



Python Programming for Data Science

Basic Data Structures: Lists, Tuples, Sets, Dictionaries

Quanskill

Giới thiệu về Data Structures

Lists - Danh sách

Tuples - Bộ dữ liệu

Sets - Tập hợp

Dictionaries - Từ điển

So sánh và Lựa chọn Data Structure

Kết luận

Giới thiệu về Data Structures

- **Lưu trữ dữ liệu:** Tổ chức và quản lý datasets lớn
- **Xử lý dữ liệu:** Thao tác hiệu quả với nhiều loại dữ liệu khác nhau
- **Machine Learning:** Chuẩn bị features và labels cho mô hình
- **Data Analysis:** Phân tích, thống kê và visualization

Trong bài này chúng ta sẽ học

4 cấu trúc dữ liệu cơ bản nhất của Python được sử dụng hàng ngày trong Data Science

Data Structure	Ordered	Mutable	Duplicates
List	o	o	o
Tuple	o	x	o
Set	x	o	x
Dictionary	o	o	x (keys)

- **Ordered:** Thứ tự các phần tử được duy trì
- **Mutable:** Có thể thay đổi sau khi tạo
- **Duplicates:** Cho phép các phần tử trùng lặp

Lists - Danh sách

List là cấu trúc dữ liệu có thứ tự, có thể thay đổi và cho phép duplicate values.

Tạo Lists

```
1 # Create empty list
2 empty_list = []
3 print(empty_list) # Output: []
4
5 # Create list with data
6 student_scores = [85, 92, 78, 96, 88]
7 print(student_scores) # Output: [85, 92, 78, 96, 88]
8
9 # Mixed data types
10 mixed_data = ["Alice", 25, 85.5, True]
11 print(mixed_data) # Output: ['Alice', 25, 85.5, True]
```

Giải thích: Lists sử dụng dấu ngoặc vuông [] và các phần tử cách nhau bởi dấu phẩy.

Index và Slicing

```
1 # Dataset of house prices (in thousands)
2 house_prices = [250, 180, 320, 150, 280, 400, 200]
3
4 # Access by index (starts from 0)
5 print(house_prices[0])    # Output: 250 (first element)
6 print(house_prices[3])    # Output: 150 (fourth element)
7 print(house_prices[-1])   # Output: 200 (last element)
8
9 # Slicing [start:end:step]
10 print(house_prices[1:4])  # Output: [180, 320, 150]
11 print(house_prices[:3])   # Output: [250, 180, 320]
12 print(house_prices[4:])   # Output: [280, 400, 200]
13 print(house_prices[::-2]) # Output: [250, 320, 280, 200]
```

Giải thích: Index âm đếm từ cuối về. Slicing giúp lấy một phần của list.

Thao tác với Lists

```
1 # Customer data
2 customers = ["John", "Mary", "Bob"]
3
4 # Adding elements
5 customers.append("Alice")          # Add at end
6 print(customers) # Output: ['John', 'Mary', 'Bob', 'Alice']
7
8 customers.insert(1, "David")       # Insert at position 1
9 print(customers) # Output: ['John', 'David', 'Mary', 'Bob', 'Alice']
10
11 # Removing elements
12 customers.remove("Bob")          # Remove first occurrence
13 print(customers) # Output: ['John', 'David', 'Mary', 'Alice']
14
15 removed = customers.pop()        # Remove and return last
16 print(removed) # Output: Alice
17 print(customers) # Output: ['John', 'David', 'Mary']
18
19 del customers[1]                 # Delete by index
20 print(customers) # Output: ['John', 'Mary']
```

Phân tích dữ liệu bán hàng

```
1 # Daily sales data for a week
2 daily_sales = [1200, 800, 1500, 900, 1800, 2200, 1600]
3 days = ["Mon", "Tue", "Wed", "Thu", "Fri", "Sat", "Sun"]
4
5 # Basic statistics
6 total_sales = sum(daily_sales)
7 average_sales = total_sales / len(daily_sales)
8 max_sales = max(daily_sales)
9 min_sales = min(daily_sales)
10
11 print(f"Total sales: ${total_sales}")           # Output: Total sales: $10000
12 print(f"Average sales: ${average_sales}")       # Output: Average sales: $1428.57
13 print(f"Max sales: ${max_sales}")              # Output: Max sales: $2200
14 print(f"Min sales: ${min_sales}")              # Output: Min sales: $800
15
16 # Find best selling day
17 best_day_index = daily_sales.index(max_sales)
18 print(f"Best day: {days[best_day_index]}")      # Output: Best day: Sat
```

Tuples - Bộ dữ liệu

Tuple là cấu trúc dữ liệu có thứ tự, **KHÔNG** thể thay đổi và cho phép duplicate values.

