

# Q2 (k-means)

姓名：王依睿

学号：1552651

## a) k-means

1. 引入 `sklearn` 和 `KMeans`

```
import sklearn
from sklearn.cluster import KMeans
```

2. k-means聚类并计算其silhouette系数

```
# k-means聚类并计算其silhouette系数
def kmeans(n_clusters):
    X = StandardScaler().fit_transform(data.T)
    # Incorrect number of clusters
    labels = KMeans(n_clusters = n_clusters).fit_predict(X)
    silhouette_score = metrics.silhouette_score(X, labels,
metric='euclidean')
    return silhouette_score
```

3. 选择不同的k进行k-means聚类并计算其silhouette系数

```
# 选择不同的k进行k-means聚类并计算其silhouette系数
def vip_kmeans():
    n_clusters_list = pd.Series(range(2, int(math.sqrt(vipno_total / 2)
+ 40)))
    vip_kmeans_df = pd.DataFrame(columns=['silhouette_score'], index =
n_clusters_list)
    for n_clusters in n_clusters_list:
        silhouette_score = kmeans(n_clusters)
        vip_kmeans_df['silhouette_score'][n_clusters] =
silhouette_score
    return vip_kmeans_df
```

## 结果

标准化后

**silhouette\_score**

<b>2</b>	0.48722
<b>3</b>	0.258304
<b>4</b>	0.409811
<b>5</b>	-0.0225924
<b>6</b>	-0.164264
<b>7</b>	0.341519
<b>8</b>	0.109789
<b>9</b>	0.184277
<b>10</b>	0.201069
<b>11</b>	-0.154637
<b>12</b>	0.279303
<b>13</b>	0.171036
<b>14</b>	-0.281389
<b>15</b>	-0.149187
<b>16</b>	0.256569
<b>17</b>	0.187516
<b>18</b>	-0.0967663

**19** 0.0399323

**20** 0.190137

**21** 0.238392

**22** 0.253588

**23** 0.276265

**24** -0.14351

**25** -0.019352

**26** 0.0447384

**27** 0.23375

**28** -0.17327

**29** 0.121717

**30** -0.2243

**31** 0.13538

**32** 0.0375936

**33** -0.271416

**34** -0.150653

**35** 0.18909

**36** 0.162469

36	0.102409
37	0.124659
38	0.189629
39	0.107688
40	0.167317
41	0.0933236
42	-0.0477688
43	-0.152667
44	-0.19394
45	0.195354
46	-0.0799995
47	-0.170349
48	-0.158859
49	-0.0380165
50	0.0929955

标准化前

silhouette_score	
2	0.941759

<b>3</b>	0.864978
<b>4</b>	0.806563
<b>5</b>	0.73279
<b>6</b>	0.302269
<b>7</b>	0.582366
<b>8</b>	0.295611
<b>9</b>	0.503593
<b>10</b>	0.264648
<b>11</b>	0.320026
<b>12</b>	0.108558
<b>13</b>	0.314193
<b>14</b>	0.221744
<b>15</b>	0.325281
<b>16</b>	0.284054
<b>17</b>	0.24598
<b>18</b>	0.312843
<b>19</b>	0.0630744
<b>20</b>	0.000000

<b>20</b>	0.263849
<b>21</b>	0.163162
<b>22</b>	0.305543
<b>23</b>	0.21772
<b>24</b>	0.242047
<b>25</b>	0.260464
<b>26</b>	0.252603
<b>27</b>	0.10152
<b>28</b>	0.306447
<b>29</b>	0.12048
<b>30</b>	0.0659082
<b>31</b>	0.128939
<b>32</b>	-0.148261
<b>33</b>	0.0438414
<b>34</b>	0.0896596
<b>35</b>	0.0911414
<b>36</b>	0.0965219
<b>37</b>	0.18304

38	0.0655652
39	0.0567768
40	0.276519
41	0.0947079
42	0.0679171
43	0.134885
44	0.132513
45	0.0357655
46	0.116735
47	0.118276
48	0.106903
49	0.0428625
50	0.0765129

