

1. MODELLE FÜR ZÄHLDATEN

1.1 WAHRSCHEINLICHKEITSMODELLE

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

$P(A) \geq 0$	Wa. von A
$P(\Omega) = 1$	gesamte Wa.

$P(A \cup B)$	$P(A) + P(B) - P(A \cap B)$	A oder B	[1]
$P(A \cap B)$		A und B	[2]
$P(A^c) P(\bar{A})$	$1 - P(A)$	nicht A	[3]
$P(A \setminus B)$	$P(A) - P(A \cap B)$	A ohne B	[4]



- Disjunkt = getrennte Elemente [5]
- $P(A \cap B) = 0 = \{ \} | P(A|B) = 0 | P(A \cup B) = P(A) + P(B)$

1.2 WAHRSCHEINLICHKEITEN BERECHNEN

- Laplace-Modell = alle Elementarereignisse sind gleich wahrsch

$$P(A) = \frac{\text{Anzahl günstige } \omega}{\text{Anzahl mögliche } \omega}$$

Bsp. Fairer Würfel:

- Wa. eine 6 zu Würfeln: $1/6$
- Wa. eine gerade Zahl zu Würfeln $3/6 = 1/2$

- mit Gegenereignis rechnen \rightarrow z.T. weniger Rechenaufwand

1.3 UNABHÄNGIGKEIT

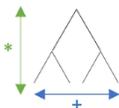
- Auftreten von A beeinflusst Wa. von B nicht \rightarrow gleicher Baum

$$P(A \cap B) = P(A) * P(B) \quad P(A) = P(A|B)$$

Bsp. Wa. zweimal hintereinander eine 6 zu Würfeln $\frac{1}{6} * \frac{1}{6}$

1.4 ABHÄNGIGKEIT / BEDINGTE WAHRSCHEINLICHKEIT

- Wa. beeinflussen sich \rightarrow ungleicher Baum
- Bedingte Wa. von A wenn B eingetreten: $P(A|B)$
- $P(A|B) \neq P(A)$



- Satz von Bayes: $P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A) * P(A)}{P(B)}$

\rightarrow richtiger Weg & korrektes Ergebnis dividiert durch alle Wege, welche das gleiche Ergebnis hervorrufen

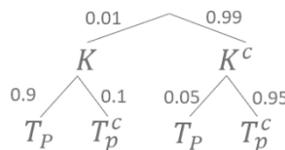
- Satz der totalen Wa. (zwei Ereignisse)
- $P(B) = P(B|A) * P(A) + P(B|A^c) * P(A^c)$
- bei mehr Ereignissen: $P(B) = \sum_{i=1}^k P(B|A_i)P(A_i)$

Bsp. Wa., dass wenn Test positiv (T_p) auch tatsächlich die Krankheit (K) vorhanden ist \rightarrow Gesucht: $P(D|T)$

Gegeben: $P(K) = 0.01 \rightarrow P(K^c) = 0.99$
 $P(T_p|K) = 0.9 \rightarrow P(T_p^c|K) = 0.1$
 $P(T_p^c|K^c) = 0.95 \rightarrow P(T_p|K^c) = 0.05$

Satz von Bayes: $P(K|T_p) = \frac{P(T_p|K) * P(K)}{P(T_p)}$ gesucht

Satz tot. Wa. $P(T_p) = P(T_p|K) * P(K) + P(T_p|K^c) * P(K^c)$
 $= 0.9 * 0.01 + 0.05 * 0.99 = 0.0585 \rightarrow$ pos. Vorhersagewert
 $P(K|T_p) = \frac{0.9 * 0.01}{0.0585} = 0.1538 \approx 15.4\%$



1.5 ODDS, LOG-ODDS, ODDS-RATIO

- Odds = wie viel Mal wahrsch. dass A eintritt statt A^c
- $odds(A) = \frac{P(A)}{1 - P(A)} = \frac{P(A)}{P(A^c)}$ wenn odds \uparrow dann auch log-odds \uparrow

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

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

Risikofaktor	Ja	Nein
Krankheit	a	b
Nein	c	d

- Odds Ratio (OR) = Chancenverhältnis, Zahl, welche etw. über die Stärke eines Zusammenhangs aussagt

$$OR = \frac{odds(A|B)}{odds(A|B^c)} = \frac{a+d}{c+b}$$

$$P(A) = \frac{a}{a+c} \rightarrow \text{Wa. zu erkranken mit Risikofaktor}$$

$$P(B) = \frac{b}{b+c} \rightarrow \text{Wa. zu erkranken ohne Risikofaktor}$$

- Risk difference (RD) = Wirksamkeit eines bestimmten Faktors
- $RD = P(A|B) - P(A|B^c)$

- Risk Ratio (RR) = Wa. durch einen bestimmten Faktor vergrößert / verkleinert wird

$$RR = \frac{P(A|B)}{P(A|B^c)}$$

Bsp. Wie hoch ist das Chancenverhältnis an einer Krankheit (K) zu erkranken, wenn unter Einfluss eines Risikofaktors (RF) steht? (z.B. Lungenkrebs, Raucher) \rightarrow Gesucht: OR

Gegeben: $P(K|RF) = 0.54 \quad P(K|RF^c) = 0.18$

$$RD = P(K|RF) - P(K|RF^c) = 0.54 - 0.18 = 0.36$$

$$RR = \frac{P(K|RF)}{P(K|RF^c)} = \frac{0.54}{0.18} = 3$$

$$odds(K|RF) = \frac{P(K|RF)}{P(K^c|RF)} = \frac{0.54}{1-0.54} = 1.17$$

$$odds(K|RF^c) = \frac{P(K|RF^c)}{P(K^c|RF^c)} = \frac{0.18}{1-0.18} = 0.22$$

$$OR = \frac{odds(K|RF)}{odds(K|RF^c)} = \frac{1.17}{0.22} = 5.33$$

\rightarrow Wa. Krank zu werden ist mit diesem RF \sim 5-mal grösser

2. DISKRETE WAHRSCHEINLICHKEITSVERTEILUNGEN

2.1 ZUFALLSVARIABLE

- Diskrete Wa.verteilung \rightarrow nur für bestimmte Werte definiert x nimmt endliche, abzählbare Werte an
- Kumulative Verteilungsfkt. \rightarrow Summe der Wa. welche \leq Wert x sind \rightarrow stetig steigend, höchster Wert ist 1 \rightarrow springt an den Stellen, die zum Wertebereich gehören
- Zufallsvariable = Zufallsexperiment = als Ergebnis eine Zahl hat
- Grossbuchstabe = Funktion | | Kleinbuchstabe = konkreter Wert

2.2 KENNZAHLEN EINER VERTEILUNG

- Erwartungswert = mittlere Lage einer Verteilung

$$E(X) = \sum x * P(X = x)$$

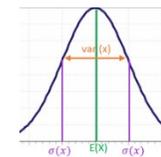
\rightarrow Zahl₁ * Wa₁ + Zahl₂ * Wa₂ + ...

- Varianz = Abweichung vom Erwartungswert

$$Var(X) = \sum (x - E(X))^2 * P(X = x)$$

- Standardabweichung = Streuung der Verteilung

$$\sigma(X) = \sqrt{Var(X)}$$



2.3 BINOMIALVERTEILUNG

- Beschreibung des Eintreffens / Nicht-Eintreffens eines bestimmten Ereignisses
- Binomialkoeffizient (TR: $nCr(n, x)$, menu;5;3) = auf wie viele Arten man x Dinge auf n Plätzen anordnen kann (Lottozahlen)

$$\binom{n}{x} = \frac{n!}{x!(n-x)!} \quad \binom{n}{2} = \frac{n(n-1)}{2} \rightarrow \# \text{ Händeschütteln}$$

