

# MATLAB EXPO

November 13–14, 2024 | Online

---

## The CLASSIX Story

### Developing the same algorithm in MATLAB and Python simultaneously

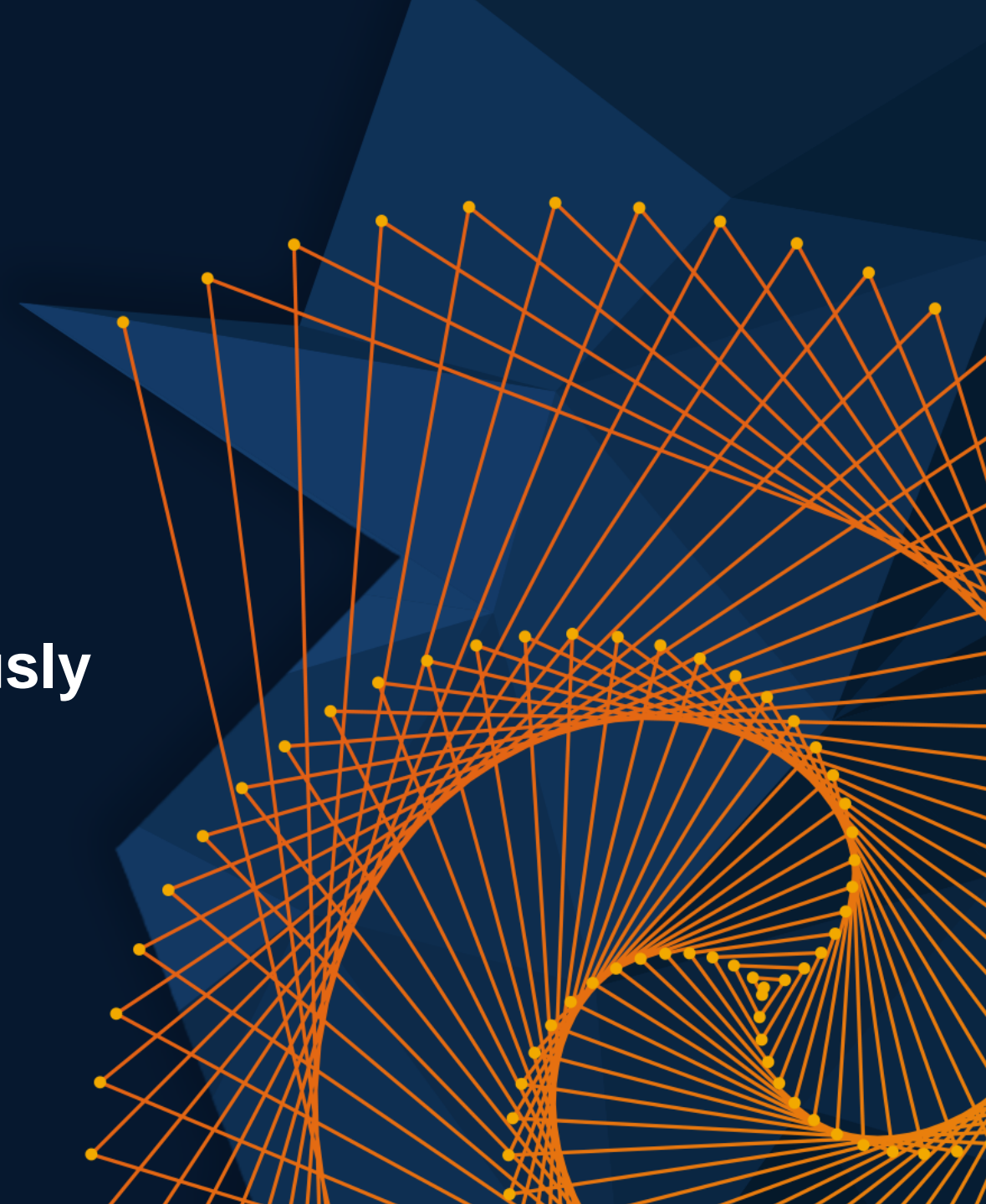
*Stefan Güttel*

*University of Manchester*



*Mike Croucher*

*MathWorks*

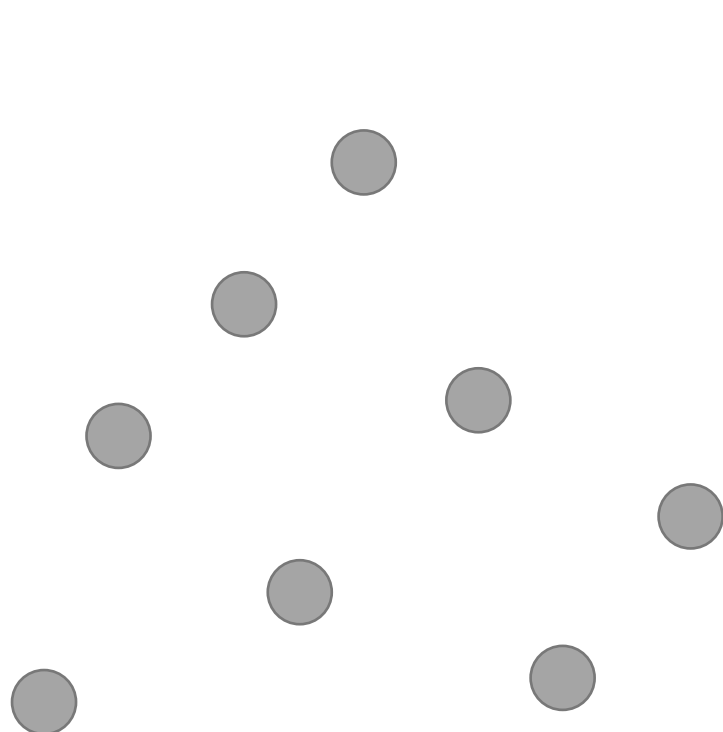


A data clustering method that is

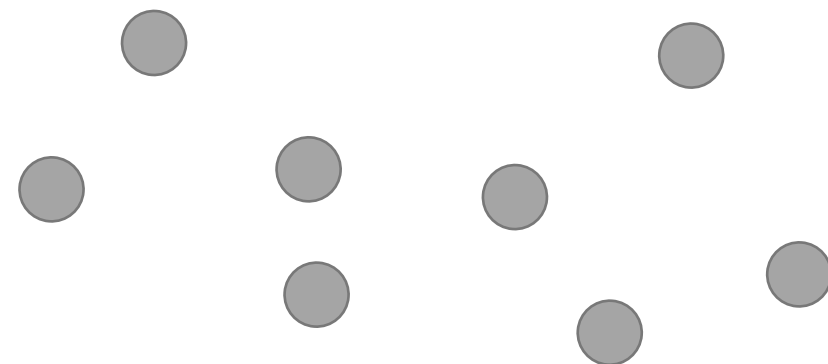
- fast
- memory efficient
- conceptually very simple
- non-iterative and 100% deterministic
- easy to tune with two hyperparameters
- computing fully explainable clustering results

Currently two implementations

- Python <https://github.com/nla-group/classix>
- MATLAB <https://github.com/nla-group/classix-matlab>



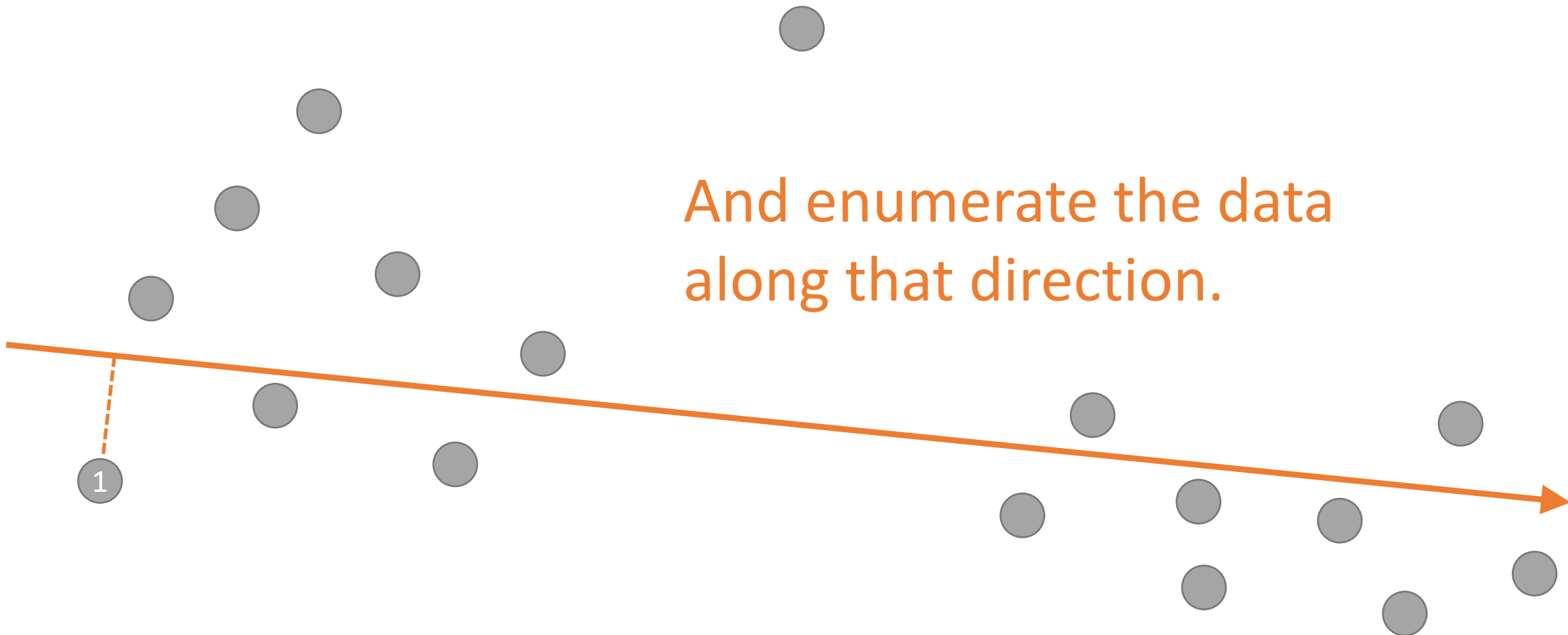
Here are some data points we want to cluster.

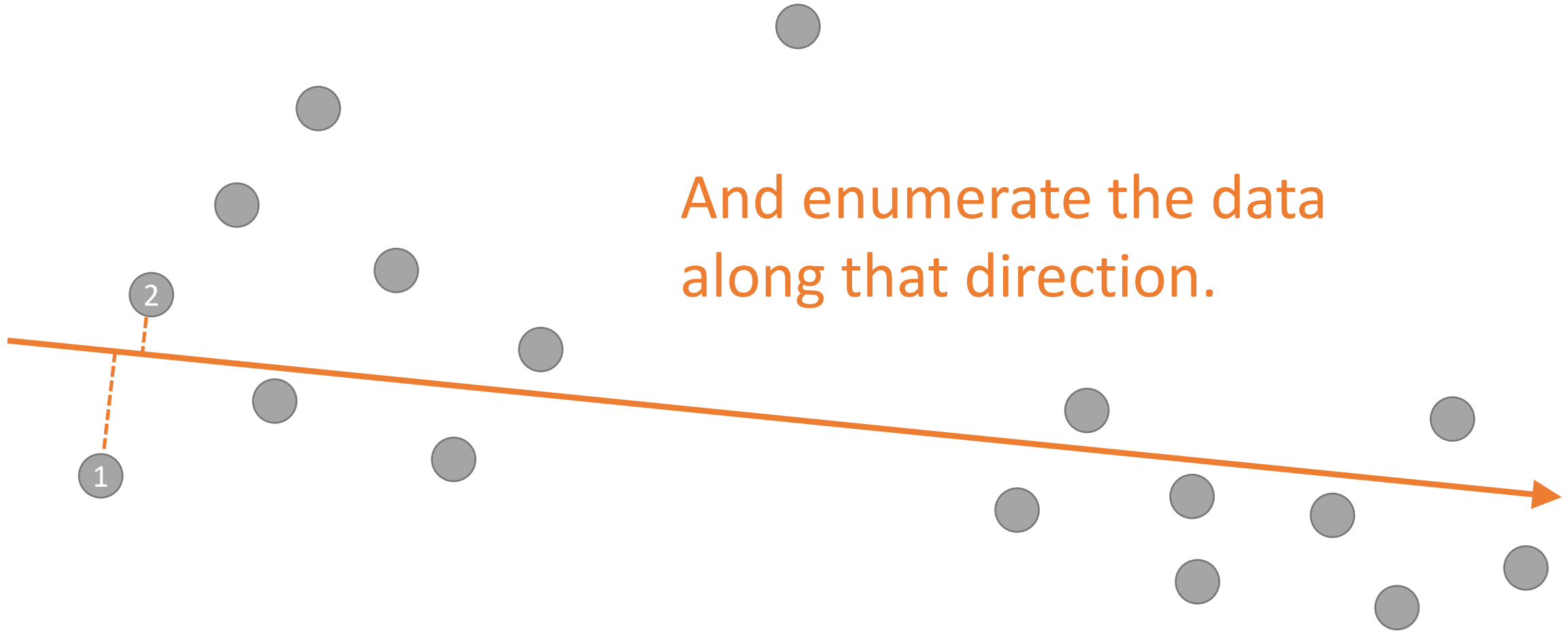


We compute the first principal component.



And enumerate the data  
along that direction.



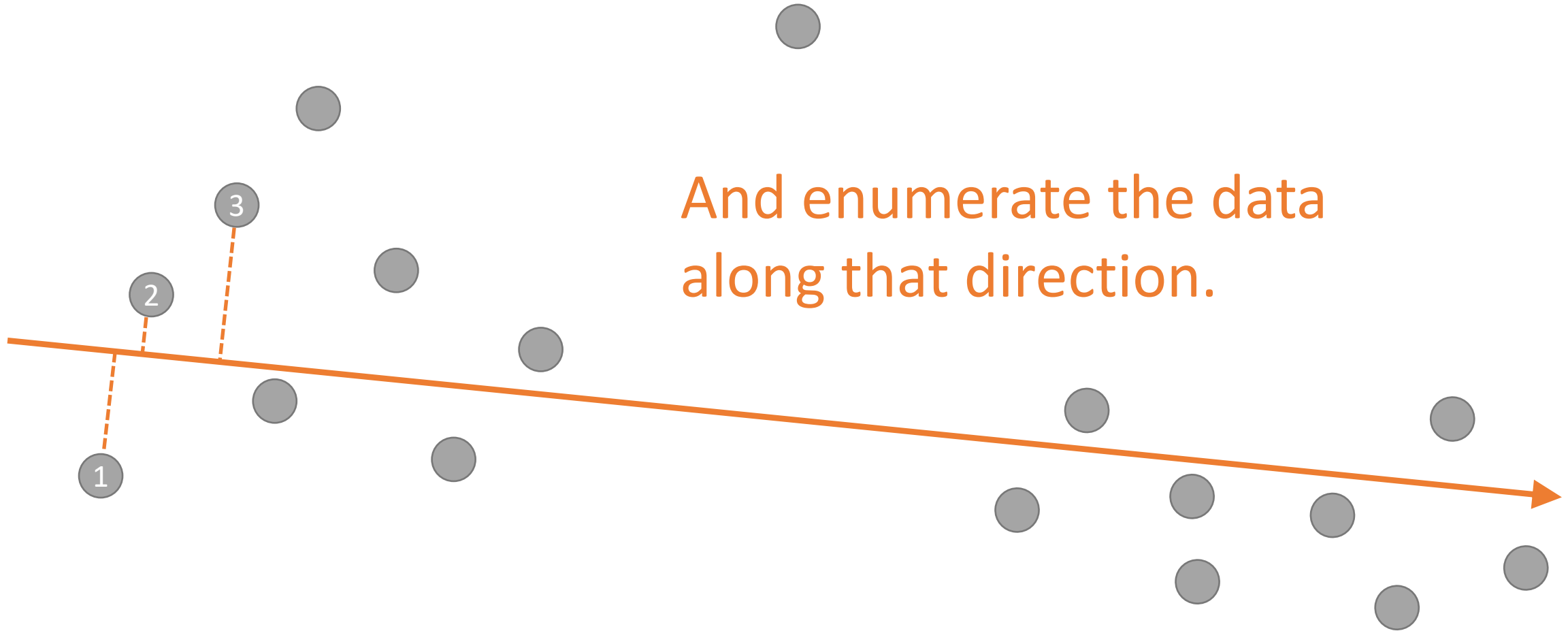


And enumerate the data along that direction.



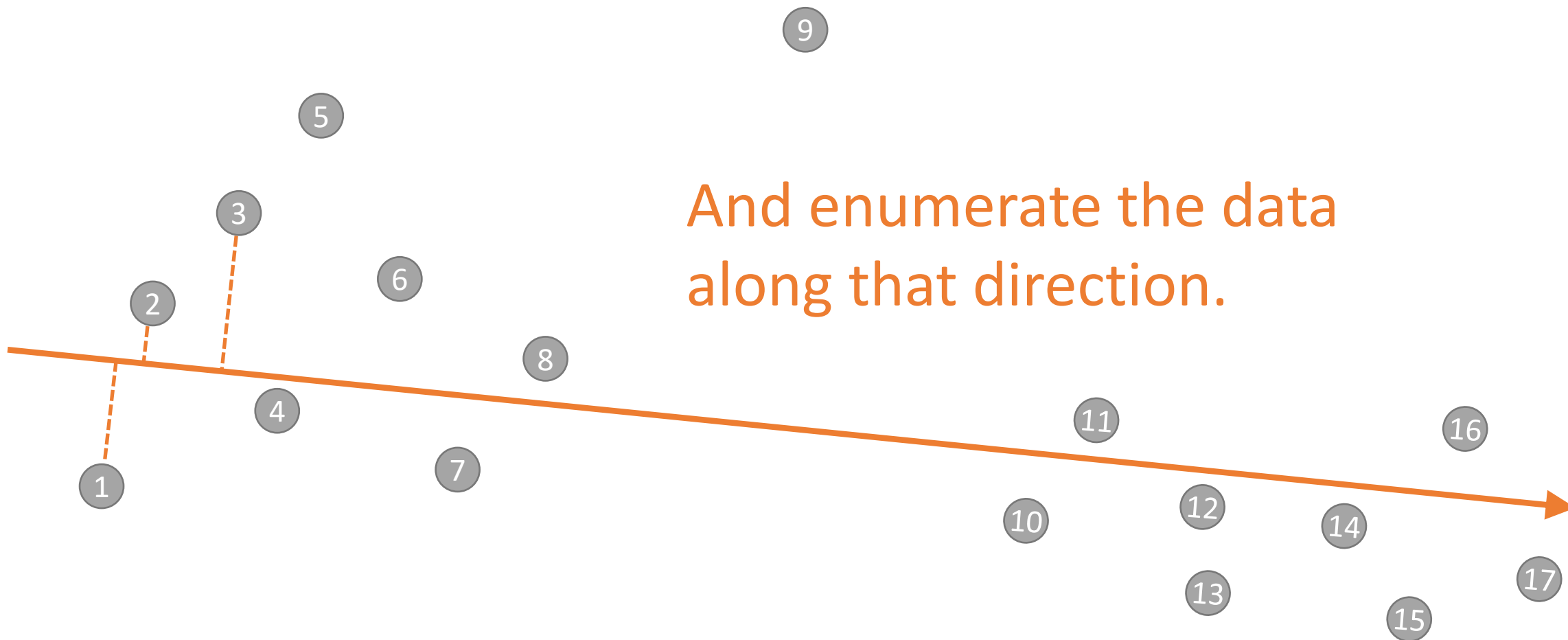
# CLASSIX

Fast and Explainable Clustering



And enumerate the data along that direction.





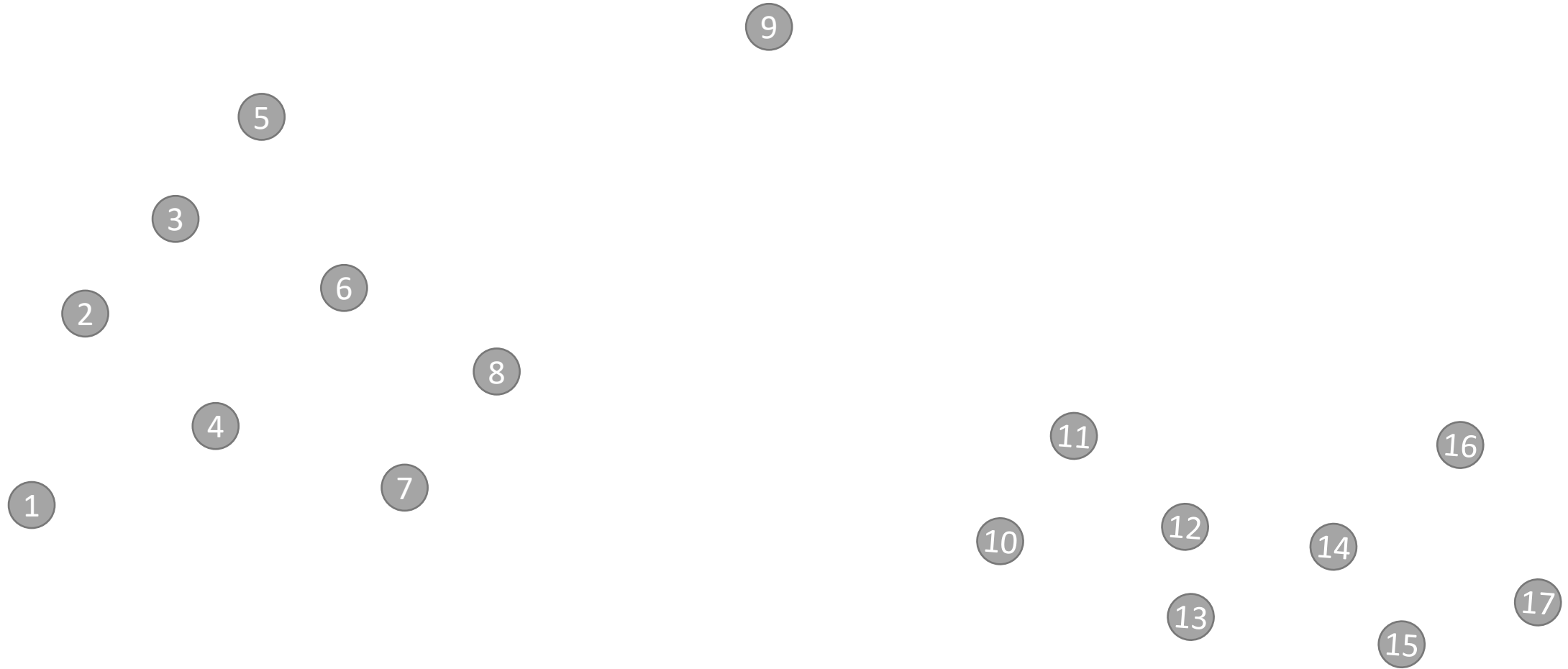
And enumerate the data along that direction.

