

Prof. R. Rojas

Mustererkennung, WS17/18

Übungsblatt 1

Boyan Hristov, Matrikelnummer: 4980301

25. Oktober 2017

Link zum Git Repository: <https://github.com/BoyanH/Freie-Universitaet-Berlin/tree/master/MachineLearning/Homework1>

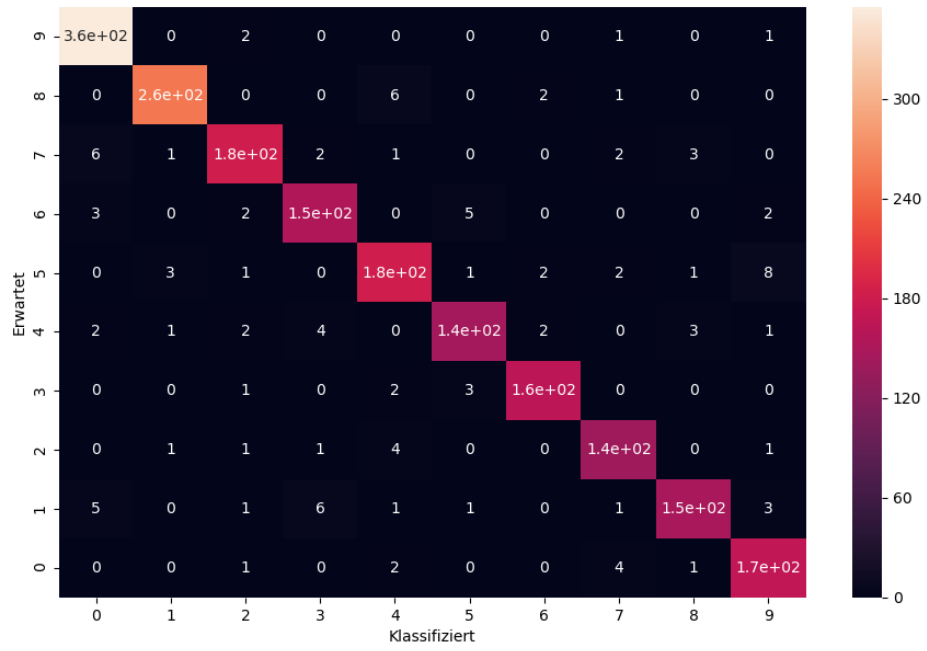
Beschreibung

Ziffererkennung mit kNN-Classifer. Dabei wurde eine simple lineare kNN Klassifizierung ohne kd-Baum benutzt. Die Eingabedaten (Grayscale Pixel von Bilder) wurden als Vektoren in n-Dimensionalen Raum (n=256) betrachtet, ohne weitere Verarbeitung von dem Merkmalen.

