

# ASTA

*The ASTA team*

## Contents

<b>1</b>	<b>Software</b>	<b>4</b>
1.1	<b>Rstudio</b> . . . . .	4
1.2	<b>R</b> basics . . . . .	4
1.3	<b>R</b> extensions . . . . .	5
1.4	<b>R</b> help . . . . .	5
<b>2</b>	<b>Data</b>	<b>6</b>
2.1	Data example . . . . .	6
2.2	Data example (continued) - variables and format . . . . .	6
2.3	Data types . . . . .	7
<b>3</b>	<b>Population and sample</b>	<b>7</b>
3.1	Aim of statistics . . . . .	7
3.2	Selecting <b>randomly</b> . . . . .	7
<b>4</b>	<b>Variable grouping and frequency tables</b>	<b>8</b>
4.1	Binning . . . . .	8
4.2	Tables . . . . .	8
4.3	2 factors: Cross tabulation . . . . .	9
<b>5</b>	<b>Graphics</b>	<b>9</b>
5.1	Bar graph . . . . .	9
5.2	The Ericksen data . . . . .	11
5.3	Histogram (quantitative variables) . . . . .	12
<b>6</b>	<b>Summary of quantitative variables</b>	<b>13</b>
6.1	Measures of center of data: Mean and median . . . . .	13
6.2	Measures of variability of data: range, standard deviation and variance . . . . .	13
6.3	Calculation of mean, median and standard deviation using <b>R</b> . . . . .	13
6.4	A word about terminology . . . . .	14
6.5	The empirical rule . . . . .	14
6.6	Percentiles . . . . .	15
6.7	Median, quartiles and interquartile range . . . . .	15

<b>7</b>	<b>More graphics</b>	<b>15</b>
7.1	Box-and-whiskers plots (or simply box plots)	15
7.2	2 quantitative variables: Scatter plot	17
<b>8</b>	<b>Appendix</b>	<b>20</b>
8.1	Recoding variables	20
<b>9</b>	<b>Point and click plotting</b>	<b>21</b>
9.1	matplotlib	21
<b>10</b>	<b>Probability of events</b>	<b>21</b>
10.1	The concept of probability	21
10.2	Actual experiment	21
10.3	Another experiment	22
10.4	Definitions	23
10.5	Theoretical probabilities of two events	23
10.6	Conditional probability	24
10.7	Conditional probability and independence	25
10.8	Discrete distribution	25
<b>11</b>	<b>Distribution of general random variables</b>	<b>26</b>
11.1	Probability distribution	26
11.2	Population parameters	27
11.3	Expected value (mean) for a discrete distribution	27
11.4	Variance and standard deviation for a discrete distribution	28
11.5	The binomial distribution	28
11.6	Distribution of a continuous random variable	29
11.7	Density function	30
11.8	Normal distribution	31
<b>12</b>	<b>Distribution of sample statistic</b>	<b>35</b>
12.1	Estimates and their variability	35
12.2	Distribution of sample mean	36
<b>13</b>	<b>Point and interval estimates</b>	<b>38</b>
13.1	Point and interval estimates	38
13.2	Point estimators: Bias	38
13.3	Point estimators: Consistency	39
13.4	Point estimators: Efficiency	39

13.5	Notation . . . . .	39
13.6	Confidence Interval . . . . .	39
13.7	Confidence interval for proportion . . . . .	40
13.8	General confidence intervals for proportion . . . . .	41
13.9	Confidence Interval for mean - normally distributed sample . . . . .	42
13.10	$t$ -distribution and $t$ -score . . . . .	42
13.11	Calculation of $t$ -score in <b>R</b> . . . . .	43
13.12	Example: Confidence interval for mean . . . . .	44
13.13	Example: Plotting several confidence intervals in <b>R</b> . . . . .	45
<b>14</b>	<b>Determining sample size</b>	<b>46</b>
14.1	Sample size for proportion . . . . .	46
14.2	Sample size for mean . . . . .	47
<b>15</b>	<b>Data collection</b>	<b>47</b>
15.1	Data collection . . . . .	47
15.2	Data collection . . . . .	48
<b>16</b>	<b>Population and sample</b>	<b>48</b>
16.1	Population and sample . . . . .	48
<b>17</b>	<b>Example: United States presidential election, 1936</b>	<b>49</b>
17.1	Example: United States presidential election, 1936 . . . . .	49
17.2	Example: United States presidential election, 1936 . . . . .	49
17.3	Example: United States presidential election, 1936 . . . . .	49
<b>18</b>	<b>Example: Bullet holes of honor</b>	<b>50</b>
18.1	Example: Bullet holes of honor . . . . .	50
18.2	Example: Bullet holes of honor . . . . .	50
<b>19</b>	<b>Theory: Biases / sampling</b>	<b>50</b>
19.1	Biases . . . . .	50
19.2	Sampling . . . . .	51
<b>20</b>	<b>Data wrangling</b>	<b>51</b>
20.1	Data wrangling . . . . .	51

# 1 Software

## 1.1 Rstudio

- Make a folder on your computer where you want to keep files to use in **Rstudio**. **Do NOT use Danish characters æ, ø, å** in the folder name (or anywhere in the path to the folder).
- Set the working directory to this folder: **Session -> Set Working Directory -> Choose Directory** (shortcut: Ctrl+Shift+H).
- Make the change permanent by setting the default directory in: **Tools -> Global Options -> Choose Directory**.

## 1.2 R basics

- Ordinary calculations:

```
4.6 * (2 + 3)^4
```

```
## [1] 2875
```

- Make a (scalar) object and print it:

```
a <- 4  
a
```

```
## [1] 4
```

- Make a (vector) object and print it:

```
b <- c(2, 5, 7)  
b
```

```
## [1] 2 5 7
```

- Make a sequence of numbers and print it:

```
s <- 1:4  
s
```

```
## [1] 1 2 3 4
```

- Note: A more flexible command for sequences:

```
s <- seq(1, 4, by = 1)
```

- **R** does elementwise calculations:

```
a * b
```

```
## [1] 8 20 28
```

```
a + b
```

```
## [1] 6 9 11
```

```
b ^ 2
```

```
## [1] 4 25 49
```

- Sum and product of elements:

```
sum(b)
```

```
## [1] 14
```

```
prod(b)
```

```
## [1] 70
```

### 1.3 R extensions

- The functionality of **R** can be extended through libraries or packages (much like plugins in browsers etc.). Some are installed by default in **R** and you just need to load them.
- To install a new package in **Rstudio** use the menu: **Tools -> Install Packages**
- You need to know the name of the package you want to install. You can also do it through a command:

```
install.packages("mosaic")
```

- When it is installed you can load it through the `library` command:

```
library(mosaic)
```

- This loads the `mosaic` package which has a lot of convenient functions for this course (we will get back to that later). It also prints a lot of info about functions that have been changed by the `mosaic` package, but you can safely ignore that.

### 1.4 R help

- You get help via `?<command>`:

```
?sum
```

- Use `tab` to make **Rstudio** guess what you have started typing.
- Search for help:

```
help.search("plot")
```

- You can find a cheat sheet with the **R** functions we use for this course here.
- Save your commands in a file for later usage:
  - Select history tab in top right pane in **Rstudio** .
  - Mark the commands you want to save.
  - Press To **Source** button.

## 2 Data

### 2.1 Data example

Data: Magazine Ads Readability

- Thirty magazines were ranked by educational level of their readers.
- Three magazines were **randomly** selected from each of the following groups:
  - Group 1: highest educational level
  - Group 2: medium educational level
  - Group 3: lowest educational level.
- Six advertisements were **randomly** selected from each of the following nine selected magazines:
  - Group 1: [1] Scientific American, [2] Fortune, [3] The New Yorker
  - Group 2: [4] Sports Illustrated, [5] Newsweek, [6] People
  - Group 3: [7] National Enquirer, [8] Grit, [9] True Confessions
- So, the data contains information about a total of 54 advertisements.

### 2.2 Data example (continued) - variables and format

- For each advertisement (54 cases), the data below were observed.
- **Variable names:**
  - WDS = number of words in advertisement
  - SEN = number of sentences in advertisement
  - 3SYL = number of 3+ syllable words in advertisement
  - MAG = magazine (1 through 9 as above)
  - GROUP = educational level (1 through 3 as above)
- Take a look at the data from within **Rstudio**:

```
magAds <- read.delim("https://asta.math.aau.dk/datasets?file=magazineAds.txt")
head(magAds)
```

```
##   WDS SEN X3SYL MAG GROUP
## 1 205  9   34   1     1
## 2 203 20   21   1     1
## 3 229 18   37   1     1
## 4 208 16   31   1     1
## 5 146  9   10   1     1
## 6 230 16   24   1     1
```

- Variable names are in the top row. They are not allowed to start with a digit, so an X has been prefixed in X3SYL.

## 2.3 Data types

### 2.3.1 Quantitative variables

- The measurements have numerical values.
- Quantitative data often comes about in one of the following ways:
  - **Continuous variables:** measurements of e.g. waiting times in a queue, revenue, share prices, etc.
  - **Discrete variables:** counts of e.g. words in a text, hits on a webpage, number of arrivals to a queue in one hour, etc.
- Measurements like this have a well-defined scale and in **R** they are stored as the type **numeric**.
- It is important to be able to distinguish between discrete count variables and continuous variables, since this often determines how we describe the uncertainty of a measurement.

### 2.3.2 Categorical/qualitative variables

- The measurement is one of a set of given categories, e.g. sex (male/female), social status, satisfaction score (low/medium/high), etc.
- The measurement is usually stored (which is also recommended) as a **factor** in **R**. The possible categories are called **levels**. Example: the levels of the factor “sex” is male/female.
- Factors have two so-called scales:
  - **Nominal scale:** There is no natural ordering of the factor levels, e.g. sex and hair color.
  - **Ordinal scale:** There is a natural ordering of the factor levels, e.g. social status and satisfaction score. A factor in **R** can have a so-called **attribute** assigned, which tells if it is ordinal.

## 3 Population and sample

### 3.1 Aim of statistics

- Statistics is all about “saying something” about a population.
- Typically, this is done by taking a random sample from the population.
- The sample is then analysed and a statement about the population can be made.
- The process of making conclusions about a population from analysing a sample is called **statistical inference**.

### 3.2 Selecting randomly

- For the magazine data:
  - First we select **randomly** 3 magazines from each group.
  - Then we select **randomly** 6 ads from each magazine.
  - An important detail is that the selection is done completely at **random**, i.e.
    - \* each magazine within a group have an equal chance of being chosen and
    - \* each ad within a magazine have an equal chance of being chosen.
- In the following it is a fundamental requirement that the data collection respects this principle of randomness and in this case we use the term **sample**.
- More generally:
  - We have a **population** of objects.
  - We choose completely at random  $n$  of these objects, and from the  $j$ th object we get the measurement  $y_j$ ,  $j = 1, 2, \dots, n$ .

- The measurements  $y_1, y_2, \dots, y_n$  are then called a **sample**.
- If we e.g. are measuring the water quality 4 times in a year then it is a bad idea to only collect data in fair weather. The chosen sampling time is not allowed to be influenced by something that might influence the measurement itself.

## 4 Variable grouping and frequency tables

### 4.1 Binning

- The function `cut` will divide the range of a numeric variable in a number of equally sized intervals, and record which interval each observation belongs to. E.g. for the variable `X3SYL` (the number of words with more than three syllables) in the magazine data:

```
# Before 'cutting':
magAds$X3SYL[1:5]
```

```
## [1] 34 21 37 31 10
```

```
# After 'cutting' into 4 intervals:
syll <- cut(magAds$X3SYL, 4)
syll[1:5]
```

```
## [1] (32.2,43] (10.8,21.5] (32.2,43] (21.5,32.2] (-0.043,10.8]
## Levels: (-0.043,10.8] (10.8,21.5] (21.5,32.2] (32.2,43]
```

- The result is a **factor** and the labels are the interval end points by default. Custom ones can be assigned through the `labels` argument:

```
labs <- c("few", "some", "many", "lots")
syll <- cut(magAds$X3SYL, 4, labels = labs) # NB: this overwrites the 'syll' defined above
syll[1:5]
```

```
## [1] lots some lots many few
## Levels: few some many lots
```

```
magAds$syll <- syll # Adding a new column to the dataset
```

### 4.2 Tables

- To summarize the results we can use the function `tally` from the `mosaic` package (remember the package **must be loaded** via `library(mosaic)` if you did not do so yet):

```
tally( ~ syll, data = magAds)
```

```
## syll
## few some many lots
## 26 14 10 4
```

- In percent:



```
tally( ~ syll, data = magAds, format = "percent")
```

```
## syll
## few some many lots
## 48.1 25.9 18.5 7.4
```

- Here we use an **R** formula (characterized by the “tilde” sign ~) to indicate that we want this variable from the dataset `magAds` (without the tilde it would look for a global variable called `syll` and use that rather than the one in the dataset).

### 4.3 2 factors: Cross tabulation

- To make a table of all combinations of two factors we use `tally` again:

```
tally( ~ syll + GROUP, data = magAds)
```

```
##      GROUP
## syll  1  2  3
## few   8 11  7
## some  4  2  8
## many  3  5  2
## lots  3  0  1
```

- Relative frequencies (in percent) columnwise:

```
tally( ~ syll | GROUP, data = magAds, format = "percent")
```

```
##      GROUP
## syll  1  2  3
## few 44.4 61.1 38.9
## some 22.2 11.1 44.4
## many 16.7 27.8 11.1
## lots 16.7  0.0  5.6
```

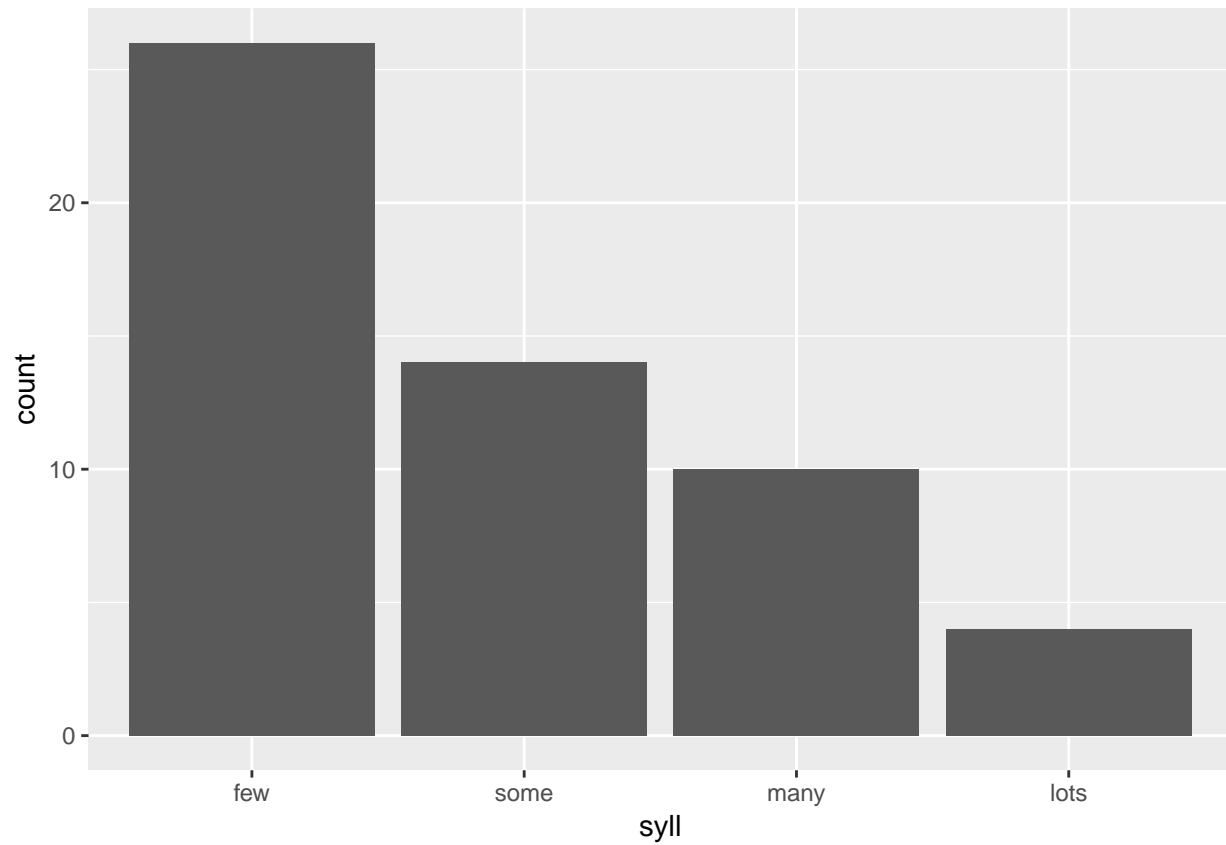
- So, the above table shows e.g. how many percentage of the advertisements in group 1 that have ‘few’, ‘some’, ‘many’ or ‘lots’ words with more than 3 syllables.

