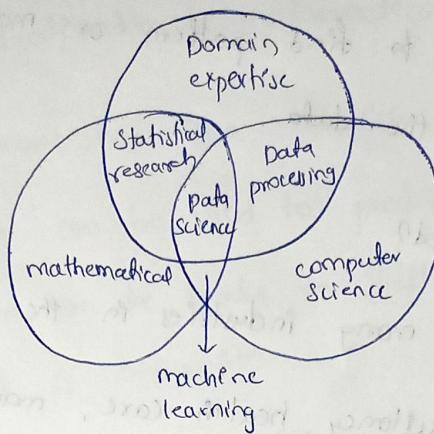


# What is Data Science?

①



- Data science combines maths & statistics, specialized programming, advanced analytics, AI & ML with specific subject matter expertise to uncover actionable insights hidden in an organization's data. These insights can be used to guide decision making & planning.
- Data science is a combination of multiple disciplines that uses statistics, data analysis & machine learning to analyze data to extract knowledge & insights from it.
- Data science is about data gathering, data analysis & decision making.
- Data science is about finding patterns in data through analysis & make future predictions by using DS companies are able to make.

1) Best decision

2) Predictive analysis

3) Pattern discovery to find pattern or may be hidden information in the data.

Where DS needed?

→ DS is used in many industries in the world today.

→ Eg:- Banking, consultancy, health care, manufacturing etc.

→ Data is a collection of raw facts.

→ Raw facts:- raw facts may be text / images.

→ Information:- which has a meaning for the data or

→ Information:- which has a meaning for the data or which collectively carries logical data. Eg:- Hindisong.

→ Data:- which does not have meaning.

Eg:- song, text.

Data Science Television:

The Quant shop:-

→ The DS television show we live in a data driven world, where massive computer analysing the trace of all human activity in our quest to predict the future.

- Data Science is a rapidly emerging discipline at the intersection of statistics, machine learning, data visualization & mathematical modeling.
- The Quantshop is a television show about data & how data can be used to predict the future. We show how data can be made to talk & explore the limits of what is knowable.
- This is a program for anyone who has ever made a bet or thought about tomorrow for anyone interested in how the mathematical or computational models that rule the world are conceived built & tested and how well they perform.

### Structure

- Each of the eight programs in all practical reason are built around a particular real world prediction challenge we watch as a student team comes to grips with the problem learning along with them as they build a forecasting model.
- They make their predictions & watch along with them to see if they are right or wrong. This first season features episodes on:

## 1) Finding Miss Universe:-

(4)

- We can computational models predict who will win a beauty contest? Is beauty just subjective or can algorithm tell who is the fairest one of all?

## 2) Modeling the movies:

- Which movie will gross the most on Christmas day? which actors deserve can data reveal the next Hollywood star?

## 3) Winning the body pool:

- How accurately can we predict Juniors weight before they are born? How can data clarify environmental risks to developing pregnancies?

## 4) The art of the Action:-

- How many millions will a particular J.W Turner painting fetch at auction? Can computers have an artistic sense of what's worth buying?

## 5) White Christmases:

- What places will wake up to a snowy Xmas this year? And how can you tell one month in advance?

(5)

### 6) Predicting the playoffs:-

- Which football team will the Super bowl? Can Google's page rank algorithm pick the winner on the field as accurately on the web?

### 7) The Ghoul pool:-

- Death comes to all men but when can we apply actuarial models to celebrate a dude who shall live & who shall die?

### 8) Playing the market:-

- Speculators get rich when right and tomorrow's price & poor when wrong how to equant predict?

Computer Science, Data Science, Real Science.

### Data & Method Centrism:-

- Scientist are data driven (Analysis of data).
- Computer scientists are algorithm driven.
- Real scientists spend anomaly (very large in size, quality or extent) amounts of efforts collecting data to answer their questions of interest.

## Concern about result:-

- Real scientists care about answers they analyse data to discover something about how the world works.
- Good science care about whether the result make sense because they care about what the answer means.
- Computer scientists worry about producing fake-looking numbers.

## Robustness:-

- Robustness is nothing but quality or condition of being strong.
- Real scientists are comfortable with the idea that data has errors.
- Computer scientists are not comfortable.
- Scientists think a lot of possible sources of bias or error in their data.
- Good programmers use strong datatype and parsing methodologies to guard against formatting errors.

## Perception:-

- Nothing is completely true or false in science while everything is either true/false in computer science.

## Asking interesting questions from Data :-

(7)

- What things might you be able to learn from a given data set.
- What do you or your people really want to know about the world.
- What will it mean to you once you find out?

