

# Data manipulation, Data Filtering

src: dpoint  
yt channel

dataframe is a collection of series.

df.dtypes

col name

datatype

df.describe()

↳ applied only on numerical & columns

↓

it describes

# mean, std, count,  
min, 25%, 50%, 75%,  
max

dtype:

> df.dtypes == 'object'

Ex:

col1

False

→ not an object dtype

col2

True

→ object dtype

> (df.dtypes == 'object').dtypes → # bool dtype('bool')

> df.dtypes[objectindex]

↓

df.dtypes == 'object'

→ { col 5 object  
col 6 object  
dtype: object

> df.dtypes[objectindex].index

only prints true values

→ index(['col 5', 'col 6'], dtype = 'object')

Accessing columns:  
 $df['col1']$  # for ~~passing~~ accessing a single value column  
 $df[['col1', 'col2', \dots]]$  # for accessing 1 or more columns.

Accessing Rows:

$df['col1'][0]$  → first value of col1

$df[0]$  → first row

$df[0:]$  → all rows

$df[0:n]$

↳ upper bound.

lower bound

→ ~~df['col1']~~

$df['newcol'] = 'dpoint'$

↳ assigned to every row in newcol.

→ ~~pd~~ pd.Categorical( $df['col6']$ )

→ like levels in R lang

→  $df['col6'].unique()$

→  $df['col6'].nnull()$

# Visualization

imports  $\rightarrow$  numpy  
pandas  
seaborn  
matplotlib.pyplot

[ %matplotlib inline ]

$\hookrightarrow$  To display inline plots for Jupyter notebooks

from pandas.plotting import scatter\_matrix.

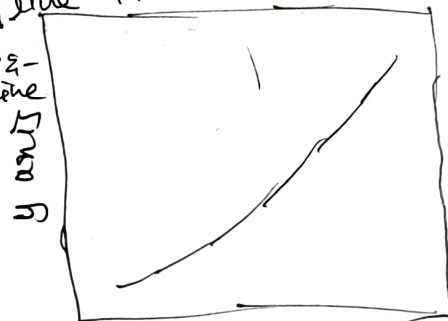
from numpy.random import randn, randint, uniform, sample

$n = [1, 2, 3, 4, 5]$

$y = [a^2 \text{ for } a \text{ in } n] \rightarrow \text{set of squares of } n.$   
[1, 4, 9, 16, 25]

plt.plot(x, y, 'r-')  
plt.xlabel('x axis')  
plt.ylabel('y axis')  
plt.title('Title is')

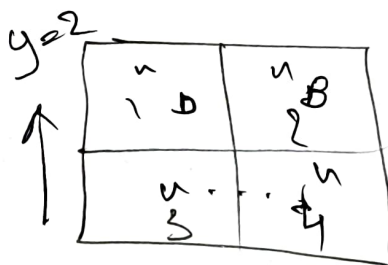
r - simple straight line  
-- st line  
•  $\rightarrow$  red line  
r-o  $\rightarrow$  2-line



$x^2$   
 $\downarrow \downarrow \rightarrow$  red line plot.  
red line

plt.subplot( $\frac{1}{n} \frac{1}{y} \frac{1}{n}$ )

(2, 2, 1)  $\rightarrow$  D  
2  $\rightarrow$  B



$\rightarrow x = 2$

# Object oriented drawing

```
fig = plt.figure()
```

↓  
object.

```
axes = fig.add_axes([1, 1, 1, 1]) or ([0.5, 0.5, 0.8, 0.8])
```

```
axes.plot(x, y, 'b')
```

(x, y, 'g')

```
axes.set_xlabel('x axis')
```

blue

```
axes.set_ylabel('y axis')
```

} } }



both these plots  
are on same canvas

in figure() → ~~edgecolor~~ = 'r' <sup>red</sup> ~~border~~  
line width → line width of edgecolor

→ if axes are not set properly <sup>they</sup> plots may overlap.

```
df = pd.DataFrame(randn(1000), index=np.linspace(1, 1000, 1000),
```

columns=['value'])

```
ts = pd.Series(randn(1000), index=np.linspace(1, 1000, 1000))
```

```
df.describe()
```