

Python recap

Objects → Mutable, $a = 5$
 Objects → Immutable, $a = (1, 3, 5)$, dictionaries

Char. of objects : { type : number, string, list, ...
 value : data value contained
 identity : related with CS

Attributes : data, functions associated with them. Name of attrs follows the name of obj, follow a dot

↳ { Data attributes : Data attached to an object.
 Methods : Function attached to an object. Typically it performs a function on the obj.

Example - import numpy as np
 $x = np.array([1, 3, 5])$ → Data attribute : $x.shape$: #elmts in x
~~np.array([1, 3, 5])~~ → no operation
 \rightarrow Method : $x.mean()$: mean of values
 there is operation

Modules : Libraries of code, eg math, numpy, ...

dir(str) gives a list of all available methods for a string

help(name.upper) gives info about the method upper applied to the string name.

Numbers : Integer, floating point, complex.

$$2 + 5, \quad 2 * 5, \quad 2 ** 5 = 2^5, \quad 6 // 7, \quad 6 // 7 = \text{parte entera de } 6/7$$

- In the interactive mode , - is the value of the last value returned.

$$\text{math.factorial}(6) = 6!$$

- input random , random.choice([]) gives a randomness of []

- Booleans . True, False, and, or, not, < , <=, ==, >=, !=

In lists, == compares identity of list, whereas is compares whether the objects are the same .

Sequences . lists, tuples, range objects, strings

S[0] , = location 1st , S[0:2] = list w/ loc. {from 2nd}

- Lists : Mutable sequences , strings are not .

S = [2, 4, 6, 8] . ~~when~~ S.append(10) = [2, 4, 6, 8, 10] . \checkmark new S

T = [10, 14, 16] , S+T = [2, 4, ..., 16]

S.reverse() , S = [10, 8, ..., 2]

S.sort() , S = [8, 2, 10, 6, 4]

(There is a function also for this: sorted())

len(S) = 5 (length) ~~long~~

(2)

- Tuples: ~~Immut~~ immutable sequences.

$T = (1, 3, 5, 7)$, can ~~not~~ $+$ and access as $T[0] = 1$.

$x = 72.15$ ~~not~~ coordinate $= (x, y)$ is a tuple.
 $y = 75.2$

I can do $(c1, c2, c3, c4) = T$, and then $c1 = 1$
 $c2 = 3$

Say now that $P = [(0, 0), (5, 9), \dots, (4, 21)]$ list of tuples:

Can do: $\left\{ \begin{array}{l} \text{for } (a, b) \text{ in } P : \\ \quad \text{print}(a, b) \end{array} \right. \quad \left(\begin{array}{l} \text{for } a, b \text{ in } P \\ \quad \text{also works} \end{array} \right)$

$c = (2)$ is an integer, $c = (2,)$ is a tuple w/ 1 object

$T.\text{count}(7)$ counts the number of 7 in T .

$$\text{sum}(T) = 1 + 3 + 5 + 7$$

Range: ~~Immut~~, it is a type!!

range(5) gives range(0, 5). To see the content do

$$\text{list}(\text{range}(5)) = [0, \dots, 4]$$

$$\text{list}(\text{range}(1, 6)) = [1, \dots, 5]$$

$$\text{list}(\text{range}(1, 13, 2)) = [1, 3, 5, \dots, 11]$$

Strings: Immutable sequences of characters

$S = "Python"$, $\text{len}(S) = 6$, $S[0] = 'P'$, $S[-1] = 'n'$

$S[0:3] = "Pyt"$ $S[-3:0] = "hon"$ (last 3 char)

"j" in S : True

"hello" + "word" = "helloworld"

$\text{name} = "Tina Fey"$, $\text{name.replace("T", "t")}$ = "tina Fey"

$\text{names} = \text{name.split(" ")}$ (splits when there is space)

$\text{names} = ["Tina", "Fey"]$

"tone".upper() = 'TONE'

lower

" ".join(list of words) = word₁ word₂ word₃ ...

Sets: Unordered collections of hashable objects
means that ~~the~~ you can use them for mutable obj.

Frozen set : immutable

Set : mutable

$\text{ids} = \text{set}([1, \dots, 7])$ turns the list [1, -7] into a set w/ 7 elements

$\{1, 2, 3, 4, 5, 6, 7\}$

$\text{ids.add(8)} = \{1, \dots, 8\}$

(3)

`ids.pop()` removes a random element of `ids`.

Say now `ids = set(range(10))`, `ids = {0, ..., 9}`

`males = set([1, 3, 6, 7])`

`females = ids - males`

`everyone = male | female`

`everyone & set([1, 2, 3])`

`females = {0, 2, 4, 5}`

(`| := U`, `- := \setminus`)
`& := \cap`)

`word = "antidisestablishmentarianism"`

~~ask about~~ `letters = set(word)`,

`key len(letters)`. gives the # of unique letters.

`x.issubset(y)` determines whether $x \subseteq y$.

Dictionaries: Maps from key objects to value objects.

Consists of pairs Key: Value
 \uparrow
 immutable

↳ curly braces

`age = { "Tim": 29, "Jim": 31, "Pam": 27, "Sam": 35 }`

`age["Pam"] = 27`

`age["Tim"] = age["Tim"] + 1`

~~#~~ += 1

`names = age.keys()`, it is a "dict-key" type object.

