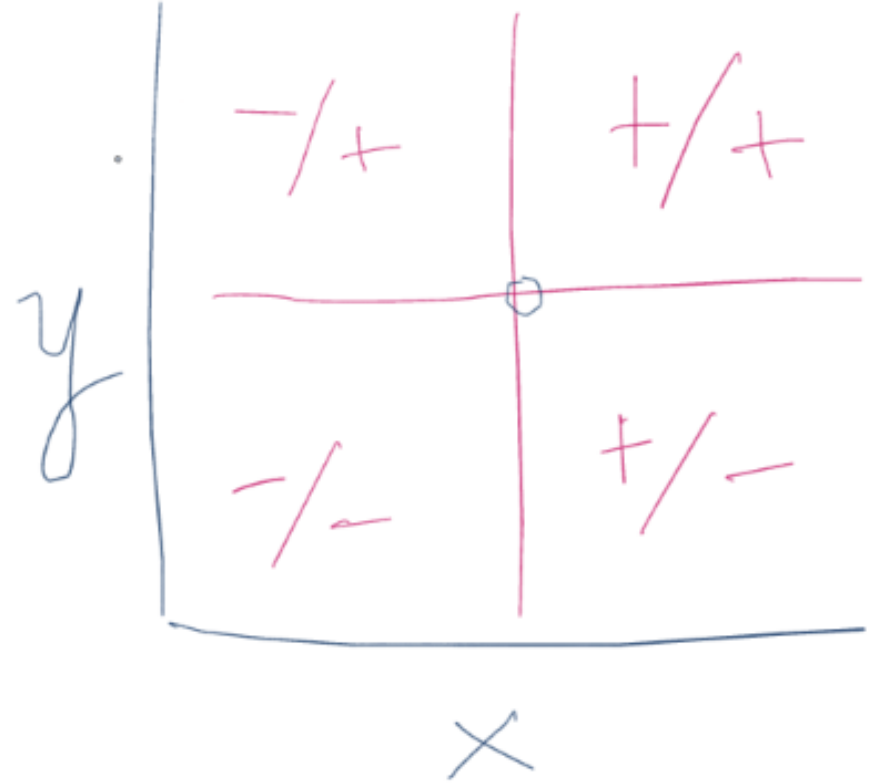# Module 4A
# Supervised Machine Learning

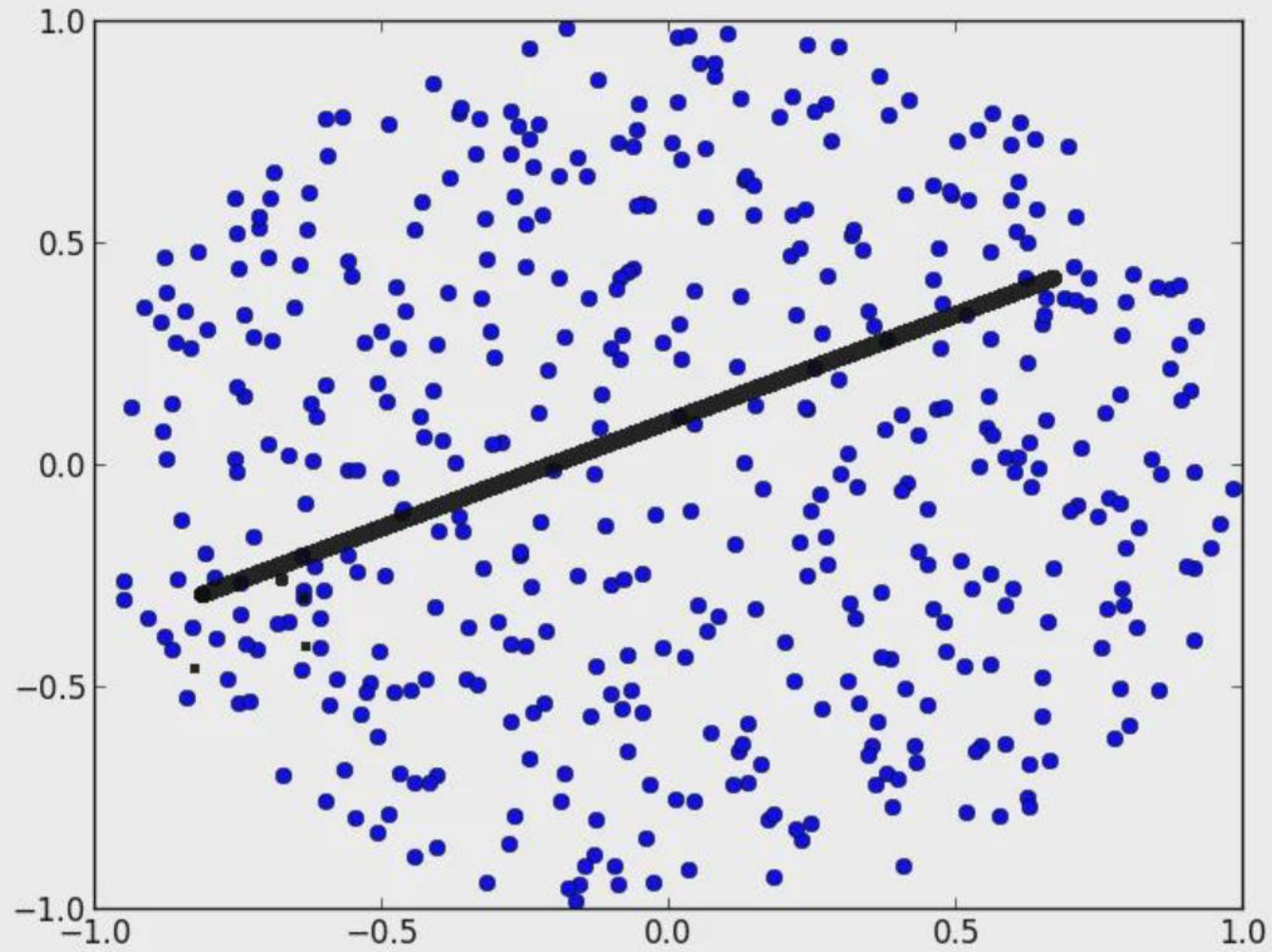Different flavors of REGRESSION and General Linear Models

