



# B+ Trees

## CS1: Java Programming Colorado State University

Slides by Wim Bohm, Russ Wakefield and Sudipto Ghosh

# Motivations

Many times you want to minimize the disk accesses while doing a search.

A binary search tree allows one key per node.

A B+ tree allows as many values that will fit on a page.



# Differences between BST and B+

- ◆ B+ is a balanced tree
  - Same distance for every path through the tree
  - Not true for BST
- ◆ B+ tree keeps data only in the leaf nodes
  - The index nodes are simply for traffic control
  - BST has data with the keys
- ◆ B+ tree has many keys per node
  - BST has 1



# Objectives

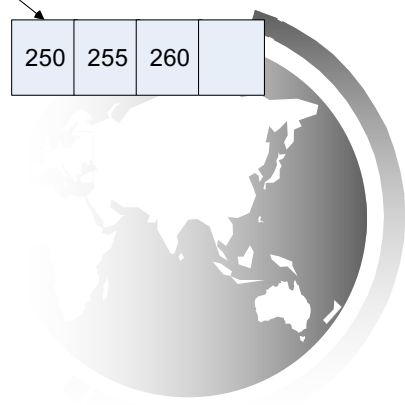
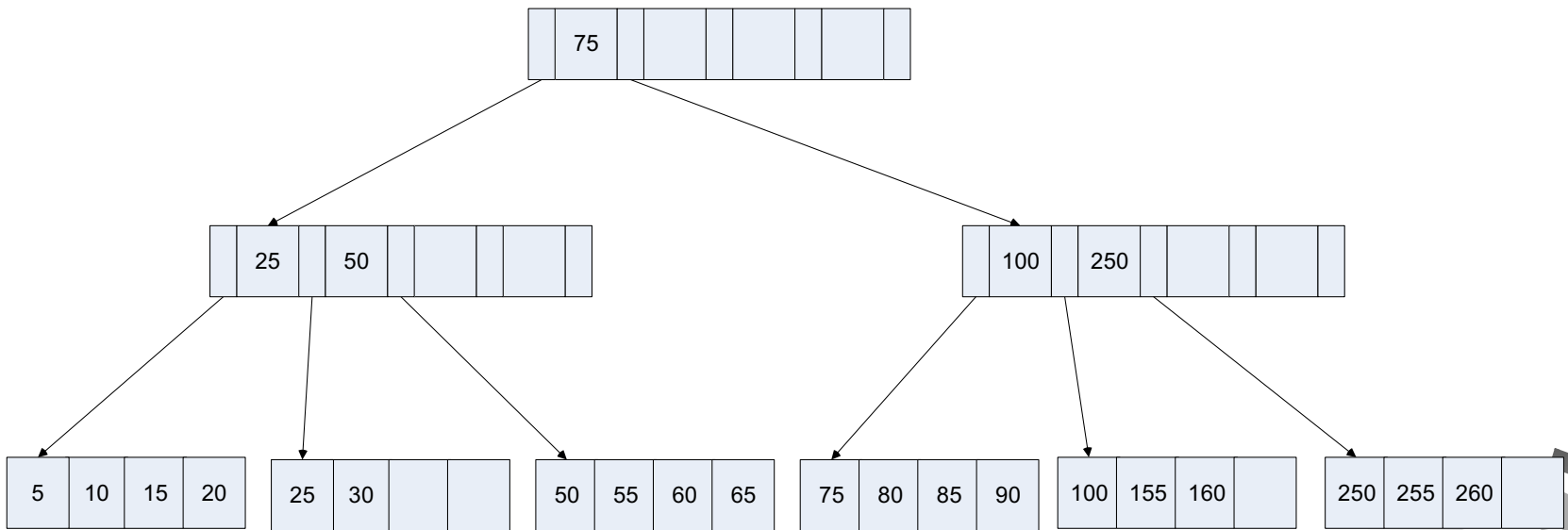
- To understand the basic structure of a B+ Tree
- To understand how the keys are stored and how it allows for efficient searching
- To understand the insertion algorithm
- To understand the deletion algorithm
- To see how rotation minimizes the cost of insertion and deletion algorithms

# What is a B+ tree?

- ◆ Dynamic structure that gracefully adjusts to changes in the file.
- ◆ Most widely used structure because it adjusts well to changes and supports both equality and range queries.
- ◆ Balanced tree in which the
  - Internal nodes direct the search
  - Leaf nodes contain the data entries.
- ◆ Leaf nodes organized into a doubly linked list.
- ◆ Allows us to easily traverse the leaf pages in either direction.



# Example of a B+ tree



# Advantages / Disadvantages

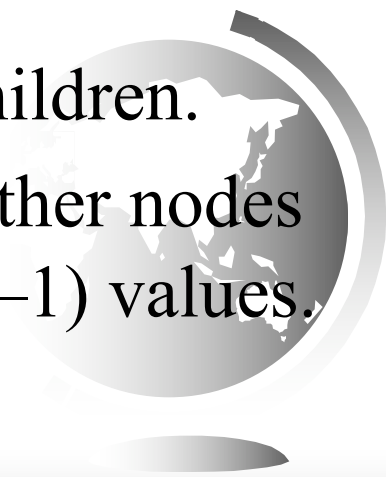
- ◆ Advantage of B<sup>+</sup> tree index files:
  - automatically reorganizes itself with small, local, changes, in the face of insertions and deletions.
  - Reorganization of entire file is not required to maintain performance.
- ◆ (Minor) disadvantage of B<sup>+</sup> trees:
  - extra insertion and deletion overhead, space overhead.
- ◆ Advantages of B<sup>+</sup> trees outweigh disadvantages
  - B<sup>+</sup> trees are used extensively



# B<sup>+</sup> Tree Index Files

A B<sup>+</sup> tree is a rooted tree satisfying the following properties:

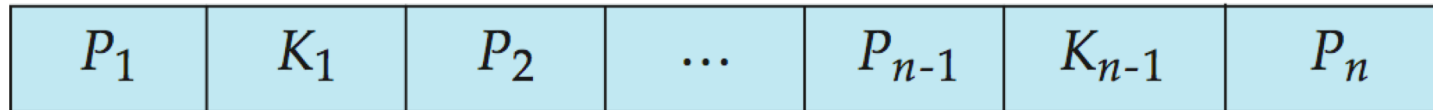
- ◆ All paths from root to leaf are of the same length
- ◆ Each node that is not a root or a leaf has between  $\lceil n/2 \rceil$  and  $n$  children.
- ◆ A leaf node has between  $\lceil (n-1)/2 \rceil$  and  $n-1$  values
- ◆ Special cases:
  - If the root is not a leaf, it has at least 2 children.
  - If the root is a leaf (that is, there are no other nodes in the tree), it can have between 0 and  $(n-1)$  values.





# B<sup>+</sup> Tree Node Structure

## ◆ Typical node



- $K_i$  are the search-key values
- $P_i$  are pointers to children (for non-leaf nodes) or pointers to records or buckets of records (for leaf nodes).

## ◆ The search-keys in a node are ordered

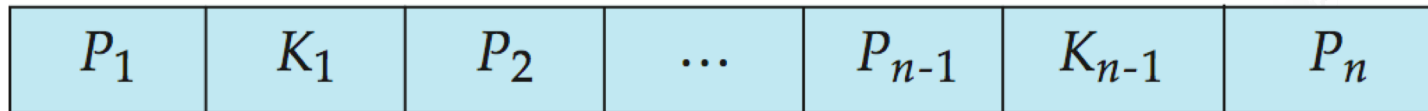
$$K_1 < K_2 < K_3 < \dots < K_{n-1}$$

Initially assume no duplicate keys, address duplicates later.

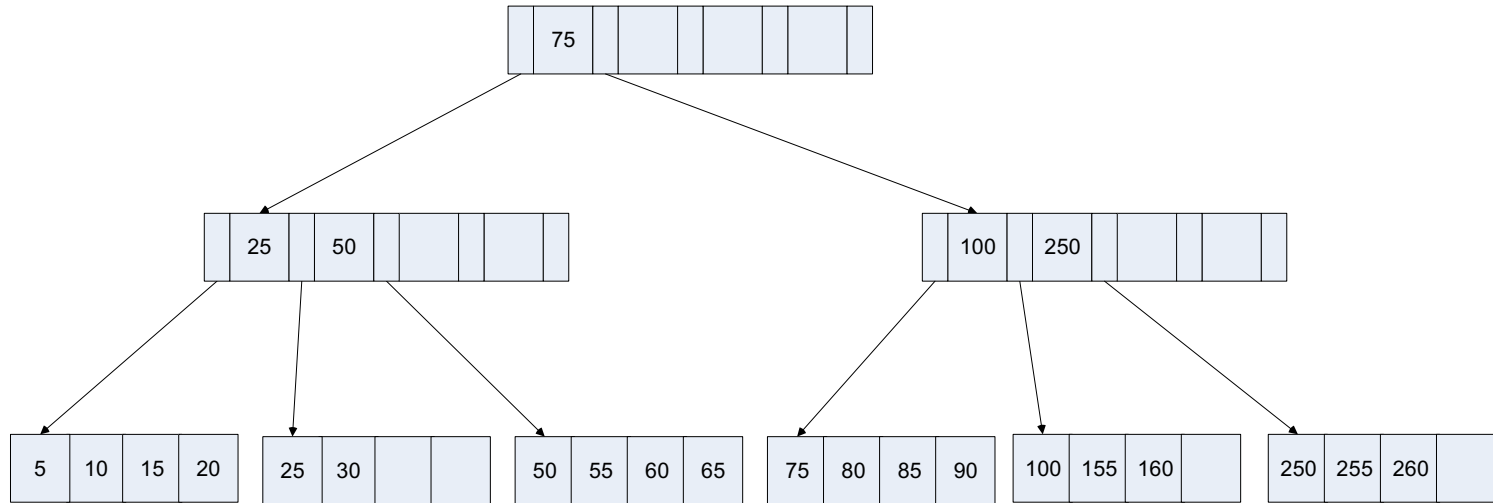


# Non-Leaf Nodes in B<sup>+</sup>Trees

- ◆ Non-leaf nodes form a multi-level sparse index on the leaf nodes.
- ◆ For a non-leaf node with  $n$  pointers:
  - All the search-keys in the subtree to which  $P_1$  points are less than  $K_1$
  - For  $2 \leq i \leq n - 1$ , all the search-keys in the subtree to which  $P_i$  points have values greater than or equal to  $K_{i-1}$  and less than  $K_i$
  - All the search-keys in the subtree to which  $P_n$  points have values greater than or equal to  $K_{n-1}$



# Example of B<sup>+</sup>tree



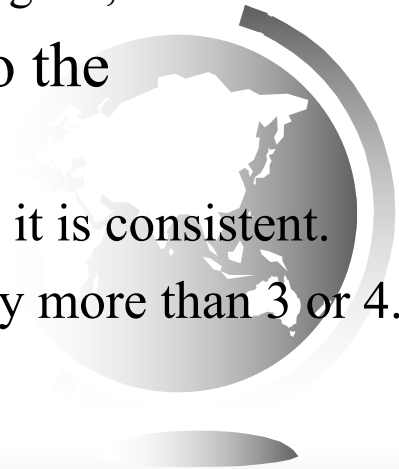
B<sup>+</sup>-tree for *instructor* file ( $n = 5$ )

- ◆ Leaf nodes must have between 2 and 4 values ( $\lceil (n-1)/2 \rceil$  and  $n - 1$ , with  $n = 5$ ).
- ◆ Non-leaf nodes other than root must have between 3 and 5 children ( $\lceil n/2 \rceil$  and  $n$  with  $n = 5$ ).
- ◆ Root must have at least 2 children.



# Observations about B<sup>+</sup>trees

- ◆ Operations (insert, delete) on the tree keep it balanced.
- ◆  $\log_f N$  cost where  $f$ =fanout,  $N$  = # of leaf pages.
- ◆ Minimum occupancy of 50% is guaranteed for each node except the root node if the deletion algorithm we will present is used.
- ◆ Each node contains  $m$  entries where  $d \leq m \leq 2d$  entries
  - $d$  is the degree, which was 2 in the previous slide
  - In practice, deletes just delete the data entry because files usually grow, not shrink.
- ◆ Search for a record is just a traversal from the root to the appropriate leaf.
  - This is the height of the tree – because the tree is balanced, it is consistent.
  - Because of the high fan-out, the height of a B<sup>+</sup> tree is rarely more than 3 or 4.



# Queries on B<sup>+</sup>Trees

- ◆ Find record with search-key value  $V$ .
  1.  $C = \text{root}$
  2. While  $C$  is not a leaf node {
    1. Let  $i$  be least value such that  $V \leq K_i$ .
    2. If no such exists, set  $C = \text{last non-null pointer in } C$
    3. Else { if ( $V = K_i$ ) Set  $C = P_{i+1}$  else set  $C = P_i$  }  
} //Now  $C$  is in leaf node containing  $K_i$
  3. Let  $i$  be least value such that  $K_i = V$
  4. If there is such a value  $i$ , follow pointer  $P_i$  to the desired record.
  5. Else no record with search-key value  $k$  exists.



