

Der Salzburger Jedermannlauf
oder: Die Anwendung der Statistik für
Spionagezwecke

Ferdinand Österreicher

FB Mathematik der Universität Salzburg

Lehrer/innen/fortbildungstag "West"
Innsbruck, 17. April 2008

Inhaltsverzeichnis

1	Einleitung	5
2	Punktschätzer	11
2.1	Schätzung des Parameters eines Alternativexperiments (A1)	11
2.2	Schätzung der Anzahl der Elemente einer Grundgesamtheit (A2)	13
2.3	Ausblick: Zur Korrektur einer Verfälschung bei Punktschätzern	18
3	Anhang: Stochastische Grundlagen	21
3.1	Rechenregeln für Erwartungswert und Varianz	21
3.2	Die Binomialverteilung	22
3.3	Die Inverse Hypergeometrische Verteilung	23

Kapitel 1

Einleitung

Die **Statistik** zerfällt in die beiden Teilgebiete

- **Beschreibende Statistik** und
- **Beurteilende Statistik**.

Während die mathematischen Kenntnisse und Fertigkeiten, welche für die beschreibende Statistik erforderlich sind, im Wesentlichen die vier Grundrechnungsarten sind, beruht die beurteilende Statistik auf stochastischen Modellen und setzt daher entsprechende Kenntnisse aus Wahrscheinlichkeitsrechnung voraus.

Die **Beurteilende Statistik** umfasst folgende beiden Teilgebiete: Das

- **Testen von Hypothesen** und das
- **Schätzen von Parametern**.

Das Testen von Hypothesen ist - auf den Punkt gebracht - "die stochastische Form des indirekten Schlusses" und, nach Auffassung des Autors, im Regelfall schwieriger zu vermitteln als das Schätzen von Parametern. Hilfreich für das Verständnis des Testens sind Kenntnis des Prinzips der Rechtssprechung im Strafrecht und/oder Vertrautheit mit einander (jedenfalls teilweise) widersprechend wissenschaftlichen Modellen, wie sie insbesondere in der Geschichte der Astronomie und der Physik anzutreffen sind.

Das Schätzen von Parametern zerfällt in die beiden Teilgebiete

- **Punktschätzer** und
- **Intervallschätzer** (insbesondere Konfidenz- oder Vertrauensintervalle).

Das Thema dieses Referats sind Punktschätzer für zwei verschiedene Anwendungssituationen.

”Anwendungsbeispiel” 1: Schätzung der Wahrscheinlichkeit für den Ausfall \perp eines Reißnagels

Man werfe einen bestimmten Reißnagel auf eine der beiden unten genauer beschriebenen Arten 100 mal, notiere der Reihe nach die Resultate der einzelnen Würfe und stelle schließlich die Folge der relativen Häufigkeiten

$$h_n = \frac{\text{Anzahl der ersten } n \text{ Würfe mit Ausfall } \perp}{n}, \quad n \in \{1, \dots, 100\}$$

graphisch dar.

Man benutze dabei einen Becher und eine ”Paschlwiese” und unterscheide folgende zwei Arten des Werfens:

(a) Der Reißnagel wird im Becher einige Male geschüttelt, wobei die Öffnung mit der Hand verdeckt wird. Anschließend wird der Becher verkehrt auf die Wiese gekippt, sodass der Reißnagel ohne Rollen auf der Wiese zu liegen kommt.

(b) Der Reißnagel wird im Becher einige Male geschüttelt, wobei die Öffnung mit der Hand verdeckt wird. Anschließend wird der Reißnagel bei schräg gehaltenem Becher so auf die Wiese gerollt, dass er deren Wand nicht berührt.

Man berücksichtige ferner, dass

- (i) stets dieselbe Person wirft und eine andere protokolliert,
- (ii) die Auswertung der Ergebnisse und deren graphische Darstellung erst nach Beendigung beider Experimente erfolgt.

Die relative Häufigkeit h_{100} ist der naheliegende Schätzer für die Wahrscheinlichkeit p des Ausfalls \perp .

Während es in diesem ”Anwendungsbeispiel” einen naheliegenden Schätzer für den unbekannten Parameter p gibt, gibt es Anwendungssituationen, für welche einem eine Fülle von Schätzern einfallen, sodass es notwendig ist, darüber nachzudenken, was einen ”guten Schätzer” ausmacht. Eine solche Anwendungssituation tritt beim ”Salzburger Jedermannlauf” auf. Sie hat - vom mathematischen Standpunkt betrachtet - dieselbe Struktur wie eine Anwendungssituation im Rahmen der Spionage im 2. Weltkrieg.

Dieses Referat hat zwei Hauptziele:

- den ”besten Schätzer” für die Anwendungssituation 2 zu finden und -
- im Zusammenhang damit -

· zwei Gütekriterien für Schätzer anzugeben, die es erlauben, Schätzer zu vergleichen.

Die Grundstruktur des Schätzens ist in folgender schematischen Darstellung widergegeben.

Schätzung	=	Wert des Parameters	+	Verfälschung	+	Zufallsschwankung
-----------	---	---------------------	---	--------------	---	-------------------

Es ist eine alte Tradition¹, das Schätzen von Parametern mit dem Zielschießen zu vergleichen, wobei folgende Entsprechungen gelten²:

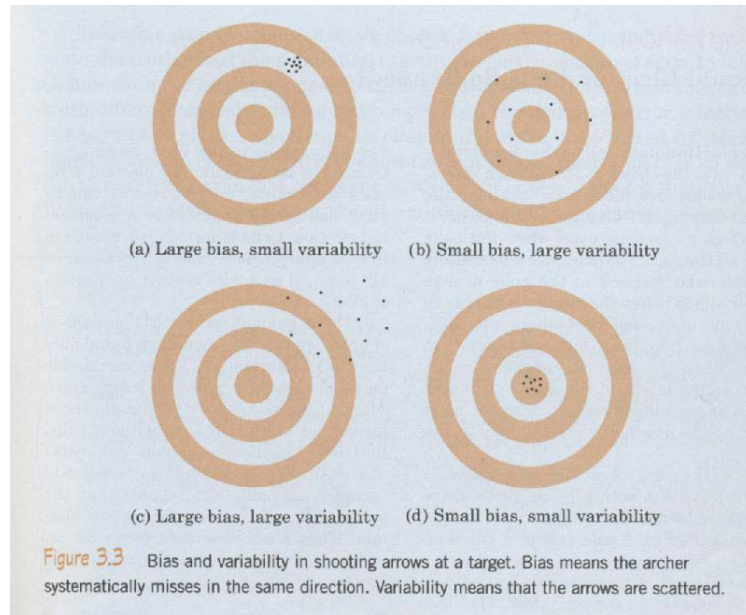
Schätzen	Zielschießen
Parameter	das Zentrum einer Zielscheibe
Verfälschung (<i>bias</i>)	Systematische Abweichung
unverfälscht (<i>unbiased</i>)	keine systematische Abweichung
kleine Varianz	hohe Präzision

Der grundsätzliche Unterschied besteht jedoch darin, dass die Zielscheibe beim Zielschießen sichtbar und daher bekannt, der Parameter beim Schätzen jedoch unbekannt ist. Daher entspricht ein Schätzer eher einer intelligenten Lenkwaffe, welche sich ihr Ziel selbst sucht, als der Kugel eines Gewehrs.

