

Metody Inżynierii Wiedzy

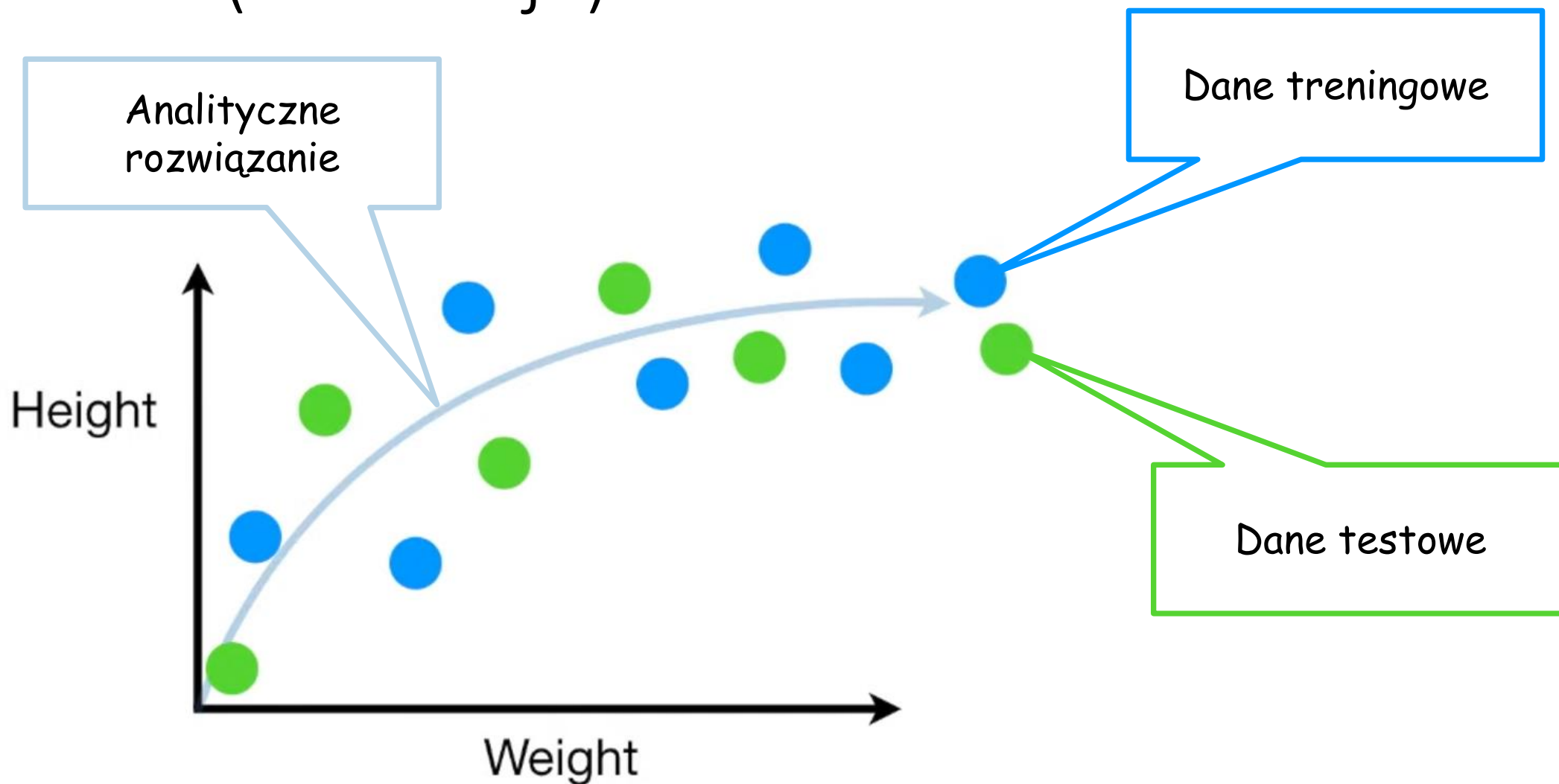
Łączenie różnych modeli w celu działania zespołowego

Dr inż. Michał Majewski

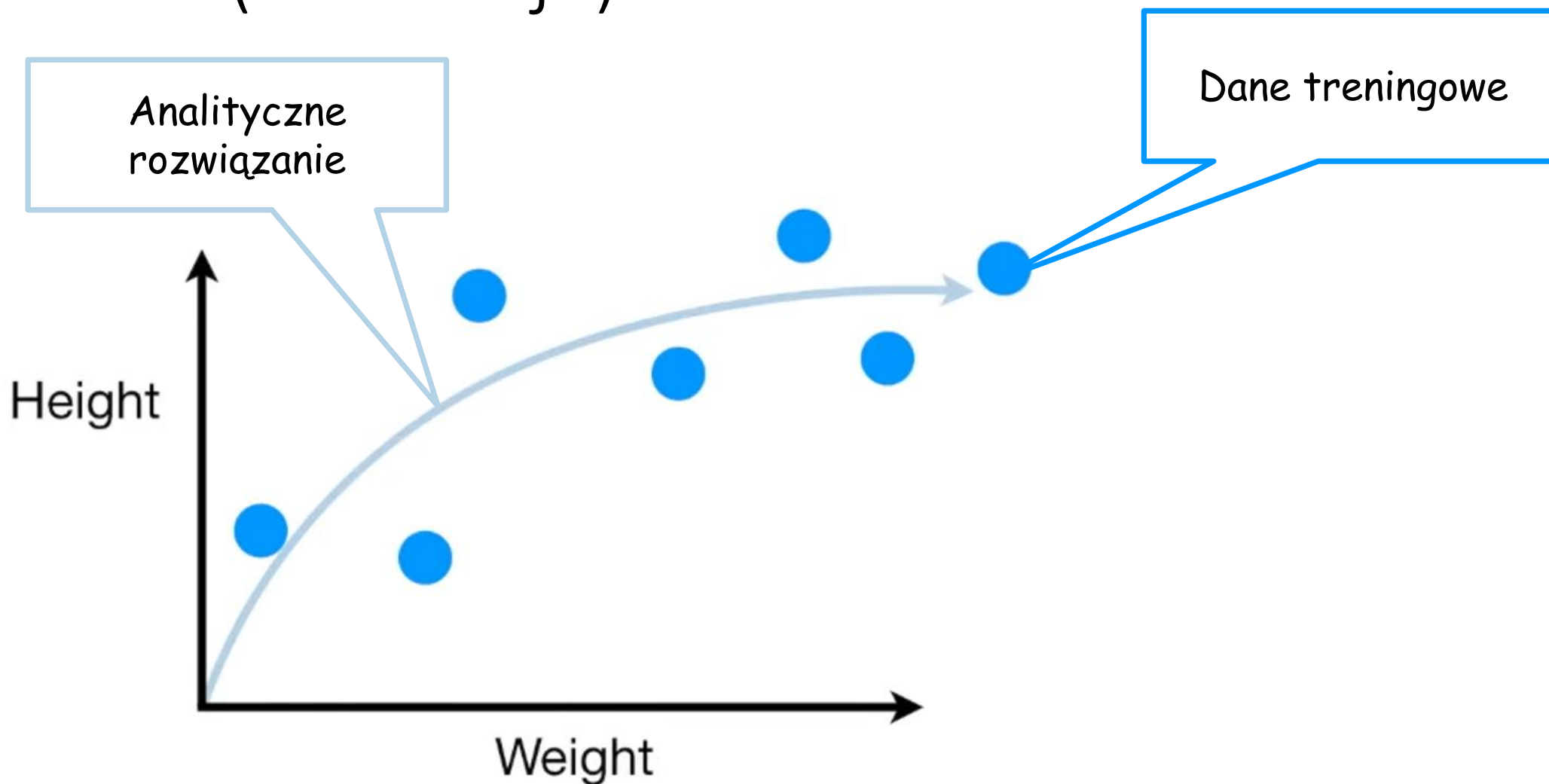
mmajew@pjawstk.edu.pl

materiały: *ftp(public) : //mmajew/MIW*

Bias (błąd, stronniczość) *Variance* (wariancja)

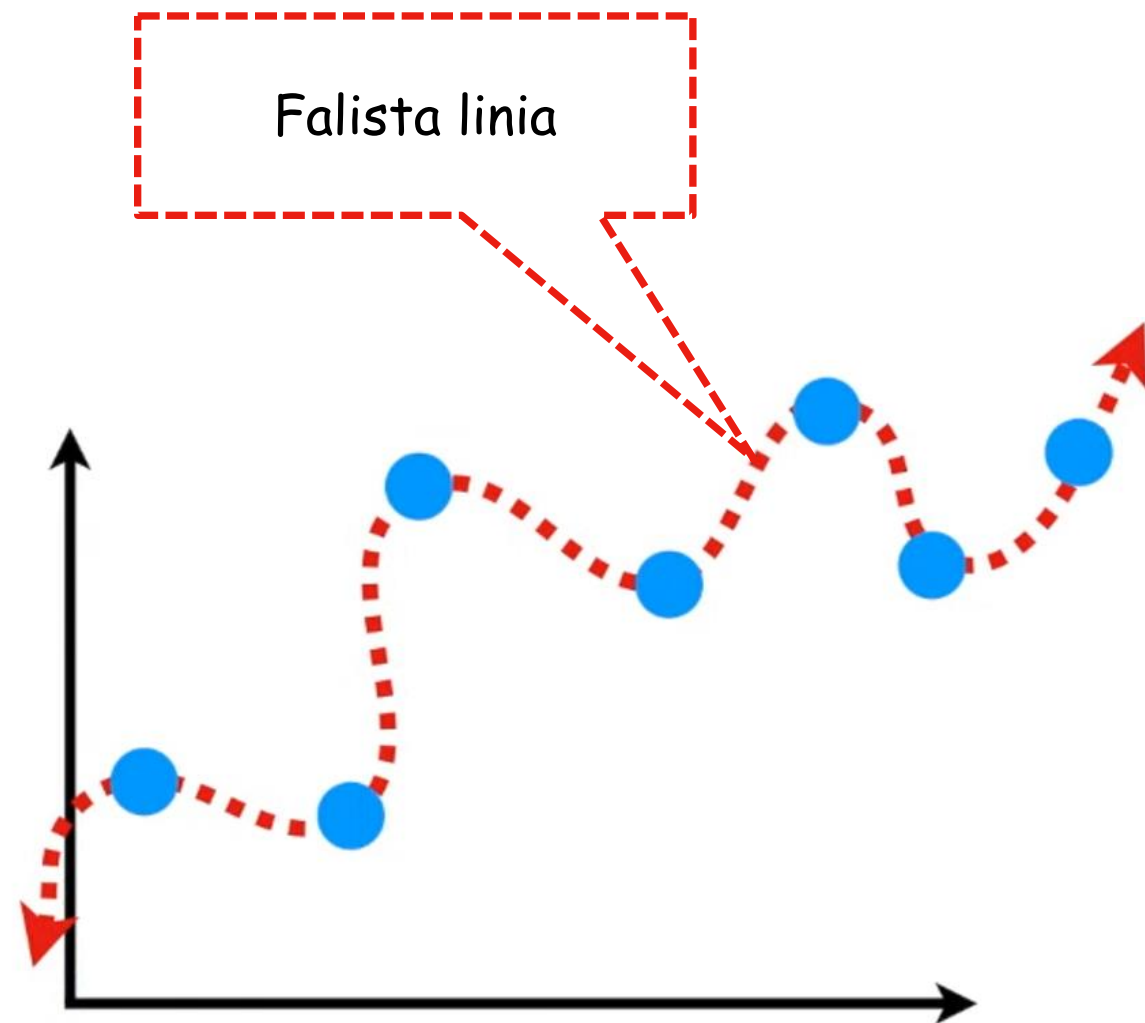
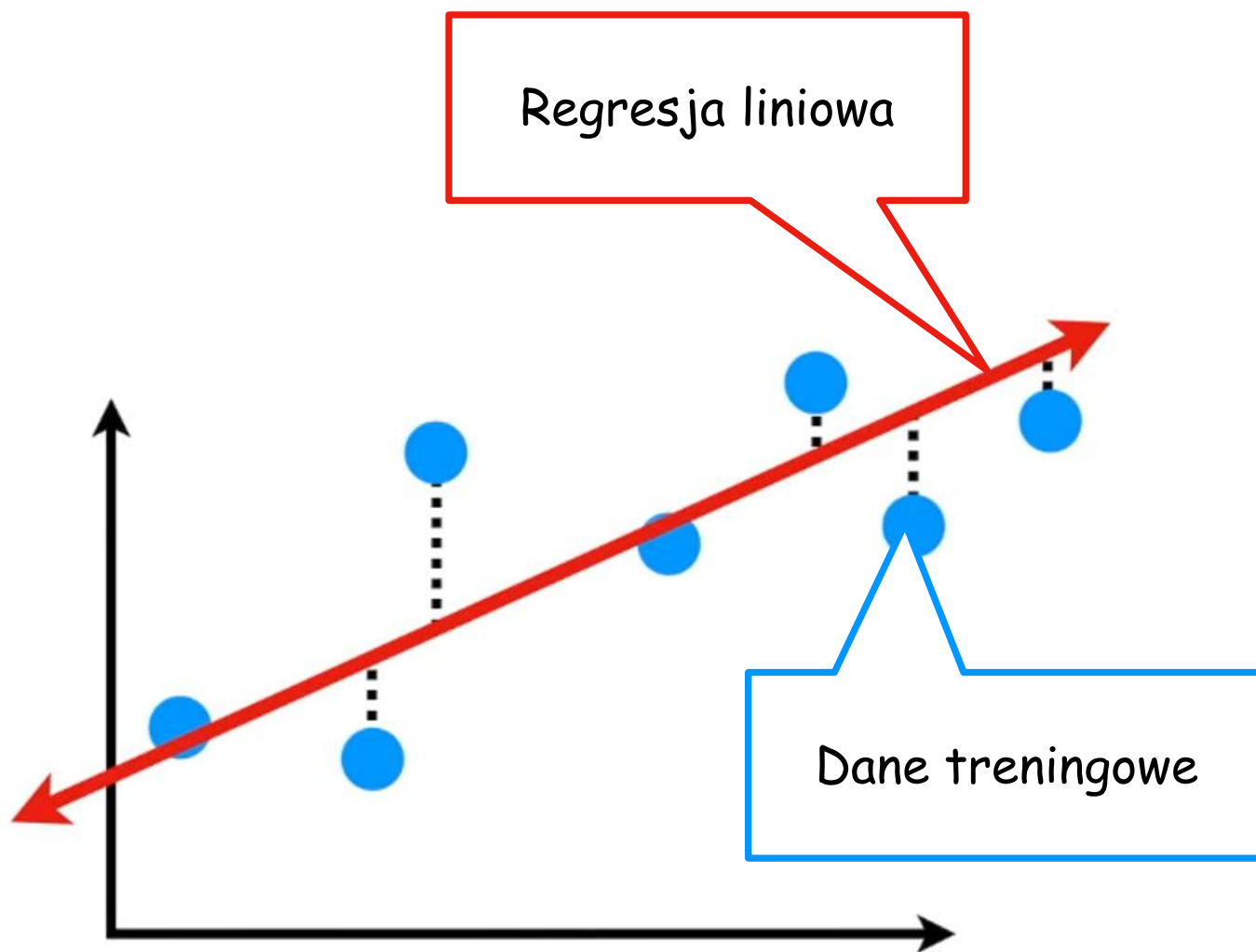


Bias (błąd, stronniczość) *Variance* (wariancja)



Bias (błąd, stronniczość) *Variance* (wariancja)

Bias - systematyczny błąd np. gdy model jest zbyt uproszczony



Bias (błąd, stronniczość) *Variance* (wariancja)

wariancja jak bardzo różnią się prognozy modelu, gdy jest uczony na różnych podzbiorach danych treningowych

Regresja liniowa

Falista linia

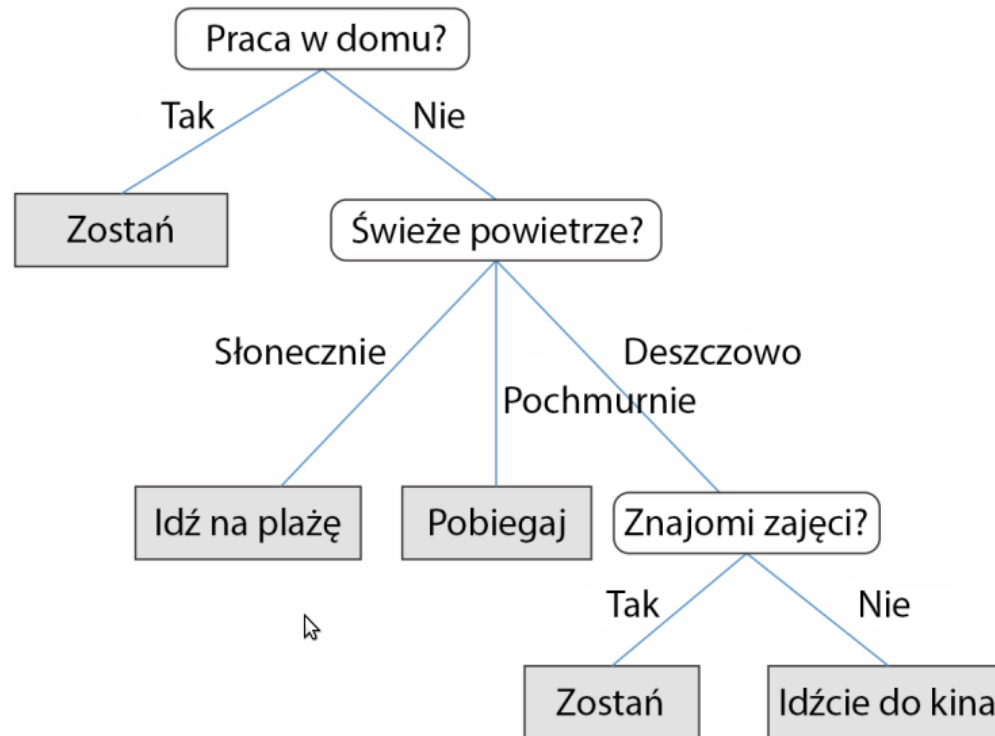
vs.

