



ACC-HPC June 2025

INTRODUCTION TO MACHINE LEARNING

Dr Kiran Waghmare
CDAC Mumbai

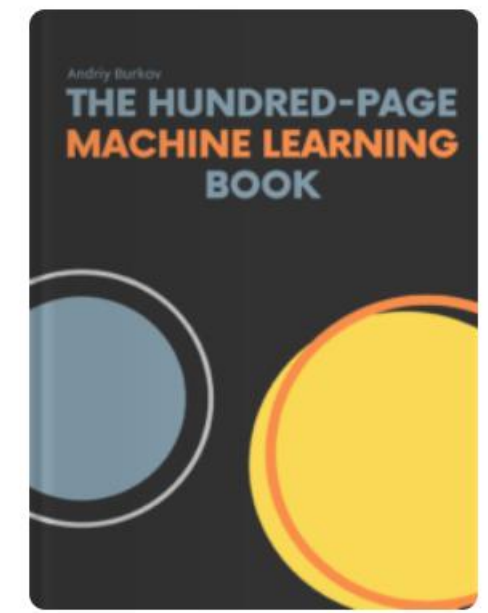
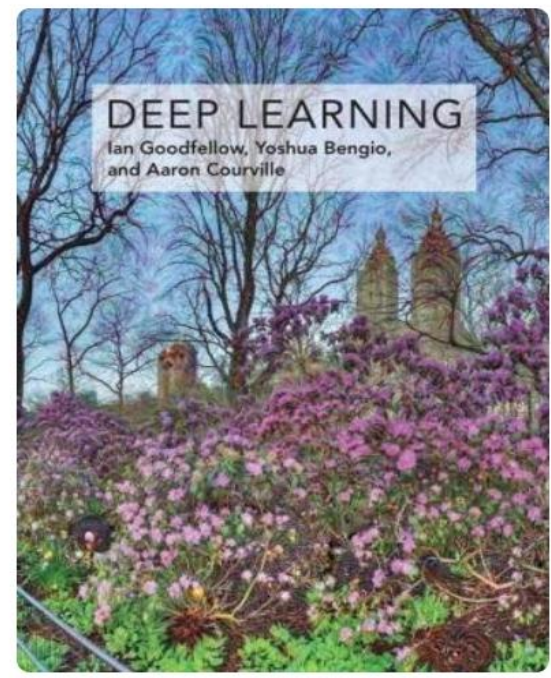
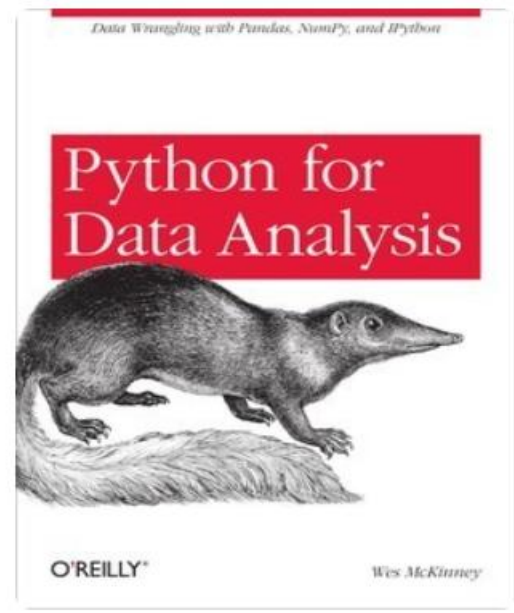
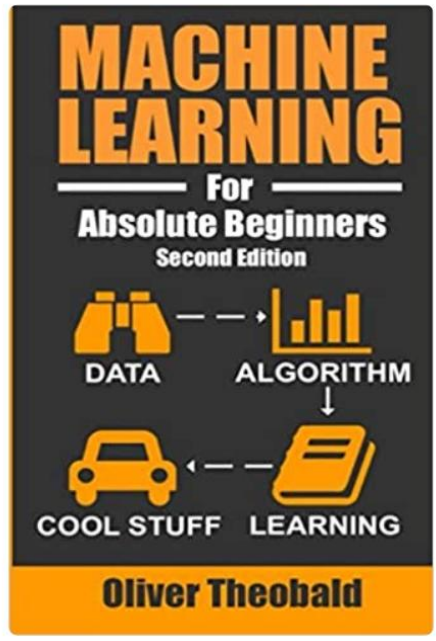
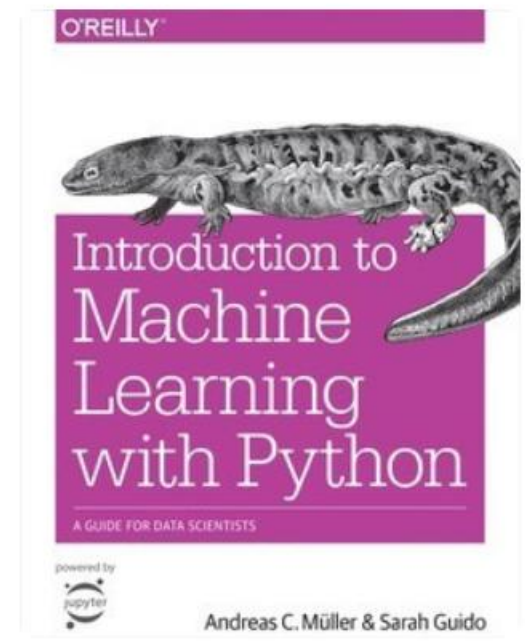
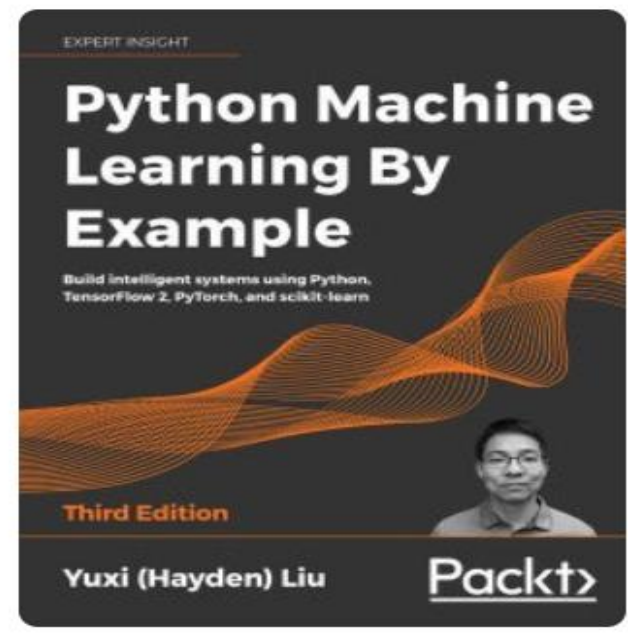
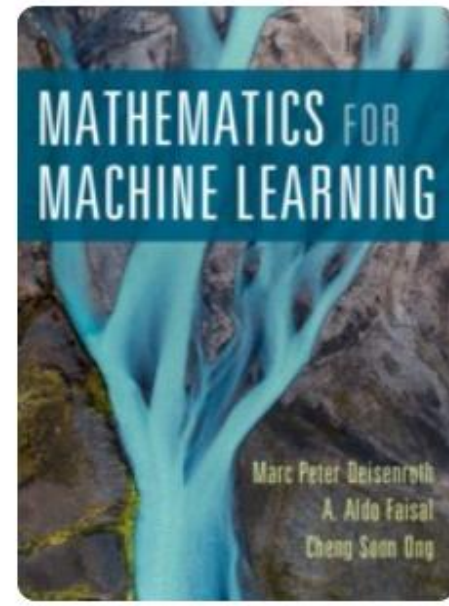
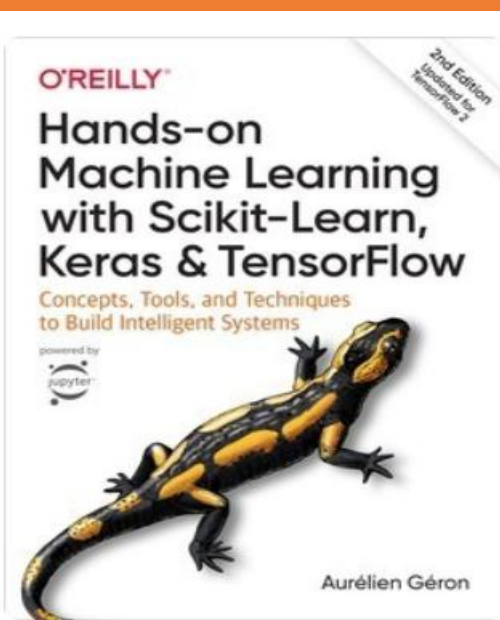
CDAC Mumbai | Kiran Waghmare

Machine Learning >



Machine Learning Roadmap





Agenda

What is Machine Learning?