## 5 Graphics

### 5.1 Bar graph

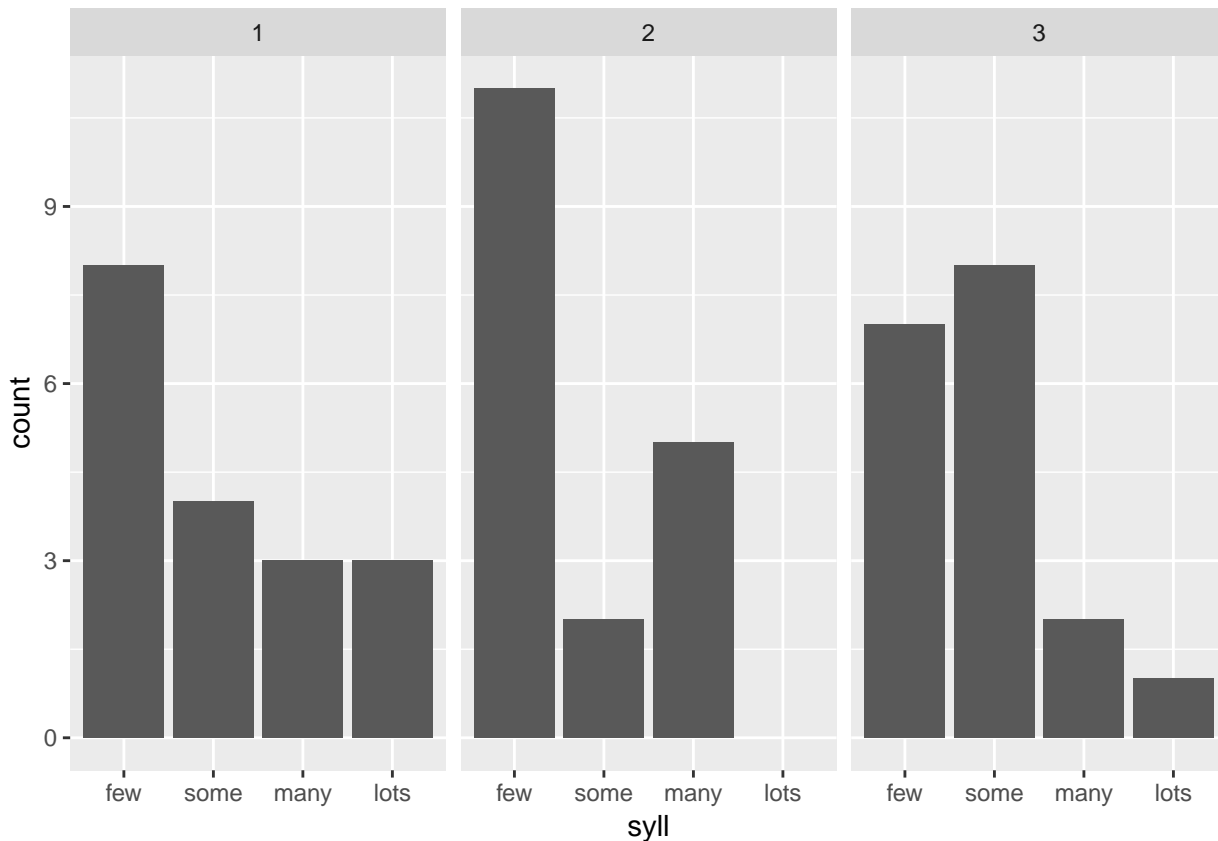
- To create a bar graph plot of table data we use the function `gf_bar` from `mosaic`. For each level of the factor a box is drawn with the height proportional to the frequency (count) of the level.

```
gf_bar( ~ syll, data = magAds)
```



- The bar graph can also be split by group:

```
gf_bar( ~ syll | GROUP, data = magAds)
```



## 5.2 The Ericksen data

- Description of data: Ericksen 1980 U.S. Census Undercount.
- This data contains the following variables:
  - **minority**: Percentage black or Hispanic.
  - **crime**: Rate of serious crimes per 1000 individuals in the population.
  - **poverty**: Percentage poor.
  - **language**: Percentage having difficulty speaking or writing English.
  - **highschool**: Percentage aged 25 or older who had not finished highschool.
  - **housing**: Percentage of housing in small, multiunit buildings.
  - **city**: A factor with levels: **city** (major city) and **state** (state or state-remainder).
  - **conventional**: Percentage of households counted by conventional personal enumeration.
  - **undercount**: Preliminary estimate of percentage undercount.
- The Ericksen data has 66 rows/observations and 9 columns/variables.
- The observations are measured in 16 large cities, the remaining parts of the states in which these cities are located, and the other U.S. states.

```
Ericksen <- read.delim("https://asta.math.aau.dk/datasets?file=Ericksen.txt")
head(Ericksen)
```

```
##      name minority crime poverty language highschool housing city
## 1  Alabama    26.1   49     19      0.2         44     7.6 state
## 2  Alaska     5.7    62     11      1.7         18    23.6 state
## 3  Arizona   18.9   81     13      3.2         28     8.1 state
```

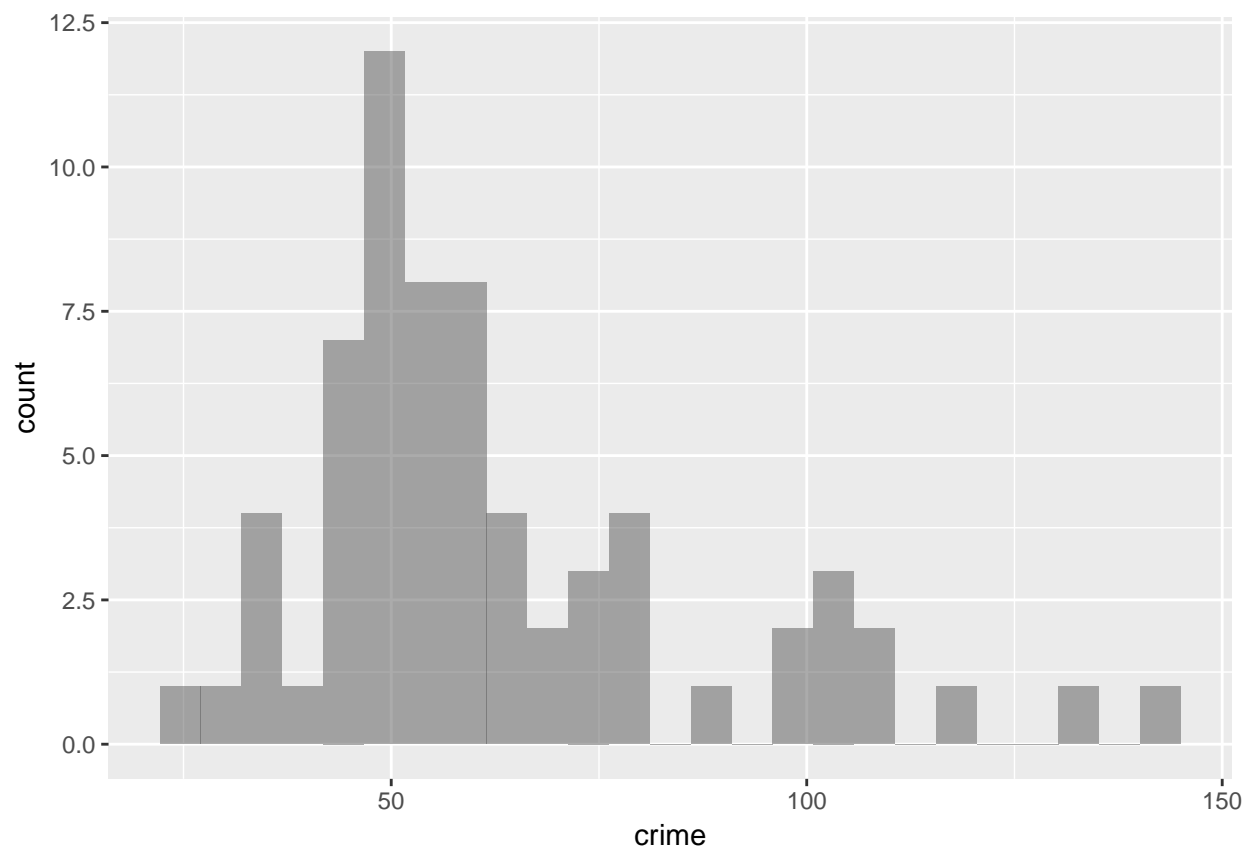
```
## 4   Arkansas      16.9   38    19    0.2    44    7.0 state
## 5 California.R   24.3   73    10    5.0    26   11.8 state
## 6   Colorado     15.2   73    10    1.2    21    9.2 state
##   conventional  undercount
## 1             0     -0.04
## 2            100     3.35
## 3             18     2.48
## 4             0     -0.74
## 5             4     3.60
## 6            19     1.34
```

- Want to make a histogram for crime rate - how?

### 5.3 Histogram (quantitative variables)

- How to make a histogram for some variable  $x$ :
  - Divide the interval from the minimum value of  $x$  to the maximum value of  $x$  in an appropriate number of equal sized sub-intervals.
  - Draw a box over each sub-interval with the height being proportional to the number of observations in the sub-interval.
- Histogram of crime rates for the Ericksen data

```
gf_histogram( ~ crime, data = Ericksen)
```



## 6 Summary of quantitative variables

### 6.1 Measures of center of data: Mean and median

- We return to the magazine ads example (WDS = number of words in advertisement). A number of numerical summaries for WDS can be retrieved using the `favstats` function:

```
favstats( ~ WDS, data = magAds)
```

```
## min Q1 median Q3 max mean sd n missing
## 31 69 96 202 230 123 66 54 0
```

- The observed values of the variable WDS are  $y_1 = 205, y_2 = 203, \dots, y_n = 208$ , where there are a total of  $n = 54$  values. As previously defined this constitutes a **sample**.
- **mean** = 123 is the **average** of the sample, which is calculated by

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

We may also call  $\bar{y}$  the **(empirical) mean** or the **sample mean**.

- **median** = 96 is the 50-percentile, i.e. the value that splits the sample in 2 groups of equal size.
- An important property of the **mean** and the **median** is that they have the same unit as the observations (e.g. meter).

### 6.2 Measures of variability of data: range, standard deviation and variance

- The **range** is the difference of the largest and smallest observation.
- The **(empirical) variance** is the average of the squared deviations from the mean:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2.$$

- **sd** = **standard deviation** =  $s = \sqrt{s^2}$ .
- Note: If the observations are measured in meter, the **variance** has unit meter<sup>2</sup> which is hard to interpret. The **standard deviation** on the other hand has the same unit as the observations.
- The standard deviation describes how much data varies around the (empirical) mean.

### 6.3 Calculation of mean, median and standard deviation using R

The mean, median and standard deviation are just some of the summaries that can be read of the `favstats` output (shown on previous page). They may also be calculated separately in the following way:

- Mean of WDS:

```
mean( ~ WDS, data = magAds)
```

```
## [1] 123
```

- Median of WDS:

```
median( ~ WDS, data = magAds)
```

```
## [1] 96
```

- Standard deviation for WDS:

```
sd( ~ WDS, data = magAds)
```

```
## [1] 66
```

We may also calculate the summaries for each group (variable `GROUP`), e.g. for the mean:

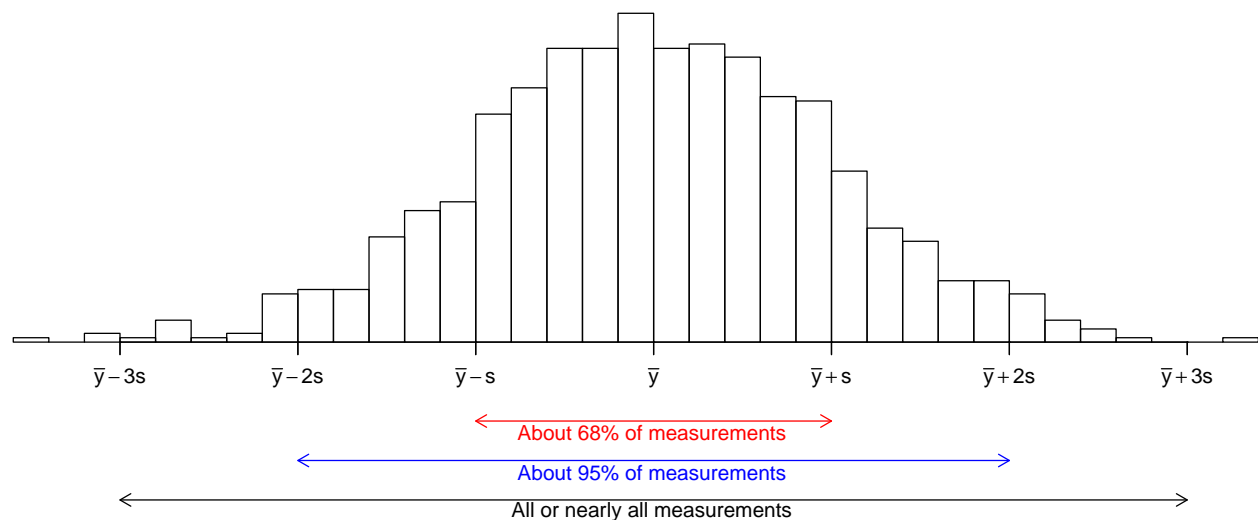
```
mean( ~ WDS | GROUP, data = magAds)
```

```
## 1 2 3  
## 140 121 106
```

## 6.4 A word about terminology

- **Standard deviation:** a measure of variability of a population or a sample.
- **Standard error:** a measure of variability of an estimate. For example, a measure of variability of the sample mean.

## 6.5 The empirical rule



If the histogram of the sample looks like a bell shaped curve, then

- about 68% of the observations lie between  $\bar{y} - s$  and  $\bar{y} + s$ .
- about 95% of the observations lie between  $\bar{y} - 2s$  and  $\bar{y} + 2s$ .
- All or almost all (99.7%) of the observations lie between  $\bar{y} - 3s$  and  $\bar{y} + 3s$ .

## 6.6 Percentiles

- The  $p$ th percentile is a value such that about  $p\%$  of the population (or sample) lies below or at this value and about  $(100 - p)\%$  of the population (or sample) lies above it.

### 6.6.1 Percentile calculation for a sample:

- First, sort data in increasing order. For the WDS variable in the magazine data:

$$y_{(1)} = 31, y_{(2)} = 32, y_{(3)} = 34, \dots, y_{(n)} = 230.$$

Here the number of observations is  $n = 54$ .

- Find the 5th percentile (i. e.  $p = 5$ ):
  - The observation number corresponding to the 5-percentile is  $N = \frac{n \cdot p}{100} = 2.7$ .
  - That is, the 5-percentile must lie between the observations  $y_{(k)} = 32$  and  $y_{(k+1)} = 34$ , where  $k = 2 < N < 3$ .
  - Let  $d = N - k = 0.7$ . One of several methods for estimating the 5-percentile:

$$y_{(k)} + d(y_{(k+1)} - y_{(k)}) = 32 + 0.7 \cdot 2 = 33.4$$

## 6.7 Median, quartiles and interquartile range

Recall

```
favstats( ~ WDS, data = magAds)
```

```
## min Q1 median Q3 max mean sd n missing
## 31 69 96 202 230 123 66 54 0
```

- 50-percentile = 96 is the **median** and it is a measure of the center of data.
- 0-percentile = 31 is the **minimum** value.
- 25-percentile = 69 is called the **lower quartile** (Q1). Median of lower 50% of data.
- 75-percentile = 202 is called the **upper quartile** (Q3). Median of upper 50% of data.
- 100-percentile = 230 is the **maximum** value.
- **Interquartile Range (IQR)**: a measure of variability given by the difference of the upper and lower quartiles:  $202 - 69 = 133$ .

## 7 More graphics

### 7.1 Box-and-whiskers plots (or simply box plots)

How to draw a box-and-whiskers plot:

- Box:
  - Calculate the median, lower and upper quartiles.
  - Plot a line by the median and draw a box between the upper and lower quartiles.
- Whiskers:
  - Calculate interquartile range and call it IQR.

- Calculate the following values:
  - \*  $L = \text{lower quartile} - 1.5 \cdot \text{IQR}$
  - \*  $U = \text{upper quartile} + 1.5 \cdot \text{IQR}$
- Draw a line from lower quartile to the smallest measurement, which is larger than  $L$ .
- Similarly, draw a line from upper quartile to the largest measurement which is smaller than  $U$ .
- Outliers: Measurements smaller than  $L$  or larger than  $U$  are drawn as circles.

