# Data Science Workshop-1 ITER, SOA University

Centre for Data Science SOA University



1/62



## pandas

#### Introduction

Imagine you're an explorer in a vast digital world, where every piece of information you need is spread across countless islands. Some of these islands are small spreadsheets, others are massive databases, and some even contain messy, tangled webs of data. But you, the brave data explorer, have a secret tool: a magical panda.

#### Your Secret Tool

This panda is no ordinary bear; it has special powers that can help you wrangle data, make sense of large tables, and turn confusing information into clear, valuable insights. With just a little guidance, this magical panda will help you explore, clean, and analyze the data landscapes you encounter.

#### pandas

- pandas derived from 'panel data'. In econometrics this is a term for 'tabular data'
- The biggest difference is that pandas is designed for working with tabular or heterogeneous data.
- Basic data structures in pandas
  - series: One-dimensional array-like object containing a sequence of values (of similar types to NumPy types) of the same type and an associated array of data labels, called its index.
  - DataFrame: a two-dimensional data structure that holds data like a two-dimension array or a table with rows and columns.





#### Contents in Pandas





#### Contents

- Series
- Dictionary to pandas series and data frame
- Selecting rows and columns
- Re-indexing, dropping entries, filter
- Arithmetic and data alignment
- Summarizing and computing descriptive statistics
- Values counts, covariance and correlation





#### Series

```
import pandas as pd
x=pd.Series([1,2,3])
Х
     2
      3
dtype: int64
x.array
<PandasArray>
[1, 2, 3]
Length: 3, dtype: int64
x.index
RangeIndex(start=0, stop=3, step=1)
y = pd.Series([4, 7, -5], index=["x", 'y', 'z'])
У
Х
      4
У
     -5
dtype: int64
```

# Basics of indexing

- The index on the left and the values on the right.
- If we will not specify an index for the data, a default one consisting of the integers 0 through N - 1 (where N is the length of the data) is created.
- You can get the array representation and index object of the Series via its array and index attributes
- want to create a Series with an index identifying each data point with a label



# Selecting and filtering

- Like numpy arrays we can use labels in the index when selecting single values or a set of values.
- NumPy-like operations, such as filtering with a Boolean array, scalar multiplication, or applying math functions, will preserve the index-value link

## Dictionary to pandas series

- Series can be thought as a fixed-length, ordered dictionary, as it is a mapping of index values to data values.
- If we have data contained in a Python dictionary, you can create a Series from it by passing the dictionary to the series method
- A Series can be converted back to a dictionary with its to\_dict method

```
import pandas as pd
sdata = {"Adishree": 2241020023, "Sushree": 2241019629,
         "Anuska": 2241019579, "Omkar": 2241019335,
         'Goswami': 2241019331}
S = pd.Series(sdata)
Adishree
            2241020023
Sushree
            2241019629
Anuska
           2241019579
Omkar
            2241019335
Goswami
            2241019331
dtvpe: int64
S.to_dict()
{'Adishree': 2241020023,
 'Sushree': 2241019629,
 'Anuska': 2241019579,
 'Omkar': 2241019335.
 'Goswami': 2241019331}
```



## Dictionary to pandas series

- when you are only passing a dictionary, the index in the resulting Series will respect the order of the keys
- if we want to see them in a different order, that can be done by passing a list containing the dictionary keys in the desired order to the index keyword argument

```
L=['Rishita','Goswami','Sweta','Omkar','Anuska']
M=pd.Series(sdata,index=L)
M
```

```
Rishita NaN
Goswami 2.241019e+09
Sweta NaN
Omkar 2.241019e+09
Anuska 2.241020e+09
dtype: float64
```





# Missing values in pandas

- As no registration number for "Some students" was found, it appears as NaN (Not a Number), which is considered in pandas to mark missing or NA values
- isna and notna functions in pandas should be used to detect missing data

```
pd.isna(M)

Rishita True
Goswami False
Sweta True
Omkar False
Anuska False
dtype: bool
```

```
Rishita False
Goswami True
Sweta False
Omkar True
```

pd.notna(M)

Omkar True Anuska True dtype: bool



# Missing values in pandas

• Automatically aligns by index label in arithmetic operations:

```
data1={'Odisha':17,'Bengal':15}
data2={'Odisha':4,'Delhi':7}
d1=pd.Series(data1)
d2=pd.Series(data2)
d1+d2
```

```
Bengal NaN
Delhi NaN
Odisha 21.0
dtype: float64
```



#### Data frame creation

- A DataFrame represents a rectangular table of data and contains an ordered, named collection of columns, each of which can be a different data type
- The DataFrame has both a row and column index; it can be thought of as a dictionary of Series all sharing the same index

:	State	СМ	Pop
0	Delhi (NCT)	Atishi Marlena	4.37
1	Gujarat	Bhupendra Patel	10.43
2	Maharashtra	Eknath Shinde	13.10
3	Odisha	Mohan Charan Majhi	5.20
4	Uttarakhand	Pushkar Singh Dhami	3.20



#### head and tail method

- The head method selects only the first five rows
- The tail returns the last five rows.

