

Daten, die auf der Erde liegen – auf Spurensuche im Supermarkt

ANDREAS EICHLER, MÜNSTER UND WOLFGANG RIEMER, KÖLN

Zusammenfassung: Die Analyse von Kassensbons, die in Supermärkten achtlos weggeworfen werden, ermöglicht eine detektivische Datenanalyse, die mit unterschiedlichen Schwerpunkten in allen Schulstufen möglich ist. In diesem Artikel wird Schritt für Schritt die Rekonstruktion des Mikrokosmos eines Supermarkts aus zunächst leblosen Zahlenkolonnen nachgezeichnet – vom Konsumentenverhalten bis hin zum Umsatz des Supermarkts. Ein 'historisches' Nachwort zur Verwendung von Kassenzetteln im Stochastikunterricht zur Zeit der Euro-Umstellung beschließt den Artikel.

1 Einführung

Sie liegen buchstäblich auf der Erde herum, zusammengeknüllt, weggeworfen oder im Einkaufswagen zurückgelassen: Nach einem prüfenden Blick haben die Bons, die man an der Kasse eines Supermarkts erhalten hat, offenbar keinerlei Wert mehr für deren kurzfristige Besitzer. Doch für angehende Statistiker aller Schulstufen können sie vielfältige Ansätze für die Analyse realer Daten geben. Durch unterschiedliche Untersuchungsperspektiven, das Aufstellen unterschiedlicher Modelle sowie die flexible Verwendung statistischer Methoden entsteht so Stück für Stück aus zunächst leblosen Zahlenkolonnen der Mikrokosmos eines Supermarkts.

Die Idee, Kassensbons auszuwerten, stammt vom zweiten Autor (2003, 2006) und wurde dort als Projekt vorgeschlagen. Einige der vielen Möglichkeiten, diese von der Erde gesammelten Daten auszuwerten, sollen in dieser Arbeit behandelt werden. Dabei geht zunächst der Blick in die unteren Klassen (auch in die Primarstufe), in der Fragen nach dem Konsumentenverhalten anhand eindimensionaler Datensätze thematisiert werden können. Der Bogen reicht dann bis zum Untersuchen von Zusammenhängen durch die Auswertung zweidimensionaler Datensätze, mit denen die Finanzlage eines Supermarktes aufgedeckt und fachlich die Datenanalyse mit der Wahrscheinlichkeitsrechnung verbunden wird. Die Arbeit ist insgesamt als Pool von Anregungen zu verstehen. Die Beispiele werden exemplarisch ausgeführt, ohne diese allerdings als fertigen Unterrichtsvorschlag für spezielle Klassenstufen auszubauen.

2 Datensammlung und erste Einsichten

2.1 Der Supermarkt

Der Supermarkt, um den es hier geht, liegt in Braunschweig und dort in einer gutbürgerlichen Gegend. Es ist ein Markt, der Kunden aus der umliegenden Gegend bedient und von Montag bis Samstag in der Zeit von 8 bis 20 Uhr geöffnet ist. Viele Kunden kommen zu Fuß oder mit dem Rad, der Parkplatz ist relativ klein. Der Supermarkt unterscheidet sich damit deutlich von den großen Märkten in der Peripherie der Stadt mit großen Parkplätzen, geschaffen für einen großen Wocheneinkauf und Angeboten, die weit über Nahrungsmittel hinausreichen.

So ist das Angebot des Supermarkts zwar eher auf gehobene Bedürfnisse ausgelegt, bedient aber neben den üblichen Haushaltswaren wie Reinigungsmittel oder Pflegebedarf in erster Linie Nahrungsmittel und Getränke. Neben Angestellten, die ausschließlich für die Bestückung der Obst- und Gemüseabteilung oder für die Fleisch- und Käsetheke zuständig sind, bedienen 14 Verkäuferinnen bzw. Verkäufer¹ die vier Kassen des Supermarkts.

2.2 Die Datensammlung

Ein Kassensbon ist ein Blatt voller Daten (siehe Abb. 1). Einige der darauf enthaltenen Daten sind unmittelbar einsichtig, das sind die (meisten) Artikel-Bezeichnungen und zugehörigen Preise. Die Angaben zur Preis-Summe aller Waren, zum Hin- und Rückgeld sowie zur Mehrwertsteuerauszeichnung sind ebenso einsichtig. Dann folgen allerdings im unteren Bereich (siehe Pfeil) des Bons Angaben, die erst nach dem Betrachten mehrerer Kassensbons decodiert werden können.² So bezeichnet

- die erste vierstellige Nummer (Zahl) die Nummer des Kunden an einer bestimmten Kasse,
- die zweite vierstellige Nummer den speziellen Supermarkt,
- die erste dreistellige Nummer die Kasse und
- die zweite dreistellige Nummer die Verkäuferin.

Das Datum und die Uhrzeit des Einkaufs sind schließlich ebenfalls noch auf dem Kassensbon enthalten. Was hier profan erzählt wird, ist schon Teil des Detektivspiels „dem Supermarkt auf der Spur“.

Die Zahlen bzw. Ziffern, die hier teilweise den im Mathematikunterricht der Sekundarstufen eher vereinzelt behandelten Kodieraspekt berühren (vgl. z.B. Herget 1994), enthalten eine zunächst unbekannte Botschaft, die entschlüsselt werden muss. Dieses für jeden Supermarkt unterschiedliche Rätsel zu lösen, ist überhaupt der Ansporn gewesen, mehr als einen Bon zu betrachten, und hat schließlich die Jagd auf das Wegwerfprodukt Kassenzettel ausgelöst.



Abb. 1: Ein Kassenzettel des Supermarkts

Seit dem 6. November 2006 sind insgesamt bisher 372 Bons unterschiedlicher Tage und unterschiedlicher Tageszeiten zusammengekommen.

2.3 Die Güte der Daten

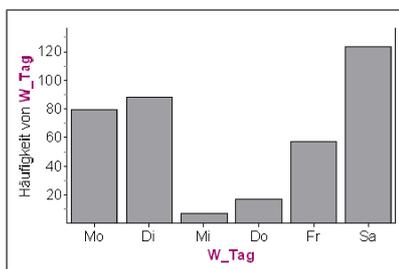


Abb. 2: Verteilung der Kassenzettel über die Wochentage

Da die Kassenzettel von einer Einzelperson bei diversen eigenen Einkäufen in dem Supermarkt gesammelt wurden, kann natürlich nur bedingt von einer Zufallsauswahl oder repräsentativen Auswahl gesprochen werden. Das liegt zum einen an der Zeit der

eigenen Einkäufe und der zu diesen Zeiten verfügbaren Bons. Abbildung 2 zeigt die Verteilung der gesammelten Bons über die Wochentage.

Auch die Verteilung der Kassenzettel über die Öffnungszeiten des Marktes ist nicht repräsentativ. Insbesondere die Randzeiten sind deutlich unterrepräsentiert (siehe Abb. 3).

