

## Schlussbemerkungen

- Dieses neu entwickelte Excel-Tabellenblatt befreit den Stochastikunterricht von aufwendigen und ungenauen Verfahren. Sie werden ersetzt durch ein vielfältig anwendbares experimentell-numerisches Vorgehen, das für Schüler deutlich leichter zu verstehen und anzuwenden ist. Neue Möglichkeiten eröffnen sich bezüglich der Aufgabengestaltung und der Vermittlung der Inhalte. Es wird Raum geschaffen, um weiterführende Themen in Angriff nehmen zu können (z.B. Markoff-Ketten, Chi-Quadrat-Test).
- Auch Graphikfähige TR oder CAS könnten den klassischen Lehrgang mit Tabellen, Algebra und Transformationen durch ein experimentell-numerisches Vorgehen ersetzen. Dazu bedarf es aber eines ganz anderen Ansatzes.
- Ich verzichte seit Jahren auf tabellierte Wahrscheinlichkeitswerte und arbeite nur noch mit dem vorgestellten Excelblatt. Fast alle Schüler können mittlerweile zu Hause auf Excel zugreifen. Notfalls können sie für die Hausaufgaben die Rechner der Schulbibliothek nutzen.

In Kurs- und Abiturarbeiten werden 4 Notebooks mit der Exceldatei bereitgestellt, über die die jeweils benötigten Werte abgerufen werden können. Dabei lösen die Schüler die Aufgaben so weit wie möglich im Heft ohne PC. Die gewünschten Werte können dann gegen Ende – für alle Aufgaben zusammen – ermittelt werden.

- Das Tabellenblatt ist ausbaubar, so dass es auch im Hochschulbereich verwendet werden kann. Einige stochastische Funktionen wie die  $\chi^2$ -, t- und Fisher-Verteilung sind aus diesem Grund bereits eingearbeitet worden.
- Über meine Homepage kann die vorgestellte Exceldatei "stochastik.xls" kostenlos bezogen werden.

### Anschrift des Verfassers

Stefan Bartz  
Jakobstraße 16  
54636 Meckel  
[www.stefanbartz.de](http://www.stefanbartz.de)

---

1 Aus didaktischen Gründen werden die Begriffe "Binomialfunktion" für  $Bin_{n,p}(X=x)$  und "kumulierende Binomialfunktion" für  $Bin_{n,p}(X \leq x)$ , statt der für Schüler verwirrenden Begriffe "Binomialverteilung" und "Verteilungsfunktion der Binomialverteilung" verwendet. Die Schreibweisen knüpfen an die der Analysis an: Funktionsname mit 3 Buchstaben,  $x$  als Variablenname und in Klammern, Scharparameter als Index, also  $Bin_{n,p}(x)$  analog zu  $\sin_t(x)$ . Entsprechend wird bei den übrigen Verteilungen  $Hyp_{n,R,N}(x)$ ,  $Poi_{\mu}(x)$  und  $Nor_{\mu,\sigma}(x)$  verfahren.

---

## Stochastische Algorithmen I

PETER EICHELSBACHER, BOCHUM

**Zusammenfassung:** Mit Hilfe der Grundbegriffe der Wahrscheinlichkeitsrechnung lassen sich sogenannte stochastische Algorithmen auch im Schulbereich untersuchen. In einer ersten Dar-

stellung klären wir den Unterschied zu deterministischen Lösungswegen, erinnern an ein Kartenspiel und stellen den Quicksort-Algorithmus und seine randomisierte Variante vor.

### Einleitung

Was ist ein stochastischer Algorithmus? Sehr einfach beantwortet ist es ein Algorithmus, bei dem der Zufall mit ins Spiel kommt. Was könnte das Interesse sein, einen bekannten, nicht vom Zufall beeinflussten Algorithmus (man nennt dies einen deterministischen Algorithmus) zu randomisieren?

Zwei Motivationen sollen hier besprochen werden.