Applications

Types of ML

ML Workflow

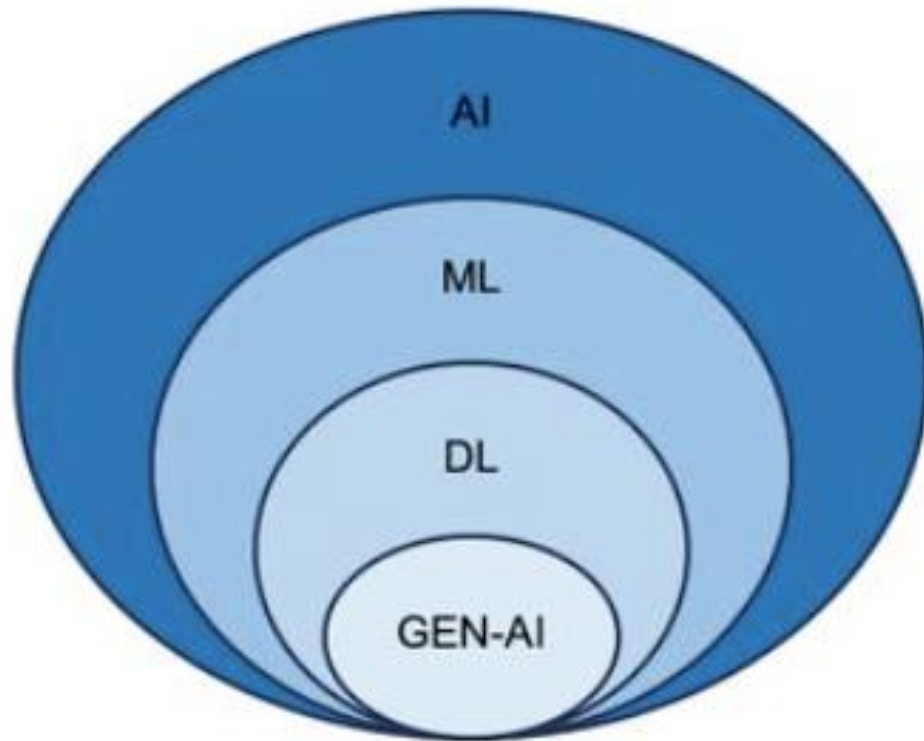
Data Pre-processing

Feature Engineering

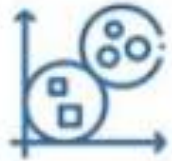
Evaluating ML Models



Overview on AI, ML, DL & GEN-AI



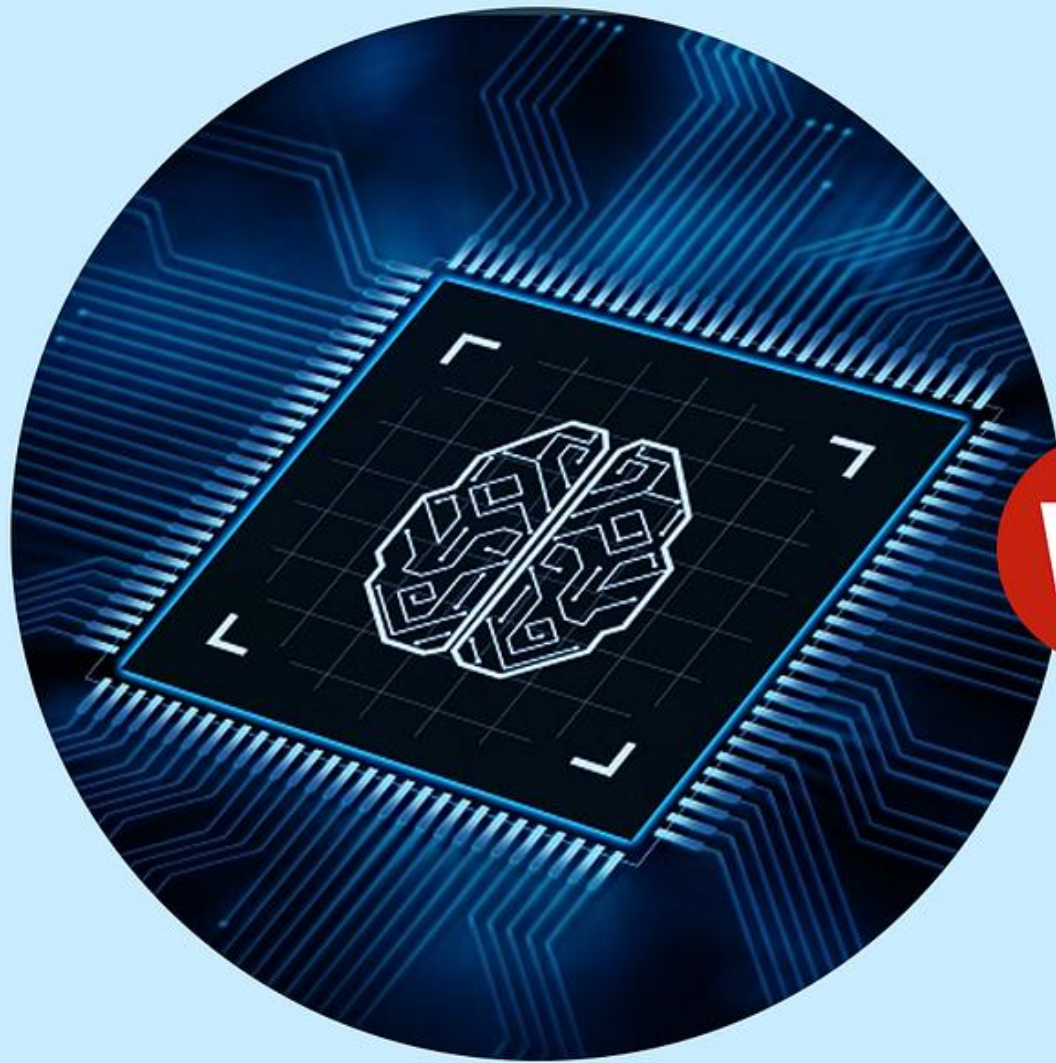
- **AI**
 - A broad field
 - Creating machines to mimic human intelligence
- **Machine learning (ML)**
 - Subset of AI
 - Learn patterns from data
 - Make decisions/predictions
- **Deep learning (DL)**
 - Subset of ML
 - Using neural networks
 - Learn complex patterns
- **Generative AI (GEN-AI)**
 - Subset of DL
 - Generating new content



Comparison

- **Traditional Programming**





Data

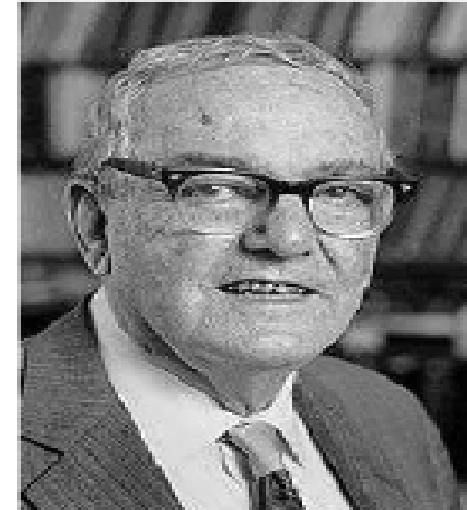
Vs



Information

Machine Learning

- **Herbert Alexander Simon:**
“Learning is any process by which a system improves performance from experience.”
- “Machine Learning is concerned with computer programs that automatically improve their performance through experience. “



Herbert Simon

[Turing Award](#) 1975

[Nobel Prize in Economics](#) 1978

What is Machine Learning?

