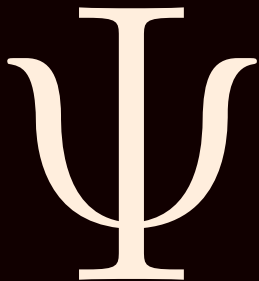


Metody Sztucznej Inteligencji

Wykład 2: **Zadanie klasyfikacji**

dr inż. Paweł Ksieniewicz
Katedra Systemów i Sieci Komputerowych

13 marca 2019





Uczenie nadzorowane



Uczenie nadzorowane

↪ Regresja:

Uczenie nadzorowane

↪ Regresja:

- × predykcja **liczby rzeczywistej** powiązanej z zadany wzorcem (wektorem cech).

Uczenie nadzorowane

↪ Regresja:

- × predykcja **liczby rzeczywistej** powiązanej z zadany wzorcem (wektorem cech).

↪ Klasyfikacja:

Uczenie nadzorowane

↪ Regresja:

- × predykcja **liczby rzeczywistej** powiązanej z zadany wzorcem (wektorem cech).

↪ Klasyfikacja:

- × predykcja **liczby całkowitej** (etykiety) powiązanej z zadany wzorcem.

Przykład

ZWIERZĘ	jajorodne?	ma łuski?	jadowite?	zmiennociepne?	ile nóg?	czy gad?
<i>kobra</i>	1	1	1	1	0	1
<i>grzechotnik</i>	1	1	1	1	0	1
<i>boa dusiciel</i>	0	1	0	1	0	1
<i>kurczak</i>	1	1	0	1	2	0
<i>gupik</i>	0	1	0	0	0	0
<i>Andrzej</i>	0	0	0	0	2	0
<i>zebra</i>	0	0	0	0	4	0
<i>pyton</i>	1	1	0	1	0	1
<i>aligator</i>	1	1	0	1	4	1

	<i>kobra</i>	<i>grzechotnik</i>	<i>boa dusiciel</i>	<i>kura</i>	<i>gupik</i>	<i>Andrzej</i>	<i>zebra</i>	<i>pyton</i>	<i>aligator</i>
<i>kobra</i>	—	0.00 ₁	3.14 ₄	2.70 ₃	3.79 ₆	4.65 ₇	5.11 ₈	2.41 ₂	3.43 ₅
<i>grzechotnik</i>	0.00 ₁	—	3.14 ₄	2.70 ₃	3.79 ₆	4.65 ₇	5.11 ₈	2.41 ₂	3.43 ₅
<i>boa dusiciel</i>	3.14 ₄	3.14 ₅	—	2.36 ₃	2.12 ₂	3.43 ₇	4.04 ₈	2.01 ₁	3.17 ₆
<i>kurczak</i>	2.70 ₄	2.70 ₅	2.36 ₃	—	3.17 ₆	3.79 ₇	3.98 ₈	1.22 ₁	1.22 ₂
<i>gupik</i>	3.79 ₆	3.79 ₇	2.12 ₁	3.17 ₄	—	2.70 ₂	3.43 ₅	2.92 ₃	3.81 ₈
<i>Andrzej</i>	4.65 ₇	4.65 ₈	3.43 ₃	3.79 ₄	2.70 ₂	—	1.22 ₁	3.98 ₅	3.98 ₆
<i>zebra</i>	5.11 ₇	5.11 ₈	4.04 ₅	3.98 ₄	3.43 ₂	1.22 ₁	—	4.51 ₆	3.79 ₃
<i>pyton</i>	2.41 ₃	2.41 ₄	2.01 ₂	1.22 ₁	2.92 ₆	3.98 ₇	4.51 ₈	—	2.45 ₅
<i>aligator</i>	3.43 ₄	3.43 ₅	3.17 ₃	1.22 ₁	3.81 ₇	3.98 ₈	3.79 ₆	2.45 ₂	—
<i>czy gad?</i>	1	1	1	0	0	0	0	1	1

Podejście minimalno-odległościowe

Podejście minimalno-odległościowe

↪ Najbardziej banalnym rozwiązaniem jest **najbliższy sąsiad**.

