

BÀI TOÁN TREE EDIT DISTANCE

Toán Tổ Hợp và Lý Thuyết Đồ Thị

1 Lý thuyết cơ bản về Tree Edit Distance

1.1 Định nghĩa

Tree Edit Distance là khoảng cách tối thiểu để biến đổi một cây thành cây khác thông qua một chuỗi các phép biến đổi cơ bản.

Cho hai cây có gốc T_1 và T_2 , Tree Edit Distance $\delta(T_1, T_2)$ là số phép biến đổi tối thiểu cần thiết để biến T_1 thành T_2 .

1.2 Các phép biến đổi cơ bản

Tree Edit Distance sử dụng ba phép biến đổi cơ bản:

1. **Insertion (Chèn)**: Chèn một nút mới vào cây
2. **Deletion (Xóa)**: Xóa một nút khỏi cây
3. **Substitution (Thay thế)**: Thay đổi nhãn của một nút

Mỗi phép biến đổi có một chi phí (cost) được định nghĩa:

- $\gamma(a \rightarrow b)$: Chi phí thay thế nhãn a bằng nhãn b
- $\gamma(a \rightarrow \lambda)$: Chi phí xóa nút có nhãn a
- $\gamma(\lambda \rightarrow b)$: Chi phí chèn nút có nhãn b

Trong đó λ biểu diễn nút rỗng.

1.3 Tính chất quan trọng

- Tree Edit Distance thỏa mãn tính chất của một metric
- $\delta(T_1, T_2) \geq 0$ với đẳng thức xảy ra khi $T_1 = T_2$
- $\delta(T_1, T_2) = \delta(T_2, T_1)$ (tính đối xứng)
- $\delta(T_1, T_3) \leq \delta(T_1, T_2) + \delta(T_2, T_3)$ (bất đẳng thức tam giác)

1.4 Ứng dụng

- **Bioinformatics:** So sánh cấu trúc phân tử RNA/DNA
- **Compiler Design:** Phân tích sự khác biệt giữa các AST
- **Web Mining:** So sánh cấu trúc HTML/XML
- **Image Processing:** Phân tích cấu trúc hình ảnh
- **Natural Language Processing:** So sánh cấu trúc cú pháp

2 Mô tả bài toán

Đề bài: Cho hai cây có gốc T_1 và T_2 với các nút được gán nhãn. Tìm khoảng cách chỉnh sửa tối thiểu (Tree Edit Distance) giữa hai cây này bằng các phương pháp:

1. **Backtracking:** Duyệt tất cả các khả năng
2. **Branch-and-bound:** Cắt tỉa các nhánh không tối ưu
3. **Divide-and-conquer:** Chia nhỏ bài toán
4. **Dynamic programming - Quy hoạch động:** Tối ưu hóa bằng lưu trữ kết quả

Input:

- Hai cây T_1 và T_2 với các nút được gán nhãn
- Hàm chi phí cho các phép biến đổi

Output:

- Khoảng cách chỉnh sửa tối thiểu
- Chuỗi các phép biến đổi tối ưu (tùy chọn)

3 Ý tưởng và giải pháp

3.1 Phân tích bài toán

Tree Edit Distance là một bài toán tối ưu hóa có thể được giải quyết bằng nhiều cách tiếp cận khác nhau:

- **Độ phức tạp:** Bài toán có độ phức tạp cao do không gian tìm kiếm lớn
- **Cấu trúc con tối ưu:** Bài toán có tính chất cấu trúc con tối ưu
- **Chồng chéo bài toán con:** Các bài toán con có thể lặp lại nhiều lần

Phương pháp	Độ phức tạp	Bộ nhớ	Tối ưu
Backtracking	$O(3^{n+m})$	$O(h)$	Có
Branch-and-bound	$O(3^{n+m}/k)$	$O(h)$	Có
Divide-and-conquer	$O(n \cdot m \cdot (n + m))$	$O(h)$	Có
Dynamic Programming	$O(n^2 \cdot m^2)$	$O(n \cdot m)$	Có

Bảng 1: So sánh các phương pháp giải Tree Edit Distance

3.2 So sánh các phương pháp

Trong đó n, m là số nút của hai cây, h là chiều cao cây, k là hệ số cắt tỉa.

4 Thuật toán chi tiết

4.1 1. Phương pháp Backtracking

Backtracking duyệt tất cả các khả năng biến đổi có thể:

```

1 Algorithm: BacktrackingTED(T1, T2, i, j)
2 Input: T1, T2 - hai cây, i, j - chỉ số nút hiện tại
3 Output: Khoảng cách chỉnh sửa tối thiểu
4
5 1. If T1[i] is empty and T2[j] is empty:
6     Return 0
7 3. If T1[i] is empty:
8     Return cost_insert(T2[j]) + BacktrackingTED(T1, T2, i, j-1)
9 5. If T2[j] is empty:
10    Return cost_delete(T1[i]) + BacktrackingTED(T1, T2, i-1, j)
11 7.
12 8. // Thu 3 phép biến đổi
13 9. substitute = cost_substitute(T1[i], T2[j]) + BacktrackingTED(T1, T2,
    i-1, j-1)
14 10. delete_op = cost_delete(T1[i]) + BacktrackingTED(T1, T2, i-1, j)
15 11. insert_op = cost_insert(T2[j]) + BacktrackingTED(T1, T2, i, j-1)
16 12.
17 13. Return min(substitute, delete_op, insert_op)
18
19 Time Complexity:  $O(3^{(n+m)})$ 
20 Space Complexity:  $O(h)$  vì h là chiều cao cây

