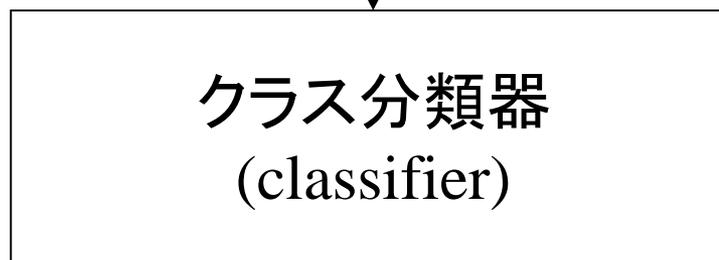


# 様々な決定木とその構成方法

和田 俊和

T1	T2	T3	T4	目標属性
1	0	1	1	1
1	0	1	1	1
1	1	1	1	1
1	1	1	0	0
1	0	1	0	0
1	1	0	1	0
1	0	0	1	0
1	1	0	1	0
0	1	0	1	1
0	0	1	1	1
0	1	0	1	1
0	1	0	1	1
0	0	0	1	1
0	0	1	0	0
0	1	0	0	0
0	0	1	0	0

T1, T2, T3, T4 の値



出力

(目標属性に一致させる)

# 分類ルール

テスト

T1	T2	T3	T4	目標属性 (Objective)
1	0	1	1	1
1	0	1	1	1
1	1	1	1	1
1	1	1	0	0
1	0	1	0	0
1	1	0	1	0
1	0	0	1	0
1	1	0	1	0
0	1	0	1	1
0	0	1	1	1
0	1	0	1	1
0	1	0	1	1
0	0	0	1	1
0	0	1	0	0
0	1	0	0	0
0	0	1	0	0

If T1=1

Then

If T3=1

Then

If T4=1

Then 目標属性=1

Else 目標属性=0

Else

目標属性=0

Else

If T4=1

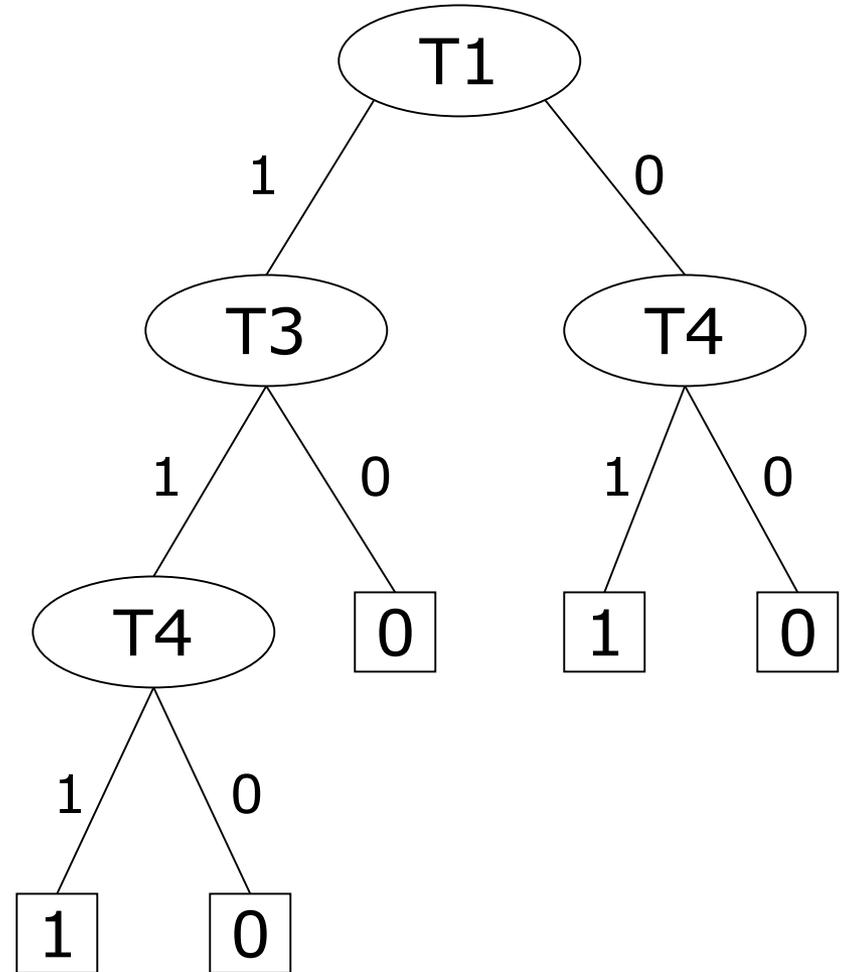
Then 目標属性=1

Else 目標属性=0

# 決定木(decision tree)

If T1=1  
Then  
    If T3=1  
    Then  
        If T4=1  
        Then 目標属性=1  
        Else 目標属性=0  
    Else  
        目標属性=0  
Else  
    If T4=1  
    Then 目標属性=1  
    Else 目標属性=0

