



# Advanced Data Science

**Introduction to Numpy, Pandas and  
Matplotlib**

**Session 2**

**Kiran Waghmare**

Program Manager  
C-DAC Mumbai

# Agenda

- **Numpy**
- **Pandas**
- **Matplotlib**

## **Unit II:**

### **Introduction to NumPy**

Brief Python Data Types;  
NumPy arrays;  
Data Types;  
Array Functions;  
Universal Functions;  
Aggregations;  
Broadcasting;  
Fancy Functions;  
Sorting arrays:  
Partial Sort;  
Structures arrays;

### **Data Manipulation using Pandas;**

Handling Missing Data;  
Hierarchical Indexing;

### **Visualization Using Matplotlib.**

Positive Indexing

0	1	2	3
---	---	---	---

One\_d

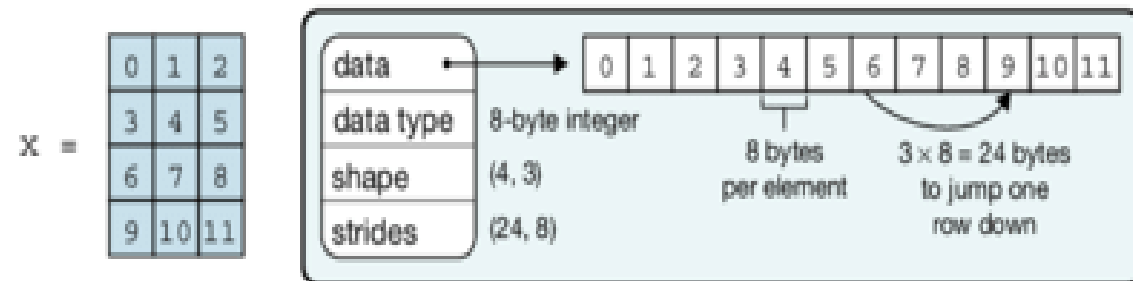
10	20	30	40
----	----	----	----

Negative Indexing

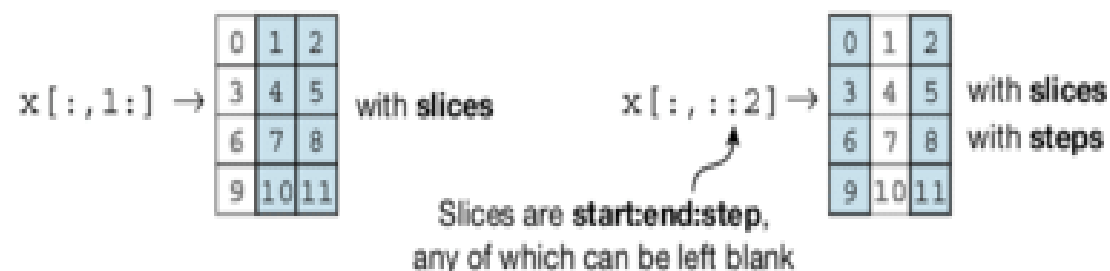
-4	-3	-2	-1
----	----	----	----



## a Data structure



## b Indexing (view)



## c Indexing (copy)

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

10	11
----	----

 with masks

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

1	5
---	---

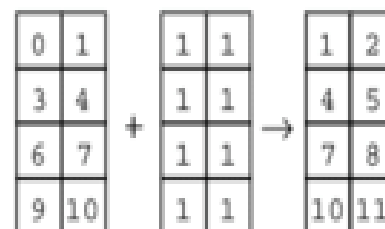
 with arrays

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

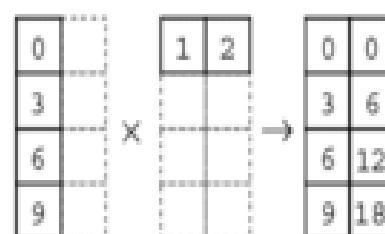
4	3
7	6

 with arrays with broadcasting

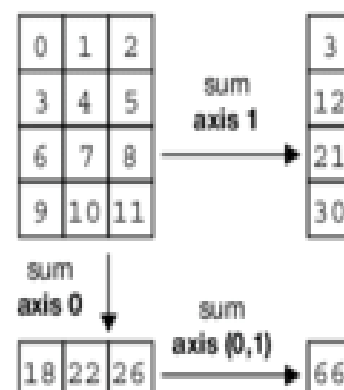
## d Vectorization



## 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]])
```

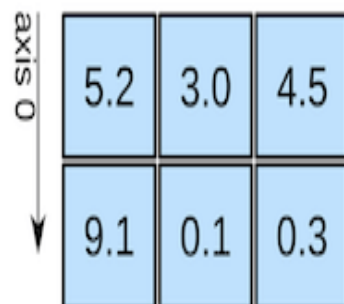
# 1D array



axis 0 →

shape: (4,)

# 2D array

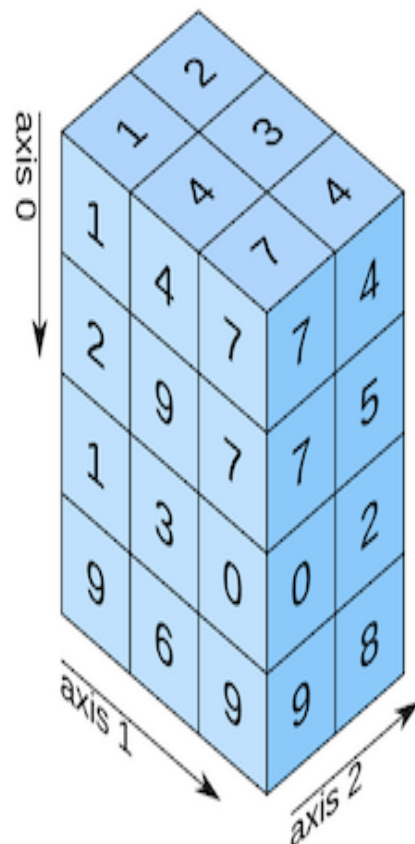


axis 0 ↓

axis 1 →

shape: (2, 3)

# 3D array



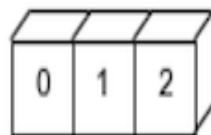
axis 0 ↓

axis 1 →

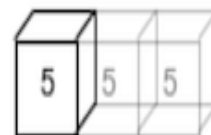
axis 2 ↗

shape: (4, 3, 2)

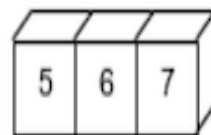
`np.arange(3)+5`



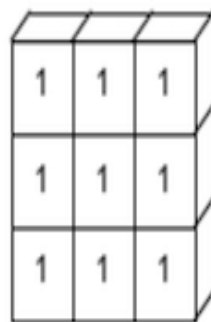
+



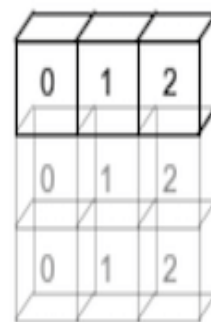
=



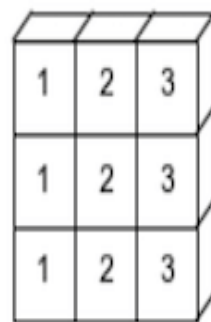
`np.ones((3,3))+np.arange(3)`



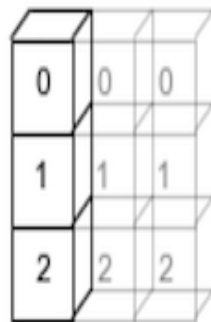
+



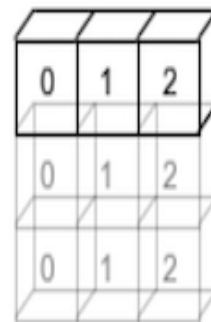
=



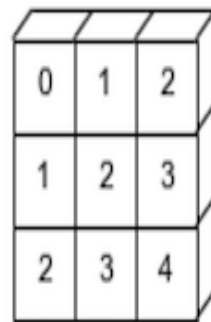
