

Systemy Sztucznej Inteligencji

Predykcja cukrzycy u pacjentów z wykorzystaniem miękkiego KNN.

Marian Dorosz

Wojciech Brożek

Oskar Fojcik

Grupa 2/3

9 maja 2022

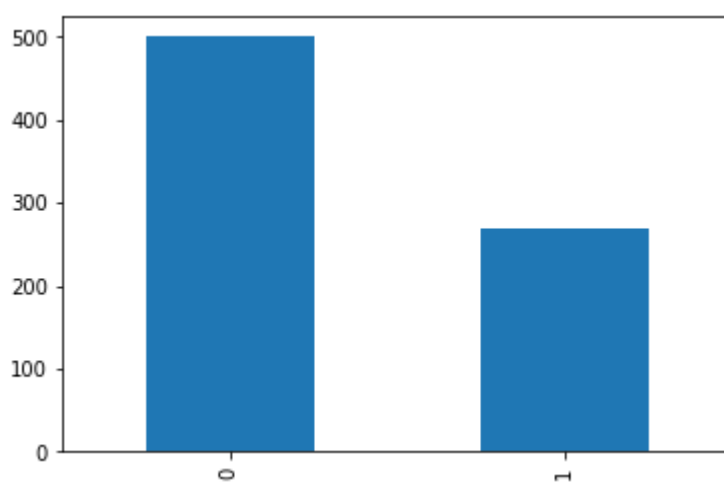
Część I

Opis programu

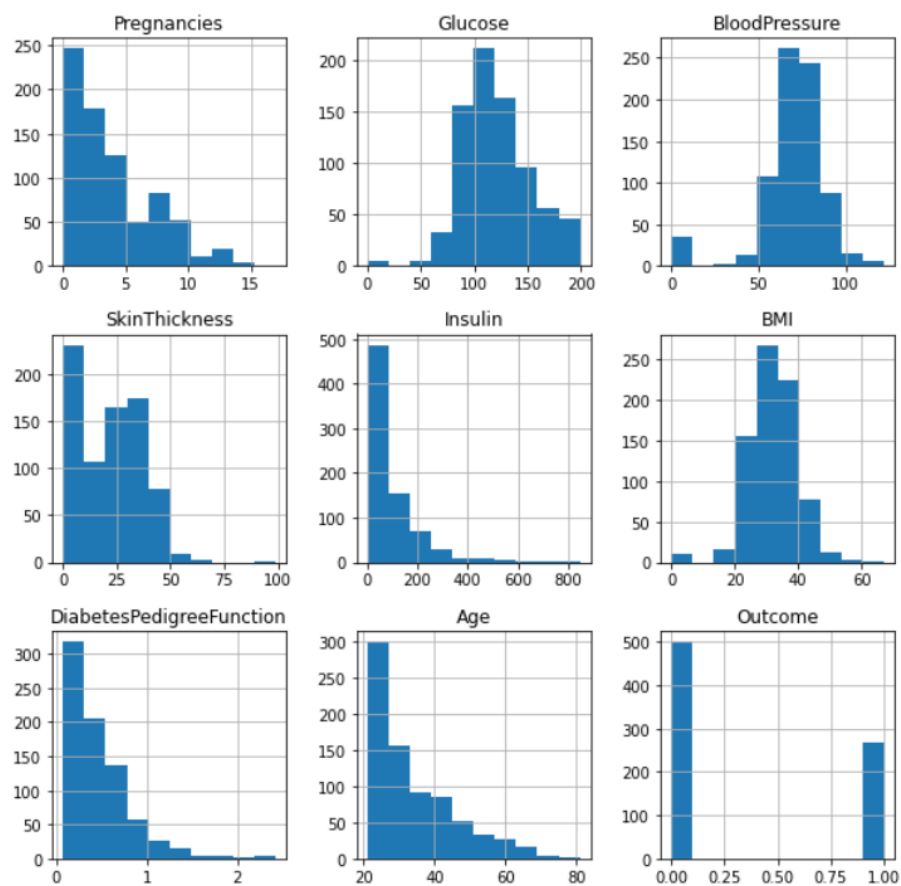
Program to implementacja prostego algorytmu uczenia maszynowego - KNN, który wykorzystując zbiory miękkie decyduje, czy dana osoba ma cukrzycę lub nie.

