

Prof. R. Rojas

Mustererkennung, WS17/18

Übungsblatt 4

Boyan Hristov, Nedeltscho Petrov

15. November 2017

Link zum Git Repository: <https://github.com/BoyanH/FU-MachineLearning-17-18/tree/master/Solutions/Homework4>

Fischer Klassifikation

Score (Output des Programs) bzw. Analyse

Das ist die Ausgabe des Programs und damit auch das Score

```
1 Best score for seed=879: 0.930510314875
2 Worst score for seed530: 0.880564603692
3 Score: 0.9305103148751357
4 Score with linear regression: 0.9055374592833876
```

Wie wir hier sehen, ist das Score auch sehr davon abhängig, wie man die Daten in Test und Train Sets splitet. Es ist deswegen so, weil einige Emails sehr untypische spam Emails bzw nicht spam Emails sind. Deswegen ist unser Score zwischen 88

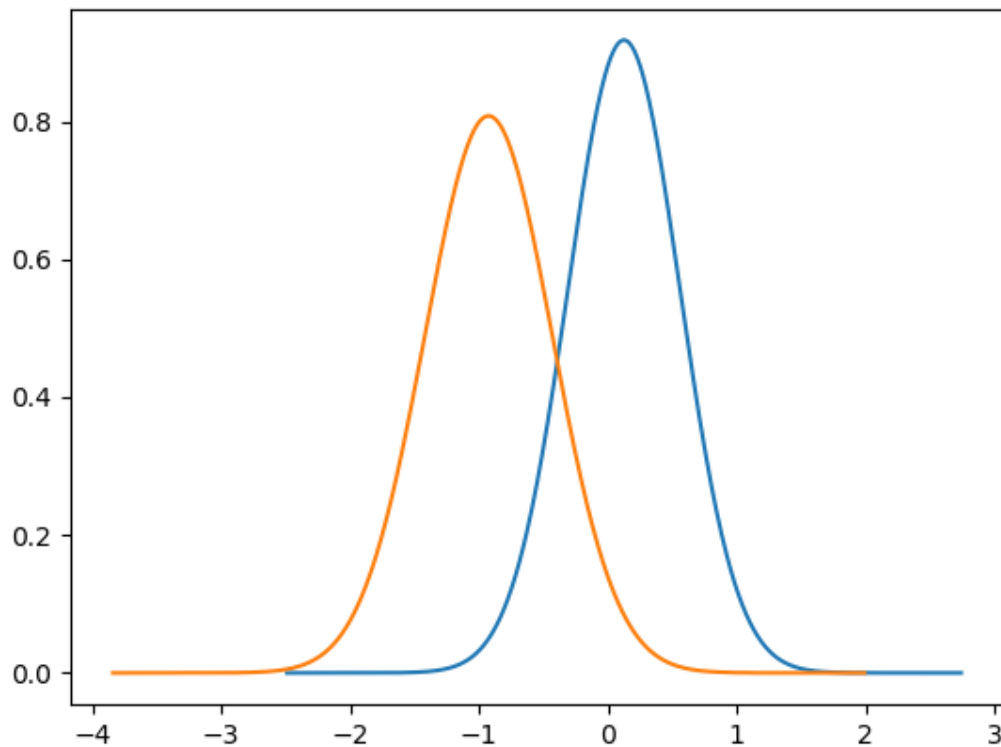
Bei linearen Regression haben wir etwas schlechteres Ergebniss mit dem Datensatz bekommen (lineare Regression von scikit-learn wurde benutzt). Das ist deswegen so, wie in der Vorlesung erklärt, weil die Kovarianzmatritzen in dem Fall für Fischer Klassifizierung besser geeignet sind. Anders gesagt, die Streuung der Daten von je Klasse liegen teilweise othogonal zu einander.

Aus performance-sichten haben wir nicht erkennbar bessere Ergebnisse bekommen (eine weitere lineare Methode, ist schnell genug).

Der Datensatz wurde im Jahre 1999 veröffentlicht, die Ergebnisse von 1998 zeigen circa 70% Da wir aber deutlich später das entwickelt haben, denken wir nicht, dass unsere Ergebnisse ausreichend sind. Wie schon auf der Webseite des Datensatzes erleutert wurde, sind falsch als Spam erkannte Emails ganz schädlich und 71% in dem Fall schon viel.

Plot

Auf dem Plot haben wir die Dichtefunktion, projiziert auf dem Fischer-Vektor berechnet. Wie man sieht, überlappen schon die Klassen sich einigermaßen. Man könnte also auch bessere Ergebnisse bekommen, könnte aber auch viel schlimmer sein.



