# INTRODUCTION TO UNSUPERVISED LEARNING

# TYPES OF MACHINE LEARNING

| | |
|---|---|
| **SUPERVISED** | Data points have known outcome |
| **UNSUPERVISED** | Data points have unknown outcome |

# TYPES OF MACHINE LEARNING

| | |
|---|---|
| **SUPERVISED** | Data points have known outcome |
| **UNSUPERVISED** | Data points have unknown outcome |

# TYPES OF UNSUPERVISED LEARNING

**CLUSTERING**     Identify unknown structure in data
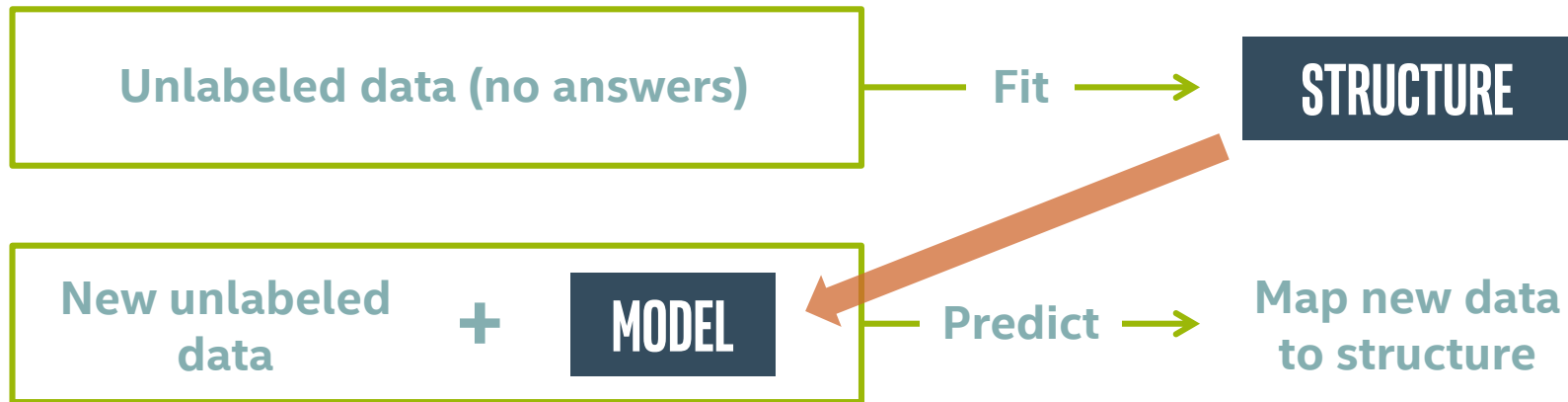
# TYPES OF UNSUPERVISED LEARNING
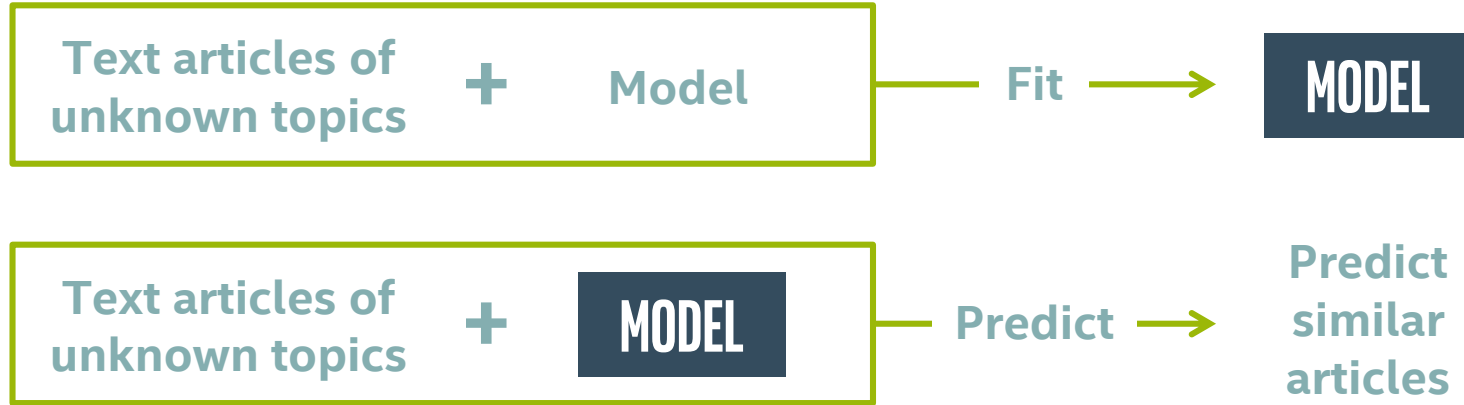
| CLUSTERING | Identify unknown structure in data |
|---|---|
| **DIMENSIONALITY REDUCTION** | Use structural characteristics to simplify data |

# UNSUPERVISED LEARNING OVERVIEW

Unlabeled data (no answers) → Fit → **STRUCTURE**

New unlabeled data + **MODEL** → Predict → Map new data to structure

# CLUSTERING: FINDING DISTINCT GROUPS

Text articles of unknown topics **+** Model → Fit → **MODEL**

Text articles of unknown topics **+** **MODEL** → Predict → Predict similar articles

# DIMENSIONALITY REDUCTION: SIMPLIFYING STRUCTURE

High resolution images + Model — Fit → **MODEL**

High resolution images + **MODEL** — Predict → Compressed images

# INTRODUCTION TO UNSUPERVISED LEARNING

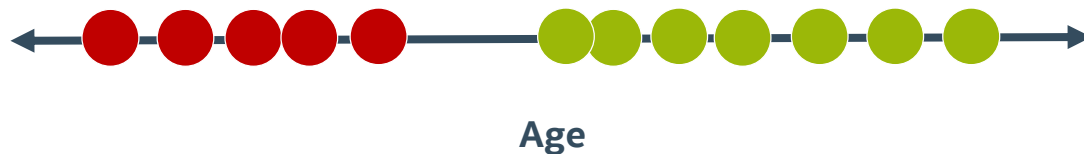**Users of a web application:**

- One feature (age)



**Age**

# INTRODUCTION TO UNSUPERVISED LEARNING

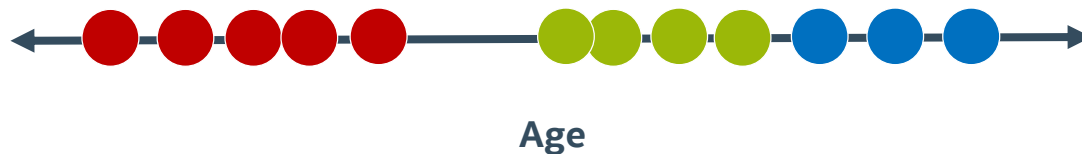**Users of a web application:**

- One feature (age)

- Two clusters



**Age**

# INTRODUCTION TO UNSUPERVISED LEARNING
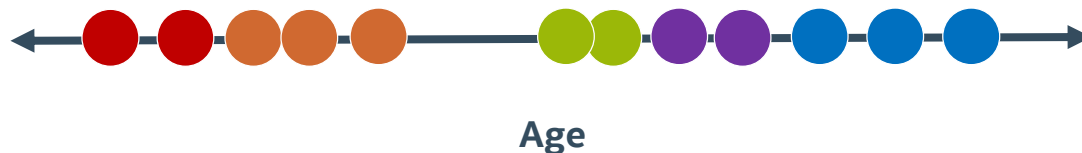
**Users of a web application:**

- One feature (age)

- Three clusters
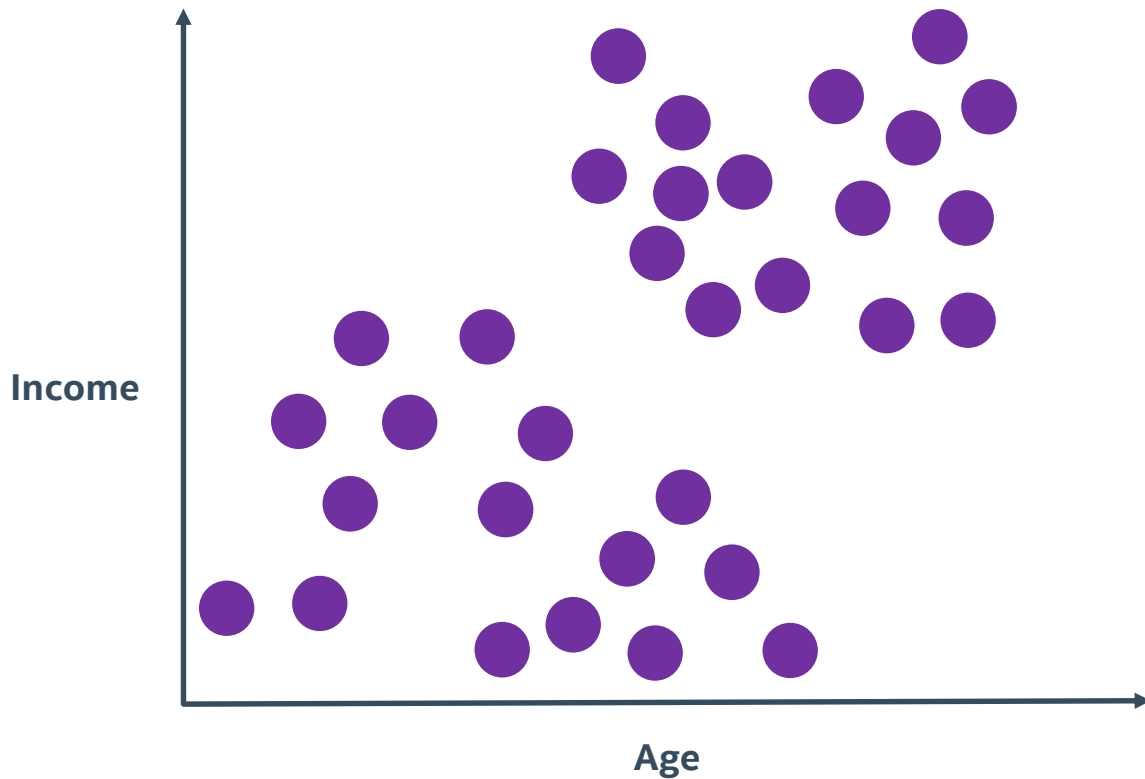
**Age**

# INTRODUCTION TO UNSUPERVISED LEARNING

**Users of a web application:**

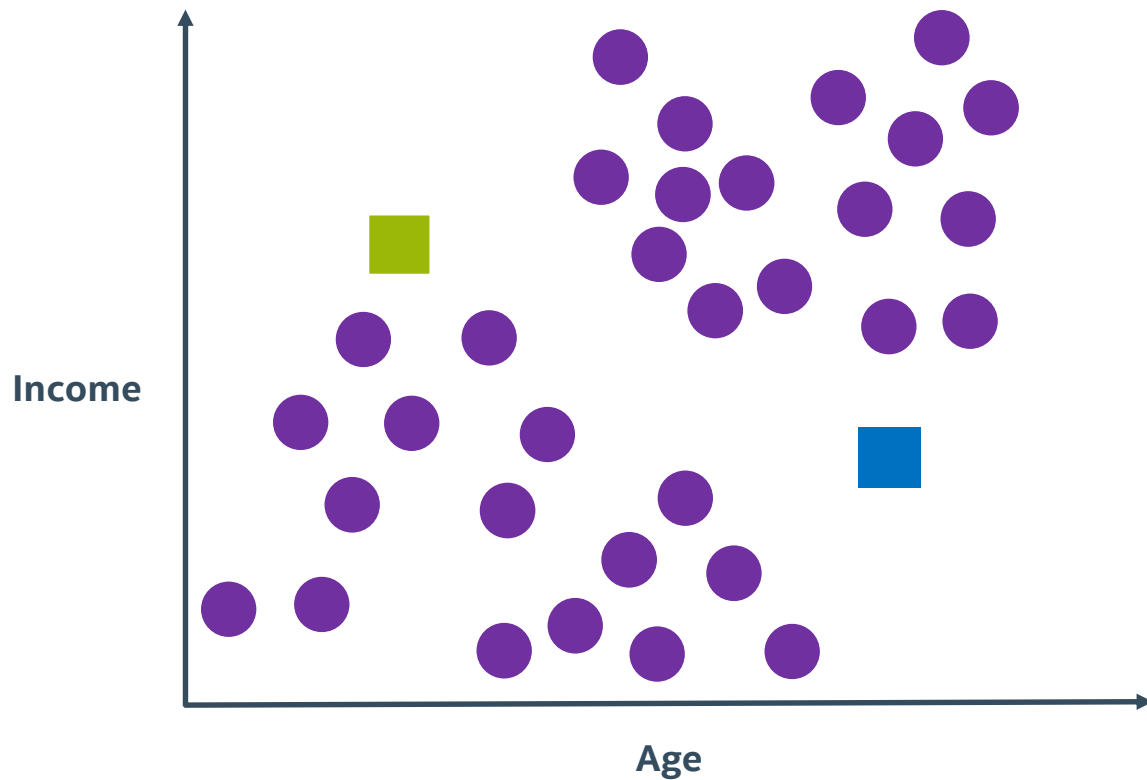- One feature (age)

- Five clusters

**Age**

# K-MEANS ALGORITHM
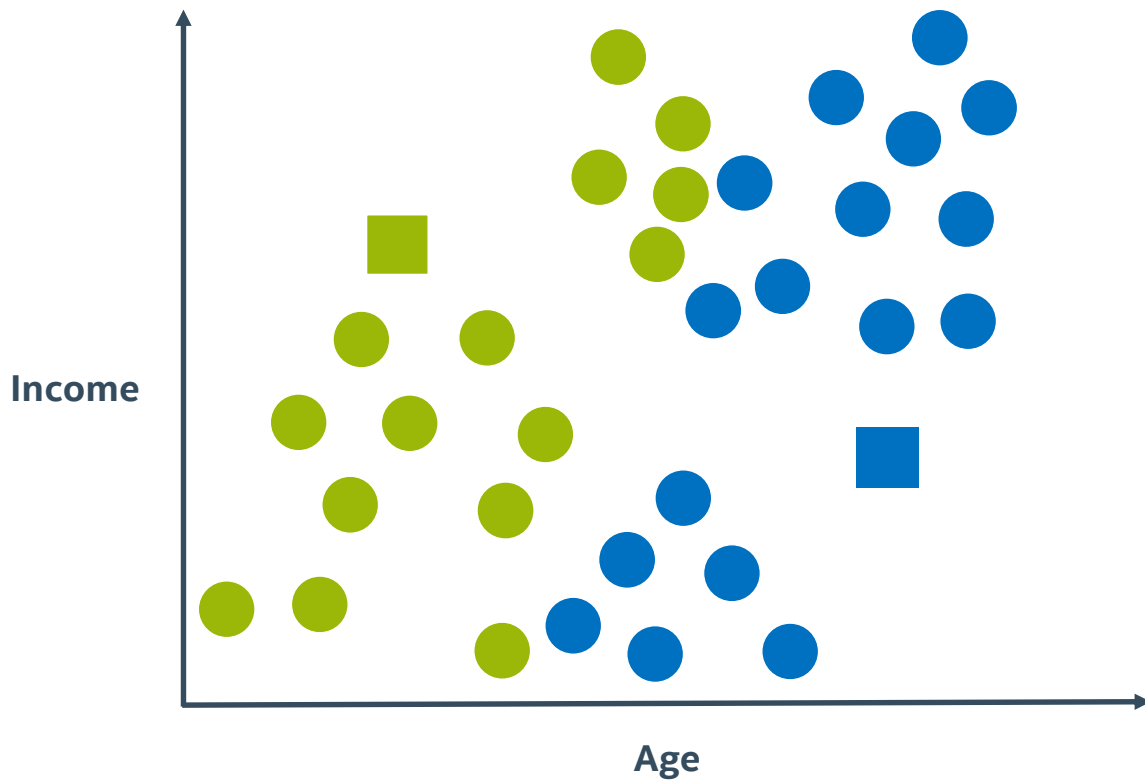
**K = 2 (find two clusters).**
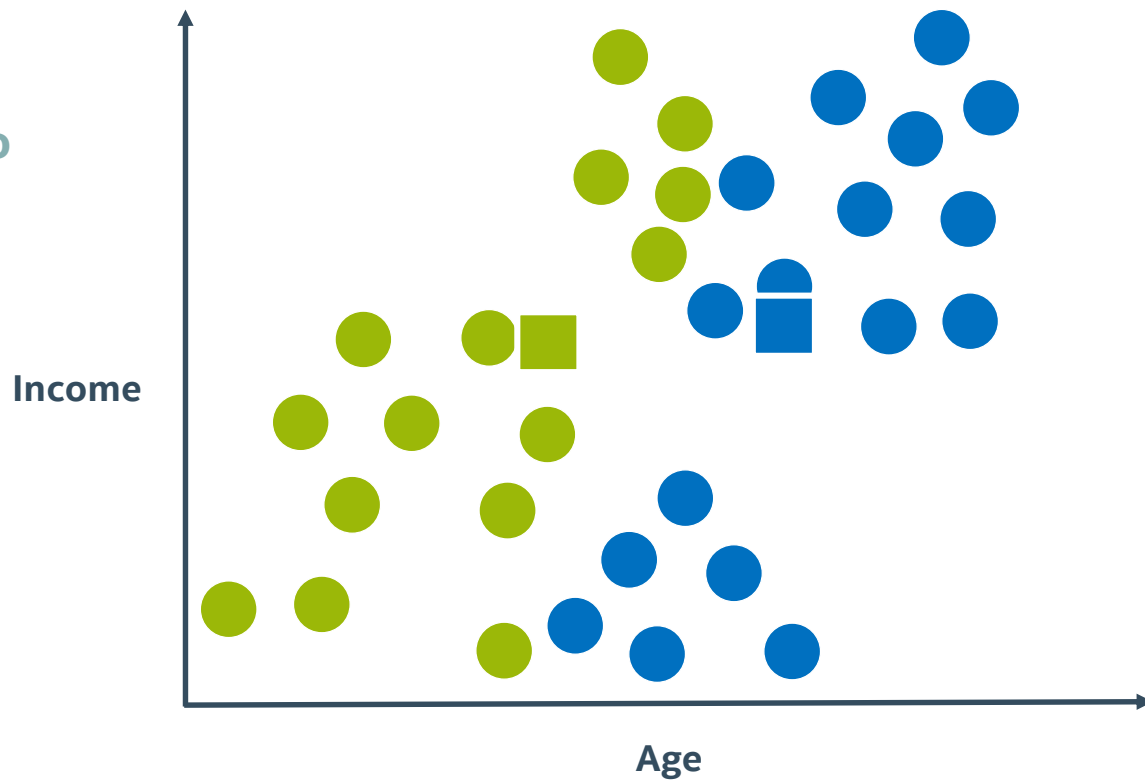
# K-MEANS ALGORITHM