# Queries on B<sup>+</sup>Trees (Cont.)

- ◆ If there are  $K$  search-key values in the file, the height of the tree is no more than  $\lceil \log_{\lceil n/2 \rceil}(K) \rceil$ . ( $n$  = number of indices / node)
- ◆ A node is generally the same size as a disk block, typically 4 kilobytes
  - and  $n$  is typically around 100 (40 bytes per index entry).
- ◆ With 1 million search key values and  $n = 100$ 
  - at most  $\log_{50}(1,000,000) = 4$  nodes are accessed in a lookup.
- ◆ Contrast this with a balanced binary tree with 1 million search key values — around 20 nodes are accessed in a lookup
  - above difference is significant since every node access may need a disk I/O, costing around 20 milliseconds

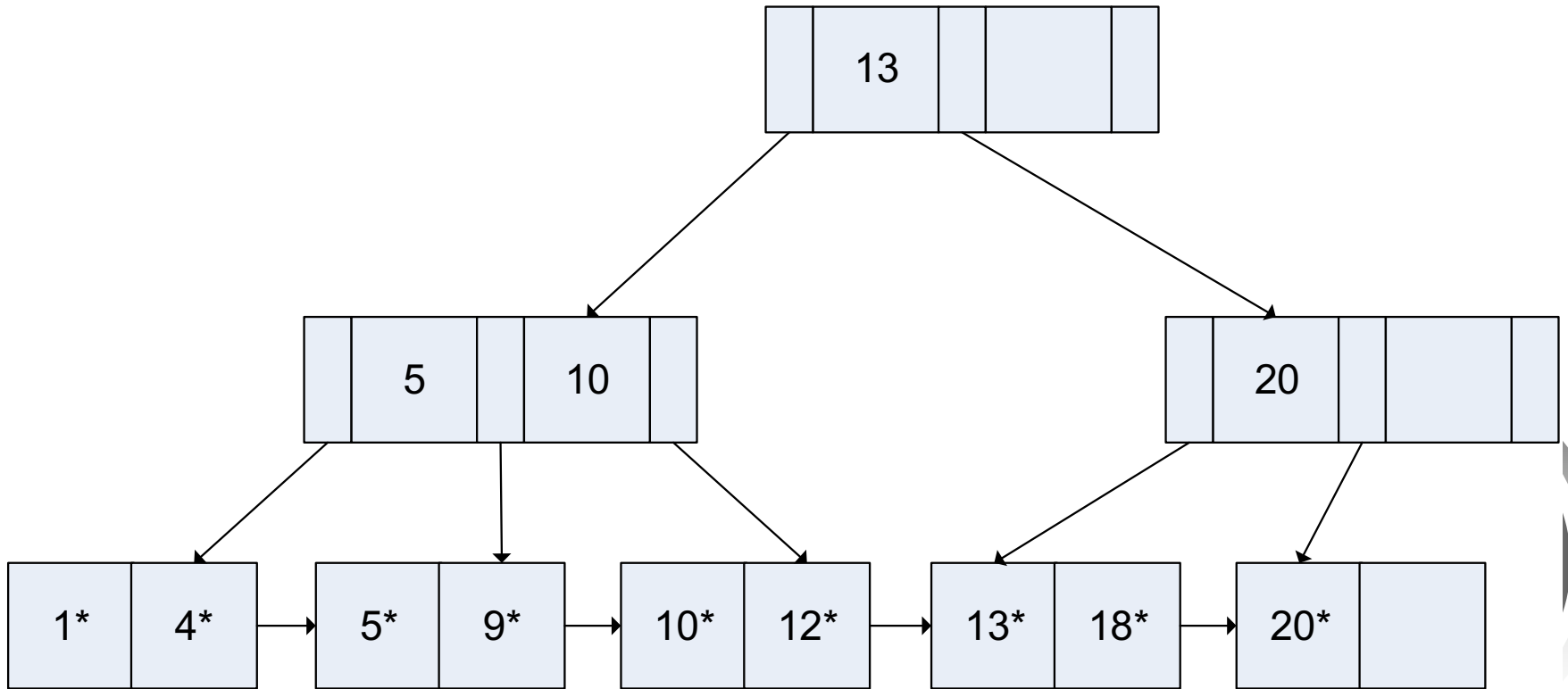


# Insertion Algorithm

Leaf page full?	Index page full?	Action
No	No	Place the record in sorted position in the appropriate leaf page
Yes	No	<ol style="list-style-type: none"><li>1. Split the leaf page including the inserted record.</li><li>2. Place Middle Key in the index page in sorted order.</li><li>3. Left leaf page contains records with keys below the middle key.</li><li>4. Right leaf page contains records with keys equal to or greater than the middle key.</li></ol>
Yes	Yes	<ol style="list-style-type: none"><li>1. Split the leaf page including the inserted record.</li><li>2. Records with keys <math>&lt;</math> middle key go to the left leaf page.</li><li>3. Records with keys <math>\geq</math> middle key go to the right leaf page.</li><li>4. Split the index page.</li><li>5. Keys <math>&lt;</math> middle key go to the left index page.</li><li>6. Keys <math>&gt;</math> middle key go to the right index page.</li><li>7. The middle key goes to the next (higher level) index.</li></ol> <p>If the next level index page is full, continue splitting the index pages.</p>

# Insertion

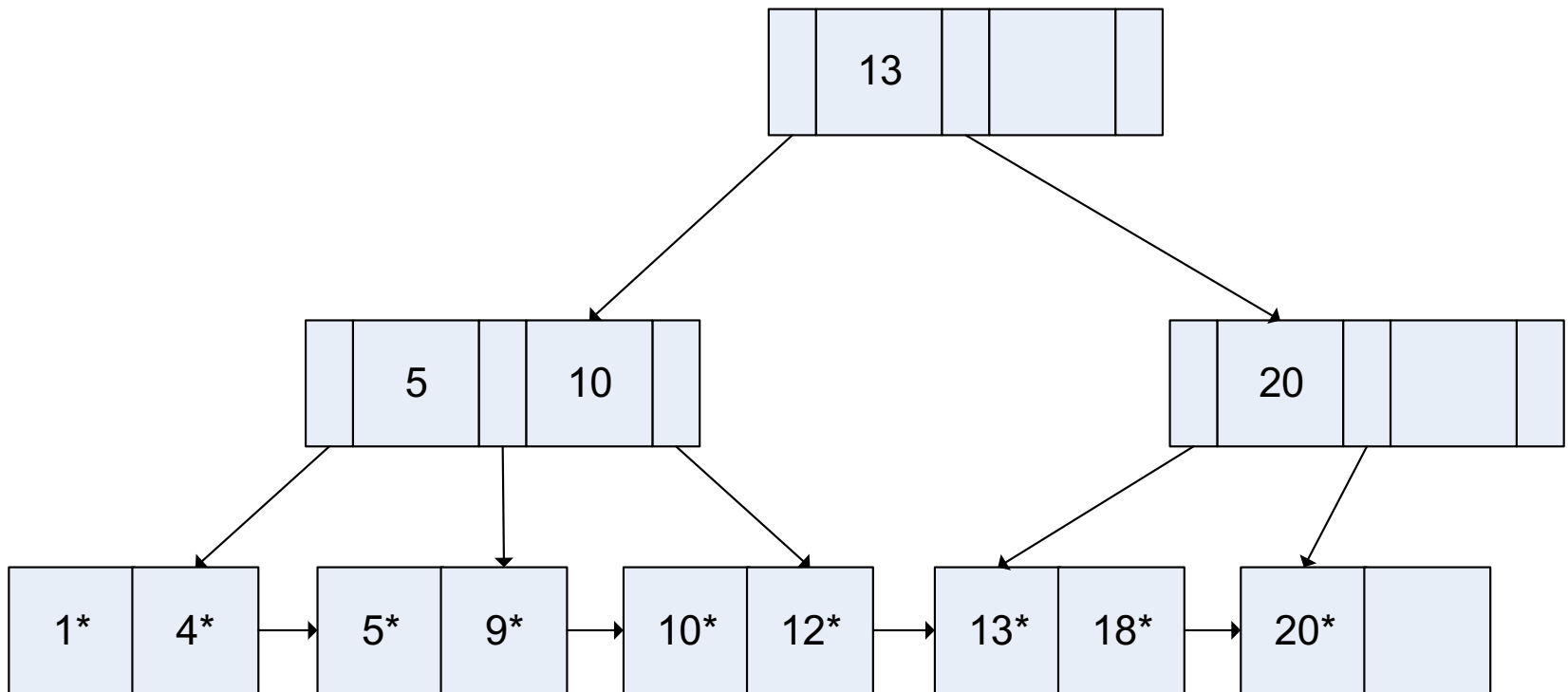
Examples of insertion with B+ tree with minimum size = 1.  
Starting with a tree looking like this:





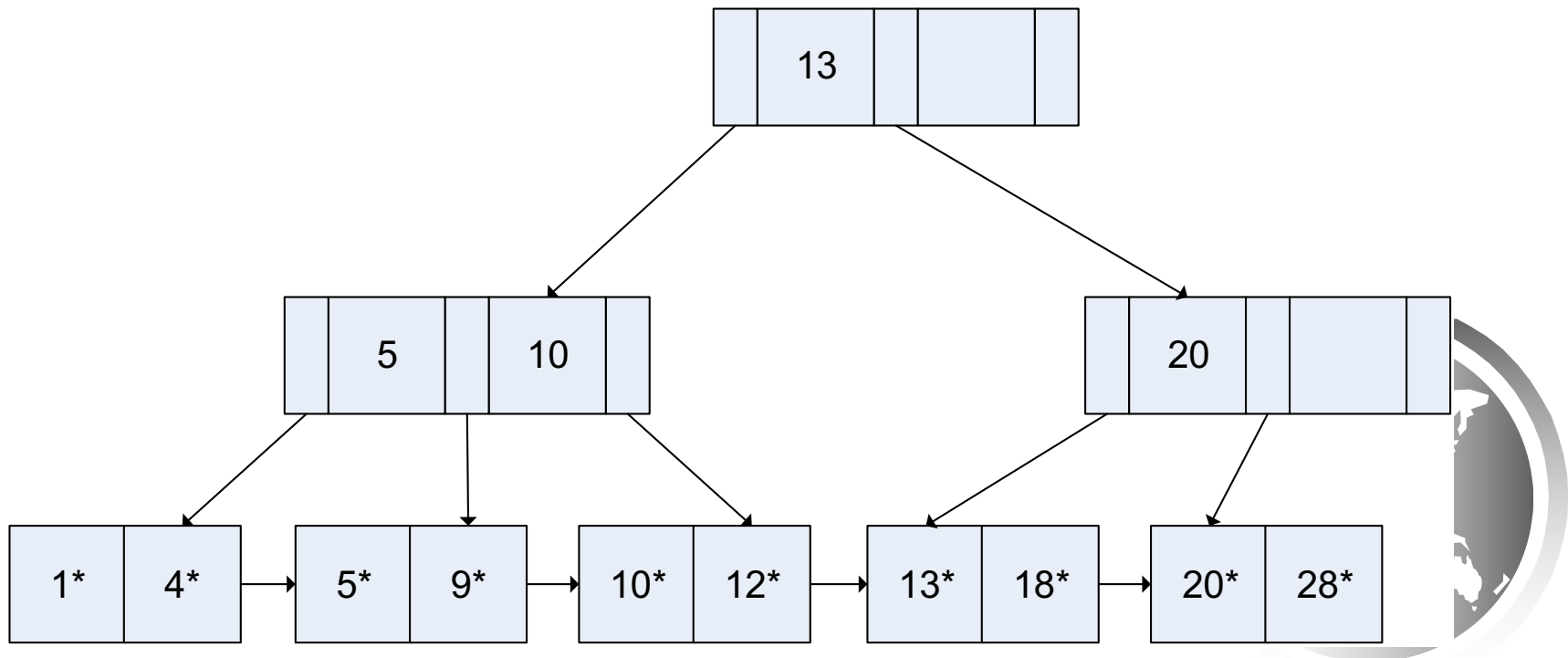
# Insert 28

Insert 28 into the below tree:



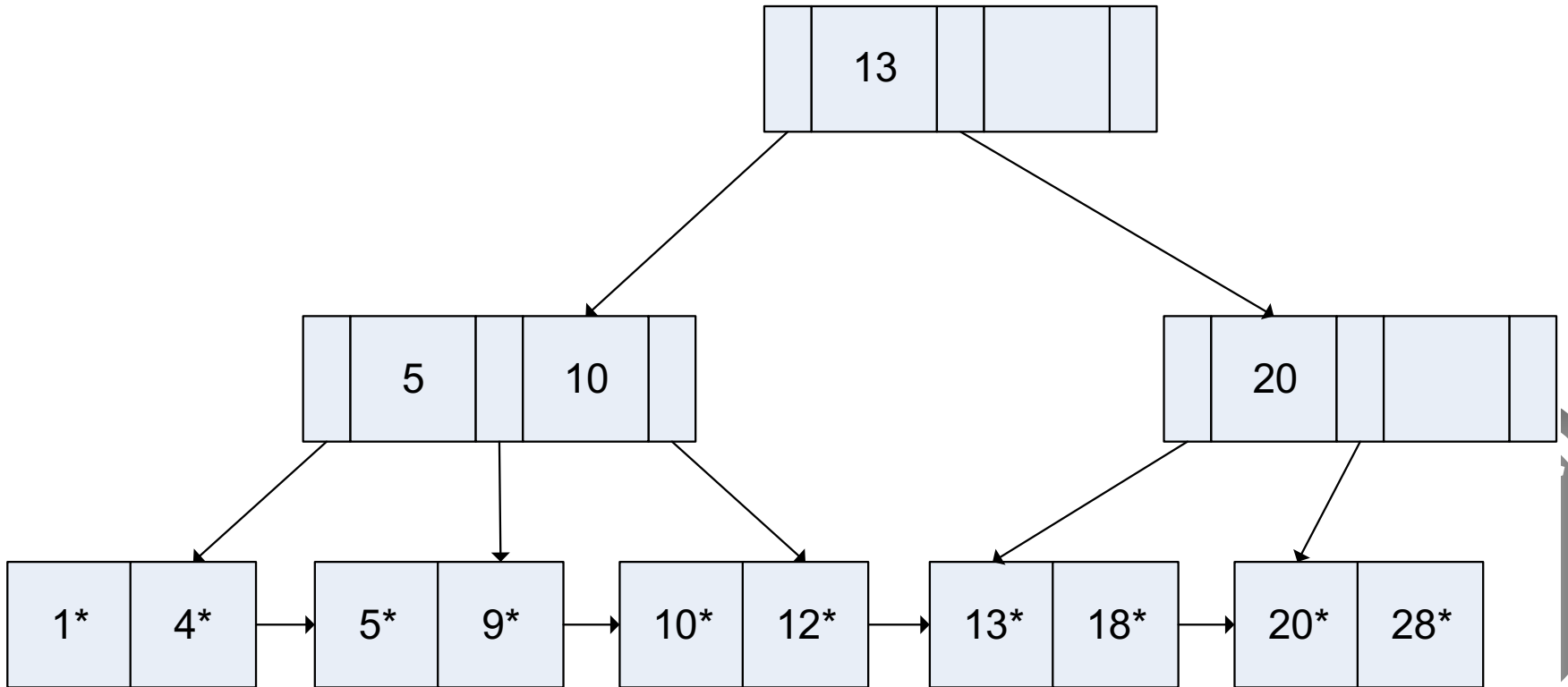
# Insert 28

- Our first insertion has an index of 28.
- We look at the leaf node to see if there is room.
- Finding an empty slot, we place the index in node in sorted order.



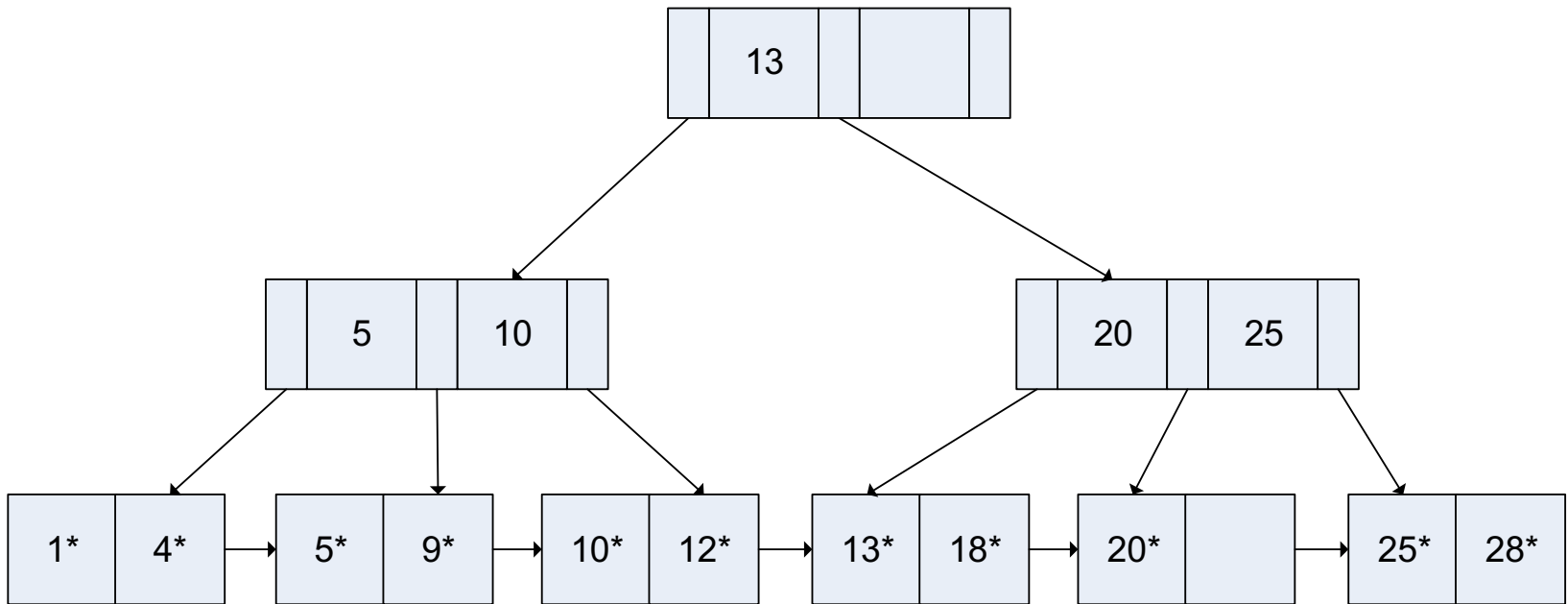
# Insert 25

Insert 25 in the below tree:



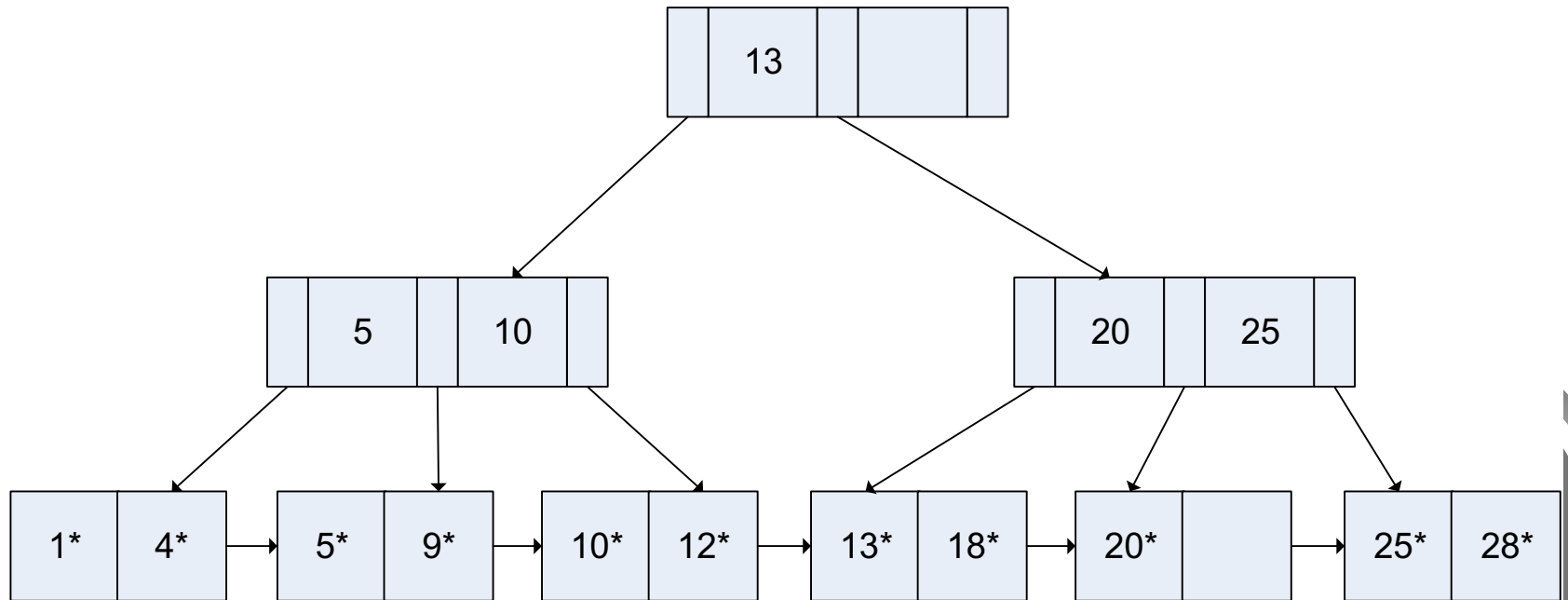
# Insert 25

- Our next insertion is at 25.
- We look at the leaf node it would go in and find there is no room.
- We split the node, and roll the middle value to the index mode above it.



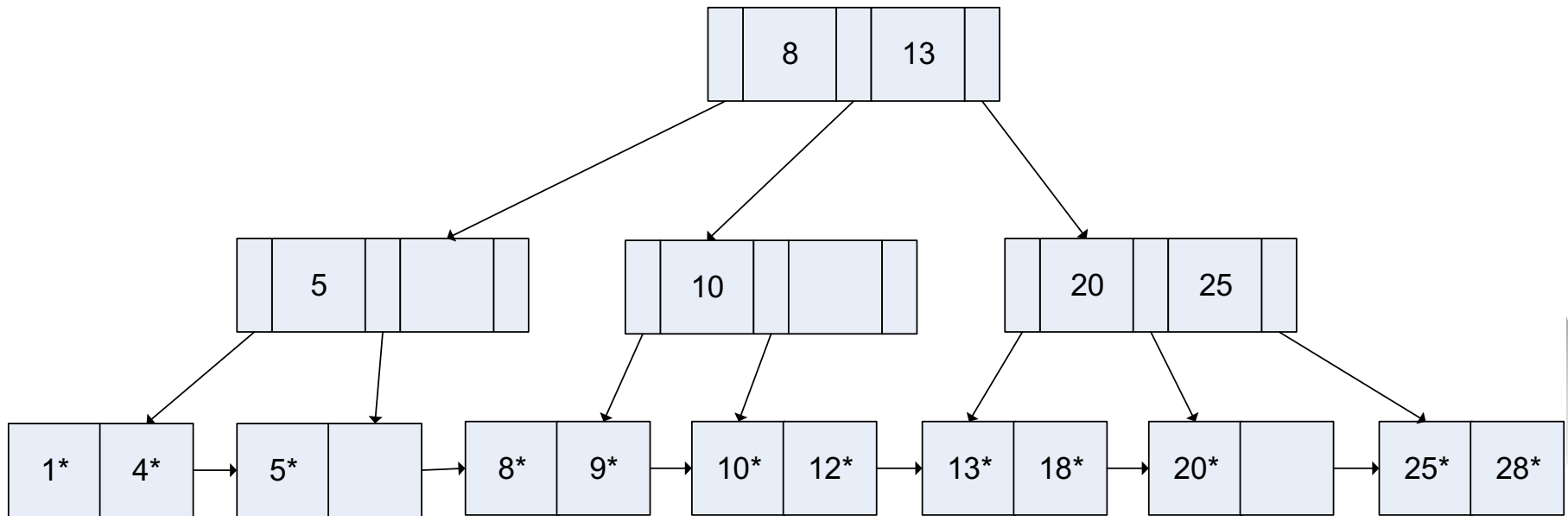
# Insert 8

Insert 8 in the below tree:



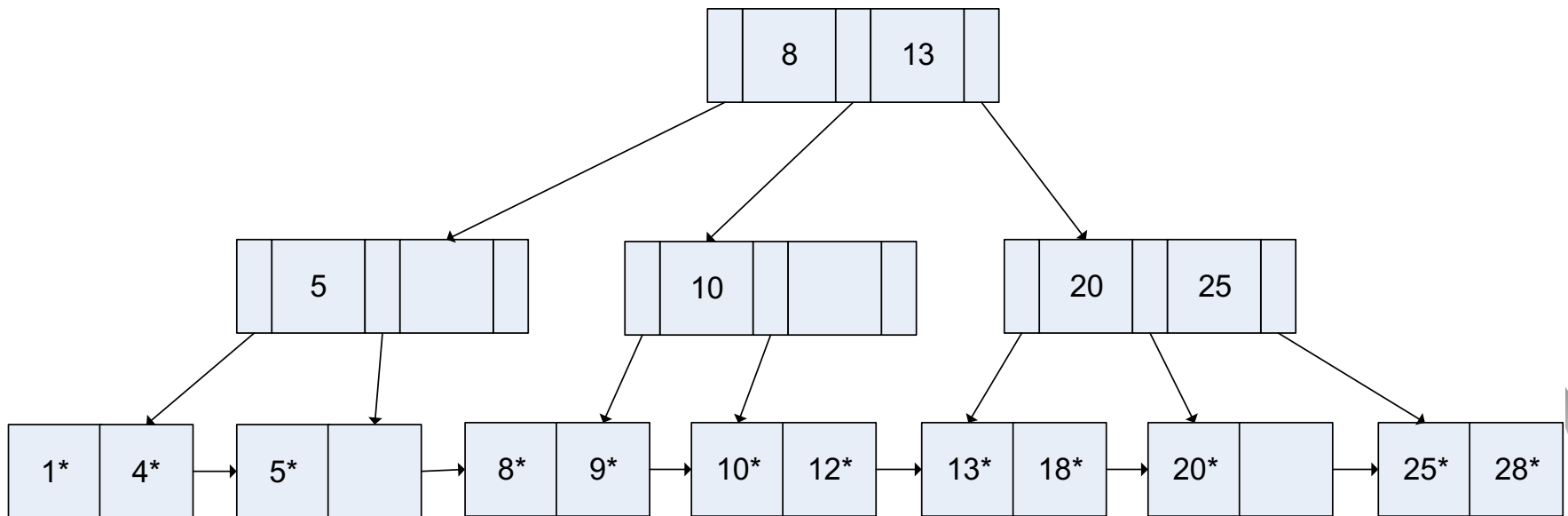
# Insert 8

- Our next case occurs when we want to add 8.
- The leaf node is full,
- so we split it and attempt to roll the index to the index node.
- It is full, so we must split it as well.



