



# データマイニング

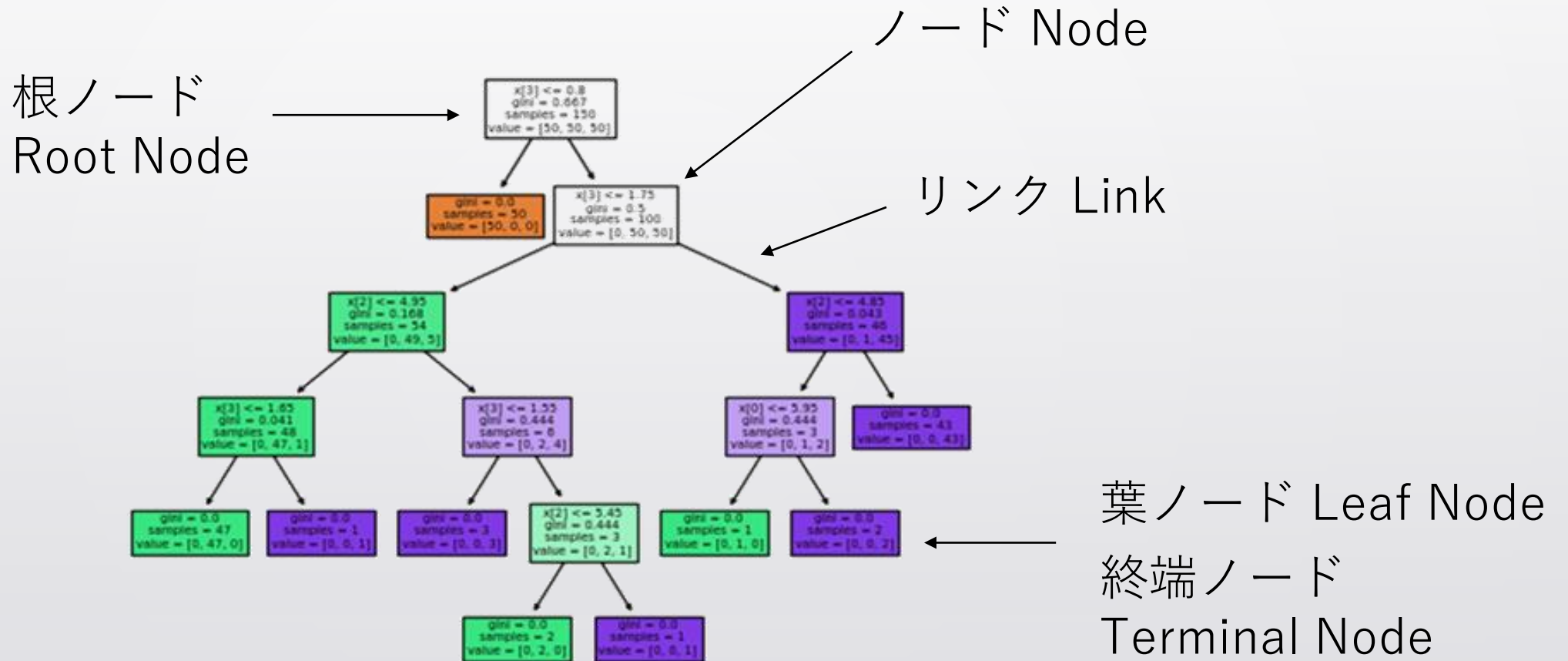
## Data Mining

### 10: 分類⑤ Classification

土居 裕和 Hirokazu Doi

長岡技術科学大学 Nagaoka University of Technology

# 決定木 Decision Tree



# 分割統治法 Divide and Conquer Induction

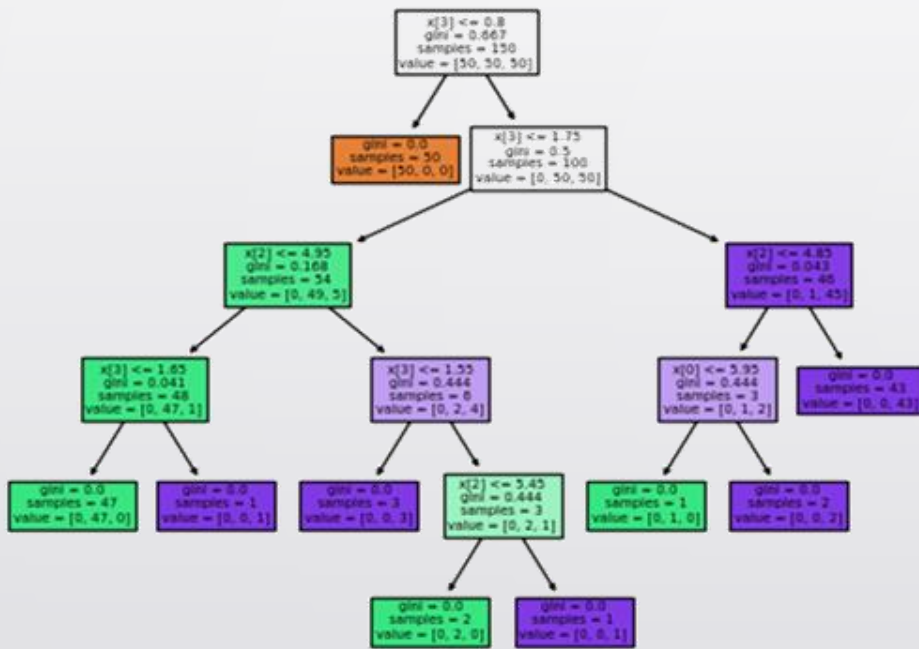
出来るだけ誤りなくデータを分類できる2  
分割基準でデータを分類する

Classify data by dichotomous criteria that minimizes  
false classification rate

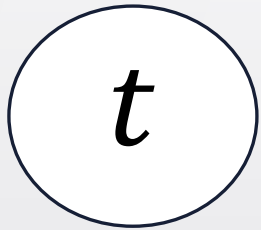


この手続きを繰り返すことで決定木を成長  
させる

Grow decision tree by repeating this procedure



## ノードのクラスの決定 Decision of the Class of a Node



$\{(x_i, t_i)\}$  : 学習データ  $x_i$  とクラスラベル  $t_i$  の集合

$C_j$  : クラス  $j$

$N_j$  : クラス  $j$  に属するデータの数

$N(t)$  : ノード  $t$  に属するデータの数

$N_j(t)$  : ノード  $t$  に属するクラス  $j$  のデータの数

## ベイズの定理 Bayes Theorem

$$P(H \cap C) = P(H|C)P(C)$$

$$P(H \cap C) = P(C|H)P(H)$$

$$P(C|H)P(H) = P(H|C)P(C)$$

$$P(C|H) = \frac{P(H|C)P(C)}{P(H)}$$



## ノードのクラス決定 Decision of the Class of a Node

ノード $t$ に属するデータが、クラス $C_j$ に属する事後確率を最大化するクラス $C_j$ を、ノード $t$ のクラスとする

Designate class  $C_j$  as the class of node  $t$  so that posterior probability of data in node  $t$  belonging to class  $C_j$  is maximized

