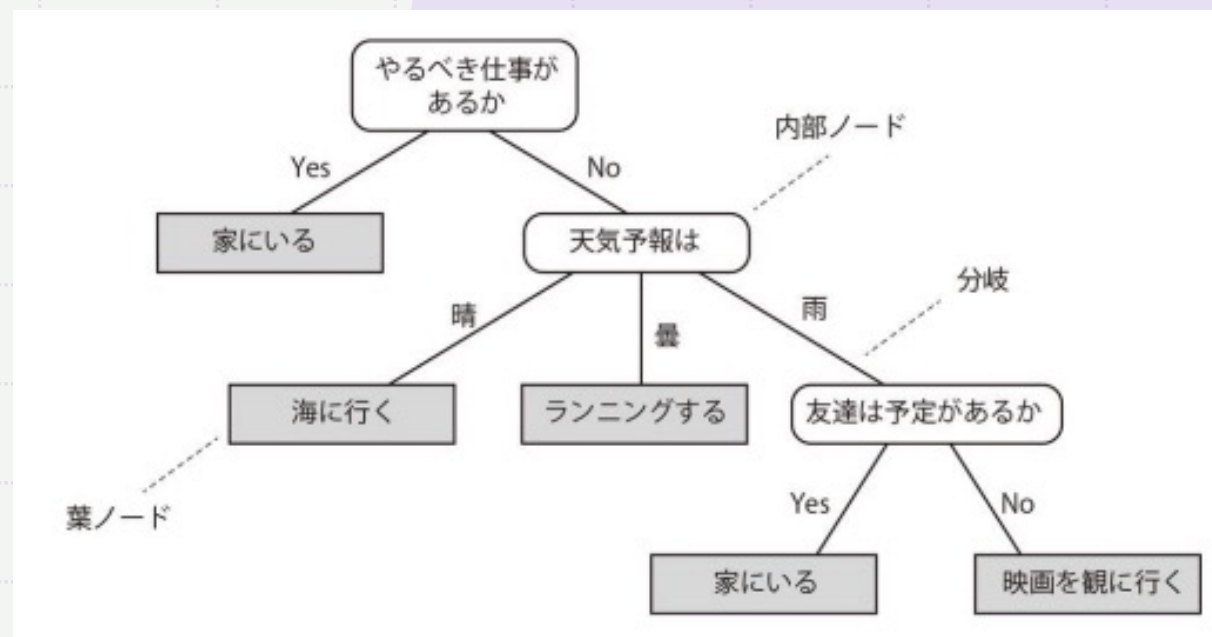


# 決定木:Decision Tree

# 概要

- ある日の行動を決定する。  
「トレーニングデータの特徴量に基づいて質問の答えを学習する。」



ちょっとわかりずらいので、、、

# 概要

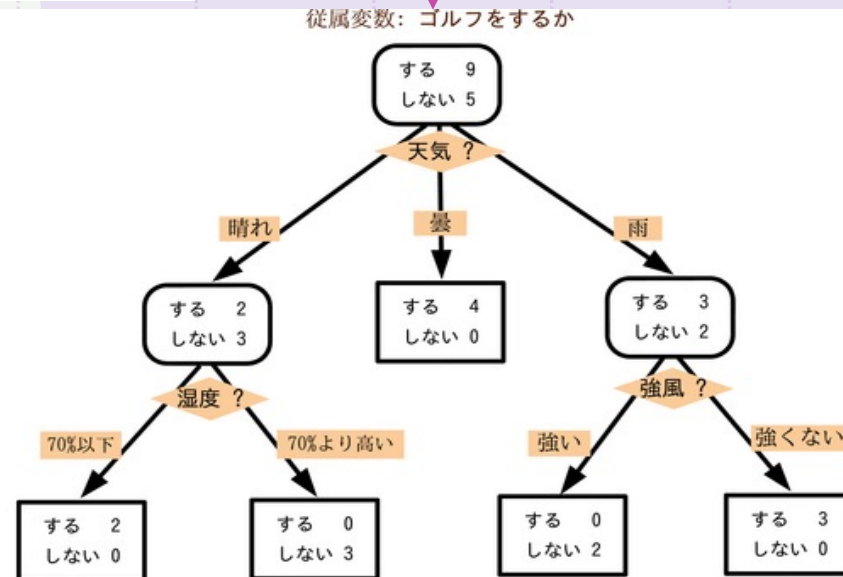
- ある日の行動を決定する。

「トレーニングデータの特徴量に基づいて質問の答えを学習する。」

ゴルフクラブ来場状況

独立変数				従属変数
天気	気温(℃)	湿度(%)	風が強い	ゴルフをする
晴れ	29	85	強くない	しない
晴れ	27	90	強い	しない
曇	28	78	強くない	する
雨	21	96	強くない	する
雨	20	80	強くない	する
雨	18	70	強い	しない
曇	18	65	強い	する
晴れ	22	95	強くない	しない
晴れ	21	70	強くない	する
雨	24	80	強くない	する
晴れ	24	70	強い	する
曇	22	90	強い	する
曇	27	75	強くない	する
雨	22	80	強い	しない

アルゴリズムを考案  
してみましょう



# 概要

Step1 :

