

## Numpy

Numerical Python (Numpy) is a powerful Python library for numerical computing, especially when working with large arrays and matrices of numbers.

→ Import numpy as np

ar = np.array([1, 3, 5, 7, 9])

ar1

one dimension

array([1, 3, 5, 7, 9])

ar2 = np.array([[1, 3, 5, 7, 9], [2, 4, 6, 8, 10]])

two dimensions (like matrix)

ar3 = np.array([[[1, 2], [3, 4]], [[5, 6], [7, 8]]])

3d array

Eg :-

ar = np.array([1, 3, 5, 7, 9])

ar

//array([1, 3, 5, 7, 9])

ar.ndim

//

ar = np.array([1, 3, 5, 7, 9], ndim=3)

ar

//array([[[1, 3, 5, 7, 9]]])

ar.ndim

//3

→ to change the dimension

## NumPy Attributes

np.dtype → datatype

np.shape → give the size (no of rows, no of columns)

np.size → how much memory the array is consuming

np.sum → sum

np.sum → total memory consumed by all the elements

Example,

Import numpy as np

ar1 = np.array([1, 2, 3, 4, 5, 6, 7])

ar2 = np.array([[1, 2, 3], [4, 5, 6]])

ar3 = np.array([[[1, 2, 3], [4, 5, 6]], [[3, 4, 5], [6, 7, 8]]])

ar4 = np.array([1, 2, 3, 4], ndmin=4)

ar4.sum

//  
1

ans.sum

//  
3

ans.shape

//  
3, 2

ans.shape

//  
4, 1, 2

ans.dtype

//dtype('float64')

ar1.size

//  
4

comlin  
Date

## np.1. Numpy

//  
0

arr = np.array([1, 2, 2, 3, 3])

arr.dtype

//dtype('float64')

arr.sum

Creating an Array

np.zeros(3) → Create array with zeros.

np.ones(3) → Create array with ones.

np.empty(3) → Create an array but elements may be anything (garbage value).

np.eye(3) → Creating an identity matrix (square matrix) 3x3, 4x4, ...

np.diag(3) → It creates a square matrix.

np.arange(start, end, step) → Some as range function of Python

np.linspace(0, 10, num=5) → It creates an array with 'num=' elements.

start  
stop  
end

Import numpy as np

→ ar = np.zeros(3)

ar

Import numpy as np

//array([0., 0., 0.])

ar.dtype

//dtype('float64')

→ ar = np.zeros((3, 4))

ar

//array([[0., 0., 0., 0.],

[0., 0., 0., 0.],

[0., 0., 0., 0.]])



arr[0:10:2]

arr[1]

// array([2, 6, 10, 14, 18])

arr[1:8:3]

// array([4, 10, 16])

arr = np.linspace(5, 12, 6)  
arr  
// array([5., 6.4, 7.8, 9.2, 10.6, 12.])

Array slicing, slicing

Interact with arr

arr = np.array([2, 4, 6, 8, 10, 12, 14, 16, 18, 20])

arr[0:7]

// array([2, 4, 6, 8])

arr[4]

arr[1:-4]

arr[14]

arr[4:3]

arr[44]

arr[4:-3]

arr[-4]

arr[0:2]

// array([[2, 4, 6, 8],  
[1, 3, 5, 7]])

arr2[0:3:2]

// array([[2, 4, 6, 8],  
[11, 22, 33, 44]])

arr2[:, 0:2]

// array([[2, 4],  
[1, 3],  
[11, 22],  
[10, 20]])

arr[4:-3]

arr  
// array([6, 8, 10, 12, 14, 16, 18])



-128 to 127

```
ar = np.array([1,2,3,4], 'l')
```

```
ar.dtype
```

```
// dtype('int32')
```

```
ar = np.array([1,2,3,127], 'b')
```

```
ar.dtype
```

```
// dtype('int8')
```

```
ar = np.array([1,2,3,4], 't')
```

```
ar.dtype
```

```
// dtype('int32')
```

```
ar
```

```
// array([1,2,3,127], dtype='int8')
```

```
ar = np.array([1,2,3,129], 'b')
```

```
ar
```

```
// array([1,2,3,-127], dtype='int8')
```

```
ar
```

```
ar = np.array([1,2,3,4], 'f')
```

```
ar.dtype
```

```
// dtype('float32')
```

```
ar
```

```
ar = np.array([1,2,3,4], 'F')
```

```
ar.dtype
```

```
// dtype('complex64')
```

```
ar
```

```
// array([1.+0.j, 2.+0.j, 3.+0.j, 4.+0.j], dtype=complex64)
```

```
ar = np.array([1,2,3,4], 'D')
```

```
ar.dtype
```

```
// dtype('uint8')
```

```
// array([1,2,3,4], dtype='uint8')
```

```
ar
```

```
// array([1,2,3,4], dtype='uint8')
```

```
ar
```

```
ar = np.array([1,2,3,4], 'O')
```

```
ar
```

```
// array([1,2,3,4], dtype='Object')
```

```
ar = np.array([1,2,3,4], 'O')
```

```
ar
```

```
// array([1,2,3,4], dtype='Object')
```

```
ar=np.array([1,2,3,4], 's')
ar.dtype
// dtype('S1')

ar
// array(['b'1', b'2', b'3', b'4'], dtype='S1')

ar=np.array(['a','b','c','d'],'U')
ar.dtype
// dtype('U1')

ar
// array(['a', 'b', 'c', 'd'], dtype='U1')

ar=np.array([1,2,3,0], 'bool')
ar.dtype
// dtype('bool')

ar
// array([True, True, True, False])
```

$\rightarrow$  Path in geschr.

