# Basics of k-means clustering

## CLUSTER ANALYSIS IN PYTHON

**Shaumik Daityari**
Business Analyst

# Why k-means clustering?

- A critical drawback of hierarchical clustering: runtime

- K means runs significantly faster on large datasets

:

k

# Step 1: Generate cluster centers

scipy  k- means clustering
1. generate cluster centers
(1)  kmeans

```
kmeans(obs, k_or_guess, iter, thresh, check_finite)
```

- `obs` : standardized observations   whiten

- `k_or_guess` : number of clusters   cluster

- `iter` : number of iterations (default: 20)   (default : 20)

- `thres` : threshold (default: 1e-05)
  k- means
  . defualt   0.00001

- `check_finite` : whether to check if observations contain only finite numbers (default: True)

NaN
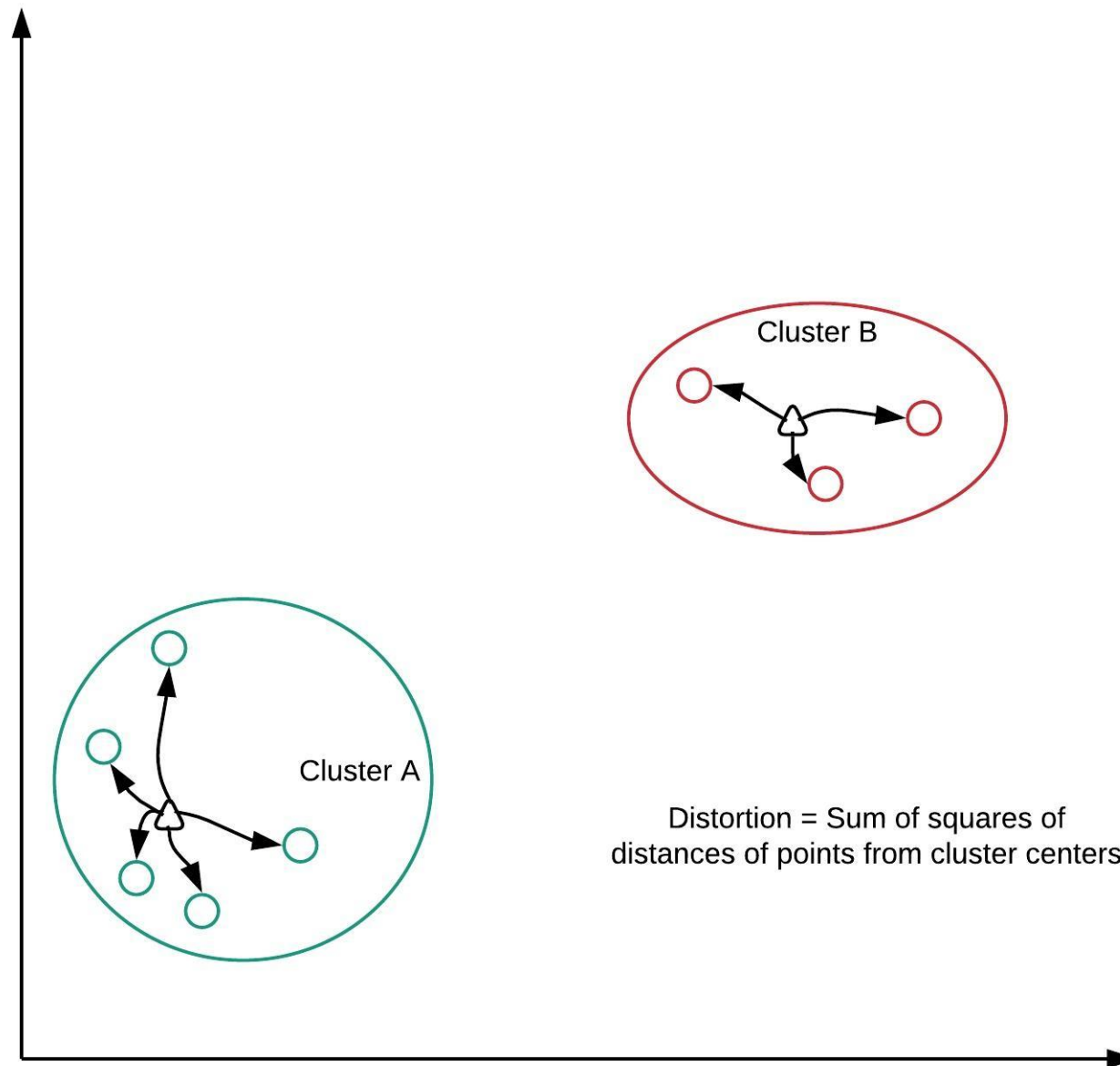
Returns two objects: cluster centers, distortion( default   : True) NaN   point

k- means   cluster center   return
cluster center   code book

k- means cluster   k- menas

# How is distortion calculated?

$$J = \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} \| \mathbf{x}_n - \mu_k \|^2$$



Cluster B

Cluster A

Distortion = Sum of squares of
distances of points from cluster centers

# Step 2: Generate cluster labels

```
vq(obs, code_book, check_finite=True)
```

vq        cluster label

3

- **obs** : standardized observations        whiten method

- **code_book** : cluster centers        kmeans method

- **check_finite** : whether to check if observations contain only finite numbers (default: True)

check_finite        NaN

Returns two objects: a list of cluster labels, a list of distortions

"code book index"        return

# A note on distortions

- `kmeans` returns a single value of distortions

- `vq` returns a list of distortions.

