

# Statistik

D-UWIS & D-ERDW & D-AGRL

Sommersemester 2005

basierend auf Skript von H.-R. Künsch

Peter Bühlmann  
Seminar für Statistik  
ETH Zürich



# Inhaltsverzeichnis

<b>1</b>	<b>Einführung (Stahel, Kap. 1)</b>	<b>7</b>
<b>2</b>	<b>Modelle für Zählraten</b>	<b>9</b>
2.1	Einleitung (Stahel, Kap. 4.1)	9
2.2	Diskrete Wahrscheinlichkeit (Stahel, Kap. 4.2, 4.6)	9
2.2.1	Der Begriff der Unabhängigkeit	11
2.3	Zufallsvariable (Stahel, Kap. 4.3, 4.4)	11
2.4	Binomialverteilung (Stahel Kap. 5.1)	13
2.5	Kennzahlen einer Verteilung (Stahel Kap. 5.3)	14
2.5.1	Kumulative Verteilungsfunktion	16
2.6	Poissonverteilung (Stahel Kap. 5.2)	17
2.6.1	Poisson-Approximation der Binomial-Verteilung	17
<b>3</b>	<b>Statistik für Zählraten</b>	<b>19</b>
3.1	Drei Grundfragestellungen der Statistik (Stahel Kap. 7.1)	19
3.2	Schätzung, statistischer Test und Vertrauensintervall bei Binomial-Verteilung (Stahel Kap. 7.2, 8.2, 9.1, 9.2)	19
3.2.1	(Punkt-)Schätzung	20
3.2.2	Statistischer Test	20
3.2.3	Vertrauensintervall	24
3.3	Schätzung, Test und Vertrauensintervall bei Poisson-Verteilung (Stahel, Kap. 7.2, 8.1, 9.1)	25
3.3.1	(Punkt-)Schätzung	25
3.3.2	Statistischer Test	25
3.3.3	Vertrauensintervall	26
<b>4</b>	<b>Modelle und Statistik für Messdaten</b>	<b>27</b>
4.1	Einleitung	27
4.2	Deskriptive Statistik (Stahel, Kap. 2 und 3.1, 3.2)	28
4.2.1	Kennzahlen	28
4.2.2	Grafische Methoden	29
4.3	Stetige Zufallsvariablen und Wahrscheinlichkeitsverteilungen (Stahel, Kap. 6.1 – 6.4, 11.2)	31
4.3.1	(Wahrscheinlichkeits-)Dichte	32
4.4	Wichtige stetige Verteilungen (Stahel, Kap. 6.2, 6.4, 6.5, 11.2)	33
4.4.1	Uniforme Verteilung	34
4.4.2	Exponential-Verteilung	34
4.4.3	Normal-Verteilung (Gauss-Verteilung)	35
4.4.4	Transformationen	36

4.4.5	Analogien zwischen Modellen und Daten . . . . .	38
4.4.6	Überprüfen der Normalverteilungs-Annahme . . . . .	38
4.5	Funktionen von Zufallsvariablen, Fehlerfortpflanzung (Stahel, Kap. 6.8 – 6.11)	39
4.6	Statistik für eine Stichprobe (Stahel, Kap. 8.3 – 8.5, 9.3) . . . . .	41
4.6.1	(Punkt-) Schätzungen . . . . .	41
4.6.2	Tests für $\mu$ . . . . .	42
4.6.3	Vertrauensintervall für $\mu$ . . . . .	45
4.6.4	Tests für $\mu$ bei nicht-normalverteilten Daten . . . . .	45
4.7	Tests bei zwei unabhängigen Stichproben (Stahel, Kap. 8.8) . . . . .	47
4.7.1	Gepaarte und ungepaarte Stichproben . . . . .	47
4.7.2	Gepaarte Tests . . . . .	48
4.7.3	Ungepaarte Tests . . . . .	48
4.7.4	Zwei-Stichproben t-Test bei gleichen Varianzen . . . . .	48
4.7.5	Weitere Zwei-Stichproben-Tests . . . . .	50
<b>5</b>	<b>Regression</b>	<b>53</b>
5.1	Korrelation und empirische Korrelation . . . . .	53
5.1.1	Die empirische Korrelation . . . . .	53
5.2	Einfache lineare Regression . . . . .	54
5.2.1	Das Modell der einfachen linearen Regression . . . . .	54
5.2.2	Parameterschätzungen . . . . .	55
5.2.3	Tests und Konfidenzintervalle . . . . .	56
5.2.4	Das Bestimmtheitsmass $R^2$ . . . . .	58
5.2.5	Allgemeines Vorgehen bei einfacher linearer Regression . . . . .	59
5.2.6	Residuenanalyse . . . . .	59
5.3	Multiple lineare Regression . . . . .	61
5.3.1	Das Modell der multiplen linearen Regression . . . . .	62
5.3.2	Parameterschätzungen und t-Tests . . . . .	63
5.3.3	Der F-Test . . . . .	64
5.3.4	Das Bestimmtheitsmass $R^2$ . . . . .	64
5.3.5	Residuenanalyse . . . . .	64
5.3.6	Strategie der Datenanalyse: ein abschliessendes Beispiel . . . . .	64

## Aufbau der Vorlesung

Die Vorlesung behandelt zuerst die Wahrscheinlichkeitsrechnung und Statistik für diskrete Variablen, welche Werte zum Beispiel in  $\{0, 1\}$ , in  $\mathbb{N}_0 = \{0, 1, 2, \dots\}$  oder in  $\mathbb{Z} = \{\dots, -1, 0, 1, \dots\}$  annehmen.

Danach werden die erarbeiteten Konzepte auf stetige Variablen übertragen, mit Wertebereichen, mit Wertebereichen zum Beispiel in  $\mathbb{R}$  oder  $[0, 1]$ . Deshalb ist der Aufbau leicht repetitiv, was sich aber in vorigen Jahren gut bewährt hat.

Schlussendlich wird auf komplexere Modellierung anhand der multiplen Regressions-Analyse eingegangen.



# Kapitel 1

## Einführung (Stahel, Kap. 1)

Die Bedeutung der Statistik liegt, für viele Wissenschaften, in der Fähigkeit

*verallgemeinernde Schlüsse von Daten ("Stichproben") auf zukünftige Daten oder umfassendere Populationen zu machen.*

Insbesondere wird dabei berücksichtigt, dass

*alle Daten gewissen Schwankungen unterworfen sind.*

Um dies zu quantifizieren benützt man

*Modelle und Gesetze der Wahrscheinlichkeitstheorie.*

Der Zufall gehorcht gewissen Gesetzen aus der Wahrscheinlichkeitstheorie, welche genauso zuverlässig sind wie andere Naturgesetze. Ob die Welt tatsächlich zufällig ist, oder ob der Zufall nur eine Terminologie ist für all diejenigen deterministischen Faktoren, die wir unmöglich alle in Betracht ziehen können, ist für viele Betrachtungen unwesentlich.





# Kapitel 2

## Modelle für Zählraten

### 2.1 Einleitung (Stahel, Kap. 4.1)

Ein **Wahrscheinlichkeitsmodell** beschreibt, welche Ergebnisse eines Experiments möglich sind und welche Chancen die verschiedenen Ergebnisse haben. Wir behandeln hier diskrete Wahrscheinlichkeitsmodelle, bei denen die möglichen Ergebnisse endlich oder “abzählbar” (z.B. in  $\mathbb{N}_0 = \{0, 1, 2, \dots\}$ ) sind. Ein Wahrscheinlichkeitsmodell erlaubt zudem mittels Simulation andere mögliche Daten zu erzeugen und so zu differenzieren, welche Schwankungen noch plausibel sind und welche nicht.

Die Experimente sind hier immer als **Zufallsexperiment** zu verstehen:

Zufallsexperiment =

Experiment, dessen Ausgang auch nicht von einem Orakel exakt vorhersagbar ist

### 2.2 Diskrete Wahrscheinlichkeit (Stahel, Kap. 4.2, 4.6)

Für die Beschreibung von Zufallsexperimenten benützen wir ein Wahrscheinlichkeitsmodell. Dieses besteht aus den folgenden Komponenten:

- **Grundraum**  $\Omega$
- **Elementarereignis**  $\omega$
- **Wahrscheinlichkeit**  $P$

Der Grundraum und die Elementarereignisse sind wie folgt:

$$\Omega = \underbrace{\{\text{mögliche Elementarereignisse } \omega\}}_{\text{mögliche Ausgänge/Resultate}}$$

*Beispiel:* 2-maliges Werfen einer Münze

$\Omega = \{KK, KZ, ZK, ZZ\}$  wobei  $K = \text{“Kopf”}$  und  $Z = \text{“Zahl”}$  bezeichnet

Elementarereignis: zum Beispiel  $\omega = KZ$

Unter einem **Ereignis**  $A$  versteht man eine Teilmenge von  $\Omega$ :

$$\text{Ereignis } A \subset \Omega$$

*Beispiel (Forts.):*  $A = \{\text{genau 1-mal Kopf}\} = \{KZ, ZK\}$ .

Die Operationen der Mengenlehre (Komplement, Vereinigung, Durchschnitt) haben eine natürliche Interpretation in der Sprache der Ereignisse.

$$\begin{aligned}A \cup B &\Leftrightarrow A \text{ **oder** } B, \text{ wobei das "oder" nicht-exklusiv ist ("oder/und")} \\A \cap B &\Leftrightarrow A \text{ **und** } B \\A^c &\Leftrightarrow \text{**nicht** } A\end{aligned}$$

*Beispiel:*  $A =$  morgen scheint die Sonne,  $B =$  morgen regnet es.

$A \cup B$  bedeutet: morgen scheint die Sonne oder morgen regnet es (und dies kann auch bedeuten, dass morgen die Sonne scheint und dass es morgen regnet);  $A \cap B$  bedeutet: morgen scheint die Sonne und morgen regnet es;  $A^c$  bedeutet: morgen scheint die Sonne nicht.

Eine **Wahrscheinlichkeit** ordnet jedem Ereignis  $A$  eine Wahrscheinlichkeit  $P(A)$  zu. Dabei sind die folgenden drei grundlegenden Regeln (Axiome von Kolmogorov) erfüllt:

1. Die Wahrscheinlichkeiten sind immer nicht-negativ:  $P(A) \geq 0$
2. Das sichere Ereignis  $\Omega$  hat Wahrscheinlichkeit eins:  $P(\Omega) = 1$
3.  $P(A \cup B) = P(A) + P(B)$  falls  $A \cap B = \emptyset$ , d.h. für alle Ereignisse, die sich gegenseitig ausschließen.

Weitere Regeln können daraus abgeleitet werden, z.B.

$$\begin{aligned}P(A^c) &= 1 - P(A), \\P(A \cup B) &= P(A) + P(B) - P(A \cap B).\end{aligned}$$

Im diskreten Fall ist eine Wahrscheinlichkeit festgelegt durch die Wahrscheinlichkeiten der Elementarereignisse  $P(\{\omega\})$ :

$$P(A) = \sum_{\omega \in A} P(\{\omega\}).$$

*Beispiel (Forts.)* Für  $A =$  genau einmal Kopf  $= \{KZ, ZK\}$  hat man  $P(A) = P(KZ) + P(ZK) = 1/4 + 1/4 = 1/2$ .

Im Wesentlichen werden in der Wahrscheinlichkeitstheorie die Wahrscheinlichkeiten gewisser Ereignisse  $A$  festgelegt (auf Grund von Plausibilitäten, Symmetrieüberlegungen, wissenschaftlichen Theorien, Fachwissen und Daten) und daraus die Wahrscheinlichkeiten von gewissen anderen Ereignissen  $B$  aus den obigen Gesetzen hergeleitet.

(Die Statistik geht umgekehrt vor: aus Daten, d.h. aus der Information, dass gewisse Ereignisse eingetreten sind, versucht man Rückschlüsse auf ein unbekanntes Wahrscheinlichkeitsmodell (unbekannte Wahrscheinlichkeiten) zu machen).

#### **Interpretation von Wahrscheinlichkeiten:**

- Idealisierung der relativen Häufigkeiten bei vielen unabhängigen Wiederholungen (**frequentistisch**)
- Mass für den Glauben, dass ein Ereignis eintreten wird" (**Bayes'sch**)

### 2.2.1 Der Begriff der Unabhängigkeit

Wenn man die Wahrscheinlichkeiten  $P(A)$  und  $P(B)$  kennt, lässt sich im Allgemeinen daraus nicht  $P(A \cap B)$  berechnen.

Wenn zwischen den Ereignissen  $A$  und  $B$  kein kausaler Zusammenhang besteht (d.h. es gibt keine gemeinsamen Ursachen oder Ausschliessungen), dann definiert man:

$$A \text{ und } B \text{ heissen (stochastisch) unabhängig} \Leftrightarrow P(A \cap B) = P(A)P(B).$$

Im Falle von unabhängigen Ereignissen, vereinfachen sich viele Situation: insbesondere sehen wir, dass in diesem Fall  $P(A \cap B)$  aus  $P(A)$  und  $P(B)$  berechnet werden kann. In der Praxis wird oftmals aufgrund von Plausibilitätsüberlegungen **deklariert**, dass zwei Ereignisse unabhängig sind.

Bei mehreren Ereignissen  $A_1, \dots, A_n$  bedeutet Unabhängigkeit, dass zum Beispiel

$$\begin{aligned} P(A_1 \cap A_2) &= P(A_1)P(A_2), \\ P(A_1 \cap A_2 \cap A_3) &= P(A_1)P(A_2)P(A_3). \end{aligned}$$

Die allgemeine Formel für unabhängige Ereignisse lautet:

$$P(A_{i_1} \cap \dots \cap A_{i_k}) = P(A_{i_1}) \cdots P(A_{i_k}) \text{ für jedes } k \leq n \text{ und jedes } 1 \leq i_1 < \dots < i_k \leq n.$$

### 2.3 Zufallsvariable (Stahel, Kap. 4.3, 4.4)

Oft sind mit einem Zufallsexperiment Zahlenwerte verknüpft, d.h. zu jedem Elementarereignis (Ergebnis)  $\omega$  gehört ein Zahlenwert  $X(\omega) = x$ .

*Beispiel:* Wert einer gezogenen Jass-Karte

$$\begin{aligned} \omega = \text{As} &\mapsto X(\omega) = 11 \\ \omega = \text{König} &\mapsto X(\omega) = 4 \\ &\vdots \\ \omega = \text{Sechs} &\mapsto X(\omega) = 0 \end{aligned}$$

In obigen Beispiel ist also  $X(\cdot)$  eine Funktion. Allgemein definiert man:

Eine **Zufallsvariable**  $X$  ist eine **Funktion**:

$$\begin{aligned} X : \quad \Omega &\rightarrow \mathbb{R} \\ \omega &\mapsto X(\omega) \end{aligned}$$

Die Funktion  $X(\cdot)$  ist nicht zufällig; das Argument  $\omega$  hingegen schon.

Die Notation  $X$  (oder auch  $Y, Z, \dots$ ) ist eher unüblich für die Bezeichnung einer Funktion: wir werden aber sehen, dass man mit Zufallsvariablen manchmal ähnlich rechnen kann wie mit gewöhnlichen Variablen  $x$  (oder auch  $y, z, \dots$ ).

Je nach Ausgang des Zufallsexperiments (d.h. von  $\omega$ ) erhalten wir einen anderen Wert  $x = X(\omega)$ : der Wert  $x$  ist **eine Realisierung** der Zufallsvariablen  $X$ . Eine Realisierung einer Zufallsvariablen ist also das Ergebnis eines Zufallsexperiments (welches mit einer Zahl beschrieben werden kann).

Eine Zufallsvariable  $X$  heisst *diskret*, falls der Wertebereich  $W = W_X$  (Menge der möglichen Werte von  $X$ ) diskret ist, d.h. abzählbar (man kann die möglichen Werte durchnummerieren). Zum Beispiel ist  $W = \{0, 1, \dots, 10\}$  endlich und deshalb diskret;  $W = \mathbb{N}_0 = \{0, 1, 2, \dots\}$  ist zwar unendlich, aber trotzdem diskret;  $W = \mathbb{R}$  ist nicht diskret (sondern kontinuierlich). Wir werden in diesem Kapitel bloss diskrete Zufallsvariablen genauer diskutieren.

## Wahrscheinlichkeitsverteilung einer Zufallsvariablen

Die Werte einer Zufallsvariablen  $X$  (die möglichen Realisationen von  $X$ ) treten mit gewissen Wahrscheinlichkeiten auf. Diese sind wie folgt definiert:

$$\begin{aligned} & \text{Wahrscheinlichkeit, dass } X \text{ den Wert } x \text{ annimmt} \\ &= P(X = x) = P(\{\omega; X(\omega) = x\}) \\ &= \sum_{\omega; X(\omega)=x} P(\omega). \end{aligned}$$

*Beispiel (Forts):*  $X =$  Wert einer gezogenen Jass-Karte

$$\begin{aligned} & \text{Wahrscheinlichkeit für Zahl 4} = P(X = 4) \\ &= P(\{\omega; \omega = \text{ein König}\}) \\ &= P(\text{Eicheln-König}) + P(\text{Rosen-König}) + P(\text{Schellen-König}) + P(\text{Schilten-König}) \\ &= 4/36 = 1/9. \end{aligned}$$

Die "Liste" von  $P(X = x)$  für alle möglichen Werte  $x$  heisst (diskrete) **(Wahrscheinlichkeits-) Verteilung** der (diskreten) Zufallsvariablen  $X$ . Zu einer Zufallsvariablen  $X$  gehört immer eine (Wahrscheinlichkeits-) Verteilung und umgekehrt:

$$\text{Zufallsvariable } X \Leftrightarrow \text{(Wahrscheinlichkeits-) Verteilung}$$

Es gilt immer

$$\sum_{\text{alle möglichen } x} P(X = x) = 1.$$

*Beispiel (Forts):*  $X =$  Wert einer gezogenen Jass-Karte  
Die Wahrscheinlichkeitsverteilung von  $X$  ist

$$\begin{aligned} P(X = 11) &= 1/9 \\ P(X = 10) &= 1/9 \\ P(X = 4) &= 1/9 \\ P(X = 3) &= 1/9 \\ P(X = 2) &= 1/9 \\ P(X = 0) &= 4/9 \end{aligned}$$

Wenn man nur an der Zufallsvariablen  $X$  interessiert ist, kann man den zu Grunde liegenden Raum  $\Omega$  vergessen, man braucht nur die Verteilung von  $X$ .

## 2.4 Binomialverteilung (Stahel Kap. 5.1)

Wir betrachten die Situation wo es um das Messen der Anzahl Erfolge (oder Misserfolge) geht. Solche Anwendungen treten z.B. auf bei der Qualitätskontrolle, Erfolg/Misserfolg bei Behandlungen (medizinisch, biologisch) oder auch bei Glücksspielen.

*Beispiel:* Werfen einer Münze

Es wird eine Münze geworfen, welche dann zufällig entweder auf Kopf (K) oder Zahl (Z) fällt.

Betrachte die Zufallsvariable  $X$  mit Werten in  $W = \{0, 1\}$ , welche folgendes beschreibt:

$$\begin{aligned} X = 0 & \quad \text{falls Münze auf Z fällt,} \\ X = 1 & \quad \text{falls Münze auf K fällt.} \end{aligned}$$

Die Wahrscheinlichkeitsverteilung von  $X$  kann durch einen einzelnen Parameter  $\pi$  beschrieben werden:

$$P(X = 1) = \pi, \quad P(X = 0) = 1 - \pi, \quad 0 \leq \pi \leq 1.$$

Fall die Münze fair ist, so ist  $\pi = 1/2$ .

### **Bernoulli( $\pi$ )-Verteilung:**

Eine Zufallsvariable  $X$  mit Werten in  $W = \{0, 1\}$  heisst Bernoulli( $\pi$ )-verteilt, falls

$$P(X = 1) = \pi, \quad P(X = 0) = 1 - \pi, \quad 0 \leq \pi \leq 1.$$

Die Bernoulli-Verteilung ist eine triviale Mathematisierung um das Eintreffen oder Nicht-Eintreffen eines Ereignis zu beschreiben.

*Beispiel (Forts.):*  $n$ -maliges Werfen einer Münze

Wir betrachten  $X =$  Anzahl K bei  $n$  unabhängigen Münzwürfen. Es ist klar, dass der Wertebereich von  $X$  die Menge  $W = \{0, 1, \dots, n\}$  ist.  $X$  lässt sich auch darstellen als Summe von unabhängigen Bernoulli-Zufallsvariablen

$$\begin{aligned} X &= \sum_{i=1}^n X_i, \\ X_i &= \begin{cases} 1 & \text{falls K im } i\text{-ten Wurf} \\ 0 & \text{falls Z im } i\text{-ten Wurf.} \end{cases} \end{aligned}$$

Die Verteilung von  $X$  im obigen Beispiel lässt sich analytisch ausrechnen. Falls  $X_1, \dots, X_n$  alle unabhängig und Bernoulli( $\pi$ )-verteilt sind, so gilt zum Beispiel:

$$\begin{aligned} P(X = 0) &= P(\text{alle } X_1 = \dots = X_n = 0) = (1 - \pi)^n, \\ P(X = 1) &= P(\text{ein } X_i = 1 \text{ und alle anderen } X_j = 0) = \binom{n}{1} \pi (1 - \pi)^{n-1}. \end{aligned}$$

Allgemein gilt die Formel der Binomial-Verteilung.

### Binomial( $n, \pi$ )-Verteilung:

Eine Zufallsvariable  $X$  mit Werten in  $W = \{0, 1, \dots, n\}$  heisst Binomial( $n, \pi$ )-verteilt, falls

$$P(X = x) = \binom{n}{x} \pi^x (1 - \pi)^{n-x}, \quad x = 0, 1, \dots, n$$

wobei  $0 \leq \pi \leq 1$  der Erfolgsparameter der Verteilung ist.

(Dabei ist  $\binom{n}{x}$  der Binomialkoeffizient, der angibt, auf wie viele Arten man  $x$  Erfolge und  $n - x$  Misserfolge anordnen kann).

Wie bereits in obigem Beispiel motiviert, beschreibt  $X$  die Anzahl Erfolge/Misserfolge (Eintreten eines bestimmten Ereignis) bei  $n$  **unabhängigen** Versuchen. Das Prädikat “unabhängig” ist wesentlich für die Korrektheit der Binomialverteilung.

