Problem A:

Nguyen S. Nguyen
UID. 004870721

**1** $k=2$. initial center points are $X_2$ and $X_4$.
- Use Euclidean distance

Run first iteration of k-means.

a) Center $X_2$

$X_2 \to X_1 : \sqrt{(2-1)^2 + (1-1)^2 + (5-1)^2} = \sqrt{1^2 + 4^2} = \sqrt{17}$

$X_2 \to X_3 : \sqrt{(6-2)^2 + (7-1)^2 + (2-5)^2} = \sqrt{4^2 + 6^2 + 3^2} = \sqrt{61}$

$X_2 \to X_4 : \sqrt{(9-2)^2 + (8-1)^2 + (7-5)^2} = \sqrt{7^2 + 7^2 + 2^2} = \sqrt{102}$

$X_2 \to X_5 : \sqrt{(3-2)^2 + (5-1)^2 + (1-5)^2} = \sqrt{1^2 + 4^2 + 4^2} = \sqrt{33}$

$X_2 \to X_6 : \sqrt{(5-2)^2 + (3-1)^2 + (7-5)^2} = \sqrt{3^2 + 2^2 + 3^2} = \sqrt{22}$

Center $X_4$:

$X_4 \to X_1 : \sqrt{(9-1)^2 + (8-1)^2 + (7-1)^2} = \sqrt{8^2 + 7^2 + 6^2} = \sqrt{149}$

$X_4 \to X_2 : \sqrt{(9-2)^2 + (8-1)^2 + (7-5)^2} = \sqrt{7^2 + 7^2 + 2^2} = \sqrt{102}$

$X_4 \to X_3 : \sqrt{(9-6)^2 + (8-7)^2 + (7-2)^2} = \sqrt{3^2 + 1^2 + 5^2} = \sqrt{35}$

$X_4 \to X_5 = \sqrt{(9-3)^2 + (8-5)^2 + (7-1)^2} = \sqrt{6^2 + 3^2 + 6^2} = \sqrt{81}$

$X_4 \to X_6 = \sqrt{(9-5)^2 + (8-3)^2 + (7-2)^2} = \sqrt{4^2 + 5^2 + 5^2} = \sqrt{66}$

b) | First cluster: $X_2, X_1, X_6, X_5$ | | Second cluster: $X_4, X_3$ |

pick points closest to initial points

c) First cluster $= \left[ \dfrac{(1+2+5+3)}{4}, \dfrac{(1+1+3+5)}{4}, \dfrac{(1+5+1+2)}{4} \right] = (2.75, 2.5, 2.25)$

Second cluster $= \left[ \dfrac{(6+9)}{2}, \dfrac{(7+8)}{2}, \dfrac{(7+2)}{2} \right] = (7.5, 7.5, 4.5)$

<u>2.</u> There are $k^N$ ways to partition N data points into a cluster. For each iteration, we produce a new clustering based only on the old clustering.

1. If old clustering is the same as new, then next clustering will be the same.

2. If new clustering is different from old, then

The algorithm repeats until there is no further change in cost. Otherwise, you would have some clustering which has a lower cost then itself which is possible.

Problem B Birch:   Branching Factor = 2     $D = \sqrt{\dfrac{2n(SS) - 2(LS)^2}{(n(n-1))}}$

# of entry in leaf = 2

Cluster threshold (diameter) = 2     $CF = (n, LS, SS)$

Table:

| Point | Coordinate |
|-------|-----------|
| X1 | (3,4) |
| X2 | (0,2) |
| X3 | (1,1) |
| X4 | (1,5) |

Insert Point $x_1$ (3,4) into New leaf

$$CF_1 = \langle 1, (3,4), (9,16)\rangle$$

$x_1 = (3,4)$

Insert Point $x_2$ (0,2) into $CF_1$ leaf temporarily. Calculate $CF_1$ again.

$$CF_1' = \langle 2, \overset{LS}{(3,6)}, \overset{SS}{(9,20)}\rangle$$

$$D = \sqrt{\dfrac{[2nSS_x - 2LS_x^2] + [2nSS_y - 2LS_y^2]}{n(n-1)}}$$

$$= \sqrt{\dfrac{[2(2)(9) - 2(3)^2] + [2(2)(20) - 2(6)^2]}{(2)(2-1)}} = \sqrt{\dfrac{18 + 8}{2}} = 3.60$$

Did not meet diameter threshold. Split into new leaf

$$CF_1 = \langle 1, (3,4), (9,16)\rangle \quad CF_2 = \langle 1, (0,2), (0,4)\rangle$$

$x_1 = (3,4)$     $x_2 = (0,2)$

Insert $x_3$ (1,1)

$\quad\hookrightarrow$ Calculate distance from $CF_1$ to $x_3$

$$\sqrt{(3-1)^2+(4-1)^2} = \sqrt{2^2+3^2} = \sqrt{13} = 3.60$$

$\quad\hookrightarrow$ Calculate distance from $CF_2$ to $x_3$

$$\sqrt{(0-1)^2+(2-1)^2} = \sqrt{1^2+1^2} = \sqrt{2} = \underline{1.41}$$

$x_3$ will insert temporarily to $x_3$ b/c it's the shortest distance.

$\quad CF_2' = \langle 2, (1,3), (1,5) \rangle$

$$D = \sqrt{\frac{[(2)(2)(1)-2(1)^2]+[2(2)(5)-2(3)^2]]}{2(2-1)}} \quad \sqrt{\frac{2+2}{2}} = \sqrt{2} = \underline{1.41}$$

Diameter is still within threshold. Will not split into new 1cluster.

$\boxed{CF_1 = \langle 1, (3,4), (9,16) \rangle}$ $\boxed{CF_2 = \langle 2, (1,3), (1,5) \rangle}$

$\qquad\qquad$ | $\qquad\qquad\qquad$ |

$\quad\bigcirc x_1 = (3,4)\qquad\qquad\bigcirc x_2 = (0,2)$
$\qquad\qquad\qquad\qquad\qquad\qquad x_3 = (1,1)$

insert $x_4$ (1,5)

$\quad\hookrightarrow$ Calculate distance from $CF_1$ to $x_4$

$$\sqrt{(3-1)^2+(4-5)^2} = \sqrt{2^2+1^2} = \sqrt{5} = \underline{2.24}$$

$\quad\hookrightarrow$ Calculate distance from $CF_2$ to $x_4$.

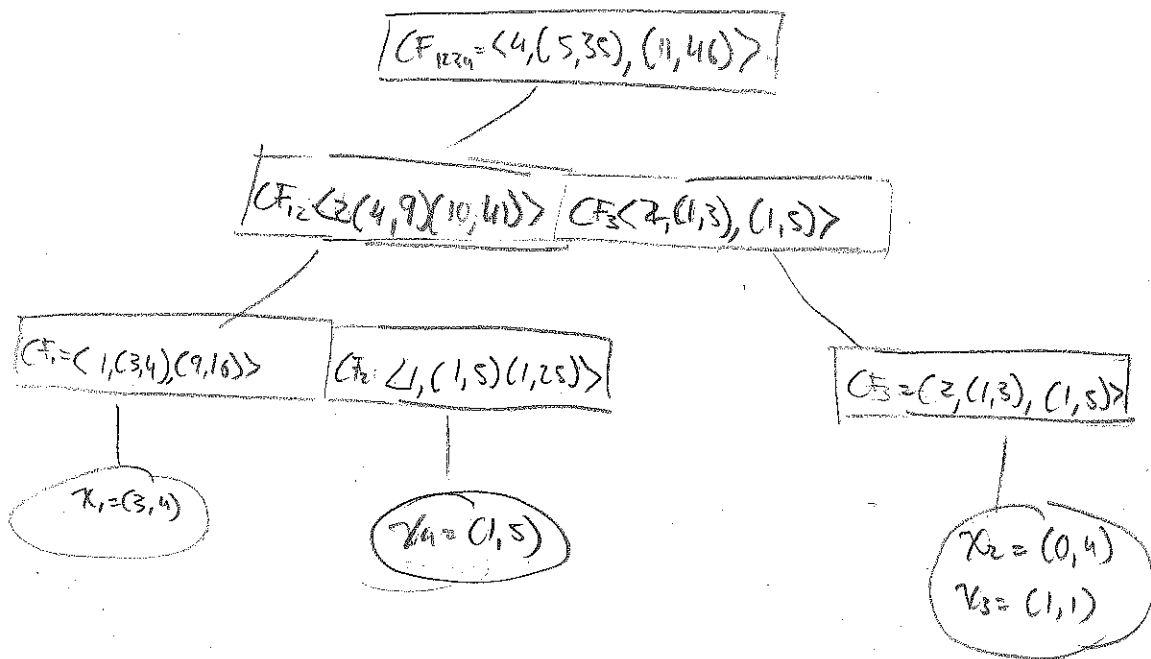$\quad "CF_2$ centroid $= \left(\dfrac{0+1}{2}, \dfrac{2+1}{2}\right) = (1.5, 1)$

$$\sqrt{(1-1.5)^2+(1-5)^2} = \sqrt{(0.5)^2+4^2} = \sqrt{16.25} = 4.03$$

insert to $CF_1$ temporarily b/c it's closer.