Część II

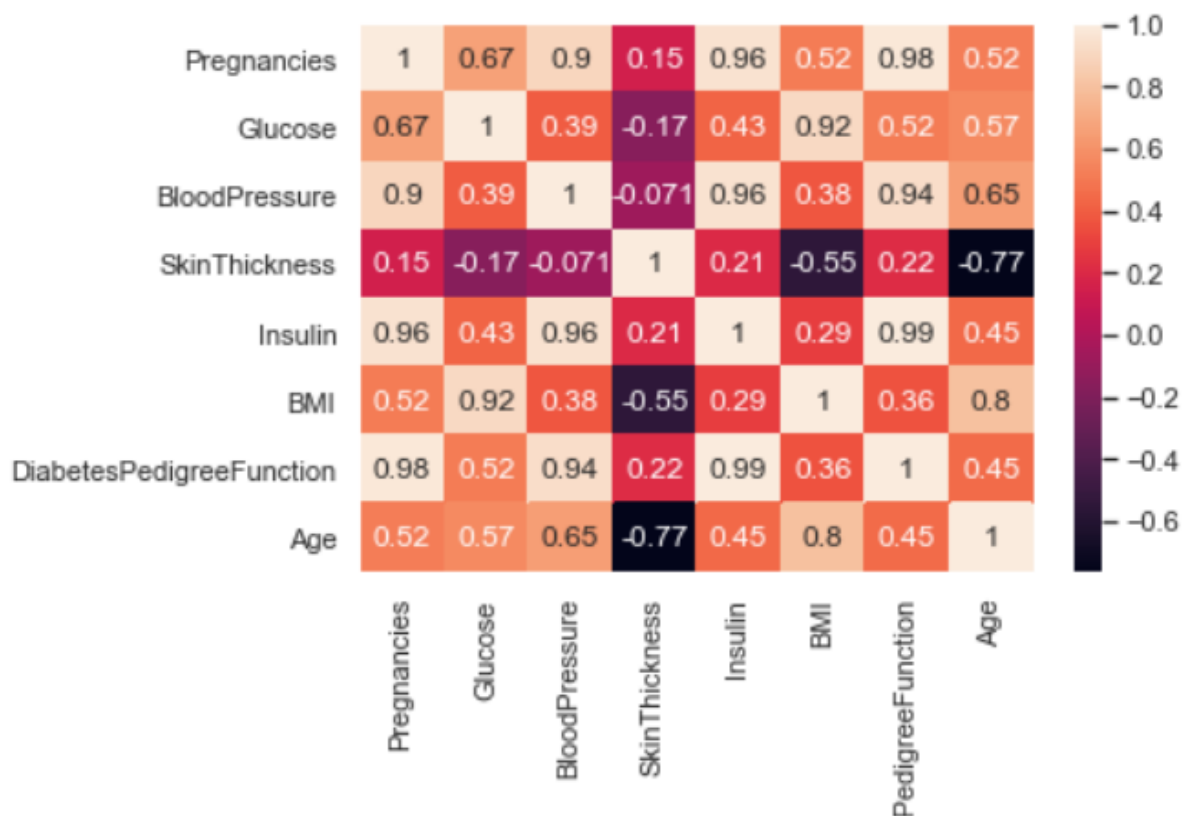
Analiza danych



Rysunek 1: Wykres przedstawiający rozkład osób zdrowych oraz chorych na cukrzycę; gdzie 0 to osoba zdrowa, 1 to osoba chora



Rysunek 2: Wykres przedstawiający jak rozdzielone są dane w poszczególnych klasach



Rysunek 3: Wykres przedstawiający korelacje pomiędzy poszczególnymi danymi im jaśniejsze pole tym większa korelacja

Wnioski:

- Istnieje wysoka korelacja pomiędzy funkcją wykrywającą cukrzycę a poziomem insuliny, glukozy oraz ciśnieniem krwi.
- Funkcja wykrywająca cukrzycę prawdopodobnie mylnie klasyfikuje kobiety w ciąży jako chore na cukrzycę.
- Osoby młode mają statystycznie większe BMI.
- Osoby o wysokim ciśnieniu krwi mają zwykle wysoki poziom glukozy.

Opis działania

Informacje ogólne

Cały program wykorzystuje trzy klasy statyczne, które odpowiadają za:

- Przetwarzanie danych
- Tworzenie zbioru miękkiego oraz obliczanie przynależności do klasy w oparciu o zadany zbiór miękkie
- Klasyfikację przy pomocy algorytmu KNN

Klasa `ProcessingData`

Ta klasa odpowiada za przetwarzanie danych. Posiada ona trzy funkcje:

- `splitSet(x,k)` - funkcja ta zwraca zbiór danych podzielony w określonej proporcji. Przyjmuje ona dwa argumenty - zbiór danych oraz w jakim stosunku ma go podzielić.

