

# NumPy - Tổng hợp các phương thức thông dụng và nâng cao cho Khoa học Dữ liệu và Trí tuệ Nhân tạo

duckAI soạn cho Hùng

11:13 AM +07, Thứ Năm, 24/07/2025

## Mục lục

<b>1</b>	<b>Giới thiệu</b>	<b>3</b>
1.1	Mục tiêu học tập . . . . .	3
1.2	Yêu cầu kiến thức . . . . .	3
<b>2</b>	<b>Các hàm cơ bản để khởi tạo mảng và sinh dữ liệu</b>	<b>3</b>
2.1	np.array . . . . .	3
2.2	np.zeros . . . . .	4
2.3	np.ones . . . . .	4
2.4	np.eye . . . . .	4
2.5	np.arange . . . . .	5
2.6	np.linspace . . . . .	5
2.7	np.random . . . . .	6
<b>3</b>	<b>Thuộc tính của mảng</b>	<b>6</b>
3.1	shape . . . . .	6
3.2	ndim . . . . .	7
3.3	size . . . . .	7
3.4	dtype . . . . .	7
<b>4</b>	<b>Đánh chỉ số và cắt lát</b>	<b>8</b>
4.1	Truy cập phần tử cụ thể . . . . .	8
4.2	Truy xuất hàng/cột . . . . .	8
4.3	Lọc điều kiện . . . . .	8
4.4	np.where . . . . .	9
4.5	np.isnan, np.isinf . . . . .	9
4.6	np.clip . . . . .	10
<b>5</b>	<b>Các phép toán và thống kê</b>	<b>10</b>
5.1	np.sum . . . . .	10
5.2	np.mean, np.std . . . . .	11
5.3	np.min / np.max . . . . .	11
5.4	np.median, np.percentile . . . . .	12
5.5	np.corrcoef . . . . .	12
5.6	np.cov . . . . .	13

5.7	np.histogram . . . . .	13
<b>6</b>	<b>Thao tác hình dạng và Reshape</b>	<b>14</b>
6.1	reshape . . . . .	14
6.2	flatten . . . . .	14
6.3	transpose (T) . . . . .	14
6.4	np.concatenate . . . . .	15
6.5	np.vstack, np.hstack . . . . .	15
6.6	np.unique . . . . .	16
6.7	np.vectorize . . . . .	16
<b>7</b>	<b>Sắp xếp và Tìm kiếm</b>	<b>17</b>
7.1	np.sort . . . . .	17
7.2	np.argsort . . . . .	17
7.3	np.argmax / np.argmin . . . . .	17
<b>8</b>	<b>Broadcasting và Đại số tuyến tính</b>	<b>18</b>
8.1	Broadcasting . . . . .	18
8.2	np.dot . . . . .	18
8.3	np.linalg . . . . .	19
8.4	np.einsum . . . . .	19
<b>9</b>	<b>Xử lý dữ liệu lớn</b>	<b>20</b>
9.1	np.save, np.load . . . . .	20
<b>10</b>	<b>Ví dụ thực hành từ Kaggle</b>	<b>20</b>
10.1	Bài 1: Chuẩn hóa dữ liệu nhiệt độ . . . . .	21
10.2	Bài 2: Tìm ngày doanh thu cao nhất . . . . .	21
10.3	Bài 3: Tính toán ma trận . . . . .	22
10.4	Bài 4: Phân tích dữ liệu phân loại . . . . .	22
10.5	Bài 5: Xử lý dữ liệu thiếu . . . . .	22
10.6	Bài 6: Phân tích phân phối dữ liệu . . . . .	23
10.7	Bài 7: Tính toán tensor . . . . .	23

# 1 Giới thiệu

Tài liệu này được thiết kế như một tài liệu nghiên cứu chuyên sâu về thư viện NumPy, một công cụ cốt lõi trong việc xử lý dữ liệu số và phát triển các ứng dụng trong Khoa học Dữ liệu và Trí tuệ Nhân tạo (AI). NumPy cung cấp các cấu trúc dữ liệu hiệu quả như mảng đa chiều (ndarray) và các hàm toán học tối ưu hóa, giúp tăng tốc độ tính toán và hỗ trợ các thuật toán phức tạp. Tài liệu nhằm mục đích phục vụ các học viên, từ người mới bắt đầu đến các nhà nghiên cứu, với các ví dụ thực tiễn và hướng dẫn chi tiết, bao gồm cả các bài tập lấy cảm hứng từ nền tảng Kaggle. Mỗi phần đều được xây dựng với mục tiêu cung cấp nền tảng lý thuyết kết hợp thực hành, đảm bảo học viên có thể áp dụng kiến thức vào các dự án thực tế.

## 1.1 Mục tiêu học tập

- Hiểu và sử dụng các phương thức cơ bản để khởi tạo và thao tác mảng NumPy. - Thành thạo các kỹ thuật thống kê, đại số tuyến tính và xử lý dữ liệu nâng cao. - Áp dụng NumPy trong các kịch bản thực tế như phân tích dữ liệu, học máy và mô phỏng.

## 1.2 Yêu cầu kiến thức

- Đã học ngôn ngữ lập trình Python cơ bản.

# 2 Các hàm cơ bản để khởi tạo mảng và sinh dữ liệu

## 2.1 np.array

**Chức năng:** Phương thức này chuyển đổi danh sách, tuple hoặc các cấu trúc iterable thành mảng NumPy (ndarray), hỗ trợ tính toán vector hóa với hiệu suất cao, rất quan trọng trong các ứng dụng khoa học dữ liệu.

**Cú pháp:** `np.array(object, dtype=None, copy=True)`

**Tham số:**

- **object:** Dữ liệu đầu vào (danh sách, tuple, iterable).
- **dtype:** Kiểu dữ liệu (ví dụ: `int32`, `float64`, `bool`), được suy ra tự động nếu không chỉ định.
- **copy:** Giá trị Boolean, xác định liệu có sao chép dữ liệu (`True`) hay sử dụng tham chiếu (`False`).

**Ví dụ:** Tạo mảng 1D và 2D với kiểu `float64`.

---

```
1 >>> import numpy as np
2 # khởi tạo mảng 1D NumPy với các phần tử [1.0, 2.0, 3.0] sử dụng
   kiểu float64
3 # hỗ trợ tính toán vector hóa hiệu quả trong phân tích số học
4 >>> arr1 = np.array([1, 2, 3], dtype=float)
5 # tạo mảng 2D NumPy [[1.0, 2.0], [3.0, 4.0]] với kiểu float64
6 # phù hợp cho các tính toán ma trận trong khoa học dữ liệu
7 >>> arr2 = np.array([[1, 2], [3, 4]], dtype=float)
8 >>> print(arr1, "\n", arr2)
```

---

## 2.2 np.zeros

**Chức năng:** Tạo mảng chứa toàn số 0, thường được sử dụng để khởi tạo mảng trước khi thực hiện các phép tính phức tạp, đặc biệt trong các thuật toán học máy như khởi tạo ma trận trọng số.

**Cú pháp:** `np.zeros(shape, dtype=float, order='C')`

**Tham số:**

- **shape:** Tuple chỉ kích thước (ví dụ: (2,3) cho ma trận 2x3).
- **dtype:** Kiểu dữ liệu (mặc định là `float64`).
- **order:** Quy định cách lưu trữ dữ liệu theo hàng ('C') hoặc cột ('F').