Calculate $CF_1 = \langle 2, (4,9), (10,41) \rangle$

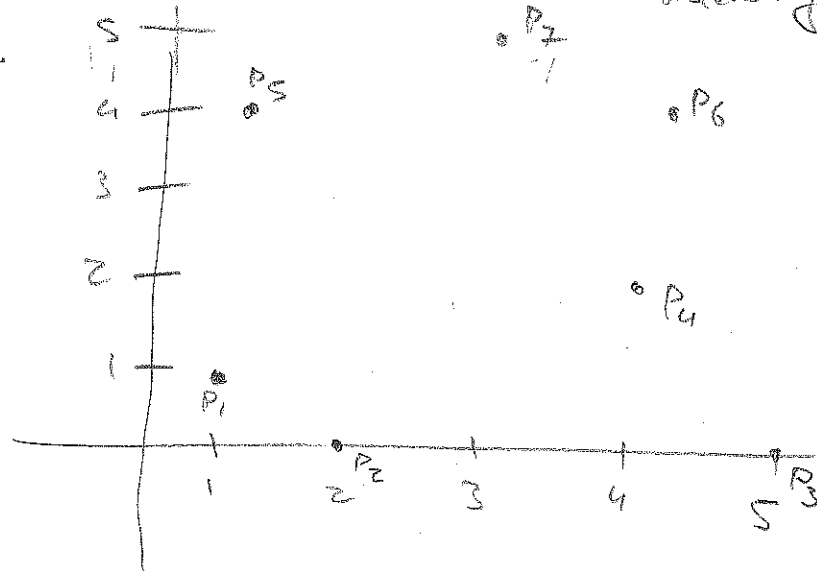$$D = \sqrt{\frac{\left[(2)(2)(10) - 2(4)^2\right] + \left[2(2)(41) - (2)(9)^2\right]}{|2(2-1)|}} = \sqrt{\frac{8+2}{2}} = 2.24$$

Did not meet threshold, split into new leaf but already max out # of entry in leaf. Need to create new branch.



$CF_{1234} = \langle 4, (5,35), (11,41) \rangle$

$CF_{12} \langle 2, (4,9)(10,41) \rangle$  $CF_3 \langle 2, (1,3), (1,5) \rangle$

$CF_1 = \langle 1, (3,4), (9,16) \rangle$  $CF_2 \langle 1, (1,5)(1,25) \rangle$  $CF_5 = \langle 2, (1,3), (1,5) \rangle$

$x_1 = (3,4)$

$x_4 = (1,5)$

$x_2 = (0,4)$
$x_3 = (1,1)$

# Problem C: DBSCAN

$\varepsilon = 2$, min Pts = 2, Points are processed in ascending order of IDs.

| ID | Coordinate |
|----|-----------|
| 1  | (1,1)     |
| 2  | (2,0)     |
| 3  | (5,0)     |
| 4  | (4,2)     |
| 5  | (1,4)     |
| 6  | (4,4)     |
| 7  | (3,5)     |
| 8  | (5,5)     |



## $P_1$ Process: Use Euclidean distances

$P_1 \to P_2 : \sqrt{(1-2)^2 + (1-0)^2} = \sqrt{2} = 1.41$

$P_1 \to P_5 = \sqrt{(1-1)^2 + (1-4)^2} = \sqrt{9} \quad 3$

$P_1 \to P_7 = \sqrt{(1-3)^2 + (1-5)^2} = \sqrt{2^2 + 4^2} = \sqrt{20} = 4.47$

$P_1 \to P_6 = \sqrt{(1-4)^2 + (1-4)^2} = \sqrt{3^2 + 3^2} = \sqrt{18} = 4.24$

$P_1 \to P_4 = \sqrt{(1-4)^2 + (1-2)^2} = \sqrt{3^2 + 1^2} = \sqrt{10} = 3.16$

$P_1 \to P_8 = \sqrt{(1-5)^2 + (1-5)^2} = \sqrt{4^2 + 4^2} = \sqrt{32} = 5.66$

$P_1 \to P_3 = \sqrt{(1-5)^2 + (1-0)^2} = \sqrt{4^2 + 1^2} = \sqrt{17} = 4.12$

$\underline{P_1 \text{ is border point b/c } P_2 \text{ is within } \varepsilon ps.}$

## $P_2$ Process:

$P_2 \to P_1 = (\text{Visited already}) = \sqrt{2} = 1.41$

$P_2 \to P_3 = \sqrt{(2-5)^2 + (0-0)^2} = \sqrt{9} = 3$

$P_2 \to P_4 = \sqrt{(2-4)^2 + (0-2)^2} = \sqrt{2^2 + 2^2} = \sqrt{8} = 2.83$

$P_2 \to P_5 = \sqrt{(2-1)^2 + (0-4)^2} = \sqrt{1^2 + 4^2} = \sqrt{17} = 4.12$

$P_2 \to P_6 = \sqrt{(2-4)^2 + (0-4)^2} = \sqrt{2^2 + 4^2} = \sqrt{20} = 4.47$

$P_2 \to P_7 = \sqrt{(2-3)^2 + (0-5)^2} = \sqrt{1^2 + 5^2} = \sqrt{26} = 5.09$

$P_2 \to P_8 = \sqrt{(2-5)^2 + (0-5)^2} = \sqrt{3^2 + 5^2} = \sqrt{34} = 5.83$

$\underline{P_2 \text{ is border point b/c } P_1 \text{ is within } \varepsilon ps.}$

## $P_3$ Process

$P_3 \to P_1 =$ (Visited already) = $\qquad = \sqrt{17} = 4.12$  (Visited)

$P_3 \to P_2 =$ (Visited Already) $\qquad = \sqrt{9} = 3$  (Visited)

$P_3 \to P_4 = \sqrt{(5-4)^2 + (0-2)^2} = \sqrt{1^2 + 2^2} = \sqrt{5} = 2.26$

$P_3 \to P_5 = \sqrt{(5-1)^2 + (0-4)^2} = \sqrt{4^2 + 4^2} = \sqrt{32} = 5.66$

$P_3 \to P_6 = \sqrt{(5-4)^2 + (0-4)^2} = \sqrt{1^2 + 4^2} = \sqrt{17} = 4.12$

$P_3 \to P_7 = \sqrt{(5-3)^2 + (0-5)^2} = \sqrt{2^2 + 5^2} = \sqrt{29} = 5.38$

$P_3 \to P_8 = \sqrt{(5-5)^2 + (0-5)^2} = \sqrt{5^2} = 5$

$P_3$ is outlier, form

## $P_4$ Process

$P_4 \to P_1 =$ (Visited Already) $= \sqrt{3^2 + 1^2} = \sqrt{10} = 3.16$  (Visited)

$P_4 \to P_2 =$ (Visited Already) $\sqrt{2^2 + 2^2} = \sqrt{8} = 2.83$  (Visited)

$P_4 \to P_3 =$ (Visited Already) $= \sqrt{1^2 + 2^2} = \sqrt{5} = 2.24$  (Visited)

$P_4 \to P_5 = \sqrt{(4-1)^2 + (2-4)^2} = \sqrt{3^2 + 2^2} = \sqrt{13} = 3.61$

$P_4 \to P_6 = \sqrt{(4-4)^2 + (2-4)^2} = \sqrt{2^2} = 2$

$P_4 \to P_7 = \sqrt{(4-3)^2 + (2-5)^2} = \sqrt{1+9} = \sqrt{10} = 3.16$

$P_4 \to P_8 = \sqrt{(4-5)^2 + (2-5)^2} = \sqrt{1^2 + 3^2} = \sqrt{10} = 3.16$

$P_4$ is border point.

