20/4/23

# Pandas Python
## (from corey schafer)

used to
read tabular data.

## [Introduction]

* import Pandas as pd — General Convention of importing

* df = pd. read-csv ('path /to /file . txt') — read into the file

* df (short for Dataframe). Print the file !

* df. shape — columns x rows of the file (is shape of the file.

* df.info() — information of all the datatypes of the columns along with the shape of the data frame

datatypes → int64, float64, object
(some of 'em)          ∟ string type.
                          (generally)

pd. set_option ('display. man_columns', 36)
* pd. set-option ('display. man-rows', 20)

a man of 20 rows is displayed.

a man of '36' columns is displayed.

df. head() ——→ get first 5 rows.
* df. head(7) ——→ get first 7 rows.
df. tail() ——→ get last 5 rows.
df. tail(7) ——→ get the last 7 rows.

## [Dataframe & Series Datatypes]

Dataframe → 2-D Datastructure.

Dataframe is $111^{er}$ to a Dictionary.

eg people = {
        "name" : ["x", "y", "z"]
    } "role" : ["A", "B", "C"]

X→A
y→B
z→C

people["name"] ⟶ ["x", "y", "z"]

to convert this Dictionary into Dataframe we can do

~~df: DataFrame~~    df = pd·DataFrame (people)

*

print(df) ▽

| | name | role |
|---|---|---|
| 0 | X | A |
| 1 | Y | B |
| 2 | Z | C |

* type(df) = .pandas·core·frame·

→ indexes *  [INDEXES]    DataFrame·

print(df["name"])          type(df["name"])
       ↓                          ↓
                           pandas·core·series·Series

| 0 | X |
|---|---|
| 1 | Y |
| 2 | Z |

Series → 1-D Datastructure |||er to Lists but
              a lot more functionality.
                                        (or)
                                        Array.

We can Also write   df·name (or) df·role but

*  df["name"]  (or) df["role"]  is more preferred.
      (✓)              (✓)

Printing out Multiple Columns

        df[    ]
            ↳ pass in '1' column (or) a list of Columns.

*  df[ ["name", "role"] ]

Now, we've seen how to access columns,

How to Access rows?
    → for that we use iloc & loc.

→ Pass in a list (or) single index (row indexes)
→ of indexes

df.iloc [ 0 , 0 ]
necessary argument *

→ Pass in the column index (or) list of column indexes
optional argument *

→ Used for slicing the Dataframe with integer-indexes.

* df.iloc [ [0,1] , 1 ]

↓

| 0 | A |
|---|---|
| 1 | B |

* df.iloc [ 0 ]

→ prints out all the columns of row index 0.

i.e.

| | name | role |
|---|---|---|
| 0 | X | A |

* df.iloc [ [1,2] , [0,1] ]

↓

| | name | role |
|---|---|---|
| 1 | Y | B |
| 2 | Z | C |

But we can use custom-named indexes

| 0 | 0 1 2 |
| 1 | |
| 2 | ↖ regular |

df.loc ——→ similar to df.iloc

* df.iloc [ [1,2] , ["role", "name"] ]

↓

| | role | name |
|---|---|---|
| 1 | B | Y |
| 2 | Z | Z |

| | name | role |
|---|---|---|
| 0 | | |
| 1 | | ↖ |
| 2 | | custom-named |

* df.columns → shows (or) returns a list of all the columns in the Dataframe df.

↓

[ "name" , "role" ]

We can Also use slicing

df.loc [[0,1,2], ["name", "role"]]

can be
written as → * df.loc [0:2, "name":"role"]

df.iloc [0:2, 0:1]

* While slicing in pandas    0:2

* Both are inclusive.

(✓) Pandas    → 0:2 → 0,1,2
   generally → 0:2 → 0,1

## Indexes

We can set a column to act as a Index [Also
Indexes need not be unique according to pandas]
we can do it by

* df.set_index ("column_name")
  df

But this does not actually change df to    set index
as `column_name` we can do that by

df = df.set_index ("column_name")
          (or)
df.set_index ("column_name", inplace = True).

df.index ——→ shows or displays the index column.
*                of the DataFrame df.

