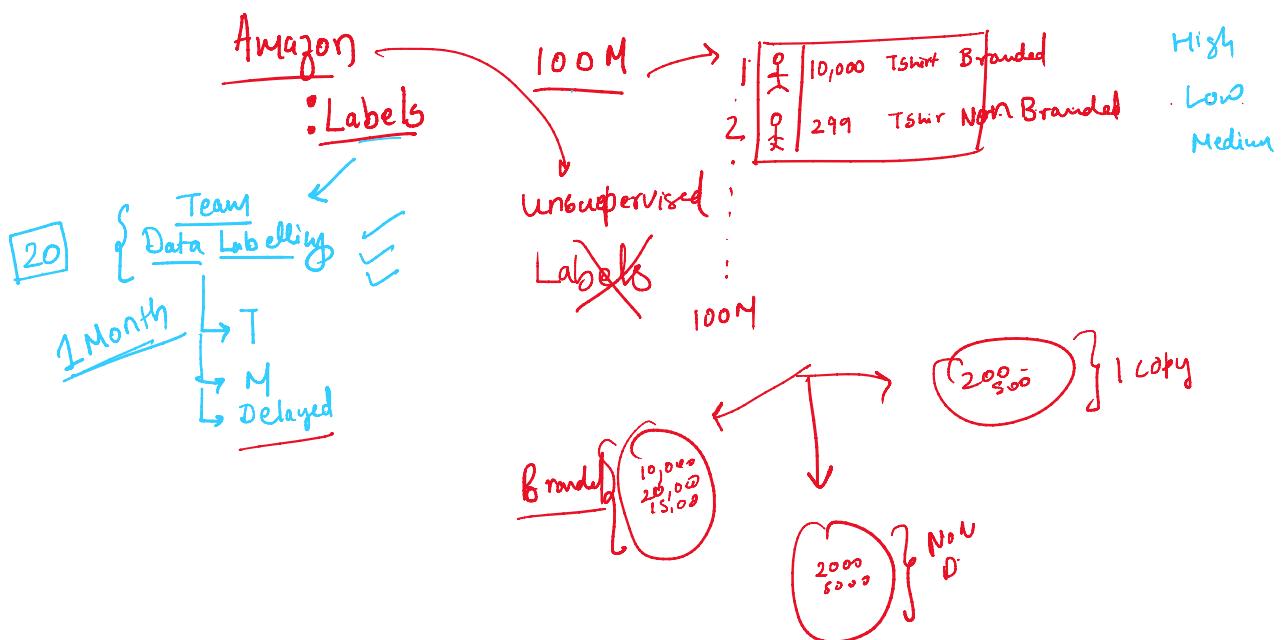
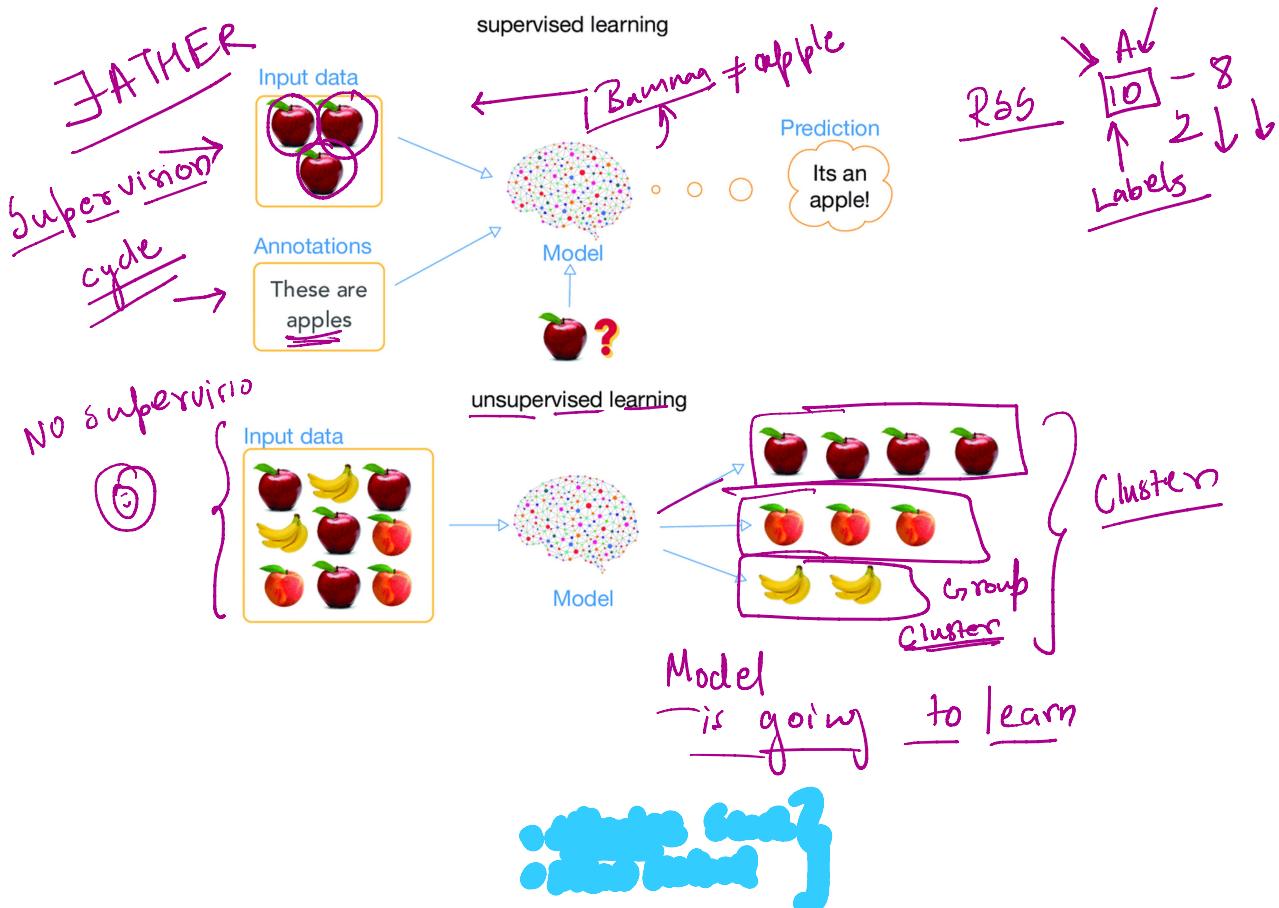