**K = 2, Randomly assign cluster centers.**

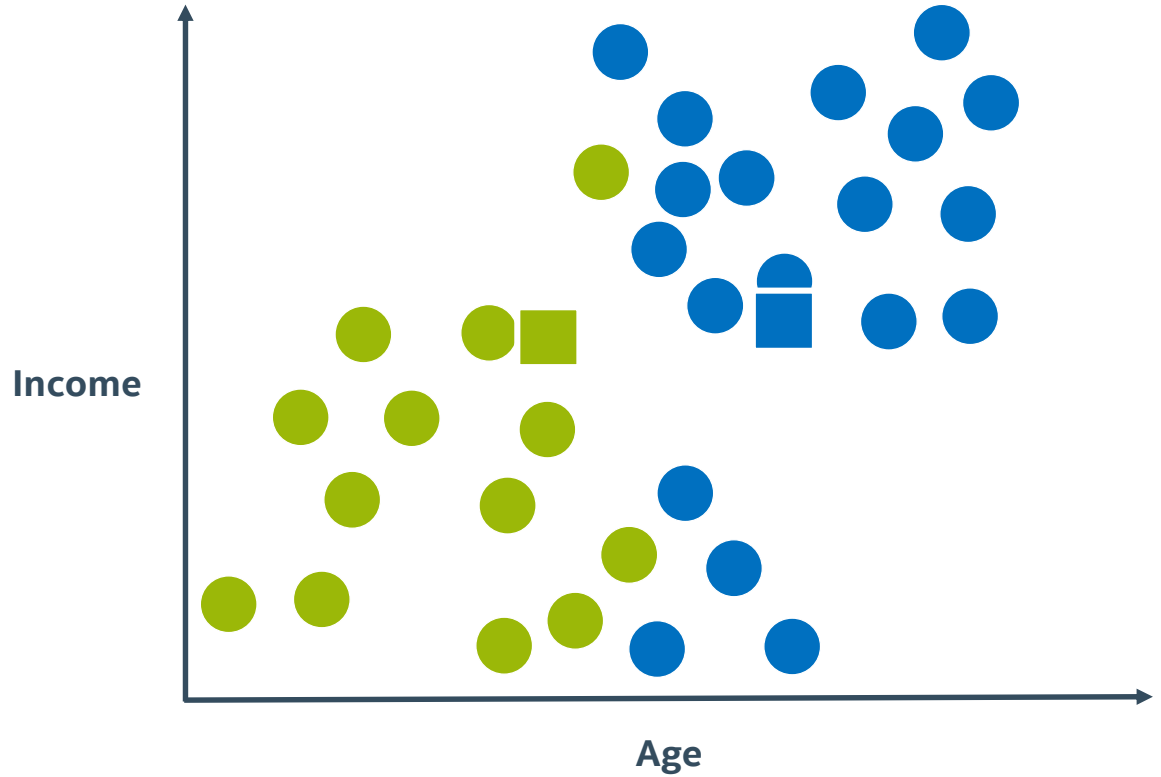# K-MEANS ALGORITHM

**K = 2, Each point belongs to closest center.**

# K-MEANS ALGORITHM

**K = 2, Each point belongs to closest center.**
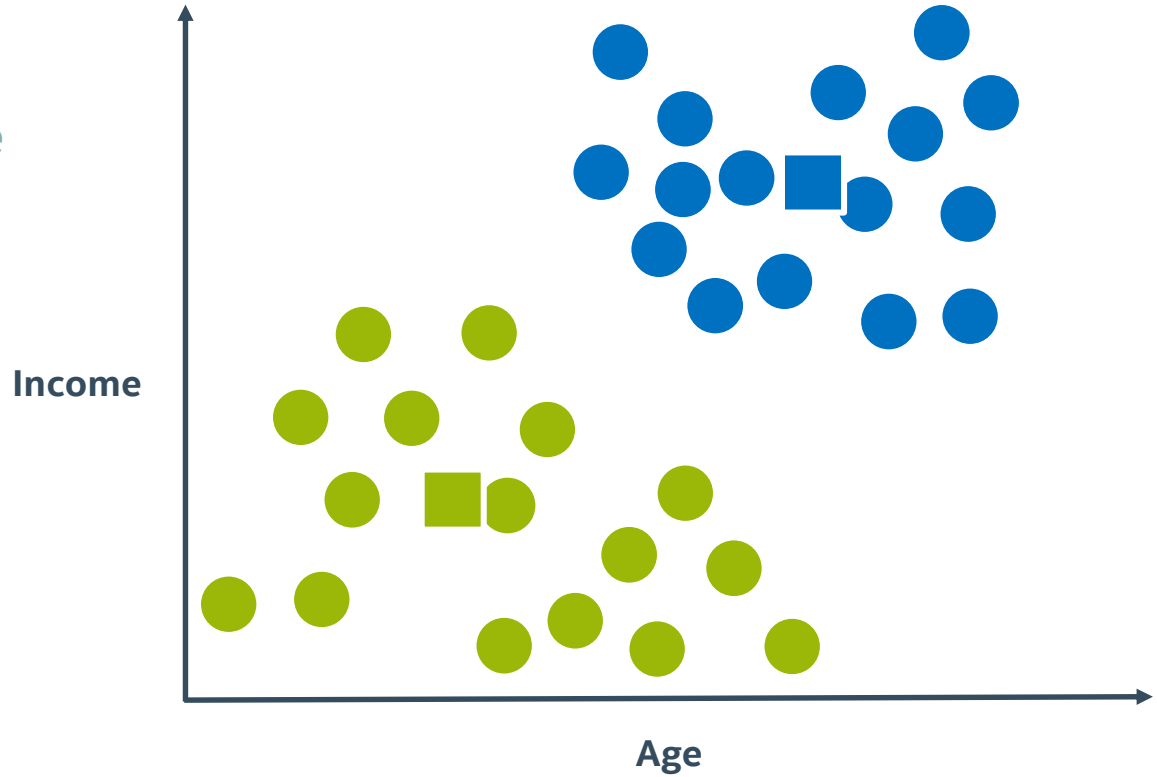


Income

Age

# K-MEANS ALGORITHM

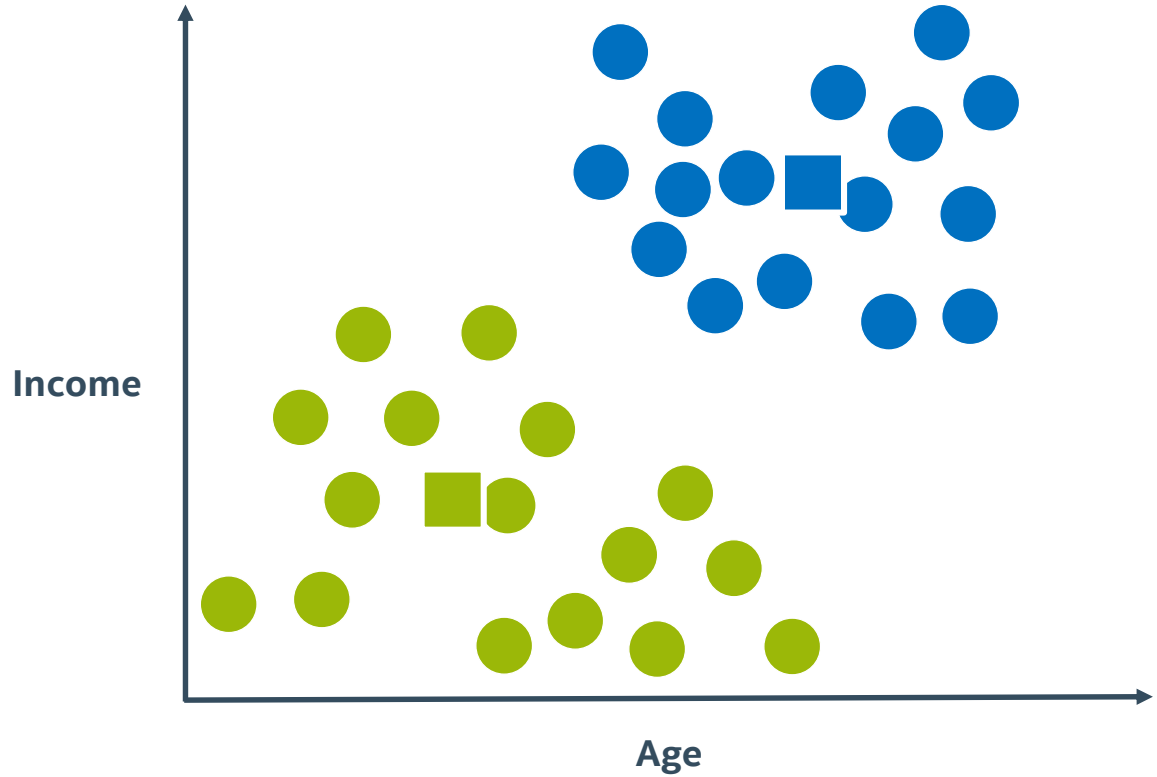**K = 2, Move each center to cluster's mean.**

# K-MEANS ALGORITHM

**K = 2, Points don't change**
**→ Converged.**
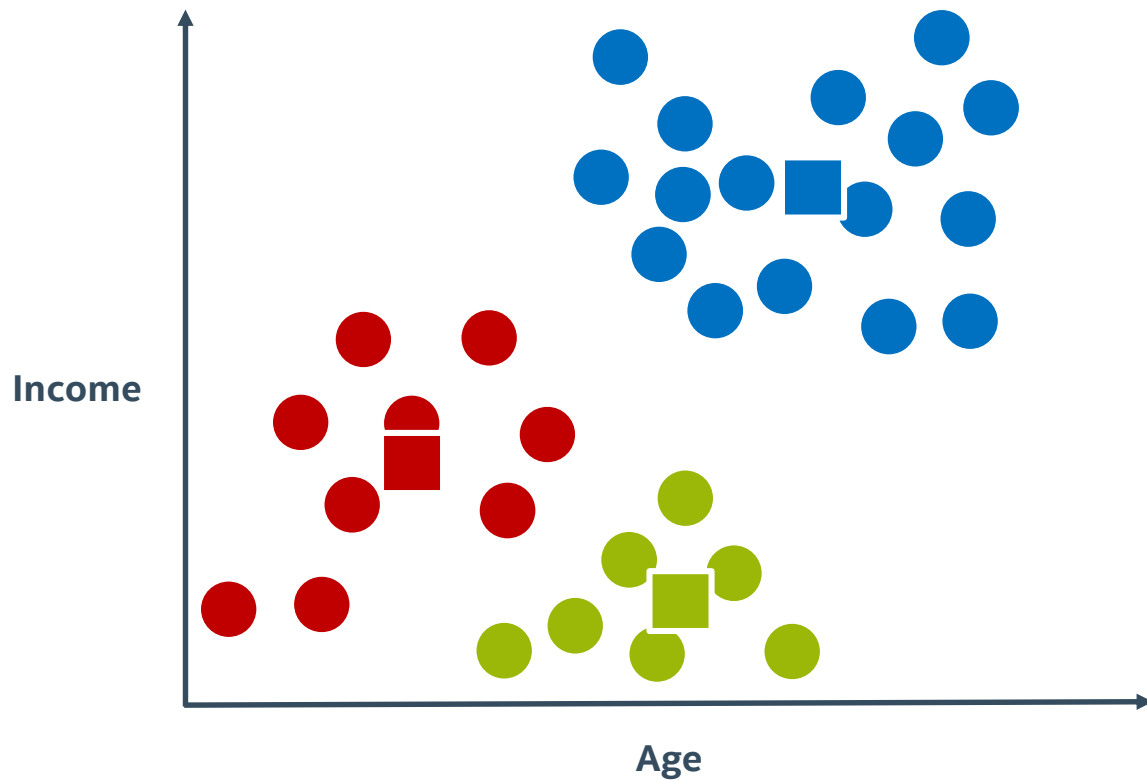
# K-MEANS ALGORITHM

**K = 2, Each point belongs to closest center.**



Income

Age

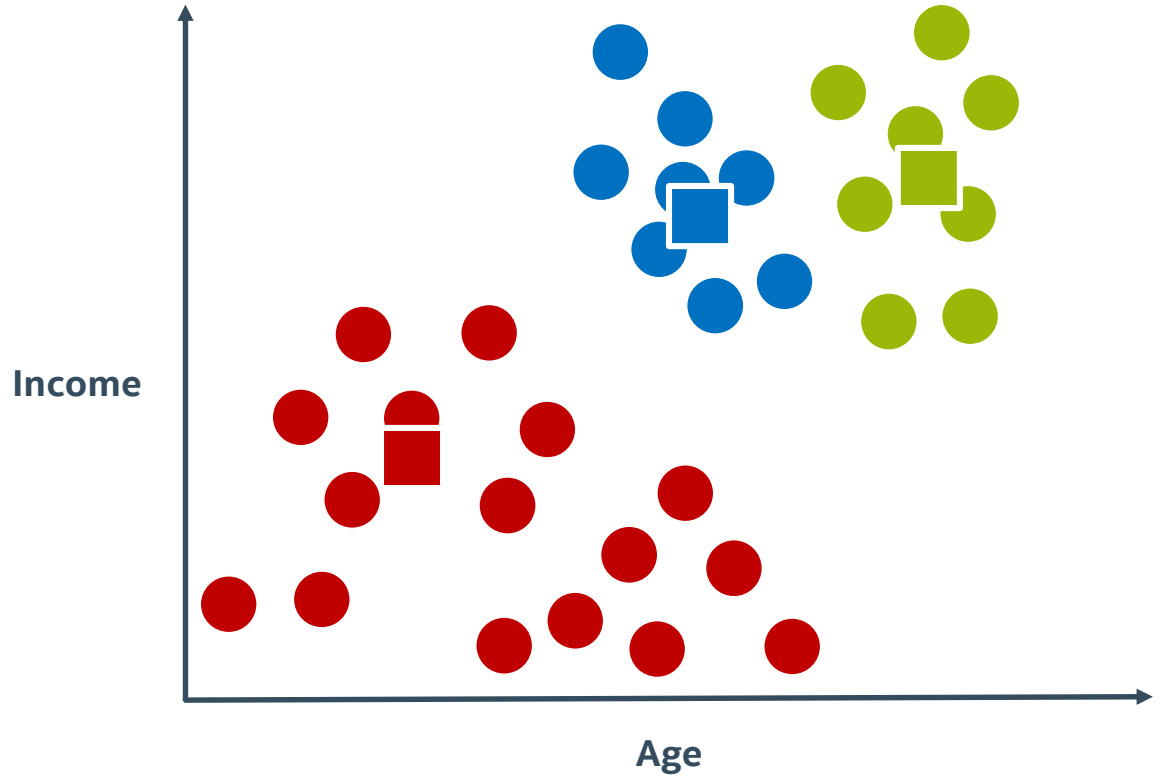# K-MEANS ALGORITHM

**K = 3**

# K-MEANS ALGORITHM

**K = 3, Results depend on initial cluster assignment.**

Income

Age

# WHICH MODEL IS THE RIGHT ONE?

# WHICH MODEL IS THE RIGHT ONE?

- **Inertia:** sum of squared distance from each point ($x_i$) to its cluster ($C_k$)

$$\sum_{i=1}^{n} (x_i - C_k)^2$$

- Smaller value corresponds to tighter clusters

- Other metrics can also be used

# WHICH MODEL IS THE RIGHT ONE?

**Initiate multiple times, take model with the best score.**

# WHICH MODEL IS THE RIGHT ONE?

**Inertia = 12.645**

Income

Age

# WHICH MODEL IS THE RIGHT ONE?

Inertia = 12.943

Income

Age

# WHICH MODEL IS THE RIGHT ONE?

Inertia = 13.112

# SMARTER INITIALIZATION OF K-MEANS CLUSTERS

# SMARTER INITIALIZATION OF K-MEANS CLUSTERS

**Pick one point at random as initial point.**

# SMARTER INITIALIZATION OF K-MEANS CLUSTERS

**Pick next point with 1/distance² probability.**

Income

Age

# SMARTER INITIALIZATION OF K-MEANS CLUSTERS

**Pick next point with 1/distance² probability.**

# SMARTER INITIALIZATION OF K-MEANS CLUSTERS

**Pick next point with 1/distance² probability.**



Income

Age

# SMARTER INITIALIZATION OF K-MEANS CLUSTERS

**Assign clusters.**



Income

Age

# CHOOSING THE RIGHT NUMBER OF CLUSTERS

# CHOOSING THE RIGHT NUMBER OF CLUSTERS

- Sometimes the question has a K

# CHOOSING THE RIGHT NUMBER OF CLUSTERS

- Sometimes the question has a K

- Clustering similar jobs on 4 CPU cores (K=4)

# CHOOSING THE RIGHT NUMBER OF CLUSTERS

- Sometimes the question has a K

- Clustering similar jobs on 4 CPU cores (K=4)

- A clothing design in 10 different sizes to cover most people (K=10)

# CHOOSING THE RIGHT NUMBER OF CLUSTERS

- Sometimes the question has a K

- Clustering similar jobs on 4 CPU cores (K=4)

- A clothing design in 10 different sizes to cover most people (K=10)

- A navigation interface for browsing scientific papers with 20 disciplines (K=20)

# CHOOSING THE RIGHT NUMBER OF CLUSTERS

- Inertia measures distance of point to cluster

# CHOOSING THE RIGHT NUMBER OF CLUSTERS

- Inertia measures distance of point to cluster

- Value decreases with increasing K as long as cluster density increases

# K-MEANS: THE SYNTAX

**Import the class containing the clustering method.**

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

# K-MEANS: THE SYNTAX

**Import the class containing the clustering method.**

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

**Create an instance of the class.**

```
kmeans = KMeans(n_clusters=3,

                    init='k-means++')
```

# K-MEANS: THE SYNTAX

Import the class containing the clustering method.

