

クラス分類問題
Classification Problem
〔 教師つき学習
Supervised Learning 〕

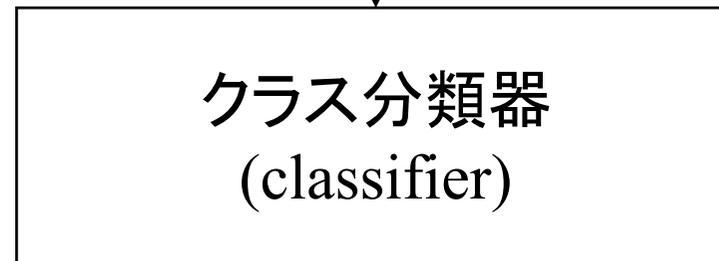
決定木

クラス分類器 (Classifier) の様々な手法

- 決定木
- Boosting
- Support Vector Machine
- k-nearest neighbor
- Naïve Bayes
- ニューラルネットワーク
- Hidden Markov Model

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



目標属性の値

T_i: 遺伝子発現量の増減, 遺伝子の局在, SNP, ハプロタイプ 等
 目標属性: ある表現型が出たか否か, 疾患か否か 等

テスト

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

If T1=1
Then

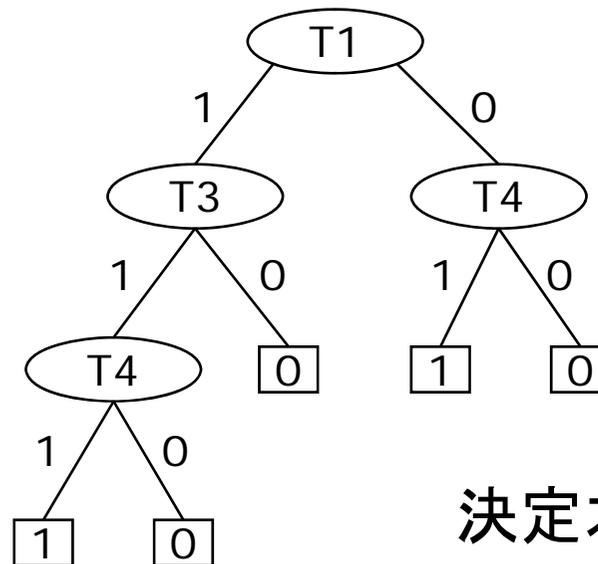
If T3=1
Then

If T4=1
Then 目標属性=1
Else 目標属性=0

Else
目標属性=0

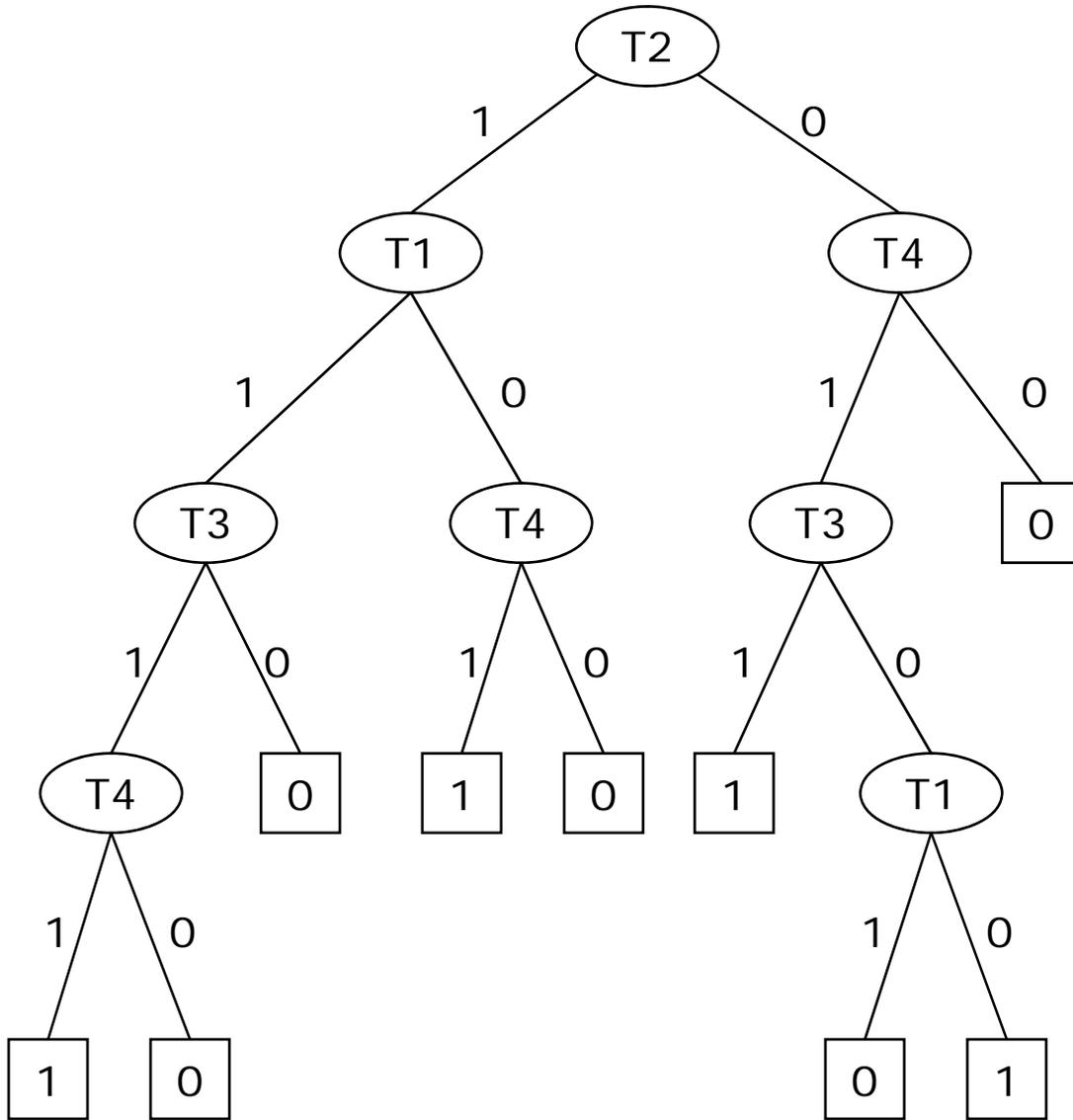
Else

If T4=1
Then 目標属性=1
Else 目標属性=0



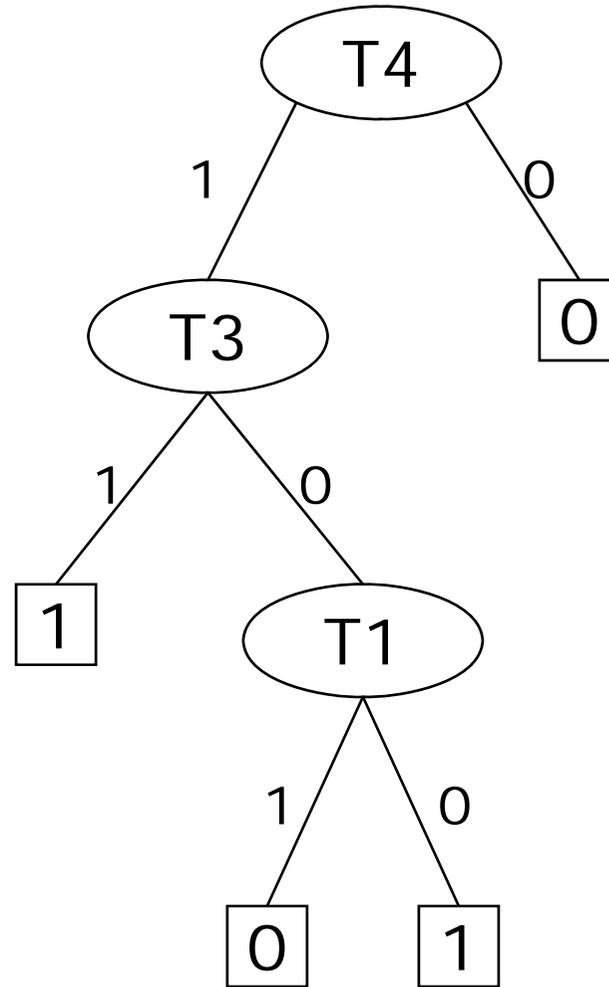
決定木(decision tree)