3]:	frame.head()					<pre>frame.tail()</pre>			
3]:		State	СМ	Рор	[5]:		State	СМ	Рор
	0	Delhi (NCT)	Atishi Marlena	4.37		0	Delhi (NCT)	Atishi Marlena	4.37
	1	Gujarat	Bhupendra Patel	10.43		1	Gujarat	Bhupendra Patel	10.43
	2	Maharashtra	Eknath Shinde	13.10		2	Maharashtra	Eknath Shinde	13.10
	3	Odisha	Mohan Charan Majhi	5.20		3	Odisha	Mohan Charan Majhi	5.20
	4	Uttarakhand	Pushkar Singh Dhami	3.20		4	Uttarakhand	Pushkar Singh Dhami	3.20





## Naming the columns

- If you specify a sequence of columns, the DataFrame's columns will be arranged in that order
- If you pass a column that isn't contained in the dictionary, it will appear with missing values in the result
- data.columns gives us all the columns in the data

```
[7]: frame1=pd.DataFrame(data,columns=["CM","Pop","State","Loan_amt"])
frame1
```

	CM	Pop	State	Loan_amt
0	Atishi Marlena	4.37	Delhi (NCT)	NaN
1	Bhupendra Patel	10.43	Gujarat	NaN
2	Eknath Shinde	13.10	Maharashtra	NaN
3	Mohan Charan Majhi	5.20	Odisha	NaN
4	Pushkar Singh Dhami	3.20	Uttarakhand	NaN

```
[9]: frame1.columns
[9]: Index(['CM', 'Pop', 'State', 'Loan_amt'], dtype='object')
```



## Retrieving a column

 To retrieve a column, either use dictionary like syntax or column name as a attribute

```
[13]:
                                         frame1.CM
[11]: frame1['State']
                                                     Atishi Marlena
[11]: 0
         Delhi (NCT)
                                  [13]:
             Gujarat
                                                    Bhupendra Patel
          Maharashtra
                                                      Eknath Shinde
              Odisha
                                                Mohan Charan Majhi
          Uttarakhand
                                               Pushkar Singh Dhami
     Name: State, dtype: object
                                          Name: CM, dtype: object
```



# Updating a column

- Columns can be modified by assignment.
- The empty Governor column could be assigned a scalar value or an array of values

[17]:		ame1['Loan_amt']=23 ame1			
[17]:		СМ	Pop	State	Loan_amt
	o	Atishi Marlena	4.37	Delhi (NCT)	23
	1	Bhupendra Patel	10.43	Gujarat	23
	2	Eknath Shinde	13.10	Maharashtra	23
	3	Mohan Charan Majhi	5.20	Odisha	23
	4	Pushkar Singh Dhami	3.20	Uttarakhand	23
	•	r dankar anign briann	3.20	Ottarakriand	23
[19]:	fr	ame1['Loan_amt']=[3			
[19]: [19]:	fr	ame1['Loan_amt']=[3			
	fr	ame1['Loan_amt']=[3 ame1	.4,7.8	,6.7,9.8,5.3	1
	fr fr	ame1['Loan_amt']=[3 ame1	.4,7.8 Pop	,6.7,9.8,5.3 <b>State</b>	Loan_amt
	fr fr	ame1['Loan_amt']=[3 rame1 CM Atishi Marlena	.4,7.8 Pop 4.37	State  Delhi (NCT)	Loan_amt



Uttarakhand

3.20

Pushkar Singh Dhami

## Updating a data frame

- When you are assigning lists or arrays to a column, the value's length must match the length of the DataFrame
- If you assign a Series, its labels will be realigned exactly to the data frames index, inserting missing values in any index values not present

```
[31]: frame2=pd.DataFrame(data,columns=["CM","Pop","State","Loan_amt"])
  val=pd.Series([3.4,6.7,9.8],index=[0,1,3])
  frame2['Loan_amt']=val
  frame2
```

[31]:		СМ	Pop	State	Loan_amt
	0	Atishi Marlena	4.37	Delhi (NCT)	3.4
	1	Bhupendra Patel	10.43	Gujarat	6.7
	2	Eknath Shinde	13.10	Maharashtra	NaN
	3	Mohan Charan Majhi	5.20	Odisha	9.8
	4	Pushkar Singh Dhami	3.20	Uttarakhand	NaN





# creating new col and deleting

- Assigning a column that doesn't exist will create a new column.
- The del method can then be used to remove this column

