

情報意味論(3) 決定木と過学習

櫻井彰人
慶應義塾大学理工学部

2014/10/19
スライド 40, 42 を修正
74以降順序を修正

1

決定木

- 復習になる方へ、ご容赦を。
- 決定木を道具に、機械学習アルゴリズム共通の課題を説明します
 - 過学習
 - バイアス
 - オッカムの剃刀 etc.

2

機械学習の材料

- 訓練データ・事例、学習データ・事例、
 - 事例=instance=sample
 - ある(一般には未知の)確率分布に従って生成される
 - 訓練データ=独立に生成された事例の集合
- 仮説集合
- 希望の結果との差を示す数値
 - 誤差、誤り率、コスト
- 未知のデータ
 - 学習結果の能力を評価するデータ
 - 訓練データとは異なる

3

機械学習の手段

- 仮説集合
 - (機械学習の)答えの候補=仮説
 - 一つの決定木 = 一つの仮説
 - 作りうる決定木の集合 = 仮説集合
- 学習過程
 - 仮説を一つとり、
 - 訓練データをうまく説明するかどうかを調べ
 - 満足が行く仮説であれば、それを答えとする。
 - 不満であれば、上記を繰り返す。

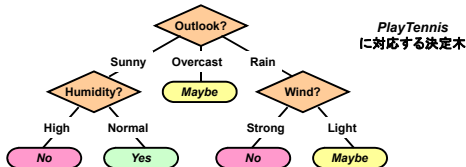
4

決定木 Decision Trees

- 分類器 Classifiers である
 - 事例: 属性 attribute (または特徴 feature) のベクトル+ラベル
- 内節 Internal Nodes: 属性、または属性値のテスト
 - 典型的: 属性 or 等しいかどうかのテスト (e.g., "Wind = ?")
 - その他 不等式や様々なテストが可能
- 枝 Branches: 枝を選ぶ条件である属性値 (テストのときはテストの結果)
 - 一対一対応 (e.g., "Wind = Strong", "Wind = Light")
- 葉 Leaves: 割当てた分類結果 (分類クラスのラベル Class Labels)

Day	Outlook	Temp	Humidity	Wind	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Cloudy	Hot	High	Weak	Yes
D4	Rainy	Mild	High	Weak	Yes
D5	Rainy	Cool	Normal	Weak	Yes
D6	Rainy	Cool	Normal	Strong	No
D7	Cloudy	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rainy	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Cloudy	Mild	High	Strong	Yes
D13	Cloudy	Hot	Normal	Weak	Yes
D14	Rainy	Mild	High	Strong	No

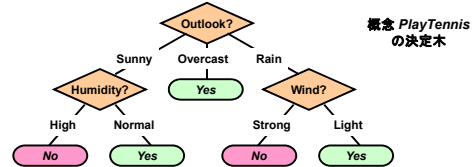
Adapted from Mitchell, 1997



5

決定木はブール関数

- 決定木はブール関数
 - 表現力: 任意のブール関数 (リテラルは属性変数のテスト)が表現可能
 - なぜ?
 - ・ 決定木のあらゆる論理関数は、Disjunctive Normal Form (DNF) でかける
 - ・ 下記の決定木: $(Sunny \wedge Normal-Humidity) \vee Overcast \vee (Rain \wedge Light-Wind)$



- 概念を表現するブール関数の例
 - \wedge, \vee, \oplus (XOR)
 - $(A \wedge B) \vee (C \wedge \neg D \wedge E)$
 - m -of- n

6

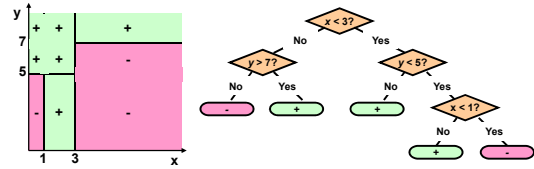
決定木と判別境界

- 事例は、多くの場合、離散属性値で表現される
 - 勿論、連続値も扱う拡張がある
 - 典型的な型
 - 名義・名辞 nominal ($\{\text{red, yellow, green}\}$)
 - 離散化・量子化 quantized ($\{\text{low, medium, high}\}$)
 - 数値の取り扱い
 - 離散化 discretization, ベクトル量子化 **vector quantization**: 閾値を用いて分割する
- ex. U. M. Fayyad and K. B. Irani, Multi-Interval Discretization of Continuous-Valued Attributes for Classification Learning, Proc. 13th IJCAI (1993).

7

決定木と判別境界

- 連続値をそのまま用いることもある
 - しかし、勿論、枝分かれする必要がある
 - 各ノードである閾値より大きい小さいかで分ける
 - 閾値は、学習時に決める
 - この場合、クラス間の判別境界は、超平面で構成されることになる
- 例: 軸並行な方形によって事例空間を分割する



8

学習過程 = 仮説出力過程

- 一般に、学習過程は、仮説を出力する(仮説を仮説空間から選んでくる)過程である。
 - 一回だけ出力する
 - 予め分かる回数だけ出力する
 - 無限回出力する
- 仮説空間が有限な場合(仮説空間に含まれる仮説の個数が有限の場合)
 - 全部を試みて、最適なものを出力する
 - 一部を試みて、その中で、最適なものを出力する
 - 多くの場合、全部を調べる時間がない
- 仮説空間が無限の場合
 - 一部を試してみるしかない
 - 最良なものを見つける方法がない限り、候補を無限に出し続ける
- いずれにせよ、探索順序が問題

9

仮説の選択順序

- 仮説を無限回出力する
 - 見かけ上は、一個の仮説を出力して終了
 - 指定した停止基準を満たした場合、終了とする
- なぜ無限回出力するか?
 - 最適解が、繰り返し計算の極限でしか求まらない
- 求める順序が問題
 - 段々「良く」なって行って欲しい
 - (何らかの基準に従い)好きなところで停止できる
 - 実際は、必ずしもそうではない。
- バイアス
 - どの順序をとるにせよ、一般には、見つけうる解(近似解)に片寄りが生じる。これを学習バイアスという。

10

決定木の学習

- 仮説集合は有限集合
 - 離散変数だけの時。ある葉に至る路上では、同じ属性は2度現れない。
 - 離散変数に対して。閾値は無言個ありうるが、異なる結果を出すものは有限個しかない。
- しかし、全数チェックはできない
 - 膨大すぎる。
- どうする?

