

i.i.d. = independent and identically distributed

Wahrscheinlichkeit (Wa., bsp. fairer Würfel)

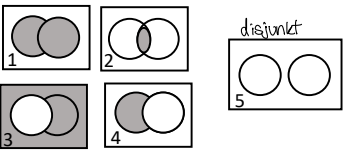
Ω	Grundraum	{1,2,3,4,5,6}
ω	Elementarereignis	{6}
A	Ereignis ($A \in \Omega$)	{2,4,6}
ϕ	Leere Menge	{}

Rechenoperationen

$P(A)$	≥ 0	Wa. von A
$P(\Omega)$	$= 1$	gesamte Wa.
$P(A \cup B)$	$P(A) + P(B) - P(A \cap B)$	1 oder/und B
$P(A \cap B)$	$= P(A B) \cdot P(B)$	[2] A und B (geschnitten)
$P(A^c)$	$1 - P(A)$	[3] nicht A
$P(A B)$	$P(A) - P(A \cap B)$	[4] A ohne B

Disjunkt [5]

$P(A \cap B) = 0$ (Schnittmenge gleich leere Menge)
 $P(A|B) = 0$
 $P(A \cup B) = P(A) + P(B)$ (Wa. addieren)



Laplace

Jedes Ereignis ist gleich wahrscheinlich!

$$P(E) = \frac{g}{m} \left(\frac{\# \text{ gesuchte Möglichkeiten}}{\# \text{ alle Möglichkeiten}} \right)$$

Bsp.: fairer Würfel
 Wa. eine 6 zu würfeln: 1/6
 Wa. eine gerade Zahl zu würfeln: 3/6

Unabhängigkeit

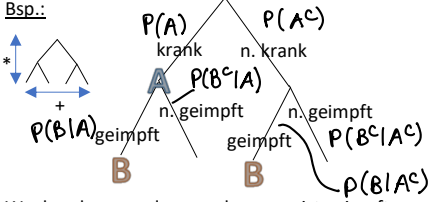
Die Wahrscheinlichkeiten beeinflussen sich gegenseitig nicht! (Gleichmässiger Baum)
Bsp.: Wa. zweimal hintereinander eine 6 zu würfeln
 $(1/6 * 1/6)$
 $P(A \cap B) = P(A) * P(B)$
 $P(A) = P(A|B)$

Abhängigkeit/ bedingte Wahrscheinlichkeit

Die Wahrscheinlichkeiten beeinflussen sich gegenseitig! (Ungleichmässiger Baum)
Bsp.: Wa. krank zu werden ist abhängig von der Impfung.
 $P(A|B) = P(A \cap B) / P(B)$

Formel für Wa., dass A eintritt, wenn B schon eingetroffen ist. (A gegeben B)

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



Wa. krank zu werden gegeben man ist geimpft
 $P(A|B) = P(A \cap B) / P(B)$; A=krank, B=geimpft;
 $P(B) = P(\text{krank}) * P(\text{geimpft}) + P(\text{n. krank}) * P(\text{geimpft})$
 $P(A \cap B) = P(\text{krank}) * P(\text{geimpft})$

Satz von Bayes

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

$P(A) = 0 \rightarrow P(A|B) = 0$

$P(A) > 0 \rightarrow P(A|B) = [0; 1] \rightarrow P(A|B) \neq P(A)$

TR: $S_X = \sigma$ $odds(A^c) = \frac{1}{odds(A)}$

Odds

Wie viel Mal wahrscheinlicher ist A als nicht A?
(Bsp.: Münzwurf 1:1 Wa. Kopf zu werfen; fairer Würfel 1:5 Wa. eine 6 zu würfeln) $P(A^c) = 1 - P(A)$

$$odds(A) = \frac{P(A)}{P(A^c)} = \frac{P(A)}{1 - P(A)}$$

$$odds(A|B = 1) = \frac{P(A|B = 1)}{1 - P(A|B = 1)}$$

Log-Odds

$$\ln(odds(A)) = \ln\left(\frac{P(A)}{P(A^c)}\right) = \ln\left(\frac{P(A)}{1 - P(A)}\right)$$

Odds-Ratio (OR)

$\rightarrow A > A^c \Rightarrow \ln odds > 0$

Wie stark hängen zwei Merkmale zusammen?
 (Für absolute Wa.s oder für bedingte Wa.s)

$$OR = \frac{odds(A)}{odds(B)}; OR = \frac{odds(A|B)}{odds(A|B^c)}$$

Bsp.: Wa. zu Erkranken sinkt um x wenn man geimpft ist. $OR = x = \frac{odds(\text{krank}|\text{nicht geimpft})}{odds(\text{krank}|\text{geimpft})}$

Stichprobengröße:
 $n \geq 4 \cdot \frac{\hat{\sigma}^2}{\sigma^2}$
 \rightarrow Breite 95% -VI
 $= 2 \cdot \sigma$

Macht: Wa, dass Beob. in Verwerfungsbereich fällt, obwohl H_0 wahr ist
 \rightarrow Verwerfungsbereich $\{0, \dots, 15\} \rightarrow$ Wahre Wa. = P_A
 = Wah. H_0 zu verwerfen $P_{PA}(X \leq 15)$

\rightarrow kann nur berechnet werden, wenn genaue Verteilung der Teststatistik unter H_0 bekannt

- Fehler 1. Art $\downarrow \Rightarrow$ Fehler 2. Art $\uparrow \Rightarrow$ Macht \downarrow
- höheres $\alpha =$ größerer Verwerfungsbereich

Auf 10% verwerfen \Rightarrow nicht unbedingt auf 1% vom

DISKRETE VERTEILUNG (abzählbar, mit „Punkt“-Wahrscheinlichkeiten rechnen möglich; $P(X=x)$, PDF)
Kennzahlen einer Verteilung

X: Zufallsvariable, x: konkreter Wert von X

Erwartungswert $\epsilon(x)$: Beschreibt die mittlere Lage der Verteilung (Mittelwert bzw. Durchschnitt)

$$\epsilon(x) = \sum_{x \in W} x * P(X = x); W = \text{Wertebereich von } x$$

Varianz $\text{Var}(x)$: Beschreibt die mittlere quadratische Abweichung eines Zufallwertes von ihrem Erwartungswert.

