

# Eksploracja danych

## KLASTERYZACJA I SEGMENTACJA

***Wojciech Waloszek***

*wowal@eti.pg.gda.pl*

***Teresa Zawadzka***

*tegra@eti.pg.gda.pl*

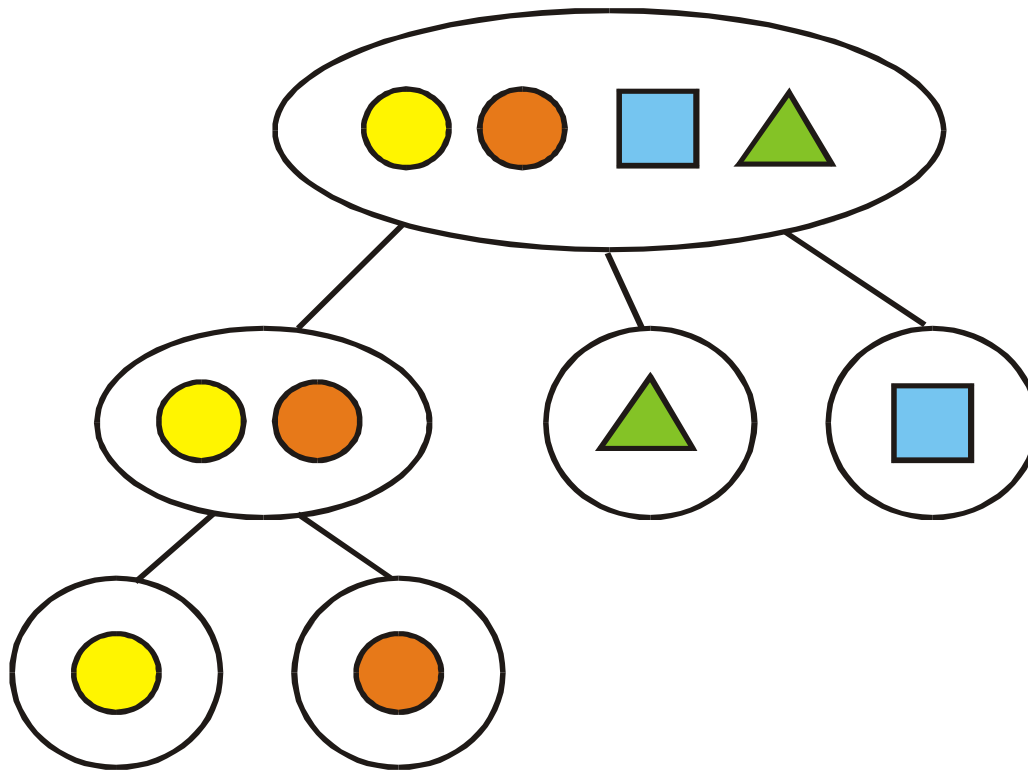
*Katedra Inżynierii Oprogramowania  
Wydział Elektroniki, Telekomunikacji i Informatyki  
Politechnika Gdańska*



# COBWEB

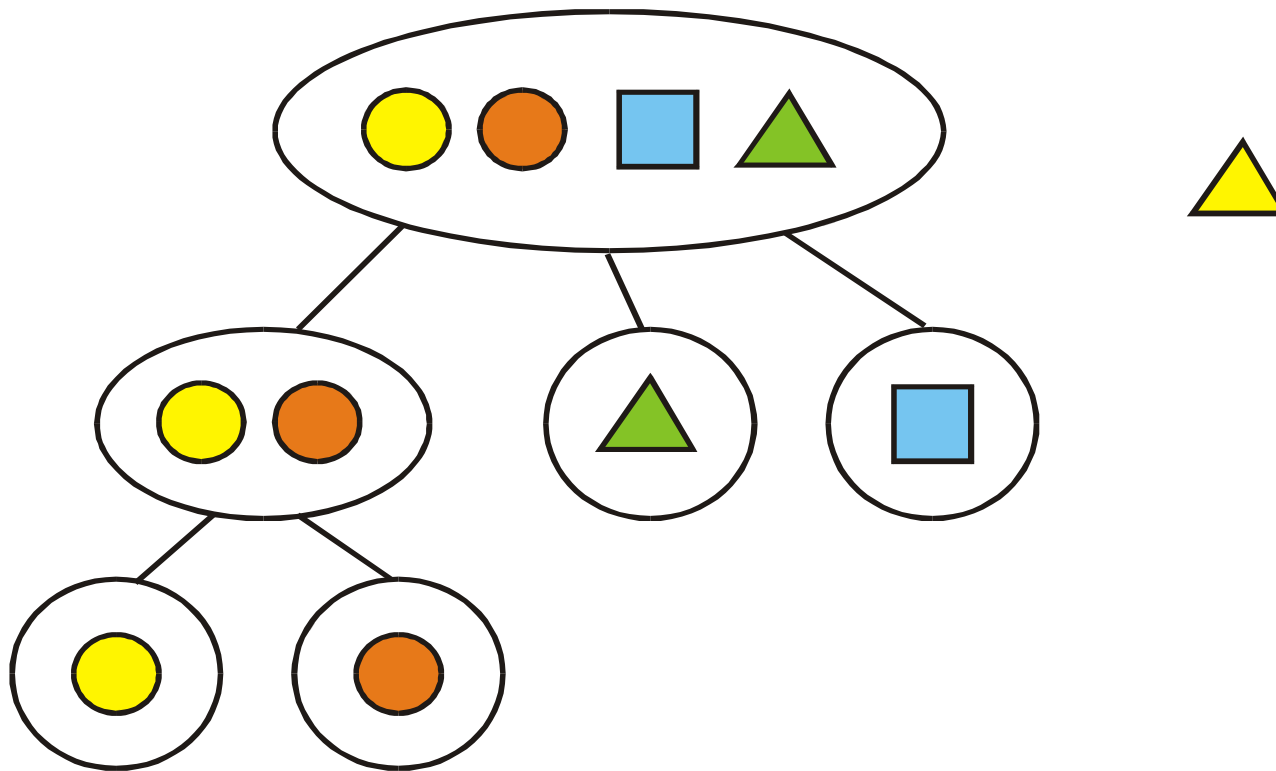
- Klasteryzacja polega na wygenerowaniu pewnych zasad podziału danych na zasadzie podobieństwa,
- Algorytm COBWEB jest przykładem algorytmu grupującego (tworzącego klastry),
- Algorytm COBWEB generuje hierarchiczny podział danych w postaci drzewa.

