

Final

Lecture - 1  
Math

20-06-2022  
Monday

FP Growth Algorithm

① Works with through divide and conquer approach.

Min sup = 3 W

Header Table

Step-1

f, a, c, d, g, i, m, p, b,  
l, o, h, j, w, k, s, e, n

Item	Frequency
f	4
a	3
c	4
d	1
g	1
i	1
m	3
p	3
b	3
l	2
o	2
h	1
j	1
w	1
k	1
s	1
<del>e</del>	<del>1</del>
e	1
h	1



L

Item	Frequency
f	4
c	4
a	3
b	3
m	3
p	3

↑ Frequent  
Pattern

(Sorted in descending  
order) W

\* f, c, a, b, m, p (order)  
for next step.



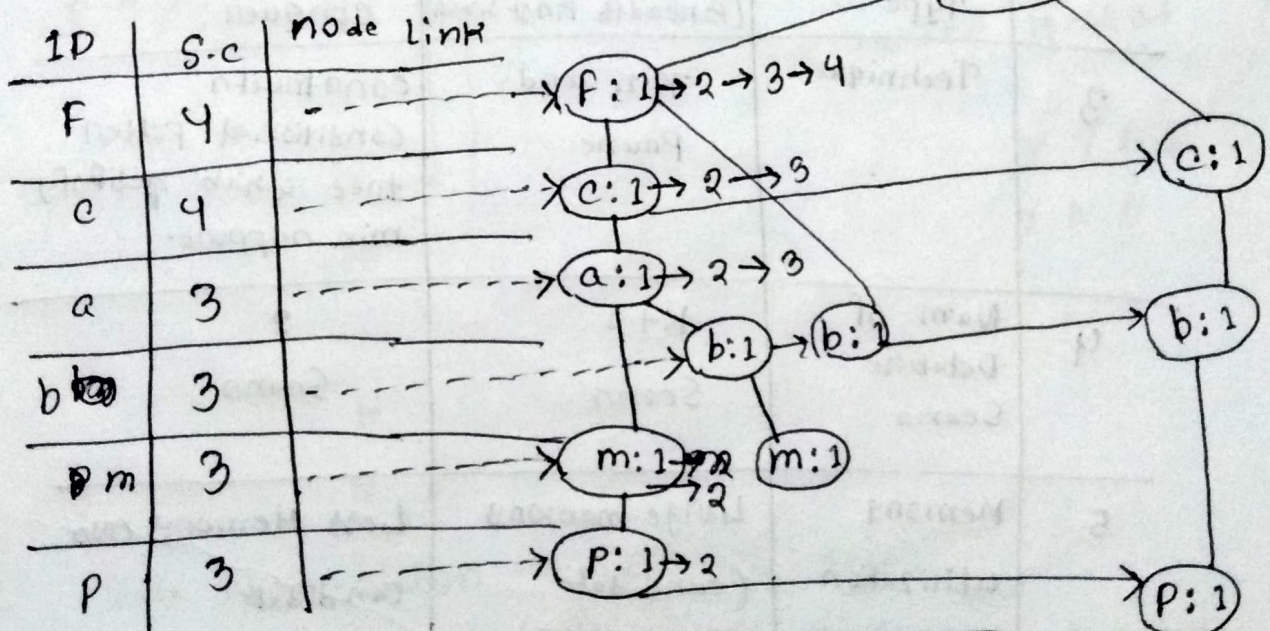
## Step - 2

f, c, a, b, m, p

ID	Item	Ordered Frequent items	Item	Conditional Pattern base
T <sub>1</sub>	f, a, c, d, g, i, m, p	f, c, a, m, p	p	{f, c, a, m: 2}, {c, b: 1}
T <sub>2</sub>	a, b, c, f, l, m, o	f, c, a, b, m	m	{f, c, a: 2}, {f, c, a, b: 1}
T <sub>3</sub>	b, f, h, j, o, w	f, b	b	{f, c, a: 1}, {f: 1}, {c: 1}
T <sub>4</sub>	b, k, c, s, p	c, b, p	a	{f, c: 3}
T <sub>5</sub>	a, f, c, e, l, p, m, n	f, c, a, m, p	c	{f: 3}

Don't need to write it for ~~on~~ (exam).