11

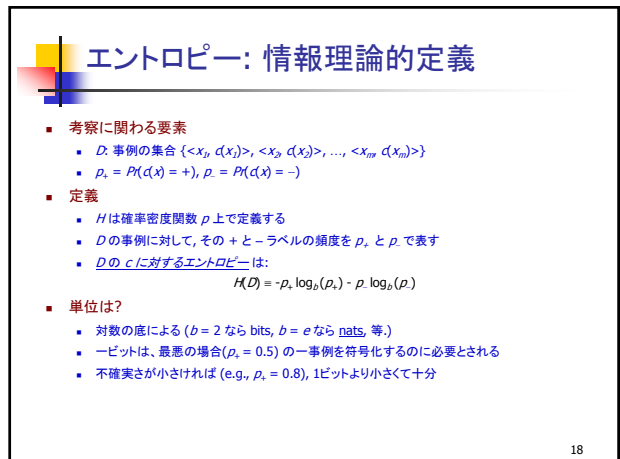
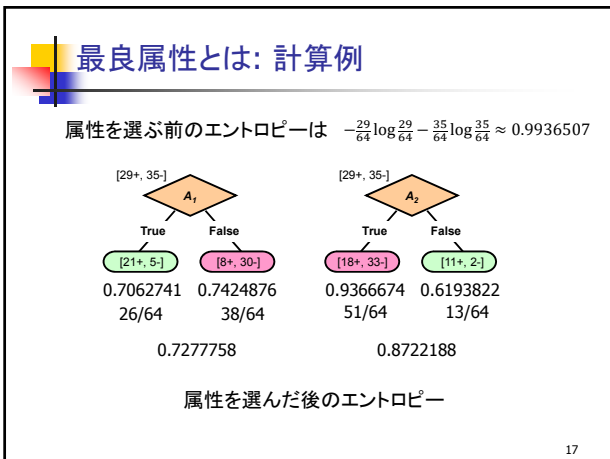
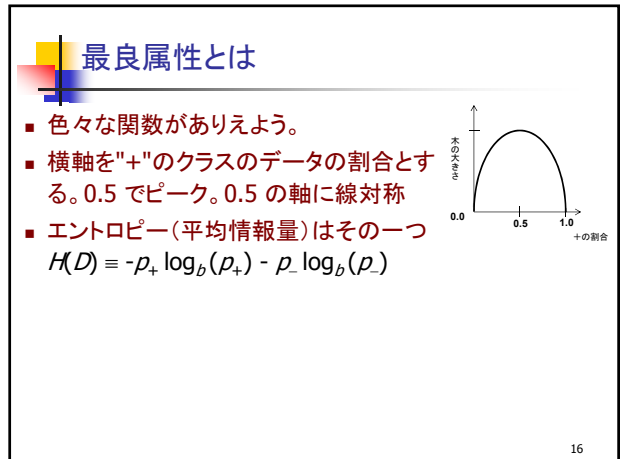
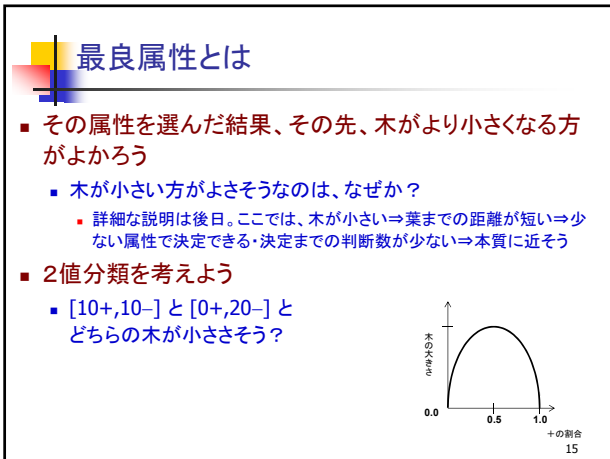
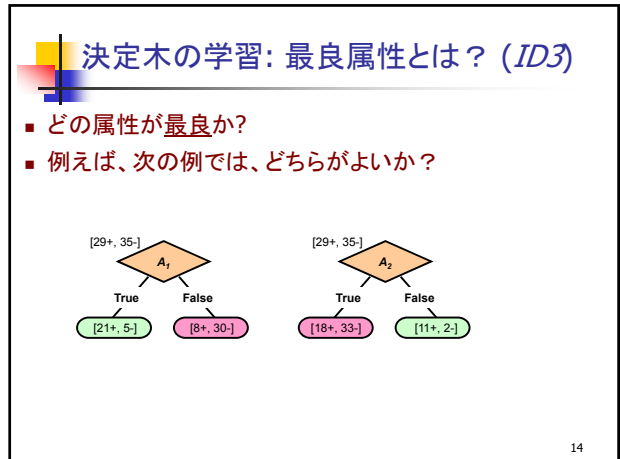
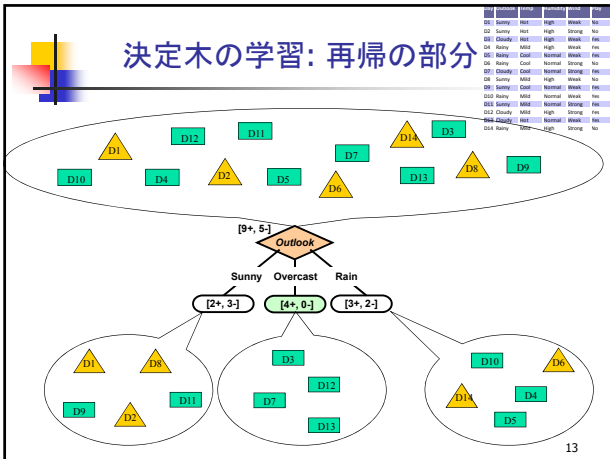
決定木の学習: トップダウン帰納 (ID3)

- アルゴリズム *Build-DT* (*Examples, Attributes*)
 - 部分木に再帰的に適用される
 - Examples: 事例の部分集合, Attributes: 属性の部分集合
- ```

IF Examples の label が同一 THEN RETURN (その label を付した葉節)
ELSE
 IF Attributes が空集合 THEN RETURN (多数派 label を付した葉節)
 ELSE
 最良属性 A を根節として選ぶ。以下で作る木を子とする木を作り、値とする。
 FOR A のそれぞれの値 v
 条件 A = v に対応した、根節からの枝を作成する
 IF {x ∈ Examples | x.A = v} = ∅
 THEN 多数派 label を付した葉節を作成
 ELSE Build-DT({x ∈ Examples | x.A = v}, Attributes - {A})

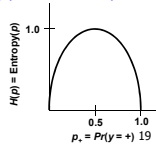
```

12



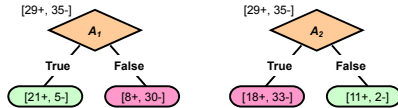
# エントロピー: 直感的説明

- 不確かさ・不明瞭さの尺度; 不確かなほど大きくなる値
  - 計る対象(量)
    - 純粋さ purity: 事例集合が、ただ一つのラベルをもつ状態に、どれだけ近いか
    - 不純さ impurity (乱雑さ disorder): ラベルがまったく分からない状態にどれだけ近いか
  - 尺度: エントロピー
    - 正の相関: 不純さ impurity, 不確かさ uncertainty, 不規則さ irregularity, 驚き surprise
    - 負の相関: 純粋さ purity, 確かさ certainty, 規則性 regularity, 冗長さ redundancy
- 例
  - 簡単のため,  $H = \{0, 1\}$ , ある分布  $P(x)$  に従うと仮定
    - (2個より多い)離散的なクラスラベルでも同様
    - さらに連続確率変数でもよい: 微分エントロピー (和を積分にしただけ)
  - $y$  に関して最も純粋: 次のいずれかの場合
    - $P(x=0) = 1, P(x=1) = 0$
    - $P(x=1) = 1, P(x=0) = 0$
  - 純粋さが最も少ない確率分布は?
    - $P(x=0) = 0.5, P(x=1) = 0.5$
    - 最大: 不純さ/不確かさ/不規則性/驚き
    - エントロピーの性質: 凹関数(「上向きに凸」)



# 情報量増分: 情報理論的定義

- 属性値に基づく分割
  - 復習:  $D$  の分割 partition は、和集合が  $D$  となるような排他的部分集合の集合
  - 目標: 属性  $A$  の属性値に基づく分割により削減される不確かさ・不純性を計る
- 定義
  - 属性  $A$  に関する  $D$  の情報量増分 は、 $A$  を用いた分割によるエントロピー減少分の期待値:
 
$$Gain(D, A) = H(D) - \sum_{v \in values(A)} \frac{|D_v|}{|D|} H(D_v) = \frac{1}{|D|} \left( |D| H(D) - \sum_{v \in values(A)} |D_v| H(D_v) \right)$$
  - 但し  $D_v$  は  $\{x \in D \mid x.A = v\}$ , すなわち,  $D$  中の事例で属性  $A$  の値が  $v$  であるものの集合
  - 補足:  $A$  による分割によって生じる部分集合  $D_v$  の大きさに従ってエントロピーの大きさを調整
    - エントロピー値は、「集合の要素一個あたりの情報量」となっているため
- どちらの属性を使うのがいい?