$$P(C_j|t) = \frac{P(t|C_j)P(C_j)}{P(t)}$$

$$\text{ノード } t \text{ のクラス} = \operatorname{argmax}_j P(C_j|t)$$

Class of node  $t$

## ノードのクラス決定 Decision of the Class of a Node

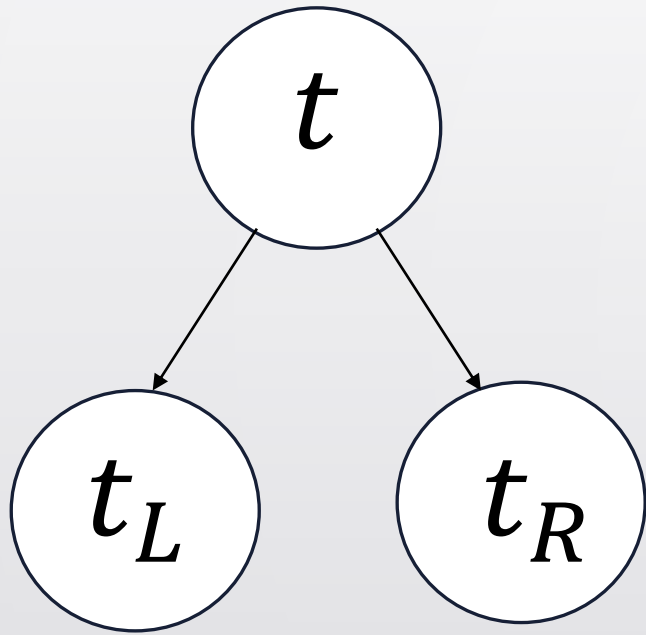
$$P(t|C_j)P(C_j) = \frac{N_j(t)}{N_j} \times \frac{N_j}{N} = \frac{N_j(t)}{N} \quad P(t) = \frac{N(t)}{N}$$

$$P(C_j|t) = \frac{P(t|C_j)P(C_j)}{P(t)} = \frac{N_j(t)}{N} \times \frac{N}{N(t)} = \frac{N_j(t)}{N(t)}$$

$$\text{ノード } t \text{ のクラス} = \operatorname{argmax}_j P(C_j|t) = \operatorname{argmax}_j \frac{N_j(t)}{N(t)}$$

Class of node  $t$

## ノードの不純度 Impurity of a Node



$$\Delta \text{Impurity} = \text{Impurity}_{\text{before}} - \text{Impurity}_{\text{after}}$$

$$\text{Impurity}_{\text{before}} = \text{Impurity}(t)$$

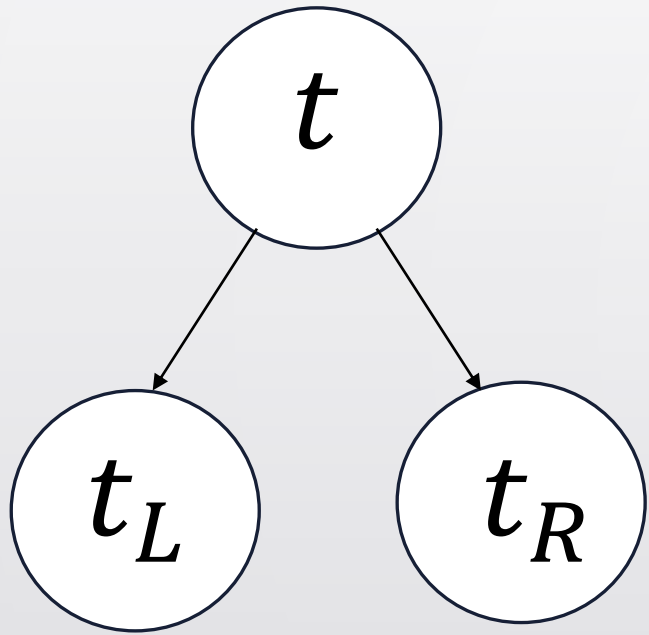
$$\text{Impurity}_{\text{after}} = P(t_L)\text{Impurity}(t_L) + P(t_R)\text{Impurity}(t_R)$$

分割前後の不純度の減少が最大になるようにデータを2分割する

Divide data so as to maximize the decrease of impurity  
 $\Delta \text{Impurity}$  after division

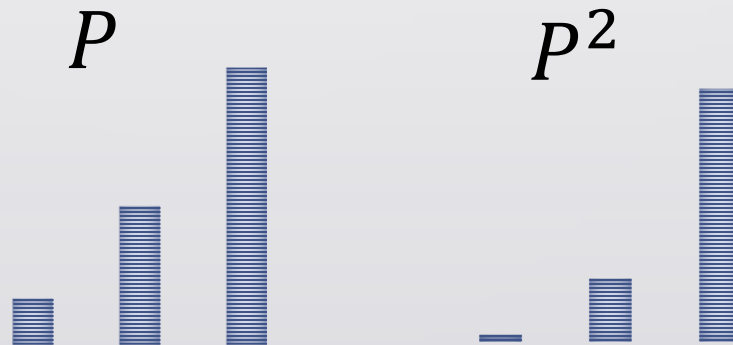


## ジニ係数 Gini Index

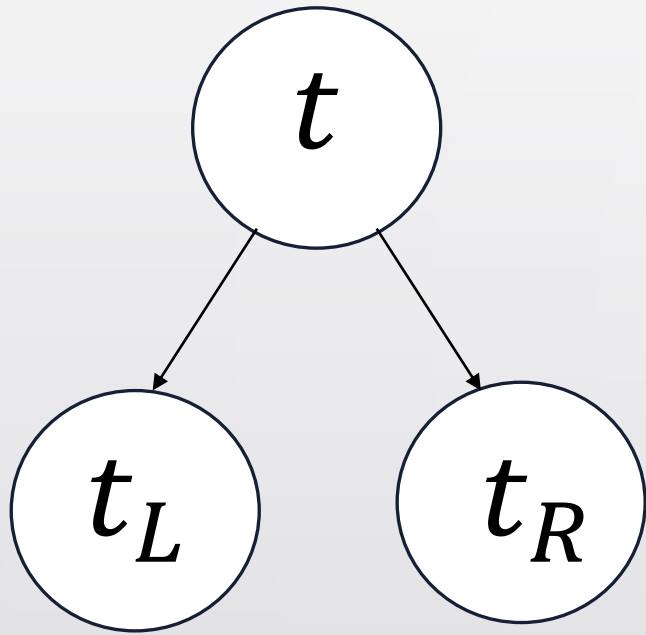


$$\text{Gini Index} = 1 - \sum_1^K P(C_j|t)^2$$

$P(C_j|t)$  のクラス間の違いが二乗により強調される  
Inter-class difference in  $P(C_j|t)$  is pronounced when squared



