

# Data Science Life Cycle – Tools and Technologies

**M.S.Mohan Sivam, B.E.,M.Tech.,M.B.A.,**

Technical Lead – Bigdata Infrastructure

Tata Consultancy Services – Chennai

# Agenda

- What is Data Science?
- Bigdata
- Data Science Life Cycle
- Tools and Technologies
- Top Algorithms in Data Science Use cases.
- Skillset required for Data scientist.

# What is Data Science?

“Data science is an inter-disciplinary field that uses scientific methods, processes, algorithms and systems to extract knowledge and insights from many structural and unstructured data.”

- Wikipedia

# What is BigData?



Put simply, big data is larger, more complex data sets, especially from new data sources.



These data sets are so voluminous that traditional data processing software just can't manage them.



But these massive volumes of data can be used to address business problems you wouldn't have been able to tackle before.



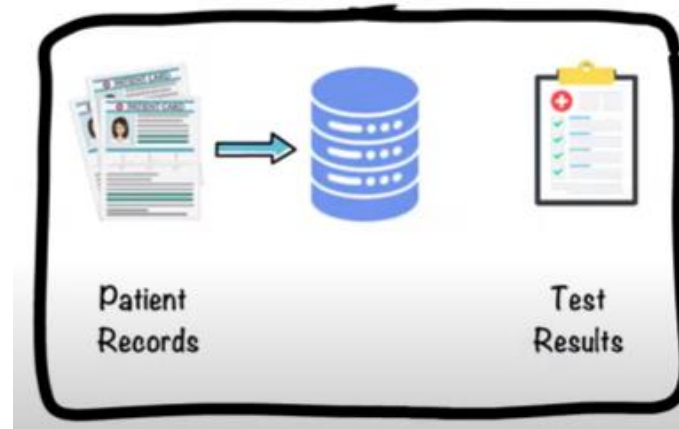
90% of the Data is Generated in last 2 year.