# 例

- 概念 PlayTennis 用の訓練事例

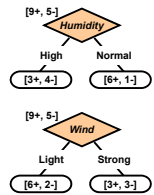
| Day | Outlook  | Temperature | Humidity | Wind   | PlayTennis? |
|-----|----------|-------------|----------|--------|-------------|
| 1   | Sunny    | Hot         | High     | Light  | No          |
| 2   | Sunny    | Hot         | High     | Strong | No          |
| 3   | Overcast | Hot         | High     | Light  | Yes         |
| 4   | Rain     | Mild        | High     | Light  | Yes         |
| 5   | Rain     | Cool        | Normal   | Light  | Yes         |
| 6   | Rain     | Cool        | Normal   | Strong | No          |
| 7   | Overcast | Cool        | Normal   | Strong | Yes         |
| 8   | Sunny    | Mild        | High     | Light  | No          |
| 9   | Sunny    | Cool        | Normal   | Light  | Yes         |
| 10  | Rain     | Mild        | Normal   | Light  | Yes         |
| 11  | Sunny    | Mild        | Normal   | Strong | Yes         |
| 12  | Overcast | Mild        | High     | Strong | Yes         |
| 13  | Overcast | Hot         | Normal   | Light  | Yes         |
| 14  | Rain     | Mild        | High     | Strong | No          |

- $ID3 = Build\text{-}DT$  但し  $Gain(\cdot)$  を使用
- $ID3$  の動きを追ってみよう

# ID3による PlayTennis 決定木作成 [1]

- 根節の属性を選ぶ

| Day | Outlook  | Temperature | Humidity | Wind   | PlayTennis? |
|-----|----------|-------------|----------|--------|-------------|
| 1   | Sunny    | Hot         | High     | Light  | No          |
| 2   | Sunny    | Hot         | High     | Strong | No          |
| 3   | Overcast | Hot         | High     | Light  | Yes         |
| 4   | Rain     | Mild        | High     | Light  | Yes         |
| 5   | Rain     | Cool        | Normal   | Light  | Yes         |
| 6   | Rain     | Cool        | Normal   | Strong | No          |
| 7   | Overcast | Cool        | Normal   | Strong | Yes         |
| 8   | Sunny    | Mild        | High     | Light  | No          |
| 9   | Sunny    | Cool        | Normal   | Light  | Yes         |
| 10  | Rain     | Mild        | Normal   | Light  | Yes         |
| 11  | Sunny    | Mild        | Normal   | Strong | Yes         |
| 12  | Overcast | Mild        | High     | Strong | Yes         |
| 13  | Overcast | Hot         | Normal   | Light  | Yes         |
| 14  | Rain     | Mild        | High     | Strong | No          |



- 事前(無条件)分布:  $9+, 5-$

- $H(D) = -(9/14) \log(9/14) - (5/14) \log(5/14) \text{ bits} = 0.94 \text{ bits}$
- $H(D, Humidity = High) = -(3/7) \log(3/7) - (4/7) \log(4/7) = 0.985 \text{ bits}$
- $H(D, Humidity = Normal) = -(6/7) \log(6/7) - (1/7) \log(1/7) = 0.592 \text{ bits}$
- $Gain(D, Humidity) = 0.94 - ((3/7) * 0.985 + (1/7) * 0.592) = 0.151 \text{ bits}$
- 同様に,  $Gain(D, Wind) = 0.94 - ((8/14) * 0.811 + (6/14) * 1.0) = 0.048 \text{ bits}$

$$Gain(D, A) = H(D) - \sum_{v \in values(A)} \frac{|D_v|}{|D|} H(D_v)$$

# ID3による PlayTennis 決定木作成 [2]

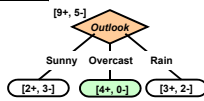
- 根節の属性を選ぶ

| Day | Outlook  | Temperature | Humidity | Wind   | PlayTennis? |
|-----|----------|-------------|----------|--------|-------------|
| 1   | Sunny    | Hot         | High     | Light  | No          |
| 2   | Sunny    | Hot         | High     | Strong | No          |
| 3   | Overcast | Hot         | High     | Light  | Yes         |
| 4   | Rain     | Mild        | High     | Light  | Yes         |
| 5   | Rain     | Cool        | Normal   | Light  | Yes         |
| 6   | Rain     | Cool        | Normal   | Strong | No          |
| 7   | Overcast | Cool        | Normal   | Strong | Yes         |
| 8   | Sunny    | Mild        | High     | Light  | No          |
| 9   | Sunny    | Cool        | Normal   | Light  | Yes         |
| 10  | Rain     | Mild        | Normal   | Light  | Yes         |
| 11  | Sunny    | Mild        | Normal   | Strong | Yes         |
| 12  | Overcast | Mild        | High     | Strong | Yes         |
| 13  | Overcast | Hot         | Normal   | Light  | Yes         |
| 14  | Rain     | Mild        | High     | Strong | No          |

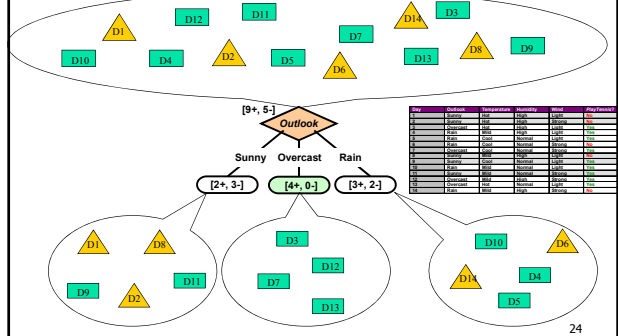
- $Gain(D, Humidity) = 0.151 \text{ bits}$
- $Gain(D, Wind) = 0.048 \text{ bits}$
- $Gain(D, Temperature) = 0.029 \text{ bits}$
- $Gain(D, Outlook) = 0.246 \text{ bits}$

- 次の属性を選ぶ(部分木の根節)

- (葉への道の上で)属性を使いきるか葉の純粋度 = 100% になるまで続ける
- 純粋度 = 100% は、一つのラベルしかないということ
- ところで  $Gain(D, A) < 0$  となりうるか?



# 再掲: 決定木の学習: 再帰的繰り返し



## ID3による PlayTennis 決定木作成 [3]

- 次の属性の選択 (部分木の根節)

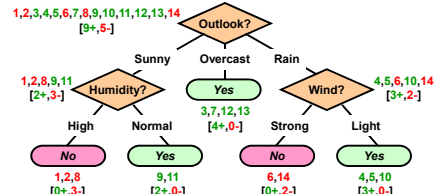
| Day | Outlook  | Temperature | Humidity | Wind   | PlayTennis? |
|-----|----------|-------------|----------|--------|-------------|
| 1   | Sunny    | Hot         | High     | Light  | No          |
| 2   | Sunny    | Hot         | High     | Strong | No          |
| 3   | Overcast | Hot         | High     | Light  | Yes         |
| 4   | Rain     | Mild        | High     | Light  | Yes         |
| 5   | Rain     | Cool        | Normal   | Light  | Yes         |
| 6   | Rain     | Cool        | Normal   | Strong | No          |
| 7   | Overcast | Cool        | Normal   | Strong | Yes         |
| 8   | Sunny    | Mild        | High     | Light  | No          |
| 9   | Sunny    | Cool        | Normal   | Light  | Yes         |
| 10  | Rain     | Mild        | Normal   | Light  | Yes         |
| 11  | Sunny    | Mild        | Normal   | Strong | Yes         |
| 12  | Overcast | Mild        | High     | Strong | Yes         |
| 13  | Overcast | Hot         | Normal   | Light  | Yes         |
| 14  | Rain     | Mild        | High     | Strong | No          |

- 約束:  $0 \log(0/a) = 0$
- $Gain(D_{Sunny}, Humidity) = 0.97 - (3/5) * 0 - (2/5) * 0 = 0.97 \text{ bits}$
- $Gain(D_{Sunny}, Wind) = 0.97 - (2/5) * 1 - (3/5) * 0.92 = 0.02 \text{ bits}$
- $Gain(D_{Sunny}, Temperature) = 0.57 \text{ bits}$
- トップダウン再帰
  - 離散値属性しかないなら,  $O(n)$  回分割をすれば終了 ( $n$  は属性数)
  - 木のレベルそれぞれで, 訓練データを一回スキャン (なぜ?)