Podejście minimalno-odległościowe

- ↪ Najbardziej banalnym rozwiązaniem jest **najbliższy sąsiad**.
- ↪ Jest to tak zwane **uczenie leniwe**.

Podejście minimalno-odległościowe

- ↪ Najbardziej banalnym rozwiązaniem jest **najbliższy sąsiad**.
- ↪ Jest to tak zwane **uczenie leniwe**.
- ↪ Uczenie polega wprost na zapamiętaniu wszystkich przypadków (wzorców, obiektów) ze zbioru uczącego.

Podójście minimalno-odległójściowe

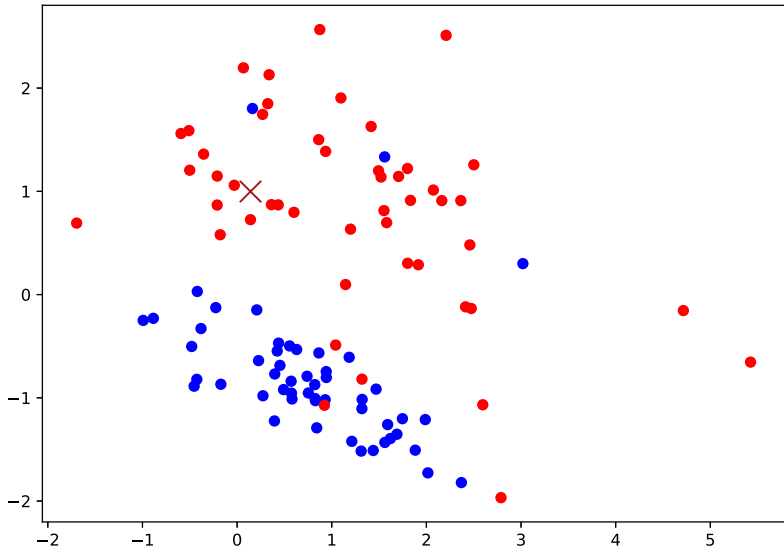
- ↪ Najbardziej banalnym rozwiązaniem jest **najbliŹszy sąsiad**.
- ↪ Jest to tak zwane **uczenie leniwe**.
- ↪ Uczenie polega wprost na zapamiętaniu wszystkich przypadków (wzorców, obiektów) ze zbioru uczącego.
- ↪ Predykcja składa się z całych dwóch kroków:

Podjęcie minimalno-odległościowe

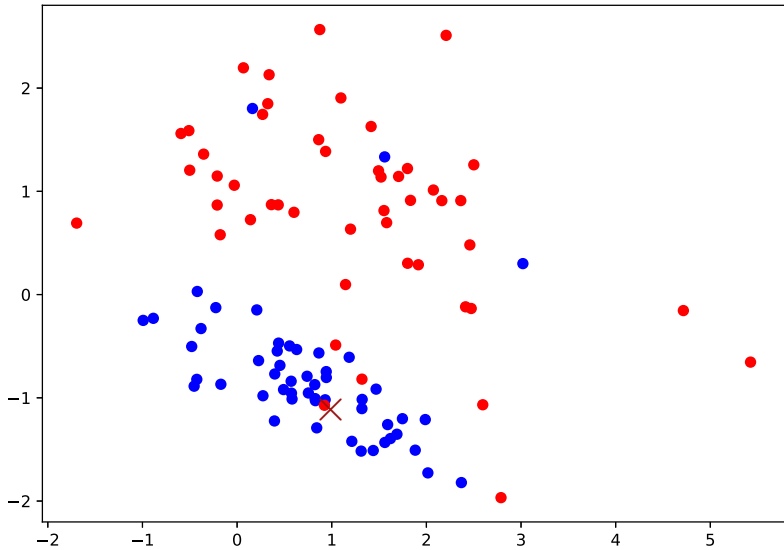
- ↪ Najbardziej banalnym rozwiązaniem jest **najbliższy sąsiad**.
- ↪ Jest to tak zwane **uczenie leniwe**.
- ↪ Uczenie polega wprost na zapamiętaniu wszystkich przypadków (wzorców, obiektów) ze zbioru uczącego.
- ↪ Predykcja składa się z całych dwóch kroków:
 - × znajdź najbliższego sąsiada,

