

課題 3

「文字列処理と動的計画法」

鈴木 潤, 吉留 崇

2019年度プログラミング演習A

本課題で学ぶこと

■ 編集距離

◆ 文字列処理 ← 課題 1

◆ 再帰的手続き ← 課題2

◆ 動的計画法

- ・ 問題3-0: 編集距離の解説
- ・ 問題内にも解説あり

編集距離 (Edit Distance)

2つの文字列がどのくらい似ているかを表す指標

(例) 「マテオ」、「マンガ」

どちらが「マリオ」に似ている？

マ	リ	オ	マ	リ	オ	「マテオ」の方が 似ている
マ	テ	オ	マ	ン	ガ	

編集距離
(後で定義)

1

2

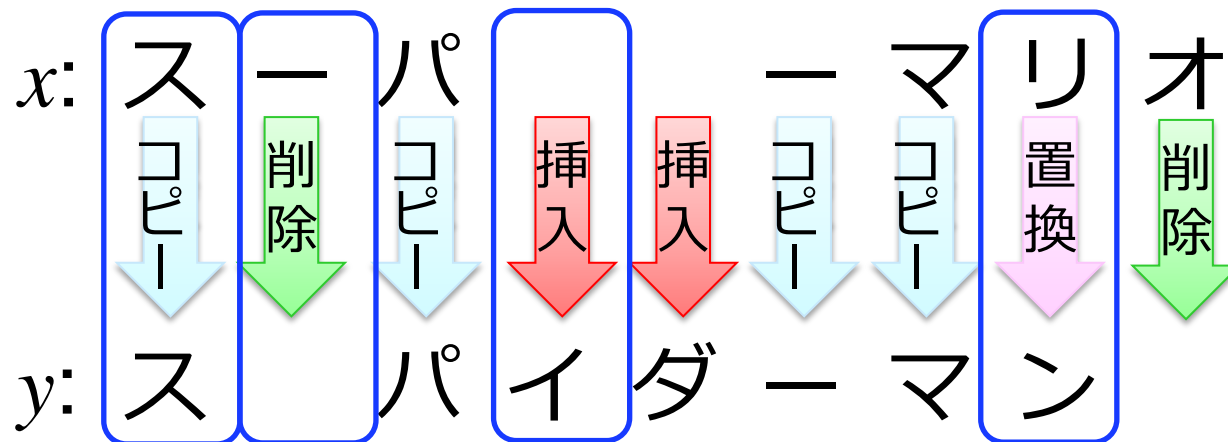
スペルチェッカーやDNAの配列の解析に応用

編集距離 (Edit Distance)

文字列 x , y の編集距離を求めるために、編集操作を行う

文字列 x にコピー、置換、挿入、削除を行い、
文字列 y に等しい文字列を得る操作

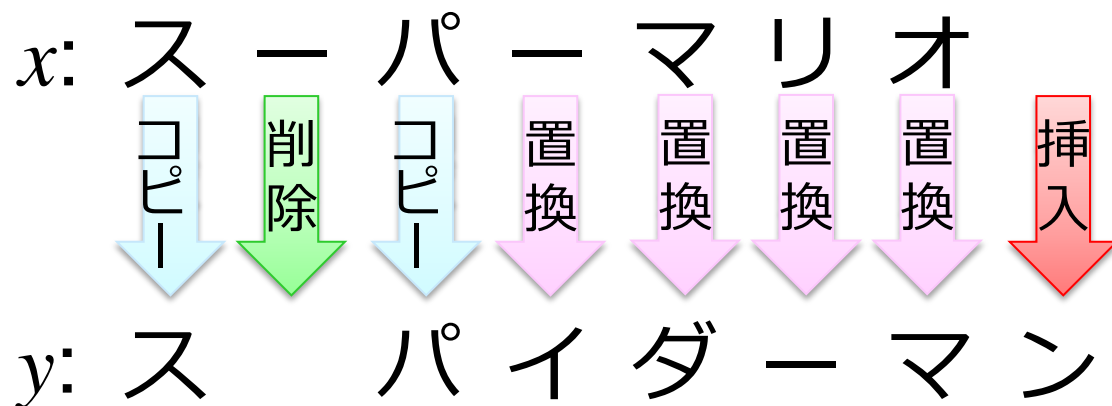
(例) x : スーパーマリオ, y : スパイダーマン



編集距離 (Edit Distance)

編集操作は、複数存在

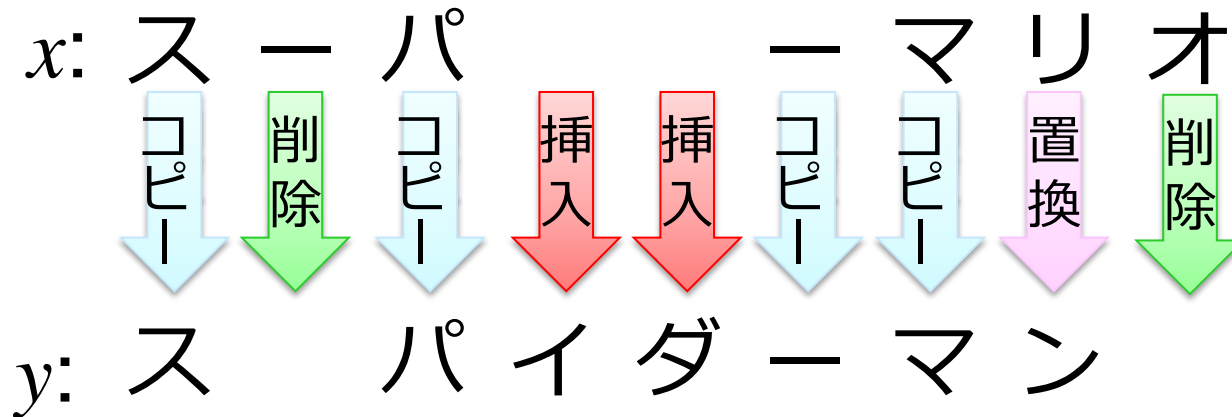
(例) x : スーパーマリオ, y : スパイダーマン



編集距離 (Edit Distance)

文字列 x を文字列 y に変換する時の「削除」, 「挿入」, 「置換」の最小回数 (「コピー」はカウントしない)

(例) x : スーパーマリオ, y : スパイダーマン



編集距離 = 5

編集距離の計算方法

文字列に対する数学記号

	1	2	3	4	5	6	7	8	9
x :	a	b	b	a	b	a	a	a	b

文字列 x の i 番目の文字： x_i 例) $x_4 = a$

文字列 x の先頭から i 番目までの文字列： X_i

例) $X_4 = abba$



接頭辞

長さ 0 の文字列： ε

編集距離の計算方法

文字列 $x = \langle x_1, x_2, \dots, x_m \rangle$, $y = \langle y_1, y_2, \dots, y_n \rangle$ の
接頭辞 X_i, Y_j の編集距離: $c_{i,j}$

$$c_{i,j} = \begin{cases} \max(i, j) & (i = 0 \text{ または } j = 0 \text{ の時}) \\ \min(c_{i-1, j-1} + d(x_i, y_j), \\ \quad c_{i-1, j} + 1, \\ \quad c_{i, j-1} + 1) & (\text{その他}) \end{cases}$$

$$d(x_i, y_j) = \begin{cases} 1 & (x_i \neq y_j) \\ 0 & (x_i = y_j) \end{cases}$$

* クロネッカーのデルタ
ではない事に注意

編集距離の計算方法

$$c_{i,j} = \min \left(\begin{array}{l} c_{i-1,j-1} + d(x_i, y_j), \\ c_{i-1,j} + 1, \\ c_{i,j-1} + 1 \end{array} \right), \quad \begin{cases} 1 & (x_i \neq y_j) \\ 0 & (x_i = y_j) \end{cases}$$

置換
削除
挿入

* クロネッカーのデルタではない事に注意

(例) $x=ab$ 、 $y=cd$

$c_{2,2} \langle ab, cd \rangle$ の編集距離は、以下の距離の最小値

$c_{1,1} \langle a, c \rangle + d(b, d)$

a b
↓ 置換
c d

$c_{1,2} \langle a, cd \rangle + 1$

a b
↓ 削除
c d

$c_{2,1} \langle ab, c \rangle + 1$

a b
↓ 挿入
c d

編集距離の計算方法

$$c_{i,j} = \max(i, j) \quad (i = 0 \text{ または } j = 0 \text{ の時})$$

$$c_{i,0} \langle X_i, \varepsilon \rangle = i, \quad \text{削除}$$

$x :$ **a** **b**
 ↓ ↓
 削除 削除
 $y :$

$$c_{0,j} \langle \varepsilon, Y_j \rangle = j, \quad \text{挿入}$$

$x :$
 ↓ ↓
 挿入 挿入
 $y :$ **c** **d**

編集距離の計算方法

$$c_{i,j} = \min(c_{i-1,j-1} + d(x_i, y_j), \begin{cases} 1 (x_i \neq y_j) & \text{置換} \\ 0 (x_i = y_j) & \text{コピー} \end{cases}, c_{i-1,j} + 1, c_{i,j-1} + 1)$$

