

The background is a dark blue gradient. It features several vertical white lines of varying lengths. Scattered throughout are small squares in teal, orange, and pink. Some squares are solid, while others are white outlines.

DATA Preprocessing

The way to Machine Learning

Abdenour Bouziane
Baaziz Mohammed Chaker

CONTENTS OF THIS SESSION

Here's what you'll find in this [slides](#) :

1. Introduction to Python
2. Introduction to Data Preprocessing :
 - a. Data cleaning
 - b. Data transformation
 - c. Data reduction
 - d. Data discretization
3. Hands-on Examples
 - Importing libraries: [pandas](#), [numpy](#)
 - Loading data: CSV, Excel, etc.
 - Handling missing values:
 - Data transformation:
 - Data reduction:
 - Data discretization:
4. Conclusion and Next Steps

Why Python for Data Preprocessing?

- Python has a rich ecosystem of libraries specifically designed for data manipulation and analysis, such as `pandas`, `NumPy`, and `scikit-learn`.
- It offers `easy-to-understand syntax`, making it accessible for beginners without sacrificing power and flexibility.
- Python's popularity in the `data science` community means there is extensive documentation, tutorials, and community support available.

Basic Python Syntax

1. Variables and Data Types:

```
main.py  +
1  # Variables
2  x = 10
3  name = "Abdou"
4  is_student = True
5
6  # Data types
7  integer_num = 10
8  float_num = 3.14
9  string_text = "Hello, world!"
10 boolean_value = True
```

Basic Python Syntax

2. Basic Arithmetic Operations:

```
main.py +
1  # Arithmetic operations
2  a = 10
3  b = 5
4  # Addition
5  sum_result = a + b
6  # Subtraction
7  difference_result = a - b
8  # Multiplication
9  product_result = a * b
10 # Division
11 division_result = a / b
12 # Modulus (remainder)
13 modulus_result = a % b
14 # Exponentiation
15 exponentiation_result = a ** b
16
```

Basic Python Syntax

3. Basic Data Structures:

```
main.py +
1 # Lists
2 my_list = [1, 2, 3, 4, 5]
3
4 # Tuples
5 my_tuple = (1, 2, 3)
6
7 # Dictionaries
8 my_dict = {"name": "John", "age": 30, "city": "New York"}
9
```

