

Kapitel 2 Suchverfahren

Gegeben sei eine Menge von Objekten, die durch (eindeutige) Schlüssel charakterisiert sind. Aufgabe von Suchverfahren ist die Suche bestimmter Objekte anhand ihres Schlüssels.

Bisher: zwei Methoden der Speicherung solcher Objekte

- sequentielle Speicherung \rightsquigarrow sortiertes Array (binäre Suche)
- verkettete Speicherung \rightsquigarrow lineare Liste

Zeitaufwand im schlechtesten Fall für eine Menge von n Elementen;

Operation	sequentiell gespeichert (sortiertes ARRAY)	verkettet gespeichert (lineare Liste)
Suche Objekt mit gegebenem Schlüssel	$O(\log n)$	$O(n)$
Einfügen an bekannter Stelle	$O(n)$	$O(1)$
Entfernen an bekannter Stelle	$O(n)$	$O(1)$

Im folgenden werden wir Verfahren behandeln, die sowohl die Suche als auch das Einfügen und Entfernen "effizient" unterstützen (bessere Laufzeitkomplexität als $O(n)$).

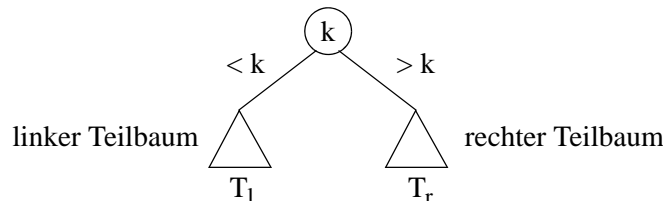
2.1 Binäre Suchbäume

Ziel: Die Operationen *Suchen*, *Einfügen* und *Entfernen* sollen alle in $O(\log n)$ Zeit durchgeführt werden.

Ansatz: Organisation der Objektmenge als Knoten eines binären Baumes.

Definition:

Ein binärer Baum heißt **binärer Suchbaum**, wenn für jeden seiner Knoten die **Suchbaumeigenschaft** gilt, d.h. alle Schlüssel im linken Teilbaum sind kleiner, alle Schlüssel im rechten Teilbaum sind größer als der Schlüssel im Knoten:

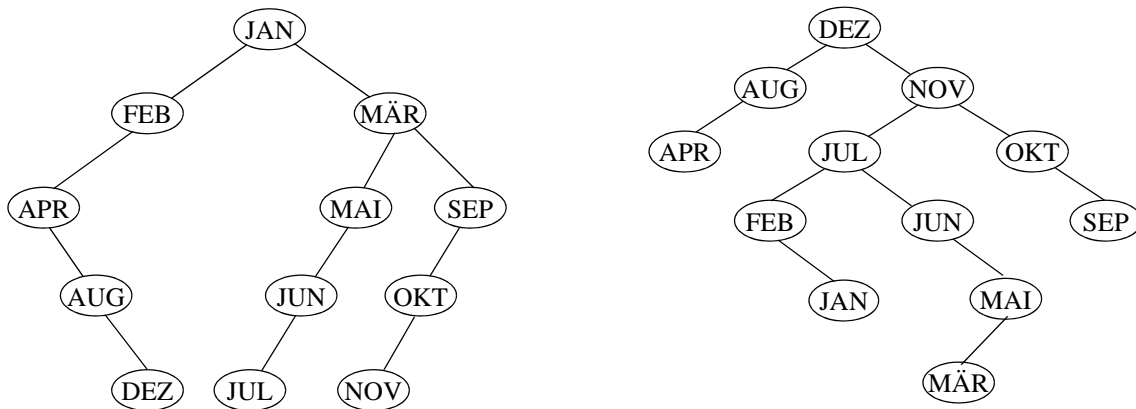


Anmerkungen:

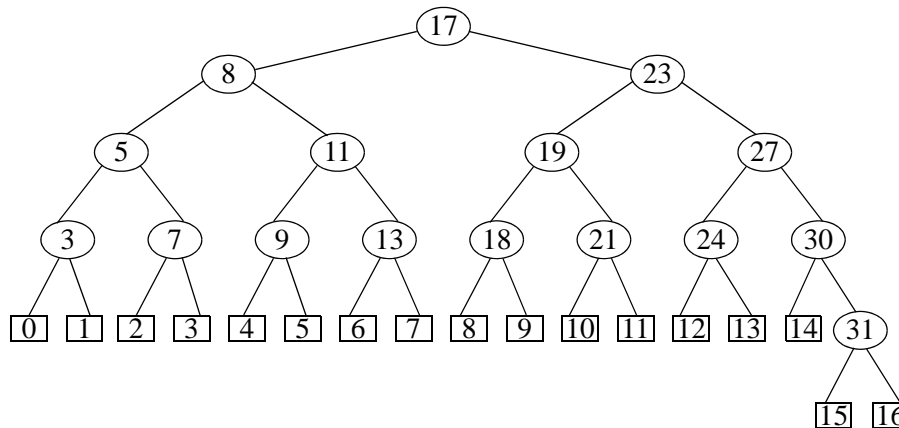
- Zur Organisation einer gegebenen Menge von n Schlüsseln gibt es eine große Anzahl unterschiedlicher binärer Suchbäume.
- Der *inorder*-Durchlauf eines binären Suchbaumes generiert die (aufsteigend) sortierte Folge der gespeicherten Schlüssel.

Beispiele:

Zwei verschiedene binäre Suchbäume über den Monatsnamen:



Der Entscheidungsbaum zur binären Suche ist ein binärer Suchbaum:



2.1.1 Allgemeine binäre Suchbäume

```

class BinaryNode
{
    int key;
    BinaryNode left;
    BinaryNode right;

    BinaryNode(int k) { key=k; left = null; right = null; }
}

public class BinarySearchTree
{
    BinaryNode root;
    public BinarySearchTree () { root = null; }
    // ... Implementierung der Methoden insert, find, delete,...
}

```

Die Methode *insert* (*int* *x*) der Klasse *BinarySearchTree* fügt einen gegebenen Schlüssel *x* ein, falls dieser noch nicht im Baum enthalten ist:

```
// Methoden der Klasse BinarySearchTree
public void insert (int x) throws Exception
{ root = insert (x, root);}

protected BinaryNode insert (int x, BinaryNode t) throws Exception
{   if ( t == null) t = new BinaryNode (x);
    else if (x < t.key)
        t.left = insert (x, t.left);
    else if (x > t.key)
        t.right = insert (x, t.right);
    else
        throw new Exception ("Inserting twice");
    return t;
}
```

Die überladene Methode *insert* (*x*, *t*) liefert eine Referenz auf die Wurzel des Teilbaums zurück, in den *x* eingefügt wurde.

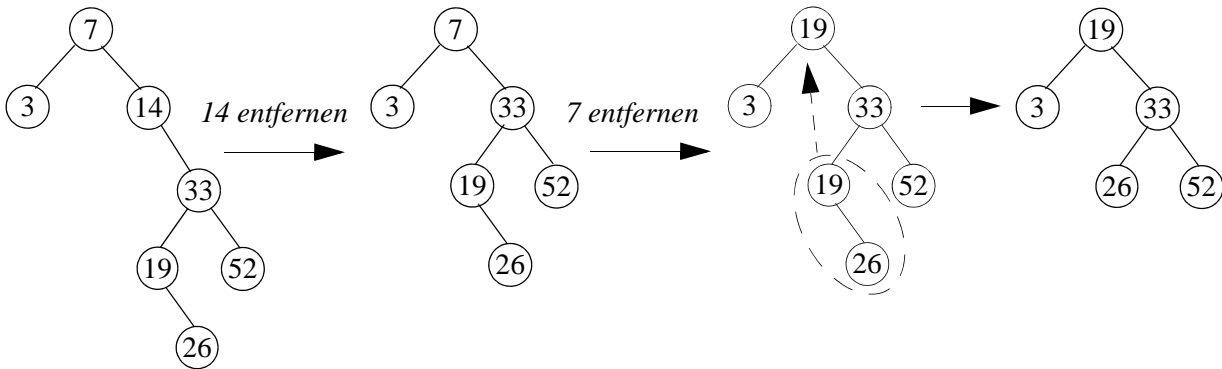
Die folgende Methode *find* (*x*) sucht in einem binären Suchbaum den Schlüssel *x* und liefert diesen zurück, falls er gefunden wurde. Andernfalls wird eine Ausnahmebehandlung durchgeführt.

```
// Methoden der Klasse BinarySearchTree
public int find (int x) throws Exception
{ return find (x, root).key;}

protected BinaryNode find (int x, BinaryNode t) throws Exception
{ while ( t != null)
    { if ( x < t.key) t = t.left;
      else if ( x > t.key) t = t.right;
      else return t;
    }
    throw new Exception ("key not found");
}
```

Das Entfernen eines Schlüssels ist etwas komplexer, da auch Schlüssel in inneren Knoten des Baumes betroffen sein können und die Suchbaumstruktur aufrecht erhalten werden muß.

Hierzu zunächst ein Beispiel:

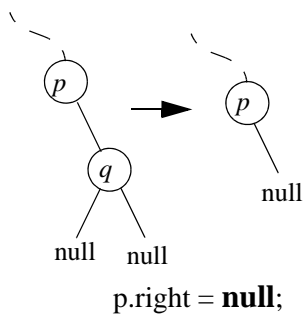


Allgemein treten zwei grundsätzlich verschiedene Situationen auf:

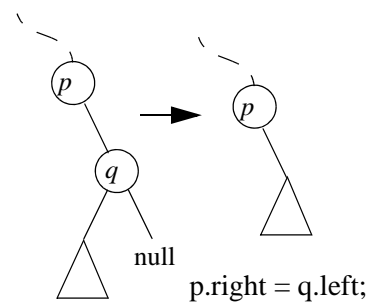
Fall 1: der Knoten q mit dem zu entfernenden Schlüssel besitzt höchstens einen Sohn (Blatt oder Halbblatt)

Fall 2: der Knoten q mit dem zu entfernenden Schlüssel besitzt zwei Söhne (innerer Knoten)

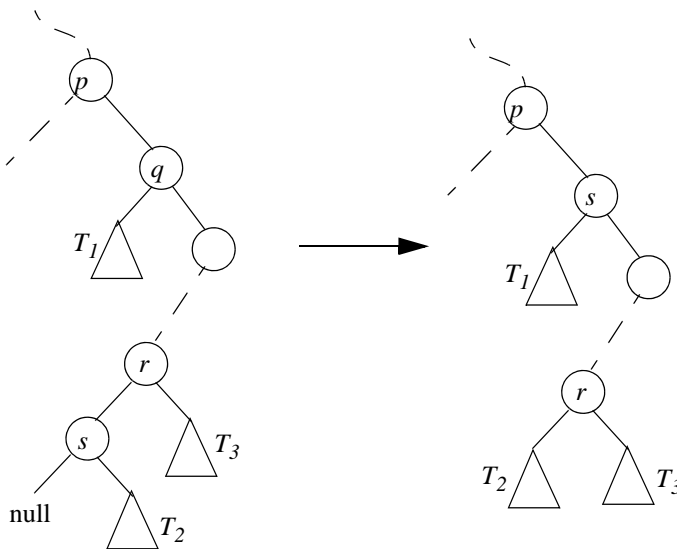
Fall 1: a) der Knoten q besitzt keinen Sohn



b) der Knoten q besitzt einen linken Sohn (rechts \rightarrow symmetrisch)