# BigData (Eg – Health care)

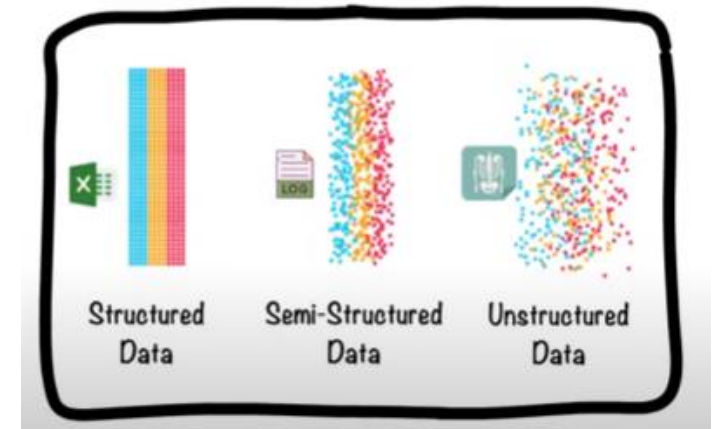
Volume



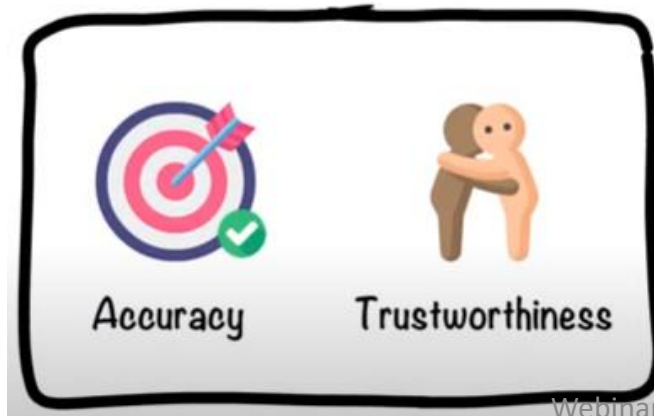
Velocity



Variety



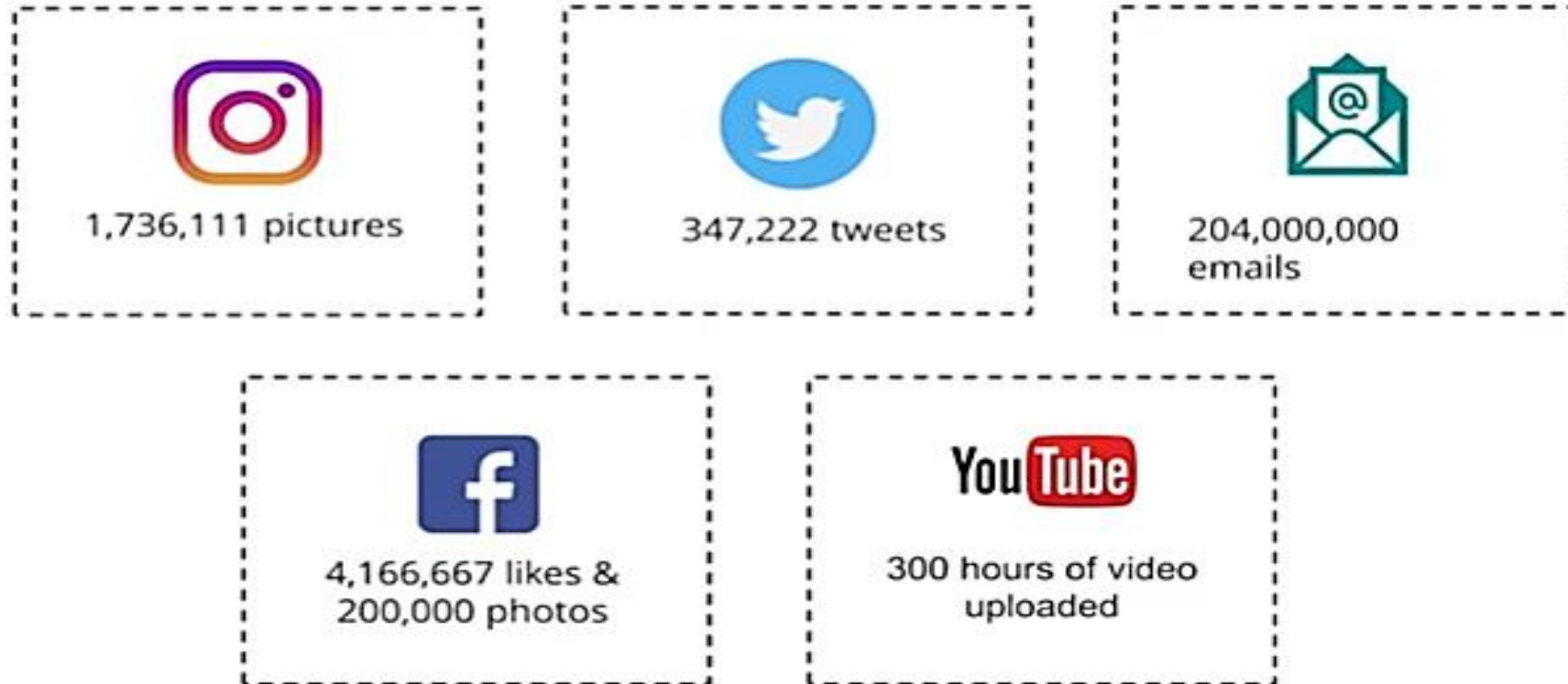
Veracity



Value

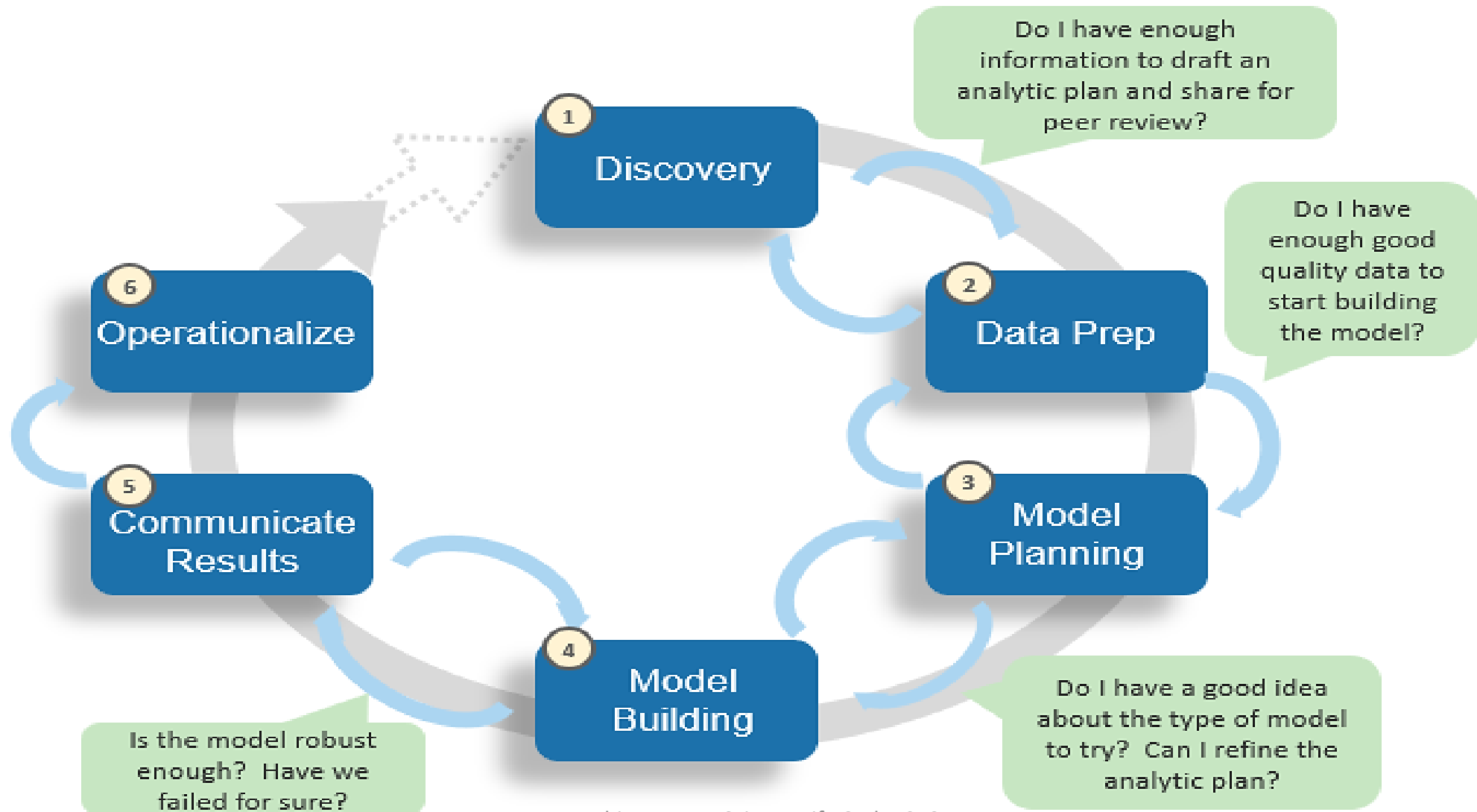


# Every minute of the Day!



**But How do we store and process this BigData!?**

# Data Science Life Cycle



# Data Science Tools





# Contd..



HADOOP



