

# Beschreibende Statistik

<p><b>Absolute Häufigkeiten</b></p> $H = \sum_1^n h_i$	<p><b>Relative Häufigkeiten</b></p> $F = \sum_1^m f_i, \quad F(x) = \frac{H(x)}{n}$	<p><b>Begriffe</b></p> <ul style="list-style-type: none"> <li>• <math>\Omega</math> = Grundgesamtheit</li> <li>• <math>n</math> = Anzahl Objekte</li> <li>• <math>X</math> = Stichprobenwerte</li> <li>• <math>a</math> = Ausprägungen</li> <li>• <math>h</math> = Absolute Häufigkeit</li> <li>• <math>f</math> = Relative Häufigkeit</li> <li>• <math>H</math> = Kumulative Absolute Häufigkeit</li> <li>• <math>F</math> = Kumulative Relative Häufigkeit</li> </ul> <div data-bbox="1469 264 2119 528"> </div>	
<p><b>Kennwerte (Lagemasse)</b></p> <ul style="list-style-type: none"> <li>• Quantil <math>i = [n \cdot q], Q = x_i = x_{[n \cdot q]}</math></li> <li>• Interquartilsabstand <math>IQR = Q_3 - Q_1</math></li> <li>• Modus <math>x_{mod}</math> = Häufigste Wert</li> </ul>			
<p><b>Arithmetisches Mittel</b></p> $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \sum_{i=1}^m a_i \cdot f_i$	<p><b>Median</b></p> $\begin{cases} x_{[\frac{n+1}{2}]} & n \text{ ungerade} \\ 0.5 \cdot (x_{[\frac{n}{2}]} + x_{[\frac{n}{2}+1]}) & n \text{ gerade} \end{cases}$	<p><b>Boxplot</b></p> <ul style="list-style-type: none"> <li>• <math>Q_1, Q_2 = x_{med}, Q_3</math></li> <li>• <math>IQR = Q_3 - Q_1</math></li> <li>• Untere Antenne <math>x_u</math>: <math>u = \min[Q_1 - 1.5 \cdot IQR, Q_1]</math></li> <li>• Obere Antenne <math>x_o</math>: <math>o = \max[Q_3 + 1.5 \cdot IQR, Q_3]</math></li> <li>• Ausreisser: <math>x_i &lt; x_u \vee x_i &gt; x_o</math></li> </ul> <div data-bbox="1832 679 2119 954"> </div>	
<p><b>Stichprobenvarianz <math>s^2</math> (Streuemasse)</b></p> $s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \overline{x^2} - \bar{x}^2, \quad (s_{kor})^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ $(s_{kor})^2 = \frac{n}{n-1} \cdot s^2$			
<p><b>Standardabweichung <math>s</math> (Streuemasse)</b></p> $s = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} = \sqrt{\overline{x^2} - \bar{x}^2}, \quad s_{kor} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$		<p><b>PMF + CDF</b></p> <div data-bbox="1010 1078 1559 1422"> </div>	<p><b>PDF + CDF</b></p> <div data-bbox="1597 1078 2119 1422"> </div>

### Nicht klassierte Daten (PMF und CDF)

Die absolute Häufigkeit kann als Funktion  $h: \mathbb{R} \rightarrow \mathbb{R}$  bezeichnet werden.

$$h_i$$

Die relative Häufigkeit kann als Funktion  $f: \mathbb{R} \rightarrow \mathbb{R}$  bezeichnet werden.

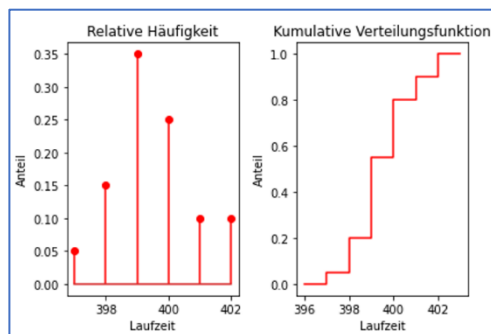
$$f_i = \frac{h_i}{n}$$

### Klassenbildung (Faustregeln)

- Die Klassen sollten gleich breit gewählt werden
- Die Anzahl der Klassen sollte zwischen 5 und 20 liegen, jedoch  $\sqrt{n}$  nicht überschreiben.

### Beispiel:

$a_i$	397	398	399	400	Total
$h_i$	1	3	7	5	16
$f_i$	$\frac{1}{16}$	$\frac{3}{16}$	$\frac{7}{16}$	$\frac{5}{16}$	1
$H_i$	1	4	11	16	
$F_i$	$\frac{1}{16}$	$\frac{4}{16}$	$\frac{11}{16}$	$\frac{16}{16}$	



### Klassierte Daten (PDF und CDF)

Die absolute Häufigkeitsdichtefunktion erhält man, indem der Wert der absoluten Häufigkeit  $h_i$ , durch die Klassenbreite (Säulenbreite)  $d_i$  geteilt wird.