```

Listing 1: Thuật toán Backtracking cho Tree Edit Distance

4.2 2. Branch-and-bound

Branch-and-bound cải thiện backtracking bằng cách cắt tỉa các nhánh không tối ưu:

```

1 Algorithm: BranchBoundTED(T1, T2, current_cost, best_cost)
2 Input: T1, T2 - hai cây, current_cost - chi phí hiện tại, best_cost -
    chi phí tốt nhất
3 Output: Khoảng cách chỉnh sửa tối thiểu
4
5 1. // Tính lower bound

```

```

6 2. lower_bound = current_cost + EstimateLowerBound(remaining_T1,
   remaining_T2)
7 3.
8 4. If lower_bound >= best_cost:
9 5.     Return INFINITY // Cat tia nhanh nay
10 6.
11 7. If T1 is empty and T2 is empty:
12 8.     best_cost = min(best_cost, current_cost)
13 9.     Return current_cost
14 10.
15 11. // Thu cac phep bien doi
16 12. For each operation in {substitute, delete, insert}:
17 13.     new_cost = current_cost + cost(operation)
18 14.     result = BranchBoundTED(new_T1, new_T2, new_cost, best_cost)
19 15.     best_cost = min(best_cost, result)
20 16.
21 17. Return best_cost
22
23 Function EstimateLowerBound(T1, T2):
24     // Uoc luong chi phi toi thieu con lai
25     return max(|T1| - |T2|, 0) * min_delete_cost +
26            max(|T2| - |T1|, 0) * min_insert_cost

```

Listing 2: Thuật toán Branch-and-bound

4.3 3. Divide-and-conquer

Chia bài toán thành các bài toán con nhỏ hơn:

```

1 Algorithm: DivideConquerTED(T1, T2)
2 Input: T1, T2 - hai cay
3 Output: Khoảng cách chỉnh sửa tối thiểu
4
5 1. If T1 is empty:
6 2.     Return SumInsertCost(T2)
7 3. If T2 is empty:
8 4.     Return SumDeleteCost(T1)
9 5.
10 6. // Chia cay thanh goc va cac cay con
11 7. root1 = root(T1), subtrees1 = children(T1)
12 8. root2 = root(T2), subtrees2 = children(T2)
13 9.
14 10. min_cost = INFINITY
15 11.
16 12. // Thu thay the goc
17 13. For each possible alignment of subtrees:
18 14.     cost = cost_substitute(root1, root2)
19 15.     For each pair (sub1, sub2) in alignment:
20 16.         cost += DivideConquerTED(sub1, sub2)
21 17.     min_cost = min(min_cost, cost)
22 18.
23 19. // Thu xoa goc cua T1
24 20. cost = cost_delete(root1)
25 21. For each subtree in subtrees1:
26 22.     cost += DivideConquerTED(subtree, T2)
27 23. min_cost = min(min_cost, cost)

```

```

28 24.
29 25. // Thu chen goc cua T2
30 26. cost = cost_insert(root2)
31 27. For each subtree in subtrees2:
32 28.     cost += DivideConquerTED(T1, subtree)
33 29. min_cost = min(min_cost, cost)
34 30.
35 31. Return min_cost

```

Listing 3: Thuật toán Divide-and-conquer

4.4 4. Dynamic Programming

Phương pháp hiệu quả nhất sử dụng quy hoạch động:

```

1 Algorithm: DynamicProgrammingTED(T1, T2)
2 Input: T1, T2 - hai cay
3 Output: Khoảng cách chỉnh sửa tối thiểu
4
5 1. // Khoi tao bang DP
6 2. n = |T1|, m = |T2|
7 3. dp[0..n][0..m] = INFINITY
8 4.
9 5. // Truong hop co so
10 6. dp[0][0] = 0
11 7. For i = 1 to n:
12 8.     dp[i][0] = dp[i-1][0] + cost_delete(T1[i])
13 9. For j = 1 to m:
14 10.    dp[0][j] = dp[0][j-1] + cost_insert(T2[j])
15 11.
16 12. // Dien bang DP
17 13. For i = 1 to n:
18 14.    For j = 1 to m:
19 15.        // Thay the
20 16.        dp[i][j] = min(dp[i][j],
21 17.                        dp[i-1][j-1] + cost_substitute(T1[i], T2[j]))
22 18.
23 19.        // Xoa
24 20.        dp[i][j] = min(dp[i][j],
25 21.                        dp[i-1][j] + cost_delete(T1[i]))
26 22.
27 23.        // Chen
28 24.        dp[i][j] = min(dp[i][j],
29 25.                        dp[i][j-1] + cost_insert(T2[j]))
30 26.
31 27. Return dp[n][m]
32
33 Time Complexity:  $O(n \cdot m)$  v i t h u t t o n Zhang-Shasha
34 Space Complexity:  $O(nm)$ 