*Beispiel:* Spermasexing (Tages-Anzeiger 6.12.2000)

Geschlechts-Beeinflussung von Kuhkälbern mit einer Methode, die Spermasexing genannt wird: deren Ziel ist, ein weibliches Kalb zu züchten. In einem Testlauf wurden zwölf Kühe mit Spermien besamt, die optisch nach dem Y-Chromosom sortiert wurden (d.h. mit der Spermasexing-Methode). Da die Methode nicht hundertprozentig sicher ist, können wir das als Zufallsexperiment auffassen. Sei  $X =$  Anzahl weiblicher gezüchteter Kuhkälber. Eine vernünftiges Modell ist dann:

$$X \sim \text{Binomial}(12, \pi),$$

wobei  $\pi$  unbekannt ist. Effektiv beobachtet wurden  $x = 11$  weiblich gezüchtete Kuhkälber: d.h.  $X = x = 11$  wurde tatsächlich **realisiert**.

Eigenschaften der Binomialverteilung (siehe Abb. 2.1):  $P(X = x)$  ist maximal wenn  $x$  gleich dem ganzzahligen Teil von  $(n + 1)\pi$  ist, und auf beiden Seiten von diesem Wert nehmen die Wahrscheinlichkeiten monoton ab. Wenn  $n\pi(1 - \pi)$  nicht allzu klein ist, ist die Verteilung glockenförmig.

## 2.5 Kennzahlen einer Verteilung (Stahel Kap. 5.3)

Eine beliebige (diskrete) Verteilung kann vereinfachend zusammengefasst werden durch 2 Kennzahlen, den **Erwartungswert**  $\mathcal{E}(X)$  und die **Varianz**  $\text{Var}(X)$  (oder **Standardabweichung**  $\sigma(X) = \sqrt{\text{Var}(X)}$ ).

Der Erwartungswert beschreibt die mittlere Lage der Verteilung und ist wie folgt definiert:

$$\mathcal{E}(X) = \sum_{x \in W_x} x P(X = x), \quad W_x = \text{Wertebereich von } X.$$

Die Varianz beschreibt die Variabilität der Verteilung, und die Standardabweichung deren Streuung:

$$\begin{aligned} \text{Var}(X) &= \sum_{x \in W_x} (x - \mathcal{E}(X))^2 P(X = x) \\ \sigma(X) &= \sqrt{\text{Var}(X)}. \end{aligned}$$

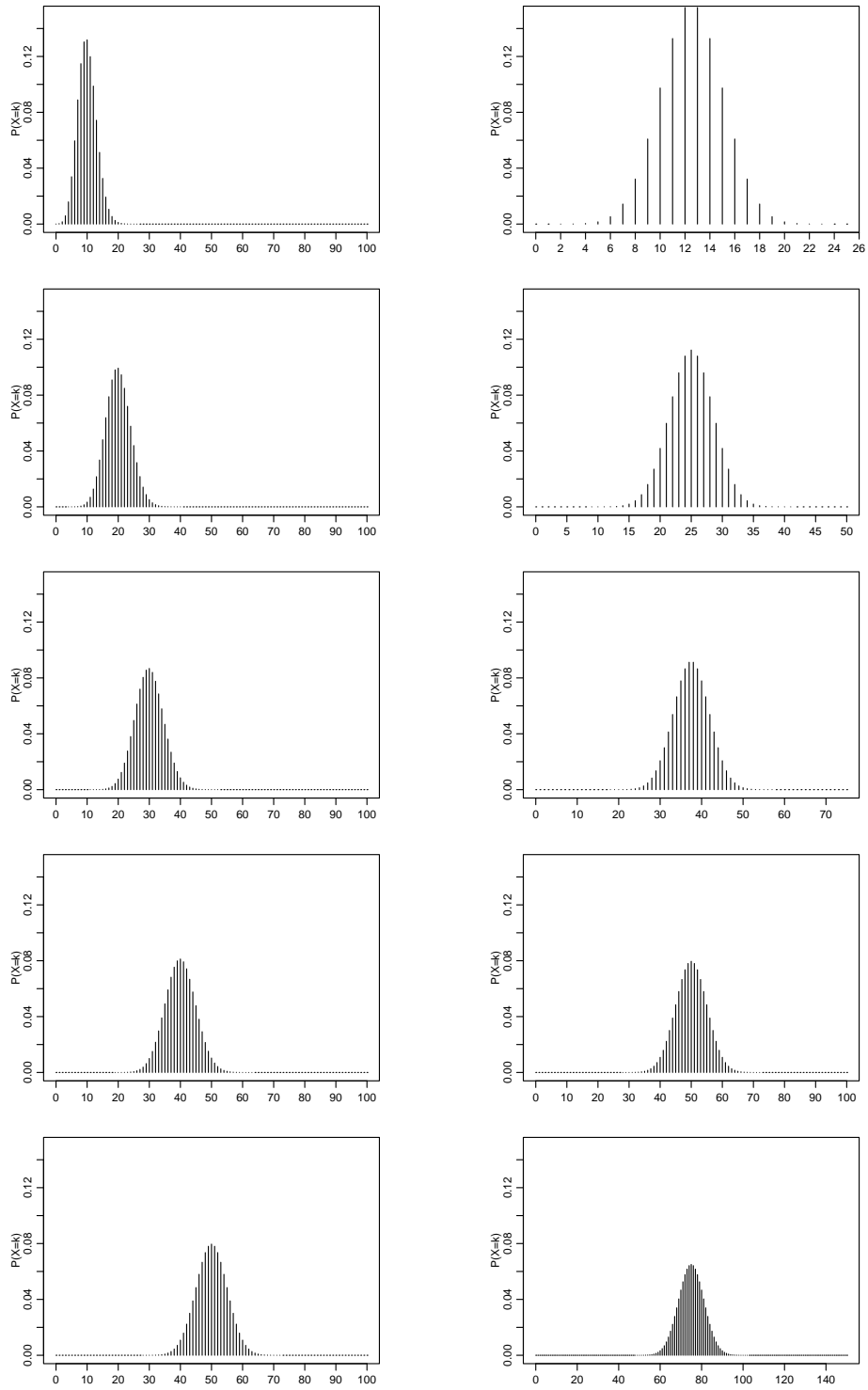


Abbildung 2.1: Die Binomialwahrscheinlichkeiten  $P(X = x)$  als Funktion von  $x$  für verschiedene  $n$ 's und  $\pi$ 's. Links ist  $n = 100$  und  $\pi = 0.1, 0.2, 0.3, 0.4, 0.5$  und rechts ist  $\pi = 0.5$  und  $n = 25, 50, 75, 100, 150$ .

Die Standardabweichung hat dieselbe Dimension wie die Messdaten: wird z.B.  $X$  in Metern ( $m$ ) gemessen, so besitzt  $\text{Var}(X)$  die Dimension Quadratmeter ( $m^2$ ) und  $\sigma(X)$  wiederum Meter ( $m$ ).

*Beispiel:* Sei  $X \sim \text{Bernoulli}(\pi)$ .

Dann:

$$\begin{aligned} \mathcal{E}(X) &= 0 \cdot P(X = 0) + 1 \cdot P(X = 1) = \pi, \\ \text{Var}(X) &= (0 - \mathcal{E}(X))^2 P(X = 0) + (1 - \mathcal{E}(X))^2 P(X = 1) = \pi^2(1 - \pi) + (1 - \pi)^2 \pi \\ &= \pi(1 - \pi), \\ \sigma(X) &= \sqrt{\pi(1 - \pi)}. \end{aligned}$$

Allgemeiner gilt für die Binomial-Verteilung (beachte, dass  $\text{Bernoulli}(\pi) = \text{Binomial}(1, \pi)$ ):

$$\begin{aligned} X &\sim \text{Bernoulli}(n, \pi), \\ \mathcal{E}(X) &= n\pi, \quad \text{Var}(X) = n\pi(1 - \pi), \quad \sigma(X) = \sqrt{n\pi(1 - \pi)}. \end{aligned}$$

### 2.5.1 Kumulative Verteilungsfunktion

Manchmal ist es geeigneter, wenn man statt der "Liste"  $P(X = x)$  (für alle  $x$ ) die sogenannte **kumulative Verteilungsfunktion** angibt:

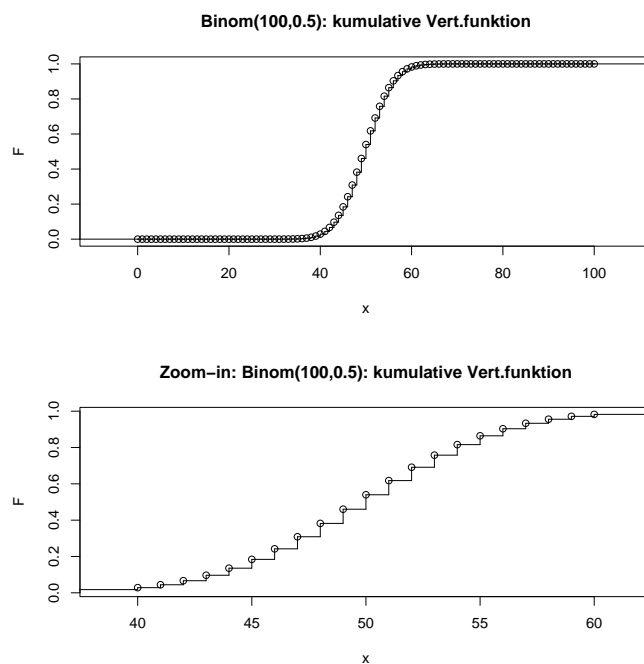


Abbildung 2.2: Kumulative Verteilungsfunktion  $F(\cdot)$  für  $X \sim \text{Binomial}(100, 0.5)$ . Unten: zoom-in für die Werte  $x \in [40, 60]$ .



$$F(x) = P(X \leq x) = \sum_{k \leq x} P(X = k).$$

Die Funktion  $F(\cdot)$  ist monoton wachsend (allerdings nicht strikt monoton) und es gilt:

$$F(-\infty) = 0, \quad F(+\infty) = 1.$$

Siehe auch Abbildung 2.2. Die Kenntnis der "Liste"  $P(X = x)$  (für alle  $x$ ) ist äquivalent zur Kenntnis der Funktion  $F(\cdot)$ , denn aus dem einen kann man das andere berechnen. Zum Beispiel gilt für  $X$  mit Wertebereich  $W_X = \{0, 1, \dots, n\}$ :  $P(X = x) = F(x) - F(x - 1)$  ( $x = 1, 2, \dots, n$ ) und  $P(X = 0) = F(0)$ .

## 2.6 Poissonverteilung (Stahel Kap. 5.2)

Der Wertebereich der Binomial( $n, \pi$ )-Verteilung ist  $W = \{0, 1, \dots, n\}$ . Falls eine Zufallsvariable nicht im vornherein einen beschränkten Wertebereich hat, so bietet sich für Zähl-daten die Poisson-Verteilung an.

Eine Zufallsvariable  $X$  mit Werten in  $\mathbb{N}_0 = \{0, 1, 2, \dots\}$  heisst Poisson( $\lambda$ )-verteilt, falls

$$P(X = x) = \exp(-\lambda) \frac{\lambda^x}{x!} \quad (x = 0, 1, 2, \dots)$$

wobei  $\lambda > 0$  ein Parameter der Verteilung ist.

Die Poisson-Verteilung ist die Standardverteilung für unbeschränkte **Zähl-daten**.

*Beispiele:* Die Poisson( $\lambda$ )-Verteilung kann bei folgenden Anwendungen gebraucht werden um die Verteilung einer Zufallsvariablen  $X$  zu modellieren:

$X =$  Anzahl Schadenmeldungen eines Versicherten pro Jahr

$X =$  Anzahl spontaner Ereignisse in einer Nervenzelle während einer Sekunde  
(via Transmitterfreisetzung an einer Synapse)

Die Kennzahlen sind wie folgt: für  $X \sim \text{Poisson}(\lambda)$ :

$$\mathcal{E}(X) = \lambda, \quad \text{Var}(X) = \lambda, \quad \sigma(X) = \sqrt{\lambda}.$$

### 2.6.1 Poisson-Approximation der Binomial-Verteilung

Betrachte  $X \sim \text{Binomial}(n, \pi)$  und  $Y \sim \text{Poisson}(\lambda)$ . Falls  $n$  gross und  $\pi$  klein mit  $\lambda = n\pi$ , dann:

$$P(X = x) = \binom{n}{x} \pi^x (1 - \pi)^{n-x} \approx P(Y = x) = \exp(-\lambda) \frac{\lambda^x}{x!} \quad (x = 0, 1, \dots, n).$$

Das heisst: für grosse  $n$  und kleine  $\pi$ :  $\text{Binomial}(n, \pi) \approx \text{Poisson}(\lambda)$  für  $\lambda = n\pi$ . Mit anderen Worten: die Poisson-Verteilung kann interpretiert werden als Verteilung für **seltene Ereignisse bei vielen unabhängigen Versuchen** (selten für einen einzelnen Fall, die Gesamt-Anzahl kann trotzdem gross sein).



# Kapitel 3

## Statistik für Zählraten

### 3.1 Drei Grundfragestellungen der Statistik (Stahel Kap. 7.1)

Die Statistik befasst sich mit dem Schliessen von einer oder mehreren Beobachtungen auf einen (oder mehrere) Parameter in einem wahrscheinlichkeitstheoretischen Modell.

*Beispiel (Forts.):* Sei  $x = 11$  die effektive Anzahl weiblicher gezüchteter Kuhkälber beim Spermasexing (vgl. Kapitel 2.4). Wir fassen  $x = 11$  auf als **Realisierung** der Zufallsvariablen  $X \sim \text{Binom}(12, \pi)$ . Wir möchten jetzt Schlüsse ziehen von der Beobachtung  $x = 11$  auf den unbekannt Parameter  $\pi$ .

**1. Grundfragestellung:** Welches ist der zu den Beobachtungen plausibelste Parameterwert? Die Antwort auf diese 1. Grundfrage heisst **(Punkt-)Schätzung**.

**2. Grundfragestellung:** Sind die Beobachtungen kompatibel (statistisch vereinbar) mit einem vorgegebenen Parameterwert? Die Antwort auf diese 2. Grundfrage heisst **statistischer Test**.

**3. Grundfragestellung:** Welche Parameterwerte sind mit den Beobachtungen kompatibel (statistisch vereinbar)? Die Antwort auf diese 3. Grundfrage heisst **Konfidenzintervall** oder **Vertrauensintervall**. Das Konfidenzintervall ist allgemeiner und informativer als ein statistischer Test.

*Beispiel (Forts.):* Im Beispiel von Spermasexing können die Grundfragen so lauten:

1. Welches ist der plausibelste Wert  $\pi$  (zu der Beobachtung  $x = 11$ )?
2. Ist die Beobachtung  $x = 11$  kompatibel mit  $\pi = 0.7$ ?
3. Welcher Bereich (Intervall) für den Parameter  $\pi$  ist mit der Beobachtung  $x = 11$  kompatibel?

### 3.2 Schätzung, statistischer Test und Vertrauensintervall bei Binomial-Verteilung (Stahel Kap. 7.2, 8.2, 9.1, 9.2)

Wir betrachten folgende Situation: gegeben ist eine Beobachtung  $x$ , welche als Realisierung von  $X \sim \text{Binomial}(n, \pi)$  aufgefasst wird. Wir möchten Schlüsse ziehen über den unbekannt Parameter  $\pi$ .

### 3.2.1 (Punkt-)Schätzung

Eine Schätzung für  $\pi$  kann pragmatisch hergeleitet werden. Da  $\mathcal{E}(X) = n\pi$  (siehe Kapitel 2.5) gilt:  $\pi = \mathcal{E}(X)/n$ . Der Wert  $n$  (Anzahl unabhängiger Versuche) ist als bekannt vorausgesetzt: die einzige Unbekannte ist dann  $\mathcal{E}(X)$ . Eine pragmatisch motivierte Schätzung ist dann:  $\widehat{\mathcal{E}(X)} = x (= \text{Beobachtung})$ , d.h. man setzt die Beobachtung gleich dem Erwartungswert. Somit ergibt sich:

$$\hat{\pi} = x/n.$$

### 3.2.2 Statistischer Test

*Beispiel:* Es wird 100-mal eine Münze geworfen.

Betrachte  $X = \text{Anzahl Kopf (K)}$  bei 100 Würfeln. Es ist vernünftig, das Modell zu benutzen:  $X \sim \text{Binomial}(100, \pi)$ . Beobachtet (realisiert) wurde  $x = 58$ . Wir möchten testen, ob die Münze fair ist, d.h. ob  $\pi = 1/2$ .

#### Motivation

Im obigen Beispiel stellen wir die folgende Überlegung an. Wir nehmen einmal an, dass die Münze fair ist, d.h. dass  $\pi = 1/2$ , und berechnen die Wahrscheinlichkeiten für “unplausible” Ereignisse von der Form  $\{X \geq c\}$  für “grosse” Werte  $c$ . Die Absicht dabei ist zu quantifizieren, ob die beobachtete Anzahl  $x = 58$  bereits zu einem “unplausiblen” Ereignis gehört (so dass man schliessen müsste, dass die Münze nicht mehr fair ist, das heisst  $\pi > 1/2$ ). Die folgende Tabelle liefert die Zahlen für  $X \sim \text{Binomial}(100, 1/2)$ .

	$c = 52$	$c = 53$	$c = 54$	$c = 55$	$c = 56$	$c = 57$	$c = 58$	$c = 59$	$c = 60$
$P(X \geq c)$	0.382	0.309	0.242	0.184	0.136	0.097	0.067	0.044	0.028

Typischerweise deklariert man ein Ereignis als “unplausibel” falls dessen Wahrscheinlichkeit weniger oder gleich 5% beträgt. In unserem Beispiel sehen wir, dass die Beobachtung  $x = 58$ , welche zu dem Ereignis  $X \geq 58$  gehört, eine zugehörige Wahrscheinlichkeit von 6.7% hat und deshalb immer noch als plausibel eingestuft wird. Das heisst, dass die Beobachtung  $x = 58$  durchaus noch als plausibel bei einer fairen Münze eingestuft werden kann. Hätte man aber 59-mal Kopf beobachtet, so würde man dies als nicht mehr sehr plausibel bei einer fairen Münze einstufen: die zugehörige Wahrscheinlichkeit ist mit 4.4% bereits eher klein. (Natürlich ist die Grenze, welche durch eine Wahrscheinlichkeit von 5% gegeben ist, willkürlich. Später werden wir dies mit dem sogenannten P-Wert charakterisieren, siehe unten.

#### Formales Vorgehen

Ein statistischer Test für den Parameter  $\pi$  im Modell  $X \sim \text{Binomial}(n, \pi)$  ist wie folgt aufgebaut.

1. Spezifiziere die sogenannte **Nullhypothese**  $H_0$ :

$$H_0: \pi = \pi_0,$$

und (anhand der Problemstellung) eine sogenannte **Alternative**  $H_A$ :

$$\begin{aligned} H_A : \quad & \pi \neq \pi_0 \text{ (zwei-seitig)} \\ & \pi > \pi_0 \text{ (ein-seitig nach oben)} \\ & \pi < \pi_0 \text{ (ein-seitig nach unten)}. \end{aligned}$$

*Beispiel (Forts.):* Es wird 100-mal eine Münze geworfen.

Betrachte  $X =$  Anzahl Kopf (K) bei 100 Würfeln mit dem Modell  $X \sim \text{Binomial}(100, \pi)$ . Wir möchten folgende Frage untersuchen: Ist die Münze fair? Oder ist sie unfair, dass sie eher zu oft auf Kopf (K) landet? Im Formalismus des statistischen Tests heisst das:  $H_0 : \pi = \pi_0 = 1/2$  und  $H_A : \pi > \pi_0 = 1/2$ . Insbesondere die Wahl der Alternative muss bezüglich der interessierenden Fragestellung getroffen werden.

2. Lege das sogenannte **Signifikanzniveau**  $\alpha$  fest. Typischerweise wählt man  $\alpha = 0.05$  (5%) oder auch  $\alpha = 0.01$  (1%).

3. Bestimme den sogenannten **Verwerfungsbereich**  $K$ . Qualitativ zeigt  $K$  in Richtung der Alternative:

$$\begin{aligned} K &= [0, c_u] \cup [c_o, n] && \text{falls } H_A : \pi \neq \pi_0, \\ K &= [c, n] && \text{falls } H_A : \pi > \pi_0, \\ K &= [0, c] && \text{falls } H_A : \pi < \pi_0. \end{aligned}$$

Quantitativ wird  $K$  so berechnet, dass

$$P_{H_0}(X \in K) = \underbrace{P_{\pi_0}}_{\text{von Binomial}(n, \pi_0)}(X \in K) \stackrel{\approx}{\leq} \alpha. \quad (3.1)$$

*Beispiel (Forts.):* Bei 100-maligem Münzwurf.

$H_0 : \pi = 1/2$  (d.h.  $\pi_0 = 1/2$ ) und  $H_A : \pi > 1/2$ . Für  $\alpha = 0.05$  haben wir in der Tabelle oben gesehen, dass  $K = [59, 100]$  ist.

4. Erst jetzt betrachte, ob die Beobachtung  $x$  in den Verwerfungsbereich  $K$  fällt:

falls ja: so verwerfe  $H_0$  (die Alternative ist dann "signifikant")

falls nein: belasse  $H_0$  (was nicht heisst, dass deswegen  $H_0$  statistisch bewiesen ist).

Diese Art der Test-Entscheidung beruht auf dem Widerspruchs-Prinzip: ein statistischer Nachweis/Beweis gelingt bloss dann, wenn die Nullhypothese verworfen werden kann. Diese Art von wissenschaftlicher Evidenz-Findung wurde bereits von Aristoteles in der Antike propagiert.

*Beispiel (Forts.):* Bei 100-maligem Münzwurf.

Da  $x = 58$  effektiv beobachtet wurde:  $H_0$  wird belassen. Das heisst, dass es keine statistische Evidenz gibt (auf dem Signifikanzniveau  $\alpha = 0.05$ ), dass die Münze zu Gunsten von Kopf (K) gefälscht ist.

*Beispiel (Forts.):* Beim Spermasexing (vgl. Kapitel 2.4) wurden  $x = 11$  weibliche Kuhkälber gezüchtet von insgesamt 12 Kälbern. Der Hersteller der Methode behauptet, dass die Erfolgswahrscheinlichkeit grösser als 70% ist. Der Test ist wie folgt.

Modell:  $X \sim \text{Binomial}(12, \pi)$

$H_0 : \pi = \pi_0 = 0.7$

$H_A : \pi > \pi_0 = 0.7$

Signifikanzniveau: wir wählen  $\alpha = 0.05$

Verwerfungsbereich:  $P_{\pi=0.7}(X \in K) \stackrel{\approx}{\leq} 0.05 \rightsquigarrow K = \{12\}$

Entscheid:  $H_0$  wird belassen, d.h. die Aussage des Herstellers ist nicht signifikant

### Fehler 1. und 2. Art

Bei einem statistischen Test treten 2 Arten von Fehlern auf.