# Drzewo COBWEB



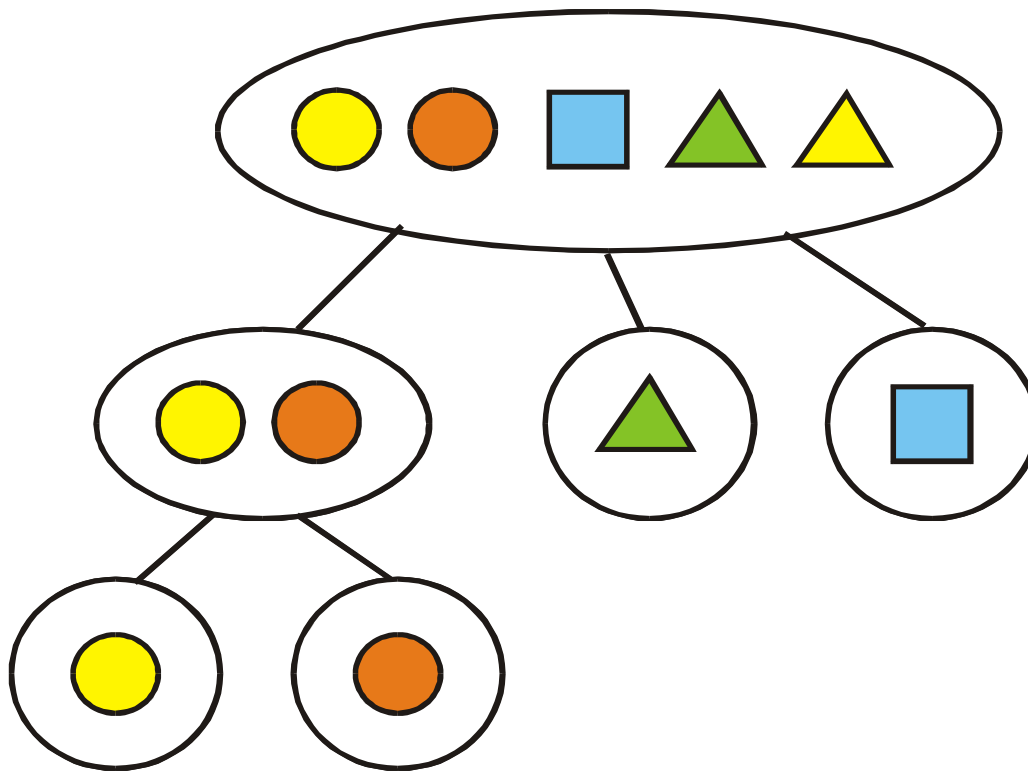
# COBWEB – grupowanie

Założmy, że chcemy przeprowadzić proces grupowania dla kolejnego przykładu:



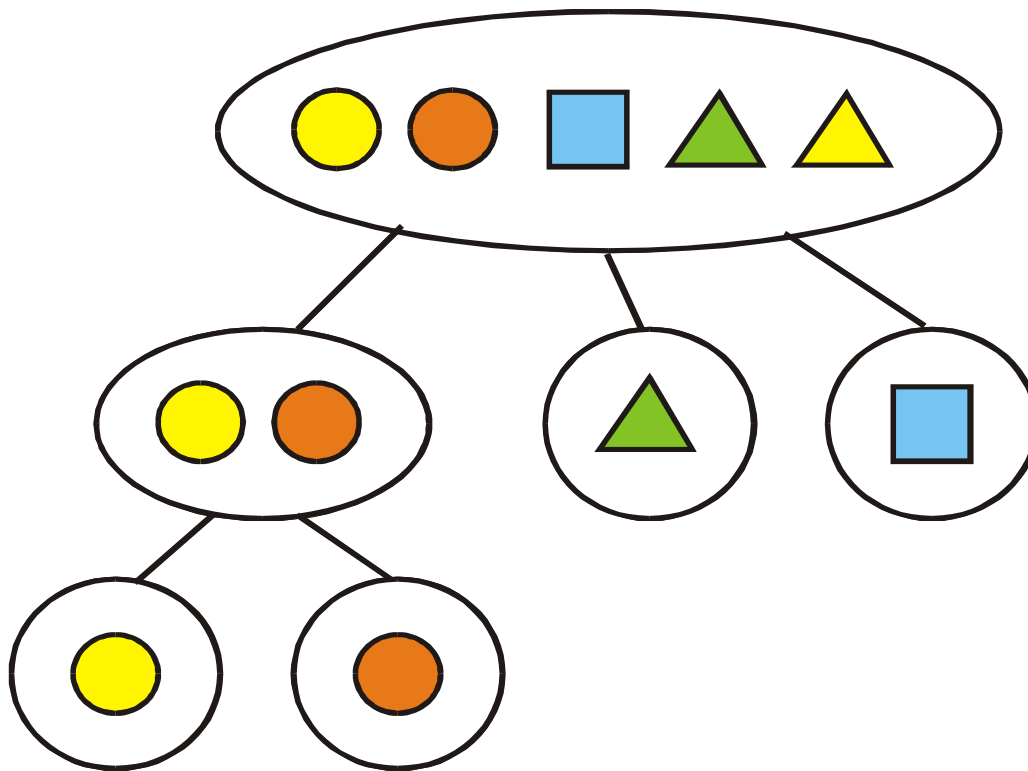
# COBWEB – grupowanie

COBWEB dokłada przykład do głównej kategorii.



# COBWEB – grupowanie

Następnie próbuje różnych możliwości dalszego grupowania.

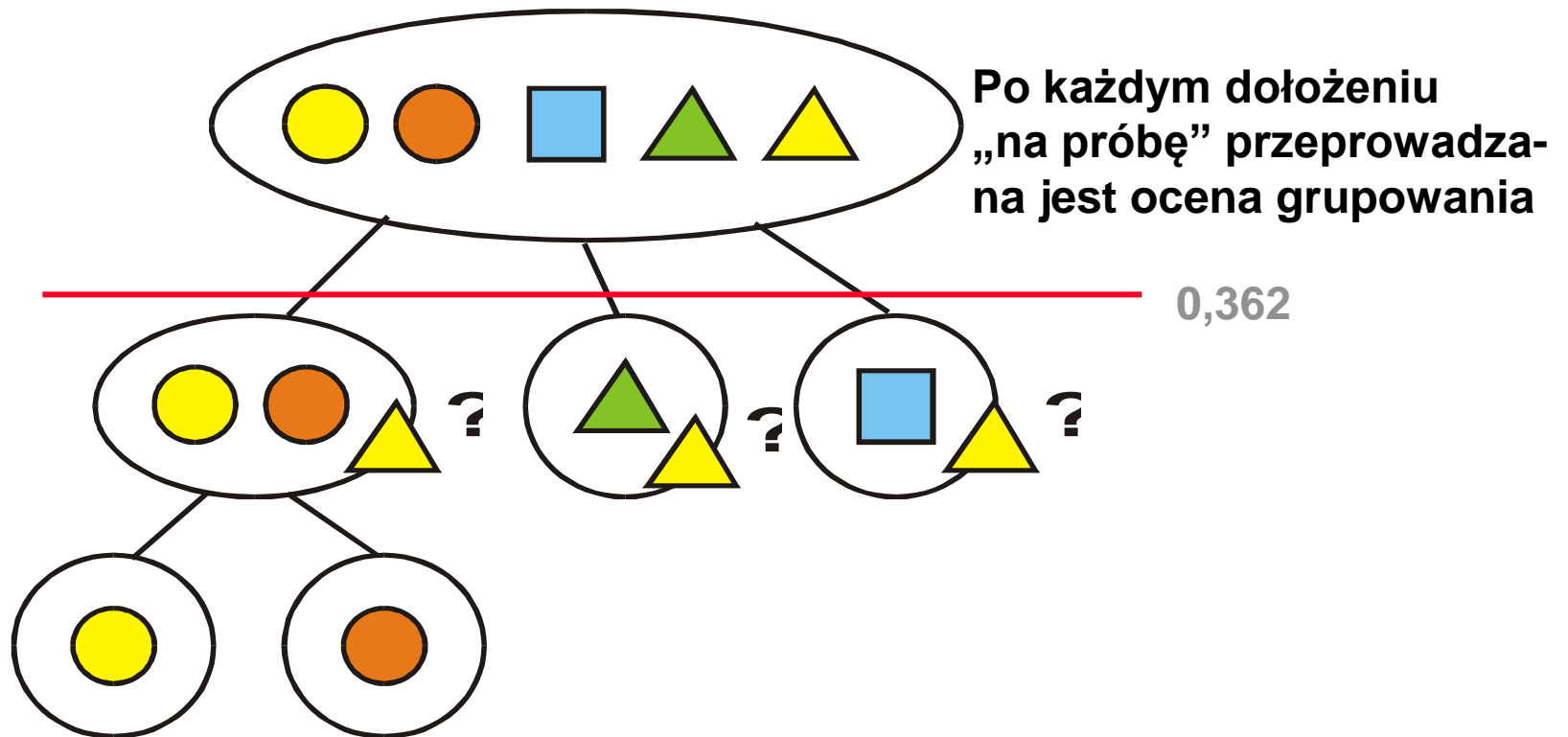


# COBWEB – operacje

- Na każdym z poziomów drzewa algorytm może wykonać jedną z operacji:
  - zaliczyć przykład do istniejącej podkategorii,
  - utworzyć osobną podkategorię dla przykładu,
  - dokonać podziału podkategorii i zaliczyć przykład do jednej z jej kategorii potomnych,
  - połączyć dwie podkategorie i zaliczyć przykład do podkategorii połączonej,
- O tym, którą operację wykonać, decyduje wartość funkcji oceny jakości grupowania.