削除挿入

$$c_{i,0} = i, \quad \text{削除} \quad c_{0,j} = j, \quad \text{挿入}$$

再帰関数による実装

動的計画法による実装

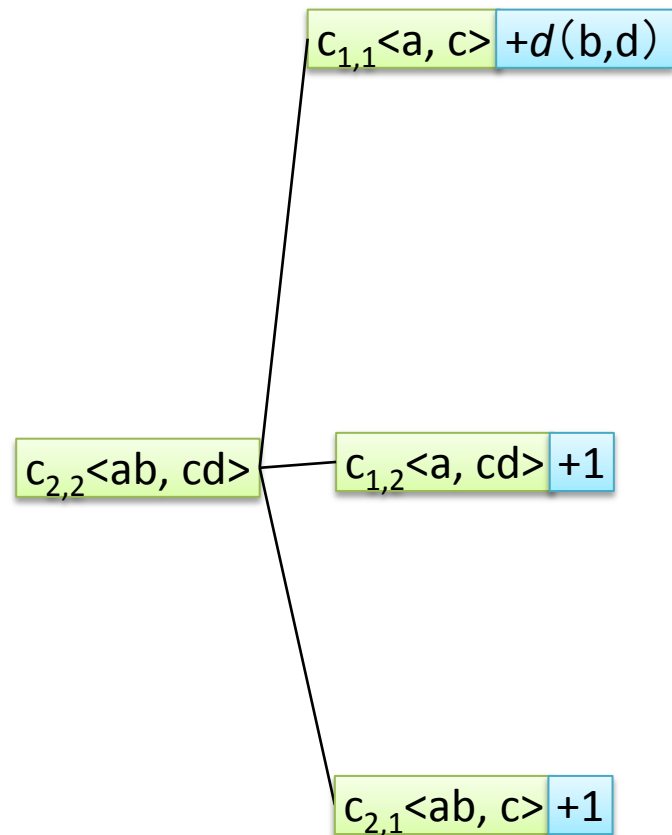
再帰的に計算

$c_{2,2} \langle ab, cd \rangle$

$$c_{i,j} \langle X_i, Y_j \rangle = \min(c_{i-1,j-1} + d(x_i, y_j), \\ c_{i-1,j} + 1, \\ c_{i,j-1} + 1)$$

$$c_{i,0} \langle X_i, \varepsilon \rangle = i, \quad c_{0,j} \langle \varepsilon, Y_j \rangle = j$$

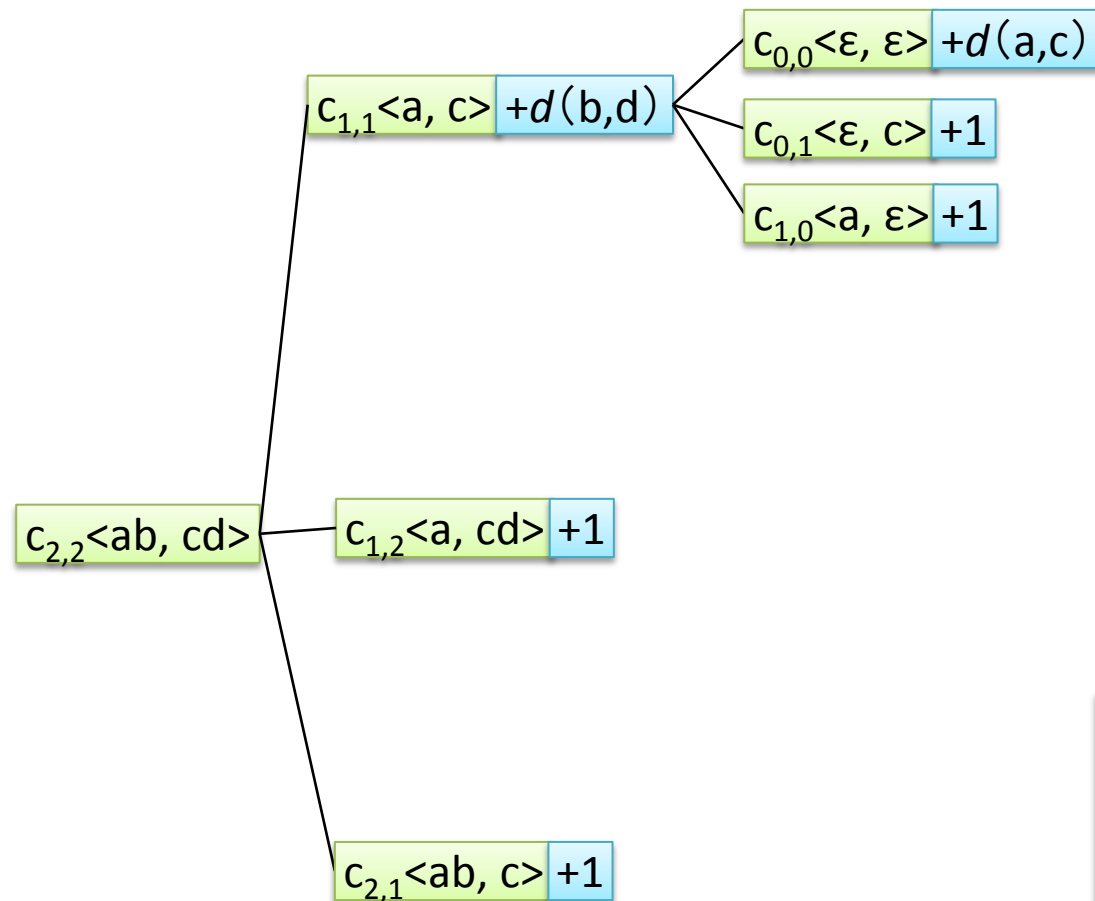
再帰的に計算



$$c_{i,j} \langle X_i, Y_j \rangle = \min(c_{i-1,j-1} + d(x_i, y_j), \\ c_{i-1,j} + 1, \\ c_{i,j-1} + 1)$$

$$c_{i,0} \langle X_i, \varepsilon \rangle = i, \quad c_{0,j} \langle \varepsilon, Y_j \rangle = j$$

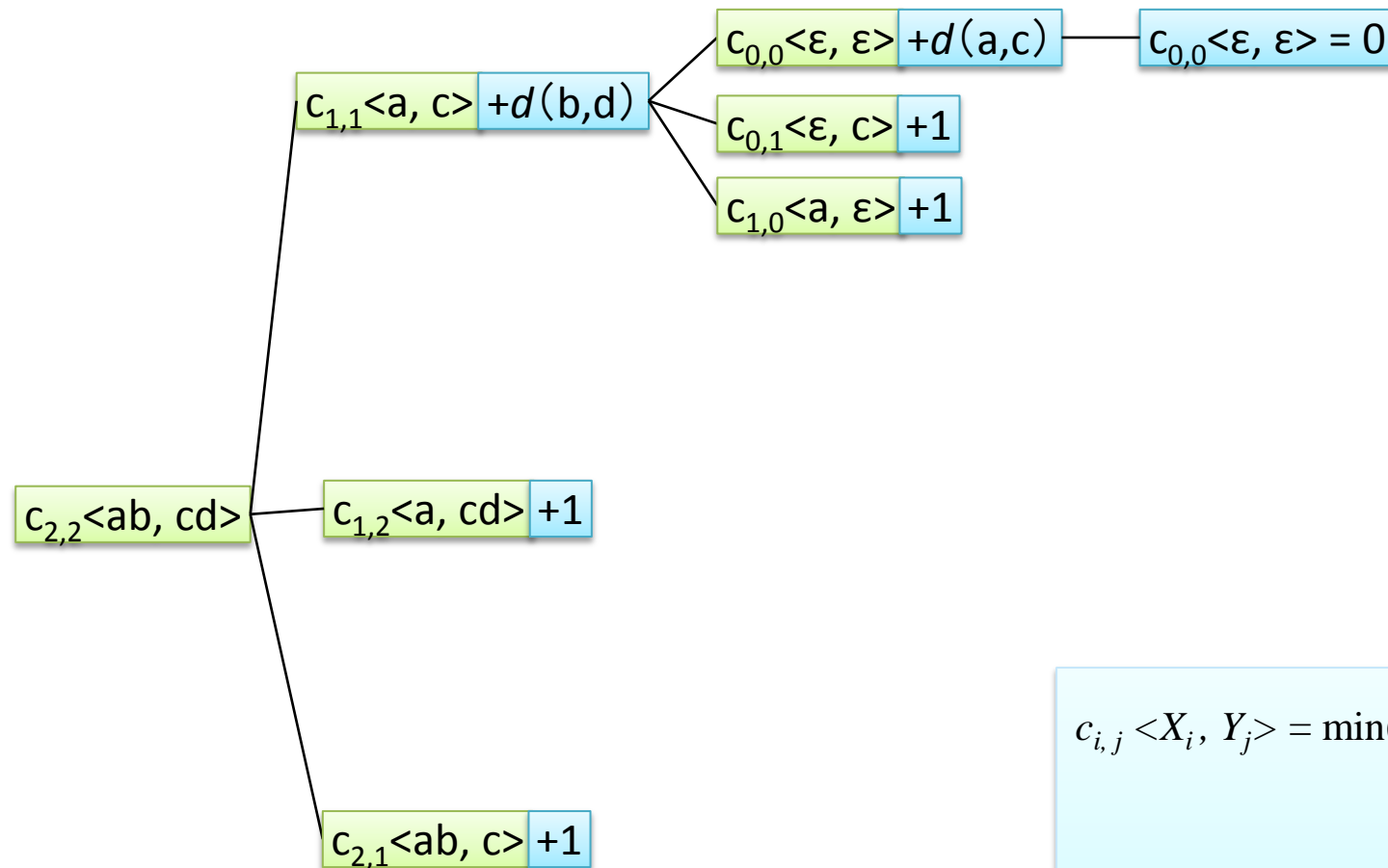
再帰的に計算



$$c_{i,j} \langle X_i, Y_j \rangle = \min(c_{i-1,j-1} + d(x_i, y_j), \\ c_{i-1,j} + 1, \\ c_{i,j-1} + 1)$$