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

Create an instance of the class.

```
kmeans = KMeans(n_clusters=3,

                init='k-means++')
```

← **final number of clusters**

# K-MEANS: THE SYNTAX

**Import the class containing the clustering method.**

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

**Create an instance of the class.**

```
kmeans = KMeans(n_clusters=3,

                init='k-means++')
```

kmeans++
cluster
initiation

# K-MEANS: THE SYNTAX

**Import the class containing the clustering method.**

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

**Create an instance of the class.**

```
kmeans = KMeans(n_clusters=3,

                   init='k-means++')
```

**Fit the instance on the data and then predict clusters for new data.**

```
kmeans = kmeans.predict(X1)
y_predict = kmeans.predict(X2)
```

# K-MEANS: THE SYNTAX

**Import the class containing the clustering method.**

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

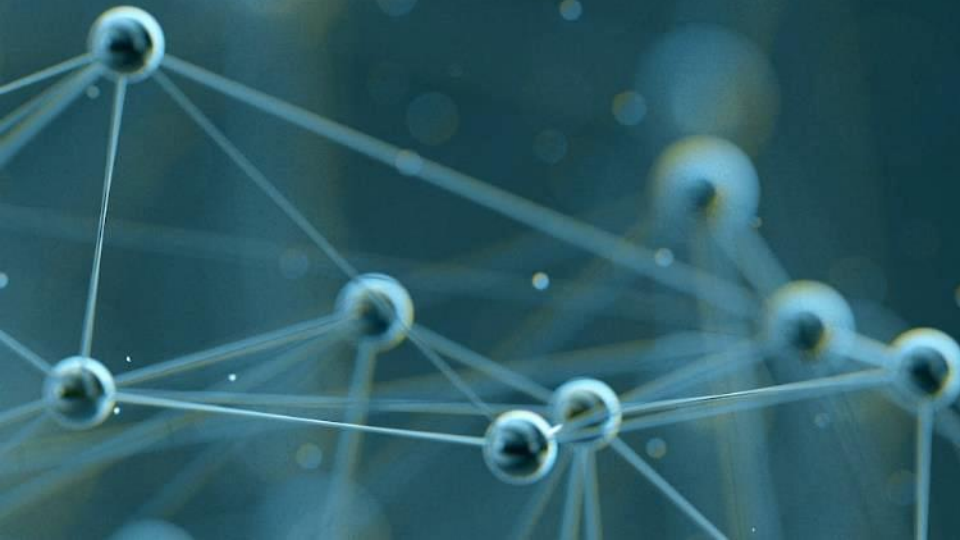**Create an instance of the class.**

```
kmeans = KMeans(n_clusters=3,

                  init='k-means++')
```

**Fit the instance on the data and then predict clusters for new data.**

```
kmeans = kmeans.predict(X1)
y_predict = kmeans.predict(X2)
```
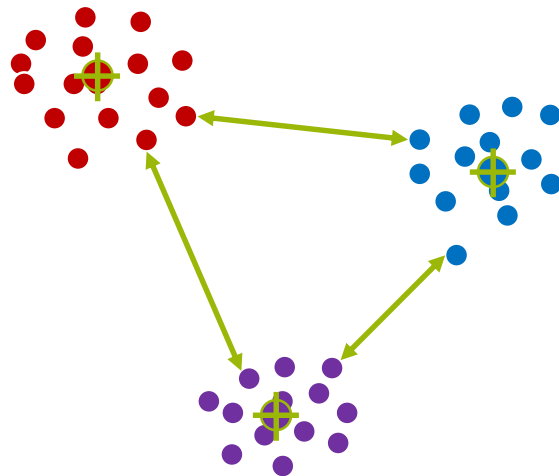
**Can also be used in batch mode with MiniBatchKMeans.**
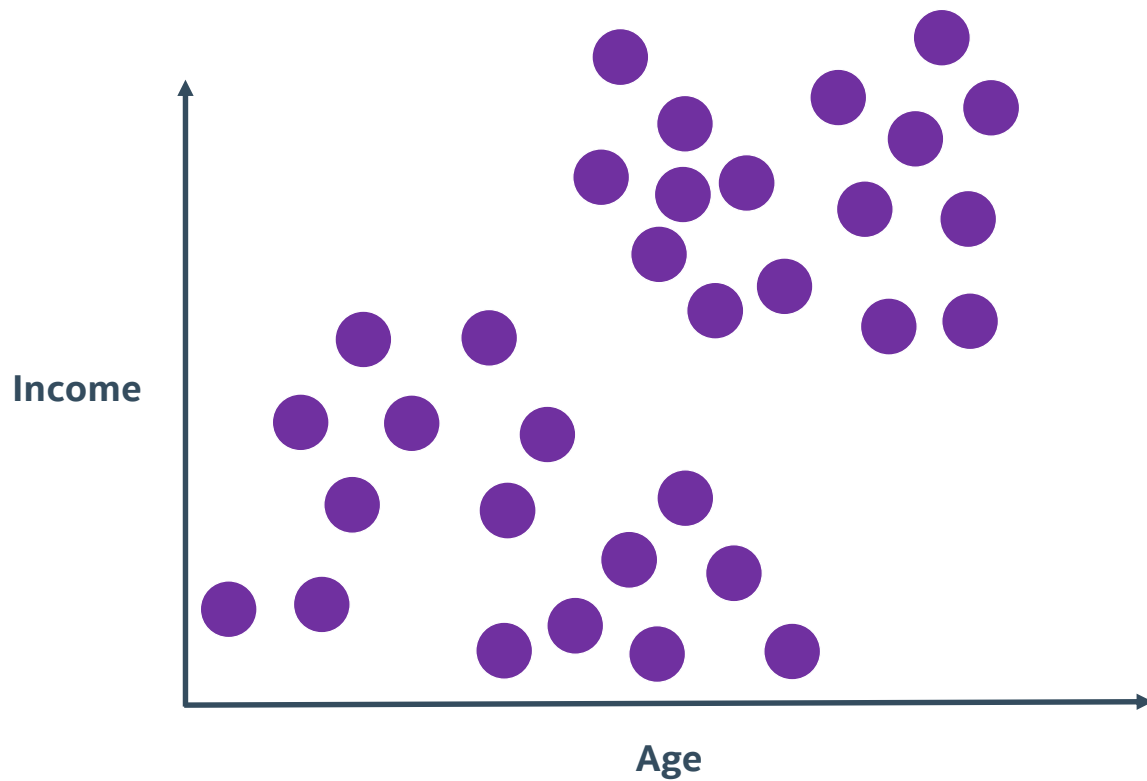
# DISTANCE METRICS
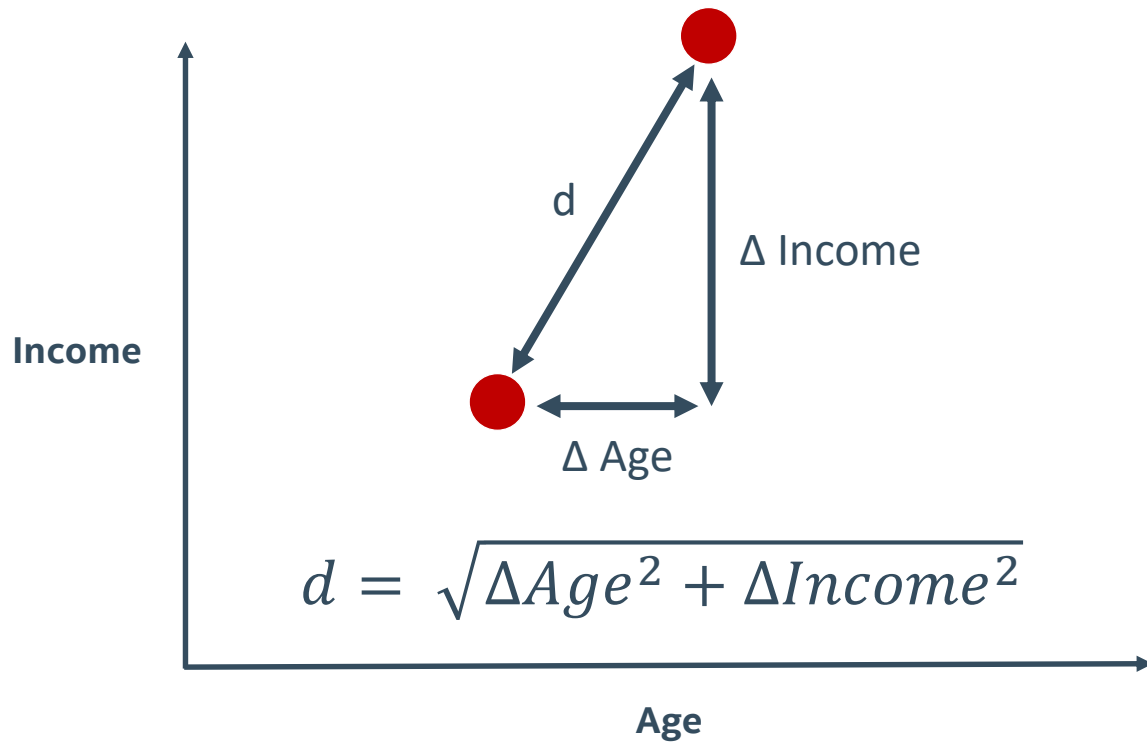
# DISTANCE METRIC CHOICE

- Choice of distance metric is extremely important to clustering success

- Each metric has strengths and most appropriate use-cases...

- ...but sometimes choice of distance metric is also based on empirical evaluation

# EUCLIDEAN DISTANCE
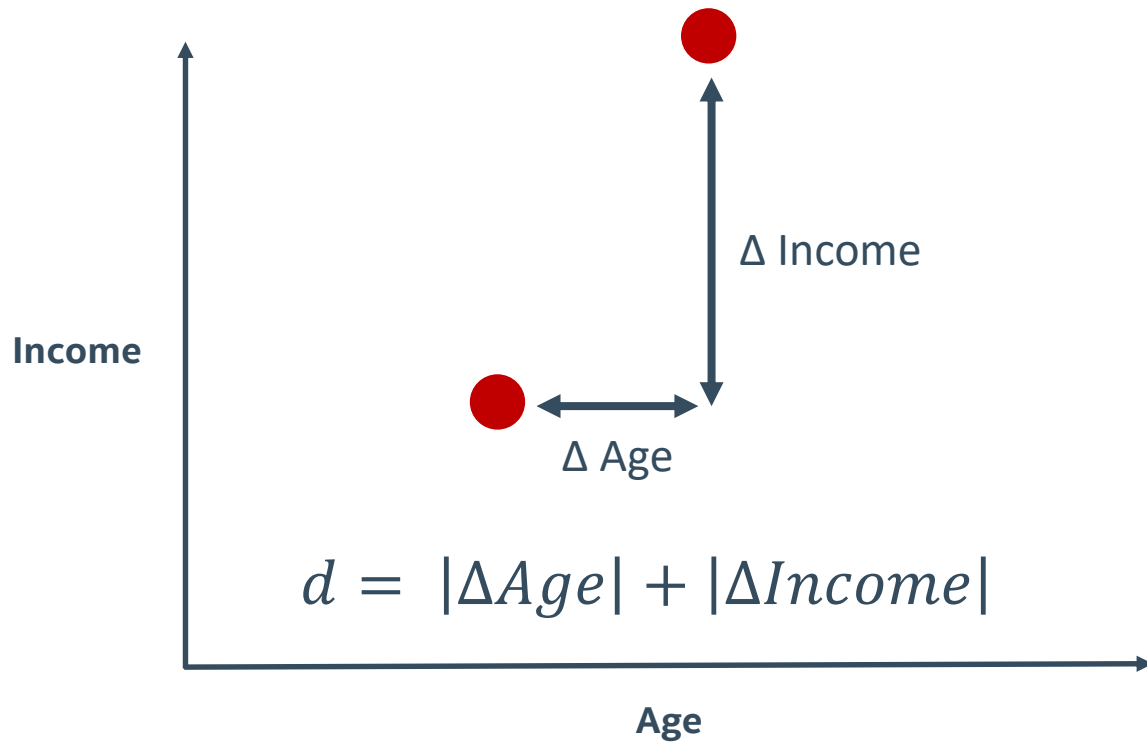
# EUCLIDEAN EUCLIDEAN DISTANCE (L2 DISTANCE)



d

Δ Income

Income

Δ Age

$$d = \sqrt{\Delta Age^2 + \Delta Income^2}$$

Age

# MANHATTAN DISTANCE (L1 OR CITY BLOCK DISTANCE)

Income

Δ Income

Δ Age

$$d = |\Delta Age| + |\Delta Income|$$

Age

# COSINE DISTANCE



$$\cos(\theta) = \frac{A \cdot B}{\|A\|\|B\|}$$

Income

Age

# COSINE DISTANCE

Income

$\theta$

$$\cos(\theta) = \frac{A \cdot B}{\|A\|\|B\|}$$

Age

# EUCLIDEAN VS COSINE DISTANCE

- Euclidean is useful for coordinate based measurements

# EUCLIDEAN VS COSINE DISTANCE

- Euclidean is useful for coordinate based measurements

- Cosine is better for data such as text where location of occurrence is less important

# EUCLIDEAN VS COSINE DISTANCE

- Euclidean is useful for coordinate based measurements

- Cosine is better for data such as text where location of occurrence is less important

- Euclidean distance is more sensitive to curse of dimensionality (see lesson 12)

# JACCARD DISTANCE

**Applies to sets (like word occurrence)**

- **Sentence A**: "I like chocolate ice cream."

- set A = {I, like, chocolate, ice, cream}

- **Sentence B**: "Do I want chocolate cream or vanilla cream?"

- set B = {Do, I, want, chocolate, cream, or, vanilla}

$$1 - \frac{A \cap B}{A \cup B} = 1 - \frac{len(shared)}{len(unique)}$$

# JACCARD DISTANCE

**Applies to sets (like word occurrence)**

- **Sentence A**: "I like chocolate ice cream."

- set A = {**I**, like, **chocolate**, ice, **cream**}

- **Sentence B**: "Do I want chocolate cream or vanilla cream?"