Tạo Tuples

```
1 # Create empty tuple
2 empty_tuple = ()
3 print(empty_tuple) # Output: ()
4
5 # Create tuple with data
6 coordinates = (10.5, 20.3)
7 print(coordinates) # Output: (10.5, 20.3)
8
9 # Single element tuple (note the comma!)
10 single_tuple = (42,)
11 print(single_tuple) # Output: (42,)
12
13 # Without parentheses (tuple packing)
14 student_info = "Alice", 20, "Computer Science"
15 print(student_info) # Output: ('Alice', 20, 'Computer Science')
16 print(type(student_info)) # Output: <class 'tuple'>
```

Giải thích: Tuples sử dụng dấu ngoặc tròn () hoặc chỉ cần dấu phẩy.

Truy cập và Unpacking

```
1 # RGB color values
2 red_color = (255, 0, 0)
3 green_color = (0, 255, 0)
4 blue_color = (0, 0, 255)
5
6 # Access by index (same as lists)
7 print(red_color[0])      # Output: 255
8 print(red_color[1:])     # Output: (0, 0)
9
10 # Tuple unpacking (very useful!)
11 r, g, b = red_color
12 print(f"Red: {r}, Green: {g}, Blue: {b}")
13 # Output: Red: 255, Green: 0, Blue: 0
14
15 # Multiple assignment using tuples
16 name, age, grade = "Bob", 22, "A"
17 print(f"Student: {name}, Age: {age}, Grade: {grade}")
18 # Output: Student: Bob, Age: 22, Grade: A
19
20 # Swapping variables
21 x, y = 10, 20
22 x, y = y, x # Swap using tuple unpacking
23 print(f"x: {x}, y: {y}") # Output: x: 20, y: 10
```

Dữ liệu tọa độ GPS và Database records

```
1 # GPS coordinates for store locations
2 store_locations = [
3     ("Store A", 40.7128, -74.0060), # New York
4     ("Store B", 34.0522, -118.2437), # Los Angeles
5     ("Store C", 41.8781, -87.6298)   # Chicago
6 ]
7
8 # Processing GPS data
9 for store_data in store_locations:
10     name, lat, lon = store_data # Unpack tuple
11     print(f"{name}: Latitude {lat}, Longitude {lon}")
12
13 # Output:
14 # Store A: Latitude 40.7128, Longitude -74.006
15 # Store B: Latitude 34.0522, Longitude -118.2437
16 # Store C: Latitude 41.8781, Longitude -87.6298
17
18 # Function returning multiple values as tuple
19 def calculate_stats(numbers):
20     return min(numbers), max(numbers), sum(numbers)/len(numbers)
21
22 data = [10, 20, 30, 40, 50]
23 min_val, max_val, avg_val = calculate_stats(data)
24 print(f"Min: {min_val}, Max: {max_val}, Average: {avg_val}")
25 # Output: Min: 10, Max: 50, Average: 30.0
```

Sets - Tập hợp

Set là collection **không có thứ tự**, có thể thay đổi và **KHÔNG cho phép duplicate**.

Tạo Sets

```
1 # Create empty set
2 empty_set = set() # Note: {} creates empty dict, not set!
3 print(empty_set) # Output: set()
4
5 # Create set with data
6 unique_ages = {25, 30, 35, 25, 30} # Duplicates will be removed
7 print(unique_ages) # Output: {25, 30, 35}
8
9 # Convert list to set (remove duplicates)
10 customer_ids = [101, 102, 101, 103, 102, 104]
11 unique_customers = set(customer_ids)
12 print(unique_customers) # Output: {101, 102, 103, 104}
13
14 # Create set from string
15 letters = set("hello")
16 print(letters) # Output: {'h', 'e', 'l', 'o'}
```

Giải thích: Sets tự động loại bỏ duplicates và không duy trì thứ tự.