$$c_{i,0} \langle X_i, \varepsilon \rangle = i, \quad c_{0,j} \langle \varepsilon, Y_j \rangle = j$$

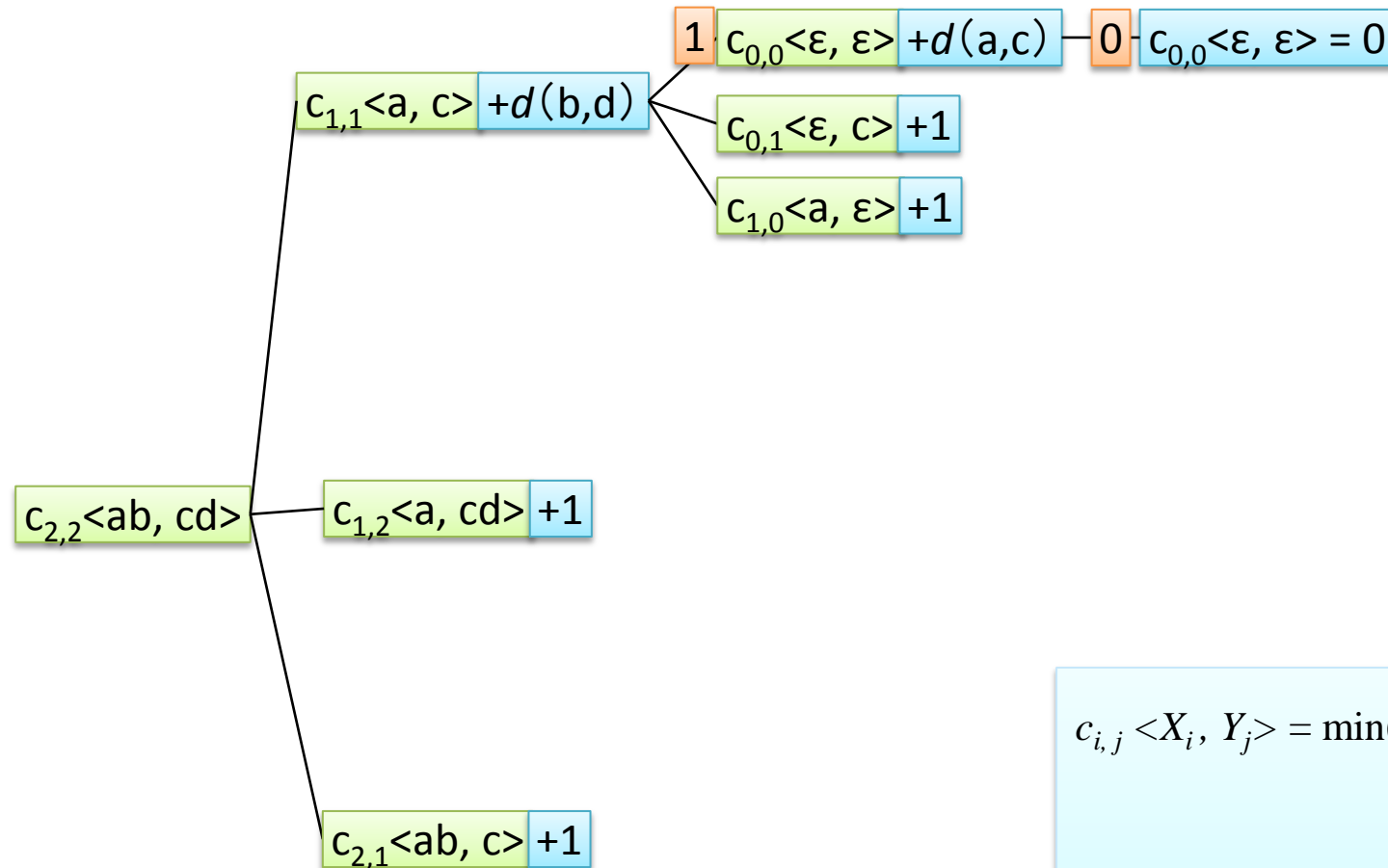
再帰的に計算



$$c_{i,j} \langle X_i, Y_j \rangle = \min(c_{i-1,j-1} + d(x_i, y_j), \\ c_{i-1,j} + 1, \\ c_{i,j-1} + 1)$$

$$c_{i,0} \langle X_i, \varepsilon \rangle = i, \quad c_{0,j} \langle \varepsilon, Y_j \rangle = j$$

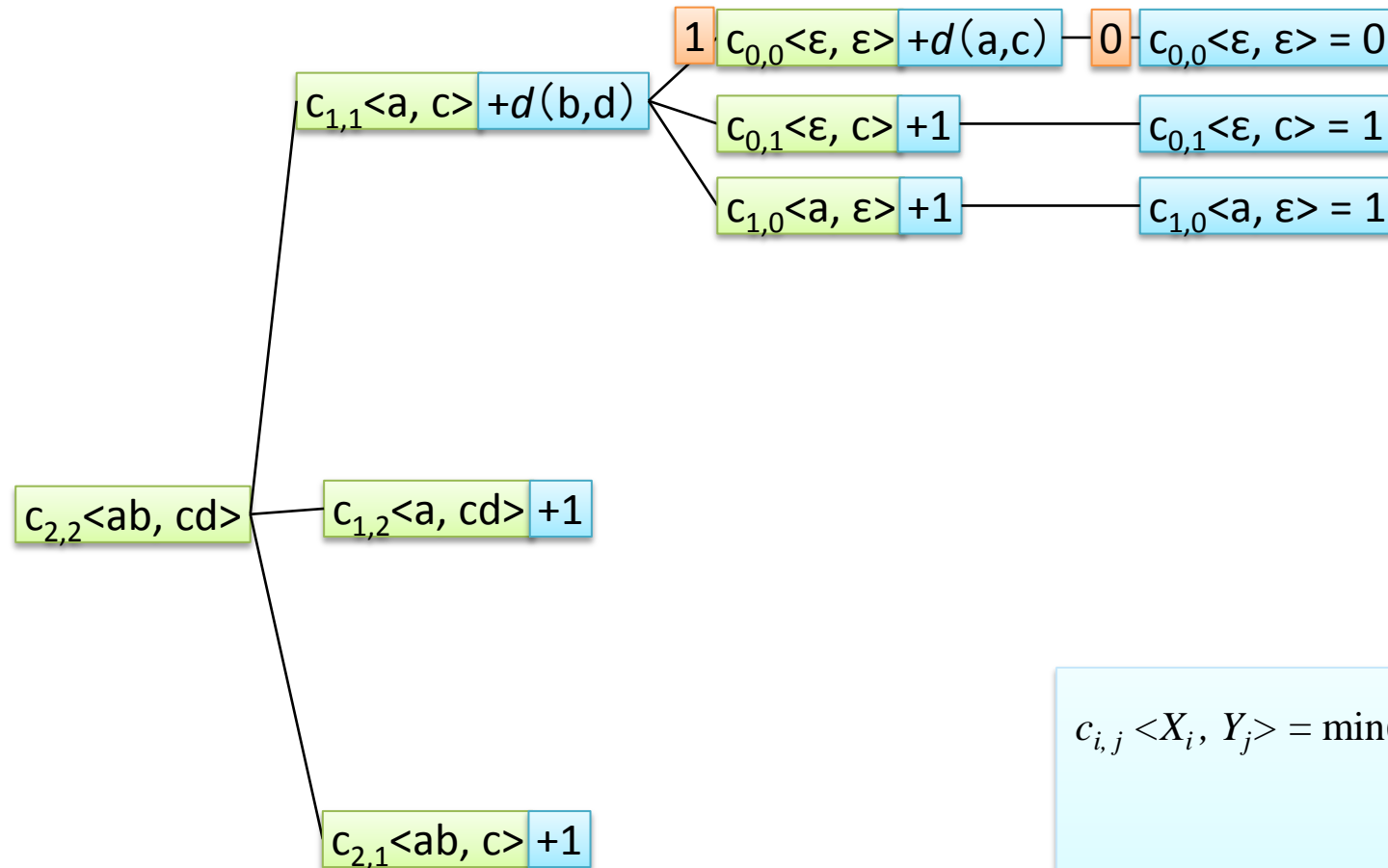
再帰的に計算



$$c_{i,j} \langle X_i, Y_j \rangle = \min(c_{i-1,j-1} + d(x_i, y_j), c_{i-1,j} + 1, c_{i,j-1} + 1)$$

$$c_{i,0} \langle X_i, \epsilon \rangle = i, \quad c_{0,j} \langle \epsilon, Y_j \rangle = j$$