**Ví dụ:** Tạo ma trận 2x3 chứa số 0.

---

```
1 >>> zerosArr = np.zeros((2, 3), dtype=int)
2 # tao ma tran 2x3 day so 0 su dung kieu integer
3 # thuong dung de khoi tao mang trong thuat toan nhu ma tran trong
  so mang nho dan
4 >>> print(zerosArr)
5 # hien thi ma tran ket qua: [[0 0 0], [0 0 0]]
```

---

## 2.3 np.ones

**Chức năng:** Tạo mảng chứa toàn số 1, hữu ích trong việc khởi tạo mảng hoặc thực hiện các phép nhân ma trận trong các bài toán toán học và học máy.

**Cú pháp:** `np.ones(shape, dtype=float, order='C')`

**Tham số:** Tương tự như `np.zeros`.

**Ví dụ:** Tạo ma trận 3x2 chứa số 1 kiểu `int`.

---

```
1 >>> onesArr = np.ones((3, 2), dtype=int)
2 # tao ma tran 3x2 day so 1 su dung kieu integer
3 # thuong dung lam diem bat dau cho cac phép cong tích luy hoac
  thao tac quy mo
4 >>> print(onesArr)
5 # hien thi ma tran ket qua: [[1 1], [1 1], [1 1]]
```

---

## 2.4 np.eye

**Chức năng:** Tạo ma trận đơn vị với các phần tử trên đường chéo chính bằng 1 và các phần tử khác bằng 0, đóng vai trò quan trọng trong đại số tuyến tính và các phép biến đổi ma trận.

**Cú pháp:** `np.eye(N, M=None, k=0, dtype=float)`

**Tham số:**

- **N:** Số hàng.
- **M:** Số cột (mặc định bằng N).
- **k:** Vị trí đường chéo (0: chính, dương: trên, âm: dưới).

- **dtype:** Kiểu dữ liệu.

**Ví dụ:** Tạo ma trận đơn vị 3x3.

---

```

1 >>> eyeArr = np.eye(3, dtype=int)
2 # tao ma tran don vi 3x3 su dung kieu integer
3 # mot ma tran co ban trong dai so tuyen tinh dung cho cac bien
  doi danh thuc
4 >>> print(eyeArr)
5 # hien thi ma tran ket qua: [[1 0 0], [0 1 0], [0 0 1]]

```

---

## 2.5 np.arange

**Chức năng:** Tạo mảng 1D chứa dãy số cách đều trong khoảng [start, stop), thường được sử dụng để tạo các chuỗi số cho mô phỏng hoặc phân tích dữ liệu.

**Cú pháp:** np.arange(start, stop, step, dtype=None)

**Tham số:**

- **start:** Giá trị bắt đầu (bao gồm, mặc định 0).
- **stop:** Giá trị kết thúc (không bao gồm).
- **step:** Khoảng cách giữa các phần tử (mặc định 1).
- **dtype:** Kiểu dữ liệu.

**Ví dụ:** Tạo dãy số chẵn từ 0 đến 10.

---

```

1 >>> rangeArr = np.arange(0, 11, 2)
2 # tao mang 1D voi cac so chan [0, 2, 4, 6, 8, 10] voi buoc nhay 2
3 # huu ich cho viec tao day so trong mo phong hoac phan tich thoi
  gian
4 >>> print(rangeArr)
5 # hien thi mang ket qua: [0 2 4 6 8 10]

```

---

## 2.6 np.linspace

**Chức năng:** Tạo mảng 1D chứa một số lượng cố định (num) các giá trị chia đều trong khoảng [start, stop], phù hợp cho việc tạo các điểm lấy mẫu đều trong các mô hình toán học.

**Cú pháp:** np.linspace(start, stop, num=50, endpoint=True, dtype=None)

**Tham số:**

- **start:** Giá trị đầu.
- **stop:** Giá trị cuối (bao gồm nếu endpoint=True).
- **num:** Số phần tử.
- **endpoint:** Giá trị Boolean, xác định có bao gồm stop không.
- **dtype:** Kiểu dữ liệu.

**Ví dụ:** Tạo 5 số chia đều từ 0 đến 1.

---

```
1 >>> linArr = np.linspace(0, 1, 5)
2 # tao 5 so chia deu [0.0, 0.25, 0.5, 0.75, 1.0] giua 0 va 1
3 # ly tuong de tao mau deu trong mo hinh hoa toan hoc hoac ve bieu
  do
4 >>> print(linArr)
5 # hien thi mang ket qua: [0.    0.25 0.5   0.75 1.   ]
```

---

## 2.7 np.random

**Chức năng:** Tạo dữ liệu ngẫu nhiên, một công cụ quan trọng trong mô phỏng, học máy, và khởi tạo trọng số cho các mô hình AI.

**Các phương thức con:**

- `np.random.rand(shape)`: Số ngẫu nhiên trong  $[0,1)$ .
- `np.random.randn(shape)`: Số ngẫu nhiên theo phân phối chuẩn ( $\text{mean}=0$ ,  $\text{std}=1$ ).
- `np.random.randint(low, high, size)`: Số nguyên ngẫu nhiên trong  $[\text{low}, \text{high})$ .

**Tham số:**

- **shape/size**: Kích thước mảng.
- **low, high**: Khoảng giá trị cho `randint`.

**Ứng dụng:** Mô phỏng dữ liệu, chia tập train/test trong học máy.

**Ví dụ:** Sinh dữ liệu ngẫu nhiên.

---

```
1 >>> randArr = np.random.rand(2, 3)
2 # tao mang 2x3 voi cac gia tri ngau nhien phan phoi deu giua 0 va
  1
3 # thuong dung de khoi tao trong so trong mang nho dan hoac mo
  phong Monte Carlo
4 >>> normArr = np.random.randn(5)
5 # tao 5 so ngau nhien tu phan phoi chuan (mean=0, std=1)
6 # huu ich de tao du lieu gia mo phong tieng noi thuc te
7 >>> intArr = np.random.randint(1, 10, 4)
8 # tao 4 so nguyen ngau nhien giua 1 va 9 (khong ke 10)
9 # ap dung trong thuat toan ngau nhien hoac phat trien game
10 >>> print(randArr, "\n", normArr, "\n", intArr)
11 # hien thi ket qua ngau nhien cua ba mang
```

---

## 3 Thuộc tính của mảng

### 3.1 shape

**Chức năng:** Trả về một tuple biểu diễn kích thước của mảng (số hàng, số cột, v.v.), giúp kiểm tra cấu trúc trước khi thực hiện các phép biến đổi hoặc tính toán.

**Cú pháp:** `a.shape`

**Ứng dụng:** Kiểm tra kích thước mảng trong các tác vụ tiền xử lý dữ liệu.

**Ví dụ:** Kiểm tra kích thước ma trận 3x4.

---

```
1 >>> a = np.zeros((3, 4))
2 # khoi tao ma tran 3x4 day so 0
3 # diem bat dau thuong dung cho cac phep tinh ma tran trong dai so
  tuyen tinh
4 >>> print(a.shape)
5 # tra ve kích thước mảng dưới dạng tuple (3, 4), chỉ 3 hàng và 4
  cột
```

---

## 3.2 ndim

**Chức năng:** Trả về số chiều của mảng (1D, 2D, 3D, v.v.), hỗ trợ phân biệt loại mảng để áp dụng các phép xử lý phù hợp.

