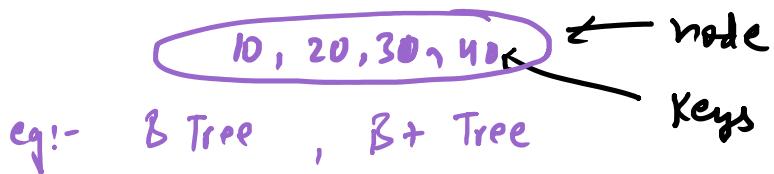


External Searching :- To locate or retrieve a record which is stored in a secondary memory like disk.

Multimay Search Tree :- A M-way search tree contains multiple keys in a node, which leads to reduction of the overall height of the tree.



B Tree :- Balanced Tree

A B Tree of order 'm' is a multiway search tree with the following properties:-

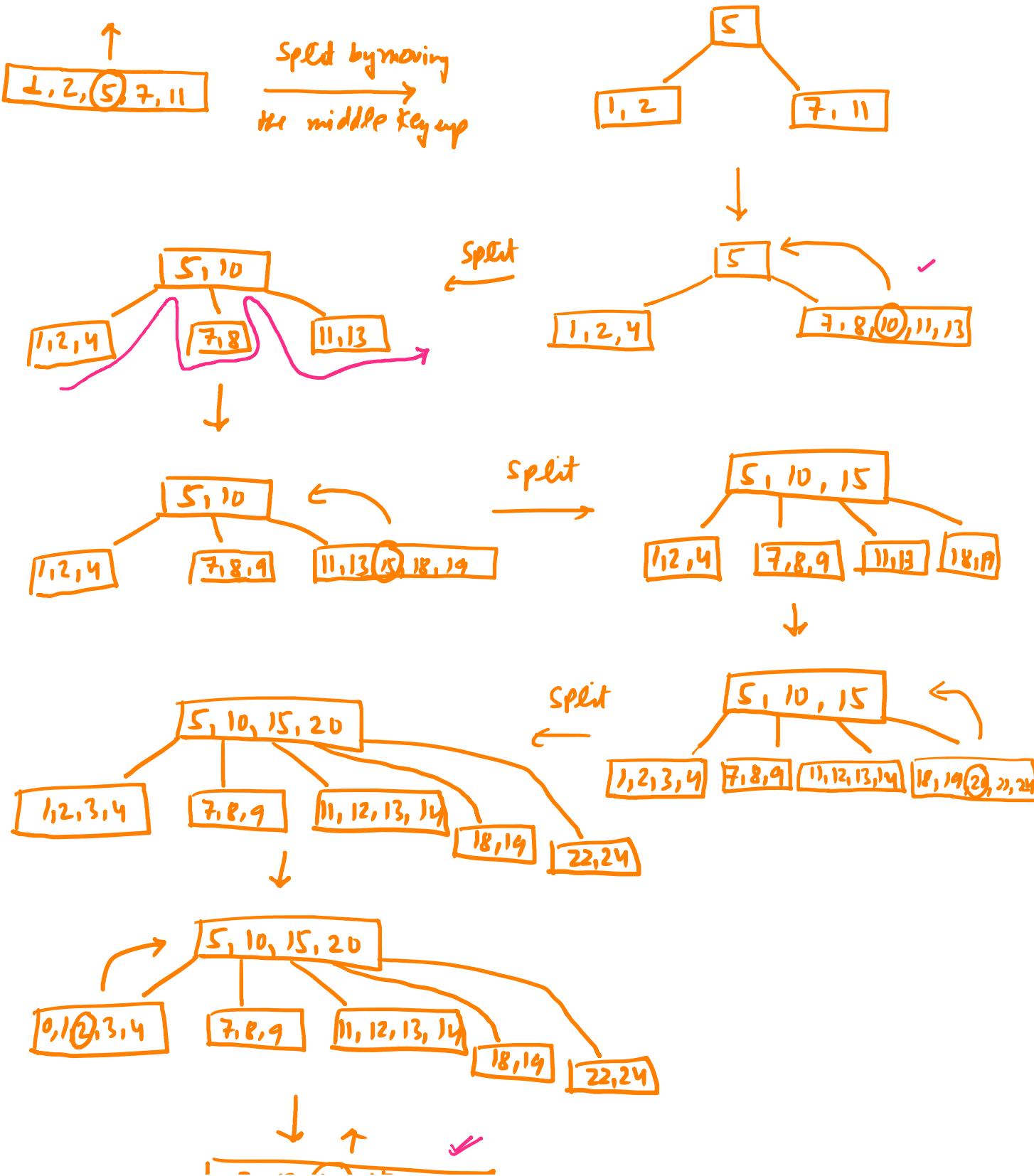
- 1) The root can have atleast 1 key & atmost $m-1$ keys
- 2) All nodes (except the root) can have atleast $\lfloor m/2 \rfloor$ keys and atmost $m-1$ keys.
- 3) All leaves are at the same depth.
- 4) If a node has 't' keys, then it will have $t+1$ children
- 5) Keys of a node are sorted in ascending order.

