

Basics of hierarchical clustering

CLUSTER ANALYSIS IN PYTHON



Shaumik Daityari
Business Analyst

Creating a distance matrix using linkage

clustering

```
scipy.cluster.hierarchy.linkage(observations,  
                                method='single', cluster  
                                metric='euclidean', metric :  
                                optimal_ordering=False : 2D  
)
```

- `method` : how to calculate the proximity of clusters
- `metric` : distance metric
- `optimal_ordering` : order data points

Which method should use?

- single: based on two closest objects
- complete: based on two farthest objects
- average: based on the arithmetic mean of all objects
- centroid: based on the geometric mean of all objects
- median: based on the median of all objects
- ward: based on the sum of squares

method

cluster

cluster

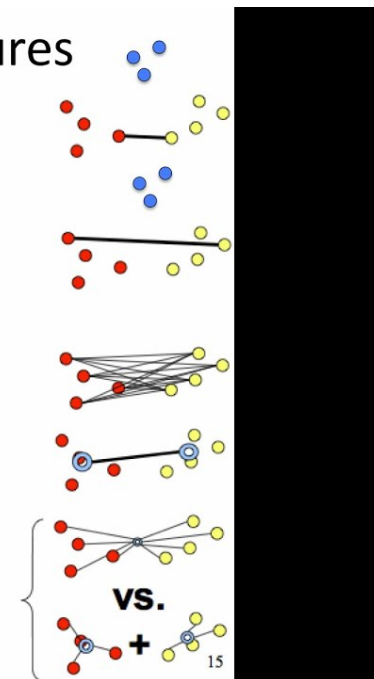
cluster centr

method

cluster proximities()

Cluster distance measures

- Single link: $D(c_1, c_2) = \min_{x_1 \in c_1, x_2 \in c_2} D(x_1, x_2)$
 - distance between closest elements in clusters
 - produces long chains $a \rightarrow b \rightarrow c \rightarrow \dots \rightarrow z$
- Complete link: $D(c_1, c_2) = \max_{x_1 \in c_1, x_2 \in c_2} D(x_1, x_2)$
 - distance between farthest elements in clusters
 - forces "spherical" clusters with consistent "diameter"
- Average link: $D(c_1, c_2) = \frac{1}{|c_1|} \frac{1}{|c_2|} \sum_{x_1 \in c_1} \sum_{x_2 \in c_2} D(x_1, x_2)$
 - average of all pairwise distances
 - less affected by outliers
- Centroids: $D(c_1, c_2) = D\left(\left(\frac{1}{|c_1|} \sum_{x \in c_1} \vec{x}\right), \left(\frac{1}{|c_2|} \sum_{x \in c_2} \vec{x}\right)\right)$
 - distance between centroids (means) of two clusters
- Ward's method: $TD_{c_1 \cup c_2} = \sum_{x \in c_1 \cup c_2} D(x, \mu_{c_1 \cup c_2})^2$
 - consider joining two clusters, how does it change the total distance (TD) from centroids?



4) Ward 연결법 (Ward's method)

군집 간 정보의 손실을 최소화하는 군집화

군집 내 편차들의 제곱합을 고려하여 군집 내 거리 (within cluster distance, ESS)를 최소화.

비슷한 크기의 군집을 생성하는 경향

Create cluster labels with fcluster

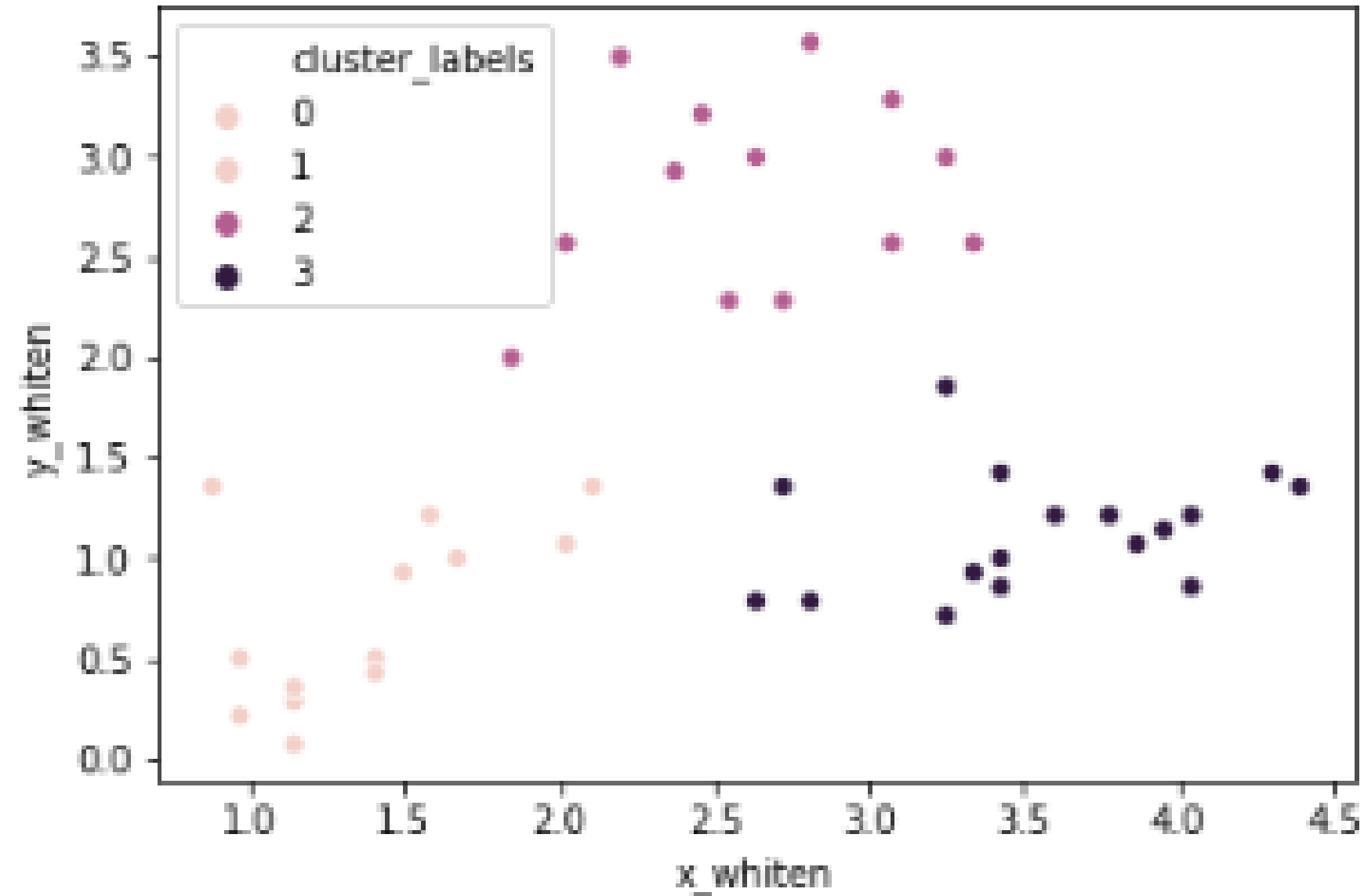
cluster label

```
scipy.cluster.hierarchy.fcluster(distance_matrix, 3, cluster_label_type='maxcluster', num_clusters=3, criterion='maxcluster')
```