`age["Tom"] = 50` adds a new entry in the dict

!!! In names "Tom" will be included! "View-object"

`ages = age.values()`

`Tom in age` : True.

`age.items()` # return a dict-key

object, essentially a list of tuples
[(name, age)]

Caution !! (Dynamic typing):

`x = 3`

`y = x`

`x = 2`

Result:

$\begin{cases} x = 2 \\ y = 3 \end{cases}$

Because integers are immutable objects

Strings also immutable

Tuples also immutable

`x = [1, 2, 3]`

`y = x`

`x[0] = 7`

Result

`x = [7, 2, 3]`

`y = [7, 2, 3]`



Lists are mutable objects

Dicts also mutable

To keep y as before, two possibilities

$\left. \begin{array}{l} y = \text{list}(x) \\ y = x[:] \end{array} \right\}$

In both cases it creates a new object.

(4)

Remark : $\left\{ \begin{array}{l} == \text{ equality element-wise} \\ \text{is } \text{ equality as objects} . \end{array} \right.$

$x = [1, 2, 3]$

$y = [1, 2, 3]$

$x == y$: True

$x \text{ is } y$: False

$x = [1, 2, 3]$

$y = x$

$x \text{ is } y$: True

Statements : Return, import, pass (do nothing, just in some cases for syntactical reasons)

if ---

elif ---

else ---

For : $\left. \begin{array}{l} \text{for } n \text{ in range}(5): \\ \quad \text{print}(n) \end{array} \right\}$

0
1
2
3
4

names = ["Pepe", "Juan", ...]

$\left. \begin{array}{l} \text{for name in names:} \\ \quad \text{print}(*name) \end{array} \right\}$

"Pepe"
"Juan"
"

$\left. \begin{array}{l} \text{for i in range(len(names)):} \\ \quad \text{print}(names[i]) \end{array} \right\}$

"Pepe"
"Juan"

In other
languages one does
(because it tests
own list is
not possible)

Can iterate over a dictionary as follows:

```
age = { "Juan": 55, "Pepe": 20, ... }
```

```
for name in age:  
    print(name, age[name])
```

} Juan 55
Pepe 20

* Can also do ... in age.keys()

Or if I want the names to be showed alphabetically, ... in sorted(age.keys()).

sorted(— , reverse=True) will order the list in reverse alphabetically and

List comprehension: $g = [a + \#2 \text{ for } a \text{ in range}(10) \text{ if } a \% 2 == 0]$.

Reading and writing files

```
for line in open("oro.txt"):  
    print(line)
```

soy
un
oro

w/ extra spaces.
This is because in the
txt file it actually says

Soy\n
un\n
oro .

This can be avoided using rstrip() method:

```
for line in open("oro.txt"):  
    line = line.rstrip()  
    print(line)
```

soy
un
oro

(5)

`split("")` method takes a string and divides it whenever a " " .

Create a txt file :

```
F = open("output.txt", "w")
F.write("Python\n")
F.close()
```

w+ : write & read
a+ : append

↑ write

↑ to close

for i in range(1, 11) :

F.write("The square of " + str(i) + " is " + str(i**2) + "\n")

F.read() returns a single string with the whole text (also \n).

F.readlines() returns a list with entries the lines of ~~the~~ the file
(also with \n, since they are strings)

```
Q = open("lol.txt", "r")
```

linea = Q.readlines(16)

print(linea)

linea = Q.readlines(5)

print(ln)

linea = Q.readlines(9)

print(ln)

lol.txt

This is line 1
This is line 2
----- 3

This is line 1

This
is line 2 .

Eg: If a file has

Pedro,	22,	1.78
Jorge,	25,	1.62
:	:	:

how to extract this data? Use `split(",")` !

~~Reading Data~~

~~Reading Data~~

```
D = open("datos.txt", "r")
Z = []
for line in D:
    currentline = line.split(",")
    Z.append(currentline)
```

(even better, line.rstrip().split(","))

Functions

inputs
↓
def add(a,b):
 mysum = a+b
 return mysum.

def add_and_sub(a,b):
 mysum = a+b
 mydiff = a-b
 return (mysum, mydiff)

⚠ Warning: For mutable objects, a function ^{might} changes its value!

def modify(mylist)
 mylist[0] = 10
 return (mylist)

L = [1, 3, 5, 7, 9]
M = modify(L)
M = [10, 3, 5, 7, 9]
L = [10, 3, 5, 7, 9]
=====

!!!

Example: Function that creates a password:

import random

(random.choice([....]) selects randomly an element of the list)

def password(length):

pw = str()

chars = "abcd...yz01...9"

for i in range(length):

pw = pw + random.choice(chars)

Scope rules: It might be that at different points of a program there are several ~~two~~ functions or variables with the same name (eg inside a function and outside). Python follows

L local
E enclosing function
G global
B built-in

- Current function
- Function that called the current function, if any
- Module that the function was defined
- Python's built-in name space

If Python cannot resolve a name, error!