# COBWEB – grupowanie

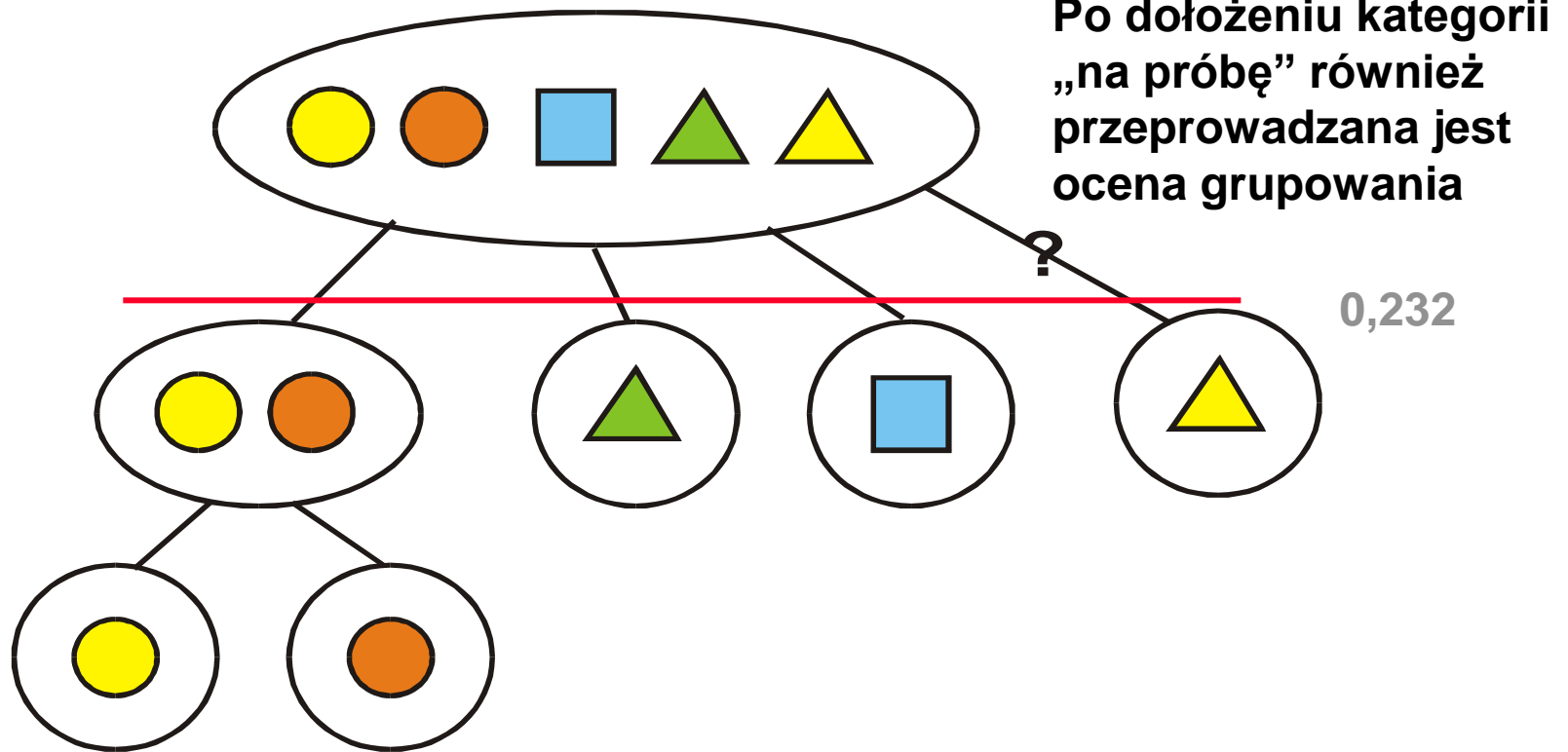
Algorytm dokłada „na próbę” przykład do każdej podkategorii.





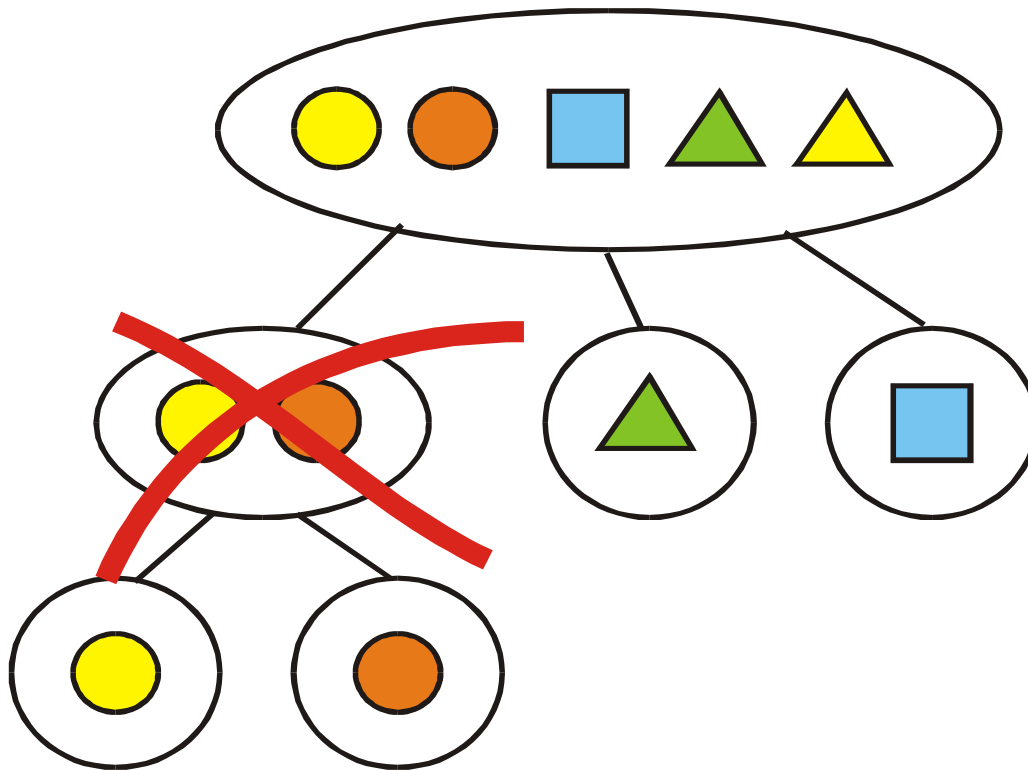
# COBWEB – grupowanie

Algorytm próbuje też stworzyć dla przykładu osobną podkategorię.