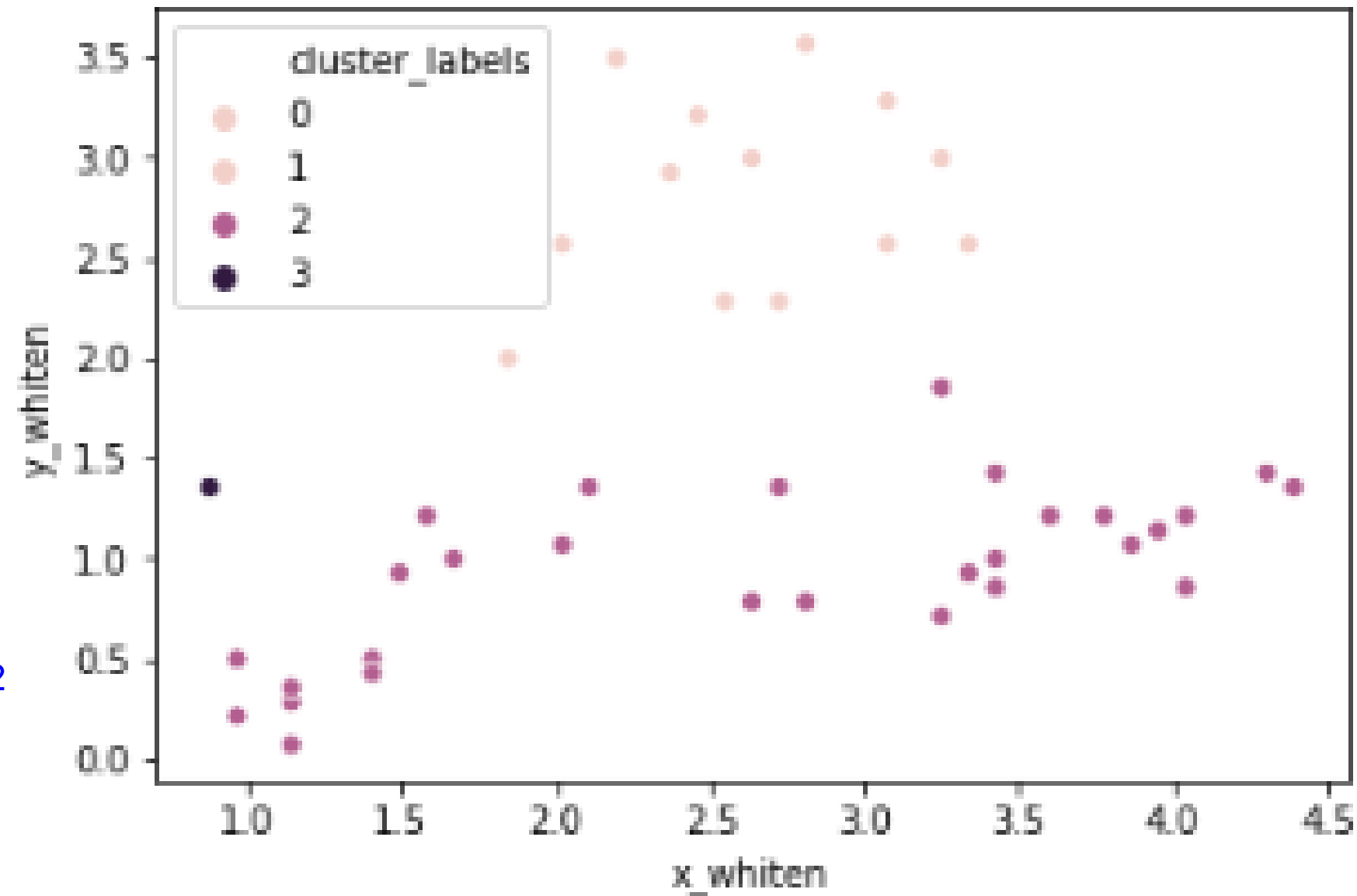
- `distance_matrix` : output of `linkage()` method
- `num_clusters` : number of clusters
- `criterion` : how to decide thresholds to form clusters `criterion` `maxcluster`

