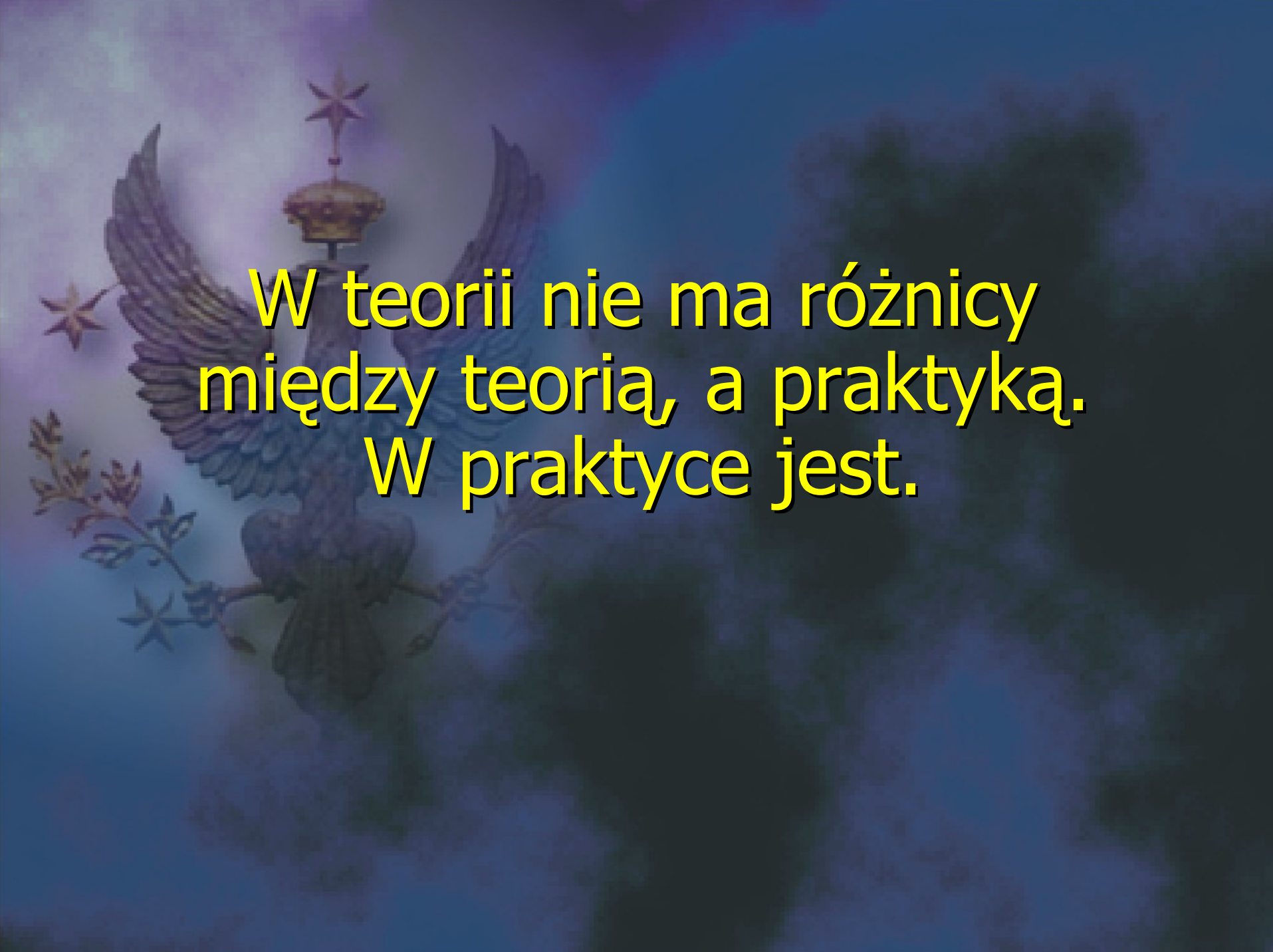


The background of the slide features a faint, stylized version of the Polish coat of arms, which is a white eagle with its wings spread, perched on a golden crown. The eagle is set against a light blue and white background with some floral and star motifs. The overall background of the slide is a dark blue gradient with a subtle pattern of trees.

