

# Chương 3: Giới thiệu về Dữ liệu

## Contents

- Giới thiệu về Dữ liệu
- 3.1 Dữ liệu là gì?
- 3.2 Các nguồn dữ liệu phổ biến
- 3.3 Cách làm việc với dữ liệu trong Python
- 3.4 Làm sạch dữ liệu
- 3.5 Tổng hợp dữ liệu
- 3.6 Xuất dữ liệu
- 3.7 Các loại dữ liệu nâng cao
- 3.8 Các công cụ và phần mềm hỗ trợ làm việc với dữ liệu
- 3.9 Các phương pháp phân tích dữ liệu
- 3.10 Thực hành nâng cao
- 3.11 Các thách thức trong làm việc với dữ liệu
- 3.12 Tương lai của dữ liệu
- 3.13 Thực hành
- 3.14 Tổng kết chương 3

## Giới thiệu về Dữ liệu

Dữ liệu là yếu tố trung tâm trong mọi dự án phân tích và trực quan hóa. Chương này sẽ giúp bạn hiểu các khái niệm cơ bản về dữ liệu, cách làm việc với dữ liệu trong Python và chuẩn bị dữ liệu cho trực quan hóa.

## 3.1 Dữ liệu là gì?

Dữ liệu là tập hợp các giá trị được thu thập để nghiên cứu hoặc phân tích. Dữ liệu có thể tồn tại dưới nhiều dạng khác nhau, như:

**Dữ liệu số (Numerical):** Số liệu có thể đo lường (ví dụ: chiều cao, cân nặng). **Dữ liệu phân loại (Categorical):** Dữ liệu thể hiện các nhóm hoặc loại (ví dụ: giới tính, màu sắc). **Dữ liệu theo thời gian (Time-series):** Dữ liệu được thu thập theo thời gian (ví dụ: giá cổ phiếu hàng ngày).

## 3.2 Các nguồn dữ liệu phổ biến

Dữ liệu có thể được lấy từ nhiều nguồn khác nhau, bao gồm:

**Tập tin:** CSV, Excel, JSON, v.v. **Cơ sở dữ liệu:** MySQL, PostgreSQL, MongoDB, v.v. **API:** Các dịch vụ cung cấp dữ liệu qua giao thức HTTP (ví dụ: Google Maps API, Twitter API). **Web scraping:** Thu thập dữ liệu từ các trang web.

## 3.3 Cách làm việc với dữ liệu trong Python

Python cung cấp nhiều thư viện mạnh mẽ để xử lý và phân tích dữ liệu. Trong chương này, chúng ta sẽ tập trung vào hai thư viện chính: **Pandas** và **NumPy**.

**3.3.1 Thư viện Pandas** Pandas là thư viện phổ biến nhất để thao tác với dữ liệu dạng bảng (DataFrame).

**Cài đặt Pandas:**

```
pip install pandas
```

**Đọc dữ liệu từ tệp CSV:**

```
import pandas as pd

# Đọc dữ liệu từ tệp CSV
data = pd.read_csv('example.csv')

# Hiển thị 5 dòng đầu tiên
print(data.head())
```

### Giải thích:

- `pd.read_csv('example.csv')`: Đọc dữ liệu từ tệp `example.csv` và lưu vào biến `data` dưới dạng `DataFrame`.
- `data.head()`: Hiển thị 5 dòng đầu tiên của `DataFrame`, giúp bạn xem sơ qua dữ liệu.

### Một số thao tác cơ bản với Pandas:

```
# Hiển thị thông tin tổng quan về dữ liệu
print(data.info())

# Hiển thị thống kê cơ bản của dữ liệu
print(data.describe())

# Lọc dữ liệu theo điều kiện
filtered_data = data[data['column_name'] > 50]
```

### Giải thích:

- [`data.info\(\)`](#): Hiển thị thông tin tổng quan về `DataFrame` như số lượng dòng, cột và kiểu dữ liệu.
- `data.describe()`: Tính toán các thống kê cơ bản (mean, std, min, max) cho dữ liệu số trong `DataFrame`.
- `data[data['column_name'] > 50]`: Lọc các dòng có giá trị trong cột `column_name` lớn hơn 50.

### 3.3.2 Thư viện NumPy NumPy hỗ trợ làm việc với các mảng số học hiệu quả.

#### Cài đặt NumPy:

```
pip install numpy
```

#### Tạo và thao tác với mảng:

```
import numpy as np

# Tạo mảng
array = np.array([1, 2, 3, 4, 5])

# Tính toán cơ bản
print("Mean:", np.mean(array))
print("Standard Deviation:", np.std(array))
```

**Giải thích:**

- `np.array([1, 2, 3, 4, 5])`: Tạo một mảng NumPy từ danh sách Python.
- `np.mean(array)`: Tính giá trị trung bình của mảng.
- `np.std(array)`: Tính độ lệch chuẩn của mảng.