25

## ID3による PlayTennis 決定木作成 [4]

| Day | Outlook  | Temperature | Humidity | Wind   | PlayTennis? |
|-----|----------|-------------|----------|--------|-------------|
| 1   | Sunny    | Hot         | High     | Light  | No          |
| 2   | Sunny    | Hot         | High     | Strong | No          |
| 3   | Overcast | Hot         | High     | Light  | Yes         |
| 4   | Rain     | Mild        | High     | Light  | Yes         |
| 5   | Rain     | Cool        | Normal   | Light  | Yes         |
| 6   | Rain     | Cool        | Normal   | Strong | No          |
| 7   | Overcast | Cool        | Normal   | Strong | Yes         |
| 8   | Sunny    | Mild        | High     | Light  | No          |
| 9   | Sunny    | Cool        | Normal   | Light  | Yes         |
| 10  | Rain     | Mild        | Normal   | Light  | Yes         |
| 11  | Sunny    | Mild        | Normal   | Strong | Yes         |
| 12  | Overcast | Mild        | High     | Strong | Yes         |
| 13  | Overcast | Hot         | Normal   | Light  | Yes         |
| 14  | Rain     | Mild        | High     | Strong | No          |



26

## 適用範囲を広げるには

- これまでのアルゴリズムでの仮定. その克服
  - 離散 出力
    - 実数値出力も可能
    - Regression trees [Breiman et al, 1984]
  - 離散 入力
    - 量子化の方法あり
    - 内節の等式テストの代わりに不等式を使用する (以前の方形の例)
- 規模の拡大
  - 大規模データベース (VLDB) からの知識発見やデータマイニング (KDD) では重要
  - 長所: 多くの事例を対象とするよいアルゴリズムあり
  - 弱点: あまりに多い属性を扱うのは難しい
- あとと助かる他の耐性
  - ノイズのあるデータ (分類ノイズ classification noise = ラベルの間違い; 属性ノイズ attribute noise = 不正確または低精度のデータ) への耐性
  - 欠測値への耐性

27

## Wekaでの例

- Wekaの紹介スライドにあるとおりです。

28

## Rにおける決定木

- Rには, 決定木関連のパッケージとして, tree, rpart, 及び rpart を多変量回帰木 (multivariate regression trees) に拡張させた mvpart がある。

29

## 分類木の例 (tree)

```
library(tree)
data(iris)
(iris.tr <- tree(Species ~., data = iris))
plot(iris.tr, type = "u"); text(iris.tr)
(iris.tr1 <- snip.tree(iris.tr, nodes = c(12, 7)))
plot(iris.tr1, type = "u"); text(iris.tr1)
```

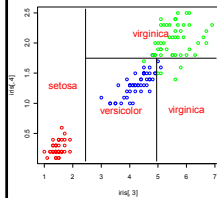
deviance =  $-2 \sum n_i \log p_i$

- root 150 329.600 setosa ( 0.33333 0.33333 0.33333 )
- Petal.Length < 2.45 50 0.000 setosa ( 1.00000 0.00000 0.00000 ) \*
- Petal.Length > 2.45 100 138.600 versicolor ( 0.00000 0.50000 0.50000 )
- Petal.Length < 1.75 54 33.320 versicolor ( 0.00000 0.90741 0.09259 )
- Petal.Length < 4.95 46 9.721 versicolor ( 0.00000 0.97917 0.02083 )
- Sepal.Length < 5.15 5 5.004 versicolor ( 0.00000 0.80000 0.20000 ) \*
- Sepal.Length > 5.15 43 0.000 versicolor ( 0.00000 1.00000 0.00000 ) \*
- Petal.Length > 4.95 6 7.638 virginica ( 0.00000 0.33333 0.66667 ) \*
- Petal.Length > 1.75 46 9.635 virginica ( 0.00000 0.02174 0.97826 ) \*
- Petal.Length < 4.95 6 5.407 virginica ( 0.00000 0.16667 0.83333 ) \*
- Petal.Length > 4.95 40 0.000 virginica ( 0.00000 0.00000 1.00000 ) \*

30

## 分類木の例 (tree)

```
iris.label<-c("S", "C", "V")[iris[, 5]]
plot(iris[,3],iris[,4],type="n")
text(iris[,3],iris[,4],labels=iris.label)
partition.tree(iris.tr1,add=T,col=2,cex=1.5)
```



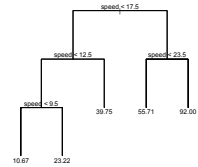
```
iris.color<-c("red","blue","green")[iris[, 5]]
plot(iris[,3],iris[,4],col=iris.color)
partition.tree(iris.tr1,add=T,col=2,cex=1.5)
```

31

## 回帰木の例 (tree)

```
> library(tree)
> data(cars)
> cars.tr<-tree(dist~speed,data=cars)
> print(cars.tr)
node), split, n, deviance, yval
* denotes terminal node
1) root 50 32540.0 42.98
2) speed < 17.5 31 8307.0 29.32
4) speed < 12.5 15 1176.0 18.20
8) speed < 9.5 6 277.3 10.67 *
9) speed > 9.5 9 331.6 23.22 *
5) speed > 12.5 16 3535.0 39.75 *
3) speed > 17.5 19 9016.0 65.26
6) speed < 23.5 14 2847.0 55.71 *
7) speed > 23.5 5 1318.0 92.00 *
> plot(cars.tr,type="u")
> text(cars.tr)
> plot(cars.tr,type="u")
> text(cars.tr)
>
```

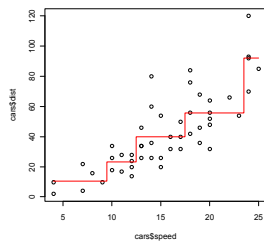
```
Library(tree)
data(cars)
cars.tr<-tree(dist~speed,data=cars)
print(cars.tr)
plot(cars.tr,type="u")
text(cars.tr)
plot(cars.tr,type="u")
text(cars.tr)
```



32

## 回帰木の例 (tree)

```
> plot(cars$speed,cars$dist)
> partition.tree(cars.tr,add=T,col=2)
```

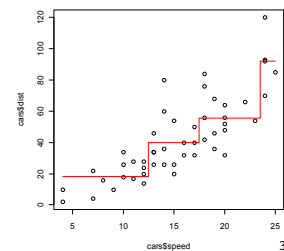
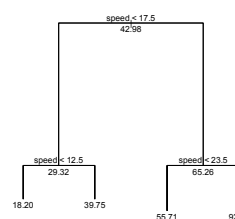


33

## 回帰木の例 (tree)

```
(cars.tr1<-prune.tree(cars.tr,best=4))
plot(cars.tr1); text(cars.tr1,all=T)

plot(cars$speed,cars$dist)
partition.tree(cars.tr1,add=T,col=2)
```



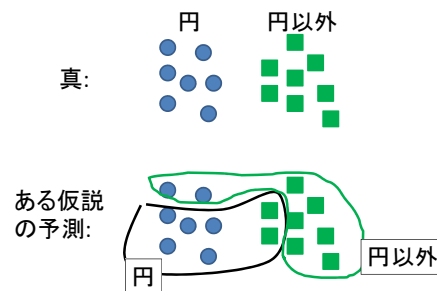
34

## 仮説の評価

- 仮説 (モデルとも言います。文脈によります) の良さ・悪さを評価する。
  - 仮説は、使う
  - 何らかの意味で「精度」や「信頼性」の高い仮説を用いたい。
  - どんな評価方法があるのか？

35

## 学習結果の評価 PrecisionとRecallの前に



36

### TP, TN, FP, FN

真:

ある仮説の予測:

TP: True Positive  
 TN: True Negative  
 FP: False Positive  
 FN: False Negative

結果 仮説による予測

円 TP FP TN 円以外

37

### Confusion matrix

|        |   | 真値                            |                                                |                                  |
|--------|---|-------------------------------|------------------------------------------------|----------------------------------|
|        |   | P                             | N                                              |                                  |
| 仮説の予測値 | P | TP<br>(True Positive)         | FP<br>(False Positive)                         | Precision = $\frac{TP}{TP + FP}$ |
|        | N | FN<br>(False Negative)        | TN<br>(True Negative)                          |                                  |
|        |   | Recall = $\frac{TP}{TP + FN}$ | Accuracy = $\frac{TP + TN}{TP + FP + TN + FN}$ |                                  |

38

### 両者のTradeoff と F-measure

$$F = \frac{1}{\frac{1}{2} \left( \frac{1}{precision} + \frac{1}{recall} \right)}$$

39

### Confusion matrix

|        |   | 真値                            |                                                |                                  |
|--------|---|-------------------------------|------------------------------------------------|----------------------------------|
|        |   | P                             | N                                              |                                  |
| 仮説の予測値 | P | TP<br>(True Positive)         | FP<br>(False Positive)                         | Precision = $\frac{TP}{TP + FP}$ |
|        | N | FN<br>(False Negative)        | TN<br>(True Negative)                          |                                  |
|        |   | Recall = $\frac{TP}{TP + FN}$ | Accuracy = $\frac{TP + TN}{TP + FP + TN + FN}$ |                                  |