Remark: Argument: Object that is passed to a function as its input when the function is called

Parameter: Variable that is used in the function definition to refer to that argument

(n=360)

↑ ↑
variable object

(name, label of obj)

Example: Look how it works depending on mutable / immutable objects:

```
def update(n, x):  
    n = 2  
    x.append(4)  
    print('update:', n, x)  
  
def main():  
    n = 1  
    x = [0, 1, 2, 3]  
    print('main:', n, x)  
    update(n, x)  
    print('main:', n, x)
```

main

Returns

main: 1 [0, 1, 2, 3]

update: 2 [0, 1, 2, 3, 4]

main: 1 [0, 1, 2, 3, 4]



n does not change its original value,
but x does

Object-oriented programming. Creating new object type ~~extending~~ extending another one.

Class (classes). This is analogous to define a function, but now this allows us to

define new methods, so new . something().

Eg.: Want to create a method ~~that~~ remove_min which removes the minimum elmt of a list:

```
class MyList(list):  
    def remove_min(self):  
        self.remove(min(self))
```

x = [2, 5, 10, 1]

y = MyList(x)

* dir(x) = (all methods for lists)

dir(y) = () + remove_min

y.remove_min()

1 disappeared from y

(7)

Numpy

Provides new data type.

import numpy as np

zv = np.zeros(5) # zero vector of dim 5

zm = np.zeros((5, 3)) # zero matrix of dim 5x3

np.empty

x = np.array([-7, -1])

x[1] = -1

A = np.array([[1, 2], [3, 4]]). A[0, 0] = 9

A.transpose()

A[:, 1] = array([2, 4])

column

A[0, :] = array([1, 2])

$$\boxed{\Delta} : [1, 2] + [3, 4] = [1, 2, 3, 4]$$

$$\text{np.array}([1, 2]) + \text{np.array}([3, 4]) = \text{array}([4, 6]).$$

||
z,

$$z_1 + z_2 = \text{array}([2, 3]).$$

~~BR[0, 1] = 22~~

v = np.array([3, 5, 8, 16, 24]).

v[[1, 2]] = array([5, 7])

list of indices

I can use booleans as array in the following sense:

$z_1 = \text{np.array}([1, 3, 5, 7, 9])$

$z_1 > 6 \rightarrow$ gives me an array $([\text{False}, \text{False}, \text{False}, \text{True}, \text{True}])$

Then I can do

$z_1[\underline{z_1 > 6}] \rightarrow$ gives an array $([7, 9])$, ie, returns an array
only with the True entries

This is
a "logical array"

Warning! Δ : Sliding an array using colon $v[3:5]$ gives a view of
the object, but does not create a new object.

Using new indexing does create a copy:

$z_1 = \text{np.array}([1, 3, 5, 7, 9])$ $w = \text{array}(3, 3, 5)$
 $w = z_1[0:3]$ $z_1 = (3, 3, 5, 7, 9) \quad !!$
 $w[0] = 3$

z_1 as before

$\text{ind} = \text{np.array}([0, 1, 2]) \quad \# \text{ indexes}$

$w = z_1[\text{ind}]$

$w[0] = 3$

$w = (3, 3, 5)$

$z_1 = (1, 3, 5, 7, 9)$

Better: z_1 as before

$w = \text{np.copy}(z_1[0:3])$
 $w[0] = 3$

(8)

`np.linspace(0, 100, 11)` gives an array of 11 points spaced equally between 0 & 100, i.e., $0, 10, 20, \dots$

`np.logspace(1, 4, 4)` gives an array of 4 points spaced "log. equally" bet.
~~10¹ & 10⁴~~, i.e., $10^1, 10^2, 10^3, 10^4$.

$A = np.array([[[123], [456]]])$, $A.shape : (2, 3)$
 $A.size : 6$

$X = np.random.random(10)$ ~~randomly chooses 10 num in [0, 1]~~

`np.any(x > 0.9)` \rightarrow True ~~is any of its elts > 0.9?~~

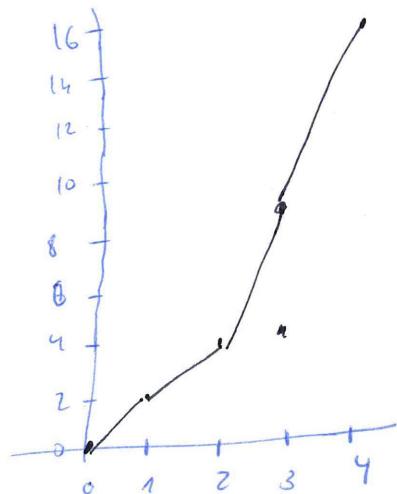
`np.all(x >= 0.1)` \rightarrow ~~All~~ Are all of its elmts ≥ 0.1 ?