Die nachstehende Abbildung ist [5], Chapter 3: "What Do Samples Tell Us?" entnommen.

¹Diese geht vermutlich auf den englischen Astronomen *John Frederick William Herschel* (1792 – 1871) zurück, der sich in einem Artikel aus dem Jahre 1869 mit der Genauigkeit beim Bogenschießen beschäftigt hat. Von ihm stammen übrigens auch die Begriffe *Positiv*, *Negativ* und *Schnappschuss* in der Photographie.

²Man vergleiche dazu die Ausführungen über *bias* und *variability* in Section 3.4 "Towards Statistical Inference" in [6].



Ein Schätzer $\hat{\theta}$ ist eine Zufallsvariable, deren Formel den unbekannten Parameter θ selbstverständlich nicht enthält, deren Verteilung von diesem jedoch abhängt. Da die Gestalt der Verteilung von $\hat{\theta}$ maßgeblich durch Erwartungswert $E(\hat{\theta})$ und Varianz $V(\hat{\theta})$ bestimmt sind, sind unsere Gütekriterien durch diese beiden Größen bestimmt. Gemäß obiger schematischer Darstellung ist es naheliegend, die *Verfälschung* (den *bias*)³ des Schätzers $\hat{\theta}$ durch die Größe

$$b(\theta) = E_{\theta}(\hat{\theta}) - \theta$$

zu definieren.⁴ Ein Schätzer heißt *unverfälscht* oder *erwartungstreu* (*unbiased*), wenn $b(\theta) \equiv 0$ ist, oder, gleichbedeutend, wenn gilt

$$E_{\theta}(\hat{\theta}) = \theta \text{ für alle möglichen Parameter } \theta.$$

³Der Begriffe *Verfälschung* (*bias*) stammt - wie viele andere Begriffe beim Schätzen von Parametern - vom englischen Statistiker *Sir Ronald Aylmer Fisher* (1890 – 1962). Die systematische Untersuchung der sogenannten *unverfälschten oder erwartungstreuen Schätzer* (*unbiased estimators*) wurde in der Folge insbesondere vom schwedischen Mathematiker *Carl H. Cramér* (1893 – 1985) vorangetrieben.

⁴Das Subskript in $E_{\theta}(\hat{\theta})$ drückt aus, dass die Verteilung des Schätzers $\hat{\theta}$ vom Parameter θ abhängt.

Erstes Ziel wird es also sein, (a) von einem gegebenen Schätzer nachzuprüfen, ob er erwartungstreu ist, bzw. (b) - wenn möglich - einen erwartungstreuen Schätzer zu finden.

Für den Fall, dass wir für eine bestimmte Anwendungssituation mehrere erwartungstreue Schätzer zur Verfügung haben, werden wir von diesen naturgemäß jenen Schätzer ausfindig zu machen versuchen, der die kleinstmögliche Varianz besitzt.

Aus den bisherigen Überlegungen geht hervor, dass Erwartungswert und Varianz einer Zufallsvariablen wichtige Größen sind. Eine knappe einschlägige Zusammenstellung findet man in Abschnitt 3.1. Die wahrscheinlichkeitstheoretischen Voraussetzung für die beiden von uns betrachteten Anwendungssituationen findet man in den Abschnitten 3.2 und 3.3.

Kapitel 2

Punktschätzer

2.1 Schätzung des Parameters eines Alternativexperiments (A1)

Das "Anwendungsbeispiel" 1 hat folgende Grundstruktur:¹

Wir betrachten ein Zufallsexperiment mit zwei möglichen Ausgängen, die wir "*Erfolg*" und "*Misserfolg*" nennen. Die Wahrscheinlichkeit eines Erfolgs sei für jede Durchführung des Experiments gleich $p \in (0, 1)$. Es werden n solcher Zufallsexperimente unabhängig voneinander durchgeführt, die Ergebnisse X_i , $i \in \{1, \dots, n\}$ der einzelnen Experimente festgestellt und protokolliert. Dabei sei

$$X_i = 1 \quad \text{oder} \quad X_i = 0,$$

je nachdem, ob beim i -ten Experiment Erfolg oder Misserfolg eintritt.² Bezeichne $S_n = \sum_{i=1}^n X_i$ die beobachtete Anzahl der Erfolge. Dann ist der zugehörige Anteil $h_n = \frac{S_n}{n}$ der Erfolge ein Schätzer für die Wahrscheinlichkeit p .

Behauptung 1: Der Schätzer h_n ist (a) erwartungstreu und (b) besitzt die Varianz $\frac{p(1-p)}{n}$.

¹Auf tatsächliche Anwendungen wird im Beitrag "Einführung in die Praxis und Theorie der Stichproben" in [9] und ansatzweise auch in [10], Abschnitt 2.1.3 eingegangen.

²Für rationale $p \in (0, 1)$ gibt es ein Urnenmodell: Aus einer Urne, welche s schwarze und w weiße - und, von der Farbe abgesehen - gleichartige Kugeln enthält ($N = s + w$) werden n Kugeln mit Zurücklegen gezogen. $X_i = 1(0)$, je nachdem, ob beim i -ten Zug eine schwarze (weiße) Kugel gezogen wird, $i \in \{1, \dots, n\}$.