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



テストの選択の順番を変更すると決定木の大きさも変化

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

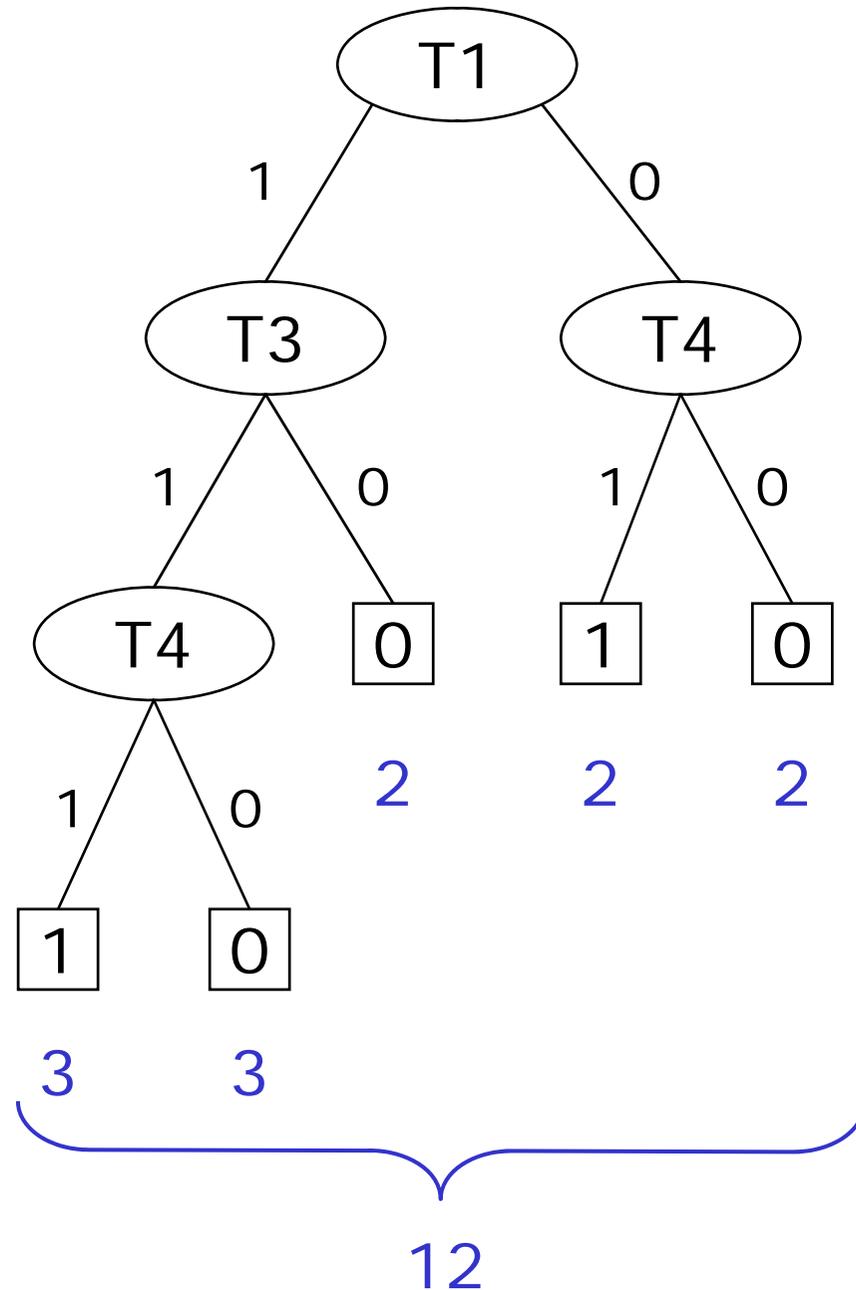


各葉ノードにおいて、到達するレコードがすべて正であるか、すべて負であるかのいずれかが成り立つ決定木を生成することを考える

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



NP問題

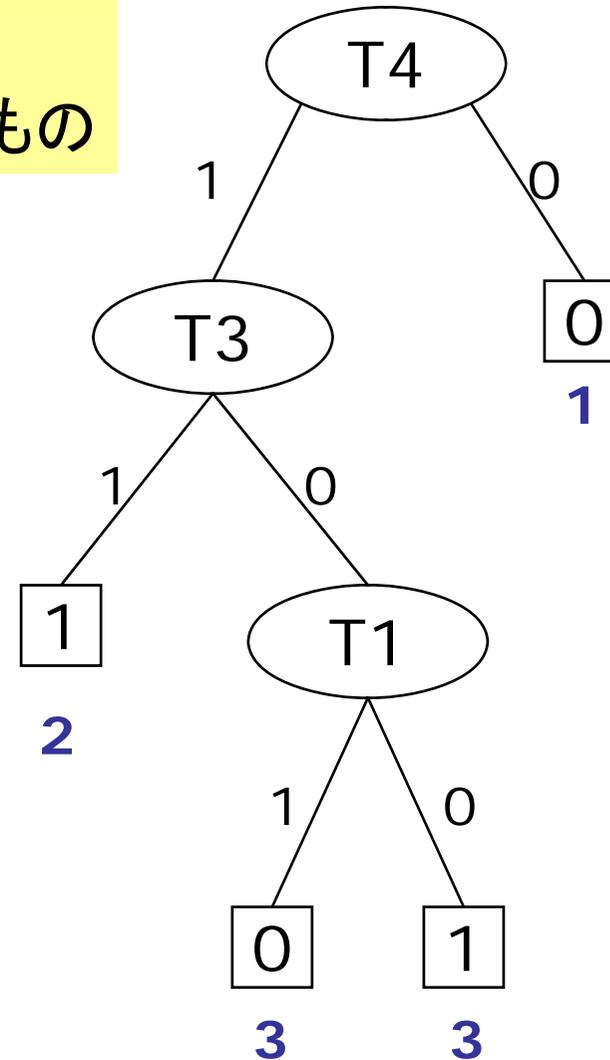
NP問題

- 正しいか否かを判定できる決定問題で、非決定性チューリングマシンが問題のサイズに関する多項式時間で計算可能な問題ある。
- 非決定性チューリングマシンとは、問題が正しいことを示す例に導いてくれる神託 (Oracle, ヒント) が与えられると、多項式時間で計算してくれる抽象的計算機。ヒントに依存して非決定的に動作する。
- 実際には都合のよいヒントを与えることは難しい。すべての解の可能性を並列的に生成してみて(このステップが手に負えない計算時間となる)、正しい例をさがす。
- 「コストが閾値以下の決定木が存在するか否か」はNP問題。
閾値以下のコストの決定木をヒントとして与えられれば、コストに比例する時間(多項式時間)で判定できるので。
- Bin Packing 問題もNP問題。
容器への分け方をヒントとして与えてもらえば、容器のサイズに納まっているか否かは積木の数に比例する時間で計算できるので。

コストが10以下の決定木は存在するか？

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

ヒントは
答えそのもの



$$2 + 3 + 3 + 1 = 9 < 10$$

NP完全問題

NP完全問題

- NP問題のひとつで、いかなるNP問題も多項式時間の手間をかければNP完全問題を解くことに帰着できる。
- すると、もしNP完全問題が多項式時間で解けるのであれば(しかし誰も肯定も否定もできていない計算機科学最大の問題)、どのようなNP問題も多項式時間で計算可能。
- 最初のNP完全問題は Cook により 1972年に提示された。
すべてのNP問題が多項式時間の手間で命題論理式が充足可能か否かを判定する問題(SATと呼ぶ)へ帰着できることを示した。
 - 非決定性チューリング機械の動作をSATにより表現できるように工夫。講義では証明する時間がない。興味がある人は以下の本の17~44ページを読んでください。
Michael R. Garey and David S. Johnson: *Computers and Intractability - A Guide to the Theory of NP-Completeness* (1979)
- SAT を多項式時間で他のNP問題 X を解くことに帰着できれば、その X もNP完全問題になる(いかなるNP問題も多項式時間で X の解決に帰着できるので)。
このようにして SAT と同格のNP完全問題が多数認定される。
- 例 Bin Packing, 「コストが閾値以下の決定木が存在するか否か」
- 多項式時間で解けそうにない問題が多いことを知る。

NP困難問題

NP困難問題

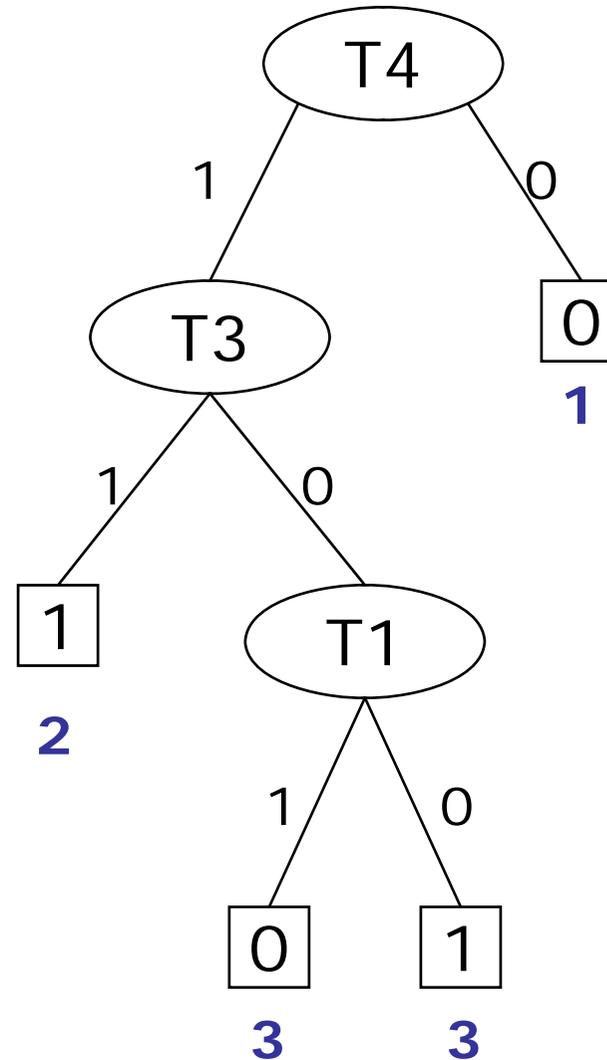
- NP問題のように正しいか否かを判定できる決定問題とは限らない。いかなるNP問題も多項式時間の手間を追加すれば、NP困難問題を解くことに帰着できるという性質を持つ。
- たとえば「コストが最小の決定木を求めよ」という最適化問題は、真偽を判定する決定問題でないが、NP困難問題である。
- なぜなら最小のコストが求めれば、NP完全問題「コストが閾値以下の決定木が存在するか否か」は直ぐに判定でき、任意のNP問題が多項式時間の手間をかけて、この最適化問題を解くことに帰着できるからである。
- NP完全問題とNP困難問題は混同されるので注意してほしい。特に、NP完全問題を最適化問題の形にしたときNP困難問題になる。

決定問題
最適化問題

コストが10以下の決定木は存在するか？
コスト最小の決定木は？

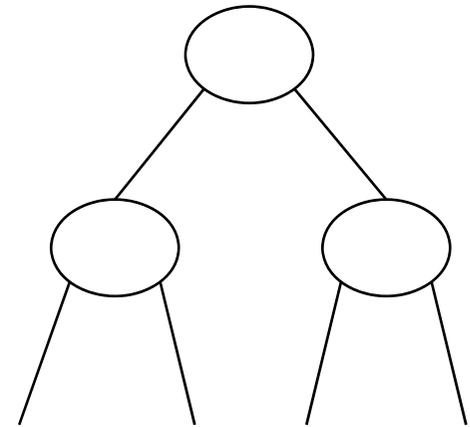
NP完全
NP困難

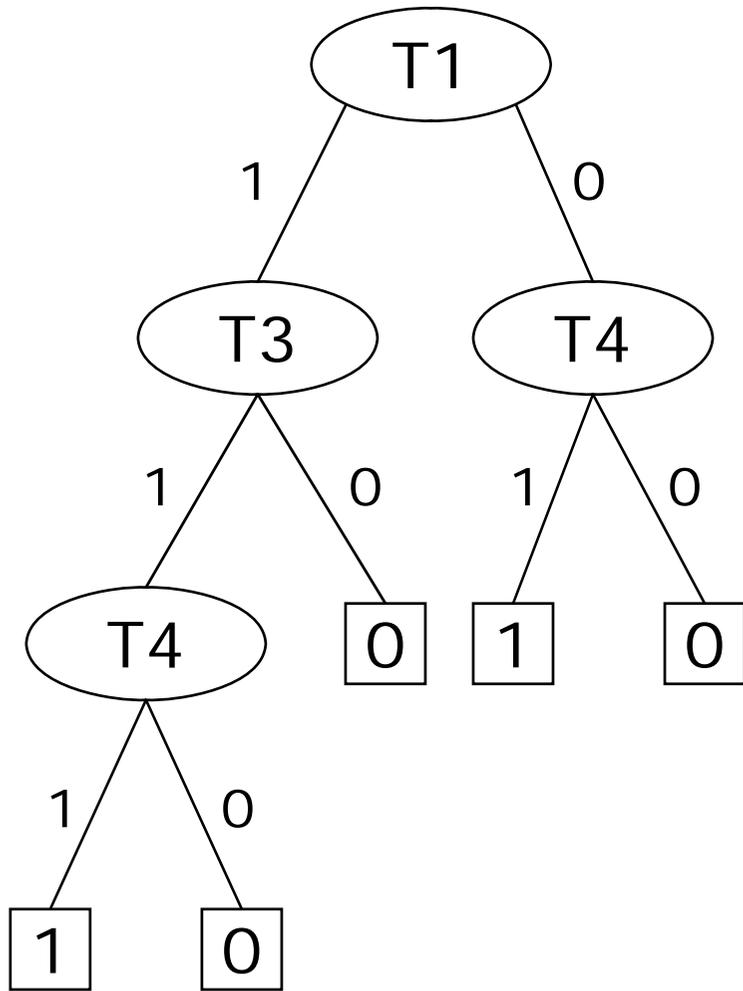
T1	T2	T3	T4	目標属性
1	0	1	1	1
1	0	1	1	1
1	1	1	1	1
0	0	1	1	1
0	1	0	1	1
0	1	0	1	1
0	1	0	1	1
0	0	0	1	1
1	1	0	1	0
1	0	0	1	0
1	1	0	1	0
0	0	1	0	0
0	0	1	0	0
1	1	1	0	0
1	0	1	0	0
0	1	0	0	0



