# PANDAS.

Date ..............

Attributes & imp fouctionality

# Pandas is built-on Numpy and Matplot lib.

①  • head()                    ⑦  • shape # (row, col) tuple
②  • info()
③  • describe() # numerical values.
④  • values. # values {In the form of 2d Array]
⑤  • columns # column names
⑥  • index # rows names

# SORTING and SUBSETTING

①  • sort_values ("col_name", ascending = false)

   • sort_values (["col_name1","col_name2"]).

②  dogs['name']

   dogs [['name', dogs 'breed']]

   outer                        Inner square brackets are
   sbort subsetting             list of column names
   the desired /col"
   from df.

   dogs [['name', 'breed']] {, ascending = [True, false]}

Spiral

filtering → condition -

dogs [ dogs [ "height_cm" ] > 50 ]

dogs [ dogs [ "breed" ] == "Labrador" ] - ⍺

dogs [ dogs [ "dob" ] < "2015 - 01 - 01" ]

Y    M    D

is_lab =           ⍺

is_brown = dogs [ "color" ] == "Brown" )

dogs [ is_lab & is_brown ].

Subsetting using .isin()

• isin()

is_black_or_brown = dogs [ "color" ].isin ( [ "Black", "Brown" ] )

New Colours in pandas:

dogs [ "height_m" ] = dogs [ "height_cm" ] / 100 .

print (dogs ).

dogs["bmi"] = dogs["weight_kg"]/dogs["height_m"]

## Multiple Manipulations.

## SUMMARY STATISTICS.

~~mean~~

① dog["height_cm"].mean().     ⑦ . std()

② . median()                    ⑧ . ~~@~~ sum().

③ . mode()                      ⑨ . quantile ()

④ . min()          a list not   ⑩ ~~agg~~ . agg()
                   a single number

⑤ . max()              ⎰ ⑪ . cumsum() , . cummax()
                       ⎱    . cummin(), . cumprod

⑥ . var()

We can also get the summary statistics for
date values.

dogs["dob"].min().      dog["dob"].max()

## AGG method

def pct30 (column) :      → 30% of data on col<sup>n</sup>.
        column . quantile (0.3)

→    dogs["weight_kg"] . agg (pct30)

dogs[["weight_kg", "height_cm"]].agg (pct30)

```python
def pct40(col):          40%
    col.quantile(0.4)
```

```python
dog["weight"].agg([pct30, pct40])
```

```python
dog["weight"].cumsum()
```

# Dropping Duplicate values

```python
vet_visits.drop_duplicates(subset="name")
```

```python
uv=vet_visits.drop_duplicates(subset=['name','breed'])
```

```python
print(uv)
```

## To Count:

```python
unique_dogs["breed"].value_counts()
```

```python
unique_dogs["breed"].value_counts(sort=True)
```

## Proportions

```python
unique_dogs["breed"].value_counts(normalize=True)
```

# Group by / Grouped Summary Statistics

Date .............

dogs. groupby ("color")["weight_kg"]. mean()
↓
the columns
to groupby
on

↓
Df

dogs. groupby ("color")["weight_kg"]. agg([min, max, sum])

dogs. groupby (["color", "breed"]).["kg"]. mean()

dogs. groupby (["color", "breed"]).[["kg", "cm"]]. mean().

## Pivot Tables

## Groupby to pivot Table.

dogs. grouby("color")["weight_kg"]. mean()

↗ performance
for sum
dogs. pivot_table (values = "weight_kg"     statistics
index = "color")

to be group by

```python
dogs.pivot_table(values = "weight_kg",
                 index = "color",
                 aggfunc = np.median)
```

```python
dogs.pivot_table(values = "weight_kg", "color"
                 index = "color",
                 aggfunc = [np.mean, np.median])
```

```python
dogs.groupby(["color", "breed"])["weight_kg"]
    .mean()
```

```python
dogs.pivot_table(values = "weight_kg",
                 index = "color",
                 columns = "breed").
```

# When there are missing values (NaNs)
in pivot table

```python
dogs.pivot_table(values = "weight_kg",
                 index = "color",
                 columns = "breed",
                 fill_value = 0,
                 margin = True)
```

mean of all except the filled_values/
                              NaNs

# EXPLICIT

## INDEXES.

dogs_ind = dogs.set_index ("name").

### for Removing the index.

dogs_ind.reset_index().

### Dropping an index.

dogs_ind.reset_index(drop=True)
↓
drop the -
~~value~~ index
entirely removing
that index/columns.