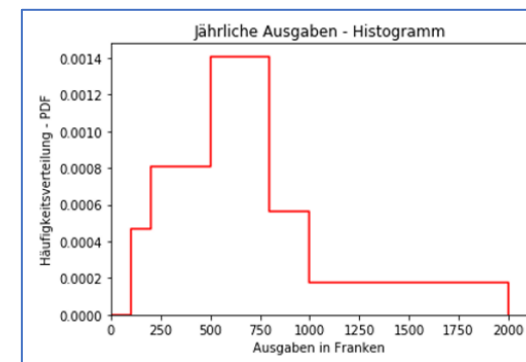
$$h(x) = \frac{h_i}{d_i}$$

Die relative Häufigkeitsdichtefunktion (PDF)  $f: \mathbb{R} \rightarrow [0,1]$  erhält man aus der absoluten Häufigkeitsdichtefunktion, indem man den Wert durch die Stichprobengröße  $n$  teilt.

$$PDF = f(x) = \frac{h(x)}{n}$$

### Beispiel:

Klassen	100 – 200	200 – 500	500 – 800	800 – 1000	Total
$h_i$	35	182	317	84	618
$f_i$	$\frac{35}{618}$	$\frac{182}{618}$	$\frac{317}{618}$	$\frac{84}{618}$	Area = 1
$d_i$	100	300	300	200	
$h(x)$	$\frac{35}{100}$	$\frac{182}{300}$	$\frac{317}{300}$	$\frac{84}{200}$	
$f(x)$	$\frac{35}{100 \cdot 618}$	$\frac{182}{300 \cdot 618}$	$\frac{317}{300 \cdot 618}$	$\frac{84}{200 \cdot 618}$	



# Deskriptive Statistik

<p><b>Varianz</b> <math>s_x^2, s_y^2</math></p> $(s_x)^2 = \overline{x^2} - \bar{x}^2, \quad (s_y)^2 = \overline{y^2} - \bar{y}^2$ <p><b>Kovarianz</b> <math>s_{xy}</math></p> $s_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}), \quad s_{xy} = \overline{xy} - \bar{x} \cdot \bar{y}$	<p><b>Abkürzungen</b></p> $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ $\overline{xy} = \frac{1}{n} \sum_{i=1}^n x_i \cdot y_i$	<p><b>Varianz (Ränge)</b> <math>(s_{rg(x)})^2, (s_{rg(y)})^2</math></p> $(s_{rg(x)})^2 = \overline{rg(x)^2} - (\overline{rg(x)})^2, \quad (s_{rg(y)})^2 = \overline{rg(y)^2} - (\overline{rg(y)})^2$ <p><b>Kovarianz (Ränge)</b> <math>s_{rg(xy)}</math></p> $s_{rg(xy)} = \overline{rg(xy)} - \overline{rg(x)} \cdot \overline{rg(y)} = \overline{rg(xy)} - \frac{(n+1)^2}{4}$																					
<p>Der <b>Korrelationskoeffizient (Pearson)</b> <math>r_{xy}</math></p> $r_{xy} = \frac{s_{xy}}{s_x \cdot s_y} = \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{\sqrt{\overline{x^2} - \bar{x}^2} \cdot \sqrt{\overline{y^2} - \bar{y}^2}}$ <p>Ist der Korrelationskoeffizient <math>r_{xy}</math></p> <ul style="list-style-type: none"> <li>• <math>r_{xy} \approx 1 \rightarrow</math> starker positiver linearer Zusammenhang</li> <li>• <math>r_{xy} \approx -1 \rightarrow</math> starker negativer linearer Zusammenhang</li> <li>• <math>r_{xy} \approx 0 \rightarrow</math> Keine lineare Korrelation</li> </ul>	<p><b>Korrelationskoeffizient (Spearman)</b> <math>r_{sp}</math></p> $r_{sp} = \frac{s_{rg(xy)}}{s_{rg(x)} \cdot s_{rg(y)}} = \frac{\overline{rg(xy)} - \overline{rg(x)} \cdot \overline{rg(y)}}{\sqrt{\overline{rg(x)^2} - (\overline{rg(x)})^2} \cdot \sqrt{\overline{rg(y)^2} - (\overline{rg(y)})^2}}$ <p>Vereinfachte Formel, sofern <i>alle Ränge unterschiedlich</i> sind</p> $r_{sp} = 1 - \frac{6 \cdot \sum_{i=1}^n d_i^2}{n \cdot (n^2 - 1)}, \quad \text{mit } d_i = rg(x_i) - rg(y_i)$ <p><b>Ränge</b></p> <p>Der Rang <math>rg(x_i)</math> des Stichprobenwertes <math>x_i</math> ist definiert als der Index von <math>x_i</math> in der nach der Grösse geordneten Stichprobe.</p> <table border="1" data-bbox="1077 1098 2078 1209"> <thead> <tr> <th><math>i</math></th> <th>1</th> <th>2</th> <th>3</th> <th>4</th> <th>5</th> <th>6</th> </tr> </thead> <tbody> <tr> <td><math>x_i</math></td> <td>23</td> <td>27</td> <td>35</td> <td>35</td> <td>42</td> <td>59</td> </tr> <tr> <td><math>rg(x_i)</math></td> <td>1</td> <td>2</td> <td>3.5</td> <td>3.5</td> <td>5</td> <td>6</td> </tr> </tbody> </table>		$i$	1	2	3	4	5	6	$x_i$	23	27	35	35	42	59	$rg(x_i)$	1	2	3.5	3.5	5	6
$i$	1	2	3	4	5	6																	
$x_i$	23	27	35	35	42	59																	
$rg(x_i)$	1	2	3.5	3.5	5	6																	
<p><u>Bemerkungen</u></p> <p>Auch wenn zwischen zwei Grössen eine Korrelation besteht, so muss das noch lange nicht einen <i>kausalen Zusammenhang</i> bedeuten. Man spricht von <i>Scheinkorrelation</i>.</p>																							
<p><b>Graphische Darstellung</b></p> <ul style="list-style-type: none"> <li>• Form linear / gekrümmt</li> <li>• Richtung positiver / negativer Zusammenhang</li> <li>• Stärke starke / schwache Streuung</li> </ul>	<p><b>Bivariate Daten</b> (Merkmale)</p> <ul style="list-style-type: none"> <li>• 2x kategoriell Kontingenztabelle + Mosaikplot</li> <li>• 1x kategoriell + 1x metrisch Boxplot oder Striptchart</li> <li>• 2x metrisch Streudiagramm</li> </ul>																						