- set B = {Do, **I**, want, **chocolate**, **cream**, or, vanilla}

$$1 - \frac{A \cap B}{A \cup B} = 1 - \frac{3}{9}$$

# DISTANCE METRICS: THE SYNTAX

**Import the general pairwise distance function.**

```
from sklearn.metrics import pairwise_distances
```

# DISTANCE METRICS: THE SYNTAX

**Import the general pairwise distance function.**

```
from sklearn.metrics import pairwise_distances
```

**Calculate the distances.**

```
dist = pairwise_distances(X,Y,
                          metric='euclidean')
```

# DISTANCE METRICS: THE SYNTAX

**Import the general pairwise distance function.**

```
from sklearn.metrics import pairwise_distances
```
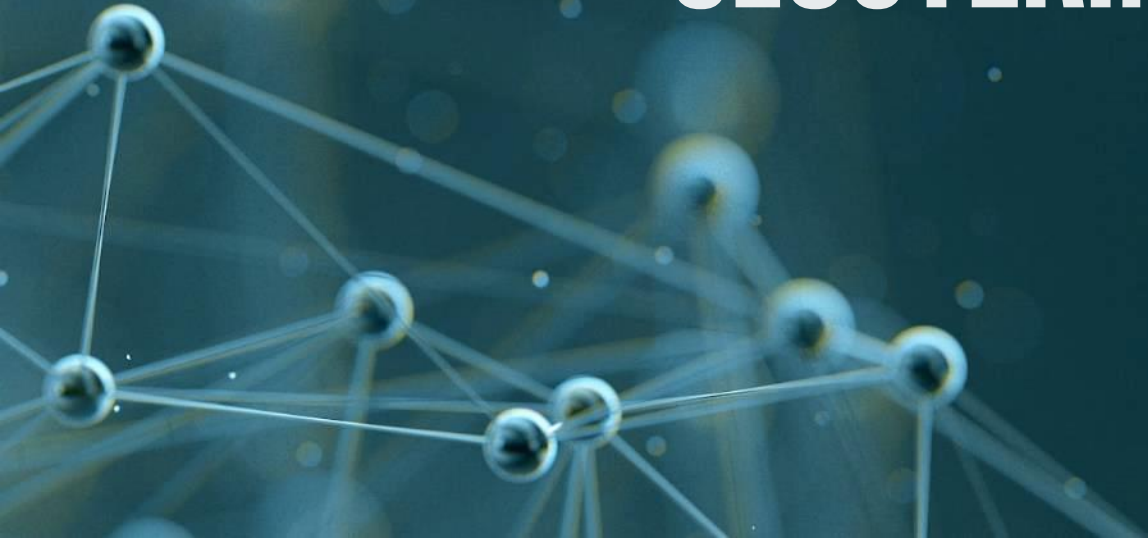
**Calculate the distances.**

```
dist = pairwise_distances(X,Y,

                  metric='euclidean')
```

← **distance metric choice**

# DISTANCE METRICS: THE SYNTAX

**Import the general pairwise distance function.**

```
from sklearn.metrics import pairwise_distances
```

**Calculate the distances.**

```
dist = pairwise_distances(X,Y,

                          metric='euclidean')
```

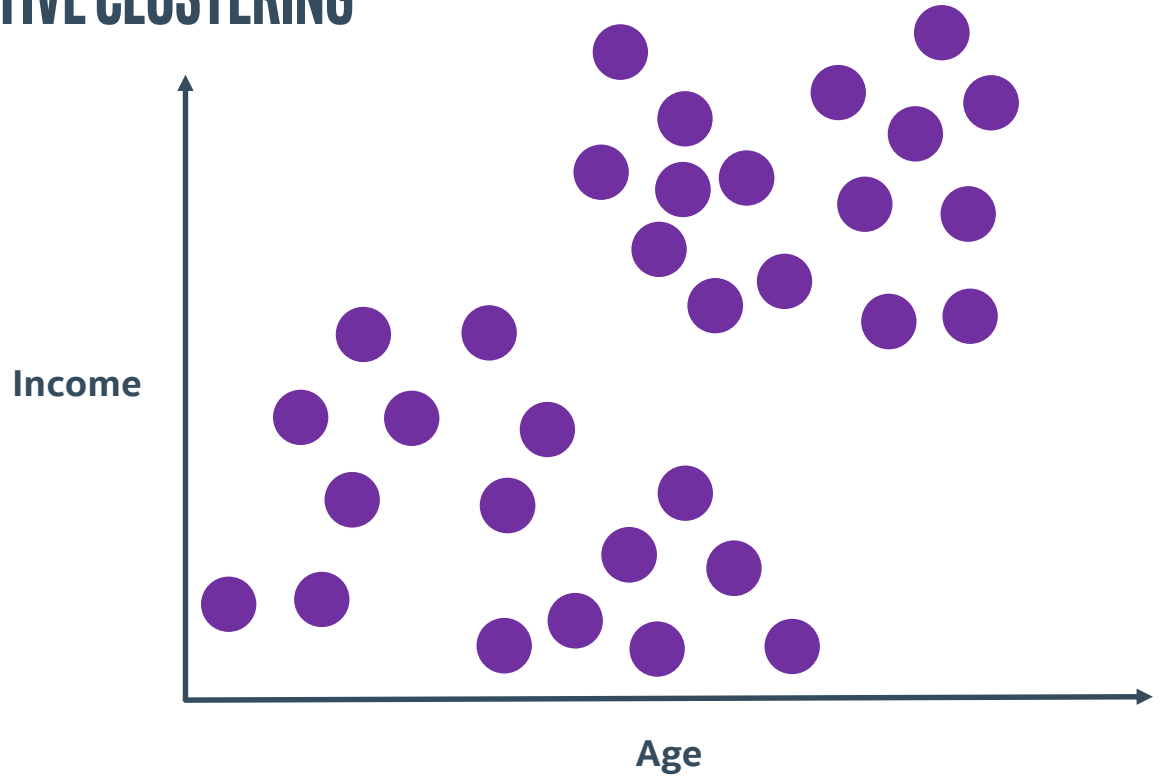**Other distance metric choices are: cosine, manhattan, jaccard, etc.**

# DISTANCE METRICS: THE SYNTAX

**Import the general pairwise distance function.**

```
from sklearn.metrics import pairwise_distances
```

**Calculate the distances.**

```
dist = pairwise_distances(X,Y,

                    metric='euclidean')
```

**Other distance metric choices are: cosine, manhattan, jaccard, etc.**

**Distance metric methods can also be imported specifically, e.g.:**

```
from sklearn.metrics import euclidean_distances
```
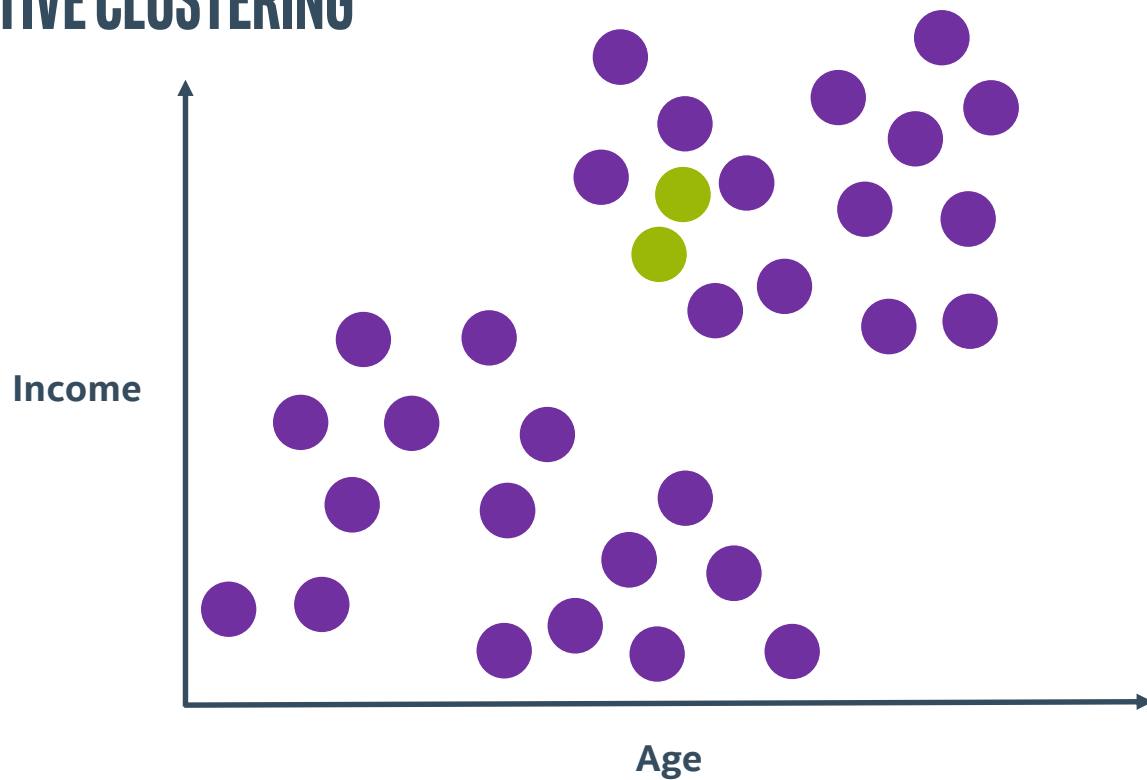
# HIERARCHICAL AGGLOMERATIVE CLUSTERING

