

Course 3:

Prepare Data for Exploration

That will cover

- How data is generated
- Different formats
- Analyze bias and credibility
- What is "Clean data"
- Data bases
- Extract own data - Spreadsheets + SQL
- Basics of data org.
- Protecting data

Week 1

~ Data is generated all time, every minute **millions text** and **hundreds of millions of emails** are sent. **Millions of online searches** are made and **videos are viewed**. Those numbers are only growing.

"Every piece of information is data"

All data is usually generated as a result of our activity in the world.

- Interviews, Observations, Forms, Questionnaires, Surveys, Cookies.

Factors to consider while collecting data:

- How the data will be collected

FP
SP
TP
P

- • Choose data sources

Sample,
Population

- • Decide what data to use

- • How much to collect

- Select the right Data type

Historical
Data

- Determine timeframe

• Buzz Words...

First Party Data:

Data collected by an individual or a group using their own resources.

- Preferred: As you collect only those which you know for exactly what.

Second Party Data:

Data collected by a group directly from its audience and then sold.

Third Party Data:

Data collected from outside sources who did not collect it directly.

- Might not be reliable.

Sample:

Sample is a part of population that is representative of the population.

Discrete data:

Can be counted and have limited number of values.

Continuous data:

Can be measured and can have almost any numeric value.

Internal Data:

Data that lives inside a company's own systems

External Data:

Lives outside of company's database

Nominal & Ordinal data.

Structured & Unstructured

— —
Mostly the data being generated is Unstructured, Audio, Video, Images, emails, social media text,

Data Model:

A Model that is used for organizing data elements and how they relate to one another.

Data elements = Attributes.

Structured

- Defined DTypes
- Mostly quantitative
- Easy - Organize
 - Search
 - Analyze
- Stored - RDBMS
 - Warehouse
- Contained in Rows & Cols

Unstructured

- Varied DTypes
- Mostly qualitative
- Difficult to search
- Provides more freedom for analysis
- Stored - Data Lakes
NoSQL
Warehouses
- Can't be put Row & Cols.

"Data models help us keep data consistent and they give us map of how data is organized"

What is Data Modelling?

It is the process of creating diagrams, that visually represent how data is organized and structured.

'Visual representations are called Data Models.'

'A blueprint of a house'

'Each level of data modeling has different level of detail'.

CONCEPTUAL

LOGICAL

PHYSICAL

① Conceptual:

Gives a high-level view of your data structure - How you want data to interact across an organization.

② Logical:

Focuses on the technical details of the model such as relationships, attributes & entities.

(3) Physical:

Show in actual that how the database is built.

Here you lay-out how each database will be put in place and DB, APPs, Features will interact in specific detail.

→ Data Modeling Techniques:

Two common:

ERD

(High level)

&

UML

(detailed)

* Data Type:

A specific kind of data attribute that tells what kind of value the data is.

Use - Record + Field ^{Row &} for Col.

W I D E
DATA

Data in which every data subject has a single row with multiple columns to hold the values of various attributes of the subject.

L

Data in which each row is one time point per subject, so each subject will have data in multiple rows.

O

D
A
T
A

N

(Only one column to add for new attribute)

G

~~WIDE~~

~~NAME~~ ~~NAME~~ ~~NAME~~

Arush 2011 120

Wide

Name	YR-2011	YR-2012	YR-2013
Arush	120	130	170
Baei	120	135	140
Heena	110	115	116

Long

Name	Years	Height
Arush	2011	120
Arush	2012	130
Arush	2013	170
Baei	2011	120
Baei	2012	135
Baei	2013	140
Heena	2011	110
Heena	2012	115
Heena	2013	116

Data Transformation:

A process of changing data's format, structure or values.

Doing this - you know before your analysis?

- Adding, copying, replicating
- Deleting fields or records
- Standardizing names of variables
- Renaming, Moving, combining
- Joining.

Why to transform the data?

Better
ORGANIZE

Better
COMPAIBILITY
(Multiple system
can use)

Data
MIGRATION
(one system to another)

DATA
MERGING

Better
ENHANCEMENT
(can be displayed with
more detailed fields)

Data
COMPARISON