集合Cのすべてのデータが同一クラスならクラスノードを作り停止. それ以外なら,  
属性の選択基準により一つの属性 A を選択し判別ノードを作る.

Step2 :

属性Aの属性値によりCを部分集合C1,C2, . . . , Cnに分けてノードを作り属性値の枝を張る.

Step3 :

ノードCi ( $1 \leq i \leq n$ ) について, アルゴリズムを再起的に適用する.

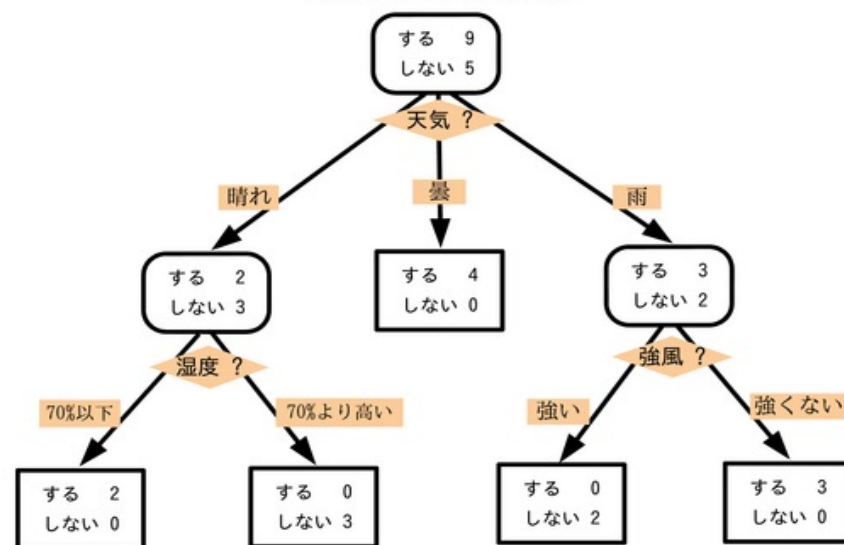
- ある日の行動を決定する。

「トレーニングデータの特徴量に基づいて質問の答えを学習する。」

ゴルフクラブ来場状況

天気	独立変数			従属変数 ゴルフをするか
	気温(℃)	湿度(%)	風が強い	
晴れ	29	85	強くない	しない
晴れ	27	90	強い	しない
曇	28	78	強くない	する
雨	21	96	強くない	する
雨	20	80	強くない	する
雨	18	70	強い	しない
曇	18	65	強い	する
晴れ	22	95	強くない	しない
晴れ	21	70	強くない	する
雨	24	80	強くない	する
晴れ	24	70	強い	する
曇	22	90	強い	する
曇	27	75	強くない	する
雨	22	80	強い	しない