$$\text{Var}(x) = \sum_{x \in W} (x - \epsilon(x))^2 * P(X = x)$$

Standardabweichung $\sigma(x)$: Beschreibt die Streuung der Verteilung bzw. die Breite

$$\sigma(x) = \sqrt{\text{Var}(x)}$$

Zufallsvariablen

Rechenregeln

$$\epsilon(X+Y) = \epsilon(X) + \epsilon(Y)$$

Varianz:

$$\epsilon(a * X) = a * \epsilon(X)$$

1. Konstanten=0

$$\epsilon(a+bX) = a + b\epsilon(X)$$

2. Faktoren hoch 2

$$\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$$

3. Minus wird Plus

$$\text{Var}(a * X) = a^2 * \text{Var}(X)$$

$$\text{Var}(2X+9 * Y) = 4\text{Var}(X) + \text{Var}(Y)$$

$$\text{Var}(a+bX) = b^2 * \text{Var}(X)$$

$$\sigma(X \pm Y) = \sqrt{\sigma^2(X) + \sigma^2(Y)}$$
 falls x und y unabhängig

$$\sigma(X \pm Y) = \sqrt{\sigma^2(X) + \sigma^2(Y) \pm 2\text{Cov}(X, Y)}$$
 falls abhängig

$$\text{Cov}(X, Y) = \epsilon((X - \epsilon(X))(Y - \epsilon(Y)))$$

$$\sigma(a+bX) = b * \sigma(X)$$

$$q(a+b \cdot x) = a + b \cdot q(x)$$

Berechnen von Vertrauensintervall für π

für Binomialverteilung durch Normalapproximation

$$I = \frac{x}{n} \pm \Phi^{-1} \left(1 - \frac{\alpha}{2} \right) * \sqrt{\frac{x}{n} \left(1 - \frac{x}{n} \right) \frac{1}{n}}$$

für $\alpha = 0.05$: $I = \frac{x}{n} \pm 1.96 * \sqrt{\frac{x}{n} \left(1 - \frac{x}{n} \right) \frac{1}{n}}$ \rightarrow Breite VI abhängig von $\frac{1}{\sqrt{n}}$

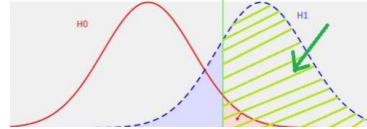
VI für gesamte erwartete Anzahl von Erfolgen:

$$I = n\pi_0 \pm \Phi^{-1} \left(1 - \frac{\alpha}{2} \right) * \sqrt{n\pi_0(1 - \pi_0)}$$

für $\alpha = 0.05$: $I = n\pi_0 \pm 1.96 * \sqrt{n\pi_0(1 - \pi_0)}$

\rightarrow VI umfasst alle Werte π_0 , für die ein (zweiseitiger) Binomialtest mit $H_0: \pi = \pi_0$ nicht verwirft

Fehler
 (hier: $\pi A > \pi_0$; grüne Linie: Signifikanzlevel α , definiert Verwerfungsbereich K)



\rightarrow größere Erfolgswa. unter H_0 :

Verteilung verschiebt sich nach rechts
 \rightarrow Verwerfungsbereich wird potentiell kleiner
 \rightarrow Macht nimmt ab

Rot: Fehler 1. Art ($\leq \alpha$); H_0 wird verworfen obwohl sie stimmen würde, $P(X \in K)$

Blau: Fehler 2. Art ($\leq \beta$); H_0 wird akzeptiert obwohl sie falsch ist, $P(X \notin K) \rightarrow$ (Wa, dass mit $p(x)$ Wert im VI liegt)

Grün: Macht ($\leq 1 - \beta$); Wahrscheinlichkeit eine bestimmte Alternativhypothese zu erkennen, wenn diese stimmt, $P(X \in K)_{\pi A}$, kann man nur berechnen wenn πA bekannt! Je grösser der Fehler 1. Art, desto grösser die Macht.

Macht einseitig vs. zweiseitig:

$H_0: \pi < \pi_0 \Rightarrow$ wird kleinere Erfolgswa. π_0 mit größerer Macht erkennen als zweiseitiger, für größere Erfolgswa. aber quasi keine Macht

$H_0: \pi > \pi_0 \Rightarrow$ andersrum

\Rightarrow ob Macht (einseitig) > Macht (zweiseitig) hängt von konkreter Alternative ab

P-Wert:

Wa, unter gültiger H_0 das Ergebnis oder etwas noch extremeres (im Sinne der Alternative) zu erhalten

• verwerfe H_0 falls $p \leq \alpha$

• Gewinnwa. π_0 unter Nullhypothese ist nicht im x VI für Gewinnwa. $\rightarrow p < 1 - x$

Einseitiger Test:

eine Seite blind
 man sucht nur grösser als oder kleiner als
Grosse Macht

Zweiseitiger Test:

sieht beide Seiten
 man sucht grösser als und kleiner als
Kleine Macht

Grundfragen der Statistik

1. **Plausibler Parameter π ?**

Punktschätzung für Erfolgswahrscheinlichkeit π

a) Momentenmethode

$$\pi = \frac{x}{n} = \frac{\epsilon(x)}{n} = \frac{\# \text{beobachtet}}{\# \text{unabhängige Versuche}} \cdot \text{erwartet}$$

b) Maximum likelihood Methode

Gleichung aufstellen ($P(A) * P(B) * \dots$), logarithmieren falls hilfreich, ableiten und gleich null setzen, nach π auflösen; genauer, aufwendiger; nur wissen, dass es das gibt

2. **Sind Werte von π mit den Daten vereinbar?**

Vertrauensintervall VI

\rightarrow Die Werte von π_0 bei denen H_0 nicht verworfen wird bilden ein $(1 - \alpha)$ -Vertrauensintervall.

\rightarrow Ein $(1 - \alpha)$ -VI enthält den wahren Wert mit einer Wa. von $(1 - \alpha)$.

Nur falls n gross ist, kann eine Normalapproximation gemacht werden! **Faustregel:** $n\pi > 5$ und $n(1 - \pi) > 5$

95% VI ist im 99% VI enthalten

Gewinnwa. ist in $1 - \alpha - \text{VI}_2$
 \rightarrow auf α nicht verwerfen
 $p > \alpha$

$X \sim \text{Bin}(n, \pi) \quad Y \sim \text{Bin}(l, \pi) \Rightarrow$ unabhängig (+gleichverteilung)
 $\Rightarrow Y+X = \text{Bin}(l+n, \pi)$

Bernoulli nicht miteinander addieren $\Rightarrow n > 1$ nicht mehr Ber

Summe von zwei unabhängigen Poissonvert. ist Poissonvert.
 $X \sim \text{Pois}(a) \quad Y \sim \text{Pois}(b) \quad Z = X+Y \quad Z \sim \text{Pois}(a+b)$

Binomialverteilung $X \sim \text{Bin}(n, \pi)$

Anzahl Erfolge (x) bei n unabhängigen Versuchen mit Erfolgswahrscheinlichkeit π .

$P(X = x) = \binom{n}{x} \pi^x (1 - \pi)^{n-x}$

$E(x) = n\pi$

$\sigma^2(x) = \text{Var}(x) = n\pi(1-\pi) \rightarrow$ am größten, wenn $\pi = 0.5$

Hypothesentest Binomialverteilung

1. Modell: X = # Erfolge von n Versuchen

$X \sim \text{Bin}(n, \pi)$

2. Nullhypothese

$H_0: \pi_0 = ?$

$H_A: \pi_A \neq \pi_0, \pi_A > \pi_0, \pi_A < \pi_0$

Einseitiger oder zweiseitiger Test?