Fall 2:



Programmcode:
s. Methode 'deleteMin'

Die folgende Methode *delete(x)* entfernt einen Schlüssel *x* aus einem binären Suchbaum. Hierbei wird eine Hilfsmethode *findMin* verwendet, die den Knoten des kleinsten Schlüssels des rechten Teilbaumes (eines inneren Knotens) bestimmt und eine weitere Hilfsmethode *deleteMin*, die diesen Knoten entfernt:

```
// Methoden der Klasse BinarySearchTree
public void delete(int x) throws Exception
{ root = delete (x, root);}

protected BinaryNode delete (int x, BinaryNode t) throws Exception
{ if (t == null)
    throw new Exception (“x does not exist (delete)”);
  if (x < t.key) t.left = delete (x, t.left);
  else if (x > t.key) t.right = delete (x, t.right);
  else if (t.left != null && t.right != null)// x is in inner node t
    { t.key = findMin (t.right).key;
      t.right = deleteMin (t.right);
    }
  else // x is in leaf or in semi-leaf node, reroot t
    t = (t.left != null) ? t.left : t.right;
  return t;
}

protected BinaryNode findMin (BinaryNode t) throws Exception
{ if ( t == null)
    throw new Exception (“key not found (findMin)”);
  while ( t.left != null) t = t.left;
return t;
}

protected BinaryNode deleteMin (BinaryNode t) throws Exception
{ if (t == null)
    throw new Exception (“key not found (deleteMin)”);
  if (t.left != null)
    t.left = deleteMin (t.left);
  else
    t = t.right;
  return t;
}
```

Zum Verständnis: Der Fall des Entfernens aus einem inneren Knoten ($t.right \neq \text{null} \ \&\& \ t.left \neq \text{null}$) wird auf eine ‘Blatt-’ ($t.right == \text{null} \ \&\& \ t.left == \text{null}$) oder eine ‘Halbblatt-Situation’ ($t.right \neq \text{null} \ || \ t.left \neq \text{null}$) zurückgeführt, indem der zu löschende Schlüssel durch den kleinsten der größeren Schlüssel ersetzt wird. Dieser Schlüssel befindet sich stets in einem Blatt oder einem Halbblatt, das daraufhin aus dem Baum entfernt wird.

Laufzeitanalyse der Algorithmen *insert*, *find* und *delete*

Alle drei Methoden sind auf einen einzigen, bei der Wurzel beginnenden Pfad des Suchbaumes beschränkt.

☛ der **maximale Aufwand** ist damit $O(h)$, wobei h die Höhe des Baumes ist.

Für die Höhe binärer Bäume mit n Knoten gilt:

- Die **maximale Höhe** eines binären Baumes mit n Knoten ist n . (☛ Lineare Liste)
- Seine **minimale Höhe** ist $\lceil \log_2(n+1) \rceil$

Begründung: Für eine gegebene Anzahl n von Knoten haben die sogenannten vollständig ausgeglichenen binären Bäume minimale Höhe. In einem vollständig ausgeglichenen binären Baum müssen alle Levels bis auf das unterste vollständig besetzt sein. Die maximale Anzahl n von Knoten in einem vollständig ausgeglichenen binären Baum der Höhe h ist:

$$n = \sum_{i=0}^{h-1} 2^i = 2^h - 1$$

$$\text{☛ } h_{\min} = \lceil \log_2(n+1) \rceil.$$

Bemerkung: Es gilt: $\lceil \log_2(n+1) \rceil = \lfloor \log_2(n) \rfloor + 1 \quad \forall n \in \mathbb{N}$.

Eine aufwendige Durchschnittsanalyse ergibt unter den beiden Annahmen:

- der Baum ist nur durch Einfügungen entstanden und
- alle möglichen Permutationen der Eingabereihenfolge sind gleichwahrscheinlich

einen mittleren Wert $h_{\text{Ø}} = 2 \cdot \ln 2 \cdot \log n \approx 1,386 \cdot \log n$.

Der **durchschnittliche Zeitbedarf** für Suchen, Einfügen und Entfernen ist damit $O(\log n)$.

Kritikpunkt am naiven Algorithmus zum Aufbau binärer Suchbäume und damit an der Klasse so erzeugter binärer Suchbäume:

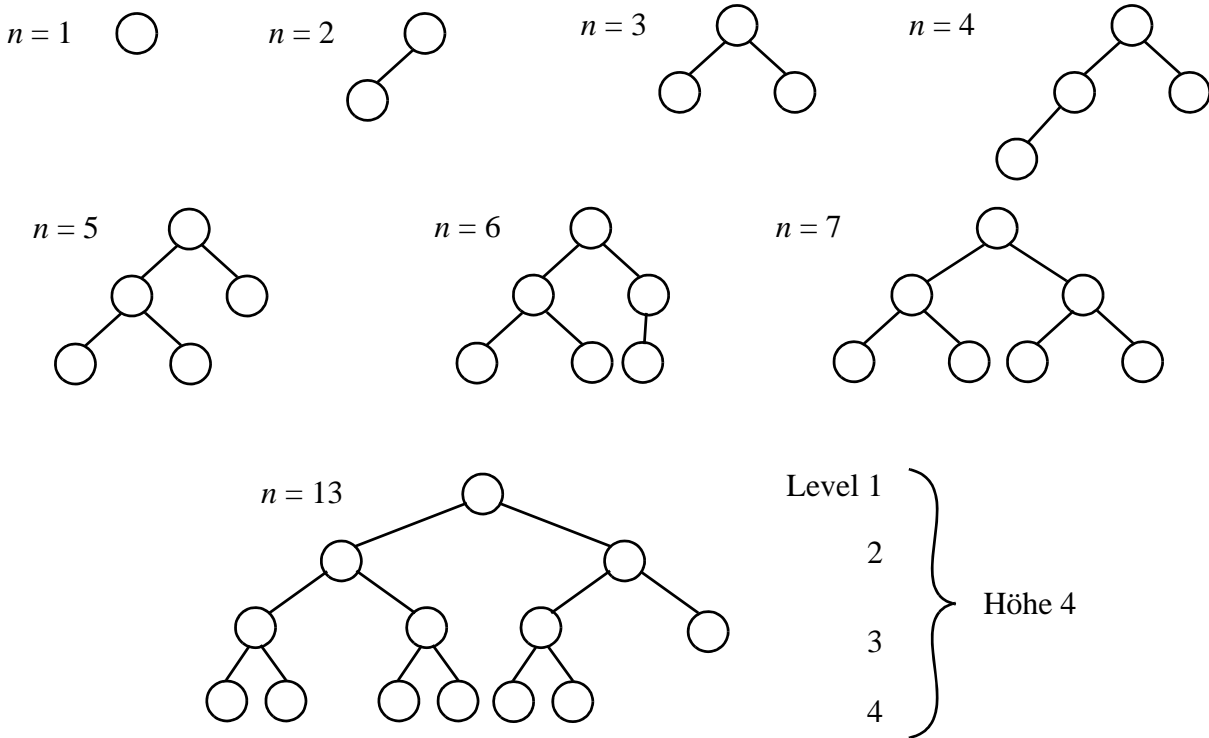
☛ im **worst-case** ist der Aufwand aller drei Operationen $O(n)$.

2.1.2 Vollständig ausgeglichene binäre Suchbäume

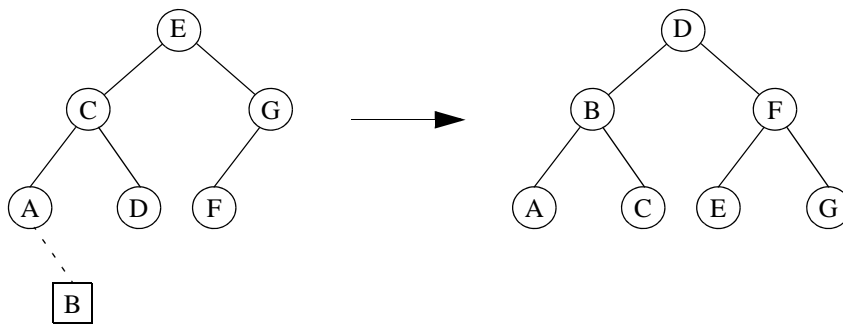
Die minimale Höhe $\lceil \log_2(n+1) \rceil$ unter allen binären Suchbäumen besitzen die **vollständig ausgeglichenen binären Suchbäume**, d.h. binäre Suchbäume, bei denen alle Levels bis auf das unterste vollständig besetzt sind.

☛ **optimaler Zeitaufwand für Suchoperationen**

Vollständig ausgeglichene binäre Bäume für verschiedene Knotenzahlen n :



Beispiel: Einfügen von Schlüssel 'B' in einen bestehenden Baum



Problem: Der Baum muß beim Einfügen von B **vollständig reorganisiert** werden.

☛ **Einfügezeit im schlechtesten Fall:** $O(n)$.

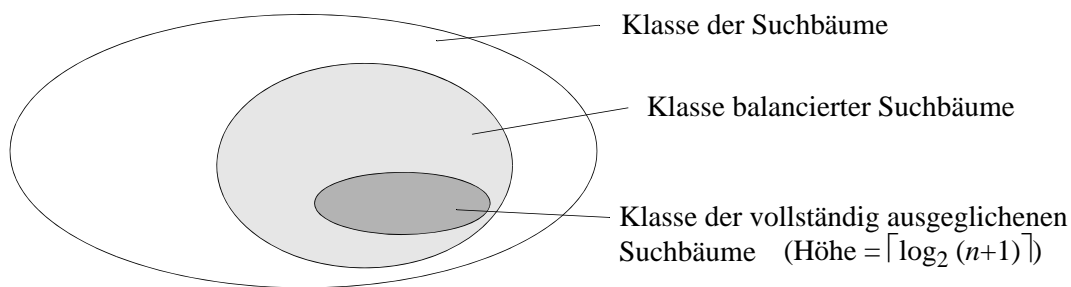
☛ Auswahl einer Kompromißlösung mit den Eigenschaften:

- die Höhe des Baumes ist im schlechtesten Fall $O(\log n)$.
- Reorganisationen bleiben auf den Suchpfad zum einzufügenden bzw. zu entfernenden Schlüssel beschränkt und sind damit im schlechtesten Fall in $O(\log n)$ Zeit ausführbar.

Definition:

Eine Klasse von Suchbäumen heißt **balanciert**, falls:

- $h_{max} = O(\log n)$
- die Operationen *Suchen*, *Einfügen* und *Entfernen* sind auf einen Pfad von der Wurzel zu einem Blatt beschränkt und benötigen damit im schlechtesten Fall $O(\log n)$ Zeit.