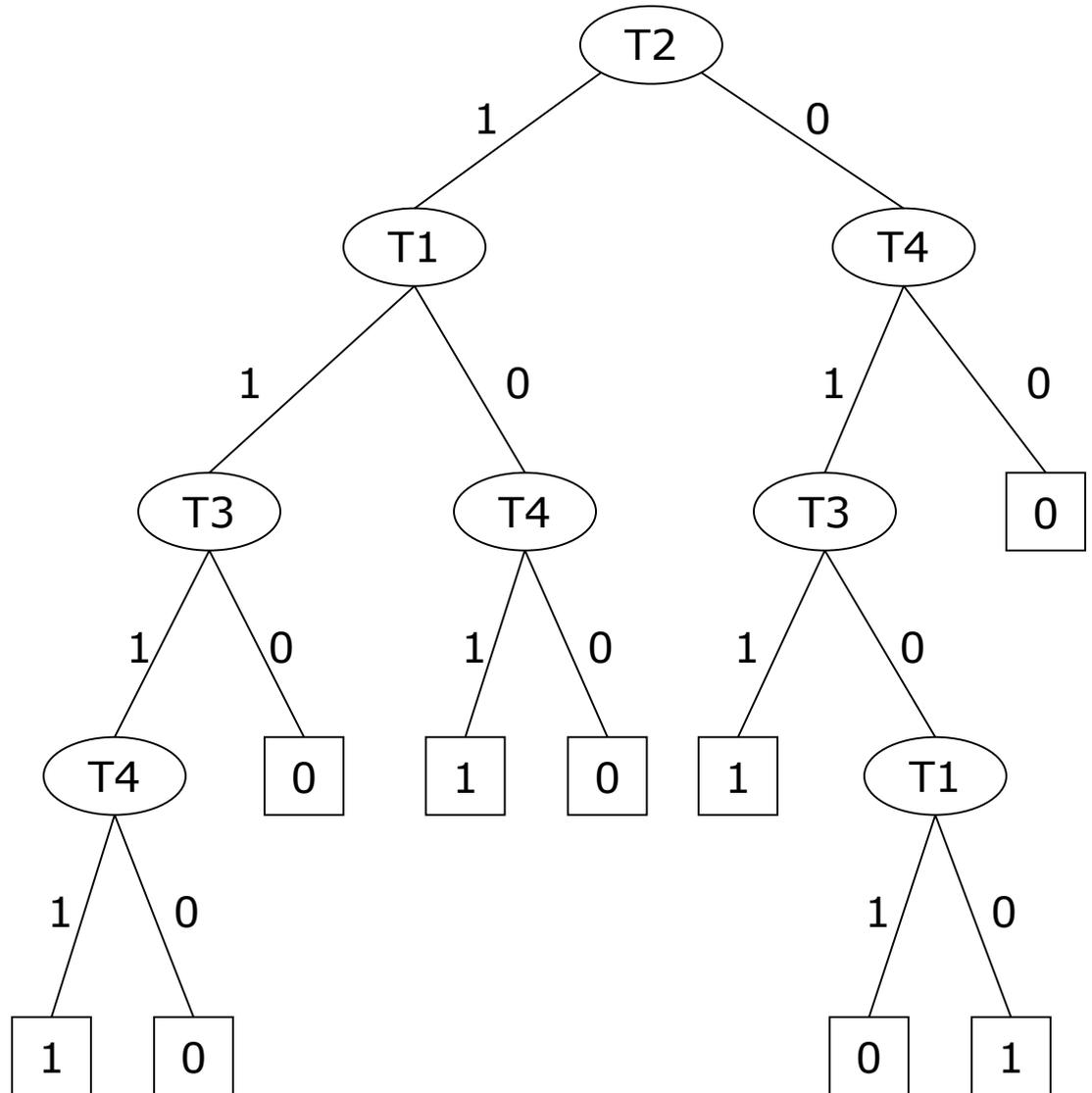
ルールによる表現



木構造表現

# テストの順番を変えると異なる決定木が得られる

T1	T2	T3	T4	目標属性
1	1	1	1	1
1	1	1	0	0
1	1	0	1	0
1	1	0	1	0
0	1	0	1	1
0	1	0	1	1
0	1	0	1	1
0	1	0	0	0
1	0	1	1	1
1	0	1	1	1
0	0	1	1	1
0	0	0	1	1
1	0	0	1	0
0	0	1	0	0
1	0	1	0	0
0	0	1	0	0



# 用語

属性(attribute)

テスト属性 (test attributes)

目標属性(Objective attribute)

T1	T2	T3	T4	Ao
----	----	----	----	----

1	1	1	1	1
1	1	1	0	0

レコード(record)