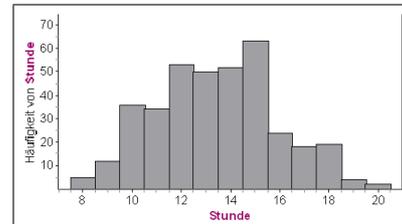


Abb. 3: Verteilung der Kassenzettel über die Stunden

Ein anderes Manko der Datengüte ist die Frage, wer überhaupt seinen Kassenzettel annimmt (wird er gleich an der Kasse gelassen, landet er nämlich in einem unerreichbaren Korb), wer ihn annimmt und kurze Zeit später achtlos zurücklässt und wer diesen mitnimmt. Rein qualitativ, aus den dreimonatigen Beobachtungen herrührend, ist kein Muster erkennbar gewesen, ob zum Beispiel erst bei besonders großen Einkäufen der Bon zur Kontrolle mitgenommen wird. Eine ebenso qualitative Einschränkung gibt es dabei: Kartenzahler scheinen den Bon zusammen mit dem Einzugsbeleg für die Karte eher mitzunehmen als fallenzulassen. Ein Indiz dafür könnte das geringe Aufkommen an solchen Bons zu sein, auf denen die Kartenzahlung erkennbar dokumentiert ist. Andererseits gibt es keinen harten Beleg dafür, wie viele Kunden mit Karte bezahlen, sondern eher einen, der in der Literatur unter dem Stichwort „Verfügbarkeitsheuristik“ (Tversky und Kahnemann 1973) bekannt geworden ist. Diese Verfügbarkeitsheuristik bedeutet in diesem Zusammenhang die subjektiv empfundene Häufigkeit (bzw. Wahrscheinlichkeit) von Kunden, die mit Karte bezahlen, denn schließlich ist man überzeugt davon, dass zumindest in der eigenen Schlange vor der Kasse prinzipiell alle mit einer (langsamen) Karte bezahlen, deren halb defekter Magnetstreifen mehrere Anläufe verursacht.

3 Wer kauft was – eine einfache Datenanalyse

Eine einfache, mit kleinen Datensätzen bereits in der Primarstufe anwendbare Analyse betrifft das Konsumentenverhalten. Was wird eigentlich vorwiegend in solch einem Supermarkt gekauft? Ein Teildatensatz

von 64 Kassenbons ergibt dabei folgendes Bild (siehe Tabelle 1):

| Warenart | Umsatz (in Euro) | Anteil am Gesamtumsatz (in %) |
|--------------------|---------------------|-------------------------------|
| Fleisch | 126,48 | 13,9 |
| Obst & Gemüse | 123,21 | 13,5 |
| Sonstiges | 122,91 | 13,5 |
| Süßigkeiten | 77,77 | 8,5 |
| Käse | 68,00 | 7,5 |
| Alkohol | 63,23 | 7,0 |
| Backwaren | 37,85 | 4,1 |
| Milch | 35,40 | 3,9 |
| Fertiggerichte | 31,55 | 3,5 |
| Getränke | 30,34 | 3,3 |
| Zigaretten | 28,00 | 3,1 |
| Backzutaten | 25,73 | 2,8 |
| Eier | 25,68 | 2,8 |
| Joghurt | 24,71 | 2,7 |
| Konserven (Gemüse) | 14,59 | 1,6 |
| Gewürze | 14,31 | 1,6 |
| Sanitärbedarf | 13,29 | 1,5 |
| Zeitschriften | 12,60 | 1,4 |
| Nudeln & Reis | 9,83 | 1,1 |
| Kaffee & Tee | 9,68 | 1,1 |
| Knabberwaren | 8,38 | 1,0 |
| Körperpflege | 5,58 | 0,6 |
| Konserven (Obst) | 3,37 | 0,4 |

Tabelle 1: Verteilung der Umsätze auf einzelne Produktklassen

Eine Datenaufbereitung in diesem Sinne berührt das elementare Umgehen mit dem Größenbereich Geld. Aus der Sicht der Datenanalyse ist hier der Aspekt der Klassierung interessant, der eine fundamentale Aufgabe jeglicher (Ein)Ordnung von Daten ist. So hat ein Supermarkt unüberschaubar viele Waren im Angebot. Unter welchen Kategorien oder eben Klassen soll man diese erfassen? Hier muss man vorher überlegen und kommt möglicherweise zu der Einsicht, die die Statistik insgesamt durchzieht: Es gibt eigentlich keine „richtigen“ oder „falschen“ Klasseneinteilungen, sondern nur sinnvolle und weniger sinnvolle. Sinnvoll bedeutet hier allein, dass man die Fragen, die man an die Daten stellt, durch eine mehr oder minder willkürliche Klasseneinteilung mehr oder weniger erschöpfend beantworten kann. Auch die oben gegebene Klasseneinteilung ist willkürlich, eine Vielzahl anderer Einteilungen ist möglich. Je nach Fragestellungen könnte man andere, auch erheblich einfachere Klasseneinteilungen vornehmen, wie etwa *Getränke*, *Nahrungsmittel*, *Anderes* oder *Frischprodukte*, *konservierte Produkte* etc.

Eine Deutung oder Interpretation der Daten, ein

grundsätzlicher Schritt des statistischen Denkens (Pfannkuch und Wild 1999), ist bei diesem Datensatz schwierig. Aus der Auswahl der Kassenbons ergibt sich etwa ein vermutlich zu hoher Anteil an Backzutaten und Süßigkeiten, da einige Kassenbons aus der Vorweihnachtszeit stammen und einen erheblichen Anteil an den jeweiligen Gesamtsummen haben. Ebenso wird man an den Tagen vor dem Wochenende einen erhöhten Umsatz an Alkoholika entdecken.

Sicher sind weitere Interpretationen anhand der Datentabelle möglich. Interessant und propädeutisch für die Statistik wichtig wird es aber dann, wenn man *vor* der Datenaufbereitung (aber nach einer Klasseneinteilung) eine wiederum fundamentale statistische Methode angewendet hat: nämlich die Hypothesenbildung bezüglich der Verteilung der Waren innerhalb des Gesamtbudgets. Dann hat man möglicherweise Argumentationsbedarf: Stimmen die Daten mit den Hypothesen überein? Ab welcher Abweichung geht man davon aus, dass die Hypothese falsch war? Gut wäre es dabei, einen weiteren Datensatz zur Verfügung zu haben, der ein erstes Erkennen der Variabilität statistischer Daten (Pfannkuch und Wild 1999) ermöglicht. An solch einem Beispiel lässt sich die Hypothesenbildung sowie die Beibehaltung bzw. Ablehnung von Hypothesen vorbereiten.

4 Dem Einzelumsatz auf der Spur

Spannend wird es, wenn man die Daten im Sinne des Umsatzes des Supermarkts hin- und herwendet. Auch das geht bezogen auf untere Klassenstufen (quasi) eindimensional.

Alle folgenden Bearbeitungen sind sinnvoll nur mit dem Rechner zu bearbeiten. Eine sehr geeignete Software, die unmittelbar die Datenanalyse unterstützt, ist *fathom*, das seit einiger Zeit in deutscher Sprache auf dem Markt ist. Der ungeheure Vorteil ist die unmittelbare, sehr einfach handhabbare Verfügbarkeit von grafischen Datenaufbereitungen in allen schulrelevanten Repräsentationen. Der Nachteil ist, etwa gegenüber Excel, die (noch) geringe Verbreitung. Alle Grafiken in dieser Arbeit sind mit *fathom* erstellt.

4.1 Der Umsatz

4.1.1 Der Umsatz pro Kunden

Nimmt man alle Kassenbons ($n = 372$), dann ergibt sich eine leicht linkssteile Verteilung der Einzelumsätze (vgl. Abb. 4). Der Begriff „Umsatz“ in der

Grafik und ebenso im Folgenden bezeichnet stets den Einzelumsatz pro Kunden.

