

Systemy Sztucznej Inteligencji

dokumentacja projektu systemu przewidywania pogody

Hubert Bojda, Dawid Gala, grupa 4/7

17 czerwca 2024

Opis programu

System sztucznej inteligencji który klasyfikuje za pomocą klasyfikatora KNN dane pobrane z pogodowego API. Dane są odpowiednio przygotowane poprzez zredukowanie ilości klas abstrakcji do najliczniejszych czyli 'Foggy', 'Overcast', 'Clear'. Zdecydowaliśmy się na takie klasy, ponieważ są one dominujące w naszej bazie danych. Następnie dane zostają przetasowane oraz podzielone na zbiór danych testowych i treningowych. W następnym kroku dane są przetwarzane przez klasyfikator.

Instrukcja obsługi

Nasz program należy uruchomić z wykorzystaniem środowiska Jupyter. Nasz projekt został utworzony w programie PyCharm, który zawiera wbudowany serwer Jupyter.

Opis działania

Klasyfikator KNN, czyli algorytm k najbliższych sąsiadów, służy do klasyfikacji i prognozowania wartości na podstawie zmiennej określonej w kolumnie decyzyjnej w bazie danych. Algorytm porównuje wartości w kolumnach objaśniających dane zjawisko z wartościami zmiennych, które są zawarte w zbiorze uczącym. Zawiera on informacje o k najbliższych obserwacjach ze zbioru uczącego. Ważnym aspektem przy tworzeniu klasyfikatora jest dobór odpowiedniej metryki obliczającej dystans między obserwacjami zbioru uczącego a zbioru treningowego. Najpopularniejsze metryki to: euklidesowa, mińkowskiego lub manhattan. Dzięki kolejnym iteracjom podział danych jest poprawiany względem zadanej metryki. Algorytm przenosi dane pomiędzy klasami tak, aby wariancja wewnątrz każdej klasy była jak najmniejsza.

Wzory

1. Obliczanie odległości euklidesowej:

$$d(x_i, x_j) = \sqrt{\sum_{k=1}^n (x_{i,k} - x_{j,k})^2} \quad (1)$$

gdzie:

- x_i i x_j to dwa punkty w przestrzeni n -wymiarowej,
- $x_{i,k}$ i $x_{j,k}$ to k -te współrzędne tych punktów.

2. Wybór k najbliższych sąsiadów:

Po obliczeniu odległości euklidesowych dla wszystkich punktów w zbiorze treningowym, wybieramy k najbliższych sąsiadów, czyli tych z najmniejszymi wartościami odległości.

3. Głosowanie większościowe:

Klasa przypisywana nowemu punktowi jest najczęściej występującą klasą wśród k najbliższych sąsiadów. Głosowanie większościowe można opisać wzorem:

$$y = \arg \max_{c \in C} \sum_{i=1}^k I(y_i = c) \quad (2)$$

gdzie:

- y to przewidywana klasa dla nowego punktu,
- C to zbiór wszystkich możliwych klas,
- I to funkcja wskaźnikowa, która przyjmuje wartość 1, gdy $y_i = c$, w przeciwnym razie 0,
- y_i to klasa przypisana do i -tego najbliższego sąsiada.

Algorytm

Poniżej przedstawiony jest algorytm dzielenia zbioru danych na zbiór danych testowych i treningowych:

Algorithm 1: Pseudokod funkcji `my_train_test_split`

Require: X , y , $test_size$ (domyślnie 0.3), $random_state$ (domyślnie *None*)

- 1: **if** $random_state$ nie jest *None* **then**
- 2: Ustaw ziarno losowości na $random_state$
- 3: **end if**
- 4: $test_size = \text{int}(test_size * \text{len}(X))$
- 5: $shuffled_indices$ = permutacja indeksów X
- 6: $indices_train = shuffled_indices[test_size:]$
- 7: $indices_test = shuffled_indices[:test_size]$
- 8: $X_train = X.\text{iloc}[indices_train]$
- 9: $X_test = X.\text{iloc}[indices_test]$
- 10: $y_train = y.\text{iloc}[indices_train]$
- 11: $y_test = y.\text{iloc}[indices_test]$
- 12: **return** $X_train, X_test, y_train, y_test$