**2.1.3 AVL-Bäume**

(Adelson-Velskij und Landis (1962))

AVL-Bäume sind ein Beispiel für eine Klasse balancierter Suchbäume.

Definition:

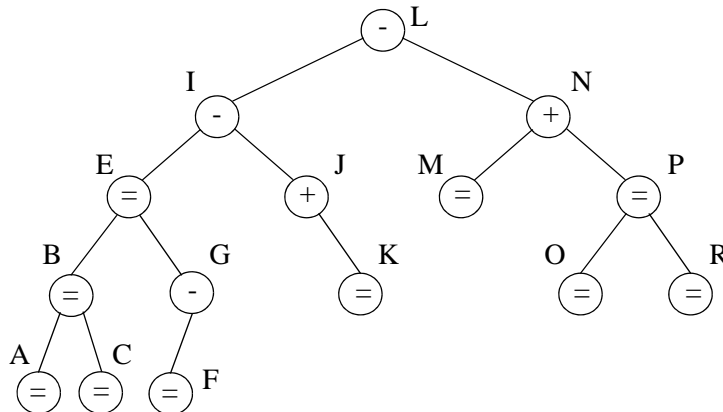
Ein binärer Suchbaum heißt **AVL-Baum**, falls für die beiden Teilbäume T_r und T_l der Wurzel gilt:

- $|h(T_r) - h(T_l)| \leq 1$
- T_r und T_l sind ihrerseits AVL-Bäume. (*rekursive Definition*)

Der Wert $h(T_r) - h(T_l)$ wird als **Balancefaktor** (BF) eines Knotens bezeichnet. Er kann in einem AVL-Baum nur die Werte -1, 0 oder 1 (dargestellt durch -, = und +) annehmen.

Mögliche Strukturverletzungen durch Einfügungen bzw. Entfernungen von Schlüsseln erfordern **Rebalancierungsoperationen**.

Beispiel für einen AVL-Baum:

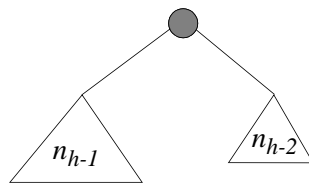


Behauptung: Die minimale Höhe $h_{min}(n)$ eines AVL-Baumes mit n Schlüsseln ist $\lceil \log_2(n + 1) \rceil$. Dies folgt aus der Tatsache, daß ein AVL-Baum minimaler Höhe einem vollständig ausgeglichenen binären Suchbaum entspricht.

Behauptung: Die maximale Höhe $h_{max}(n)$ eines AVL-Baumes mit n Schlüsseln ist $O(\log n)$.

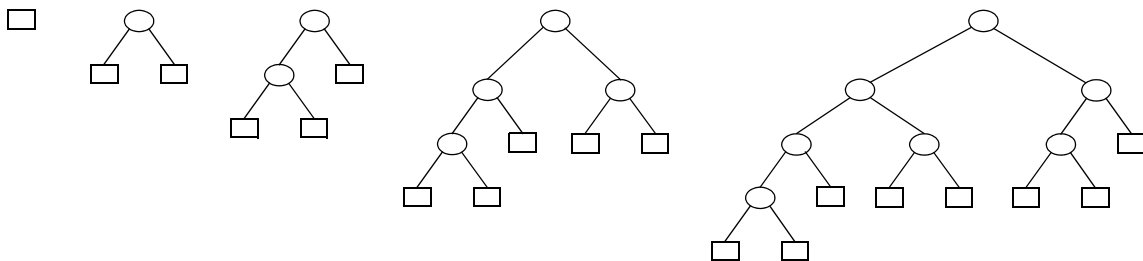
Beweis: Die maximale Höhe wird realisiert von sogenannten *minimalen AVL-Bäumen*. Dies sind AVL-Bäume, die für eine gegebene Höhe h die minimale Anzahl von Schlüsseln abspeichern.

Minimale AVL-Bäume haben bis auf Symmetrie die folgende Gestalt:



Sei n_h die minimale Anzahl von Schlüsseln in einem AVL-Baum der Höhe h . Dann gilt: $n_h = n_{h-1} + n_{h-2} + 1$ für $h \geq 2$ und $n_0 = 0, n_1 = 1$.

Die minimalen AVL-Bäume der Höhen $h = 0, 1, 2, 3, 4$ haben bis auf Symmetrie folgende Gestalt:



Die Rekursionsgleichung für n_h erinnert an die Definition der *Fibonacci-Zahlen*:

$$fib(n) = \begin{cases} n & \text{für } n \leq 1 \\ fib(n-1) + fib(n-2) & \text{für } n > 1 \end{cases}$$

h	0	1	2	3	4	5	6
n_h	0	1	2	4	7	12	20
$fib(h)$	0	1	1	2	3	5	8

Hypothese: $n_h = fib(h + 2) - 1$

Beweis: Induktion über h

Induktionsanfang: $n_0 = fib(2) - 1 = 1 - 1 = 0 \checkmark$

Induktionsschluß: $n_{h+1} = n_h + n_{h-1} + 1$
 $= 1 + fib(h + 2) - 1 + fib(h + 1) - 1$
 $= fib(h + 3) - 1$

Hilfssatz: $fib(n) = \frac{1}{\sqrt{5}} \cdot (\varphi_1^n - \varphi_2^n)$ mit $\varphi_1 = \frac{1 + \sqrt{5}}{2}$, $\varphi_2 = \frac{1 - \sqrt{5}}{2} \approx -0,618$.

Der Beweis erfolgt durch vollständige Induktion. (wird hier übergangen)

Für jede beliebige Schlüsselanzahl $n \in \mathbb{N}$ gibt es ein eindeutiges $h_{max}(n)$ mit:

$$n_{h_{max}(n)} \leq n < n_{h_{max}(n) + 1}.$$

Mit obiger Hypothese folgt hieraus: $n + 1 \geq fib(h_{max}(n) + 2)$.

Durch Einsetzen erhalten wir:

$$n + 1 \geq \frac{1}{\sqrt{5}} \cdot \left(\varphi_1^{h_{max}(n) + 2} - \underbrace{\varphi_2^{h_{max}(n) + 2}}_{< 0,62 < \sqrt{5}/2} \right) \geq \frac{1}{\sqrt{5}} \cdot \varphi_1^{h_{max}(n) + 2} - \frac{1}{2}$$

Und damit: $\frac{1}{\sqrt{5}} \cdot \varphi_1^{h_{max}(n) + 2} \leq n + \frac{3}{2}$.

Durch Auflösen nach $h_{max}(n)$ ergibt sich: $\log_{\varphi_1}(\frac{1}{\sqrt{5}}) + h_{max}(n) + 2 \leq \log_{\varphi_1}(n + \frac{3}{2})$.

Und für $h_{max}(n)$: $h_{max}(n) \leq \log_{\varphi_1}(n + \frac{3}{2}) - (\log_{\varphi_1}(\frac{1}{\sqrt{5}}) + 2) \leq$

$$\leq \log_{\varphi_1}(n) + const = \log_{\varphi_1}(2) \cdot \log_2(n) + const$$

$$\left(\log_b(a) = \frac{\log_c(a)}{\log_c(b)} \right) \quad \leftarrow \frac{\ln 2}{\ln \varphi_1} \cdot \log_2(n) + const \approx 1,44 \cdot \log_2(n) + const$$

Für große Schlüsselanzahlen ist die Höhe eines AVL-Baumes somit um maximal 44% größer als die des vollständig ausgeglichenen binären Suchbaumes.

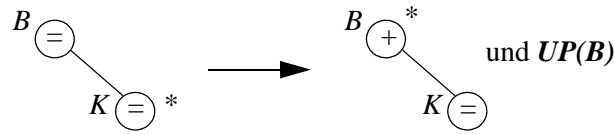
Also gilt die Behauptung: $h_{max}(n) = O(\log n)$.

Einfügen von Schlüsseln: $Insert(k)$

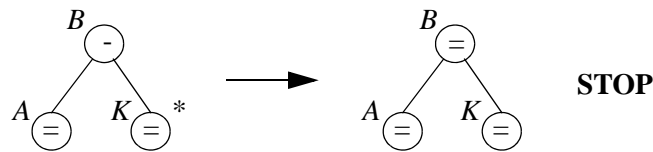
(In der folgenden Beschreibung sind symmetrische Fälle nicht dargestellt.)

Der Schlüssel k wird in einen neuen Sohn K des Knotens B eingefügt.

Fall 1: B ist ein Blatt

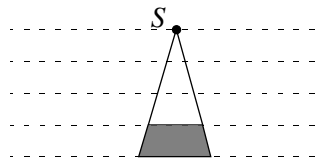


Fall 2: B hat einen linken Sohn



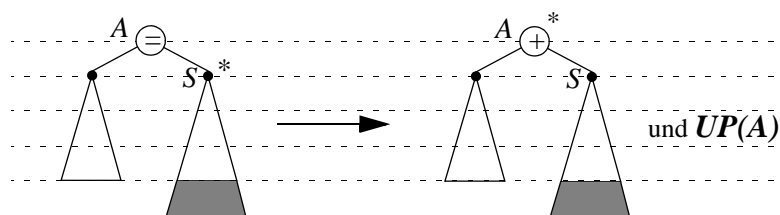
Die Methode $UP(S)$ wird aufgerufen für einen Knoten S , dessen Teilbaum in seiner Höhe um 1 gewachsen ist. S ist die Wurzel eines korrekten AVL-Baumes.

▮▮▮ mögliche Strukturverletzung durch einen zu hohen Teilbaum!



Fall 1: der Vater von S hat BF '='

1.1 der Vater von S ist nicht die Wurzel



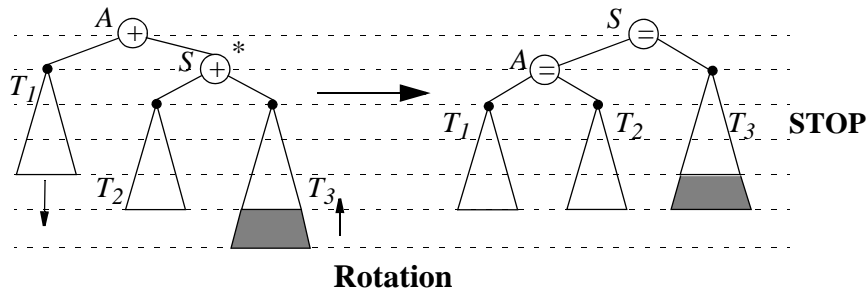
1.2 der Vater von S ist die Wurzel: \rightarrow dieselbe Transformation und **STOP**.

Fall 2: der Vater von S hat BF '+' oder '-' und S ist die Wurzel des kürzeren Teilbaumes.

▮▮▮ In beiden Fällen wird der BF im Vater zu '=' und **STOP**.