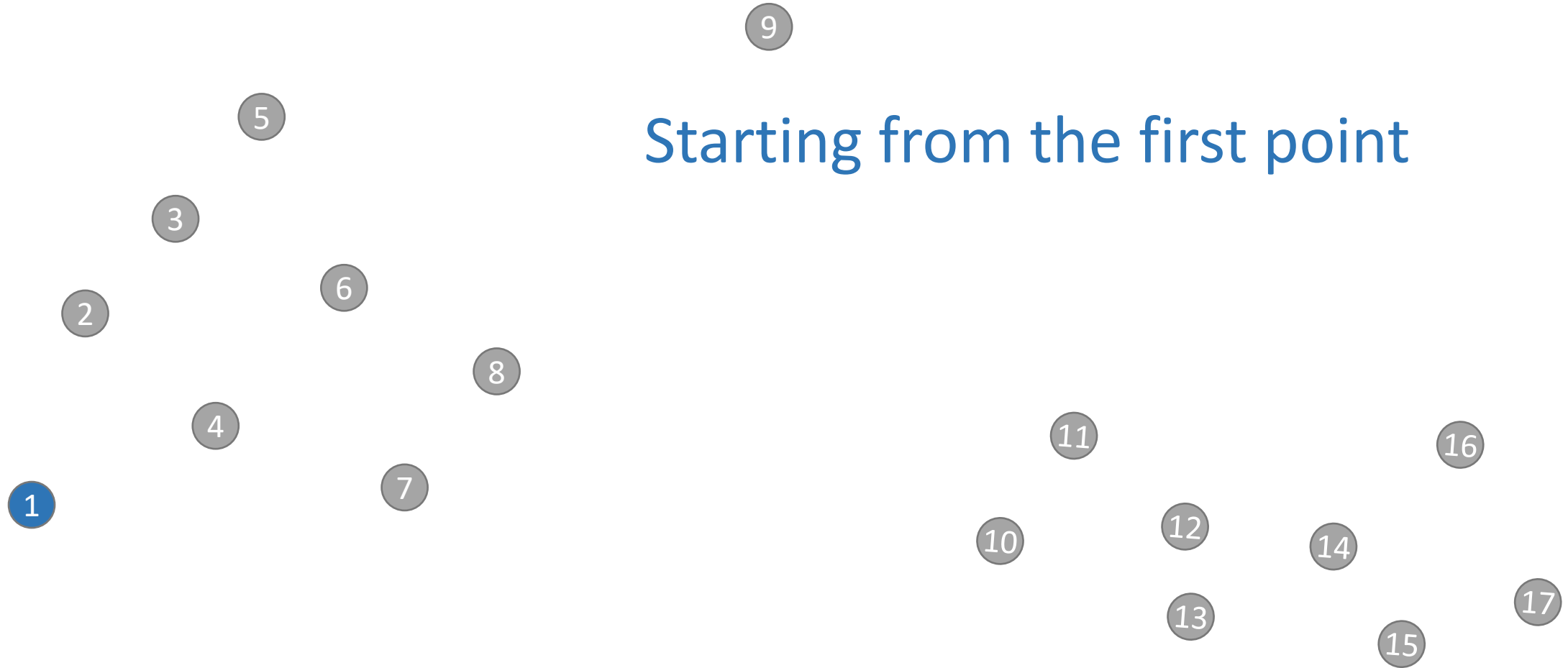


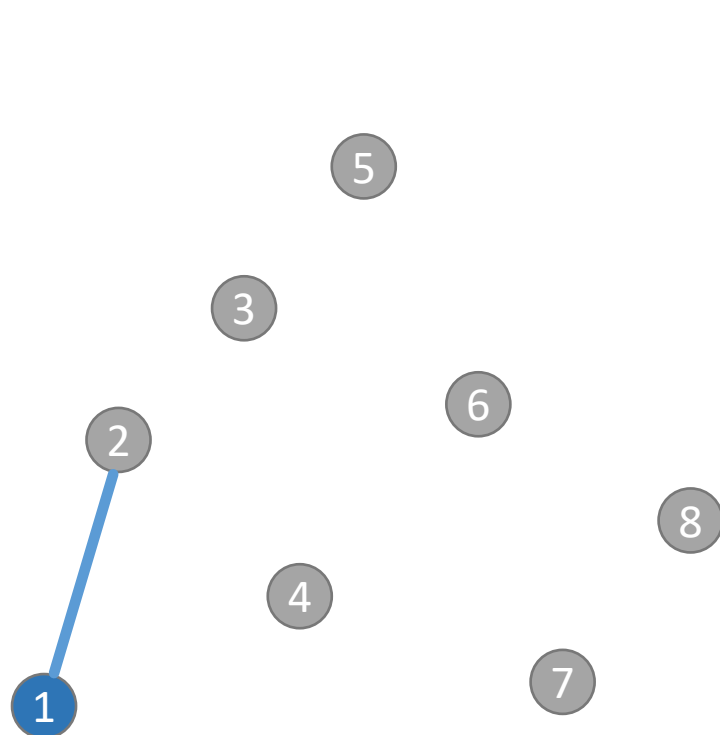


# CLASSIX

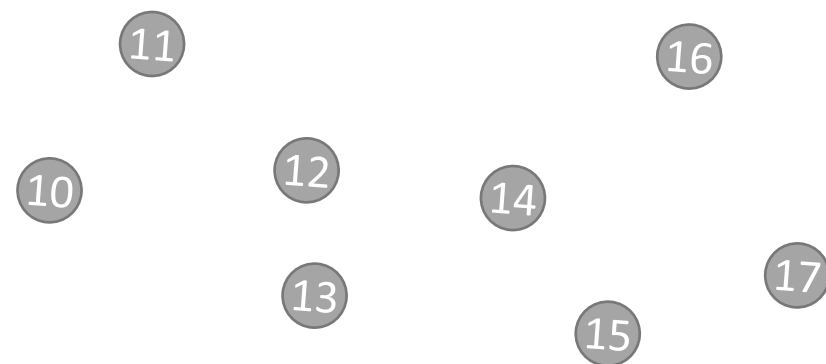
Fast and Explainable Clustering





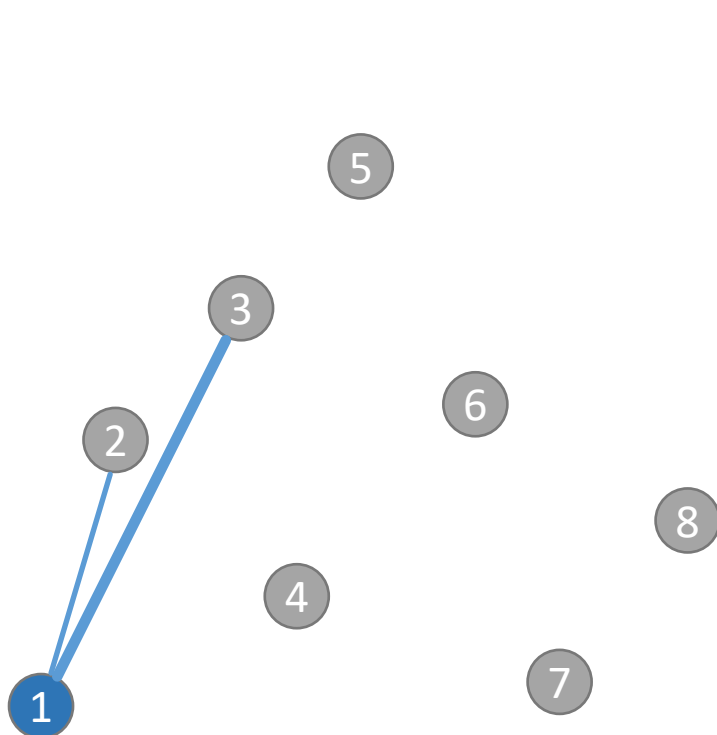


Starting from the first point we compute distances to the following points.

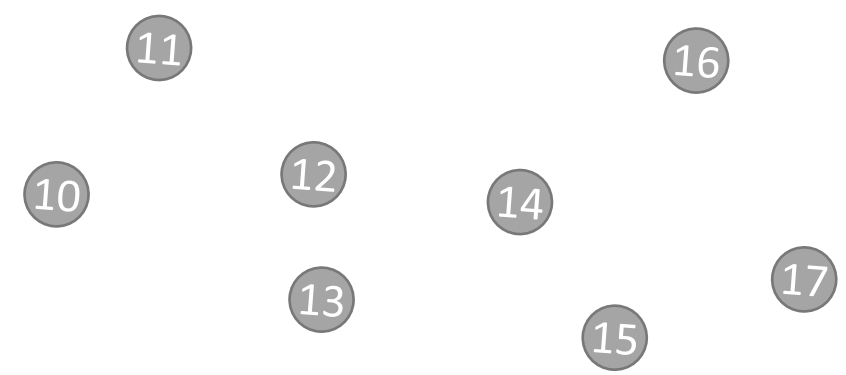


# CLASSIX

Fast and Explainable Clustering

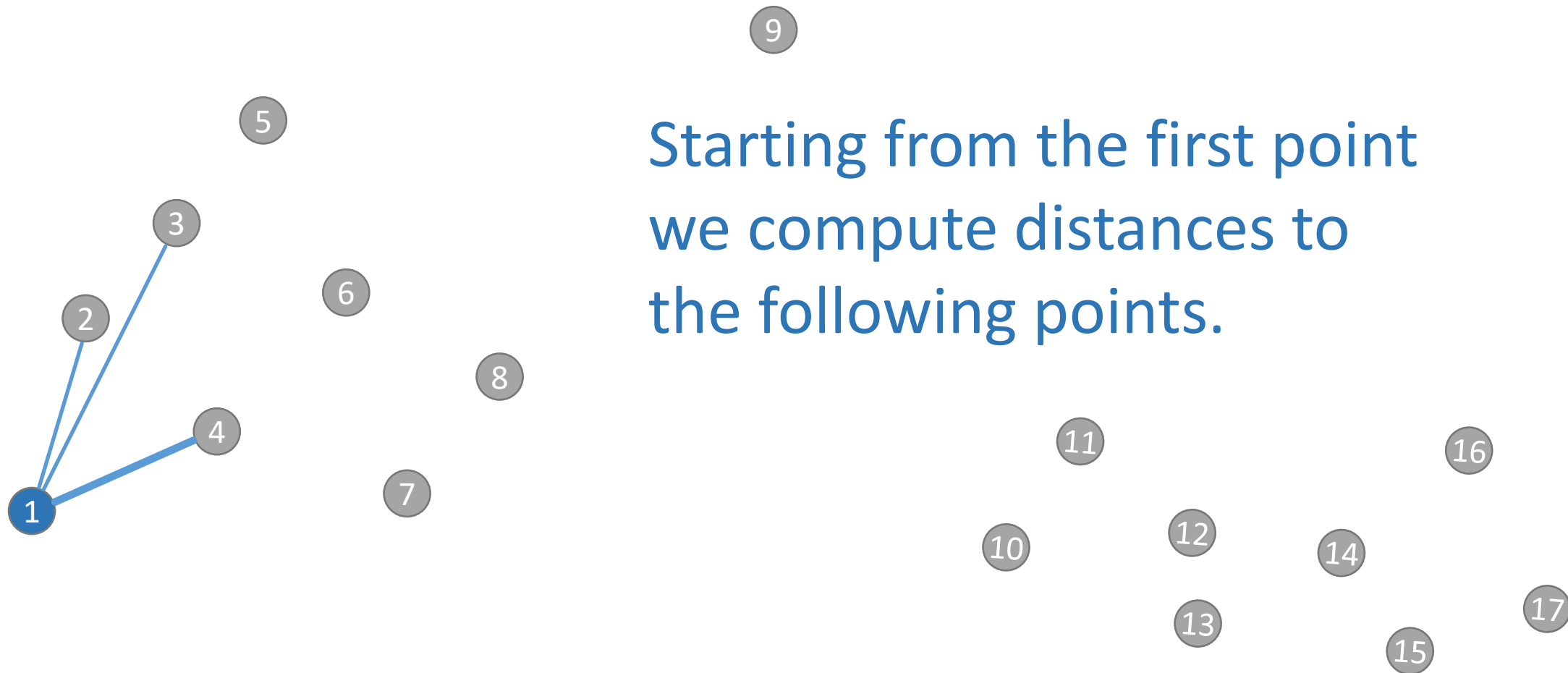


Starting from the first point we compute distances to the following points.



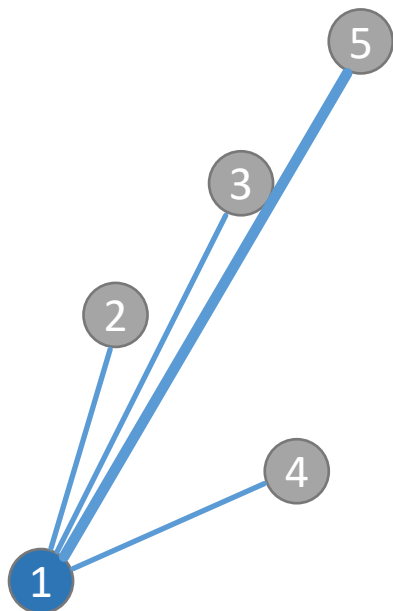
# CLASSIX

Fast and Explainable Clustering



# CLASSIX

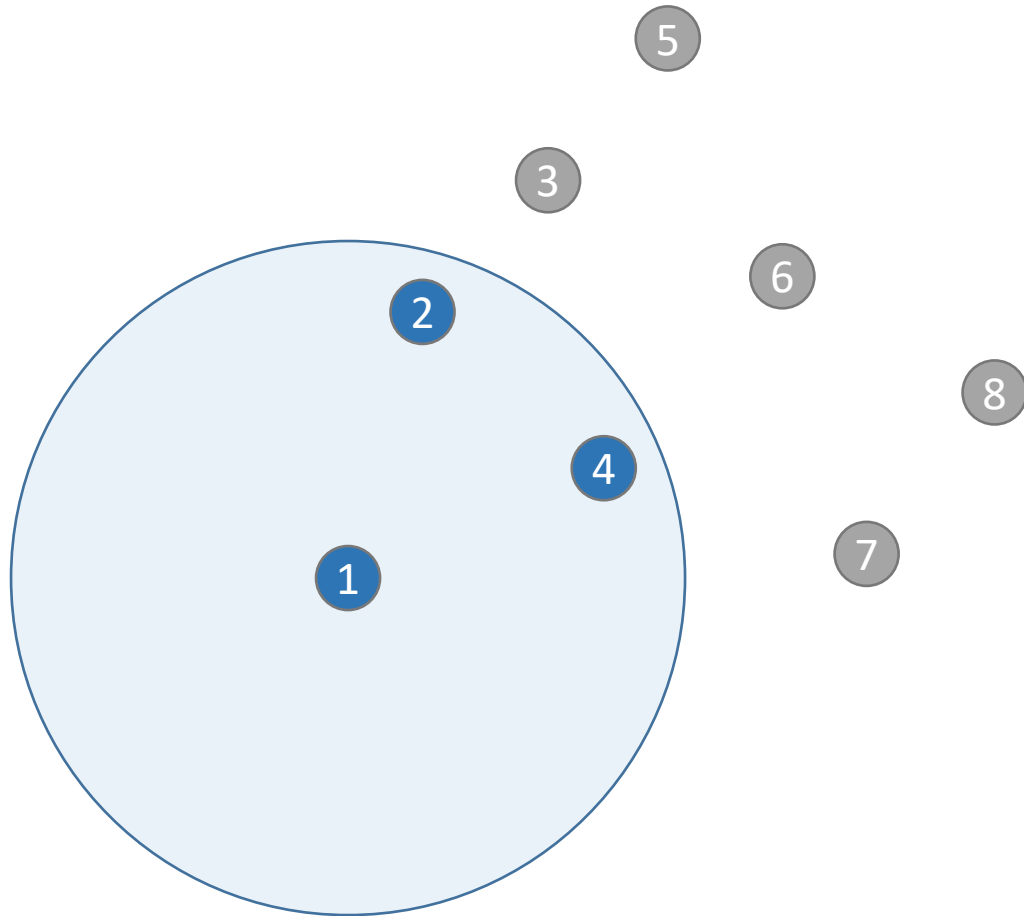
Fast and Explainable Clustering



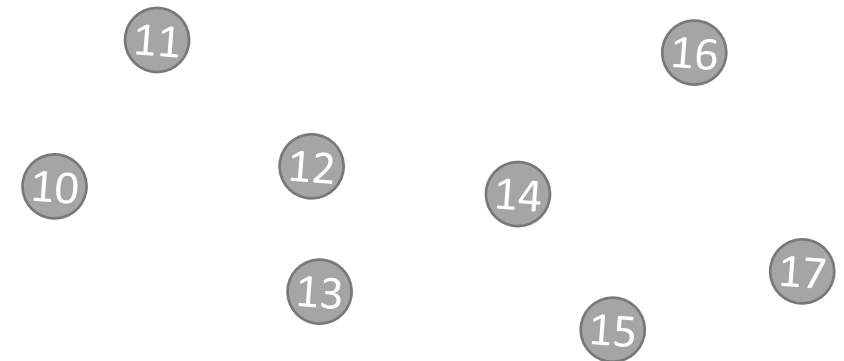
Starting from the first point we compute distances to the following points.

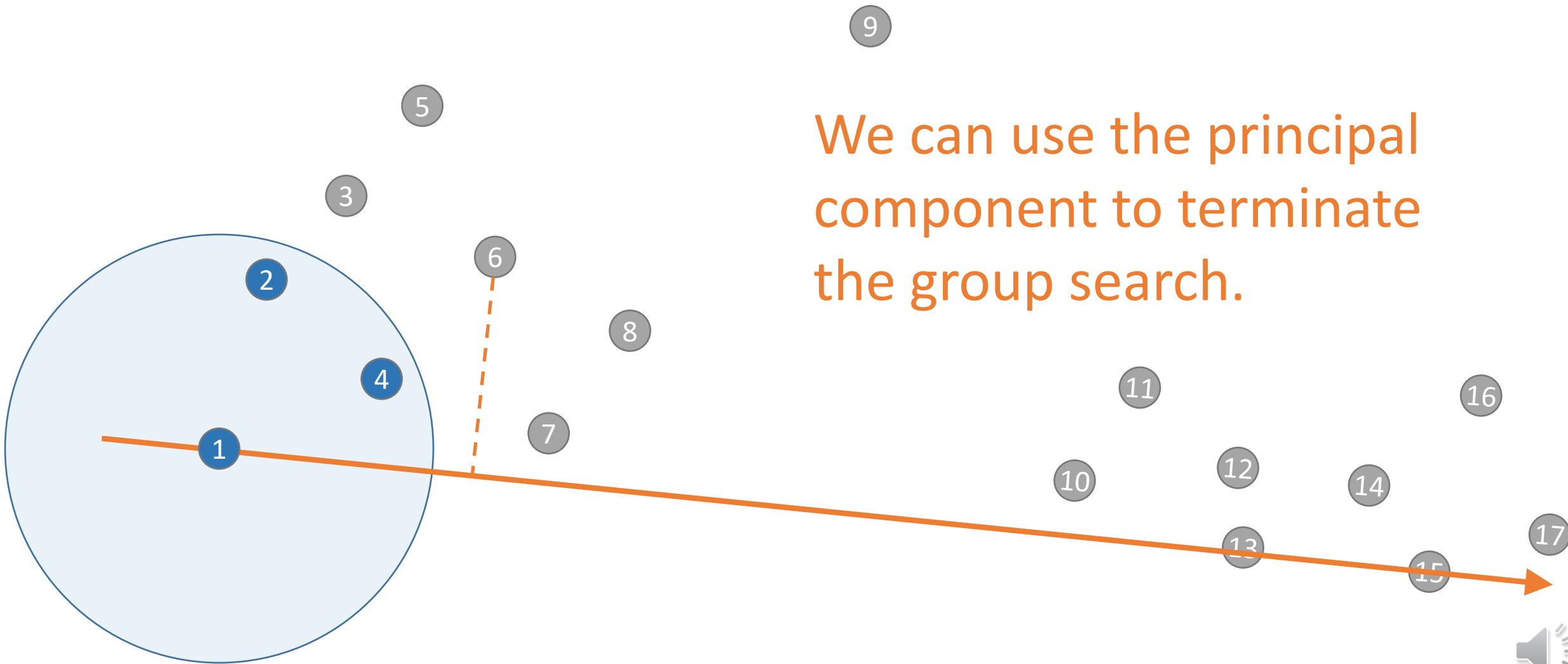
# CLASSIX

Fast and Explainable Clustering



All points within a predefined radius become part of a group.