```

Listing 4: Thuật toán Dynamic Programming

5 Code C++

```

1 #include <iostream>
2 #include <vector>
3 #include <string>
4 #include <algorithm>
5 #include <climits>
6 #include <queue>
7 #include <unordered_map>
8
9 using namespace std;
10
11 // Cau truc bieu dien nut cay
12 struct TreeNode {
13     string label;
14     vector<TreeNode*> children;
15     int id; // ID duy nhat cho moi nut
16
17     TreeNode(string lbl, int node_id) : label(lbl), id(node_id) {}
18
19     ~TreeNode() {
20         for (auto child : children) {
21             delete child;
22         }
23     }
24 };
25
26 class TreeEditDistance {
27 private:
28     // Ham chi phi
29     int cost_substitute(const string& a, const string& b) {
30         return (a == b) ? 0 : 1;
31     }
32
33     int cost_delete(const string& a) {
34         return 1;
35     }
36
37     int cost_insert(const string& b) {
38         return 1;
39     }
40
41     // Chuyen cay thanh danh sach post-order
42     void getPostOrder(TreeNode* root, vector<TreeNode*>& post_order) {
43         if (!root) return;
44
45         for (auto child : root->children) {
46             getPostOrder(child, post_order);
47         }
48         post_order.push_back(root);
49     }
50
51     // Tinh kích thước cây con tại mỗi nút
52     void computeSubtreeSizes(const vector<TreeNode*>& post_order,
53                             vector<int>& subtree_sizes) {
54         subtree_sizes.resize(post_order.size());
55
56         for (int i = 0; i < post_order.size(); i++) {

```

```

57         subtree_sizes[i] = 1; // Chu the nut hien tai
58
59         for (auto child : post_order[i]->children) {
60             for (int j = 0; j < i; j++) {
61                 if (post_order[j] == child) {
62                     subtree_sizes[i] += subtree_sizes[j];
63                     break;
64                 }
65             }
66         }
67     }
68 }
69
70 public:
71     // 1. Phuong phap Backtracking
72     int backtrackingTED(TreeNode* t1, TreeNode* t2) {
73         if (!t1 && !t2) return 0;
74         if (!t1) return cost_insert(t2->label) +
75             sumInsertCost(t2->children);
76         if (!t2) return cost_delete(t1->label) +
77             sumDeleteCost(t1->children);
78
79         int min_cost = INT_MAX;
80
81         // Thay the nut goc
82         int substitute_cost = cost_substitute(t1->label, t2->label);
83         substitute_cost += computeChildrenCost(t1->children, t2->
children);
84         min_cost = min(min_cost, substitute_cost);
85
86         // Xoa nut goc cua t1
87         int delete_cost = cost_delete(t1->label);
88         for (auto child : t1->children) {
89             delete_cost += backtrackingTED(child, t2);
90         }
91         min_cost = min(min_cost, delete_cost);
92
93         // Chen nut goc cua t2
94         int insert_cost = cost_insert(t2->label);
95         for (auto child : t2->children) {
96             insert_cost += backtrackingTED(t1, child);
97         }
98         min_cost = min(min_cost, insert_cost);
99
100         return min_cost;
101     }
102
103     // 2. Phuong phap Branch-and-bound
104     int branchBoundTED(TreeNode* t1, TreeNode* t2) {
105         int best_cost = INT_MAX;
106         branchBoundHelper(t1, t2, 0, best_cost);
107         return best_cost;
108     }
109
110 private:
111     void branchBoundHelper(TreeNode* t1, TreeNode* t2,

```

```

112         int current_cost, int& best_cost) {
113     // Uoc luong lower bound
114     int lower_bound = current_cost + estimateLowerBound(t1, t2);
115
116     if (lower_bound >= best_cost) {
117         return; // Cat tia
118     }
119
120     if (!t1 && !t2) {
121         best_cost = min(best_cost, current_cost);
122         return;
123     }
124
125     if (!t1) {
126         best_cost = min(best_cost, current_cost + sumInsertCost({t2
127     }));
128     return;
129     }
130
131     if (!t2) {
132         best_cost = min(best_cost, current_cost + sumDeleteCost({t1
133     }));
134     return;
135     }
136
137     // Thu cac phep bien doi
138     // Thay the
139     int new_cost = current_cost + cost_substitute(t1->label, t2->
140     label);
141     branchBoundChildrenHelper(t1->children, t2->children, new_cost,
142     best_cost);
143
144     // Xoa
145     for (auto child : t1->children) {
146         branchBoundHelper(child, t2, current_cost + cost_delete(t1->
147     label), best_cost);
148     }
149
150     // Chen
151     for (auto child : t2->children) {
152         branchBoundHelper(t1, child, current_cost + cost_insert(t2->
153     label), best_cost);
154     }
155     }
156
157     int estimateLowerBound(TreeNode* t1, TreeNode* t2) {
158     int size1 = countNodes(t1);
159     int size2 = countNodes(t2);
160     return abs(size1 - size2);
161     }
162
163 public:
164     // 3. Phuong phap Divide-and-conquer
165     int divideConquerTED(TreeNode* t1, TreeNode* t2) {
166     if (!t1 && !t2) return 0;
167     if (!t1) return sumInsertCost({t2});