## WHY INDEXING?

# Indexing makes subsetting easier

# ~~Index~~ values in the index don't
need to be unique.

# Multi-level indexes
## a.k.a hierarchical Indexes

Date ..............

dogs_ind3 = dogs.set_index (["color", "breed"])

## Subsetting with outer level

dog_ind3. loc [["lab", "Chihuahua"]]

## Subsetting with inner levels with a
## list of tuples.

dog_ind3. loc [('lab', 'brown'), ('Chihua', 'Tau')]

## Sorting by index values

dogs_ind3. sort_index ()

~~dogs_ind3~~

## Controlling sort-index

dogs_ind3. sort_index (level = ["color", "breed"],

ascending = [True, False]

*Spiral*

# Now we have 2 problems

① Index values are just data, and storing data in multiple forms confuses with my/our brain.

② It also violates "tidy data" principles

normal tabular data

③ plus need to learn two syntaxes

④ I would suggest not to use but knowing about it is better [: It helps in reading other peoples code some might find it useful:

# SLICING AND SUBSETTING
## WITH .loc and .iloc

→ Sorting the index before slicing.

dogs_srt = dogs.set_index(["breed", "color
                            sort_index()

# Slicing at the outer index level

Date ..............

```
dog_srt.loc["Chow Chow" : "Poodle"]
```
↓ It is included
← specifying the index values →

for **inner index level** : List of Tuples.
```
dog-srt.loc["Tan" : "Grey"] . (X)
```
↓
Wrong approach

```
dog-srt.loc[("lab" ; "Brown") : ("Schnauzer", Grey)
```

## Slicing columns

```
dog-srt.loc[:, "name" : "height_cm"]
```

mixing both slicing for rows as well as columns

```
dogs-srt.loc[("lab" ; "brown") : ("sch", grey),
"name" : "height"]
```

Spiral

Index Slicing can be ~~~~~~~~
Imp when it comes to
Dates.

dogs = dogs.set_index("DOB").sort_index()

So now, I can directly slice
the data using dates/ partial dates.

dogs.loc["2014":"2016"]
                        → included.

dogs.iloc[2:5, 1:4]
                    → final values are
                          not included

## Working with pivot tables

dog_height_bybreed-vs_color = dog_pack.pivot_table(
                        "height", index = "breed",
                        columns = "color")

dhbbvc.loc["chow chow": "poodle"]

dhbvc.mean(axis="index")
                    → rows

# Visualizing data

```python
import matplotlib.pyplot as plt.

dog_pack ["height_cm"]. hist().
plt. show().
```

## Bins:

```python
dog_pack ["height_cm"]. his(bins=5).
    plt. show().
```

## BAR PLOT

```python
avg_weight_by_breed = dog_pack.groupby("breed")
                            ["Weight_by"].
                        mean().

avg_weight_by_breed.plot(kind="bar",
            title="Mean Weight by dog breed")

plt. show().
```

# Line plots.

```
sully. plot( x= "date", y= "weight", kind= "line")
    plt.show()
of.
```

```
sully. plot (x= "date", y= "weight", kind= "line",
plt.show()      rot= 45)
                the x axis labels
    rotating at 45° to read easily.
```

## SCATTER PLOT { rel" b/w two numeric variables}

```
sully. dogpack (x= "height_cm", y= "weight",
                kind ='scatter')

    plt. show()
```

## LAYERING PLOTS { one plot on top of other plot}.

```
dog-pack [dog-pack["sex"] == "F"]["height"].hist()

dog.pack [dog-pack["sex"] == "M"]["height"].hist()
```

plt. legend (["F"; "M"])

plt. show()

for Transparency

dog-pack["sex"] == "F"]["height"].hist(alpha=0.7)

dog-pack["sex"] == "M"]["height"].hist(alpha=0.7) ..

plt. legend (["F"; "M"])

plt. show().

## MISSING VALUES.

Detecting missing values

dogs. isna().

dogs. isna(). any() → # In columns.

dogs. isna(). sum()

# Removing Missing Values

```
dogs. dropna ()
     Replacing missing values
dogs. fillna (0) .
```

## Creating dataframes

```
my_dict = { "key1" : value1,

            "key2" : value2,

            "key3 : value3 }
```

```
my.dict = { "title" : "Charlotte's Web",
            "author": "E.B. White,"
            "published" : 1952

          }
```

```
pd. DataFrame (list_of_dicts).
```

## DICTIONARY OF LISTS - By Column

```
pd. Dataframes.
```

# Reading & Writing CSV

Date ..............

```
# new_dogs = pd.read_csv("new_dogs.csv")

new_dogs.to_csv("new_dogs_with_bmi.csv").
```

---