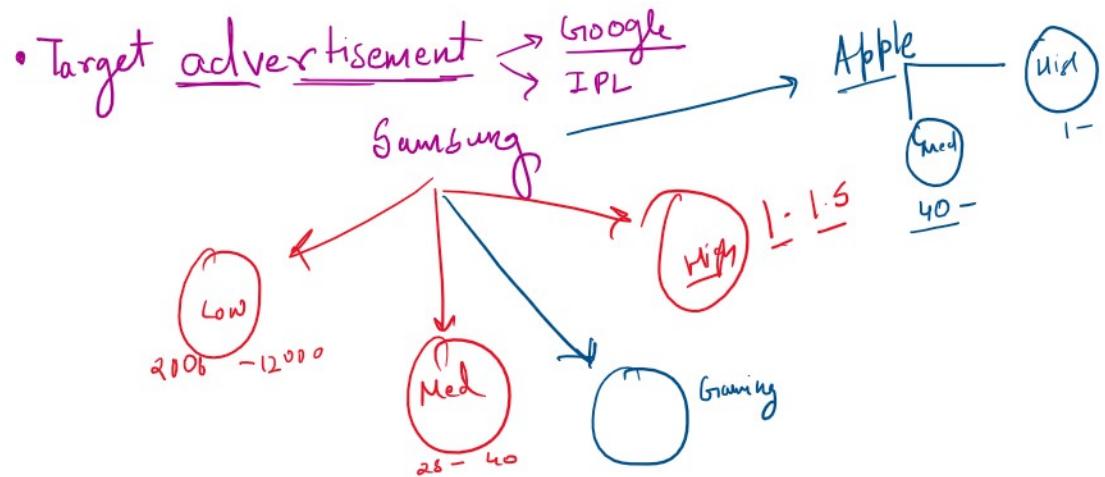
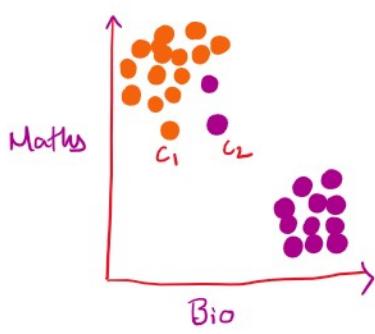


