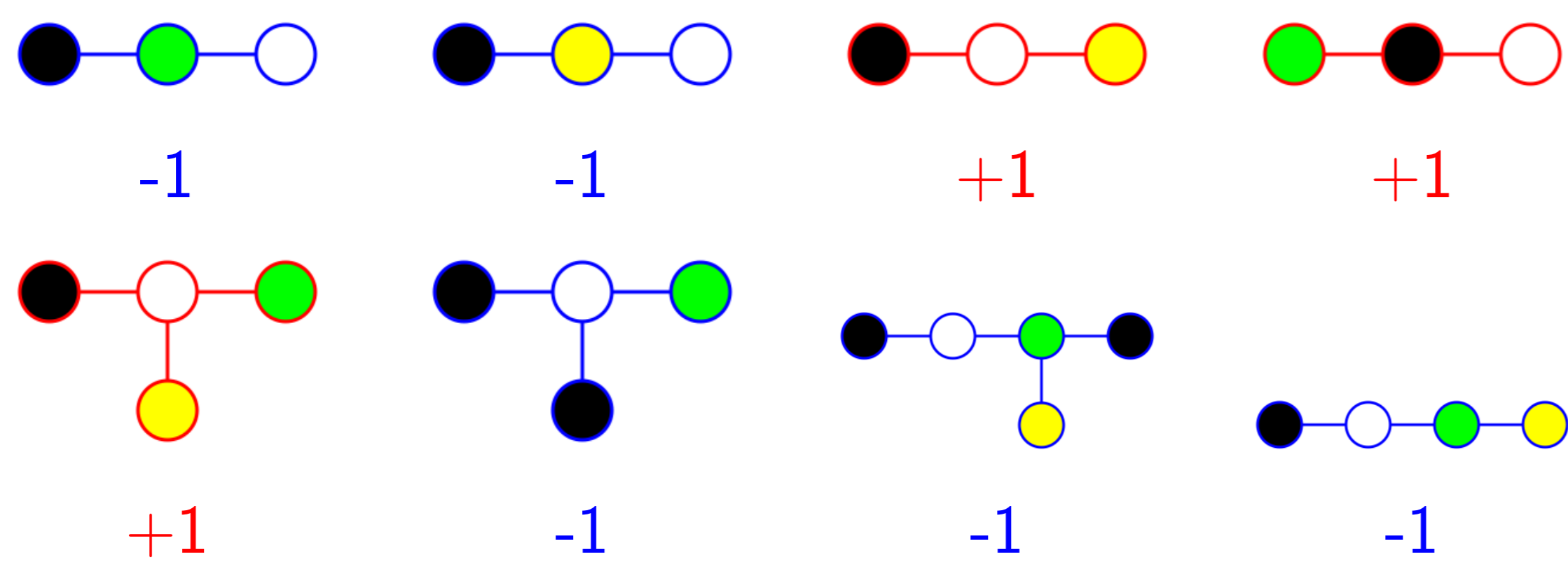


# 全部分グラフ指示子に基づく決定木勾配ブースティング

北海道大学 修士2年 横山侑政

## グラフの教師あり分類問題

入力 クラス ( $y = +1$  or  $-1$ ) のわかるグラフ  $N$  個



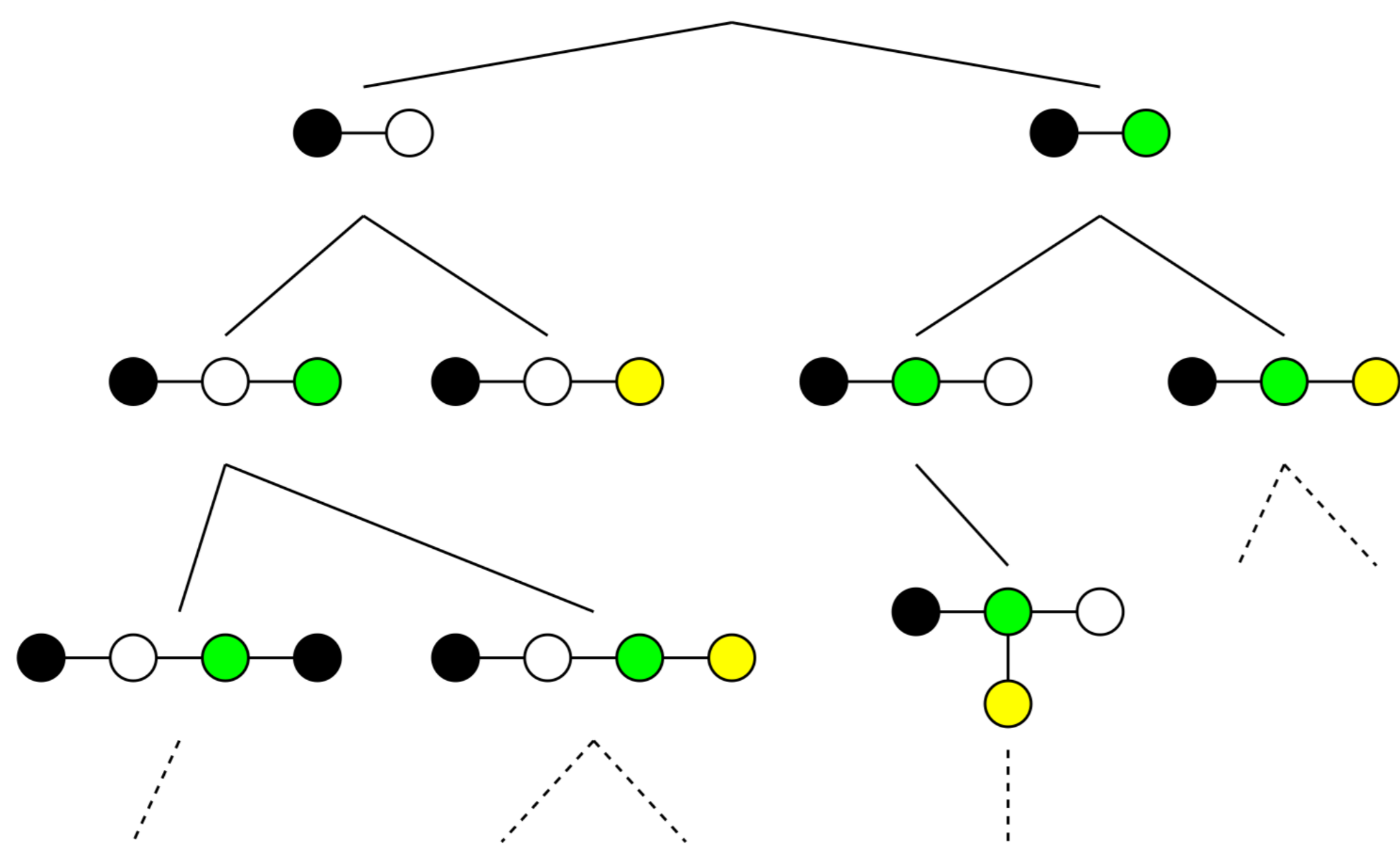
出力 グラフ  $G$  が  $+1$  か  $-1$  を予測する予測モデル  $F$

特徴量 部分グラフ有無

	$x_1$	$x_2$	$x_3$	$x_4$	...
$G_1$	1	0	1	0	...
$G_2$	0	0	0	1	...
$G_3$	1	1	0	0	...
$G_4$	1	1	0	1	...