45]:		ame1['Capital']='Bh ame1	ubanes	war'		
45]:		СМ	Pop	State	Loan_amt	Capital
	0	Atishi Marlena	4.37	Delhi (NCT)	3.4	Bhubaneswar
	1	Bhupendra Patel	10.43	Gujarat	7.8	Bhubaneswar
	2	Eknath Shinde	13.10	Maharashtra	6.7	Bhubaneswar
	3	Mohan Charan Majhi	5.20	Odisha	9.8	Bhubaneswar
	4	Pushkar Singh Dhami	3.20	Uttarakhand	5.3	Bhubaneswar
47]:	de	1 frame1['Capital']				
49]:	fr	ame1				
_	fr		Рор	State	Loan_amt	
49]: 49]:	fr 0	ame1		<b>State</b> Delhi (NCT)	Loan_amt	
_		ame1	Рор		_	
_	o	ame1  CM  Atishi Marlena	<b>Pop</b> 4.37	Delhi (NCT)	3.4	
_	0	came1 CM Atishi Marlena Bhupendra Patel	<b>Pop</b> 4.37 10.43	Delhi (NCT) Gujarat	3.4 7.8	





#### **Nested Dictionaries**

 If a nested dictionary is passed to the DataFrame, pandas will interpret the outer dictionary keys as the columns, and the inner keys as the row indices

	India	USA	Japan
2015	2.1	18.40	4.4
2020	2.7	21.06	5.0
2025	5.0	25.00	5.4





- Creating DataFrame from dict of lists, dict of dicts
- From 2-dim numpy arrays

```
#not more than 2-d
a=np.array([[1,2],[3,4],[5,6]])
a=pd.DataFrame(a,index=['r1','r2','r3'],columns=['c1','c2'])
#a.columns=['c1','c2']
#a.index=['r1','r2','r3']
a
```

```
r1 1 2
r2 3 4
r3 5 6
```



Creating Pandas DataFrame from lists of lists.

```
data = [['Virat', 35], ['Dhoni', 42], ['Shubman', 24]]
df = pd.DataFrame(data, columns=['Name', 'Age'])
df
```

	Name	Age
0	Virat	35
1	Dhoni	42
2	Shubman	24





Creating Dataframe from list of dicts

```
        a
        b
        c
        d

        0
        1
        2
        3
        NaN

        1
        10
        20
        30
        7.0
```



Creating DataFrame from Dictionary of series

	one	two
а	10	10
b	20	20
С	30	30
d	40	40





## Naming index and column

We can use name attribute for naming the index and columns

Country	India	USA
Year		
2020	3.8	13.5
2025	5.0	15.0



## Index objects in pandas

- Index is an immutable sequence (ordered multiset) used for indexing DataFrame and Series
- standard python indexing and slicing notation works for index objects
- create an index object using pd.Index() method
- NumPy attributes like shape, size, ndim and dtype works for index objects
- Immutable is the key difference between numpy array and indices
- many set operations like union, intersection, symmetric\_difference can be used for index





# Index objects in pandas

```
import pandas as pd
                                                       print(I.shape)
df=pd.DataFrame(['a','b','c'],index=[1,2,3])
T=df.index
                                                       print(I.size)
print(I[2])
                                                       print(I.dtype)
3
                                                       print(I.ndim)
labels = pd.Index(np.arange(1,4))
print('labels', labels)
pd.Series([1.5, -2.5, 0], index=labels)
                                                       (3,)
Int64Index([1, 2, 3], dtype='int64')
     1.5
1
                                                       int64
     -2.5
2
     0.0
dtype: float64
I1=pd.Index([1,2,3,4,5,6])
I2=pd.Index([3,4,5,6,7,8])
print(I1.union(I2))
print(I1.intersection(I2))
print(I1.symmetric difference(I2))
print(I1.append(I2))
Int64Index([1, 2, 3, 4, 5, 6, 7, 8], dtype='int64')
Int64Index([3, 4, 5, 6], dtype='int64')
Int64Index([1, 2, 7, 8], dtype='int64')
```



Int64Index([1, 2, 3, 4, 5, 6, 3, 4, 5, 6, 7, 8], dtype='int64')

**∢** ∄ ▶

# Reindexing

- reindex is used to create a new object with the values rearranged
- DataFrame.reindex(labels,index,columns, axis, method, copy, level, fill\_value, limit, tolerance)
- labels and index works same: New index labels

```
import pandas as pd
import numpy as no
frame = pd.DataFrame(np.arange(9).reshape((3, 3)),
                     index=["i", "ii", "iii"],columns=["A", "B", "C"])
frame
   A B C
frame2 = frame.reindex(["i", "iii", "ii"])
frame?
   A B C
frame2 = frame.reindex(index=["i", "iii", "ii"])
frame2
   A B C
```



ii 3 4 5

# Reindexing

- columns: labels for the columns
- axis: Axis to target. Can be either the axis name ('index', 'columns') or number (0, 1)
- fill\_value: Value to use for missing value

```
frame2 = frame.reindex(["ii","A","D", "i"],axis=0)
frame2 = frame.reindex(columns=["A", "B"])
                                                         frame?
frame2
    A B
                                                           i 0.0
                                                                1.0 2.0
                                                         frame2 = frame.reindex(["ii","A","D", "i"],axis=0,fill value=7)
frame2 = frame.reindex(["C","A","D", "B"],axis=1)
                                                         frame2
frame2
                                                             A B C
                                                          ii 3 4 5
                                                          D 7 7 7
```