Poniżej przedstawione są algorytmy klasyfikatora KNN:

Algorithm 2: Opis algorytmu KNN

Require: $X_train, y_train, X_test, k_neighbors$

- 1: $predictions \leftarrow []$
- 2: **for** x **in** każdy punkt w X_test **do**
- 3: Oblicz odległości między punktem x a danymi treningowymi X_train .
- 4: Posortuj indeksy odległości w celu znalezienia k najbliższych sąsiadów.
- 5: Wybierz etykiety odpowiadające indeksom k najbliższych sąsiadów.
- 6: Znajdź najczęściej występującą etykietę (głosowanie większościowe).
- 7: Dodaj przewidywaną etykietę do listy przewidywań.
- 8: **end for**
- 9: **return** $predictions$

Algorithm 3: Funkcja obliczająca dokładność

Require: $X_train, y_train, X_test, y_test, k_neighbors$

- 1: $predictions \leftarrow$ wywołaj funkcję $predict(X_train, y_train, X_test, k_neighbors)$
- 2: Oblicz dokładność przewidywań poprzez porównanie przewidywanych etykiet z rzeczywistymi etykietami y_test .
- 3: **return** dokładność

Zbiór danych

Zbiór danych weatherHistory.csv pochodzi ze stony Kaggle.com.

Plik posiada 96453 rekordów i 12 kolumn. Jej schemat jest przedstawiony poniżej.

Tabela 1: Dane pogodowe z pliku CSV

<i>FormattedDate</i>	<i>Summary</i>	<i>PrecipType</i>	<i>Temperature(C)</i>
2006-04-18 07:00:00.000 +0200	Clear	rain	8.6889
2006-04-12 04:00:00.000 +0200	Overcast	rain	6.6222
2006-04-11 19:00:00.000 +0200	Foggy	rain	8.8

Tabela 2: Dane pogodowe z pliku CSV

<i>ApparentTemperature(C)</i>	<i>Humidity</i>	<i>WindSpeed(km/h)</i>	<i>WindBearing(degrees)</i>
8.6889	0.93	1.4329	290.0
2.6167	0.93	25.0355	0.0
5.2944	0.99	26.5006	339.0

Tabela 3: Dane pogodowe z pliku CSV

<i>Visibility(km)</i>	<i>Pressure(millibars)</i>	<i>DailySummary</i>
5.8443	1012.96	Partly cloudy until night.
6.118	1003.68	Foggy overnight and breezy in the morning.
2.6565	1004.99	Foggy in the evening.

Przygotowanie danych

Przed rozpoczęciem tasowania danych i ich normalizacji trzeba było odpowiednio dopasować zbiór danych do klasyfikatora. W pierwszej kolejności zostały usunięte kolumny z wartościami typu String, tego typu kolumn nie da się przetworzyć w sposób matematyczny. Jedynie została kolumna *Summary*, która jest naszą kolumną decyzyjną

Listing 1: Usuwanie kolumn

```
1 del df['Formatted Date']
2 del df['Apparent Temperature (C)']
3 del df['Precip Type']
4 del df['Loud Cover']
5 del df['Daily Summary']
6 del df['Humidity']
7 del df['Wind Bearing (degrees)']
```

Aby zredukować ilość klas abstrakcji zamieniliśmy wszystkie rekordy z podsumowaniem wskazującym na "overcast, clear, foggy", jednak z lekko innym opisem. Dzięki temu udało się zredukować dane do trzech najbardziej według nas znaczących klas abstrakcji.

Usuwanie rekordów w których występują opisy które nie pasują do klas abstrakcji.