## Google N-gram:-

- Google uses data to improve search result and provide fresh access to out of print books.
- N-gram is a set of sequential of N-words ~~test~~
- The google N-gram displays user to selected word or passes (N-gram) in a graph that shows how those passes have occurred in a corpus.
- Google N-gram corpora are made up of scanned books available in google books.

## Properties of Data:-

- One purpose of DS is to structure data, may be it is interpretable and easy to work with.

→ Data can be characterized into 3 groups.

(8)

- 1) Structured & unstructured
- 2) Quantitative & categorical
- 3) Big Data & little data.

### 1). Structured & Unstructured data.

#### Unstructured

→ Unstructured data is unorganized but we must organise the data for analysis process.

→ e.g., collection of tweets from Twitter.

→ Our first step is generally to build a matrix structure

→ A bag of words model

will construct a matrix with a row of each tweet and a column for each frequently used vocabulary words.



→ The no. of times tweet ( $i$ ) word ( $j$ )

#### Structured

→ Structured data is organized and easier to work with the data.

→ Data is often represented by a matrix; where the rows of a matrix represent distinct items or records and the columns represent distinct properties of items or records.

□	□	□	□
□	□	□	□
□	□	□	□

## 2) Quantitative & categorical data:-

(a)

### Quantitative

- Quantitative Data consists of numerical data (or) values like height and weight.
- The data can be incorporated directly into algebraic formulae and mathematical models (or) displayed in conventional graphs and charts.
- Quantitative data can be analysed using various statistical techniques such as mean, median, mode, variance & standard deviation etc.

### Categorical

- Categorical data consists of cables describing the properties of the objects under investigation like gender, hair, colour and occupation.
- On other hand the categorical data can represent with a specific data or ranking.
- The difference between them may not consistent or measurable.

### 3) Big data & Small data:-

Big data	Small data.
<p>→ Big data refers to extremely large and complex data sets that are beyond the capacity of the traditional data processing and analysis tools, to handle effectively.</p> <p>→ Big data can exhibit inconsistencies in data quality format or structure. This variability can make it challenging to work with and analyze the data effectively.</p>	<p>→ Little data involves relatively small amounts of information typically ranging from a few kilobytes to giga bytes. It can be stored and processed on a single computer or small server.</p> <p>→ Analysis of little data can be performed using common data analysis and visualization techniques.</p>

### Developing scoring system:-

→ Scores are functions that map the features of each entity to a numerical values of merit the basic approaches for building effective scoring system and evaluating them.

## Gold standards & proxies:

(11)

- In data science, a gold standard is a set of labels or answers that we trust to be correct.
- In the original formulation of BMI, the gold standard was, the body that percentages carefully measured on a small number of subjects of course, such measurements are subject to some error, but by defining these values to be the gold standard for fitness.
- The presence of a gold standard provides a rigorous or best way to develop a good scoring system.
- We can use curve fitting technique like linear regression to weight the input features so as to best approximate the "right answers" on the gold standard instances.
- But it can be hard to find real gold standards proxies are easier to find data that should correlate well with the desired but unobtained ground truth.

- BMI was designed to be a proxy for body fat percentage. It is easily computer from height & weight and does a pretty good job correlating with body fat.
- Proxies are particularly good when evaluating scoring / ranking systems.

### Score Vs Ranking:-

- Rankings are permutative ordering n entities by me it, generally constructed by sorting the output for some scoring systems. popular examples of ranking, Rating system includes:

#### (1) Football / Basketball top twenty:-

- Press agencies, generally rank the top college sports teams by aggregating the votes of coaches or sports writers.
- Typically each voter provides their own personal ranking of the top 20 teams, and each team gets awarded more points the higher they appear on the voters lists.

→ Summing up the points from each voter gives a total score for each team and sorting these scores defines the ranking. (13)

## (2) University Academic Rankings:

- The magazine called US News and World Report publishes annual rankings of the top American colleges and universities.
- Their methodology is proprietary and changes each year, presumably to motivate people to buy the news ranking.
- But it is generally a score produced from statistics like faculty/student-ratio, acceptance ratio, the standardized test score of its students and applications.
- Since scores and rankings are duals of each other, which produces or provides or more meaningful?
- Representation of the data? As in any comparison the best answer is that it depends on factors like:-
  - i) will the numbers be presented in isolation?
  - ii) what is underlying distribution of score?
  - iii) Do you care about the extremes or middle?

## Recognizing Good Scoring Functions:-

→ Good scoring functions are good because they  
are easily interpretable and generally believable.  
Here we review the properties of statistics which  
point in these direction.

### (1) Easily Computable:-

→ Good statistics can be easily described and  
presented. BMI is a excellent example, it contains  
only two parameters and is evaluated using  
only simple algebra.