i 0 1 2

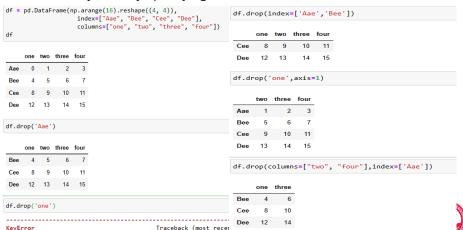
# Reindexing

- method: use for filling holes in reindexed DataFrame.
- Note: this is only applicable to DataFrames/Series with a monotonically increasing/decreasing index.
- None (default): don't fill gaps, ffill: Propagate last valid observation forward to next valid, bfill: Use next valid observation to fill gap

```
frame1 = pd.Series(["A", "B", "C"],
                                       frame1.reindex(np.arange(6),
                   index=[0, 2, 4])
                                                       method="ffill")
frame1
                                            А
dtype: object
                                       dtype: object
frame1.reindex(np.arange(6))
                                       frame1.reindex(np.arange(6),
                                                       method="bfill")
     NaN
                                               А
     NaN
     NaN
                                            NaN
```

## Dropping entries from an Axis

 Remove rows or columns by specifying label names and corresponding axis, or by directly specifying index or column names



#### Indexing, selection, filter in Series

 Series indexing works analogously to NumPy array indexing, except you can use the Series's index values instead of only integers

```
print(obj[["b", "a", "d"]])
obj = pd.Series(np.arange(4.),
             index=["a", "b", "c", "d"l)
                                              1.0
                                       b
obi
                                       a 0.0
    0.0
                                              3.0
   1.0
                                       dtype: float64
    2.0
    3.0
dtvpe: float64
                                       print(obj[[1, 3]])
print('obj["b"]=',obj["b"])
                                             1.0
print('obj[1]=',obj[1])
                                              3.0
                                       dtype: float64
obi["b"]= 1.0
obi[1]= 1.0
                                       print(obj[obj < 2])</pre>
print(obj[2:4])
                                              0.0
    2.0
                                              1.0
    3.0
                                       dtype: float64
dtvpe: float64
```

#### Indexing in DataFrame

- Indexing into a DataFrame retrieves one or more columns either with a single value or sequence
- it is more preferable to use loc and iloc methods

```
data['c1']

r1 0
r2 4
r3 8
r4 12
Name: c1, dtype: int32
```

```
    c1
    c2
    c3
    c4

    r1
    0
    1
    2
    3

    r2
    4
    5
    6
    7

    r3
    8
    9
    10
    11

    r4
    12
    13
    14
    15
```

```
c1 c3

r1 0 2

r2 4 6

r3 8 10
```



# Indexing in DataFrame

ata[:	2]			
c1	c2	c3	c4	
<b>r1</b> 0	1	2	3	
_				
r <b>2</b> 4	5	6	1	
ata>5				
	с1	c2	с3	c4
r1 Fa	lse	False	False	False
	lse	False	True	True
Fa				



True

True

True

True

## loc and iloc in pandas

- loc and iloc for label-based and integer-based indexing
- loc[] is primarily label based, but may also be used with a boolean array.
- Single label. Note this returns the row as a Series
- List of labels. Note using [[]] returns a DataFrame

```
import pandas as pd
df = pd.DataFrame([[1, 2], [4, 5], [7, 8]],
        index=['Cobra', 'Viper', 'RatSnake'],
        columns=['max_speed', 'length'])
df
```

Name:	Viper,	dtyp	e: int	64
df.loc	[['Vipe	er','	Cobra'	]]
	max_sp	eed l	ength	
Viper		4	5	
Cobra		1	2	

df.loc['Viper']

max\_speed

	max_speed	length
Cobra	1	2
Viper	4	5
RatSnake	7	8

## loc and iloc in pandas

- Slice with labels for row and single label for column.
- Boolean list with the same length as the row axis
- Conditional that returns a boolean Series
- Conditional that returns a boolean Series with column labels specified

```
df.loc['Cobra':'Viper', 'max_speed']
                                      df.loc[df['length'] > 6]
Cohra
Viper
                                                max_speed length
Name: max speed, dtype: int64
                                       RatSnake
df.loc[[False, False, True]]
                                      df.loc[df['length']>6,'max speed']
         max speed length
                                      RatSnake
RatSnake
                                      Name: max speed, dtype: int64
```

Data Science Workshop-1

### loc and iloc in pandas

• iloc: integer-location based indexing for selection by position

```
c1 c2 c3 c4

r1 0 1 2 3

r2 4 5 6 7

r3 8 9 10 11

r4 12 13 14 15
```



