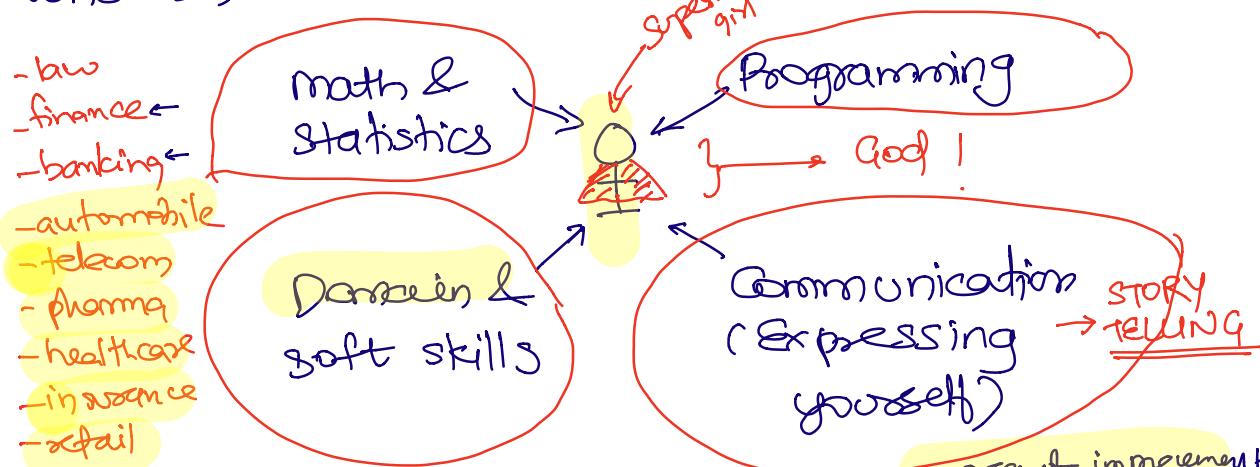


- Introduction to Data Science
- Introduction to Bigdata
- Setting up Python on your machine

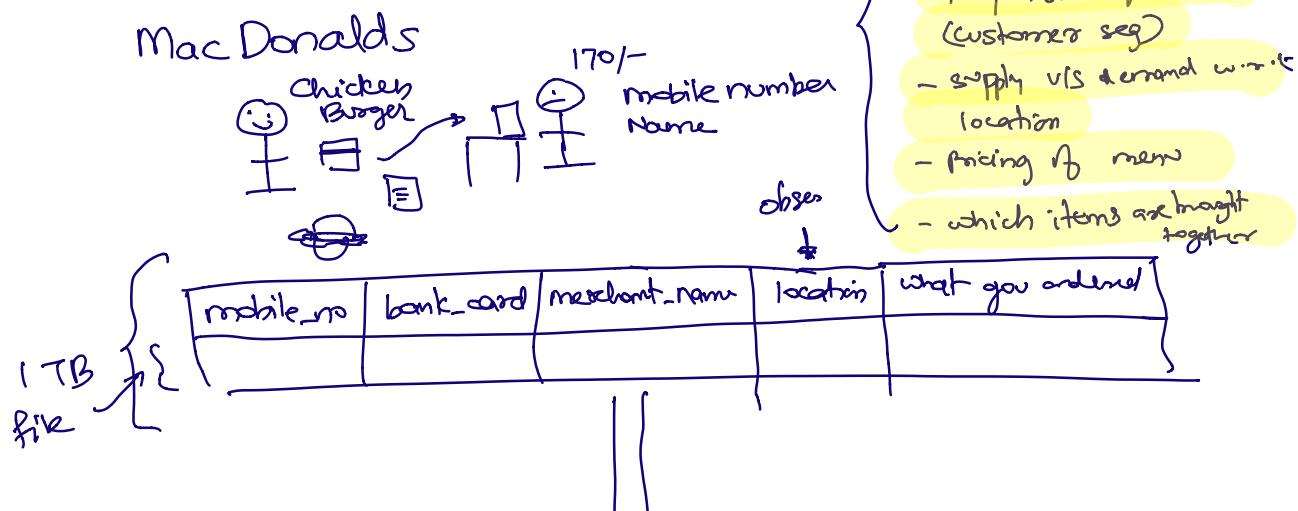
What is Data Science ?

Set of Practices that enable Data Scientist to get **INFORMED DECISIONS**

Who is a Data scientist ?



