

# A Training on Python for Data Analysis

From Fundamentals to Machine Learning

---

Nahiyan Bin Noor, MS

July 22, 2025

Data Analyst - Intermediate  
Institute for Digital Health & Innovation  
University of Arkansas for Medical Sciences

Introduction to Python

Generating a Sample Dataset

Exploratory Data Analysis (EDA)

Data Visualization

Machine Learning Models

# Introduction to Python

---

# What is Python?

- A **high-level, general-purpose** programming language.

# What is Python?

- A **high-level, general-purpose** programming language.
- Renowned for its **simple, clean syntax** that emphasizes readability.

# What is Python?

- A **high-level, general-purpose** programming language.
- Renowned for its **simple, clean syntax** that emphasizes readability.
- **Interpreted**, not compiled, which makes development faster and more interactive.

# What is Python?

- A **high-level, general-purpose** programming language.
- Renowned for its **simple, clean syntax** that emphasizes readability.
- **Interpreted**, not compiled, which makes development faster and more interactive.
- It's not just for data! Python is used for web development, automation, scientific computing, and much more.

# Why Python for Data Analysis?

Python has become the de-facto language for data science for several key reasons:

1. **Vast Ecosystem of Libraries:** This is its superpower. Specialized libraries provide powerful, pre-built functionality for nearly any data task.



# Why Python for Data Analysis?

Python has become the de-facto language for data science for several key reasons:

1. **Vast Ecosystem of Libraries:** This is its superpower. Specialized libraries provide powerful, pre-built functionality for nearly any data task.
2. **Ease of Learning:** Its straightforward syntax allows you to focus on solving problems rather than getting stuck on complex language rules.

# Why Python for Data Analysis?

Python has become the de-facto language for data science for several key reasons:

1. **Vast Ecosystem of Libraries:** This is its superpower. Specialized libraries provide powerful, pre-built functionality for nearly any data task.
2. **Ease of Learning:** Its straightforward syntax allows you to focus on solving problems rather than getting stuck on complex language rules.
3. **Large, Active Community:** If you have a problem, chances are someone has solved it and shared the solution online.

# Why Python for Data Analysis?

Python has become the de-facto language for data science for several key reasons:

1. **Vast Ecosystem of Libraries:** This is its superpower. Specialized libraries provide powerful, pre-built functionality for nearly any data task.
2. **Ease of Learning:** Its straightforward syntax allows you to focus on solving problems rather than getting stuck on complex language rules.
3. **Large, Active Community:** If you have a problem, chances are someone has solved it and shared the solution online.
4. **Integration and Versatility:** You can handle the entire data workflow—from data acquisition to machine learning—all in one language.

# Our Core Toolkit for Today

- **Pandas:** The essential library for data manipulation. It provides the **DataFrame**, a structure for handling tabular data.
  - *Our workhorse for importing, cleaning, filtering, and aggregating data.*

# Our Core Toolkit for Today

- **Pandas:** The essential library for data manipulation. It provides the **DataFrame**, a structure for handling tabular data.
  - *Our workhorse for importing, cleaning, filtering, and aggregating data.*
- **Matplotlib & Seaborn:** The primary libraries for creating a wide range of static and interactive visualizations.
  - *How we'll create histograms, bar charts, and scatter plots.*

# Our Core Toolkit for Today

- **Pandas:** The essential library for data manipulation. It provides the **DataFrame**, a structure for handling tabular data.
  - *Our workhorse for importing, cleaning, filtering, and aggregating data.*
- **Matplotlib & Seaborn:** The primary libraries for creating a wide range of static and interactive visualizations.
  - *How we'll create histograms, bar charts, and scatter plots.*
- **Scikit-learn & Statsmodels:** The go-to libraries for machine learning and statistical modeling.
  - *The tools we'll use to build our predictive and inferential models.*