Drzewa decyzyjne

Marcin S. Szczuka

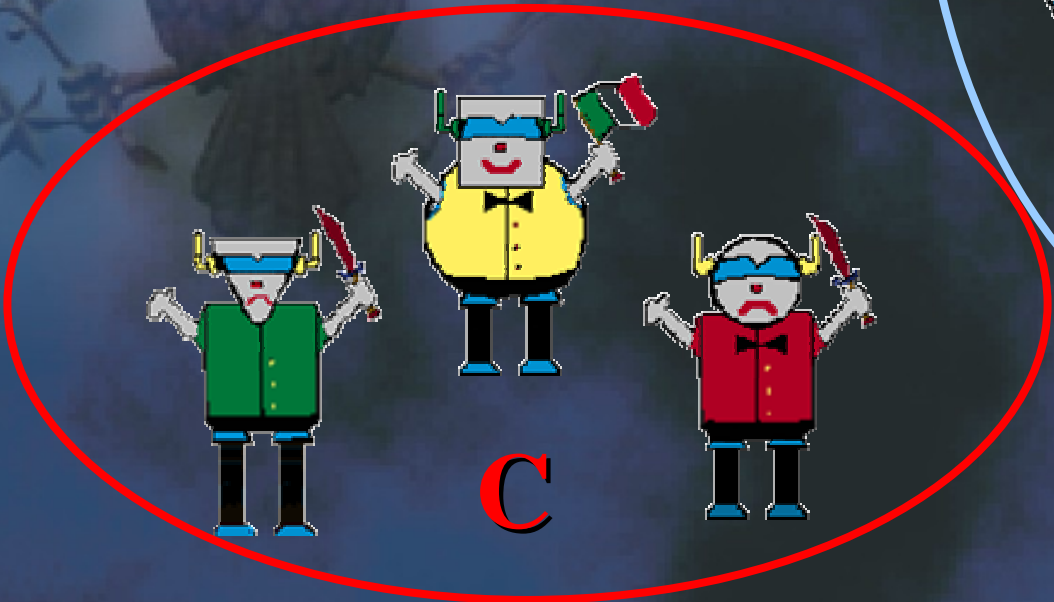
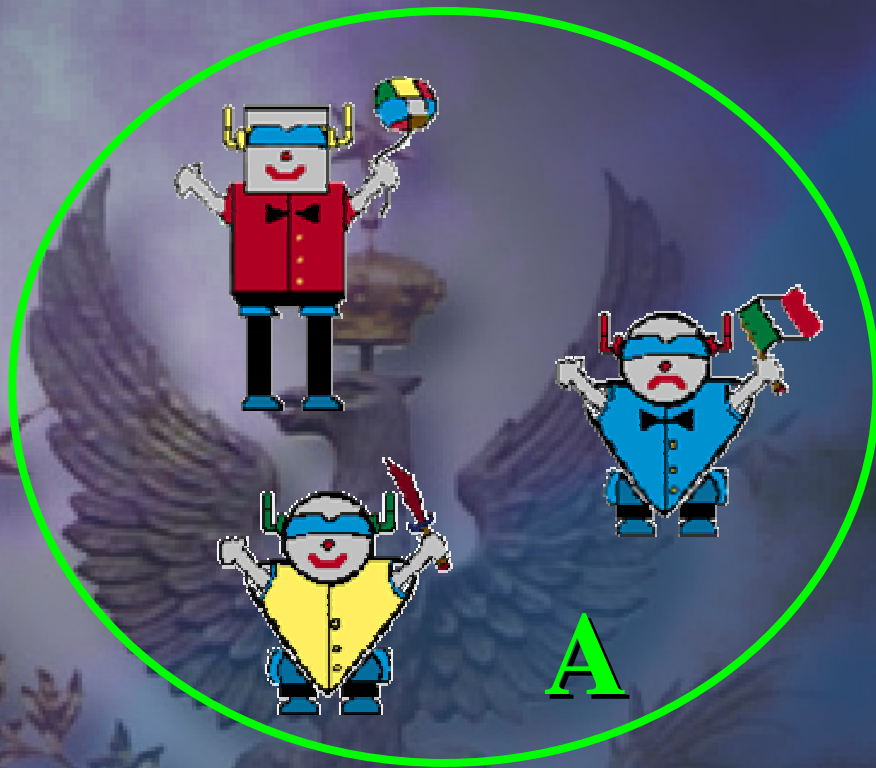
Wykład 3 – część 1

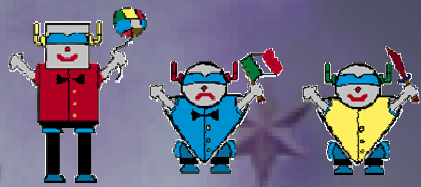


W teorii nie ma różnicy
między teorią, a praktyką.
W praktyce jest.

