

# Exercise Sheet 1: Python Basics

This first exercise sheet tests the basic functionalities of the Python programming language in the context of a simple prediction task. We consider the problem of predicting health risk of subjects from personal data and habits. We first use for this task a decision tree



adapted from the webpage <http://www.refactorthis.net/post/2013/04/10/Machine-Learning-tutorial-How-to-create-a-decision-tree-in-RapidMiner-using-the-Titanic-passenger-data-set.aspx>. For this exercise sheet, you are required to use only pure Python, and to not import any module, including numpy. In exercise sheet 2, the nearest neighbor part of this exercise sheet will be revisited with numpy.

## Classifying a single instance (15 P)

---

- Create a function that takes as input a tuple containing values for attributes (smoker,age,diet), and computes the output of the decision tree.
- Test your function on the tuple `('yes',31,'good')`,

```
def riskCalc(person):
    if ( person[0] == 'yes'):
        if ( person[1] < 30 ):
            return 'Less Risk'
        elif ( person[1] > 30 ):
            return 'More Risk'
    elif ( person[0] == 'no'):
        if ( person[2] == 'good'):
            return 'Less Risk'
        elif ( person[2] == 'poor' ):
            return 'More Risk'

print riskCalc(('yes',31,'poor'))
```

More Risk

## Reading a dataset from a text file (10 P)

---

The file `health-test.txt` contains several fictitious records of personal data and habits.

- Read the file automatically using the methods introduced during the lecture.
- Represent the dataset as a list of tuples.

```
def getTuple(line):
    L = str.split(line[:-1],',')
    L[1]= int(L[1])
    return tuple(L)

def readTuplesFromFile(file):
    tupleList=[]
    for line in open (file, 'r'):
        tupleList += [getTuple(line)]
    return tupleList
print(readTuplesFromFile('health-test.txt'))
```

```
[('yes', 21, 'poor'), ('no', 50, 'good'), ('no', 23, 'good'), ('yes', 45, 'poor'), ('yes', 51, 'good'), ('no', 60, 'good'), ('no', 15, 'poor'), ('no', 18, 'good')]
```

## Applying the decision tree to the dataset (15 P)

- Apply the decision tree to all points in the dataset, and compute the percentage of them that are classified as "more risk".

```
risk=[]
moreRisk=[]
def decisionTree():
    global risk,moreRisk
    data = readTuplesFromFile('health-test.txt')

    risk = [ riskCalc(person) for person in data ]

    moreRisk = filter (lambda val: val == 'More Risk', risk)

decisionTree()
print ( '%.2f %%' % (float(len(moreRisk)) / float(len(risk))*100))
```

37.50 %

## Learning from examples (10 P)

Suppose that instead of relying on a fixed decision tree, we would like to use a data-driven approach where data points are classified based on a set of training observations manually labeled by experts. Such labeled dataset is available in the file `health-train.txt`. The first three columns have the same meaning than for `health-test.txt`, and the last column corresponds to the labels.

- Write a procedure that reads this file and converts it into a list of pairs. The first element of each pair is a triplet of attributes, and the second element is the label.

```
def dataPairs(file):
    tuples = readTuplesFromFile(file)
    pairList = []
    for tup in tuples:
        pairList += [ (tup[:-1],tup[len(tup)-1])]
    return pairList

print dataPairs('health-train.txt')
```

```
[('yes', 54, 'good'), 'less'), (('no', 55, 'good'), 'less'), (('no', 26, 'good'), 'less'), (('yes', 40, 'good'), 'more'), (('yes', 25, 'poor'), 'less'), (('no', 13, 'poor'), 'more'), (('no', 15, 'good'), 'less'), (('no', 50, 'poor'), 'more'), (('yes', 33, 'good'), 'more'), (('no', 35, 'good'), 'less'), (('no', 41, 'good'), 'less'), (('yes', 30, 'poor'), 'more'), (('no', 39, 'poor'), 'more'), (('no', 20, 'good'), 'less'), (('yes', 18, 'poor'), 'less'), (('yes', 55, 'good'), 'more')]
```

## Nearest neighbor classifier (25 P)

We consider the nearest neighbor algorithm that classifies test points following the label of the nearest neighbor in the training data. For this, we need to define a distance function between data points. We define it to be

$$d(a,b) = (a[0] \neq b[0]) + ((a[1] - b[1]) / 50.0)^2 + (a[2] \neq b[2])$$

where `a` and `b` are two tuples corresponding to the attributes of two data points.