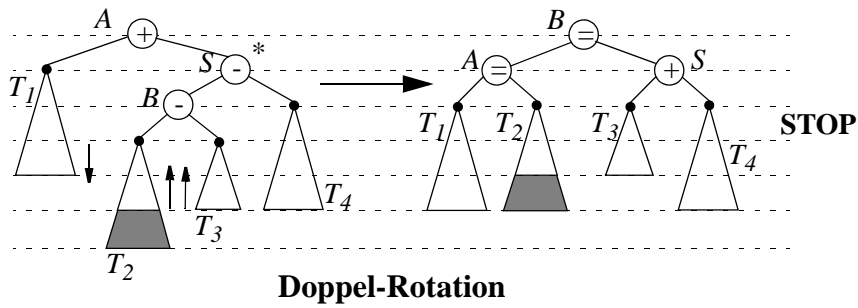
Fall 3: der Vater von S hat BF '+' oder '-' und S ist die Wurzel des höheren Teilbaumes.

3.1 der Vater von S hat BF '+' und S hat BF '+':



3.2 der Vater von S hat BF '+' und S hat BF '-'.

z.B.: B hat BF '-'.



der Fall B hat BF '+' wird analog gehandhabt.

3.3 der Vater von S hat BF '-' und S hat BF '-': → symmetrisch zu 3.1

3.4 der Vater von S hat BF '-' und S hat BF '+': → symmetrisch zu 3.2

Beim Einfügen genügt eine einzige Rotation bzw. Doppelrotation um eine Strukturverletzung zu beseitigen. Wir werden sehen, daß dies beim Entfernen nicht genügt.

Entfernen von Schlüsseln: Delete(k)

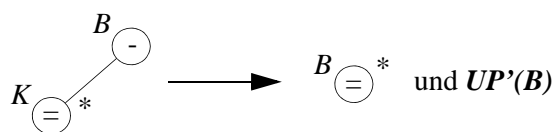
(In der folgenden Beschreibung sind symmetrische Fälle nicht dargestellt.)

Der Knoten K mit Schlüssel k wird entfernt.

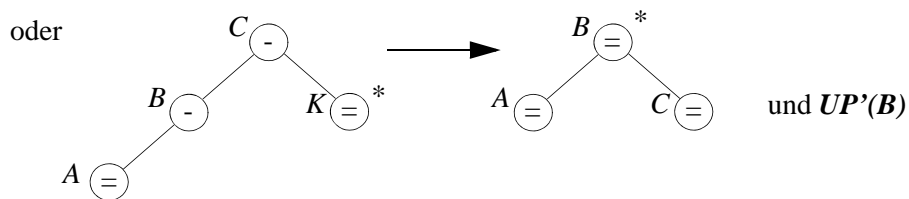
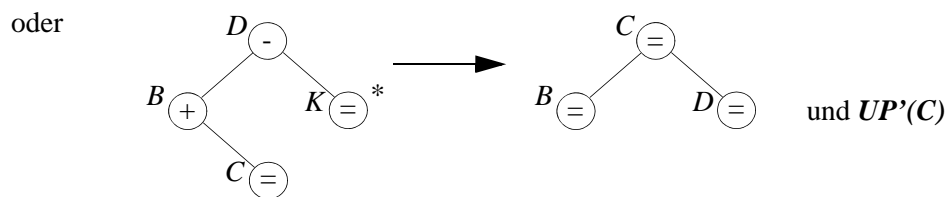
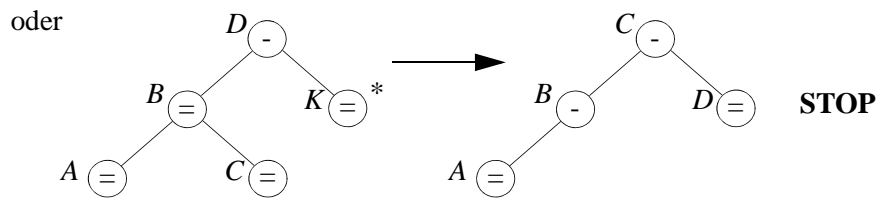
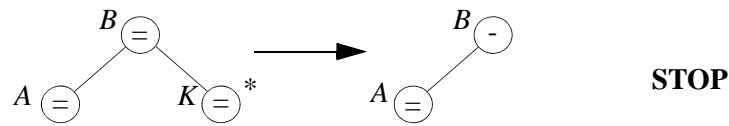
Fall 1: K hat höchstens einen Sohn

1.1 K ist ein Blatt

1.1.1 K hat keinen Bruder

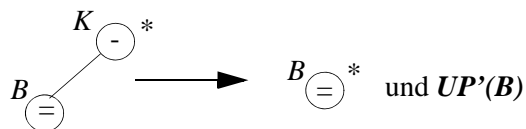


1.1.2 K hat einen Bruder



1.2 K hat genau einen Sohn

1.2.1 K hat einen linken Sohn



1.2.2 K hat einen rechten Sohn: \rightarrow symmetrisch

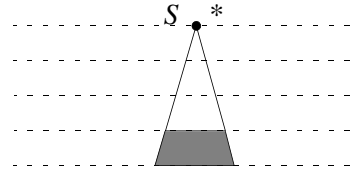
Fall 2: K ist ein innerer Knoten (hat zwei Söhne)

Man bestimme in dem AVL-Baum den **kleinsten** Schlüssel s , der **größer als k** ist. s ist in einem Halbblatt S . Ersetze k durch s und entferne den Schlüssel s .

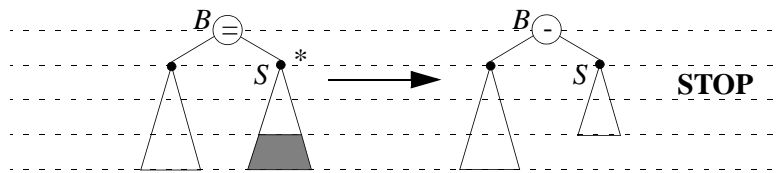
☛ damit haben wir Fall 2 auf Fall 1 zurückgeführt.

Die **Methode UP'(S)** wird aufgerufen für einen Knoten S , dessen Teilbaum in seiner Höhe um 1 reduziert ist. Der Teilbaum mit Wurzel S ist ein korrekter AVL-Baum.

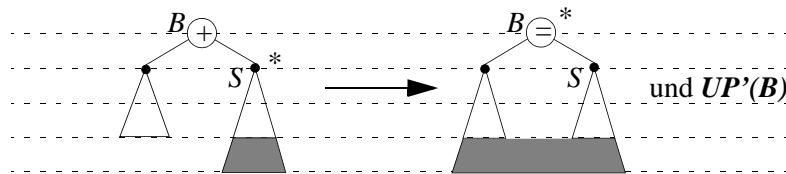
→ **mögliche Strukturverletzung** durch einen zu niedrigen Teilbaum!



Fall 1: der Vater von S hat BF '='.



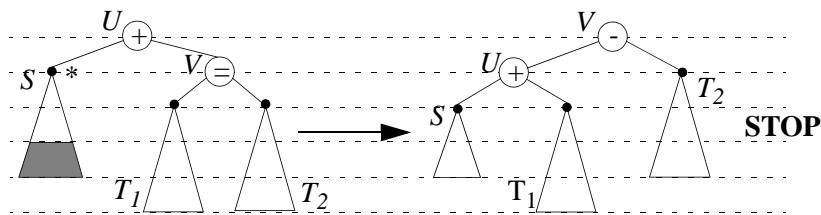
Fall 2: der Vater von S hat BF '+' oder '-' und S ist die Wurzel des höheren Teilbaums.



Fall 3: der Vater von S hat BF '+' oder '-' und S ist die Wurzel des kürzeren Teilbaums.

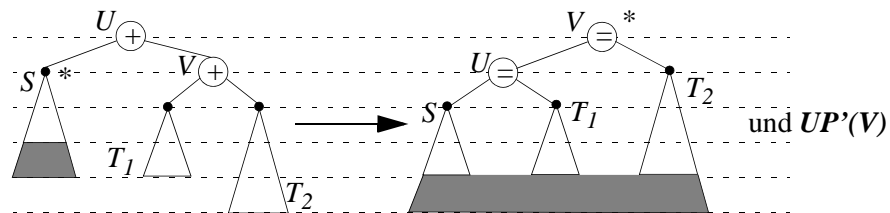
3.1 der Vater von S hat BF '+'

3.1.1 der Bruder von S hat BF '='



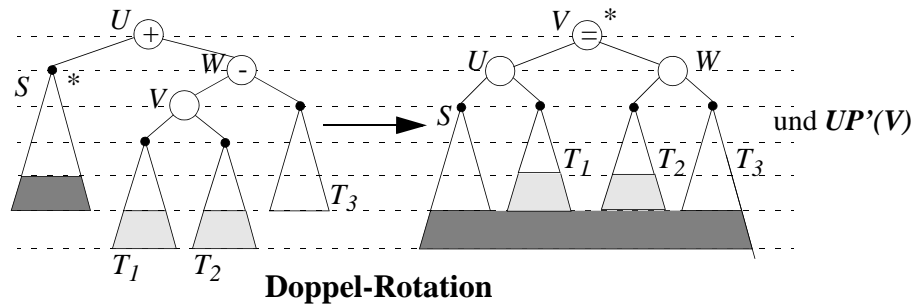
Rotation

3.1.2 der Bruder von S hat BF '+'



Rotation

3.1.3 der Bruder von S hat $BF' -'$.



Mindestens einer der beiden Bäume T_1 und T_2 hat die durch den hell schraffierten Bereich angegebene Höhe.

3.2 der Vater von S hat $BF' -'$: \implies symmetrisch zu 3.1.

Fall 4: S ist die Wurzel. \implies **STOP**

Im Falle von Entferne-Operationen wird eine mögliche Strukturverletzung also nicht notwendigerweise durch eine einzige Rotation bzw. Doppelrotation beseitigt. Im schlechtesten Fall muß auf dem Suchpfad bottom-up vom zu entfernenden Schlüssel bis zur Wurzel auf jedem Level eine Rotation bzw. Doppelrotation durchgeführt werden.

Korollar: Die **AVL-Bäume** bilden eine Klasse **balancierter Bäume**.

2.2 B-Bäume

Sei n die Anzahl der Objekte und damit der Datensätze. Wir nehmen nun an, daß das Datenvolumen zu groß ist, um im Hauptspeicher gehalten zu werden, z.B. $n = 10^6$.

\implies Datensätze auf externen Speicher auslagern, z.B. Plattenspeicher.

Beobachtung: Der Plattenspeicher wird als Menge von **Blöcken** (mit einer Größe im Bereich von 1 - 4 KByte) betrachtet, wobei ein Block durch das Betriebssystem jeweils komplett in den Hauptspeicher übertragen wird. Diese Übertragungseinheiten werden auch als **Seiten** des Plattenspeichers bezeichnet.