# HIERARCHICAL AGGLOMERATIVE CLUSTERING

**Find closest pair, merge into a cluster.**

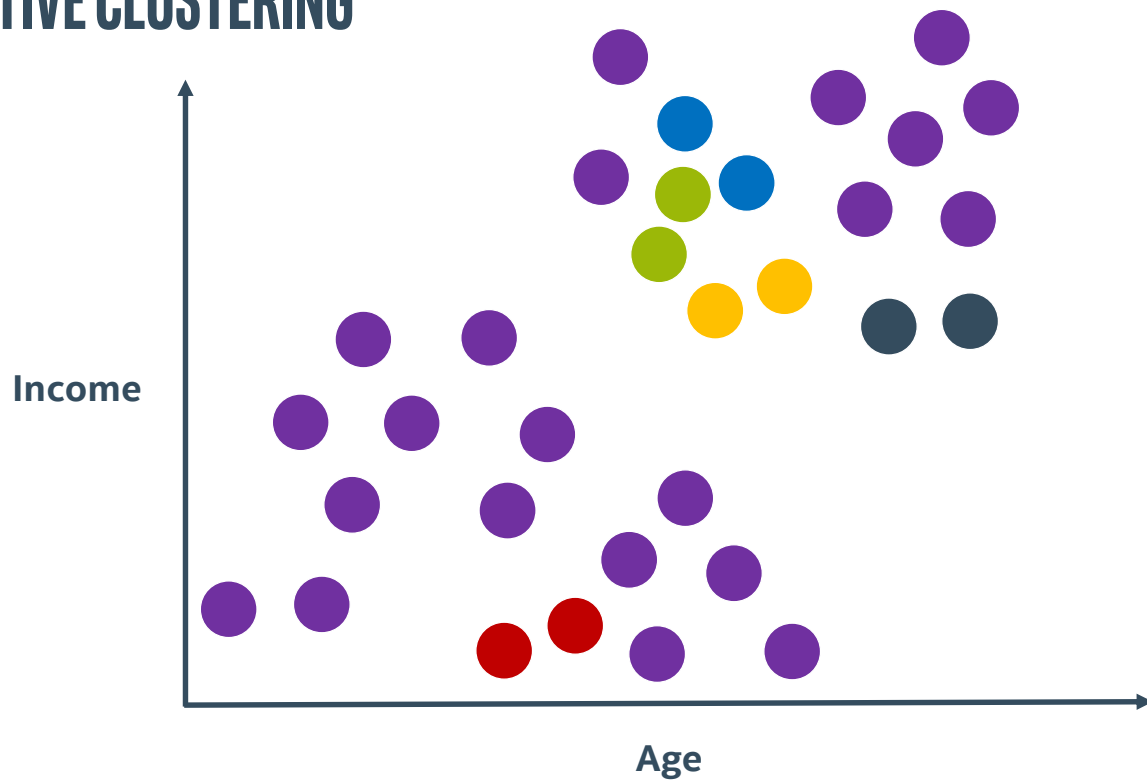# HIERARCHICAL AGGLOMERATIVE CLUSTERING

**Find next closest pair and merge.**

# HIERARCHICAL AGGLOMERATIVE CLUSTERING
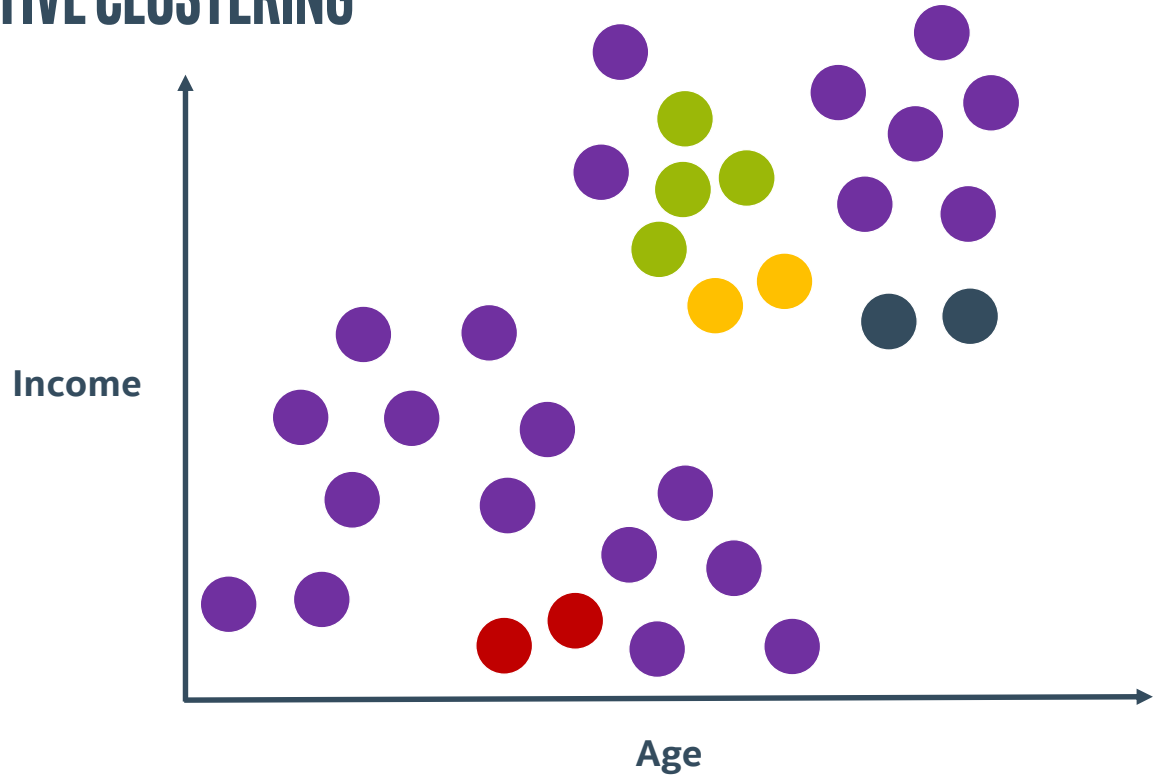
**Find next closest pair and merge.**

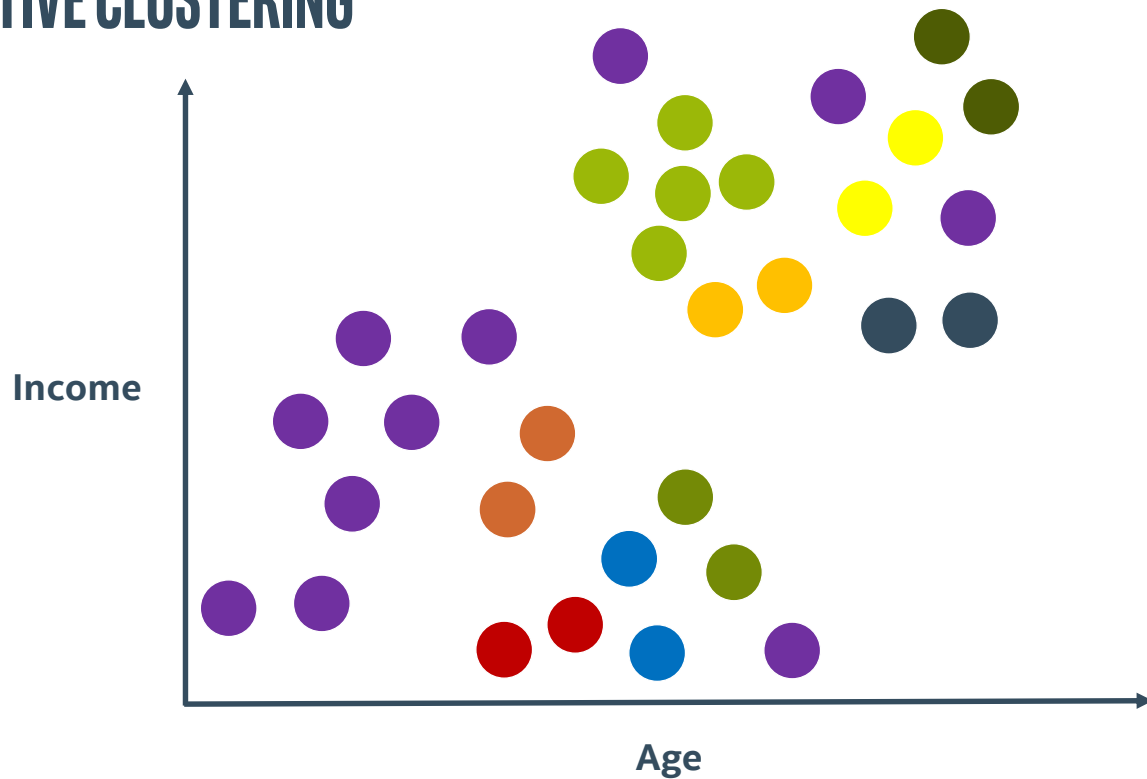# HIERARCHICAL AGGLOMERATIVE CLUSTERING

**Keep merging closest pairs.**

Income

Age

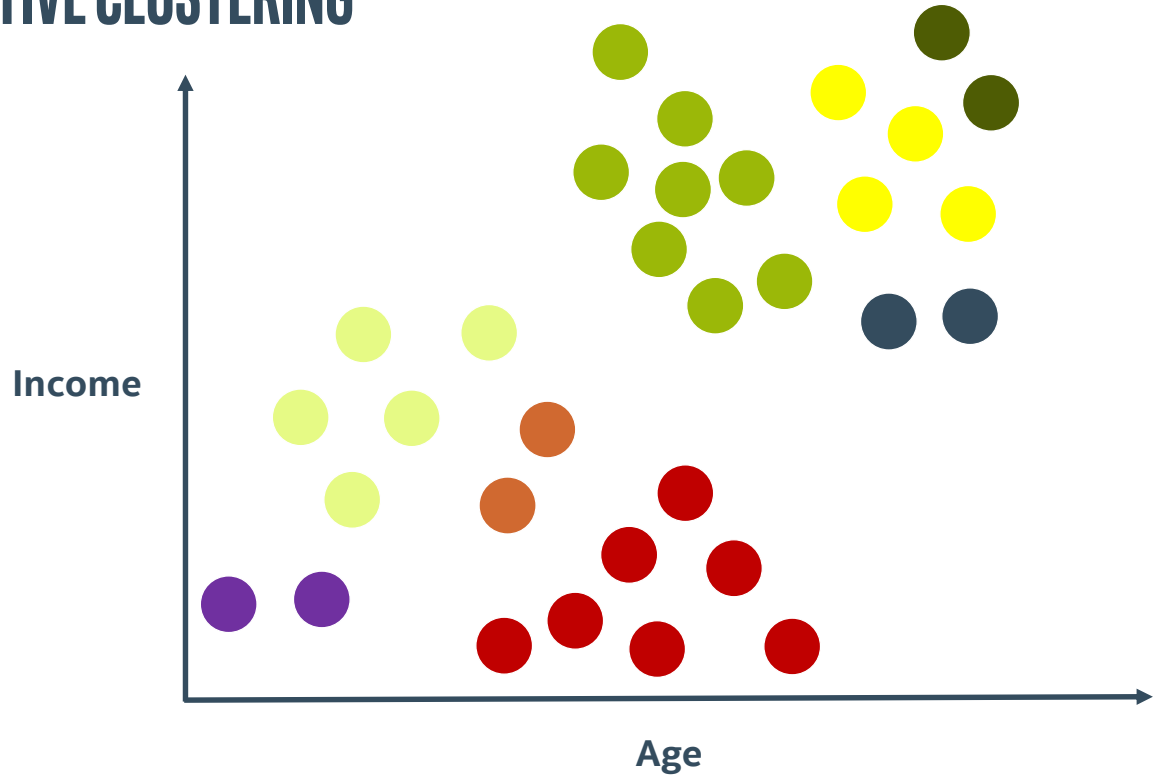# HIERARCHICAL AGGLOMERATIVE CLUSTERING

**If the closest pair is two clusters, merge them.**

# HIERARCHICAL AGGLOMERATIVE CLUSTERING

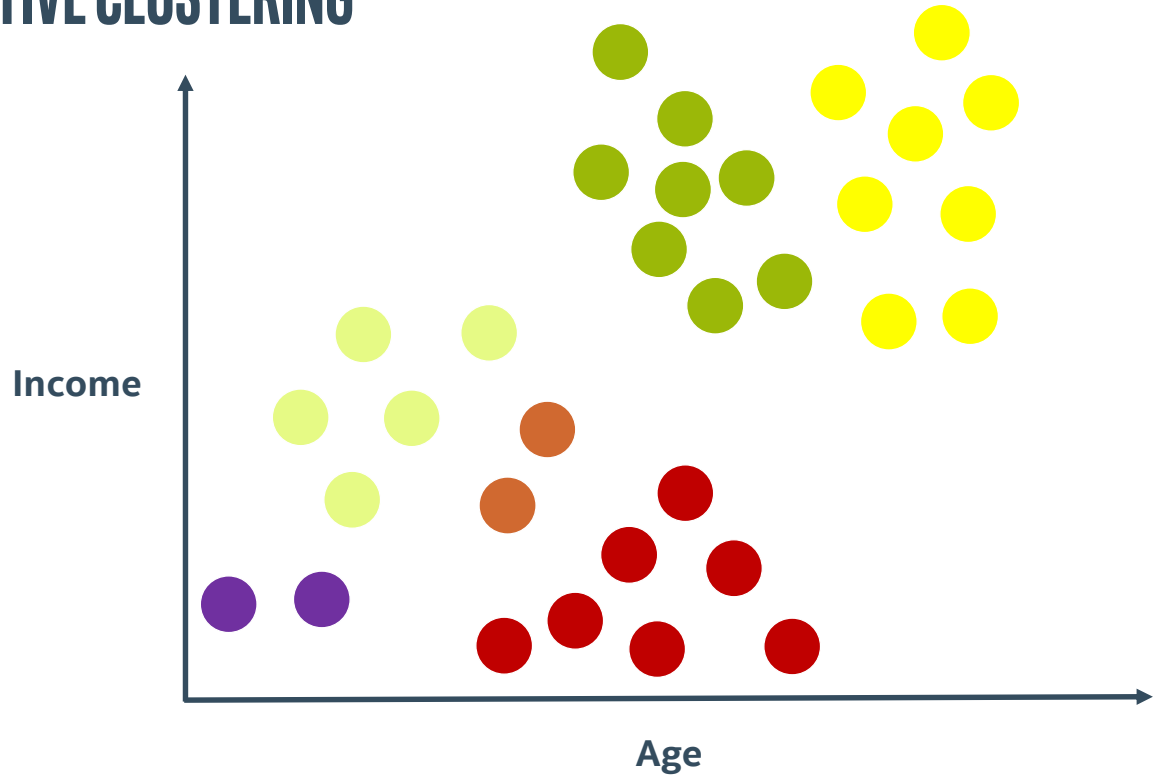**Keep merging closest pairs and clusters.**