- [Arthur Samuel, 1959]
 - Field of study that gives computers
 - the ability to learn without being explicitly programmed
- [Kevin Murphy] algorithms that
 - automatically detect patterns in data
 - use the uncovered patterns to predict future data or other outcomes of interest
- [Tom Mitchell] algorithms that
 - improve their performance (P)
 - at some task (T)
 - with experience (E)

The concept of learning in a ML system

- Learning = Improving with experience at some task
 - Improve over task T ,
 - With respect to performance measure, P
 - Based on experience, E .

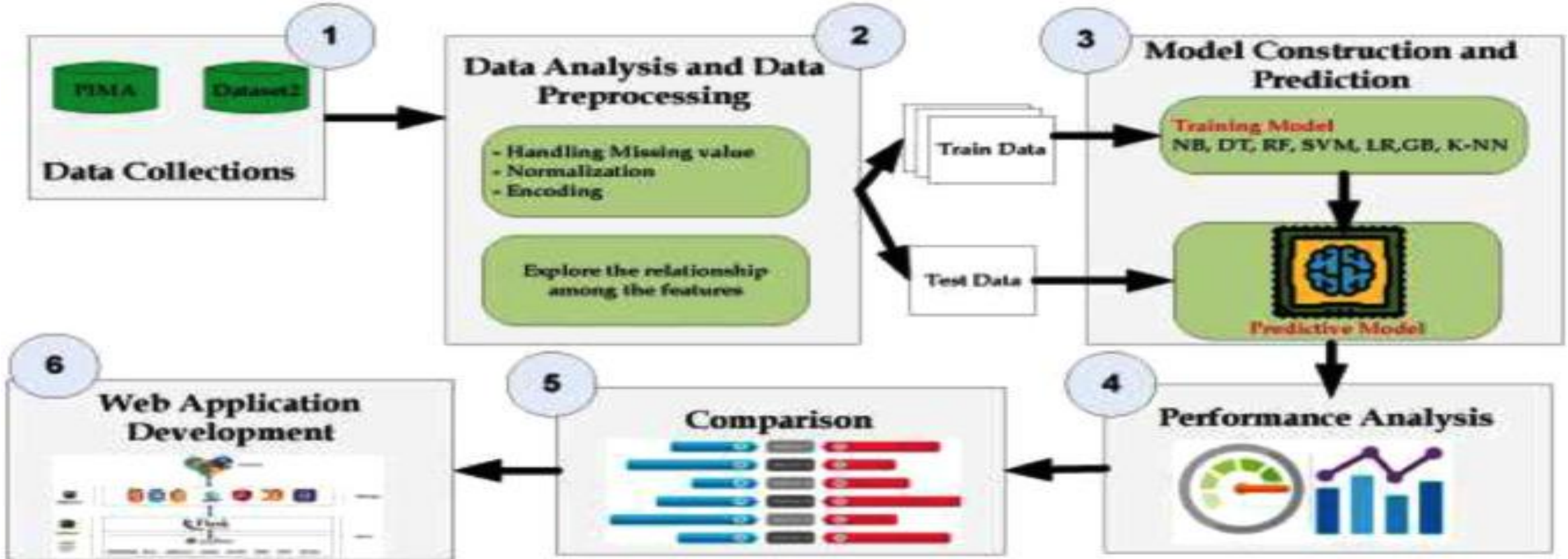
Definition

A computer program is said to learn from **experience E** with respect to some class of **tasks T** and performance **measure P**, if its performance at tasks T , as measured by P , **improves** with experience E .

Examples

- i) Handwriting recognition learning problem
 - Task T : Recognising and classifying handwritten words within images
 - Performance P : Percent of words correctly classified
 - Training experience E : A dataset of handwritten words with given classifications
- ii) A robot driving learning problem
 - Task T : Driving on highways using vision sensors
 - Performance measure P : Average distance traveled before an error
 - training experience: A sequence of images and steering commands recorded while observing a human driver
- iii) A chess learning problem
 - Task T : Playing chess
 - Performance measure P : Percent of games won against opponents
 - Training experience E : Playing practice games against itself

Industry Case Study : 1



Definition

A computer program which learns from experience is called a *machine learning program* or simply a *learning program*. Such a program is sometimes also referred to as a *learner*.

What is Machine Learning?

- If you are a Scientist



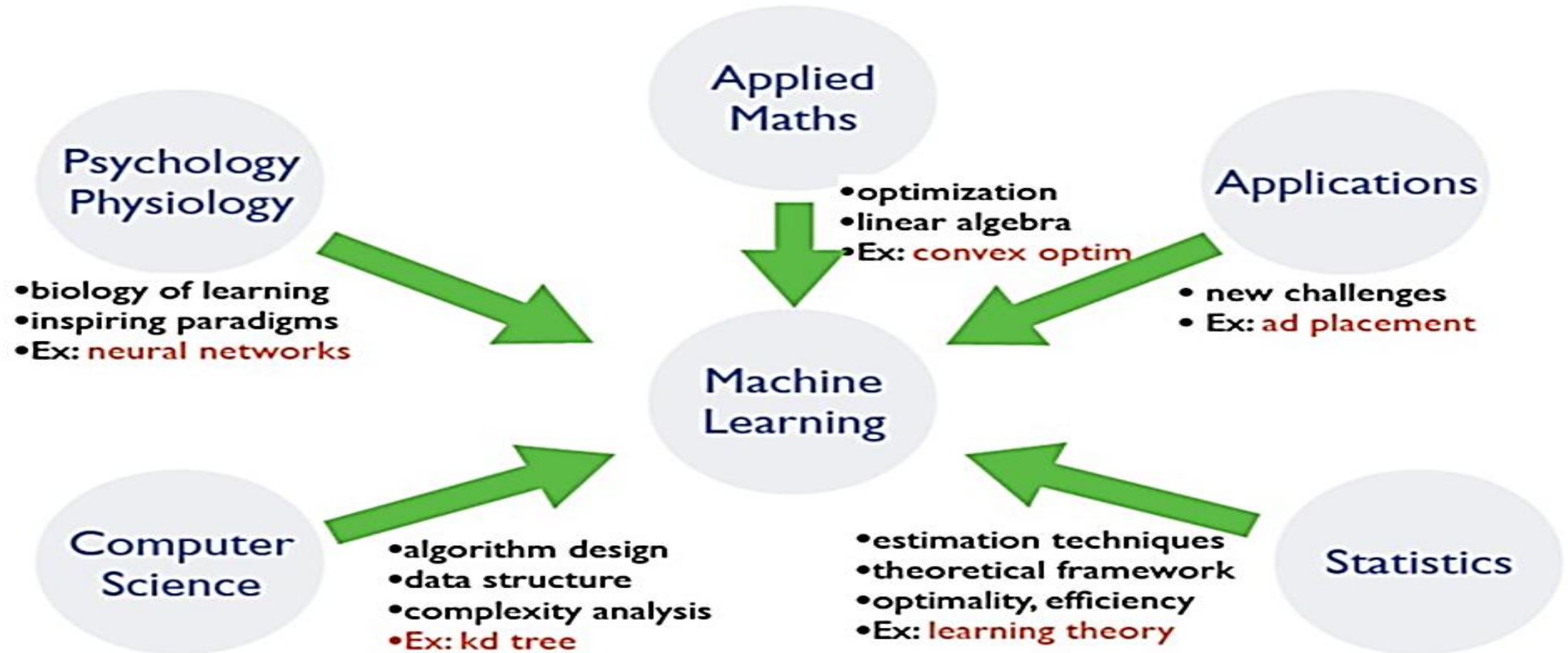
Why Study Machine Learning?

