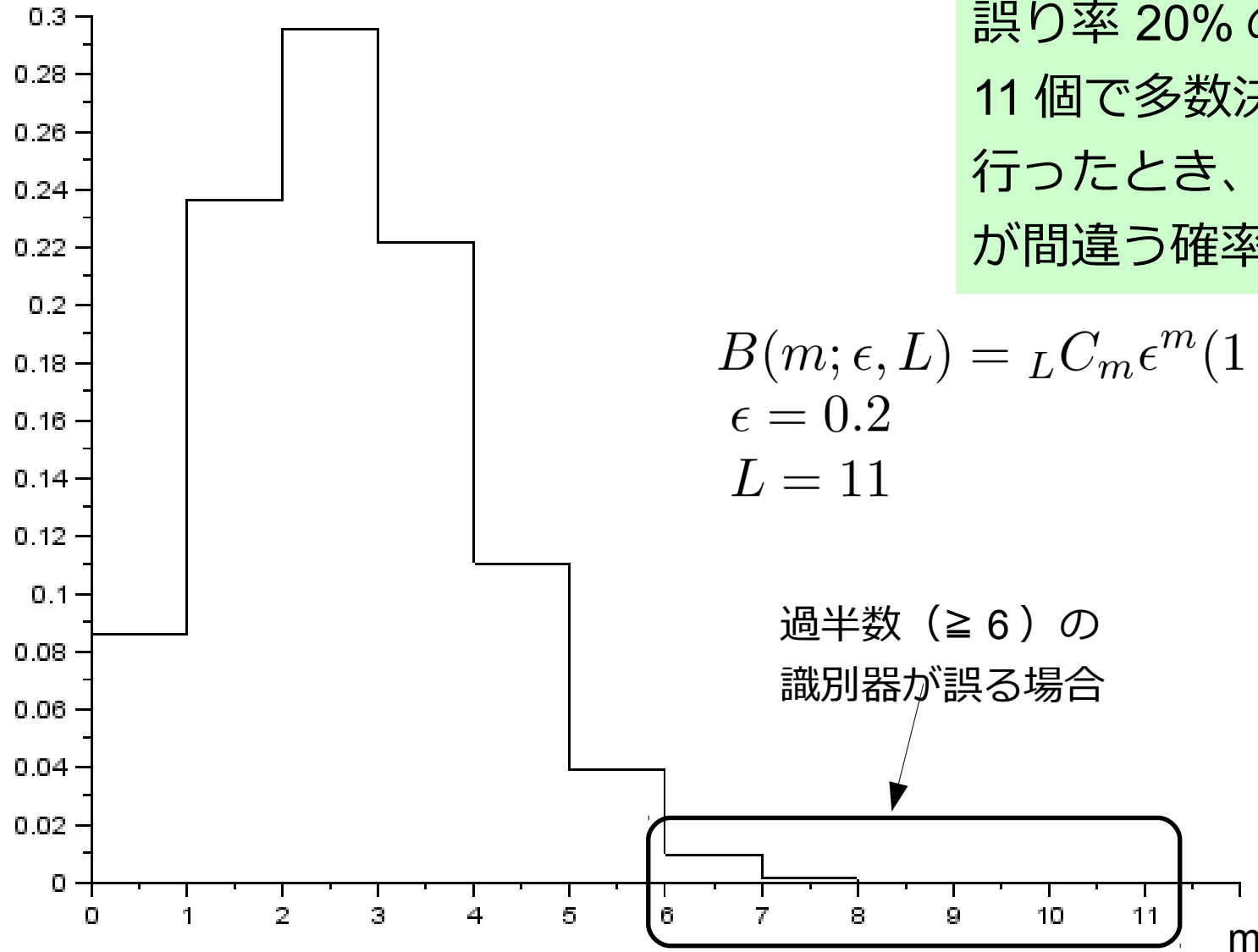


9. アンサンブル学習

- アンサンブル学習とは
 - 分類器を複数組み合わせ、それらの結果を統合することで個々の分類器よりも性能を向上させる方法
- アイディア
 - 訓練例集合から全く独立に L 個の分類器 (誤り率 ϵ , 誤りは独立) を作成
 - m 個の分類器が誤る確率は二項分布 $B(m; \epsilon, L)$
 - $\epsilon < 0.5$ のとき、 $m > L/2$ となる B は小さい値