**Cú pháp:** `a.ndim`

**Ứng dụng:** Xác định số chiều trong các tác vụ đa chiều như xử lý hình ảnh.

**Ví dụ:** Kiểm tra số chiều của ma trận 2D.

---

```
1 >>> a = np.array([[1, 2], [3, 4]])
2 # định nghĩa mảng 2D với 2 hàng và 2 cột
3 # đại diện cho cấu trúc ma trận trong dữ liệu đa chiều
4 >>> print(a.ndim)
5 # tra về số chiều: 2, chỉ rằng đây là cấu trúc 2D
```

---

## 3.3 size

**Chức năng:** Trả về tổng số phần tử trong mảng, giúp đánh giá dung lượng dữ liệu trước khi xử lý.

**Cú pháp:** `a.size`

**Ứng dụng:** Kiểm tra số lượng phần tử trong các bài toán phân tích dữ liệu lớn.

**Ví dụ:** Đếm tổng số phần tử.

---

```
1 >>> a = np.ones((2, 3))
2 # tao ma tran 2x3 day so 1
3 # thuong dung lam co so de tích lũy giá trị trong xử lý dữ liệu
4 >>> print(a.size)
5 # tra về tổng số phần tử: 6 (2 hàng * 3 cột)
```

---

## 3.4 dtype

**Chức năng:** Trả về kiểu dữ liệu của các phần tử trong mảng, hỗ trợ kiểm tra hoặc ép kiểu dữ liệu trong quá trình xử lý.

**Cú pháp:** `a.dtype`

**Ứng dụng:** Đảm bảo tính nhất quán kiểu dữ liệu trong các phép tính số học.

**Ví dụ:** Kiểm tra kiểu dữ liệu của mảng số thực.

---

```
1 >>> arr = np.array([1.2, 3.4])
2 # tao mảng với các giá trị thực
```

---

```
3 # thường dùng cho dữ liệu liên tục trong tính toán khoa học
4 >>> print(arr.dtype)
5 # trả về kiểu dữ liệu: float64, chỉ đoạn do chính xác 64-bit
```

---

## 4 Đánh chỉ số và cắt lát

### 4.1 Truy cập phần tử cụ thể

**Chức năng:** Lấy hoặc sửa đổi giá trị của phần tử tại vị trí xác định, một thao tác cơ bản trong việc thao tác dữ liệu.

**Cú pháp:** `a[i, j]` (cho mảng 2D).

**Ứng dụng:** Truy xuất hoặc cập nhật dữ liệu trong các tập dữ liệu đa chiều.

**Ví dụ:** Lấy phần tử tại hàng 0, cột 1.

---

```
1 >>> a = np.array([[1, 2, 3], [4, 5, 6]])
2 # định nghĩa ma trận 2x3 với các giá trị integer
3 # đại diện cho tập dữ liệu mà mọi phần tử có thể được truy cập
  riêng lẻ
4 >>> print(a[0, 1])
5 # lấy phần tử tại hàng 0, cột 1, trả về: 2
```

---

### 4.2 Truy xuất hàng/cột

**Chức năng:** Trích xuất hàng, cột hoặc vùng con của mảng bằng kỹ thuật cắt lát, hỗ trợ phân tích dữ liệu chi tiết.

**Cú pháp:**

- `a[:, j]`: Cột thứ `j`.
- `a[i, :]`: Hàng thứ `i`.
- `a[i:j, k:l]`: Vùng con từ hàng `i` đến `j-1`, cột `k` đến `l-1`.

**Ứng dụng:** Tách dữ liệu để phân tích hoặc huấn luyện mô hình.

**Ví dụ:** Lấy cột đầu và hai hàng đầu.

---

```
1 >>> a = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
2 # tạo ma trận 3x3 đại diện cho tập dữ liệu nhỏ
3 # hữu ích để trích xuất tập con trong phân tích dữ liệu
4 >>> print(a[:, 0])
5 # trích xuất cột đầu: [1 4 7], chọn tất cả hàng cho cột 0
6 >>> print(a[0:2, :])
7 # trích xuất hai hàng đầu: [[1 2 3], [4 5 6]], sử dụng cắt lát
```

---

### 4.3 Lọc điều kiện

**Chức năng:** Lấy các phần tử thỏa mãn điều kiện logic, một kỹ thuật quan trọng trong việc lọc dữ liệu theo ngưỡng.



**Cú pháp:** `a[condition]`

**Ứng dụng:** Lọc dữ liệu trong các tác vụ phân tích thời tiết hoặc tài chính.

**Ví dụ:** Lọc nhiệt độ lớn hơn 30.

---

```
1 >>> temps = np.random.randint(20, 40, 10)
2 # tao mang 10 so nguyen ngau nhien giua 20 va 39
3 # mo phong du lieu doc nhiet do tu cam bien
4 >>> hotTemps = temps[temps > 30]
5 # loc mang chi giu nhung nhiet do tren 30 do
6 # phép tính thông dụng trong phân tích dữ liệu khi hau
7 >>> print(hotTemps)
8 # tra ve mang cac gia tri vuot qua 30, vi du: [32 35 38]
```

---

## 4.4 np.where

**Chức năng:** Trả về chỉ số của phần tử thỏa mãn điều kiện hoặc thay thế giá trị dựa trên điều kiện, hỗ trợ xử lý dữ liệu thiếu và tối ưu hóa.

**Cú pháp:** `np.where(condition[, x, y])`

**Tham số:**

- **condition:** Biểu thức logic.
- **x, y:** Giá trị thay thế khi điều kiện đúng (x) hoặc sai (y). Nếu không có x, y, trả về chỉ số.

**Ứng dụng:** Xử lý dữ liệu thiếu, thay thế giá trị, hoặc tìm vị trí trong phân tích dữ liệu lớn.

**Ví dụ:** Thay thế giá trị  $< 0$  bằng 0.

---

```
1 >>> data = np.array([-1, 2, -3, 4])
2 # tao mang voi cac gia tri am va duong
3 # dai dien cho tap du lieu can duoc chuan hoa
4 >>> result = np.where(data < 0, 0, data)
5 # thay the tat ca gia tri am bang 0, giu nguyen gia tri duong
6 # ky thuat dung de bao dam du lieu khong am trong tai chinh
7 >>> print(result)
8 # tra ve mang da sua: [0 2 0 4]
9 >>> indices = np.where(data < 0)
10 # tra ve chi so noi dieu kien (data < 0) dung
11 # huu ich de xac dinh vi tri bat thuong trong tap du lieu
12 >>> print(indices)
13 # tra ve: [0 2] (chi so cua -1 va -3)
```

---

## 4.5 np.isnan, np.isinf

**Chức năng:** Kiểm tra giá trị thiếu (`np.isnan`) hoặc vô cực (`np.isinf`), rất cần thiết trong việc làm sạch dữ liệu thực tế.

**Cú pháp:** `np.isnan(a)`, `np.isinf(a)`

**Ứng dụng:** Xử lý dữ liệu thiếu hoặc lỗi trong các tập dữ liệu lớn từ các nguồn không đồng nhất.

**Ví dụ:** Kiểm tra và loại bỏ giá trị NaN.

---

```

1 >>> data = np.array([1, np.nan, 3, np.inf, 5])
2 # tao mang voi gia tri NaN (thieu) va vo cuc
3 # mo phong du lieu thuc te co loi co the xay ra
4 >>> nanMask = np.isnan(data)
5 # tao bo loc boolean danh dau gia tri NaN
6 # can thiet de lam sach du lieu trong phan tich thong ke
7 >>> infMask = np.isinf(data)
8 # tao bo loc boolean danh dau gia tri vo cuc
9 # huu ich de phat hien tran so trong tinh toan so hoc
10 >>> print("NaN:", nanMask)
11 # tra ve: [False  True False False False]
12 >>> print("Inf:", infMask)
13 # tra ve: [False False False  True False]
14 >>> cleanData = data[~np.isnan(data) & ~np.isinf(data)]
15 # loc bo ca gia tri NaN va vo cuc
16 # buoc tien xu ly chuan trong ong day hoc may
17 >>> print(cleanData)
18 # tra ve: [1. 3. 5.]