The Time is Ripe

- Algorithms
 - Many basic effective and efficient algorithms available.
- Data
 - Large amounts of on-line data available.
- Computing
 - Large amounts of computational resources available.

Domains include in ML

Where does ML fit in?



Google



Google Search

I'm Feeling Lucky



Google

Q AUSTRALIA NEWS

Q australia news - Google Search

Q australia news today

Q australia news now

Q australia news latest

Q australia news covid

Q australia news code

Q australia news headlines



NETFLIX

Category Codes 2022

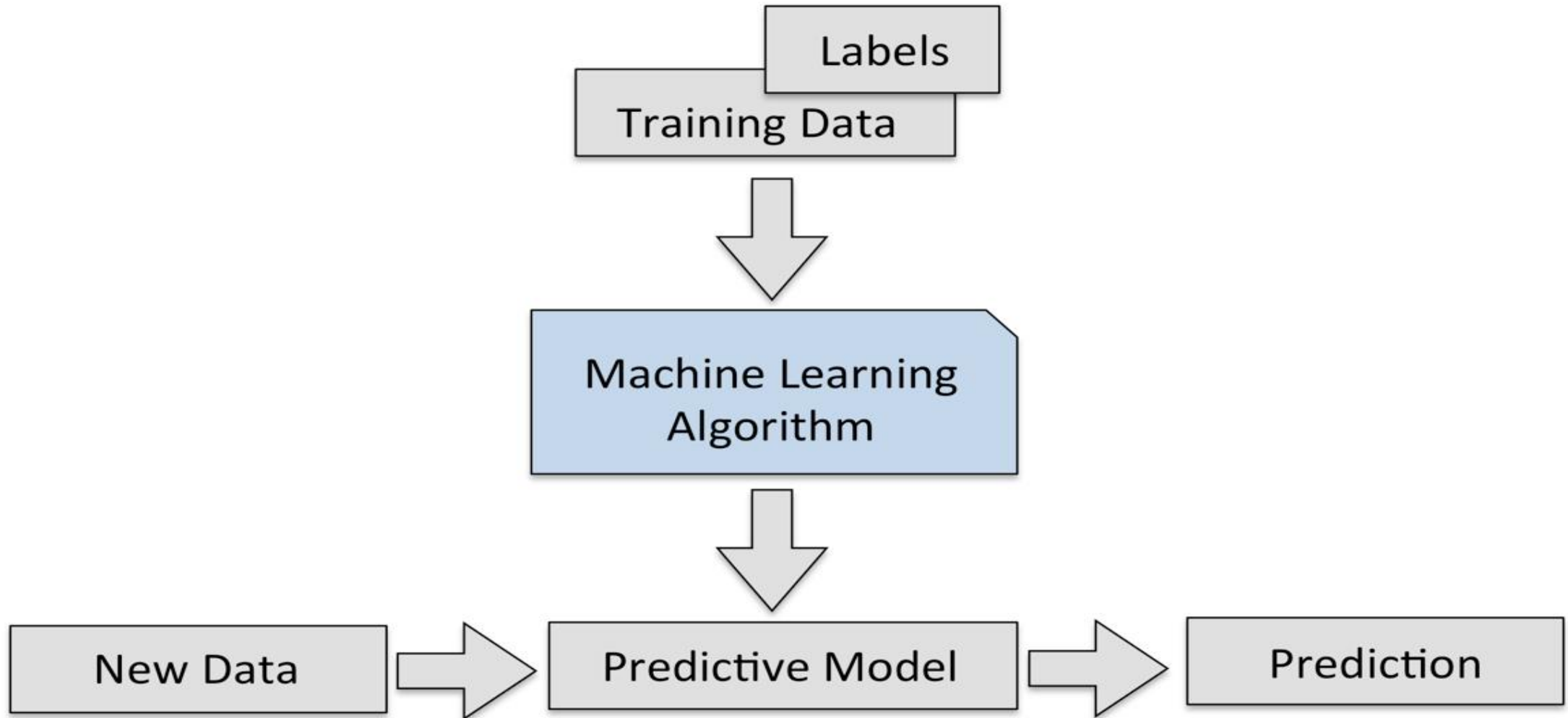
NETFLIX

How does Machine Learning Work?

Input Data → Analyze Data → Find Patterns → Prediction → Stores the Feedback



Machine Learning Model

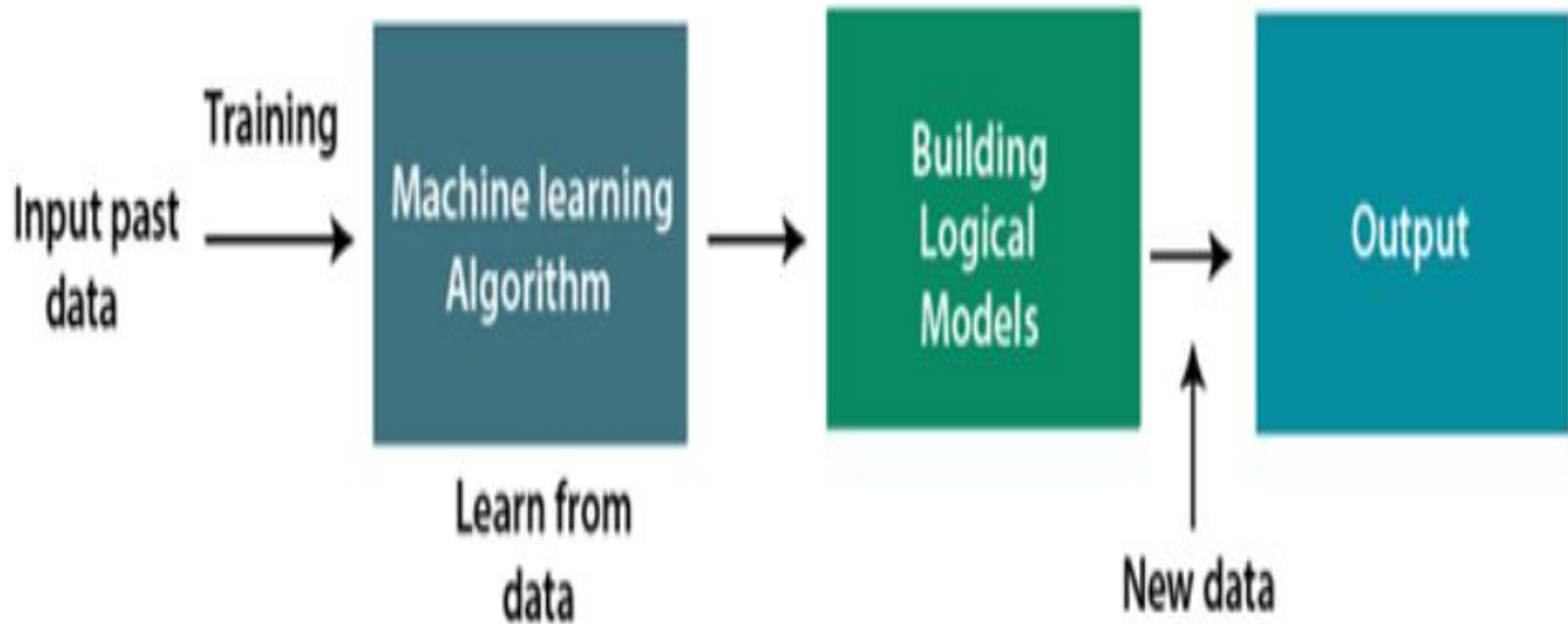


Features of Machine Learning:

- Machine learning **uses data to detect various patterns** in a given dataset.
- It can **learn from past data** and improve automatically.
- It is a **data-driven technology**.
- Machine learning is much **similar to data mining** as it also deals with a huge amount of data.
- **Following are some key points that show the importance of Machine Learning:**
 - **Rapid increment** in the production of data
 - **Solving complex problems**, which are difficult for a human
 - **Decision-making in various sectors** including finance
 - **Finding hidden patterns and extracting useful information** from data.