Listing 2: Usuwanie rekordów

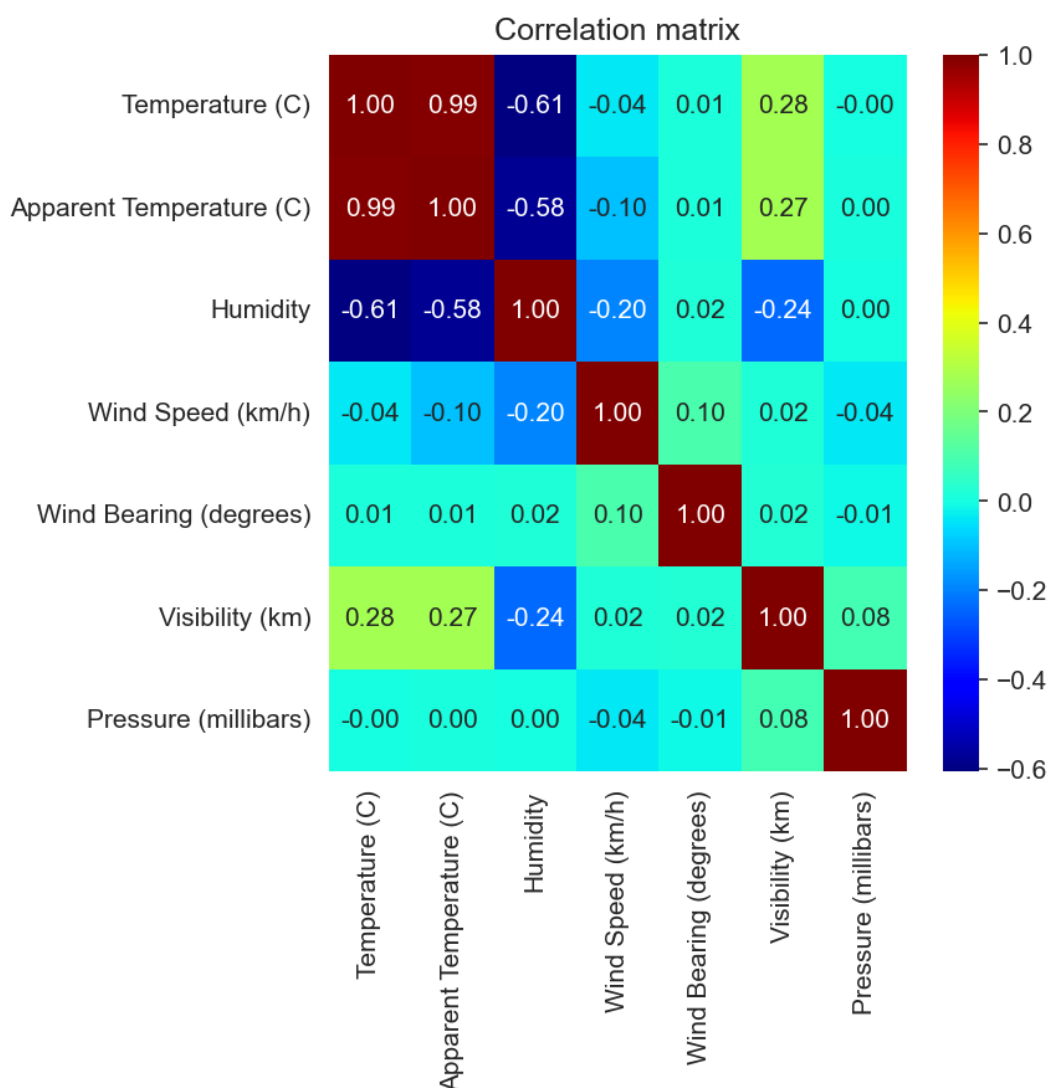
```
1 df = df.drop(df[df['Summary'] == 'Breezy and Mostly Cloudy'].index)
2 df = df.drop(df[df['Summary'] == 'Breezy and Partly Cloudy'].index)
3 df = df.drop(df[df['Summary'] == 'Humid and Mostly Cloudy'].index)
4 df = df.drop(df[df['Summary'] == 'Breezy and Overcast'].index)
5 df = df.drop(df[df['Summary'] == 'Humid and Partly Cloudy'].index)
6 df = df.drop(df[df['Summary'] == 'Windy and Foggy'].index)
7 df = df.drop(df[df['Summary'] == 'Windy and Overcast'].index)
```

Reszta usuniętych rekordów znajduje się w sekcji z pełnym kodem.

Dzięki tym zmianom, w naszej bazie zostały najważniejsze dla nas dane, a sama baza skurczyła się do 34635 wierszy, co znacznie skróci nasze czasy obliczeń.