Matplotlib & Pyplot

`import matplotlib.pyplot as plt`

`plt.tplot([0, 1, 4, 9, 16])` \rightarrow

(we ; at the end ~~to~~ to not show the object)
~~ie the name of the object~~

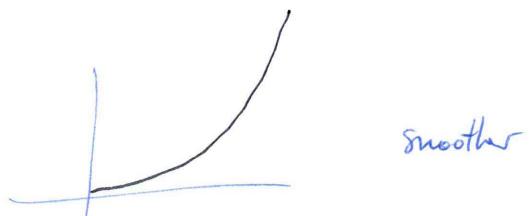


i.e. it plots the piecewise linear function
 passing through $(0,0), (1,1), (2,4), (3,9), (4,16)$

$x = np.linspace(0, 10, 21)$

$y = x**2$

`plt.plot(x, y)` \sim gives



smoother

$x = \text{as before}$

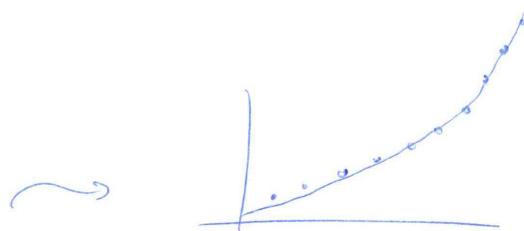
$$y^1 = x^{**} 2.0$$

$$y^2 = x^{**} 1.5$$

`plt.plot(x, y1, "bo-",)`

blue nodes
as circles
solid line

$\circ = \text{circle}$
 $b = \text{blue}$
 $g = \text{green}$
 $r = \text{red}, \text{ etc}$



`plt.plot(x, y1, "bo-", linewidth=2, markersize=4)`

↑
line

↑
nodes



Example:

$x = \text{np.linspace}(0, 10, 20)$

$$y^1 = x^{**} 2.0$$

$$y^2 = x^{**} 1.5$$

`plt.plot(x, y1, "bo-", linewidth=2, markersize=12, label="$x^{3.24}$")`

`plt.plot(x, y2, "gs-", linewidth=2, markersize=12, label="$x^{3x.55}$")`

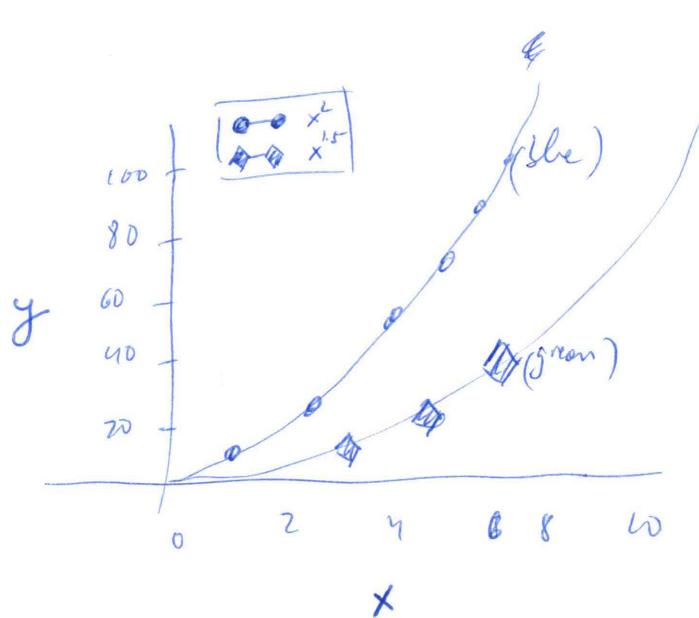
`plt.xlabel("x")`

`plt.ylabel("y")`

`plt.axis([0.5, 10.5, -5, 105])`

`plt.legend(loc="upper left")`

`plt.savefig("myplot.pdf")`



(9)

Remark : plt. loglog (x, y, "b0-") returns a plot with ^{both axes} log scale.

plt. semilogx () ————— with log scale in x axes

plt. semilog y () ————— y

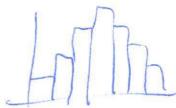
np. logspace (-1, 1, 20) gives log-eq. space pts from 10^{-1} to 10^1 .

Histograms:

x = np. random. normal (size = 1000)

generates an array of 1000 numbers
get from the standard normal distribution
 $N(0, 1)$.

plt. hist (x); # returns



plt. hist (x, normed = True); # instead of plotting the number of observation, it plots
density # the proportion of observations

plt. hist (x, norm = True, bins = np. linspace (-5, 5, 21)) # set differences bkt of
0.5, in 20 ~~buckets~~
bins

x = np. random. gamma (2, 3) $\sim \gamma(2, 3)$

plt. hist (x, bins = 30, cumulative = True)

gives the cumulative histogram,

(histtype = "step")

appearance of the histogram

Can put many plots into one:

plt.figure()

plt.subplot(221) # 22 indicates 2x2

plt.hist(x, ...)

plt.subplot(222)

{

223

224

Random processes

import random

random.choice([0,1]) # returns 0 or 1 randomly

{ or any list, or a range object , eg range(1,7) for 1,...,6 .

random.choice (random.choice ([range(1,7), range(1,9), range(1,11)])) # encodes the situation where you have 3 dice of 6, 8 & 10 faces, you choose one at random and you roll the die.

Example : Roll the dice a number of times and plot the results. (say 100 times)

rolls = []

for k in range(100) :

 rolls.append (random.choice ([1,...6]))

plt.hist(rolls, bins = np.linspace(0.5, 6.5, 7));

Example: What probability distribution follows the random variable consisting of the sum of rolling 10 dices?

