

Prof. R. Rojas

Mustererkennung, WS17/18

Übungsblatt 3

Boyan Hristov, Nedeltscho Petrov

7. November 2017

Link zum Git Repository: <https://github.com/BoyanH/FU-MachineLearning-17-18/tree/master/Solutions/Homework3>

Klassifikation mit Gauss Verteilung

Score (Output des Programs) bzw. Analyse

Das ist die Ausgabe des Programs und damit auch das Score

```
1 Score 3 vs 5: 0.923312883436
2 Score 3 vs 7: 0.929712460064
3 Score 3 vs 8: 0.843373493976
4 Score 5 vs 7: 0.931596091205
5 Score 5 vs 8: 0.874233128834
6 Score 7 vs 8: 0.837060702875
7 Score 3 vs 5 vs 7 vs 8: 0.838810641628
8 Score all: 0.840558046836
9 Score all with train data: 0.969551501852
```

Leider haben wir sehr komisches Score bei der letzten Hausaufgabe bekommen und können das jetzige anhand von selbsterstellte Daten von linearen Regression vergleichen. Mit den Daten aus der Musterlösung, die im Tutorium vorgestellt wurde, wurden wir sagen, dass die jetzige Lösung mit Gauss Verteilung nicht besser (sogar schlechter) ist, als diese mit linearen Regression (da war die Unterscheidung zwischen 5 und 7 mit 94 + % Genauigkeit).

Wenn man alle Ziffer vergleicht, bekommt man ein Score von 84.0558046836% Genauigkeit mit den Testdaten und 96.9551501852% Genauigkeit mit den Traindaten. D.h. wir lernen immer noch ganz viel äuswendig"bzw. es gibt relativ viel überfitting".

Das Programm funktioniert ungefähr so schnell wie bei der Klassifikation mit linearen Regression. Die Implementierung war aber schwieriger und kostet mehr Speicher damit es performant ist (Zentren, Konvergenzmatritzen, pseudo-invertierte Konvergenzmatritzen und Determinanten). Bei der Implementierung mit linearen Regression musste man nur die "Hat"Matritzen speichern.

Es ist aber leichter für mehrere Klassen zu implementieren. Man muss nicht separat die Ergebnissen vom mehreren Klassifikatoren vergleichen um zwischen mehrere Klassen zu unterscheiden. Man kann direkt

die Wahrscheinlichkeit, dass ein Punkt zu einer Klasse gehört berechnen, und dann direkt die Klasse mit der Höchsten Wahrscheinlichkeit wählen.

Erklärung interessanteren Teilen des Codes

Die Fit Methode ist diesmal nicht besonders interessant, da wir dort nur die Zentren berechnen (Alle Punkte aufsummieren und durch ihre Anzahl teilen) und die Konvergenzmatrixberechnungsfunktion aufrufen. Diese wäre aber interessant, da wir die vektorisiert haben und die Hauptimplementierung ausgelagert haben.

Die Summe und die Teilung durch der Anzahl kann man hier auch sehen. Sonst haben wir `numpy.vectorize` benutzt, um die Berechnung der einzelnen Summaren zu beschleunigen.

```
1 self.covariance_matrix[label] = np.vectorize(GaussianClassifier.covariance_for_point,
    signature='(m),(n)->(m,m)')(
2     points_per_label[label], self.centers[label]).sum(axis=0) / len(
    points_per_label[label])
```

Wir nutzen hier die Formel aus der Vorlesung. Das einzigste Unterschied hier ist, dass wir den ersten Vektor transponieren und nicht der zweiten. Das ist so, da wir am Ende eine Matrix bekommen wollen und Numpy interpretiert eine 1-dimensionale Liste als eine Matritze der Form $1 \times n$ und nicht $m \times 1$ (Vektor) als wir diese betrachten wollen und umgekehrt für das zweite.

```
1 def covariance_for_point(point, center):
2     return np.matrix(point - center, dtype=np.float64).T.dot(np.matrix(point - center,
    dtype=np.float64))
```

Bei der Berechnung der Wahrscheinlichkeit haben wir auch die Formel aus der Vorlesung benutzt. Wichtig wäre hier, dass es zu eine Determinante = 0 kommen könnte, wenn wir linear abhängige Features haben. Das ist oft der Fall bei manchen Ziffern, da wir die Pixels als Features betrachten und diese oft abhängig sind. Z.b wenn eine Ziffer ein Horizontales Strich hat, dann falls ein Pixel grauer ist, sind mit größtem Wahrscheinlichkeit auch diese daneben (links und rechts) auch grauer.