# Kombinatorik

<p><b>Fakultät</b></p> $n! = 1 \cdot 2 \cdot \dots \cdot n = \prod_{k=1}^n k$	<p><b>Binomialkoeffizient</b></p> <p>Wie viele Möglichkeiten gibt es <math>k</math> Objekte aus einer Gesamtheit von <math>n</math> Objekten auszuwählen.</p> $\binom{n}{k} = \frac{n!}{(n-k)! \cdot k!}$																
<p><b>Systematik</b></p> <ul style="list-style-type: none"> <li><math>k</math> Anzahl Stellen</li> <li><math>n</math> Anzahl Optionen pro Stelle</li> </ul>	<table border="1"> <thead> <tr> <th colspan="2">Variation (mit Reihenfolge)</th> <th colspan="2">Kombination (ohne Reihenfolge)</th> </tr> <tr> <th>Mit Wiederholung</th> <th>Ohne Wiederholung</th> <th>Mit Wiederholung</th> <th>Ohne Wiederholung</th> </tr> </thead> <tbody> <tr> <td><math>n^k</math></td> <td><math>\frac{n!}{(n-k)!}</math></td> <td><math>\binom{n+k-1}{k}</math></td> <td><math>\binom{n}{k}</math></td> </tr> <tr> <td>Zahenschloss</td> <td>Schwimmwettkampf</td> <td>Zahnarzt</td> <td>Lotto</td> </tr> </tbody> </table>	Variation (mit Reihenfolge)		Kombination (ohne Reihenfolge)		Mit Wiederholung	Ohne Wiederholung	Mit Wiederholung	Ohne Wiederholung	$n^k$	$\frac{n!}{(n-k)!}$	$\binom{n+k-1}{k}$	$\binom{n}{k}$	Zahenschloss	Schwimmwettkampf	Zahnarzt	Lotto
Variation (mit Reihenfolge)		Kombination (ohne Reihenfolge)															
Mit Wiederholung	Ohne Wiederholung	Mit Wiederholung	Ohne Wiederholung														
$n^k$	$\frac{n!}{(n-k)!}$	$\binom{n+k-1}{k}$	$\binom{n}{k}$														
Zahenschloss	Schwimmwettkampf	Zahnarzt	Lotto														
<p><b>Variation mit Wiederholung</b> (Zahenschloss)</p> <p>Wie viele Möglichkeiten gibt es bei einem Zahenschloss (0 – 9) mit 6 Zahlenkränzen?</p> $n = 10, \quad k = 6$ $n^k = 10^6$	<p><b>Variation ohne Wiederholung</b> (Schimmwettkampf)</p> <p>Bei einem Schwimmwettkampf starten 10 Teilnehmer. Wie viele mögliche Platzierungen der ersten drei Plätze (Podest) gibt es?</p> $n = 10, \quad k = 3$ $\frac{n!}{(n-k)!} = \frac{10!}{(10-3)!} = \frac{10!}{(7)!}$																
<p><b>Kombination mit Wiederholung</b> (Zahnarzt)</p> <p>3 Spielzeuge werden aus 5 Töpfen gezogen. Jeder Topf ist mit einer (unterschiedlichen) Art von Spielzeug befüllt.</p> <p>Wie viele Möglichkeiten hat das Kind?</p> $n = 5, \quad k = 3$ $\binom{n+k-1}{k} = \binom{5+3-1}{3} = \binom{7}{3}$	<p><b>Kombination ohne Wiederholung</b> (Lotto)</p> <p>Wie gross sind die Chancen beim Lotto 6 aus 49 Zahlen richtig zu ziehen?</p> <p>Jede Zahl ist nur einmal vorhanden und die Zahlen werden nicht zurückgelegt. Die Reihenfolge in der gezogen wird spielt keine Rolle.</p> $n = 49, \quad k = 6$ $\binom{n}{k} = \binom{49}{6}$																

## Ideen

- Berechnung durch Aufteilung in mehrere Kombinationen
- Berechnung über Inverse
- Prozente = Wahrscheinlichkeit / Gesamt-Wahrscheinlichkeit

## Wahrscheinlichkeit, dass ein Ereignis $x$ -Mal Auftritt

Beim Rommé spielt man mit  $110$  Karten: *sechs* davon sind *Joker*. Zu Beginn eines Spiels erhält jeder Spieler genau  $12$  Karten.