Man erkennt die oben bereits gegebene Charakterisierung des Supermarkts: „Großeinkäufe“ mit einem Umsatz über 50 Euro sind selten (zumindest in der Stichprobe), das Maximum liegt bei rund 72 Euro. Der Modalwert liegt bei der gewählten Klasseneinteilung bei 10 Euro, die Hälfte der Kunden tätigt Einkäufe zwischen etwa 10 und 20 Euro.

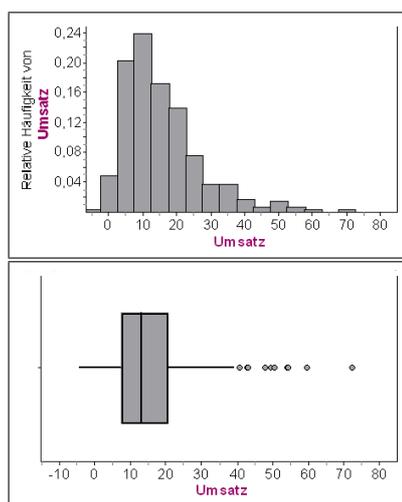


Abb. 4: Verteilung der Einzelumsätze

Unter der Annahme, dass es sich bei der Stichprobe um eine zumindest annähernd repräsentative Stichprobe handelt (siehe oben), erhält man über das arithmetische Mittel ($\bar{x} \approx 15,3$) eine Schätzung hinsichtlich des Umsatzes pro Kunden bei einer Standardabweichung von etwa 11 (Euro). Diese Schätzung wird im weiteren Verlauf der Analyse wichtig.

Mit Hilfe der (sinnvollen) Annahme, dass die Mittelwerte der Verteilung des Umsatzes pro Kunden in Stichproben von je 372 Kassenbons annähernd normalverteilt sind, ergibt sich nach dem Standardverfahren zur Berechnung von Konfidenzintervallen ($\bar{x} \pm \frac{s}{\sqrt{n}}$, mit $\bar{x} \approx 15$ und $s \approx 11$ aus der Stichprobe mit $n = 372$ Kassenbons) das Intervall $I_\mu = [14,2; 16,4]$ (Euro) für den tatsächlichen Umsatz (μ) pro Kunden (Konfidenzniveau: $\alpha = 0,95$).

4.1.2 Der Umsatz zu verschiedenen Zeitpunkten

Ein weiterer Aspekt bei der Untersuchung der Kassenbons könnte von der Frage ausgehen, ob sich beispielsweise an verschiedenen Wochentagen oder auch zu verschiedenen Tageszeiten Unterschiede in den Umsätzen ergeben. Diese Aspekte werden im Folgenden zunächst qualitativ, über den grafikgesteuerten Vergleich der Verteilung sowie den numerisch

gesteuerten Vergleich von Kennzahlen der Verteilung betrachtet. Anschließend folgt mit der Beurteilung mit Hilfe des Anpassungstest ein „hartes“ statistisches Verfahren.

Bezüglich der Tageszeiten (vgl. Abb. 5) ergibt sich ein ähnliches Bild der Umsatzverteilung pro Kunden, wie es sich bereits in der Verteilung bezüglich aller Kassenbons gezeigt hat. Die Aufteilung in Tageszeiten bedeutet hier eine Drittelung der Öffnungszeiten (Morgen: 8-12 Uhr, Mittag: 12-16 Uhr, Abend: 16-20 Uhr). Methodisch betreibt man hier eine Clustering der Daten und damit einen ersten, qualitativen Schritt in Richtung der Beurteilung von Abhängigkeiten zwischen zwei Merkmalen. Die Verfahren bleiben allerdings weiterhin die gleichen, die bei der Analyse eindimensionaler Daten verwendet werden.

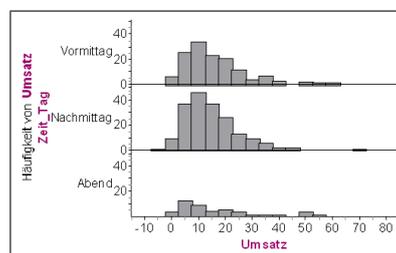


Abb. 5: Umsätze pro Kunden nach Tageszeiten

Ebenso ergibt sich das gleiche Bild, wenn man als Clustering der Gesamtdaten die Wochentage nimmt (vgl. Abb. 6). Im Histogramm sind hier die Tage Mittwoch und Donnerstag auf Grund ihres geringen Vorkommens im Datensatz ausgenommen worden.

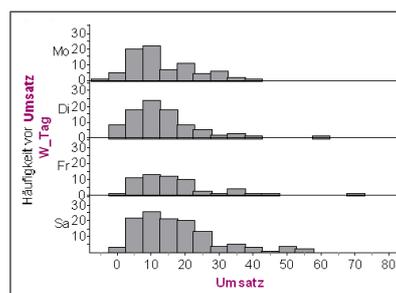


Abb. 6: Umsätze pro Kunden nach Tageszeiten

Im Sinne des Aufbaus des statistischen Denkens sind in diesem Schritt der Datenanalyse mehrere Aspekte wichtig.

Zunächst ist in nahezu allen realen Datensätzen die Frage nach Abhängigkeiten von Interesse und hier speziell die Frage, ob die Umsätze pro Kunden von der Tageszeit oder dem Tag selber abhängig sind. Diese Frage wird hier mit einem zunächst qualitativen Verfahren, der Clustering der Daten und deren Visualisierung, aufgenommen.

Als zweiter wichtiger Aspekt ergibt sich der Vergleich von Verteilungen, der als fundamentaler Schritt im Aufbau statistischen Wissens eingeschätzt wird (Makar und Confrey 2004). Der Vergleich kann hier zunächst über die Betrachtung der Form der Verteilung, dann aber über die Betrachtung der charakteristischen Kennzahlen, also der verfügbaren Lage- und Streuparameter erfolgen (vgl. Tabelle 2).

| Zeitpunkt | $\bar{x}_{0,5}$ | \bar{x} | SW | $Q_{0,5}$ | σ |
|-----------|-----------------|-----------|------|-----------|----------|
| Morgen | 13,0 | 15,7 | 59,6 | 13,7 | 11,2 |
| Mittag | 12,4 | 14,5 | 76,9 | 11,9 | 10,1 |
| Abend | 10,5 | 16,0 | 55,1 | 17,2 | 17,2 |
| Montag | 13,1 | 9,9 | 42,7 | 14,1 | 9,6 |
| Dienstag | 13,2 | 11,6 | 59,6 | 10,1 | 9,8 |
| Freitag | 16,5 | 14,5 | 72,0 | 12,3 | 12,4 |
| Samstag | 17,4 | 15,0 | 55,1 | 13,7 | 12,0 |

Tabelle 2: Kennzahlen zur Umsatzverteilung: Median, arithmetisches Mittel, Spannweite, Quartilsabstand und Standardabweichung

Eine wichtige Erweiterung zum Vergleich der Verteilungen ist die Simulation. Dabei wird nebenbei ein wichtiger Aspekt des statistischen Denkens berührt, nämlich die Erkenntnis der Variabilität statistischer Daten (Pfannkuch und Wild 1999). Nimmt man an, dass die Umsätze pro Person unabhängig von den Zeitpunkten der Messung sind, so kann man zufällige Stichproben aus dem Gesamtdatensatz ziehen (mit dem gleichen Umfang, der für verschiedene Messzeitpunkte vorhanden ist). Diese simulierten Stichproben sollten dann wiederum eine der realen Stichprobe ähnliche Verteilung aufweisen. Durch die Simulation wird nun gerade die Erkenntnis der Variabilität gefördert. Jede Simulation ergibt eine neue, den anderen zwar ähnliche, im Detail dennoch unterschiedliche Verteilung (vgl. Abb. 7). Man kann aber trotz der *Variabilität* der Einzelverteilungen das Muster, die immer ähnliche Form der Verteilungen, erkennen.

