

Inteligência Artificial

Relatório 3 - Árvores de Decisão

Grupo 8

Luís Pinto, n.º 201704025
Sónia Almeida, n.º 201811293
Miguel Lançóis, n.º 201506342

30 de Abril de 2019

Conteúdo

1	Introdução	3
2	Algoritmos para árvores de decisão	3
2.1	ID3	4
2.2	Cart	4
2.3	Mars	4
2.4	Chaid	4
3	Implementação	5
4	Resultados	5
4.1	Restaurant	5
4.2	Weather	6
4.3	Iris	7
5	Conclusão	9
6	Bibliografia/webgrafia	9

1 Introdução

Neste trabalho, vamos explorar a criação de árvores de decisão com base em algoritmos estudados em aula. Estas árvores são geradas através de ditos algoritmos e conjuntos de dados obtidos previamente.

NOTA: Sempre que aparecerem parênteses retos em frente a alguma palavra, é uma referência à bibliografia/webgrafia onde os números dentro dos parênteses representam os links na secção 6.

Árvores de Decisão

- **O que são?**

Estes grafos são árvores que nos permitem determinar o resultado final de um dado conjunto de atributos. Cada nó, que não uma folha, representa um atributo ao qual estão associadas arestas que representam o valor associado a certo atributo. Se o nosso problema fosse gerar uma árvore de decisão para lançamento de moedas, o nó raiz seria o atributo "Qual face saíu?", as arestas "Cara/Coroa" e os nós filhos "True/False", dependendo do resultado final pretendido, por exemplo.

- **Para que servem?**

Com árvores de decisão, torna-se possível organizar informação de forma mais compacta e eficaz para evitar ter em consideração todos os pequenos aspetos ao tomar decisões. Se o problema em questão for descobrir quais as melhores características de uma galinha para obter ovos mais saborosos, podemos arranjar, por exemplo, um conjunto de dados com os atributos: Cór das penas, Cór dos olhos, Tamanho, Comida. Com estes dados, podemos gerar uma árvore de decisão que nos permitirá escolher as melhores galinhas para obter ovos mais saborosos, por exemplo.

2 Algoritmos para árvores de decisão

Existem vários algoritmos que permitem gerar árvores de decisão a partir de uma série de dados. Existem dois tipos de árvores de decisão:

- As árvores de classificação (= classification trees)
- As árvores de regressão (= regression trees)

As **árvores de classificação** permitem prever a classe da variável desejada: é prevista a classe da variável. Para este tipo de árvore, queremos maximizar a homogeneidade de um subconjunto construído a partir de um conjunto maior. Existem várias formas de medir esta homogeneidade, as mais famosas sendo a entropia de Shannon e o critério de Gini. Estas duas métricas são usadas, respetivamente, nos algoritmos ID3 (Iterative Dichotomiser 3) e CART (Classification And Regression Tree).

As **árvores de regressão**, elas, permitem prever uma quantidade real como um tempo de espera por exemplo ou um preço: neste caso, é previsto um valor numérico. Neste caso, queremos maximizar a variância entre classes para que os valores sejam as mais distantes possível.

Os algoritmos que vão ser apresentados no que segue são os seguintes: o ID3 (Iterative Dichotomiser 3), o CART (Classification And Regression Tree), o MARS (Multivariate adaptive regression spline) e o CHAID (CHi-squared Automatic Interaction Detector) embora existam dezenas de algoritmos para árvores de decisão. [4 e 6]

2.1 ID3

O ID3 usa a função de entropia de Shannon para construir a sua árvore. Em termodinâmica, a entropia designa uma grandeza extensiva, ou seja, que é proporcional a uma grandeza que caracteriza um sistema, como a sua massa por exemplo, ao contrário de uma grandeza intensiva que não depende do tamanho do sistema, como a massa volúmica. A entropia caracteriza a desordem de um sistema. Quanto mais entropia, maior a desordem. Assim, um sistema ordenado tem um valor de entropia nulo. Em informática, a entropia de Shannon caracteriza a falta de informação de um sistema. Quanto menos informação, maior o valor de entropia. Este valor de entropia, para uma variável aleatória V com valores $v_k, k \in 1, \dots, n$, com probabilidade $P(v_k)$, é calculado da seguinte maneira: [7]

$$H(V) = \sum_{k=1}^n P(v_k) \log_2 \frac{1}{P(v_k)} = - \sum_{k=1}^n P(v_k) \log_2 P(v_k) \text{ e } H(V) \in [0, 1]$$

