

# グラフ構造データを 対象とした 解釈可能決定集合の拡張

5416063 松山航太

機械学習の解釈性研究

Interpretable  
Decision Sets

グラフ構造データ  
への適用



Graph\_Interpretable  
Decision Sets  
(本手法)

H, Lakkaraju, et al : 『Interpretable Decision Sets: A Joint Framework for Description and Prediction』 ,KDD '16(2016)

1. ルール集合の獲得

	○○	××	...	class
D1	No	Yes	...	c1
D2	Yes	No	...	c2
...	...	...	...	...

○○=Y and ××=Y then c1  
○○=N then c2  
××=N and △△=N then c2  
...

2. 適切な部分集合の獲得

○○=Y and ××=Y then c1  
××=N and △△=N then c2

データセット  
 $D$

クラス相関ルール (if-thenルール) を用いて  
ルールセットを獲得

ルール集合  
 $S$

評価関数  $f(R, D)$  に従って  $R$  を選択  
$$R^* = \underset{R \in S}{\operatorname{argmax}} (f(R, D))$$

解釈可能  
決定集合  
 $R$

# 評価関数

データセット

$$R^* = \underset{R \subseteq S}{\operatorname{argmax}} (f(R, D))$$

出力となる  
選択ルール集合

$$f(R, D) = \sum_{i=1}^7 \lambda_i f_i(R, D)$$

クラス相関ルールの  
全体集合Sの  
部分集合

重み  
( $\lambda_i \geq 0$ )

各観点からの評価

# 集合サイズ (小さいほど高評価)

$$f_1(R, D) = |S| - \text{size}(R)$$

$$\begin{array}{l} \text{size}(R) \\ \parallel \\ 3 \end{array} \left\{ \begin{array}{l} \bigcirc\bigcirc = Y \text{ and } \times\times = Y \text{ then } c1 \\ \times\times = N \text{ and } \triangle\triangle = N \text{ then } c2 \\ \triangle\triangle = N \text{ then } c2 \end{array} \right.$$

# ルールの本体部の長さ (短いほど高評価)

$$f_2(R, D) = L_{max} \cdot |S| - \sum_{r \in R} \text{length}(r)$$

○○=Y and ××=Y then c1

××=N and △△=N then c2

△△=N } then c2

$$\begin{aligned} \text{length}(r) \quad \sum_{r \in R} \text{length}(r) &= 2 + 2 + 1 \\ &= 5 \end{aligned}$$

# ルールの説明範囲の重複度合い（小さいほど高評価）

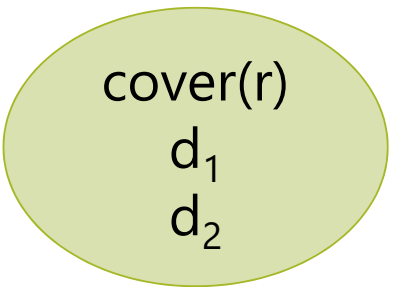
$$f_3(R, D) = N \cdot |S|^2 - \sum_{\substack{r_i, r_j \in R \\ i \leq j \\ c_i = c_j}} \text{overlap}(r_i, r_j) \quad f_4(R, D) = N \cdot |S|^2 - \sum_{\substack{r_i, r_j \in R \\ i \leq j \\ c_i \neq c_j}} \text{overlap}(r_i, r_j)$$

$$\text{overlap}(r, r') = \text{cover}(r) \cap \text{cover}(r')$$

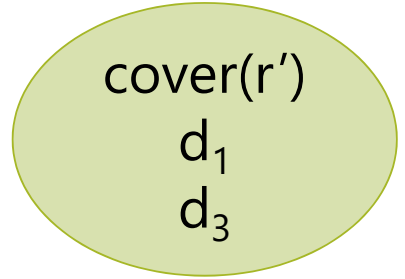
データセットD

	○○	××	...	<del>class</del>
d <sub>1</sub>	Y	Y	...	<del>c1</del>
d <sub>2</sub>	Y	N	...	<del>c2</del>
d <sub>3</sub>	N	Y	...	<del>c2</del>

ルールr  
○○=Y then c1



ルールr'  
××=Y then c2



$$\text{overlap}(r, r') = 1$$

d<sub>1</sub>



# クラスの網羅性 (大きいほど高評価)

$$f_5(R, D) = \sum_{c' \in C} 1 \quad (\exists r = (s, c) \in R \text{ such that } c = c')$$