- Binomialverteilung: Anz. Erfolge in einer Serie von gleichartigen, unabhängigen Versuchen, die jeweils genau zwei mögliche Ergebnisse haben (= Bernoulli-Prozesse)

$$X \sim Bin(n, \pi) \quad P(X = x) = \binom{n}{x} * \pi^x * (1 - \pi)^{n-x} \text{ (binompdf)}$$

x = # Erfolge | n = # Versuche | π = # Wa. für Erfolg

TR: menu;5;5;A \rightarrow pdf | | menu;5;5;B \rightarrow cdf

$$P[X = x] = \text{binompdf}(n, \pi, x) \text{ «genau x Treffer»}$$

$$P[X \leq x] = \text{binomcdf}(n, \pi, x) \text{ «höchstens x»}$$

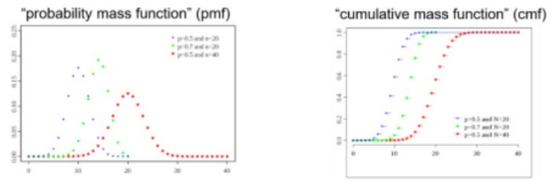
$$P[X < x] = \text{binomcdf}(n, \pi, x - 1)$$

$$P[X \geq x] = 1 - \text{binomcdf}(n, \pi, x - 1) \text{ «mindestens x»}$$

$$P[X > x] = 1 - \text{binomcdf}(n, \pi, x)$$

$E(X) = n * \pi$	$Var(X) = n\pi(1 - \pi)$
$\sigma(X) = \sqrt{Var(X)} = \sqrt{n * \pi * (1 - \pi)}$	

- Streuung $Var(x)$ wächst mit $n \rightarrow$ aber immer langsamer
- Festes n ist $Var_{max}(x)$ bei $\pi = 0.5$
- $P(X=x)$ (pmf = probability mass fkt) ist maximal wenn $x \approx n * \pi$
- Wenn n gross, sind Wa. nur um $n * \pi$ herum gross, sonst klein
- Wenn $n * \pi * (1 - \pi)$ nicht zu klein \rightarrow Verteilung symmetrisch



2.4 BERNOULLIVERTeilUNG

- Erfolg oder Misserfolg bei nur einem Versuch
- $P(X = 1) = \pi \quad || \quad P(X = 0) = 1 - \pi \quad || \quad 0 \leq \pi \leq 1$

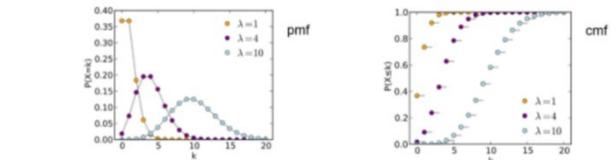
$E(X) = \pi \quad Var(X) = \pi(1 - \pi) = \sigma^2$

- Unabhängige Zufallsvariable
- $X \sim Bin(n, \pi) \quad | \quad Y \sim Ber(\pi) \rightarrow X + Y \sim bin(n + 1, \pi)$

2.5 POISSONVERTEILUNG

- Seltene Ereignisse werden in einem vorgegebenen Zeitraum gezählt (Bspw. Unfallraten, radioaktiver Zerfall) \rightarrow keine klare Obergrenze vorhanden
- $X \sim Pois(\lambda) \quad P(X = x) = \frac{\lambda^x}{x!} * e^{-\lambda} \quad (x \in \{0, 1, \dots, \infty\})$
- λ = durchschnittliches auftreten pro Zeit
 ∞ = einfach eine grosse Zahl im TR eintippen
- Poissonverteilungen addierbar, wenn unabhängig voneinander $\rightarrow E((X + Y) \quad | \quad E(a * X) = a * E(X)$
 - Poisson-Approximation: für grosse n und kleine π gilt: $Bin(n, \pi) \approx Pois(\lambda)$
 - TR: **poissCdf** (λ , untere Grenze, obere Grenze) (menu;5;5; k)

$E(X) = \lambda \quad Var(X) = \lambda \quad \sigma(X) = \sqrt{\lambda}$



2.6 UNIFORME VERTEILUNG

- Alle Ereignisse sind gleich wahrsch.
- $X \sim Unif(n) \quad P(X = x) = \frac{1}{n} \quad (x \in \{0, 1, \dots, n\})$

$E(X) = \frac{n+1}{2} \quad Var(X) = \frac{(n+1)(n-1)}{12} = \sigma^2$

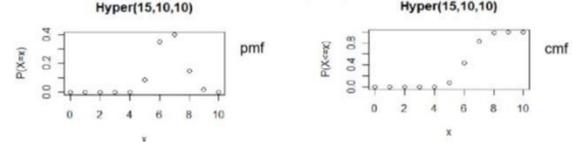
2.7 HYPERGEOMETRISCHE VERTEILUNG

- Gewinnwa. nicht konst., da abhängig was davor gezogen wurde \rightarrow vernachlässigbar, wenn sehr grosse Lostrommel und wenige Ziehungen

$X \sim hyper(N, n, m) \quad P(X = x) = \frac{\binom{m}{x} \binom{N-m}{n-x}}{\binom{N}{n}} = \frac{\text{«günstige»}}{\text{«mögliche»}}$

TR (gespeicherte Funktion): **hyp(x,N, m, n)**
 Bezogen auf Kugelziehen aus Urne: $N = \#$ Anz
 $m = \#$ markierte $n =$ gezogen $| \quad x =$ markiert von gezogenen

$E(X) = \frac{n * m}{N}$
 $Var(X) = \frac{n * m * (N - m)(N - n)}{N^2 * (N - 1)} = \sigma^2$



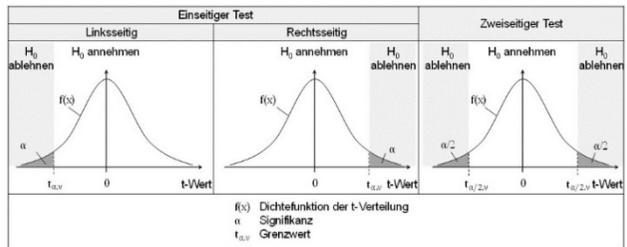
3. GRUNDFRAGESTELLUNGEN DER STATISTIK

3.1 WELCHES IST DER ZU DEN BEOBACHTUNGEN PLAUSIBELSTE PARAMETERWERT π ?

- Antwort: Punkt-Schätzung \rightarrow Schätzwert mit Dach: $\hat{\pi}$
- Momenten-Methode (MM): Beobachtung = $E(X)$
- Maximum-Likelihood-Methode (MLE):
 - Welches ist die höchste Wa. basierend auf Messdaten?
 $\frac{dP(X=x)}{d\pi} = \frac{d}{d\pi} (\binom{n}{x} \pi^x * (1 - \pi)^{n-x})$
 - Nullstellen der Ableitung suchen \rightarrow Extrema
 - Für eine Beobachtung: $\hat{\pi} = \frac{x}{n}$
 - π soll so gewählt werden, dass am besten zur Beobachtung passt
 Bsp. Münze wird 3-mal geworfen Ergebnis: ZKK, Würfe sind unabhängig voneinander, $p =$ Wa. für K, $MLE(p) = 0.33 \rightarrow 1/3$

3.2 SIND DIE BEOBACHTUNGEN KOMPATIBEL (STATISTISCH VEREINBAR) MIT VORGEGEBENEM PARAMETER π ?

- Antwort: statistischer Test
- 1) Modell \rightarrow passende Wa.-Verteilung und X festlegen
 - 2) H_0 (Nullhypothese) & H_A (Alternativhypothese) aufstellen
 - $H_0: \pi = \pi_0$ kein Effekt / Veränderung
 - $H_A: \pi \neq \pi_0$ zweiseitig
 - $H_A: \pi > \pi_0$ einseitig nach oben
 - $H_A: \pi < \pi_0$ einseitig nach unten
 - 3) Teststatistik (T) = # Erfolge bei n Versuchen
 - unter Annahme, dass H_0 stimmt
 - T ist eine Funktion von X
 - 4) Signifikanzniveau α festlegen (meist: 0.05 || 0.01) = Stichprobendaten, welche von der Annahme abweichen, dass Annahme verworfen wird
 - wenn α grösser wird vergrössert sich Verwerfungsbereich d.h. Macht nimmt tendenziell zu
 - 5) Verwerfungsbereich K = enthält alle Ereignisse, welche «abnormal» genug sind \rightarrow bei Binomialtest mit binomcdf
 - Gesucht: kleinste Zahl c, für die gilt: $P(X \geq c) \leq \alpha$