Mit der statistischen Methode des Anpassungstests, der – wenn überhaupt – ein Thema am Ende der Schulstochastik ist, kann man die qualitative Untersuchung schließlich statistisch absichern. Dazu kann man beispielsweise die Cluster verschiedener Wochentage heranziehen und untersuchen, ob Unterschiede in den Verteilungen hinsichtlich der verschiedenen Mess-Zeitpunkte innerhalb des Zufälligen liegen, also etwa mit dem χ^2 -Anpassungstest sich jeweils die Hypothese der identischen Verteilung beibehalten lässt (die folgenden Schritte sind nur in Kürze beschrieben. Sie dienen nicht dazu, das Vorgehen des Anpassungstests zu erlernen, sondern sollen allein das Nachvollziehen in Kenntnis dieses

Tests ermöglichen).

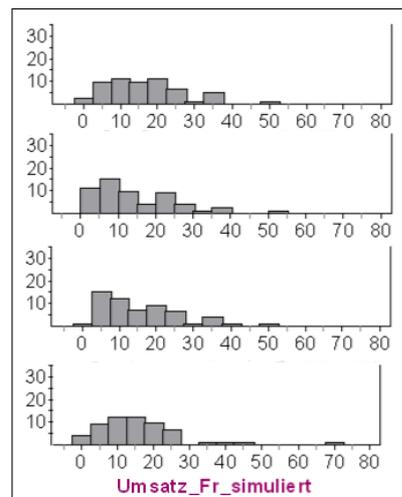


Abb. 7: Simulation von vier Freitagen (bzw. 57 Umsätzen)

Die Frage, die bei dem χ^2 -Anpassungstest (hier in der Ausführungsvariante von Fisz 1973) beurteilt werden soll, ist die folgende: Passt die empirisch ermittelte Verteilung der Umsätze eines bestimmten Wochentags mit der ebenfalls empirisch ermittelten Gesamtverteilung der Umsätze überein oder sind die Abweichungen der beiden Verteilungen voneinander überzufällig? Dazu werden beispielsweise hinsichtlich des Freitags die vorhandenen 57 Daten in 5 Klassen mit 10 Daten und ein Randklasse mit 7 Daten aufgeteilt. Die Häufigkeit dieser Klassen (k_i) ist $H_{fr}(k_i) = 10$ für $i = 1, \dots, 5$ und $H_{fr}(k_i) = 7$ für $i = 6$; vgl. Tabelle 3). Als Vergleichsmaßstab gilt die Verteilung aller Umsätze. Mit der gleichen Klassierung ergeben sich die absoluten Häufigkeiten $H_{ges}(k_i)$. Daraus ergibt sich – ebenfalls mit der gleichen Klassierung – die theoretische Häufigkeit $H_{theo}(k_i)$ der Umsätze (wenn man die relative Häufigkeit hinsichtlich der Gesamtstichprobe als Wahrscheinlichkeit für die Klasse setzt) der 6 Klassen durch:

$$H_{theo}(k_i) = \frac{H_{ges}(k_i)}{n_{ges}} \cdot n_{fr} = \frac{H_{ges}(k_i)}{372} \cdot 57$$

Mit diesen Vorgaben lassen sich die einzelnen Summanden der χ^2 -Testverteilung mit

$$t_i = \frac{(H_{theo}(k_i) - H_{fr}(k_i))^2}{H_{theo}(k_i)}$$

bestimmen, die Summe berechnen und mit Hilfe der χ^2 -Verteilung die Anpassungshypothese beurteilen. In diesem Fall ergibt sich als Testgröße

$$t = \sum_{i=1}^6 t_i \approx 5,37$$

Der Vergleich mit der χ^2 -Verteilung (mit 5 Freiheitsgraden) ergibt hinsichtlich des Wertes ihrer Verteilungsfunktion: $F_{\chi^2_{r=5}}(5, 37) \approx 0,627$, d.h. die Hypothese der identischen Verteilung kann nicht abgelehnt werden.

| Klasse (i) | $H_{fr}(k_i)$ | $H_{ges}(k_i)$ | $H_{theo}(k_i)$ | t_i |
|----------------|---------------|----------------|-----------------|-------|
| 1 | 10 | 83 | 12,7 | 0,58 |
| 2 | 10 | 68 | 10,4 | 0,02 |
| 3 | 10 | 70 | 10,7 | 0,05 |
| 4 | 10 | 37 | 5,67 | 3,31 |
| 5 | 10 | 82 | 12,6 | 0,52 |
| 6 | 7 | 32 | 4,9 | 0,90 |
| Summe | 57 | 372 | 57 | 5,37 |

Tabelle 3: Intervalle zum Anpassungstest

Man erkennt an der Formel für die Testgröße t , dass eine große Abweichung der theoretischen Häufigkeit von der in der Stichprobe vorhandenen (etwa in der Klasse 4, siehe Tab. 3) einen großen Summanden hinsichtlich der Testgröße t ergibt. Insgesamt erkennt man ebenso, dass die für die reale Verteilung am Freitag durch die Wahl der Klassengrenzen künstlich erzeugten annähernden Gleichverteilung ebenso annähernd in der theoretischen Verteilung auftritt.

Ebenso wie bei der Verteilung für den Freitag kann auch für die übrigen Wochentage keine Abweichung von der Verteilung aller Umsätze pro Kunden statistisch nachgewiesen werden. Das ist ein erstaunliches Resultat, das möglicherweise wiederum auf die Charakteristik des Supermarkts, kein Ziel größerer Wochenendeinkäufe zu sein, zurückzuführen ist.

Der beschriebene Ansatz der Datenanalyse lässt sich weiter ausbauen. Mögliche Untersuchungsspekte könnten beispielsweise die Beurteilung der Abhängigkeit des Umsatzes pro Kunden zu den Merkmalen „Verkäuferin“ oder auch „Kasse“ sein.

5 Dem Gesamtumsatz auf der Spur

Einen wesentlichen Anteil daran, die Spannung bei der Sammlung der Kassensbons aufrecht zu erhalten, hat die Untersuchung des Gesamtumsatzes ausgemacht, der sich aus den Umsätzen der vier Kassen ergibt. Eine weitere Klasse statistischer Methoden kommt bei der Analyse dieser Daten ins Spiel, bestehend aus den Verfahren der Regression und Korrelation. Um einen zumindest kleinen Eindruck der genannten Spannung zu vermitteln, werden im ersten Abschnitt die ersten Erkenntnisse in chronologischer Reihenfolge dargestellt und erst anschließend eine umfassendere Analyse vorgenommen.

5.1 Die ersten Schritte

Nach der Analyse der ersten Kassensbons-Daten (vgl. Abb. 1) war nicht eindeutig geklärt, was die Kodierungen im unteren Teil des Bons tatsächlich bedeuten. Insbesondere war zunächst die Zuordnung der ersten vierziffrigen Codes zur Nummer des Kunden (periodisch von 1 bis 9999) sowie die Zuordnung der ersten dreiziffrigen Codes zur Kasse unklar. Die sich allmählich herauschälende Hypothese in diesem Sinne führte beschreibend zur Darstellung der Kundennummern an einer Kasse in Form einer Punktwolke bzw. einer Zeitreihe (vgl. Abb. 8). Die Skalierung der Zeitachse beschreibt die Öffnungszeit in Stunden seit dem 6. November 2006 (täglich 12 Stunden ohne Sonn- und Feiertage).