Podjęcie minimalno-odległościowe

- ↪ Najbardziej banalnym rozwiązaniem jest **najbliższy sąsiad**.
- ↪ Jest to tak zwane **uczenie leniwe**.
- ↪ Uczenie polega wprost na zapamiętaniu wszystkich przypadków (wzorców, obiektów) ze zbioru uczącego.
- ↪ Predykcja składa się z całych dwóch kroków:
 - × znajdź najbliższego sąsiada,
 - × zwróć etykietę odnalezionego najbliższego sąsiada.



	<i>kobra</i>	<i>grzechotnik</i>	<i>boa dusiciel</i>	<i>kura</i>	<i>gupik</i>	<i>Andrzej</i>	<i>zebra</i>	<i>pyton</i>	<i>aligator</i>
<i>kobra</i>	—	0.00 ₁	3.14 ₄	2.70 ₃	3.79 ₆	4.65 ₇	5.11 ₈	2.41 ₂	3.43 ₅
<i>grzechotnik</i>	0.00 ₁	—	3.14 ₄	2.70 ₃	3.79 ₆	4.65 ₇	5.11 ₈	2.41 ₂	3.43 ₅
<i>boa dusiciel</i>	3.14 ₄	3.14 ₅	—	2.36 ₃	2.12 ₂	3.43 ₇	4.04 ₈	2.01 ₁	3.17 ₆
<i>kurczak</i>	2.70 ₄	2.70 ₅	2.36 ₃	—	3.17 ₆	3.79 ₇	3.98 ₈	1.22 ₁	1.22 ₂
<i>gupik</i>	3.79 ₆	3.79 ₇	2.12 ₁	3.17 ₄	—	2.70 ₂	3.43 ₅	2.92 ₃	3.81 ₈
<i>Andrzej</i>	4.65 ₇	4.65 ₈	3.43 ₃	3.79 ₄	2.70 ₂	—	1.22 ₁	3.98 ₅	3.98 ₆
<i>zebra</i>	5.11 ₇	5.11 ₈	4.04 ₅	3.98 ₄	3.43 ₂	1.22 ₁			
<i>pyton</i>	2.41 ₃	2.41 ₄	2.01 ₂	1.22 ₁	2.92 ₆	3.98 ₇			
<i>aligator</i>	3.43 ₄	3.43 ₅	3.17 ₃	1.22 ₁	3.81 ₇	3.98 ₈			
<i>czy gad?</i>	1	1	1	0	0	0	0	1	1



K najbliższych sąsiadów — k NN

K najbliższych sąsiadów — k NN

↪ Jest to nadal **uczenie leniwe**.

K najbliższych sąsiadów — k NN

↪ Jest to nadal **uczenie leniwe**.

↪ Predykcja nadal składa się z całych dwóch kroków:

K najbliższych sąsiadów — k NN

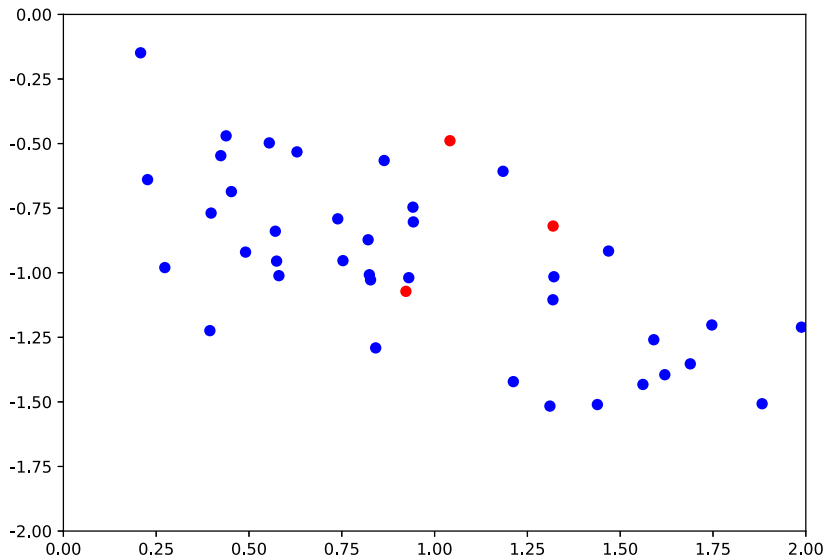
- ↪ Jest to nadal **uczenie leniwe**.
- ↪ Predykcja nadal składa się z całych dwóch kroków:
 - × znajdź k najbliższych sąsiadów,

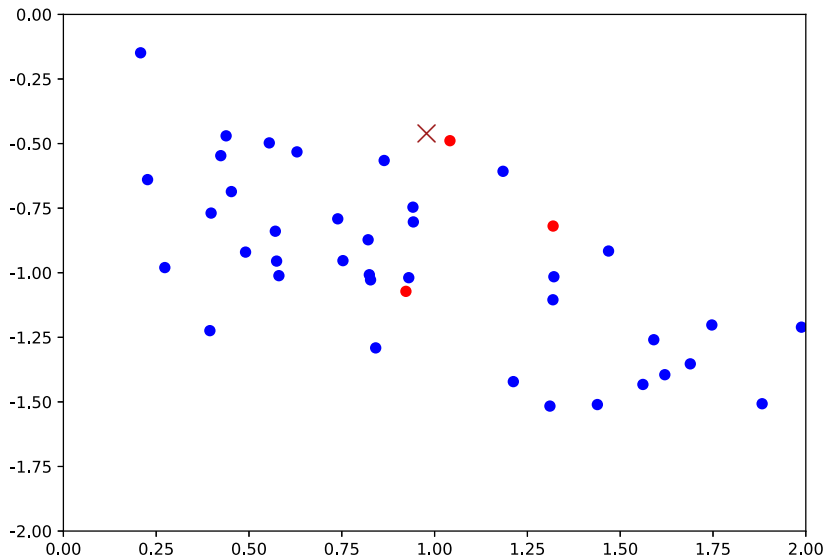
K najbliższych sąsiadów — k NN

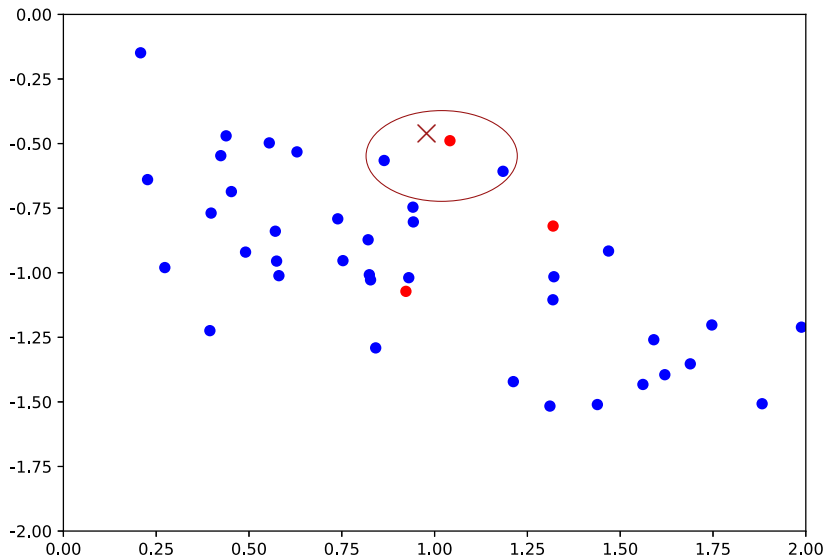
↪ Jest to nadal **uczenie leniwe**.