Beweis: (a) ist eine unmittelbare Konsequenz von (1) in Abschnitt 3.1 und Abschnitt 3.2:

$$E_p(h_n) = E_p\left(\frac{S_n}{n}\right) = \frac{1}{n} \cdot E_p(S_n) = \frac{1}{n} \cdot np = p \quad \forall p \in (0, 1).$$

(b) ist eine unmittelbare Konsequenz von (2) in Abschnitt 3.1 und Abschnitt 3.3:

$$V_p(h_n) = V_p\left(\frac{S_n}{n}\right) = \left(\frac{1}{n}\right)^2 \cdot V_p(S_n) = \frac{1}{n^2} \cdot np(1-p) = \frac{p(1-p)}{n} \quad \forall p \in (0, 1). \quad \square$$

In diesem Anwendungsbeispiel muss man ein wenig nachdenken, um - neben dem naheliegenden Schätzer h_n - weitere erwartungstreue Schätzer zu finden: h_n ist

$$h_n = \sum_{i=1}^n \frac{1}{n} X_i$$

eine gewichtete Summe der einzelnen Zufallsvariablen X_i . Es ist daher naheliegend, anstelle der Gewichte $\frac{1}{n}$ allgemeine Gewichte α_i , $i \in \{1, \dots, n\}$ zu verwenden. In diesem Zusammenhang gilt folgende Behauptung.

Behauptung 2: Sei

$$W_n = \{ \boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n) \in [0, 1]^n : \sum_{i=1}^n \alpha_i = 1 \}$$

die Menge aller Wahrscheinlichkeitsverteilungen auf der Menge $\{1, \dots, n\}$. Dann definiert

$$\hat{p}_{\boldsymbol{\alpha}} = \sum_{i=1}^n \alpha_i \cdot X_i, \quad \boldsymbol{\alpha} \in W_n$$

eine Familie von Schätzern, für welche gilt

- (a) alle Schätzer dieser Familie sind erwartungstreu und
- (b) für die Varianz jedes Schätzer $\hat{p}_{\boldsymbol{\alpha}}$, $\boldsymbol{\alpha} \in W_n$ gilt

$$\frac{p(1-p)}{n} \leq V(\hat{p}_{\boldsymbol{\alpha}}) \leq p(1-p),$$

wobei

- (i) die untere Schranke genau dann gilt, wenn $\hat{p}_\alpha = h_n$ ist und
- (ii) die obere Schranke genau dann, wenn gilt $\hat{p}_\alpha = X_i$, $i \in \{1, \dots, n\}$.

Beweis: Der hübsche Beweis, welcher die üblichen Rechenregeln für Erwartungswert und Varianz von Summen von Zufallsvariablen benützt, unterbleibt hier.

2.2 Schätzung der Anzahl der Elemente einer Grundgesamtheit (A2)

Den Umfang einer durchnummerierten Grundgesamtheit zu schätzen, ist ein eindrucksvolles Beispiel für eine Situation, bei welcher eine Vielzahl vernünftiger Schätzer zur Verfügung steht. Aus dieser ist ein Schätzverfahren auszuwählen, welches "möglichst genau" ist. Wir gehen dabei von folgendem Anwendungsbeispiel aus.³

Anwendungsbeispiel 2: Der Salzburger Jedermannlauf

Der "Salzburger Jedermannlauf" ist eine Breitensportveranstaltung, welche seit Jahren am Nationalfeiertag in der Landeshauptstadt Salzburg organisiert wird. Die Teilnehmer erhalten Startnummern $1, 2, \dots, N$, wobei N die uns unbekannte Teilnehmerzahl ist. Nach Beendigung des Laufes erfolgt eine Verlosung von 25 Preisen. Dafür werden 25 Startnummern zufällig und ohne Zurücklegen aus der Menge der Startnummern jener Läufer/innen gezogen, welche den Lauf beenden.

Von der Veranstaltung im Jahr 1994 wurden uns folgende 24 Nummern übermittelt.

616, 1436, 737, 11, 1133, 1003, 705, 139, 614, 665, 1057, 1076,
1070, 1075, 1382, 1384, 1394, 776, 650, 8, 688, 1065, 269, 195.

Versuchen Sie aus dieser Information die Teilnehmerzahl zu schätzen, indem Sie von der selbstverständlich nicht ganz realistischen Annahme ausgehen, dass die Startnummern lückenlos vergeben werden und alle Personen, welche eine Startnummer besitzen, auch den Lauf beenden und an der Verlosung teilnehmen.

³Auf weitere ähnlich strukturierte Anwendungen wird in [12] eingegangen.

Diese Idealisierung lässt sich durch folgendes Urnenmodell beschreiben.

Urnenmodell: Gegeben sei eine Urne mit N Jetons, welche von 1 bis N durchnummeriert sind. Der Parameter N sei unbekannt. n ($\leq N$) Jetons werden zufällig und ohne Zurücklegen gezogen. x_1, \dots, x_n seien die Nummern der gezogenen Jetons. N ist mit Hilfe der beobachteten Werte x_1, \dots, x_n zu schätzen.

Punktschätzer beim Ziehen ohne Zurücklegen

Beispiel: Jemand wählt eine Zahl $N \in \{10, \dots, 80\}$ und bestückt die Urne mit N Jetons. Wir sollen aufgrund einer Stichprobe vom Umfang $n = 5$ die gewählte Zahl N schätzen. Die Nummern der gezogenen Jetons sind

$$x_1 = \dots, x_2 = \dots, x_3 = \dots, x_4 = \dots, x_5 = \dots.$$

Wir haben folgende Schätzer, das sind Vorschriften, den Schätzwert zu bestimmen, erarbeitet.

1. Der Mittelschätzer

Aufgrund des Gesetzes der Großen Zahlen approximiert das Stichprobenmittel

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

für hinreichend große n den Erwartungswert $\mu = E(X_i) = \frac{N+1}{2}$ der Zufallsvariablen. D.h. es ist $\bar{X}_n \cong \frac{N+1}{2}$ und daher $N \cong 2\bar{X}_n - 1$. Diese für die Momentenmethode, die auf den englischen Statistiker *Karl Pearson* (1857 – 1936) zurückgeht, typische Überlegung motiviert den Mittelschätzer

$$\bar{N} = 2\bar{X}_n - 1.$$

Der entsprechende Schätzwert ist also $2 \cdot \bar{5} - 1 = \dots$.

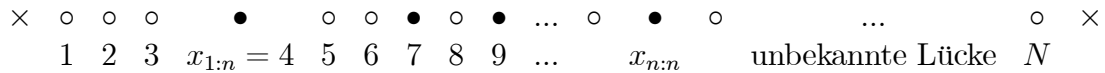
Im Folgenden bezeichnen $X_{1:n} < X_{2:n} < \dots < X_{n:n}$ die der Größe nach geordneten Stichprobenwerte. $X_{i:n}$ heißt dabei die *i-te Ordnungst Statistik*, $i \in \{1, \dots, n\}$.