Case - study :



Insurance companies DS \Rightarrow statistical model for hospitals

A hospital chain XYZ group wants to understand in which area of Mumbai he shall open a new hospital?

- finding strategic locations in the given area to open hospitals
- finding strategic locations to install a mobile tower for better service.
- Gym fitness centers

DATA PRODUCT

→ Senator ✓/S
→ Senator ✓/S

mark zuckerberg (fB) } Data Privacy.
Sunder P. (Google)

Is my data exposed?
Are you keeping a saltine watch on me?
Are you able to download my files from device?
Are you able to read my email?
Are you spying me?
Are you capable to spy me?

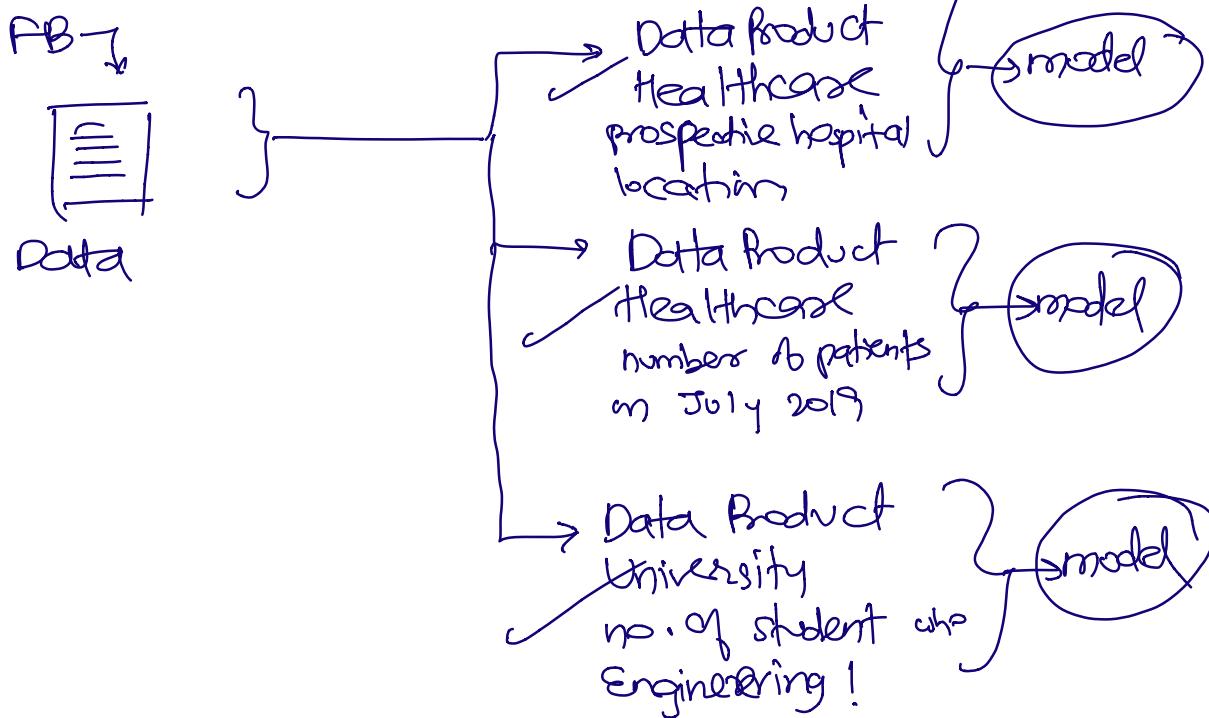
COMPROMISING USER'S
PRIVACY

FB → ?
Google → ? NO!