↪ Predykcja nadal składa się z całych dwóch kroków:

- × znajdź k najbliższych sąsiadów,
- × zwróć najczęściej pojawiającą się etykietę wśród odnalezionych najbliższych sąsiadów.







	<i>kobra</i>	<i>grzechotnik</i>	<i>boa dusiciel</i>	<i>kura</i>	<i>gupik</i>	<i>Andrzej</i>	<i>zebra</i>	<i>pyton</i>	<i>aligator</i>
<i>kobra</i>	—	0.00 ₁	3.14 ₄	2.70 ₃	3.79 ₆	4.65 ₇	5.11 ₈	2.41 ₂	3.43 ₅
<i>grzechotnik</i>	0.00 ₁	—	3.14 ₄	2.70 ₃	3.79 ₆	4.65 ₇	5.11 ₈	2.41 ₂	3.43 ₅
<i>boa dusiciel</i>	3.14 ₄	3.14 ₅	—	2.36 ₃	2.12 ₂	3.43 ₇	4.04 ₈	2.01 ₁	3.17 ₆
<i>kurczak</i>	2.70 ₄	2.70 ₅	2.36 ₃	—	3.17 ₆	3.79 ₇	3.98 ₈	1.22 ₁	1.22 ₂
<i>gupik</i>	3.79 ₆	3.79 ₇	2.12 ₁	3.17 ₄	—	2.70 ₂	3.43 ₅	2.92 ₃	3.81 ₈
<i>Andrzej</i>	4.65 ₇	4.65 ₈	3.43 ₃	3.79 ₄	2.70 ₂	—	1.22 ₁	3.98 ₅	3.98 ₆
<i>zebra</i>	5.11 ₇	5.11 ₈	4.04 ₅	3.98 ₄	3.43 ₂	1.22 ₁			
<i>pyton</i>	2.41 ₃	2.41 ₄	2.01 ₂	1.22 ₁	2.92 ₆	3.98 ₇			
<i>aligator</i>	3.43 ₄	3.43 ₅	3.17 ₃	1.22 ₁	3.81 ₇	3.98 ₈			
<i>czy gad?</i>	1	1	1	0	0	0	0	1	1

~~~~~

Ile powinno wynosić  $k$ ?

~~~~~

Ile powinno wynosić k ?

↪ Powinno być **nieparzyste**, dla ewidentnego wyniku głosowania.

Ile powinno wynosić k ?

- ↪ Powinno być **nieparzyste**, dla ewidentnego wyniku głosowania.
- ↪ Im **mniejsze** k , tym mniejsza próbka, a więc większy wpływ obserwacji odstających.

Ile powinno wynosić k ?

- ↪ Powinno być **nieparzyste**, dla ewidentnego wyniku głosowania.
- ↪ Im **mniejsze** k , tym mniejsza próbka, a więc większy wpływ obserwacji odstających.
- ↪ Im **większe** k , tym bardziej upraszczamy klasyfikację do wyboru **klasy dominującej według prawdopodobieństwa a priori**.

Ile powinno wynosić k ?

- ↪ Powinno być **nieparzyste**, dla ewidentnego wyniku głosowania.
- ↪ Im **mniejsze** k , tym mniejsza próbka, a więc większy wpływ obserwacji odstających.
- ↪ Im **większe** k , tym bardziej upraszczamy klasyfikację do wyboru **klasy dominującej według prawdopodobieństwa a priori**.
- ↪ W większości pakietów **domyślnie wynosi 5**.

Ile powinno wynosić k ?