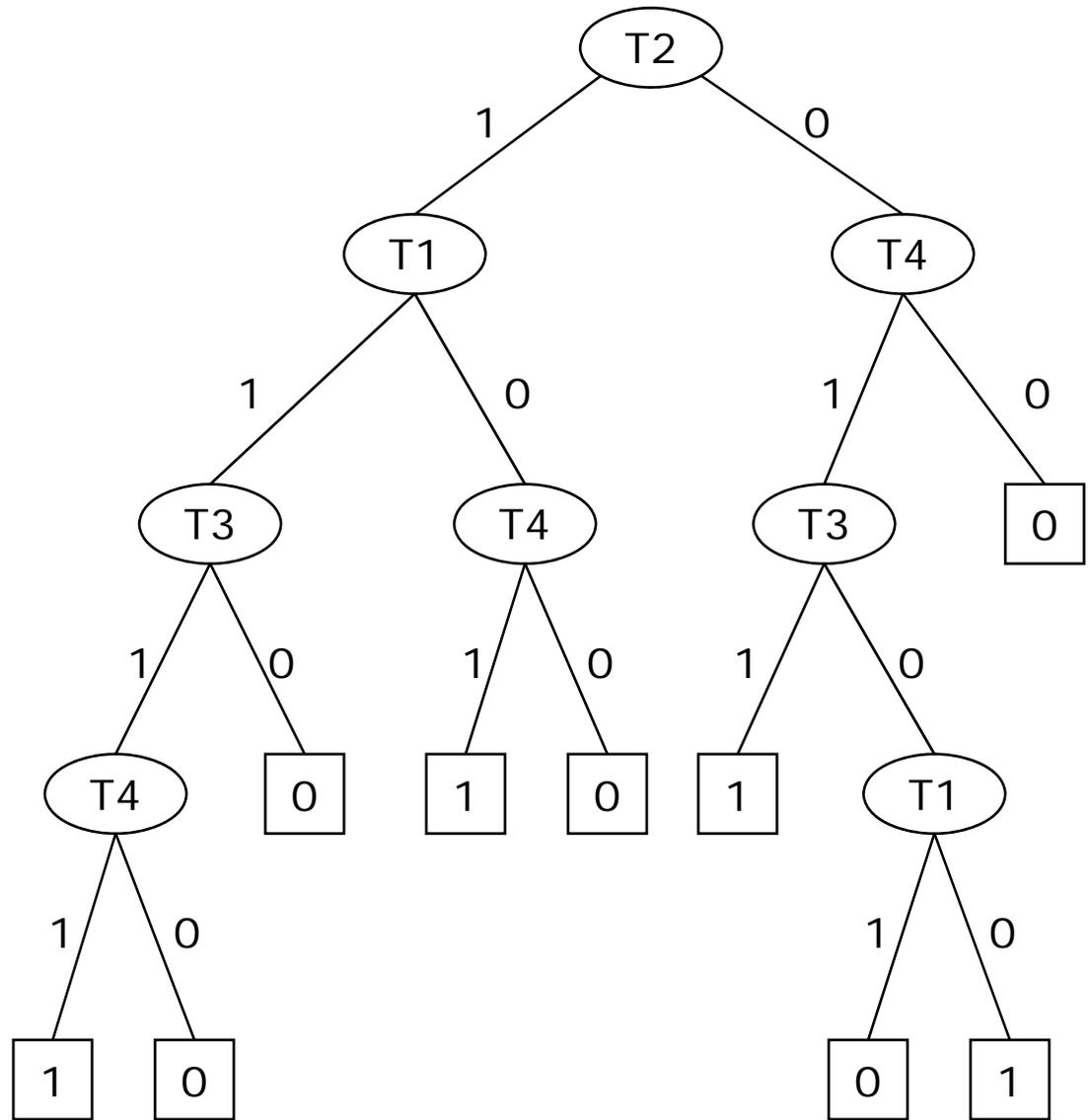
近似的解法 貪欲(greedy)アルゴリズム

- 「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





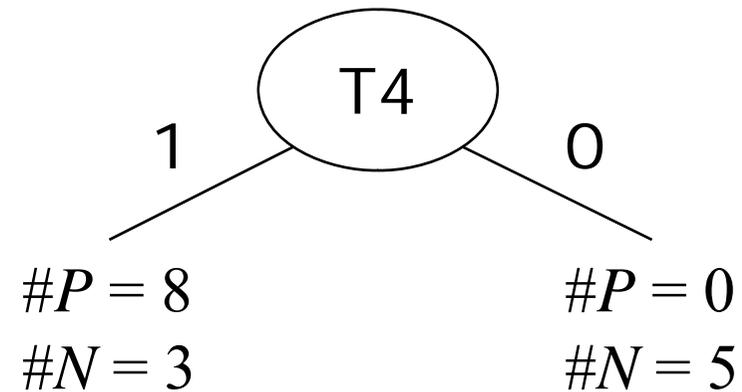
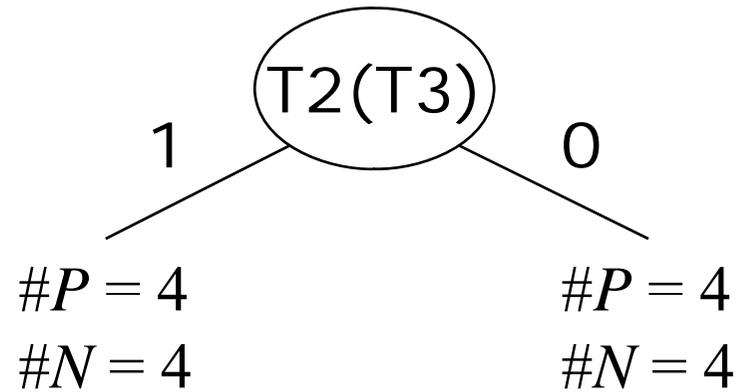
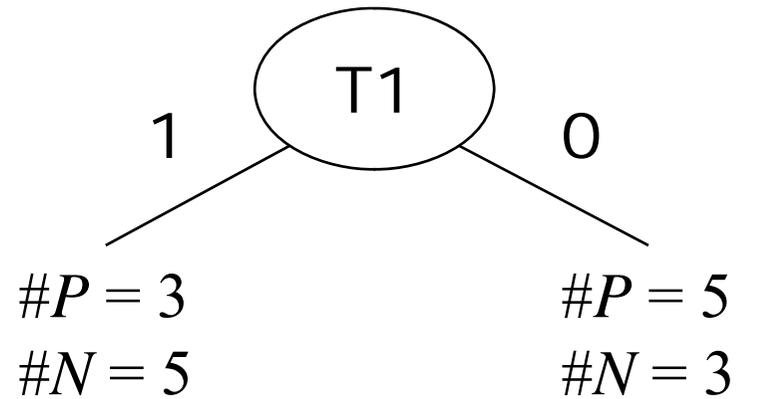
$3+3+2+2+2=12$



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

#P: 正レコードの数 #N: 負レコードの数

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



$$\begin{aligned} \#P + \#N &= n \\ \#P &= m \end{aligned}$$

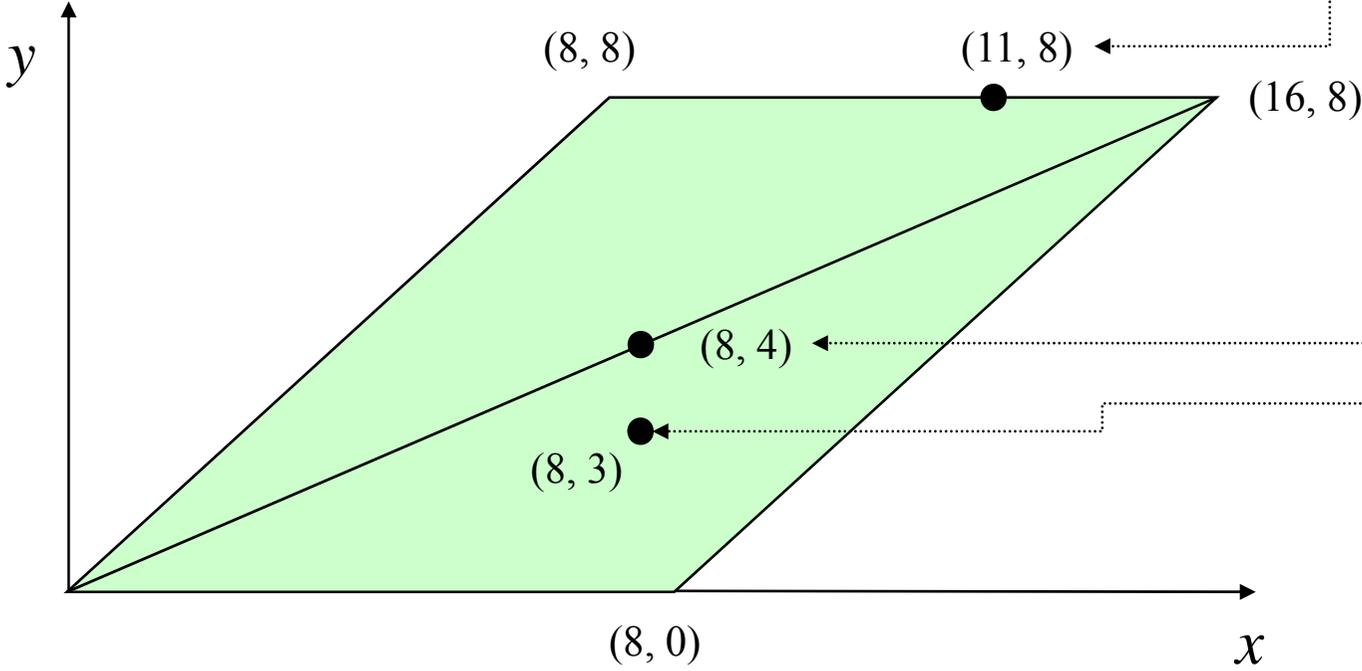
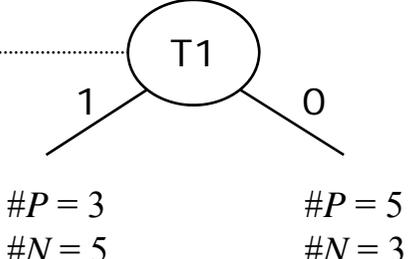
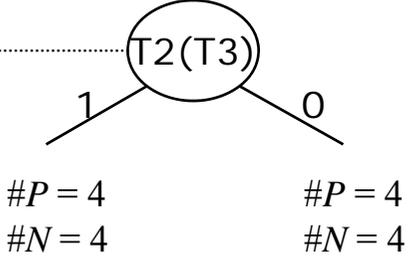
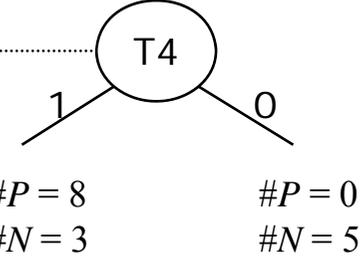
n, m は test の選択に依存しないので定数と考える。選択した test が決めるのは変数の x と y

