**Mohammad Saleh**
**GRA at PHY WSU**
**m.saleh@cern.ch**
**fa9520@wayne.edu**
**Jan 25, 2017**

**CSC 5825**
**Homework 1**

# Solutions:

**Question 1:**

**a.** Let start the comparison by looking at the advantages of OCR over scanning and faxing it: – The OCR will be able to save the document as a characters (A, a, b). This will make the documents editable. While scanning will give you an image only and you cannot change any text in the document. So if you looking to edit or copy the text the OCR will be more preferable in this case.

– The disadvantage of using OCR is the accuracy. As the document may include different fonts and hand writing, the OCR will not be accurate where the scan will show the accurate document. So if the details in the document is the important concern then the scan will be more preferable.

– May be one other advantage of the optical reader is the size of the machine as the OCR can be smaller than the scanner but it will actually depend on the model of the scanner and the OCR. In case your are constrained to small area, probably the optical reader will be more preferable.

**b.** The OCR will fail in many situation: As we mentioned before many fonts are available and some of them can be cursive which will be very hard for the OCR to e able to detect the characters. Also the OCR is matching the character with a template pixel by pixel. If we have our character with low resolution (small number of pixel), the OCR will fail as it will not have enough pixel information.

The barcode reader is still used because it very reliable as the barcode is a combination of line and the barcode reader will not fail to estimate the difference in the distance between whose line. Also the technology used in the barcode is simple and efficient and money wise it is cheaper than the OCR. So in case you need to tag a merchandise you don't need to assign characters to be read by OCR later on. The barcode will be enough to tag the merchandise.

**Question 2:**

According to the equation of the circle eq.1, we have the radius and the center coordinates parameters that define any circle. The parameters for the hypothesis circle will be in between the most specified and the most general hypothesis in this sense will be only talking about changing the radius of the circle. The ellipse will be more useful as in the circle case we are constrained with the size of the radius and the variance across the x-axis and the y-axis will be the same where in the ellipse you have two different sizes of the axes (major and minor) and you can see the correlation between the x-axis and the y-axis more clearly. In addition to the major and minor axes we can also use a rotation angle for ellipse to include desired point (as example figure1) rather than the circle which will not have difference if you do rotation (isotropic). To generalize for K>2, we can use N binary hypothesis, with problems less than using triangles as using ellipse we will have more control on the K class using the two-axes and rotation angle.

$$(x - x_0)^2 + (y - y_0)^2 = R^2 \tag{1}$$

**Question 3:**

The answer for this question depends mainly on the application but we can discuss the different possible cases. S is the closest to the positive (doesn't make false positive) and G is close to the negative (doesn't make false negative). In the case where false positive and false negative have the same cost, we can choose h halfway between S and G. Whereas if false positive costs more than false negative, we choose h closer to S. Also, if false negative costs more than false positive, we choose h closer to G. As we said early it is application dependent.
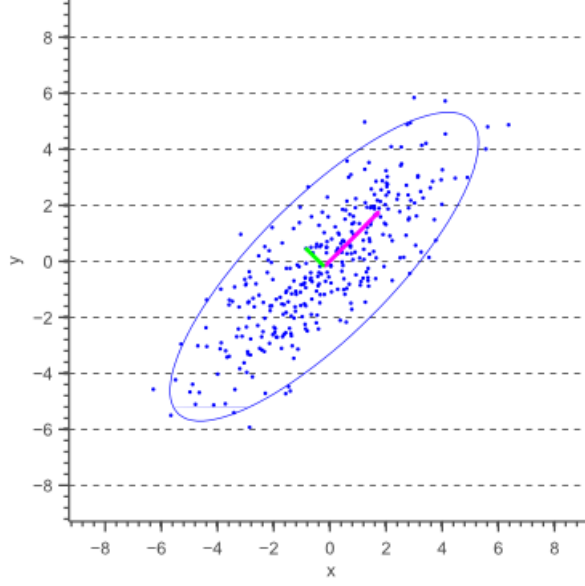
Figure 1: an example of particle correlation represented by an ellipse rotated by some angle.

**Question 4:**

For given training set, we set S, most specified hypothesis, and G, most generalized hypothesis. Then we choose our h between S and G but we are not sure if h closer to S or G. SO the region between S and G depend on our intuitive. So the best queries by supervisor will be for instances between S and Q. In this case the region between S and G will be smaller as more instances will be available and we will be able to set larger boundaries on S and G.