## Step - 3



Tree Generation

f → c → a → b → m → p

Am:

# Decision Tree

Ex: 8.1 (Dataset - Table 8.1) Book page - 338

Sol<sup>n</sup>: Formula: Entropy( $E$ ) =  $-\sum_j P(j|E) \log_2 P(j|E)$   
 $\rightarrow$  percentage of distribution

Decision column  $\rightarrow$  class: buys-computer

Yes - 9  
 No - 5 } 14

$S = [9+, 5-]$

Information gain

$$\begin{aligned} \text{Info}(D) &= -\frac{9}{14} \log_2 \left(\frac{9}{14}\right) - \frac{5}{14} \log_2 \left(\frac{5}{14}\right) \\ &= -0.64 \times (-0.64) - 0.36 \times (-1.48) \\ &= 0.99 \text{ bits} \end{aligned}$$

Main  
 Entropy

Now, calculating the expected information Re-  
 requirement for each attributes:

For "age":

Youth = 5

middle-aged = 4

Senior = 5

$$\begin{aligned}
 \text{Info}_{\text{age}}(D) &= \frac{5}{14} \times \left( -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) + \\
 &\quad \frac{4}{14} \times \left( -\frac{4}{4} \log_2 \frac{4}{4} \right) + \frac{5}{14} \times \left( -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right) \\
 &= 0.36 \times (0.0533 + 0.44) + 0.28 \times 0 + \\
 &\quad 0.36 \times (0.44 + 0.53) \\
 &= 0.69 \text{ bits}
 \end{aligned}$$

Now,  $\text{Gain}(\text{age}) = \text{Info}(D) - \text{Info}_{\text{age}}(D)$

Now  $= 0.94 - 0.69$

$= 0.25 \text{ bits}$  (overall information gain with respect to age) ✓

For "Income":

High = 4, medium = 6, low = 4

$$\begin{aligned}
 \text{Info}_{\text{income}}(D) &= \frac{4}{14} \times \left( -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} \right) + \\
 &\quad \frac{6}{14} \times \left( -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} \right) + \\
 &\quad \frac{4}{14} \times \left( -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right) \\
 &= 0.28 \times (0.5 + 0.5) + 0.43 (0.39 + 0.53) \\
 &\quad + 0.28 (0.31 + 0.5) \\
 &= 0.902 \text{ bits}
 \end{aligned}$$

Now,  $\text{Gain}(\text{age}) = \text{Info}(D) - \text{info}_{\text{income}}(D)$

$= 0.94 - 0.902$

$= 0.038$

$= 0.03 \text{ bits}$



For "Student":

yes = 7, no = 7

$$\begin{aligned}\text{Info}_{\text{student}}(D) &= \frac{7}{14} \times \left( -\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} \right) + \\ &\quad \frac{7}{14} \times \left( -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} \right) \\ &= 0.5 (0.19 + 0.40) + 0.5 (0.52 + 0.46) \\ &= 0.785 \text{ bits}\end{aligned}$$

$$\begin{aligned}\text{Now, Gain}(\text{Student}) &= \text{Info}(D) - \text{Info}_{\text{student}}(D) \\ &= 0.94 - 0.785 \\ &= 0.15 \text{ bits}\end{aligned}$$

For "Credit-rating":

