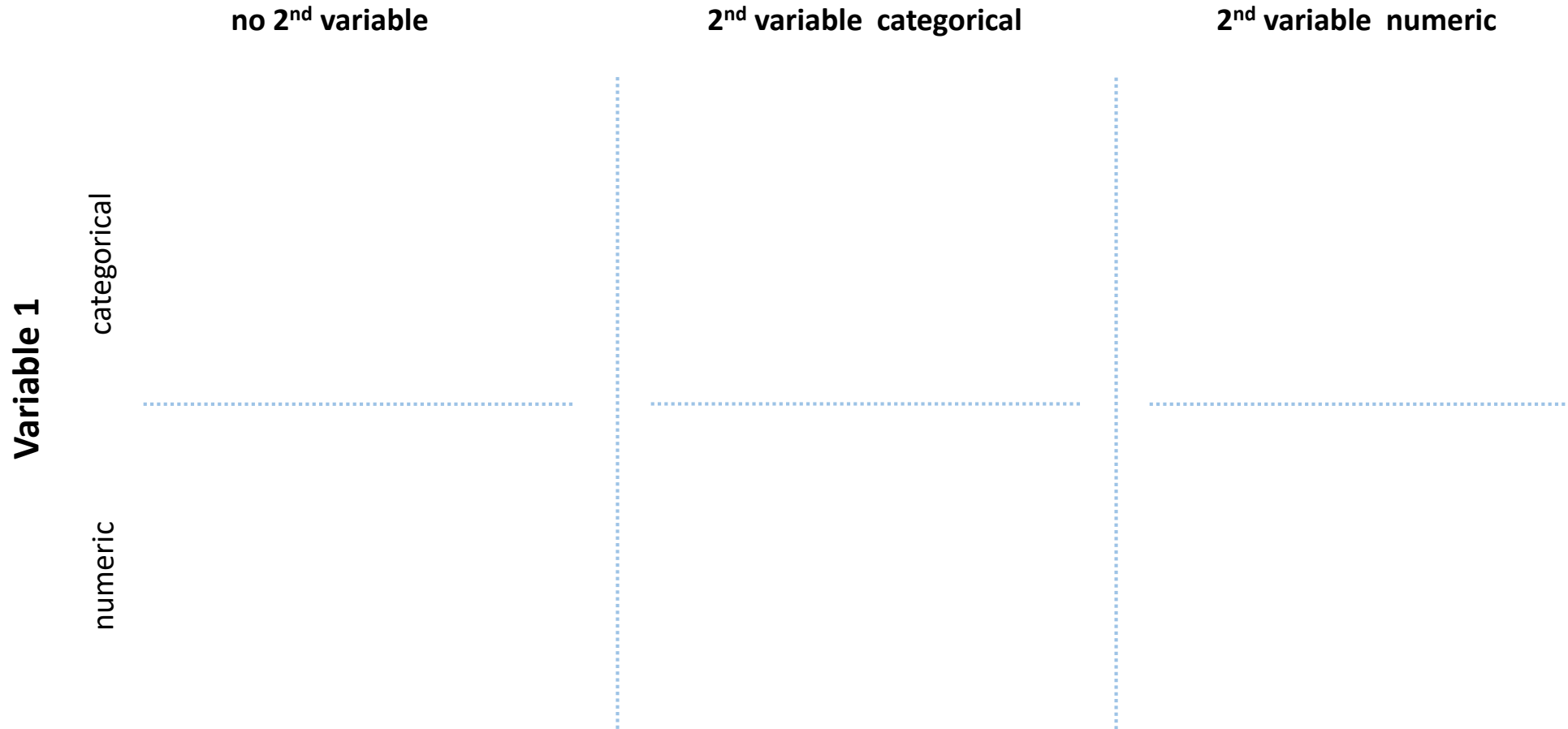
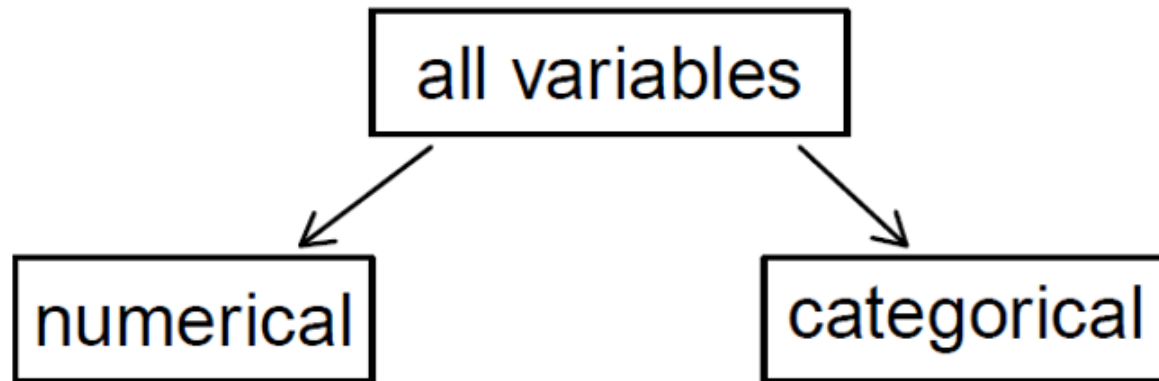


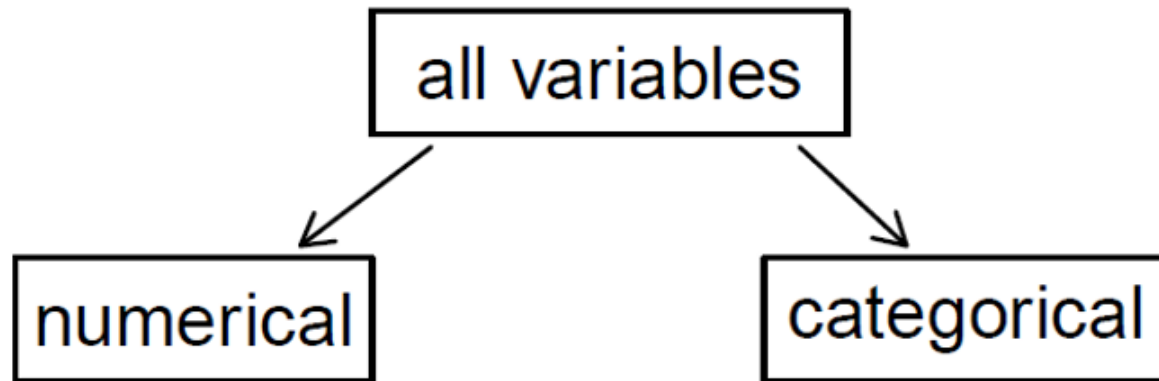
# Exploratory Data Analysis

## Graphical Tools

Prof. Dr. Gero Szepannek  
Statistics, Business Mathematics & Machine Learning  
Stralsund University of Applied Sciences





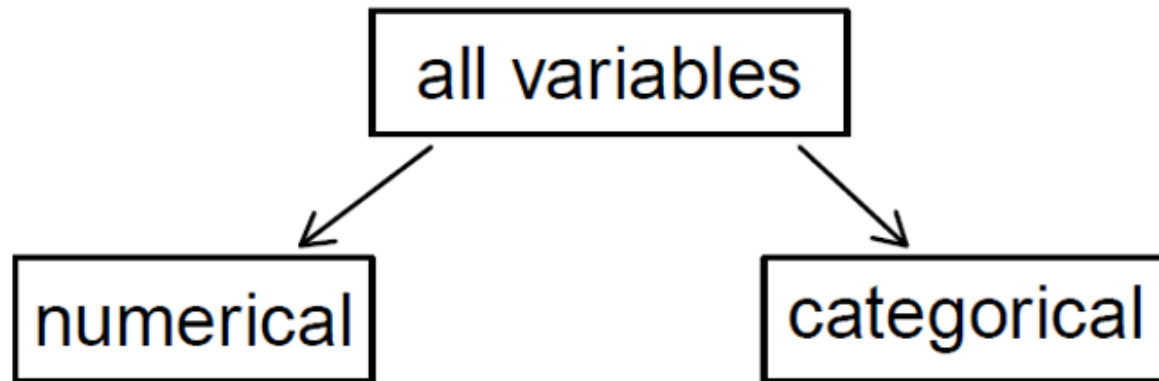
**Examples:**

- Size
- Number of eyes

**Example:**

- Colour



**Examples:**

- Size
- Number of eyes

**Example:**

- Colour



**Different analysis methodology!**



Variable 1

categorical

no 2<sup>nd</sup> variable!