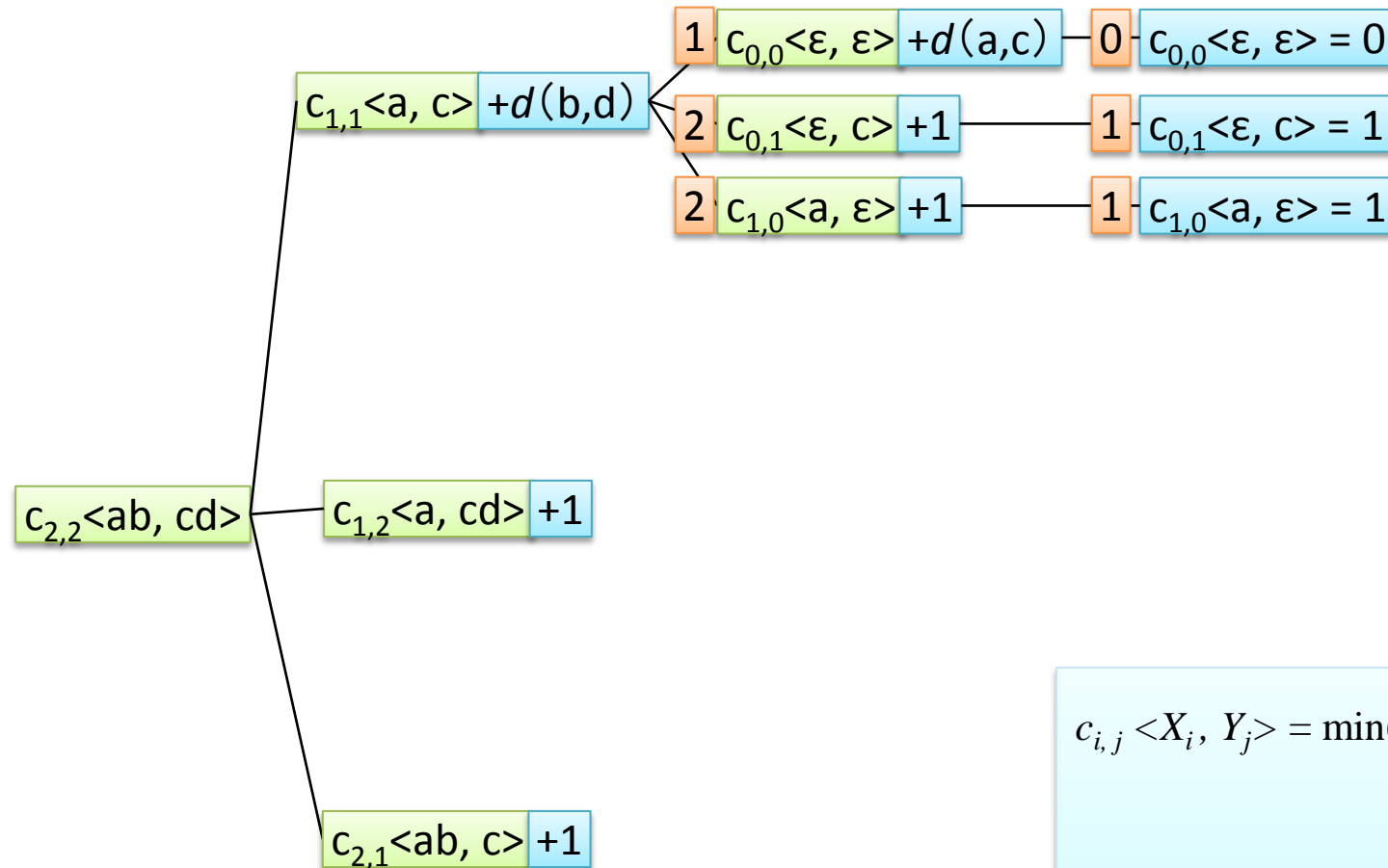
再帰的に計算



$$c_{i,j} \langle X_i, Y_j \rangle = \min(c_{i-1,j-1} + d(x_i, y_j), c_{i-1,j} + 1, c_{i,j-1} + 1)$$

$$c_{i,0} \langle X_i, \epsilon \rangle = i, \quad c_{0,j} \langle \epsilon, Y_j \rangle = j$$

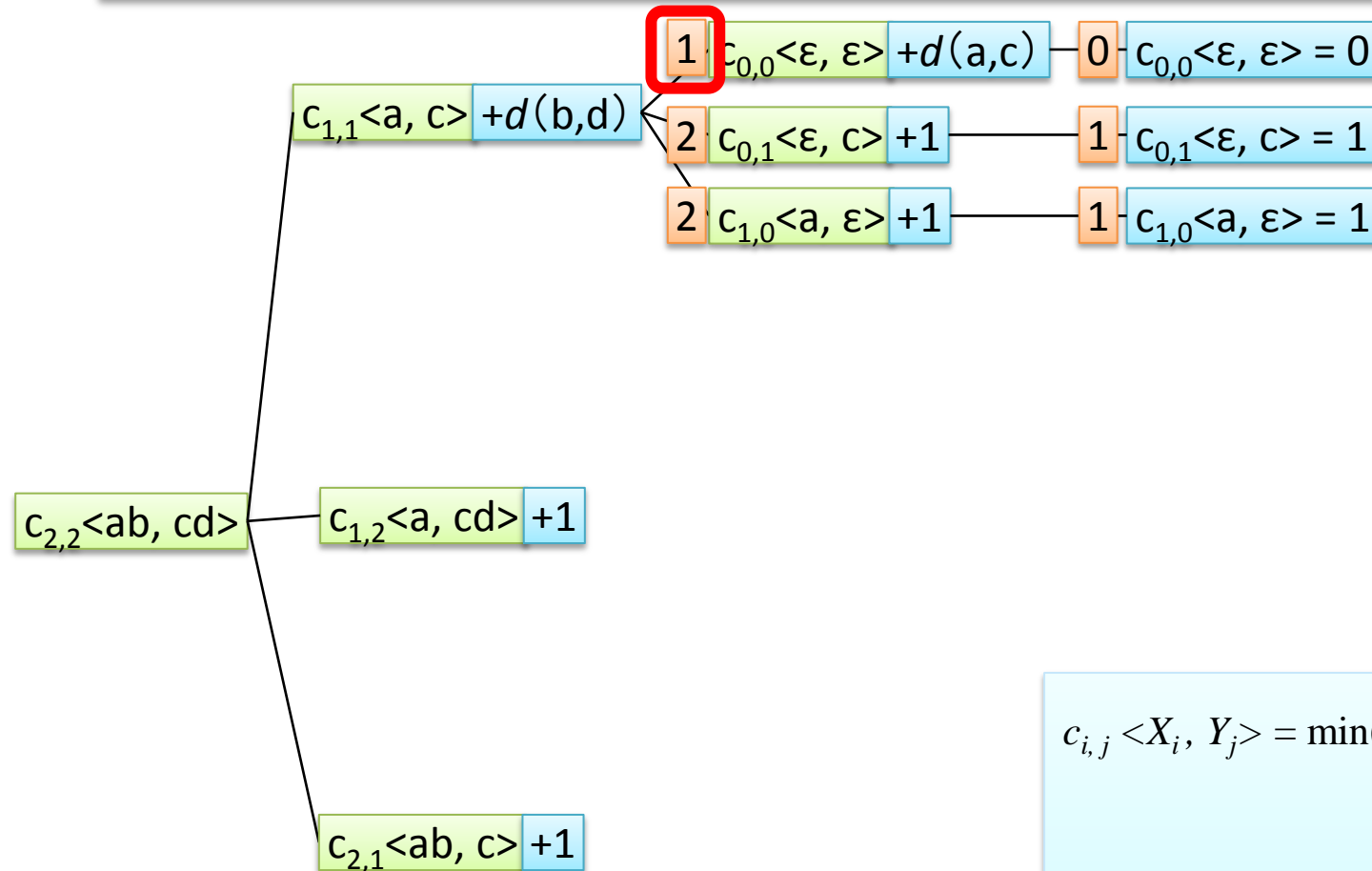
再帰的に計算



$$c_{i,j} \langle X_i, Y_j \rangle = \min(c_{i-1,j-1} + d(x_i, y_j), \\ c_{i-1,j} + 1, \\ c_{i,j-1} + 1)$$

$$c_{i,0} \langle X_i, \epsilon \rangle = i, \quad c_{0,j} \langle \epsilon, Y_j \rangle = j$$