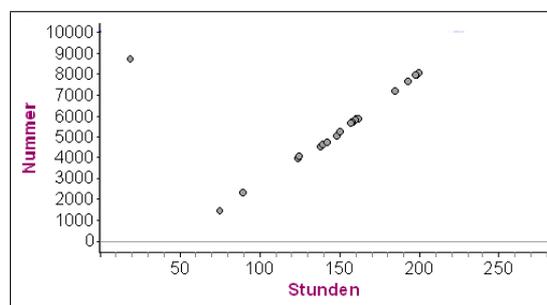


Abb. 8: Zeitreihe der Kundennummern an Kasse 1

Aus einer solchen Darstellung wurden folgende Hypothesen gewonnen:

1. Der erste vierziffrige Code beschreibt die Kundennummer, fortlaufend und sich periodisch wiederholend von 1 bis 9999 (ohne diese Hypothese wäre der linke Punkt in der Wolke kaum zu erklären gewesen)
2. Die Kundennummer kann gut durch eine lineare Funktion (in Abhängigkeit von der Öffnungszeit) beschrieben werden.

Mit der zweiten Hypothese geht die Möglichkeit einer Prognose einher, zu welchem Zeitpunkt der 10000. Kunde durch die Kasse 1 geht, wann also der Umsprung von 9999 auf 1 erfolgt. Diese Prognose ist einerseits ebenfalls eine Hypothese, andererseits aber auch ein Kriterium für die bisher aufgestellten Hypothesen gewesen. Die Schätzung des „Umsprungs“ sowie die Einpassung der Regressionsgeraden (hier noch schlicht nach Augenmaß) sowie die Begrenzung der Kundennummer ($Nummer < 10000$, $Nummer \in \mathbf{N}$) sind in der folgenden Abbildung (Abb. 9) aufgenommen.

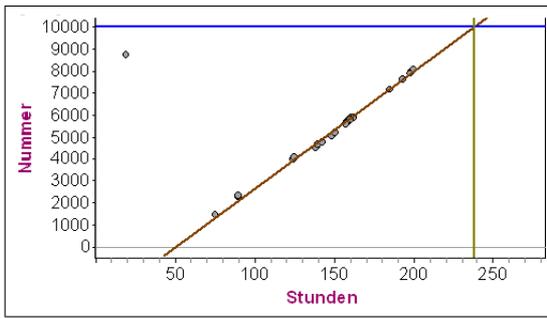


Abb. 9: Zeitreihe der Kundennummern an Kasse 1

Ergebnis dieser ersten Prognose: nach (ungefähr) 234 Stunden (nach dem 6.11.2006, 8 Uhr), also (ungefähr) nach 19 Tagen und 6 Stunden Öffnungszeit, erfolgt der Umsprung. Bezogen auf den Startpunkt bezieht sich die Prognose des Umsprungs auf Dienstag, den 28. November 2006, gegen 14 Uhr. Ab diesem Zeitpunkt sollten dann, so die erweiterte Hypothese, die Punkte der Zeitreihe durch eine Gerade mit identischer Steigung gegenüber der vorangegangenen Periode beschreibbar sein.

Abbildung 10 zeigt Entwicklung der Zeitreihe nach der Sammlung der nächsten Kassensbons bis kurze Zeit nach dem erwarteten Umsprung.

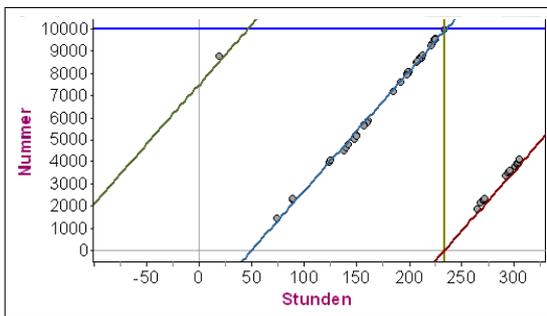


Abb. 10: Zeitreihe der Kundennummern an Kasse 1

Tatsächlich hat – von links kommend – der zuletzt gesammelte Kassensbon, der sehr nahe an dem prognostizierten Termin liegt (Dienstag, 28.11.2006, 13.45 Uhr), die Nummer 9946. Ebenso erkennt man zumindest den prognostizierten Umsprung sowie die Periodizität der Regression. Die eingepasste Gerade für das unten dargestellte Intervall hat die Gleichung (mit gerundeten Achsenabschnitten sowie mit der verkürzten Bezeichnung t für die Variable *Stunden*)

$$g(t) := \begin{cases} 54 \cdot t + 7400 & \text{für } 0 \leq t < 48,15 \\ 54 \cdot t - 2600 & \text{für } 48,15 \leq t < 233,34 \\ 54 \cdot t - 12600 & \text{für } 233,34 \leq t < 330 \end{cases}$$

Die weiteren Schritte der Datenanalyse nach diesem Heureka-Effekt sind im Folgenden dargestellt.

5.2 Der Umsatz der einzelnen Kassen

Für die mathematisch gesteuerte Regression sowie die Beurteilung der linearen Abhängigkeit der Kunden-Nummern von der Öffnungszeit kann man zunächst die Umsprünge vermeiden und die Nummern so transformieren, dass sie monoton wachsen (vgl. Abb. 11). Die Variable *Nummer_trans* bezeichnet hier damit die Anzahl der Kunden seit Untersuchungsbeginn.

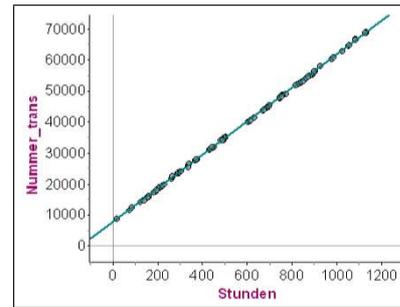


Abb. 11: Zeitreihe zu Kasse 1

Nach der Methode der kleinsten Quadrate für die Regression ist die Gerade mit der Gleichung

$$g(t) = 53,97 \cdot t + 7370$$

optimal. Der Korrelationskoeffizient für diese Gerade ist gerundet 1, also ist die Güte der linearen Abhängigkeit annähernd optimal. An den Datenpunkten selbst kann man in diesem Fall erkennen, dass der angezeigte aber gerundete Wert des Korrelationskoeffizienten von 1 *nicht* bedeutet, dass alle Punkte auf einer Geraden liegen. Die Abweichungen sind allerdings gering (vgl. Abb 12).