3. Teststatistik

T = # Erfolge von n Versuchen unter der Annahme, dass H_0 stimmt.

$T \sim \text{Bin}(n, \pi_0)$

4. Signifikanzniveau

$\alpha = ?$ (meist 5% $\rightarrow \alpha = 0.05$) $\rightarrow P(x \leq k) \leq 0.05$

5. Verwerfungsbereich

$\pi_A \neq \pi_0 \quad P(T \leq C_U) \leq \alpha/2 \text{ und } P(T \geq C_O) \leq \alpha/2$

$\pi_A > \pi_0 \quad P(T \geq C_O) \leq \alpha$

$\pi_A < \pi_0 \quad P(T \leq C_U) \leq \alpha$

Normalapproximation von C_U/C_O :

Zweiseitig, $\alpha = 0.05$

$C \approx n\pi_0 \pm 1.96 \sqrt{n\pi_0(1-\pi_0)}$

Einseitig, $\alpha = 0.05: \pi_A > \pi_0: +1.64 \rightarrow$ Obergrenze n

$\pi_A < \pi_0: -1.64 \rightarrow$ Untergrenze 0

6. Testentscheid

Liegt der beobachtete Wert im Verwerfungsbereich?

JA $\rightarrow H_0$ wird verworfen

NEIN $\rightarrow H_0$ wird nicht verworfen (aber auch nicht bewiesen)

Bernoulliverteilung $Y \sim \text{Ber}(\pi)$

Erfolg oder Misserfolg bei nur einem Versuch? (Spez. Fall der Binom. mit $n=1$)

$P(X = x) = \binom{1}{x} \pi^x (1 - \pi)^{1-x}$

$P(X = 1) = \pi$

$P(X = 0) = 1 - \pi$

$E(x) = \pi$

$\sigma^2(x) = \text{Var}(x) = \pi(1-\pi)$

Uniforme Verteilung $Y \sim \text{unif}(n)$

Alle Ereignisse sind gleich wahrscheinlich, n ist die Anzahl möglicher Ereignisse.

$P(X = x) = \frac{1}{n}$

$E(x) = \frac{(n+1)}{2}$

$\sigma^2(x) = \text{Var}(x) = \frac{(n+1)(n-1)}{12}$

Hypergeometrische Verteilung $Y \sim \text{hyper}(n, N, m)$

Modell: Urne mit N Kugeln, m sind markiert; es werden n gezogen und x davon sind markiert.

Ohne zurücklegen!!

$P(X = x) = \frac{\binom{m}{x} \binom{N-m}{n-x}}{\binom{N}{n}} = \begin{matrix} \text{günstig} \\ \text{möglich} \end{matrix}$

$E(x) = \frac{n \cdot m}{N}$

$\sigma^2(x) = \text{Var}(x) = \frac{n \cdot m \cdot (N-n)(N-m)}{N^2 \cdot (N-1)}$

Poissonverteilung $Y \sim \text{Pois}(\lambda)$

Seltene Ereignisse und ihr vorkommen in einem bestimmten Zeitraum

λ : durchschnittliches Auftreten pro Zeit

$P(X = x) = e^{-\lambda} \cdot \frac{\lambda^x}{x!}$

$E(x) = \lambda$

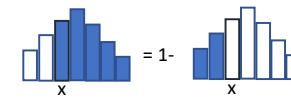
$\sigma^2(x) = \text{Var}(x) = \lambda$

Summen von Poissonverteilten Zufallsvariablen: Wenn $X \sim \text{Pois}(\lambda_x)$ und $Y \sim \text{Pois}(\lambda_y)$ unabhängig sind, dann ist $X+Y \sim \text{Pois}(\lambda_x + \lambda_y)$

Prinzip „bis und mit“ = „1-ohne“:

$P(X \leq x) = 1 - P(X > x) = 1 - P(X \geq x+1)$

$P(X \geq x) = 1 - P(X < x) = 1 - P(X \leq x-1)$



Für große n kann man Binomial mit Poisson annähern: $\lambda = n \cdot \pi$

Maximum likelihood $\pi = \frac{x}{n}$

$H_A: p > p_0 \Rightarrow$ Für Verwerfungsbereich $P(x \geq k) \leq \alpha$

Verteilungen von diskreten Zufallsvariablen

In Versicherungen gängige Vert.
Ereignisse die beliebig oft eintreten können

	Binomialverteilung	Poissonverteilung	Allgemeine Formeln
Beschreibung	X = Anzahl „Erfolge/Treffer“ in n unabhängigen Versuchen (p = Erfolgswahrscheinlichkeit)	X = Anzahl „Erfolge“ pro Intervall (Zeit, Fläche usw.),	
Darstellung	$X \sim \text{Bin}(n, p)$	$X \sim \text{Poisson}(\lambda)$	
Werte, die X annehmen kann (diskret -> abzählbar)	<ul style="list-style-type: none"> abzählbare, endliche Anzahl Ereignisse 0, 1, 2, 3, 4 ..., n 	<ul style="list-style-type: none"> abzählbare, unendliche Anzahl Ereignisse 0, 1, 2, 3, ... 	<ul style="list-style-type: none"> abzählbar ..., -2, -1, 0, 1, 2, ...
Wahrscheinlichkeitsfunktion f(x)	$f(x) = P[X = x] = \binom{n}{x} \cdot p^x \cdot (1-p)^{n-x}$	$f(x) = P[X = x] = e^{-\lambda} \cdot \frac{\lambda^x}{x!}$	Selber herausfinden, ohne Formel... z.B. bei Jasskarten: $P[\text{König}] = P[X=4] = \frac{4}{36} = \frac{1}{9}$
Kumulative Verteilungsfunktion F(x)	$F(x) = P[X \leq x] = \sum_{k=0}^x \binom{n}{k} \cdot p^k \cdot (1-p)^{n-k}$ <p>Wichtig: $F(n) = P[X \leq n] = 1$</p>	$F(x) = P[X \leq x] = \sum_{k=0}^x e^{-\lambda} \cdot \frac{\lambda^k}{k!}$	$F(x) = P[X \leq x] = \sum_{k \leq x} P[X = k]$
Erwartungswert E(X)	$E(X) = n \cdot p$	$E(X) = \lambda$	$E(X) = \sum_{i=1}^n k_i \cdot P[X = k_i]$
Varianz Var(X)	$\text{Var}(X) = n \cdot p \cdot (1-p)$	$\text{Var}(X) = \lambda$	$\text{Var}(X) = \sum_{i=1}^n (k_i - E(X))^2 \cdot P[X = k_i]$
Standardabweichung $\sigma(X)$	$\sigma(X) = \sqrt{\text{Var}(X)}$	$\sigma(X) = \sqrt{\text{Var}(X)}$	$\sigma(X) = \sqrt{\text{Var}(X)}$
Approximatives zweiseitiges 95% Vertrauensintervall	$I \approx \frac{x}{n} \pm 1.96 \cdot \sqrt{\frac{\frac{x}{n} \cdot (1 - \frac{x}{n})}{n}}$ <p>Vertrauensintervall für die Erfolgswahrkeit p.</p>	$I \approx x \pm 1.96\sqrt{x}$	
Approx. einseitiges 95% Vertrauensintervall	$\left[-\infty, \frac{x}{n} + 1.64 \sqrt{\frac{\frac{x}{n} \cdot (1 - \frac{x}{n})}{n}} \right] \text{ für } H_A: p < p_0$ $\left[\frac{x}{n} - 1.64 \sqrt{\frac{\frac{x}{n} \cdot (1 - \frac{x}{n})}{n}}, \infty \right] \text{ für } H_A: p > p_0$	Vertrauensintervall für λ .	