2. Der Medianschätzer: Der Einfachheit halber sei $n = 2m+1$ mit $m \in \mathbb{N}_0$. Ersetzt man nun beim Mittelschätzer das Stichprobenmittel \bar{X}_n durch den Stichprobenmedian $X_{m+1:2m+1}$, so erhält man den Medianschätzer

$$\hat{N}_{m+1:2m+1} = 2 \cdot X_{m+1:2m+1} - 1.$$

Der zugehörige Schätzwert ist also: $2 \cdot \dots - 1 = \dots$.

3. Die Lückenmethode⁴:



Motivation: Wegen $X_{n:n} \leq N$ schätzt das Stichprobenmaximum den Parameter N in aller Regel zu kurz. Um diesen Bias zu korrigieren, ersetzen wir die Länge $N - X_{n:n}$ der unbekannten "Lücke" durch

(a) die bekannte Länge $X_{1:n} - 1$ der ersten Lücke und erhalten somit gemäß $N = X_{n:n} + N - X_{n:n} \cong X_{n:n} + X_{1:n} - 1$ den Lückenschätzer

$$\tilde{N}_1 = X_{n:n} + X_{1:n} - 1.$$

Der entsprechende Schätzwert ist also $\dots + \dots - 1 = \dots$.

(b) die durchschnittliche Länge

$$\frac{1}{n} [X_{1:n} - 1 + \sum_{i=2}^n (X_{i:n} - X_{i-1:n} - 1)] = \frac{1}{n} X_{n:n} - 1$$

der bekannten Lücken und erhalten somit gemäß $N = X_{n:n} + N - X_{n:n} \cong X_{n:n} + \frac{1}{n} X_{n:n} - 1 = \frac{n+1}{n} X_{n:n} - 1$ den Maximumschätzer

$$\hat{N}_{n:n} = \frac{n+1}{n} X_{n:n} - 1.$$

Der entsprechende Schätzwert ist also $\frac{6}{5} \cdot \dots - 1 = \dots$.

Anmerkung 1: Der Medianschätzer und der Maximumschätzer gehören der folgenden Familie von Schätzern an

$$\hat{N}_{i:n} = \frac{n+1}{i} \cdot X_{i:n} - 1, \quad i \in \{1, \dots, n\}.$$

Dabei gilt gemäß Abschnitt 3.3 $X_{i:n} \sim IH_{i,N,n}$ oder, in Worten, die i -te Ordnungsstatistik $X_{i:n}$ ist gemäß der *Inversen Hypergeometrischen Verteilung*

⁴In der nachstehenden Abbildung sind die beobachteten Werte durch \bullet gekennzeichnet. So ist beispielsweise $x_{1:n} = 4$.

mit den Parametern i, N und n verteilt. Deren Erwartungswert und Varianz sind, wie ebenfalls in Abschnitt 3.3 gezeigt wird,

$$\begin{aligned} (1) \quad E(X_{i:n}) &= i \cdot \frac{N+1}{n+1} \\ (2) \quad V(X_{i:n}) &= \frac{i(n+1-i)}{n+1} \cdot \frac{N+1}{n+1} \cdot \frac{N-n}{n+2}. \end{aligned}$$

Folgerungen: $X_{i:n}$ schätzt offensichtlich zu kurz. Aus (1) lässt sich aber leicht ein "unverfälschter" Schätzer für N ermitteln. Indem man nämlich $X_{i:n}$ mit dem Faktor $\frac{n+1}{i}$ multipliziert und schließlich 1 abzieht, erhält man den Schätzer $\hat{N}_{i:n}$ der oben angegebenen Familie. Dessen Erwartungswert ist wegen Abschnitt 3.1, (i) und (1) tatsächlich

$$\begin{aligned} E\left(\frac{n+1}{i}X_{i:n} - 1\right) &= \frac{n+1}{i} \cdot E(X_{i:n}) - 1 \\ &= \frac{n+1}{i} \cdot i \cdot \frac{N+1}{n+1} - 1 \\ &= N. \end{aligned}$$

Von allen Schätzern der Familie ist der Maximumschätzer $\hat{N}_{n:n}$ derjenige mit der kleinsten Varianz. Wegen Abschnitt 3.1, (ii) und (2) gilt nämlich

$$\begin{aligned} V\left(\frac{n+1}{i}X_{i:n} - 1\right) &= \left(\frac{n+1}{i}\right)^2 V(X_{i:n}) \\ &= \frac{(n+1)^2}{i^2} \cdot \frac{i(n+1-i)}{n+1} \cdot \frac{N+1}{n+1} \cdot \frac{N-n}{n+2} \\ &= \left(\frac{n+1}{i} - 1\right) \cdot (N+1) \cdot \frac{N-n}{n+2} \\ &\geq \left(\frac{n+1}{n} - 1\right) \cdot (N+1) \cdot \frac{N-n}{n+2} \\ &= \frac{N+1}{n} \cdot \frac{N-n}{n+2} \\ &= V\left(\frac{n+1}{n}X_{n:n} - 1\right). \end{aligned}$$

Mit Hilfe von tiefgründigen Methoden lässt sich übrigens zeigen, dass $\hat{N}_{n:n}$ unter allen erdenklichen erwartungstreuen Schätzern jener mit kleinster Varianz ist.⁵ Siehe dazu beispielsweise [11], Abschnitt 1.2.5.

⁵Diese Methoden wurden im Zeitraum 1945–1960 von den Statistikern *Calyampudi R. Rao* (1920–), *David H. Blackwell* (1919–), *Erich L. Lehmann* (1917–) und *Henry Scheffé* (1907–1977) entwickelt.

Fallstudie: Anwendung der Statistik für Spionagezwecke⁶

Mindestens einmal in der Geschichte der Statistik wurde ein statistisches Verfahren für Spionagezwecke verwendet; und zwar im 2. Weltkrieg von den Alliierten zur Schätzung der deutschen Waffenproduktion.