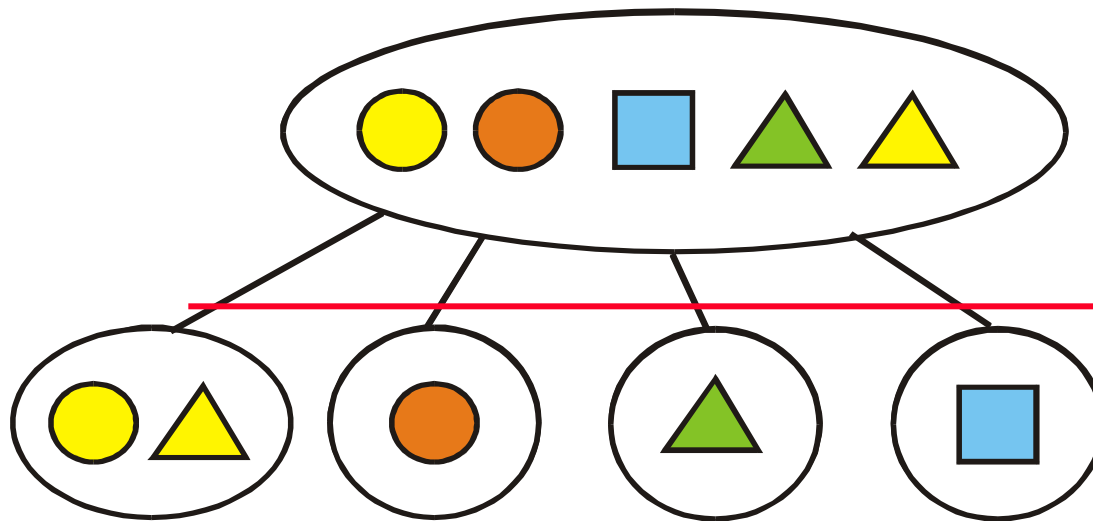
# COBWEB – podział

COBWEB próbuje też usunąć jedną z podkategorii i dodać przykład do jednej z jej kategorii potomnych.



# COBWEB – podział

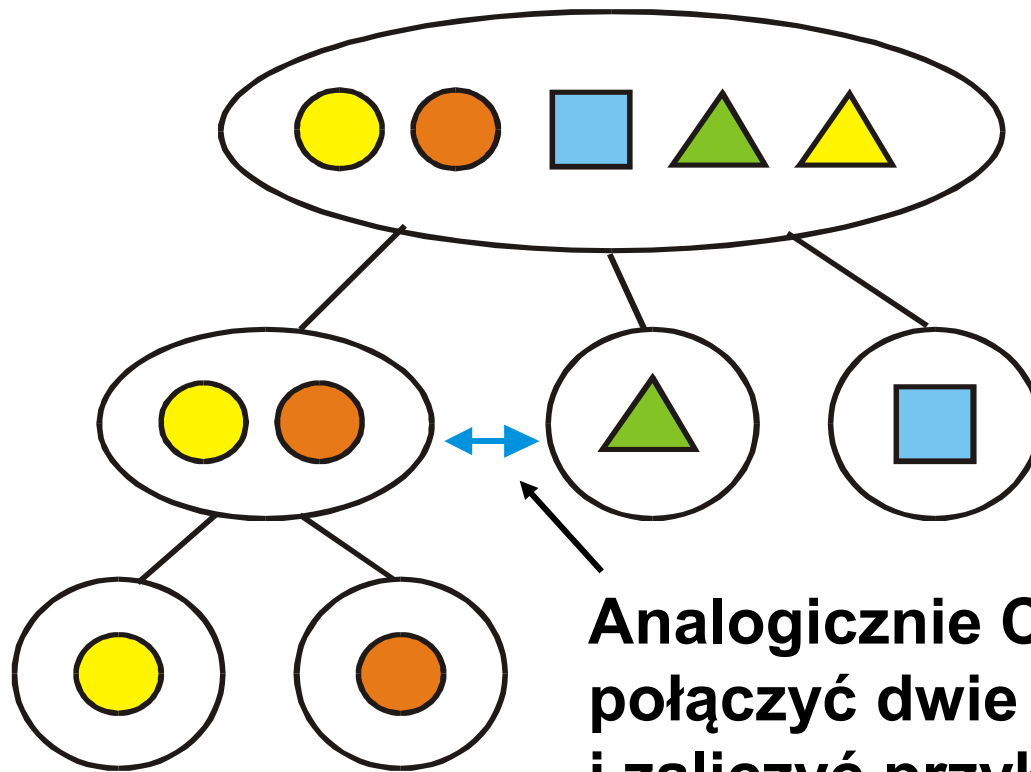
COBWEB próbuje też usunąć jedną z podkategorii i dodać przykład do jednej z jej kategorii potomnych.



Oczywiście po tych posunięciach (nadal wykonywanych „na próbę”) obliczamy ocenę grupowania