**Fehler 1. Art:** Fälschliches Verwerfen von  $H_0$ , obwohl  $H_0$  richtig ist.

**Fehler 2. Art:** Fälschliches Beibehalten von  $H_0$ , obschon die Alternative zutrifft.

Der Fehler 1. Art wird als "schlimmer" betrachtet: er wird direkt kontrolliert mittels der

Konstruktion eines Tests: die Formel (3.1) besagt:

$$P(\text{Fehler 1. Art}) = P_{H_0}(X \in K) \stackrel{\approx}{\leq} \alpha.$$

Das Signifikanzniveau kontrolliert also die Wahrscheinlichkeit für eine Fehler 1. Art. Es gilt aber auch:

$P(\text{Fehler 2. Art})$  wird grösser falls  $\alpha$  kleiner gewählt wird.

Die Wahl von  $\alpha$  steuert also einen Kompromiss zwischen Fehler 1. und 2. Art. Weil man aber primär einen Fehler 1. Art vermeiden will, wählt man  $\alpha$  klein, z.B.  $\alpha = 0.05$ .

*Beispiel (Forts.):* Beim Spermasexing nehmen wir einmal an, dass in Tat und Wahrheit der Parameter  $\pi = 0.8 \in H_A$  ist (die Spezifikationen des Tests sind wie oben:  $H_0 : \pi = 0.7$ ,  $H_A : \pi > 0.7$  und  $\alpha = 0.05$ ). Da der Verwerfungsbereich  $K = \{12\}$  ist (siehe oben), gilt dann:

$$P(\text{Test behält } H_0 \text{ bei, obschon } \pi = 0.8) = P_{\pi=0.8}(X \leq 11) = 1 - P_{\pi=0.8}(X = 12) = 0.93.$$

Das heisst, dass ein Fehler 2. Art (unter der Annahme dass  $\pi = 0.8$ ) mit grosser Wahrscheinlichkeit auftritt. Das ist natürlich enttäuschend, aber unvermeidlich bei der kleinen Anzahl von 12 Versuchen. Beachte, dass die Wahrscheinlichkeit für einen Fehler 1. Art  $\stackrel{\approx}{\leq} 0.05$  ist.

### Der P-Wert

Die Entscheidung eines Tests mit "Verwerfen" oder "Beibehalten" der Nullhypothese  $H_0$  ist abhängig von der etwas willkürlichen Wahl des Signifikanzniveaus  $\alpha$ . Mathematisch bedeutet dies, dass der Verwerfungsbereich  $K = K(\alpha)$  abhängig von der Wahl von  $\alpha$  ist.

Man kann sich einfach überlegen, dass qualitative Folgendes gilt:

Verwerfungsbereich  $K = K(\alpha)$  wird kleiner mit kleiner werdendem  $\alpha$ ,

weil bei kleinem  $\alpha$  die Wahrscheinlichkeit für einen Fehler 1. Art klein ist (und dies wird natürlich erreicht, wenn es tendenziell schwierig ist eine Nullhypothese  $H_0$  zu verwerfen - sprich Verwerfungsbereich  $K$  klein). Umgekehrt gilt natürlich auch, dass  $K = K(\alpha)$  grösser wird mit wachsendem  $\alpha$ . Dies impliziert: es gibt ein Signifikanzniveau, bei dem die Nullhypothese  $H_0$  "gerade noch" verworfen wird.

Der P-Wert ist definiert als das kleinste Signifikanzniveau, bei dem die Nullhypothese  $H_0$  (gerade noch) verworfen wird

Der P-Wert kann folgendermassen gerechnet werden: die Beobachtung  $X = x$  kommt auf die Grenze des Verwerfungsbereichs  $K = K(\text{P-Wert})$  mit Signifikanzniveau = P-Wert zu liegen; siehe auch Abbildung 3.1.

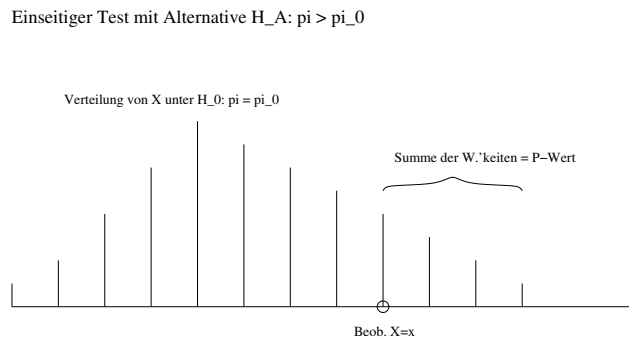


Abbildung 3.1: Schematische Darstellung des P-Werts bei einer einseitigen Alternative  $H_A: \pi > \pi_0$ .

Der P-Wert liefert mehr Information als bloss die Test-Entscheidung bei einem vorbestimmten Signifikanzniveau  $\alpha$  (z.B.  $\alpha = 0.05$ ). Insbesondere gilt aufgrund der Definition des P-Werts:

verwerfe  $H_0$  falls P-Wert  $\leq \alpha$   
 belasse  $H_0$  falls P-Wert  $> \alpha$ .

Zusätzlich zu dieser Entscheidungsregel quantifiziert der P-Wert wie signifikant eine Alternative ist (d.h. wie gross die Evidenz ist für das Verwerfen von  $H_0$ ). Sprachlich wird manchmal wie folgt übersetzt:

P-Wert  $\approx 0.05$  : schwach signifikant  
 P-Wert  $\approx 0.01$  : signifikant  
 P-Wert  $\approx 0.001$  : stark signifikant  
 P-Wert  $\leq 10^{-4}$  : äusserst signifikant

*Beispiel (Forts.):* Beim Spermasexing betrachten wir die Null-Hypothese  $\pi = 0.7$  und die Alternative  $\pi > 0.7$ . Beobachtet wurde  $x = 11$ , aufgefasst als Realisierung von  $X \sim \text{Binomial}(12, \pi)$ . Der P-Wert ist dann:

$$P_{\pi=0.7}(X \geq 11) = P_{\pi=0.7}(X = 11) + P_{\pi=0.7}(X = 12) = 0.085.$$

Wie wir bereits früher gesehen haben, liefert dies kein Verwerfen von  $H_0$  auf dem Signifikanzniveau  $\alpha = 0.05$ . (Wenn man - aus irgendwelchen Gründen - im voraus das Signifikanzniveau  $\alpha = 0.09$  gewählt hätte, so könnte man  $H_0$  auf diesem Signifikanzniveau  $\alpha = 0.09$  verwerfen).

### 3.2.3 Vertrauensintervall

Informativer als ein statistischer Test ist ein sogenanntes Vertrauensintervall (auch Konfidenzintervall genannt). Es liefert eine Antwort auf die 3. Grundfragestellung von Kapitel 3.1: Welche Werte von  $\pi$  sind mit der Beobachtung  $x$  kompatibel (statistisch vereinbar).

Ein Vertrauensintervall  $I$  zum Niveau  $1 - \alpha$  besteht aus allen Parameterwerten, die im Sinne des statistischen Tests zum Signifikanzniveau  $\alpha$  mit der Beobachtung verträglich sind (üblicherweise nimmt man den zweiseitigen Test). Mathematisch heisst dies:

$$I = \{\pi_0; \text{Nullhypothese } H_0 : \pi = \pi_0 \text{ wird belassen}\}. \quad (3.2)$$

Diese Beziehung stellt eine Dualität zwischen Tests und Vertrauensintervall dar.

Die Berechnung kann grafisch, oder mit einer Tabelle erfolgen. Falls  $n$  "gross" ist, so kann die sogenannte Normalapproximation (siehe Kap. 4.5) benützt werden. Letztere ergibt folgendes approximatives Konfidenzintervall  $I$  zum Niveau  $1 - \alpha = 0.95$  für den unbekannt Parameter  $\pi$ :

$$I \approx \frac{x}{n} \pm 1.96 \sqrt{\frac{x}{n} \left(1 - \frac{x}{n}\right) \frac{1}{n}} \quad (3.3)$$

Das Vertrauensintervall  $I = I(x)$  hängt von der Beobachtung ab. Wenn man anstelle der Beobachtung die zugehörige Zufallsvariable  $X$  einsetzt, so ist  $I(X)$  zufällig und hat die Eigenschaft:

$$P(\pi \in I(X)) \stackrel{\approx}{>} 1 - \alpha.$$

Dies kann so interpretiert werden, dass der wahre Parameter  $\pi$  mit Wahrscheinlichkeit  $1 - \alpha$  im Konfidenzintervall  $I$  enthalten ist.

*Beispiel (Forts.):* Beim Spermasexing erhält man für ein zweiseitiges Konfidenzintervall zum Niveau  $1 - \alpha = 0.95$  mittels einer Tabelle oder dem Computer für die Berechnung von (3.2):

$$I = (0.615, 0.998)$$

Das heisst, dass der wahre "Zucht"-Parameter  $\pi$  mit einer Wahrscheinlichkeit von 95% in  $I$  liegt. Es besteht also auf Grund der kleinen Stichprobe grosse Unsicherheit, wie erfolgreich die Methode bei langfristigem Einsatz tatsächlich sein wird. Die Näherungsformel in (3.3) ist für dieses Beispiel nicht besonders gut, weil  $n = 12$  eher klein ist. Man erhält mit (3.3):

$$I \approx (0.760, 1.073)$$

Der rechte Endpunkt ist natürlich zu gross, denn der Parameter  $\pi$  ist ja kleiner oder gleich 1.



### 3.3 Schätzung, Test und Vertrauensintervall bei Poisson-Verteilung (Stahel, Kap. 7.2, 8.1, 9.1)

Wir betrachten folgende Situation: gegeben ist eine Beobachtung  $x$ , welche als Realisierung von  $X \sim \text{Poisson}(\lambda)$  aufgefasst wird. Wir möchten Schlüsse ziehen über den unbekannt Parameter  $\lambda$ .

#### 3.3.1 (Punkt-)Schätzung

Da  $E(X) = \lambda$  (siehe Kapitel 2.6), und unter Verwendung der pragmatischen Schätzung von  $E(X)$  mittels der Beobachtung  $x$ , erhält man die Schätzung:

$$\hat{\lambda} = x.$$

#### 3.3.2 Statistischer Test

Ein statistischer Test für den Parameter  $\lambda$  im Modell  $X \sim \text{Poisson}(\lambda)$  erfolgt völlig analog zu der Konstruktion bei der Binomial-Verteilung in Kapitel 3.2.2.

1. Spezifiziere die **Nullhypothese**  $H_0$ :

$$H_0 : \lambda = \lambda_0,$$

und (anhand der Problemstellung) eine **Alternative**  $H_A$ :

$$\begin{aligned} H_A : \quad & \lambda \neq \lambda_0 \text{ (zwei-seitig)} \\ & \lambda > \lambda_0 \text{ (ein-seitig nach oben)} \\ & \lambda < \lambda_0 \text{ (ein-seitig nach unten)}. \end{aligned}$$

2. Lege das **Signifikanzniveau**  $\alpha$  fest, zum Beispiel  $\alpha = 0.05$ .

3. Bestimme den **Verwerfungsbereich**  $K$ . Qualitativ zeigt  $K$  in Richtung der Alternative:

$$\begin{aligned} K &= [0, c_u] \cup [c_o, \infty) && \text{falls } H_A : \lambda \neq \lambda_0, \\ K &= [c, \infty) && \text{falls } H_A : \lambda > \lambda_0, \\ K &= [0, c] && \text{falls } H_A : \lambda < \lambda_0. \end{aligned}$$

Quantitativ wird  $K$  so berechnet, dass

$$P_{H_0}(X \in K) = \underbrace{P_{\lambda_0}}_{\text{von Poisson}(\lambda_0)} (X \in K) \stackrel{\approx}{\leq} \alpha.$$

4. Erst jetzt betrachte, ob die Beobachtung  $x$  in den Verwerfungsbereich  $K$  fällt:  
falls ja: verwerfe  $H_0$ ;  
falls nein: belasse  $H_0$ .

Die Konzepte wie Fehler 1. und 2. Art oder P-Wert sind identisch wie in Kapitel 3.2.2.

### 3.3.3 Vertrauensintervall

Das Vertrauensintervall  $I$  zum Niveau  $1 - \alpha$  (siehe auch Kapitel 3.2.3) besteht aus allen Werten  $\lambda$ , die beim zugehörigen Test akzeptiert werden. Manchmal kann auch die folgende Approximation benutzt werden für ein zweiseitiges Konfidenzintervall zum Niveau  $1 - \alpha = 0.95$ :

$$I = I(x) \approx x \pm 1.96\sqrt{x}.$$

*Beispiel:* Im Jahr 1992 gab es  $x = 554$  Tote bei Verkehrsunfällen in der Schweiz. Wir fassen diese Beobachtung auf als eine Realisierung von  $X \sim \text{Poisson}(\lambda)$ . Die Schätzung ist dann  $\hat{\lambda} = 554$  und das approximative Vertrauensintervall ist  $I = I(x) \approx (507.9, 600.1)$ .

# Kapitel 4

## Modelle und Statistik für Messdaten

### 4.1 Einleitung

In vielen Anwendungen hat man es nicht mit Zähl-, sondern mit Messdaten zu tun, bei denen die Werte im Prinzip kontinuierlich sind. Zur Illustration betrachten wir zwei Datensätze. Beim ersten werden zwei Methoden zur Bestimmung der latenten Schmelzwärme von Eis verglichen. Wiederholte Messungen der freigesetzten Wärme beim Uebergang von Eis bei  $-0.72^\circ\text{C}$  zu Wasser bei  $0^\circ\text{C}$  ergaben die folgenden Werte (in cal/g):

Methode A	79.98	80.04	80.02	80.04	80.03	80.03	80.04	79.97	80.05	80.03
Methode A	80.02	80.00	80.02							
Methode B	80.02	79.94	79.98	79.97	79.97	80.03	79.95	79.97		

Obwohl die Messungen mit der grösstmöglichen Sorgfalt durchgeführt und alle Störeinflüsse ausgeschaltet wurden, variieren die Messungen von Fall zu Fall. Wir werden diese Variationen innerhalb der Messreihen als zufällig modellieren, das heisst wir interpretieren diese Werte als Realisierungen von Zufallsvariablen. Wir werden dann die Frage beantworten, ob die Unterschiede zwischen den Methoden ebenfalls als zufällig angesehen werden können, oder ob ein systematischer Unterschied plausibler ist, der auch in der ganzen Population, d.h. in weiteren Messungen, bestehen bleibt. Im letzteren Fall werden wir dann noch zusätzlich angeben, wie gross der systematische Unterschied etwa ist.

Im zweiten Beispiel wurde bei 11 Individuen die Aggregation von Blutplättchen vor und nach dem Rauchen einer Zigarette gemessen. Die folgenden Daten geben den Anteil aggregierter Blutplättchen (in Prozent) nach einer Stimulation an.

Individuum	1	2	3	4	5	6	7	8	9	10	11
Vorher	25	25	27	44	30	67	53	53	52	60	28
Nachher	27	29	37	56	46	82	57	80	61	59	43

Wieder variieren die Werte in einer nicht vorhersehbaren Art. Diesmal handelt es sich jedoch weniger um Messfehler, sondern um Variation zwischen Individuen (vermutlich gäbe es auch noch eine gewisse Variation beim gleichen Individuum, wenn der Test wiederholt würde). Die Aggregation bei diesen 11 Individuen ist meistens, aber nicht immer nach dem Rauchen höher, und die Fragestellung lautet, ob es sich hier um einen zufälligen Effekt handelt, der auf die spezifische Stichprobe beschränkt ist, oder ob man dieses Resultat

auf eine grössere Population verallgemeinern kann. Im letzteren Fall möchte man wieder angeben, wie gross die mittlere Zunahme etwa ist.

## 4.2 Deskriptive Statistik (Stahel, Kap. 2 und 3.1, 3.2)

Bei einer statistischen Analyse ist es wichtig, nicht einfach blind ein Modell anzupassen oder ein statistisches Verfahren anzuwenden. Die Daten sollten immer mit Hilfe von geeigneten grafischen Mitteln dargestellt werden, da man nur auf diese Weise unerwartete Strukturen und Besonderheiten entdecken kann. Kennzahlen können einen Datensatz grob charakterisieren. Im Folgenden werden die Daten mit  $x_1, \dots, x_n$  bezeichnet.

### 4.2.1 Kennzahlen

Häufig will man die Verteilung der Daten numerisch zusammenfassen. Dazu braucht man mindestens zwei Kenngrössen, eine für die Lage und eine für die Streuung. Die bekanntesten solchen Grössen sind das *arithmetische Mittel* für die Lage,

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

und die *empirische Standardabweichung* für die Streuung,

$$s_x = \sqrt{\text{var}} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}.$$

(der Nenner  $n - 1$ , anstelle von  $n$ , ist mathematisch begründet und hat die Eigenschaft, dass kein “systematischer” Fehler auftritt).

Alternative Kenngrössen sind der *Median* als Lagemass und die *Quartilsdifferenz* als Streuungsmass. Diese werden mit Hilfe von sogenannten Quantilen definiert.

### Quantil

Das *empirische  $\alpha$ -Quantil* ist anschaulich gesprochen der Wert, bei dem  $\alpha \times 100\%$  der Datenpunkte kleiner und  $(1 - \alpha) \times 100\%$  der Punkte grösser sind.

Zur formalen Definition führen wir die geordneten Werte ein:

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}.$$

Das empirische  $\alpha$ -Quantil ist dann gleich

$$\begin{aligned} & \frac{1}{2}(x_{(\alpha n)} + x_{(\alpha n + 1)}) \quad \text{falls } \alpha \cdot n \text{ eine ganze Zahl ist,} \\ & x_{(k)} \quad \text{wobei } k = \text{nächstgrössere ganze Zahl von } \alpha \cdot n; \quad \text{falls } \alpha \cdot n \text{ keine ganze Zahl ist.} \end{aligned}$$

Der (empirische) Median ist das empirische 50%-Quantil: d.h., es markiert die “mittlere” Beobachtung und ist also ein Mass für die Lage der Daten.

Die Quartilsdifferenz ist gleich

$$\text{empirisches 75\%-Quantil} - \text{empirisches 25\%-Quantil}$$

und ist ein Streuungsmass für die Daten.

Median und Quartilsdifferenz haben den Vorteil, dass sie robust sind: das heisst, dass sie viel weniger stark durch extreme Beobachtungen beeinflusst werden können als arithmetisches Mittel und Standardabweichung.

*Beispiel:* Messung der Schmelzwärme von Eis mit Methode A

Aufgrund von  $n = 13$  Messungen: das arithmetische Mittel  $\bar{x} = 80.02$  und die Standardabweichung  $s_x = 0.024$ . Ferner, für  $n = 13$  ist  $0.25n = 3.25$ ,  $0.5n = 6.5$  und  $0.75n = 9.75$ . Damit ist das 25%-Quantil gleich  $x_{(4)} = 80.02$ , der Median gleich  $x_{(7)} = 80.03$  und das 75%-Quantil gleich  $x_{(10)} = 80.04$ .

## Standardisierung

Durch Verschiebung und Skalierung der Werte kann man erreichen, dass zwei oder mehrere Datensätze die gleiche Lage und Streuung haben. Insbesondere kann man einen Datensatz so standardisieren, dass das arithmetische Mittel gleich Null und die Standardabweichung gleich 1 ist. Dies erreicht man mittels der linear transformierten Variablen

$$z_i = \frac{x_i - \bar{x}}{s_x} \quad (i = 1, \dots, n).$$

Alle Aspekte einer Verteilung, die bei einer Verschiebung oder Skalierung unverändert bleiben, machen die Form der Verteilung aus. Dazu gehört insbesondere die Schiefe (Asymmetrie) der Verteilung, für die es auch Kennzahlen gibt.

### 4.2.2 Grafische Methoden

Einen Überblick über die auftretenden Werte ergibt das *Histogramm*. Um ein Histogramm zu zeichnen, bildet man Klassen  $(c_{k-1}, c_k]$  und berechnet die Häufigkeiten  $h_k$ , d.h. die Anzahl Werte in diesem Intervall. Dann trägt man über den Klassen Balken auf, deren Fläche *proportional* zu  $h_k$  ist.

Beim Boxplot hat man ein Rechteck, das vom empirischen 25%- und vom 75%-Quantil begrenzt ist, und Linien, die von diesem Rechteck bis zum kleinsten- bzw. grössten "normalen" Wert gehen (per Definition ist ein normaler Wert höchstens 1.5 mal die Quartilsdifferenz von einem der beiden Quartile entfernt). Zusätzlich gibt man noch Ausreisser durch Sterne und den Median durch einen Strich an. Der Boxplot ist vor allem dann geeignet, wenn man die Verteilungen einer Variablen in verschiedenen Gruppen (die im allgemeinen verschiedenen Versuchsbedingungen entsprechen) vergleichen will; siehe Abbildung 4.1.

Die *empirische kumulative Verteilungsfunktion*  $F_n(\cdot)$  ist eine Treppenfunktion, die links von  $x_{(1)}$  gleich null ist und bei jedem  $x_{(i)}$  einen Sprung der Höhe  $\frac{1}{n}$  hat (falls ein Wert mehrmals vorkommt, ist der Sprung ein Vielfaches von  $\frac{1}{n}$ ). In andern Worten:

$$F_n(x) = \frac{1}{n} \text{Anzahl}\{i \mid x_i \leq x\}.$$

Abbildung 4.2 zeigt die empirische kumulative Verteilungsfunktion für die Messungen der Schmelzwärme von Eis mit Methode A.

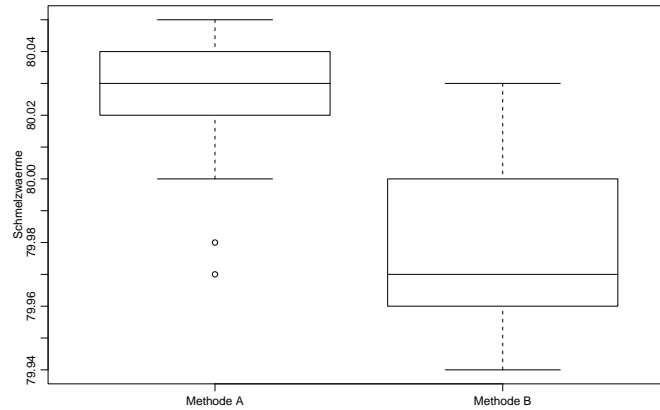


Abbildung 4.1: Boxplots für die zwei Methoden zur Bestimmung der Schmelzwärme von Eis.