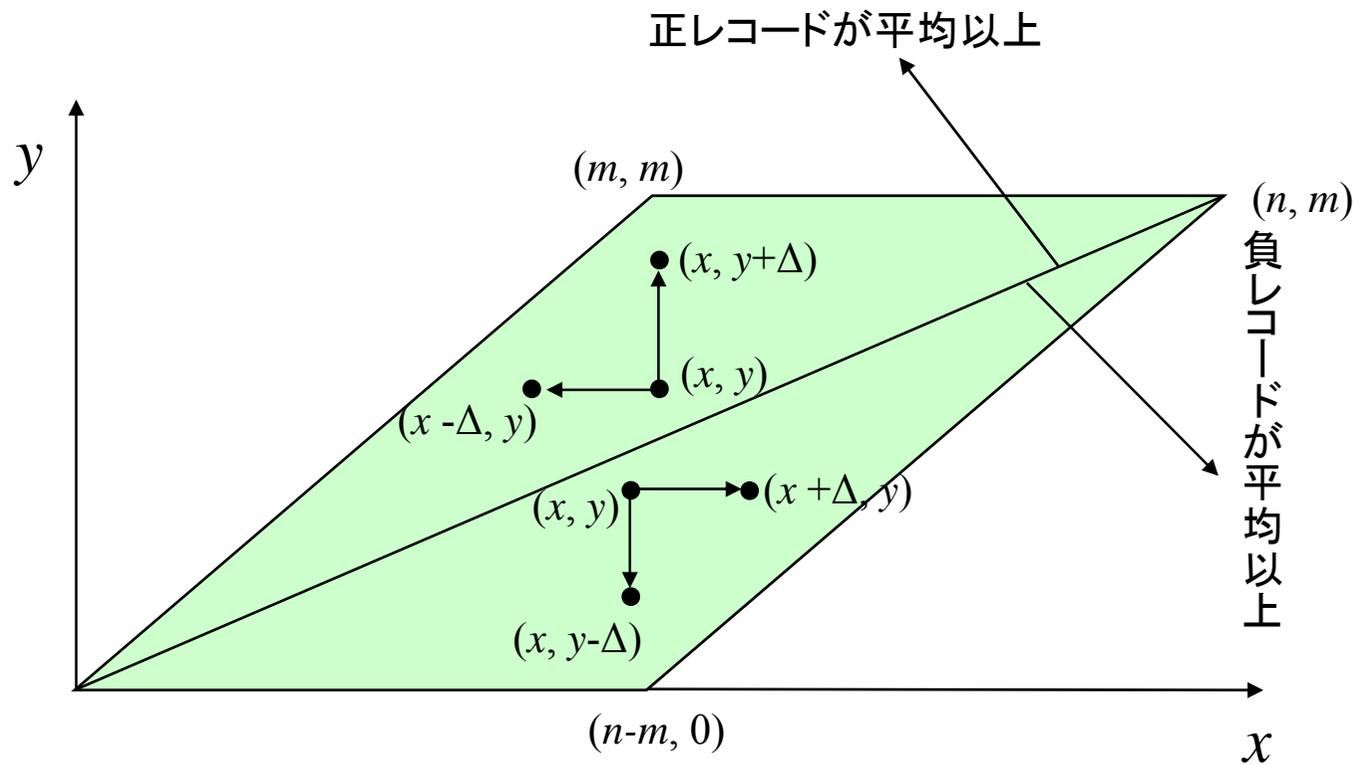
test

1 0

$$\begin{aligned} \#P + \#N &= x \\ \#P &= y \end{aligned}$$

$$\begin{aligned} \#P + \#N &= n - x \\ \#P &= m - y \end{aligned}$$





評価基準 $\varphi(x, y)$ が満たすべき条件

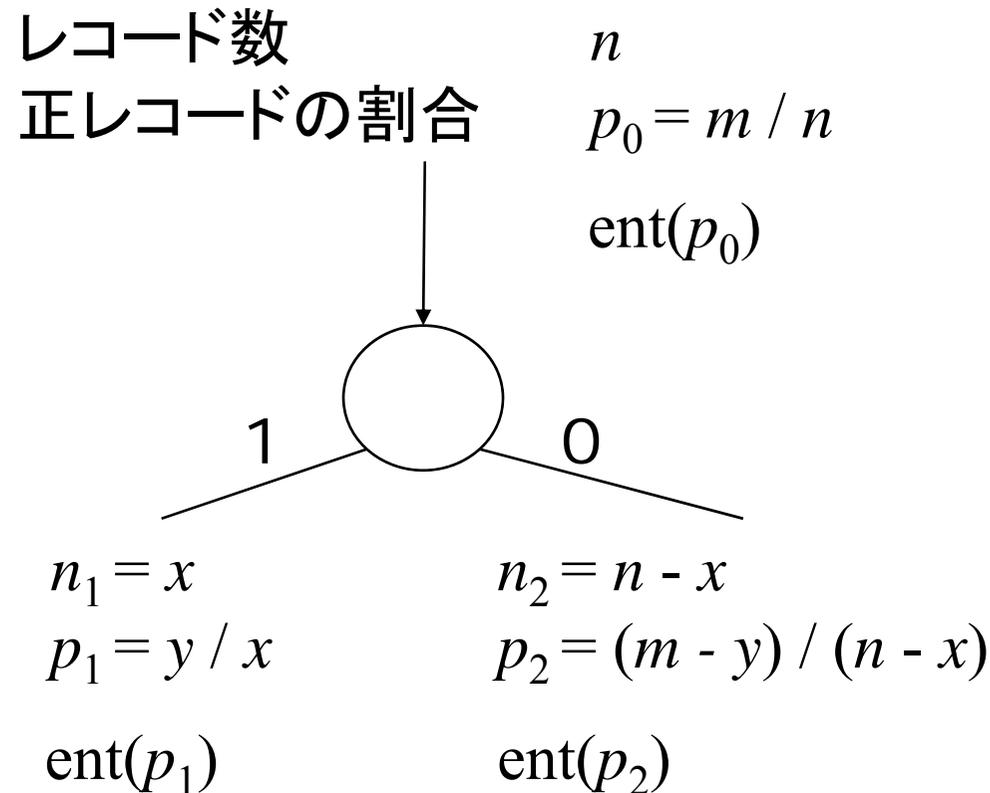
$\varphi(x, y)$ は $m/n = y/x$ のとき最小

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

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

$\Delta > 0$

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



Entropy Gain / Information Gain

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

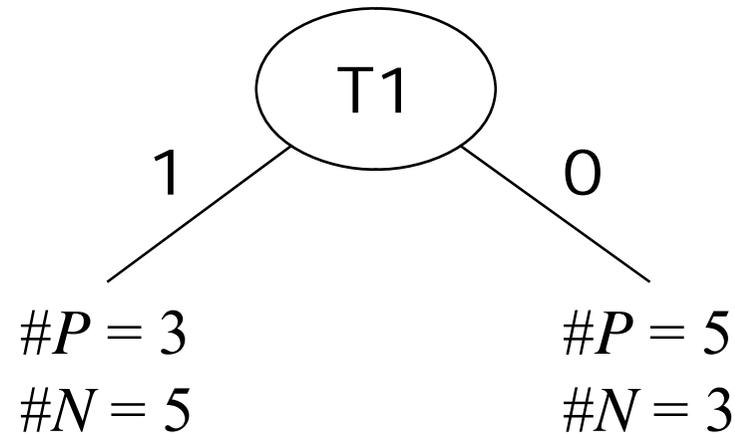
$$\#P = 8$$

$$\#N = 8$$

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

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

$$= 1$$



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

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

$$= 0.95444$$

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

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

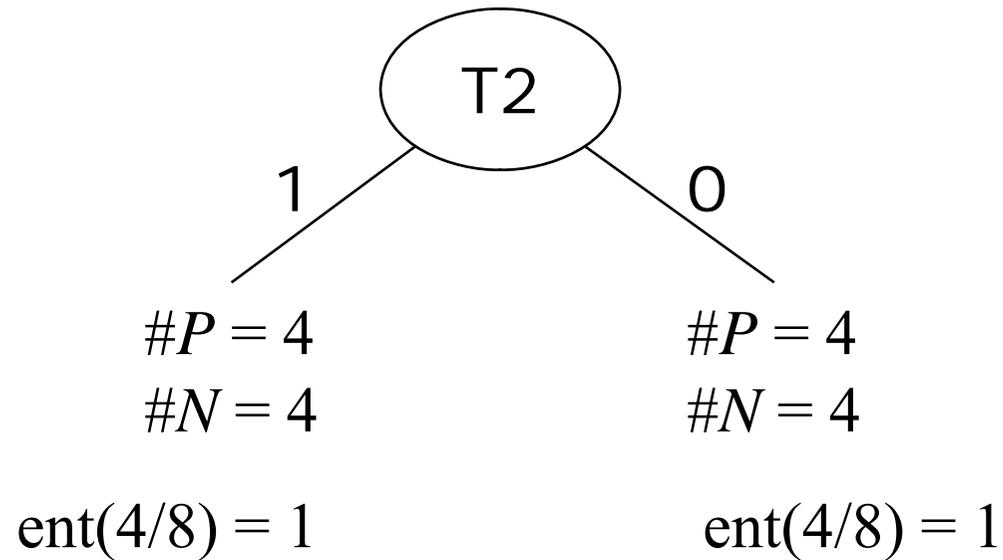
$$= 1 - 0.95444$$

$$= 0.04556$$

$$\#P = 8$$

$$\#N = 8$$

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



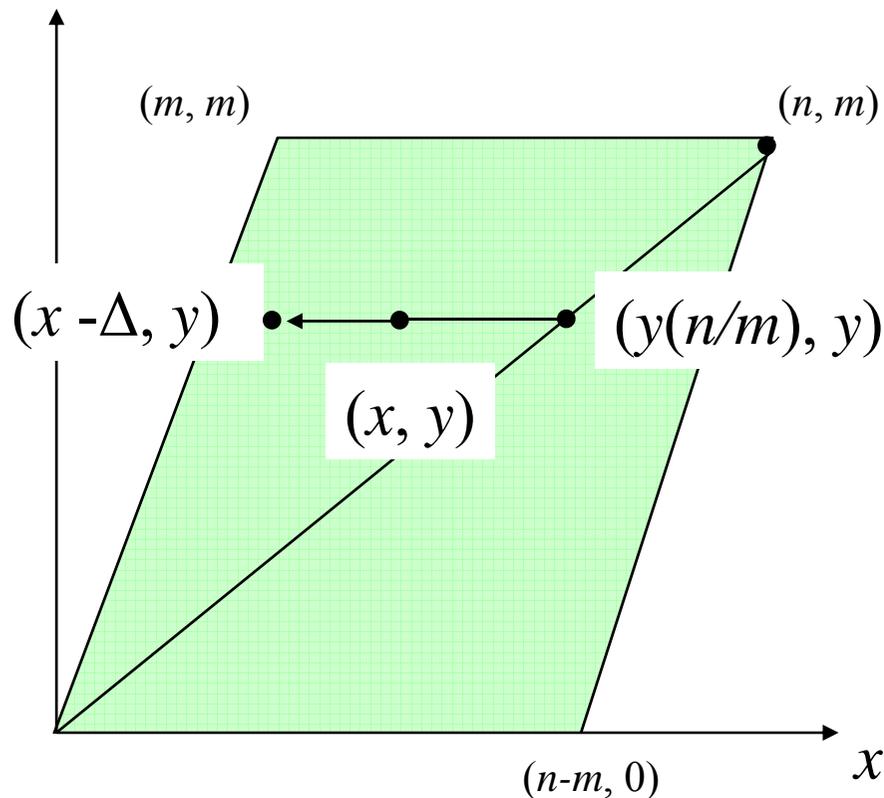
$$\begin{aligned} \text{Entropy Gain} &= \text{ent}(8/16) - (8/16)\text{ent}(4/8) - (8/16)\text{ent}(4/8) \\ &= 0 \end{aligned}$$

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

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



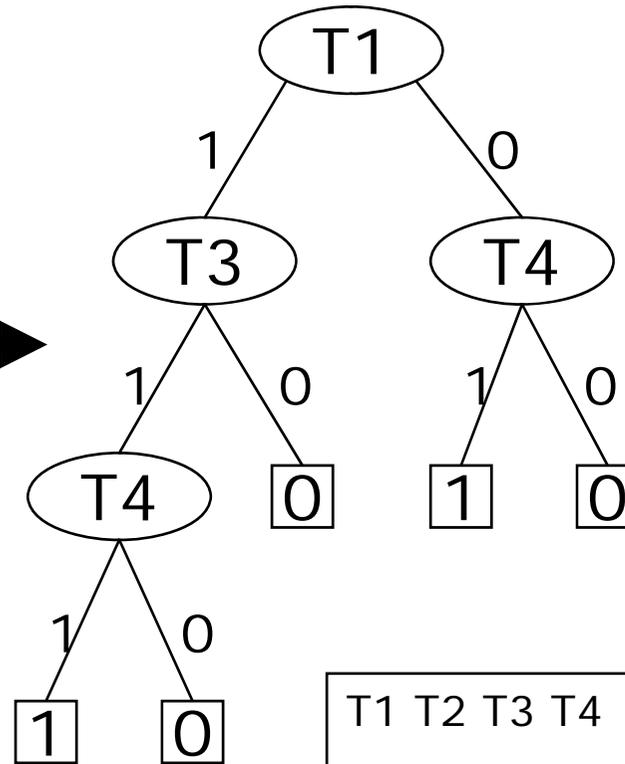
$$\begin{aligned} \text{Ent}(x,y) & \leq \max \{ \text{Ent}(x-\Delta, y), \text{Ent}(y(n/m), y) \} \\ & \leq \text{Ent}(x-\Delta, y) \end{aligned}$$

最小値

未知のテストレコードに対する正答率

訓練(training)レコード

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
<u>1</u>	<u>1</u>	<u>0</u>	<u>1</u>	<u>0</u>
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



正答率 = 4/5

T1	T2	T3	T4	目標属性	決定木の予測
1	1	1	1	1	1
1	0	0	1	0	0
<u>1</u>	<u>1</u>	<u>0</u>	<u>1</u>	<u>1</u>	<u>0</u>
0	1	1	1	1	1
0	0	0	0	0	0

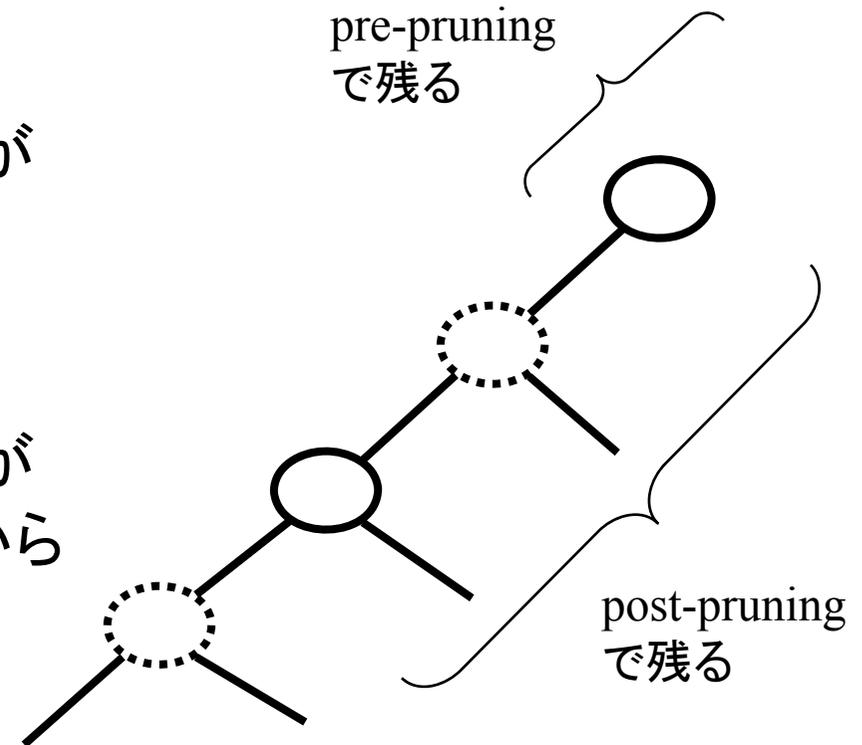
テストレコードが訓練レコードと異なる場合

Overfitting 問題

- 決定木を十分に大きくすれば、
訓練レコードに関しては正答率を改善できる
- しかし巨大な木のテストレコードに対する正答率は
必ずしも高くない
- 訓練レコードに対して過剰に適合 (overfitting) した状態

