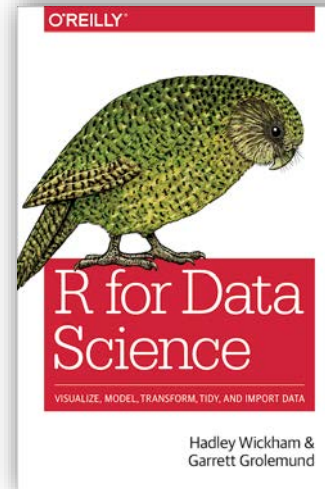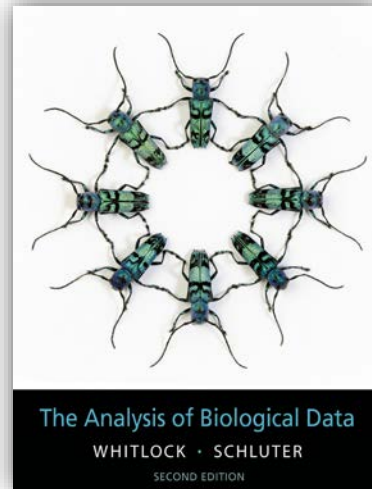# Data Science in Bioinformatics

Palle Villesen & Thomas Bataillon
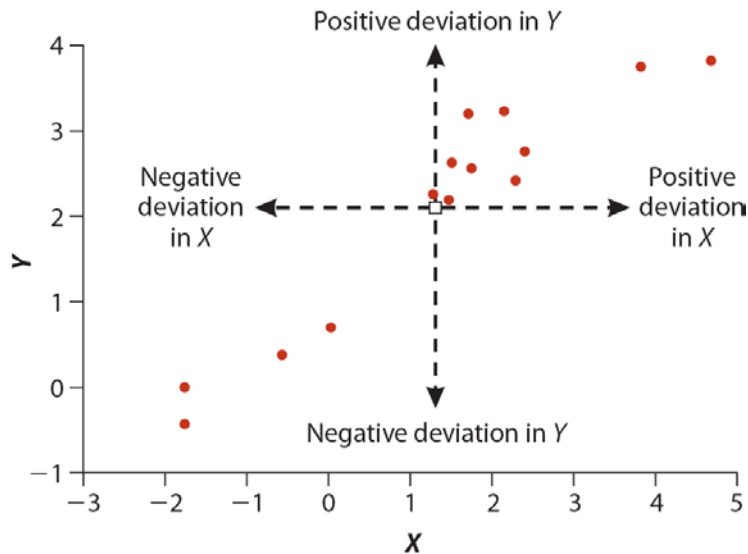
# Outline for week 11

- Chapter 16+17
  - Open discussion & Exercises
- Thursday
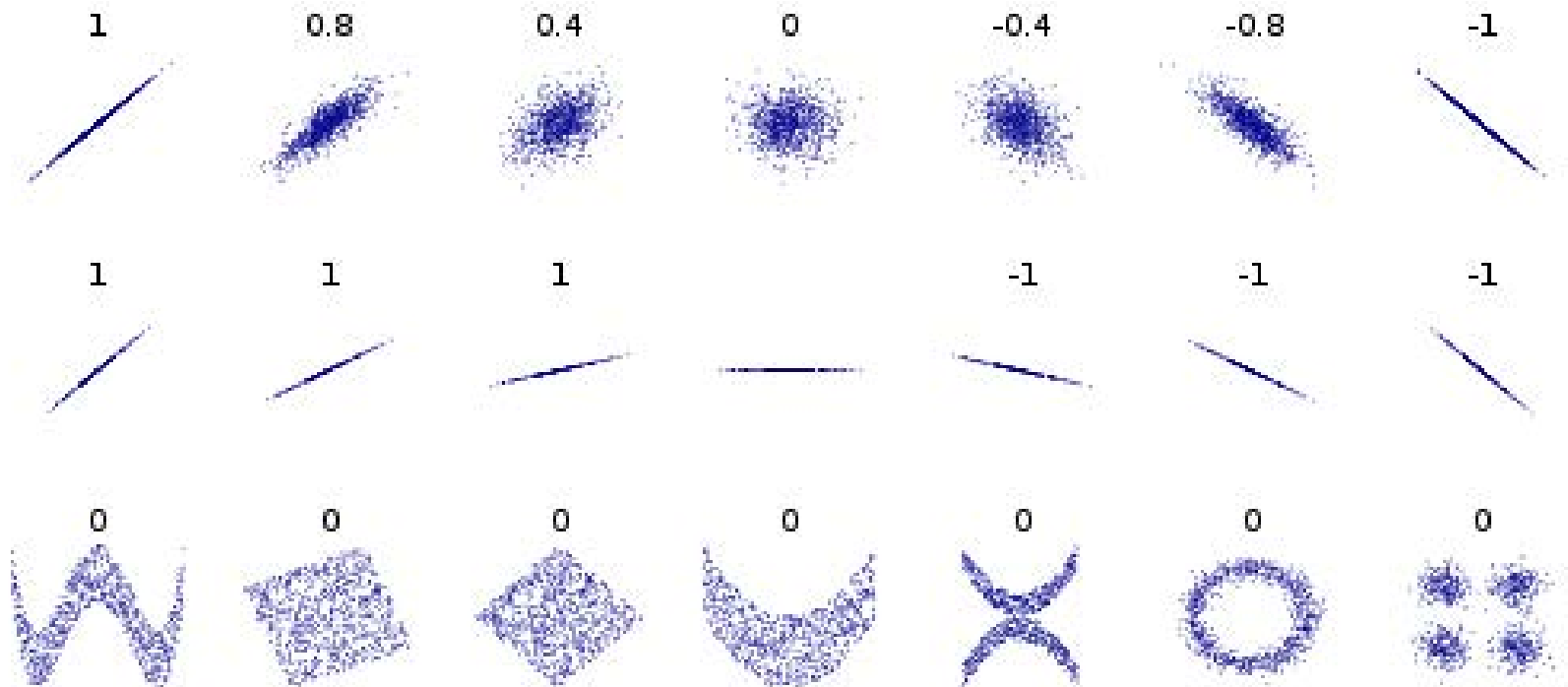  - Who will be responsible

