# Who am I ?

# Beginner → Intermediate #(Python)

( started with Mathematics 2008.
↓
GATE/CAT
↓
Data Science )

~15 yrs

## AKASH PUSHKAR CHARAN
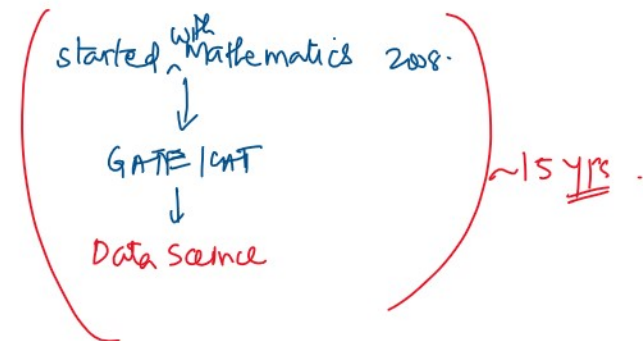### Quick witted | Tech-Savvy | Observer

- Machine Learning Enthusiast
- IIT Kanpur Alumnus
- Working as Lead Data Scientist with Accenture Strategy & Consulting

An instructor by passion and data scientist by profession who always tries to look at the problem in a different way. I started teaching when I was in first year of my graduation which amounts to **10+** years of experience.

*"It's not who I am underneath but what I do that defines me"*
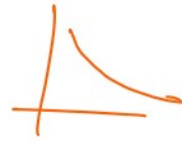
PROVIDE LINKEDIN PROFILE

https://www.linkedin.com/in/akash-pushkar-04642925/

# Setting the ground rules & expectations

## Ground rules :

① Sheer ==focus==, ==attention== & ==interaction==

Mobile phone → DND

doubts
↳ ask questions
— Two windows ⎡→ after break
⎣→ end of the session.⎦

② **Revision is the key to success** -
↳ ==Always revise all the notes (OneNote) & notebooks before== you attend the next session

③ 4 times a week.
In case, you happen to miss a class, ==pls cover the topics== using ·==LMS== .

④ Real life application ⇒ Industrial examples.
→ (learning/experience from my past projects)

⑤ **Everyone will code.**

⑥ Assessment questions ⎡→ HW
⎢→ QIC : Question in class
⎢→ Reading assignment ✓
⎣→ ==Case-studies==

⑦ Breaks # (10-15mins)

# Course Structure

**Pre-requisites:** Python Basics →
- Conditional statements
- Looping statements
- Functions / Lambda functions

**I** {

**Module #①** NumPy } For data manipulation activities ⇒ **EDA: Exploratory Data Analysis**

**Module #②** Pandas
↳ 2.0 has been released.

**Module #③** Data visualization using **Matplotlib** →
- MATLAB style
- Object Oriented style

# Storytelling with data
↳ seaborn
↳ Plotly / Altair ⇒ Enterprise visualization

**Module #④** Probability & statistics & a bit of maths : Foundation for ML / DS.
- Descriptive
- Inferential

**Module #⑤** Introduction to ML
- Supervised
- Unsupervised.

What is data manipulation?



**What is Data Manipulation?**                                    IntelliPaat

Data Manipulation is the process converting data into a format that is easy to process and is more organized

- Data extraction from multiple sources
- Manipulating data using Python
- Organized and readable information

Supply chain # Demand Forecasting

Sales data

| Date | Item | Quanity | Price | Sales | Site ID | Cust ID |
|------|------|---------|-------|-------|---------|---------|
| 13 Jun | SKU001 | 10 × | 100 | 1000 | | |
| 12 Jun | SKU002 | 5 × | 50 | 250 | | |
| 8 May | SKU003 | 2 × | 10 | 20 | | |

Item Master

| Item | Description | UOM | Item categ | Weigh | volume | Brand |
|------|-------------|-----|------------|-------|--------|-------|
| SKU001 | Biscuit | Ea | | | | |
| SKU002 | Shampo | Ea | | | | |
| SKU003 | Choco | Pack | | | | |

Site/Location Master

| Site ID | family name | City | State | Country | Lat | Long | Postal Code | Region |
|---------|-------------|------|-------|---------|-----|------|-------------|--------|
| 001 | | | | | | | | |
| 002 | | | | | | | | |
| 003 | | | | | | | | |
| 004 | | | | | | | | |

| Cust Id |
|---------|
| Customer Master |

Facts Table ⟶ Transaction Table (Sales data)

Dimension Table ⟶ Master Files (Item | site | customer)

# Why Data Manipulation?

There could be multiple problems with the data some of those are:

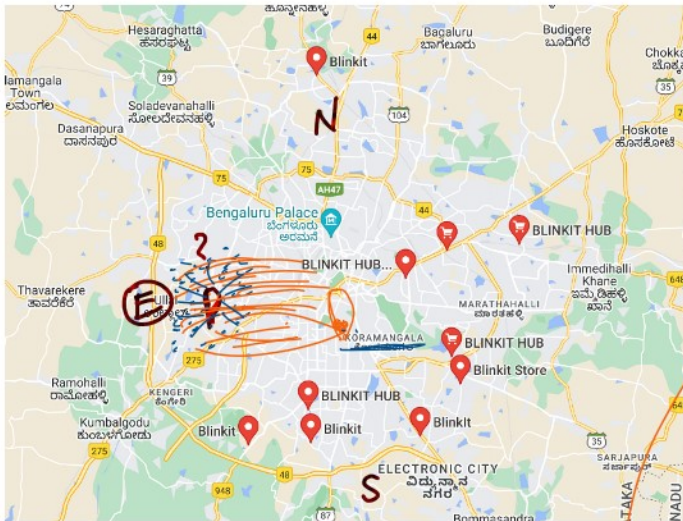| Missing Values | Incorrect Format | Different Units | Unnecessary Data |

Since the data is accumulated from multiple sources some values in a row might be missing. This could be due to multiple reasons such as equipment not being available, user not willing to share data etc.

---

## Missing Values Treatment

Location | Site Master

| Site ID | City | Address/State | Country | Region | Postal code | Lat. | Long. |
|---------|------|---------------|---------|--------|-------------|------|-------|
| LOC001 | BLR | KA | IN | South | 560xxx | — | — |
| LOC002 | " | " | | South | " | — | — |
| LOC003 | " | ✗ | | East 560055 | | ✗ | ✗ |
| | | | | | | | |
| LOC100 | " | " | | West | 560023 | | |

How will you treat missing Lat & Long

a) Drop the row(s)/column(s)
b) Impute with mean/median/mode.

---



KPI

Cost/order → [₹ 10 to ₹ 20]     100%

Finance/Account Segment

Site 003 → Site Address → Postal code → (Lat) & (Long)

Brick & Mortar

site #003 & Electricity + Lease (Rental)

Sales | Site ID

Site

Finance | Operational | Marketing | Sales

Input → Black Box → output

**Data Science is an enabler !**

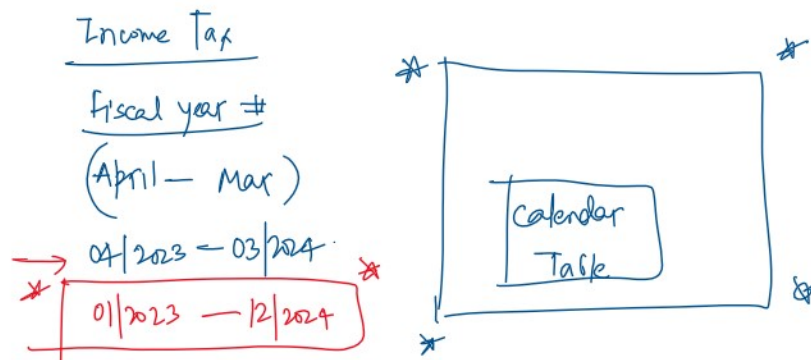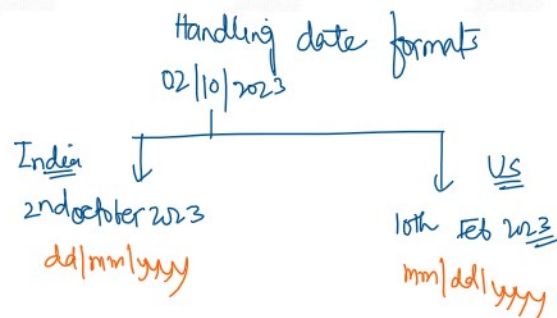## Why Data Manipulation?

There could be multiple problems with the data some of those are:

| Missing Values | Incorrect Format | Different Units | Unnecessary Data |

Sometimes data from different sources might be formatted differently such as having different date formats or formats of currency etc.

Handling date formats

02/10/2023

India → 2nd october 2023 → dd/mm/yyyy

US → 10th Feb 2023 → mm/dd/yyyy

Income Tax

fiscal year #

(April — Mar)

→ 04/2023 — 03/2024

01/2023 — 12/2024

Calendar Task

There could be multiple problems with the data some of those are:

| Missing Values | Incorrect Format | Different Units | Unnecessary Data |
|---|---|---|---|

Different sources might also have different Units of measurements such as temperature being measured in Celsius, Fahrenheit and Kelvin or distance being measured in Miles and Kilometers etc.

*kg and litres*

---

Sometimes we have a large dataset with columns that contain values that are not relevant to the tasks that you are performing. For example Some datasets have unique id columns which are not important.

*Feature engineering*
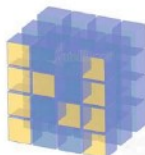
*Variables / Columns & importance*

*↓*

*Rest everything we drop.*

---

## What is NumPy?

The NumPy library is a very popular Python library and the abbreviation is "Numerical Python". The purpose of NumPy library is to do scientific computation and apply them to python applications.

NumPy

**How will NumPy help in Data Science?**

1. First of all, It is a open source Python library

2. It is fast because it is written in C and Python

3. In Python, there is no in-built array capabilities

4. You can use List as an alternative for arrays, but NumPy is better. But how is it better? We will discuss that now.