$H_A: \pi \neq \pi_0$ $K = [0, c_u] \cup [c_o, n]$ zweiseitig $\rightarrow \frac{\alpha}{2}$
 $H_A: \pi > \pi_0$ $K = [c, n]$ einseitig
 $H_A: \pi < \pi_0$ $K = [0, c]$ einseitig

- Einseitiger Test = nur Abweichungen in eine Richtung von H_0 detektiert \rightarrow nicht grosse Abweichung nötig, damit detektiert
 - Macht gross
 - unsicher, dass die Abweichung, in die eine Richtung von der H_0 relevant ist \rightarrow einseitiger Test bevorzugt
- Zweiseitiger Test = Abweichungen in beide Richtungen von H_0
 - müssen aber gross sein, damit detektiert \rightarrow Macht ist klein
 - zweiseitiger Verwerfungsbereich mit $\frac{1}{2}$ Signifikanzniveau α \rightarrow Immer eine untere und eine obere Grenze vorhanden
 - Spezifität nicht sehr hoch \rightarrow eignet sich, wenn man grosse/kleine Gewinn-Wa. erkennen will

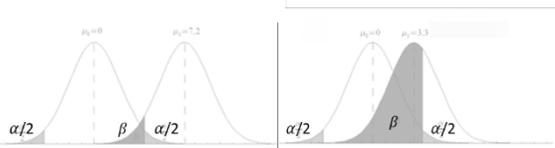
- Normalapproximation (nur für Überschlagsrechnungen)
 - Falls n gross und $\pi \approx 0.5$
 - gut falls $n * \pi > 5$ und $n * (1 - \pi) > 5$
 - je grösser # Versuche \rightarrow bessere Normalapproximation
- $\pi \neq \pi_0 \quad K = [0, c_u] \cup [c_o, n] \quad c = n\pi_0 \pm 1.96\sqrt{n\pi_0(1-\pi_0)}$
 TR (gesp. Fkt.): napp_1o(n, π_0); napp_1u $\frac{\alpha}{2}$ pro Seite
 \rightarrow o = aufrunden, u = abrunden
- $\pi > \pi_0 \quad K = [c, n] \quad c = n\pi_0 + 1.64\sqrt{n\pi_0(1-\pi_0)}$
 TR (gesp. Fkt.): napp_2(n, π_0); aufrunden
- $\pi < \pi_0 \quad K = [0, c] \quad c = n\pi_0 - 1.64\sqrt{n\pi_0(1-\pi_0)}$
 TR (gesp. Fkt.): napp_3(n, π_0); abrunden

6) Testentscheid: Beobachteter Wert im Verwerfungsbereich?

- Ja = H_0 verwerfen
- Nein = H_0 nicht verworfen aber auch nicht bewiesen
- Fehlertypen
 1. Art = α -Fehler = H_0 verwerfen, obwohl H_0 korrekt
 2. Art = β -Fehler = H_0 beibehalten, obwohl H_A korrekt
 - α beschränkt Mögl. für Fehler 2. Art \rightarrow grösser falls α kleiner
 - wenn Fehler 1. Art kleiner wird, wird Fehler 2. Art grösser
 - wenn Fehler 2. Art grösser wird, wird Fehler 1. Art kleiner

$\alpha = 0.05$	H_0 wahr	H_0 falsch
H_0 keep	Spezifität $\rightarrow 1 - \alpha$	Fehler 2. Art $\rightarrow \beta$
H_0 discard	Fehler 1. Art $\rightarrow \alpha$	Sensitivität / Macht $\rightarrow 1 - \beta$

- Macht = Wa. H_A zu entdecken, falls H_A richtig ist
 $1 - P(\text{Fehler 2. Art})$
 - Je grösser Fehler 1. Art (α), desto grösser die Macht
 - wird bei zweiseitigem Test meist kleiner (ausser einseitiger detektiert auf falsche Seite \rightarrow Macht = 0)
- $H_0 \uparrow | H_A \downarrow | \text{Macht} \downarrow \quad || \quad \alpha \downarrow | \beta \uparrow | \text{Macht} \downarrow$



Bsp. Berechnung Macht
 $\rightarrow [0,8] \cup [18,19] | p_A = 0.85 | n = 19$
 Behauptung: Macht = 0.199
 binomcdf für beide Verwerfungsbereiche (=0.00002 & =0.198)
 Macht = Addieren der Werte = 0.199 \rightarrow Behauptung richtig

- P-Wert = Wa. unter Gültigkeit von H_0 das beobachtete Ergebnis oder ein extremeres zu erhalten
 = das kleinste α , bei dem H_0 gerade noch verworfen wird
 - Verwerfe H_0 falls P-Wert $\leq \alpha$
 - Belasse H_0 falls P-Wert $> \alpha$

- je kleiner der P-Wert, desto signifikanter Berechnung
 \rightarrow Ergebnis spricht gegen H_0
 - wenn x grösser wird, dann wird p-Wert kleiner (Rest konst.)
- Bsp. Einseitiger Binomialtest (p-Wert bestimmen)
 $\rightarrow \pi = 0.55 | H_A: \pi > \pi_0 | n = 10 | x = 3 |$
 Behauptung: p=0.204
 - zuerst binompdf (für genau 3 Treffer) = 0.0746 = Wa. dafür
 - p-Wert \rightarrow extremere Ereignisse, d.h. $x \geq 3 \rightarrow$
 binomcdf = 0.97 \rightarrow d.h. Behauptung falsch

3.3 WELCHE PARAMETERWERTE π SIND MIT DEN BEOBACHTUNGEN KOMPATIBEL?

- Antwort: Konfidenzintervall / Vertrauensintervall (VI) = Bereich, der mit einer gewissen Wa. den Parameter eine Verteilung einer Zufallsvariablen einschliesst
- Schliesst wahren Parameter mit $P(1 - \alpha)$ ein
- Werte von π_0 bei denen H_0 nicht verworfen wird
- Ein $(1-\alpha)$ -Vertrauensintervall enthält den wahren Parameter mit Wahrscheinlichkeit $1-\alpha$
- VI = $\{\pi_0 ; \text{Nullhypothese } (H_0: \pi = \pi_0) \text{ wird belassen}\}$

Bei $\alpha = 0.05$: $95\text{-VI} = \frac{x}{n} \pm 1.96 * \sqrt{\frac{x}{n} * (1 - \frac{x}{n}) * \frac{1}{n}}$

- TR (gesp. Fkt.): na_95o(n, x) oder na_95u(n, x)
- zweiseitiges 95%-VI umfasst bei einer Binomvert. alle Werte p_0 , für die ein zweiseitiger Binomialtest mit $H_0: p=p_0$ nicht verwirft \rightarrow mit 95%-Wa. ist korrekt
 - Normalapproximation: für 95%-VI (falls n gross) \rightarrow Faustregel
 - Parameter in VI, wenn nein \rightarrow verwerfen
 - \sqrt{n} - Gesetz = mit n -mal so vielen Beobachtungen wird das VI um den Faktor \sqrt{n} kleiner (weniger breit)
 \rightarrow für eine halb so grosse Standardabweichung (Streuung) braucht man 4-mal so viele Daten Bsp. $n = 4$, d. h. $\frac{VI}{\sqrt{4}}$
 - doppelte Genauigkeit bei viermal so vielen Daten

3.4 VERGLEICHE

- Hypothesentest (gut für internationale Feststellungen)
 - \oplus klares Prozedere, klare Aussage über Fehler 1. & 2.Art
 - \ominus wie deutlich verworfen? wie gross wahrer Parameter?
- P-Wert
 - \oplus klar ob und wie deutlich verworfen wird
 - \ominus keine Aussage Fehler 1&2, wie gross wahrer Parameter?
- 95%-VI (Unsicherheit vermindern \rightarrow beste Art)
 - \oplus klar ob und wie deutlich verworfen, klar wie gross der Parameter etwa ist
 - \ominus keine klare Aussage über Fehler 1. & 2. Art