*Note: Whiskers are minimum and maximum of the observations that are not deemed to be outliers.*

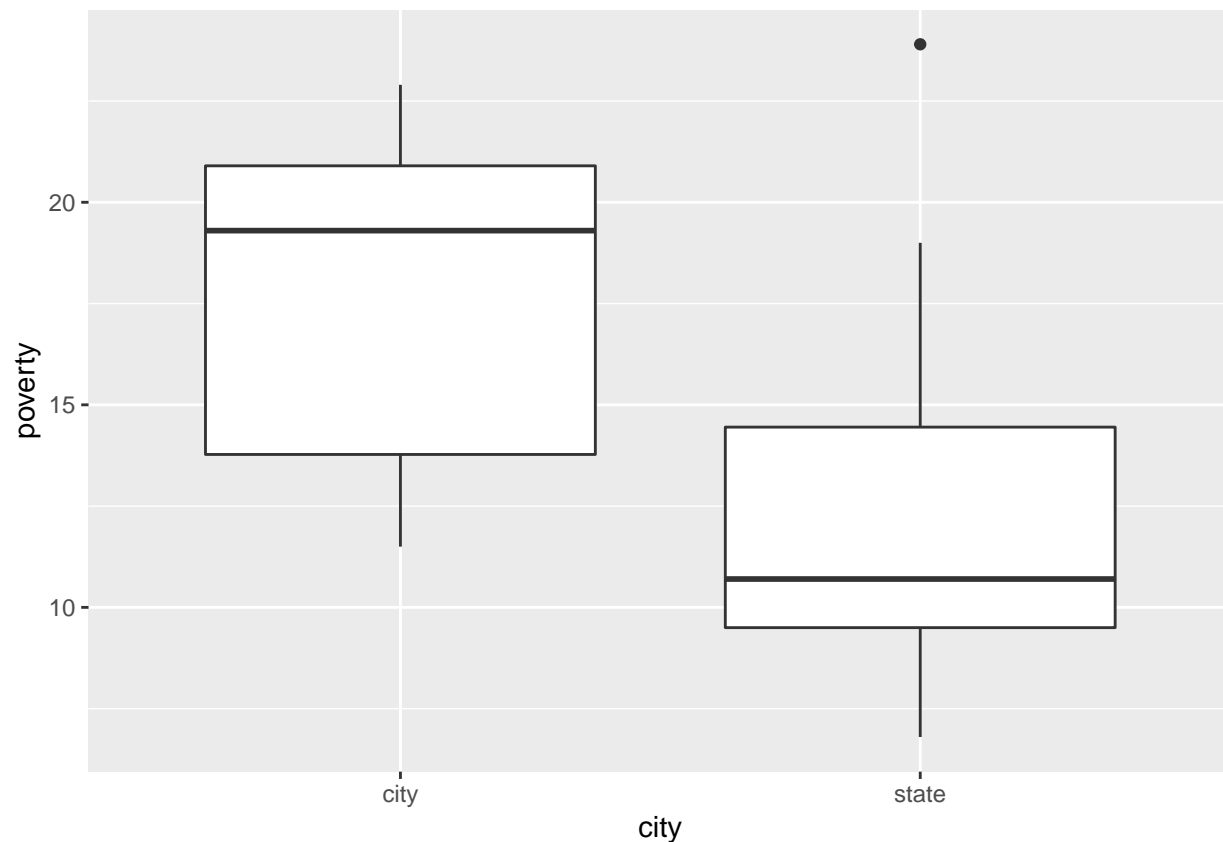
### 7.1.1 Boxplot for Ericksen data

Boxplot of the poverty rates separately for cities and states (variable `city`):

```
favstats(poverty ~ city, data = Ericksen)
```

```
##   city min  Q1 median Q3 max mean  sd  n missing
## 1 city 11.5 13.8   19 21 23  18 4.0 16     0
## 2 state 6.8  9.5   11 14 24  12 3.7 50     0
```

```
gf_boxplot(poverty ~ city, data = Ericksen)
```



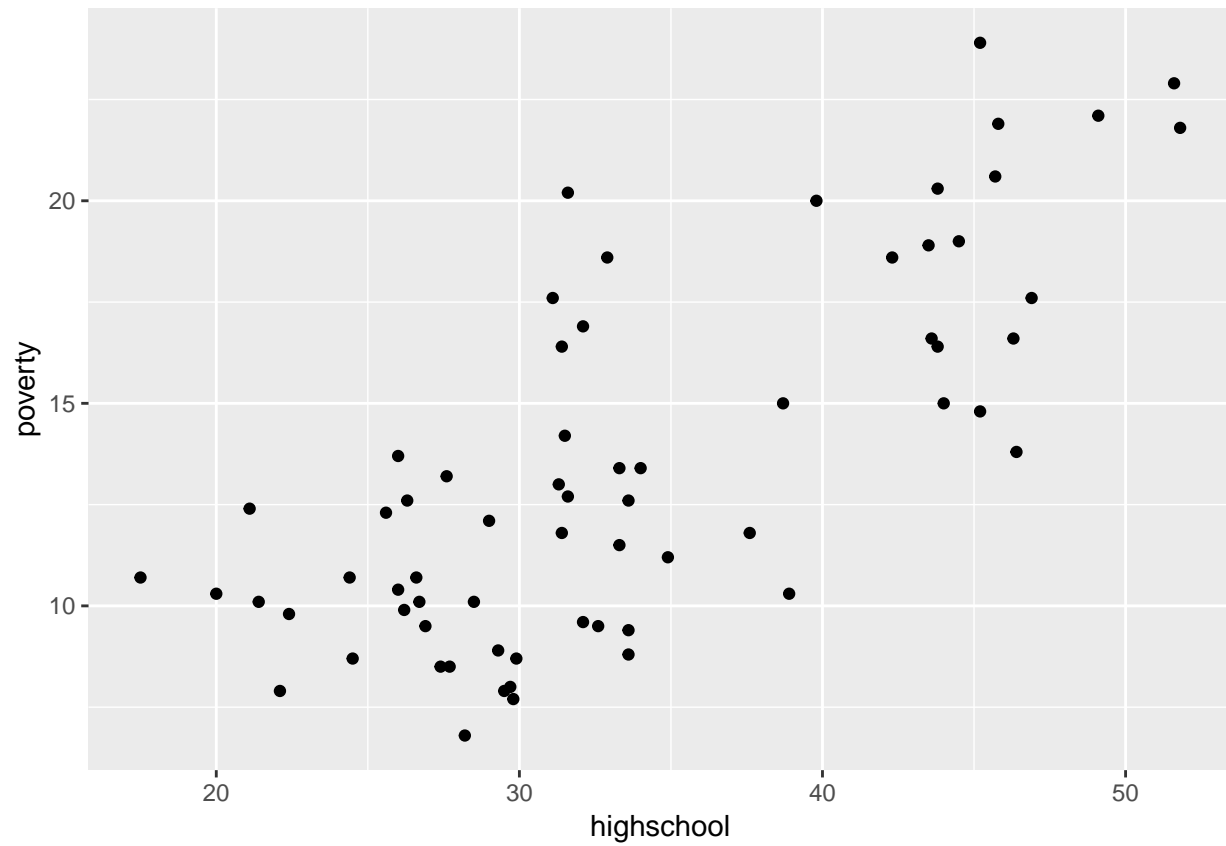
- There seems to be more poverty in the cities.
- A single state differs noticeably from the others with a high poverty rate.



## 7.2 2 quantitative variables: Scatter plot

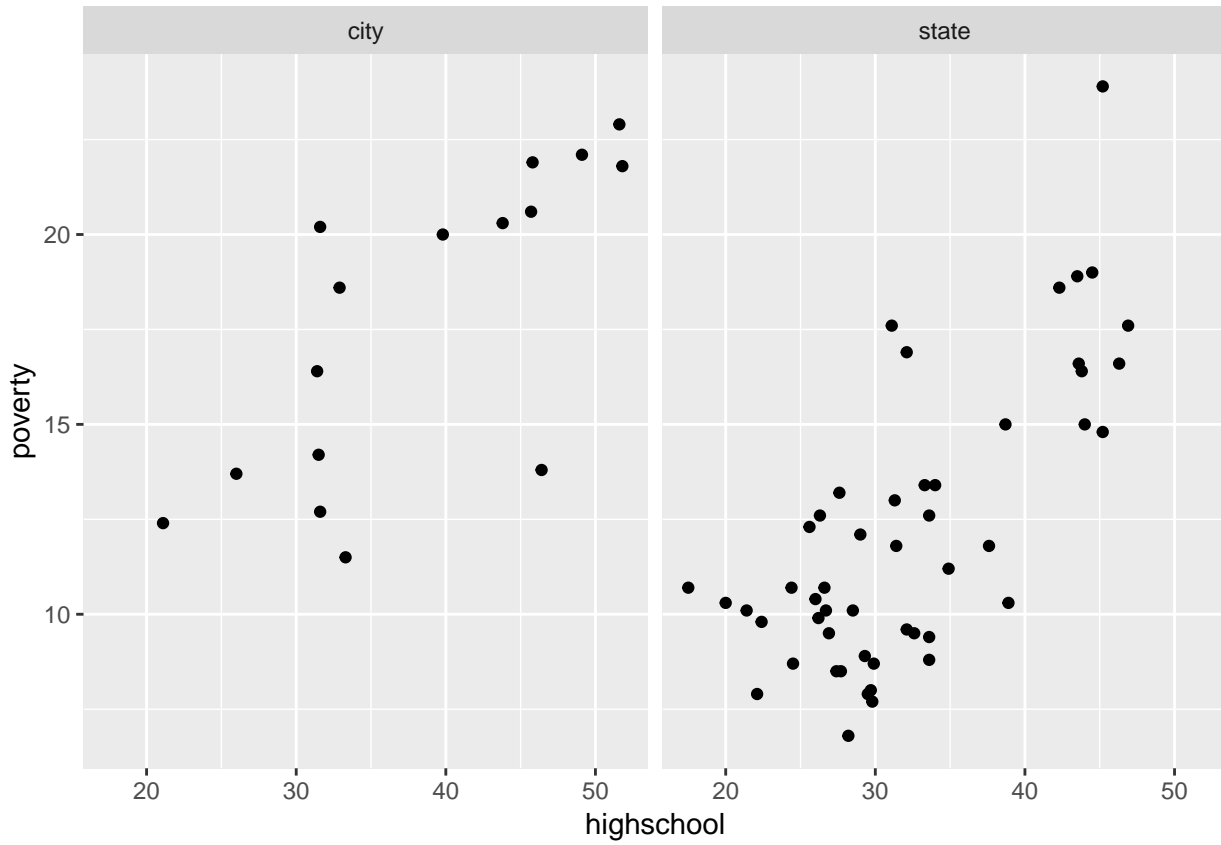
For two quantitative variables the usual graphic is a scatter plot:

```
gf_point(poverty ~ highschool, data = Ericksen)
```

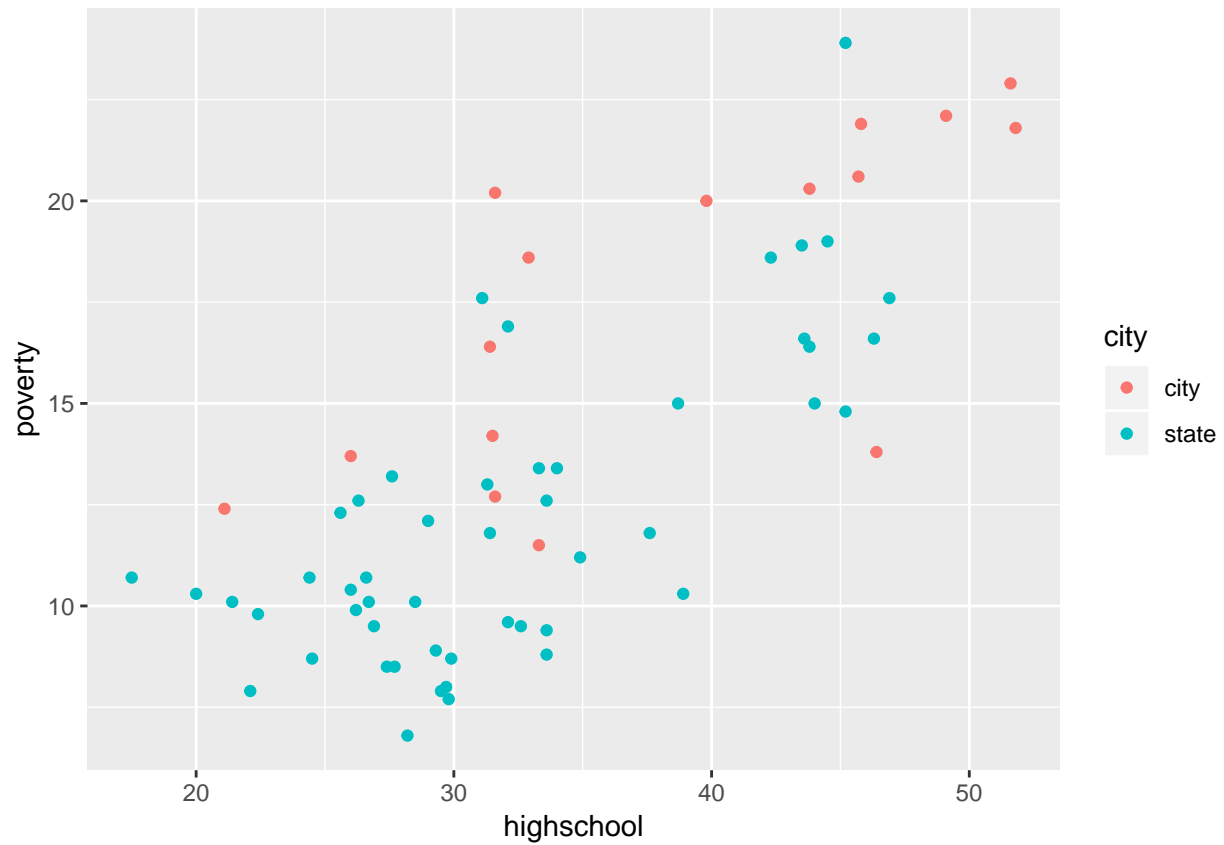


This can be either split or coloured according to the value of city:

```
gf_point(poverty ~ highschool | city, data = Ericksen)
```

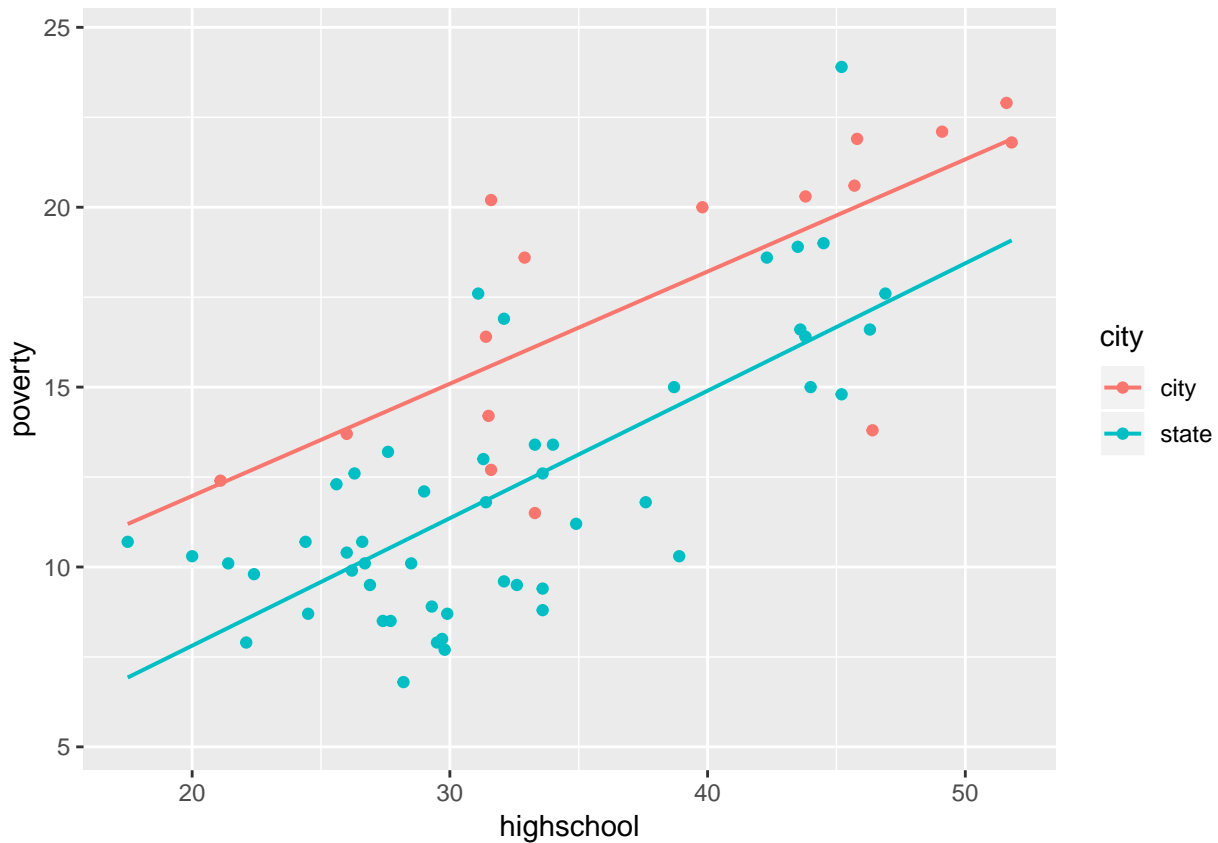


```
gf_point(poverty ~ highschool, col = ~city, data = Ericksen)
```



If we want a regression line along with the points we can do:

```
gf_point(poverty ~ highschool, col = ~city, data = Ericksen) %>% gf_lm()
```



## 8 Appendix

### 8.1 Recoding variables

- The function `factor` will directly convert a vector to be of type `factor`. E.g.:

```
head(magAds$GROUP)
```

```
## [1] 1 1 1 1 1 1
```

```
f <- factor(magAds$GROUP)
magAds$GROUP <- f
head(magAds$GROUP)
```

```
## [1] 1 1 1 1 1 1
## Levels: 1 2 3
```

- Custom labels for the levels can also be used:

```
f <- factor(magAds$GROUP,
           levels = c("1", "2", "3"),
           labels = c("high", "medium", "low"))
magAds$GROUP <- f
head(magAds$GROUP)
```

```
## [1] high high high high high high
## Levels: high medium low
```

- In this way the numbers are replaced by more informative labels describing the educational level.

## 9 Point and click plotting

### 9.1 mplot

- If `mosaic` is loaded and the package `manipulate` is installed you can construct plots using point and click using the function `mplot`.
- You simply use `mplot` on your dataset and answer the question and then you can change things by pressing the settings button (cog wheel) in the top left of the plot window.

```
mplot(Ericksen)
```

- In the end you can press “Show expression” to get the code for the plot.

## 10 Probability of events

### 10.1 The concept of probability

- Experiment: Measure the waiting times in a queue where we note 1, if it exceeds 2 minutes and 0 otherwise.
- The experiment is carried out  $n$  times with results  $y_1, y_2, \dots, y_n$ . There is **random variation** in the outcome, i.e. sometimes we get a 1 other times a 0.
- **Empirical probability** of exceeding 2 minutes:

$$p_n = \frac{\sum_{i=1}^n y_i}{n}.$$

- **Theoretical probability** of exceeding 2 minutes:

$$\pi = \lim_{n \rightarrow \infty} p_n.$$

- We try to make inference about  $\pi$  based on a sample, e.g. “Is  $\pi > 0.1$ ?” (“do more than 10% of the customers experience a waiting time in excess of 2 minutes?”).
- Statistical inference is concerned with such questions when we only have a finite sample.

### 10.2 Actual experiment

- On February 23, 2017, a group of students were asked how long time (in minutes) they waited in line last time they went to the canteen at AAU’s Copenhagen campus:

```
y_canteen <- c(2, 5, 1, 6, 1, 1, 1, 1, 3, 4, 1, 2, 1, 2, 2, 2, 4, 2, 2, 5, 20, 2, 1, 1, 1, 1)
x_canteen <- ifelse(y_canteen > 2, 1, 0)
x_canteen
```

```
## [1] 0 1 0 1 0 0 0 0 1 1 0 0 0 0 0 0 1 0 0 1 1 0 0 0 0 0
```

- Empirical probability of waiting more than 2 minutes:

```
p_canteen <- sum(x_canteen) / length(x_canteen)
p_canteen
```

```
## [1] 0.27
```

- Question: Is the population probability  $\pi > 1/3$ ?
- Notice: One student said he had waited for 20 minutes (we doubt that; he was trying to make himself interesting. Could consider ignoring that observation).