kmeans
    vq

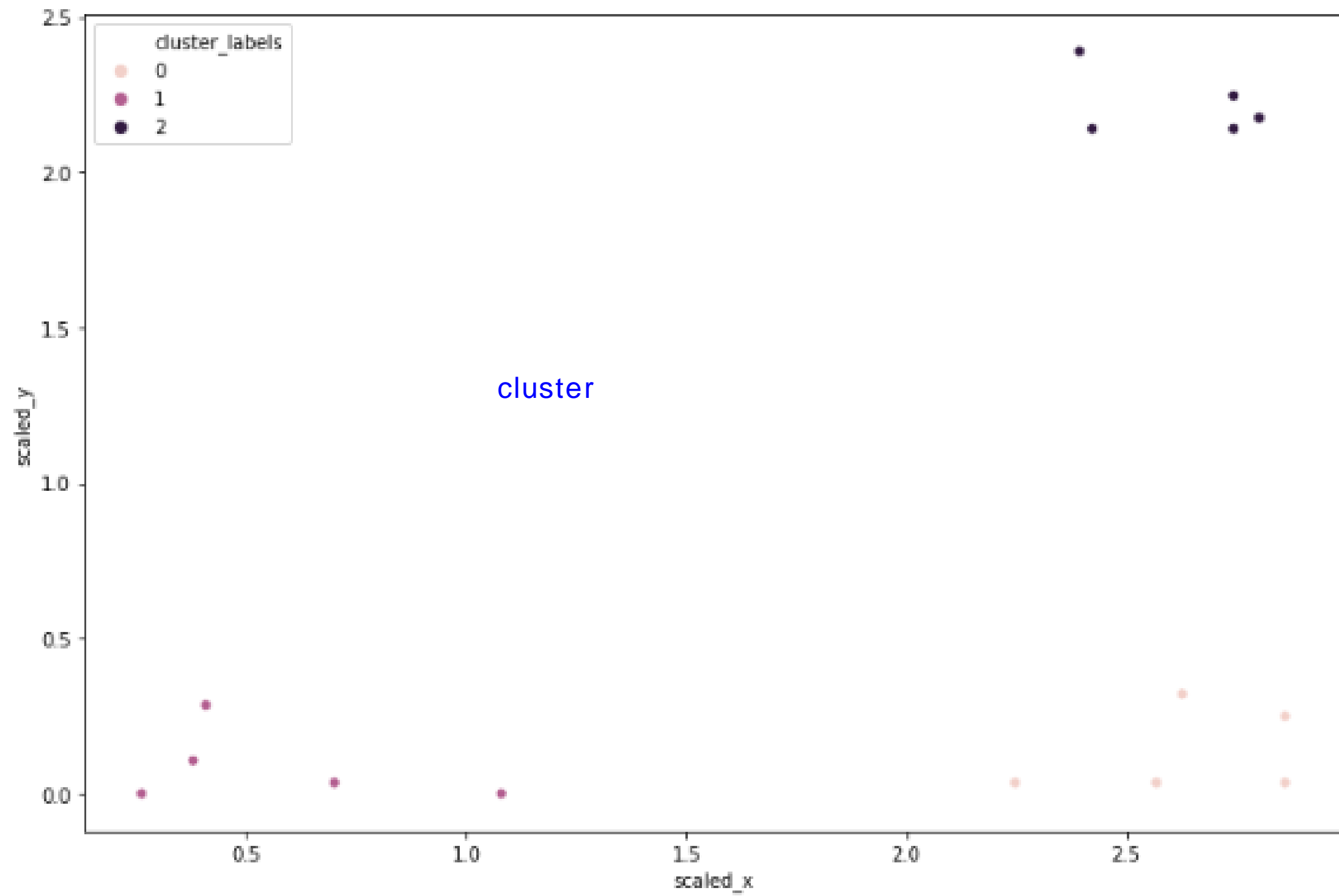vq                                    kmeans                    .

# Running k-means

```python
# Import kmeans and vq functions
from scipy.cluster.vq import kmeans, vq
```

```python
# Generate cluster centers and labels
cluster_centers, _ = kmeans(df[['scaled_x', 'scaled_y']], 3)
df['cluster_labels'], _ = vq(df[['scaled_x', 'scaled_y']], cluster_centers)
```