```

---

## 4.6 np.clip

**Chức năng:** Giới hạn giá trị trong mảng trong khoảng `[a_min, a_max]`, giúp loại bỏ các giá trị ngoại lai trong dữ liệu.

**Cú pháp:** `np.clip(a, a_min, a_max)`

**Tham số:**

- **a:** Mảng đầu vào.
- **$a_{min}, a_{max}$ :** *Giới hạn dưới và trên*. **Ứng dụng:** Chuẩn hóa dữ liệu trong xử lý tín hiệu hoặc hình ảnh.

**Ví dụ:** Giới hạn giá trị trong `[0, 10]`.

---

```

1 >>> data = np.array([15, 5, -3, 12])
2 # tao mang voi cac gia tri ngoai khoang mong muon [0, 10]
3 # dai dien cho du lieu cam bien chua duoc chuan hoa
4 >>> clippedData = np.clip(data, 0, 10)
5 # gioi han tat ca gia tri trong khoang [0, 10], cat bo cuc tri
6 # ky thuat dung de tranh ngoai lai trong xu ly anh hoac tin hieu
7 >>> print(clippedData)
8 # tra ve: [10  5  0 10]

```

---

## 5 Các phép toán và thống kê

### 5.1 np.sum

**Chức năng:** Tính tổng các phần tử theo trục hoặc toàn mảng, một phép toán cơ bản trong phân tích dữ liệu số.

**Cú pháp:** `np.sum(a, axis=None, dtype=None)`

**Tham số:**

- **a:** Mảng đầu vào.
- **axis:** Trục tính tổng (None: toàn bộ, 0: cột, 1: hàng).
- **dtype:** Kiểu dữ liệu kết quả.

**Ứng dụng:** Tính tổng doanh thu hoặc điểm số trong các tập dữ liệu kinh doanh.

**Ví dụ:** Tính tổng doanh thu theo ngày và tổng cộng.

---

```

1 >>> salesData = np.array([[100, 200], [300, 400]])
2 # định nghĩa ma trận 2x2 đại diện cho dữ liệu bán hàng hai ngày
  va hai sản phẩm
3 >>> print(np.sum(salesData, axis=0))
4 # tính tổng theo cột (tổng dọc), trả về [400 600]
5 # hữu ích để tổng hợp tổng ngày trên các sản phẩm
6 >>> print(np.sum(salesData))
7 # tính tổng tất cả phần tử trong ma trận, trả về 1000
8 # phép tính thông dụng để tính tổng doanh thu tổng thể

```

---

## 5.2 np.mean, np.std

**Chức năng:** Tính trung bình (**np.mean**) và độ lệch chuẩn (**np.std**), các chỉ số thống kê quan trọng trong phân tích dữ liệu.

**Cú pháp:** **np.mean**(a, axis=None), **np.std**(a, axis=None)

**Ứng dụng:** Chuẩn hóa z-score và phân tích độ biến thiên trong học máy.

**Ví dụ:** Chuẩn hóa z-score dữ liệu ngẫu nhiên.

---

```

1 >>> xData = np.random.randn(1000)
2 # tạo 1000 số ngẫu nhiên từ phân phối chuẩn
3 # mô phỏng tập dữ liệu lớn cho phân tích thống kê
4 >>> meanVal = np.mean(xData)
5 # tính giá trị trung bình của tập dữ liệu
6 # đo lường khuynh hướng trung tâm trong thống kê mô tả
7 >>> stdVal = np.std(xData)
8 # tính độ lệch chuẩn, chỉ rõ độ phân tán của dữ liệu
9 # liên quan đến chuẩn hóa trong học máy
10 >>> zScores = (xData - meanVal) / stdVal
11 # chuẩn hóa dữ liệu thành z-score với trung bình ~0 và độ lệch
   chuẩn ~1
12 # bước tiền xử lý cho nhiều mô hình thống kê
13 >>> print(zScores.mean(), zScores.std())
14 # kiểm tra chuẩn hóa, dự kiến giá trị gần 0 và 1

```

---

## 5.3 np.min / np.max

**Chức năng:** Tìm giá trị nhỏ nhất (**np.min**) hoặc lớn nhất (**np.max**) trong mảng, hỗ trợ xác định phạm vi dữ liệu.

**Cú pháp:** **np.min**(a, axis=None), **np.max**(a, axis=None)

**Ứng dụng:** Xác định cực trị như nhiệt độ cao nhất hoặc giá cổ phiếu thấp nhất.

**Ví dụ:** Tìm nhiệt độ thấp nhất và cao nhất.

---

```

1 >>> temps = np.array([25, 30, 28, 32, 27])
2 # tao mang du lieu doc nhiet do bang Celsius
3 # dai dien cho tap du lieu thoi gian nho
4 >>> print(np.min(temps), np.max(temps))
5 # tra ve nhiet do thap nhat va cao nhat: 25 32
6 # huu ich de xac dinh khoang trong giam sat moi truong

```

---

## 5.4 np.median, np.percentile

**Chức năng:** Tính trung vị (`np.median`) và phần trăm (`np.percentile`), cung cấp các chỉ số phân vị quan trọng trong thống kê.

**Cú pháp:** `np.median(a, axis=None)`, `np.percentile(a, q, axis=None)`

**Tham số:**

- **q:** Phần trăm (ví dụ: 90 cho percentile 90).

**Ứng dụng:** Phân tích phân phối và xác định ngưỡng bất thường trong dữ liệu.

**Ví dụ:** Tính trung vị và percentile 90 của doanh thu.

---

```

1 >>> revData = np.array([120, 100, 150, 130, 170])
2 # dinh nghia mang du lieu doanh thu cua mot doanh nghiep nho
3 >>> medianVal = np.median(revData)
4 # tinh gia tri trung vi (130), do luong trung tam manh me
5 >>> p90Val = np.percentile(revData, 90)
6 # tinh percentile 90 (164), chi nguong 10% doanh thu cao nhat
7 >>> print(medianVal, p90Val)
8 # tra ve: 130 164, huu ich de phat hien ngoai lai trong phan tich
   tai chinh

```

---

## 5.5 np.corrcoef

**Chức năng:** Tính hệ số tương quan Pearson giữa các biến, hỗ trợ phân tích mối quan hệ tuyến tính trong dữ liệu.

**Cú pháp:** `np.corrcoef(x, y=None)`

**Tham số:**

- **x:** Mảng 1D hoặc 2D (hàng là biến, cột là quan sát).
- **y:** Mảng 1D bổ sung (tùy chọn).

**Ứng dụng:** Đánh giá mối quan hệ giữa các đặc trưng trong học máy.

**Ví dụ:** Tính tương quan giữa doanh thu và chi phí quảng cáo.

---

```

1 >>> salesData = np.array([100, 120, 150, 130, 170])
2 # dai dien cho du lieu doanh thu nam trong nam ky
3 >>> adsData = np.array([50, 60, 80, 70, 90])
4 # dai dien cho chi phi quang cao trong cung ky
5 >>> corrMatrix = np.corrcoef(salesData, adsData)
6 # tinh ma tran he so tuong quan Pearson
7 # danh gia moi quan he tuyen tinh giua doanh thu va quang cao

```

---

```
8 >>> print(corrMatrix)
9 # tra ve ma tran 2x2, voi duong cheo 1 va gia tri tuong quan
   ngoai duong cheo
```

---

## 5.6 np.cov

**Chức năng:** Tính ma trận hiệp phương sai giữa các biến, hỗ trợ phân tích độ biến thiên trong các tập dữ liệu đa biến.

**Cú pháp:** `np.cov(m, bias=False)`

**Tham số:**

- **m:** Mảng 2D (hàng là biến, cột là quan sát).
- **bias:** Giá trị Boolean, chuẩn hóa theo N (True) hoặc N-1 (False).