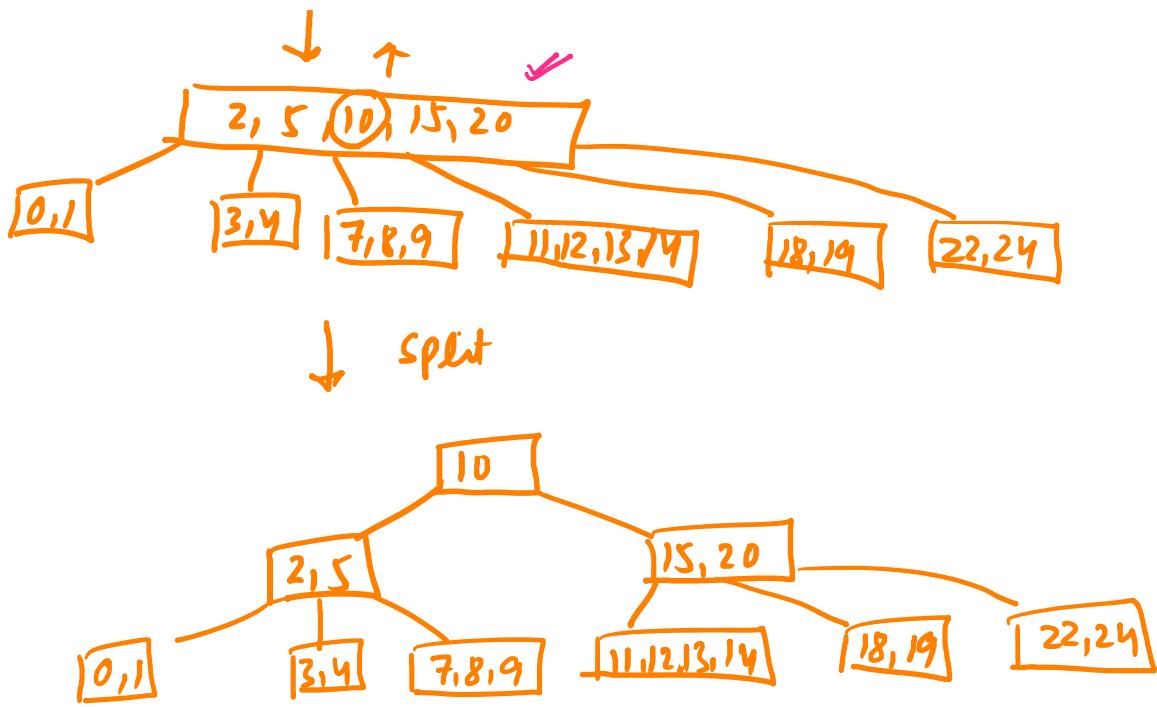


- 6) All the keys in the subtrees pointed by S_1, S_2, \dots, S_{t+1} ,
 $S_1 < K_1 < S_2 < K_2 < S_3 < \dots < S_t < K_t < S_{t+1}$

Insert the following keys in a B Tree of Order 5.

1, 7, 5, 2, 11, 4, 8, 13, 10, 15, 19, 9, 18, 24, 3, 12, 14, 20, 21, 0





Deletion :-

1) For a key in a non-leaf node \rightarrow swap the key with its inorder successor & then delete the inorder successor.

2) For a key in a leaf node

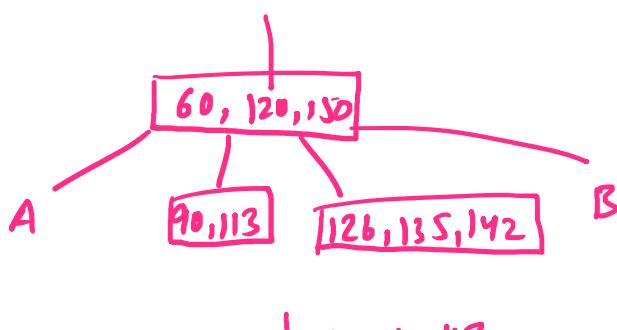
- (a) If after deletion of the key Property 2 holds then no issues
- (b) But if Property 2 does not hold then there is a issue.
Keys in that leaf note $< \lfloor m/2 \rfloor$, in case of more than

Solutions to 2-(b) :-

1) Borrowing

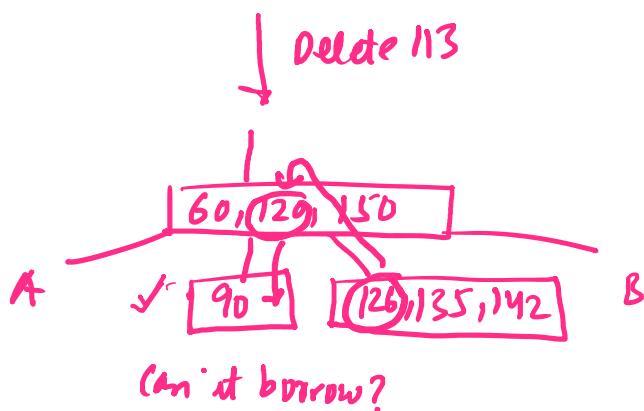
Assume $m=5$

1 nodes or a
non-root node



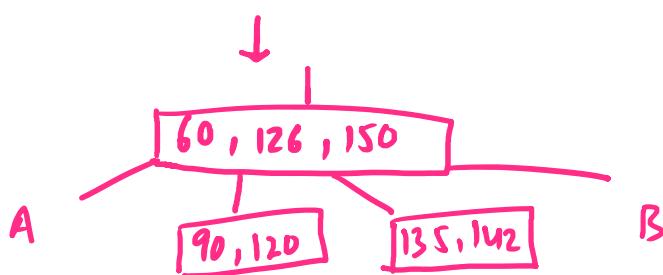
2 options:-

i) First Re-align

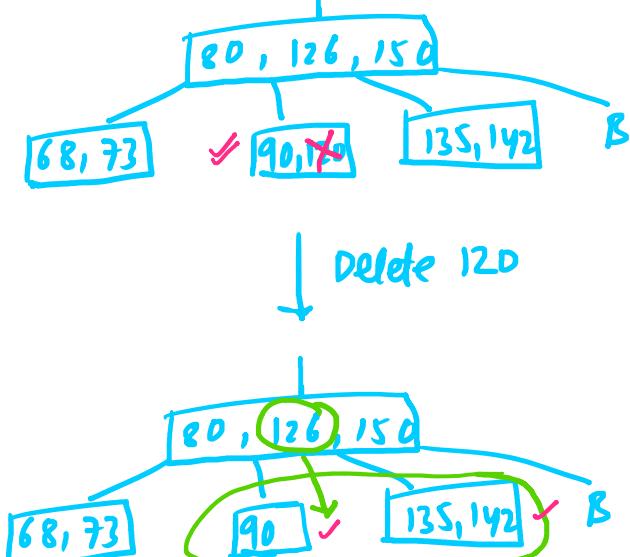


1 options:-

- 1) First key of right sibling node via parent
- 2) Last key of left sibling node via parent

