

CS 240 – Data Structures and Data Management

Module 4: Dictionaries - Enriched

T. Biedl E. Kondratovsky M. Petrick O. Veksler

Based on lecture notes by many previous cs240 instructors

David R. Cheriton School of Computer Science, University of Waterloo

Winter 2022

Outline

- 1 More Balanced Search Trees
 - Balanced binary search trees
 - Scapegoat Trees
 - Amortized analysis
 - Analysis of scapegoat trees

Outline

- 1 More Balanced Search Trees
 - Balanced binary search trees
 - Scapegoat Trees
 - Amortized analysis
 - Analysis of scapegoat trees

Overview of balanced binary search trees

We will see numerous variants of binary search trees.

The operations then have the following run-times:

- $\Theta(\log n)$ worst-case time (**AVL-trees**)
- $\Theta(\log n)$ amortized time (**Scapegoat trees**) and no rotations.
- $\Theta(\log n)$ expected time (**Treaps**)
- $\Theta(\log n)$ expected time (**Skip lists**) and space is smaller. (It's not even a tree.)
- $\Theta(\log n)$ amortized time (**Splay trees**) and space is smaller, and can handle biased requests.

(We will see “rotations”, “amortized” and “biased requests” later.)

Overview of balanced binary search trees

General strategy for balanced binary search tree:

- Use a binary search tree, but impose *structural condition*
- Argue that structural condition implies $O(\log n)$... height (where ... might be worst-case / avg-case / expected)
- With this, search takes $O(\log n)$... time
- Explain how to do insert and delete so that structural condition continues to hold.
 - ▶ This must be done so that run-time stays $O(\text{height})$
 - ▶ With this, insert/delete takes $O(\log n)$... time

Outline

- 1 More Balanced Search Trees
 - Balanced binary search trees
 - Scapegoat Trees
 - Amortized analysis
 - Analysis of scapegoat trees

Scapegoat trees

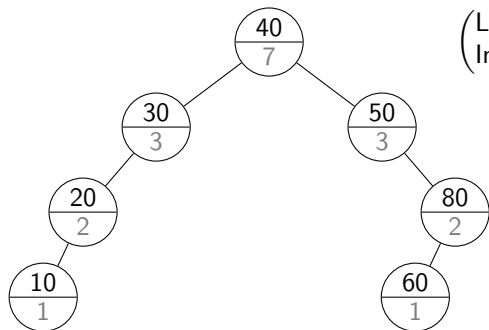
- Can we have balanced binary search trees *without* rotations?
(A later application will need such a tree.)
- This sounds impossible—we must sometimes restructure the tree.
- **Idea:** Rather than doing a small local change, occasionally do a large (near-global) rebuild.
- With the right setup, this will lead to $O(\log n)$ height and $O(\log n)$ *amortized* time for all operations.

Scapegoat trees

Fix a constant α with $\frac{1}{2} < \alpha < 1$. A **scapegoat tree** is a binary search tree where any node v with a parent satisfies

$$v.size \leq \alpha \cdot v.parent.size.$$

(Lower number = subtree-size.)
(In our examples, $\alpha = \frac{2}{3}$.)

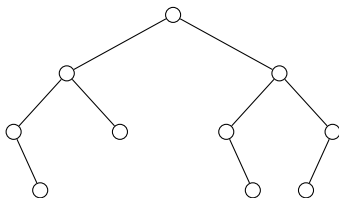


- $v.size$ needed during updates \rightsquigarrow must be stored
- Any subtree is a constant fraction smaller \rightsquigarrow height $O(\log n)$.

Scapegoat tree operations

- *search*: As for a binary search tree. $O(\text{height}) = O(\log n)$.
- For *insert* and *delete*, occasionally restructure a subtree into a **perfectly balanced tree**:

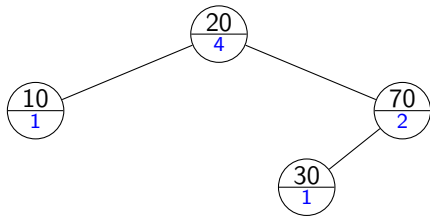
$$|\text{size}(z.\text{left}) - \text{size}(z.\text{right})| \leq 1 \quad \text{for all nodes } z.$$