1. Algorithmen sind erfunden worden für Situationen *endlicher* Probleme, etwa das Sortieren von Gegenständen nach einer Ordnung, z.B. das Sortieren einer Liste von Zahlen, das Prüfen zweier Zahlenreihen auf Gleichheit, das Auffinden

eines optimalen Wertes einer Funktion auf einer endlichen Menge. Algorithmen suchen dabei möglichst effizient alle Möglichkeiten durch, und daher muss Endlichkeit gefordert werden. Doch schon einfachste Fragestellungen der obigen Art für endliche Situationen können zu Laufzeiten der Algorithmen führen, die z.B. die Lebensdauer der Erde bei weitem übersteigt, und benötigen Speicherplatz, der auch bei modernen Rechnern nicht zur Verfügung steht. Nun versucht man sein Glück mit randomisierten Algorithmen, um die Laufzeit zu verkürzen und den Speicherbedarf zu verkleinern.

2. Algorithmen sollen bei allen denkbaren Eingaben korrekt und schnell arbeiten. Häufig gibt es aber Eingabewerte, die zu einer sehr ungünstigen Laufzeit führen. Die Laufzeit-Analyse eines Algorithmus muss aber solche *worst case* Eingaben berücksichtigen, auch wenn viele Eingaben zu einem guten Verhalten des Algorithmus führen.

Beim *Auswürfeln* eines Eingabewertes würde man nun erwarten, dass sich der Algorithmus mit hoher Wahrscheinlichkeit gut verhält. Der stochastische Algorithmus sollte nun für alle Eingaben eine Laufzeit liefern, die mit hoher Wahrscheinlichkeit gutartig ist. Ungünstige Verläufe mögen mit geringer Wahrscheinlichkeit eintreten. Liegt ein Algorithmus vor, der eine sehr große Menge möglicher Fälle zu untersuchen hat, liegt die Hoffnung im randomisierten Algorithmus darin, dass das Erkunden eines kleinen Teils der großen Menge mit großer Wahrscheinlichkeit schon eine relevante Information erbringt.

Stochastische Algorithmen werden also in der Regel entweder die Lösung nicht immer mit Sicherheit finden oder nicht mit absoluter Präzision. Wir wollen ein Beispiel geben, um das Wort *Randomisierung* zu entdecken: Ein verrückter Millionär sendet an 20 Personen einen Brief mit dem Text:

*Diesen Text haben weitere 19 Personen erhalten. Wählen Sie aus, ob Sie mir zurückschreiben oder nicht. Kommt bei mir genau ein Brief an, so erhält der Absender eine Million Euro. Falls keiner oder mehr als einer eintrifft, so erhält niemand etwas. Der Kontakt untereinander ist verboten.*

Was ist die beste Strategie? Wenn es eine beste Strategie gibt, so kann man davon ausgehen, dass jeder der 19 Mitspieler diese Strategie unabhängig von den anderen verwendet. Es gibt aber nur die Strategie, einen Brief abzuschicken oder nicht. Beide Strategien haben aber nach den genannten Spielregeln keinen Erfolg!

Was passiert nun bei einer Randomisierung? Wir schicken mit Wahrscheinlichkeit  $p$  mit  $p \in (0, 1)$  einen Brief ab. Dann gewinnt man mit Wahrscheinlichkeit  $f(p) = p(1-p)^{19}$  eine Million. Der Wert dieser Funktion in  $p$  ist maximal bei  $p_0 = \frac{1}{20}$ , wie man leicht sieht: Es ist

$$f'(p) = (1-p)^{19} - 19p(1-p)^{18}$$

und  $f'(p) = 0$  genau dann wenn

$$(1-p)^{19} = 19p(1-p)^{18}, \text{ also } 1 = 20p.$$

Wenn wir uns eine Zufallsmaschine bauen, die mit Wahrscheinlichkeit  $\frac{1}{20}$  entscheidet, so schicken wir im positiven Fall den Brief ab, sonst nicht. Der randomisierte Algorithmus ist jedem deterministischen Algorithmus überlegen. Zugegeben, dieses Beispiel ist etwas konstruiert, aber es hebt auf einfache Weise mögliche Vorteile eines randomisierten Algorithmus gegenüber einem deterministischen hervor. In diesem Artikel wollen wir uns mit dem Sortieren einer langen Liste von Zahlen befassen. Hier ist eine randomisierte Version eines Sortier-Algorithmus unser Fokus.