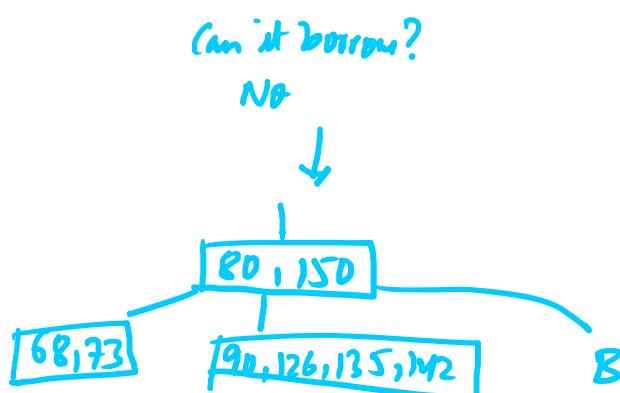


2) (Consolidation :-)



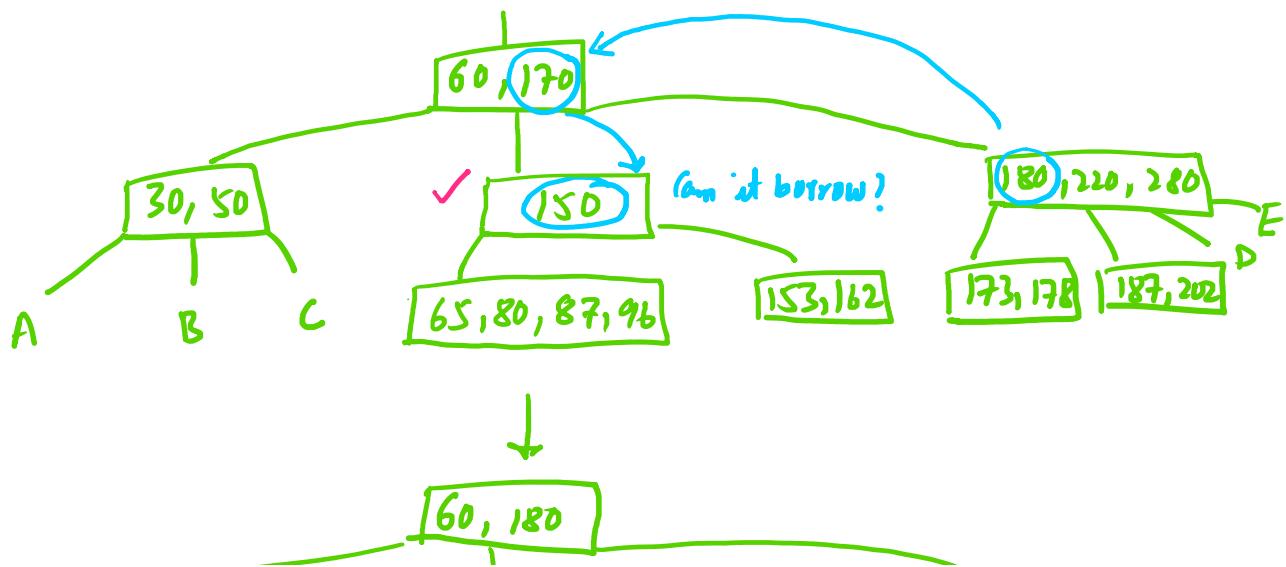
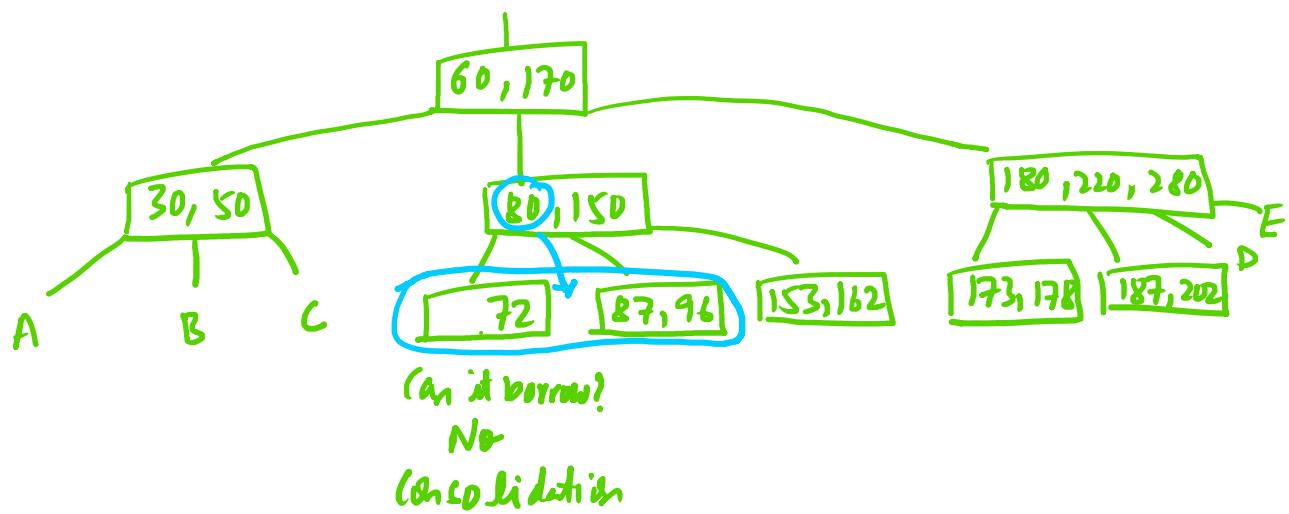
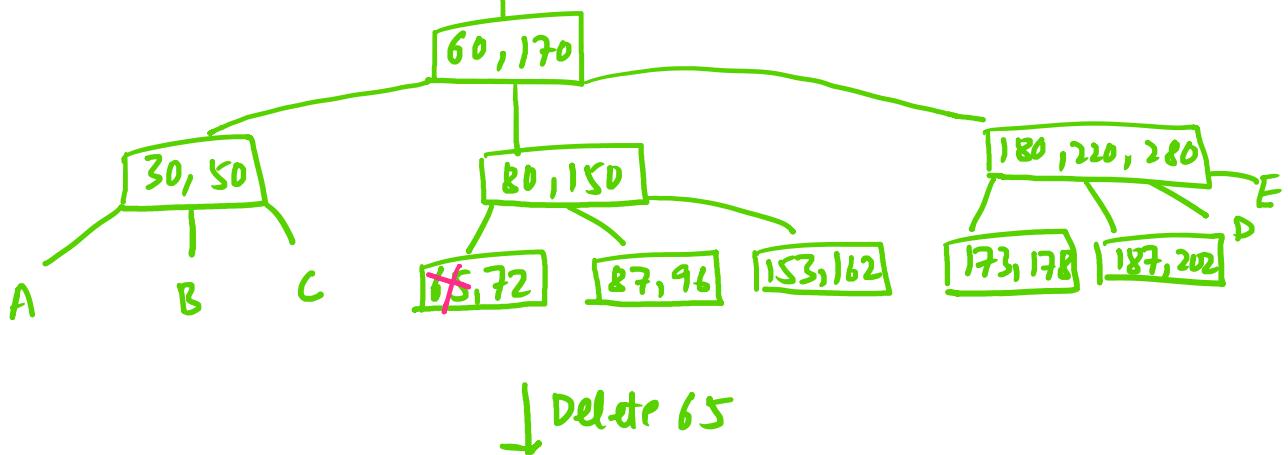
2 options:-