- \* It is the process of gathering, managing, and analyzing data efficiently.

\* Plays a fundamental step in the data processing pipeline.

\* It involves the creation, updating, transfer, and cleaning of raw data from diverse external sources into structured files or storage mechanisms, where it awaits further processing and analysis.

Microbathing

A hybrid wave pattern

June 1970 - small batches 107 data,  
Processed frequently

# The Data Lakes Weekflow

- \* Data source identification - Identify and assess the data source, understand the data format, structure and access method

Type Amphibian - tree lany

- x Batch

Real-time

- using a hybrid or a  
sequential, enabling robust  
processing.

→ Batch Processing

- \* Data Cleaning
  - \* Data Validation
  - \* Data Transformation
  - \* Data Feeding
  - \* Data Monitoring

Blossom - we can't

- Data as collected and  
processed at scheduled intervals  
backups

Data collection: Data from various sources

- Characteristics Latency, Data size, Processing, Scheduling, Fault.

## Example of Batch Processing

### Architecture

Data Sources  
(Databases, logs, CSV)

- 1) Payroll System
- 2) Bank Statement Generation
- 3) E-commerce Order Processing
- 4) Credit Card Billing

20

25

30

35

40

45

50

55

60

65

70

75

80

85

90

95

100

105

110

115

120

125

130

135

140

145

150

155

160

165

170

175

180

185

190

195

200

205

210

215

220

225

230

235

240

245

250

255

260

265

270

275

280

285

290

295

300

305

310

315

320

325

330

335

340

345

350

355

360

365

370

375

380

385

390

395

400

405

410

415

420

425

430

435

440

445

450

455

460

465

470

475

480

485

490

495

500

505

510

515

520

525

530

535

540

545

550

555

560

565

570

575

580

585

590

595

600

605

610

615

620

625

630

635

640

645

650

655

660

665

670

675

680

685

690

695

700

705

710

715

720

725

730

735

740

745

750

755

760

765

770

775

780

785

790

795

800

805

810

815

820

825

830

835

840

845

850

855

860

865

870

875

880

885

890

895

900

905

910

915

920

925

930

935

940

945

950

955

960

965

970

975

980

985

990

995

1000

Adv

Efficiency at Scale

Automated Scheduling

Error Handling and Retry Logic

Resource Optimization

Reduced manual intervention

Consistency and Standardization

Limitations

Delayed Results

Resource Spikes

Complex Debugging

Limited Flexibility

Data Staleness

Storage Layer  
(HDFS, S3, etc.)

Batch Processing  
(Spark, Flink, Hadoop)

Proteected Data  
(DW, DBS, Reports)

B (Tool) Reports  
(Gralogan, ELK, etc.)

Common Batch Processing Tool

Apache Spark

Apache Flink

Talend

Aws Glue

## Stream Processing

Date Inspection - Data is continuously received from  
Source like Kafka, Sensors, Log, API etc.

Context Delivery - Result are pushed immediately to dashboards,  
dashboards, alerting systems, or API's.

### Characteristics

Volume, Data size, Processing, Use case, Latency tolerance

Fast

Example

- 1) Fraud detection in Banking
- 2) Real-Time Price matching (Uber, Ola)
- 3) Real-Time Analytics for E-commerce
- 4) Social media feeds and Notifications,
- 5) Log Monitoring and Alerting
- 6) Online Gaming Leaderboards.
- 7) Smart Traffic Systems

## Architectures

Continuous Data Sources

Data ingestion



Data processing engine



Alerting System  
(Email, SMS, Push, etc.)

Downstream DB

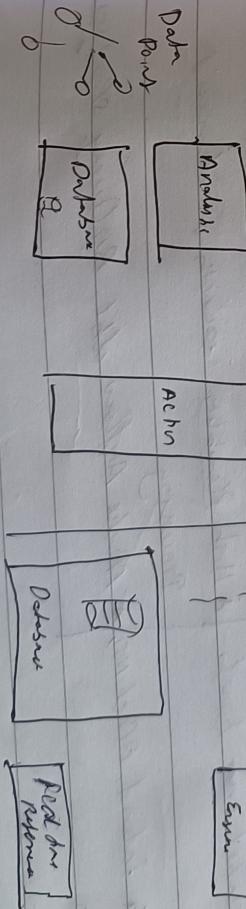
(NoSQL, DW)