$y_s = []$

for rep. in range (500000):

$j = 0$

for k in range (10):

$x = \text{random.choice}([1, \dots, 6])$

$j = j + x$

$y_s.append(j)$

`plt.hist(y_s); # normal distribution, (Central Limit Theorem)`

mp.random.random(5) creates a list w/ 5 numbers chosen from the uniform dist in [0,1]
((5,3)) ————— matrix w/ 5×3

mp.random.normal(0,1) # chooses a number from $N(0,1)$

(0,1,(6,8)) # ~~creates~~ a 6×8 matrix of random numbers.

mp.random.randint(1,7) # chooses a random integer 1, ..., 6.

$A = \text{mp.array}(\text{some matrix})$

np.sum(A) # sum all entries

sum the elements in every column (so over the rows)

np.sum(A, axis=0) # sum the elements in every row (so over the columns)

np.sum(A, axis=1) # sum the elements in every column (so over the rows)

Example (The previous example revisited):

$X = \text{mp.random.randint}(1,7, (100000, 10))$

$Y = \text{np.sum}(X, axis=1)$

`plt.hist(Y)`.

How to measure time?

import time

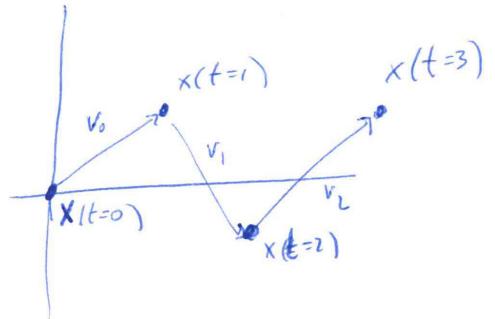
start_time = time.clock()

end_time = time.clock()

total_time = end_time - start_time

Random walk : A walker moves in \mathbb{R}^2 randomly at times $t=0, 1, 2, 3, \dots$.

$$x(t=k) = x(t=0) + \sum v_i.$$



np.cumsum (matrix , axis= 0 or 1) # returns a matrix with partial sum in any column/row

np.concatenate (a1, a2) # concatenate two np arrays

Example: Plot a 5-iteration random walker starting from $(0,0) \in \mathbb{R}^2$:

$X_0 = \text{np.array}([0, 0])$

$\delta X = \text{np.random.normal}(0, 1, (2, 5))$

$X = \text{np.concatenate}((X_0, \text{np.cumsum}(\delta X, \text{axis}=1), \text{axis}=1)$

`plt.plot(X[0], X[1], "r--")`

`plt.savefig("rw.pdf")`

More about strings of characters

- F.read() creates a unique string with all the text, also \n. How to remove those? A copy form is to use F.replace("\n", "").
- ⚠ Since strings are immutable, this does not change F, one has to do F=F.replace(...).
- ⚠ In some computers, there are \r or also \n\r\n. Better do F=F.replace("\r", "") as well.
- If you create a function and want to create a "help" entry, this can be done by creating a "docstring": after def ... : , write """ This is the manual ... """.

Remark: A safer way to read files is

with open("dna.txt", "r") as F:

seg = F.read

Count words from a text

from collections import Counter

text = "hola me llamo Pablo Casado y he perdido cinco claves en un año".

word_counts = Counter(text.split(" ")) # creates a dict with keys the different words and values the frequency

Remark: To turn text into a string with no capital letters or punctuation, do

text = text.lower()

skip = [".", ",", ";", ":", "!", "?", "(", ")"]

for ch in skip: ⚡ text = text.replace(ch, "")

• Read a book and save result

with open("book.txt", "r", encoding="utf8") as current_file:

text = current_file.read()

text = text.replace("\n", "").replace("\r", "")

• Find fragment in a book:

ind = text.find("What's in a name?")

print(text[ind:ind+100]) # prints the 100 char after the start of the sentence.

• Navigate file directories

import os

os.listdir("./Descargas") # list files in ./Descargas

where you already are

what if I have many directories English, German, ... each of them w/ books of different authors
(as new folder)

Shakespeare... and inside those the author .txt's files with the books? Can do

for language in os.listdir("./Books"):

for author in os.listdir("./Books/" + language):

for title in os.listdir("./Books/" + language + "/" + author):

inputfile = "./Books/" + language + "/" + author + "/" + title

(forgot the indentation)

- Create table of data with Pandas

import pandas as pd

table = pd.DataFrame (columns = ("name", "age"))

table.loc[1] = "James", 22

table.loc[2] = "Pinte", 30

Result:		
	name	age
1	James	22
2	Pinte	30

table.head() } # gives the first / last 5 lines
table.tail()

Remark: "hole".capitalize() returns "Hole"

table.name # gives a list with all entries of the column 'name'.

table = pd.read_csv("data.csv", index_col=0) # reads a csv file

More numpy:

• Distance bet. pts: ~~def~~

import numpy as np

p1 = np.array([1,1])

p2 = np.array([4,4])

p2 - p1 \rightarrow array([3,3])

np.power(p2-p1, 2) # does $(3^2, 3^2)$, to return array([9,9])

np.sqrt(np.sum(np.power(p2-p1, 2))).