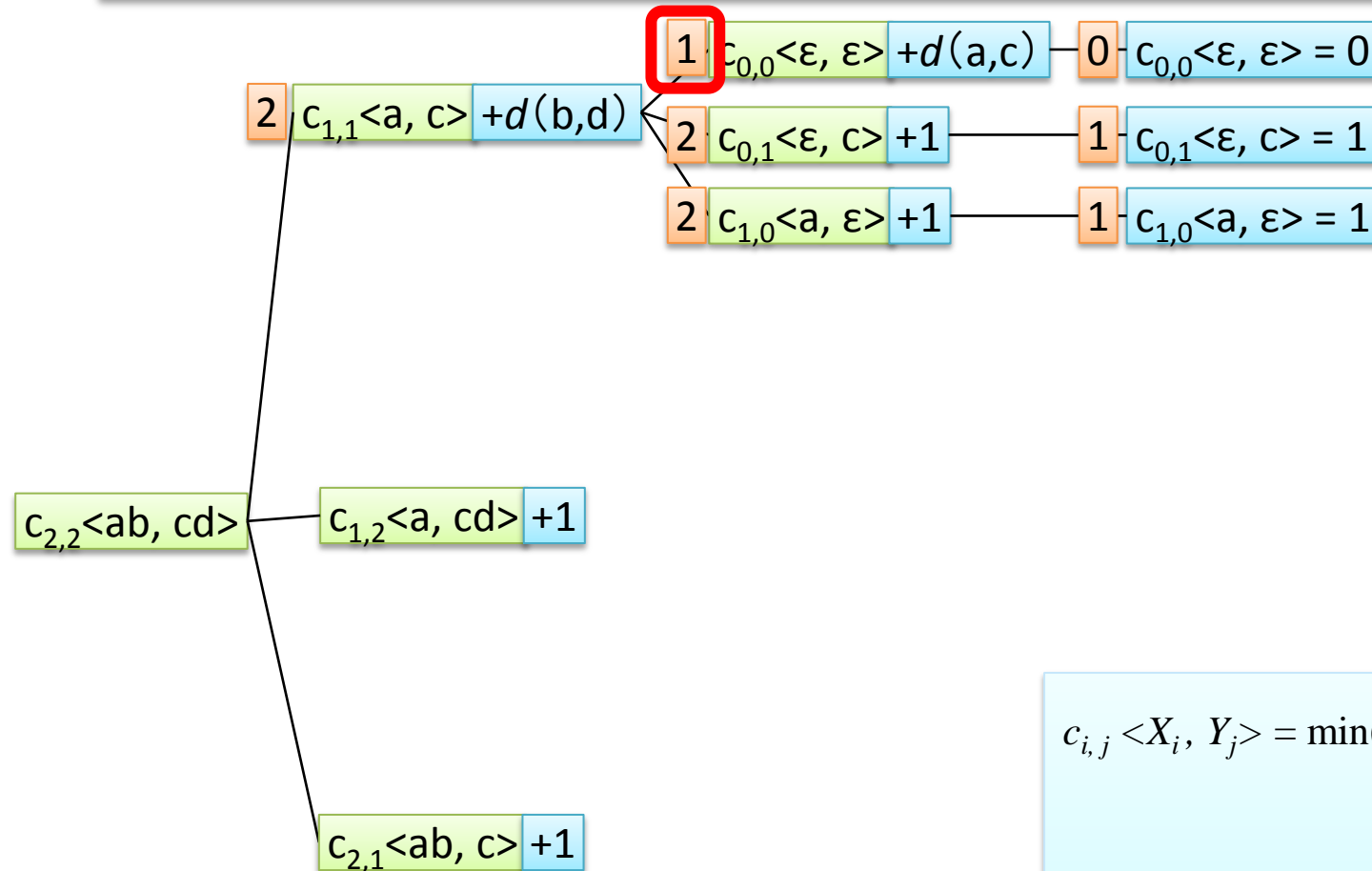
再帰的に計算



$$c_{i,j} \langle X_i, Y_j \rangle = \min(c_{i-1,j-1} + d(x_i, y_j), c_{i-1,j} + 1, c_{i,j-1} + 1)$$

$$c_{i,0} \langle X_i, \epsilon \rangle = i, \quad c_{0,j} \langle \epsilon, Y_j \rangle = j$$

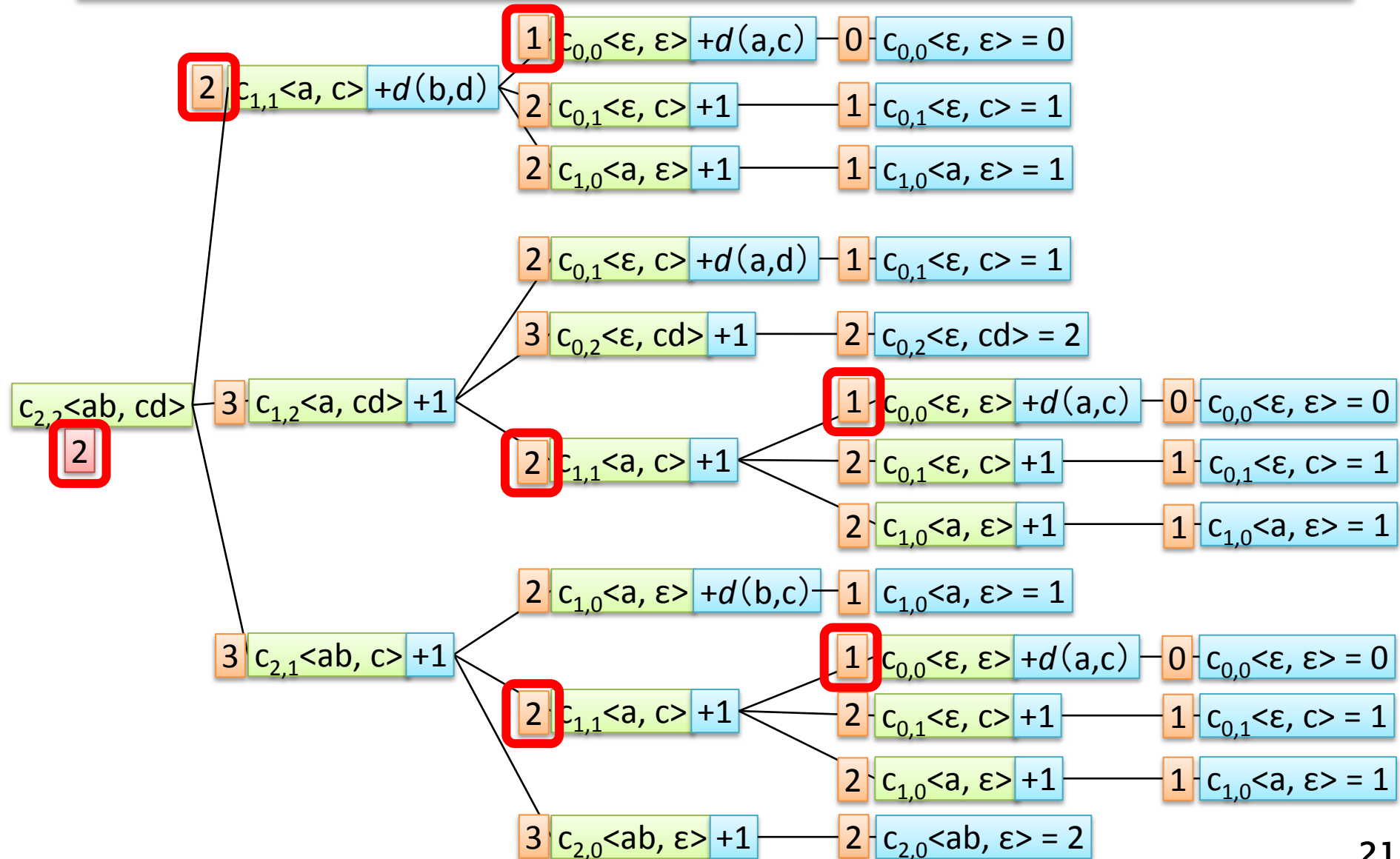
再帰的に計算



$$c_{i,j} <X_i, Y_j> = \min(c_{i-1,j-1} + d(x_i, y_j), c_{i-1,j} + 1, c_{i,j-1} + 1)$$