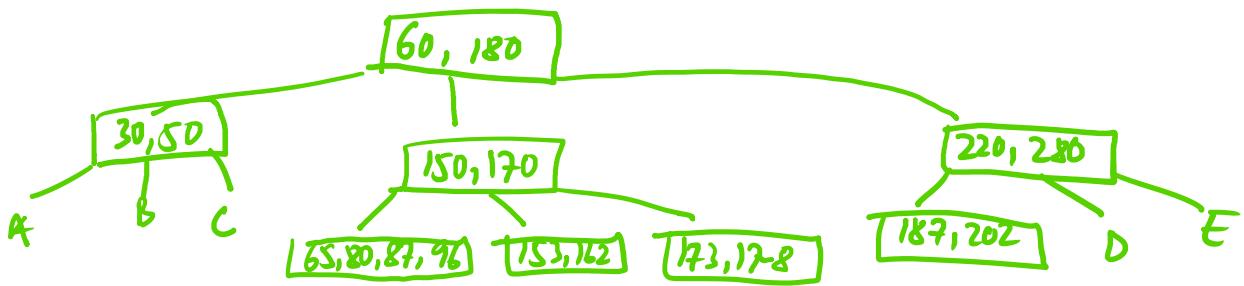
- 1) Combine with the right sibling node by moving the middle key down.
- 2) Combine with the left sibling node by moving the middle key down



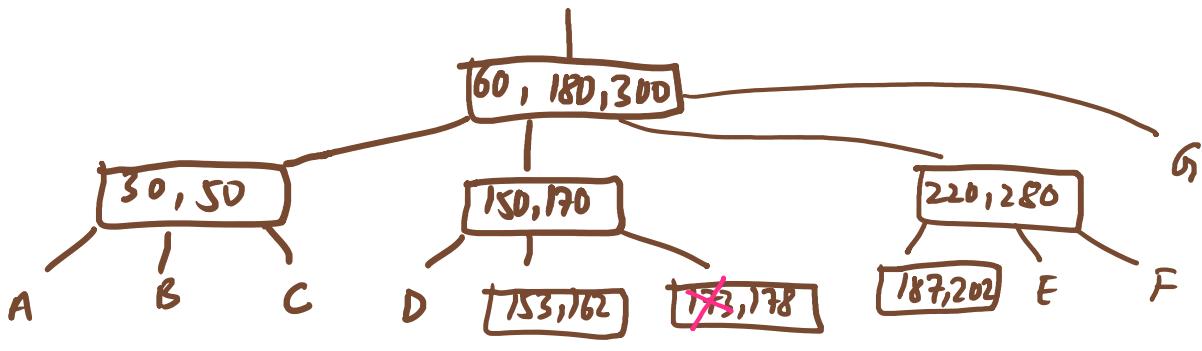
3) Consolidation 4 Borrowing

$B \rightarrow C \rightarrow \dots$
 i $\underbrace{B \rightarrow C}_{\text{parent}(i)}$