Overfitting の回避方法1

- pre-pruning 法
決定木を構成中に Entropy Gain が小さいテスト属性は生成しない
- post pruning 法
決定木を構成後に Entropy Gain が小さいテスト属性を木の葉ノードからボトムアップに除く
- 実験的には Overfitting をある程度回避できると言われている



○ Entropy Gain が
大きいテスト属性

⊙ Entropy Gain が
小さいテスト属性

このような例題を作ってみてください

Overfitting の回避方法2

小さい決定木を生成するアプローチ

Entropy Gain を最大化するテストを効率的に計算できるか？

- 論理積でレコード集合を分割

$$(T_1 = 1) \wedge (T_2 = 1) \wedge \dots \wedge (T_n = 1)$$

NP困難 branch-and-bound 法

- テスト属性 T の領域に3つ以上の異なる離散値がある場合

例 血液型 $\in \{A, O\}$ ($\subseteq \{A, B, AB, O\}$)

k 個の離散値のとき $O(k)$ 、つまり k に線形比例する計算時間

- テスト属性が数値の場合

例 $21 \leq \text{年齢} \leq 44$

多項式時間

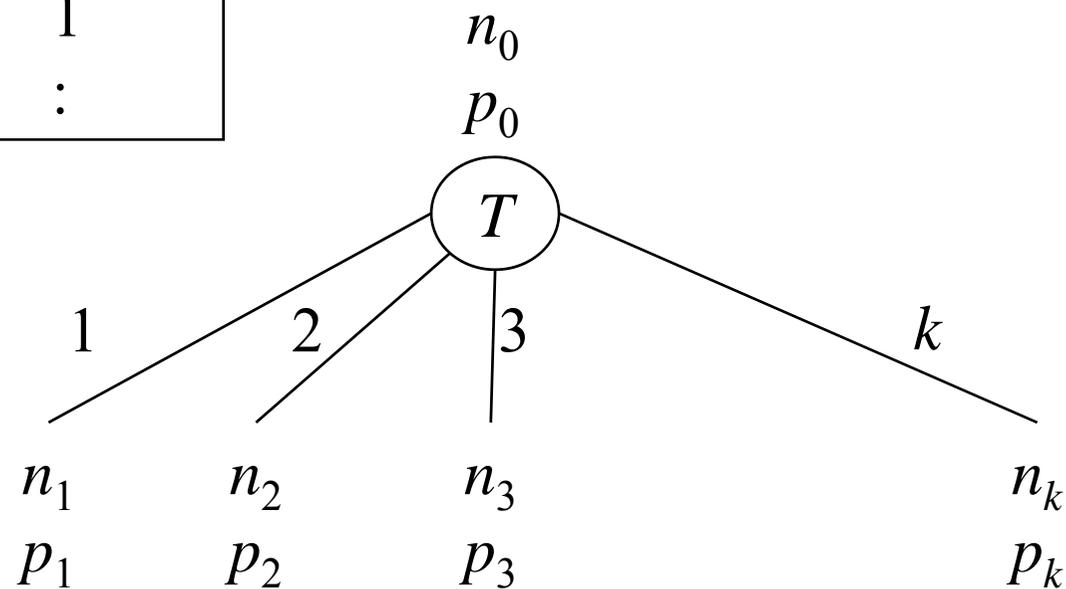
テスト属性 T の領域に3つ以上の異なる離散値がある場合

- T の領域を $\{1, 2, \dots, k\}$

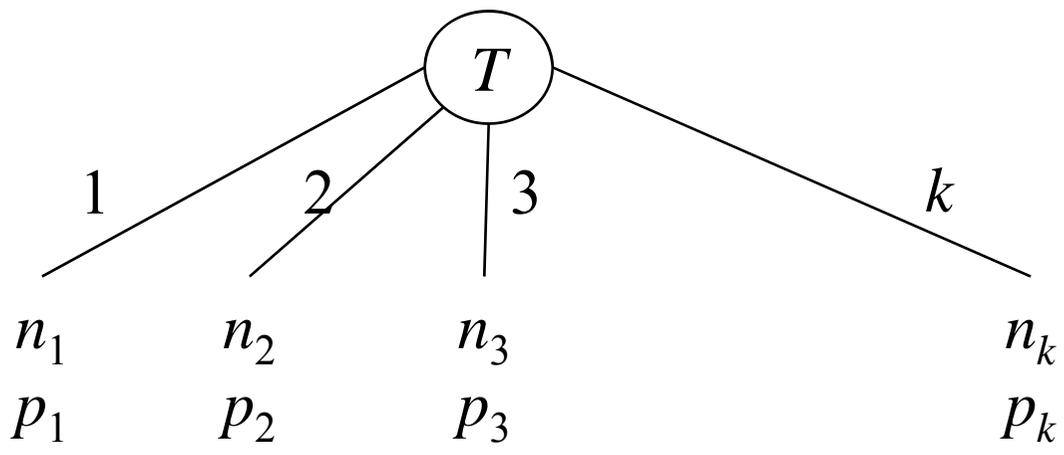
T ...	目標属性A
2	1
3	0
1	1
⋮	⋮

各属性値ごとに分割する方式

例えば $k=1000$ のときは
巨大な木になる



$$\text{ent}(p_0) - \sum (n_i / n_0) \text{ent}(p_i)$$

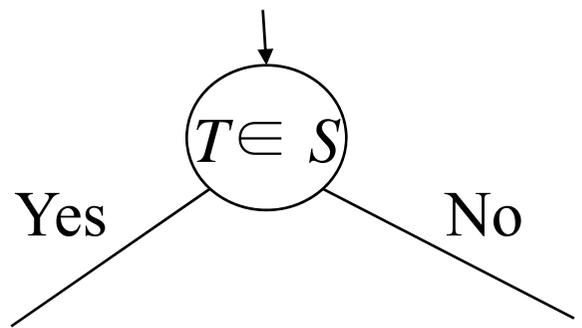


一般性を失うことなく
 $p_1 \geq p_2 \geq \dots \geq p_k$
 と仮定

$\#P + \#N = n$
 $\#P = m$

属性値の部分集合で分割する方式

$\{1, 2, \dots, k\}$ の部分集合 S



データ数
 $\Sigma\{n_i \mid i \in S\} = x$
 正のデータ数
 $\Sigma\{p_i n_i \mid i \in S\} = y$

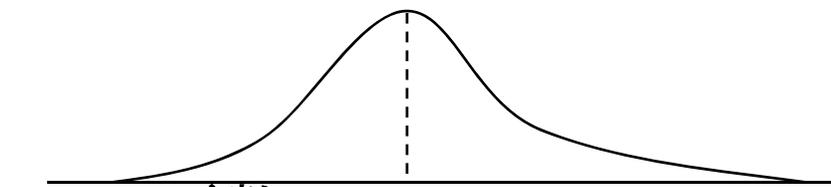
$\Sigma\{n_i \mid i \notin S\}$
 $\{p_i n_i \mid i \notin S\}$

定理
 ある j が存在して
 $S = \{1, \dots, j\}$ または $\{j, \dots, k\}$
 が Entropy Gain を最大化

遺伝子の状態

目標属性が数値の場合の拡張

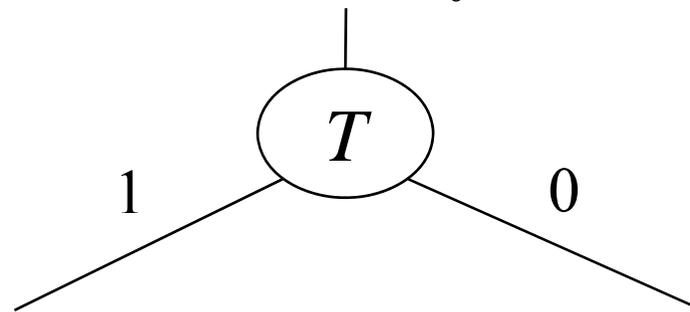
T_1	T_2	...	目標属性 A (血圧)
1	0		125
0	1		145
1	1		113
⋮	⋮	⋮	



レコード数 = n

目標属性値の総和 $\sum t[A] = m$

平均値 = $m / n = p_0$



目標属性の例としては
癌細胞を1/10にするのに
必要な放射線量や
抗がん剤の量など

レコード数 = x

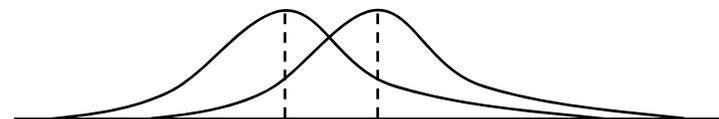
$\sum \{t[A] \mid t[T]=1\} = y$

平均値 = $y / x = p_1$

レコード数 = $n - x$

平均値 = $(m - y) / (n - x)$

= p_2



クラス間分散 $\text{Var}(x, y) = x (p_1 - p_0)^2 + (n - x) (p_2 - p_0)^2$
 クラス間分散を最大化する T_1 をみつきたい

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

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

- Entropy Gain を最大化する方法がクラス間分散を最大化するのに使える
- 決定木のように、レコード分割を繰り返した木構造を回帰木 (Regression Tree) と呼ぶ

付録

「コストが閾値以下の決定木が存在するか否か」
はNP完全問題

「コストが閾値以下の決定木が存在するか否か」がNP完全問題であることを証明する。

L. Hyafil and R. L. Rivest: “Constructing Optimal Binary Decision Trees is NP-complete,”
Information Processing Letters, Vol. 5, No. 1, May 1976

証明のロードマップ

SAT は、多項式時間の手間をかければ EXACT COVER by 3-SETS というNP完全問題を解くことに帰着できる。

この EXACT COVER by 3-SETS を多項式時間の変形により、「コストが閾値以下の決定木が存在するか否か」を解くことに帰着できることを示す。

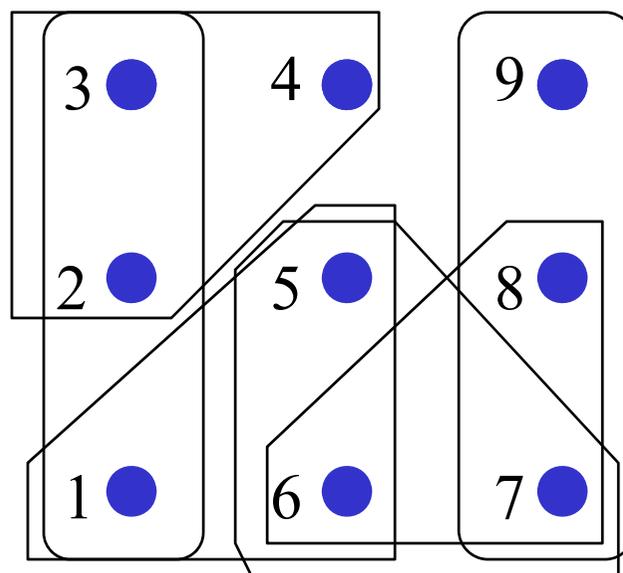
EXACT COVER BY 3-SETS 問題 (NP完全問題)

3の倍数の個数だけ元を含む有限集合 X と、
3個の元を含む X の部分集合の集まり $S = \{T_1, T_2, \dots\}$ が与えられたとき、次の条件を満たす部分集合 $S_1 \subseteq S$ は存在するか？