Idee: Konstruiere eine Baumstruktur mit der Eigenschaft:

1 Knoten des Baumes $\hat{=}$ 1 Seite (bzw. mehreren Seiten) des Plattenspeichers

Für binäre Baumstrukturen bedeutet das:

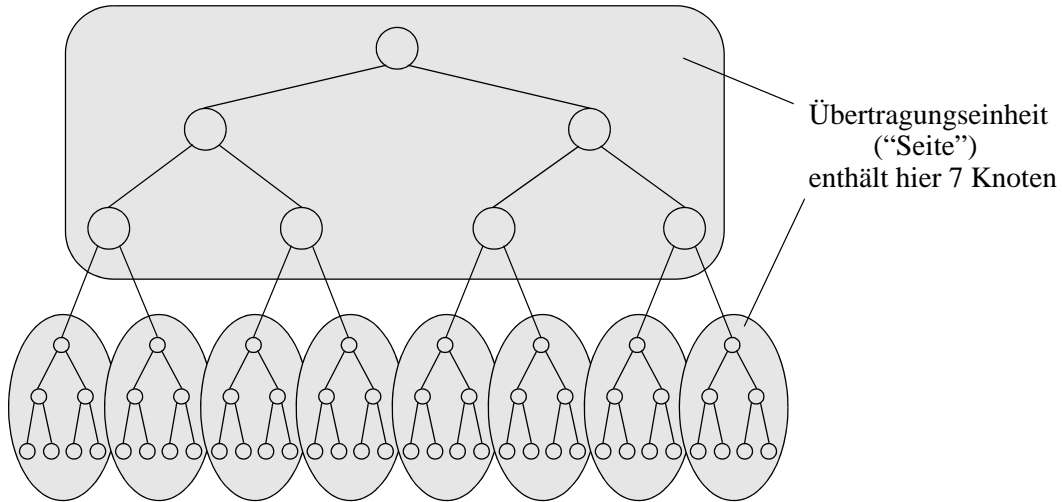
left oder right-Zeiger folgen $\hat{=}$ 1 Plattenspeicherzugriff

Beispiel:

Sei die Anzahl der Datensätze: $n = 10^6$

$$\implies \log_2(10^6) = \log_2(10^3)^2 = 2 \cdot \log_2(10^3) \approx 20 \text{ Plattenspeicherzugriffe}$$

Idee: Zusammenfassen mehrerer binärer Knoten zu einer Seite



Beispiel für einen B-Baum:

99 Knoten in einer Seite \implies 100-fache Verzweigung:

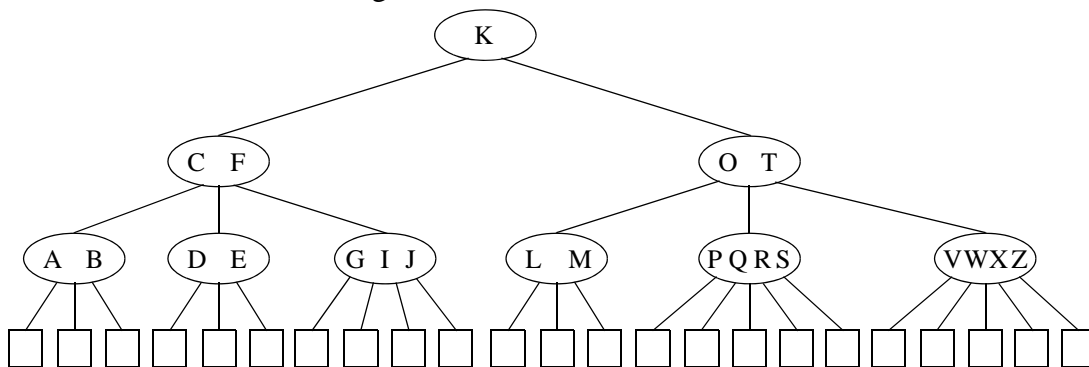
$$\implies \log_{100}(10^6) = 3 \text{ Plattenspeicherzugriffe.}$$

(reduziert auf 2, falls die Wurzelseite immer im Hauptspeicher liegt)

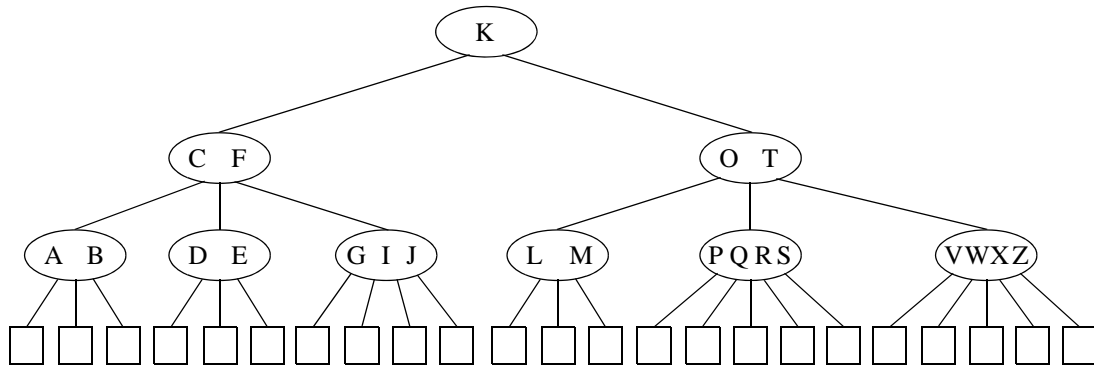
Definition: B-Baum der Ordnung m (Bayer und McCreight (1972))

- (1) Jeder Knoten enthält höchstens $2m$ Schlüssel.
- (2) Jeder Knoten außer der Wurzel enthält mindestens m Schlüssel.
- (3) Die Wurzel enthält mindestens einen Schlüssel.
- (4) Ein Knoten mit k Schlüsseln hat genau $k+1$ Söhne.
- (5) Alle Blätter befinden sich auf demselben Level.

Beispiel: für einen B-Baum der Ordnung 2



Beispiel: derselbe B-Baum nochmals zum Üben



Berechnung der maximalen Höhe h_{max} eines B-Baumes der Ordnung m mit n Schlüsseln:

- Level 1 hat $k_1 = 1$ Knoten
- Level 2 hat $k_2 \geq 2$ Knoten
- Level 3 hat $k_3 \geq 2(m+1)$ Knoten
- ...
- Level $h+1$ hat $k_{h+1} \geq 2(m+1)^{h-1}$ (äußere, leere) Knoten

Ein B-Baum mit n Schlüsseln teilt den Wertebereich der Schlüssel in $n + 1$ Intervalle. Es gilt also

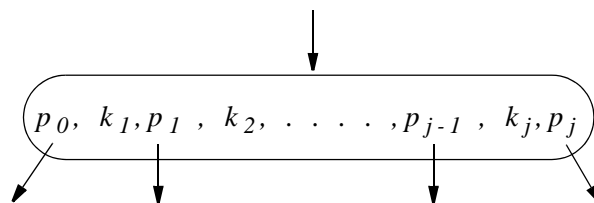
$$k_{h+1} = n + 1 \geq 2(m+1)^{h-1}, \text{ d.h. } h \leq 1 + \log_{m+1} \left(\frac{n+1}{2} \right).$$

Da die Höhe immer ganzzahlig ist, und da diese Rechnung minimalen Füllgrad annimmt, folgt:

$$h \leq \left\lceil \log_{m+1} \left(\frac{n+1}{2} \right) \right\rceil + 1$$

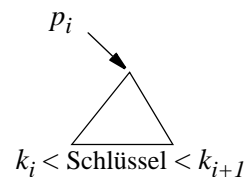
Beobachtung: Jeder Knoten (außer der Wurzel) ist mindestens mit der Hälfte der möglichen Schlüssel gefüllt. Die Speicherplatzausnutzung beträgt also mindestens 50%!

Allgemeine Knotenstruktur:



wobei: $k_1 < k_2 < \dots < k_j$ und $m \leq j \leq 2m$;

p_i zeigt auf den Teilbaum mit Schlüsseln zwischen k_i und k_{i+1} .



Anmerkung: Da die Schlüssel in jedem Knoten eines B-Baumes aufsteigend sortiert sind, kann ein Knoten nach Übertragung in den Hauptspeicher binär durchsucht werden.

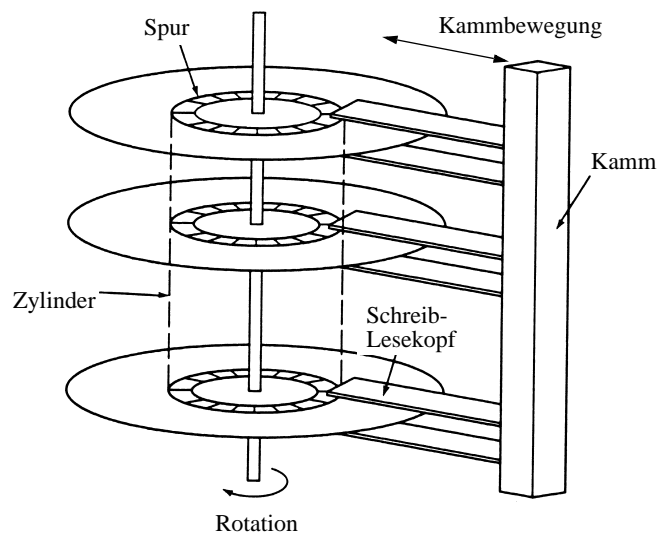
```
class Entry { public int key;
             public Page son;
             ... // Konstruktor..
             }
```

```
class Page { public int numberOfEntries;
            public Page firstSon;
            public Entry [] entries;
            // Im Konstruktor initialisieren mit new Entry [2*m+1];
            ...
            }
```

```
class Btree { protected Page root;
             protected int m;
             ...
             }
```

Welche Ordnung m von B-Bäumen ist auf realen Rechnern und Plattenspeichern günstig?

Hierzu betrachten wir zunächst den physischen Aufbau eines Magnetplattenspeichers. Dieser besteht aus einer Reihe übereinanderliegender rotierender Magnetplatten, die in Zylinder, Spuren und Sektoren unterteilt sind. Der Zugriff erfolgt über einen Kamm mit Schreib-/Leseköpfen, der quer zur Rotation bewegt wird.



Der Seitenzugriff erfolgt nun in mehreren Phasen. Etwas vereinfacht lässt sich die Zugriffszeit in die Zeiten für die folgenden Phasen zerlegen:

Phase	Hard Disk Drive		Solid State Disk	
	Ursache	Zeit	Ursache	Zeit
Positionierungszeit (PZ)	Kammbewegung	8 ms	Addressberechnung	0.1 ms
Latenzzeit (LZ)	Warten auf Sektor	4 ms	---	---
Übertragungszeit (ÜZ)	Übertragung der Daten	$3.3 \cdot 10^{-6}$ ms / Byte	Übertragung der Daten	$1.0 \cdot 10^{-5}$ ms / Byte

→ **Zugriffszeit für eine Seite:** $PZ + LZ + \ddot{U}Z \cdot (\text{Seitengröße})$.

Sei die Größe eines Schlüssels durch α Bytes und die eines Zeigers durch β Bytes gegeben.