2024

### loc and iloc in pandas

```
data.iloc[[1,3]]
                               data.iloc[[True,True,False,False]]
     c1
         c2 c3 c4
                                  c1 c2 c3 c4
  r2
     12
         13
                 15
                                r2
                                      5
                               data.iloc[[True,True,False,False],
 data.iloc[0:2,1:3]
                                       [True, True, False, False]]
     c2
         c3
                                  c1 c2
  r1
  r2
      5
          6
data.at['r2','c3'] data.iat[1,2]
                                6
6
```



## Indexing in dataframe

Table 5.4: Indexing options with DataFrame

Туре	Notes
df[column]	Select single column or sequence of columns from the DataFrame; special case conveniences: Boolean array (filter rows), slice (slice rows), or Boolean DataFrame (set values based on some criterion)
df.loc[rows]	Select single row or subset of rows from the DataFrame by label
df.loc[:, cols]	Select single column or subset of columns by label
df.loc[rows,	Select both row(s) and column(s) by label
df.iloc[rows]	Select single row or subset of rows from the DataFrame by integer position
df.iloc[:,	Select single column or subset of columns by integer position
df.iloc[rows,	Select both row(s) and column(s) by integer position
df.at[row, col]	Select a single scalar value by row and column label
<pre>df.iat[row, col]</pre>	Select a single scalar value by row and column position (integers)

reindex method Select either rows or columns by labels

#### Experiement loc and iloc methods on series object



39 / 62

8

### Integer indexing pitfalls

#### Let us consider an example

·	
<pre>ser = pd.Series(np.arange(3.)) ser</pre>	<pre>ser = pd.Series(np.arange(3.),index=['a','b','c']) ser</pre>
0 0.0 1 1.0 2 2.0 dtype: float64	a 0.0 b 1.0 c 2.0 dtype: float64
ser.iloc[-1]	ser[-1]
2.0	2.0
ser[-1]	ser.iloc[-1]
ValueError	7 2.0

If you have an axis index containing integers, data selection will always be label oriented.

2024

## NumPy ufuncs

 NumPy ufuncs (element-wise array methods) also work with pandas objects

```
frame = pd.DataFrame(np.random.standard normal((4, 3)),
                                                       np.abs(frame)
                  columns=list("bde"),
                  index=["India", "USA", "Japan", "China"])
frame
                                                          India
                                                                 0.680518
                                                                              0.592259
                                                                                          0.446467
 India 0.680518 0.592259 0.446467
                                                          USA
                                                                 0 478302 0 529713
                                                                                          1 430172
  USA -0.478302 0.529713 -1.430172
                                                                  1.246007 0.617735
                                                                                          0.864200
Japan 1.246007 0.617735 -0.864200
                                                        Japan
 China -1 488501 -0 085932 2 205331
                                                                 1.488501
                                                                              0.085932 2.205331
                                                        China
```

# Function Application and mapping

- Apply a function along an axis of the DataFrame
- either the DataFrame's index (axis=0) or the DataFrame's columns (axis=1)

```
def f1(x):
df = pd.DataFrame([[4, 9]] * 3,
                                df.apply(np.sum, axis=0)
                                                                      return x.max() - x.min()
               columns=['A', 'B'])
                                                                  df.apply(f1, axis=1)
                                      12
df
                                      27
                                dtype: int64
                                                                  dtype: int64
  A B
                                df.apply(np.sum, axis=1)
                                                                  def f1(x):
                                                                      return x.max() - x.min()
                                      13
                                                                  df.apply(f1, axis=0)
                                      13
                                      13
                                                                  dtvpe: int64
                                dtype: int64
```

• create a 4\*4 matrix, where each row contains Student's name, registration no, sec and cgpa.

```
D={'Name':['A',"B","C","D"],'Reg_no':[1,2,3,4],
'Sec':['i','iii','iii','iv'],'CGPA':[7,8,9,10]} Find the name, registration no
of the student with highest cgpa. (Use np.max on cgpa column) What
is the row and column number of the cell with the highest CGPA?
(Convert dataframe to numpy array and see which indices has
maximum value)
```



• create a 4\*4 matrix, where each row contains Student's name, registration no, sec and cgpa.

```
D={'Name':['A',"B","C","D"],'Reg_no':[1,2,3,4], 'Sec':['i','iii','iii','iv'],'CGPA':[7,8,9,10]} Find the name, registration no of the student with highest cgpa.(Use np.max on cgpa column) What is the row and column number of the cell with the highest CGPA? (Convert dataframe to numpy array and see which indices has maximum value)
```

**Ans**  $df.loc[df['CGPA'] == np.max(df.CGPA),['Name','Reg_no']]$ row,  $col = np.where(df.to_numpy() == np.max(df.CGPA))$ 