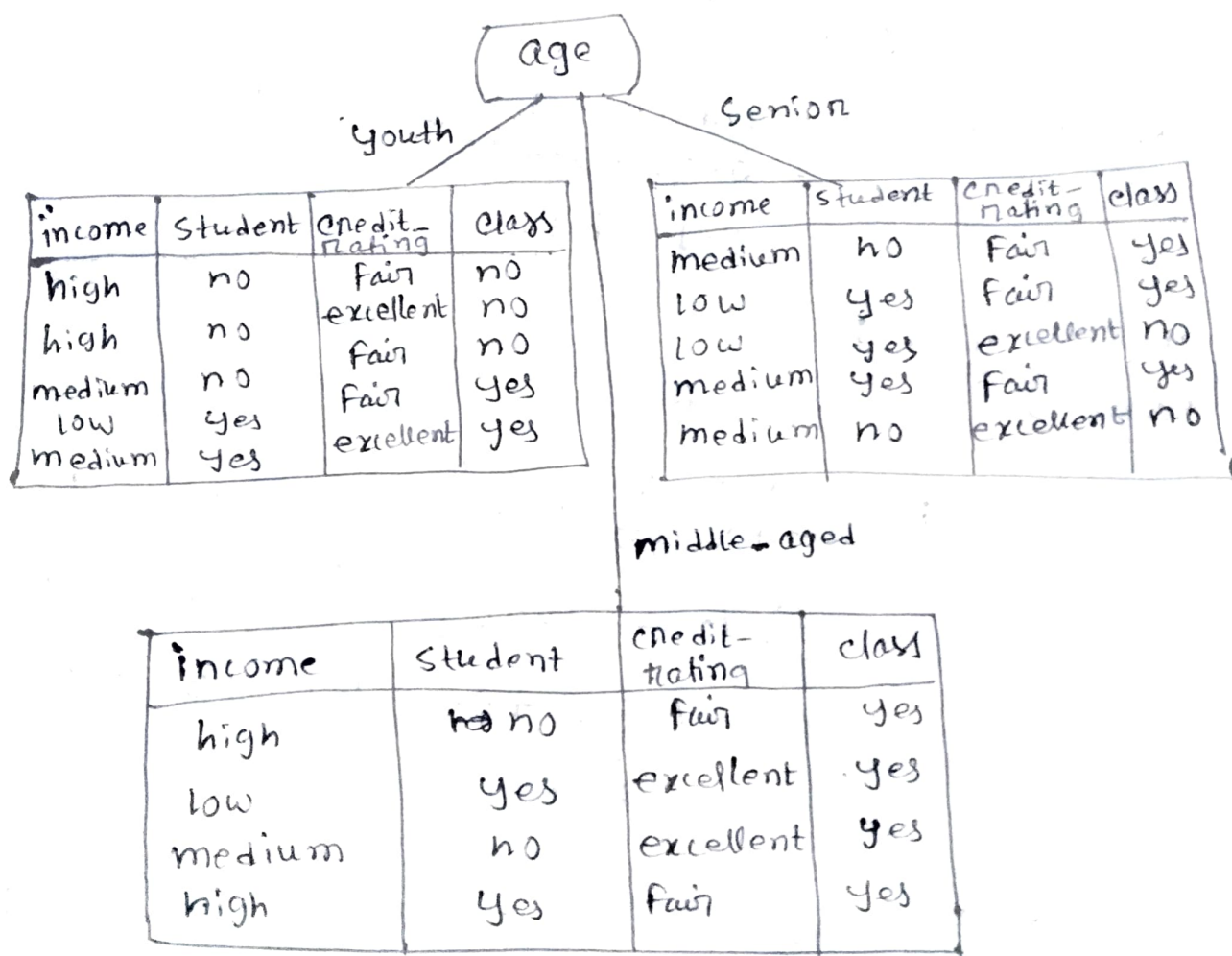
Fair = 8, excellent = 6

$$\begin{aligned}\text{Info}_{\text{credit-rating}}(D) &= \frac{8}{14} \left( -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \right) + \\ &\quad \frac{6}{14} \left( -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} \right) \\ &= 0.57 (0.31 + 0.5) + 0.43 (0.5 + 0.5) \\ &= 0.89 \text{ bits}\end{aligned}$$

$$\begin{aligned}\text{Now, Gain}(\text{Credit-rating}) &= \text{Info}(D) - \text{Info}_{\text{credit-rating}}(D) \\ &= 0.94 - 0.89 \\ &= 0.05 \text{ bits}\end{aligned}$$

We can see "age" has the highest information gain among the attributes, it is selected as the splitting attribute.

