**Q3** **(a)**

Clan.

C1 :- buys - Computer = "Yes"

C2 :- buys - Computer = "no"

_Data to be classified:-_

$X = $ (age $<=30$, Income = medium, Student = yes, Credit - rating = Fair)

$P(C_1) = P($buys - comp = "Yes"$) = \frac{9}{14}$

$= 0.643$

$P($buys - comp = "no"$) = \frac{5}{14}$

$= 0.357$

— Compute $P(X|C_i)$ for each class.

$P($age = "$<=30$" | buys - comp = "Yes"$) = \frac{2}{9}$

$= 0.22$

$P($age = "$<=30$" | buy - comp = "no"$) = \frac{3}{5} = 0.6$

$P($income = "medium" | buys - comp = "Yes"$) = \frac{4}{9} = 0.444$

$P($incom = "medium" | buys - comp = "no"$) = \frac{2}{5} = 0.4$

$P($Student = "Yes" | buys - comp = "Yes"$) = \frac{6}{9} = 0.667$

$P($Student = "Yes" | buys - comp = "no"$) = \frac{1}{5} = 0.2$

$P($credit - rating = "fair" | buys - comp = "Yes"$) = \frac{6}{9} = 0.667$

$P($credit - rating = "fair" | buys - comp = "no"$) = \frac{2}{5} = 0.4$

$X = (age <= 30, income = medium, Student = yes,$
$credit - rating = fair).$

$P(X|c_i) : P(x|buy\_comp = 'Yes') =$
$$= 0.222 \times 0.444 \times 0.667$$
$$\times 0.667$$
$$= 0.044$$

$P(x| buys\_comp = 'no') = 0.6 \times 0.4 \times 0.2 \times 0.4 = 0.019$

$P(x|c_i) \ast P(c_i) : P(x| buys\_comp = 'Yes') \ast P(buys\_comp = 'Yes')$

$$= 0.028$$

$P(x|buys\_comp = 'no') \ast P(buys\_comp = 'no') = 0.007$

→ Therefore, $X$ belongs to class ("buys-comp = Yes")

[b] Limitation of Naive Bayes

— Naives Bayesian predication requires each requires each conditional prob. be non-zero.

— otherwise the predicted prob. will be zero.

$$P(x|c_i) = \prod_{k-1}^{n} P(x_k|c_i)$$

— eg, Suppose a dataset with 1000 tuples, income = low(0) income = medium (990), & income = high(10).

— Use laplacian Correction

○ · Add 1 to each Case.

Prob (income = low) = 1/1003
Prob (income = med) = 991/1003
Prob (incu = high) = 11/1003

— The "Correction" prob estimate are close to thier "UnCorrected" Counterparts.

## Q2 b

— K-means, each Cluster is represented by the Centre of the Cluste.

— Given k, the k-means algorithm is implemented in 4 8kps.

1. Partition object into k non-empty Subsets.

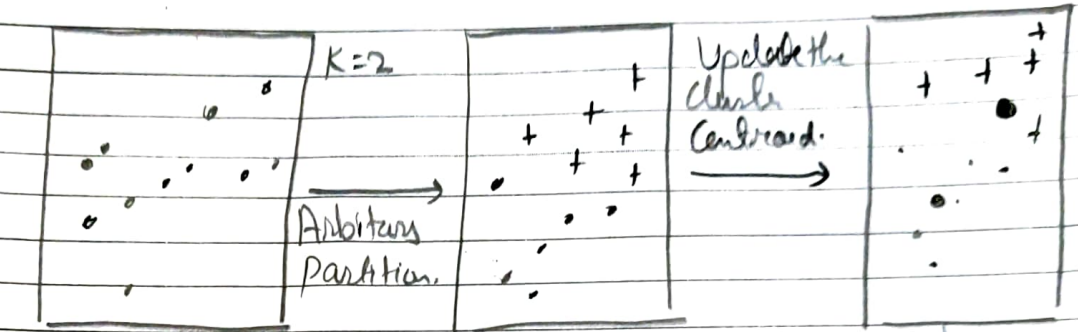2. Compute Seed points as the Centroids of the Cluster of the Current partitioning

3. Assign each object to the Cluster with the nearest Seed point.

4. Go back to Step 2, Stop when the assignment does not change.
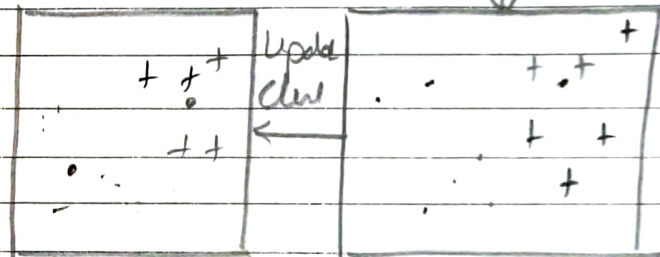
→ Example.



The initial data set.

— Partition object into k non empty ∨ Subset.

↑ Loop if needed.

Reassign obj.

— Limitation.

- K-mean Clustering Algorithm has limitation.

1. It requires to specify the no. of cluster (k) in advance.

2. It can't handle noisy data od outliea.

3. It is not Suitable to identify cluster with non-convex shapes.

**[Q2]** **[a]** Decision Tree Algorithm. [DTA].

— DTA belongs to the family of Supervised learning algorithm.

— The decision Criteria are different for classification & regression tree.

— Decision Tree use multiple algorithm to decide to split a node into 2 or more Sub-node.

— The decision Tree split selects the split all nodes available variable & then select the split which result in most homogenous Sub-node.

— The algorithm Selects is also based on the type of target Variable

ID3 , C4.5, CART ( Classification & Regression Tree) CHAID, MARS.