

Inferring a Graph from Path Frequency



Tatsuya Akutsu^{1,2} & Daiji Fukagawa²

¹ Institute for Chemical Research, Kyoto Univ., Japan

² Graduate School of Informatics, Kyoto Univ., Japan

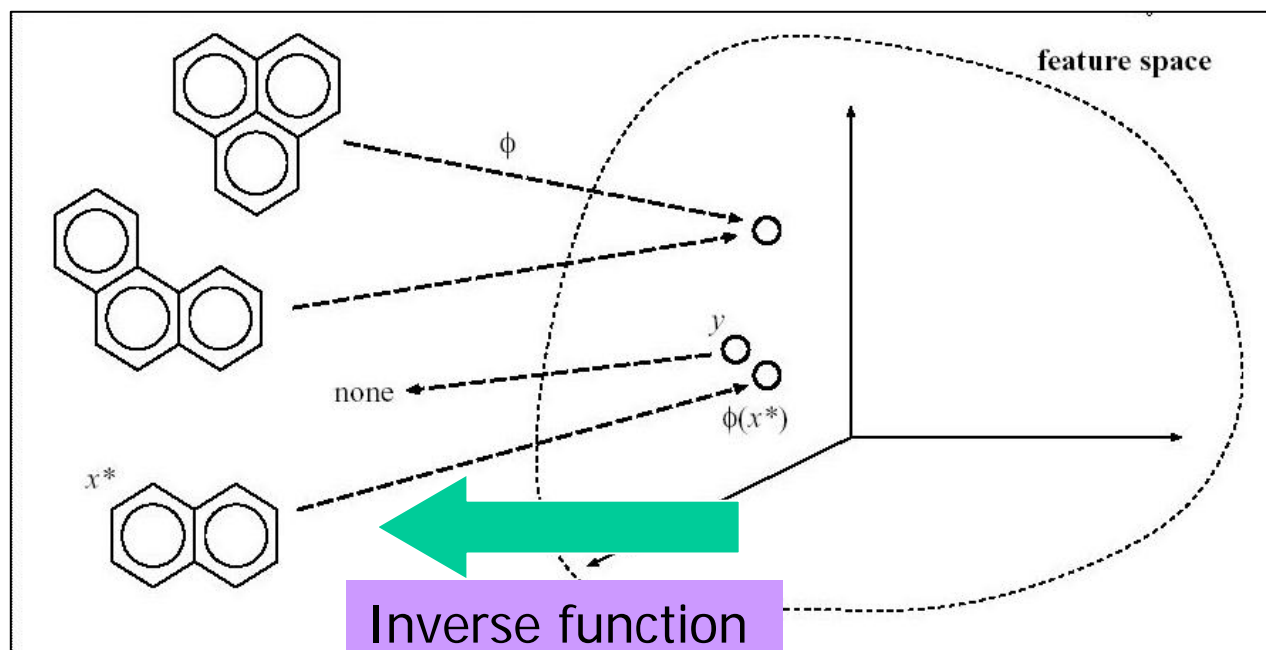
Outline

- Introduction
 - String inference from spectrum feature (SISF)
 - Graph inference from path frequency (GIPF)
 - Optimization versions (SISF-M, GIPF-M)
- Algorithms for special cases
- Complexity results
 - Strong NP-completeness of GIPF in general case
(Reduction from 3-PARTITION)
- Conclusion

Motivation

- Kernel methods (e.g. Support Vector Machine) have been applied to various problems. In kernel methods,
 Data (sequences, chemical compounds,...) → **Feature vector**
- This work: we consider reverse direction, i.e.,
 Feature vector → **Data**

May be useful for designing new sequences/chemical compounds

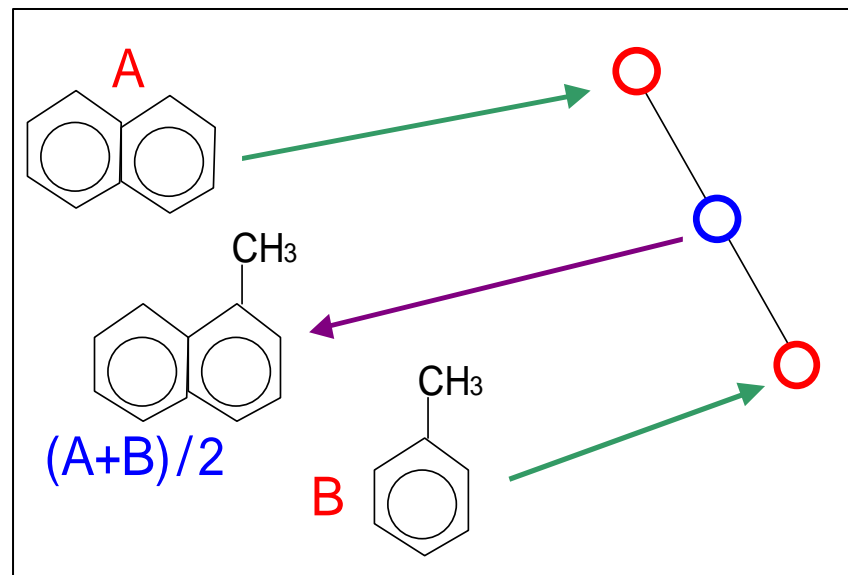


Motivation (continued)

Potential application:

drug design

- For example, design a **new compound** which is the **middle** of known compounds A and B



Related work

- Kernel PCA + regression [Bakir, Weston, Scölkopf 2004]
- Graph pre-image [Bakir, Zien, Tsuda 2004]
- But, no complexity studies

graph inference problem

Graph inference from path frequency

Given **path frequency vector** v , infer the **original graph** whose feature vector (=path frequency) is equal to v (or closest to v).

(length) (path: occurrence)

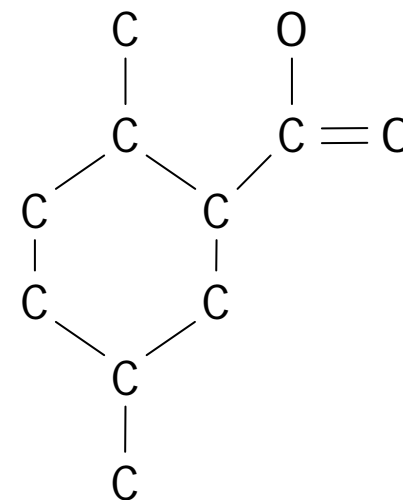
0 --- C x 9, O x 2

1 --- CC x 18, CO x 2, OC x 2

2 --- CCC x 24, CCO x 2, OCC x 2, OCO x 2

3 --- CCCC x 26, CCCCO x 4, OCCCC x 4

: : : : :



Spectrum Feature for Strings [Leslie *et al.* 02]

For a string S ,

- Spectrum feature of level k is a frequency vector of all possible k -grams.

e.g. spectrum feature of level 2 for 'aababb':

