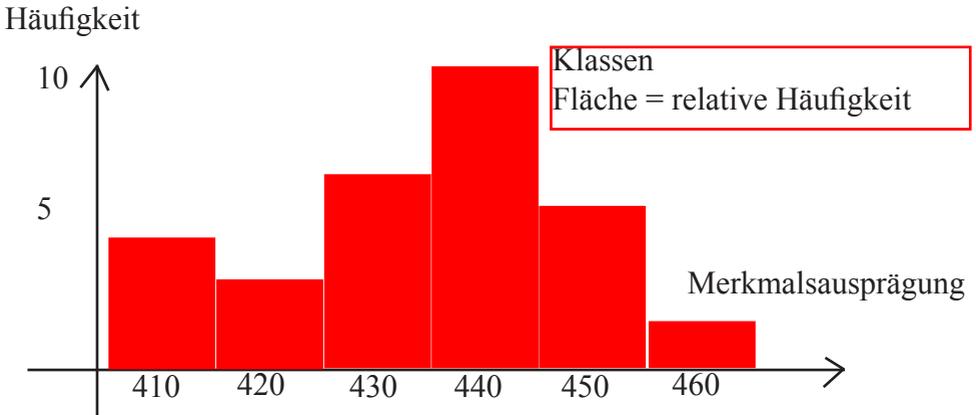


<b>Merkmalsausprägung <math>x_i</math></b>	410	420	430	440	450	460
<b>Häufigkeit <math>n_i</math> (wieviel mal)</b>	4	3	6	10	5	2



TI - 89:

APPS -> DATA MATRIX EDITOR -> NEW

Typ : Data

C1 : Merkmalsausprägung x;

C2 : alle  $n_i$  Werte;

F5 ->

Calculation Types : one Var

x : C1

Frequent Category: Yes

Frequ : C2

**bei relativer Häufigkeit:**

y Achsen Höhe =

Höhe der Klasse in % oder Merkmalsausprägung

Breite der Klasse in x

$$y_i = a + b * x_i$$

$$\bar{y} = a + b * \bar{x}$$

**Stichprobenumfang n**

**absolute Häufigkeit  $n_i$  :  $n_1 = 4, n_2 = 3...$**

$$\text{relative Häufigkeit } h(x_i) = \frac{n_i}{n} = \frac{4}{30}$$

$$\text{empirischer Mittelwert } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i =$$

$$\frac{1}{30} * (4 * 410 + 3 * 420 + 6 * 430 + 10 * 440 + 5 * 450 + 2 * 460)$$

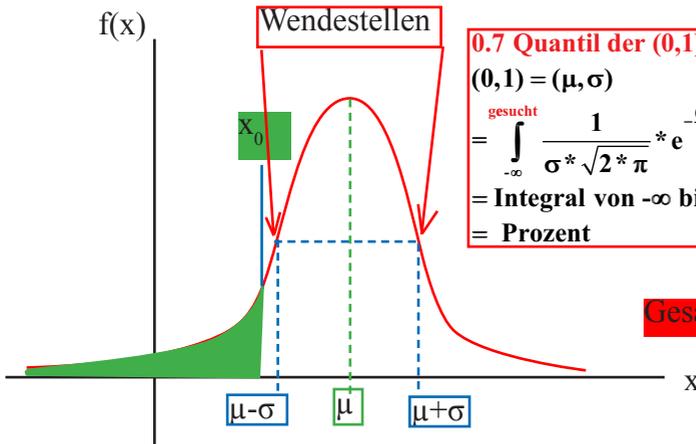
$$\text{bei k Klassen } \bar{x} = \frac{1}{n} \sum_{i=1}^k n_i * x_i \quad x_i = \text{Stabmitten}$$

$$\text{empirische Standardabweichung } s = \sqrt{\frac{1}{n-1} * \sum_{i=1}^k n_i * (\bar{x} - x_i)^2}$$

$$\text{Varianz} = s^2$$

Dichtefunktion, Glockenkurve, Gaussche Kurve oder Normalverteilung:

Fläche entspricht der relativen Häufigkeit:  
Integration von  $-\infty$  bis  $x_0$



**0.7 Quantil der (0,1)-Normverteilung:**

$(0,1) = (\mu, \sigma)$

gesucht

$$= \int_{-\infty}^{\text{gesucht}} \frac{1}{\sigma * \sqrt{2 * \pi}} * e^{-\frac{(x-\mu)^2}{2 * \sigma^2}} * dx = 0.7$$

= Integral von  $-\infty$  bis **gesucht** = Fläche von 0.7  
= **Prozent**

**Gesamtfläche immer 1**

$$f(x) = \frac{1}{\sigma * \sqrt{2 * \pi}} * e^{-\frac{(x-\mu)^2}{2 * \sigma^2}}$$