○○=Y and ××=Y then c1  
○○=N and △△=Y then c2

c1, c2  
両方 = 2



○○=Y and ××=Y then c1  
○○=Y and △△=Y then c1

c1  
のみ = 1



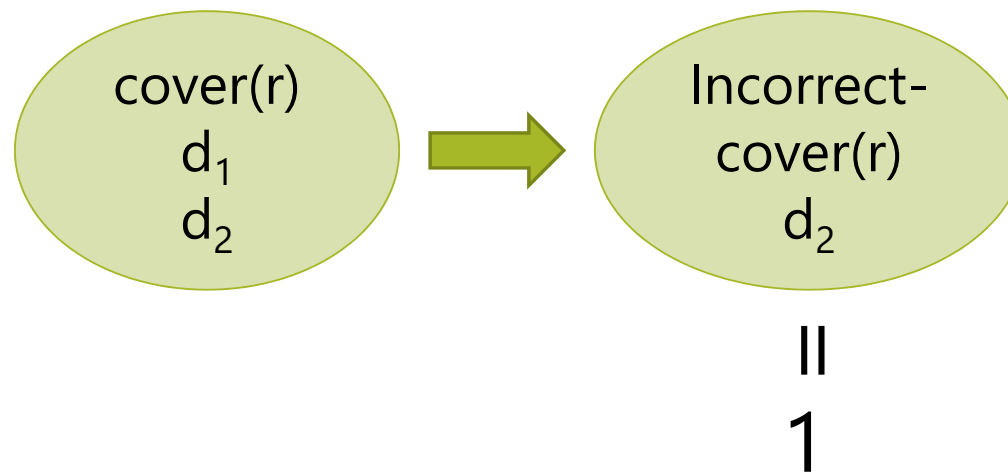
# 誤分類される事例数 (小さいほど高評価)

$$f_6(R, D) = N \cdot |S| - \sum_{r \in R} |\text{incorrect-cover}(r)|$$

	○○	××	...	class
d <sub>1</sub>	Y	Y	...	Y
d <sub>2</sub>	Y	N	...	N
d <sub>3</sub>	N	N	...	N

ルール $r$

○○=Y then c1



# 事例の被覆度 (大きいほど高評価)

$$f_7(R, D) = \sum_{(x,y) \in D} 1(|\{r | (x,y) \in \text{correct-cover}(r)\}| \geq 1)$$

データセットD

	○○	××	...	class
d <sub>1</sub>	Y	Y	...	c1
d <sub>2</sub>	Y	N	...	c2
d <sub>3</sub>	N	N	...	c3

r<sub>1</sub>: ○○ = Y then c1  
r<sub>2</sub>: ×× = N then c2

correct-  
cover(r<sub>1</sub>)  
d<sub>1</sub>

correct-  
cover(r<sub>2</sub>)  
d<sub>2</sub> d<sub>3</sub>

$$f_7(R) = |\{d_1, d_2, d_3\}| = 3$$

$$R^* = \underset{R \subseteq S}{\operatorname{argmax}} (f(R, D))$$

### 劣モジュラ最大化問題

$f_1(R, D)$ から $f_7(R, D)$ は全て劣モジュラ関数  
既存のアルゴリズム(SLS等)で計算可能

劣モジュラ関数： $f(S \cup \{i\}) - f(S) \geq f(T \cup \{i\}) - f(T)$

台集合： $V, S, T \subseteq V, S \subseteq T, i \notin Y$

## IDS(従来手法)

	○○	××	...	class
D1	No	Yes	...	Y
D2	Yes	No	...	N
...	...	...	...	...

グラフデータの  
トランザクション化

○○=Y and ××

○○=N then c2

××

△△=N then c1

...