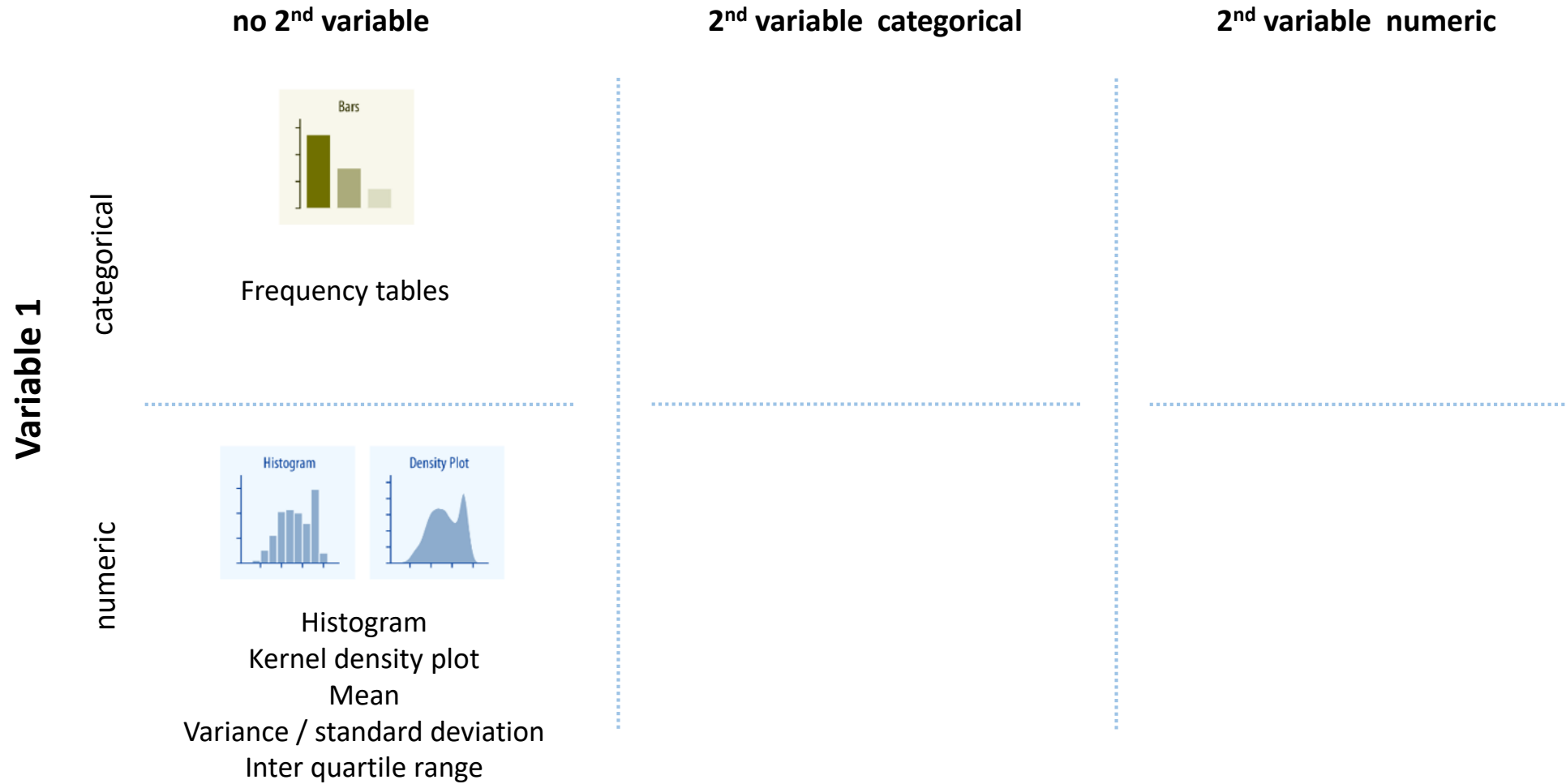


numeric



2<sup>nd</sup> variable categorical

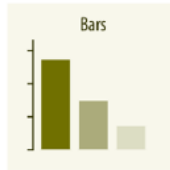
2<sup>nd</sup> variable numeric



Variable 1

categorical

no 2<sup>nd</sup> variable

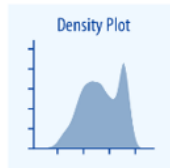
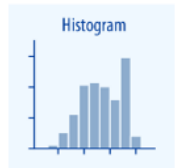