<u>aa</u> x 1
<u>ab</u> x 2
<u>ba</u> x 1
<u>bb</u> x 1

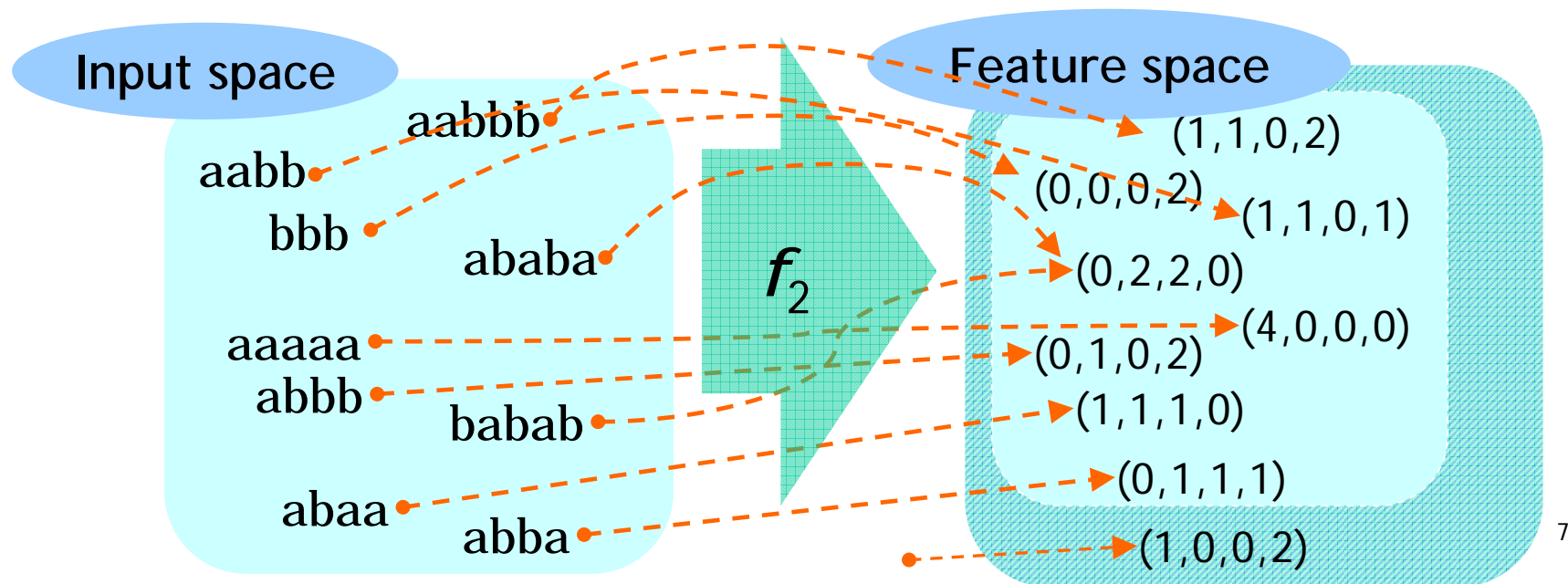
$$\rightarrow f_2(\text{'aababb'}) = (1, 2, 1, 1)$$

Spectrum Feature for Strings [Leslie *et al.* 02]

f_K : **mapping** from input space to feature space

$$f_K(s) = \left(occ(t, s) \right)_{t \in \Sigma^K} \quad \begin{array}{l} K: \text{level } (>0) \\ \Sigma: \text{alphabet} \end{array}$$

where $occ(t, s)$ is # of occurrences of a substring t in a string s .



Problem 1

SISF: *String Inference from Spectrum Feature*

Input: an integer K , feature vector $\mathbf{v} = (v_t)_t \quad K$

Output: a string s which, if it exists, satisfies $f_K(s) = \mathbf{v}$,
otherwise "no solution."

Ex) $\Sigma = \{a, b\}$, $K=2$, $\mathbf{v} = (v_{aa}, v_{ab}, v_{ba}, v_{bb}) = (1, 1, 0, 2)$

$$|s| = (v_{aa} + v_{ab} + v_{ba} + v_{bb}) + (K-1) = (1+1+0+2) + 1 = 5$$

$$f_2('aaaaa') = (4, 0, 0, 0)$$

$$f_2('aaaab') = (3, 1, 0, 0)$$

: : : :

$$f_2('aabbb') = (1, 1, 0, 2)$$

: : : :

$$f_2('bbbbbb') = (0, 0, 0, 4)$$

solution: $s = 'aabbb'$

Solutions may not be unique:

$\mathbf{v} = (1, 1, 1, 1)$ 4 solutions

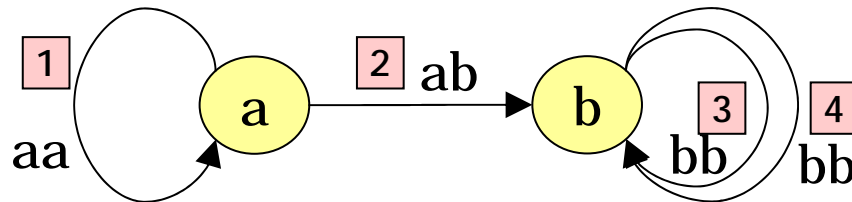
$\mathbf{v} = (1, 2, 2, 1)$ 12 solutions

Linear time algorithm for SISF

Reduction to **Eulerian graph problem** [Pevzner]

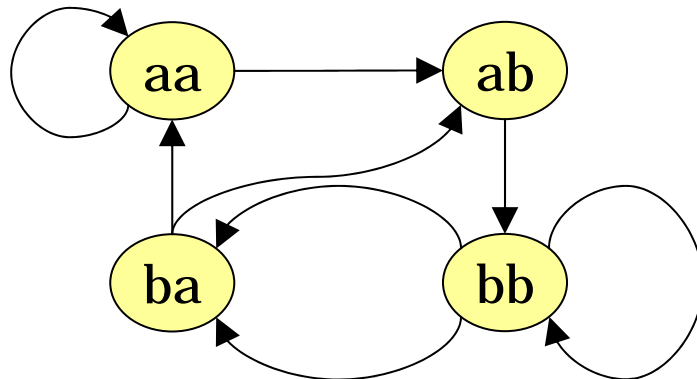
$f_k(s) = v$ for some s G_v has a Eulerian path

Ex.1) $K=2, v=(v_{aa}, v_{ab}, v_{ba}, v_{bb})=(1, 1, 0, 2)$



solution
 \xrightarrow{aabb}
 1 aa
 2 ab
 3 bb
 4 bb

Ex.2) $K=3, v=(v_{aaa}, v_{aab}, v_{aba}, v_{abb}, v_{baa}, v_{bab}, v_{bba}, v_{bbb})=(1, 1, 0, 1, 1, 1, 2, 1)$



2 solutions
 $\left\{ \begin{array}{l} bbaabbbab \\ bbaaabbab \end{array} \right.$