$1 - \alpha = \text{specificity} = \text{TNR} = \frac{TN}{FP + TN}$   
 $\alpha = \text{FPR} = \frac{FP}{FP + TN}$

$\beta = \text{FNR} = \frac{FN}{FN + TP}$   
 $1 - \beta = \text{sensitivity} = \text{TPR} = \frac{TP}{FN + TP}$

第一種の過誤  
 第二種の過誤  
 帰無仮説は、「陽性でない」(陽性であることを示したくないから)  
 第一種の過誤=棄却した(陽性だと言った)が、それは誤り  
 第二種の過誤=受理した(陰性だと言った)が、それは誤り

AUC = 0.97  
TPR = 0.98  
FPR = 0.27

AUC = 0.97  
Sensitivity = 0.97  
Specificity = 0.73

40

### ROC curve

- Receiver operating characteristics
  - "ROC"という用語はレーダが開発された当初、操作盤上にあつたノブの名
    - <http://www.math-koubou.jp/stata/files/r12/est006.pdf>

41

### ROC curve

|        |   | 真値                            |                                                |                                  |
|--------|---|-------------------------------|------------------------------------------------|----------------------------------|
|        |   | P                             | N                                              |                                  |
| 仮説の予測値 | P | TP<br>(True Positive)         | FP<br>(False Positive)                         | Precision = $\frac{TP}{TP + FP}$ |
|        | N | FN<br>(False Negative)        | TN<br>(True Negative)                          |                                  |
|        |   | Recall = $\frac{TP}{TP + FN}$ | Accuracy = $\frac{TP + TN}{TP + FP + TN + FN}$ |                                  |

ROC curve ("Receiver Operating Characteristics")

$\alpha = \text{FPR} = \frac{FP}{FP + TN}$   
 $1 - \alpha = \text{specificity}$

AUC = 0.97  
TPR = 0.98  
FPR = 0.27

AUC = 0.97  
Sensitivity = 0.97  
Specificity = 0.73

42

## 訓練誤差と汎化誤差

- 訓練誤差・誤率: 訓練データ学習データに対する、仮説出力値の、真の出力値に対する誤差・誤率
  - 簡単に数えることができる。
- 未知データに対する誤差・誤率: (訓練データと同じ母集団から、同じ方法で抽出した)未知データに対する、仮説出力値の、真の出力値に対する誤差・誤率。テストエラーともいう。

43

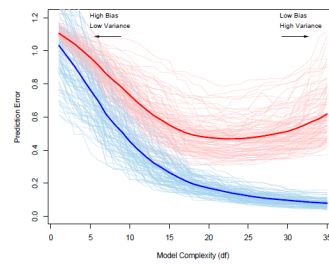


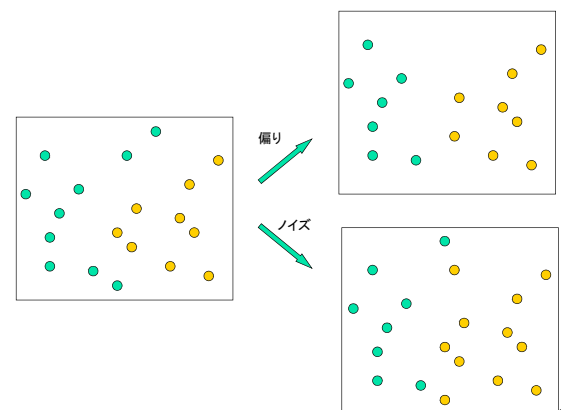
FIGURE 7.1. Behavior of test sample and training sample error as the model complexity is varied. The light blue curves show the training error  $E_T$ , while the light red curves show the conditional test error  $E_{TT}$  for 100 training sets of size 50 each, as the model complexity is increased. The solid curves show the expected test error  $E[E_{TT}]$  and the expected training error  $E[E_T]$ .

Elements of Statistical Learning 44

## 過学習

- over-learning とか over-training と呼ばれる
- 学習すべきでないものまで、学習してしまう
- 学習すべきでないもの
  - 学習データに含まれる偏り
    - 無限集合(真の概念が含む事例は無限個ある)の有限部分集合であるため、かならず、偏りがある。
  - 学習データに含まれる誤り
    - 現実データにはノイズがある。分類クラスにも属性値にもノイズは存在する。
- 学習してしまう
  - 学習能力が高いから
  - 調節可能なパラメータ数が多い

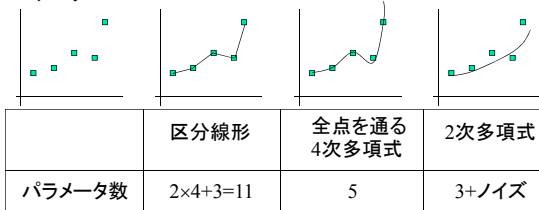
45



46

## 再掲: 関数近似の例(ノイズ)

データ



多分過学習?

多分過学習

47

## 関数近似の applet

- 分かりきったことかもしれませんが、デモプログラムを用いて実験してみると、よく分かります。

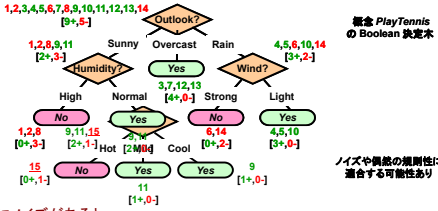
<http://www.mste.uiuc.edu/users/exner/java/f/leastquares/>

48



## 決定木における過学習: 例

- 既出例: 帰納した木



- 訓練事例にノイズがあると

- 事例 15: <Sunny, Hot, Normal, Strong, >
  - この例は実は noisy である。すなわち、正しいラベルは +
  - 以前に作成した木は、これを、誤分類する
- 決定木はどのように更新されるべきか (incremental learning を考える)?
- 新しい仮説  $h = T$  の性能は  $h = T$  より悪く なる と予想される (ノイズに騙されているから!)<sup>49</sup>

## 帰納学習における過学習

- 定義

- 仮説  $h$  が訓練データ集合  $D$  を過学習する (〜に overfits する) というのは、もし他の仮説  $h'$  で  $error_D(h) < error_D(h')$  であるが  $error_{test}(h) > error_{test}(h')$  となるものがあること
- 原因: 訓練事例が少なすぎる (あまりにも少ないデータに基づく判断); ノイズ; 単なる偶然

- 過学習に対応するには?

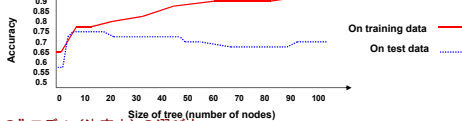
- 予防策
  - 過学習が発生する前に対応する
  - 重要な relevant 属性 (i.e., モデルにとって有用なものを)のみを用いる
    - 注意: 鶏と卵の問題; 重要性 relevance を予測する尺度が必要
- 回避策
  - 問題が起こりそうときに、脇をすりぬける
  - テスト集合を確保しておき、仮説  $h$  がその上で悪くなりそうときに、学習を停止する
- 泳がせ策
  - 問題は発生するにまかせ、発生を検出し、その後回復する
  - モデルを作ってみて、過学習に寄与する要素を発見・除去する (刈る prune)

50

## 決定木学習: 過学習の予防と回避

- 過学習にどう立ち向かうか?