What is Machine Learning Model?

- **Definition:**
 - Machine Learning is a concept which allows the machine
 - **to learn from examples and experience,**
 - and that **too without being explicitly programmed.**
- Machine Learning algorithms are an evolution of **normal algorithms.**
- They make your **programs “smarter”, by allowing them to automatically learn** from the data you provide.
- The algorithm is mainly divided into:
 - **Training Phase**
 - **Testing phase**



Training Phase : Mango example

- You take a randomly selected specimen of mangoes from the market (**training data**),
- make a table of all the **physical characteristics** of each mango,
 - like color, size, shape, grown in which part of the country,
 - sold by which vendor, etc (**features**),
 - along with the sweetness, juiciness, ripeness of that mango (**output variables/classifier/predictor**).
- You feed this data to the machine learning algorithm (classification/regression), and it learns a model of the correlation between an average mango's physical characteristics, and its quality.

Machine Learning End Product



Ordinary
System



With AI



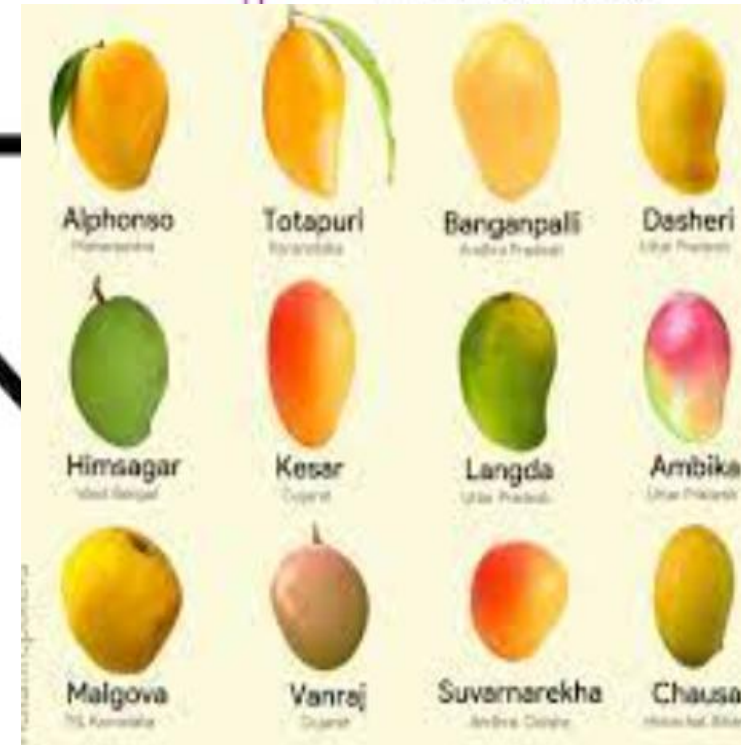
Machine
Learning



Learns



Predicts



Testing Phase

- Next time when you go shopping, you will **measure the characteristics of the mangoes** which you are purchasing(test data)and **feed it to the Machine Learning algorithm**.
- It will **use the model** which was **computed earlier to predict**
 - if the mangoes are sweet, ripe and/or juicy.
- The algorithm may internally **use the rules**, similar to the one you manually wrote earlier (for eg, a decision tree).
- Finally, you can now shop for mangoes with great confidence, without worrying about the details of how to choose the best mangoes.

Conclusion as an Algorithm

- You know what! you can make your algorithm **improve over time** (reinforcement learning) so that it **will improve its accuracy** as it gets trained on more and more training dataset.
- In case it makes a wrong prediction it will update its rule by itself.
- The best part of this is, you can use the **same algorithm to train different models**.
- You can create one each for predicting the quality of apples, grapes, bananas, or whatever you want.

Machine Learning Libraries – Comparison Chart

Category	Library	Primary Purpose	Key Features
Numerical Computing	NumPy	Fast numerical operations	Arrays, linear algebra, mathematical functions
Data Manipulation	Pandas	Data cleaning & analysis	DataFrames, merging, filtering, grouping
Visualization	Matplotlib	Basic plotting	Line, bar, scatter, histograms
	Seaborn	Statistical visualization	Heatmaps, boxplots, pairplots
Machine Learning	Scikit-Learn	Classical ML algorithms	Preprocessing, ML models, pipelines, metrics

What is Scikit-Learn?

- **Main concept**

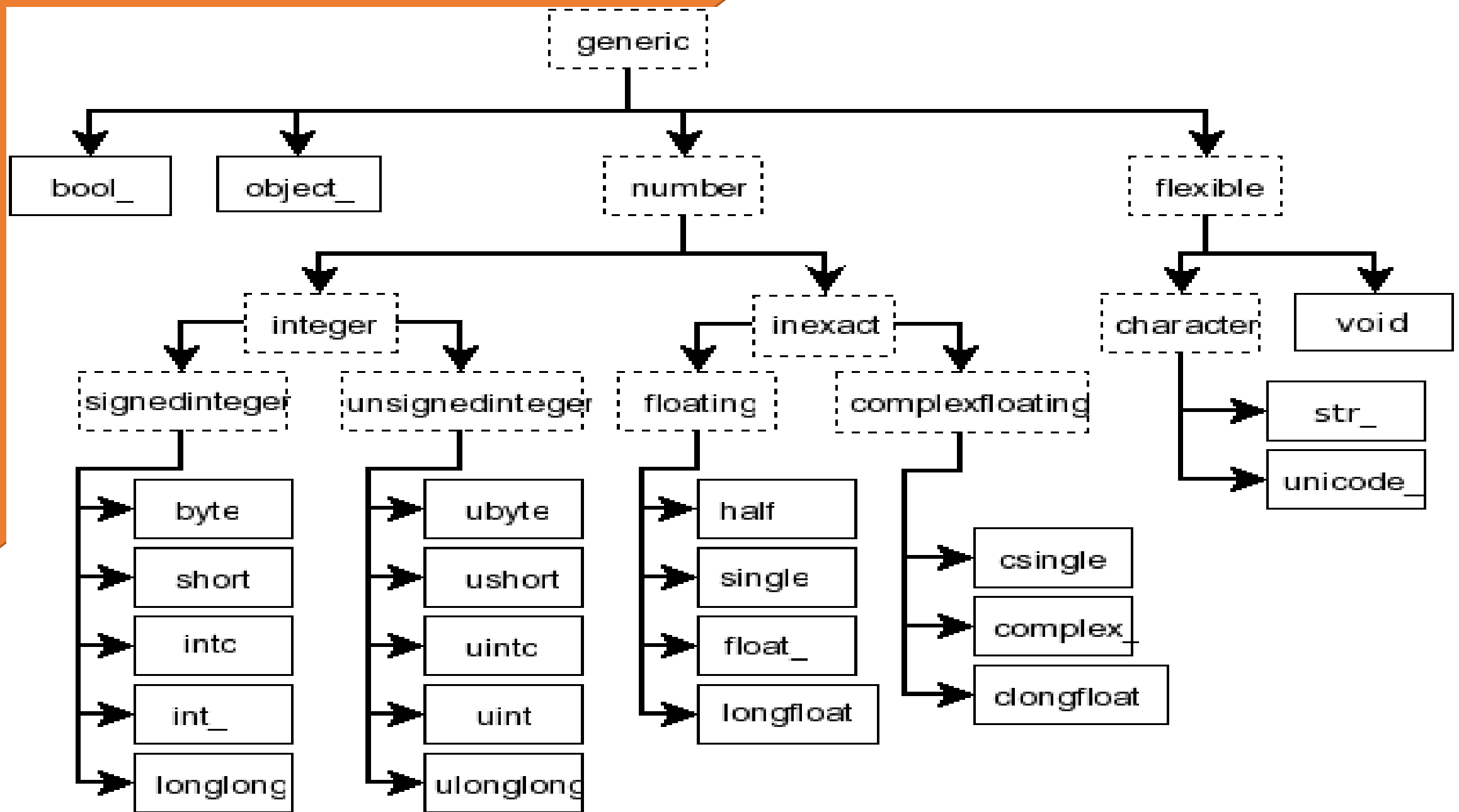
- A powerful open-source machine learning library for Python.
- Built on NumPy, SciPy, and Matplotlib.
- Provides simple, consistent APIs for building ML models.
- Ideal for beginners and production-grade prototypes.

