

2022年度 知識工学 レポート課題

氏名:
学生番号:

2022年5月29日

1 課題1

1.1 課題

決定木学習モデルを download し, windows または linux にインストールして付属のデモファイルについて解説をする. ない場合は, 簡単な例を作成して動作説明をする (解説 blog などでも利用されているデータを参考にして良い)

さらに動かした結果を付与し説明を加えよ. いくつかのデモが付与されているが1つで良い.

1.2 手法

RuleQuest Research[1] の配布する, See5/C5.0 のデモプログラムを実行する.
付属のデモファイル `anneal.{data,names,test}` を使用した.

1.3 デモファイル

1.3.1 `anneal.names`

`anneal.names` はスキーマファイルである. 後述する `anneal.data`, `anneal.test` のクラス, 各属性の型を定義する.

```
1 1,2,3,5,U.  
2  
3 family:      discrete 12.|GB,GK,GS,TN,ZA,ZF,ZH,ZM,ZS.  
4 product-type: C, H, G.  
5 steel:       R,A,U,K,M,S,W,V.  
6 carbon:      continuous.  
7 hardness:    continuous.
```

上記が `anneal.names` の冒頭の抜粋である.

1 行目はクラスを定義する. 与えるデータが 1, 2, 3, 5, U の 5 つのクラスに分類される.

3 行目以降は各属性について定義する. 3 行目は `family` という属性値について定義する. `family` は最大で 12 種類の離散値を取ることができる. | 以降の記述はコメントであることに注意されたい.

4,5 行目はそれぞれ `product-type`, `steel` という属性値について定義する. これらはそれぞれ C, H, G と R, A, U, K, M, S, W, V のうちのいずれかの値を取ると明示的に指定している.

6,7 行目はそれぞれ `carbon`, `hardness` という属性値について定義する. これらは連続値 (実数値) を取ることができる.

1.3.2 anneal.data, anneal.test

anneal.data は学習データ、anneal.test はテストデータであり、同じ形式を取る。

```
anneal.data
1 TN,C,A,00,00,N/A,N/A,3,000,N,N/A,N/A,N/A,N/A,N/A,N/A,N/A,C,N/A,N/A,N/A,N/A,N/A,
2   N/A,N/A,N/A,N/A,N/A,SHEET,0.600,1220.0,0761,N/A,0000,N/A,5
3 N/A,C,A,00,00,N/A,S,2,000,N/A,N/A,E,N/A,N/A,N/A,N/A,N/A,N/A,N/A,N/A,N/A,N/A,N/A,
4   N/A,N/A,N/A,N/A,N/A,SHEET,0.699,0610.0,0762,N,0000,N/A,3
5 TN,C,N/A,00,00,N/A,A,1,000,N/A,N/A,N/A,N/A,N/A,N/A,N/A,N/A,N/A,N/A,N/A,N/A,Y,N/A,
6   N/A,N/A,N/A,N/A,N/A,SHEET,0.500,0609.9,0612,N/A,0000,N/A,5
7 N/A,C,A,00,60,T,N/A,N/A,000,N/A,N/A,G,N/A,N/A,N/A,N/A,B,Y,N/A,N/A,N/A,N/A,N/A,
8   N/A,N/A,N/A,N/A,N/A,SHEET,0.800,0356.1,0762,N/A,0000,N/A,3
9 ZS,C,A,00,00,N/A,S,3,000,N,N/A,E,N/A,N/A,N/A,N/A,N/A,N/A,N/A,N/A,N/A,N/A,N/A,N/A,
10  N/A,N/A,N/A,N/A,N/A,COIL,1.000,0075.0,0000,N/A,0000,N/A,3
```

上記が anneal.data の冒頭の抜粋である（適当に改行を加えた）。

各行が学習データであり、コンマで区切られたデータが並んでいる。1 列目から順に anneal.names に定義した属性が順番に並んでいる。そして末尾にクラスが与えられる。

1.4 実行結果

実行結果の全体は付録 A に添付している。ここでは一部を説明するに留める。また、詳細については <https://www.rulequest.com/see5-unix.html#CLASSIFIERS> に記載されている。

11 行目には学習データの要約が記載されている。デモプログラムでは 400 ケースまでの制限があるため、400 ケース 38 属性のデータが読み込まれたことを示している。