# Correlation and regression

- Is two variables independent or not?
- How much of the variation in Y is explained by the variation in X?
- What is the standard error?
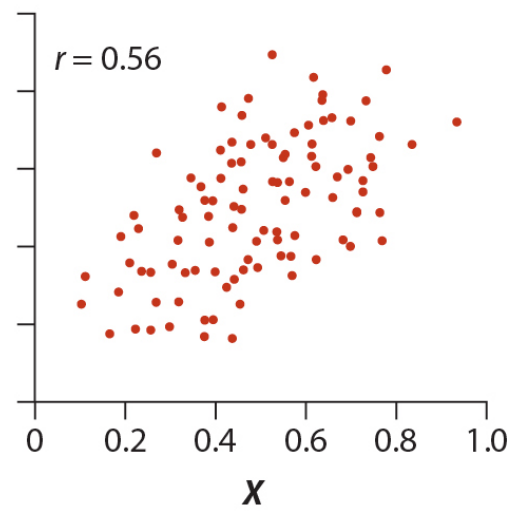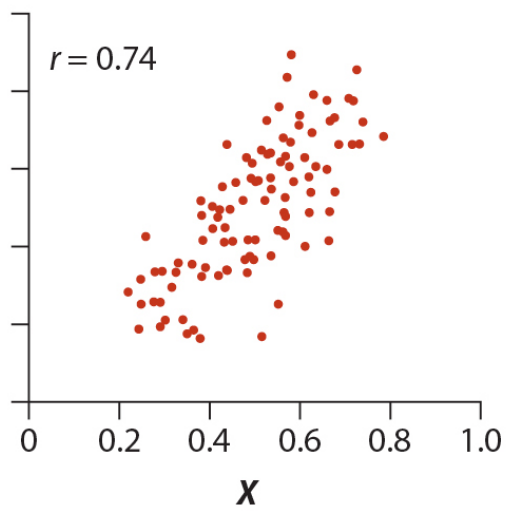- How strong is the effect?
- How do we test it?

# Pearson correlation

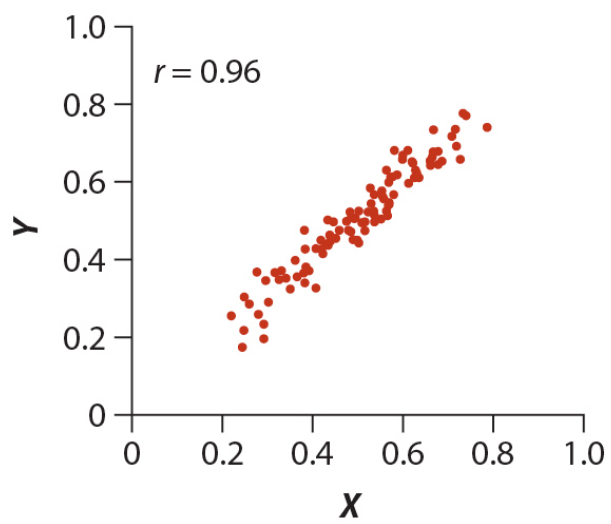$$r_{xy} = \frac{\displaystyle\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{(n-1)s_x s_y} = \frac{\displaystyle\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\displaystyle\sum_{i=1}^{n}(x_i - \bar{x})^2 \sum_{i=1}^{n}(y_i - \bar{y})^2}},$$