- Do this at the *highest* node where the size-condition of scapegoat trees is violated

Scapegoat Tree Insertion Example

Example:



Scapegoat tree insertion

```
scapegoatTree::insert(k, v)
1.   $z \leftarrow \text{BST}::\text{insert}(k, v)$ 
2.   $S \leftarrow$  stack initialized with  $z$ 
3.  while ( $p \leftarrow z.\text{parent} \neq \text{NIL}$ )    // update sizes, get path
4.      increase  $p.\text{size}$ 
5.       $S.\text{push}(p)$ 
6.       $z \leftarrow p$ 
7.  while ( $S.\text{size} \geq 2$ )                // size-condition violated?
8.       $p \leftarrow P.\text{pop}()$ 
9.      if ( $p.\text{size} < \alpha \cdot \max\{p.\text{left}.\text{size}, p.\text{right}.\text{size}\}$ )
10.         rebuild subtree at  $p$  into perfectly balanced tree
11.         return
```

- Rebuilding at p (line 10) can be done in $O(p.\text{size})$ time (exercise).
- This restores scapegoat tree (we rebuild at the highest violation).

Outline

- 1 More Balanced Search Trees
 - Balanced binary search trees
 - Scapegoat Trees
 - **Amortized analysis**
 - Analysis of scapegoat trees

Detour: Amortized analysis

As for dynamic arrays and lazy deletion, we have the following pattern:

- usually the operation is fast,
- the occasional operation is quite slow.

The worst-case run-time bound here would not reflect that overall this works quite well.

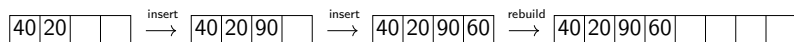
Instead, try to find an **amortized run-time bound**: A bound that holds if we add the bounds up over all operations.

$$\sum_{i=1}^k T^{\text{actual}}(\mathcal{O}_i) \leq \sum_{i=1}^k T^{\text{amort}}(\mathcal{O}_i).$$

(where $\mathcal{O}_1, \dots, \mathcal{O}_k$ is any feasible sequence of operations, $T^{\text{actual}}(\cdot)$ is the actual run-time, and $T^{\text{amort}}(\cdot)$ is the amortized run-time (or an upper bound for it).

Detour: Amortized analysis

For dynamic arrays, some ad-hoc methods work.



- Direct argument:
 - ▶ $n/2$ fast inserts takes $\Theta(1)$ time each.
 - ▶ Then one slow insert takes $\Theta(n)$.
 - ▶ Averaging out therefore $\Theta(1)$ per operation.
 - ▶ This is doing math with asymptotic notation - dangerous.
- Explicitly define $T^{\text{amort}}(\cdot)$ and verify.
 - ▶ Set time units such that $T^{\text{actual}}(\text{insert}) \leq 1$ and $T^{\text{actual}}(\text{resize}) \leq n$.
 - ▶ Define $T^{\text{amort}}(\text{insert}) = 3$ and $T^{\text{amort}}(\text{resize}) = 0$.

$$\text{Verify } \sum_{i=1}^k T^{\text{actual}}(\mathcal{O}_i) \leq \sum_{i=1}^k T^{\text{amort}}(\mathcal{O}_i).$$

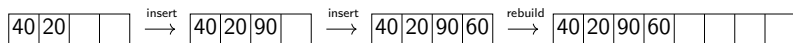
Potential function method

Usually we need more systematic methods.

- **Potential function:** A function $\Phi(\cdot)$ that depends on the current status of the data structure.
 - ▶ E.g.: $\Phi(i) = \max\{0, 2 \cdot \text{size} - \text{capacity}\}$ for dynamic arrays.
 - ▶ “ i ” = operations $\mathcal{O}_1, \dots, \mathcal{O}_i$ have been executed.
- Potential function must satisfy: $\Phi(0) = 0$, $\Phi(i) \geq 0$ for all i .
 - ▶ Can verify this for dynamic-array function above.
- Define $T^{\text{amort}}(\mathcal{O}_i) = T^{\text{actual}}(\mathcal{O}_i) + \Phi(i) - \Phi(i-1)$
 - ▶ Often we just write $T^{\text{amort}}(\mathcal{O}) = T^{\text{actual}}(\mathcal{O}) + \Phi^{\text{after}} - \Phi^{\text{before}}$

Lemma: This satisfies $\sum_i T^{\text{actual}}(\mathcal{O}_i) \leq \sum_i T^{\text{amort}}(\mathcal{O}_i)$.