9.1 なぜ性能が向上するのか



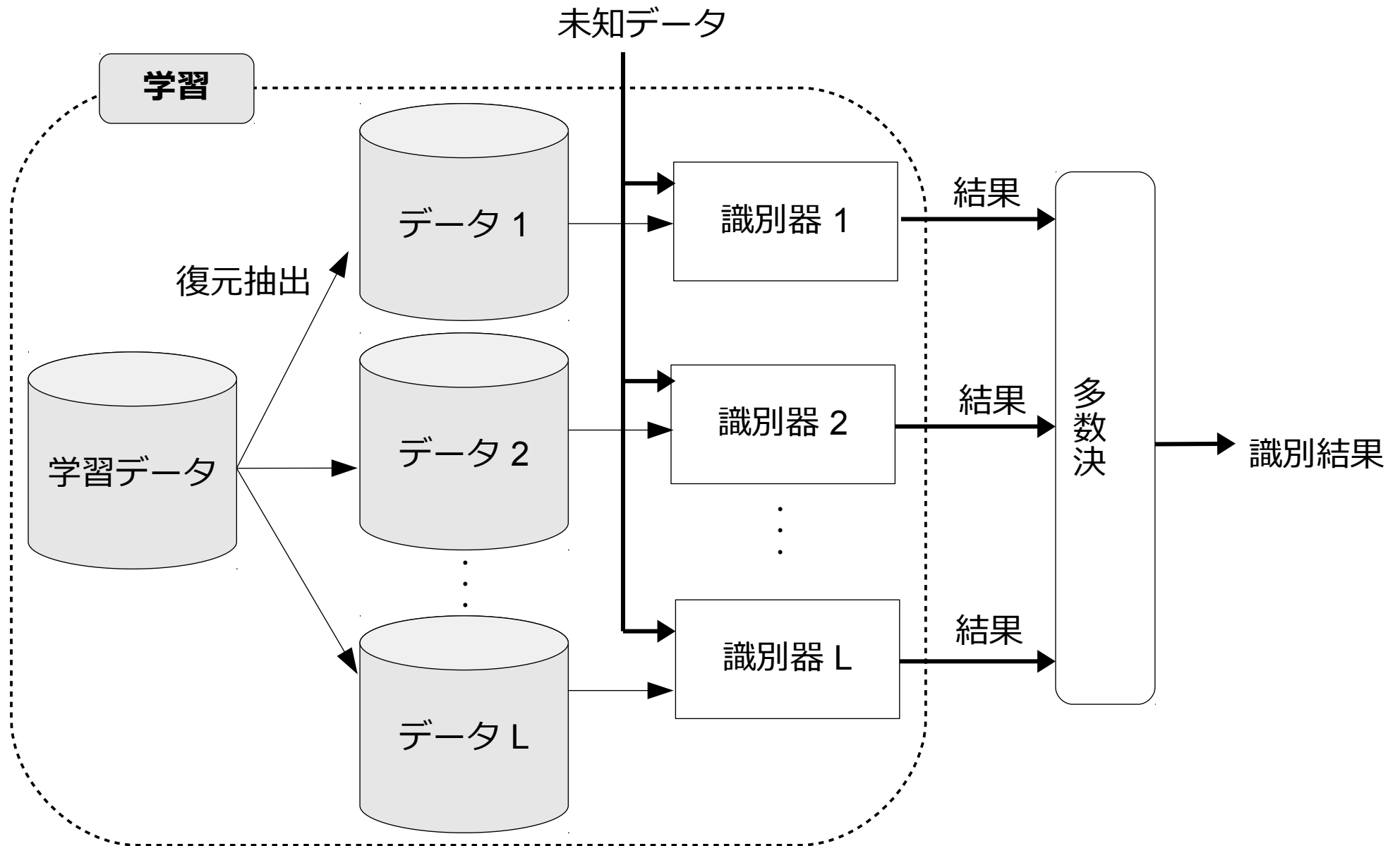
誤り率 20% の分類器
11 個で多数決統合を
行ったとき、6 個以上
が間違ふ確率は **1.2%**

