

# Data Mining

## Home work 09

### Machine Learning Part:2

Aqeel Labash  
**Lecturer:** Jaak Vilo

06 April 2016

## First Question

The article focus on classification. And here is the list

- No matter what algorithm you pick it's consist from three parts : Representation, Evaluation, Optimization.
- **Representation:** The feature we want to use in away computer can handle.
- **Evaluation or objective or scoring function:** used to distinguish between good and bad classifiers.
- **Optimization :** a method that search for best scoring classifiers.
- Generalization is what we want. "if there are 100,000 words in the dictionary, the spam filter described above has 2100,000 possible different inputs."
- Always test your module on data different than the train data.
- data by itself is not enough for generalization knowledge is also required to know which module to apply." A corollary of this is that one of the key criteria for choosing a representation is which kinds of knowledge are easily expressed in it. For example, if we have a lot of knowledge about what makes examples similar in our domain, instance-based methods may be a good choice. If we have knowledge about probabilistic dependencies, graphical models are a good fit. And if we have knowledge about what kinds of preconditions are required by each class, "IF . . . THEN . . ." rules may be the best option."
- **Over fitting** could be decomposed to Bias and Variance." Bias is a learner's tendency to consistently learn the same wrong thing. Variance is the tendency to learn random things irrespective of the real signal."
- Some ways to fight over fitting is: regularization term, cross validation, statistical significance.
- In most cases over fitting doesn't happen because of noise.
- intuitions used for low dimension usually don't work with high dimension problems." our intuitions, which come from a three-dimensional world, often do not apply in high-dimensional ones."
- Numbers in theory might not be applicable and not always correct." consider the space of Boolean functions of  $d$  Boolean variables. If there are  $e$  possible different examples, there are  $2^e$  possible different functions, so since there are  $2^d$  possible examples, the total number of functions is  $2^{2^d}$ . And even for hypothesis spaces that are "merely" exponential, the bound is still very loose, because the union bound is very pessimistic. For example, if there are 100 Boolean features and the hypothesis space is decision trees with up to 10 levels, to guarantee  $\delta = \epsilon = 1\%$  in the bound above we need half a million examples. But in practice a small fraction of this suffices for accurate learning."
- **Good Features:** are those which independent from each other and highly correlated with the class.
- Awesome feature is not necessarily provided by the data we might have to build it." Often, the raw data is not in a form that is amenable to learning, but you can construct features from it that are."
- The algorithm with more data wins." As a rule of thumb, a dumb algorithm with lots and lots of data beats a clever one with modest amounts of it. (After all, machine learning is all about letting data do the heavy lifting.)"

- The more data we have, the more complex the classifier.
- "Variable size learners can in principle learn any function given sufficient data, but in practice they may not, because of limitations of the algorithm (for example, greedy search falls into local optima) or computational cost."
- Ensemble learning give better results usually.
- if the classifier is simpler doesn't mean it's more accurate.
- being able to represent the data doesn't mean that module can learn it."For example, standard decision tree learners cannot learn trees with more leaves than there are training examples."
- correlation doesn't mean one cause another.it's just an observation on the data that might have cause relation.

## Second Question

For this task I used the following Code :

```

1
2 import matplotlib.pyplot as plt
3
4
5 # In [2]:
6
7 original = {}
8 TotalPositive = 0
9 TotalNegative = 0
10 with open('data.class', 'r') as f:
11     f = f.readlines()
12     for line in f:
13         line = line.split()
14         if line[1] == 'T':
15             original[line[0]] = True
16             TotalPositive += 1
17         else:
18             original[line[0]] = False
19             TotalNegative += 1
20     roc1 = []
21     with open('roc1.txt', 'r') as f:
22         f = f.readlines()
23         for line in f:
24             roc1.append(line.strip())
25     roc2 = []
26     with open('roc2.txt', 'r') as f:
27         f = f.readlines()
28         for line in f:
29             roc2.append(line.strip())
30     roc3 = []
31     with open('roc3.txt', 'r') as f:
32         f = f.readlines()
33         for line in f:
34             roc3.append(line.strip())
35     roc4 = []
36     with open('roc4.txt', 'r') as f:
37         f = f.readlines()
38         for line in f:
39             roc4.append(line.strip())
40     rocperfect = []
41     for x in original.keys():
42         if original[x]:
43             rocperfect = [x] + rocperfect
44         else:
45             rocperfect.append(x)
46
47
48 # In [3]:
49
50 def GetTPFP(k, dataset):
51     TP = 0
52     FP = 0
53     TN = 0
54     FN = 0