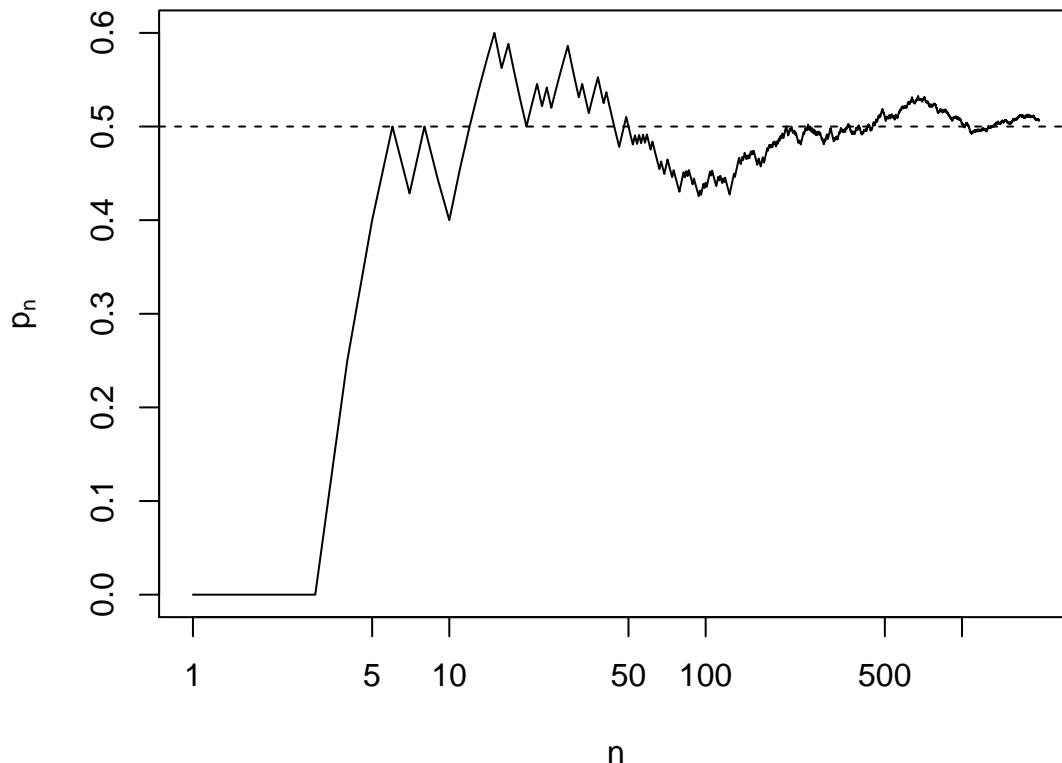
### 10.3 Another experiment

- John Kerrich, a South African mathematician, was visiting Copenhagen when World War II broke out. Two days before he was scheduled to fly to England, the Germans invaded Denmark. Kerrich spent the rest of the war interned at a camp in Hald Ege near Viborg, Jutland, and to pass the time he carried out a series of experiments in probability theory. In one, he tossed a coin 10,000 times. His results are shown in the following graph.
- Below,  $x$  is a vector with the first 2,000 outcomes of John Kerrich's experiment (0 = tail, 1 = head):

```
head(x, 10)
```

```
## [1] 0 0 0 1 1 1 0 1 0 0
```

- Plot of the empirical probability  $p_n$  of getting a head against the number of tosses  $n$ :



(The horizontal axis is on a log scale).

## 10.4 Definitions

- **Sample space:** All possible outcomes of the experiment.
- **Event:** A subset of the sample space.

We conduct the experiment  $n$  times. Let  $\#(A)$  denote how many times we observe the event  $A$ .

- **Empirical probability** of the event  $A$ :

$$p_n(A) = \frac{\#(A)}{n}.$$

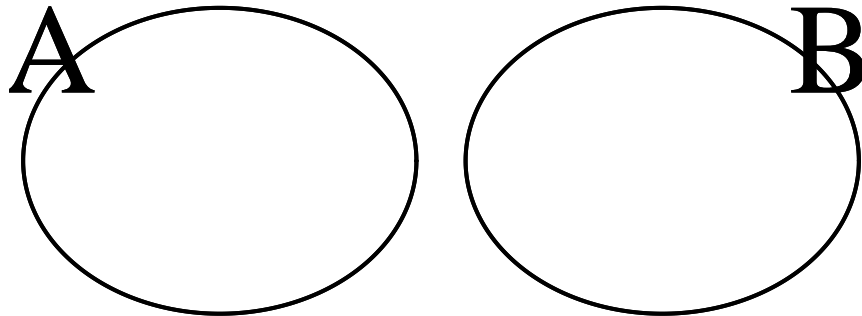
- **Theoretical probability** of the event  $A$ :

$$P(A) = \lim_{n \rightarrow \infty} p_n(A)$$

- We always have  $0 \leq P(A) \leq 1$ .

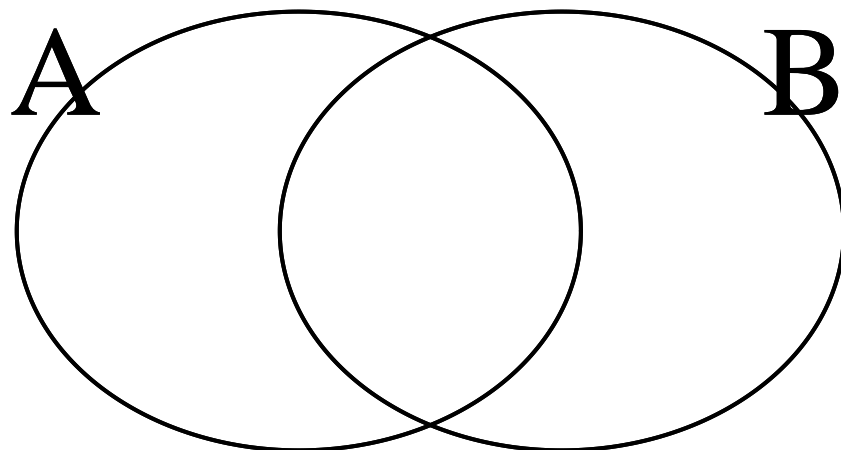
## 10.5 Theoretical probabilities of two events

- If the two events  $A$  and  $B$  are **disjoint** (non-overlapping) then
  - $\#(A \text{ and } B) = 0$  implying that  $P(A \text{ and } B) = 0$ .
  - $\#(A \text{ or } B) = \#(A) + \#(B)$  implying that  $P(A \text{ or } B) = P(A) + P(B)$ .



- If the two events  $A$  and  $B$  are **not disjoint** then the more general formula is

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B).$$



## 10.6 Conditional probability

- Say we consider two events  $A$  and  $B$ . Then the **conditional probability** of  $A$  given (or conditional on) the event  $B$  is written  $P(A | B)$  and is defined by

$$P(A | B) = \frac{P(A \text{ and } B)}{P(B)}.$$

- The above probability can be understood as: “how probable  $A$  is if we know that  $B$  has happened”.

---

### 10.6.1 Example with magazine data:

```
magAds <- read.delim("https://asta.math.aau.dk/datasets?file=magazineAds.txt")

# Create two new factors 'words' and 'education':
magAds$words <- cut(magAds$WDS, breaks = c(31, 72, 146, 230), include.lowest = TRUE)
magAds$education <- factor(magAds$GROUP, levels = c(1, 2, 3), labels = c("high", "medium", "low"))

library(mosaic)
tab <- tally(~ words + education, data = magAds)
tab
```

```
##           education
## words      high medium low
##  [31,72]      4      6   5
##   (72,146]      5      6   8
##  (146,230]      9      6   5
```

- The event  $A = \{\text{words} = (146, 230]\}$  (the ad is a “difficult” text) has empirical probability

$$p_n(A) = \frac{9 + 6 + 5}{54} = \frac{20}{54} \approx 37\%.$$

- Say we only are interested in the probability of a “difficult” text (event  $A$ ) for high education magazines, i.e. conditioning on the event  $B = \{\text{education} = \text{high}\}$ . Then the empirical conditional probability can be calculated from the table:

$$p_n(A | B) = \frac{9}{4 + 5 + 9} = \frac{9}{18} = 0.5 = 50\%.$$

- The conditional probability of  $A$  given  $B$  may theoretically be expressed as

$$\begin{aligned} P(A | B) &= P(\text{words} = (146, 230] | \text{education} = \text{high}) \\ &= \frac{P(\text{words} = (146, 230] \text{ and } \text{education} = \text{high})}{P(\text{education} = \text{high})}, \end{aligned}$$

which translated to empirical probabilities (substituting  $P$  with  $p_n$ ) will give



$$\begin{aligned}
p_n(A | B) &= \frac{p_n(\text{words} = (146, 230] \text{ and education} = \text{high})}{p_n(\text{education} = \text{high})} \\
&= \frac{\frac{9}{54}}{\frac{4+5+9}{54}} \\
&= \frac{9}{4 + 5 + 9} \\
&= 50\%
\end{aligned}$$

as calculated above.

## 10.7 Conditional probability and independence

- If information about  $B$  does not change the probability of  $A$  we talk about independence, i.e.  $A$  is **independent** of  $B$  if

$$P(A | B) = P(A) \Leftrightarrow P(A \text{ and } B) = P(A)P(B)$$

The last relation is symmetric in  $A$  and  $B$ , and we simply say that  $A$  and  $B$  are **independent events**.

- In general the events  $A_1, A_2, \dots, A_k$  are independent if

$$P(A_1 \text{ and } A_2 \text{ and } \dots \text{ and } A_k) = P(A_1)P(A_2) \cdots P(A_k).$$

### 10.7.1 Magazine data revisited

- Recall the empirical probabilities calculated above:

$$p_n(A) = 37\% \quad \text{and} \quad p_n(A | B) = 50\%.$$

- These indicate (we cannot say for sure as we only consider a finite sample - we will later see how to test for this) that the theoretical probability

$$P(A) \neq P(A | B)$$

and hence that knowledge about  $B$  (high education level) may convey information about the probability of  $A$  (the ad containing a “difficult” text).

## 10.8 Discrete distribution

### 10.8.1 Example: Magazine data

```
# Table with the percentage of ads in each combination of the levels of 'words' and 'education'
tab <- tally( ~ words + education, data = magAds, format = "percent")
round(tab, 2) # Round digits
```

```
##           education
## words      high medium low
## [31,72]     7.4  11.1  9.3
## (72,146]    9.3  11.1 14.8
## (146,230] 16.7  11.1  9.3
```

- The 9 disjoint events above (corresponding to combinations of **words** and **education**) make up the whole sample space for the two variables. The empirical probabilities of each event is given in the table.

## 10.8.2 General discrete distribution

- In general:
  - Let  $A_1, A_2, \dots, A_k$  be a subdivision of the sample space into pairwise disjoint events.
  - The probabilities  $P(A_1), P(A_2), \dots, P(A_k)$  are called a **discrete distribution** and satisfy

$$\sum_{i=1}^k P(A_i) = 1.$$

---

## 10.8.3 Example: 3 coin tosses

- **Random/stochastic variable:** A function  $Y$  that translates an outcome of the experiment into a number.
- Possible outcomes in an experiment with 3 coin tosses:
  - 0 heads (TTT)
  - 1 head (HTT, THT, TTH)
  - 2 heads (HHT, HTH, THH)
  - 3 heads (HHH)
- The above events are disjoint and make up the whole sample space.
- Let  $Y$  be the number of heads in the experiment:  $Y(TTT) = 0, Y(HTT) = 1, \dots$
- Assume that each outcome is equally likely, i.e. probability  $1/8$  for each event. Then,
  - $P(\text{no heads}) = P(Y = 0) = P(TTT) = 1/8$ .
  - $P(\text{one head}) = P(Y = 1) = P(HTT \text{ or } THT \text{ or } TTH) = P(HTT) + P(THT) + P(TTH) = 3/8$ .
  - Similarly for 2 or 3 heads.
- So, the distribution of  $Y$  is

Number of heads, $Y$	0	1	2	3
Probability	1/8	3/8	3/8	1/8

# 11 Distribution of general random variables

## 11.1 Probability distribution

- We are conducting an experiment where we make a quantitative measurement  $Y$  (a random variable), e.g. the number of words in an ad or the waiting time in a queue.
- In advance there are many possible outcomes of the experiment, i.e.  $Y$ 's value has an uncertainty, which we quantify by the **probability distribution** of  $Y$ .
- For any interval  $(a, b)$ , the distribution states the probability of observing a value of the random variable  $Y$  in this interval:

$$P(a < Y < b), \quad -\infty < a < b < \infty.$$

- $Y$  is **discrete** if we can enumerate all the possible values of  $Y$ , e.g. the number of words in an ad.
- $Y$  is **continuous** if  $Y$  can take any value in a interval, e.g. a measurement of waiting time in a queue.

### 11.1.1 Sample

We conduct an experiment  $n$  times, where the outcome of the  $i$ th experiment corresponds to a measurement of a random variable  $Y_i$ , where we assume

- The experiments are **independent**
- The variables  $Y_1, \dots, Y_n$  have the **same distribution**

## 11.2 Population parameters

- When the sample size grows, then e.g. the mean of the sample,  $\bar{y}$ , will stabilize around a fixed value,  $\mu$ , which is usually unknown. The value  $\mu$  is called the **population mean**.
- Correspondingly, the standard deviation of the sample,  $s$ , will stabilize around a fixed value,  $\sigma$ , which is usually unknown. The value  $\sigma$  is called the **population standard deviation**.
- Notation:
  - $\mu$  (mu) denotes the population mean.
  - $\sigma$  (sigma) denotes the population standard deviation.

---

Population	Sample
$\mu$	$\bar{y}$
$\sigma$	$s$

---

### 11.2.1 Distribution of a discrete random variable

- Possible values for  $Y$ :  $\{y_1, y_2, \dots, y_k\}$ .
- The **distribution** of  $Y$  is the probabilities of each possible value:  $p_i = P(Y = y_i)$ ,  $i = 1, 2, \dots, k$ .
- The distribution satisfies:  $\sum_{i=1}^k p_i = 1$ .

## 11.3 Expected value (mean) for a discrete distribution

- The **expected value** or **(population) mean** of  $Y$  is

$$\mu = \sum_{i=1}^k y_i p_i$$

- An important property of the expected value is that it has the same unit as the observations (e.g. meter).

### 11.3.1 Example: number of heads in 3 coin flips

- Recall the distribution of  $Y$  (number of heads):

---

y (number of heads)	0	1	2	3
$P(Y = y)$	1/8	3/8	3/8	1/8

---

- Then the expected value is

$$\mu = 0\frac{1}{8} + 1\frac{3}{8} + 2\frac{3}{8} + 3\frac{1}{8} = 1.5.$$

Note that the expected value is not a possible outcome of the experiment itself.

## 11.4 Variance and standard deviation for a discrete distribution

- The **(population) variance** of  $Y$  is

$$\sigma^2 = \sum_{i=1}^k (y_i - \mu)^2 p_i$$

- The **(population) standard deviation** is  $\sigma = \sqrt{\sigma^2}$ .
- Note: If the observations have unit meter, the **variance** has unit meter<sup>2</sup> which is hard to interpret. The **standard deviation** on the other hand has the same unit as the observations (e.g. meter).

### 11.4.1 Example: number of heads in 3 coin flips

The distribution of the random variable ‘number of heads in 3 coin flops’ has variance

$$\sigma^2 = (0 - 1.5)^2 \frac{1}{8} + (1 - 1.5)^2 \frac{3}{8} + (2 - 1.5)^2 \frac{3}{8} + (3 - 1.5)^2 \frac{1}{8} = 0.75.$$

and standard deviation

$$\sigma = \sqrt{\sigma^2} = \sqrt{0.75} = 0.866.$$

## 11.5 The binomial distribution

- The **binomial distribution** is a discrete distribution
- The distribution occurs when we conduct a success/failure experiment  $n$  times with probability  $\pi$  for success. If  $Y$  denotes the number of successes it can be shown that

$$p_Y(y) = P(Y = y) = \binom{n}{y} \pi^y (1 - \pi)^{n-y},$$

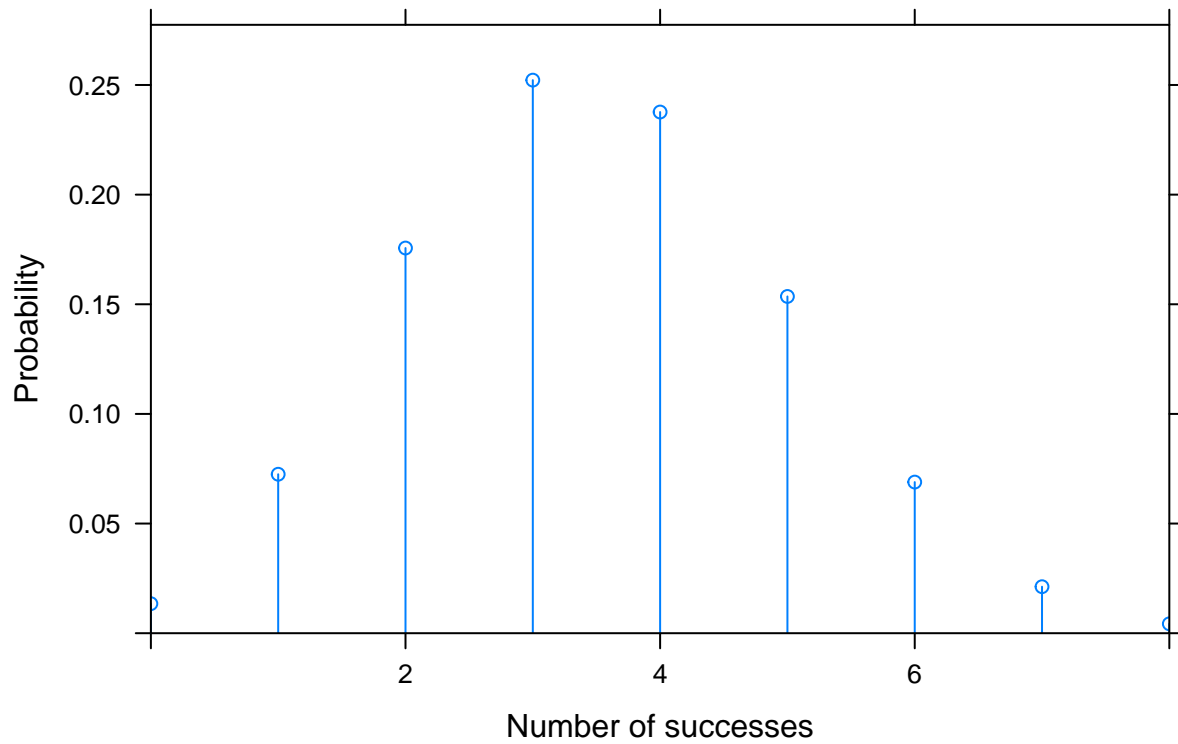
where  $\binom{n}{y} = \frac{n!}{y!(n-y)!}$  and  $m!$  is the product of the first  $m$  integers.

