

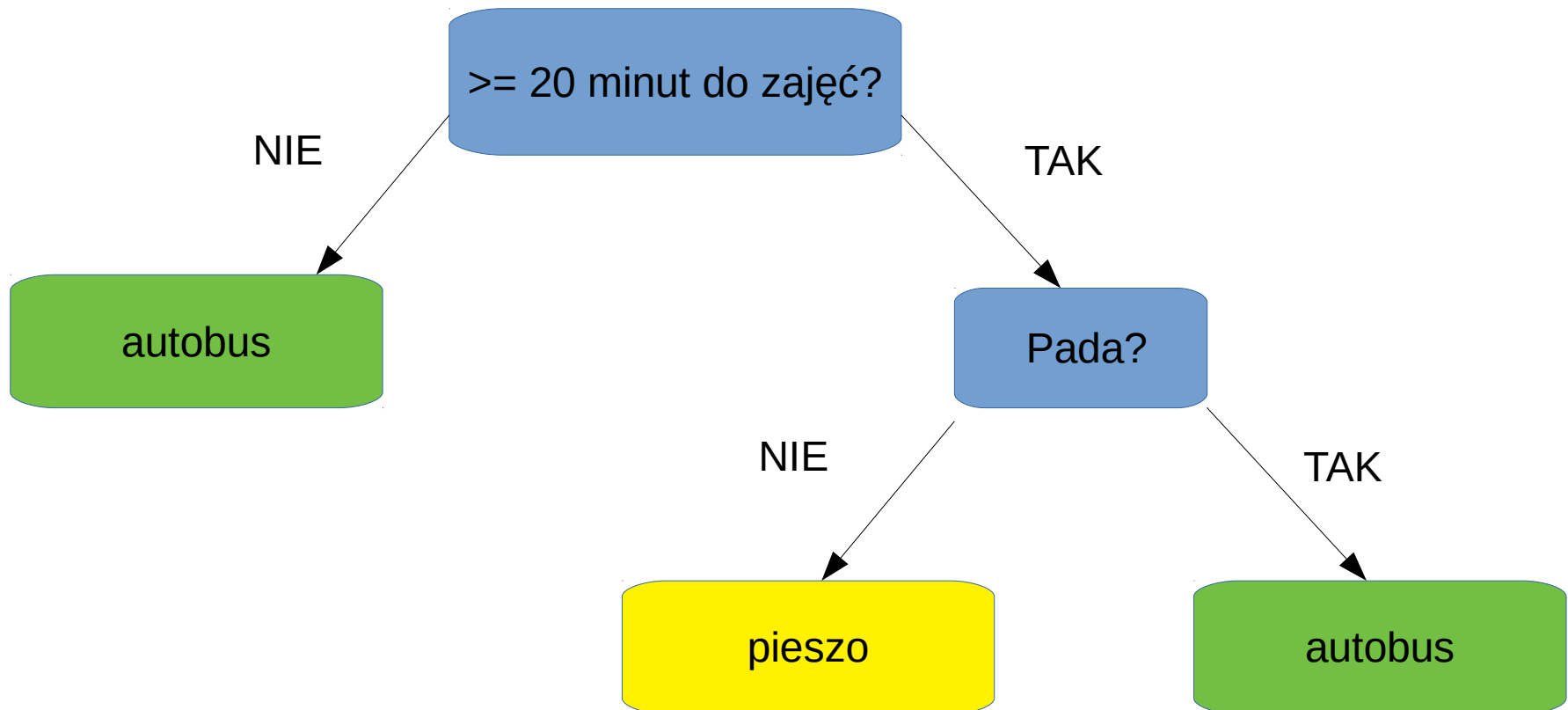
Drzewa decyzyjne

Kamil Michalak

A decorative blue wavy line that spans the width of the slide, positioned at the bottom. It has a smooth, undulating shape, with a darker blue area on the left and a lighter blue area on the right, separated by a thin white line.

Czym są drzewa decyzyjne?

- Decyzja/klasyfikacja na podstawie atrybutów.
- Testy/pytania w węzłach, odpowiedzi w liściach
- Jak dotrzeć na uczelnię?

