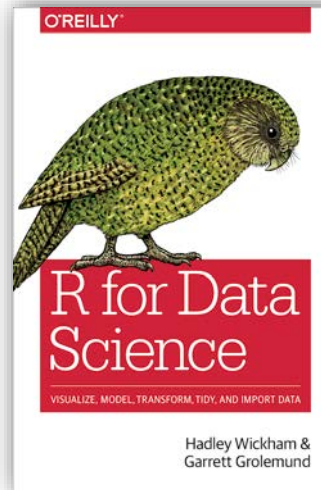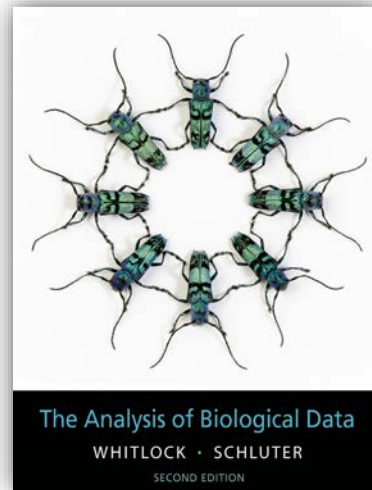# Data Science in Bioinformatics
## week.02.remember.your.statistics.class

Palle Villesen & Thomas Bataillon

# Outline for week 02

- Any questions from last week ?
- Basic statistics
  - Most data are samples
  - Describing data
    - Descriptive statistics of a sample
    - Displaying data
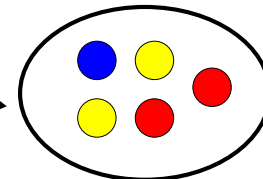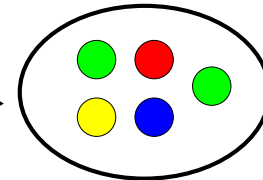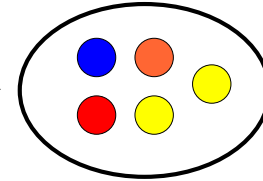  - Sampling with uncertainty
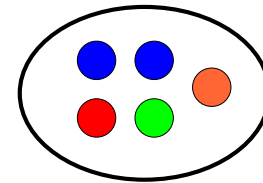    - Sampling distributions