Thao tác cơ bản với Sets

```
1 # Website visitors
2 monday_visitors = {"user1", "user2", "user3", "user4"}
3 tuesday_visitors = {"user3", "user4", "user5", "user6"}
4
5 # Add and remove elements
6 monday_visitors.add("user7")
7 print(monday_visitors)    # Output: {'user1', 'user2', 'user3', 'user4', 'user7'}
8
9 monday_visitors.remove("user1") # Raises error if not found
10 print(monday_visitors)   # Output: {'user2', 'user3', 'user4', 'user7'}
11
12 monday_visitors.discard("user999") # No error if not found
13 print(monday_visitors)   # Output: {'user2', 'user3', 'user4', 'user7'}
14
15 # Check membership
16 print("user2" in monday_visitors)    # Output: True
17 print("user5" in monday_visitors)    # Output: False
18
19 # Set length
20 print(len(monday_visitors))    # Output: 4
```

Phép toán tập hợp

```
1 # Customer segments
2 premium_customers = {"Alice", "Bob", "Charlie", "David"}
3 active_customers = {"Bob", "Diana", "Eve", "Charlie"}
4
5 # Union - all customers
6 all_customers = premium_customers | active_customers
7 print(all_customers)
8 # Output: {'Alice', 'Bob', 'Charlie', 'David', 'Diana', 'Eve'}
9
10 # Intersection - customers in both segments
11 both_segments = premium_customers & active_customers
12 print(both_segments) # Output: {'Bob', 'Charlie'}
13
14 # Difference - premium customers who are not active
15 premium_only = premium_customers - active_customers
16 print(premium_only) # Output: {'Alice', 'David'}
17
18 # Symmetric difference - customers in only one segment
19 exclusive_customers = premium_customers ^ active_customers
20 print(exclusive_customers) # Output: {'Alice', 'David', 'Diana', 'Eve'}
```

Phân tích dữ liệu khách hàng

```
1 # E-commerce customer behavior analysis
2 email_subscribers = {"cust1", "cust2", "cust3", "cust4", "cust5"}
3 purchase_customers = {"cust2", "cust4", "cust6", "cust7", "cust8"}
4 mobile_app_users = {"cust1", "cust3", "cust6", "cust9", "cust10"}
5
6 # Find customers who subscribed but never purchased
7 subscribed_no_purchase = email_subscribers - purchase_customers
8 print(f"Subscribed but no purchase: {subscribed_no_purchase}")
9 # Output: Subscribed but no purchase: {'cust1', 'cust3', 'cust5'}
10
11 # Find highly engaged customers (all three activities)
12 highly_engaged = email_subscribers & purchase_customers & mobile_app_users
13 print(f"Highly engaged customers: {highly_engaged}")
14 # Output: Highly engaged customers: set()
15
16 # Total unique customers across all platforms
17 total_customers = email_subscribers | purchase_customers | mobile_app_users
18 print(f"Total unique customers: {len(total_customers)}")
19 # Output: Total unique customers: 10
```

Dictionaries - Từ điển

Dictionary lưu trữ dữ liệu dưới dạng key-value pairs, có thứ tự (Python 3.7+) và có thể thay đổi.

Tạo Dictionaries

```
1 # Create empty dictionary
2 empty_dict = {}
3 print(empty_dict) # Output: {}
4
5 # Create dictionary with data
6 student_grades = {
7     "Alice": 85,
8     "Bob": 92,
9     "Charlie": 78,
10    "Diana": 96
11 }
12 print(student_grades)
13 # Output: {'Alice': 85, 'Bob': 92, 'Charlie': 78, 'Diana': 96}
14
15 # Using dict() constructor
16 product_prices = dict(laptop=999, mouse=25, keyboard=75)
17 print(product_prices) # Output: {'laptop': 999, 'mouse': 25, 'keyboard': 75}
```

Giải thích: Keys phải unique và immutable (string, number, tuple). Values có thể là bất kỳ data type nào.

Truy cập và thay đổi dữ liệu

```
1 # Sales data by region
2 sales_by_region = {
3     "North": 15000,
4     "South": 12000,
5     "East": 18000,
6     "West": 14000
7 }
8
9 # Access values by key
10 print(sales_by_region["North"]) # Output: 15000
11 print(sales_by_region.get("North")) # Output: 15000
12 print(sales_by_region.get("Central", 0)) # Output: 0 (default value)
13
14 # Modify existing value
15 sales_by_region["North"] = 16000
16 print(sales_by_region["North"]) # Output: 16000
17
18 # Add new key-value pair
19 sales_by_region["Central"] = 10000
20 print(sales_by_region)
21 # Output: {'North': 16000, 'South': 12000, 'East': 18000, 'West': 14000, 'Central': 10000}
22
23 # Remove key-value pair
24 del sales_by_region["Central"]
25 removed_value = sales_by_region.pop("West", 0)
26 print(f"Removed West sales: {removed_value}") # Output: Removed West sales: 14000
```