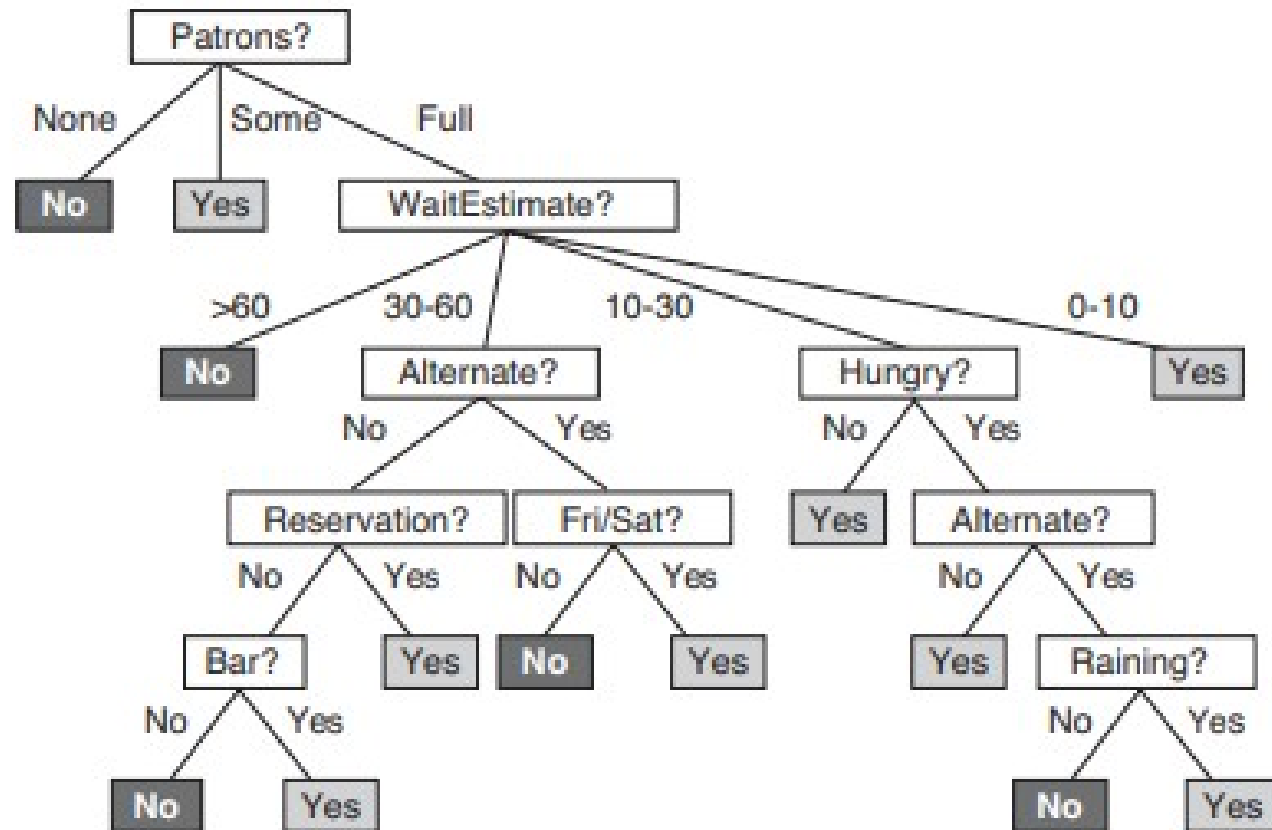


Przykład (*AI: A Modern Approach*)

Czy zaczekać na stolik w restauracji?

- *Alternate*: czy jest inna restauracja w pobliżu.
- *Bar* : czy restauracja ma bar, w którym można zaczekać.
- *Fri/Sat*: czy jest piątek lub sobota.
- *Hungry*: czy jesteśmy bardzo głodni.
- *Patrons*: ile osób jest w restauracji (*None/Some/Full*).
- *Price*: przedział cenowy (\$, \$\$, \$\$\$).
- *Raining*: czy pada.
- *Reservation*: czy mamy rezerwację.
- *Type*: rodzaj restauracji (*French/Italian/Thai/burger*).
- *WaitEstimate*: czas oczekiwania oszacowany przez właściciela (0–10 minut, 10–30, 30–60, lub >60)

Przykład (*AI: A Modern Approach*)