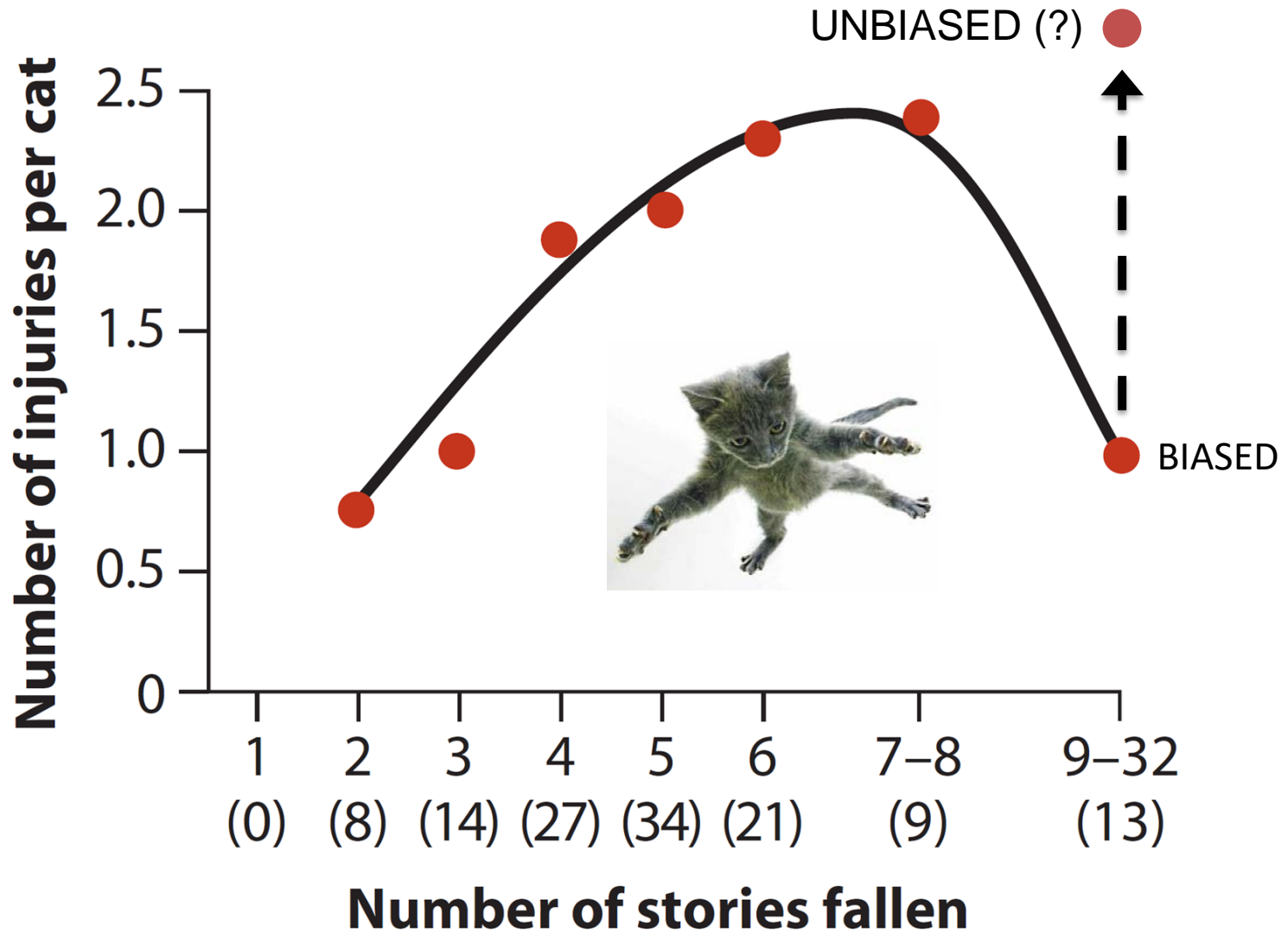# Population    VS    sample

**UNKNOWN**

- Virtually infinite
- Parameters
- Probability distribution

**KNOWN**

- n obs: $x_1, x_2, \ldots x_n$
- Parameter estimates
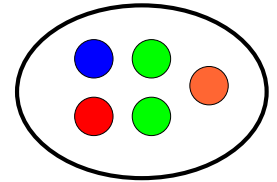- Sampling distribution
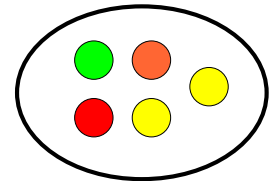
# Biased sample: "High-rise syndrome"