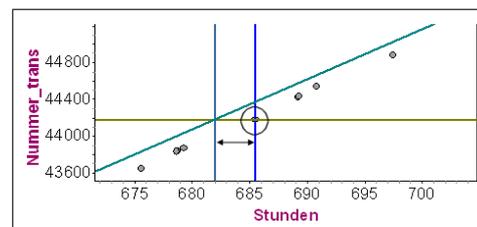


Abb. 12: Ausschnitt aus der Zeitreihe der Anzahl der Kunden an Kasse 1

In diesem Fall ist die Residuenbestimmung in Richtung der Zeitachse aussagekräftiger als die übliche Bestimmung in vertikaler Richtung, also in Richtung der Achse der transformierten Kunden-Nummern. So gibt die Bestimmung der horizontalen Abweichungen an, wie weit zeitlich die Abweichungen der realen von den idealen Kundennummern sind. Für das markierte Datum (vgl. Abb 12) gilt etwa, dass diese

Kundennummer gegenüber der idealen Kundennummer (eingepasste Gerade) um rund dreieinhalb Stunden zu spät durch die Kasse gegangen ist. Betrachtet man alle Daten in dieser Weise, so ergibt sich im Mittel eine zeitliche Abweichung von rund 4 Stunden und 40 Minuten (vgl. auch Abb. 13 zu allen Residuen in diesem Sinne).

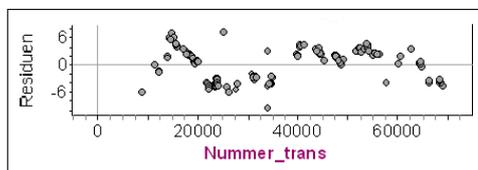


Abb. 13: Zeitliche Abweichungen der Anzahl der Kunden an Kasse 1 von der eingepassten Regressionsgeraden

Wie ist nun die Gerade zu interpretieren? Aus der Geradengleichung gewinnt man die Aussage, dass pro Stunde rund 54 Kunden bedient werden. Auch an dieser Stelle wäre es sinnvoll, solch einen Durchlaufkoeffizienten vorab schätzen zu lassen. So ist der mittlere Durchlauf – trotz Kartenzahler, voller Einkaufswagen, mühsam die Cent suchende Kunden, fehlender Auswägungen von Frischgemüse und den, an Kasse 1 allerdings seltenen, Leerlaufzeiten – recht hoch. Es ist möglicherweise ein Durchlauf, in dem sich wiederum die Struktur des Supermarkts zeigt, in der die großen Wocheneinkäufe eher selten sind (siehe oben). Aus dem Durchlauf pro Stunde ergibt sich weiterhin, dass im Mittel etwa alle 185 Stunden bzw. alle 15 Tage und 4 Stunden 10 000 Kunden die Kasse 1 des Supermarkts passieren.

| Kasse | Gerade Korrelationskoeffizient | Durchlauf pro Stunde, entspricht 10000 Kunden in |
|-------|---|--|
| 2 | $g_2(t) = 46 \cdot t - 3100$ $r = 1$ | 46,0 18 d 1 h |
| 3 | $g_3(t) = 25,9 \cdot t + 4100$ $r = 1$ | 25,9 32 d 2 h |
| 4 | $g_4(t) = 14,4 \cdot t + 4200$ $r = 1$ | 14,4 57 d 10 h |

Tabelle 7: Kunden-Durchläufe

Für die anderen Kassen, die bei Bedarf der Kasse 1 zugeschaltet werden, ergeben sich nach dem gleichen Analysemuster die in Tabelle 7 angeführten Geradenanpassungen bzw. Durchlaufkoeffizienten (vgl. auch Abb. 14).

Man erkennt an den Werten in der Tabelle 7 deutlich die übergreifende Strategie, die Kassen in der Reihenfolge ihrer Nummerierung zur Hauptkasse, der Kasse 1, hinzuschalten.

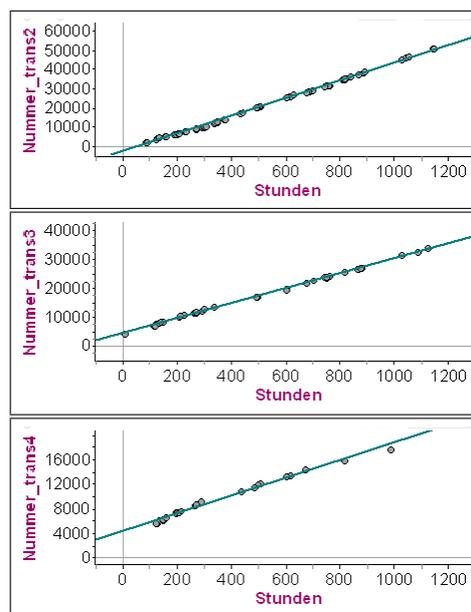


Abb. 14: Zeitreihe der Kundennummern an den Kassen 2, 3 und 4

5.3 Der Umsatz des Supermarkts

Sind sowohl der Umsatz pro Kunden als auch der Kundendurchsatz bekannt, so kann man natürlich unmittelbar die Höhe der Kasseneinnahmen des Supermarkts in bestimmten Zeitabschnitten berechnen. Während hier der Kundendurchsatz sehr genau bestimmt werden kann, ist für den Umsatz pro Kunden ein Intervall gegeben, das bei angemessener Güte der Daten mit hoher Wahrscheinlichkeit den realen mittleren Umsatz pro Kunde enthält.

Es ergibt sich pro Tag im Mittel der Durchsatz von $12 \cdot (53,97 + 46 + 25,9 + 14,4) \approx 1683$ Kunden. Legt man weiterhin das mit dem Standardverfahren bestimmte Konfidenzintervall zum mittleren Umsatz pro Kunden, so erhält man für den Umsatz pro Tag in Euro das Intervall $[23900; 27600]$ (gerundet). Insgesamt ergibt sich

| Zeitspanne | Einnahmen in Euro | |
|-----------------|-------------------|--------------|
| | untere Grenze | obere Grenze |
| Tag | 23900 | 27600 |
| Woche | 167300 | 193200 |
| Monat (26 Tage) | 621500 | 717700 |
| Jahr (2007) | 7242300 | 8364400 |

Tabelle 8: Einnahmen des Supermarkts

5.4 Kalkulation eines Supermarktes

Aus dem Kassenumsatz die Kalkulation des Supermarkts zu rekonstruieren, ist ein komplexes Problem und würde den Umfang des Artikels sprengen. Es sollen nur kurze Hinweise gegeben werden, die projektartig weiter verfolgt werden könnten, aber weit über die hier erfolgte Datenanalyse hinausgehen.

Ein Supermarkt hat das verständliche Ziel, Gewinn zu machen. Von den Einnahmen an der Kasse abzurechnende Kosten entstehen durch

- die monatliche Miete; solche Informationen kann man aus dem Netz beziehen.
- Nebenkosten, die schwer abzuschätzen sind (extensives Heizen im Winter, Kühlung von Waren etc.).
- das Personal. Das monatliche Gehalt einer Verkäuferin bzw. eines Verkäufers ist (im Mittel) auch im Netz verfügbar⁴. Zuzurechnen wären die vom Arbeitgeber zur Hälfte zu bezahlenden Sozialabgaben. Zudem müsste man eine Schätzung machen, wie viele Angestellte ein Supermarkt abgesehen von denjenigen an den Kassen beschäftigt.
- die an den Staat gehende Mehrwertsteuer.
- natürlich durch den Einkauf der Waren. Kann man nach Abzug der an den Staat gehenden Mehrwertsteuer die Miet-, Neben- und Personalkosten bestimmen, so hat man die Möglichkeit, den Aufschlag des Supermarkts zu bestimmen, ab dem ein Gewinn (ohne Beachtung der Steuer) entsteht.