**Ứng dụng:** Phân tích độ biến thiên giữa các đặc trưng trong học máy.

**Ví dụ:** Tính hiệp phương sai giữa doanh thu và chi phí.

---

```
1 >>> salesData = np.array([100, 120, 150, 130, 170])
2 # chua du lieu doanh thu cho nam quan sat
3 >>> adsData = np.array([50, 60, 80, 70, 90])
4 # chua du lieu chi phi quang cao cho cung ky
5 >>> covMatrix = np.cov([salesData, adsData])
6 # tinh ma tran hiệp phương sai, do luong do bien dong giua doanh
   thu va quang cao
7 # buoc quan trong trong phân tích thống kê đa biến
8 >>> print(covMatrix)
9 # tra ve ma tran 2x2 voi phương sai tren duong cheo va hiệp
   phương sai ngoai duong cheo
```

---

## 5.7 np.histogram

**Chức năng:** Tính biểu đồ tần suất của dữ liệu, một công cụ quan trọng trong việc hình dung phân phối dữ liệu.

**Cú pháp:** `np.histogram(a, bins=10, range=None, density=False)`

**Tham số:**

- **a:** Mảng đầu vào.
- **bins:** Số lượng hoặc danh sách các khoảng.
- **range:** Tuple (min, max) của dữ liệu.
- **density:** Giá trị Boolean, chuẩn hóa thành xác suất.

**Ứng dụng:** Phân tích phân phối dữ liệu trong các nghiên cứu thống kê.

**Ví dụ:** Tạo histogram cho dữ liệu ngẫu nhiên.

---

```
1 >>> data = np.random.randn(1000)
2 # tao 1000 mau ngau nhien tu phân phối chuẩn
3 # dai dien cho tap du lieu lon de phân tích phân phối
```

```

4 >>> histData, binsData = np.histogram(data, bins=10)
5 # tính histogram với 10 khoảng, đếm số lần xuất hiện trong mỗi
    khoảng
6 # công cụ cơ bản để hiểu phân phối dữ liệu
7 >>> print("Frequencies:", histData)
8 # tra về số lần xuất hiện trong mỗi khoảng
9 >>> print("Bins:", binsData)
10 # tra về các khoảng định nghĩa histogram

```

---

## 6 Thao tác hình dạng và Reshape

### 6.1 reshape

**Chức năng:** Thay đổi hình dạng của mảng mà không làm thay đổi dữ liệu, hỗ trợ chuẩn bị dữ liệu cho các mô hình học máy.

**Cú pháp:** `a.reshape(new_shape)`

**Tham số:**

- **`new_shape` :** *Tuple* hoặc *nguyên* (*-1* tùy ý). **Ứng dụng:** Biến đổi dữ liệu để phù hợp với các lớp đầu vào của mạng nơ-ron.

**Ví dụ:** Chuyển vector thành ma trận 2x3.

```

1 >>> a = np.arange(6)
2 # tạo mảng 1D [0, 1, 2, 3, 4, 5] sử dụng arange
3 # đây là số thường dùng làm cơ sở để thay đổi hình dạng
4 >>> print(a.reshape((2, 3)))
5 # thay đổi mảng thành ma trận 2x3 [[0 1 2], [3 4 5]]
6 # biến đổi thông dụng để điều chỉnh dữ liệu với lớp đầu vào mô
    hình

```

---

### 6.2 flatten

**Chức năng:** Chuyển mảng đa chiều thành mảng 1D, hỗ trợ chuẩn bị dữ liệu cho các thuật toán yêu cầu đầu vào vector.

**Cú pháp:** `a.flatten()`

**Ứng dụng:** Chuẩn bị dữ liệu đầu vào cho các mô hình học máy đơn giản.

**Ví dụ:** Làm phẳng ma trận 2D.

```

1 >>> a = np.array([[1, 2], [3, 4]])
2 # định nghĩa ma trận 2x2
3 # đại diện cho tập dữ liệu nhỏ cần được tuyến tính hóa
4 >>> print(a.flatten())
5 # chuyển ma trận thành mảng 1D [1 2 3 4]
6 # hữu ích để đưa vào thuật toán cần đầu vào vector

```

---

### 6.3 transpose (T)

**Chức năng:** Chuyển vị ma trận bằng cách đổi hàng thành cột, một thao tác quan trọng trong đại số tuyến tính.

**Cú pháp:** `a.T` hoặc `np.transpose(a)`

**Ứng dụng:** Điều chỉnh ma trận trong các bài toán học máy hoặc xử lý hình ảnh.

**Ví dụ:** Chuyển vị ma trận 2x3.

---

```
1 >>> a = np.array([[1, 2, 3], [4, 5, 6]])
2 # tao ma tran 2x3
3 # thuong dung de dai dien du lieu voi 2 mau va 3 dac trung
4 >>> print(a.T)
5 # chuyen vi ma tran thanh 3x2 [[1 4], [2 5], [3 6]]
6 # bien doi can thiet cho cac phép tính ma tran hoac canh chinh
   dac trung
```

---

## 6.4 np.concatenate

**Chức năng:** Ghép nhiều mảng thành một theo trục được chỉ định, hỗ trợ kết hợp dữ liệu từ nhiều nguồn.

**Cú pháp:** `np.concatenate((a1, a2, ...), axis=0)`

**Tham số:**

- **(a1, a2, ...):** Tuple các mảng cần ghép.
- **axis:** Trục ghép (0: hàng, 1: cột).

**Ứng dụng:** Kết hợp tập dữ liệu trong các dự án tiền xử lý dữ liệu lớn.

**Ví dụ:** Ghép hai mảng theo hàng.

---

```
1 >>> a1 = np.array([[1, 2], [3, 4]])
2 # dinh nghia ma tran 2x2, dai dien cho tap con du lieu
3 >>> a2 = np.array([[5, 6]])
4 # dinh nghia ma tran 1x2, tap con khac de ghép lại
5 >>> print(np.concatenate((a1, a2), axis=0))
6 # ghép theo hàng, tra ve [[1 2], [3 4], [5 6]]
7 # phuong phap de hop nhat du lieu theo chieu doc trong tien xu ly
```

---

## 6.5 np.vstack, np.hstack

**Chức năng:** Ghép mảng theo chiều dọc (`vstack`) hoặc chiều ngang (`hstack`), hỗ trợ tổ chức dữ liệu theo cấu trúc mong muốn.