Hierarchical clustering with ward method

ward method
cluster
seaborn plot 가
label 0 가 cluster가



Hierarchical clustering with single method



single method

cluster
cluster

가 가

single mothod

clustering
cluster

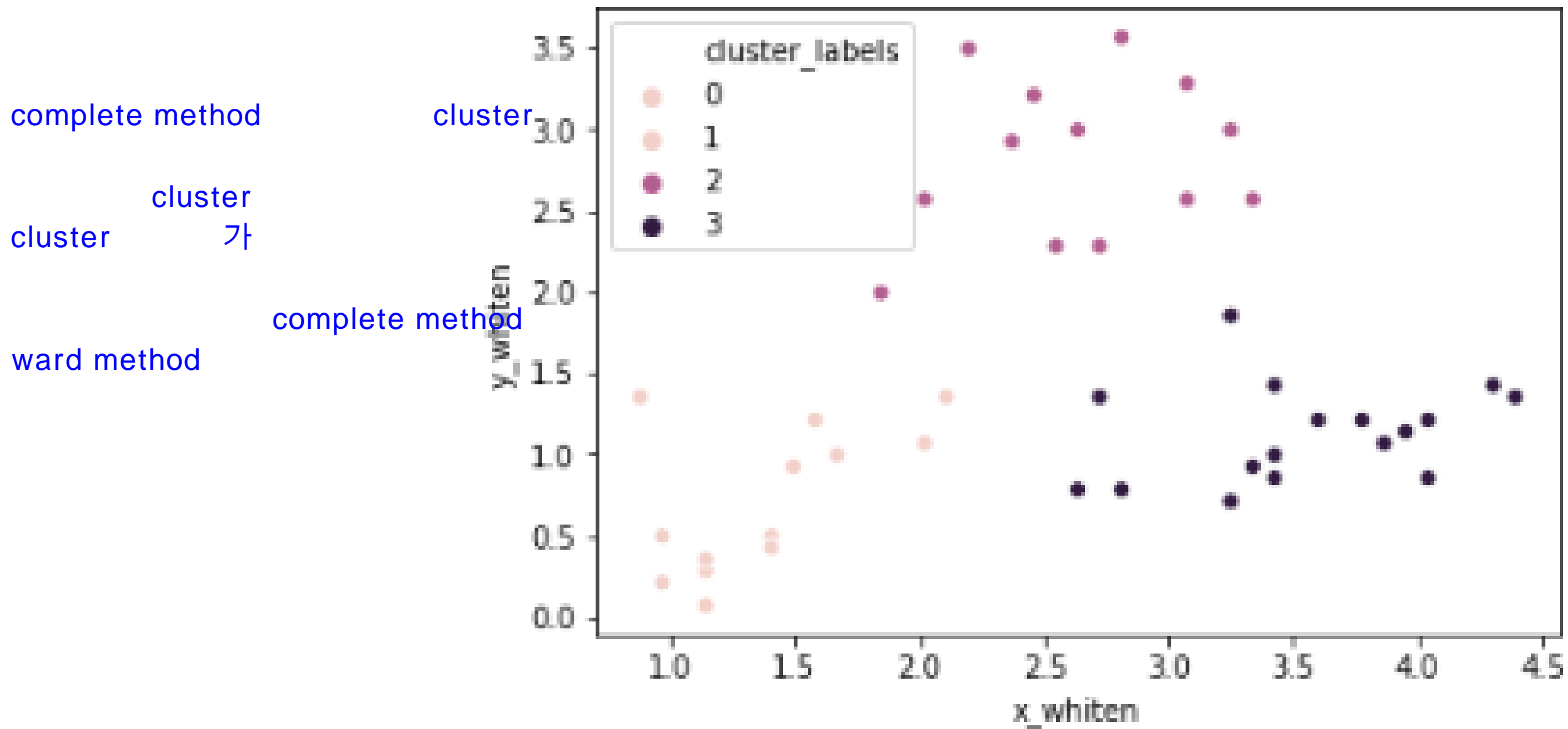
1

3

가

2

Hierarchical clustering with complete method



Final thoughts on selecting a method

- No one right method for all
- Need to carefully understand the distribution of data

가

1.

2. data

Let's try some exercises

CLUSTER ANALYSIS IN PYTHON

Visualize clusters

CLUSTER ANALYSIS IN PYTHON



Shaumik Daityari
Business Analyst

Why visualize clusters?

- Try to make sense of the clusters formed
- An additional step in validation of clusters
- Spot trends in data

cluster ?

- cluster center
- cluster 가 가
-

cluster

An introduction to seaborn

- `seaborn`: a Python data visualization library based on `matplotlib`
- Has better, easily modifiable aesthetics than matplotlib! plotting
- Contains functions that make data visualization tasks easy in the context of data analytics
- Use case for clustering: `hue` parameter for plots seaborn scatter cluster label
cluster

Visualize clusters with matplotlib

```
from matplotlib import pyplot as plt
```

```
df = pd.DataFrame({'x': [2, 3, 5, 6, 2],  
                  'y': [1, 1, 5, 5, 2],  
                  'labels': ['A', 'A', 'B', 'B', 'A']}) cluster  
  
colors = {'A': 'red', 'B': 'blue'} dictionary  
  
df.plot.scatter(x='x',  
               y='y',  
               c=df['labels'].apply(lambda x: colors[x])) lambda c  
  
plt.show()
```

Visualize clusters with seaborn

```
from matplotlib import pyplot as plt
import seaborn as sns
```

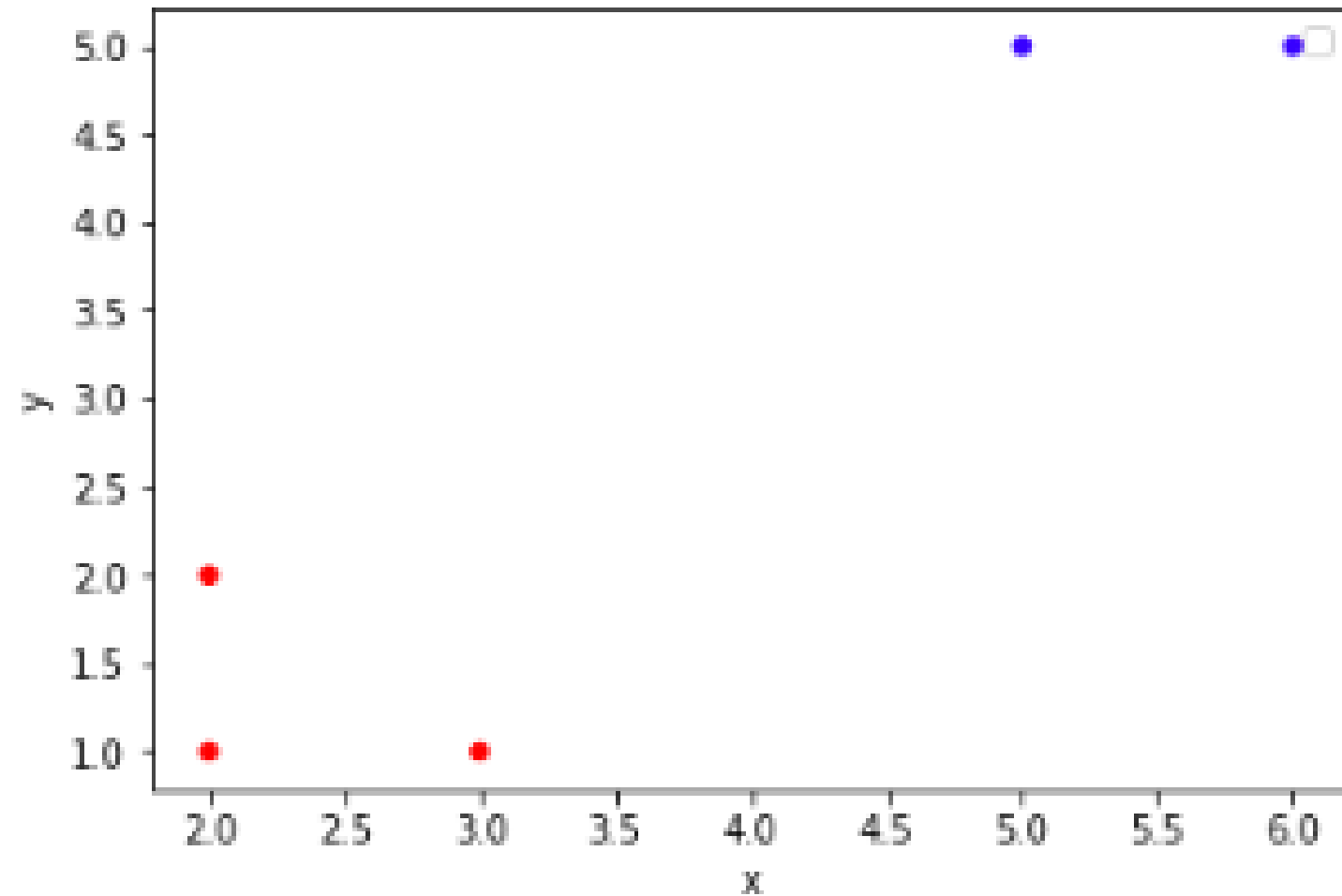
```
df = pd.DataFrame({'x': [2, 3, 5, 6, 2],
                   'y': [1, 1, 5, 5, 2],
                   'labels': ['A', 'A', 'B', 'B', 'A']})

sns.scatterplot(x='x',
                y='y',
                hue='labels', label='cluster label', lw=2)

plt.show()
```