- Find the mode of a list (most repeated value)

import scipy.stats as ss

votes = [1, 2, 3, 3, 1, 2, 1, 3, 3, 3]

mode, count = ss.mstats.mode(votes)

$\left\{ \begin{array}{l} \text{value} \\ \text{it appears} \end{array} \right\}$
 number of times

Submatrix

Given a np.array A (matrix), one can get a submatrix by doing $A[1:2, 0:2]$ eg.
 the first column is $A[:, 0]$, second column is $A[:, 1]$, etc.

$A.shape[0]$ returns #rows; $A.shape[1]$ returns #columns.

Arrange

np.arange(0, 10, 2, [dtype=int]) #return an array [0, 2, 4, 6, 8].

np.arange(10).reshape(5, 2) gives the matrix $\begin{pmatrix} 0 & 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 & 9 \end{pmatrix}$

np.repeat(m, k) gives an array (m, m, m, \dots, m) K times

Sort (Ordering) a list

a = np.array([-2, 5, 8, -4, 3, 6])

np.sort(a) returns the ordered list

np.argsort(a) returns a list with the indices corresponding to the sorted elements, i.e.

if $a = (x_0, \dots, x_n)$, and $np.argsort(a) = (x_{i_1}, x_{i_3}, x_{i_0}, x_{i_2}, \dots)$ then

np.argsort(a) returns $(i_4, i_3, i_0, i_2, \dots)$

• A list with the 3 largest numbers of a can be got by doing

$\text{ind} = \text{np.argsort}(a)$

~~del~~

$\alpha[\text{ind}[0:2]]$

Example (KNN algorithm): Given two sets of pts $A = \{x_1, \dots, x_m\}$, $B = \{y_1, \dots, y_n\}$, and a pt $p \in \mathbb{R}^n$, a fraction which determines to what set the pt p is the closest, based on the nearest $K \in \mathbb{N}$ points in $A \cup B$.

def distance(p, q): (easy)

def find_nearest_nbhds(p, points, K)

Find the K nearest neighbors and return indices in points

distances = np.zeros(points.shape[0])

for i in range(len(distances)):

 distances[i] = distance(p, points[i])

ind = np.argsort(distances)

return ind[1:K]

def Knn_predict(p, points, outcomes, K=5)

points is supposed to be $A \cup B$ and

ind = find_nearest_nbhds(p, points, K)

outcomes [$\underbrace{0, 0, \dots, 0}_m, \underbrace{1, \dots, 1}_m$]

mode, count = ss.stats.mode(outcomes[ind])

return mode.

- `ss.norm(mu, sigma).rvs(n, m)` returns a $n \times m$ matrix with entries chosen at random $\sim N(\mu, \sigma)$

Enumerate

`seasons = ["spring", "autumn", ...]`

`list(enumerate(seasons))` returns $\{(0, "spring"), (1, "autumn"), \dots\}$

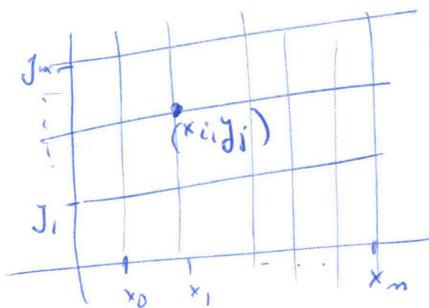
\uparrow
index \uparrow
element

- Meshgrid : Let $x = (x_1, \dots, x_n)$, $y = (y_1, \dots, y_m)$

`xx, yy = np.meshgrid(x, y)` returns two matrices:

$$xx = \left(\begin{matrix} x_i \\ \vdots \\ x_n \end{matrix} \right)_{i=1 \dots n}^{j=1 \dots m}$$

$$yy = \left(\begin{matrix} y_j \\ \vdots \\ y_m \end{matrix} \right)_{i=1 \dots n}^{j=1 \dots m}$$



let $A = \left(\begin{matrix} (x_i, y_j) \\ \vdots \\ (x_n, y_m) \end{matrix} \right)_{i=1 \dots n}^{j=1 \dots m}$

(so xx has repeated rows, yy has repeated columns)

More on Pandas

input pandas as pd

Two main objects: Series & Data Frame

→
 1-dim array-like 2-dim array-like

Series

x = pd.Series([6, 8, 3, 6])

indices	values
1	6
0	6
1	8
2	3
3	6

Can define your own index labels:

x = pd.Series([6, 8, 3, 6], index = ["q", "w", "e", "r"]) →

q	6
w	8
e	3
r	6

Can access to the value by doing x["w"] or x[[{"w"}, {"r}]].

• If age = { "Tom": 50, ... } is a dict, get a pd series by doing

x = pd.Series(age) → Tom 50
 : : .

• Data Frame: Before we saw one way. Can also turn a dict into Data Frame; by using lists in the values.

data = { "name": ["Tom", ...], "age": [50, ...], "ZIP": ["06008", ...] }

x = pd.DataFrame(data, columns = ["name", "age", "ZIP"])

	name	age	ZIP
0	Tom	50	06008
1			
2			
3			

Can recover a single column by doing $x["name"]$ or $x.name$.