- ↪ Powinno być **nieparzyste**, dla ewidentnego wyniku głosowania.
- ↪ Im **mniejsze** k , tym mniejsza próbka, a więc większy wpływ obserwacji odstających.
- ↪ Im **większe** k , tym bardziej upraszczamy klasyfikację do wyboru **klasy dominującej według prawdopodobieństwa a priori**.
- ↪ W większości pakietów **domyślnie wynosi 5**.
- ↪ Najlepszy byłby dobór eksperymentalny przez **zbiór walidujący**.

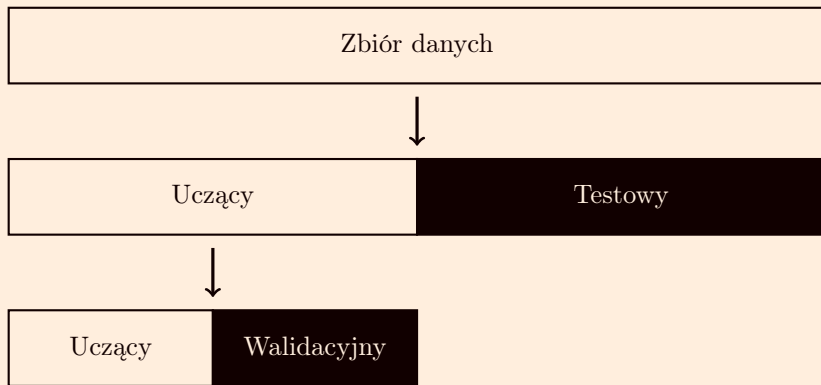
Zbiór danych

Zbiór danych



Uczący

Testowy



Wady i zalety kNN

Wady i zalety kNN

- 👍 Uczenie jest nieprzyzwoicie wręcz szybkie. (*bo go nie ma*)

Wady i zalety kNN

- 👍 Uczenie jest nieprzyzwoicie wręcz szybkie. *(bo go nie ma)*
- 👍 Nie wymaga żadnej matematyki do wyjaśnienia. *(o ile nie skupiamy się na definicji odległości)*

Wady i zalety kNN

- 👍 Uczenie jest nieprzyzwoicie wręcz szybkie. *(bo go nie ma)*
- 👍 Nie wymaga żadnej matematyki do wyjaśnienia. *(o ile nie skupiamy się na definicji odległości)*
- 👎 Jest szalenie nieefektywny pamięciowo. *(musi zapisać w pamięci wszystko)*

Wady i zalety kNN

- 👍 Uczenie jest nieprzyzwoicie wręcz szybkie. *(bo go nie ma)*
- 👍 Nie wymaga żadnej matematyki do wyjaśnienia. *(o ile nie skupiamy się na definicji odległości)*
- 👎 Jest szalenie nieefektywny pamięciowo. *(musi zapisać w pamięci wszystko)*
- 👎 Jest przygnębiająco nieefektywny obliczeniowo. *(im większe k i im większa liczba wzorców, tym dłużej szuka)*