```



```

162     if (!t2) return sumDeleteCost({t1});
163
164     int min_cost = INT_MAX;
165
166     // Thay the goc
167     int substitute_cost = cost_substitute(t1->label, t2->label);
168     substitute_cost += computeChildrenCost(t1->children, t2->
children);
169     min_cost = min(min_cost, substitute_cost);
170
171     // Xoa goc t1
172     int delete_cost = cost_delete(t1->label);
173     for (auto child : t1->children) {
174         delete_cost += divideConquerTED(child, t2);
175     }
176     min_cost = min(min_cost, delete_cost);
177
178     // Chen goc t2
179     int insert_cost = cost_insert(t2->label);
180     for (auto child : t2->children) {
181         insert_cost += divideConquerTED(t1, child);
182     }
183     min_cost = min(min_cost, insert_cost);
184
185     return min_cost;
186 }
187
188 // 4. Phuong phap Dynamic Programming (Zhang-Shasha)
189 int dynamicProgrammingTED(TreeNode* t1, TreeNode* t2) {
190     vector<TreeNode*> post1, post2;
191     getPostOrder(t1, post1);
192     getPostOrder(t2, post2);
193
194     if (post1.empty() && post2.empty()) return 0;
195     if (post1.empty()) return post2.size();
196     if (post2.empty()) return post1.size();
197
198     int n = post1.size();
199     int m = post2.size();
200
201     vector<int> size1, size2;
202     computeSubtreeSizes(post1, size1);
203     computeSubtreeSizes(post2, size2);
204
205     // Bang DP chinh
206     vector<vector<int>> dp(n + 1, vector<int>(m + 1, 0));
207
208     // Khoi tao
209     for (int i = 1; i <= n; i++) {
210         dp[i][0] = dp[i-1][0] + cost_delete(post1[i-1]->label);
211     }
212     for (int j = 1; j <= m; j++) {
213         dp[0][j] = dp[0][j-1] + cost_insert(post2[j-1]->label);
214     }
215
216     // Dien bang DP

```

```

217     for (int i = 1; i <= n; i++) {
218         for (int j = 1; j <= m; j++) {
219             // Xoa
220             dp[i][j] = dp[i-1][j] + cost_delete(post1[i-1]->label);
221
222             // Chen
223             dp[i][j] = min(dp[i][j],
224                             dp[i][j-1] + cost_insert(post2[j-1]->label
225 ));
226
227             // Thay the (neu la cay con)
228             if (isAncestor(post1[i-1], post1) && isAncestor(post2[j
229 -1], post2)) {
230                 dp[i][j] = min(dp[i][j],
231                                 dp[i-1][j-1] + cost_substitute(post1[i
232 -1]->label,
233                                                         post2[j
234 -1]->label));
235             }
236         }
237     }
238     return dp[n][m];
239 }
240
241 private:
242     // Cac ham ho tro
243     int sumInsertCost(const vector<TreeNode*>& nodes) {
244         int total = 0;
245         for (auto node : nodes) {
246             total += cost_insert(node->label);
247             total += sumInsertCost(node->children);
248         }
249         return total;
250     }
251
252     int sumDeleteCost(const vector<TreeNode*>& nodes) {
253         int total = 0;
254         for (auto node : nodes) {
255             total += cost_delete(node->label);
256             total += sumDeleteCost(node->children);
257         }
258         return total;
259     }
260
261     int computeChildrenCost(const vector<TreeNode*>& children1,
262                             const vector<TreeNode*>& children2) {
263         // Su dung quy hoach dong don gian cho danh sach cac cay con
264         int n = children1.size();
265         int m = children2.size();
266
267         if (n == 0) return sumInsertCost(children2);
268         if (m == 0) return sumDeleteCost(children1);
269
270         vector<vector<int>> dp(n + 1, vector<int>(m + 1, 0));

```

```

269     for (int i = 1; i <= n; i++) {
270         dp[i][0] = dp[i-1][0] + sumDeleteCost({children1[i-1]});
271     }
272     for (int j = 1; j <= m; j++) {
273         dp[0][j] = dp[0][j-1] + sumInsertCost({children2[j-1]});
274     }
275
276     for (int i = 1; i <= n; i++) {
277         for (int j = 1; j <= m; j++) {
278             dp[i][j] = min({
279                 dp[i-1][j] + sumDeleteCost({children1[i-1]}),
280                 dp[i][j-1] + sumInsertCost({children2[j-1]}),
281                 dp[i-1][j-1] + divideConquerTED(children1[i-1],
children2[j-1])
282             });
283         }
284     }
285
286     return dp[n][m];
287 }
288
289 void branchBoundChildrenHelper(const vector<TreeNode*>& children1,
290                               const vector<TreeNode*>& children2,
291                               int current_cost, int& best_cost) {
292     int children_cost = computeChildrenCost(children1, children2);
293     best_cost = min(best_cost, current_cost + children_cost);
294 }
295
296 int countNodes(TreeNode* root) {
297     if (!root) return 0;
298     int count = 1;
299     for (auto child : root->children) {
300         count += countNodes(child);
301     }
302     return count;
303 }
304
305 bool isAncestor(TreeNode* node, const vector<TreeNode*>& post_order)
306 {
307     // Kiem tra xem node co phai la to tien cua cac nut khac khong
308     // Day la ham don gian hoa, can cai dat chinh xac hon
309     return true;
310 }
311 public:
312     // Ham tao cay mau de test
313     TreeNode* createSampleTree1() {
314         TreeNode* root = new TreeNode("A", 1);
315         TreeNode* b = new TreeNode("B", 2);
316         TreeNode* c = new TreeNode("C", 3);
317         TreeNode* d = new TreeNode("D", 4);
318         TreeNode* e = new TreeNode("E", 5);
319
320         root->children = {b, c};
321         b->children = {d};
322         c->children = {e};