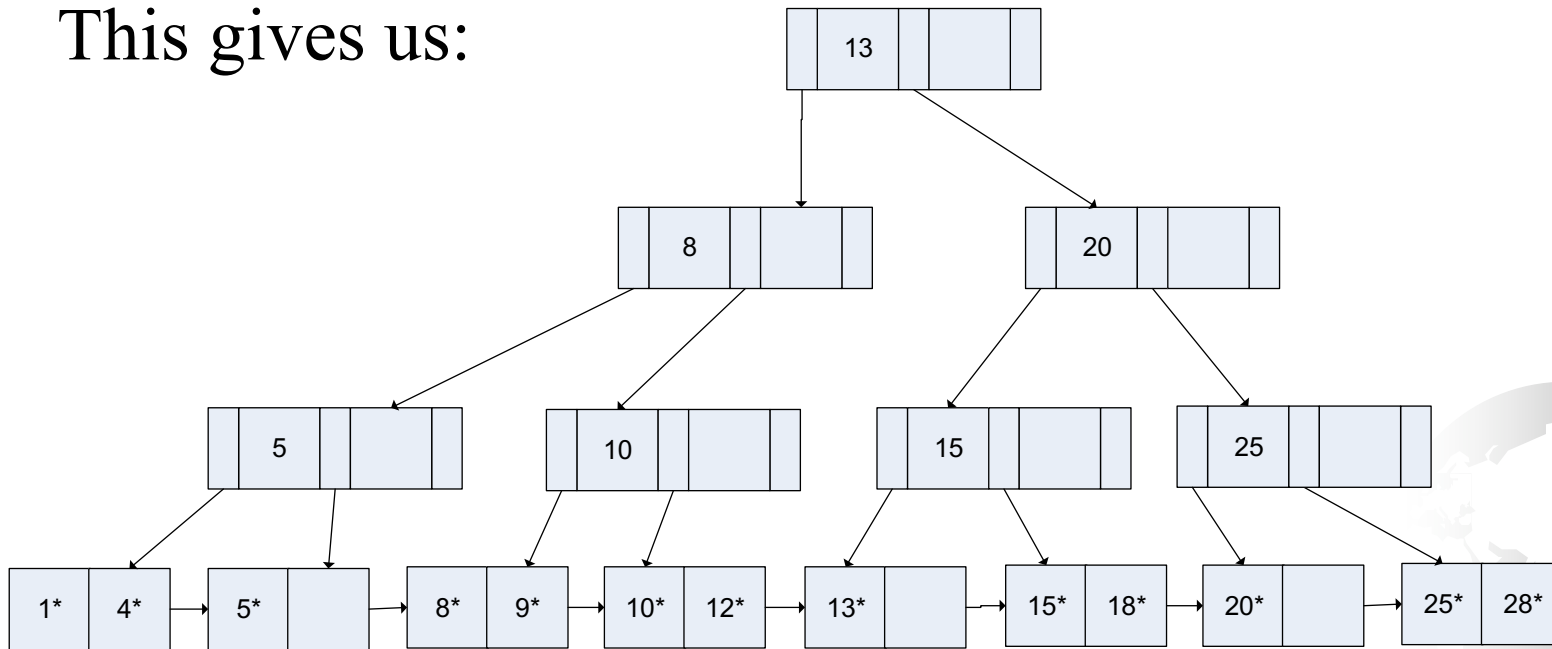
# Insert 15

Insert 15 into the below tree:



# Insert 15

- Our last case occurs when we want to add 15.
- This is going to result in the root node being split.
- The leaf node is full, as are the two index nodes above it.
- This gives us:



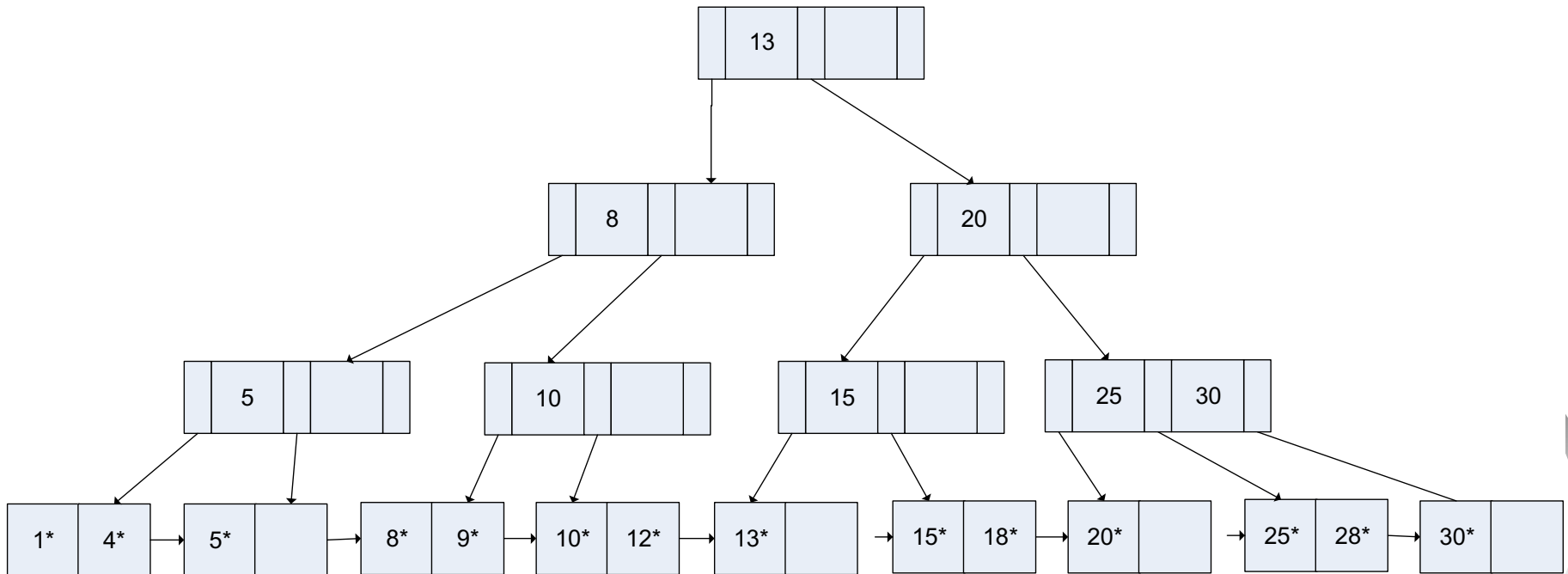


# Delete Algorithm

Leaf page less than 1/2?	Index page less than 1/2?	Action
No	No	<ol style="list-style-type: none"><li>1. Delete the record from the leaf page.</li><li>2. Arrange keys in ascending order to fill void.</li><li>3. If the key of the deleted record appears in the index pages, use the next key to replace it.</li></ol>
Yes	No	<ol style="list-style-type: none"><li>1. Combine the leaf page and its sibling.</li><li>2. Change the index page to reflect the change.</li></ol>
Yes	Yes	<ol style="list-style-type: none"><li>1. Combine the leaf page and its sibling.</li><li>2. Adjust the index page to reflect the change.</li><li>3. Combine the index page with its sibling</li><li>4. Delete entry from parent</li></ol> <p>Continue combining index pages until you reach a page with the correct fill factor or you reach the root page.</p>

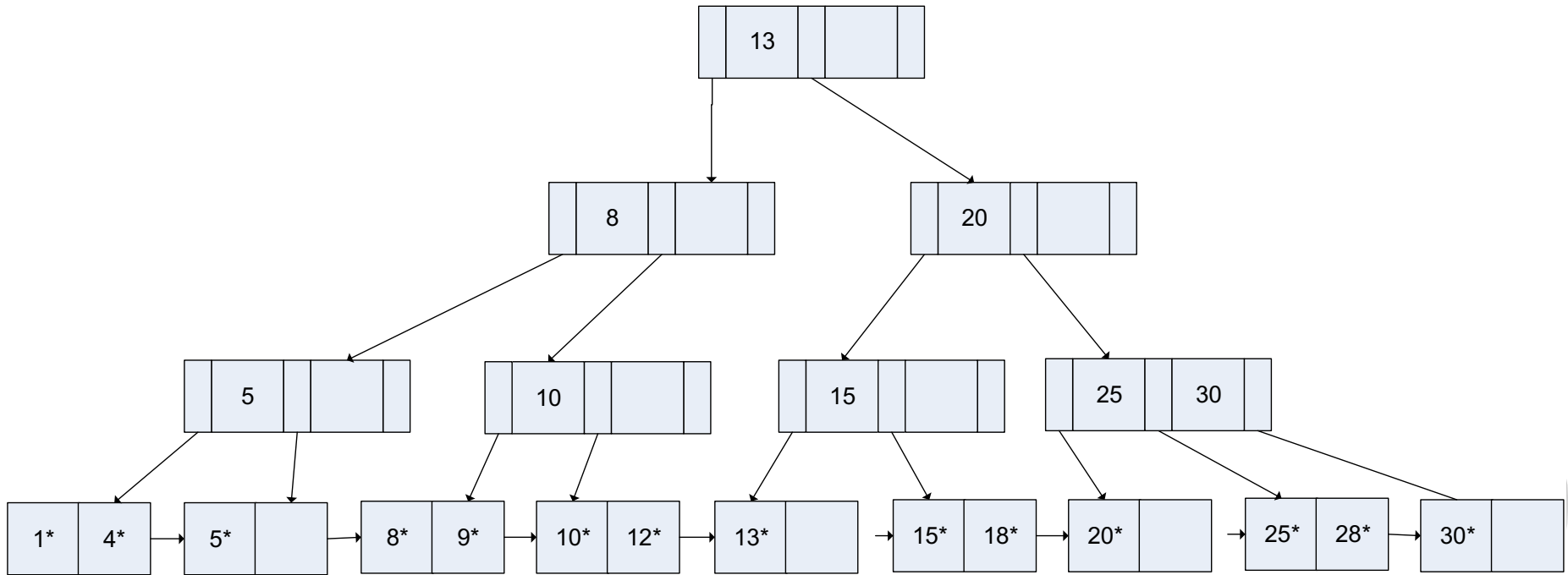
# Delete

Let's take our tree from the insert example with a minor modification. We have added 30 to give us an index node with 2 indexes in it:



# Delete 18

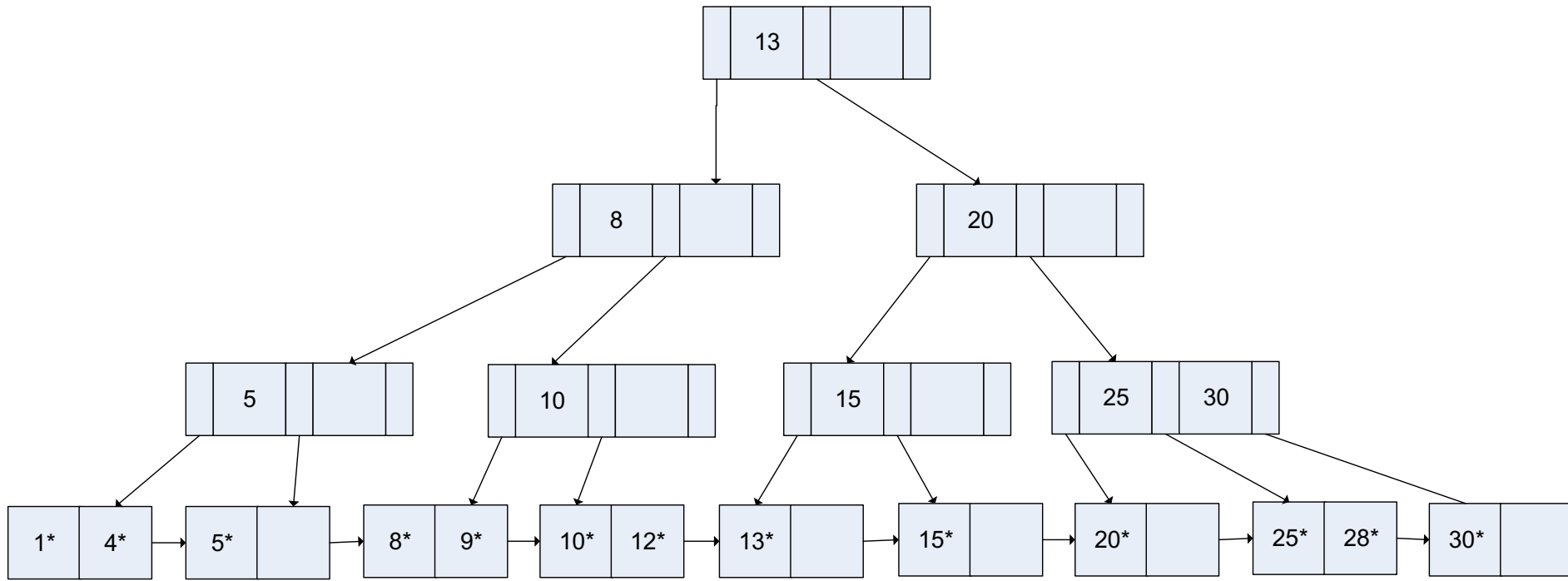
Delete 18 from the below tree:



# Delete 18

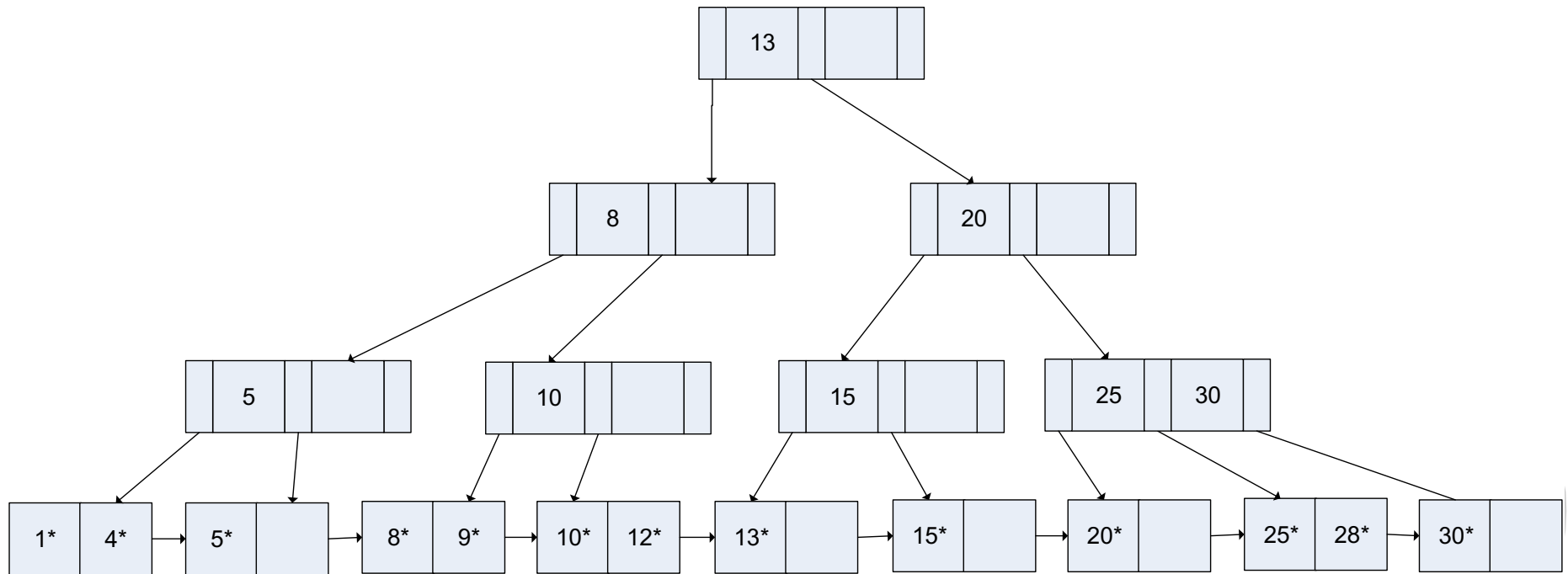
Our first delete is of 18.

Simplest case is that it is not an index and in a leaf node that deleting it will not take you below 1/2.



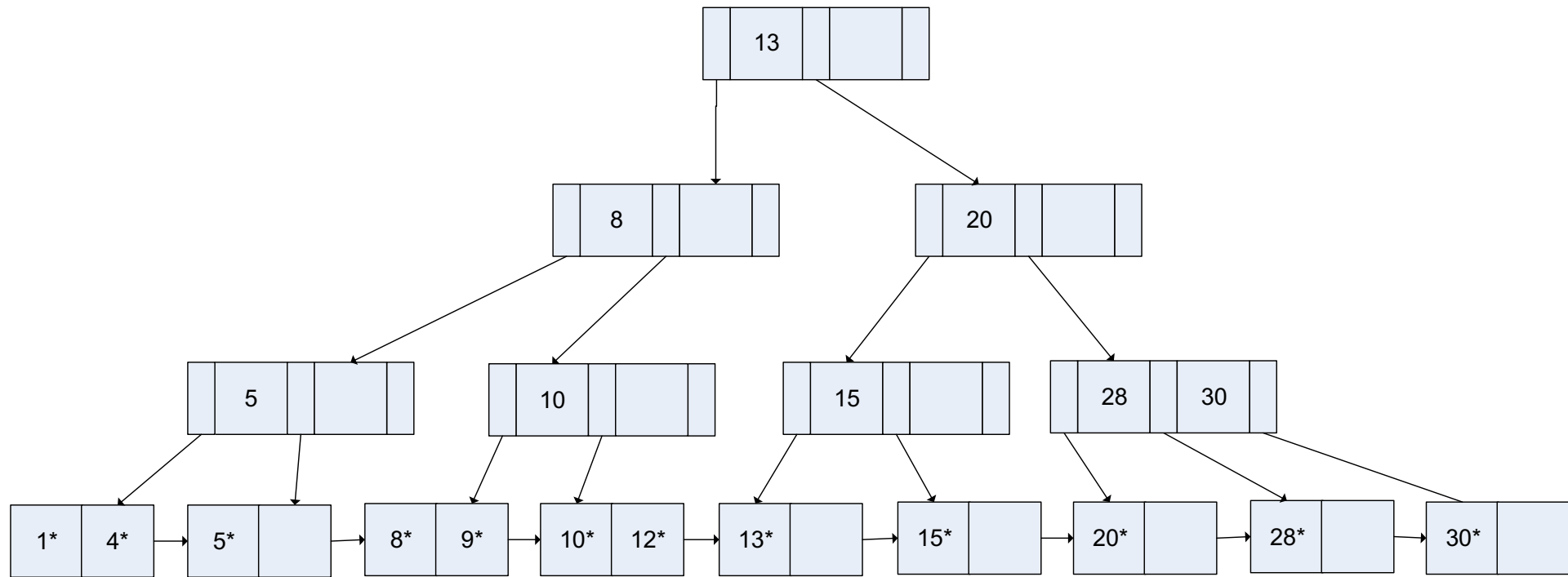
# Delete 25

Delete 25 from the below tree:



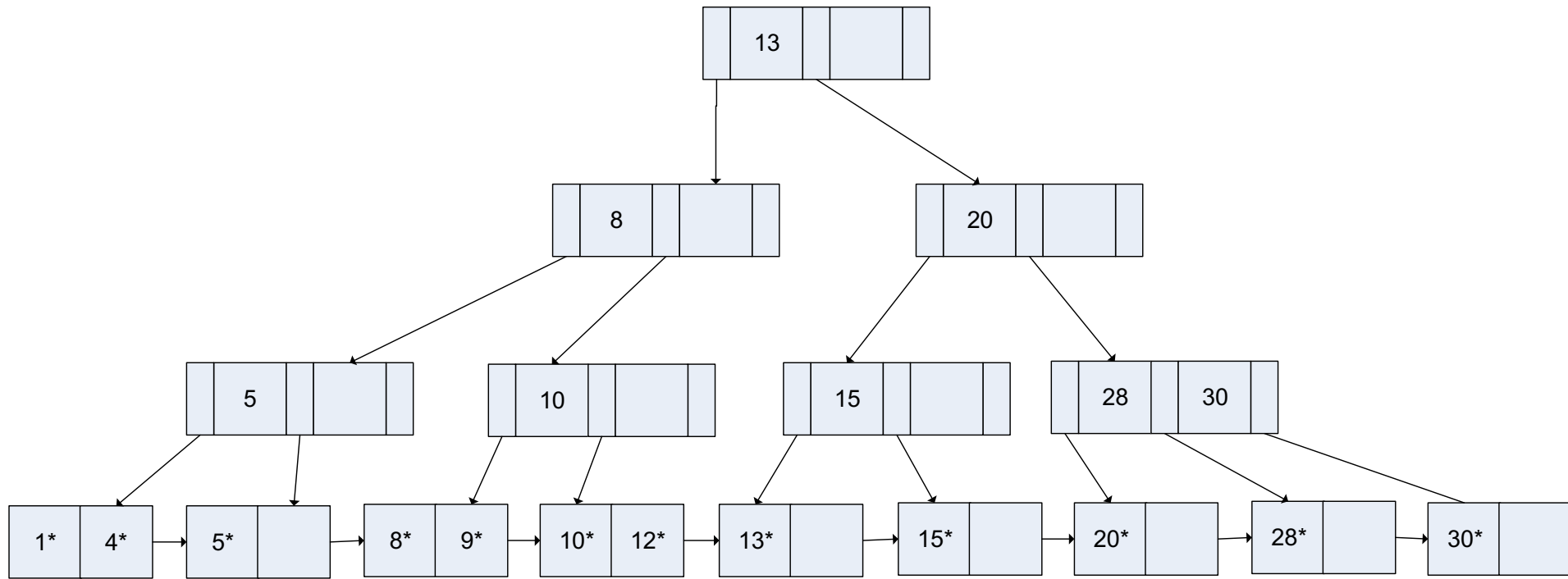
# Delete 25

This delete is similar, except the index appears in an index node. In that case, the next index replaces the one in the index node. Let's delete 25 from the previous slide.



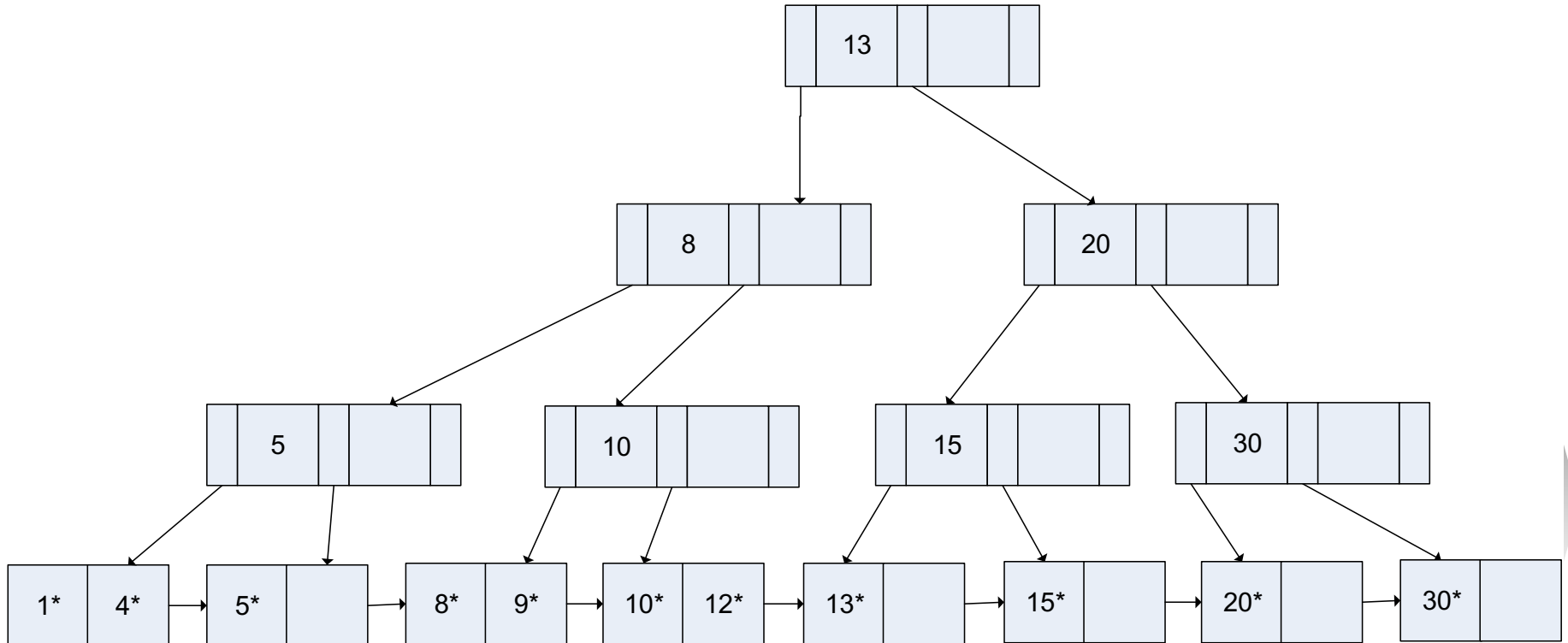
# Delete 28

Delete 28 from the below tree:



# Delete 28

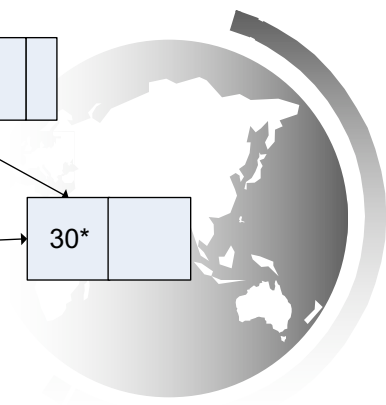
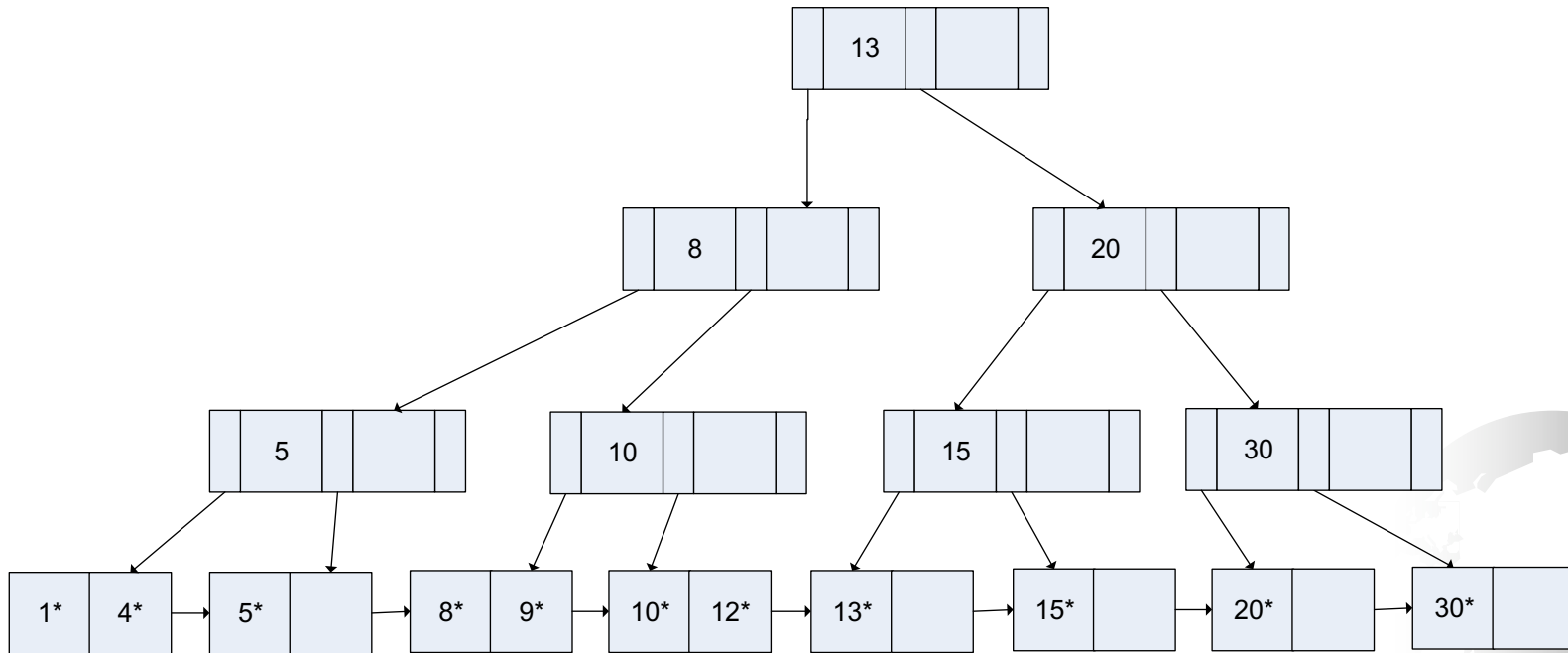
Our next case takes the node below d. Let's delete 28. We combine the leaf page (in our case it is empty) with its sibling and update the index appropriately. That gives us:





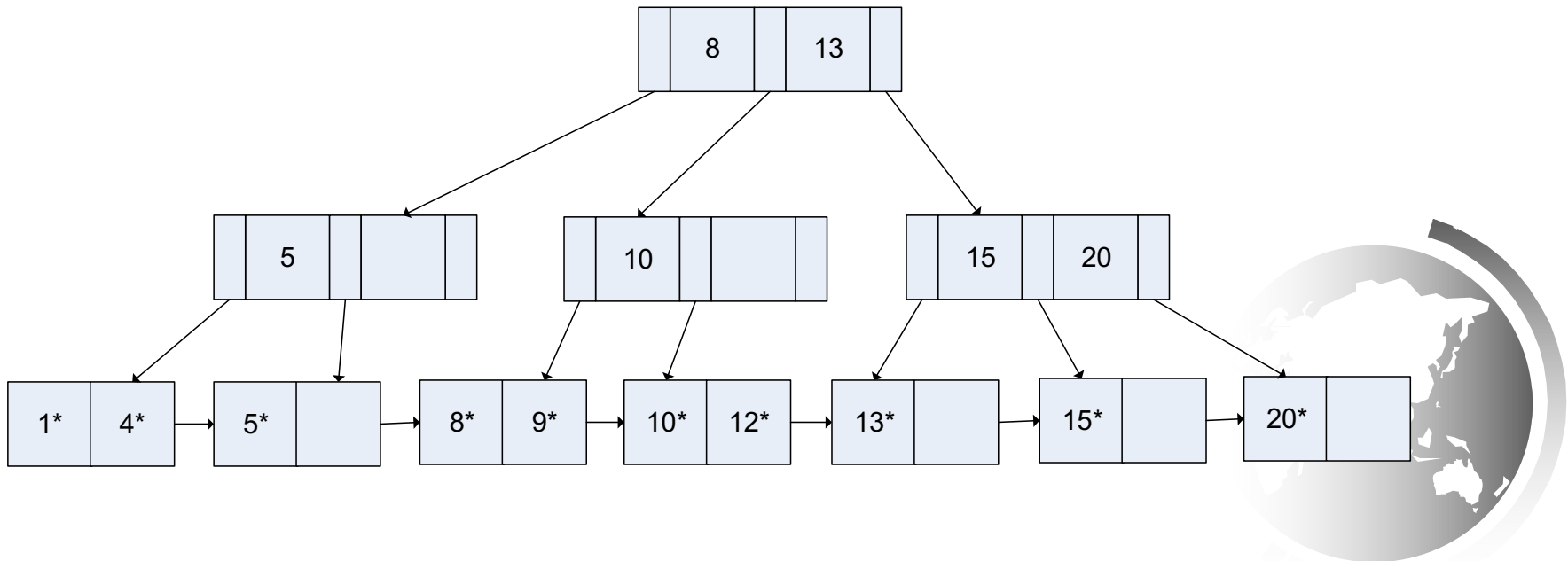
# Delete 30

Delete 30 from the below tree:



# Delete 30

- Next we delete 30.
- This takes us below d for the index.
- We combine the indexes, which has the effect of taking the index above below d.
- This continues to the root.



# Delete 30

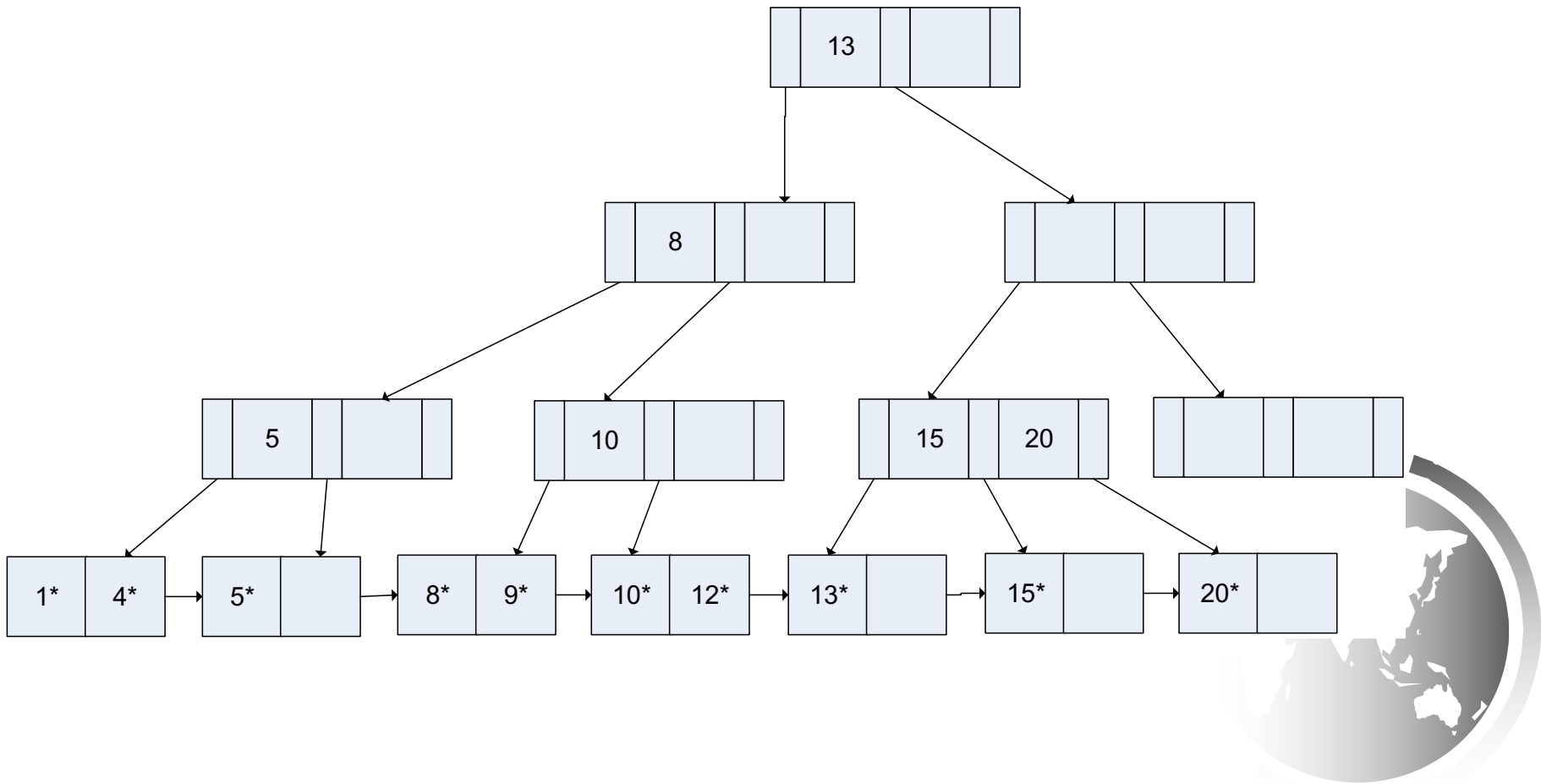
Woah. That seemed like magic. What process got us to that?

Ok – let's go through it.

1. When we deleted 30, it took the data entry node that 30 was in below d.
2. Now we have to merge with the sibling.
3. When we merge – it's to the sibling on the left, which means pointer in the index above is no longer valid.
4. We remove it, (which leaves it less than d), pull down the index from above and merge the index node with its sibling.