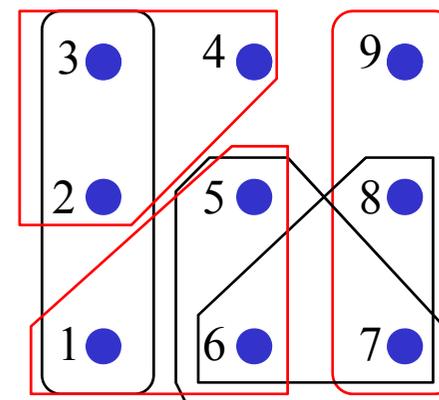
- $\cup \{T \mid T \in S_1\} = X$
 X の任意の元は、 S_1 のどれか一つに含まれる
- S_1 の元は互いに交わらない
 任意の $i \neq j$ について $T_i \cap T_j = \emptyset$

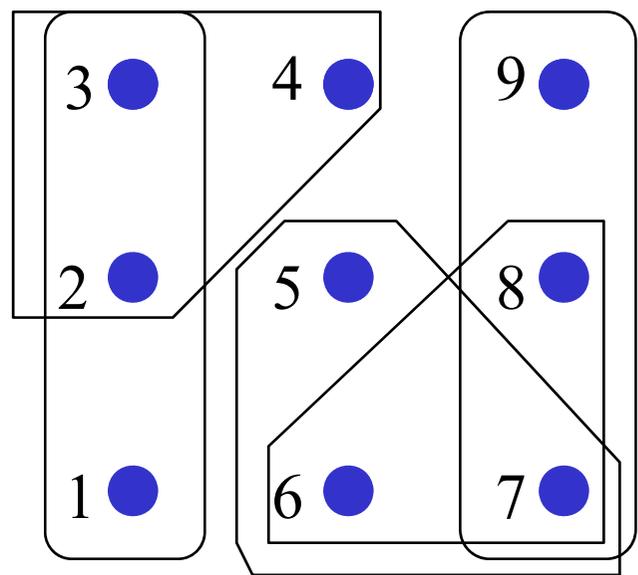
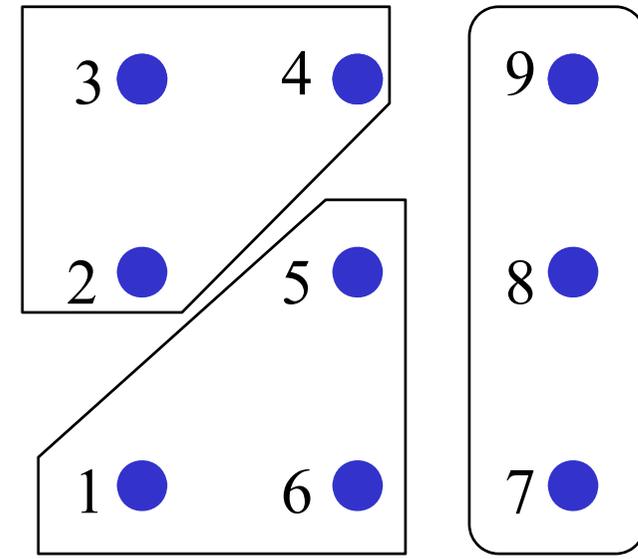
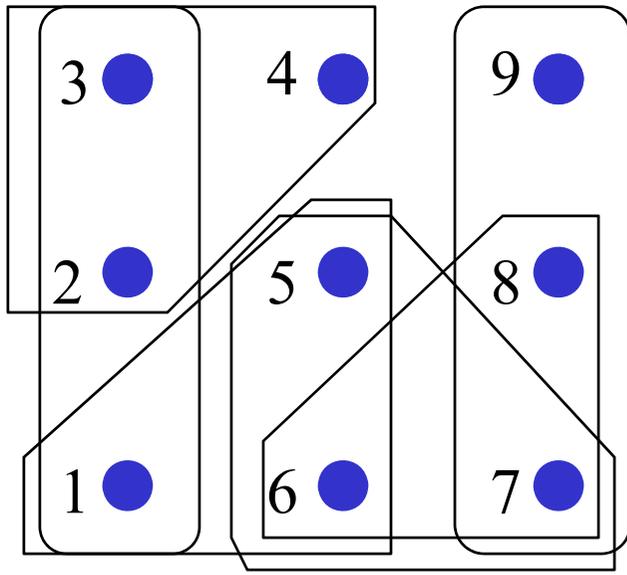
例

$X = \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$
 $S = \{\{1, 2, 3\}, \{2, 3, 4\},$
 $\{1, 5, 6\}, \{5, 6, 7\},$
 $\{6, 7, 8\}, \{7, 8, 9\}\}$



ヒント





EXACT COVER BY 3-SETS
解がない

テスト属性

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

属性(attribute) フィーチャー

レコード データ タップル

レコード t の属性 A の値を $t[A]$ と表記

目標値 目標属性の値から一つ選択. 例えば 1

目標属性が A 、目標値が 1 のとき、

- $t[A]=1$ となるレコード t を **正(positive)**
- $t[A]\neq 1$ となるレコード t を **負(negative)**

EXACT COVER BY 3-SETS

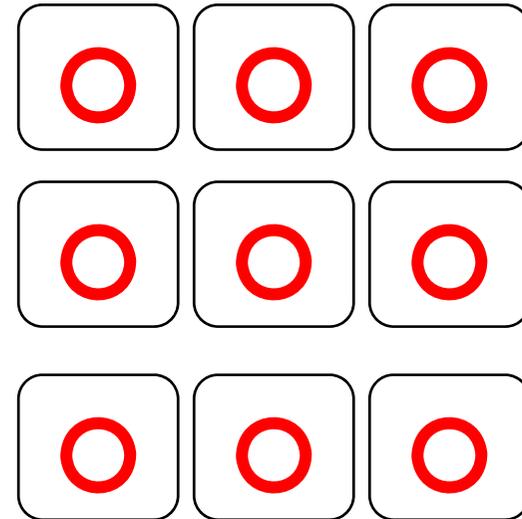
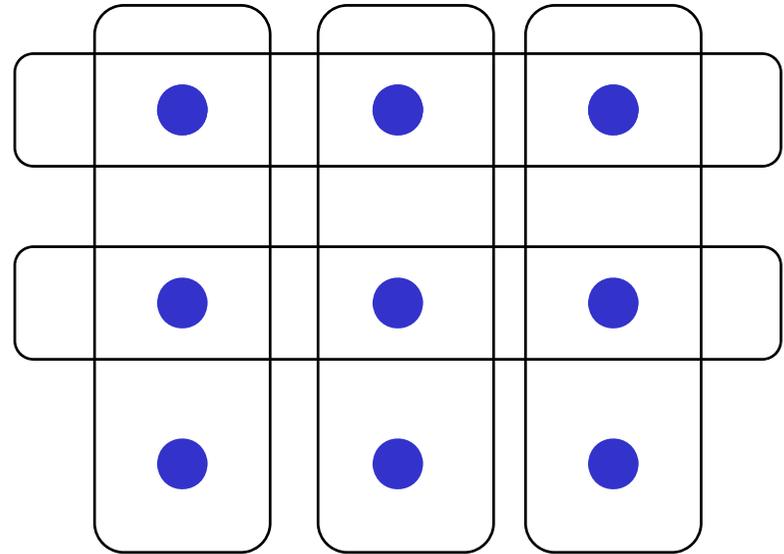
$X, S = \{T1, T2, \dots\}$

● 正レコード X の元

□ テスト属性

○ 負レコード Y の元

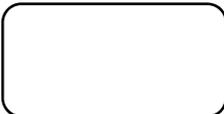
決定木を作成するために
負のレコードを追加する

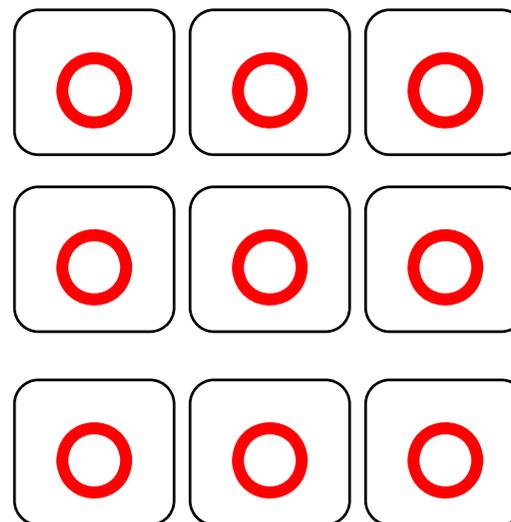
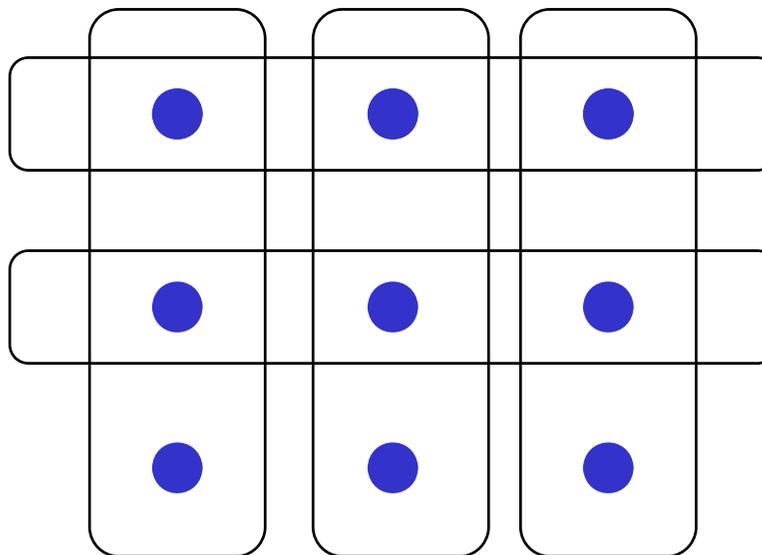


例

● 正レコード $t[A]=1$ X の元

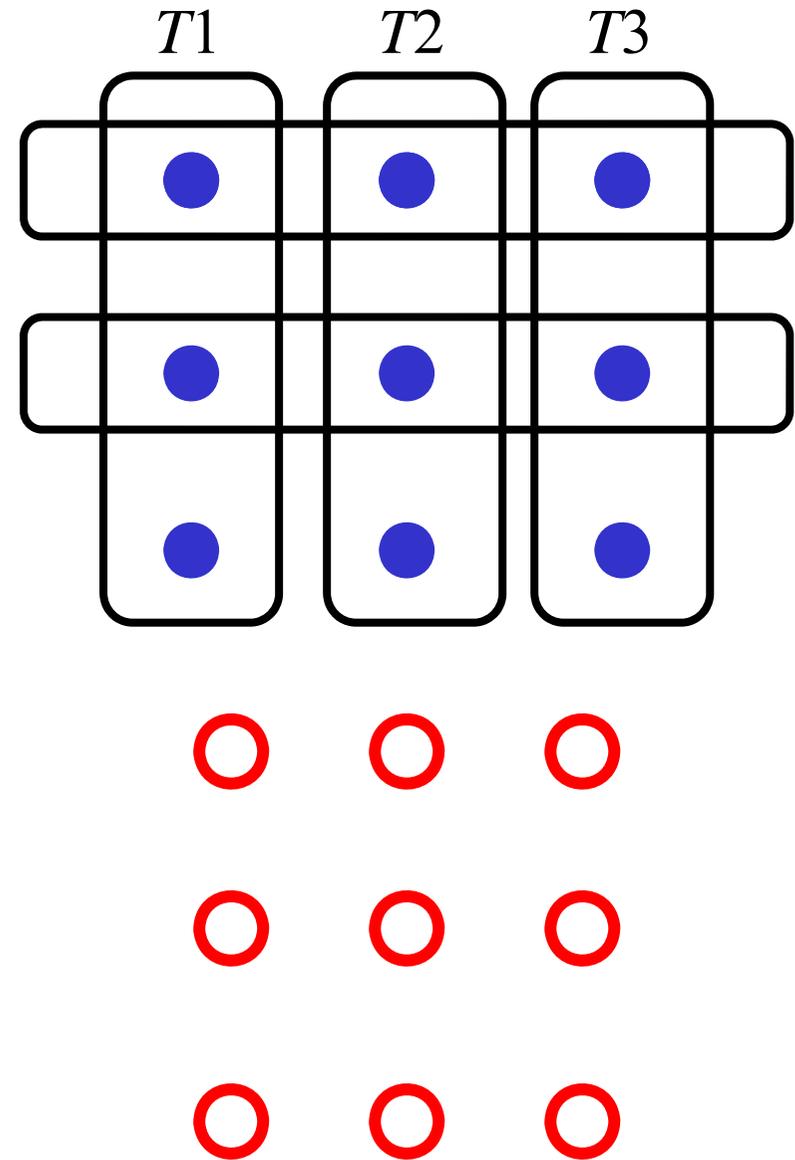
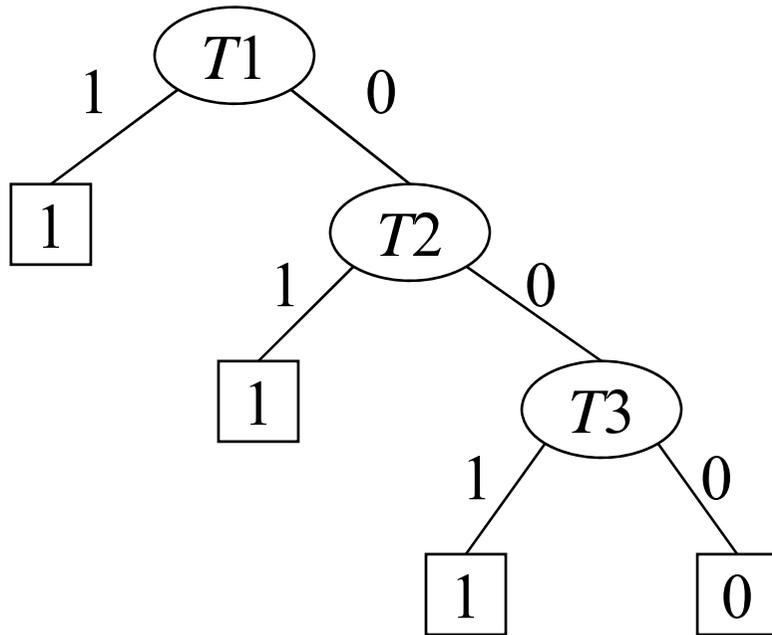
○ 負レコード $t[A]=0$ Y の元