```

```

323
324     return root;
325 }
326
327 TreeNode* createSampleTree2() {
328     TreeNode* root = new TreeNode("A", 1);
329     TreeNode* c = new TreeNode("C", 2);
330     TreeNode* b = new TreeNode("B", 3);
331     TreeNode* e = new TreeNode("E", 4);
332     TreeNode* f = new TreeNode("F", 5);
333
334     root->children = {c, b};
335     c->children = {e, f};
336
337     return root;
338 }
339
340 // Ham in cay
341 void printTree(TreeNode* root, string indent = "", bool isLast =
342 true) {
343     if (!root) return;
344
345     cout << indent;
346     if (isLast) {
347         cout << "          ";
348         indent += "          ";
349     } else {
350         cout << "          ";
351         indent += "          ";
352     }
353     cout << root->label << endl;
354
355     for (int i = 0; i < root->children.size(); i++) {
356         bool last = (i == root->children.size() - 1);
357         printTree(root->children[i], indent, last);
358     }
359 };
360
361 // Ham main de test
362 int main() {
363     TreeEditDistance ted;
364
365     // Tao hai cay mau
366     TreeNode* tree1 = ted.createSampleTree1();
367     TreeNode* tree2 = ted.createSampleTree2();
368
369     cout << "=== TREE EDIT DISTANCE DEMO ===" << endl << endl;
370
371     cout << "C y 1:" << endl;
372     ted.printTree(tree1);
373     cout << endl;
374
375     cout << "C y 2:" << endl;
376     ted.printTree(tree2);
377     cout << endl;

```

```

378
379 // Test cac phuong phap
380 cout << "=== K T Q U C C PH NG PH P ===" << endl;
381
382 cout << "1. Backtracking: " << ted.backtrackingTED(tree1, tree2) <<
endl;
383 cout << "2. Branch-and-bound: " << ted.branchBoundTED(tree1, tree2)
<< endl;
384 cout << "3. Divide-and-conquer: " << ted.divideConquerTED(tree1,
tree2) << endl;
385 cout << "4. Dynamic Programming: " << ted.dynamicProgrammingTED(
tree1, tree2) << endl;
386
387 // Giai phong bo nho
388 delete tree1;
389 delete tree2;
390
391 return 0;
392 }

```

Listing 5: Cài đặt đầy đủ bằng C++

6 Code Python

```

1 class TreeNode:
2     """ L p b i u d i n n t c y """
3     def __init__(self, label, node_id=0):
4         self.label = label
5         self.id = node_id
6         self.children = []
7
8     def add_child(self, child):
9         self.children.append(child)
10
11     def __str__(self):
12         return f"TreeNode({self.label})"
13
14 class TreeEditDistance:
15     """ L p g i i b i t o n T r e e E d i t D i s t a n c e """
16
17     def __init__(self):
18         self.memo = {} # Cho dynamic programming
19
20     def cost_substitute(self, a, b):
21         """Chi ph thay t h n h n a b n g n h n b"""
22         return 0 if a == b else 1
23
24     def cost_delete(self, a):
25         """Chi ph x a n t c n h n a"""
26         return 1
27
28     def cost_insert(self, b):
29         """Chi ph c h n n t c n h n b"""
30         return 1

```

```

31
32 def get_post_order(self, root):
33     """Duyệt cây theo thứ tự post-order"""
34     if not root:
35         return []
36
37     result = []
38     for child in root.children:
39         result.extend(self.get_post_order(child))
40     result.append(root)
41     return result
42
43 def count_nodes(self, root):
44     """Mức số nút trong cây"""
45     if not root:
46         return 0
47     count = 1
48     for child in root.children:
49         count += self.count_nodes(child)
50     return count
51
52 def sum_insert_cost(self, nodes):
53     """Tổng chi phí chèn danh sách các nút"""
54     total = 0
55     for node in nodes:
56         if node:
57             total += self.cost_insert(node.label)
58             total += self.sum_insert_cost(node.children)
59     return total
60
61 def sum_delete_cost(self, nodes):
62     """Tổng chi phí xóa danh sách các nút"""
63     total = 0
64     for node in nodes:
65         if node:
66             total += self.cost_delete(node.label)
67             total += self.sum_delete_cost(node.children)
68     return total
69
70 # 1. Phương pháp Backtracking
71 def backtracking_ted(self, t1, t2):
72     """
73     Giải Tree Edit Distance bằng Backtracking
74
75     Args:
76         t1, t2: Hai cây con số sinh
77
78     Returns:
79         int: Không có chênh lệch tối thiểu
80     """
81     if not t1 and not t2:
82         return 0
83     if not t1:
84         return self.cost_insert(t2.label) + self.sum_insert_cost(t2.
children)
85     if not t2:

```

```

86         return self.cost_delete(t1.label) + self.sum_delete_cost(t1.
children)
87
88         min_cost = float('inf')
89
90         # Thay t h n t g c
91         substitute_cost = self.cost_substitute(t1.label, t2.label)
92         substitute_cost += self.compute_children_cost(t1.children, t2.
children)
93         min_cost = min(min_cost, substitute_cost)
94
95         # X a n t g c c a t1
96         delete_cost = self.cost_delete(t1.label)
97         for child in t1.children:
98             delete_cost += self.backtracking_ted(child, t2)
99         min_cost = min(min_cost, delete_cost)
100
101         # Ch n n t g c c a t2
102         insert_cost = self.cost_insert(t2.label)
103         for child in t2.children:
104             insert_cost += self.backtracking_ted(t1, child)
105         min_cost = min(min_cost, insert_cost)
106
107         return min_cost
108
109     # 2. Ph ng ph p Branch-and-bound
110     def branch_bound_ted(self, t1, t2):
111         """
112         G i i Tree Edit Distance b ng Branch-and-bound
113
114         Args:
115             t1, t2: Hai c y c n s o s nh
116
117         Returns:
118             int: K h o n g c c h c h nh s a t i t h i u
119         """
120         self.best_cost = float('inf')
121         self._branch_bound_helper(t1, t2, 0)
122         return self.best_cost
123
124     def _branch_bound_helper(self, t1, t2, current_cost):
125         """H m quy cho branch-and-bound"""
126         # c l ng lower bound
127         lower_bound = current_cost + self.estimate_lower_bound(t1, t2)
128
129         if lower_bound >= self.best_cost:
130             return # C t t a nh nh n y
131
132         if not t1 and not t2:
133             self.best_cost = min(self.best_cost, current_cost)
134             return
135
136         if not t1:
137             cost = current_cost + self.sum_insert_cost([t2])
138             self.best_cost = min(self.best_cost, cost)
139             return

```

```

140
141     if not t2:
142         cost = current_cost + self.sum_delete_cost([t1])
143         self.best_cost = min(self.best_cost, cost)
144         return
145
146     # Thay thế node t1 bằng node t2
147     # Thay thế t1
148     new_cost = current_cost + self.cost_substitute(t1.label, t2.
label)
149     children_cost = self.compute_children_cost(t1.children, t2.
children)
150     self.best_cost = min(self.best_cost, new_cost + children_cost)
151
152     # Xóa node t1
153     for child in t1.children:
154         self._branch_bound_helper(child, t2,
155                                   current_cost + self.cost_delete(t1.
label))
156
157     # Chèn node t2
158     for child in t2.children:
159         self._branch_bound_helper(t1, child,
160                                   current_cost + self.cost_insert(t2.
label))
161
162     def estimate_lower_bound(self, t1, t2):
163         """ Tính toán lower bound cho branch-and-bound """
164         size1 = self.count_nodes(t1)
165         size2 = self.count_nodes(t2)
166         return abs(size1 - size2)
167
168     # 3. Phân hoạch Divide-and-conquer
169     def divide_conquer_ted(self, t1, t2):
170         """
171         Giải bài toán Tree Edit Distance bằng Divide-and-conquer
172
173         Args:
174             t1, t2: Hai cây nhị phân
175
176         Returns:
177             int: Kết quả chi phí tối thiểu
178         """
179         if not t1 and not t2:
180             return 0
181         if not t1:
182             return self.sum_insert_cost([t2])
183         if not t2:
184             return self.sum_delete_cost([t1])
185
186         min_cost = float('inf')
187
188         # Thay thế node gốc
189         substitute_cost = self.cost_substitute(t1.label, t2.label)
190         substitute_cost += self.compute_children_cost(t1.children, t2.
children)

```



```

191     min_cost = min(min_cost, substitute_cost)
192
193     # X a g c t1
194     delete_cost = self.cost_delete(t1.label)
195     for child in t1.children:
196         delete_cost += self.divide_conquer_ted(child, t2)
197     min_cost = min(min_cost, delete_cost)
198
199     # Ch n g c t2
200     insert_cost = self.cost_insert(t2.label)
201     for child in t2.children:
202         insert_cost += self.divide_conquer_ted(t1, child)
203     min_cost = min(min_cost, insert_cost)
204
205     return min_cost
206
207 # 4. Ph n g p h p Dynamic Programming
208 def dynamic_programming_ted(self, t1, t2):
209     """
210     G i i Tree Edit Distance b n g Dynamic Programming (Zhang-
211     Shasha)
212
213     Args:
214         t1, t2: Hai c y c n s o s n h
215
216     Returns:
217         int: K h o n g c c h c h n h s a t i t h i u
218     """
219     if not t1 and not t2:
220         return 0
221     if not t1:
222         return self.count_nodes(t2)
223     if not t2:
224         return self.count_nodes(t1)
225
226     # Ch u y n i c y t h n h d n g post-order
227     post1 = self.get_post_order(t1)
228     post2 = self.get_post_order(t2)
229
230     n, m = len(post1), len(post2)
231
232     # K h i t o b n g DP
233     dp = [[0] * (m + 1) for _ in range(n + 1)]
234
235     # T r n g h p c s
236     for i in range(1, n + 1):
237         dp[i][0] = dp[i-1][0] + self.cost_delete(post1[i-1].label)
238     for j in range(1, m + 1):
239         dp[0][j] = dp[0][j-1] + self.cost_insert(post2[j-1].label)
240
241     # i n b n g DP
242     for i in range(1, n + 1):
243         for j in range(1, m + 1):
244             # X a
245             cost_del = dp[i-1][j] + self.cost_delete(post1[i-1].
label)

