

Pandas

# pandas

Import numpy as np

Import pandas as pd

```
data = pd.Series([0.25, 0.5, 0.75, 1.0])
```

data.

```
put 0 0.25
    1 0.50
    2 0.75
    3 1.00
```

→ data.values

```
array([0.25, 0.5, 0.75, 1.0])
```

→ data.index

Range index (start=0, stop=4, step=2)

→ data = pd.Series([1, 2, 3, 4], index=['a', 'b', 'c', 'd'])

→ data.values

→ data.index → Index values are a, b, c, d.

→ Pandas are capable for non-contiguous or nonsequential indices



→ Population-dict = {  
'California': 32222,  
'Texas': 2234,  
'New York': 34222,  
'Florida': 349999}.

Population = pd.Series(Population-dict).  
population

→ Data can be a scalar, which is repeated to fill the specified index.

→ pd.Series(5) index = [100, 200, 300].

output  
100 5

200 5

300 5

Data can be a dictionary, in which index denotes

→ pd.Series({2: 'a', 1: 'b', 3: 'c'}) to the sorted dict keys

output  
1 b

2 a

3 c.

## Imp in pandas

Series → one-D Array →

Horrible indices.

Dataframe → 2-D Array →

both Horrible row indices,

Horrible column names

-> area-dict = { 'California': 32432, 'Texas': 695662,  
                  'Newyork': 294532, 'Florida': 39454,  
                  'Gmina': 149995 }

$\rightarrow \text{area} = \text{pd.Series}(\text{area\_dict})$ .

Creating a 2-D object containing this info using DataFrame.

```
States = pd.read_csv('states.csv')  
States = States[{'population': population,  
                 'area': area}]
```

→ States-index

Dataframe can be

Thought as a generalized  
3dian of 2-D numpy  
array where both the  
rows and columns have  
a generalized index for  
accessing the data.

→ From a list of dict. Any list of dictionaries can be made into a dataframe.

`Data = [{'a': i, 'b': 2*i} for i in range(3)]`

`Pd.DataFrame(Data)`

Output	a	b
0	0	0
1	1	2
2	2	4

→ From a dictionary or Series objects.

`Pd.DataFrame({'population': population, 'area': area})`

→ From a 2-D Numpy Array.

Given a 2-D array of data, we can create a DataFrame with any specified column and index names. If omitted, an integer index will be used for each:

`Pd.DataFrame(np.random.rand(3,2), columns = ['fore', 'back'], index = ['a', 'b', 'c'])`



→ From a Numpy Structured Array.

A Pandas Data Frame operates much like a structured array, and can be created directly from one.

→  $A = np.zeros(3, \text{dtype}=[('A', 'i8'), ('B', 'f8')]).$

Output:  $\text{array}([0, 0.0), (0, 0.0), (0, 0.0)], \text{dtype}=[('A', 'i8'), ('B', 'f8')]).$

→ pd.DataFrame(A).

### The Pandas Index Object

→ This Index object is an interesting struct itself, and it can be thought of either as an immutable array or as an ordered set.

ind = pd.Index([2, 3, 5, 7, 11]).

→ Print(ind.size, ind.shape, ind.ndim, ind.dtype)

→ indA = pd.Index([1, 3, 5, 7, 9]).

indB = pd.Index([2, 3, 5, 7, 11]).

→ indA & indB # Intersection.

→ indA | indB # Union.

→ indA ^ indB # Symmetric difference.



→ Indexers : loc, iloc and ix.

data = pd.Series(['a', 'b', 'c'], index=[1, 3, 5])  
data

loc attribute allows indexing and slicing that always referencle the implicit index.

data.loc[2]

→ 'a'

→ data.loc[1:3]

output

1	a
3	b

→ loc attribute allows indexing and slicing that always referencles the implicit python-style index.

→ data.iloc[2]

'b'

→ data.iloc[1:3]

3	b
5	c

# Handling missing Data

- Numpy does provide some special aggregations that will ignore these missing values.

Import Numpy as np

Import pandas as pd

```
Vals1 = np.array([1, np.nan, 3, 4])
```

Vals2.

- NAN: Missing numerical data

```
Vals2 = np.array([1, np.nan, 3, 4]).
```

Vals2.dtype

↳ Everything or any operation done with Nan gives out to be NAN.

- np.nansum(Vals2), np.nanmin(Vals2), np.nanmax(Vals2)

NAN and None in pandas

= = = = =

```
Pd.Series([1, np.nan, 2, None]).
```

```
n = pd.Series(range(2), dtype=int)
```

None

## A Operating on null Values

→ Pandas treats None and NaN as essentially interchangeable for indicating missing or null values.

Several useful methods are:

isnull()

Generate a Boolean mask indicating missing values.

notnull()

opposite of isnull().

dropna()

Return a filtered version of the data.

fillna()

Return a copy of the data with missing values filled or imputed.

## A Detecting null Values

→ Two useful methods for detecting null data:  
isnull() or notnull().

→ data = pd.Series([1, np.nan, 'Hello', None]).

→ data.isnull().

→ data[data.notnull()].



## \* Dropping null values.

- dropna() [which removes NA values].
- fillna() (which fills in NA values).

## \* Dropping null values in Data Frame

- df = pd.DataFrame([[1, np.nan, 2],  
[2, 3, 5],  
[np.nan, 4, 6]]).
- df.

dropna() will drop all the null values present.

→ df.dropna() → will drop all values from columns and rows.

→ df.dropna(axis=1 (columns)) → dropping all null values from column itself.

→ df[3] = np.nan

df → creating new column and value of nan value

→ df.dropna(axis=1 (columns)) how='all').



removing all the rows  
and columns containing  
NA values.