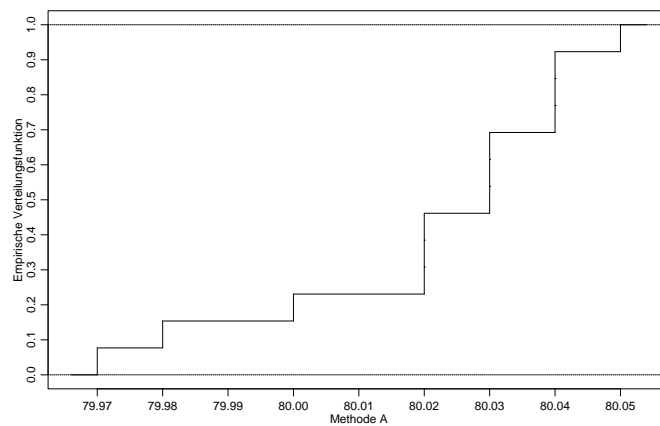


Abbildung 4.2: Empirische kumulative Verteilungsfunktion der Messungen der Schmelzwärme von Eis mit Methode A.

## Mehrere Variablen

Wenn wir zwei verschiedene Größen messen, d.h. wenn die Daten von der Form  $(x_1, y_1), \dots, (x_n, y_n)$  sind, interessiert man sich in erster Linie für die Zusammenhänge und Abhängigkeiten zwischen den Variablen. Diese kann man aus dem *Streudiagramm* ersehen, welches die Daten als Punkte in der Ebene darstellt: Die  $i$ -te Beobachtung entspricht dem Punkt mit Koordinaten  $(x_i, y_i)$ . Die Abbildung 4.3 zeigt das Streudiagramm für die Werte “vorher” und “nachher” bei der Blutplättchen-Aggregation. Man sieht einen klaren monotonen Zusammenhang, Individuen haben also eine Tendenz zu starker, bzw. schwacher Aggregation, unabhängig vom Rauchen.

Für die numerische Zusammenfassung der Abhängigkeit ist die **empirische Korrelation**  $r$  (oder auch mit  $\hat{\rho}$  bezeichnet) am gebräuchlichsten:

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

Die empirische Korrelation ist eine dimensionslose Zahl zwischen -1 und +1. Das Vorzeichen von  $r$  misst die Richtung und der Betrag die Stärke des linearen Zusammenhangs

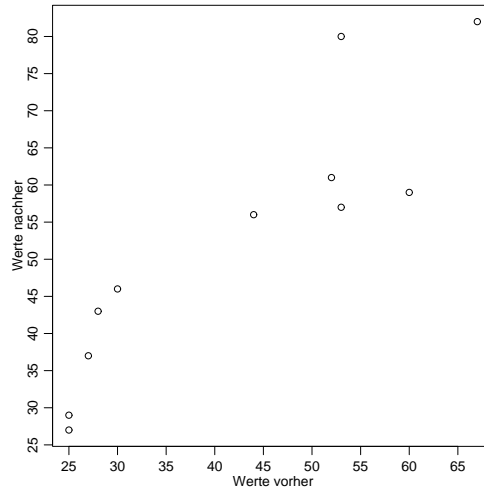


Abbildung 4.3: Streudiagramm der Blutplättchen-Aggregation vor und nach dem Rauchen einer Zigarette.

zwischen den beiden Variablen. Im Fall der Aggregation von Blutplättchen ist die empirische Korrelation gleich 0.9, was den Eindruck vom Streudiagramm bestätigt. Man sollte jedoch nie  $r$  berechnen, ohne einen Blick auf das Streudiagramm zu werfen, da ganz verschiedene Strukturen den gleichen Wert von  $r$  ergeben können.

Weitere Ausführungen sind in Kapitel 5.1 zu finden.

### 4.3 Stetige Zufallsvariablen und Wahrscheinlichkeitsverteilungen (Stahel, Kap. 6.1 – 6.4, 11.2)

Eine Zufallsvariable  $X$  heisst stetig, wenn deren Wertebereich  $W_X$  kontinuierlich ist; z.B.  $W_x = \mathbb{R}$ ,  $\mathbb{R}^+$  oder  $[0, 1]$ .

In Kapitel 2.3 hatten wir gesehen, dass die Wahrscheinlichkeitsverteilung einer diskreten Zufallsvariablen beschrieben werden kann, indem man die “Punkt”-Wahrscheinlichkeiten  $P(X = x)$  für alle möglichen  $x$  im Wertebereich angibt. Für eine stetige Zufallsvariable  $X$  gilt jedoch:

$$P(X = x) = 0 \text{ für alle } x \in W_X.$$

Dies impliziert, dass wir die Wahrscheinlichkeitsverteilung von  $X$  nicht mittels der Angaben von “Punkt”-Wahrscheinlichkeiten beschreiben können.

Die Wahrscheinlichkeitsverteilung einer stetigen Zufallsvariablen  $X$  kann jedoch beschrieben werden, indem man die Wahrscheinlichkeiten für alle Intervalle  $(a, b]$  ( $a < b$ ) angibt:

$$P(X \in (a, b]) = P(a < X \leq b).$$

Diese Information ist auch in der kumulativen Verteilungsfunktion  $F(x) = P(X \leq x)$  enthalten, denn:

$$P(a < X \leq b) = F(b) - F(a).$$

Zusammenfassend heisst dies, dass die Wahrscheinlichkeitsverteilung einer stetigen Zufallsvariablen  $X$  durch die kumulative Verteilungsfunktion beschrieben werden kann.

### 4.3.1 (Wahrscheinlichkeits-)Dichte

Die Idee einer "Punkt"-Wahrscheinlichkeit  $P(X = x)$  kann im Infinitesimalen auch für stetige Zufallsvariablen formuliert werden.

Die (Wahrscheinlichkeits-)Dichte  $f(\cdot)$  ist definiert als Ableitung der kumulativen Verteilungsfunktion:

$$f(x) = F'(x).$$

Damit erhalten wir folgende Interpretation:

$$P(x < X \leq x + h) \approx hf(x) \quad \text{falls } h \text{ klein ist.}$$

Die Begründung dafür ist:

$$P(x < X \leq x + h)/h = (F(x + h) - F(x))/h \approx f(x)$$

wobei die letzte Approximation aus der Definition einer Ableitung folgt.

Weil  $F(x) = \int_{-\infty}^x f(y)dy$  die Stammfunktion der Dichte ist, gelten die folgenden Eigenschaften:

1.  $f(x) \geq 0$  für alle  $x$  (da  $F(\cdot)$  monoton wachsend ist)
2.  $P(a < X \leq b) = \int_a^b f(x)dx$
3.  $\int_{-\infty}^{\infty} f(x)dx = 1$  (wegen 2.)

### Kennzahlen von stetigen Verteilungen

Der Erwartungswert  $\mathcal{E}(X)$  und die Standardabweichung  $\sigma_X$  einer stetigen Zufallsvariablen  $X$  haben dieselbe Bedeutung wie im diskreten Fall in Kapitel 2.5; einzig die Berechnungen sehen etwas anders aus. Es gelten:

$$\begin{aligned}\mathcal{E}(X) &= \int_{-\infty}^{\infty} xf(x)dx, \\ \text{Var}(X) &= \int_{-\infty}^{\infty} (x - \mathcal{E}(X))^2 f(x)dx, \quad \sigma_X = \sqrt{\text{Var}(X)}.\end{aligned}$$

In der frequentistischen Interpretation ist der Erwartungswert eine Idealisierung des arithmetischen Mittels der Werte einer Zufallsvariablen (bei vielen Wiederholungen).

Für den Erwartungswert der transformierten Zufallsvariablen  $Y = g(X)$ , wobei  $g : \mathbb{R} \rightarrow \mathbb{R}$  eine Transformation ist, gilt:

$$\mathcal{E}(Y) = \mathcal{E}(g(X)) = \int_{-\infty}^{\infty} g(x)f(x)dx.$$

Damit erhält man, dass

$$\text{Var}(X) = \mathcal{E}((X - \mathcal{E}(X))^2).$$



Die obige Formel gilt auch für diskrete Zufallsvariablen. Die folgenden **Rechenregeln** (auch für diskrete Zufallsvariablen) erweisen sich oft als nützlich: für beliebige  $a, b \in \mathbb{R}$ ,

$$\begin{aligned}\mathcal{E}[a + bX] &= a + b\mathcal{E}[X], \\ \text{Var}(X) &= \mathcal{E}[X^2] - (\mathcal{E}[X])^2, \\ \text{Var}(a + bX) &= b^2\text{Var}(X).\end{aligned}$$

Die **Quantile** (einer Verteilung von  $X$ )  $q(\alpha)$  ( $0 < \alpha < 1$ ) sind wie folgt definiert:

$$P(X \leq q(\alpha)) = \alpha.$$

Das heisst:

$$F(q(\alpha)) = \alpha \Leftrightarrow q(\alpha) = F^{-1}(\alpha).$$

Dies kann auch so interpretiert werden, dass  $q(\alpha)$  der Punkt ist, so dass die Fläche von  $-\infty$  bis  $q(\alpha)$  unter der Dichte  $f(\cdot)$  gleich  $\alpha$  ist. Siehe auch Abbildung 4.4. Das 50%-Quantil heisst der **Median**.

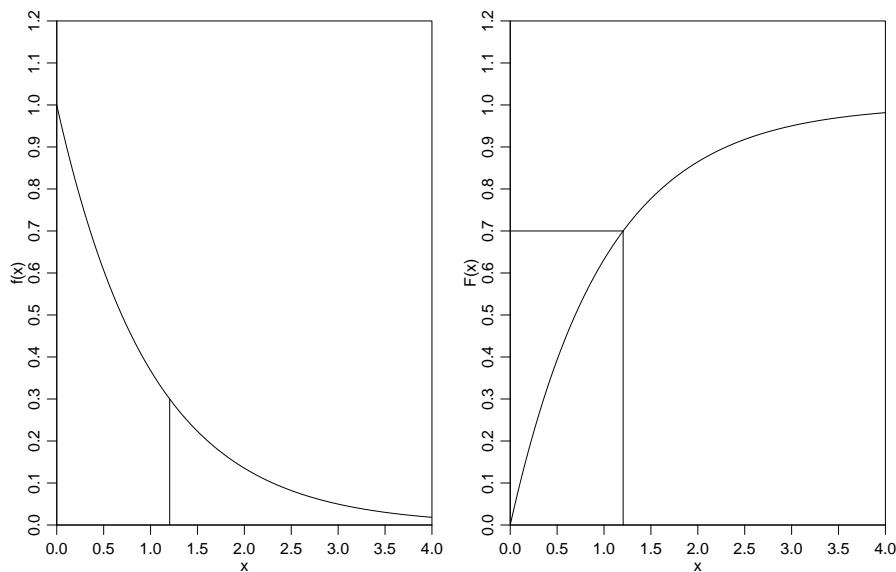


Abbildung 4.4: Illustration des 70%-Quantils  $q(0.7)$ . Links ist die Dichte gezeichnet, mit der Fläche von  $-\infty$  (oder hier auch 0) bis  $q(0.7)$  gleich 0.7. Rechts die kumulative Verteilungsfunktion und  $q(0.7)$  als Wert der Umkehrfunktion an der Stelle 0.7.

#### 4.4 Wichtige stetige Verteilungen (Stahel, Kap. 6.2, 6.4, 6.5, 11.2)

Wir haben in Kapitel 4.3 gesehen, dass wir die Wahrscheinlichkeitsverteilung einer stetigen Zufallsvariablen mit der kumulativen Verteilungsfunktion  $F(\cdot)$  oder der Dichte  $f(\cdot)$  charakterisieren können.

#### 4.4.1 Uniforme Verteilung

Die Uniforme Verteilung tritt auf bei Rundungsfehlern und als Formalisierung der völligen “Ignoranz”.

Eine Zufallsvariable  $X$  mit Wertebereich  $W_X = [a, b]$  heisst Uniform( $[a, b]$ ) verteilt, falls

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

Die Dichte ist also konstant auf dem Intervall  $[a, b]$ . Das heisst, dass die gleiche Wahrscheinlichkeit vorliegt auf dem ganzen Wertebereich  $W_X = [a, b]$ , deshalb der Name “uniform”.

Die zugehörige kumulative Verteilungsfunktion ist

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

Für  $X \sim \text{Uniform}([a, b])$  sind die Kennzahlen wie folgt:

$$\begin{aligned} \mathcal{E}(X) &= (a+b)/2, \\ \text{Var}(X) &= (b-a)^2/12, \quad \sigma_X = \sqrt{\text{Var}(X)}. \end{aligned}$$

#### 4.4.2 Exponential-Verteilung

Die Exponential-Verteilung ist das einfachste Modell für Wartezeiten auf Ausfälle.

*Beispiel:* Ionenkanäle

In Membranen von Muskel- und Nerven-Zellen gibt es viele Kanäle wo Ionen fließen können, falls der Kanal offen ist. Simple kinetische Modelle motivieren, dass die Offenzeit eines Kanals mit der Exponential-Verteilung modelliert werden kann.

Eine Zufallsvariable  $X$  mit Wertebereich  $W_X = \mathbb{R}^+ = [0, \infty)$  heisst Exponential-verteilt mit Parameter  $\lambda \in \mathbb{R}^+$  ( $\text{Exp}(\lambda)$ ) falls

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

Die zugehörige kumulative Verteilungsfunktion ist

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

Die Dichte und kumulative Verteilungsfunktion für  $\lambda = 1$  sind in Abbildung 4.4 zu sehen.

Für  $X \sim \text{Exp}(\lambda)$  sind die Kennzahlen wie folgt:

$$\begin{aligned} \mathcal{E}(X) &= 1/\lambda, \\ \text{Var}(X) &= 1/\lambda^2, \quad \sigma_X = \sqrt{\text{Var}(X)}. \end{aligned}$$

Einen Zusammenhang zwischen der Exponential- und Poisson-Verteilung ergibt sich wie folgt. Wenn die Zeiten zwischen den Ausfällen eines Systems Exponential( $\lambda$ )-verteilt sind, dann ist die Anzahl Ausfälle in einem Intervall der Länge  $t$  Poisson( $\lambda t$ )-verteilt.

### 4.4.3 Normal-Verteilung (Gauss-Verteilung)

Die Normal-Verteilung (manchmal auch Gauss-Verteilung genannt) ist die häufigste Verteilung für Messwerte.

*Beispiel:* Die Lichtmessungen eines “White Dwarf Sterns” können modelliert werden als Realisierungen von Normal-verteilten Zufallsvariablen.

Eine Zufallsvariable  $X$  mit Wertebereich  $W_X = \mathbb{R}$  heisst Normal-verteilt mit Parametern  $\mu \in \mathbb{R}$  und  $\sigma^2 \in \mathbb{R}^+$  ( $\mathcal{N}(\mu, \sigma^2)$ ) falls

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

Die zugehörige kumulative Verteilungsfunktion  $F(\cdot)$  ist nicht explizit darstellbar: allgemein gilt aber  $F(x) = \int_{-\infty}^x f(y)dy$ .

Für  $X \sim \mathcal{N}(\mu, \sigma^2)$  sind die Kennzahlen wie folgt:

$$\begin{aligned}\mathcal{E}(X) &= \mu, \\ \text{Var}(X) &= \sigma^2, \quad \sigma_X = \sqrt{\text{Var}(X)}.\end{aligned}$$

Das heisst, dass die Parameter  $\mu$  und  $\sigma^2$  eine natürliche Interpretation als Erwartungswert und Varianz der Verteilung haben. Drei Normalverteilungen mit verschiedenen Werten von  $\mu$  und  $\sigma$  sind in Abbildung 4.5 dargestellt.

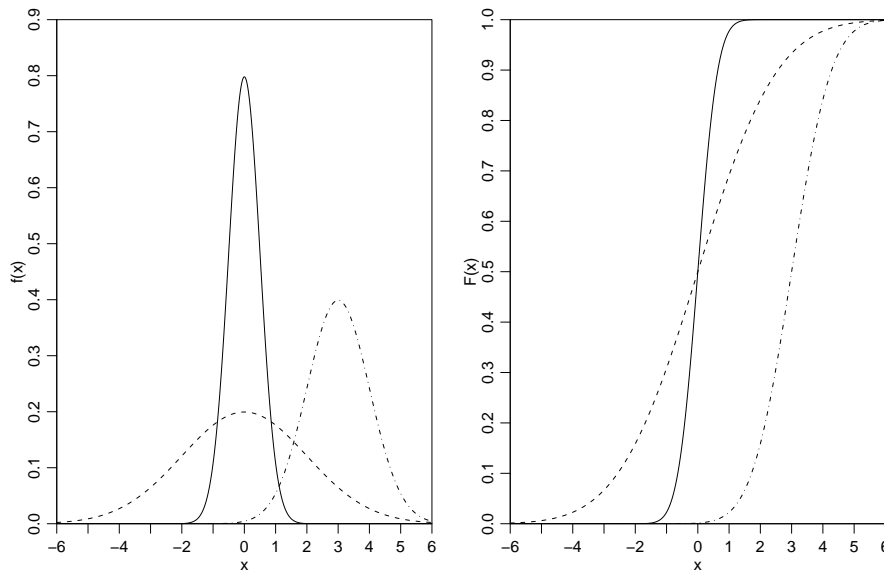


Abbildung 4.5: Dichten (links) und kumulative Verteilungsfunktionen (rechts) der Normalverteilungen mit  $\mu = 0, \sigma = 0.5$  (ausgezogen),  $\mu = 0, \sigma = 2$  (gestrichelt) und  $\mu = 3, \sigma = 1$  (Strich-Punkte).

## Die Standard-Normalverteilung

Die Normal-Verteilung mit  $\mu = 0$  und  $\sigma^2 = 1$  heisst Standard-Normalverteilung. Deren Dichte und kumulative Verteilungsfunktion werden wie folgt bezeichnet:

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right),$$
$$\Phi(x) = \int_{-\infty}^x \varphi(y) dy.$$

Die Werte von  $\Phi(\cdot)$  sind tabelliert. Wir werden unten sehen, dass eine Normal-Verteilung  $\text{cal}N(\mu, \sigma^2)$  immer in eine Standard-Normalverteilung transformiert werden kann. Deshalb werden die Werte von  $\Phi(\cdot)$  genügen um Wahrscheinlichkeiten und Quantile einer allgemeinen  $\mathcal{N}(\mu, \sigma^2)$ -Verteilung zu berechnen.

### 4.4.4 Transformationen

Es kann nützlich sein, eine stetige Zufallsvariable  $X$  zu transformieren:

$$Y = g(X),$$

wobei  $g : \mathbb{R} \rightarrow \mathbb{R}$  eine Transformation ist.

#### Lineare Transformationen

Wir betrachten hier den Fall einer linearen Transformation

$$g(x) = a + bx \quad (a, b \in \mathbb{R}).$$

Für  $Y = g(X)$  gilt dann (siehe auch Kapitel 4.3.1):

$$\begin{aligned} \mathcal{E}(Y) &= \mathcal{E}(a + bX) = a + b\mathcal{E}(X), \\ \text{Var}(Y) &= \text{Var}(a + bX) = b^2 \text{Var}(X), \quad \sigma_Y = b\sigma_X. \end{aligned} \tag{4.1}$$

Überdies gelten für  $b > 0$ :

$$\begin{aligned} \alpha - \text{Quantil von } Y &= q_Y(\alpha) = a + bq_X(\alpha), \\ f_Y(y) &= \frac{1}{b} f_X\left(\frac{y-a}{b}\right). \end{aligned} \tag{4.2}$$

#### Standardisieren einer Zufallsvariablen

Wir betrachten eine stetige Zufallsvariable  $X$ . Wir können  $X$  immer linear transformieren, so dass die transformierte Zufallsvariable Erwartungswert = 0 und Varianz = 1 hat. Dies geschieht wie folgt: betrachte die lineare Transformation

$$g(x) = -\frac{\mathcal{E}(X)}{\sigma_X} + \frac{1}{\sigma_X}x$$

und bilde die Transformierte

$$Z = g(X) = \frac{X - \mathcal{E}(X)}{\sigma_X}.$$

Mit Hilfe der Regeln in (4.1) gilt dann:  $\mathcal{E}(Z) = 0$ ,  $\text{Var}(Z) = 1$ .

Falls  $X \sim \mathcal{N}(\mu, \sigma)$ , so ist die standardisierte Zufallsvariable

$$Z = \frac{X - \mu}{\sigma} \sim \mathcal{N}(0, 1).$$

Dies folgt aus der Regel in (4.2) für linear transformierte Dichten. Diese lineare Transformation überführt also eine Normal-Verteilung in eine Standard-Normalverteilung. (Allgemeiner gilt: eine linear transformierte Normalverteilung ist wiederum eine Normalverteilung. Diese Eigenschaft, dass man mit linearen Transformationen innerhalb der Verteilungsklasse bleibt ist eine spezielle Eigenschaft der Normalverteilung und im Allgemeinen nicht richtig).

*Beispiel:* Berechnung von Wahrscheinlichkeiten bei  $\mathcal{N}(\mu, \sigma^2)$ .

Wir betrachten  $X \sim \mathcal{N}(2, 4)$  und möchten  $P(X \leq 5)$  berechnen. Man geht dann wie folgt vor.

$$P(X \leq 5) = P\left(\frac{X - 2}{\sqrt{4}} \leq \frac{5 - 2}{\sqrt{4}}\right) = P(Z \leq 3/2),$$

wobei  $Z \sim \mathcal{N}(0, 1)$ . Somit hat man:

$$P(X \leq 5) = P(Z < 3/2) = 0.933,$$

wobei man den numerischen Wert mittels einer Tabelle oder Computer bestimmt.

Mit der Regel (4.2) findet man für das 95%-Quantil ((4.2) angewendet mit  $a = -\mu/\sigma = 2/2 = 1$  und  $b = 1/\sigma = 1/2$ ):

$$q_X(0.95) = (q_Z(0.95) + 1) \cdot 2 = (\Phi^{-1}(0.95) + 1) \cdot 2 = 5.290.$$

## Nichtlineare Transformationen

Wenn  $g : \mathbb{R} \rightarrow \mathbb{R}$  eine beliebige Transformation ist, so sind viele der oben diskutierten Eigenschaften komplizierter. Die folgende Formel kann aber nützlich sein. Für eine stetige Zufallsvariable  $X$  mit Dichte  $f_X(\cdot)$  gilt:

$$\mathcal{E}(g(X)) = \int_{-\infty}^{\infty} g(x) f_X(x) dx.$$

Siehe auch Kapitel 4.3.1.