$$B(m; \epsilon, L) = {}_L C_m \epsilon^m (1 - \epsilon)^{L-m}$$

$\epsilon = 0.2$
 $L = 11$

過半数 (≥ 6) の
識別器が誤る場合

9.2 バギング



9.2 バギング

- 特徴

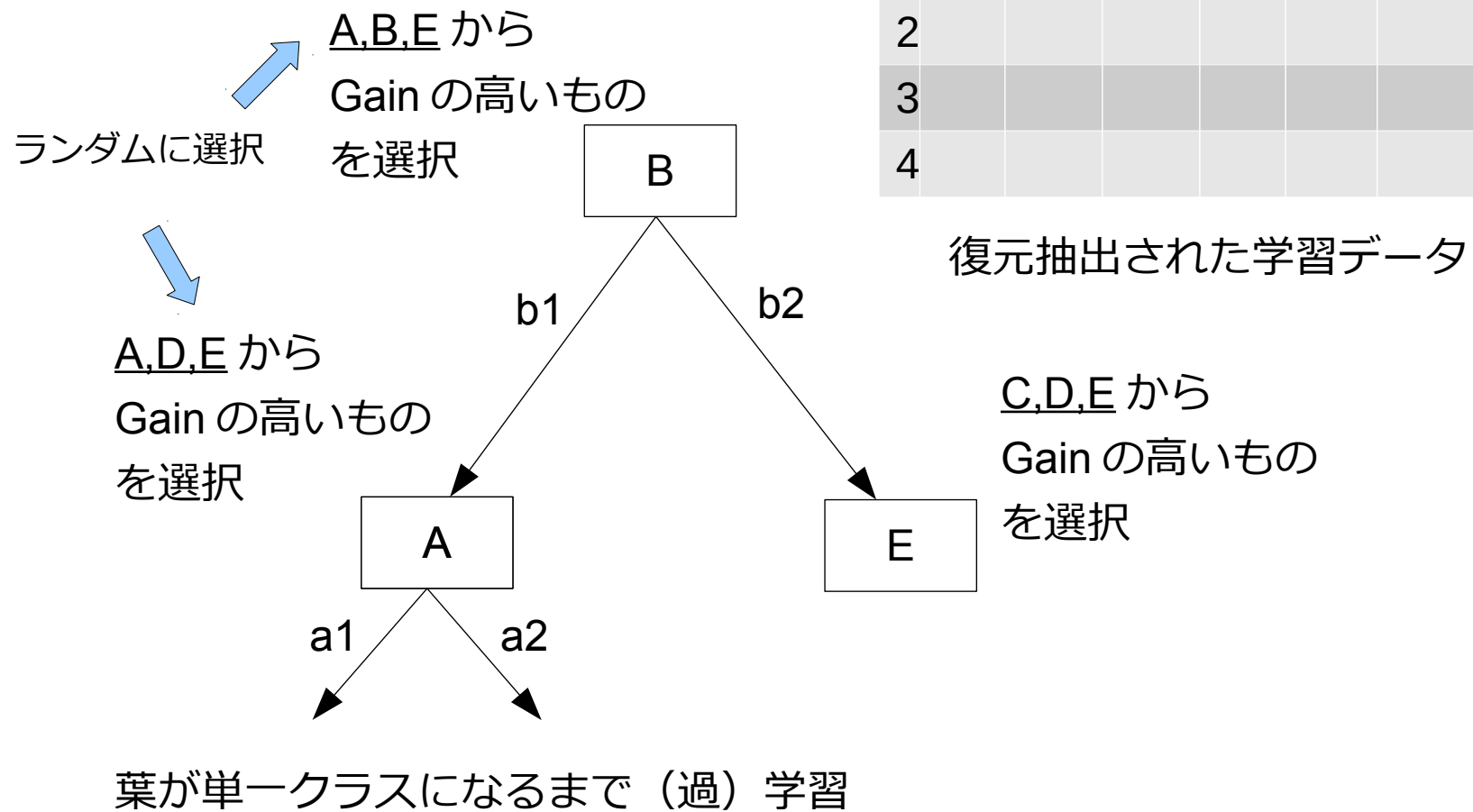
- 訓練例から復元抽出することで、元のデータと同じサイズの独立なデータ集合を作成する。

n 回行って、あるデータが抽出されない確率： $(1 - \frac{1}{n})^n$

$n \rightarrow \infty$ で
約 0.368

- 各々のデータに対して同じアルゴリズムで分類器を作成する
 - アルゴリズムは不安定 (学習データの違いに敏感) の方がよい
 - 例) 枝刈りをしない決定木
- 結果の統合は多数決

9.3 ランダムフォレスト



9.3 ランダムフォレスト

- 特徴
 - バギングと同様、学習データは復元抽出
 - 識別器作成に使用できる特徴をランダムに制限することで、各抽出データ毎に全く異なった識別器ができる
 - 識別器は意図的に過学習させる

Weka の RandomForest

The screenshot shows the Weka Explorer interface with the 'Classify' tab selected. The classifier chosen is 'RandomForest' with parameters '-I 3 -K 0 -S 1 -print -num-slots 1'. The 'Test options' section shows 'Use training set' selected. The 'Result list' on the left shows three entries for 'trees.RandomForest' at different times. The 'Classifier output' pane displays the output of the first two trees. Two red arrows point to the output of the two trees, with a handwritten red note in Japanese.