Frequency tables

2<sup>nd</sup> variable categorical

2<sup>nd</sup> variable numeric

numeric



Histogram

Kernel density plot

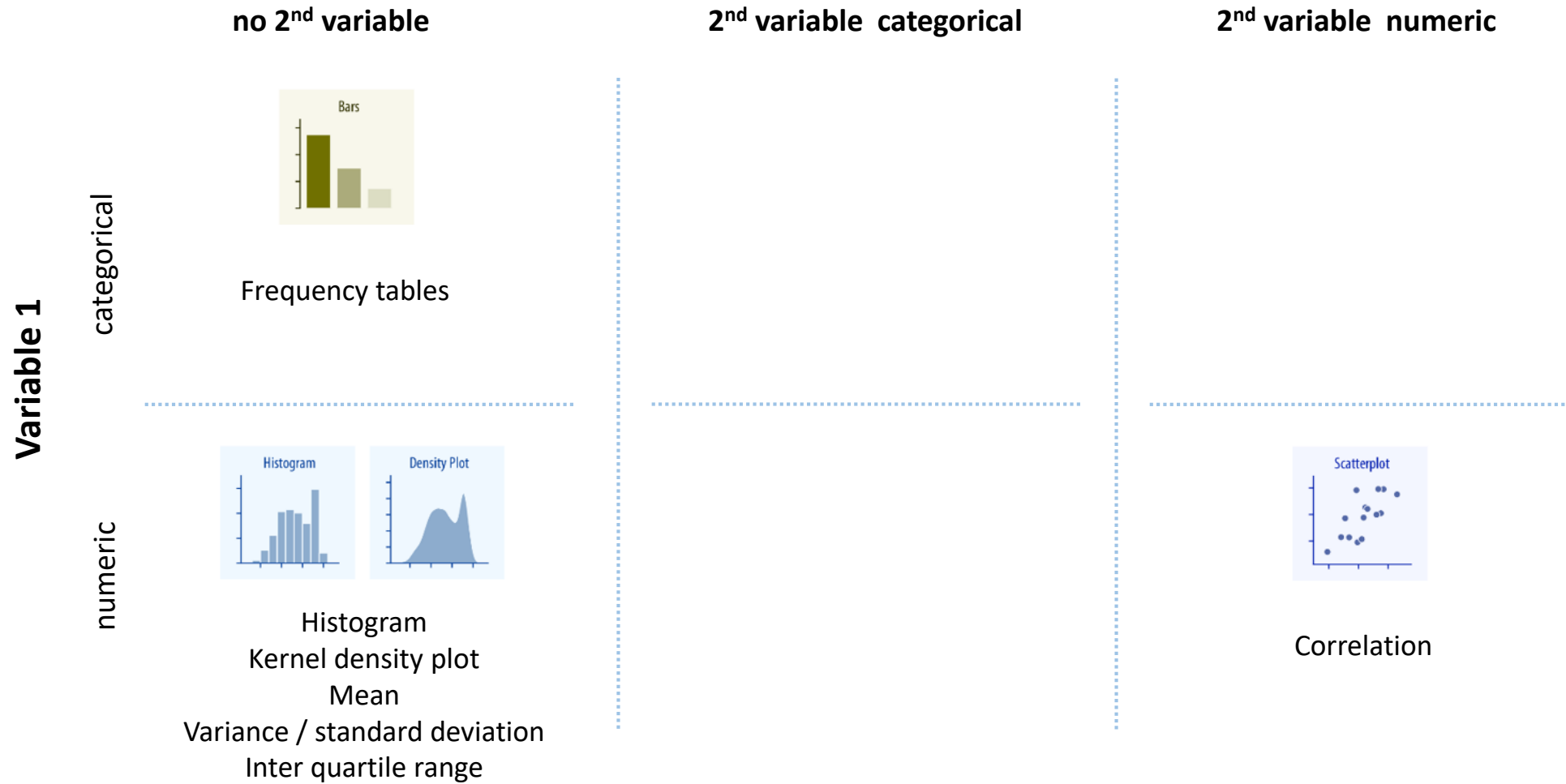
Mean

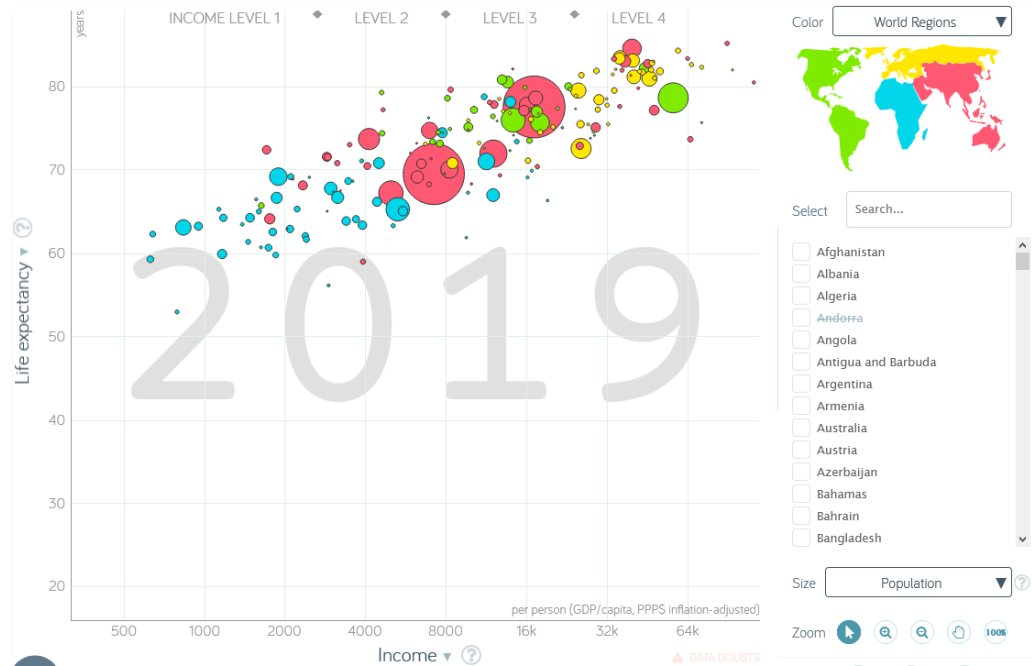
Variance / standard deviation

Inter quartile range





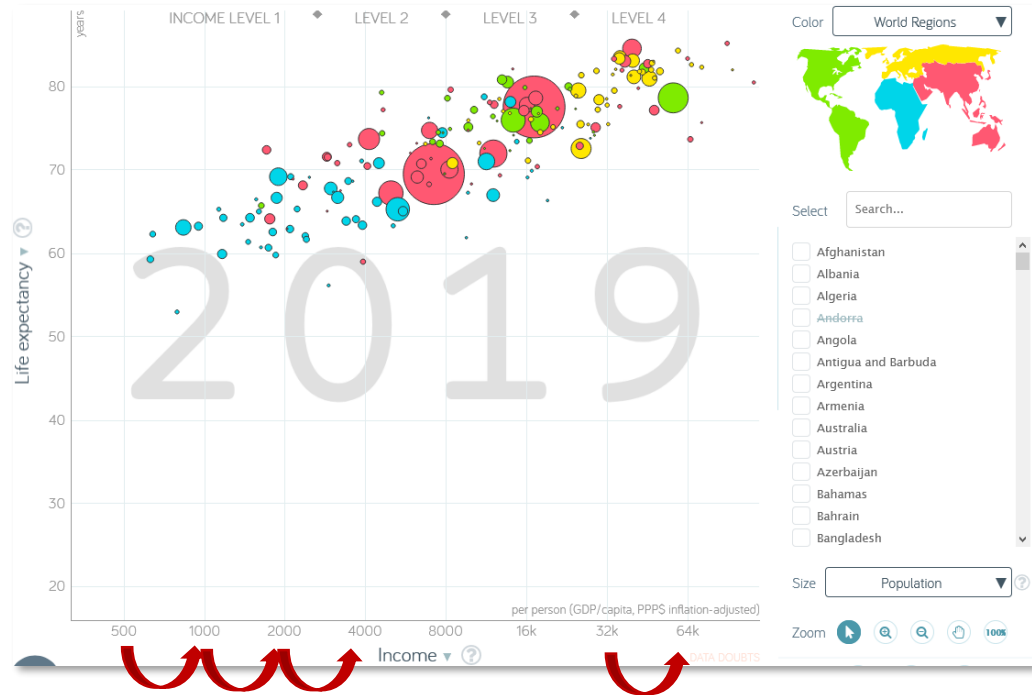




# GAPMINDER

<https://www.gapminder.org/>





x2

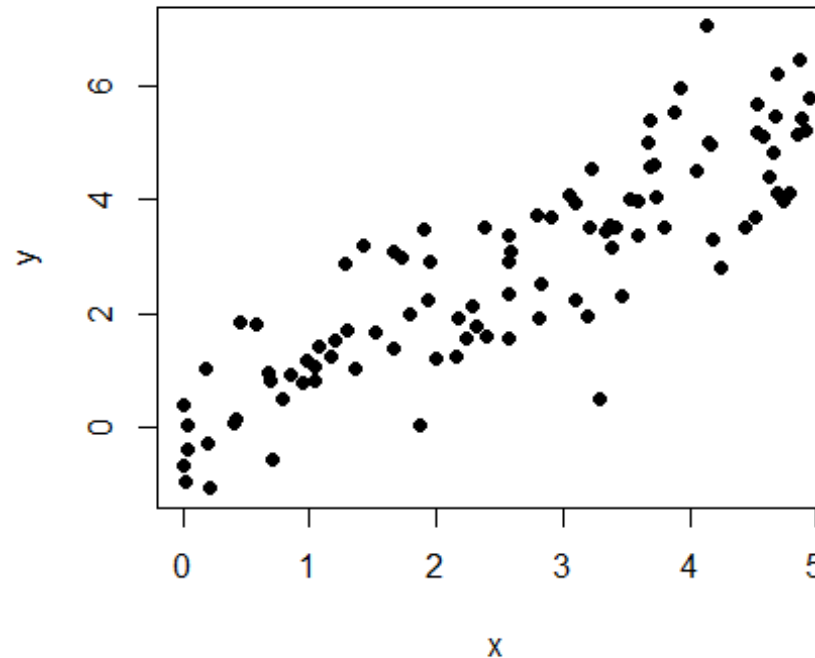
- Note: x-axis logarithmically scaled...
- This is often done if the differences are huge

A statistical measure to quantify the dependency between two numeric variables is given by the **coefficient of correlation  $\rho$**  :

**Interpretation:**

$-1 \leq \rho \leq 1$  where:

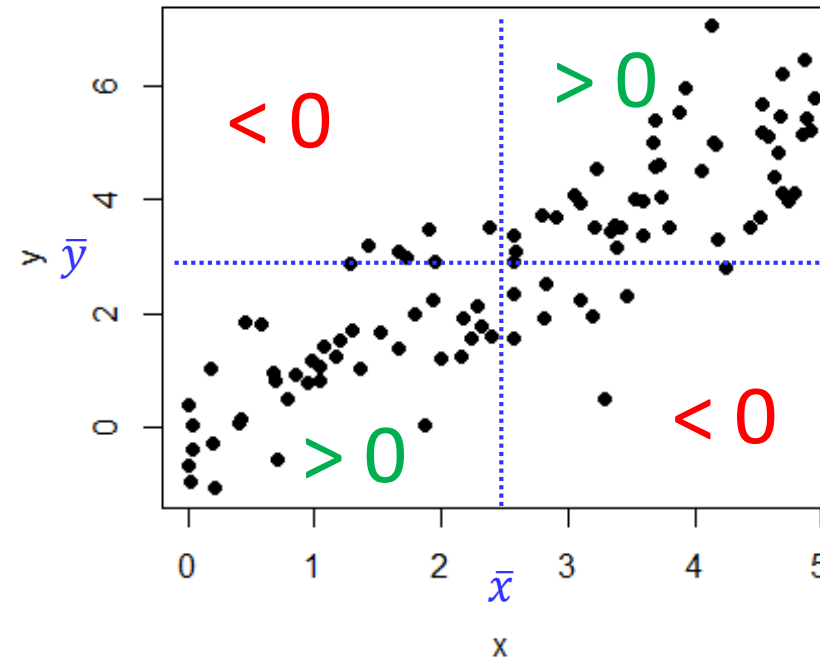
Correlation	Interpretation
$\rho > 0$	Positive dependency
$\rho = 0$	No (linear) dependency
$\rho < 0$	Negative dependency



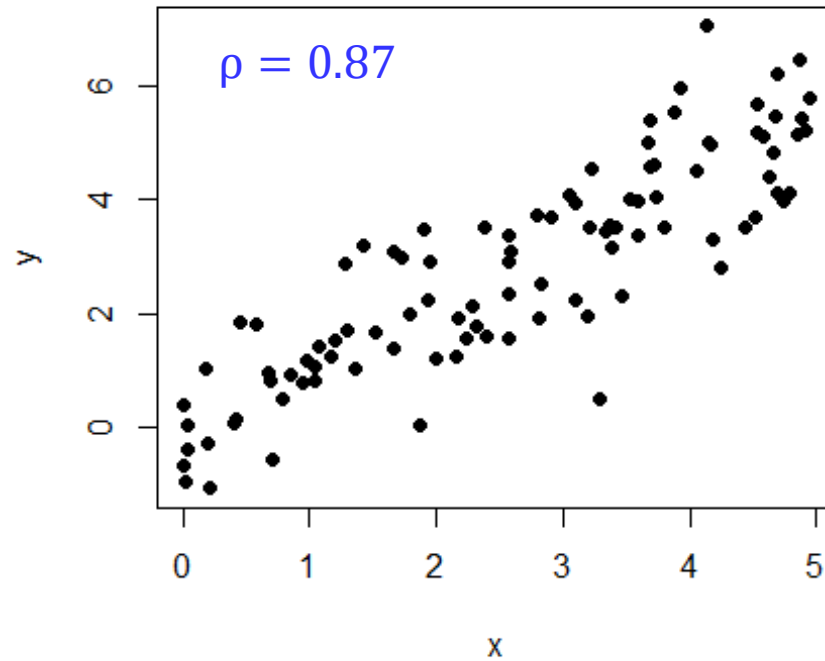
...comparison with the mean of x and y

$$\rho = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{j=1}^n (x_j - \bar{x})^2 \cdot \sum_{j=1}^n (y_j - \bar{y})^2}}$$

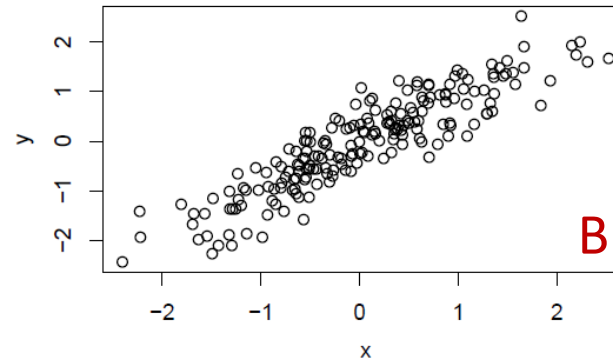
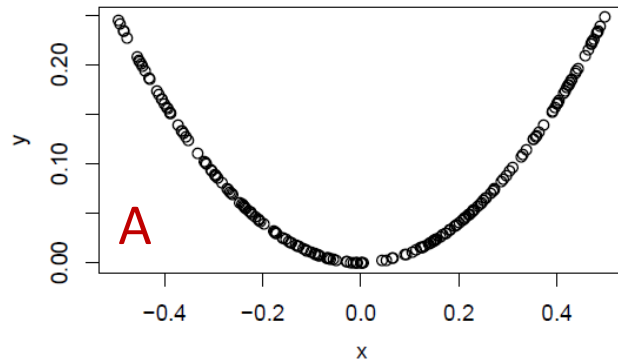
...only for scaling  $-1 \leq \rho \leq 1$



Correlation	Interpretation
$\rho > 0$	Positive dependency
$\rho = 0$	No (linear) dependency
$\rho < 0$	Negative dependency

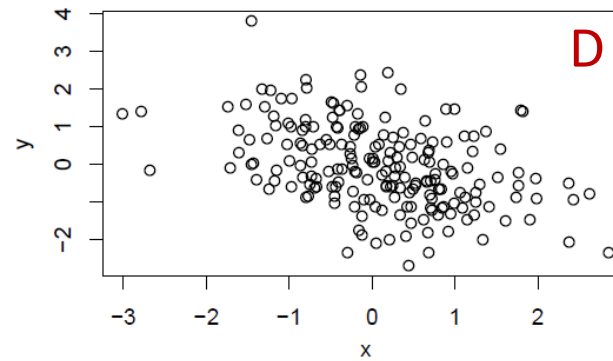
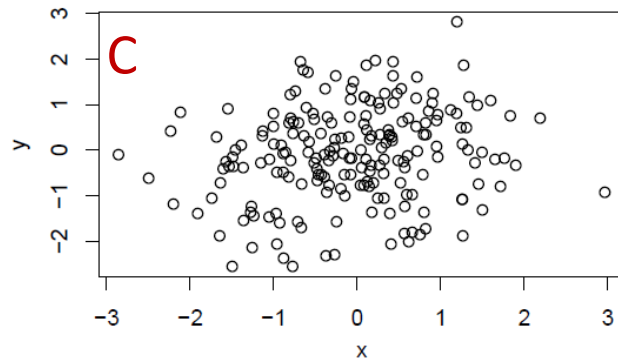