Dane testowe

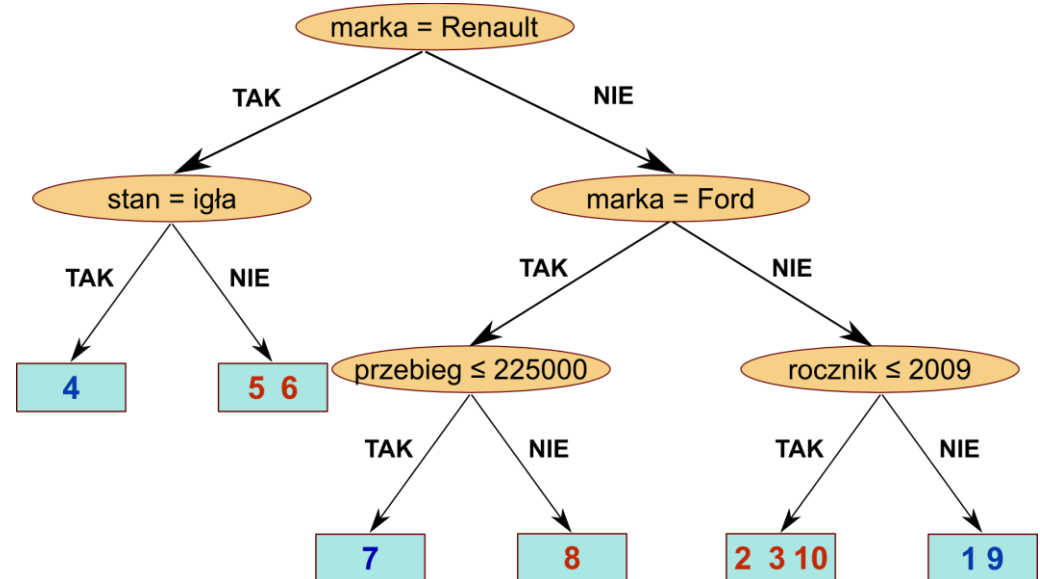
nadmierne dopasowanie
overfitting

Drzewo decyzyjne

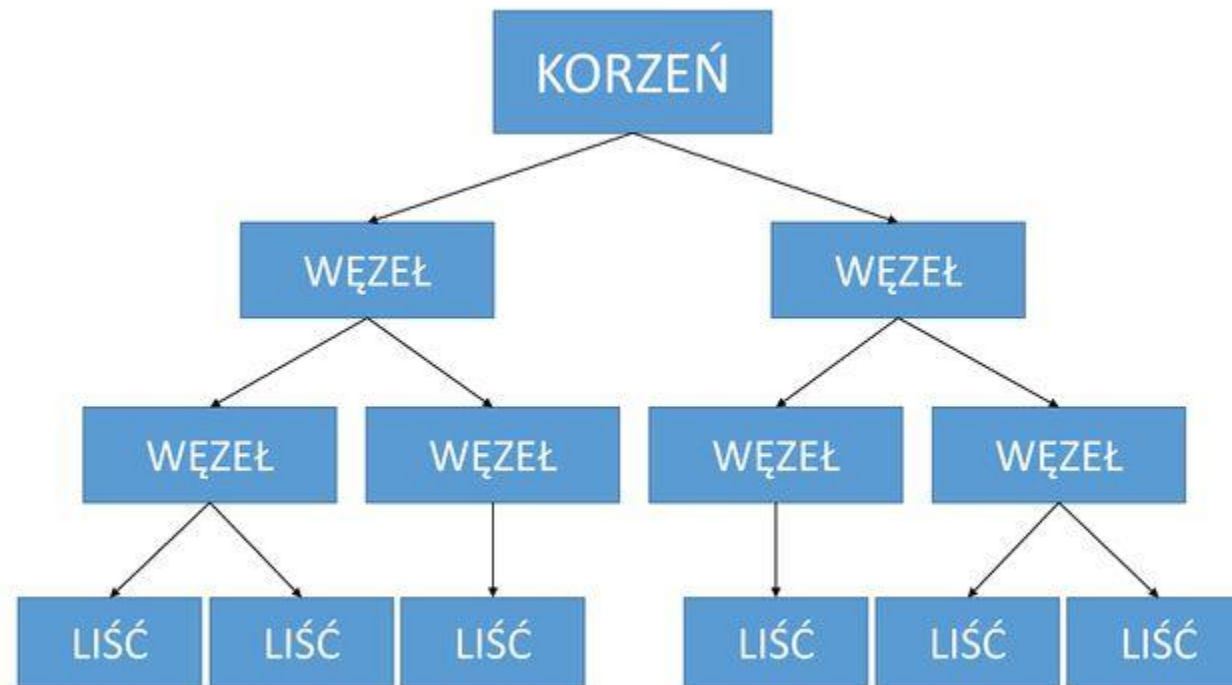
Drzewo klasyfikacji



Drzewo regresyjne



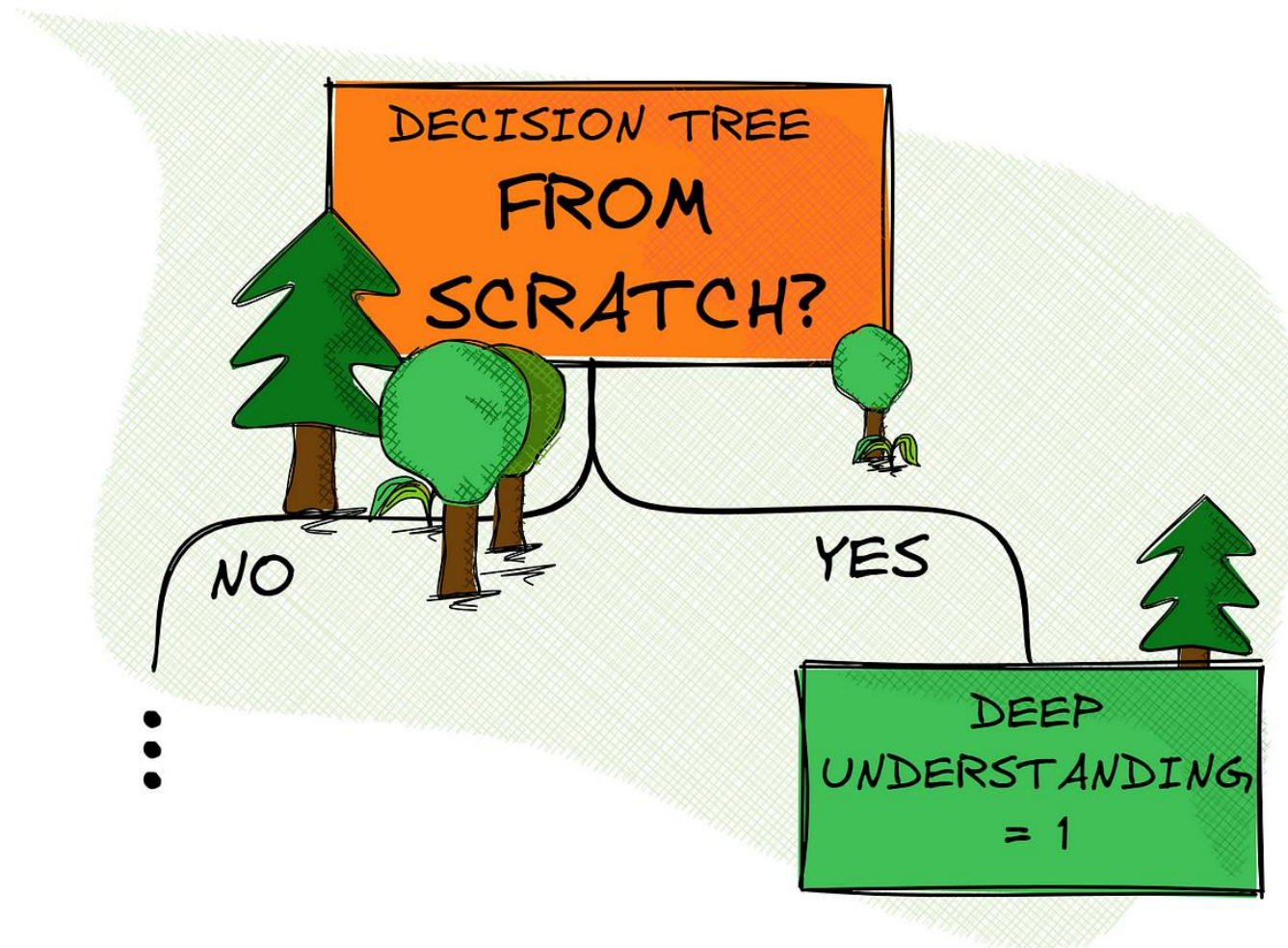
Drzewo decyzyjne



Drzewo decyzyjne – jak zbudować

Dane wejściowe

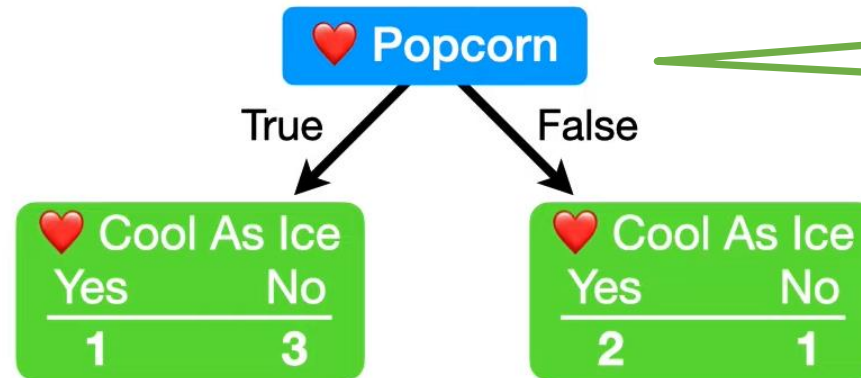
Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No



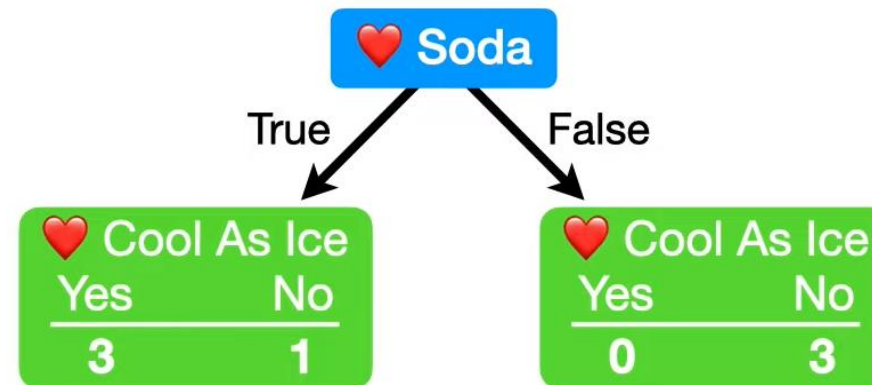
Drzewo decyzyjne – jak zbudować

Dane wejściowe

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No



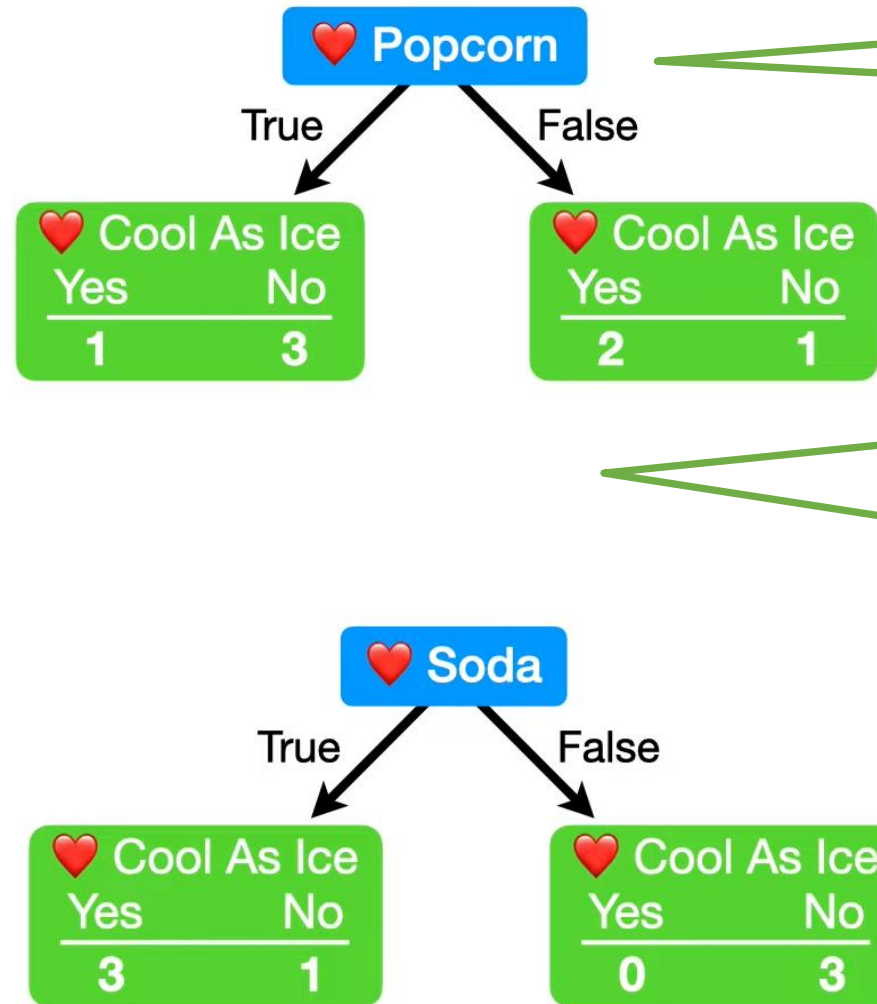
Szukamy węzła startowego (korzenia)



Drzewo decyzyjne – jak zbudować

Dane wejściowe

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No



Szukamy węzła startowego (korzenia)

Kryteria rozgałęzienia (miara zanieczyszczeń):
> Entropia
> **Wskaźnik Giniego**
> Błąd klasyfikacji

Porównanie wskaźników zanieczyszczeń – wykład 5

Entropia

Dla wszystkich niepustych klas $p(i|t) \neq 0$:

$$I_H(t) = - \sum_{i=1}^c p(i|t) \log_2 p(i|t)$$

- Wyrażenie $p(i|t)$ oznacza proporcję pomiędzy próbkami należącymi do klasy i w danym węźle t ,
- Entropia będzie wynosiła 0, jeśli wszystkie próbki w węźle będą należały do tej samej klasy,
- Maksymalną wartość osiągnie wtedy, gdy będziemy mieli do czynienia z jednorodnym rozkładem klas,
- Poprzez kryterium entropii próbujemy zmaksymalizować wzajemne informacje w drzewie.