Difference → Supervised and Unsupervised
Labels without Labels

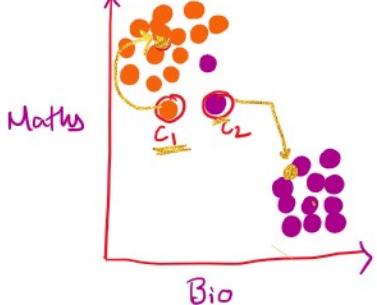




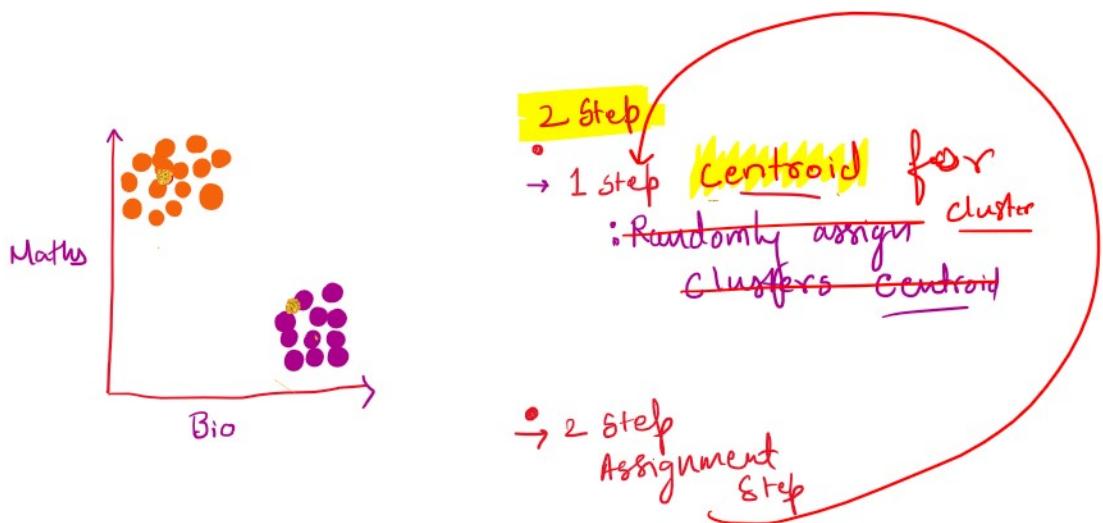
K Means



- 2 Step
 - 1 step
: Randomly assign clusters Centroid
- \Rightarrow 2 Step Assignment Step

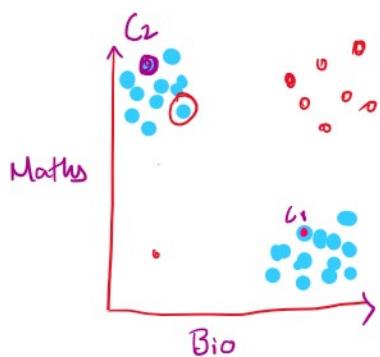


- 2 Step
 - 1 step Centroid for cluster
: Randomly assign clusters Centroid
- \Rightarrow 2 Step Assignment Step



- Pick the first centroid point (C_1) randomly.
- Compute distance of all points in the dataset from the selected centroid. The distance of x_i point from the farthest centroid can be computed by

From <<https://towardsdatascience.com/understanding-k-means-k-means-and-k-medoids-clustering-algorithms-ad9c9fbf47ca>>



1. Randomly Select 1 Point from the given Point
- ② dist is higher highest would be your next centroid
- ③ calc the dist to iB nearest centroid