Data: Dane wejściowe zbiór danych x , liczba k

Result: $x[:n]$, $x[n:]$

$n = \text{int}(\text{len}(x)*k);$

return $x[:n]$, $x[n:]$;

Algorithm 1: Podział danych wejściowych na dwa zbiory.

- `shuffle(x)` - funkcja ta zwraca dane przetasowane przy pomocy algorytmu **Fisher-Yates shuffle**. Przyjmuje ona jako argument zbiór danych.

Data: Dane wejściowe zbiór danych x

Result: Przetasowany zbiór x

$i = \text{len}(x) - 1;$

while $i > 0$ **do**

$j = \text{random}(i, \text{len}(x));$

 swap($x[i]$, $x[j]$);

end

return x ;

Algorithm 2: Algorytm tasowania danych.

- `normalize(x)` - ta funkcja normalizuje zbiór danych wykorzystując algorytm normalizacji min-max. Wynikiem jej działania jest znormalizowany zbiór danych. Przyjmuje ona jeden argument i jest nim zbiór danych.

Data: Dane wejściowe zbiór danych x

Result: Znormalizowany zbiór danych x

Pobierz kolumny pomijając kolumnę 'Outcome' i zapisz je w $values$

Stwórz listę nazw kolumn zbioru $values$ i przypisz je do zmiennej $cols$

foreach col *in* $cols$ **do**

 Pobierz wszystkie wiersze dla danej kolumny i przypisz je do zmiennej $data$;

 Pobierz maksymalną wartość ze zmiennej $data$ i przypisz do zmiennej max ;

 Pobierz minimalną wartość ze zmiennej $data$ i przypisz do zmiennej min ;

for $(r = 0; r < len(x); r++)$ **do**

$x[r, col] = (x[r, col] - min) / (max - min)$;

end

end

return x

Algorithm 3: Algorytm normalizacji danych (min-max).

Klasa `soft_set`

Zadaniem tej klasy jest zbudowanie zbioru miękkiego dla zadanego wektora danych. W jej skład wchodzi następujące funkcje:

- `build_soft_set(x)` - ta funkcja służy do zbudowania zbioru miękkiego dla zadanej jako argument kolekcji danych. Jej wynikiem jest słownik z wartościami 0-1. Samo budowanie słownika odbywa się następująco: algorytm wylicza średnie wartości dla zadanej kolumny, a następnie sprawdza ile wartości w kolumnie jest poniżej i powyżej średniej dla osób zdrowych i niezdrowych. W oparciu o te dane uzupełnia słownik wartościami.

Data: Dane wejściowe zbiór danych x

Result: Zbiór miękkie w postaci słownika $soft_set$

Utwórz słownik $soft_set$;

Utwórz klucz w słowniku: $soft_set[0]$ i przypisz do niego nowy słownik;

Utwórz klucz w słowniku: $soft_set[1]$ i przypisz do niego nowy słownik;

foreach key **in** $soft_set$ **do**

$data1 = x$;

 Pobierz każdy wiersz, dla którego kolumna 'Outcome' == key ;

 Z $data1$ oraz $data2$ usuń kolumnę 'Outcome';

foreach key **in** $data1.columns$ **do**

$mean =$ średnia wartości całej kolumny $data1[col]$;

$temp1 = 0$;

$temp2 = 0$;

end

foreach val **in** $data2[col]$ **do**

if $val < mean$ **then**

$temp1++$;

else

$temp2++$;

end

end

if $temp1 < temp2$ **then**

$soft_set[key][col] = 0$;

else

$soft_set[key][col] = 1$;

end

end

return $soft_set$

Algorithm 4: Algorytm budujący zbiór miękkie.

```
{0: {'Pregnancies': 0, 'Glucose': 0, 'BloodPressure': 1, 'SkinThickness': 1, 'Insulin': 0, 'BMI': 0, 'DiabetesPedigreeFunction': 0, 'Age': 0}, 1:
{'Pregnancies': 1, 'Glucose': 1, 'BloodPressure': 1, 'SkinThickness': 1, 'Insulin': 0, 'BMI': 1, 'DiabetesPedigreeFunction': 0, 'Age': 1}}
```

Rysunek 4: Przykładowy słownik reprezentujący zbiór miękkie.