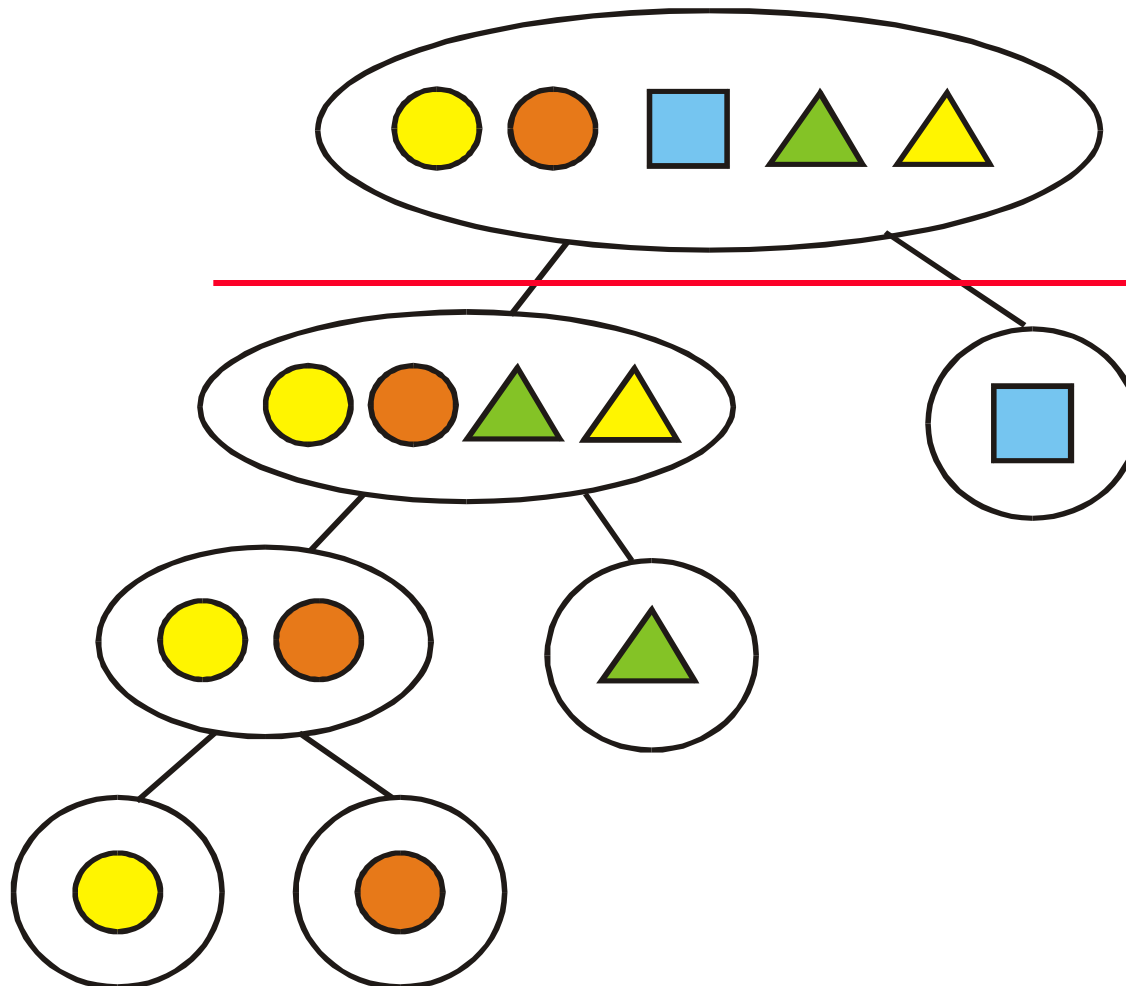
0,326

# COBWEB – scalenie



**Analogicznie COBWEB może połączyć dwie kategorie i zaliczyć przykład do kategorii połączonej**

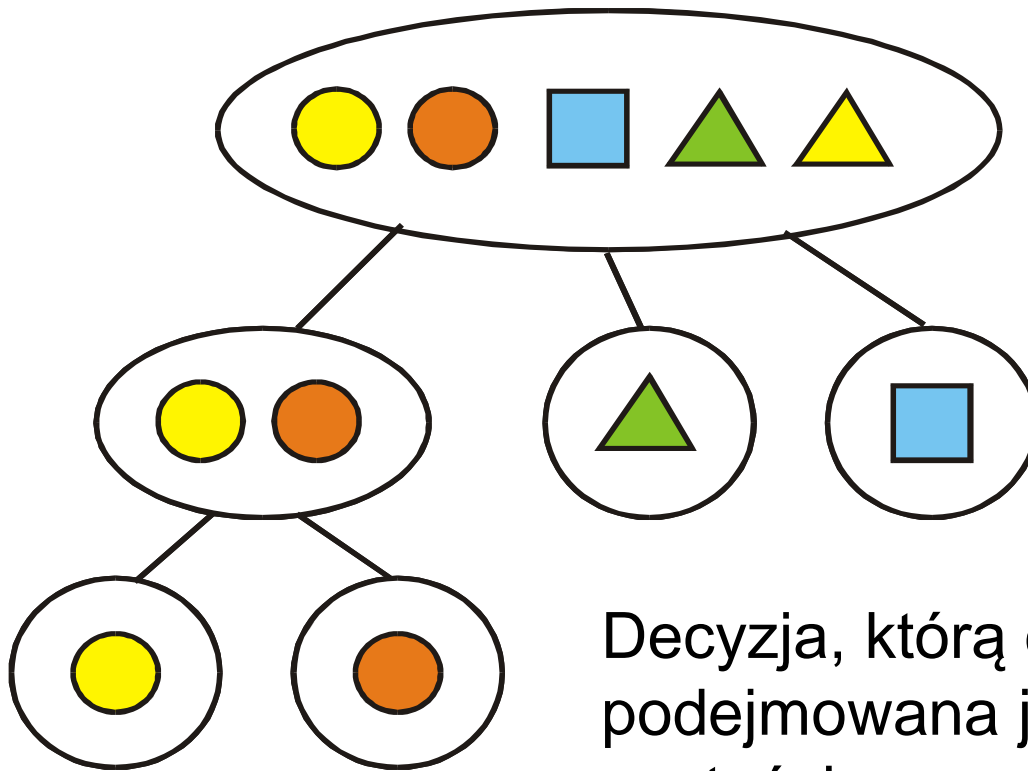
# COBWEB – scalenie



0,428

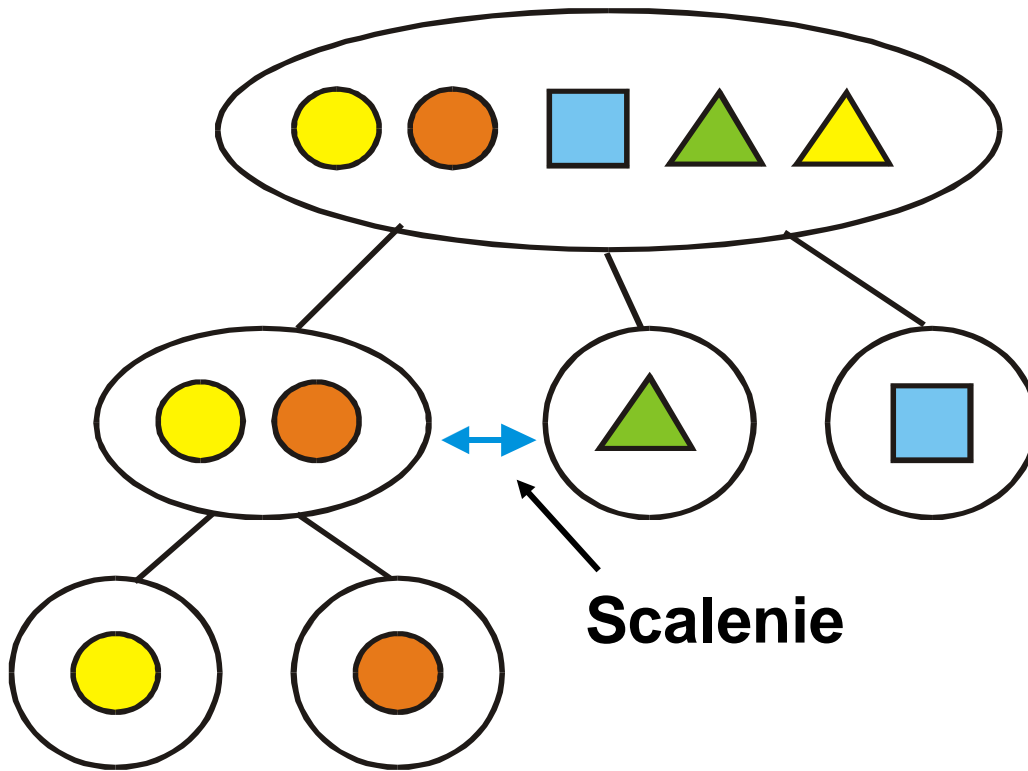
Po tym, ostatnim już,  
posunięciu „na  
próbę” obliczamy  
ocenę grupowania