- 予防策
  - 重要な属性を選択 (i.e., 決定木では有用)
  - 重要な予測: 属性を filter する, または 部分集合選択
- 回避策
  - 検証集合 validation set を抜き出しおき,  $h$  の予測精度 がそれに対し悪化したら学習を停止



- “最良の”モデル (決定木) の選び方

- 上述: 性能を測定するにあたって、訓練データとそれとは別の検証データを用いる
- 別法: 最小記述長 Minimum Description Length (MDL):
  - 最小化せよ:  $size(h = T) + size(\text{誤分類 misclassifications } (h = T))$

51

## 決定木学習: 過学習の予防と回避

- 基本的なアプローチが2つある

- Pre-pruning (回避): 木を作成する途中で木の生長を止める。信頼性ある選択をするに十分なデータはないと判断されたとき
- Post-pruning (回復): 木を一杯まで構築し節を削除する。削除するのは、十分な証拠がないとみなされるもの

- 枝刈りすべき部分木を評価する方法

- Cross-validation: 仮説の有用性を評価するために、予めデータをとりおく (Mitchell 第4章)
- 統計的検定: 観測された規則性が偶然起こったものとして捨ててよいかどうかをテストする (Mitchell 第5章)
- 最小記述長 Minimum Description Length (MDL)
  - 仮説  $T$  の複雑度の増加分は、単に(説明しようとしているデータの)例外を記憶するのに必要な記述量より大きいか小さいか?
    - Tradeoff: モデルを記述する versus 残差誤差を記述する

52

## Reduced-Error Pruning

- Post-Pruning, Cross-Validation Approach

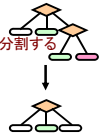
- 所与のデータを 訓練データ training set と 検証データ Validation set に分割する

- 関数  $Prune(T, node)$

- 引数  $node$  を根節とする部分木を除去
- 引数  $node$  を葉節とする (そこにある事例には多数派のラベルを付与)

- アルゴリズム Reduced-Error-Pruning ( $D$ )

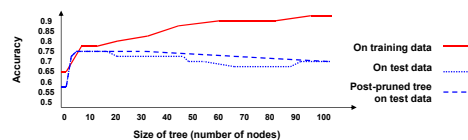
- $D$  を分割する。  $D_{train}$  (訓練 training / "growing"),  $D_{validation}$  (検証 validation / "pruning")
- $D_{train}$  に  $ID3$  を適用して、完全な木  $T$  を作る
- UNTIL  $D_{validation}$  で計測した精度が悪化する DO
- FOR  $T$  中のそれぞれの内節 candidate
- Temp[candidate] ← Prune ( $T$ , candidate)
- Accuracy[candidate] ← Test (Temp[candidate],  $D_{validation}$ )
- $T \leftarrow T$  = Temp 中で Accuracy が最良のもの
- RETURN (pruneしおえた)  $T$



53

## Reduced-Error Pruning の効果

- Reduced-Error Pruning によるテスト誤差の減少



- 節を刈ることによってテスト誤差が減少する
- 注:  $D_{validation}$  は  $D_{train}$  と  $D_{test}$  のどちらとも異なる

- 賛成論と批判論

- 賛成: 最も正確な  $T$  ( $T$  の部分木) のうちで最小のものが生成できる
- 批判:  $T$  を作るのにわざわざデータ量を減らしている
  - $D_{validation}$  をとりおけるだけの余裕があるか?
  - データ量が十分でなければ、誤差をなおさら大きくする ( $D_{train}$  が不十分)

54

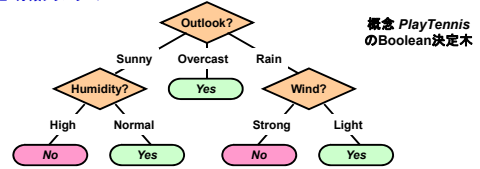
## Rule Post-Pruning

- しばしば用いられる方法
  - これもよく知られた overfitting 対応策
  - C4.5でその亜種が用いられた。C4.5は ID3の派生・後継。
- アルゴリズム Rule-Post-Pruning ( $D$ )
  - $D$ から  $T$ を生成 (ID3を使用) - 可能な限り  $D$ に適合するまで成長させる (過学習も許す)
  - $T$ を等価な規則集合に変換 (根節から葉節へ道一つにつき1規則)
  - それぞれの規則を、独立に、条件をどれでも、推定精度が改善する限り、除去することにより刈り込む(一般化する)
  - 刈り込んだ規則をソートする
    - 推定精度に従ってソートする
    - 列に並べて、 $D_{test}$ に適用する

55

## 決定木を規則に変換する

- 規則の構文
  - 左辺: 条件 (属性の等式テスト上の連言複素形 conjunctive formula)
  - 右辺: 分類クラスラベル

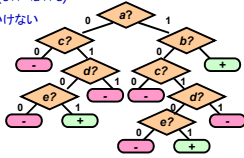


- 例
  - IF (Outlook = Sunny) ^ (Humidity = High) THEN PlayTennis = No
  - IF (Outlook = Sunny) ^ (Humidity = Normal) THEN PlayTennis = Yes
  - ...

56

## 決定木における重複

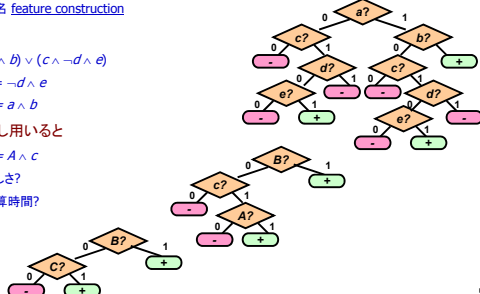
- 決定木: 表現上の短所
  - 決定木は、一番簡単な表現というわけではない
  - ポイント: 属性を重複 replication させる必要がある場合がある
- 属性重複の例
  - e.g., Disjunctive Normal Form (DNF):  $(a \wedge b) \vee (c \wedge \neg d \wedge e)$
  - (どちらかの) 連言は部分木として重複させないといけない
- 部分解
  - 新しい属性を作る
  - 別名 constructive induction (CI)
  - Mitchellの第10章参照



57

## 少しだけ: 決定木の構成的帰納法

- 新しい属性の合成
  - 一つの "+ 節" に到る直前の二つの属性の連言から新しい属性を合成 synthesize する
  - 別名 feature construction
- 例
  - $(a \wedge b) \vee (c \wedge \neg d \wedge e)$
  - $A = \neg d \wedge e$
  - $B = a \wedge b$
- 繰り返し用いると
  - $C = A \wedge c$
  - 正しさ?
  - 計算時間?



58

## 決定木: 他の話題

- 他の機械学習に共通する課題

59

## 連続値属性

- 連続値属性を扱う2つの方法
  - 離散化
    - 実数値属性を、予め、いくつかの範囲に分ける
    - e.g.,  $\{high = Temp > 35^\circ C, med = 10^\circ C < Temp \leq 35^\circ C, low = Temp \leq 10^\circ C\}$
  - 内節を分けるのに、閾値を用いる
    - e.g.,  $A \leq a$ によって二つの部分集合  $A \leq a$ と  $A > a$ ができる
    - この離散化に際して、情報増分が同様に計算される
- 情報増分を最大にする分割はどうやって得るか?
  - FOR 連続値属性  $A$ のそれぞれ
    - 事例  $\{x \in D\}$ を  $x.A$ に従って、分割する
    - FOR 異なったラベルを持つ  $A$ の値の順序対  $(l, u)$ それぞれ
    - 閾値の候補として、中点 mid-point の情報増分を評価, i.e.,  $D_{A \leq (l+u)/2} D_{A > (l+u)/2}$
  - 例
 

|              |    |    |    |    |    |    |    |
|--------------|----|----|----|----|----|----|----|
| $A = Length$ | 10 | 15 | 21 | 28 | 32 | 40 | 50 |
| Class:       | -  | +  | +  | -  | +  | +  | -  |

    - 閾値のチェック:  $Length \leq 12.5? \leq 24.5? \leq 30? \leq 45?$

60

## 多値属性に伴う問題

- 問題
  - もしある属性が多値であると、 $Gain(\cdot)$ はそれを選びやすい(なぜ?)
  - 例えば、日付(2007/11/01等)を属性として用いることを想像してみればわかる!
- 一つのアプローチ:  $GainRatio$ を  $Gain$ の代わりに使用

$$Gain(D, A) = H(D) - \sum_{v \in values(A)} \left[ \frac{|D_v|}{|D|} \cdot H(D_v) \right]$$

$$GainRatio(D, A) = \frac{Gain(D, A)}{SplitInformation(D, A)}$$

$$SplitInformation(D, A) = - \sum_{v \in values(A)} \left[ \frac{|D_v|}{|D|} \log \frac{|D_v|}{|D|} \right]$$