# Biased sample

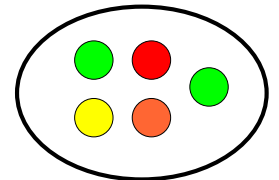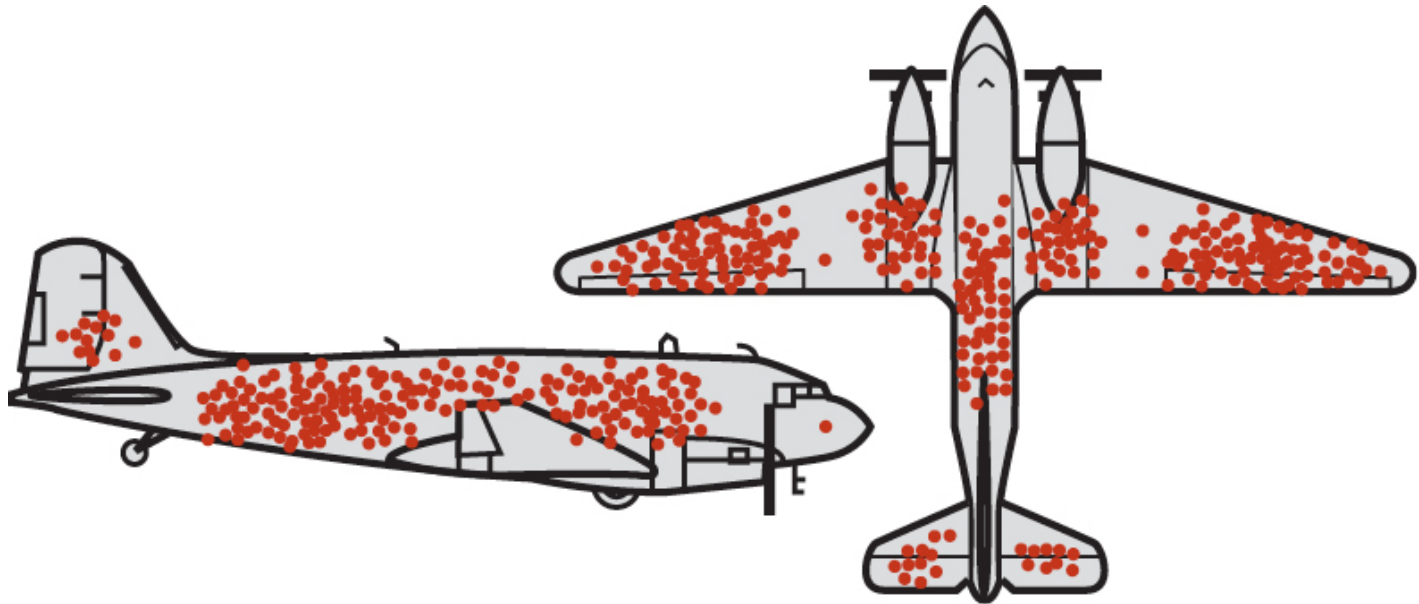- What if the blue balls are more slippery than the others? ...

# Bullet holes in returning planes in WWII

## How would you enforce the planes?

https://xkcd.com/1827/

# Displaying data

# What is best?

# What is best?

# What is best?

# Displaying data

- If is really difficult in ggplot, then it is probably a bad idea

# Describing data

- Arithmetic mean = average
- Standard deviation (compares all observations with mean)
- Median = 50% quantile
- IQR (75% quantile – 25% quantile)

- Present a case where one is really bad

# Life expectancy is a mean

# Standard deviation

$$s = \sqrt{\frac{\sum(Y_i - \bar{Y})^2}{n - 1}}$$

# Median absolute deviation

$$\mathrm{MAD} = \mathrm{median}(\ |X_i - \mathrm{median}(X)|\ ),$$

# Exercise

- Make a small dataset with 3 outliers:
  ```
  x = c(70:90,1000,1100,1200)
  ```
- Calculate mean and median
  ```
  mean(x)
  ```
  ```
  median(x)
  ```
- Calculate sd, iqr and mad
  ```
  sd(x)
  ```
  ```
  quantile(x)
  ```
  ```
  mad(x)
  ```
- Conclusions?