kmeans          cluster center
vq          cluster label

```python
# Plot clusters
sns.scatterplot(x='scaled_x', y='scaled_y', hue='cluster_labels', data=df)
plt.show()
```
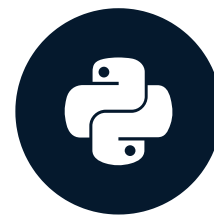
# Next up: exercises!
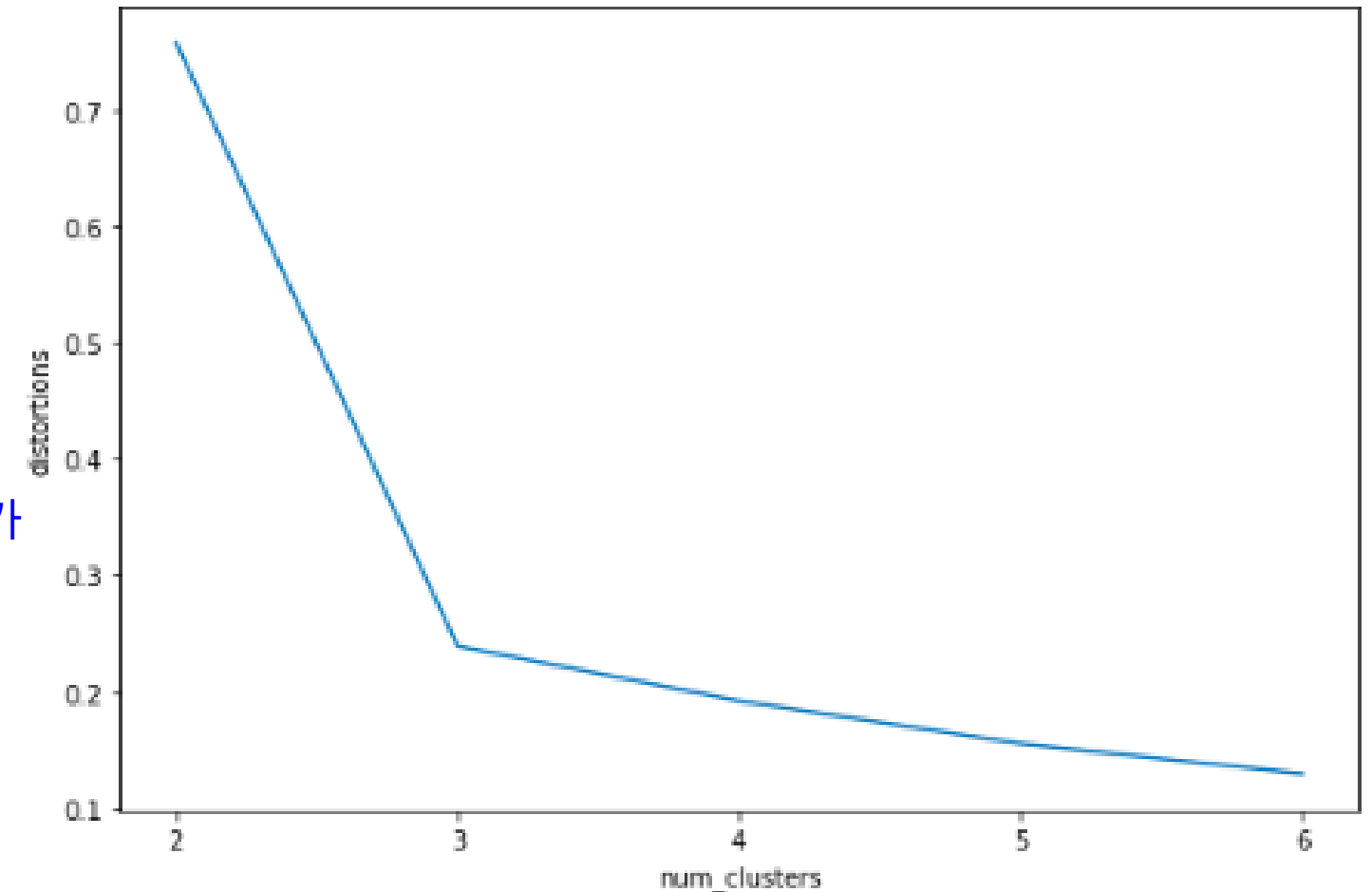
# How many clusters?

## CLUSTER ANALYSIS IN PYTHON

**Shaumik Daityari**
Business Analyst

k- means clustering      cluster

# How to find the right k?

- No *absolute* method to find right number of clusters (k) in k-means clustering

- Elbow method



k- means clustering                                    cluster

                                    ,

- >  elbow plot              dataset            cluster
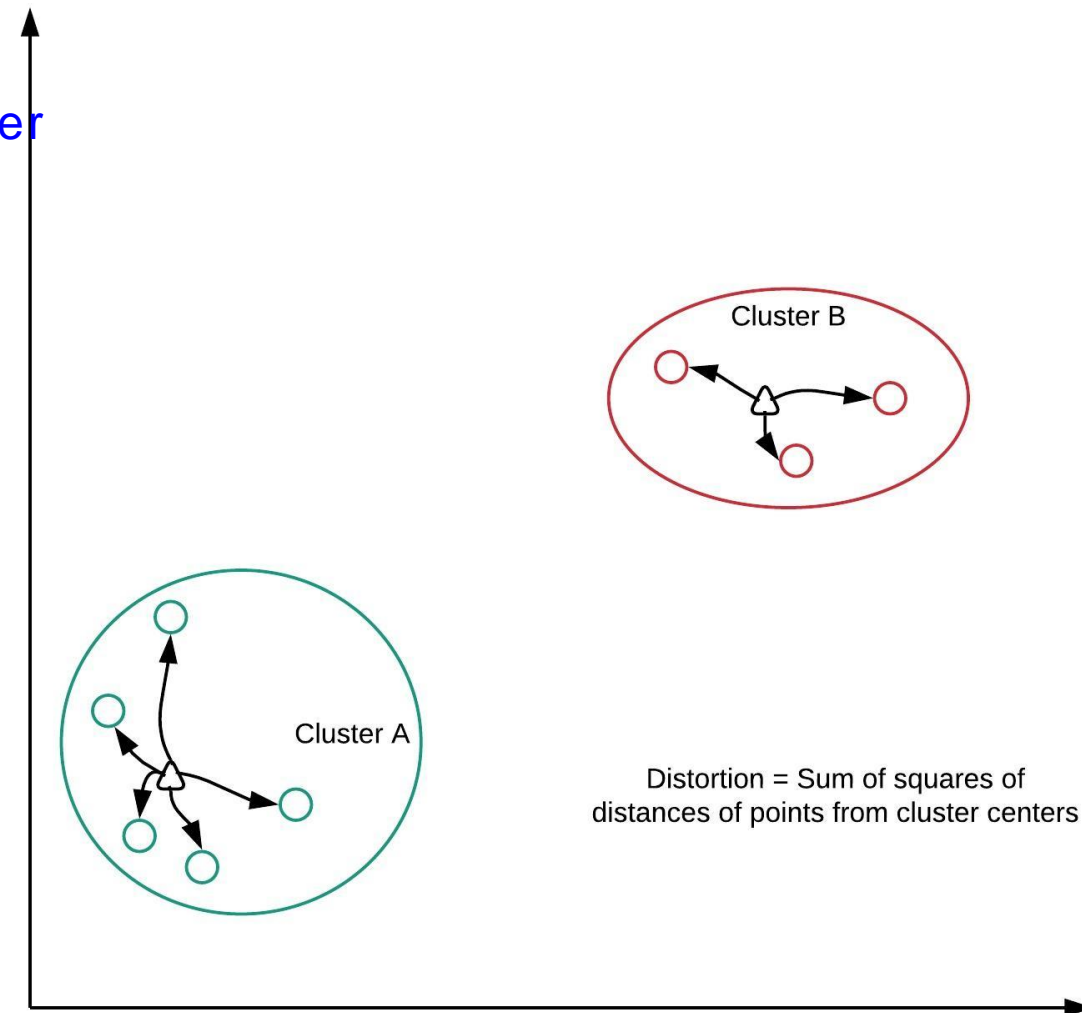
# Distortions revisited