### Decision Tree



Gain ratio

Formula: 1)  $\text{Split Info}_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right)$

2)  $\text{Gain Ratio}(A) = \frac{\text{Gain}(A)}{\text{Split Info}_A(D)}$

Ex: 8.2 computation of gain ratio for the attribute "income".

$$\begin{aligned} \text{Split Info}_{\text{income}}(D) &= - \frac{4}{14} \times \log_2 \left( \frac{4}{14} \right) - \frac{6}{14} \times \log_2 \left( \frac{6}{14} \right) \\ &\quad - \frac{4}{14} \times \log_2 \left( \frac{4}{14} \right) \\ &= 0.52 + 0.52 + 0.52 \\ &= 1.56 \end{aligned}$$

From "ex: 8.1" we found  $\text{Gain}(\text{income}) = 0.03$ .

Therefore,  $\text{Gain Ratio}(\text{income}) = \frac{0.03}{1.56} = 0.019$  Ans:

For "age":

$$\begin{aligned} \text{Split Info}_{\text{income}}(\text{age})(D) &= - \frac{5}{14} \times \log_2 \frac{5}{14} - \frac{4}{14} \times \log_2 \left( \frac{4}{14} \right) \\ &\quad - \frac{5}{14} \times \log_2 \frac{5}{14} \\ &= 0.53 + 0.52 + 0.53 \\ &= 1.58 \end{aligned}$$

From "ex: 8.1" we found  $\text{Gain}(\text{age}) = 0.25$

Therefore,  $\text{Gain Ratio}(\text{age}) = \frac{0.25}{1.58} = 0.16$

For "Student":

$$\text{Split Info}_{\text{student}}(D) = -\frac{7}{14} \times \log_2\left(\frac{7}{14}\right) - \frac{7}{14} \times \log_2\left(\frac{7}{14}\right)$$

$$\begin{aligned} &= 0.5 + 0.5 \\ &= 1 \end{aligned}$$

Previously we found  $\text{gain}(\text{student}) = 0.15$

$$\text{Therefore, GainRatio}(\text{student}) = \frac{0.15}{1} = 0.15$$

For "credit-rating":

$$\text{Split Info}_{\text{credit-rating}}(D) = -\frac{8}{14} \times \log_2\left(\frac{8}{14}\right) - \frac{6}{14} \times \log_2\left(\frac{6}{14}\right)$$

$$= 0.46 + 0.52$$

$$= 0.98$$

Previously we found  $\text{gain}(\text{credit-rating}) = 0.05$ .

$$\text{Therefore, GainRatio}(\text{credit-rating}) = \frac{0.05}{0.98}$$

$$= 0.05$$

Ans:



Gini Index

Formulas:

- 1)  $Gini(D) = 1 - \sum_{i=1}^m p_i^2$
- 2)  $Gini_A(D) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2)$
- 3)  $\Delta Gini(A) = Gini(D) - Gini_A(D)$

**Ex: 8.3**

Impurity of D:

$$Gini(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2$$

$$= 0.439$$

Calculating Gini Index for each attribute.

For "income":

possible  
considering each of the splitting subsets. ~~can~~  
consider the subset {low, medium}. This would result  
in 10 tuples in partition  $D_1$  satisfying the condition  
"income"  $\in$  {low, medium}. The remaining four tuples  
of D would be assigned to partition  $D_2$ .