**relative Häufigkeit = Fläche unter der Kurve = 1:**

$f(x)$  hat das absolute Maximum bei  $x = \mu$

$$f(\mu) = \frac{1}{\sigma * \sqrt{2 * \pi}} \text{ d.h. der Maximalwert wird kleiner, wenn } \sigma \text{ erhöht wird}$$

$f(x)$  hat Wendestellen bei  $x = \mu \pm \sigma$   $f'(x)=0$   $f''(x) \neq 0$

$$\text{Mittelwert} = \bar{x} = \mu = E(X) = \int_{-\infty}^{\infty} x * f(x) * dx$$

$$\text{Standardabweichung} = s = \sigma = V(X) = \sqrt{\int_{-\infty}^{\infty} (x - E(X))^2 * f(x) * dx}$$

$$\text{Varianz} = v = \sqrt{V(X)}$$

**relative Häufigkeit = Fläche unter der Kurve = 1:**

**Ti 89 hat Numerikprobleme:**

**nSolve = Solve**

**nInt =  $\int$**

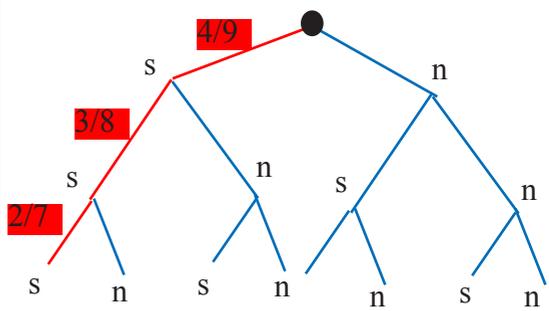


**Zufallsversuch** -> Das Werfen einer Münze  
**Stichprobenraum  $\Omega$**  -> Die Menge aller möglichen Ausgänge, oder Ergebnisse eines Zufallsversuches  
**Ausfall  $\omega$**  -> Ein Ergebnis  
**Ereignis** -> Teilmenge des Stichprobenraumes  $\Omega$   
**Beispiel Würfel:**  
 $\omega_1 = \text{"Kopf"}$      $\omega_2 = \text{"Zahl"}$   
 $\Omega = \{\text{"Kopf"}, \text{"Zahl"}\}$   
**Ereignis**  $A = \{2, 4, 6\}$   
**Assiationsregel**  $p(A) = \sum_{\omega \in A} p(\omega)$   
**Gegeneignis**  $= p(\bar{A}) = 1 - p(A)$   
**Laplace – Gerät** -> Zufallsgerät, bei dem alle Ausfälle dieselbe Wahrscheinlichkeit besitzen  
 $p(A) = \frac{\text{Anzahl günstige Fälle}}{\text{Anzahl mögliche Fälle}}$

**Mehrstufige Zufallsversuche:**

In einem Baum gilt:  
 Die Wahrscheinlichkeit eines Pfades ist gleich dem Produkt aller Wahrscheinlichkeiten längs des Pfades

Reisegruppe mit 4 Schmugglern und 5 ehrliche Leute. Wie gross ist die Wahrscheinlichkeit, dass von 3 Stichproben 3 Schmuggler ertappt werden?  
 Schmuggler = s  
 nicht Schmuggler = n



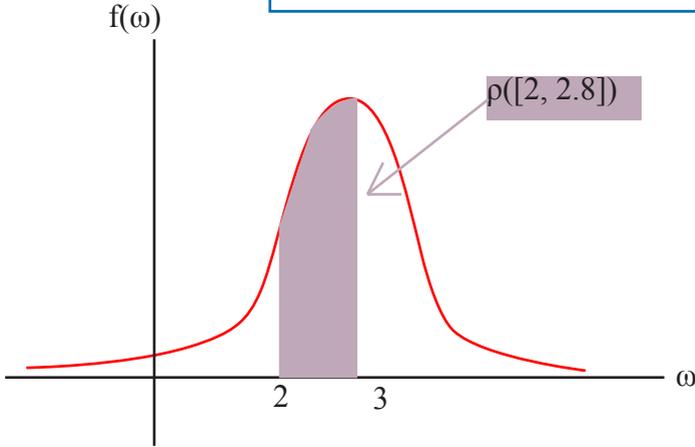
3 mal Schmuggler hintereinander=  
 $(4/9) * (3/8) * (2/7) = 1/21$

**Wahrheitstabelle:**

$\omega$	sss	ssn	snn	sns.....
$p(\omega)$	1/21	...	...	...

