

Victor Leonardo Mascarenhas Soares Horta

ÁRVORE BINÁRIA DE BUSCA ÓTIMA (OBST)

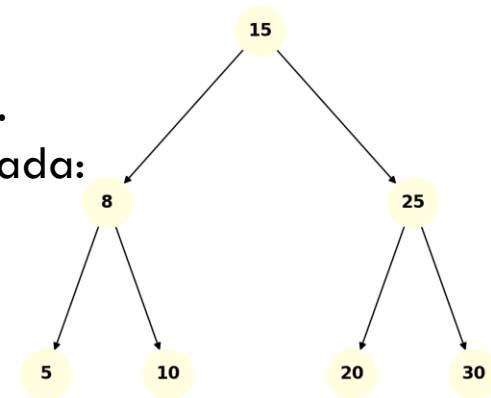


INTRODUÇÃO

O que é uma árvore binária de busca(BST)?

Uma **árvore binária de busca** é um grafo $G = (V, E)$ tal que:

1. G é **conexo**.
2. G é **acíclico**.
3. Existe uma raiz r que define uma estrutura **hierárquica**.
4. Cada nó tem no máximo 2 filhos (grau ≤ 3 se considerar o pai).
5. Existe uma **ordem total** $<$ sobre as chaves em V , e ela é respeitada:
 1. Subárvore esquerda = elementos menores.
 2. Subárvore direita = elementos maiores.



Do ponto de vista de **uso**, ela permite buscas rápidas: cada comparação descarta metade dos nós possíveis.

PROBLEMA

O **Problema da Árvore Binária de Busca Ótima (OBST)** consiste em determinar a estrutura de uma árvore binária de busca que **minimiza o custo esperado das operações de busca**, dado que não apenas conhecemos o conjunto de chaves, mas também as **probabilidades associadas a cada chave**.

Formalmente, seja $K = \{k_1, k_2, \dots, k_n\}$ um conjunto ordenado de chaves e p_i a probabilidade de busca da chave k_i . O objetivo é construir uma BST de modo a reduzir:

$$E[\text{custo}] = \sum_{i=1}^n p_i \cdot \text{profundidade}(k_i)$$

onde $\text{profundidade}(k_i)$ representa a posição do nó k_i na árvore (raiz com profundidade 1).

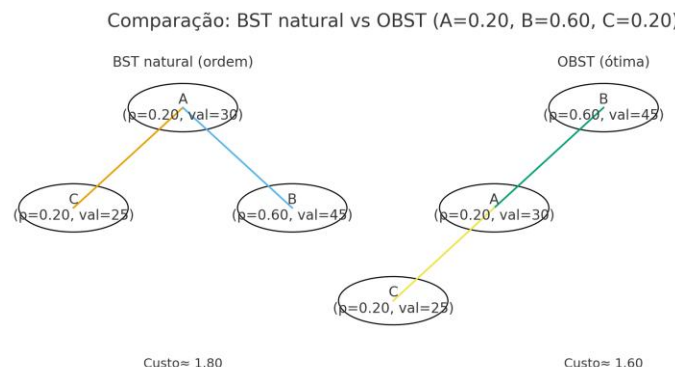
Exemplo intuitivo

Se temos 3 chaves:

A: 20% das buscas e valor 30

B: 60% das buscas e valor 45

C: 20% das buscas e valor 25



Em uma árvore **sem considerar probabilidades**, A ficaria na raiz.

Mas como B é a mais provável (60%), a **OBST coloca B como raiz**, para reduzir o custo médio.

APLICAÇÕES

Esse problema não é só teórico. Ele aparece em várias aplicações:

Em **compiladores**, para organizar tabelas de símbolos.

Em **bancos de dados**, quando muitas consultas são feitas em alguns registros específicos.

Em **compressão e recuperação de informação**, onde acessamos dados com frequências diferentes.

E, no dia a dia?

