

Earle Aguilar  
804 501 476  
11/10/2017

## 1. K-means

a.

	Dataset 1	Dataset 2	Dataset 3
Purity	1.0	0.868	0.78
NMI	1.0	0.464	0.170

b. Strengths

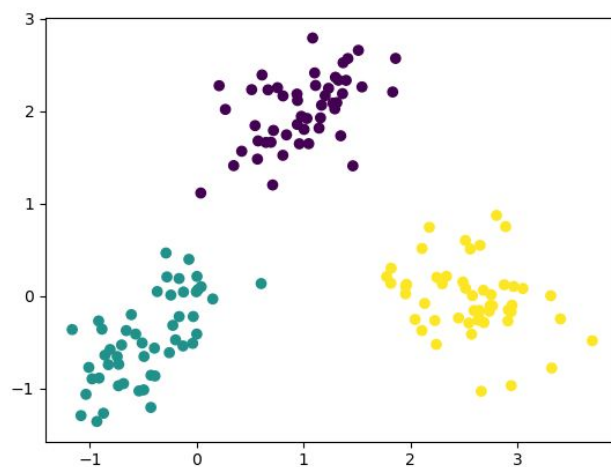
i. Efficient

c. Weaknesses

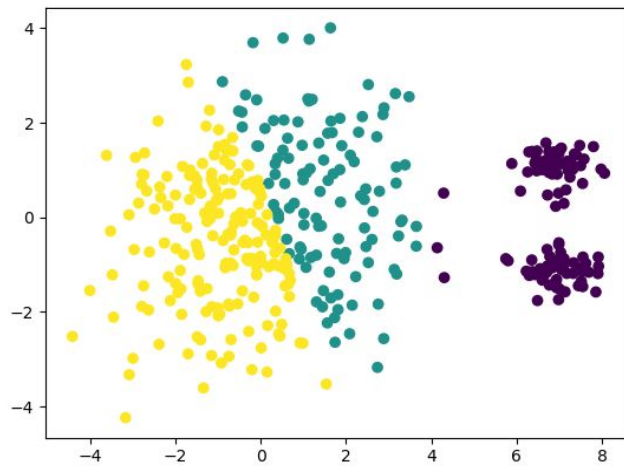
i. Need to specify K in advance.

ii. Sensitive to noisy data and outliers.

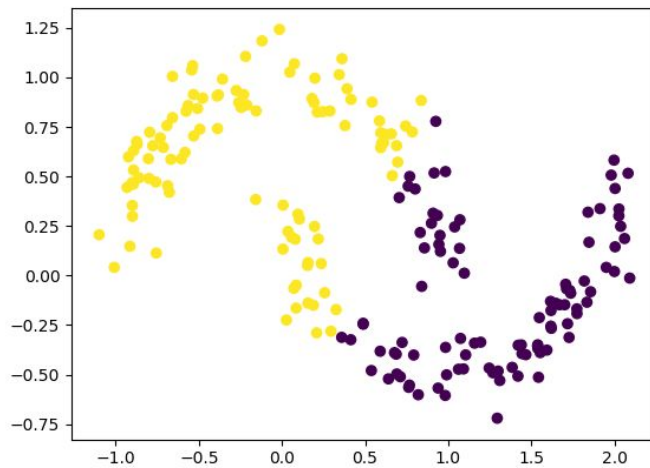
iii. Can terminate at a local optima.



d.



e.



f.

## 2. DBSCAN

a.

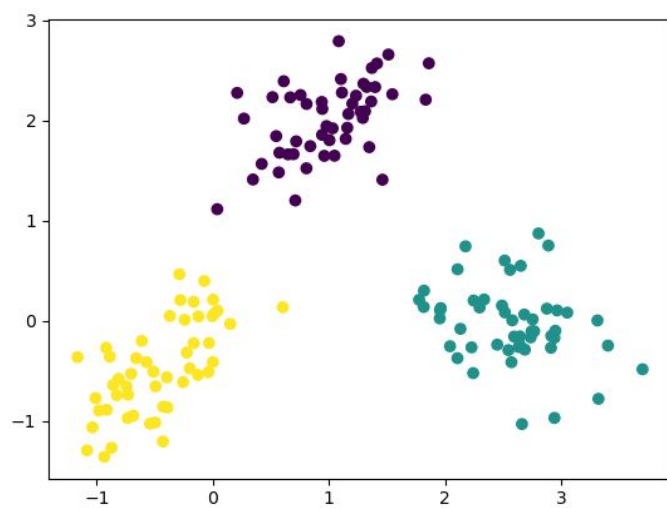
	Dataset 1	Dataset 2	Dataset 3
Purity	1.0	0.953	1.0
NMI	1.0	1.0	1.0