# **Review:** Correlation:

- Measures the amount/degree of linear association between two **numerical** variables

- Estimate the degree to which variables **covary**
  - With no attempt to interpret the causality of the association

Example: arm length and leg length covary together (individuals with longer arms often have longer legs) but they are influenced by other underlying variables **not** each other (longer legs do not cause longer arms)
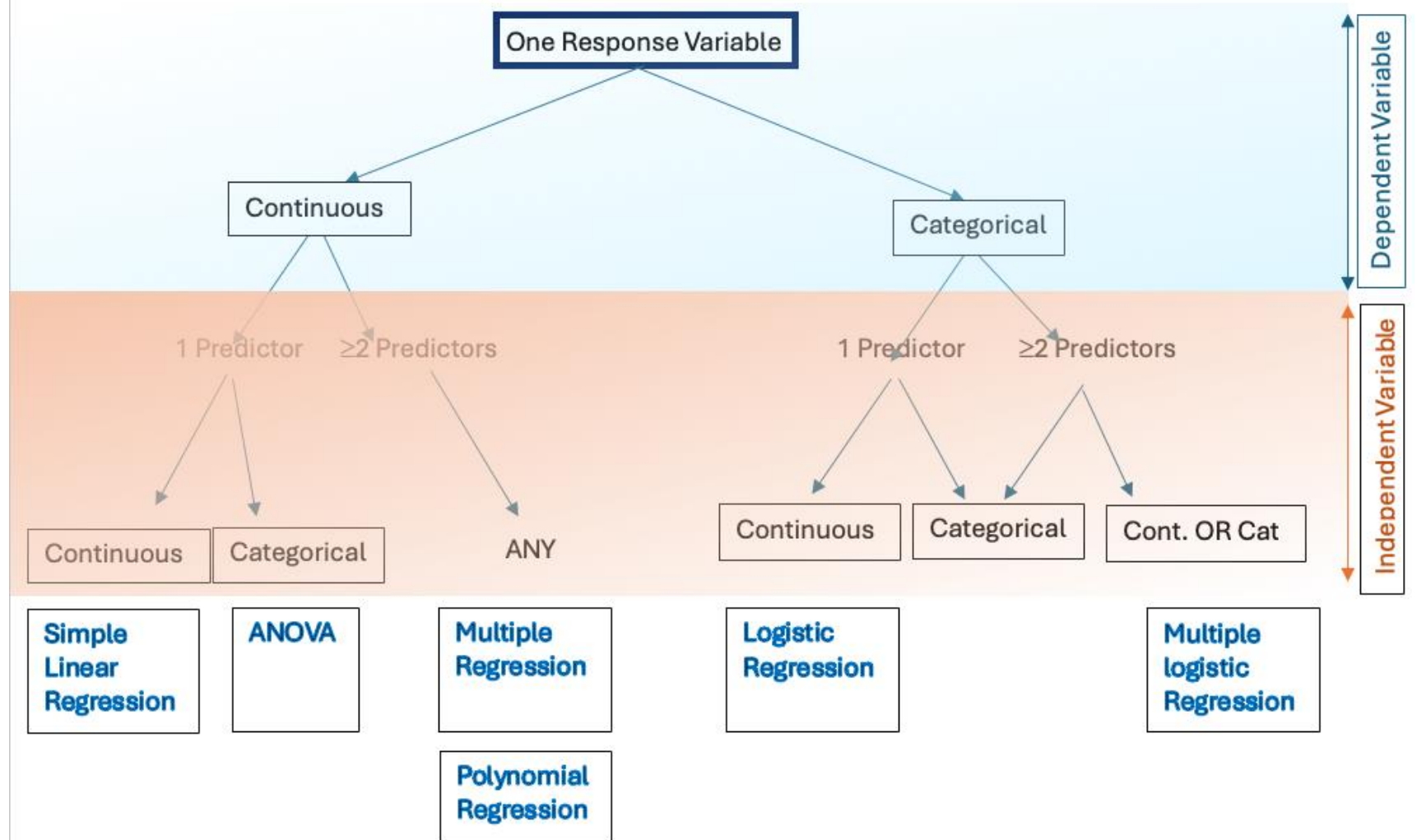
# Scientists be like