NOSQL DATA  
BASES



ETL TOOLS



DEVOPS



MACHINE  
LEARNING



ARTIFICIAL  
INTELLIGENCE



DATA  
VISUALIZATION



ETC.,

# Top Algorithms in Data Science Use cases.



Decision Tree



Random  
Forest



Association  
Rule Mining



Linear  
Regression



K-Means  
Clustering

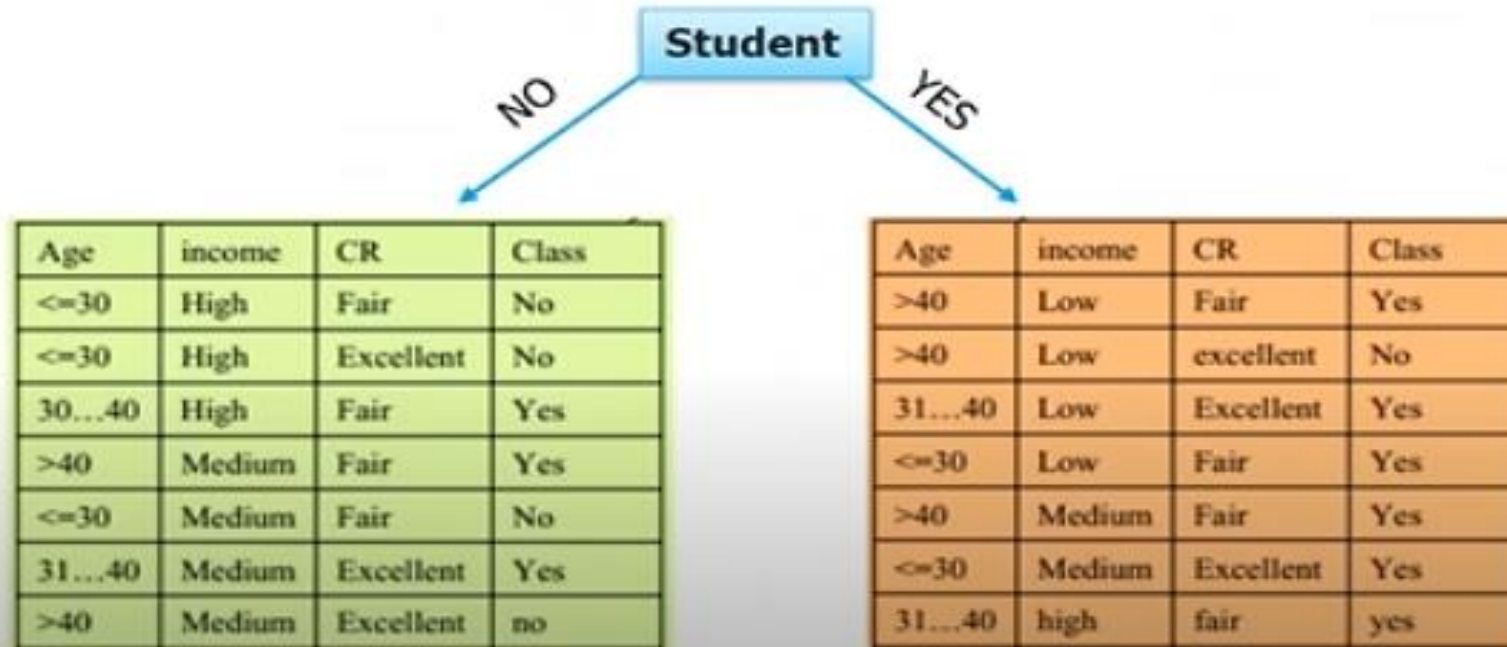
# Decision Tree

rec	Age	Income	Student	Credit_rating	Buys_computer
r1	<=30	High	No	Fair	No
r2	<=30	High	No	Excellent	No
r3	31...40	High	No	Fair	Yes
r4	>40	Medium	No	Fair	Yes
r5	>40	Low	Yes	Fair	Yes
r6	>40	Low	Yes	Excellent	No
r7	31...40	Low	Yes	Excellent	Yes
r8	<=30	Medium	No	Fair	No
r9	<=30	Low	Yes	Fair	Yes
r10	>40	Medium	Yes	Fair	Yes
r11	<=30	Medium	Yes	Excellent	Yes
r12	31...40	Medium	No	Excellent	Yes
r13	31...40	High	Yes	Fair	Yes
r14	>40	Medium	No	Excellent	No