We can use the principal component to terminate the group search.

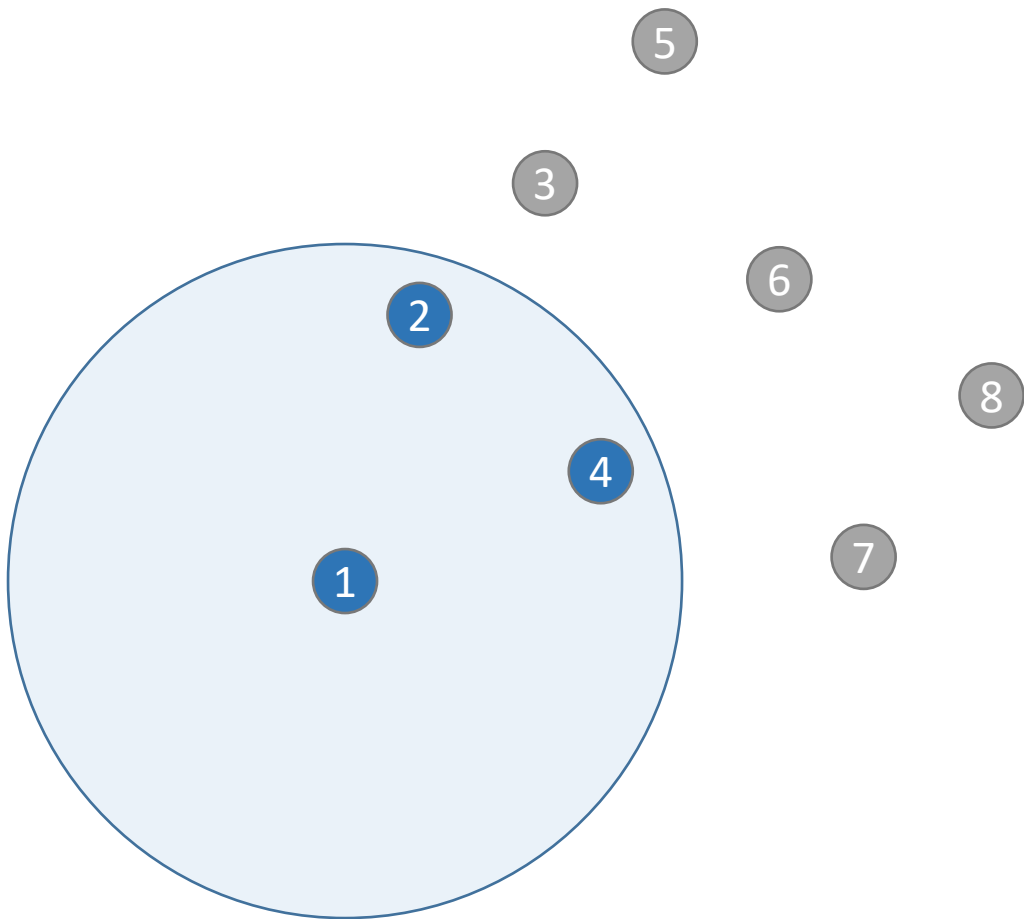




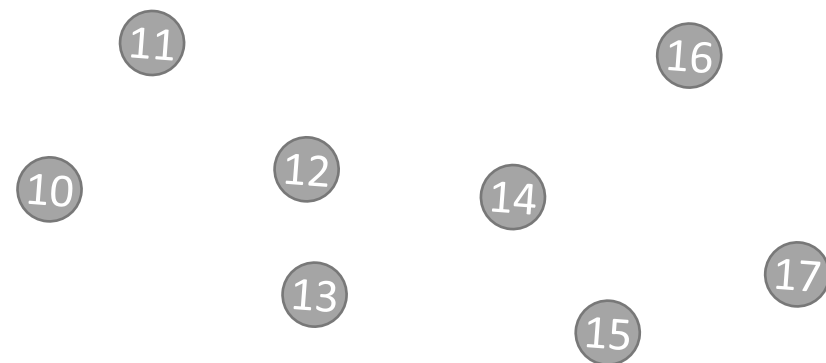


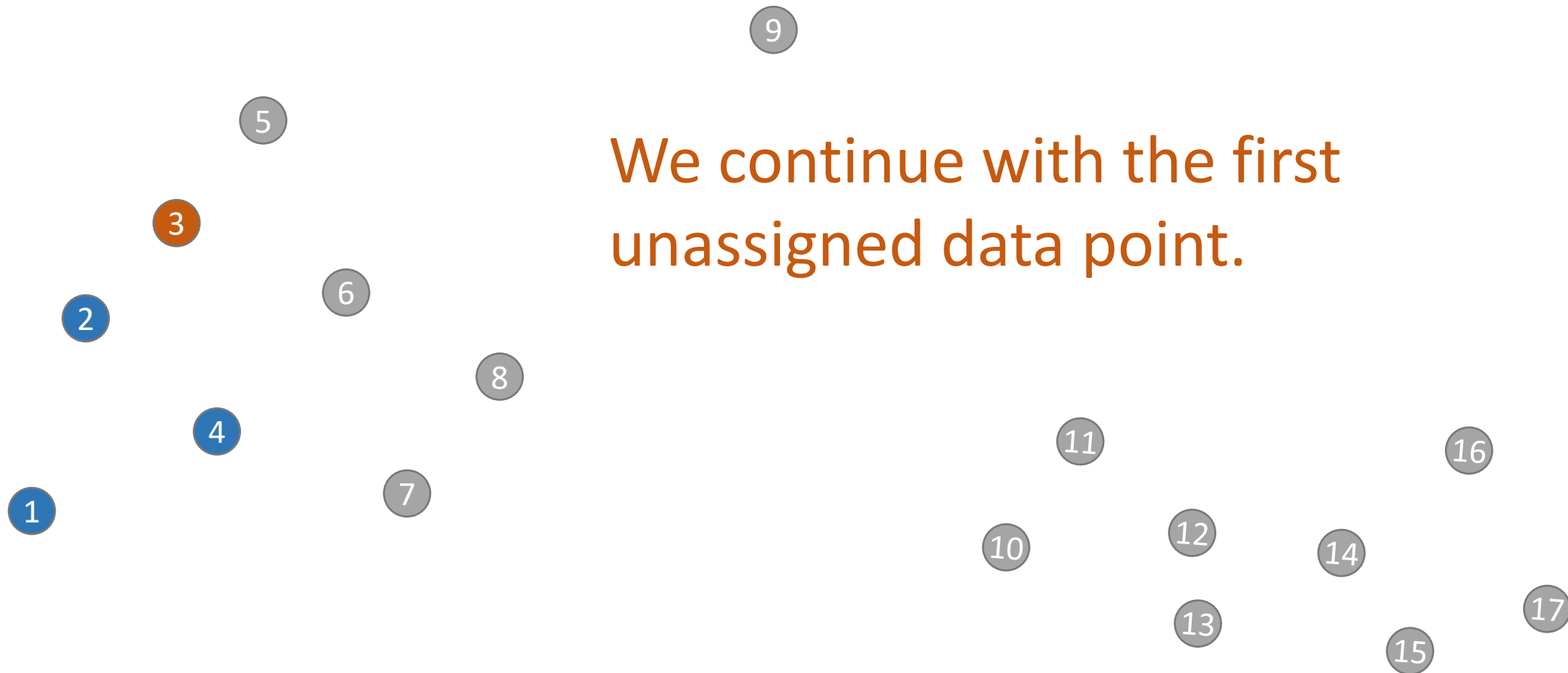
# CLASSIX

Fast and Explainable Clustering



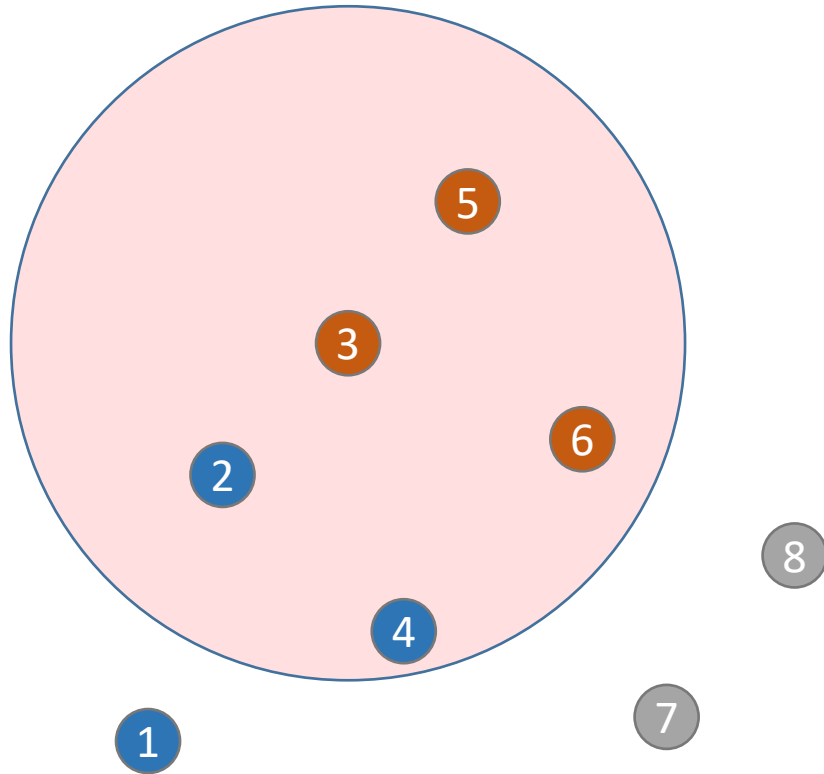
No need for range query data structures to implement this.



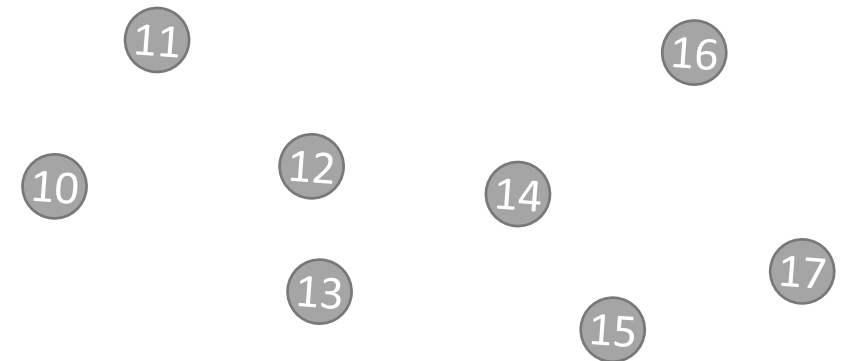


We continue with the first unassigned data point.



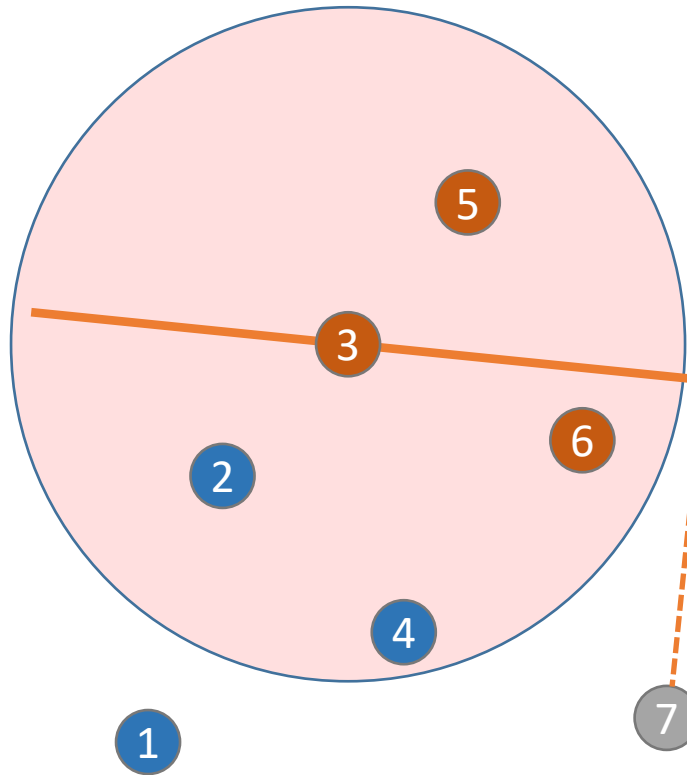


And group all data points within the predefined radius.



# CLASSIX

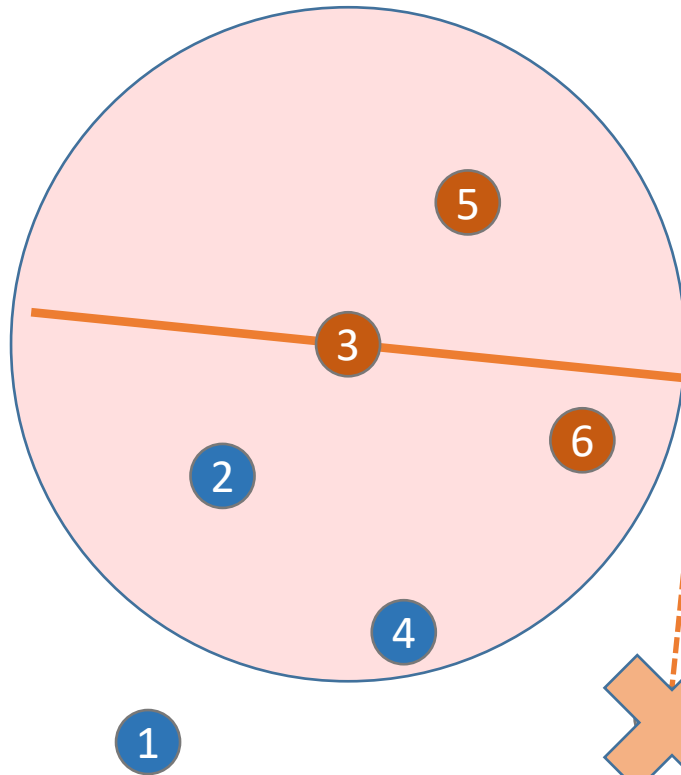
Fast and Explainable Clustering



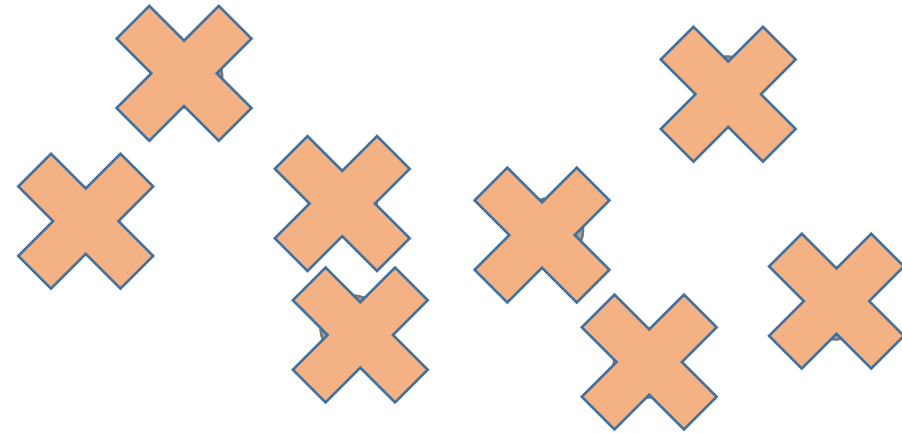
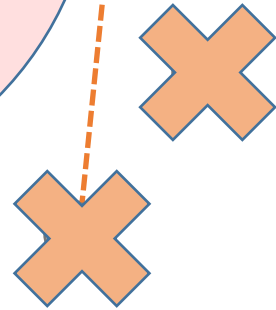
Using the same trick as before,  
point 7 and all the following ones  
do not need to be considered.

# CLASSIX

Fast and Explainable Clustering

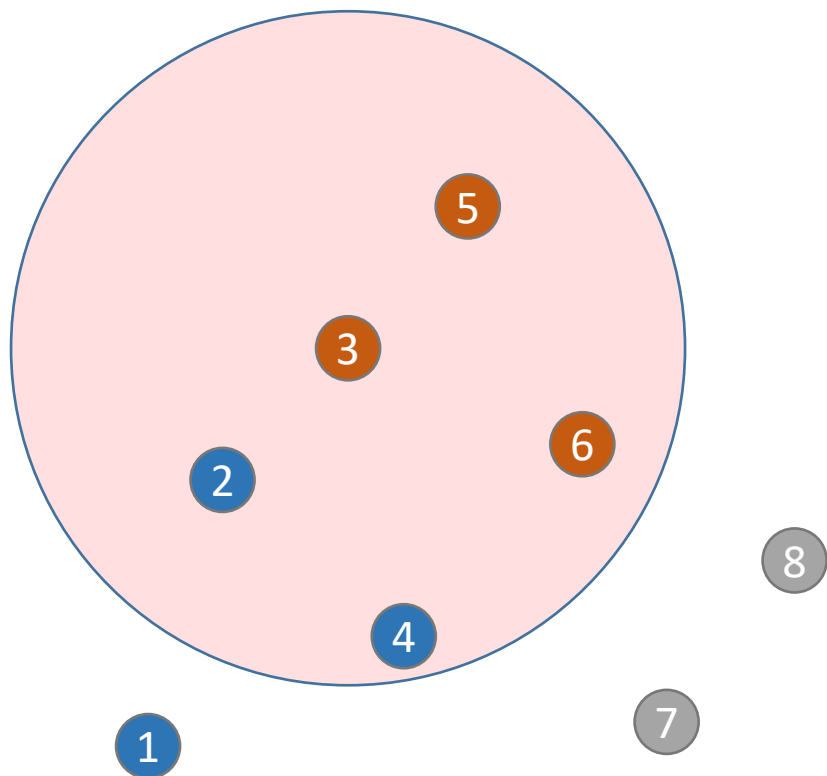


Using the same trick as before,  
point 7 and all the following ones  
do not need to be considered.



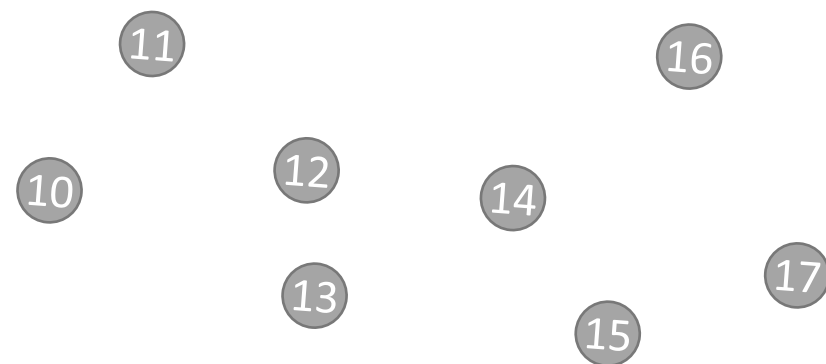
# CLASSIX

Fast and Explainable Clustering



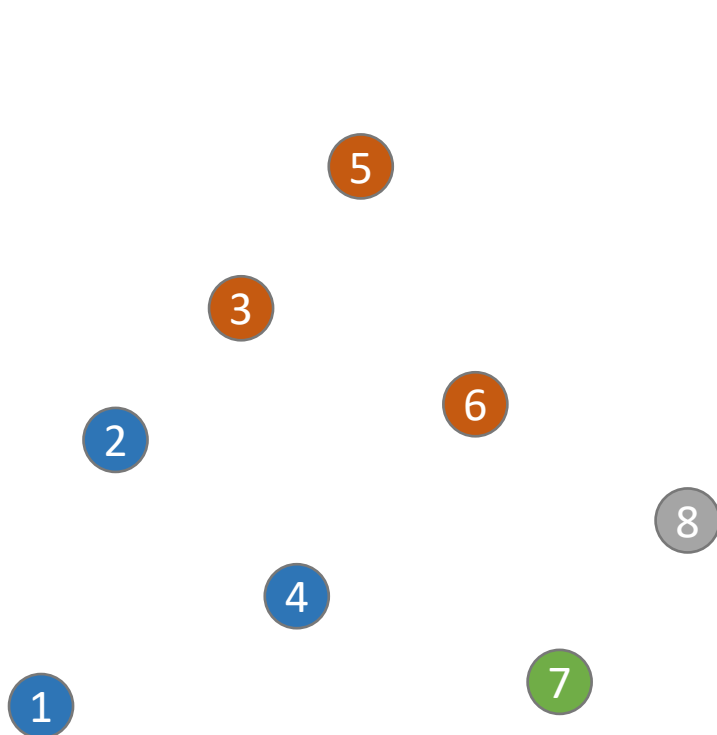
9

This early search termination keeps the number of distance calculations low.