## Erinnerung an ein Kartenspiel: Patience

Erinnerung soll hierbei bedeuten, dass wir das folgende Patience-Spiel in Eichelsbacher (2003) bereits einmal vorgestellt haben: Wir benötigen ein Kartenspiel von  $N$  Karten, die "linear geordnet" sind, d.h. für je zwei Karten lässt sich entscheiden, welche die höherwertige ist. Beispielsweise kann man von einem Kartenspiel nur die Herzkarten entnehmen und diese mit der Ordnung

$$A_s < 2 < 3 < \dots < B < D < K$$

versehen. Die Karten werden gemischt und auf einen Stapel gelegt. Von diesem Stapel wird nun sukzessive eine nach der anderen Karte von oben entnommen und in einem von mehreren Stapeln abgelegt. Die Regel dabei ist, dass eine Karte *nicht* auf einen Stapel gelegt werden kann, wenn sie höher als dessen oberste Karte ist. Gibt es keinen Stapel, auf den eine neu gezogene Karte gelegt werden kann, muss ein neuer Stapel angefangen werden. Ziel ist es, das Spiel mit so wenig Stapeln wie möglich zu beenden.

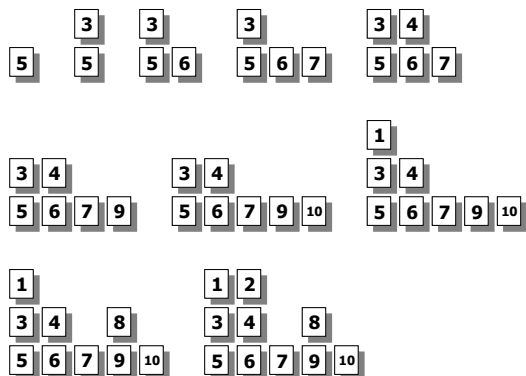
Wir wollen nun annehmen, das Kartenspiel sei gut gemischt (wie viele Mischvorgänge dafür notwendig sind, steht auf einem anderen Blatt und kann hier nicht im Detail erläutert werden).

Wir probieren nun, unser einfaches Patience-Spiel zu spielen. Wie kann man erreichen, dass zum Schluss so wenig Stapel wie möglich ausgelegt sind? Eine Strategie liegt auf der Hand: Wir legen jede neue Karte auf den erstmöglichen Stapel. Diese Strategie wollen wir die *gierige Strategie* nennen. Schauen wir sie uns anhand eines Beispiels an:

**Beispiel:** Der gemischte Kartenstapel von  $N = 10$  Karten sei gegeben durch

**5 3 6 7 4 9 10 1 8 2**

Spielen wir damit Patience gemäß der gierigen Strategie, erhalten wir die folgende Abfolge von Stapeln:



Wir haben unser Spiel also mit fünf Stapeln beendet. Aber wie gut ist das? Eichelsbacher (2003) leitet her, dass in dem oben dargestellten Patience-Spiel die gierige Strategie optimal ist (und mehr). Die geschilderte Variante des Patience-Spiels wird auch *Patience-Sortieren* genannt, und zwar aus folgendem Grund: Am Ende des Spiels liegt Karte 1 oben im ersten Stapel, wenn man Karte 1 wegnimmt, liegt Karte 2 oben auf einem der Stapel, wenn man diese wegnimmt, liegt Karte 3 oben auf einem der Stapel usw.

Mit Hilfe dieses Vorgehens kann man also die Karten von Hand sortieren. Wir haben also bereits einen spielerischen Sortier-Algorithmus gefunden. Wir hatten weiter in Eichelsbacher (2003) nach einer möglichen Computersimulation des obigen Patience-Spiels gefragt. Dieser muss ein mathematisches Modell zugrunde liegen. Hierbei ist die erste Frage, die es zu klären gilt: Was bedeutet es, dass wir das Kartenspiel zu Beginn gut mischen?