13-38 行目には学習結果が決定木として示されている。上から順に、hardness の属性値が 75 より大きいかわそれ以下かに分け、大きければクラス U に分類し、それ以下であれば更に分類する。strength の属性値が 375 より大きいかわそれ以下かに分け、大きく、かつ steel の属性値が U, K, W, V, M, N/A のいずれかであればクラス 2 に、R, A のどちらかであればクラス 3 に、S であればクラス 1 に分類する。strength の属性値が 375 以下で、かつ family の属性値が TN であればクラス 5 に、ZS であればクラス 3 に、N/A であれば更に分類。と続く。

また、分類クラスの後ろの括弧付きの数字は、(そのクラスに分類されたケース数/分類が誤っていたケース数)を示している。

41-68 行目には学習データに対する決定木の評価が示されている。Size は分類されたインスタンス数が 0 でない決定木の葉の個数を、Errors は実際の（与えられた）クラスと誤って分類されたインスタンスの個数及び割合を示している。

50-57 行目は混同行列と呼ばれる行列であり、行が実際のクラス、列が分類されたクラスを表している。すなわち対角成分が正しく分類されたインスタンス数を、非対角成分が誤って分類されたインスタンス数を表す。

59-68 行目は各属性が学習データの分類をするときに使われたインスタンス数の割合を示している。すべてのインスタンスの hardness が使われ、分類が完了しなかった 95% のインスタンスの strength が分類に使われた、といった具合である。

74-89 行目はテストデータに決定木を適用したときの混同行列である。

2 課題 2

2.1 課題

売れているノートパソコンの人気機種の属性を調べる

- Web サイトの人気順

- 上位 1-10 は正例（売れた）
 - 下位 11 以降は負例（売れないと仮定）
1. moodle に C5.0 用のファイルを置いたので実行してどういう木になったか分析
 2. 元々のファイル notePC_rank.xlsx から属性を 1 つ選び、追加して分析せよ

2.2 課題 2.1

実行結果は付録 B に示す。

12-20 行目の決定木から分類結果を表にまとめると表 1 のようになる。

表 1 ノート PC の売れ方と価格・CPU スコアの関係

| CPU スコア | ≤ 9019 | > 9019 | |
|---------|--------|---------|----------|
| 価格 | | ≤ 11550 | > 11550 |
| N/A | 負例 (2) | | |
| ≤ 49784 | 負例 (2) | | |
| > 49784 | 正例 (4) | 正例 (3) | |
| ≤ 68020 | | | |
| > 68020 | | 負例 (4) | 正例 (3/1) |

ここからは大まかに価格が 49784 円より高いほうが売れ行きがよく、価格が高いが CPU スコアが 9019 から 11550 の間であるとあまり売れないと読み取れる。

2.3 課題 2.2

ディスプレイサイズの属性を加え、notepc.data, notepc.names を 付録 C のようにした。

このときの C5.0 の実行結果は 8 行目の属性数を除き 2.2 節と同じであった。

参考文献

- [1] RuleQuest Research Data Mining Tools, <https://www.rulequest.com/index.html>, 2022 年 5 月 19 日アクセス。

付録 A anneal の決定木の学習結果

```

1  anneal.output
2  C5.0 [Release 2.11a]          Thu May 19 20:54:42 2022
3  -----
4
5  Options:
6      Application './C50Release2/Data/anneal'
7
8      ** This demonstration version cannot process **
9      ** more than 400 training or test cases.      **
10
11 Read 400 cases (38 attributes) from ./C50Release2/Data/anneal.data
12
13 Decision tree:

```

```

14
15 hardness > 75: U (21/1)
16 hardness <= 75:
17 :...strength > 375:
18 :...steel in {U,K,W,V}: 2 (0)
19 :   steel = N/A: 2 (6)
20 :   steel in {R,A}: 3 (5)
21 :   steel = M: 2 (6)
22 :   steel = S: 1 (2)
23 strength <= 375:
24 :...family = TN: 5 (35)
25   family = ZS: 3 (19)
26   family = N/A:
27 :...enamelability in {3,4,5}: 3 (0)
28   enamelability = 1: 3 (3)
29   enamelability = 2: 2 (5)
30   enamelability = N/A:
31 :...surface-quality in {D,E,F,G}: 3 (240)
32   surface-quality = N/A:
33 :...condition in {A,X}: 3 (0)
34   condition = N/A: 3 (24)
35   condition = S:
36 :...temper_rolling = -: 2 (0)
37   temper_rolling = N/A: 2 (32)
38   temper_rolling = T: 3 (2)
39
40
41 Evaluation on training data (400 cases):
42
43         Decision Tree
44         -----
45         Size      Errors
46
47         13      1( 0.2%)  <<
48
49
50         (a)  (b)  (c)  (d)  (e)  <-classified as
51         ----  ---  ---  ---  ---  -----
52         2                                (a): class 1
53         49                                (b): class 2
54         293                               (c): class 3
55         35                                (d): class 5
56         20                                (e): class U
57
58
59 Attribute usage:
60
61         100% hardness
62         95% strength
63         91% family
64         77% enamelability
65         75% surface-quality
66         15% condition
67         9% temper_rolling
68         5% steel
69

```

```

70
71      ** This demonstration version cannot process **
72      ** more than 400 training or test cases.      **
73
74 Evaluation on test data (400 cases):
75
76      Decision Tree
77      -----
78      Size      Errors
79
80      13      6( 1.5%)  <<
81
82
83      (a)  (b)  (c)  (d)  (e)  <-classified as
84      ----  ---  ---  ---  ---
85      2          3          (a): class 1
86          38          (b): class 2
87          311         1    (c): class 3
88          25          (d): class 5
89          2          18    (e): class U
90
91
92 Time: 0.0 secs

```

付録 B notepc_s の決定木の学習結果

```

1
2 C5.0 [Release 2.11a] Thu May 19 23:35:59 2022
3 -----
4
5 Options:
6 Application `notepc_s'
7
8 Read 18 cases (5 attributes) from notepc_s.data
9
10 Decision tree:
11
12 price = N/A: 0 (2)
13 price <= 49784: 0 (2)
14 price > 49784:
15 :...cpu_score <= 9019: 1 (4)
16   cpu_score > 9019:
17   :...price <= 68020: 1 (3)
18     price > 68020:
19     :...cpu_score <= 11550: 0 (4)
20       cpu_score > 11550: 1 (3/1)
21
22
23 Evaluation on training data (18 cases):
24
25      Decision Tree
26      -----
27      Size      Errors
28

```

```

29         6      1( 5.6%)   <<
30
31
32         (a)   (b)   <-classified as
33         ----  ----
34         8      1      (a): class 0
35             9      (b): class 1
36
37
38     Attribute usage:
39
40         100%  price
41         78%   cpu_score
42
43
44     Time: 0.0 secs

```

付録C notepc の.data と.names

```

1  ----- notepc.data -----
2  Dell,80282.0,10088,8,1.06,13.3,0
3  NEC,111183.0,13803,8,2.1,15.6,0
4  Dell,37980.0,4146,4,1.59,14,0
5  Lenovo,49588.0,11275,8,1.5,14,0
6  Lenovo,N/A,16190,8,2.2,17.3,0
7  Lenovo,71060.0,10088,8,1.7,15.6,0
8  Dell,96122.0,10521,16,1.714,15.6,0
9  マウスコンピューター,N/A,7951,8,1.59,15.6,0
10 Dynabook,127380.0,10521,16,1.94,15.6,0
11 NEC,89700.0,6132,8,N/A,15.6,1
12 Dell,56026.0,7149,8,1.83,15.6,1
13 マウスコンピューター,129800.0,12580,16,2.13,15.6,1
14 Dell,49980.0,6324,8,1.78,15.6,1
15 マイクロソフト,80500.0,7951,8,1.11,12.4,1
16 Dell,64676.0,13851,8,N/A,14,1
17 Lenovo,56430.0,11275,8,1.45,14,1
18 Lenovo,75460.0,13803,16,1.5,14,1
19 Dell,64980.0,10088,8,1.78,15.6,1

```

```

1  ----- notepc.names -----
2  0,1.
3
4  maker:      discrete 6.
5  price:      continuous.
6  cpu_score:   continuous.
7  memory:     continuous.
8  weight:     continuous.
9  display:    continuous.

```