## P5 Process

$P_5 \to P_1 = (\text{Visited}) \qquad = \sqrt{3^2} \qquad = 3 \qquad (\text{Visited})$

$P_5 \to P_2 = (\text{Visited}) \qquad = \sqrt{1+4^2} = \sqrt{17} = 4.12 \qquad (\text{Visited})$

$P_5 \to P_3 = (\text{Visited}) \qquad = \sqrt{4^2+4^2} = \sqrt{32} = 5.66 \qquad (\text{Visited})$

$P_5 \to P_4 = (\text{Visited}) \qquad = \sqrt{3^2+2^2} = \sqrt{13} = 3.6 \qquad (\text{Visited})$

$P_5 \to P_6 = \sqrt{(1-4)^2 + (4-4)^2} = \sqrt{3^2+0^2} = \sqrt{3^2} = 3$

$P_5 \to P_7 = \sqrt{(1-3)^2 + (5-4)^2} = \sqrt{2^2+1^2} = \sqrt{5} = 2.24$

$P_5 \to P_8 = \sqrt{(1-5)^2 + (4-5)^2} = \sqrt{4^2+1^2} = \sqrt{17} = 4.12$

**P5 is outlier.**

## P6 Process

$P_6 \to P_1 = (\text{Visited}) \qquad = \sqrt{3^2+3^2} = \sqrt{18} = 4.24$

$P_6 \to P_2 = (\text{Visited}) \qquad \sqrt{(2)^2+4^2} = \sqrt{20} = 4.47$

$P_6 \to P_3 = (\text{Visited}) \qquad \sqrt{1^2+4^2} = \sqrt{12} = 4.12$

$P_6 \to P_4 = (\text{Visited}) \qquad = \sqrt{2^2} = 2$

$P_6 \to P_5 = \sqrt{(\text{Visited})} \qquad = \sqrt{(3)^2} = 3$

$P_6 \to P_7 = \sqrt{(4-3)^2 - (4-5)^2} = \sqrt{1^2+1^2} = \sqrt{2} = \underline{1.41}$

$P_6 \to P_8 = \sqrt{(4-5)^2 + (4-5)^2} = \sqrt{1^2+1^2} = \sqrt{2} = \underline{1.41}$

**P6 is core point. Merge with P4.**

## $P_7$ Process

$P_7 \rightarrow P_1 =$ (Visited) $= \sqrt{2^2 + 4^2} = \sqrt{20} = 4.47$

$P_7 \rightarrow P_2 =$ (Visited) $= \sqrt{1^2 + 5^2} = \sqrt{26} = 5.09$

$P_7 \rightarrow P_3 =$ (Visited) $= \sqrt{2^2 + 5^2} = \sqrt{29} = 5.38$

$P_7 \rightarrow P_4 :$ (Visited) $= \sqrt{1^2 + 3^2} = \sqrt{10} = 3.16$

$P_7 \rightarrow P_5 =$ (Visited) $= \sqrt{2^2 + 1^2} = \sqrt{5} = 2.24$

$P_7 \rightarrow P_6 =$ (Visited) $\sqrt{1^2 + 1^2} = \sqrt{2} = 1.41$

$P_7 \rightarrow P_8 = \sqrt{(3-5)^2 + (5-5)^2} = \sqrt{2^2 + 0} = \sqrt{2^2} = 2$

$\underline{P_7 \text{ is Core Point . Merge with } P_6, P_4.}$

## $P_8$ Process

$P_8 \rightarrow P_1 =$ (Visited) $= \sqrt{4^2 + 4^2} = \sqrt{32} = 5.66$

$P_8 \rightarrow P_2 =$ (Visited) $= \sqrt{3^2 + 5^2} = \sqrt{34} = 5.83$

$P_8 \rightarrow P_3 =$ (Visited) $= \sqrt{5^2} = 5$

$P_8 \rightarrow P_4 =$ (Visited) $= \sqrt{1^2 + 5^2} = \sqrt{10} = 3.16$

$P_8 \rightarrow P_5 =$ (Visited) $= \sqrt{4^2 + 1^2} = \sqrt{17} = 4.12$

$P_8 \rightarrow P_6 =$ (Visited) $= \sqrt{2} = 1.41$

$P_8 \rightarrow P_7 =$ (Visited) $= \sqrt{2^2} = 2$

$\underline{P_8 \text{ is core point, merge with } P_7, P_6, P_4}$

Since $P_1, P_2$ are not connected
to a core point, then they
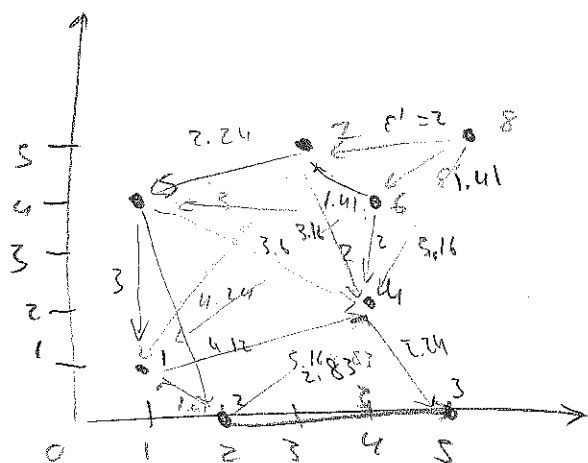will be converted to outliers there
is only one cluster that formed.

$C_1 = \{ P_7, P_6, P_4 \}$

| ID | Coordinate | |
|---|---|---|
| 1 | (1,1) | border point → noise |
| 2 | (2,0) | border point → noise |
| 3 | (5,0) | noise |
| 4 | (4,2) | border point |
| 5 | (1,4) | noise |
| 6 | (4,4) | core point |
| 7 | (3,5) | core point |
| 8 | (5,5) | core point |

initially border point
but later
↓ converted to noise

**Problem D**  OPICS (using calculations in problem C, DBSCAN)

$\varepsilon = 3$, min pts = 2, Run algorithm, starting from point #8 (5,5)

| ZD | coordinate |
|----|-----------|
| 1 | (1,1) |
| 2 | (2,0) |
| 3 | (5,0) |
| 4 | (4,2) |
| 5 | (1,4) |
| 6 | (4,4) |
| 7 | (3,5) |
| 8 | (5,5) |



Starting at #8 (core distance = 2)

- get neighbors: Point 6, 7

  ↳ output Point 8 to ordered list. ⊖

  ↳ Pt 6, pt 7 = core distance $\varepsilon' = 2$     Comparing to Point 8

  1. Seeds = empty priority queue

     ↳ update priority queue()  →

     | Point 6 |
     |---------|
     | Point 7 |

     } assume lower ID is higher priority since both points are tie with same core distance

     ↳ Look at Point 6 (core distance = 1.41)

        ↳ get Neighbors: Point 7, Point 4 (Point 8) points already looked at

     ↳ output Point 6 to ordered list

        ↳ update priority queue() →     Compare to Point 6.

        | Point 7 | $\varepsilon' = 1.41$ from point 6 |
        |---------|----------|
        | Point 4 | $\varepsilon = 2$ |
        | Point 5 | $\varepsilon = 3$ |

        ↳ look at Point 7 (core distance = 2)

           ↳ get Neighbors: (Point 5, Point 8) ← already looked at

           ↳ output point 7 to ordered list

           | Point 4 at 2 | b/c point 4 reachable distance is still smaller. |
           |------|------|
           | Point 5 | |

           ↳ update priority Queue() →

           distance = 2.24

↳ Look at Point 4 (core distance = 2.24)
  ↳ Get Neighbors : Pt 3, 6, 2   →already looked at
  ↳ output Point 4 to ordered list
    ↳ update priority queue()

| Point 3 |
| Point 5 |  } 2.24
| Point 2 |  } 2.83

} same distance but lower if it is higher priority.

↳ look at Point 3 (core distance = 3)
  ↳ Get Neighbors: Pt 3, ④, → already looked at
  ↳ output Point 3 to ordered list
    ↳ update priority queue()

| Point 5 |
| Point 2 |  b/c distance is still 2.83

↳ Look at Point 5 (core distance = 3)
  ↳ Neighbors: Pt ④, 1, ⑥  → already looked at
  ↳ output Point 5 to ordered list

| Point 2 |
| Point 1 |  ε = 3

↳ Look at Point 2 (core distance = 2.83)
  ↳ Neighbors: Pt 1, 3, 4  → already looked at

| point 1 |

↳ look at Point 1
  ↳ output point 1 to ordered list

↳

| Point 8 |
| Point 6 |
| Point 7 |
| Point 4 |
| Point 3 |
| Point 5 |
| Point 2 |
| Point 1 |

Problem E:

1) — Calculate information of class label.