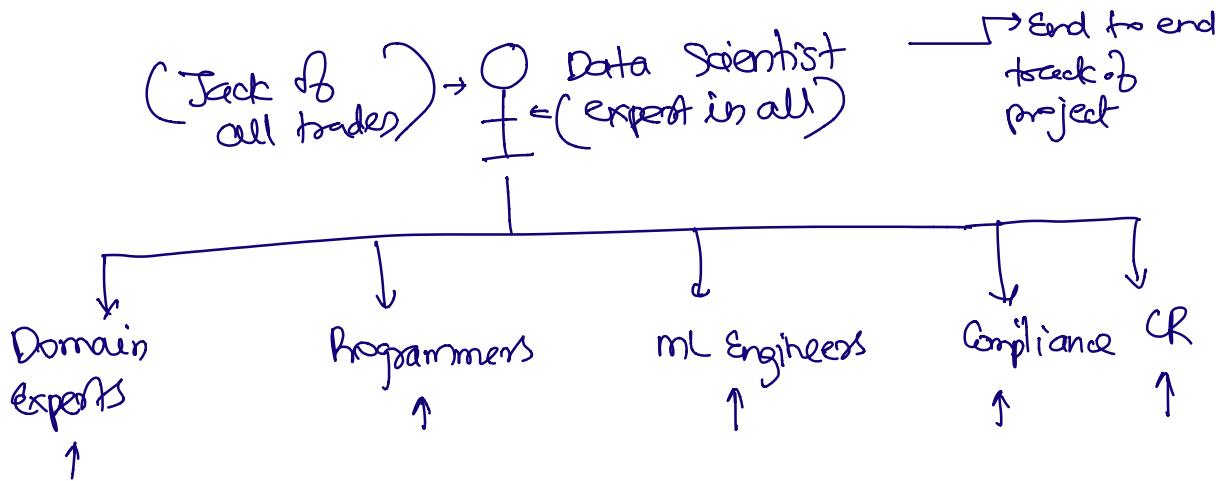
Google } never sell/share/use
 FB } your data directly } → they sell
STATISTICS

Age 18 - 35 } → Chicken Burgers } → Mumbai, Andheri

18 - 22 } → veg Burgers } → Mumbai, Dadar



file | equation | algo.



Set of Practices (Data Science) (Guidelines)

- ① Understand the use-case. (Are you the best fit for this use-case)
- ② Relate your use-case with domain for creative insights.
- ③ Gather relevant data / Data Product for the use-case.
 - Is data available with customer?
 - Is customer willing to give plain data?
 - Is the data-sharing satisfies all compliance standards? HIPAA | ISO 27001 | ISO 9001:2015

- ④ Apply Statistical analysis over the data.
- ⑤ Apply mL / DL / Dm algo. to create data product (model)
- ⑥ viz. the result.
- ⑦ Present your results / findings to stakeholders.

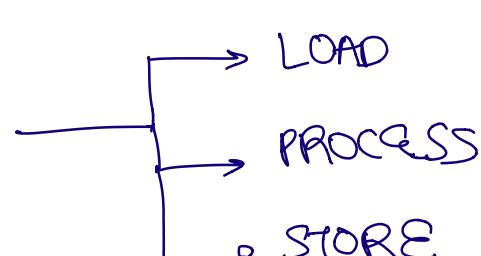
Bigdata

what is Bigdata ?

Bigdata is a TERM. Its not a tool, technique, technology -

Bigdata is all about the ability of your sw, framework, architecture (sw, hw) to HANDLE the data.

Handle refers to



? } if your app /
slow / arch. fails
you are facing
**BIGDATA
PROBLEM.**

BIGDATA IS A PROBLEM !

507 MB → 201201India.txt

Notepad → It tried its best but FAILED TO LOAD the data.