部分グラフの有無の計算

探索木を考えて  
深さ優先探索



親ノードは  
子ノードの  
部分グラフ

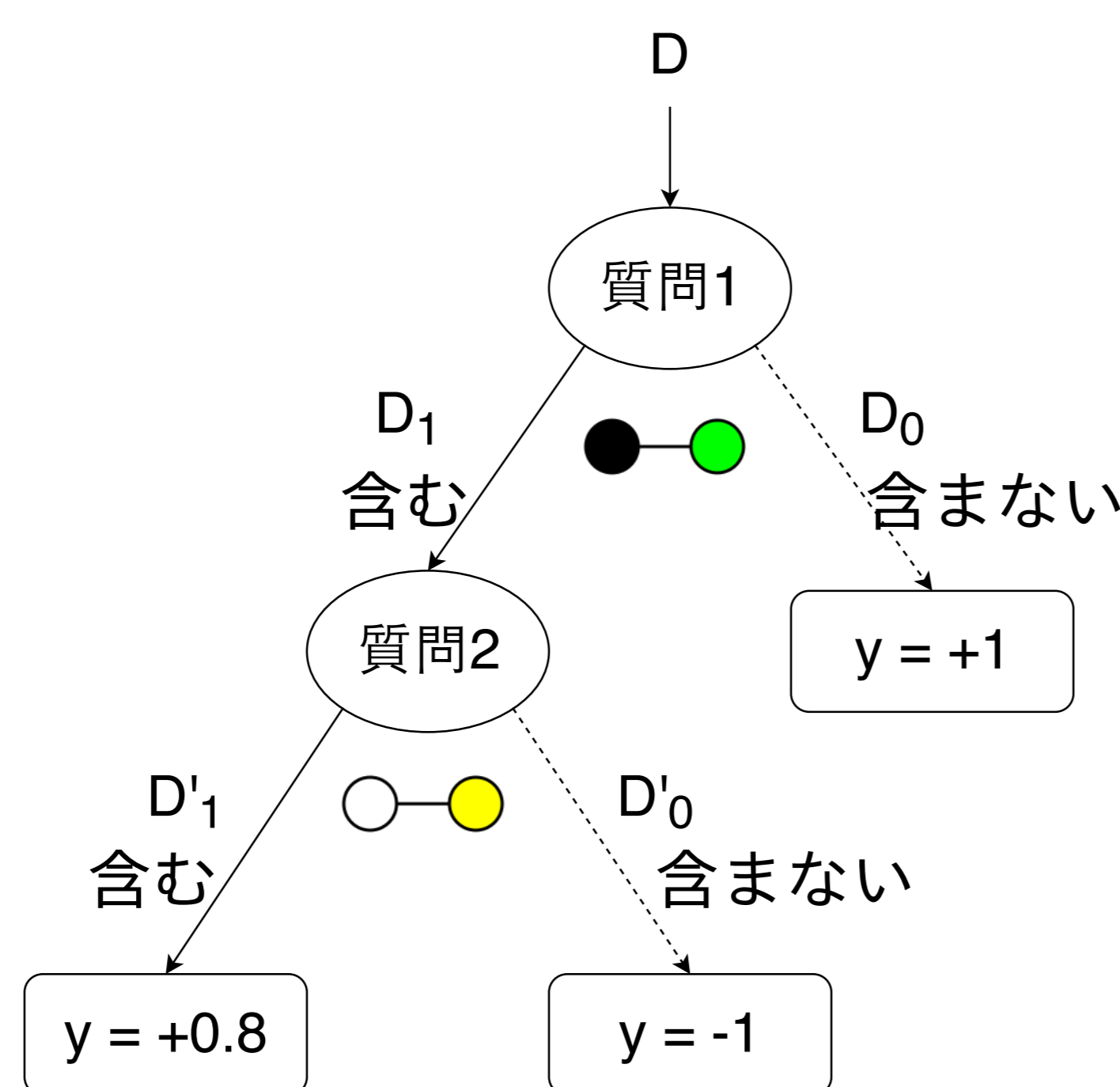
問題点 部分グラフの数が膨大  
全ての部分グラフの有無を計算するのは困難

対処法 陽には全ての部分グラフを計算せず  
決定木の勾配ブースティングモデル  $F$  を学習する

## 決定木の勾配ブースティング

決定木

入力データに  
内部ノードで質問をして  
葉ノードで予測値を返す



勾配ブースティング

決定木に加法的に学習する

$$F(G) = T_0(G) + sT_1(G) + sT_2(G) + sT_3(G) + \dots \quad (1)$$

$T_k$  は残差  $r_i$  を予測する回帰木

$$r_i = \frac{\partial L(y_i, F_{k-1}(G_i))}{\partial F} \quad (2)$$

$s$  はシュリンク係数,  $L$  は損失関数

## 内部ノードの学習

内部ノードは二乗誤差和を最小にする部分グラフを学習する

$$\arg \min_{x_j \in X} [\text{SSE}(D_1(x_j)) + \text{SSE}(D_0(x_j))] \quad (3)$$

$X$ : 全ての部分グラフの集合 (数が膨大で解けない)

$D_1(x_j) : \{G_i : x_j \text{ を含む}\}, D_0(x_j) : \{G_i : x_j \text{ を含まない}\}$

$\text{SSE}(D)$ : sum of square error,  $\min_m \sum_{G_i \in D} (r_i - m)^2$   
データ集合  $D$  の残差の二乗誤差の合計

## 部分グラフの有無の性質

部分グラフの有無の性質

部分グラフ  $x_j$  と その拡大グラフ  $x'_j$  を考えたとき

あるグラフ  $G_i$  に 部分グラフ  $x_j$  が無いならば 部分グラフ  $x'_j$  も無い

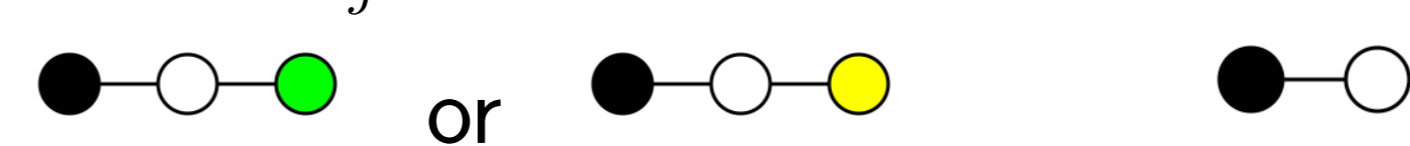


数式  $x_j \subset x'_j$  として  $G_i \not\supset x_j \Rightarrow G_i \not\supset x'_j$