- `get_membership(x, v)` - ta funkcja sprawdza do jakiej klasy przynależy dany wektor. Przyjmuje ona jako argumenty zbiór miękki oraz wektor danych. Jej wynikiem jest liczba zmiennoprzecinkowa.

Data: Dane wejściowe zbiór miękki *soft_set*, wektor danych *v*

Result: Liczba zmiennoprzecinkowa

result = nowy słownik;

```
foreach key in soft_set do
    result[key] = 1;
    foreach key2 in soft_set[key2] do
        result[key] += soft_set[key][key2] * v[key2]
    end
end
```

end

Zwróć największą wartość przypisaną do jednego z kluczy słownika.;

Algorithm 5: Algorytm tasowania danych.

Klasa `soft_KNN`

Ta klasa to implementacja algorytmu KNN, który do głosowania wykorzystuje zbiory miękkie. Składa się ona z następujących funkcji:

- `calc_euclidian_distance(val1, val2)` - ta funkcja zwraca jednowymiarową odległość euklidesową dla dwóch wartości. Innymi słowy implementuje wzór $|val1 - val2|$, gdzie val_i to wartość.

Data: Dane wejściowe liczby *val1*, *val2*

Result: Liczba zmiennoprzecinkowa

`return abs(val1 - val2);`

Algorithm 6: Obliczanie metryki euklidesowej.

- `clustering(x, v, k)` - ta funkcja jest implementacją algorytmu KNN, z wykorzystaniem elementów klasy *soft_set*. Jej zadaniem jest ocena, czy wprowadzone dane należą do osoby zdrowej lub chorej. KNN jest leniwym algorytmem uczenia nadzorowanego, więc dane muszą mieć kolumnę, dzięki której jest możliwe jednoznaczne zidentyfikowanie klasy obiektu. Sam algorytm polega na wyliczeniu k sąsiadów, obok których znajduje się badany element. W tym celu wykorzystuje się metrykę (np.: Euklidesową), która określa jak blisko danego sąsiada znajduje się badany element. Klasa, do której zostanie przypisany element jest zależna od tego, jak wielu sąsiadów danego typu posiada. Jeżeli np.: większość sąsiadów badanego obiektu to cukrzycy, to będzie to oznaczać, że osoba jest najprawdopodobniej cukrzykiem.

Data: Dane wejściowe zbiór danych x , wektor danych $sample$, liczba całkowita k

Result: Nazwa klasy, do której należy próbka

$soft_x = soft_set.build_soft_set(x);$

$distances = [];$

for ($i = 0; i < len(x), i++$) **do**

$temp_x = soft_set.build_soft_set(soft_x, x[i]);$

$temp_sample = soft_set.build_soft_set(soft_x, sample);$

$distances.append(soft_KNN.calc_euclidian_distance(temp_x, temp_sample));$

end

$temp =$ kopia x ;

Do zmiennej $temp$ dodaj kolumnę 'distance' i przypisz jej tablicę $distances$;

Przesortuj tabelę $temp$ zgodnie z wartościami kolumny 'distances';