In wieviel Prozent aller möglichen Fälle sind darunter **genau zwei Joker**?

$$\frac{\binom{6}{2} \cdot \binom{104}{10}}{\binom{110}{12}}$$

In wieviel Prozent aller möglichen Fälle ist darunter **mindestens ein Joker**?

$$1 - \frac{\binom{104}{12}}{\binom{110}{12}}$$

Von  $100$  Glühbirnen sind genau *drei defekt*. Es werden nun  $6$  Glühbirnen zufällig ausgewählt.

Wie viele Möglichkeiten gibt es, wenn sich **mindestens eine defekte** Glühbirne in der Auswahl befinden soll?

$$\binom{100}{6} - \binom{97}{6} = 203'880'032$$

Mit wie viel Prozent Chancen ist bei einer Auswahl von  $6$  Glühbirnen **keine defekt**?

$$\frac{\binom{97}{6}}{\binom{100}{6}}$$

Sind in mehr als 60% aller Fälle von vier (nicht gleichaltrigen) Geschwistern mindestens zwei im gleichen Monat geboren?

$$1 - \frac{12 \cdot 11 \cdot 10 \cdot 9}{12^4}$$

Auf wie viele Arten lassen sich  $10$  Bücher in ein Regal reihen?

$$n = 10, \quad k = 10$$

$$\frac{n!}{(n-k)!} = 10!$$

## Aufteilung in mehrere Kombinationen

Wie viele Worte lassen sich aus den Buchstaben des Wortes ABRAKADABRA bilden? (Nur Worte in denen alle Buchstaben vorkommen!)

$$A = 5x, \quad B = 2x, \quad R = 2x, \quad D = 1x, \quad K = 1x$$

$$\binom{11}{5} \cdot \binom{6}{2} \cdot \binom{4}{2} \cdot \binom{2}{1} \cdot \binom{1}{1} = 83160$$

# Elementare Wahrscheinlichkeitsrechnung

Ergebnisraum  $\Omega$ : Menge aller möglichen Ergebnisse des Zufallsexperiments. Zähldichte  $\rho: \Omega \rightarrow [0,1]$ , die jedem Ereignis seine Wahrscheinlichkeit zuordnet.

Für jedes Ereignis aus  $\Omega$  gleichwahrscheinlich ist, wird  $(\Omega, P)$  *Laplace-Raum* genannt.

$$P(M) = \frac{|M|}{|\Omega|}$$

Zwei Ereignisse  $A$  und  $B$  heißen **stochastisch unabhängig**, falls

$$P(A \cap B) = P(A) \cdot P(B)$$

Zwei Zufallsvariablen  $X: \Omega \rightarrow \mathbb{R}$  und  $Y: \Omega \rightarrow \mathbb{R}$  heißen **stochastisch abhängig**, falls

$$P(X = x, Y = y) = P(X = x) \cdot P(Y = y), \quad \text{für alle } x, y \in \mathbb{R}$$

Für **stochastisch unabhängige** Zufallsvariablen  $X$  und  $Y$  gilt

$$E(X \cdot Y) = E(X) \cdot E(Y), \quad V(X + Y) = V(X) + V(Y)$$

**Kenngrossen** (Varianz und Erwartungswert)

$$E(X + Y) = E(X) + E(Y), \quad E(\alpha X) = \alpha E(X)$$

$$V(X) = E(X^2) - E(X)^2 = \left[ \sum_{x \in \mathbb{R}} P(X = x) \cdot x^2 \right] - E(X)^2$$

$$V(\alpha X + \beta) = \alpha^2 \cdot V(X), \quad S(X) = \sqrt{V(X)}$$

Wahrscheinlichkeit eines Ereignisses  $B$  mit Vorbedingung  $A$

$$P(B|A) = \frac{P(B \cap A)}{P(A)}$$

**Multiplikationssatz**

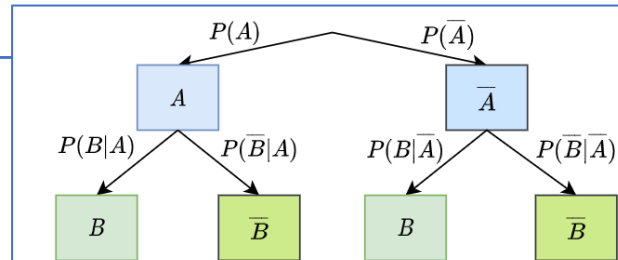
$$P(A \cap B) = P(A) \cdot P(B|A) = P(B) \cdot P(A|B)$$

Satz von der **Totalen Wahrscheinlichkeit**

$$P(B) = P(A) \cdot P(B|A) + P(\bar{A}) \cdot P(B|\bar{A})$$

Satz von **Bayes**

$$P(A|B) = \frac{P(A) \cdot P(B|A)}{P(B)}$$



	A	$\bar{A}$	$\Sigma$
B	$P(A \cap B)$	$P(\bar{A} \cap B)$	$P(B)$
$\bar{B}$	$P(A \cap \bar{B})$	$P(\bar{A} \cap \bar{B})$	$P(\bar{B})$
$\Sigma$	$P(A)$	$P(\bar{A})$	$P(\Omega)$