- **Key Features**

- Easy-to-use APIs for ML operations.
- Wide range of ML algorithms: regression, classification, clustering, etc.
- Preprocessing tools: scaling, encoding, feature extraction.
- Model evaluation tools: metrics, cross-validation.
- Pipelines: automate workflow from pre-processing → modeling.
- Built-in toy datasets: Iris, Digits, Wine, Breast Cancer.

Numpy Library

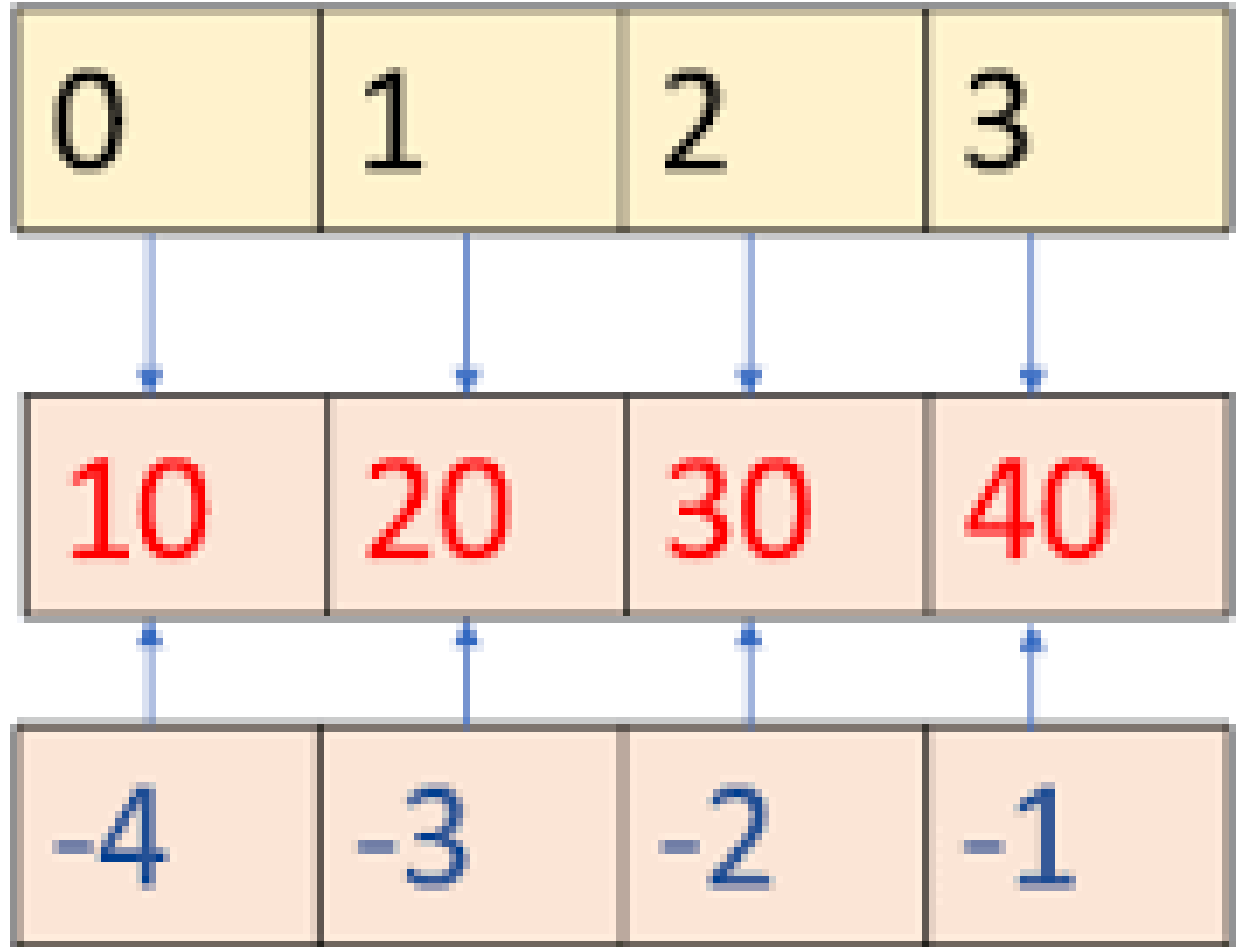
- Provides the foundation for **scientific computing in Python**.
- Supports **multi-dimensional arrays (ndarrays)** that are faster and more efficient than Python lists.
- Enables vectorized operations, reducing the need for loops and making code faster.
- Offers **powerful tools** for linear algebra, statistics, and mathematical operations.
- Forms the core dependency for many **ML libraries** like Pandas, Scikit-Learn, TensorFlow, PyTorch.
- Provides **broadcasting**, allowing operations on arrays of different shapes.
- Handles **large datasets** efficiently, crucial for machine learning workflows.
- Highly optimized using C/C++ backend, giving significant performance boosts.



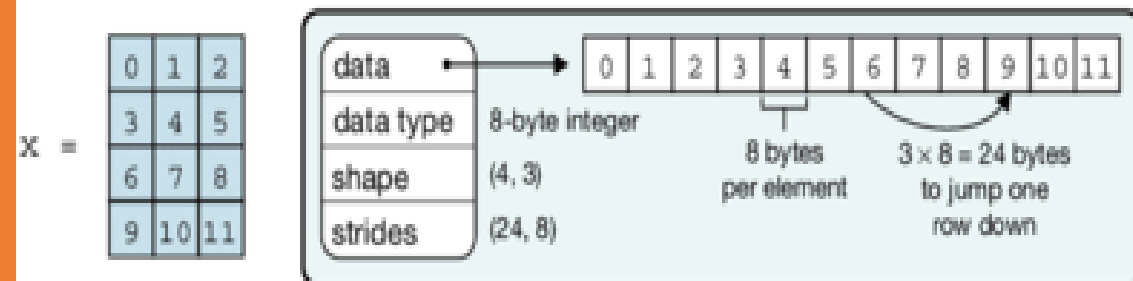
Positive Indexing

One_d

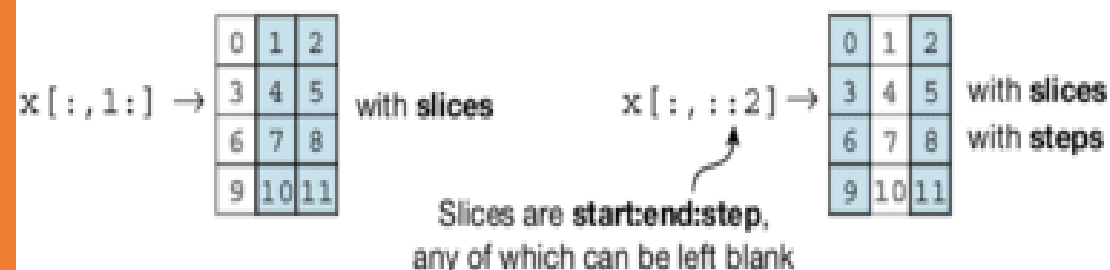
Negative Indexing



a Data structure



b Indexing (view)



c Indexing (copy)

$x[1, 2] \rightarrow 5$ with scalars $x[x > 9] \rightarrow$

10	11
----	----

 with masks

$x[\begin{bmatrix} 0 & 1 \\ 1 & 2 \end{bmatrix}] \rightarrow [x[0, 1], x[1, 2]] \rightarrow$