- Distortion: sum of squared distances of points from cluster centers

- Decreases with an increasing number of clusters

- Becomes zero when the number of clusters equals the number of points

- Elbow plot: line plot between cluster centers and distortion



data point    cluster center

cluster    0    0

point

cluster

Cluster B

Cluster A

Distortion = Sum of squares of distances of points from cluster centers

# Elbow method

- Elbow plot: plot of the number of clusters and distortion

- Elbow plot helps indicate number of clusters present in data

-                               cluster    k- means clustering
    x      cluster      y                   elbow plot      (cluster      1~data point      )
-           plot    cluster

# Elbow method in Python

```python
# Declaring variables for use
distortions = []

num_clusters = range(2, 7)
```

```python
# Populating distortions for various clusters
for i in num_clusters:
    centroids, distortion = kmeans(df[['scaled_x', 'scaled_y']], i)
    distortions.append(distortion)
```

k- means

```python
# Plotting elbow plot data
elbow_plot_data = pd.DataFrame({'num_clusters': num_clusters,
                                'distortions': distortions})

sns.lineplot(x='num_clusters', y='distortions',
             data = elbow_plot_data)
plt.show()
```
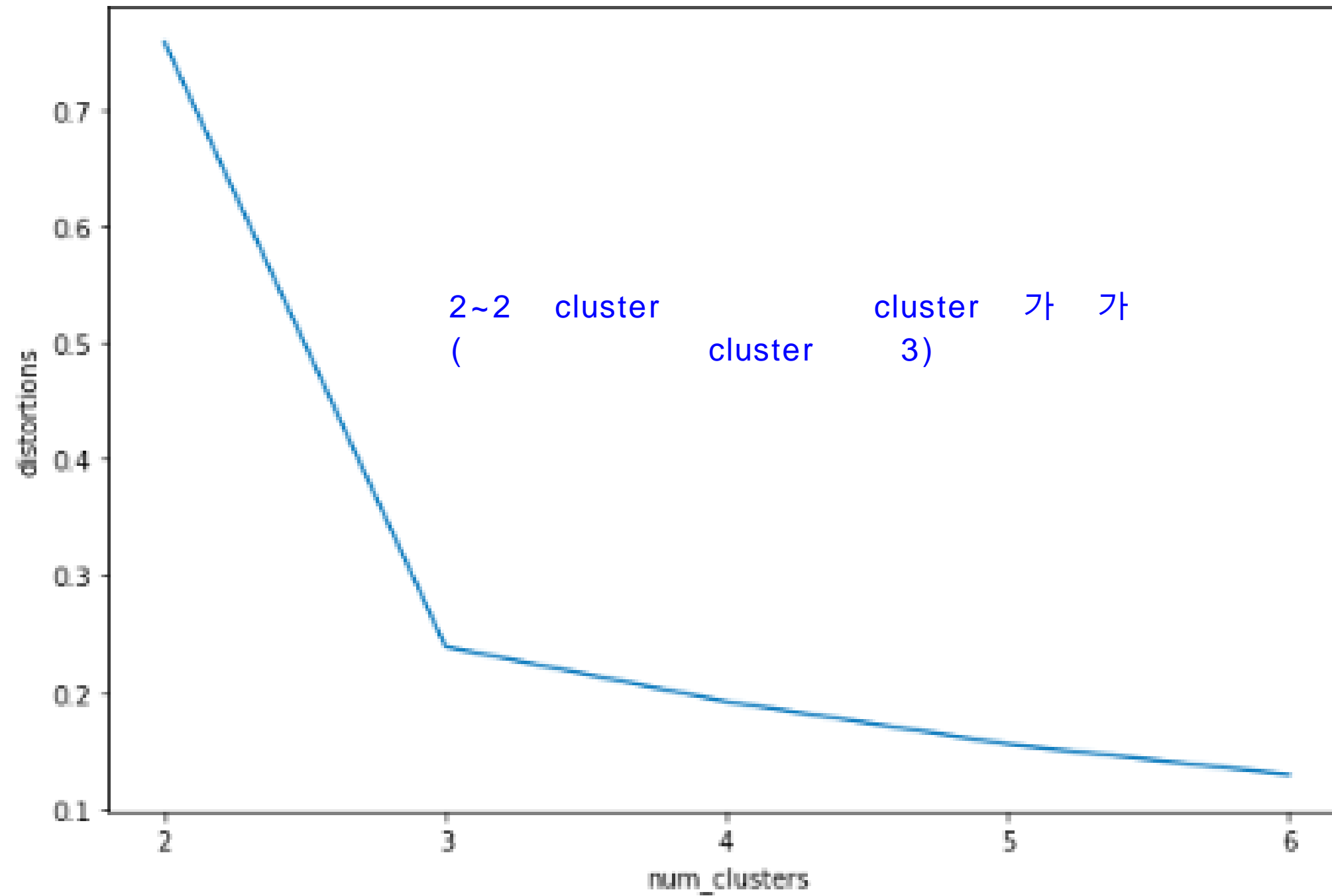
df