6 Didaktische Anmerkungen

In der Analyse von Daten, die auf dem Boden liegen, steckt einiges didaktisches Potenzial. Zunächst ist die Analyse der Daten in unterschiedlichen Anforderungsstufen in verschiedenen Klassenstufen möglich, prinzipiell vom Ende der Primarstufe bis in die Sekundarstufe II. Die Spannweite liegt inhaltlich von einfachen Häufigkeitsbestimmungen bis hin zu Verfahren der beurteilenden Statistik wie Anpassungstest. Innerhalb einzelner Untersuchungsschritte kann die Schwierigkeit variiert werden. So sind an manchen Stellen wie etwa dem Vergleich von Umsätzen pro Kunden an bestimmten Tagen mit dem allgemeinen Umsatz pro Kunden verschiedene Methoden einsetzbar, für die Sekundarstufe I eine qualitative Untersuchung von Verteilungen durch Betrachtung von Mittelwerten und Streuung, für die Sekundarstufe evtl. mit dem Einsatz von inferenzstatistischen Verfahren.

Ein Unterrichtsprojekt zu Kassensbons enthält weiterhin alle Stufen, die eine „komplette statistische Untersuchung“ (Biehler, Hartung 2006, 53) aufweist: „*Problemstellung – Planung der Erhebung – Datenerhebung – Auswertung – Interpretation – Schluss-*

folgerung und Ergebnisbericht“. Es ist zudem über die Betrachtung der Kalkulation eines Supermarkts noch erheblich erweiterbar.

Versteht man die Datenanalyse als *Kompetenz*, die Schüler in ihrer Schulzeit erwerben sollen, so kann man die Aufteilung des statistischen Denkens, die Wild und Pfannkuch (1999) formuliert haben, als Auflistung von Teilkompetenzen verstehen:

- Erkennen der Notwendigkeit von Daten
- Flexible Aufbereitung von Daten
- Erkennen der Variabilität von Daten
- Erkennen und Beschreiben von Mustern in den Daten
- Verbindung mit dem Sachkontext.

Im Sinne dieser Teilaspekte des statistischen Denkens soll die Datenanalyse zu den Kassensbons noch einmal rückblickend in Kürze betrachtet werden.

Erkennen der Notwendigkeit von Daten: Für einen Statistiker ist die Formulierung dieser Stufe wahrscheinlich banal, für Schüler zunächst einmal nicht. Hier ist die Erfahrung aufzubauen, dass über das Vergöbern von Datensätzen Beurteilungen von Situationen erst einmal über „Stammtischmeinungen“ bzw. naive Schätzungen hinauskommen.

In dem vorgestellten Projekt wird die Notwendigkeit einer vergrößerten Datenbasis in mehreren Phasen deutlich. Zunächst braucht man mehr Daten (Kassensbons), um die Kodierungen des Kassenzettels zu verstehen. Glaubt man zu verstehen (hypothetisch), so braucht man mehr Daten, um seine Hypothesen zu bestätigen (etwas die Kodierung der Kundennummer). Will man weitere Untersuchungen machen (etwa zum Umsatz pro Kunden aufgeteilt nach Verkäuferinnen), so braucht man unter Umständen noch mehr Daten und kann dabei die Grenzen der Aussagefähigkeit eines begrenzten Datensatzes erkennen.

Flexible Aufbereitung: Die flexible Datenaufbereitung durchzieht die explorative Datenanalyse der Kassenzettel. So werden die Daten nach unterschiedlichen Gesichtspunkten grafisch dargestellt und analysiert, etwa als Gesamtdatensatz oder geclustert nach Aspekten wie Einzeltagen oder Tageszeiten.

Die eingesetzten Methoden können breit gestreut werden. So kann die Analyse eher qualitative Untersuchungen verschiedener Diagramme (Boxplots oder Säulendiagramme), oder eher quantitativ-statistische

Verfahren wie etwa den Anpassungstest umfassen. Die flexible Datenaufbereitung hat dabei zum Ziel, in möglichst vielfältiger Weise Muster in den Daten zu erkennen und beispielgebunden Möglichkeiten und Grenzen einzelner statistischer Methoden zu erfahren.

Erkennen der Variabilität von Daten: Dieser Aspekt wird in der Analyse der Kassensbons eher unterschwellig behandelt. Dennoch ergibt sich etwa bei der Analyse der nach Wochentagen geclusterten Umsätze pro Kunden stets die Frage, ob die Unterschiede in den Verteilungen auf eben der Variabilität der Daten beruhen oder ob systematische Unterschiede bestehen. Einen etwas expliziteren Zugang zum Erkennen der Variabilität kann dagegen die in Kapitel 4.1.2 simulierte Verteilung von Umsätzen pro Person an vier virtuellen Freitagen leisten.

Erkennen und Beschreiben von Mustern in den Daten: Die Suche nach Mustern in den Daten ist ein Kernstück der Datenanalyse. Hier geht es einerseits um die qualitative Beschreibung von Mustern, etwa der Verteilung von Umsätzen pro Kunden. Die statistische Analyse mit Hilfe des Anpassungstests hilft bei der Beurteilung, ob ein Muster (in den Umsätzen pro Kunden) mit den in den Daten steckenden empirischen Phänomenen wie etwa der Umsatzverteilung pro Kunden an bestimmten Tagen übereinstimmt. Schließlich lässt sich der Durchlauf von Kunden an den Kassen in fast idealer Form mit Hilfe eines Musters, einer Regressionsgeraden, beschreiben und die Güte des Musters mit Hilfe des Korrelationskoeffizienten beurteilen.

Verbindung mit dem Sachkontext: Die Verbindung der Datenanalyse mit dem Sachkontext ist bei realen Daten fast von selbst gegeben. So ist etwa die Beurteilung der rein statistisch erzeugten Ergebnisse untrennbar mit der über die mathematische Verarbeitung der Daten hinausgehende Betrachtung der Charakteristik eines Supermarkts (Kap. 2.1) oder die daran anschließende Beurteilung der Datengüte verbunden (Kap. 2.3).

7 Eurodebatte und die Exponentialverteilung, ein historischer Nachtrag

Als der zweitgenannte Autor im April 2001 mit seinen Schülern erstmals auf Kassenzettel-Jagd ging, dachte noch niemand an den Euro.

Im April 2002 war er schon 4 Monate als der berühmte „Teuro“ unter uns. Und zu allem Ärger hingen da seit Wochen diese Plakate „im ALDI“,

die weismachen wollten, dass bei den Lebensmittel-discountern alles billiger geworden ist, erinnern Sie sich? Wieder so eine geschickte Werbemasche? Hat man doch oft gehört. Und überall gelesen. Kurz vor der Euromstellung Preise rauf, um dann im Frühjahr ein bisschen billiger werden zu können. Und satte Preiserhöhungen unter dem Strich.