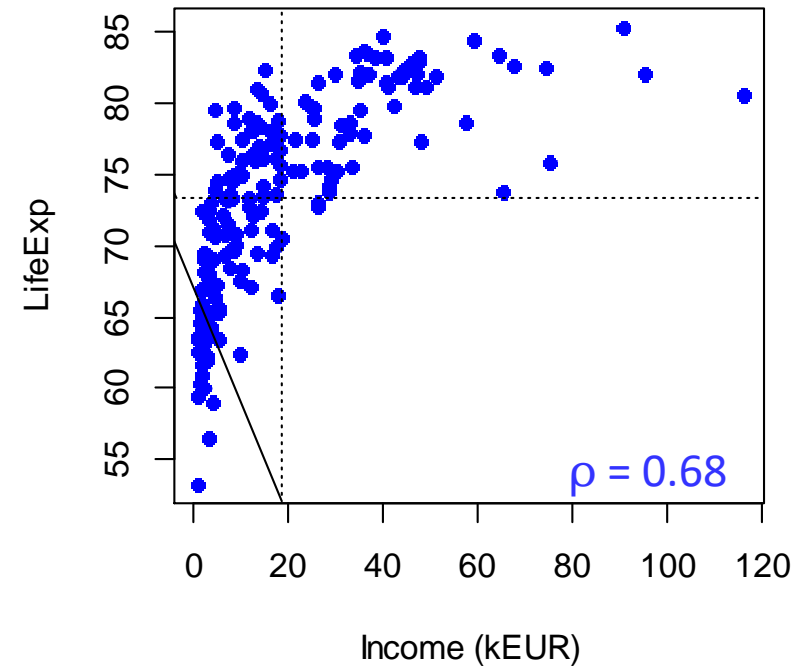
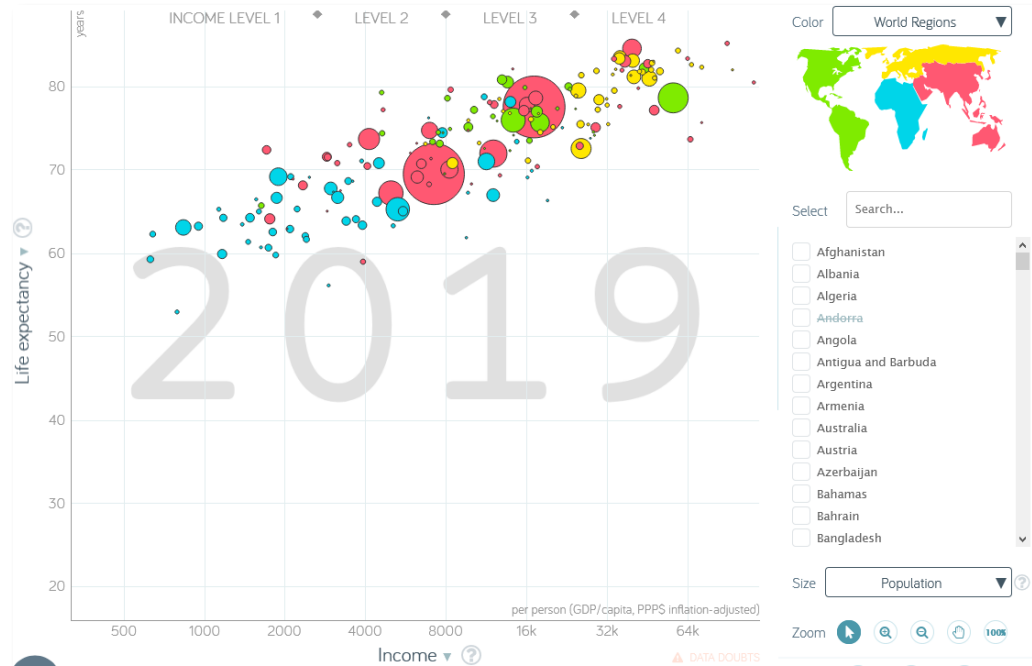
But note: No correlation (i.e.  $\rho = 0$ )  $\nRightarrow$  independence (!) (just: no linear dependency).



Which is the corresponding plot to a correlation of:

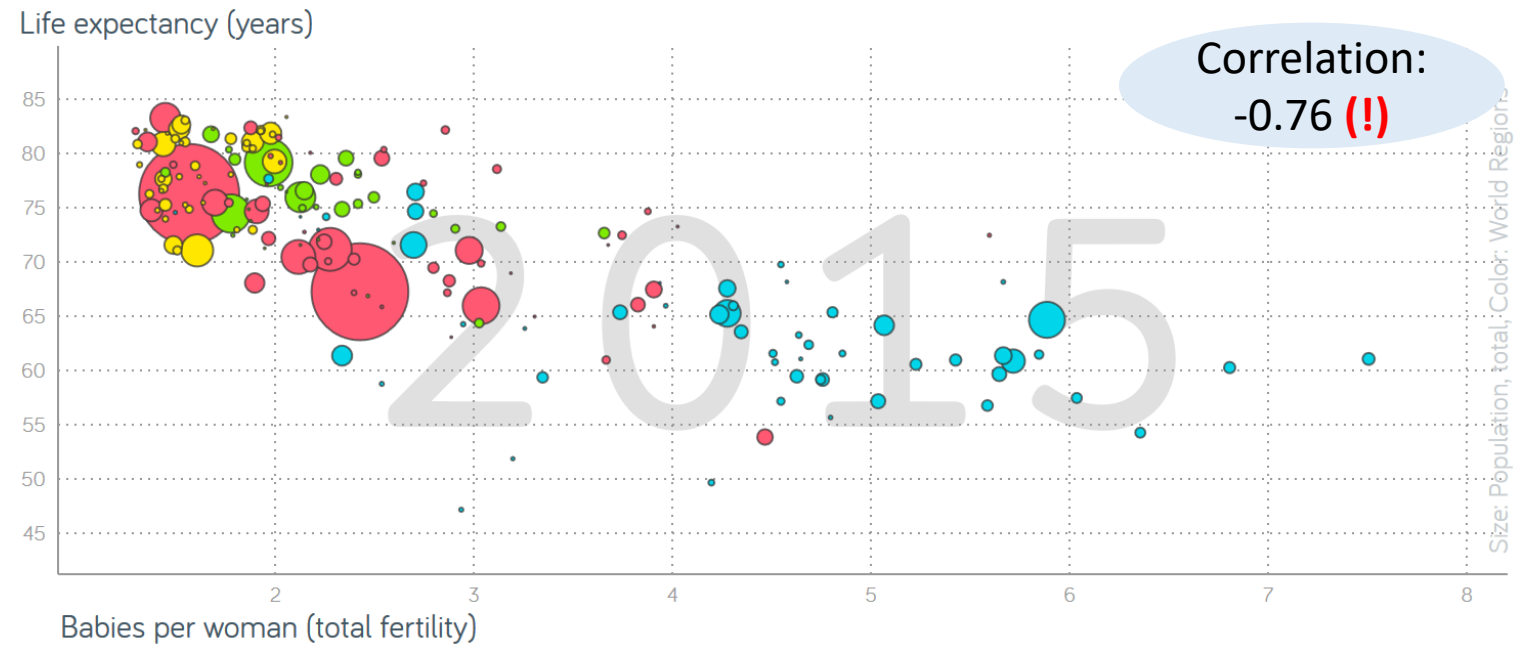
- ☐ 0.9
- ☐ -0.42

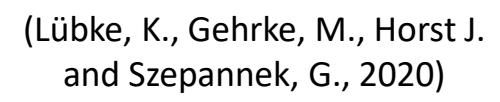
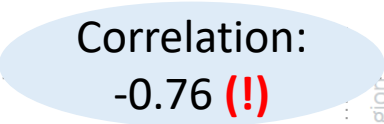




$\rho = 0.68 \rightarrow$  Strong positive dependency  
btw income and life expectancy!



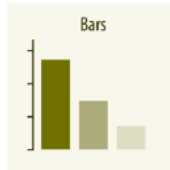




Variable 1

categorical

no 2<sup>nd</sup> variable



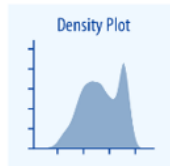
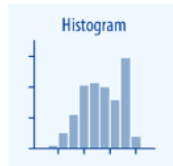
Frequency tables

2<sup>nd</sup> variable categorical



2<sup>nd</sup> variable numeric

numeric



Histogram

Kernel density plot

mean

Variance / standard deviation

Inter quartile range

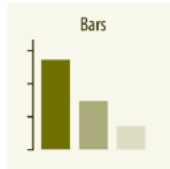


Correlation

Variable 1

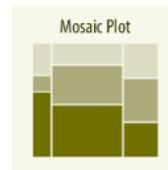
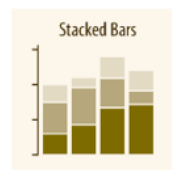
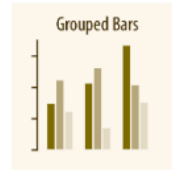
categorical

no 2<sup>nd</sup> variable



Frequency tables

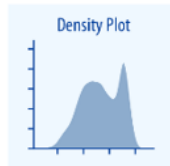
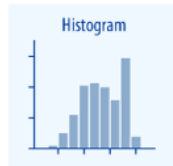
2<sup>nd</sup> variable categorical



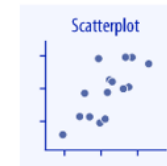
Contingency tables  
 $\chi^2$  Test  
Cramer's V

2<sup>nd</sup> variable numeric

numeric



Histogram  
Kernel density plot  
mean  
Variance / standard deviation  
Inter quartile range



