

HW #4:

1.

Calculating Purity:

Output Label 1: 5 2's

Output Label 2: 5 3's 1 1

Output Label 3: 3 1's, 1 4, 1 1

Output Label 4: 4 4's.

$$\text{Purity} = (5+5+3+4) / 20 = .85$$

Precision, Recall, F-Measure

$$\text{TP: } 4+3+2+1+3+2+1+3+2+1+4+3+2+1 = 32$$

$$\text{FP: } 4 + 4 = 8$$

$$\text{FN: } 5+4=9$$

$$\text{Precision} = \text{TP}/(\text{TP}+\text{FP}) = 32/40 = .8$$

$$\text{Recall} = \text{TP}/(\text{TP}+\text{FN}) = 32/41 = .78$$

$$\text{F-measure} = 2 * \text{Precision} * \text{Recall} / (\text{Precision} + \text{Recall}) = .78987$$

NMI:

$$P(\text{cluster } 1) = 1/4$$

$$P(\text{cluster } 2) = 1/4$$

$$P(\text{cluster } 3) = 1/4$$

$$P(\text{cluster } 4) = 1/4$$

Used wolfram alpha to calculate H

$$H(\text{Clusters}) = 2$$

$$P(\text{Class } 1) = 1/4$$

$$P(\text{Class } 2) = 3/10$$

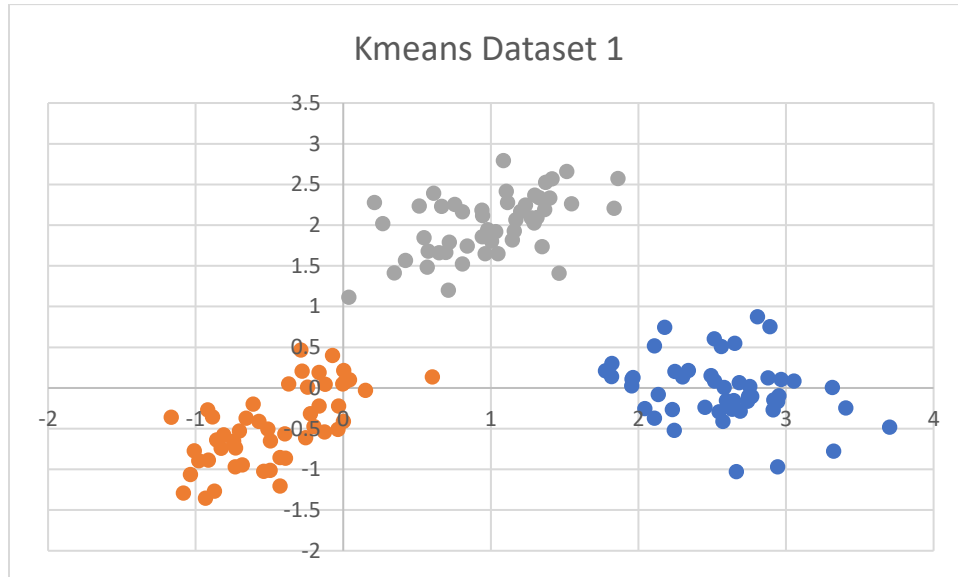
$$P(\text{Class } 3) = 1/4$$

$$P(\text{Class } 4) = 1/5$$

$$H(\text{Classes}) = 1.9855$$

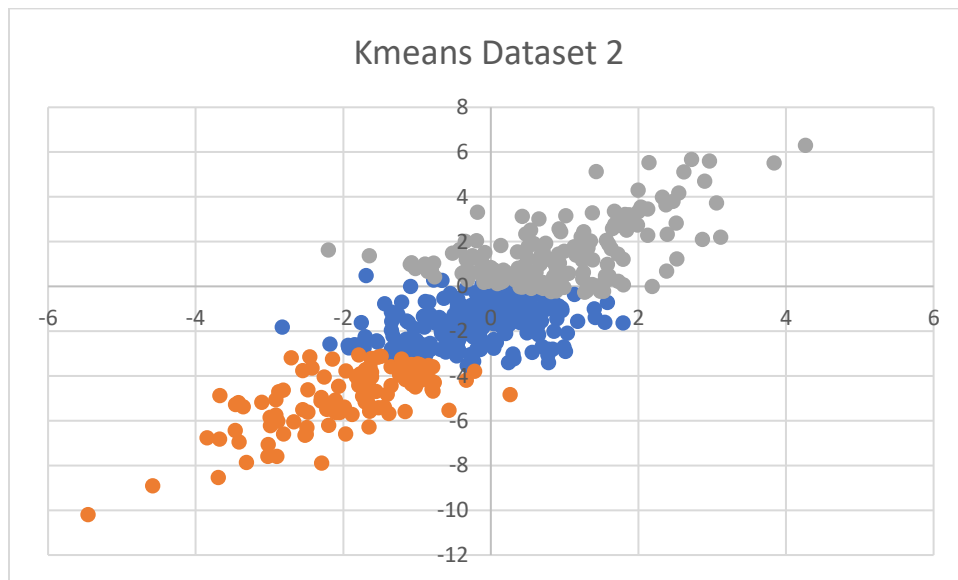
$$I = 5/20 * \log(20 * 5/16) + 4/20 * \log(20 * 4/16) + 1/20 * \log(20 * 1/16) + 1/20 * \log(20 * 1/16) + 4/20 * \log(20 * 4/16) + 5/20 * \log(20 * 5/16) = 2.2829$$