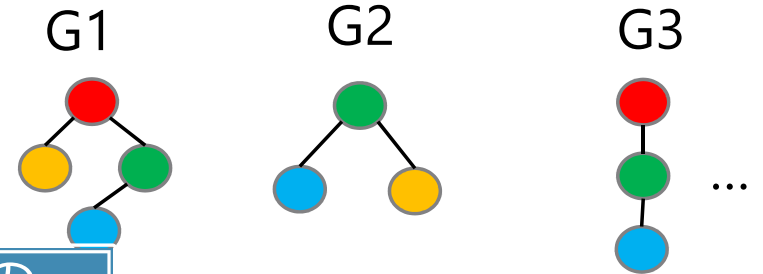
グラフを考慮した  
評価関数の拡張

○○=Y and ××

××

△△=N then c2

## GIDS(本研究)



$g_1=Y$  and  $g_2=Y$  then c1

$g_3=Y$  then c2

$g_2=Y$  and  $g_3=Y$  then c2

...

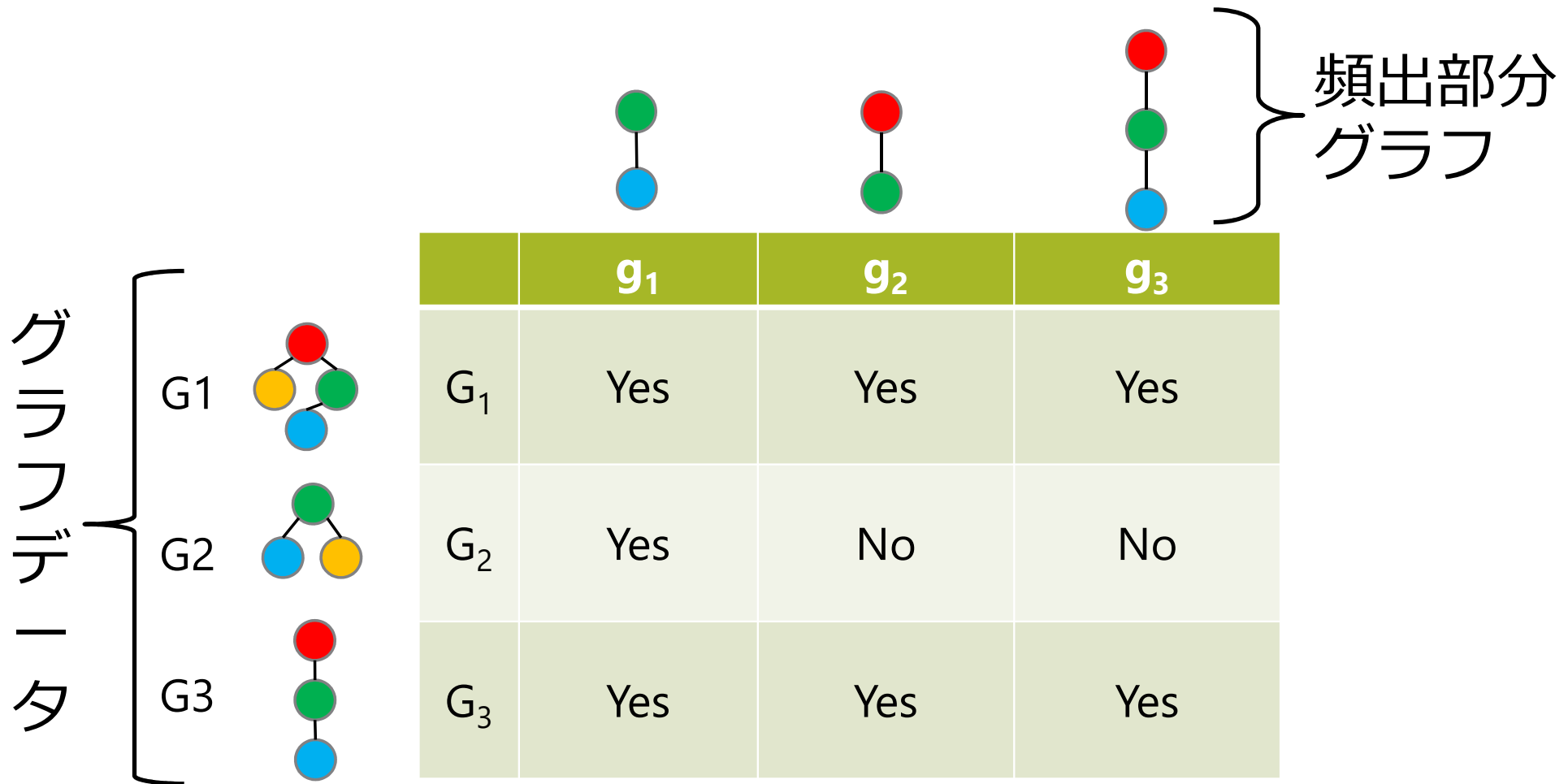
$g_1=Y$  and  $g_2=Y$  then c1

$g_3=Y$  then c2

1.  
ルール集合  