$$c_{i,0} <X_i, \epsilon> = i, \quad c_{0,j} <\epsilon, Y_j> = j$$

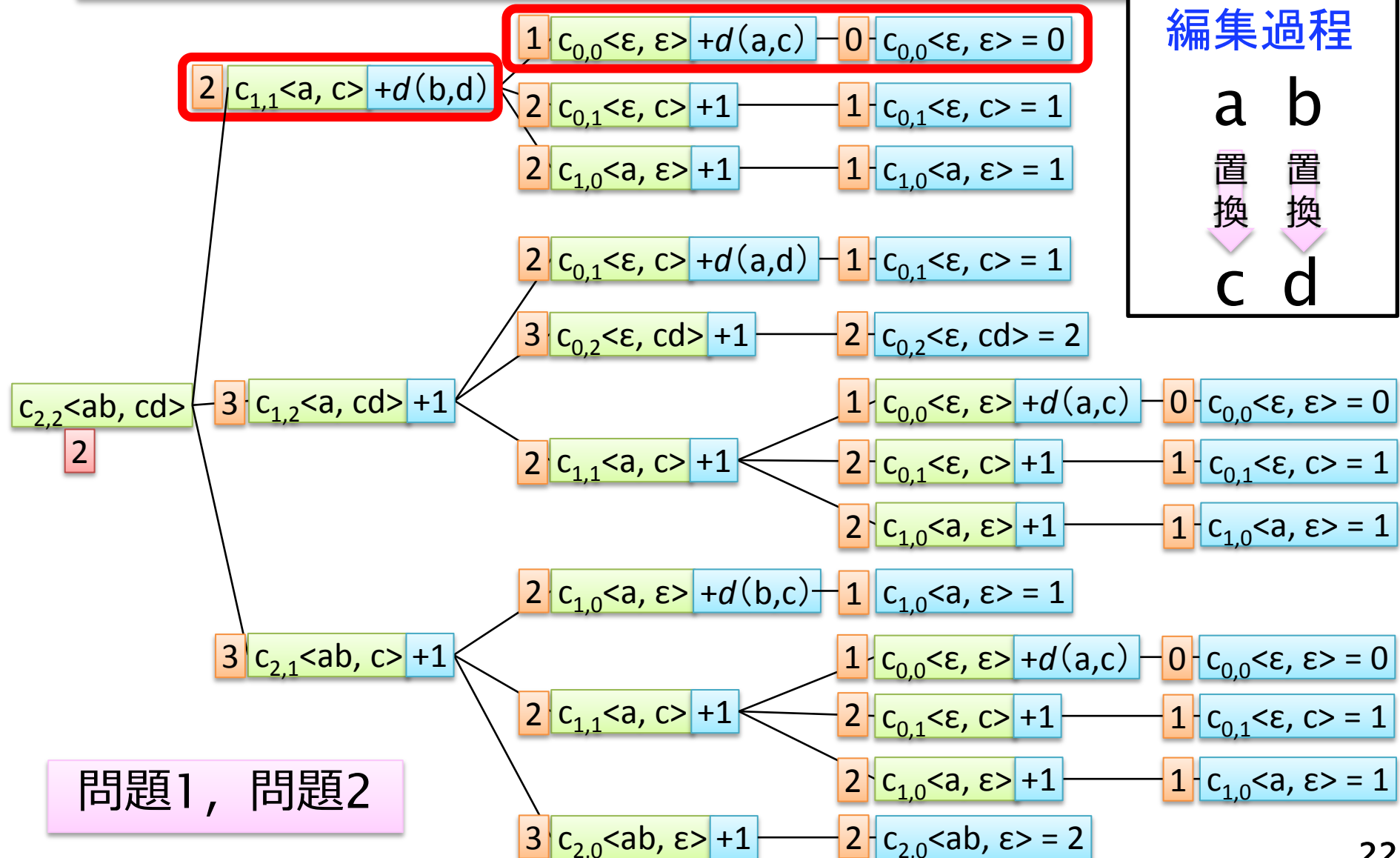
再帰的に計算



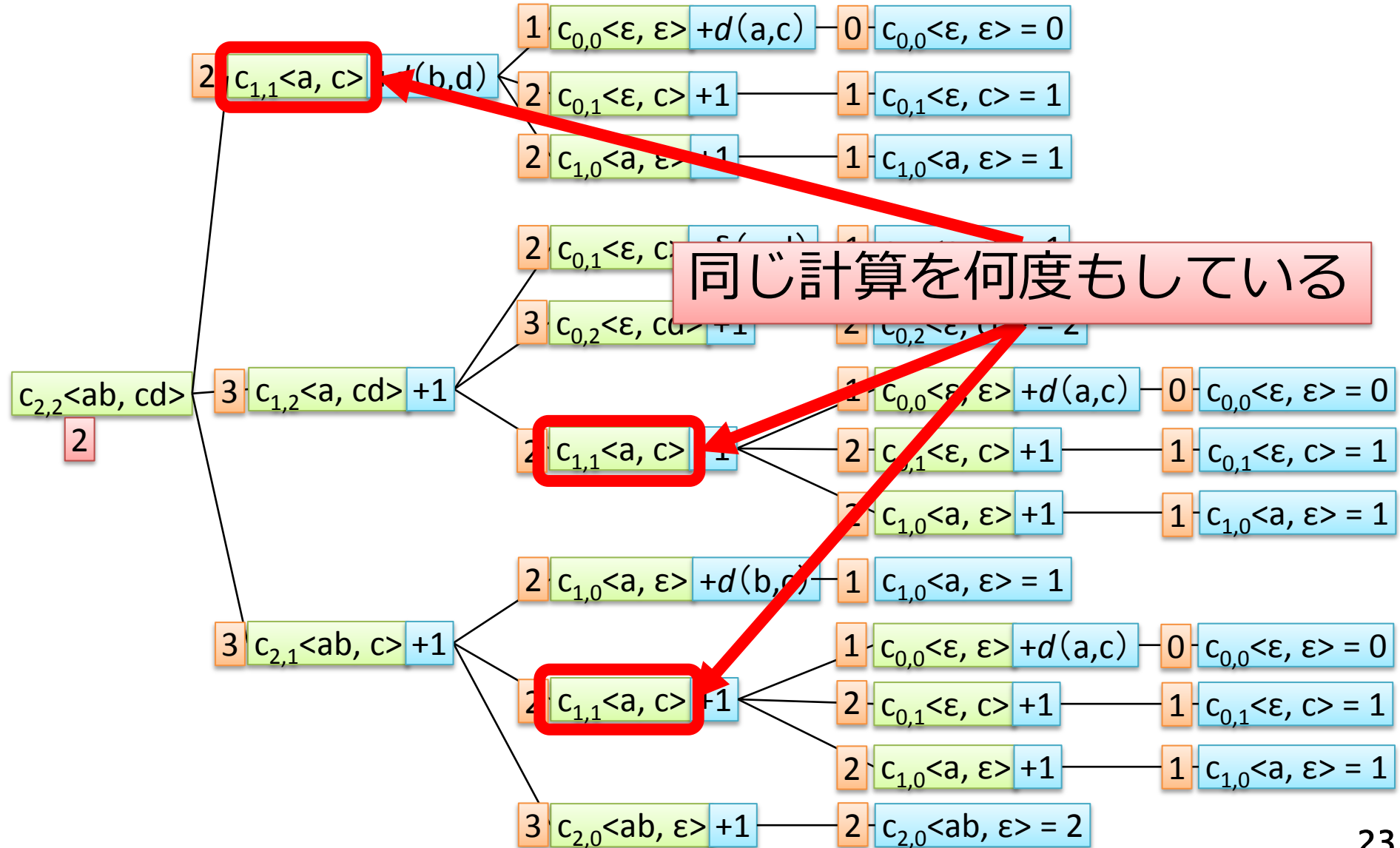
再帰的に計算

編集距離を
与える
編集過程

a b
置 置
換 換
↓ ↓
c d



再帰的計算の問題点



再帰的計算の問題点

- 同じ計算を何度も行うため計算量大

→ メモ化 (memoization) による対策

一度計算したものをメモとして記憶, 次に計算するときはそのメモを呼び出す

動的計画法(Dynamic Programing, DP)

$c_{i,j} < ab, cd >$

	ϵ	c	d
ϵ	$C_{0,0}$	$C_{0,1}$	$C_{0,2}$
a	$C_{1,0}$	$C_{1,1}$	$C_{1,2}$
b	$C_{2,0}$	$C_{2,1}$	$C_{2,2}$

$$c_{i,j} < X_i, Y_j > = \min(c_{i-1,j-1} + d(x_i, y_j), \\ c_{i-1,j} + 1, \\ c_{i,j-1} + 1)$$

$$c_{i,0} < X_i, \epsilon > = i, c_{0,j} < \epsilon, Y_j > = j$$

$C_{0,0}$ 、 $C_{0,1}$ 、 $C_{1,0}$ が計算済み
ならば、 $C_{1,1}$ の値分かる

動的計画法(Dynamic Programming, DP)

