

WSI Decision Tree Classifier sprawozdanie

Małgorzata Grzanka

05.05.2024

Spis treści

1	Wstęp	1
2	Skrypt przygotowujące dane dla modelu	2
2.1	Zbiór danych	2
2.2	Zasada działania skryptu	2
3	Sposób działania drzewa decyzyjnego ID3	3
3.1	Ogólna zasada działania	3
3.2	Dzielenie danych	3
3.3	Wybór kolumny, po której się dzieli	3
3.4	Kryterium stopu	3
4	Hiperparametry	4
5	Wpływ głębokości na dokładność predykcji	4
5.1	Sposób badania wpływu	4
5.2	Wyniki	5
5.3	Wnioski	7
5.3.1	Wpływ głębokości na dokładność modelu	7
5.3.2	Wpływ głębokości na pozostałe elementy oceny modelu	7

1 Wstęp

Celem zadania była implementacja drzewa decyzyjnego tworzonego algorytmem ID3 z ograniczeniem maksymalnej głębokości drzewa. Do analizy działania stworzonego modelu wykorzystam zbiór danych Cardio Vascular Disease Detection z kaggle, gdzie polem do predykcji jest pole *cardio*. Dodatkowym elementem zadania było stworzenie skryptu przygotowującego dane dla modelu.

2 Skrypt przygotowujące dane dla modelu

2.1 Zbiór danych

Zadany zbiór danych zawierał następujące pola:

- wiek (age) : int (days)
- wzrost (height) : int (cm)
- płeć (gender) : 1: first gender, 2: second gender
- ciśnienie skurczowe (ap_hi) : int
- ciśnienie rozkurczowe (ap_lo) : int
- cholesterol (cholesterol) : 1: normal, 2: above normal, 3: well above normal
- glikaza (gluc) : 1: normal, 2: above normal, 3: well above normal
- palenie (smoke) : binary (0/1)
- spożywanie alkoholu (alco) : binary (0/1)
- aktywność fizyczna (active) : binary (0/1)

Przedstawione pola to cechy pacjenta, a zadaniem modelu jest przewidzenie na ich podstawie, czy pacjent ma chorobę serca, czy nie (pole cardio). Mamy zatem do czynienia z klasyfikacją binarną. Dane przechowywane są w pliku csv.

2.2 Zasada działania skryptu

Skrypt spełnia powyższe cechy:

- Eliminuje niepoprawne dane ze zbioru danych (NaNs albo ujemne, kiedy nie powinno być żadnych ujemnych - np. waga):
 - Usuwa całą kolumnę danych, jeśli jej ilość niepoprawnych danych jest większa niż jakiś zadany *drop_column_percentage* wszystkich danych.
 - Wypełnia niepoprawne dane średnią z danej kolumny, jeśli ilość jej niepoprawnych danych jest mniejsza niż jakiś zadany *drop_column_percentage*, ale większa niż inny, mniejszy zadany *fill_column_percentage*.
 - Usuwa wszystkie próbki z niepoprawnymi danymi w tej kolumnie, jeśli jest ich mniej niż jakiś zadany *fill_column_percentage*.
- Eliminuje outliers'ów (odstające dane) ze zbioru danych. Jako outliers'y traktowane są dane, które są więcej niż 4 odchylenia standardowe od średniej wartości w kolumnie.

3 Sposób działania drzewa decyzyjnego ID3

3.1 Ogólna zasada działania

Drzewo decyzyjne służy do klasyfikacji binarnej, czyli do przewidywania wartości binarnych na podstawie zadanych cech (kolumn/pól) ze zbioru danych. Działa na zasadzie wielokrotnego dzielenia zbioru danych na mniejsze podzbiory na podstawie obecnych w zbiorze cech. Wybierane są te cechy, które po podzieleniu zapewniają najlepszą pewność w przypisaniu określonej przewidywanej wartości do próbki, która znalazła się w tym podziorze.

3.2 Dzielenie danych

Dzielenie danych odbywa się za pomocą zadania pytania Tak / Nie. W przypadku dzielenia po kolumnie zawierającej dane binarne, jest to bardzo proste. Dane, których wartość w tej kolumnie wynosi 1 zostają przydzielone do "lewego" podprzedziału, a dane, których wartość w tej kolumnie wynosi 0, zostają przydzielone do "prawego" podprzedziału. W przypadku danych ciągłych, należy dokonać zamienienia danych na dane binarne. W moim modelu polega to na rozpatrzeniu warunku, czy wartość danej próbki w tej kolumnie jest mniejsza lub równa jakiejś zadanej wartości (wartość ta jest wzięta z określonym krokiem z wszystkich możliwych wartości w tej kolumnie). Następnie podział następuje podobnie jak przy danych binarnych - jeśli warunek jest spełniony, próbka zostaje przydzielona do "lewego" podprzedziału, a jeśli nie - do "prawego".