# Distribution of the data: boxplots

# Distribution of the data: histograms

# Outliers

- Outliers:
  - Influence mean and sd

- Median and iqr are more robust
  - Median is a single data point
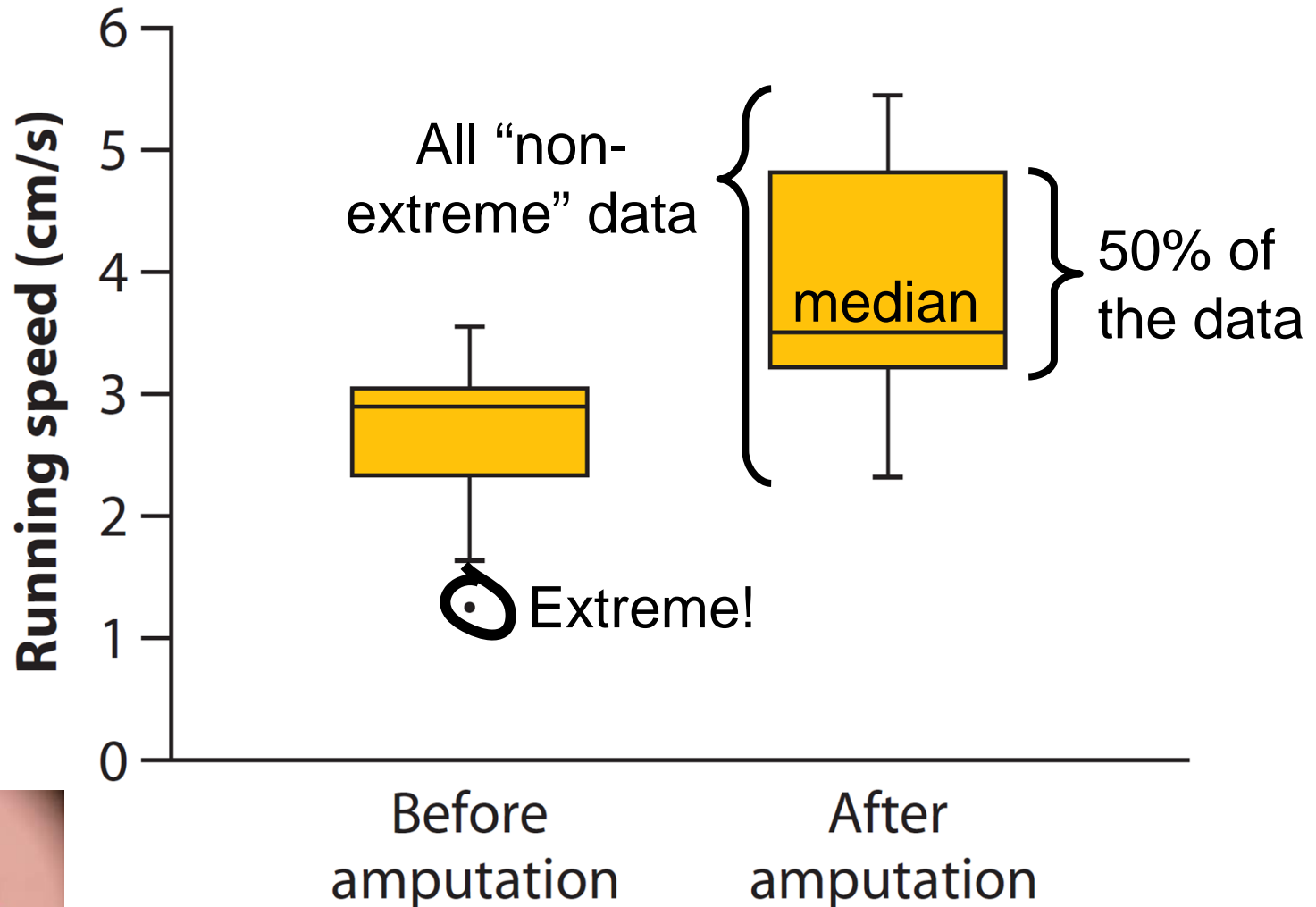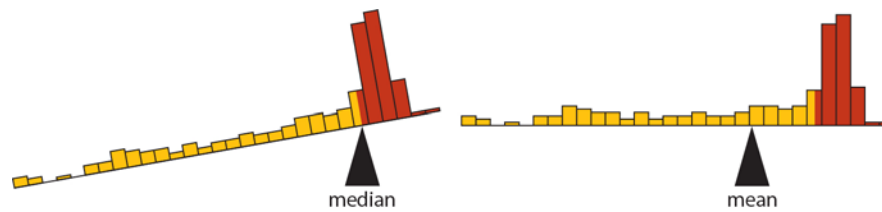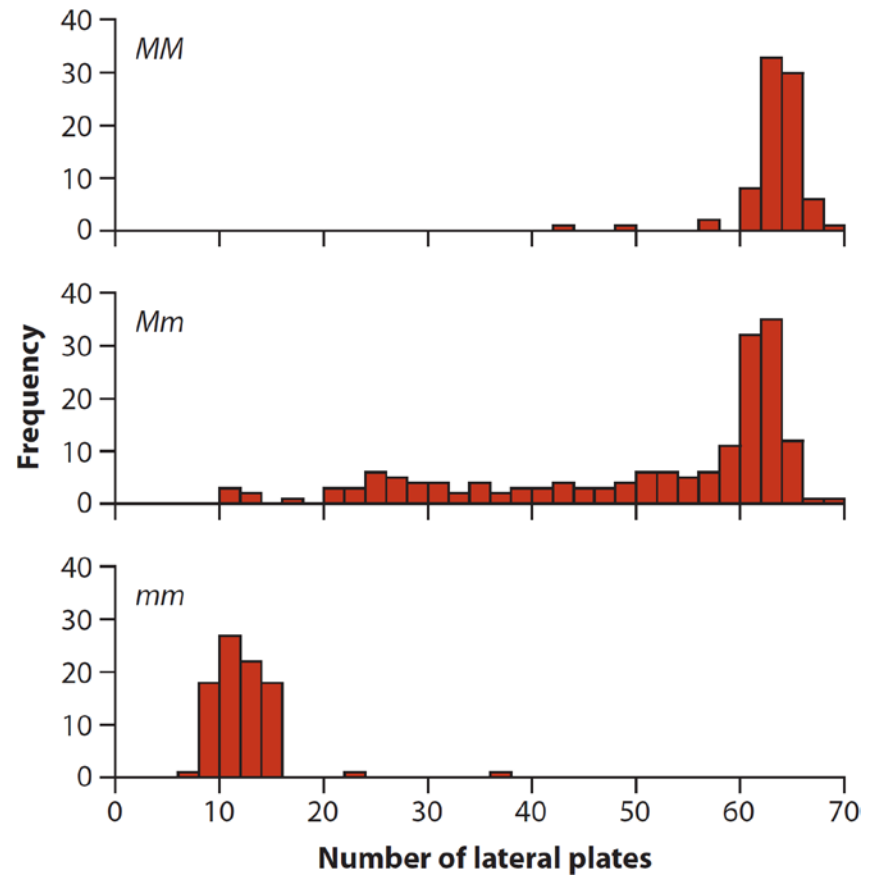  - iqr only use half of the data