Jedes deutsche Kriegsgerät, ob V-2-Rakete, Panzer oder Autoreifen, war während des Produktionsprozesses mit einer Seriennummer versehen worden. War beispielsweise die Gesamtanzahl der bis zu einem bestimmten Zeitpunkt hergestellten Mark-I-Panzer gleich N , so besaß jeder dieser Panzer eine Seriennummer zwischen 1 und N . Nun wurden den Alliierten im Verlauf der Kriegshandlungen einige Nummern

$$1 \leq X_{1:n} < \dots < X_{n:n} \leq N$$

bekannt (entweder dadurch, dass Panzer zerstört oder erbeutet oder dass einschlägige Dokumente erbeutet wurden). Das Verfahren, das ursprünglich zur Schätzung von N angewendet wurde, war der Lückenschätzer

$$\begin{aligned} \tilde{N}_2 &= X_{n:n} + \frac{1}{n-1} \sum_{i=2}^n (X_{i:n} - X_{i-1:n} - 1) = X_{n:n} + \frac{1}{n-1} (X_{n:n} - X_{1:n}) - 1 \\ &= \frac{n}{n-1} X_{n:n} - \frac{1}{n-1} X_{1:n} - 1. \end{aligned}$$

Nach Ende des Krieges, als die Dokumente des deutschen Kriegsministeriums zugänglich wurden, fand man, dass die Schätzwerte für die Waffenproduktion, die auf statistischen Methoden beruhten, weit zuverlässiger waren, als jene, denen andere Informationen zugrunde lagen. So lag beispielsweise der mittels Seriennummern-Schätzer erhaltene Schätzwert 3400 für die bis 1942 erzeugten deutschen Panzer dem tatsächlichen Wert sehr nahe. Der "offizielle" Schätzwert der Alliierten, der auf Informationen beruhte, welche vom Geheimdienst und von Spionageaktivitäten stammten, war hingegen mit 18000 weit überhöht. Fehler dieser Größenordnung, vielfach in der offenbar sehr effektiven "Nazi-Propaganda" begründet, waren nicht ungewöhnlich. Lediglich das sehr objektive Seriennummern-Verfahren war gegenüber derartig verfälschenden Einflüssen unempfindlich!

⁶Eine freie Übersetzung von Case Study 5.4.1 in [4]

2.3 Ausblick: Zur Korrektur einer Verfälschung bei Punktschätzern

In den vorangehenden Ausführungen über Punktschätzer und in [10], Abschnitt 1.4.1 hatten die naheliegenden Kandidaten für Schätzer der entsprechenden Parameter, nämlich

$$X_{n:n} \text{ für } N \quad \text{und} \quad \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \text{ für } \sigma^2$$

die Tendenz, den jeweiligen Parameter zu unterschätzen. In den beiden vorliegenden Fällen liegt also gemäß der nachstehenden schematischen Darstellung eine negative Verfälschung vor:

Schätzung	=	Wert des Parameters	+	Verfälschung	+	Zufallsschwankung
-----------	---	---------------------	---	--------------	---	-------------------

In beiden Fällen war es jedoch möglich, diese Verfälschung durch eine geeignete Modifikation dieser Kandidaten zu beheben. Das Resultat sind die entsprechenden unverfälschten Schätzer⁷

$$\frac{n+1}{n} \times X_{n:n} - 1 \text{ für } N \quad \text{und} \quad S_n^2 = \frac{n}{n-1} \times \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \text{ für } \sigma^2.$$

Motivation: Schießt man mit einem Gewehr auf eine Zielscheibe, so wird die systematische Abweichung durch den Höhenverlust der Kugel infolge (a) der Schwerkraft und (b) des Luftwiderstandes hervorgerufen. Im Folgenden sei in knapper Form auf die Korrektur des durch die Schwerkraft hervorgerufenen Höhenverlusts eingegangen.

Angenommen,
 die Kugel verlässt den Gewehrlauf mit einer Geschwindigkeit von v Metern pro Sekunde,
 die Zielscheibe ist vom Ende des Gewehrlaufes x_0 Meter entfernt und
 das Ende des Gewehrlaufes und die Mitte der Zielscheibe befinden sich auf gleicher Höhe, nämlich auf der Höhe von y_0 Metern.

⁷Hinsichtlich der Begründung, warum beim Schätzer S_n^2 für σ^2 nicht durch n , sondern durch $n-1$ dividiert wird, siehe auch [8].

2.3. AUSBLICK: ZUR KORREKTUR EINER VERFÄLSCHUNG BEI PUNKTSCHÄTZERN¹⁹

Ferner sei $g \cong 9.81 \text{ m/sec}^2$ die Erdbeschleunigung.

Bezeichnen nun

α den, in Bogenlänge angegebenen, *Anstell-* oder *Abgangswinkel*⁸ des Gewehrs,

x die horizontale Entfernung der Kugel vom Ende des Gewehrlaufes in Richtung Zielscheibe und

$y_\alpha(x)$ ihre Höhe als Funktion von x ,

so ist letztere aufgrund von Galileis Fallgesetz gleich

$$y_\alpha(x) = y_0 + \frac{1}{2 \cos^2(\alpha)} x (\sin(2\alpha) - \frac{g}{v^2} \cdot x) .$$

Ist der *Anstellwinkel* gleich $\alpha = 0$, so ergibt sich demgemäß

$$y_0(x) = y_0 - \frac{g}{2v^2} \cdot x^2$$

und somit der Höhenverlust $\frac{g}{2v^2} \cdot x^2$. Dieser lässt sich dadurch vermeiden, dass man den *Anstellwinkel* α so wählt, dass $y_\alpha(x_0) = y_0$ gilt. Dies wird offensichtlich dadurch erreicht, dass man den letzten Term in der Formel für $y_\alpha(x)$ gleich Null setzt. Dementsprechend hat man den Anstellwinkel gleich

$$\alpha = \frac{1}{2} \arcsin\left(\frac{g}{v^2} \cdot x_0\right)$$

zu wählen. Dabei wird mit $\arcsin(\cdot)$ der Hauptwert des Arcus-Sinus be-

⁸Dies ist der von der Visierlinie und der sogenannten *Laufseelenachse* eingeschlossene Winkel.

Zur weiteren Begriffsklärung:

Die *Visierlinie* ist die durch *Kimme*, *Korn* und Zielobjekt bestimmte Gerade. Sie ist im vorliegenden Fall waagrecht.

Die *Laufseelenachse* ist die durch den Gewehrlauf bestimmte Gerade.

Die *Kimme* ist eine V-förmige Kerbe am Visier des Gewehrs.

Das *Korn* ist der zugehörige komplementäre Teil des Visiers. Es hat demnach die Form \wedge .

zeichnet.