Details zur Implementierung

In der Fit Methode haben wir beide Kovarianzmatritzen berechnet, als auch die Medians. Anhand der Formel von der Vorlesung haben wir dann unser Fischer Vektor berechnet. Nachdem, haben wir den Mittelpunkt beider Zentren auf dem Fischer-Vektor projiziert, damit wir dem später für die Klassifizierung benutzen können.

```
1  def fit(self, X_train, y_train):
2      X_a, X_b = FisherClassifier.split_in_classes(X_train, y_train)
3      cov_mat_a = np.cov(X_a, rowvar=False, bias=True)
4      cov_mat_b = np.cov(X_b, rowvar=False, bias=True)
5      center_a = np.array(X_a, dtype=np.float64).mean(0)
6      center_b = np.array(X_b, dtype=np.float64).mean(0)
7
8      alpha = np.linalg.pinv(cov_mat_a + cov_mat_b).dot(center_a - center_b)
9      alpha_normalized = alpha / np.linalg.norm(alpha)
10     self.alpha = alpha_normalized
11
12     # to determine whether a point belongs to class a or class b we need a threshold
13     # on the 1 dimensional space. This one is the projected point between the 2 centers
```

```

14         self.threshold = self.project_point((center_a + center_b) / 2)
16         self.plot_probability_distribution(center_a, center_b, X_a, X_b)

```

Die eigentliche Klassifizierung ist ganz simpel. Wir projizieren das neue Punkt auf dem Fischer-Vektor und schauen, ob es vor oder nach unser in Fit definierten Threshold ist.

```

1     def predict_single(self, x):
2         # project x into alpha (AKA Fisher's vector)
3         projected = self.project_point(x)
4         return projected < self.threshold

```

Die ganze Plot Funktionalität ist nicht so interessant, außer 2 relevante Teile. Wir haben die folgende Funktion benutzt, um die Dichtefunktion zu berechnen (da bei Fischer Normalverteilung eine Voraussetzung ist)

```

1 @staticmethod
2     def get_density_function(center, covariance):
3         return lambda x: math.e ** (
4             (-1/2) * ((x - center) / covariance) ** 2
5             ) / (covariance * math.sqrt((2*math.pi)))

```

Dabei brauchen wir aber die Kovarianz. Die wird von alle projizierten Punkten so berechnet.

```

1 @staticmethod
2     def get_covariance_for_projected(points, center):
3         vectorized_sq_distances_sum = np.vectorize(lambda x, m: (x - m)**2)
4         square_distances_sum = np.sum(vectorized_sq_distances_sum(points, center))
5         return math.sqrt(square_distances_sum / len(points))

```

Code in FisherClassifier.py

```

1 from Classifier import Classifier
2 from Parser import get_data_set
3 import numpy as np
4 from sklearn.linear_model import LinearRegression
5 import math
6 from random import random
7 import matplotlib.pyplot as plt
8 from scipy.stats import multivariate_normal
9
10
11 class FisherClassifier(Classifier):
12     def __init__(self, Ximport csv
13 import numpy as np
14 import os
15 from sklearn.model_selection import train_test_split
16
17
18 def parse_data():
19     file_name = os.path.join(os.path.dirname(__file__), './Dataset/spambase.data')
20     csv_file = open(file_name, 'rt')
21     reader = csv.reader(csv_file, delimiter=',', quoting=csv.QUOTE_NONE)
22     data = []
23
24     for row in reader:
25         filtered = list(filter(lambda x: x != '', row))
26         data.append(list(map(lambda x: float(x), filtered)))
27
28     return data

```

```

31 def get_points_and_labels_from_data(data):
32     points = np.array(list(map(lambda x: x[:-1], data)), dtype=np.float64)
33     labels = np.array(list(map(lambda x: int(x[-1]), data)))
34
35     return points, labels
36
37
38 def get_data_set(seed):
39     data = parse_data()
40     X, y = get_points_and_labels_from_data(data)
41     # for determined results we use a seed for random_state, so that data is always split
42     X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.8, test_size
43                                                         =0.2,
44                                                         random_state=seed)
45
46     return X_train, X_test, y_train, y_test
47
48 _train, y_train):
49     self.alpha = None
50     self.threshold = None
51     self.fit(X_train, y_train)
52
53 @staticmethod
54 def split_in_classes(X_train, y_train):
55     split_X = ([], [])
56
57     for idx, x in enumerate(X_train):
58         current_label = y_train[idx]
59         split_X[current_label].append(x)
60
61     return split_X
62
63 @staticmethod
64 def get_density_function(center, covariance):
65     return lambda x: math.e ** (
66         (-1/2) * ((x - center) / covariance) ** 2
67     ) / (covariance * math.sqrt((2*math.pi)))
68
69 @staticmethod
70 def get_covariance_for_projected(points, center):
71     vectorized_sq_distances_sum = np.vectorize(lambda x, m: (x - m)**2)
72     square_distances_sum = np.sum(vectorized_sq_distances_sum(points, center))
73     return math.sqrt(square_distances_sum / len(points))
74
75 def plot_class(self, points, center):
76     projected_center = self.project_point(center)
77     projected_points = list(map(lambda x: self.project_point(x), points))
78     covariance = FisherClassifier.get_covariance_for_projected(projected_points,
79 projected_center)
80     density_a = FisherClassifier.get_density_function(projected_center, covariance)
81     plot_distance = 5000
82     y_of_plot = [density_a(float(x) / 100) for x in range(-plot_distance, plot_distance
83 )]
84
85     x_to_plot = [float(x) / 100 for x in range(-plot_distance, plot_distance)]
86     first_decent = None
87     last_decent = None
88
89     for idx, y in enumerate(y_of_plot):
90         if y > 0.01 and first_decent is None:
91             first_decent = idx
92         elif y <= 0.01 and first_decent is not None:
93             last_decent = idx
94             break
95
96     beauty_margin = (last_decent - first_decent)
97     start = int(first_decent - beauty_margin/2)
98     end = int(last_decent + beauty_margin/2)