43 / 62

• create a 4\*4 matrix, where each row contains Student's name, registration no, sec and cgpa.

```
D = \{\text{'Name':} [\text{'A',"B'',"C'',"D''}], \text{'Reg\_no':} [1,2,3,4], \\ \text{'Sec':} [\text{'i','iii','iii','iv'}], \text{'CGPA':} [7,8,9,10] \} \text{ Find the name, registration no of the student with highest cgpa.} (Use np.max on cgpa column) What is the row and column number of the cell with the highest CGPA?} (Convert dataframe to numpy array and see which indices has maximum value)}
```

**Ans**  $df.loc[df['CGPA'] == np.max(df.CGPA),['Name','Reg_no']]$ row,  $col = np.where(df.to_numpy() == np.max(df.CGPA))$ 

Ocunt the number of missing values in each column of df. Which column has the maximum number of missing values?(create a function to count nan values in a column, use df.apply and that function)



• create a 4\*4 matrix, where each row contains Student's name, registration no, sec and cgpa.

```
D = \{\text{'Name':} [\text{'A',"B'',"C'',"D''}], \text{'Reg\_no':} [1,2,3,4], \\ \text{'Sec':} [\text{'i','iii','iii','iv'}], \text{'CGPA':} [7,8,9,10] \} \text{ Find the name, registration no of the student with highest cgpa.} (Use np.max on cgpa column) What is the row and column number of the cell with the highest CGPA?} (Convert dataframe to numpy array and see which indices has maximum value)}
```

• Count the number of missing values in each column of df. Which column has the maximum number of missing values?(create a function to count nan values in a column, use df.apply and that function) Ans missings\_each\_col = df.apply(lambda x: x.isnull().sum()) missings\_each\_col.argmax()

Data Science Workshop-1

• Create a 4\*4 matrix, entries from np.arange(16). Reverse the entries

of the matrix. The matrix will look like 
$$\begin{bmatrix} 15 & 14 & 13 & 12 \\ 11 & 10 & 9 & 8 \\ 7 & 6 & 5 & 4 \\ 3 & 2 & 1 & 0 \end{bmatrix}$$



44 / 62

• Create a 4\*4 matrix, entries from np.arange(16). Reverse the entries

of the matrix. The matrix will look like  $\begin{bmatrix} 15 & 14 & 13 & 12 \\ 11 & 10 & 9 & 8 \\ 7 & 6 & 5 & 4 \\ 3 & 2 & 1 & 0 \end{bmatrix}$ 

```
x=pd.DataFrame(np.arange(16).reshape((4,4)))
x.loc[::-1,::-1]
```

Oreate a 4\*4 matrix where each entries are from standard normal distribution. Add a column to this data frame, which contains row wise minimums.



• Create a 4\*4 matrix, entries from np.arange(16). Reverse the entries

```
of the matrix. The matrix will look like \begin{bmatrix} 15 & 14 & 13 & 12 \\ 11 & 10 & 9 & 8 \\ 7 & 6 & 5 & 4 \\ 3 & 2 & 1 & 0 \end{bmatrix}
```

```
x=pd.DataFrame(np.arange(16).reshape((4,4)))
x.loc[::-1,::-1]
```

Create a 4\*4 matrix where each entries are from standard normal distribution. Add a column to this data frame, which contains row wise minimums.

```
df=pd.DataFrame(np.random.randn(16).reshape((4,4)))
print(df)
row_max=df.apply(np.max,axis=1)
df['row_max']=row_max
df
```





Create a 4\*4 matrix entries from standard normal distribution. Add another row where all entries are NaN. Replace missing values in 'col1' and 'col2' columns with their respective mean.





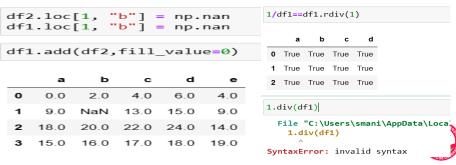
## Arithmetic and Data Alignment

- Arithmetic operation between objects that have different indexes, if any index pairs are not the same, the respective index in the result will be the union of the index pairs.
- Using the add method on df1, I pass df2 and an argument to fill\_value, which substitutes the passed value for any missing values in the operation

f)	-											а	ь	c	d	e
	а	ь	c	d							0	0.0	2.0	4.0	6.0	NaN
o	О	1	2	3							1	9.0	11.0	13.0	15.0	NaN
1	4	5	6	7							2	18.0	20.0	22.0	24.0	NaN
2	8	9	10	11												
1f2		pd	. Da	taF	ran			2(20).re: ist("ab	shape((4, cde"))	5)),	<b>3</b>	NaN L.add	NaN (df2,	NaN fill	NaN _valu	
1152	2				ran					5)),	_					
1152		Ŀ		c	ran d		. co			5)),	_	L.add	(df2,	fill_	_valu	e=0)
 	a	E		e 2	d	•	- - 4			5)),	dfi	L.add a	(df2,	fill <sub>.</sub>	_valu	e=0)
0	2 a	1	) (	e 2 7	<b>d</b> 3 8	<b>e</b> 4	+ 1			5)),	df1	L.add a 0.0	(df2,	fill_c	_valu d 6.0 15.0	e=0) e