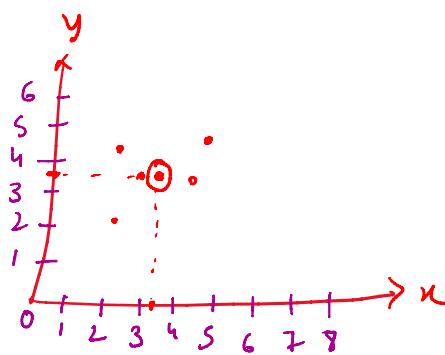
$$E_C = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

x_1, y_1 x_2, y_2

Mean

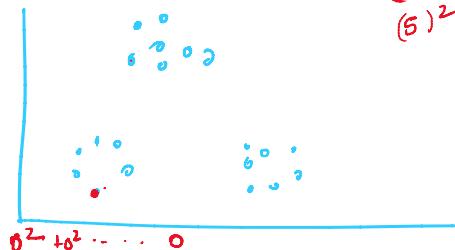
$$\begin{cases} x = \frac{2+2+3+4+5+5}{6} \\ = 3.3 \end{cases}$$

$$y = \frac{2+3+5+4+5+3+5+4+7}{9} \\ = 3.6$$



K Means
Centroid
↓
No of center

$SSD_{C_1} + SSD_{C_2} > V_2$



$$\frac{(10)^2}{(5)^2}$$

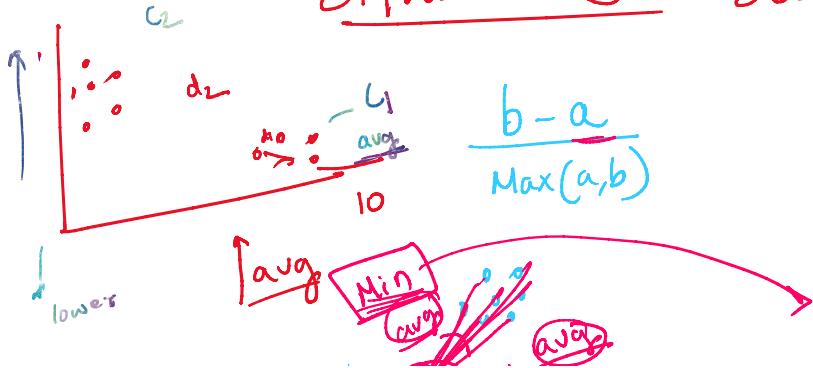
Elbow → sum of square distance
Cluster $K=1 \rightarrow SSD = V_1$
 $K=2 \rightarrow SSD = V_2$
 $K=3 \vdots$
 $SSD = 0$