# Generating a Sample Dataset

---

# We start by generating a synthetic dataset.

```
1  import pandas as pd
2  import numpy as np
3  import random
4
5  NUM_ROWS = 50000
6  sexes = ['Female', 'Male']
7  races = ['White', 'Black or African American', 'Asian', 'Other']
8
9  # Create a dictionary of data (simplified for slide)
10 data = {
11     'Sex': [random.choice(sexes) for _ in range(NUM_ROWS)],
12     'AgeOnIndexDate': np.random.randint(18, 85, size=NUM_ROWS),
13     'FirstRace': [random.choice(races) for _ in range(NUM_ROWS)],
14     'Depression': np.random.choice([0, 1], size=NUM_ROWS, p=[0.7, 0.3]),
15     'ChronicPain': np.random.choice([0, 1], size=NUM_ROWS, p=[0.6, 0.4]),
16     'ElixhauserScore': np.random.randint(0, 20, size=NUM_ROWS),
17     'NumberOfEdVisits': np.random.randint(0, 15, size=NUM_ROWS)
18 }
19 df = pd.DataFrame(data)
20 df.to_csv('synthetic_patient_data.csv', index=False)
21 print("Data generated and saved.")
```



# Exploratory Data Analysis (EDA)

---

## Load the data and get a high-level overview.

---

```
1 import pandas as pd
2 df = pd.read_csv('synthetic_patient_data.csv')
3
4 # Display first 5 rows
5 print("--- df.head() ---")
6 print(df.head())
7
8 # Display dataset info
9 print("\n--- df.info() ---")
10 df.info()
```

---

## Zoom in on specific subsets of the data.

---

```
1  # Select a single column
2  ages = df['AgeOnIndexDate']
3
4  # Select multiple columns
5  demographics = df[['Sex', 'AgeOnIndexDate', 'FirstRace']]
6
7  # Filter rows based on a condition
8  # Patients with an Elixhauser Score greater than 15
9  high_risk_patients = df[df['ElixhauserScore'] > 15]
10
11 # Filter with multiple conditions (AND: &)
12 # Female patients with chronic pain
13 female_chronic_pain = df[(df['Sex'] == 'Female') & (df['ChronicPain'] == 1)]
```

---

## Calculate summary statistics for columns.

---

```
1  # Get summary statistics for all numerical columns
2  # Provides count, mean, std, min, max, and percentiles.
3  print(df.describe())
4
5  # Count occurrences in a categorical column
6  print("\n--- Value Counts for FirstRace ---")
7  print(df['FirstRace'].value_counts())
```

---

# Use groupby to split-apply-combine.

---

```
1  # Group by one column and calculate a single metric
2  # What is the average age for each sex?
3  avg_age_by_sex = df.groupby('Sex')['AgeOnIndexDate'].mean()
4  print(avg_age_by_sex)
5
6
7  # Group by a column and apply multiple aggregations
8  # Get stats for ElixhauserScore for each race
9  race_summary = df.groupby('FirstRace')['ElixhauserScore'].agg(
10     ['mean', 'max', 'count']
11 )
12 print(race_summary)
```

---

# Data Visualization

---

# Use a histogram to see the spread of a numerical variable.

---

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 sns.set_theme(style="whitegrid")
4
5 # Histogram of patient ages
6 plt.figure(figsize=(10, 6))
7 sns.histplot(data=df, x='AgeOnIndexDate', bins=30, kde=True)
8 plt.title('Distribution of Patient Age')
9 plt.xlabel('Age')
10 plt.ylabel('Frequency')
11 plt.show()
```

---

## Use a count plot to see the frequency of categories.

---

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 sns.set_theme(style="whitegrid")
4
5 # Count plot of patients by race
6 plt.figure(figsize=(10, 6))
7 sns.countplot(data=df, y='FirstRace',
8               order=df['FirstRace'].value_counts().index)
9 plt.title('Number of Patients by Race')
10 plt.xlabel('Count')
11 plt.ylabel('Race')
12 plt.show()
```

---



# Machine Learning Models

---

# Convert categorical data to numbers and split into train/test sets.

---

```
1  from sklearn.model_selection import train_test_split
2  from sklearn.preprocessing import StandardScaler
3
4  # Convert categorical variables into dummy/indicator variables
5  df_processed = pd.get_dummies(df, drop_first=True)
6
7  # Define Features (X) and Target (y)
8  y = df_processed['ReceivedMOUD']
9  X = df_processed.drop('ReceivedMOUD', axis=1)
10
11 # Split data into training and testing sets
12 X_train, X_test, y_train, y_test = train_test_split(
13     X, y, test_size=0.3, random_state=42, stratify=y)
14
15 # Scale numerical features
16 scaler = StandardScaler()
17 X_train = scaler.fit_transform(X_train)
18 X_test = scaler.transform(X_test)
```