Organizar prateleiras de supermercado

Menus de aplicativos

Agendas

O que é mais usado fica em destaque.

COMO O ALGORITMO RESOLVE

A solução do problema usa **programação dinâmica**. Não vou entrar nos detalhes, mas a ideia geral é esta:

1. Divide em subproblemas

Analisa intervalos menores de chaves (subárvores).

Ex.: $\{A, B, C, D\} \rightarrow$ primeiro A, B, C, D sozinhos \rightarrow depois pares $(A, B) \rightarrow$ depois trios (A, B, C) , etc.

2. Usa tabelas auxiliares (Programação Dinâmica)

$e[i, j]$:custo mínimo do intervalo.

$w[i, j]$:soma das probabilidades.

$root[i, j]$:raiz ótima daquele intervalo.

Evita recalcular problemas já resolvidos.

3. Testa raízes possíveis

$$e[i, j] = \min_{r=i}^j (e[i, r-1] + e[r+1, j] + w[i, j])$$

4. Reconstrói a árvore ótima

Após preencher as tabelas, identifica a raiz global.

Usa $root[i, j]$ para reconstruir recursivamente a árvore final.

COMO O ALGORITMO RESOLVE

“O custo de um intervalo de chaves é o menor valor que consigo ao testar cada possível raiz, somando o custo da esquerda, o custo da direita e o peso do intervalo.”

EXEMPLO

Suponha que temos um sistema de busca de nomes em uma lista ordenada de quatro pessoas:

Ana, Bruno, Carla e Diego.

As probabilidades de cada nome ser pesquisado são:

Ana: **35%**

Bruno: **25%**

Carla: **25%**

Diego: **15%**

Além disso, existe uma pequena chance de buscas por nomes que não estão na lista, distribuída igualmente entre os espaços entre as chaves.

Como organizar esses nomes em uma Árvore Binária de Busca de forma que o número médio de comparações para encontrar um nome seja o menor possível?

Em outras palavras, queremos construir a **Árvore Binária de Busca Ótima (OBST)** para esse conjunto de chaves e probabilidades.

CÓDIGO EM PYTHON

```
from typing import List, Any
import sys

# -----
# Validação de p e q
# -----
def validate_probabilities(p: List[float], q: List[float]) -> None:
    if len(q) != len(p) + 1:
        raise ValueError("q deve ter tamanho len(p)+1 (há n+1 gaps para n chaves).")
    if any(x < 0 for x in p + q):
        raise ValueError("Probabilidades negativas não são permitidas.")
    s = sum(p) + sum(q)
    if abs(s - 1.0) > 1e-6:
        raise ValueError(f"As probabilidades devem somar 1.0 (atual: {s:.6f}).")
```


CÓDIGO EM PYTHON

```
34 # Núcleo: OBST (DP clássica)
35 # -----
36
37 #Implementa a OBST por programação dinâmica e devolve três tabelas:
38 def optimal_bst(p: List[float], q: List[float], keys: List[Any]):
39     n = len(p)
40     if len(keys) != n:
41         raise ValueError("keys e p devem ter o mesmo tamanho.")
42     #menor custo esperado para as chaves do intervalo i..j;
43     e = [[0.0]*(n+2) for _ in range(n+2)]
44     #peso (soma das probabilidades p e q) no intervalo i..j;
45     w = [[0.0]*(n+2) for _ in range(n+2)]
46     #índice da raiz ótima escolhida para i..j.
47     root = [[0]*(n+2) for _ in range(n+2)]
48
49     #Inicializa os casos-base (intervalo vazio i..i-1):
50     #custo = q[i-1], peso = q[i-1].
51     #Isto modela a chance de "não encontrar" entre as chaves.
52     for i in range(1, n+2):
53         e[i][i-1] = q[i-1]
54         w[i][i-1] = q[i-1]
55     for l in range(1, n+1):
56         for i in range(1, n-l+2):
57             j = i + l - 1
58             #Inicializa o custo do intervalo com infinito
59             e[i][j] = float("inf")
60             #Atualiza o peso do intervalo i..j a partir do peso i..j-1:
61             #adiciona p[j-1] (prob. da chave j) e q[j] (gap após j).
62
63             w[i][j] = w[i][j-1] + p[j-1] + q[j]
64             melhor_r, melhor_custo = None, float("inf")
65             #Testa cada chave r em i..j como candidata a raiz do intervalo.
66             for r in range(i, j+1):
```

