

Clustering Metrics



Week 07 - Day 04

K-Means

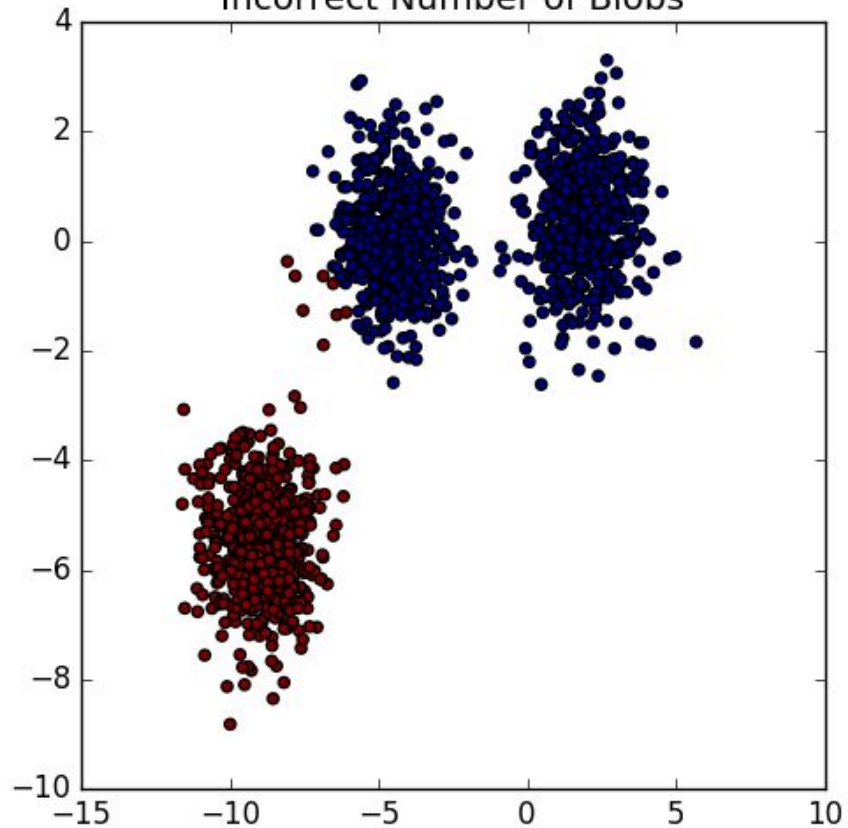
1. Based on distances
2. Uses centroids
3. Iterations
4. K is an input

Problems

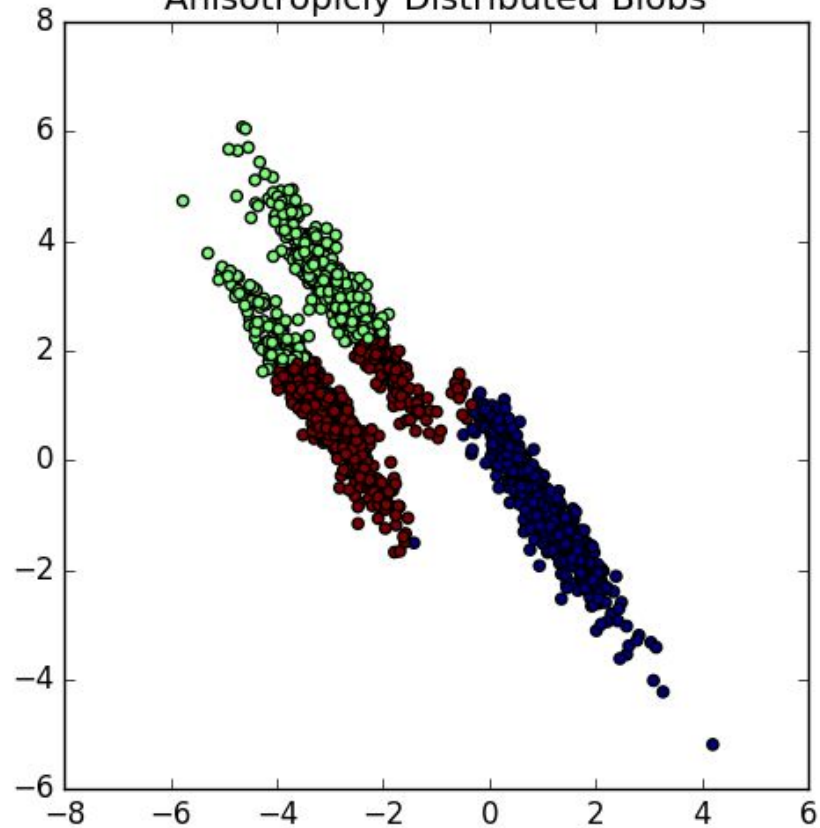
1. K is an input
2. Outliers
3. Irregular shapes