A partir da entropia, é avaliado o ganho de informação, que é calculado da seguinte forma: determina-se a falta de informação, entre 0 e 1, e subtrai-se a entropia. Se a falta de informação for total, o valor de ganho de informação é 1 - entropia. Isto acontece quando o atributo não nos traz informação nenhuma. Isto acontece com o atributo Class no restaurante, que contém 6 Yes e 6 No: não nos traz informação nenhuma. Este algoritmo é baseado sobre um conjunto de classes iniciais. Para que o modelo seja geral que chegue, devemos ter cuidado para que essas classes não influenciem de mais o modelo: queremos evitar overfitting e ter um modelo muito refinado para os dados que temos. Nesse caso, o modelo é demasiado específico e não é válido para outro tipo de dados. [7] Todavia, queremos obter um modelo genérico e não demasiado específico. Outra limitação do ID3 é que só é usado um atributo de cada vez para tomar uma decisão. Classificar dados contínuos pode ser computacionalmente caro. [5]

2.2 Cart

O CART usa o índice de Gini para classificar os seus dados. O índice Gini é muito utilizado para caracterizar as desigualdades dentro de uma população. Este índice varia entre 0 e 1. Se o índice tiver valor 0, a igualdade é perfeita, se tiver valor 1, a desigualdade é inteira (por exemplo, só uma pessoa é que possui as riquezas todas). Em informática, este valor é calculado multiplicando o valor da probabilidade que um elemento seja escolhido com a probabilidade que este elemento esteja mal classificado. [4]

2.3 Mars

Para agrupar os dados, é calculada uma soma ponderada de regressões lineares. [3]

2.4 Chaid

Este algoritmo é muito utilizado no domínio de marketing para selecionar um grupo de consumidores e prever as consequências das suas escolhas sobre outras variáveis. [2] É baseado sobre o teste de Bonferroni [1], mais especificamente sobre o ajusto do critério de rejeição controlando a probabilidade de ocorrer um certo erro.

3 Implementação

Linguagem: Para a implementação do código decidimos utilizar a linguagem *Python 3.7*.

Estrutura de Dados: Na implementação do nosso código, utilizamos dicionários e listas para manipulação da informação. Ambas estas estruturas de dados são de fácil manipulação e essa foi a razão principal para as termos escolhido ao invés de outras estruturas. Os "dataset", ficheiros em formato csv, são lidos e convertidos para um dicionário, utilizando o "ID" de cada linha como chave para o seu conteúdo. De seguida, construímos um dicionário com todos os valores para cada atributo presente no "dataset" que servirá de guia para a construção da árvore de decisão. São utilizadas listas quando se pretende analisar o conteúdo de um ou dois atributos isoladamente.

Para armazenamento da árvore de decisão é utilizada a classe "TreeNode" que é constituída pelo atributo *att*, com o valor do nó correspondente, e por um dicionário *children* onde, para cada chave correspondente ao valor do atributo, se irá referenciar o nó seguinte da árvore.

Organização do código: O código inicia-se com uma chamada da função *main()* onde se faz a leitura do ficheiro csv, *get_dictionary()*, e, de seguida, obtém-se o dicionário com os valores de cada atributo com a função *branching()*. Com os valores obtidos é, finalmente, chamado o algoritmo ID3 que irá retornar a árvore de decisão.

Para a implementação do ID3 foi necessário a implementação de várias funções auxiliares, tais como:

- *is_all_the_same_value()*: permite indicar se todos os exemplos restantes apresentam o mesmo valor da classe final;
- *most_common_value()*: retorna o valor mais comum da classe final;
- *get_next_att()*: calculando a entropia dos valores do atributo, retorna o atributo que deverá ser escolhido de seguida para a árvore.

4 Resultados

4.1 Restaurant

```

1 runfile('D:/FEUP/IA/TP/Trabalho_3/tp3(3).py', wdir='D:/FEUP/IA/TP/Trabalho_3')
2 Please insert the name of the file with the data: restaurant.csv
3 time_make_dic = 0.0
4 time_dic = 5.41792106628418
5 time_parent_branching = 0.0
6 time_ent = 0.0
7 time_branching = 0.0
8 time_ID3 = 0.0
9 <Pat>
10     Full:
11         <Price>
12             $:
13                 <Rain>
14                     No:
15                         <Bar>
16                             No:
17                                 <Fri>