```

```

95     plt.plot(x_to_plot[start:end], y_of_plot[start:end])

97     def plot_probability_distribution(self, center_a, center_b, points_a, points_b):
98         self.plot_class(points_a, center_a)
99         self.plot_class(points_b, center_b)
100        plt.show()

102     def project_point(self, x):
103         return x.dot(self.alpha)

105     def fit(self, X_train, y_train):
106         X_a, X_b = FisherClassifier.split_in_classes(X_train, y_train)
107         cov_mat_a = np.cov(X_a, rowvar=False, bias=True)
108         cov_mat_b = np.cov(X_b, rowvar=False, bias=True)
109         center_a = np.array(X_a, dtype=np.float64).mean(0)
110         center_b = np.array(X_b, dtype=np.float64).mean(0)

112         alpha = np.linalg.pinv(cov_mat_a + cov_mat_b).dot(center_a - center_b)
113         alpha_normalized = alpha / np.linalg.norm(alpha)
114         self.alpha = alpha_normalized

116         # to determine whether a point belongs to class a or class b we need a threshold
117         # on the 1 dimensional space. This one is the projected point between the 2 centers
118         self.threshold = self.project_point((center_a + center_b) / 2)

120         self.plot_probability_distribution(center_a, center_b, X_a, X_b)

122     def predict(self, X):
123         return list(map(lambda x: self.predict_single(x), X))

125     def predict_single(self, x):
126         # project x into alpha (AKA Fisher's vector)
127         projected = self.project_point(x)
128         return projected < self.threshold

131 # max_score = 0
132 # min_score = 100
133 # best_seed = 0
134 # worst_seed = 0

136 # for i in range(1000):
137 #     X_train, X_test, y_train, y_test = get_data_set(i)
138 #     classifier = FisherClassifier(X_train, y_train)
139 #     score = classifier.score(X_test, y_test)
140 #     if score > max_score:
141 #         max_score = score
142 #         best_seed = i
143 #     if score < min_score:
144 #         min_score = score
145 #         worst_seed = i
146 #
147 # print('Best score for seed={}: {}'.format(best_seed, max_score))
148 # print('Worst score for seed={}: {}'.format(worst_seed, min_score))

150 X_train, X_test, y_train, y_test = get_data_set(879)
151 classifier = FisherClassifier(X_train, y_train)
152 score = classifier.score(X_test, y_test)
153 print('Score: {}'.format(classifier.score(X_test, y_test)))

155 lm = LinearRegression()
156 y_train_modified = list(map(lambda x: 1 if x == 1 else -1, y_train))
157 lm.fit(X_train, y_train_modified)
158 prediction = np.array(list(map(lambda x: 1 if x > 0 else 0, lm.predict(X_test)))), dtype=np.
    float64)
159 score = np.mean(prediction == np.array(y_test, dtype=np.float64))
160 print('Score with linear regression: {}'.format(score))

```

Code in Parser.py

```
1 import csv
2 import numpy as np
3 import os
4 from sklearn.model_selection import train_test_split

7 def parse_data():
8     file_name = os.path.join(os.path.dirname(__file__), './Dataset/spambase.data')
9     csv_file = open(file_name, 'rt')
10    reader = csv.reader(csv_file, delimiter=',', quoting=csv.QUOTE_NONE)
11    data = []

13    for row in reader:
14        filtered = list(filter(lambda x: x != '', row))
15        data.append(list(map(lambda x: float(x), filtered)))

17    return data

20 def get_points_and_labels_from_data(data):
21     points = np.array(list(map(lambda x: x[:-1], data)), dtype=np.float64)
22     labels = np.array(list(map(lambda x: int(x[-1]), data)))

24     return points, labels

27 def get_data_set(seed):
28     data = parse_data()
29     X, y = get_points_and_labels_from_data(data)
30     # for determined results we use a seed for random_state, so that data is always split
31     X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.8, test_size
32                                                         =0.2,
33                                                         random_state=seed)

34     return X_train, X_test, y_train, y_test
```