

Extra Material for P1W3D3

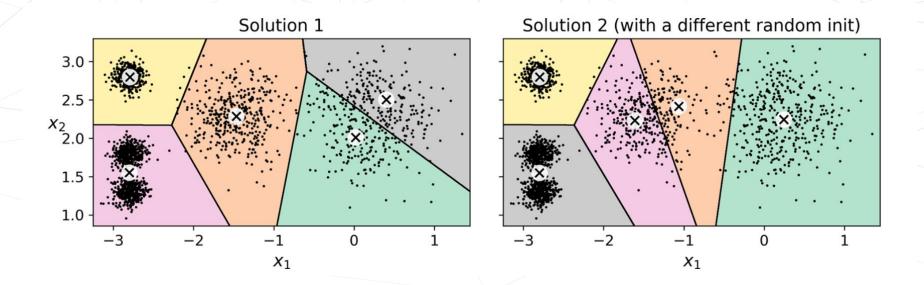
P1 W3 D3 AM GAUSSIAN MIXTURE MODEL

KMeans - Advantages

- Guaranteed to converge in a finite number of steps (usually quite small iterations).
- One of the fastest clustering algorithm.
- Easy to understand.

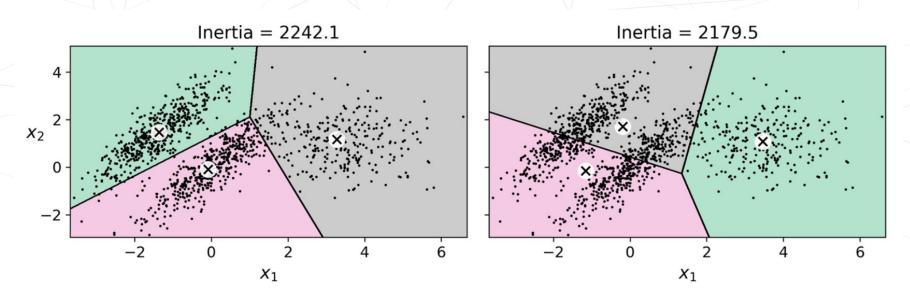
KMeans - Disadvantages (1)

Depends on the centroid initialization.



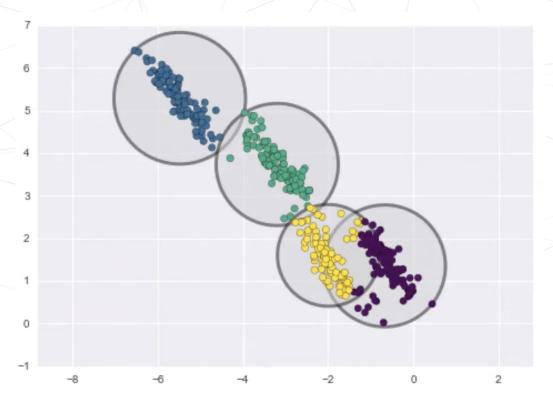
KMeans - Disadvantages (2)

 KMeans does not behave very well when the clusters have varying sizes, different densities, or non spherical shapes.



KMeans - Disadvantages (3)

KMeans is a hard-clustering technique.



Problem : News Clustering (1)

- Cluster 1 : news related to Sport.
- Cluster 2 : news related to Technology.





Problem : News Clustering (2)

Which cluster for this news?

FIFA Bakal Uji Teknologi Offside Semi-Otomatis di Piala Arab 2021

jun | CNN Indonesia





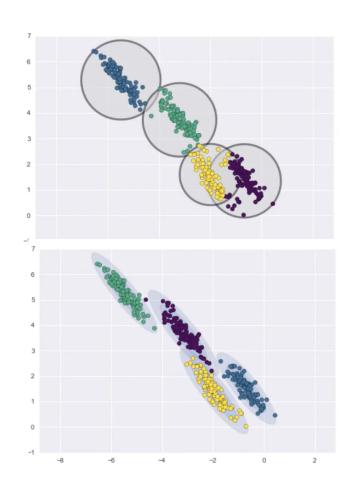




KMeans vs GMM

- KMeans:
 - → Good for spherical-shape
 - → Hard clustering

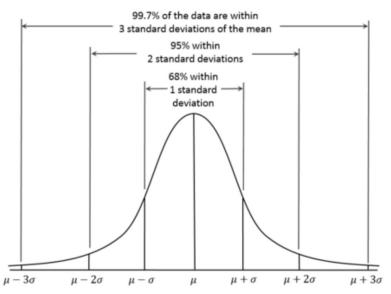
- GMM:
 - → Good for non spherical-shape
 - → Soft clustering

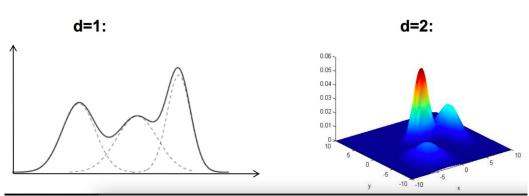


Gaussian Mixture Model (1)