Correlation



Mosaicplot

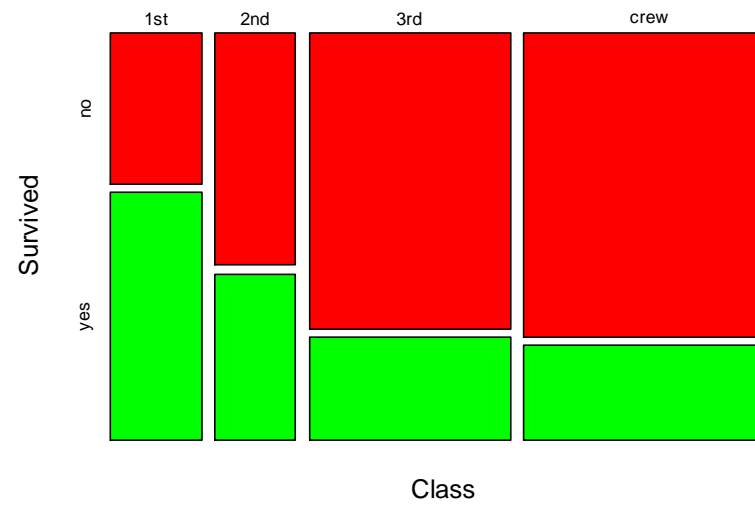
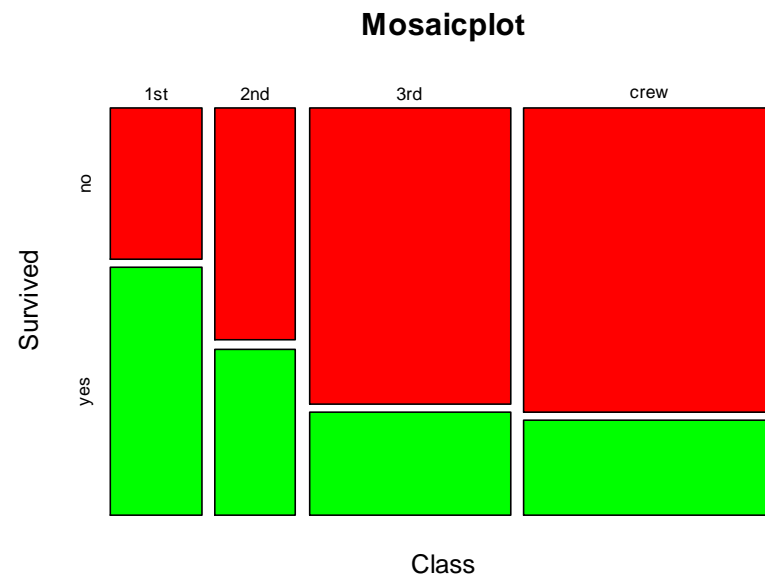


Figure taken from: <https://www.geo.de/geolino/mensch/10493-rtkl-geschichte-die-letzte-nacht-auf-der-titanic>

$P(\text{survived}|\text{class})$  $P(\text{class})$



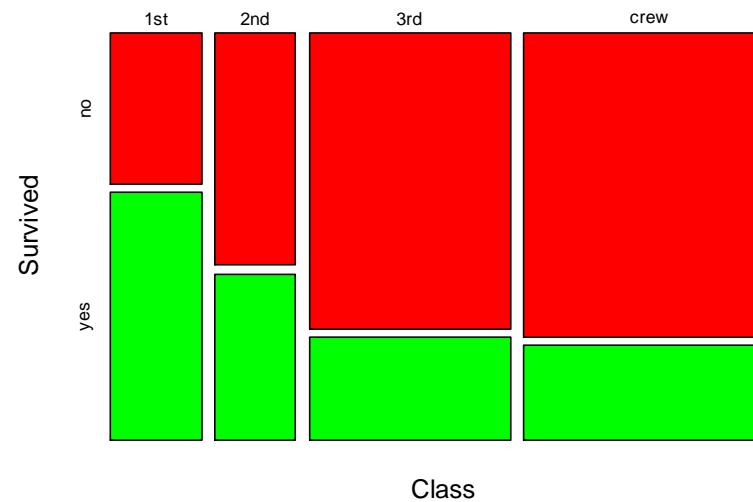
## Absolute frequencies

	1 <sup>st</sup>	2 <sup>nd</sup>	3 <sup>rd</sup>	crew	total
no	123	166	528	679	1496
yes	201	118	181	211	711
total	324	284	709	890	2207

## Conditional frequencies

	1 <sup>st</sup>	2 <sup>nd</sup>	3 <sup>rd</sup>	crew	total
no	0,37963	0,58451	0,74471	0,76292	1496
yes	0,62037	0,41549	0,25529	0,23708	711
total	324	284	709	890	2207

Mosaicplot



observed

			$\Sigma$
Frau	50	20	
Mann	10	20	
$\Sigma$			





observed

			$\Sigma$
Frau	50	20	
Mann	10	20	
$\Sigma$			

expected

			$\Sigma$
Frau			
Mann			
$\Sigma$			

**What counts could we expect if gender and preference were independent?**



observed

			$\Sigma$
Frau	50	20	
Mann	10	20	
$\Sigma$			

-

expected

			$\Sigma$
Frau			
Mann			
$\Sigma$			

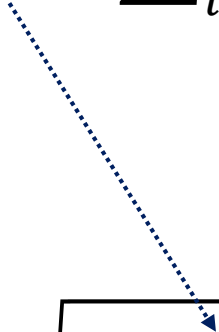
)<sup>2</sup>



expected

			$\Sigma$
Frau			
Mann			
$\Sigma$			

$$\chi^2 = \sum_{i,j} \frac{(o_{ij} - e_{ij})^2}{e_{ij}}$$

$$V = \sqrt{\frac{\chi^2/n}{\min(c-1, r-1)}}$$


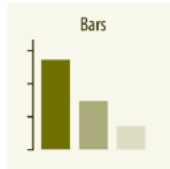
n: #observations  
c/r: #columns/rows of the table

**$0 \leq V \leq 1$  measures the dependency between two categorical variables.**

Variable 1

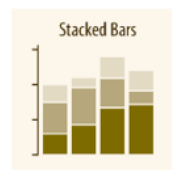
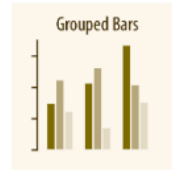
categorical

no 2<sup>nd</sup> variable



Frequency tables

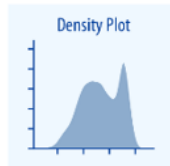
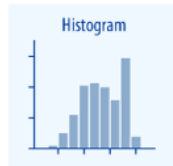
2<sup>nd</sup> variable categorical



Contingency tables  
 $\chi^2$  Test  
Cramer's V

2<sup>nd</sup> variable numeric

numeric



Histogram  
Kernel density plot  
mean  
Variance / standard deviation  
Inter quartile range



Correlation

Variable 1

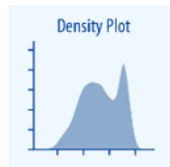
categorical

numeric

no 2<sup>nd</sup> variable



Frequency tables



Histogram

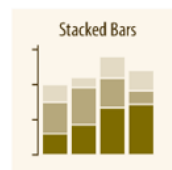
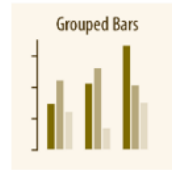
Kernel density plot

mean

Variance / standard deviation

Inter quartile range

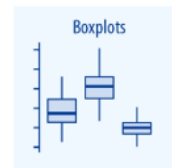
2<sup>nd</sup> variable categorical



Contingency tables

$\chi^2$  Test

Cramer's V

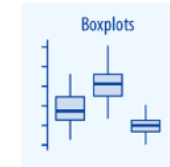


Boxplots

Mean | group

Variance | group

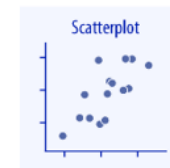
2<sup>nd</sup> variable numeric



Boxplots

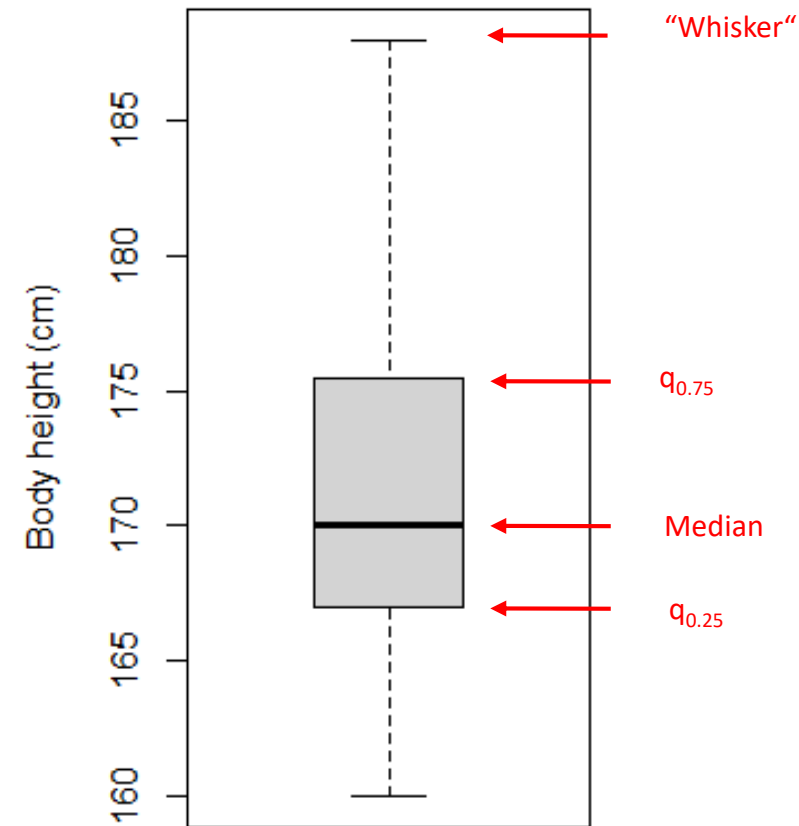
Mean | group

Variance | group



Correlation

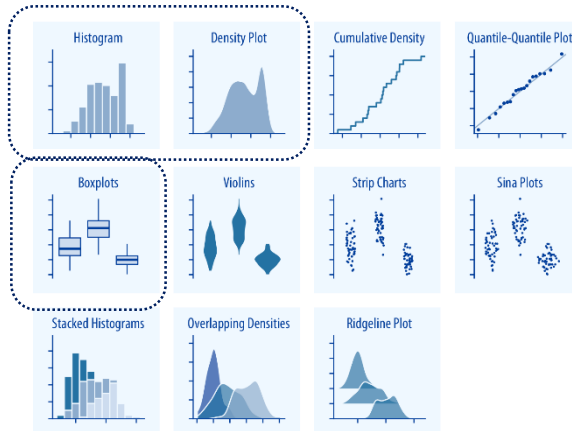
- Visualization of the distribution based on five key measures: min, max,  $q_{0.75}$ ,  $q_{0.25}$  & median.
- The box covers 50% of the data.
- ‘Whiskers’:  $q_{0.75}/q_{0.25} \pm \frac{3}{2} IQR$  or the maximum/minimum.
- ...Points outside the whiskers are suspicious to be outliers.