Macierz korelacji

Macierz korelacji - jest to macierz która reprezentuje korelacje pomiędzy danymi w naszym zbiorze danych. Wartości w każdej komórce pochodzą z przedziału $<-1,1>$. Jeśli wartości są bliskie 1 to oznacza że wartości w kolumnach są silnie skorelowane dodatnio np. gdy jedna zmienna wzrasta to druga też. Natomiast gdy wartości są bliskie -1 to oznacza że wartości są silnie skorelowane ujemnie np. gdy jedna zmienna maleje to druga rośnie. Gdy wartość jest bliska 0 to zmienne nie są w ogóle skorelowane.



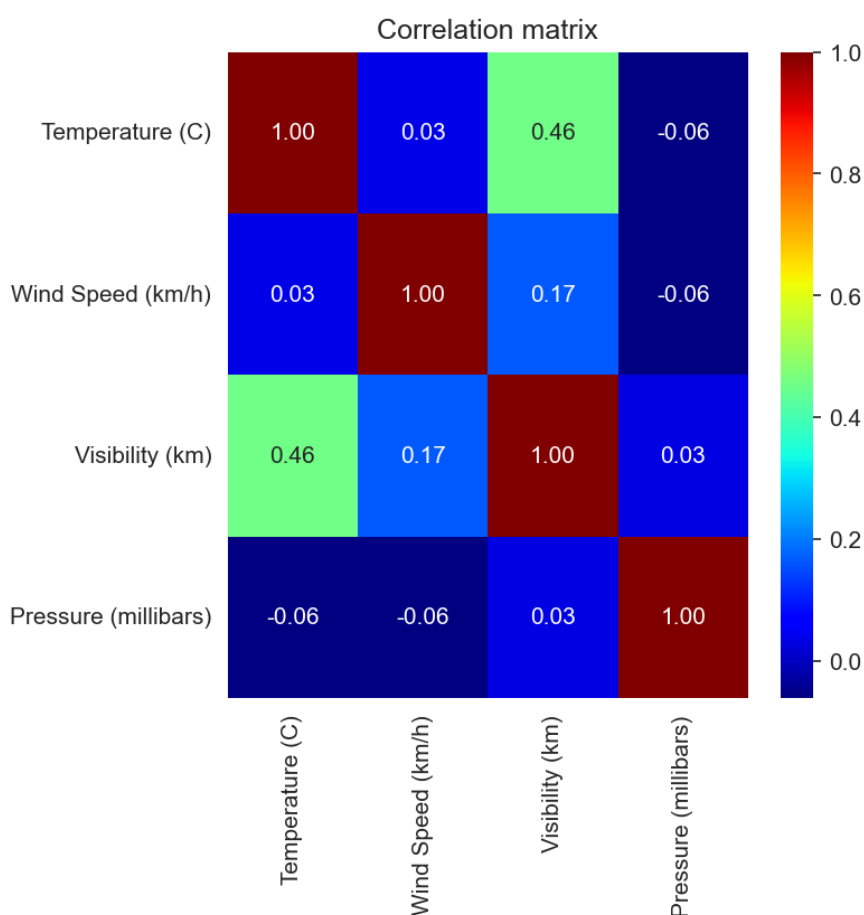
Rysunek 1: Macierz korelacji pełnego zbioru danych

Analiza macierzy korelacji

Powyższa macierz przedstawia zależność pomiędzy zmiennymi dla całej bazy danych liczącej 96453 rekordów.

- Kolumna $Temperature(C)$ oraz $ApparentTemperature(C)$ są silnie dodatnio skorelowane. Nie jest to zaskoczeniem ponieważ kiedy rzeczywista temperatura rośnie lub maleje to odczuwalna również.
- Ciekawym zjawiskiem jest zestawienie kolumn $Temperature(C)$ i $ApparentTemperature(C)$ wraz z kolumną $Humidity$. Na *Rysunek1* wartość korelacji wynosi odpowiednio -0.61 i -0.58. Co oznacza, że gdy temperatura rośnie to wilgotność spada.
- Podsumowując, analizę można w matematyczny sposób przedstawić wpływ różnych zmiennych w naszym obserwowanym zjawisku. Tego typu wizualizacje pomagają wyobrazić i zrozumieć jak czynniki mają na siebie wpływ w obserwowanym doświadczeniu a które nie.