→ **Seitengröße** $\approx 2m(\alpha + \beta)$

→ **Zugriffszeit** pro Seite = $PZ + LZ + \ddot{U}Z \cdot 2m(\alpha + \beta) = a + b \cdot m$

mit: $a = PZ + LZ$ und $b = 2(\alpha + \beta) \cdot \ddot{U}Z$.

Andererseits ergibt sich für die **interne Verarbeitungszeit** pro Seite bei binärer Suche:

$c \cdot \log_2(m) + d$ für Konstanten c und d .

Die **gesamte Verarbeitungszeit pro Seite** ist damit: $a + b \cdot m + c \cdot \log_2(m) + d$.

Die maximale Anzahl von Seiten auf einem Suchpfad eines B-Baumes mit n Schlüsseln ist:

$$h_{max} = \left\lceil \log_{m+1} \left(\frac{n+1}{2} \right) \right\rceil + 1 = f \cdot \frac{\log_2 \left(\frac{n+1}{2} \right)}{\log_2(m)} \quad \text{für eine Konstante } f.$$

Die **maximale Suchzeit** MS ist damit gegeben durch die Funktion:

$$MS(m) = g \cdot \left(\frac{a+d}{\log_2(m)} + \frac{b \cdot m}{\log_2(m)} + c \right) \quad \text{mit } g = f \cdot \log_2 \left(\frac{n+1}{2} \right).$$

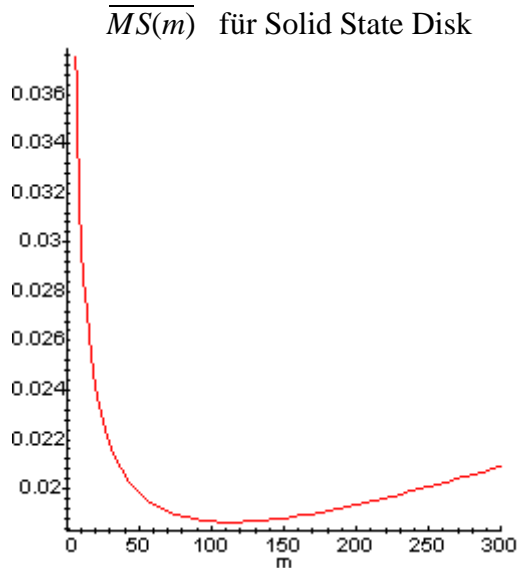
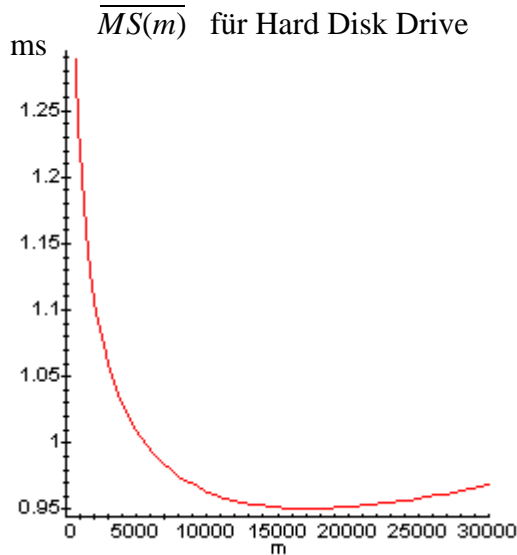
Für die obigen Konstanten setzen wir beispielhaft die folgenden Werte (eines HDD) ein:

$$a = 0,012s, \quad a + d \approx a = 0,012s = 12 \text{ ms}$$

$$\alpha = 8, \quad \beta = 4 \rightarrow b = 24 \cdot 3,3 \cdot 10^{-6} \text{ ms} = 79,2 \cdot 10^{-6} \text{ ms}.$$

Damit ist die folgende Funktion $\overline{MS(m)}$ zu minimieren:

$$\overline{MS(m)} = \left(\frac{12}{\log_2(m)} + \frac{0,0000792 \cdot m}{\log_2(m)} \right) ms \rightarrow \text{MIN.}$$



Die maximale Suchzeit ist somit nahezu minimal für $16000 \leq m \leq 19000$. Dies entspricht in obigem Beispiel einer „optimalen“ Seitengröße von 432 KByte. Für die exemplarischen Werte der Solid-State-Disk ergibt sich eine „optimale“ Seitengröße von 3,0 KB (für $m = 125$).

Einfügen in B-Bäumen:

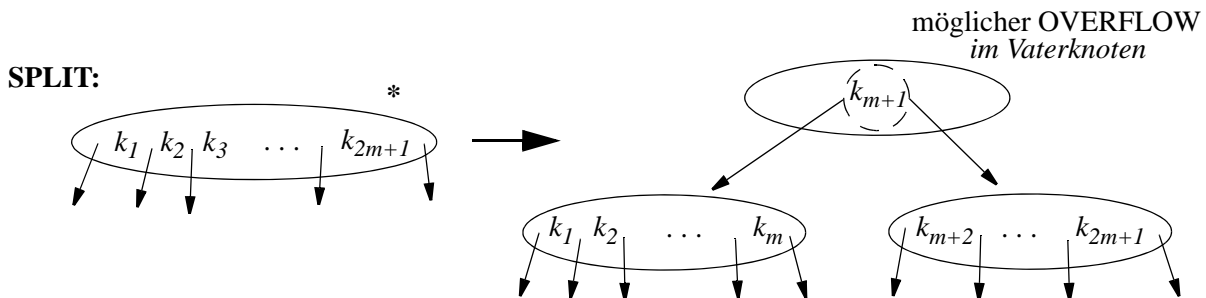
Zunächst wird der Knoten K gesucht, in den der neue Schlüssel einzufügen ist. Dieser ist stets ein Blatt. Der Schlüssel wird in K an der entsprechenden Stelle eingefügt.

Sei s die Anzahl der Schlüssel in K nach dem Einfügen:

Fall 1: $s \leq 2m$: \rightarrow **STOP**

Fall 2: $s = 2m + 1$: \rightarrow **OVERFLOW**

Ein **OVERFLOW** wird behandelt durch Aufspalten (**SPLIT**) des Knotens. Dies kann einen **OVERFLOW** im Vaterknoten zur Folge haben. Auf diese Weise kann sich ein **OVERFLOW** bis zur Wurzel des Baumes fortsetzen. Falls die Wurzel gesplittet wird wächst die Höhe des Baumes um 1.



Entfernen aus B-Bäumen:

Der zu entfernende Schlüssel k wird im Baum gesucht und aus dem gefundenen Knoten K gelöscht. Falls K kein Blatt ist, wird der entfernte Schlüssel durch den Schlüssel p ersetzt, der der kleinste Schlüssel im Baum ist, der größer als k ist. Sei P der Knoten, in dem p liegt. Dann ist P ein Blatt, aus dem p nun entfernt wird. Auf diese Weise wird der Fall "Löschen in einem inneren Knoten" auf den Fall "Löschen in einem Blatt" zurückgeführt.

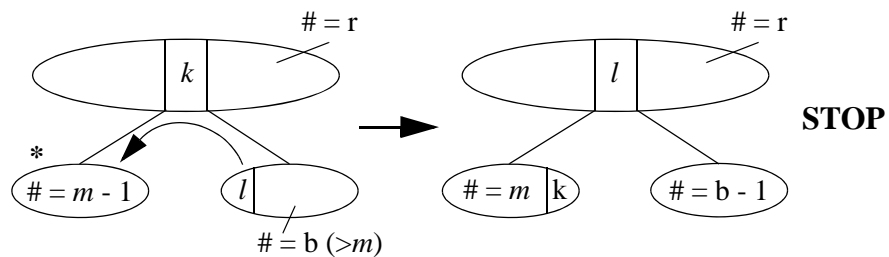
Sei s die Anzahl der Schlüssel in K nach dem Entfernen:

Fall 1: $s \geq m$: \rightarrow **STOP**

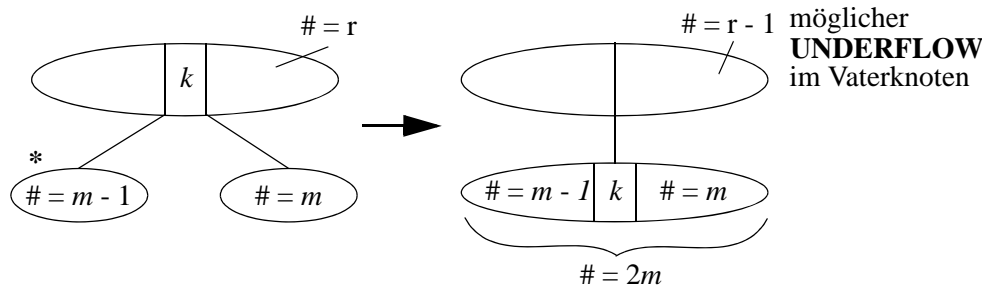
Fall 2: $s = m - 1$: \rightarrow **UNDERFLOW**

UNDERFLOW-Behandlung:

Fall 1: Der Bruder hat $\geq m+1$ Schlüssel: \rightarrow **Ausgleichen** mit dem Bruder



Fall 2: Der Bruder hat m Schlüssel: \rightarrow **Verschmelzen** mit dem Bruder unter Hinzunahme des trennenden Vaterschlüssels
 \rightarrow möglicher **UNDERFLOW** im Vaterknoten.



So kann sich auch die UNDERFLOW-Behandlung bis zur Wurzel des Baumes fortsetzen. Wird aus der Wurzel des Baumes der letzte Schlüssel entfernt, so wird dieser gelöscht; die Höhe des Baumes verringert sich damit um 1.

Korollar: Die B-Bäume bilden eine Klasse balancierter Suchbäume.

Aufwandsabschätzung für die durchschn. Anzahl von Aufspaltungen (Splits) pro Einfügung:

Bezeichne s die durchschnittliche Anzahl der Knoten-Splits pro Einfügung:

Modell:

Ausgehend von einem leeren Baum wird ein B-Baum der Ordnung m für die Schlüssel k_1, k_2, \dots, k_n durch n aufeinanderfolgende Einfügungen konstruiert.

Sei t die Gesamtzahl der Aufspaltungen bei der Konstruktion dieses B-Baumes

$$\rightarrow s = t/n.$$

Sei p die Anzahl der Knoten (Seiten) des B-Baumes mit $p \geq 3$

$$\rightarrow t < p - 1 \text{ (genauer: } t \leq p - 2)$$

Für die Anzahl n der Schlüssel in einem B-Baum der Ordnung m mit p Knoten gilt:

$$n \geq 1 + (p - 1) \cdot m.$$

$$\rightarrow p-1 \leq \frac{n-1}{m} \rightarrow s = \frac{t}{n} < \frac{p-1}{n} \leq \frac{1}{n} \cdot \frac{n-1}{m} < \frac{1}{m} \text{ wobei im allg.: } 16000 \leq m \leq 19000 \text{ (s.o.).}$$

Es gilt: $t \leq p - 2 \rightarrow t < p - 1$ für $p \geq 3$

Denn:

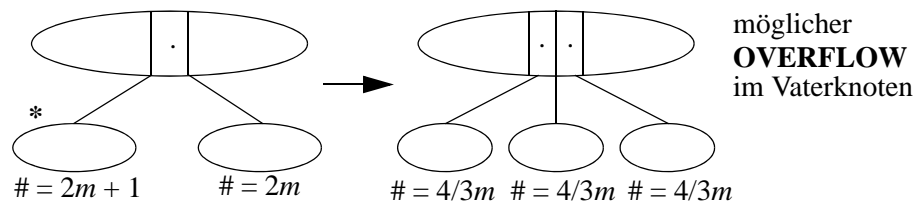
- bei jeder Aufspaltung wird mindestens ein zusätzlicher Knoten geschaffen
- Aufspalten der Wurzel \rightarrow 2 zusätzliche Knoten
- wenn ein B-Baum mehr als einen Knoten hat, dann ist die Wurzel mindestens einmal aufgespalten worden
- dem ersten Knoten des B-Baums geht keine Aufspaltung voraus.