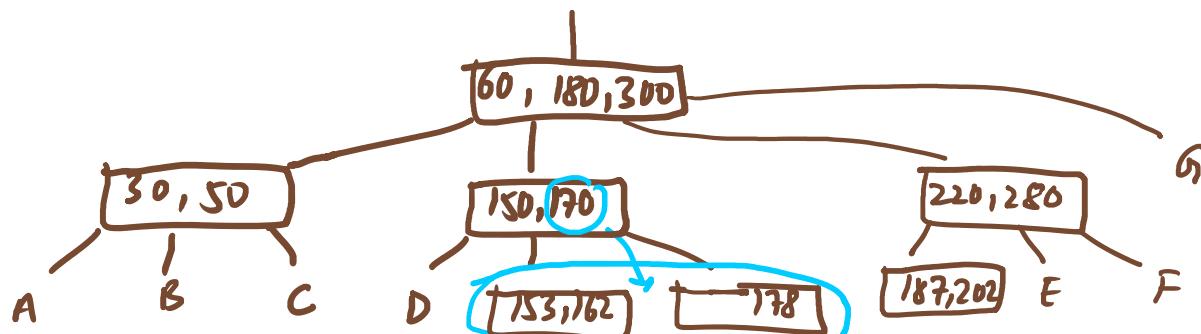




4) Double Consolidation :-



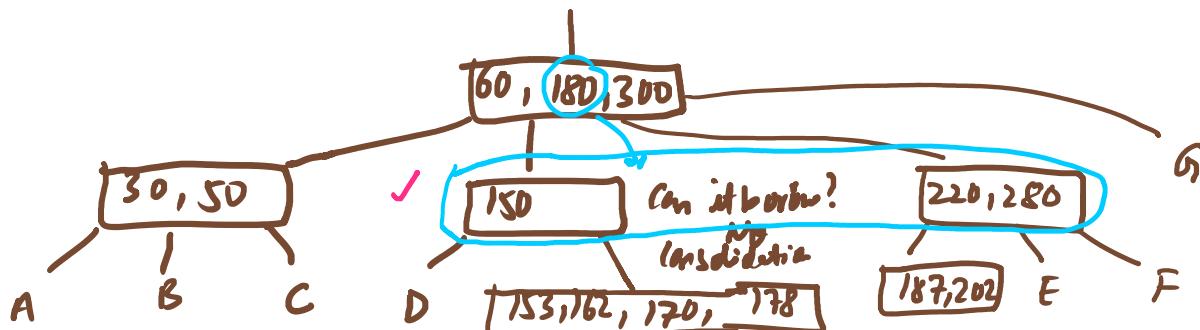
↓ Delete 171



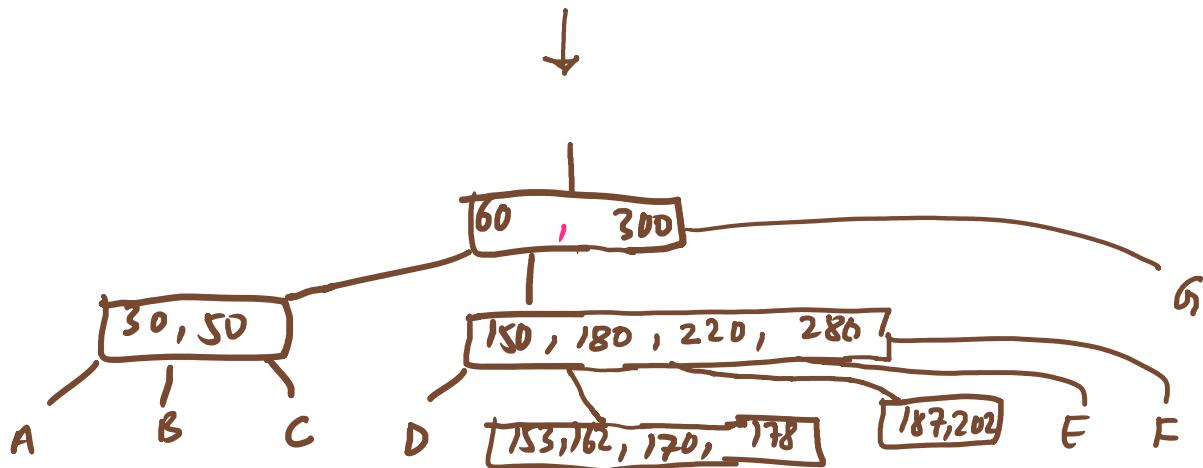
Can it borrow?



No
Consolidation



A ~ ~ D [155, 162, 170, 178]

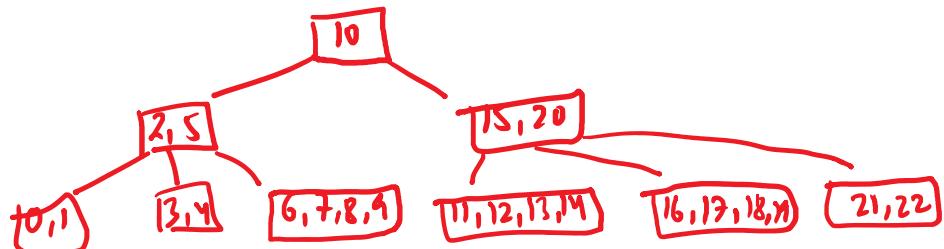


B + Tree :-

Drawback of a B tree :-

- 1) (Key, Value) pair is present in all the nodes (internal & external)

↑ ↑
Index Data



This will increase the no. of levels thereby increasing the height.

- 2) Searching of multiple Keys is time consuming

e.g.: - 1, 3 what if the leaves are connected.

B + Tree :- Two types of nodes

- 1) Internal Nodes - Only keys are stored → index information
- 2) External Node - Actual data or record will be stored.