の獲得

2.  
適切な  
部分集合の  
獲得

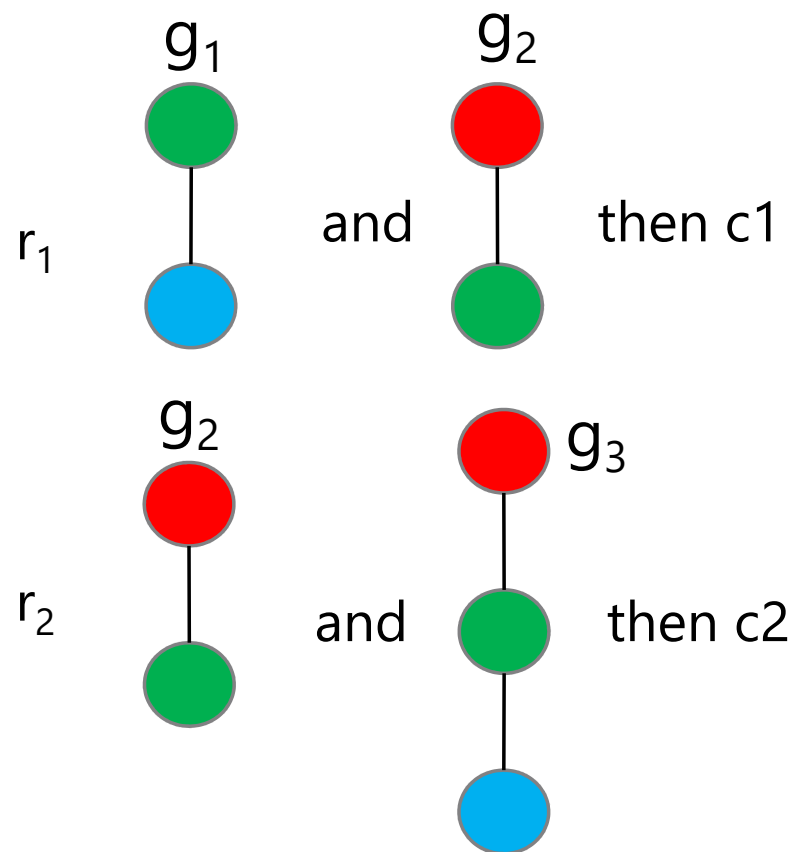
# グラフ構造データのトランザクション化



# 評価関数の拡張

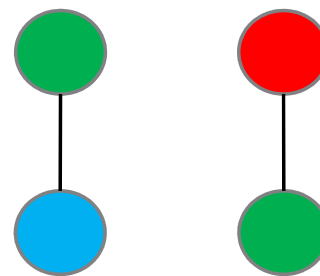
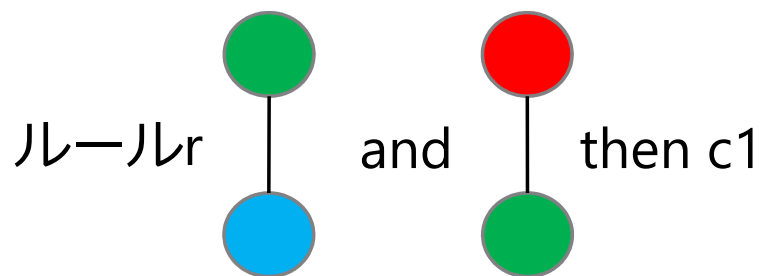
評価関数をルール中の  
グラフの形状を考慮するよう拡張

$g_1=Y$  and  $g_2=Y$  then  $c_1$   
 $g_2=Y$  and  $g_3=Y$  then  $c_2$   
...