**Cú pháp:** `np.vstack((a1, a2, ...))`, `np.hstack((a1, a2, ...))`

**Ứng dụng:** Kết hợp dữ liệu trong các tác vụ tiền xử lý hoặc phân tích đa chiều.

**Ví dụ:** Ghép theo chiều dọc và ngang.

---

```
1 >>> a1 = np.array([[1, 2], [3, 4]])
2 # dinh nghia ma tran 2x2, phan cua tap du lieu lon hon
3 >>> a2 = np.array([[5, 6]])
4 # dinh nghia ma tran 1x2, phan khac de ghép vào
5 >>> print(np.vstack((a1, a2)))
6 # ghép theo chiều dọc, tra ve [[1 2], [3 4], [5 6]]
7 # huu ích để thêm hàng mới trong du lieu thoi gian
8 >>> print(np.hstack((a1, a2.T)))
```

```
9 # ghép theo chiều ngang sau khi chuyển vì a2
10 # tra về [[1 2 5], [3 4 6]], lý tưởng để kết hợp đặc trưng
```

---

## 6.6 np.unique

**Chức năng:** Trả về các giá trị duy nhất trong mảng, hỗ trợ phân tích tần suất và loại bỏ trùng lặp trong dữ liệu.

**Cú pháp:** `np.unique(ar, returnCounts=False)`

**Tham số:**

- **ar:** Mảng đầu vào.
- **returnCounts:** Giá trị Boolean, trả về số lần xuất hiện.

**Ứng dụng:** Phân tích dữ liệu phân loại hoặc chuẩn bị dữ liệu cho mô hình.

**Ví dụ:** Tìm giá trị duy nhất và tần suất.

---

```
1 >>> data = np.array([1, 2, 2, 3, 1, 4])
2 # tạo mảng với các giá trị trùng lặp
3 # mô phỏng tập dữ liệu do đo nhiều lần
4 >>> uniqueVals, counts = np.unique(data, returnCounts=True)
5 # xác định các giá trị duy nhất và tần suất xuất hiện
6 # bước quan trọng trong phân tích dữ liệu danh mục
7 >>> print(uniqueVals, counts)
8 # tra về: [1 2 3 4] [2 2 1 1], chỉ số lần xuất hiện của mỗi giá trị
```

---

## 6.7 np.vectorize

**Chức năng:** Áp dụng hàm tùy chỉnh lên từng phần tử của mảng, tăng cường khả năng xử lý dữ liệu linh hoạt.

**Cú pháp:** `np.vectorize(pyfunc, otypes=None)`

**Tham số:**

- **pyfunc:** Hàm Python cần áp dụng.
- **otypes:** Kiểu dữ liệu đầu ra.

**Ứng dụng:** Tùy chỉnh xử lý dữ liệu trong các dự án học máy phức tạp.

**Ví dụ:** Áp dụng hàm bình phương.

---

```
1 >>> def squareFunc(x):
2 # định nghĩa hàm tính bình phương của một số
3 # biến đổi đơn giản để trình bày
4     return x**2
5 >>> vecSquare = np.vectorize(squareFunc)
6 # chuyển hàm thành vector để áp dụng lên mọi phần tử mảng
7 # tăng hiệu suất so với vòng lặp Python trên dữ liệu lớn
8 >>> data = np.array([1, 2, 3])
9 # tạo mảng mẫu để kiểm tra hàm vector hóa
```



```
10 >>> print(vecSquare(data))
11 # áp dụng hàm bình phương, trả về [1 4 9]
12 # hữu ích cho các biến đổi phần tử trong tiền xử lý dữ liệu
```

---

## 7 Sắp xếp và Tìm kiếm

### 7.1 np.sort

**Chức năng:** Sắp xếp mảng theo thứ tự tăng dần, hỗ trợ tổ chức dữ liệu trong các tác vụ phân tích.

**Cú pháp:** `np.sort(a, axis=-1, kind='quicksort')`

**Tham số:**

- **axis:** Trục sắp xếp (-1: trục cuối, None: làm phẳng).
- **kind:** Thuật toán ('quicksort', 'mergesort', 'heapsort').

**Ứng dụng:** Sắp xếp dữ liệu số hoặc chuỗi trong các bài toán tối ưu hóa.

**Ví dụ:** Sắp xếp doanh thu tăng dần.

---

```
1 >>> revData = np.array([120, 100, 150, 130])
2 # định nghĩa mảng dữ liệu doanh thu
3 # đại diện cho dữ liệu tại chính các sắp xếp
4 >>> print(np.sort(revData))
5 # sắp xếp mảng theo thu từ thấp đến: [100 120 130 150]
6 # phép tính thông dụng để xếp hạng hoặc phân tích xu hướng
```

---

### 7.2 np.argsort

**Chức năng:** Trả về chỉ số của phần tử khi sắp xếp tăng dần, hỗ trợ sắp xếp đồng bộ nhiều mảng.

**Cú pháp:** `np.argsort(a, axis=-1)`

**Ứng dụng:** Sắp xếp liên kết dữ liệu trong các tập dữ liệu phức tạp.

**Ví dụ:** Tìm chỉ số của doanh thu khi sắp xếp.

---

```
1 >>> revData = np.array([120, 100, 150, 130])
2 # chưa dữ liệu doanh thu theo thu từ ban đầu
3 # cần chỉ số sắp xếp để đồng bộ dữ liệu liên quan
4 >>> print(np.argsort(revData))
5 # trả về chỉ số [1 0 3 2] để sắp xếp mảng
6 # giúp đồng bộ sắp xếp các mảng liên kết
```

---

### 7.3 np.argmax / np.argmin

**Chức năng:** Trả về chỉ số của phần tử lớn nhất (`np.argmax`) hoặc nhỏ nhất (`np.argmin`), hỗ trợ tìm vị trí cực trị.

**Cú pháp:** `np.argmax(a, axis=None), np.argmin(a, axis=None)`

**Ứng dụng:** Xác định vị trí cực trị trong các chuỗi thời gian hoặc dữ liệu kinh doanh.

**Ví dụ:** Tìm ngày có doanh thu cao nhất.

---

```

1 >>> revData = np.array([120, 100, 150, 130])
2 # đại diện cho dữ liệu doanh thu trong bốn ngày
3 # vị trí giá trị cao nhất là mục tiêu phân tích
4 >>> print(np.argmax(revData))
5 # trả về chỉ số 2, chỉ vị trí của giá trị lớn nhất 150
6 # hữu ích để xác định ngày hiệu suất cao nhất

```

---

## 8 Broadcasting và Đại số tuyến tính

### 8.1 Broadcasting

**Chức năng:** Thực hiện phép toán trên các mảng có kích thước khác nhau bằng cách tự động mở rộng kích thước, tối ưu hóa tính toán trong NumPy.

**Ứng dụng:** Chuẩn hóa dữ liệu hoặc áp dụng phép toán trên toàn mảng mà không cần vòng lặp.

**Ví dụ:** Chuẩn hóa ma trận bằng cách trừ trung bình mỗi cột.

---

```

1 >>> data = np.array([[1, 2], [3, 4], [5, 6]])
2 # định nghĩa ma trận 3x2 đại diện cho tập dữ liệu với hai đặc
   #   trung
3 >>> meanVals = np.mean(data, axis=0)
4 # tính trung bình theo cột [3 4]
5 # bước để đưa dữ liệu về trung tâm trong chuẩn hóa thống kê
6 >>> normalizedData = data - meanVals
7 # áp dụng broadcasting để trừ trung bình khỏi mọi hàng
8 # trả về ma trận đã chuẩn hóa cho phân tích tiếp theo
9 >>> print(normalizedData)
10 # trả về: [[-2 -2], [0 0], [2 2]], chỉ độ lệch từ trung bình cột

```

---

### 8.2 np.dot

**Chức năng:** Thực hiện nhân ma trận hoặc tích vô hướng, một phép toán cốt lõi trong đại số tuyến tính và học máy.

**Cú pháp:** `np.dot(a, b)`

**Ứng dụng:** Tính toán trong mạng nơ-ron hoặc biến đổi hình học.

**Ví dụ:** Nhân hai ma trận 2x2.

---

```

1 >>> A = np.array([[1, 2], [3, 4]])
2 # định nghĩa ma trận 2x2 A, ma trận hệ số trong hệ phương trình
   #   tuyến tính
3 >>> B = np.array([[5, 6], [7, 8]])
4 # định nghĩa ma trận 2x2 B, đối tượng thu hai để nhân
5 >>> print(np.dot(A, B))
6 # thực hiện nhân ma trận, trả về [[19 22], [43 50]]
7 # phép tính trung tâm trong biến đổi tuyến tính và mạng nơ-ron

```

---

## 8.3 np.linalg

**Chức năng:** Cung cấp các hàm đại số tuyến tính như nghịch đảo ma trận, định thức, giá trị riêng, và giải hệ phương trình, hỗ trợ các bài toán toán học nâng cao.