**Question 5:**

**a.** The likelihood ratio is defined as $p(\mathbf{x}|C_i)/p(\mathbf{x}|C_j)$, The discriminant function discriminates between two classes, defined as the ratio of class conditional density using Baye's rule: $g(\mathbf{x}) = P(C_i|x)/P(C_j|x)$, i$\neq$j, where it will choose $C_i$ if $g(\mathbf{x}) > 1$ and $C_j$ otherwise. Also from Baye's rule: $P(C_k|x) = p(x|C_k) * P(C_k)/p(x)$. Implies that $g(x) = \frac{p(x|C_i)}{p(x|C_j)} * \frac{P(C_i)}{P(C_j)}$, which include the liklihood ratio.

**b.** The same analogy here where the log odds ratio is defined as $\log \frac{P(C_i|x)}{P(C_j|x)}$. Also the discriminating function will be log formulated, $g(x) = \log \frac{P(C_i|x)}{P(C_j|x)}$, which choose $C_i$ if $g(\mathbf{x}) > 0$ (since log(1)=0) and $C_j$ otherwise. Also from Baye's rule: $P(C_k|x) = p(x|C_k) * P(C_k)/p(x)$. Implies that $g(x) = \log \frac{p(x|C_i)}{p(x|C_j)} + \log \frac{P(C_i)}{P(C_j)}$
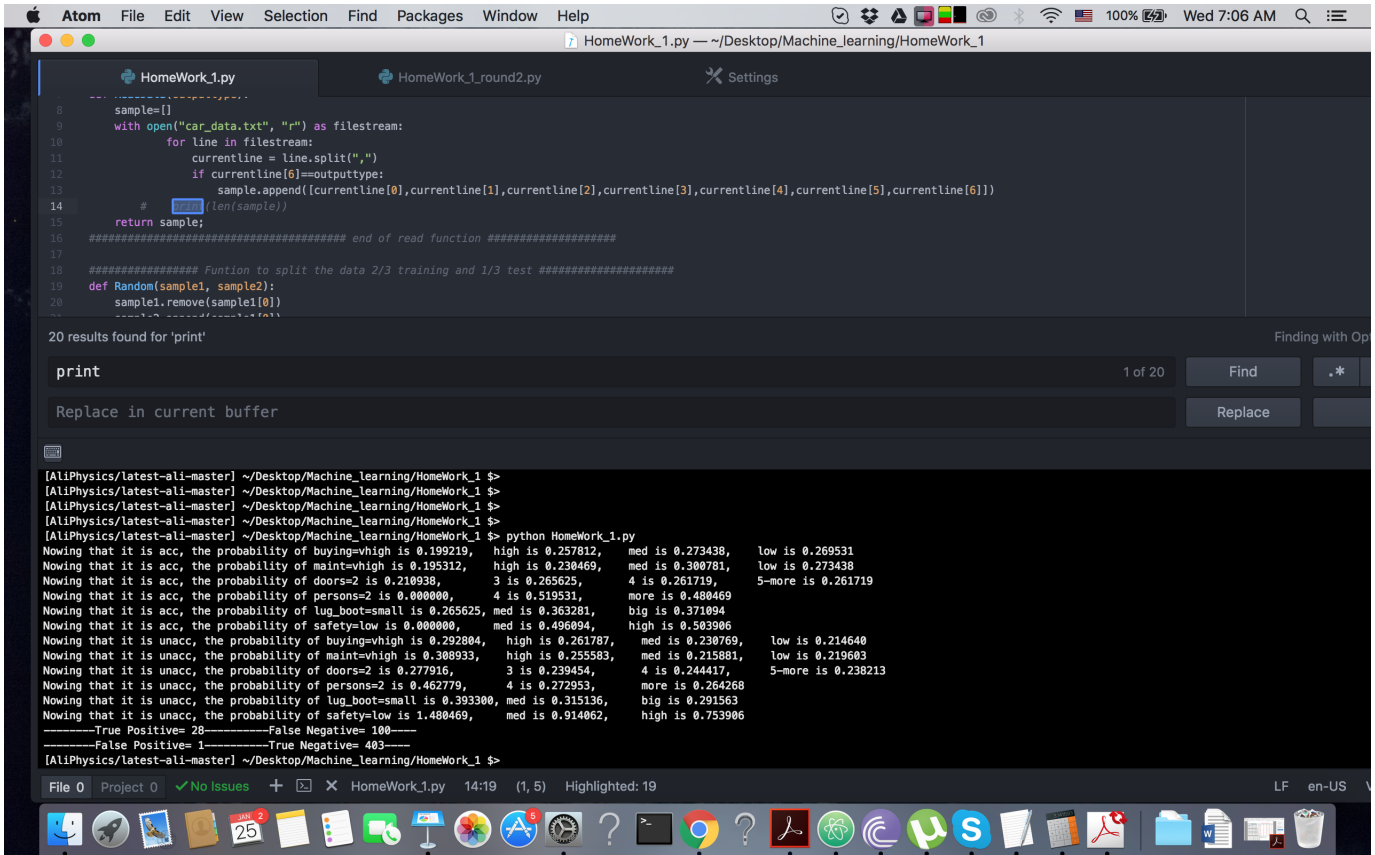
Figure 2: snapshot of my code after compiling. It shows all the requested parameters for the homework.