Błąd klasyfikacji

Błąd klasyfikacji możemy określić jako:

$$I_E(t) = 1 - \max\{p(i|t)\}$$

- Jest to kryterium przydatne do przycinania,
- Nie jest zalecane do rozwijania drzewa, ponieważ wykazuje mniejszą czułość na zmiany w rozkładzie prawdopodobieństwa klas wewnątrz węzła.

Wskaźnik Giniego

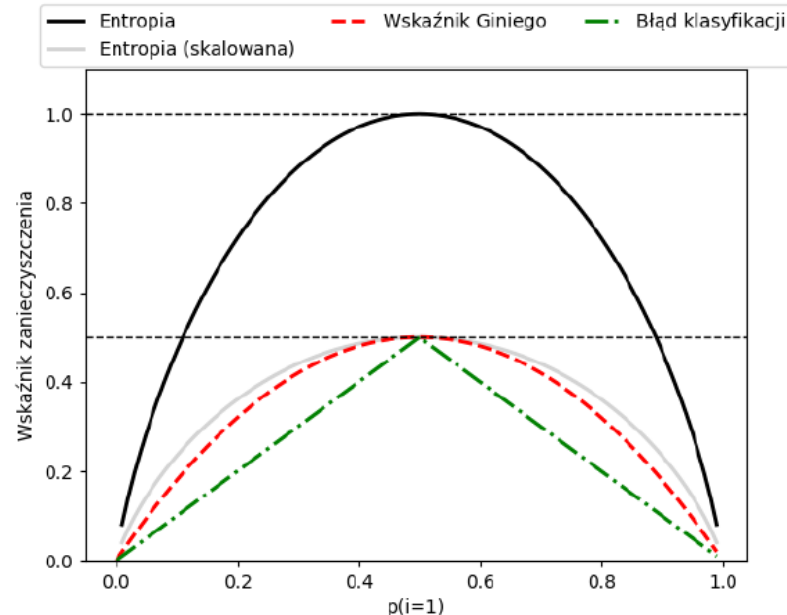
Wskaźnik Giniego możemy interpretować jako kryterium służące do minimalizowania prawdopodobieństwa nieprawidłowej klasyfikacji:

$$I_G(t) = \sum_{i=1}^c p(i|t)(1 - p(i|t)) = 1 - \sum_{i=1}^c p(i|t)^2$$

- Podobnie jak w przypadku entropii, wskaźnik Giniego uzyskuje największą wartość, gdy klasy są między sobą idealnie wymieszane; np. dla binarnej konfiguracji klas ($c = 2$):

$$I_G(t) = 1 - \sum_{i=1}^c p(i|t)^2 = 0,5$$

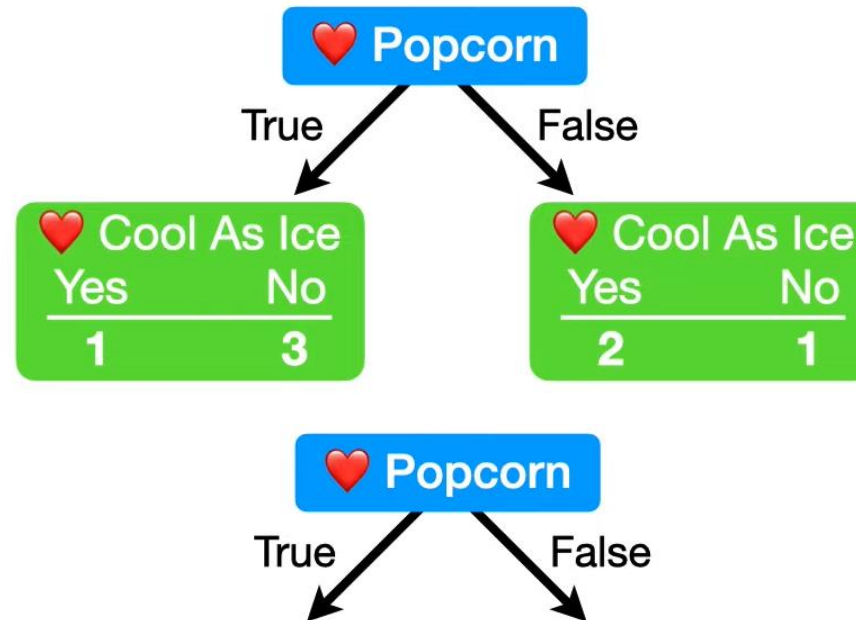
- Wskaźnik Giniego i entropia generują zazwyczaj podobne wyniki,
- Zamiast różnych kryteriów zanieczyszczeń, lepiej jest eksperymentować z różnymi wartościami granicy przycinania.



Drzewo decyzyjne – jak zbudować

Dane wejściowe

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No



Kryteria rozgałęzienia:
Wskaźnik Giniego I_G
Gini impurity

$$I_{G,pop,false} = 1 - \left(\frac{2}{2+1}\right)^2 - \left(\frac{1}{2+1}\right)^2 = 0.444$$

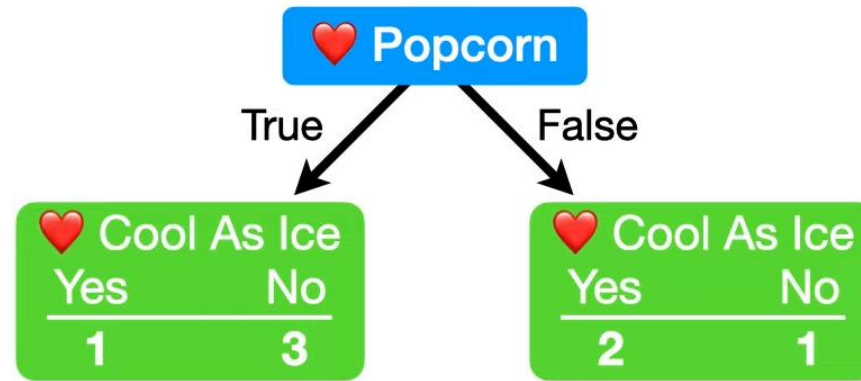
$$I_{G,pop,true} = 1 - \left(\frac{1}{1+3}\right)^2 - \left(\frac{3}{1+3}\right)^2 = 0.375$$

$$I_{G,popcorn} = \left(\frac{4}{4+3}\right) 0.375 + \left(\frac{3}{4+3}\right) 0.444 = \mathbf{0.405}$$

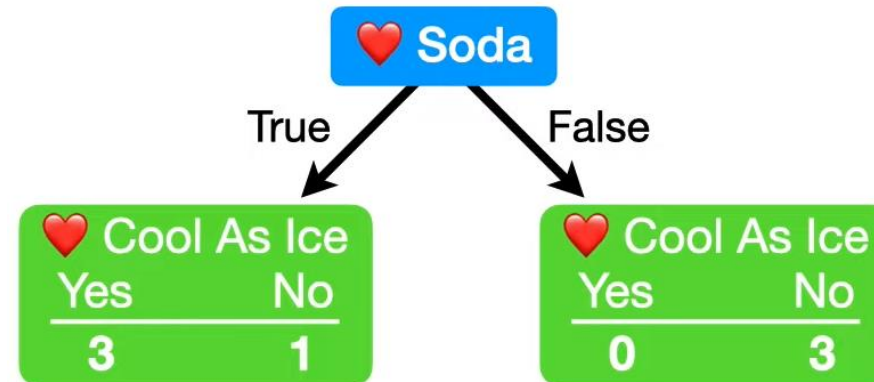
Drzewo decyzyjne – jak zbudować

Dane wejściowe

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No



$$I_{G,\text{popcorn}}=0.405$$



$$I_{G,\text{soda}}=0.214$$

Kryteria rozgałęzienia:
Wskaźnik Giniego I_G
Gini impurity

Drzewo decyzyjne – jak zbudować

Dane wejściowe

Sortujemy

Kryteria rozgałęzienia:
Wskaźnik Giniego IG
Gini impurity

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No

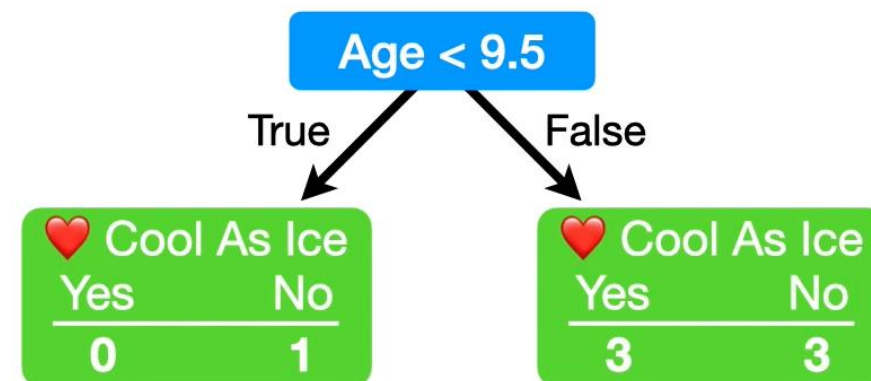
Określamy
granice w połowie

Drzewo decyzyjne – jak zbudować

Dane wejściowe

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No

Kryteria rozgałęzienia:
Wskaźnik Giniego I_G
Gini impurity



$$I_{G,9.5,true} = 1 - \left(\frac{0}{1}\right)^2 - \left(\frac{1}{1}\right)^2 = 0$$

$$I_{G,9.5,false} = 1 - \left(\frac{3}{3+3}\right)^2 - \left(\frac{3}{3+3}\right)^2 = 0.5$$

$$I_{G,9.5} = \left(\frac{1}{1+6}\right) 0 + \left(\frac{6}{1+6}\right) 0.5 = \mathbf{0.429}$$

Drzewo decyzyjne – jak zbudować

Dane wejściowe

Obliczamy IG dla każdej granicy

Kryteria rozgałęzienia:
Wskaźnik Giniego IG
Gini impurity

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No

0.429
0.343
0.476
0.476
0.343
0.429

Age < 9.5

True

False

♥ Cool As Ice	
Yes	No
0	1

♥ Cool As Ice	
Yes	No
3	3

$$I_{G,9.5,true} = 1 - \left(\frac{0}{1}\right)^2 - \left(\frac{1}{1}\right)^2 = 0$$

$$I_{G,9.5,false} = 1 - \left(\frac{3}{3+3}\right)^2 - \left(\frac{3}{3+3}\right)^2 = 0.5$$

$$I_{G,9.5} = \left(\frac{1}{1+6}\right) 0 + \left(\frac{6}{1+6}\right) 0.5 = \mathbf{0.429}$$