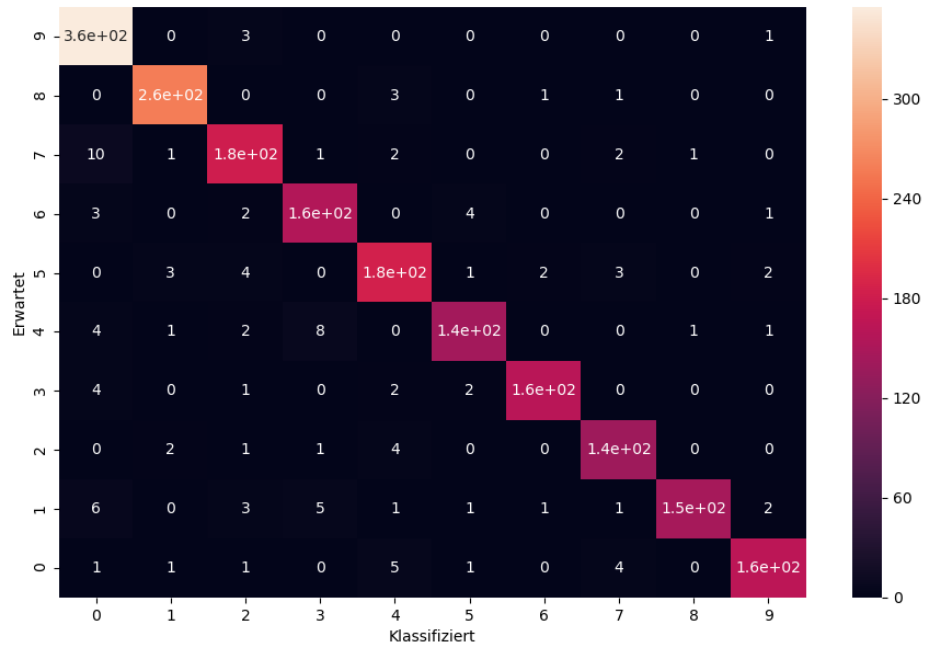
Fehlerrate

Das ist die Ausgabe des Programs und damit auch die Fehlerrate

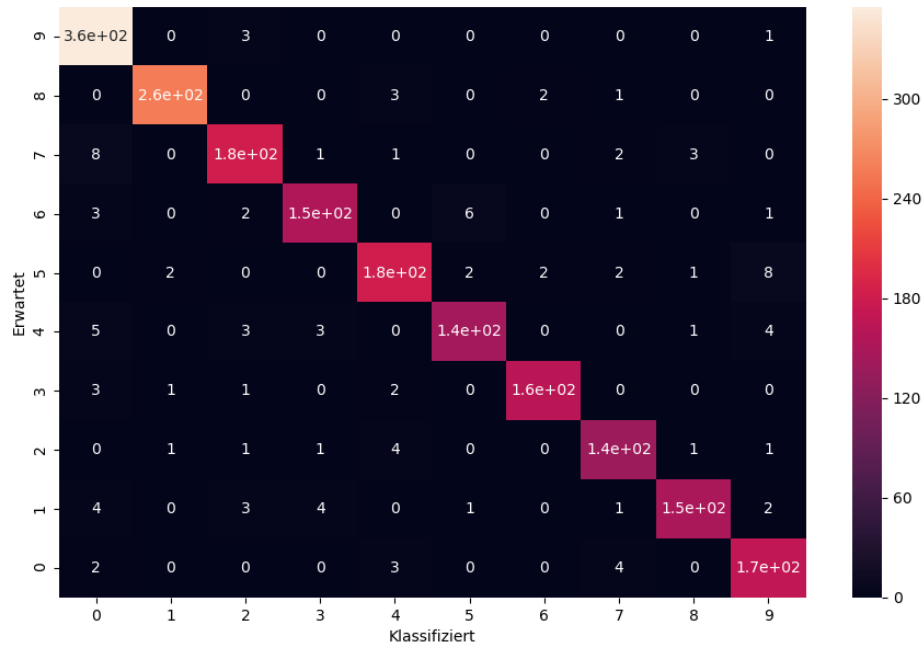
```
1 Error rate for k=1 is 5.630293971101146%
2 Error rate for k=2 is 5.879422022919781%
3 Error rate for k=3 is 5.5306427503736915%
```



$k=1$



$k=2$



k=3

Code

```

1 import os
2 import math

4 import seaborn as sn
5 import pandas as pd
6 import matplotlib.pyplot as plt

8 dir_path = os.path.dirname(os.path.realpath(__file__))
9 trainingSetData = []

11 def extractDataFromLine(line):
12     line = line.replace(' \n', '') # clear final space and new line chars

14     return list(map(float, line.split(' '))); # map line to a list of floats

16 def train():
17     explicitPathRead = os.path.join(dir_path, './Dataset/train')
18     f = open(explicitPathRead, 'r')

20     for line in f:
21         line = line.replace(' \n', '') # clear final space and new line chars
22         currentDigitData = extractDataFromLine(line)
23         trainingSetData.append(currentDigitData)

25     # TODO: maybe construct a kd-tree out of data

27 def test(knnCount):
28     explicitPathRead = os.path.join(dir_path, './Dataset/test')
29     f = open(explicitPathRead, 'r')

31     errors = 0
32     testsCount = 0

```

```

33 confusionsMatrix = [[0 for x in range(10)] for y in range(10)] # 10 by 10 matrix

34
35 for line in f:
36     currentLineData = extractDataFromLine(line)
37     result = classify(currentLineData, knnCount)
38     expected = str(int(currentLineData[0]))
39
40     confusionsMatrix[int(expected)][int(result)] += 1
41
42     testsCount += 1
43     if result != expected:
44         errors += 1
45
46 errorRate = errors / testsCount
47 printConfusionsMatrix(confusionsMatrix, knnCount)
48 print("Error rate for k={0} is {1}%".format(knnCount, errorRate * 100))
49
50 def classify(digitData, knnCount):
51     kNN = getKNN(trainingSetData, digitData, knnCount)
52     labelsOfKNN = list(map(getLabel, kNN)) # intentionally as integers, to make mapping easier
53     digitRepetitions = [0 for x in range(10)]
54
55     for label in labelsOfKNN:
56         digitRepetitions[label] += 1
57
58     return str(digitRepetitions.index(max(digitRepetitions)))
59
60
61 def getKNN(trainingSetData, digitData, k):
62     knn = [] # distance to the kNN
63     knnToDataMapper = {}
64
65     for trainingDigitData in trainingSetData:
66         currentDistance = getDistanceBetweenPoints(getCoords(digitData), getCoords(
            trainingDigitData))
67
68         if len(knn) < k or max(knn) > currentDistance:
69             knn.append(currentDistance)
70             knn.sort()
71             knnToDataMapper[currentDistance] = trainingDigitData
72             knn = knn[:k]
73
74     return list(map(lambda x: knnToDataMapper[x], knn))
75
76
77 def getCoords(data):
78     return data[1:]
79
80 def getLabel(data):
81     return int(data[0])
82
83 def getDistanceBetweenPoints(a, b):
84     if len(a) != len(b):
85         raise Exception('Points must be in the same n-dimensional space to calculate distance!')
86
87     squares = 0
88
89     for i in range(len(a)):
90         squares += math.pow(a[i] - b[i], 2)
91
92     return math.sqrt(squares)
93
94 def printConfusionsMatrix(matrix, k):
95     explicitImgPath = os.path.join(dir_path, './Plots/confusion_matrix_for_k_{0}.png'.format(k))
96     digits = [str(x) for x in range(10)];
97
98     df_cm = pd.DataFrame(matrix, index = digits,

```

```
99 columns = digits.reverse() )
101 plt.figure(figsize = (11,7))
102 heatmap = sn.heatmap(df_cm, annot=True)
104 heatmap.set(xlabel='Klassifiziert', ylabel='Erwartet')
106 plt.savefig(explicitImgPath, format='png')
108 train()
109 test(1)
110 test(2)
111 test(3)
```