Damit wir die 0 Determinante umgehen (wegen Division mit 0), haben wir $2 * \pi * ||X||$ ausgelagert und benutzen später das Minimum von die und eine kleine Delta.

```
1 def get_possibility_for_class(self, point_class, point):
2     two_pi_det = 2 * math.pi * self.covariance_matrix_det[point_class]
3     left_side = 1 / max(0.2, math.sqrt(two_pi_det))
4     right_side = math.e**(-0.5 * (point - self.centers[point_class]).T.
5         dot(self.covariance_matrix_pinv[point_class]).dot(point -
        self.centers[point_class]))
7     return left_side * right_side
```

Die Voraussage machen wir, in dem wir die Wahrscheinlichkeit für je Klasse berechnen und dann diese mit höchsten nehmen.

```
1 possibilities = list(map(lambda x: self.get_possibility_for_class(x, point), self.classes))
2 winning_index = possibilities.index(max(possibilities))
3 return self.classes[winning_index]
```

Vollständiges Code

```
1 from Classifier import Classifier
2 from Parser import *
3 import numpy as np
4 import math

7 class GaussianClassifier(Classifier):
8     @staticmethod
9     def covariance_for_point(point, center):
10         # calculate covariance for a single point (well not really, but a single summar
11         # thingy)
12         # idea of this method is to easily vectorize it with np.vectorize
13
14         # implementation specific:
15         # in the formula, the first vector should be transposed and the second not
16         # this is only done, because we need to receive a matrix at the end
17         # with the way numpy handles single vectors, we actually need to transpose the
18         # first one
19         # and not the second in order to do that what we are used to in math
20         # (numpy treats a 1-dimensional array as 1xn matrix and not as nx1 as we want to)
21         return np.matrix(point - center, dtype=np.float64).T.dot(np.matrix(point - center,
22         dtype=np.float64))

23 def __init__(self, train_data, classes = [x for x in range(10)]):
24     """
25     :param classes: list of classes the classifier should train itself to distinguish
26     (e.g [3,5] for 3 vs 5 classifier) default is all digits
27     :param trainData:
28     :param trainLabels:
29     :param testData:
30     :param testLabels:
31     """
32
33     self.centers = {}
34     self.covariance_matrix = {}
35     self.covariance_matrix_det = {}
36     self.covariance_matrix_pinv = {}
37     (train_labels, train_points) = get_labels_and_points_from_data(train_data, classes)
38     self.classes = classes
39     self.fit(train_labels, train_points)

40 def fit(self, train_labels, train_points):
41     assert(len(train_labels) == len(train_points))
42     points_per_label = {}

43     # sort points in a dictionary, separated by classes
44     # eg {3: [first 256 dimension vector, second 256 dimensional vector, etc.], 5: ...
45     # ...}
46     for idx, point in enumerate(train_points):
47         current_label = train_labels[idx]
48         if current_label not in points_per_label:
49             points_per_label[current_label] = [point]
50         else:
51             points_per_label[current_label].append(point)

52     # then for each class, find the centroid and the covariance matrix
53     # for optimization reasons, we also save the inverse of the covariance matrix and
54     # it's determinant
55     for label in points_per_label:
56         # average of all points from the current class (with axis 0, so row-wise
57         # average)
58         self.centers[label] = np.array(points_per_label[label], dtype=np.float64).mean
59         (0)
60         # calculate covariance matrix using vectorization (see covariance_for_point
61         # static method)
```

```

58         # using the formula  $1/n * (\sum_i ((point - center)(point - center)^T))$ 
59         self.covariance_matrix[label] = np.vectorize(GaussianClassifier.
covariance_for_point, signature='(m),(n)->(m,m)')(
60             points_per_label[label], self.centers[label]).sum(axis=0) / len(
points_per_label[label])
61         # also calculate and save determinant and pseudo-inverse of matrix for
performance reasons
62         self.covariance_matrix_det[label] = np.linalg.det(self.covariance_matrix[label]
])
63         self.covariance_matrix_pinv[label] = np.linalg.pinv(self.covariance_matrix[
label])

64     def predict(self, X):
65         return list(map(lambda x: self.predict_single(x), X))

66     def predict_single(self, point):
67         possibilities = list(map(lambda x: self.get_possibility_for_class(x, point), self.
classes))
68         winning_index = possibilities.index(max(possibilities))
69
70         return self.classes[winning_index]

71     def get_possibility_for_class(self, point_class, point):
72         # using the formula from the lecture, calculate the probability for a point with
coordinates to
73         # be part of a class

74         # only important thing here is that  $2 * \pi * \det(covariance\_matrix)$  can be zero
75         # (in case the covariance_matrix doesn't have a full rank (when we have identical
values for some features
76         # this can often be the case because of the white pixels at the edges)),
77         # so we use
78         # np.nextafter to replace any zeros with a reaaaaly small float (because of
DivideByZero exceptions...)

79         two_pi_det = 2 * math.pi * self.covariance_matrix_det[point_class]
80         left_side = 1 / max(0.2, math.sqrt(two_pi_det))
81         right_side = math.e**(-0.5 * (point - self.centers[point_class]).T.
82             dot(self.covariance_matrix_pinv[point_class]).dot(point -
self.centers[point_class]))

83         return left_side * right_side

84 train_data = parse_data_file('./Dataset/train')
85 test_data = parse_data_file('./Dataset/test')

86 three_vs_five = GaussianClassifier(train_data, [3,5])
87 (three_vs_five_test_labels, three_vs_five_test_data) = get_labels_and_points_from_data(
test_data, [3,5])
88 print("Score 3 vs 5: {}".format(three_vs_five.score(three_vs_five_test_data,
three_vs_five_test_labels)))

89 three_vs_seven = GaussianClassifier(train_data, [3,7])
90 (three_vs_seven_test_labels, three_vs_seven_test_data) = get_labels_and_points_from_data(
test_data, [3,7])
91 print("Score 3 vs 7: {}".format(three_vs_seven.score(three_vs_seven_test_data,
three_vs_seven_test_labels)))

92 three_vs_eight = GaussianClassifier(train_data, [3,8])
93 (three_vs_eight_test_labels, three_vs_eight_test_data) = get_labels_and_points_from_data(
test_data, [3,8])
94 print("Score 3 vs 8: {}".format(three_vs_eight.score(three_vs_eight_test_data,
three_vs_eight_test_labels)))

95 five_vs_seven = GaussianClassifier(train_data, [5,7])
96 (five_vs_seven_test_labels, five_vs_seven_test_data) = get_labels_and_points_from_data(
test_data, [5,7])