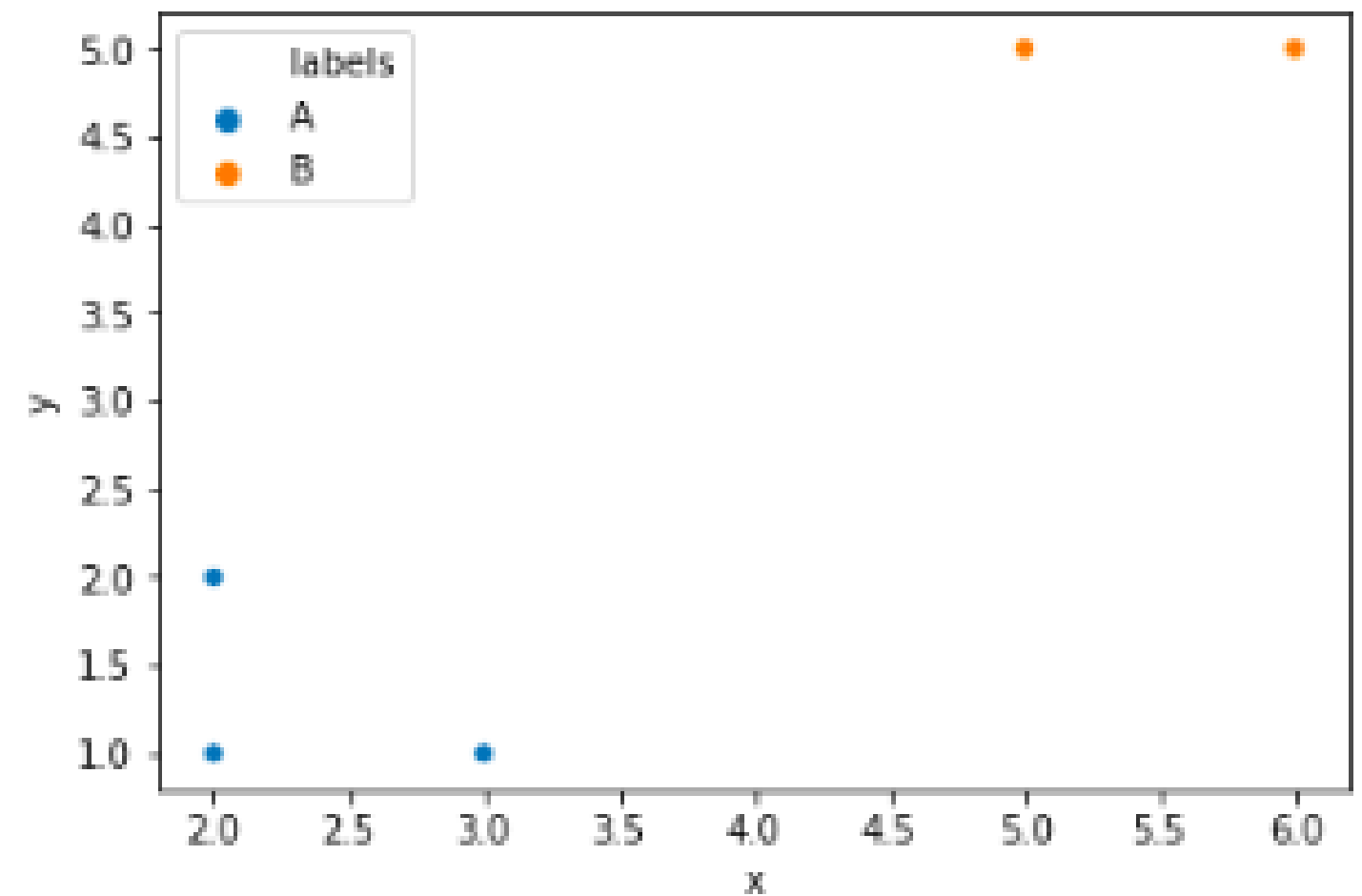
Comparison of both methods of visualization

MATPLOTLIB PLOT



1. dataframe cluster label
2. seaborn

SEABORN PLOT



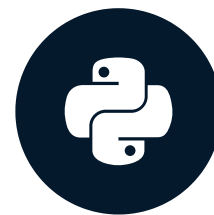
seaborn
가 cluster

Next up: Try some visualizations

CLUSTER ANALYSIS IN PYTHON

How many clusters?