$$\text{NMI} = I/\sqrt{H(\text{Cluster}) * H(\text{Class})} = 1.145$$



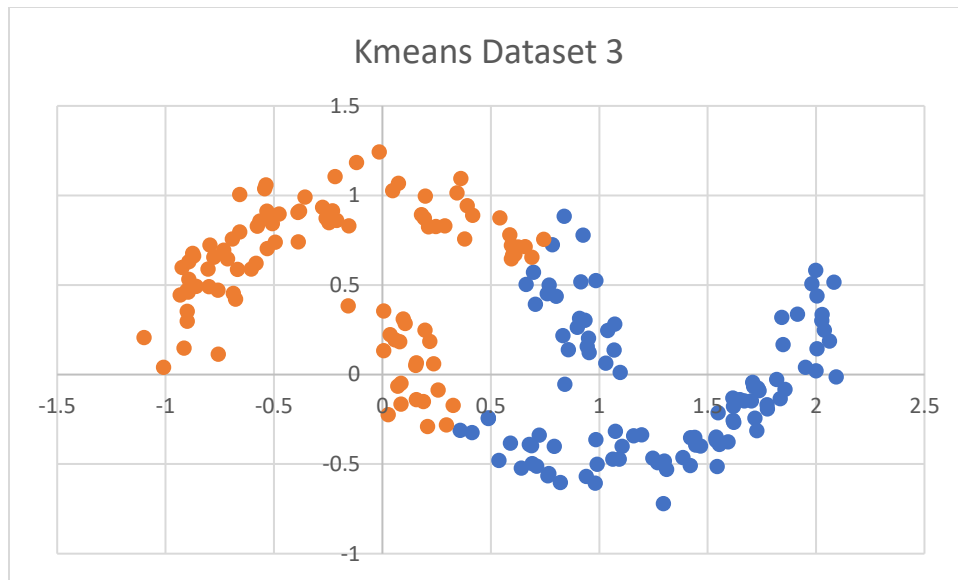
Purity: 1.0

NMI: 1.0



Purity: .764

NMI: 0.04685

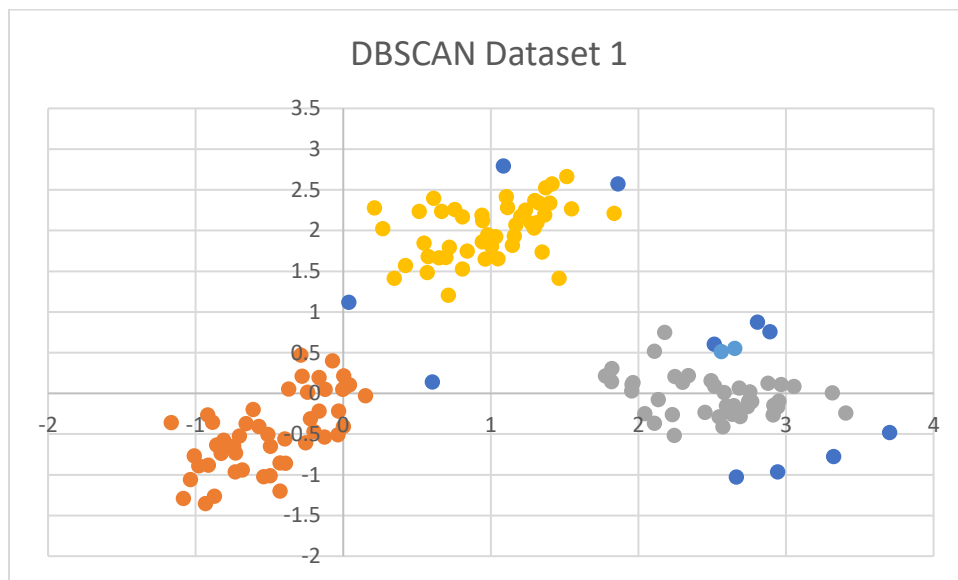


Purity: 0.76

NMI: .145025

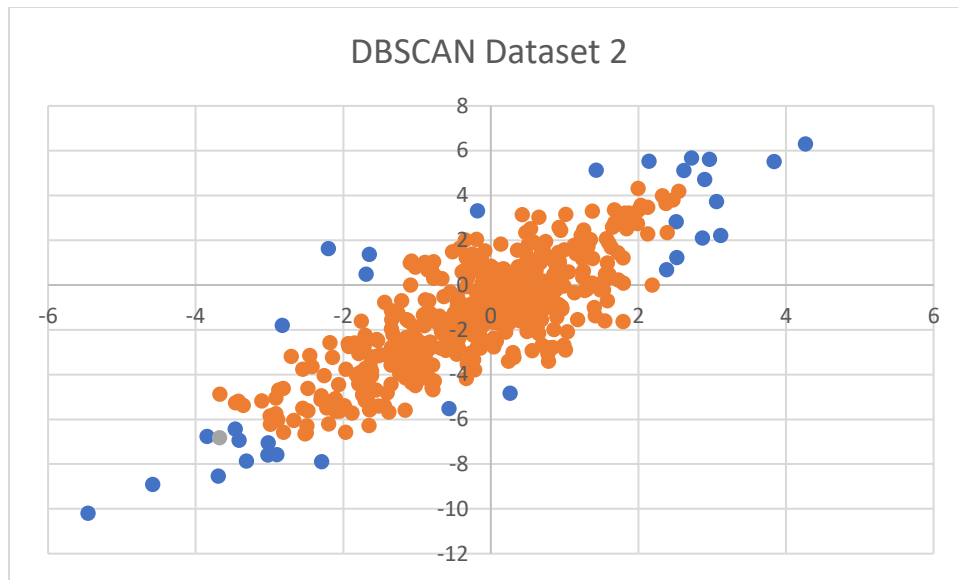
Why KMeans?

KMeans is simple, efficient, and is flexible and the time complexity is pretty good. The problem is that KMeans cannot deal with non-spherical clusters very well, clusters that do not have uniform size or density clusters, and you have to specify the number of clusters before the algorithm starts.