**Spezielle Verteilung** (Summary)

*Diskret*

$$E(X) = \sum_{x \in \mathbb{R}} f(x) \cdot x$$

$$V(X) = \sum_{x \in \mathbb{R}} f(x) \cdot (x - E(X))^2$$

*Stetig*

$$E(X) = \int_{-\infty}^{\infty} f(x) \cdot x \, dx$$

$$V(X) = \int_{-\infty}^{\infty} f(x) \cdot (x - E(X))^2 \, dx$$

## Spezielle Verteilungen

<p><b>Bernoulli-Verteilung</b> (Einmaliges zurücklegen)</p> <p>Bernoulli-Experimente sind Zufallsexperimente mit nur zwei möglichen Ergebnissen. Wir bezeichnen diese Ergebnisse mit 1 und 0.</p> $P(X = 1) = p, \quad P(X = 0) = 1 - p = q$ <ol style="list-style-type: none"> <li><math>E(X) = E(X^2) = p</math></li> <li><math>V(X) = p \cdot (1 - p)</math></li> </ol>	<p><b>Approximation</b> durch die Normalverteilung</p> <ul style="list-style-type: none"> <li><b>Binomialverteilung:</b> <math>\mu = np, \sigma^2 = npq</math></li> <li><b>Poissonverteilung:</b> <math>\mu = \lambda, \sigma^2 = \lambda</math></li> </ul> $P(a \leq X \leq b) = \sum_{x=a}^b P(X = x) \approx \Phi_{\mu, \sigma} \left( b + \frac{1}{2} \right) - \Phi_{\mu, \sigma} \left( a - \frac{1}{2} \right)$ <p><b>Faustregel</b> Die Approximation (Binomialverteilung) kann verwendet werden, wenn <math>npq &gt; 9</math></p> <p>Für grosses <math>n</math> (<math>n \geq 50</math>) und kleiner <math>p</math> (<math>p \leq 0.1</math>) kann Binomial- durch die Poisson-Verteilung approximiert werden</p> $B(n, p) \approx Poi(n \cdot p)$	
<p>Eine <b>Hypergeometrische</b> Verteilung kann durch eine <b>Binomialverteilung</b> angenähert werden, wenn <math>n \leq \frac{N}{20}</math></p> $H(N, M, N) \approx B\left(n, \frac{M}{N}\right)$		
<p><b>Hypergeometrische Verteilung</b> (Ohne zurücklegen)</p> <ul style="list-style-type: none"> <li><math>N</math> = Objekte gesamthaft</li> <li><math>M</math> = Objekte einer bestimmten Sorte</li> <li><math>n</math> = Stichprobengrösse</li> <li><math>x</math> = Merkmalsträger</li> </ul>	<p><b>Binomialverteilung</b> (Mit zurücklegen)</p> <ul style="list-style-type: none"> <li><math>n</math> = Anzahl Wiederholungen</li> <li><math>p</math> = Wahrscheinlichkeit für ein Ergebnis 1</li> <li><math>q = 1 - p</math></li> </ul>	<p><b>Poisson Verteilung</b></p> <ul style="list-style-type: none"> <li><math>\lambda</math> = Rate</li> </ul>
$P(X = x) = \frac{\binom{M}{x} \cdot \binom{N-M}{n-x}}{\binom{N}{n}}$	$P(X = x) = \binom{n}{x} \cdot p^x \cdot q^{n-x}$	$P(X = x) = \frac{\lambda^x}{x!} \cdot e^{-\lambda}, \quad \lambda > 0$
<p>Schreibweise: <math>X \sim H(N, M, n)</math></p> <ol style="list-style-type: none"> <li><math>\mu = E(X) = n \cdot \frac{M}{N}</math></li> <li><math>\sigma^2 = V(X) = n \cdot \frac{M}{N} \cdot \left(1 - \frac{M}{N}\right) \cdot \frac{N-n}{N-1}</math></li> <li><math>\sigma = S(X) = \sqrt{V(X)}</math></li> </ol>	<p>Schreibweise: <math>X \sim B(n; p)</math></p> <ol style="list-style-type: none"> <li><math>\mu = E(X) = np</math></li> <li><math>\sigma^2 = V(X) = npq</math></li> <li><math>\sigma = S(X) = \sqrt{npq}</math></li> </ol>	<p>Schreibweise <math>X \sim Poi(\lambda)</math></p> <ol style="list-style-type: none"> <li><math>\mu = E(X) = \lambda</math></li> <li><math>\sigma^2 = V(X) = \lambda</math></li> <li><math>\sigma = S(X) = \sqrt{\lambda}</math></li> </ol>

## Spezielle Verteilungen

Bei einer *stetigen* Zufallsvariable  $X$  lässt sich die Verteilungsfunktion als Integral einer Funktion  $f$  darstellen