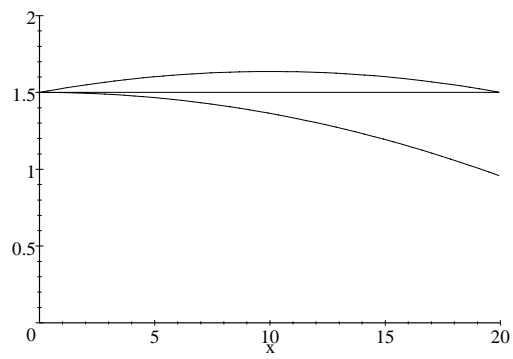


Abbildung der *Visierlinie* und der Bahnkurven der Kugel bei *Anstellwinkel* $\alpha = 0$ und bei adjustiertem *Anstellwinkel*

Kapitel 3

Anhang: Stochastische Grundlagen

3.1 Rechenregeln für Erwartungswert und Varianz

Es sei X eine Zufallsvariable mit endlichem Wertebereich $W_X = \{\omega_1, \dots, \omega_m\} \subset \mathbb{R}$ und Wahrscheinlichkeitsverteilung $P(X = \omega_j)$, $j \in \{1, \dots, m\}$. Dann heißen die Größe

$$E(X) = \sum_{j=1}^m \omega_j \cdot P(X = \omega_j)$$

der *Erwartungswert* und

$$V(X) = E[(X - E(X))^2] = \sum_{j=1}^m (\omega_j - E(X))^2 \cdot P(X = \omega_j)$$

die *Varianz* von X .

Rechenregeln: Seien $a, b \in \mathbb{R}$, $a \neq 0$ und $Y = aX + b$ die Zufallsvariable mit Wertebereich $W_Y = aW_X + b$ und Verteilung $P(Y = a\omega_j + b) = P(X = \omega_j)$, $j \in \{1, \dots, m\}$. Dann gelten

$$E(Y) = aE(X) + b \tag{i}$$

und

$$V(Y) = a^2 V(X). \tag{ii}$$

Für die Berechnung der Varianz von X sind folgende äquivalente Darstellungen sehr nützlich

$$\begin{aligned} V(X) &= E(X^2) - (E(X))^2 \\ &= E[X(X - E(X))] - E(X)E(X - 1) \\ &= E[X(X + E(X))] - E(X)E(X + 1). \end{aligned} \quad (\text{iii})$$

3.2 Die Binomialverteilung

Binomialverteilung: Seien $n \in \mathbb{N}$ und $p \in (0, 1)$. Dann heißt die durch

$$P(S_n = k) = \binom{n}{k} p^k (1-p)^{n-k}, \quad k \in \{0, \dots, n\}$$

gegebene Verteilung *Binomialverteilung* mit den Parametern n und p .

Bezeichnung: Die Zufallsvariable S_n mit dem Wertebereich $\{0, \dots, n\}$ und der obigen Verteilung heißt Binomialverteilt mit den Parametern n und p (kurz: $S_n \sim B_{n,p}$).

Erwartungswert und Varianz: Es gelten

$$E(S_n) = np \quad \text{und} \quad V(S_n) = np(1-p).$$

Herleitung: Entscheidend ist die Rekursionsformel

$$k \cdot P(S_n = k) = np \cdot P(S_{n-1} = k-1), \quad k \in \{1, \dots, n\},$$

wobei $S_{n-1} \sim B_{n-1,p}$, welche auf dem Sachverhalt

$$k \cdot \binom{n}{k} = n \cdot \binom{n-1}{k-1}, \quad k \in \{1, \dots, n\}$$

beruht. Für den Erwartungswert gilt demnach

$$\begin{aligned} E(S_n) &= \sum_{k=0}^n k \cdot P(S_n = k) \\ &= \sum_{k=1}^n k \cdot P(S_n = k) = np \sum_{k=1}^n P(S_{n-1} = k-1) = np, \end{aligned}$$

denn die Summe

$$\sum_{k=1}^n P(S_{n-1} = k-1) = \sum_{k-1=0}^{n-1} P(S_{n-1} = k-1)$$

ist - als Summe aller Wahrscheinlichkeiten der $B_{n-1,p}$ -verteilten Zufallsvariablen S_{n-1} - gleich 1.

Nun zur Varianz:

$$\begin{aligned} E[(S_n - 1)S_n] &= \sum_{k=2}^n (k-1)k \cdot P(S_n = k) \\ &= np \sum_{k-1=1}^{n-1} (k-1)P(S_{n-1} = k-1) \\ &= E(S_n) E(S_{n-1}) . \end{aligned}$$

Die Varianz ist daher im Hinblick auf die mittlere Darstellung in Abschnitt 3.1, (iii)

$$\begin{aligned} V(S_n) &= E[(S_n - 1)S_n] - E(S_n) E(S_n - 1) \\ &= E(S_n) E(S_{n-1}) - E(S_n) E(S_n - 1) \\ &= E(S_n) (1 - E(S_n) + E(S_{n-1})) \\ &= np(1 - (1 + n - 1)p + (n - 1)p) \\ &= np(1 - p) . \quad \square \end{aligned}$$

3.3 Die Inverse Hypergeometrische Verteilung

Verteilung der i -ten Ordnungsstatistik beim Ziehen ohne Zurücklegen: Seien $N \in \mathbb{N}$, $n \in \{1, \dots, N\}$ und $i \in \{1, \dots, n\}$. Dann ist die i -te Ordnungsstatistik gemäß der durch

$$P(X_{i:n} = k) = \frac{\binom{k-1}{i-1} \binom{N-k}{n-i}}{\binom{N}{n}}, \quad k \in \{i, \dots, i + N - n\},$$

gegebenen Verteilung verteilt. Diese heißt *Inverse Hypergeometrische Verteilung*¹ mit den Parametern i , N und n .

¹Eine andere Bezeichnung ist *Hypergeometrische Wartezeitverteilung*.

Herleitung: Der kleinstmögliche Wert von $X_{i:n}$ ist i , wobei $X_{i:n} = i$ genau dann zutrifft, wenn sich unter den Nummern der n gezogenen Jetons die i kleinsten Nummern befinden.

Der größtmögliche Wert von $X_{i:n}$ ist $i + N - n$, wobei $X_{i:n} = i + N - n = N - (n - i)$ genau dann zutrifft, wenn sich unter den Nummern der n gezogenen Jetons die $n - i$ größten Nummern befinden.