No ground truth

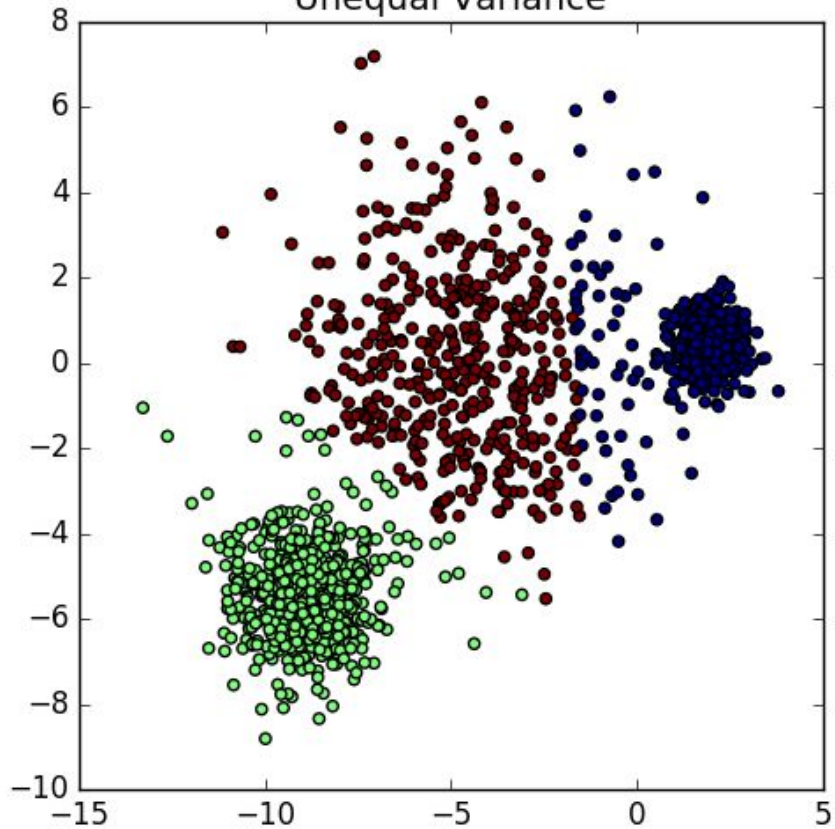
Incorrect Number of Blobs



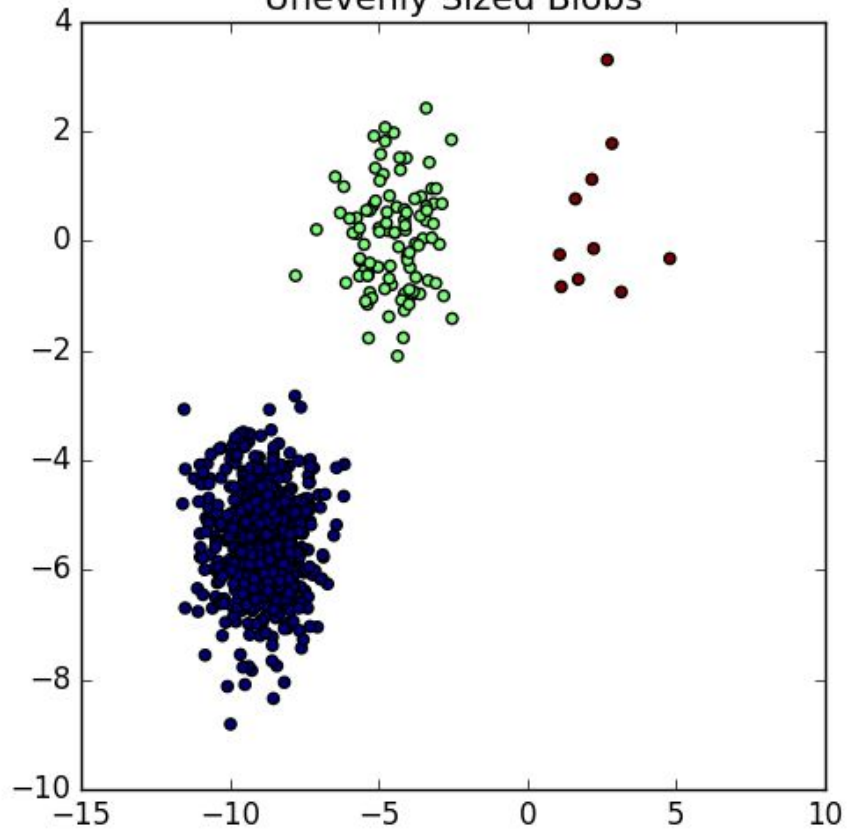
Anisotropically Distributed Blobs



Unequal Variance



Unevenly Sized Blobs



Visual inspection is crucial!

Clustering Metrics

Inertia, Silhouette

"Many indices (more than 30) has been published in the literature for finding the right number of clusters in a dataset."

Inertia

Average squared distance between
each point and its centroid

Similar to MSE

(mean squared error)