Eine oft gebrauchte Verteilung ist die **Lognormal-Verteilung**. Wenn  $X \sim \mathcal{N}(\mu, \sigma^2)$ , dann heisst  $Y = \exp(X)$  Lognormal-verteilt mit Parametern  $\mu \in \mathbb{R}$  und  $\sigma^2 \in \mathbb{R}^+$ . Die Lognormal-Verteilung ist nicht mehr symmetrisch und es gilt:  $\mathcal{E}(Y) = \exp(\mu + \sigma^2/2)$ .

### 4.4.5 Analogien zwischen Modellen und Daten

Zufallsvariablen und Verteilungen beschreiben die Population, d.h. was mit welcher Wahrscheinlichkeit passieren könnte. Daten  $x_1, \dots, x_n$  interpretieren wir als Realisierungen von Zufallsvariablen  $X_1, \dots, X_n$  (man könnte auch die  $n$  Daten als  $n$  Realisierungen von einer Zufallsvariablen  $X$  interpretieren; die Schreibweise mit mehreren Zufallsvariablen hat jedoch Vorteile, siehe Abschnitt 4.5).

Aus Daten können wir Rückschlüsse auf die zugrunde liegende Verteilung ziehen. Insbesondere haben alle Grössen, die für Zufallsvariablen definiert sind, ein Gegenstück für Datensätze gemäss folgender Tabelle. Die empirischen Grössen sind Schätzungen für die theoretischen Grössen. Diese werden mit wachsendem Stichprobenumfang  $n$  immer genauer.

Daten	Population (Modell)
Histogramm	Dichte
empirische kumulative Verteilungsfkt.	theoretische kumulative Verteilungsfkt.
empirische Quantile	theoretische Quantile
Arithmetisches Mittel	Erwartungswert
empirische Standardabweichung	theoretische Standardabweichung.

### 4.4.6 Überprüfen der Normalverteilungs-Annahme

Oft wollen wir überprüfen ob eine Verteilung ein brauchbares Modell für einen Datensatz darstellt. Das heisst, wir wollen überprüfen, ob ein Datensatz  $x_1, \dots, x_n$  als Realisierungen von einer Zufallsvariablen  $X$  mit einer Modell-Verteilung (z.B. mit einer kumulativen Verteilungsfunktion  $F(\cdot)$ ) aufgefasst werden kann.

Im Prinzip kann man das Histogramm der empirischen Daten mit der Dichte der Modell-Verteilung vergleichen. Oft sieht man aber die Abweichungen oder Übereinstimmungen besser, wenn man stattdessen die Quantile benutzt.

#### Q-Q Plot

Die Idee des Q-Q-Plot (Quantil-Quantil Plot) ist die empirischen Quantile gegen die theoretischen Quantile der Modell-Verteilung zu plotten. Konkret: plote für  $\alpha = 0.5/n, 1.5/n, \dots, (n - 0.5)/n$  die theoretischen Quantile der Modell-Verteilung  $q(\alpha)$  auf der x-Achse gegen die empirischen Quantile, welche den geordneten Beobachtungen  $x_{[1]} < x_{[2]} < \dots < x_{[n]}$  entsprechen, auf der y-Achse. Wenn die Beobachtungen gemäss der Modell-Verteilung erzeugt wurden, sollten diese Punkte ungefähr auf der Winkelhalbierenden  $y = x$  liegen.

#### Normal-Plot

Meist will man nicht eine spezifische Verteilung, sondern eine ganze Klasse von Verteilungen prüfen, also zum Beispiel die Klasse der Normalverteilungen mit beliebigem  $\mu$  und  $\sigma$ .

Ein Q-Q Plot wo die Modell-Verteilung gleich der Standard-Normalverteilung  $\mathcal{N}(0, 1)$  ist heisst Normal-Plot.

Falls die Daten Realisierungen von  $X \sim \mathcal{N}(\mu, \sigma^2)$  sind, so gilt für die Quantile von  $X$ :

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

siehe (4.2). Wenn man also einen Normal-Plot macht, so sollten die Punkte im Normal-Plot ungefähr auf der Geraden  $\mu + \sigma \cdot x$  liegen. Abbildung 4.6 zeigt zwei Normal-Plots: einmal

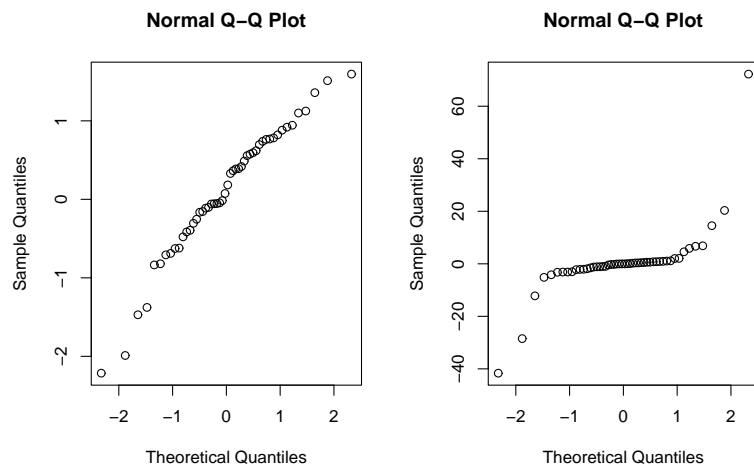


Abbildung 4.6: Links: Normal-Plot für 50 Realisierungen von  $\mathcal{N}(0, 1)$ . Rechts: Normal-Plot für 50 Realisierungen von Cauchy-Verteilung (sehr langschwänzig).

wo die Daten-generierende Verteilung eine Normal-Verteilung ist, und eine Situation wo die Daten von einer sehr langschwänzigen Verteilung erzeugt sind. Weitere Illustrationen sind in Abbildung 11.2. in Stahel ersichtlich.

## 4.5 Funktionen von Zufallsvariablen, Fehlerfortpflanzung (Stahel, Kap. 6.8 – 6.11)

In den meisten Anwendungen hat man es nicht mit einer, sondern mit mehreren Zufallsvariablen zu tun. Üblicherweise misst man die gleiche Größe mehrmals (man hat mehrere Individuen, oder man wiederholt die Messungen).

Die Messungen  $x_1, x_2, \dots, x_n$  fassen wir als Realisierungen der Zufallsvariablen  $X_1, \dots, X_n$  auf. Diese Notation ist oft bequemer als die Interpretation, dass die Messungen als  $n$  unabhängige Realisierungen einer Zufallsvariablen  $X$  sind. Oft sind wir an Funktionen der Zufallsvariablen  $X_1, \dots, X_n$  interessiert:

$$Y = g(X_1, \dots, X_n),$$

wobei  $g : \mathbb{R}^n \rightarrow \mathbb{R}$ , und  $Y$  wiederum eine Zufallsvariable ist. Wir betrachten hier vor allem die Funktion

$$\bar{X}_n = n^{-1} \sum_{i=1}^n X_i.$$

Beachte die folgende Verbindung: wenn  $x_i$  Realisierungen der Zufallsvariablen  $X_i$  sind, so ist das arithmetische Mittel der Daten  $\bar{x}_n = n^{-1} \sum_{i=1}^n x_i$  eine Realisierung der Zufallsvariablen  $\bar{X}_n$ .

Wir sind hier an der Verteilung der Zufallsvariablen  $\bar{X}_n$  interessiert (die Kenntnis dieser Verteilung ist ein wichtiges Werkzeug um später Statistik aufgrund von arithmetischen Mitteln von Daten zu machen). Dazu machen wir üblicherweise die folgende Annahme.

### Die i.i.d. Annahme

Oft treffen wir die Annahme, dass die Zufallsvariablen  $X_1, \dots, X_n$  **unabhängig** voneinander sind und dass alle **dieselbe** Verteilung haben. Dafür benutzt man die Notation:

$$X_1, \dots, X_n \text{ i.i.d.}$$

Die Abkürzung i.i.d. steht für: independent, identically distributed.

*Beispiel:*

$$X_1, \dots, X_n \text{ i.i.d. } \sim \mathcal{N}(\mu, \sigma^2)$$

bedeutet, dass alle  $X_i$ 's voneinander unabhängig sind und alle dieselbe Normal-Verteilung  $\mathcal{N}(\mu, \sigma^2)$  haben.

### Kennzahlen und Verteilung von $\bar{X}_n$

Wir nehmen in diesem Abschnitt an, dass

$$X_1, \dots, X_n \text{ i.i.d. } \sim \text{kumulative Verteilungsfkt. } F.$$

Wegen dem zweiten "i" in i.i.d. hat jedes  $X_i$  dieselbe Verteilung und dieselben Kennzahlen:  $\mathcal{E}(X_i) = \mu$ ,  $\text{Var}(X_i) = \sigma_X^2$ .

Zwei Kennzahlen von  $\bar{X}_n$  sind dann:

$$\begin{aligned} \mathcal{E}(\bar{X}_n) &= \mu, \\ \text{Var}(\bar{X}_n) &= \frac{\sigma_X^2}{n}. \end{aligned}$$

Der Erwartungswert von  $\bar{X}_n$  ist also gleich demjenigen einer einzelnen Zufallsvariablen  $X_i$ , die *Varianz nimmt jedoch ab falls  $n$  gross ist*. Daraus folgt das folgende Gesetz:

**Gesetz der Grossen Zahlen:** falls  $X_1, \dots, X_n$  i.i.d., dann

$$\bar{X}_n \longrightarrow \mu \quad (n \rightarrow \infty).$$

*Beispiel:*  $n$ -maliges Werfen eines Würfels

Betrachte  $X_i =$  Augenzahl im  $i$ -ten Wurf. Dann

$$\begin{aligned} \bar{X}_n &= \text{durchschnittliche Augenzahl bei } n \text{ Würfeln} \\ \longrightarrow \mu &= \mathcal{E}(X_i) = (1 + 2 + 3 + 4 + 5 + 6)/6 = 3.5 \quad (n \rightarrow \infty). \end{aligned}$$



Das heisst, bei sehr vielen Würfeln wird die durchschnittliche Augenzahl nahe bei 3.5 sein.

Die Verteilung von  $\bar{X}_n$  ist im allgemeinen schwierig anzugeben. Ein Spezialfall ist der folgende:

$$\bar{X}_n \sim \mathcal{N}(\mu, \sigma_X^2/n) \text{ falls } X_1, \dots, X_n \text{ i.i.d. } \sim \mathcal{N}(\mu, \sigma_X^2).$$

Falls die einzelnen  $X_i$ 's nicht normal-verteilt sind, so gilt erstaunlicherweise die obige Verteilungs-Formel immer noch approximativ. Dies liefert der folgende berühmte Satz.

**Zentraler Grenzwertsatz:** falls  $X_1, \dots, X_n$  i.i.d. , dann

$$\bar{X}_n \approx \mathcal{N}(\mu, \sigma_X^2/n),$$

wobei die Approximation im Allgemeinen besser wird mit grösserem  $n$ . Überdies ist auch die Approximation besser, je näher die Verteilung von  $X_i$  bei der Normal-Verteilung  $\mathcal{N}(\mu, \sigma_X^2)$  ist.

Ebenso ist die standardisierte Zufallsvariable

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma_X}$$

ungefähr  $\mathcal{N}(0, 1)$  verteilt.

## 4.6 Statistik für eine Stichprobe (Stahel, Kap. 8.3 – 8.5, 9.3)

Wir betrachten Daten  $x_1, \dots, x_n$  welche als Realisierungen von  $X_1, \dots, X_n$  i.i.d. aufgefasst werden. Zwei Kennzahlen der Zufallsvariablen  $X_i$  sind:  $\mathcal{E}(X_i) = \mu$  und  $\text{Var}(X_i) = \sigma_X^2$ . Typischerweise sind diese (und andere) Kennzahlen unbekannt, und man möchte Rückschlüsse aus den Daten darüber machen.

*Beispiel:* Blutplättchen-Aggregation (siehe Abschnitt 4.1)

Die Blutplättchen-Aggregation ist ein Beispiel eines sogenannten *gepaarten Vergleichs*, wo man bei jeder Versuchseinheit eine Grösse unter zwei verschiedenen Bedingungen misst. Von Interesse ist, ob ein systematischer Unterschied bezüglich der Aggregation vor und nach dem Rauchen einer Zigarette besteht. Um dies zu untersuchen bilden wir die *Differenzen*:  $x_i = \text{Aggregation "nachher"} - \text{Aggregation "vorher"}$  ( $i = 1, \dots, 11$ ), und wir haben somit eine (uns interessierende) Stichprobe.

### 4.6.1 (Punkt-) Schätzungen

Die (Punkt-) Schätzungen für den Erwartungswert und die Varianz sind:

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

Beachte dass die Schätzer hier als Funktionen der Zufallsvariablen  $X_1, \dots, X_n$  geschrieben sind: insbesondere sind  $\hat{\mu}$  und  $\hat{\sigma}_X^2$  selbst wieder Zufallsvariablen (die Verteilungseigenschaften von  $\hat{\mu}$  wurden in Abschnitt 4.5 diskutiert). Mit der Interpretation, dass die Daten  $x_i$  Realisierungen der Zufallsvariablen  $X_i$  sind, sind die realisierten Schätzer gleich dem arithmetischen Mittel und der empirischen Varianz der Daten.

#### 4.6.2 Tests für $\mu$

*Beispiel:* Blutplättchen-Aggregation (Forts.)

Wir wollen testen, ob ein systematischer Unterschied zwischen Aggregation “nachher” und Aggregation “vorher” besteht. Da  $x_i$  gerade die Differenz der Aggregationen zwischen “nachher” und “vorher” ist, betrachten wir das folgende Test-Problem:

$$H_0 : \mu = 0, \quad H_A : \mu \neq 0.$$

Um auf den Parameter  $\mu$  zu testen, machen wir vorerst einmal die Annahme, dass

$$X_1, \dots, X_n \text{ i.i.d. } \mathcal{N}(\mu, \sigma_X^2). \quad (4.3)$$

Eine Abschwächung dieser Annahme wir später diskutiert.

#### Der z-Test

Wir nehmen an, dass die Daten  $x_1, \dots, x_n$  Realisierungen von (4.3) sind. Überdies machen wir die Annahme, dass  $\sigma_X^2$  bekannt ist.

Der z-Test für den Parameter  $\mu$  ist dann wie folgt.

1. Spezifiziere die Nullhypothese  $H_0 : \mu = \mu_0$  und die Alternative  $H_A : \mu \neq \mu_0$  (oder “<” oder “>”).

2. Lege das Signifikanzniveau  $\alpha$  fest (z.B.  $\alpha = 0.05$ ).

3. Betrachte die **Teststatistik**  $\bar{X}_n$ . Unter der Nullhypothese gilt (siehe Abschnitt 4.5):

$$\bar{X}_n \sim \mathcal{N}(\mu_0, \sigma_X^2/n).$$

Der Verwerfungsbereich für die Teststatistik  $\bar{X}_n$  für die 2-seitige Alternative  $H_A : \mu \neq \mu_0$  ist dann

$$K = (-\infty, \mu_0 - \Phi^{-1}(1 - \alpha/2)\sigma_X/\sqrt{n}] \cup [\mu_0 + \Phi^{-1}(1 - \alpha/2)\sigma_X/\sqrt{n}, \infty).$$

Somit folgt mit einfacher Umformung, dass

$$P_{H_0}[\bar{X}_n \in K] = P_{\mu_0}[|\bar{X}_n - \mu_0| > \frac{\sigma_X}{\sqrt{n}}\Phi^{-1}(1 - \frac{\alpha}{2})] = \alpha,$$

das heisst, die Wahrscheinlichkeit eines Fehlers 1. Art ist gerade gleich dem Signifikanzniveau  $\alpha$ .

4. Verwerfe  $H_0$  falls das arithmetische Mittel der Daten  $\bar{x}_n \in K$  (ansonsten belasse  $H_0$ ).

Zusammenfassend ist der z-Test wie folgt:

$$\begin{aligned} \text{verwerfe } H_0, \text{ falls } & \left| \frac{\sqrt{n}(\bar{x}_n - \mu_0)}{\sigma_X} \right| > \Phi^{-1}(1 - \alpha/2) \quad \text{bei } H_A : \mu \neq \mu_0, \\ & \frac{\sqrt{n}(\bar{x}_n - \mu_0)}{\sigma_X} < -\Phi^{-1}(1 - \alpha) \quad \text{bei } H_A : \mu < \mu_0, \\ & \frac{\sqrt{n}(\bar{x}_n - \mu_0)}{\sigma_X} > \Phi^{-1}(1 - \alpha) \quad \text{bei } H_A : \mu > \mu_0. \end{aligned}$$

Im Unterschied zu den Tests in Kapitel 3.2.2 basiert der z-Test auf *mehreren* Beobachtungen. Diese werden aber mittels einer realisierten Teststatistik  $\bar{x}_n$ , oder auch in der standardisierten Form

$$z = \frac{\sqrt{n}(\bar{x}_n - \mu_0)}{\sigma_X}$$

zusammengefasst (eine Funktion der Daten). Ansonsten sind die Konzepte genau gleich wie in Kapitel 3.2.2: insbesondere brauchen wir für die Bestimmung des Verwerfungsbereichs  $K$  die Verteilung der Zufallsvariablen unter der Nullhypothese  $H_0 : \mu = \mu_0$ :

$$Z = \frac{\sqrt{n}(\bar{X}_n - \mu_0)}{\sigma_X} \sim \mathcal{N}(0, 1).$$

### Der t-Test

Wir vorhin nehmen wir an, dass die Daten Realisierungen von (4.3) sind. In der Praxis ist die Annahme, dass  $\sigma_X$  bekannt ist, oftmals unrealistisch. Wir können aber die Schätzung

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

benutzen. Dies führt aber zu einer zusätzlichen Unsicherheit, welche berücksichtigt werden muss.

Die Teststatistik beim t-Test ist

$$t = \frac{\sqrt{n}(\bar{x}_n - \mu_0)}{\hat{\sigma}_X},$$

und deren Verteilung unter der Nullhypothese  $H_0 : \mu = \mu_0$  ist

$$T = \frac{\sqrt{n}(\bar{X}_n - \mu_0)}{\hat{\sigma}_X} \sim t_{n-1},$$

wobei  $t_{n-1}$  eine sogenannte t-Verteilung mit  $n - 1$  Freiheitsgraden ist.

Die  $t_\nu$ -Verteilung ist eine symmetrische Verteilung um 0, welche langschwänziger ist als die Standard-Normalverteilung  $\mathcal{N}(0, 1)$ . Für  $T \sim t_\nu$  gilt:

$$\begin{aligned} \mathcal{E}(T) &= 0 \\ \text{Var}(T) &= \frac{\nu}{\nu - 2}. \end{aligned}$$

Für grosse  $\nu$  ist  $t_\nu$  ähnlich zu  $\mathcal{N}(0, 1)$ : insbesondere strebt die  $t_\nu$ -Verteilung gegen die Standard-Normalverteilung  $\mathcal{N}(0, 1)$  falls  $\nu \rightarrow \infty$ . Abbildung 4.7 zeigt die Dichte einer  $t_5$ -Verteilung.

Zusammenfassend ist der  $t$ -Test wie folgt:

$$\begin{aligned} \text{verwerfe } H_0 : \mu = \mu_0, \text{ falls } & |t| = \left| \frac{\sqrt{n}(\bar{x}_n - \mu_0)}{\hat{\sigma}_X} \right| > t_{n-1, 1-\alpha/2} \quad \text{bei } H_A : \mu \neq \mu_0, \\ & t = \frac{\sqrt{n}(\bar{x}_n - \mu_0)}{\hat{\sigma}_X} < -t_{n-1, 1-\alpha} \quad \text{bei } H_A : \mu < \mu_0, \\ & t = \frac{\sqrt{n}(\bar{x}_n - \mu_0)}{\hat{\sigma}_X} > t_{n-1, 1-\alpha} \quad \text{bei } H_A : \mu > \mu_0, \end{aligned}$$

wobei  $t_{n-1, \alpha}$  das  $\alpha$ -Quantile der  $t_{n-1}$ -Verteilung bezeichnet. Dieses Quantil ist tabelliert (siehe z. B. Stahel, Tabelle 8.5.g, p. 187) oder kann mittels Computer berechnet werden. Es ist etwas grösser als das Quantil der Normalverteilung und ergibt daher einen etwas kleineren Verwerfungsbereich. Für grosse  $n$  ist der Unterschied allerdings minim (da  $t_{n-1} \approx \mathcal{N}(0, 1)$  falls  $n$  gross).

Abbildung 4.7 illustriert den Verwerfungsbereich des  $t$ -Tests bei  $n = 6$  Beobachtungen. Der P-Wert bei 2-seitiger Alternative  $H_A : \mu \neq \mu_0$  kann wie folgt berechnet werden:

$$P - \text{Wert} = 2 \left( 1 - F_{t_{n-1}} \left( \frac{\sqrt{n}|\bar{x}_n - \mu_0|}{\hat{\sigma}_X} \right) \right),$$

wobei  $F_{t_{n-1}}$  die kumulative Verteilungsfunktion der  $t$ -Verteilung mit  $n - 1$  Freiheitsgraden bezeichnet.

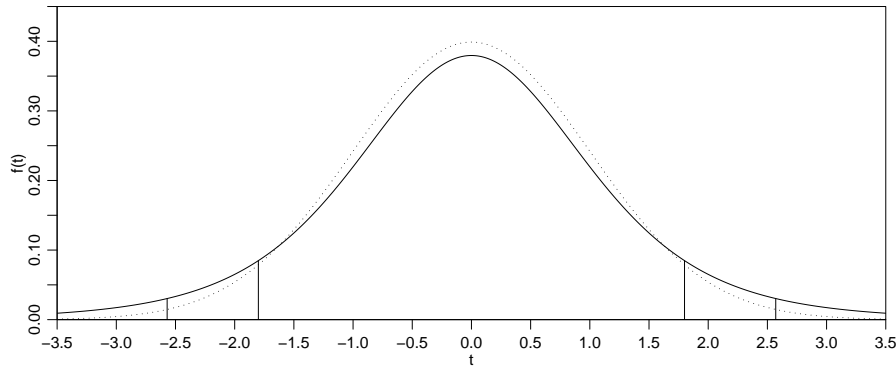


Abbildung 4.7: Dichte der  $t$ -Verteilung mit 5 Freiheitsgraden (ausgezogen) und der Normalverteilung (gestrichelt). Die Flächen ausserhalb der vertikalen Linien sind gleich 2.5% bzw. gleich dem halben P-Wert für einen hypothetischen Datensatz mit  $\sqrt{6}|\bar{x}_6 - \mu_0|/\hat{\sigma}_X = 1.8$ .