従属変数: ゴルフをするか



# 属性の選択基準

情報利得(Information gain:IG)を最大にする

$$IG(D_p, f) = I(D_p) - \sum_{j=1}^m \frac{N_j}{N_p} I(D_j)$$

ここで、 $f$ は分割を行う特徴量であり、 $D_p$ は親のデータセット、 $D_j$ は $j$ 番目の子ノードのデータセットである。 $I$ は**不純度**(impurity)<sup>※35</sup>を数値化したものであり、 $N_p$ は親ノードのデータ点の総数、 $N_j$ は $j$ 番目の子ノードのデータ点の個数である。このように、情報利得は「親ノードの不純度」と「子ノードの不純度の合計」との差にすぎない。つまり、子ノードの不純度が低いほど、情報利得は大きくなる。

Step1 :

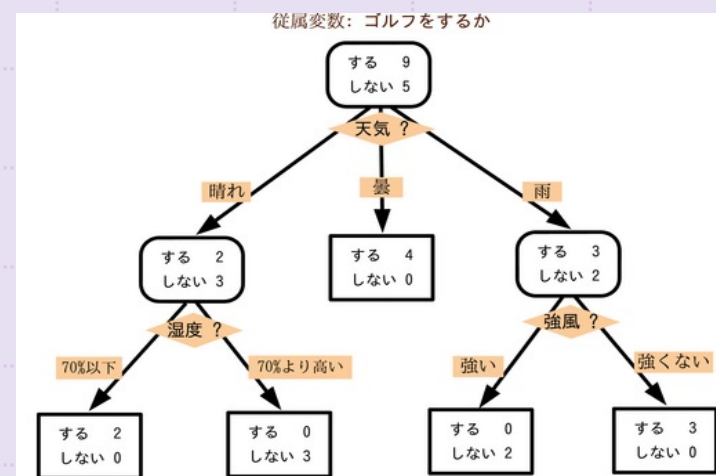
集合Cのすべてのデータが同一クラスならクラスノードを作り停止。それ以外なら、**属性の選択基準**により一つの属性Aを選択し判別ノードを作る。

Step2 :

属性Aの属性値によりCを部分集合C1,C2,..., Cnに分けてノードを作り属性値の枝を張る。

Setp3 :

ノードCi ( $1 \leq i \leq n$ ) について、アルゴリズムを再起的に適用する。



# 属性の選択基準

## 不純度の指標

$$I_H(t) = - \sum_{i=1}^c p(i|t) \log_2 p(i|t)$$

entropy

$$I_G(t) = \sum_{i=1}^c p(i|t)(1 - p(i|t)) = 1 - \sum_{i=1}^c p(i|t)^2$$

Gini 不純度

$$I_E(t) = 1 - \max \{p(i|t)\}$$

分類誤差

$$IG(D_p, f) = I(D_p) - \sum_{j=1}^m \frac{N_j}{N_p} I(D_j)$$

ここで、 $f$ は分割を行う特徴量であり、 $D_p$ は親のデータセット、 $D_j$ は $j$ 番目の子ノードのデータセットである。 $I$ は不純度 (impurity) <sup>※35</sup> を数値化したものであり、 $N_p$ は親ノードのデータ点の総数、 $N_j$ は $j$ 番目の子ノードのデータ点の個数である。このように、情報利得は「親ノードの不純度」と「子ノードの不純度の合計」との差にすぎない。つまり、子ノードの不純度が低いほど、情報利得は大きくなる。