- Write a function that retrieves for a test point the nearest neighbor in the training set, and classifies the test point accordingly.
- Test your function on the tuple `('yes', 31, 'good')`

```

closestPoint= tuple()
def distance(a,b):
    dist = (a[0]!=b[0])+(((a[1]-b[1])/50.0)**2)+(a[2]!=b[2])
    return dist

def NearestNeighbourClassify(a):
    distVect=[]
    trainPairs = dataPairs('health-train.txt')
    for pair in trainPairs:
        distVect+= [distance(a,pair[0])]
    return trainPairs[distVect.index(min(distVect))][1]

point = ('yes',31,'good')
tuple((point,NearestNeighbourClassify(point)))

```

```
(( 'yes', 31, 'good'), 'more')
```

- Apply both the decision tree and nearest neighbor classifiers on the test set, and find the data point(s) for which the two classifiers disagree, and with which probability it happens.

```

decisionTree()
riskDecisionTree = [ 'more' if (riskVal =='More Risk') else 'less' for riskVal in risk]

riskNN = []
testPoints = readTuplesFromFile('health-test.txt')
for testPoint in testPoints:
    riskNN += [NearestNeighbourClassify(testPoint)]

mismatchCount=0
mismatchList=[]
for i in range(len(riskNN)):
    if ( riskDecisionTree[i] != riskNN[i]):
        mismatchList += [testPoints[i]]
        mismatchCount+=1

print mismatchList

print mismatchCount/float(len(riskNN))

```

```

[('yes', 51, 'good')]
0.125

```

One problem of simple nearest neighbors is that one needs to compare the point to predict to all data points in the training set. This can be slow for datasets of thousands of points or more. Alternatively, some classifiers train a model first, and then use it to classify the data.

## Nearest mean classifier (25 P)

---

We consider one such trainable model, which operates in two steps:

(1) Compute the average point for each class, (2) classify new points to be of the class whose average point is nearest to the point to predict.

For this classifier, we convert the attributes smoker and diet to real values (for smoker: yes=1.0 and no=0.0, and for diet: good=0.0 and poor=1.0), and use the modified distance function:

$$d(a,b) = (a[0]-b[0])**2 + ((a[1]-b[1])/50.0)**2 + (a[2]-b[2])**2$$

We adopt an object-oriented approach for building this classifier.

- Implement the methods `train` and `predict` of the class `NearestMeanClassifier`.

```

class NearestMeanClassifier:

    val = {'yes':1.0,'no':0.0,'good':0.0,'poor':1.0}
    moreCentroid=(0,0,0)
    lessCentroid=(0,0,0)
    moreCount=0
    lessCount=0
    # Training method that takes as input a dataset
    # and produces two internal vectors corresponding
    # to the mean of each class.

    def calcDistance(self,a,b):
        return (self.val[a[0]]-b[0])**2+((a[1]-b[1])/50.0)**2+(self.val[a[2]]-b[2])**
2

    def train(self,dataset):
#        global moreCentroid,lessCentroid,moreCount,lessCount
        for datapoint in dataPairs(dataset):
            if ( datapoint[1] == 'more'):
                self.moreCentroid=(self.moreCentroid[0]+self.val[datapoint[0][0]],sel
f.moreCentroid[1]+datapoint[0][1],self.moreCentroid[2]+self.val[datapoint[0][2]])
                self.moreCount+=1
            elif ( datapoint[1] == 'less' ):
                self.lessCentroid=(self.lessCentroid[0]+self.val[datapoint[0][0]],sel
f.lessCentroid[1]+datapoint[0][1],self.lessCentroid[2]+self.val[datapoint[0][2]])
                self.lessCount+=1

            self.moreCentroid=(self.moreCentroid[0]/float(self.moreCount),self.moreCentro
id[1]/float(self.moreCount),self.moreCentroid[2]/float(self.moreCount))
            self.lessCentroid=(self.lessCentroid[0]/float(self.lessCount),self.lessCentro
id[1]/float(self.lessCount),self.lessCentroid[2]/float(self.lessCount))
        # Prediction method that takes as input a new data
        # point and predicts it to belong to the class with
        # nearest mean.
    def predict(self,x):
        l=[]
        l += [self.calcDistance(x,self.moreCentroid)]
        l += [self.calcDistance(tuple(x),self.lessCentroid)]
        index = l.index(min(l))
        if ( index == 0):
            return 'more'
        elif ( index ==1 ):
            return 'less'

```

- Build an object of class `NearestMeanClassifier` , train it on the training data, and print the mean

vector for each class.

```
c = NearestMeanClassifier()
c.train('health-train.txt')
print c.predict(('no',31,'poor'))
print tuple((c.moreCentroid,'more'))
print tuple((c.lessCentroid,'less'))
```

```
more
((0.5714285714285714, 37.142857142857146, 0.5714285714285714), 'more')
((0.3333333333333333, 32.111111111111114, 0.2222222222222222), 'less')
```

- Predict the test data using the nearest mean classifier and print all test examples for which all three classifiers (decision tree, nearest neighbor and nearest mean) agree.

```
riskNM = []
testPoints= readTuplesFromFile('health-test.txt')
for testPoint in testPoints:
    riskNM += [c.predict(testPoint)]

for i in range(len(riskNN)):
    if ( riskDecisionTree[i] == riskNN[i] and riskNN[i] == riskNM[i]):
        print testPoints[i]
```

```
('no', 50, 'good')
('no', 23, 'good')
('yes', 45, 'poor')
('no', 60, 'good')
('no', 15, 'poor')
('no', 18, 'good')
```