CLUSTER ANALYSIS IN PYTHON



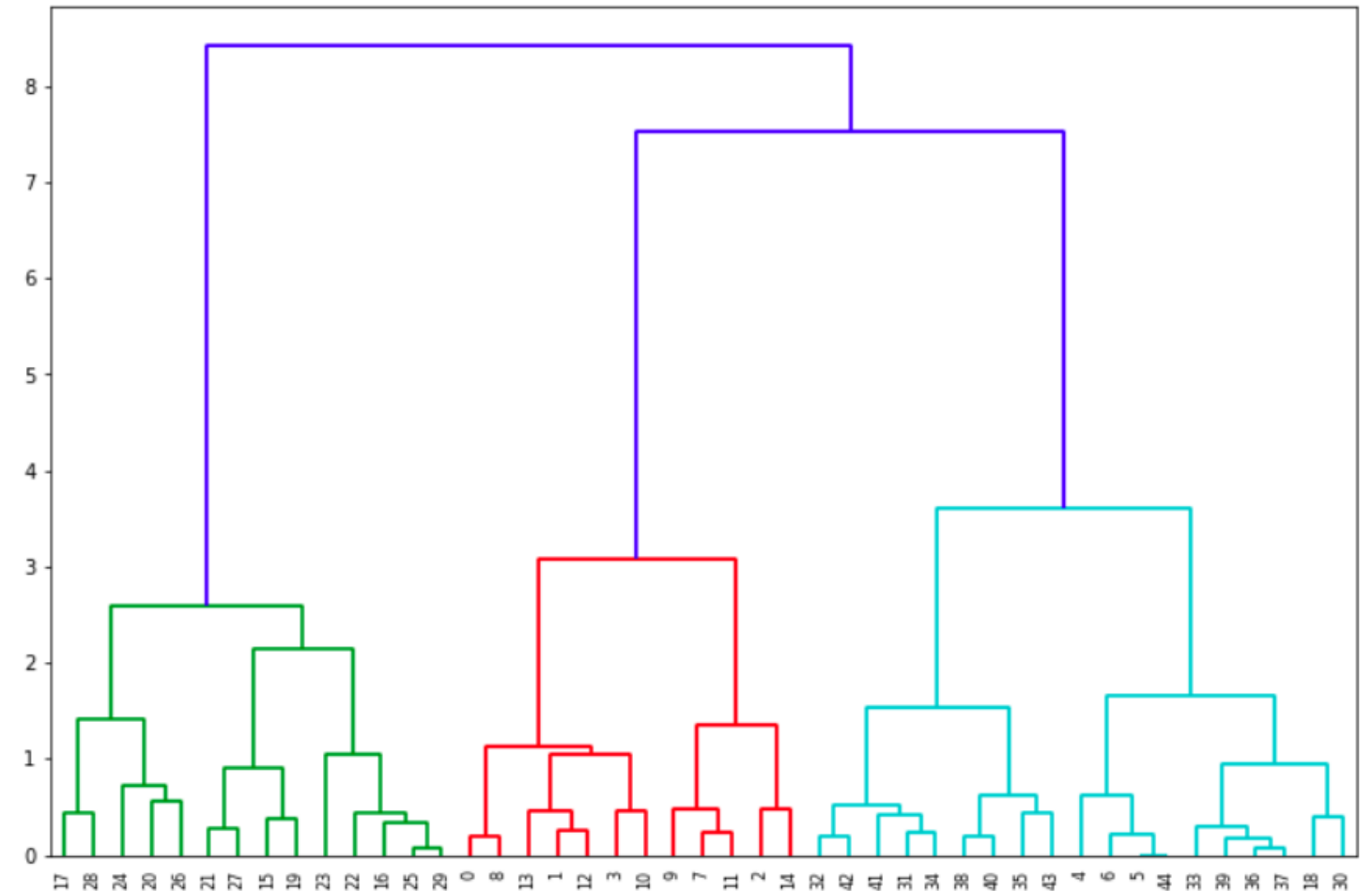
cluster가

Shaumik Daityari
Business Analyst

Introduction to dendrograms

- Strategy till now - decide clusters on visual inspection
- Dendrograms help in showing progressions as clusters are merged
- A dendrogram is a branching diagram that demonstrates how each cluster is composed by branching out into its child nodes

