

Systemy Sztucznej Inteligencji

Dokumentacja Projektu

Porównanie algorytmu KNN oraz Naiwnego Klasyfikatora Bayesa przy klasyfikacji odręcznie
pisanym cyfr.

Piotr Skowroński gr. 3/6

Krzysztof Czuba gr. 4/7

Jakub Poreda gr. 3/6

27 czerwca 2023

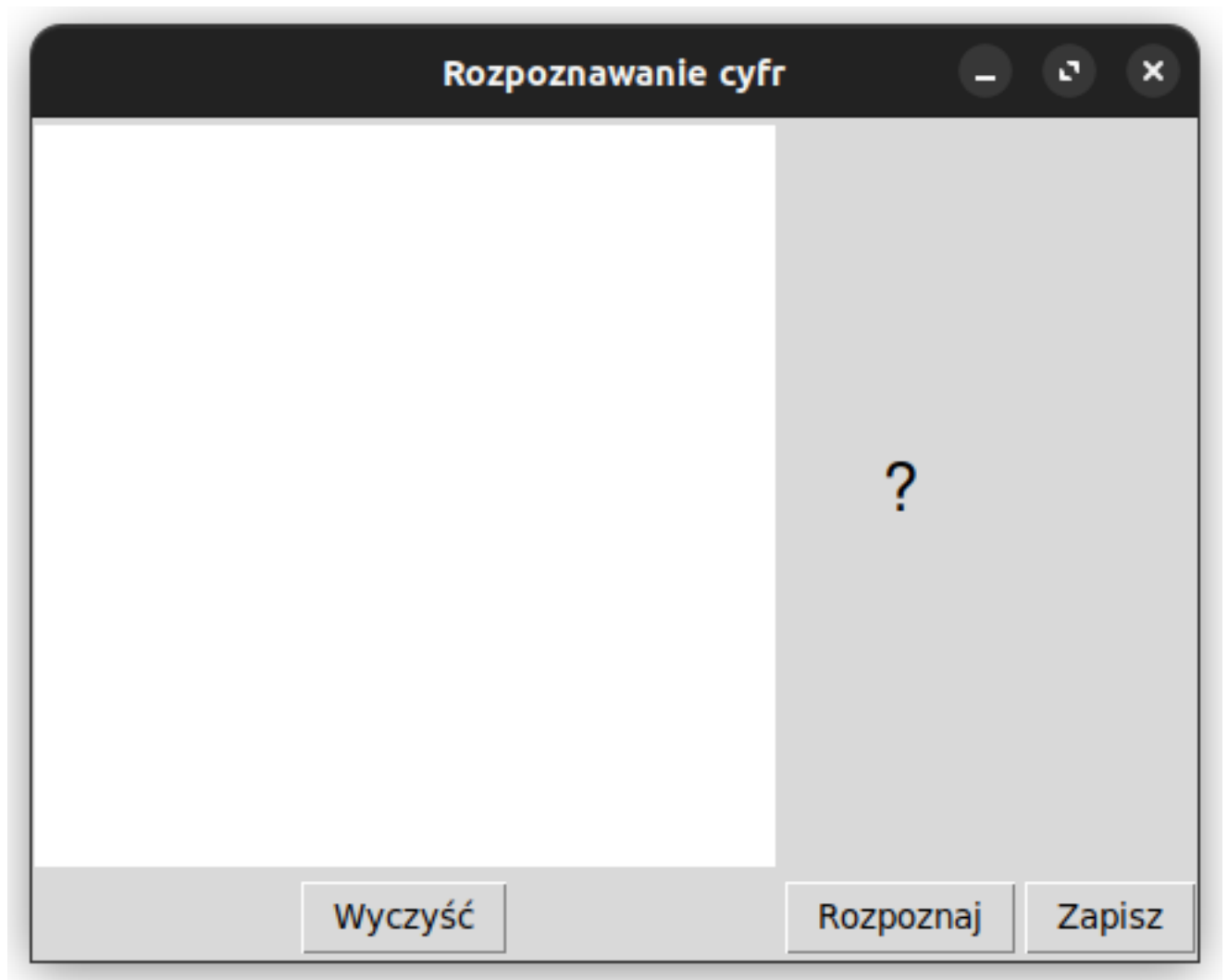
Spis treści

1	Wstęp	3
1.1	Opis programu	3
1.2	Użyte biblioteki	4
1.3	Baza danych	4
2	Opis działania	6
2.1	Normalizacja danych	6
2.2	Algorytm k najbliższych sąsiadów	6
2.3	Metryka odległości	7
2.4	Pseudokod algorytmu kNN	7
2.5	Wyniki klasyfikatora kNN	7
2.6	Algorytm Naiwnego Bayesa	9
2.7	Pseudokod algorytmu Naiwnego Bayesa	10
2.8	Wyniki klasyfikatora Naiwnego Bayesa	10
2.9	Implementacja	11
3	Wnioski	13
3.1	Porównanie wyników klasyfikatorów	13
3.2	Potencjał rozwoju	13
4	Pełen kod aplikacji	14
4.1	main.py	14
4.2	bayes.py	16
4.3	knn.py	17
4.4	test.py	17
4.5	utils.py	19