← Quersumme = 1

Die Funktion  $f(\omega)$  entspricht also dem Begriff der Dichtfunktion = Wahrscheinlichkeit



**Bei Binominalverteilung:**

$$E(X) = n * p$$

$$V(X) = n * p * (1 - p)$$

**6\* Würfel, 6 verschiedene Zahlen:**

1-2-3-4-5-6      Pfad egal

$$1 * \frac{5}{6} * \frac{4}{6} * \frac{3}{6} * \frac{2}{6} * \frac{1}{6}$$

**bei 12 Würfeln genau 6\* eins:**

$$p = \binom{12}{6} * \binom{1}{2}^6 * \binom{1}{2}^6 \quad \binom{n}{k} * p^k * (1-p)^{n-k}$$

**Aus 12 Plätzen 6 wählen**

**Geordnete Stichproben ohne Zurücklegen:**

Aus  $n$  Objekten sind  $k$  herauszugreifen und in einer Folge anzuordnen.

Wieviele Möglichkeiten gibt es??

Anzahl Möglichkeiten des Versuches = **Anzahl Pfade =  $x$**

$$p(\Omega) = 1$$

$$\frac{x * (n - k)!}{n!} = 1$$

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

-> Es gibt  $\frac{n!}{(n - k)!}$  Möglichkeiten aus  $n$  Objekten  $k$  herauszugreifen und anzuordnen

**Permutationen :**

Eine Anordnung von  $n$  Kugeln (allgemein : Elemente) in einer bestimmten Reihenfolge heisst Permutation. Für die Anzahl der möglichen Permutationen gilt dann:

**1. Alle  $n$  Kugeln sind voneinander verschieden:**

$$p(n) = n!$$

**2. Ungeordneter Fall (Reihenfolge egal)**

totale Anzahl ungeordnete Falle  $x * k! = \frac{n!}{(n - k)!}$

$$x = \frac{n!}{(n - k)! * k!} = \binom{n}{k} \text{ lies "n tief k" (Binomenalkoeffizient)}$$

-> Es gibt  $\binom{n}{k}$  Möglichkeiten aus einer Menge mit  $n$  Elementen eine Teilmenge von  $k$  Elementen herauszugreifen

**3. Unter den  $n$  Kugeln befinden sich jeweils  $n_1, n_2, \dots, n_k$  einander gleiche:**

$$p(n; n_1, \dots, n_k) = \frac{n!}{n_1! * n_2! * \dots * n_k!}$$

**günstige Fälle:**

10 Transistoren, 4 defekt, 3 nehmen

die Wahrscheinlichkeit, dass 2 defekt sind unter den 3 genommenen:

$$\binom{4}{2} * \binom{6}{1} = \binom{4 \text{ defekte}}{2 \text{ davon}} * \binom{6 \text{ gute}}{1 \text{ genommen}}$$

bei mehreren Stufen:  $\binom{n}{k} * \binom{n_2}{k_2}$

bei  $n=k$ :  $n!$

bei Mannschaften:  $/2$

Beispiel Geordnete Stichprobe ohne Zurücklegung:

Aus  $n$  Objekten sind  $k$  herauszugreifen und in einer Folge anzuordnen.

Beispiel ungeordnete Stichprobe ohne Zurücklegung:

Wieviele Möglichkeiten gibt es, in einem Verein mit 70 Mitgliedern einen fünfköpfigen Vorstand zu wählen?

Oder:

Im Lotto einen viere erzielen....

Oder:

Frau Maier will ihre 5 Kinder in einer Reihe anordnen. Auf wieviel Arten kann sie das tun?

Es gibt  $p(3)=3! = 6$  verschiedene Möglichkeiten, 3 verschiedenfarbige Kugeln anzuordnen:



In einer Urne befinden sich 5 Kugeln, 3 weiße und 2 rote. Sie lassen sich auf

$$p(5; 3,2) = \frac{5!}{3!2!} = 10$$

verschiedene Arten anordnen

Beim Zahlenlotto werden aus 45 Zahlen 6 gezogen:

**1. Wieviele Möglichkeiten gibt es, einen 4-er zu erzielen?**

**mache 2. Stufen Versuch:**

1. Stufe 4 Zahlen aus Urne 1 ziehen  $\binom{6}{4}$  Möglichkeiten

2. Stufe 2 Zahlen aus Urne 2 ziehen  $\binom{39}{2}$  Möglichkeiten

**Gesamt :**

$$\binom{6}{4} * \binom{39}{2} \text{Möglichkeiten} = \frac{6!}{(6-4)! * 4!} * \frac{39!}{(39-2)! * 2!}$$

**2. Wie gross ist die Wahrscheinlichkeit für einen Vierer?**

a) entweder über Bäume

$$b) p(\text{"vierer"}) = \frac{\text{Anzahl der günstigen Fälle}}{\text{Anzahl der möglichen Fälle}} = \frac{\text{Anzahl Vierer}}{\text{Anzahl der möglichen Fälle}}$$

$$= \frac{\binom{6}{4} * \binom{39}{2} \text{Möglichkeiten}}{\binom{45}{6}}$$

**Produktregel:**

In einem Zweistufenversuch seien  $n_1$  die Anzahl der Möglichkeiten wie die erste Stufe ausgehen kann und  $n_2$  die Anzahl der Möglichkeiten für die zweite Stufe.