Der Sinn eines solchen Mischens ist es offenbar, die Karten aus einer möglicherweise bekannten Anfangssituation so durcheinander zu bringen, dass die Reihenfolge nicht mehr vorhersagbar ist. In mathematischer Sprechweise erhalten wir durch Mischen eine Permutation unserer Karten, oder, wenn wir die Karten mit 1 bis  $N$  durchnummerieren, eine Permutation der Zahlen  $1, \dots, N$ . Die Menge aller Permutationen von  $\{1, \dots, N\}$  heißt *symmetrische Gruppe* und wird mit  $S_N$  bezeichnet.

Ein Kartenspiel können wir als gut gemischt bezeichnen, wenn wir nach dem Mischen keine Idee mehr haben, welche Permutation vorliegt, d.h. wenn alle Permutationen  $\sigma \in S_N$  gleich wahrscheinlich sind. Nun hat  $S_N$   $N!$  viele Elemente, d.h. wir dürfen das Kartenspiel als gut gemischt bezeichnen, wenn für die Situation nach dem

Mischen gilt:

$$P(\sigma \text{ liegt vor}) = P(\sigma) = \frac{1}{N!}.$$

Wie lässt sich nun am praktischsten eine Permutation  $\sigma \in S_N$  simulieren? Formal lässt sich eine Permutation  $\sigma$  als eine Folge von  $N$  Paaren

$$(i, \sigma(i)) \text{ mit } 1 \leq i \leq N \text{ und } 1 \leq \sigma(i) \leq N$$

aufschreiben. Hierbei bezeichnet  $\sigma(i)$  den Platz an dem die Ziffer  $i$  in der Permutation  $\sigma$  erscheint. Die Permutation

$$\sigma = \{(1, 2), (2, 1), (3, 3)\} \in S_3$$

bedeutet, dass Karte 2 nach dem Mischen auf Platz 1, Karte 1 auf Platz 2 und Karte 3 auf Platz 3 gelandet ist. Manchmal schreibt man dafür kürzer

$$\sigma = (2 \ 1 \ 3) \in S_3$$

Nun ist es offenbar gar nicht wichtig, dass die  $\sigma(i) \in \{1, \dots, N\}$  sind. Wichtig ist allein, dass für  $i \neq j$  auch  $\sigma(i) \neq \sigma(j)$  ist. Dann kann man die  $i \in \{1, \dots, N\}$  in der auf- oder absteigenden Reihenfolge der  $\sigma(i)$  aufschreiben und erhält eine zufällige Permutation. Dies ist beispielsweise gewährleistet, wenn man für jedes  $i$  eine Zufallsvariable  $\sigma(i)$  zieht, die unabhängig von allen anderen und uniform im Intervall  $[0, 1]$  verteilt ist (eine solche ist leicht auf einem Computer zu erzeugen).

Hierbei bemerke man, dass für zwei uniform auf  $[0, 1]$  verteilte Zufallsvariablen  $X, Y$  mit Wahrscheinlichkeit eins  $X \neq Y$  gilt. (Dies ist eine theoretische Bemerkung; natürlich kann der Computer beim Generieren nur endlich viele Dezimalen einer Zufallszahl darstellen. Es kommt also vor, dass er zwei gleiche Realisierungen erstellt, aber sehr selten.) Es wird also stets für  $i \neq j$  entweder  $\sigma(i) < \sigma(j)$  oder umgekehrt  $\sigma(j) < \sigma(i)$  gelten. Ein Rezept zum Erzeugen einer Zufallspermutation lautet genauer also so:

- Erzeuge unabhängige, uniform in  $[0, 1]$  verteilte Zufallszahlen  $\sigma(1), \dots, \sigma(N)$ .
- Ordne  $\{1, \dots, N\}$  so, dass 
$$\sigma(i_1) < \sigma(i_2) < \dots < \sigma(i_N)$$
 gilt. Hierbei ist  $(i_1, \dots, i_N)$  die neue Anordnung der Menge  $\{1, \dots, N\}$ .
- Dann ist  $(i_1, \dots, i_N) \in S_N$  eine zufällige Permutation.