# エントロピー Entropy

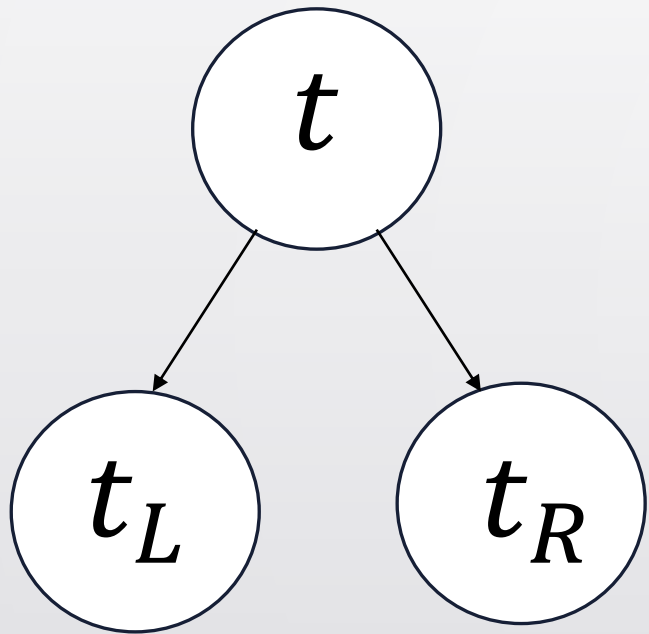


ノード  $t$  に含まれるクラスがばらついている  
Class of data belonging to node  $t$  is not uniform



ノード  $t$  の不純度が高いと  $P(C_j|t)$  が小さくなる  
 $P(C_j|t)$  gets small when impurity of node  $t$  is large

# エントロピー — Entropy

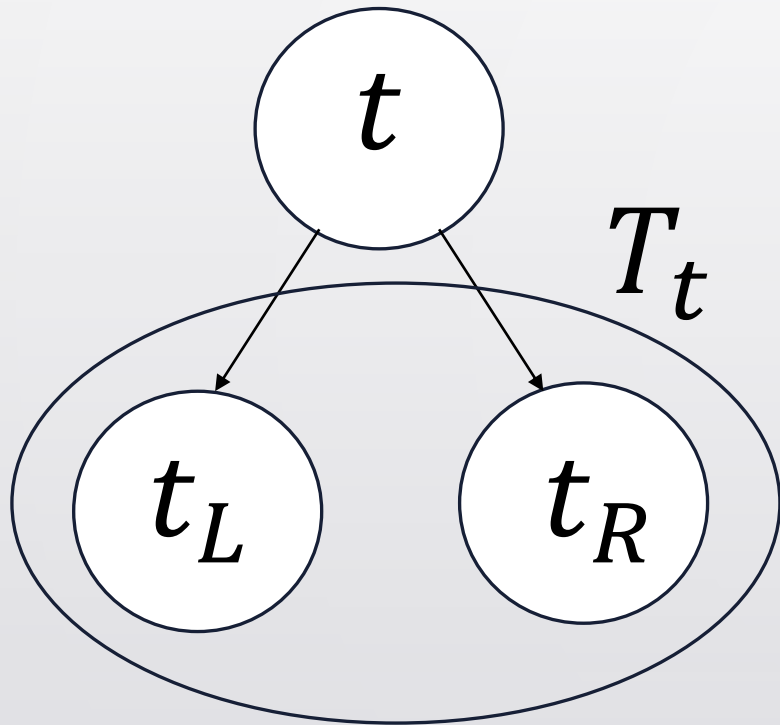


$$Entropy = - \int p \log(p) = E \left[ \log \frac{1}{p} \right]$$

確率が低い事象が起こると大きくなる  
Increases when an event with low-probability occurs

$$Impurity = - \sum_1^K P(C_j|t) \log P(C_j|t)$$

# 木の剪定 Pruning Decision Tree

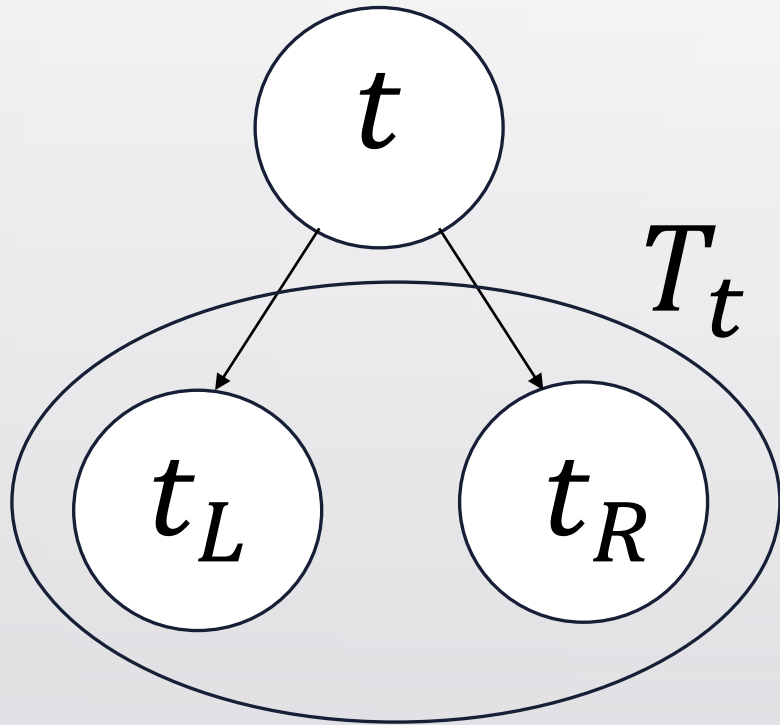


ノード  $t$  を終端ノードとすべきかどうか  
= ノード  $t$  の分岐  $T_t$  を削除すべきかどうか

The problem of whether node  $t$  should be determined as a terminal node equals to whether branch  $T_t$  of node  $t$  should be removed or not.