- Expected value:  $\mu = n\pi$ .
- Variance:  $\sigma^2 = n\pi(1 - \pi)$ .
- Standard deviation:  $\sigma = \sqrt{n\pi(1 - \pi)}$ .

*# The binomial distribution with  $n = 10$  and  $\pi = 0.35$ :*

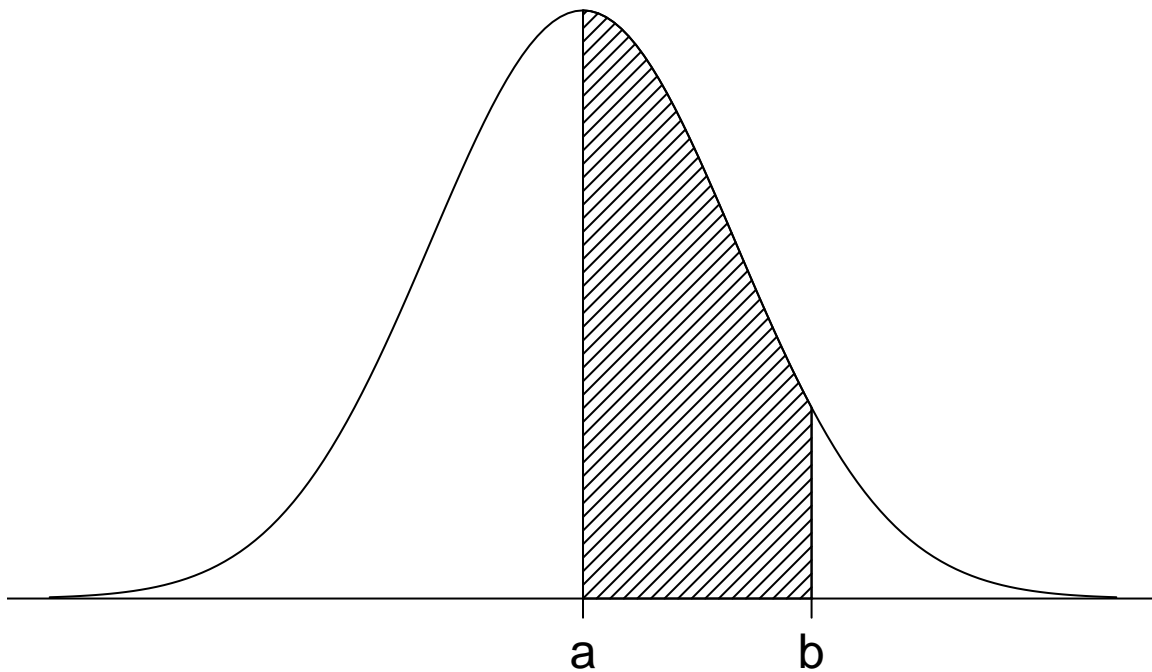
```
plotDist("binom", size = 10, prob = 0.35,
        ylab = "Probability", xlab = "Number of successes", main = "binom(n = 10, prob = 0.35)")
```

### binom(n = 10, prob = 0.35)



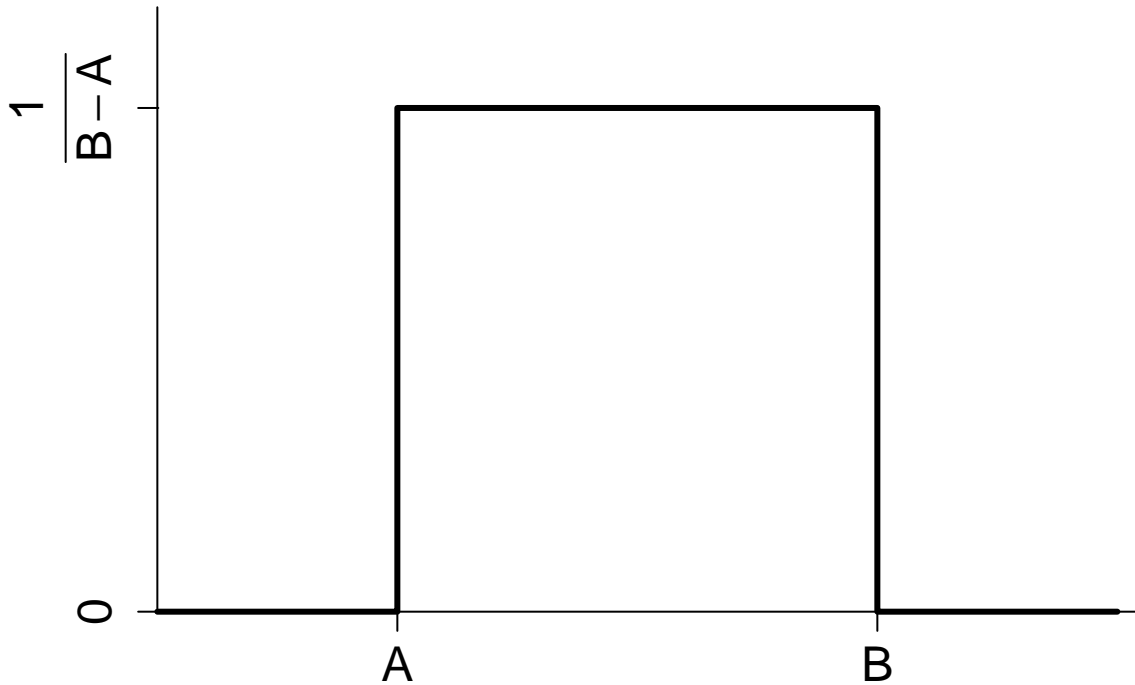
### 11.6 Distribution of a continuous random variable

- The distribution of a continuous random variable  $Y$  is characterized by the so-called probability density function  $f_Y$ .



- The area under the graph of the probability density function between  $a$  and  $b$  is equal to the probability of an observation in this interval.
- $f_Y(y) \geq 0$  for all real numbers  $y$ .
- The area under the graph for  $f_Y$  is equal to 1.
- For example the **uniform distribution** from  $A$  to  $B$ :

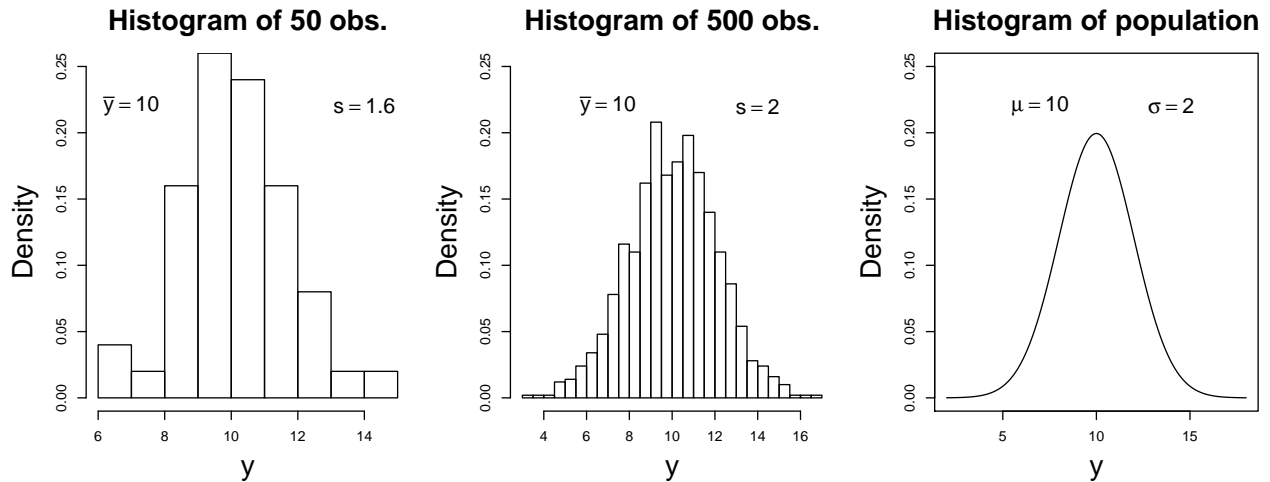
$$f_Y(y) = \begin{cases} \frac{1}{B-A} & A < y < B \\ 0 & \text{otherwise} \end{cases}$$



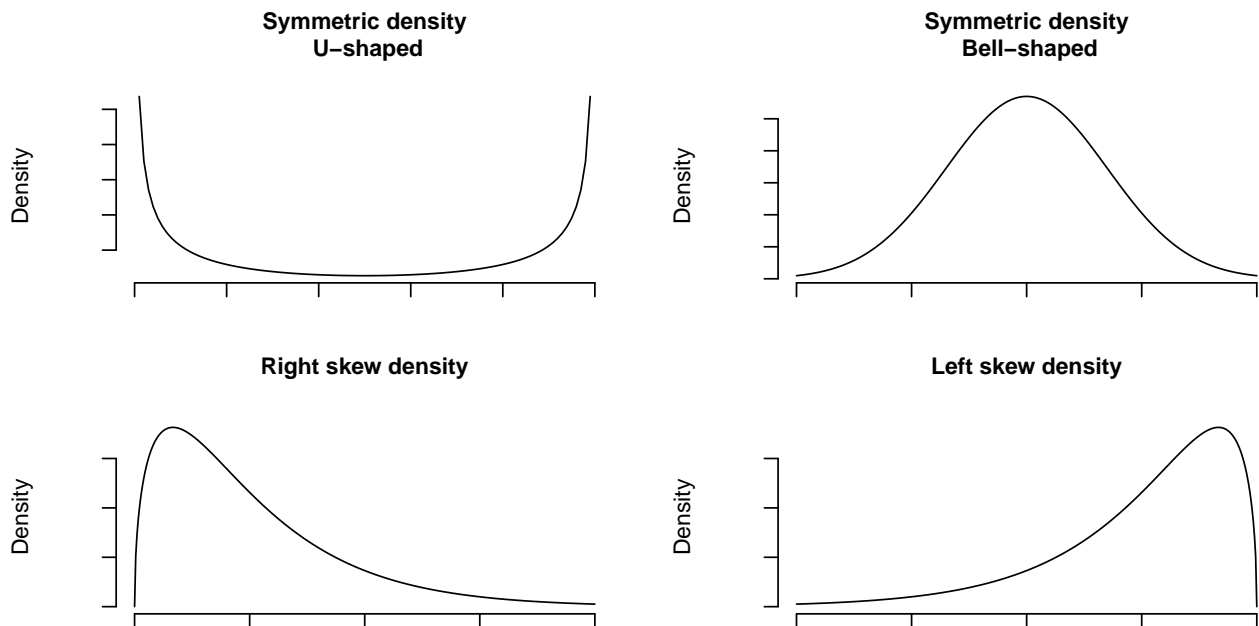
## 11.7 Density function

### 11.7.1 Increasing number of observations

- Another way to think about the density is in terms of the histogram.
- If we draw a histogram for a sample where the area of each box corresponds to the relative frequency of each interval, then the total area will be 1.
- When the number of observations (sample size) increase we can make a finer interval division and get a more smooth histogram.
- We can imagine an infinite number of observations, which would produce a nice smooth curve, where the area below the curve is 1. A function derived this way is also what we call the **probability density function**.



### 11.7.2 Density shapes



## 11.8 Normal distribution

- The normal distribution is a continuous distribution determined by two parameters:
  - $\mu$ : the **mean** (expected value), which determines where the distribution is centered.
  - $\sigma$ : the **standard deviation**, which determines the spread of the distribution about the mean.
- The distribution has a bell-shaped probability density function:

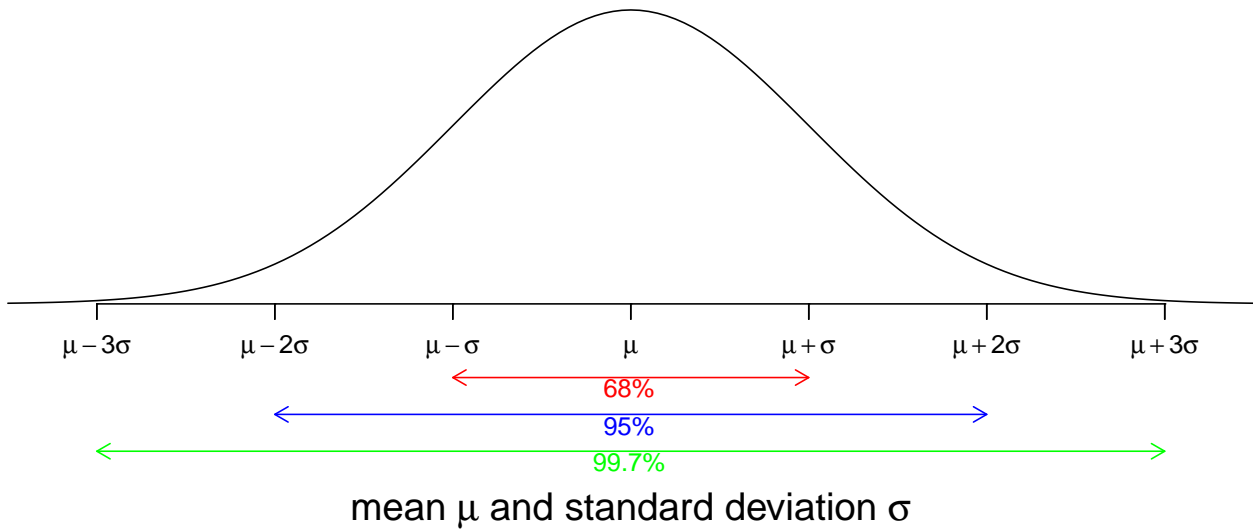
$$f_Y(y; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y - \mu)^2\right)$$

- When a random variable  $Y$  follows a normal distribution with mean  $\mu$  and standard deviation  $\sigma$ , then we write  $Y \sim \text{norm}(\mu, \sigma)$ .

- We call  $\text{norm}(0, 1)$  the **standard normal distribution**.

### 11.8.1 Reach of the normal distribution

#### Density of the normal distribution



Interpretation of standard deviation:

- $\approx 68\%$  of the population is within 1 standard deviation of the mean.
- $\approx 95\%$  of the population is within 2 standard deviations of the mean.
- $\approx 99.7\%$  of the population is within 3 standard deviations of the mean.

### 11.8.2 Normal $z$ -score

- If  $Y \sim \text{norm}(\mu, \sigma)$  then the corresponding  $z$ -score is

$$Z = \frac{Y - \mu}{\sigma} = \frac{\text{observation} - \text{mean}}{\text{standard deviation}}$$

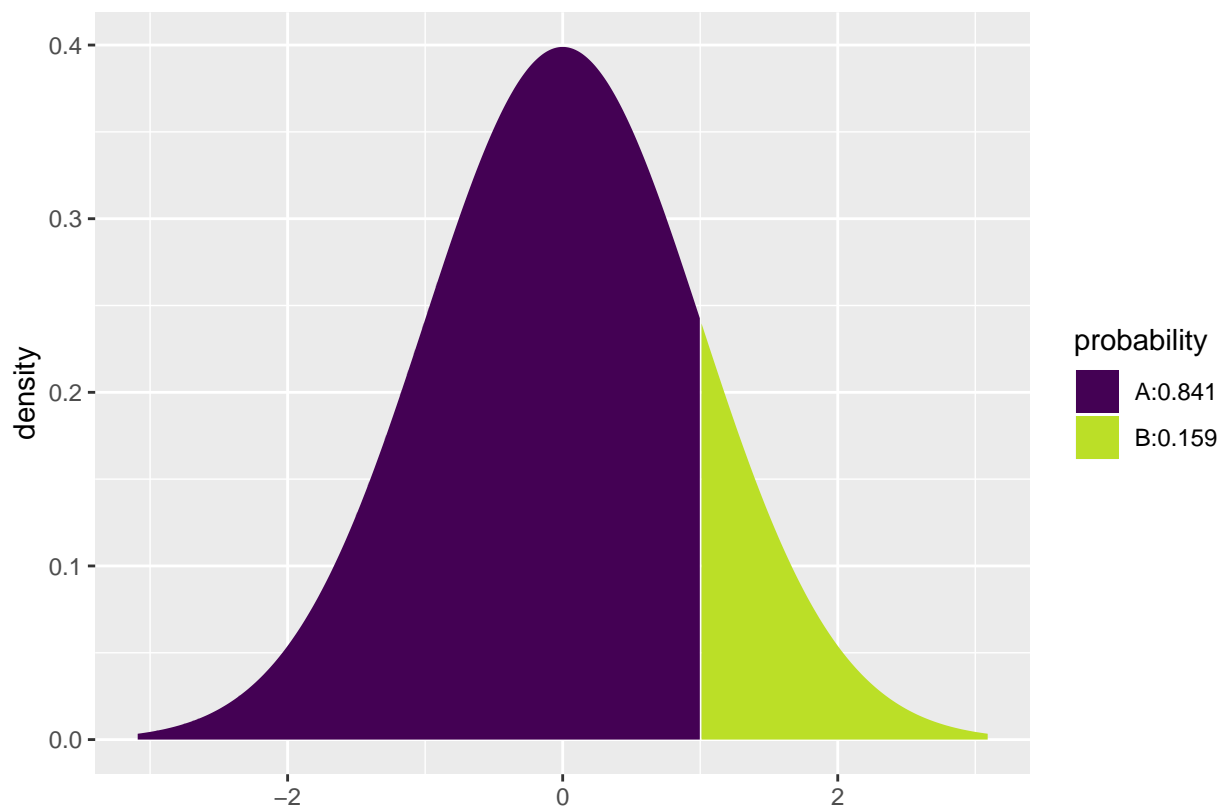
- I.e.  $Z$  counts the number of standard deviations that the observation lies away from the mean, where a negative value tells that we are below the mean.
- We have that  $Z \sim \text{norm}(0, 1)$ , i.e.  $Z$  has zero mean and standard deviation one.
- This implies that
  - $Z$  lies between  $-1$  and  $1$  with probability  $68\%$
  - $Z$  lies between  $-2$  and  $2$  with probability  $95\%$
  - $Z$  lies between  $-3$  and  $3$  with probability  $99.7\%$
- It also implies that:
  - The probability of  $Y$  being between  $\mu - z\sigma$  and  $\mu + z\sigma$  is equal to the probability of  $Z$  being between  $-z$  and  $z$ .