Problem 2

SISF-M: SISF with the Minimum Error

Input: an integer K , feature vector $v = (v_t)_t \quad K$

Output: a string s which minimizes the distance between $f_K(s)$ and v

例) $K=2, v=(v_{aa}, v_{ab}, v_{ba}, v_{bb})=(1, 2, 0, 1)$

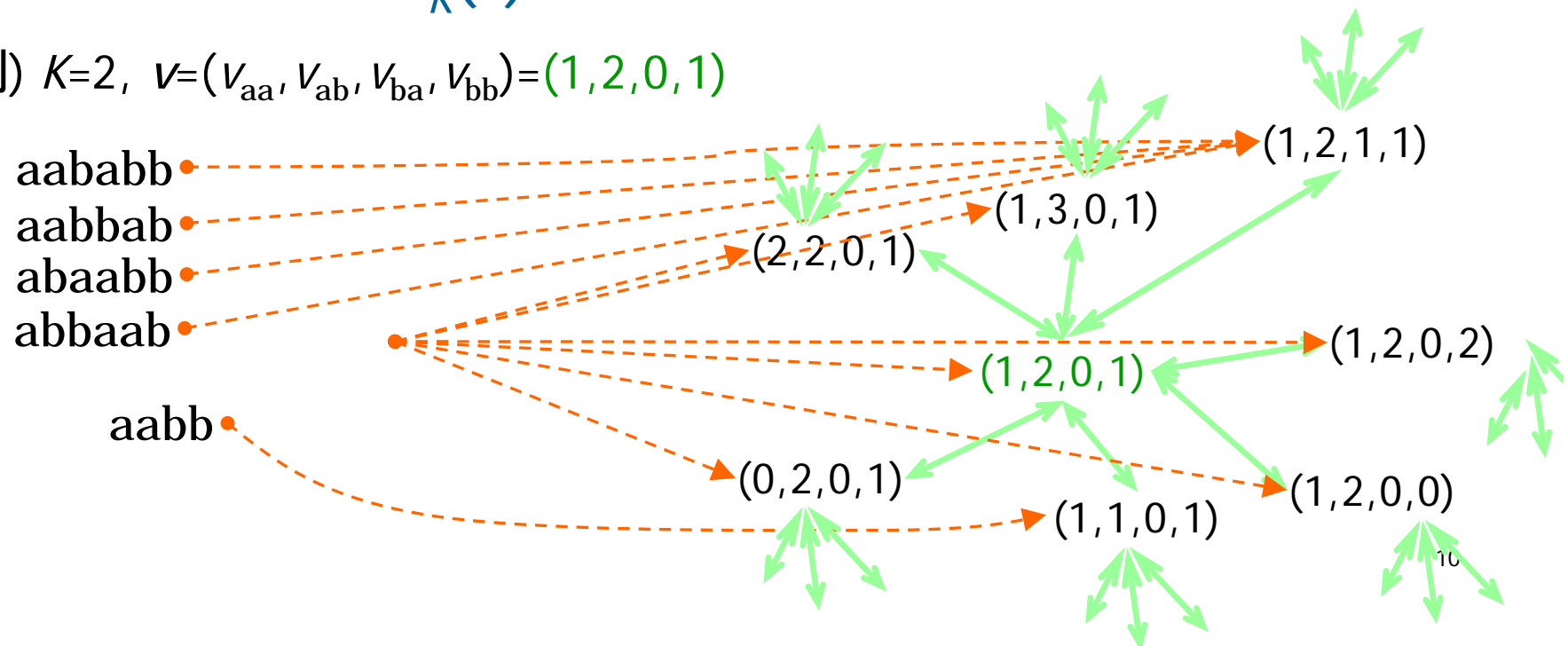
aababb

aabbab

abaabb

abbaab

aabb



Algorithm for SISF-M

- It seems difficult to apply Eulerian path technique.
- Thus, we employ another approach based on Dynamic programming.
- The algorithm is a special case of the graph inference algorithm.

Feature vector for graphs: Path Frequency (c.f. Marginalized Graph Kernel)

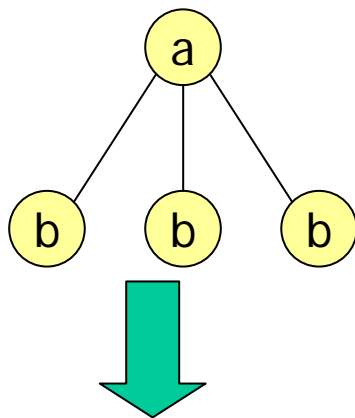
$$\text{Feature vector } f_K: G \rightarrow \mathbb{N}^K$$

$$f_K(G) = \left(\text{occ}(t, G) \right)_{t \in \Sigma^{\leq K}}$$

K : level (>0)
 Σ : alphabet
 G : graph

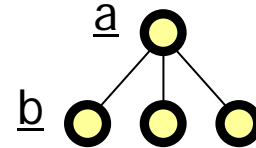
where $\text{occ}(t, G)$ is # of occurrences of paths labeled with t in a graph G

e.g.) $K=1$



$$f_1(G) = (1, 3, 0, 3, 3, 0)$$

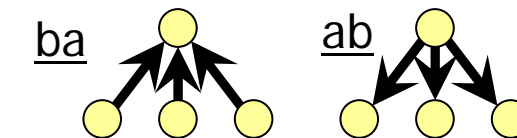
$k=0$



$$\text{occ}(a, G) = 1$$

$$\text{occ}(b, G) = 3$$

$k=1$



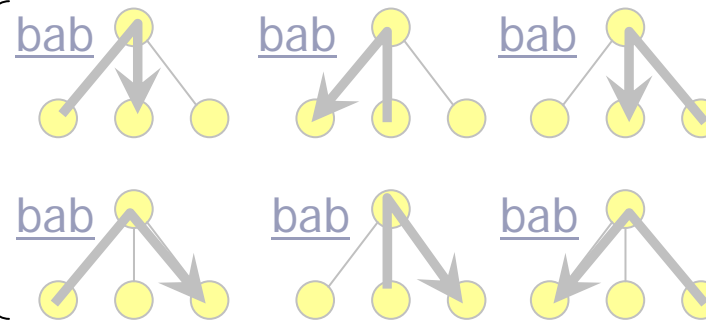
$$\text{occ}(aa, G) = 0$$

$$\text{occ}(ab, G) = 3$$

$$\text{occ}(ba, G) = 3$$

$$\text{occ}(bb, G) = 0$$

$k=2$



:

Problem 3 & 4

GIPF: *Graph Inference from Path Frequency*

Input: an integer K , feature vector $v = (v_t)_t \quad K$

Output: a graph G which satisfies, if it exists, $f_K(G) = v$
otherwise "no solution."

GIPF-M: *GIPF with Minimum Error*

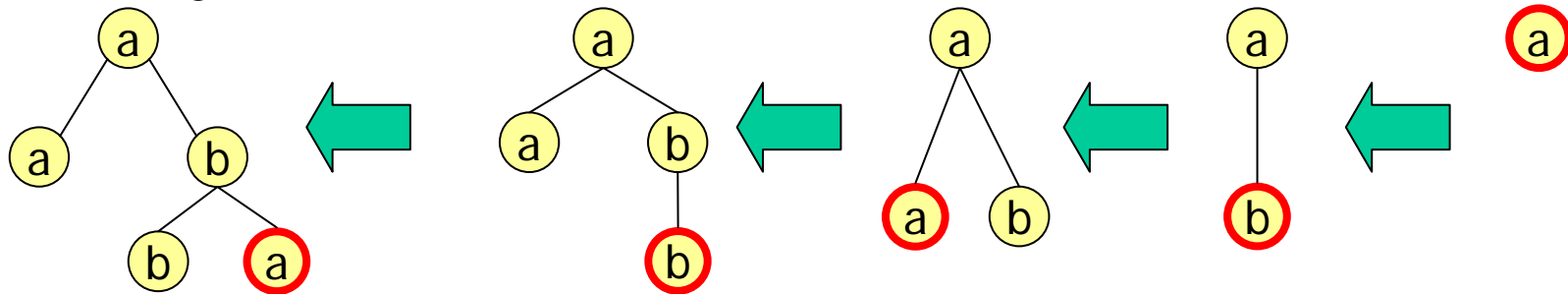
Input: an integer K , feature vector $v = (v_t)_t \quad K$

Output: a graph G , which minimizes the distance between
 $f_K(G)$ and v

Dynamic Programming for restricted GIPF

(1) Trees, $K=1$, fixed

Any tree can be constructed by inserting a leaf one by one.



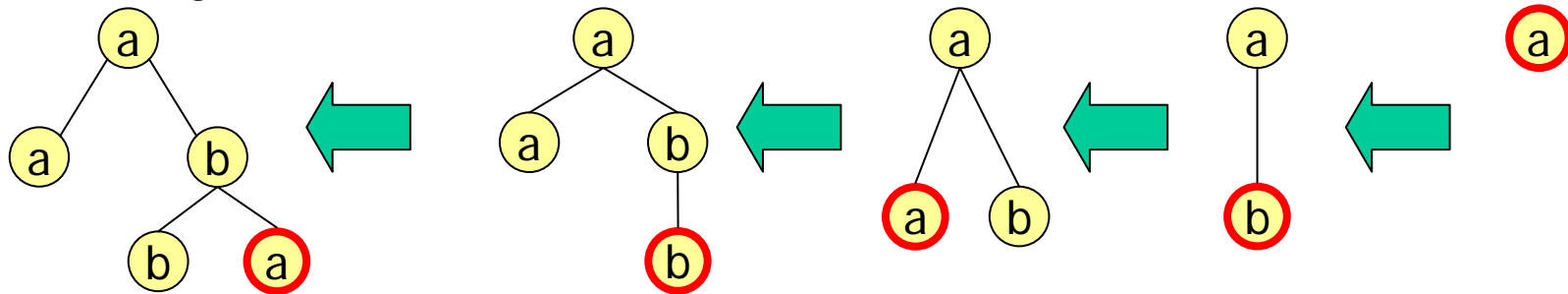
$D(v)=1$ iff. There exists a tree T s.t. $f_v(T)=v$

$$\begin{aligned}
 D(n_a, n_b, n_{aa}, n_{ab}, n_{ba}, n_{bb}) = 1 &\Leftrightarrow \\
 &(D(n_a - 1, n_b, n_{aa} - 2, n_{ab}, n_{ba}, n_{bb}) = 1) \vee \\
 &(D(n_a - 1, n_b, n_{aa}, n_{ab} - 1, n_{ba} - 1, n_{bb}) = 1) \vee \\
 &(D(n_a, n_b - 1, n_{aa}, n_{ab} - 1, n_{ba} - 1, n_{bb}) = 1) \vee \\
 &(D(n_a, n_b - 1, n_{aa}, n_{ab}, n_{ba}, n_{bb} - 2) = 1)
 \end{aligned}$$