CÓDIGO EM PYTHON

```
66     for r in range(i, j+1):
67
68         #Calcula o custo se r for raiz:
69         #custo da subárvore esquerda e[i][r-1] +
70         # custo da subárvore direita e[r+1][j] +
71         #w[i][j] (todo mundo do intervalo fica um nível mais fundo quando escolhemos uma raiz).
72
73         custo = e[i][r-1] + e[r+1][j] + w[i][j]
74         if custo < melhor_custo:
75             melhor_custo, melhor_r = custo, r
76     e[i][j] = melhor_custo
77     root[i][j] = melhor_r
78     return e, w, root
```

CÓDIGO EM PYTHON

```
80 # -----
81 # Ordens (percursos) a partir da tabela root
82 #Essas três funções percorrem recursivamente os intervalos indicados pela tabela root para gerar, sem construir a árvore, as listas de chaves em pré-ordem, em-ordem e
83 # pós-ordem.
84 # -----
85 def preorder_from_root(root, keys, i, j):
86     if j < i:
87         return []
88     r = root[i][j]
89     return [keys[r-1]] + preorder_from_root(root, keys, i, r-1) + preorder_from_root(root, keys, r+1, j)
90
91 def inorder_from_root(root, keys, i, j):
92     if j < i:
93         return []
94     r = root[i][j]
95     return inorder_from_root(root, keys, i, r-1) + [keys[r-1]] + inorder_from_root(root, keys, r+1, j)
96
97 def postorder_from_root(root, keys, i, j):
98     if j < i:
99         return []
100     r = root[i][j]
101     return postorder_from_root(root, keys, i, r-1) + postorder_from_root(root, keys, r+1, j) + [keys[r-1]]
```

CÓDIGO EM PYTHON

```
103 # -----
104 # Programa principal
105 # -----
106 def main():
107     # Exemplo visível de entrada:
108     keys = ["Ana", "Bruno", "Carla", "Diego"]
109     p = [0.35, 0.25, 0.25, 0.15]
110     q = [0.02, 0.02, 0.02, 0.02, 0.02]
111
112     print(">>> EXEMPLO UTILIZADO:")
113     print(f"    keys = {keys}")
114     print(f"    p    = {p}")
115     print(f"    q    = {q}\n")
116
117     # Normaliza p+q para 1.0 (mantém proporções) – evita erro de soma
118     total = sum(p) + sum(q)
119     if abs(total - 1.0) > 1e-9:
120         esc = 1.0 / total
121         p = [round(x*esc, 10) for x in p]
122         q = [round(x*esc, 10) for x in q]
```

CÓDIGO EM PYTHON

```
124 # Validação final
125 validate_probabilities(p, q)
126
127 print(">>> Executando OBST (programação dinâmica...)")
128 e, w, root = optimal_bst(p, q, keys)
129 n = len(keys)
130
131 # Custo
132 cost_dp = e[1][n]
133 print(f">>> Custo mínimo esperado (DP): {cost_dp:.6f}\n")
134
135 # Raízes por intervalo (decisões do algoritmo)
136 print(">>> Raízes escolhidas por intervalo:")
137 for L in range(1, n+1):
138     for i in range(1, n-L+2):
139         j = i + L - 1
140         r = root[i][j]
141         print(f"    Intervalo [{i},{j}] → chaves {keys[i-1:j]} | raiz = '{keys[r-1]}'")
142
143 # Ordem escolhida (percursos)
144 pre = preorder_from_root(root, keys, 1, n)
145 ino = inorder_from_root(root, keys, 1, n)
146 post = postorder_from_root(root, keys, 1, n)
147
148 print("\n>>> Ordem escolhida (percursos da OBST):")
149 print("    Pré-ordem (raiz → esquerda → direita):", pre)
150 print("    Em-ordem (esquerda → raiz → direita):", ino, "    (deve sair ordenada)")
151 print("    Pós-ordem (esquerda → direita → raiz):", post)
152
153 print("\n>>> FIM 📌 A OBST minimiza o número médio de comparações considerando as probabilidades.")
154
```