WICHTIG: Bernoulliverteilung $X \sim \text{Bernoulli}(p)$ ist ein Spezialfall der Binomialverteilung für $n=1$ (d.h. $X \sim \text{Binom}(1, p)$). Es wird nur ein einziger Versuch beschrieben, die Wahrscheinlichkeit für einen Erfolg ist gleich p und die Wahrscheinlichkeit für einen Misserfolg ist gleich 1-p.

WICHTIG: APPROXIMATIONEN: Für $n = \text{gross}$ ist es mühsam, ohne Computer die kumulative Wahrscheinlichkeit $P[X \leq x]$ der Binomialverteilung zu berechnen. Deshalb wird meistens mit anderen Verteilungen approximiert:

- Normalapproximation (für $n = \text{gross}$, sehr häufig verwendet):** $X \approx \text{Normal}(\mu, \sigma^2)$; $\mu = n \cdot p$; $\sigma^2 = n \cdot p \cdot (1-p)$

$$P[X \leq x] \approx P\left[Z \leq \frac{x-\mu}{\sigma}\right] = \Phi\left(\frac{x-\mu}{\sigma}\right) \rightarrow z\text{-Tabelle}$$
- Poissonapproximation (für $n = \text{gross}$ und $p = \text{klein}$):** $X \approx \text{Poisson}(\lambda)$; $\lambda = n \cdot p$

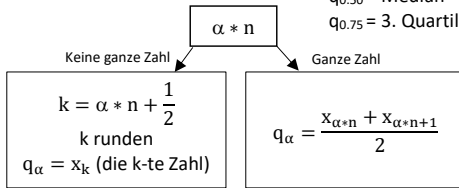
! Median ist nicht immer in Stichprobe

STETIGE VERTEILUNG (Wertebereich für x kontinuierlich, keine „Punkt“-Wahrscheinlichkeiten! $P(X=x)=0$ für alle $x \in W_x$, nur $P(X \leq x)$ möglich, Fläche unter der Kurve, Intervalle!!)

α -Quantil (empirisch)

Wert q_α , bei dem mindestens $\alpha * 100\%$ der Datenpunkte kleiner und $(1-\alpha) * 100\%$ grösser als q_α sind. **ACHTUNG!** Daten der Reihenfolge nach sortieren!
 $x_1 \leq x_2 \leq \dots \leq x_n$

$q_{0,25} = 1.$ Quartil
 $q_{0,50} =$ Median
 $q_{0,75} = 3.$ Quartil



Inter-Quartile-Range: (~68% aller Werte)

IQR = $q_{0,75} - q_{0,25}$

IQR und Median: robust gegen Ausreisser

$P(X \leq q_\alpha) = \alpha$

Kennzahlen

Für die Lage:

Arithmetisches Mittel = Durchschnitt ($=\mu$)

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

→ nicht robust gegen Ausreisser

Für die Streuung („Breite“):

Empirische Standardabweichung ($=\sigma$)

$$S_x = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2}$$

Korrelation: $y = a \cdot x \rightarrow a$ ist Steigung im Streudiagramm zw. x und y

↳ misst wie stark Punkte auf einer Linie liegen

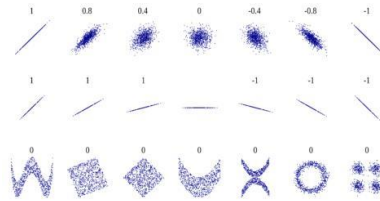
Gerade: Korrelation 1 | kein linearer Zusammenhang: Korrelation 0

Für linearen Zusammenhang: (zwischen X und Y, zwei stetigen Zufallsvariablen) $Y = a + bX$

$E(Y) = a + bE(X)$; $Var(Y) = b^2 Var(X)$; $q_y = a + bq_x$

wenn x normiert, y auch; wenn nicht, keine Aussage möglich

Korrelation [-1,1]: (KEIN kausaler Zusammenhang!!!)



Stärke und Richtung von linearen Abhängigkeiten, sagt nichts über Steigung aus!

$cor=0$; kein linearer Zusammenhang, aber evt. Sonst

Empirische Korrelation:

$$r = \frac{S_{xy}}{S_x S_y} \quad S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \frac{1}{n-1}$$

Gesetz der grossen Zahlen (GGZ)

$X_1, X_2 \dots X_n \sim F$ i.i.d (independent and identical distributed; F: irgendeine Verteilung)

$$E(\bar{x}_n) = \mu \quad \sigma_{\bar{x}_n} = \frac{\sigma_x}{\sqrt{n}} \quad Var(\bar{x}_n) \neq Var(x_i)$$

Wurzel-n-Gesetz: Für eine halb so grosse Standardabweichung (Streuung) braucht man 4mal so viele Daten! (gilt für Std.abw. von Mittelwerten, nicht einzelne Messungen)

Zentraler Grenzwertsatz (ZGWS)

$X_1, X_2 \dots X_n \sim F$ i.i.d \rightarrow aus GGZ: $\bar{x}_n \sim N(\mu, \frac{\sigma_x^2}{n})$

$S_n \sim N(n * \mu, n * \sigma_x^2)$

Bsp.: $n=100$ Spiele, $E(x)=1/3$, $Var(x)=28.6$

S_n =totale Gewinne (Standardisieren **S.6**)

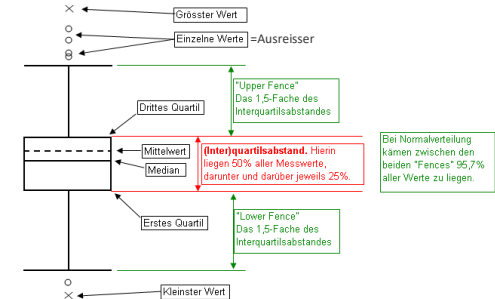
$S_n \sim N(100 * 1/3, 100 * 28.6) = N(33,2860)$