```

```

55 for i in range (3000):
56 #Identified True
57 if i<k:
58 TP+=original[dataset[i]]
59 else:
60 TN+=original[dataset[i]]
61 #(TP,FP,TN,FN)
62 #(F11,F01,F10,F00)
63 #return (TP,TotalPositive-TP,TN,TotalNegative-TN)
64 return (float (TP)/float (TotalPositive),float (k-TP)/float (TotalNegative))
65
66
67
68 # In [4]:
69
70 roc1cm=[]
71 roc2cm=[]
72 roc3cm=[]
73 roc4cm=[]
74 rocperfectcm=[]
75 print 'processing roc1'
76 for i in range (3000):
77 roc1cm.append(GetTPFP(i,roc1))
78 print 'processing roc2'
79 for i in range (3000):
80 roc2cm.append(GetTPFP(i,roc2))
81 print 'processing roc3'
82 for i in range (3000):
83 roc3cm.append(GetTPFP(i,roc3))
84 print 'processing roc4'
85 for i in range (3000):
86 roc4cm.append(GetTPFP(i,roc4))
87 print 'processing rocperfect'
88 for i in range (3000):
89 rocperfectcm.append(GetTPFP(i,rocperfect))
90 print 'Finished'
91
92
93 # In [5]:
94
95 plt.figure('roc1.jpg')
96 plt.plot([x[1] for x in roc1cm],[x[0] for x in roc1cm], 'ro')
97 plt.ylabel('TPR')
98 plt.xlabel('FPR')
99 plt.title('Roc1')
100 #plt.show()
101 plt.savefig('roc1.jpg')
102
103
104 # In [6]:
105
106 plt.figure('roc2.jpg')
107 plt.plot([x[1] for x in roc2cm],[x[0] for x in roc2cm], 'ro')
108 plt.ylabel('TPR')
109 plt.xlabel('FPR')
110 plt.title('Roc2')
111 #plt.show()
112 plt.savefig('roc2.jpg')
113
114
115 # In [7]:
116
117 plt.figure('roc3.jpg')
118 plt.plot([x[1] for x in roc3cm],[x[0] for x in roc3cm], 'ro')
119 plt.ylabel('TPR')
120 plt.xlabel('FPR')
121 plt.title('Roc3')
122 #plt.show()
123 plt.savefig('roc3.jpg')
124
125
126 # In [8]:
127
128 plt.figure('roc4.jpg')
129 plt.plot([x[1] for x in roc4cm],[x[0] for x in roc4cm], 'ro')
130 plt.ylabel('TPR')

```

```

131 plt.xlabel('FPR')
132 plt.title('Roc4')
133 #plt.show()
134 plt.savefig('roc4.jpg')
135
136
137 # In[9]:
138
139 plt.figure('rocperfect.jpg')
140 plt.plot([x[1] for x in rocperfectcm],[x[0] for x in rocperfectcm], 'ro')
141 plt.ylabel('TPR')
142 plt.xlabel('FPR')
143 plt.title('Rocperfect')
144 #plt.show()
145 plt.savefig('rocperfect.jpg')
146
147
148 # In[10]:
149
150 AUCROC1=0
151 AUCROC2=0
152 AUCROC3=0
153 AUCROC4=0
154 AUCROCPERFECT=0
155
156 for i in range(1,3000):
157 AUCROC1+= (roc1cm[i][1] - roc1cm[i-1][1])*roc1cm[i][0]
158 AUCROC2+= (roc2cm[i][1] - roc2cm[i-1][1])*roc2cm[i][0]
159 AUCROC3+= (roc3cm[i][1] - roc3cm[i-1][1])*roc3cm[i][0]
160 AUCROC4+= (roc4cm[i][1] - roc4cm[i-1][1])*roc4cm[i][0]
161 AUCROCPERFECT+= (rocperfectcm[i][1] - rocperfectcm[i-1][1])*rocperfectcm[i][0]
162
163 print AUCROC1,AUCROC2,AUCROC3,AUCROC4,AUCROCPERFECT

```

What I did in the previous is simply calculated TP and from it I calculated TPR & FPR.

$$TPR = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{FP + TN} = 1 - specificity$$

the confusion matrix and from it I calculated TPR,FPR then draw them.And here is the figures :