# Contd..

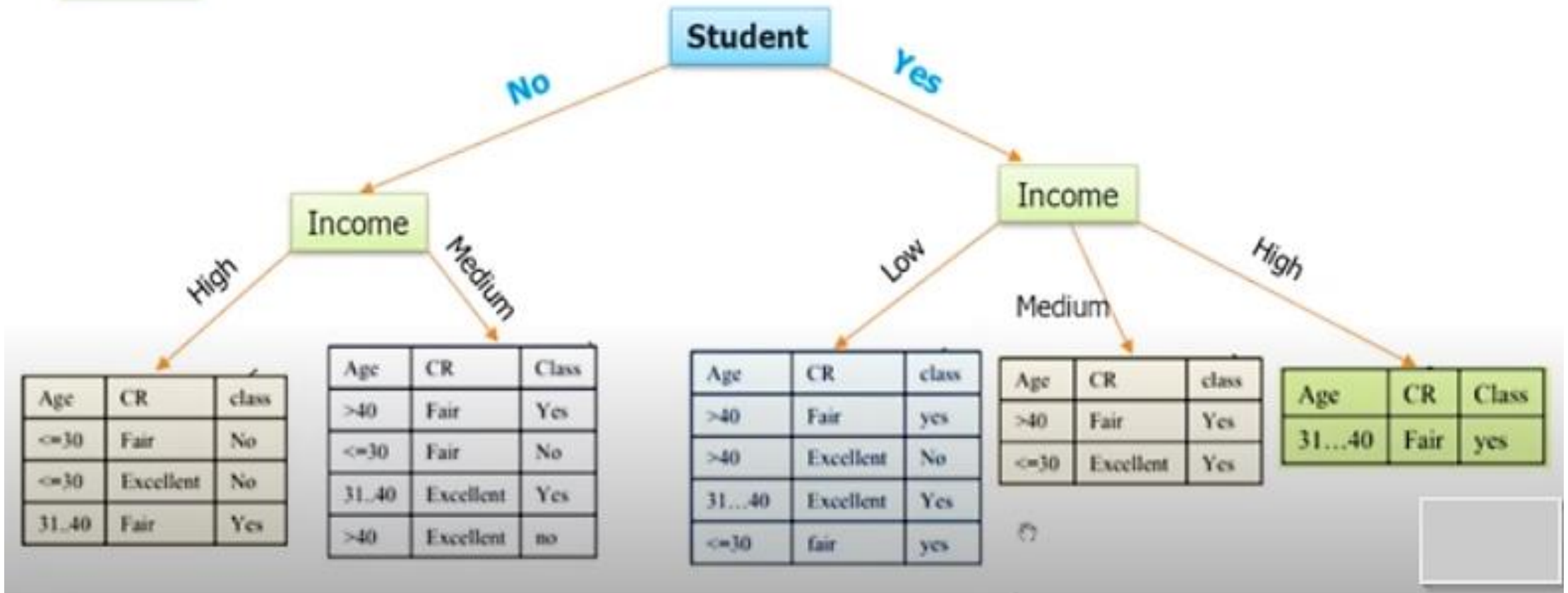
- Step 1

Step-1



# Contd..

Step-2



# Contd..

## Classification rules :

- 1.  $\text{student}(\text{no}) \wedge \text{income}(\text{high}) \wedge \text{age}(\leq 30) \Rightarrow \text{buys\_computer}(\text{no})$
- 2.  $\text{student}(\text{no}) \wedge \text{income}(\text{high}) \wedge \text{age}(31 \dots 40) \Rightarrow \text{buys\_computer}(\text{yes})$
- 3.  $\text{student}(\text{no}) \wedge \text{income}(\text{medium}) \wedge \text{CR}(\text{fair}) \wedge \text{age}(> 40) \Rightarrow \text{buys\_computer}(\text{yes})$
- 4.  $\text{student}(\text{no}) \wedge \text{income}(\text{medium}) \wedge \text{CR}(\text{fair}) \wedge \text{age}(\leq 30) \Rightarrow \text{buys\_computer}(\text{no})$
- 5.  $\text{student}(\text{no}) \wedge \text{income}(\text{medium}) \wedge \text{CR}(\text{excellent}) \wedge \text{age}(> 40) \Rightarrow \text{buys\_computer}(\text{no})$
- 6.  $\text{student}(\text{no}) \wedge \text{income}(\text{medium}) \wedge \text{CR}(\text{excellent}) \wedge \text{age}(31 \dots 40) \Rightarrow \text{buys\_computer}(\text{yes})$
- 7.  $\text{student}(\text{yes}) \wedge \text{income}(\text{low}) \wedge \text{CR}(\text{fair}) \Rightarrow \text{buys\_computer}(\text{yes})$
- 8.  $\text{student}(\text{yes}) \wedge \text{income}(\text{low}) \wedge \text{CR}(\text{excellent}) \wedge \text{age}(31 \dots 40) \Rightarrow \text{buys\_computer}(\text{yes})$
- 9.  $\text{student}(\text{yes}) \wedge \text{income}(\text{low}) \wedge \text{CR}(\text{excellent}) \wedge \text{age}(> 40) \Rightarrow \text{buys\_computer}(\text{no})$
- 10.  $\text{student}(\text{yes}) \wedge \text{income}(\text{medium}) \Rightarrow \text{buys\_computer}(\text{yes})$
- 11.  $\text{student}(\text{yes}) \wedge \text{income}(\text{high}) \Rightarrow \text{buys\_computer}(\text{yes})$