Bsp. Einseitiger Binomialtest (Verwerfungsbereich bestimmen)
 $\rightarrow H_0: \pi = 0.6 | H_A: \pi > 0.6 | n = 26 \alpha = 0.05$
 Behauptung: $K[21,26] \rightarrow K = [c, n]$ (aus Theorie), weil $H_A: \pi > \pi_0$
 binomcdf(26, 0.6, 21, 26) = 0.021 < 0.05 \rightarrow d.h. K richtig gewählt
 3. Zahl \rightarrow untere Grenze (hier c), 4. Zahl \rightarrow obere Grenze (hier n)

4. DESKRIPTIVE STATISTIK

4.1 KENNZAHLEN

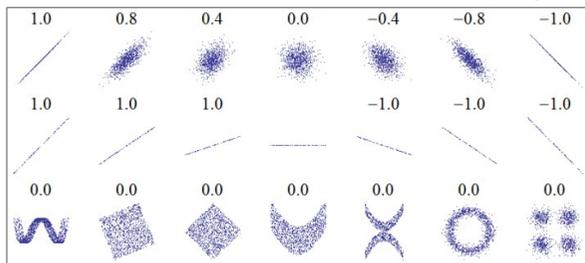
- Kennzahlen charakterisieren einen Datensatz
- Kennzahlen für Lage
 - Arithmetisches Mittel $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
 - Median $q_{0.5}$
- Kennzahlen für Streuung
 - Streuung \rightarrow empirische Standardabweichung
$$s_x = \sqrt{var} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$
 - IQR (=Inter-Quartile-Range)
- α -Quantil (empirisch, $\alpha = \%$)
 - q_α = Wert aus Daten \rightarrow bei dem mindestens $\alpha * 100\%$ der Datenpunkte kleiner und $(1 - \alpha) * 100\%$ grösser als q_α sind.
 - Daten müssen der Reihenfolge nach sortiert sein
 - $\alpha * n$ berechnen ($q_{0.25} \rightarrow$ 1. Quartil, $q_{0.75} \rightarrow$ 3. Quartil)
 - ganze Zahl: $q_\alpha = \frac{1}{2}(x_{\alpha*n} + x_{\alpha*n+1})$
 - keine ganze Zahl: $k = \alpha * n + \frac{1}{2} \rightarrow$ Runden zu Zahl
 - $q_\alpha = x_k$ (k-te Zahl)
 - empirische Median = $q_{0.5}$
 - mittlere Beobachtung, robust gegen Ausreisser
 - Zahl aus Stichprobe, wenn $\alpha * n$ keine ganze Zahl
 - wenn ganze Zahl \rightarrow Mittelwert aus zwei Beobachtungen aus der Stichprobe
 - Quartilsdifferenz (IQR)
 - Streuungsmass für Daten, robust gegen Ausreisser
$$IQR = q_{0.75} - q_{0.25}$$
- Kovarianz = Mass für gegenseitige Abhängigkeit \rightarrow misst ob Zahlen in die gleiche Richtung gehen

$$Cov(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$$
 - x, y = stetige Zufallsvariablen, μ = Erwartungswert
- Korrelation = Stärke einer statistischen Beziehung von zwei Variablen $Corr(X, Y) = \rho_{XY} = \frac{Cov(X, Y)}{\sigma_X \sigma_Y}$
 - dimensionslose normierte Zahl $\in [-1, 1]$
 - $Corr(X, Y) = 1 \rightarrow Y = a + bX \quad b > 0$
 - $Corr(X, Y) = -1 \rightarrow Y = a - bX \quad b > 0$
 - $Corr(X, Y) = 0 \rightarrow X$ und Y unabhängig voneinander

• Empirische Korrelation

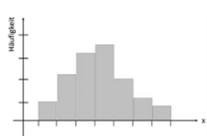
$$\hat{\rho} = r = \frac{s_{xy}}{s_x s_y} \quad || \quad s_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

- dimensionslose Zahl $\in [-1, 1]$
- Vorzeichen gibt Richtung an, Betrag gibt die Stärke des linearen Zusammenhangs zwischen zwei Variablen an
- Korrelationskoeffizient = Grad des lin. Zusammenhangs

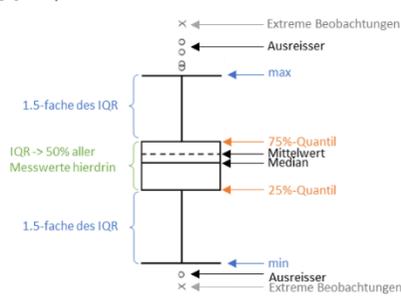


4.2 GRAPHISCHE METHODEN

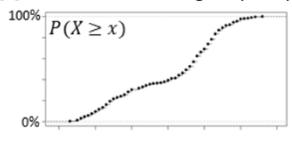
[1] Histogramm



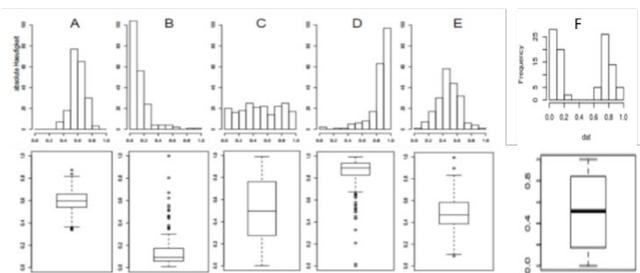
[3] Boxplot



[2] Kumulative Verteilungsfkt. (ECDF)



• Zusammenhang Histogramm – Boxplot



• Hypothese

- konfirmatorisch: langsam, formell, «bestätigen»
- explorativ: schnell, informell

4.3 ZUFALLSVARIABLEN (X=Zufallsvariable, x=Wert)

- $\sum P(X = x_i) = 1 \rightarrow$ wenn Elemente mehrfach gezählt werden > 1
- $P(X=x)$: Wa., dass ZV genau den Wert x produziert
- $P(X \geq x)$: Wa., dass ZV mindestens Wert x oder höher produziert
- $P(X \leq x)$: Wa., dass ZV mindestens Wert x oder tiefer produziert
- Funktion einer Zufallsvariable
- Lineare Transformation $g(x) = a + b * x$ für $Y = a + b * X$

$$E(Y) = a + b * E(X) \quad | \quad Var(Y) = b^2 * Var(X)$$

$$\sigma(Y) = |b| * \sigma(X) = \sqrt{b^2 * Var(X)}$$

- mehrere x vorhanden:
- $E(Y) = a + b_1 * E(X_1) + b_2 * E(X_2)$
- $Var(Y) = (b_1)^2 * Var(X_1) + (b_2)^2 * Var(X_2)$
- Quantil: $q_y^\alpha = a + b * q_x^\alpha$ falls $b > 0$

- i.i.d = independent, identically, distributet $\rightarrow X_1, \dots, X_n$ i. i. d
- Zufallsvariablen unabhängig und haben dieselbe Verteilung
- Funktion mehrerer Zufallsvariablen $y = g(x_1, x_2, \dots, x_n)$
- $E(X_1 + X_2) = E(X_1) + E(X_2) \rightarrow$ gilt immer
- $Var(X_1 + X_2) = Var(X_1) + Var(X_2)$ wenn X_1, X_2 unabhängig
- Gesetz der Grossen Zahlen (GGZ) = relative beobachtete Häufigkeit nähert sich mit steigendem Stichprobenumfang seinem Erwartungswert \rightarrow Einpendeln
- je grösser die Stichprobe, desto genauer die Schätzung
- Std.abweichung von $\bar{X}_n =$ Standard-fehler des arith. Mittels
- $E(\bar{X}_n) = \mu \quad | \quad Var(\bar{X}_n) = \frac{\sigma_x^2}{n} \quad | \quad \sigma(\bar{X}_n) = \frac{\sigma_x}{\sqrt{n}} \bar{X}_n \rightarrow \mu (n \rightarrow \infty)$
- Zentraler Grenzwertsatz (ZGS) = Mittelwert einer beliebigen Verteilung nähert sich mit zunehmendem Stichprobenumfang der Normalverteilung an
- je grösser Stichprobe, desto näher an Normalverteilung
- $(\bar{X}_n) \sim N(\mu_x, \frac{\sigma_x^2}{n}) \quad | \quad S_n = \sum_{i=1}^n x_n \quad | \quad S_n \sim N(n * \mu_x, n * \sigma_x^2)$