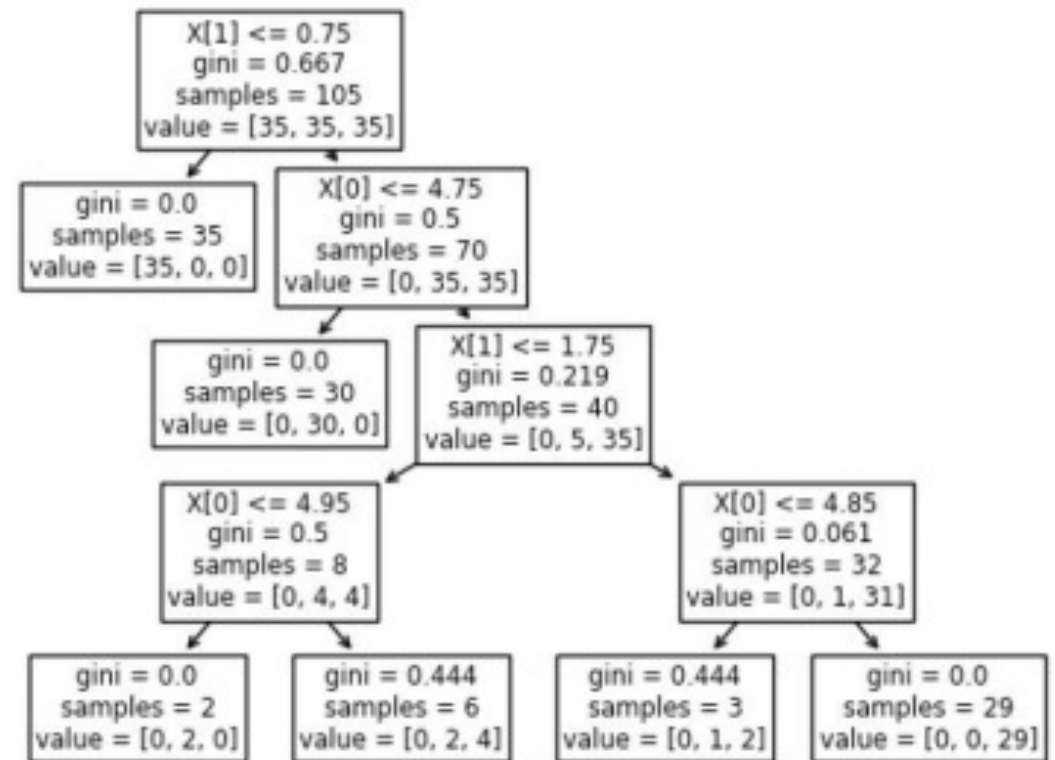
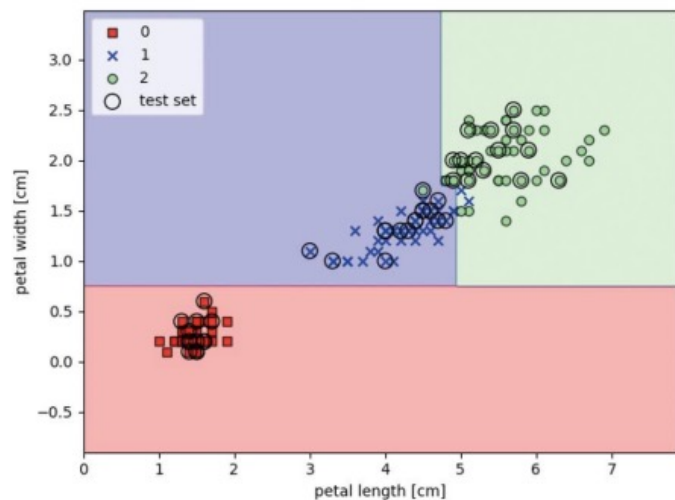
多くの実装で  
エントロピー or Gini不純度  
が用いられている