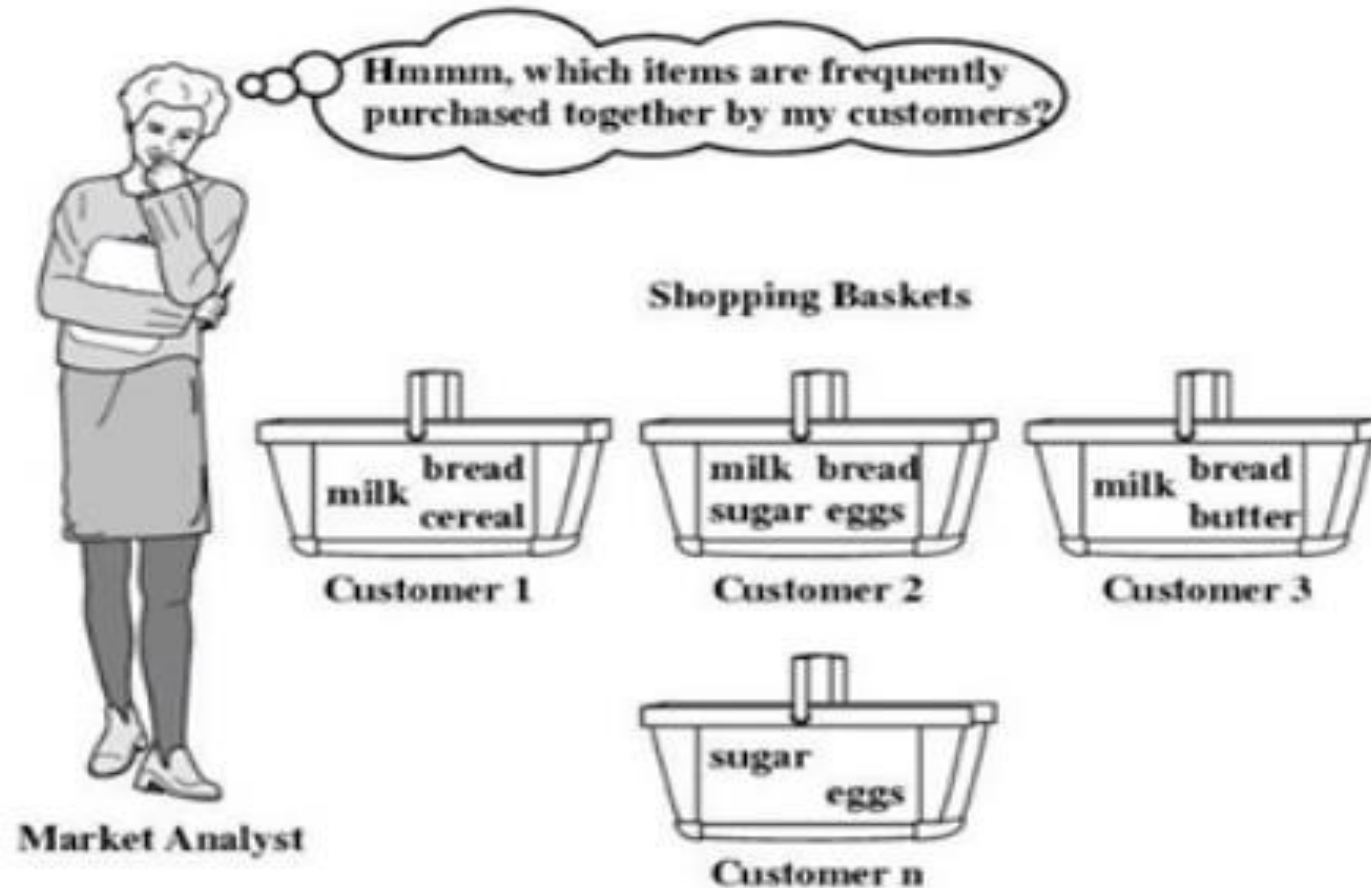


# Random Forest



- Suppose you are indecisive about watching a movie, “Kaatru Veliyidai”
- You can do one of the following
  1. Either you ask your best friend, whether you will like the movie or not.
  2. Or you can ask your Group of Friends.

# Association Rule Mining





# Contd..

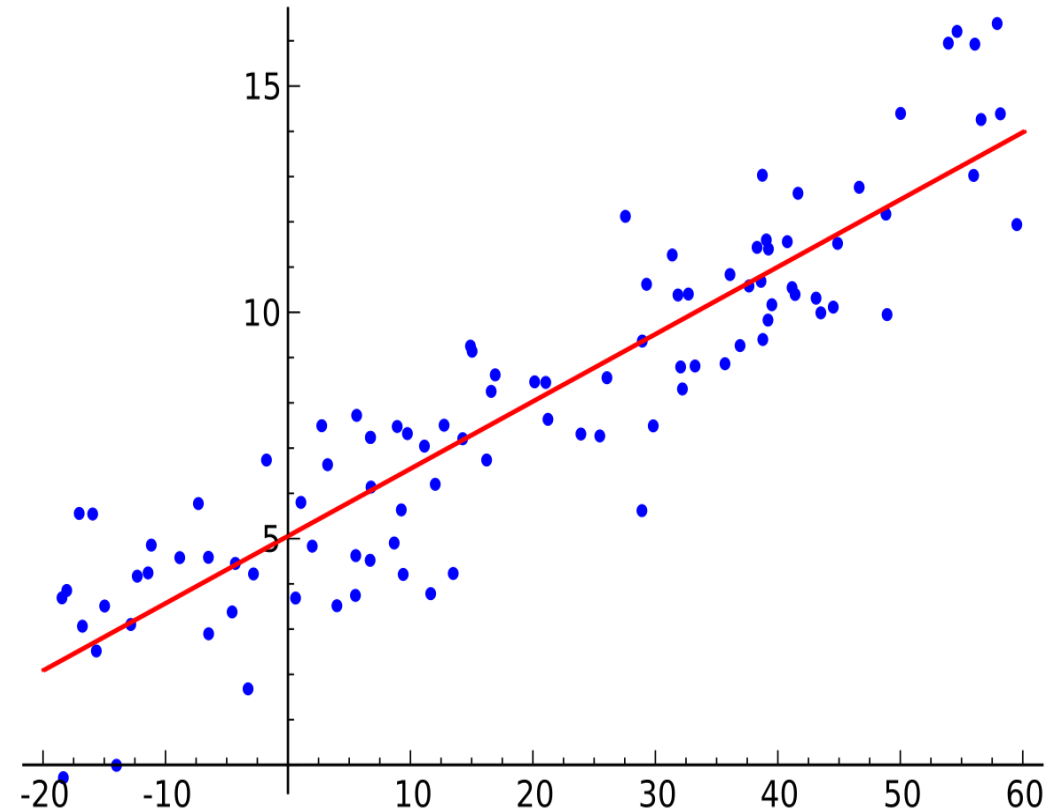
- **Association Rule Mining** is a Popular and well Researched method for discovering interested relations between variables in large Data.
- The Rule found in the sales data of a Super Market would indicate that **if a Customer buys Onions and Tomato together**, he or she is likely to also buy a Chicken.