レコード  $t$  の属性  $A$  の値を  $t[A]$  と表記

**目標値** 目標属性の値の1つ. 例えば 1

目標属性  $A_o$ 、目標値が 1 のとき、

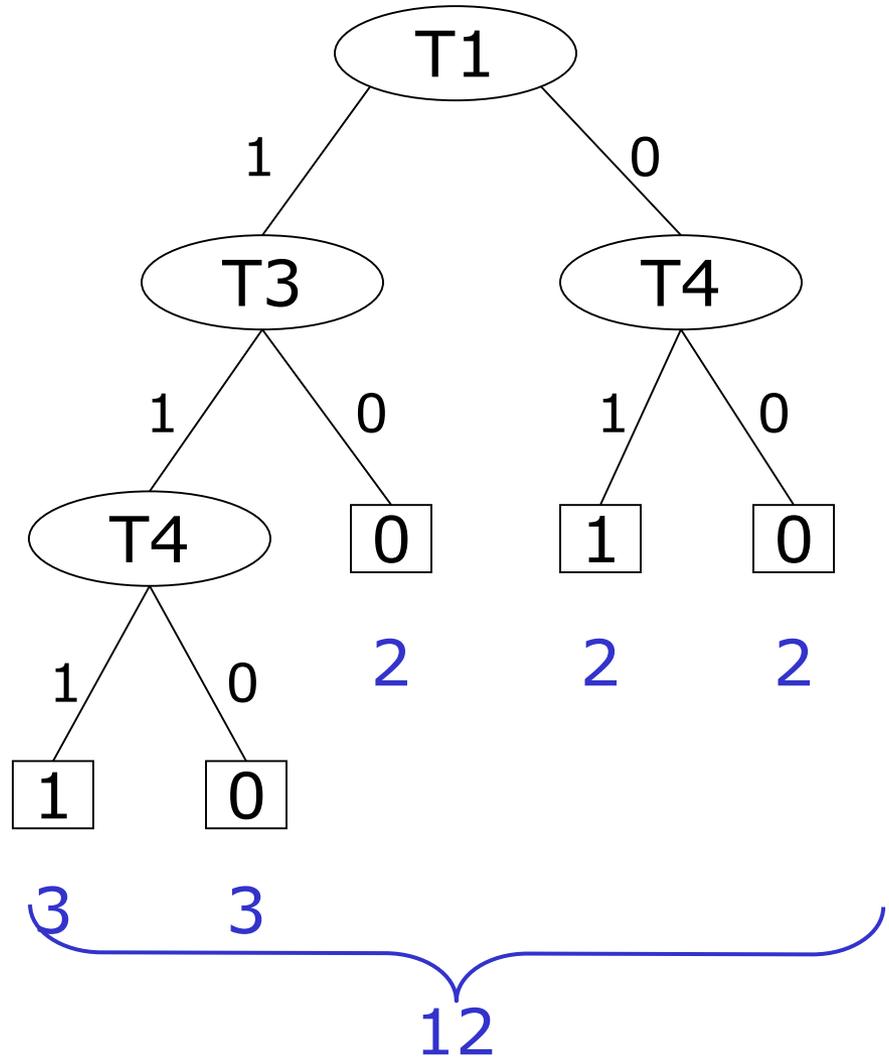
- $t[A_o]=1$  となる  $t$  を **正(positive)**
- $t[A_o]\neq 1$  となる  $t$  を **負(negative)**