### Arithmetic and Data Alignment

- If data in both corresponding DataFrame locations is missing the result will be missing.
- df1.add(df2) is equivalent to df2.radd(df1), r stands for reverse, the arguments are reversed
- Experiment with sub, rsub, div, rdiv, floordiv, rfloordiv, mul, rmul, pow, rpow



## Sorting

- To sort lexicographically by row or column label, use the sort\_index method, which returns a new, sorted object
- you can sort by rows or columns by specifying axis key word

```
frame.sort_index(axis='index')

d a b c

one 4 5 6 7

three 0 1 2 3
```

frame.sort\_index(axis='columns')

d a b c three 0 1 2 3 one 4 5 6 7

a b c d
three 1 2 3 0



2024

## Sorting

- The data is sorted in ascending order by default but can be sorted in descending order by the ascending key word
- sort\_values is used to sort a Series and DataFrame
- Any missing values are sorted to the end of the Series by default
- Missing values can be sorted to the start instead by using the na\_position argument to be first

```
frame.sort_index(axis="columns", ascending=False)
                                    obj.sort values()
    d c b a
three 0 3 2 1
                                    dtype: int64
obj = pd.Series([4, 7, -3, 2])
                                    obj.sort values(na_position="first")
obj
                                        NaN
                                        NaN
                                       -3.0
dtype: int64
                                        2.0
```

## Sorting

- While sorting a DataFrame, we can use the column name as the criteria to sort
- either pass a column name in the parentheses or multiple names in a column in side the parentheses

```
frame = pd.DataFrame(
                          frame.sort values("b") frame.sort_values(["a", "b"])
    {"a": [-4, 7, 3, 2],
     "b": [0, 1, 0, 1]})
frame
                              a b
    a b
                           3 2 1
```

## ranking

- Compute numerical data ranks 1 through n(the length)
- The smallest gets rank one and largest gets rank n
- Equal values are assigned a rank that is the average of the ranks of those values



### ranking

- If same values appear more than once, Ranks can also be assigned according to the order in which they're observed in the data
- we can also rank in descending order
- ranks for DataFrame can be computed either for rows or columns

```
obj=pd.Series([3,2,3,1])
                                                            frame.rank(axis='index')
                             : frame = pd.DataFrame(
obj.rank(method='first')
                                   {"b": [4.3, 7, -3, 2],
                                    "a": [0, 1, 0, 1],
0
     3.0
                                                                3.0 1.5 2.0
                                    "c": [-2, 5, 8, -2.5]})
     2.0
     4.0
                               frame
     1.0
                                                               2.0 3.5 1.0
dtype: float64
                                                            frame.rank(axis='columns')
obj.rank(ascending=False)
                                  4.3 0 -2.0
     1.5
     3.0
                                                                3.0 2.0 1.0
                                2 -3.0 0 8.0
     1.5
     4.0
                                                                1.0
                                                                   2.0 3.0
                                  2.0 1 -2.5
dtype: float64
                                                                30 20 10
```

## Summary statistics

 numpy methods like sum(), mean(), cumsum(), can be applied by specifying axis labels

```
df.sum(axis=0)
: df=pd.DataFrame(
                                                                   df.cumsum()
     np.random.randn(9).reshape((3,3)))
                                           3.053371
  df
                                           1.975640
                                                                               0
                                                                                          1
                                         -0.550882
                                    dtype: float64
                                                                       0.893613
                                                                                  0.684894
                                                                                             -0.7194
                                                                       1.143361
                                                                                  1.826619
                                                                                              1.2067
                                    df.sum(axis=1)
   0 0.893613 0.684894 -0.719420
                                                                       3.053371
                                                                                  1.975640
                                                                                             -0.5508
   1 0.249749 1.141725 1.926191
                                          0.859087
                                           3.317664
                                    1
  2 1.910010 0.149021 -1.757654
                                           0.301377
                                                                   df.cumsum(axis=1)
                                    dtype: float64
: df.sum()
                                                                               O
                                                                                          1
                                    df.mean(axis=1)
                                                                       0.893613
                                                                                             0.85908
                                                                                  1 578506
      3.053371
                                           0.286362
      1.975640
                                           1.105888
                                    1
                                                                       0.249749
                                                                                  1.391473
                                                                                             3.31766
     -0.550882
                                           0.100459
                                                                       1.910010
                                                                                  2.059031
                                                                                             0.30137
  dtype: float64
                                    dtvpe: float64
```