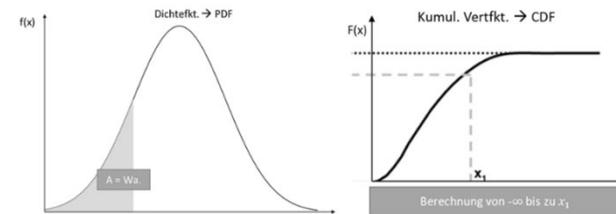
5. STETIGE VERTEILUNGEN

- Zufallsvariable X ist stetig, der Wertebereich kontinuierlich
- Bsp. Grösse messen \rightarrow unendlich viele Nachkommastellen
- Pkt.-Wa. $P(X = x) = 0 \quad | \quad P(X \in (a, b]) = P(a < x \leq b)$
- Kumulative Verteilungsfunktion (CDF) = Wa. dass das Ergebnis X kleiner oder gleich dem Wert x
- $F(x) = P(X \leq x)$ weil $P(a < x \leq b) = F(b) - F(a)$
- \rightarrow «kleiner als», Monotoner Anstieg von 0 bis 1

- Wahrscheinlichkeitsdichte $f(x) = F'(x)$ (PDF)
- $f(x) \geq 0$ für alle x (weil F monoton steigend)
- \rightarrow keine einzelnen Punkte vorhanden, sondern über unendlich viele Werte ein Integral \rightarrow Wa. der Dichtefunktion
- $P(a < x \leq b) = F(b) - F(a) = \int_a^b f(x) dx$
- $\int_{-\infty}^{\infty} f(x) dx = 1$
- Nullverteilung (null distribution) = Stichprobenverteilung der Teststatistik gemäss H_0
- Quantile q
- $P(x \leq q(\alpha)) = \alpha \quad | \quad F(q(\alpha)) = \alpha \quad | \quad q(\alpha) = F^{-1}(\alpha)$

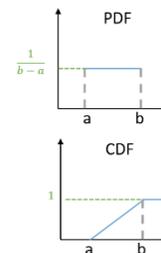
$$E(X) = \int_{-\infty}^{\infty} x * f(x) dx$$

$$Var(x) = \int_{-\infty}^{\infty} (x - E(x))^2 * f(x) dx = E(X^2) - (E(X))^2 = \sigma^2$$



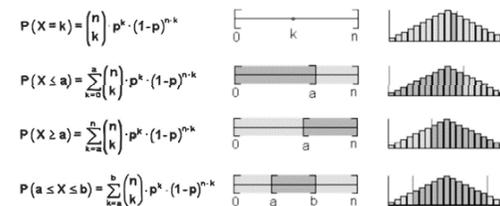
5.1 UNIFORME VERTEILUNG

- Völlige Ignoranz \rightarrow Dichte konstant auf dem Intervall \rightarrow auf ganzem Wertebereich gleiche Wa. $\rightarrow X \sim Unif([a, b])$
- PDF $\rightarrow P(X = x)$
- $f(x) = \begin{cases} \frac{1}{b-a} & \text{für } a \leq x \leq b \\ 0 & \text{falls nicht} \end{cases}$
- CDF $\rightarrow P(X \geq x) = 1 - P(X \leq x)$
- $F(x) = \begin{cases} 0 & \text{falls } x < a \\ \frac{x-a}{b-a} & \text{falls } a \leq x \leq b \\ 1 & \text{falls } x > b \end{cases}$



$$E(X) = \frac{a+b}{2} \quad Var(x) = \frac{(b-a)^2}{12} = \sigma^2$$

$$Median m(x) \rightarrow 0.5 = \frac{m-a}{b-a}$$



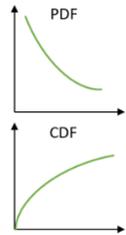
5.2 EXPONENTIALVERTEILUNG

- Zeitdauer für seltene Ereignisse → Warten ohne Gedächtnis
Wa. verändert sich nie (Bsp. Wartezeiten auf Ausfälle)
 $X \sim \text{Exp}(\lambda) \quad P(X = x) = \frac{\lambda^x}{x!} * e^{-\lambda}$
- Pois(x) hat eine bestimmte Anzahl an Events, bei Exp(x) ist es die Dauer zwischen zwei Events
- wenn die Zeiten zwischen den Ausfällen eines Systems exponential-verteilt sind, dann ist die # Ausfälle im Intervall der Länge t Poisson-verteilt

PDF

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

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



$E(x) = \frac{1}{\lambda} \quad \text{Var}(x) = \frac{1}{\lambda^2} \quad m(x) \rightarrow 0.5 * 1 - e^{-\lambda * m}$

5.3 NORMALVERTEILUNG / GAUSSVERTEILUNG

- Zufällige Messfehler, Summe von unabhängigen, gleichverteilten Zufallsvariablen $X \sim N(\mu, \sigma^2)$
- PDF
$$f(x) = \frac{1}{\sigma * \sqrt{2\pi}} * \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$
- CDF: keine konkrete Angabe → TABELLE

$E(x) = \mu \quad \text{Var}(x) = \sigma^2$
 Kurvenhöhepunkt Breite

5.4 STANDARDNORMALVERTEILUNG

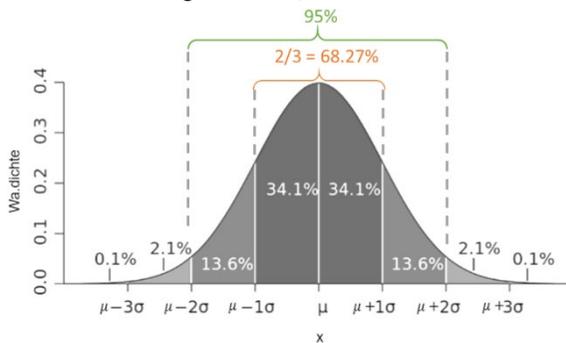
- Standardnormalverteilung $Z \sim N(0,1)$
- PDF → $\varphi(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$
- CDF → $\Phi(x) = \int_{-\infty}^x \varphi(y) dy$ (Symmetrieachse = x-Achse)
- Standardisierung = Transformation einer Zufallsvariablen, so dass $E(x) = 0 \mid \text{Var}(x) = 1 \mid \sigma(x) = 1$ gilt → Erreichung durch linear transformierte Variablen
$$g(x) = \frac{x-E(x)}{\sigma_x} = -\frac{E(x)}{\sigma_x} + \frac{1}{\sigma_x} x = a + b x$$

$$Z = g(X) = \frac{X-E(X)}{\sigma_x}$$
- Standardisieren: $P(X \leq x) = P\left(z \leq \frac{x-\mu}{\sigma}\right) = \Phi$

- 1.) Gegeben: Stelle(x) → cdf gesucht, nur für Wa. kleiner als
 $P(X \geq \mu) \rightarrow$ Ausgangslage
 $P(X \leq -x) = 1 - P(X \leq x) = P(X \geq x)$
 $z = \frac{x-\mu}{\sigma} \rightarrow$ in TABELLE nachschauen
- 2.) Gegeben: Wa. P → Stelle(x) gesucht, Wa. in TABELLE nachschauen = $z \rightarrow x = z * \sigma + \mu$

Bsp. Gegeben: $X \sim N(\mu = -85, \sigma^2 = 0.09)$
 Behauptung: $P(X \leq -85.608) = 0.0091$
 $\sigma = \sqrt{0.09} = 0.3 \mid z = \frac{x-\mu}{\sigma} = \frac{(-85.608)-(-85)}{0.3} = -2.86 \rightarrow$ Tabelle
 Wert: 0.9909, weil $P(X \geq \mu)$ sein soll → $1 - 0.9909 = 0.0091$

- Prozentzahlen: wie viele Messwerte darin enthalten
- Normalverteilung erkennbar, wenn QQ-Plot = Gerade ergibt



5.5 LOGNORMAL-VERTEILUNG

- Für positive Zufallsvariablen verwendet → wenn Transformation auf eine Normalverteilung führt = X-Lognormal-verteilt
 $Y \sim N(\mu, \sigma^2), X = \exp(Y) \rightarrow \mu \in \mathbb{R}, \sigma^2 \in \mathbb{R}^2$
- Lognormal-Verteilung nicht mehr symmetrisch
 $E(X) = \exp\left(\mu, \frac{\sigma^2}{2}\right) > \exp(E(Y))$