$c_{i,j} < ab, cd >$

	ϵ	c	d
ϵ	$C_{0,0}$	$C_{0,1}$	$C_{0,2}$
a	$C_{1,0}$	$C_{1,1}$	$C_{1,2}$
b	$C_{2,0}$	$C_{2,1}$	$C_{2,2}$

$$c_{i,j} < X_i, Y_j > = \min(c_{i-1,j-1} + d(x_i, y_j), \\ c_{i-1,j} + 1, \\ c_{i,j-1} + 1)$$

$$c_{i,0} < X_i, \epsilon > = i, c_{0,j} < \epsilon, Y_j > = j$$

$C_{1,1}$ 、 $C_{1,2}$ 、 $C_{2,1}$ が計算済み
ならば、 $C_{2,2}$ の値分かる

動的計画法(Dynamic Programming, DP)

$c_{i,j} < ab, cd >$

	ϵ	c	d
ϵ	$c_{0,0}$	$c_{0,1}$	$c_{0,2}$
a	$c_{1,0}$	$c_{1,1}$	$c_{1,2}$
b	$c_{2,0}$	$c_{2,1}$	$c_{2,2}$

$$c_{i,j} < X_i, Y_j > = \min(c_{i-1,j-1} + d(x_i, y_j), \\ c_{i-1,j} + 1, \\ c_{i,j-1} + 1)$$

$$c_{i,0} < X_i, \epsilon > = i, c_{0,j} < \epsilon, Y_j > = j$$

ボトムアップに $c_{i,j}$ を求める

↓
編集距離 $c_{2,2}$ 得られる

同じ計算は1度のみ \Rightarrow 計算量減

動的計画法(Dynamic Programming, DP)

$c_{i,j} < ab, cd >$

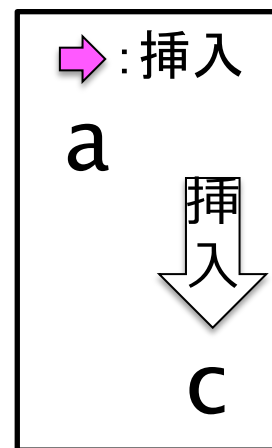
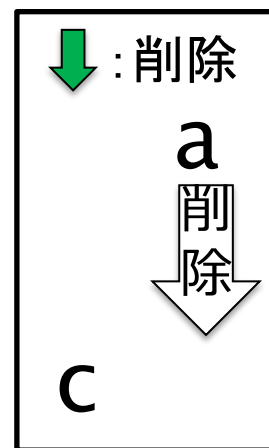
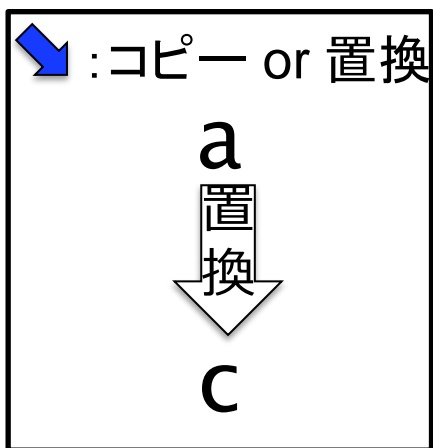
	ϵ	c	d
ϵ	$c_{0,0}$	$c_{0,1}$	$c_{0,2}$
a	$c_{1,0}$	$c_{1,1}$	$c_{1,2}$
b	$c_{2,0}$	$c_{2,1}$	$c_{2,2}$

$$c_{i,j} < X_i, Y_j > = \min(\boxed{c_{i-1,j-1} + d(x_i, y_j)},$$

$$\boxed{c_{i-1,j} + 1},$$

$$\boxed{c_{i,j-1} + 1})$$

$$c_{i,0} < X_i, \epsilon > = i, c_{0,j} < \epsilon, Y_j > = j$$



動的計画法による計算

$c_{i,j} \langle ab, cd \rangle$

	ϵ	c	d
ϵ			
a			
b			

$$c_{i,j} \langle X_i, Y_j \rangle = \min(c_{i-1,j-1} + \delta(x_i, y_j), \\ c_{i-1,j} + 1, \\ c_{i,j-1} + 1)$$

$$c_{i,0} \langle X_i, \epsilon \rangle = i, c_{0,j} \langle \epsilon, Y_j \rangle = j$$

動的計画法による計算

$c_{i,j} < ab, cd >$

	ϵ	c	d
ϵ	0 \rightarrow 1 \rightarrow 2		
a			
b			

$$c_{i,j} < X_i, Y_j > = \min(c_{i-1,j-1} + \delta(x_i, y_j), c_{i-1,j} + 1, c_{i,j-1} + 1)$$

$$c_{i,0} < X_i, \epsilon > = i, c_{0,j} < \epsilon, Y_j > = j$$

$$c_{0,0} < \epsilon, \epsilon > = 0, c_{0,1} < \epsilon, c > = 1, c_{0,2} < \epsilon, cd > = 2$$

挿入

動的計画法による計算

$c_{i,j} < ab, cd >$

	ϵ	c	d
ϵ	0	1	2
a	1		
b	2		

$$c_{i,j} < X_i, Y_j > = \min(c_{i-1,j-1} + \delta(x_i, y_j), c_{i-1,j} + 1, c_{i,j-1} + 1)$$

$$c_{i,0} < X_i, \epsilon > = i, c_{0,j} < \epsilon, Y_j > = j$$

$$c_{0,0} < \epsilon, \epsilon > = 0, c_{1,0} < a, \epsilon > = 1, c_{2,0} < ab, \epsilon > = 2$$

削除

動的計画法による計算

$c_{i,j} < ab, cd >$

		$0 + d(a, c) = 0 + 1 = 1$	c	d
ϵ	0	1	2	
a	1	1		
b	2			

Diagram illustrating the dynamic programming table for sequence alignment between "ab" and "cd". The table shows the cost $c_{i,j}$ for aligning the first i characters of "ab" with the first j characters of "cd". The cells are labeled with their values, and arrows indicate the optimal alignment path from $(0,0)$ to $(1,1)$. The path is: $(0,0) \rightarrow (1,1)$ (blue arrow), $(1,1) \rightarrow (1,2)$ (green arrow), and $(1,2) \rightarrow (2,2)$ (pink arrow). The cost values are: $c_{0,0}=0$, $c_{0,1}=1$, $c_{0,2}=2$, $c_{1,0}=1$, $c_{1,1}=1$, $c_{1,2}=2$. The optimal alignment is "a" aligned with "c" and "b" aligned with "d".

$$c_{i,j} < X_i, Y_j > = \min(c_{i-1,j-1} + d(x_i, y_j), c_{i-1,j} + 1, c_{i,j-1} + 1)$$

$$c_{i,0} < X_i, \epsilon > = i, c_{0,j} < \epsilon, Y_j > = j$$

$$c_{1,1} < a, c > = \min(c_{0,0} < \epsilon, \epsilon > + d(a, c), \text{ 置換} \\ c_{0,1} < \epsilon, c > + 1, \text{ 削除} \\ c_{1,0} < a, \epsilon > + 1) \text{ 挿入}$$

動的計画法による計算

$c_{i,j} < ab, cd >$

	ϵ	$1 + d(a, d) = 1 + 1 = 2$	d
ϵ	0	1	2
a	1	1	2
b	2		

Diagram illustrating the dynamic programming table for the edit distance between "ab" and "cd". The table shows the cost of operations (insertion, deletion, substitution) for each character pair. The values are calculated as follows:

- $c_{0,0} = 0$ (Base case)
- $c_{0,1} = 1$ (Insertion of 'a')
- $c_{0,2} = 2$ (Insertion of 'b')
- $c_{1,1} = 1$ (Substitution of 'a' for 'a')
- $c_{1,2} = 2$ (Substitution of 'a' for 'd')
- $c_{2,2} = 2$ (Substitution of 'b' for 'd')

Arrows indicate the transitions used to calculate the values:

- From $c_{0,1}$ to $c_{1,1}$ (Insertion of 'a')
- From $c_{0,2}$ to $c_{1,2}$ (Insertion of 'b')
- From $c_{1,1}$ to $c_{1,2}$ (Substitution of 'a' for 'd')
- From $c_{0,1}$ to $c_{1,2}$ (Deletion of 'a')
- From $c_{0,2}$ to $c_{1,2}$ (Deletion of 'b')

$$c_{i,j} < X_i, Y_j > = \min(c_{i-1,j-1} + d(x_i, y_j), c_{i-1,j} + 1, c_{i,j-1} + 1)$$

$$< X_i, \epsilon > = i, c_{0,j} < \epsilon, Y_j > = j$$

$$c_{1,2} < a, cd > = \min(c_{0,1} < \epsilon, c > + d(a, d), \text{置換 } c_{0,2} < \epsilon, cd > + 1, \text{削除 } c_{1,1} < a, c > + 1) \text{ 挿入}$$

動的計画法による計算

$c_{i,j} \langle ab, cd \rangle$

	ϵ	c	d
ϵ		1	2
a	1	1	2
b	2	2	

$1+d(b,c)$
 $= 1+1=2$

$1+1=2$

$2+1=3$

$$c_{i,j} \langle X_i, Y_j \rangle = \min(c_{i-1,j-1} + d(x_i, y_j), c_{i-1,j} + 1, c_{i,j-1} + 1)$$

$$c_{i,0} \langle X_i, \epsilon \rangle = i, c_{0,j} \langle \epsilon, Y_j \rangle = j$$

$$c_{2,1} \langle ab, c \rangle = \min(c_{1,0} \langle a, \epsilon \rangle + d(b, c), \text{置換 } c_{1,1} \langle a, c \rangle + 1, \text{削除 } c_{2,0} \langle ab, \epsilon \rangle + 1) \text{挿入}$$