Wenn man für jeden Ausgang der ersten Stufe alle Möglichkeiten der zweiten Stufe offen hat, so berechnet sich die totale Anzahl der Möglichkeiten des Gesamtversuchs als Produkt:

$$n_1 * n_2$$

Zufallsvariablen und ihre verteilung:

$$\sum_{i=1}^n p(X = x_i) * x_i$$

Wahrscheinlichkeit, Wert

Bsp.

Beim Würfeln wird bei geraden Zahlen das doppelte ausgezahlt, bei ungeraden verliert man das dreifache der Zahl

$X(\omega)$ =Gewinn bei Ausfall  $\omega$

Verteilung von X:

$X=x_i$	-3	4	-9	8	-15	12
$p(X=x_i)$	1/6	1/6	1/6	1/6	1/6	1/6

$$E(X) = (-3) * \frac{1}{6} + 4 * \frac{1}{6} - 9 * \frac{1}{6} + 8 * \frac{1}{6} - 15 * \frac{1}{6} + 12 * \frac{1}{6}$$

$$= -\frac{1}{2} \text{ Durchschnittlicher Verlust pro Spiel ist 0.5 Fr., wenn man oft spielt}$$

Unter einem **n-stufigen Bernoulli-Versuch** versteht man die n-Fache Durchführung eines Zufalsversuchs, dessen Stichprobenraum nur aus 2 Elementen besteht:

$$\Omega_0 = \{0,1\}$$

Die 1 wird als Erfolg bezeichnet und

$p = p(1)$  als Erfolgswahrscheinlichkeit

Sei X die Anzahl der Erfolge in einem n-Stufigen Bernoulliversuch.

Dann gilt für die Verteilungsfunktion

(k plätze aus n) bsp. von 12 Münzenwürfen (n), 3\*Kopf(k) :

$$p(X = x_i) = p(X = k) = \binom{n}{k} * p^k * (1-p)^{n-k}$$

oder bei mehreren  $\sum_{k=\text{anfang}}^{\text{ende}} \binom{n}{k} * p^k * (1-p)^{n-k}$

Bsp :

Eine Münze wird 120 mal geworfen. Wie gross ist die Wahrscheinlichkeit, dass zwischen 48 und 60 mal Kopf ergibt?

$$\sum_{k=48}^{60} \binom{120}{k} * 0.5^k * (1-0.5)^{120-k}$$

**Beschreibende Statistik :**

Merkmalsausprägung  $x_i$

relative Häufigkeit  $\frac{n_i}{n}$

empirischer Mittelwert  $\bar{x} = \frac{1}{n} * \sum_{i=1}^k n_i * x_i$

empirische Varianz  $s^2 = \frac{1}{n-1} * \sum_{i=1}^k n_i * (\bar{x} - x_i)^2$

**Wahrscheinlichkeitsrechnung :**

Wert  $x_i$  einer Zufallsgrösse X

$p(X = x_i)$ , X nimmt Wert  $x_i$  an

$E(X) = \sum_{k=1}^n p(X = x_i) * x_i$ , falls nur n Werte

$\int_{-y}^y x * f(x) * dx$ , falls die Werte von X alle

relle Zahlen durchlaufen

Nimmt eine Zufallsvariable X nur endlich viele Werte an mit den Wahrscheinlichkeiten

$p(X = x_i)$  und dem Erwartungswert  $m = E(X)$ , dann :

Varianz der Zufallsvariable X :

$V(X) = \sum_{i=1}^n (x_i - m)^2 * p(X = x_i)$   $s = \sqrt{V(X)}$  heisst Standardabweichung von X

$V(X) = n * p * (1 - p)$

Anzahl Stufen in diesem Bernoulliveruch n

Erfolgswahrscheinlichkeit p

Füher n - stufigen Versuch durch - sagen wir N mal.