Drzewa decyzyjne

- q Drzewo – skierowany, acykliczny graf planarny.
- q Liść – wierzchołek bez wychodzących krawędzi.
- q Wierzchołki wewnętrzne (nie liście) reprezentują testy.
- q Każda krawędź odpowiada wynikowi odpowiedniego testu.
- q Liście odpowiadają decyzjom (wg. hipotezy).



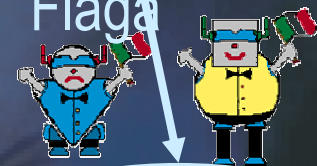


Trzyma

Balon

Miecz

Flaga



Wzrost

Uśmiech

Kolor

Wysoki

Niski

Tak

Nie

Zieleń

Czerw.

Niebieski

Żółty

A

B

C

A

C

C

B

A





Trzyma

Balon

Miecz

Flaga

Wzrost

Uśmiech

Kolor

Wysoki

Niski

Tak

Nie

Zieleń

Czerw.

Niebieski

Żółty

A

B

C

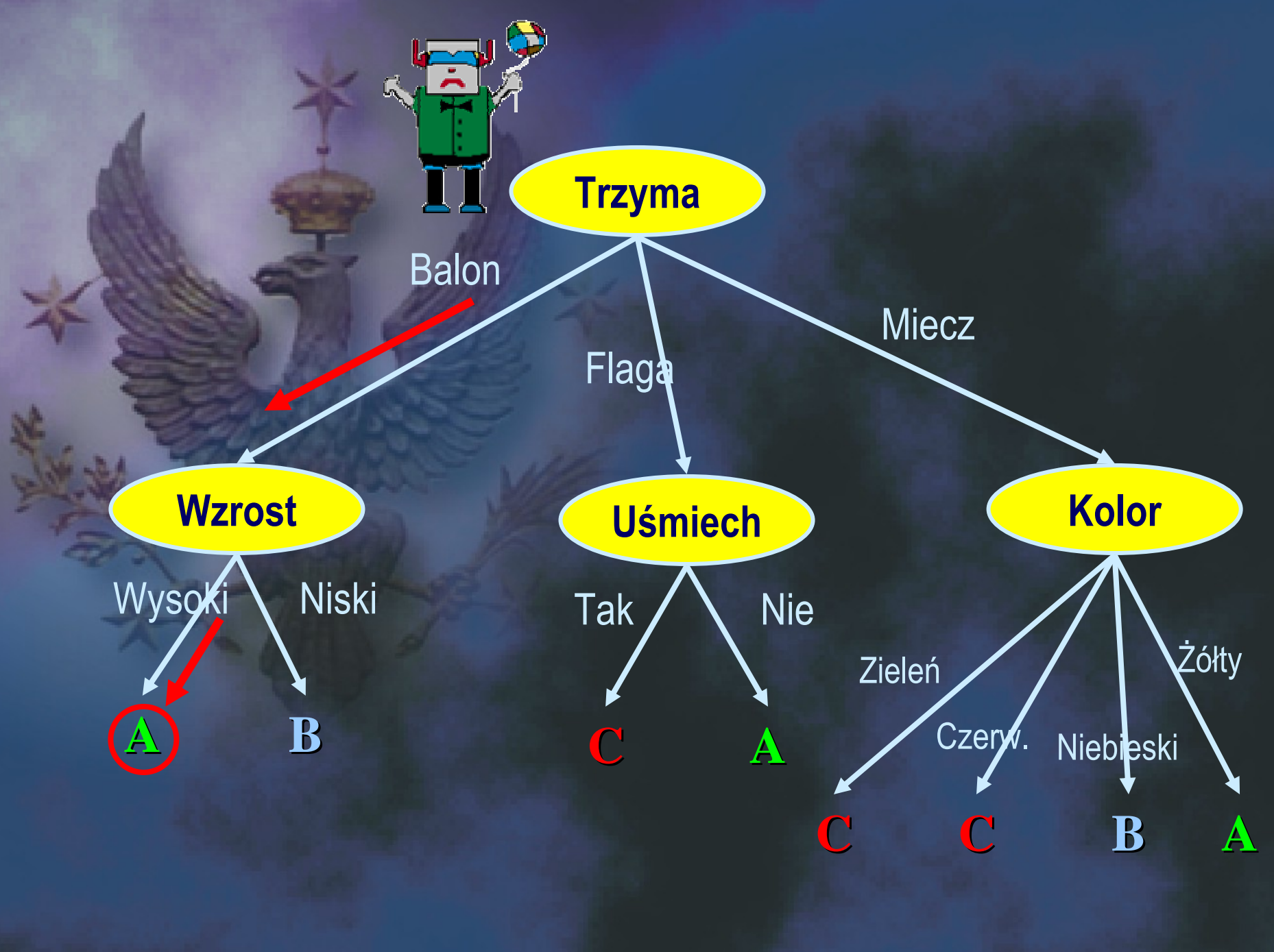
A

C

C

B

A



The background of the slide features a faint, stylized coat of arms of Poland, which includes a white eagle with wings spread, a crown on its head, and a sword above the crown. The eagle is set against a blue background with a subtle floral pattern.