1 Wstęp

1.1 Opis programu

Celem programu jest klasyfikowanie odręcznie pisanych cyfr przez użytkownika. Aplikacja zawiera proste GUI, które pozwala użytkownikowi narysować cyfrę na płótnie, a następnie po wciśnięciu przycisku 'Rozpoznaj' program klasyfikuje cyfrę za pomocą jednego z klasyfikatorów w celu rozpoznania narysowanej cyfry. Otrzymane wyniki klasyfikacji są wyświetlane w bloku po prawej stronie interfejsu użytkownika.



Rysunek 1: Wygląd aplikacji



Rysunek 2: Rozpoznawanie

1.2 Użyte biblioteki

Program korzysta z następujących zewnętrznych bibliotek:

- Pillow
 - Do transformacji zapisanych cyfr na macierz
 - Do transformacji narysowanej cyfry na macierz
- numpy
- seaborn

1.3 Baza danych

Baza danych składa się z 1000 obrazów cyfr narysowanych przez nas (100 dla każdej cyfry). Każdy piksel obrazu jest reprezentowany w skali szarości (ma wartość od 0 do 255, gdzie 0 to biały, a 255 to czarny kolor). Obrazy są przechowywane w formacie png.



Rysunek 3: Przykładowy obraz z bazy danych



Rysunek 4: Ten sam obraz po zmianie rozdzielczości na 28x28 pikseli

2 Opis działania

2.1 Normalizacja danych

Piksele wczytanego obrazu są konwertowane na skalę szarości tj. 0 - kolor biały, 255 - kolor czarny oraz są normalizowane do przedziału od 0 do 1 co ułatwia modelowi dopasowanie cyfr. Wzór na normalizację pojedynczego piksela:

$$z_i = \frac{x_i - \min(x)}{\max(x) - \min(x)}$$

gdzie:

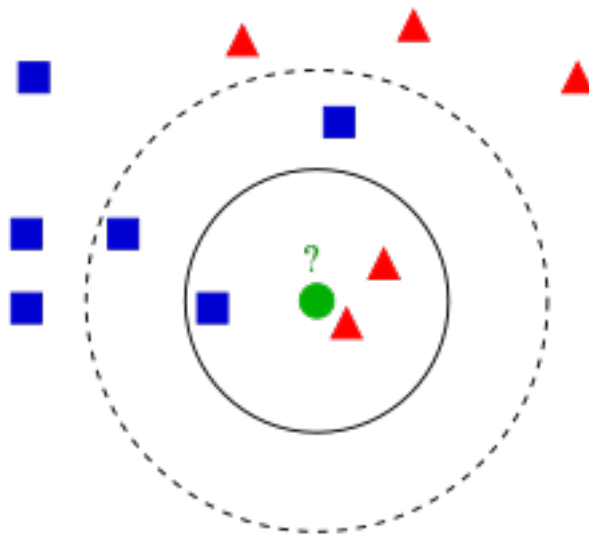
- z_i - znormalizowany piksel
- x_i - piksel
- x - zbiór wszystkich pikseli

Ostatecznie wzór ma postać:

$$z_i = \frac{x_i}{255}$$

2.2 Algorytm k najbliższych sąsiadów

Klasyfikator kNN to jedna z ważniejszych nieparametrycznych metod klasyfikacji. W tej metodzie klasyfikowany obiekt przydzielamy do tej klasy, do której należy większość z k sąsiadów.



Rysunek 5: Przykład klasyfikacji metodą kNN

W przypadku $k=3$ (mniejszy okrąg), zielona kropka zostanie zakwalifikowana do czerwonych trójkątów. W przypadku $k=5$ (większy okrąg) - do niebieskich kwadratów.

2.3 Metryka odległości

Użyta została odległość Minkowskiego określona wzorem:

$$L_m(x, y) = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{\frac{1}{p}}$$

gdzie:

- L_m - odległość między punktami x i y
- x, y - punkty w przestrzeni n wymiarowej
- x_i, y_i - i -ta współrzędna punktów x i y
- p - parametr określający rodzaj metryki

Przetestowaliśmy skuteczność klasyfikatora kNN dla liczby sąsiadów $k \in \{1, 3, 5, 7, 9, 11, 13, 15\}$ oraz dla wartości parametru $p \in \{1, 2, 3\}$.