-> Grundversuch wird n \* N mal durchgeführt

-> habe ca. n \* N \* p Erfolge erzielt

-> durchschnittliche Anzahl Erfolge pro Bernoulliversuch

» n \* p wird ums genauer, je grösser N

W = Die Menge aller möglichen Ausgänge W = {0,1}

X = Anzahl Erfolge

X =  $x_i$  Verteilung von X

P( $x = x_i$ ) = Wahrscheinlichkeit

E(X) = Durchschnittliche Anzahl Erfolge, wenn der Bernoulliversuch sehr

oft durchgeführt wird

= n \* p

man berechne die Wahrscheinlichkeit dafür, dass bei 120 Würfeln mit einem Laplace-Würfel 19 oder 20 oder 21 mal die sechs erscheint:

$$n * p * (1 - p) > 9$$

**1.Methode:**

$$p(a \leq X \leq b) = \sum_{\text{Anfang}}^{\text{Ende}} \binom{n}{k} * p^k * (1-p)^{n-k} = \sum_{n=19}^{21} \binom{120}{k} * \left(\frac{1}{6}\right)^k * \left(\frac{5}{6}\right)^{120-k}$$

**2.Methode :**

$$\mu = n * p = 120 * \frac{1}{6} = 20$$

$$\sigma = \sqrt{n * p * (1-p)} = \sqrt{120 * \frac{1}{6} * \frac{5}{6}} = 4.0824$$

$$p(a \leq X \leq b) = \frac{1}{\sigma * \sqrt{2 * \pi}} * \int_{18.5}^{21.5} e^{-\frac{(x-\mu)^2}{2 * \sigma^2}} * dx$$

**Boissonverteilung** nur, wenn n\* gleicher Versuch mit gleicher Ausgangslage (Zurücklegung).

**Anzahl Pfade**  $\binom{4}{2}$  4\*gezogen und 2\*Ass

$$Y = X * \beta + \epsilon$$

Die binominalverteilte Zufallsvariable X mit

$$n * p * (1-p) > 9 :$$

$$\mu = n * p$$

und

$$\sigma = \sqrt{n * p * (1-p)}$$

ist für den Fall, wo  $n * p * (1-p) > 9$  ist, angenähert normalverteilt.

Eine Zufallsvariable X heisst normalverteilt, wenn sie jede reelle Zahl als Wert annehmen kann und für jedes Intervall [a,b] gilt:

$$p(a \leq X \leq b) = \frac{1}{\sigma * \sqrt{2 * \pi}} * \int_a^b e^{-\frac{(x-\mu)^2}{2 * \sigma^2}} * dx$$

für eine normalverteilte Zufallsvariable X gilt:

$$E(X) = \mu = n * p \quad V(X) = \sigma^2$$

$$n * p * (1-p) \leq 9 :$$

**Faustregel:** ( $n > 10$  und  $p < 0.05$ ) oder ( $n \geq 100$  und  $p \leq 0.1$ )

Die Poissonverteilung ist, wenn  $n * p * (1-p) \leq 9$ ,

Die Wahrscheinlichkeit sehr klein ist, aber die Anzahl der betrachteten Ereignisse sehr gross. Man spricht von "**seltenen Ereignissen**".

Die binominalverteilte Zufallsvariable X mit  $E(X) = \mu = n * p$  gilt:

$$V(X) = \mu$$

$p(X) = k$  strebt mit wachsendem n einer festen Zahl zu:

$$p = \lim_{x \rightarrow \infty} \binom{n}{k} * p^k * (1-p)^{n-k} = \frac{\mu^k}{k!} * e^{-\mu}$$

X = k		0	1	2	3	...	k	...	∞
-----		-	-	-	-----	-	-	-	-
p(X = k)					$\frac{\mu^k}{k!} * e^{-\mu}$				

→ Verteilung

Bsp.

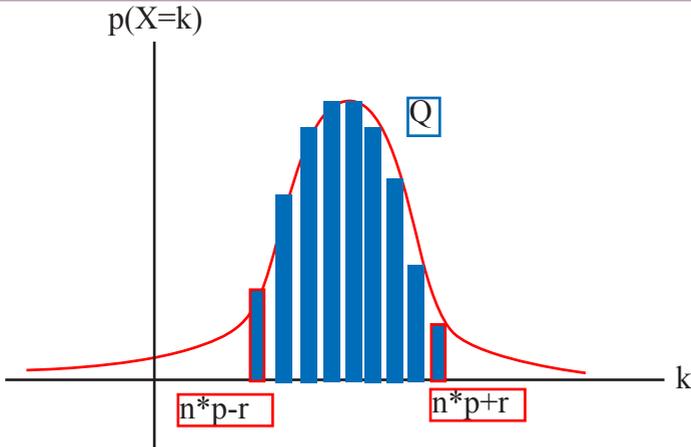
Eine fabrik produziert Schrauben, die mit einer Wahrscheinlichkeit von  $p = 0.005$  defekt sind. Wie gross ist die Wahrscheinlichkeit, in einer Schachtel von 2000 Schrauben 10 oder mehr defekte zu finden?