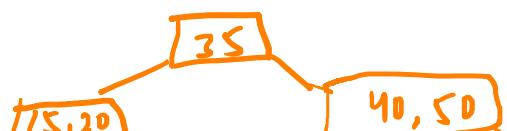
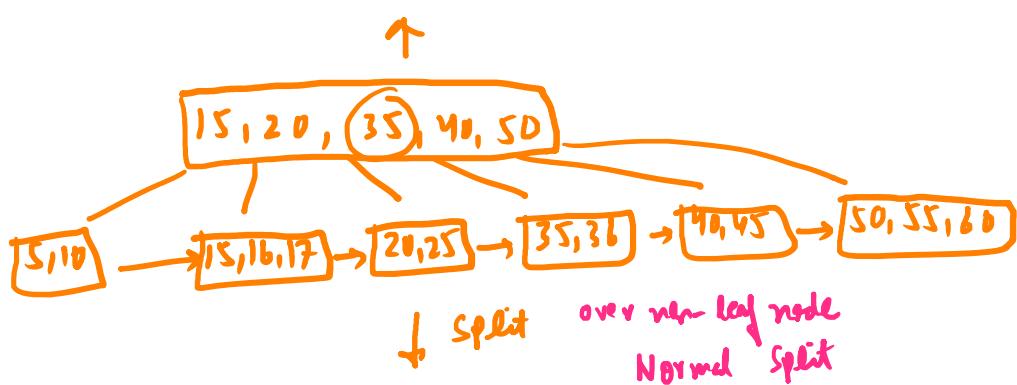
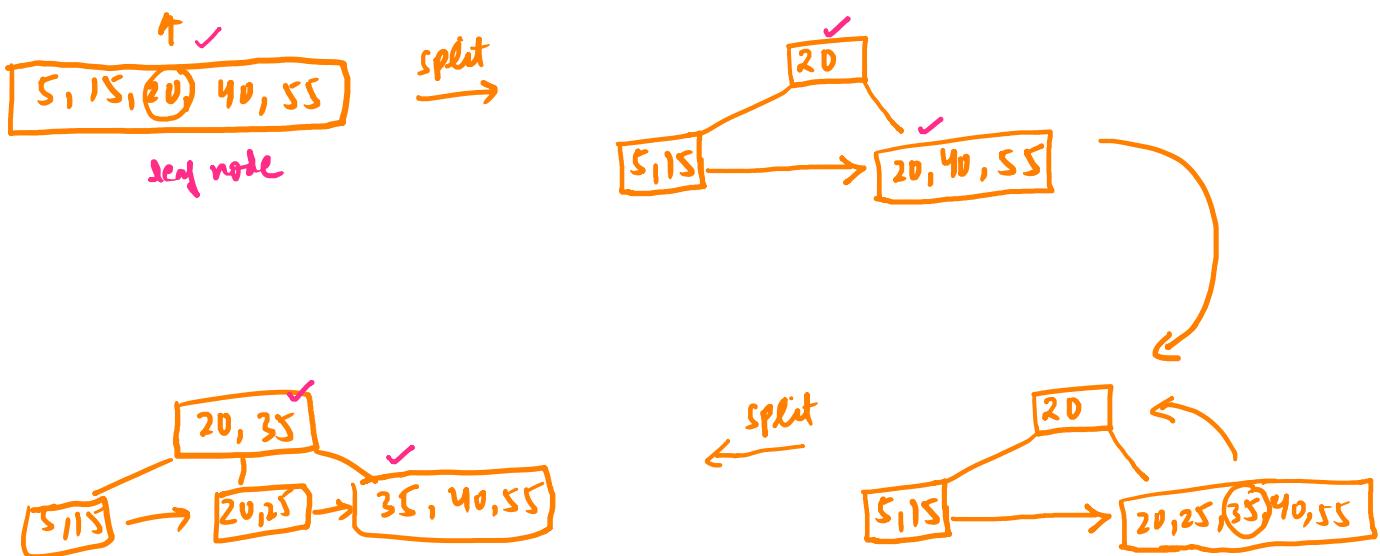
All the external nodes are connected via linked list.

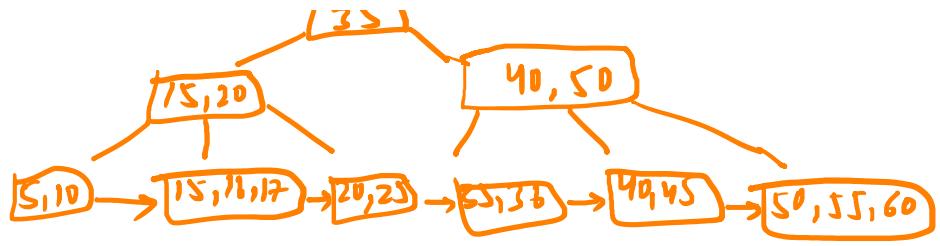
B+Tree Insertion :-

- B+Tree of order 'm' will follow all the properties of a B Tree of order 'm'.
- Whenever a split is performed over a leaf then retain a copy of that middle element in the right child leaf node.

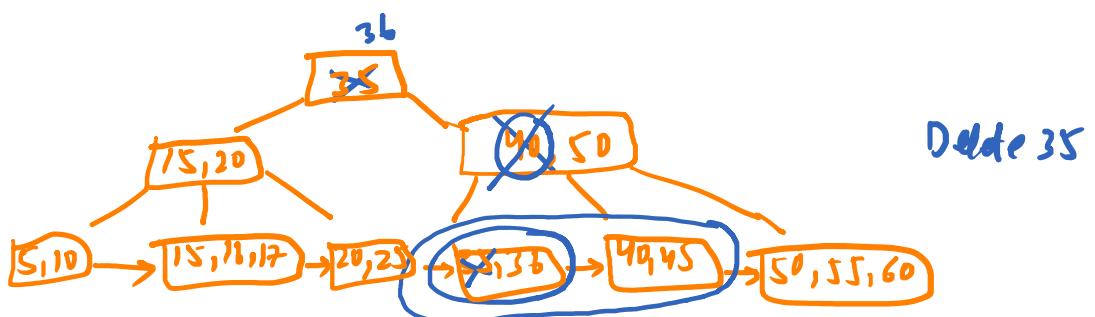
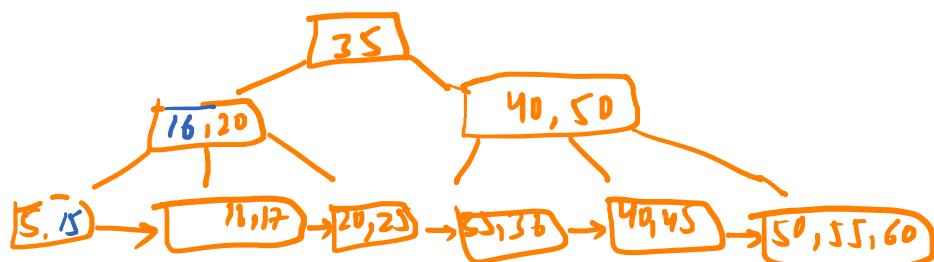
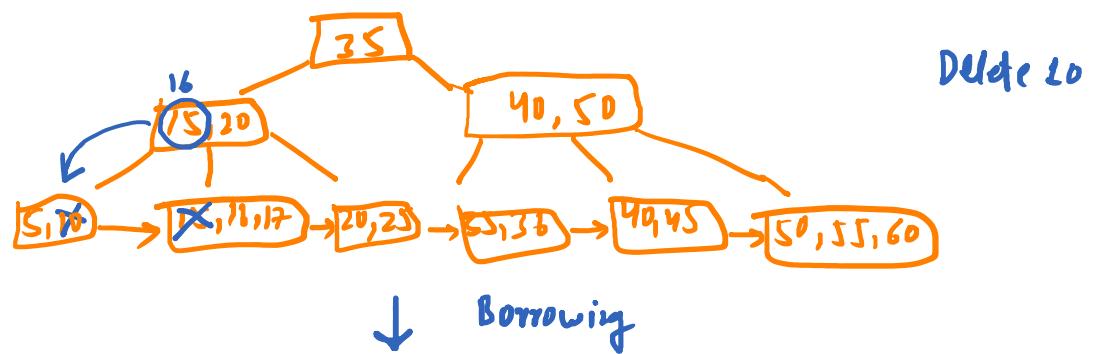
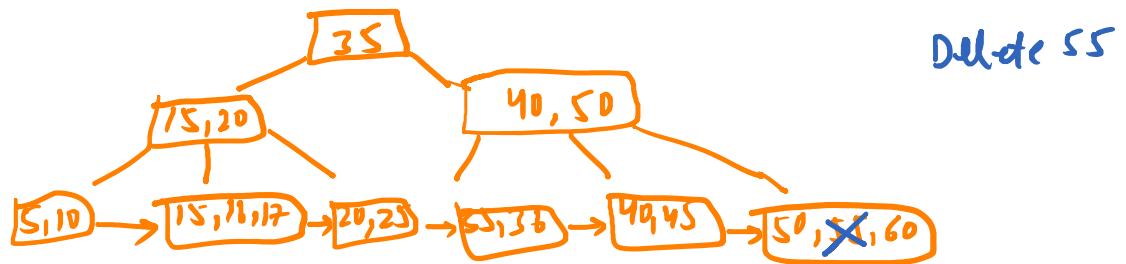
Insert the following keys in a B+Tree of Order 5.