6. STATISTIK FÜR STICHPROBEN

Vorzeichen-Test, Wilcoxon-Test, Mann-Whitney-Test

- (Punkt-)Schätzungen
$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i \mid \hat{\sigma}_X^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

- T-Test besser als Z-Test

Test	Annahmen			
	σ_x bekannt	$x_i \sim N$	Symmetrisch	i.i.d.
z	X	X	X	X
t		X	X	X

6.1 Z-TEST

- 1.) Basierend auf mehreren Beobachtungen $Z \sim N(\mu, \sigma_x^2)$
 $\sigma_x =$ bekannt $\mid \sigma_{\bar{X}_n} =$ berechnen $\mid X_i =$ kontinuierlich
- 2.) Hypothese betrachten
 $H_0: \mu_0 = ? \mid H_A = \mu_A \neq \mu_0, \mu_A > \mu_0, \mu_A < \mu_0$
Einseitiger oder zweiseitiger Test?
- 3.) Teststatistik
$$Z = \frac{\bar{x}_n - \mu_0}{\frac{\sigma_{\bar{X}_n}}{\sigma_x}} = \frac{\sqrt{n} * (\bar{x}_n - \mu_0)}{\sigma_x} = \frac{\text{beobachtet-erwartet}}{\text{Standardfehler}} \mu \rightarrow \text{meist } 0$$
- 4.) Signifikanzniveau & Verwerfungsbereich → TABELLE
 $H_A: \mu_A \neq \mu_0 \quad K = (-\infty, -z_{1-\frac{\alpha}{2}}] \cup [z_{1-\frac{\alpha}{2}}, \infty)$
 $H_A: \mu_A < \mu_0 \quad K = (-\infty, -z_{1-\alpha}]$
 $H_A: \mu_A > \mu_0 \quad K = [z_{1-\alpha}, \infty)$
Bsp. $Z_{0.985}, 0.985$ in Tabelle suchen → Spalte und Zeile = Grenze
- 5.) Testentscheid & Vertrauensintervall
Beobachteter Wert (in Z einsetzen) im Verwerfungsbereich?
Ja → H_0 verwerfen \mid Nein → H_0 nicht verwerfen

$H_A: \mu_A \neq \mu_0 \quad VI = \left[\bar{X}_n \pm z_{1-\frac{\alpha}{2}} * \frac{\sigma_x}{\sqrt{n}} \right]$
 $H_A: \mu_A < \mu_0 \quad VI = \left(-\infty, \bar{X}_n - z_{1-\alpha} * \frac{\sigma_x}{\sqrt{n}} \right]$
 $H_A: \mu_A > \mu_0 \quad VI = \left[\bar{X}_n - z_{1-\alpha} * \frac{\sigma_x}{\sqrt{n}}, \infty \right)$

6.2 T-TEST

1.) $T \sim N(\mu, \sigma_x^2)$

$\sigma_x =$ unbekannt, $\sigma_{\bar{x}}$ = geschätzt, $X_i =$ kontinuierlich

$$E(T) = 0 \mid Var(T) = \frac{n}{n-2} \mid \hat{\sigma}_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2$$

2.) Hypothese betrachten

$H_0: \mu_0 = ? \mid H_A = \mu_A \neq \mu_0, \mu_A > \mu_0, \mu_A < \mu_0$
Einseitiger oder zweiseitiger Test?

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

$$T = \frac{\bar{x}_n - \mu_0}{\frac{\hat{\sigma}_{\bar{x}_n}}{\sqrt{n}}} = \frac{\sqrt{n}(\bar{x}_n - \mu_0)}{\hat{\sigma}_x} = \frac{\text{beobachtet-erwartet}}{\text{Standardfehler}} \mid \sigma_{\bar{x}_n} = \frac{\hat{\sigma}_x}{\sqrt{n}}$$

$\mu_0 \rightarrow$ meist 0

P-Wert (bei $H_A: \mu \neq \mu_0$) $P(|T| > |t|) = 2 * (1 - F_{t_{n-1}}(\frac{\sqrt{n}|\bar{x}_n - \mu_0|}{\hat{\sigma}_x}))$

4.) Signifikanzniveau & Verwerfungsbereich → TABELLE

Quantile in Tabelle: $T \sim t_{13}, P(T \leq t) = 0.95 \rightarrow T_{13}, Sp 0.95$
 $T \sim t_{13}, P(X \leq 2.659) \rightarrow T_{13}, 2.659$ suchen & Spalte ablesen

$H_A: \mu_A \neq \mu_0 \quad K = (-\infty, -t_{n-1; 1-\frac{\alpha}{2}}] \cup [t_{n-1; 1-\frac{\alpha}{2}}, \infty)$

$H_A: \mu_A < \mu_0 \quad K = (-\infty; -t_{n-1; 1-\frac{\alpha}{2}}]$

$H_A: \mu_A > \mu_0 \quad K = [t_{n-1; 1-\frac{\alpha}{2}}; \infty)$

- Je kleiner degrees of freedom (df) bei der Verteilung t_{df} , desto wahrscheinlicher sind Werte mit grossem Absolutbetrag

5.) Testentscheid & Vertrauensintervall

- Beobachteter Wert (in T einsetzen) im Verwerfungsbereich?
Ja → H_0 verwerfen

Nein → H_0 nicht verwerfen: falls $|\frac{\sqrt{n} * |\bar{x}_n - \mu_0|}{\hat{\sigma}_x}| < t_{n-1; 1-\frac{\alpha}{2}}$

$H_A: \mu_A \neq \mu_0 \quad VI = [\bar{X}_n \pm t_{n-1; 1-\frac{\alpha}{2}} * \frac{\hat{\sigma}_x}{\sqrt{n}})$

$H_A: \mu_A < \mu_0 \quad VI = (-\infty, \bar{X}_n + t_{n-1; 1-\frac{\alpha}{2}} * \frac{\hat{\sigma}_x}{\sqrt{n}})$

$H_A: \mu_A > \mu_0 \quad VI = [\bar{X}_n + t_{n-1; 1-\frac{\alpha}{2}} * \frac{\hat{\sigma}_x}{\sqrt{n}}, \infty)$

• Faustregel für minimale Stichprobengrösse n

$n \geq 4 * \frac{\hat{\sigma}_x^2}{\delta^2} \mid \delta = \frac{1}{2}$ von der totalen Breite des VI

6.3 ZWEI STICHPROBEN: GEPAART / UNGEPAART

• Gepaarte Stichprobe:

Direkte Zuordnung → Stichprobengrösse dieselbe

Bsp. Zwillinge, Messung vorher/nachher

1.) Differenz bilden $u_i = x_i - y_i \rightarrow (i = 1, \dots, n)$ x_i, y_i sind abhängig

2.) Teststatistik: $T_u = \frac{\sqrt{n} * (u_i - 0)}{\hat{\sigma}_u}$

Falls Daten normalverteilt → T-Test

• Ungepaarte Stichprobe

Paarung nicht möglich

x_i, y_i unabhängig

$m \neq n \rightarrow$ oft, nicht immer

Bsp. Messungen mit Methode A oder B, aber nicht mit beiden gleichzeitig → kein eindeutiger Zusammenhang

Gepaart	Ungepaart
Gleichgrosse Stichproben	Unterschiedlich
Klare Zuordnung	Keine Zuordnung
Differenz der Paare	
Mehr Macht	Weniger Macht

6.4 ZWEI STICHPROBEN T-TEST MIT GLEICHEN VARIANZEN

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

Arithmetisches Mittel: $\bar{X}_n = \frac{1}{n} \sum X_i \mid \bar{Y}_m = \frac{1}{m} \sum Y_i$

Empirische Varianz: $\hat{\sigma}_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

1.) x_1, \dots, x_n i.i.d. $\sim N(\mu_x, \sigma^2) \mid y_1, \dots, y_m$ i.i.d. $\sim N(\mu_y, \sigma^2)$