SAÍDA

```
PS C:\Users\masca> & C:/Users/masca/AppData/Local/Microsoft/WindowsApps/python3.11.exe "c:/Users/masca/OneDrive/Documents/Mestrado 2025,1/obst_completo.py"

>>> Executando OBST (programação dinâmica)...
>>> Custo mínimo esperado (DP): 2.036364

>>> Raízes escolhidas por intervalo:
Intervalo [1,1] → chaves ['Ana'] | raiz = 'Ana'
Intervalo [2,2] → chaves ['Bruno'] | raiz = 'Bruno'
Intervalo [3,3] → chaves ['Carla'] | raiz = 'Carla'
Intervalo [4,4] → chaves ['Diego'] | raiz = 'Diego'
Intervalo [1,2] → chaves ['Ana', 'Bruno'] | raiz = 'Ana'
Intervalo [2,3] → chaves ['Bruno', 'Carla'] | raiz = 'Bruno'
Intervalo [3,4] → chaves ['Carla', 'Diego'] | raiz = 'Carla'
Intervalo [1,3] → chaves ['Ana', 'Bruno', 'Carla'] | raiz = 'Bruno'
Intervalo [2,4] → chaves ['Bruno', 'Carla', 'Diego'] | raiz = 'Carla'
Intervalo [1,4] → chaves ['Ana', 'Bruno', 'Carla', 'Diego'] | raiz = 'Bruno'

>>> Ordem escolhida (percursos da OBST):
Pré-ordem (raiz → esquerda → direita): ['Bruno', 'Ana', 'Carla', 'Diego']
Em-ordem (esquerda → raiz → direita): ['Ana', 'Bruno', 'Carla', 'Diego'] (deve sair ordenada)
Pós-ordem (esquerda → direita → raiz): ['Ana', 'Diego', 'Carla', 'Bruno']

>>> FIM - A OBST minimiza o número médio de comparações considerando as probabilidades.
```

CONCLUSÃO

- **O OBST** é forma de organizar uma BST para **minimizar o custo médio** de busca, usando **probabilidades** de acerto (p) e de falha (q).
- **Como decide:** avalia subproblemas com tabelas $e[i,j]$ (custo), $w[i,j]$ (peso) e $root[i,j]$ (raiz ótima), escolhendo a estrutura de **menor profundidade esperada**.
- **Por que importa:** quando o padrão de buscas é **conhecido/estável**, reduz comparações em média → **desempenho previsível**.
- **No nosso exemplo (Ana 0.35, Bruno 0.25, Carla 0.25, Diego 0.15):** a OBST escolhe **Bruno como raiz**, pois **equilibra** os acessos e diminui o custo total (mesmo Ana sendo a mais provável isoladamente).
- **Quando usar:** tabelas de símbolos, índices/consulta em BD, dicionários, autocompletar; quando **não usar:** cenários muito dinâmicos (prefira árvores balanceadas/splay).
- **Essência: probabilidade guiando estrutura** → menos comparações em média, respeitando a ordem da BST.

OBRIGADO!

O repositório no GitHub contém:

- Slides em PDF
- Código em Python
- Dados de exemplo
- README.md com link para vídeo.