- $SplitInformation$ :  $c = |values(A)|$  に、ほぼ、比例  
i.e., 多くの値をもつ属性にハンディを負わせる
  - e.g., 仮定:  $c_j = c_{total} = n$  として  $c_j = 2$
  - $SplitInformation(A_j) = \log(n)$ ,  $SplitInformation(A_j) = 1$
  - もし  $Gain(D, A) = Gain(D, A_j)$  とすると,  $GainRatio(D, A) \ll GainRatio(D, A_j)$
- すなわち,  $GainRatio(\cdot)$ を用いれば, (分岐数が少ない方への)選択バイアスが表現できる

61

## 補足: Gini index

- もう一つの分割の指標
  - $n$  は分類・クラスの個数
  - $Gini(D)$  は  $D$  内の分布が偏れば偏るほど、すなわち、pure になるほど小さくなる

$$Gini(D) = \sum_{i \neq j} p_i p_j = 1 - \sum_{i=1}^n p_i^2$$

$$GiniGain(D, A) = Gini(D) - \sum_{v \in values(A)} \left[ \frac{|D_v|}{|D|} \cdot Gini(D_v) \right]$$

62

## コスト付き属性

- 応用分野毎
  - 医療: 体温検査のコストは 1000円; 血液検査 1500円; 生検 50000円
    - また検査の侵襲性・無侵襲性も考慮する必要あり
    - 患者へのリスクも(e.g., 羊水検査)
  - 他のコスト
    - サンプリング時間: e.g., ロボットのソーナー(レンジファインダー, etc.)
    - 人工物, 生体へのリスク(どんな情報を収集するか)
    - 関連する分野(e.g., 断層装置): 非破壊検査
- 低い期待コストでいかに consistent な木を作るか?
  - 一つのアプローチ: 情報増分  $gain$  を **コスト正規化増分  $Cost-Normalized-Gain$**  で置き換える
  - 正規化関数の例
    - [Nunez, 1998]:

$$Cost-Normalized-Gain(D, A) = \frac{Gain^2(D, A)}{Cost(D, A)}$$

$$Cost-Normalized-Gain(D, A) = \frac{2^{Gain(D, A)} - 1}{(Cost(D, A) + 1)^w} \quad w \in [0, 1]$$

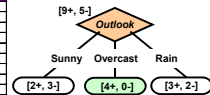
但し  $w$  はコストの重要性を定める

63

## 欠測値: 属性値が不明

- 問題: 属性  $A$  の値がない事例があるとうどうなるか?
  - しばしば, 訓練時やテスト時に, 必ずしも全ての属性値が入手できるとは限らない
  - 例: 医療診断
    - < Fever = true, Blood-Pressure = normal, ..., Blood-Test = ? ... >
    - 値は、本当になかったり, またあっても信頼度が低かったりする
  - 欠測値: 訓練時 versus 分類時
    - 訓練時: ある  $x \in D$  について  $A$  の値が与えられていないとき  $Gain(D, A)$  を評価する
    - 分類時:  $A$  の値を知らずに, 新しい事例  $x$  を分類する
- 解:  $Gain(D, A)$  の計算の中に推測を入れる

| Day | Outlook  | Temperature | Humidity | Wind   | PlayTennis? |
|-----|----------|-------------|----------|--------|-------------|
| 1   | Sunny    | Hot         | High     | Light  | No          |
| 2   | Sunny    | Hot         | High     | Strong | No          |
| 3   | Overcast | Hot         | High     | Light  | Yes         |
| 4   | Rain     | Mild        | High     | Light  | Yes         |
| 5   | Rain     | Cool        | Normal   | Light  | Yes         |
| 6   | Rain     | Cool        | Normal   | Strong | No          |
| 7   | Overcast | Cool        | Normal   | Strong | Yes         |
| 8   | Sunny    | Mild        | ???      | Light  | No          |
| 9   | Sunny    | Cool        | Normal   | Light  | Yes         |
| 10  | Rain     | Mild        | Normal   | Light  | Yes         |
| 11  | Sunny    | Mild        | Normal   | Strong | Yes         |
| 12  | Overcast | Mild        | High     | Strong | Yes         |
| 13  | Overcast | Hot         | Normal   | Light  | Yes         |
| 14  | Rain     | Mild        | High     | Strong | No          |



64

## 欠測値: 対応策

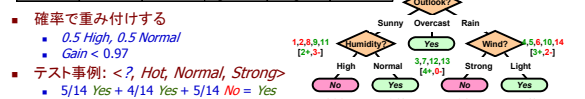
- 訓練事例はとにかく使用する. 木を(根節から)辿りつつ作っていくとき
  - 考慮すべき属性のどれについても, 事例中でもし値が知られていないなら, それを推測する
  - その推測は, 今いる節に割当てられた事例の知られている値に基づく
- $x.A$  の最もありそうな値を推測する
  - 第一案: 節  $n$  で属性  $A$  をテストするなら,  $n$  を通る事例の  $A$  の値でもっとも多いものを用いる
  - 第二案 [Mingers, 1989]: 節  $n$  で属性  $A$  をテストするなら,  $n$  を通る事例で  $x$  と同じクラスラベルをもつものの  $A$  の値でもっとも多いものを用いる
- 推測値を分散させる
  - 両隣け: 値の分布に従い, 推測値を分散させる
  - $x.A$  の可能な値  $v_j$  の分布に比例して確率  $p_j$  を割当てる[Quinlan, 1993]
    - 木の子孫に,  $x$  の内の  $p_j$  分を割当てる. データ数に 3.7個などという値が出現する
    - これを用いて  $Gain(D, A)$  or  $Cost-Normalized-Gain(D, A)$  を計算する
- どのアプローチにおいても, 新事例も同様に分類する

65

## 欠測値: 例

- $x.A$  の最もありそうな値を予測する
  - 第一案: Humidity = Normal
  - 第二案: Humidity = High (No 事例はすべて High)
  - (最も Gain の大きなものはどちらだろうか? High: Gain = 0.97, Normal: Gain < 0.97)

| Day | Outlook  | Temperature | Humidity | Wind   | PlayTennis? |
|-----|----------|-------------|----------|--------|-------------|
| 1   | Sunny    | Hot         | High     | Light  | No          |
| 2   | Sunny    | Hot         | High     | Strong | No          |
| 3   | Overcast | Hot         | High     | Light  | Yes         |
| 4   | Rain     | Mild        | High     | Light  | Yes         |
| 5   | Rain     | Cool        | Normal   | Light  | Yes         |
| 6   | Rain     | Cool        | Normal   | Strong | No          |
| 7   | Overcast | Cool        | Normal   | Strong | Yes         |
| 8   | Sunny    | Mild        | ???      | Light  | No          |
| 9   | Sunny    | Cool        | Normal   | Light  | Yes         |
| 10  | Rain     | Mild        | Normal   | Light  | Yes         |
| 11  | Sunny    | Mild        | Normal   | Strong | Yes         |
| 12  | Overcast | Mild        | High     | Strong | Yes         |
| 13  | Overcast | Hot         | Normal   | Light  | Yes         |
| 14  | Rain     | Mild        | High     | Strong | No          |



- 確率が重み付けする
  - 0.5 High, 0.5 Normal
  - Gain < 0.97
- テスト事例: < ? , Hot, Normal, Strong >
  - 5/14 Yes + 4/14 Yes + 5/14 No = Yes

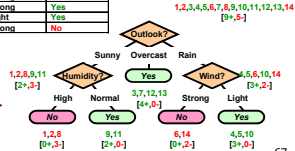
66

## 欠測値: 例

- $X:A$  の最もありそうな値を予測する
  - 第一案: Humidity = Normal
  - 第二案: Humidity = High (No 事例はすべて High)
  - (最も Gain の大きなものはどうだろうか? High: Gain = 0.97, Normal: Gain < 0.97)

| Day | Outlook  | Temperature | Humidity | Wind   | PlayTennis? |
|-----|----------|-------------|----------|--------|-------------|
| 1   | Sunny    | Hot         | High     | Light  | No          |
| 2   | Sunny    | Hot         | High     | Strong | No          |
| 3   | Overcast | Hot         | High     | Light  | Yes         |
| 4   | Rain     | Mild        | High     | Light  | Yes         |
| 5   | Rain     | Cool        | Normal   | Light  | Yes         |
| 6   | Rain     | Cool        | Normal   | Strong | No          |
| 7   | Overcast | Cool        | Normal   | Strong | Yes         |
| 8   | Sunny    | Mild        | ???      | Light  | No          |
| 9   | Sunny    | Cool        | Normal   | Light  | Yes         |
| 10  | Rain     | Mild        | Normal   | Light  | Yes         |
| 11  | Sunny    | Mild        | Normal   | Strong | Yes         |
| 12  | Overcast | Mild        | High     | Strong | Yes         |
| 13  | Overcast | Hot         | Normal   | Light  | Yes         |
| 14  | Rain     | Mild        | High     | Strong | No          |