$$\sum d_1^2 + d_2^2 \dots d_n^2 \quad K=n$$

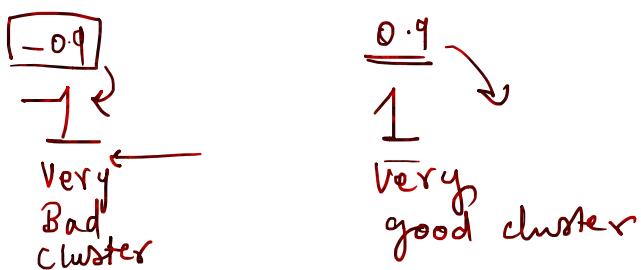
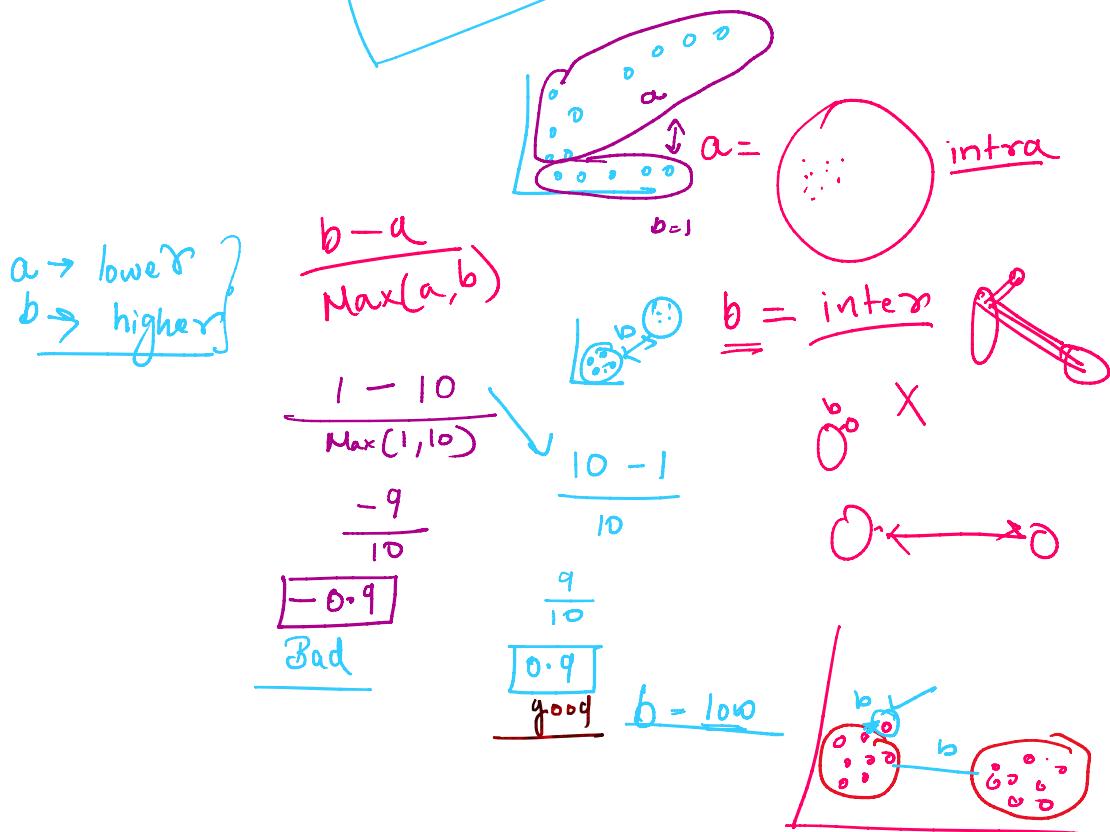
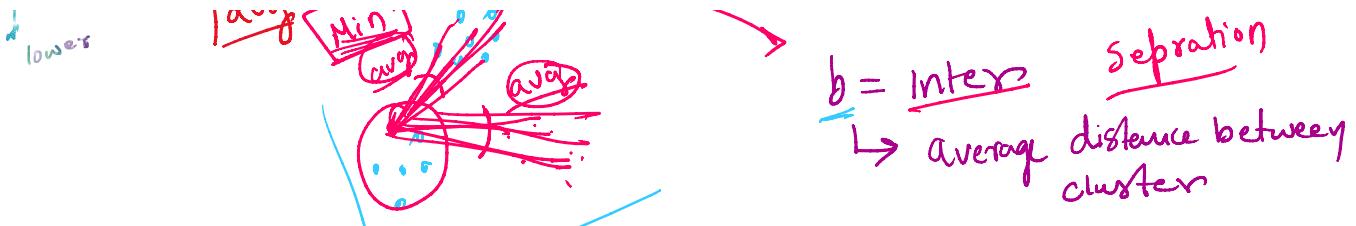
{ Silhouette } → PCA LDA T5
Hierarchical