X = Anzahl defekte Schrauben aus 2000 - Stufigem Bernoulli Versuch

$$m = n * p = 2000 * 0.005$$

$$p(X \geq 10) \text{ mind. 10 defekte} = 1 - p(X < 10)$$

$$= 1 - \sum_{k=0}^9 \frac{m^k}{k!} * e^{-m}$$



**Ungleichungen :**

**Nie beidseitiges Dividieren oder multiplizieren mit negativen Zahlen!**

$$-b < a < b \rightarrow a^2 < b^2$$

$$a * p^2 + b * p + c \leq 0$$

**Rechung :**

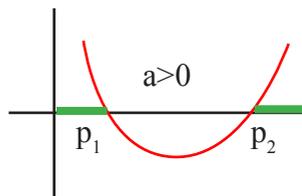
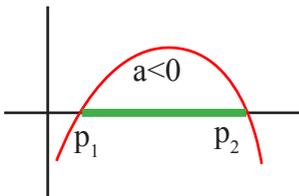
$$a * p^2 + b * p + c = 0$$

ergibt  $p_1, p_2$

$$n * p - 2 * \sqrt{n * p * (1-p)} \leq X \leq n * p + 2 * \sqrt{n * p * (1-p)} \quad | *n \quad | -p$$

$$-\frac{2}{n} * \sqrt{n * p * (1-p)} \leq \frac{X}{n} - p \leq 2 * \sqrt{\frac{p * (1-p)}{n}} \quad | \wedge^2$$

$$\left(\frac{X}{n} - p\right)^2 \leq \frac{4 * n * p * (1-p)}{n^2}$$



## Vertrauenswahrscheinlichkeit, Erfahrungswahrscheinlichkeit:

Die  $X$  - Werte liegen also in einem Intervall  $[n \cdot p - r, n \cdot p + r]$  für einen geeigneten Radius  $r$ . Man sucht nun ein  $r$  so, dass die Wahrscheinlichkeit dafür, dass  $X$  im Bereich  $n \cdot p - r \leq X \leq n \cdot p + r$  liegt, gerade  $Q$  ist.

Algebraisch:

$$P(m - r \leq X \leq m + r) = Q$$

$$E(X) = n \cdot p$$

$$V(X) = n \cdot p \cdot (1 - p)$$

**Achtung:**

Fonfidenzintervall  $\frac{X}{n} = \text{Prozent}$

$$p_1 \leq p \leq p_2$$

**n geeignet? e. ungünstig?**

1. Konfidentintervall bestimmen

2. ungünstigste Fälle von  $p_2$   $n$  bestimmen  $\rightarrow$  ev.  $n$  zu klein

Wir lassen nicht jede Zahl für  $Q$  zu, sonder nur:

$$P(m - 2 \cdot s \leq X \leq m + 2 \cdot s) = F\left(\frac{(n \cdot p + 2 \cdot s) - n \cdot p}{s}\right) - F\left(\frac{(n \cdot p - 2 \cdot s) - n \cdot p}{s}\right) = 0.9544 \text{ aus Liste}$$

$$P(m - 3 \cdot s \leq X \leq m + 3 \cdot s) = F\left(\frac{(n \cdot p + 3 \cdot s) - n \cdot p}{s}\right) - F\left(\frac{(n \cdot p - 3 \cdot s) - n \cdot p}{s}\right) = 0.9974 \text{ aus Liste}$$

$X$  und  $n$  ist bekannt ergibt nach lösen der Ungleichung:

$$\left(\frac{X}{n} - p\right)^2 \leq m \cdot \frac{p \cdot (1 - p)}{n} \quad \text{dabei ist } m = 4 \text{ bei } Q = 0.954 \quad m = 9 \text{ bei } Q = 0.997$$

$$\left(\frac{X}{n} - p\right)^2 - m \cdot \frac{p \cdot (1 - p)}{n} = 0 \quad \text{ergibt} \quad p_1 = 0.0029 \quad p_2 = 0.0428$$

bsp.

Von 1000 werkstücken einer Sendung erwiesen sich 30 als defekt.

-(Rechnung mit  $n = 1000$  und  $X = 30$ )

$$p_1 = 0.0029 \quad p_2 = 0.0428$$

**Achtung, um sicher zu sein Radius vergrößern!**  $\rightarrow p_1 = 0.002, p_2 = 0.043$

zu **95.4%** kann man sicher sein, dass zwischen 2% und 4.3% der Werkstücke kapputt sind!

$$0.002 \leq p \leq 0.043$$

**Allgemein:**

Eine Stichprobe sei entstanden durch einen  $n$  - stufigen Bernoulliversuch mit unbekannter Erfolgswahrscheinlichkeit  $p$  und  $X$  sei die Anzahl der erreichten Erfolge.