2.4 Pseudokod algorytmu kNN

Data: Dane wejściowe: zbiór treningowy *train_data* zbiór testowy *test_data* liczba sąsiadów k metryka p

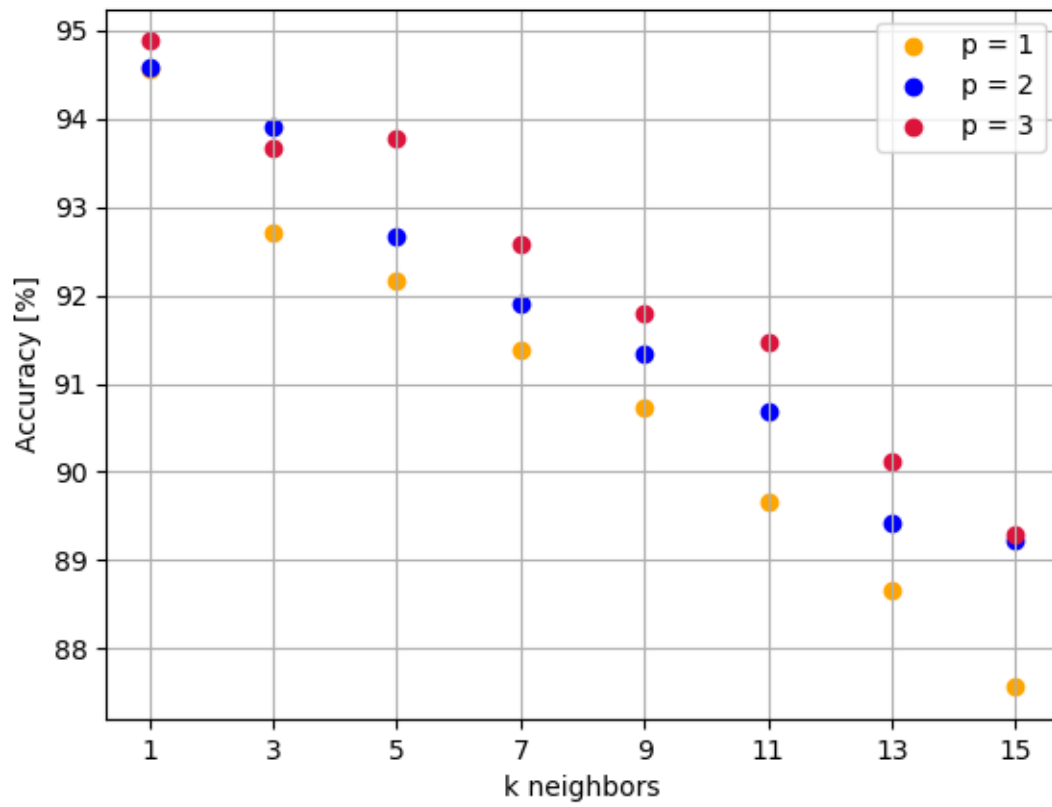
Result: Zbiór testowy z przewidzianymi etykietami

```
foreach test_instance in test_data do
    distances = [];
    foreach train_instance in train_data do
        distance = point_distance(test_instance, train_instance, p);
        distances.append((train_instance, distance));
    end
    sorted_distances = sort(distances, by = distance);
    k_nearest_neighbors = sorted_distances[: k];
    test_instance.set_predicted_label(predicted_label);
end
return test_data;
```