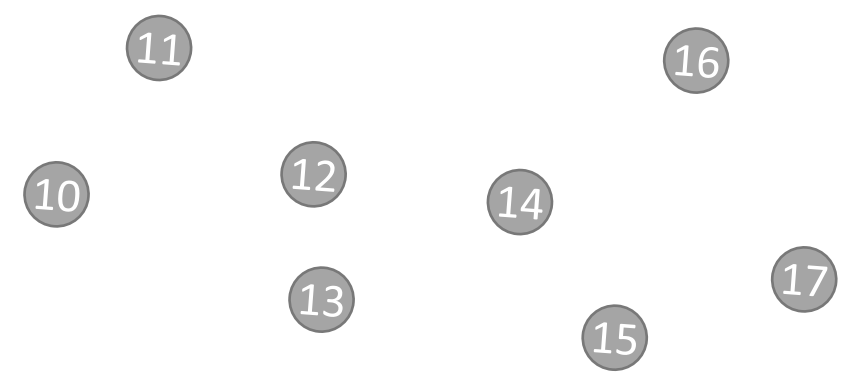


# CLASSIX

Fast and Explainable Clustering



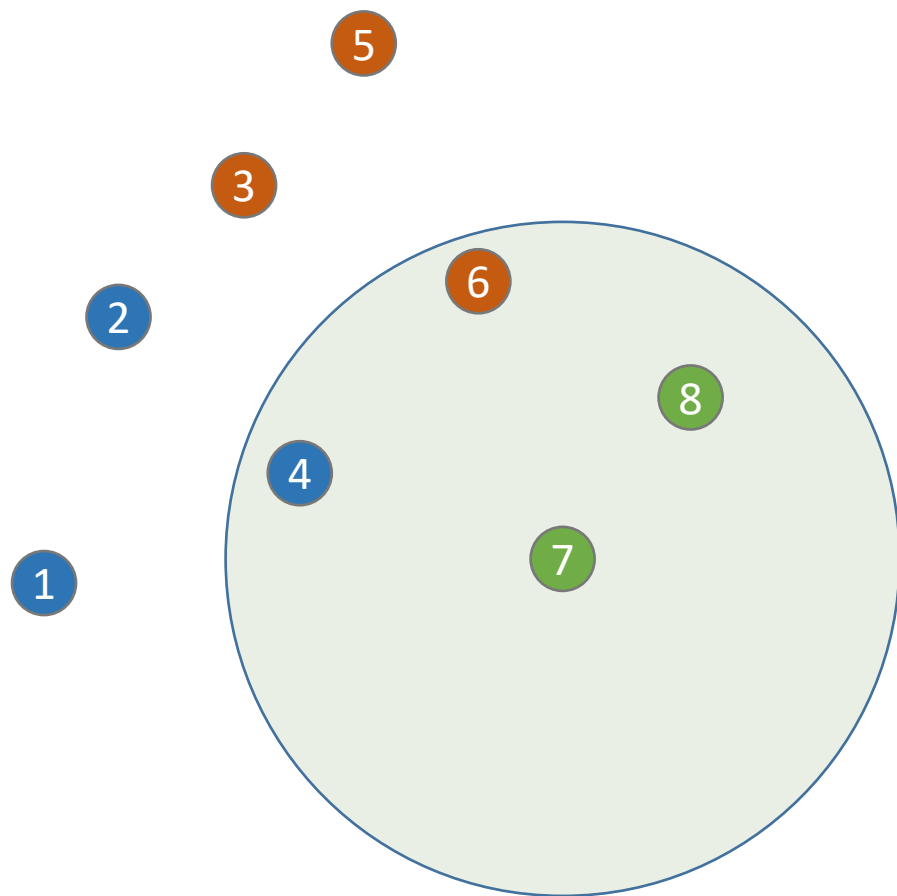
We continue grouping the data points in their order.



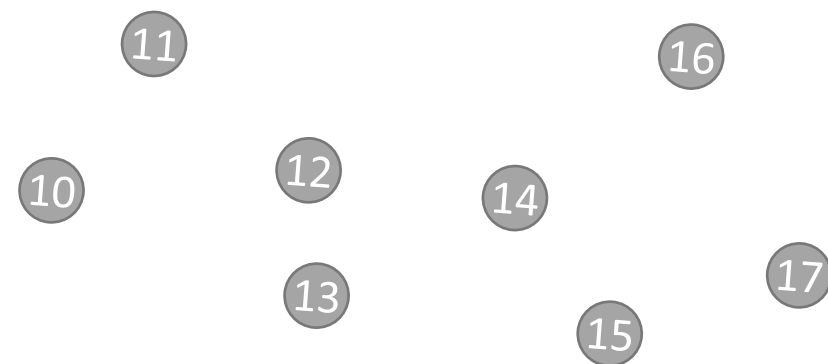


# CLASSIX

Fast and Explainable Clustering

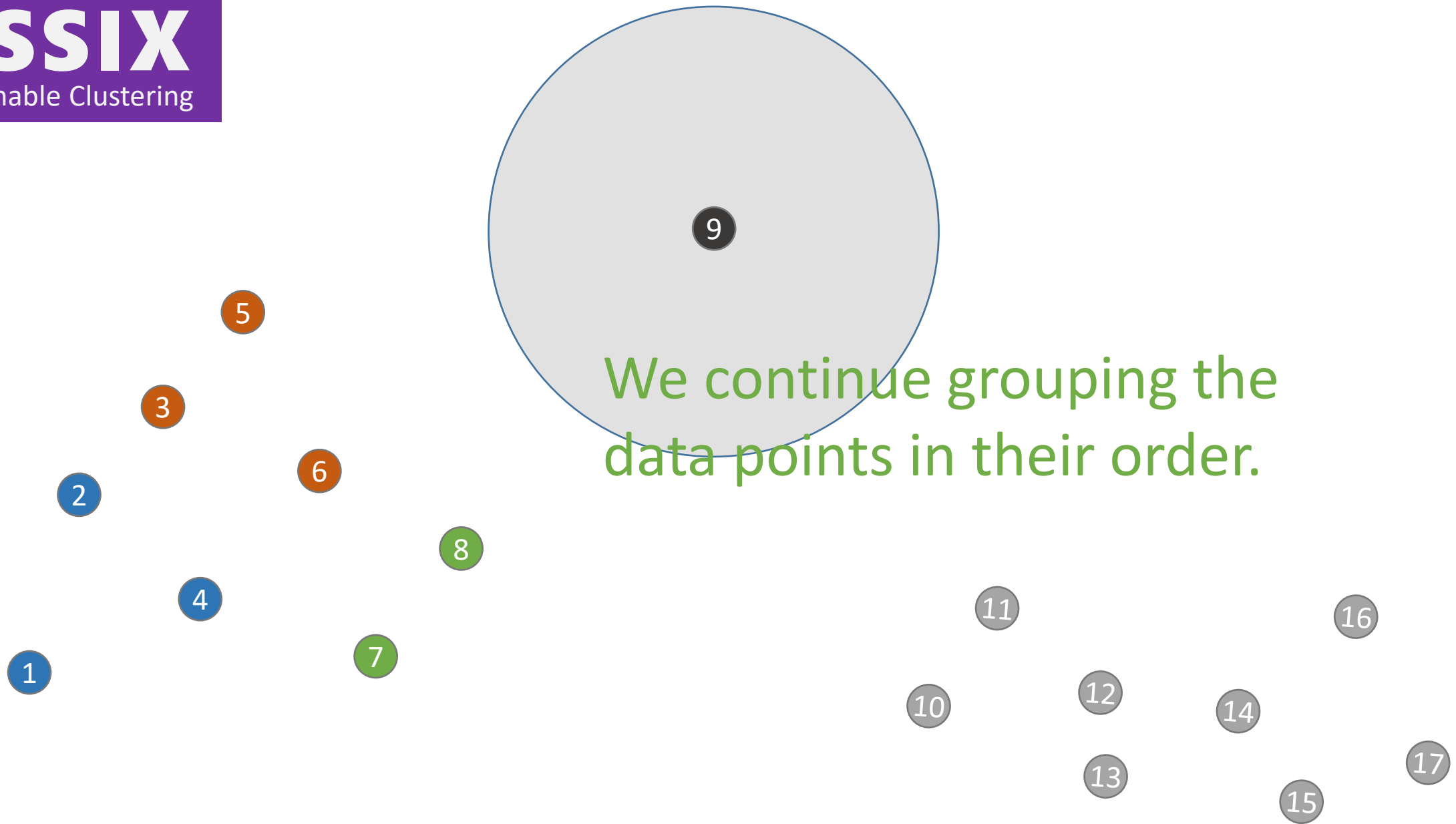


We continue grouping the data points in their order.



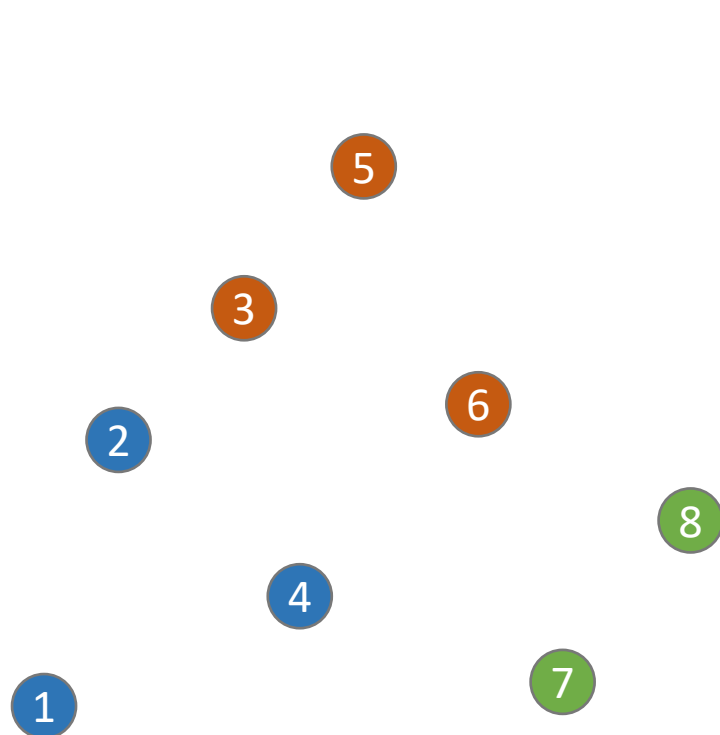
# CLASSIX

Fast and Explainable Clustering



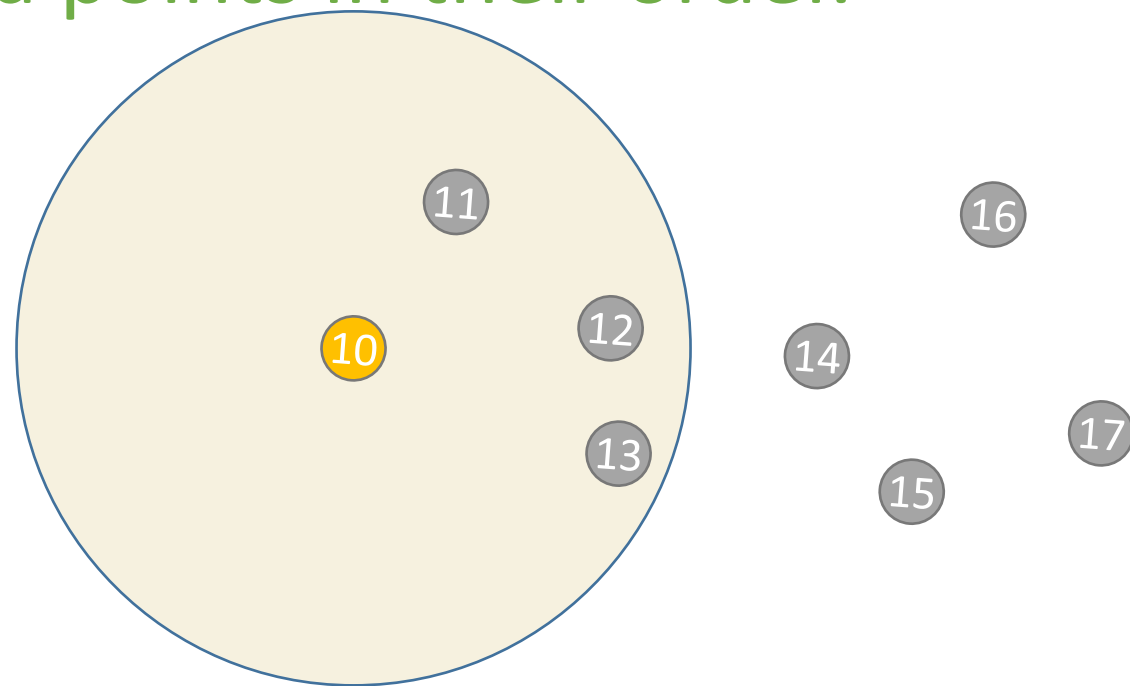
# CLASSIX

Fast and Explainable Clustering



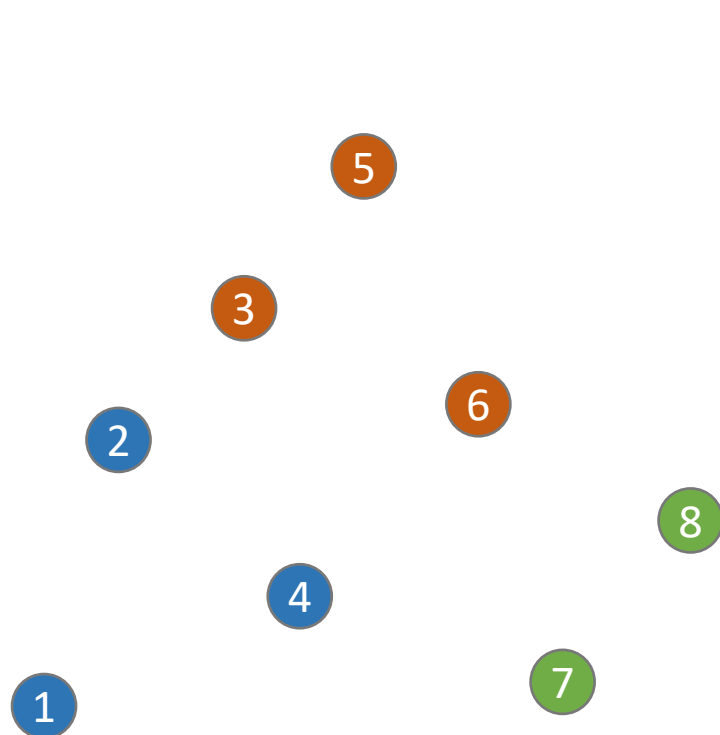
9

We continue grouping the data points in their order.



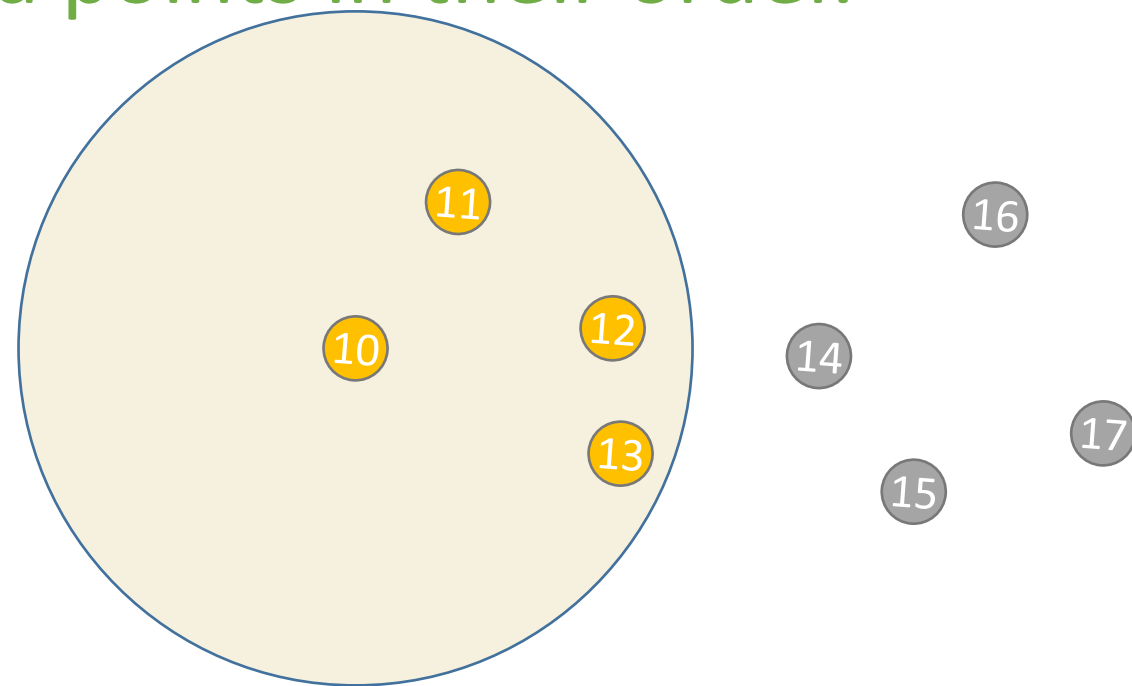
# CLASSIX

Fast and Explainable Clustering



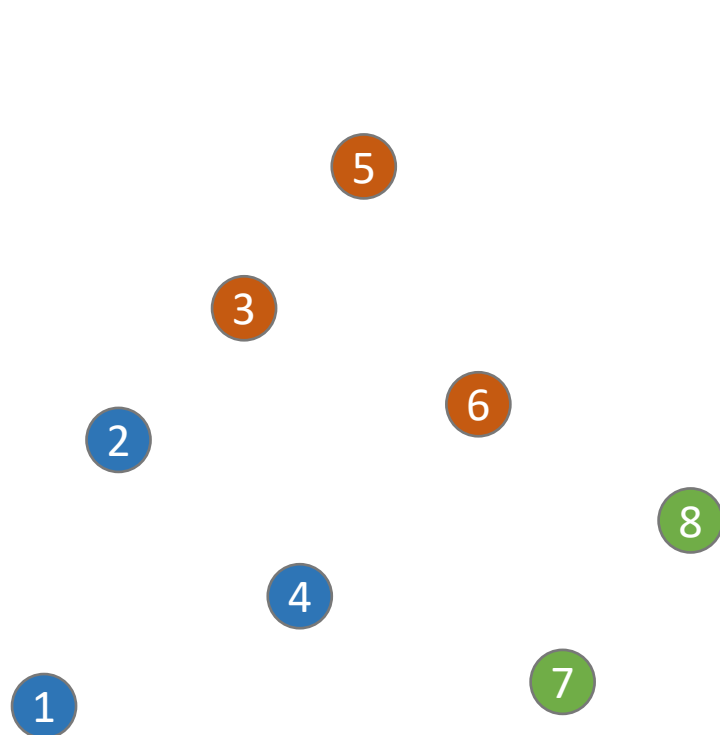
9

We continue grouping the data points in their order.

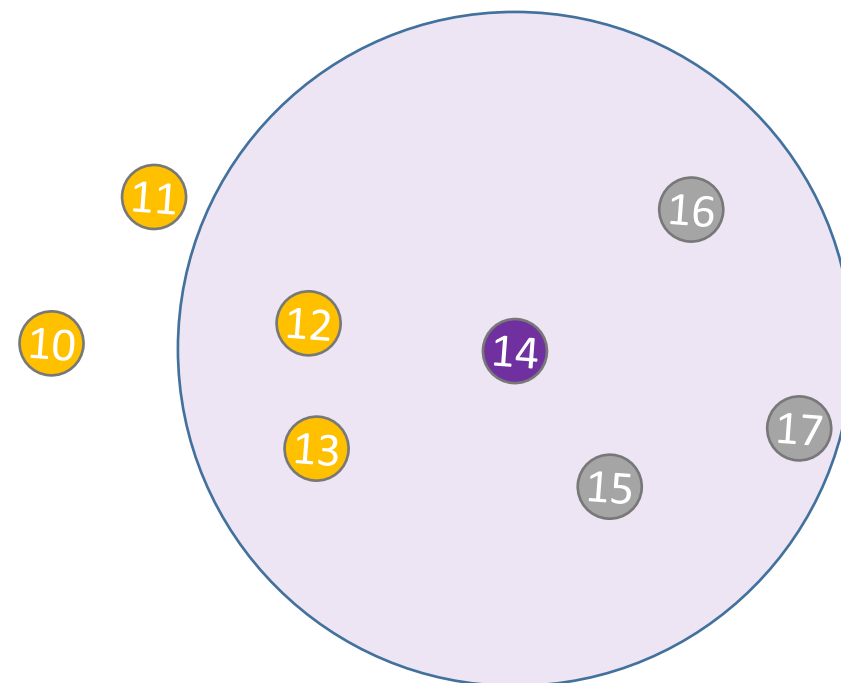


# CLASSIX

Fast and Explainable Clustering

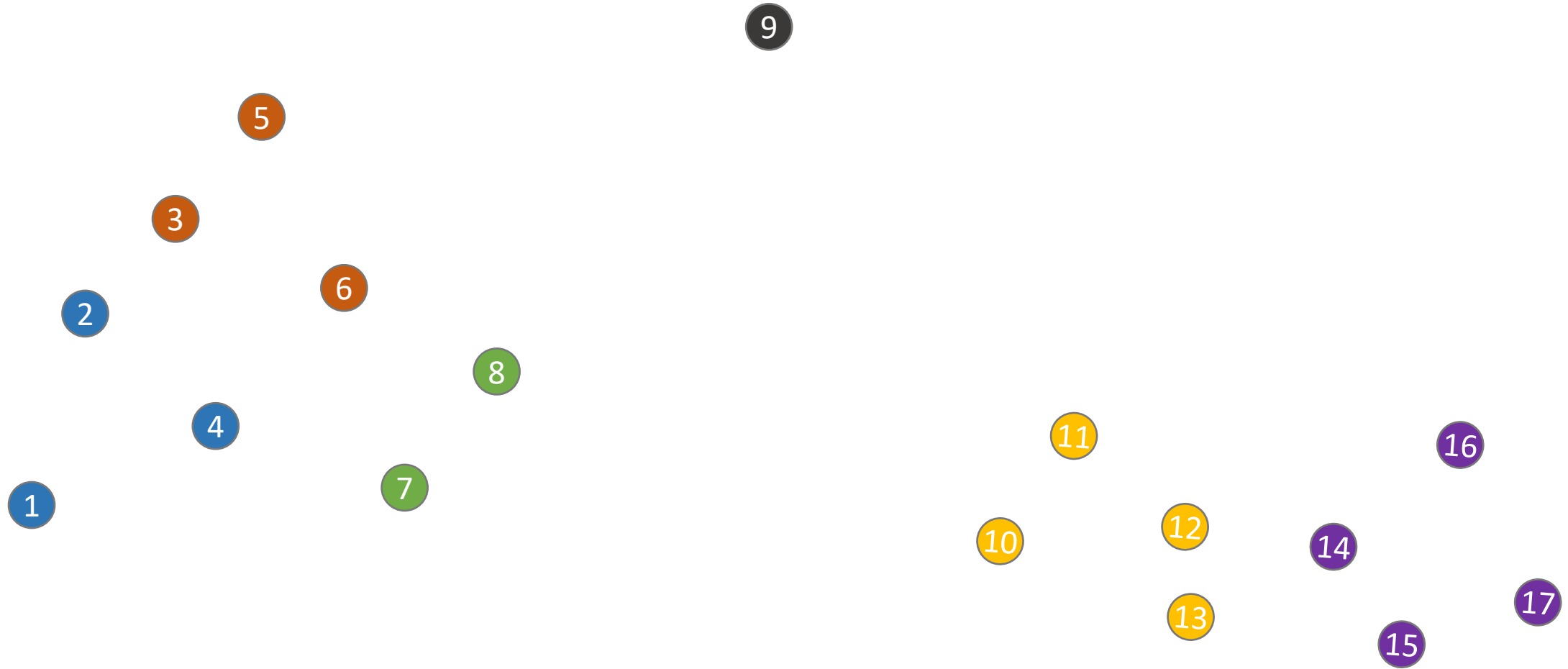


We continue grouping the data points in their order.