動的計画法による計算

$c_{i,j} \langle ab, cd \rangle$

	ϵ	c	d
ϵ	0		2
a	1	1	2
b	2	2	2

$1+d(b,d)$
 $= 1+1=2$

$2+1=3$

$2+1=3$

$$c_{i,j} \langle X_i, Y_j \rangle = \min(c_{i-1,j-1} + d(x_i, y_j), c_{i-1,j} + 1, c_{i,j-1} + 1)$$

$$c_{i,0} \langle X_i, \epsilon \rangle = i, c_{0,j} \langle \epsilon, Y_j \rangle = j$$

$$c_{2,2} \langle ab, cd \rangle = \min(c_{1,1} \langle a, c \rangle + d(b, d), \text{置換 } c_{1,2} \langle a, cd \rangle + 1, \text{削除 } c_{2,1} \langle ab, c \rangle + 1) \text{挿入}$$

動的計画法による計算

$c_{i,j} < ab, cd >$

	ϵ	c	d
ϵ	0	1	2
a	1	1	2
b	2	2	2

$$c_{i,j} < X_i, Y_j > = \min(c_{i-1,j-1} + d(x_i, y_j), \\ c_{i-1,j} + 1, \\ c_{i,j-1} + 1)$$

$$c_{i,0} < X_i, \epsilon > = i, c_{0,j} < \epsilon, Y_j > = j$$



ここが編集距離

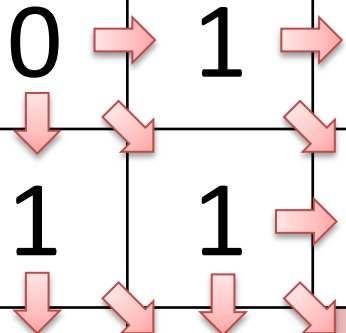
問題3, 問題4

動的計画法による計算

$c_{i,j} < ab, cd >$

編集操作の求め方

	ϵ	c	d
ϵ	0	1	2
a	1	1	2
b	2	2	2

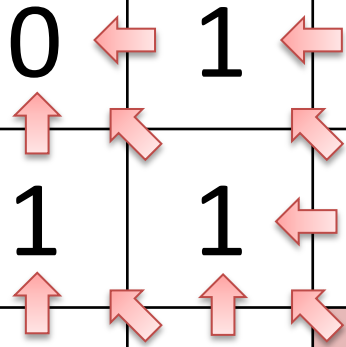


値を埋める際に使ったデータの
流れの矢印を表示

動的計画法による計算

$c_{i,j} < ab, cd >$

	ϵ	c	d
ϵ	0	1	2
a	1	1	2
b	2	2	2



編集操作の求め方

矢印を180°回転

動的計画法による計算

$c_{i,j} \langle ab, cd \rangle$

	ϵ	c	d
ϵ	0	1	2
a	1	1	2
b	2	2	2

編集操作の求め方

右下から0の部分までたどって
いったときの矢印が編集操作
(ただし, 編集操作が1つに決
まるとは限らない)

a b
↓置換 ↓置換
c d
R R

問題6

課題3：全部で7問

- 問題1 (5点) [記述問題]再帰による編集距離の計算 (手書き可)
- 問題2 (10点) 再帰関数による実装
- 問題3 (10点) [記述問題]動的計画法による編集距離の計算 (手書き可)

2週目の面接のみ受付

- 問題4 (10点) 動的計画法による実装
- 問題5 (15点) [プログラム+記述問題]再帰呼び出しにおけるメモ化
- 問題6 (10点) 編集操作の表示
- 問題7 (10点) スpellチェッカーの作成

4週目の面接のみ受付

問題1

$X=\text{"on"}$, $Y=\text{"so"}$ に対して, $c_{i,j}$ の漸化式を手作業で再帰的に適用することにより $c_{2,2}$ の値を求めよ.

◎ $c_{2,2}$ を求める途中経過を省略することなく, すべてレポートに記すこと(課題説明スライドと同様な木構造を使って記すこと).

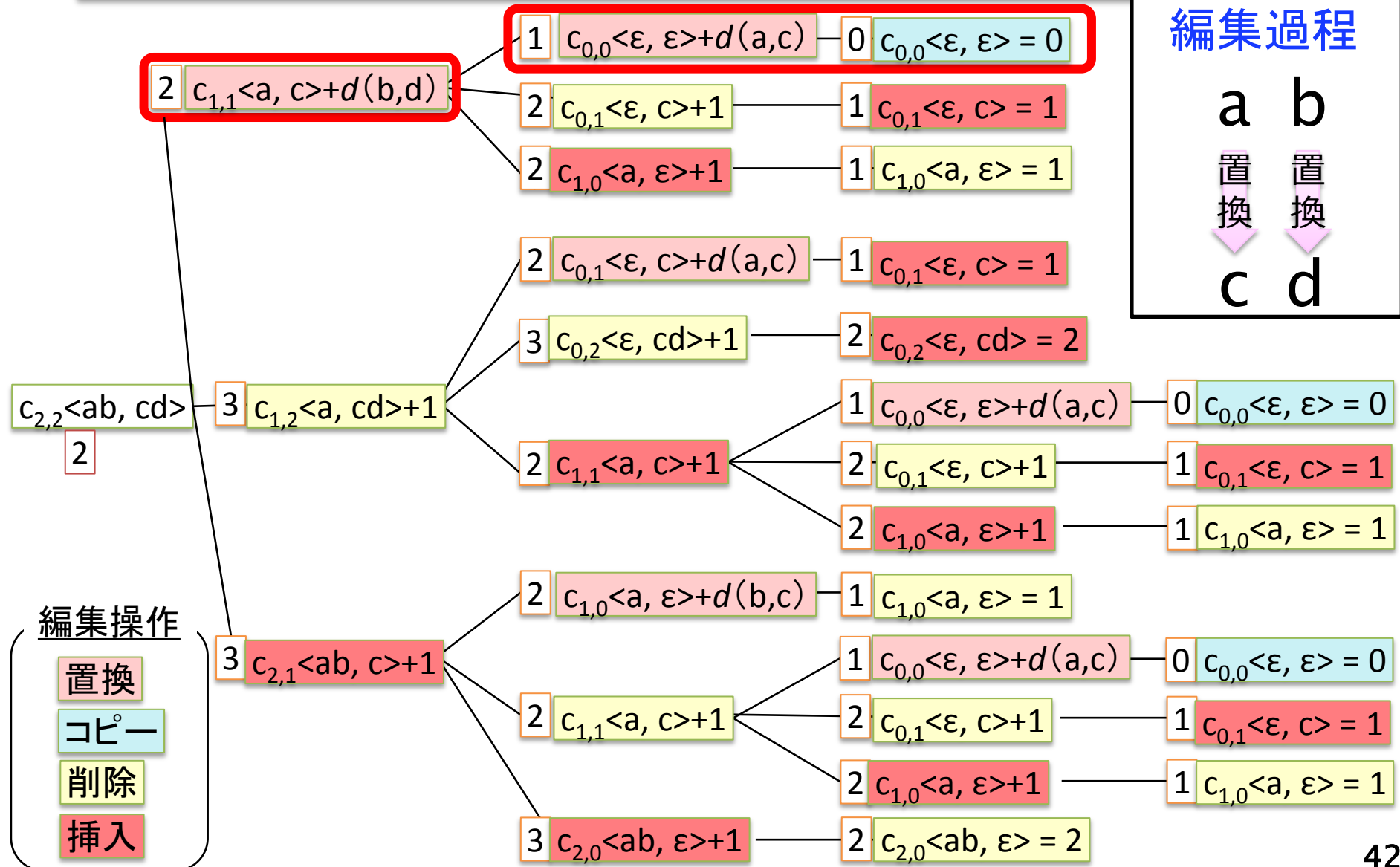
◎ **注意**: 以下の**3つすべて**を記せ

- $c_{i,j}$ の値を求めるのに使った接頭辞の編集距離の値とその添字.
- 編集操作が分かるように明示せよ.
- 編集距離を与える編集過程をすべて示せ.

問題1

編集距離を
与える
編集過程

a b
置 置
換 換
↓ ↓
c d



面接を受けるときの注意事項

- **C言語の用語や関数の動作をあらかじめ調べてくること**
プログラムのソースコードに説明を書いてかまわないので、面接時に説明できるようにすること
- **課題に関する問題文をきちんと読むこと**
問題文を読むことが解答への近道になることもあります
- **2週目の面接時に、途中だったとしても、取り組んでいる問題を持ってきましょう**
困っていることやわからないことがあれば、面接の時にアドバイスをするので、途中のプログラムでも遠慮なく持ってきましょう