# Regression:

- Statistics is about <u>prediction</u>

- Used to **predict** value of one numerical variable from the value of another
  - predicting dependent/response variable, Y from independent/predictor X

- **Linear** regression assumes that the relationship between X and Y can be described by a line
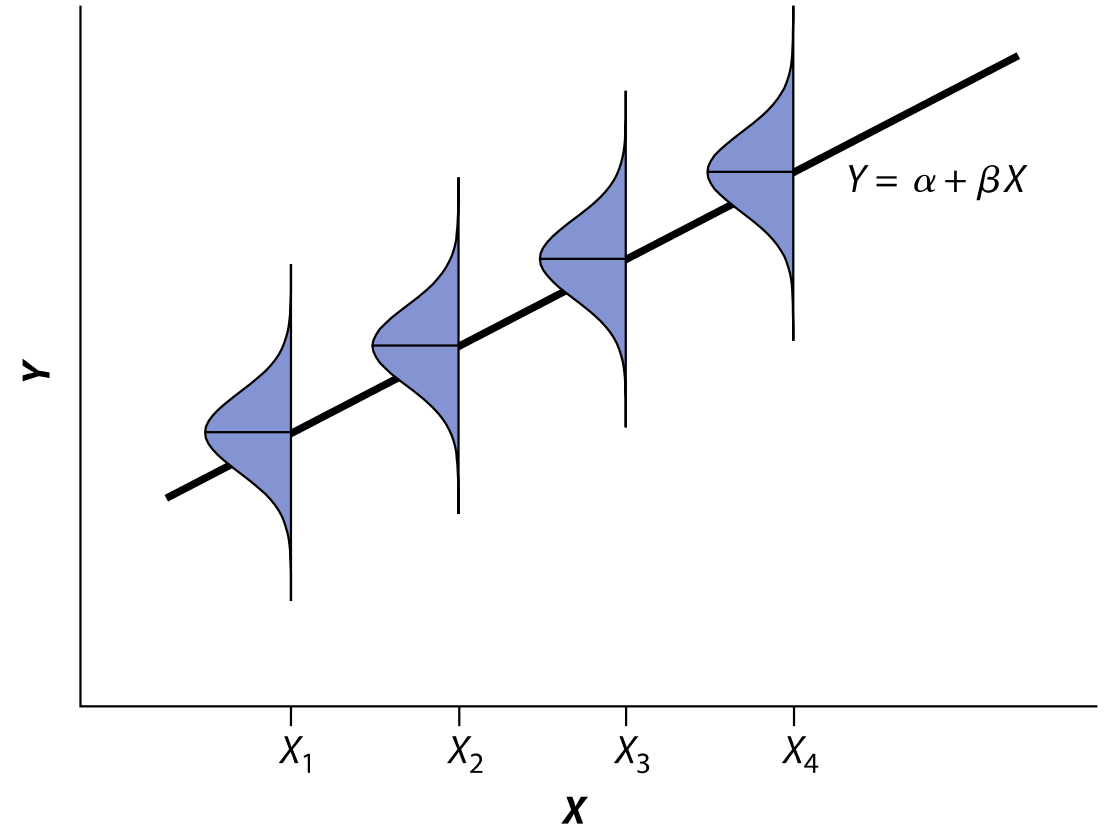  - Fits a straight line to a (messy) scatterplot