2024

#### statistics

 idxmin and idxmax, return the index value where the minimum or maximum values are attained

```
import pandas as pd
import numpy as np
obj=pd.DataFrame(np.random.
                  randn(30).reshape((5,6)))
obj
0 -0.528642 0.438604 -1.850550 -0.465469 -1.754782 -0.878637
   -0.992590 -0.253095 -1.229630
                                1.249513 -0.929660
    0.383950 -0.677413 0.096648 -0.521331 -1.078097
                                                    0.383043
    2 341758 -1 092688
                       0.132866
                                 0.352772
                                          0.946163
                                                    0.990012
4 -0.589614 -0.436097 1.995528
                                3.191171
```

```
obj.idxmax()
dtype: int64
obj.idxmax(axis=1)
a
dtype: int64
obj.idxmax(axis=0)
```

#### Describe

obj.describe()

	0	1	2	3	4	5
count	5.000000	5.000000	5.000000	5.000000	5.000000	5.000000
mean	0.122972	-0.404138	-0.171028	0.761331	-0.189888	-0.038386
std	1.338654	0.566310	1.483189	1.537825	1.525303	0.734693
min	-0.992590	-1.092688	-1.850550	-0.521331	-1.754782	-0.878637
25%	-0.589614	-0.677413	-1.229630	-0.465469	-1.078097	-0.469812
50%	-0.528642	-0.436097	0.096648	0.352772	-0.929660	-0.216534
75%	0.383950	-0.253095	0.132866	1.249513	0.946163	0.383043
max	2.341758	0.438604	1.995528	3.191171	1.866934	0.990012

#### Exercise

Consider a column containing non-numeric data and then apply describe method to that. a) first apply to a non-numeric series b)a data frame with one column non-numeric

# Summary statisitcs and other methods

count	Number of non-NA values
describe	Compute set of summary statistics
min, max	Compute minimum and maximum values
argmin, argmax	Compute index locations (integers) at which minimum or maximum value is obtained, respectively; not available on DataFrame objects
id×min, id×ma×	Compute index labels at which minimum or maximum value is obtained, respectively
quantile	Compute sample quantile ranging from 0 to 1 (default: 0.5)
sum	Sum of values
mean	Mean of values
median	Arithmetic median (50% quantile) of values
mad	Mean absolute deviation from mean value
prod	Product of all values
var	Sample variance of values
std	Sample standard deviation of values
skew	Sample skewness (third moment) of values
kurt	Sample kurtosis (fourth moment) of values
cumsum	Cumulative sum of values
cummin,	Cumulative minimum or maximum of values, respectively
cumprod	Cumulative product of values

### Covariance and correlation

single observed value of dependent variable mean of all values of independent variable

single observed value of independent variable

total count of property sample values

$$\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})$$

Cov(X,Y) =

mean of all values of independent variable

n-1

ш

population count minus one (Bessel's Correction)

 $\alpha$ 

$$Cor(X,Y) = \frac{Cov(X,Y)}{2}$$



$$Cor(X,Y) = rac{\sum (x_i - ar{x})(y_i - ar{y})}{\sqrt{\sum (x_i - ar{x})^2 (y_i - ar{y})^2}}$$
 standard

(SOA)

Data Science Workshop-1

deviations

2024

# Syntax

- series\_name1.cov(series\_name2)
- series\_name1.corr(series\_name2)
- df.cov() for covariance and df.corr() for corelation among columns of the data frame.



### unique values, counts

- function is unique, which gives you an array of the unique values in a Series
- value\_counts computes a Series containing value frequencies
- isin performs a vectorized set membership check



### Covariance

obj[0].cov(obj[0]) 0.8261131446271064

obj.cov()

	0	1	2	3	4	5
0	0.826113	0 101669	0.610734	1 129590	-0 377262	0.070390
_	0.020110					
1	0.101669	0.770697	0.475538	1.436232	-0.082309	0.466032
2	-0.610734	0.475538	1.565414	0.505303	0.374347	1.454724
3	1.128580	1.436232	0.505303	5.724417	0.946254	1.015754
4	-0.377262	-0.082309	0.374347	0.946254	1.410920	1.457691
5	-0.979280	0.466032	1.454724	1.015754	1.457691	2.544626



### Correlation

```
obj[0].corr(obj[0])
```

1.0

obj.corr()

	0	1	2	3	4	5
0	1.000000	0.127417	-0.537054	0.518976	-0.349440	-0.675422
1	0.127417	1.000000	0.432941	0.683781	-0.078932	0.332784
2	-0.537054	0.432941	1.000000	0.168800	0.251889	0.728877
3	0.518976	0.683781	0.168800	1.000000	0.332959	0.266141
4	-0.349440	-0.078932	0.251889	0.332959	1.000000	0.769312
5	-0.675422	0.332784	0.728877	0.266141	0.769312	1.000000



### isin function

```
obj = pd.Series(["c", "a", "d", "a", "a", "b", "b", "c", "c"])
uniques = obj.unique()
uniques
array(['c', 'a', 'd', 'b'], dtype=object)
obj.value_counts()
ä
Б
a
dtype: int64
mask=obj.isin(["b", "c"])
mask
0
       True
      False
ż
      False
      False
      False
       True
       True
ń
       True
dtype: bool
obj[mask]
      ъ
ĕ
      Б
```