```

```

18                                     No: No (1)
19                                     Yes: Yes (1)
20                                     Yes: Yes (1)
21                                     Yes: No (1)
22                                     $$: No (4)
23                                     $$$: No (2)
24                                     None: No (2)
25                                     Some: Yes (4)
26 time_print = 0.0

```

4.2 Weather

```

1 runfile('D:/FEUP/IA/TP/Trabalho_3/tp3(3).py', wdir='D:/FEUP/IA/TP/Trabalho_3')
2 Please insert the name of the file with the data: weather.csv
3 time_make_dic = 0.0
4 time_dic = 4.520612716674805
5 time_parent_branching = 0.0
6 time_ent = 0.0
7 time_branching = 0.0
8 time_ID3 = 0.0
9 <Temp>
10     (64.0, 69.0):
11         <Weather>
12             overcast: yes (1)
13             rainy:
14                 <Humidity>
15                     (65.0, 70.0): no (1)
16                     (70.0, 80.0): no (1)
17                     (80.0, 90.0): yes (1)
18                     (90.0, 95.0): no (1)
19                     (95.0, 97.0): no (1)
20             sunny: yes (2)
21     (69.0, 72.0):
22         <Weather>
23             overcast: yes (2)
24             rainy:
25                 <Humidity>
26                     (65.0, 70.0): no (1)
27                     (70.0, 80.0): no (1)
28                     (80.0, 90.0): no (1)
29                     (90.0, 95.0): no (1)
30                     (95.0, 97.0): yes (1)
31             sunny: yes (1)
32     (72.0, 75.0):
33         <Weather>
34             overcast: yes (1)
35             rainy: no (1)
36             sunny: no (1)
37     (75.0, 83.0):
38         <Humidity>
39             (65.0, 70.0): yes (3)
40             (70.0, 80.0): yes (2)
41             (80.0, 90.0): yes (1)
42             (90.0, 95.0): no (1)
43             (95.0, 97.0): yes (3)

```

```

44         (83.0, 86.0):
45             <Weather>
46                 overcast: yes (1)
47                 rainy: no (1)
48                 sunny: no (1)
49 time_print = 0.0

```

4.3 Iris

```

1  runfile('D:/FEUP/IA/TP/Trabalho_3/tp3(3).py', wdir='D:/FEUP/IA/TP/Trabalho_3')
2  Please insert the name of the file with the data: iris.csv
3  time_make_dic = 0.0
4  time_dic = 4.952691316604614
5  time_parent_branching = 0.0
6  time_ent = 0.0
7  time_branching = 0.0
8  time_ID3 = 0.006543159484863281
9  <petallength>
10     (1.0, 1.5): Iris-setosa (23)
11     (1.5, 3.9):
12         <sepalwidth>
13             (4.3, 5.0):
14                 <sepalwidth>
15                     (2.0, 2.7): Iris-versicolor (1)
16                     (2.7, 3.0): Iris-setosa (8)
17                     (3.0, 3.1): Iris-setosa (8)
18                     (3.1, 3.4): Iris-setosa (6)
19                     (3.4, 5.4): Iris-setosa (2)
20             (5.0, 5.6):
21                 <sepalwidth>
22                     (2.0, 2.7): Iris-versicolor (5)
23                     (2.7, 3.0): Iris-setosa (17)
24                     (3.0, 3.1): Iris-setosa (1)
25                     (3.1, 3.4): Iris-setosa (1)
26                     (3.4, 5.4): Iris-setosa (15)
27             (5.6, 6.1):
28                 <sepalwidth>
29                     (2.0, 2.7): Iris-versicolor (1)
30                     (2.7, 3.0): Iris-versicolor (1)
31                     (3.0, 3.1): Iris-setosa (2)
32                     (3.1, 3.4): Iris-setosa (2)
33                     (3.4, 5.4): Iris-setosa (2)
34             (6.1, 6.6): Iris-setosa (27)
35             (6.6, 8.9): Iris-setosa (27)
36     (3.9, 4.7):
37         <sepalwidth>
38             (4.3, 5.0): Iris-virginica (1)
39             (5.0, 5.6): Iris-versicolor (5)
40             (5.6, 6.1): Iris-versicolor (15)
41             (6.1, 6.6): Iris-versicolor (8)
42             (6.6, 8.9): Iris-versicolor (3)
43     (4.7, 5.4):
44         <sepalwidth>
45             (2.0, 2.7):
46                 <sepalwidth>