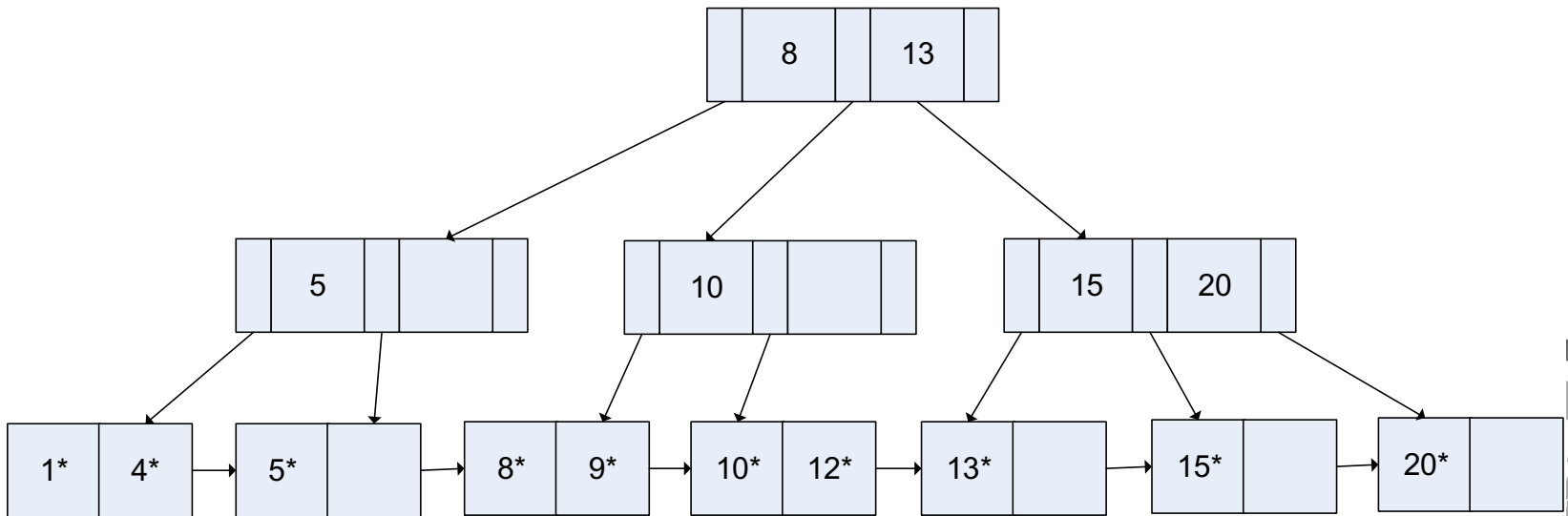


# Delete 30



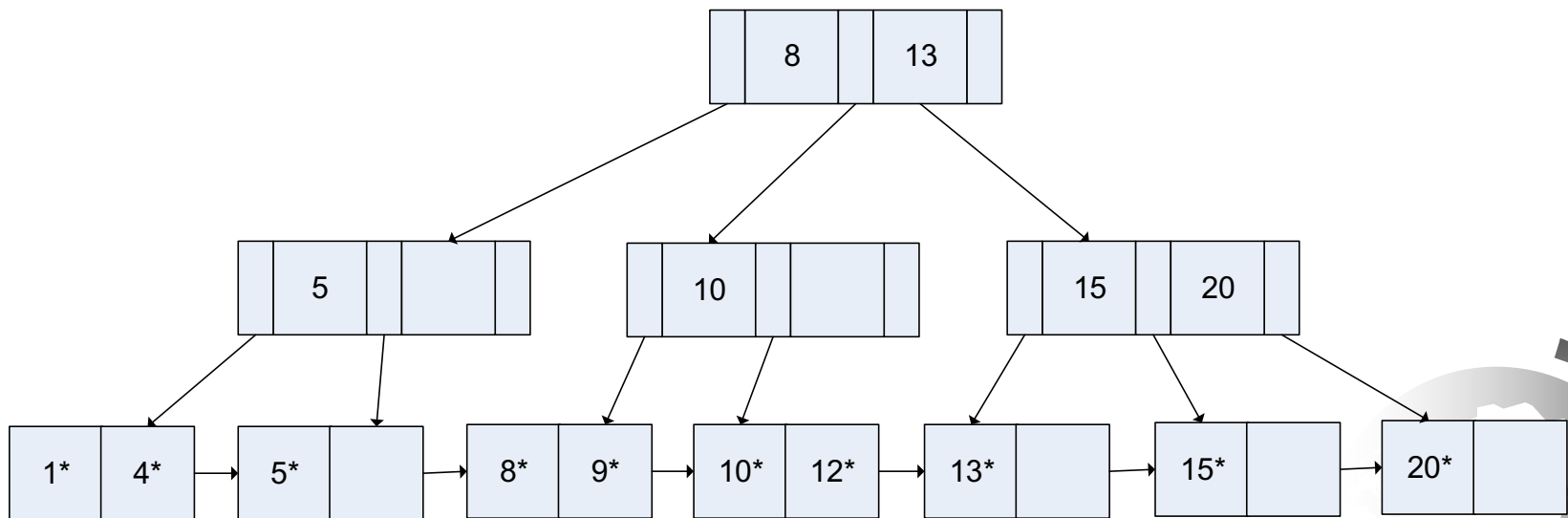
# Delete 30

Repeating the process gets us back to



# Delete 5

Delete 5 from the below tree:

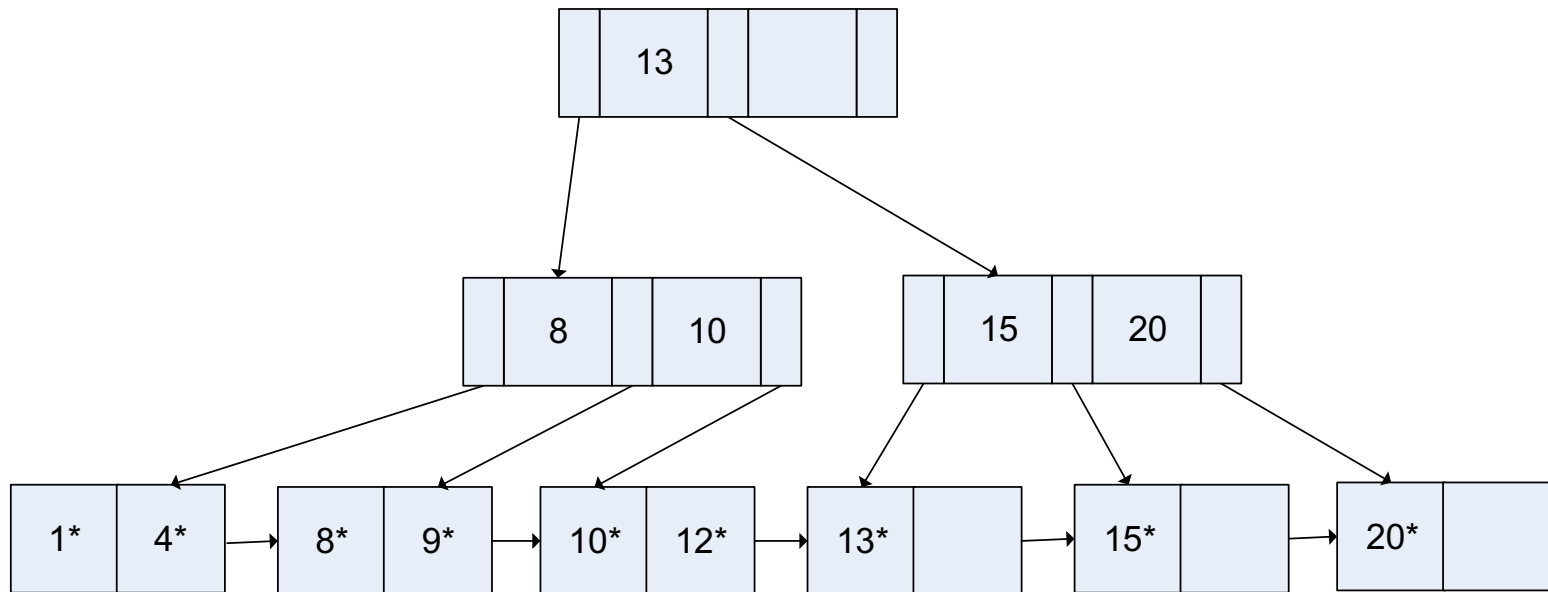


# Delete 5

Our last example deletes 5.

This takes the node and the index above it below d.

We remove the leaf node and combine the index with its neighbor.

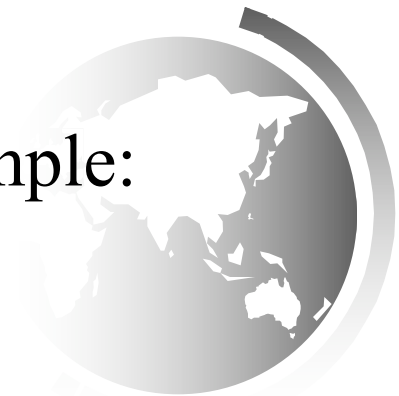


# Rotation

It is also possible to rebalance a tree to reduce the number of splits – called rotation.

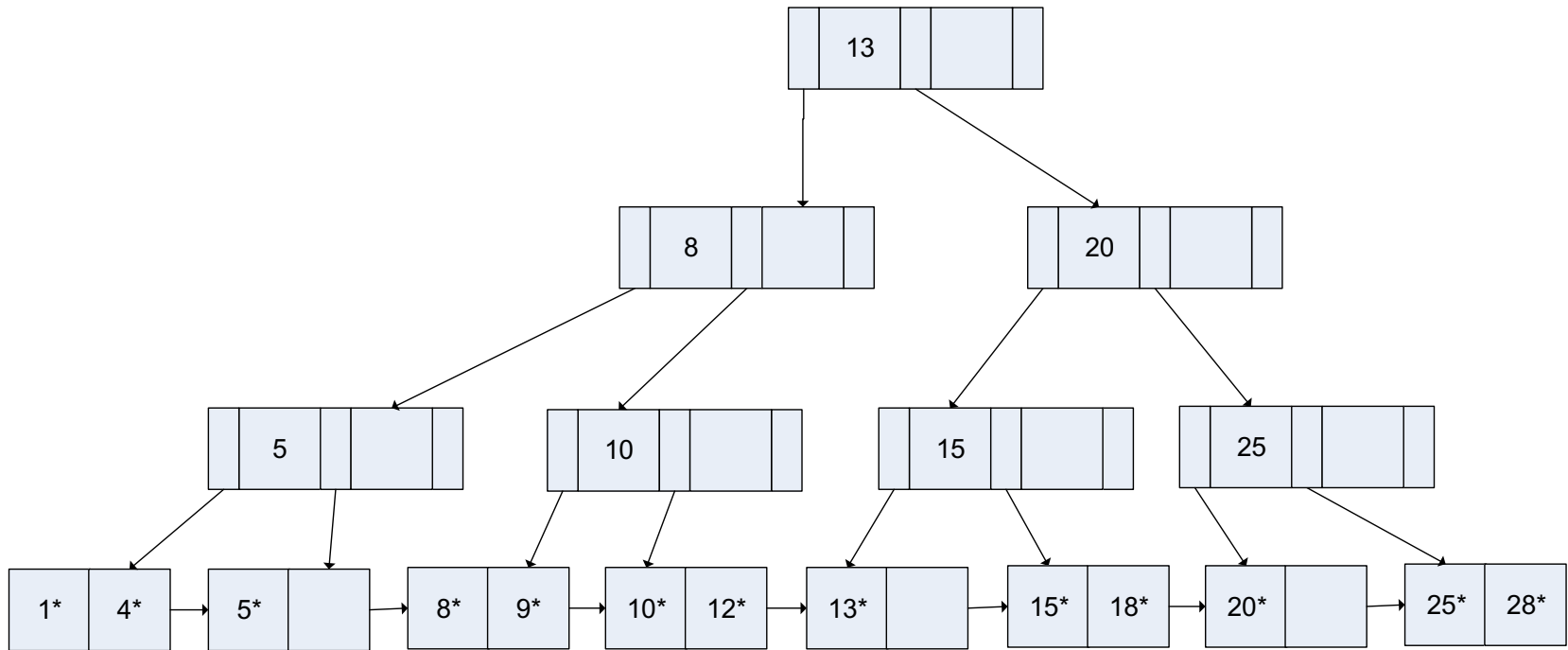
If you are trying to insert, and a leaf page is full, but its sibling isn't – you can move an index to a sibling and avoid splitting.

Let's go back to a tree from our insert example:





# Add 3

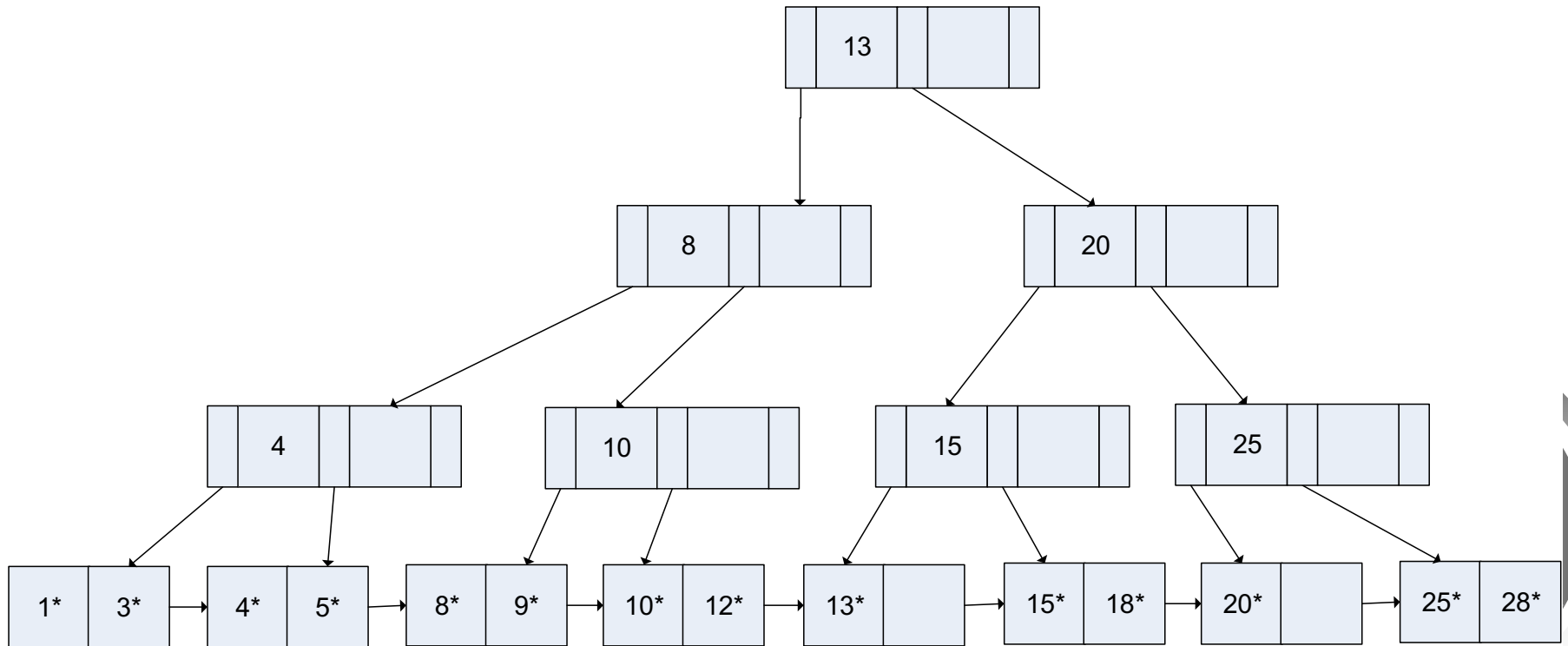


# Add 3 (Rotation)

We check the sibling to see if it has room.