**Ambient (room) temperature** ⟶ **plant growth**

<span style="color:#7a3b10"><u>A common point of confusion:</u> Pearson's correlation is the slope of the best-fit line **_after_ _standardizing_ both** variables. The regression slope is the **unstandardized** version that tells you how much Y changes per unit X.</span>
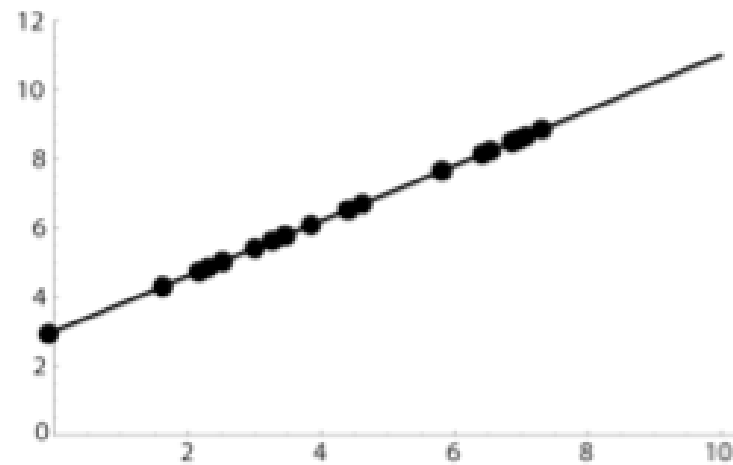
- Other kinds:
  - Lasso
    - Variable selection (weighting a predictor variable by 0)
  - Ridge
    - Allows analysis even in the face of **collinearity**
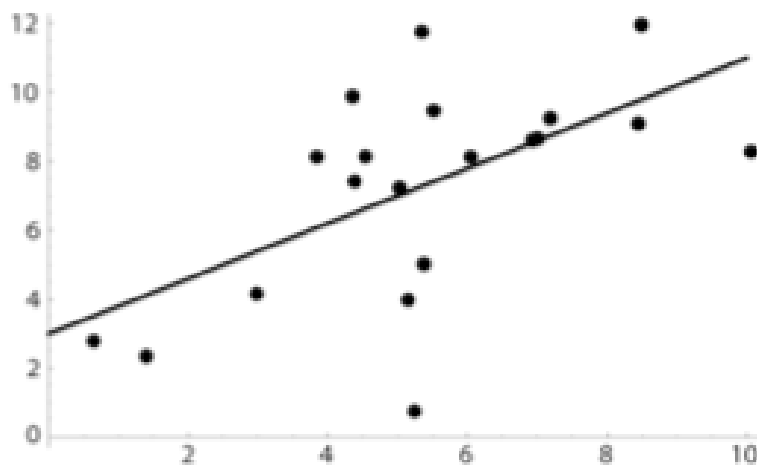
# Regression:

- **Linear** regression assumes that the relationship between X and Y can be described by a line
  - Fits a straight line to a (messy) scatterplot

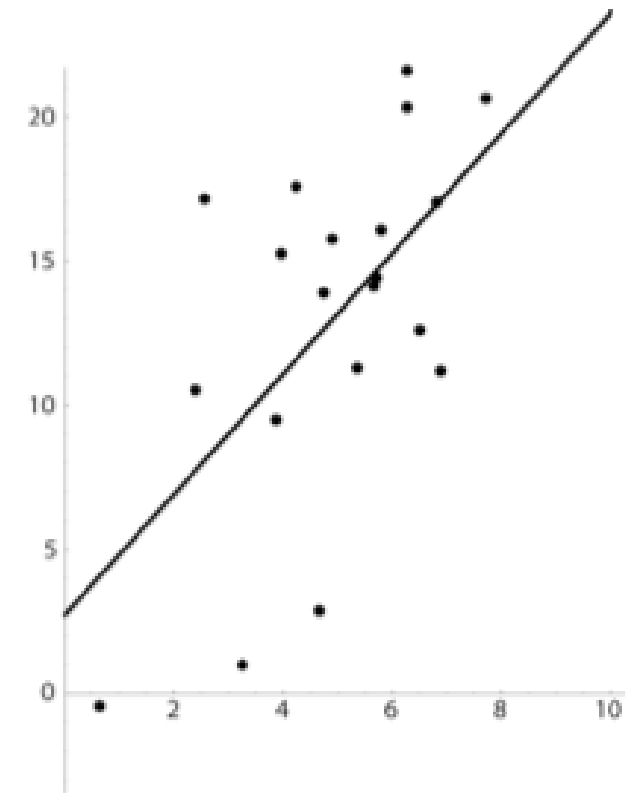- **Homoscedasticity**: Y is normally distributed with equal variance for all values of X

$$Y = \alpha + \beta X$$

**correlation vs regression**



$r = 1$;  $Y = 3 + 0.8X$          $r = 0.6$;  $Y = 3 + 0.8X$          $r = 0.6$;  $Y = 3 + 2X$