Classifier output

All the base classifiers:

```
RandomTree
=====
windy = TRUE
| humidity = high : no (5/0)
| humidity = normal
| | outlook = sunny : yes (1/0)
| | outlook = overcast : yes (0/0)
| | outlook = rainy : no (2/0)
windy = FALSE
| outlook = sunny : no (2/0)
| outlook = overcast : yes (1/0)
| outlook = rainy : yes (3/0)

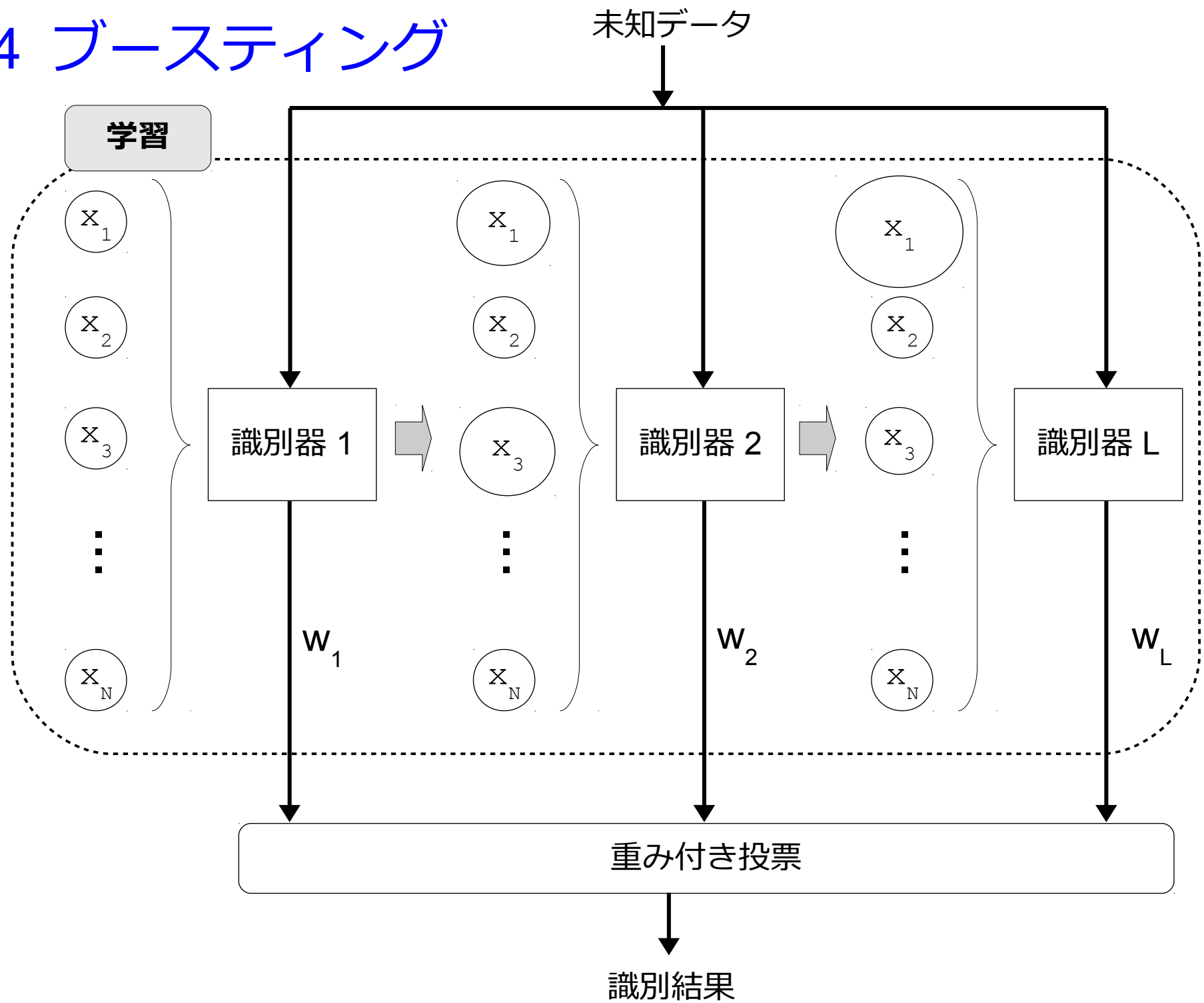
Size of the tree : 11

RandomTree
=====
humidity = high
| outlook = sunny : no (3/0)
| outlook = overcast : yes (2/0)
| outlook = rainy
| | windy = TRUE : no (3/0)
| | windy = FALSE : yes (1/0)
humidity = normal : yes (5/0)

Size of the tree : 8
```

異なった決定木が
学習されている

9.4 ブースティング



9.4 ブースティング

- 特徴
 - 逐次的に相補的な分類器を作成
 - 以前の分類器が誤った事例に重みを付けて次の分類器を学習
 - 学習アルゴリズムが重みに対応していない場合は、重みに比例した数を復元抽出
 - 結果は分類器に対する重み付き投票