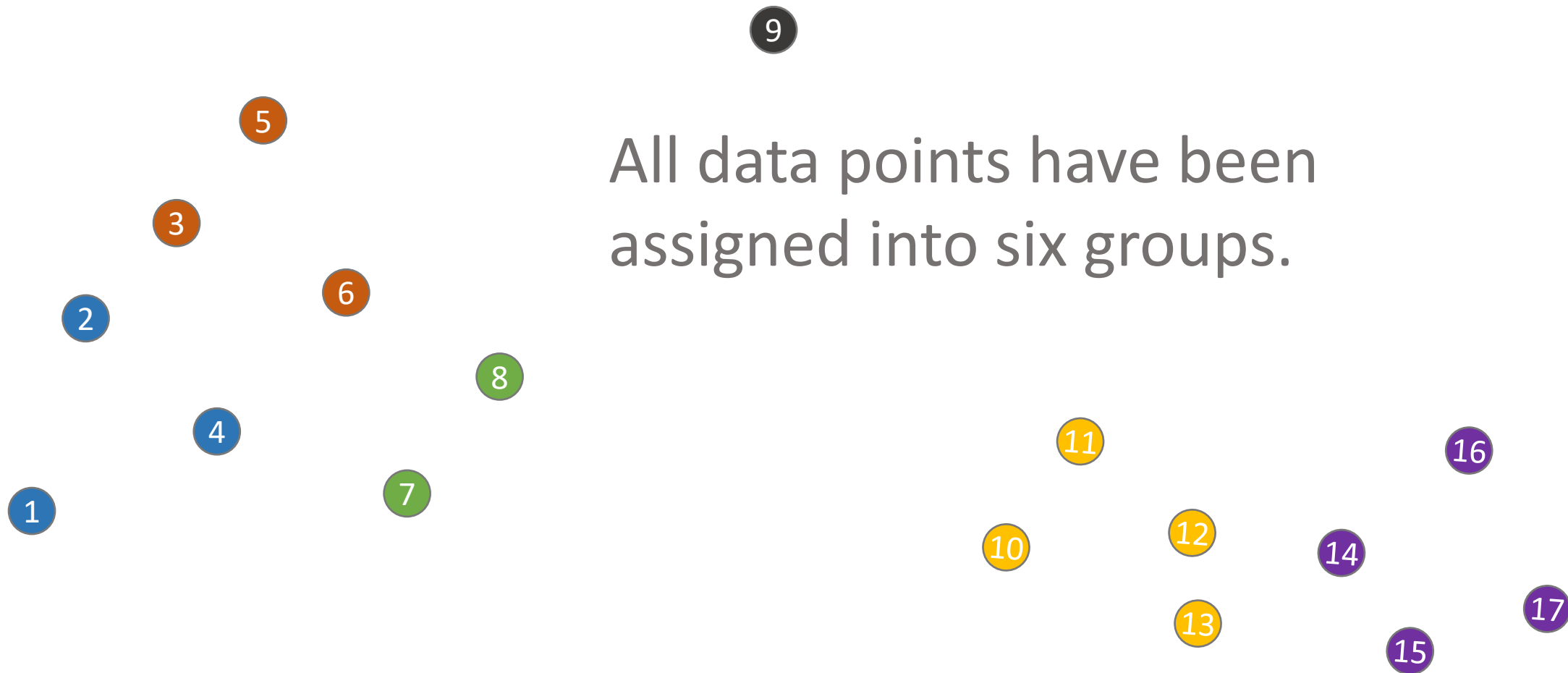
# CLASSIX

Fast and Explainable Clustering

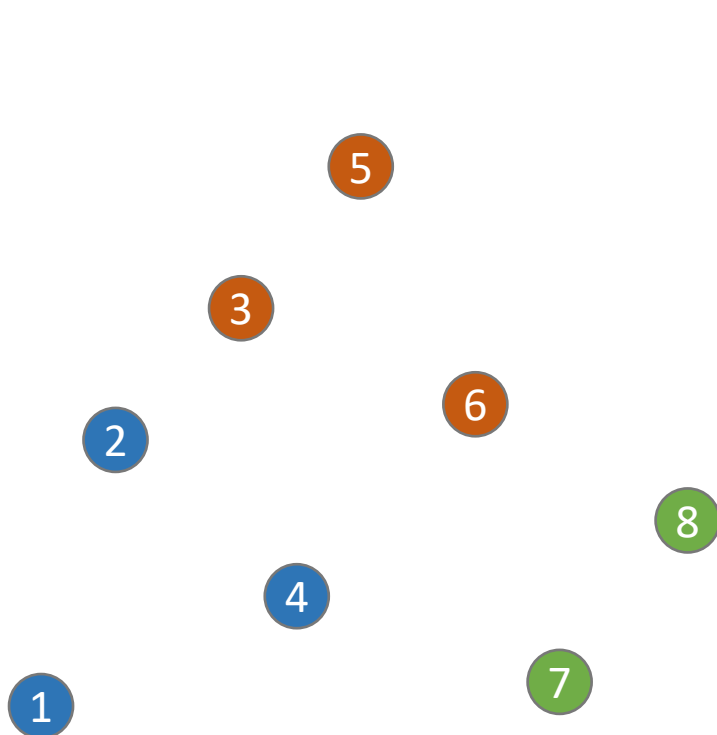


# CLASSIX

Fast and Explainable Clustering

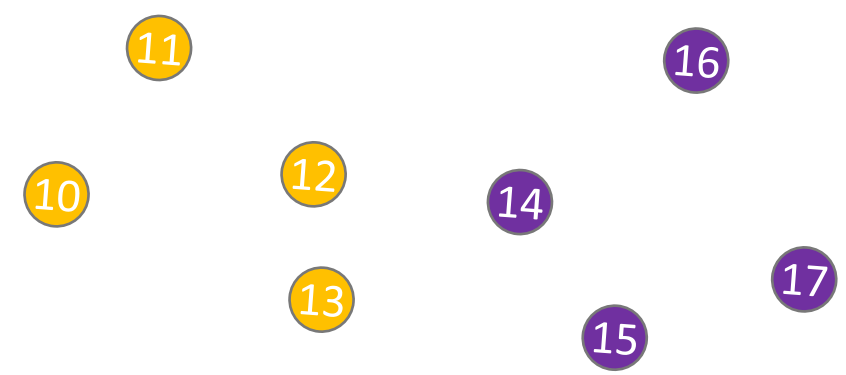


All data points have been assigned into six groups.



9

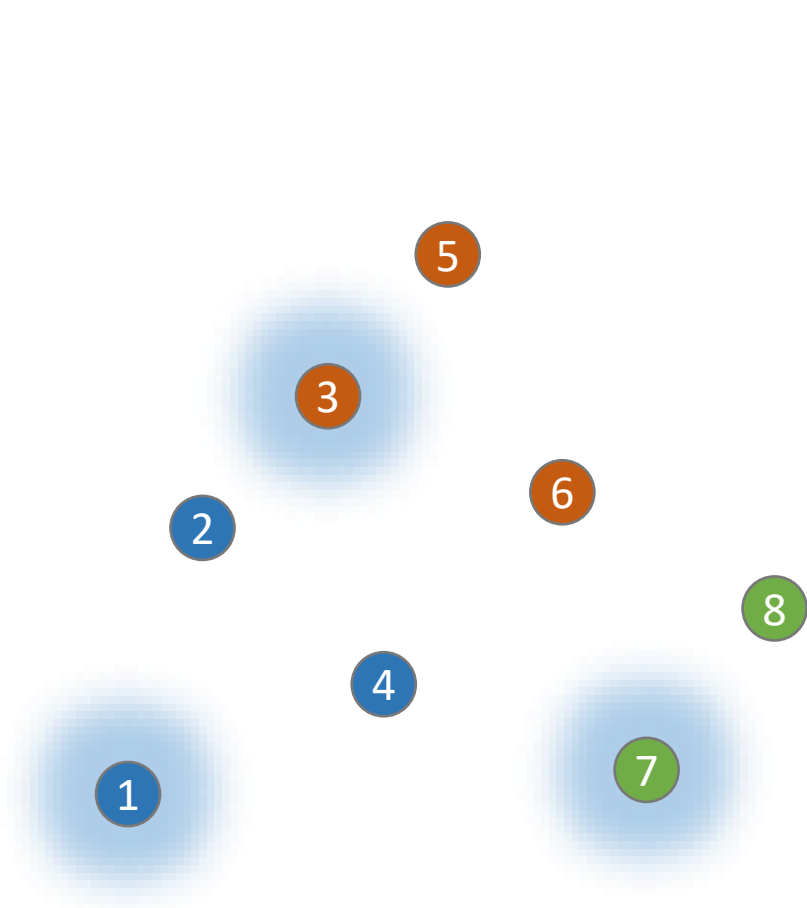
Now we just have to merge these groups into clusters.



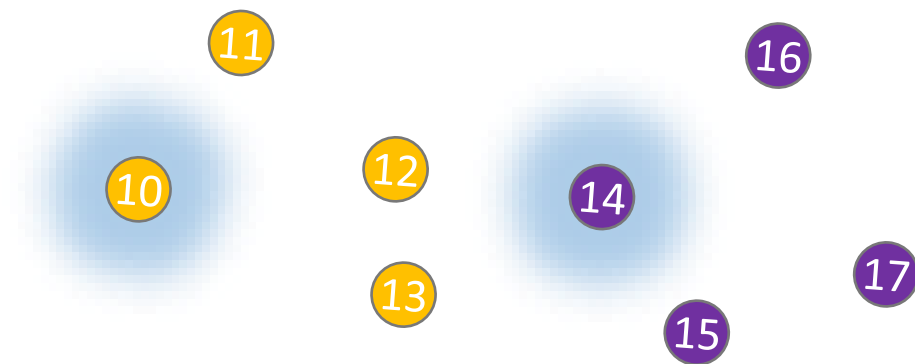


# CLASSIX

Fast and Explainable Clustering

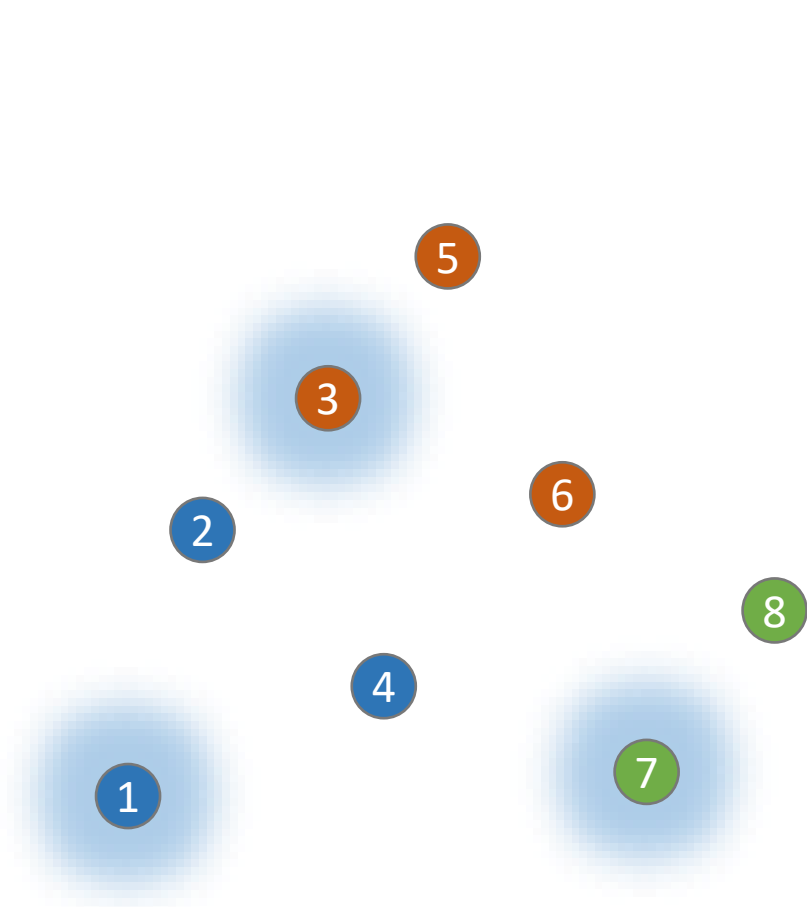


This can be done efficiently by using the starting points of each group.

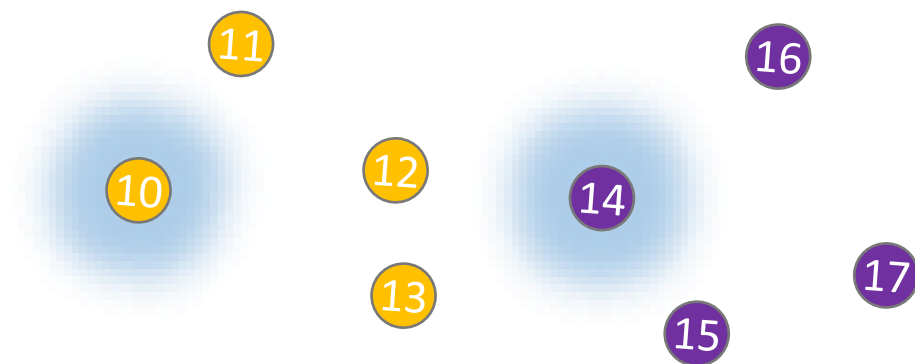


# CLASSIX

Fast and Explainable Clustering

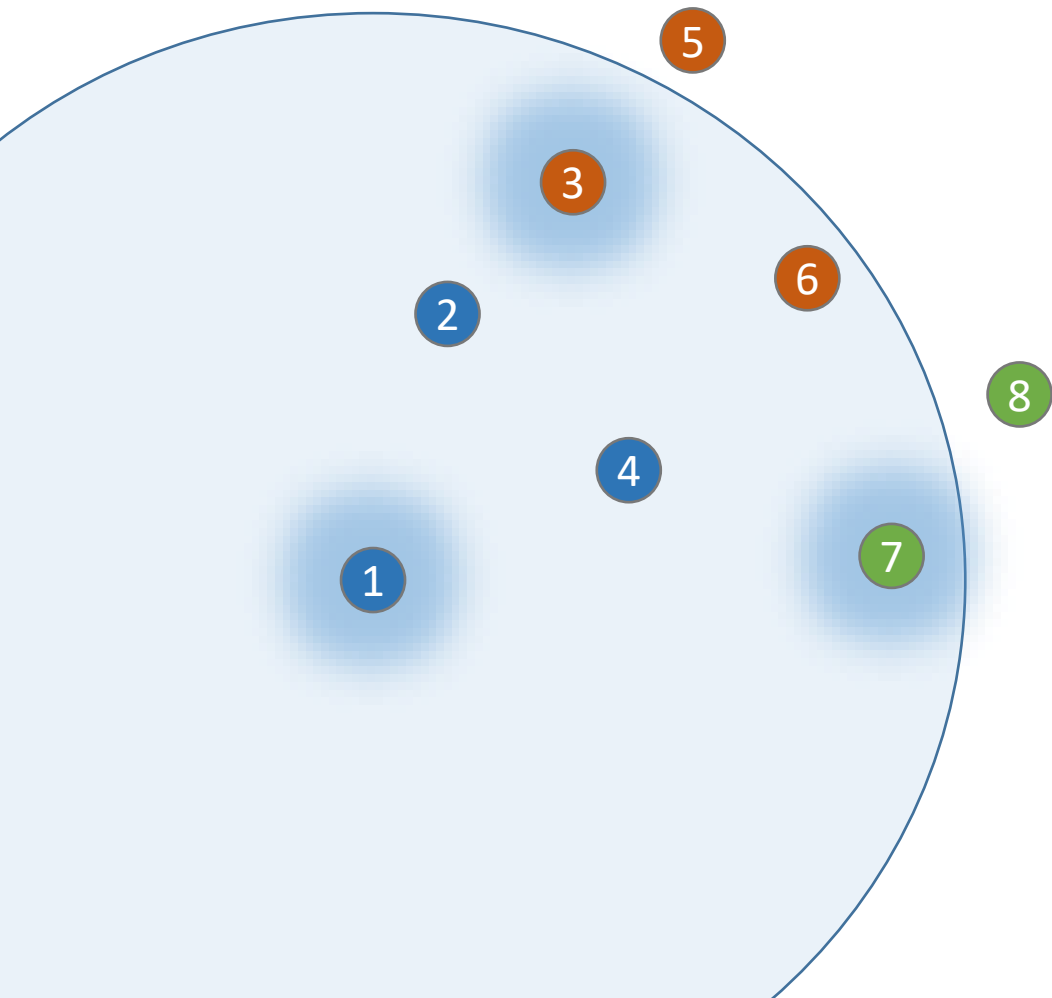


The starting points are already sorted and we can again use early search termination.

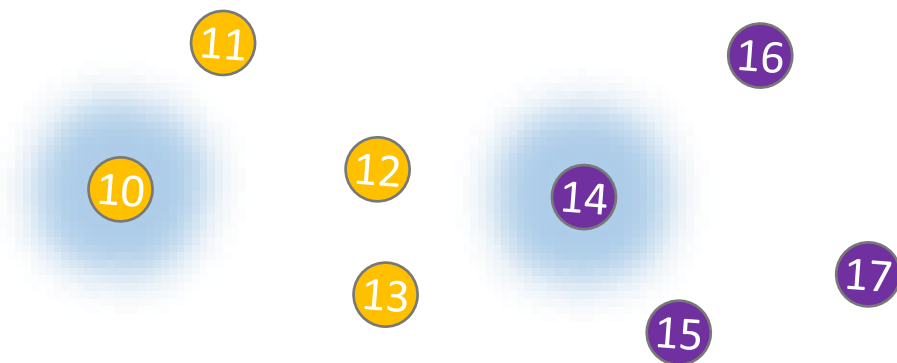


# CLASSIX

Fast and Explainable Clustering

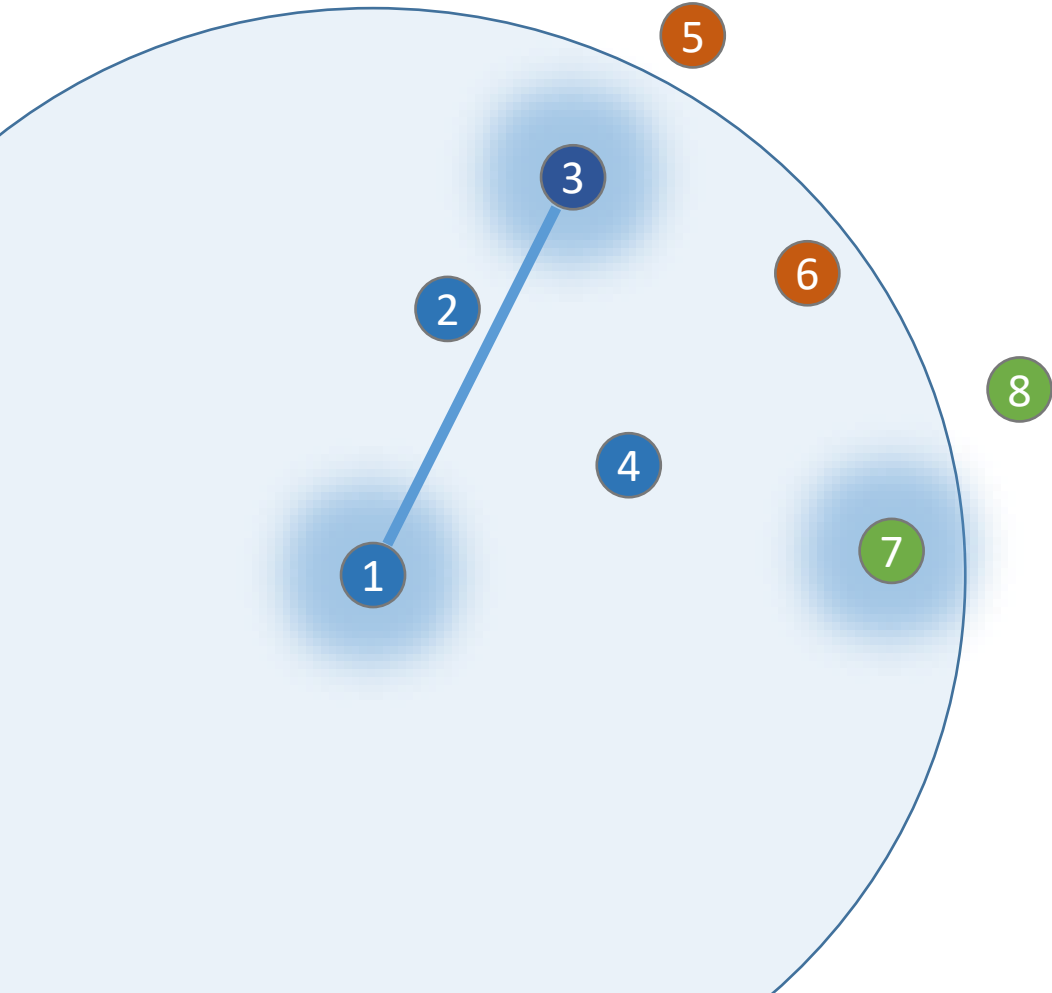


We merge two starting points within 1.5x the predefined radius.