ノード  $t$  の分岐 Branch of node  $t$

## 木のコスト Cost of Tree



ノード  $t$  の分岐 Branch of node  $t$

$$R_{\alpha}(t) = \underline{R(t)} + \alpha$$

ノード  $t$  における誤り率 Error rate at node  $t$

$$R_{\alpha}(T) = \sum_{t \in T} R_{\alpha}(t) = R(T) + \alpha |\tilde{T}|$$

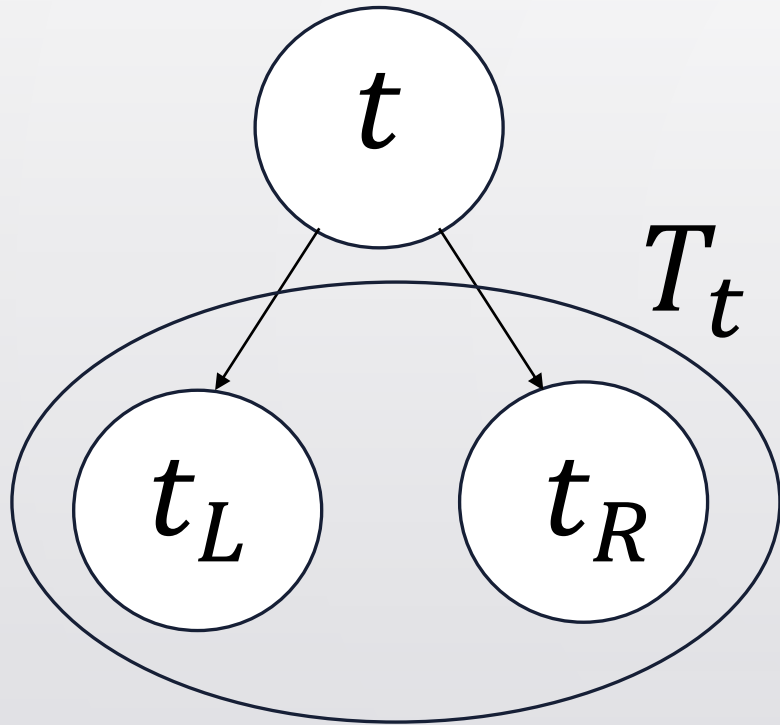
$|\tilde{T}|$  : 終端ノードの数 Number of terminal node

小さい木で高い正答率を達成したい

Aims to achieve high accuracy by decision tree with smaller size



## 木のコスト Cost of Tree



ノード  $t$  の分岐 Branch of node  $t$

$$R_\alpha(t) = R(t) + \alpha$$

$$R_\alpha(T_t) = R(T_t) + \alpha |\tilde{T}_t|$$

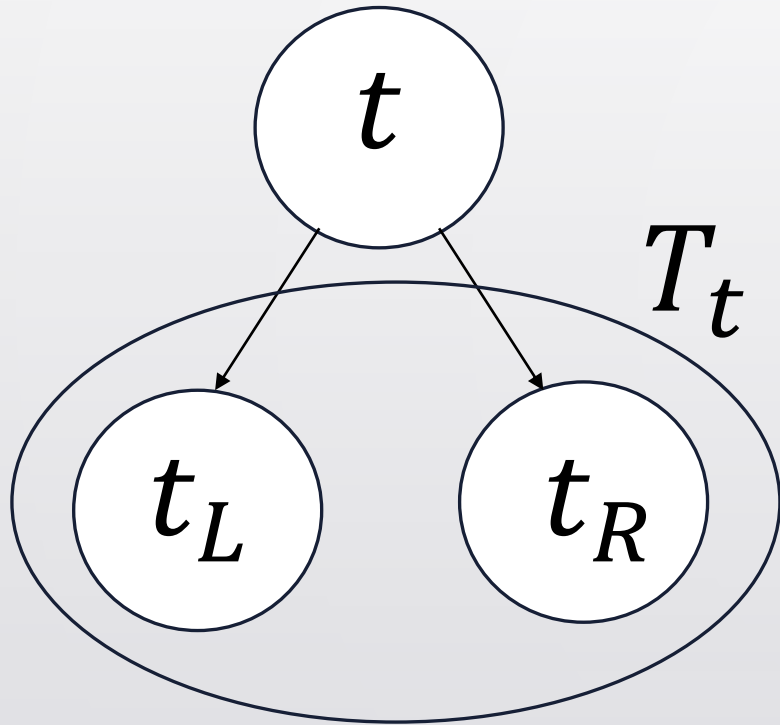
$T_t$  は複数のノードを含むことに注意

Note that  $T_t$  includes multiple nodes

$R_\alpha(T_t) \geq R_\alpha(t)$  なら  $T_t$  を削除する

Remove  $T_t$  if  $R_\alpha(T_t) \geq R_\alpha(t)$

## 木のコスト Cost of Tree



ノード  $t$  の分岐 Branch of node  $t$

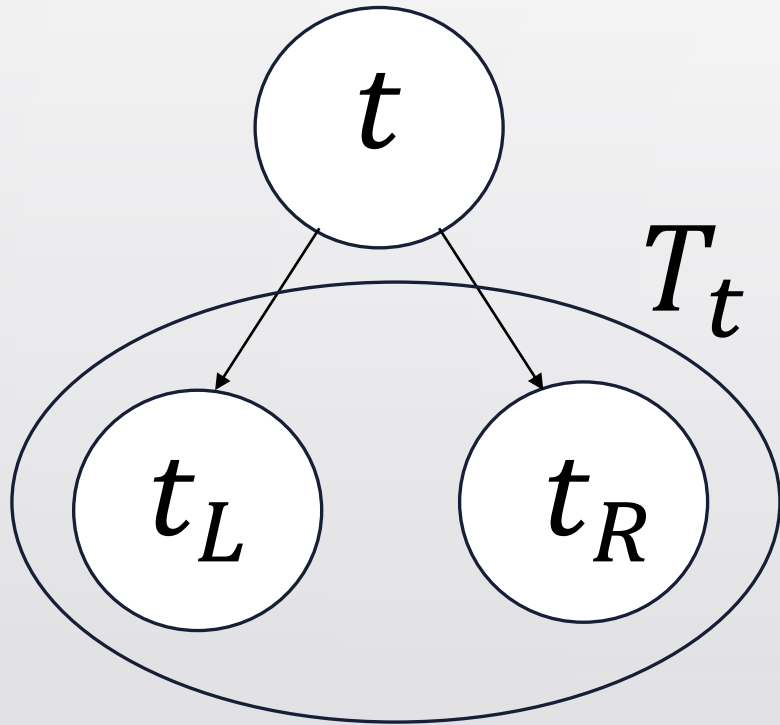
$$R(T_t) + \alpha |\tilde{T}_t| \geq R(t) + \alpha$$



$$\frac{R(t) - R(T_t)}{|\tilde{T}_t| - 1} \leq \alpha \text{ なら } T_t \text{ を削除する}$$

$$\text{Remove } T_t \text{ if } \frac{R(t) - R(T_t)}{|\tilde{T}_t| - 1} \leq \alpha$$

## 木のコスト Cost of Tree



ノード  $t$  の分岐 Branch of node  $t$

$$\frac{R(t) - R(T_t)}{|\tilde{T}_t| - 1} \leq \alpha \text{ なら } T_t \text{ を削除する}$$

