

July 12

① csv - separated by commas (,)
② tsv - " " tab (---)

Python-tutorial 10

Pandas : (Panel Data System, a library).

• file Read Write operation :

import pandas as pd

draft1 = pd.read_csv(--- address ---)

draft1 / or draft1.head()

tabular form

	a	b	c	d	message
0	1	2	3	4	hello
1	5	6	7	8	world
2	9	10	11	12	foo

by default first line as heading (row & column)

type(draft1)

pandas.core.frame.DataFrame

head() heading + top-5 data sets

tail() " + bottom-5 data sets

head(8) " + top-8 data sets

tail(18) " + bottom 18 data sets

draft2 = pd.read_table('test.tsv')

draft2.head()

test test1 test2 test3 test4

seeth seeth seeth seeth seeth seeth

kumar kumar kumar kumar kumar kumar

(*) if you want to read csv file in tsv format:

```
draft1 = pd.read_csv('test.tsv')
```

```
draft1.tail()
```

(:: no space) ↓

test\ttest1\ttest2\ttest3\ttest4

0

seedh\tseedh\tseedh\tseedh\tseedh\tseedh

1

kumar\tkumar\tkumar\tkumar\tkumar\tkumar

but if we do like ↓:

```
draft1 = pd.read_csv('test.tsv', sep='\\t')
```

```
draft1.tail()
```

test test1 test2 test3 test4

seedh seedh seedh seedh seedh seedh

kumar kumar kumar kumar kumar kumar

Similarly, for tsv in csv (:: sep = ',')

(*) if you don't want to consider first row as header instead it will be part of data
(:: header = None) → & by default heading is [0, 1, 2, 3, 4, 5]

(*) if you want to give your own heading:

```
df = pd.read_csv('ex2.csv',
```

```
names = ['M', 'A', 'D', 'H', 'U',  
         'P', 'I'])
```

df

(if there are extra column then NaN)

	M	A	D	H	U	P	I
0	a	b	c	d	e	NaN	NaN
1	f	g	h	i	j	NaN	NaN
2	k	l	m	n	o	NaN	NaN

(*) Multilevel indexing (when you want to store data inside data).
hierarchical indexing

`parsed = pd.read_csv('csv_mindex.csv')`

parsed

	key1	key2	value1	value2
(default index)	0	one	a	1
	1	one	b	3
	2	one	c	5
	3	two	a	7
	4	two	b	9
	5	two	c	11

`parsed = pd.read_csv('csv_mindex.csv',
index_col = ['key1', 'key2'])`

	key1	key2	value1	value2
		a	1	2
	one	b	3	4
		c	5	6
		d	7	8
	two	e	9	10
		f		

(because key1 & key2 became a column in index and default index is removed.)

③ xlsx - excel

(*) if you want to read some data & skip some:
(\therefore skip rows = [0, 2, 3])

(*) if there is some null data in our data set & you want to know (True/False):

```
result = pd.read_csv('ex5.csv')
```

result True/False

```
pd.isnull(result)
```

a b c d e

0 False False False False True null-data

1 False False True False False

2 False False False False False

(*) Excel - Sheets :-

```
draft3 = pd.read_excel(...xlsx)
```

or draft3 = pd.read_excel(...xlsx,

(for particular sheet) Sheet_name = 'Sheet2')

(by default first sheet of excel)

(*) accessing data from any website (in tabular form)
PANDAS-HTML won't be able to access data present in other than tabular form:

```
url = "http://www. .... .html"
```

```
BB_data = pd.read_html(url)
```

BB_data

[...] = [dataframe1, df2, ...]

type(BB_data) → list

multiple tables present on website are merged in a list

data-frame \rightarrow collection of rows & columns.
($m \times 1$), ($1 \times n$) \rightarrow series

BB_data[0] \rightarrow First dataframe of list

df = BB_data[0]

df.head(10) \rightarrow first 10 rows

[df['column1 say'] \rightarrow extracting data from required columns
(Type \downarrow Series)

[df[['col1', 'col2', 'col15']] \rightarrow multiple-columns
(Type \downarrow list) (: pass as list)

* when you want to know all column names:

df.columns (: heading)

\rightarrow Index(['M', 'A', 'D', 'H', 'U', 'P', 'I', 'N'])

* data type of columns:

titanic_train.dtypes

passangerId	int64
Name	object (/string)
Age	float64
Sex	object

* Some Questions:

- ① read data in df format
- ② print first five & last two rows
- ③ change column name
- ④ select one column & print its data type
- ⑤ print data dtype
- ⑥ check if any data is null

* [it is possible to convert list into series & vice-versa.]

* (skipping random rows is possible & not column - first we have to read all data & then filter it)

① `df = Pd.read_csv(r'https://---.csv', sep=',')`

② `df.head()`

② `df.tail(2)`

③ `df.rename(columns = {'age': 'Age'})`

check other way too

④ `df['age'].dtype`

⑤ `df.dtypes`

⑥ `df.isnull().sum()` ✓

⑦ `(, skiprows = [2, 3, 8])`

✗