$p(i|t)$  は、特定のノード  $t$  においてクラス  $i$  に所属しているデータ点の割合を表す。



# 特徴

P.88 ..... 決定木による決定領域



## 意味解釈可能性 (Interpretability)

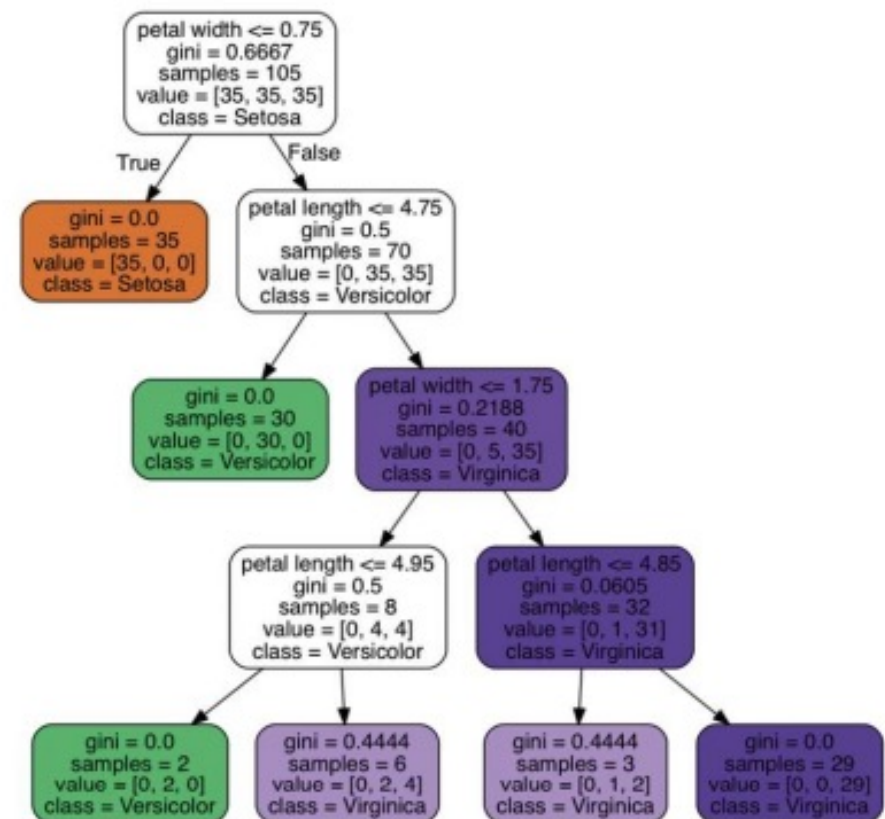
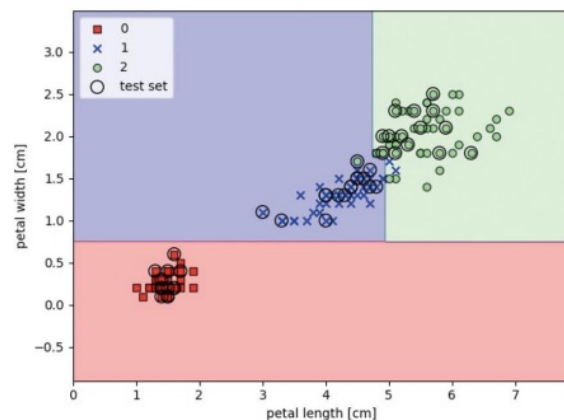
得られた結果の意味解釈ができるかどうか？

# 意味解釈可能性 (Interpretability)

得られた結果の意味解釈ができるかどうか？

P.90 .....分割条件を示す決定木の画像

P.88 .....決定木による決定領域





# Random Forest

決定木で多数決？（アンサンブル学習）

1. サイズ  $n$  のランダムなブートストラップ標本を復元抽出<sup>※41</sup> する（訓練データセットから  $n$  個のデータ点をランダムに選択する）。
2. ブートストラップ標本から決定木を成長させる。各ノードで以下の作業を行う。
  - 2.1  $d$  個の特徴量をランダムに非復元抽出<sup>※42</sup> する。
  - 2.2 たとえば情報利得を最大化することにより、目的関数に従って最適な分割となる特徴量を使ってノードを分割する。
3. 手順 1 ～ 2 を  $k$  回繰り返す。
4. 決定木ごとの予測をまとめ、**多数決**に基づいてクラスラベルを割り当てる。多数決については、第 7 章で詳しく説明する。