### 11.8.3 Calculating probabilities in the standard normal distribution

- The function `pnorm` always outputs the area to the left of the  $z$ -value (quantile/percentile) we give as input (variable `q` in the function), i.e. it outputs the probability of getting a value less than  $z$ . The first argument of `pnorm` denotes the distribution we are considering.

```
# For a standard normal distribution the probability of getting a value less than 1 is:  
left_prob <- pnorm("norm", q = 1, mean = 0, sd = 1)
```



```
left_prob
```

```
## [1] 0.84
```

- Here there is a conflict between **R** and the textbook, since in the book we always consider right probabilities in the normal distribution. Since the total area is 1 and we have the left probability we easily get the right probability:

```
right_prob <- 1 - left_prob  
right_prob
```

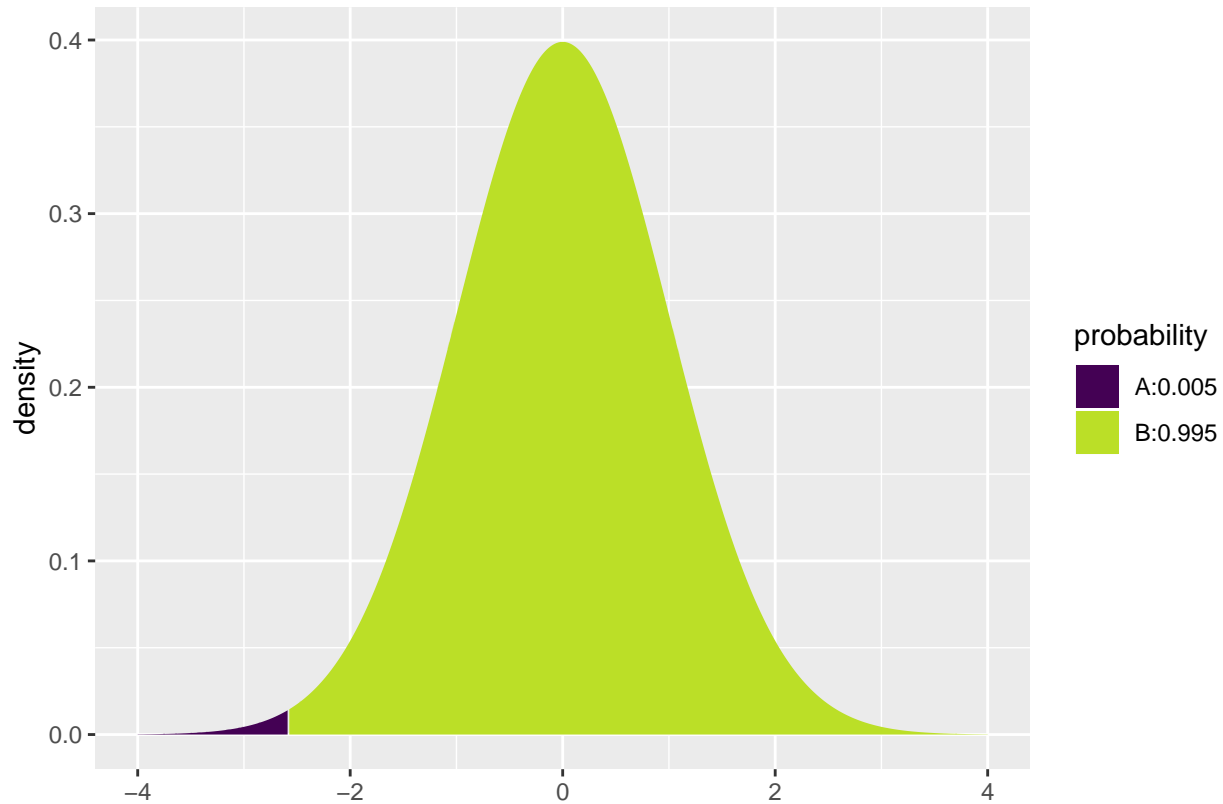
```
## [1] 0.16
```

- For  $z = 1$  we have a right probability of  $p = 0.1587$ , so the probability of an observation between  $-1$  and  $1$  is  $1 - 2 \cdot 0.1587 = 0.6826 = 68.26\%$  due to symmetry.

#### 11.8.4 Calculating $z$ -values (quantiles) in the standard normal distribution

- If we have a probability and want to find the corresponding  $z$ -value we again need to decide on left/right probability. The default in **R** is to find the left probability, so if we want the  $z$ -value with e.g. 0.5% probability to the left we get:

```
left_z <- qdist("norm", p = 0.005, mean = 0, sd = 1, xlim = c(-4, 4))
```

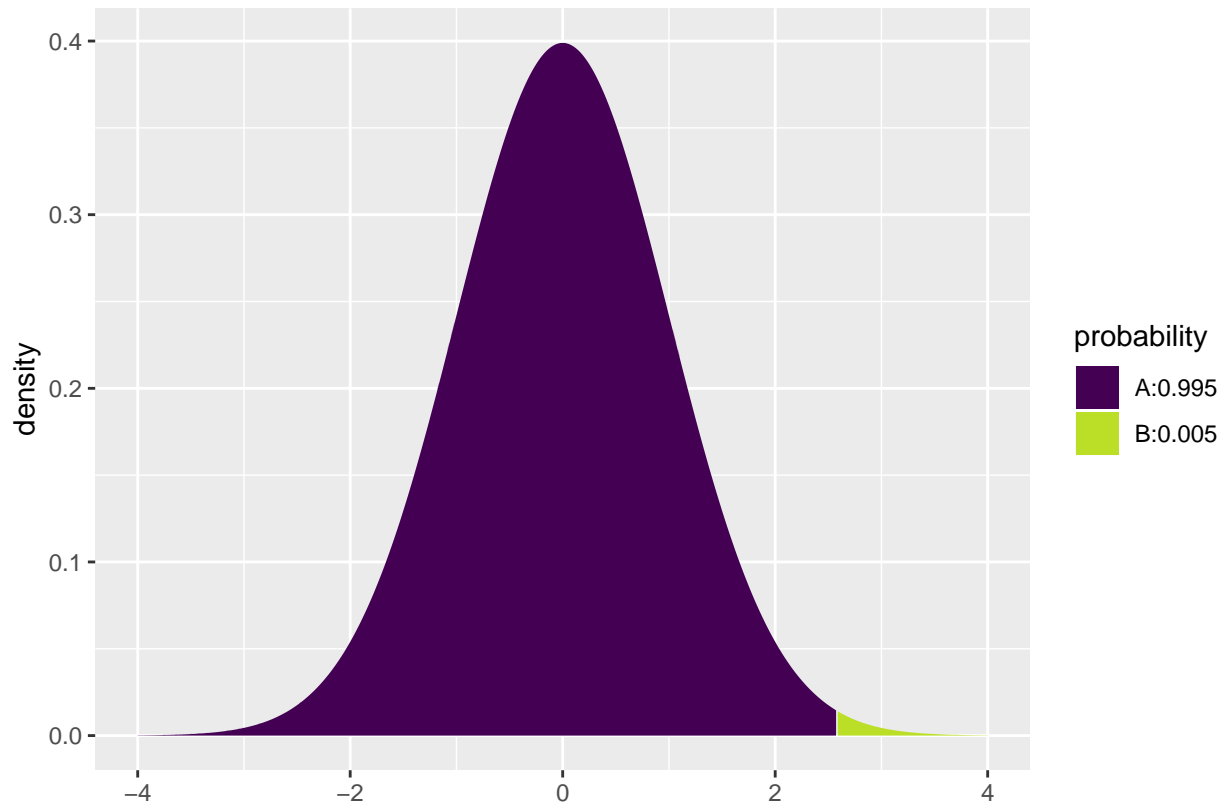


```
left_z
```

```
## [1] -2.6
```

- However, in all the formulas in the course we follow the textbook and consider  $z$ -values for a given right probability. E.g. with 0.5% probability to the right we get:

```
right_z <- qdist("norm", p = 1-0.005, mean = 0, sd = 1, xlim = c(-4, 4))
```



```
right_z
```

```
## [1] 2.6
```

- Thus, the probability of an observation between  $-2.576$  and  $2.576$  equals  $1 - 2 \cdot 0.005 = 99\%$ .

### 11.8.5 Example

The Stanford-Binet Intelligence Scale is calibrated to be approximately normal with mean 100 and standard deviation 16.

What is the 99-percentile of IQ scores?

- The corresponding  $z$ -score is  $Z = \frac{IQ-100}{16}$ , which means that  $IQ = 16Z + 100$ .
- The 99-percentile of  $z$ -scores has the value 2.326 (can be calculated using `qdist`).
- Then, the 99-percentile of IQ scores is:

$$IQ = 16 \cdot 2.326 + 100 = 137.2.$$

- So we expect that one out of hundred has an IQ exceeding 137.

## 12 Distribution of sample statistic

### 12.1 Estimates and their variability

We are given a sample  $y_1, y_2, \dots, y_n$ .

- The sample mean  $\bar{y}$  is the most common estimate of the population mean  $\mu$ .
- The sample standard deviation,  $s$ , is the most common estimate of the population standard deviation  $\sigma$ .

We notice that there is an uncertainty (from sample to sample) connected to these statistics and therefore we are interested in describing their **distribution**.

## 12.2 Distribution of sample mean

- We are given a sample  $y_1, y_2, \dots, y_n$  from a population with mean  $\mu$  and standard deviation  $\sigma$ .
- The sample mean

$$\bar{y} = \frac{1}{n}(y_1 + y_2 + \dots + y_n)$$

then has a distribution where

- the distribution has mean  $\mu$ ,
- the distribution has standard deviation  $\frac{\sigma}{\sqrt{n}}$  (also called the **standard error**), and
- when  $n$  grows, the distribution approaches a normal distribution. This result is called **the central limit theorem**.

### 12.2.1 Central limit theorem

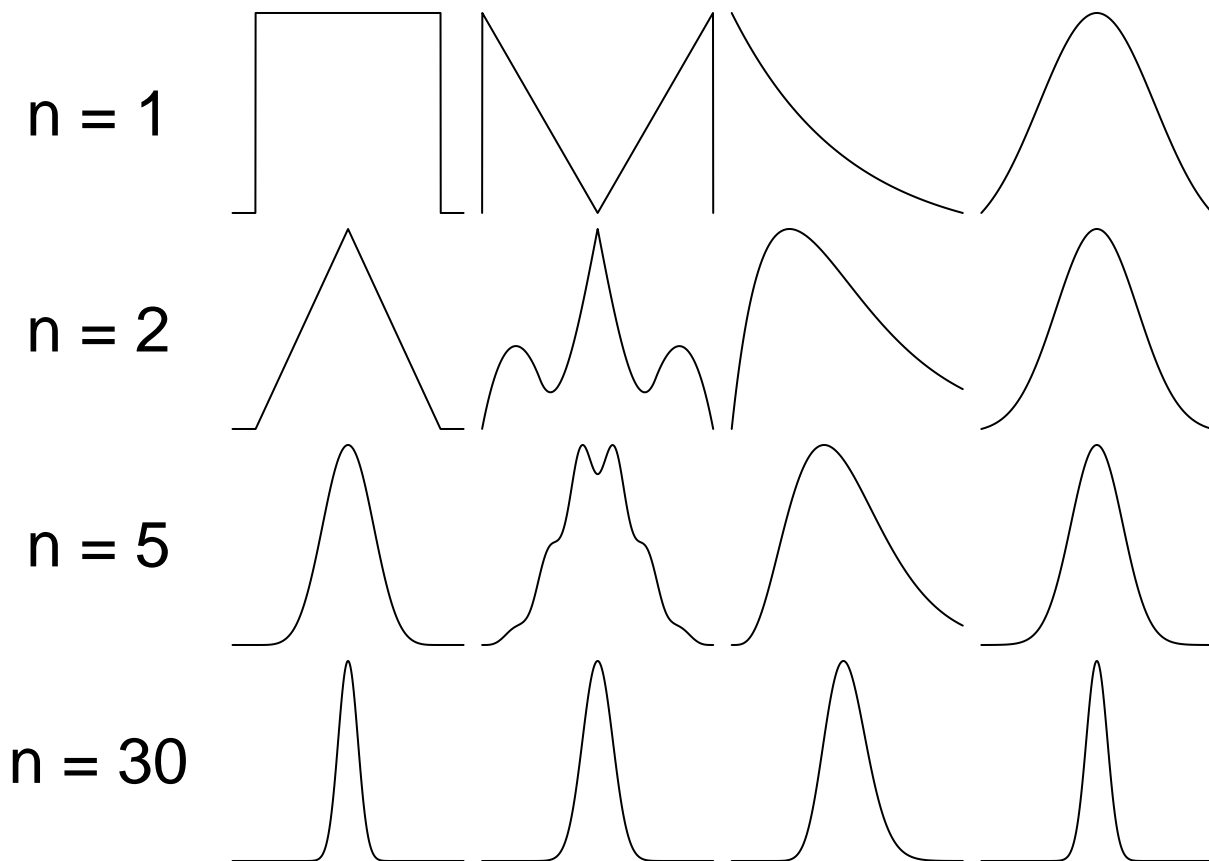
- The points above can be summarized as

$$\bar{y} \approx \text{norm}\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

i.e.  $\bar{y}$  is approximately normally distributed with mean  $\mu$  and standard error  $\frac{\sigma}{\sqrt{n}}$ .

- When our sample is sufficiently large (such that the above approximation is good) this allows us to make the following observations:
  - We are 95% certain that  $\bar{y}$  lies in the interval from  $\mu - 2\frac{\sigma}{\sqrt{n}}$  to  $\mu + 2\frac{\sigma}{\sqrt{n}}$ .
  - We are almost completely certain that  $\bar{y}$  lies in the interval from  $\mu - 3\frac{\sigma}{\sqrt{n}}$  to  $\mu + 3\frac{\sigma}{\sqrt{n}}$ .
- This is not useful when  $\mu$  is unknown, but let us rephrase the first statement to:
  - We are 95% certain that  $\mu$  lies in the interval from  $\bar{y} - 2\frac{\sigma}{\sqrt{n}}$  to  $\bar{y} + 2\frac{\sigma}{\sqrt{n}}$ , i.e. we are directly talking about the uncertainty of determining  $\mu$ .

### 12.2.2 Illustration of CLT



- Four different population distributions ( $n=1$ ) of  $y$  and corresponding sampling distributions of  $\bar{y}$  for different sample sizes. As  $n$  increases the sampling distributions become narrower and more bell-shaped.

