

# ML lab2 KNearestNeighbor

0416037 李家安

## 1. Result

```
1 resubstation, algo=brute, dist=manhattan
2 1.0
3 --- 0.3311185036791992 seconds ---
4 resubstation, algo=kd_tree, dist=manhattan
5 1.0
6 --- 0.045838832298649414 seconds ---
7 resubstation, algo=ball_tree, dist=manhattan
8 1.0
9 --- 0.12276124954223633 seconds ---
10
11 resubstation, algo=brute, dist=euclidean
12 1.0
13 --- 0.3285883778751953 seconds ---
14 resubstation, algo=kd_tree, dist=euclidean
15 1.0
16 --- 0.0344395637512287 seconds ---
17 resubstation, algo=ball_tree, dist=euclidean
18 1.0
19 --- 0.12823184858398438 seconds ---
20
21 kFold, foldNum=12, algo=brute, dist=manhattan
22 0.623313577672
23 --- 0.35594725688825684 seconds ---
24 kFold, foldNum=12, algo=kd_tree, dist=manhattan
25 0.621873837833
26 --- 0.10834942626953125 seconds ---
27 kFold, foldNum=12, algo=ball_tree, dist=manhattan
28 0.621681586283
29 --- 0.16851843781171875 seconds ---
30
31
32 kFold, foldNum=12, algo=brute, dist=euclidean
33 0.612786148828
34 --- 0.48883493461688887 seconds ---
35 kFold, foldNum=12, algo=kd_tree, dist=euclidean
36 0.612282687753
37 --- 0.10797286833638371 seconds ---
38 kFold, foldNum=12, algo=ball_tree, dist=euclidean
39 0.618215858863
40 --- 0.17551788221435547 seconds ---
41
42 resubstation, algo=brute, dist=mahalanobis
43 1.0
44 --- 1.7785589738529785 seconds ---
45 kFold, foldNum=12, algo=brute, dist=mahalanobis
46 0.677833988625
47 --- 2.841188955387807 seconds ---
48
49 resubstation, algo=brute, dist=cosDist
50 0.567782758477
51 --- 473.8676188468833 seconds ---
52 kFold, foldNum=12, algo=brute, dist=cosDist
53 0.458921873238
54 --- 438.4166786287428 seconds ---
```

## 2. Environment

Ubuntu 16.04

Python 3.5

## 3. Language and Library

Python 3.5

Pandas

NumPy

SciPy

SciKit-learn

## 4. How to use it

```
$ python3 wine.py
```

## 5. Code

在程式中我先對資料做處理，將 dataset 以網路的方式抓下來，用 pandas 轉換為 DataFrame 的形式，再把 quality 這個 column 切出來分成 feature 和 target，同時對 feature 做 normalization。

定義兩個函式 resubstitution 和 kfold 分別做 resubstitution 跟 kfold，分別用不同的距離跟算法當參數輸入。

resubstation 的部分，以距離跟算法建立 KNeighborsClassifier 物件 nbr，nbr.fit() 來訓練，然後用 nbr.predict 來預測，最後用 confusion matrix 測量精準度。過程中我們沒有對 data 做任何的切割，也就是直接拿原始資料做訓練以及預測，符合 resubstation 的採樣方式。

kfold 的部分，以 kFold object 將 dataset 分割為 12 份，用 for loop 分別計算每次預測的精準度，加總後再除以 12 算出平均精準度。

運算時將 ball\_tree, kd\_tree, brute 與 manhattan, euclidean 交叉丟入兩函式運算，以及將 brute 與 mahalanobis, cosDist 交叉丟入函式運算，就可以得到各類模型的精確度了。