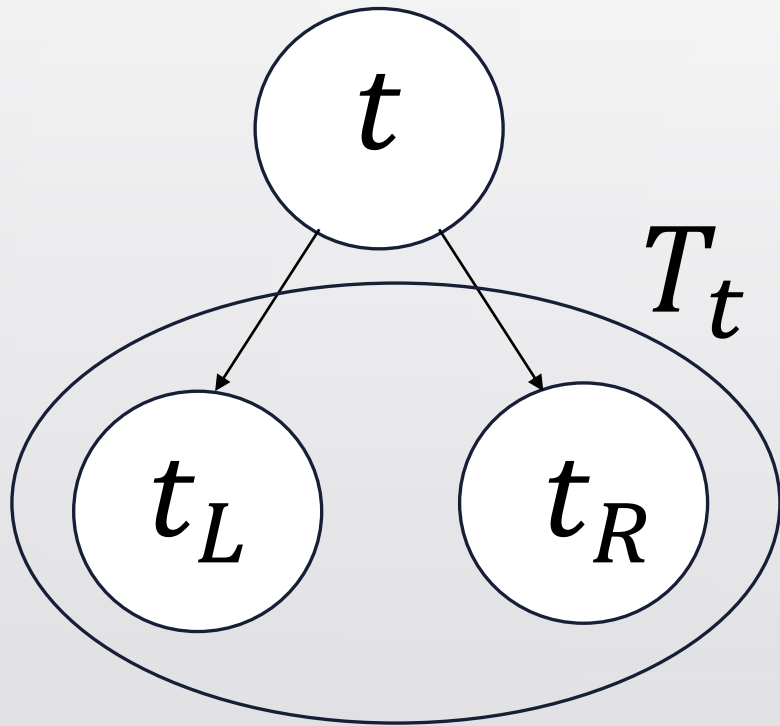


$$g(t) = \frac{R(t) - R(T_t)}{|\tilde{T}_t| - 1}$$

$g(t)$  は剪定するかどうかの判断基準を与える

$g(t)$  gives a criteria to decide whether to prune branch

# 剪定のアルゴリズム Algorithm of Pruning



- 全ての内部ノードについて  $g(t)$  を計算する  
Compute  $g(t)$  for all the internal nodes
- $g(t)$  が最小のノードの分岐を削除する  
Remove branch of the node with minimum  $g(t)$
- この手続きを繰り返す  
Repeat the procedure above until certain criteria is met

ノード  $t$  の分岐 Branch of node  $t$

# アンサンブル学習 Ensemble Learning

## ノーフリーランチ定理 No Free Lunch Theorem

あらゆる問題に対して最適な分類器は存在しない

There is no classifier that shows best performance for every classification problem