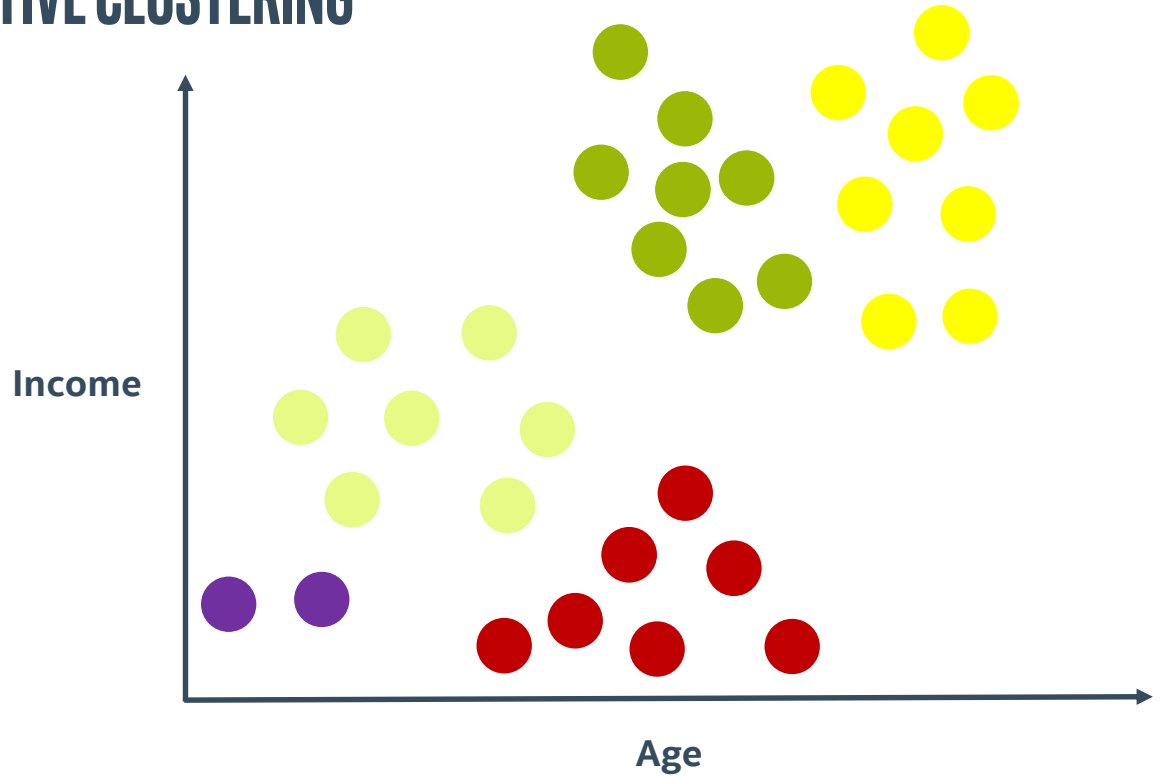# HIERARCHICAL AGGLOMERATIVE CLUSTERING

**Keep merging closest pairs and clusters.**

# HIERARCHICAL AGGLOMERATIVE CLUSTERING
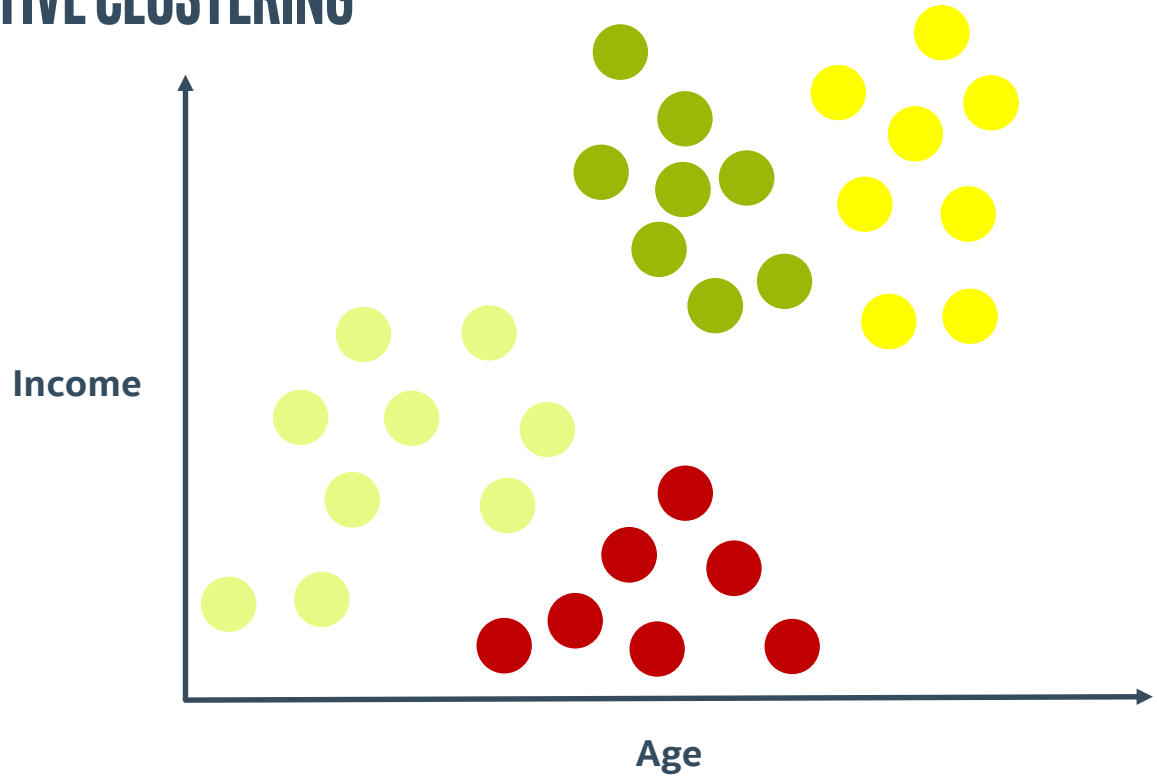
**Current number
of clusters = 6.**

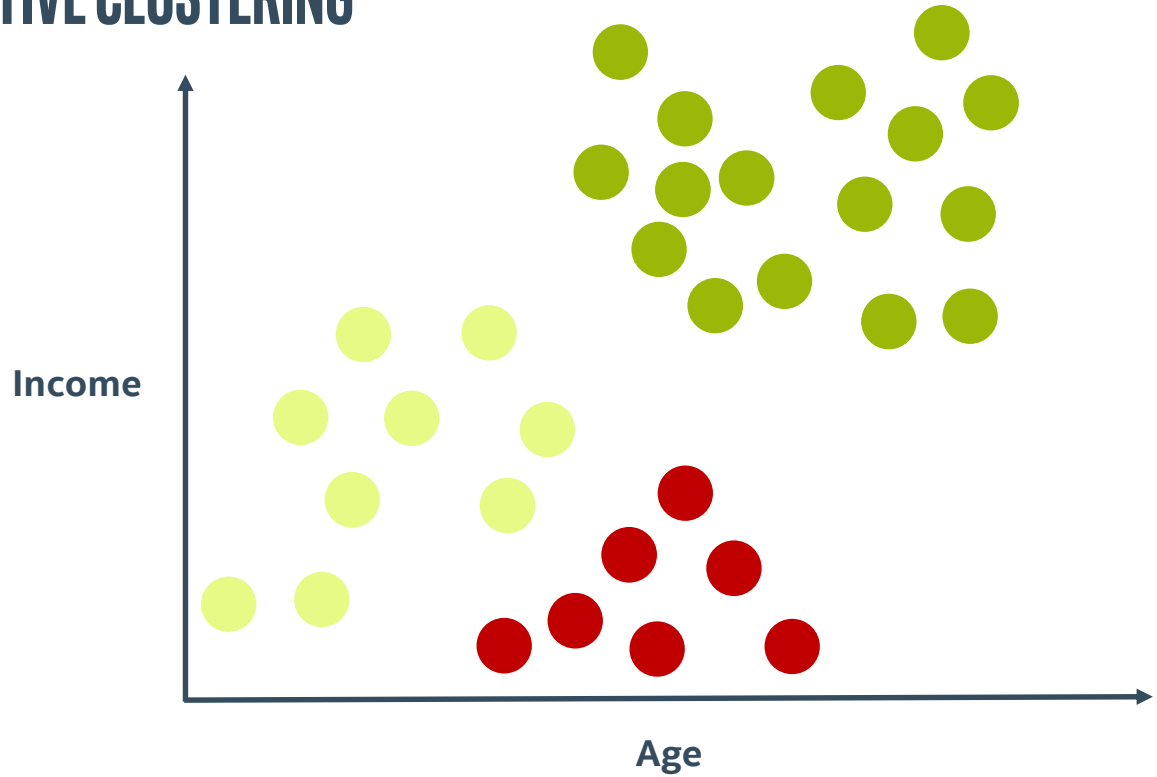# HIERARCHICAL AGGLOMERATIVE CLUSTERING

**Current number
of clusters = 5.**

# HIERARCHICAL AGGLOMERATIVE CLUSTERING
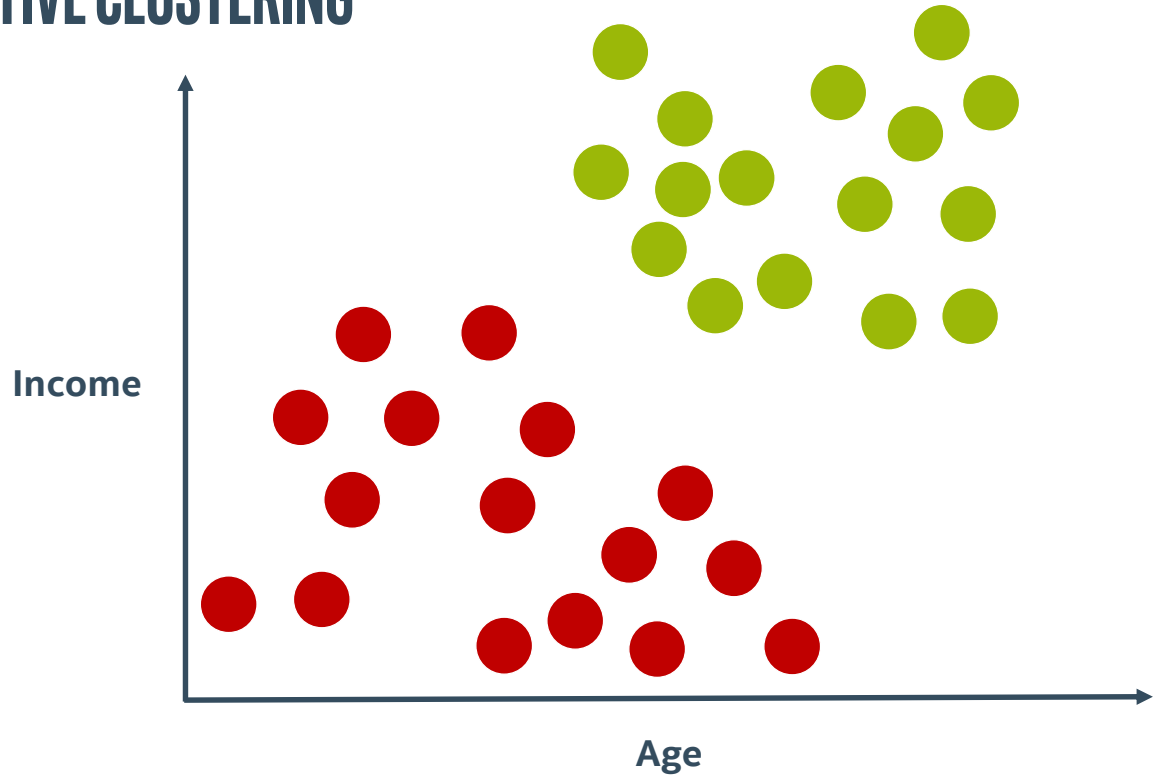
**Current number of clusters = 4.**

# HIERARCHICAL AGGLOMERATIVE CLUSTERING
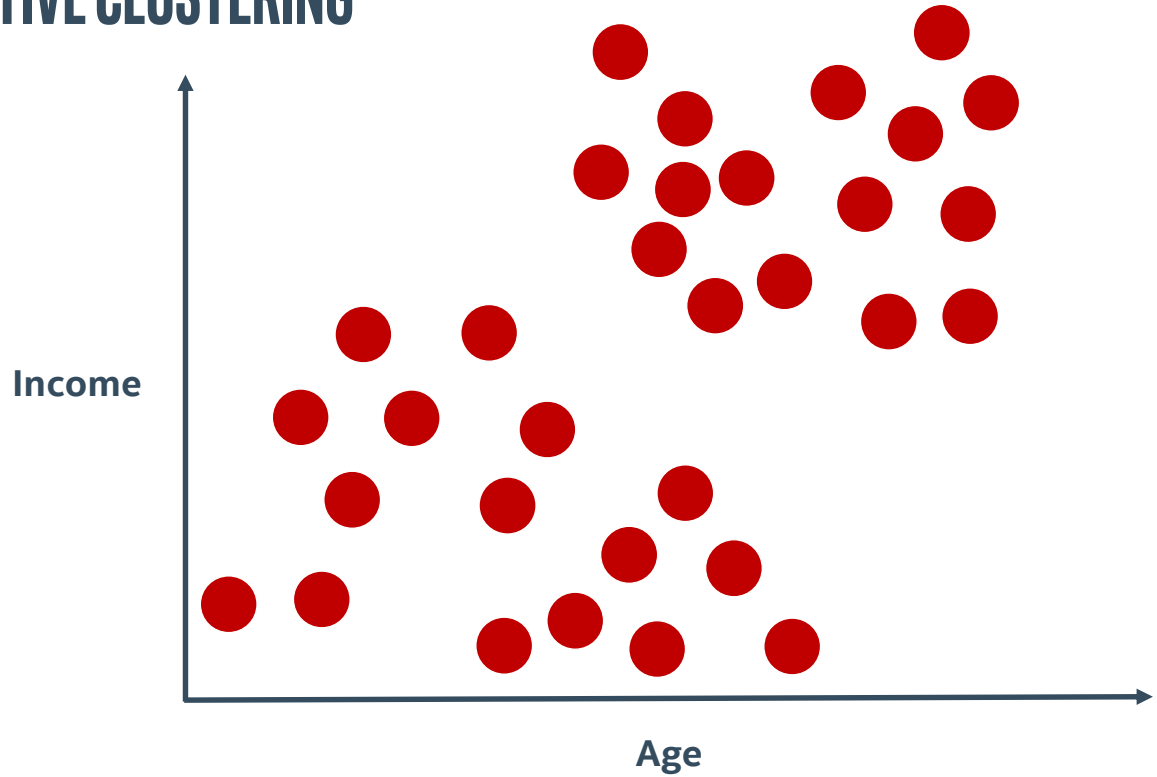
**Current number of clusters = 3.**

# HIERARCHICAL AGGLOMERATIVE CLUSTERING

**Current number of clusters = 2.**

# HIERARCHICAL AGGLOMERATIVE CLUSTERING

**Current number
of clusters = 1.**

# AGGLOMERATIVE CLUSTERING STOPPING CONDITIONS

**CONDITION 1**
The correct number of clusters is reached

# AGGLOMERATIVE CLUSTERING STOPPING CONDITIONS
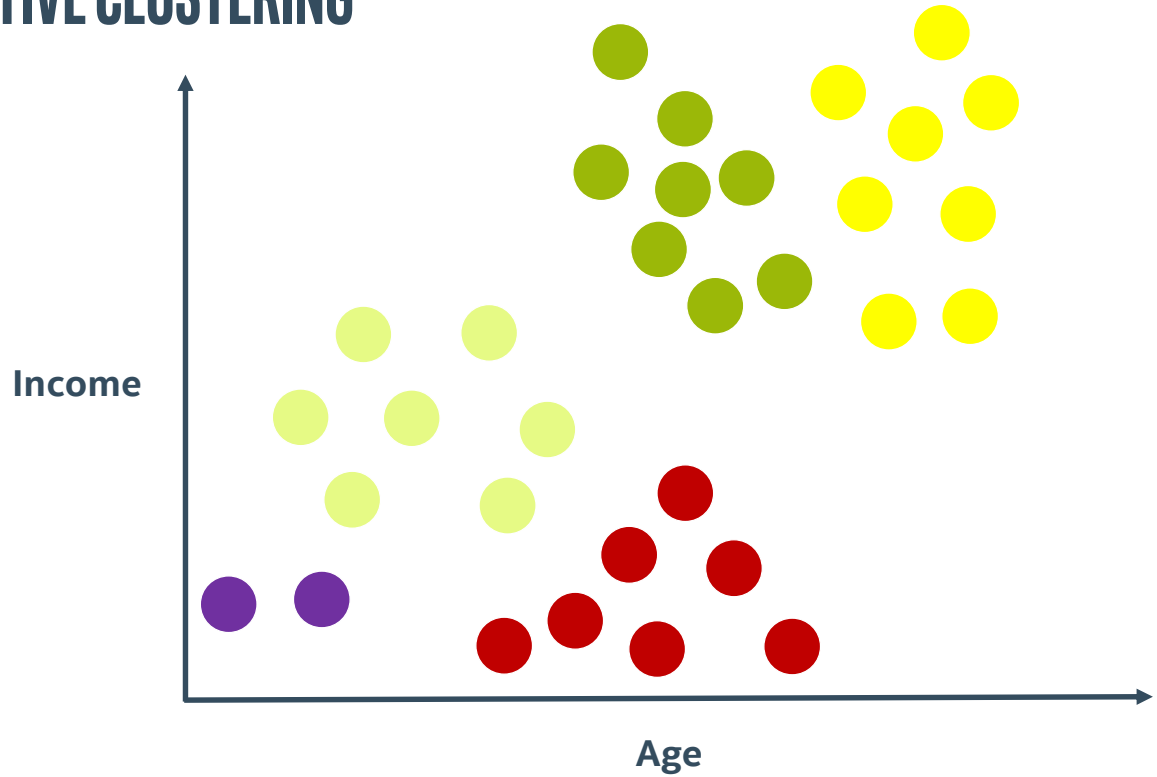
**CONDITION 1**     The correct number of clusters is reached

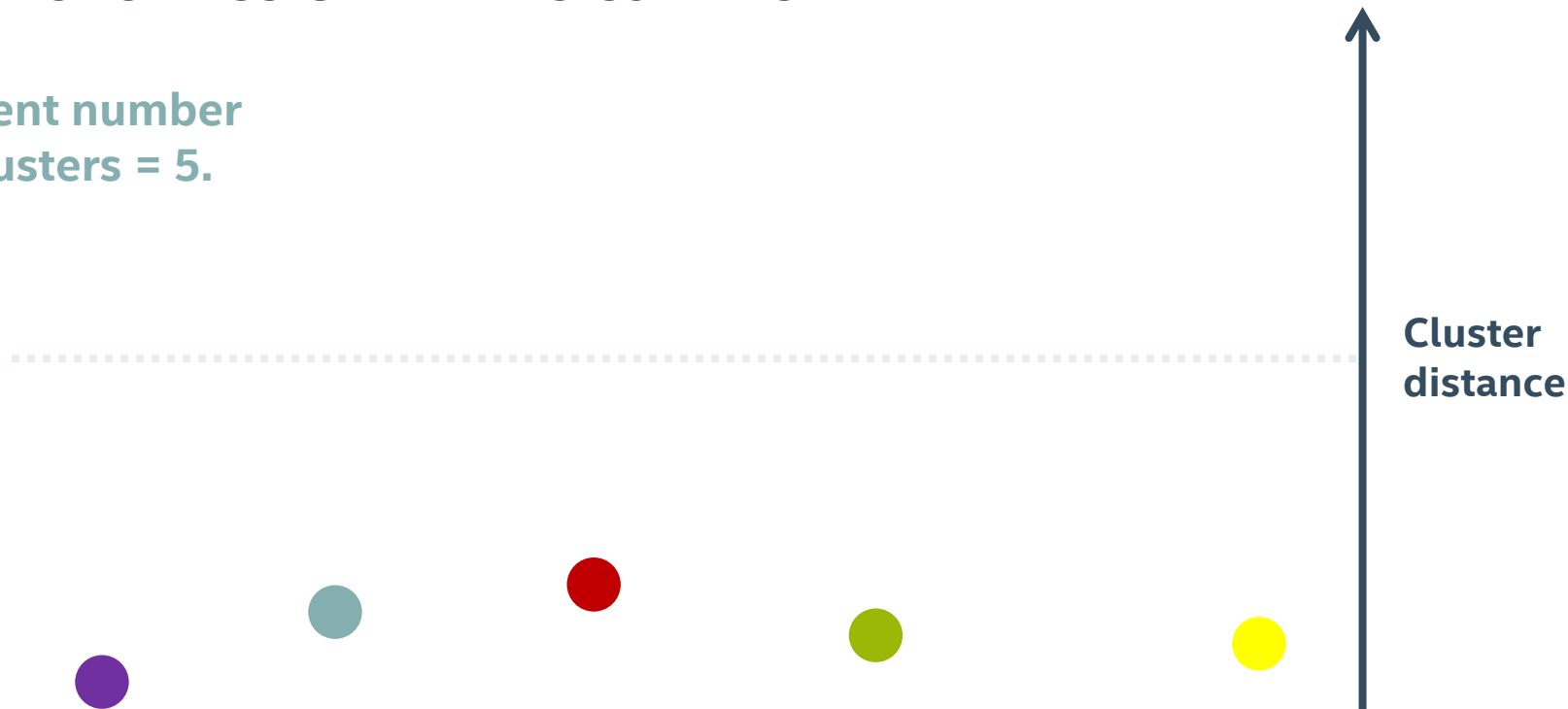**CONDITION 2**     Minimum average cluster distance reaches a set value