The Gini Index value computed based on this partitioning is:

$$Gini_{income \in \{low, medium\}}(D)$$

$$= \frac{10}{14} Gini(D_1) + \frac{4}{14} Gini(D_2)$$

$$= \frac{10}{14} \left(1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2\right) + \frac{4}{14} \left(1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2\right)$$

$$= 0.443$$

$$= Gini_{income \in \{high\}}(D)$$

Ans:

Subset {low, high}:

$$\text{Gini}_{\text{income}} \{ \text{low, high} \} (D)$$

$$= \frac{8}{14} \text{Gini}(D_1) + \frac{6}{14} \text{Gini}(D_2)$$

$$= \frac{8}{14} \left( 1 - \left( \frac{2}{4} \right)^2 - \left( \frac{2}{4} \right)^2 \right) +$$

$$\frac{6}{14} \left( 1 - \left( \frac{5}{8} \right)^2 - \left( \frac{3}{8} \right)^2 \right) + \frac{6}{14} \left( 1 - \left( \frac{4}{6} \right)^2 - \left( \frac{2}{6} \right)^2 \right)$$

$$= 0.4588$$

$$= \text{Gini}_{\text{income}} \{ \text{medium} \} (D)$$

Subset {medium, high}:

$$\text{Gini}_{\text{income}} \{ \text{medium, high} \} (D)$$

$$= \frac{10}{14} \text{Gini}(D_1) + \frac{4}{14} \text{Gini}(D_2)$$

$$= \frac{10}{14} \left( 1 - \left( \frac{6}{10} \right)^2 - \left( \frac{4}{10} \right)^2 \right) + \frac{4}{14} \left( 1 - \left( \frac{3}{4} \right)^2 - \left( \frac{1}{4} \right)^2 \right)$$

$$= 0.450$$

$$= \text{Gini}_{\text{income}} \{ \text{low} \} (D)$$

From the calculation we can see the best binary split for attribute "income" is on {low, medium} (or {high}). Because it minimizes the Gini index.

{low, medium} ✓  
 {low, high} ✓  
~~{high, medium}~~  
 {medium, high}

For "age":

- ✓ {youth, middle-age}
  - ✓ {youth, senion}
  - ✓ {middle-age, senion}
- Subsets

Subset {youth, middle-age}

$$\begin{aligned} \text{Gini}_{\text{age}} &\in \{\text{youth, middle-age}\} (D) \\ &= \frac{10}{14} \text{Gini}(D_1) + \frac{4}{14} \text{Gini}(D_2) \\ &= \frac{10}{14} \left( 1 - \left(\frac{6}{10}\right)^2 - \left(\frac{4}{10}\right)^2 \right) + \frac{4}{14} \left( 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 \right) \\ &= 0.457 \\ &= \text{Gini}_{\text{age}} \{\text{senion}\} (D) \end{aligned}$$

Subset {youth, senion}

$$\begin{aligned} \text{Gini}_{\text{age}} &\in \{\text{youth, senion}\} (D) \\ &= \frac{10}{14} \text{Gini}(D_1) + \frac{4}{14} \text{Gini}(D_2) \\ &= \frac{10}{14} \left( 1 - \left(\frac{5}{10}\right)^2 - \left(\frac{5}{10}\right)^2 \right) + \frac{4}{14} \left( 1 - \left(\frac{4}{4}\right)^2 \right) \\ &= 0.357 \\ &= \text{Gini}_{\text{age}} \in \{\text{middle-age}\} (D) \end{aligned}$$



Subset {middle age, Senior}

$$\text{Gini}_{\text{age}} \in \{\text{middle age, Senior}\} (D)$$

$$= \frac{9}{14} \text{Gini}(D_1) + \frac{5}{14} \text{Gini}(D_2)$$

$$= \frac{9}{14} (1 - (\frac{4}{9})^2) + \frac{5}{14} (1 - (\frac{3}{5})^2 - (\frac{2}{5})^2)$$

$$= \frac{9}{14} (1 - (\frac{7}{9})^2 - (\frac{2}{9})^2) + \frac{5}{14} (1 - (\frac{2}{5})^2 - (\frac{3}{5})^2)$$

$$= 0.3036$$

$$= \text{Gini}_{\text{age}} \in \{\text{youth}\} (D)$$

Let's consider "Student":

binary attribute

$$\text{Gini}_{\text{student}}(D) = \frac{7}{14} (1 - (\frac{6}{7})^2 - (\frac{1}{7})^2) + \frac{7}{14} (1 - (\frac{4}{7})^2)$$