# COBWEB – przykład

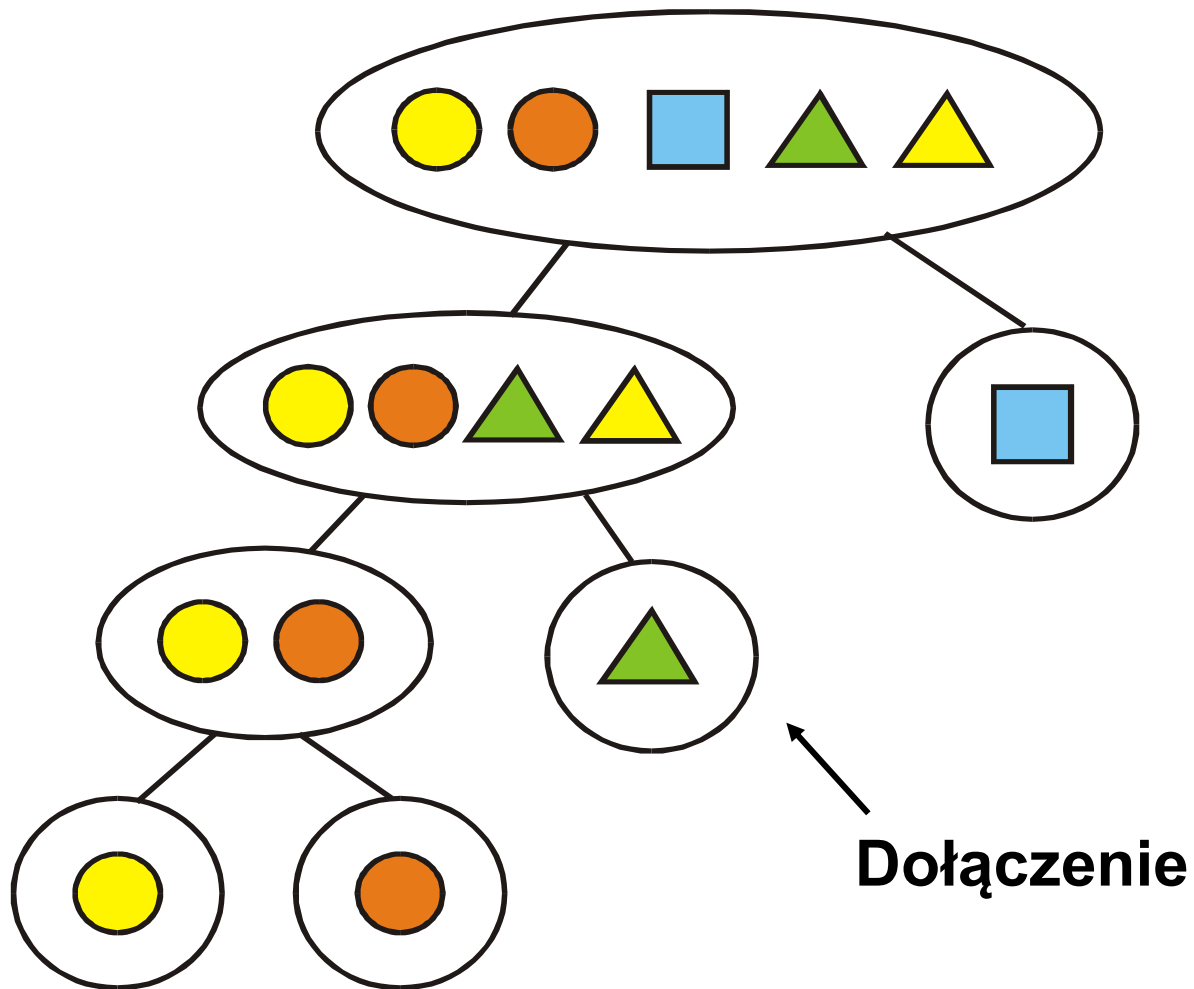


Decyzja, którą operację wykonać podejmowana jest na podstawie wartości oceny grupowania na każdym poziomie drzewa.

