# Data Science Fundamentals

**Duration: 16 hours (2 days)**

**Pre-requisites**

- Basic Python programming skills

- Familiarity with Excel or CSV data formats

- Fundamental understanding of statistics (mean, median, variance)

- Optional: Prior exposure to pandas or NumPy

---

**Course Outcomes**

By the end of this 2-day workshop, participants will:

- Understand the data science lifecycle from data to insights

- Perform data cleaning, transformation, and analysis using Python

- Visualize data using matplotlib and seaborn

- Build and evaluate simple machine learning models

- Apply end-to-end data science workflows on real datasets

---

**Day 1: Data Exploration & Analysis**

1. **Introduction to Data Science**

    o What is Data Science?

    o Data Science lifecycle: Define → Acquire → Prepare → Analyze → Act

    o Tools of the trade (Python, Jupyter, pandas, scikit-learn)

2. **Python for Data Science Refresher**

    o Lists, dictionaries, functions, loops

    o Working with Jupyter Notebooks

3. **Data Manipulation with pandas**

    o Loading CSV, Excel files

    o DataFrames: filter, sort, group, merge

    o Handling missing values, duplicates

4. **Exploratory Data Analysis (EDA)**

   o Descriptive statistics (mean, median, std)

   o Value counts, correlation, outliers

   o Hands-on EDA with a real dataset

5. **Data Visualization**

   o Line, bar, pie, scatter, histogram

   o Using matplotlib and seaborn

   o Customizing plots for reports

6. **Hands-on Lab: Titanic Dataset Analysis**

   o Clean, explore, and visualize insights from Titanic dataset

---

**Day 2: Data Science in Action (Model Building)**

7. **Introduction to Machine Learning**

   o Supervised vs. Unsupervised learning

   o Classification vs. Regression

8. **Feature Engineering**

   o Encoding categorical variables

   o Scaling numerical values

   o Splitting data into train/test sets

9. **Model Building with scikit-learn**

   o Train a simple Logistic Regression model

   o Evaluate with accuracy, precision, recall

   o Use Confusion Matrix and ROC curve

10. **Model Improvement Techniques**

    o Cross-validation

    o Hyperparameter tuning (GridSearchCV)

    o Avoiding overfitting

11. **Real-World Use Case: Predict Customer Churn**

- Load and analyze a churn dataset
- Build and evaluate a classification model

12. **Deployment & Next Steps**
- Saving model with joblib
- Basic intro to deployment using Streamlit