### 12.2.3 Example

- Body Mass Index (BMI) of people in Northern Jutland (2010) has mean  $\mu = 25.8 \text{ kg/m}^2$  and standard deviation  $4.8 \text{ kg/m}^2$ .
- A random sample of  $n = 100$  costumers at a burger bar had an average BMI given by  $\bar{y} = 27.2$ .
- If “burger bar” has “no influence” on BMI (and the sample is representative of the population/people in Northern Jutland), then

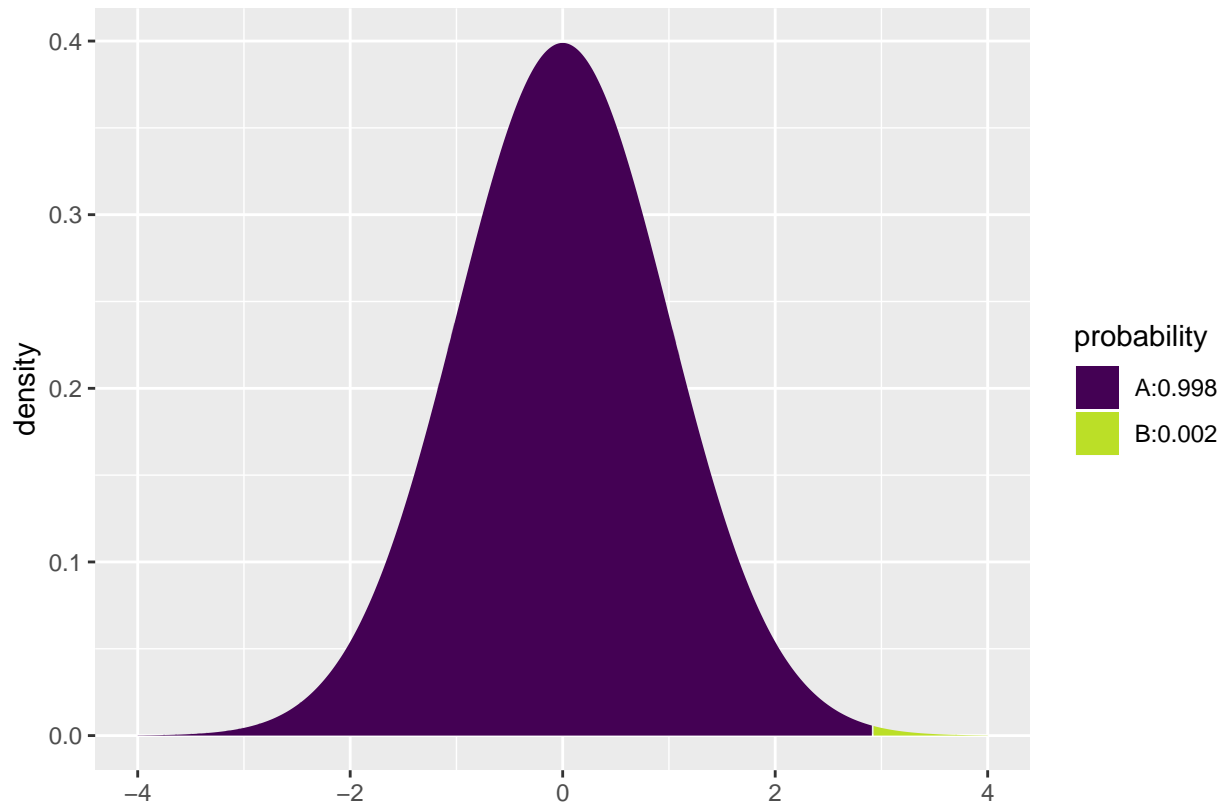
$$\bar{y} \approx \text{norm}\left(\mu, \frac{\sigma}{\sqrt{n}}\right) = \text{norm}(25.8, 0.48).$$

- For the actual sample this gives the observed  $z$ -score

$$z_{obs} = \frac{27.2 - 25.8}{0.48} = 2.92$$

- Recalling that the  $z$ -score is (here approximately) standard normal, the probability of getting a higher  $z$ -score is:

```
1 - pdist("norm", mean = 0, sd = 1, q = 2.92, xlim = c(-4, 4))
```



## [1] 0.0018

- Thus, it is highly unlikely to get a random sample with such a high  $z$ -score. This indicates that costumers at the burger bar has a mean BMI, which is higher than the population mean.

## 13 Point and interval estimates

### 13.1 Point and interval estimates

- We want to study hypotheses for population parameters, e.g. the mean  $\mu$  and the standard deviation  $\sigma$ .
  - If  $\mu$  is e.g. the mean waiting time in a queue, then it might be relevant to investigate whether it exceeds 2 minutes.
- Based on a sample we make a **point estimate** which is a guess of the parameter value.
  - For instance we have used  $\bar{y}$  as an estimate of  $\mu$  and  $s$  as an estimate of  $\sigma$ .
- We often want to supplement the point estimate with an **interval estimate** (also called a **confidence interval**). This is an interval around the point estimate, in which we are confident (to a certain degree) that the population parameter is located.
- The parameter estimate can then be used to investigate our hypothesis.

### 13.2 Point estimators: Bias

- If we want to estimate the population mean  $\mu$  we have several possibilities e.g.
  - the sample mean  $\bar{y}$

- the average  $y_T$  of the sample upper and lower quartiles
- Advantage of  $y_T$ : Very large/small observations have little effect, i.e. it has practically no effect if there are a few errors in the data set.
- Disadvantage of  $y_T$ : If the distribution of the population is skewed, i.e. asymmetrical, then  $y_T$  is **biased**, meaning that in the long run this estimator systematically over or under estimates the value of  $\mu$ .
- Generally we prefer that an estimator is **unbiased**, i.e. its distribution is centered around the true parameter value.
- Recall that for a sample from a population with mean  $\mu$ , the sample mean  $\bar{y}$  also has mean  $\mu$ . That is,  $\bar{y}$  is an unbiased estimate of the population mean  $\mu$ .

### 13.3 Point estimators: Consistency

- From previous lectures we know that the standard error of  $\bar{y}$  is  $\frac{\sigma}{\sqrt{n}}$ ,
  - i.e. the standard error decrease when the sample size increase.
- In general an estimator with this property is called **consistent**.
- $y_T$  is also a consistent estimator, but has a variance that is greater than  $\bar{y}$ .

### 13.4 Point estimators: Efficiency

- Since the variance of  $y_T$  is greater than the variance of  $\bar{y}$ ,  $\bar{y}$  is preferred.
- In general we prefer the estimator with the smallest possible variance.
  - This estimator is said to be **efficient**.
- $\bar{y}$  is an efficient estimator.

### 13.5 Notation

- The symbol  $\hat{\cdot}$  above a parameter is often used to denote a (point) estimate of the parameter. We have looked at an
  - estimate of the population mean:  $\hat{\mu} = \bar{y}$
  - estimate of the population standard deviation:  $\hat{\sigma} = s$
- When we observe a 0/1 variable, which e.g. is used to denote yes/no or male/female, then we will use the notation

$$\pi = P(Y = 1)$$

for the proportion of the population with the property  $Y = 1$ .

- The estimate  $\hat{\pi} = (y_1 + y_2 + \dots + y_n)/n$  is the relative frequency of the property  $Y = 1$  in the sample.

### 13.6 Confidence Interval

- The general definition of a confidence interval for a population parameter is as follows:
  - A **confidence interval** for a parameter is constructed as an interval, where we expect the parameter to be.
  - The probability that this construction yields an interval which includes the parameter is called the **confidence level** and it is typically chosen to be 95%.
  - (1-confidence level) is called the **error probability** (in this case  $1 - 0.95 = 0.05$ , i.e. 5%).
- In practice the interval is often constructed as a symmetric interval around a point estimate:

- **point estimate** ± **margin of error**
- Rule of thumb: With a margin of error of 2 times the standard error you get a confidence interval, where the confidence level is approximately 95%.
- I.e: **point estimate** ± **2 x standard error** has confidence level of approximately 95%.

## 13.7 Confidence interval for proportion

- Consider a population with a distribution where the probability of having a given characteristic is  $\pi$  and the probability of not having it is  $1 - \pi$ .
- When *no/yes* to the characteristic is denoted 0/1, i.e.  $y$  is 0 or 1, the distribution of  $y$  have a standard deviation of:

$$\sigma = \sqrt{\pi(1 - \pi)}.$$

That is, the standard deviation is not a “free” parameter for a 0/1 variable as it is directly linked to the probability  $\pi$ .

- With a sample size of  $n$  the standard error of  $\hat{\pi}$  will be (since  $\hat{\pi} = \frac{\sum_{i=1}^n y_i}{n}$ ):

$$\sigma_{\hat{\pi}} = \frac{\sigma}{\sqrt{n}} = \sqrt{\frac{\pi(1 - \pi)}{n}}.$$

- We do not know  $\pi$  but we insert the estimate and get the **estimated standard error** of  $\hat{\pi}$ :

$$se = \sqrt{\frac{\hat{\pi}(1 - \hat{\pi})}{n}}.$$

- The rule of thumb gives that the interval

$$\hat{\pi} \pm 2\sqrt{\frac{\hat{\pi}(1 - \hat{\pi})}{n}}$$

has confidence level of approximately 95%. I.e., before the data is known the random interval given by the formula above has approximately 95% probability of covering the true value  $\pi$ .

### 13.7.1 Example: Point and interval estimate for proportion

- Now we will have a look at a data set concerning votes in Chile. Information about the data can be found here.

```
Chile <- read.delim("https://asta.math.aau.dk/datasets?file=Chile.txt")
```

- We focus on the variable `sex`, i.e. the gender distribution in the sample.

```
library(mosaic)
tally(~ sex, data = Chile)
```

```
## sex
##   F   M
## 1379 1321
```



```
tally( ~ sex, data = Chile, format = "prop")
```

```
## sex
##   F   M
## 0.51 0.49
```

- Unknown population proportion of females (F),  $\pi$ .
  - Estimate of  $\pi$ :  $\hat{\pi} = \frac{1379}{1379+1321} = 0.5107$
  - Rule of thumb :  $\hat{\pi} \pm 2 \times se = 0.5107 \pm 2\sqrt{\frac{0.5107(1-0.5107)}{1379+1321}} = (0.49, 0.53)$  is an approximate 95% confidence interval for  $\pi$ .
- 

### 13.7.2 Example: Confidence intervals for proportion in R

- **R** automatically calculates the confidence interval for the proportion of females when we do a so-called hypothesis test (we will get back to that later):

```
prop.test( ~ sex, data = Chile, correct = FALSE)
```

```
##
## 1-sample proportions test without continuity correction
##
## data:  Chile$sex [with success = F]
## X-squared = 1, df = 1, p-value = 0.3
## alternative hypothesis: true p is not equal to 0.5
## 95 percent confidence interval:
##  0.49 0.53
## sample estimates:
##      p
## 0.51
```

- The argument `correct = FALSE` is needed to make **R** use the “normal” formulas as on the slides and in the book. When `correct = TRUE` (the default) a mathematical correction which you have not learned about is applied and slightly different results are obtained.
- 

## 13.8 General confidence intervals for proportion

- Based on the central limit theorem we have:

$$\hat{\pi} \approx N\left(\pi, \sqrt{\frac{\pi(1-\pi)}{n}}\right)$$

if  $n\hat{\pi}$  and  $n(1-\hat{\pi})$  both are at least 15.

- To construct a confidence interval with (approximate) confidence level  $1 - \alpha$ :

- 1) Find the so-called **critical value**  $z_{crit}$  for which the upper tail probability in the standard normal distribution is  $\alpha/2$ .

- 2) Calculate  $se = \sqrt{\frac{\hat{\pi}(1-\hat{\pi})}{n}}$
- 3) Then  $\hat{\pi} \pm z_{crit} \times se$  is a confidence interval with confidence level  $1 - \alpha$ .

### 13.8.1 Example: Chile data

Compute for the **Chile** data set the 99% and 95%-confidence intervals for the probability that a person is female:

- For a 99%-confidence level we have  $\alpha = 1\%$  and
  - 1)  $z_{crit} = \text{qdist}(\text{"norm"}, 1 - 0.01/2) = 2.576$ .
  - 2) We know that  $\hat{\pi} = 0.5107$  and  $n = 2700$ , so  $se = \sqrt{\frac{\hat{\pi}(1-\hat{\pi})}{n}} = 0.0096$ .
  - 3) Thereby, a 99%-confidence interval is:  $\hat{\pi} \pm z_{crit} \times se = (0.4859, 0.5355)$ .
- For a 95%-confidence level we have  $\alpha = 5\%$  and
  - 1)  $z_{crit} = \text{qdist}(\text{"norm"}, 1 - 0.05/2) = 1.96$ .
  - 2) Again,  $\hat{\pi} = 0.5107$  and  $n = 2700$  and so  $se = 0.0096$ .
  - 3) Thereby, we find as 95%-confidence interval  $\hat{\pi} \pm z_{crit} \times se = (0.4918, 0.5295)$  (as the result of `prop.test`).

## 13.9 Confidence Interval for mean - normally distributed sample

- When it is reasonable to assume that the population distribution is normal we have the **exact** result

$$\bar{y} \sim \text{norm}\left(\mu, \frac{\sigma}{\sqrt{n}}\right),$$

i.e.  $\bar{y} \pm z_{crit} \times \frac{\sigma}{\sqrt{n}}$  is not only an approximate but rather an exact confidence interval for the population mean,  $\mu$ .

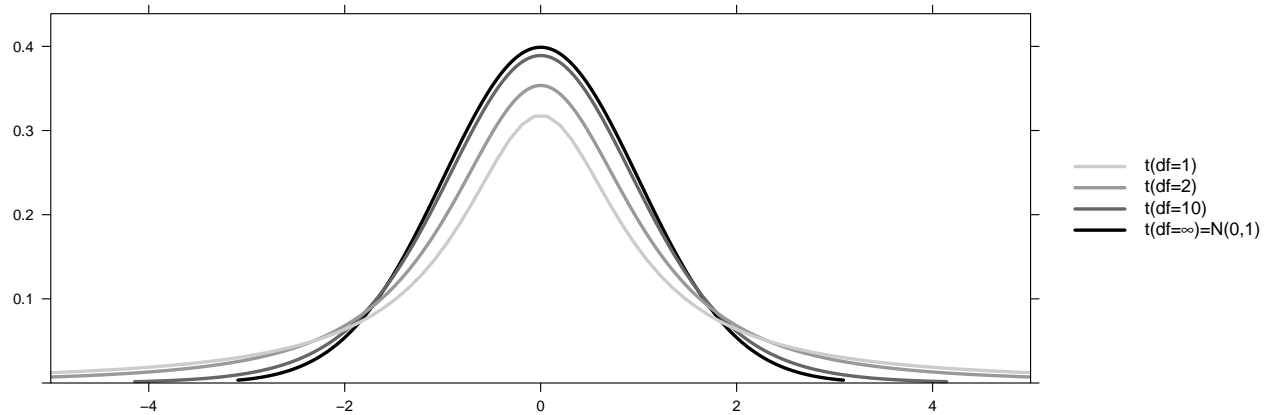
- In practice we **do not know**  $\sigma$  and instead we are forced to apply the sample standard deviation  $s$  to find the **estimated standard error**  $se = \frac{s}{\sqrt{n}}$ .
- This extra uncertainty, however, implies that an exact confidence interval for the population mean  $\mu$  cannot be constructed using the  $z$ -score.
- Luckily, an exact interval can still be constructed by using the so-called  **$t$ -score**, which apart from the confidence level depends on the **degrees of freedom**, which are  $df = n - 1$ . That is the confidence interval now takes the form

$$\bar{y} \pm t_{crit} \times se.$$

### 13.10 $t$ -distribution and $t$ -score

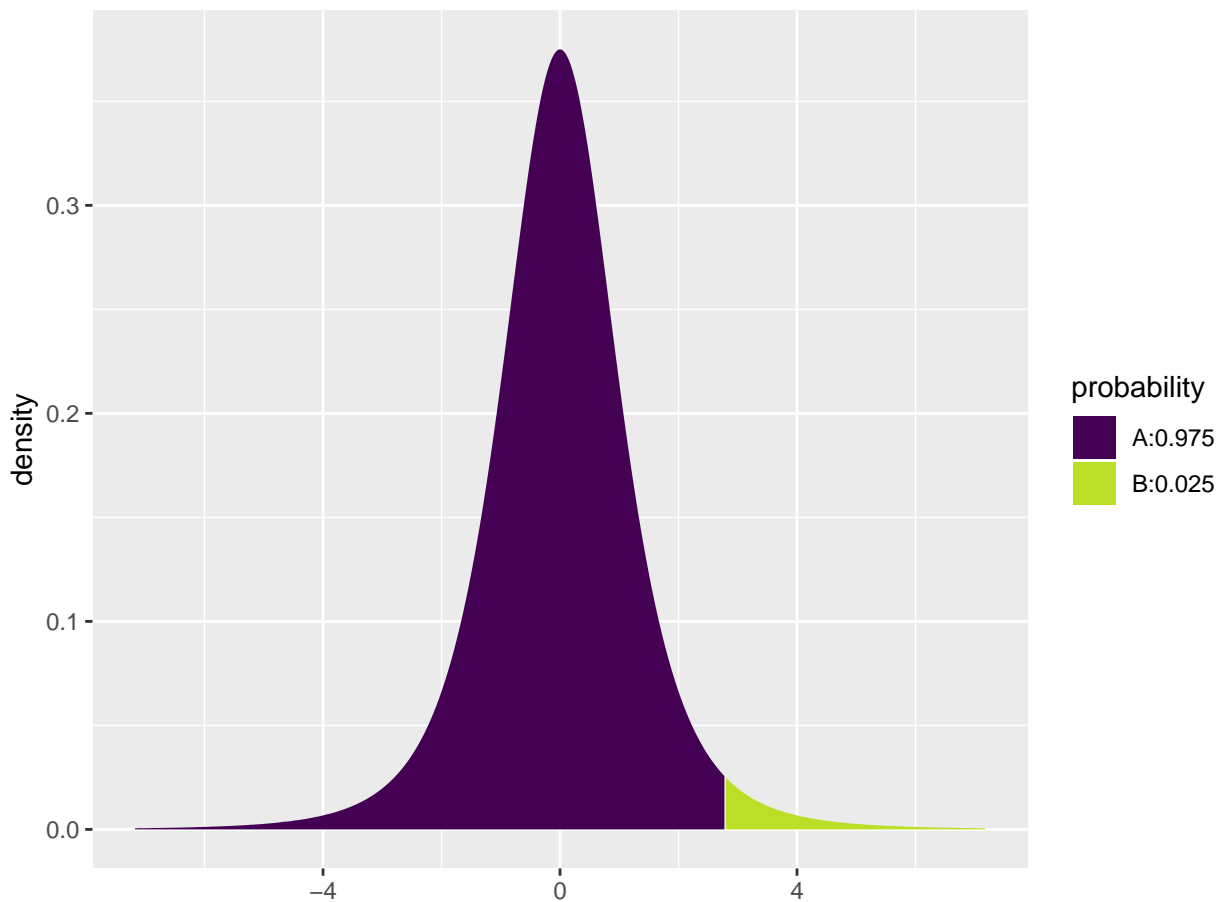
- Calculation of  $t$ -score is based on the  **$t$ -distribution**, which is very similar to the standard normal  $z$ -distribution:
  - it is symmetric around zero and bell shaped, but
  - it has “heavier” tails and thereby
  - a slightly larger standard deviation than the standard normal distribution.
  - Further, the  $t$ -distribution’s standard deviation decays as a function of its **degrees of freedom**, which we denote  $df$ .
  - and when  $df$  grows the  $t$ -distribution approaches the standard normal distribution.