3.3 Wybór kolumny, po której się dzieli

Cecha, po której dzielony jest zbiór danych lub jakiś jego kolejny podzbiór, powinna zapewniać najlepszą pewność w przypisaniu określonej przewidywanej wartości do próbki, która znalazła się w tym podziorze. Aby znaleźć taką cechę, można zastosować na przykład entropię lub gini impurity. W moim modelu zdecydowałam się na tę drugą opcję. Dla każdego możliwego podziału liczone jest gini impurity i wybierane jest to, które ma najmniejszą wartość (nieczystość podziału jest możliwie najmniejsza).

3.4 Kryterium stopu

Dzielenie zbioru na kolejne podzbiory w moim algorytmie kończy się, gdy:

- podzbiór osiągnie jakąś zadaną głębokość
- liczba próbek w podziorze jest mniejsza, niż zadana maksymalna liczba próbek
- wszystkie próbki w podziorze są jednej klasy (tzw. liść czysty)

W dalszej analizie skupię się jedynie na wpływie maksymalnej głębokości na dokładność precyzji, traktując pozostałe hiperparametry jako zmienne stałe.

4 Hiperparametry

Mój model przyjmuje następujące hiperparametry:

- maksymalna głębokość - maksymalna ilość podzbiorów w jednej gałęzi
- minimalna liczba próbek - ile minimalnie próbek może być w podzbiorze (ustawiona na stałe na 1)
- limit wartości dyskretnych - ile maksymalnie może być różnych wartości dyskretnych w kolumnie numerycznej (stosowana, aby przyspieszyć trenowanie modelu)
- krok podziału wartości dyskretnych - jeśli liczba unikalnych wartości dyskretnych jest większa niż limit, do warunków podziału brane są wartości generowane z zadany krokiem
- kryterium oceny - może być gini albo entropia (entropia niestety jeszcze nie zaimplementowana)

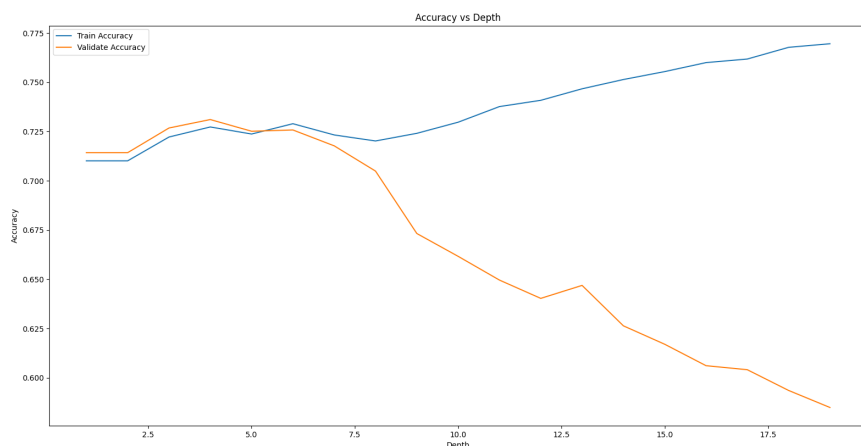
5 Wpływ głębokości na dokładność predykcji

Zbadam wpływ głębokości na dokładność predykcji modelu. W tym celu, minimalna liczba próbek w podzbiorze zostanie ustawiona na 1, aby rozpatrywać jedynie sam wpływ głębokości. Bez zmian pozostają także hiperparametry: limit wartości dyskretnych (300), krok podziału wartości dyskretnych (100) oraz kryterium oceny (gini).

5.1 Sposób badania wpływu

Dla zbadania wpływu głębokości, podzieliłam zbiór danych na 3 podzbiory - zbiór treningowy, zbiór walidacyjny i zbiór testowy. Model wytrenuję na zbiorze treningowym z różnymi głębokościami w przedziale 1-19, a następnie użyję go do predykcji zbiorów treningowego oraz walidacyjnego, aby znaleźć optymalną głębokość. Następnie, dokonam ostatecznej predykcji na zbiorze testowym.

5.2 Wyniki



Rysunek 1: Wpływ głębokości na dokładność predykcji zbioru treningowego i walidacyjnego

Jak widać na powyższym wykresie, dokładność predykcji dla zbioru walidacyjnego rośnie do głębokości równej około 4 - wtedy osiąga największą wartość równą około 73%. Następnie maleje aż do wartości poniżej 60% skuteczności dla głębokości 19.