`np.arange(3).reshape((3,1))+np.arange(3)`



+



=



# 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&lt;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]: &lt;&gt; (</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>&lt;</code>	Smaller than
<code>&gt;</code>	Greater than
<code>&lt;=</code>	Smaller than or equal
<code>&gt;=</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

# 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] &gt; 0.6]</code>	Rows where the column col is greater than 0.6
<code>df[(df[col] &gt; 0.6) &amp; (df[col] &lt; 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

# Importing Data

Operator	Description
<code>pd.read_csv(filename)</code>	From a CSV file
<code>pd.read_table(filename)</code>	From a delimited text file (like TSV)
<code>pd.read_excel(filename)</code>	From an Excel file
<code>pd.read_sql(query, connection_object)</code>	Read from a SQL table/database
<code>pd.read_json(json_string)</code>	Read from a JSON formatted string, URL or file.
<code>pd.read_html(url)</code>	Parses an html URL, string or file and extracts tables to a list of dataframes
<code>pd.read_clipboard()</code>	Takes the contents of your clipboard and passes it to <code>read_table()</code>
<code>pd.DataFrame(dict)</code>	From a dict, keys for columns names, values for data as lists

# Exporting Data

Operator	Description
<code>df.to_csv(filename)</code>	Write to a CSV file
<code>df.to_excel(filename)</code>	Write to an Excel file
<code>df.to_sql(table_name, connection_object)</code>	Write to a SQL table
<code>df.to_json(filename)</code>	Write to a file in JSON format

- **Unit II:**
- **Introduction to NumPy**
- **Brief Python Data Types;**
- **NumPy arrays;**
- **Data Types;**
- **Array Functions;**
- **Universal Functions;**
- **Aggregations;**
- **Broadcasting;**
- **Fancy Functions;**
- **Sorting arrays:**
- **Partial Sort;**
- **Structures arrays;**
- **Data Manipulation using Pandas;**
- **Handling Missing Data;**
- **Hierarchical Indexing;**
- 
- **Visualization Using Matplotlib.**