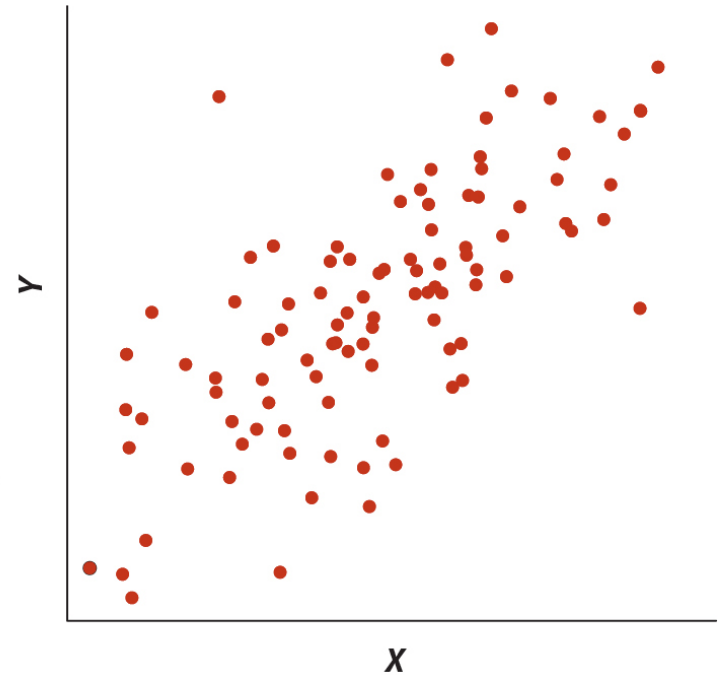# Correlation

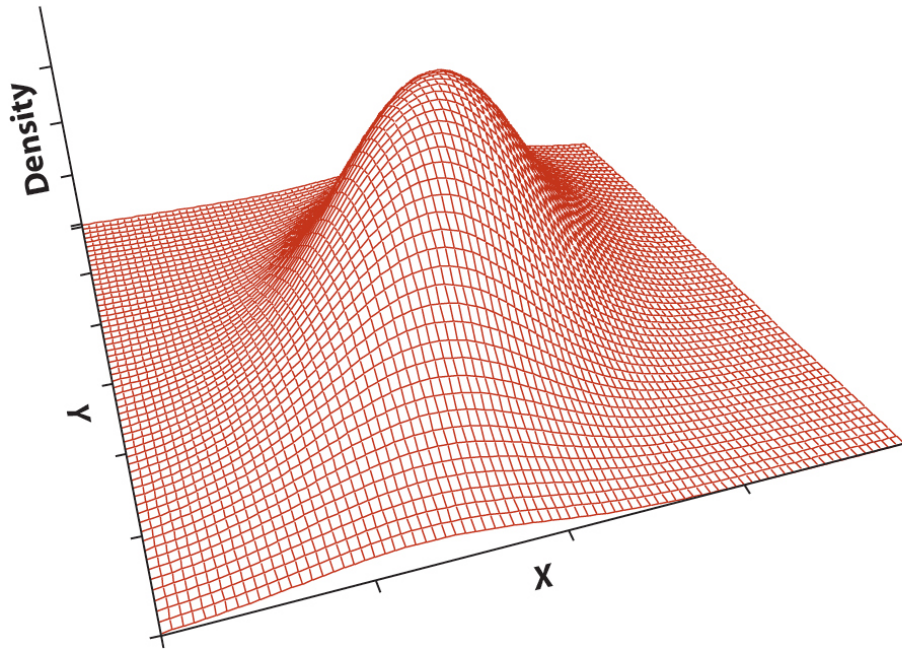# Sample size and correlation

- $SE_r = \sqrt{\dfrac{(1-r^2)}{n-2}}$

- **Confidence interval**
- Correlation coefficient is not normally distributed
- We convert it to "z" that is normal – calculates 95% CI and converts back from z to r
- Or we bootstrap the confidence interval

# Testing correlation

- H0: the population correlation is 0
- HA: the population correlation is not 0
- t = r/SE, Use t distribution with n-2 df
- cor.test
- **Alternative way**
- Permutation test
  - Sample shuffle x and y, calculate r
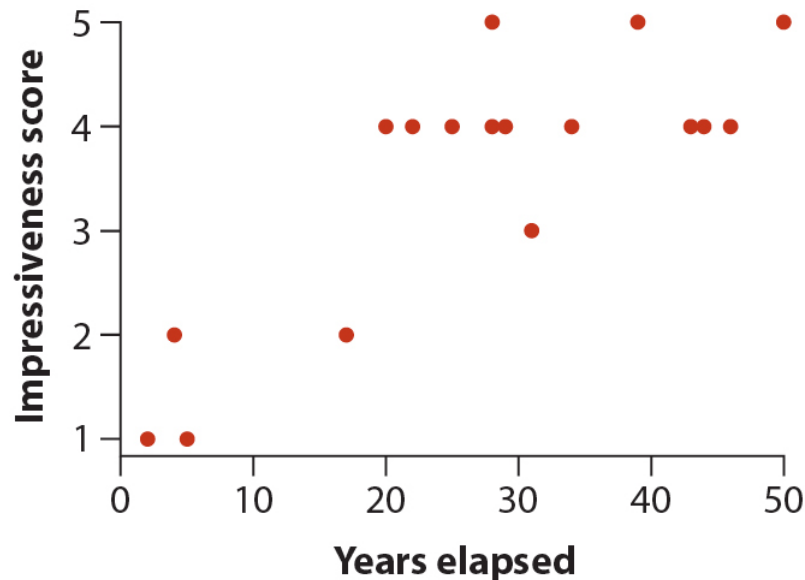  - How many are more extreme than observed r?
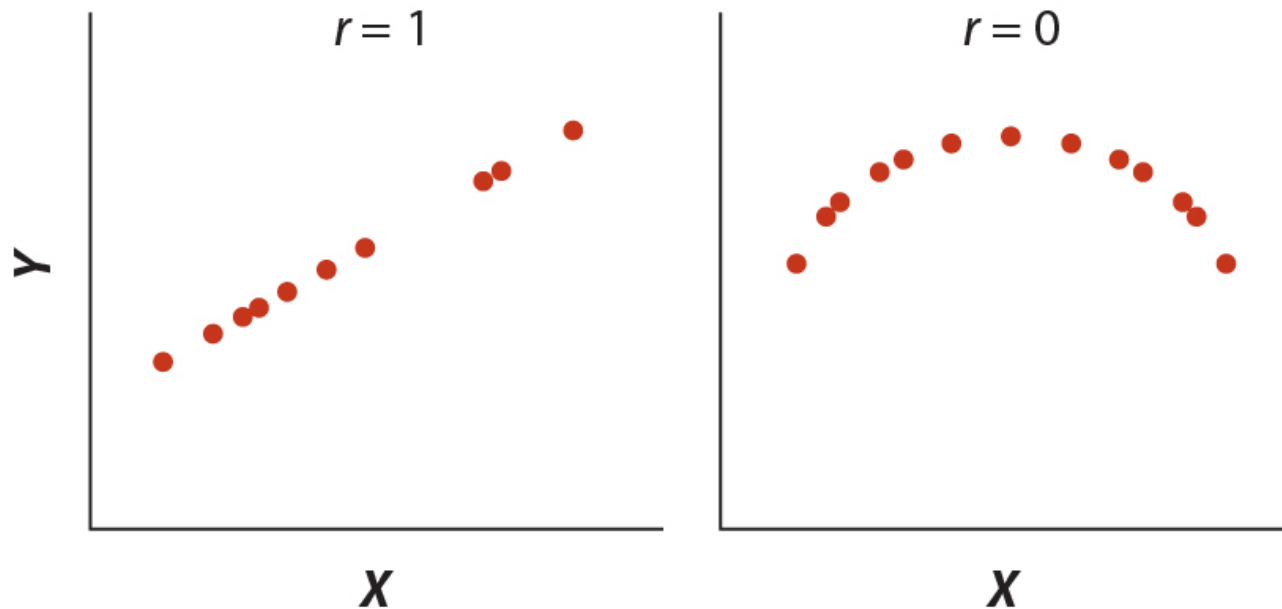
# Assumptions
# bivariate normal distribution