Drzewo decyzyjne – jak zbudować

Dane wejściowe

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No

Szukamy minimum

0.429

0.343

0.476

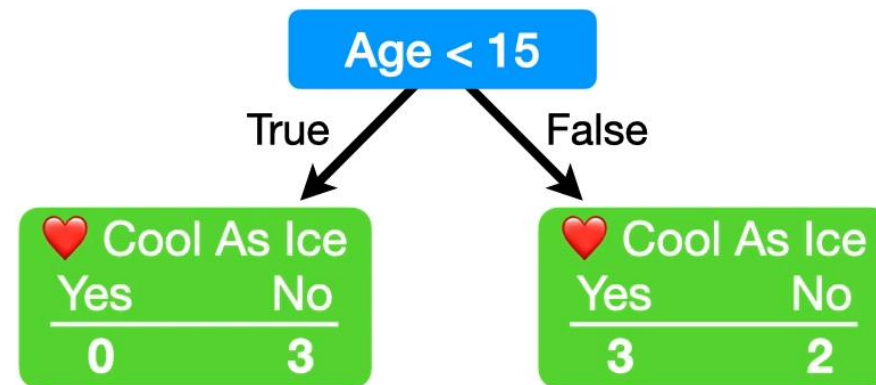
0.476

0.343

0.429

Kryteria rozgałęzienia:
Wskaźnik Giniego IG
Gini impurity

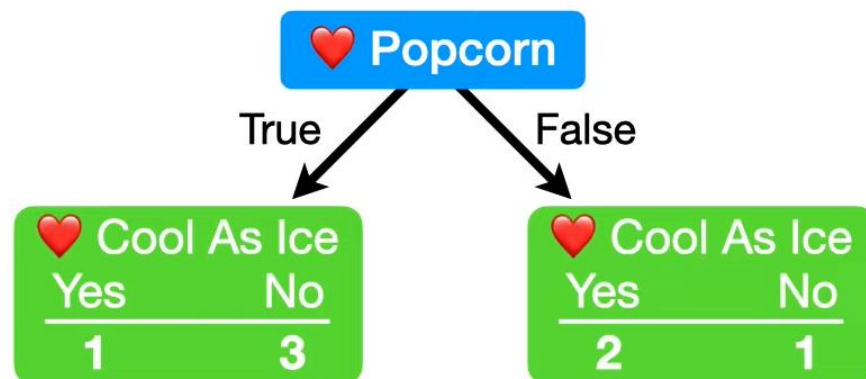
W tym przykładzie
wybieramy 15 lub 44



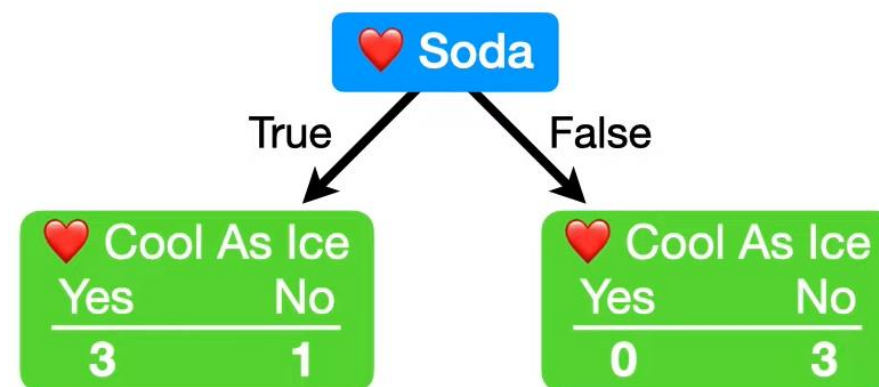
Drzewo decyzyjne – jak zbudować

Dane wejściowe

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No

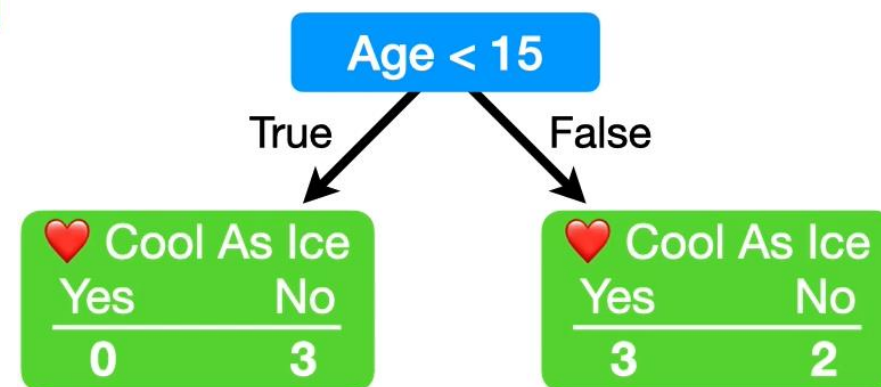


$$I_{G, \text{popcorn}} = 0.405$$



$$I_{G, \text{soda}} = 0.214$$

Kryteria rozgałęzienia:
Wskaźnik Giniego I_G
Gini impurity

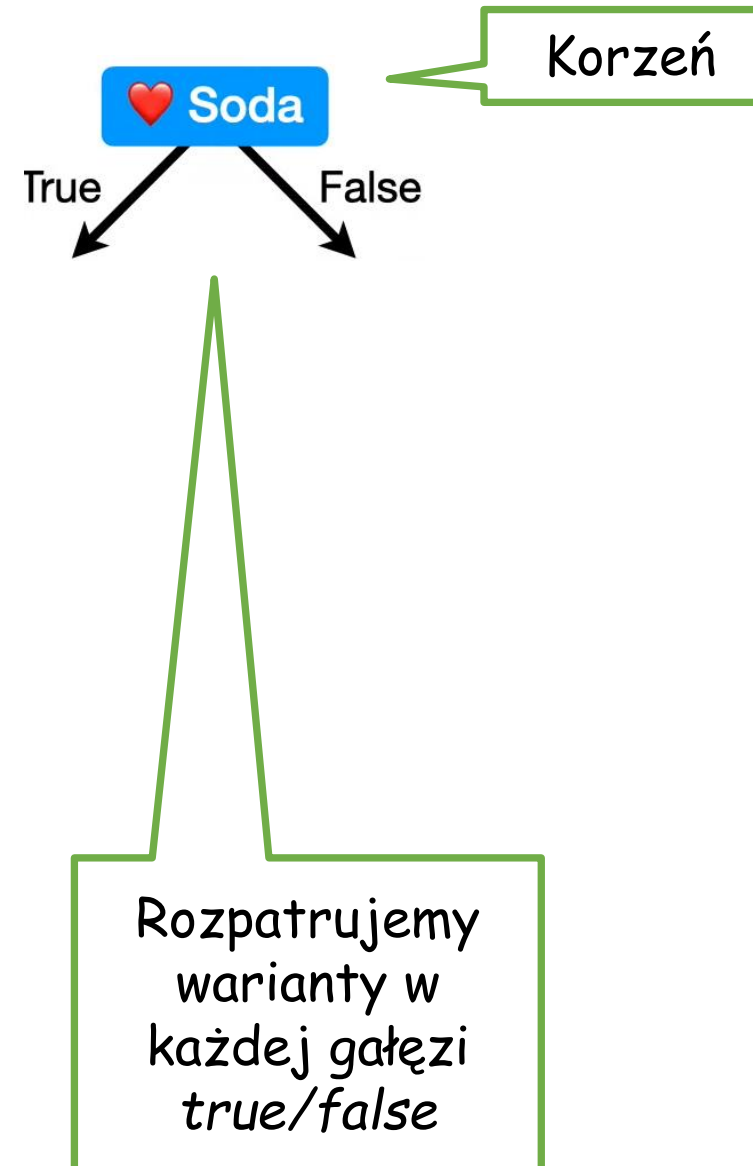


$$I_{G, 15} = 0.343$$

Drzewo decyzyjne – jak zbudować

Dane wejściowe

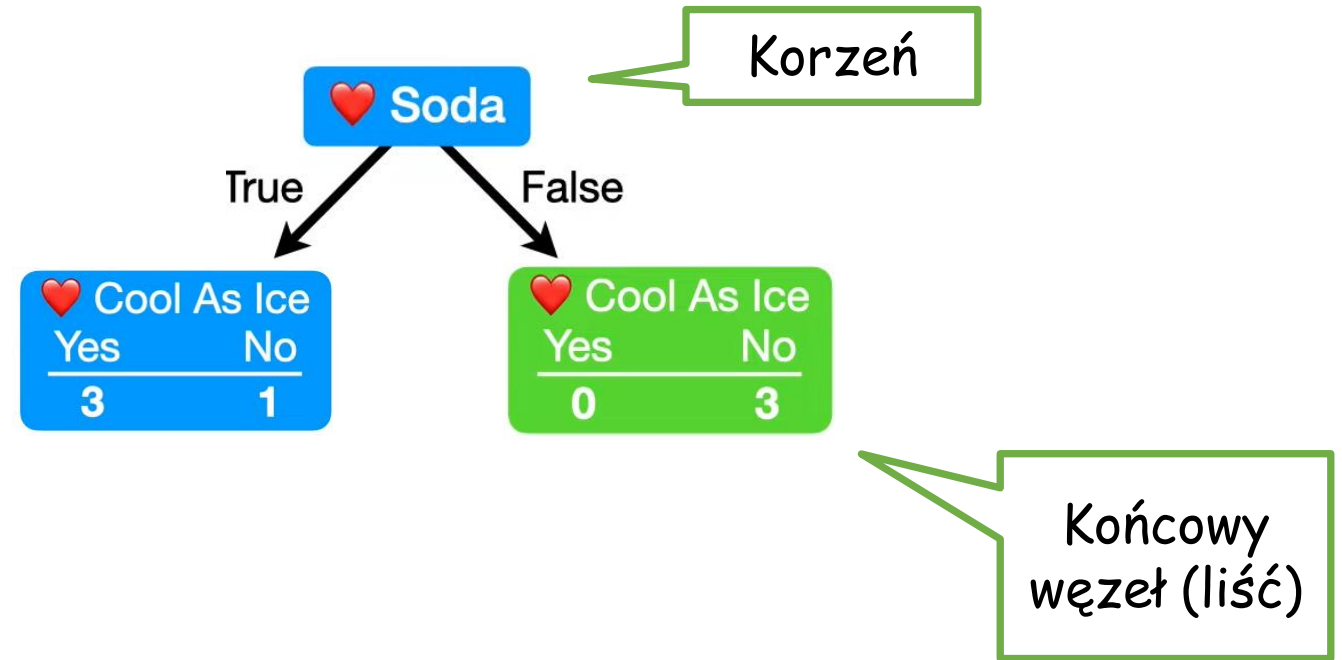
Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No



Drzewo decyzyjne – jak zbudować

Dane wejściowe

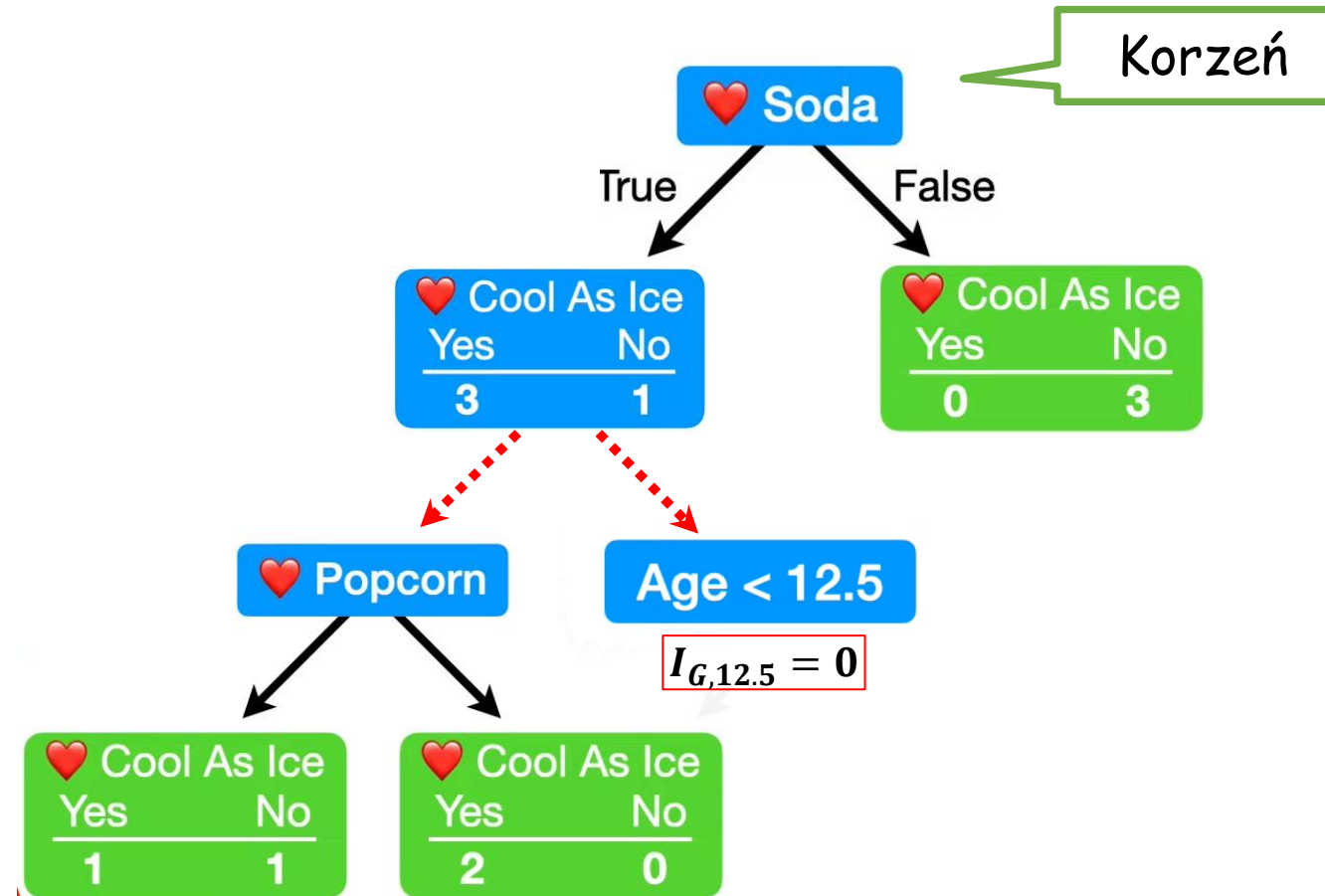
Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No



Drzewo decyzyjne – jak zbudować

Dane wejściowe

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No



$$I_{G,pop,true} = 1 - \left(\frac{1}{1+1}\right)^2 - \left(\frac{1}{1+1}\right)^2 = 0.5$$

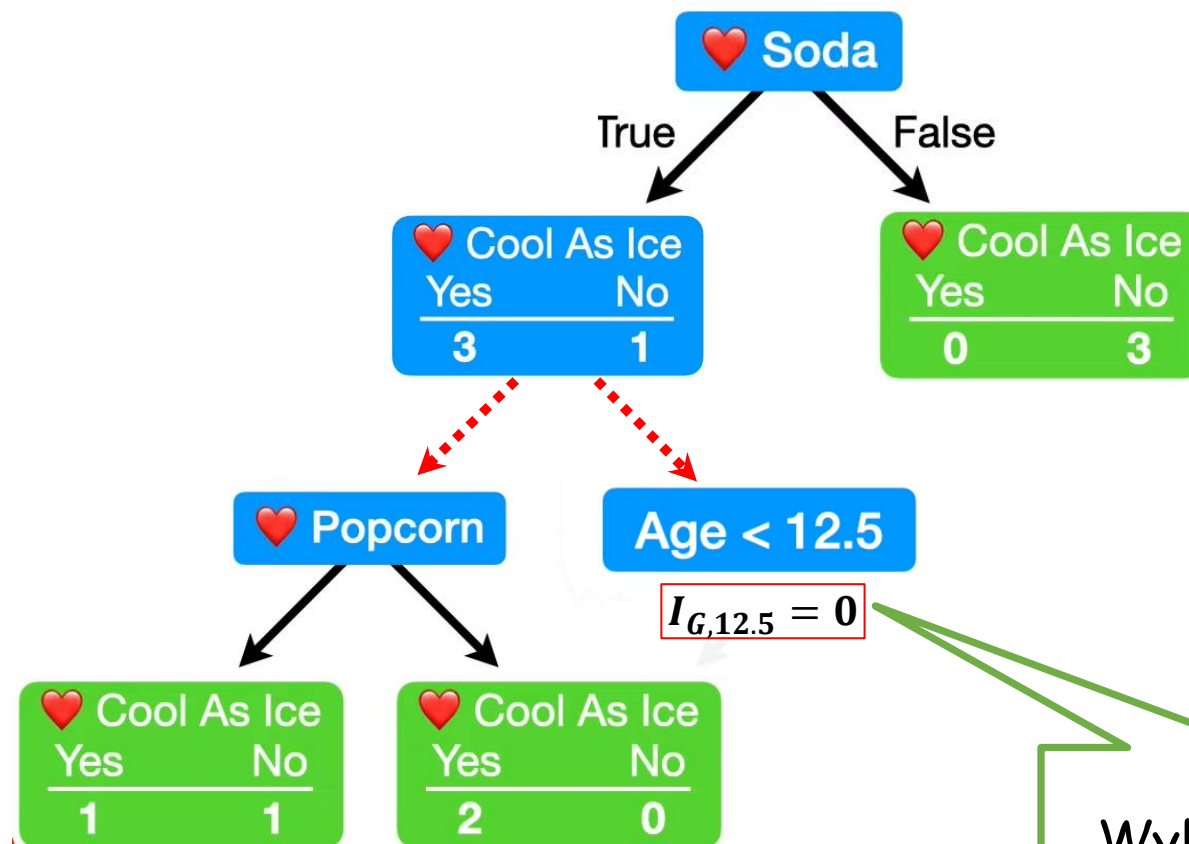
$$I_{G,pop,false} = 1 - \left(\frac{2}{0+2}\right)^2 - \left(\frac{0}{0+2}\right)^2 = 0$$

$$I_{G,popcorn} = \left(\frac{2}{2+2}\right) 0.0 + \left(\frac{2}{2+2}\right) 0.5 = 0.25$$

Drzewo decyzyjne – jak zbudować

Dane wejściowe

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No



$$I_{G,pop,true} = 1 - \left(\frac{1}{1+1}\right)^2 - \left(\frac{1}{1+1}\right)^2 = 0.5$$

$$I_{G,pop,false} = 1 - \left(\frac{2}{0+2}\right)^2 - \left(\frac{0}{0+2}\right)^2 = 0$$

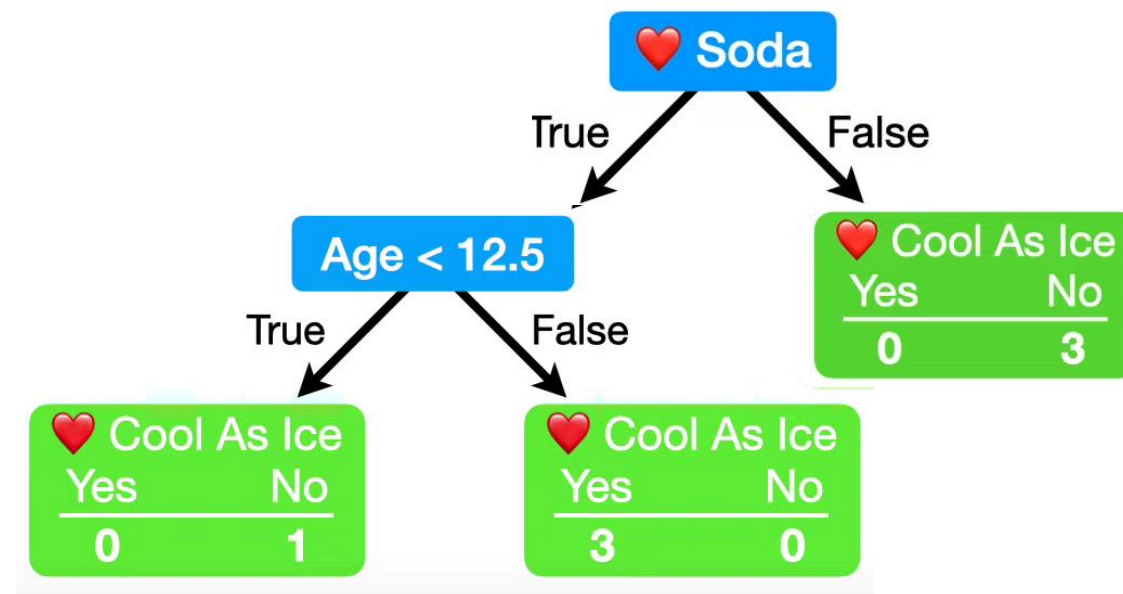
$$I_{G,popcorn} = \left(\frac{2}{2+2}\right) 0.0 + \left(\frac{2}{2+2}\right) 0.5 = 0.25$$