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(u) \cdot du$$

Liegt eine beliebige Normalverteilung  $N(\mu, \sigma)$  vor, muss *standardisiert* werden. Statt ursprünglichen Zufallsvariablen  $X$  betrachtet man die Zufallsvariable

$$U = \frac{X - \mu}{\sigma}$$

### Gauss-Verteilung oder Normalverteilung

Die stetige Zufallsvariable  $X$  folgt der *Normalverteilung* mit den Parametern,  $\mu, \sigma \in \mathbb{R}, \sigma > 0$ , wenn sie folgende Dichtefunktion hat:

$$\varphi_{\mu, \sigma}(x) = \frac{1}{\sqrt{2\pi} \cdot \sigma} \cdot e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$

Schreibweise:  $X \sim N(\mu; \sigma)$

Ist  $\mu = 0$  und  $\sigma = 1$ , so spricht man von der *Standardnormalverteilung*. Ihre Dichtefunktion wird einfach mit  $\varphi$  bezeichnet; sie ist gegeben durch

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{1}{2}x^2}$$

### Die Verteilungsfunktion der Normalverteilung

Die kumulative Verteilungsfunktion (CDF) von  $\varphi_{\mu, \sigma}(x)$  wird mit  $\Phi_{\mu, \sigma}(x)$  bezeichnet. Sie ist definiert durch

$$\Phi_{\mu, \sigma}(x) = P(X \leq x) = \int_{-\infty}^x \varphi_{\mu, \sigma}(t) dt = \frac{1}{\sqrt{2\pi} \cdot \sigma} \cdot \int_{-\infty}^x e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2} dt$$

### Erwartungswert und Varianz der Normalverteilung

Für eine Zufallsvariable  $X \sim N(\mu; \sigma)$  gilt

$$E(X) = \mu, \quad V(X) = \sigma^2$$

### Zentraler Grenzwertsatz

Für eine Folge  $X_1, X_2, \dots, X_n$  von Zufallsvariablen definieren wir die *n-te Summe*  $S_n$  und das *arithmetische Mittel*  $\bar{X}_n$ .

Haben alle Zufallsvariablen denselben Erwartungswert  $E(X_i) = \mu$  und dieselbe Varianz  $V(X_i) = \sigma^2$  so folgt

$$E(S_n) = n \cdot \mu, \quad V(S_n) = n \cdot \sigma^2, \quad E(\bar{X}_n) = \mu, \quad V(\bar{X}_n) = \frac{\sigma^2}{n} = \frac{1}{n^2} \cdot V(S_n)$$

Sind die Zufallsvariablen alle identisch  $N(\mu, \sigma)$  verteilt, so sind die Summe  $S_n$  und das arithmetische Mittel  $\bar{X}_n$  wieder normalverteilt mit...

- $S_n$ :  $N(n \cdot \mu, \sqrt{n} \cdot \sigma)$
- $\bar{X}_n$ :  $N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$

Verteilungsfunktion  $F_n(u)$  der dazugehörigen standardisierten Zufallsvariable

$$U_n = \frac{((X_1 + X_2 + \dots + X_n) - n\mu)}{\sqrt{n} \cdot \sigma} = \frac{(\bar{X} - \mu)}{\frac{\sigma}{\sqrt{n}}}$$

Konvergiert für  $n \rightarrow \infty$  gegen die Verteilungsfunktion  $\Phi(u)$  der Standardnormalverteilung:

$$\lim_{n \rightarrow \infty} F_n(u) = \Phi(u) = \frac{1}{\sqrt{2\pi}} \cdot \int_{-\infty}^u e^{-\frac{1}{2}t^2} dt$$



# Methode der kleinsten Quadrate

## Lineare Regression

Gegeben sind Datenpunkte  $(x_i; y_i)$  mit  $1 \leq i \leq n$ . Die *Residuen / Fehler*  $\epsilon_i = g(x_i) - y_i$  dieser Datenpunkte sind Abstände in  $y$ -Richtung zwischen  $y_i$  und der Geraden  $g$ . Die *Ausgleichs- oder Regressiongerade*, sei diejenige Gerade für die, die Summe der quadrierten Residuen  $\sum_{i=1}^n \epsilon_i^2$  am kleinsten ist.

### Regressionsgerade

Die *Regressionsgerade*  $g(x) = mx + d$  mit den Parametern  $m$  und  $d$  ist die Gerade, für welche die *Residualvarianz*  $s_\epsilon^2$  minimal ist.

$$\text{Steigung: } m = \frac{s_{xy}}{s_x^2}, \quad y\text{-Achsenabschnitt: } d = \bar{y} - m\bar{x}, \quad s_\epsilon^2 = s_y^2 - \frac{s_{xy}^2}{s_x^2}$$

### Bestimmtheitsmass

Die Totale Varianz setzt sich zusammen aus der Residualvarianz und der Varianz der prognostizierten Werte