---

# Predicting a binary outcome: 'ReceivedMOUD'.

---

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.metrics import classification_report, accuracy_score
3
4 # Train the model
5 log_reg = LogisticRegression(random_state=42, max_iter=1000)
6 log_reg.fit(X_train, y_train)
7
8 # Make predictions and evaluate
9 y_pred = log_reg.predict(X_test)
10 print(f"Accuracy: {accuracy_score(y_test, y_pred):.4f}")
11 print(classification_report(y_test, y_pred))
```

---

# Using a more powerful model for prediction.

---

```
1  from sklearn.ensemble import RandomForestClassifier
2
3  # Train the model
4  rf_clf = RandomForestClassifier(n_estimators=100, random_state=42, n_jobs=-1)
5  rf_clf.fit(X_train, y_train)
6
7  # Make predictions and evaluate
8  y_pred_rf = rf_clf.predict(X_test)
9  print(f"Accuracy: {accuracy_score(y_test, y_pred_rf):.4f}")
10 print(classification_report(y_test, y_pred_rf))
```

---

# Goal: Understand the effect of a predictor on a continuous outcome.

---

```
1 import statsmodels.api as sm
2
3 # Use the original, unscaled dataframe for easier interpretation
4 # Y = continuous outcome, X = predictors
5 Y_lin = df['NumberOfEdvisitUniqueEncounter']
6 X_lin = df[['totalMorphineDose', 'AgeOnIndexDate']]
7
8 # We must add a constant (intercept) to our predictors
9 X_lin = sm.add_constant(X_lin)
10
11 # Fit the Ordinary Least Squares (OLS) model
12 model_lin = sm.OLS(Y_lin, X_lin).fit()
13
14 # Print the detailed summary
15 print(model_lin.summary())
```

# Goal: Model a count outcome, like number of ED visits.

---

```
1  import statsmodels.api as sm
2
3  # Y = count outcome, X = predictors
4  Y_count = df['NumberOfEdvisitUniqueEncounter']
5  X_count = df[['ChronicPain', 'AgeOnIndexDate', 'Depression']]
6  X_count = sm.add_constant(X_count)
7
8  # Fit the Poisson model
9  # We use .GLM (Generalized Linear Model) with the Poisson family
10 poisson_model = sm.GLM(Y_count, X_count, family=sm.families.Poisson()).fit()
11
12 # Print the summary
13 print(poisson_model.summary())
```

---

# Goal: Model a count outcome with overdispersion.

---

```
1 import statsmodels.api as sm
2
3 # Use the same Y and X as the Poisson model
4 Y_count = df['NumberOfEdvisitUniqueEncounter']
5 X_count = df[['ChronicPain', 'AgeOnIndexDate', 'Depression']]
6 X_count = sm.add_constant(X_count)
7
8 # Fit the Negative Binomial model
9 # A good alternative to Poisson if the variance of the count is
10 # much larger than its mean (overdispersion).
11 neg_bin_model = sm.GLM(Y_count, X_count,
12                         family=sm.families.NegativeBinomial()).fit()
13
14 # Print the summary
15 print(neg_bin_model.summary())
```

---

# Modeling a binary outcome for causal inference.

---

```
1 import pandas as pd
2 import statsmodels.api as sm
3
4 # Define the dependent (Y) and independent (X) variables
5 Y = df['ReceivedMOUD'] # Binary outcome (0 or 1)
6 X = df[['ChronicPain', 'AgeOnIndexDate', 'Depression']]
7
8 # Add a constant (intercept) to the model
9 X = sm.add_constant(X)
10
11 # Instantiate and fit the Logit model
12 logit_model = sm.Logit(Y, X).fit()
13
14 # Print the detailed statistical summary
15 print(logit_model.summary())
```

---



# OLS vs Poisson vs Negative Binomial Regression

Feature	OLS	Poisson	Negative Binomial
Outcome Type	Continuous	Count	Count (overdispersed)
Data Example	Blood pressure, BMI	Number of ER visits	Number of ER visits
Assumes	Normality, constant variance	Mean = Variance	Variance = Mean
Common Pitfall	Poor with skewed data	Poor fit with overdispersion	Better fit for overdispersed data
Interpretation	Change in mean outcome	Change in event rate (IRR)	Change in event rate (IRR)
Model Output	Linear coefficient ( $\beta$ )	Incidence Rate Ratio (IRR)	Incidence Rate Ratio (IRR)

# Questions?

All codes and training materials will be found here:  
<https://github.com/Nahiyan140212/RTEC-Training>