The expression of the density function is of slightly complicated form and will not be stated here, instead the  $t$ -distribution is plotted below for  $df = 1, 2, 10$  and  $\infty$ .



### 13.11 Calculation of $t$ -score in R

```
qdist("t", p = 1 - 0.025, df = 4)
```



```
## [1] 2.8
```

- We seek the quantile (i.e. value on the x-axis) such that we have a given **right tail** probability. This is the critical t-score associated with our desired level of confidence.
- To get e.g. the  $t$ -score corresponding to a right tail probability of 2.5 % we have to look up the 97.5 % quantile using `qdist` with  $p = 1 - 0.025$  since `qdist` looks at the area to the **left hand side**.
- The degrees of freedom are determined by the sample size; in this example we just used  $df = 4$  for illustration.
- As the  $t$ -score giving a right probability of 2.5 % is 2.776 and the  $t$ -distribution is symmetric around 0, we have that an observation falls between -2.776 and 2.776 with probability  $1 - 2 \cdot 0.025 = 95$  % for a  $t$ -distribution with 4 degrees of freedom.

### 13.12 Example: Confidence interval for mean

- We return to the dataset `Ericksen` and want to construct a 95% confidence interval for the population mean  $\mu$  of the variable `crime`.

```
Ericksen <- read.delim("https://asta.math.aau.dk/datasets?file=Ericksen.txt")
stats <- favstats( ~ crime, data = Ericksen)
stats
```

```
## min Q1 median Q3 max mean sd n missing
## 25 48 55 73 143 63 25 66 0
```

```
qdist("t", 1 - 0.025, df = 66 - 1, plot = FALSE)
```

```
## [1] 2
```

- I.e. we have
  - $\bar{y} = 63.06$
  - $s = 24.89$
  - $n = 66$
  - $df = n - 1 = 65$
  - $t_{crit} = 2.$
- The confidence interval is  $\bar{y} \pm t_{crit} \frac{s}{\sqrt{n}} = (56.94, 69.18)$
- All these calculations can be done automatically by **R**:

```
t.test( ~ crime, data = Ericksen, conf.level = 0.95)
```

```
##
## One Sample t-test
##
## data: crime
## t = 20, df = 60, p-value <2e-16
## alternative hypothesis: true mean is not equal to 0
## 95 percent confidence interval:
## 57 69
## sample estimates:
## mean of x
## 63
```

### 13.13 Example: Plotting several confidence intervals in R

- We shall look at a built-in R dataset `chickwts`.
- `?chickwts` yields a page with the following information

An experiment was conducted to measure and compare the effectiveness of various feed supplements on the growth rate of chickens. Newly hatched chicks were randomly allocated into six groups, and each group was given a different feed supplement. Their weights in grams after six weeks are given along with feed types.

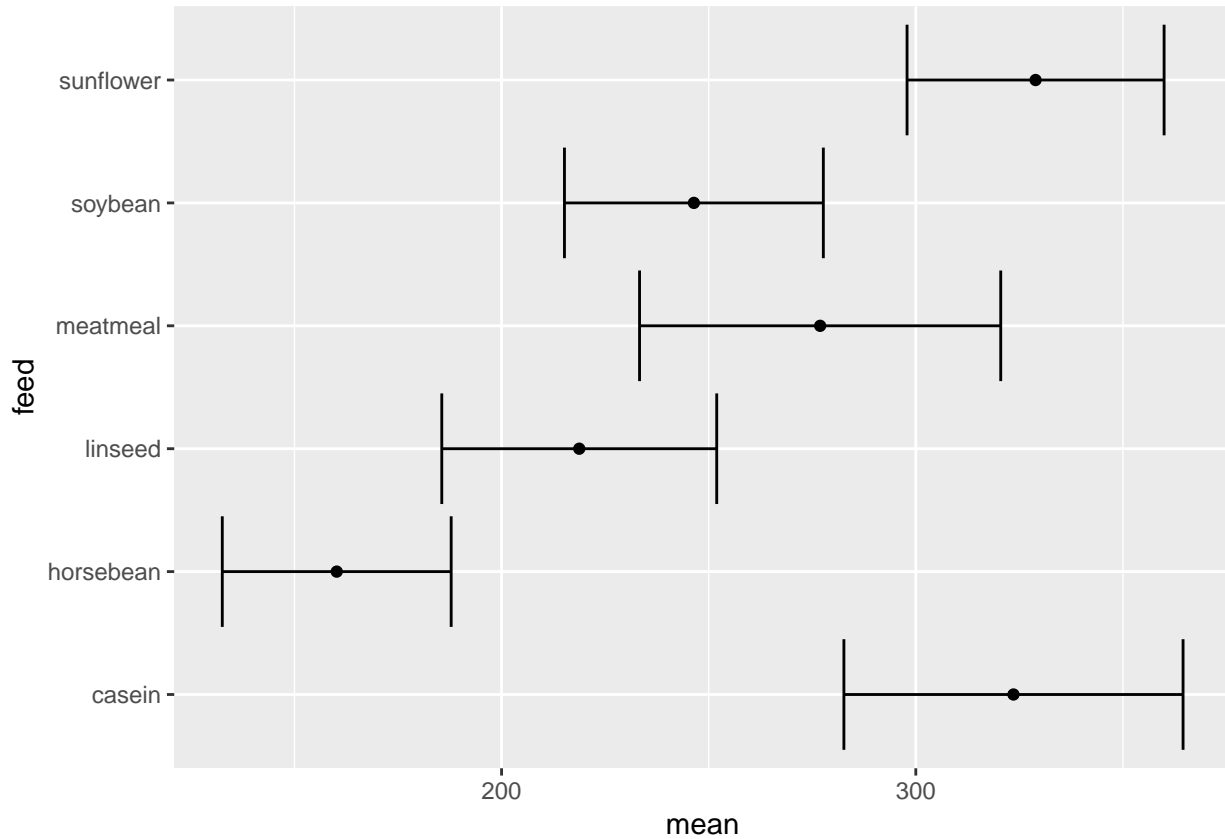
- `chickwts` is a data frame with 71 observations on 2 variables:
  - `weight`: a numeric variable giving the chick weight.
  - `feed`: a factor giving the feed type.
- Calculate a confidence interval for the mean weight for each feed separately; the confidence interval is from lower to upper given by `mean±tscore * se`:

```
cwei <- favstats( weight ~ feed, data = chickwts)
se <- cwei$sd / sqrt(cwei$n) # Standard errors
tscore <- qdist("t", p = .975, df = cwei$n - 1, plot = FALSE) # t-scores for 2.5% right tail probability
cwei$lower <- cwei$mean - tscore * se
cwei$upper <- cwei$mean + tscore * se
cwei[, c("feed", "mean", "lower", "upper")]
```

```
##      feed mean lower upper
## 1 casein  324   283   365
## 2 horsebean 160   133   188
## 3 linseed  219   186   252
## 4 meatmeal 277   233   321
## 5 soybean  246   215   278
## 6 sunflower 329   298   360
```

- We can plot the confidence intervals as horizontal line segments using `gf_errorbarh`:

```
gf_errorbarh(feed ~ lower + upper, data = cwei) %>%
  gf_point(feed ~ mean)
```



## 14 Determining sample size

### 14.1 Sample size for proportion

- The confidence interval is of the form point estimate  $\pm$  estimated margin of error.
- When we estimate a proportion the margin of error is

$$M = z_{crit} \sqrt{\frac{\pi(1-\pi)}{n}},$$

where the critical  $z$ -score,  $z_{crit}$ , is determined by the specified confidence level.

- Imagine that we want to plan an experiment, where we **want to achieve a certain margin of error**  $M$  (and thus a specific width of the associated confidence interval).
- If we solve the equation above we see:
  - If we choose sample size  $n = \pi(1-\pi)\left(\frac{z_{crit}}{M}\right)^2$ , then we obtain an estimate of  $\pi$  with margin of error  $M$ .
- If we do not have a good guess for the value of  $\pi$  we can use the worst case value  $\pi = 50\%$ . The corresponding sample size  $n = \left(\frac{z_{crit}}{2M}\right)^2$  ensures that we obtain an estimate with a margin of error, which is at the *most*  $M$ .

#### 14.1.1 Example

- Let us choose  $z_{crit} = 1.96$ , i.e the confidence level is 95%.

- How many voters should we ask to get a margin of error, which equals 1%?
- Worst case is  $\pi = 0.5$ , yielding:

$$n = \pi(1 - \pi) \left( \frac{z_{crit}}{M} \right)^2 = \frac{1}{4} \left( \frac{1.96}{0.01} \right)^2 = 9604.$$

- If we are interested in the proportion of voters that vote for “socialdemokratiet” a good guess is  $\pi = 0.23$ , yielding

$$n = \pi(1 - \pi) \left( \frac{z_{crit}}{M} \right)^2 = 0.23(1 - 0.23) \left( \frac{1.96}{0.01} \right)^2 = 6804.$$

- If we instead are interested in “liberal alliance” a good guess is  $\pi = 0.05$ , yielding

$$n = \pi(1 - \pi) \left( \frac{z_{crit}}{M} \right)^2 = 0.05(1 - 0.05) \left( \frac{1.96}{0.01} \right)^2 = 1825.$$

## 14.2 Sample size for mean

- The confidence interval is of the form point estimate  $\pm$  estimated margin of error.
- When we estimate a mean the margin of error is

$$M = z_{crit} \frac{\sigma}{\sqrt{n}},$$

where the critical  $z$ -score,  $z_{crit}$ , is determined by the specified confidence level.

- Imagine that we want to plan an experiment, where we **want to achieve a certain margin of error**  $M$ .
- If we solve the equation above we see:
  - If we choose sample size  $n = \left( \frac{z_{crit}\sigma}{M} \right)^2$ , then we obtain an estimate with margin of error  $M$ .
- Problem: We usually do not know  $\sigma$ . Possible solutions:
  - Based on similar studies conducted previously, we make a qualified guess at  $\sigma$ .
  - Based on a pilot study a value of  $\sigma$  is estimated.

## 15 Data collection

### 15.1 Data collection

- Getting numbers to report is easy
- Getting sensible and trustworthy numbers to report is orders of magnitude more difficult

Ronald Fisher (1890-1962):

To consult the statistician after an experiment is finished is often merely to ask him to conduct a post mortem examination. He can perhaps say what the experiment died of.

Said about Fisher:

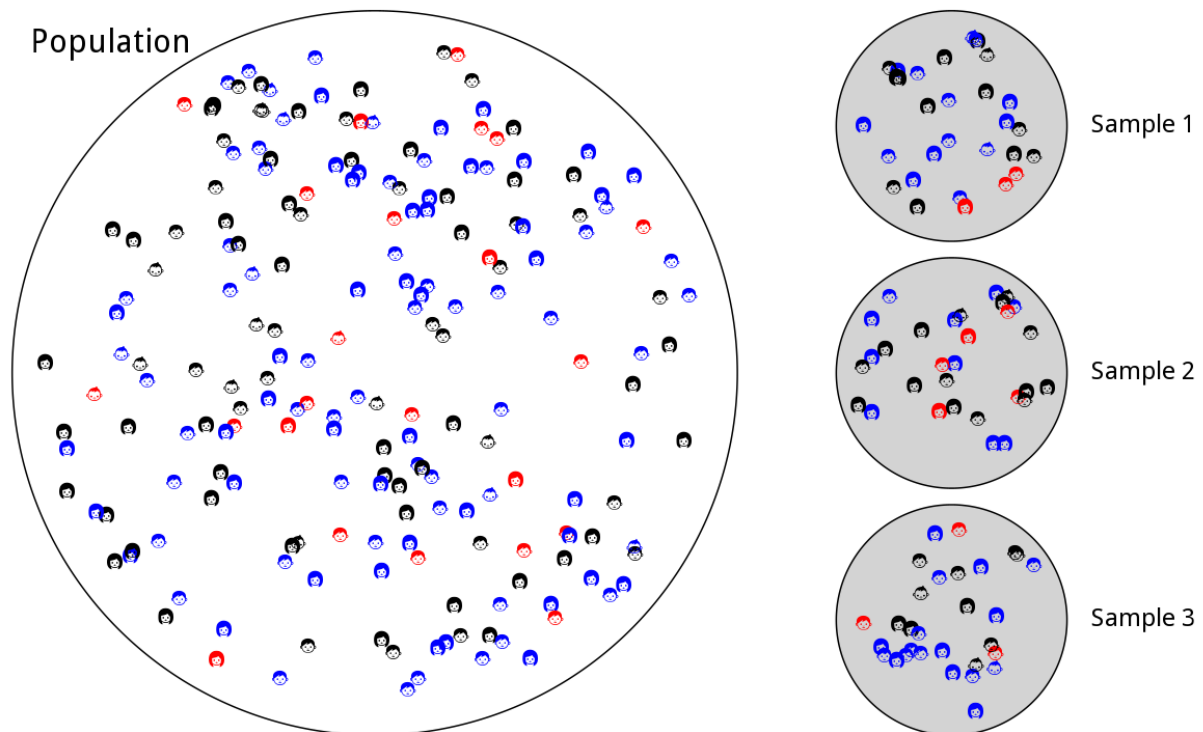
- Anders Hald (1913-2007), Danish statistician: *“a genius who almost single-handedly created the foundations for modern statistical science”*
- Bradley Efron (b. 1938): *“the single most important figure in 20th century statistics”*

## 15.2 Data collection

- Competences, ideally:
  - Statistics, both conceptually and analyses
  - Data wrangling (loading data; right format for analyses, tables, figures; ...)
  - Visualizations
  - Knowledge about subject (best with access to experts)
- Not just downloading a spreadsheet!
  - Population vs sample
  - Descriptives of the sample (e.g. mean)
  - Statistical inference about population (how close is sample's mean to population's mean)
- Do collect and analyze data, but know about pitfalls and limitations in generalisability!

## 16 Population and sample

### 16.1 Population and sample



Sample 3 of size  $n = 30$ :

shape	color	n_sample	p_sample	p_pop	p_diff
baby	black	2	0.07	0.04	-0.02
baby	blue	1	0.03	0.04	0.01
baby	red	0	0.00	0.01	0.01



shape	color	n_sample	p_sample	p_pop	p_diff
man	black	5	0.17	0.12	-0.04
man	blue	8	0.27	0.22	-0.04
man	red	3	0.10	0.08	-0.02
woman	black	3	0.10	0.23	0.13
woman	blue	8	0.27	0.22	-0.05
woman	red	0	0.00	0.02	0.02

- Descriptive vs statistical inference.

## 17 Example: United States presidential election, 1936

### 17.1 Example: United States presidential election, 1936

(Based on Agresti, this and this.)

- Current president: Franklin D. Roosevelt
- Election: Franklin D. Roosevelt vs Alfred Landon (Republican governor of Kansas)
- Literary Digest: magazine with history of accurately predicting winner of past 5 presidential elections

### 17.2 Example: United States presidential election, 1936

- Literary Digest poll ( $\hat{\pi}$  and  $1 - \hat{\pi}$ ): Landon: 57%; Roosevelt: 43%
- Actual results ( $\pi$  and  $1 - \pi$ ): Landon: 38%; Roosevelt: 62%
- Sampling error: 57%-38% = 19%
  - Practically all of the sampling error was the result of **sample bias**
  - Poll size of > 2 mio. individuals participated – extremely large poll

### 17.3 Example: United States presidential election, 1936

- Mailing list of about 10 mio. names was created
  - Based on every telephone directory, lists of magazine subscribers, rosters of clubs and associations, and other sources
  - Each one of 10 mio. received a mock ballot and asked to return the marked ballot to the magazine
- “respondents who returned their questionnaires represented only that subset of the population with a relatively intense interest in the subject at hand, and as such constitute in no sense a random sample ... it seems clear that the minority of anti-Roosevelt voters felt more strongly about the election than did the pro-Roosevelt majority” (*The American Statistician*, 1976)
- Biases:
  - Selection bias
    - \* List generated towards middle- and upper-class voters (e.g. 1936 and telephones)
    - \* Many unemployed (club memberships and magazine subscribers)
  - Non-response bias
    - \* Only responses from 2.3/2.4 mio out of 10 million people
    - \* Cannot force people to participate: but mail may be junk (phone, interviews, online, pay/paid, ...)

## 18 Example: Bullet holes of honor

### 18.1 Example: Bullet holes of honor

(Based on this.)

- World War II
- Royal Air Force (RAF), UK
  - Lost many planes to German anti-aircraft fire
- Armor up!
  - Where?
  - Count up all the bullet holes in planes that returned from missions
    - \* Put extra armor in the areas that attracted the most fire

### 18.2 Example: Bullet holes of honor

- Hungarian-born mathematician Abraham Wald:
  - If a plane makes it back safely with a bunch of bullet holes in its wings: holes in the wings aren't very dangerous
    - \* **Survivorship bias**
  - Armor up the areas that (on average) don't have any bullet holes
    - \* They never make it back, apparently dangerous

## 19 Theory: Biases / sampling

### 19.1 Biases

Agresti section 2.3:

- Sampling/selection bias
  - Probability sampling: each sample of size  $n$  has same probability of being sampled
    - \* Still problems: undercoverage, groups not represented (inmates, homeless, hospitalized, ...)
  - Non-probability sampling: probability of sample not possible to determine
    - \* E.g. volunteer sampling
- Response bias
  - E.g. poorly worded, confusing or even order of questions
  - Lying if think socially unacceptable
- Non-response bias
  - Non-response rate high; systematic in non-responses (age, health, believes)

## 19.2 Sampling

Agresti section 2.4:

- Random sampling schemes:
  - Simple sampling: each possible sample of equal size equally probable
  - Systematic sampling
  - Stratified sampling
  - Cluster sampling
  - Multistage sampling
  - ...

## 20 Data wrangling

### 20.1 Data wrangling

This will be illustrated with two specific cases.

The material is on Moodle.