Example: Dynamic arrays



- Potential function $\Phi(i) = \max\{0, 2 \cdot \text{size} - \text{capacity}\}$
- As before set time units such that

$$T^{\text{actual}}(\text{insert}) \leq 1 \text{ and } T^{\text{actual}}(\text{resize}) \leq n.$$

- *insert* increases *size*, does not change *capacity*

$$\Rightarrow \Delta\Phi = \Phi^{\text{after}} - \Phi^{\text{before}} \leq 2 - 0 = 2$$

$$T^{\text{amort}}(\text{insert}) \leq 1 + 2 - 0 = 3 \in O(1)$$

- *rebuild* happens only if $\text{size} = \text{capacity} = n$

$$\Rightarrow \Phi^{\text{before}} = 2n - n = n.$$

$$\Rightarrow \Phi^{\text{after}} = 2n - 2n = 0 \text{ since new capacity is } 2n.$$

$$T^{\text{amort}}(\text{insert}) \leq n + 0 - n = 0 \in O(1)$$

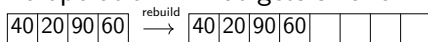
Result: The amortized run-time of dynamic arrays is $O(1)$.

Potential function method

How to find a suitable potential function?

(No recipe, but some guidelines.)

- Study the expensive operation: What gets smaller?



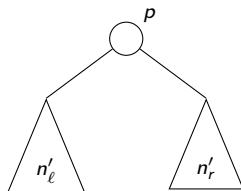
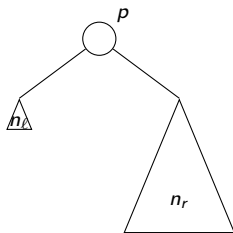
- ▶ Dynamic arrays: *rebuild* increases capacity.
We want the potential function to get smaller.
So potential function should have term “ $-capacity$.”
- Study condition $\Phi(\cdot) \geq 0$ and $\Phi(0) = 0$.
 - ▶ Dynamic arrays: Usually have $capacity \leq 2 \cdot size$.
So usually $2 \cdot size - capacity \geq 0$,
 - ▶ We added a $\max\{0, \dots\}$ term so that also $\Phi(0) = 0$.
- Compute the amortized time and see whether you get good bounds.
- Rinse, lather, repeat.

Outline

- 1 More Balanced Search Trees
 - Balanced binary search trees
 - Scapegoat Trees
 - Amortized analysis
 - Analysis of scapegoat trees

Amortized analysis of scapegoat trees

- Expensive operation: Rebuild subtree at p .



- Claim:** If we rebuild at p , then $|n_r - n_\ell| \geq (2\alpha - 1)n_p$.
Proof:

- Idea:** Potential function should involve $\sum_v |v.\text{left.size} - v.\text{right.size}|$.

Amortized analysis of scapegoat trees

- Use $\Phi(i) = c \cdot \sum_v \max\{|v.\text{left} - v.\text{right}| - 1, 0\}$ for some constant c .
- *insert* and *delete* increases contribution at ancestors by at most 1 and does not increase other contributions.

$$\begin{aligned} T^{\text{amort}}(\text{insert}) &= T^{\text{actual}}(\text{insert}) + \Phi_{\text{after}} - \Phi_{\text{before}} \\ &\leq \log n + c \#\{\text{ancestors}\} \in O(\log n) \end{aligned}$$

- *rebuild* decreases contribution at p by $(2\alpha - 1)n_p$ and does not increase other contributions.

$$\begin{aligned} T^{\text{amort}}(\text{rebuild}) &= T^{\text{actual}}(\text{rebuild}) + \Phi_{\text{after}} - \Phi_{\text{before}} \\ &\leq n_p + c(-(2\alpha - 1)n_p) \end{aligned}$$

With $c = 1/(2\alpha - 1)$, this is at most 0 and *rebuild* is free.

Result: All operations have amortized run-time in $O(\log n)$.