- $s_y^2$  Totale Varianz
- $s_{\hat{y}}^2$  prognostizierte (erklärte) Varianz
- $s_\epsilon^2$  Residualvarianz

$$s_y^2 = s_\epsilon^2 + s_{\hat{y}}^2$$

Das *Bestimmtheitsmass*  $R^2$  beurteilt die globale Anpassungsgüte einer Regression über den Anteil der prognostizierten Varianz  $s_{\hat{y}}^2$  an der totalen Varianz  $s_y^2$

$$R^2 = \frac{s_{\hat{y}}^2}{s_y^2}$$

Das *Bestimmtheitsmass*  $R^2$  entspricht dem Quadrat des *Korrelationskoeffizienten*

$$R^2 = \frac{s_{xy}^2}{s_x^2 \cdot s_y^2} = (r_{xy})^2$$

### Residuenquadrate

- $y_i$ : beobachtete  $y$ -Werte
- $\hat{y}_i$ : prognostizierte bzw. erklärte  $y$ -Werte
- $\epsilon$ : Residuen (oder auch Fehler)

$$\sum_{i=1}^n \left( \frac{y_i - g(x_i)}{\epsilon_i} \right)^2 = \sum_{i=1}^n \left( \frac{y_i - \hat{y}_i}{\epsilon_i} \right)^2$$

### Kleinste Quadrate (KQM)

Die Parameter  $m$  und  $q$  werden mit der Matrix  $A$  berechnet

$$A = \begin{pmatrix} x_1 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{pmatrix}, \quad A^T \cdot A \cdot \begin{pmatrix} m \\ q \end{pmatrix} = A^T \cdot \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

### Linearisierungsfunktionen

Ausgangsfunktion	Transformation
$y = q \cdot x^m$	$\log(y) = \log(q) + m \cdot \log(x)$
$y = q \cdot m^x$	$\log(y) = \log(q) + \log(m) \cdot x$
$y = q \cdot e^{m \cdot x}$	$\ln(y) = \ln(q) + m \cdot x$
$y = \frac{1}{q + m \cdot x}$	$V = q + m \cdot x; V = \frac{1}{y}$
$y = q + m \cdot \ln(x)$	$y = q + m \cdot U; u = \ln(x)$
$y = \frac{1}{q \cdot m^x}$	$\log\left(\frac{1}{y}\right) = \log(q) + \log(m) \cdot x$

## Schliessende Statistik

<p><b>Erwartungstreue Schätzfunktion</b></p> <p>Eine Schätzfunktion <math>\Theta</math> eines Parameters <math>\theta</math> heisst <i>erwartungstreu</i>, wenn</p> $E(\Theta) = \theta$	<p><b>Effizienz Schätzfunktion</b></p> <p>Gegeben sind zwei <i>erwartungstreue</i> Schätzfunktionen <math>\Theta_1</math> und <math>\Theta_2</math> desselben Parameters <math>\theta</math>. Man nennt <math>\Theta_1</math> <i>effizienter als</i> <math>\Theta_2</math>, falls</p> $V(\Theta_1) < V(\Theta_2)$	<p><b>Konsistenz Schätzfunktion</b></p> <p>Eine Schätzfunktion <math>\Theta</math> heisst <i>konsistent</i>, wenn</p> $E(\Theta) \rightarrow \theta \text{ und } V(\Theta) \rightarrow 0 \text{ für } n \rightarrow \infty$
<p>Grundgesamtheit mit Erwartungswert <math>\mu</math>, Varianz <math>\sigma^2</math> und Zufallsstichprobe <math>X_1, X_2, X_3</math>. Die folgende Schätzfunktion ist gegeben.</p> $\Theta_1 = \frac{1}{3} \cdot (2X_1 + X_2)$		
<p>Ist diese Schätzfunktion <i>erwartungstreu</i> (Parameter: <math>\mu</math>)?</p> $E(\Theta_1) = E\left(\frac{1}{3} \cdot (2X_1 + X_2)\right) = \frac{1}{3} \cdot (2E(X_1) + E(X_2))$ $E(\Theta_1) = \frac{1}{3} \cdot (2\mu + \mu) = \frac{3\mu}{3} = \mu$ <p>Da <math>E(\Theta_1) = \mu</math> ist die Funktion erwartungstreu.</p>	<p>Berechne die <i>Effizienz</i> der Schätzfunktion (Parameter: <math>\sigma^2</math>):</p> $V(\Theta_1) = V\left(\frac{1}{3} \cdot (2X_1 + X_2)\right) = \frac{1}{9} \cdot V(2X_1 + X_2) = \frac{1}{9} \cdot (V(2X_1) + V(X_2))$ $V(\Theta_1) = \frac{1}{9} \cdot (4 \cdot V(X_1) + V(X_2)) = \frac{1}{9} \cdot (4\sigma^2 + \sigma^2) = \frac{5\sigma^2}{9}$	
<p><b>Likelihood-Funktion</b></p> <p>Wir betrachten eine Zufallsvariable <math>X</math> und ihre Dichte (PDF)</p> $f_x(x \theta)$ <p>Welche von <math>x</math> und einem oder mehreren Parametern <math>\theta</math> abhängig sind. Für eine Stichprobe vom Umfang <math>n</math> mit <math>x_1, \dots, x_n</math> nennen wir die vom Parameter <math>\theta</math> abhängige Funktion ... die Likelihood-Funktion der Stichprobe.</p> $L(\theta) = f_x(x_1 \theta) \cdot f_x(x_2 \theta) \cdot \dots \cdot f_x(x_n \theta)$		<p><b>Vorgehen – Likelihood Funktion</b></p> <ol style="list-style-type: none"> <li>1. Likelihood-Funktion bestimmen</li> <li>2. Maximalstelle der Funktion bestimmen <ul style="list-style-type: none"> <li>• (Partielle) Ableitung <math>L'(\theta) = 0</math></li> </ul> </li> </ol>
<p><b>Erwartungswert</b> (Funktion, Wert)</p> $\bar{X} = \frac{1}{n} \cdot \sum_{i=1}^n X_i, \quad \hat{\mu} = \bar{x} = \frac{1}{n} \cdot \sum_{i=1}^n x_i$	<p><b>Varianz</b> (Funktion, Wert)</p> $S^2 = \frac{1}{n-1} \cdot \sum_{i=1}^n (X_i - \bar{X})^2, \quad \hat{\sigma}^2 = s^2 = \frac{1}{n-1} \cdot \sum_{i=1}^n (x_i - \bar{x})^2$	