2.) Hypothese betrachten

$H_0: \mu_x = \mu_y \mid H_A = \mu_x \neq \mu_y, \mu_x > \mu_y, \mu_x < \mu_y$

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

$$S_{Pool}^2 = \frac{1}{n+m-2} ((n-1)\hat{\sigma}_x^2 + (m-1)\hat{\sigma}_y^2) = \frac{1}{n+m-2} (\sum_{i=1}^n (X_i - \bar{X}_n)^2 + \sum_{i=1}^m (Y_i - \bar{Y}_m)^2)$$

$$T = \frac{(\bar{X}_n - \bar{Y}_m) - (\mu_x - \mu_y)}{S_{Pool} * \sqrt{\frac{1}{n} + \frac{1}{m}}} \mu_x, \mu_y \rightarrow \text{meist } 0, \sqrt{S_{Pool}^2} !$$

4.) Signifikanzniveau & Verwerfungsbereich → TABELLE

$H_A: \mu_x \neq \mu_y \quad K = (-\infty; -t_{n+m-2; 1-\frac{\alpha}{2}}] \cup [t_{n+m-2; 1-\frac{\alpha}{2}}; \infty)$

$H_A: \mu_x < \mu_y \quad K = (-\infty; -t_{n+m-2; 1-\alpha}]$

$H_A: \mu_x > \mu_y \quad K = [t_{n+m-2; 1-\alpha}; \infty)$

5.) Testentscheid: Beobachteter Wert in Verwerfungsbereich?

- Ja → H_0 verwerfen | Nein → H_0 nicht verworfen
- $(\mu_x - \mu_y)$ ist häufig unter H_0 gleich 0 → wird weggelassen
- $S_{Pool}^2 = \hat{\sigma}^2$
- Bei gepaarten Stichproben auch ungepaarter Test durchführbar → umgekehrt nicht möglich

6.5 WELCH-TEST; ZWEI STICHPR. T-TEST UNGLEICHE VARIANZEN

- Varianzen in beiden Gruppen müssen nicht gleich sein → aber kleinere Macht → an Prüfung gleiche Varianzen annehmen

6.6 MULTIPLES TESTEN – BONFERRONI-KORREKTUR

- Liste von wirklich "wichtigen" Resultaten generieren
- Anpassung der Signifikanzniveaus für mehrere zusammenhängende Tests
 $\alpha_{adj} = \frac{\alpha}{j} \mid \text{korr. } \alpha - \text{Niveau für jeden Test} = \frac{\text{gewünschtes } \alpha}{\# \text{ Test}}$
- Wenn mehrere Tests durchgeführt werden bei verschiedenen Gruppen gibt es eine Alpha-Fehler-Kumulierung
Alpha-Fehler-Risiko: $1 - (1 - \alpha)^j$

Bsp. 4 Gruppen Mittelwerte vergleichen

alle $\alpha = 5\% \rightarrow 1 - (1 - 0.05)^6 = 0.265$

(6 Tests, weil alle miteinander vergleichen)

d.h auch wenn kein signifikanter Unterschied im Mittelwert würde zu 26.5% ein Test fälschlicherweise signifikant werden
Daher Bonferroni-Korrektur anwenden, damit das Gesamt α für alle Tests zusammen bei max. 5% liegt

$\alpha_{adj} = \frac{0.05}{6} = 0.0083 \rightarrow$ für jeden Test $\alpha = 0.0083$ einsetzen

- Nachteil: korrigiert α -Niveau stark nach unten → Risiko für Fehler 2. Art steigt

6.7 HYPOTHESETEST VERGLEICHE

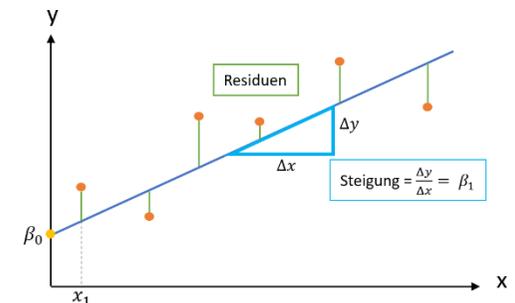
Test	Was wird getestet	Bsp.
Binomialtest	Anteil in Gruppe gleich π_0 ?	Wirksamkeit von Medikament
t-Test, eine Stichprobe	Erwartungswert in Gruppe gleich μ_0 ?	Füllmenge in Getränkeflaschen
t-Test, 2 gepaarte Stichproben	Erwartungswert in beiden Gruppen gleich?	Reaktionszeit von Haupt- und Nebenhand
t-Test, 2 ungepaarte Stichproben	Erwartungswert in beiden Gruppen gleich?	Aktivität von Gen XY bei Gesunden und Kranken

7. REGRESSION

Multiple-Linear-Regression

- Voraussagen basierend auf Messdaten machen / gerade fitten
- $Y = \beta_0 + \beta_1 * x + \varepsilon \mid \varepsilon \sim N(0, \sigma^2)$
- $\varepsilon =$ zufälliger Fehler / Fehler-Variable / Rausch-Terme
Zusammenhang zwischen der erklärenden und der Ziel- Variablen nicht exakt
- $X_i =$ erklärende Variable, deterministisch, wurde bestimmt
- $Y_i =$ Zielvariable, Zufallsvariable von ε abhängig
- Residuen = Vertikaler Abstand zwischen Punkt und Geraden
 $R_i = Y_i - (\beta_0 + \beta_1 * x_i) (i = 1, \dots, n)$
- Einfache lineare Regression: nur eine x-Variable vorkommend
- Linear = wenn Parameter β in der Ableitung verschwindet

$\varepsilon_i \sim N(0, \sigma^2) \mid E(\varepsilon_i) = 0 \mid Var(\varepsilon_i) = \sigma^2$



7.1 PARAMETERSCHÄTZUNGEN

METHODE DER KLEINSTEN QUADRATE

- Schätzungen haben keinen systematischen Fehler, sondern streuen nur um den wahren Parameter herum
- Methode der kleinsten Quadrate und Maximum Likelihood Methode sind äquivalent
- β_0 und β_1 sind Minimierer von $\sum_{i=1}^n (Y_i - (\beta_0 + \beta_1 * x_i))^2$
- Approximation $R_i \approx \epsilon_i$
- Parameterschätzung

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (Y_i - \bar{Y}_n)(x_i - \bar{x}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \quad | \quad \hat{\beta}_0 = \bar{Y}_n - \hat{\beta}_1 * \bar{x}_n$$
- Varianzschätzung $\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n R_i^2$
- Standardfehler $s.e.(\hat{\beta}_1) = \sqrt{\widehat{Var}(\hat{\beta}_1)} = \frac{\hat{\sigma}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x}_n)^2}}$
- Teststatistik $T = \frac{\hat{\beta}_1 - 0}{s.e.(\hat{\beta}_1)} = \frac{\text{beobachtet-erwartet}}{\text{geschätzter Standardfehler}}$
- Verwerfungsbereich $K = (-\infty, -t_{n-2; 1-\frac{\alpha}{2}}] \cup [t_{n-2; 1-\frac{\alpha}{2}}, \infty)$
- Vertrauensintervall

$$VI = [\hat{\beta}_1 - s.e.(\hat{\beta}_1)t_{n-2; 1-\frac{\alpha}{2}}; \hat{\beta}_1 + s.e.(\hat{\beta}_1)t_{n-2; 1-\frac{\alpha}{2}}]$$
- Faustregel 95%-VI $[\hat{\beta}_1 - 2 * s.e.(\hat{\beta}_1); \hat{\beta}_1 + 2 * s.e.(\hat{\beta}_1)]$
- Bestimmtheitsmass = wie gut stimmt das Regressionsmodell mit den tatsächlichen Beobachtungen überein

$$R^2 = \frac{SS_R}{SS_Y} = \hat{\rho}_{Y\hat{Y}}^2 \rightarrow \text{je näher } R^2 \text{ an } 1 \text{ desto besser Modell}$$

- Regressionsgerade anpassen $\rightarrow \beta_0$ und β_1 berechnen
- T-Test \rightarrow hat x einen Einfluss auf Y? $\rightarrow H_0: \beta_1 = 0$ wenn nicht signifikant \rightarrow x hat keinen Einfluss auf Y
- T-Test \rightarrow Regression durch Nullpunkt? $\rightarrow H_0: \beta_0 \neq 0$ wenn nicht signifikant \rightarrow kleineres Modell mit Regression durch Nullpunkt benutzen
- Vertrauensintervall β_0 und β_1 & Bestimmtheitsmass R^2
- Modell-Voraussetzungen mit Residuenanalyse überprüfen

7.2 GENERALISED LINEAR MODELS (GLMS)

- Zusammenhang zwischen erklärenden Variablen und Parametern einer Verteilung feststellen
- Störgrösse bei klassischen linearen Modellen normalverteilt bei GLMs nicht notwendig
- Binomialregression / logistische Regression $Y \sim Bin(n, p(x))$

$$P(x) = \frac{e^{\beta_0 + \beta_1 * x}}{1 + e^{\beta_0 + \beta_1 * x}} \quad | \quad \log\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + \beta_1 * x$$

Bsp. Bei welcher Dosis ist Genesungswahrscheinlichkeit 80%?
- Poissonregression $Y \sim Pois(\lambda(x))$

$$\lambda(x) = e^{\beta_0 + \beta_1 * x} \quad | \quad \log(\lambda(x)) = \beta_0 + \beta_1 * x$$

Bsp. Morgen -5°C, 95%-Quantil der Unfälle morgen?