Wybieramy wiek 12.5 jako kolejne rozgałęzienie

Drzewo decyzyjne – jak zbudować

Dane wejściowe

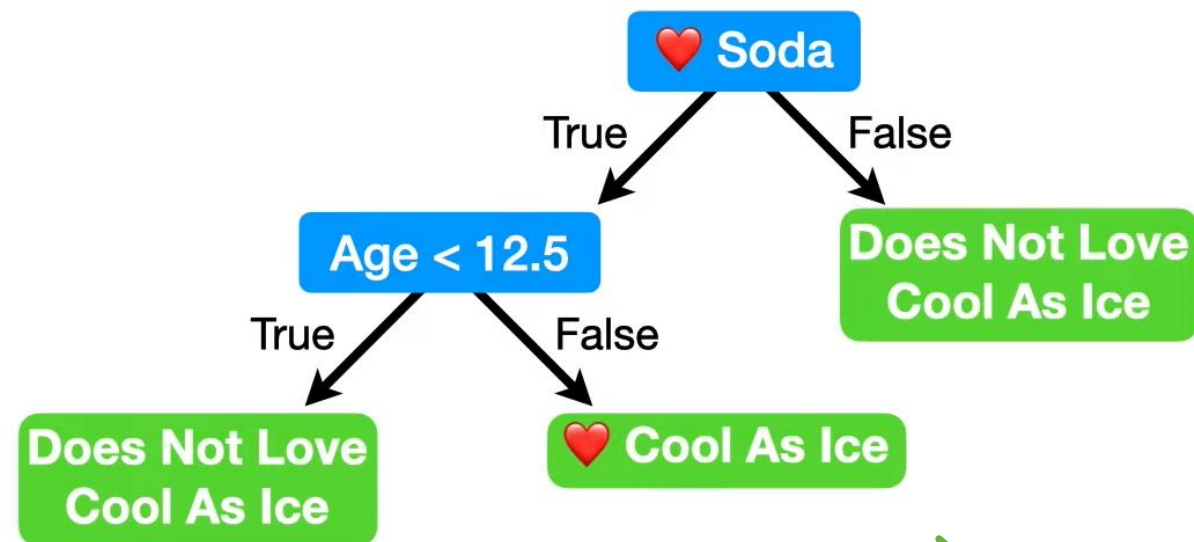
Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No



Drzewo decyzyjne – jak zbudować

Dane wejściowe

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No

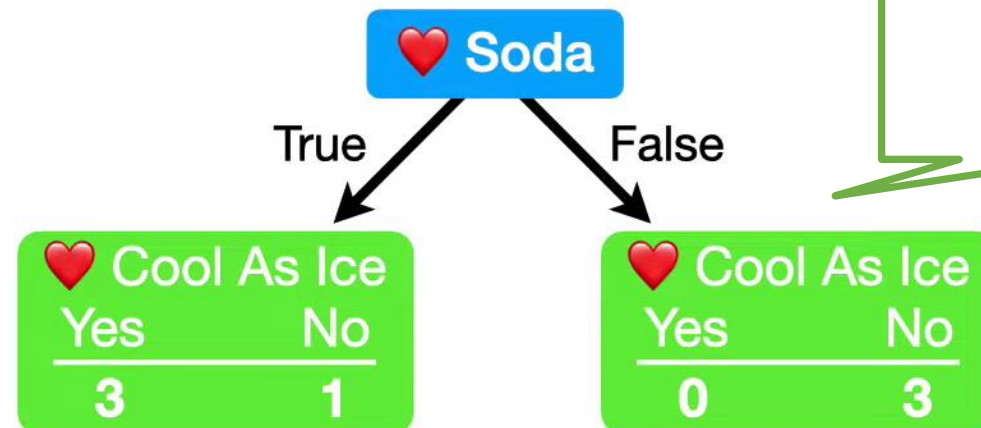
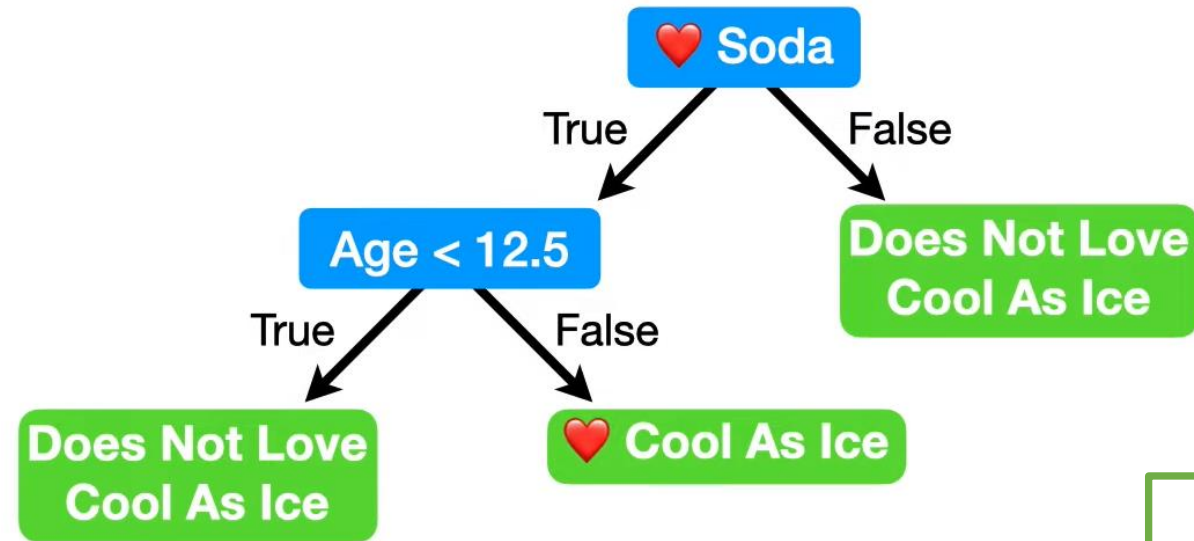


Finalne drzewo decyzyjne

Drzewo decyzyjne – jak zbudować

Dane wejściowe

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No

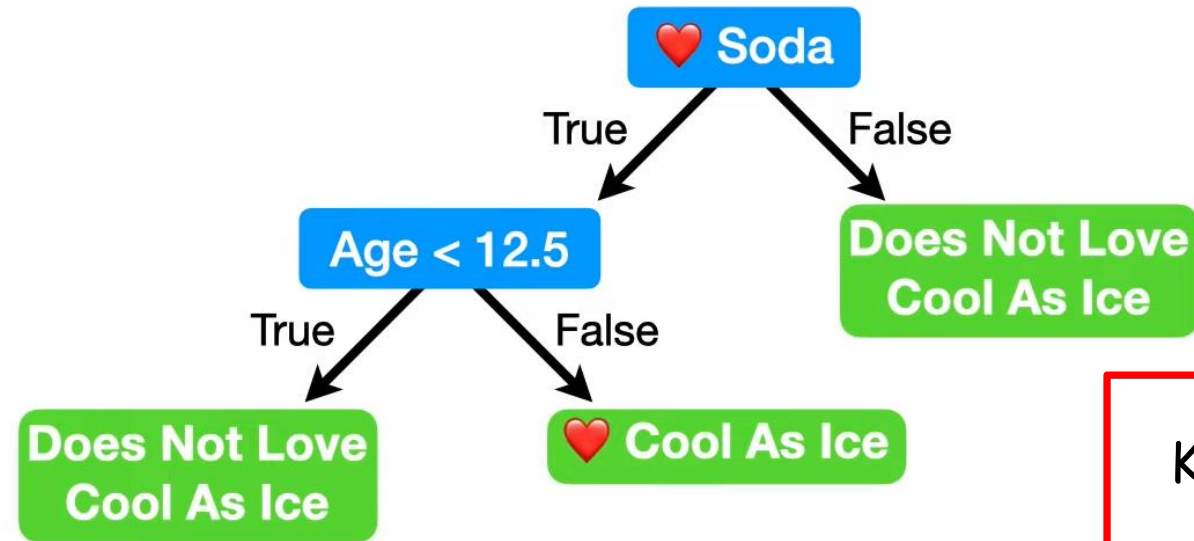


Uproszczone
drzewo
decyzyjne

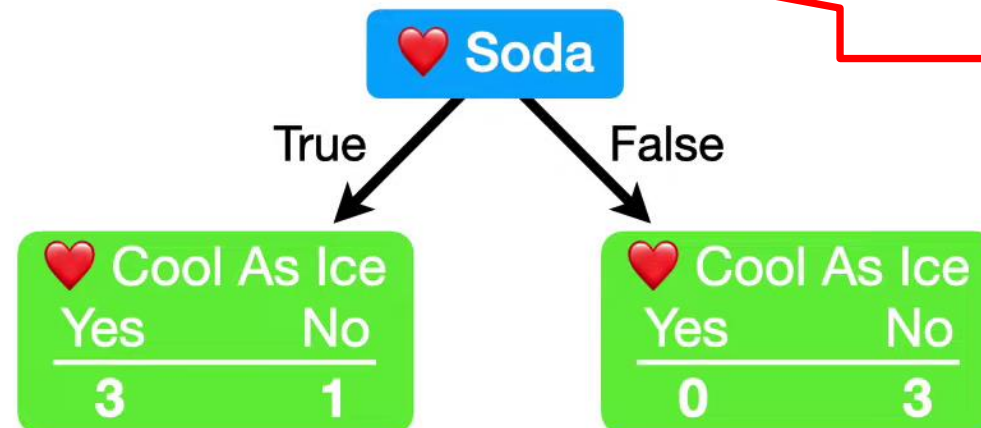
Drzewo decyzyjne – jak zbudować

Dane wejściowe

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No



Które lepsze?
Walidacja
krzyżowa



Drzewo decyzyjne – brak danych

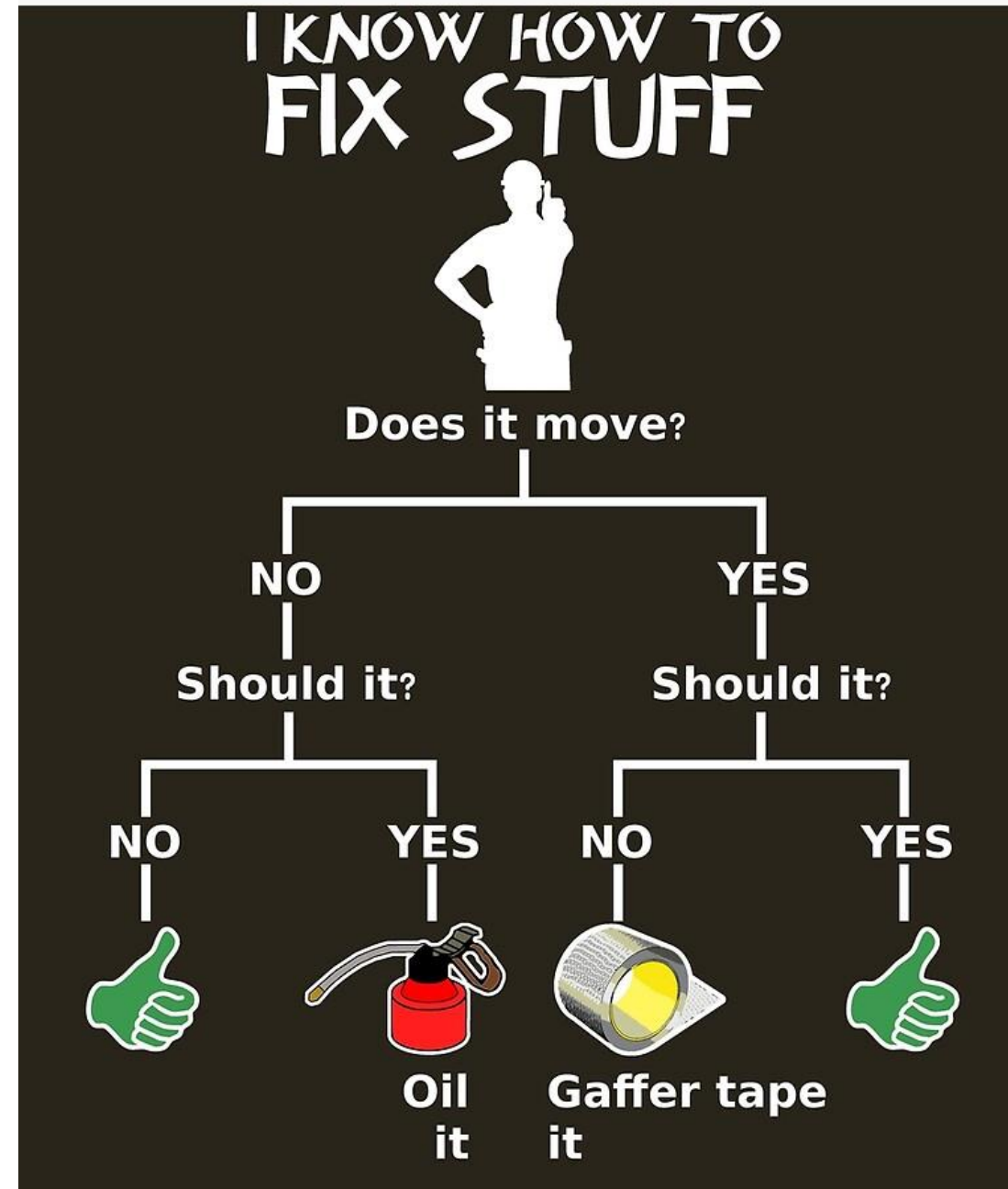
Dane wejściowe

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	?	Yes
No	Yes	35	Yes
Yes	?	38	Yes
Yes	No	50	No
No	No	83	No

- ❖ Usuwanie rekordów z brakującymi danymi
- ❖ Uzupełnianie brakujących danych
- ❖ Użycie algorytmów imputacji
- ❖ Znaczniki brakujących danych
- ❖ Użycie algorytmów uczenia maszynowego do prognozowania brakujących wartości

Kodujemy

public/mmajew/MIW/05/
00_decision tree.py

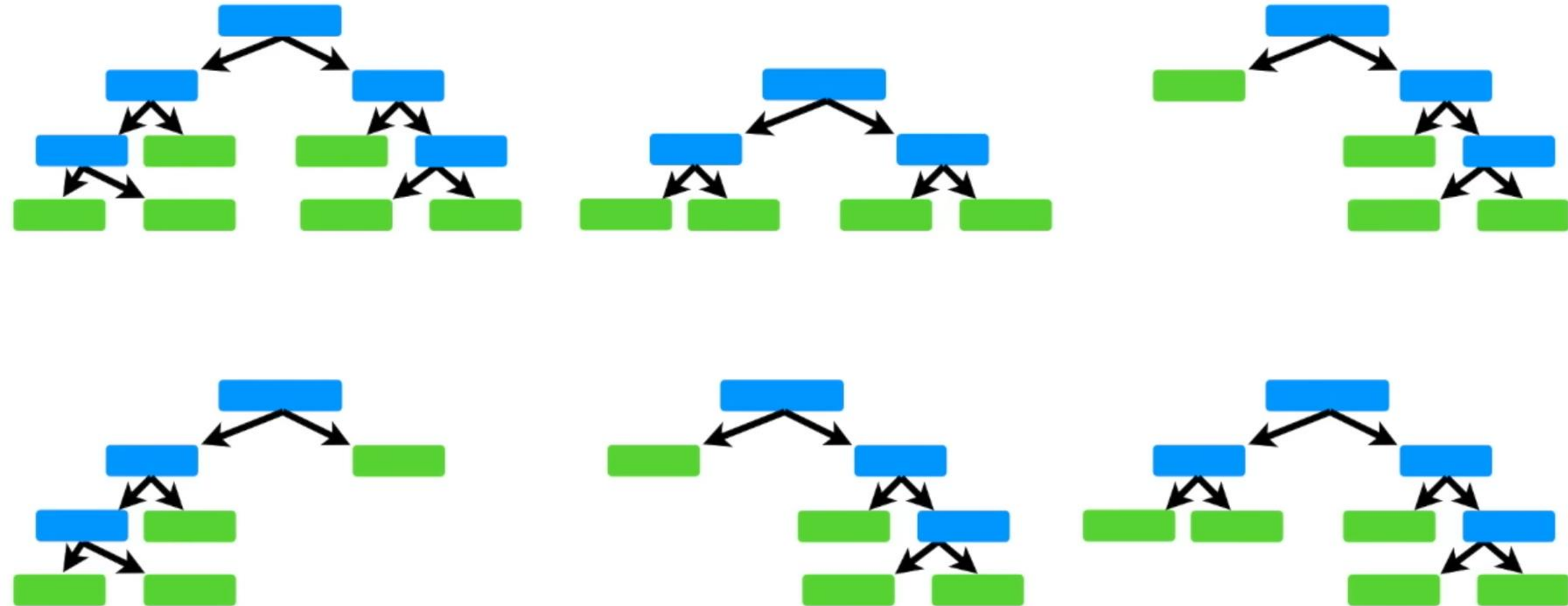


Losowy las

Losowy las to **zbiór drzew decyzyjnych**, które są **szkolone niezależnie** od siebie na losowych podzbiorach danych i/lub losowych podzbiorach cech, a **następnie agregowane** w celu uzyskania stabilnego i **wydajnego modelu klasyfikacji** lub regresji.