2~2   cluster               cluster
(                 cluster       3)

# Final thoughts on using the elbow method

- Only gives an indication of optimal _k_ (numbers of clusters)

- Does not always pinpoint how many _k_ (numbers of clusters)

- Other methods: average silhouette and gap statistic

- elbow                     cluster
                     k

ex) elbow

-

= >                                cluster

# Next up: exercises

## CLUSTER ANALYSIS IN PYTHON

# Limitations of k-means clustering

## CLUSTER ANALYSIS IN PYTHON

**Shaumik Daityari**
Business Analyst

k-

# Limitations of k-means clustering

- How to find the right _K_ (number of clusters)?

- Impact of seeds

- Biased towards equal sized clusters

k- means clustering                 clustering                         runtime

,

-                ,              cluster     k                  . elbow method                    k
- k- means clustering                    seed     clustering
-                      :                    cluster

# Impact of seeds

seed    cluster

Initialize a random seed

```
from numpy import random

random.seed(12)
```

cluster center
        cluster                                    .
=>                                  k- means clustering

seed    numpy
        1D array
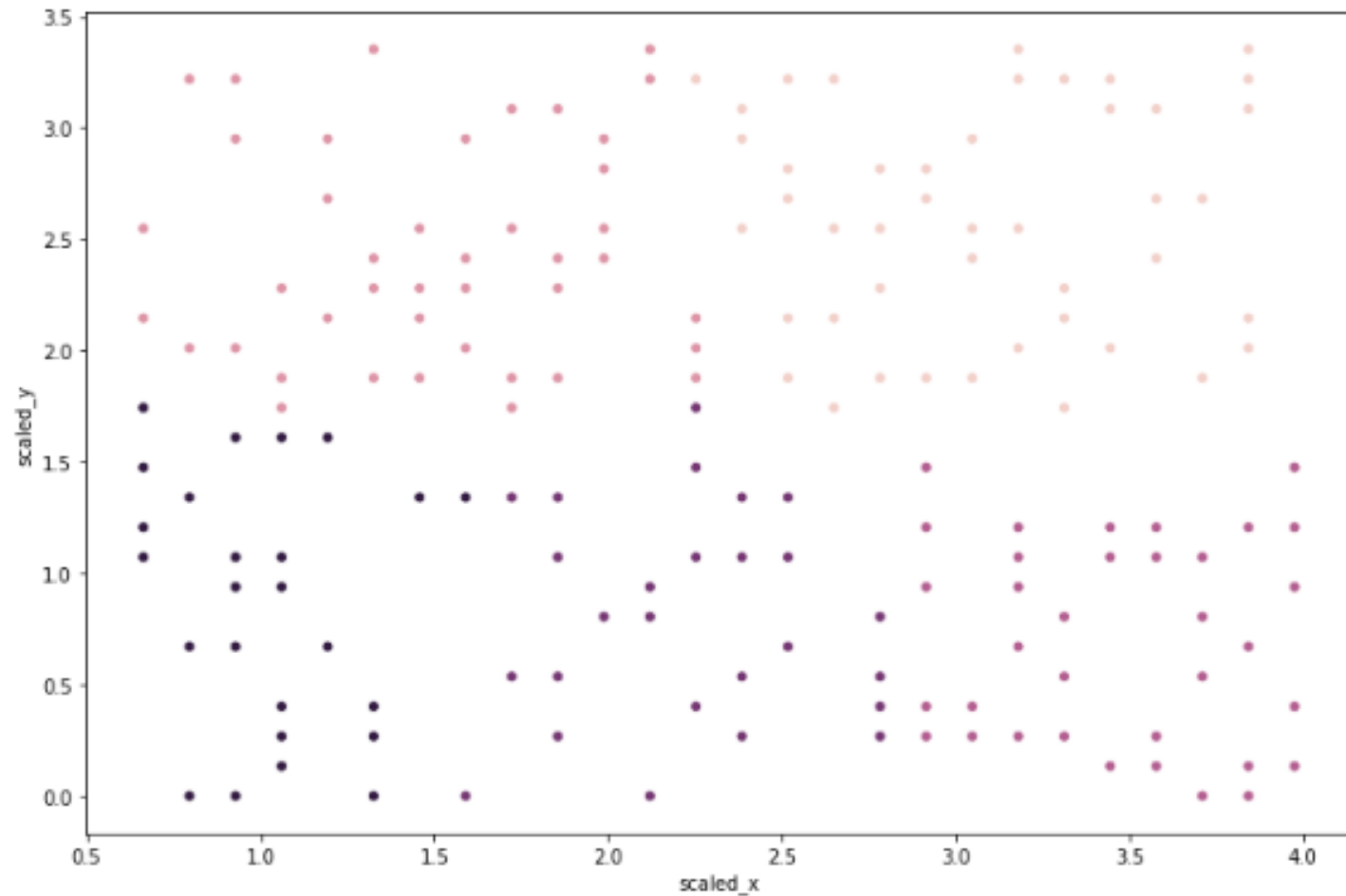                            seed                    k- means clustering
            200        point                    5        cluster              .
        cluster