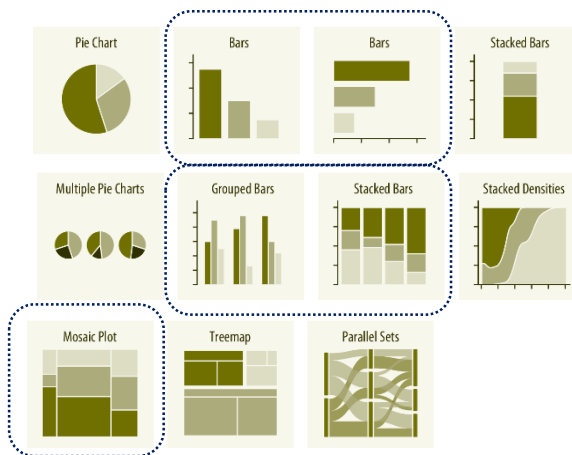
## Amounts



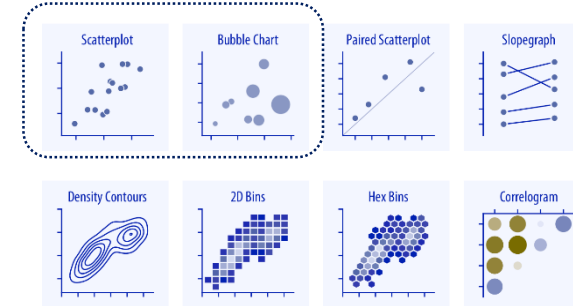
## Distributions



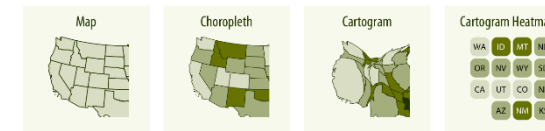
## Proportions



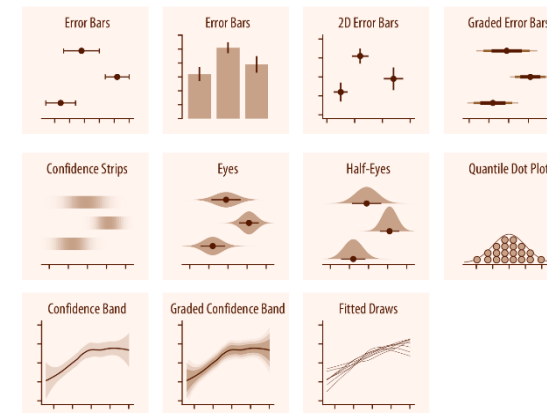
## Dependency

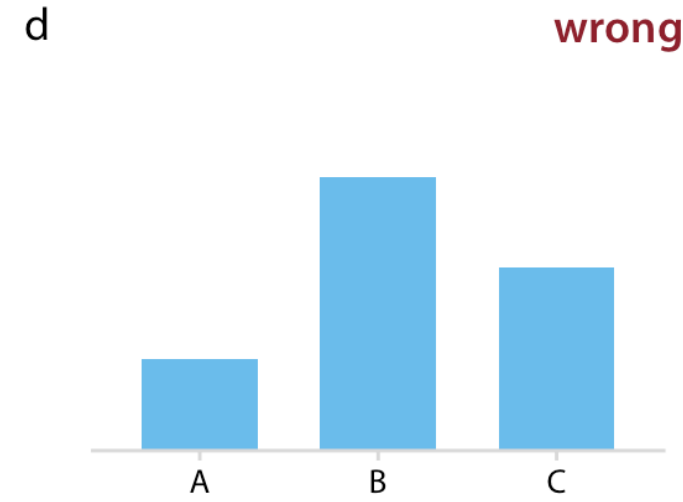
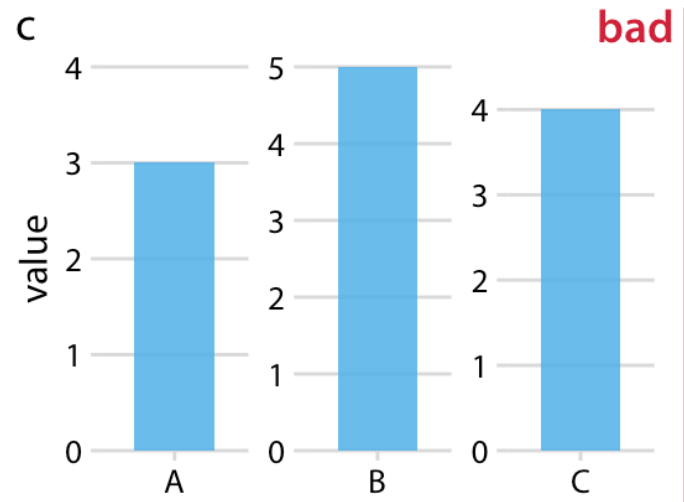
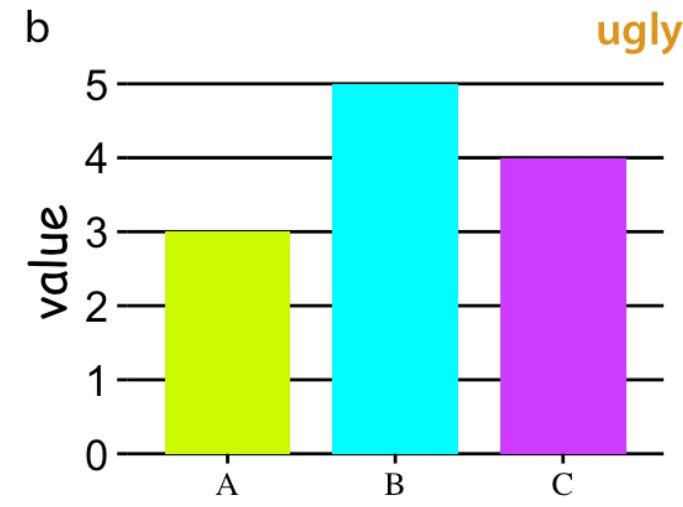
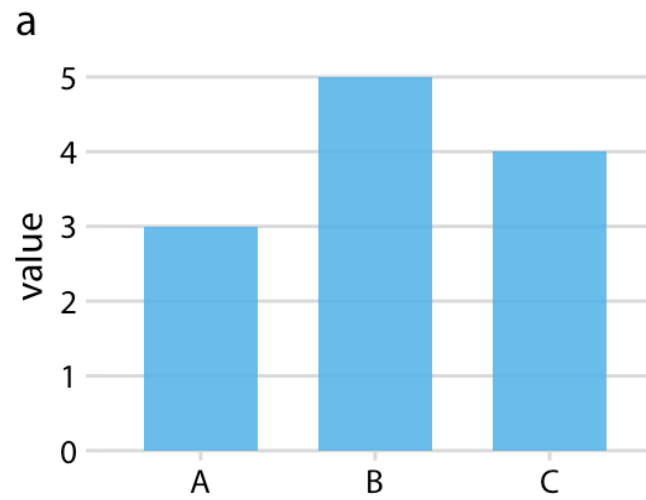


## Geodata



## Uncertainty

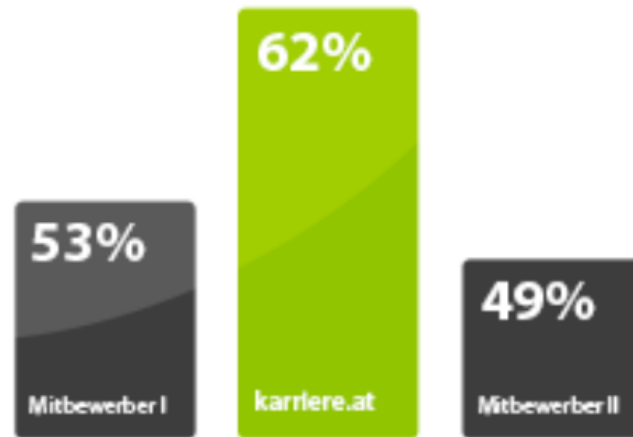






## Höchste Bekanntheit

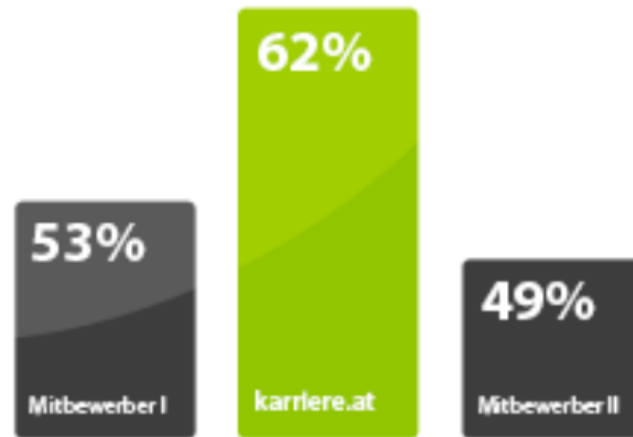
Fast 2/3 der Arbeitnehmer kennen **karriere.at**.  
Im Mitbewerbsvergleich ist das spitze.