∴ for NOTEPAD, it's a BIGDATA

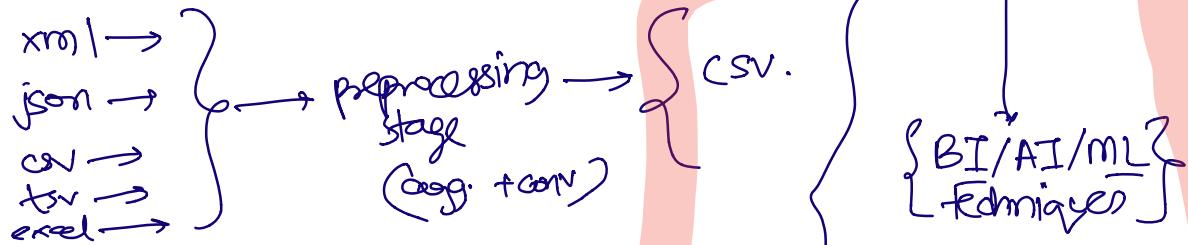
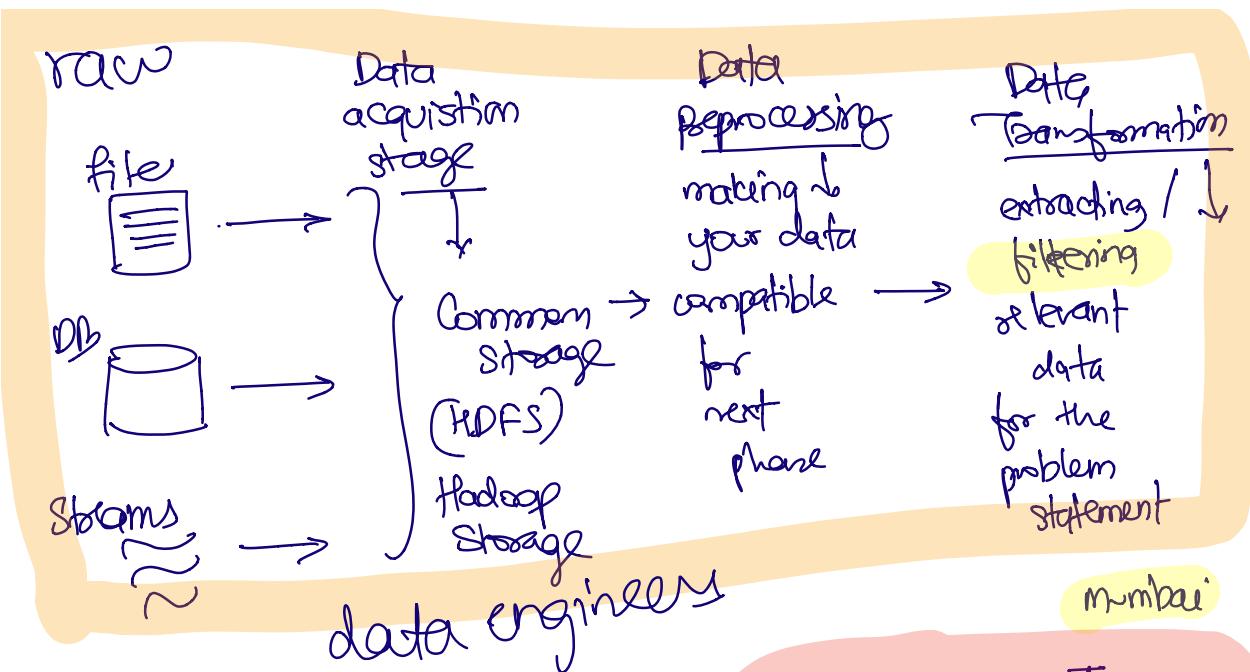
Wordpad → It tried, had some latency, but LOADED THE DATA

∴ for WORDPAD, it's a NORMAL DATA

Bigdata → Hadoop { Java based distributed storage & processing framework.

my capability to store and process the data is directly proportional to my underlying HW's capability.

- Apache Hadoop
 - Apache Spark
 - Apache Storm
- } Data processing
Data Transformation



DATA

LAKE

ELT

Data Engineers \Rightarrow data acq, data preprocessing, data conversion, data agg, data transformation
 \Rightarrow Hadoop, Spark, Storm, Kafka, Solr, Java, Python, Elasticsearch

Data Analyst \Rightarrow create statistical inference out of the data \rightarrow (Excel (Advanced), Tableau, SPSS, SAS, R, Python)

Business Analyst \Rightarrow create statistical inference out of the data w.r.t. domain.

ML Engineers \Rightarrow the one who make data product.
(Python \rightarrow numpy, sciPy, pandas, sklearn, tensorflow, keras, nttle, spaCy, beautifulsoup, openCV, etc)

Data Scientists \Rightarrow THE GOD (smiley face)
Master/Jack of all trades on all tools & technologies
with expertise in domain & knowledge & Good presentation & soft skills.
Handling the above team + interacting with customers.

① Python

- \rightarrow Basics of Python
- \rightarrow Numpy / SciPy
- \rightarrow Pandas
- \rightarrow matplotlib & seaborn (Viz)
- \rightarrow sklearn for ML

② Statistics

③ Exploratory Data Analysis (EDA)

④ Machine Learning (limited)

⑤ Bigdata using R, Spark (basics)

Which programming language I must learn
to become Analyst / DS ?

- ① Python 3
4
- ② R
- ③ Scala
- ④ Go lang

SQL

Excel techniques
(Business)

Jupyter → Python Hands-on exercises .

How to create a local lab environment
in your machine?

Anaconda Python