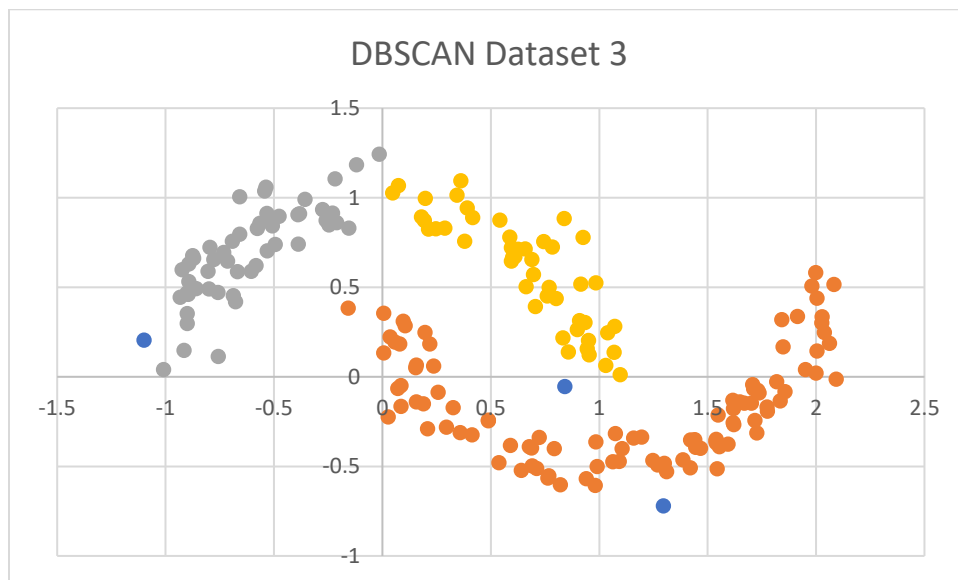
Purity: 0.926666

NMI: 0.9745



Purity: 0.708

NMI: 0.00525345

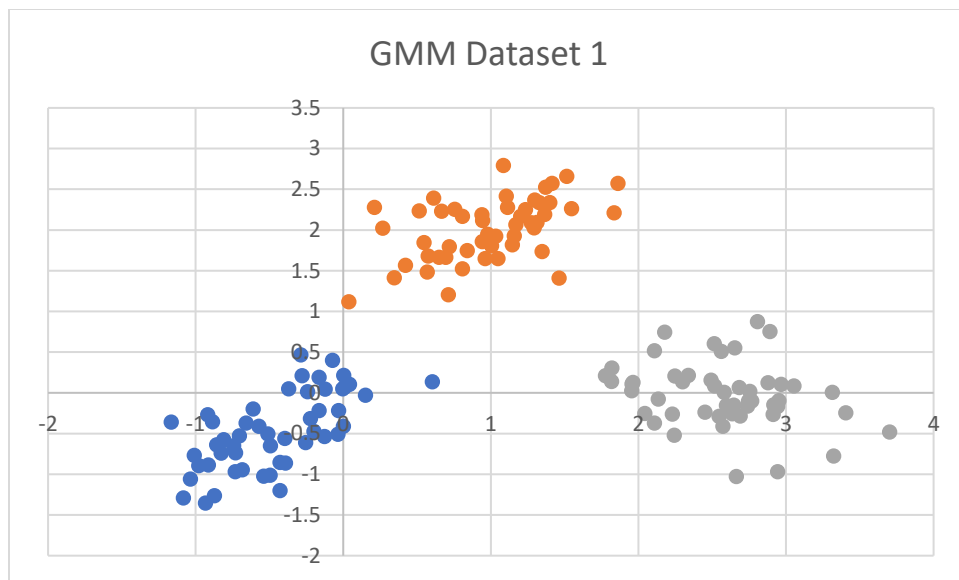


Purity: 0.985

NMI: 0.81735

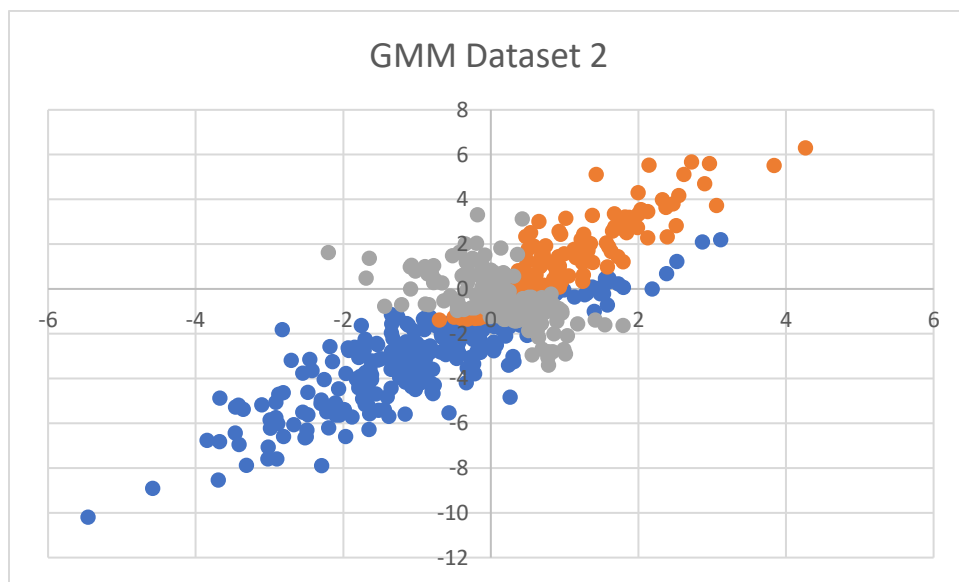
Why DBSCAN?

It can do clusters of arbitrary shapes and it is robust to noise. The problem is that it isn't very good if there aren't enough connected regions of sufficiently high density.



