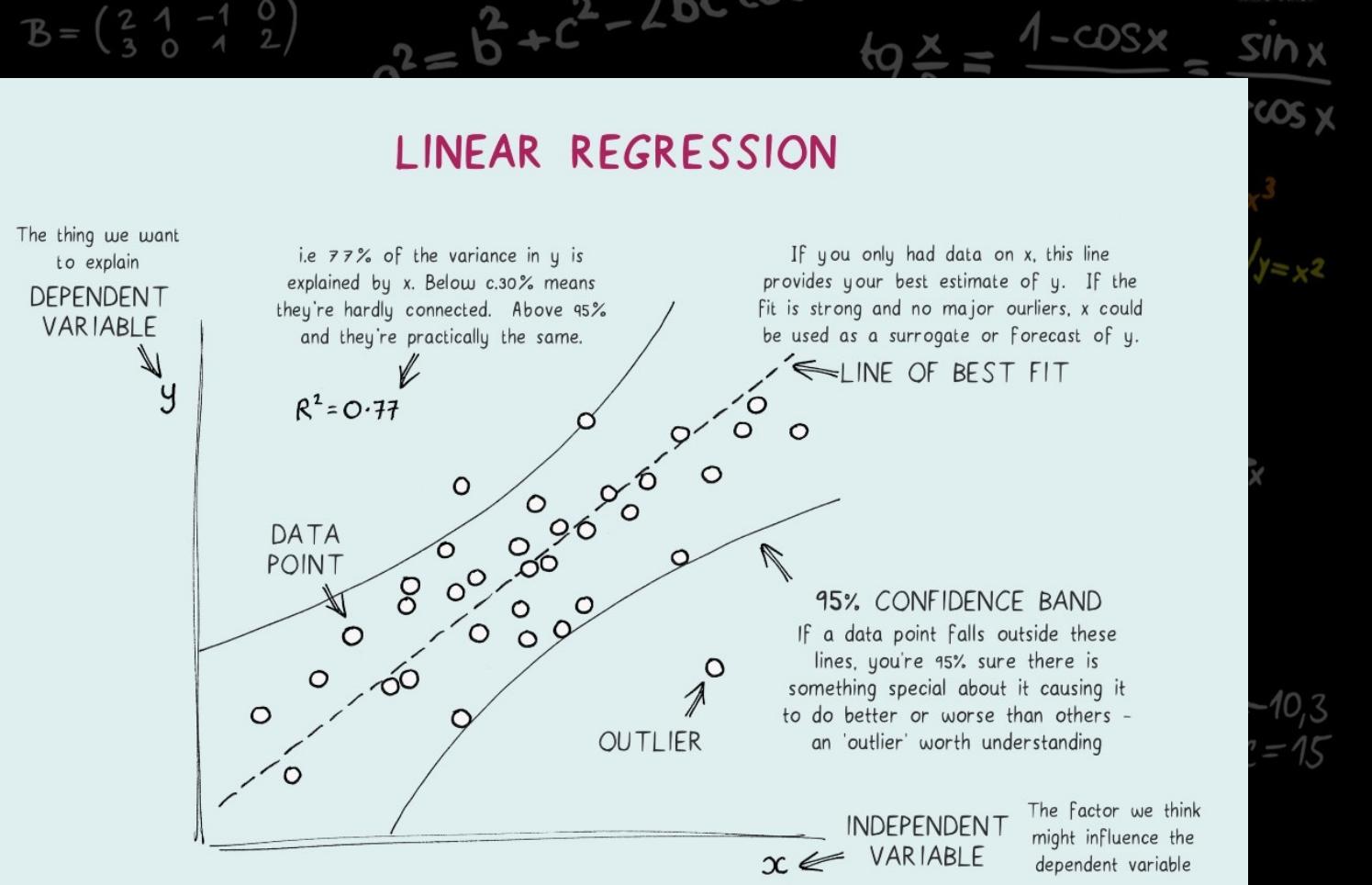
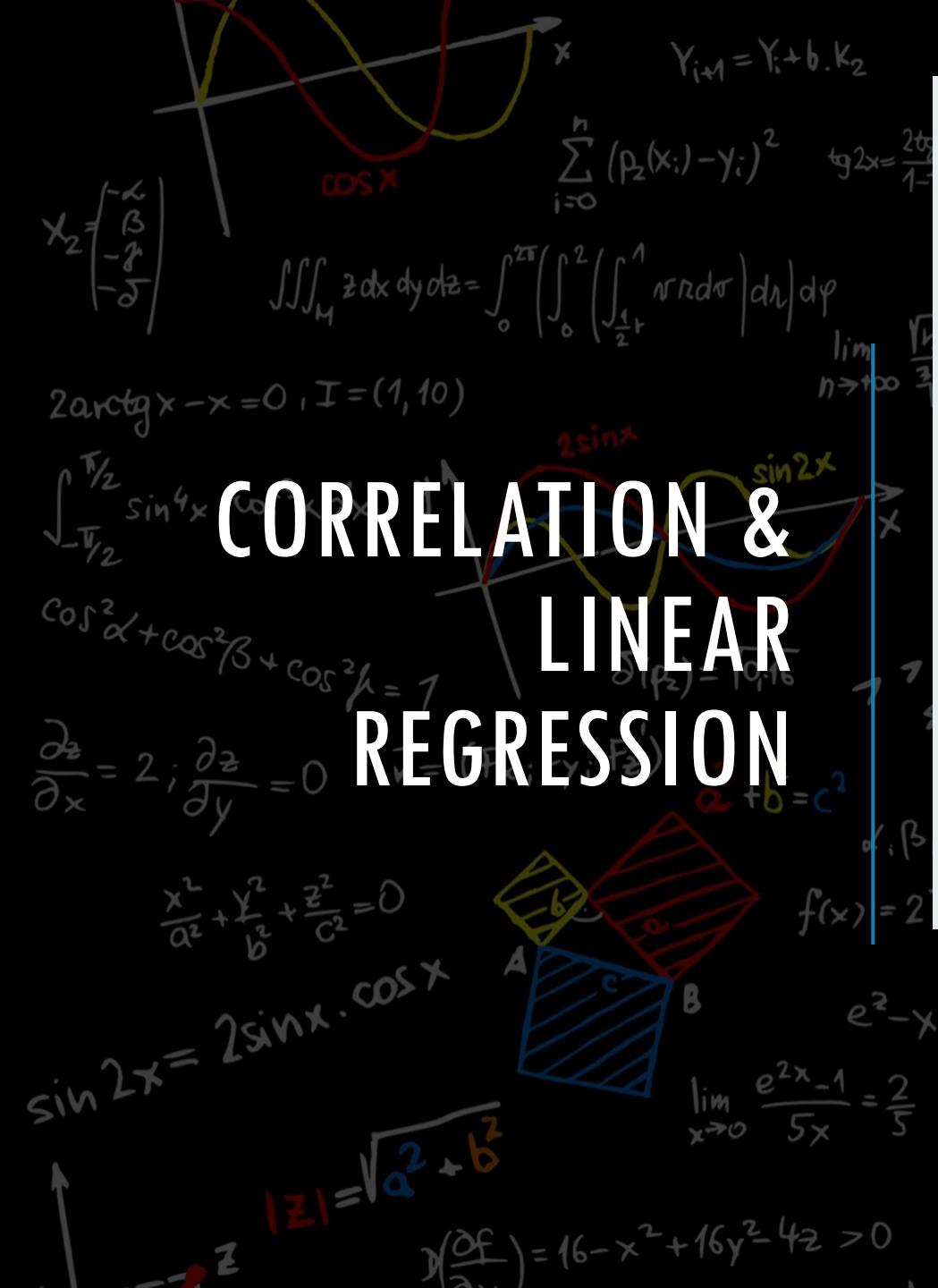
Linear regression is a very common approach to model the relationship between two interval/ratio variables.