Association Rules

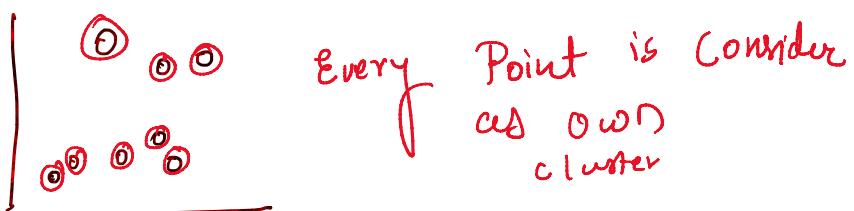
Silhouette Score

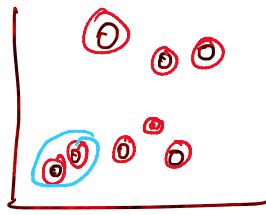


$a = \text{avg dist within}$
cluster cohesion
 $b = \text{inter}$ separation

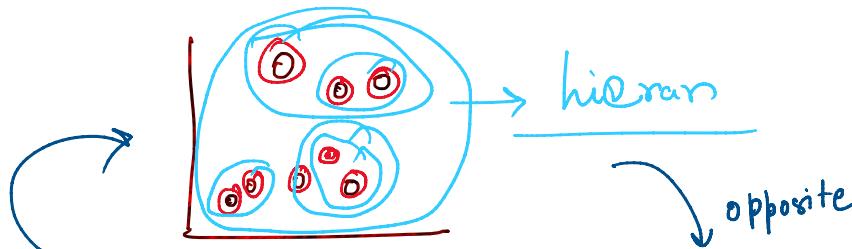


Hierarchical Clustering

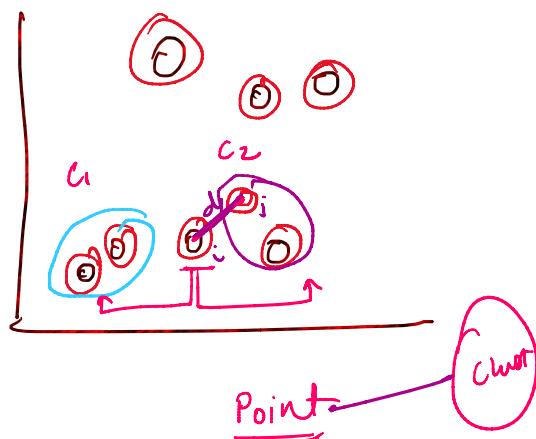




You Merge the
Cluster which
are closer



* Agglomerative *



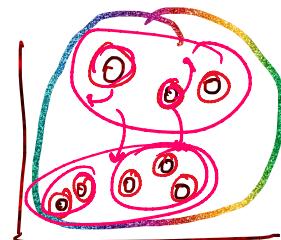
distance between
a Point and a
cluster } Linkage

Single linkage ($\min(d_1, d_2)$)

Average linkage ($\text{avg}(d_1, d_2)$)

complete linkage ($\max(d_1, d_2)$)

Divisive ⇒ All Points
belong
to
one
clust.



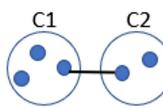
Dividing the
cluster

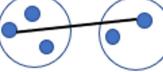
complete linkage ($\max(d_1, d_2)$)

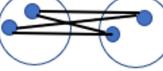
Ward
$$\frac{\sum \text{Sim}(P_i, P_j)}{|C_1| |C_2|}$$

$$\frac{d_1 + d_2 + d_3}{3}$$



Single Linkage  $D(C1, c2) = \min D(C1, C2)$
Minimum distance between data points in clusters

Complete Linkage  $D(C1, c2) = \max D(C1, C2)$
Maximum distance between data points in clusters

Average Linkage  $D(C1, c2) = \text{Average distances of all pairs in clusters}$

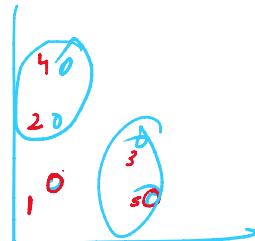
Centroid method  Minimum distance between centroids of the clusters

$$\left\{ \frac{d_1^4 + d_2^2 + d_3^3}{3} \right\}$$

Ward's method 

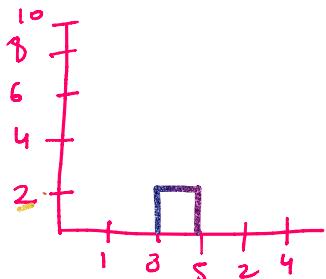
Minimum variance to minimize the total within-cluster variance

	1	2	3	4	5
1	0	-			
2	9	0			
3	3	7	0		
4	6	5	9	0	
5	11	10	2	8	0



1 step
find the closest points

	1	2	3	4	5
1	0	-			
2	9	0			
3	3	7	0		
4	6	5	9	0	
5	11	10	2	8	0



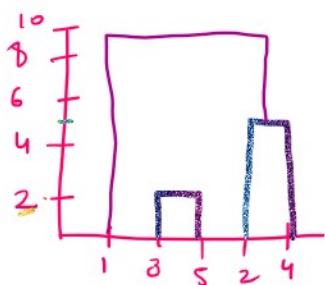
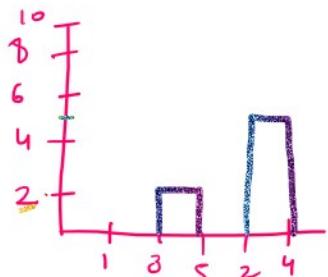
	3	5	1	2	4
3.5	0				

2nd dist
of
Point and
cluster.

