

What is Hierarchical Clustering?

Let's say we have the below points and we want to cluster them into groups:

initial points

We can assign each of these points to a separate cluster:

multiple clusters

Now, based on the similarity of these clusters, we can combine the most similar clusters together and repeat this process until only a single cluster is left:

single cluster

We are essentially building a hierarchy of clusters. That's why this algorithm is called hierarchical clustering. I will discuss how to decide the number of clusters in a later section. For now, let's look at the different types of hierarchical clustering.

Types of Hierarchical Clustering

There are mainly two types of hierarchical clustering:

Agglomerative hierarchical clustering

Divisive Hierarchical clustering

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
```

Bad key "text.kerning_factor" on line 4 in
C:\Users\91920\anaconda3\lib\site-packages\matplotlib\mpl-data\stylelib_classic_tes
t_patch.mplstyle.
You probably need to get an updated matplotlibrc file from
<https://github.com/matplotlib/matplotlib/blob/v3.1.3/matplotlibrc.template>
or from the matplotlib source distribution

```
In [2]: data = pd.read_csv('C:/Users/91920/Downloads/Wholesale customers data.csv')
data.head()
```

```
Out[2]:
```

	Channel	Region	Fresh	Milk	Grocery	Frozen	Detergents_Paper	Delicassen
0	2	3	12669	9656	7561	214	2674	1338
1	2	3	7057	9810	9568	1762	3293	1776
2	2	3	6353	8808	7684	2405	3516	7844
3	1	3	13265	1196	4221	6404	507	1788
4	2	3	22615	5410	7198	3915	1777	5185

```
In [3]: data.info()

<class 'pandas.core.frame.DataFrame'>
```

```

RangeIndex: 440 entries, 0 to 439
Data columns (total 8 columns):
#   Column              Non-Null Count  Dtype
---  -
0   Channel              440 non-null    int64
1   Region               440 non-null    int64
2   Fresh                440 non-null    int64
3   Milk                 440 non-null    int64
4   Grocery              440 non-null    int64
5   Frozen               440 non-null    int64
6   Detergents_Paper     440 non-null    int64
7   Delicassen           440 non-null    int64
dtypes: int64(8)
memory usage: 27.6 KB

```

There are multiple product categories – Fresh, Milk, Grocery, etc. The values represent the number of units purchased by each client for each product. Our aim is to make clusters from this data that can segment similar clients together. We will, of course, use Hierarchical Clustering for this problem.

```

In [4]: from sklearn.preprocessing import normalize
data_scaled = normalize(data)
data_scaled = pd.DataFrame(data_scaled, columns=data.columns)
data_scaled.head()

```

