# Carbon_Emission_Prediction (AICTE_Internship - June 2025)

## Week 2:-

### Overview:

This phase aims to conduct an in-depth exploratory data analysis (EDA) to understand the structure, distribution, and relationships within the country-level dataset. This step is critical in identifying key trends, patterns, and anomalies before applying predictive modeling techniques to estimate $CO_2$ emissions.

**Data exploration**, also known as exploratory data analysis (EDA), provides numerous advantages that are essential for the success of any data-driven project. One of the most important benefits is that it allows analysts to develop a deep understanding of the dataset by examining its structure, distributions, and key characteristics. This includes identifying central tendencies, variability, missing values, and potential anomalies. By visualizing data through charts, histograms, scatter plots, and correlation matrices, patterns and relationships between variables become clearer, enabling more informed decision-making during the modeling phase.

Another major advantage is that data exploration helps uncover hidden trends, outliers, and inconsistencies that may otherwise go unnoticed. Recognizing these issues early allows for corrective actions such as data transformation, normalization, or removal of problematic records, which improves model robustness and accuracy. Additionally, EDA helps in feature selection by highlighting which variables are most relevant or redundant, reducing dimensionality and preventing overfitting in machine learning models.

Furthermore, data exploration supports better communication of insights to stakeholders. Visualizations and descriptive statistics provide intuitive summaries of complex datasets, making findings more accessible to non-technical audiences. Overall, EDA serves as a critical foundation for building reliable, interpretable, and high-performing models by ensuring that the underlying data is well-understood before any algorithm is applied.

## Challenges faced:

Having been exposed to proper **Machine Learning** only recently, I was dumbfounded by the next phase that came after Data Cleaning – **Data Exploration**.

There was just too much to learn at once, making my mind lag at times, but with repeated studying and even asking ChatGPT for help, I was able to understand many of the things and even acquire some more relevant knowledge on this subject.

Another common challenge is the overwhelming number of syntax and the ecosystem for libraries available. The libraries like **Numpy, Pandas, Seaborn, and Matplotlib** had so much in them that I struggled to remember.

I also struggled with setting up environments, managing dependencies, and dealing with issues like version conflicts. Especially when I was suddenly slapped with a screen full of red & incomprehensible errors, which made me grit my teeth in anger. Most of these errors happened due to version conflicts that were simply resolved using **pip commands**.

**Conclusion:**

All in all, it was a very satisfying journey, making me step into a whole new world of learning and experience. Although the path was thorny and full of challenges, but I achieved a satisfying sense of accomplishment after reaching the week end.