

1 Drzewa decyzyjne w teorii decyzji

W teorii decyzji drzewo decyzyjne jest drzewem decyzji i ich możliwych konsekwencji (stanów natury). Zadaniem drzew decyzyjnych może być zarówno stworzenie planu, jak i rozwiązanie problemu decyzyjnego.

Metoda drzew decyzyjnych jest szczególnie przydatna w problemach decyzyjnych z licznymi, rozgałęziającymi się wariantami oraz w przypadku podejmowania decyzji w warunkach ryzyka.

2 Drzewa decyzyjne w uczeniu maszynowym

Drzewa decyzyjne w uczeniu maszynowym służy do wyodrębniania wiedzy z zestawu przykładów (patrz eksploracja danych). Zakładamy, że posiadamy zestaw przykładów: obiektów opisanych przy pomocy atrybutów, którym przyporządkowujemy jakąś decyzję (patrz tabela decyzyjna).

Przykład: Chcemy zautomatyzować proces przyjmowania kandydatów na praktyki w dużej firmie. Posiadamy setki przykładów z przeszłości, chcemy wydobyć z nich reguły decyzyjne. Atrybuty Wykształcenie, Języki obce, Doświadczenie i Ogólne wrażenie są kodowane skalą od 1 do 5.

Wiek	Płeć	Wykształcenie	Języki	Doświadczenie	Prezentacja	Przyjęty
25	m	2	4	1	4	nie
22	k	4	3	4	2	nie
21	m	4	5	5	4	tak
29	m	1	3	2	3	nie

Table 1:

Na podstawie tabeli decyzyjnej stworzymy drzewo, którego węzłami są poszczególne atrybuty, gałęziami wartości odpowiadające tym atrybutom, a liście tworzą poszczególne decyzje. Na podstawie przykładowych danych wygenerowano następujące drzewo:

Drzewo w takiej postaci odzwierciedla w jaki sposób były na podstawie atrybutów były podejmowane decyzje klasyfikujące (dla uproszczenia połączono niektóre gałęzie). Zaletą tej reprezentacji jest jej czytelność dla człowieka. W prosty sposób można przekształcić ją do reprezentacji regułowej.

3 Algorytm tworzenia drzewa ID3

3.1 Wstęp

Algorytm tworzenia drzew decyzyjnych ID3 jest jednym z prostszych algorytmów ale zarazem daje on dość dobre wyniki. Celem jest oczywiście stworzenie drzewa, które za pomocą wartości atrybutów przyjmowanych przez elementy podzieli nam dziedzinę na klasy równoważności (w domyśle podgrupy mające taką samą wartość jednej ze zmiennych). Oczywiście w typowym przypadku możliwości stworzenia takiego drzewa będzie wiele. Więc ustalamy dodatkowy cel, jakim będzie minimalna wysokość drzewa, liczona jako największa odległość od korzenia do liścia. Algorytm ID3 zawsze (jeżeli to możliwe) stworzy nam drzewo decyzyjne. Natomiast nie zawsze jest to drzewo optymalnej wielkości. Algorytm ID3 jest algorytmem zachłannym, decyzje o rozbudowie drzewa są podejmowane na podstawie przybliżonej oceny każdego z wariantów jakie możemy przyjąć w danym kroku. Raz podjęta decyzja nie jest już zmieniana - nie jest to algorytm adaptacyjny. Przyjrzyjmy się jego działaniu.

3.2 W dużym skrócie

Dopuki każdy z liści drzewa nie jest homogeniczny (zmienne decyzyjne jego elementów nie są jednakowe) powtarzaj:

- Wybierz ten spośród nieużytych jeszcze atrybutów, który minimalizuje średnią entropię (opis średniej entropii dalej)
- Rozwin niehomogeniczne liście względem wybranego atrybutu.

3.3 Formula Entropii

Entropia jest miarą z teorii informacji, charakteryzująca czystość i homogeniczność zbioru atrybutów.

3.3.1 Dane

- nb, liczba instancji w lisciu b
- nbc, liczba instancji w lisciu b należących do klasy c. $nbc \leq nb$
- nt, całkowita liczba instancji we wszystkich liściach

3.3.2 Prawdopodobieństwo

$$P_b = \frac{n_{bc}}{n_b}$$

- Jeżeli wszystkie instancje w grupie są klasyfikowane pozytywnie, wtedy $P_b = 1$ (liść homogeniczny pozytywnie)
- Jeżeli wszystkie instancje w grupie są klasyfikowane negatywnie, wtedy $P_b = 0$ (liść homogeniczny negatywnie)

3.3.3 Entropia

$$Entropia = Sum(c) \left(-\frac{n_{bc}}{n_b} \log_2 \left(\frac{n_{bc}}{n_b} \right) \right)$$

- Entropia jest zerowa jeżeli zbiór jest idealnie homogeniczny
- Entropie wynosi 1 jeżeli zbiór jest idealnie niehomogeniczny ze względu na atrybut (tzn. nie jest on dzielony na żadne podgrupy przez ten atrybut)

3.3.4 Średnia entropia

$$ŚredniaEntropia = Sum(b) \left(\frac{n_b}{n_t} \right) * [Sum(c) \left(-\frac{n_{bc}}{n_b} \log_2 \left(\frac{n_{bc}}{n_b} \right) \right)]$$

3.4 Minimalizacja entropii = Minimalizacja wysokości drzewa ???

Ogólne założenie jest aby tworzyć drzewa decyzyjne optymalnej wielkości ale z praktycznego punktu widzenia nie jest to uzasadnione ze względu na duży koszt obliczeniowy. W zastępstwie korzystamy z przybliżonych procedur tworzenia małych, ale niekoniecznie najmniejszych drzew decyzyjnych.

3.5 Algorytm w pseudokodzie

1. Zainicjuj drzewo (wszystkie elementy przypisane do korzenia, który jest zarazem liściem)
2. Dopuki nie wszystkie liście są homogeniczne powtarzaj
 - Jeżeli nie ma nieużytych atrybutów -> **koniec z błędem**
 - Oblicz średnią entropię dla nieużytych jeszcze atrybutów
 - Wybierz ten atrybut, który minimalizuje średnią entropię (dla którego wyliczony w poprzednim punkcie wskaźnik jest najmniejszy)
 - Rozwijaj niehomogeniczne liście względem wybranego atrybutu (liść staje się węzłem, do którego są "przyczepione" liście powstałe z podziału tego liścia na liście zawierające każdą z przyjmowanych przez wybrany atrybut wartości)
3. Wypisz drzewo -> **poprawny koniec**