Dynamic Programming for restricted GIPF

(1) Trees, $K=1$, fixed (cont'd.)

Any tree can be constructed by inserting a leaf one by one.



$D(v)=1$ iff. There exists a tree T s.t. $f_v(T)=v$

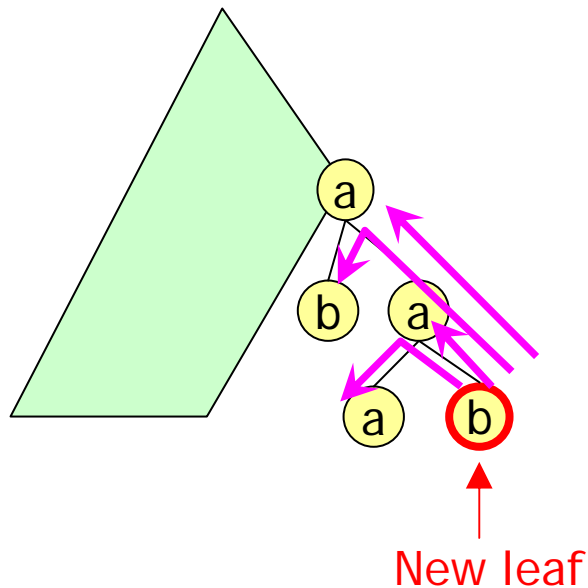
(Theorem) GIPF for trees is solved in polynomial time in n (the size of tree) if $K=1$ and a fixed alphabet.

→ GIPF-M is also solved by searching in this table.

Dynamic Programming for restricted GIPF

(2) Trees, fixed K , ,

- Extension of DP for $K=1$ (not straightforward)
- More complicated data structure than $K=1$.

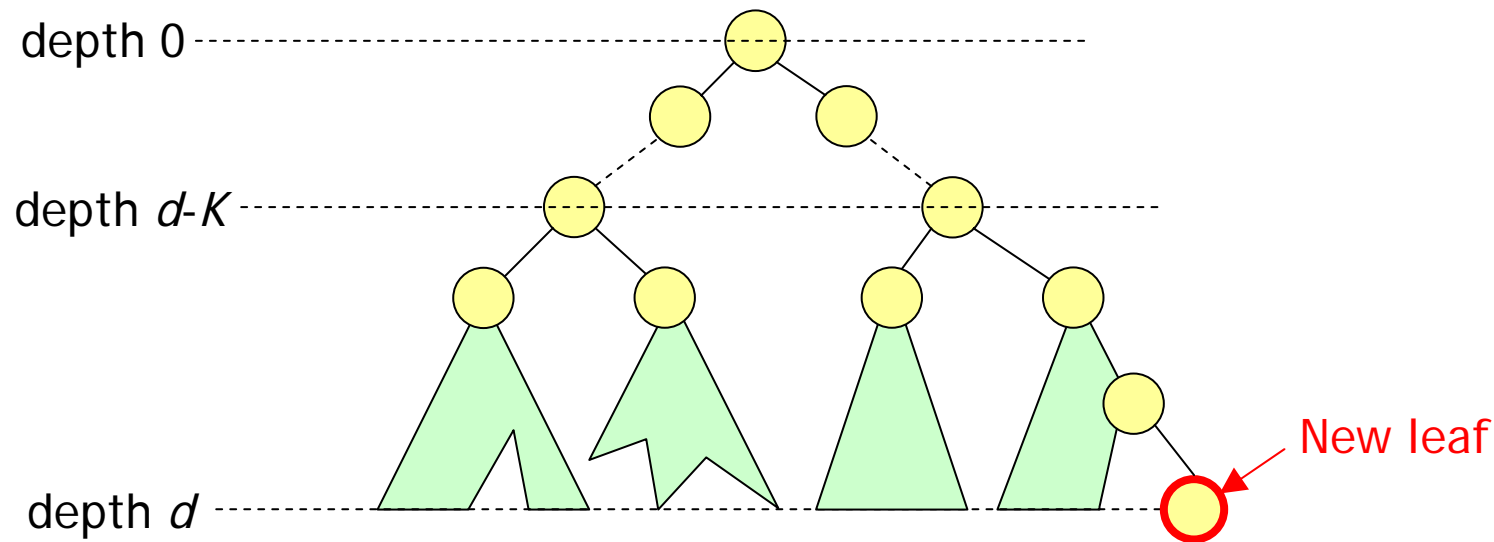


When a new leaf is added, much more new paths appear
 $\rightarrow O(k)$ new paths

Dynamic Programming for restricted GIPF

(2) Trees, fixed K , , (cont'd)

- Extended DP table: $D(v, e, d)$
 - v : feature vector, e : paths around leaves, d : depth
- K , , are fixed, so is the size of e .



(Theorem) GIPF (and GIPF-M) for trees is solved in polynomial time in n if K , , are all fixed.

Strong NP-completeness of GIPF

- GIPF is strongly NP-complete even if the underlying graph is a tree and $K=3$

(We improved the result from that in the proceedings, where this result was shown for non-tree graphs)

- Reduction from 3-PARTITION problem

Reduction from: 3-PARTITION, $P=(X, w, B)$

- $X=\{x_1, \dots, x_{3m}\}$, $w(x_i)=w_i$, $w_j=Bm$

Reduction into: GIPF, $Q=(v, K)$

- $v=\{a, b, c, d, y\}$ $\{x_1, \dots, x_{3m}\}$ $\{A_1, \dots, A_m\}$
- $K=3$
- $v(s)=O(\text{poly}(m+B))$ for every $s \in v$
- $|\{s \in v \mid v(s) \text{ is non-zero}\}| = O(\text{poly}(m+B))$

Strong NP-completeness of GIPF

(Theorem) GIPF is strongly NP-complete, even if $K=3$ and the underlying graph is a tree.