*Beispiel (Forts.)* Blutplättchen-Aggregation (siehe Abschnitt 4.1)

Wir betrachten die Differenzen  $x_i = \text{Aggregation "nachher"} - \text{Aggregation "vorher"}$  ( $i = 1, \dots, 11$ ) und fassen diese auf als i.i.d. Realisierungen von  $\mathcal{N}(\mu, \sigma_X^2)$ . Die interessierende Nullhypothese und Alternative ist  $H_0 : \mu = \mu_0 = 0$  und  $H_A : \mu > \mu_0 = 0$ . Die realisierte Teststatistik ist

$$\frac{\sqrt{n}(\bar{x}_n - \mu_0)}{\hat{\sigma}_X} = 4.27$$

und das relevante Quantil für  $\alpha = 0.05$  ist  $t_{10;0.95} = 1.812$ . Die Test-Entscheidung ist also: verwerfe  $H_0$  auf dem 5% Signifikanz-Niveau. Der P-Wert ist

$$P_{H_0}[T > 4.27] = 1 - F_{10}(4.27) = 0.00082.$$

Dies bedeutet, dass der Einfluss von dem Rauchen einer Zigarette bezüglich der Blutplättchen-Aggregation hoch signifikant ist.

Wenn man anstelle der einseitigen die zweiseitige Alternative  $H_A : \mu \neq \mu_0 = 0$  spezifiziert, so sind die Resultate wie folgt: das relevante Quantil für  $\alpha = 0.05$  ist  $t_{10;0.975} = 2.23$ . Die Test-Entscheidung bleibt dieselbe: verwerfe  $H_0$  auf dem 5% Signifikanz-Niveau. Der P-Wert ist

$$P_{H_0}[|T| > 4.27] = 2(1 - F_{10}(4.27)) = 0.0016,$$

### 4.6.3 Vertrauensintervall für $\mu$

Analog wie bei Zählraten in Kapitel 3.2.3 besteht das Vertrauensintervall aus denjenigen Werten  $\mu$ , bei denen der entsprechende Test nicht verwirft.

Wir nehmen wiederum an, dass die Daten Realisierungen von (4.3) sind. Dies führt dann auf die folgenden zweiseitigen Vertrauensintervalle (die dazugehörigen Tests sind zweiseitig mit Alternative  $H_A : \mu \neq \mu_0$ ) zum Niveau  $1 - \alpha$ :

$$\begin{aligned} \bar{x}_n \pm \Phi^{-1}(1 - \alpha/2) \frac{\sigma_X}{\sqrt{n}} & \text{ falls } \sigma_X \text{ bekannt,} \\ \bar{x}_n \pm t_{n-1, 1-\alpha/2} \frac{\hat{\sigma}_X}{\sqrt{n}} & \text{ falls } \sigma_X \text{ unbekannt.} \end{aligned}$$

*Beispiel (Forts.):* Aggregation von Blutplättchen

Wir haben 10 Freiheitsgrade und  $t_{10,0.975} = 2.23$ . Das zweiseitige Konfidenzintervall für die Erhöhung der Blutplättchen-Aggregation nach dem Rauchen einer Zigarette ist somit (%-ige Zunahme)

$$I = 10.27 \pm 2.23 \cdot 7.9761/\sqrt{11} = [4.91, 15.63].$$

Insbesondere ist die Null nicht im Intervall  $I$ : das heisst, der Wert  $\mu = 0$  ist nicht mit den Daten kompatibel (was wir bereits vom t-Test (siehe oben) wissen).

### 4.6.4 Tests für $\mu$ bei nicht-normalverteilten Daten

Der z- und t-Test sind optimal falls die Daten Realisierungen von normalverteilten Zufallsvariablen sind wie in (4.3). Optimalität bedeutet hier, dass dies die Tests sind, welche die beste Macht (-Kurve) haben, siehe unten.

Wir betrachten hier die allgemeinere Situation, wo die Daten Realisierungen sind von

$$X_1, \dots, X_n \text{ i.i.d. ,} \tag{4.4}$$

wobei  $X_i$  eine beliebige Verteilung hat. Wir bezeichnen mit  $\mu$  einen Lageparameter der Verteilung (z.B.  $\mu =$  Median der Verteilung von  $X_i$ ). Die Nullhypothese ist von der Form  $H_0 : \mu = \mu_0$ .

## Die Macht eines Tests

Wir haben in Kapitel 3.2.2 gesehen, dass es zwei Fehler für einen Test gibt:

Fehler 1. Art = fälschliches Verwerfen von  $H_0$ , obschon  $H_0$  richtig ist,

und

Fehler 2. Art( $\mu$ ) = (fälschliches) Beibehalten von  $H_0$  falls  $\mu \in H_A$  richtig ist.

Die Wahrscheinlichkeit für einen Fehler 1. Art ist gerade gleich  $\alpha$ ; beim Fehler 2. Art ( $\mu$ ) betrachtet man oft die Macht:

$$\text{Macht}(\mu) = 1 - P(\text{Fehler 2. Art}(\mu)) = P(\text{Verwerfen von } H_0 \text{ falls } \mu \text{ stimmt}).$$

Für  $\mu \in H_A$  kann man die Macht( $\mu$ ) interpretieren als die Chance, dass man richtigerweise  $H_A$  entdeckt falls  $\mu \in H_A$  stimmt. Für eine Teststatistik  $T$  und einen dazugehörigen Verwerfungsbereich  $K$  gilt dann:

$$\begin{aligned} P_{\mu_0}(T \in K) &= \alpha, \\ \text{Macht}(\mu) &= P_{\mu}(T \in K). \end{aligned}$$

## Der Vorzeichen-Test

Wir betrachten die Situation wo die Daten Realisierungen von (4.4) sind, wobei die einzelnen  $X_i$  nicht-normalverteilt sind. Wir bezeichnen hier mit  $\mu = \text{Median}$  der Verteilung von  $X_i$ ; im Falle einer symmetrischen Verteilung ist  $\mu = \mathcal{E}(X_i)$ .

Der Vorzeichen-Test benutzt die folgende Teststatistik:

$$V = \text{Anzahl } X_i\text{'s mit } (X_i > \mu_0).$$

Beachte dass  $V = \text{Anzahl positiver Vorzeichen von } (X_i - \mu_0)$  was den Namen des Tests erklärt.

Wir betrachten die Nullhypothese  $H_0 : p = P(X_i > \mu_0) = 1/2$  und die Alternative  $H_A : p \neq 1/2$  (oder einseitige Versionen). Somit ist unter der Nullhypothese  $H_0$  die Teststatistik  $V$  folgendermassen verteilt:

$$V \sim \text{Binomial}(n, 1/2),$$

und der Vorzeichentest wird somit zum Test für den Parameter  $p$  bei einer Binomialverteilung.

Der Vorzeichentest stimmt immer, falls die Daten Realisierungen von (4.4) sind: das heisst, die Wahrscheinlichkeit für einen Fehler 1. Art ist kontrolliert durch  $\alpha$  bei beliebiger Verteilung der  $X_i$ 's. Für den z- und t-Test stimmt dies nicht: wegen dem Zentralen Grenzwertsatz wird aber die Wahrscheinlichkeit für einen Fehler 1. Art approximativ kontrolliert durch  $\alpha$ , zumindest falls  $n$  gross ist.

Die Macht des z- oder t-Tests wird aber im Allgemeinen schnell schlecht, wenn die  $X_i$ 's in (4.4) nicht mehr normalverteilt sind. Deshalb hat bei nicht-normalverteilten Daten (nicht in guter Approximation normalverteilt) der Vorzeichentest oftmals eine höhere Macht als

der z- oder t-Test. Ein Nachteil des Vorzeichentests ist, dass er die Information nicht ausnützt, um wieviel die  $X_i$  von dem Wert  $\mu_0$  abweichen (siehe die Definition der Teststatistik  $V$  oben).

*Beispiel (Forts.):* Blutplättchen-Aggregation

Die Nullhypothese ist  $H_0 : \mu = \mu_0 = 0$ . Die realisierte Teststatistik ist dann  $v = 10$  und der P-Wert bei einseitiger Alternative  $H_A : \mu > \mu_0 = 0$  ist 0.005 (beim t-Test war der P-Wert = 0.00082).

## Der Wilcoxon-Test

Der Wilcoxon-Test ist ein Kompromiss, der keine Normalverteilung voraussetzt wie der t-Test und die Information der Daten besser ausnützt als der Vorzeichen-Test.

Die Voraussetzung für den Wilcoxon-Test ist: die Daten sind Realisierungen von (4.4) wobei die Verteilung der  $X_i$ 's stetig und symmetrisch ist (symmetrische Dichte um  $\mu = \mathcal{E}(X_i)$ ). Der P-Wert bei ein- oder zwei-seitiger Alternative kann mittels Computer berechnet werden.

Der Wilcoxon-Test ist in den allermeisten Fällen vorzuziehen: er hat in vielen Situationen oftmals bessere Macht als der t- und als der Vorzeichen-Test. Nur falls die Daten sehr gut mit einer Normalverteilung beschrieben werden ist der t-Test für gute Datenanalyse "vollumfänglich tauglich": diese Annahme oder Bedingung kann man z.B. mit dem Normal-Plot (siehe Kap. 4.4.6) grafisch überprüfen.

*Beispiel (Forts.):* Blutplättchen-Aggregation

Die Nullhypothese ist  $H_0 : \mu = \mu_0 = 0$ . Der P-Wert bei einseitiger Alternative  $H_A : \mu > \mu_0 = 0$  ist 0.002528.

## 4.7 Tests bei zwei unabhängigen Stichproben (Stahel, Kap. 8.8)

Oft möchte man einen Vergleich zweier Methoden (Gruppen, Versuchsbedingungen, Behandlungen) hinsichtlich der Lage der Verteilung machen.

### 4.7.1 Gepaarte und ungepaarte Stichproben

In allen Anwendungen ist neben der Auswertung auch die korrekte Planung des Versuches wichtig. Man muss sicherstellen, dass eventuelle Unterschiede tatsächlich durch die verschiedenen Methoden und nicht durch eine andere Störgrösse verursacht sind. Die beiden wichtigsten Prinzipien dazu sind *Blockbildung* und *Randomisierung*.

Randomisierung bedeutet hier, dass man die Reihenfolge der Versuche und die Zuordnung von Versuchseinheit zu Versuchsbedingung zufällig wählt: man hat dann Beobachtungen (realisierte Zufallsvariablen)

$x_1, x_2, \dots, x_n$  unter Versuchsbedingung 1,

$y_1, y_2, \dots, y_m$  unter Versuchsbedingung 2.

Im Allgemeinen ist  $m \neq n$ , aber nicht notwendigerweise. Bei solch **zufälliger Zuordnung von verschiedenen Versuchseinheiten zu zwei verschiedenen Versuchsbedingungen** spricht man von einer **ungepaarten Stichprobe**.

*Beispiel:*

Zufällige Zuordnung von 100 Testpatienten zu Gruppe der Grösse 60 mit Medikamenten-Behandlung und zu anderer Gruppe der Grösse 40 mit Placebo-Behandlung.

*Beispiel:*

Datensatz zu latenter Schmelzwärme von Eis in Kapitel 4.1.

Andererseits liegt eine **gepaarte Stichprobe** vor, wenn **beide Versuchsbedingungen an derselben Versuchseinheit eingesetzt** werden. Die Daten sind dann von der folgenden Struktur:

$$\begin{aligned}x_1, \dots, x_n &\text{ unter Versuchsbedingung 1,} \\y_1, \dots, y_n &\text{ unter Versuchsbedingung 2.}\end{aligned}$$

Notwendigerweise ist dann die Stichprobengrösse  $n$  für beide Versuchsbedingungen dieselbe.

*Beispiel:*

Datensatz zu Blutplättchen-Aggregation, siehe Kapitel 4.1.

## 4.7.2 Gepaarte Tests

Bei der Analyse von gepaarten Vergleichen arbeitet man mit den Differenzen innerhalb der Paare,

$$u_i = x_i - y_i \quad (i = 1, \dots, n),$$

welche wir als Realisierungen von i.i.d. Zufallsvariablen  $U_1, \dots, U_n$  auffassen. Kein Unterschied zwischen den beiden Versuchsbedingungen heisst dann einfach  $E[U_i] = 0$  (oder auch  $\text{Median}(U_i) = 0$ ). Tests dafür sind in Kapitel 4.6 beschrieben. Dabei ist zu beachten, dass die vorausgesetzte Symmetrie für die Verteilung von  $U_i$  beim Wilcoxon-Test immer gilt unter der Nullhypothese, dass  $X_i$  und  $Y_i$  dieselbe Verteilung haben.

## 4.7.3 Ungepaarte Tests

Bei ungepaarten Stichproben hat man Daten  $x_1, \dots, x_n$  und  $y_1, \dots, y_m$  (siehe Kapitel 4.7.1), welche wir als Realisierungen der folgenden Zufallsvariablen auffassen:

$$\begin{aligned}X_1, \dots, X_n &\text{ i.i.d. ,} \\Y_1, \dots, Y_m &\text{ i.i.d. ,}\end{aligned} \tag{4.5}$$

wobei auch alle  $X_i$ 's von allen  $Y_j$ 's unabhängig sind.

## 4.7.4 Zwei-Stichproben t-Test bei gleichen Varianzen

Das einfachste Problem lässt sich unter folgender Annahme an (4.5) lösen:

$$\begin{aligned}X_1, \dots, X_n \text{ i.i.d.} &\sim \mathcal{N}(\mu_X, \sigma^2), \\Y_1, \dots, Y_m \text{ i.i.d.} &\sim \mathcal{N}(\mu_Y, \sigma^2).\end{aligned} \tag{4.6}$$



Die interessierende Null-Hypothese ist

$$H_0 : \mu_X = \mu_Y.$$

Der Zwei-Stichproben t-Test (bei gleichen Varianzen) verwirft dann die Nullhypothese  $H_0 : \mu_X = \mu_Y$ , falls

$$|T| = \frac{|\bar{X}_n - \bar{Y}_m|}{S_{pool} \sqrt{1/n + 1/m}} > t_{n+m-2, 1-\alpha/2} \text{ bei Alternative } H_A : \mu_X \neq \mu_Y,$$

$$T = \frac{\bar{X}_n - \bar{Y}_m}{S_{pool} \sqrt{1/n + 1/m}} > t_{n+m-2, 1-\alpha} \text{ bei Alternative } H_A : \mu_X > \mu_Y,$$

$$T = \frac{\bar{X}_n - \bar{Y}_m}{S_{pool} \sqrt{1/n + 1/m}} < t_{n+m-2, 1-\alpha} \text{ bei Alternative } H_A : \mu_X < \mu_Y.$$

Dabei ist

$$S_{pool}^2 = \frac{1}{n+m-2} \left( \sum_{i=1}^n (X_i - \bar{X}_n)^2 + \sum_{i=1}^m (Y_i - \bar{Y}_m)^2 \right)$$

die gepoolte Schätzung für die gemeinsame Varianz  $\sigma^2$ . Die Wahl des Nenners in der Teststatistik  $T$  ergibt sich aus

$$\text{Var}(\bar{X}_n - \bar{Y}_m) = \sigma^2 \left( \frac{1}{n} + \frac{1}{m} \right). \quad (4.7)$$

Beweis von (4.7):

1.  $\bar{X}_n$  und  $\bar{Y}_m$  sind unabhängig, weil alle  $X_i$ 's von allen  $Y_j$ 's unabhängig sind.

2. Wegen der Unabhängigkeit von  $\bar{X}_n$  und  $\bar{Y}_m$  gilt:

$$\text{Var}(\bar{X}_n - \bar{Y}_m) = \text{Var}(\bar{X}_n) + \text{Var}(-\bar{Y}_m) = \text{Var}(\bar{X}_n) + \text{Var}(\bar{Y}_m).$$

3.  $\text{Var}(\bar{X}_n) = \sigma^2/n$  und  $\text{Var}(\bar{Y}_m) = \sigma^2/m$ .

Somit ist mit Schritt 2:  $\text{Var}(\bar{X}_n - \bar{Y}_m) = \sigma^2(1/n + 1/m)$ .  $\square$

Die Herleitung des Zwei-Stichproben t-Tests ist wie folgt. Man ersetzt die unbekannte Differenz  $\mu_X - \mu_Y$  durch die Schätzung  $\bar{X}_n - \bar{Y}_m$  und beurteilt, ob diese Schätzung "nahe bei" 0 liegt ("weit weg von" 0 würde Evidenz für  $H_A$  bedeuten). Dies wird so quantifiziert, dass man durch die Wurzel der geschätzten Varianz dividiert, und dies als Teststatistik benutzt:

$$\begin{aligned} T &= \frac{\bar{X}_n - \bar{Y}_m}{\sqrt{\widehat{\text{Var}}(\bar{X}_n - \bar{Y}_m)}} \\ &= \frac{\bar{X}_n - \bar{Y}_m}{S_{pool} \sqrt{1/n + 1/m}}. \end{aligned}$$

Unter der Annahme (4.6) und der Null-Hypothese  $\mu_X = \mu_Y$  gilt dann:

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

Somit kommt man zu der oben angegebenen Entscheidungsregel, analog zum t-Test für eine Stichprobe, siehe Kapitel 4.6.2.

*Beispiel:* Schmelzwärme von Eis, siehe Kapitel 4.1.

Die Null-Hypothese sei  $H_0 : \mu_X = \mu_Y$  und wir betrachten die Alternative  $H_A : \mu_X \neq \mu_Y$ . Die Kennzahlen des Datensatzes sind:  $\bar{x}_{13} = 80.021$ ,  $\bar{y}_8 = 79.979$ ,  $s_{pool}^2 = 7.2 \cdot 10^{-4}$ . Damit hat die Testgröße den Wert 3.47 was deutlich grösser ist als das 97.5% Quantil  $t_{19, 0.975} = 2.093$ .

## 4.7.5 Weitere Zwei-Stichproben-Tests

### Zwei-Stichproben t-Test bei ungleichen Varianzen

Anstelle der Annahme in (4.6) gelte:

$$\begin{aligned} X_1, \dots, X_n \text{ i.i.d.} &\sim \mathcal{N}(\mu_X, \sigma_X^2), \\ Y_1, \dots, Y_m \text{ i.i.d.} &\sim \mathcal{N}(\mu_Y, \sigma_Y^2). \end{aligned}$$

Die Verallgemeinerung des Zwei-Stichproben t-Tests für ungleiche Varianzen  $\sigma_X^2 \neq \sigma_Y^2$  ist in der Literatur zu finden und in vielen statistischen Programmen implementiert.

### Zwei-Stichproben Wilcoxon-Test (Mann-Whitney Test)

Die Voraussetzungen für den Zwei-Stichproben Wilcoxon-Test, manchmal auch Mann-Whitney Test genannt, bezüglich (4.5) sind wie folgt:

$$\begin{aligned} X_1, \dots, X_n \text{ i.i.d.} &\sim \text{beliebige kumulative Verteilungsfunktion } F(\cdot), \\ Y_1, \dots, Y_m \text{ i.i.d.} &\sim F(\cdot - \delta). \end{aligned}$$

Dies bedeutet, dass die Verteilung von  $Y_j$  die um  $\delta$  verschobene Verteilung von  $X_i$  ist, denn:  $P(Y_j \leq x + \delta) = F_Y(x + \delta) = F_X(x + \delta - \delta) = F_X(x) = P(X_i \leq x)$ .

Die Berechnung des P-Werts eines Zwei-Stichproben Wilcoxon-Tests kann mittels Computer erfolgen. Aus den gleichen Gründen wie im Fall einer Stichprobe (siehe Kapitel 4.6.4) ist der Wilcoxon-Test im Allgemeinen dem t-Test vorzuziehen.

## 4.8\* Versuchsplanung (Stahel, Kap. 14.1 - 14.2)

Genauso wichtig wie die Auswertung der Daten sind Überlegungen, wie man die Daten gewinnen soll. Wir haben bisher vor allem Vergleiche zwischen zwei Behandlungen besprochen (gepaart oder ungepaart). Man soll dabei nie eine neue Behandlung vergleichen mit Resultaten für die Standardbehandlung aus früheren Studien. Es braucht immer eine **Kontrollgruppe** in der gleichen Studie, die sich möglichst wenig von der Gruppe mit der neuen Behandlung unterscheidet. Dann stellt sich natürlich die Frage, wie man die Zuordnung zu den beiden Gruppen durchführen soll. Im gepaarten Fall muss man analog entscheiden, in welcher Reihenfolge man die beiden Behandlungen durchgeführt werden. Systematische Unterschiede zwischen den Gruppen, bzw. systematische Effekte der Reihenfolge kann man am besten vermeiden, wenn man die Zuordnung zufällig macht (sogenannte **Randomisierung**). Zufällig heisst dabei nicht willkürlich, sondern mit Hilfe von Zufallszahlen.

Ein weiterer wichtiger Punkt ist, dass das Experiment wenn möglich **doppelblind** sein soll. Das heisst, dass weder die Person, welche die Behandlung durchführt oder beurteilt, noch die Versuchsperson die Gruppenzugehörigkeit kennen. Dies ist nötig, um Nebeneffekte auszuschalten (nicht die Behandlung wirkt, sondern der mit der Behandlung verbundene Aufwand).

Nicht immer ist ein randomisiertes, doppelblindes Experiment möglich (aus ethischen oder praktischen Gründen). Dies erschwert die Auswertung und Interpretation unter Umständen gewaltig, weil man Störeffekte praktisch nicht ausschliessen kann. Ein bekanntes Beispiel ist der Zusammenhang zwischen Rauchen und Lungenkrebs, der lange umstritten war, weil die genetische Veranlagung und der Effekt des Lebensstils nicht auszuschliessen waren.



# Kapitel 5

## Regression

### 5.1 Korrelation und empirische Korrelation

Die gemeinsame Verteilung von abhängigen Zufallsvariablen  $X$  und  $Y$  ist i.A. kompliziert, und man begnügt man sich oft mit einer **vereinfachten** Kennzahl zur Beschreibung der Abhängigkeit. Die Kovarianz und Korrelation zwischen  $X$  und  $Y$  sind wie folgt definiert:

$$\begin{aligned}\text{Cov}(X, Y) &= E[(X - \mu_X)(Y - \mu_Y)] \quad (\text{Kovarianz}) \\ \text{Corr}(X, Y) &= \rho_{XY} = \text{Cov}(X, Y) / (\sigma_X \sigma_Y) \quad (\text{Korrelation}),\end{aligned}$$

wobei  $\sigma_X = \sqrt{\text{Var}(X)}$ , und analog für  $\sigma_Y$ .