The method assumes that there is a linear relationship between the dependent variable and the independent variable, and finds a best fit model for this relationship.

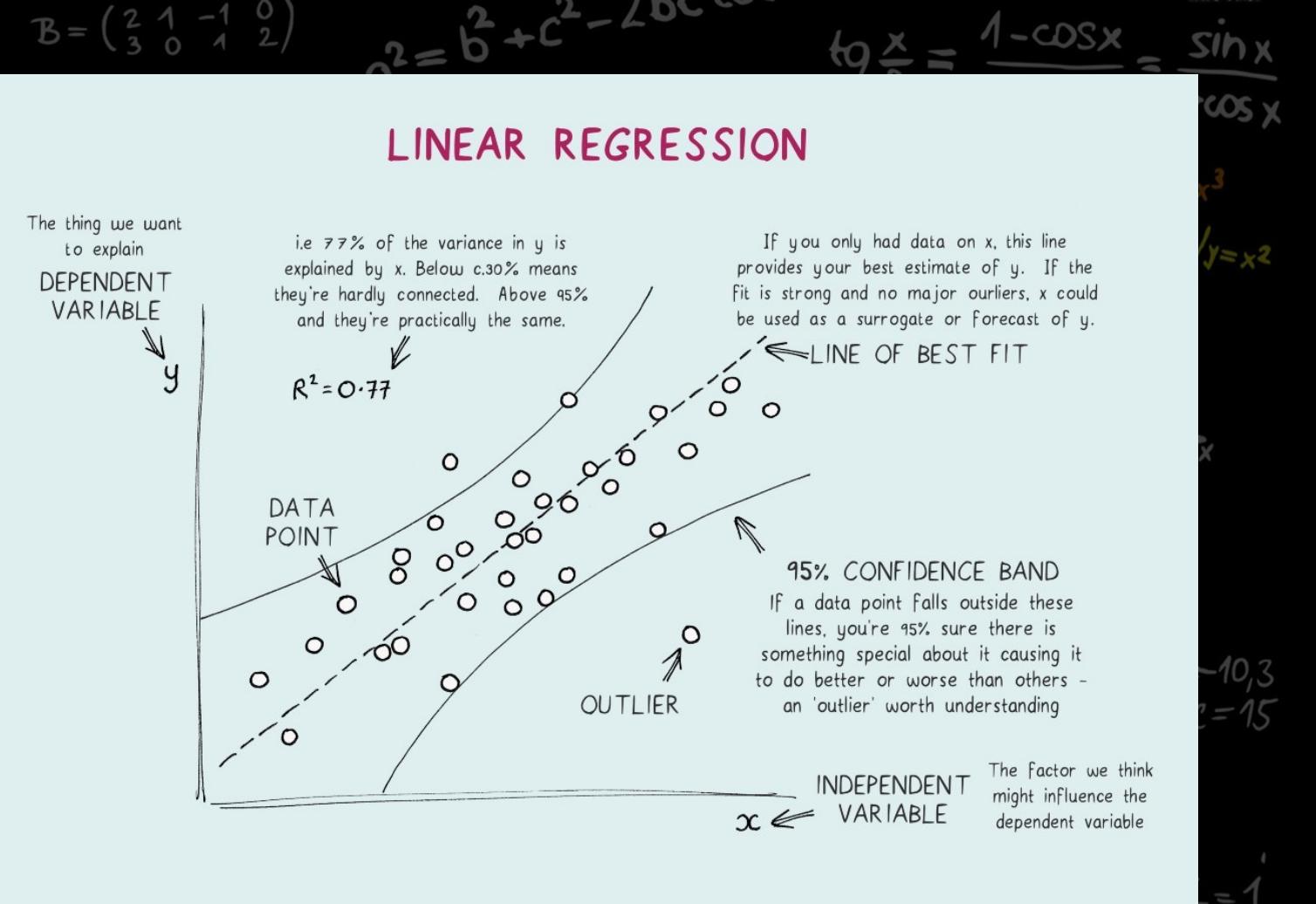
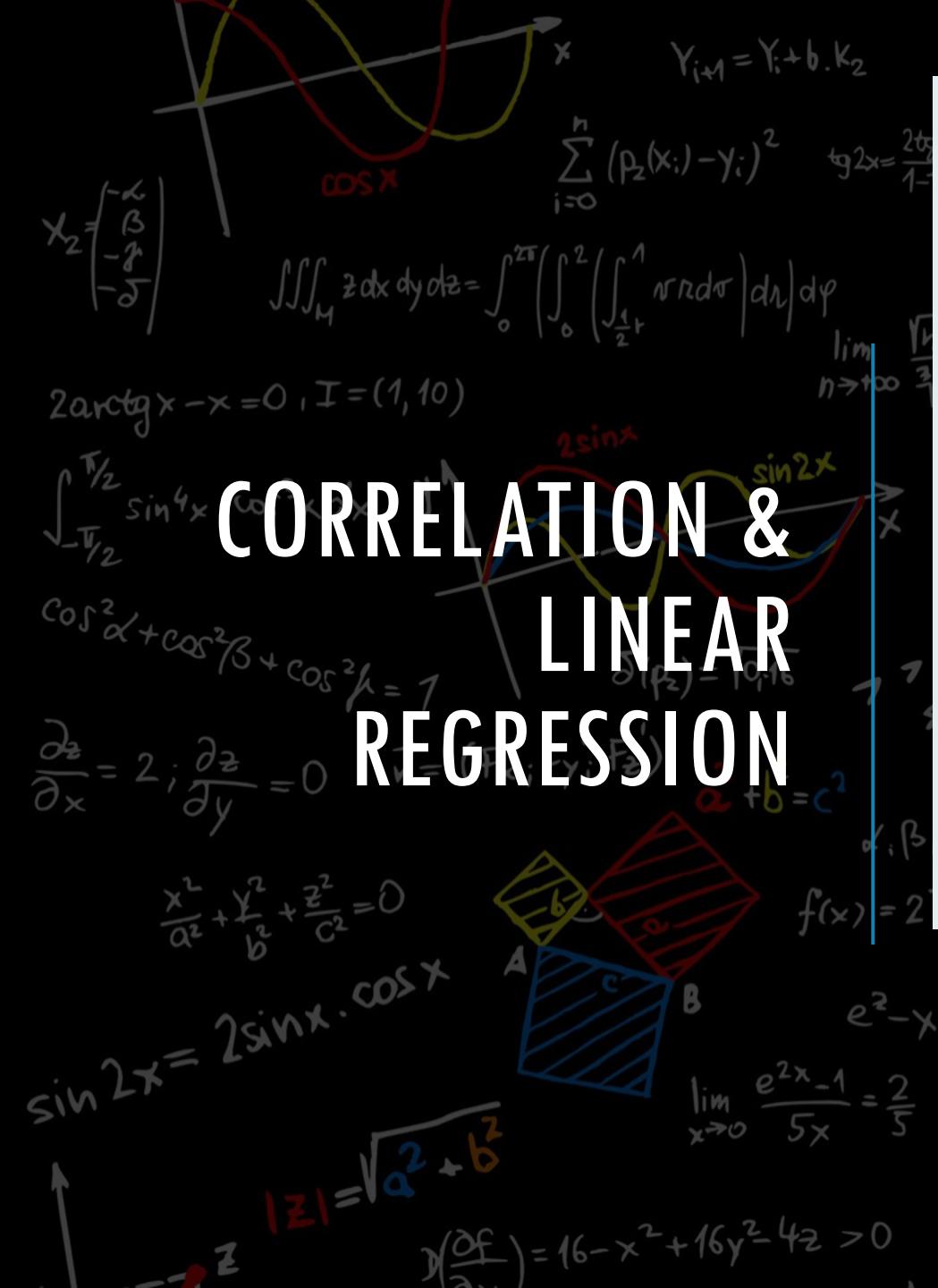
The outcome of linear regression includes estimating the intercept and the slope of the linear model.