Poniżej znajduje się macierz korelacji dla zredukowanego zbioru danych, który zawiera 34635 rekordów.



Rysunek 2: Macierz korelacji zbioru danych liczący 34635 rekordów

Implementacja

Projekt składa się z jednego pliku Jupyter o nazwie *projekt.ipynb*. Środowisko Jupyter jest popularne dla analityków danych i inżynierów wykorzystujących uczenie maszynowe. Jupyter pozwala na wykonywanie kodu krok po kroku poprzez podział na bloki. Ta funkcjonalność umożliwia analizowanie kodu i jego prototypowanie. Wspiera biblioteki typu Matplotlib lub Seaborn które pozwalają na interaktywną wizualizację danych.

Tasowanie i normalizacja danych

Listing 3: Podział danych na zbiór testowy i treningowy

```
1 y = df.pop("Summary")
2 X = df
3
4 def my_train_test_split(X, y, test_size=0.3, random_state=None):
5
6     if random_state is not None:
7         np.random.seed(random_state)
8
9     test_size = int(test_size * len(X))
10
11     shuffled_indices = np.random.permutation(X.index)
12     indices_train = shuffled_indices[test_size:]
13     indices_test = shuffled_indices[:test_size]
14
15     X_train = X.iloc[indices_train]
16     X_test = X.iloc[indices_test]
17     y_train = y.iloc[indices_train]
18     y_test = y.iloc[indices_test]
19
20     return X_train, X_test, y_train, y_test
21
22 X_train, X_test, y_train, y_test = my_train_test_split(X, y, test_size
    =0.3)
```

Powyższy kod odpowiada za podział zbioru danych na zbiory testowe i treningowe

- Na początku zbiór danych jest dzielony na X - cechy i y - etykiety
- Funkcja *my_train_test_split* przyjmuje jako argument zbiór cech, etykiet oraz opcjonalnie parametr rozmiaru zbioru testowego i ziarno losowości.
- Zwracane są 2 zbiory cech i 2 zbiory etykiet określone jako testowe i treningowe.

Listing 4: Klasyfikator KNN

```
1 class KNN:
2     def __init__(self, k_neighbors=3):
3         self.k_neighbors = k_neighbors
4
5     def __distance_metric(self, point, data):
6         return np.sqrt(np.sum((data - point) ** 2, axis=1))
7
8     def fit(self, X_train, y_train):
9         self.X_train = X_train
10        self.y_train = y_train
11
12    def predict(self, X_test):
13        predictions = []
14
15        for x in X_test.values:
16
17            distances = self.__distance_metric(x, self.X_train)
18
19            indices = np.argsort(distances)[:self.k_neighbors]
20
21            labels = [self.y_train.iloc[i] for i in indices]
22            result = max(set(labels), key=labels.count)
23
24            predictions.append(result)
25
26        return predictions
27
28    def score(self, X_test, y_test):
29        y_pred = self.predict(X_test)
30        accuracy = sum(y_test == y_pred) / len(y_test)
31        return accuracy
```

Powyższy kod przedstawia implementację klasyfikatora KNN

- Klasyfikator jest zaimplementowany jako klasa.
- Konstruktor pobiera parametr K najbliższych sąsiadów, domyślnie przypisanych jest 3 sąsiadów.
- Metoda *distance_metric* za pomocą metryki euklidesowej wyznacza odległość pomiędzy wektorami zmiennych.
- Metoda *fit* pobiera treningowy zbiór cech i etykiet i przypisuje do pól w klasie KNN.
- Metoda *predict* jest odpowiedzialna za wyznaczenie odległości pomiędzy danymi zbioru *X_test* i *X_train* i odpowiednie przewidywanie danych. Metoda zwraca etykietę, która najczęściej występowała wśród najbliższych sąsiadów.
- Metoda *score* pobiera zbiory testowe cech i etykiet a następnie porównuje przewidywane etykiety z tymi rzeczywistymi. Metoda zwraca stosunek poprawnie sklasyfikowanych punktów.