Sei nun $k \in \{i, \dots, i + N - n\}$. Die Wahrscheinlichkeit jedes Ereignisses $E_k = \{X_{i:n} = k\}$ wird vermittle des Laplace-Prinzips bestimmt:

$$P(E_k) = \frac{\text{Anzahl der für das Ereignis } E_k \text{ günstigen Fälle}}{\text{Anzahl aller möglichen Fälle}},$$

wobei wir uns der Bequemlichkeit halber vorstellen, dass wir der Urne alle n Jetons mit einem Griff entnehmen.²

Die Anzahl aller möglichen Fälle ist gleich der Anzahl der n -elementigen Teilmengen einer N -elementigen Menge und somit $\binom{N}{n}$.

Das Ereignis $\{X_{i:n} = k\}$ ist durch folgende drei Eigenschaften bestimmt:

- (1) von den $k - 1$ Nummern $1, \dots, k - 1$ werden genau $i - 1$ gezogen,
- (2) eines der Jetons hat die Nummer k ,
- (3) von den $N - k$ Nummern $k + 1, \dots, N$ werden genau $n - i$ gezogen,

wobei jede der möglichen Anordnungen für (1) mit jeder der möglichen Anordnungen für (3) zu kombinieren ist.

Die Anzahl der für das Ereignis $\{X_{i:n} = k\}$ günstigen Fälle ist daher

$$|\{X_{i:n} = k\}| = \binom{k-1}{i-1} \cdot \binom{1}{1} \cdot \binom{N-k}{n-i} = \binom{k-1}{i-1} \binom{N-k}{n-i}$$

und die gesuchte Wahrscheinlichkeit somit

$$P(X_{i:n} = k) = \frac{\binom{k-1}{i-1} \binom{N-k}{n-i}}{\binom{N}{n}}. \quad \square$$

Bezeichnung: Die Zufallsvariable $X_{i:n}$ mit dem Wertebereich $\{i, \dots, i + N - n\}$ und der obigen Verteilung heißt *Invers Hypergeometrischverteilt* mit den Parametern i , N und n (kurz: $X_{i:n} \sim IH_{i,N,n}$).

²Wie bei der Hypergeometrischen Verteilung stellt sich auch hier heraus, dass die Wahrscheinlichkeit nicht davon abhängt, ob man der Urne die n Jetons mit einem Griff oder nacheinander entnimmt.

Erwartungswert und Varianz: Es gelten

$$\begin{aligned} E(X_{i:n}) &= i \cdot \frac{N+1}{n+1} \\ V(X_{i:n}) &= \frac{i(n+1-i)}{n+1} \cdot \frac{N+1}{n+1} \cdot \frac{N-n}{n+2}. \end{aligned}$$

Herleitung: Entscheidend ist die Rekursionsformel

$$k \cdot P(X_{i:n} = k) = i \frac{N+1}{n+1} \cdot P(X_{i+1:n+1} = k+1), \quad k \in \{i, \dots, i+N-n\},$$

wobei $X_{i+1:n+1} \sim IH_{i+1, N+1, i+1}$. Für den Erwartungswert gilt demnach

$$E(X_{i:n}) = \sum_{k=i}^{i+N-n} k \cdot P(X_{i:n} = k) = i \frac{N+1}{n+1} \sum_{k=i}^{i+N-n} P(X_{i+1:n+1} = k+1) = i \frac{N+1}{n+1},$$

denn die Summe

$$\sum_{k=i}^{i+N-n} P(X_{i+1:n+1} = k+1) = \sum_{k+1=i+1}^{i+1+(N+1)-(n+1)} P(X_{i+1:n+1} = k+1)$$

ist - als Summe aller Wahrscheinlichkeiten der $IH_{i+1, N+1, i+1}$ -verteilten Zufallsvariablen $X_{i+1:n+1}$ - gleich 1.

Nun zur Varianz: Es ist

$$\begin{aligned} E[(X_{i:n} + 1)X_{i:n}] &= \sum_{k=i}^{i+N-n} (k+1)k \cdot P(X_{i:n} = k) \\ &= k \frac{N+1}{s+1} \sum_{k+1=i+1}^{i+1+(N+1)-(n+1)} (k+1)P(X_{i:n} = k+1) \\ &= E(X_{i:n}) E(X_{i+1:n+1}). \end{aligned}$$

Wegen

$$\begin{aligned} E(X_{i+1:n+1}) - E(X_{i:n} + 1) &= (i+1) \frac{N+2}{n+2} - i \frac{N+1}{n+1} - 1 \\ &= \frac{N+2}{n+2} - 1 - i \left(\frac{N+1}{n+1} - 1 - \left(\frac{N+2}{n+2} - 1 \right) \right) \\ &= \frac{N-n}{n+2} - i \left(\frac{N-n}{n+1} - \frac{N-n}{n+2} \right) \\ &= \frac{N-n}{n+2} \left(1 - \frac{i}{n+1} \right) \end{aligned}$$

ist die Varianz daher im Hinblick auf die letzte Darstellung in Abschnitt 3.1, (iii) gleich

$$\begin{aligned}
 V(X_{i:n}) &= E[(X_{i:n} + 1)X_{i:n}] - E(X_{i:n})E(X_{i:n} + 1) \\
 &= E(X_{i:n})E(X_{i+1:n+1}) - E(X_{i:n})E(X_{i:n} + 1) \\
 &= E(X_{i:n})[E(X_{i+1:n+1}) - E(X_{i:n} + 1)] \\
 &= i \frac{N+1}{n+1} \cdot \frac{N-n}{n+2} \left(1 - \frac{i}{n+1}\right) \\
 &= i \left(1 - \frac{i}{n+1}\right) \frac{N+1}{n+1} \frac{N-n}{n+2}. \quad \square
 \end{aligned}$$

Literaturverzeichnis

[1] **Bücher**

- [2] *Freedman, D., Pisani, R. and R. Purves*: Statistics. Norton & Co., New York 2007
- [3] *Krämer, W.*: Statistik verstehen: Eine Gebrauchsanweisung. Campus Verlag, Frankfurt - New York, 1999
- [4] *Larsen, R.J. and M.L. Marx*: An Introduction to Mathematical Statistics and its Applications. Pearson Prentice Hall, Upper Saddle River, New Jersey 2006
- [5] *Moore, D.S.*: Statistics: Concepts and Controversies. W.H. Freeman & Co., New York 2001
- [6] *Moore, D.S. and G.P. McCabe*: Introduction to the Practice of Statistics. W.H. Freeman & Co., New York 2004