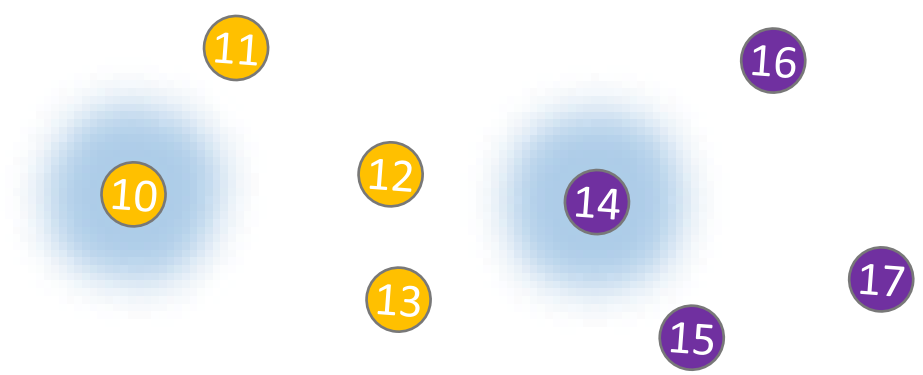


# CLASSIX

Fast and Explainable Clustering

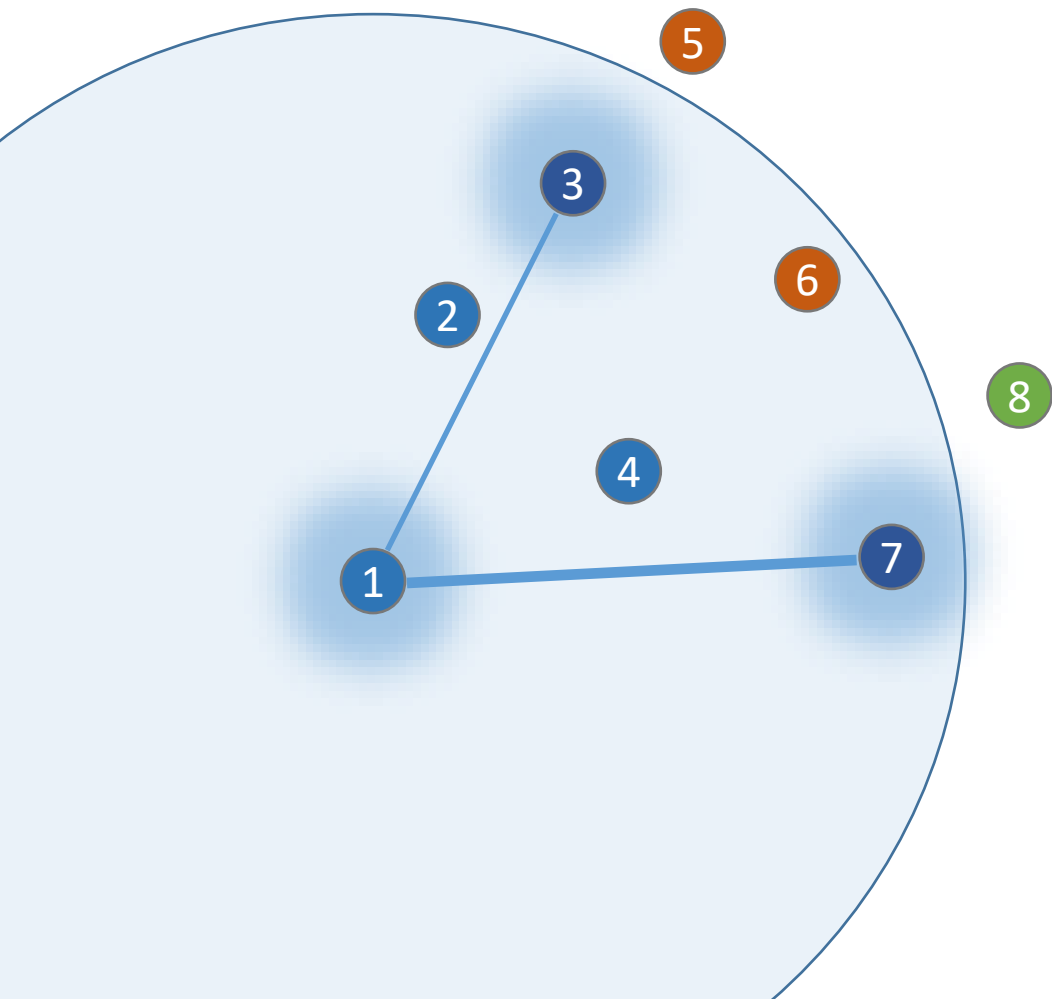


We merge two starting points within  $1.5x$  the predefined radius.

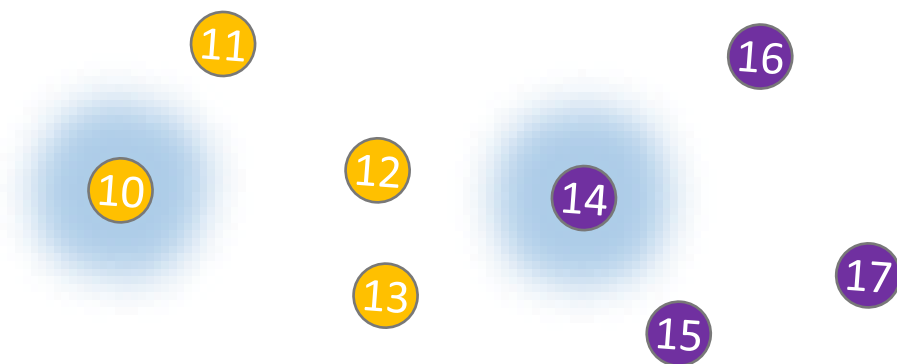


# CLASSIX

Fast and Explainable Clustering

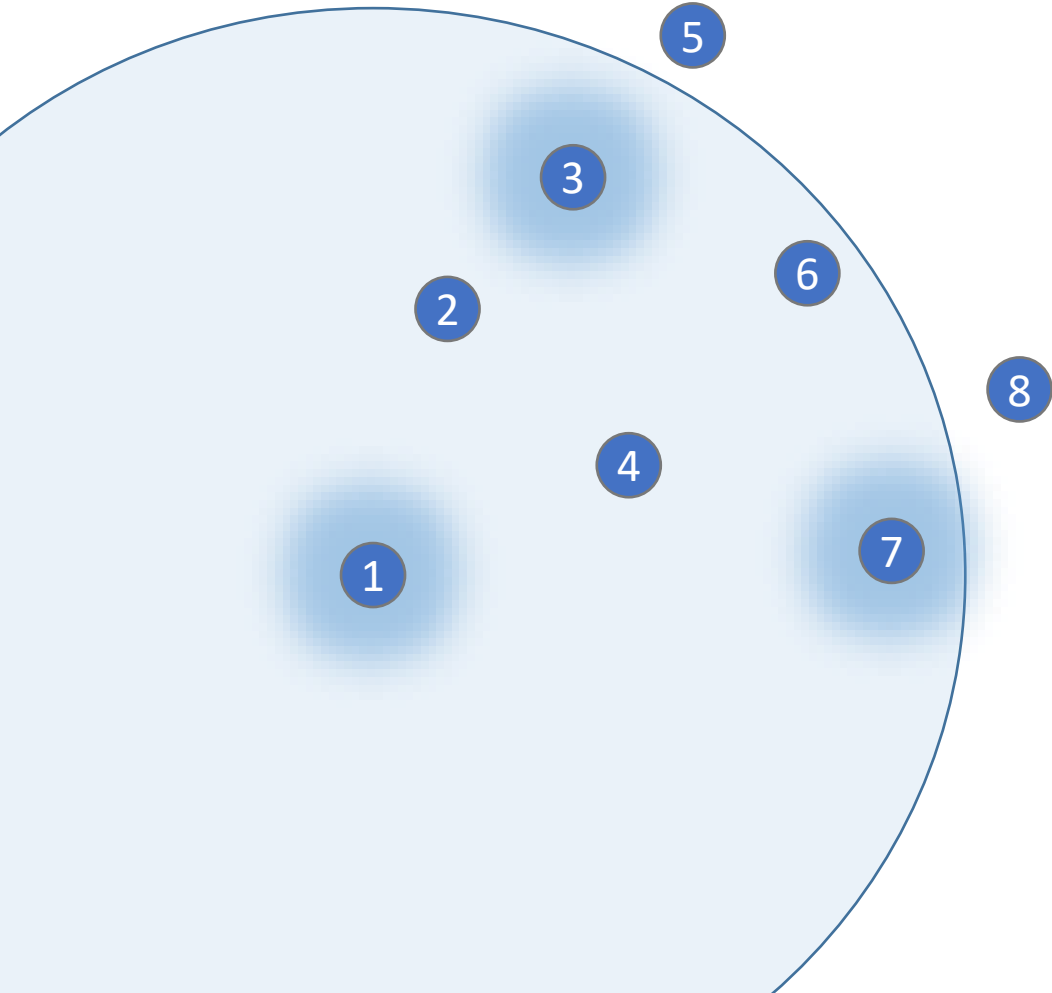


We merge two starting points within  $1.5x$  the predefined radius.

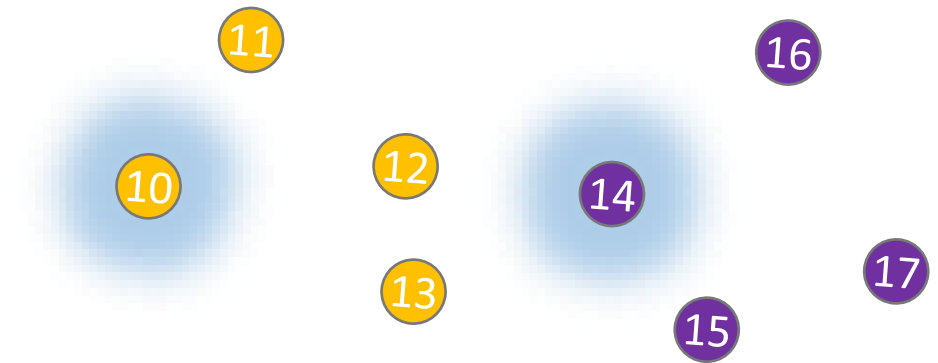


# CLASSIX

Fast and Explainable Clustering

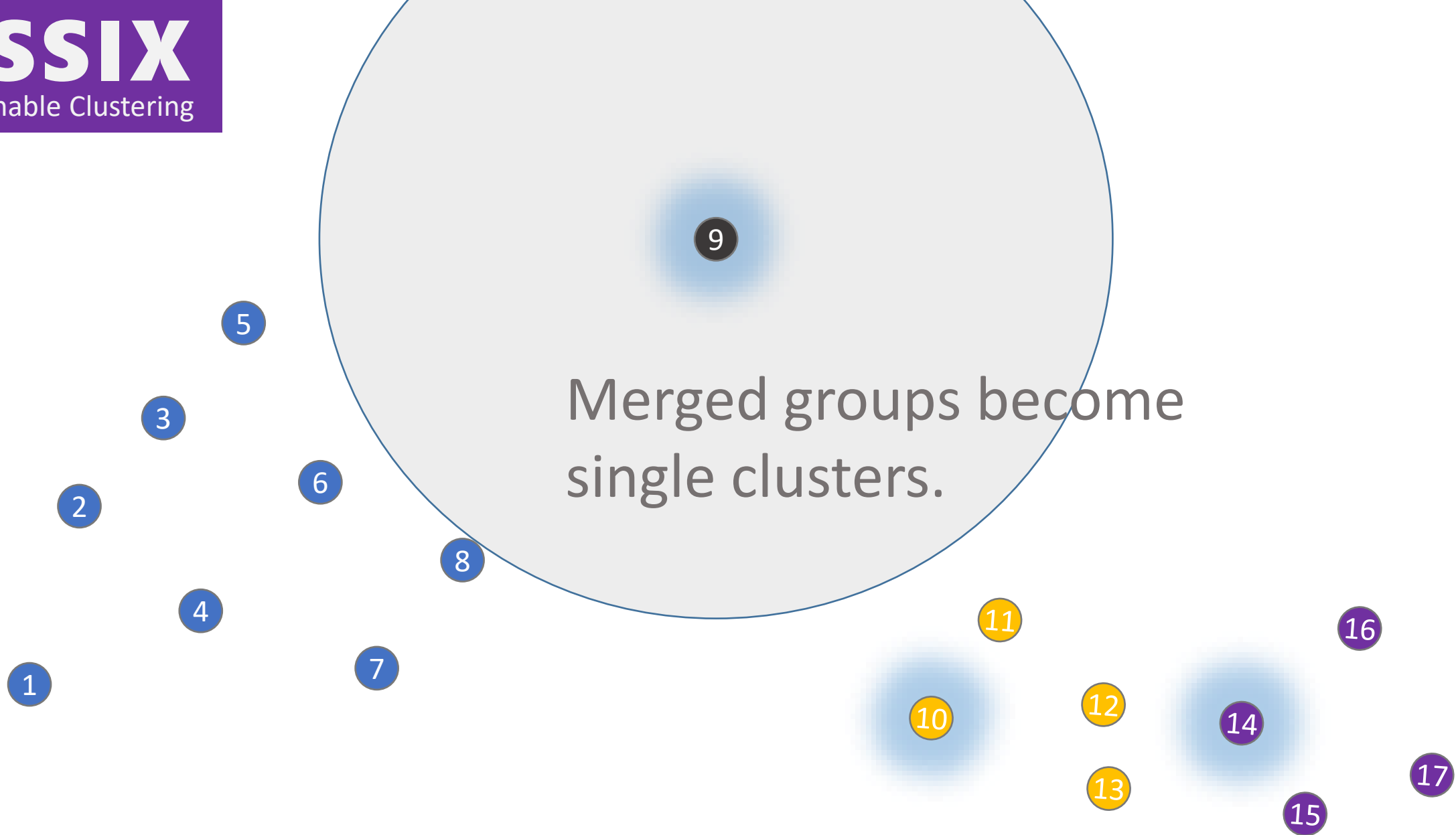


Merged groups become single clusters.



