

Twoenie drzewa - algorytm rekurencyjny.

- dzielimy obszar $X = \{(x_1, x_2, \dots, x_p)\}$ na prostokąty parami wartości R_j ; na poziomie: jeden prostokąt X
- jeśli mamy pewien podzest X na R_1, R_2, \dots, R_j , to wybieramy wszystkie możliwe podziały dowolnego z prostokątów R_k otrzymane przez hiperpłaszczyzny $\{x_j < s\}$;

$$R^{(1)}(j, s) = R_k \cap \{(x_1, \dots, x_p) : x_j < s\}$$

$$R^{(2)}(j, s) = R_k \cap \{(x_1, \dots, x_p) : x_j \geq s\}$$



• wybieramy najlepszy podział - ten taki, który minimalizuje:

np. dla ~~drzewa~~ poprzedniej regresji

$$\sum_{i: x^{(i)} \in R^{(1)}(j, s)} (y_i - \hat{y}_{R^{(1)}})^2 + \sum_{i: x^{(i)} \in R^{(2)}(j, s)} (y_i - \hat{y}_{R^{(2)}})^2$$

(+ suma dla pozostałych prostokątów)

dla klasyfikacji:

$$\hat{p}_{mk} = \frac{\# \text{ punktów w } R_m, \text{ które są klasy } k}{\# \text{ punktów w } R_m}$$

$$\hat{p}_{m1} + \hat{p}_{m2} + \dots + \hat{p}_{mK} = 1$$

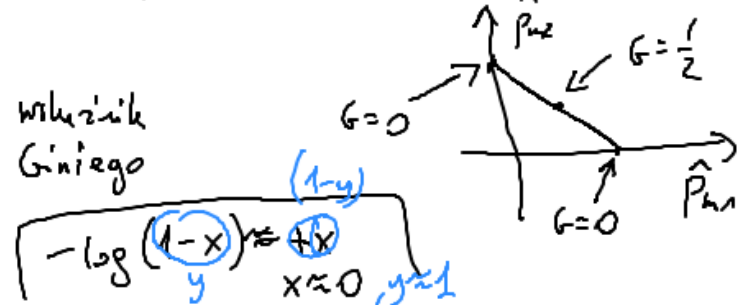
minimalizujemy
 odwołujemy
 z obu stron
 po wszystkich
 m z wagami
 od dowolnie punktów

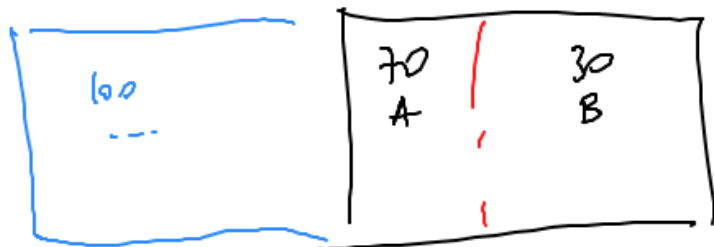
$$E = 1 - \max_k \hat{p}_{mk}$$

$$G = \sum_{k=1}^K \hat{p}_{mk} (1 - \hat{p}_{mk}) = 1 - \sum_{k=1}^K \hat{p}_{mk}^2$$

$$D = - \sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk} \quad \text{- entropia}$$

wskaznik
 Ginięgo



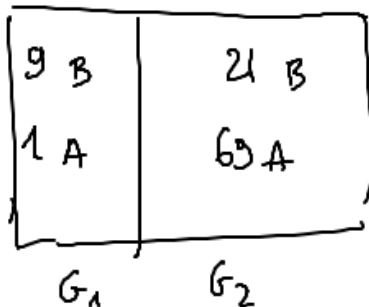


$$G = 1 - \left(\left(\frac{30}{100} \right)^2 + \left(\frac{70}{100} \right)^2 \right) = 1 - 0.09 - 0.49 = 0.42$$

średnia: $\frac{100}{100} \cdot G = 0.42$

podział:

0.5



$$G_1 = 1 - \left(\left(\frac{9}{10} \right)^2 + \left(\frac{1}{10} \right)^2 \right) = 1 - 0.81 - 0.01 = 0.18$$

$$G_2 = 1 - \left(\left(\frac{21}{90} \right)^2 + \left(\frac{69}{90} \right)^2 \right) \approx 0.36$$

~~średnia: $G_1 + G_2 =$~~

średnia: $0.9 \cdot G_2 + 0.1 \cdot G_1 \approx 0.32 + 0.02 = 0.34$
0.45 0.05

Tracimy dane z pomocą algorytmu zachłanego tak długo, jak się da, lub zatrzymujemy się po osiągnięciu pewnej wielkości drzewa (np. liczba liści, głębokość, luba głębokość i liści). Następnie przycinamy.

Przycinanie

"cost complexity pruning"

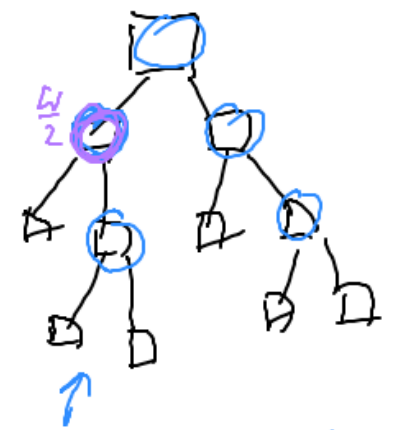
(przycinanie węzła korzeni (?!))

$\alpha \geq 0$

dla regresji:

minimalizujemy

$$(*) \quad \underbrace{\sum_{m=1}^{|T|} \sum_{i: x_i \in R_m} (y_i - \hat{y}_{R_m})^2}_{S} + \underbrace{\alpha |T|}_{\text{liczba liści w drzewie}}$$



dla każdego węzła
wycinamy kawałek zmiłci
 α , który lepiej zaczyna
się z optymalności przycinania
drzewa w tym węzle

po wszystkich możliwych poddrzewach T naszego drzewa.

- zaczynamy od wyjściowego drzewa T_1
- rozbieramy najmniejsze α , dla którego (*) jest minimalna dla pełnego
właściwego poddrzewa T_{A+1} drzewa $T_A \rightarrow T_{A+1}$
- porównujemy dla drzewa T_{A+1} i T_A

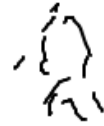
Np. dla węzła 0:
suma S się zmniejsza po
jego obcięciu o pewną liczbę
 w
liczba liści zmniejszamy się o 2
 $(w - 2\alpha) \leq 0 \quad \alpha \geq \frac{w}{2}$

(*) zmniejszy się $0 \rightarrow (w - 2\alpha) \leq 0 \quad \alpha \geq \frac{w}{2}$

Jaką 2 wybrać?

$L=0$
długo
opóźnienia

L duże
Jakość
kwalit.



- 1) Treningowy drzewo używając wielu danych (trainingowy)
i szukamy i sprawdzamy drzewo na pozostałych danych (testowy)
- 2) Użył wchodzący kategoriej