**Question 6:**

Below is the syntax for all the parts (a, b, c, d). I did add comments in order to make it more clear. I did also include a snapshot for the compiled code. Sorry incase of not perfectly organised code. I am mainly using C++ in my research but I wrote the syntax in PYTHON to get familiar with the syntax (of course same logic as in C++).

```python
import random
import math
from decimal import *
getcontext().prec = 6

################### Funtion to read the data from the text file ######################
def ReadData(outputtype):
    sample=[]
    with open("car_data.txt", "r") as filestream:
            for line in filestream:
                    currentline = line.split(",")
                    if currentline[6]==outputtype:
                            sample.append([currentline[0],currentline[1],currentline[2],currentline[3],
        #       print(len(sample))
    return sample;
######################################### end of read function #####################

################## Funtion to split the data 2/3 training and 1/3 test ####################
def Random(sample1, sample2):
```

3

```
87         sample1.remove(sample1[0])
88         sample2.append(sample1[0])
89         while (float(len(sample1))/float(len(sample2)))>2.0:
90             import random
91             random=random.randint(0,len(sample1)-1)
92     #         print(random)
93             sample2.append(sample1[random])
94             sample1.remove(sample1[random])
95     #         print(len(sample2))
96         return;
97     ################ end of random function  ####################
98
99     ################ Funtion to optain the conditional Probability using Naive Bayes classifier ##
100    ################## We could done this function more general and loop based but for now it is m
101    def Byes(sample,sampleoutput,Buying,Maint,Doors,Persons,Boot,Safety):
102        for i in range(0,len(sample)):
103            if(sample[i][6]==sampleoutput):
104                if(sample[i][0]=='vhigh'):
105                    Buying[0]+=1.0
106                if(sample[i][0]=='high'):
107                    Buying[1]+=1.0
108                if(sample[i][0]=='med'):
109                    Buying[2]+=1.0
110                if(sample[i][0]=='low'):
111                    Buying[3]+=1.0
112                if(sample[i][1]=='vhigh'):
113                    Maint[0]+=1.0
114                if(sample[i][1]=='high'):
115                    Maint[1]+=1.0
116                if(sample[i][1]=='med'):
117                    Maint[2]+=1.0
118                if(sample[i][1]=='low'):
119                    Maint[3]+=1.0
120                if(sample[i][2]=='2'):
121                    Doors[0]+=1.0
122                if(sample[i][2]=='3'):
123                    Doors[1]+=1.0
124                if(sample[i][2]=='4'):
125                    Doors[2]+=1.0
126                if(sample[i][2]=='5more'):
127                    Doors[3]+=1.0
128                if(sample[i][3]=='2'):
129                    Persons[0]+=1.0
130                if(sample[i][3]=='4'):
131                    Persons[1]+=1.0
132                if(sample[i][3]=='more'):
133                    Persons[2]+=1.0
134                if(sample[i][4]=='small'):
135                    Boot[0]+=1.0
136                if(sample[i][4]=='med'):
137                    Boot[1]+=1.0
138                if(sample[i][4]=='big'):
139                    Boot[2]+=1.0
140                if(sample[i][5]=='low'):
```

```
141                    Safety[0]+=1.0
142              if(sample[i][5]=='med'):
143                    Safety[1]+=1.0
144              if(sample[i][5]=='high'):
145                    Safety[2]+=1.0
146      return;
147  ################# end of Naive Bayes classifier function  #####################
148
149
150
151
152  testsample=[]
153  testsample1=[]
154  buyingProb=[0]*4
155  maintProb=[0]*4
156  doorsProb=[0]*4
157  personProb=[0]*3
158  bootProb=[0]*3
159  safetyProb=[0]*3
160  buyingProb1=[0]*4
161  maintProb1=[0]*4
162  doorsProb1=[0]*4
163  personProb1=[0]*3
164  bootProb1=[0]*3
165  safetyProb1=[0]*3
166  multiply=[0]*6
167  multiply1=[0]*6
168  classifier=[0]*4
169  accsample=ReadData("acc\n")    ###### store the acc output cases in accsample
170  unccsample=ReadData("unacc\n") ###### store the unacc output cases in unccsample
171  Random(accsample,testsample1)    ###### random 2/3 training set stored in accsample while test s
172  Random(unccsample,testsample)    ###### random 2/3 training set stored in unccsample while test s
173  testsample.extend(testsample1)    ###### merging the training and testing sample with the same p
174  unccsample.extend(accsample)
175
176  Byes(unccsample,"acc\n",buyingProb1,maintProb1,doorsProb1,personProb1,bootProb1,safetyProb1)
177  ######### trainning the classifier (+ive case)
178  Byes(unccsample,"unacc\n",buyingProb,maintProb,doorsProb,personProb,bootProb,safetyProb)
179  ######### trainning the classifier (-ive case)
180
181  x=sum(buyingProb)
182  y=sum(buyingProb1)
183
184  for i in range(0,4):  ######### normalizing the probabilities
185      buyingProb[i]=buyingProb[i]/x
186      maintProb[i]=maintProb[i]/x
187      doorsProb[i]=doorsProb[i]/x
188      buyingProb1[i]=buyingProb1[i]/y
189      maintProb1[i]=maintProb1[i]/y
190      doorsProb1[i]=doorsProb1[i]/y
191  for i in range(0,3):
192      personProb[i]=personProb[i]/x
193      bootProb[i]=bootProb[i]/x
194      safetyProb[i]=safetyProb[i]/y
```