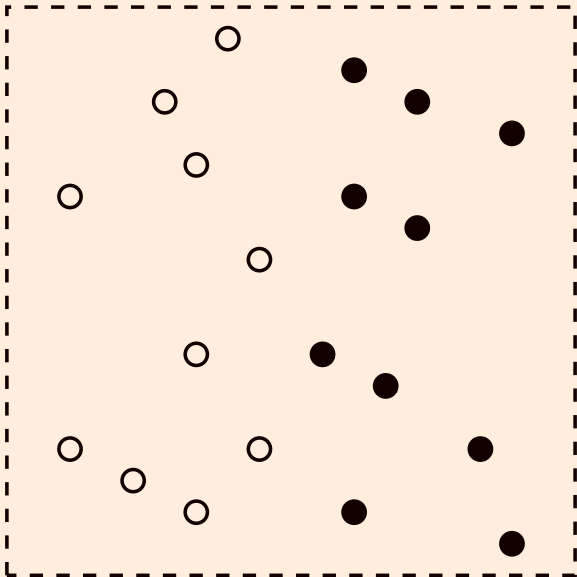
Wady i zalety kNN

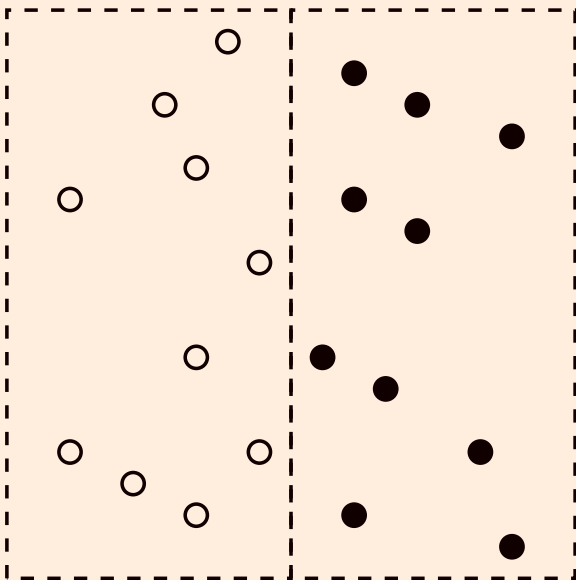
- 👍 Uczenie jest nieprzyzwoicie wręcz szybkie. *(bo go nie ma)*
- 👍 Nie wymaga żadnej matematyki do wyjaśnienia. *(o ile nie skupiamy się na definicji odległości)*
- 👎 Jest szalenie nieefektywny pamięciowo. *(musi zapisać w pamięci wszystko)*
- 👎 Jest przegnąbiająco nieefektywny obliczeniowo. *(im większe k i im większa liczba wzorców, tym dłużej szuka)*
- 👎 **NIE MÓWI NAM NIC O PROBLEMIE.** *(nie budujemy w nim żadnego modelu)*

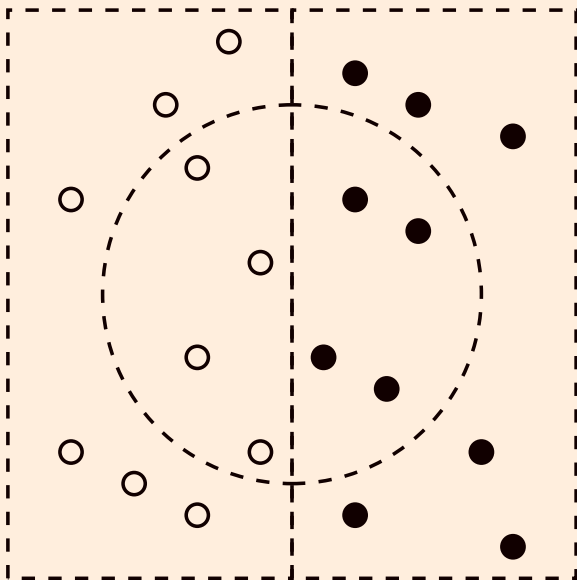


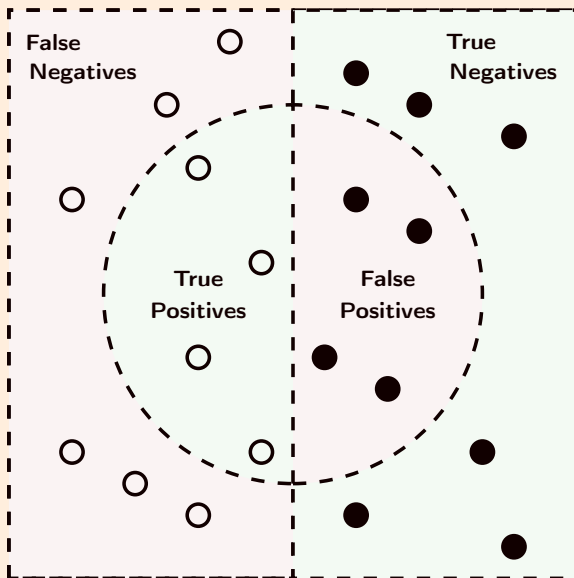
Miary oceny jakości











Macierz pomyłek dla problemu binarnego

	Wynik predykcji		
Etykieta	Prawda	Fałsz	Suma
Prawda	TP	FN	P
Fałsz	FP	TN	N
Suma	P'	N'	