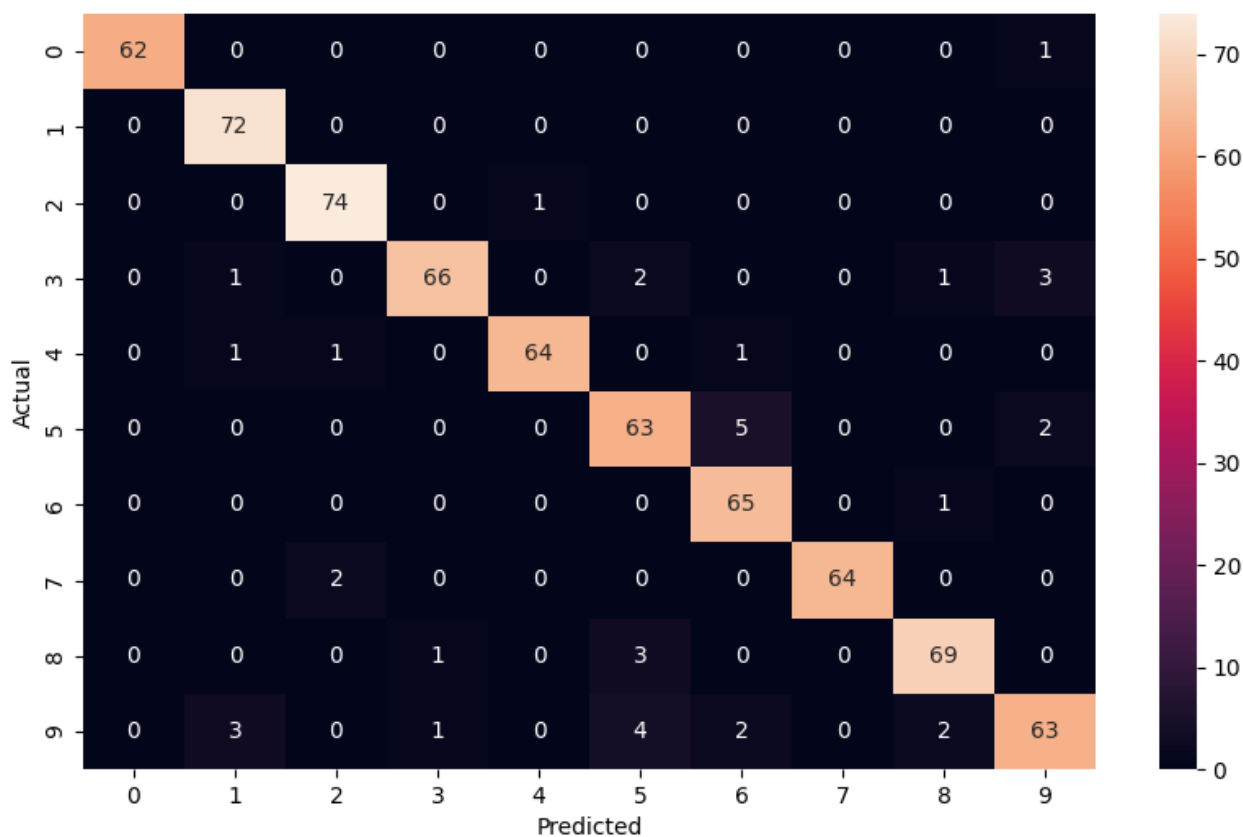
Algorithm 1: Algorytm k najbliższych sąsiadów.

2.5 Wyniki klasyfikatora kNN

Po przeanalizowaniu różnych wartości parametrów k i p , algorytm kNN klasyfikuje cyfry z największą dokładnością na poziomie 94.9% dla pary parametrów $k = 2$ i $p = 3$.



Rysunek 6: Zależność dokładności klasyfikacji od liczby sąsiadów k i parametru p



Rysunek 7: Macierz błędów algorytmu kNN dla $k = 1$ i $p = 3$

2.6 Algorytm Naiwnego Bayesa

Klasyfikator Naiwnego Bayesa jest to klasyfikator probabilistyczny, oparty na twierdzeniu Bayesa. Zakłada on niezależność cech. Naiwny klasyfikator bayesowski predykuje klasę nowego obiektu w zbiorze na podstawie prawdopodobieństwa warunkowego.

$$p(D_i | X_1, X_2, \dots, X_n)$$

gdzie:

D_i - klasa obiektu

X_1, X_2, \dots, X_n - cechy obiektu

Prawdopodobieństwo wystąpienia elementu w danym zbiorze

$$p(D_i) = \frac{|D_i|}{N}$$

gdzie:

$|D_i|$ - moc zbioru D_i

N - liczba wszystkich elementów w zbiorze

2.7 Pseudokod algorytmu Naiwnego Bayesa

Data: Dane wejściowe: Dane D Nowy punkt NP

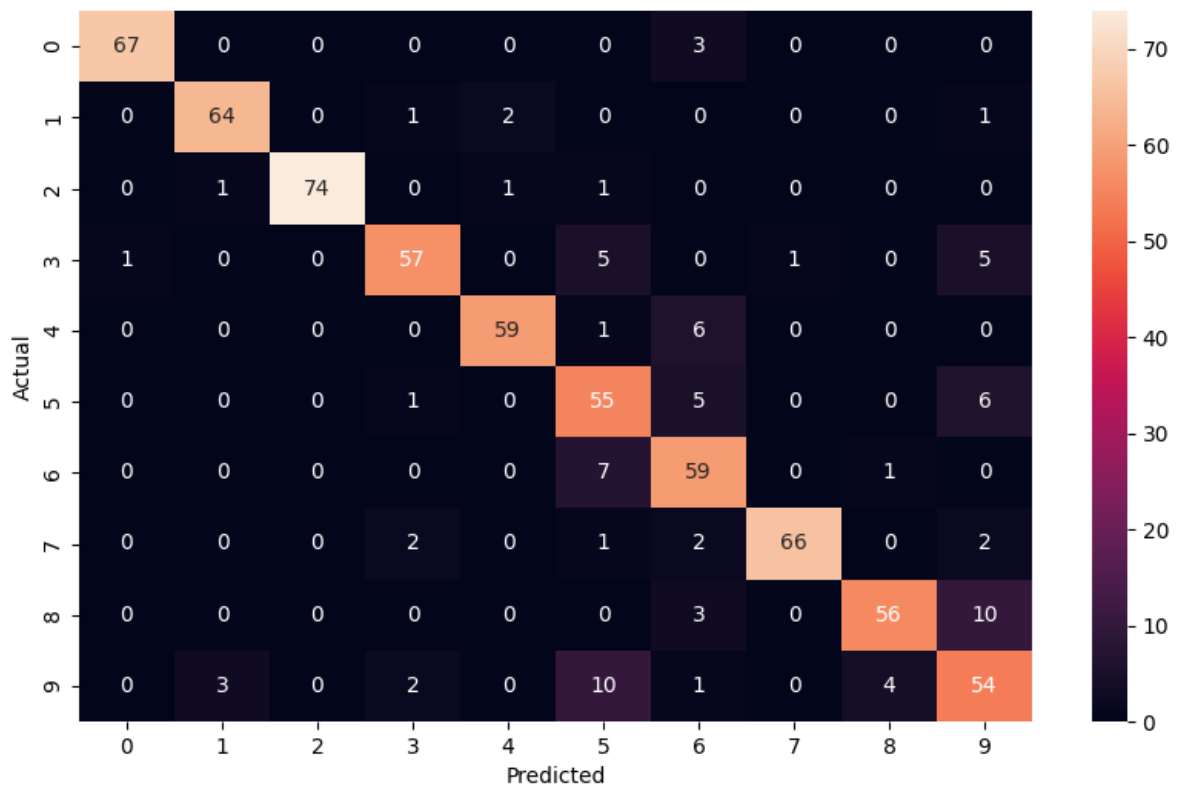
Result: Przewidziana klasa dla NP

```
probabilities = [];  
foreach  $c$  in classes do  
    | probabilities.append(probability of class  $c$ );  
end  
chosen_class = None;  
max_probability = -1;  
foreach  $p, c$  in zip(probabilities, classes) do  
    | if  $p > max\_probability$  then  
    |     | max_probability =  $p$ ;  
    |     | chosen_class =  $c$ ;  
    | end  
end  
return chosen_class
```

Algorithm 2: Algorytm Naiwnego Bayesa.

2.8 Wyniki klasyfikatora Naiwnego Bayesa

Po przeprowadzeniu kilkunastu testów algorytm Naiwnego Bayesa uzyskuje dokładność na poziomie 88.7%.



Rysunek 8: Macierz błędów algorytmu Naiwnego Bayesa

2.9 Implementacja

1. main.py - główny plik programu
2. knn.py - plik zawierający implementację algorytmu kNN
3. bayes.py - plik zawierający implementację naiwnego klasyfikatora Bayesa
4. utils.py - plik zawierający funkcje pomocnicze (np. wczytywanie danych)
5. test.py - plik zawierający funkcje testujące klasyfikatory

```

1 # Algorytm kNN
2 class KNN:
3     def __init__(self, k: int = 3, p: int = 2):
4         self.k = k
5         self.p = p
6
7     @staticmethod
8     def minkowski_distance(X1, X2, p) -> float:

```

```

9         return sum([abs((a - b) ** p) for a, b in zip(X1, X2)]) ** (1 /
10                     p)
11
12     def fit(self, X, Y) -> None:
13         self.X_train = X
14         self.Y_train = Y
15
16     def predict(self, X) -> list[float]:
17         return [self.predict_point(x) for x in X]
18
19     def predict_point(self, x) -> float:
20         distances = [
21             KNN.minkowski_distance(x, x_train, self.p)
22             for x_train in self.X_train
23         ]
24         k_indices = np.argsort(distances)[: self.k]
25         k_nearest_labels = [self.Y_train[i] for i in k_indices]
26         most_common = Counter(k_nearest_labels).most_common(1)
27         return most_common[0][0]
28
29 # Algorytm Naiwnego Bayesa
30 class NaiveBayesClassifier:
31     def fit(self, X, y):
32         self.classes = np.unique(y)
33         self.num_classes = len(self.classes)
34         self.num_features = X.shape[1]
35         self.class_probs = np.zeros(self.num_classes)
36         self.feature_probs = np.zeros((self.num_classes, self.
37             num_features))
38         for i, c in enumerate(self.classes):
39             X_c = X[y == c]
40             self.class_probs[i] = len(X_c) / len(X)
41             self.feature_probs[i] = np.mean(X_c, axis=0)
42
43     def predict(self, X):
44         y_pred = np.zeros(X.shape[0], dtype=np.str_)
45         for i, x in enumerate(X):
46             posterior_probs = []
47             for j in range(self.num_classes):
48                 prior = np.log(self.class_probs[j])
49                 likelihood = np.sum(
50                     np.log(self.compute_feature_prob(self.feature_probs[
51                         j], x))
52                 )
53                 posterior = prior + likelihood
54                 posterior_probs.append(posterior)
55             y_pred[i] = self.classes[np.argmax(posterior_probs)]
56         return y_pred
57
58     def compute_feature_prob(self, feature_prob, x):
59         epsilon = 1e-9
60         return feature_prob * x + (1 - feature_prob) * (1 - x + epsilon)
61
62 # Funkcja klasyfikująca narysowaną cyfrę
63 def predict_digit(img):

```

```
61     img = img.resize((28, 28))
62     img = img.convert("L")
63     img = np.array(img).flatten()
64     img = np.invert(img)
65     img = img / 255
66     x_train, y_train = read_digits("imgs")
67     x_train = x_train / 255
68     model = KNeighborsClassifier(n_neighbors=1, p=3)
69     model.fit(x_train, y_train)
70     y_pred = model.predict(np.array([img]))[0]
71     return y_pred
```

3 Wnioski

3.1 Porównanie wyników klasyfikatorów

W badaniu porównującym algorytm kNN i naiwny klasyfikator bayesowski w rozpoznawaniu pisma odręcznego stwierdzono, że k-NN osiągnął lepsze wyniki i dokładność klasyfikacji niż naiwny klasyfikator bayesowski. Wykorzystano zbiór danych, który został podzielony na dane treningowe (75%) i dane testowe (25%), gdzie każda próbka była reprezentowana jako obraz 28 x 28 pikseli przedstawiający odręczne cyfry. Naiwny klasyfikator bayesowski uzyskał dokładność na poziomie około 88.7%, podczas gdy kNN osiągnął dokładność na poziomie około 94.9%. Wynika z tego, że k-NN jest bardziej skutecznym narzędziem do rozpoznawania cyfr pisma odręcznego niż naiwny klasyfikator bayesowski, zwłaszcza w przypadku większej liczby danych treningowych i większego wymiaru wejściowego. Naiwny klasyfikator bayesowski ma jednak pewne zalety, takie jak krótszy czas potrzebny na dopasowanie i testowanie. W porównaniu do tego, klasyfikator kNN wymagał znacznie więcej czasu. Warto jednak zauważyć, że k-NN może być dobrym wyborem dla mniejszych zbiorów danych i problemów o mniejszej liczbie wymiarów.

3.2 Potencjał rozwoju

Projekt rozpoznawania cyfr przy użyciu uczenia maszynowego ma potencjał rozwoju w wielu różnych obszarach: - Zwiększenie dokładności rozpoznawania: Można podjąć działania mające na celu zwiększenie dokładności rozpoznawania cyfr. Można eksperymentować z różnymi algorytmami uczenia maszynowego, takimi jak głębokie sieci neuronowe, konwolucyjne sieci neuronowe (CNN).

- Rozpoznawanie innych zestawów danych: Oprócz rozpoznawania cyfr można rozważyć rozwinięcie projektu w celu rozpoznawania innych zestawów danych. Można badać możliwość rozpoznawania liter, symboli matematycznych lub innych obiektów.

- Adaptacja do innych języków i pism: Projekty rozpoznawania cyfr można rozszerzyć na inne języki i pisma. Na przykład, można badać możliwość rozpoznawania cyfr w języku chińskim, arabskim, japońskim lub innych językach, które posiadają swoje własne znaki i pismo. Odpowiednie dane treningowe i dostosowanie modelu będą kluczowe w takim przypadku.