- How to ~~view~~ change the indices 0, 1, ...? By doing

$x.reindex(["Person1", ...])$

- If x, y are pd Series of the same dim, then $x+y$ sums the entries w/ the same ~~both~~ labels, and returns NaN if a label appears only in one of x or y .

- Read files: Either is it in csv or a txt with dots separated by , 's, do

$x = pd.read_csv("data.csv")$

To add a column from another file,

$x["new_column"] = pd.read_csv("new-data.txt")$

⚠ $x["new_column"]$ adds a column, $x.loc[7]$ adds a new row

$x.iloc[0:10]$ slices the first 10 rows

$x.iloc[0:10, 1:3]$ ————— and 2nd, 3rd columns.

$x.columns$ # lists the name of the columns.

- Compute correlations

Let x be a pd.DataFrame object w/ columns some attributes of some rows of data.

Would like to compare two such attributes

$\text{corr_x} = \text{pd.DataFrame.corr}(x)$ # compute correlation matrix!

Graphical plotting:

```
import matplotlib.pyplot as plt
plt.figure(figsize=(10,10))
plt.pcolor(x)      # actually plots the corrmat
plt.colorbar()     # the legend for the colors
plt.savefig("...pdf")
```

Remark. $x.\text{transpose}()$ gives the actual transpose of the table.

- Spectral co-clustering: Attempt to find clusters of related data (eg words from books, and find clusters by language).

```
from sklearn.cluster.bicluster import SpectralCoclustering
model = SpectralCoclustering(n_clusters=6, random_state=0)
model.fit(corr_x)
```

model.rows_- # returns a $(n\text{-clusters} \times \# \text{rows in } x)$ -matrix with True/False
 \uparrow
 {the data}

```
np.sum(model.rows == , axis=1) # returns an array with the # of elts in each cluster
```

Examples with Pandas

```
birddata = pd.read_csv("bird_tracking.csv")
```

```
- ix = birddata.bird_name == "Eric" # returns the indices corresponding to the value "Eric" in the column "bird-name".
```

```
x,y = birddata.longitude [ix], birddata.latitude [ix] # returns arrays with the columns long, lat corresponding to ix, i.e Eric.
```

```
plt.figure(figsize=(7,7))
```

```
plt.plot(x,y,".")
```

```
bird_names = pd.unique(birddata.bird_name) # returns an list array w/ the unique entries of the column bird-name (ie without repetition)
```

Example:

```
bird_names (as before)
```

```
for bird in bird_names:
```

```
    ix = birddata.bird_name == bird
```

```
    x,y = birddata.longitude [ix], birddata.latitude [ix]
```

```
    plt.plot(x,y,".")
```

```
plt.xlabel("longitude"); plt.ylabel("latitude")
```

```
plt.legend(loc="lower right")
```

Useful trick : How to ignore the NaN's of a list? (This eg gives error if you want to plot): let x be an array with numbers and some NaN's. Then

```
ind = np.isnan(x)           # ind returns an array with True if a value is NaN
plt.hist(x[~ind])          # ~ind is the complement of ind, so we only plot the "False" values
                           # in ind.
```

- Pandas can also plot histograms :

```
bird_data.speed - 2d. plot (kind = 'hist', range = [0, 30])
plt.xlabel ("2D speed")
```

Good point of this : Do not have to deal with NaN's.

- Handling time with Datetime

```
import datetime
```

```
datetime.datetime.today()      # returns a datetime-type obj: datetime(2020, 6, 18, 11, 43, 50)
                               #           y M D H M S
```

```
time_1 = datetime.datetime()
```

```
time_2 = _____
```

time_2 - time_1 → returns the time it has passed (timedelta object)

• Say we have a string "2013-08-15 00:18:08+00". How to read off the date from this?

```
datetime.datetime.strptime(x[:-3], "%Y-%m-%d %H:%M:%S")
```

• How to add a list to a pandas DataFrame?

`birddata["newcolumn"] = pd.Series(x, index = birddata.index)`

↑ ↑
the existing DDF the list

• How many days have passed? I can divide timedelta objects: ~~divide~~

~~timedelta~~

`time_2 - time_1 / delta datetime.timedelta(days=1)`

Graphs with NetworkX

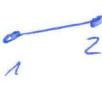
`import networkx as nx`

`g = nx.Graph() # creates an empty graph`

`g.add_node(1) # creates a node labelled w/ 1`

`g.add_nodes_from([2,3,"u","v"])`

`g.nodes() # show the nodes`

`g.add_edge(1,2) # adds a node` 

`g.add_edges_from([(1,5),("u","w")]) # A nodes are added automatically!`

`g.remove_node(2)`

`g.remove_nodes_from([4,5])`

`g.number_of_nodes()`

Plotting a graph

import matplotlib.pyplot as plt

nx.draw(G, with_labels=True, node_color="lightblue", edge_color="grey")

G.degree() # returns a dict-like object w/ the list of nodes & degrees

Erdős - Rényi random graphs: Parameters (N, p): $N = \# \text{nodes}$, $p = \text{probability that there is a graph for every pair of nodes}$

from scipy.stats import bernoulli

bernoulli.rvs(p=0.2) # bernoulli distr, returns 0 or 1 w/ prob p

G = nx.Graph(); N=20; p=0.2

G.add_nodes(range(N))

for nod1 in G.nodes(): # can be omitted

 for nod2 in G.nodes:

 if bernoulli.rvs(p=p) == True:

 G.add_edge(nod1, nod2)

(creates a directed graph, if not directed should add (and $\text{nod1} < \text{nod2}$)

Degree distribution

def deg_distr(G):

 deg_sequence = [d for n, d in G.degree]

 plt.hist(deg_sequence, histtype="step")

 plt.xlabel("..."); plt.ylabel("..."); plt.title("...").