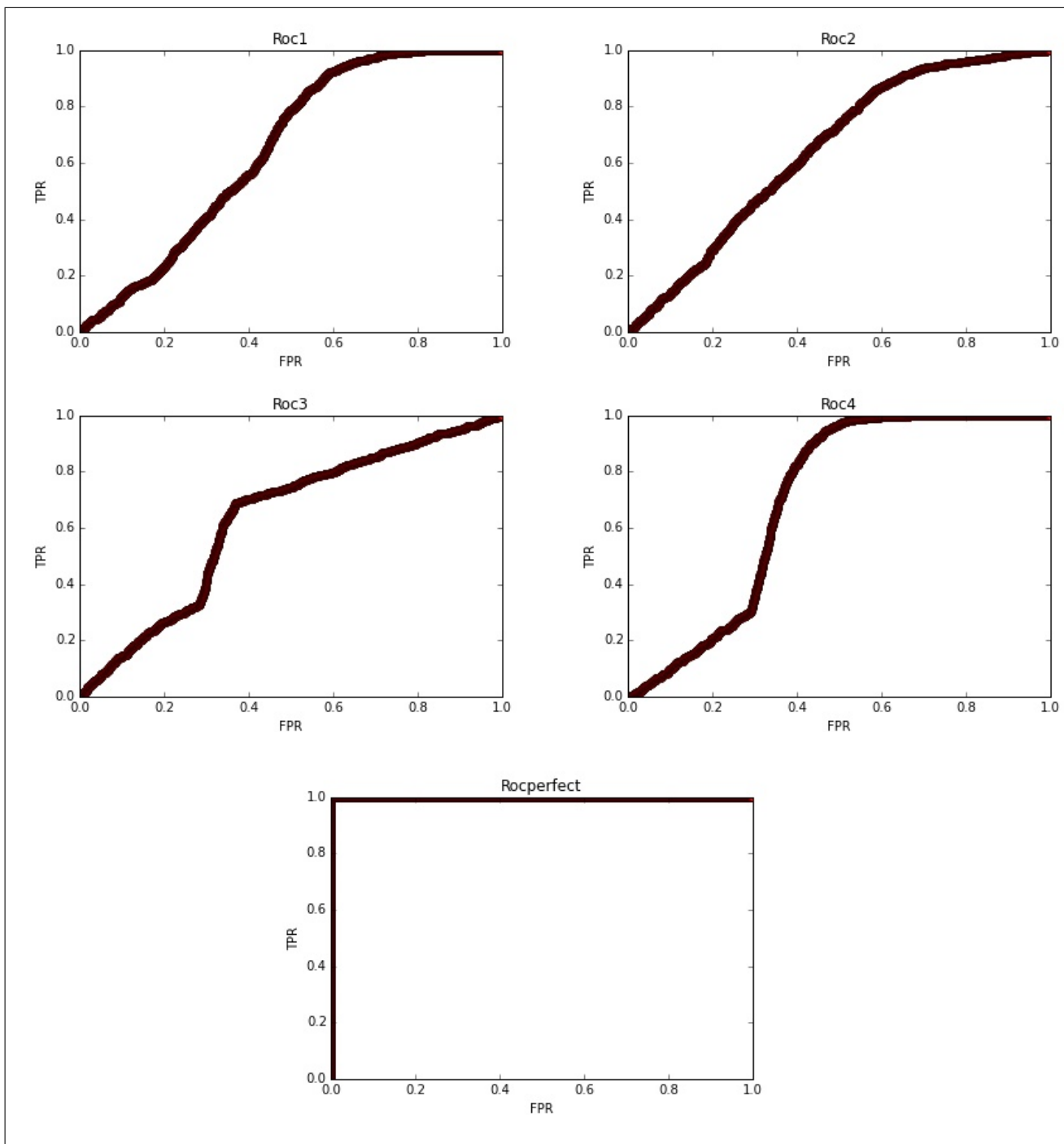


Figure 1: ROC for all the predictions.

We can notice from the previous figures that Roc4 give a good prediction for TP and minimum FP compared to others.