 テスト属性



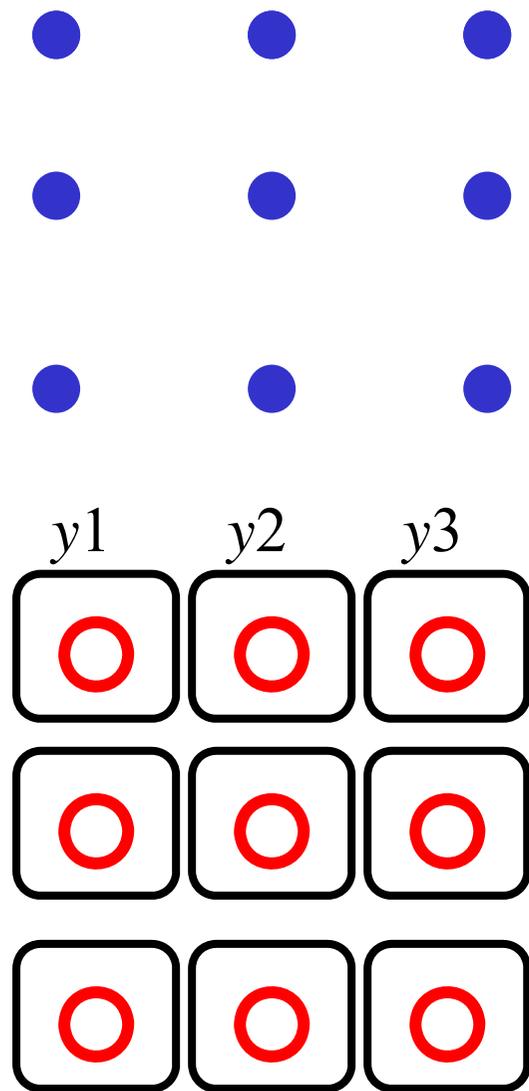
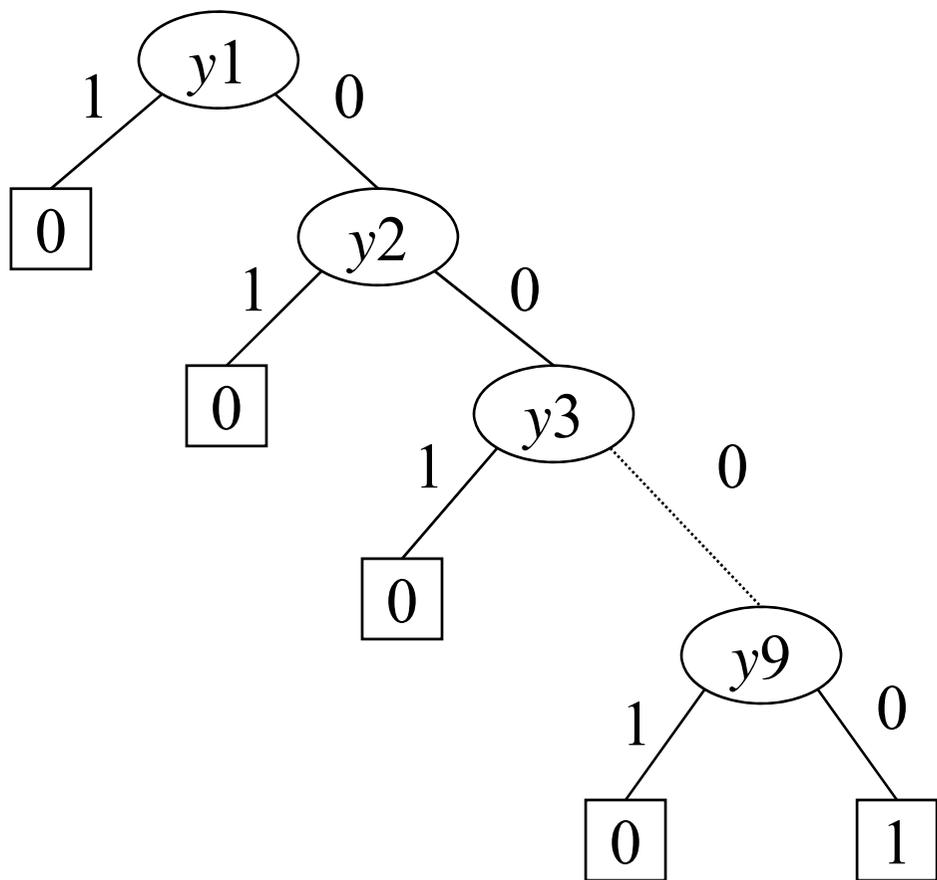
正レコードと負レコードを
決定木で分離する戦略1

$S = \{T1, T2, \dots\}$ のテスト属性
を使って正レコードを包含する



正レコードと負レコードを
決定木で分離する戦略2

$Y = \{y_1, y_2, \dots\}$ のテスト属性を
すべて使って負レコードを包含する



正レコードと負レコードを 分離する決定木の構成法

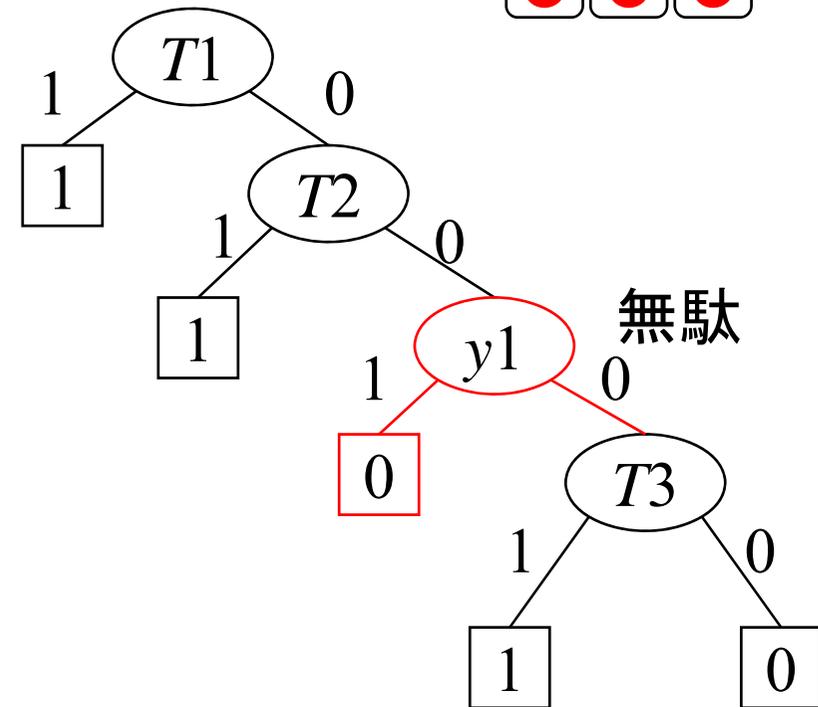
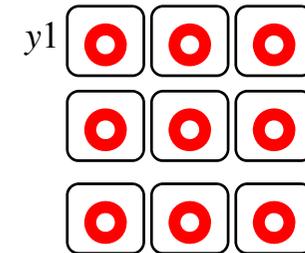
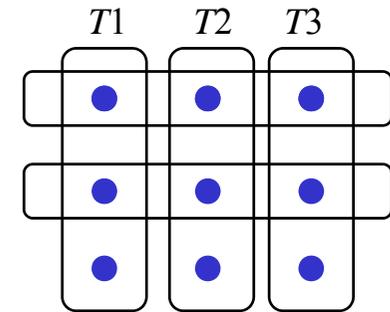
戦略1

正レコードをすべて包含する
 S の部分集合をテスト属性として使う
(Y のテスト属性を使うのは無駄)

戦略2

負レコードを包含する Y の
テスト属性をすべて使う
(S のテスト属性を使うのは無駄)

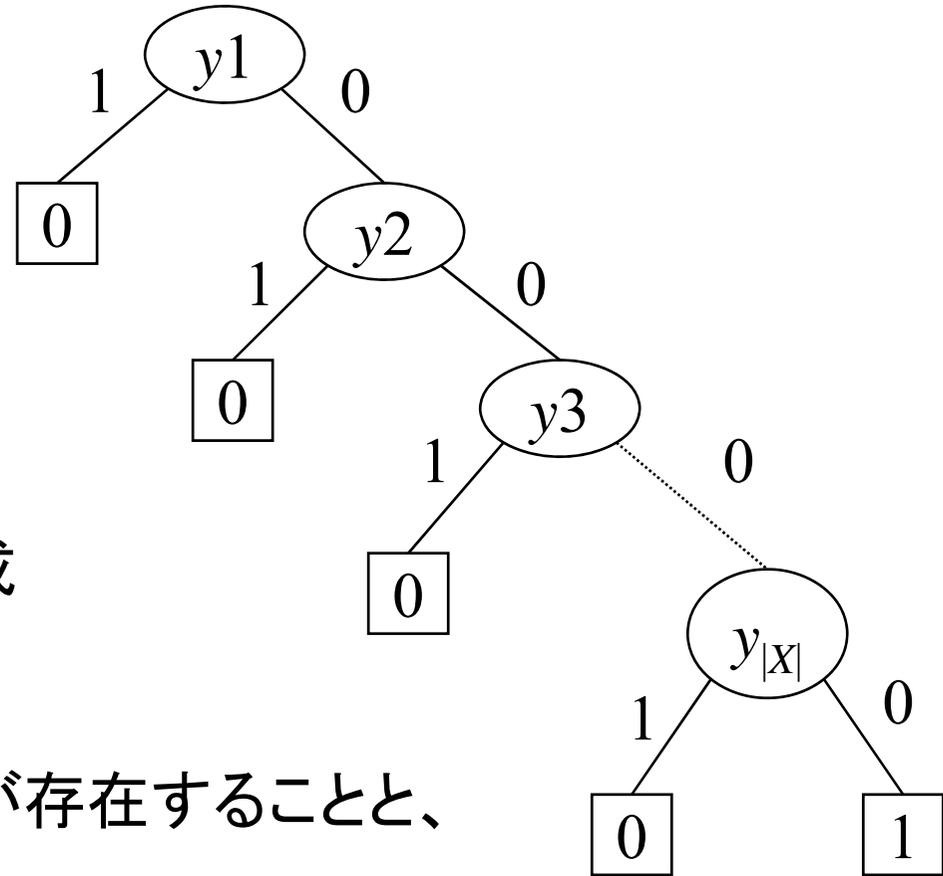
したがって、戦略1もしくは戦略2が
コスト最小の木を生成する



戦略2のコスト

$$1+2+\dots+|X|+|X|$$

は戦略1のコストより大きい



戦略1がコスト最小の木を生成

EXACT COVER BY 3-SETSが存在すること、
戦略1よりコストが

$$1+2+\dots+|X|/3 + |X|/3 \quad (= w \text{ とおく})$$

となる決定木が存在すること

(コストが w 以下の決定木が存在すること) は同値.