Basic Python Syntax

4. Functions:

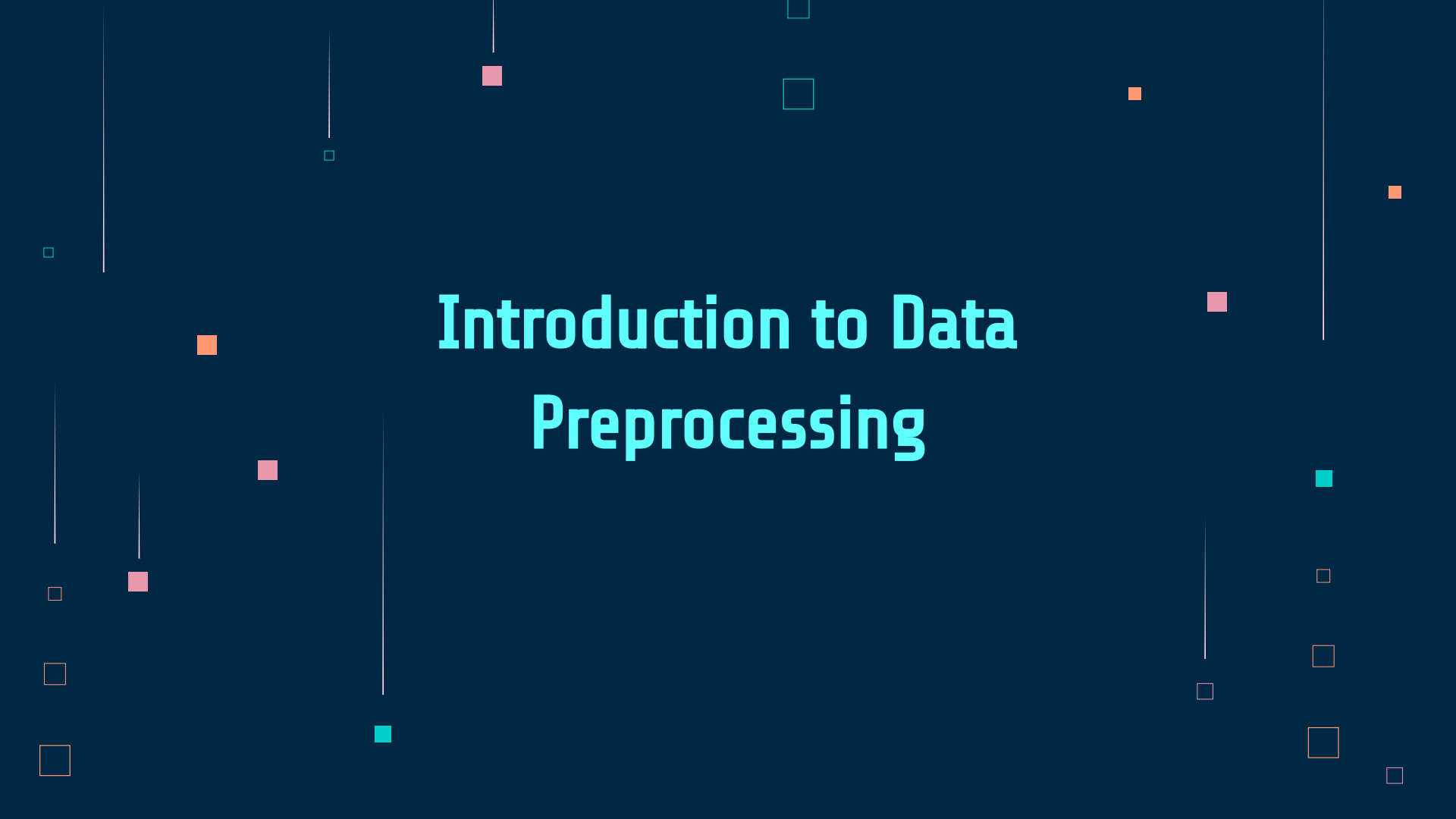
```
main.py +
1 # Function definition
2 def greet(name):
3     return "Hello, " + name + "!"
4
5 # Function call
6 greeting = greet("Alice")
7 print(greeting)
```

Basic Python Syntax

5. Control Flow Statements:

```
main.py +
1 # If-else statement
2 x = 10
3 if x > 5:
4     print("x is greater than 5")
5 else:
6     print("x is not greater than 5")
7
8 # For loop
9 for i in range(5):
10     print(i)
11
12 # While loop
13 counter = 0
14 while counter < 5:
15     print(counter)
16     counter += 1
17
```


Introduction to Data Preprocessing



What is Data Preprocessing ?

Data preprocessing is a crucial step in the data analysis pipeline that involves **cleaning**, **transforming**, and preparing raw data into a format suitable for further analysis. It aims to improve the quality of data, **enhance** its interpretability, and facilitate the performance of machine learning **algorithms**.

Why is Data Preprocessing important ?

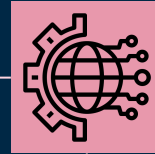
- **Improves Data Quality:** Raw data often contains errors, outliers, and missing values. Preprocessing helps identify and rectify these issues, ensuring that the data is accurate and reliable.
- **Enhances Interpretability:** Preprocessed data is easier to understand and interpret, making it more accessible for analysis and decision-making.
- **Prepares Data for Analysis:** Many machine learning algorithms require data to be in a specific format. Preprocessing ensures that data meets these requirements, enabling the application of various analytical techniques.

Steps involved in Data Preprocessing

Data Cleaning



Data
Transformation



Data Exploration



Data
Discretization



1. Data Exploration

Data exploration is the process of examining and understanding the structure, patterns, and relationships within a dataset before performing further analysis. It involves summarizing the data using descriptive statistics and visualizing it through charts and graphs.

```
1  import pandas as pd
2
3  # Load data
4  data = pd.read_csv('data.csv')
5
6  # Data Exploration
7  summary_stats = data.describe()
8  data.head()
9  data.tail()
10 data.columns
11 data.dtypes
12 print(data['data column'].describe())
13 print(data['column'])
14 data.shape()
15
16
```