← Different correlation; same slope          Same correlation; different slope →
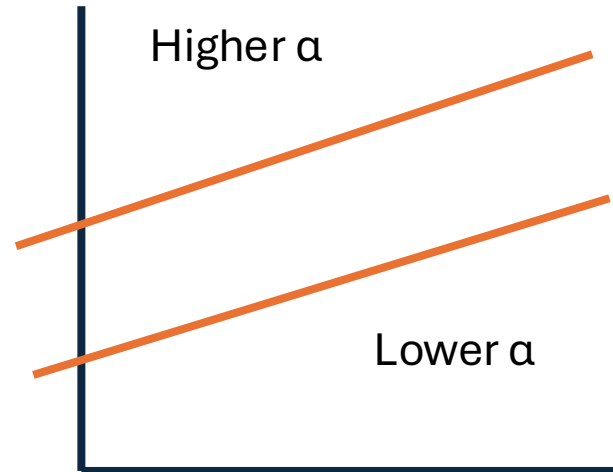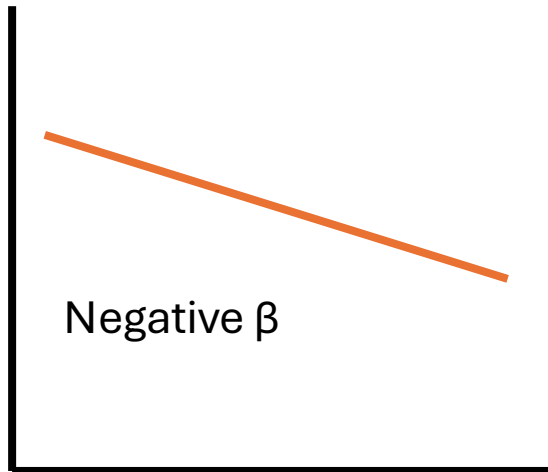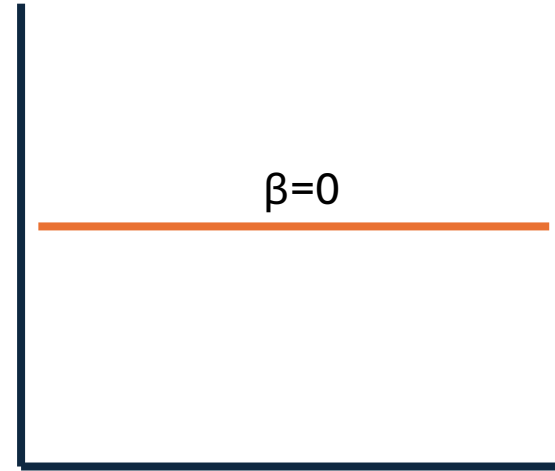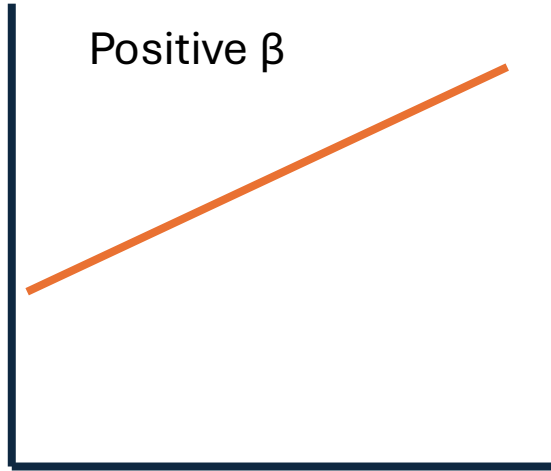
# The parameters of linear regression

$$Y = \boxed{\alpha} + \enclose{circle}{\beta}X + \varepsilon_1$$

intercept

Slope*

# Regression Overview



Positive β

β=0

Negative β

Higher α

Lower α

# Estimating a regression line

$$Y = a + bX + \varepsilon_1$$

## Quick Review:

$$Z = \frac{\overline{Y} - \mu}{\sigma_{\overline{Y}}} = \boxed{\frac{\overline{Y} - \mu}{\sigma / \sqrt{n}}}$$

$$t = \frac{\overline{Y} - \mu}{SE_{\overline{Y}}} = \boxed{\frac{\overline{Y} - \mu}{s / \sqrt{n}}}$$

$$\text{F-value} = \frac{\boxed{\mathbf{MS_{group}}}}{\boxed{\mathbf{MS_{error}}}} \quad \begin{array}{l} \textbf{SIGNAL} \\[1em] \textbf{NOISE} \end{array}$$

$$r = \frac{\boxed{\sum (X - \overline{X})(Y - \overline{Y})}}{\boxed{\sqrt{\sum (X - \overline{X})^2} \sqrt{\sum (Y - \overline{Y})^2}}} = \frac{Co \, \mathrm{var} \, iance(X,Y)}{s_x s_y}$$

# (Ordinary) Least Squares:

- Best fitting line through a scatterplot
  - Line that minimized spread of y values

- Minimize $SS_{residuals}$
  - Measurement of how much the line's predicted $y_i$ deviate from actual data values



