

# MA/CSSE 473

## Day 29

Optimal BST  
Conclusion

Then student  
questions about the  
exam



Dynamic Programming Example

**OPTIMAL BINARY SEARCH TREES**



## Recap: Optimal Binary Search Trees

- Suppose we have  $n$  distinct data keys  $K_1, K_2, \dots, K_n$  (in increasing order) that we wish to arrange into a Binary Search Tree
- This time the expected number of probes for a successful or unsuccessful search depends on the shape of the tree and where the search ends up
- Guiding principle for optimization?
- This discussion follows Reingold and Hansen, *Data Structures*. **An excerpt on optimal static BSTS is posted on Moodle.** I use  $a_i$  and  $b_i$  where Reingold and Hansen use  $\alpha_i$  and  $\beta_i$ .



## What contributes to the expected number of probes?

- Frequencies, depth of node
- For successful search, number of probes is one more than the depth of the corresponding internal node
- For unsuccessful, number of probes is equal to the depth of the corresponding external node

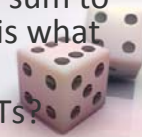


## Optimal BST Notation

- Keys are  $K_1, K_2, \dots, K_n$ , in internal nodes  $x_1, x_2, \dots, x_n$
- Let  $v$  be the value we are searching for
- For  $i = 1, \dots, n$ , let  $a_i$  be the probability that  $v$  is key  $K_i$
- For  $i = 1, \dots, n-1$ , let  $b_i$  be the probability that  $K_i < v < K_{i+1}$ 
  - Similarly, let  $b_0$  be the probability that  $v < K_1$ , and  $b_n$  the probability that  $v > K_n$
  - Each  $b_i$  is associated with external node  $y_i$

- Note that 
$$\sum_{i=1}^n a_i + \sum_{i=0}^n b_i = 1$$

- We can also just use *frequencies* instead of *probabilities* when finding the optimal tree (and divide by their sum to get the probabilities if we ever need them). That is what we will do in an example.
- Should we try exhaustive search of all possible BSTs?



## What not to measure

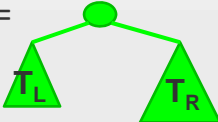
- What about external path length and internal path length?
- These are too simple, because they do not take into account the frequencies.
- We need *weighted* path lengths.



## Weighted Path Length

$$C(T) = \sum_{i=1}^n a_i [1 + \text{depth}(x_i)] + \sum_{i=0}^n b_i [\text{depth}(y_i)]$$

**Note:**  $y_0, \dots, y_n$  are the external nodes of the tree

- If we divide this by  $\sum a_i + \sum b_i$ , we get the expected number of probes.
- We can also define C recursively:
- $C(\square) = 0$ . If  $T =$   , then

$C(T) = C(T_L) + C(T_R) + \sum a_i + \sum b_i$ , where the summations are over all  $a_i$  and  $b_i$  for nodes in T

- It can be shown by induction that these two definitions are equivalent (a homework problem).



## Example

- Frequencies of vowel occurrence in English
- : A, E, I, O, U
- a's: 32, 42, 26, 32, 12
- b's: 0, 34, 38, 58, 95, 21
- Draw a tree (with E as root), and see which is best. (sum of a's and b's is 390).



## Strategy

- We want to minimize the weighted path length
- Once we have chosen the root, the left and right subtrees must themselves be optimal EBSTs
- We can build the tree from the bottom up, keeping track of previously-computed values



## Intermediate Quantities

- Cost: Let  $C_{ij}$  (for  $0 \leq i \leq j \leq n$ ) be the cost of an optimal tree (not necessarily unique) over the frequencies  $b_i, a_{i+1}, b_{i+1}, \dots, a_j, b_j$ . Then
- $C_{ii} = 0$ , and 
$$C_{ij} = \min_{i < k \leq j} (C_{i,k-1} + C_{kj}) + \sum_{t=i}^j b_t + \sum_{t=i+1}^j a_t$$
- This is true since the subtrees of an optimal tree must be optimal
- To simplify the computation, we define
- $W_{ii} = b_i$ , and  $W_{ij} = W_{i,j-1} + a_j + b_j$  for  $i < j$ .
- Note that  $W_{ij} = b_i + a_{i+1} + \dots + a_j + b_j$ , and so
- $C_{ii} = 0$ , and 
$$C_{ij} = W_{ij} + \min_{i < k \leq j} (C_{i,k-1} + C_{kj})$$
- Let  $R_{ij}$  (root of best tree from  $i$  to  $j$ ) be a value of  $k$  that minimizes  $C_{i,k-1} + C_{kj}$  in the above formula



## Code

```
# initialize the main diagonal
for i in range(n + 1):
    R[i][i] = i
    W[i][i] = b[i]
    # Draw this cell of the table in the given window.
    drawSquare(i, i, W[i][i], C[i][i], R[i][i], win, indent, squareSize)
# Now populate each of the n upper diagonals:
for d in range(1, n+1): # fill in this diagonal
    # The previous diagonals are already filled in.
    for i in range(n - d + 1):
        j = i + d; # on the dth diagonal, j - i = d
        opt = i + 1 # until we find a better one
        for k in range(i+2, j+1):
            if C[i][k-1]+C[k][j] < C[i][opt-1]+C[opt][j]:
                opt = k
        R[i][j] = opt
        W[i][j] = W[i][j-1] + a[j] + b[j]
        C[i][j] = C[i][opt-1] + C[opt][j] + W[i][j]
        # Draw this cell of the table in the given window.
        drawSquare(i, j, W[i][j], C[i][j], R[i][j], win, indent, squareSize)
```



## Results

R00: 0	R01: 1	R02: 2	R03: 2	R04: 3	R05: 4
W00: 0	W01: 66	W02: 146	W03: 230	W04: 357	W05: 390
C00: 0	C01: 66	C02: 212	C03: 418	C04: 754	C05: 936
	R11: 1	R12: 2	R13: 3	R14: 3	R15: 4
	W11: 34	W12: 114	W13: 198	W14: 325	W15: 358
	C11: 0	C12: 114	C13: 312	C14: 624	C15: 798
		R22: 2	R23: 3	R24: 4	R25: 4
		W22: 38	W23: 122	W24: 249	W25: 282
		C22: 0	C23: 122	C24: 371	C25: 532
			R33: 3	R34: 4	R35: 4
			W33: 58	W34: 185	W35: 218
			C33: 0	C34: 185	C35: 346
				R44: 4	R45: 5
				W44: 95	W45: 128
				C44: 0	C45: 128
					R55: 5
					W55: 21
					C55: 0

**How to  
construct the  
optimal tree?**

**Analysis of the  
algorithm?**

- Constructed by diagonals, from main diagonal upward
- What is the optimal tree?



## Running time

- Most frequent statement is the comparison  
if  $C[i][k-1]+C[k][j] < C[i][opt-1]+C[opt][j]$ :
- How many times  
does it execute:

$$\sum_{d=1}^n \sum_{i=0}^{n-d} \sum_{k=i+2}^{i+d} 1$$

```
simplify(sum(sum(sum(1,k=i+2..i+d),i=0..n-d),d=1..n));
```

$$-\frac{1}{6}n + \frac{1}{6}n^3$$