# COBWEB – przykład

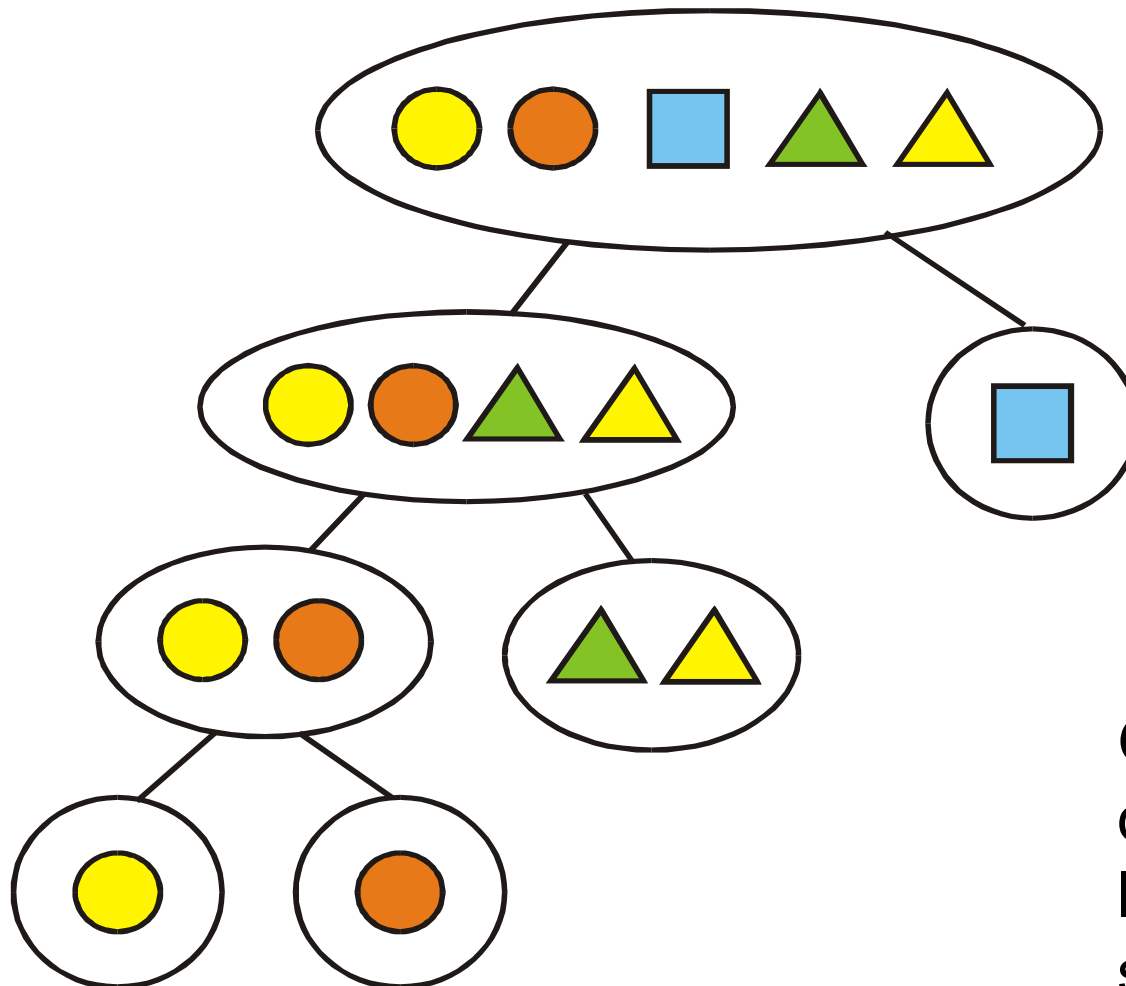


# COBWEB – przykład



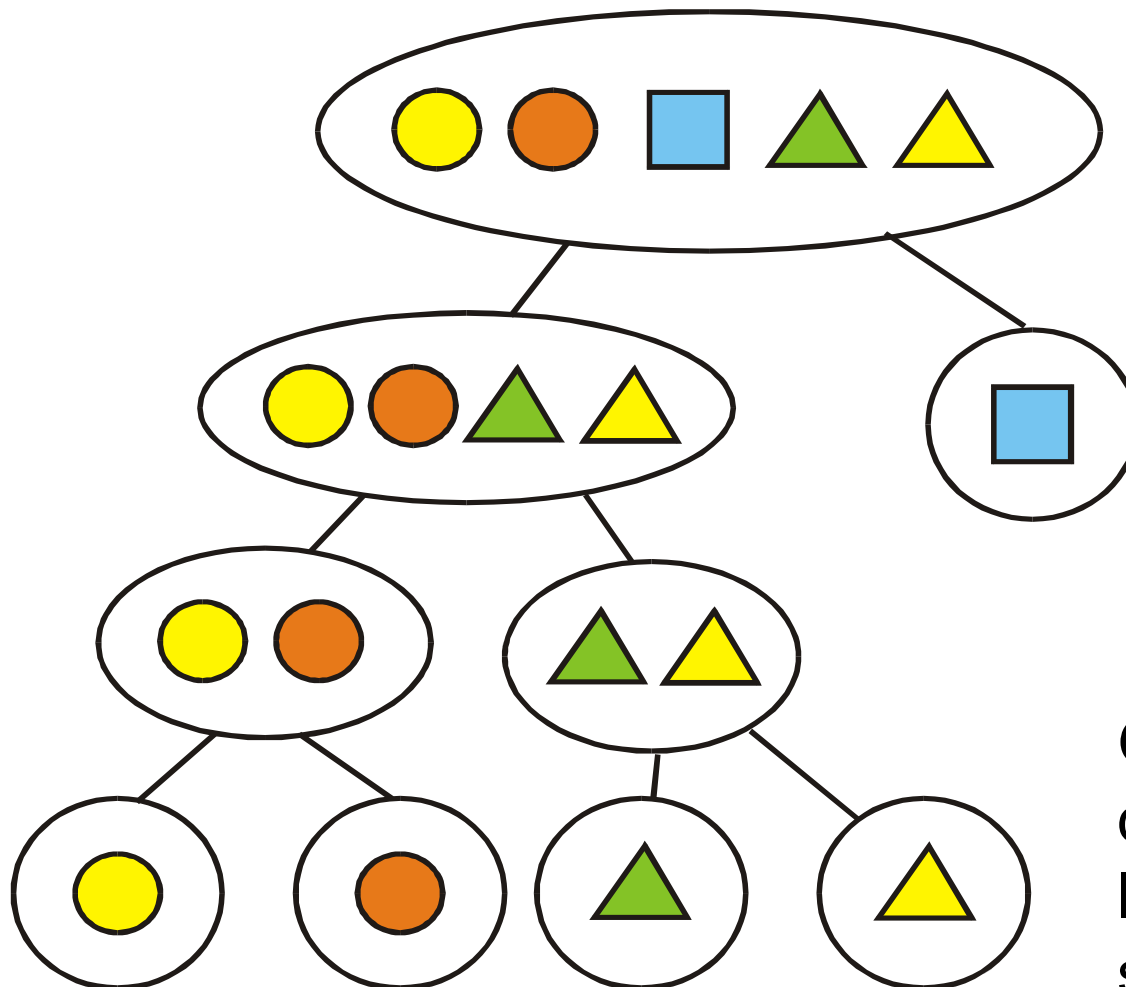


# COBWEB – przykład



Gdy przykład dołączany jest do liścia, generowane są dwa nowe liście

# COBWEB – przykład



Gdy przykład dołączany jest do liścia, generowane są dwa nowe liście

# COBWEB – ocena grupowania

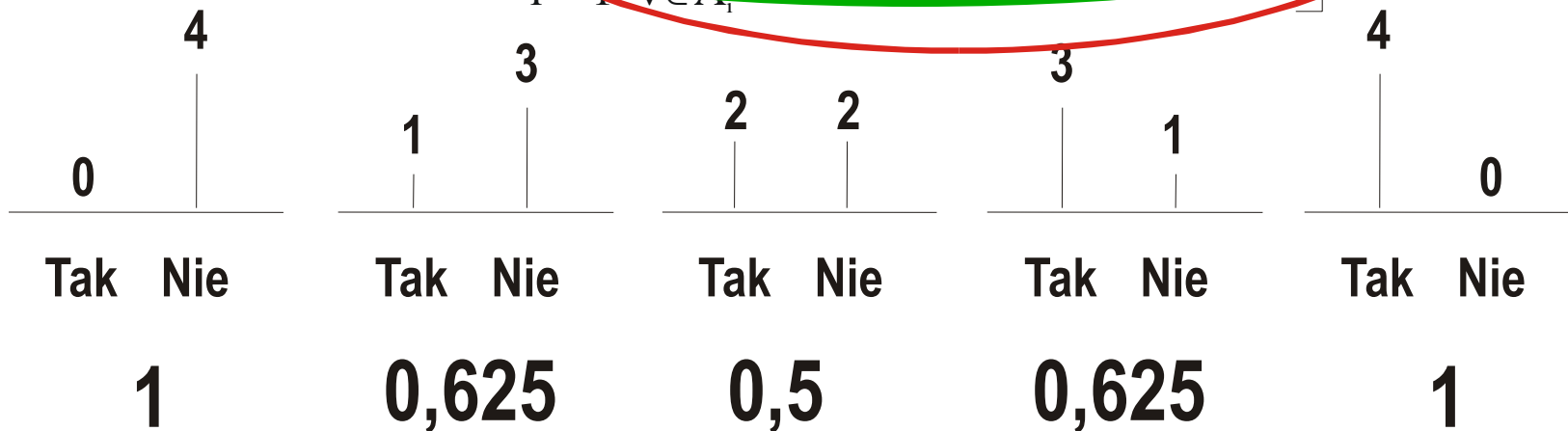
$$f(h) = \frac{1}{|C_h|} \sum_{d \in C_h} \Pr_{x \in \Omega}(h(x)=d).$$

związłość

$$\left[ \sum_{i=1}^n \sum_{v \in A_i} \Pr_{x \in \Omega}(a_i(x)=v | h(x)=d) \right]^2 -$$

$$= \sum_{i=1}^n \sum_{v \in A_i} \Pr_{x \in \Omega}(a_i(x)=v)^2$$

precyzja

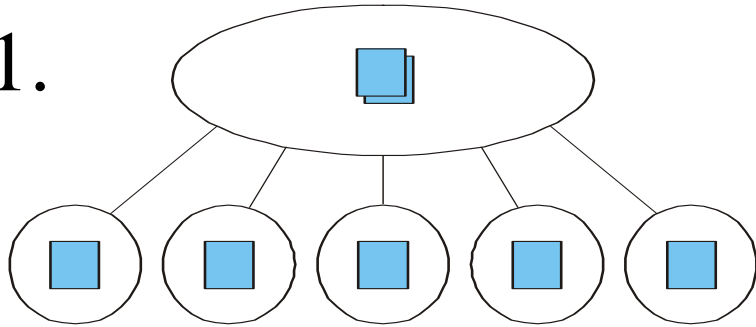


# COBWEB – zalety

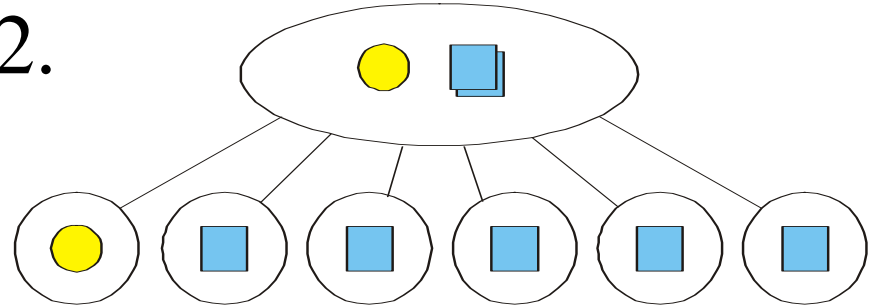
- COBWEB może być używany w trybie inkrementacyjnym,
- Produkowany jest hierarchiczny podział przykładów,
- Funkcja oceny jakości grupowania może być dostosowana do grupowania niestandardowych typów atrybutów lub nawet niestandardowych przykładów

# COBWEB – wada

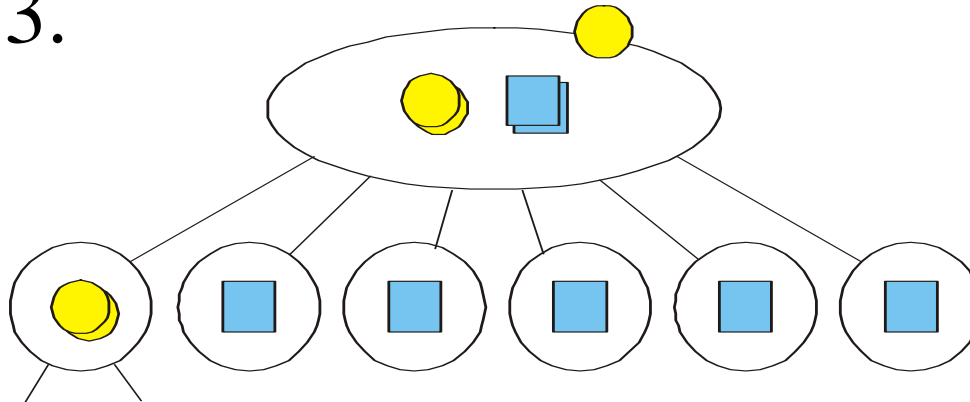
1.