Utworzenie drzewa decyzyjnego

```
function DECISION-TREE-LEARNING(examples, attributes, parent_examples) returns  
a tree  
  
  if examples is empty then return PLURALITY-VALUE(parent_examples)  
  else if all examples have the same classification then return the classification  
  else if attributes is empty then return PLURALITY-VALUE(examples)  
  else  
     $A \leftarrow \operatorname{argmax}_{a \in \text{attributes}} \text{IMPORTANCE}(a, \text{examples})$   
    tree  $\leftarrow$  a new decision tree with root test A  
    for each value  $v_k$  of A do  
       $\text{exs} \leftarrow \{e : e \in \text{examples} \text{ and } e.A = v_k\}$   
      subtree  $\leftarrow$  DECISION-TREE-LEARNING(exs, attributes – A, examples)  
      add a branch to tree with label (A =  $v_k$ ) and subtree subtree  
  return tree
```

Zbiór danych

Example	Input Attributes										Goal
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	<i>WillWait</i>
x₁	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>French</i>	<i>0–10</i>	<i>y₁ = Yes</i>
x₂	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Thai</i>	<i>30–60</i>	<i>y₂ = No</i>
x₃	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Some</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Burger</i>	<i>0–10</i>	<i>y₃ = Yes</i>
x₄	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Thai</i>	<i>10–30</i>	<i>y₄ = Yes</i>
x₅	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>Full</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>French</i>	<i>>60</i>	<i>y₅ = No</i>
x₆	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$</i>	<i>Yes</i>	<i>Yes</i>	<i>Italian</i>	<i>0–10</i>	<i>y₆ = Yes</i>
x₇	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>None</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Burger</i>	<i>0–10</i>	<i>y₇ = No</i>
x₈	<i>No</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$</i>	<i>Yes</i>	<i>Yes</i>	<i>Thai</i>	<i>0–10</i>	<i>y₈ = Yes</i>
x₉	<i>No</i>	<i>Yes</i>	<i>Yes</i>	<i>No</i>	<i>Full</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Burger</i>	<i>>60</i>	<i>y₉ = No</i>
x₁₀	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>Italian</i>	<i>10–30</i>	<i>y₁₀ = No</i>
x₁₁	<i>No</i>	<i>No</i>	<i>No</i>	<i>No</i>	<i>None</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Thai</i>	<i>0–10</i>	<i>y₁₁ = No</i>
x₁₂	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Burger</i>	<i>30–60</i>	<i>y₁₂ = Yes</i>

Figure 18.3 Examples for the restaurant domain.