4 Pełen kod aplikacji

4.1 main.py

```
1 import os
2 import sys
3 import tkinter as tk
4 from datetime import datetime
5
6 import numpy as np
7 from PIL import Image, ImageGrab # Pillow
8 from sklearn.naive_bayes import GaussianNB
9 from sklearn.neighbors import KNeighborsClassifier # scikit-learn
10
11 DIRNAME = os.path.dirname(os.path.dirname(__file__))
12 sys.path.append(DIRNAME)
13
14 from digit_recognition.bayes import NaiveBayesClassifier
15 from digit_recognition.knn import KNN
16 from digit_recognition.utils import read_digits
17
18
19 class Settings:
20     WIDTH = HEIGHT = 300
21     FONTSIZE = 20
22
23
24 def predict_digit(img):
25     img = img.resize((28, 28))
26     img = img.convert("L")
27     img = np.array(img).flatten()
28     img = np.invert(img)
29     img = img / 255
30     x_train, y_train = read_digits("imgs")
31     x_train = x_train / 255
32     model = KNeighborsClassifier(n_neighbors=1, p=3)
33     model.fit(x_train, y_train)
34     y_pred = model.predict(np.array([img]))[0]
35     return y_pred
36
37
38 class App(tk.Tk):
39     def __init__(self):
40         tk.Tk.__init__(self)
41         tk.Tk.title(self, "Rozpoznawanie cyfr")
42         self.x = self.y = 0
43         self.canvas = tk.Canvas(
44             self,
45             width=Settings.WIDTH,
46             height=Settings.HEIGHT,
47             bg="white",
48             cursor="cross",
49         )
50         self.label = tk.Label(
```

```

51         self, text="?", font=("Helvetica", Settings.FONTSIZE)
52     )
53     self.classify_button = tk.Button(
54         self, text="Rozpoznaj", command=self.classify_handwriting
55     )
56     self.clear_button = tk.Button(
57         self, text="Wyczyść", command=self.clear_all
58     )
59     self.save_button = tk.Button(
60         self, text="Zapisz", command=self.save_to_file
61     )
62     self.canvas.grid(
63         row=0,
64         column=0,
65         pady=2,
66         sticky=tk.W,
67     )
68     self.label.grid(row=0, column=1, pady=2, padx=2)
69     self.classify_button.grid(row=1, column=1, pady=2, padx=2)
70     self.clear_button.grid(row=1, column=0, pady=2)
71     self.save_button.grid(row=1, column=2, pady=2)
72     self.canvas.bind("<B1-Motion>", self.draw_lines)
73
74     def clear_all(self):
75         self.canvas.delete("all")
76
77     def get_canvas_image(self) -> Image:
78         x = self.winfo_rootx() + self.canvas.winfo_x()
79         y = self.winfo_rooty() + self.canvas.winfo_y()
80         x1 = x + self.canvas.winfo_width()
81         y1 = y + self.canvas.winfo_height()
82         return ImageGrab.grab().crop((x, y, x1, y1))
83
84     def classify_handwriting(self):
85         im = self.get_canvas_image()
86         digit = predict_digit(im)
87         self.label.configure(text=str(digit))
88
89     def save_to_file(self):
90         folder = "imgs"
91         filename = datetime.today().strftime("%d-%m-%Y %Hh%Mm%Ss") + ".png"
92         im = self.get_canvas_image()
93         im.save(os.path.join(folder, filename))
94
95     def draw_lines(self, event):
96         self.x = event.x
97         self.y = event.y
98         r = 20
99         self.canvas.create_oval(
100             self.x - r,
101             self.y - r,
102             self.x + r,
103             self.y + r,
104             fill="black",

```

```

105         outline="black",
106     )
107
108
109 def main() -> int:
110     app = App()
111     app.mainloop()
112     return 0
113
114
115 if __name__ == "__main__":
116     raise SystemExit(main())

```

4.2 bayes.py

```

1  import numpy as np
2
3
4  class NaiveBayesClassifier:
5      def fit(self, X, y):
6          self.classes = np.unique(y)
7          self.num_classes = len(self.classes)
8          self.num_features = X.shape[1]
9          self.class_probs = np.zeros(self.num_classes)
10         self.feature_probs = np.zeros((self.num_classes, self.num_features))
11
12         for i, c in enumerate(self.classes):
13             X_c = X[y == c]
14             self.class_probs[i] = len(X_c) / len(X)
15             self.feature_probs[i] = np.mean(X_c, axis=0)
16
17     def predict(self, X):
18         y_pred = np.zeros(X.shape[0], dtype=np.str_)
19         for i, x in enumerate(X):
20             posterior_probs = []
21             for j in range(self.num_classes):
22                 prior = np.log(self.class_probs[j])
23                 likelihood = np.sum(
24                     np.log(self.compute_feature_prob(self.feature_probs[j], x))
25                 )
26                 posterior = prior + likelihood
27                 posterior_probs.append(posterior)
28             y_pred[i] = self.classes[np.argmax(posterior_probs)]
29         return y_pred
30
31     def compute_feature_prob(self, feature_prob, x):
32         epsilon = 1e-9
33         return feature_prob * x + (1 - feature_prob) * (1 - x + epsilon)
34
35     def accuracy_score(self, y_test, y_pred):
36         return sum([y_test[i] == y_pred[i] for i in range(len(y_test))]) / len(y_test)