**Note:**I added a plot for the perfect file how the perfect classifier would look like. The following code calculate approximation of the area under curve. We calculate the sum of all rectangles (point by point) and so on.

```

1 AUCROC1=0
2 AUCROC2=0
3 AUCROC3=0
4 AUCROC4=0
5
6 for i in range(1,3000):
7     AUCROC1+= (roc1cm[i][1] - roc1cm[i-1][1]) * roc1cm[i][0]
8     AUCROC2+= (roc2cm[i][1] - roc2cm[i-1][1]) * roc2cm[i][0]
9     AUCROC3+= (roc3cm[i][1] - roc3cm[i-1][1]) * roc3cm[i][0]
10    AUCROC4+= (roc4cm[i][1] - roc4cm[i-1][1]) * roc4cm[i][0]
11 print AUCROC1,AUCROC2,AUCROC3,AUCROC4

```

The previous code give the following result :

$$AUC1 = 0.650940277346, AUC2 = 0.647270924831, AUC3 = 0.629960231006, AUC4 = 0.696844993141$$

## Third Question

To charactraize the previous plots.I would explain what is happening.When we start we predict 0 or 1 positive so  $TP+FP=1$  at best cases, while  $k$  is growing  $TP+FP$  is growing as well.At the end we predict all the data as True.In that point we have  $TPR = 1, FPR = 1$ .

The previous figures For this task I used the hint for this question so I subtracted FPR from TPR and then tried to find the max value. I used the following code :

```
1 def FindK(dataset):
2     k=0
3     lst=[]
4     for i in range(3000):
5         lst.append({'x':i, 'y':dataset[i][0]-dataset[i][1]})
6
7     lst.sort(key=lambda x: x['y'], reverse=True)
8     return lst[0]['x']
9     print FindK(roc1cm),FindK(roc2cm),FindK(roc3cm),FindK(roc4cm)
10    roc1best = FindK(roc1cm)
11    roc2best = FindK(roc2cm)
12    roc3best = FindK(roc3cm)
13    roc4best = FindK(roc4cm)
14    rocperfectbest= FindK(rocperfectcm)
15    print roc1best,roc2best,roc3best,roc4best,rocperfectbest
```

the previous code output the following data :

$$roc1best : 2179, roc2best : 2089, roc3best : 1500, roc4best : 1996, rocperfectbest : 1215$$

To get better view I drawn the data and pointed out those points using the following code:

```
1 plt.figure('roc1_best.jpg')
2 plt.plot([x[1] for x in roc1cm],[x[0] for x in roc1cm], 'ro')
3 plt.ylabel('TPR')
4 plt.xlabel('FPR')
5 plt.title('Roc1 Best Point')
6 plt.axvline(x=roc1cm[roc1best][1])
7 plt.axhline(y=roc1cm[roc1best][0])
8 #plt.show()
9 plt.savefig('roc1_best.jpg')
10
11
12 # In [24]:
13
14 plt.figure('roc2_best.jpg')
15 plt.plot([x[1] for x in roc2cm],[x[0] for x in roc2cm], 'ro')
16 plt.ylabel('TPR')
17 plt.xlabel('FPR')
18 plt.title('Roc2 Best Point')
19 plt.axvline(x=roc2cm[roc2best][1])
20 plt.axhline(y=roc2cm[roc2best][0])
21 #plt.show()
22 plt.savefig('roc2_best.jpg')
23
24
25 # In [25]:
26
27 plt.figure('roc3_best.jpg')
28 plt.plot([x[1] for x in roc3cm],[x[0] for x in roc3cm], 'ro')
29 plt.ylabel('TPR')
30 plt.xlabel('FPR')
31 plt.title('Roc3 Best Point')
32 plt.axvline(x=roc3cm[roc3best][1])
33 plt.axhline(y=roc3cm[roc3best][0])
34 #plt.show()
35 plt.savefig('roc3_best.jpg')
36
37
38 # In [27]:
39
40 plt.figure('roc4_best.jpg')
```

```

41 plt.plot([x[1] for x in roc4cm],[x[0] for x in roc4cm], 'ro')
42 plt.ylabel('TPR')
43 plt.xlabel('FPR')
44 plt.title('Roc4 Best Point')
45 plt.axvline(x=roc4cm[roc4best][1])
46 plt.axhline(y=roc4cm[roc4best][0])
47 #plt.show()
48 plt.savefig('roc4_best.jpg')
49
50
51 # In [28]:
52
53 plt.figure('rocperfect_best.jpg')
54 plt.plot([x[1] for x in rocperfectcm],[x[0] for x in rocperfectcm], 'ro')
55 plt.ylabel('TPR')
56 plt.xlabel('FPR')
57 plt.title('Rocperfect Best Point')
58 plt.axvline(x=rocperfectcm[rocperfectbest][1])
59 plt.axhline(y=rocperfectcm[rocperfectbest][0])
60 #plt.show()
61 plt.savefig('rocperfect_best.jpg')