Class P: Play-Golf = "Yes"
Class N: play-Golf = "No"

$$I(p,n) = I(4,2) = -\frac{4}{6}\log_2\left(\frac{4}{6}\right) - \frac{2}{6}\log_2\left(\frac{2}{6}\right)$$

$$= \left[\frac{\ln(4/6)}{\ln(2)}\right] \times \left(\frac{-4}{6}\right) + \left[\frac{\ln(2/6)}{\ln(2)}\right] \times \left(\frac{-2}{6}\right)$$

$$= .390 + .528 \Rightarrow \boxed{I(p,n) = 0.918}$$

| Windy | $P_i$ Yes | $N_i$ No | $I(p_i,n_i)$ |
|---|---|---|---|
| True | 2 | 2 | $-\frac{2}{4}\log_2\left(\frac{2}{4}\right) - \frac{2}{4}\log_2\left(\frac{2}{4}\right) = 0.5+0.5=1$ |
| False | 2 | 0 | 0 |

$$E(windy) = \frac{4}{6}I(2,2) + \frac{2}{6}(I(2,0)) = \frac{4}{6}(1) + \frac{2}{6}(0) \Rightarrow \boxed{E(windy) = .667}$$

$$Gain(windy) = I(p,n) - E(windy) = 0.918 - .667 = .251$$

$$\boxed{Gain(windy) = .251}$$

| Humidity | $P_i$ Yes | $n_i$ No | $I(p,n)$ |
|---|---|---|---|
| normal | 3 | 1 | $-\frac{3}{4}\log_2\left(\frac{3}{4}\right) - \frac{1}{4}\log_2\left(\frac{1}{4}\right) = .311 + .5 = .811$ |
| high | 1 | 1 | $-\frac{1}{2}\log_2\left(\frac{1}{2}\right) - \frac{1}{2}\log_2\left(\frac{1}{2}\right) = 0.5+0.5=1$ |

$$E(humidity) = \frac{4}{6}I(3,1) + \frac{2}{6}I(1,1) = \frac{4}{6}(.811) + \frac{2}{6}(1) \Rightarrow \boxed{E(humidity) = .874}$$

$$Gain(humidity) = I(p,n) - E(humidity) = 0.918 - .874 \Rightarrow \boxed{Gain(humidity) = .044}$$

| Outlook | $P_i$ Yes | $n_i$ No | $I(p_i,n_i)$ |
|---|---|---|---|
| overcast | 2 | 0 | 0 |
| rain | 2 | 1 | $-\frac{2}{3}\log_2\left(\frac{2}{3}\right) - \frac{1}{3}\log_2\left(\frac{1}{3}\right) = .390 + .528 = .918$ |
| sunny | 0 | 1 | 0 |

$$E(outlook) = \frac{2}{6}(0) + \frac{3}{6}(.918) + \frac{1}{6}(0) \Rightarrow \boxed{E(outlook) = .459}$$

$$Gain(outlook) = I(p,n) - E(outlook) = 0.918 - .459 = .459$$

$$\boxed{Gain(outlook) = .459}$$

| Temperature | $P_i$ | $n_i$ | $I(P_i, n_i)$ |
|---|---|---|---|
| mild | 4 | 1 | $-\frac{4}{5}\log_2(\frac{4}{5}) - \frac{1}{5}\log_2(\frac{1}{5}) = .257 + .464 = .721$ |
| cool | 0 | 1 | 0 |

$$E(Temperature) = \frac{5}{6}(.721) + 1(0) \Rightarrow \boxed{E(Temperature) = .601}$$

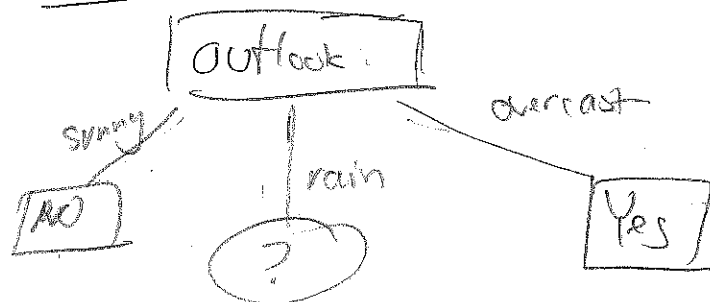$$Gain(Temperature) = I(p,n) - E(Temperature) = .918 - .601 = .317$$

$$\boxed{Gain(Temperature) = .317}$$

| Day | $P_i$ | $n_i$ | $I(P_i, n_i)$ |
|---|---|---|---|
| 07-20 | 1 | 0 | 0 |
| 07-21 | 1 | 0 | 0 |
| 07-22 | 1 | 0 | 0 |
| 07-23 | 0 | 1 | 0 |
| 07-26 | 0 | 1 | 0 |
| 07-30 | 1 | 0 | 0 |

$$\boxed{E(Day) = 0}$$

$$Gain(Day) = I(p,n) - E(Day) = .918 - 0 \Rightarrow \boxed{Gain(Day) = .918}$$

Decision Tree:



Looking only at outlook= rain for next step

| temp. | $P_i$ | $n_i$ | $I(P_i, n_i)$ |
|---|---|---|---|
| mild | 2 | 0 | 0 |
| cool | 0 | 1 | 0 |

$$E(temp) = \frac{2}{3}(0) + \frac{1}{3}(0) = 0$$
$$Gain(temp) = .918$$

| humidity | $P_i$ | $n_i$ | $I(P_i, n_i)$ |
|---|---|---|---|
| high | 1 | 0 | 0 |
| normal | 1 | 1 | 1 |

$$E(humidity) = \frac{1}{3}(0) + \frac{2}{3}(1) = \frac{2}{3} = .667$$
$$Gain(humidity) = .918 - .667 = .251$$

| Windy | $P_i$ | $n_i$ | $I(P_i, n_i)$ |
|---|---|---|---|
| true | 1 | 1 | 1 |
| false | 1 | 0 | 0 |

$$E(humidity) = .667$$
$$Gain(humidity) = .251$$

**Outlook**
- sunny → No
- overcast → Yes
- rain → **Temperature**
  - mild → Yes — pure leaves
  - cool → NO — pure leaves, no need to further expansion.

Days initially had the hight Information gain but were very shallow. The decision tree would have been shallow and broad which would have not been a good decision tree.

2) Using Naive Bayes Classification, predict:

| Temperature | Yes | No | Outlook | Yes | No | Humidity | Yes | No | Windy | Yes | No |
|---|---|---|---|---|---|---|---|---|---|---|---|
| mild | 3/4 | 1 | overcast | 2 | 0 | normal | 3 | 1 | False | 2 | 0 |
| cool | 0 | 1 | sunny | 0 | 1 | high | 1 | 1 | True | 2 | 2 |
|  |  |  | rain | 2 | 1 |  |  |  |  |  |  |
| mild | 4/4 | 1/2 | overcast | 2/4 | 0/2 | normal | 3/4 | 1/2 | False | 2/4 | 0/4 |
| cool | 0/4 | 1/2 | sunny | 0/4 | 1/2 | high | 1/4 | 1/4 | True | 2/4 | 2/4 |
|  |  |  | rain | 2/4 | 1/2 |  |  |  |  |  |  |

P( Temperature = mild, Outlook=rain, Humidity=high, Windy=False | Play Golf= Yes )=

$$1 \times 0.5 \times .25 \times .5 = \boxed{0.0625} \times (c_i(\text{"Yes"})) = \boxed{.04167}$$
$$= .667$$

P( Temperature = mild, Outlook= rain, humidity=high, Windy= True | Play Golf= No )=

$$0.5 \times 0.5 \times 0.5 \times 0 = \boxed{0 \text{ probability of not play golf}}$$
$$\times (c_i(\text{"No"})) \quad \text{with given condition.}$$
$$= 2/6$$

3) Naive Bayes takes less data and greatly reduce computation cost, only count the class distribution.

3) Naive Bayes takes less data and greatly reduce computation cost, only count the class distribution.

Programming Assignment F:

The purpose of this programming assignment is to determine the top three features from Wine Dataset using random forests method. Furthermore, we will investigate the popular clustering methods (K-Means, Birch, Agglomerative Clustering, and DBSCAN) using those top three features that we determined from random forests method. We will run each clustering method for 1, 3, 5, 11 clusters, and provide a 3D graph visualization as well as runtime for each iteration. The programming assignment is done in Python 3.6 using sklearn package (http://scikit-learn.org/ ). The hardware includes a MacBook Pro with Intel core i7 and 16GB of RAM with an operating system of OSX Version 10.12.3.

1) For the random forest method, I used RandomForestClassifier ensemble with n_estimators= 1000 and criterion = 'entropy' as parameters. The parameter, n_estimators, is the number of trees that will be generated in the algorithm. The higher number of trees increases the accuracy of determining the importance of the features. Random Forest essentially takes an average of the all values in the trees. In this case, I have given 1000 trees, which might have been more than enough, in order to accurately determine the top three features. Therefore, the algorithm returns all the features with their respective value below. The top 3 features are: **Flavanoids, Proline, Color Intensity.**

| Flavanoids | 0.17783608295298423 |
| Proline | 0.15316211932034413 |
| Color Intensity | 0.13859670997925955 |
| OD280/OD315 of diluted wines | 0.13543375734648969 |
| Alcohol | 0.11724624000008846 |
| Hue | 0.081493384905736613 |
| Total phenols | 0.05901134355487616 |
| Malic Acid | 0.030775215634327207 |
| Magnesium | 0.028231940198704893 |
| Alcalinity of Ash | 0.025352049995492284 |
| Proanthocyanins | 0.019412659155907064 |
| Ash | 0.013472501555337116 |
| Nonflavanoid phenols | 0.0099182386198699863 |

2) The first clustering method, K-Means, will take parameter n_clusters = 3 against the top 3 features as determined above. In order to generate the run time of K-Means fitting, I use another python library call timeit to output the running time. I start the timer before the fitting function and stop after the fitting function. Running the algorithm 5 different times:

| Run 1 | 0.016372540994780138 |
| Run 2 | 0.015201788017293438 |
| Run 3 | 0.015945635008392856 |
| Run 4 | 0.013141681993147358 |
| Run 5 | 0.011502273991936818 |

The average running time is about **14.42 milliseconds** to fit 3 clusters in K-means method. Below is the 3d graph visualization with 3 clusters.