- cluster
- cluster cluster
-



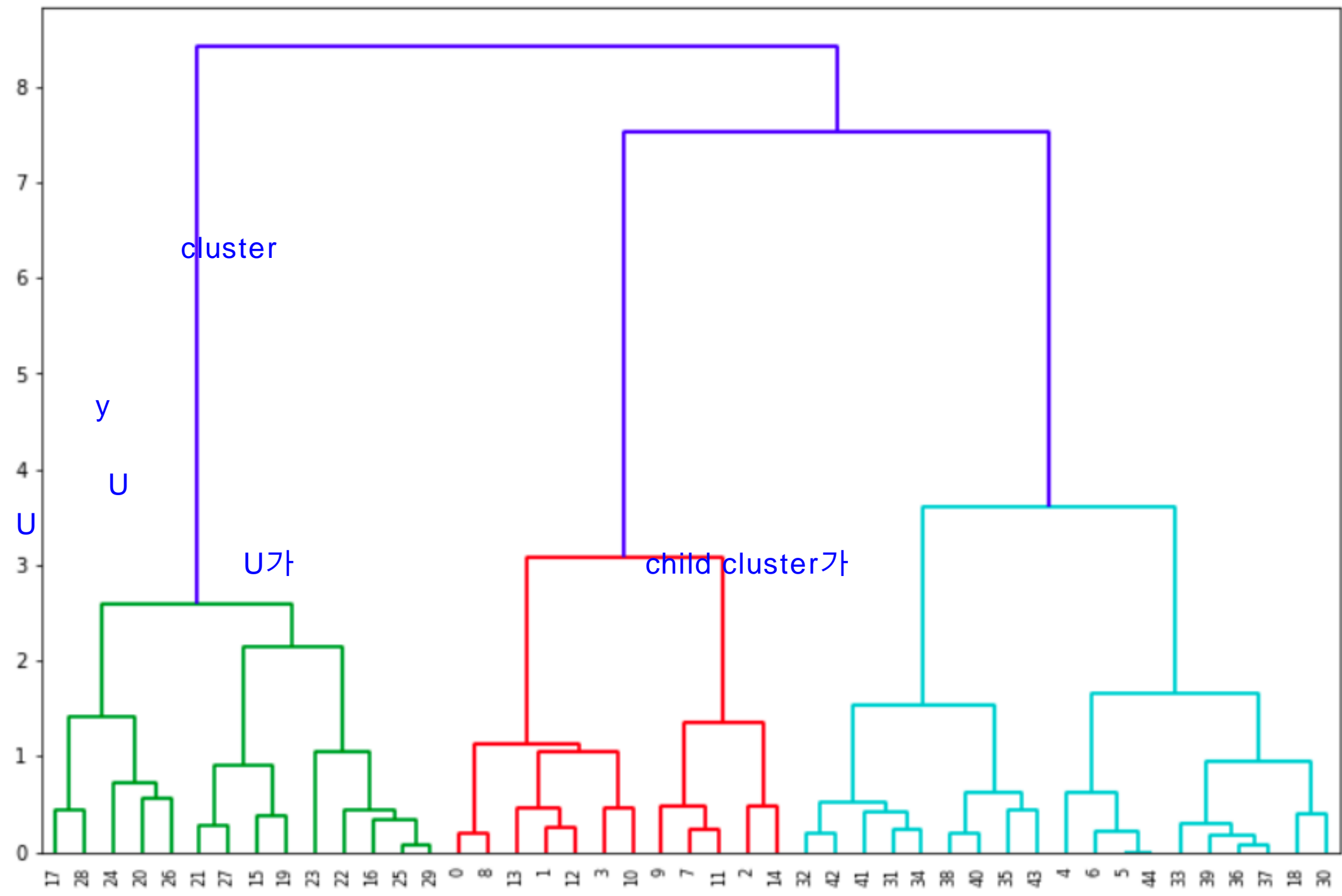
Create a dendrogram in SciPy

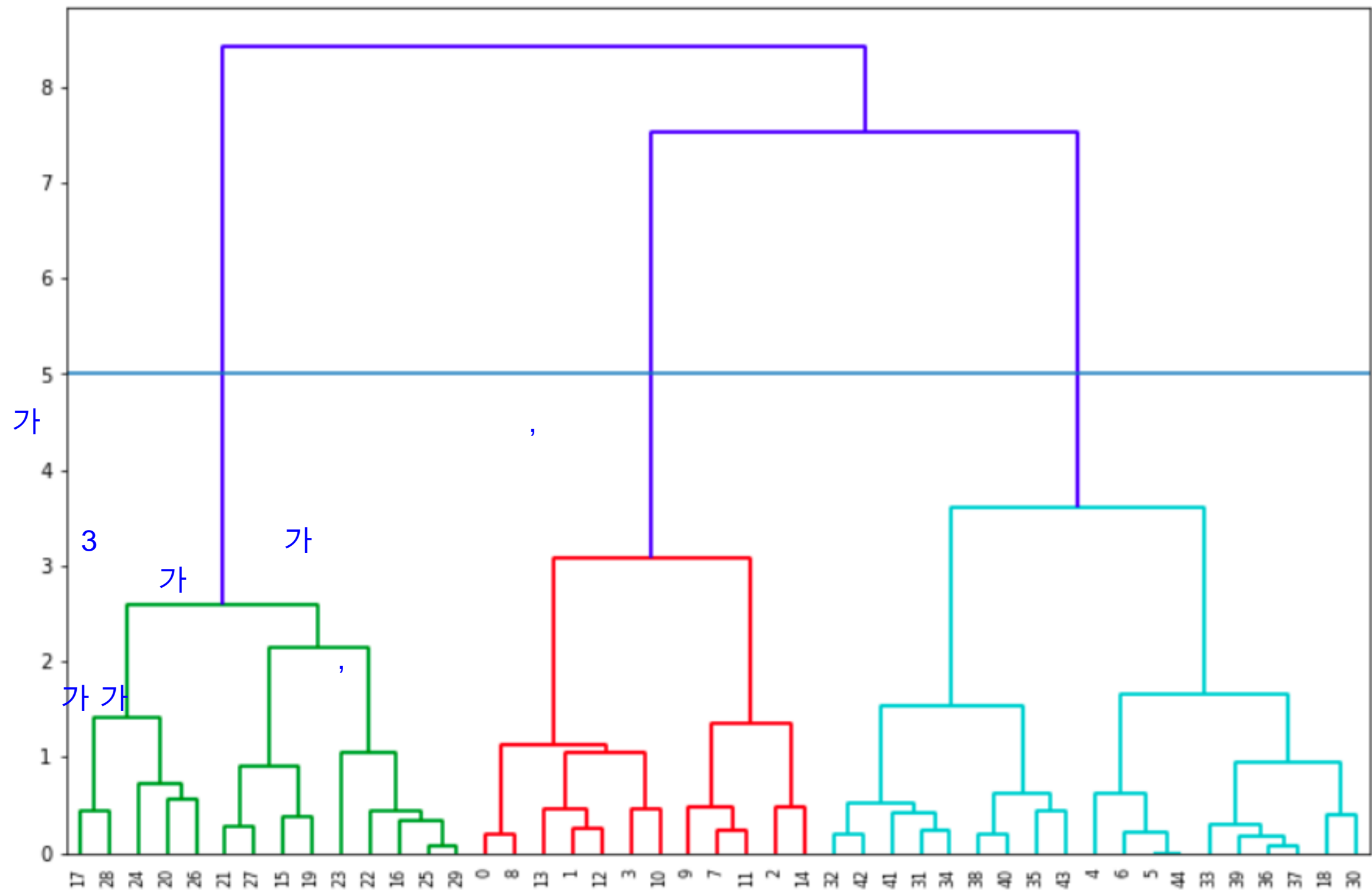
```
from scipy.cluster.hierarchy import dendrogram
```

```
Z = linkage(df[['x_whiten', 'y_whiten']], linkage
            method='ward',
            metric='euclidean')

dn = dendrogram(Z)           dendrogram           linkage           plot
plt.show()
```