Các method quan trọng

```
1 # Employee information
2 employees = {
3     "emp001": {"name": "Alice", "dept": "IT", "salary": 75000},
4     "emp002": {"name": "Bob", "dept": "Sales", "salary": 65000},
5     "emp003": {"name": "Charlie", "dept": "IT", "salary": 80000}
6 }
7
8 # Get all keys, values, items
9 print(employees.keys())
10 # Output: dict_keys(['emp001', 'emp002', 'emp003'])
11
12 print(list(employees.values())[0])
13 # Output: {'name': 'Alice', 'dept': 'IT', 'salary': 75000}
14
15 # Iterate through dictionary
16 for emp_id, info in employees.items():
17     print(f"ID: {emp_id}, Name: {info['name']}, Salary: ${info['salary']}")
18
19 # Output:
20 # ID: emp001, Name: Alice, Salary: $75000
21 # ID: emp002, Name: Bob, Salary: $65000
22 # ID: emp003, Name: Charlie, Salary: $80000
```

Phân tích dữ liệu bán hàng

```
1 # Product sales data
2 sales_data = {
3     "products": ["laptop", "mouse", "keyboard", "monitor"],
4     "quantities": [50, 200, 150, 75],
5     "prices": [999, 25, 75, 299],
6     "categories": ["electronics", "accessories", "accessories", "electronics"]
7 }
8
9 # Calculate total revenue by product
10 for i in range(len(sales_data["products"])):
11     product = sales_data["products"][i]
12     revenue = sales_data["quantities"][i] * sales_data["prices"][i]
13     print(f"{product}: ${revenue}")
14
15 # Output:
16 # laptop: $49950
17 # mouse: $5000
18 # keyboard: $11250
19 # monitor: $22425
20
21 # Count products by category
22 category_count = {}
23 for category in sales_data["categories"]:
24     category_count[category] = category_count.get(category, 0) + 1
25
26 print(category_count) # Output: {'electronics': 2, 'accessories': 2}
```

So sánh và Lựa chọn Data Structure

List - Sử dụng khi

- Cần thứ tự và cho phép duplicates
- Thường xuyên thêm/xóa elements
- Time series data, sequences

Tuple - Sử dụng khi

- Dữ liệu không thay đổi (coordinates, RGB values)
- Return multiple values từ functions
- Dictionary keys

Set - Sử dụng khi

- Loại bỏ duplicates
- Membership testing (in operator)
- Set operations (union, intersection)

Dictionary - Sử dụng khi

- Key-value relationships
- Fast lookups by key
- Grouping và counting data

Kết hợp tất cả Data Structures

```
1 # Customer transaction data
2 customers = {
3     "customer001": {
4         "name": "Alice Johnson",
5         "transactions": [150, 200, 75, 300], # List of amounts
6         "location": (40.7128, -74.0060), # Tuple for coordinates
7         "purchase_categories": {"electronics", "clothing", "books"} # Set
8     },
9     "customer002": {
10        "name": "Bob Smith",
11        "transactions": [100, 250, 180],
12        "location": (34.0522, -118.2437),
13        "purchase_categories": {"electronics", "sports", "books"}
14    }
15 }
16
17 # Analysis
18 for customer_id, data in customers.items():
19     name = data["name"]
20     total_spent = sum(data["transactions"]) # List operation
21     avg_transaction = total_spent / len(data["transactions"])
22     lat, lon = data["location"] # Tuple unpacking
23     num_categories = len(data["purchase_categories"]) # Set size
24
25     print(f"Customer: {name}")
26     print(f"Total spent: ${total_spent}, Average: ${avg_transaction:.2f}")
27     print(f"Location: ({lat}, {lon}), Categories: {num_categories}")
28     print("---")
```

Kết luận

Những gì chúng ta đã học

- **Lists:** Ordered, mutable, allows duplicates - dùng cho sequences
- **Tuples:** Ordered, immutable, allows duplicates - dùng cho fixed data
- **Sets:** Unordered, mutable, no duplicates - dùng cho unique values
- **Dictionaries:** Ordered, mutable, unique keys - dùng cho key-value pairs

Tips cho Data Science

- Sử dụng Lists cho time series và numerical data
- Tuples cho coordinates, fixed records
- Sets cho unique values và set operations
- Dictionaries cho data grouping và fast lookups

Questions?