→ müssen unabhängig und gleich verteilt sein

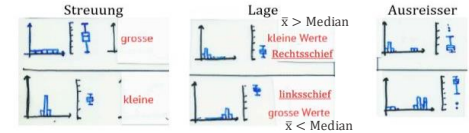
→ Für große n ist Approximation besser

Graphische Methoden

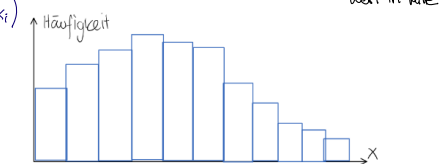
Boxplot



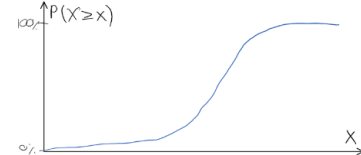
Hauptmerkmale zur Diskussion, bzw. zum Vergleich:



Histogramm



Empirische kumulative Verteilungsfunktion (ECDF)



Kontinuierliche Verteilungsfunktionen

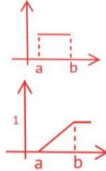
Wahrscheinlichkeitsdichte: $f(x)$

Kummulative Verteilungfnk.: $F(x)$ (Integral von $f(x)$)

Uniforme Verteilung $X \sim \text{Unif}([a,b])$

Jeder Wert ist gleich wahrscheinlich

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{falls } a \leq x \leq b \\ 0 & \text{sonst} \end{cases}$$



$$P(x \leq X) F(x) = \begin{cases} 0 & \text{falls } x < a \\ \frac{x-a}{b-a} & \text{falls } a \leq x \leq b \\ 1 & \text{falls } x > b \end{cases}$$

$$\mathcal{E}(X) = \frac{a+b}{2}, \quad \text{Var}(X) = \frac{(b-a)^2}{12}, \quad \sigma_X = \frac{b-a}{\sqrt{12}}$$

Exponentialverteilung $X \sim \text{Exp}(\lambda)$

Warten ohne Gedächtnis (Bsp.: Radioaktiver Zerfall)

- Einfaches Modell für Wartezeiten auf Ausfälle
- Wenn die Zeit zwischen den Ausfällen Exponentialverteilt sind (λ), dann ist die Anzahl Ausfälle in einem Intervall t Poisson-verteilt ($t\lambda$).

$$f(x) = \begin{cases} \lambda \exp(-\lambda x) & \text{falls } x \geq 0 \\ 0 & \text{sonst} \end{cases}$$



$$F(x) = \begin{cases} 1 - e^{-\lambda x} & \text{falls } x \geq 0 \\ 0 & \text{falls } x < 0 \end{cases}$$



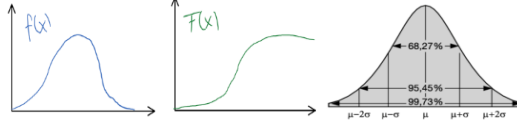
$$\mathcal{E}(X) = \frac{1}{\lambda}, \quad \text{Var}(X) = \frac{1}{\lambda^2}, \quad \sigma_X = \frac{1}{\lambda}$$

- $P(T > t) = 1 - P(T \leq t) = e^{-\lambda t}$
- $P(T > t + s | T > s) = \frac{P(T > t+s \text{ und } T > s)}{P(T > s)} = \frac{P(T > t+s)}{P(T > s)} = \frac{e^{-\lambda(t+s)}}{e^{-\lambda s}} = e^{-\lambda t} = P(T > t)$
"Es spielt keine Rolle, ob man schon s Sekunden gewartet hat"

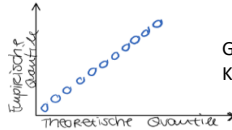
Normalverteilung $X \sim N(\mu, \sigma^2)$

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$F(x) = ??? \rightarrow$ Standardnormalverteilung



Prüfen ob Normalverteilt: QQ-Plot (S.9)



Gerade: Normalverteilt
Krumm: NICHT n.v.

Für z-Test: $\bar{X}_n = N(\mu, \sigma_{\bar{X}_n}^2) \rightarrow N(\mu, \lambda) = \frac{\bar{X}_n - \mu}{\sigma_{\bar{X}_n}}$

\Rightarrow Für t-Test nicht weil σ geschätzt ist

Streuung Werte: Standardabweichung s.d.

Streuung arith. Mittel: "s.e.d."

\rightarrow Wenn Steigung kleiner \Rightarrow kleinere Varianz

Standardnormalverteilung $Z \sim N(0,1)$

Mittelwert μ wird auf 0 gesetzt (Symmetrieachse=x-Achse)

Standardabweichung σ wird auf 1 gesetzt ($x[-1,1]$)

$$f(x) = \varphi(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$$

$$F(x) = \Phi(x) = \int_{-\infty}^x \varphi(x) dx$$

TABELLE!!

Von Normalverteilt zu Standardnormalverteilt:

Standardisieren (oder Normieren):

$$x_n = \frac{x - \mu}{\sigma} \quad \text{disnormieren: } x = x_n \sigma + \mu$$

Bsp.: $X \sim N(2,2^2) \rightarrow P(X \leq 5) = ? (=0.93)$

$$P(X \leq 5) = P\left(Z \leq \frac{5-\mu}{\sigma}\right) = P\left(Z \leq \frac{5-2}{2}\right) = P\left(Z \leq \frac{3}{2}\right)$$

$$P\left(Z \leq \frac{3}{2}\right) = \Phi(1.5) = (\text{Tabelle}) = 0.93$$

Wichtig: Quadrat!! Wir rechnen mit σ , nicht mit σ^2 ! Also zuerst Wurzelziehen!

Normalvert.: Median = Erwartungswert
 \hookrightarrow symmetrisch am Erwartwert

\rightarrow Menschl. Leben keine Ex-Vert., da Wk. zu Sterben im Alter zunimmt

Stichprobengröße: $n \geq 4 \cdot \frac{\hat{\sigma}_x^2}{\sigma^2}$

$\sigma^2 = \frac{1}{2}$ Breite VI

10 Pers. Sprint
1x mit Koffein
1x ohne
→ äußere Faktoren können variieren, intrinsische sind gleich

10 Pers. 1x Sprint
5 mit Koffein
5 ohne
→ nicht genau 1 kann 1 zugeordnet werden

Statistische Tests (für Erwartungswert μ)

Z-Test basierend auf mehreren Beobachtungen

σ_x : bekannt; $\hat{\sigma}_{\bar{x}_n}$: berechnen

X_i : kontinuierlich, normalverteilt; Gepaarte Variablen!

1. Modell: $X_1, \dots, X_n \sim N(\mu, \sigma_x^2)$

2. Nullhypothese

$H_0: \mu_0 = ?$

$H_A: \mu_A \neq \mu_0, \mu_A > \mu_0, \mu_A < \mu_0$

Einseitiger oder zweiseitiger Test?