Es bleibt also festzuhalten: das Patience-Spiel bietet eine Möglichkeit, die Karten zu sortieren. Wenn Sie sich einen größeren Stapel von bereits

korrigierten Klausuren vornehmen und ihn alphabetisch sortieren wollen, so spielen Sie doch einfach Patience damit, es erfolgt eine (schnelle ?) alphabetische Sortierung. Interessant ist, dass man herausgefunden hat, dass man diesen Sortier-Algorithmus implementieren kann mit

$$N \cdot \log N - N \cdot \log \log N + O(N)$$

Vergleichen (zweier Werte) maximal. Des Weiteren hat kein anderer Algorithmus ein besseres *worst case* Verhalten. Dies findet man in Fredman (1975). Es bezeichnet hier  $O(N)$  eine Größe, die durch ein  $cN$  mit  $c$  konstant nach oben abgeschätzt werden kann.

Die Simulation einer Mischung des Kartendecks kann am Rechner mit Hilfe der uniform verteilten Zufallszahlen erfolgen, allerdings benötigt man hier intern bereits einen Sortieralgorithmus!

## Quicksort

Der Algorithmus mit Hilfe des Patience-Spiels ist natürlich nicht eine Vorgehensweise, auf die man in natürlicher Art und Weise stoßen würde. Nach der spielerischen Variante, die ein Sortieren ermöglicht, und deren Computer-Implementierung selbst einen Sortier-Algorithmus benötigt, stellen wir nun den gebräuchlichsten Sortier-Algorithmus vor, den so genannten Quicksort-Algorithmus. Auf ihn kann man kommen, seine Analyse soll hier mittels der Methoden der Wahrscheinlichkeitsrechnung erfolgen. Die Laufzeit-Analyse einer randomisierten Form wird dabei recht einfach erfolgen können.

Gegeben sei eine Menge  $S$  von  $n$  paarweise verschiedenen Zahlen. Sie soll der Größe nach sortiert werden. Zunächst wählt man ein Element  $y$  aus. Man kann zum Beispiel immer das erste Element nehmen, man kann es aber auch zufällig auswählen. Das ausgewählte Element  $y$  nennt man *Pivot-Element*. Man vergleicht dieses Pivot-Element nun mit allen anderen Zahlen und partitioniert die Zahlen in zwei Mengen; die eine umfasst alle Zahlen, die kleiner sind als  $y$  (genannt  $S_1$ ), die andere alle Zahlen, die größer als  $y$  sind (genannt  $S_2$ ). Man benötigt dazu  $n-1$  Vergleiche. Man gebe  $S_1$ ,  $y$  und  $S_2$  in dieser Reihenfolge aus.

Nun wende man den beschriebenen Algorithmus (Auswahl des Pivot-Elements, Partitionierung, Ausgabe) rekursiv auf  $S_1$  und  $S_2$  an (zumindest solange die Mengen mehr als ein Element besitzen). Wählt man  $y$  fest (zum Beispiel immer das erste Element), sprechen wir vom *deterministischen Quicksort-Algorithmus*, wählen wir  $y$  zufällig nach der Gleichverteilung auf  $S$ , so sprechen wir vom *randomisierten Quicksort-Algorithmus*

Das Maß für die so genannte *Laufzeit* des Algorithmus ist nun die Gesamtanzahl der benötigten Vergleiche von zwei Zahlen.  $X_n$  bezeichne im Folgenden die Anzahl der von Quicksort durchgeführten Vergleiche. Dies ist im Fall des randomisierten Quicksort-Algorithmus eine Zufallsvariable, daher die Bezeichnung  $X_n$ .

Wir wollen zunächst diskutieren, in welchen Situationen der Algorithmus besonders schnell bzw. besonders langsam den Sortier-Vorgang erledigt. Man wähle zum Beispiel als Pivot-Element immer das erste Element. Jetzt ist die Laufzeit von Quicksort im schlimmsten Fall von der Größenordnung  $n^2/2$ , wenn die vorliegende Liste von Zahlen bereits sortiert ist. Denn dann benötigen wir im ersten Schritt  $n-1$  Vergleiche, im zweiten Schritt  $n-2$  Vergleiche und im letzten Schritt noch einen Vergleich, also insgesamt

$$(n-1) + (n-2) + (n-3) + \dots + 1 = n(n-1)/2$$

Vergleiche. Dieser *worst case* kann auch beim randomisierten Quicksort auftreten, nämlich genau dann, wenn der Zufall stets als Pivot-Element das kleinste oder das größte Element auswählt. Auf der anderen Seite ist Quicksort besonders schnell, falls stets der mittlere Wert, der *Median*, der zu sortierenden Mengen gewählt wird (deterministisch oder zufällig). Nur kommt hier schon die Schwäche der deterministischen Version zum Vorschein: soll das Pivot-Element der Median sein, so müssen wir das mittlere Element finden. Es ist das mittlere Element einer sortierten Sequenz von Elementen.