- 確率で重み付けする
  - 0.5 High, 0.5 Normal
  - Gain < 0.97
- テスト事例: <?, Hot, Normal, Strong>
  - 1/3 Yes + 1/3 Yes + 1/3 No = Yes
  - 5/14 Yes + 4/14 Yes + 5/14 No = Yes



67

## ところで学習とは

68

## Induction (帰納)

- OED (Oxford English Dictionary) によれば
  - the process of inferring a general law or principle from the observations of particular instances
  - これは、inductive inference のこととする
  - inductive reasoning は: the process of reassigning a probability (or credibility) to a law or proposition from the observation of particular events

69

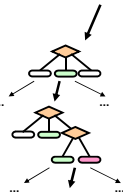
## 帰納とは(2)

- 帰納とは
  - データに潜在する規則性を得ること
  - 物体の落下の実験データ → 万有引力の法則
  - 惑星の公転運動 → 楕円運動、面積速度一定、調和
- 帰納で得た規則の正しさはどう測るか

70

## ID3 による仮説空間探索

- 探索問題
  - 探索の対象は 決定木全部の空間, すなわちブール関数をすべて表現可能な空間
    - Pros: 表現力; 柔軟性
    - Cons: 計算量; 巨大, 意味の分からない木も含む
  - 目的: もっともよい決定木を見出す (最小な consistent な木)
  - 障害: この木を見出す問題は NP-hard
  - Tradeoff
    - heuristics の使用 (探索の案内役としての目子)
    - 貪欲 greedy アルゴリズムの使用
    - すなわち, バックトラックなしの山登り hill-climbing (gradient "descent") ...
- 統計的学習
  - 事例の部分集合  $D_0$  の統計的用量  $p_0, p$  に基づく決定
  - ID3 では, 全てのデータを使用
  - ノイズのあるデータに対してロバスト



71

## ID3 の帰納バイアス

- 探索におけるヒューリスティックは帰納バイアスである
  - $H$  は  $X$  の冪集合 (全部分集合の集合)
  - ⇒ 帰納バイアスなしと言っよいか? いや, そうではない...
    - 短い木への選好 (終了条件から) がある
    - 情報量増分が高い属性を根節に近いうちという選好がある
    - Gain(): ID3 の帰納バイアスを体現するヒューリスティック関数
  - ID3 の帰納バイアス
    - ある仮説への選好をヒューリスティック関数に表現している
    - 比較して: 仮説空間  $H$  を制限すること (命題論理の正規形に基づく制限: k-CNF, etc.)
- 短い木を好むこと
  - データに適合する木の中で最短のものを選ぶ
  - オッカムの剃刀バイアス: 観測を説明する最短の仮説をとれ

学習時に用いる、データ以外の仮定。それにより、こちらの仮説がより良い、この仮説はとらない、この仮説はとるといことが決まる  
これがないと、データを説明する仮説が多数(無限に)あって、結論が得られない  
いや、しかし、データ以外の情報を使って、仮説を選択するのはまずいのではないか?  
もし、バイアスが避けえないとしたら、どういふバイアスがよいのか?

72

## 学習とバイアス

- **バイアス**: 仮説間に順位があるとき、その順位
  - 同時に複数個の仮説をみたときの、選好順位
  - 一度に一個ずつ見るときの、探索順序
- データに適合する仮説は、一般に、多量にあるので、学習するには**バイアスが必要**
  - 仮説を一個選択するのではなく、複数個の仮説を用いる場合でも、「データに適合する仮説をすべて用いる」のではない限り、**バイアスが必要**である。

学習:

データ → 仮説

いや、しかし、データ以外の情報を使って、仮説を選択するのはまずいのではないかと?

もし、バイアスが避けえないとしたら、どういうバイアスがよいのか?

73

## Occam の剃刀

- 人口に膾炙しているのは
  - Entities should not be multiplied beyond necessity.
- Bertrand Russell によれば
  - It is vain to do with more what can be done with fewer.
- 最も普通の解釈
  - Among the theories that are consistent with the observed phenomena, one should select the simplest theory.

74

## Isaac Newton の言葉

- We are to admit no more causes of natural things than such as are both true and sufficient to explain the appearances. To this purpose the philosophers say that Nature does nothing in vain, and more is in vain when less will serve; for **Nature is pleased with simplicity**, and affects not the pomp of superfluous causes.

75

## オッカムの剃刀: ある選好バイアス

- 帰納バイアス2つ: 選好バイアス preference biases と言語バイアス language biases
  - **選好バイアス**
    - ・ 学習アルゴリズムに(普通は暗黙的に)組み込まれている
    - ・ 言い換えれば: 探索順序の規定
  - **言語バイアス**
    - ・ 知識(仮説)の表現に(普通は暗黙的に)組み込まれている
    - ・ 言い換えれば: 探索空間の制限
    - ・ 別名 制限バイアス
- オッカムの剃刀 Occam's Razor: 賛成意見
  - 短い仮説の方が、長い仮説に比べ、個数が少ない
    - ・ 例えば、ビット列で考えれば、長さ  $n$  のものは  $n+1$  のものに比べ半数,  $n \geq 0$ .
    - ・ 短い仮説が、もしデータにぴったり合ったとしたら、偶然とは考え難い
      - ・ 短い仮説は、個数が少ないので、説明できる現象の数が少ない
    - ・ 長い仮説 (例: 200 個の節を持つ木, かつ  $|D| = 100$ ) の場合には、偶然である可能性が高い
      - ・ いずれかの木がデータにぴったり合う。どれに合うかは偶然であるが、どれかに合うこと自体は当然。
  - 得るものと捨てたもの
    - ・ 他の条件が同一であれば、複雑なモデルの汎化能力は単純なモデルほどではない
    - ・ あとになってもっと柔軟な(微調整可能な)モデルが必要になることはないとは仮定

76

## オッカムの剃刀と決定木: 二つの問題

- オッカムの剃刀 Occam's Razor: 反対意見
  - 仮説空間  $H$  に依存して  $size(h)$  が決まる。同じ  $k$  でも  $H$  が異なると  $size(h)$  が異なる。
  - 「小ささ」を選好することへの疑問: 「少ない」ことは正当化にならない
- オッカムの剃刀 Occam's Razor は Well-Defined か?
  - 内部の知識表現 knowledge representation によってどの  $h$  が「短い」かがきまる --- 恣意的?
    - ・ 例えば、テスト "(Sunny ^ Normal-Humidity) v Overcast v (Rain ^ Light-Wind)" は一語?
  - 答: 表現言語を固定; 十分長いところでは、長い仮説は、内部表現によらず、やっぱり長い
    - ・ 反論: 答えになっていない。実際には「短い仮説」に関する議論が重要
- 「短い仮説」であって、どうして他の「小さい仮説空間」ではないのか?
  - 小さい仮説集合を定義する方法はいくつもある。
  - 選好バイアスで用いる size が何であって、適当に基準  $S$  を選べば  $size(h)$  をその限界内に制限することができる (i.e., " $S$  に合致する木のみ受理する")
    - ・ e.g., 節の個数が素数であって、文字 "Z" で始まる属性を用いている木
    - ・ なぜ「小さな木であって、(例えば)  $A_0, A_1, \dots, A_{11}$  を順番にテストするもの」ではないのか?
    - ・  $size(h)$  に基づいて小さな仮説集合を定義することに、特別の意味があるのか?
  - 参考: Chapter 6, Mitchell's Machine Learning

77

## エピクロスの多説明原理

- ギリシャの哲学者 Epicurus
  - If more than one theory is consistent with the observations, keep all theories (Principle of Multiple Explanations).
- その一つの理由: 一つを他から選り出す理由がない

78