したがって、EXACT COVER by 3-SETS を多項式時間の変形により、「コストが閾値以下の決定木が存在するか否か」を解くことに帰着できる (証明終)

付録 以下の低利の証明のスケッチ

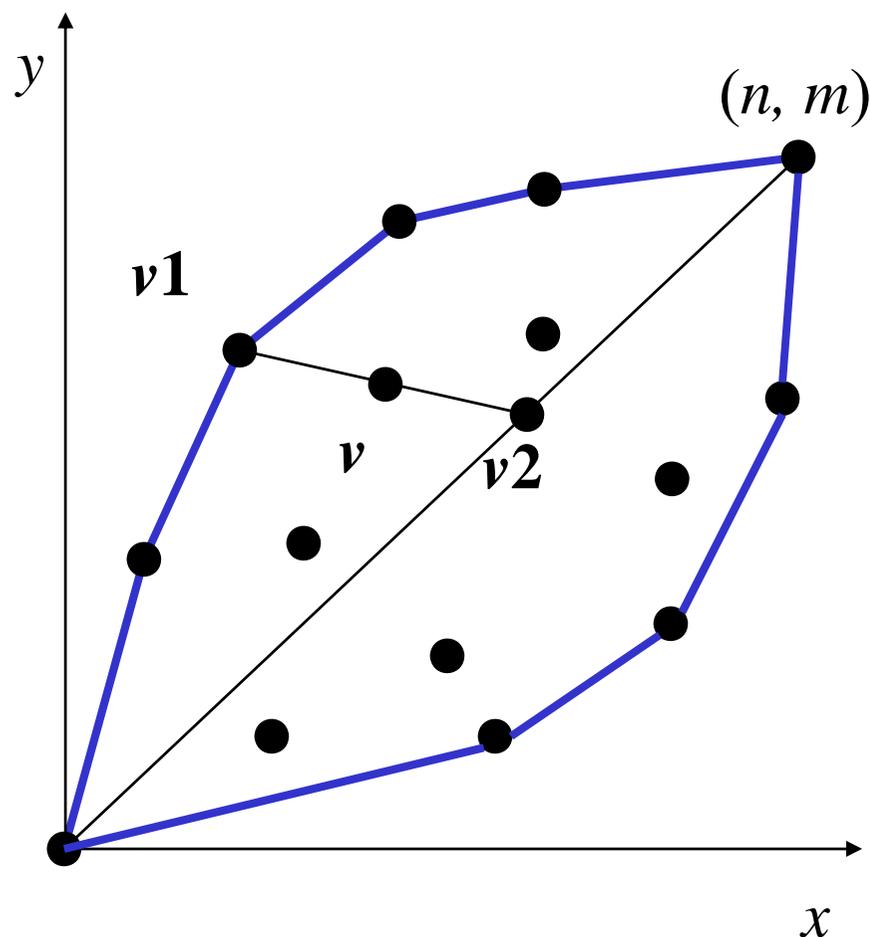
定理 ある j が存在して $S = \{1, \dots, j\}$ または $\{j, \dots, k\}$ が Entropy Gain を最大化

Shinichi Morishita. On Classification and Regression. *Proc. of First Int'l Conf. on Discovery Science, Lecture Notes in Artificial Intelligence*, Springer, Vol. 1532, pages 40-57 (1998)

<http://mlab.cb.k.u-tokyo.ac.jp/~moris/paper/ds98.pdf.gz>

各テスト $T \in S$ に 2次元座標の点 (x, y) (スタンプ点と呼ぶ) を対応させる

凸包 (Convex Hull): すべてのスタンプ点を含む最小の凸多角形 (凸多角形 P : P の任意の2点を結ぶ直線が P 内を通過)

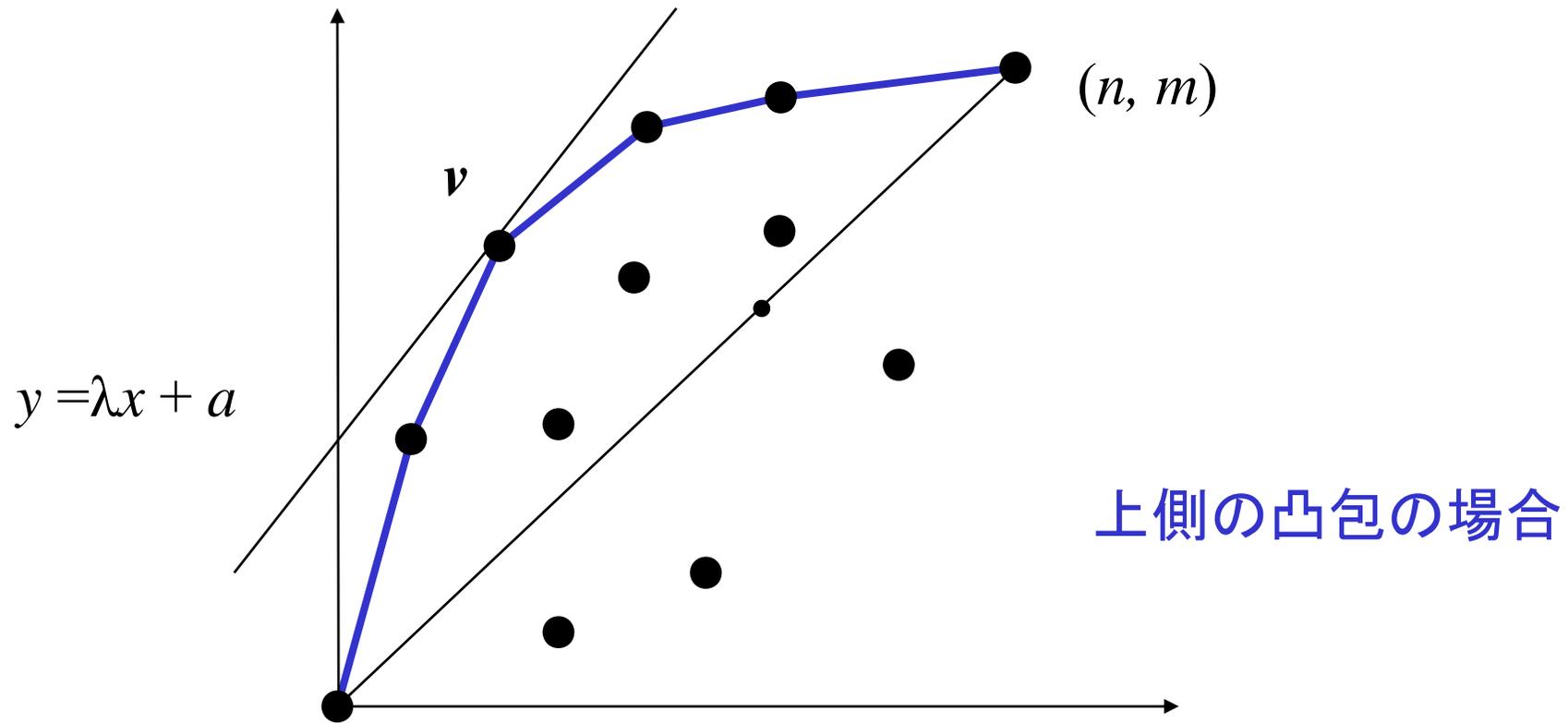


凸包内の任意のスタンプ点 v に対して、 v を内分点として持つ凸包上のスタンプ点 $v1$ と対角線上の点 $v2$ が存在する(演習問題)。

$$\text{Ent}(v) \leq \max \{ \text{Ent}(v1), \text{Ent}(v2) \}$$

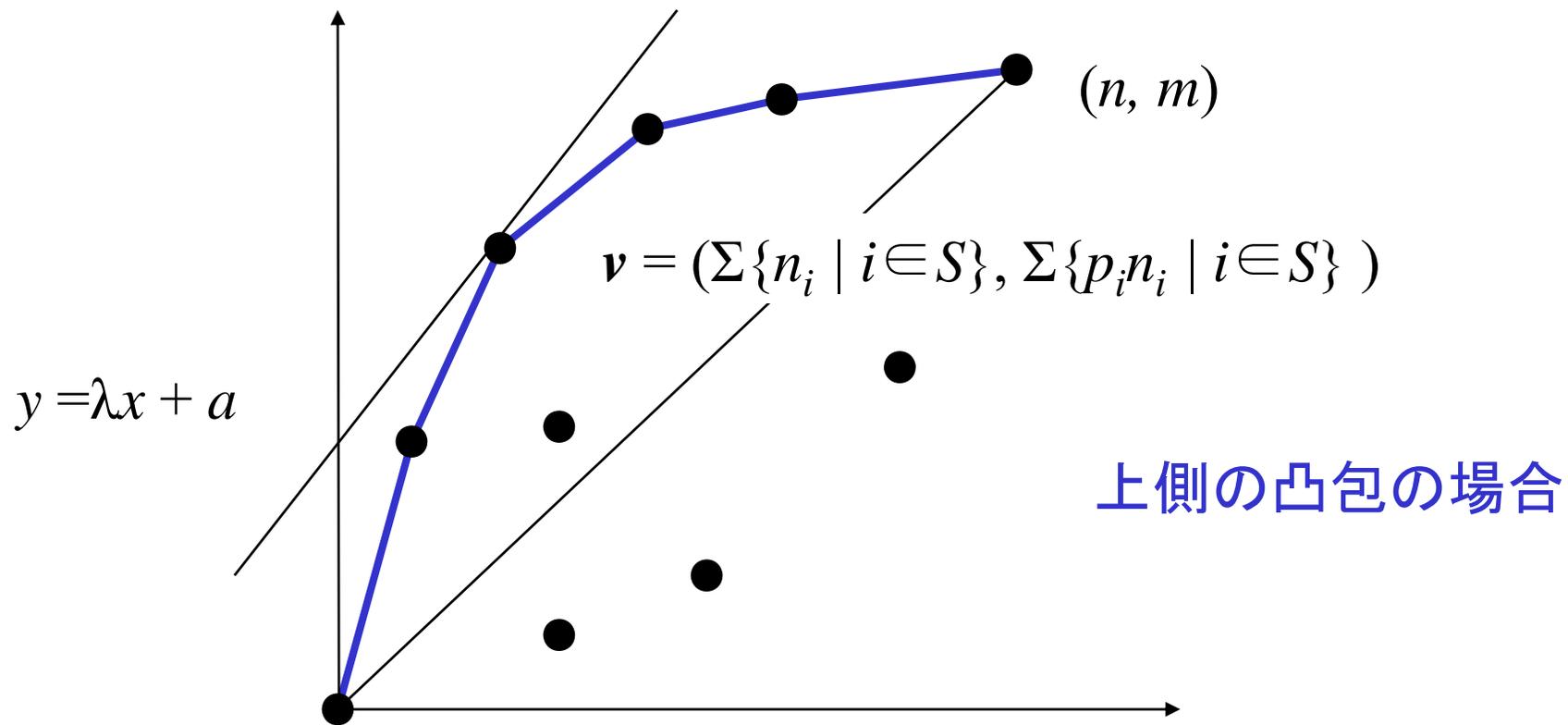
$\text{Ent}(v2)$ は最小値なので
 $\text{Ent}(v) \leq \text{Ent}(v1)$

Entropy Gain を最大化するスタンプ点が凸包上に存在



上側の凸包上のスタンプ点 v に対して、ある傾き λ の直線が存在し、全点中で v は y 切片 a を最大化

下側の凸包上のスタンプ点 v に対しては、ある傾き λ の直線が存在し、全点中で v は y 切片 a を最小化



$$\begin{aligned}
 a &= y - \lambda x \\
 &= \sum\{p_i n_i \mid i \in S\} - \lambda \sum\{n_i \mid i \in S\} \\
 &= \sum\{(p_i - \lambda) n_i \mid i \in S\}
 \end{aligned}$$

$p_1 \geq p_2 \geq \dots \geq p_k$ であるので
 最大化するのは $p_j - \lambda \geq 0$ と $p_{j+1} - \lambda < 0$ を満たす $S = \{1, \dots, j\}$
 ただし $p_1 - \lambda < 0$ のときは $S = \emptyset$, $p_k - \lambda > 0$ のときは $S = \{1, \dots, k\}$.