The second clustering method, Birch, takes in parameter n_clusters = 3 against the top 3 features as determined above, branching factor default to 50 and threshold (diameter of sub-cluster) default to 0.5.  I start the timer before the fitting function and stop after the fitting function. Running the algorithm 5 times:

| Run 1 | 0. 010433707997435704 |
| Run 2 | 0. 009979043010389432 |
| Run 3 | 0. 008869687997503206 |
| Run 4 | 0. 008747225016122684 |
| Run 5 | 0. 009066089987754822 |

The average running time is about **9.42 milliseconds** to fit 3 clusters in Birch method. Below is the 3d graph visualization with 3 clusters.

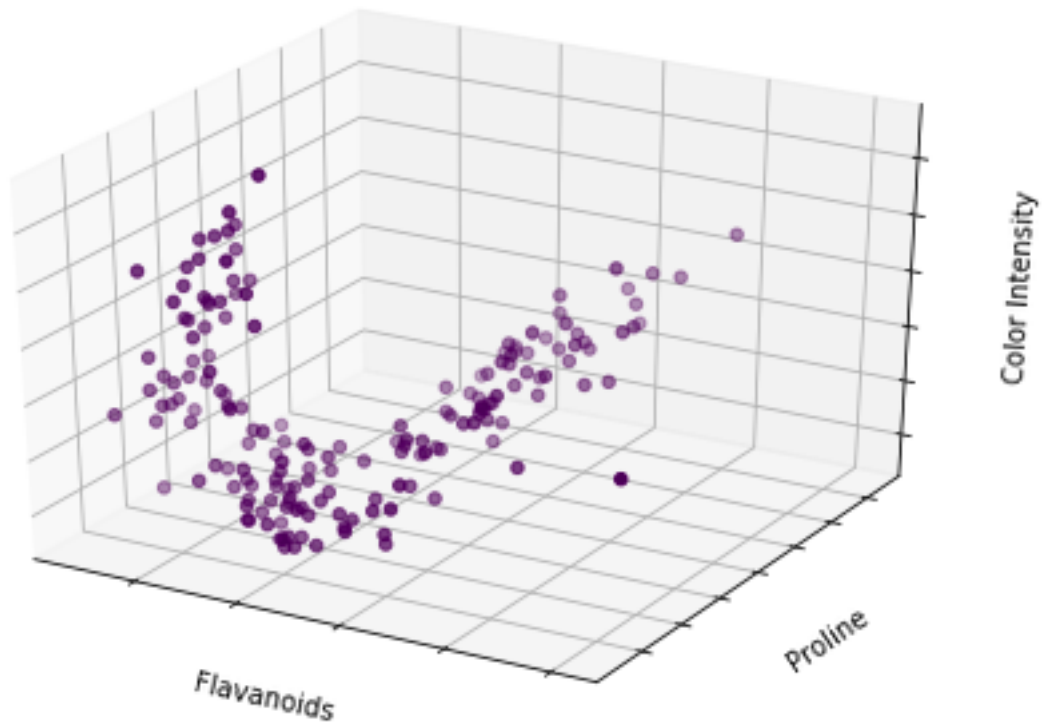The third clustering method, Agglomerative Clustering, will take parameter n_clusters = 3 against the top 3 features, affinity is set to default which is Euclidean, and linkage that is set to ward (minimizes the variance of the clusters being merged) by default. I start the timer before the fitting function and stop after the fitting function. Running the algorithm 5 times:

| Run 1 | 0. 0011187860218342394 |
| Run 2 | 0. 001062514988007024 |
| Run 3 | 0. 001056721986969933 |
| Run 4 | 0. 0011835710029117763 |
| Run 5 | 0. 001036215981002897 |

The average running time is about **1.09 milliseconds** to fit 3 clusters in Agglomerative method. Below is the 3d graph visualization with 3 clusters.
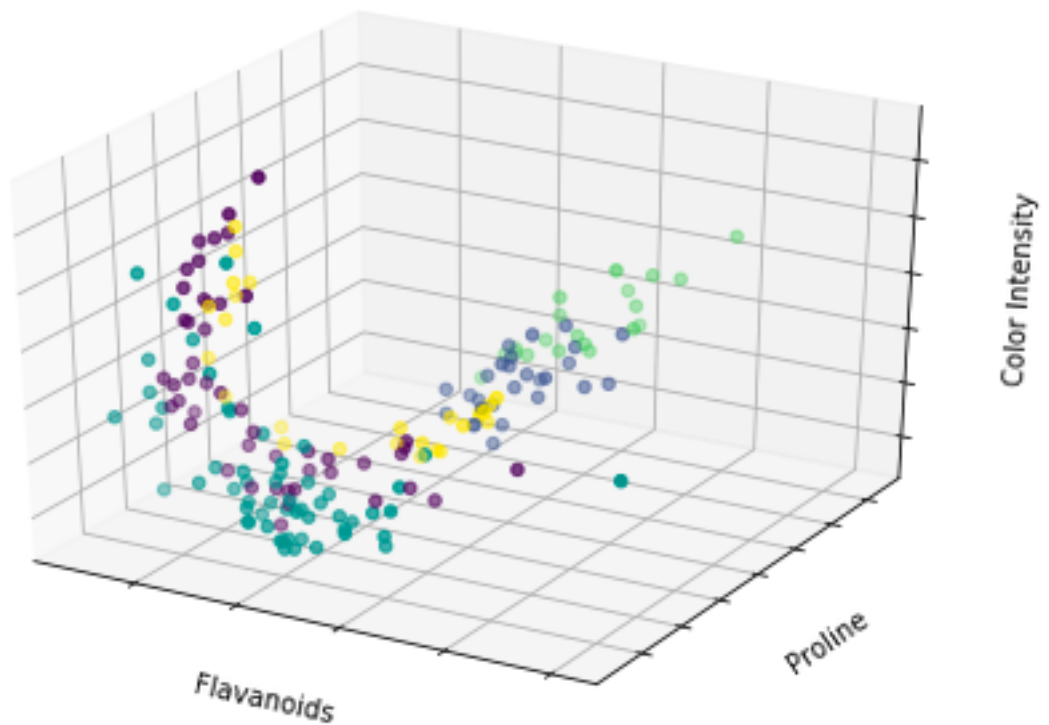
3) K-Means

Running K-means method against 1 cluster for the top 3 features as determined above. I start the timer before the fitting function and stop after the fitting function. Running the algorithm 5 times:

n clusters = 1

| Run 1 | 0. 0036769699945580214 |
|-------|------------------------|
| Run 2 | 0. 003603240998927504 |
| Run 3 | 0. 003719937987625599 |
| Run 4 | 0. 003697048989124596 |
| Run 5 | 0. 004235491011058912 |

The average running time is about **3.78 milliseconds** to fit 1 cluster in K-Means method, which is faster than fitting 3 clusters. Intuitively, running K-means with 1 cluster is much faster because it does not require the algorithm to reexamine every point against the new centers and to recalculate the mean center points. Below is the 3d graph visualization with 1 clusters.

Running K-means method against 5 cluster for the top 3 features. I start the timer before the fitting function and stop after the fitting function. Running the algorithm 5 times:

n clusters = 5

| Run 1 | 0. 015278347011189908 |
|-------|-----------------------|
| Run 2 | 0. 01614384798449464 |
| Run 3 | 0. 014840518997516483 |
| Run 4 | 0. 014069013006519526 |
| Run 5 | 0. 018971159995999187 |