# 復元抽出と非復元抽出



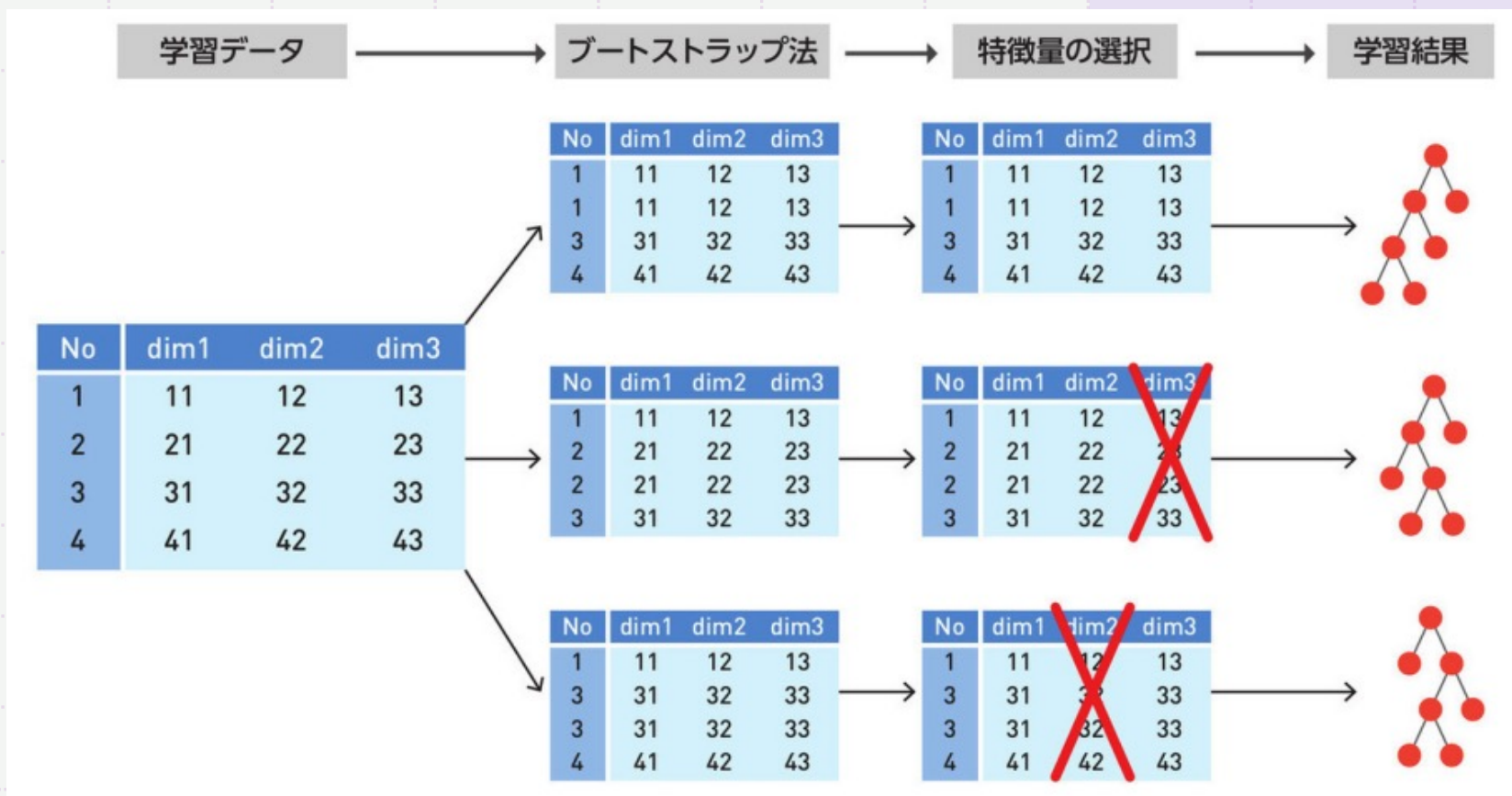
## 復元抽出と非復元抽出

復元抽出 (sampling with replacement) と非復元抽出 (sampling without replacement) という言葉になじみがないかもしれないので、簡単な思考実験をしてみよう。くじ引きゲームで、つぼの中から数字を適当に引くとしよう。0、1、2、3、4 の 5 つの数字の入ったつぼを用意し、1 回につき 1 つの数字を引く。1 回目につぼから特定の数字を引く確率は 5 分の 1 である。さて、非復元抽出では、1 回ごとにつぼに数字を戻さない。したがって、次の回に残りの数字の中から特定の数字を引く確率は前の回によって決まる。たとえば、つぼの中に残っている数字が 0、1、2、4 であるとすれば、次の回に 0 を引く可能性は 4 分の 1 になる。

