

Klasifikacija besedil z jedrnimi metodami

Tomaž Tomažič (63100281)

14. junij 2015

1 Uvod

V nalogi sem uporabili SVM za klasifikacijo besedil da bi pri tem spoznal pomen jeder.

2 Podatki

Podatke smo si morali sami poiskati v treh kategorijah. Shranil sem bloge. Izbral sem si dve podobni kategoriji Windows blog (A) in Bing blog (B). Za drugacno temo pa sem izbral Zoella blog (C), ki je popularen blog za zenske, licenje in podobne reci.

3 Metode

Potrebno je bilo napisati jedro za podporne vektorje implementirane v sklearn knjižnjici. Jedro deluje na principu merjenja razdalje s kompresijo, kot je bilo podano v navodilih naloge. Razdalje sem obrnil tako da sem odstel vse vrednosti od največje v matriki razdalj. Napovedi sem naredil za vse tri kombinacije med kategorijami: A in B, A in C, B in C. Napoved sem naredil tako, da sem se naucil model na 80% podatkih, ostalih 20% pa sem uporabil za napoved verjetnosti razreda.

4 Rezultati

Rezultate sem preveril kar z metodo ostrega očesa, ko sem izpisal napovedane vrjetnosti. Rezultati so bili nekoliko boljsi ko sem napovedoval cisto razlicne razrede (A-C in B-C) v primerjavi z podobnimi razredi (A-B). Menim da bi bilo to mozno se precej izboljsati z iskanjem pravih parametrov.

5 Izjava o izdelavi domače naloge

Domačo nalogo in pripadajoče programe sem izdelal sam.

Priloge

A Programska koda

```
1 import zlib
2 import copy
3 import numpy as np
4 from sklearn.svm import SVC
5 from sklearn.metrics import roc_auc_score, accuracy_score, auc
6
7 class SVM:
8     def __init__(self, data):
9         self.learner = SVC(C = 1.1, kernel=self.string_kernel, probability=True)
10
11         self.data = data
12
13     def Z(self, a):
14         return len(zlib.compress(a.encode()))
15
16     def distance(self, a, b):
17         Z = self.Z
18         a = self.data[int(a)]
19         b = self.data[int(b)]
20
21         ab = a+b
22         ba = b+a
23         return (Z(ab) - Z(a)) / Z(a) + (Z(ba) - Z(b)) / Z(b)
24
25     def string_kernel(self, X, Y):
26         M = np.zeros((len(X), len(Y)))
27         for i, a in enumerate(X):
28             for j, b in enumerate(Y):
29                 M[i, j] = self.distance(a[0], b[0])
30         return M.max() - M
31
32     def fit(self, X, Y):
33         return self.learner.fit(X, Y)
34
35     def predict(self, X):
36         return self.learner.predict_proba(X)
37
38 data_a = [ open("source/windows/"+str(i)+".txt", "rt").read() for i in range(1,
39             21)]
40 data_b = [ open("source/bing/"+str(i)+".txt", "rt").read() for i in range(1,
41             21)]
42 data_c = [ open("source/zoella/"+str(i)+".txt", "rt").read() for i in range(1,
43             21)]
44
45 def test(data_a, data_b):
46     da = data_a[:-4]
47     db = data_b[:-4]
48     data_ab = copy.deepcopy(da)
49     data_ab.extend(db)
50
51     Xab = np.vstack(np.r_[0:len(data_ab)])
52     Yab = np.array([0]*len(da) + [1]*len(db))
53
54     s = SVM(data_ab)
```

```
51     s.fit(Xab, Yab)
    s.data.extend(data_a[-4:])
53     s.data.extend(data_b[-4:])

55     return s.predict(np.vstack(np.r_[len(da)*2:len(data_a)*2]))

57 print(test(data_a, data_b))
print(test(data_a, data_c))
59 print(test(data_b, data_c))
```