The average running time is about **15.864 milliseconds** to fit 5 clusters in K-Means method, which is slightly slower from what I observed for 3 clusters. Interestingly, running against 5 clusters should be slower than running 3 clusters, however, the average running time is comparable to running 3 clusters. One explanation could be that running K-Means against 5 clusters did not require the algorithm to recalculate the centers again, or had minor recalculation of the mean centers. I also noticed that K-means method was not able to handle non-convex clusters correctly. Below is the 3d graph visualization with 5 clusters.
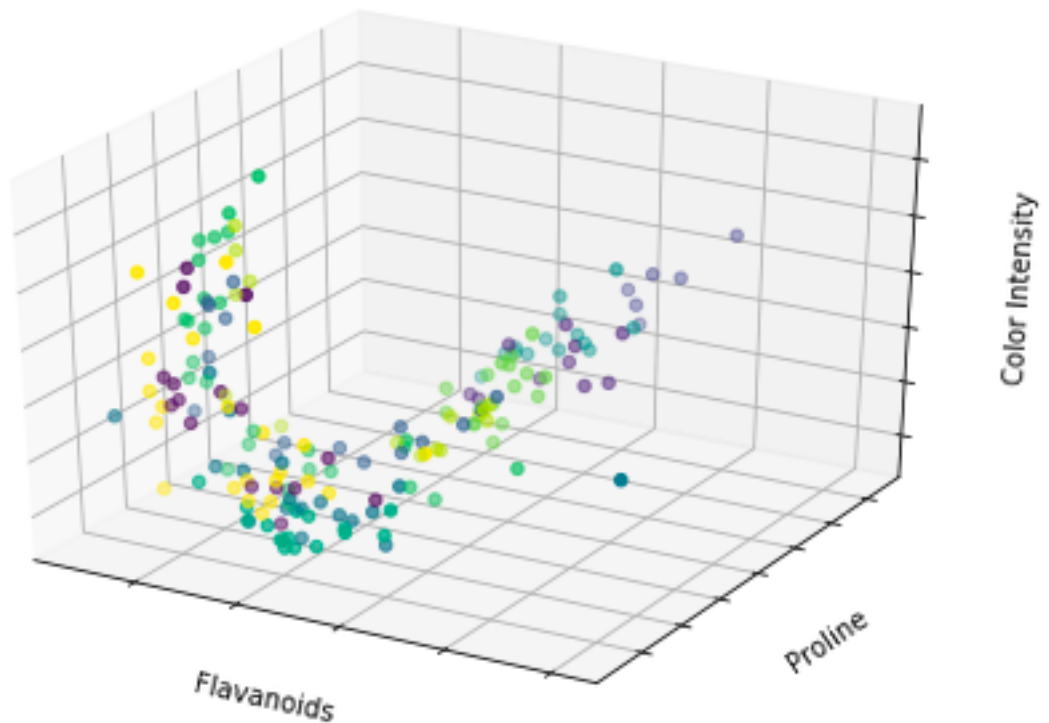
Running K-means method against 11 cluster for the top 3 features as determined above. I start the timer before the fitting function and stop after the fitting function. Running the algorithm 5 times:

n clusters = 11

| Run 1 | 0. 017634788004215807 |
|-------|------------------------|
| Run 2 | 0. 021763442986411974 |
| Run 3 | 0. 021575298014795408 |
| Run 4 | 0. 013702736992854625 |
| Run 5 | 0. 02184795998618938 |

The average running time is about **19.3 milliseconds** to fit 11 clusters in K-Means method, which is significantly slower than 3 clusters as expected. Intuitively, this makes sense since the more clusters require more reexamining of mean points and more recalculating of the new means points. Below is the 3d graph visualization with 5 clusters.
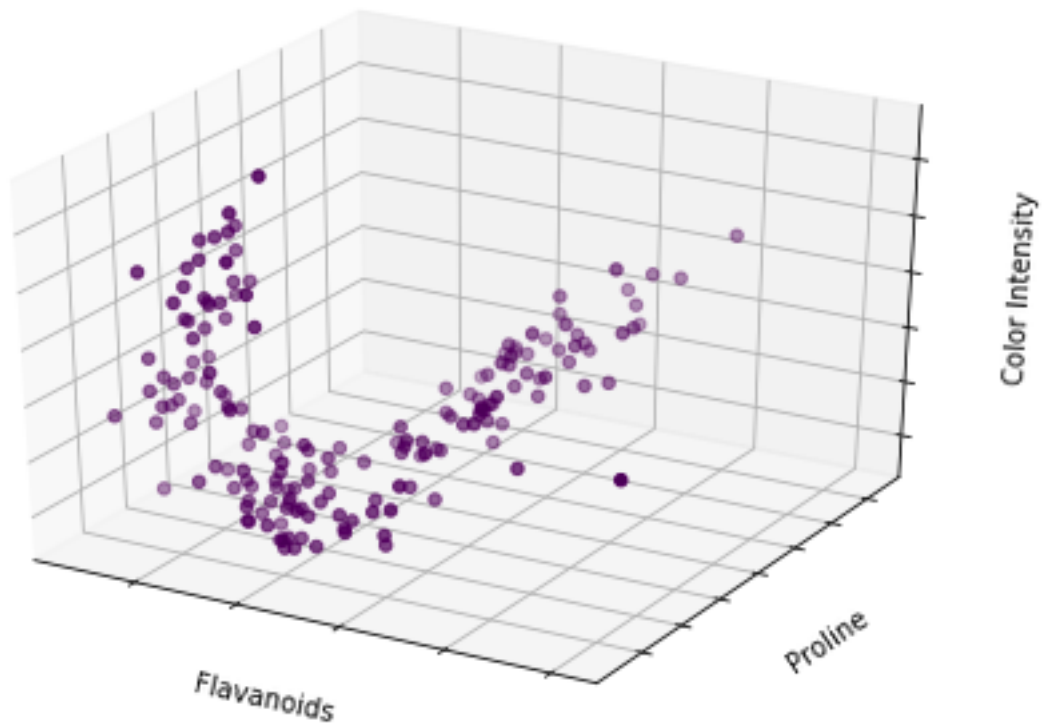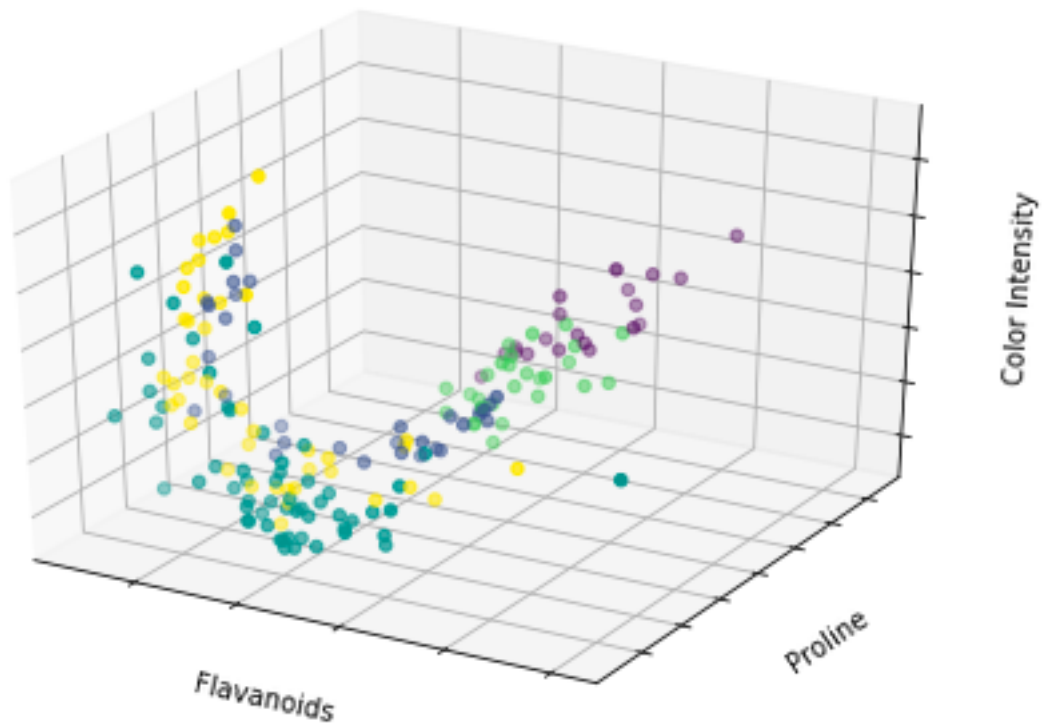
Birch

Running Birch method against 1 cluster for the top 3 features as determined above. I start the timer before the fitting function and stop after the fitting function. Running the algorithm 5 times:

n clusters = 1

| Run 1 | 0. 006875534018035978 |
|-------|------------------------|
| Run 2 | 0. 0070446610043291 |
| Run 3 | 0. 006853063008747995 |
| Run 4 | 0. 006931835989234969 |
| Run 5 | 0. 007208497991086915 |

The average running time is about **6.98 milliseconds** to fit 1 cluster which is faster than 3 cluster. Below is the 3d graph visualization with 1 clusters.