Testy

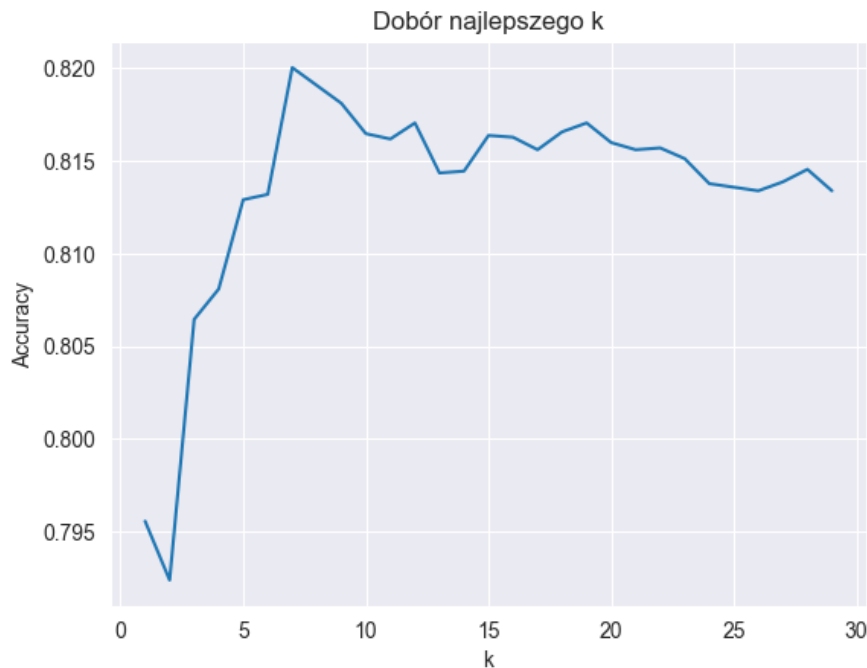
Pierwsza faza testów polegała na odpowiednim zredukowaniu ilości klas w zbiorze danych na potrzeby projektu. Nasza pierwotna baza zawiera sporą ilość rekordów oraz klas co powoduje wysoką złożoność obliczeniową modelu.

Po przygotowaniu zbioru danych można było zacząć testy modelu.

Listing 5: Testowanie

```
1 list_of_accuracies = []
2 for i in range(1,30):
3     my_knn = KNN(i)
4     my_knn.fit(X_train, y_train)
5     list_of_accuracies.append(my_knn.score(X_test, y_test))
```

Powyższy kod jest odpowiedzialny za utworzenie tablicy w której będą umieszczane wyniki dokładności naszego modelu w zależności od liczby k sąsiadów. W powyższym kodzie pobierzemy wyniki dokładności z zakresu $[1, 30]$. Dzięki wykorzystanym danym możemy zbadać dla jakiego k nasz model radzi sobie najlepiej. Z pomocą posłuży nam biblioteka Matplotlib, która rysuje wykresy.



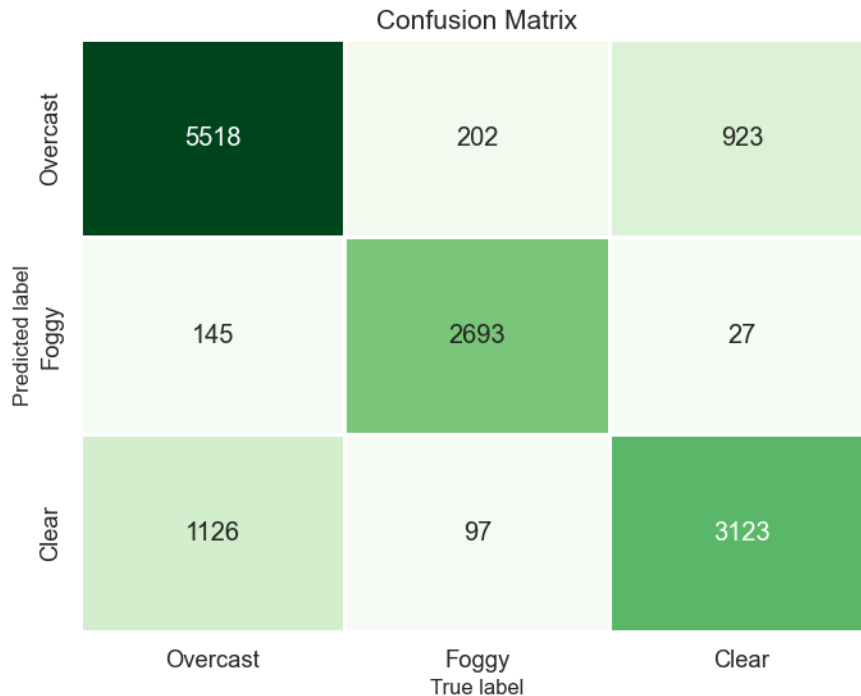
Rysunek 3: Wykres, który przedstawia zależność k sąsiadów od accuracy