Ist  $n \cdot p \cdot (1 - p) > 9$ , so bilden alle jene  $p$ , welche die Ungleichung

$$\left(\frac{X}{n} - p\right)^2 \leq m \cdot \frac{p \cdot (1 - p)}{n} \text{ erfüllen,}$$

das Vertrauensintervall für  $p$  zur Vertrauenswahrscheinlichkeit  $Q$  oder kurz

$Q\%$  - Vertrauensintervall oder **Konfidenzintervall**.

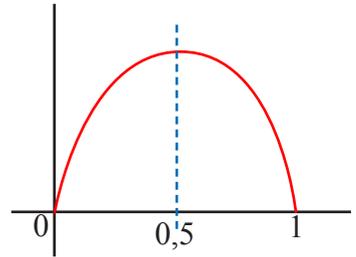
**Der notwendige Stichprobenumfang:**

"Genauigkeit von 3% mit einer Vertrauenswahrscheinlichkeit von 95,4%"

$$\left| \frac{X}{n} - p \right| \leq 0.03 \quad \left( \frac{X}{n} - p \right)^2 \leq 0.03^2$$

$$\left( \frac{X}{n} - p \right)^2 \leq \frac{4 * p * (1-p)}{n} \leq 0.03^2$$

$$\frac{4 * p * (1-p)}{0.03^2} \leq n \quad n(p) \quad \text{mit} \quad 0 \leq p \leq 1$$



$$n(p) = \frac{4 * p * (1-p)}{0.03^2} = \text{Parabel} \quad \text{"schlechtester Fall" bei } 0,5 \text{ weil}$$

$$\frac{4 * 0,5 * (1-0,5)}{0,03^2} \approx 1111,11$$

Für eine Genauigkeit von 3% und einer Vertrauenswahrscheinlichkeit von 95,4% sind mindestens 1112 Befragungen nötig.

**Regression**

**univariate Statistik = nur ein Merkmal wird untersucht**

**bivariate Statistik = zwei Merkmale werden untersucht**

**multivariate Statistik = mehr als 2 Merkmale werden untersucht**

**Abweichung**  $\varepsilon_i = y_i - F(x_i) \quad y_i = \beta_0 + \beta_1 * x_i + \varepsilon_i$

$$SSE = \sum_{i=1}^n (y_i - f(x_i))^2$$

**in Matrizenform:**

$$\begin{pmatrix} 1 & x_i \end{pmatrix} * \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} + \varepsilon_i = 1 * \beta_0 + x_i * \beta_1 + \varepsilon_i$$

$$SSE = \sum_{i=1}^n y_i - (\beta_0 + x_i * \beta_1)^2 = \sum_{i=1}^n \varepsilon_i^2 \quad |\varepsilon|^2 \text{ muss minimal sein!!!}$$

$$\frac{\partial SSE}{\partial \beta_0} = -2 * \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 * x_i)) = 0 \quad |:-2$$

$$\sum_{i=1}^n y_i = n * \beta_0 + \beta_1 * \sum_{i=1}^n x_i$$

$$\frac{\partial SSE}{\partial \beta_1} = -2 * \sum_{i=1}^n x_i * (y_i - (\beta_0 + \beta_1 * x_i)) = 0 \quad |:-2$$

$$\sum_{i=1}^n y_i * x_i = \beta_0 * \sum_{i=1}^n x_i + \beta_1 * \sum_{i=1}^n x_i^2$$

Mit Matrizen dargestellt :

$$\text{linke Seite : } \begin{pmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i * y_i \end{pmatrix} = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{pmatrix} * \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix} \quad \mathbf{X^T} \quad \mathbf{Y}$$

$$\begin{aligned} \text{rechte Seite : } & \begin{pmatrix} n * b_0 + b_1 * \sum_{i=1}^n x_i \\ b_0 * \sum_{i=1}^n x_i + b_1 * \sum_{i=1}^n x_i^2 \end{pmatrix} = \begin{pmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{pmatrix} * \begin{pmatrix} b_0 \\ b_1 \end{pmatrix} \\ & = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{pmatrix} * \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \dots & \dots \\ 1 & x_n \end{pmatrix} * \begin{pmatrix} b_0 \\ b_1 \end{pmatrix} \quad \mathbf{X^T} \quad \mathbf{X} \quad \mathbf{b} \end{aligned}$$

$$\rightarrow \mathbf{X^T * Y = X^T * X * b}$$

$$\mathbf{b = (X^T * X)^{-1} * X^T * Y}$$

TI-89 :

$\{x_1, x_2, \dots, x_n\}$  STO A

$\{y_1, y_2, \dots, y_n\}$  STO B

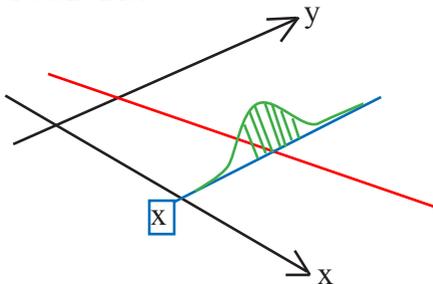
LinReg A, B

QuadReg A, B

ExpReg A, B

danach ShowStat

Damit das Werkzeug "lineare Regression" sinnvoll angewendet werden kann auf Datenpunkte  $(x,y)$ , sollte das Residuum  $\varepsilon = y - (\beta_1 + \beta_2 * x)$  and der Stelle  $x$  normalverteilt sein mit Mittelwert 0 und einer Varianz  $\sigma^2$ , welche unabhängig von der Stelle  $x$  ist.