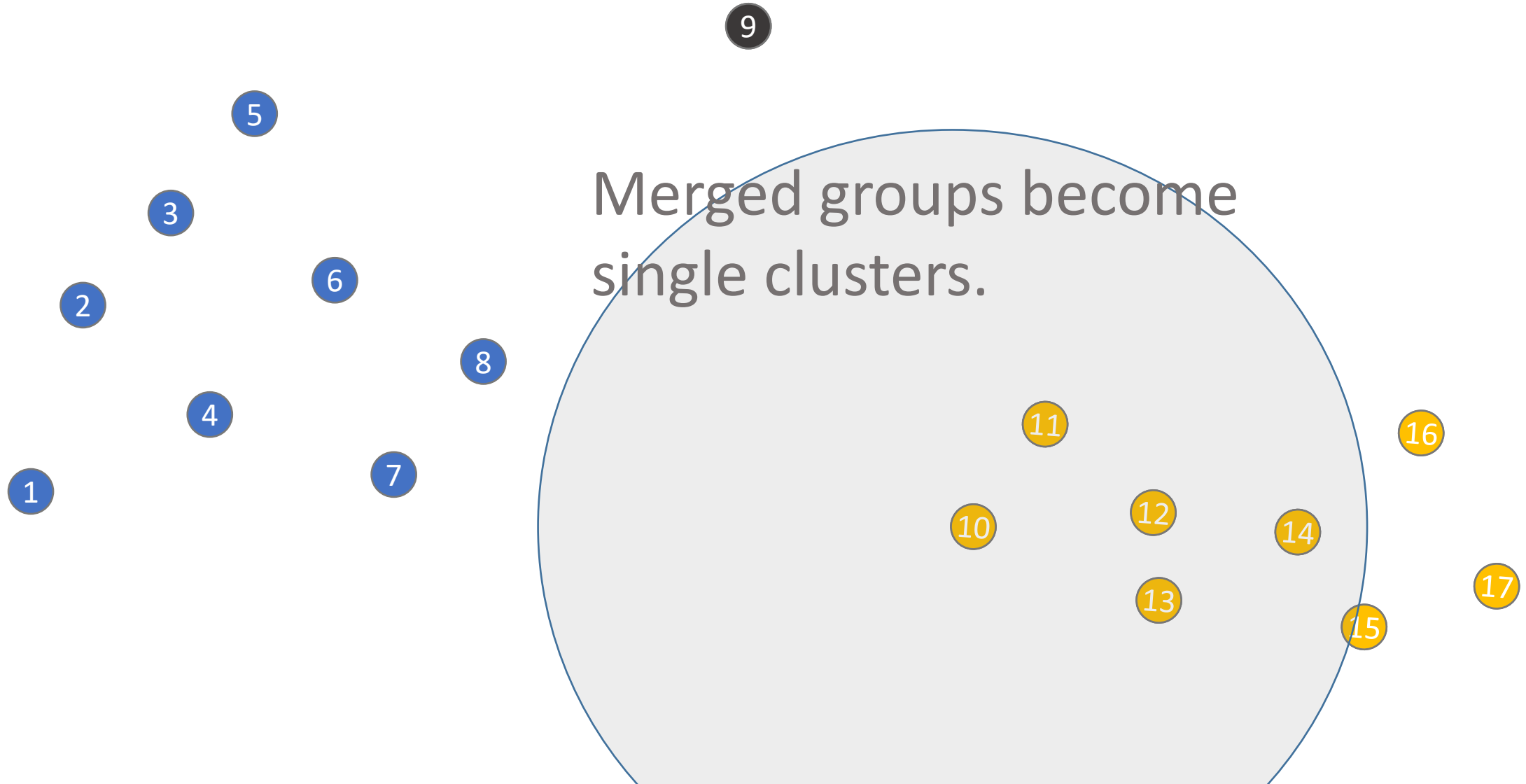
# CLASSIX

Fast and Explainable Clustering

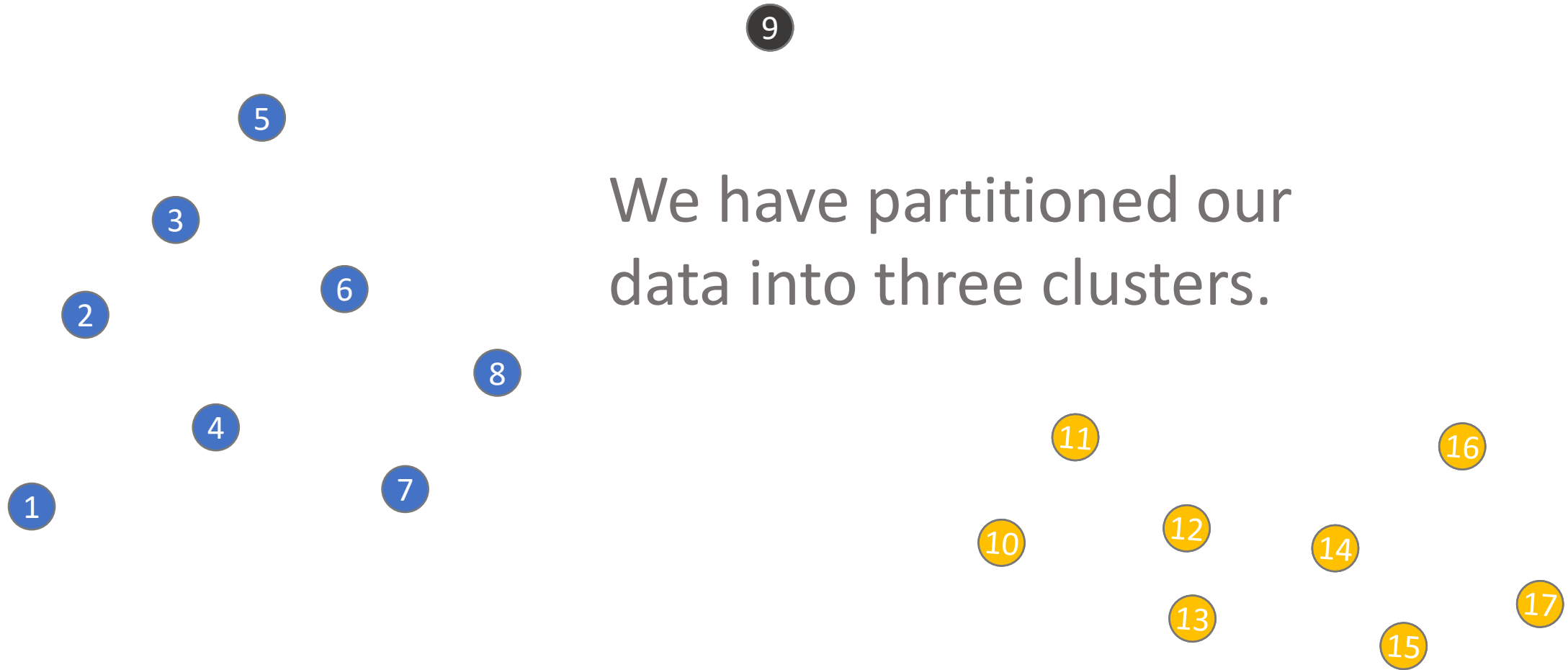


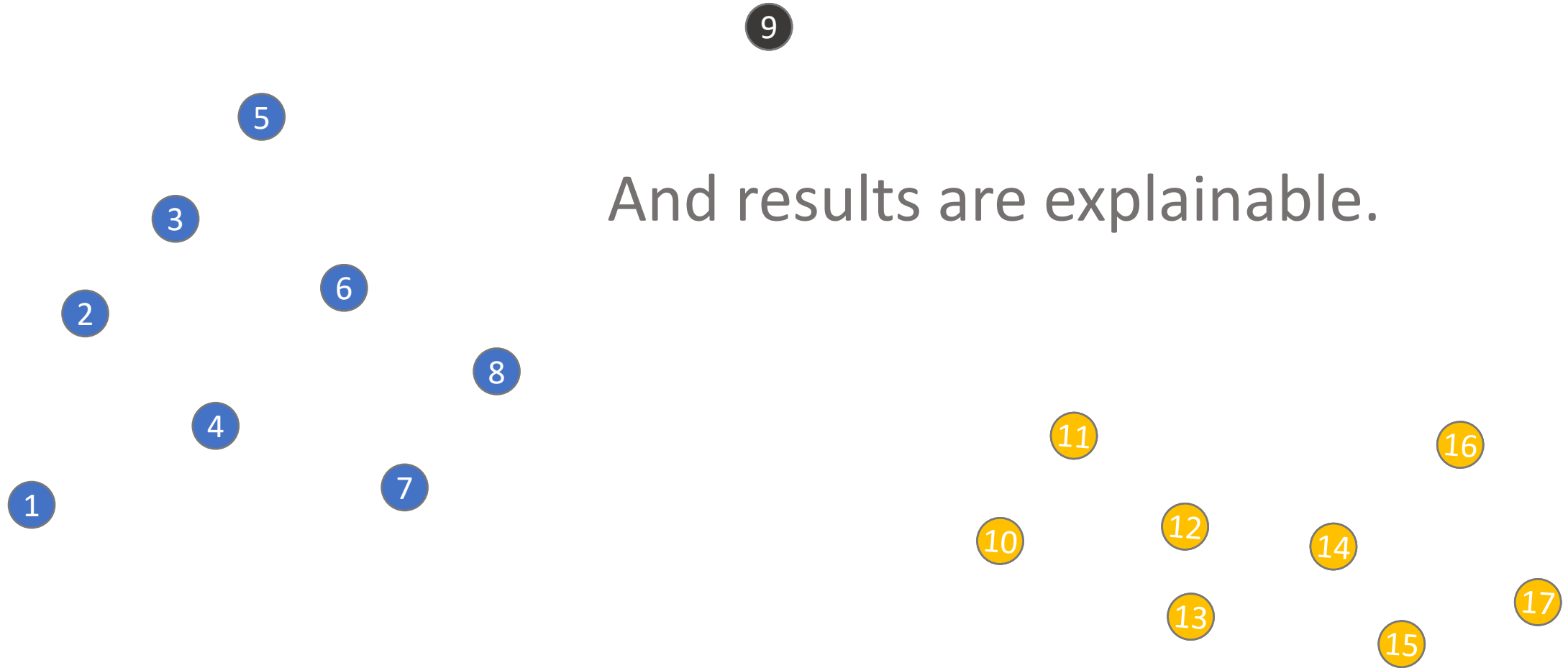
# CLASSIX

Fast and Explainable Clustering



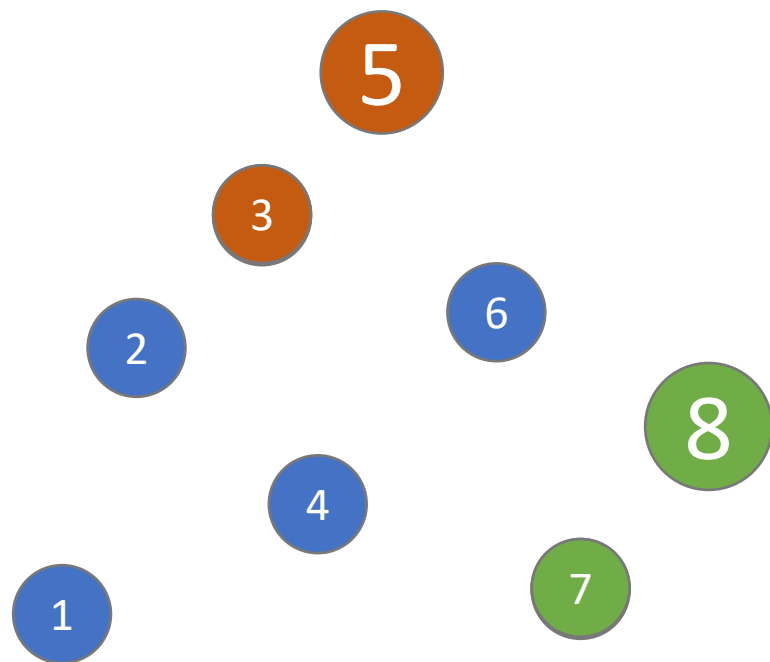






# CLASSIX

Fast and Explainable Clustering



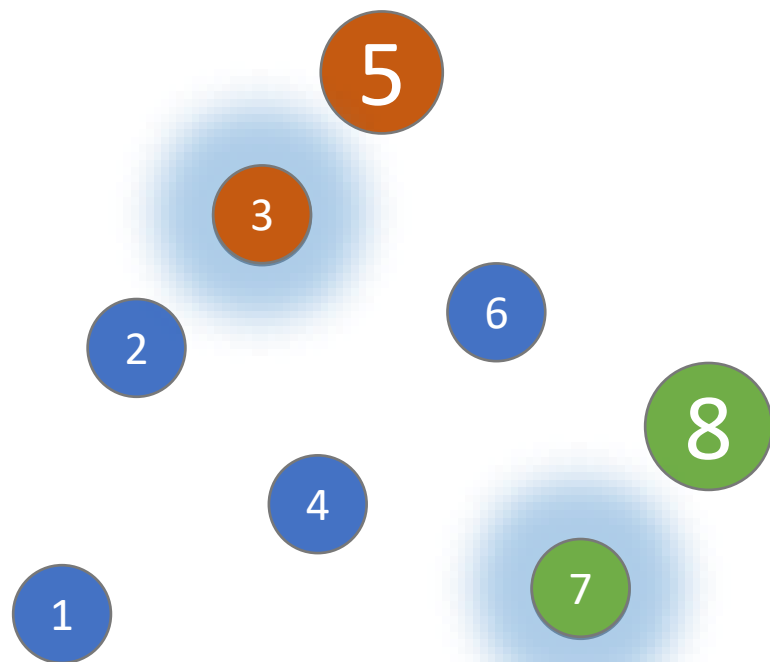
9

Point 5 and 8 are in the same cluster because they are connected by a path of nearby starting points.



# CLASSIX

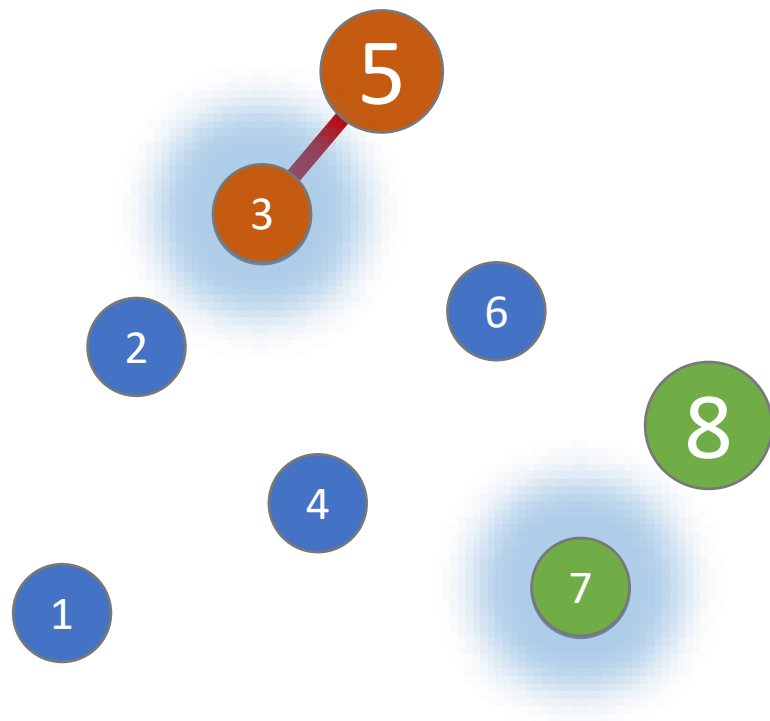
Fast and Explainable Clustering



9

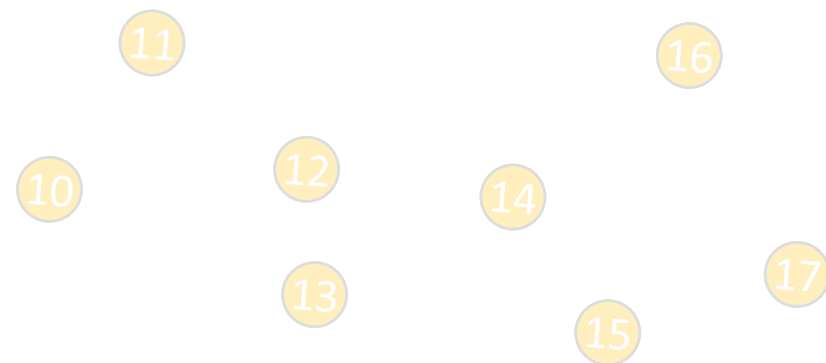
Point 5 and 8 are in the same cluster because they are connected by a path of nearby starting points.





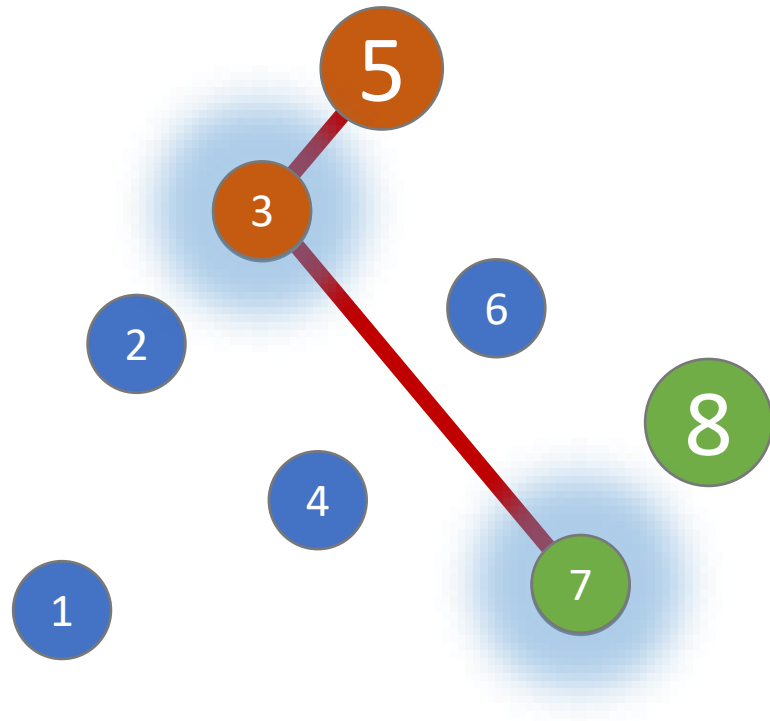
9

Point 5 and 8 are in the same cluster because they are connected by a path of nearby starting points.



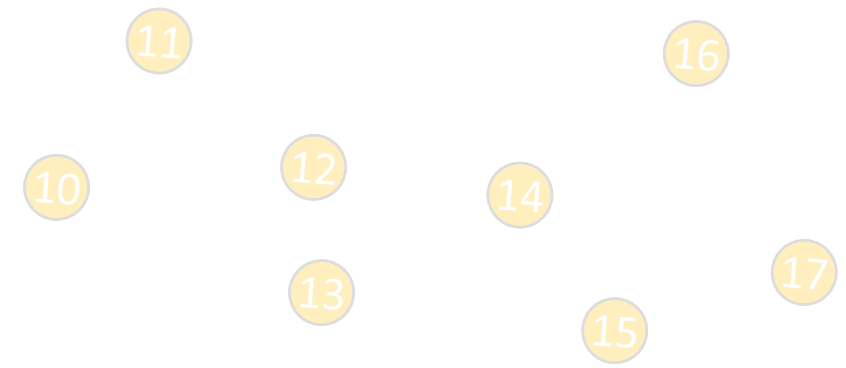
# CLASSIX

Fast and Explainable Clustering



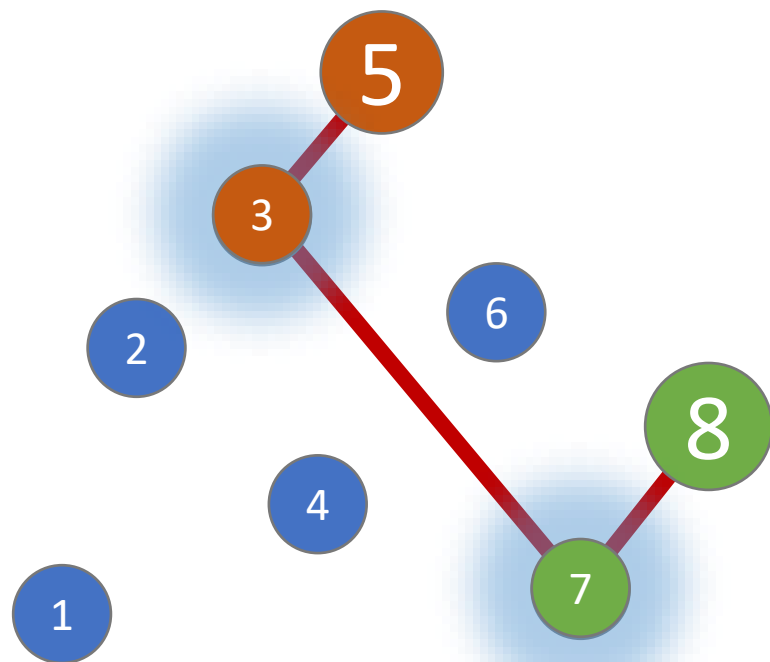
9

Point 5 and 8 are in the same cluster because they are connected by a path of nearby starting points.



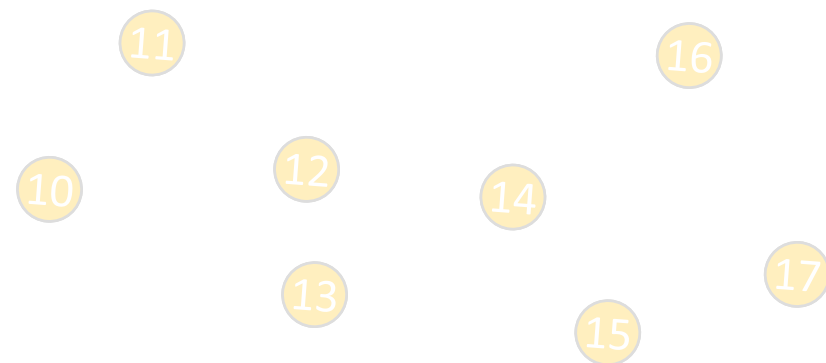
# CLASSIX

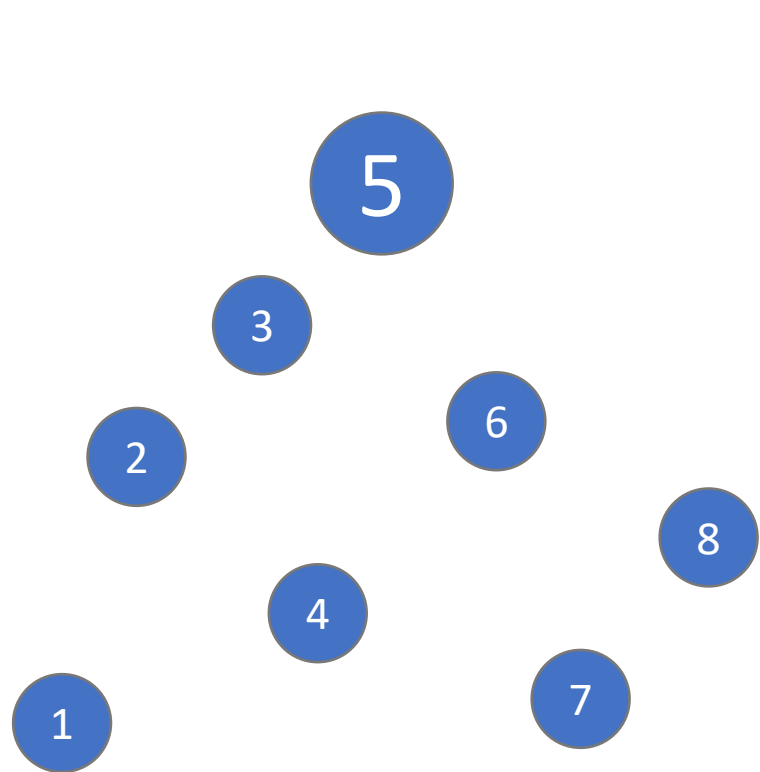
Fast and Explainable Clustering



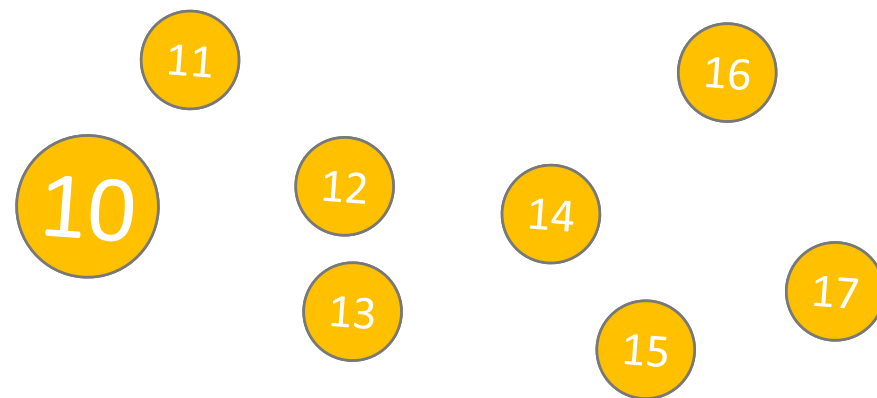
9

Point 5 and 8 are in the same cluster because they are connected by a path of nearby starting points.





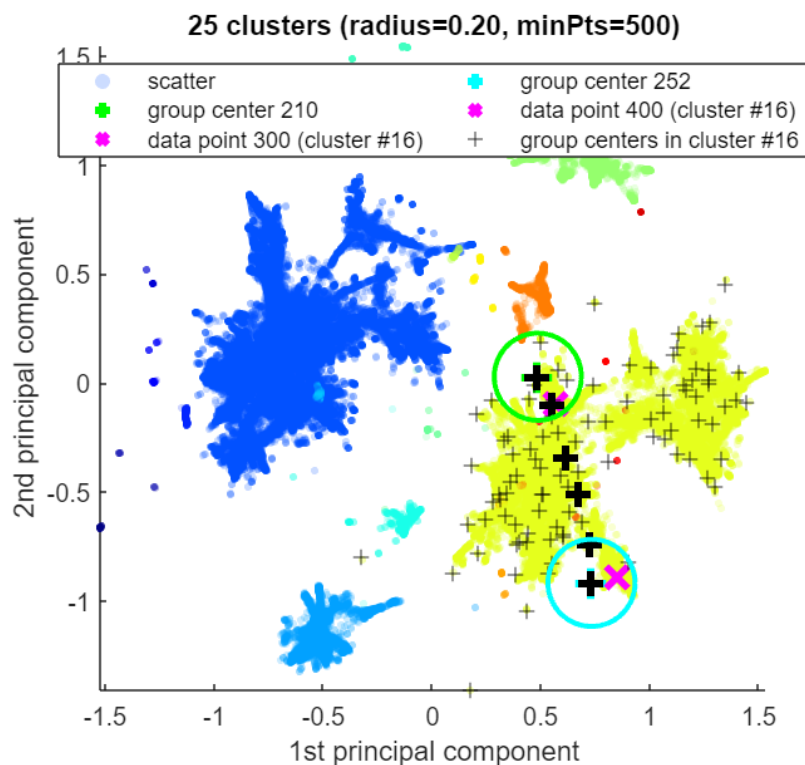
And there is no such path between points 5 and 10, hence different clusters.