1	5
---	---

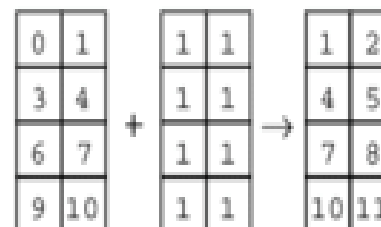
 with arrays

$x[\begin{bmatrix} 1 & 1 & 0 \\ 2 & & \end{bmatrix}] \rightarrow x[\begin{bmatrix} 1 & 1 & 1 & 0 \\ 2 & 2 & 1 & 0 \end{bmatrix}] \rightarrow$

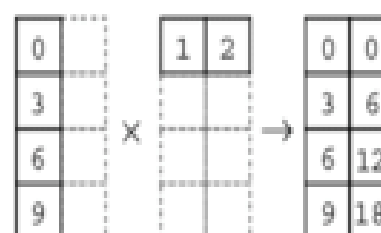
4	3
7	6

 with arrays with broadcasting

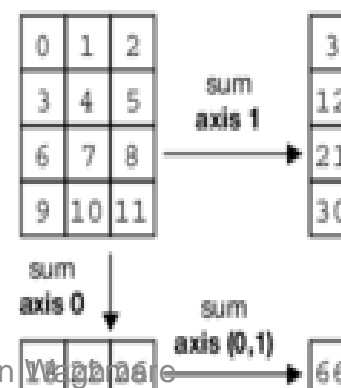
d Vectorization



e Broadcasting



f Reduction



g Example

```
In [1]: import numpy as np
```

```
In [2]: x = np.arange(12)
```

```
In [3]: x = x.reshape(4, 3)
```

```
In [4]: x
```

```
Out[4]:
```

```
array([[ 0,  1,  2],
       [ 3,  4,  5],
       [ 6,  7,  8],
       [ 9, 10, 11]])
```

```
In [5]: np.mean(x, axis=0)
```

```
Out[5]: array([4.5, 5.5, 6.5])
```

```
In [6]: x = x - np.mean(x, axis=0)
```

```
In [7]: x
```

```
Out[7]:
```

```
array([[ -4.5,  -4.5,  -4.5],
       [ -1.5,  -1.5,  -1.5],
       [  1.5,   1.5,   1.5],
       [  4.5,   4.5,   4.5]])
```

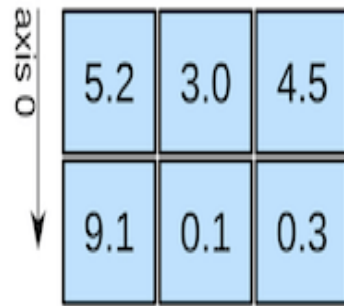
1D array



axis 0 →

shape: (4,)

2D array

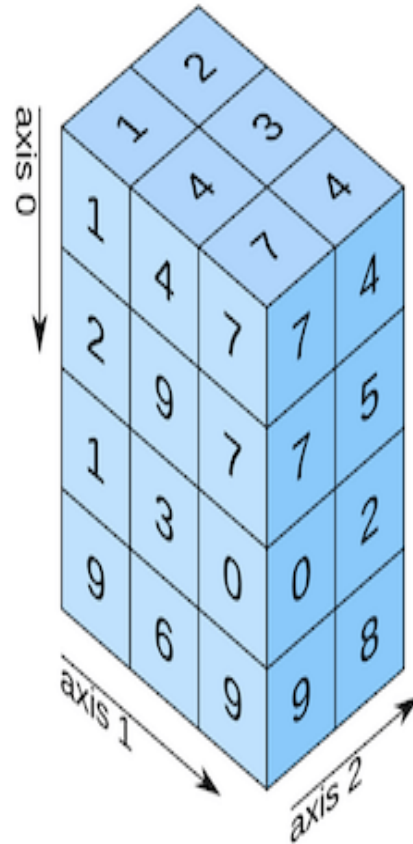


axis 0 ↓

axis 1 →

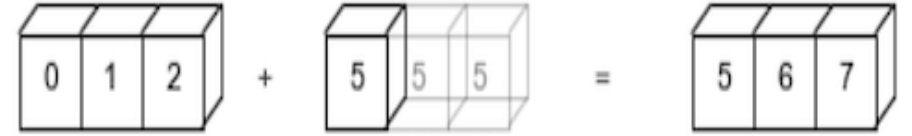
shape: (2, 3)

3D array

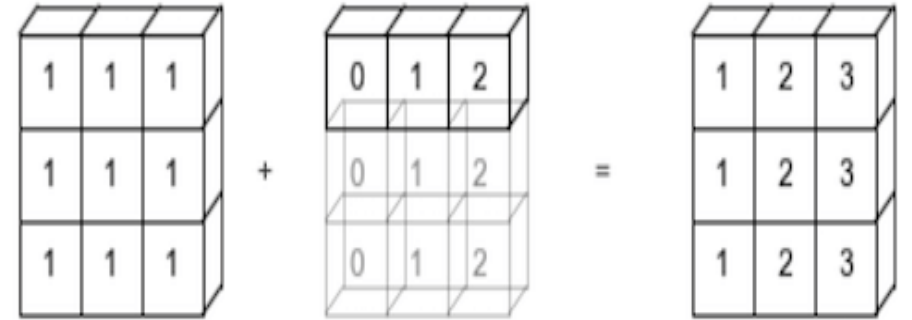


shape: (4, 3, 2)

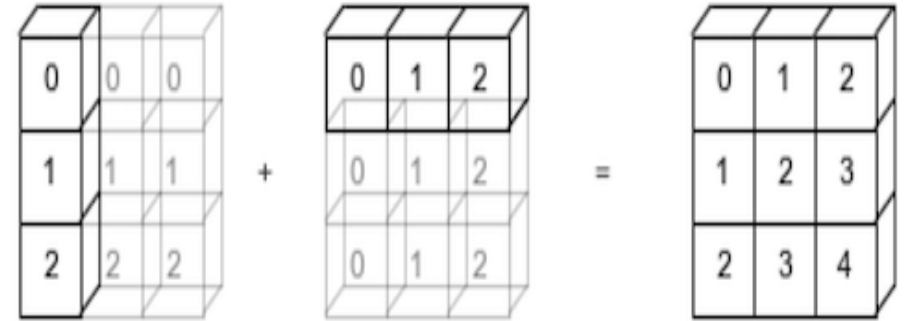
`np.arange(3)+5`



`np.ones((3, 3))+np.arange(3)`



`np.arange(3).reshape((3, 1))+np.arange(3)`



NumPy Basics

NumPy Basics

Operator	Description
<code>np.array([1,2,3])</code>	1d array
<code>np.array([(1,2,3),(4,5,6)])</code>	2d array
<code>np.arange(start,stop,step)</code>	range array

Placeholders

Operator	Description
<code>np.linspace(0,2,9)</code>	Add evenly spaced values btw interval to array of length
<code>np.zeros((1,2))</code>	Create and array filled with zeros
<code>np.ones((1,2))</code>	Creates an array filled with ones
<code>np.random.random((5,5))</code>	Creates random array
<code>np.empty((2,2))</code>	Creates an empty array

Array

Operator	Description
<code>array.shape</code>	Dimensions (Rows,Columns)
<code>len(array)</code>	Length of Array
<code>array.ndim</code>	Number of Array Dimensions
<code>array.dtype</code>	Data Type
<code>array.astype(type)</code>	Converts to Data Type
<code>type(array)</code>	Type of Array

Copying/Sorting

Operator	Description
<code>np.copy(array)</code>	Creates copy of array
<code>other = array.copy()</code>	Creates deep copy of array
<code>array.sort()</code>	Sorts an array
<code>array.sort(axis=0)</code>	Sorts axis of array

Array Manipulation

Adding or Removing Elements

Operator	Description
<code>np.append(a,b)</code>	Append items to array
<code>np.insert(array, 1, 2, axis)</code>	Insert items into array at axis 0 or 1
<code>np.resize((2,4))</code>	Resize array to shape(2,4)
<code>np.delete(array,1,axis)</code>	Deletes items from array