# 決定木の最小化

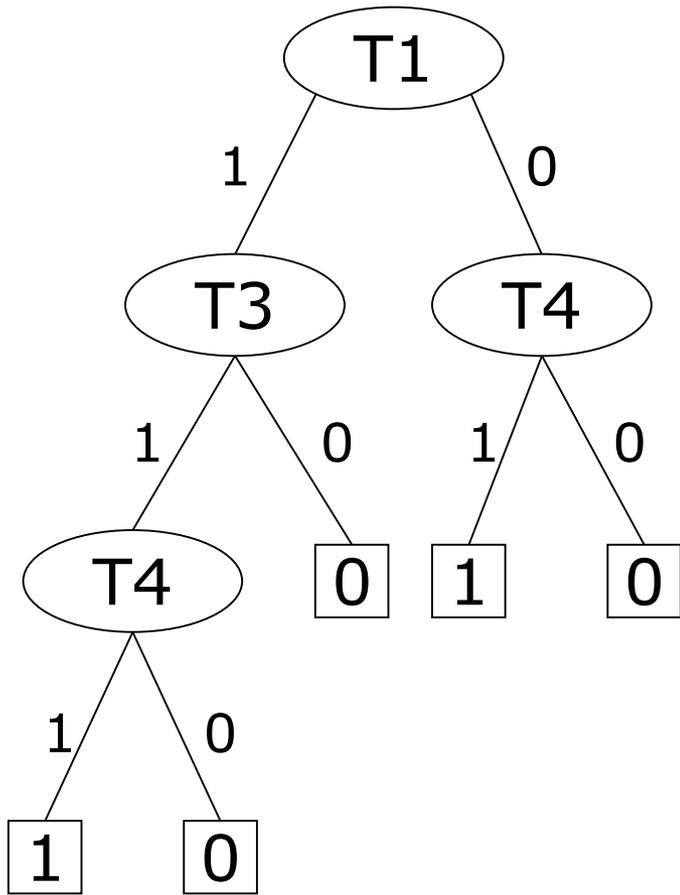
# 決定木のコスト

## 決定木の評価

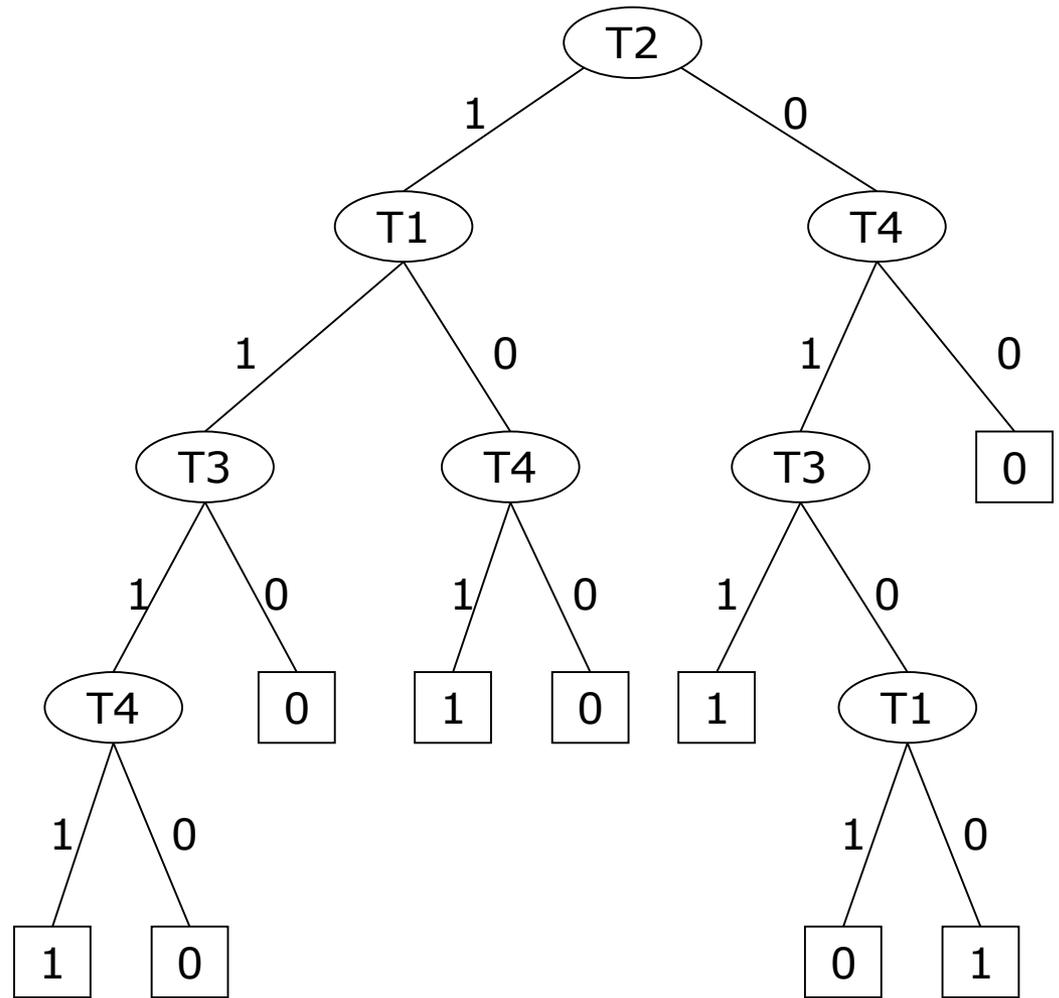
- 根ノードから葉ノード  $x$  に至るテストの数  $\text{cost}(x)$
- 決定木のコスト  
 $= \sum \{ \text{cost}(x) \mid x \text{ は葉ノード} \}$
- 決定木のコストは小さいほうが望ましい



