



Prof. R. Rojas

Mustererkennung, WS17/18

Übungsblatt 2

Boyan Hristov, Nedeltscho Petrov

31. Oktober 2017

Link zum Git Repository: <https://github.com/BoyanH/Freie-Universitaet-Berlin/tree/master/MachineLearning/Homework2>

Score

Das ist die Ausgabe des Programs und damit auch das Score

```
1 Score for 3 vs 7: 100.0%
2 Score for 3 vs 8: 100.0%
3 Score for 5 vs 7: 100.0%
4 Score for 5 vs 8: 100.0%
5 Score for 7 vs 8: 100.0%
```

Lineare Regression

Wir benutzen die folgende Formel für lineare Regression mit Least Squares

$$\vec{B} = (X^T X)^{-1} X^T \vec{y} \quad (1)$$

wobei \vec{B} die Lösung von der folgenden Gleichung ist: (im besten Fall, falls Matrix mit vollem Rank)

$$y_i = B_0 + B_i X_{i1} + B_i X_{i2} + \dots + B_n X_{in} \quad (2)$$

dabei ist X die Eingabematrix mit je Zeile ein Punkt im n-Dimensionalen Raum. Damit wir diese Gleichung lösen können, mit B_0 als 1. Glied, müssen wir erstmal eine Spalte mit Einsen einfügen.

```
1 ones = np.ones((len(self.trainData), 1), dtype=float)
2 X = np.append(ones, self.trainData, axis = 1)
```

Weiter, um die 2 Klassen gleichwertig in der linearen Regression zu betrachten, dürfen wir natürlich nicht gleich die Labels benutzen, sondern diese erstmal normalisieren (so zu sagen). Dabei mapen wir die eine Labels zu -1 und die andere zu 1. Z.B wenn wir zwischen 3 und 5 unterscheiden wollen, wird das je Label 3 zu -1 und je 5 zu 1. Das passiert folgendermaßen

```

1 def normalizeLabels(self, labels):
2     return list(map(lambda x: -1 if int(x) == self.classA else 1, labels))

```

Und damit unsere fit Methode

```

1 def fit(self):
2     ones = np.ones((len(self.trainData), 1), dtype=float)
3     X = np.append(ones, self.trainData, axis = 1)
4     xtxInversed = LinearRegressionClassifier.pseudoInverse(X.T.dot(X))
5     normalizedLabels = self.normalizeLabels(self.trainLabels)
6     self.beta = xtxInversed.dot(X.T).dot(normalizedLabels)

```

Pseudoinverse

Wir benutzen die Moore-Penrose Pseudoinverse. Nach der Formel gilt:

$$A^+ = \lim_{\delta \rightarrow 0} A^*(AA^* + \delta E)^{-1} \quad (3)$$

Wobei A^* die adjungierte Matrix ist. Da wir aber in dem Bereich der reellen Zahlen uns befinden, ist die adjungierte Matrix identisch zu der transponierten. Deswegen gilt:

$$A^+ = \lim_{\delta \rightarrow 0} A^T(AA^T + \delta E)^{-1} \quad (4)$$

Damit unsere Implementierung

```

1 @staticmethod
2 def pseudoInverse(X):
3     delta = np.nextafter(np.float16(0), np.float16(1)) # as close as we can get to lim delta
4     # -> 0
5     pseudoInverted = X.T.dot(np.linalg.inv(X.dot(X.T) + delta * np.identity(len(X))))

```

Vollständige Implementierung (Ohne Classifier Klasse)

```

1 import pandas as pd
2 import numpy as np
3 from operator import itemgetter
4 import matplotlib.pyplot as plt
5 import os
6 from Classifier import Classifier

7
8 import seaborn as sn
9 import pandas as pd
10 import matplotlib.pyplot as plt

11
12 class LinearRegressionClassifier(Classifier):

13
14
15     # B = (X^T X)^-1 X^T y
16
17     # hat H = X * B^ = ` X (X^T X)^-1 X^T
18
19     @staticmethod