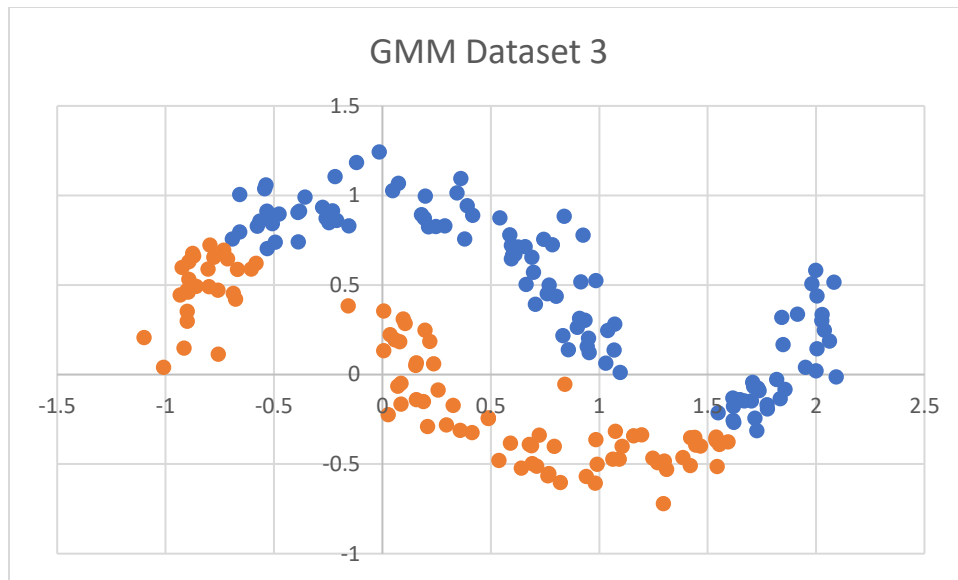
Purity: 1.0

NMI: 1.0



Purity: 0.764

NMI: 0.075603



Purity: 0.69

NMI: .0759484

Why GMM?

Probably has the best accuracy and is better at clusters of differing sizes and shapes. The problem of GMM is that the computation time is much longer and intensive compared to the other clustering methods.