Wynikowe drzewo

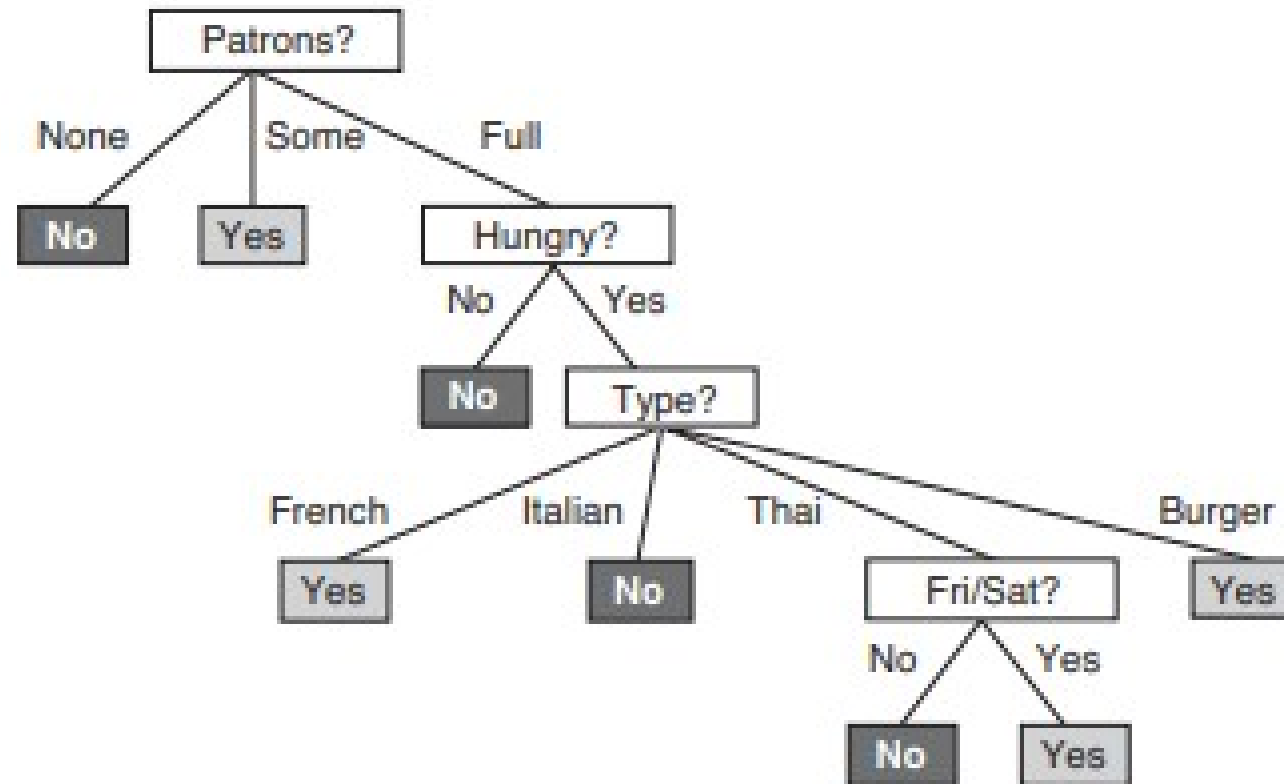


Figure 18.6 The decision tree induced from the 12-example training set.

Wybór cechy do testowania

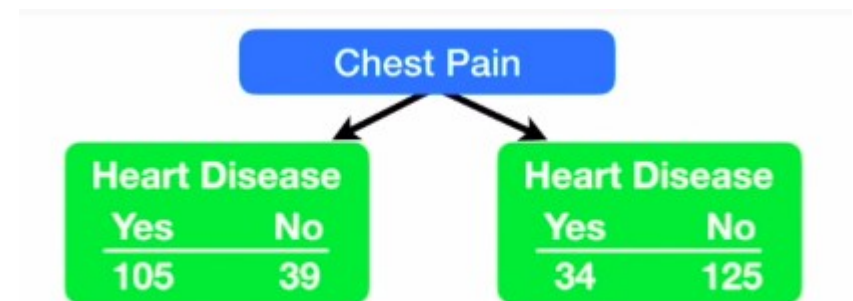
- Sposób 1: maksymalizowanie spadku entropii
- Entropia – niepewność wystąpienia danego zdarzenia w danej chwili, jeśli zdarzenie jest pewne to entropia wynosi 0.
- Wzór na entropię zmiennej losowej V o zbiorze wartości $\{v_1, \dots, v_n\}$

$$H(V) = - \sum_{k=1}^n P(v_k) \cdot \log_2(P(v_k))$$

Wybór cechy do testowania

- Sposób 2: „Gini impurity”

Chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
No	No	No	No
Yes	Yes	Yes	Yes
Yes	Yes	No	No
Yes	No	???	Yes
etc...	etc...	etc...	etc...



Wybór cechy do testowania

- Sposób 2: „Gini impurity”



- $G(L) = 1 - P(\text{„Yes"})^2 - P(\text{„No"})^2 = 0.395$
- $G(R) = 0.336$
- $G(\text{„Chest Pain"}) = \text{średnia ważona z L i R}$
- Wybieramy cechę z najmniejszą wartością G

Chest Pain

Heart Disease

Yes	No
7	26

Heart Disease

Yes	No
6	76

$$G(\text{„Chest Pain”}) = 0.29$$

Good Circ.