Def: Let G be a graph. The adjacent matrix of G is $A = (a_{ij})$, $m \times m$ matrix,
 $m = \# \text{nodes}$, and $a_{ij} = \begin{cases} 0, & \text{if } \nexists \text{ edge bet. } v_i \text{ & } v_j \\ 1, & \exists \end{cases}$

- * Ofc G can be reconstructed from A . To do this w/ Python

```
g = nx.to_networkx_graph(A)
```

- ## • Connected components of a graph

How to get the largest conn weight? Two ways:

(1) $g = \text{max_connected_component_subgraphs}(G)$. -- next -- ()

(2) $g = \max(\text{max connected components}(G), \text{key} = \text{len})$

Statistics with Python

```
import numpy as np
```

— *scipy as ss*

— matplotlib.pyplot as plt

$x = 10 * \text{ss. uniform.rvs}(\text{size} = n)$

creates an array w/ m values chosen from
the uniform distribution $[0, 10)$.

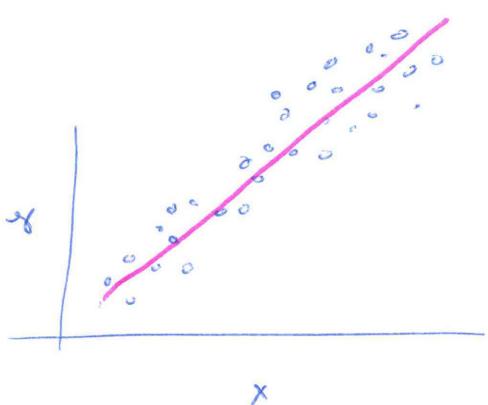
$y = 5 + 2 * x + \text{ss.norm.rvs}(\text{loc} = 0, \text{scale} = 1, \text{size} = 50)$

noise form $\sim N(0, 1)$

Can plot this by doing

```

plt.figure()
plt.plot(x,y, "o", ms=5)
xx = np.array([0,10])
plt.plot(xx, 5+2*xx)
plt.xlabel("x"); plt.ylabel("y")
    }
```



- least squares problem: Given $Y = \beta_0 + \beta_1 X$ and data $\{(x_i, y_i)\}$, want to build $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 \hat{x}$ which minimizes $\sum e_i^2$, where $e_i = \hat{y}_i - y_i$.

```
import statsmodels.api as sm
```

```
mod = sm.OLS(y,x)
```

```
est = mod.fit()
```

```
print (est.summary())
```

```
X = sm.add_constant(x)
```

```
mod = sm.OLS(y,X)
```

```
est = mod.fit()
```

```
print (est.summary())
```

} for a estimation for $Y = \beta \cdot X$

} for $Y = \beta_0 + \beta_1 X$

• scikit-learn for linear regression

$x1 = 10 * \text{ss.uniform.rvs}(\text{size}=500)$

$x2 = \underline{\hspace{10em}}$

$y = 5 + 2 * x1 - 1 * x2 + \text{ss.norm.rvs}(\text{loc}=0, \text{scale}=1, \text{size}=500)$

$X = \text{np.stack}([x1, x2], \text{axis}=1)$ # takes two arrays and stacks them together, i.e., concatenates them, so you get a matrix. In this case as columns ($\text{axis}=1$)

from mpl_toolkits.mplot3d import Axes3D
 $\text{fig} = \text{plt.figure}()$
 $\text{ax} = \text{fig.add_subplot}(111, \text{projection}='3d')$
 $\text{ax.scatter}(X[:, 0], X[:, 1], y, c=y)$

} Produces a 3d graphic

from sklearn.linear_model import LinearRegression

$\text{lm} = \text{LinearRegression}(\text{fit_intercept}=\text{True})$

$\text{lm.fit}(X, y)$

lm.intercept_- # returns β_0

$\text{lm.coef}_- [0]$ # — β_1

$\underline{\hspace{10em}} [1]$ # — β_2

• How to "predict" values?

$X0 = \text{np.array}([2, 4]).\text{reshape}(1, -1)$

$\text{lm.predict}(X0)$

- (Assessing model accuracy) : To evaluate the performance of the regression model, one has to quantify how much the predictions agree w/ the observed data, e.g. mean square error. Out of some data, we want to split it randomly ~~and~~, see how they perform differently and compare results.

```
from sklearn.model_selection import train_test_split
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.5, random_state=1)
```

```
lm = LinearRegression(fit_intercept=True) # for the  $\beta_0$  term
```

```
lm.fit(X_train, y_train) # get the model
```

```
lm.score(X_test, y_test)
```

