



# Python Cơ Bản Cho Data Science

Kiểu Dữ Liệu Cơ Bản: String, Int và Float

---

Quanskill

Giới Thiệu

String (Chuỗi)

Integer (Số Nguyên)

Float (Số Thập Phân)

Chuyển Đổi Kiểu Dữ Liệu

Ứng Dụng Thực Tế

Tóm Tắt và Luyện Tập

# **Giới Thiệu**

---

- Trong Data Science, chúng ta xử lý nhiều loại dữ liệu khác nhau
- Hiểu kiểu dữ liệu giúp chọn đúng phương pháp phân tích
- Python có các kiểu dữ liệu cơ bản phù hợp cho từng mục đích

## Ba Kiểu Dữ Liệu Cơ Bản

- **String (str)**: Văn bản, nhãn, danh mục
- **Integer (int)**: Số nguyên, đếm, chỉ số
- **Float**: Số thập phân, đo lường, tính toán

Python cung cấp hàm `type()` để kiểm tra kiểu dữ liệu:

```
1 # Check data types in Python
2 name = "Alice"
3 age = 25
4 height = 1.75
5
6 print("Type of name:", type(name))
7 print("Type of age:", type(age))
8 print("Type of height:", type(height))
```

## Output:

```
Type of name: <class 'str'>
Type of age: <class 'int'>
Type of height: <class 'float'>
```

# **String (Chuỗi)**

---

- String (chuỗi) là kiểu dữ liệu văn bản
- Được bao quanh bởi dấu ngoặc kép "" hoặc đơn ''
- Rất quan trọng trong phân tích dữ liệu văn bản (text analytics)

## Ví Dụ Trong Data Science

- Tên khách hàng: "Nguyễn Văn A"
- Danh mục sản phẩm: "Electronics", "Clothing"
- Bình luận khách hàng: "Sản phẩm rất tốt!"
- Địa chỉ email: "user@example.com"

```
1 # Creating strings - different ways
2 customer_name = "John Doe"
3 product_category = 'Electronics'
4 review = """This product is amazing!
5 I would definitely recommend it."""
6
7 print("Customer:", customer_name)
8 print("Category:", product_category)
9 print("Review:", review)
```

## Output:

Customer: John Doe

Category: Electronics

Review: This product is amazing!

I would definitely recommend it.

```
1 # String concatenation for data processing
2 first_name = "Maria"
3 last_name = "Garcia"
4
5 # Method 1: Using + operator
6 full_name = first_name + " " + last_name
7 print("Full name:", full_name)
8
9 # Method 2: Using f-strings (recommended)
10 customer_id = 12345
11 message = f"Welcome {full_name}! Your ID is {customer_id}"
12 print("Message:", message)
```

## Output:

Full name: Maria Garcia

Message: Welcome Maria Garcia! Your ID is 12345

```
1 # Useful string methods for data cleaning
2 email = " JOHN.DOE@GMAIL.COM "
3
4 # Clean and standardize the email
5 clean_email = email.strip().lower()
6 print("Original:", f"{email}")
7 print("Cleaned:", clean_email)
8
9 # Extract information
10 domain = clean_email.split("@")[1]
11 print("Domain:", domain)
12
13 # Check content
14 print("Contains gmail:", "gmail" in clean_email)
```

## Output:

```
Original: ' JOHN.DOE@GMAIL.COM '
Cleaned: john.doe@gmail.com
Domain: gmail.com
Contains gmail: True
```

# Integer (Số Nguyên)

---

- Integer (int) là số nguyên: ..., -2, -1, 0, 1, 2, ...
- Không có phần thập phân
- Dùng để đếm, chỉ số, ID, ranking