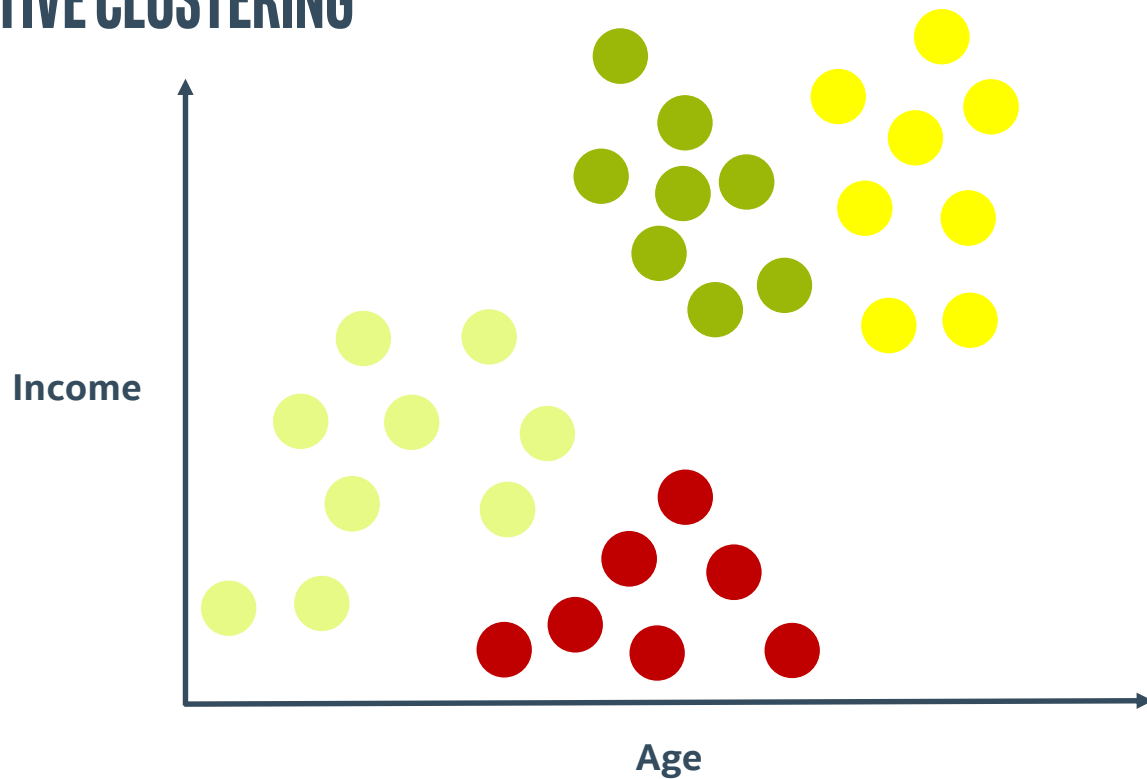# HIERARCHICAL AGGLOMERATIVE CLUSTERING

**Current number of clusters = 5.**

# HIERARCHICAL AGGLOMERATIVE CLUSTERING

**Current number
of clusters = 5.**

**Cluster
distance**

# HIERARCHICAL AGGLOMERATIVE CLUSTERING

**Current number
of clusters = 4.**

# HIERARCHICAL AGGLOMERATIVE CLUSTERING

**Current number
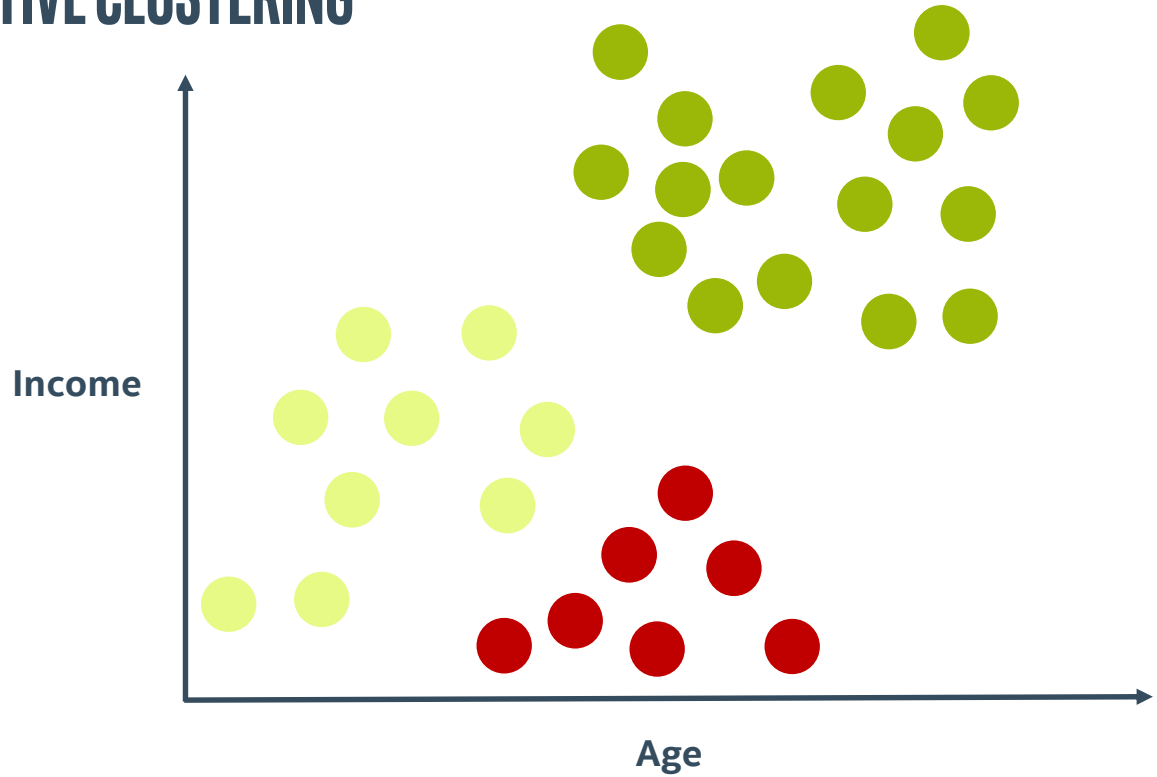of clusters = 4.**

**Cluster
distance**

# HIERARCHICAL AGGLOMERATIVE CLUSTERING

**Current number
of clusters = 3.**

# HIERARCHICAL AGGLOMERATIVE CLUSTERING

**Current number of clusters = 3.**
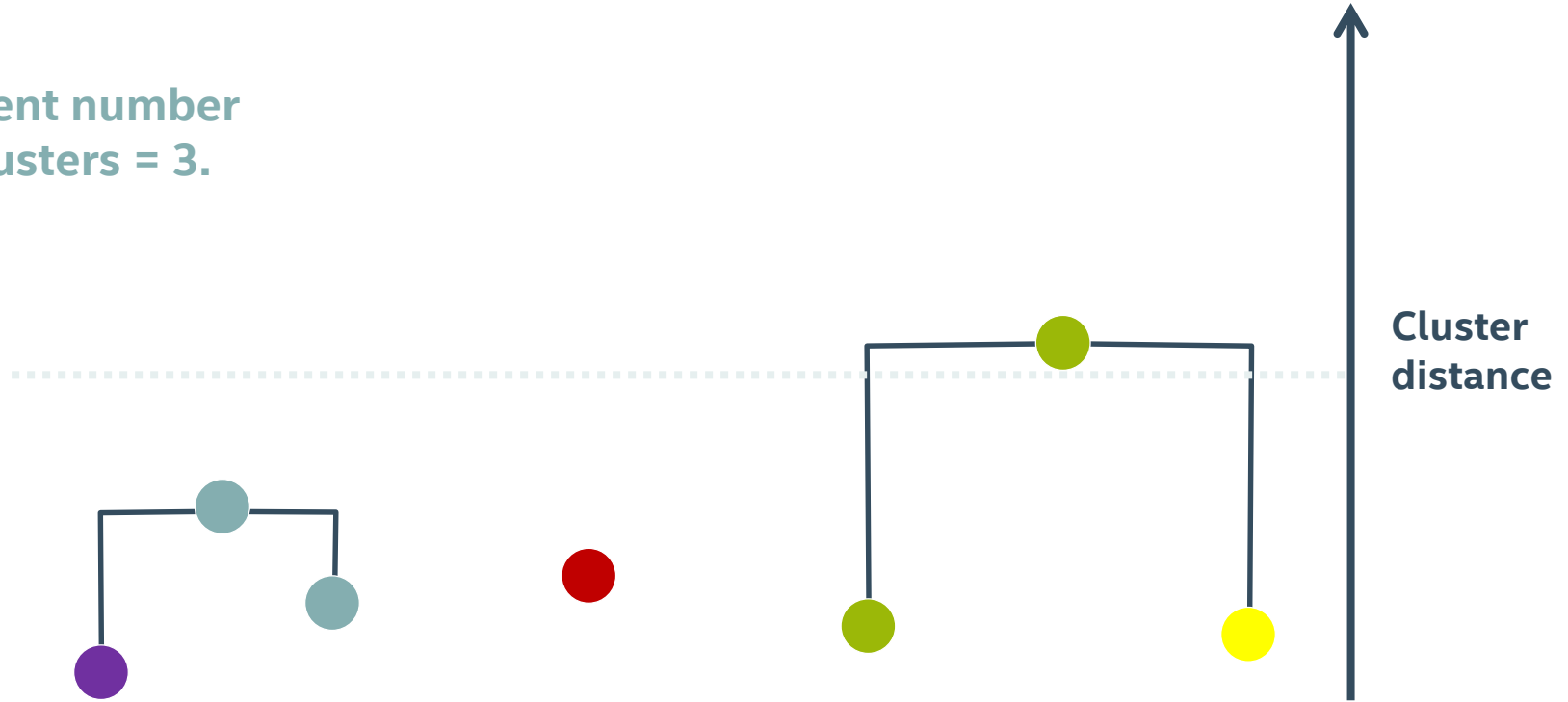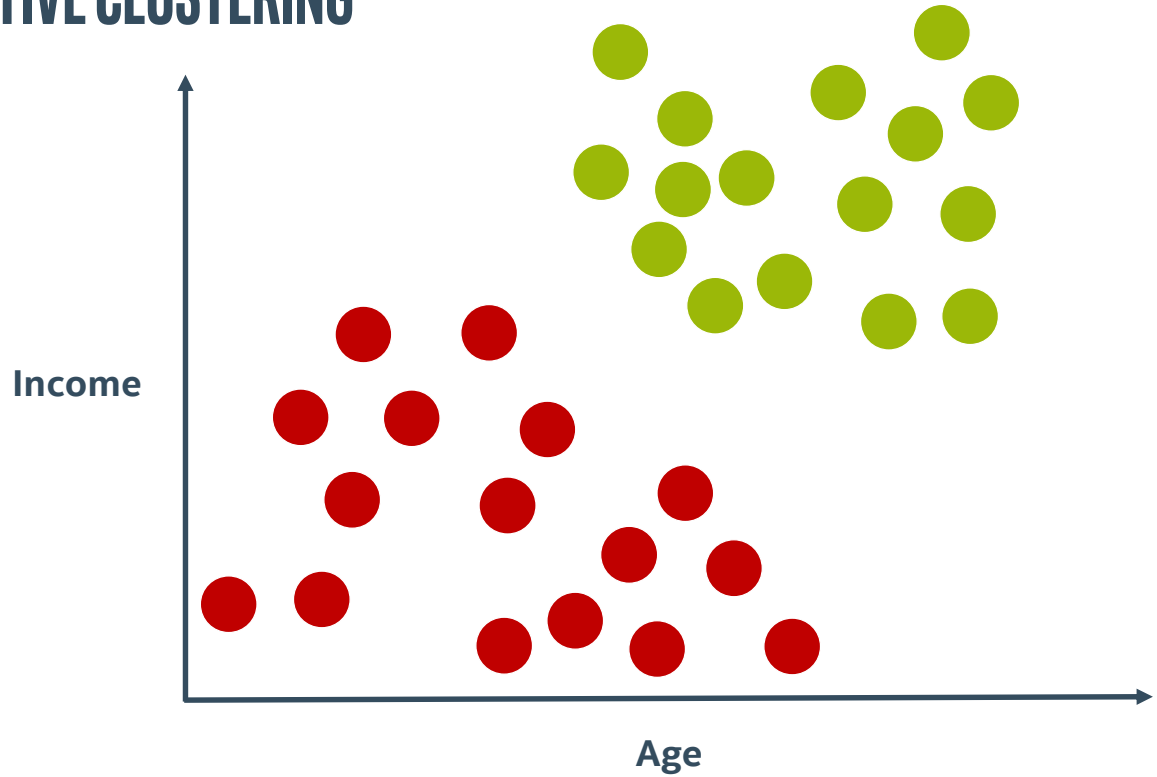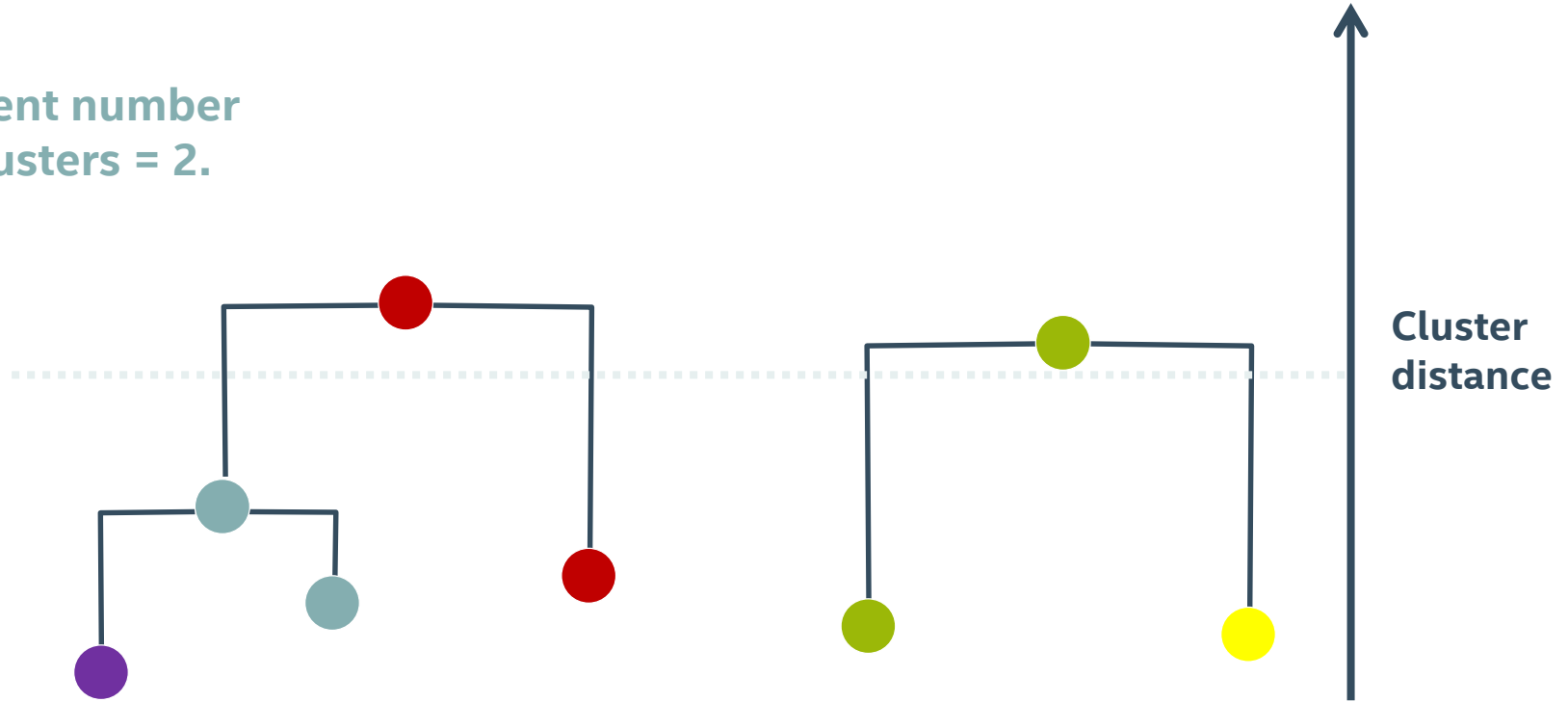
**Cluster distance**

# HIERARCHICAL AGGLOMERATIVE CLUSTERING

**Current number
of clusters = 2.**

# HIERARCHICAL AGGLOMERATIVE CLUSTERING
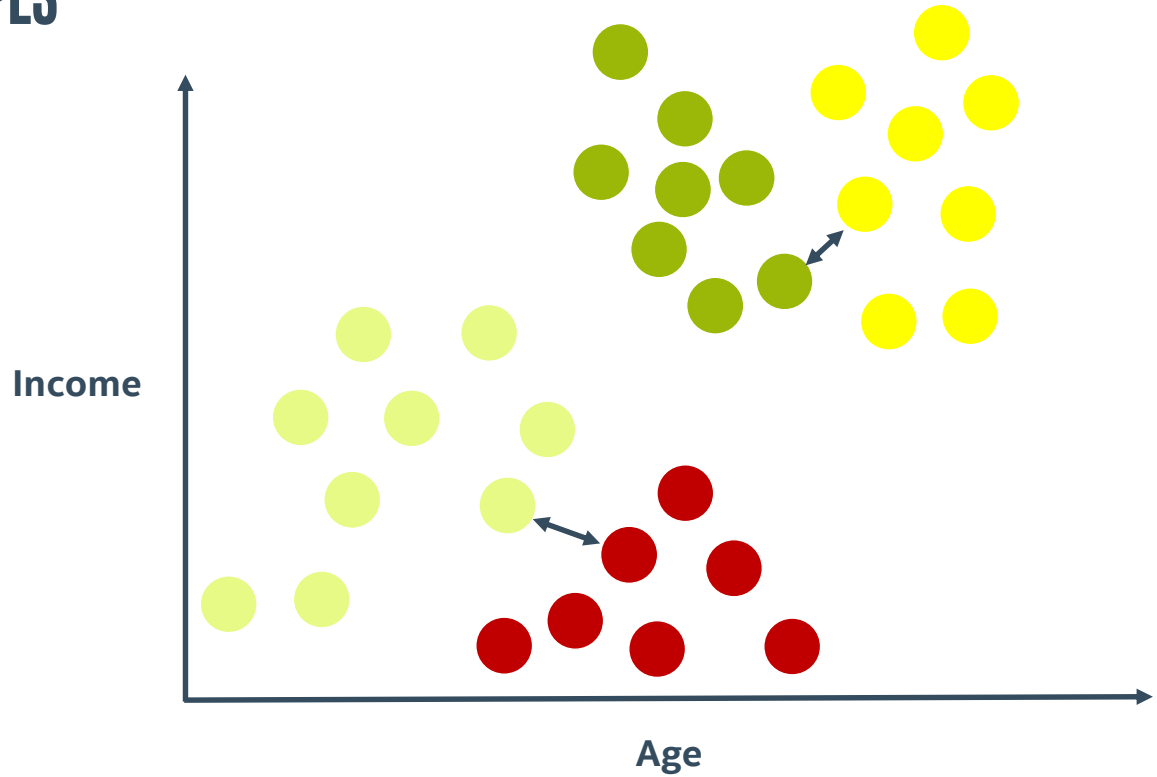
**Current number
of clusters = 2.**
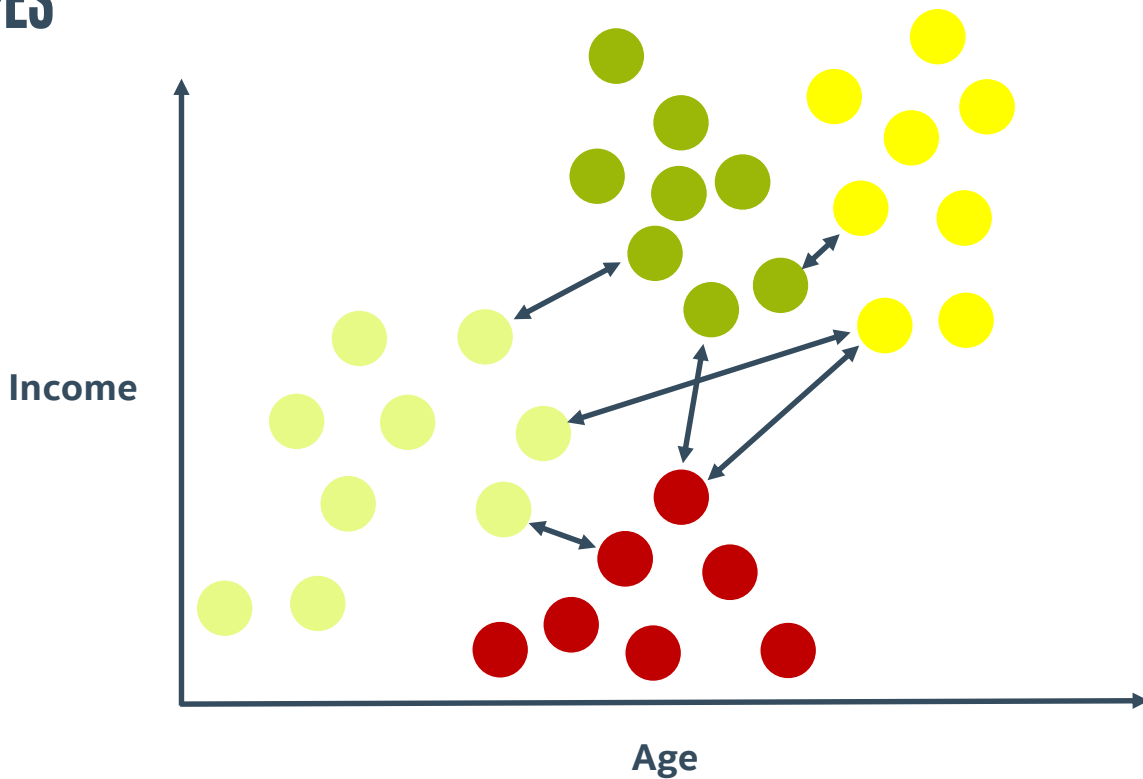
**Cluster
distance**

# HIERARCHICAL LINKAGE TYPES

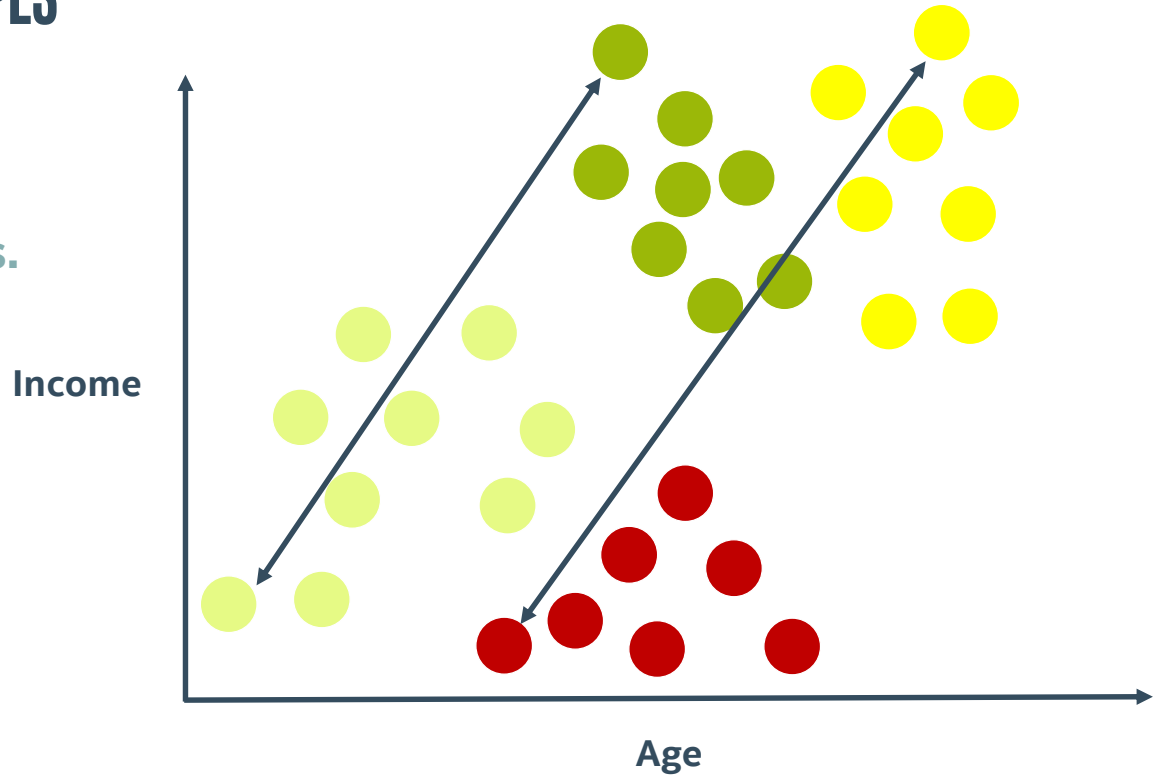**Single linkage: minimum pairwise distance between clusters.**