Die Stochastik-Kurse am Heinrich-Mann-Gymnasium in Köln gingen erneut auf Zetteljagd und das Stöbern in dem alten Zettelhaufen begann von vorne. Man musste auf den alten (zum Glück noch nicht entsorgten) und den neuen Kassenzetteln gleiche Waren finden. Die ersten Ergebnisse zeigen: Brechbohnen wurden teurer (+), Toilettenpapier billiger (–). Was sich hinter „Wacholderbauch“ (–) verbirgt, blieb zunächst ungeklärt (1 € = 1,95583 DM).

| Ware | März 2001 [März 2002] | Diff. | Vorz. |
|--------------------|----------------------------|--------|-------|
| Backofen Pommes | 1,32DM (0,67€) [0,86€] | 0,19€ | + |
| Bandnudeln | 1,29DM (0,66€) [0,65€] | –0,01€ | – |
| Coca Cola 2l | 2,79DM (1,43 €) [1,27€] | –0,16€ | – |
| Fischstäbchen | 2,49DM (1,27 €) [1,39€] | 0,12€ | + |
| Brechbohnen | 1,39DM (0,71€) [0,79€] | 0,08€ | + |
| Gouda | 4,17DM (2,13€) [2,27€] | 0,14€ | + |
| Klarspüler | 1,99DM (1,02€) [1,17€] | 0,15€ | + |
| Toilettenpapier | 4,98DM (2,55€) [2,49€] | –0,06€ | – |
| Wachholderbauch | 1,99DM (1,02€) [0,99€] | –0,03€ | – |
| Waffeleier | 1,29DM (0,66€) [0,69€] | 0,03€ | + |

Tabelle 9: kleine Stichprobe vom Umfang 10

Und mit dieser Stichprobe war man beim einfachsten Verfahren der beurteilenden Statistik, beim Vorzeichen-test, gelandet. Und zu unserer Verblüffung – und allen Unkenrufen zum Trotz – mussten wir feststellen, dass tatsächlich von den 360 verglichenen Preisen 213 (das sind 59 Prozent) gesenkt worden waren.

Genauere Untersuchungen zerrten dann den tieferen Grund für diese vielen Preissenkungen ans Tageslicht: Nicht eine barmherzige Kundenfreundlichkeit, sondern die sublimen Psychologie der „Schwellenpreise“: Wenn früher ein Produkt 1,99 DM gekostet hätte, müsste es nun wegen des Umrechnungskurses 1,01 € kosten. Wenn man daraus 0,99 € macht, hat ALDI 2 Cent „verschenkt“, kann

aber wieder einen psychologisch günstigen Schwellenpreis anbieten. Diese Vielzahl minimaler Preissenkungen wurden durch massivere Preiserhöhungen bei einzelnen Produkten ausgeglichen. Unsere „ALDI-Kassenzettelforschung“ ergab insgesamt eine (recht moderate) Preiserhöhung von ca. 2 Prozent während eines Jahres. Wer mochte, konnte darauf guten Gewissens reimen: „Willst du keinen Teuro haben, musst du dich bei ALDI laben“.

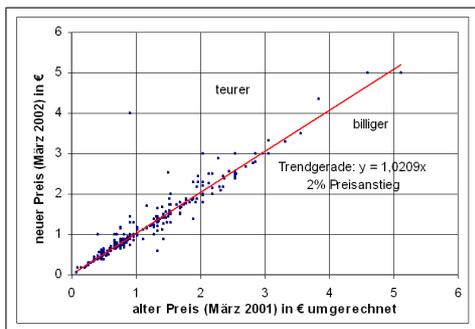


Abb. 15: Alte und neue Preise im Vergleich

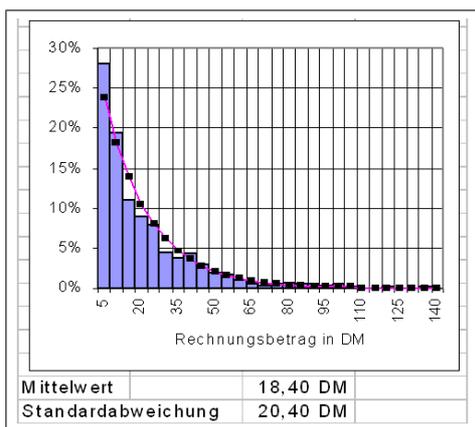


Abb. 16: Exponentialverteilung der Rechnungsbeträge (Säulen) mit virtuellen Kunden

Eine weitere überraschende Entdeckung machten Paul, David und Wladi, die inzwischen alle Mathematik und Informatik studieren: Die Rechnungsbeträge sind nicht, wie man vielleicht erwarten würde, normalverteilt: Unabhängig von der Supermarktkette erwiesen sie sich als exponentialverteilt (vgl. Abb. 16). Das gilt für Aldi, für Lidl, für Plus und selbst für Kaiser's. Das „Verrückte“ ist, dass man genau diese Verteilung erhält, wenn man im Rahmen einer Modellrechnung „virtuelle Kunden“ beim Einkauf durch den Laden schlendern lässt, die jede Ware, an der sie vorbeikommen, rein zufällig und unabhängig von

dem, was schon im Korb ist, mit einer festen Wahrscheinlichkeit hinzupacken.

Anmerkungen

- ¹ Die korrekte Bezeichnung für diesen Berufsstand ist Einzelhandelskauffrau bzw. Einzelhandelskaufmann. Statt dieser etwas sperrigen Bezeichnung wird kurz von Verkäuferinnen gesprochen und hier vereinfachend auch nur die weibliche Form verwendet. Im Gegensatz dazu wird bei den Kunden nur die männliche Form verwendet.
- ² Die Supermärkte haben verschiedene Vorgehensweisen bei den Angaben auf ihren Kassensbons. Welche Kodierungen von einzelnen Märkten aufgeführt werden, kann daher nicht allgemein angegeben werden. Die De-Codierung ist damit immer wieder eine neue Aufgabe.
- ³ Vgl <http://www.boeckler.de/cps/rde/xchg/SID-3D0AB75D-FEC5218E/hbs/hs.xml/550.html>.

Literatur

- Biehler, R., Hartung, R. (2006). Die Leitidee Daten und Zufall. In Blum, W. Druke-Noe, C., Hartungs, R., Köller, O. (Hrsg.), Bildungsstandards Mathematik: konkret. Berlin: Cornelsen-Scriptor, S. 51-80.
- Eichler, A. (2007) „Geld weg — Arzt weg!“ — Was ist dran am Ärzteprotest? Praxis der Mathematik 49(13), S. 20-26.
- Fisz, M. (1973) Wahrscheinlichkeitsrechnung und mathematische Statistik. Berlin: VEB.
- Herget, W. (1994) Artikelnummern und Zebrastrifen, Balkencode und Prüfziffern – Mathematik und Informatik im Alltag. In Blum, W., Henn, W., Klika, M., Maass, J. (Hrsg.), Materialien fuer einen realitätsbezogenen Mathematikunterricht. Bd. 1. Hildesheim: Franzbecker. 1994. S. 69-84.
- Makar, K. und Confrey, J. (2004). Secondary teachers' statistical reasoning in comparing two groups. In Ben-Zvi, D., und Garfield, J. B. (Hrsg.), The challenge of developing statistical literacy, reasoning, and thinking (S. 353-373). Dordrecht: Kluwer.
- Pfannkuch, M., Wild, C. (1999). Statistical Thinking in Empirical Enquiry. In: International Statistical Review 67(3) 1999, S. 223-248.
- Riemer, W. (2003). Stochastik. Stuttgart: Klett-Verlag, S.154.
- Riemer, W. (2006). Lambacher-Schweizer 6 NRW. Stuttgart: Klett-Verlag, S.164.
- Tversky, A. und Kahnemann, D. (1973). Availability: A heuristic for judging frequency and probability. Cognitive Psychology 5, S. 207-232.

Anschriften der Verfasser

Andreas Eichler
 Institut für Didaktik der Mathematik
 Westfälische Wilhelms-Universität Münster
 Fliegerstraße 21
 48149 Münster
a.eichler@uni-muenster.de
 Wolfgang Riemer
 Studienseminar für Lehrämter an Schulen
 Claudiusstraße 1
 50678 Köln