

H30 Ⅱ

- (1) 非負自然数 k は高々 $\log(n+1)$ コの bit を用いて 2進法表記できる。
 また、含める 1 の個数が k コの要素は nCk (2) (3) に存在しないため、これから
 それぞれ $0 \sim nCk-1$ までの相異なるインデックスを割り当てれば $\log nCk$ (2) の bit
 を用いて 2進法表記できる。
 以上2つの情報をあわせて、 \mathbb{Z} を k も含めて符号化するための符号長は高々

$$\log(n+1) + \log nCk \quad (\text{bit}) \quad \text{で十分} \quad \square$$

(2) 題意より

$$\begin{aligned} f(\theta) &= \sum_{t: x_t=1} (y_t - \theta)^2 \\ &= \sum_{t \in \{1, 2, 8, 9\}} (y_t - \theta)^2 = 3(1-\theta)^2 + (0-\theta)^2 = 4\theta^2 - 6\theta + 3 \\ &= 4\left(\theta - \frac{3}{4}\right)^2 + \frac{3}{4} \quad \text{from } \theta = \frac{3}{4} \end{aligned}$$

(3) $L(\mathbb{Z}) = \log(n+1) + \log nCk$ (但し、 k は \mathbb{Z} に含める bit 数)

• $i=1$ $(y_0^{(1)}, y_1^{(1)}) = (1000, 1110)$ より
 $\Delta(11y) = L(1000) + L(1110)$
 $= \log 5 + \log 4 + \log 5 + \log 4 = \log 400$

• $i=2$ $(y_0^{(2)}, y_1^{(2)}) = (11000, 110)$ より

$$\begin{aligned} \Delta(21y) &= L(11000) + L(110) \\ &= \log 6 + \log 10 + \log 4 + \log 3 = \log 720 \end{aligned}$$

• $i=3$ $(y_0^{(3)}, y_1^{(3)}) = (110, 11000)$ より

$$\begin{aligned} \Delta(31y) &= L(110) + L(11000) \\ &= \log 4 + \log 3 + \log 6 + \log 10 = \log 720 \end{aligned}$$

$$\cdot i=4 \quad (y_0^{(4)}, y_1^{(4)}) = (1100, 1100) \text{ である}$$

$$\begin{aligned} \Delta(i|y) &= L(1100) + L(1100) \\ &= \log 5 + \log 6 + \log 5 + \log 6 = \log 900 \end{aligned}$$

よって) $\Delta(i|y)$ を最小にする i は $i=1$ である. ($\Delta(1|y) < \Delta(2|y) = \Delta(3|y) < \Delta(4|y)$)

(4) $i^*=1$ にて) S を $y_1^{(1)}, y_1^{(0)}$ に分割した後、各々の x_2, x_3, x_4, y の値を置き換えることにより、

$y_1^{(1)}$					$y_1^{(0)}$				
t	x_2	x_3	x_4	y	t	x_2	x_3	x_4	y
1	0	0	1	1	4	1	1	0	1
2	1	0	0	1	5	0	1	1	0
3	0	1	1	1	6	0	1	0	0
4	1	1	0	0	7	0	0	1	0

同様に $i=2, 3, 4$ について $\Delta(i|y_1^{(1)}), \Delta(i|y_0^{(1)})$ を計算すると以下のようになる。

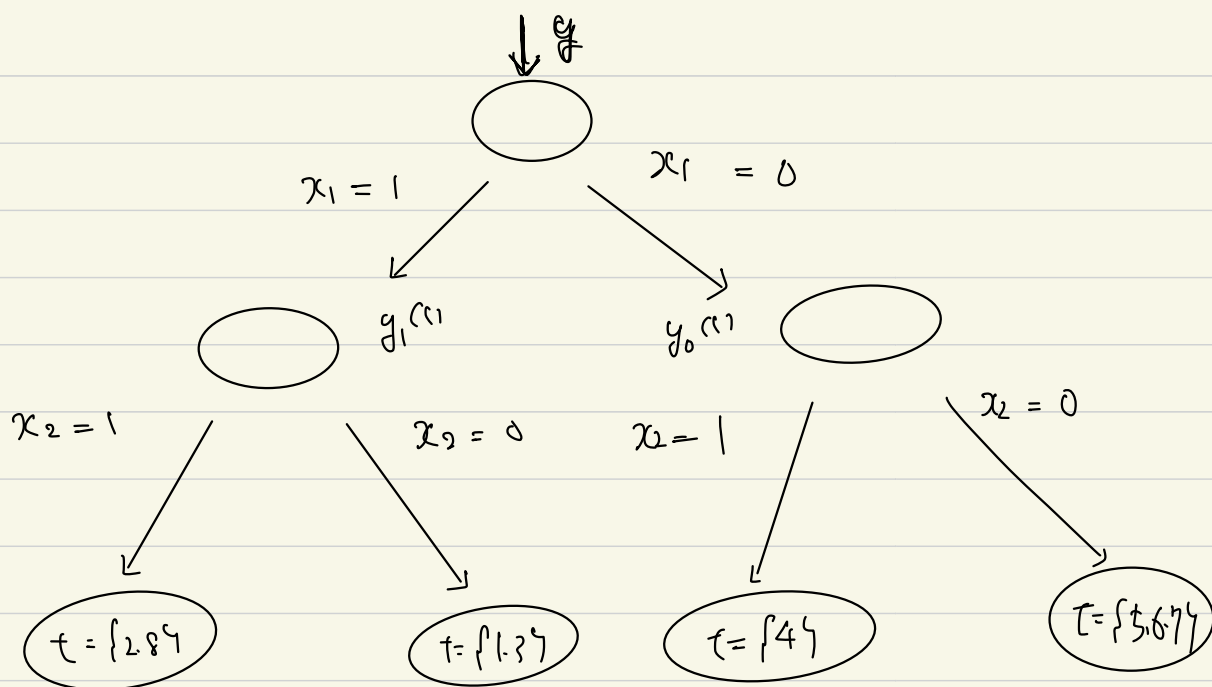
$$\left\{ \begin{array}{ll} \Delta(2|y_1^{(1)}) = L(10) + L(11) = \log 18 & \Delta(2|y_0^{(1)}) = L(1) + L(000) = \log 8 \\ \Delta(3|y_1^{(1)}) = L(10) + L(11) = \log 18 & \Delta(3|y_0^{(1)}) = L(100) + L(0) = \log 24 \\ \Delta(4|y_1^{(1)}) = L(11) + L(10) = \log 18 & \Delta(4|y_0^{(1)}) = L(00) + 2L(10) = \log 18 \end{array} \right.,$$