$classes = \{0: 0, 1: 0\};$

for ($i = 0; i < k; i++$) **do**

$classes[temp[i].Outcome] += 1;$

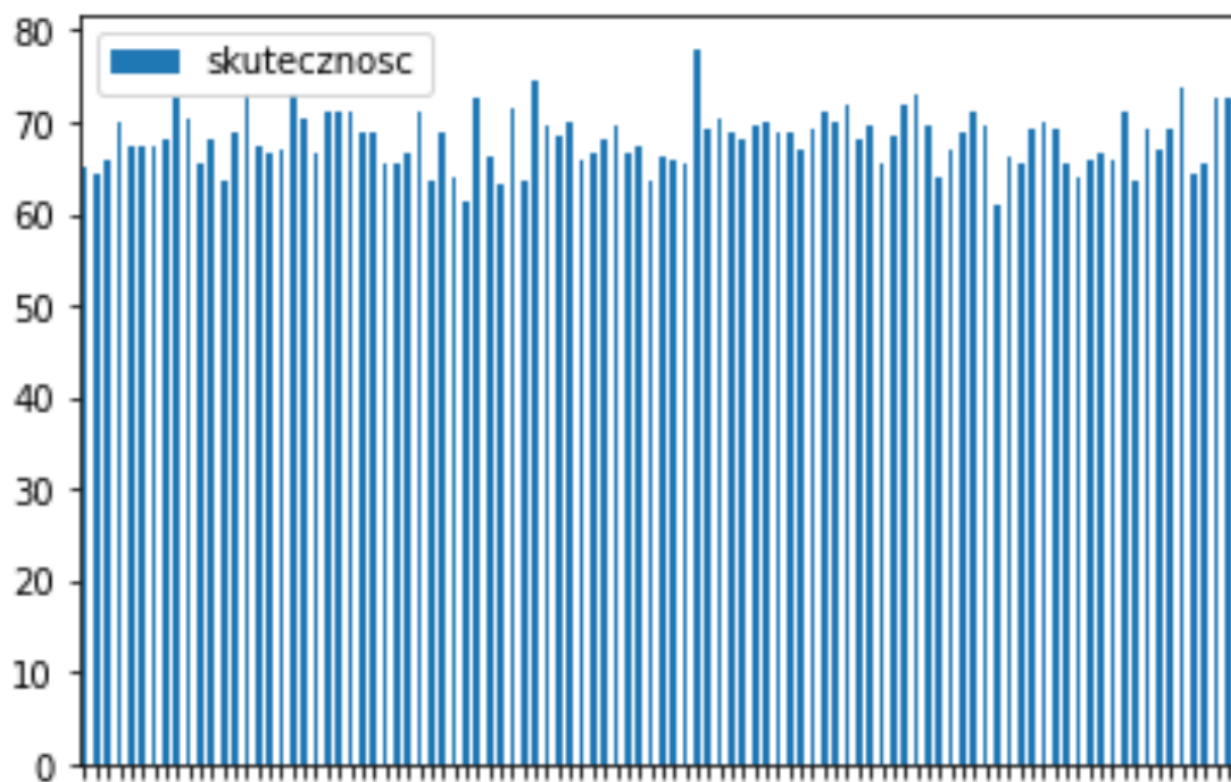
end

Zwróć klucz z największą przypisaną wartością;

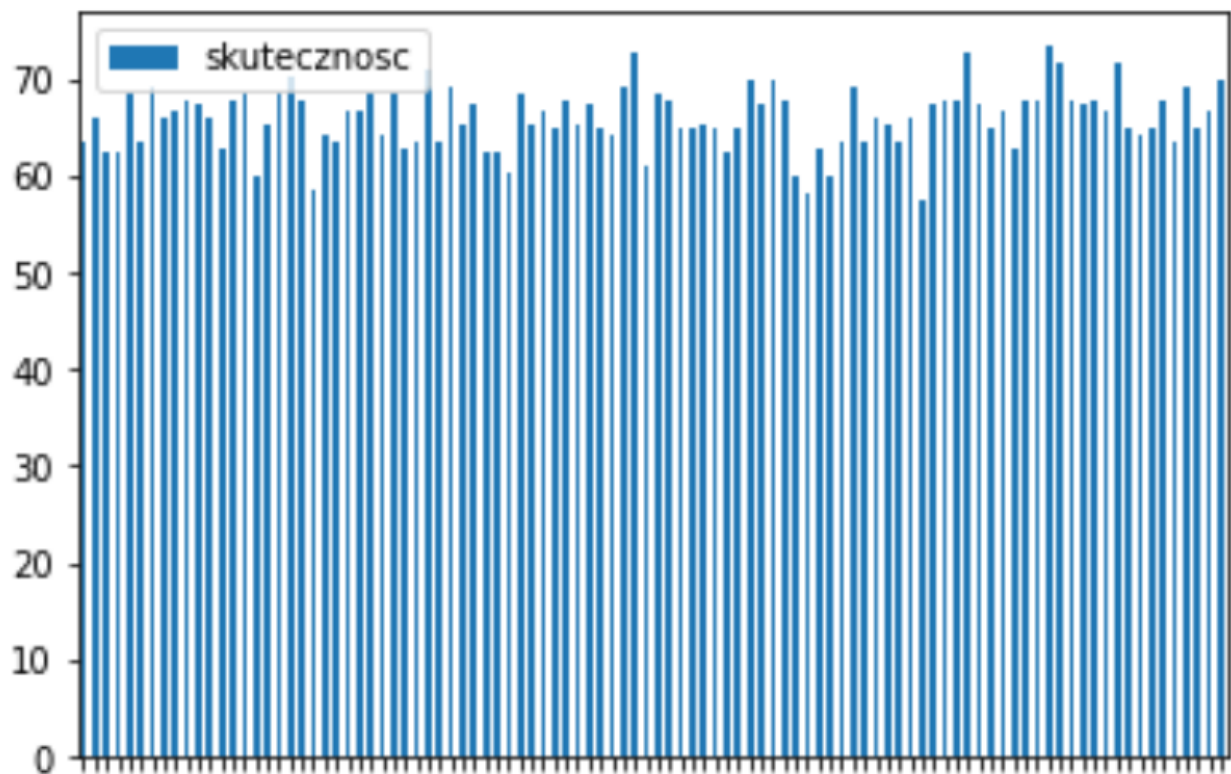
Algorithm 7: Klasteryzacja danych.

