

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

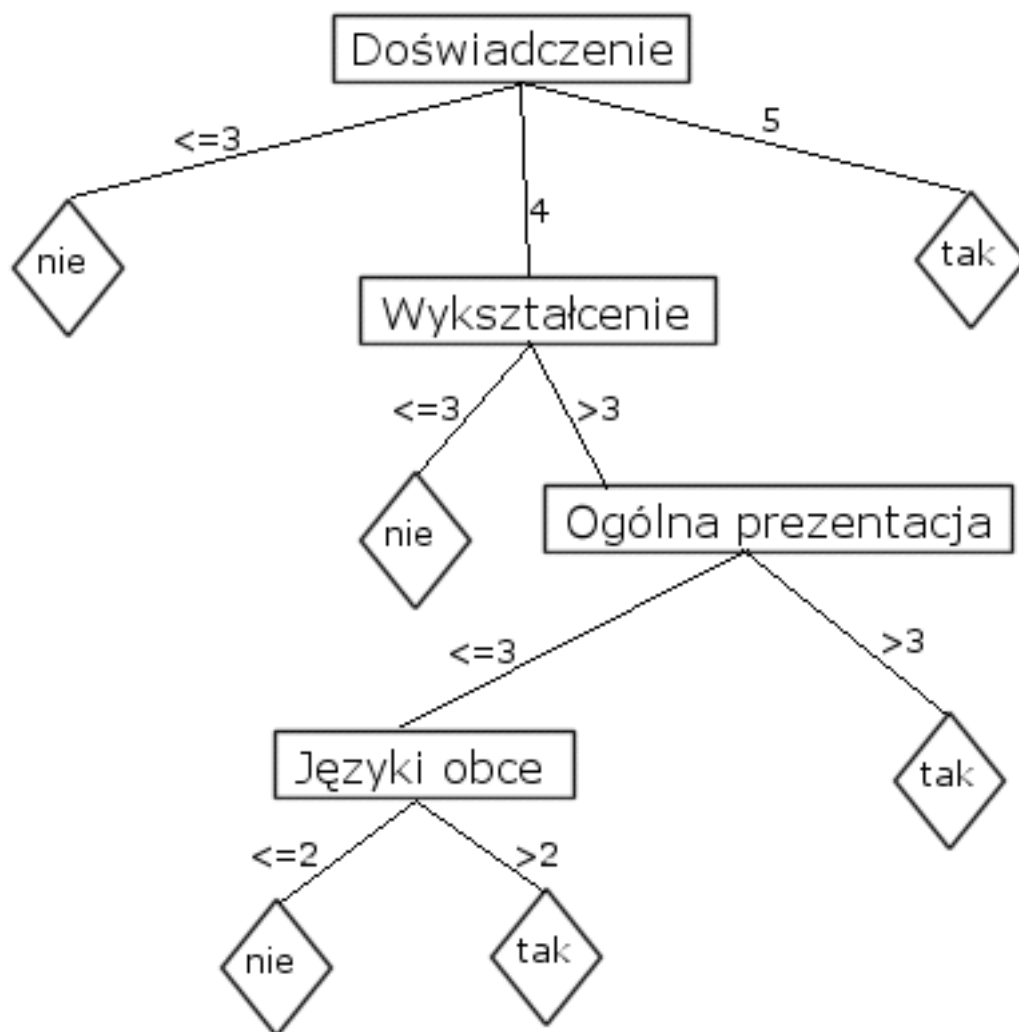
Drzewa decyzyjne w uczeniu maszynowym służy do wyodrębniania wiedzy z zestawu przykładów (patrz eksploracja danych). Zakładamy, że posiadamy zestaw przykładów: obiektów opisanych przy pomocy atrybutów, którym przyporządkowujemy jakąś decyzję (patrz tabela decyzyjna).

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 W dużym skrócie

Dopuki każdy z liści drzewa nie należy do tej samej klasy równoważności (nie jest homogeniczny - jego zmienne decyzyjne nie są jednakowe) powtarzaj:

- Wybierz niehomogeniczny liść
- Zamień ten liść na węzeł testowy dzielący ten podzbiór na tak niehomogeniczne podzbiory jak to możliwe, zgodnie z wyliczeniem entropii

3.2 Trochę dokładniej

- Oblicz entropie dla każdego z atrybutów, które chcemy wykorzystać do tworzenia następnego poziomu drzewa. Wybierz najlepszy z nich (ten który minimalizuje entropię)

3.3 Minimalizacja entropii?

Generalnie 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.

3.4 Wybieranie atrybutu do obliczania entropii

Najważniejszym punktem algorytmu ID3 jest wybór atrybutu do testowania na każdym liściu drzewa.

Procedura:

- Należy sprawdzić jak atrybut dzieli elementy należące do liścia
- Minimalizuj średnią entropię (oblicz entropie dla każdego z węzłów, dla każdego z atrybutów i wybierz ten z najmniejszą entropią)

3.5 Formula Entropii

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

3.5.1 Dane

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

3.5.2 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)

3.5.3 Entropia

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

- 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)

3.5.4 Średnia entropia

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