## Ví Dụ Trong Data Science

- Số lượng sản phẩm: 150
- ID khách hàng: 12345
- Tuổi: 25
- Số lần click: 1250
- Ranking: 1, 2, 3, ...

```
1 # Creating integers for data analysis
2 total_customers = 1000
3 new_signups = 45
4 churned_customers = 12
5
6 print("Total customers:", total_customers)
7 print("Data type:", type(total_customers))
8
9 # Calculate retention
10 active_customers = total_customers + new_signups - churned_customers
11 print("Active customers:", active_customers)
12
13 # Calculate percentages (will be float)
14 churn_rate = churned_customers / total_customers * 100
15 print("Churn rate:", churn_rate, "%")
```

## Output:

Total customers: 1000  
Data type: <class 'int'>  
Active customers: 1033  
Churn rate: 1.2 %

```
1 # Integer operations for data analysis
2 sales_q1 = 250
3 sales_q2 = 300
4 sales_q3 = 275
5 sales_q4 = 425
6
7 # Basic arithmetic
8 total_sales = sales_q1 + sales_q2 + sales_q3 + sales_q4
9 average_sales = total_sales // 4 # Integer division
10 remainder = total_sales % 4      # Remainder
11
12 print("Total sales:", total_sales)
13 print("Average (integer division):", average_sales)
14 print("Remainder:", remainder)
15
16 # Comparison
17 best_quarter = max(sales_q1, sales_q2, sales_q3, sales_q4)
18 print("Best quarter sales:", best_quarter)
```

## Output:

```
Total sales: 1250
Average (integer division): 312
Remainder: 2
Best quarter sales: 425
```

# Float (Số Thập Phân)

---

- Float là số thập phân (floating point number)
- Có phần nguyên và phần thập phân
- Dùng cho các phép đo lường chính xác, tỷ lệ phần trăm

## Ví Dụ Trong Data Science

- Giá sản phẩm: 29.99
- Tỷ lệ conversion: 0.035 (3.5%)
- Chiều cao: 1.75 (mét)
- Điểm rating: 4.5/5.0
- Temperature: 25.6°C

```
1 # Creating floats for precise calculations
2 product_price = 29.99
3 discount_rate = 0.15    # 15%
4 tax_rate = 0.08        # 8%
5
6 print("Original price: $", product_price)
7 print("Data type:", type(product_price))
8
9 # Calculate final price
10 discount_amount = product_price * discount_rate
11 discounted_price = product_price - discount_amount
12 tax_amount = discounted_price * tax_rate
13 final_price = discounted_price + tax_amount
14
15 print("Discount amount: $", round(discount_amount, 2))
16 print("Final price: $", round(final_price, 2))
```

## Output:

Original price: \$ 29.99  
Data type: <class 'float'>  
Discount amount: \$ 4.5  
Final price: \$ 27.59

```
1 # Float operations for statistical analysis
2 temperatures = [23.5, 25.1, 22.8, 26.3, 24.7]
3
4 # Calculate statistics
5 total_temp = sum(temperatures)
6 count = len(temperatures)
7 average_temp = total_temp / count
8
9 print("Temperature readings:", temperatures)
10 print("Total:", total_temp)
11 print("Count:", count)
12 print("Average:", round(average_temp, 2), " °C")
13
14 # Find min and max
15 min_temp = min(temperatures)
16 max_temp = max(temperatures)
17 range_temp = max_temp - min_temp
18
19 print("Min temperature:", min_temp, " °C")
20 print("Max temperature:", max_temp, " °C")
21 print("Temperature range:", range_temp, " °C")
```

## Output:

Temperature readings: [23.5, 25.1, 22.8, 26.3, 24.7]