```

The previous code output the following figures :

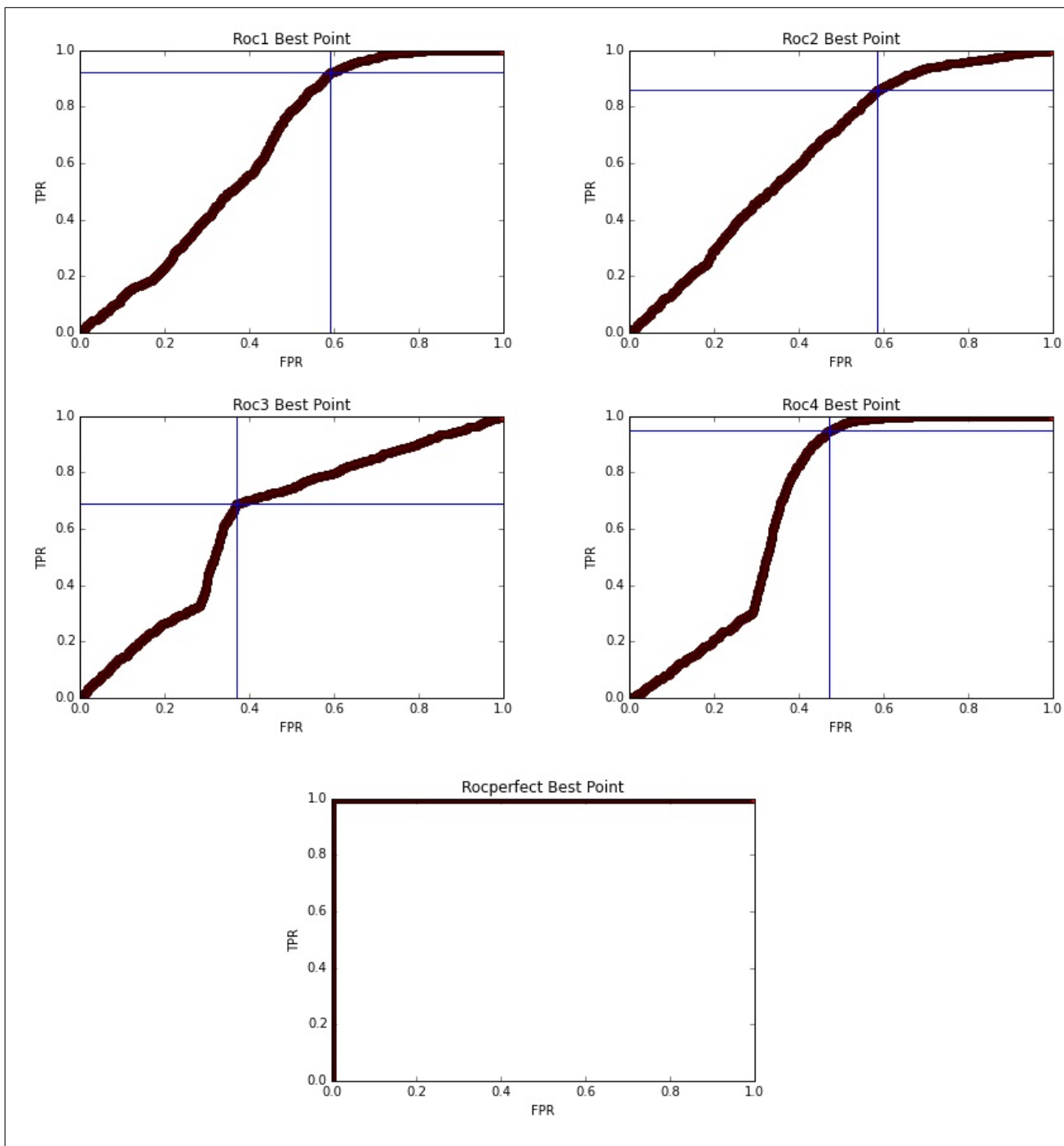


Figure 2: ROC for all the predictions pointing the best split point.

From the previous figure we can see that at the "Rocperfect" we can't see the blue lines that's because they overlap with the data itself.

## Fourth Question

## Fifth Question

## Sixth Question

## References

[1] receiver operating Characteristics