Running Birch method against 5 clusters for the top 3 features as determined above. I start the timer before the fitting function and stop after the fitting function. Running the algorithm 5 times:

n clusters = 5

| Run 1 | 0. 007144284987589344 |
|-------|------------------------|
| Run 2 | 0. 006923697976162657 |
| Run 3 | 0. 006931835989234969 |
| Run 4 | 0. 006867500982480124 |
| Run 5 | 0. 007241435989271849 |

The average running time is about **5.65 milliseconds** to fit 5 cluster. Interestingly, comparing the running time in this iteration to running 1 cluster iteration above, the difference in running time is minimal. This makes sense because Birch method only scans the data once and does not need to reexamine each point again.
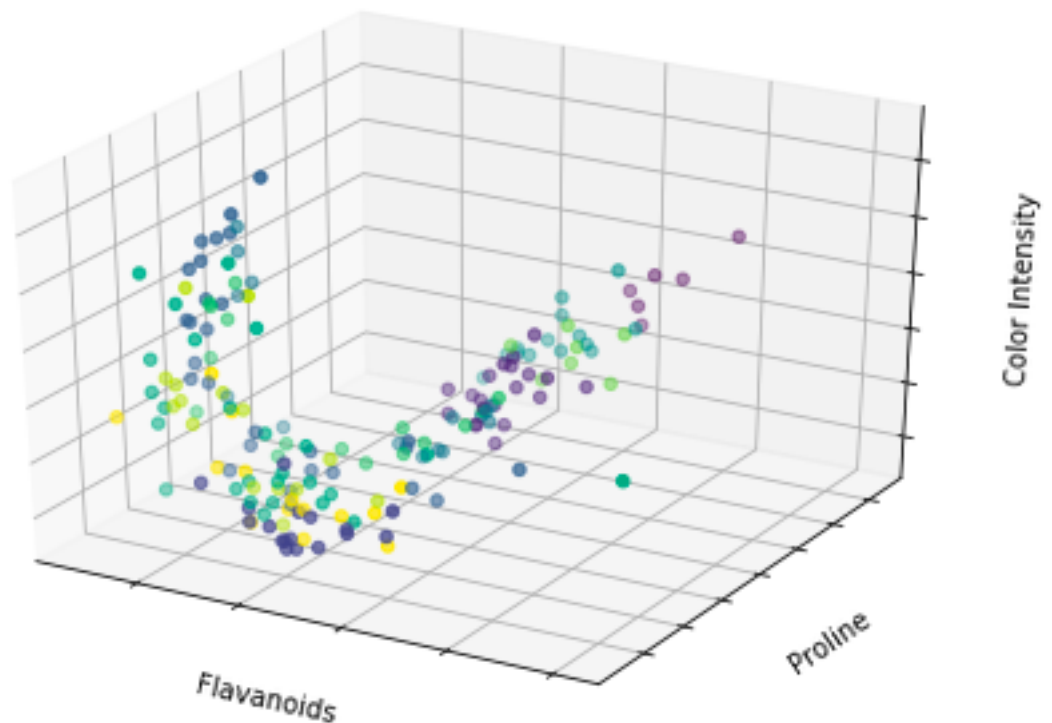
Running Birch method against 11 clusters for the top 3 features as determined above. I start the timer before the fitting function and stop after the fitting function. Running the algorithm 5 times:

n clusters = 11

| Run 1 | 0. 007426388008752838 |
|-------|------------------------|
| Run 2 | 0. 008581353002227843 |
| Run 3 | 0. 009152141021331772 |
| Run 4 | 0. 008396883000386879 |
| Run 5 | 0. 008861561014782637 |

The average running time is about **8.48 milliseconds** to fit 5 cluster. The running time for 11 clusters is slightly slower than running time for 5 clusters, however, the difference is still minimal. As I discussed from previous iteration, BIRCH only requires a single scan of data to build the CF tree in memory. Comparing to K-Means, Birch is much faster and efficient because it avoids re-examining all the points.
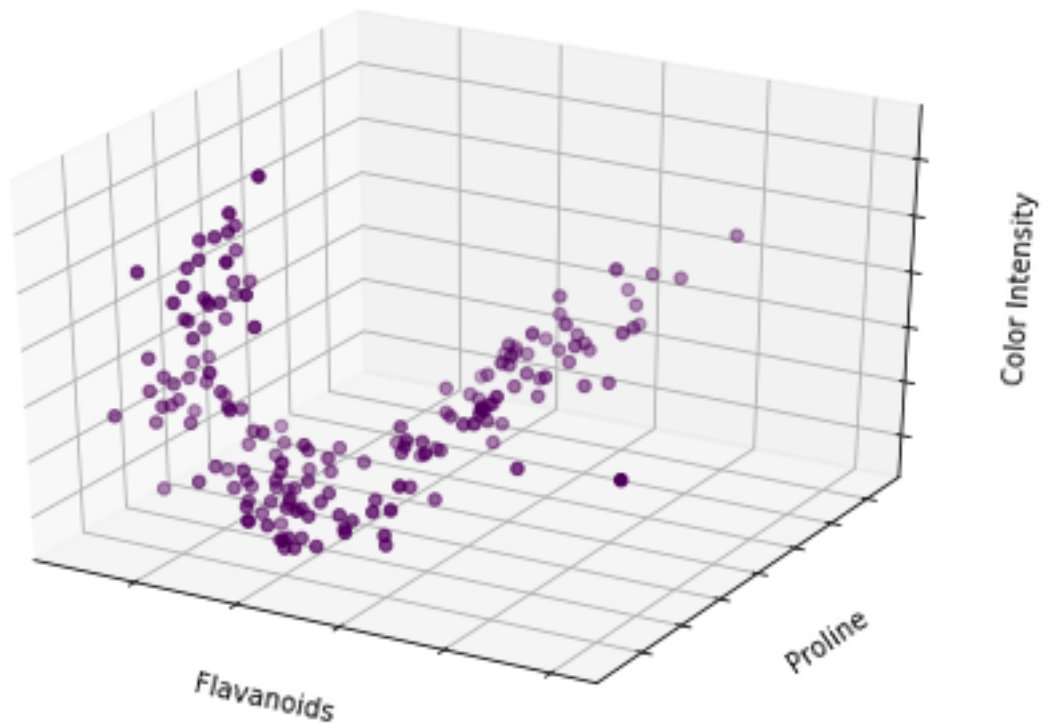
Agglomerative Clustering

Running Agglomerative Clustering method against 1 clusters for the top 3 features as determined above. I start the timer before the fitting function and stop after the fitting function. Running the algorithm 5 times:

n clusters = 1

| Run 1 | 0. 0007832249975763261 |
| Run 2 | 0. 000833658006740734 |
| Run 3 | 0. 0008580129942856729 |
| Run 4 | 0. 0008462440164294094 |
| Run 5 | 0. 0008289699908345938 |

The average running time is about **8.48 milliseconds** to fit 1 cluster. Comparing to the same method running at 3 clusters, it is very comparable.
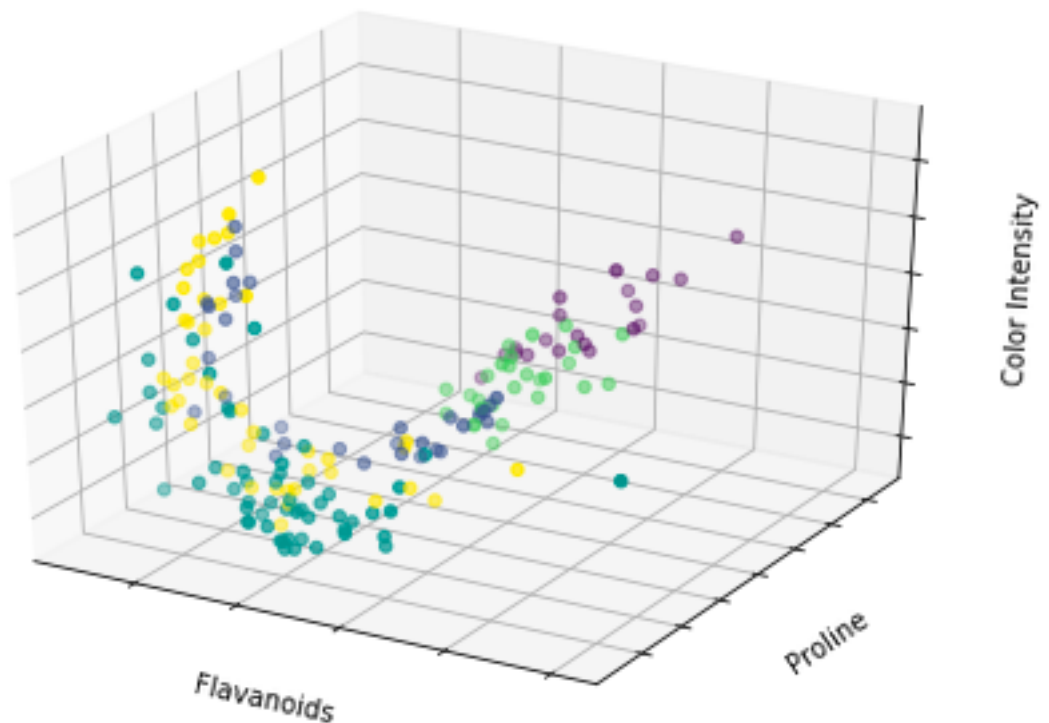
Running Agglomerative Clustering method against 5 clusters for the top 3 features as determined above. I start the timer before the fitting function and stop after the fitting function. Running the algorithm 5 times:

n clusters = 5

| Run 1 | 0. 0008039359818212688 |
|-------|------------------------|
| Run 2 | 0. 0008217980212066323 |
| Run 3 | 0. 0008288479875773191 |
| Run 4 | 0. 0008330359996762127 |
| Run 5 | 0. 0011607060150709003 |