# Other stuff

- Rank correlation
  - No assumption on distribution of x and y
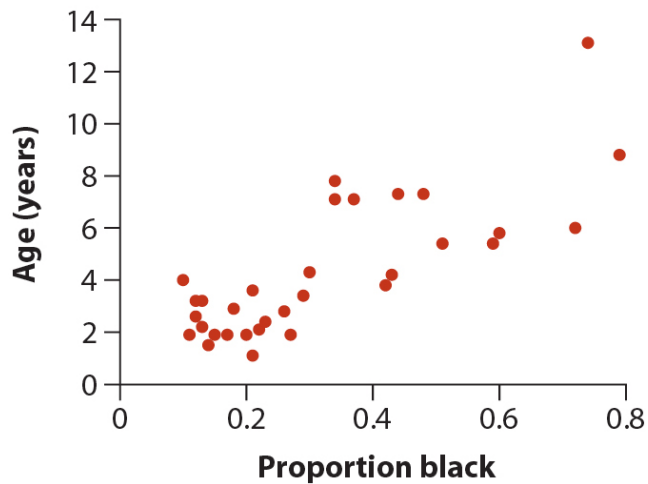  - It ranks all data and works on ranks

# Low correlation?

# Linear regression

- $Y = \beta_o + \beta_1 X \ldots$
- $Y = \beta_o + \beta_1 X + \varepsilon$

- It is a **prediction** method
- Also called statistical learning (or machine learning)
  - Prediction and/or inference
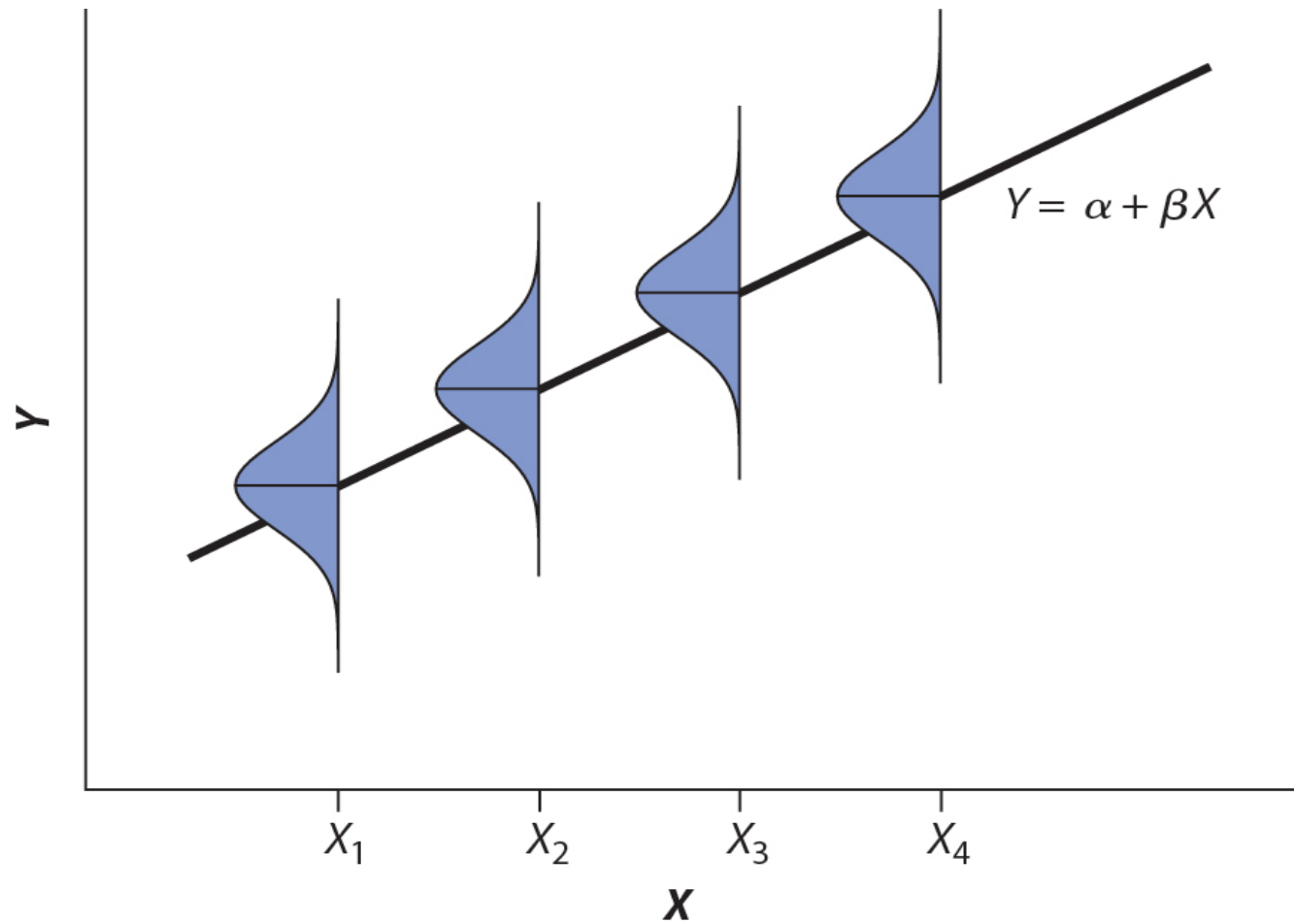- Supervised learning (we need known data to build the model)
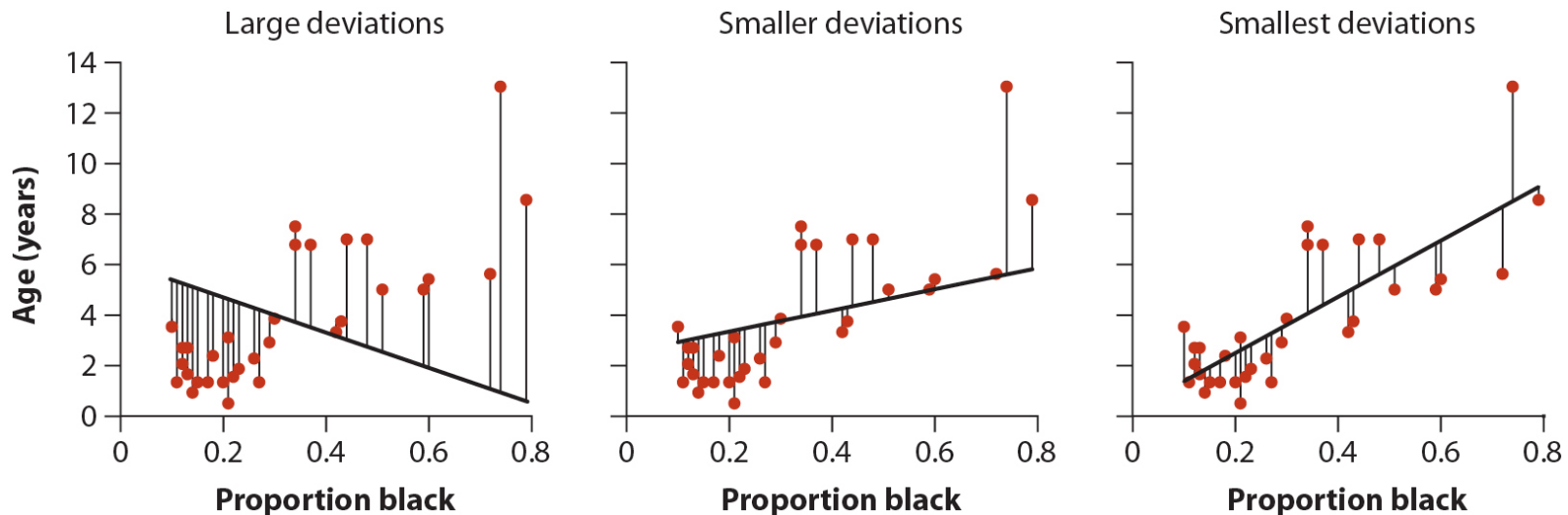
# Predict lion age from nose

# Essentials of regression

- Assumptions
- Tests & Measure of Fit
  - Parametric
  - Resampling methods
- Validation / Inspection
  - Visual
  - Test on residuals
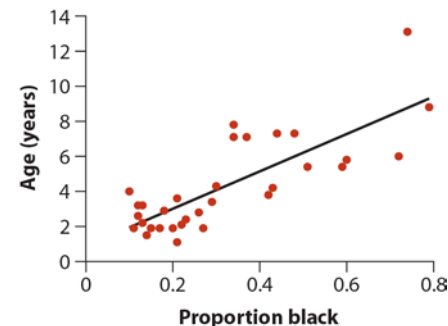  - Influence of outliers, colinearity

# Assumptions…



$Y = \alpha + \beta X$

# A visual on least squares...

Large deviations        Smaller deviations        Smallest deviations



$$\hat{\alpha} = \bar{y} - \hat{\beta}\,\bar{x},$$

$$\hat{\beta} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2}$$

$$= \frac{\mathrm{Cov}(x, y)}{\mathrm{Var}(x)}$$

$$= r_{xy}\frac{s_y}{s_x}.$$

minimum prediction error

# The model for lions

- Age = 0.88 + 10.65 * proportion black

- Standard error of slope

- Confidence interval for the slope

- Also possible by bootstrapping

- CI for slope: 7.56 < slope < 13.73

# What can we do

- We have shown that nose patterns change with age (inference)

- We can estimate how much it changes

- We can calculate how much of the variation is explained by this linear relation ship ($R^2$)