Large deviations | Smaller deviations | Smallest deviations

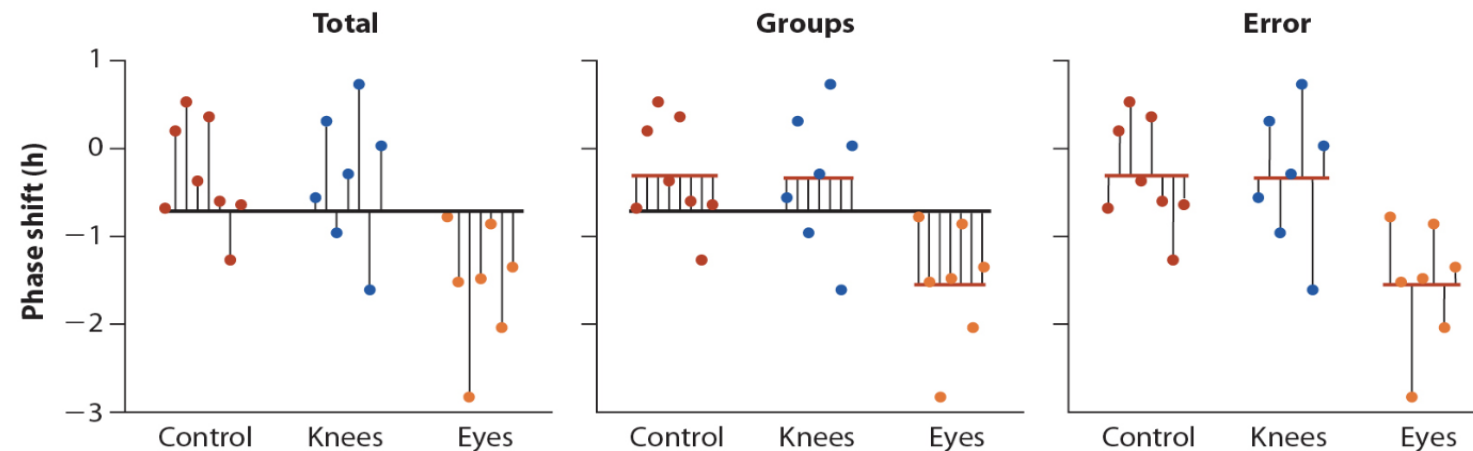**In Linear regression, the line of best fit (OLS) is the GRAND Mean in ANOVA**



Figure 20.1: Whitock and Schluter, Fig 15.1.2 – Illustrating the partitioning of sum of squares into $MS_{group}$ and $MS_{error}$ components.
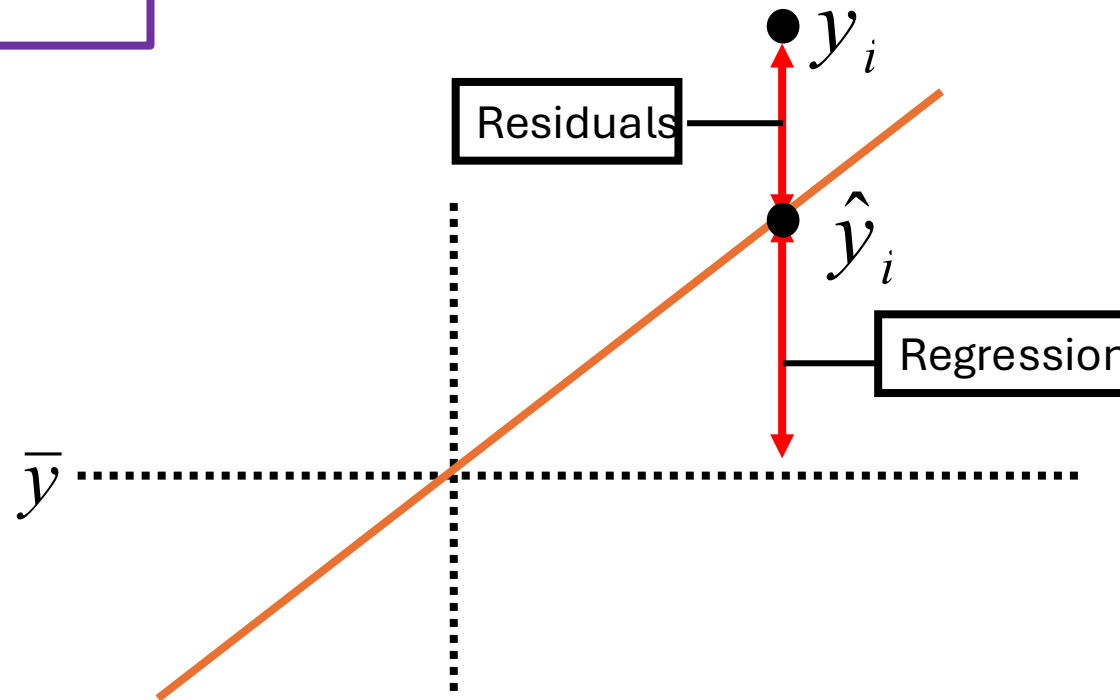
# Regression Overview