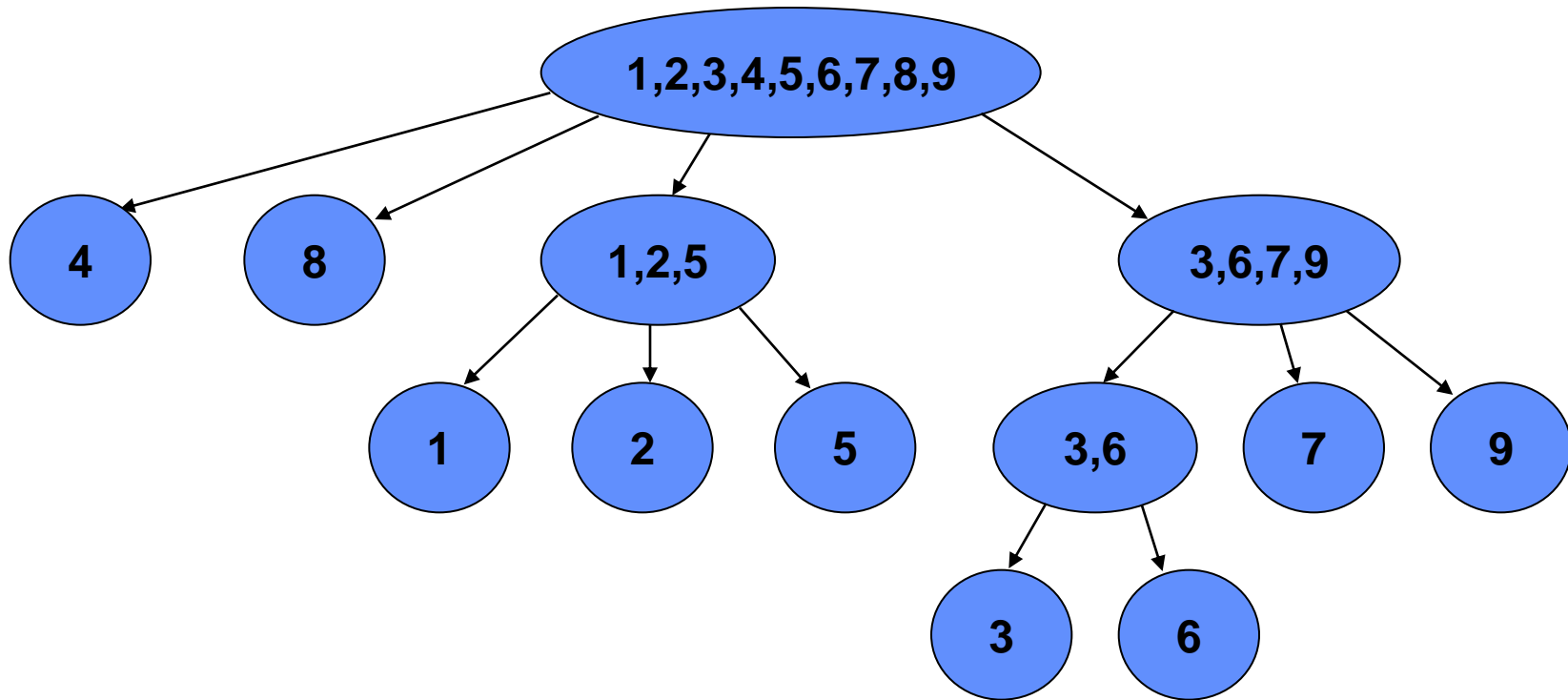
2.



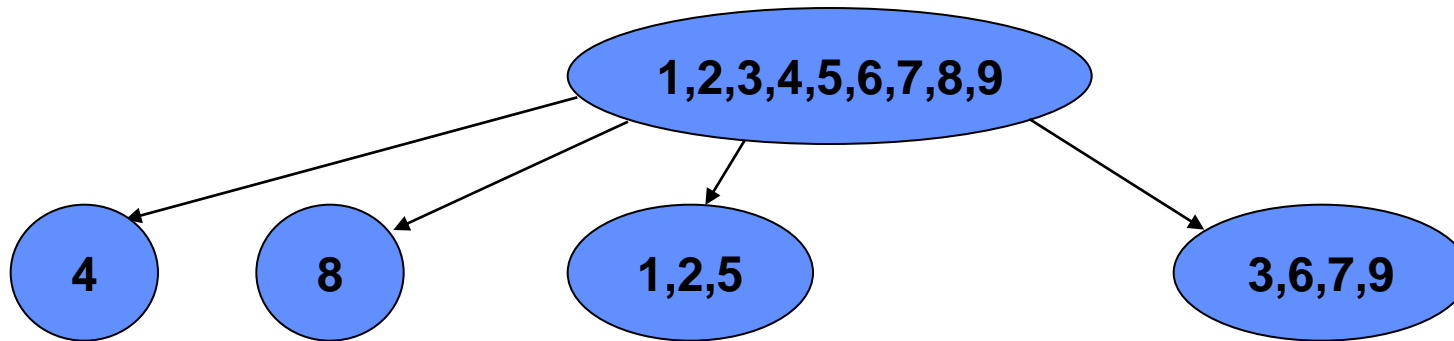
3.



# COBWEB – przykład „Z.K.”



# COBWEB – klasyfikacja



**W praktyce jako wynik grupowania traktuje się drugi poziom drzewa COBWEB licząc od korzenia**

**Takie zredukowane drzewo może być traktowane również jako klasyfikator**

# Grupowanie - przykład

S.C.	D.O.R.	Wiek	Wykształcenie	Sam.	Z.K.
S	800	32	wyższe	tak	tak
S	1200	35	średnie	tak	tak
S	700	26	podstawowe	nie	nie
M	600	45	wyższe	nie	tak
M	650	38	średnie	tak	tak
S	900	28	wyższe	nie	nie
S	1100	65	średnie	tak	nie
M	500	22	średnie	nie	nie
S	800	43	podstawowe	tak	nie



# COBWEB – podsumowanie

- COBWEB dołącza każdy przykład do drzewa kategorii, wykonując proste operacje,
- O tym, którą operację wykonać, decyduje funkcja oceny jakości grupowania według kryteriów precyzji i zwięzłości,
- Działa w trybie inkrementacyjnym, można dostosowywać sposób jego działania,
- Algorytm zachłanny, efekt jego działania mocno zależy od kolejności podawania przykładów na wejście

# Dziękujemy za uwagę

Zapraszamy na wykład:

**KLASTERYZACJA I SEGMENTACJA cz. 2**