**Cú pháp:**

- `np.linalg.inv(A)`: Nghịch đảo ma trận (yêu cầu ma trận vuông và khả nghịch).
- `np.linalg.det(A)`: Tính định thức ma trận (định lượng khả nghịch, 0 nếu không khả nghịch).
- `np.linalg.eig(A)`: Tính giá trị riêng và vector riêng (phân tích phổ ma trận).
- `np.linalg.solve(A, b)`: Giải hệ phương trình tuyến tính  $Ax = b$  (A phải khả nghịch).

**Ứng dụng:** Giải hệ phương trình, phân tích thành phần chính trong học máy.

**Ví dụ:** Giải hệ phương trình  $Ax = b$ .

---

```
1 >>> A = np.array([[3, 1], [2, 4]])
2 # định nghĩa ma trận hệ số 2x2 A cho hệ phương trình tuyến tính
3 # đại diện cho mô hình vật lý hoặc kinh tế
4 >>> b = np.array([7, 10])
5 # định nghĩa vector 2x1 b nằm bên phải của phương trình Ax = b
6 >>> x = np.linalg.solve(A, b)
7 # giải hệ Ax = b, với A là ma trận khả nghịch
8 # phương pháp dùng trong tối ưu hóa và mô phỏng vật lý
9 >>> print(x)
10 # tra về giải pháp: [2. 1.], các giá trị x thỏa mãn Ax = b
11 >>> evals, evecs = np.linalg.eig(A)
12 # tính giá trị riêng và vector riêng của A
13 # cần thiết cho phân tích phổ trong phân tích thành phần chính
14 >>> print(evals)
15 # tra về giá trị riêng, chỉ các hệ số phong to của ma trận
```

---

## 8.4 np.einsum

**Chức năng:** Thực hiện phép toán tensor hiệu quả bằng ký hiệu Einstein, cho phép định nghĩa tùy chỉnh cách các chiều được kết hợp, tối ưu hóa tính toán phức tạp.

**Cú pháp:** `np.einsum(subscripts, *operands)`

**Tham số:**

- **subscripts:** Chuỗi ký hiệu (ví dụ: `'ij,jk->ik'` cho nhân ma trận, `'ij,ij->'` cho tổng tích từng phần tử), mô tả cách các chiều được ánh xạ.
- **operands:** Các mảng đầu vào tham gia phép toán.

**Ứng dụng:** Tính toán tensor trong học sâu, đặc biệt hữu ích với dữ liệu đa chiều.

**Ví dụ:** Nhân ma trận và tính tổng theo hàng.

---

```
1 >>> A = np.array([[1, 2], [3, 4]])
2 # định nghĩa ma trận 2x2 A, mau tensor de trình bày
3 >>> B = np.array([[5, 6], [7, 8]])
```

```

4 # định nghĩa ma trận 2x2 B, đối tượng khác cho phép toán
5 >>> result = np.einsum('ij,ij->', A, B)
6 # sử dụng ký hiệu Einstein 'ij,ij->' để tính tổng tích từng phần
  tu
7 # tương đương với tích vô hướng của mảng đã đề, trả về 70 (1*5 +
  2*6 + 3*7 + 4*8)
8 # phương pháp tối ưu hóa cho các phép rút gọn tensor
9 >>> print("Sum of products:", result)
10 # trả về: 70, kết quả số của tổng tích từng phần tu
11 >>> sumRows = np.einsum('ij->i', A)
12 # tính tổng phần tử theo hàng sử dụng ký hiệu 'ij->i'
13 # giảm mảng 2D xuống mảng 1D với tổng hàng
14 >>> print(sumRows)
15 # trả về: [3 7], tổng của mỗi hàng, hữu ích để tổng hợp đặc trưng

```

---

## 9 Xử lý dữ liệu lớn

### 9.1 np.save, np.load

**Chức năng:** Lưu (np.save) và tải (np.load) mảng vào/từ file định dạng .npy, hỗ trợ lưu trữ dữ liệu số hiệu quả.

**Cú pháp:** np.save(file, arr), np.load(file)

**Tham số:**

- **file:** Tên file hoặc đối tượng file.
- **arr:** Mảng cần lưu.

**Ứng dụng:** Lưu trữ và tái sử dụng dữ liệu lớn trong các dự án học máy hoặc phân tích dữ liệu.

**Ví dụ:** Lưu và tải mảng.

```

1 >>> data = np.array([[1, 2], [3, 4]])
2 # tạo ma trận 2x2, tập dữ liệu nhỏ để lưu trữ
3 >>> np.save('dataFile.npy', data)
4 # lưu mảng NumPy vào file nhị phân .npy
5 # dạng hiệu quả để lưu trữ dữ liệu số lớn
6 >>> loadedData = np.load('dataFile.npy')
7 # tải mảng từ file .npy đã lưu
8 # giúp lấy lại dữ liệu để phân tích hoặc đào mô hình tiếp theo
9 >>> print(loadedData)
10 # hiển thị ma trận đã tải lại: [[1 2], [3 4]]

```

---

## 10 Ví dụ thực hành từ Kaggle

Phần này cung cấp các bài tập thực hành lấy cảm hứng từ nền tảng Kaggle, nhằm giúp học viên áp dụng kiến thức NumPy vào các tình huống thực tế trong Khoa học Dữ liệu và AI.

## 10.1 Bài 1: Chuẩn hóa dữ liệu nhiệt độ

**Bài toán:** Chuẩn hóa danh sách nhiệt độ về z-score để chuẩn bị dữ liệu cho phân tích thống kê.

**Lời giải:**

---

```
1 >>> import numpy as np
2 >>> temps = np.array([25, 28, 30, 27, 29, 31, 26])
3 # định nghĩa mảng dọc nhiệt độ bằng Celsius
4 # đại diện cho dữ liệu thời tiết trong một tuần
5 >>> meanTemp = np.mean(temps)
6 # tính trung bình nhiệt độ trên tập dữ liệu
7 # đo lường xu hướng trung tâm để chuẩn hóa
8 >>> stdTemp = np.std(temps)
9 # tính độ lệch chuẩn, chỉ đo biến động nhiệt độ
10 >>> zScores = (temps - meanTemp) / stdTemp
11 # chuẩn hóa nhiệt độ thành z-score
12 # biến đổi dữ liệu về phân phối chuẩn để phân tích thống kê
13 >>> print("Original temperatures:", temps)
14 >>> print("Z-scores:", zScores)
15 # tra về nhiệt độ ban đầu và z-score của chúng
16 # với trung bình ~0 và độ lệch chuẩn ~1
```