x
y
dendrogram
U
가
point
가
U가

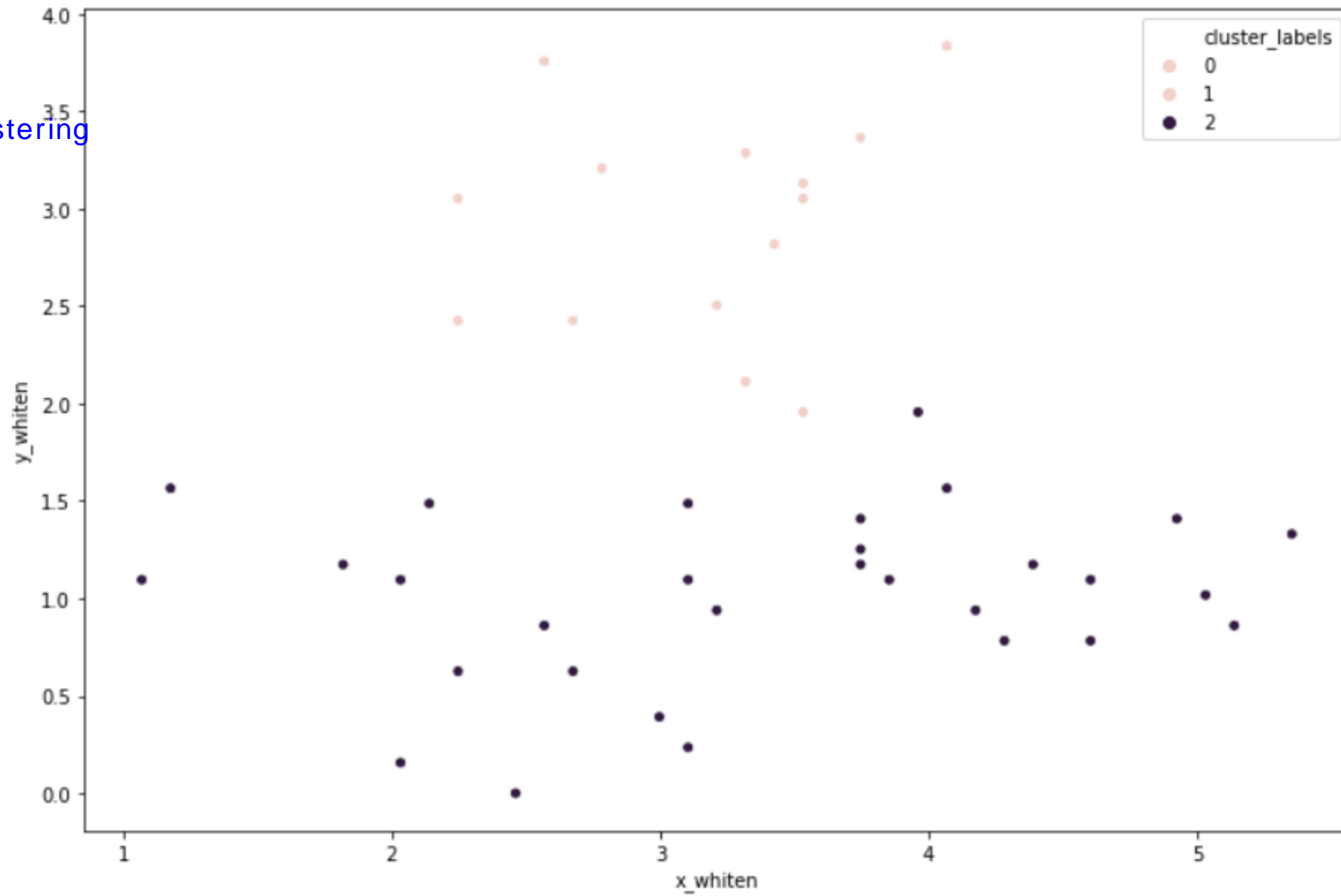




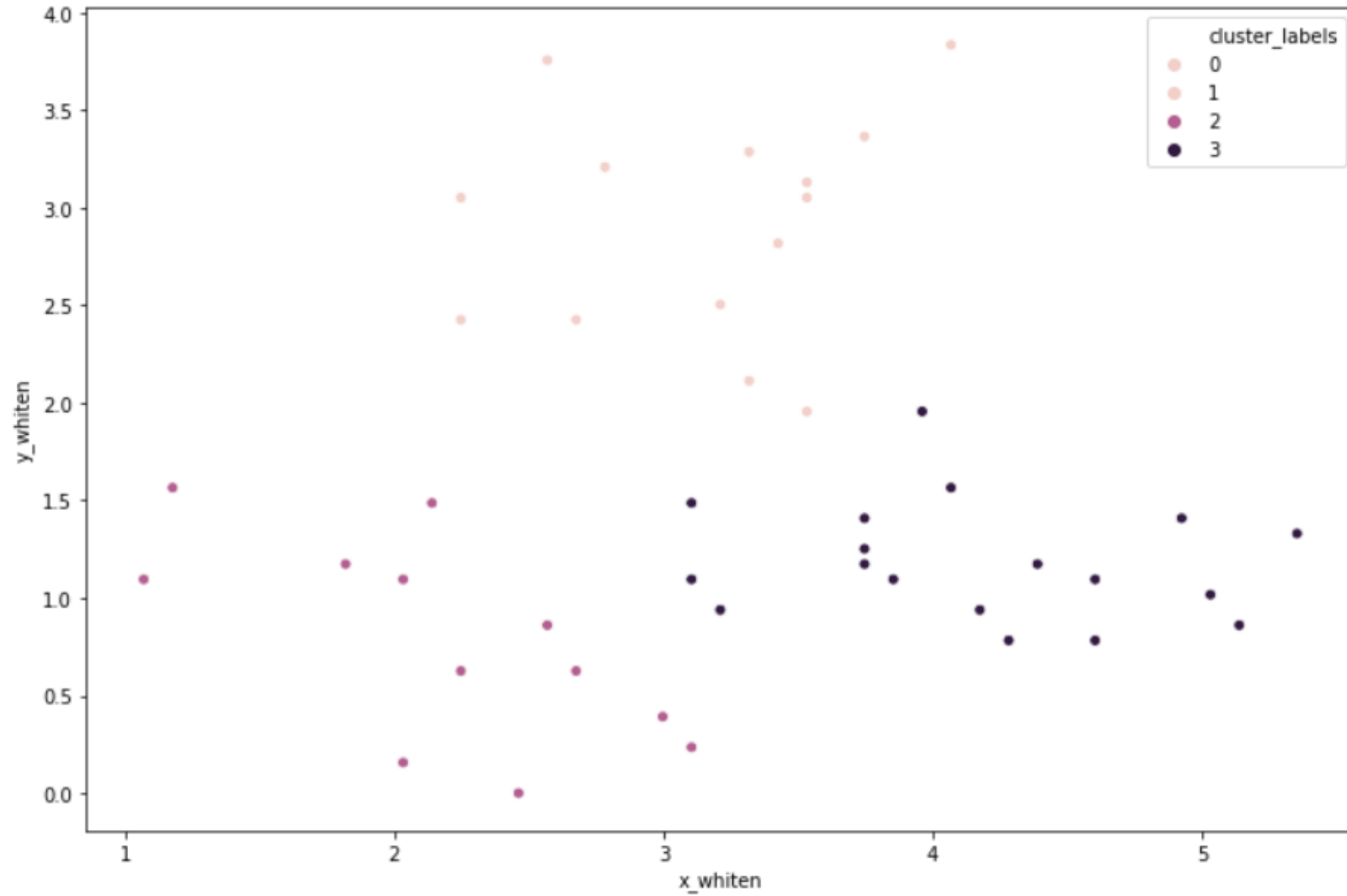
ex) point
ex) 3

2

clustering



3 cluster

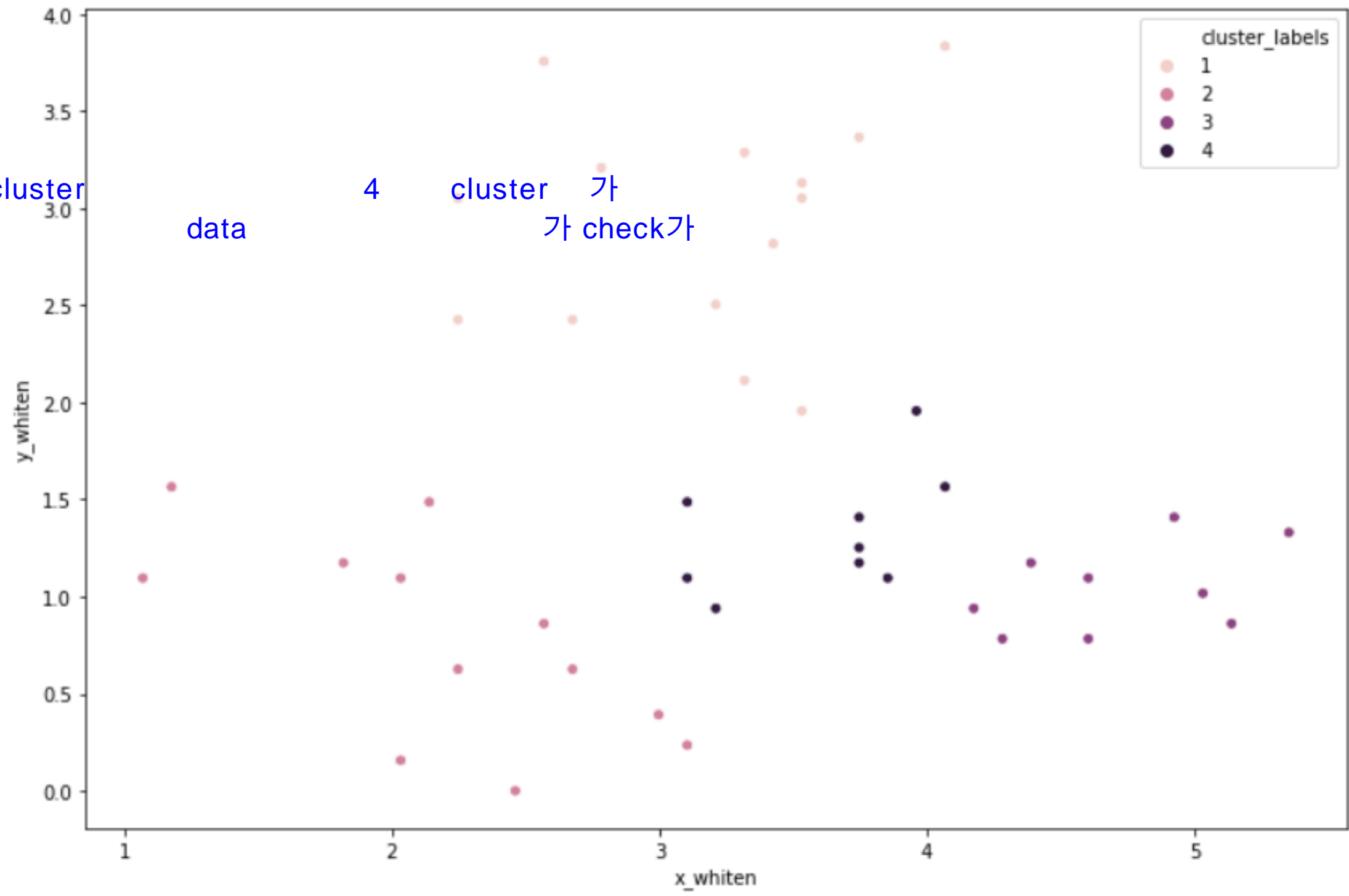


4 cluster
3 cluster
cluster

data

4

cluster 가
가 check가



Next up - try some exercises

CLUSTER ANALYSIS IN PYTHON

Limitations of hierarchical clustering

CLUSTER ANALYSIS IN PYTHON



Shaumik Daityari
Business Analyst

Measuring speed in hierarchical clustering

- `timeit` module
- Measure the speed of `.linkage()` method
- Use randomly generated points
- Run various iterations to extrapolate

- `timeit` module
- clustering process `.linkage()` method
- `XY` `.` `가`

Use of timeit module

100

```
from scipy.cluster.hierarchy import linkage