W przypadku danych treningowych, dokładność predykcji rośnie w nim wraz ze wzrostem głębokości modelu - od wartości 71% do ponad 76%.

```
->->-> Model rating for depth 1<-<-<-  
Confution matrix values (tp, fp, fn, tn): (5094, 1665, 3213, 7036)  
Model accuracy: 0.7131937911571026  
Precision: 0.7536617842876165  
Recall: 0.6132177681473456  
F-measure: 0.6762246117084827  
FPR: 0.19135731525112057  
TNR: 0.8086426847488795
```

Rysunek 2: Wyniki dla danych testowych dla głębokości 1

```
->->-> Model rating for depth 4<-<-<-  
Confution matrix values (tp, fp, fn, tn): (5918, 2283, 2389, 6418)  
Model accuracy: 0.7253057384760113  
Precision: 0.7216193147177173  
Recall: 0.712411219453473  
F-measure: 0.7169857039011388  
FPR: 0.26238363406505  
TNR: 0.7376163659349501
```

Rysunek 3: Wyniki dla danych testowych dla głębokości 4

```
->->-> Model rating for depth 5<-<-<-  
Confution matrix values (tp, fp, fn, tn): (5303, 1700, 3004, 7001)  
Model accuracy: 0.7234242709313264  
Precision: 0.757246894188205  
Recall: 0.6383772721800891  
F-measure: 0.692749836708034  
FPR: 0.1953798413975405  
TNR: 0.8046201586024595
```

Rysunek 4: Wyniki dla danych testowych dla głębokości 5

```
->->-> Model rating for depth 8<-<-<-  
Confution matrix values (tp, fp, fn, tn): (6124, 2959, 2183, 5742)  
Model accuracy: 0.6976716839134525  
Precision: 0.6742265771220962  
Recall: 0.7372095822800048  
F-measure: 0.7043128234617596  
FPR: 0.34007585335018964  
TNR: 0.6599241466498104
```

Rysunek 5: Wyniki dla danych testowych dla głębokości 8

```
->->-> Model rating for depth 19<-<-<-<
Confution matrix values (tp, fp, fn, tn): (5168, 3919, 3139, 4782)
Model accuracy: 0.5850188146754468
Precision: 0.5687245515571696
Recall: 0.6221259179005658
F-measure: 0.5942278946763252
FPR: 0.45040799908056545
TNR: 0.5495920009194345
```

Rysunek 6: Wyniki dla danych testowych dla głębokości 19

5.3 Wnioski

5.3.1 Wpływ głębokości na dokładność modelu

Spadek dokładności dla danych walidacyjnych, a wzrost dokładności dla danych treningowych przy wzroście maksymalnej głębokości wiąże się z zjawiskiem przeczenia (over-fitting). Przez zbyt dokładne podzielenie drzewa poprzez dane treningowe, model nie oddaje ogólnego trendu w korelacji między danymi, a jest jedynie dobrym odwzorowaniem podziału danych testowych. Z tego powodu, predykcje na danych testowych są dobre, ponieważ model dobrze się do nich dopasowuje. Natomiast przez to, że ogólny trend nie jest w nim uwzględniony, predykcje próbek innych niż treningowe okazują się bardzo często niepoprawne (dane walidacyjne).

Wzrost dokładności zarówno dla danych walidacyjnych, jak i treningowych przy mniejszych dokładnościach wiąże się z niedouczeniem (under-fitting). Dane są zbyt mało podzielone, aby oddawać trend w korelacji między cechami, a przewidywana klasa, więc zarówno predykcje dla danych walidacyjnych, jak i testowych, są niepoprawne.

Za optymalna wartość głębokości uznaje wartości z okolicy punktu maksymalnego wykresu danych walidacyjnych. Są to wartości 3, 4 oraz 5, z czego 4 daje najlepszą możliwą predykcję dla danych testowych - 0.725%.

5.3.2 Wpływ głębokości na pozostałe elementy oceny modelu

Do pozostałych elementów oceny modelu zaliczamy:

- precision - poprawnie przewidziane jedynki / wszystkich przewidzianych jedynek (okładność pozytywnych przewidywań modelu).
- recall - poprawnie przewidziane jedynki / wszystkie jedynki w danych (zdolność do przewidywania jedynek)
- true negative rate (TNR) - poprawnie przewidziane zera / wszystkie zera w danych (zdolność przewidywania zer)

Jak widać w zamieszczonych wyżej wynikach, wzrost maksymalnej głębokości podczas trenowania modelu spowodował spadek TNR oraz Precision oraz wzrost Recall. Maksymalne wartości Precision oraz TNR i minimalne Recall osiągnęły dla wartości z przecięcia się krzywych dokładności dla danych walidacyjnych i treningowych (głębokości 5 oraz 6). Jednak dla bardzo dużych głębokości (19), przeuczenie jest zbyt duże i wszystkie elementy są negatywne.

Ze względu na wzrost dopasowania model zaczął przewidywać więcej jedynek, przez co procent dobrze przewidzianych jedynek wzrósł (Recall). Stało się to jednak kosztem spadku dokładności w przewidywaniu zer (TNR) i wzrostu liczby niepoprawnie przewidzianych jedynek (Precision). Inaczej mówiąc, zbyt duża ilość przewidzianych jedynek sprawiła, że wiele z nich zostało przewidzianych niepoprawnie, a liczba przewidzianych zer spadła, co zmniejszyło dokładność przewidywania zer.