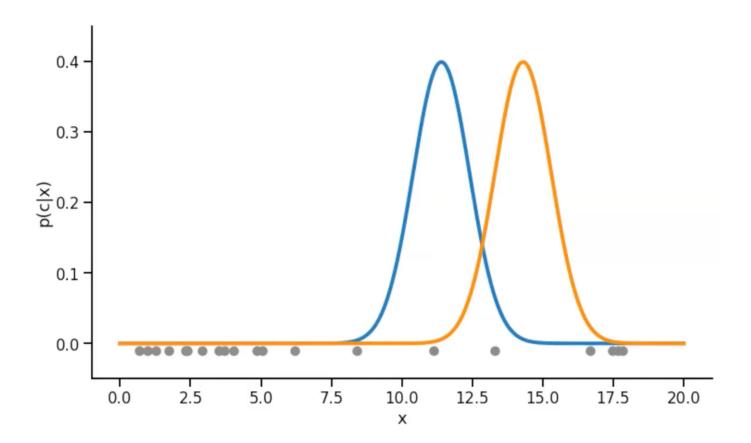
- Using probabilistic approaches.
- Based on Expectation-Maximisation algorithm.
- GMM is a model that assumes data generated from a mixture of K Gaussian (bell shape) components.

Gaussian Mixture Model (2)

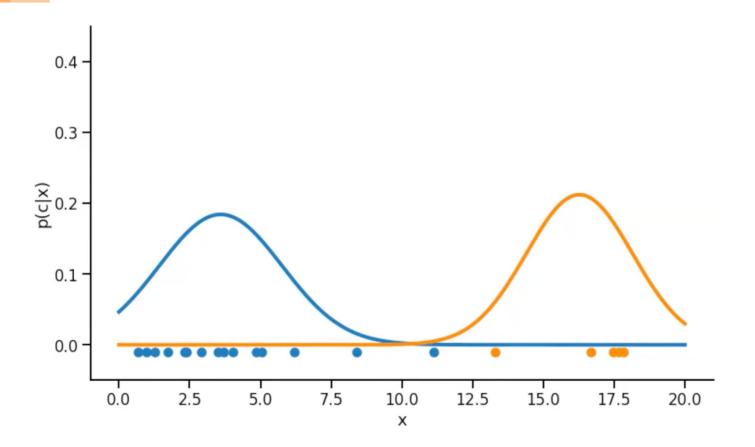




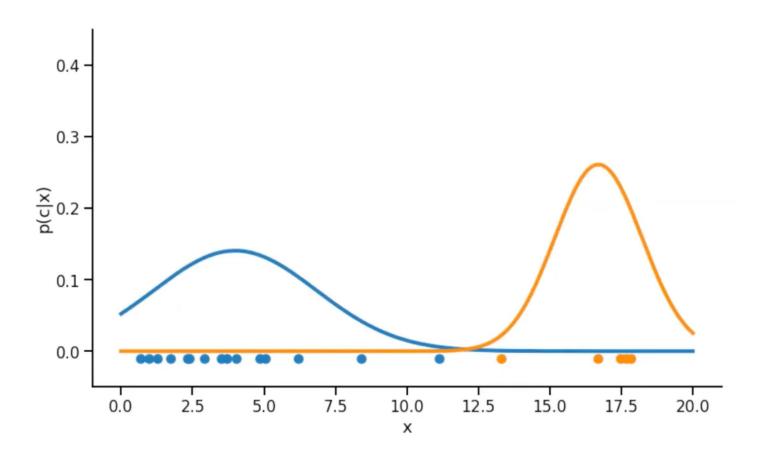
Gaussian Mixture Model - Illustration (Iteration 1)



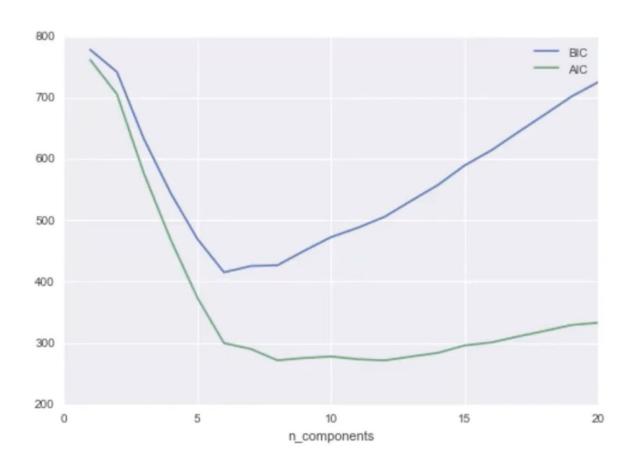
Gaussian Mixture Model - Illustration (Iteration 2)



Gaussian Mixture Model - Illustration (Iteration 3)