5

```
195        personProb1[i]=personProb1[i]/y
196        bootProb1[i]=bootProb1[i]/y
197        safetyProb1[i]=safetyProb1[i]/y
198    #
199    ######### Printing all the parameters from training as requested
200    samplekind1='acc'
201    print"Nowing_that_it_is_%s,_the_probability_of_buying=vhigh_is_%f,___high_is_%f,____med_is_%f,_
202    print"Nowing_that_it_is_%s,_the_probability_of_maint=vhigh_is_%f,___high_is_%f,____med_is_%f,_
203    print"Nowing_that_it_is_%s,_the_probability_of_doors=2_is_%f,_____3_is_%f,_____4_is_%f,___
204    print"Nowing_that_it_is_%s,_the_probability_of_persons=2_is_%f,_____4_is_%f,_____more_is_%f"
205    print"Nowing_that_it_is_%s,_the_probability_of_lug_boot=small_is_%f,_med_is_%f,_____big_is_%f"%
206    print"Nowing_that_it_is_%s,_the_probability_of_safety=low_is_%f,_____med_is_%f,_____high_is_%f"
207
208    personProb1[0]=personProb1[0]+0.2      ######### TO avoid zero probability  ##################
209    safetyProb1[0]=safetyProb1[0]+0.2       ######### TO avoid zero probability  ##################
210
211    samplekind='unacc'
212    print"Nowing_that_it_is_%s,_the_probability_of_buying=vhigh_is_%f,___high_is_%f,____med_is_%f,_
213    print"Nowing_that_it_is_%s,_the_probability_of_maint=vhigh_is_%f,___high_is_%f,____med_is_%f,_
214    print"Nowing_that_it_is_%s,_the_probability_of_doors=2_is_%f,_____3_is_%f,_____4_is_%f,___
215    print"Nowing_that_it_is_%s,_the_probability_of_persons=2_is_%f,_____4_is_%f,_____more_is_%f"
216    print"Nowing_that_it_is_%s,_the_probability_of_lug_boot=small_is_%f,_med_is_%f,_____big_is_%f"%
217    print"Nowing_that_it_is_%s,_the_probability_of_safety=low_is_%f,_____med_is_%f,_____high_is_%f"
218
219
220
221    ############### validation check and TP TN FP FN using sample tests ####################
222    for i in range(0,len(testsample)):
223        if testsample[i][0]=='vhigh':
224            multiply[0]=buyingProb[0]
225            multiply1[0]=buyingProb1[0]
226        if testsample[i][0]=='high':
227            multiply[0]=buyingProb[1]
228            multiply1[0]=buyingProb1[1]
229        if testsample[i][0]=='med':
230            multiply[0]=buyingProb[2]
231            multiply1[0]=buyingProb1[2]
232        if testsample[i][0]=='low':
233            multiply[0]=buyingProb[3]
234            multiply1[0]=buyingProb1[3]
235        if testsample[i][1]=='vhigh':
236            multiply[1]=maintProb[0]
237            multiply1[1]=maintProb1[0]
238        if testsample[i][1]=='high':
239            multiply[1]=maintProb[1]
240            multiply1[1]=maintProb1[1]
241        if testsample[i][1]=='med':
242            multiply[1]=maintProb[2]
243            multiply1[1]=maintProb1[2]
244        if testsample[i][1]=='low':