3. Teststatistik $Z \sim N(0,1)$

$$Z = \frac{\bar{X}_n - \mu_0}{\frac{\hat{\sigma}_{\bar{X}_n}}{\sqrt{n}}} = \frac{\sqrt{n}(\bar{X}_n - \mu_0)}{\sigma_x} \quad \hat{\sigma}_{\bar{X}_n} = \frac{\sigma_x}{\sqrt{n}}$$

(= beobachtet - erwartet / Standardfehler)

4. Signifikanzniveau

$\alpha = ?$ (meist 5% $\rightarrow \alpha = 0.05$)

5. Verwerfungsbereich → TABELLE!

$\mu_A \neq \mu_0$ $K = (-\infty, -\Phi^{-1}(1-\alpha/2)] \cup [\Phi^{-1}(1-\alpha/2), \infty)$

$\mu_A > \mu_0$ $K = [\Phi^{-1}(1-\alpha), \infty)$

$\mu_A < \mu_0$ $K = (-\infty, -\Phi^{-1}(1-\alpha)]$

6. Testentscheid

Liegt der beobachtete Wert (in Z einsetzen) im

Verwerfungsbereich?

JA → H_0 wird verworfen

NEIN → H_0 wird nicht verworfen (aber auch nicht bewiesen)

Vertrauensintervall für μ (α -V.I.)

$\mu_A \neq \mu_0$ $[\bar{x}_n \pm \Phi^{-1}(1-\frac{\alpha}{2}) * \frac{\sigma_x}{\sqrt{n}}]$

$\mu_A > \mu_0$ $[\bar{x}_n - \Phi^{-1}(1-\alpha) * \frac{\sigma_x}{\sqrt{n}}, \infty)$

$\mu_A < \mu_0$ $(-\infty, \bar{x}_n + \Phi^{-1}(1-\alpha) * \frac{\sigma_x}{\sqrt{n}}]$

Schreibweise: $\Phi^{-1}(1-\frac{\alpha}{2}) = Z_{1-\alpha/2}$

T-Test besser als Z-Test

Bsp. $P(Z \leq 2) \approx 0.97 > P(T \leq 2) \approx 0.95$

Vertrauensintervall: enthält unplausible Werte für Teststatistik falls Nullhypothese stimmt und plausible Parameterwerte

Gepaarter T-Test EINSTICHPROBENTEST

σ_x : unbekannt; $\hat{\sigma}_x$: geschätzt

X_i : kontinuierlich, normalverteilt; Gepaarte Variablen!

SCHÄTZEN: $\hat{\sigma}_x$ und $\hat{\sigma}_x^2 \rightarrow$ geschätzte Varianz!

$$\hat{\sigma}_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2 \quad \hat{\sigma}_x \text{ in } \sqrt{\hat{\sigma}_x^2}$$

1. Modell: $X_1, \dots, X_n \sim N(\mu, \sigma_x^2)$

2. Nullhypothese

$H_0: \mu_0 = ?$

$H_A: \mu_A \neq \mu_0, \mu_A > \mu_0, \mu_A < \mu_0$

Einseitiger oder zweiseitiger Test?

3. Teststatistik $T \sim t_{n-1}$

$$T = \frac{\bar{X}_n - \mu_0}{\frac{\hat{\sigma}_{\bar{X}_n}}{\sqrt{n}}} = \frac{\sqrt{n}(\bar{X}_n - \mu_0)}{\hat{\sigma}_x} \quad \hat{\sigma}_{\bar{X}_n} = \frac{\hat{\sigma}_x}{\sqrt{n}}$$

4. Signifikanzniveau

$\alpha = ?$ (meist 5% $\rightarrow \alpha = 0.05$)

5. Verwerfungsbereich → TABELLE!

$\mu_A \neq \mu_0$ $K = (-\infty, -t_{n-1, 1-\alpha/2}] \cup [t_{n-1, 1-\alpha/2}, \infty)$

$\mu_A > \mu_0$ $K = [t_{n-1, 1-\alpha}, \infty)$

$\mu_A < \mu_0$ $K = (-\infty, -t_{n-1, 1-\alpha}]$

6. Testentscheid

Liegt der beobachtete Wert im Verwerfungsbereich?

Beob. Wert einsetzen: $t = \frac{\sqrt{n}(\bar{x} - \mu_0)}{\hat{\sigma}_x}$

JA → H_0 wird verworfen

NEIN → H_0 wird nicht verworfen (aber auch nicht bewiesen)

Vertrauensintervall für μ (α -V.I.)

$\mu_A \neq \mu_0$ $[\bar{x}_n \pm t_{n-1, 1-\alpha/2} * \frac{\hat{\sigma}_x}{\sqrt{n}}]$

$\mu_A > \mu_0$ $[\bar{x}_n - t_{n-1, 1-\alpha} * \frac{\hat{\sigma}_x}{\sqrt{n}}, \infty)$

$\mu_A < \mu_0$ $(-\infty, \bar{x}_n + t_{n-1, 1-\alpha} * \frac{\hat{\sigma}_x}{\sqrt{n}}]$

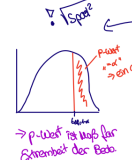
Schreibweise: $t_{df} \rightarrow$ degrees of freedom

Generell: Je kleiner df, desto wahrscheinlicher sind

Werte mit grossem Betrag

→ gepaarter 2 Stichproben T-Test:

gleiche Formel (Diff. zw. Z_u & Z_u) μ



Anzahl falscher Tests bei $\alpha = a$ wird \times Tests $\text{mg} = a \times$

Ungepaarter T-Test ZWEISTICHPROBENTEST

Annahme: $\sigma_x^2 = \sigma_y^2 \rightarrow$ gleiche Varianz!

$\bar{X}_n = \frac{1}{n} \sum x_i$ $\bar{Y}_m = \frac{1}{m} \sum y_i \rightarrow$ arithmetisches Mittel

$$S_{\text{pool}}^2 = \frac{1}{n+m-2} \left(\sum_{i=1}^n (x_i - \bar{x}_n)^2 + \sum_{i=1}^m (y_i - \bar{y}_m)^2 \right) = \frac{1}{n+m-2} ((n-1)\hat{\sigma}_x^2 + (m-1)\hat{\sigma}_y^2)$$

1. Modell: $X_1, \dots, X_n \sim N(\mu, \sigma_x^2)$

2. Nullhypothese

$H_0: \mu_0 = ?; \mu_x = \mu_y$

$H_A: \mu_A \neq \mu_0, \mu_A > \mu_0, \mu_A < \mu_0$

Einseitiger oder zweiseitiger Test?