# コストの比較



$$3+3+2+2+2=12$$



$$4+4+3+3+3+3+4+4+2=30$$

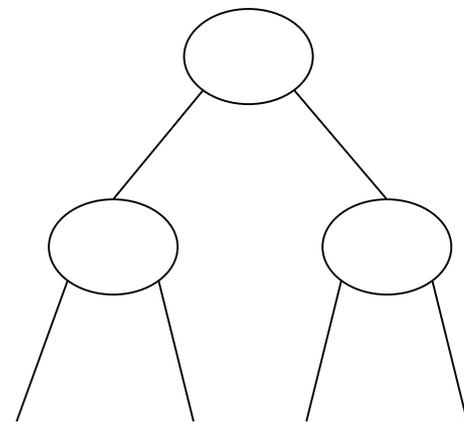
近似的解法

Greedy Algorithm

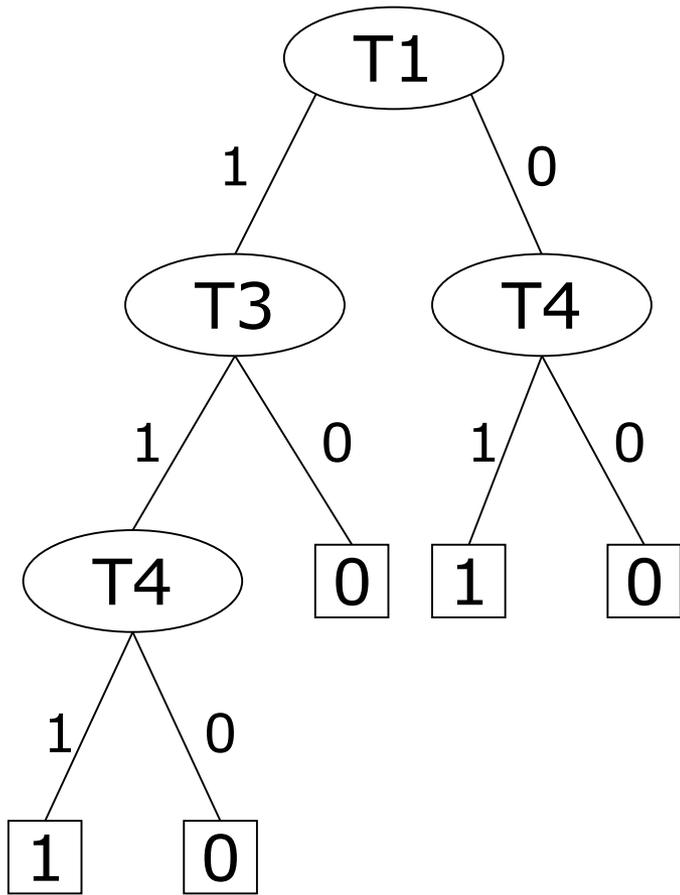
エントロピー

# 近似的解法 貪欲(greedy)アルゴリズム

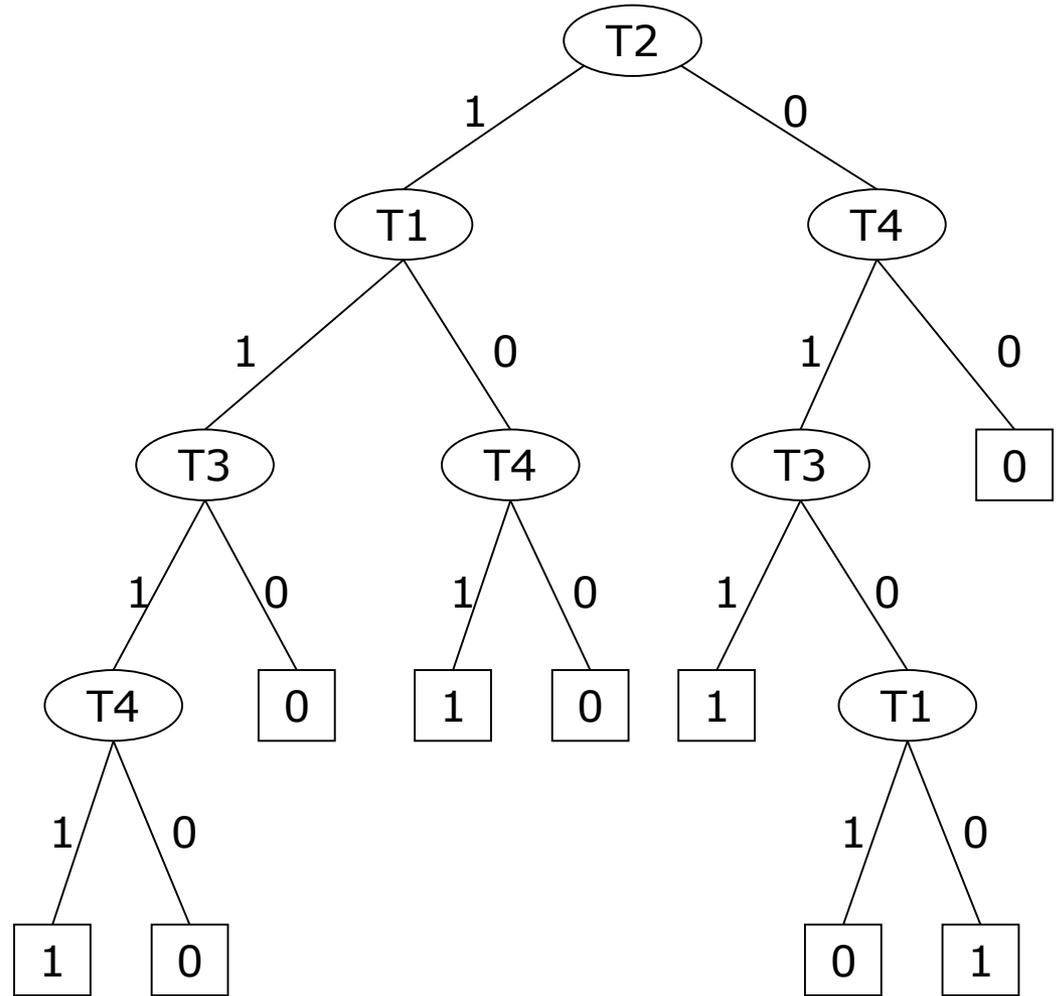
- 「最も良さそうな」テストを選択し、レコード集合を分割
- 分割したレコード集合に同じ戦略を繰り返し適用する
- 一度選択したテストは変更しない
- 「最も良さそうな」の基準を工夫



# 最初にT2ではなくT1を選ぶべき



$$3+3+2+2+2=12$$

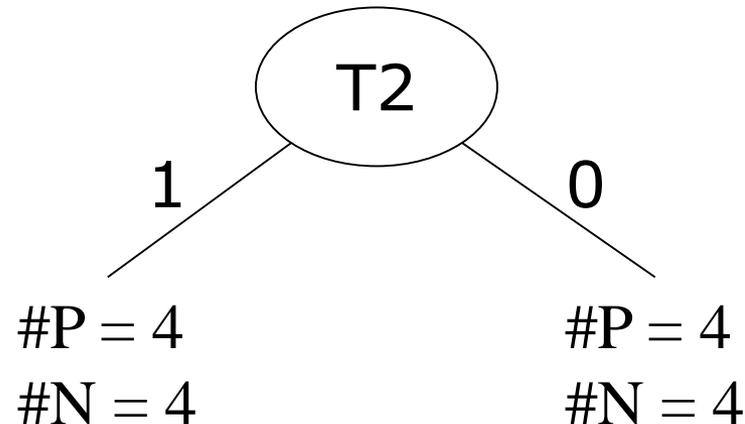
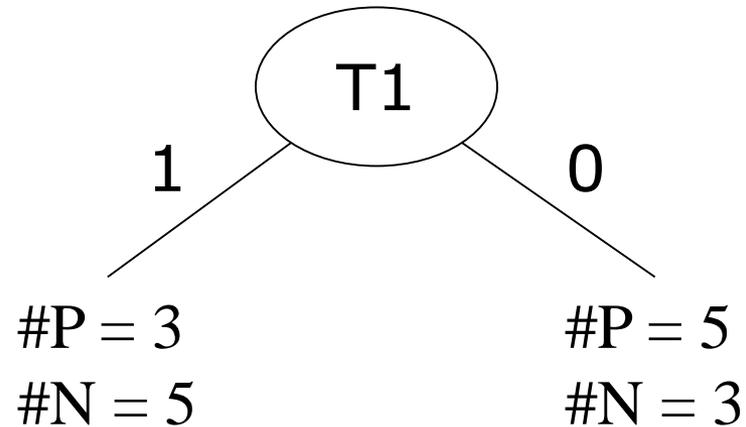