Z powyższego wykresu wynika, że nasz model radzi sobie najlepiej z $k = 7$. Jednakże w przedziale $[1, 10]$ wartości accuracy znacząco rosną i dopiero w przedziale $(10, 30]$ wykres stabilizuje się.

Natomiast wszystkie wartości dokładności oscylują w przedziale $(0.792, 0.82]$. Pokazuje to, że powyższy model oraz baza danych dobrze się zachowuje podczas klasyfikacji.

Tablica pomyłek

Tablica pomyłek jest to wizualizacja stosowana przy ocenie jakości klasyfikacji. Wiersze tablicy odpowiadają poprawnym klasom decyzyjnym a kolumny przedstawiają przewidywania klasyfikatora. Poniżej znajduje się tablica dla naszego modelu którego k wynosi 13.



Rysunek 4: Tablica pomyłek dla $k = 13$

Podsumowanie

Analizując powyższe wykresy, można wysunąć kilka wniosków na temat klasyfikatora KNN. Widać, że algorytm ten doskonale radzi sobie z danymi numerycznymi, co znajduje potwierdzenie w przypadku bazy pogodowej. Wyniki dokładności oscylują w przedziale między 0.792 a 0.82, co jasno wskazuje na skuteczność działania tego klasyfikatora. Jego zdolność do dokładnego przewidywania warunków pogodowych na podstawie historycznych danych jest zatem bardzo obiecująca.

Pełen kod aplikacji

Listing 6: Pełen kod aplikacji

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 from sklearn.metrics import confusion_matrix
6
7 df = pd.read_csv('weatherHistoryWithRoundValues.csv')
8
9 del df['Formatted Date']
10 del df['Apparent Temperature (C)']
11 del df['Precip Type']
12 del df['Loud Cover']
13 del df['Daily Summary']
14 del df['Humidity']
15 del df['Wind Bearing (degrees)']
16
17 df = df.drop(df[df['Summary'] == 'Breezy and Mostly Cloudy'].index)
18 df = df.drop(df[df['Summary'] == 'Breezy and Partly Cloudy'].index)
19 df = df.drop(df[df['Summary'] == 'Humid and Mostly Cloudy'].index)
20 df = df.drop(df[df['Summary'] == 'Breezy and Overcast'].index)
21 df = df.drop(df[df['Summary'] == 'Humid and Partly Cloudy'].index)
22 df = df.drop(df[df['Summary'] == 'Windy and Foggy'].index)
23 df = df.drop(df[df['Summary'] == 'Windy and Overcast'].index)
24 df = df.drop(df[df['Summary'] == 'Breezy and Foggy'].index)
25 df = df.drop(df[df['Summary'] == 'Windy and Partly Cloudy'].index)
26 df = df.drop(df[df['Summary'] == 'Breezy and Dry'].index)
27 df = df.drop(df[df['Summary'] == 'Dry and Mostly Cloudy'].index)
28 df = df.drop(df[df['Summary'] == 'Windy and Dry'].index)
29 df = df.drop(df[df['Summary'] == 'Humid and Overcast'].index)
30 df = df.drop(df[df['Summary'] == 'Dangerously Windy and Partly Cloudy'].
    index)
31 df = df.drop(df[df['Summary'] == 'Windy and Mostly Cloudy'].index)
32 df = df.drop(df[df['Summary'] == 'Dry and Partly Cloudy'].index)
33 df = df.drop(df[df['Summary'] == 'Windy'].index)
34 df = df.drop(df[df['Summary'] == 'Breezy'].index)
35 df = df.drop(df[df['Summary'] == 'Rain'].index)
36 df = df.drop(df[df['Summary'] == 'Partly Cloudy'].index)
37 df = df.drop(df[df['Summary'] == 'Mostly Cloudy'].index)
38 df = df.drop(df[df['Summary'] == 'Dry'].index)
39 df = df.drop(df[df['Summary'] == 'Light Rain'].index)
40 df = df.drop(df[df['Summary'] == 'Drizzle'].index)
41
42 df = df.reset_index(drop=True)
43
44 listofvalues = df['Summary'].unique().tolist()
45 for value in listofvalues:
46     count = df['Summary'].str.contains(value).sum()
47     print(value, " występuje: ", count)
48
49 plt.figure(figsize=(5, 5), dpi=150)
```