2. Data Cleaning

Identifying and handling missing values, outliers, and inconsistencies in the data.

```
main.py +  
1 import pandas as pd  
2  
3 # Load data  
4 data = pd.read_csv('data.csv')  
5  
6 # Check for missing values  
7 missing_values = data.isnull().sum()  
8  
9 # Handle missing values (e.g., fill with mean)  
10 data.fillna(data.mean(), inplace=True)  
11  
12 # Remove duplicate records  
13 data.drop_duplicates(inplace=True)  
14
```

Binning Technique

Data binning, bucketing is a data pre-processing method used to minimize the effects of small observation errors. The original data values are divided into small intervals known as bins and then they are replaced by a general value calculated for that bin. This has a smoothing effect on the input data and may also reduce the chances of overfitting in the case of small datasets

There are 2 methods of dividing data into bins:

1. Equal Frequency Binning: bins have an equal frequency.
2. Equal Width Binning : bins have equal width with a range of each bin are defined as $[\min + w]$, $[\min + 2w]$ $[\min + nw]$ where $w = (\max - \min) / (\text{no of bins})$.

Binning Technique

1. Equal frequency

Suppose we have a dataset with the following values:

Values=[5,8,12,15,18,20,22,25,28,30]

We want to perform equal frequency binning with 3 bins.

Binning Technique

1. Equal frequency

1. **Sort the Data:**

Sorted Values=[5,8,12,15,18,20,22,25,28,30]

2. **Determine Number of Observations per Bin:**

3. Total Observations=10 | Target Frequency per Bin= $10/3 \approx 3.33$

4. **Identify Bin Boundaries:**

- Bin 1: [5,12)
- Bin 2: [12,22)
- Bin 3: [22,30]

Binning Technique

1. Equal frequency

5. Assign Values to Bins:

- Values 5,8,12 belong to Bin 1.
- Values 15,18,20 belong to Bin 2.
- Values 22,25,28,30 belong to Bin 3.

6. Handle Ties:

- Ties can be handled by either assigning them to the same bin or distributing them evenly across adjacent bins, depending on the implementation.

Binning Technique

1. Equal Width

1. Calculate Range of Values:

$$\text{Range} = \text{Max}(\text{Values}) - \text{Min}(\text{Values}) = 30 - 5 = 25$$

2. Calculate Bin Width: $\text{Bin Width} = 25/3 \approx 8.33$

3. Define Bin Edges:

- Bin 1: [5,13.33)
- Bin 3: [21.67,30]

Binning Technique

1. Equal Width

4. Assign Values to Bins:

- i. Values 5,8,12 belong to Bin 1.
- ii. Values 15,18,20 belong to Bin 2.
- iii. Values 22,25,30 belong to Bin 3.

5. Handle Edge Cases:

- iv. Values equal to the minimum or maximum may be included in the bin adjacent to them or assigned to a special bin, depending on the implementation.

3. Data Transformation

Converting and standardizing data into a suitable format for analysis. This may involve encoding categorical variables, scaling numerical features, and normalizing data distributions.

main.py



```
1 from sklearn.preprocessing import StandardScaler, LabelEncoder
2
3 # Encode categorical variables
4 label_encoder = LabelEncoder()
5 data['category_encoded'] = label_encoder.fit_transform(data['category'])
6
7 # Scale numerical features
8 scaler = StandardScaler()
9 data[['numerical_feature1', 'numerical_feature2']] = scaler.fit_transform(data[['numerical_feature1', 'numerical_feature2']])
10
```

4. Data Discretization

Grouping continuous data into discrete intervals or categories.

main.py



```
1 # Bin numerical data
2 data['binned_feature'] = pd.cut(data['numerical_feature'], bins=3, labels=['Low', 'Medium', 'High'])
3 |
```

Do you have any questions?

abdenour.bouziane@ensia.edu.dz

mohammed.baaziz@ensia.edu.dz

THANKS

Resources

Python : <https://github.com/AbderrahimRezki/PythonWorkshop.git>

More about Data preprocessing :
<https://www.geeksforgeeks.org/data-preprocessing-in-data-mining/>

Github repository to work on :
<https://github.com/AbdenourBouziane/ETC-Data-Preprocessing-Workshop>