Błąd i dokładność

$$ERR = \frac{FN + FP}{TP + FN + FP + TN}$$

$$ACC = 1 - ERR$$

Czy dokładność to dobra miara?

↪ Nie znaczy przesadnie wiele, jeśli mamy do czynienia z **problemem niezbalansowanym**.

Czy dokładność to dobra miara?

- ↪ Nie znaczy przesadnie wiele, jeśli mamy do czynienia z **problemem niezbalansowanym**.
- ↪ **Najczęściej mamy do czynienia z problemem niezbalansowanym.**

Czy dokładność to dobra miara?

- ↪ Nie znaczy przesadnie wiele, jeśli mamy do czynienia z **problemem niezbalansowanym**.
- ↪ **Najczęściej mamy do czynienia z problemem niezbalansowanym.**



Jeśli klasyfikujemy bardzo rzadką chorobę (założmy, że występuje u jednego człowieka na milion) i ograniczymy nasz klasyfikator do informowania każdego pacjenta o tym że jest zdrowy, niezależnie od wyników badań, otrzymamy dokładność powyżej 99%.

Czy dokładność to dobra miara?

- ↪ Nie znaczy przesadnie wiele, jeśli mamy do czynienia z **problemem niezbalansowanym**.
- ↪ **Najczęściej mamy do czynienia z problemem niezbalansowanym.**



Jeśli klasyfikujemy bardzo rzadką chorobę (założmy, że występuje u jednego człowieka na milion) i ograniczymy nasz klasyfikator do informowania każdego pacjenta o tym że jest zdrowy, niezależnie od wyników badań, otrzymamy dokładność powyżej 99%.



Wszyscy chorzy umrą w niewiedzy i mękach.

Precision i recall

Precision i recall

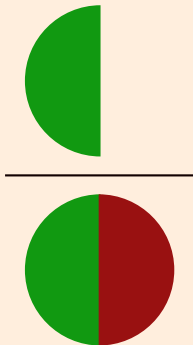
Precision

Ile z **ocenionych pozytywnie wzorców** zostało ocenionych słusznie.

Precision i recall

Precision

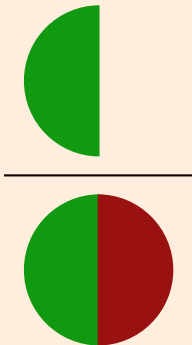
Ile z **ocenionych pozytywnie wzorców** zostało ocenionych słusznie.



Precision i recall

Precision

Ile z **ocenionych pozytywnie wzorców** zostało ocenionych słusznie.



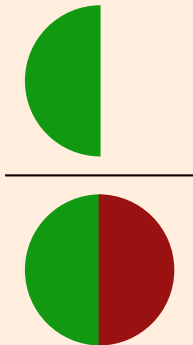
Recall

Ile z **pozytywnych wzorców** zostało ocenionych słusznie.

Precision i recall

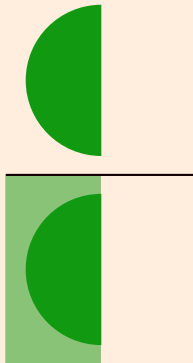
Precision

Ile z **ocenionych pozytywnie wzorców** zostało ocenionych słusznie.



Recall

Ile z **pozytywnych wzorców** zostało ocenionych słusznie.



F-score

$$F_1 = 2 * \frac{precision * recall}{precision + recall}$$