Seed: `np.array(1000, 2000)`

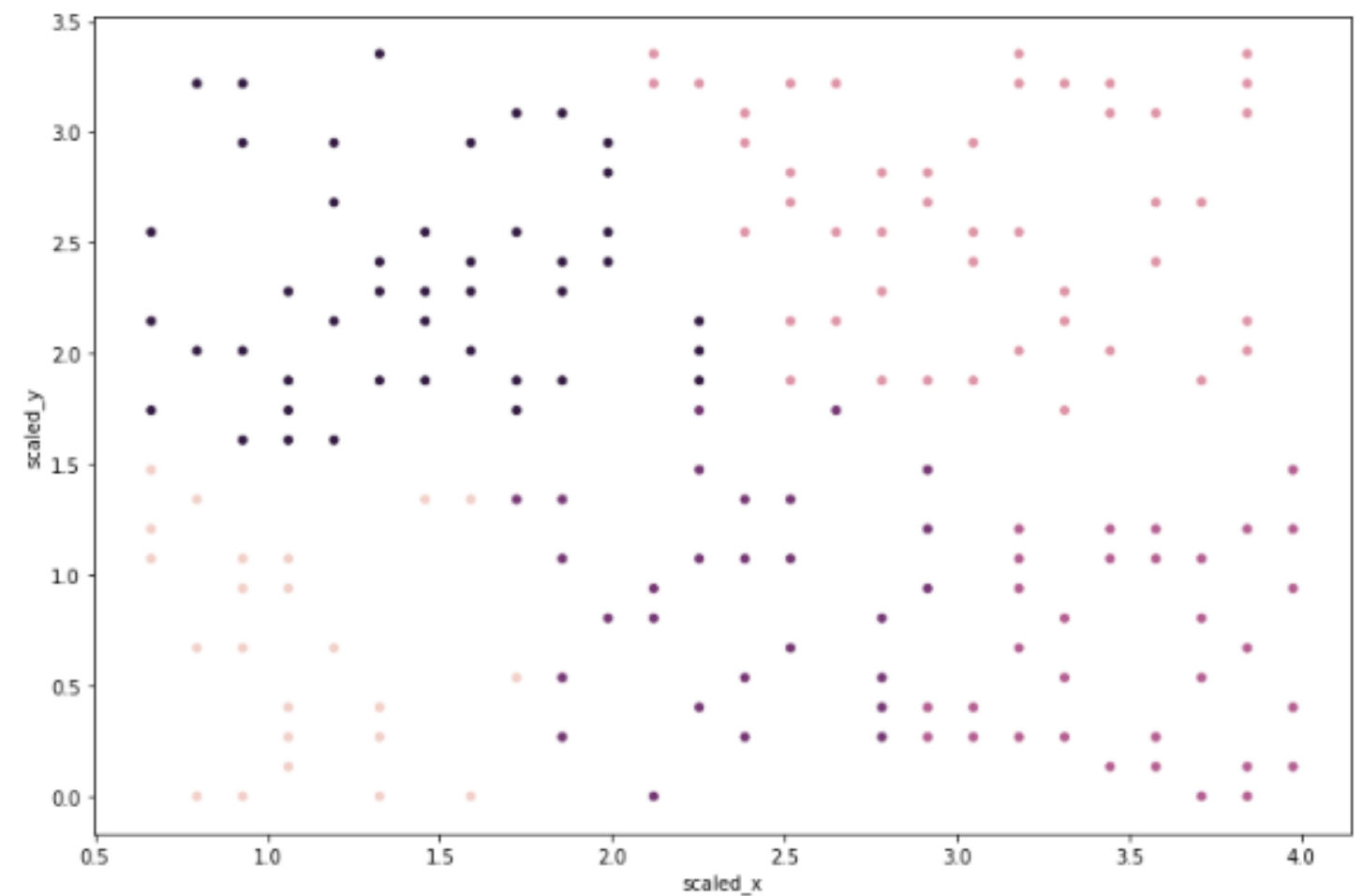Cluster sizes: 29, 29, 43, 47, 52

Seed: `np.array(1,2,3)`

Cluster sizes: 26, 31, 40, 50, 53
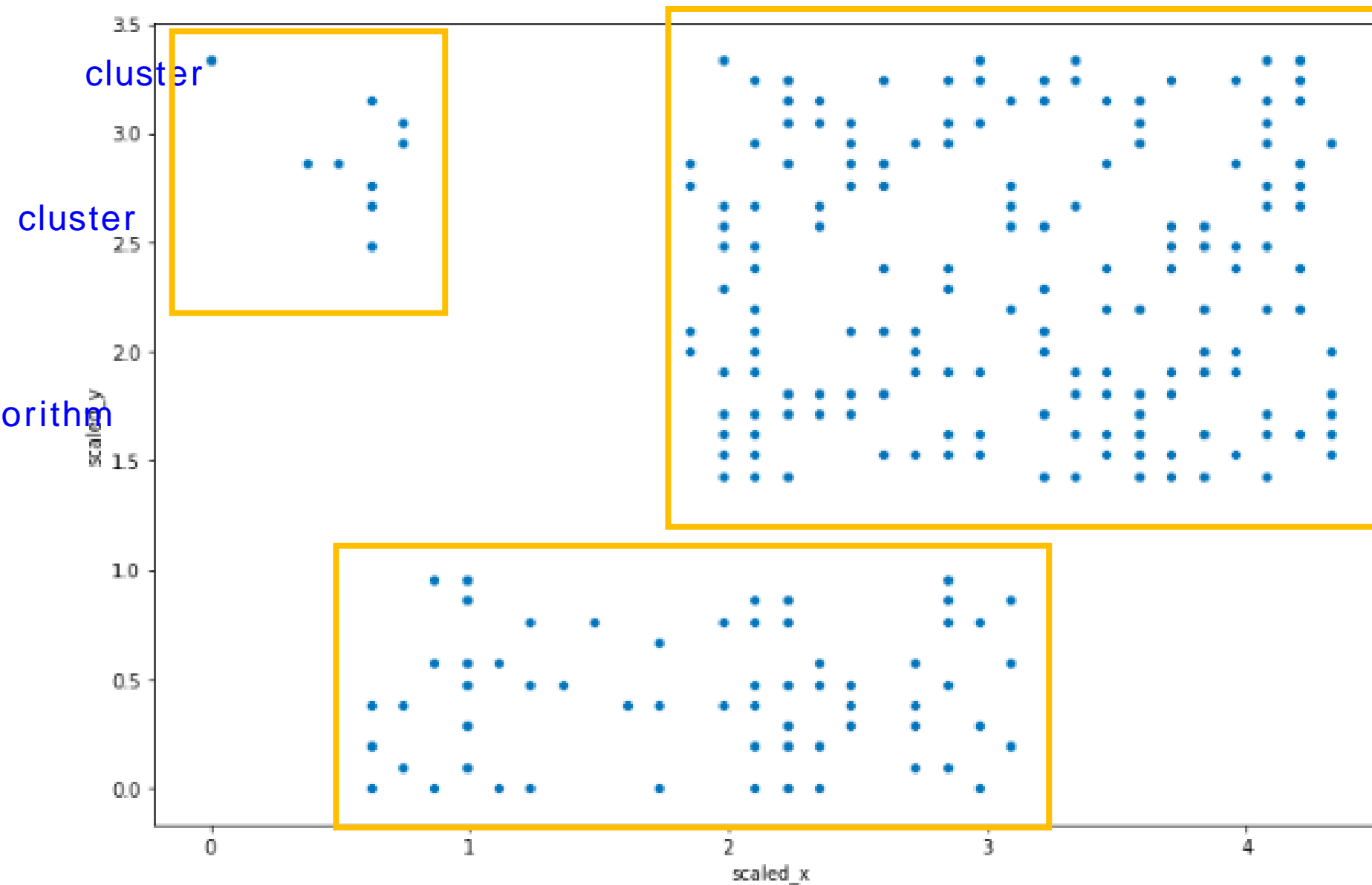
# Impact of seeds: plots