```

```

245
246         # Chèn
247         cost_ins = dp[i][j-1] + self.cost_insert(post2[j-1].
label)
248
249         # Thay thế
250         cost_sub = (dp[i-1][j-1] +
251                     self.cost_substitute(post1[i-1].label, post2[
j-1].label))
252
253         dp[i][j] = min(cost_del, cost_ins, cost_sub)
254
255     return dp[n][m]
256
257     def compute_children_cost(self, children1, children2):
258         """Tinh chi phí tối ưu biến đổi danh sách cây
con"""
259         n, m = len(children1), len(children2)
260
261         if n == 0:
262             return self.sum_insert_cost(children2)
263         if m == 0:
264             return self.sum_delete_cost(children1)
265
266         # Sắp xếp DP để tính giá cho danh sách cây con
267         dp = [[0] * (m + 1) for _ in range(n + 1)]
268
269         for i in range(1, n + 1):
270             dp[i][0] = dp[i-1][0] + self.sum_delete_cost([children1[i
-1]])
271         for j in range(1, m + 1):
272             dp[0][j] = dp[0][j-1] + self.sum_insert_cost([children2[j
-1]])
273
274         for i in range(1, n + 1):
275             for j in range(1, m + 1):
276                 delete_cost = dp[i-1][j] + self.sum_delete_cost([
children1[i-1]])
277                 insert_cost = dp[i][j-1] + self.sum_insert_cost([
children2[j-1]])
278                 match_cost = (dp[i-1][j-1] +
279                               self.divide_conquer_ted(children1[i-1],
children2[j-1]))
280
281                 dp[i][j] = min(delete_cost, insert_cost, match_cost)
282
283     return dp[n][m]
284
285     def print_tree(self, root, indent="", is_last=True):
286         """In cây đã tính giá"""
287         if not root:
288             return
289
290         print(indent, end="")
291         if is_last:
292             print(" ", end="")

```

```

293         indent += "    "
294     else:
295         print("                ", end="")
296         indent += "    "
297     print(root.label)
298
299     for i, child in enumerate(root.children):
300         last = (i == len(root.children) - 1)
301         self.print_tree(child, indent, last)
302
303     def create_sample_tree1(self):
304         """ T o c y m u 1          test """
305         root = TreeNode("A", 1)
306         b = TreeNode("B", 2)
307         c = TreeNode("C", 3)
308         d = TreeNode("D", 4)
309         e = TreeNode("E", 5)
310
311         root.add_child(b)
312         root.add_child(c)
313         b.add_child(d)
314         c.add_child(e)
315
316         return root
317
318     def create_sample_tree2(self):
319         """ T o c y m u 2          test """
320         root = TreeNode("A", 1)
321         c = TreeNode("C", 2)
322         b = TreeNode("B", 3)
323         e = TreeNode("E", 4)
324         f = TreeNode("F", 5)
325
326         root.add_child(c)
327         root.add_child(b)
328         c.add_child(e)
329         c.add_child(f)
330
331         return root
332
333     def main():
334         """ H m main          test c c ph ng ph p """
335         ted = TreeEditDistance()
336
337         # T o hai c y m u
338         tree1 = ted.create_sample_tree1()
339         tree2 = ted.create_sample_tree2()
340
341         print("=== TREE EDIT DISTANCE DEMO ===\n")
342
343         print("C y 1:")
344         ted.print_tree(tree1)
345         print()
346
347         print("C y 2:")
348         ted.print_tree(tree2)

```

```

349 print()
350
351 # Test c c ph ng ph p
352 print("=== K T Q U C C PH NG PH P ===")
353
354 print(f"1. Backtracking: {ted.backtracking_ted(tree1, tree2)}")
355 print(f"2. Branch-and-bound: {ted.branch_bound_ted(tree1, tree2)}")
356 print(f"3. Divide-and-conquer: {ted.divide_conquer_ted(tree1, tree2)}")
357
358 print(f"4. Dynamic Programming: {ted.dynamic_programming_ted(tree1, tree2)}")
359
360 if __name__ == "__main__":
361     main()

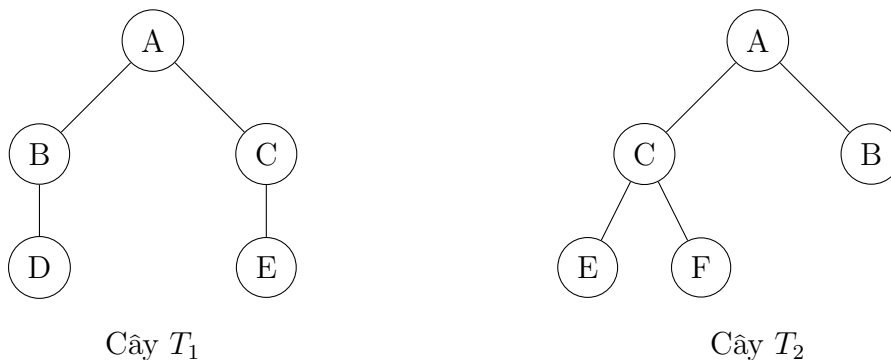
```