---

## 10.2 Bài 2: Tìm ngày doanh thu cao nhất

**Bài toán:** Xác định ngày có doanh thu cao nhất và thấp nhất trong một chuỗi dữ liệu kinh doanh.

**Lời giải:**

---

```
1 >>> salesData = np.array([120, 100, 150, 130, 170, 110, 140])
2 # đại diện cho dữ liệu doanh thu trong bảy ngày
3 # tập dữ liệu thời gian cho phân tích kinh doanh
4 >>> maxDay = np.argmax(salesData)
5 # xác định chỉ số của giá trị doanh thu cao nhất
6 # chỉ ngày có hiệu suất cao nhất
7 >>> minDay = np.argmin(salesData)
8 # xác định chỉ số của giá trị doanh thu thấp nhất
9 # chỉ ngày có hiệu suất thấp nhất
10 >>> print("Revenue:", salesData)
11 >>> print("Highest day:", maxDay, "with revenue", salesData[
    maxDay])
12 # tra về ngày và doanh thu cao nhất, ví dụ: "Highest day: 4 with
    revenue 170"
13 >>> print("Lowest day:", minDay, "with revenue", salesData[minDay
    ])
14 # tra về ngày và doanh thu thấp nhất, ví dụ: "Lowest day: 1 with
    revenue 100"
```

---

## 10.3 Bài 3: Tính toán ma trận

**Bài toán:** Tính nghịch đảo của ma trận 2x2 và nhân với vector để giải hệ phương trình tuyến tính.

**Lời giải:**

---

```
1 >>> A = np.array([[4, 7], [2, 6]])
2 # định nghĩa ma trận hệ số 2x2 A cho hệ phương trình tuyến tính
3 # đại diện cho mô hình vật lý hoặc kinh tế
4 >>> b = np.array([8, 10])
5 # định nghĩa vector 2x1 b nằm bên phải của phương trình Ax = b
6 >>> AInv = np.linalg.inv(A)
7 # tính nghịch đảo của ma trận A
8 # yêu cầu A là ma trận vuông và khả nghịch, bước để giải hệ
9 >>> result = np.dot(AInv, b)
10 # nhân nghịch đảo A với vector b để giải cho x trong Ax = b
11 # phương pháp thông dụng trong bài toán tối ưu hóa
12 >>> print("Matrix A:\n", A)
13 >>> print("Inverse A:\n", AInv)
14 >>> print("Result:", result)
15 # trả về ma trận ban đầu, nghịch đảo và giải pháp [2. 1.]
```

---

## 10.4 Bài 4: Phân tích dữ liệu phân loại

**Bài toán:** Đếm tần suất các giá trị duy nhất trong tập dữ liệu phân loại để chuẩn bị cho mô hình học máy.

**Lời giải:**

---

```
1 >>> categories = np.array(['A', 'B', 'A', 'C', 'B', 'A', 'C', 'C']
2                               ])
3 # tạo mảng dữ liệu danh mục với nhãn trùng lặp
4 # mô phỏng tập dữ liệu như phân đoạn khách hàng
5 >>> uniqueCats, counts = np.unique(categories, returnCounts=True)
6 # xác định các danh mục duy nhất và tần suất xuất hiện
7 # bước tiền xử lý để tạo đặc trưng trong học máy
8 >>> print("Categories:", uniqueCats)
9 >>> print("Frequencies:", counts)
10 # trả về: Categories: ['A' 'B' 'C'] Frequencies: [3 2 3]
# chỉ phân phối các danh mục trong tập dữ liệu
```

---

## 10.5 Bài 5: Xử lý dữ liệu thiếu

**Bài toán:** Thay thế giá trị NaN bằng trung bình của mảng để xử lý dữ liệu không hoàn chỉnh.

**Lời giải:**

---

```
1 >>> data = np.array([1, np.nan, 3, 4, np.nan])
2 # tạo mảng với giá trị thiếu (NaN)
3 # mô phỏng tập dữ liệu thực tế có độ mất mát
4 >>> meanVal = np.nanmean(data)
```

```

5 # tính trung bình, bỏ qua giá trị NaN
6 # phương pháp mạnh mẽ để xử lý dữ liệu thiếu trong thống kê
7 >>> cleanedData = np.where(np.isnan(data), meanVal, data)
8 # thay thế giá trị NaN bằng trung bình tính được
9 # kỹ thuật đơn giản để chuẩn bị dữ liệu cho phân tích
10 >>> print("Original data:", data)
11 >>> print("Cleaned data:", cleanedData)
12 # tra về dữ liệu ban đầu và dữ liệu đã được xử lý, ví dụ: Original:
    [1. nan 3. 4. nan], Cleaned: [1. 2.666 3. 4. 2.666]

```

---

## 10.6 Bài 6: Phân tích phân phối dữ liệu

**Bài toán:** Tạo histogram cho tập dữ liệu ngẫu nhiên và phân tích phân phối để đánh giá đặc điểm thống kê.

**Lời giải:**

```

1 >>> data = np.random.randn(1000)
2 # tạo 1000 mẫu ngẫu nhiên từ phân phối chuẩn
3 # đại diện cho tập dữ liệu lớn để phân tích phân phối
4 >>> histData, binsData = np.histogram(data, bins=10, density=True
    )
5 # tính histogram với 10 khoảng, chuẩn hóa thành mật độ xác suất
6 # công cụ thống kê để trình bày phân phối dữ liệu
7 >>> print("Normalized frequencies:", histData)
8 # tra về tần suất chuẩn hóa cho mỗi khoảng
9 # tổng bằng 1 do chuẩn hóa mật độ
10 >>> print("Bins:", binsData)
11 # tra về ranh giới các khoảng định nghĩa các khoảng histogram

```

---

## 10.7 Bài 7: Tính toán tensor

**Bài toán:** Sử dụng `np.einsum` để tính tổng tích các phần tử của hai ma trận, minh họa tính toán nâng cao.

**Lời giải:**

```

1 >>> A = np.array([[1, 2], [3, 4]])
2 # định nghĩa ma trận 2x2 A, mau tensor để thực hiện phép toán
3 >>> B = np.array([[5, 6], [7, 8]])
4 # định nghĩa ma trận 2x2 B, đối tượng khác cho tính toán
5 >>> result = np.einsum('ij,ij->', A, B)
6 # sử dụng ký hiệu Einstein 'ij,ij->' để tính tổng tích từng phần
    tu
7 # tương đương với tích vô hướng của mảng đã đề, tra về 70 (1*5 +
    2*6 + 3*7 + 4*8)
8 # phương pháp tối ưu hóa cho các phép rút gọn tensor trong học
    sâu
9 >>> print("Sum of products:", result)
10 # tra về: 70, kết quả của tổng tích từng phần tử

```

---

## Tài liệu Tham khảo

- NumPy Documentation: <https://numpy.org/doc/stable/>
- Học ứng dụng Kaggle: <https://www.kaggle.com/learn/>
- Scipy Lecture Notes: <https://scipy-lectures.org/>

## Hướng dẫn Sử dụng cho Học viên

- Tài liệu này được thiết kế để học viên có thể tự học hoặc sử dụng trong các khóa đào tạo về Khoa học Dữ liệu và AI.
- Mỗi ví dụ đều đi kèm chú thích chi tiết, giúp học viên hiểu rõ ý nghĩa và ứng dụng thực tế của từng phương thức.
- Học viên được khuyến khích thực hành các bài tập thực tế và mở rộng bằng cách tích hợp với các thư viện khác như Pandas, Matplotlib hoặc Scikit-learn.
- Nếu gặp lỗi biên dịch hoặc cần hỗ trợ thêm, hãy liên hệ với người hướng dẫn hoặc tham khảo các tài liệu tham khảo được liệt kê.