```

```

47         (4.3, 5.0): Iris-virginica (3)
48         (5.0, 5.6): Iris-virginica (3)
49         (5.6, 6.1): Iris-virginica (2)
50         (6.1, 6.6):
51             <petalwidth>
52                 (0.1, 0.2): Iris-versicolor (1)
53                 (0.2, 1.2): Iris-versicolor (1)
54                 (1.2, 1.5): Iris-versicolor (1)
55                 (1.5, 1.9): Iris-versicolor (1)
56                 (1.9, 3.5): Iris-virginica (1)
57         (6.6, 8.9): Iris-virginica (3)
58     (2.7, 3.0):
59         <sepallength>
60             (4.3, 5.0): Iris-virginica (8)
61             (5.0, 5.6): Iris-virginica (8)
62             (5.6, 6.1):
63                 <petalwidth>
64                     (0.1, 0.2): Iris-virginica (4)
65                     (0.2, 1.2): Iris-virginica (4)
66                     (1.2, 1.5): Iris-virginica (4)
67                     (1.5, 1.9): Iris-versicolor (1)
68                     (1.9, 3.5): Iris-virginica (4)
69             (6.1, 6.6):
70                 <petalwidth>
71                     (0.1, 0.2): Iris-virginica (4)
72                     (0.2, 1.2): Iris-virginica (4)
73                     (1.2, 1.5): Iris-versicolor (2)
74                     (1.5, 1.9): Iris-virginica (3)
75                     (1.9, 3.5): Iris-virginica (1)
76             (6.6, 8.9): Iris-versicolor (1)
77     (3.0, 3.1):
78         <petalwidth>
79             (0.1, 0.2): Iris-virginica (5)
80             (0.2, 1.2): Iris-virginica (5)
81             (1.2, 1.5): Iris-virginica (5)
82             (1.5, 1.9):
83                 <sepallength>
84                     (4.3, 5.0): Iris-virginica (3)
85                     (5.0, 5.6): Iris-virginica (3)
86                     (5.6, 6.1): Iris-virginica (2)
87                     (6.1, 6.6): Iris-virginica (1)
88                     (6.6, 8.9): Iris-versicolor (1)
89             (1.9, 3.5): Iris-virginica (2)
90     (3.1, 3.4):
91         <petalwidth>
92             (0.1, 0.2): Iris-versicolor (5)
93             (0.2, 1.2): Iris-versicolor (5)
94             (1.2, 1.5): Iris-versicolor (1)
95             (1.5, 1.9): Iris-versicolor (4)
96             (1.9, 3.5): Iris-virginica (3)
97     (3.4, 5.4): Iris-virginica (19)
98     (5.4, 7.9): Iris-virginica (30)
99 time_print = 0.0

```


5 Conclusão

Durante a realização deste trabalho, vários detalhes surgiram relativamente ao algoritmo utilizado para gerar uma árvore de decisão e, também, à árvore em si. São referidos, de seguida, alguns aspetos negativos e positivos observados neste trabalho.

Negativos:

- Ao aumentar o número de casos observados, a árvore gerada começa a demorar bastante tempo o que pode não ser muito benéfico para certas situações;
- Dados com muitos atributos diferentes, podem gerar árvores com muitos ramos o que é desvantajoso para árvores de decisão;
- Por vezes podem ser observados imensos casos e parecer que um atributo está ligado a uma dada classe mas, na verdade, pode não estar e induzir em erro.

Positivos:

- Conseguimos condensar grandes quantidades de informações em árvores de tamanhos relativamente pequenos;
- Depois de obtida a árvore, é muito fácil de verificar certas possíveis relações entre os atributos e as classes;
- Permitem-nos focar mais em dados atributos sabendo por onde começar um dado estudo mais profundo.

6 Bibliografia/webgrafia

1. https://en.wikipedia.org/wiki/Chi-square_automatic_interaction_detection
2. <https://fr.wikipedia.org/wiki/CHAID>
3. https://fr.wikipedia.org/wiki/R%C3%A9gression_multivari%C3%A9e_par_spline_adaptative
4. [https://fr.wikipedia.org/wiki/Arbre_de_d%C3%A9cision_\(apprentissage\)](https://fr.wikipedia.org/wiki/Arbre_de_d%C3%A9cision_(apprentissage))
5. <https://fr.slideshare.net/HARDIKSINGH7/id3-algorithm-56632554>
6. <https://medium.com/deep-math-machine-learning-ai/chapter-4-decision-trees-algorithms-b93975f7a1f1>
7. *Artificial Intelligence, a Modern Approach*, Russell and Norvig, cap 18.