The average running time is about **0.8895 milliseconds** to fit 5 cluster. Comparing the running time against fitting 1 cluster above, the difference is very minimal. Intuitively, this makes sense because Agglomerative clustering works by having each point as its own cluster initially and find a pair of most similar clusters and merge basing on their distance (distance between two clusters). Furthermore, Agglomerative Clustering handles non-convex clustering better than K-Means.
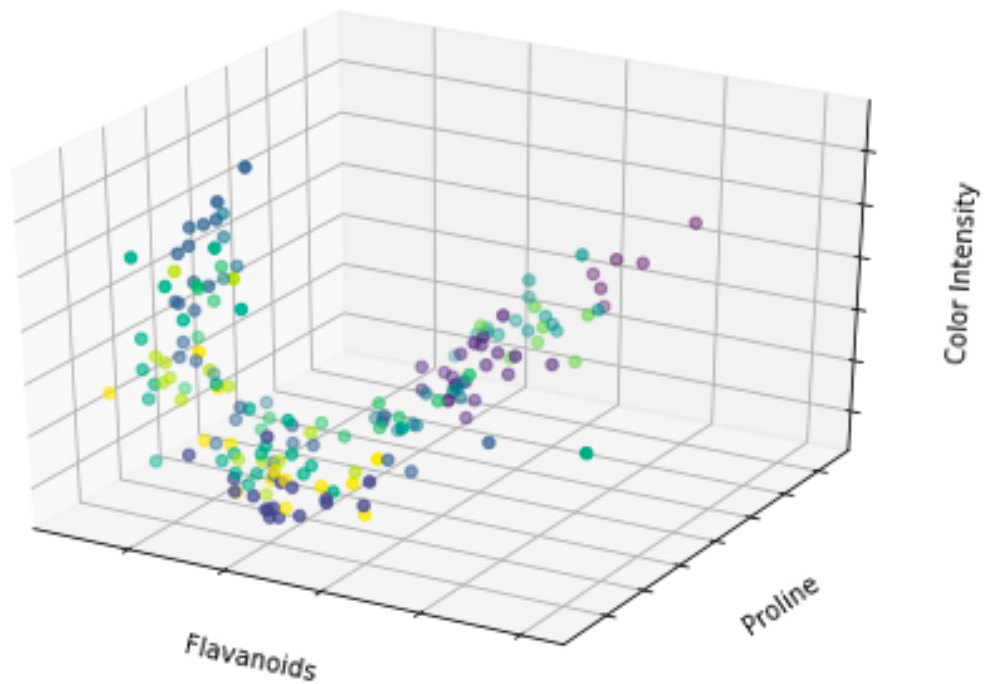
Running Agglomerative Clustering method against 11 clusters for the top 3 features as determined above. I start the timer before the fitting function and stop after the fitting function. Running the algorithm 5 times:
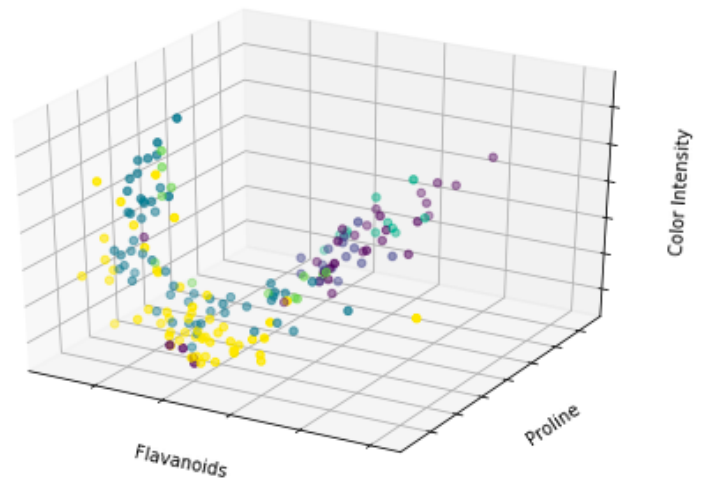
n clusters = 11

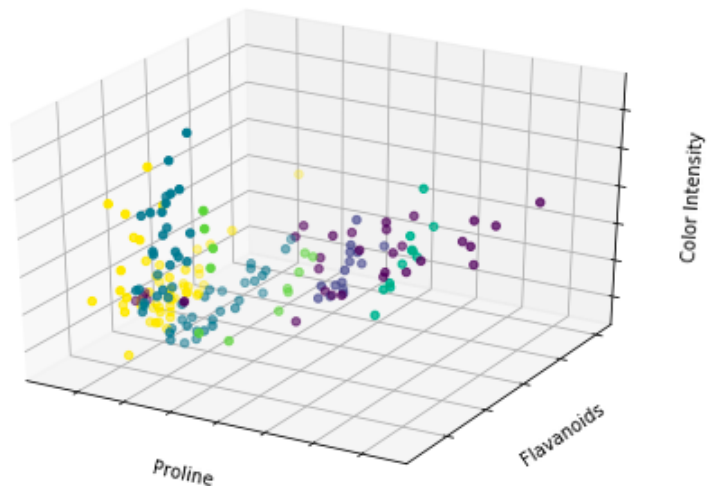| Run 1 | 0. 0008172669913619757 |
|---|---|
| Run 2 | 0. 0008151039946824312 |
| Run 3 | 0. 0008365200192201883 |
| Run 4 | 0. 0008419170044362545 |
| Run 5 | 0. 0008454299822915345 |

The average running time is about **0.83124 milliseconds** to fit 11 cluster. Again, the difference in average running time to fit 11 clusters is very minimal.

4) The nice thing about DBSCAN clustering method is that it does not require a defined number of clusters like in K-Means, Birch, and Agglomerative Clustering methods. However, it does require epsilon (radius from each point to neighbors) and the minimum neighbors in order to form clusters. However, choosing the correct epsilon and minimum sample can be difficult and important since it determines the number of clusters generated in the model. DBSCAN can also discover any arbitrary shapes.



The rule of thumb of choosing the minimum points can be derived from the number of dimensions D in the dataset, minPts >= D + 1. Since we have 3 dimensions in this dataset, any minimum sample above 4 would be ideal. Using Euclidean metric in this dataset, I have concluded that an epsilon value of 20, and minimum sample of 5 do generate the same 5 clusters 3D graph as above. Comparing the 3-D graphs from DBSCAN, Birch, and Agglomerative Clustering methods, I believe that having 5 clusters in this dataset is the most optimized number of clusters.

5) In conclusion, I used random forest method to determine the top three important features of Wine dataset which are Flavanoids, Proline, and Color Intensity. Using the top three features, I ran against popular clustering methods such as K-Means, Birch, Agglomerative Clustering, and DBSCAN.

K-Means clustering method is the first method that is used to analyze for clusters because of its simplicity and easy to use. However, K-means do have some shortfalls such that it does not detect outliers, and it does not work well with non-convex shapes. In this dataset, which contains a nonconvex shape as seen in the graphs above, K-means method was unable to cluster the shape correctly. Furthermore, K-means does not scale with large datasets. Looking at the running time comparison above, having more clusters would take longer running time since K-Means needs compare the distance again for every point against the mean center of clusters until all points are satisfied with their mean cluster centers.

BIRCH clustering method is more efficient than K-means in terms of running time. Since Birch only requires one data scan to build a CF tree and it does not require data all at once, this can reduce the running time greatly. However, there might be some running cost if the CF Tree requires splitting but this is minimal because the splitting would be done in memory. If new point is inserted, the algorithm can use depth search traverse to find the most optimal cluster to fit the point in thus avoiding having unnecessary rescans of data like in K-Means. It can also scale well with larger dataset. However, BIRCH do have some downfalls. The size of each node in CF tree can hold only a limited number of entries due to the size, and the method doesn't perform well if the shape is not spherical because it uses radius or diameter to control the boundary of a cluster.

Agglomerative Clustering method seems to be the fastest among K-Means and Birch, however it does not scale very well with large dataset due to its rather slow time complexity O(n^2) Agglomerative Clustering can produce an ordering of data points, which maybe informative for data display, and smaller clusters are generated which may be helpful for discovery. However, data points may be incorrectly grouped at early stage and there is no way to backtrack. Therefore, the result should be examined with closely consideration to ensure that the clustering makes sense. Furthermore, Agglomerative Clustering uses the distance between clusters to merge and the use of different distance metrics for measuring may generate different results. Therefore, Agglomerative may require multiple experiments with different metrics.

DBSCAN does scale well with large datasets and able to discover clusters with any arbitrary shape. It does not require a specific number of clusters, and is robust to outliers. DBSCAN requires two parameters, epsilon (the radius from the core data point), and minimum number of neighbors. However, choosing a meaningful distance threshold epsilon can be difficult. DBSCAN cannot cluster datasets with large differences in densities since the minimum of neighbors and epsilon combination cannot be chosen appropriately for all clusters. A better method similar to DBSCAN is OPTICS which

epsilon only serves the upper limit of the radius. The radius in OPTICS would be based on how dense the number of neighbors defined by the minimum neighbors given. The OPTICS method allows extraction of clusters with any arbitrary epsilon by calculating the core distance and reachability distance.