Seed: `np.array(1000, 2000)`

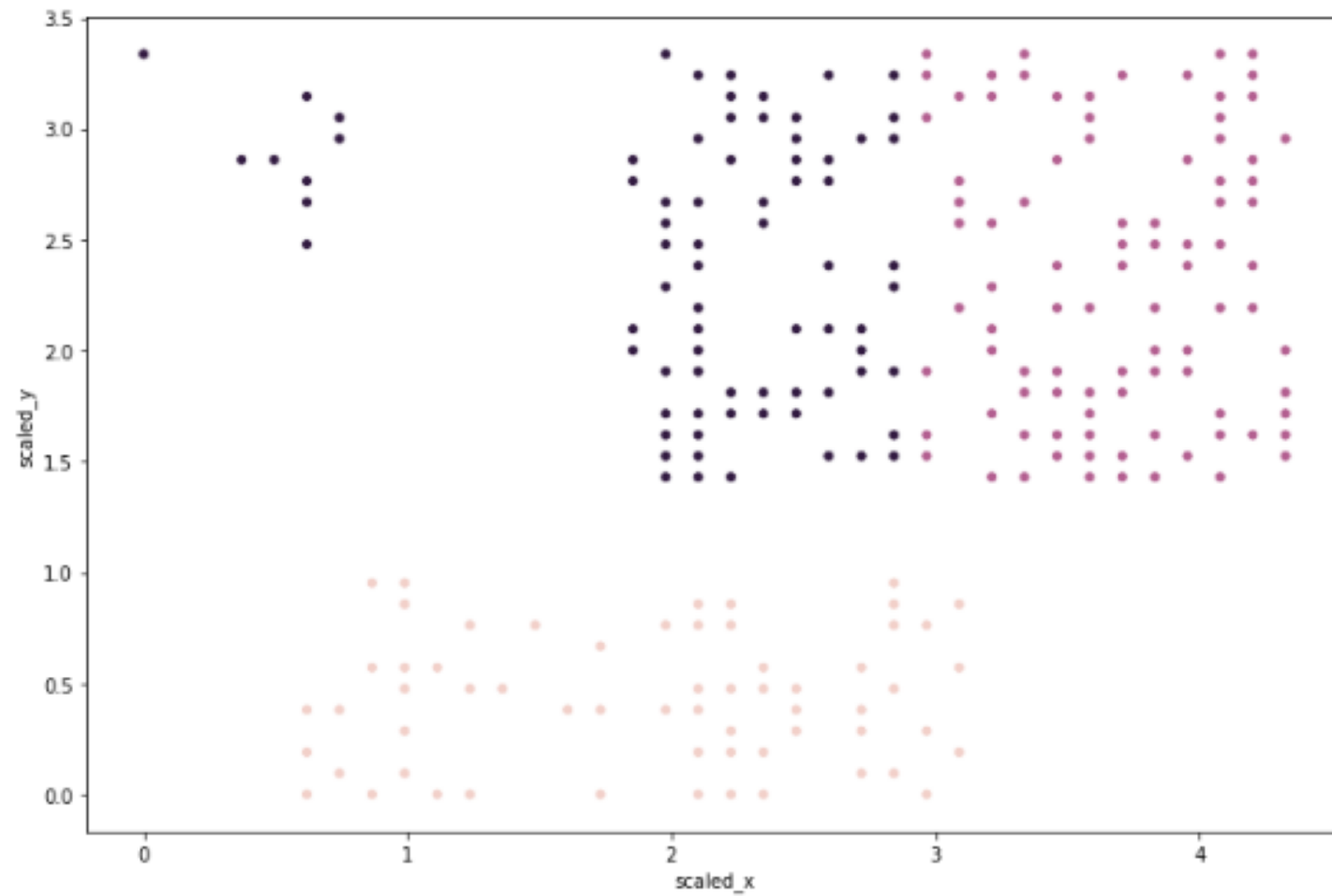Seed: `np.array(1,2,3)`

# Uniform clusters in k means
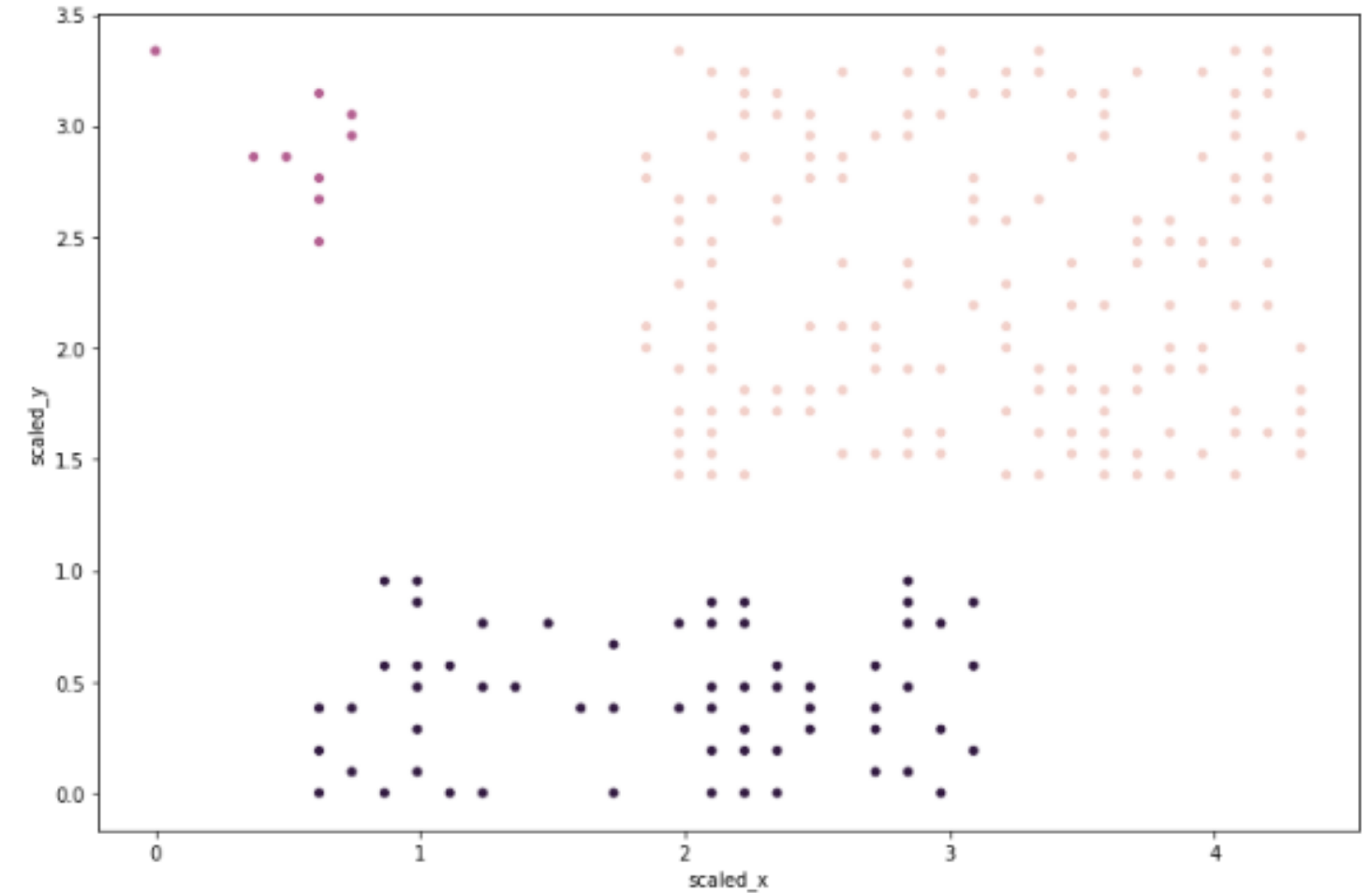
# Uniform clusters in k-means: a comparison

K-means clustering with 3 clusters

Hierarchical clustering with 3 clusters



k- means cluster                         seed

    : k- means cluster

= >            cluster

clustering                      cluster                      slide

# Final thoughts

- Each technique has its pros and cons

- Consider your data size and patterns before deciding on algorithm

- Clustering is exploratory phase of analysis

clustering

# Next up: exercises

## CLUSTER ANALYSIS IN PYTHON