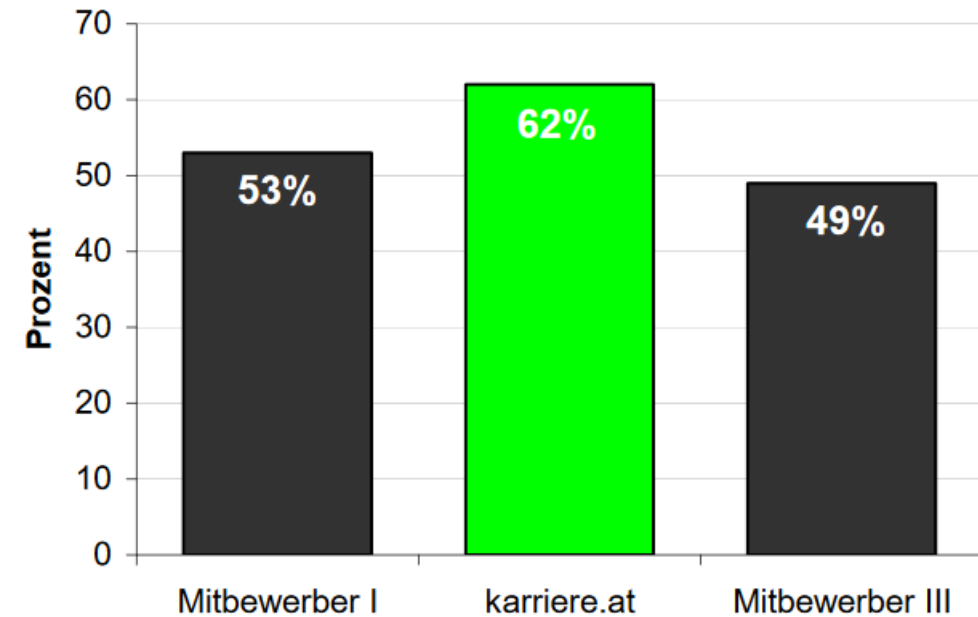
(gefunden am 12. November 2014 auf <http://www.karriere.at/hr>)

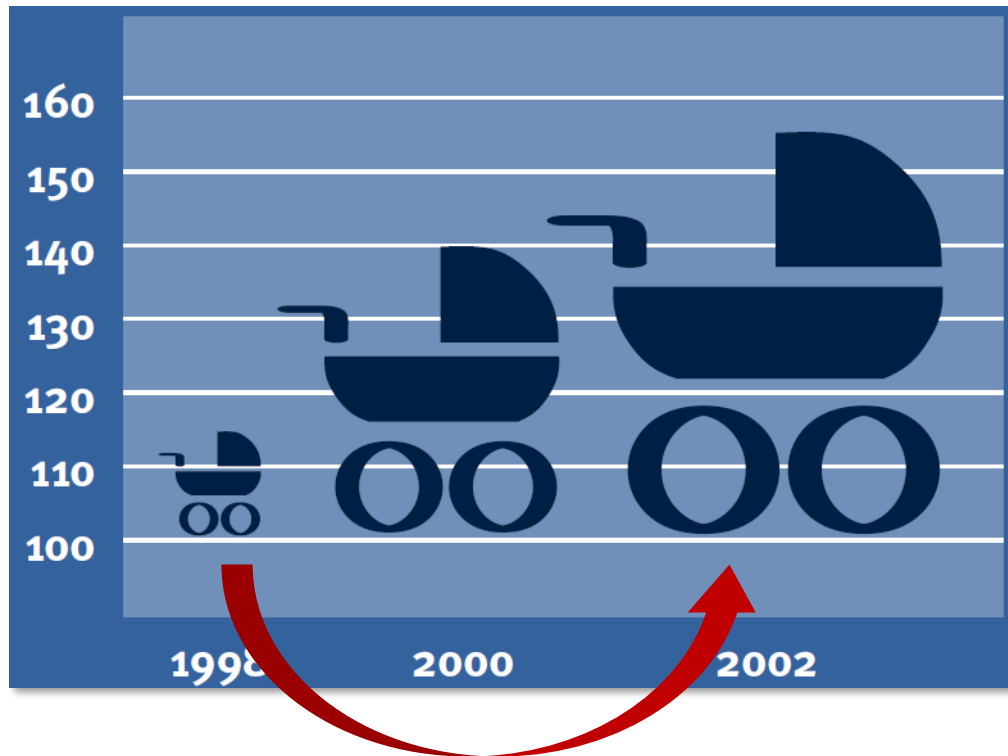
## Höchste Bekanntheit

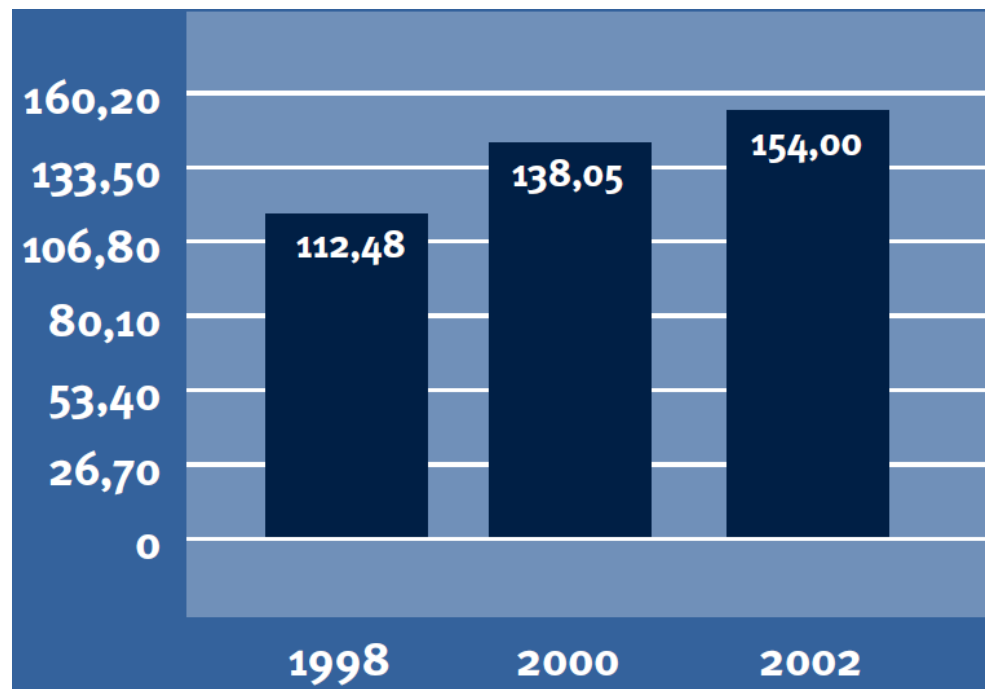
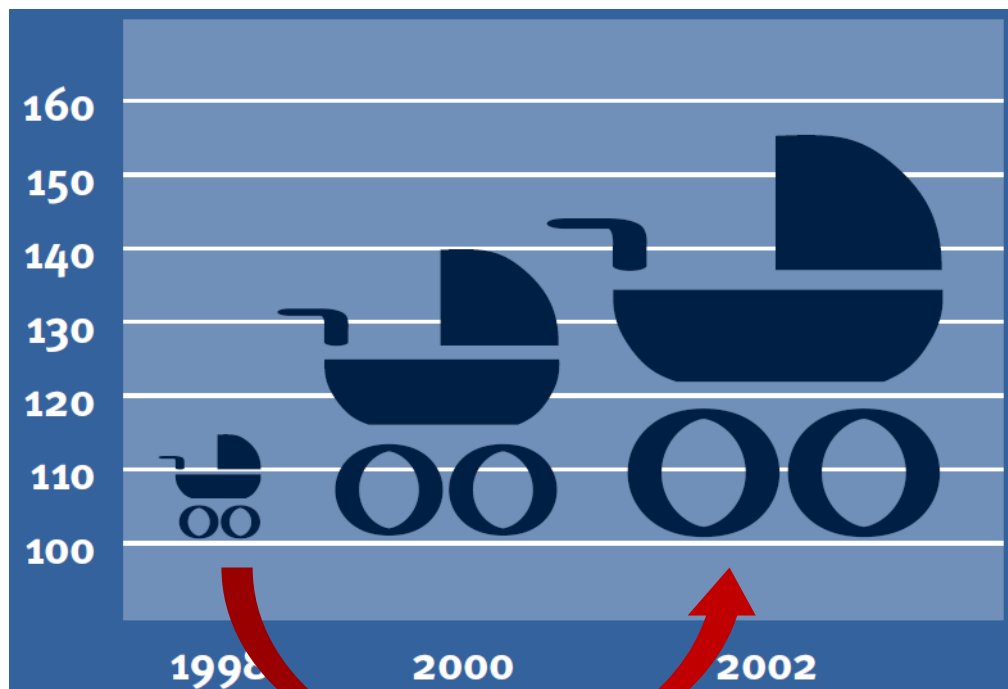
Fast 2/3 der Arbeitnehmer kennen **karriere.at**.  
Im Mitbewerbsvergleich ist das spitze.



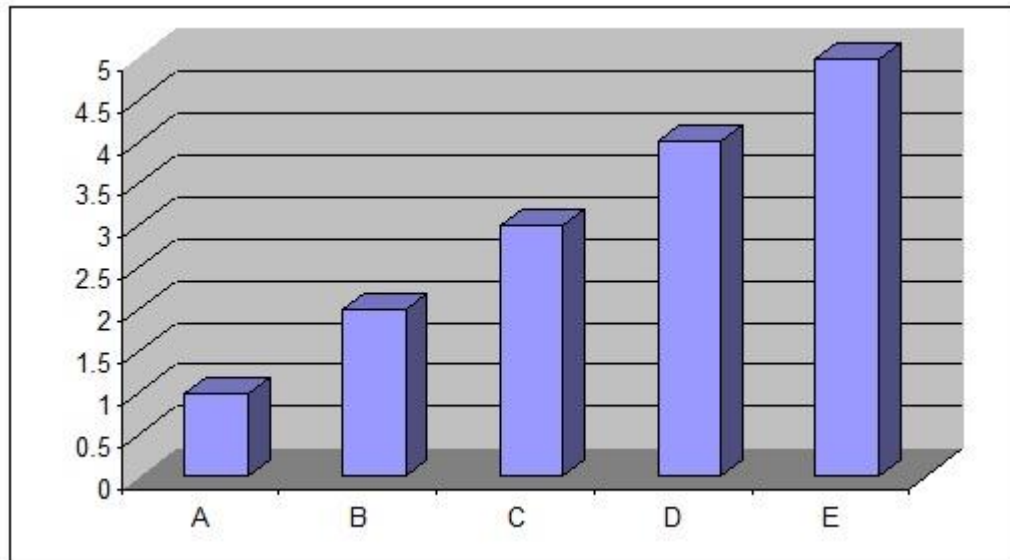
(gefunden am 12. November 2014 auf <http://www.karriere.at/hr>)