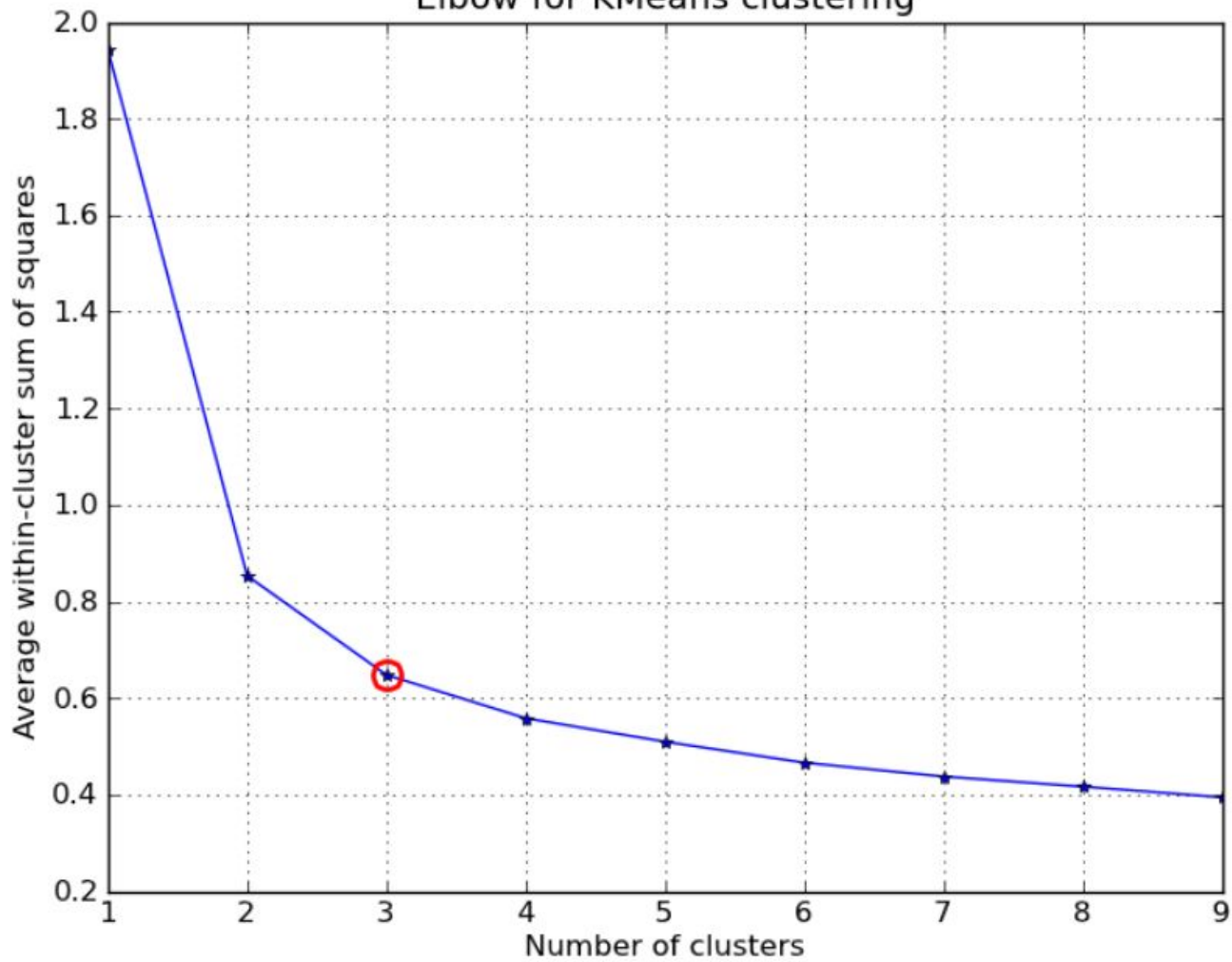
Perfect inertia

=

One cluster for each point

Inertia + elbow

Elbow for KMeans clustering



Elbow

=

Often good-enough

Silhouette Score/Coefficient

Silhouette = Cohesion + Separation

Can be calculated for every point!

Cohesion

=

“Intra-cluster distance”

Cohesion

=

Distance from point to centroid