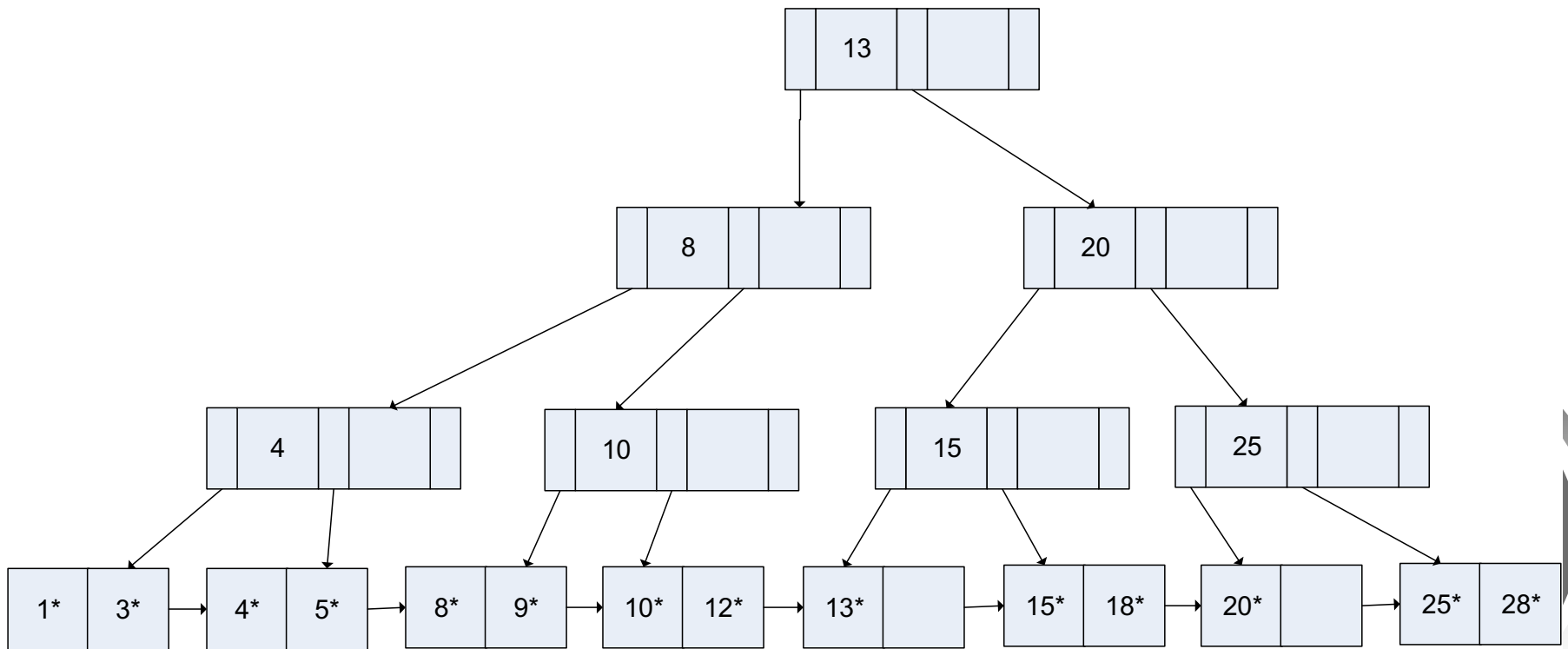
It does, so we move a record to it adjusting the index.

Now we have :

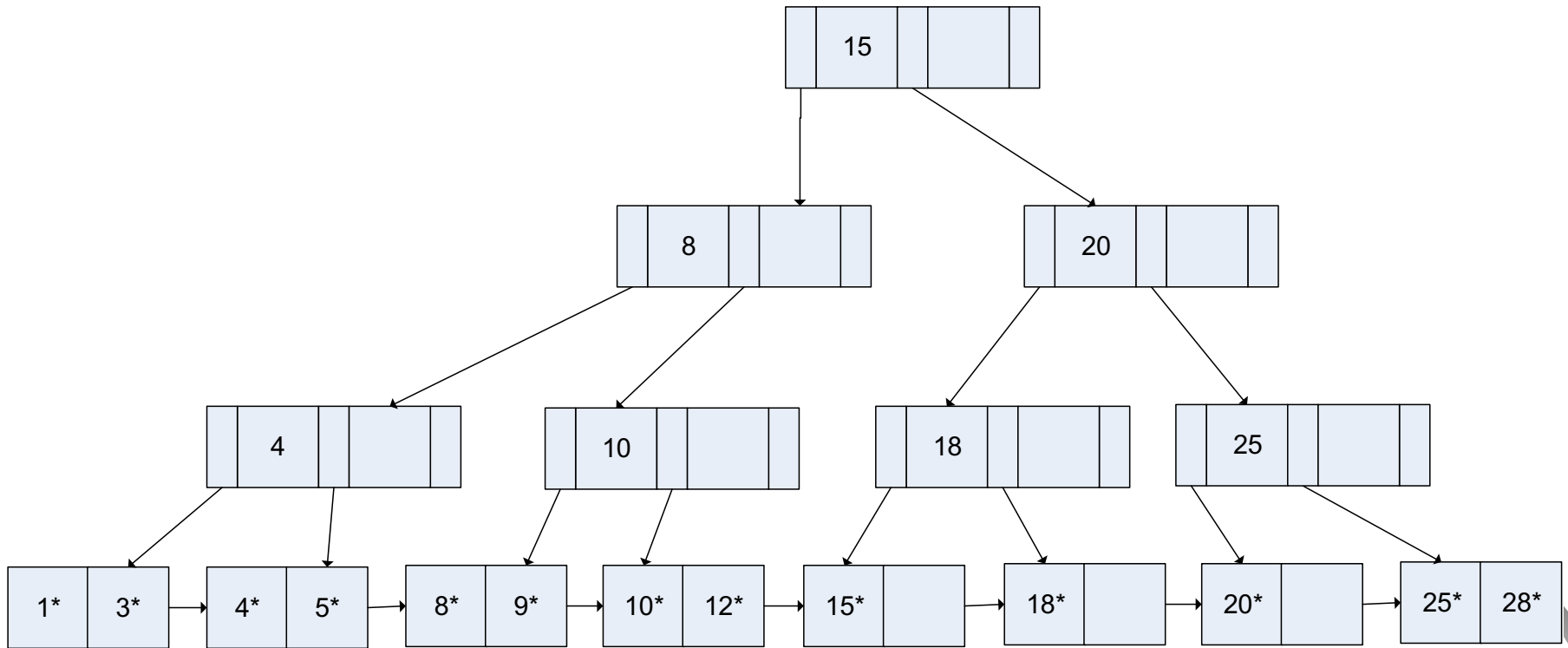


# Delete 13 (Rotation)

The same concept works with deletes.

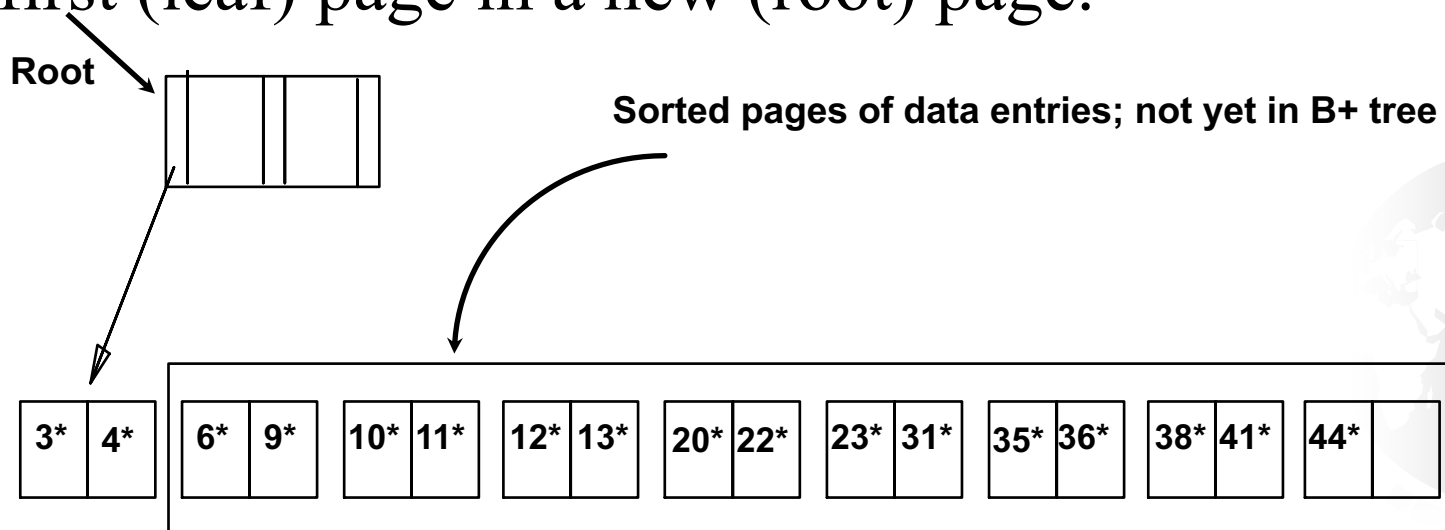


# Delete 13 (Rotation)



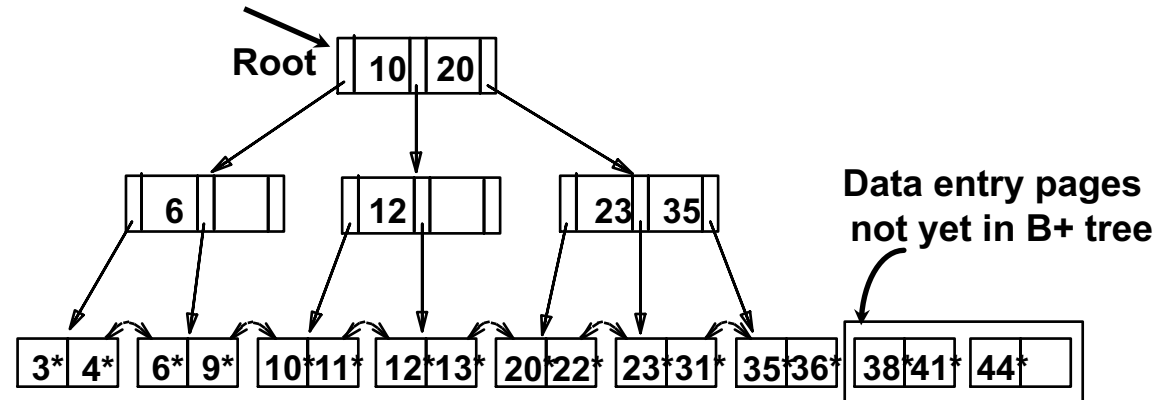
# Bulk Loading of a B+ Tree

- ◆ If we have a large collection of records, and we want to create a B+ tree on some field, doing so by repeatedly inserting records is very slow.
- ◆ Bulk Loading can be done much more efficiently.
- ◆ *Initialization*: Sort all data entries, insert pointer to first (leaf) page in a new (root) page.

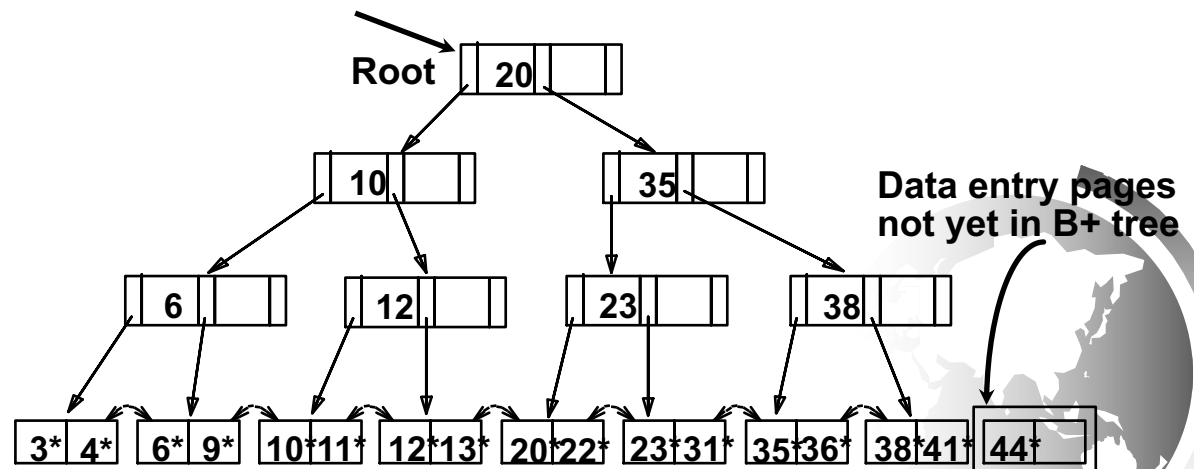


# Bulk Loading (Contd.)

◆ Index entries for leaf pages always entered into right-most index page just above leaf level.



◆ When this fills up, it splits. (Split may go up right-most path to the root.)



◆ Much faster than repeated inserts.

# Summary of Bulk Loading

- ◆ Option 1: multiple inserts.
  - Slow.
  - Does not give sequential storage of leaves.
- ◆ Option 2: Bulk Loading
  - Has advantages for concurrency control.
  - Fewer I/Os during build.
  - Leaves will be stored sequentially (and linked, of course).
  - Can control “fill factor” on pages.