b. Strengths

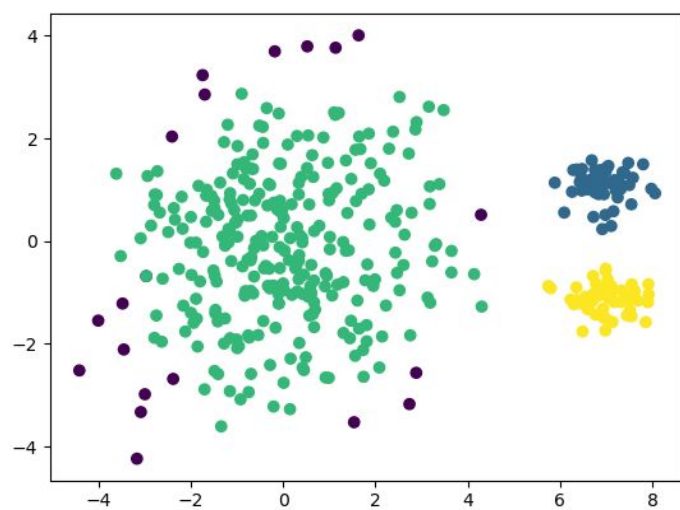
- i. Robust to noise and outliers.
- ii. Can handle clusters of different shapes and sizes.

c. Weaknesses

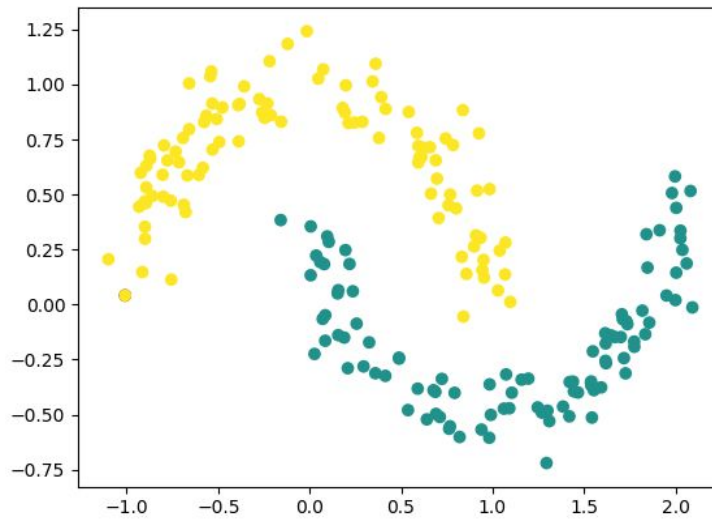
- i. Sensitive to parameters.



d.



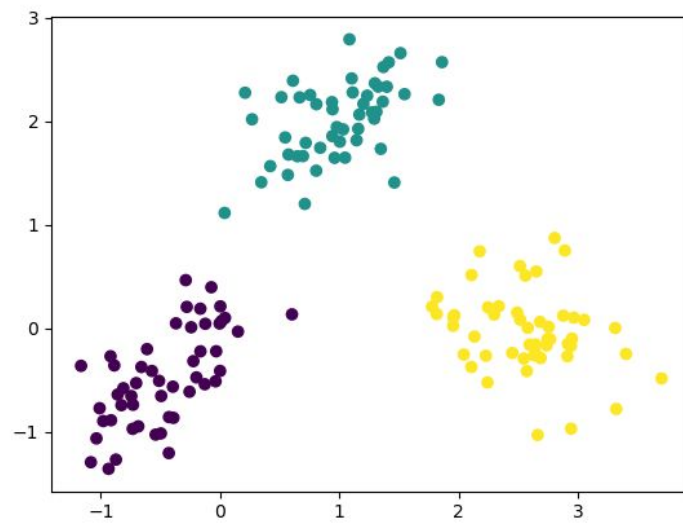
e.



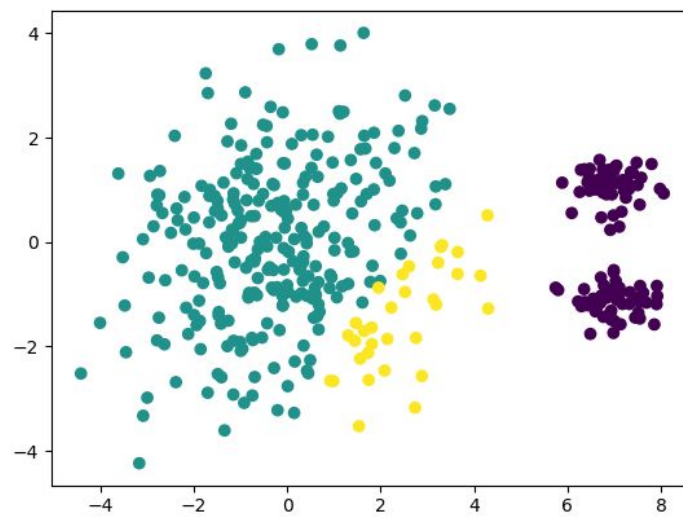
- f.
3. GMM
- a.

