

***k* nearest neighbors**

1) Given a dataset as follows:

X1	X2	Class
0.376000	0.488000	0
0.312000	0.544000	0
0.298000	0.624000	0
0.394000	0.600000	0
0.506000	0.512000	0
0.488000	0.334000	1
0.478000	0.398000	1
0.606000	0.366000	1
0.428000	0.294000	1
0.542000	0.252000	1

- Classifying the testset with 1NN, 3NN:

X1	X2	Class
0.550000	0.364000	?
0.558000	0.470000	?
0.456000	0.450000	?

0.450000	0.570000	?
----------	----------	---

2) Implement k NN from scratch in Python. The program requires 3 parameters:

- file name of trainset
- file name of testset
- number of nearest neighbors (k)

Dataset with m examples, n dimensions (attribute), c classes ($0, 1, \dots, c-1$), is in the format:

```
val_i1_a1 val_i1_a2 ... val_i1_an class_i1
val_i2_a1 val_i2_a2 ... val_i2_an class_i2
...
val_im_a1 val_im_a2 ... val_im_an class_im
```

The program reports the classification results (accuracy, confusion matrix) with different trials $k=1, 3$, etc for 5 datasets:

- Iris (**.trn**: trainset, **.tst**: testset)
- Optics (**.trn**: trainset, **.tst**: testset)
- Letter (**.trn**: trainset, **.tst**: testset)
- Leukemia (**.trn**: trainset, **.tst**: testset)
- Fp (**.trn**: trainset, **.tst**: testset)

<http://www.cit.ctu.edu.vn/~dtngchi/ml/data.tar.gz>

3) Proof of Cover-Hart's theorem:

For sufficiently large training set size m , the error rate of the 1 NN classifier is less than twice the Bayes error rate.