これに対し、復元抽出では、引いた数字を常につぼに戻すため、特定の数字を引く確率は毎回変化しない。毎回同じ数字を引く可能性もある。つまり、復元抽出では、サンプル (数字) は独立しており、共分散は 0 である。たとえば、数字を 5 回引いた結果は次のようになるかもしれない。

- ランダムな非復元抽出：2、1、3、4、0
- ランダムな復元抽出：1、3、3、4、1

# Random Forestの直感的理解



# 多数決（majority voting）の性能

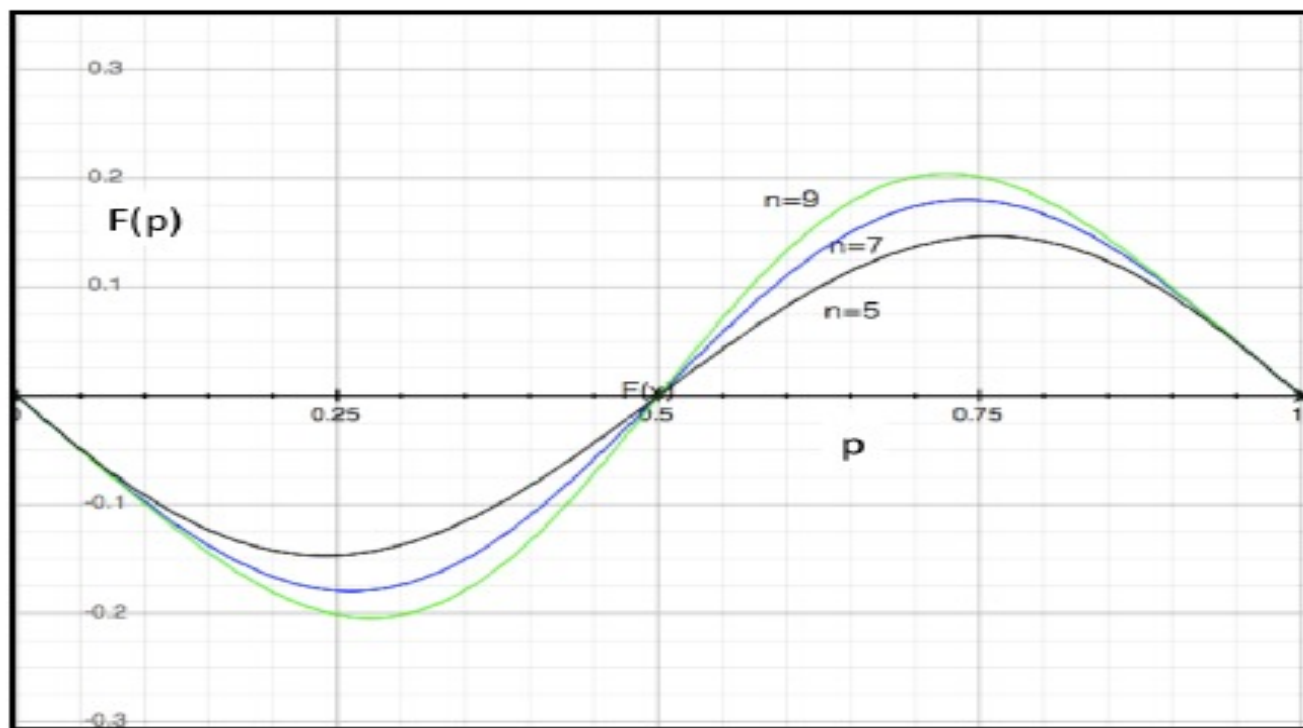


Fig. 2. Transition of consultation performance.

多数決の性能

正答率 50%の人がいくら集まっても、正答率は上がらない。

正答率 100%の人に多数決は必要ない。

一定程度の正答率の高さをもつ人が複数集まることによって多数決が有効に働く

# Theoretical Value of consultation

- Number of members to participate in the consultation

$n$

- Accuracy rate of the answer for each member

$p$

Accuracy rate of the consultation  $Q(p)$ ?

$$Q(p) = p^n + {}_nC_{n-1}p^{n-1}(1-p) + {}_nC_{n-2}p^{n-2}(1-p)^2 + \cdots + {}_nC_{\frac{n+1}{2}}p^{\frac{n+1}{2}}(1-p)^{\frac{n-1}{2}} \quad (1)$$