$$4+4+3+3+3+3+4+4+2=30$$

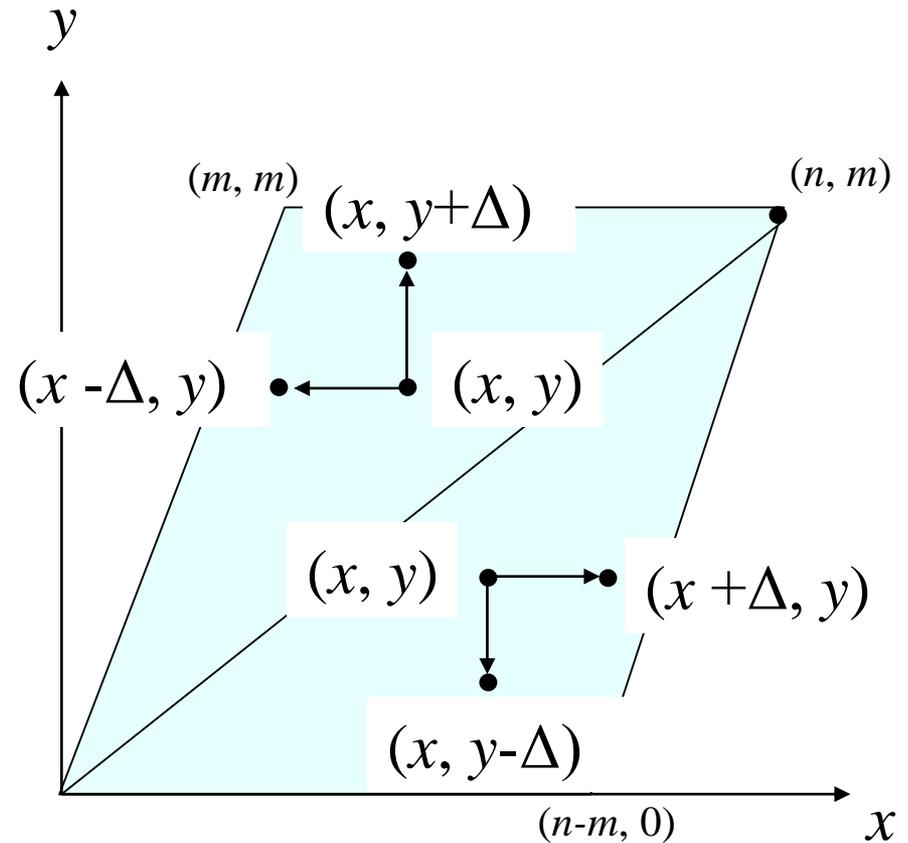
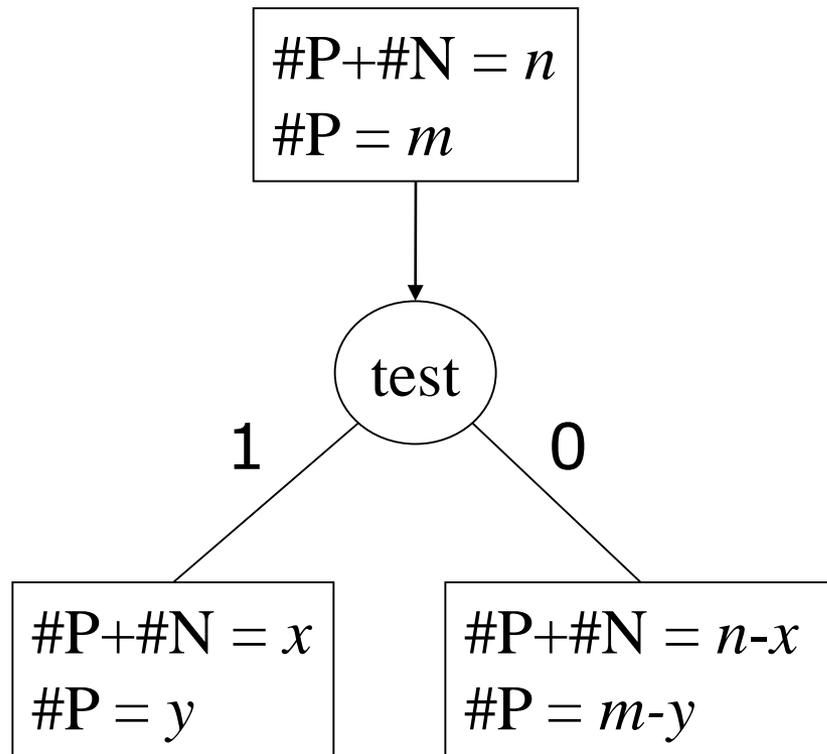
# T2 ではなく T1 を選択する評価基準は？