- Lineare Regression $Y \sim N(\mu(x), \sigma^2)$
 Einfache Regression: $\mu(x) = \beta_0 + \beta_1 x + \beta_2 x^2$
 Multiple Regression: $\mu(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$

8. QUALITÄTSTEST - RESIDUENANALYSE

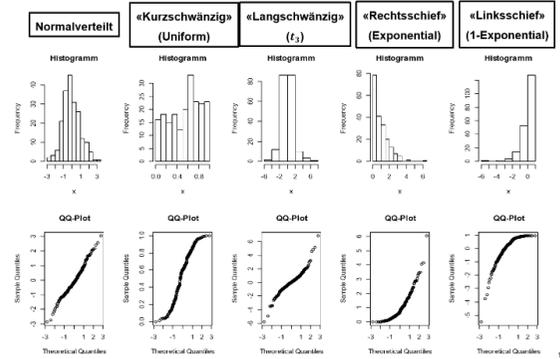
- Für die Überprüfung der Modell-Voraussetzungen für die einfache lineare Regression
- Residuen messen den vertikalen Abstand zwischen beobachtetem Punkt und gefitteter Gerade
- $E(E_i) = 0 \rightarrow E(Y_i) = \beta_0 + \beta_1 x_i \rightarrow$ kein systematischer Fehler im Modell \rightarrow Überprüfung im Tukey-Anscombe Plot
- E_1, \dots, E_n i. i. d. | $Corr(E_i, E_j) = 0 (i \neq j) \rightarrow$ Fehler müssen voneinander unabhängig sein, keine seriellen Korrelationen vorhanden
 \rightarrow Überprüfen mit QQ-Plot und Tukey-Anscombe (Fehler müssen ebenfalls gleichverteilt sein)
- E_1, \dots, E_n i. i. d. $N(0, \sigma^2) \rightarrow$ Fehler sind normalverteilt
 \rightarrow Überprüfen mit Normalplot
- Übergang zwischen Verletzen und Einhalten des Modells ist fließend \rightarrow leichte Verletzungen weniger schlimm und Ergebnisse doch noch brauchbar

8.1 QQ-PLOT

- Histogramm transformieren \rightarrow Hilfe, für Vergleich ob Modell und Verteilung funktionieren
- Q-Q-Plot (Quantil-Quantil Plot) = empirische Quantile werden gegen die theoretischen Quantile der Modell-Verteilung aufgetragen \rightarrow wenn Beobachtungen gemäss Modell, dann Punkte auf Winkelhalbierenden von $y = x$

$$\alpha = \frac{0.5}{n}, \frac{1.5}{n}, \dots, \frac{(n-0.5)}{n}$$
- Normal-Plot = Überprüfung ob Normalverteilung vorliegt

$$q(\alpha) = \mu + \sigma \Phi^{-1}(\alpha)$$

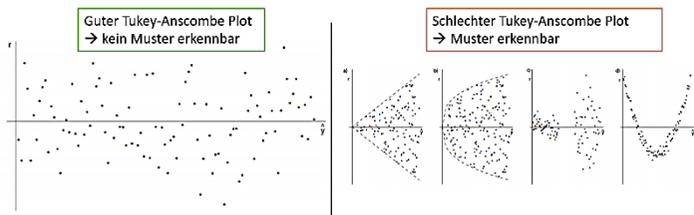


- Punkte sollen auf der Geraden mit Achsenabschnitt μ und Steigung σ liegen = dann gut

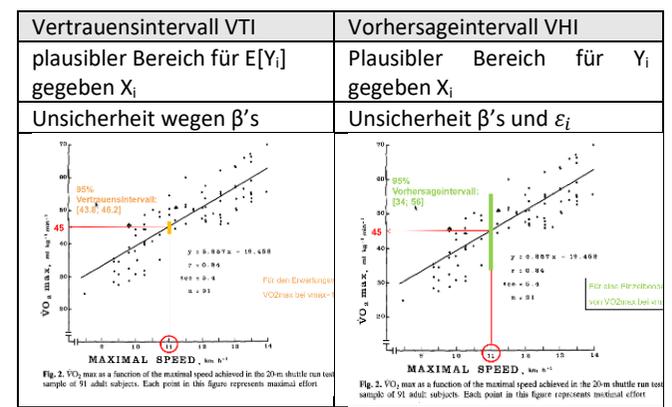
8.2 TUKEY ANSCOMBE PLOT

- Plot der Residuen r_i (Y-Achse) gegen die angepassten Werte \hat{y}_i (X-Achse)
- Idealfall: gleichmässige Streuung der Punkte um Null \rightarrow konstante Breite, uniforme Streuung innerhalb des Bandes
 Abweichungen: kegelförmiges Anwachsen der Streuung mit \hat{y}_i
- Möglicherweise kann man Zielvariable logarithmieren

$$\log(Y_i) = \beta_0 + \beta_1 x_i + \epsilon_i$$



8.3 VERTRAUENSINTERVALL UND VORHERSAGEINTERVALL VON Y



8.4 R-OUTPUT

Modell $Y_i = \beta_0 + \beta_1 x_1 + E_i ; E_i \sim N(0, \sigma^2)$ i. i. d.
 $Y_i = -19.5 + 5.9x_1 + E_i ; E_i \sim N(0, 5.4^2)$ i. i. d.

Residuals:	Min	1Q	Median	3Q	Max
	-10.2	-4.5	-0.2	4.7	12.0

Coefficients:	Estimate	Std. Error	t value	Pr(> t)
β_0 Intercept Achsenabschnitt	-19.5	4.7	-4.1	8.5e-05
β_1 x	5.9	0.4	14.3	<2e-16

Signif. Codes: 0 0.001 0.01 0.05 0.1 1

Residual standard error: 5.4 on 89 degrees of freedom
 Multiple R-squared: 0.6981, Adjusted R-squared: 0.8948
 F-statistic: 205.8 on 1 and 89 DF, p-Value <2.2e-16

- erwarteter Achsenabschnitt
- Falls man x um eine Einheit erhöht, sagt das Modell voraus, dass sich y um den Wert 5.9 (Steigung) erhöht
- Wenn es um α geht \rightarrow Estimate miteinander vergleichen
- Standardfehler von $\hat{\beta}_1 (= \sigma_{\hat{\beta}_1})$
 Approx. 95%-VI: $5.9 \pm 2 * 0.4$
 Exaktes, 95%-VI: $5.9 \pm 1.99 * 0.4$
 $\rightarrow t_{89, 0.975} = 1.99 \rightarrow$ TABELLE
- Freiheitsgrade: $df = n - 2$
- t Value: Zwischen ± 2 , wenn \uparrow , dann Fläche klein
 \rightarrow d.h. P-Wert astronomisch klein
- Beobachtete Teststatistik
 Im Test $H_0: \beta_0 = 0$ vs. $H_1: \beta_1 \neq 0$
 $\frac{5.9}{0.4} = 14.3$
- P-Wert: Angenommen $\beta_1 = 0$, wie Wa. ist es diese Beobachtung oder etwas Extremes zu erhalten.