(Proof)

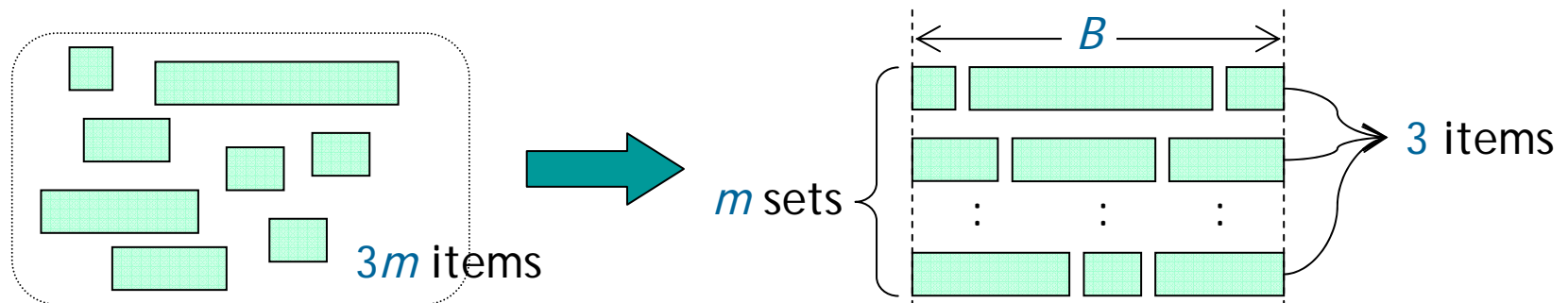
- Reduction from 3-PARTITION can be done in $\text{poly}(m+B)$ time and thus its size is bounded by $\text{poly}(m+B)$.
- GIPF Q is 'yes' 3-PARTITION P is 'yes'

3-PARTITION problem [Garey&Johnson]

3-PARTITION: strongly NP-complete problem

Input: a set $X = \{x_1, \dots, x_{3m}\}$ of $3m$ items, for each x_i a weight $w(x_i)$, s.t. $\sum_i w(x_i) = mB$

Output: 'yes' if there exist a partition A_1, \dots, A_m of X s.t. $|A_h| = 3$ and $\sum_{x \in A_i} w(x) = B$, 'no' otherwise

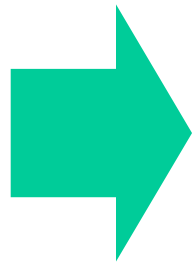
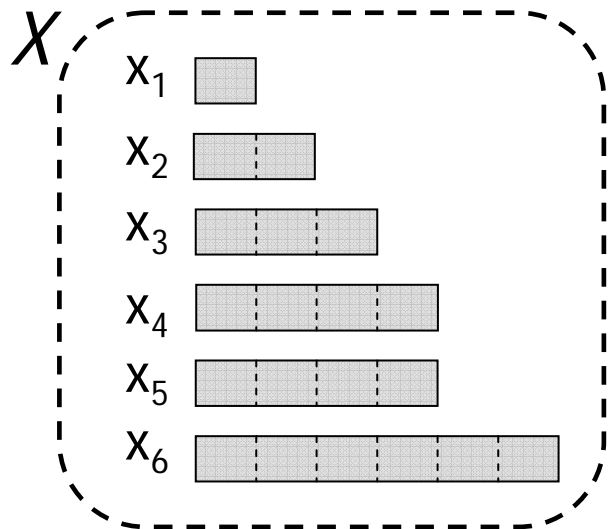


- 3-PARTITION does not have pseudo-polynomial time algorithm unless $P=NP$.
(i.e., cannot be solved in $\text{poly}(m+B)$ unless $P=NP$)
- Strongly NP-complete even if $B/4 < w(x) < B/2$

An Example of Reduction (1)

An instance of 3-PARTITION P :

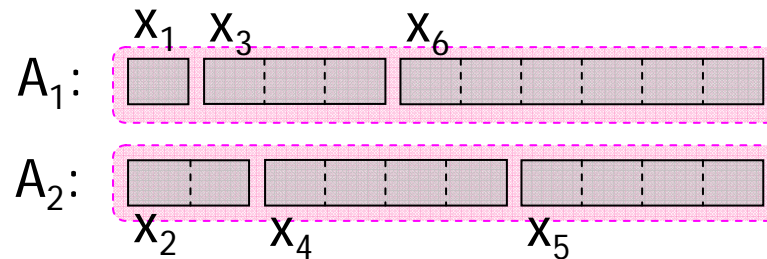
- $m=2 \rightarrow X = \{x_1, x_2, x_3, x_4, x_5, x_6\}$ ($|X|=3m=6$)
- $w=(1, 2, 3, 4, 4, 6)$, $B=10$
- $w(x_1)=1$, $w(x_2)=2$, $w(x_3)=3$, $w(x_4)=4$, $w(x_5)=4$, $w(x_6)=6$



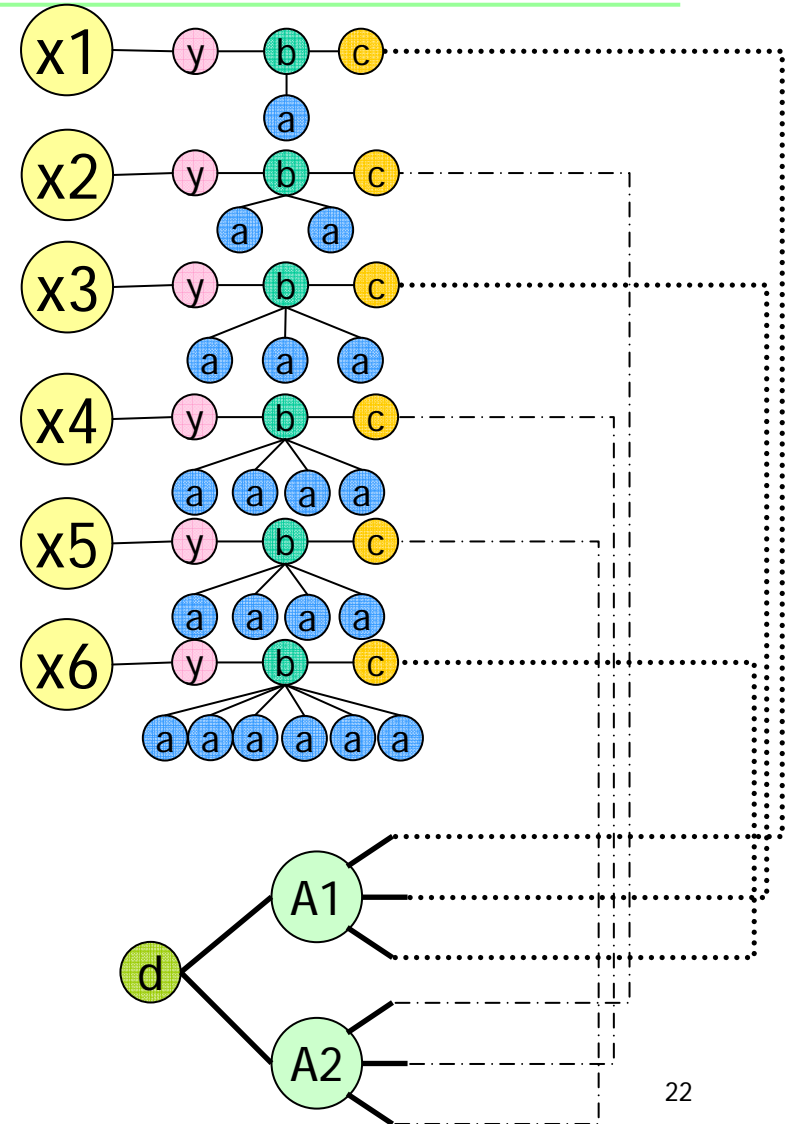
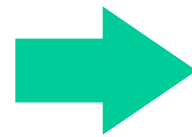
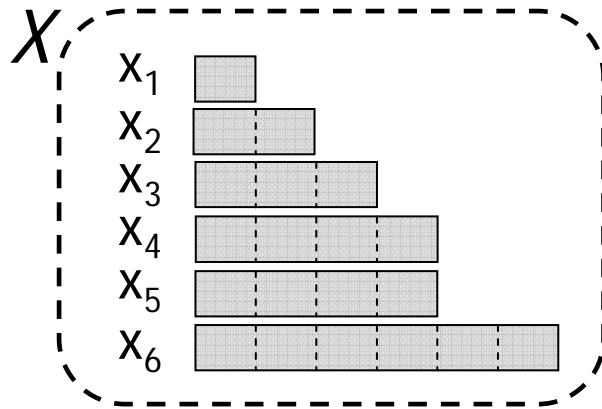
The solution for P :