$$\text{Support}(A \rightarrow B) = \frac{\text{Number of Records containing both A and B}}{\text{Amount of Records}}$$

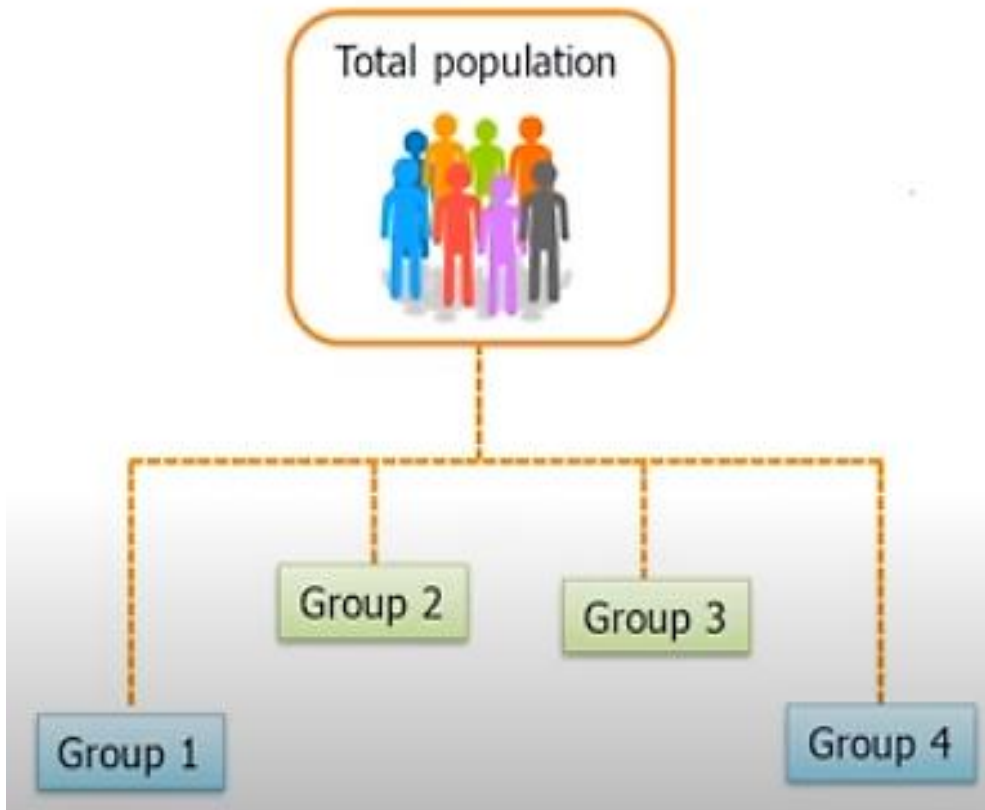
$$\text{Confidence}(A \rightarrow B) = \frac{\text{Number of Records containing both A and B}}{\text{Number of Records containing A}}$$

# Regression Analysis – Linear Regression

- Regression Analysis helps to understand how value of **Dependent variable** Changes when any one of **Independent variable** Changes, while other dependent variables are kept fixed.
- Linear Regression is the most Popular Algorithm used for **Prediction** and **Forecasting**



# K-Means Clustering



- Kmeans algorithm is an iterative algorithm that tries to partition the dataset into Kpre-defined distinct non-overlapping subgroups (clusters) where each data point belongs to only one group. It tries to make the inter-cluster data points as similar as possible while also keeping the clusters as different (far) as possible.
- The Objects Group 1 should be as similar as possible.
- But there should be much difference between objects in different Groups.
- The attributes of the Objects are allowed to determine which objects should be grouped together.

[https://github.com/msmohansivam/SKCT Webinar Hadoop](https://github.com/msmohansivam/SKCT_Webinar_Hadoop)

# Thank You

<https://in.linkedin.com/in/m-s-mohan-sivam-281305a7>

+91 9790748175