逆に言えば

あるグラフ  $G_i$  が

部分グラフ  $x'_j$  を持つならば 部分グラフ  $x_j$  も持つ



## 部分グラフの有無の計算の枝刈り

部分グラフの有無の計算の枝刈り

式(3)を解きたいが全て部分グラフの集合  $X$  が大きすぎて計算困難

すでに有無を計算した部分グラフ  $x_j$  の拡大グラフ  $x'_j$  での

二乗誤差の下限值が分かれば枝刈りできる

二乗誤差和の下限值

部分グラフの有無の性質を使うと以下がいえる

$$\min_{x'_j \in X \supset x_j} [\text{SSE}(D_1(x'_j)) + \text{SSE}(D_0(x'_j))] \quad (4)$$

$$= \min_{S \subset D_1(x_j)} [\text{SSE}(D_1(x_j) \setminus S) + \text{SSE}(D_0(x_j) \cup S)] \quad (5)$$

$X \supset x_j : \{x'_j \mid x'_j \in X, x'_j \supset x_j\}$   $x_j$  の拡大グラフの集合

二乗誤差和の下限値の計算

式(5)は組み合わせを考える必要があるが

以下の等式を示すことで計算が簡単になる

(前からの部分集合 と後ろからの部分集合 を考えればよい)

$$\min_{C \subset A} [\text{SSE}(A \setminus C) + \text{SSE}(B \cup C)] \quad (6)$$

$$= \min \left( \min_{k=2, \dots, |A|-1} [\text{SSE}(A_{\leq k}) + \text{SSE}(B \cup A_{> k})], \min_{k=2, \dots, |A|-1} [\text{SSE}(A_{> k}) + \text{SSE}(B \cup A_{\leq k})] \right) \quad (7)$$

$A = \{a_1, \dots, a_n\}, A_{\leq k} = \{a_1, \dots, a_k\}, A_{> k} = \{a_{k+1}, \dots, a_n\}$

$B = \{b_1, \dots, b_m\}$

$\text{SSE}(A) = \min_m \sum_{a_i \in A} (a_i - m)^2$  集合  $A$  の二乗誤差の合計

証明  $\text{SSE}(A \setminus C) + \text{SSE}(B \cup C)$  (8)

$$= \sum_{i=0}^n (a_i - \bar{a}')^2 - \sum_{i=0}^k (c_i - \bar{a}')^2 + \sum_{i=0}^m (b_i - \bar{b}')^2 + \sum_{i=0}^k (c_i - \bar{b}')^2 \quad (9)$$

$$= - \sum_{i=0}^k (c_i - \bar{a})^2 - \frac{(\sum_{i=0}^k (c_i - \bar{a}))^2}{n-k} + \sum_{i=0}^n (a_i - \bar{a})^2 + \sum_{i=0}^k (c_i - \bar{b})^2 - \frac{(\sum_{i=0}^k (c_i - \bar{b}))^2}{m+k} + \sum_{i=0}^m (b_i - \bar{b})^2 \quad (10)$$

$$= - \left( \frac{1}{n-k} + \frac{1}{m+k} \right) \left( \sum_{i=0}^k c_i \right)^2 + \left( \frac{2\bar{a}n}{n-k} - \frac{2\bar{b}m}{m+k} \right) \sum_{i=0}^k c_i \quad (11)$$

$$- \frac{nk}{n-k} \bar{a}^2 + \frac{mk}{m+k} \bar{b}^2 + \sum_{i=0}^n (a_i - \bar{a})^2 + \sum_{i=0}^m (b_i - \bar{b})^2 \quad (12)$$

$k$  を固定したとき上に凸な  $\sum_{i=0}^k c_i$  の二次式

$\bar{a}$ :  $A$  の平均,  $\bar{a}'$ :  $A \setminus C$  の平均,  $\bar{b}$ :  $B$  の平均,  $\bar{b}'$ :  $B \cup C$  の平均,  $k = |C|$