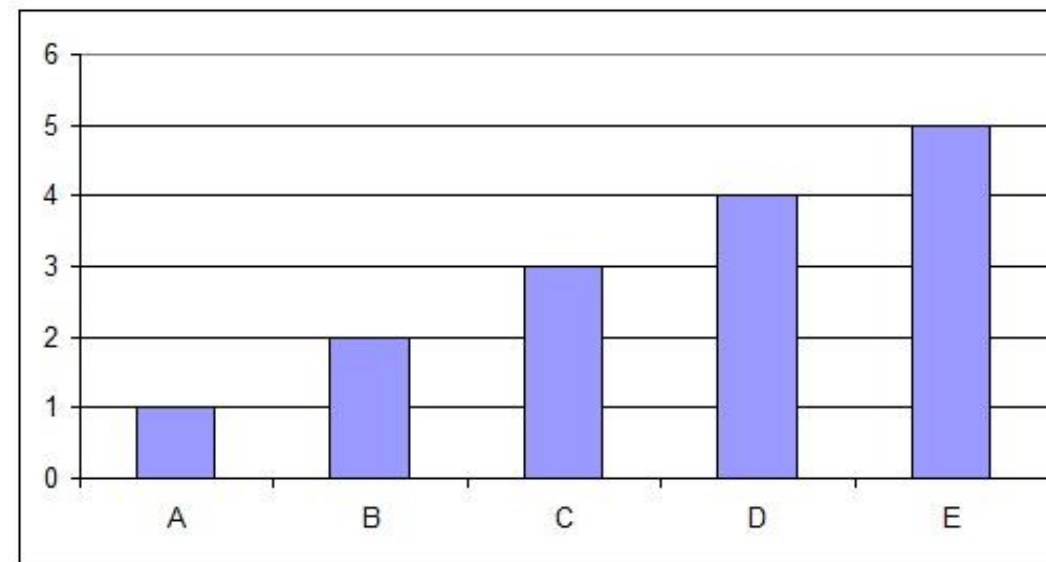
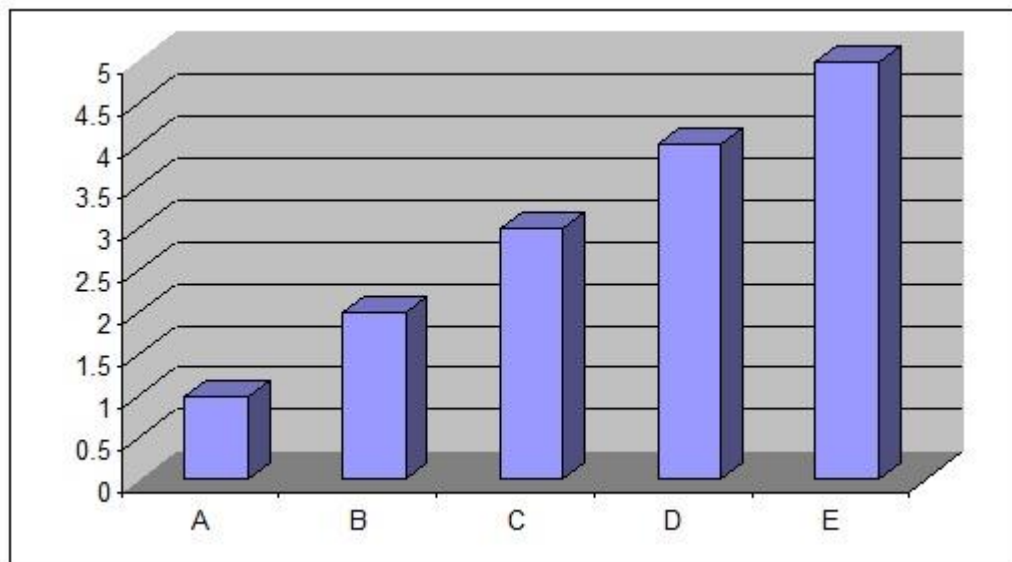




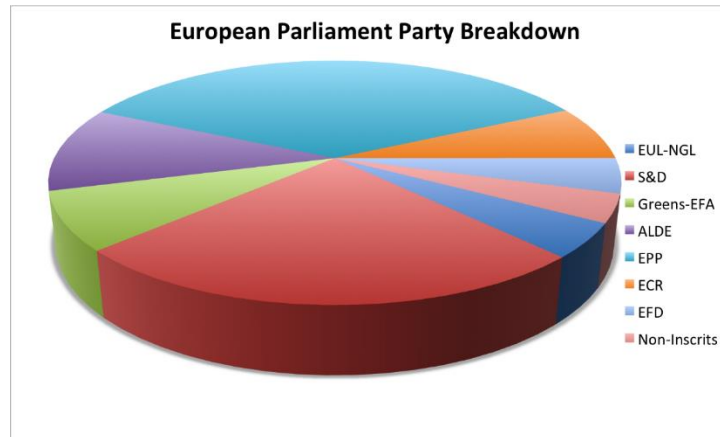
Source: [http://www.wdr.de/tv/applications/fernsehen/wissen/quarks/pdf/Q\\_Zahlen.pdf](http://www.wdr.de/tv/applications/fernsehen/wissen/quarks/pdf/Q_Zahlen.pdf)



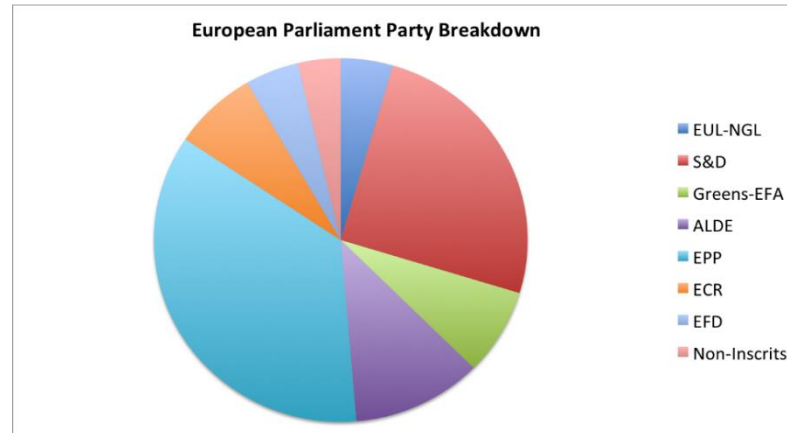
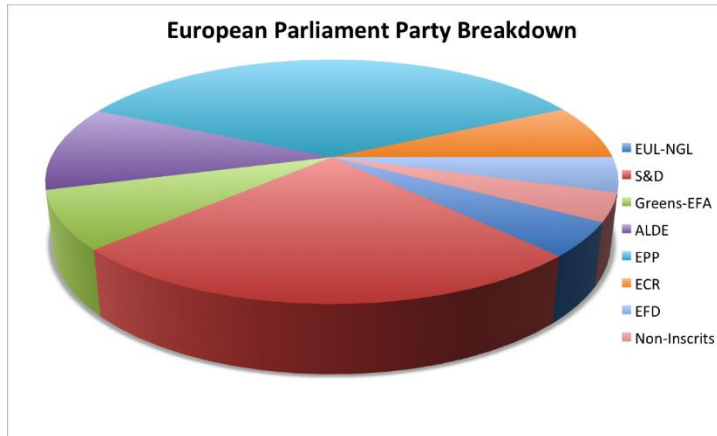
Source: <http://consultantjournal.com/blog/use-3d-charts-at-your-own-risk>



Source: <http://consultantjournal.com/blog/use-3d-charts-at-your-own-risk>

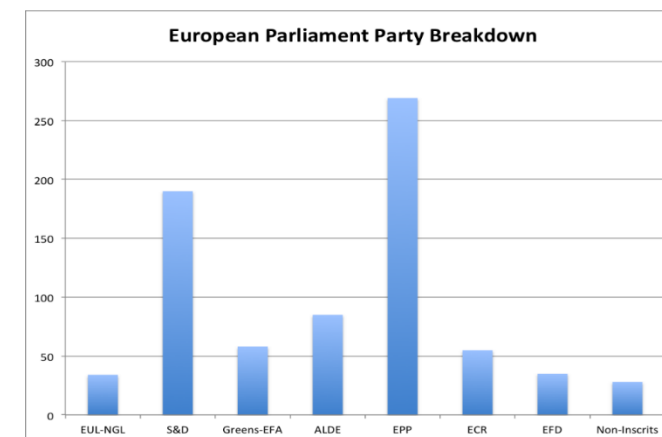
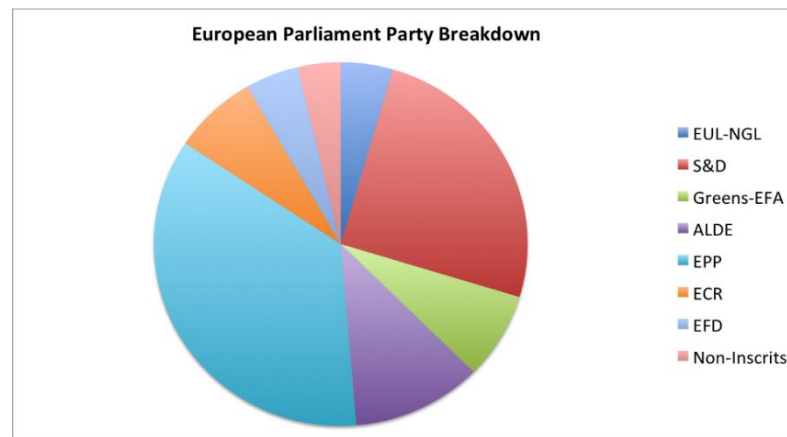
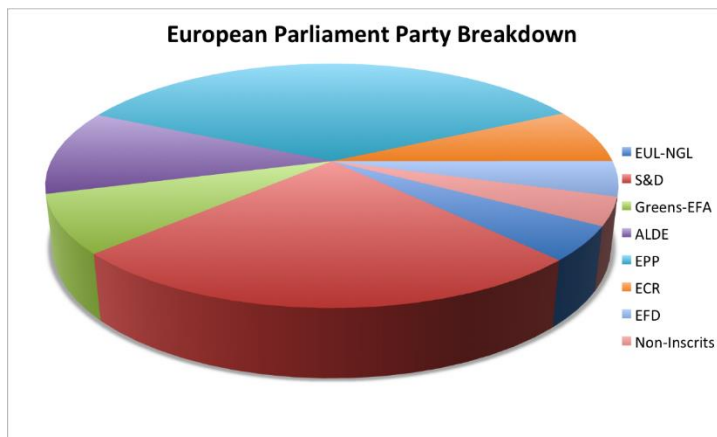


Source: <http://www.businessinsider.com/pie-charts-are-the-worst-2013-6?IR=T>



Source: <http://www.businessinsider.com/pie-charts-are-the-worst-2013-6?IR=T>





Source: <http://www.businessinsider.com/pie-charts-are-the-worst-2013-6?IR=T>



Source: A. Quatember: *Statistischer Unsinn – Wenn Medien an der Prozenzhürde scheitern*, Springer Spektrum.

