

# Pattern Recognition - KNN classifier implement

1. Store [UCI dataset - Wine Data Set](#) as pandas dataframe, and store every feature and label into variable X and y respectively.

```
1 import pandas as pd
2
3 def load_feature(filename):
4     df = pd.read_csv(filename)
5     feature = df.drop('Class',1)
6     return feature
7
8 def load_label(filename):
9     df = pd.read_csv(filename)
10    label = df.Class
11    return label
12
13 def main():
14     X = load_feature('wine.data')
15     y = load_label('wine.data')
```

```
In [38]: runfile('/home/snoopyknight/文件/knn/knn.py', wdir='/home/snoopyknight/文件/knn')
0   Alcohol  Malic acid  Ash  Alcalinity of ash  Magnesium  Total phenols \
1   14.23      1.71    2.43      15.6           127         2.80
2   13.20      1.78    2.14      11.2           100         2.65
3   13.16      2.36    2.67      18.6           101         2.80
4   14.37      1.95    2.50      16.8           113         3.85
5   13.24      2.59    2.87      21.0           118         2.80

   Flavanoids  Nonflavanoid phenols  Proanthocyanins  Color intensity  Hue \
0      3.06          0.28           2.29           5.64          1.04
1      2.76          0.26           1.28           4.38          1.05
2      3.24          0.30           2.81           5.68          1.03
3      3.49          0.24           2.18           7.80          0.86
4      2.69          0.39           1.82           4.32          1.04

   OD280/OD315 of diluted wines  Proline
0              3.92             1065
1              3.40             1050
2              3.17             1185
3              3.45             1480
4              2.93              735
```

*print(X.head())*

```
In [39]: runfile('/home/snoopyknight/文件/knn/knn.py', wdir='/home/snoopyknight/文件/knn')
0      1
1      1
2      1
3      1
4      1
Name: Class, dtype: int64
```

*print(y.head())*

2. Split training data and testing data.(based on sklearn.model\_selection.train\_test\_split)

The store them in X\_train, X\_test, y\_train, y\_test.

```
1 from sklearn.model_selection import train_test_split
2
3 def main():
4     X = load_feature('wine.data')
5     y = load_label('wine.data')
6
7     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=
8     0.5)
```

3. Calculate cosine similarity of all testing data and training data, then store them in the `cs_array(row:X_test_arr,col: X_train)`.

```
1 from sklearn.model_selection import train_test_split
2 import numpy as np
3
4 def cos_sim(X_train, X_test, y_train, y_test):
5     X_train_arr = np.array(X_train.values)
6     X_test_arr = np.array(X_test.values)
7     cs_array = cosine_similarity(X_test_arr,X_train_arr)
8     return cs_array
```

```
In [42]: runfile('/home/snoopyknight/文件/knn/knn.py', wdir='/home/snoopyknight')
[[0.99974582 0.98767391 0.99962477 ... 0.99976859 0.99965528 0.9996475 ]
 [0.99885173 0.99555036 0.99892307 ... 0.9961601  0.99565948 0.99565911]
 [0.99979148 0.99307442 0.99983312 ... 0.9980913  0.99778233 0.99774827]
 ...
 [0.99878215 0.98303509 0.99856005 ... 0.99996635 0.99999747 0.9999962 ]
 [0.99991292 0.9919165  0.99991913 ... 0.99862683 0.99836576 0.99833441]
 [0.99944678 0.99438196 0.99955029 ... 0.99721853 0.99686059 0.99680792]]

print(cs_array)
```

#### 4. knn\_classify

- choose k largest values of cosine similarity in each row. (It means that k training data nodes which is close to the current testing data node.)
- Compare labels of this k nodes and choose the most common label as the result of the prediction of the testing data

```
1 def knn_classify(X_train, X_test, y_train, y_test, k):
2     cs_array = cos_sim(X_train, X_test, y_train, y_test)
3     k_list = []
4     y_pred_list = []
5     for i in range(len(cs_array)):
6         k_list = heapq.nlargest(k, range(len(cs_array[i])), cs_array[i].take)
7         #k_list stores index of k largest cosine similarity value
8         class_list = []
9         for idx in k_list:
10             class_list.append(y_train.iloc[idx])
11             #class_list stores k prediction classes of each node.
12         a = np.array(class_list)
13         counts = np.bincount(a)
14         print(np.argmax(counts))
15         #choose the most common class as the result prediction
16         y_pred_list.append(np.argmax(counts))
```

```

17         print("=====")
18     y_pred = pd.Series(y_pred_list)
19     return y_pred

```

```

In [47]: runfile('/home/snoopyknight/文件/knn/knn.py', wdir='/home/snoopyknight/文件/knn')
[3, 3, 2, 3, 2]
3
=====
[2, 3, 2, 2, 2]
2
=====
[2, 1, 2, 2, 2]
2
=====
[2, 2, 2, 2, 2]
2
=====
[2, 2, 3, 1, 3]
2
=====
[2, 2, 2, 2, 2]
2
=====
[1, 1, 1, 1, 1]
1
=====
[3, 3, 2, 2, 3]
3

```

#### 5. Check the classifier report and accuracy

```

1 from sklearn.metrics import classification_report as clf_report
2 from sklearn.metrics import accuracy_score
3
4 y_pred = knn_classify(X_train, X_test, y_train, y_test, 5)
5 accuracy = accuracy_score(y_test, y_pred)
6 report = clf_report(y_test, y_pred)

```

```

          precision    recall  f1-score   support

     1         0.92      0.88      0.90         26
     2         0.78      0.85      0.82         34
     3         0.81      0.76      0.79         29

 avg / total         0.83      0.83      0.83         89

=====
accuracy = 0.8314606741573034

```

The code is on: [https://github.com/SnoopyKnight/knn\\_classifier](https://github.com/SnoopyKnight/knn_classifier)

#### Reference:

- <http://enginebai.logdown.com/posts/241676/knn>
- <https://pandas.pydata.org/pandas-docs/stable/>
- <http://scikit-learn.org/stable/>
- <http://www.numpy.org/>