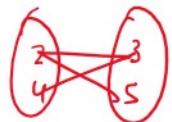
	35	1	2	4
35	0			
1	11	0		
2	10	9	0	
4	9	6	5	0

5 step

	35	1	2	4
35	0			
1	11	0		
2	10	9	0	
4	9	6	5	0



	35	24	1
35	0		
24	10	0	
1	11	9	0



$$\begin{cases} 23 - 7 \\ 25 - 10 \\ 43 - 9 \\ 45 - 8 \end{cases}$$

Association Rules - A priori

Id	Name of Product
1	I ₁ I ₂ I ₃
2	I ₂ I ₄
3	I ₂ I ₃
4	I ₁ I ₂ I ₄
5	I ₁ I ₃

1 item set
I₁
T₋

Frequency
- 6
- 7

	I ₁	I ₂	I ₃	I ₄	
Y	I ₁	I ₂	I ₄		
S	I ₁	I ₃			
b		I ₂	I ₃		
T		I ₁	I ₃		
g		I ₁	I ₂	I ₃	
q		I ₁	I ₂	I ₃	

→

I ₁	- 6
I ₂	- 7
I ₃	- 6
I ₄	- 2
I ₅	- 2

Threshold = 2

Id	Name of item
1	I ₁ I ₂ I ₅
2	I ₂ I ₄
3	I ₂ I ₃
4	I ₁ I ₂ I ₄
5	I ₁ I ₃ ✓
6	I ₂ I ₃
7	I ₁ I ₃ ✓
8	I ₁ I ₂ I ₃ I ₅ ✓
9	I ₁ I ₂ I ₃ ✓

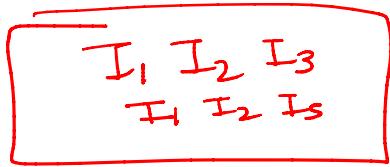
→

2 Item
freq set

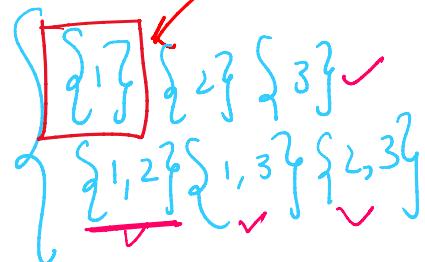
- ✓ I₁ I₂ - 4
- ✓ I₁ I₃ - 4
- ✓ I₁ I₄ - 1
- ✓ I₁ I₅ - 2
- ✓ I₂ I₃ - 4
- ✓ I₂ I₄ - 2
- ✓ I₂ I₅ - 2
- I₃ I₄ - 6
- I₃ I₅ - 1
- I₄ I₅ - 0

3 Freq ✓

I ₁ I ₂ I ₃	2
I ₁ I ₂ I ₄	1
I ₁ I ₂ I ₅	2
I ₁ I ₃ I ₄	0
I ₁ I ₃ I ₅	1
I ₁ I ₄ I ₅	0
I ₂ I ₃ I ₄	0
I ₂ I ₃ I ₅	1
I ₂ I ₄ I ₅	0
I ₃ I ₄ I ₅	0



Freq Item 3 $\rightarrow \{1, 2, 3\} \quad \{1, 2, 5\}$



✓ Support $\rightarrow 2/9$
 ✓ confidence $\rightarrow 50\%$

$$\text{Rule 1} \rightarrow \{1\} \rightarrow \{2, 3\}$$

Confidence = $\frac{\text{Support}(1, 2, 3)}{\text{Support}(1)} = \frac{2/9}{6/9} = \underline{\underline{33\%}}$

1 ↗ 2, 3 ~~33% ↗ 50%~~

Rule 2 $\{2\} \rightarrow \{1, 3\}$

$$\text{coh} = \frac{\text{Support}\{1, 2, 3\}}{\text{Support}(2)} = \frac{2/9}{7/9} = \underline{\underline{28\%}}$$

~~28% ↗ 50%~~

Rule 3 $\rightarrow \{3\} \rightarrow \{1, 2\} \rightarrow = 33.3\%$

Rule 4 \rightarrow

$$\underline{\{1, 2\}} \rightarrow \{3\} = \frac{\text{support}(1, 2, 3)}{\text{support}(1, 2)} = \frac{2/9}{4/9} = 50\%$$

50% $\geq 50\%$

Valid Rule

$$\{1, 2\} \rightarrow \{3\}$$

* PCA
* Demo Apriori
* LDA

TG