Zadanie

Mając zbiór etykietowanych przykładów skonstruuj drzewo, które najlepiej przybliży proces podejmowania decyzji dla tych przykładów.

Notacja

Każdy przykład jest opisany przez zbiór atrybutów:

$$(a_1(x), \dots, a_n(x))$$

gdzie $x \in X$ i a_1, \dots, a_n są atrybutami (cechami) takimi że:

$$A_i: X \rightarrow V_i$$

V_i nazywamy przestrzenią (wartości) atrybutu.

Np. atrybut **Kolor** ma przestrzeń

{Czerwony, Zielony, Niebieski, Żółty}

Przykład tablicy decyzyjnej

Trzyma a_1	Uśmiech a_2	Wzrost a_3	Kolor a_4	Klub d
Flaga	Nie	Niski	Niebieski	A
Miecz	Nie	Wysoki	Zielony	B
Flaga	Tak	Wysoki	Żółty	C



Budowanie drzewa – ogólna idea

• Podejście „dziel i rządź”.

q Wybierz „najlepszy” atrybut i ustaw jako test w korzeniu.

q Stwórz gałąź dla każdej wartości atrybutu. Usuń atrybut z dalszych rozważań.

q Na końcu każdej gałęzi konstruuj (rekurencyjnie) drzewo z przykładów odpowiadających tej gałęzi.



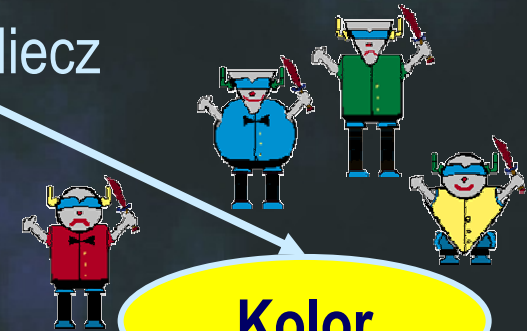
Trzyma

Balon

Miecz



Flaga



Wzrost

Uśmiech

Kolor

Wysoki

Niski

Tak

Nie

Zielony

Czerw.

Niebieski

Żółty

A

B

C

A

C

C

B

A



ID3 - Notacja

S – próbka treningowa, zbiór przykładów

$S_{a=v}$ – podzbiór przykładów mających v jako wartość na atrybucie a

d_l – decyzja dla liścia l w drzewie

a_n – test w wierzchołku wewnętrznym drzewa n

A – zbiór dostępnych atrybutów (testów)

$\text{dec}(S, d)$ – zwraca najczęstszą decyzję w S lub wartość domyślną d , gdy S pusty.

ID3 - algorytm

```
ID3( $S, A, d$ )  
create root;  
if ( $S$  pusty) or (wszystkie przykłady z tą samą decyzją) or  
  ( $A$  pusty) then  
  create-leaf  $l$ ;  $d_1 := \text{dec}(S, d)$ ;  
  return (drzewo z pojedynczym liściem  $l$ );  
endif;  
create node  $n$ ;  
 $a_n := \text{choose-attribute}(S, A)$ ; root :=  $n$ ;  
 $d := \text{dec}(S, d)$ ;  
forall  $v \in V_{a_n}$  do  
  add-subtree ID3( $S_{a_n=v}$ ,  $A - \{a_n\}$ ,  $d$ );  
end;  
return (tree);
```

ID3 – wybór atrybutu

`choose-attribute(S, A)`;

Zwraca atrybut z A , który prowadzi do najlepszego podziału S tj. atrybut będący najlepszym dyskryminatorem ze względu na decyzję.

ID3 - dyskusja

- q Prosty ID3 jest zwykle za prosty. W rzeczywistości korzysta się z jego rozszerzeń np. C4.5 i C5 .
- q Złożoność jest rozsądna – $O(|S|n \log n)$.
- q Algorytm jest dokładny (kompletny). Każda hipoteza może być skonstruowana.

Drzewa – co dalej

Na następnym wykładzie:

- q Wybieranie atrybutu dla drzewa na podstawie miary entropijnej.
- q Rozszerzenia i uzupełnienia algorytmu konstrukcji drzewa.
- q Więcej o złożoności i rozwiązywaniu problemów związanych z nią.