BIC & AIC



Algorithm Comparison

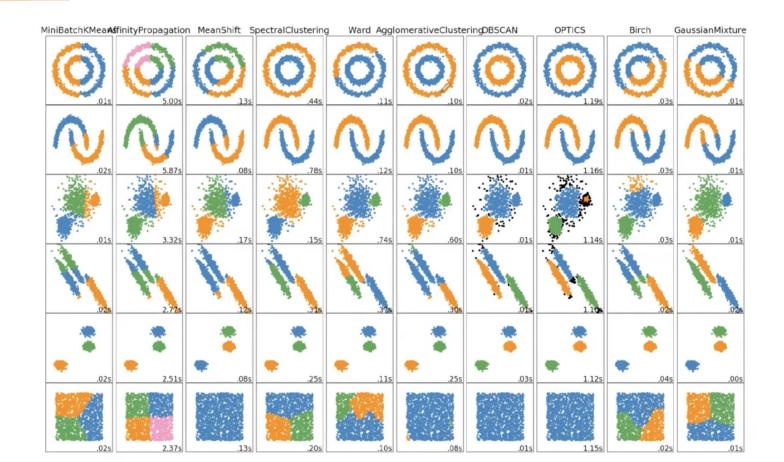




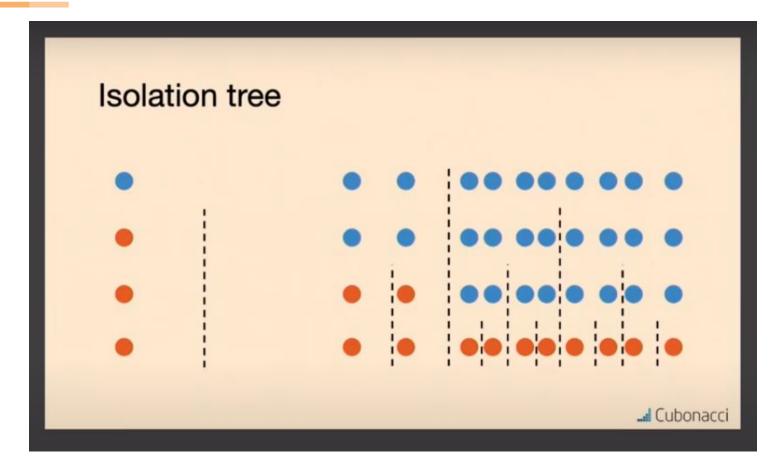
Illustration (Outlier)



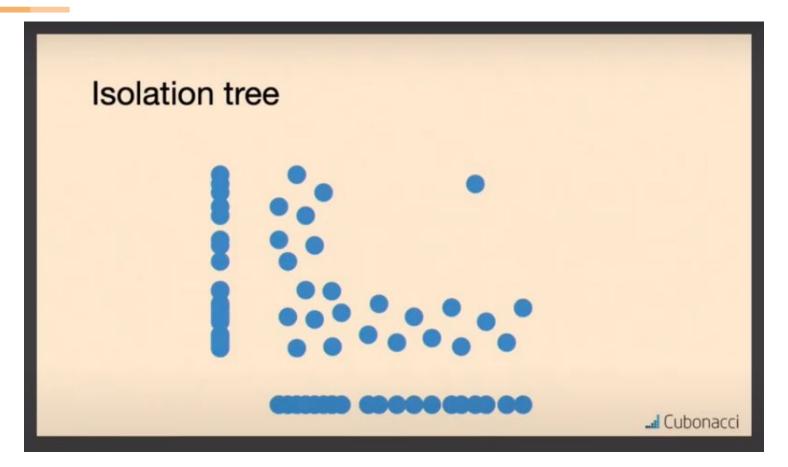
Illustration (Novelty)



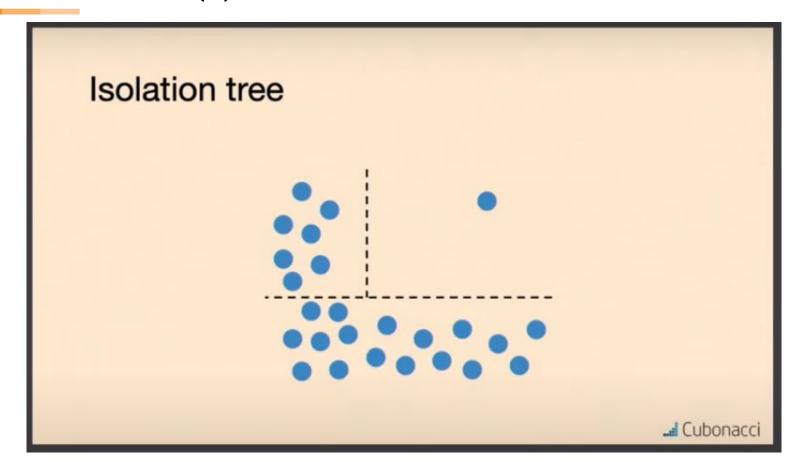
Isolation Forest (1)



Isolation Forest (2)



Isolation Forest (3)



Isolation Forest (4)

