Pandas in Python (self study).

① → Software library written for python for data manipulation & analysis.

i) → pd1 = pd. Dataframe (dict1)
                ↳ An already defined struc.
                 → sort of like an excel sheet is made.
        ↳ Import Pandas as pd.

ii) → pd1. to_csv ('friends.csv') } fn to convert data to excel sheet.
                           (transfer)
              ↳ Name of the excel sheet that will be made

iii) → pd1. head (2) → shows first 2 rows of table.
                     [useful when data has thousands of rows].
    • tail() also there
    • describe() → shows stats analysis of numerical columns.

iv) to read & edit → read1 = pd. read_csv ('read2. csv')
    excel sheet in                        ↳ to read.
    pandas           read1 ['columnName'] [0] = 50
                               ↳ to edit

② Pandas in Excel with power & functions provided by python.
Especially for data analysis. Also open source.
Generally used with numpy.

③ 2. data structures
                → Series (1D) (Just one column (or) row).
                ↳ Dataframes (2D) ( like in excel) .

→ Created a dataframe for showing operations:—

    newdf1 = pd. Dataframe (np. random. rand (334,5),
                                  index = np. arange (334))
                                       ↳ optional.

shape, describe, info (optional)
() () () extra

convert to numpyarray for other
→ newdf1.to_numpy() operation in numpy

④ Superimp func^n foundation newdf1. columns/rows = list ("ABCD").

newdf1. index     newdf1. describe()     newdf1.drop
                                         (0, axis)      commands shown

newdf1. dtypes     type (newdf1)

newdf1. sort - index (axis = x, ascending = false/true).

(or) sort -values

↳ doesn't modify original data.

⑤ warning → A value is trying - - - - - slice from Afrome.

→ shown bcz :-

of copy & view functions.

newdf2 = new df1.

newdf2 [0][0] = 93.           ⎱ newdf2 as
                              ⎰ view executed
then newdf1 also changed.       not copy..

      bcz unless you specify newdf2 = newdf1.copy()
newdf2 will act as a pointer & changes in it will
affect newdf1 too.

                                                    list
                                                  → for all
                                                  : ican be
                                                  : used.
∴ bcz of this confusion, we use newdf1. loc [0,0] = 50

eg:- 1.loc [(newdf1 ['0'] <0.3)]              ↳ loc func.
                                              to change values.
        ↳ Returns copy of all rows satisfying
                        given condition, (AND can also be used)
⑥                          in column 0,

→ iloc also there. Like loc selects rows & columns
                            wrt to their name.
        iloc → uses only index (0,1,...) not names given
                    by you.

⑦ Note :-

newdf1 = newdf1.drop [0] → changes the database

newdf1. drop [0] → shows how change will be, no change
                                        made unless '=' used
                        (OR):
                        inplace can be used.

# Pandas Masterclass.

① **Dataframe** → A table having rows & column.. (2D)
heterogeneous.

② ➤ merge, concatenate (or) reshape data. (slice, date)
Good at missing data handling.
Has powerful time series tool to work.

} Adv of Pandas.

③ ➤ Dataframe — → data
— → index
— → column.

} 3 components of dataframe.

Most convention way to create dataframe is by ndarray.

④ Diff formats from where data can be imported from:-
➤ .csv , .tsv , .pkl.
        ↓        ↓
    comma     tab
  seperated   seperated.

⑤ ➤ data = np. array ( . . . . ) → can be dict, list also.
s = pd. Series (data)
⤷ other parameter also present.

Simi for dataframe..
syntax

when extracting info. * &
[ ] → info abt single column.
[[ ]] → info abt multiple column.

⑥ **How errors happens**

➤ when multiple paragraph error given error is only one place but so many functions & classes are called that all error are shown.
You can use Traceback to see where first point is.

```
df.value_counts()    b fill → belowfill.
normalize
       to True    df.fillna('4')
       for
       percentage.         └→ fills all null
                              value by 4.
```

*① df [['Name', 'Age']]
    df ['Name'].

⑦ df [df ['Age' > 20]] → returns row.

    df ['Age' > 20] → returns true/false.

⑧ df.loc [0] → returns first row.    ┌─────────┐
                                      │   for   │
                                      │ editing. │
                                      └─────────┘

    df.loc [0, ['Name']] → returns name from
                                      first row.

    df.loc [0:2 ['Name', 'Age']]
              └→ start to end (slicing)

    iloc → only uses index
         └→ if slicing used, goes to start to end-1.
                 └→ Generally iloc & loc only used for slicing.

⑨ → iloc → only for index specified in parameter.
              Not columns, for manipulating (or) taking
              (or) editing columns use loc.

⑩ Inplace = True & drop → imp concepts.
                └→ for making changes permanent.

⑪ → while reading file into pandas, path is to be specified
     if .csv not in same folder. [ r to be mentioned
                                      before " ].

⑫ → df.iloc [0:5, -2:] → try.

** df.drop ('column name', axis=1, inplace=True).
                    └→ if dropping column give axis=1
                       (or) else won't execute.

Other func → ununique() ⟶ no parameters ⟶ list()

df.isnull().sum()
↳ imp
not null() or isnull() == false

 Page No. Date

(13) Groupby concept (imp) & little adv | rename()
& inplace()

zip ⟶ store in something
(like dict (or) list
merge() } adv to access).
eg:- dict1 { zip(---- }

both
similar, see
once

→ Group by

df.column

» df.groupby('columnname').size()

* unique
(or)
unique

df.groupby('Sex')['Survived'].mean().

| S.No. | City | Dept | Manager |
|-------|------|------|---------|
| 1 | Mumbai | Sales | Ds-M |
| 2 | Mumbai | Sales | " |
| 3 | Goa | Sales | " |
| 4 | Goa | Tech | " |

groupby('City') → 2 block
Goa & Mum

groupby('Manager') → 1 block
whole t

↳ used instead of for loop

⤴ Never
used

(14). Joining 2 different datasets.

1) ⟹ Inner Join, Outer Join, etc.

ii) ⟹ Combining df1 & df2 → [ merge() ]

↳ can only be used for
common columns.

⟹ pd.merge(df1, df2, on=
'city').

iii) ⟹ on='', how='' & other parameters can
also be passed on.

(15) ⟹ pd.concat() → Merge dateframe vertically &
horizontally. entirely.
(without particular condition like
inner (or) outer join).