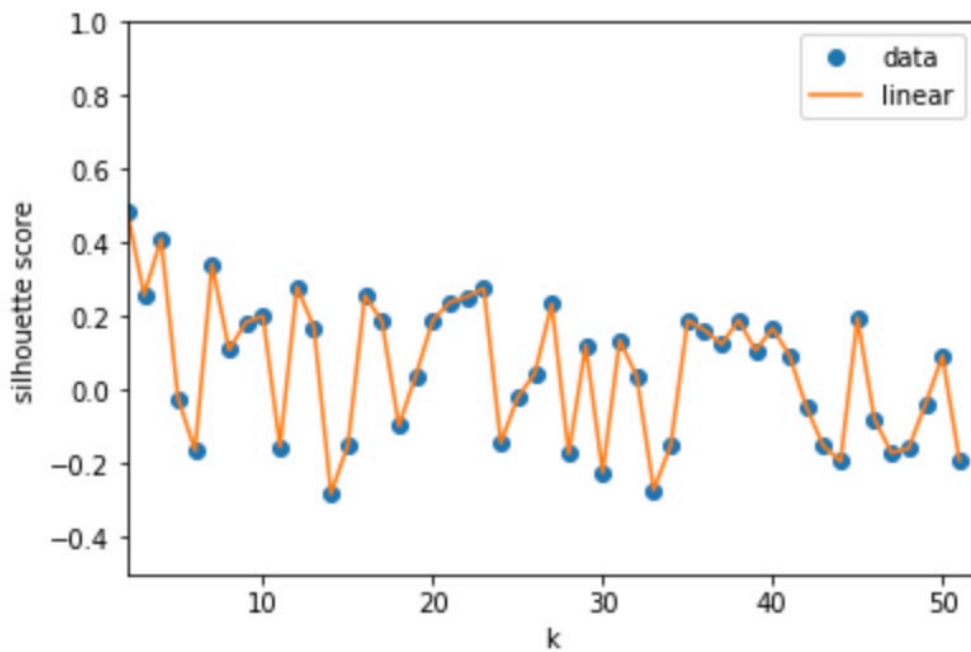
1. 利用 `matplotlib.pyplot` 以k值为横轴、Silhouette系数值为y轴，画出Silhouette系数值-k值的函数图

```
# 作图
def kmeans_plot(vip_kmeans_df):
    f = interp1d(vip_kmeans_df.index, vip_kmeans_df.silhouette_score)
    xnew = vip_kmeans_df.index.copy()
    plt.plot(vip_kmeans_df.index, vip_kmeans_df.silhouette_score, 'o',
             xnew, f(xnew), '-')
    plt.xlim([2, int(math.sqrt(vipno_total / 2) + 40)])
```

```
plt.ylim([-0.5, 1])
plt.legend(['data', 'linear'], loc='best')
plt.xlabel('k')
plt.ylabel('silhouette score')
plt.suptitle(
    "Silhouette analysis for KMeans clustering on reco data",
    fontsize=14, fontweight='bold')
plt.show()
```

标准化后

### Silhouette analysis for KMeans clustering on reco data



标准化前





[illegible]

## 2. 查看knn结果

```
hash_size: 3
k: 5
distance      vipno
0      11056  29000000549289
1      11176  1593140967467
2      11466  1595142205462
3      11766  1590141216914
4      11790  1592140611301
```

## 结论

1. 若将Silhouette系数作为评价k-means聚类质量的标准，则k-means算法倾向于将数据点聚类成较少的类别。
2. 根据聚类结果来看，k-means聚类倾向于将极大多数点分为一簇，其现实意义可能是由于大多数人的购物习惯都较为相似，数据量较少也是其中一个原因。
3. 正因为极大多数点划分为同一簇，从而使得验证Ish的knn结果时，Ish的knn查询结果与输入vipno都在同一个簇。