Linear regression can then be used as a predictive model, whereby the model can be used to predict a y value for any given x. In practice, the model shouldn't be used to predict values beyond the range of the x values used to develop the model.

CORRELATION & LINEAR REGRESSION



CORRELATION & LINEAR REGRESSION



Mathematically, a regression line assumes the form of $y = ax + b$, with a the intercept (the location where it intersects an axis) and b the slope (the steepness of a line).

Both can be used to estimate an average rate of change. The greater the magnitude of the slope, the steeper the line and the greater the rate of change.

CORRELATION & LINEAR REGRESSION



Indeed, in R, the same syntax can be used to perform both ANOVA and linear regression!

ASSESSING MODEL ASSUMPTIONS: FORMAL TESTS TO ASSESS NORMALITY OF RESIDUALS

There are formal tests to assess the normality of residuals.

Common tests include Shapiro-Wilk, Anderson-Darling, Kolmogorov-Smirnov, and D'Agostino-Pearson.

Their results are dependent on sample size.

When the sample size is large, the tests may indicate a statistically significant departure from normality, even if that departure is small.

When sample sizes are small, they won't detect departures from normality!!

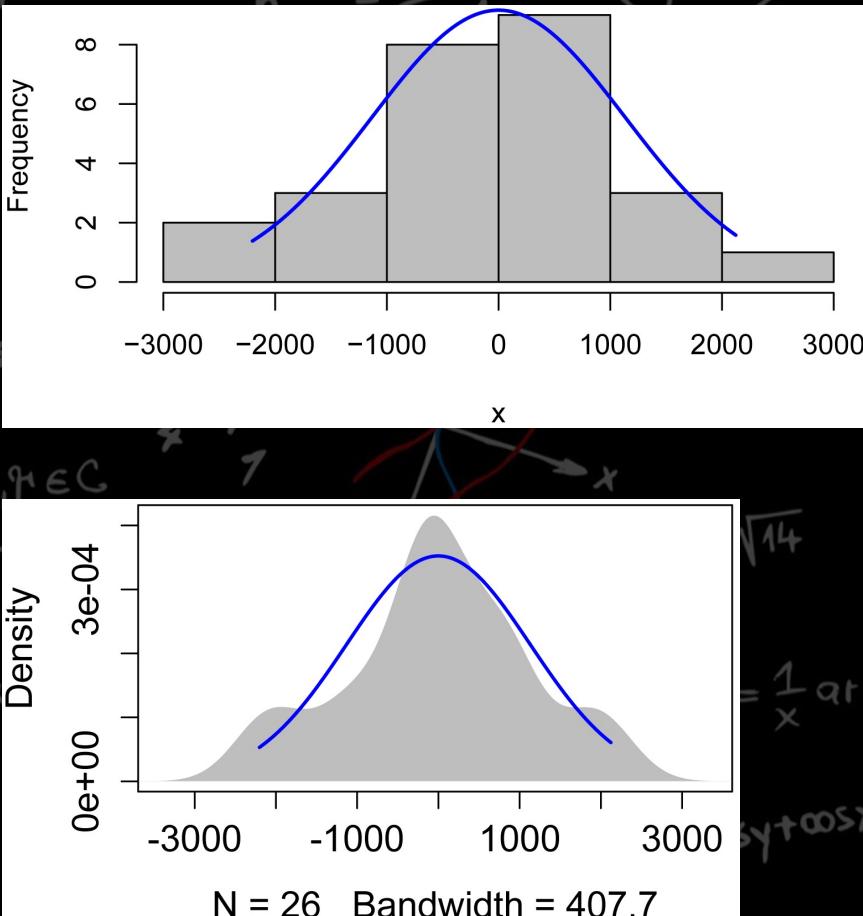
ASSESSING MODEL ASSUMPTIONS: FORMAL TESTS TO ASSESS NORMALITY OF RESIDUALS

There are no definitive guidelines as to what range of skew or kurtosis are acceptable for considering residuals to be normally distributed.

If we were forced to give advice for skewness calculations, be cautious if the absolute value is > 0.5, and consider it not normally distributed if the absolute value is > 1.

Some authors use 2 as a cut-off for normality, and others use a higher limit for kurtosis.

ASSESSING MODEL ASSUMPTIONS: FORMAL TESTS TO ASSESS NORMALITY OF RESIDUALS

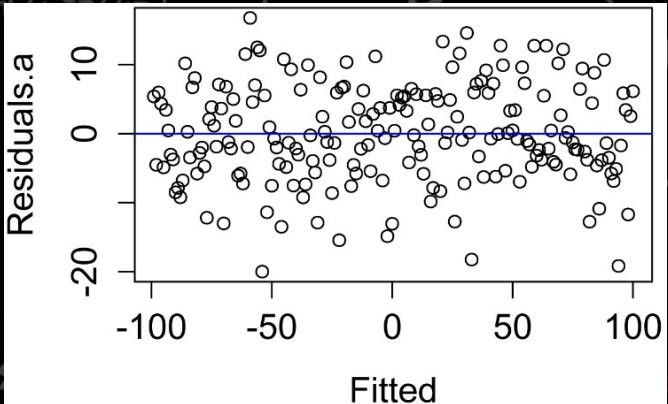


Usually, the best method to see if model residuals meet the assumptions of normal distribution and homoscedasticity are to plot them and inspect the plots visually.

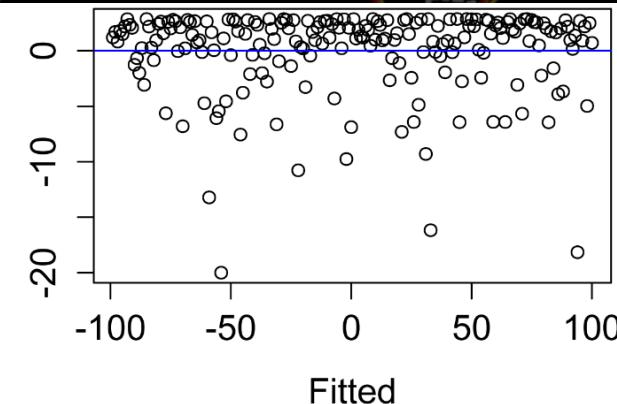
ASSESSING MODEL ASSUMPTIONS: FORMAL TESTS TO ASSESS NORMALITY OF RESIDUALS

Patterns in the plot of residuals versus fitted values can indicate a lack of homoscedasticity or that errors are not independent of fitted values.