                          are
→ When we do this we Asetting custom indexes

the use of this feature is
   Imagine we have fields like "name" & "email"
we have unique emails then that can act as a field

⇒ we do ‚ df = df.set_index ("email")

　　　　 ↳ We can Search based on email using

'loc'
　　 ＊ df.loc ("johndoe@gmail.com")
　　　　　　　　 ↳ gets all info relating to
　　　　　　　　　 the email "johndoe@gmail.com"
　　　　　　　　　 the name & stuff.

＊ Note
　 We cannot use ~~iloc~~ once we change the indexes
　　　　　　 loc[0]
　Since indexes 0,1,... do not exist since we
　changed it
　　　　　 But,
　　　　　　 we can still use iloc [0] & stuff
　　　　 ~~~~~ with iloc it works.
　　　　　　　　　 [still]

　We can reset indexes by doing
　　 ＊ df = df.reset_index()
　　　　　　(or)
　　　 df.reset_index (inplace = True)

　　　　 Now, all the integer indexes are restored.
　　　────────
　　　　 We can use ＊ schema_df.sort_index()
　　　　　　　　　　　　　 ↳ sorts all the index

　　　　　　　 (in pandas)　　 ascending=false.
　　　　　　　 can be　　　　 inplace = True
　　　　　　　 used in
　　　　　　 many places
　　　　　　 where the changes　　 ＊
　　　　　　 are temporary

　　　 df.shape, df.index, df.columns.

Note *

df.loc & df.iloc return the smaller Dataframes wrt to the values we put in them. ~~The~~ The returned Dataframes are of the type ~~pandas~~ 'pandas.core.frame. DataFrame'

Same as all the Dataframe.

(A new sub- Dataframe ~~/~~ is returned, of the original 'df'.

---

## Filtering

Consider Dataframe

|   | Name | last | email |
|---|------|------|-------|
| 0 | John | Doe | johndoe@gmail.com |
| 1 | Jane | Doe | Janedoe@gmail.com. |
| 2 | James | Brad | JamesB@gmail.com. |
| 3 | Jennifer | Anniston | JenniferAnniston@gmail.com. |

* We use Conditionals such as
  `'=='`, `'<'`, `'>'` to filter our DataFrame.

  → (df["last"] == "Doe") → filt

  ↘ this will return a Series of True/False
  (Array)           values.

[ True, True, False, False] → df[[ True, True, False, False]]

→ df [filt] will give indexes 0, 1
  ↳ it is a sub-DataFrame to the original DF
  (it is filtered.)

we can do ~~df.loc~~ df.loc [filt] instead
  ↘ this has 4 Advantages

* df.loc processes the Arrays of T/F values

we can do
  → df.loc [filt, ["Name", "email"]]
  *
  to get only Name & email fields.

Using Conditionals: '&' and '|'   'and'.

filt = (df["last"] == "Doe") & (df["name"] == "Joh...)

df.loc[filt]

↘ returns index '0'.
  Dataframe with ~

| Negation → '~'
| (Not)

df.loc[~filt]

↘ returns indexes
  '1', '2', '3'   Dataframe with

Similarly we use '|' for (or).

We can also use '>', '<', '>=', '<=' for :

— filtering.

Examples
→ high_salary = (df["salary"] > 70000)

→ df[high_salary]

→ df.loc[high_salary, ["Country", "prog_lang", "salary"]]

Now let's say we want to see if the "country" is in a set of our preferred countries

→ countries = ["USA", "UK", "India", "Germany", "Canada"]

we use ".isin()" function to check if
(Python-pandas fun1) "country" is in the array.

→ filt1 = df["country"].isin(countries)

Now let's say we want to see how many use "python" in their work.
consider "prog_lang" format as (HTML, CSS; Java; python —.

We use a (python-pandas func) ".str.contains()"
to check if substring is present inside a main string.

we need to put a additional parameter 'na = False'
to deal with NaN values in the DataFrame. ✳

(Null or None)   NaN = None = NULL
                   Basically NULL
                   values in
                   Pandas DataFrame.

~~~~~~~~~~~~~

filt 2 = df ["prog lang"].str.contains ("Python", na = false)
df.loc [filt 2].

filtering is one of the first
things we do with Data