```

```

108 print("Score 5 vs 7: {}".format(five_vs_seven.score(five_vs_seven_test_data,
    five_vs_seven_test_labels)))

110 five_vs_eight = GaussianClassifier(train_data, [5,8])
111 (five_vs_eight_test_labels, five_vs_eight_test_data) = get_labels_and_points_from_data(
    test_data, [5,8])
112 print("Score 5 vs 8: {}".format(five_vs_eight.score(five_vs_eight_test_data,
    five_vs_eight_test_labels)))

114 seven_vs_eight = GaussianClassifier(train_data, [7,8])
115 (seven_vs_eight_test_labels, seven_vs_eight_test_data) = get_labels_and_points_from_data(
    test_data, [7,8])
116 print("Score 7 vs 8: {}".format(seven_vs_eight.score(seven_vs_eight_test_data,
    seven_vs_eight_test_labels)))

118 combined = GaussianClassifier(train_data, [3, 5, 7, 8])
119 (combined_test_labels, combined_test_data) = get_labels_and_points_from_data(test_data,
    [3,5,7,8])
120 print("Score 3 vs 5 vs 7 vs 8: {}".format(combined.score(combined_test_data,
    combined_test_labels)))

122 all_digits = [x for x in range(10)]
123 all_classifier = GaussianClassifier(train_data, all_digits)
124 (all_test_labels, all_test_data) = get_labels_and_points_from_data(test_data, all_digits)
125 print("Score all: {}".format(all_classifier.score(all_test_data, all_test_labels)))

127 all_digits = [x for x in range(10)]
128 all_classifier = GaussianClassifier(train_data, all_digits)
129 (all_train_labels, all_train_data) = get_labels_and_points_from_data(train_data, all_digits
    )
130 print("Score all: {}".format(all_classifier.score(all_train_labels, all_train_data)))

```

Und das simple Parser

```

1 import csv
2 import numpy as np

5 def parse_data_file(file_name):
6     file = open(file_name, 'rt')
7     reader = csv.reader(file, delimiter=' ', quoting=csv.QUOTE_NONE)
8     data = []

10     for row in reader:
11         filtered = list(filter(lambda x: x != '', row))
12         data.append(list(map(lambda x: float(x), filtered)))

14     return data

17 def get_labels_and_points_from_data(data, classes):
18     data = list(filter(lambda x: int(x[0]) in classes, data))
19     labels = np.array(list(map(lambda x: int(x[0]), data)))
20     points = np.array(list(map(lambda x: x[1:], data)), dtype=np.float64)

22     return labels, points

```