**Regression bei multivariaten Daten :**

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix} \quad b = \begin{pmatrix} b_1 \\ b_2 \\ \dots \\ b_p \end{pmatrix} \quad X = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \dots & \dots & \dots & \dots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix}$$

Achtung : bei  $x^2$   $X = \begin{pmatrix} 1 & x_1 & x_1^2 \\ \dots & \dots & \dots \\ 1 & x_n & x_n^2 \end{pmatrix}$  bei  $\frac{1}{x}$   $X = \begin{pmatrix} 1 & x_1 & \frac{1}{x_1} \\ \dots & \dots & \dots \\ 1 & x_n & \frac{1}{x_n} \end{pmatrix}$

$$b = (X^T * X)^{-1} * X^T * Y$$

**Die Funktion kann nicht linear sein, die Regression ist es aber!**

**Trick :**

Daten  $(x_i/y_i)$  Modell  $y = a * e^{b*x}$  (Funktion aussuchen, die etwa passt)

$$y = a * e^{b*x} \quad | \ln(\ )$$

$$\ln(y) = \ln(a) + b * x \quad \rightarrow \text{ist linear}$$

$$\ln(y) = a^* + b^* * x$$

Rechne lineare Regression für transformierte Daten

Rücktransformation :  $a = e^{a^*}$  ,  $b = b^*$

**Beispiel :**

$$n = 2$$

$$y_1 = 1, \quad y_2 = 3$$

$$\hat{y}_1 = -1, \quad \hat{y}_2 = -3$$

$$\bar{y} = 2$$

$$\bar{\hat{y}} = -2$$

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{(-1+2)^2 + (-3+2)^2}{(1-2)^2 + (3-2)^2} = \frac{2}{2} = 1$$

Nur aus der Formel für  $R^2$  kann nicht geschlossen werden:

$$(R^2 = 1) \neq (y_i \approx \hat{y}_i)$$

**Wie gut ist der Fit??**

Als Mass für den Fehler :

$$SS_{\text{tot}} = \sum_{i=1}^n (y_i - \bar{y})^2 \quad \text{ist bis auf den Vorfaktor } \frac{1}{n-1} \text{ de Varianz!}$$

**Durch x verursachte Varianz von y** oder **erklärte Varianz** von y :

$$\frac{1}{n-1} * \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

Es seien Datenpunkte  $(x_{i1}, x_{i2}, \dots, x_{ip}, y_i)$  gegeben mit  $i = 1, 2, \dots, n$

für eine Regression  $\hat{y}_i = f(x_{i1}, x_{i2}, \dots, x_{ip})$  durchgeführt worden, dann

heisst der Term :

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{\text{erklärte Varianz}}{SS_{\text{tot}}} = \text{Bestimmtheitsmass oder Determinationskoeffizient}$$

der Regression.

$\bar{y}$  = empirischer Mittelwert der  $y_i$

$\hat{y}$  = empirischer Mittelwert der  $\hat{y}_i$

je besser die Näherung  $y_i \gg \hat{y}_i$ , deso näher liegt  $R^2$  bei 1.

Es seien Datenpunkte  $(x_i, y_i)$  gegeben mit  $i = 1, 2, \dots, n$  für die eine Regression

$\hat{y}_i = b_0 + b_1 * x_i$  durchgeführt worden ist, dann gilt :

$$\sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \bar{y})^2 - SSE$$

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

$$SSE = (1 - R^2) * \sum_{i=1}^n (y_i - \bar{y})^2$$

Die Qualität des Fits ist umso besser, je kleiner der SSE.

$1 - R^2$  macht den SSE klein, wenn  $R^2$  nahe bei 1

Mehr noch :

$$1 - R^2 \geq 0 \quad \text{da} \quad SSE \geq 0 \quad \text{und} \quad SS_{\text{tot}} \geq 0 \quad \quad 0 \leq R^2 \leq 1$$

**Korrelationskoeffizient :**

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

$r^2$  = Bestimmtheitsmass

r ist positiv, wenn die Regressionsgerade steigend ist

r ist negativ, wenn die Regressionsgerade fallend ist