## Least Squares:

- What are the elements of this equation?

$$SS_{residual} = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

$$\hat{y}_i = a + bx_i$$



$y_i$

Residuals

$\hat{y}_i$

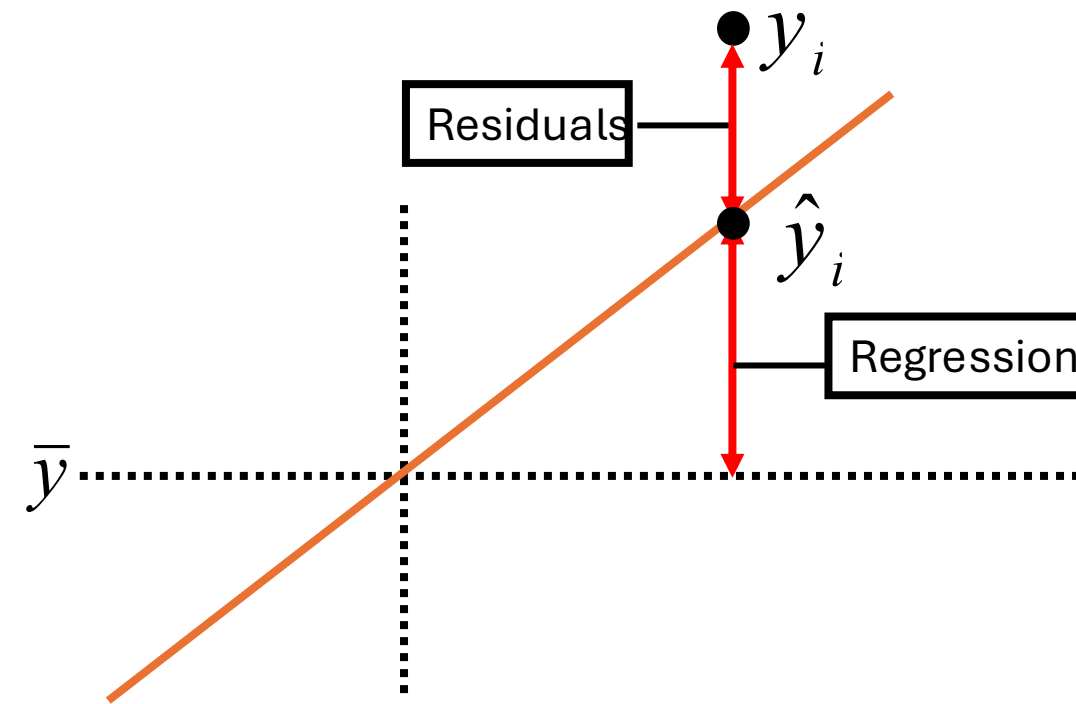Regression

$\bar{y}$

- Residuals:
  - Residuals measure the scatter of points above and below the least squares regression line
  - $MS_{residual}$ is the variance of the residuals, $residual = Y_i - \hat{Y}_i$

$$MS_{residual} = \frac{\mathring{a}\,(Y_i - \hat{Y}_i)^2}{n-2}$$



  - $MS_{regression}$ is the variance of the regression

$$MS_{regression} = \frac{\sum(\hat{Y}_i - \bar{Y})^2}{n-2}$$

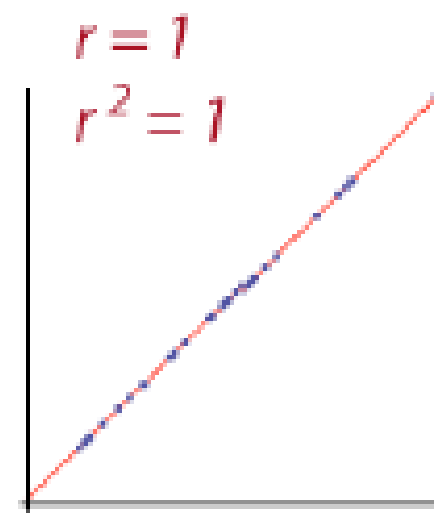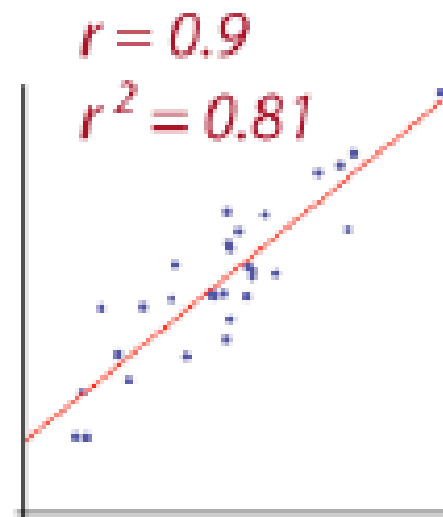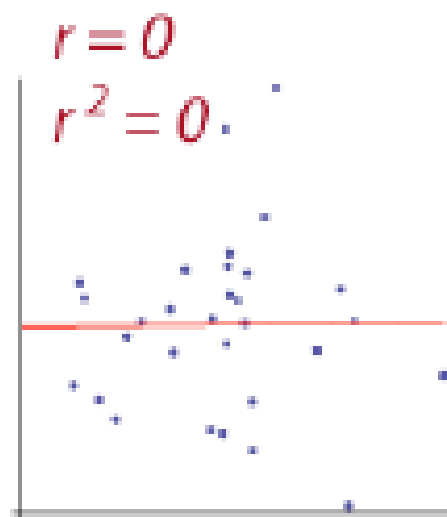  - Coefficient of determination ($r^2$) = SSR/SST