```

50 sns.heatmap(df.drop("Summary", axis=1).corr(), annot=True, cmap="jet",
    fmt=".2f")
51 plt.title("Correlation matrix")
52 plt.show()
53
54 y = df.pop("Summary")
55 X = df
56
57
58 def my_train_test_split(X, y, test_size=0.3, random_state=None):
59     if random_state is not None:
60         np.random.seed(random_state)
61
62     test_size = int(test_size * len(X))
63
64     shuffled_indices = np.random.permutation(X.index)
65     indices_train = shuffled_indices[test_size:]
66     indices_test = shuffled_indices[:test_size]
67
68     X_train = X.iloc[indices_train]
69     X_test = X.iloc[indices_test]
70     y_train = y.iloc[indices_train]
71     y_test = y.iloc[indices_test]
72
73     return X_train, X_test, y_train, y_test
74
75
76 X_train, X_test, y_train, y_test = my_train_test_split(X, y, test_size
    =0.3)
77
78 class KNN:
79     def __init__(self, k_neighbors=3):
80         self.k_neighbors = k_neighbors
81
82     def __distance_metric(self, point, data):
83         return np.sqrt(np.sum((data - point) ** 2, axis=1))
84
85     def fit(self, X_train, y_train):
86         self.X_train = X_train
87         self.y_train = y_train
88
89     def predict(self, X_test):
90         predictions = []
91
92         for x in X_test.values:
93             distances = self.__distance_metric(x, self.X_train)
94
95             indices = np.argsort(distances)[:self.k_neighbors]
96
97             labels = [self.y_train.iloc[i] for i in indices]
98             result = max(set(labels), key=labels.count)
99
100             predictions.append(result)
101
102     return predictions

```

```

103
104     def score(self, X_test, y_test):
105         y_pred = self.predict(X_test)
106         accuracy = sum(y_test == y_pred) / len(y_test)
107         return accuracy
108
109
110 list_of_accuracies = []
111 for i in range(1, 30):
112     my_knn = KNN(i)
113     my_knn.fit(X_train, y_train)
114     list_of_accuracies.append(my_knn.score(X_test, y_test))
115
116 x_data = np.arange(1, 30, 1)
117 plt.plot(x_data, list_of_accuracies)
118 plt.xlabel('k')
119 plt.ylabel('Accuracy')
120 plt.title('Choosing the best k')
121 plt.grid(True)
122 plt.show()
123
124 my_knn = KNN(13)
125 knn_predicitons = my_knn.predict(X_test)
126 cm = confusion_matrix(y_test, knn_predicitons, labels=listofvalues)
127 cm_df = pd.DataFrame(cm, index=listofvalues, columns=listofvalues)
128 plt.figure(figsize=(8, 6))
129 sns.heatmap(cm_df, annot=True, cmap='Greens', fmt='d', linewidths=1,
130             linecolor='white', cbar=False)
131 plt.title('Confusion Matrix', fontsize=15)
132 plt.ylabel('Predicted label', fontsize=12)
133 plt.xlabel('True label', fontsize=12)
134 plt.show()

```