# HIERARCHICAL LINKAGE TYPES

**Single linkage: minimum pairwise distance between clusters.**

# HIERARCHICAL LINKAGE TYPES

**Complete linkage: maximum pairwise distance between clusters.**
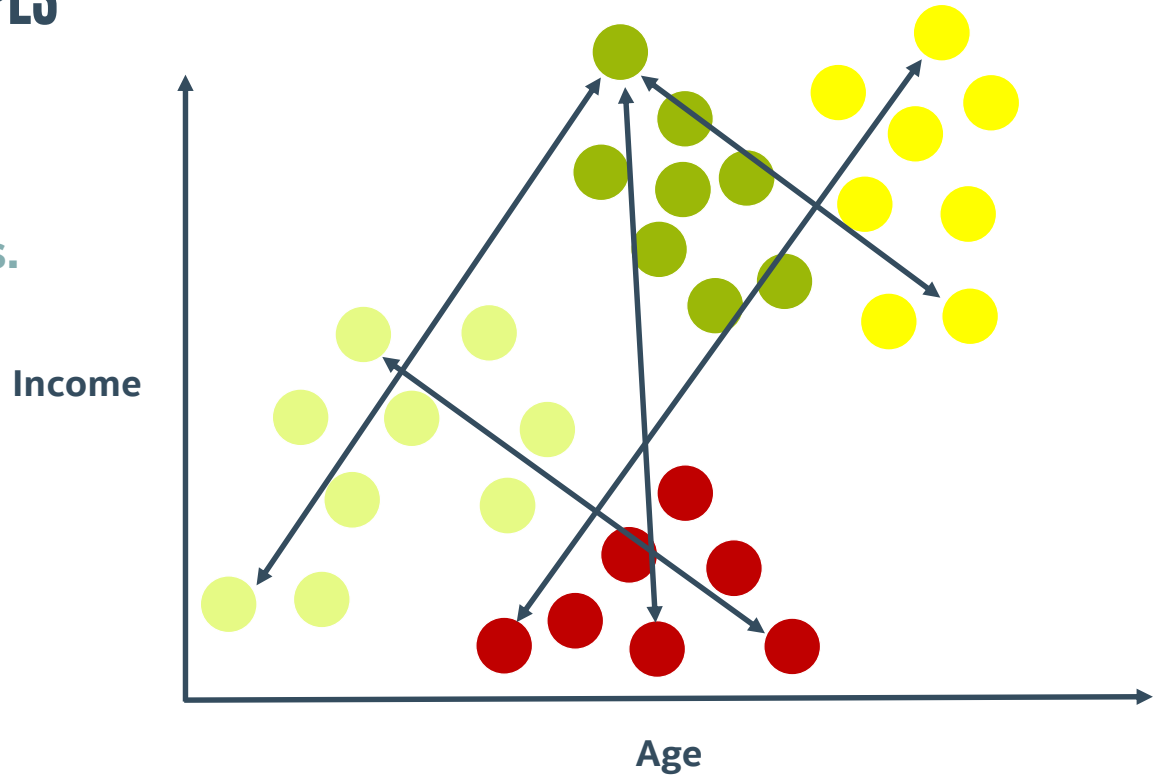
Income

Age
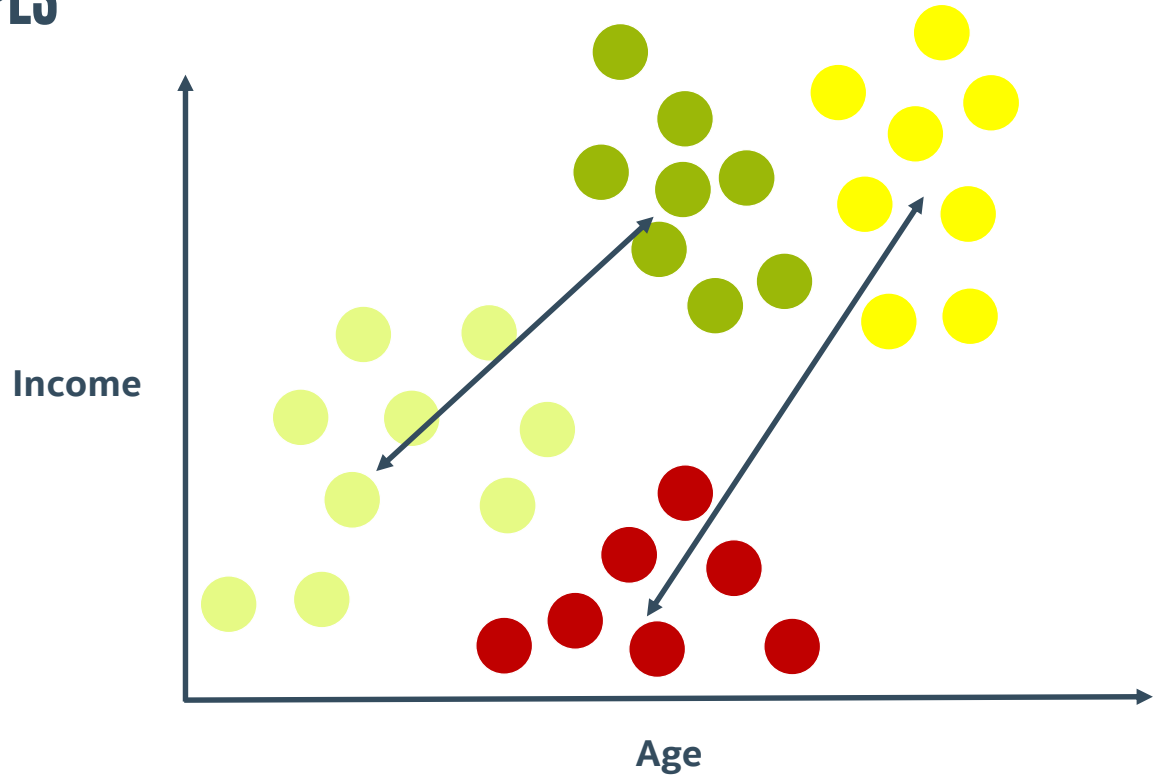
# HIERARCHICAL LINKAGE TYPES

**Complete linkage:
maximum pairwise
distance between clusters.**

# HIERARCHICAL LINKAGE TYPES

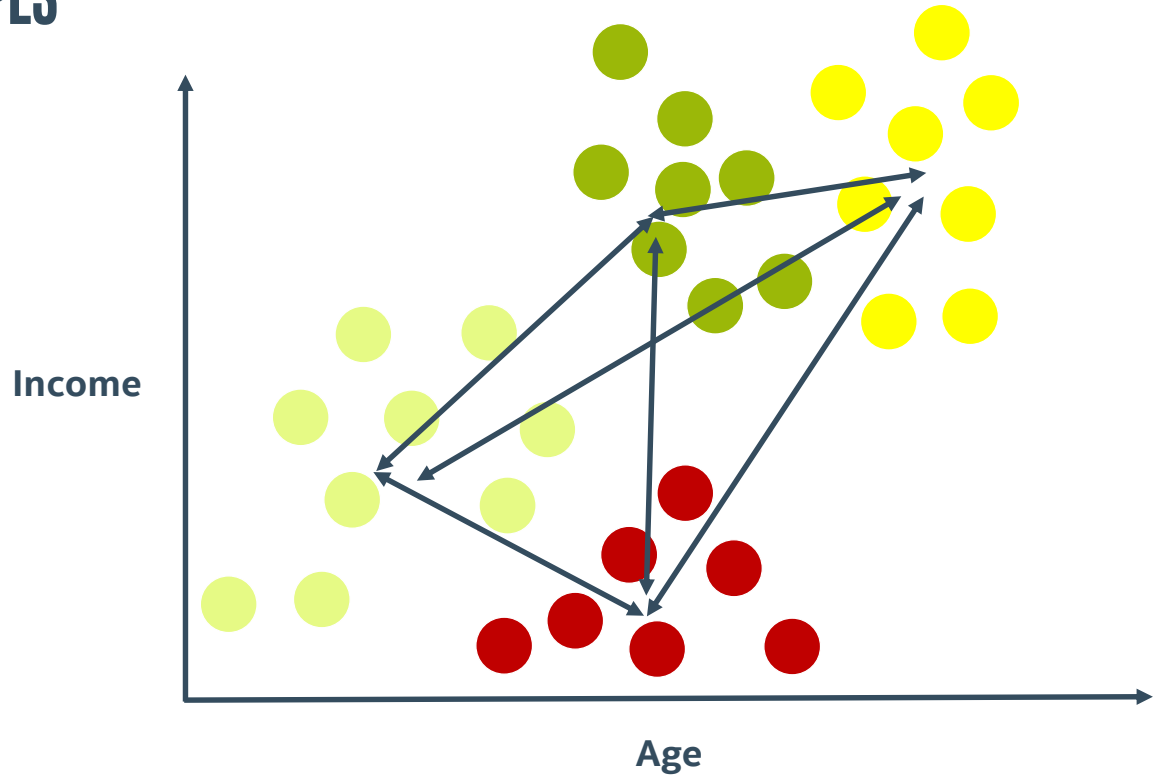**Average linkage: average pairwise distance between clusters.**
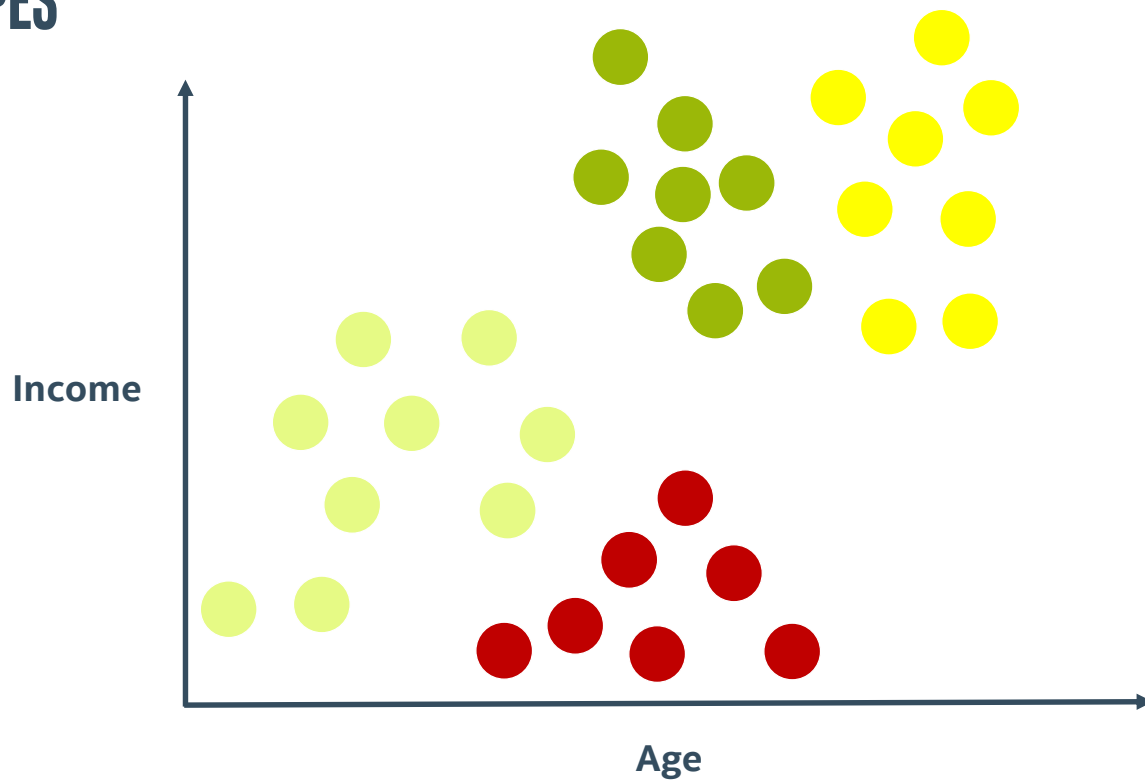


Income

Age

# HIERARCHICAL LINKAGE TYPES
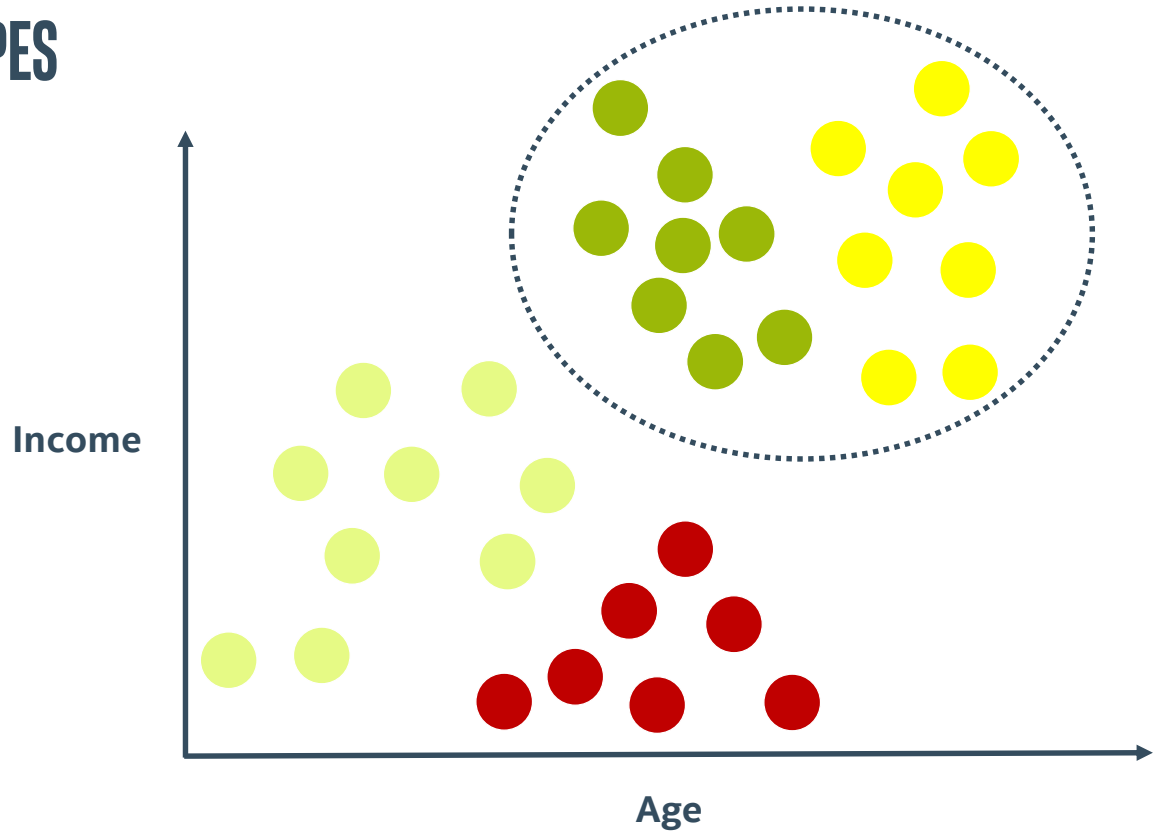
**Average linkage: average pairwise distance between clusters.**

# HIERARCHICAL LINKAGE TYPES

**Ward linkage: merge based on best inertia.**



Income

Age

# HIERARCHICAL LINKAGE TYPES

**Ward linkage: merge based on best inertia.**

# AGGLOMERATIVE CLUSTERING: THE SYNTAX

Import the class containing the clustering method.

```
from sklearn.cluster import AgglomerativeClustering
```

Create an instance of the class.

```
agg = AgglomerativeClustering(n_clusters=3,
                              affinity='euclidean',
                              linkage='ward')
```

Fit the instance on the data and then predict clusters for new data.

```
agg = agg.fit(X1)

y_predict = agg.predict(X2)
```

# AGGLOMERATIVE CLUSTERING: THE SYNTAX

**Import the class containing the clustering method.**

```
from sklearn.cluster import AgglomerativeClustering
```

**Create an instance of the class.**

```
agg = AgglomerativeClustering(n_clusters=3,

                    affinity='euclidean',

                    linkage='ward')
```

**final number of clusters**

**Fit the instance on the data and then predict clusters for new data.**

```
agg = agg.fit(X1)

y_predict = agg.predict(X2)
```

# AGGLOMERATIVE CLUSTERING: THE SYNTAX

**Import the class containing the clustering method.**

```
from sklearn.cluster import AgglomerativeClustering
```

**Create an instance of the class.**

```
agg = AgglomerativeClustering(n_clusters=3,

                    affinity='euclidean',

                    linkage='ward')
```
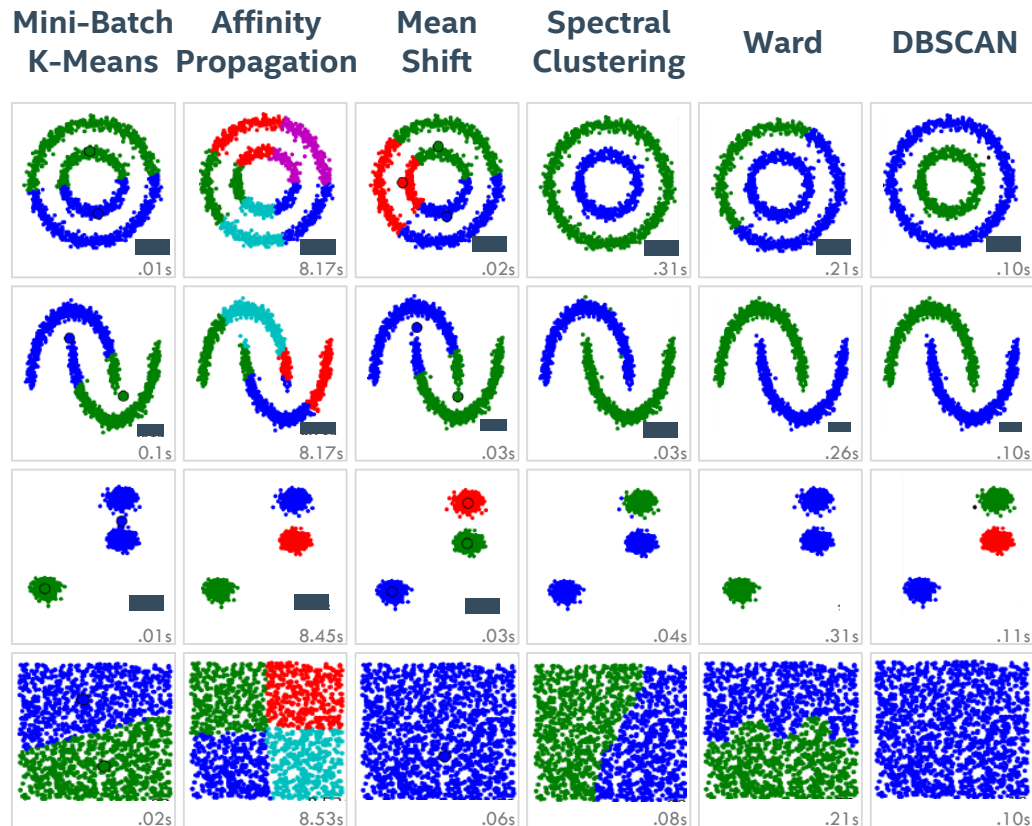
← **cluster affinity and aggregation**

**Fit the instance on the data and then predict clusters for new data.**

```
agg = agg.fit(X1)

y_predict = agg.predict(X2)
```

# OTHER TYPES OF CLUSTERING



Reference: http://scikit-learn.org/stable/auto_examples/cluster/plot_cluster_comparison.html