	Dataset 1	Dataset 2	Dataset 3
Purity	1.0	0.875	0.69
NMI	1.0	0.608	0.076

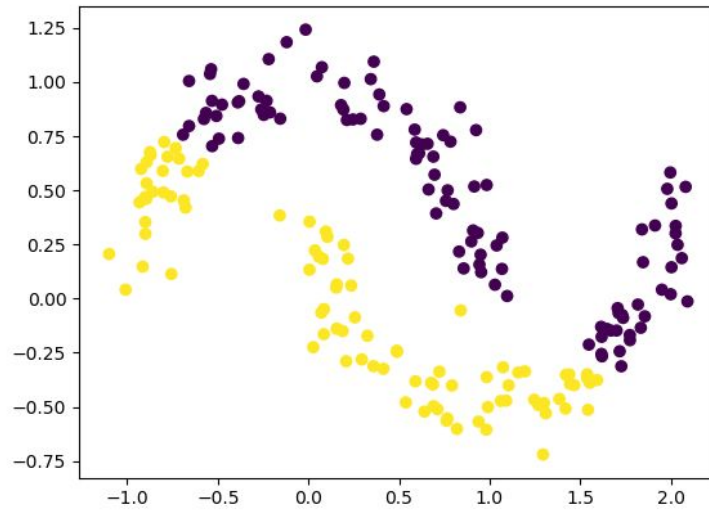
- b. Strengths
- i. Mixture models are more general than partitioning. Can have different densities and size of clusters
  - ii. Clusters may be characterized by a small number of parameters.
- c. Weaknesses
- i. Converges to local optima.
  - ii. Computationally expensive for large distributions
  - iii. Hard to estimate number of clusters
  - iv. Can only deal with spherical clusters.



d.



e.



f.

#### 4. Problem 1

		Predicted				$\text{Purity} = \frac{1}{20}(4+5+5+4)$ $= \frac{18}{20} = 0.9$
Actual		$c_1$	$c_2$	$c_3$	$c_4$	
$c_1$		0	1	4	0	
$c_2$		5	0	0	0	
$c_3$		0	5	0	0	
$c_4$		0	0	1	4	

$\begin{pmatrix} 22 \\ 22 \\ 22 \\ 2 \end{pmatrix}$ $c_1$	$\begin{pmatrix} 33 \\ 33 \\ 33 \\ 1 \end{pmatrix}$ $c_2$	$\begin{pmatrix} 1111 \\ 4 \end{pmatrix}$ $c_3$	$\begin{pmatrix} 44 \\ 44 \end{pmatrix}$ $c_4$	$TP+FP = \left(\frac{5}{2}\right) + \left(\frac{6}{2}\right) + \left(\frac{5}{2}\right) + \left(\frac{4}{2}\right)$ $TP+FP = 41$ $TP = \left(\frac{5}{2}\right) + \left(\frac{5}{2}\right) + \left(\frac{4}{2}\right) + \left(\frac{4}{2}\right)$ $TP = 32$ $FP = 9$
---	---	---	--	--

#pairs = 190

Total Neg = 190 - 41		Same Cluster	Different Cluster
= 149	Same Class	32	9
	Different Class	9	141

$FN = 9$  and  $TN = 149 - 9 = 141$

$TP = 32$

Precision =  $\frac{TP}{TP+FP} = \frac{32}{41}$

Recall =  $\frac{TP}{TP+FN} = \frac{32}{32+9}$

Precision = 0.78

Recall = 0.8

$F = \frac{2(0.78)(0.8)}{0.78+0.8} = 0.79$

a.



	C1	C2	C3	C4	sum
1	0	1	4	0	5
2	5	0	0	0	5
3	0	5	0	0	5
4	0	0	1	4	5
sum	5	6	5	4	20

$$I(\Omega, C) = \sum_{k=1}^4 \sum_{j=1}^4 \frac{|C_{k,j}|}{N} \log \left( \frac{N |C_{k,j}|}{|C_k| |W_j|} \right) =$$

$$= \frac{5}{20} \log(4) + \frac{1}{20} \log\left(\frac{2}{3}\right) + \frac{5}{20} \log\left(\frac{10}{3}\right) + \frac{4}{20} \log\left(\frac{10}{5}\right) + \frac{1}{20} \log\left(\frac{4}{5}\right) + \frac{4}{20} \log(4)$$

$$I(\Omega, C) = 1.623$$

$$H(\Omega) = 2, H(C) = 1.985$$

$$NMI(\Omega, C) = \frac{1.623}{\sqrt{2 \times 1.985}} = 0.8143$$

b.