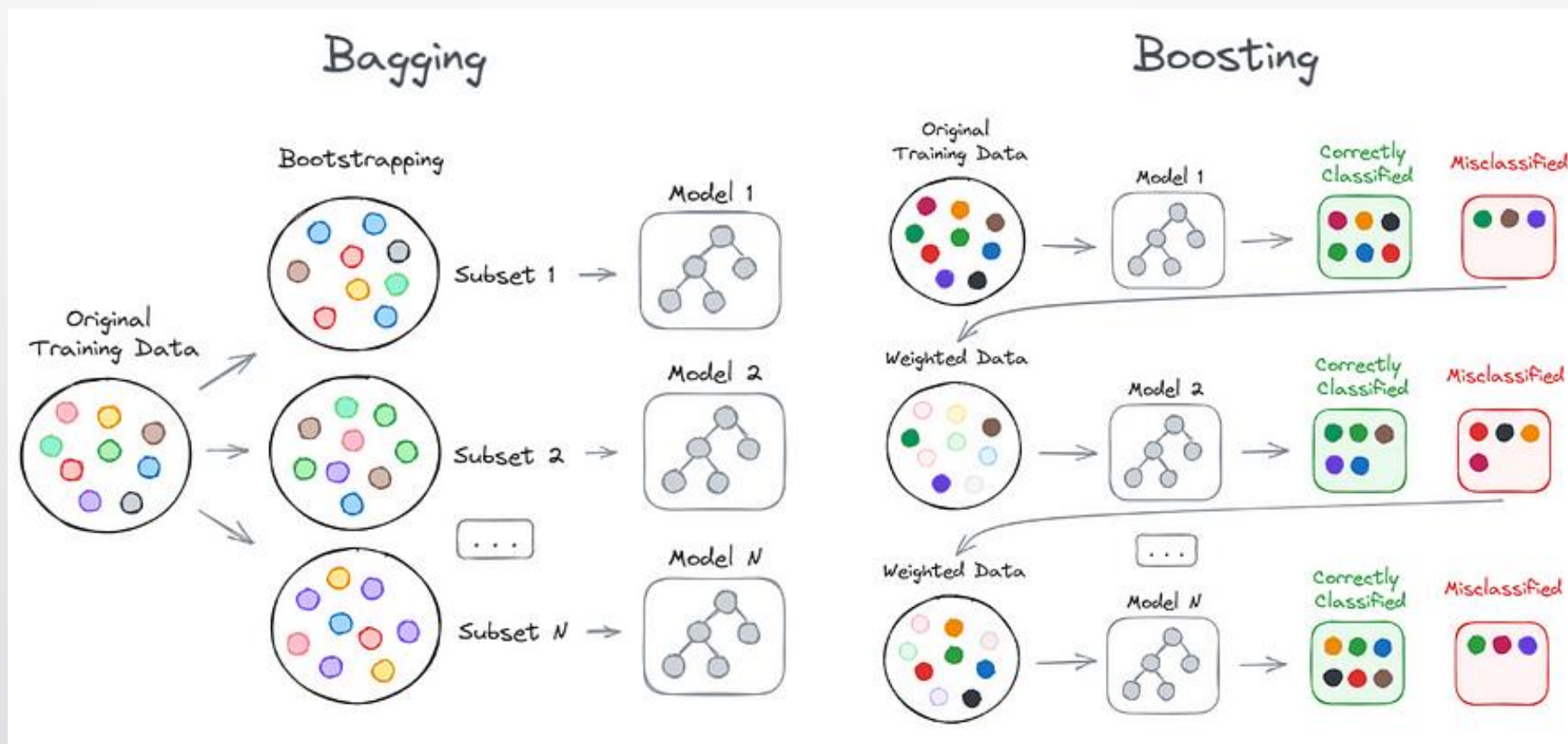


複数の弱識別器を組み合わせることで、精度のよい分類器を作る

Create superior classifier by combining multiple weak classifiers

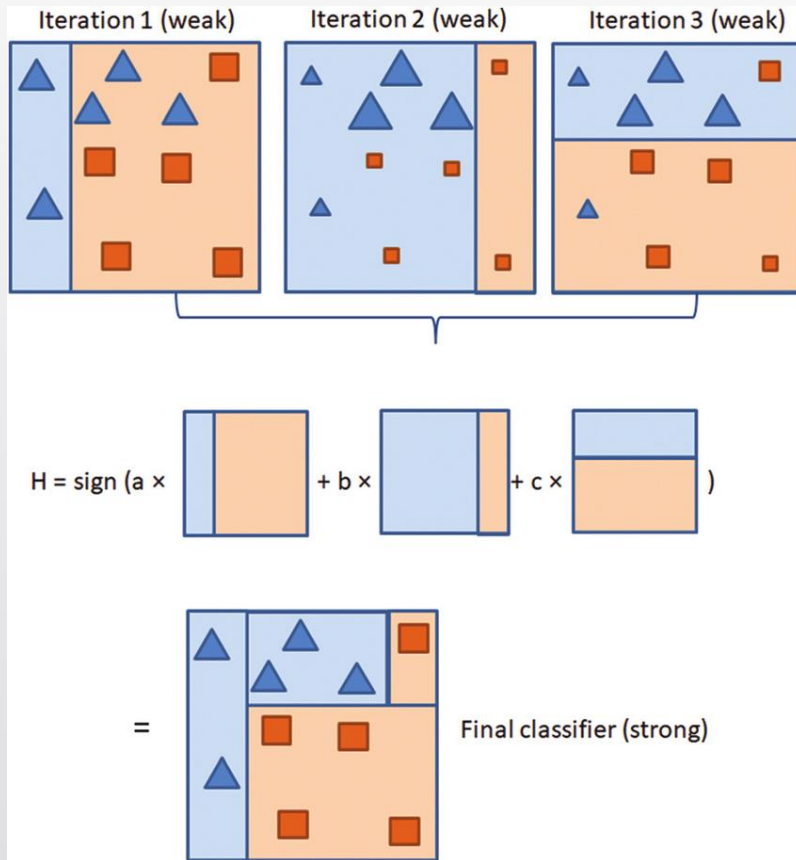


# バギングとブースティング Bagging and Boosting



<https://pub.towardsai.net/bagging-vs-boosting-the-power-of-ensemble-methods-in-machine-learning-6404e33524e6>

# アダブースト Ada(ptive)Boost(ing)

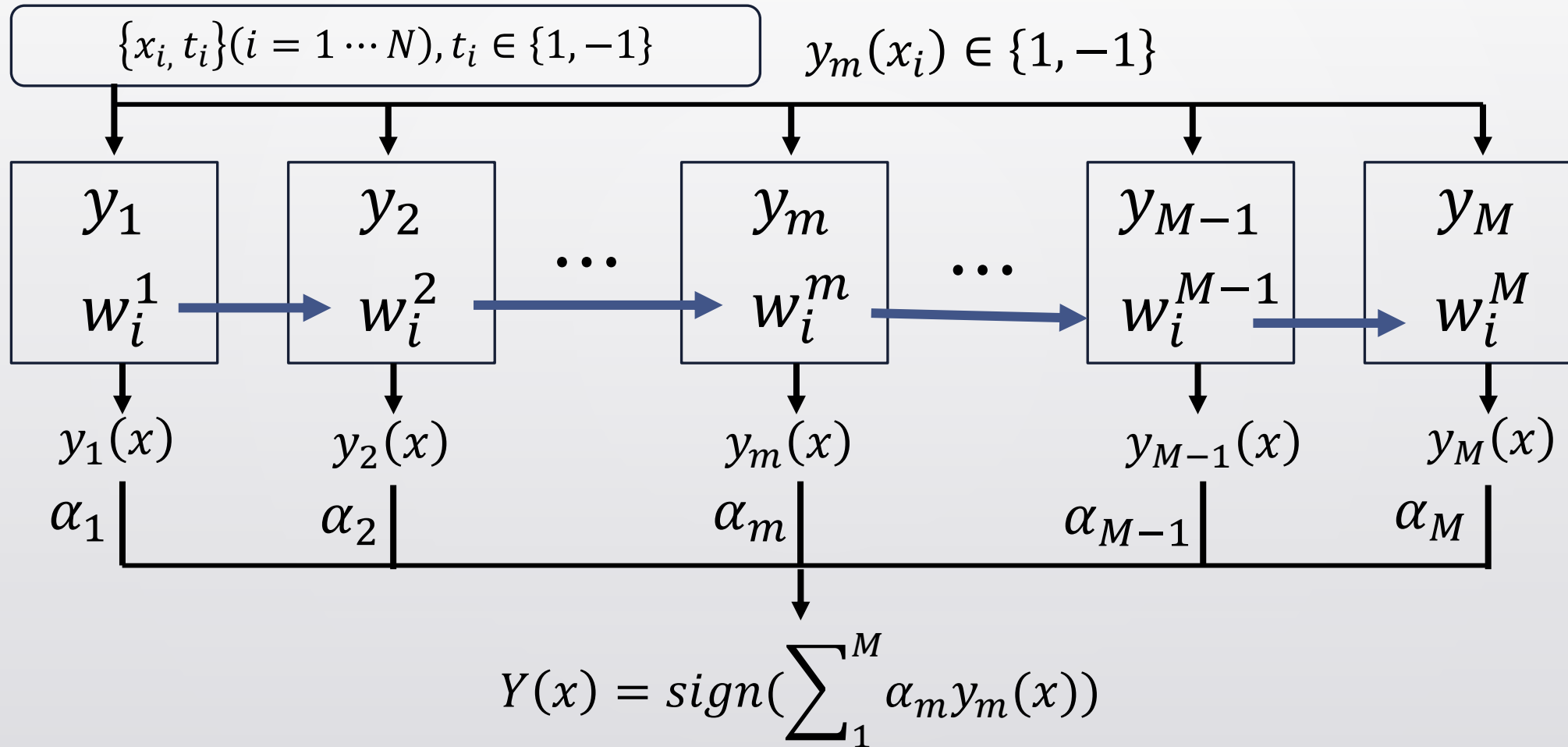


- ブースティングアルゴリズムの一つ  
One of boosting algorithms

- $t$ 個めの弱学習器が誤識別したデータの重みを大きくして、 $t + 1$ 個めの弱学習器をトレーニングする

Train  $t + 1$ -th weak learner by giving large weight to data points that  $t$ -th weak learner misclassified