$$= 0.3673$$

Similarly,  $\text{Gini}_{\text{credit rating}}(D) = 0.4285$

Reduction of impurity:

$$\text{age} \rightarrow 0.459 - 0.357 = 0.102$$

$$\text{income} \rightarrow 0.459 - 0.443 = 0.016$$

$$\text{age} \rightarrow 0.457 - 0.357 = 0.1$$

$$\text{income} \rightarrow 0.458 - 0.443 = 0.02$$

$$\text{Student} \rightarrow$$

$$\text{age} \rightarrow 0.459 - 0.357 = 0.102$$

$$\text{income} \rightarrow 0.459 - 0.443 = 0.016$$

$$\text{Student} \rightarrow 0.459 - 0.367 = 0.092$$

$$\text{Credit-rating} \rightarrow 0.459 - 0.428 = 0.031$$

Ans

## K-means clustering

→ Division - cluster

■ Supervised Algo - Divide data into different category.

■ Unsupervised Algo - Divide data into different category.

↳ Division - classification / class

• But there is functional difference → working procedure difference.

Example:

Num of obj = 6, Num of clusters = 2

Q:

No	X	Y
1	1	1
2	2	3
3	1	2
✓ 4	3	3
✓ 5	2	2
6	3	1

} chosen points

Soln:

Step-1: choose random k points and set as cluster centers.

$$C_1 = (2, 2)$$

$$C_2 = (3, 3)$$

Step-2: calculating the distance between ~~objects~~ objects into cluster centroids by using Euclidean Distance.

$$D_i = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

C1: ~~(2,2)~~ Fon (2,2) w

$$\begin{aligned} 1. D_1 &= \{(1,1), (2,2)\} \\ &= \sqrt{(2-1)^2 + (2-1)^2} \\ &= \sqrt{1+1} \\ &= 1.41 \end{aligned}$$

$$\begin{aligned} 2. D_{2_1} &= \{(2,3), (2,2)\} \\ &= \sqrt{(2-2)^2 + (2-3)^2} \\ &= \sqrt{0+1} \\ &= \underline{1} \end{aligned}$$

$$\begin{aligned} 3. D_1 &= \{(1,2), (2,2)\} \\ &= \sqrt{(2-1)^2 + (2-2)^2} \\ &= \sqrt{1+0} \\ &= 1 \end{aligned}$$

$$\begin{aligned} 4. D_1 &= \{(3,3), (2,2)\} \\ &= \sqrt{(2-3)^2 + (2-3)^2} \\ &= \sqrt{1+1} \\ &= 1.41 \end{aligned}$$

$$\begin{aligned} 5. D_1 &= \{(2,2), (2,2)\} \\ &= \sqrt{(2-2)^2 + (2-2)^2} \\ &= 0 \end{aligned}$$

$$\begin{aligned} 6. D_1 &= \{(3,1), (2,2)\} \\ &= \sqrt{(2-3)^2 + (2-1)^2} \\ &= 1.41 \end{aligned}$$

C2: Fon (3,3) w

$$\begin{aligned} 1. D_2 &= \{(1,1), (3,3)\} \\ &= \sqrt{(3-1)^2 + (3-1)^2} \\ &= 2.82 \end{aligned}$$

$$\begin{aligned} 2. D_2 &= \{(2,3), (3,3)\} \\ &= \sqrt{(3-2)^2 + (3-3)^2} \\ &= 1 \end{aligned}$$

$$\begin{aligned} 3. D_2 &= \{(1,2), (3,3)\} \\ &= \sqrt{(3-1)^2 + (3-2)^2} \\ &= \sqrt{5} = 2.23 \end{aligned}$$

$$\begin{aligned} 4. D_2 &= \{(3,3), (3,3)\} \\ &= \sqrt{(3-3)^2 + (3-3)^2} \\ &= 0 \end{aligned}$$

$$\begin{aligned} 5. D_{2_2} &= \{(2,2), (3,3)\} \\ &= \sqrt{(3-2)^2 + (3-2)^2} \\ &= 1.41 \end{aligned}$$

$$\begin{aligned} 6. D_2 &= \{(3,1), (3,3)\} \\ &= \sqrt{(3-3)^2 + (3-1)^2} \\ &= 2 \end{aligned}$$

$$C_1 = \{(1,1), (1,2), (2,2), (3,1)\}$$

$$C_2 = \{(2,3), (3,3)\}$$

Note: For same value we can  
take as per our wish.  
→ (2,2)