- We can predict the age of lions

# R² = SS(regression) / SS(total)

$$\text{TSS} = \sum_{i=1}^{n} \left( y_i - \bar{y} \right)^2$$
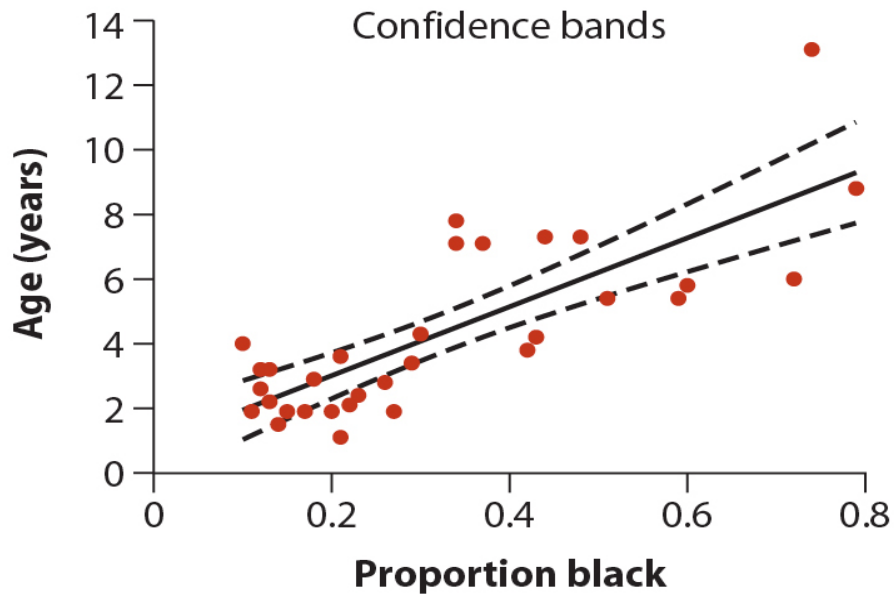
$$\text{ESS} = \sum_{i=1}^{n} \left( \hat{y}_i - \bar{y} \right)^2.$$

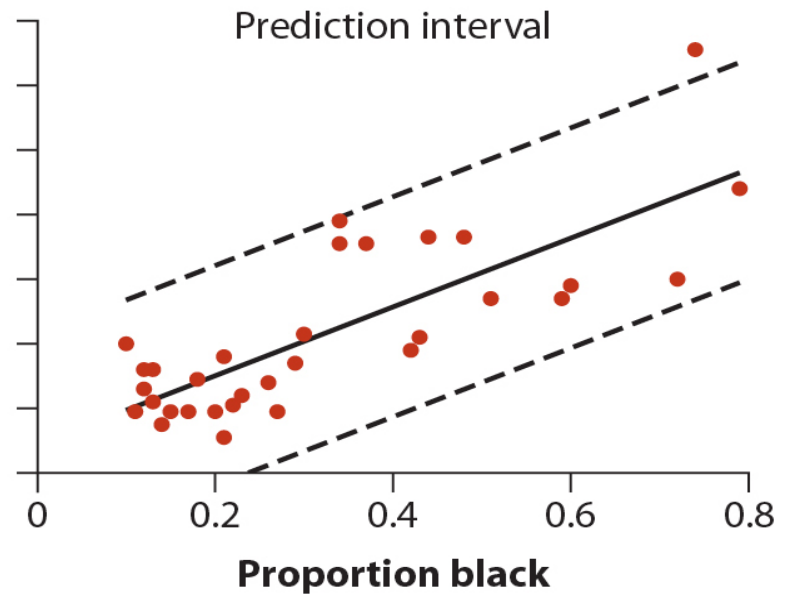$$RSS = \sum_{i=1}^{n} (\varepsilon_i)^2 = \sum_{i=1}^{n} (y_i - (\alpha + \beta x_i))^2$$

# Testing the slope

- H0: slope = 0
- HA: slope != 0
- T.test – but for any practical purposes R will do an ANOVA test
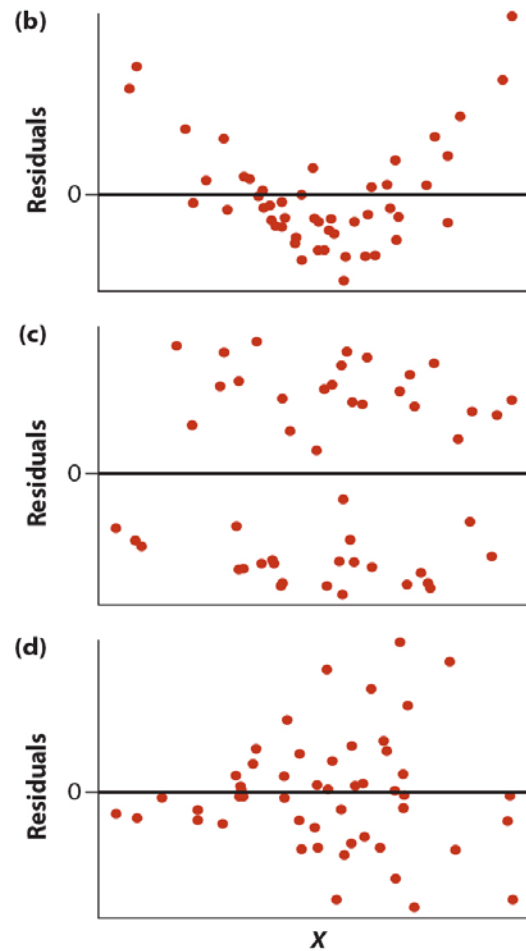- Permutation test? How?
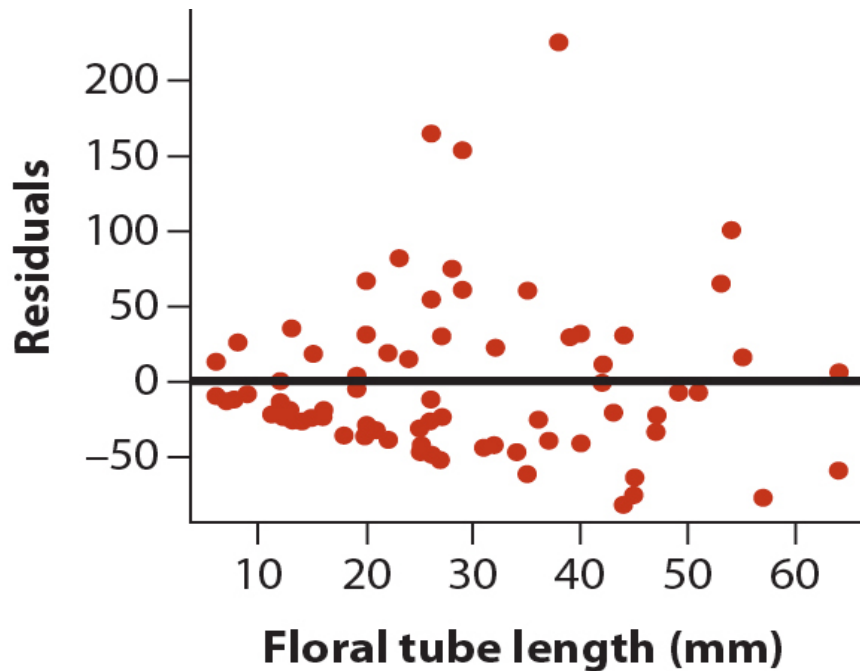
# Predictions



