

# Algorytmy Ucznia Maszynowego

## Sprawozdanie

|                  |                   |             |  |
|------------------|-------------------|-------------|--|
| Data:            | 28.03.2023        | Dzień:      | Wtorek                                       |
| Grupa:           | grupa nr 2        | Godzina:    | 18:55  |
| Nazwisko i imię: | Rydyński Grzegorz | Prowadzący: | prof. dr hab. inż. Ewa Skubalska-Rafajłowicz |
| Nr indeksu:      | 252958            | Temat:      | Sprawozdanie z programów                     |

## Spis treści

|          |  |          |
|----------|--|----------|
| <b>1</b> | <b>Wstęp</b>   | <b>1</b> |
| <b>2</b> | <b>Algorytm k najbliższych sąsiadów</b>  | <b>2</b> |
| 2.1      | Działanie algorytmu - obliczanie odległości . . . . .                          | 2        |
| <b>3</b> | <b>Testy</b>   | <b>3</b> |
| 3.1      | Seria 1 - zmiana ilości sąsiadów dla stałego zbioru danych testowych . . . . . | 3        |

## 1. Wstęp

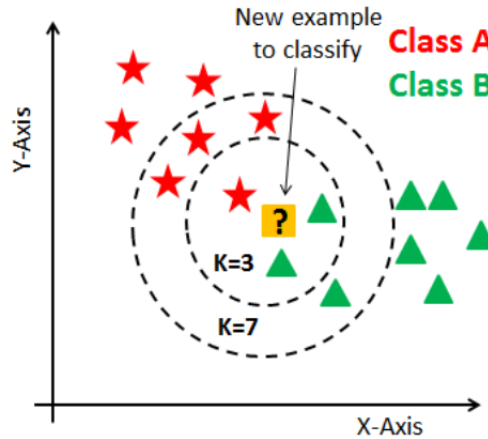
Na zajęciach projektowych należało zrealizować 3 programy realizujące zadanie klasyfikacji danych z wybranego zbioru. Należało zrealizować następujące algorytmy:

- algorytm k najbliższych sąsiadów,
- algorytm perceptronu Rosenblatta,
- wybrany algorytm klasyfikacji.

Pierwsze dwa algorytmy zostały zrealizowane na zbiorze danych Iris.

## 2. Algorytm k najbliższych sąsiadów

Algorytm k najbliższych sąsiadów jest jednym z nieparametrycznych algorytmów, wykorzystywanym do zadań regresji oraz klasyfikacji. Algorytm oblicza odległość danego punktu danych ze zbioru testowego do wszystkich pozostałych punktów danych ze zbioru treningowego. Po obliczeniu wszystkich odległości, wybierane jest k najbliższych punktów treningowych (nazywanych sąsiadami) do punktu testowego. Następnie, na podstawie ilości wystąpień danej klasy, do punktu testowego przypisywana jest klasa, która najczęściej wystąpiła wśród sąsiadów.



Rys. 1: Wizualizacja algorytmu odpowiednio dla  $k = 3$  i  $k = 7$

### 2.1. Działanie algorytmu - obliczanie odległości

Obliczanie odległości pomiędzy punktami wykonywane jest na podstawie różnych metryk długości. Najbardziej popularną metryką i metryką domyślnie wykorzystywaną przez biblioteki Scikit-Learn, (która została wykorzystana również w programie) jest metryka Euklidesowa, opisana wzorem:

$$d(p, q) = \sqrt{\sum (p_i - q_i)^2} \quad (1)$$

Do innych metryk można zaliczyć:

- metrykę Minkowskiego
- metrykę Czebyszewa
- metrykę "miejską" (metryka Manhattana)
- metrykę Minkowskiego

Z metryki Minkowskiego wynikają wszystkie powyższe metryki, ponieważ opisywana jest następująco:

$$d(p, q) = (\sum |p_i - q_i|^p)^{1/p} \quad (2)$$

Dla parametru  $p = 2$  otrzymujemy metrykę Euklidesową, a dla  $p = 1$  otrzymujemy metrykę Manhattana, którą można zapisać jako:

$$d(p, q) = \sum |p_i - q_i| \quad (3)$$

Dla parametru  $p$  dążącego do  $\infty$  otrzymujemy metrykę Czebyszewa:

$$d(p, q) = \lim_{p \rightarrow \infty} (\sum |p_i - q_i|^p)^{1/p} = \max |p_i - q_i| \quad (4)$$

analogicznie dla  $p$  dążącego do  $-\infty$  otrzymujemy:

$$d(p, q) = \lim_{p \rightarrow -\infty} (\sum |p_i - q_i|^p)^{1/p} = \min |p_i - q_i| \quad (5)$$

Implementacja funkcji wyliczającej dystans na podstawie metryki Euklidesowej wygląda następująco:

```
1  @staticmethod
2  def calculate_distance(test_row: list, train_row: list) -> float:
3      """ This method calculates the euclidean distance between data points
4      :param test_row: test_data point of type list
5      :param train_row: train_data point of type list
6      :return: euclidean distance between data points in the n-th dimension
7      """
8      euclidean_distance = 0.0
9      for i in range(len(train_row)-1):
10         euclidean_distance += (test_row[i] - train_row[i])**2
11     return sqrt(euclidean_distance)
```

Listing 1: Implementacja funkcji wyliczającej dystans w programie

### 3. Testy

Dla algorytmu zostały przeprowadzone 3 serie testów, w celu sprawdzenia jakości predykcji.

#### 3.1. Seria 1 - zmiana ilości sąsiadów dla stałego zbioru danych testowych

Pierwsza seria testów polega na zmianie ilości branych pod uwagę sąsiadów. Badany był zakres sąsiadów  $n\_neighbors = (3, 8)$