B*-Bäume

B*-Bäume verhalten sich im wesentlichen wie B-Bäume, jedoch wird bei **OVERFLOW** eines Knotens dieser nicht gleich aufgespalten, sondern wie beim UNDERFLOW der Bruder betrachtet:

Fall 1: Bruder hat $b \leq 2m-1$ Schlüssel \rightarrow **Ausgleichen** mit dem Bruder

Fall 2: Bruder hat $b = 2m$ Schlüssel \rightarrow **Verteilen** auf drei Knoten:



Definition: B*-Baum der Ordnung m

m sei ein Vielfaches von 3.

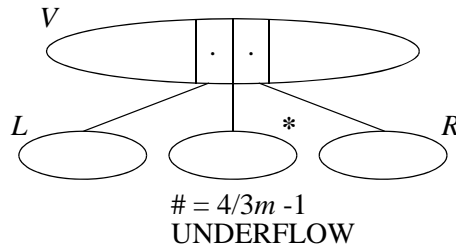
- (1) Jeder Knoten außer der Wurzel enthält mindestens $\frac{4}{3}m$, höchstens $2m$ Schlüssel.
- (2) Die Wurzel enthält mindestens einen, höchstens $\frac{8}{3}m$ Schlüssel.
- (3) Ein Knoten mit k Schlüsseln hat genau $k+1$ Söhne.
- (4) Alle Blätter befinden sich auf demselben Level.

B*-Bäume besitzen eine Speicherplatzausnutzung von mindestens 66%!

Nach einer zum B-Baum äquivalenten Berechnung ergibt sich für die maximale Höhe h_{max} eines B*-Baumes der Ordnung m mit n Schlüsseln:

$$h_{max} = \left\lceil \log_{4/3m+1} \left(\frac{n+1}{2} \right) \right\rceil + 1.$$

Entfernen in B*-Bäumen

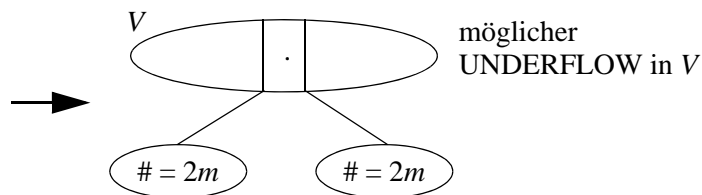


Fall 1: R hat mehr als $\frac{4}{3}m$ Schlüssel \rightarrow Ausgleichen von $*$ und R

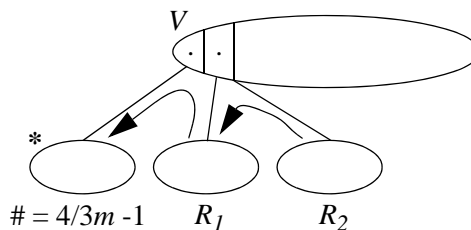
Fall 2: R hat $\frac{4}{3}m$ Schlüssel:

2.1 L hat mehr als $\frac{4}{3}m$ Schlüssel \rightarrow Ausgleichen von $*$ und L

2.2 L hat $\frac{4}{3}m$ Schlüssel:



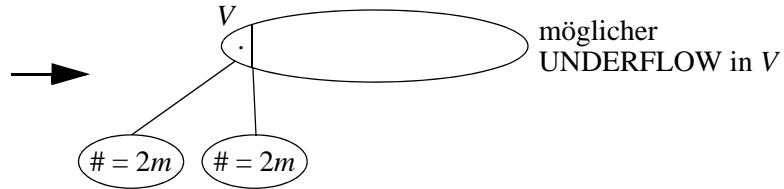
Fall 3: $*$ hat nur einen direkten Bruder, aber weitere indirekte Brüder:



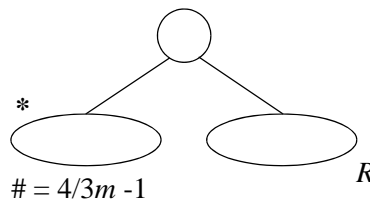
3.1 Falls R_1 mehr als $\frac{4}{3}m$ Schlüssel besitzt \rightarrow Ausgleichen von $*$ und R_1

3.2 Falls R_1 genau $\frac{4}{3}m$ Schlüssel besitzt und R_2 mehr als $\frac{4}{3}m$ Schlüssel besitzt
 → zuerst Ausgleichen von R_1 und R_2 , dann von $*$ und R_1

3.3 Falls R_1 und R_2 jeweils genau $\frac{4}{3}m$ Schlüssel besitzen



Fall 4: $*$ ist Sohn einer binären Wurzel:



4.1 R hat $\frac{4}{3}m$ Schlüssel:



4.2 R hat mehr als $\frac{4}{3}m$ Schlüssel → Ausgleichen von $*$ mit R

Korollar: Die B^* -Bäume bilden eine Klasse balancierter Suchbäume.

Bemerkung: In fast allen gängigen Datenbanksystemen sind B-Bäume, B^* -Bäume oder deren Varianten zur effizienten Ausführung von Suchoperationen implementiert.

2.3 Optimale binäre Suchbäume

Bisher: Häufigkeiten, mit denen einzelne Schlüssel gesucht werden, bleiben unberücksichtigt. → **Annahme der Gleichverteilung** der Suchhäufigkeiten

Jetzt: Wir **berücksichtigen Häufigkeiten**, mit denen einzelne Schlüssel gesucht werden
 → Je häufiger ein Schlüssel gesucht wird, desto höher im Baum soll er plaziert werden.

Anwendungsbeispiel:

→ "Tabelle" der Schlüsselworte bei Compilern (*statisch!*).

Seien die Zugriffshäufigkeiten zu den n einzufügenden Schlüsseln k_1, k_2, \dots, k_n mit p_1, p_2, \dots, p_n

($\sum_{i=1}^n p_i = 1$ nicht unbedingt erforderlich) zum Zeitpunkt der Einfügung bekannt.

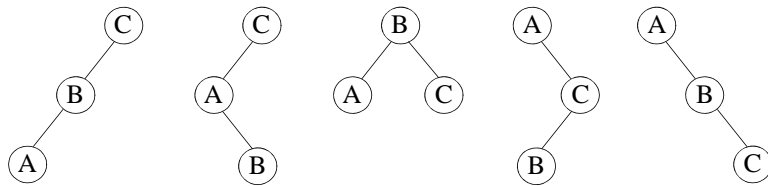
Problem: Wie finden wir unter diesen Voraussetzungen einen “optimalen” Suchbaum aus den $\binom{2n}{n} / (n + 1) \approx 4^n / (\sqrt{\pi} \cdot n^{3/2})$ verschiedenen Suchbäumen für n Schlüssel?

Beispiel:

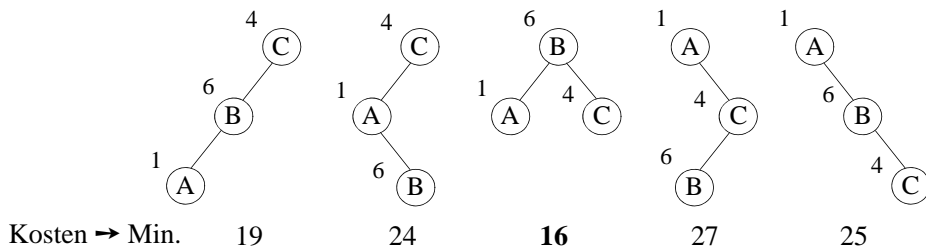
Sei $n = 3$;

Gegeben seien die Schlüssel A, B und C mit den Zugriffshäufigkeiten p_A, p_B und p_C .

→ 5 mögliche Suchbäume; welcher optimal ist, hängt vom Gewicht (den “Kosten”) des Baumes und damit von den Gewichten der einzelnen Schlüssel ab.



Konkretes Beispiel: $p_A = 1, p_B = 6$ und $p_C = 4$.



Kosten = $\alpha \cdot$ (erwartete Anzahl von Knotenzugriffen) + $\beta \cdot$ (erwartete Anzahl von Schlüsselvergleichen)
 wobei: α : Kosten für einen Knotenzugriff
 β : Kosten für einen Schlüsselvergleich

Für einen binären Baum gilt: Kosten = $(\alpha + \beta) \cdot$ (erwartete Anzahl von Schlüsselvergleichen)

Für eine allgemeine Lösung werden wir zusätzlich erfolglose Suchoperationen, d.h. die Suche von Schlüsseln, die zwischen den vorhandenen Schlüsseln liegen, berücksichtigen.

Damit ergibt sich die folgende allgemeine **Problemstellung**:

Gegeben: n Schlüssel: $k_1 < k_2 < \dots < k_n$.

$2n+1$ Gewichte (Häufigkeiten):

$$p_1, p_2, \dots, p_n \text{ und } q_0, q_1, q_2, \dots, q_n \text{ mit } w = \sum_{i=1}^n p_i + \sum_{i=0}^n q_i, \text{ wobei:}$$

$$p_i/w = \text{Wahrscheinlichkeit } (k_i = \text{Suchargument})$$

$$q_i/w = \text{Wahrscheinlichkeit } (k_i < \text{Suchargument} < k_{i+1})$$

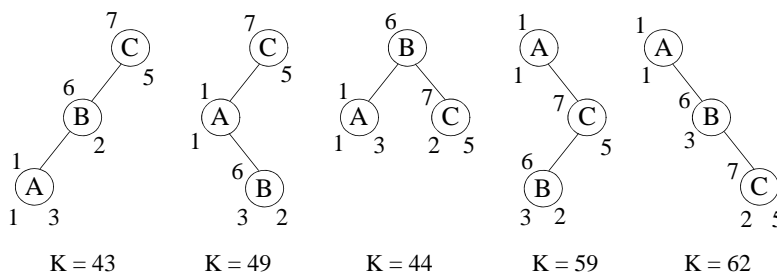
$$(\text{o.E. } k_0 = -\infty, k_{n+1} = \infty)$$

Zu konstruieren ist ein binärer Suchbaum derart, daß die zu erwartende Anzahl der bei einer Suchoperation auszuführenden Schlüsselvergleiche, ausgedrückt durch das folgende **Kostenmaß K**, minimiert wird:

$$K = \sum_{i=1}^n p_i \cdot lev(\bigcirc i) + \sum_{j=0}^n q_j \cdot (lev(\square j) - 1) \rightarrow \text{MIN.}$$

mit: $\bigcirc j$ bezeichne den Baumknoten mit Schlüssel k_j und
 $\square k$ bezeichne den äußeren Baumknoten "zwischen" $\bigcirc k$ und $\bigcirc(k+1)$.