Step-3: Recalculating the position of the centroid.

$$\text{Mean} = \left( \frac{x_1 + x_2 + \dots + x_n}{n}, \frac{y_1 + y_2 + \dots + y_n}{n} \right)$$

$$c_1 = \left( \frac{1 + 1 + 2 + 3}{4}, \frac{1 + 2 + 2 + 1}{4} \right)$$

$$\text{New } c_1 = (1.75, 1.5)$$

$$c_2 = \left( \frac{2+3}{2}, \frac{3+3}{2} \right) = (2.5, 3)$$

$$\Rightarrow \text{New } c_2 = (2.5, 3)$$

Step-4: Go back to Step 2, unless the centroids are not changing.

$c_1: \text{Fon } (1.75, 1.5)$

$$\begin{aligned} D_1 &= \sqrt{\{(1, 1), (1.75, 1.5)\}} \\ &= \sqrt{(1.75-1)^2 + (1.5-1)^2} \\ &= 0.90 \end{aligned}$$

$$\begin{aligned} D_1 &= \{(2, 3), (1.75, 1.5)\} \\ &= \sqrt{(1.75-2)^2 + (1.5-3)^2} \\ &= 1.52 \end{aligned}$$

$$\begin{aligned} D_1 &= \{(1, 2), (1.75, 1.5)\} \\ &= \sqrt{(1.75-1)^2 + (1.5-2)^2} \\ &= 0.90 \end{aligned}$$

$c_2: \text{Fon } (2.5, 3)$

$$\begin{aligned} D_2 &= \{(1, 1), (2.5, 3)\} \\ &= \sqrt{(2.5-1)^2 + (3-1)^2} = \sqrt{(2.5-1)^2 + (3-1)^2} \\ &= 2.5 \end{aligned}$$

$$\begin{aligned} D_2 &= \{(2, 3), (2.5, 3)\} \\ &= \sqrt{(2.5-2)^2 + (3-3)^2} \\ &= 0.5 \end{aligned}$$

$$\begin{aligned} D_2 &= \{(1, 2), (2.5, 3)\} \\ &= \sqrt{(2.5-1)^2 + (3-2)^2} \\ &= 1.80 \end{aligned}$$

$$D_1 = \{(3,3), (1.75, 1.5)\}$$

$$= \sqrt{(1.75-3)^2 + (1.5-3)^2}$$

$$= 1.95$$

$$D_1 = \{(2,2), (1.75, 1.5)\}$$

$$= \sqrt{(1.75-2)^2 + (1.5-2)^2}$$

$$= 0.55$$

$$D_1 = \{(3,1), (1.75, 1.5)\}$$

$$= \sqrt{(1.75-3)^2 + (1.5-1)^2}$$

$$= 1.34$$

$$D_2 = \{(3,3), (2.5, 3)\}$$

$$= \sqrt{(2.5-3)^2 + (3-3)^2}$$

$$= 0.5$$

$$D_2 = \{(2,2), (2.5, 3)\}$$

$$= \sqrt{(2.5-2)^2 + (3-2)^2}$$

$$= 1.11$$

$$D_2 = \{(3,1), (2.5, 3)\}$$

$$= \sqrt{(2.5-3)^2 + (3-1)^2}$$

$$= 2.06$$

$$\therefore C_1 = \{(1,1), (1,2), (2,2), (3,1)\}$$

$$C_2 = \{(2,3), (3,3)\} \text{ ~~(3,1)~~}$$

So, we can see the result of the two iteration are same ~~as~~ i.e. clustering outcome are same.