```

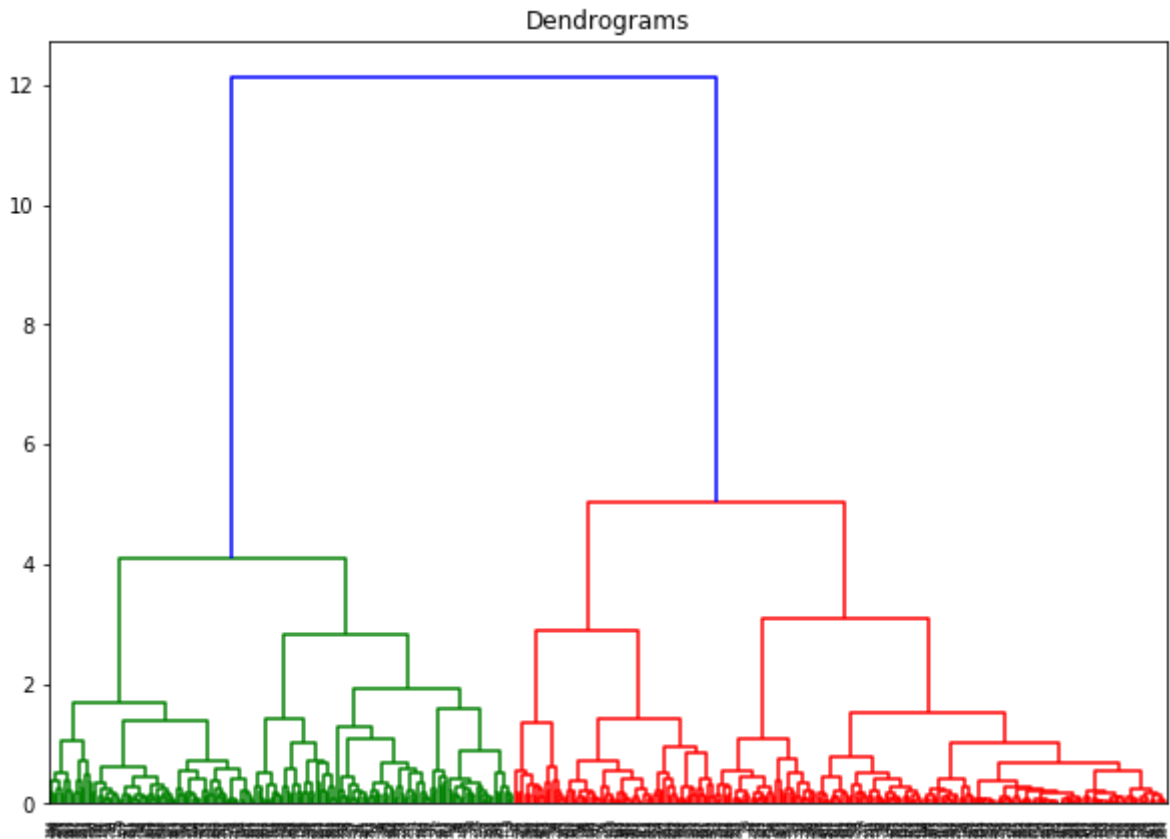
Out[4]:
   Channel  Region  Fresh  Milk  Grocery  Frozen  Detergents_Paper  Delicassen
0  0.000112  0.000168  0.708333  0.539874  0.422741  0.011965         0.149505  0.074809
1  0.000125  0.000188  0.442198  0.614704  0.599540  0.110409         0.206342  0.111286
2  0.000125  0.000187  0.396552  0.549792  0.479632  0.150119         0.219467  0.489619
3  0.000065  0.000194  0.856837  0.077254  0.272650  0.413659         0.032749  0.115494
4  0.000079  0.000119  0.895416  0.214203  0.284997  0.155010         0.070358  0.205294

```