Testy

Przeprowadzone zostały testy dla większych i mniejszych zbiorów testowych przeasowanych i znormalizowanych. Większy zawierał 537 elementów a mniejszy 375.



Rysunek 5: Wykres skuteczności dla większego zbioru treningowego i walidacyjnego. Średnia skuteczność wynosi 68.19 %



Rysunek 6: Wykres skuteczności dla mniejszego zbioru treningowego i walidacyjnego.
Średnia skuteczność wynosi 65.98 %

Pełen kod aplikacji

```

1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4
5 class ProcessingData:
6     @staticmethod
7     def splitSet(x: pd.DataFrame, k: int) -> pd.DataFrame:
8         n = int(len(x)*k)
9         xTrain = x[:n]
10        xVal = x[n:]
11        return xTrain, xVal
12
13    @staticmethod
14    def shuffle(x: pd.DataFrame) -> pd.DataFrame:
15        for i in range(len(x)-1, -1, -1):
16            j = np.random.randint(i, len(x))
17            x.iloc[i], x.iloc[j] = x.iloc[j], x.iloc[i]
18        return x
19
20    @staticmethod
21    def normalize(x: pd.DataFrame) -> pd.DataFrame:

```

```

22     values = x.loc[:, x.columns != 'Outcome']
23     columnNames=values.columns.tolist()
24     for column in columnNames:
25         data = x.loc[:,column]
26         max1 = max(data)
27         min1 = min(data)
28         for row in range(0,len(x),1):
29             x.at[row,column] = (x.at[row,column]-min1)/(max1-min1)
30     return x
31
32 class soft_set:
33
34     @staticmethod
35     def build_soft_set(x: pd.DataFrame) -> dict:
36         soft_set = {}
37         soft_set[0] = {}
38         soft_set[1] = {}
39         for key in soft_set:
40             data1 = x
41             data2 = x.loc[x['Outcome']==key]
42             data1 = data1.drop(['Outcome'], axis=1)
43             data2 = data2.drop(['Outcome'], axis=1)
44             for col in data1.columns:
45                 mean = data1[col].mean()
46                 temp1 = 0 # lower or equal to mean
47                 temp2 = 0 # greater than mean
48                 for val in data2[col]:
49                     if val < mean:
50                         temp1 += 1
51                     else:
52                         temp2 += 1
53                 if temp1 > temp2:
54                     soft_set[key][col] = 0
55                 else:
56                     soft_set[key][col] = 1
57         return soft_set
58
59     @staticmethod
60     def get_membership(soft_set: dict, vector: pd.Series) -> float:
61         vector = vector.to_dict()
62         result = dict()
63         for key in soft_set:
64             result[key] = 1
65             for key2 in soft_set[key]:
66                 result[key] += soft_set[key][key2] * vector[key2]
67         return float(result[max(result, key=result.get)])
68
69 class soft_KNN:
70     @staticmethod
71     def calc_euclidian_distance(val1: float, val2: float) -> float:
72         return np.abs(val1-val2)
73
74     @staticmethod
75     def clustering(x: pd.DataFrame, sample: pd.Series, k: int) -> str:
76         soft_x = soft_set.build_soft_set(x)

```

```
77     distances = []
78     for i in range(i, len(x)):
79         temp_x = soft_set.get_membership(soft_x, x.iloc[i])
80         temp_sample = soft_set.get_membership(soft_x, sample)
81         distances.append(soft_KNN.calc_euclidian_distance(temp_x,
82                                                             temp_sample))
82     tempdf = x.copy()
83     tempdf['distance'] = distances
84     tempdf = tempdf.sort_values(by='distance')
85     classes = {
86         0: 0,
87         1: 0
88     }
89     for i in range(k):
90         classes[tempdf.iloc[i].Outcome] += 1
91
92     return max(classes, key = classes.get)
```