T1	T2	T3	T4	目標属性
1	0	1	1	1
1	0	1	1	1
1	1	1	1	1
1	1	1	0	0
1	0	1	0	0
1	1	0	1	0
1	0	0	1	0
1	1	0	1	0
0	1	0	1	1
0	0	1	1	1
0	1	0	1	1
0	1	0	1	1
0	0	0	1	1
0	0	1	0	0
0	1	0	0	0
0	0	1	0	0

#P : 正レコードの数  
#N : 負レコードの数



# テスト選択のための評価値を定める



- test を選択すると  $(x, y)$  が定まる. 評価関数を  $\varphi(x, y)$  とおく
- 評価基準  $\varphi(x, y)$  が満たしてほしい条件

$\varphi(x, y)$  は  $m/n = y/x$  のとき最小

$\varphi(x, y) \leq \varphi(x, y+\Delta), \varphi(x, y) \leq \varphi(x-\Delta, y)$  if  $m/n < y/x$

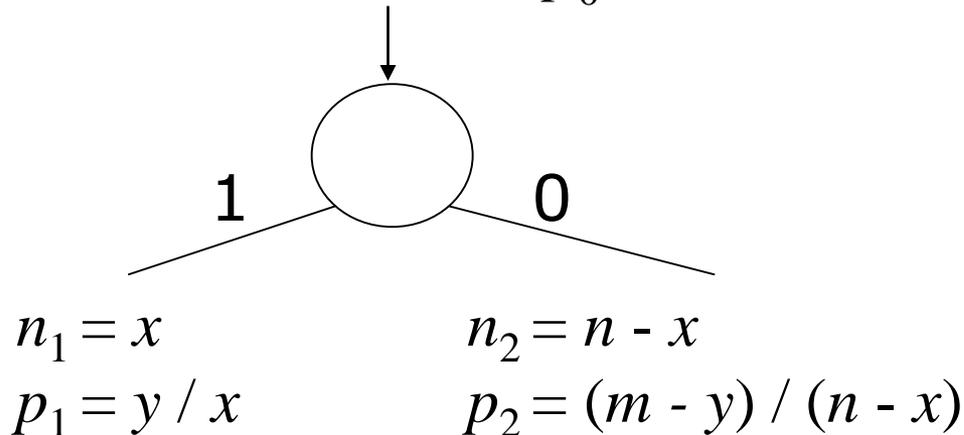
$\varphi(x, y) \leq \varphi(x, y-\Delta), \varphi(x, y) \leq \varphi(x+\Delta, y)$  if  $m/n > y/x$

$\Delta > 0$

# 評価値:エントロピーゲイン

正レコードの割合	$p = \#P / (\#P + \#N)$
負レコードの割合	$1-p = \#N / (\#P + \#N)$
エントロピー	$\text{ent}(p) = -p \log_2 p - (1-p) \log_2 (1-p)$

レコード数  $n$   
正レコードの割合  $p_0 = m / n$



Entropy Gain

$$\text{Ent}(x, y) = \text{ent}(p_0) - (n_1/n) \text{ent}(p_1) - (n_2/n) \text{ent}(p_2)$$

# T1を選んだときのエントロピーゲイン

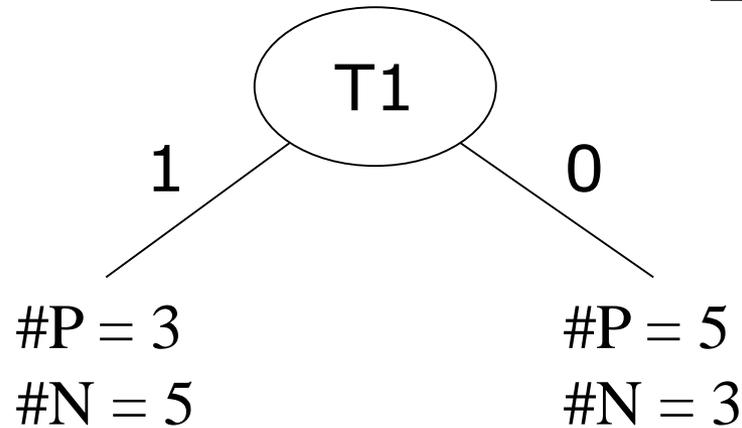
$$\#P = 8$$

$$\#N = 8$$

$$\text{ent}(8/16)$$

$$= 2 \left( - (1/2) \log_2(1/2) \right)$$

$$= 1$$



$$\text{ent}(3/8)$$

$$= - (3/8) \log_2(3/8) - (5/8) \log_2(5/8)$$

$$= 0.95444$$

$$\text{ent}(5/8) = \text{ent}(3/8)$$

$$\text{Entropy Gain} = \text{ent}(8/16) - (8/16)\text{ent}(3/8) - (8/16)\text{ent}(5/8)$$