# グラフサイズ (小さいほど高評価)

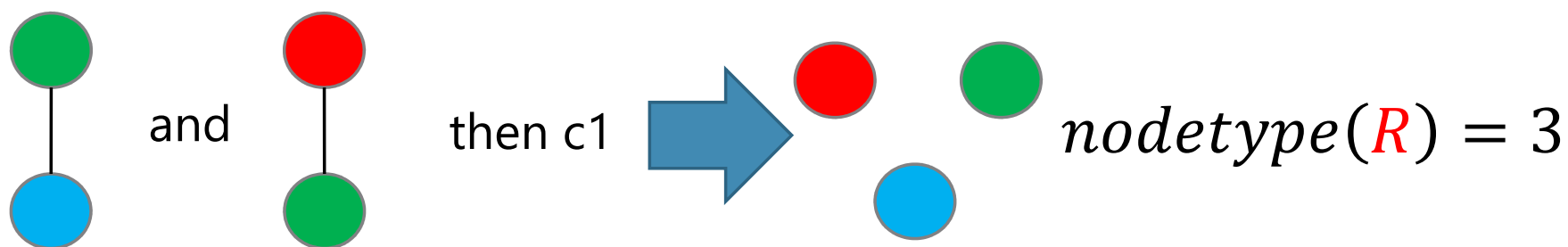
$$f_8(R, D) = \sum_{r \in R} (G_{max} - \text{graphsize}(r))$$



$$\text{graphsize}(r) = 4 + 2 = 6$$

# ノードラベルの種類数（多いほど高評価）

$$f_9(R, D) = \text{nodetype}(R)$$



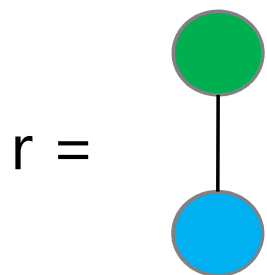


# グラフ非類似度 (大きいほど高評価)

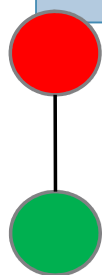
グラフ編集距離  
(Graph edit distance)

$$f_{10}(R, D) = \sum_{r \in R} (\text{graph dissimilarity}(r))$$

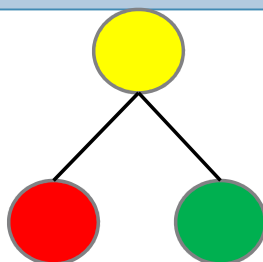
$$\text{graph dissimilarity}(r) = \frac{1}{|G|C_2} \sum_{g_i, g_j \in G} \text{GED}(g_i, g_j)$$



and



and

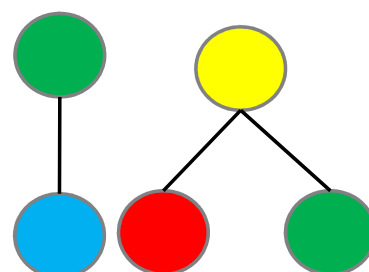
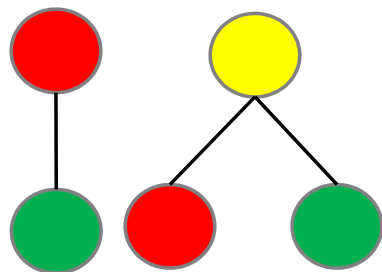
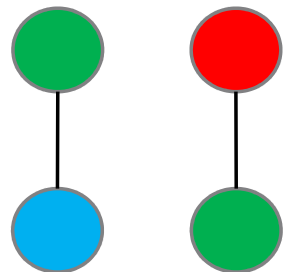