Die Korrelation  $\rho_{XY}$  ist eine dimensionslose, normierte Zahl mit Werten  $\rho_{XY} \in [-1, 1]$ .

Die Korrelation misst Stärke und Richtung der **linearen Abhängigkeit** zwischen  $X$  und  $Y$ . Es gilt

$$\begin{aligned}\text{Corr}(X, Y) &= +1 \text{ genau dann wenn } Y = a + bX \text{ für ein } a \in \mathbb{R} \text{ und ein } b > 0, \\ \text{Corr}(X, Y) &= -1 \text{ genau dann wenn } Y = a + bX \text{ für ein } a \in \mathbb{R} \text{ und ein } b < 0.\end{aligned}$$

Überdies gilt:

$$X \text{ und } Y \text{ unabhängig} \implies \text{Corr}(X, Y) = 0. \quad (5.1)$$

Die Umkehrung gilt i.A. nicht.

#### 5.1.1 Die empirische Korrelation

In Kapitel 4.2.2 und Abbildung 4.3 haben wir ein Beispiel gesehen, wo die Daten  $(x_1, y_1), \dots, (x_n, y_n)$  als Realisierungen von i.i.d. Zufallsvektoren  $(X_1, Y_1), \dots, (X_n, Y_n)$  aufgefasst werden können.

Die empirischen Korrelation ist dann

$$\widehat{\text{Corr}}(X, Y) = \hat{\rho}_{XY} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}.$$

Analog zur Korrelation gelten die folgenden Eigenschaften:

$$\begin{aligned} \hat{\rho}_{XY} &\in [-1, 1], \\ \hat{\rho}_{XY} = +1 &\Leftrightarrow y_i = a + bx_i \text{ f\"ur alle } i = 1, \dots, n, \text{ und f\"ur ein } a \in \mathbb{R} \text{ und ein } b > 0, \\ \hat{\rho}_{XY} = -1 &\Leftrightarrow y_i = a + bx_i \text{ f\"ur alle } i = 1, \dots, n, \text{ und f\"ur ein } a \in \mathbb{R} \text{ und ein } b < 0. \end{aligned}$$

## 5.2 Einfache lineare Regression

Wir betrachten das folgende Beispiel aus der Chemie. Die Dimerisation von 1,3-Butadien l\"auft nach einem Reaktionsmodell zweiter Ordnung ab und ist deshalb charakterisiert durch die Gleichung  $\frac{d}{dt}C(t) = -\kappa C(t)^2$ , wobei  $C$  hier den Partialdruck des Edukts und  $t$  die Zeit bedeutet. Diese Gleichung hat L\"osungen der Form

$$\frac{1}{C(t)} = \frac{1}{C(0)} + \kappa t.$$

Messungen zu verschiedenen Zeitpunkten im Ablauf der Reaktion sind in Abbildung 5.1 dargestellt. Die letzte Gleichung zeigt, dass der Kehrwert des Partialdruckes linear von der Zeit abh\"angen sollte. Wegen zuf\"alligen Messfehlern und kleinen systematischen Abweichungen vom einfachen Modell liegen die Punkte nicht genau auf einer Geraden.

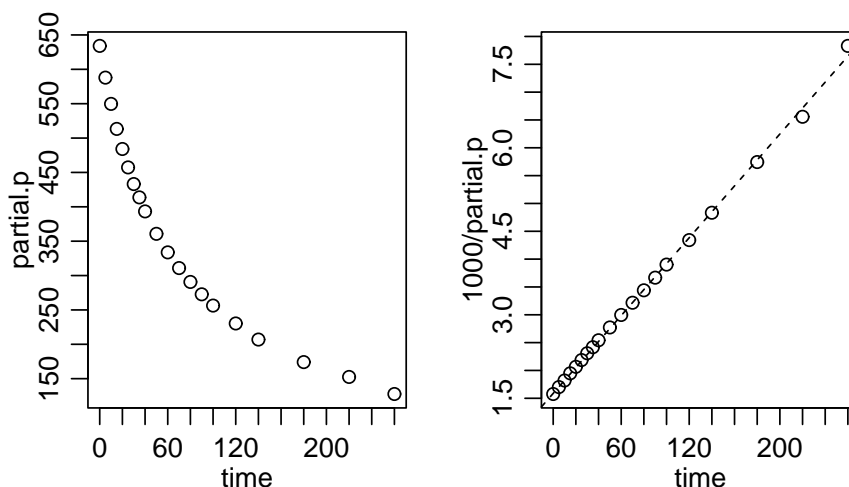


Abbildung 5.1: Partialdruck von Butadien (links) und Kehrwert  $1000 \times 1/\text{Partialdruck}$  (rechts), gegen die Zeit aufgetragen

### 5.2.1 Das Modell der einfachen linearen Regression

Im obigen Beispiel haben wir Daten

$$(x_1, y_1), \dots, (x_n, y_n),$$

wobei  $x_i$  der Zeitpunkt der  $i$ -ten Messung und  $y_i$  der Kehrwert des Partialdrucks der  $i$ -ten Messung bezeichnen. Diese fassen wir auf als Realisierungen des folgenden Modells:

$$\begin{aligned} Y_i &= h(x_i) + E_i \quad (i = 1, \dots, n), \\ E_1, \dots, E_n &\text{ i.i.d. , } \mathcal{E}(E_i) = 0, \text{ Var}(\mathcal{E}_i) = \sigma^2. \end{aligned}$$

Die  $Y$ -Variable ist die **Zielvariable** (engl: response variable) und die  $x$ -Variable ist die **erklärende Variable** oder **Co-Variable** (engl: explanatory variable; predictor variable; covariate). Die Zufallsvariablen  $E_i$  werden öfters als Fehler-Variablen oder Rausch-Terme bezeichnet. Sie besagen, dass der Zusammenhang zwischen der erklärenden und der Ziel-Variablen nicht exakt ist. Die erklärenden Variablen  $x_i$  ( $i = 1, \dots, n$ ) sind deterministisch, hingegen sind die Ziel-Variablen  $Y_i$  Zufallsvariablen (wegen den Zufalls-Variablen  $E_i$ ).

Die Modelle für die Funktion  $h(\cdot)$  sind:

$$\begin{aligned} h(x) &= \beta_0 + \beta_1 x && : \text{einfache lineare Regression,} \\ h(x) &= \beta_1 x && : \text{einfache lineare Regression durch Nullpunkt.} \end{aligned}$$

Meistens betrachten wir das allgemeinere Modell mit einem Achsenabschnitt  $\beta_0$ . Das Modell ist illustriert in Abbildung 5.2, wo für die Fehler-Variablen eine  $\mathcal{N}(0, 0.1^2)$ -Verteilung spezifiziert wurde.

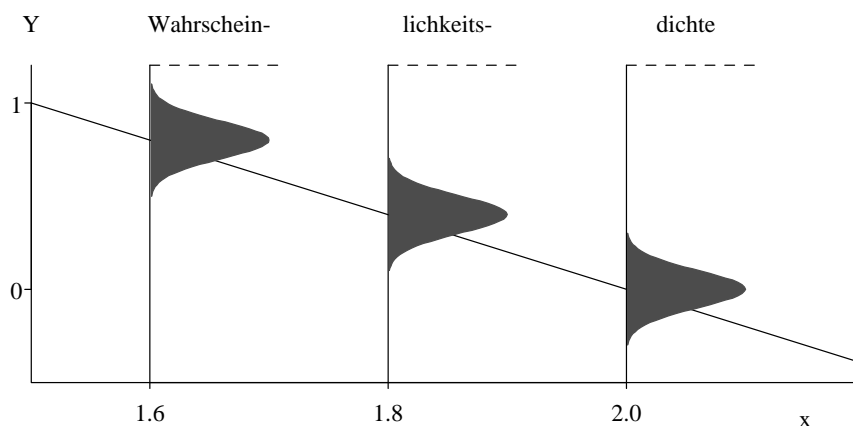


Abbildung 5.2: Veranschaulichung des Regressionsmodells  $Y_i = 4 - 2x_i + E_i$  mit  $E_i \sim \mathcal{N}(0, 0.1^2)$  für drei Beobachtungen.

## 5.2.2 Parameterschätzungen

Die unbekanntenen Modell-Parameter in der einfachen linearen Regression sind  $\beta_0$ ,  $\beta_1$  und auch die Fehlervarianz  $\sigma^2$ . Die Methode der Kleinsten-Quadrate liefert die folgenden Schätzungen:

$$\hat{\beta}_0, \hat{\beta}_1 \text{ sind Minimierer von } \sum_{i=1}^n (Y_i - (\beta_0 + \beta_1 x_i))^2.$$

Die Lösung dieses Optimierungsproblem ist eindeutig:

$$\begin{aligned} \hat{\beta}_1 &= \frac{\sum_{i=1}^n (Y_i - \bar{Y}_n)(x_i - \bar{x}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \\ \hat{\beta}_0 &= \bar{y}_n - \hat{\beta}_1 \bar{x}_n. \end{aligned}$$

Das Prinzip der Kleinsten-Quadrate liefert sogenannte erwartungstreue Schätzungen:

$$\mathcal{E}(\hat{\beta}_0) = \beta_0, \quad \mathcal{E}(\hat{\beta}_1) = \beta_1,$$

das heisst, dass die Schätzungen keinen systematischen Fehler haben (z.B. sind sie nicht systematisch zu gross was bedeuten würde, dass z.B.  $\mathcal{E}(\hat{\beta}_1) > \beta_1$ ).

Für die Schätzung von  $\sigma^2$  benützen wir das Konzept der Residuen. Falls wir Realisationen der Fehler-Terme  $E_i$  beobachten könnten, so könnten wir die empirische Varianzschätzung für  $\sigma^2$  verwenden. Hier approximieren wir zuerst die unbeobachteten Fehler-Variablen  $E_i$  durch die **Residuen**:

$$R_i = Y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \quad (i = 1, \dots, n).$$

Da  $E_i = Y_i - (\beta_0 + \beta_1 x_i)$  scheint die Approximation  $R_i \approx E_i$  vernünftig. Als Varianzschätzung benutzt man dann:

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n R_i^2. \quad (5.2)$$

Dabei ist zu beachten, dass (bei einfacher linearer Regression mit einem Achsenabschnitt  $\beta_0$ ) gilt:  $\sum_{i=1}^n R_i = 0$ . Das heisst, die Varianzschätzung in (5.2) ist wie die empirische Varianz bei einer Stichprobe (siehe Kapitel 4.6.1), ausser dass wir den Faktor  $1/(n-2)$  anstelle von  $1/(n-1)$  nehmen. Dieser Faktor entspricht der folgenden Faustregel:  $1/(n - \text{Anzahl Parameter})$ , wobei die Anzahl Parameter ohne den zu schätzenden Varianz-Parameter zu zählen ist (in unserem Falle sind dies die Parameter  $\beta_0, \beta_1$ ).

Bei einem Datensatz mit realisierten  $y_i$  ( $i = 1, \dots, n$ ) werden die Schätzungen mit den Werten  $y_i$  anstelle von  $Y_i$  gerechnet. Zum Beispiel sind dann die realisierten Residuen von der Form  $r_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)$ .

### 5.2.3 Tests und Konfidenzintervalle

Wir diskutieren hier die 2. und 3. Grundfragestellung (siehe Kapitel 3.1) im Kontext der einfachen linearen Regression. Dabei werden wir entscheidend mehr Schlussfolgerungen ziehen, als bloss eine best passende Regressionsgerade zu finden.

#### Der t-Test in Regression

Wir betrachten hier als Beispiel den folgenden Datensatz. Es wurden  $n = 111$  Messungen gemacht von mittlerer täglicher Temperatur ( $x$ -Variable) und mittlerem täglichem Ozongehalt ( $Y$ -Variable). Die Daten und die Regressionsgerade  $\hat{\beta}_0 + \hat{\beta}_1 x$  sind in Abbildung 5.2.3 ersichtlich. Die interessierende Frage in der Praxis lautet: hat die Temperatur einen Einfluss auf den Ozongehalt. Diese Frage kann man in ein Test-Problem übersetzen:

$$\begin{aligned} H_0 &: \beta_1 = 0, \\ H_A &: \beta_1 \neq 0. \end{aligned}$$

Es wird hier “per default” ein zwei-seitiger Test durchgeführt, nämlich der t-Test für die Steigung in der einfachen linearen Regression.

Wir machen hier die Annahme, dass

$$E_1, \dots, E_n \text{ i.i.d. } \mathcal{N}(0, \sigma^2). \quad (5.3)$$



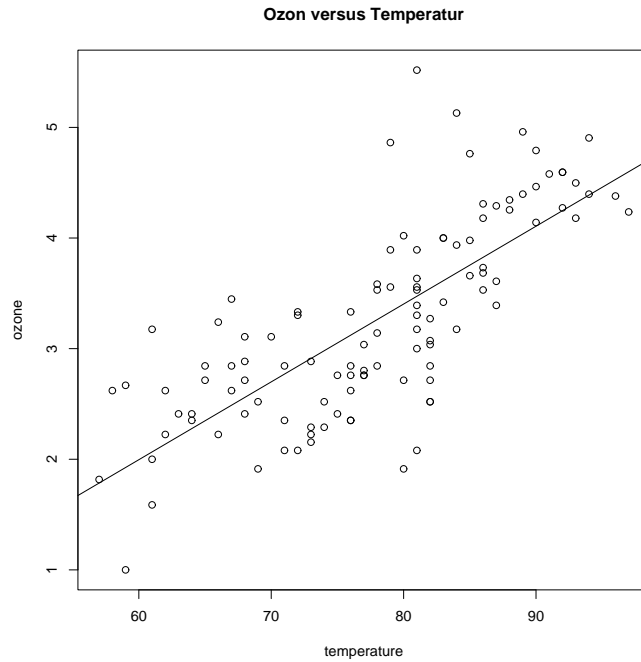


Abbildung 5.3: Streudiagramm und angepasste Regressionsgerade für den Ozon-Temperatur Datensatz.

Die Teststatistik ist

$$\frac{\hat{\beta}_1}{\widehat{s.e.}(\hat{\beta}_1)},$$

$$\widehat{s.e.}(\hat{\beta}_1) = \sqrt{\widehat{\text{Var}}(\hat{\beta}_1)} = \frac{\hat{\sigma}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}.$$

Unter der Nullhypothese und der Annahme von normalverteilten Fehlern in (5.3) gilt:

$$T \sim t_{n-2} \text{ unter } H_0 : \beta_1 = 0,$$

und der P-Wert dieses zwei-seitigen t-Test kann dann analog wie in Kapitel 4.6.2 berechnet werden (mit  $n - 2$  anstelle von  $n - 1$  Freiheitsgraden), und er wird auch von statistischer Software geliefert.

Völlig analog erhält man auch einen Test für  $H_0 : \beta_0 = 0$  bei zwei-seitiger Alternative  $H_A : \beta_0 \neq 0$ . Der entsprechende P-Wert, unter Annahme der Normalverteilung in (5.3) wird von statistischer Software geliefert.

Der Computer-Output vom R bei dem Anpassen einer einfachen linearen Regression für den Datensatz von Ozon als Funktion von Temperatur sieht wie folgt aus:

Call:

```
lm(formula = ozone ~ temperature)
```

Residuals:

```
Min      1Q  Median      3Q      Max
```

-1.49016 -0.42579 0.02521 0.36362 2.04439

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-2.225984	0.461408	-4.824	4.59e-06 ***
temperature	0.070363	0.005888	11.951	< 2e-16 ***

---

Residual standard error: 0.5885 on 109 degrees of freedom

Multiple R-Squared: 0.5672, Adjusted R-squared: 0.5632

F-statistic: 142.8 on 1 and 109 DF, p-value: < 2.2e-16

Die zweite Kolonne bei "Coefficients" beschreibt die Punktschätzer  $\hat{\beta}_i$  ( $i = 0, 1$ ); die dritte Kolonne die geschätzten Standardfehler  $\widehat{s.e.}(\hat{\beta}_i)$  ( $i = 0, 1$ ); die vierte Kolonne die Teststatistik  $\hat{\beta}_i / \widehat{s.e.}(\hat{\beta}_i)$  ( $i = 0, 1$ ), welche sich aus der zweiten dividiert durch die dritte Kolonne ergibt; die fünfte Kolonne bezeichnet den P-Wert für  $H_0 : \beta_i = 0$  und  $H_A : \beta_i \neq 0$  ( $i = 0, 1$ ). Überdies ist die geschätzte Standardabweichung für den Fehler  $\hat{\sigma}$  ersichtlich unter "Residual standard error"; die "degrees of freedom" sind gleich  $n - 2$ .

## Konfidenzintervalle

Basierend auf der Normalverteilungsannahme erhält man die folgenden zwei-seitigen Konfidenzintervalle für  $\beta_i$  ( $i = 0, 1$ ) zum Niveau  $1 - \alpha$ :

$$\begin{aligned} \hat{\beta}_0 \pm \widehat{s.e.}(\hat{\beta}_0) t_{n-2; 1-\alpha/2} & \text{ für } \beta_0, \\ \hat{\beta}_1 \pm \widehat{s.e.}(\hat{\beta}_1) t_{n-2; 1-\alpha/2} & \text{ für } \beta_1. \end{aligned}$$

### 5.2.4 Das Bestimmtheitsmass $R^2$

Die Güte eines Regressionsmodells kann mit dem sogenannten Bestimmtheitsmass  $R^2$  quantifiziert werden. Dazu betrachten wir eine Beziehungen zwischen verschiedenen Variationsquellen: mit der Bezeichnung  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$  für den Wert der angepassten Geraden beim Wert  $x_i$  ist

$$\underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{SS_Y} = \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{SS_E} + \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{SS_R}. \quad (5.4)$$

Dabei beschreibt  $SS_Y$  die totale Variation der Zielvariablen (ohne Einfluss der erklärenden Variablen  $x$ ),  $SS_E$  die Variation des Fehlers (Residuen-Quadratsumme), und  $SS_R$  die Variation, welche durch die Regression erklärt wird (Einfluss der erklärenden Variablen  $x$ ). Das Bestimmtheitsmass ist dann wie folgt definiert:

$$R^2 = \frac{SS_R}{SS_Y}, \quad (5.5)$$

und beschreibt den Anteil der totalen Variation, welche durch die Regression erklärt wird. Wegen 5.4 gilt, dass  $0 \leq R^2 \leq 1$ : falls  $R^2$  nahe bei 1 ist, so erklärt das Regressionsmodell viel der totalen Variation und ist somit gut; falls  $R^2 \approx 0$  taugt das Regressionsmodell nicht

besonders viel. Die Realisation von  $R^2$  ist im Computer-Output zu finden unter "Multiple R-squared".

Im Falle der einfachen linearen Regression gilt auch:

$$R^2 = \hat{\rho}_{XY}^2,$$

d.h.  $R^2$  ist gleich der quadrierten empirischen Korrelation.

### 5.2.5 Allgemeines Vorgehen bei einfacher linearer Regression

Grob zusammengefasst kann bei einfacher linearer Regression folgendemassen vorgegangen werden.

1. Anpassen der Regressionsgeraden; d.h. Berechnung der Punktschätzer  $\hat{\beta}_0, \hat{\beta}_1$ .
2. Testen ob erklärende Variable  $x$  einen Einfluss auf die Zielvariable  $Y$  hat mittels t-Test für  $H_0 : \beta_1 = 0$  und  $H_a : \beta_1 \neq 0$ . Falls dieser Test nicht-signifikantes Ergebnis liefert, so ist das Problem "in der vorliegenden Form uninteressant".
3. Testen ob Regression durch Nullpunkt geht mittels t-Test für  $H_0 : \beta_0 = 0$  und  $H_A : \beta_0 \neq 0$ . Falls dieser Test nicht-signifikantes Ergebnis liefert, so benützt man das kleinere Modell mit Regression durch Nullpunkt.
4. Bei Interesse Angabe von Konfidenzintervallen für  $\beta_0$  und  $\beta_1$ .
5. Angabe des Bestimmtheitsmass  $R^2$ . Dies ist in gewissem Sinne eine informellere (und zusätzliche) Quantifizierung als der statistische Test in Punkt 2.
6. Überprüfen der Modell-Voraussetzungen mittels Residuenanalyse. Dieser wichtige Schritt wird ausführlicher in Kapitel 5.2.6 beschrieben.

### 5.2.6 Residuenanalyse

Wir werden hier graphische Methoden beschreiben, basierend auf realisierten Residuen  $r_i (i = 1, \dots, n)$ , welche zur Überprüfung der Modell-Voraussetzungen für die einfache lineare Regression eingesetzt werden können. Die Modell-Voraussetzungen sind, in prioritärer Reihenfolge, die folgenden.

1.  $\mathcal{E}(E_i) = 0$ .  
Somit gilt  $\mathcal{E}(Y_i) = \beta_0 + \beta_1 x_i$ , das heisst: es gibt keinen systematischen Fehler im Modell.  
Abweichungen von dieser Annahme könnte zum Beispiel durch einen nicht-linearen Zusammenhang zwischen  $x$  und  $Y$  verursacht sein.
2.  $E_1, \dots, E_n$  i.i.d.  
Abweichungen könnte z.B. eine nicht-konstante Varianz der Fehlers sein, d.h.  $\text{Var}(E_i) = \sigma_i^2$  mit verschiedenen  $\sigma_i^2$  für  $i = 1, \dots, n$ . Eine andere Abweichung könnte durch korrelierte/abhängige Fehler verursacht sein.
3.  $E_1, \dots, E_n$  i.i.d.  $\mathcal{N}(0, \sigma^2)$ .  
Abweichungen könnte durch eine lang-schwänzige Fehlerverteilung verursacht sein.

## Der Tukey-Anscombe Plot

Der wichtigste Plot in der Residuenanalyse ist der Plot der Residuen  $r_i$  gegen die angepassten Werte  $\hat{y}_i$ , der sogenannte Tukey-Anscombe Plot.