Es gibt zwar einen randomisierten Algorithmus zur Bestimmung des Medians, der weniger Vergleiche benötigt als jeder theoretisch mögliche deterministische Algorithmus (und dieser kann an anderer Stelle einmal besprochen werden), aber die Laufzeit (die bei  $cn$  Vergleichen mit einer Konstanten  $c$  liegt) müsste bei jeder Auswahl des Pivot-Elements hinzugerechnet werden.

Wenn nun aber – rein theoretisch – die randomisierte Variante von Quicksort immer als Pivot-Element den Median liefert, so sind die beiden Mengen  $S_1$  und  $S_2$  jeweils etwa gleich groß. Dann benötigt man jeweils etwa  $n$  Vergleiche, aber die Anzahl der Iterationen (hier Halbierungen) ist etwa  $\log_2(n)$ . Dies ist einfach einzusehen:

$\log_2(n)$  ist diejenige Zahl, mit der man 2 potenzieren muss um  $n$  zu erhalten. Insgesamt ist also ein günstiger Verlauf des Quicksort-Algorithmus durch  $X_n = n \log_2(n)$  gegeben. Wir fassen zusammen: Die Anzahl der Vergleiche beim Quicksort-Algorithmus liegt

$$\text{zwischen } n \log_2(n) \text{ und } n^2/2.$$

Für die deterministische Variante von Quicksort kann man sich nun aufmachen und die *durchschnittliche Laufzeit* ausrechnen. Wir betrachten dazu den Ansatz, führen die Rechnung aber nicht durch. Ist das Pivot-Element das  $i$ -t größte Element der Liste, so haben die beiden rekursiv zu sortierenden Teilfelder  $S_1$  und  $S_2$  die Größe  $i-1$  und  $n-i$ . Dann beschreibt

$$X_n = n - 1 + \frac{1}{n} \sum_{i=1}^n (X_{i-1} + X_{n-i})$$

für  $n \geq 1$  und bei Setzung  $X_0 := 0$  das durchschnittliche Verhalten des Quicksorts. Man sieht:

$$X_n = n - 1 + \frac{2}{n} \sum_{i=1}^n X_{i-1}.$$

Eine solche Rekursionsgleichung kann man lösen. Eine verbreitete Methode ist die der erzeugenden Funktion. Wir lassen dies hier weg und teilen mit:

$$X_n = 2(n+1) \sum_{i=1}^n \frac{1}{i} - 2n - 4.$$

Wenn wir nun verwenden, dass

$$\frac{1}{j} \leq \frac{1}{x} \text{ für } x \in [j-1, j]$$

und daher

$$1 + \sum_{j=2}^n \frac{1}{j} \leq 1 + \int_1^n \frac{1}{x} dx = 1 + \log n,$$

so sehen wir, dass der Quicksort-Algorithmus in seiner deterministischen Form zum Sortieren von  $n$  Zahlen im Durchschnitt

$$n \cdot \log n + O(n)$$

Vergleiche benötigt. Wieder meint  $O(n)$  einen Ausdruck, der durch  $c n$  mit  $c$  konstant nach oben abgeschätzt werden kann. Der Algorithmus ist also im Mittel wesentlich schneller als im worst case ( $n^2$ ). Man bedenke aber, dass zum Beispiel bei Eingabe-Listen von überwiegend vorsortierten Daten die Aussage über die durchschnittliche Laufzeit wenig aussagekräftig ist! Im randomisierten Fall wird man nun aber das folgende Resultat beobachten können: *wählt man das Pivot-Element nicht deterministisch sondern zufällig, so ist bei allen(!) Eingaben die Laufzeit mit hoher Wahrscheinlichkeit gutartig.*

Zur Klärung dieser Aussage machen wir uns nun auf, im randomisierten Fall den Erwartungswert von  $X_n$  zu bestimmen.