```

In [5]: import scipy.cluster.hierarchy as shc
plt.figure(figsize=(10, 7))
plt.title("Dendrograms")
dend = shc.dendrogram(shc.linkage(data_scaled, method='ward'))

```



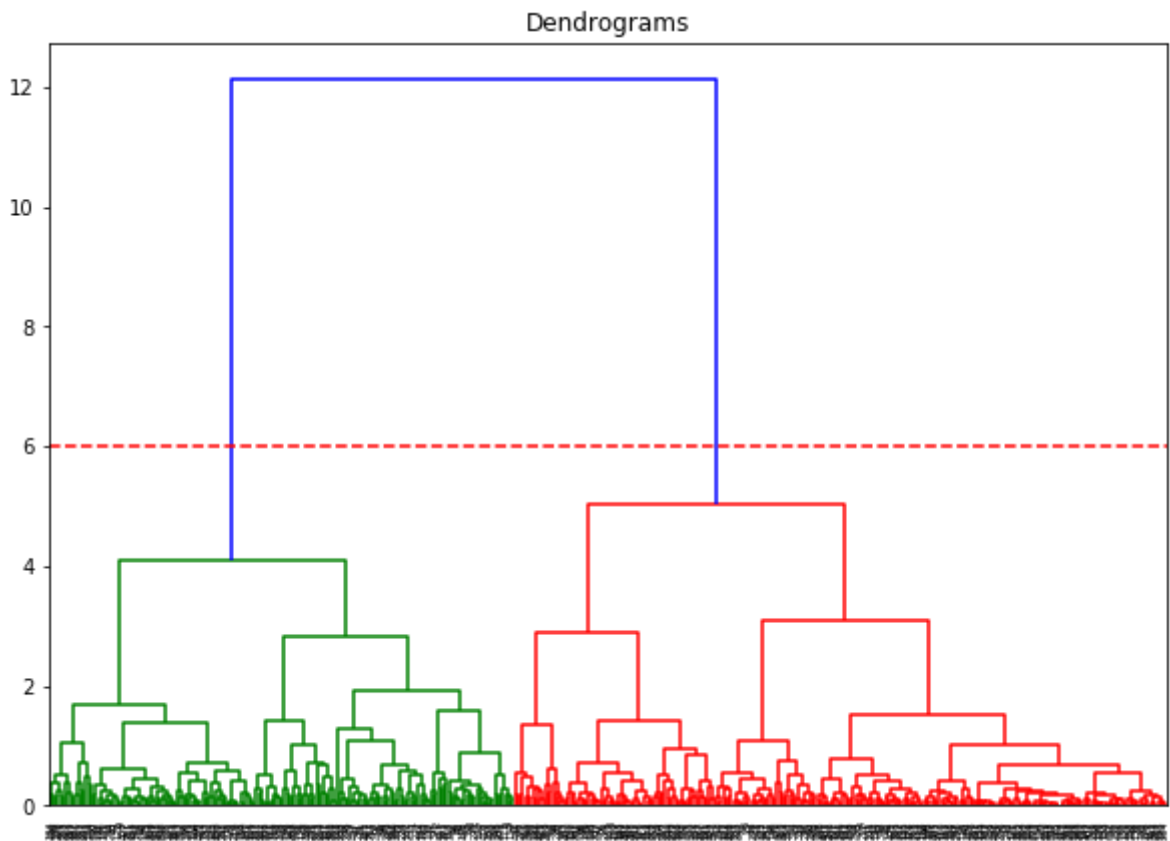
```
In [6]: shc.linkage(data_scaled, method='ward')
```

```
Out[6]: array([[8.30000000e+01, 2.82000000e+02, 1.53464710e-02, 2.00000000e+00],
               [1.50000000e+02, 3.36000000e+02, 1.80052884e-02, 2.00000000e+00],
               [2.58000000e+02, 2.73000000e+02, 1.99150754e-02, 2.00000000e+00],
               ...,
               [8.71000000e+02, 8.73000000e+02, 4.10333306e+00, 1.83000000e+02],
               [8.74000000e+02, 8.75000000e+02, 5.02655073e+00, 2.57000000e+02],
               [8.76000000e+02, 8.77000000e+02, 1.21238165e+01, 4.40000000e+02]])
```

The x-axis contains the samples and y-axis represents the distance between these samples. The vertical line with maximum distance is the blue line and hence we can decide a threshold of 6 and cut the dendrogram:

```
In [6]: plt.figure(figsize=(10, 7))
plt.title("Dendrograms")
dend = shc.dendrogram(shc.linkage(data_scaled, method='ward'))
plt.axhline(y=6, color='r', linestyle='--')
```