Precision of the predicted **mean**
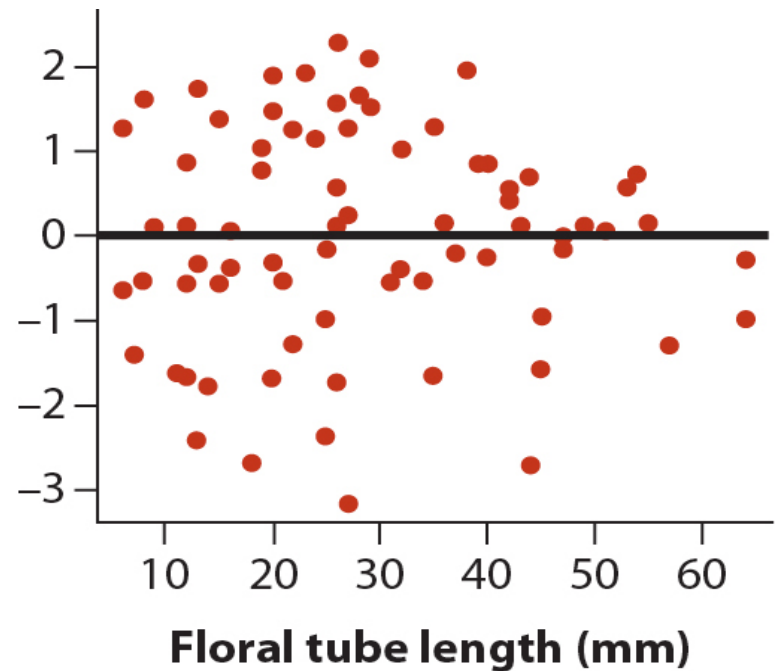
Precision of the predicted **value**
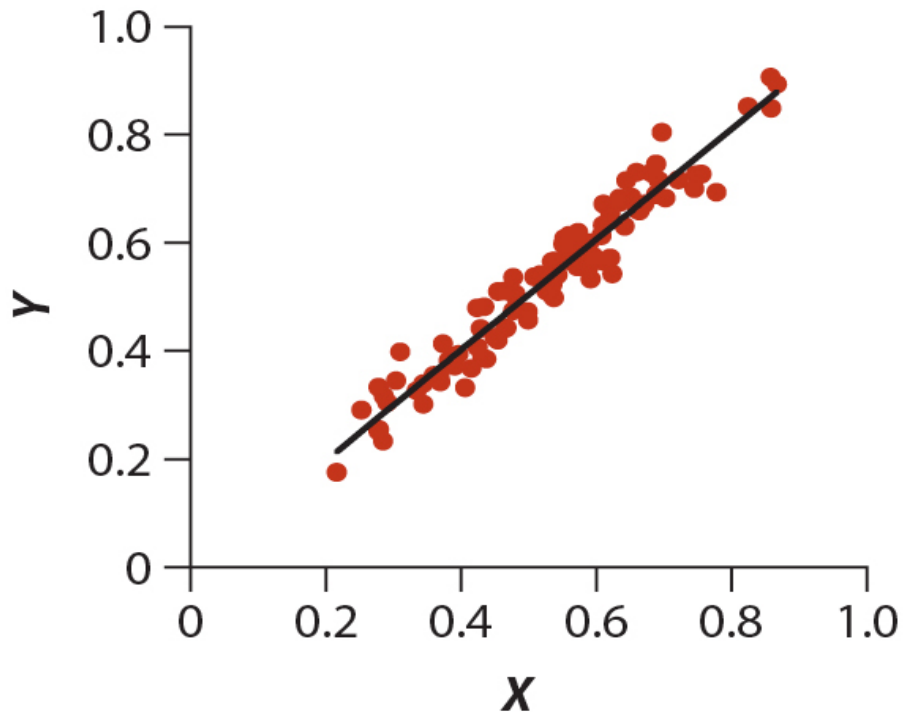
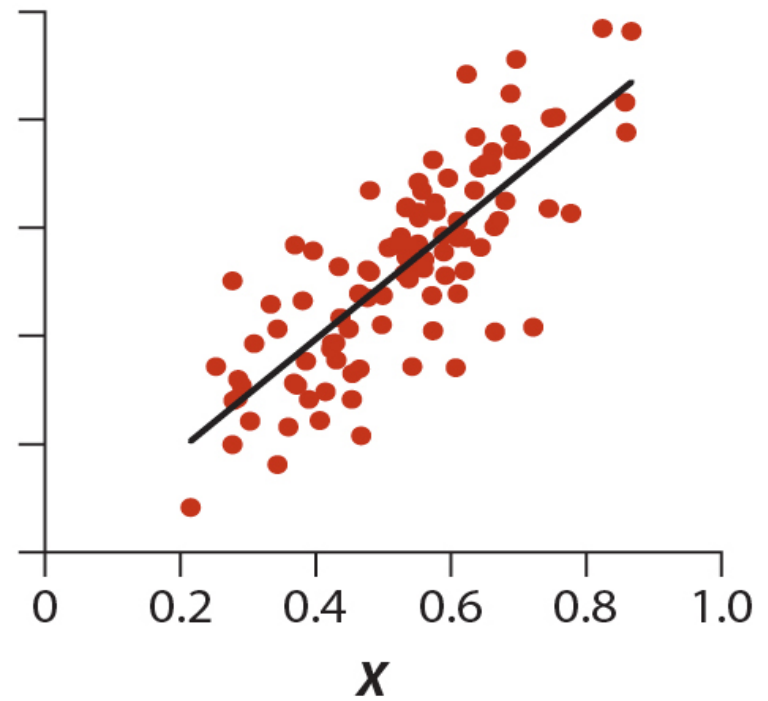# Residual plots are essential