Zeitschriften

- [7] *Engel, A.*: Statistik in der Schule: Ideen und Beispiele aus neuerer Zeit. Der Mathematikunterricht, Jahrgang 28, Heft 1 (1982), S. 57 – 85
- [8] *Diepgen, R.*: Warum nur $n - 1$ und nicht n ? Erwartungstreue - leicht gemacht. Stochastik in der Schule **19** (1999), Nr. 1, S. 10 – 13

Skripten

- [9] *Österreicher, F.*: Ausgewählte Kapitel der Statistik. LV-Unterlagen, Salzburg 1986

- [10] *Österreicher, F.*: Skriptum zur Lehrveranstaltung Statistik für Lehramt. Salzburg 2007 *)
- [11] *Österreicher, F.*: Skriptum zur Lehrveranstaltung Mathematische Statistik, Salzburg 2008 *)

Seminarunterlagen

- [12] *Österreicher, F.*: Schätzen des Umfangs von Populationen. Fortbildungsseminar, Meran 1990
- [13] *Österreicher, F.* und *M. Weiß*: Unterlagen zum Sochastikseminar: Teil 1, Salzburg 2007 *)
- [14] *Österreicher, F.* und *M. Weiß*: Unterlagen zum Sochastikseminar: Teil 2, Salzburg 2007 *)

*) verfügbar unter "<http://www.uni-salzburg.at>" > Fakultäten und Fachbereiche > Naturwissenschaftliche Fakultät > Fachbereich Mathematik > Personen > Dozenten >

Diplomarbeiten

- [15] *Weiß, M.*: Binomialverteilung und Normalapproximation: Grundlegendes und Hintergrundinformation für den Stochastikunterricht. Salzburg 1995
- [16] *Kolmberger, M.*: Statistik in der Nußschale - Ist unser Würfel fair? Diplomarbeit, Salzburg 1997
- [17] *Jandl, M.*: Computereinsatz im Stochastikunterricht. Salzburg 1997

Schulbücher für die AHS-Oberstufe

- [18] *Malle·Ramharter·Ulovec·Kandl*: Mathematik verstehen 6, öbvhtp Verlagsgesellschaft, Wien 2005
- [19] *Malle·Ramharter·Ulovec·Kandl*: Mathematik verstehen 7, öbvhtp Verlagsgesellschaft, Wien 2006
- [20] *Malle·Ramharter·Ulovec·Kandl*: Mathematik verstehen 8, öbvhtp Verlagsgesellschaft, Wien 2007

- [21] *Götz·Reichel·R.Müller·Hanisch·Hederer·Wenzel·M.Müller*: Mathematik Lehrbuch 6, öbvhtp Verlagsgesellschaft, Wien 2005
- [22] *Götz·Reichel·R.Müller·Hanisch·Hederer·Wenzel·M.Müller*: Mathematik Lehrbuch 7, öbvhtp Verlagsgesellschaft, Wien 2006
- [23] *Götz·Reichel·R.Müller·Hanisch·Hederer·Wenzel·M.Müller*: Mathematik Lehrbuch 8, öbvhtp Verlagsgesellschaft, Wien 2007
- [24] *Taschner*: Mathematik 2 - Übungs- und Lehrbuch für die 6. Klasse AHS, Oldenbourg, Wien 1999
- [25] *Taschner*: Mathematik 3 - Übungs- und Lehrbuch für die 7. Klasse AHS, Oldenbourg, Wien 2000
- [26] *Taschner*: Mathematik 4 - Übungs- und Lehrbuch für die 8. Klasse AHS, Oldenbourg, Wien 2001
- [27] *Geretschläger·Griesel·Postel*: Elemente der Mathematik 6, E. Dorner, Wien 2005
- [28] *Geretschläger·Griesel·Postel*: Elemente der Mathematik 7, E. Dorner, Wien 2006
- [29] *Geretschläger·Griesel·Postel*: Elemente der Mathematik 8, E. Dorner, Wien 2007
- [30] *Steiner·Novak*: MatheMaster 6 - Mathematik für die 6. Klasse AHS, Reniets Verlag, Wien 2005
- [31] *Steiner·Novak*: MatheMaster 7 - Mathematik für die 7. Klasse AHS, Reniets Verlag, Wien 2006
- [32] *Steiner·Novak*: MatheMaster 8 - Mathematik für die 8. Klasse AHS, Reniets Verlag, Wien 2007

Schulbücher für Höhere Lehranstalten für Wirtschaftliche Berufe

- [33] *Hanisch·Reichel·Müller·Schak*: Mathematik für HLA 4. öbvhtp, Wien 2006

Schulbücher für Handelsakademien

- [34] *Kronfellner, M. Peschek · Blasonig · Fischer · Kronfellner, J.*: Angewandte Mathematik 4. Verlag Hölder Pichler Tempsky, Wien 2001
- [35] *Schneider · Thannhauser*: Mathematik Arbeitsbuch und Aufgabensammlung einschließlich Lösungen. Band 4 für den V. Jahrgang HAK. Rudolf Trauner Verlag, Linz 1999
- [36] *Steiner · Weilharter*: Mathematik und ihre Anwendungen in der Wirtschaft. Band 4. Reniets Verlag, Wien 2008
- [37] *Hinkelmann · Böhm · Hofbauer · Metzger · Schuhäker*: Mathe mit Gewinn 1. öbvhtp, Wien 2005

Schulbücher für Höhere Technische Lehranstalten

- [38] *Schärf*: Mathematik 2 für HTL. Oldenbourg Verlag, Wien 1998
- [39] *Schärf*: Mathematik 3 für HTL. Oldenbourg Verlag, Wien 1999
- [40] *Schalk · Steiner*: Mathematik 4. Reniets Verlag, Wien 2001
- [41] *Timischl · Kaiser*: Ingenieur-Mathematik 4, Verlag E. Dorner, Wien 2005

Schulbücher aus Deutschland

- [42] *Barth · Haller*: Stochastik Leistungskurs. Ehrenwirth Verlag, München 1983
- [43] *Heigl · Feuerpfeil*: Stochastik Leistungskurs. Bayrischer Schulbuch Verlag, München 1987
- [44] *Lambacher · Schweizer*: Stochastik Leistungskurs. Ernst Klett Schulbuchverlag, Stuttgart 1988