Separation

=

“Inter-cluster distance”

Separation

=

Distance from point to closer cluster

separation - cohesion

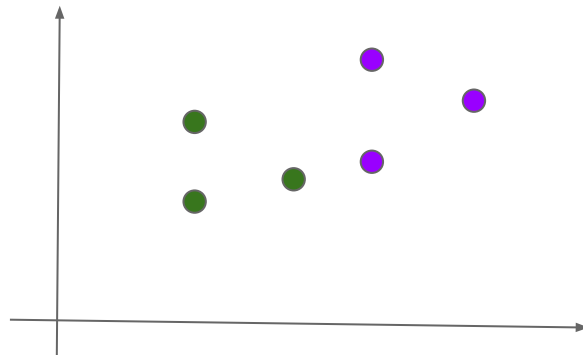
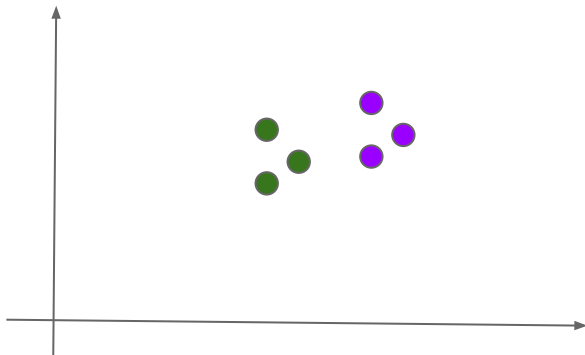
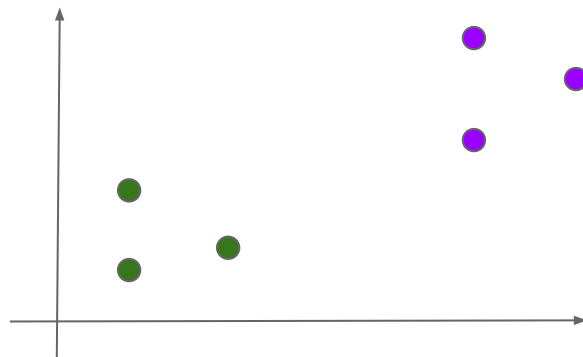
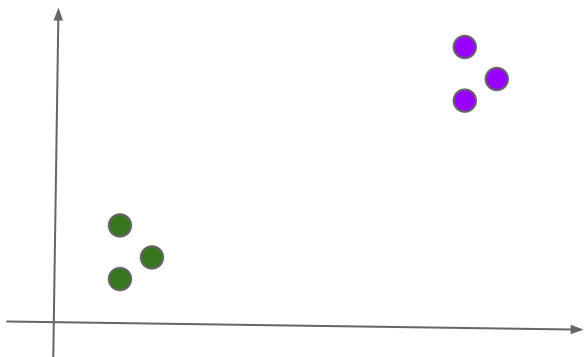
$\max(\text{separation}, \text{cohesion})$