$$A_1 = \{x_1, x_3, x_6\} \rightarrow 1+3+6=10$$

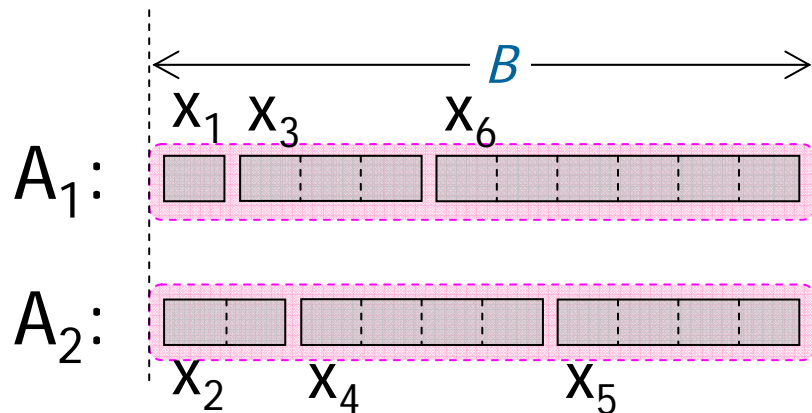
$$A_2 = \{x_2, x_4, x_5\} \rightarrow 2+4+4=10$$



An Example of Reduction (2)



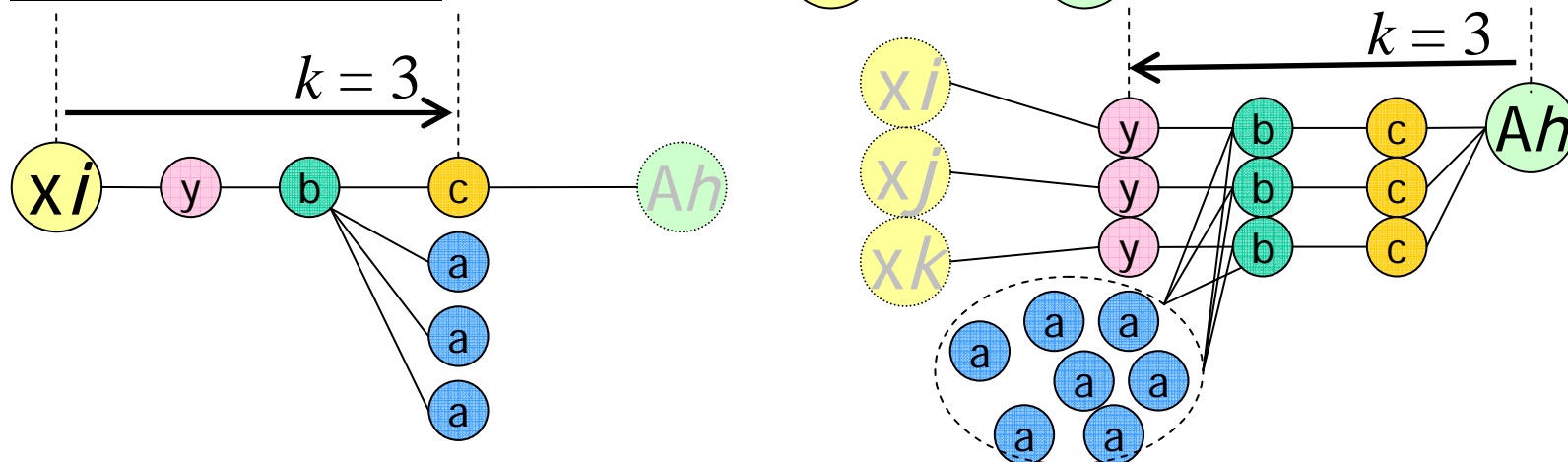
Solution for 3-PARTITION



An Example of Reduction (3)

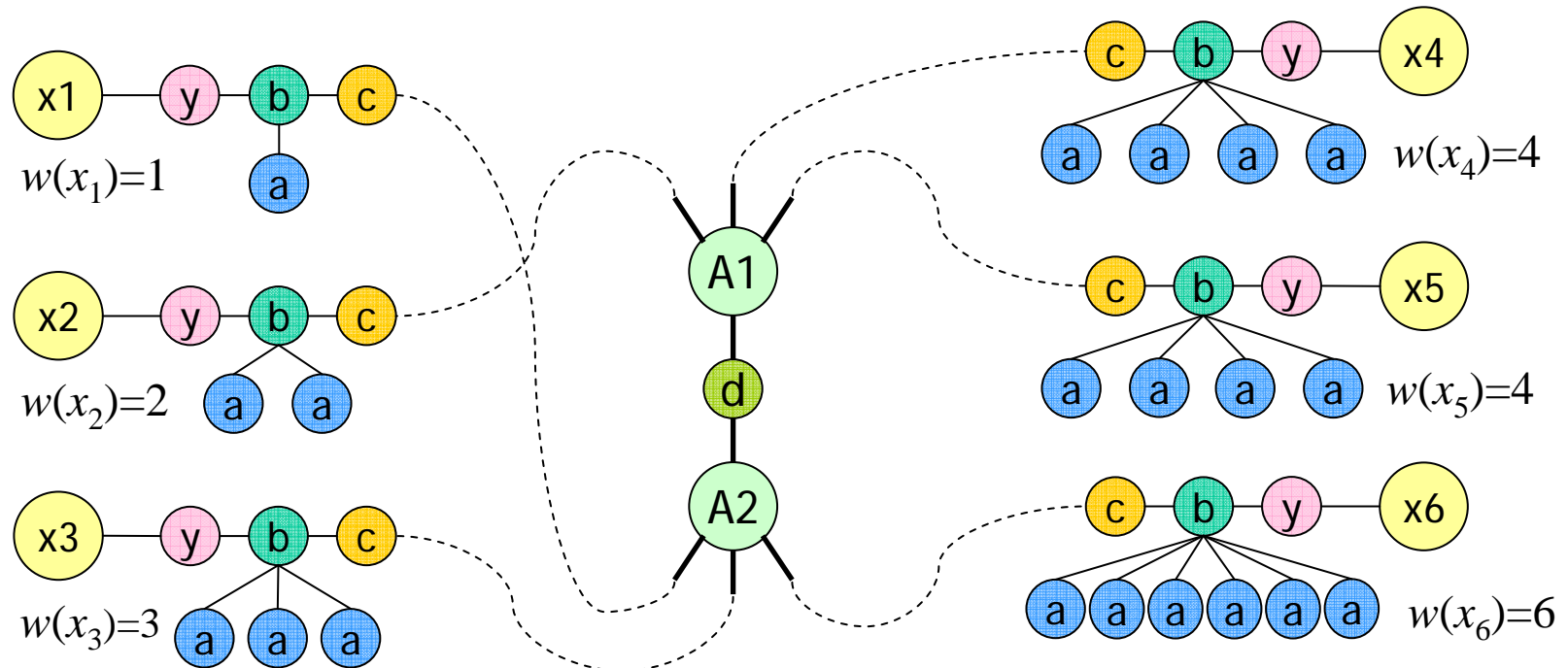
There are **two kinds** of vertices which have **unique label**. $\rightarrow x_i$'s (x_i) and A_h 's (A_h)

- x_i encodes the weight of i -th item $\rightarrow w(x_i)$
- A_h is a **matchmaker** of x_i 's, but doesn't know who matches who, because x_i and A_h are distant.