Losowy las



Dobór liczby drzew w losowym lesie może być ustalony **empirycznie przy użyciu walidacji krzyżowej** lub za pomocą **metod automatycznego strojenia hiperparametrów**, takich jak `RandomizedSearchCV` lub `GridSearchCV`, w celu znalezienia optymalnej liczby drzew, która zapewnia najlepszą wydajność na zbiorze testowym. **Typowe wartości liczby drzew** znajdują się w **zakresie od kilkudziesięciu do kilkuset**, ale ostateczny wybór zależy od charakterystyki danych i złożoności problemu klasyfikacji lub regresji.

Drzewo decyzyjne vs. Losowy las

Drzewo decyzyjne (Decision Tree Classifier)

- ✓ **Prostota interpretacji:** Drzewa decyzyjne są łatwe do interpretacji i wizualizacji
- ✓ **Obsługa zarówno danych kategorycznych, jak i numerycznych**
- ✓ **Skuteczność w danych o niskiej do średniej złożoności:** W przypadku danych o prostszej strukturze, drzewa decyzyjne mogą osiągać dobre wyniki bez konieczności zbyt skomplikowanych procedur.
- **Podatność na przeuczenie:** Drzewa decyzyjne mają tendencję do tworzenia zbyt skomplikowanych modeli, które dobrze dopasowują się do danych treningowych, ale słabo generalizują na nowe dane.
- **Niestabilność:** Małe zmiany w danych treningowych mogą prowadzić do znaczących zmian w strukturze drzewa decyzyjnego.

Losowy las (Random Forest Classifier)

- ✓ **Odporność na przeuczenie:** Losowe lasy agregują wiele drzew, co pomaga w średnich wynikach i zmniejsza wariancję modelu.
- ✓ **Skuteczność na dużych zbiorach danych:** Losowe lasy dobrze radzą sobie z dużymi zbiorami danych i wieloma cechami
- ✓ **Brak potrzeby ręcznego strojenia parametrów:** W przeciwieństwie do drzew decyzyjnych, losowe lasy rzadko wymagają ręcznego dostrojenia parametrów
- **Mniej interpretowalne:** Losowe lasy są mniej interpretowalne niż pojedyncze drzewa decyzyjne
- **Zwiększone zużycie zasobów:** W porównaniu z pojedynczym drzewem decyzyjnym, losowe lasy mogą wymagać więcej zasobów obliczeniowych ze względu na potrzebę szkolenia i predykcji wielu drzew.

Losowy las – próbki początkowe

Original Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	No	Yes	167	Yes

Bootstrapped Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

bootstrap losowa próba z powtórzeniami

Losowy las – próbki początkowe

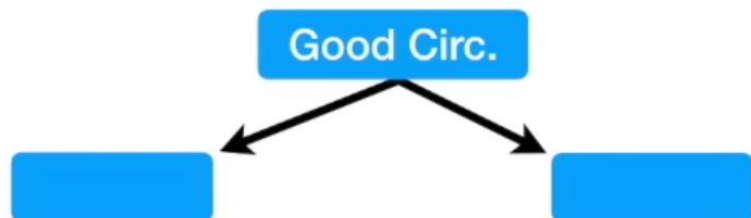
*"feature bagging" lub "feature subsampling",
losujemy losowe podzbiory cech
np. kolumny GBC i BA*

Bootstrapped Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

Losowy las – próbki początkowe

Np. po obliczeniu IG jako korzeń mamy GC...



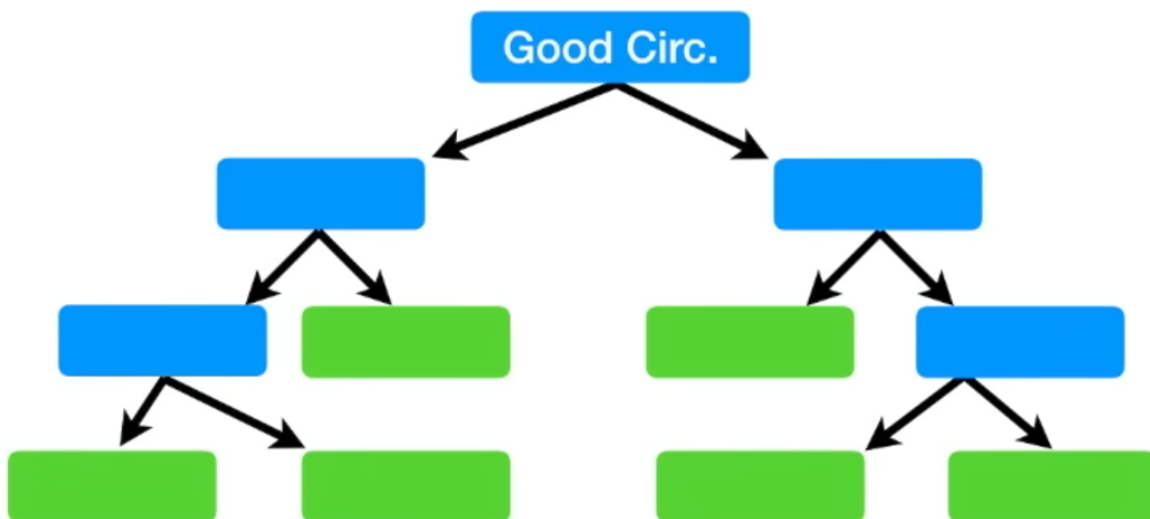
Bootstrapped Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

...i losowo wybieramy kolejne kolumny do analizy

Losowy las – próbki początkowe

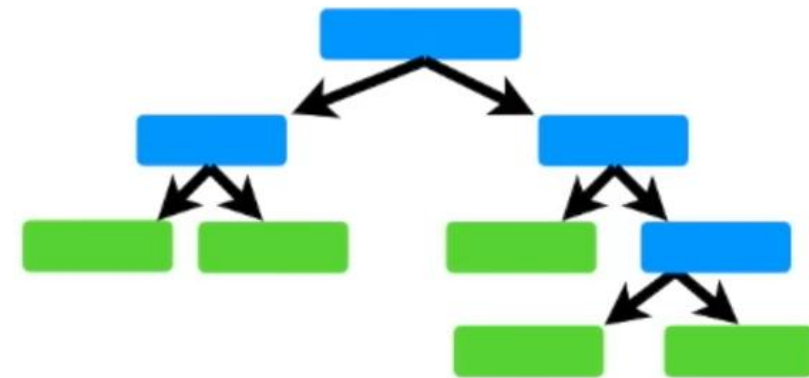
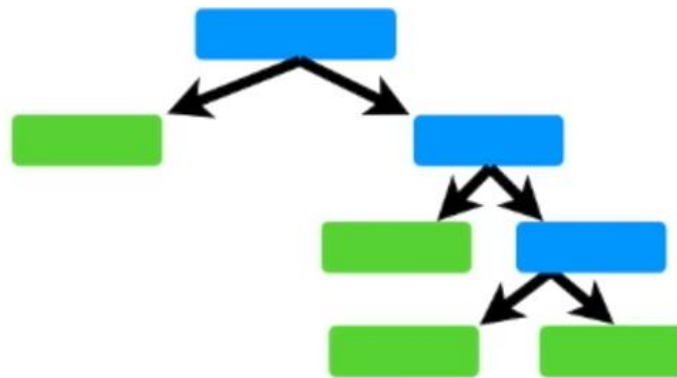
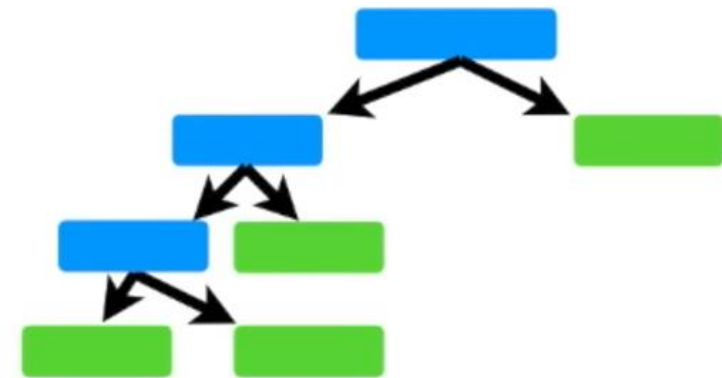
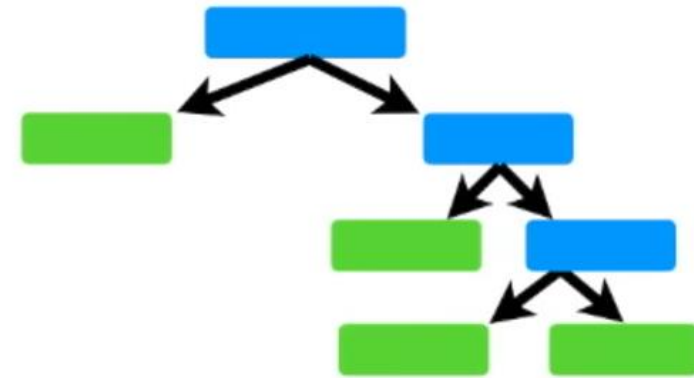
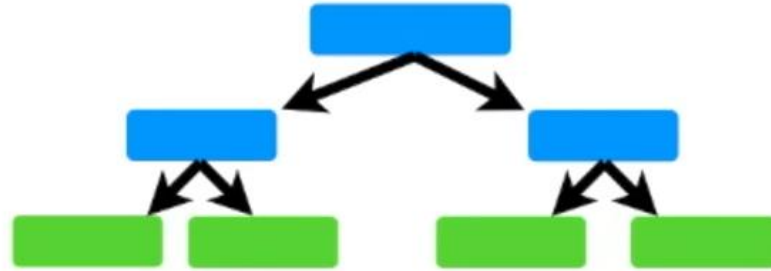
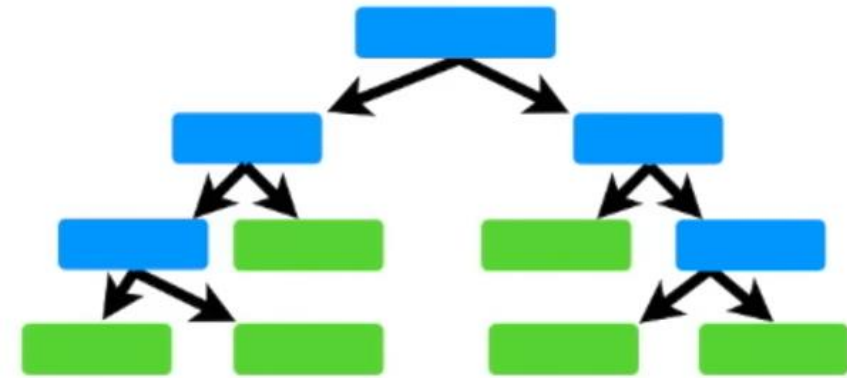
I dla danego bootstrapped data mamy drzewko



Bootstrapped Dataset

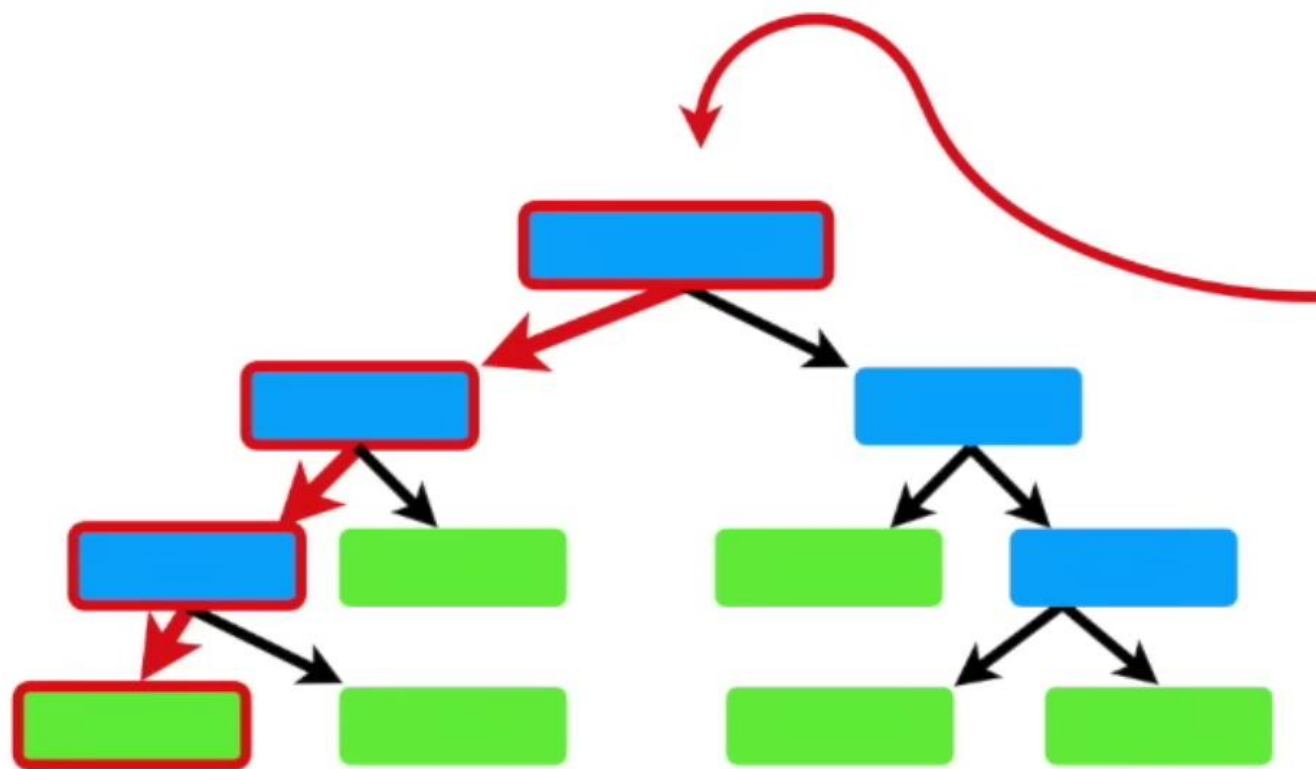
Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

Losowy las



I tak mamy losowy las, jak go użyć?