Read time

Data -> Batch

Analytics Engine

Analytics Engine



Kafka (Support Stream Processing)

Continuous (Batch window or never)

## Kafka Architecture

Producers  
Consumers

Brokers

Topic

Partitions

## ETL

[Source System]  
↓ Extract

Cloud Data warehouse + Load

[Logging Area or ETL tool]  
↓ Transform (Clean, bin, enrich)

↓ Transform (ustering, segments)

[Data warehouse or Database] + Load

[Load Structured Data] on working

↓  
To Store  
Structured  
data

ETL (Extract → Transform → Load)

Banking and Finance

Healthcare

Retail Pos Systems

Insurance

ELT use case

Streaming Analytics / IoT

E-commerce

Social Media Analytics

Real-time Fraud Detection

High Scalability with on-demand

Processing

Informatica, Talend, Alteryx, Apache Airflow, abt, BigQuery, Snowflake

partake

## SQl

### Database

// is an organized collection of Structured data

e.g. folder in a computer

- > CREATE DATABASE databases;
- > SHOW DATABASES;
- > USE databases;
- > CREATE TABLE Products (product\_id INT, product\_name VARCHAR(50), price FLOAT);
- > SHOW TABLES;
- > INSERT INTO Products (product\_id, product\_name, price) VALUES (1, "iPhone 13 Pro", 120000),  
(2, "iPhone 13", 60000);
- > SELECT \* FROM Products;
- > SELECT Product\_name FROM Products;
- > SELECT Product\_name, Price FROM Products;
- > Modify
- > SELECT \* FROM Products;
- > SELECT Product\_name, brand FROM Products;
- > SELECT \* FROM Products WHERE Price >= 50000;
- > SELECT Product\_name FROM Products WHERE brand = 'Apple'.
- > SELECT Product\_name, Price FROM Products WHERE Price BETWEEN 150000 AND 200000
- > SELECT \* FROM Products\_new;

### Adding a column

→ ALTER TABLE Products\_new ADD brand VARCHAR(40);

→ Brand is still null, so update it.

→ UPDATE Products\_new  
SET brand = 'Apple'  
WHERE Product\_id;

→ DROP TABLE Products\_new;

→ DROP DATABASE databases;

### SELECT

- SELECT brand FROM Product WHERE brand LIKE '%Apple' (Case Insensitive)
- SELECT brand FROM Product WHERE brand LIKE '%apple%' (Case Sensitive)
- SELECT brand FROM Product WHERE brand = 'apple' (Second Way)
- SELECT brand FROM Product WHERE brand LIKE 'apple' (Third Way)
- SELECT \* FROM Product WHERE price LIKE '1000000'
- SELECT \* FROM Product WHERE brand IN ('Redmi', 'Realme')

### COUNT

- SELECT COUNT(\*) FROM Product;
- SELECT COUNT(\*) FROM Product WHERE value;
- SELECT COUNT(\*) FROM Product WHERE value > 10;
- SELECT COUNT(\*) FROM Product WHERE value < 10;
- SELECT COUNT(\*) FROM Product WHERE NOT value < 10;
- SELECT COUNT(\*) FROM Product WHERE NOT value < 10;

### SUM

- SELECT SUM(price) FROM Product;
- SELECT AVG(price) FROM Product;
- SELECT MIN(price) FROM Product;
- SELECT MAX(price) FROM Product;

### GROUP BY

- SELECT COUNT(\*) FROM Product GROUP BY brand ASC;
- SELECT COUNT(\*) FROM Product GROUP BY brand DESC;

### AND, OR, NOT

- SELECT \* FROM Product WHERE Price > 100000 AND Product\_id < 10;
- SELECT \* FROM Product WHERE Price > 100000 OR Product\_id < 10;
- SELECT \* FROM Product WHERE Product\_id < 10;
- SELECT \* FROM Product WHERE NOT Product\_id < 10;

### LIMIT

- SELECT \* FROM Product LIMIT 10;
- SELECT \* FROM Product LIMIT 10 OFFSET 5;

### Average Function

Count	Min	Max	Avg
Sum			

Table 2

Specifications

Joining

ID RAM

ROM

1 6

128

1 8

128

2 16

256

2 3

128

6 6

64

7 4

128

7 8

256

7 9

512

Types of join

Inner Join

Left Join

01

a) Get RAM and ROM details of Order #33  
specifications



Inner Join



Left Join

Right Join

Full Join

01

→ SELECT Products.Product\_id, Products.Product\_name,

Specifications.RAM, Specifications.ROM FROM Product

JOIN Specifications

ON Product.Product\_id = Specifications.ID;

→ SELECT Products.ProductId, Products.ProductName, Specifications.RAM,

Specifications.ROM FROM Product LEFT JOIN Specifications

ON Product.ProductId = Specifications.ID;

UNION ALL

SELECT Products.ProductId, Products.ProductName, Specifications.RAM,

Specifications.ROM FROM Product RIGHT JOIN Specifications

ON Product.ProductId = Specifications.ID;