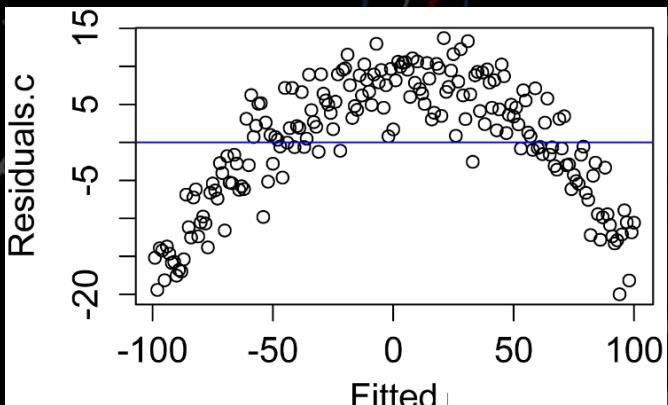
Normally distributed and homoscedastic



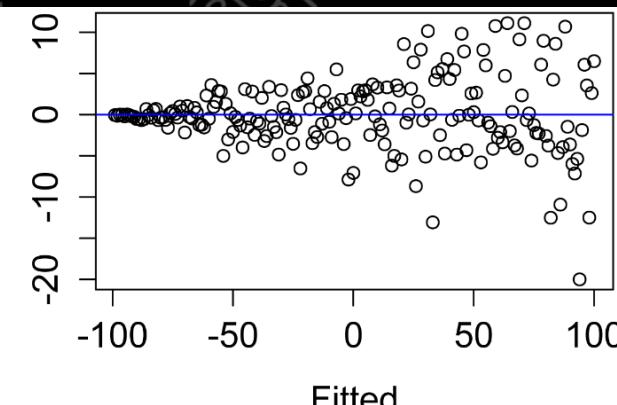
Non-normal distribution of residuals



Not independent of the fitted values



Heteroscedasticity



MODEL SELECTION: THE AKAIKE INFORMATION CRITERION

The Akaike information criterion (AIC) is a mathematical method for evaluating how well a model fits the data it was generated from. In statistics, AIC is used to compare different possible models and determine which one is the best fit for the data. AIC is calculated from:

→ The number of independent variables used to build the model.

→ The maximum likelihood estimate of the model (how well the model reproduces the data).

→ The best-fit model according to AIC is the one that explains the greatest amount of variation using the fewest possible independent variables.

MODEL SELECTION: THE AKAIKE INFORMATION CRITERION

AIC

AIC is most often used for model selection. By calculating and comparing the AIC scores of several possible models, you can choose the one that is the best fit for the data.

When testing a hypothesis, you might gather data on variables that you aren't certain about, especially if you are exploring a new idea. You want to know which of the independent variables you have measured explain the variation in your dependent variable.

A good way to find out is to create a set of models, each containing a different combination of the independent variables you have measured. These combinations should be based on:

→ Your knowledge of the study system – avoid using parameters that are not logically connected, since you can find spurious correlations between almost anything!

→ Your experimental design – for example, if you have split two treatments up among test subjects, then there is probably no reason to test for an interaction between the two treatments.

MODEL SELECTION: THE AKAIKE INFORMATION CRITERION

AIC

Once you've created several possible models, you can use AIC to compare them.

Lower AIC scores are better, and AIC penalizes models that use more parameters.

If two models explain the same amount of variation, the one with fewer parameters will have a lower AIC score and will be the better-fit model.

MODEL SELECTION: THE AKAIKE INFORMATION CRITERION

AIC

AIC determines the relative information value of the model using the maximum likelihood estimate and the number of parameters (independent variables) in the model.

$$AIC = 2K - 2\ln(L)$$

K is the number of independent variables used and L is the log-likelihood estimate (a.k.a. the likelihood that the model could have produced your observed y-values).

To compare models using AIC, you need to calculate the AIC of each model. If a model is more than 2 AIC units lower than another, then it is considered significantly better than that model.