Losowy las – nowy pacjent

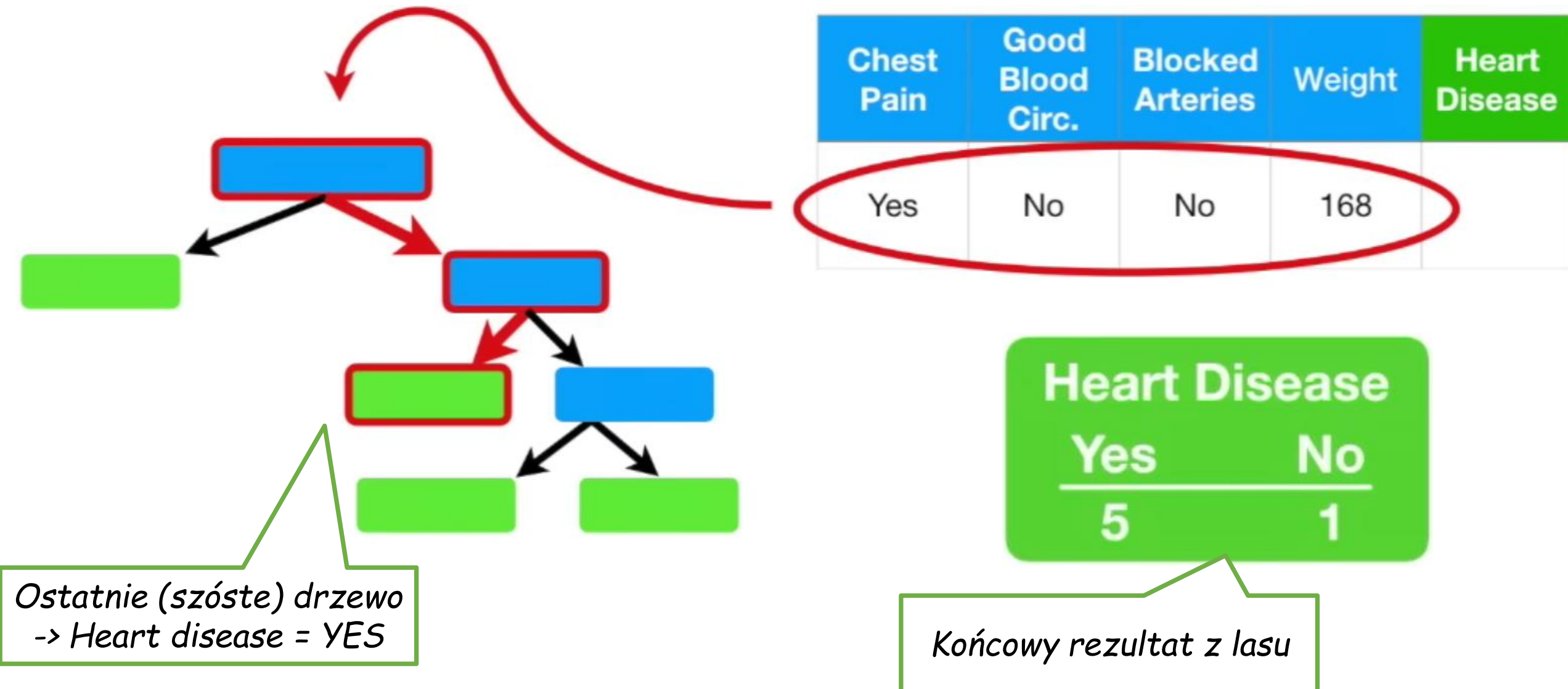


Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	No	No	168	

Heart Disease	
Yes	No
1	0

Pierwsze drzewo -> Heart disease = YES

Losowy las – nowy pacjent



Losowy las – testowanie danych

Original Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	No	Yes	167	Yes

Out-Of-Bag Dataset

Bootstrapped Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

Losowy las – na co zwrócić uwagę?

- ☐ Liczba drzew **OVER 9000!!!**
- ☐ Parametry drzewa, np.: maksymalna głębokość, minimalna liczba próbek wymaganych do podziału węzła i minimalna liczba próbek wymaganych w liściu
- ☐ Liczba cech
- ☐ Kryterium podziału
- ☐ Różnorodność danych (zwiększenie różnorodności między drzewami)
- ☐ Ocena modelu (OOB, walidacja krzyżowa)
- ☐ Interpretowalność modelu
- ☐ Skalowanie
- ☐ Monitorowanie wydajności



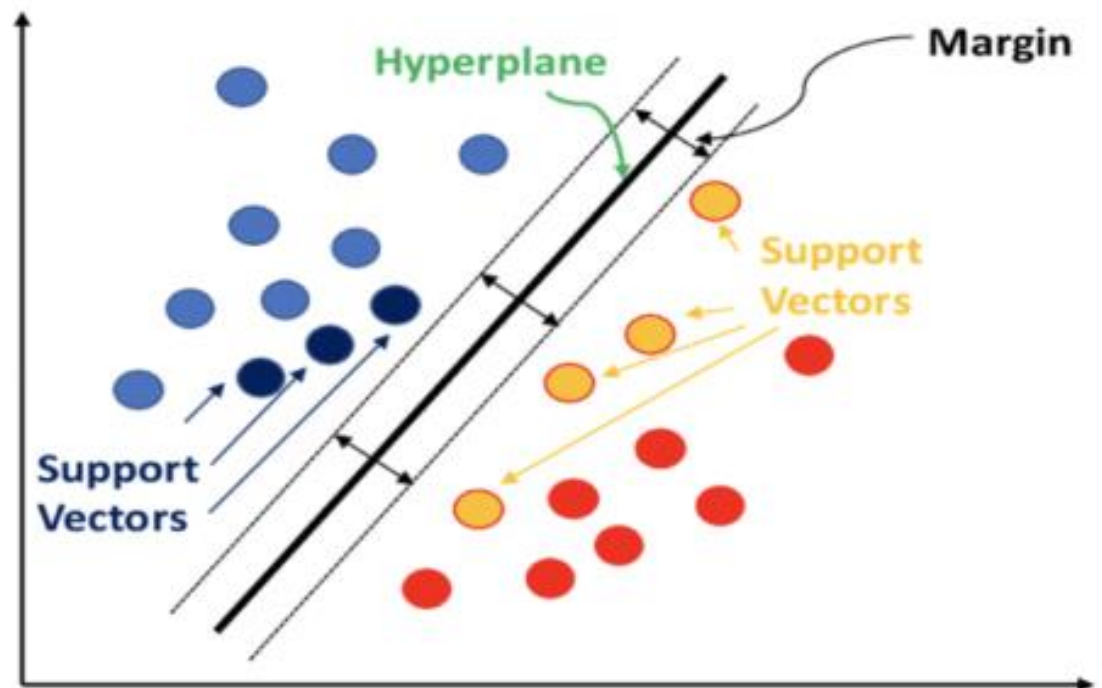
Kodujemy

public/mmajew/MIW/05/
01_random forest.py



Maszyny wektorów nośnych (Support Vector Machines, SVM)

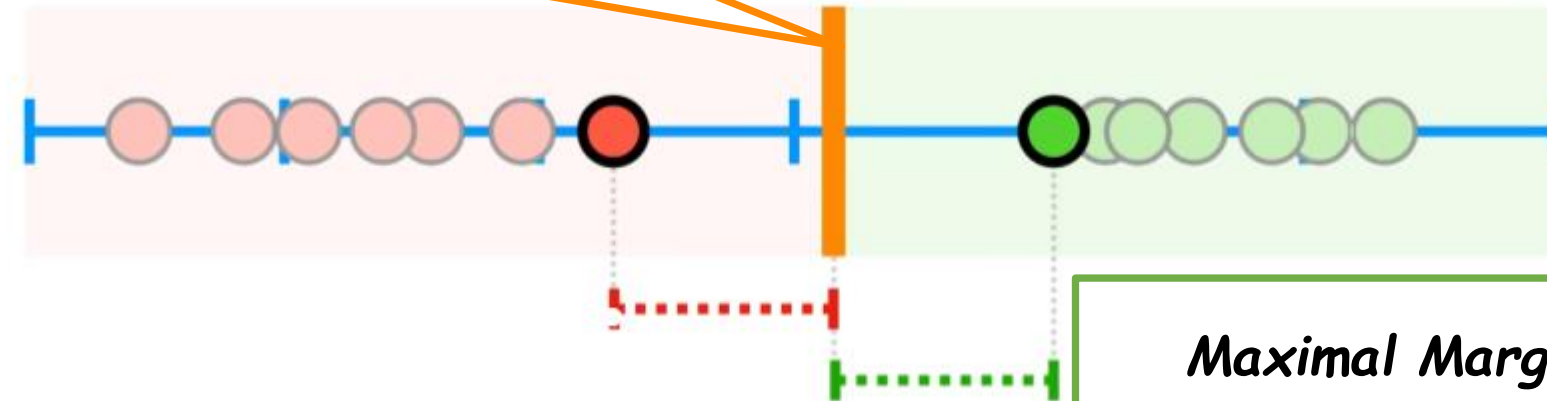
- ❖ znalezienie **optymalnego zestawu hiperpłaszczyzn**, które mogą skutecznie rozdzielić dane należące do różnych klas lub przewidywać wartości w przypadku regresji.
- ❖ **maksymalizacja marginesu między hiperpłaszczyzną a najbliższymi punktami danych** różnych klas, zwanych **wektorami nośnymi**.
- ❖ w przypadku **nieliniowych zbiorów danych**, SVM może stosować tzw. **funkcje jądra**, które pozwalają na wyznaczenie bardziej skomplikowanych granic decyzyjnych.



<https://datatron.com/wp-content/uploads/2021/05/Support-Vector-Machine.png>

Maszyny wektorów nośnych (Support Vector Machines, SVM)

Próg, wartość graniczna *threshold*

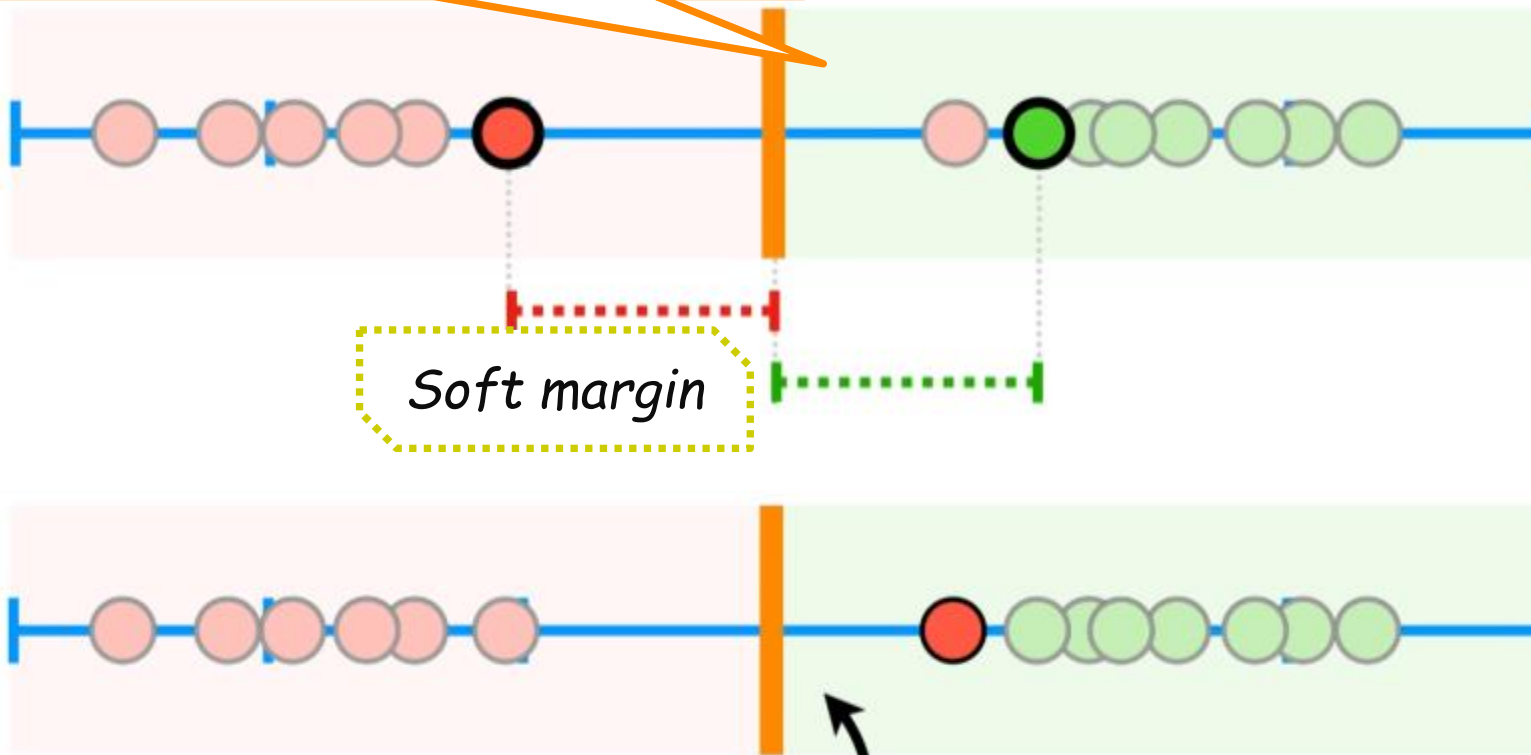


Dane odstające *outliers*

low bias, high variance

Maszyny wektorów nośnych (Support Vector Machines, SVM)

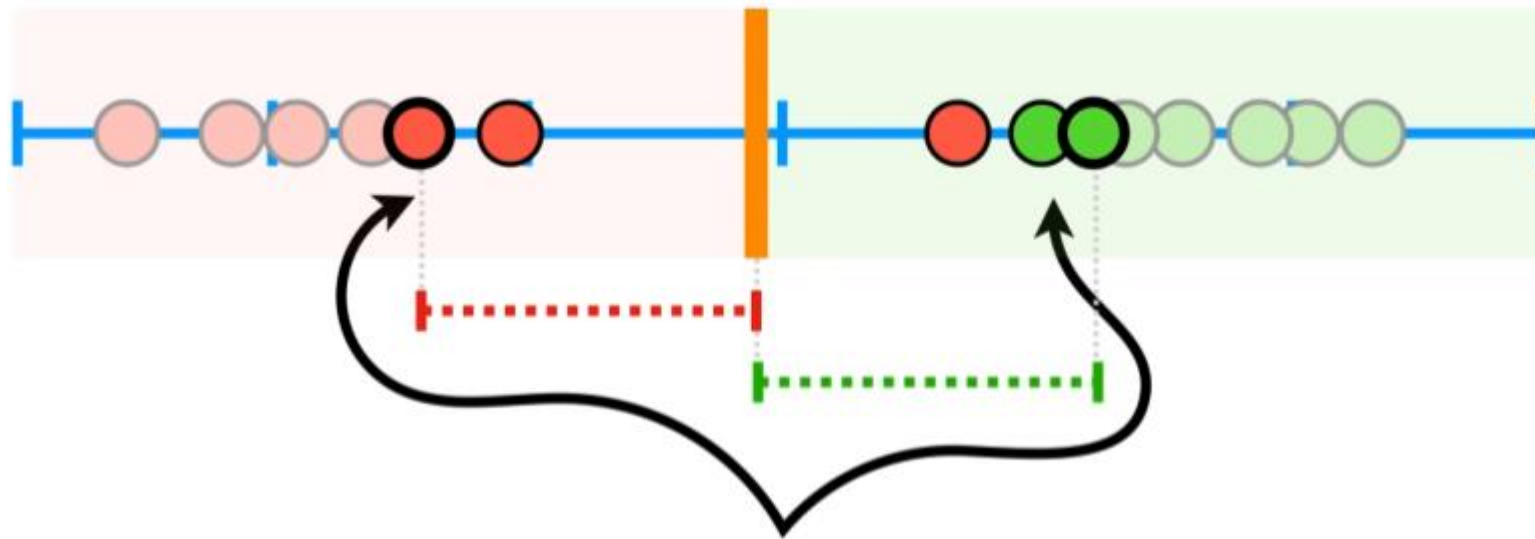
Próg, wartość graniczna *threshold*



Błąd systematyczny na danych
treningowych *bias*

*higher bias,
lower variance*

Maszyny wektorów nośnych (Support Vector Machines, SVM)