Beispiel: Seien im obigen Beispiel: $p_A = 1, p_B = 6, p_C = 7$ und $q_0 = 1, q_1 = 3, q_2 = 2, q_3 = 5$.

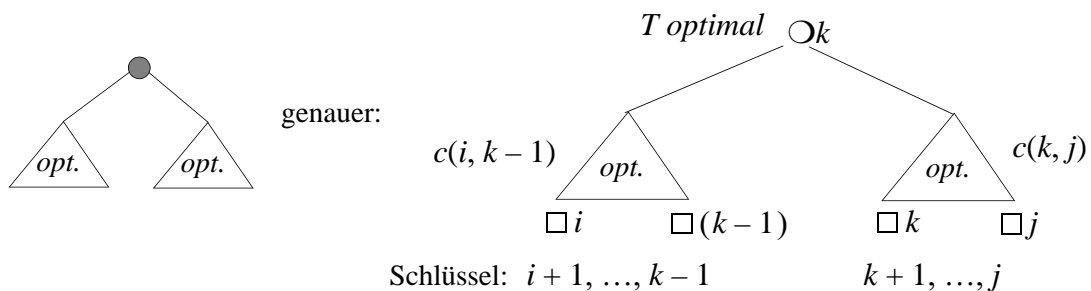


Optimalitätskriterium: Alle **Teilbäume** eines optimalen Suchbaumes sind **optimal**.

→ **Idee:** ausgehend von einzelnen Knoten als minimalen Teilbäumen werden systematisch "immer größere" optimale Teilbäume konstruiert.

"Bottom-up" - Methode:

Ein neuer optimaler Teilbaum ergibt sich aus einer geeigneten Wurzel und zwei optimalen Teilbäumen (*bereits berechnet*) dieser Wurzel.



Seien $c(i, j)$, $0 \leq i \leq j \leq n$, die Kosten eines optimalen Teilbaumes mit den Gewichten $p_{i+1}, \dots, p_j; q_i, \dots, q_j$; sei weiterhin: $w(i, j) = \sum_{k=i+1}^j p_k + \sum_{l=i}^j q_l$.

Dann können die Werte $c(i, j)$ nach folgendem Rekursionsschema berechnet werden:

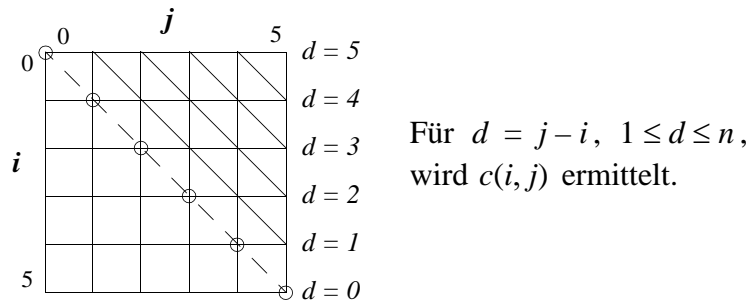
$$c(i, i) = 0 \rightarrow \text{der Baum besteht aus } \square i \text{ (leerer Baum).}$$

$$c(i, j) = w(i, j) + \min_{i < k \leq j} (c(i, k-1) + c(k, j)) \quad \text{für } i < j.$$

Denn: das minimale Gewicht des Teilbaums mit Wurzel $\circ k$ und den Grenzen i bzw. j ist gegeben durch: $w(i, j) + c(i, k-1) + c(k, j)$.

Mit Hilfe dieses Rekursionsschemas wird $c(0, n)$ und der zugehörige optimale binäre Suchbaum bestimmt.

Es gibt hierbei $\approx \frac{n^2}{2}$ Werte $c(i, j) \rightarrow$ **Platzbedarf** $O(n^2)$.



Laufzeitbedarf für diese Art der Konstruktion eines optimalen binären Suchbaumes:

Für $d = j - i, 1 \leq d \leq n,$ berechnen wir $c(i, j)$.

Wieviele Werte für $\circ k$ müssen bei der Minimum-Bestimmung betrachtet werden?

Für jedes $d, 1 \leq d \leq n,$ gibt es $n - d + 1$ verschiedene $c(i, j)$'s.

Für jedes dieser $c(i, j)$ mit $j - i = d$ gibt es d mögliche Werte für $\circ k$.

Damit gibt es für jedes $d, 1 \leq d \leq n, d \cdot (n - d + 1)$ mögliche Werte für $\circ k$.

Insgesamt $\sum_{d=1}^n d \cdot (n - d + 1) = \frac{n^3}{6} + O(n^2)$ mögliche Werte für $\circ k$.

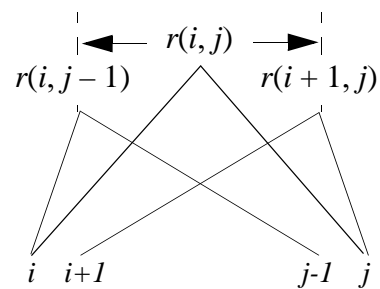
$$\left(\sum_{d=1}^n d^2 = \frac{n \cdot (n + 1) \cdot (2n + 1)}{6} \right)$$

\rightarrow **Laufzeitbedarf:** $O(n^3)$; dies ist für einen praktischen Einsatz indiskutabel!

Es gilt jedoch folgende **Monotonie-Eigenschaft:**

Sei $r(i, j)$ ein Wert für $\circ k$, für den $c(i, j)$ minimal wird.

$$\rightarrow r(i, j - 1) \leq r(i, j) \leq r(i + 1, j)$$



→ Es genügt: $c(i, j) = w(i, j) + \min_{r(i, j-1) \leq k \leq r(i+1, j)} (c(i, k-1) + c(k, j))$ zu berechnen.

Für jedes d , $1 \leq d \leq n$, müssen damit also weniger Werte für Ok untersucht werden, nämlich nur noch:

$$\sum_{\substack{d \leq j \leq n \\ i = j-d}} (r(i+1, j) - r(i, j-1) + 1) = r(n-d+1, n) - r(0, d-1) + (n-d+1) < 2n$$

→ **Laufzeitbedarf:** reduziert auf $O(n^2)$.

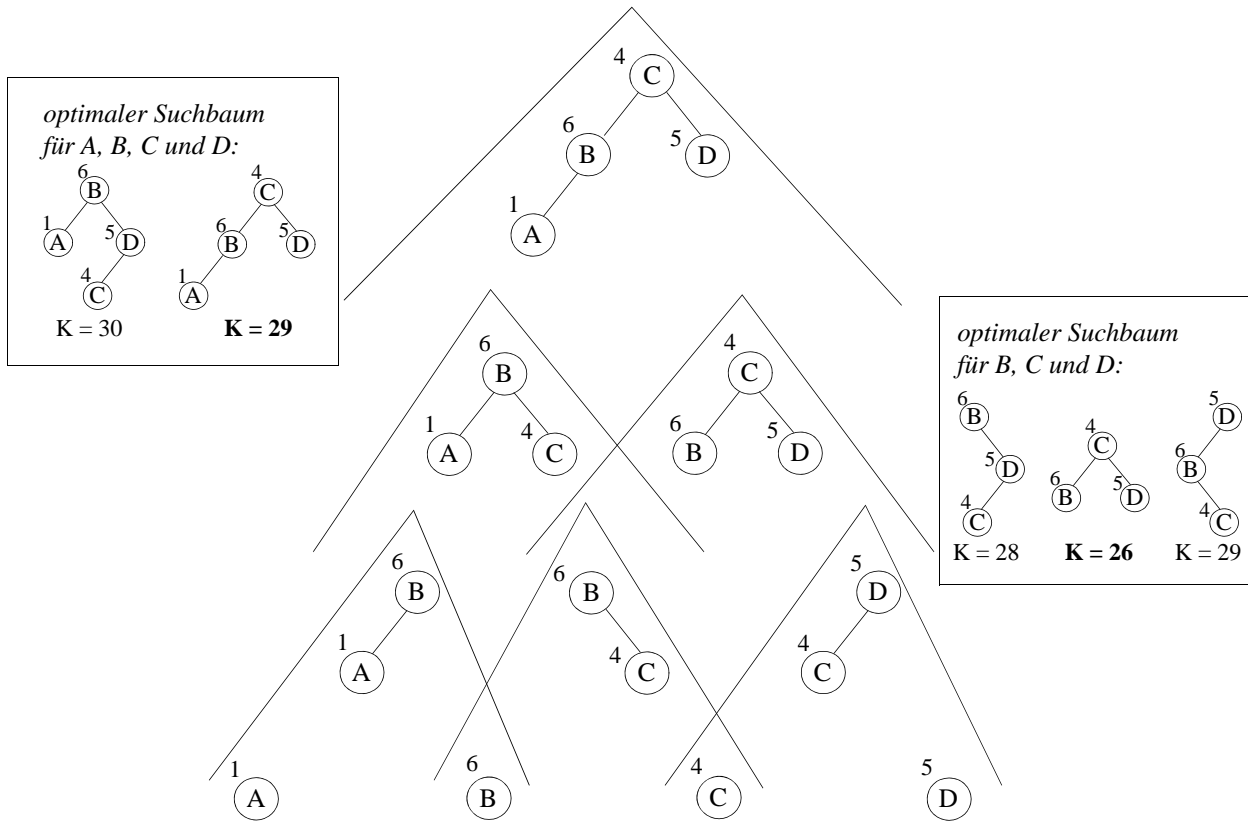
Abschließendes Beispiel:

Sei $n = 4$;

Gegeben seien die Schlüssel A, B, C und D mit den Zugriffshäufigkeiten:

$$p_A = 1, p_B = 6, p_C = 4 \text{ und } p_D = 5.$$

Sei weiterhin: $q_i = 0, 0 \leq i \leq 4$.



Wie bereits angedeutet, sind optimale binäre Suchbäume nur für statische Anwendungen geeignet (kein Einfügen oder Entfernen). Nahezu-optimale Suchbäume zur Handhabung von Zugriffshäufigkeiten und gleichzeitig Einfügen bzw. Entfernen werden später behandelt.

2.4 Hashverfahren

Bisher: Suchen mit Hilfe von **Schlüsselvergleichen**.

Jetzt: stattdessen **Adressberechnung**

→ Auswertung einer Funktion, der sogenannten **Hash-** oder **Adressfunktion**.

Vorteil: Die Suche erfolgt weitgehend unabhängig von den anderen Schlüsselwerten

→ sie ist in der Regel schneller.

(im Durchschnitt konstanter statt logarithmischer Aufwand)