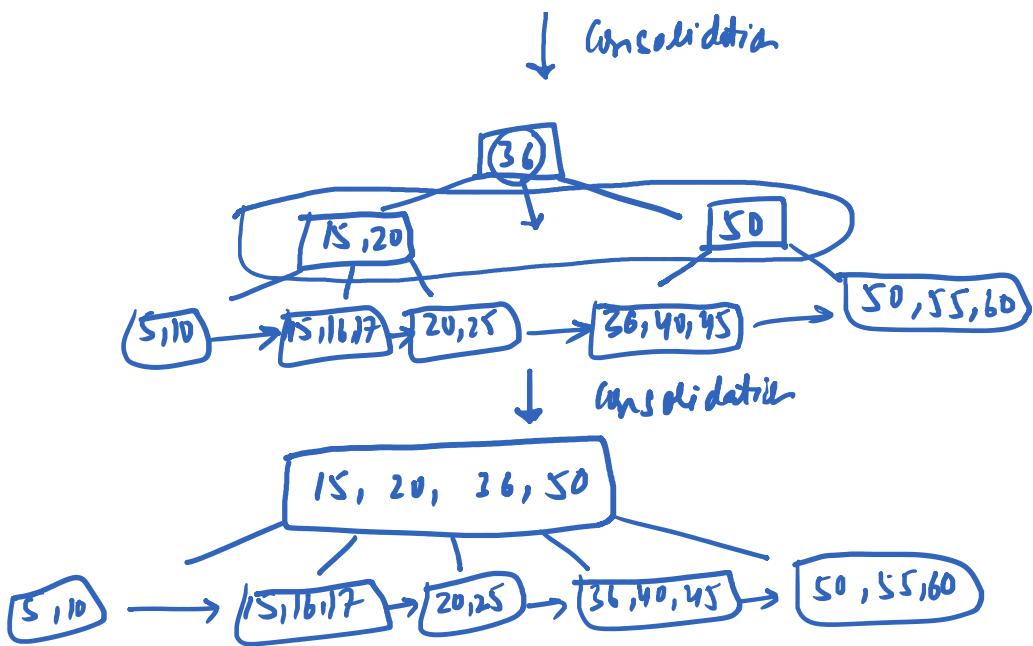
5, 40, 15, 20, 55, 35, 25, 10, 36, 50, 45, 60, 16, 12





B+ Tree Deletion :- Logic is same as B Tree





Hashing :- Time complexity to search a record
 $\Theta(1)$, $O(n)$

mapping the data to its location.

Hash function $h(x) = y$

\uparrow \uparrow
 data location

DS used in a hash table

Collision : if x_1 & x_2 are 2 data values where $x_1 \neq x_2$
 but $h(x_1) = h(x_2)$

Example :- Insert the following elements in a hash table of size 10
 20, 35, 14, 26, 9, 53

0	20
1	
2	
3	53

$$\text{Hash function} = x \% n$$

$$x \mod n$$

\uparrow \uparrow
 x n

2	
3	53
4	14
5	35
6	26
7	
8	
9	9

$\pi \bmod n$
 \uparrow
 data
 n
 \uparrow
 size of hash table

$$h(x) = x \% 10$$

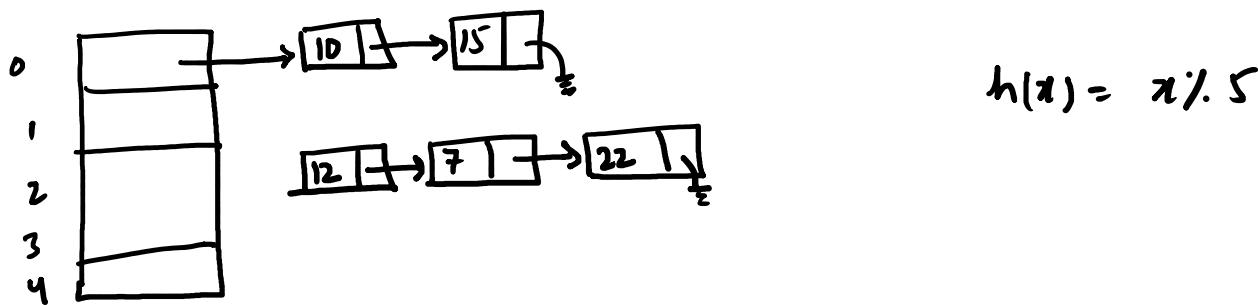
$$h(20) = 20 \% 10 = 0$$

$$\approx 35 \% 10 = 5$$

$$h(26) = 26 \% 10 = 6$$

Collision Resolution Techniques :-

- 1) Direct Chaining :- Array of linked lists
 Separate chaining



- 2) Open Addressing :- Array based implementation

Probing [Linear Probing
 Quadratic Probing]