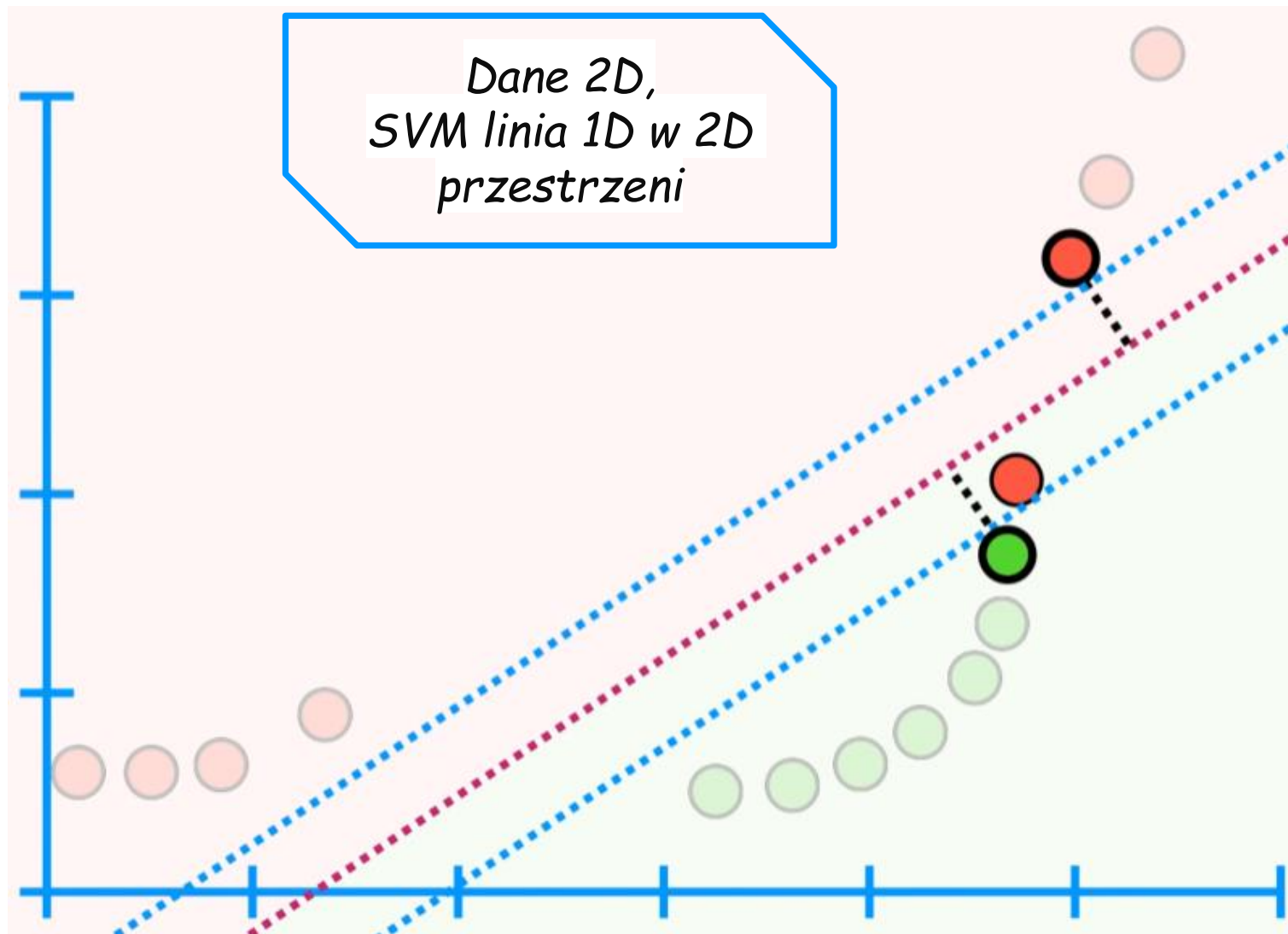


Dane na krawędzi i wewnątrz *soft margin* nazywamy wektorami nośnymi/wspierającymi *support vectors*

Walidacja krzyżowa
w celu ustalenia *soft margin*

Soft Margin Classifier
aka
Support Vector Classifier

Maszyny wektorów nośnych (Support Vector Machines, SVM)

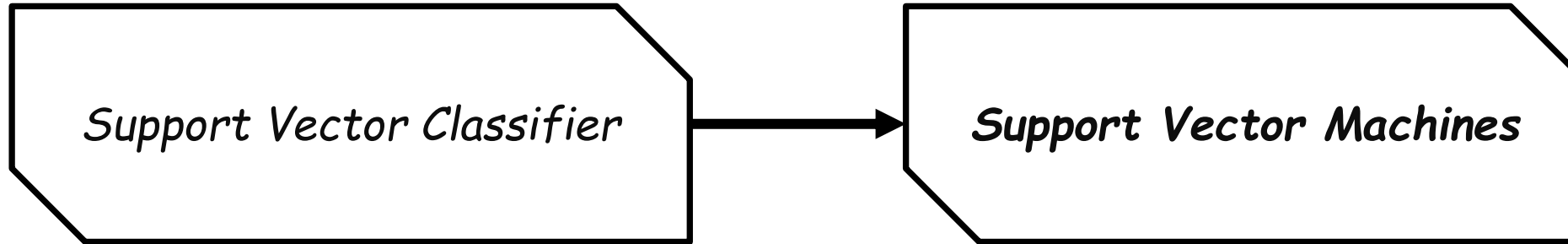


Granica decyzji,
soft margin

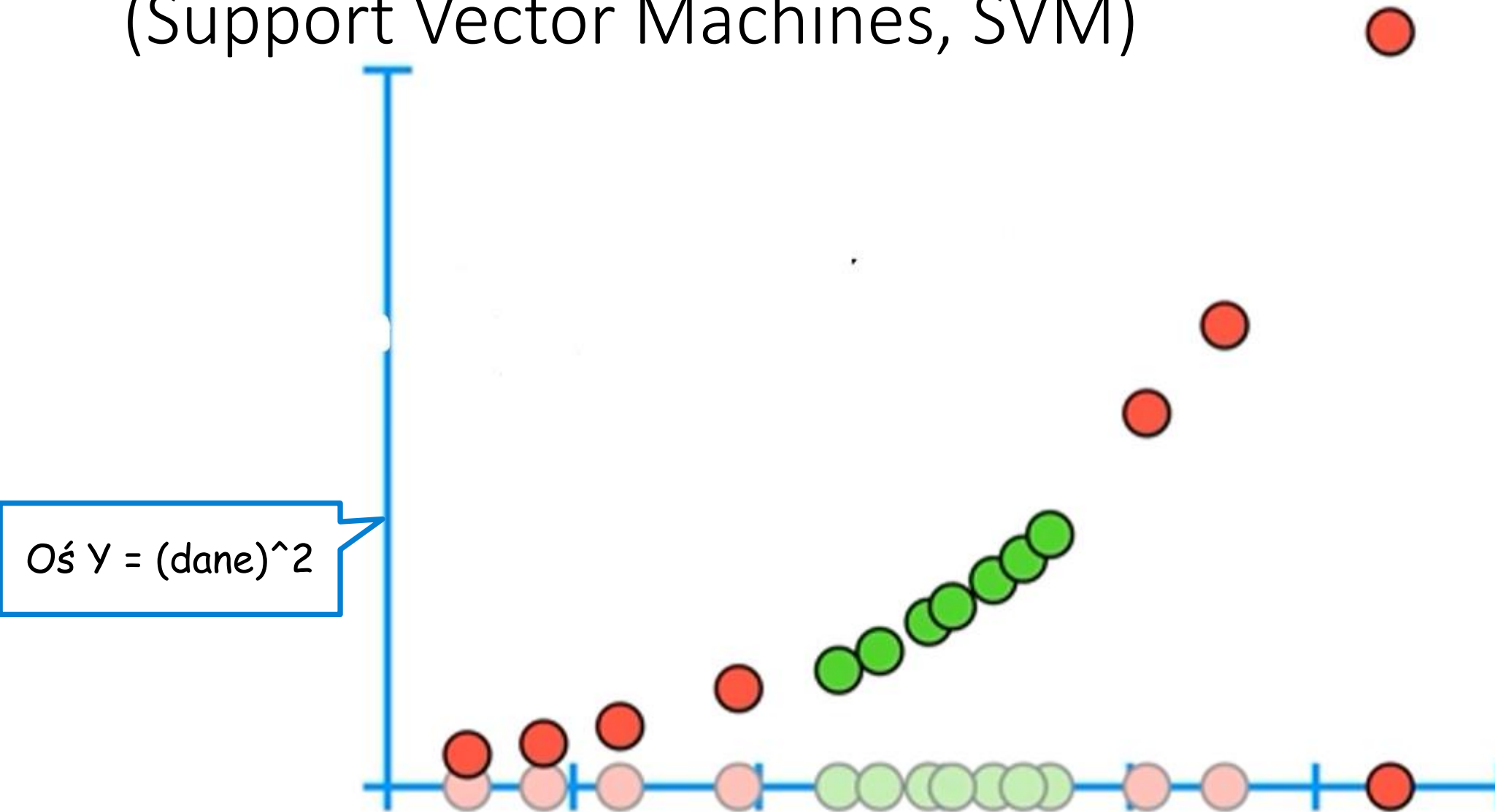
"**Overlapping**" klasy w danych treningowych nakładają się. Nie ma jednoznacznej granicy decyzyjnej (hiperpłaszczyzny)

Maszyny wektorów nośnych (Support Vector Machines, SVM)

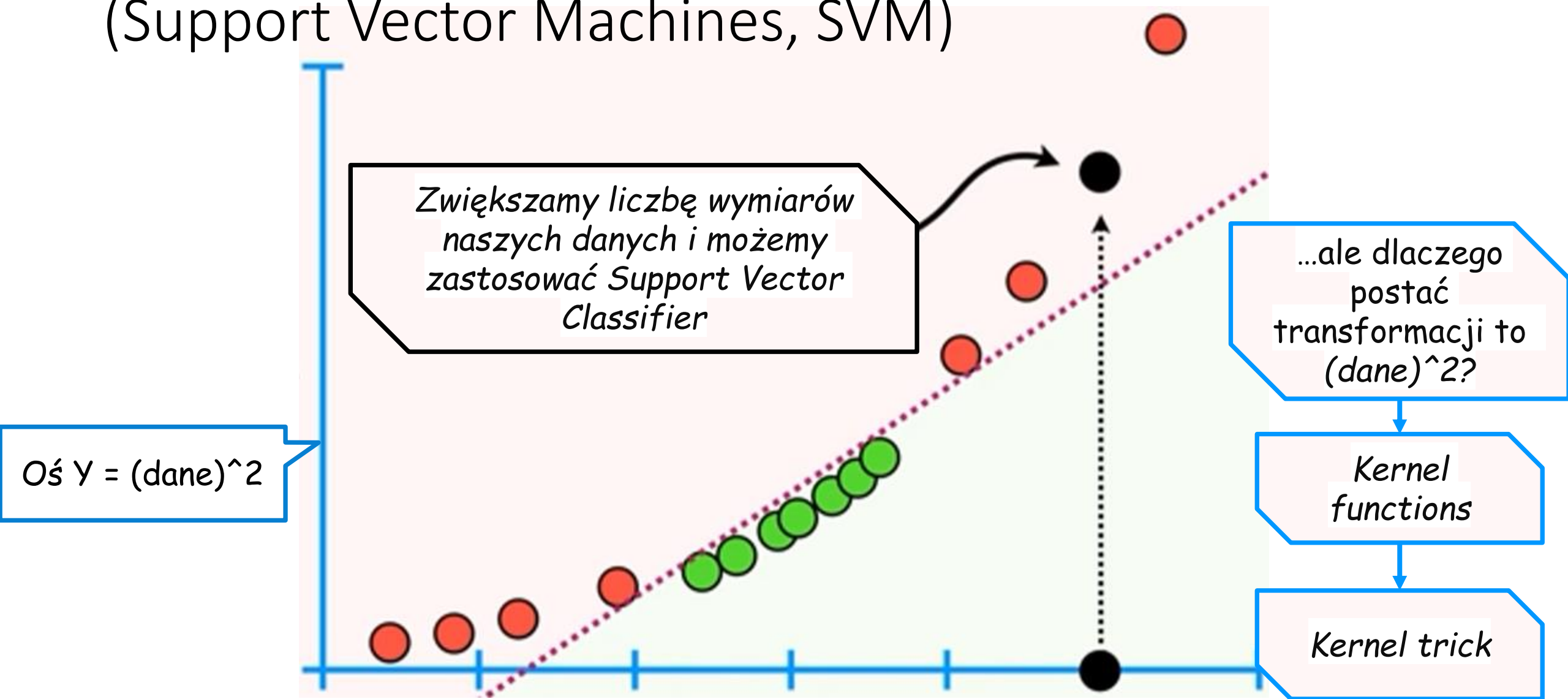
"**Overlapping**" klasy w danych
treningowych nakładają się.
Nie ma jednoznacznej granicy
decyzyjnej (hiperpłaszczyzny)



Maszyny wektorów nośnych (Support Vector Machines, SVM)



Maszyny wektorów nośnych (Support Vector Machines, SVM)



Kodujemy

public/mmajew/MIW/05/
02_SVM_log reg.py

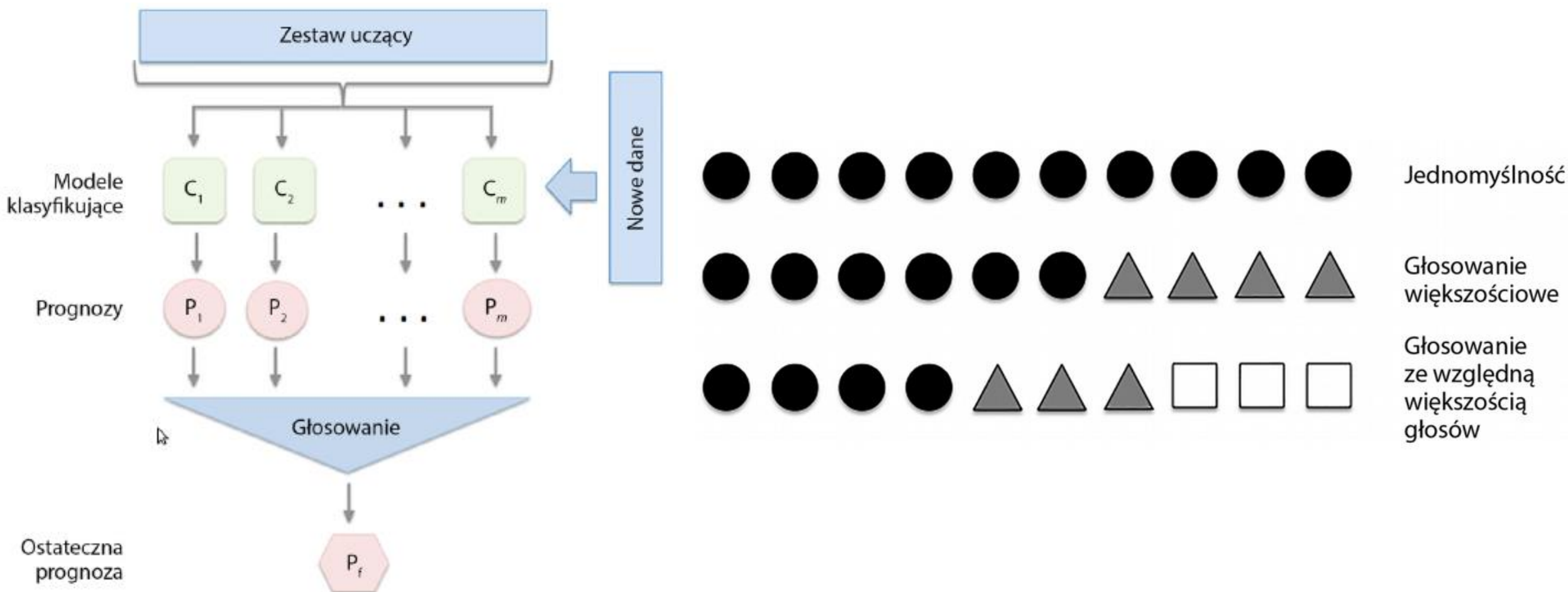


Coding
the SVM
algorithm in numpy

from sklearn
import svm

Uczenie zespołów

Celem metod zespołowych (*ensemble methods*) jest łączenie różnych klasyfikatorów w jeden meta klasyfikator wykazujący większą skuteczność uogólniania niż każdy ze składowych algorytmów.



Projekt 3 a. (na 5 pkt.)

Szkic: [public/mmajew/MIW/05/03_VotingClassifier.py](#)

Projekt wykonaj z wykorzystaniem pakietu scikit-learn

1. Stwórz zbiór danych za pomocą funkcji `make_moons(n_samples = 10000, noise = 0.4)`.
2. Rozdziel uzyskany zestaw danych na podzbiory uczący i testowy przy użyciu metody `train_test_split()`.
- 3.a Wytrenuj klasyfikatory `LogisticRegression`, `SVM` oraz `RandomForestClassifier` i połącz klasyfikatory `SVM`, `LogisticRegression` oraz `RandomForestClassifier` w **jeden zespół (`VotingClassifier`)**.
4. Oceń osiągnięte rezultaty: wyświetl dokładność trenowania i testowania, narysuj kontur decyzji na podstawie predykcji `VotingClassifier`.

Termin 30.04

Drugą część projektu (projekt 3 b. za kolejne 5 pkt podam na kolejnych zajęciach)