then c1

$$\{\text{GED}(g_1, g_2)$$

$$+ \text{GED}(g_2, g_3)$$

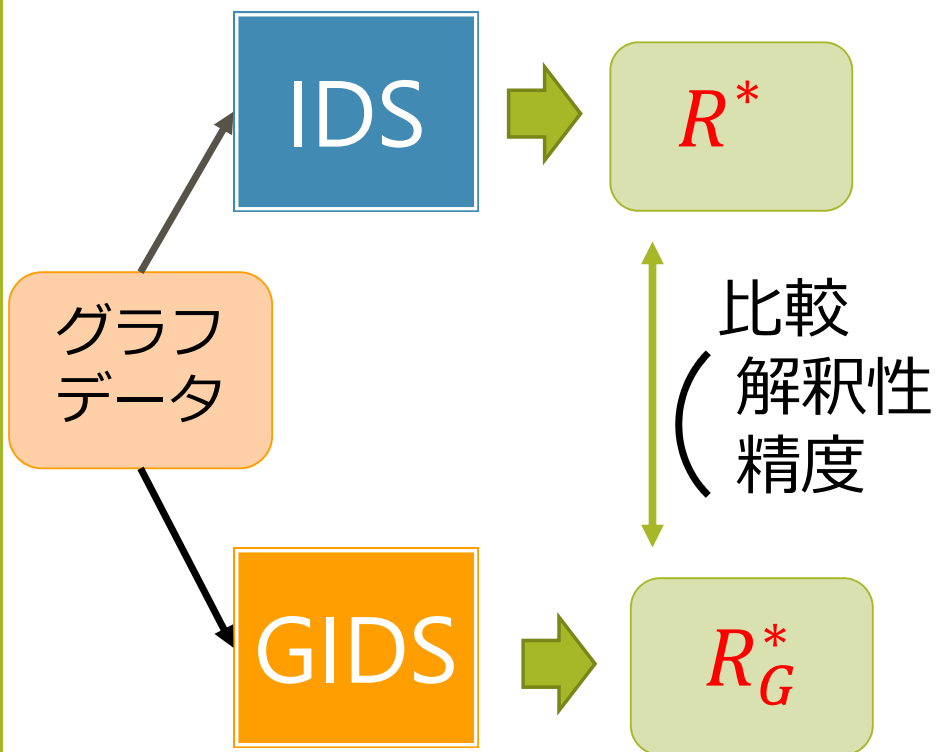
$$+ \text{GED}(g_1, g_3)\} / 3$$



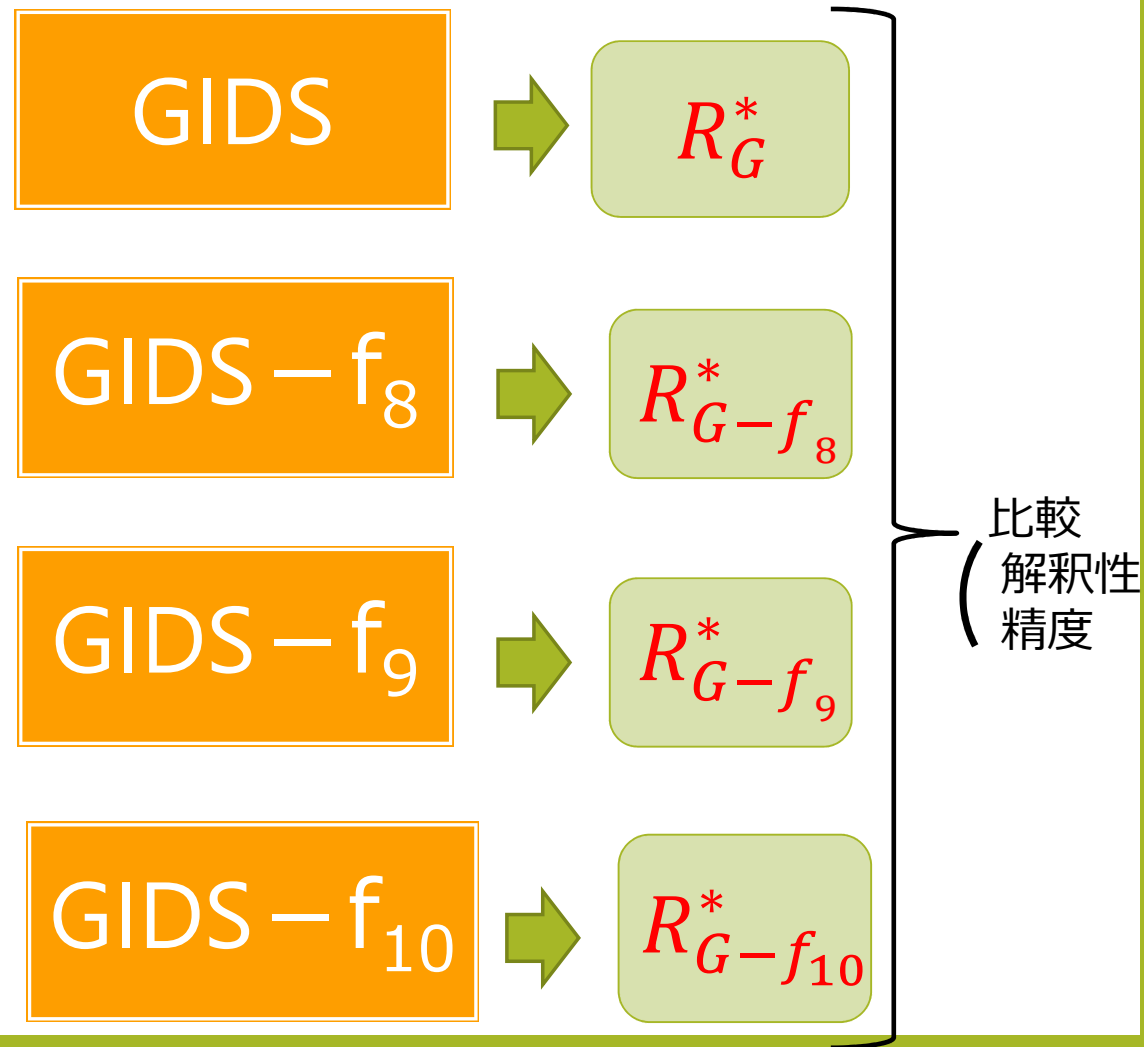
$$= \text{graph dissimilarity}(r)$$

# 実験設定

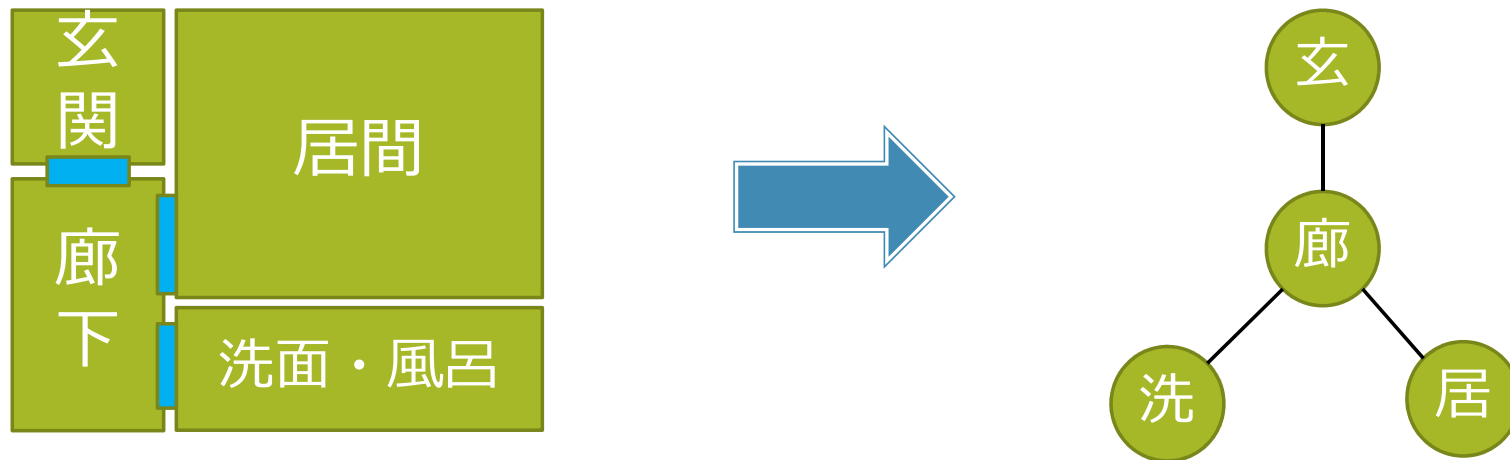
## 拡張の有無の比較



## 各関数の有無の比較



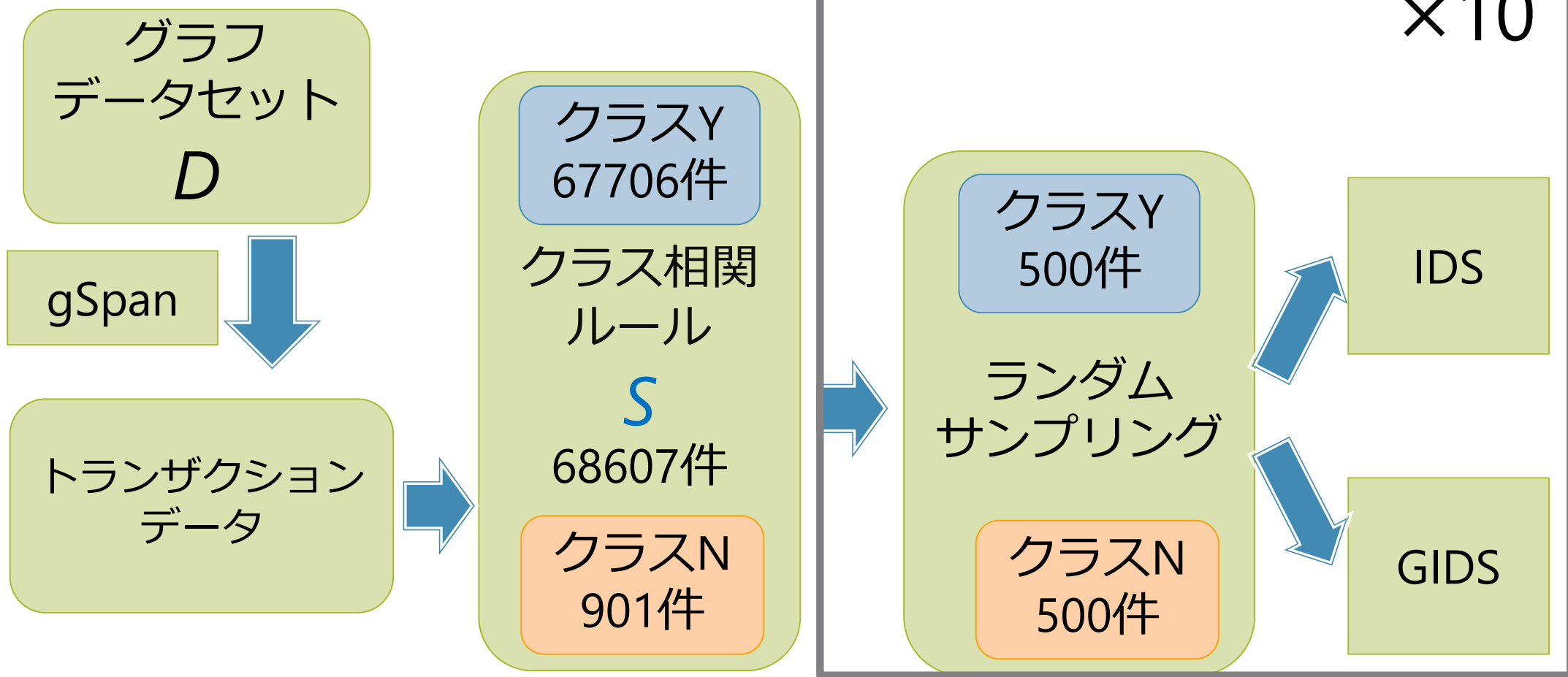
# データセット:Homes間取り図グラフ



データ数:480件 (Y,N各クラス240件ずつ)