**Theorem:** *Beim randomisierten Quicksort-Algorithmus beträgt die mittlere Laufzeit*

$$E(X_n) = 2(n+1) \left(1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n}\right) - 4n,$$

also

$$E(X_n) = 2n \log n + O(n)$$

Bevor wir den Beweis geben, hier eine schnelle Vergleichs-Diskussion: da

$$2 \log(n) \text{ ungefähr } 1.39 \log_2(n)$$

ist, muss man also im Mittel nur 39 Prozent mehr Vergleiche vornehmen als im günstigsten Fall!

**Beweis:** Eine der fundamentalen Ideen in der Stochastik zur Bestimmung des Erwartungswertes einer Zufallsvariablen ist die Zerlegung der Zufallsvariable in eine Summe einfacherer Zufallsvariablen. Der Erwartungswert ist dann die Summe der Erwartungswerte der einfacheren Zufallsvariablen (Linearität des Erwartungswertes). Für die Menge  $S$  schreiben wir nun

$$S = \{s_{(1)}, s_{(2)}, \dots, s_{(n)}\}$$

wobei dies die korrekte Ordnung von  $S$  sei, also  $s_{(i)} < s_{(j)}$  für  $i < j$ . Somit bezeichnet  $s_{(i)}$  das  $i$ -t kleinste Element. Es sei  $1_A$  die Indikatorfunktion auf der Menge  $A$ , also die Funktion, die einem  $x$  den Wert 1 zuordnet, wenn  $x \in A$ , und sonst den Wert 0 zuordnet. Dann gilt

$$X_n = \sum_{1 \leq i < j \leq n} 1_{A_{ij}},$$

wobei  $A_{ij}$  das Ereignis bezeichnet, dass es im Verlauf des Algorithmus zum Vergleich von  $s_{(i)}$  und  $s_{(j)}$  kommt. Dann ist

$$E(X_n) = \sum_{1 \leq i < j \leq n} E(1_{A_{ij}})$$

Der Erwartungswert eines Indikators  $1_A$  ist

$$E(1_A) = 1 \cdot P(A) + 0 \cdot P(A^c) = P(A)$$

Also müssen wir  $P(A_{ij})$  berechnen. Wir behaupten:

$$P(A_{ij}) = \frac{2}{j-i+1} \text{ für alle } 1 \leq i < j \leq n.$$

Dies sieht man so: das Ereignis  $A_{ij}$  tritt genau dann ein, wenn entweder  $s_{(i)}$  oder  $s_{(j)}$  das erste der Elemente  $s_{(i)}, \dots, s_{(j)}$  ist, das als Pivot-Element  $y$  herangezogen wird. Denn falls ein anderes Element dieser Teilliste gewählt würde - nennen wir es  $s_{(k)}$  mit  $i < k < j$  - dann würden  $s_{(i)}$  und  $s_{(j)}$  zwar je mit  $s_{(k)}$  verglichen werden, aber sie würden nicht miteinander verglichen werden, denn sie gehörten dann zu verschiedenen Teillisten  $S_1$  und  $S_2$ , die beim rekursiven Verfahren ja im Anschluss getrennt behandelt werden. Nun hat jedes der Elemente  $s_{(i)}, \dots, s_{(j)}$  dieselbe Wahrschein-

lichkeit, als erstes als Pivot-Element gewählt zu werden. Also folgt

$$P(A_{ij}) = \frac{2}{j-i+1}.$$

Wir müssen also – mit etwas Geduld –

$$2 \sum_{j=2}^n \sum_{i=1}^{j-1} \frac{1}{j-i+1}$$

ausrechnen bzw. umformen. Für jedes  $j$  ist die innere Summe

$$\frac{1}{j} + \frac{1}{j-1} + \dots + \frac{1}{2}, \text{ hat also die Darstellung } \sum_{i=2}^j \frac{1}{i}$$