## 3.4 Làm sạch dữ liệu

Dữ liệu thực tế thường không hoàn hảo và cần được làm sạch trước khi trực quan hóa.

**Các bước làm sạch dữ liệu:****1. Xử lý giá trị thiếu:**

- Xóa các dòng hoặc cột có giá trị thiếu:

```
data = data.dropna()
```

- Điền giá trị thay thế:

```
data['column_name'] = data['column_name'].fillna(0)
```

**2. Xử lý dữ liệu không hợp lệ:**

- Loại bỏ các giá trị không hợp lệ hoặc ngoại lệ (outliers).

**3. Chuyển đổi kiểu dữ liệu:**

- Chuyển đổi cột sang kiểu số hoặc chuỗi phù hợp:

```
data['column_name'] = data['column_name'].astype(float)
```

#### 4. Chuẩn hóa dữ liệu:

- Đổi tên cột để dễ sử dụng:

```
data.rename(columns={'OldName': 'NewName'}, inplace=True)
```

## 3.5 Tổng hợp dữ liệu

Việc tổng hợp giúp rút ra thông tin ý nghĩa từ dữ liệu. Pandas hỗ trợ các phương pháp sau:

- **Nhóm dữ liệu:**

```
grouped = data.groupby('category_column').mean()  
print(grouped)
```

#### Giải thích:

- `data.groupby('category_column')`: Nhóm dữ liệu theo cột `category_column`.
- `.mean()`: Tính giá trị trung bình của các nhóm.

#### Tạo Pivot Table:

```
pivot_table = data.pivot_table(values='value_column', index='category_column', columns=  
print(pivot_table)
```

#### Giải thích:

- `pivot_table()`: Tạo bảng tổng hợp dữ liệu, sử dụng `category_column` làm chỉ mục và `other_column` làm các cột.
- `aggfunc='sum'`: Tính tổng giá trị cho mỗi nhóm.

## 3.6 Xuất dữ liệu

Sau khi làm việc với dữ liệu, bạn có thể xuất kết quả ra các định dạng khác nhau:

```
# Xuất dữ liệu ra CSV
data.to_csv('output.csv', index=False)

# Xuất dữ liệu ra Excel
data.to_excel('output.xlsx', index=False)
```

### Giải thích:

- `data.to_csv('output.csv')`: Lưu DataFrame vào tệp CSV.
- `data.to_excel('output.xlsx')`: Lưu DataFrame vào tệp Excel.

## 3.7 Các loại dữ liệu nâng cao

- **Dữ liệu phi cấu trúc (Unstructured Data)**: Giới thiệu về dữ liệu không có cấu trúc rõ ràng như văn bản, hình ảnh, video.
- **Dữ liệu bán cấu trúc (Semi-structured Data)**: Ví dụ về dữ liệu như XML, JSON mà có cấu trúc nhưng không hoàn toàn theo định dạng bảng.

## 3.8 Các công cụ và phần mềm hỗ trợ làm việc với dữ liệu

- Giới thiệu về các công cụ như Jupyter Notebook, Google Colab, và các phần mềm như Tableau, Power BI cho việc phân tích và trực quan hóa dữ liệu.

## 3.9 Các phương pháp phân tích dữ liệu

- **Phân tích mô tả (Descriptive Analysis)**: Cách sử dụng thống kê mô tả để tóm tắt dữ liệu.
- **Phân tích dự đoán (Predictive Analysis)**: Giới thiệu về các mô hình học máy cơ bản để dự đoán xu hướng tương lai từ dữ liệu hiện tại.

## 3.10 Thực hành nâng cao

- Cung cấp một số bài tập thực hành nâng cao hơn, chẳng hạn như:
  - Sử dụng Pandas để thực hiện phân tích dữ liệu từ một tệp CSV lớn.
  - Tạo biểu đồ trực quan hóa dữ liệu bằng Matplotlib hoặc Seaborn từ dữ liệu đã xử lý.

## 3.11 Các thách thức trong làm việc với dữ liệu

- Thảo luận về các vấn đề thường gặp như dữ liệu không đầy đủ, dữ liệu không chính xác, và cách giải quyết chúng.

## 3.12 Tương lai của dữ liệu

- Một cái nhìn tổng quan về xu hướng tương lai trong lĩnh vực dữ liệu, như Big Data, Machine Learning, và AI.

## 3.13 Thực hành

Hãy thử các bài tập sau để làm quen với dữ liệu:

1. Đọc tệp data.csv và hiển thị 5 dòng đầu tiên.
2. Kiểm tra xem tệp có giá trị thiếu không và xử lý chúng.
3. Tạo một Pivot Table để tính trung bình một cột theo từng nhóm.

## 3.14 Tổng kết chương 3

**Nội dung:** Trong chương 3 này đã giới thiệu các khái niệm cơ bản và công cụ làm việc với dữ liệu trong Python. Trong chương 4, chúng ta sẽ tìm hiểu về thống kê mô tả và cách trực quan hóa các đặc điểm cơ bản của dữ liệu.