```

```

20 def pseudoInverse(X):
21     # Determining whether a matrix is invertable or not with gaus elimination
22     # (most efficient from the naive approaches) takes  $n^3$  time, which
23     # makes the whole program slower for not much precision gain
24     # therefore we simply calculate the pseudo inversed matrix every time
25
26     # if LinearRegressionClassifier.isInvertable(X):
27     #     return np.linalg.inv(X)
28
29     # calculate pseudo-inverse  $A^+$  of a matrix  $A$  ( $X$  in our case)
30     #  $A^+ = \lim_{\delta \rightarrow 0} A(AA^* + \delta E)^{-1}$  where  $A^*$  is the conjugate transpose
31     # and  $E$  is the identity matrix
32
33     # In our case, we are working with real numbers, so  $A^* = A^T$ 
34     # so the formula is  $A^T(AA^T + \delta E)^{-1}$ 
35
36     delta = np.nextafter(np.float16(0), np.float16(1)) # as close as we can get to  $\lim_{\delta \rightarrow 0}$ 
37     pseudoInverted = X.T.dot(np.linalg.inv(X.dot(X.T) + delta * np.identity(len(X))))
38
39     return pseudoInverted
40
41 @staticmethod
42 def isInvertable(X):
43     # apply gaus elimination
44     # if the matrix is transformable in row-echelon form
45     # then it is as well invertable
46
47     X = np.copy(X) # don't really change given matrix
48     m = len(X)
49     n = len(X[0])
50     for k in range(min(m, n)):
51         # Find the k-th pivot:
52         #  $i_{\max} = \max(i = k \dots m, \text{abs}(A[i, k]))$ 
53         i_max = k
54         max_value = X[k][k]
55         for i in range(k, m):
56             if X[i][k] > max_value:
57                 max_value = X[i][k]
58                 i_max = i
59         if X[i_max][k] == 0:
60             return False
61         for i in range(n):
62             temp = X[k][i]
63             X[k][i] = X[i_max][i]
64             X[i_max][i] = temp
65         # Do for all rows below pivot:
66         for i in range(k+1, m):
67             # for  $i = k + 1 \dots m$ :
68             f = X[i][k] / X[k][k]
69             # Do for all remaining elements in current row:
70             for j in range(k + 1, n):
71                 X[i][j] = X[i][j] - (X[k][j] * f)
72             X[i][k] = 0
73
74     return True
75
76 def __init__(self, trainSet, testSet, classA, classB):
77     self.classA = classA
78     self.classB = classB
79
80     trainSet = self.filterDataSet(trainSet)
81     testSet = self.filterDataSet(testSet)
82
83     self.trainData = list(map(lambda x: x[1:], trainSet))
84     self.trainLabels = list(map(itemgetter(0), trainSet))
85     self.testSet = list(map(lambda x: x[1:], testSet))
86     self.testLabels = list(map(itemgetter(0), testSet))

```

```

88     self.fit()

90 def filterDataSet(self, dataSet):
91     return list(filter(lambda x: int(x[0]) in [self.classA, self.classB], dataSet))

93 def fit(self):
94     # fill X with (1,1...,1) in it's first column to be able to get the
95     # wished  $y_i = B_0 + B_1x_{i1} + B_2x_{i2} + \dots + B_nx_{in}$ 
96     ones = np.ones((len(self.trainData), 1), dtype=float)
97     X = np.append(ones, self.trainData, axis = 1)

99     # and then used the following formula to calculate our closest possible B
100    # which solves best our least squares regression
101    #  $B = (X^T X)^{-1} X^T y$ 
102    # where  $y = (-1, 1, -1, 1, \dots, -1)$  (for example) is a vector
103    # of the labels corresponding to the given data points

105    xtInversed = LinearRegressionClassifier.pseudoInverse(X.T.dot(X))

107    # normalize y, so the two possible classes are mapped to -1 or 1
108    normalizedLabels = self.normalizeLabels(self.trainLabels)
109    self.beta = xtInversed.dot(X.T).dot(normalizedLabels)

111 def predictSingle(self, X):
112     X = np.append(np.array([1]), np.array(X), axis=0)
113     return self.classA if (X.dot(self.beta) < 0) else self.classB

115 def predict(self, X):
116     return np.array(list(map(lambda x: self.predictSingle(x), X)))

118 def test(self):
119     print('Score for {} vs {}: {}'.format(
120         self.classA, self.classB, self.score(self.testSet, self.testLabels) * 100))
121     self.printConfusionsMatrix(self.confusion_matrix(self.testSet, self.testLabels))

124 def normalizeLabels(self, labels):
125     return list(map(lambda x: -1 if int(x) == self.classA else 1, labels))

127 def printConfusionsMatrix(self, matrix):
128     explicitImgPath = os.path.join(dir_path, './Plots/confusion_matrix_for_{0}vs_{0}.png'.
129         format(
130             self.classA, self.classB))
131     digits = [str(x) for x in range(10)];

132     df_cm = pd.DataFrame(matrix, index = digits,
133         columns = digits )

135     plt.figure(figsize = (11,7))
136     heatmap = sn.heatmap(df_cm, annot=True)

138     heatmap.set(xlabel='Klassifiziert', ylabel='Erwartet')

140     plt.savefig(explicitImgPath, format='png')

142 def extractDataFromLine(line):
143     line = line.replace(' \n', '') # clear final space and new line chars
144     return list(map(float, line.split(' '))); # map line to a list of floats

146 def parseDataFromFile(file, dataArr):
147     for line in file:
148         line = line.replace(' \n', '') # clear final space and new line chars
149         currentDigitData = extractDataFromLine(line)
150         dataArr.append(currentDigitData)

153 trainSet = []

```

```

154 testSet = []

156 dir_path = os.path.dirname(os.path.realpath(__file__))
157 explicitPathTrainData = os.path.join(dir_path, './Dataset/train')
158 explicitPathTestData = os.path.join(dir_path, './Dataset/test')
159 trainFile = open(explicitPathTestData, 'r')
160 testFile = open(explicitPathTestData, 'r')

162 parseDataFromFile(trainFile, trainSet)
163 parseDataFromFile(testFile, testSet)

165 LinearRegressionClassifier(trainSet, testSet, 3, 5).test()
166 LinearRegressionClassifier(trainSet, testSet, 3, 7).test()
167 LinearRegressionClassifier(trainSet, testSet, 3, 8).test()
168 LinearRegressionClassifier(trainSet, testSet, 5, 7).test()
169 LinearRegressionClassifier(trainSet, testSet, 5, 8).test()
170 LinearRegressionClassifier(trainSet, testSet, 7, 8).test()

```