したがって、 $i=2$ (ランダム)

$i=2$

を再帰的な分割に用いる。

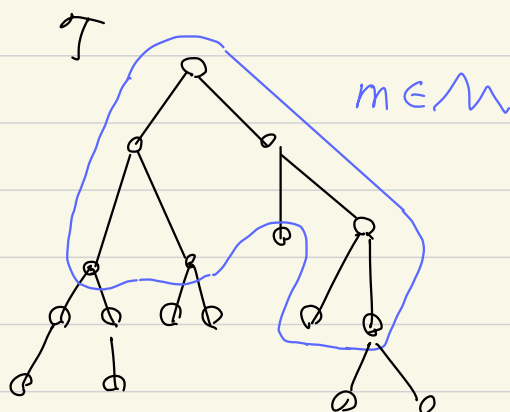
葉の深さが2になった時点で分割は終了するので、これで終了。得られる分割木は次のようになる。



(5) 機械学習分野における、いわゆる「過学習」とよばれる訓練データに過度に適合してモデルの作成が原因で、より一般的なデータにあってはいるであろう特徴量を見逃してしまうような現象が発生するためである。

この場合、モデル用訓練データは S であり、テストデータは新たに与えられるデータとなる。過学習を防ぎ、一般的なモデルに対する予測率を上げるためには、純粹に訓練データの数を増やすか、訓練データをいくつかの分割 (たとえばグループに分け、それらの部分からモデルを構築して全体のモデルを構築する) という方法が考えられる。

(6) 分割木での頂点数を N とした時、 $O(N^2)$ で動作するアルゴリズムを考える。但し、あらかじめ、木の木において、各頂点に分類されるものの番号一欄を各頂点に保存しておくものとする。…(*)



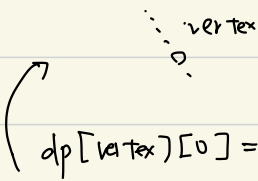
この時、各頂点において、以下の情報を保存する配列を用意する

$dp[vertex][0/1] := 0/1$ vertexから出る2分木の分割を

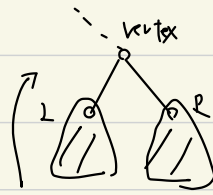
(木の木で、集合 M に属する部分木 $m \in M$ の例)

$\left\{ \begin{array}{l} 0 \dots \text{考慮しない} \\ 1 \dots \text{考慮する} \end{array} \right.$
 隣の
 $vertex$ を根とする部分木 M_{vertex} の最小コスト

図として表記すると



こゝで vertex に分類された 木の葉に対応する y を考える



さらに各分岐を行う.

$$dp[vertex][1] = \min_{\text{Lchild}} dp[Lchild][0/1] + \min_{\text{Rchild}} dp[Rchild][0/1] + C_I$$

これから、 $dp[vertex][0/1]$ の値はその子頂点から動的計画法を用いて計算することが出来る.

したがって、根を始点とする DFS (深さ優先探索) を行う. その帰りがたに $dp[vertex][0/1]$ の値を更新してゆくことで、任意の頂点に対する $dp[vertex][0/1]$ の値を計算できる

以下にその疑似コードを示す (なお \parallel (コメント) を用いて適宜説明を加えている)

```
void dfs (int u) {  
    for v ∈ child[u]  
        dfs (v)  
    // これまでのデータから dp[u] の更新
```

$$dp[u][0] = C_L + C_I + L[y_u]$$

// y_u : u を根とする部分木の葉にあたる全てのデータの y を並べた得られる2元系列

$$dp[u][1] = C_I + \min_{0 \leq i \leq 1} dp[leftchild][i] + \min_{0 \leq i \leq 1} dp[rightchild][i]$$

// u の部分木に対する問題をその子孫の部分木の問題に分割する

```
}  
dfs (root);
```

この時、もともとの問題である $m \in M$ に対する最小ペナルティの値は

$dp[root][i]$ に等しいことは明らかである. また、その部分木の復元も $\min_{0 \leq i \leq 1}$

根から順に各頂点をたどってゆき、 $dp[vertex][0] / dp[vertex][1]$ のどちらが最適解と見做さ求めることで $O(N)$ で復元可能である.

最後に、前に挙げた疑似コードの計算量が $O(N^2)$ であることを示そう、

$dp[u][0]$ の更新 (には y_u の列挙のコストがかかる) は $y'_{leftchild} + y'_{rightchild}$
+ (自身が葉の場合、 y_u) において求められる。このコストは高々 $O(N)$ であるため、取得及び
 $L[y_u]$ の計算にも $O(N)$ のコストがかかる。

$dp[u][1]$ の更新には、 $O(1)$ の算術演算のみを要する。

任意の頂点について上記が成り立ち、かつ dfs 自体は $O(N)$ で終了するので、
全体の計算量は高々 $O(N^2)$ である。 ■