## MATLAB demo

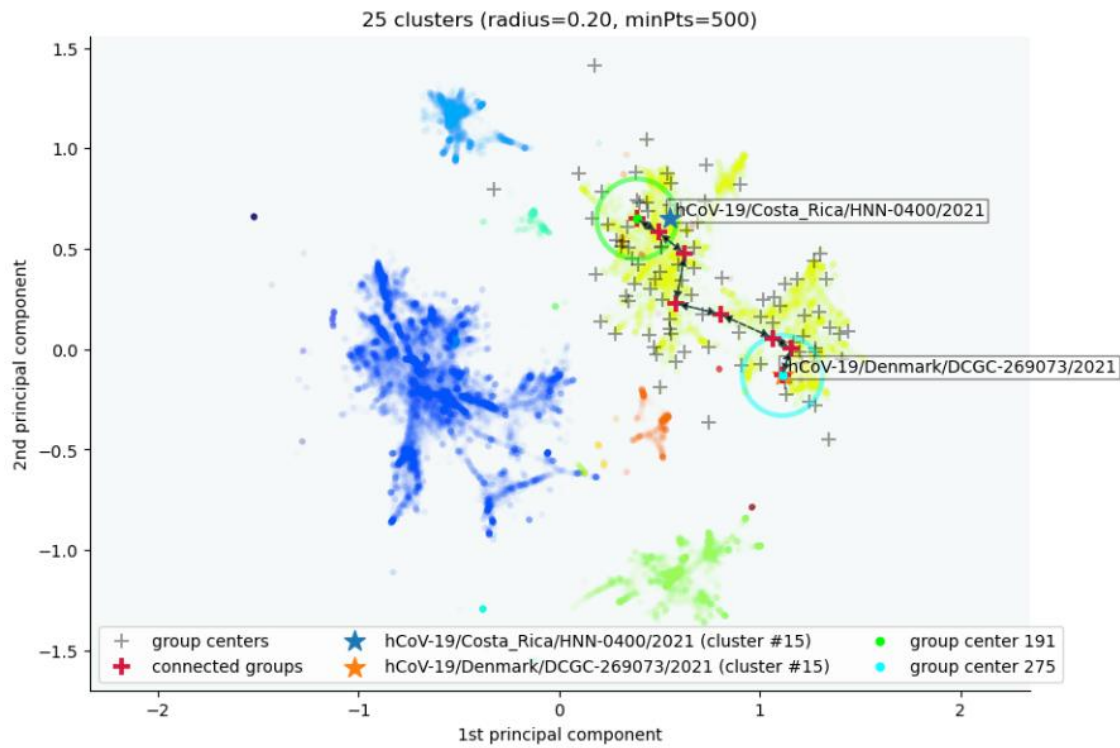
### Clustering 5.7M RNA sequences of coronavirus



	CLASSIX	DBSCAN (5%)
Parameters	(0.2, 500)	(0.1, 1)
Runtime (s)	6.7	493
Clusters	25	39
AMI	0.61	0.60

## Same in Python

### Clustering 5.7M RNA sequences of coronavirus



Data point hCoV-19/Costa\_Rica/HNN-0400/2021 is in group 191.  
Data point hCoV-19/Denmark/DCGC-269073/2021 is in group 275.  
Both groups were merged into cluster #15.

The two groups are connected via groups  
191 <-> 210 <-> 239 <-> 230 <-> 258 <-> 272 <-> 279 <-> 275.

	CLASSIX.py	HDBSCAN
Parameters	(0.2, 500)	(180000, 5)
Runtime (s)	7.7	4080
Clusters	25	4
AMI	0.61	0.59

## MATLAB vs Python

Timing comparison on URI machine learning datasets

	Dim	Size	#Classes	CLASSIX.m	CLASSIX.mex	CLASSIX.py
Banknote	4	1372	2	0.044	0.031	0.078
Dermatology	34	366	6	0.019	0.017	0.047
Ecoli	7	336	7	0.015	0.012	0.028
Glass	9	214	26	0.009	0.009	0.013
Iris	4	150	4	0.008	0.007	0.010
Seeds	7	210	3	0.013	0.010	0.029
Wine	13	178	2	0.011	0.010	0.020
Phoneme	256	4509	4	20.861	6.195	5.369
VDU Signals	2	2028780	11	1.028	1.034	2.860

## How did we get there?

A tale of developing the same algorithm  
in two languages simultaneously

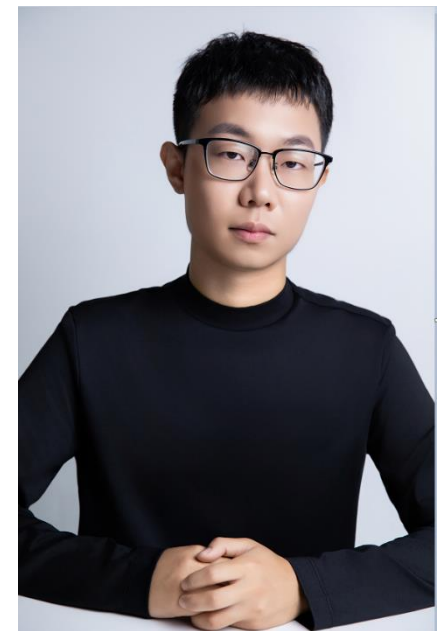
# Step 0: Original version of CLASSIX

Developed in Python with PhD student Xinye Chen between 2021–2022

First arXiv preprint in February 2022 and GitHub release

Used inefficient disjoint set data structure for keeping track of clusters

Still significantly faster than e.g. DBSCAN

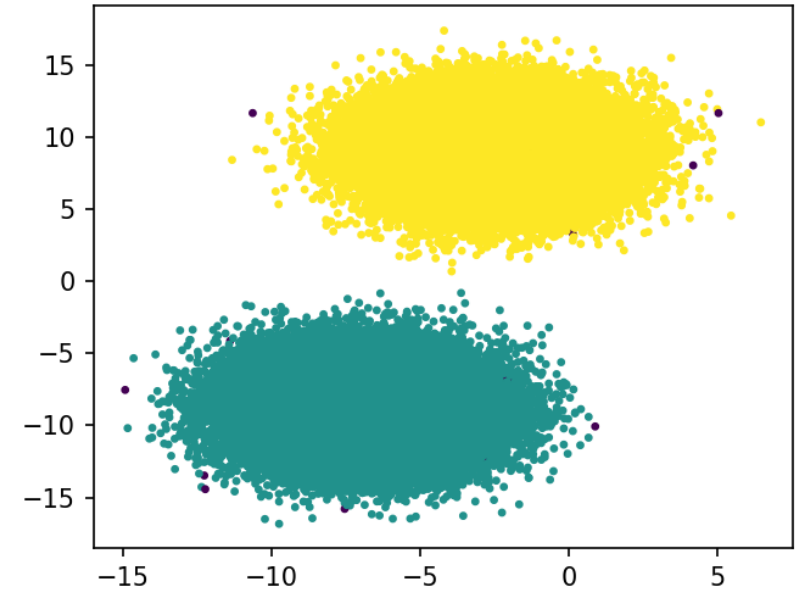


# Step 0: Original version of CLASSIX

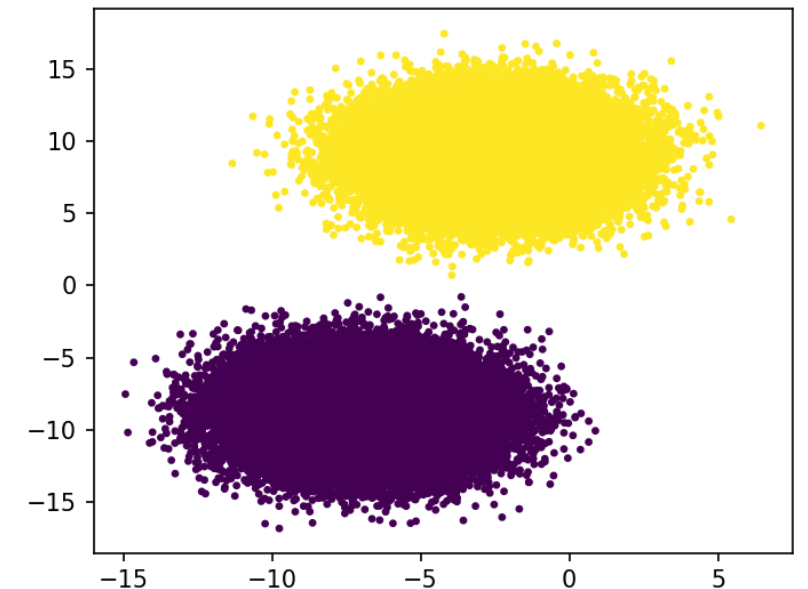
```
print('DBSCAN (sklearn 1.5.2)')
st = time()
clustering = DBSCAN(eps=3, min_samples=5).fit(X)
print(' Runtime:', time()-st, 'seconds')
print(' ARI:      ', ari(clustering.labels_, y))
```

```
print('CLASSIX.py version 0.8.8')
st = time()
clx = CLASSIX(radius=0.2, minPts=5)
clx.fit(X)
print(' Runtime:', time()-st, 'seconds')
print(' ARI:      ', ari(clx.labels_, y))
```

DBSCAN Runtime: 28.344 seconds, ARI: 1.000

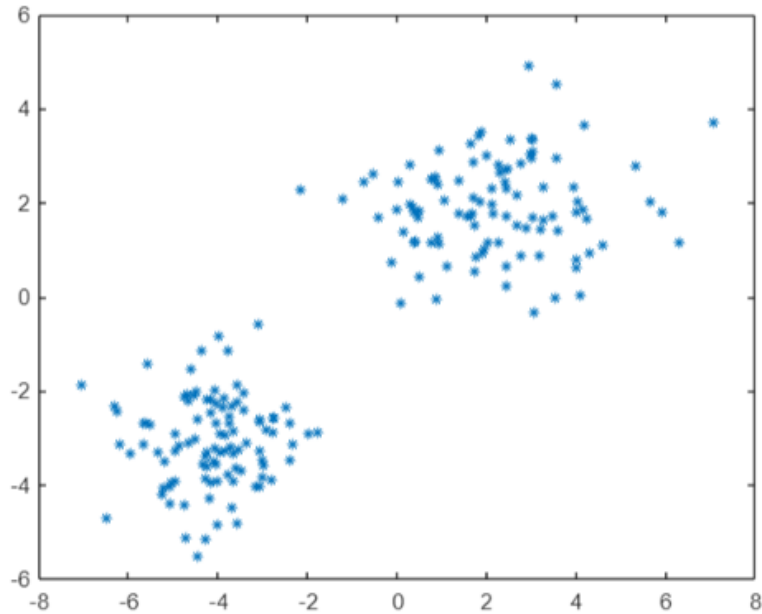


CLASSIX 0.8.8 Runtime: 3.739 seconds, ARI: 1.000



# Step 1: Using interoperability to run CLASSIX.py in MATLAB

```
rng('default') % For reproducibility
mu1 = [2 2]; % Mean of the 1st cluster
sigma1 = [2 0; 0 1]; % Covariance of the 1st cluster
mu2 = [-4 -3]; % Mean of the 2nd cluster
sigma2 = [1 0; 0 1]; % Covariance of the 2nd cluster
r1 = mvnrnd(mu1,sigma1,100);
r2 = mvnrnd(mu2,sigma2,100);
X = [r1; r2];
plot(X(:,1),X(:,2),"*",MarkerSize=5);
```



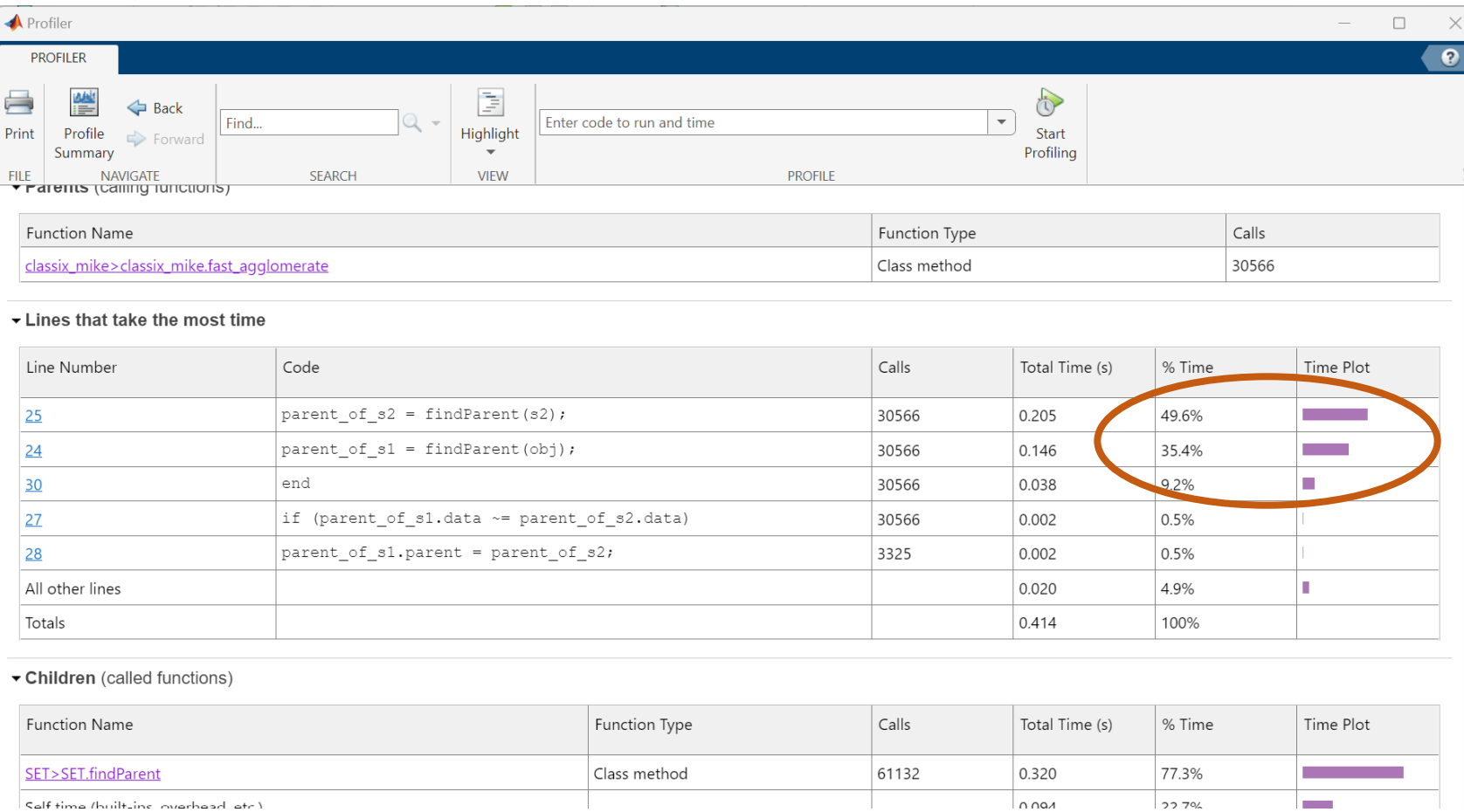
Calling CLASSIX is straightforward. We don't even need to convert the MATLAB array X to a Numpy array as it's all done automatically.

```
clx = py.classix.CLASSIX(radius=0.3, verbose=0);
clx = clx.fit(X);
clx.explain(plot=false);
```

```
CLASSIX clustered 200 data points with 2 features.
The radius parameter was set to 0.30 and minpts was set to 0.
```

[MATLAB Meets Python: Amplifying Research Impact with Cross-Platform Integration - MATLAB \(mathworks.com\)](https://mathworks.com)

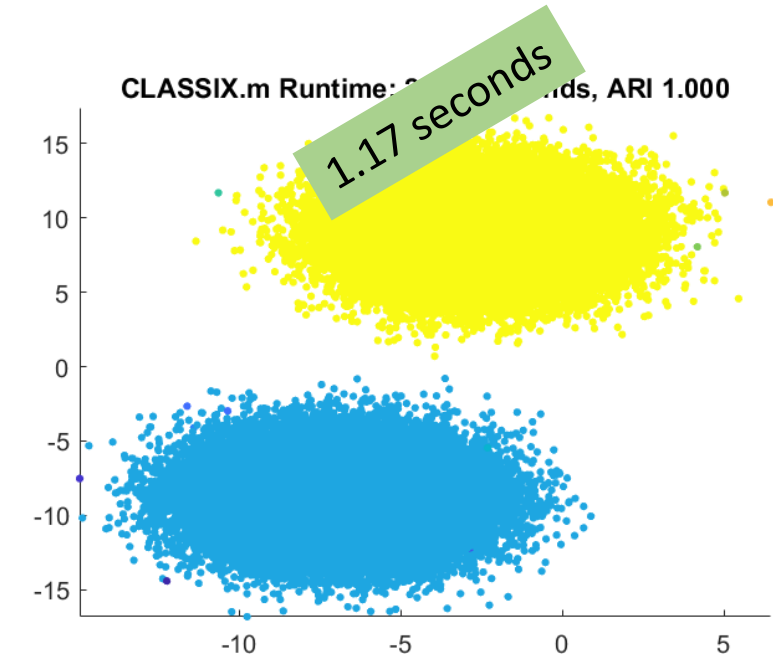
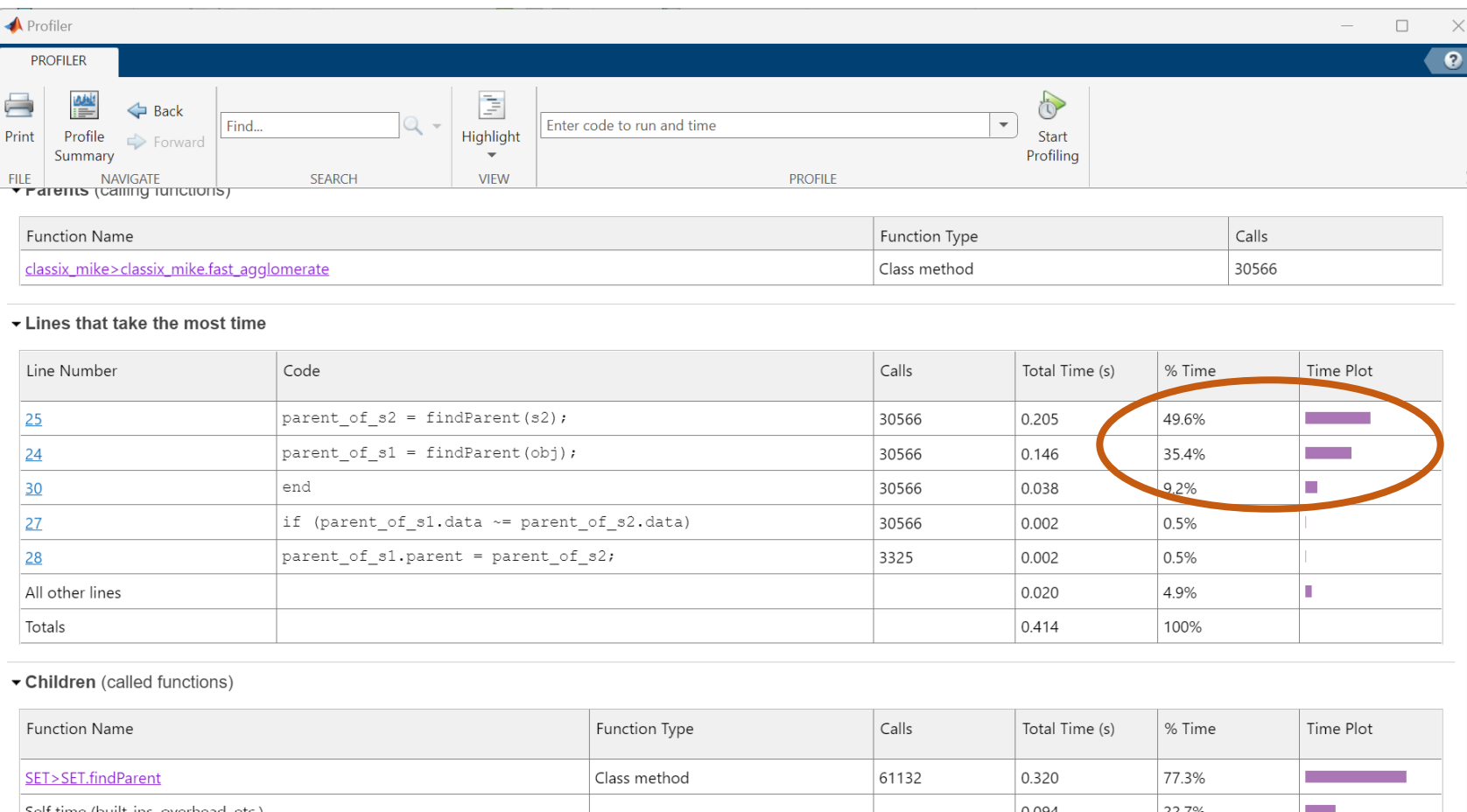
# Step 2: Write native MATLAB version based on Python original



Large amount of time spent on disjoint set structure operations!



# Step 2: Write native MATLAB version based on Python original



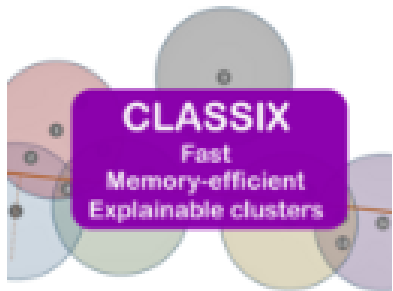
Large amount of time spent on disjoint set structure operations!

MATLAB profiler helped us to improve both MATLAB and Python versions

# Step 3: Use insights to rewrite both Python and MATLAB versions

We used what we learned in Python to improve Mike's MATLAB code and got another significant speed-up.

The MATLAB version is currently faster than the Python package



## Fast and explainable clustering with CLASSIX

Version 1.3 (23.8 MB) by Stefan Güttel

A fast and easy-to-use clustering method that provides explanations for the computed clusters.

<https://github.com/nla-group/classix-matlab>

[⊕ Follow](#)

CLASSIX.py v0.8.8  
3.74 seconds

CLASSIX.m v0.1  
2.37 seconds

not using minPts

CLASSIX.m v1.0  
1.17 seconds

CLASSIX.py v1.0.0  
3.14 seconds

CLASSIX.m v1.3  
0.85 seconds

using mex for  
submatrix-vector  
product

CLASSIX.py v1.2.5  
0.94 seconds



- Make a Python package available to MATLAB users using interoperability
- Rewrite the Python code as MATLAB code
- Iterate between Python package and MATLAB toolbox, using insights from one to drive improvements in the other. New algorithm is faster in both languages than it would have been otherwise
- BONUS: MATLAB itself is improved a little

<https://github.com/nla-group/classix-matlab>