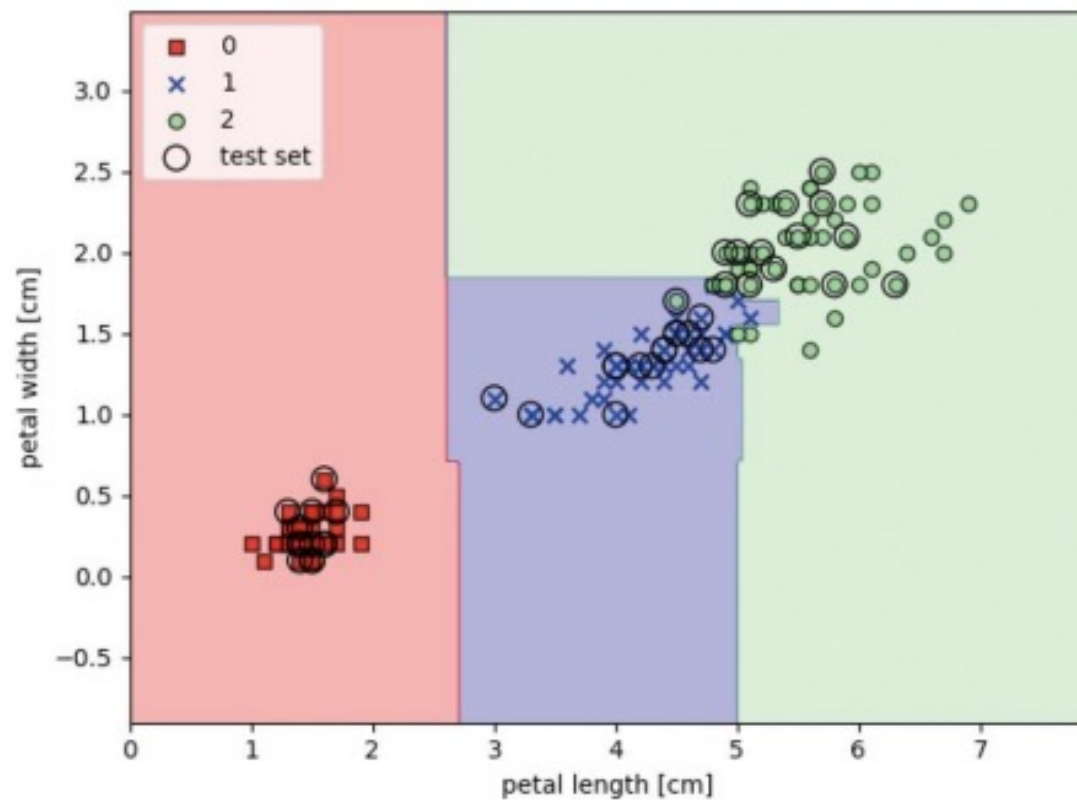
where  $n$  is odd number.

Performance of consultation  $F(p)$  is,

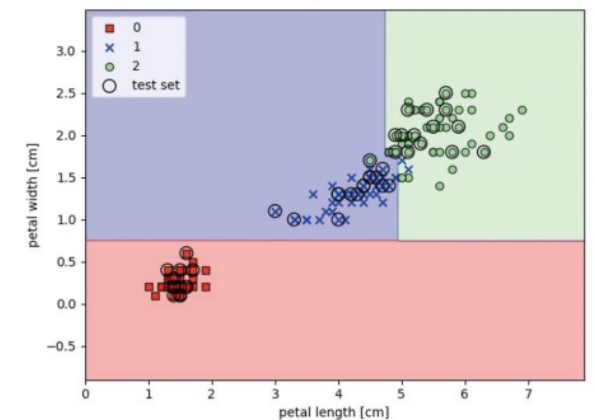
$$F(p) = Q(p) - p$$

# Random Forestの実行結果

P.93 .....ランダムフォレストによる決定領域



P.88 .....決定木による決定領域



実行結果について、バイアスとバリエーションの観点から考察してみよう。