# $R^2$ predicts the amount of variance in Y explained by the regression line

We saw this in ANOVA where $R^2$ gave 'precision' of model (i.e. Ability of the model to explain variation)
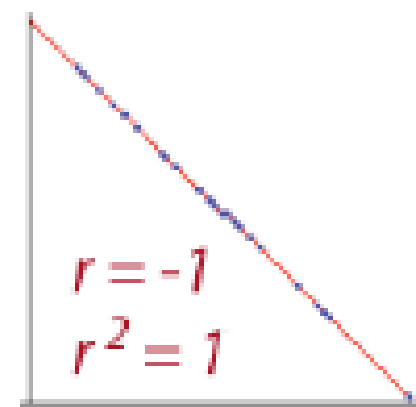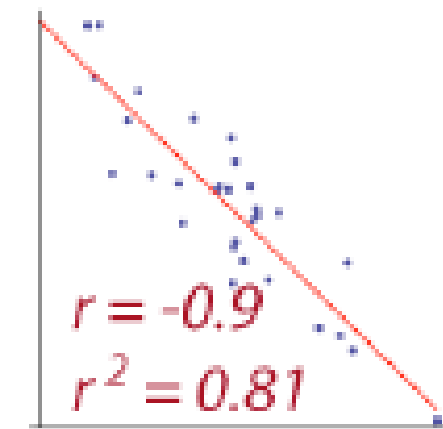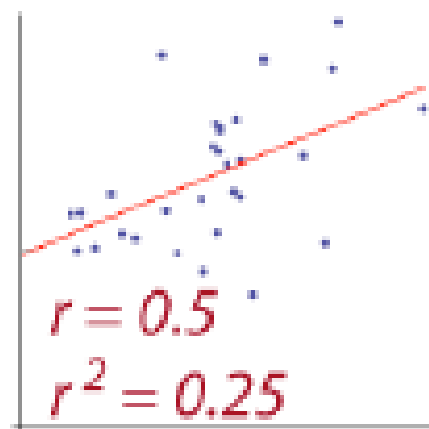
- The coefficient of determination

- Sometimes written as $r^2$

- Square of the correlation coefficient, r

$$R^2 = \frac{SS_{regression}}{SS_{Total}}$$

$r = 0$
$r^2 = 0$

$r = 0.9$
$r^2 = 0.81$

$r = 1$
$r^2 = 1$

$r = 0.5$
$r^2 = 0.25$

$r = -0.9$
$r^2 = 0.81$

$r = -1$
$r^2 = 1$

Y

Y

X

X

X

# Best estimate of slope:

b =  Sum of cross products

Sum of squares of X

$$b = \frac{\sum_{i=1}^{n} (X_i - \overline{X})(Y_i - \overline{Y})}{\sum_{i=1}^{n} (X_i - \overline{X})^2}$$

# Best estimate of slope:

b = $\dfrac{\text{Sum of cross products}}{\text{Sum of squares of X}}$

$$b = \frac{\overset{n}{\underset{i=1}{\sum}} (X_i - \overline{X})(Y_i - \overline{Y})}{\overset{n}{\underset{i=1}{\sum}} (X_i - \overline{X})^2}$$

$$r = \frac{\sum (X - \overline{X})(Y - \overline{Y})}{\sqrt{\sum (X - \overline{X})^2} \sqrt{\sum (Y - \overline{Y})^2}} = \frac{Co\operatorname{var}iance(X,Y)}{s_x s_y}$$

Denominator normalizes based on **both** X and Y variables

Denominator ONLY normalizes based on X (Independent/Explanatory variable)

18

# Finding **a**:

$$\overline{Y} = a + b\overline{X}$$

**OR**

$$a = \overline{Y} - b\overline{X}$$