# Inspecting residuals



y = a + bx

sqrt(y) = a + bx

# Measurement errors on Y...



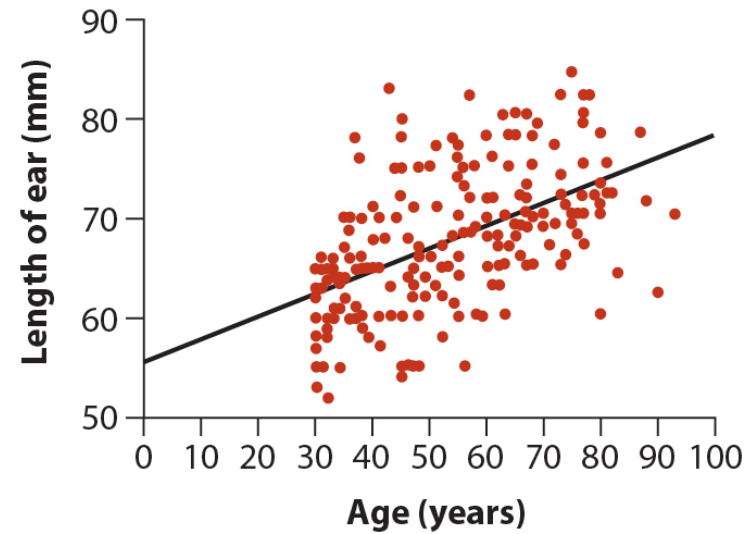No measurement error

Measurement error in Y

# Measurement errors on X



Measurement error in $X$
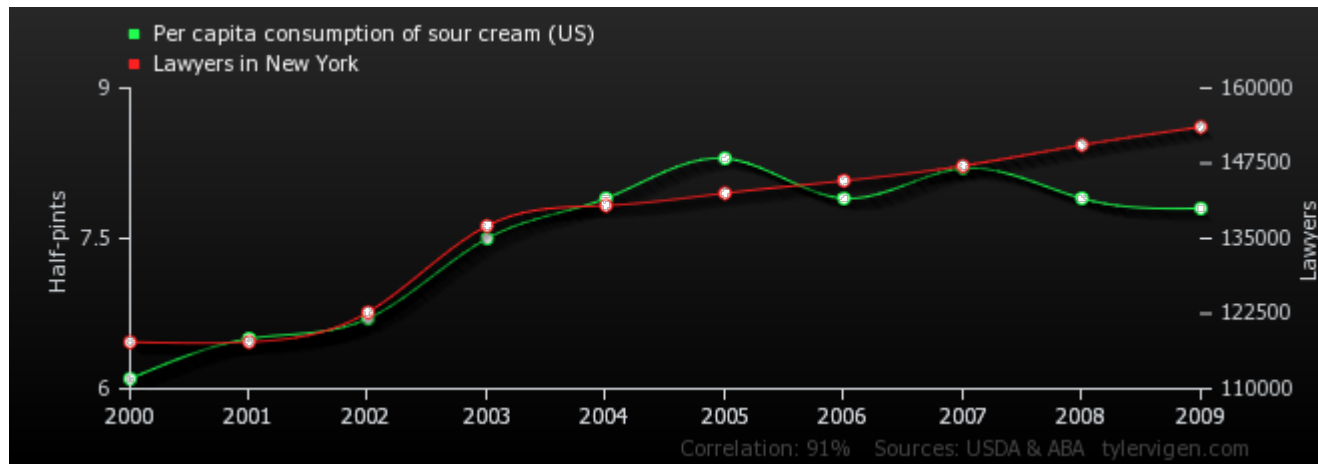
# Extrapolation

# Transformation

# How are they connected?

- What is the difference between correlation and regression?

- What are similarities?

- Sanity checks?

- Transformations?

# What does it mean?

- Correlation != causation

# Exercises

- Distribution of levels of gene expression in Orangutans
- Is evolution (dN/dS) related to gene expression?
- We will summarize all genes pr. Chromosome
- Chromosomes with high gene expression
- Chromosome with high dN/dS
- What is dN/dS?
- What do you expect? Why?

# Pet project – guns and USA

- I have taken data on gun ownership pr. State and combined with death rate (by firearms)
- Also data on countries in the world

- Are these two variables correlated?
- Is it significant?
- How strong is the effect? (slope)
- What is $R^2$ ?