An Example of Reduction (4)

3-PARTITION P : $w = (1, 2, 3, 4, 4, 6) \rightarrow$ GIPF Q : $(v, 3)$

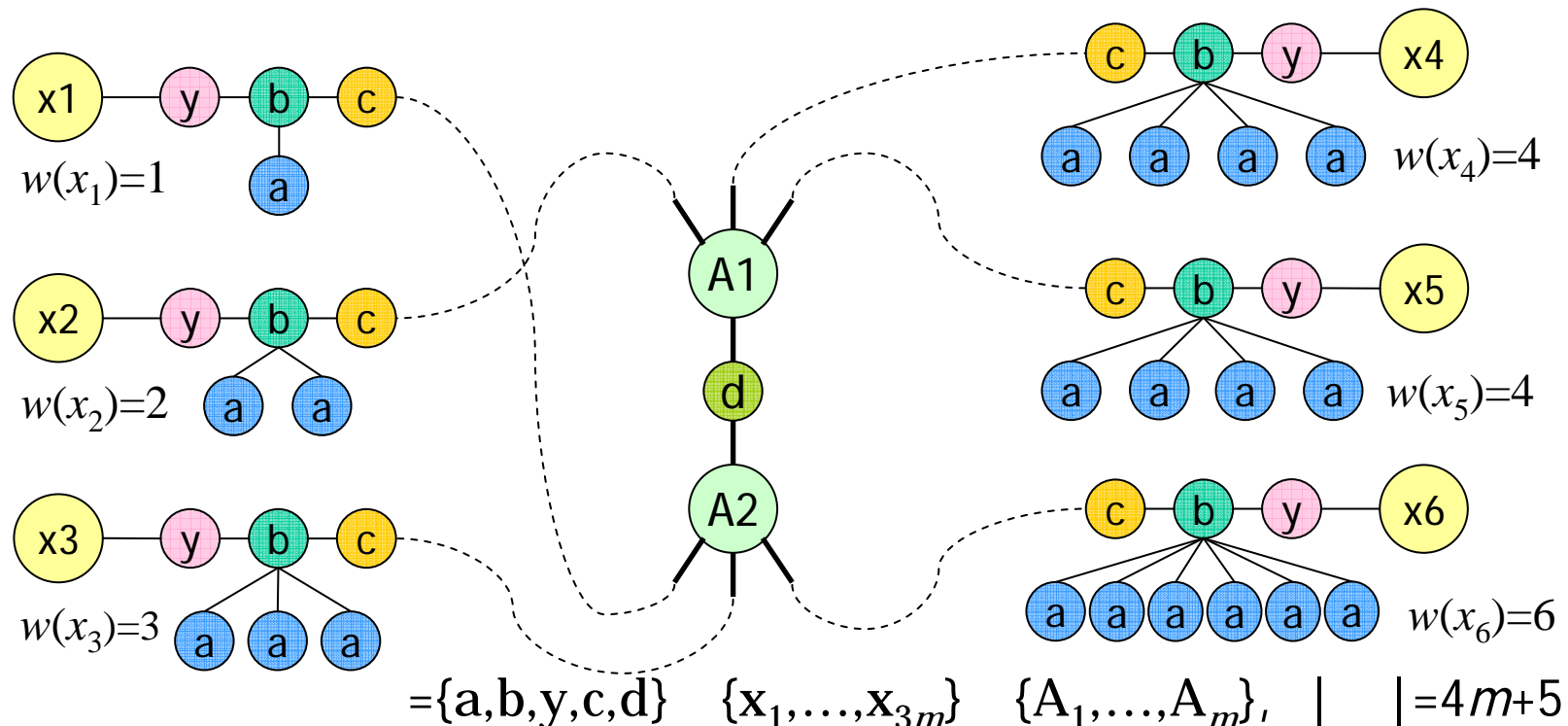


$$= \{a, b, y, c, d\} \quad \{x_1, \dots, x_{3m}\} \quad \{A_1, \dots, A_m\}, \quad | \quad | = 4m + 5$$

Feature vector specifies structures of blocks, but does not specify the connection between blocks $\{x_i\}$ and $\{A_1, A_2\}$.

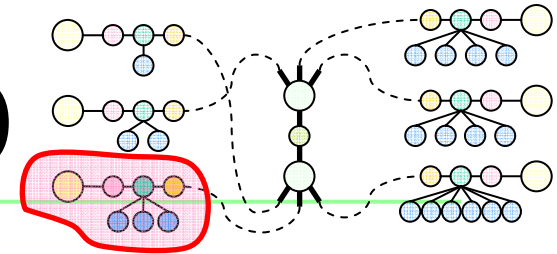
An Example of Reduction (5)

3-PARTITION P : $w = (1, 2, 3, 4, 4, 6) \rightarrow$ GIPF Q : $(v, 3)$

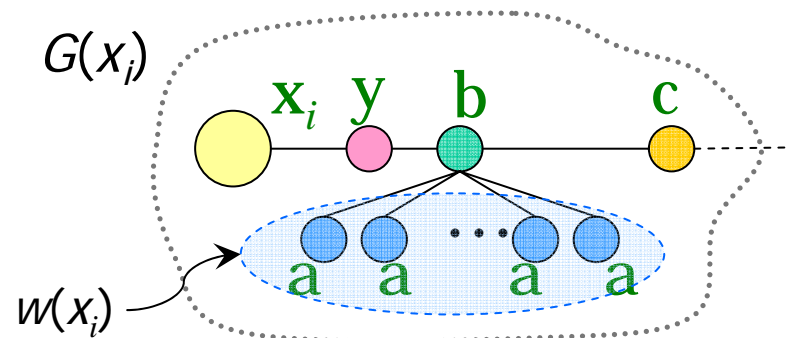


- The connection satisfying the constraints given by feature vector corresponds to a solution of 3-PARTITION.
- In this case, $\{x_1, x_3, x_6\}$ and $\{x_2, x_4, x_5\}$ correspond to a solution of 3-PARTITION.

An Example of Reduction (6)



For each $x_i; i=1,2,\dots,3m$, generate a graph $G(x_i)$:



$G(x_i)$ encodes $w(x_i)$

- A label x_i is unique in the whole graph
- # of a 's = $w(x_i)$

Paths of length 3 which determine $G(x_i)$

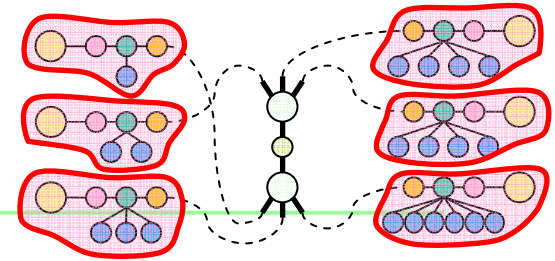
0 --- $\{(\underline{x}_i:1), (\underline{a}:w(x_i)), (\underline{y}:1), (\underline{b}:1), (\underline{c}:1)\}$

1 --- $\{(\underline{x}_i\underline{y}:1), (\underline{y}\underline{x}_i:1), (\underline{y}\underline{b}:1), (\underline{b}\underline{y}:1), (\underline{b}\underline{c}:1), (\underline{c}\underline{b}:1), (\underline{a}\underline{b}:w(x_i)), (\underline{b}\underline{a}:w(x_i))\}$

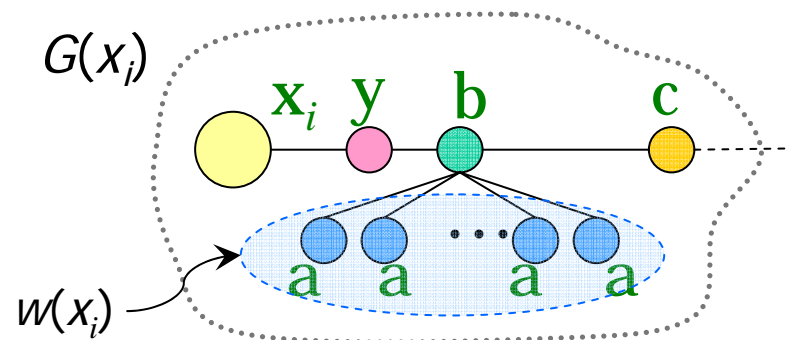
2 --- $\{(\underline{x}_i\underline{y}\underline{b}:1), (\underline{y}\underline{b}\underline{a}:w(x_i)), (\underline{y}\underline{b}\underline{c}:1), (\underline{b}\underline{y}\underline{x}_i:1), (\underline{c}\underline{b}\underline{y}:1), (\underline{c}\underline{b}\underline{a}:w(x_i)), (\underline{a}\underline{b}\underline{y}:w(x_i)), (\underline{a}\underline{b}\underline{c}:w(x_i)), (\underline{a}\underline{b}\underline{a}:w(x_i)(w(x_i)-1))\}$