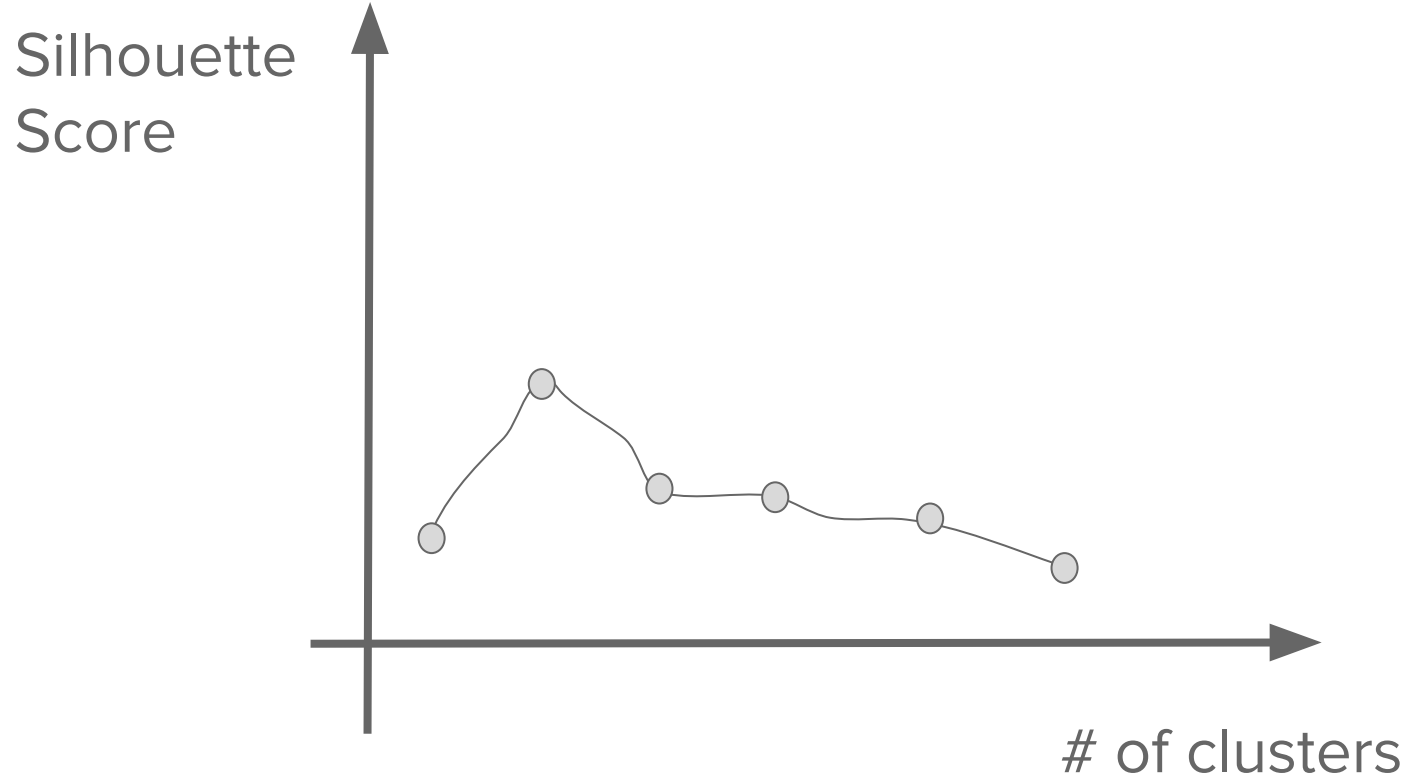
Silhouette Coefficient

=

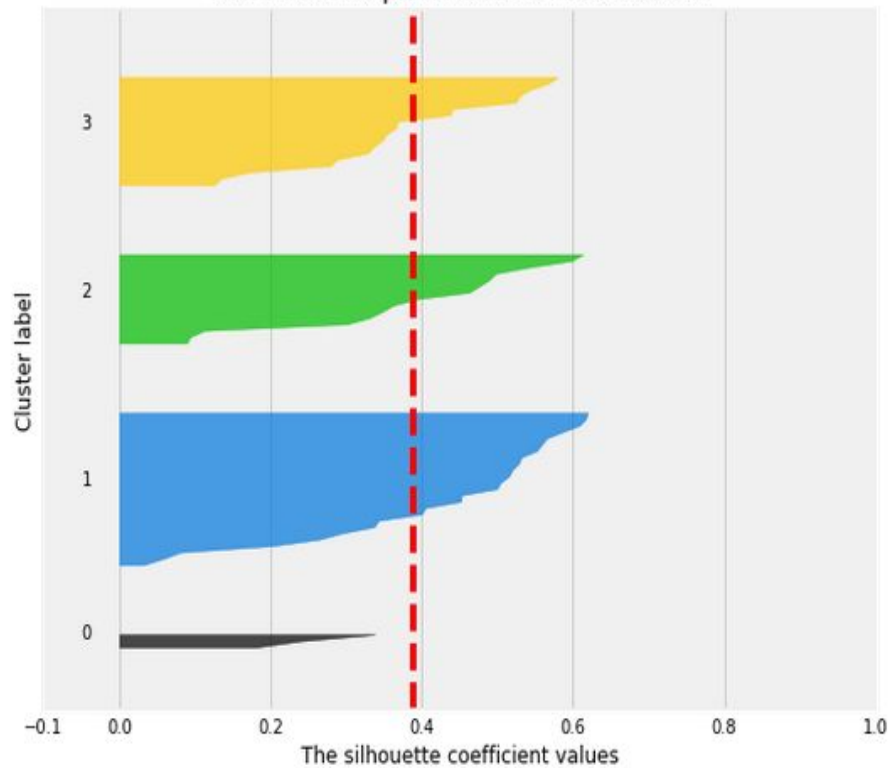
Average Silhouette

(of each point)

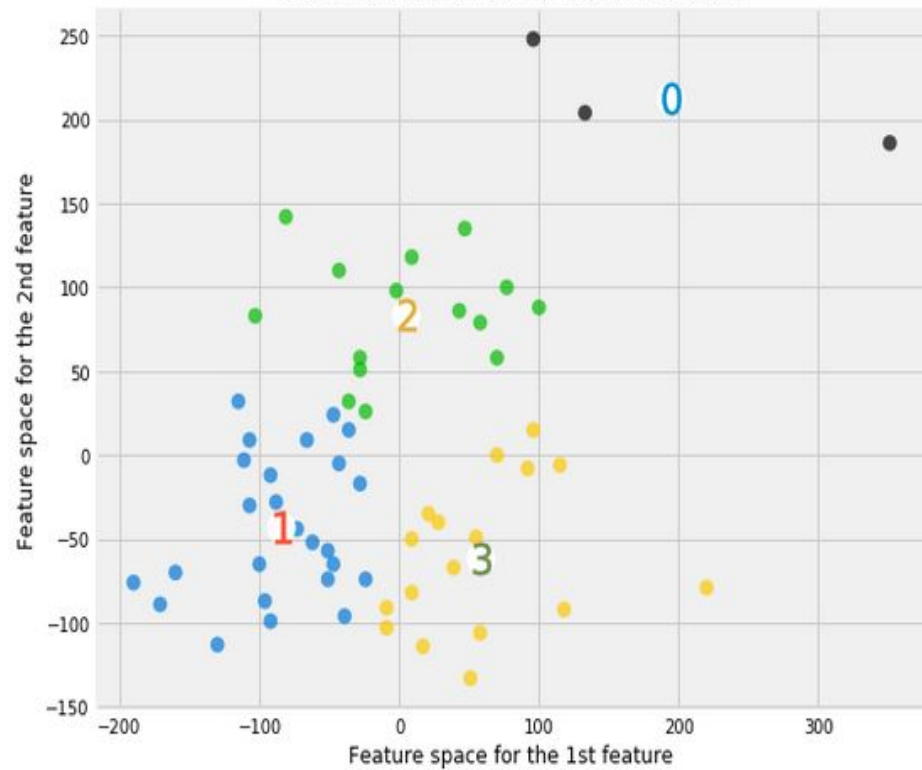




The silhouette plot for the various clusters.



The visualization of the clustered data.



Silhouette Coefficient
is good but not perfect!

Metrics
when you have
the labels

Unusual...why?

Labels = Classification!

Too easy :)

Completeness

Homogeneity

V-Measure Score

Mutual Information Score

Completeness

“Indicates that all members of a given class are assigned to the same cluster.”

[0,1]

Homogeneity

“Indicates each cluster contains only members of a single class.”

$[0,1]$

V-Measure

$$V = \frac{2 \cdot \text{homogeneity} \cdot \text{completeness}}{\text{homogeneity} + \text{completeness}}$$

[0,1]

Similar to what?

F1 score!

Mutual Information Score

$$MI(i, j) = \sum_{a, b} P(a_i, b_j) \cdot \log \left(\frac{P(a_i, b_j)}{P(a_i) \cdot P(b_j)} \right)$$