Some of the other examples are as follows:-

- 1) easily understandable.
- 2) Monotonic Interpretation of variables.
- 3) Produces generally satisfying results on outliers.
- 4) Uses systematically normalized variables.

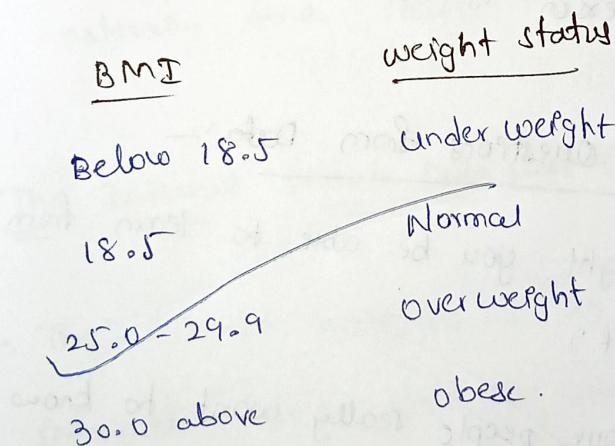
## BMI [Body Mass Index]:-

(15)

→ BMI is a measure used to access an individuals body weight in relation to their height. It is calculated by dividing a persons weight in kilograms by the square of their height in meters.

→ Formula to calculate BMI:- 
$$\frac{\text{Weight (in kg)}}{(\text{height})^2 \text{ (in meters)}}$$

→ BMI of an individual is calculated by the use of a mathematical formula it can also be estimated using weight in pounds to estimate BMI.



- BMI is commonly used to categorize individuals into different weight classes such as underweight and obese.
- BMI is used to provide a good measure of obesity. But BMI fails to provide actual information on body composition like amount of muscle bone fat etc.
- As single measure BMI is clearly not a perfect measure of health. But it's still a useful starting point for important condition that becomes more likely when a person is obese.

### Asking Interesting Questions from Data! -

- 1) What things might you be able to learn from a given data set?
- 2) What do you/your people really want to know about the world?
- 3) What will it mean to you once you find out?

## The baseball encyclopedias:-

(17)

- A large collection of information about one or many subjects often arranged alphabetically in articles in a book or set of books.
- Baseball has long had an outsize importance in the world of data science.
- What makes baseball important to data science is its extensive statistical record of play, dating back for well over a hundred years.
- Baseball is a sport of discrete events; pitchers throw balls & batters try to hit them, that naturally lends itself to informative statistics.

## The Internet movie Data base (ImDb);-

- The Internet movie db is an online db containing information and statistics about movies, TV shows and video games as well as actors, directors & other film industry professionals.

- This information includes list of cast to crew members, movie release date & box office information.
- IMDb was first published in 1990 by Weedham, a computer programmer as a group of scripts which allowed users to search a list of film credits.

### Advanced Ranking Techniques:-

- In the absence of any gold standard, these methods produce statistics which ~~are~~ are often revealing and informative.
- When several powerful techniques have been developed to compute ranking from specific type of inputs.

#### ① Elo ranking:-

- Rankings are often formed by analyzing sequence of binary comparison which arise naturally in competitive between entities.

### \* Sports Contest Results:-

- Typical sorting events
- Assume that in a football (or) chess match, pit teams 'A' & 'B' against each other.
- Only one team will win.
- Thus each match is essentially a binary comparison of merit.

### \* Implicit Comparison:-

- Suppose a student has been accepted by both universities 'A' and 'B' but opts for 'A'. This can be taken as an implicit vote. that 'A' is better than 'B'.
- The Elo system starts by rating all players equally and their incrementally adjusted to each player's score in response to the outcome of each match.

$$r'(A) = r(A) + k(s_A - e_A)$$

- where  $\rightarrow r'(A)$  &  $r(A)$  represent the previous and updated score for player 'A'.

- $k$  is a fixed parameter reflecting the maximum possible score adjustment in response to a single match.
- $s_A$  is the scoring results achieved by player A in the match under consideration. Typically  $s_A = 1$  if A won &  $s_A = -1$  if A lost.
- $\mu_A$  was the expected once competing against B if A has exactly the same skill level as B then presumably  $\mu_A = 0$ .
- ⇒ what is expected Match Score?

if  $P(A>B)$  estimates the probability that A beats B

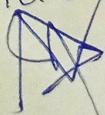
$$\mu_A = 1 \cdot P_{A>B} + (-1) (1 - P_{A>B})$$

→ If the ranking system is meaningful, this probability should be a function of the difference the score  $r(A) & r(B)$ .

The target function:-

→ we need a function  $f(x)$  that takes  $X'$  and  $Y'$

fields a probability



$$f(0) = 1/2$$

$$f(\infty) = 1$$

$$f(-\infty) = 0$$