### Vertrauensintervalle

Wir legen eine grosse Wahrscheinlichkeit  $\gamma$  fest (z.B.  $\gamma = 95\%$ ).  $\gamma$  heisst *statistische Sicherheit* oder *Vertrauensniveau*.  $\alpha = 1 - \gamma$  ist die sogenannte *Irrtumswahrscheinlichkeit*.

Dann bestimmen wir zwei Zufallsvariablen  $\theta_u$  und  $\theta_o$  so, dass sie den wahren Parameterwert  $\theta$  mit der Wahrscheinlichkeit  $\gamma$  einschliessen:

$$P(\theta_u \leq \theta \leq \theta_o) = \gamma$$

**Spezialfall: Anteilswert  $p$  einer Bernoulli-Verteilung** (Funktion, Wert)

$$\bar{X} = \frac{1}{n} \cdot \sum_{i=1}^n X_i, \quad \hat{p} = \bar{x} = \frac{1}{n} \cdot \sum_{i=1}^n x_i$$

### Intervallschätzung

Geben Sie das Vertrauensintervall für  $\mu$  an ( $\sigma^2$  unbekannt). Gegeben sind...

$$n = 10, \quad \bar{x} = 102, \quad s^2 = 16, \quad \gamma = 0.99$$

1. Verteilungstyp bestimmen

Verteilungstyp mit Param  $\mu$  und  $\sigma^2$  unbekannt  $\rightarrow$  *T-Verteilung*

2. Verteilung und Quantile berechnen

$$f = n - 1 = 9, \quad p = \frac{1 + \gamma}{2} = 0.995, \quad c = t_{(p;f)} = t_{(0.995;9)} = 3.25$$

3. Vertrauensintervall bestimmen

$$e = c \cdot \frac{S}{\sqrt{n}} = 4.111, \quad \theta_u = \bar{X} - e = 97.89, \quad \theta_o = \bar{X} + e = 106.11$$

	Verteilung der Grundgesamtheit	Param	Schätzfunktionen	Standardisierte Zufallsvariable	Verteilung / Quantile	Intervallgrenzen
1	Normalverteilung (Varianz $\sigma^2$ bekannt)	$\mu$	$\bar{X}$	$U = \frac{(\bar{X} - \mu)}{\sigma/\sqrt{n}}$	Standardnormalverteilung $c = u_p, p = \frac{1+\gamma}{2}$	$e = c \cdot \frac{\sigma}{\sqrt{n}}$ $\theta_u = \bar{X} - e, \quad \theta_o = \bar{X} + e$
2	Normalverteilung (Varianz $\sigma^2$ unbekannt und $n \leq 30$ ; sonst Fall 1 mit $s$ als Schätzwert für $\sigma$ )	$\mu$	$\bar{X}, S^2$	$T = \frac{(\bar{X} - \mu)}{s/\sqrt{n}}$	<i>t-Verteilung</i> $c = t_{(p;f=n-1)}, p = \frac{1+\gamma}{2}$	$e = c \cdot \frac{S}{\sqrt{n}}$ $\theta_u = \bar{X} - e, \quad \theta_o = \bar{X} + e$
3	Normalverteilung	$\sigma^2$	$\bar{X}, S^2$	$Z = (n-1) \frac{s^2}{\sigma^2}$	Chi-Quadrat-Verteilung $c_1 = z_{(p_1;f=n-1)}, p_1 = \frac{1-\gamma}{2}$ $c_2 = z_{(p_2;f=n-1)}, p_2 = \frac{1+\gamma}{2}$	$e = (n-1) \cdot \frac{S^2}{c_2}$ $\theta_u = \frac{e}{c_2}, \quad \theta_o = \frac{e}{c_1}$
4	Bernoulli-Verteilung mit $n\hat{p}(1-p) > 9$	$p$	$\bar{X}$ $P(X_i = 1) = p$	$U = \frac{\bar{X} - p}{\sqrt{p(1-p)/n}}$	Standardnormalverteilung näherungsweise $c = u_q, q = \frac{1+\gamma}{2}$	$e = c \cdot \sqrt{\frac{\bar{X} \cdot (1-\bar{X})}{n}}$ $\theta_u = \bar{X} - e, \quad \theta_o = \bar{X} + e$
5	Beliebig mit $n > 30$	$\mu, \sigma^2$	Wie im Fall 1 (gegebenenfalls mit $s$ als Schätzwert für $\sigma$ ) bzw. wie im Fall 3			