3 --- $\{(\underline{x}_i\underline{y}\underline{b}\underline{a}:w(x_i)), (\underline{x}_i\underline{y}\underline{b}\underline{c}:1), (\underline{c}\underline{b}\underline{y}\underline{x}_i:1), (\underline{a}\underline{b}\underline{y}\underline{x}_i:w(x_i))\}$

An Example of Reduction (7)



For each $x_i; i=1,2,\dots,3m$, generate a graph $G(x_i)$:



$G(x_i)$ encodes $w(x_i)$

- A label x_i is unique in the whole graph
- # of a 's = $w(x_i)$

In total of $G(x_1), G(x_2), \dots, G(x_{3m})$

0 --- $\{x_i:1\}$ $\{(a:mB), (y:3m), (b:3m), (c:3m)\}$

1 --- $\{x_i y:1, (y x_i:1)\}$ $\{(ab:mB), (ba:mB), (yb:3m), (by:3m), (bc:3m), (cb:3m),\}$

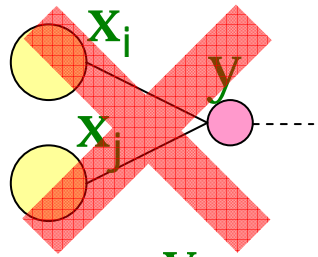
2 --- $\{x_i y b:1, (b y x_i:1)\}$ $\{(y b a:mB), (y b c:3m), (c b y:3m), (c b a:mB), (a b y:mB), (a b c:mB), (a b a: w(x_i)^2 - mB)\}$

3 --- $\{x_i y b a:w(x_i), (x_i y b c:1), (c b y x_i:1), (a b y x_i:w(x_i))\}$

Note for Uniqueness of Graph $G(x_i)$

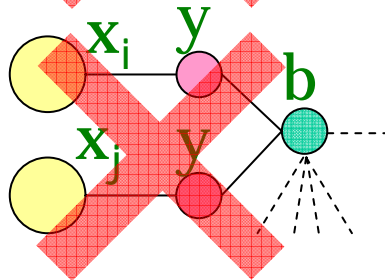
It is necessary to prove that **the set of paths uniquely determines $G(x_i)$** .

- The following cases does NOT occur:



--- Because a path ' $x_i y x_j$ ' is not given

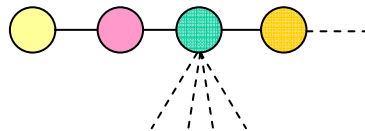
→ 1-to-1 correspondence between (x_i, y)



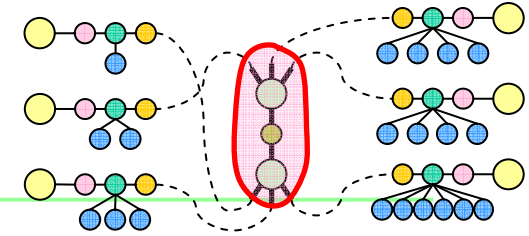
--- Because a path ' $y b y$ ' is not given

(Even if tottering (backtrack) is admitted, provable similarly by using # of ' $y b$ ')

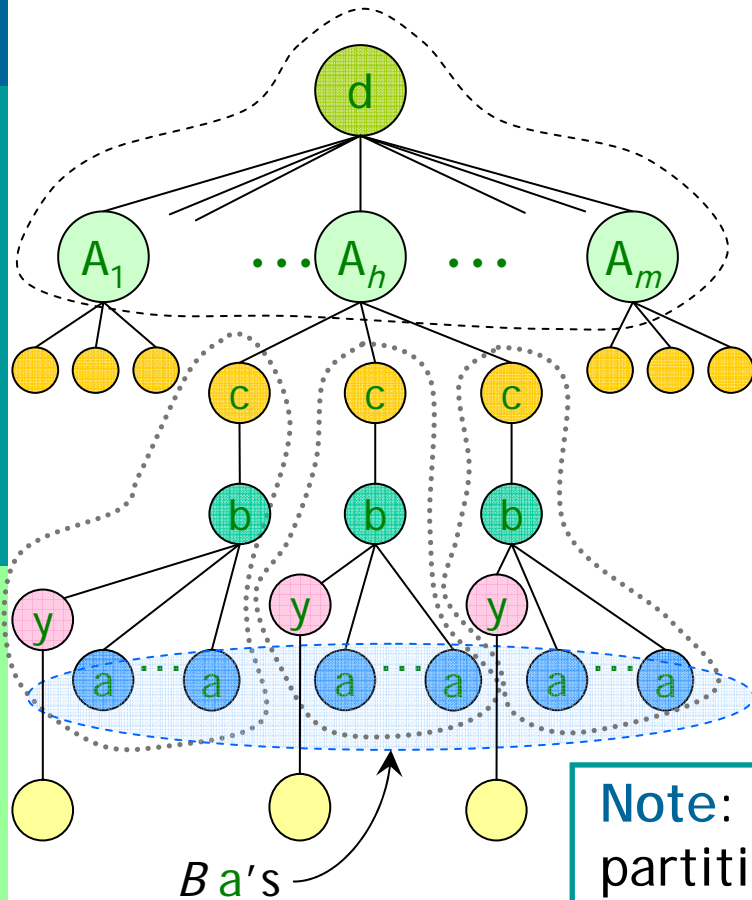
Uniqueness for a quadruple (x_i, y, b, c) can be proved in a similar way



An Example of Reduction (8)



Generation of center graph G_C



Remaining paths in G_C

- 0 --- ${}_h\{(A_h:1)\} \quad \{(d:1)\}$
- 1 --- ${}_h\{(A_hc:3m), (A_hd:m), (cA_h:3m), (dA_h:m)\}$
- 2 --- ${}_h\{(A_hcb:3), (bcA_h:3), (cA_hc:6), (cA_hd:3), (dA_hc:3)\}$
 ${}_{h,k}\{(A_hdA_k:1)\}$
- 3 --- ${}_h\{(A_hcby:3), (ybcA_h:3), (A_hcba:B), (abcA_h:B), (dA_hcb:3), (bcA_hd:3), (cA_hcb:6), (bcA_hc:6)\}$
 ${}_{h,k}\{(A_hdA_kc:3), (cA_kdA_h:3)\}$

Note: center graph is determined without knowing partition (without information about $w(x_i)$'s)

Strong NP-completeness of GIPF

⟨Theorem⟩ GIPF is strongly NP-complete, even if $K=3$ and the underlying graph is a tree.

Hardness results for other special case

(Theorem) GIPF is strongly NP-complete, even for trees of bounded degree 4 and of fixed K .

Bounded degrees ()

- Branchings for \underline{a} 's and for center \underline{d} \rightarrow Use **binary tree**

Bounded alphabets ()

- \underline{x}_i 's and \underline{A}_h 's \rightarrow Encode with fixed alphabets

In both cases, we cannot bound K by a constant

Note: if all of $\underline{a}, \underline{d}, \underline{x}_i, \underline{A}_h, K$ are fixed, then the problem can be solved in poly. time

Conclusion

- GIPF is **strongly NP-complete** even if underlying graph is a tree and $K=3$.
- GIPF (and GIPF-M) for trees is **solvable in polynomial time** by using DP, if n , K and Δ are all fixed.

Still ongoing:

- Our DP is extendable to **outer-planar graphs**
- Completeness results for more restricted cases

Future work:

- Complexity of SISF-M in general cases
- Approximation algorithm, etc.