3. Teststatistik $T \sim t_{n+m-2}$

$$T = \frac{\bar{X}_n - \bar{Y}_m}{\frac{S_{\text{pool}}}{\sqrt{\frac{1}{n} + \frac{1}{m}}}}$$

beachte: oben S_{pool}^2 , hier S_{pool}

4. Signifikanzniveau

$\alpha = ?$ (meist 5% $\rightarrow \alpha = 0.05$)

5. Verwerfungsbereich → TABELLE!

$\mu_A \neq \mu_0$ $K = (-\infty, -t_{n+m-2, 1-\alpha/2}] \cup [t_{n+m-2, 1-\alpha/2}, \infty)$

$\mu_A > \mu_0$ $K = [t_{n+m-2, 1-\alpha}, \infty)$

$\mu_A < \mu_0$ $K = (-\infty, -t_{n+m-2, 1-\alpha}]$

6. Testentscheid

Liegt der beobachtete Wert im Verwerfungsbereich?

Beobachteten Wert einsetzen

JA → H_0 wird verworfen

NEIN → H_0 wird nicht verworfen (aber auch nicht bewiesen)

GEPAART	UNGEPAART
Gleichgrosse Stichproben	Unterschiedliche Stichproben
Klare Zuordnung (rechts/links, vorher/nachher)	Keine Zuordnung
Differenz der Paare	
MEHR MACHT	WENIGER MACHT

ungepaarter T-Test: andere Formel: keine Differenz $\rightarrow Z_u$

p-Wert in Tabelle: $t_{df, x} = T$

$p = 2 \cdot (1-x)$
↳ zweiseitig

Mann-Whitney U-Test (Zweistichproben W-Test)

Annahme: $X_i \sim F$ i.i.d.
 $Y_i \sim G$ i.i.d.

Teste μ : $H_0: F=G$ (X_i und Y_i gleich verteilt)
 $H_A: F=G+\gamma$ (X_i, Y_i gleiche Form aber verschoben)

Prinzip: $H_0: F=G$

1. Absolute Beträge der beob. Werte nehmen
 2. Ränge nehmen
 3. getrennt nach positiven und negativen Ausgangswerten Rangsummen bilden
 4. mit PC U-Test machen
- falls H_0 stimmt: + bzw. - Rangsummen etwa gleich, sonst stehen sie in einem bestimmten Verhältnis.

Welch-Test: setzt nicht voraus, dass σ beider Gruppen gleich sind.

Bonferroni Korrektur Multiples Testen *Fehler 1. Art (falsch positiv) höchstens = α*

Bsp.: 1000 Tests mit $\alpha=0.05$ (H_0 sei bei jedem wahr)
 Wegen Fehler 1. Art (H_0 verwerfen obwohl wahr) werden bei etwa 5% der Tests H_0 fälschlicherweise verworfen. → SG

Ziel: für alle 1000 Tests insgesamt eine Wa. von 5% zu haben, H_0 fälschlicherweise zu verwerfen!

Die Tests werden auf dem neuen Signifikanzniveau

$\tilde{\alpha} = \frac{\alpha}{n}$ (n: Anzahl Tests) durchgeführt.

→ verringert Wa. für falsch positiv für alle Test zusammen auf < 5%

Welchen Test muss ich durchführen?



Stichprobengröße (Faustregel)

$n \geq 4 \frac{\sigma_x^2}{\delta^2}$ δ : Breite des 95%-VI

Vertrauens-/Vorhersage-Intervall

95% Vertrauensintervall < 95% Vorhersageintervall

Vertrauensintervall: Durchschnitt (viele Messungen) Gibt erwarteten Wert mit Wa. $1-\alpha$ heraus	Vorhersageintervall: Man kann bei EINER Messung etwa so viel erwarten Grössere Streuung
---	---

Wa., dass Test bei wahren H_0

kein sign. Ergebnis zeigt: $1-\alpha$

→ Wa., dass bei n Tests keiner eins zeigt: $(1-\alpha)^n$

→ Wa., dass mind. 1 Test fälschlicherweise verwirft: $1-(1-\alpha)^n$

→ Fehler folgt Normalverteilung mit Erwartungswert 0

→ für gegebene x folgt Y einer Normalverteilung mit $E(Y) = \beta_0 + \beta_1 x$

Lineare Regression

R-Output:

Bsp.: Ein Bauer notiert den Ertrag von Kartoffelfeldern in kg pro m². Er stellt fest, dass sich durch Dünger den Ertrag erhöhen lässt, ist sich jedoch nicht sicher ob es wirklich am Dünger liegt. Er versucht deshalb, den Ertrag durch eine Lineare Regression zu beschreiben. Falls es am Dünger liegt, hängt der Ertrag somit vom Grundertrag (β_0), vom Zusatztertrag durch Einsatz von x Einheiten Dünger ($\beta_1 \cdot \text{dünger}$) und von einem Standardfehler (σ) ab.

$H_0: \beta_1 = 0$ (Dünger hat keinen Einfluss)

$H_A: \beta_1$ hilft zur Erklärung des Modells (es kommt auf die Menge Dünger an)

ertrag = $\beta_0 + \beta_1 \cdot \text{dünger} + \epsilon$, $\epsilon \sim N(0, \sigma)$.

Der (unvollständige) Regressionsoutput sieht wie folgt aus:

Residuals:

Min 1Q Median 3Q Max
-86.00 -9.96 3.17 10.99 35.86

Coefficients:	β	$\sigma(\beta)$	$t(\beta)$ Teststatistik	P-Wert
	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	3.0437	6.5182	0.47	0.64
dünger	0.5471	0.0612	8.93	1.3e-11 ***

0 **** 0.001 *** 0.01 ** 0.05 ' 0.1 ' ' 1

Residual standard error: 21.6 on 46 degrees of freedom

Multiple R-squared: 0.634, Adjusted R-squared: 0.626

Meistens 2

Anz. Datenpunkte: Anz. Felder = degrees of freedom + Anz. β 's = 46+2 = 48 = n
→ $df = n - 2$

Koeffizienten: $\beta = \sigma(\beta) \cdot t(\beta)$

Bsp.: $\beta_1 = \sigma(\beta_1) \cdot t(\beta_1) = 0.0612 \cdot 8.93 = 0.5471$

95%-VI:

genau: $VI(\beta) = \beta \pm t_{df, 0.975} \cdot \sigma(\beta)$

Bsp.: $VI(\beta_1) = \beta_1 \pm t_{46, 0.975} \cdot \sigma(\beta_1) = 0.5471 \pm 2.01 \cdot 0.0612 = [0.42, 0.67]$

→ da 0 nicht in VI, H_0 wird verworfen

Approx.: $VI(\beta) = \beta \pm 2 \cdot \sigma(\beta)$

Verwerfungsbereich:

$K(\beta) = (-\infty, -t_{df, 1-\alpha/2}] \cup [t_{df, 1-\alpha/2}, \infty)$

Bsp.: für $\alpha = 0.05$ $K(\beta_1) = (-\infty, -2.01] \cup [2.01, \infty)$

→ da $t(\beta_1) = 8.93$ in K, H_0 wird verworfen

P-Wert:

$t(\beta) = t_{df, 1-P\text{-Wert}/2}$ → Tabelle: was ist P-Wert? $(1-P\text{-Wert}/2)$

Bsp.: $t(\beta_0) = t_{46, 1-P\text{-Wert}/2} = 0.47$ → Tabelle: $(1-P\text{-Wert}/2) \sim 0.68$

P-wert = $2 \cdot (1-0.68) = 0.64$

→ Wenn t-Wert weit über größtem Wert von Tabelle

in Zeile mit df ist, ist

p-wert null

→ P-Wert = α
 $t_{\text{table}} = \frac{\text{Estimate}}{\text{Std. Error}}$ → Wert in Tabelle für t_n -Verteilung
→ bei zweiseitig: verdoppeln

Erwarteter Ertrag:

Menge Dünger = x → Ertrag = $\beta_0 + \beta_1 \cdot x$

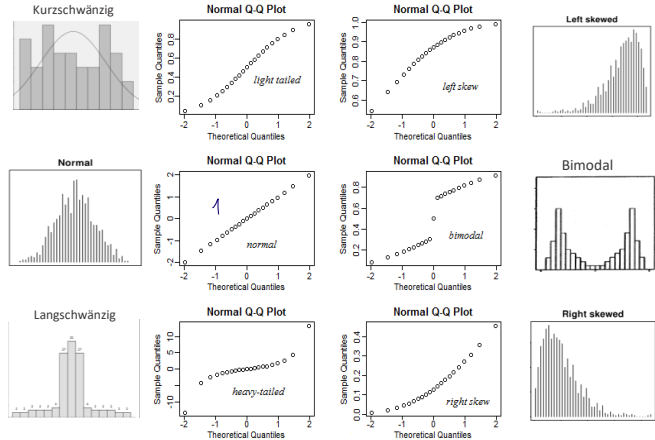
Bsp.: Bei 1.5l/m² Dünger: Ertrag = $3.0437 + 0.5471 \cdot 1.5 = 3.8643$

QQ-Plot und Histogramm (siehe auch S.5)

Ist es plausibel, dass QQ-Plot und Histogramm von denselben Daten stammen? Welche Graphiken gehören zusammen?

Zur Erinnerung: nur bei Geraden im QQ-Plot ist die Normalitätsannahme bestätigt.

uniform

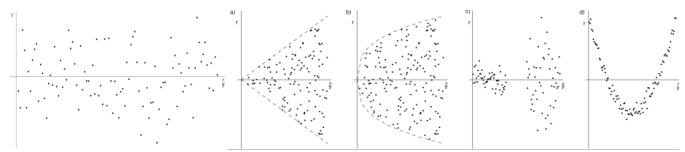


Tukey-Anscombe Plot

Ist die Fehlervarianz konstant? Passen die Daten zum Modell?

Ja:

Nein:

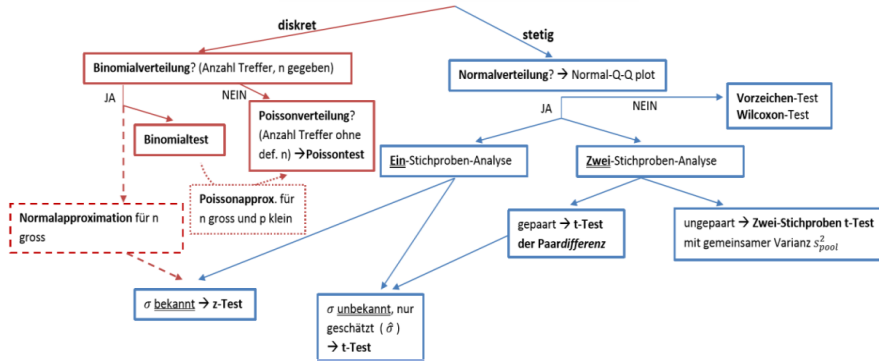


Fehlervarianz nicht konstant

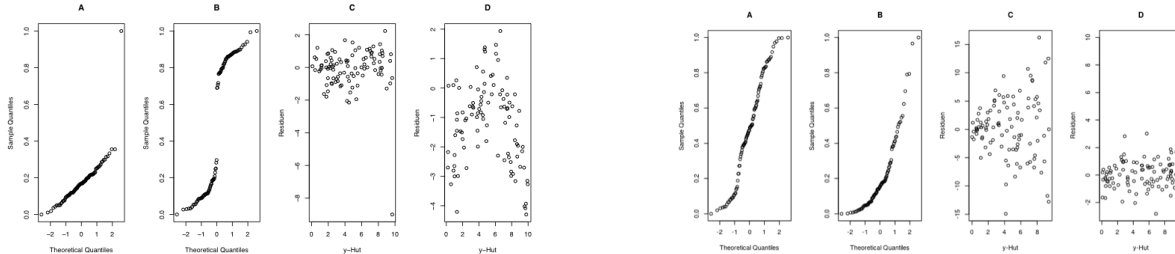
Systematischer Fehler

Residuen: vertikaler Abstand zw. beobachtetem Punkt und gefitteter Gerade

Man kann nicht alle Modellannahmen mit TA überprüfen



Zufallsvariablen Streudiagramme → arith. Mittel muss kleineren Stichprobenwert entsprechen
→ Median ($\alpha = 50\%$ - Quantil) nur wenn $\alpha \cdot n$ keine ganze Zahl ist
→ dann Zahl aus Stichprobe



- ✓ a) Die Residuen in 'A' könnten von einer Normalverteilung mit Ausreißern stammen.
- ✗ b) Die Residuen in 'B' könnten von einer Normalverteilung mit Ausreißern stammen.
- ✓ c) Plot 'C' zeigt, dass abgesehen von einigen Ausreißern die Linearitätsannahme erfüllt ist.
- ✓ d) Plot 'D' zeigt, dass die Linearitätsannahme des geschätzten Modells nicht angebracht ist.

- ✓ a) Betrachten Sie den QQ-Plot in Abbildung 'A'. Es ist plausibel, dass die Residuen von einer Uniformen Verteilung stammen.
- ✓ b) Betrachten Sie den QQ-Plot in Abbildung 'B'. Es ist plausibel, dass die Residuen von einer Rechtsschiefen Verteilung (d.h. die Dichtefunktion fällt bei grossen Zahlen langsamer ab als bei kleinen Zahlen) stammen.
- ✗ c) Betrachten Sie den TA-Plot in Abbildung 'C'. Der Plot zeigt konstante Varianz der Residuen.
- ✗ d) Betrachten Sie den TA-Plot in Abbildung 'D'. Der Plot zeigt keine auffallenden Abweichungen von den Modellannahmen.