Blocked

100/33

Chst Pn

13/102

Zostaje liściem

17/3

7/22

$$G(\text{węzeł}) = 0.2$$

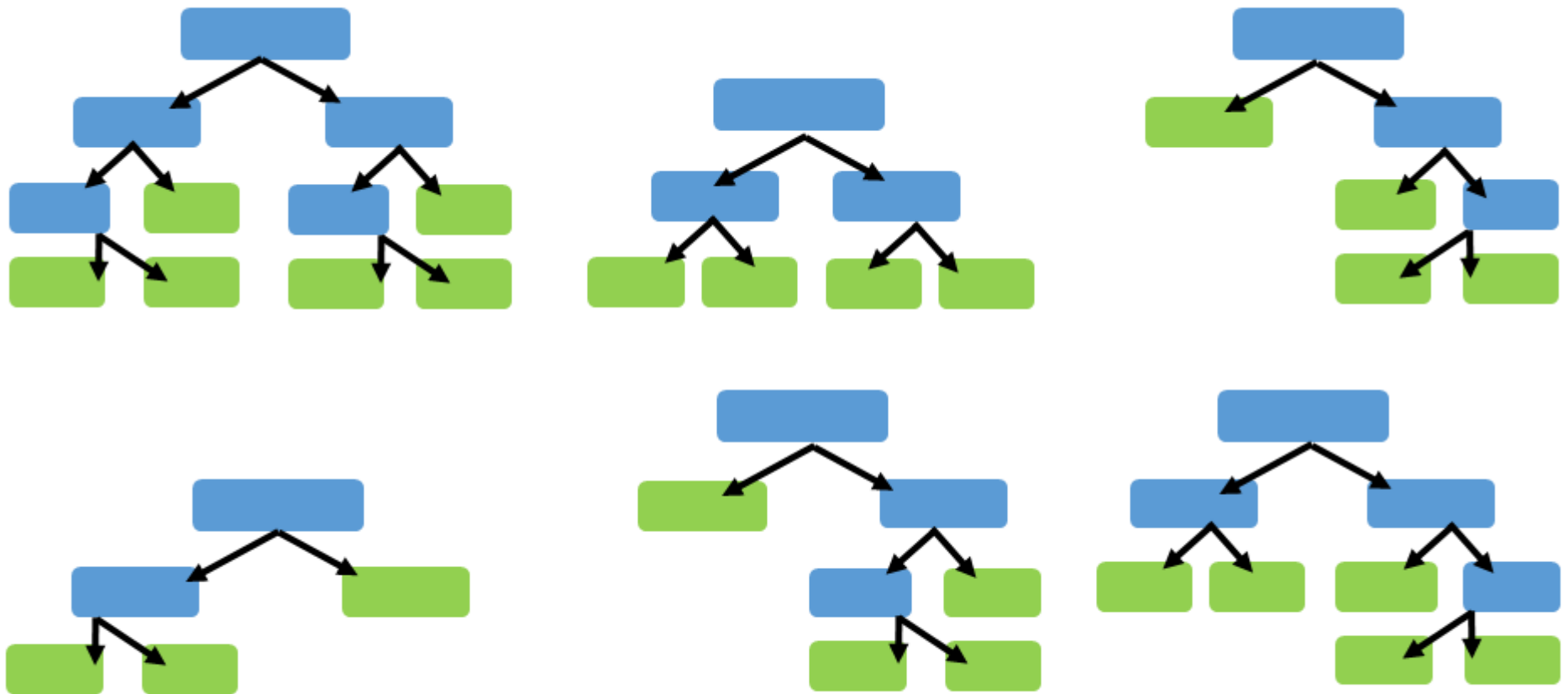
Zastosowania drzew decyzyjnych

- Medycyna – diagnozowanie pacjentów na podstawie reguł znalezionych w bazie danych
- Ocena inwestycji
- Udzielanie kredytów

Problem

- „overfitting” - drzewo działa dobrze dla zbioru uczącego, ale nie dla dowolnych danych

Las losowy (Random forest)



Las losowy (Random forest)

- Tworzymy n drzew
- Dokonujemy klasyfikacji obiektu na każdym drzewie
- Ostateczna klasyfikacja – taka jak na większości drzew.

Tworzenie lasu losowego

- Krok 1: Tworzymy „bootstrapped dataset”
Dla oryginalnego k -elementowego zbioru danych losujemy k -elementów z powtórzeniami.
- Krok 2: Tworzymy n drzew decyzyjnych.
W każdym kroku tworzenia drzewa wybieramy najlepszą cechę jedynie z podzbioru wszystkich cech.

Testowanie lasu losowego

- Bierzemy dane, które nie znalazły się w „bootstrapped dataset” („out-of-bag”)
- Dla każdego rekordu sprawdzamy czy jest dobrze klasyfikowany przez las
- Miarą skuteczności lasu jest % poprawnie sklasyfikowanych danych „out-of-bag”.

Testowanie lasu losowego

- Teraz możemy sprawdzić jak liczne podzbiory cech przy generowaniu lasu dają najlepszą skuteczność.
- Zwykle $\sim \sqrt{\text{liczba cech}}$

Dziękuję za uwagę