$$= 1 - 0.95444$$

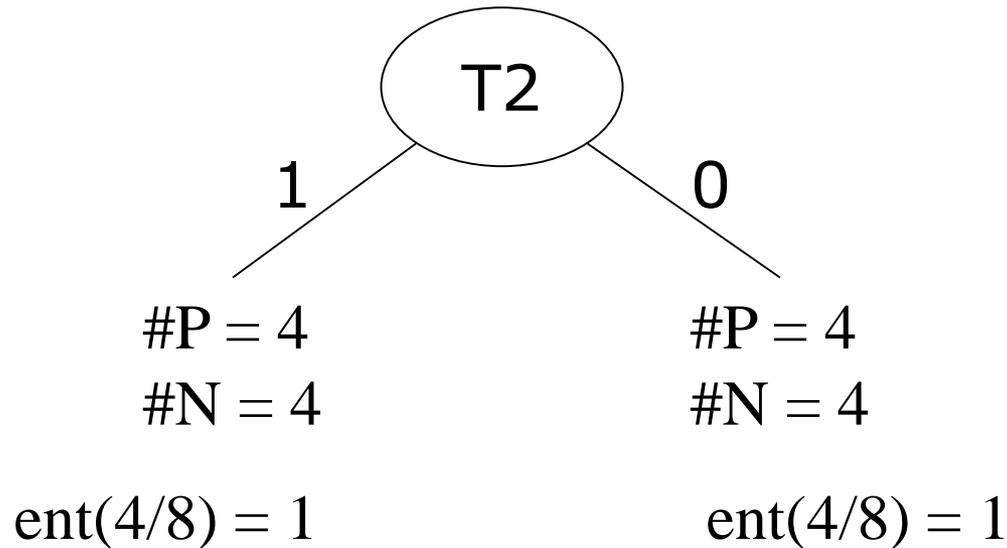
$$= 0.04556$$

# T2を選んだときのエントロピーゲイン

$$\#P = 8$$

$$\#N = 8$$

$$\text{ent}(8/16) = 1$$



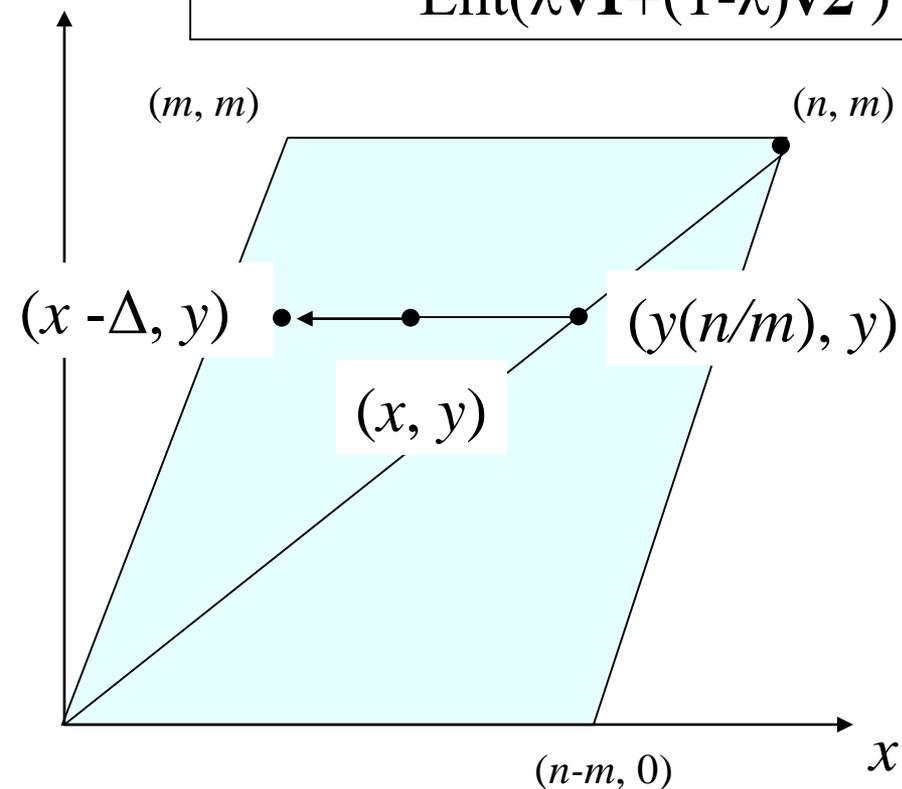
$$\begin{aligned} \text{Entropy Gain} &= \text{ent}(8/16) - (8/16)\text{ent}(4/8) - (8/16)\text{ent}(4/8) \\ &= 0 \end{aligned}$$

- $\text{Ent}(x, y)$  は  $m/n = y/x$  のとき最小
- $\text{Ent}(x, y)$  は凸関数 任意の点  $\mathbf{v1}, \mathbf{v2}$  と  $0 \leq \lambda \leq 1$  について

$$\text{Ent}(\lambda \mathbf{v1} + (1-\lambda) \mathbf{v2}) \leq \lambda \text{Ent}(\mathbf{v1}) + (1-\lambda) \text{Ent}(\mathbf{v2})$$

すると

$$\text{Ent}(\lambda \mathbf{v1} + (1-\lambda) \mathbf{v2}) \leq \max \{ \text{Ent}(\mathbf{v1}), \text{Ent}(\mathbf{v2}) \}$$



$$\text{Ent}(x, y)$$

$$\leq \max \{ \text{Ent}(x - \Delta, y), \text{Ent}(y(n/m), y) \}$$

$$\leq \text{Ent}(x - \Delta, y)$$

# 演習問題

T1	T2	T3	T4	目標属性
1	0	1	1	1
1	1	1	1	1
1	1	1	0	0
1	0	1	0	0
1	1	0	1	0
1	0	0	1	0
1	1	0	1	0
0	1	0	1	1
0	0	1	1	1
0	1	0	1	1
0	0	0	1	1
0	0	1	0	0
0	1	0	0	0
0	0	1	0	0

- 左の表のようなテスト-目標属性間の関係があるとき, エントロピーゲインに基づいて決定木を求めなさい