245            multiply[1]=maintProb[3]
246            multiply1[1]=maintProb1[3]
247        if testsample[i][2]=='2':
248            multiply[2]=doorsProb[0]
```

```python
249             multiply1[2]=doorsProb1[0]
250         if testsample[i][2]=='3':
251             multiply[2]=doorsProb[1]
252             multiply1[2]=doorsProb1[1]
253         if testsample[i][2]=='4':
254             multiply[2]=doorsProb[2]
255             multiply1[2]=doorsProb1[2]
256         if testsample[i][2]=='5more':
257             multiply[2]=doorsProb[3]
258             multiply1[2]=doorsProb1[3]
259         if testsample[i][3]=='2':
260             multiply[3]=personProb[0]
261             multiply1[3]=personProb1[0]
262         if testsample[i][3]=='4':
263             multiply[3]=personProb[1]
264             multiply1[3]=personProb1[1]
265         if testsample[i][3]=='more':
266             multiply[3]=personProb[2]
267             multiply1[3]=personProb1[2]
268         if testsample[i][4]=='small':
269             multiply[4]=bootProb[0]
270             multiply1[4]=bootProb1[0]
271         if testsample[i][4]=='med':
272             multiply[4]=bootProb[1]
273             multiply1[4]=bootProb1[1]
274         if testsample[i][4]=='big':
275             multiply[4]=bootProb[2]
276             multiply1[4]=bootProb1[2]
277         if testsample[i][5]=='low':
278             multiply[5]=bootProb[0]
279             multiply1[5]=bootProb1[0]
280         if testsample[i][5]=='med':
281             multiply[5]=bootProb[1]
282             multiply1[5]=bootProb1[1]
283         if testsample[i][5]=='high':
284             multiply[5]=bootProb[2]
285             multiply1[5]=bootProb1[2]
286         prob=multiply[0]*multiply[1]*multiply[2]*multiply[3]*multiply[4]*multiply[5]*x/1062
287         prob1=multiply1[0]*multiply1[1]*multiply1[2]*multiply1[3]*multiply1[4]*multiply1[5]*y/1062
288 #        print(prob)
289 #        print(prob1)
290         output='unacc\n'
291         if prob1>prob:
292             output='acc\n'
293         if testsample[i][6]==output:
294             if output=='acc\n':
295                 classifier[0]+=1     ######## True Positive #########
296             if output=='unacc\n':
297                 classifier[1]+=1     ######## True Negative #########
298         if testsample[i][6]!=output:
299             if output=='acc\n':
300                 classifier[2]+=1     ######## False Positive #########
301             if output=='unacc\n':
302                 classifier[3]+=1     ######## False Negative  #########
```

```python
303    print"————————True_Positive=_%d——————————False_Negative=_%d————"%(classifier[0],classifier[3])
304    print"————————False_Positive=_%d——————————True_Negative=_%d————"%(classifier[2],classifier[1])
```