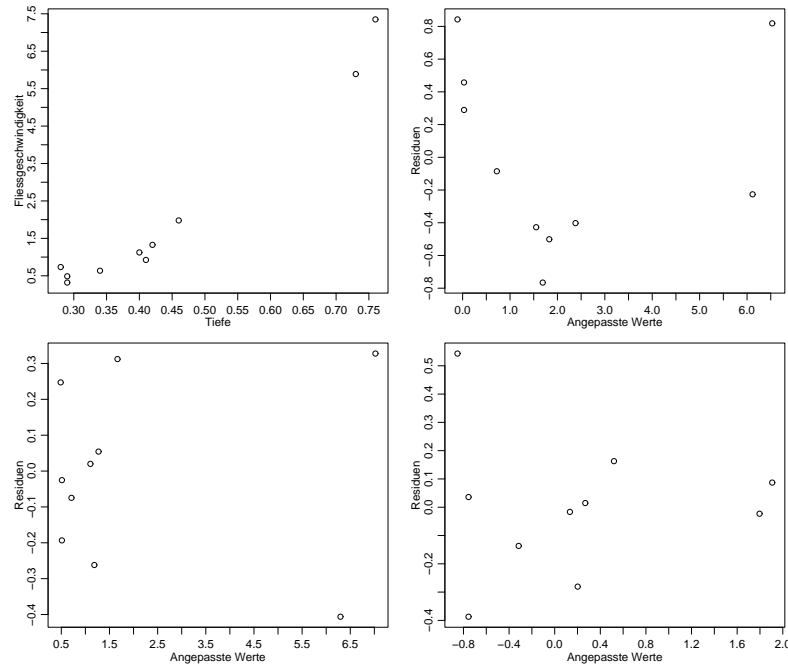


Abbildung 5.4: Streudiagramm von Tiefe und Fließgeschwindigkeit (oben links), Tukey-Anscombe Plots für einfache lineare Regression (oben rechts), für quadratische Regression (siehe Kapitel 5.3.1) (unten links) und für einfache lineare Regression mit logarithmierten Variablen  $\log(Y)$  und  $\log(x)$  (unten rechts).

Im Idealfall: gleichmässige Streuung der Punkte um Null.

Abweichungen:

- kegelförmiges Anwachsen der Streuung mit  $\hat{y}_i$

- evtl. kann man die Zielvariable logarithmieren (falls  $Y_i$ 's positiv sind), d.h. man benutzt das neue Modell

$$\log(Y_i) = \beta_0 + \beta_1 x_i + \varepsilon_i.$$

- Ausreisserpunkte

- evtl. können robuste Regressions-Verfahren verwendet werden (siehe Literatur)

- unregelmässige Struktur

- Indikation für nichtlinearen Zusammenhang

- evtl. Ziel und/oder erklärende Variablen transformieren (siehe auch das Beispiel in Abbildung 5.1).

Für den Ozon-Datensatz ist der Tukey-Anscombe Plot in Abbildung 5.5 gezeigt.

Nichtlineare Zusammenhänge können in der Praxis natürlich vorkommen: sie zeigen an, dass die Regressionsfunktion nicht korrekt ist. Abhilfe schaffen die Aufnahme zusätzlicher erklärender Variablen (z.B. quadratische Terme, siehe Kapitel 5.3.1) oder - wie bereits oben angedeutet - Transformationen der erklärenden und/oder der Ziel-Variablen. Ein einfaches

Beispiel ist in Abbildung 5.4 gezeigt, bei dem es um den Zusammenhang zwischen Tiefe und Fließgeschwindigkeit von Bächen geht. Bei einfacher Regression zeigt der Tukey-Anscombe Plot eine klare nichtlineare Struktur, die verschwindet, wenn man entweder einen quadratischen Term dazunimmt (siehe Kapitel 5.3.1) oder wenn man beide Variablen logarithmiert (d.h. einen Potenzzusammenhang anpasst mit dem Modell

$$\log(Y_i) = \beta_0 + \beta_1 \log(x_i) + \varepsilon_i \quad (i = 1, \dots, n).$$

Mit so wenigen Daten kann man zwischen diesen beiden Modellen nicht unterscheiden. Die Nichtlinearität des Zusammenhangs ist natürlich auch im ursprünglichen Streudiagramm ersichtlich, wenn man genau hinschaut. Häufig sind aber Abweichungen von der Linearität im Tukey-Anscombe Plot besser zu sehen.

### Plot bezüglich serieller Korrelation

Um die Unabhängigkeitsannahme der  $E_1, \dots, E_n$  zu überprüfen, kann der folgende Plot gemacht werden: plote  $r_i$  gegen die Beobachtungsnummer  $i$ .

Im Idealfall: gleichmässige Streuung der Punkte um Null.

Abweichungen:

- langfristiges Zonen mit durchwegs positiven oder negativen Residuen  
die Punktschätzungen sind immer noch OK, aber die Tests und Konfidenzintervalle stimmen nicht mehr evtl. Regression mit korrelierten Fehlern verwenden (siehe Literatur)

Für den Ozon-Datensatz ist der serielle Korrelations-Plot in Abbildung 5.5 gezeigt.

### Der Normalplot

Mit dem Normalplot (siehe Kapitel 4.4.6) können wir die Normalverteilungsannahme in (5.3) überprüfen.

Im Idealfall: approximativ eine Gerade

Abweichungen:

- Abweichung von einer Geraden Evtl. robuste Regression benutzen (siehe Literatur)

Für den Ozon-Datensatz ist der Normalplot in Abbildung 5.5 gezeigt.

### Das Auffinden eines guten Modells

Oftmals werden mehrere Modelle in einer Art “workflow-feedback” Prozeß betrachtet und angepasst. Man beginnt mit einem ersten Modell; dann, aufgrund von Residuenanalyse wird das Modell modifiziert. Das modifizierte Modell (immer noch als linear angenommen in evtl. transformierten Variablen) wird wiederum mit linearer Regression angepasst, und mit Residuenanalyse wird das neue Modell beurteilt. Dieses Vorgehen wird iteriert bis man ein “zufriedenstellendes” Modell gefunden und angepasst hat.

## 5.3 Multiple lineare Regression

Oftmals hat man mehrere erklärende Variablen  $x_{i,1}, \dots, x_{i,p-1}$  ( $p > 2$ ).

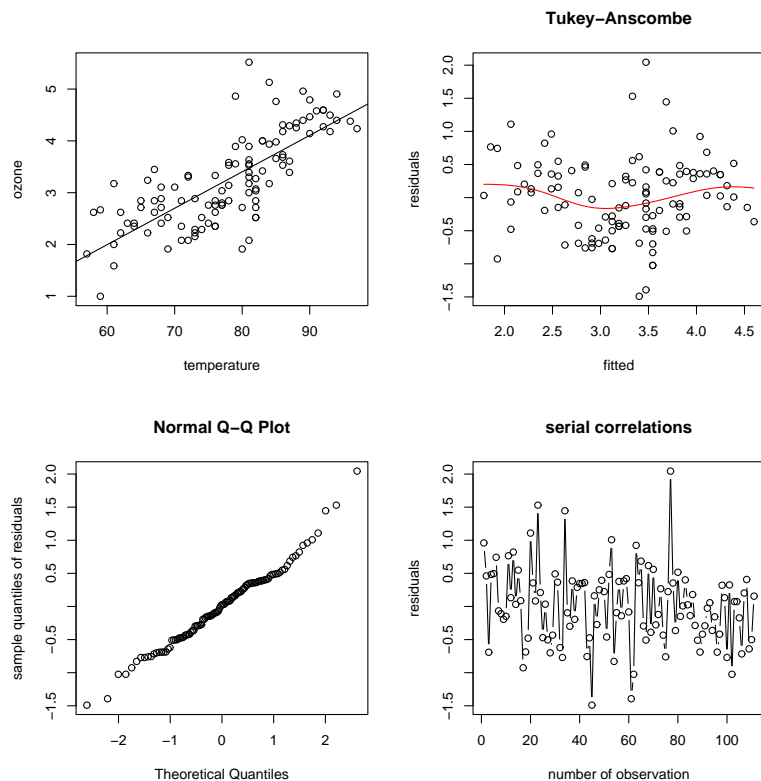


Abbildung 5.5: Ozon-Datensatz: Streudiagramm mit angepasster Regressiongerade (oben links); Tukey-Anscombe Plot (oben rechts); serieller Korrelations-Plot (unten links); Normalplot (unten rechts).

### 5.3.1 Das Modell der multiplen linearen Regression

Das Modell ist wie folgt:

$$Y_i = \beta_0 + \sum_{j=1}^{p-1} \beta_j x_{i,j} + E_i,$$

$$E_1, \dots, E_n \text{ i.i.d. , } \mathcal{E}(E_i) = 0, \text{ Var}(\mathcal{E}_i) = \sigma^2.$$

Wie bei der einfachen linearen Regression nehmen wir an, dass die erklärenden Variablen deterministisch sind. Es ist oftmals nützlich, das obige Modell in Matrix-Schreibweise darzustellen:

$$\begin{matrix} Y & = & X & \times & \beta & + & E \\ n \times 1 & & n \times p & & p \times 1 & & n \times 1 \end{matrix} \quad (5.6)$$

wobei  $X$  eine  $(n \times p)$ -Matrix ist mit Kolonnenvektoren  $(1, 1, \dots, 1)^T$ ,  $(x_{1,1}, x_{2,1}, \dots, x_{n,1})^T$  und letztendlich  $(x_{1,p-1}, x_{2,p-1}, \dots, x_{n,p-1})^T$ .

Beispiele von multipler linearer Regression sind unter anderen:

**Simple lineare Regression:**  $Y_i = \beta_0 + \beta_1 x_i + E_i$  ( $i = 1, \dots, n$ ).

$$p = 2 \quad X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}, \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}.$$

**Quadratische Regression:**  $Y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + E_i$  ( $i = 1, \dots, n$ ).

$$p = 3, \quad X = \begin{pmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{pmatrix}, \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix}.$$

Zu beachten ist, dass die Funktion quadratisch ist in den  $x_i$ 's, aber *linear* in den Koeffizienten  $\beta_j$  und deshalb ein Spezialfall des multiplen linearen Regressions Modells.

**Regression mit transformierten erklärenden Variablen:**

$Y_i = \beta_0 + \beta_1 \log(x_{i2}) + \beta_2 \sin(\pi x_{i3}) + E_i$  ( $i = 1, \dots, n$ ).

$$p = 3, \quad X = \begin{pmatrix} 1 & \log(x_{12}) & \sin(\pi x_{13}) \\ 1 & \log(x_{22}) & \sin(\pi x_{23}) \\ \vdots & \vdots & \vdots \\ 1 & \log(x_{n2}) & \sin(\pi x_{n3}) \end{pmatrix}, \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix}.$$

Wiederum, das Modell ist *linear* in den Koeffizienten  $\beta_j$ , aber nichtlinear in den  $x_{ij}$ 's.

### 5.3.2 Parameterschätzungen und t-Tests

Analog zur einfachen linearen Regression wird meist die Methode der Kleinsten Quadrate benutzt:

$$\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_{p-1} \text{ sind Minimierer von } \sum_{i=1}^n (Y_i - (\beta_0 + \beta_1 x_{i,1} + \dots + \beta_{p-1} x_{i,p-1}))^2.$$

Die eindeutige Lösung dieser Optimierung ist explizit darstellbar falls  $p < n$ :

$$\hat{\beta} = (X^T X)^{-1} X^T Y,$$

wobei  $\hat{\beta}$  den  $p \times 1$  Vektor  $(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_{p-1})^T$  bezeichnet, und  $X, Y$  wie in (5.6).

Die Schätzung der Fehlervarianz ist

$$\frac{1}{n-p} \sum_{i=1}^n R_i^2, \quad R_i = Y_i - (\hat{\beta}_0 + \sum_{j=1}^{p-1} \hat{\beta}_j x_{i,j}).$$

Unter der Annahme, dass die Fehler normalverteilt sind wie in (5.3), können auch ähnlich zur einfachen Regression t-Tests für die folgenden Hypothesen gemacht werden:

$$H_{0,j} : \beta_j = 0; \quad H_{A,j} : \beta_j \neq 0 \quad (j = 0, \dots, p-1).$$

Der wesentliche Unterschied besteht aber in der Interpretation der Parameter:

$\beta_j$  misst den linearen Effekt  
 der  $j$ -ten erklärenden Variablen auf die Zielvariable  $Y$   
**nach** Elimination der linearen Effekte  
 aller anderen Variablen auf  $Y$  ( $j = 1, \dots, p - 1$ )

Insbesondere impliziert dies, dass man die Koeffizienten  $\beta_j$  nicht einfach durch einzelne, individuelle simple lineare Regressionen von  $Y$  auf die  $j$ -te erklärende erhalten kann.

*Beispiel:* Wir betrachten  $p = 3$  und 2 erklärende Variablen. Wir nehmen an, dass die beiden erklärenden Variablen empirisch stark korreliert sind. Es kann dann durchaus geschehen, dass:

sowohl  $H_{0,1} : \beta_1 = 0$  als auch  $H_{0,2} : \beta_2 = 0$  werden nicht verworfen, obschon mindestens einer der Koeffizienten  $\beta_1$  oder  $\beta_2$  ungleich Null ist.

Um den Trugschluss zu vermeiden, dass es keine Effekt der erklärenden Variable auf die Ziel-Variable gibt, muss man den sogenannten F-Test betrachten.

### 5.3.3 Der F-Test

Der (globale) F-Test quantifiziert die Frage, ob es mindestens eine erklärende Variable gibt, welche einen relevanten Effekt auf die Zielvariable (im Sinne der linear Regression). Die folgende Nullhypothese wird beim (globalen) F-Test betrachtet:

$$H_0 : \beta_1 = \dots = \beta_{p-1} = 0$$

$$H_A : \text{mindestens ein } \beta_j \neq 0 \text{ (} j = 1, \dots, p - 1 \text{)}.$$

Der P-Wert des (globalen) F-Tests ist im Computer-Output gegeben unter "F-statistic".

### 5.3.4 Das Bestimmtheitsmass $R^2$

Das Bestimmtheitsmass  $R^2$  ist in der multiplen linearen Regression über die Formel (5.5) definiert (mit Hilfe der Zerlegung in (5.4)). Eine Interpretation im Sinne einer quadrierten Stichproben-Korrelation zwischen der Ziel-Variablen und den erklärenden Variablen lässt sich nicht mehr herstellen.

### 5.3.5 Residuenanalyse

Die Residuenanalyse geht völlig analog zu Kapitel 5.2.6. Das allgemeine Vorgehen bei multipler linearer Regression ist wie in Kapitel 5.2.5, unter Einbezug des F-Tests nach dem Schritt 1.

### 5.3.6 Strategie der Datenanalyse: ein abschliessendes Beispiel

Wir betrachten ein Beispiel wo die Asphalt-Qualität als Funktion von 6 erklärenden Variablen analysiert wird.

```
y = RUT : log("rate of rutting") = log(change of rut depth in inches
      per million wheel passes)
      ["rut" := 'Wagenspur', ausgefahrenes Geleise]
```



```

x1 = VISC : log(viscosity of asphalt)
x2 = ASPH : percentage of asphalt in surface course
x3 = BASE : percentage of asphalt in base course
x4 = RUN  : '0/1' indicator for two sets of runs.
x5 = FINES: 10* percentage of fines in surface course
x6 = VOIDS: percentage of voids in surface course

```

Die Daten sind in Abbildung 5.6 dargestellt. Die Zusammenhänge werden linearer, wenn

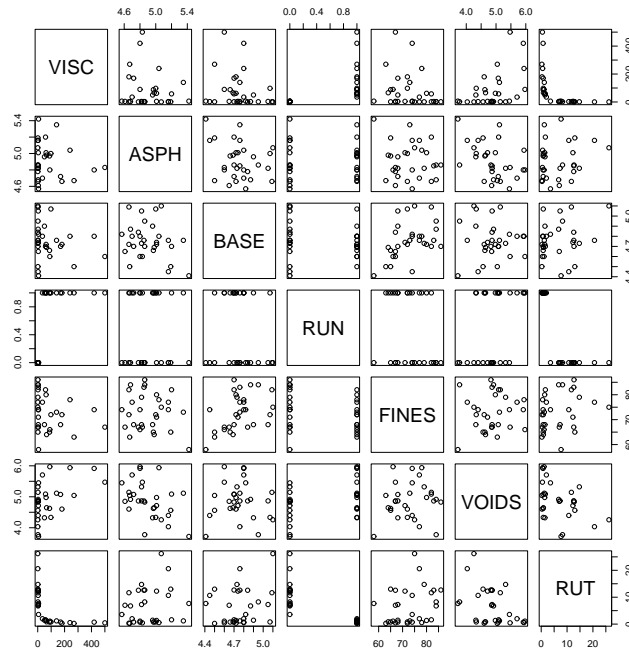


Abbildung 5.6: Paarweise Streudiagramme für den Asphalt-Datensatz. Die Zielvariable ist "RUT".

man die Zielvariable "RUT" logarithmiert und ebenfalls die erklärende Variable "VISC".

```

y = LOGRUT : log("rate of rutting") = log(change of rut depth in inches
per million wheel passes)
["rut":= 'Wagenspur', ausgefahrenes Geleise]
x1 = LOGVISC : log(viscosity of asphalt)
x2 = ASPH : percentage of asphalt in surface course
x3 = BASE : percentage of asphalt in base course
x4 = RUN : '0/1' indicator for two sets of runs.
x5 = FINES: 10* percentage of fines in surface course
x6 = VOIDS: percentage of voids in surface course

```

Die transformierten Daten sind in Abbildung 5.7 dargestellt.

Mittels R wird ein multiples lineares Modell angepasst. Der Output sieht wie folgt aus:

```

Call:
lm(formula = LOGRUT ~ ., data = asphalt1)

```

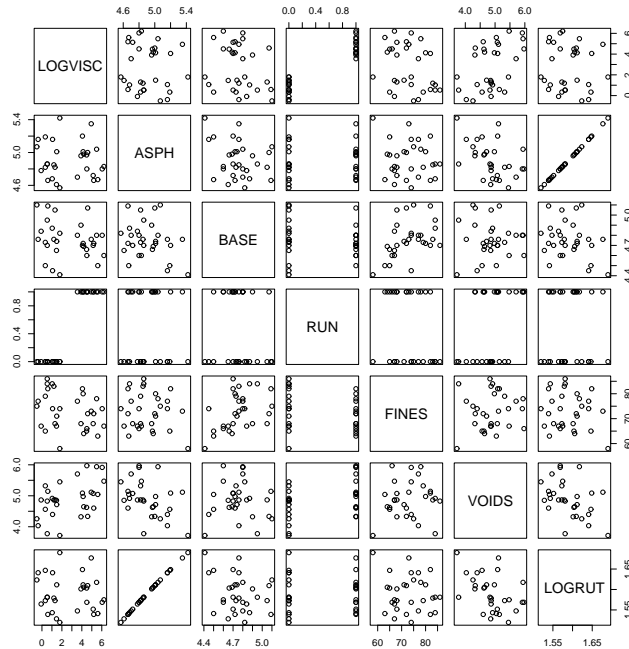


Abbildung 5.7: Paarweise Streudiagramme für den transformierten Asphalt-Datensatz. Die Zielvariable ist “LOGRUT”, die log-transformierte ursprüngliche Variable “RUT”. Die erklärende Variable “LOGVISC” ist ebenfalls die log-transformierte ursprüngliche Variable “VISC”.

Residuals:

	Min	1Q	Median	3Q	Max
	-0.48348	-0.14374	-0.01198	0.15523	0.39652

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	-5.781239	2.459179	-2.351	0.027280	*
LOGVISC	-0.513325	0.073056	-7.027	2.90e-07	***
ASPH	1.146898	0.265572	4.319	0.000235	***
BASE	0.232809	0.326528	0.713	0.482731	
RUN	-0.618893	0.294384	-2.102	0.046199	*
FINES	0.004343	0.007881	0.551	0.586700	
VOIDS	0.316648	0.110329	2.870	0.008433	**

---

Signif. codes: 0 ‘\*\*\*’ 0.001 ‘\*\*’ 0.01 ‘\*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 0.2604 on 24 degrees of freedom

Multiple R-Squared: 0.9722, Adjusted R-squared: 0.9653

F-statistic: 140.1 on 6 and 24 DF, p-value: < 2.2e-16

Wir sehen, dass die Variablen “LOGVISC”, “ASPH” und “VOID” signifikant oder sogar hoch-signifikant sind; die Variable “RUN” ist bloss schwach signifikant. Der F-Test ist hoch-signifikant, das Bestimmtheitsmass  $R^2$  sehr nahe bei 1. Die degrees of freedom sind

hier  $n - p = 24$  mit  $p = 7$ , d.h.  $n = 31$ . Die Residuenanalyse ist mittels Tukey-Anscombe und Normalplot in Abbildung 5.8 zusammengefasst: die Normalverteilungsannahme für die Fehler ist eine vernünftige Approximation. Der Tukey-Anscombe Plot zeigt etwas systematische Variation was durch Nichtlinearität induziert sein könnte; das  $R^2$  aber bereits sehr nahe bei 1 liegt, so kann man trotzdem sagen, dass die multiple lineare Regression sehr viel der totalen Variation erklären kann.

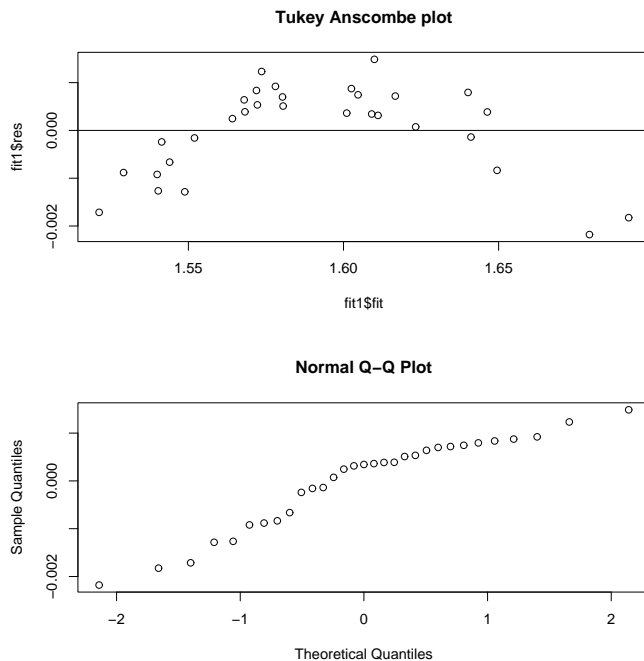


Abbildung 5.8: Tukey-Anscombe Plot (oben) und Normalplot (unten) beim Asphalt-Datensatz mit den transformierten Variablen “LOGRUT” und “LOGVISC”.

Ohne log-Transformationen, d.h. das untransformierte Modell wie in Abbildung 5.6, ist das Bestimmtheitsmass  $R^2 = 0.7278$ , also wesentlich schlechter als im transformierten Modell.