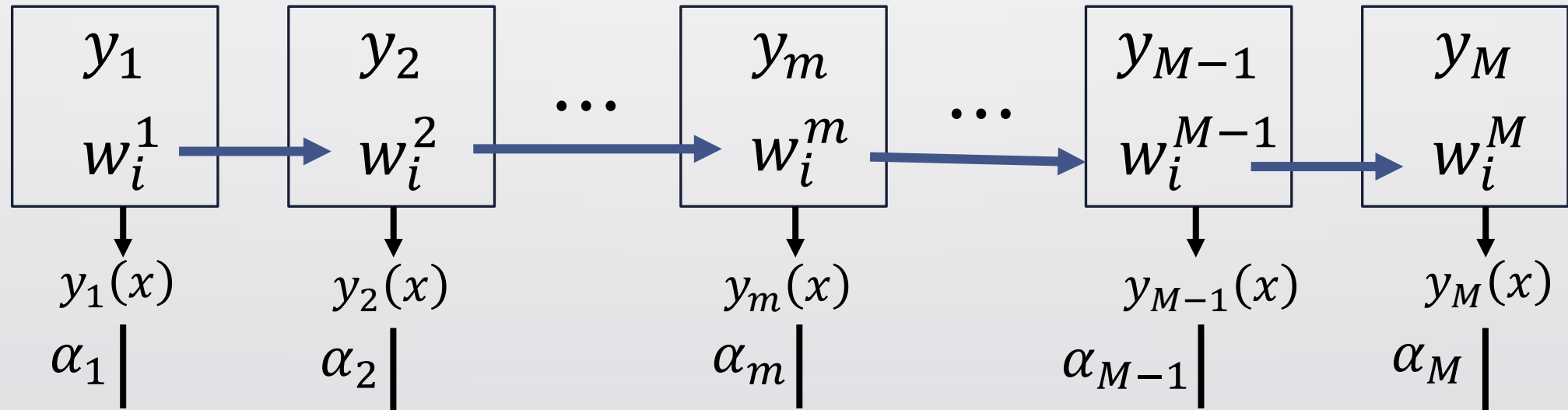
# アダブースト AdaBoost



# アダブースト AdaBoost

弱識別器 $y_m$ を逐次的に訓練する

Train weak learners  $y_m$  sequentially

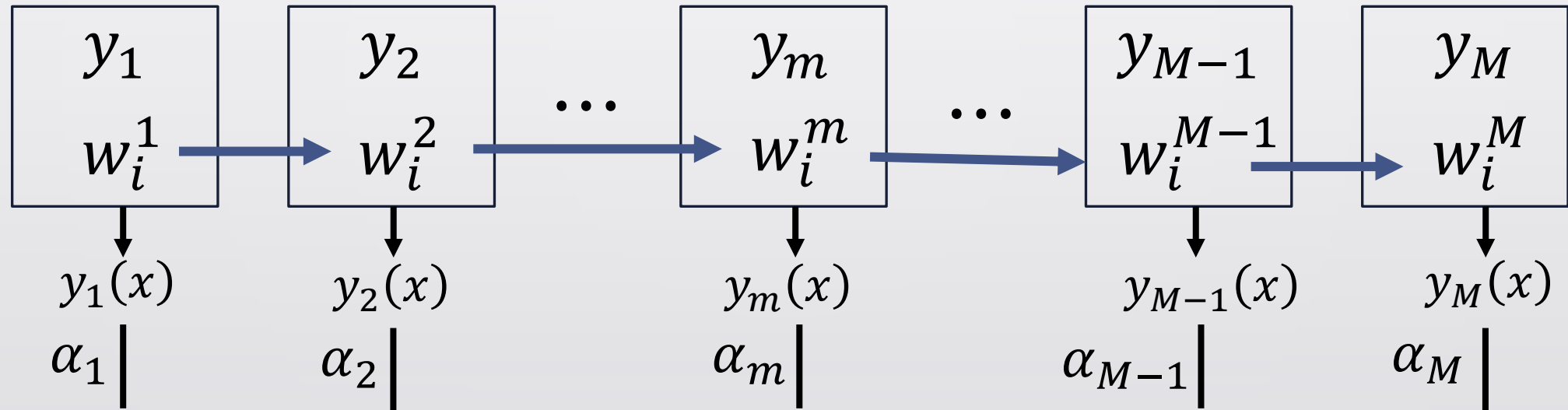




## アダブースト AdaBoost

データの重み  $w_i^m$  と弱学習器の重み  $\alpha_m$  が更新される

Update weight of data  $w_i^m$  and weight of weak learner  $\alpha_m$





## アダブースト AdaBoost

$E_m$  が最小になるよう弱識別器  $y_m$  を訓練する

Train weak learners  $y_m$  so that  $E_m$  is minimized

$$E_m = \frac{\sum_1^N w_i^m I(y_m(x_i))}{\sum_1^N w_i^m} \quad I(y_m(x_i)) = \begin{cases} 1 & (y_m(x_i) \neq t_i) \\ 0 & (y_m(x_i) = t_i) \end{cases}$$

## アダブースト AdaBoost

データの重み  $w_i^m$  と弱学習器の重み  $\alpha_m$  が更新される

Update weight of data  $w_i^m$  and weight of weak learner  $\alpha_m$

$$\alpha_m = \ln \left( \frac{1}{E_m} - 1 \right) \geq 0$$

$$w_i^{m+1} = \begin{cases} w_i^m \exp(\alpha_m) & (y_m(x_i) \neq t_i) \\ w_i^m & (y_m(x_i) = t_i) \end{cases}$$

# アダブースト AdaBoost

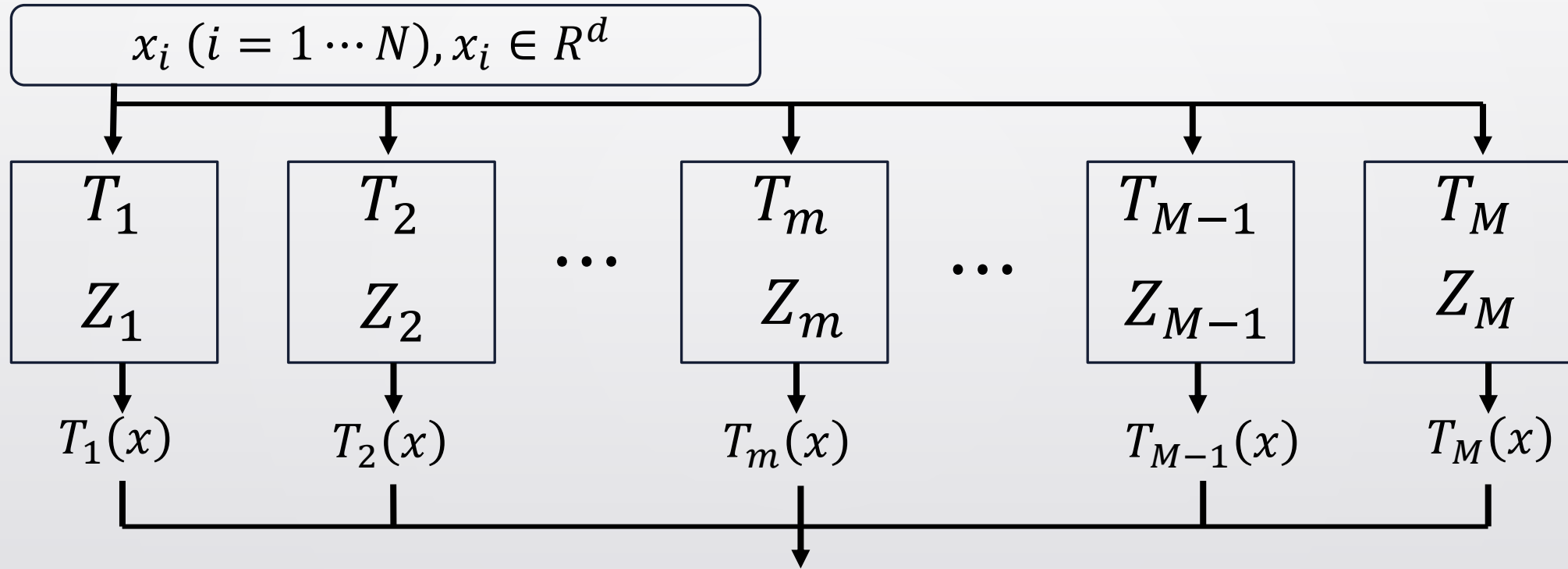
最終的な出力は、弱学習器の出力の重み付き和によって決まる

Output of AdaBoost classifier is determined by weighted sum of output of each weak classifier