import pandas as pd
import random, timeit
points = 100
df = pd.DataFrame({'x': random.sample(range(0, points), points), 0~100
                  'y': random.sample(range(0, points), points)})
%timeit linkage(df[['x', 'y']], method = 'ward', metric = 'euclidean')
interpreter
timeit
report
```

1.02 ms ± 133 µs per loop (mean ± std. dev. of 7 runs, 1000 loops each)

jupyter notebook

1.02ms

가

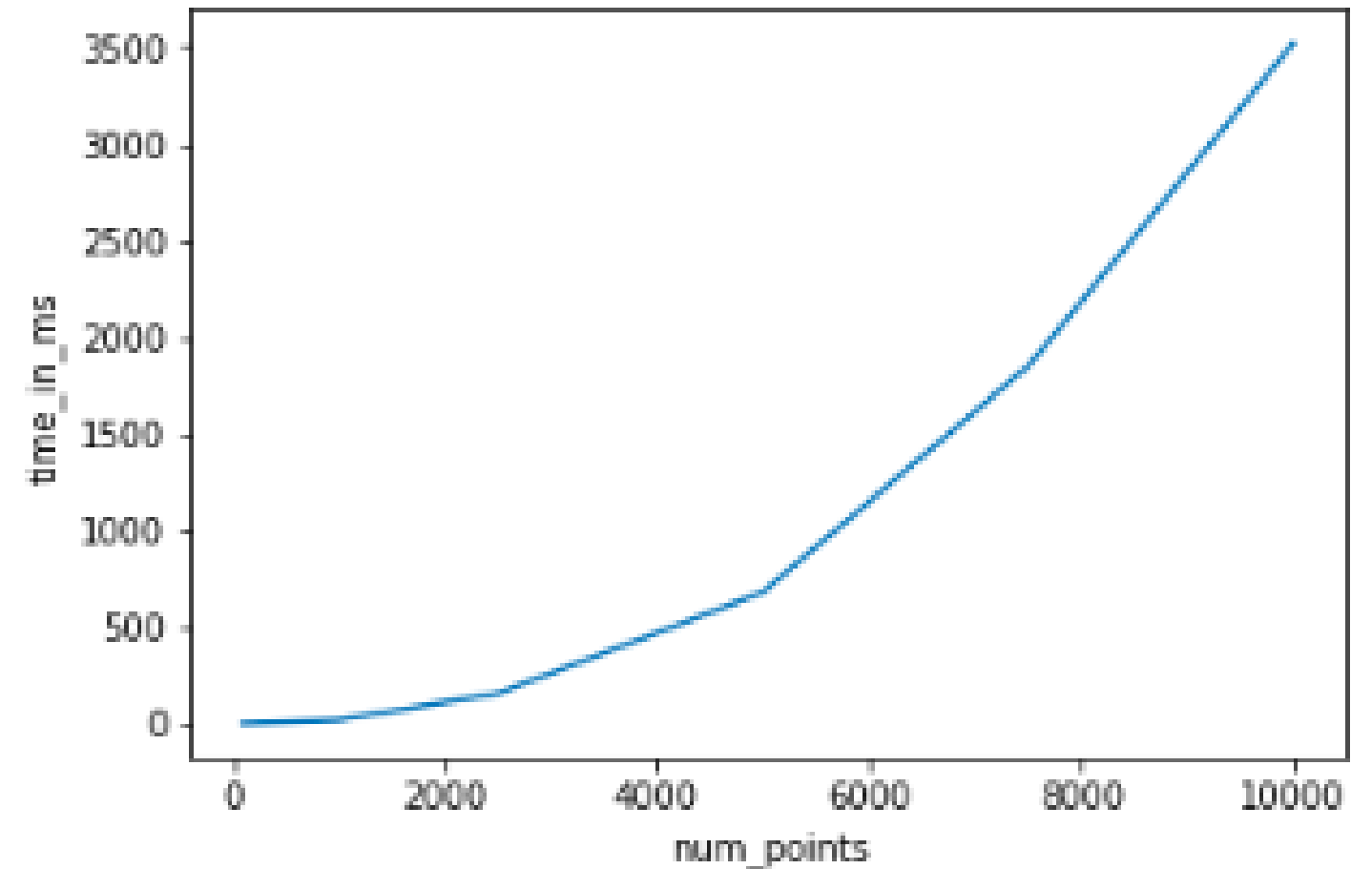
=> point ,

plot

Comparison of runtime of linkage method

- Increasing runtime with data points
- Quadratic increase of runtime
- Not feasible for large datasets

- point data point 가 가 runtime 가
- 가 가 가 2
 , , 1 ,
 data point clustering



Next up - exercises

CLUSTER ANALYSIS IN PYTHON