T, Ozaki:Extraction of Characteristic Subgraph Patterns with Support Threshold from Databases of Floor Plans, Cander2019

# 実験設定



グラフサイズ

ノード種類数

グラフ非類似度

	IDS	GIDS	$\neg f_8$	$\neg f_9$	$\neg f_{10}$
精度 (AUC)	0.52	<b><u>0.58</u></b>	<b>0.52</b>	0.58	0.58
	0.35	<b><u>0.00</u></b>	<b>0.35</b>	0.00	0.00
	0.22	<b><u>0.11</u></b>	<b>0.22</b>	0.11	0.11
	<b><u>2.95</u></b>	3.00	<b>2.95</b>	3.00	3.00
	4.20	4.20	<b>2.00</b>	4.20	4.20
	1.00	1.00	1.00	1.00	1.00
	9.95	<b><u>10.94</u></b>	<b>12.25</b>	<b>10.88</b>	10.94
ノード種類数 : $f_9$	7.80	<b><u>9.40</u></b>	<b>7.90</b>	<b>9.30</b>	9.40
非類似度 : $f_{10}$	0.26	<b><u>0.50</u></b>	<b>0.27</b>	0.50	0.50

・ グラフサイズの関数を除くと性能が大きく下がる

・ 他の関数は誤差レベルの貢献

## まとめ

- ・ IDSのグラフ構造データへの拡張  
- 精度、解釈性の向上

## 今後の課題

- ・ 追加実験
- ・  $\lambda$ の最適化
- ・ 主観評価