```
Out[6]: <matplotlib.lines.Line2D at 0x235f75f6608>
```

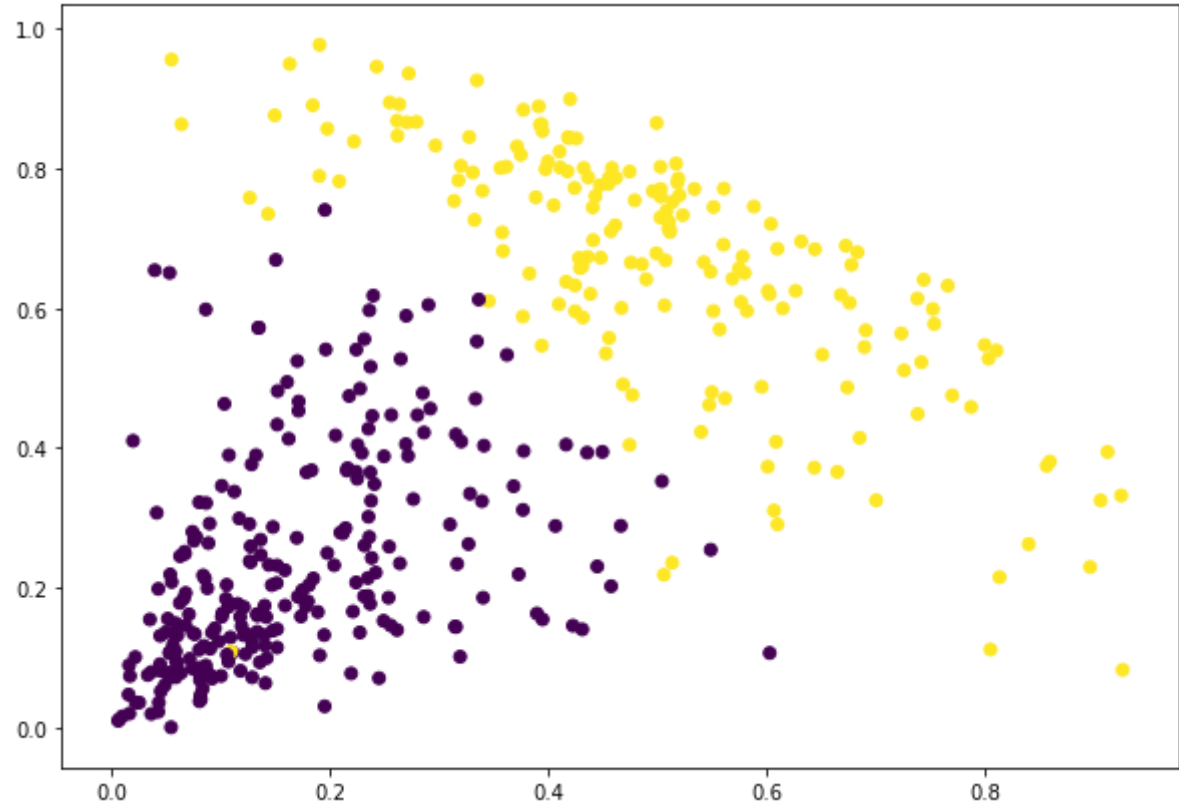


```
In [7]: from sklearn.cluster import AgglomerativeClustering
cluster = AgglomerativeClustering(n_clusters=2, affinity='euclidean', linkage='ward')
cluster.fit_predict(data_scaled)
```

```
Out[7]: array([1, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0,
0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 1, 1,
1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1, 1, 0, 1,
1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 1, 1, 0,
0, 0, 0, 0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1,
0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0,
0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1,
0, 1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1,
0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1,
0, 0, 1, 1, 0, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1, 1, 0, 1, 0,
0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0,
0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1,
1, 1, 1, 0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0,
0, 0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 1, 1, 1, 0,
0, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0,
0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 0, 1,
1, 1, 0, 1, 0, 1, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0,
0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0,
1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 1, 1, 1,
1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 1, 0, 1],
dtype=int64)
```

```
In [8]: plt.figure(figsize=(10, 7))
plt.scatter(data_scaled['Milk'], data_scaled['Grocery'], c=cluster.labels_)
```

```
Out[8]: <matplotlib.collections.PathCollection at 0x235f71e3708>
```



```
In [ ]:
```