Linear Probing :- $h(x) = (x+i) \% n$ rehashing function
 where $i = 1, 2, 3, \dots$

0	20
1	30
2	40
3	53
4	14
5	35
6	26

Insert 30

$$h(30) = 30 \% 10 = 0$$

$$re-h(30+1) = 31 \% 10 = 1$$

$\dots + 40$

collision

varat

Primary Clustering

5	35
6	26
7	
8	
9	9

Insert 40

- h(40) = 0 collision
 re- h(40+1) = 1 collision
 re- h(40+2) = 2 vacant

Primary Clustering

Quadratic Probing $re-h(x) = (x+i^2) \% n$

where $i = 1, 2, 3, \dots$

0	20
1	30
2	
3	
4	14
5	35
6	40
7	
8	
9	9

Insert 40

- h(40) = 0 collision
 h(40+1^2) = 1 collision
 h(40+2^2) = 4 collision
 h(40+3^2) = 9 collision
 h(40+4^2) = 6 vacant

0	20
1	30
2	40
3	50
4	
5	
6	
7	
8	
9	

Huffman Tree :- Used for data compression

Symbol	Code (Fixed Length)
A	000
B	001
C	010
D	011
E	100

0 0 0 0 0 1 0 1 0
A B C

3 bits per symbol

Can this be reduced?

length of message = n symbols
 size of message = $3n$

Lossless

compression

Huffman tree uses variable length codes.

Frequently occurring symbol should be given shorter codes.

Symbol	Frequency
A	2
B	1

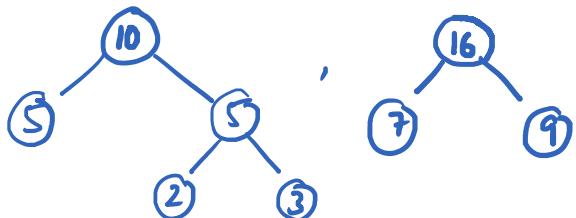
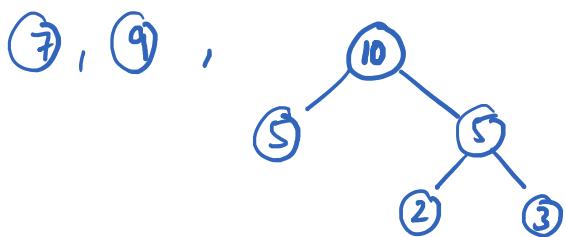
Huffman Tree Algorithm :-

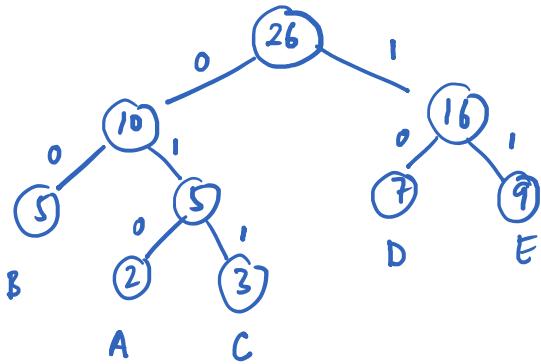
- 1) All the symbols with their frequencies shall

A	2
B	5
C	3
D	7
E	9
	<u>26</u>

- 1) All the symbols with their frequencies shall become individual nodes (or single node tree)
- 2) Select 2 trees with minimum root value.
- 3) Combine these 2 trees by making one as Left subtree & other as right subtree, and their overall root node shall have the value equivalent to the sum of individual roots.
- 4) Repeat step 2 to 3 until one tree remains & shall be the Huffman Tree.

(2), (3), (5), (7), (9)





Huffman Tree

$$\begin{aligned}\text{Average bits / symbol} &= \frac{3+2+3+2+2}{5} \\ &= 2.4 \text{ bits / symbol}\end{aligned}$$

Symbol	Frequency	Huffman Code
A	2	010
B	5	00
C	3	011
D	7	10
E	9	11

$$\begin{aligned}\text{Average bits / symbol} &= \frac{\sum_{i=1}^n f_i x_i}{\sum_{i=1}^n f_i} = \frac{2 \times 3 + 5 \times 2 + 3 \times 3 + 7 \times 2 + 9 \times 2}{26} \\ &= \frac{6 + 10 + 9 + 14 + 18}{26} \\ &= 2.19 \text{ bits / symbol}\end{aligned}$$

$\downarrow -\log_2 p_i$

$$\begin{aligned}\text{Shannon's Entropy} &= - \sum_{i=1}^n p_i \log_2 p_i \\ &= - \left[\frac{2}{26} \log_2 \frac{2}{26} + \frac{5}{26} \log_2 \frac{5}{26} + \frac{3}{26} \log_2 \frac{3}{26} + \frac{7}{26} \log_2 \frac{7}{26} + \frac{9}{26} \log_2 \frac{9}{26} \right] \\ &= 2.5 \text{ bits / symbol}\end{aligned}$$

Average information content

A A A A A

$$E = - \sum_{i=1}^n p_i \log_2 p_i$$

$$E = - p_A \log_2 p_A \quad p_A = 1$$

min value of entropy

$$\boxed{E = 0} \Rightarrow \text{no randomness}$$

A B C D E

$$E = - \sum_{i=1}^n p_i \log_2 p_i$$

$$1/5, 1/5, 1/5, 1/5, 1/5$$

$$p_i = \frac{1}{5}$$

Entropy

$$H = - \sum_{i=1}^n p_i \log_2 p_i$$

$$= -5 \left[\frac{1}{5} \log_2 \frac{1}{5} \right] \quad p_i = \frac{1}{5}$$

$$= - \log_2 \frac{1}{5} = \boxed{\log_2 5}$$

$$\text{max. value of entropy} = \log_2 n$$

Highly random