Wir erhalten

$$E(X_n) = 2 \sum_{j=2}^n \sum_{i=2}^j \frac{1}{i}$$

Nun vertauschen wir die Summationsreihenfolge:

$$E(X_n) = 2 \sum_{i=2}^n \sum_{j=i}^n \frac{1}{i}$$

Dies sieht man mit Hilfe von Abb. 1 ein:

j \ i	2	3	..	n
2	•			
3	•	•		
..				
n	•	•	..	•

Abb. 1: Doppelsummenschema

Nun ist aber die innere Summe

$$(n-i+1) \frac{1}{i} = (n+1) \frac{1}{i} - 1$$

Damit erhalten wir

$$E(X_n) = 2(n+1) \sum_{i=2}^n \frac{1}{i} - 2(n-1)$$

Und nun ist die rechte Seite gleich

$$2(n+1) \sum_{i=1}^n \frac{1}{i} - 2(n+1) - 2(n-1) = 2(n+1) \sum_{i=1}^n \frac{1}{i} - 4n$$

Die Asymptotik folgt nun wieder aus

$$1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n} = \log n + O(1)$$

Damit ist der Satz bewiesen. ■

Die mittlere Laufzeit ist etwas besser als die ange-deutete im deterministischen Fall. Eine noch bes-sere Aussage über die Qualität des randomisierten Quicksorts ist, dass  $X_n$  mit hoher Wahrschein-lichkeit in der Nähe seines Erwartungswertes liegt.

Eine solche Aussage wurde in Rösler (1991) be-wiesen und liest sich so:

Für beliebige  $\lambda, \varepsilon > 0$  existiert eine Konstante  $c(\lambda, \varepsilon) > 0$ , so dass

$$P(X_n > (1+\varepsilon)E(X_n)) \leq c(\lambda, \varepsilon)n^{-2\lambda\varepsilon}$$

wobei  $n$  eine natürliche Zahl ist. Dies ist die Rechtfertigung, dass für alle(!) Eingaben die Lauf-zeit mit hoher Wahrscheinlichkeit gutartig ist.

Nimmt man an, dass alle möglichen Eingaben, also alle Permutationen der Eingabe, gleich wahr-scheinlich sind und interessiert sich dann für den Erwartungswert der Laufzeit, so kann man mittels einer leichten Umformulierung des obigen Bewei-ses zeigen, dass auch hier die durchschnittliche Laufzeit in der Größenordnung  $n \cdot \log n$  liegt. Wir führen dies nicht aus.

Es bleibt abschließend darauf hinzuweisen, dass es viele Untersuchungen zu Sortieralgorithmen gibt. Hoare hat 1961 den hier diskutierten Quicksort eingeführt, siehe Hoare (1961). Andere Sortier-algorithmen tragen den Namen *Heapsort* oder *Mergesort* und werden unter anderem in Krengel (2000) diskutiert. Ein gut ausgearbeitetes Skriptum ist Koenig (2004).

## Literatur

- Aldous, D. und Diaconis, P. (1999): Longest in-creasing subsequences: from patience sorting to the Baik-Deift-Johansson theorem. Bull. Amer. Math. Soc. (N.S.) 36, 413–432.
- Eichelsbacher, P. und Löwe, M. (2003): Geduld und Zufall. Stochastik in der Schule 23(2), 2–6.
- Fredman, M.L. (1975): On computing the length of the longest increasing subsequence. Discrete Mathematics 11, 29–35.
- Hoare, C.A.R. (1961): Computer Journal 5, 10–15.
- König, W. (2004): Stochastische Algorithmen. Universität Leipzig.
- Krengel, U. (2000): Einführung in die Wahrschein-lichkeitstheorie und Statistik (5. Aufl.). Wies-baden: Vieweg.
- Rösler, U. (1991): A limit theorem for Quicksort. Theor. Inf. Appl. 25, 85–100.

## Anschrift des Verfassers

Peter Eichelsbacher  
Fakultät für Mathematik  
Ruhr-Universität Bochum, NA 3 / 68  
44780 Bochum  
peter.eichelsbacher@ruhr-uni-bochum.de