Combining Arrays

Operator	Description
<code>np.concatenate((a,b),axis=0)</code>	Split an array into multiple sub-arrays.
<code>np.vstack((a,b))</code>	Split an array in sub-arrays of (nearly) identical size
<code>np.hstack((a,b))</code>	Split the array horizontally at 3rd index

More

Operator	Description
<code>other = ndarray.flatten()</code>	Flattens a 2d array to 1d
<code>array = np.transpose(other)</code>	Transpose array
<code>array.T</code>	Transpose array
<code>inverse = np.linalg.inv(matrix)</code>	Inverse of a given matrix

Slicing and Subsetting

Operator	Description
<code>array[i]</code>	1d array at index i
<code>array[i,j]</code>	2d array at index[i][j]
<code>array[i<4]</code>	Boolean Indexing, see Tricks
<code>array[0:3]</code>	Select items of index 0, 1 and 2
<code>array[0:2,1]</code>	Select items of rows 0 and 1 at column 1
<code>array[:1]</code>	Select items of row 0 (equals <code>array[0:1, :]</code>)
<code>array[1:2, :]</code>	Select items of row 1
<code>[comment]: <> (</code>	<code>array[1,...]</code>
<code>array[: :-1]</code>	Reverses array

Mathematics

Operations

Operator	Description
<code>np.add(x,y)</code>	Addition
<code>np.subtract(x,y)</code>	Subtraction
<code>np.divide(x,y)</code>	Division
<code>np.multiply(x,y)</code>	Multiplication
<code>np.sqrt(x)</code>	Square Root
<code>np.sin(x)</code>	Element-wise sine
<code>np.cos(x)</code>	Element-wise cosine
<code>np.log(x)</code>	Element-wise natural log
<code>np.dot(x,y)</code>	Dot product
<code>np.roots([1,0,-4])</code>	Roots of a given polynomial coefficients

Comparison

Operator	Description
<code>==</code>	Equal
<code>!=</code>	Not equal
<code><</code>	Smaller than
<code>></code>	Greater than
<code><=</code>	Smaller than or equal
<code>>=</code>	Greater than or equal
<code>np.array_equal(x,y)</code>	Array-wise comparison

Basic Statistics

Operator	Description
<code>np.mean(array)</code>	Mean
<code>np.median(array)</code>	Median
<code>array.corrcoef()</code>	Correlation Coefficient
<code>np.std(array)</code>	Standard Deviation

More

Operator	Description
<code>array.sum()</code>	Array-wise sum
<code>array.min()</code>	Array-wise minimum value
<code>array.max(axis=0)</code>	Maximum value of specified axis
<code>array.cumsum(axis=0)</code>	Cumulative sum of specified axis

Pandas

- Provides **DataFrame** and **Series** structures for easy data handling
- Makes **data cleaning, filtering, and transformation** simple and fast
- Offers powerful **data analysis** tools (grouping, aggregation, summarization)
- Supports **reading and writing** data from multiple formats (CSV, Excel, JSON, SQL)
- Enables easy **handling of missing values**
- Integrates seamlessly with **NumPy, Matplotlib, Scikit-Learn**
- Allows **label-based indexing** and **intuitive selection** of rows/columns
- Great for **preprocessing data** before machine learning
- Used extensively for **EDA (Exploratory Data Analysis)**
- Highly optimized for performance and real-world data workloads

Pandas Library

Create Test Objects

Operator	Description
<code>pd.DataFrame(np.random.rand(20,5))</code>	5 columns and 20 rows of random floats
<code>pd.Series(my_list)</code>	Create a series from an iterable my_list
<code>df.index = pd.date_range('1900/1/30', periods=df.shape[0])</code>	Add a date index

Viewing/Inspecting Data

Operator	Description
<code>df.head(n)</code>	First n rows of the DataFrame
<code>df.tail(n)</code>	Last n rows of the DataFrame
<code>df.shape</code>	Number of rows and columns
<code>df.info()</code>	Index, Datatype and Memory information
<code>df.describe()</code>	Summary statistics for numerical columns
<code>s.value_counts(dropna=False)</code>	View unique values and counts
<code>df.apply(pd.Series.value_counts)</code>	Unique values and counts for all columns

Selection

Operator	Description
<code>df[col]</code>	Returns column with label col as Series
<code>df[[col1, col2]]</code>	Returns columns as a new DataFrame
<code>s.iloc[0]</code>	Selection by position
<code>s.loc['index_one']</code>	Selection by index
<code>df.iloc[0,:]</code>	First row
<code>df.iloc[0,0]</code>	First element of first column

Data Cleaning

Operator	Description
<code>df.columns = ['a','b','c']</code>	Rename columns
<code>pd.isnull()</code>	Checks for null Values, Returns Boolean Array
<code>pd.notnull()</code>	Opposite of <code>pd.isnull()</code>
<code>df.dropna()</code>	Drop all rows that contain null values
<code>df.dropna(axis=1)</code>	Drop all columns that contain null values
<code>df.dropna(axis=1,thresh=n)</code>	Drop all rows have have less than n non null values
<code>df.fillna(x)</code>	Replace all null values with x
<code>s.fillna(s.mean())</code>	Replace all null values with the mean
<code>s.astype(float)</code>	Convert the datatype of the series to float
<code>s.replace(1,'one')</code>	Replace all values equal to 1 with 'one'
<code>s.replace([2,3],['two', 'three'])</code>	Replace all 2 with 'two' and 3 with 'three'
<code>df.rename(columns=lambda x: x + 1)</code>	Mass renaming of columns
<code>df.rename(columns={'old_name': 'new_ name'})</code>	Selective renaming
<code>df.set_index('column_one')</code>	Change the index
<code>df.rename(index=lambda x: x + 1)</code>	Mass renaming of index

Filter, Sort, and Groupby

Operator	Description
<code>df[df[col] > 0.6]</code>	Rows where the column col is greater than 0.6
<code>df[(df[col] > 0.6) & (df[col] < 0.8)]</code>	Rows where $0.8 > \text{col} > 0.6$
<code>df.sort_values(col1)</code>	Sort values by col1 in ascending order
<code>df.sort_values(col2, ascending=False)</code>	Sort values by col2 in descending order.5
<code>df.sort_values([col1, col2], ascending=[True, False])</code>	Sort values by col1 in ascending order then col2 in descending order
<code>df.groupby(col)</code>	Returns a groupby object for values from one column
<code>df.groupby([col1, col2])</code>	Returns groupby object for values from multiple columns
<code>df.groupby(col1)[col2]</code>	Returns the mean of the values in col2, grouped by the values in col1
<code>df.pivot_table(index=col1, values=[col2, col3], aggfunc=mean)</code>	Create a pivot table that groups by col1 and calculates the mean of col2 and col3
<code>df.groupby(col1).agg(np.mean)</code>	Find the average across all columns for every unique col1 group
<code>df.apply(np.mean)</code>	Apply the function <code>np.mean()</code> across each column
<code>nf.apply(np.max, axis=1)</code>	Apply the function <code>np.max()</code> across each row

Join/Combine

Operator	Description
<code>df1.append(df2)</code>	Add the rows in df1 to the end of df2 (columns should be identical)
<code>pd.concat([df1, df2],axis=1)</code>	Add the columns in df1 to the end of df2 (rows should be identical)
<code>df1.join(df2,on=col1,how='inner')</code>	SQL-style join the columns in df1 with the columns on df2 where the rows for col have identical values. The 'how' can be 'left', 'right', 'outer' or 'inner'

Statistics

Operator	Description
<code>df.describe()</code>	Summary statistics for numerical columns
<code>df.mean()</code>	Returns the mean of all columns
<code>df.corr()</code>	Returns the correlation between columns in a DataFrame
<code>df.count()</code>	Returns the number of non-null values in each DataFrame column
<code>df.max()</code>	Returns the highest value in each column
<code>df.min()</code>	Returns the lowest value in each column
<code>df.median()</code>	Returns the median of each column
<code>df.std()</code>	Returns the standard deviation of each column