→ df.dropna(axis='rows') thresh=3).

SHOT ON REDMI  
AI DUAL CAMERAS for the new / column to keep:

- `data = pd.Series([1, np.nan, 2, None, 3], index = list('abcde'))`
- `data`
- `data.fillna(0)` → This will fill the `NAN` values as zero
  
- # Forward-Fill.
- `data.fillna(method = 'ffill')`
- we can specify a forward-fill to propagate the previous values forward.
  
- # back-fill
- `data.fillna(method = 'bfill')`
- and we can specify a back-fill to propagate the next values backward.
  
- `df.fillna(method = 'ffill', axis = 1)`.

A Multiply Indexed Series

- Representing two dimensional data within a one-d series.
  
- `index = [('california', 2010), ('california', 2010), ('Newyork', 2010), ('Newyork', 2010), ('Texas', 2010), ('Texas', 2010)]`.
  
- `populations = [33871648, 37253956, 18976457, 19378102, 20851420, 25145561]`.
  
- `pop = pd.DataFrame(populations, index = index)`
- Pop

- with this indexing scheme you can straightly index or slice the series based on this multiple index
- `pop[['California', 2010], ('Texas', 2010)]`.  
For the values of 2010 → you'll need to do some messy munging to make it happen.
- `pop[[i for i in pop, index if i[2] == 2010]]`.

## Pandas Multindex

- index = pd.MultiIndex.from\_tuples(linden)
- index
- pop = pop.reindex(linden)
- pop
- pop[:, 2010] → for accessing all data from 2010.

## MultIndex as extra dimension

- pop\_df = pop.unstack()
- pop\_df

unstack method will quickly create and convert a multiply indexed series into a conventionally indexed data frame.

- pop\_df.stack().

The stack() method provides the opposite operation.

## Combining Datasets : concat and append

Import pandas as pd.

Import numpy as np

def make\_df(cats, ind):

    data = [i + str(c) + str(i) for i in ind]  
    for c in cats]

    return pd.DataFrame(data, ind).

make\_df('ABC', range(3))

### Simple concatenation with pd.concat

Ser1 = pd.Series(['A1', 'B1', 'C1'], index=[1, 2, 3]).

Ser2 = pd.Series(['D1', 'E1', 'F1'], index=[4, 5, 6]).

pd.concat([Ser1, Ser2]).

df1 = make\_df('AB', [1, 2]).

df2 = make\_df('AB', [3, 4]).

print(df1); print(df2); print(pd.concat([df1, df2]))

df3 = make\_df('AB', [0, 2]).

df4 = make\_df('AB', [0, 1]).

print(df3); print(df4); print(pd.concat([df3, df4]))

axis=1.

## Duplicate Indices

```
n = make_df(['A', 'B'], [0, 1])
y = make_df(['A', 'B'], [2, 3])
y.index = n.index #
print(n); print(y);
print(pd.concat([n, y]))
```

## Ignoring the Order

→ print(n); print(y); print(pd.concat([n, y]), ignore\_index = True))

## Adding multiIndex Keys

→ Keys are helpful to specify a label for the data samples  
The result will be a hierarchically ordered series  
(containing the data):

→ print(pd.concat([n, y]), keys = ['n', 'y']))

## Difference

np.concatenate and  
pd.concat is that  
Pandas concatenation  
preserves indices, even  
if the result will have  
duplicate indices.

## \* concatenation with joins

d75 = make\_dt('ABC', [1, 2])

d76 = make\_dt('BCD', [3, 4])

print(d75); print(d76); print(pd.concat([d75, d76]).)

→ The join is a union of the input columns (join='outer')  
[outer = union]

→ The inner join is the join which can be used for the intersection of the columns.

Print(d75); Print(d76); print(pd.concat([d75, d76],  
join\_axes=[d75.columns]))

## \* The append method

→ Direct array concat can be made by append method using the fewer strokes only.

→ For ex:- pd.concat([d71, d72]), you can simply call d71.append(d72):

→ append method in pandas it does not modify the original object, it creates a new object with combined data.

creates a new

→ The reason is, append, index and date buffer

## \* Combining datasets: merge and join

→ pd.merge() → This is a subset of what is known as Relational Algebra which is just a formal set of rules for manipulating Relational Data.

### Three categories of join

- one-to-one join
- many-to-one join
- many-to-many join

#### (i) one-to-one join

```
d71 = pd.DataFrame({'Employee': ['Bob', 'Jake', 'Lisa', 'Sue'],
                     'group': ['Accounting', 'Engineering', 'Engineering',
                               'HR']})
```

```
d72 = pd.DataFrame({'Employee': ['Lisa', 'Bob', 'Jake', 'Sue'],
                     'Hire-date': [2004, 2008, 2012, 2014]})
```

```
d73 = pd.merge(d71, d72)
```

(ii) many-to-one join

```
d74 = pd.DataFrame({ 'group': ['Accounting', 'Engineering', 'HR'],
                     'supervisor': ['carly', 'guido', 'steve'] })
print(pd.merge(d73, d74))
pd.merge(d73, d74)
```

(iii) many-to-many joins

By performing a many-to-many join, we can recover the skills associated with any individual person:

```
d75 = pd.DataFrame({ 'group': ['Accounting', 'Accounting', 'Engineering', 'Engineering', 'HR', 'HR'],
                     'skills': ['math', 'spreadsheets', 'coding', 'linear', 'spreadsheets', 'organization'] })
```

\* `dat.unique` → returns array of unique values.

\* `dat.nunique` → gives out number of unique values

\* `.corr()` = shows correlation