```


4.3 knn.py

```

1  from collections import Counter
2
3  import numpy as np
4
5
6  class KNN:
7      def __init__(self, k: int = 3, p: int = 2):
8          self.k = k
9          self.p = p
10
11     @staticmethod
12     def minkowski_distance(X1, X2, p) -> float:
13         return sum([abs((a - b) ** p) for a, b in zip(X1, X2)]) ** (1 / p)
14
15     def fit(self, X, Y) -> None:
16         self.X_train = X
17         self.Y_train = Y
18
19     def predict(self, X) -> list[float]:
20         return [self.predict_point(x) for x in X]
21
22     def predict_point(self, x) -> float:
23         distances = [
24             KNN.minkowski_distance(x, x_train, self.p)
25             for x_train in self.X_train
26         ]
27         k_indices = np.argsort(distances)[: self.k]
28         k_nearest_labels = [self.Y_train[i] for i in k_indices]
29         most_common = Counter(k_nearest_labels).most_common(1)
30         return most_common[0][0]
31
32     def accuracy_score(self, y_test, y_pred):
33         return sum([y_test[i] == y_pred[i] for i in range(len(y_test))]) /
34             len(
35                 y_test
36             )

```

4.4 test.py

```

1  import os
2  import sys
3  from collections import defaultdict
4
5  import matplotlib.pyplot as plt
6  from sklearn.metrics import accuracy_score, confusion_matrix
7  from sklearn.model_selection import train_test_split
8  from sklearn.neighbors import KNeighborsClassifier

```

```

9
10 DIRNAME = os.path.dirname(os.path.dirname(__file__))
11 sys.path.append(DIRNAME)
12
13
14 from digit_recognition.bayes import NaiveBayesClassifier
15 from digit_recognition.knn import KNN
16 from digit_recognition.utils import read_digits
17
18
19 def test_knn(
20     X,
21     y,
22     k_list: list[int],
23     p_list: list[int],
24     filename: str,
25     n_tests: int = 100,
26     plot: bool = False,
27     colors: list[str] = [],
28 ) -> None:
29     m2 = {}
30     for p in p_list:
31         m = defaultdict(float)
32         for _ in range(n_tests):
33             for k in k_list:
34                 X_train, X_test, y_train, y_test = train_test_split(
35                     X, y, train_size=0.75
36                 )
37                 model = KNeighborsClassifier(n_neighbors=k, p=p)
38                 model.fit(X_train, y_train)
39                 y_pred = model.predict(X_test)
40                 accuracy = accuracy_score(y_test, y_pred)
41                 m[k] += 100 * accuracy / n_tests
42     m2[p] = m
43     if plot:
44         assert len(colors) == len(p_list)
45         for i, (k, v) in enumerate(m2.items()):
46             plt.scatter(k_list, v.values(), label=f"p={k}", color=colors[i])
47         plt.grid()
48         plt.xlabel("k neighbors")
49         plt.ylabel("Accuracy [%]")
50         plt.legend([f"p = {p}" for p in p_list], loc="upper right")
51         plt.xticks(k_list)
52         plt.savefig(filename)
53     else:
54         print(m2)
55
56
57 def test_bayes(X, y, n_tests: int = 100) -> None:
58     acc = 0
59     for _ in range(n_tests):
60         X_train, X_test, y_train, y_test = train_test_split(
61             X, y, train_size=0.75
62         )
63         model = NaiveBayesClassifier()

```

```
64     model.fit(X_train, y_train)
65     y_pred = model.predict(X_test)
66     accuracy = model.accuracy_score(y_test, y_pred)
67     acc += 100 * accuracy / n_tests
68     print(f"accuracy: {round(acc, 3)}")
```

4.5 utils.py

```
1  import os
2
3  import numpy as np
4  from PIL import Image
5
6
7  def read_digits(dirname: str):
8      """Read digits' images from dirname directory.
9      Return a tuple of two lists:
10     - X: list of image pixels (each pixel is a number from
11       0 (white) to 255 (black).
12     - y: list of digit labels from 0 to 9.
13     Dimensions:
14     - X: nsamples * 28 * 28 (each image is resized to 28x28 pixels).
15     - y: nsamples."""
16     X = []
17     y = []
18     for subdir, _, files in os.walk(dirname):
19         for file in files:
20             label = subdir[-1]
21             filepath = subdir + os.sep + file
22             img = Image.open(filepath)
23             img = img.resize((28, 28))
24             img = img.convert("L")
25             img = np.array(img).flatten()
26             img = np.invert(img)
27             X.append(img)
28             y.append(label)
29     return np.array(X), np.array(y)
```