# Estimating with uncertainty

# Keywords

- We want to say something about the **population**
- But we only have a **sample**
- So the **sample estimate** is different from the true value because of **sampling error**
- **The sampling distribution** is the distribution of estimate from different samples
- **The standard error** is the standard deviation of the sampling distribution
- **Confidence intervals on the estimate**
  - The 2SE rule of thumb
  - Bootstrap (chapter 19)

# All genes

```
df <- read_csv(file =
"chap04e1HumanGeneLengths.csv")
```

```
set.seed(0)
dfsub <- df %>% sample_n(size = 100)
```

- ```
  df <- read_csv(file="chap04e1HumanGeneLengths.csv")
  ```
- # Replicate figure 4.1-3
  - the sampling distribution of the mean for n=100
- # Replicate figure 4.1-4
  - the sampling distribution of the mean for n=20, n=100 and n=500
- # Calculate standard error from your samples (n=20, n=100, n=500)
- # Compare with table 4.2-1

```r
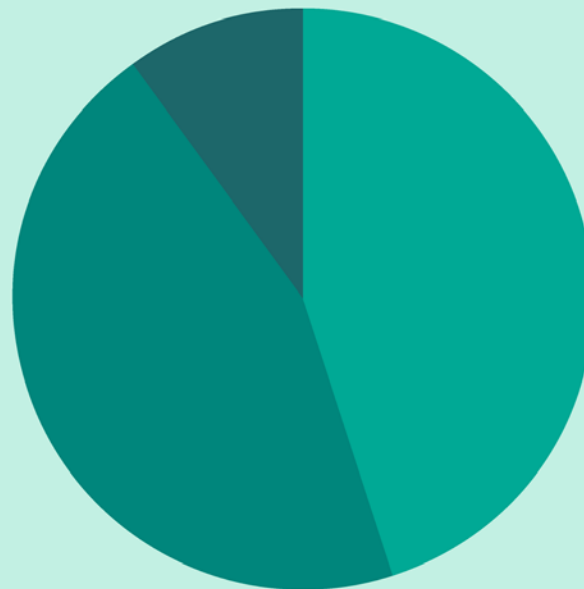set.seed(0)
r <- data.frame() %>% tbl_df()
for (i in 1:10000) {
  dfsub = df %>% sample_n(size = 100)
  r <- rbind(r, data.frame(n = nrow(dfsub), gene.mean = mean(dfsub$geneLength)))
}
```

GROUP WORK

- Splitting up tasks
- Scolding and accusing others of not doing their part
- Actual work

TRUTH FACTS