Listing 6: Cài đặt đầy đủ bằng Python

7 Ví dụ minh họa

7.1 Ví dụ 1: Cây đơn giản

Xét hai cây sau:

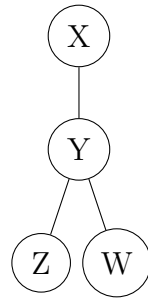


Phân tích các phép biến đổi:

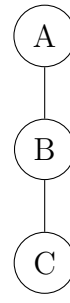
1. Giữ nguyên nút A: Chi phí = 0
2. Hoán đổi vị trí B và C: Không phát sinh chi phí trực tiếp
3. Thêm nút F vào cây con C: Chi phí = 1 (chèn F)
4. Xóa nút D: Chi phí = 1 (xóa D)

Tổng chi phí tối thiểu: $\delta(T_1, T_2) = 2$

7.2 Ví dụ 2: Trường hợp phức tạp



Cây T_3



Cây T_4

Các bước biến đổi tối ưu:

1. Thay thế $X \rightarrow A$: Chi phí = 1
2. Thay thế $Y \rightarrow B$: Chi phí = 1
3. Thay thế $Z \rightarrow C$: Chi phí = 1
4. Xóa nút W : Chi phí = 1

Tổng chi phí: $\delta(T_3, T_4) = 4$

8 Phân tích độ phức tạp

8.1 Độ phức tạp thời gian

Phương pháp	Trường hợp tốt nhất	Trường hợp trung bình	Trường hợp xấu nhất
Backtracking	$O(2^{\min(n,m)})$	$O(3^{(n+m)/2})$	$O(3^{n+m})$
Branch-and-bound	$O(2^{\min(n,m)})$	$O(3^{(n+m)/3})$	$O(3^{n+m})$
Divide-and-conquer	$O(n \cdot m)$	$O(n \cdot m \cdot \log(n+m))$	$O(n \cdot m \cdot (n+m))$
Dynamic Programming	$O(n^2 \cdot m^2)$	$O(n^2 \cdot m^2)$	$O(n^2 \cdot m^2)$

Bảng 2: Phân tích độ phức tạp thời gian chi tiết

8.2 Độ phức tạp không gian

- **Backtracking**: $O(h)$ với h là chiều cao cây
- **Branch-and-bound**: $O(h + \text{size của priority queue})$
- **Divide-and-conquer**: $O(h)$ cho call stack
- **Dynamic Programming**: $O(n \cdot m)$ cho bảng DP

9 Tối ưu hóa và cải tiến

9.1 Tối ưu hóa bộ nhớ cho DP

Dynamic Programming có thể được tối ưu hóa bộ nhớ:

```

1 // Sử dụng rolling array thay vì bảng 2D đầy đủ
2 vector<int> prev(m + 1), curr(m + 1);
3
4 for (int i = 0; i <= n; i++) {
5     for (int j = 0; j <= m; j++) {
6         if (i == 0) {
7             curr[j] = j;
8         } else if (j == 0) {
9             curr[j] = i;
10        } else {
11            curr[j] = min({
12                prev[j] + cost_delete(T1[i-1]),
13                curr[j-1] + cost_insert(T2[j-1]),
14                prev[j-1] + cost_substitute(T1[i-1], T2[j-1])
15            });
16        }
17    }
18    prev = curr;
19 }

```

Listing 7: Tối ưu hóa bộ nhớ

9.2 Cải tiến Branch-and-bound

- **Lower bound tốt hơn:** Sử dụng relaxation của bài toán con
- **Heuristic ordering:** Ưu tiên các nhánh có khả năng tối ưu cao
- **Memoization:** Lưu trữ kết quả các bài toán con đã giải

10 Kết luận

10.1 Tổng kết các phương pháp

Bài toán Tree Edit Distance được giải quyết thành công bằng bốn phương pháp:

1. **Backtracking:** Đơn giản nhưng chậm, phù hợp cho cây nhỏ
2. **Branch-and-bound:** Cải thiện backtracking, hiệu quả hơn nhưng vẫn exponential
3. **Divide-and-conquer:** Cân bằng giữa thời gian và không gian
4. **Dynamic Programming:** Hiệu quả nhất về thời gian, polynomial complexity

10.2 Lựa chọn phương pháp

- **Cây nhỏ** ($n, m < 10$): Backtracking hoặc Branch-and-bound
- **Cây trung bình** ($10 \leq n, m \leq 100$): Divide-and-conquer
- **Cây lớn** ($n, m > 100$): Dynamic Programming
- **Bộ nhớ hạn chế**: Divide-and-conquer hoặc Backtracking

10.3 Ứng dụng thực tế

Tree Edit Distance có ứng dụng rộng rãi trong:

- **Công nghệ sinh học**: So sánh cấu trúc RNA/DNA
- **Xử lý ngôn ngữ tự nhiên**: Phân tích cú pháp
- **Web mining**: So sánh cấu trúc HTML/XML
- **Compiler design**: Tối ưu hóa AST
- **Image processing**: Nhận dạng mẫu