Total: 122.4

Count: 5

Average: 24.48 °C

Min temperature: 22.8 °C

Max temperature: 26.3 °C

## **Chuyển Đổi Kiểu Dữ Liệu**

---

```
1 # Type conversion in data processing
2 user_input = "123"          # String from user input
3 price_str = "45.99"         # String from CSV file
4 count_float = 10.0          # Float from calculation
5
6 print("Original types:")
7 print("user_input:", type(user_input), "value:", user_input)
8 print("price_str:", type(price_str), "value:", price_str)
9 print("count_float:", type(count_float), "value:", count_float)
10
11 # Convert types
12 user_id = int(user_input)      # String to int
13 product_price = float(price_str) # String to float
14 item_count = int(count_float)   # Float to int
15
16 print("\nAfter conversion:")
17 print("user_id:", type(user_id), "value:", user_id)
18 print("product_price:", type(product_price), "value:", product_price)
19 print("item_count:", type(item_count), "value:", item_count)
```

```
1 # Handling conversion errors safely
2 data_inputs = ["123", "45.99", "invalid", "67.5", "abc"]
3
4 print("Processing data inputs:")
5 for i, input_value in enumerate(data_inputs):
6     print(f"\nInput {i+1}: '{input_value}'")
7
8     try:
9         # Try to convert to float first
10        as_float = float(input_value)
11        print("  As float:", as_float)
12
13        # Check if it's actually an integer
14        if as_float.is_integer():
15            as_int = int(as_float)
16            print("  As integer:", as_int)
17
18    except ValueError:
19        print("  Error: Cannot convert to number")
20        print("  Keeping as string:", input_value)
```

## Ứng Dụng Thực Tế

---

```
1 # Customer data analysis example
2 customers = [
3     {"name": "Alice Johnson", "age": 28, "spent": 1250.50},
4     {"name": "Bob Smith", "age": 35, "spent": 890.25},
5     {"name": "Carol Davis", "age": 42, "spent": 2100.75}
6 ]
7
8 print("Customer Analysis Report")
9 print("=" * 30)
10
11 total_spent = 0.0
12 total_customers = len(customers)
13
14 for customer in customers:
15     name = customer["name"]           # String
16     age = customer["age"]            # Integer
17     spent = customer["spent"]       # Float
18
19     print(f"\nCustomer: {name}")
20     print(f"Age: {age} years old")
21     print(f"Total spent: ${spent}")
22
23     total_spent += spent
24
25 average_spent = total_spent / total_customers
26 print(f"\nSummary:")
27 print(f"Total customers: {total_customers}")
28 print(f"Total revenue: ${round(total_spent, 2)}")
29 print(f"Average per customer: ${round(average_spent, 2)})")
```

## Output:

### Customer Analysis Report

---

Customer: Alice Johnson  
Age: 28 years old  
Total spent: \$1250.5

Customer: Bob Smith  
Age: 35 years old  
Total spent: \$890.25

Customer: Carol Davis  
Age: 42 years old  
Total spent: \$2100.75

Summary:  
Total customers: 3  
Total revenue: \$4241.5  
Average per customer: \$1413

## Tóm Tắt và Luyện Tập

---

**Bảng 1:** So sánh các kiểu dữ liệu cơ bản

Kiểu	Mô tả	Ví dụ
String (str)	Văn bản, ký tự	"Alice", 'Product A'
Integer (int)	Số nguyên	25, 100, -5
Float	Số thập phân	29.99, 3.14, -0.5

## Điểm Quan Trọng

- Sử dụng `type()` để kiểm tra kiểu dữ liệu
- Chuyển đổi giữa các kiểu: `int()`, `float()`, `str()`
- Xử lý lỗi khi chuyển đổi với `try-except`
- Chọn đúng kiểu dữ liệu cho từng mục đích

## Bài Tập 1: Xử Lý Thông Tin Sản Phẩm

Tạo một chương trình xử lý thông tin sản phẩm:

- Tên sản phẩm: "Laptop Dell XPS"
- Giá gốc: 1299.99
- Số lượng: 50
- Tính giá sau giảm giá 10%

```
1 # Your code here:  
2 product_name = "Laptop Dell XPS"  
3 original_price = 1299.99  
4 quantity = 50  
5 discount_rate = 0.10  
6  
7 # Calculate discounted price  
8 # Print product information  
9 # Show data types
```

- Khi nào nên sử dụng Integer thay vì Float?
- Làm thế nào để xử lý dữ liệu đầu vào không hợp lệ?
- Tại sao việc hiểu kiểu dữ liệu quan trọng trong Data Science?

**Questions?**