The diagram illustrates the AdaBoost output calculation. At the top, a horizontal line represents the summation of weighted outputs from  $M$  weak classifiers. Above this line, the outputs  $y_1(x)$ ,  $y_2(x)$ ,  $y_m(x)$ ,  $y_{M-1}(x)$ , and  $y_M(x)$  are listed. Below the line, the corresponding weights  $\alpha_1$ ,  $\alpha_2$ ,  $\alpha_m$ ,  $\alpha_{M-1}$ , and  $\alpha_M$  are shown, each connected to the line by a vertical line segment. A downward arrow from the middle of the line points to the final equation:

$$Y(x) = \text{sign}\left(\sum_1^M \alpha_m y_m(x)\right)$$

# ランダムフォレスト Random Forest



多数決投票 Majority Voting

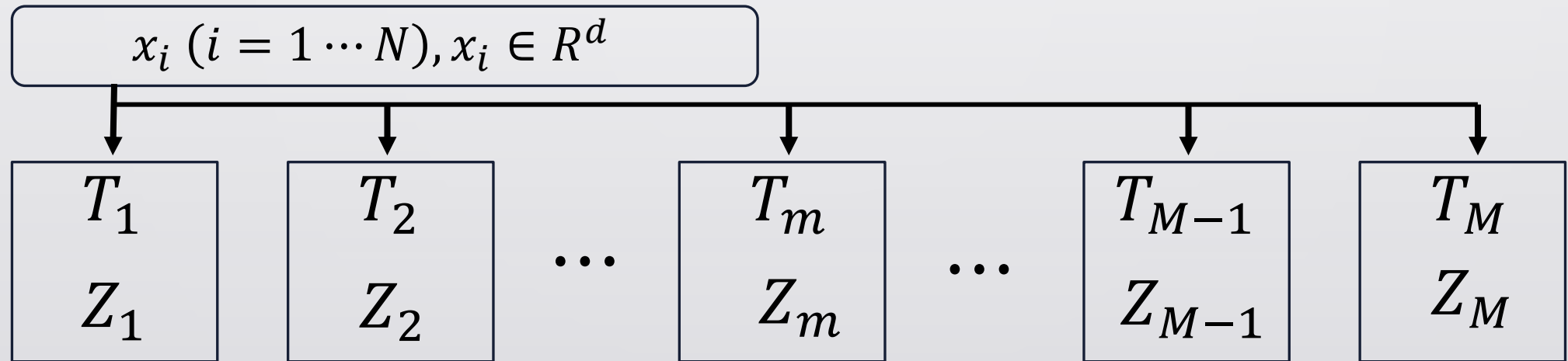
# ランダムフォレスト Random Forest

データからブートストラップサンプル $Z_m$ を抽出する

Extract bootstrapped samples  $Z_m$  from dataset

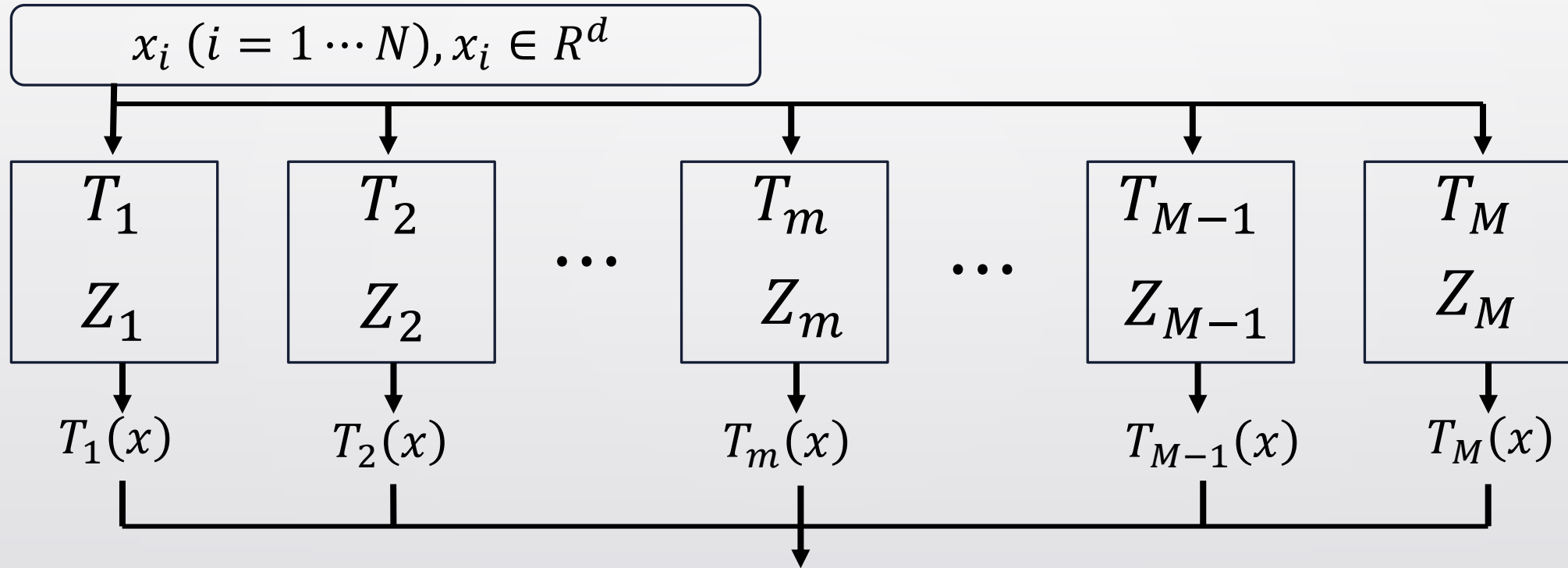
ランダムに選択した $d'$ 次元の特徴量を使って決定木 $T_m$ を構成する

Grow decision tree  $T_m$  based on randomly-selected  $d'$ -dimensional features





# ランダムフォレスト Random Forest



多数決投票 Majority Voting