

Learn Python

file extension: '.py'

Code (is stored in a variable) used to take input
↓
name = input('Your name?')
print('Hello ' + name)
prints

{ } ↗ user input
Your name? Ankit
Hello Ankit ↗

HOW PYTHON WORKS? { Python gets written as source code. We use this translation service called an interpreter to run this code, and then we'll split out some instructions for our machine, so that the code can be executed.

name = input('Hi') (takes input from user)
print('Hello' + name)
↑
gives output

print(type(6))
returns datatype
(int, float, etc)

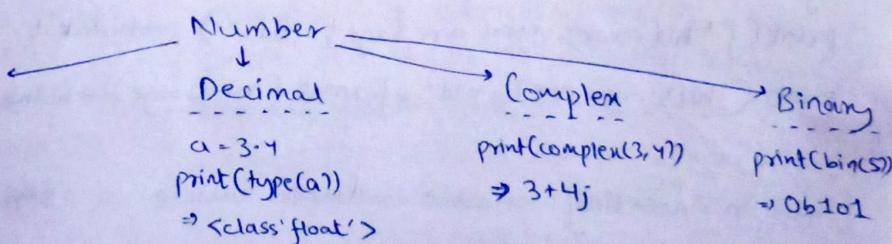
In python, $2^{**5} = 2^5 = 2^5$, $5//4$ = returns rounded off integer on division

$5 \% 3 = 2$ → returns remainder, round(3.9) = rounds off value

abs(-3) = +3 → returns absolute value

{ ex: $5//4 = 1$
{ ex: $5//4 = 1.25$

Operator Precedence: () > ** > * (or) / > + (or) -



Also, (binary → int)

print(int('0b101', 2))

=> 5

(base 10) (base of no. inside brackets)

→ a, b, c = 1, 2, 3

print(a)

=> 1

{ Assigning values to multiple variables at once }

→ print('a' + 'b') # string concatenation

(Comments)

→ print('hi') ↗ (String)
→ long_string = ""
WOW
000
///
print(long_string) } WOW
000
///

∴ We can use " " (or) " "

(or) " "

> (better for longer sentences)

→ str()
↳ converts to string

→ a = str(100)
b = int(a)
print(type(b))
↳ type class 'int'

→ Escape Sequence
weather = "It is \ kind of " sunny
print(weather)
↳ It's "kind of" sunny
take care

→ ex: name = 'Ankit'
0 1 2 3 → (index of string)
print(name[0]) → A
• (name[0:3]) → AnK [start:stop] = [start, stop]
• (name[0:5:2]) → A k t [start:stop:stepover]
• (name[1:]) → nkit prints all values after index 1
• (name[:3]) → Ank prints all values before index 3
• (name[::2]) → A K t steps by 2 (for whole string)
• (name[-1]) → t prints last letter i.e. sequence starts from end of string
• (name[::1]) → t i k n A

↳ t = given tab space
n = moves to next line
↳ it tells, "whatever comes after me, it's a string"

→ Formatting Sequence:

age = 28

print(f"hi{name}, you are {age} old.") → hiAnkit, you are 28 old.

print('hi{.you are {} old.'.format(name, age)) → hiAnkit, you are 28 old.

(helps in inserting variable values inside a string)

* String is immutable i.e. once declared, it can't be changed.

Ex: name[0] = 'Z' → error
print(name) {we can only change full string}
↳ name = 'Yellow'

→ Few features of String

que = 'to be alive'
print(que.upper()) → To BE ALIVE
• (que.capitalize()) → To be alive
• (len(que)) → 11
• que.find('be') → 3

↳ (que.replace('be', 'me'))
↳ to me alive
↳ (que) printed original of
↳ string is immutable
↳ to be alive

→ Boolean

print(bool(0)) → False
bool(1) → True

print(bool('True')) → True
print(bool('False')) → False

→ if: print('*' * 10) → 10 * * * * * * * * *

→ year = input('Birth year = ?')
↳ this input will be taken in string form

→ In python, lists are a form of arrays
↳ (mutable)

(ordered sequence of obj.'s that can be of any datatype)

ex: cart = ['a', 'b', 'c', 'd']
 1 2 3 → index

print(cart[1:3]) → [b, c]

new_cart = cart → new_cart - original one!

new_cart[0] = 'g'

print(new_cart) → [g, b, c, d]

print(cart) → [g, b, c, d]

a copy of cart is made

↳ new_cart = cart [:]

new_cart[0] = 'g'

print(new_cart) → [g, b, c, d]

print(cart) → [a, b, c, d]

(so any changes in new_cart - changes in orig. cart)

→ Matrix ≡ 2D array → ex: matrix = [[1, 2, 1], [0, 1, 0], [1, 0, 1]]
 in np: print(matrix[0][1]) → 2

→ ex: basket = [1, 2, 3, 4, 5].

basket.append(100)

basket.insert(5, 77)

print(basket)

basket.pop()

basket.remove(3)

basket.clear()

These are few methods!

• append(100) → adds 100 in last of list

• insert(5, 77) → adds '77' at index 5

• pop() → removes last element in the list

• remove(3) → removes element whose value is 3

• clear() → removes all elements in the list

If we do → new_list = basket.append(100) → None
print(new_list)

as we are just making changes in place (+e. changes in basket),
it's not returning any value

→ ex: bus = ['a', 'b', 'c', 'd', 'e']

find any ele.
print('x' in bus) → off False

similarly, we can do for string
print('i' in "Yellow Hi") → off True

counting ele.
print(bus.count('d'))

↳ 1

→ Sets are ordered collection of unique items

ex: `s = {1, 2, 3, 4, 5}`
`print(s)` → `{1, 2, 3, 4, 5}`

`s.add(100)` # to add an element '100' in set : `(set_name).add(value)`

We can also convert a list to set using 'set' keyword. Eg vice versa

ex: `list = []` ↳ `print(list(s))`
`print(set(list))`

`print(len(s))` # to find length of set → `5`
`print(1 in s)` # to find '1' in set → `True`

`s1 = s.copy()` ↳ copying set
`s.clear()` ↳ emptying set

`print(s1)` → `{1, 2, 3, 4, 5, 100}`
`print(s)` → `set()`

→ `ms = {1, 2, 3, 4, 5}` # my-set
`ys = {4, 5, 6, 7, 8, 9, 10}` # your-set

`print(ms.difference(ys))` → finds diff. btw both sets
`↳ op: {1, 2, 3}`

`print(ms.discard(5))` → removes that ele. from set
`↳ op: None`

`print(ms)` → `{1, 2, 3}` → since we have discarded 5

`print(ms.difference_update(ys))` → finds diff. both sets & modifies it
`↳ op: None`

`print(ms)` → `{1, 2, 3}`

↳ commented

`print(ms.intersection(ys))` ↳ finds & prints common elements in both
`↳ op: {4, 5}`

`print(ms & ys)`

`print(ms.isdisjoint(ys))` ↳ checks if sets are not overlapping
`↳ op: False` ↳ very diag.

`print(ms.union(ys))` ↳ writing the sets
`↳ op: {1, 2, 3, 4, 5, 6, 7, 8, 9, 10}`

`print(ms | ys)`

ex: `ms = {4, 5}`
`ys = {4, 5, 6, 7, 8, 9, 10}`

`print(ms.issubset(ys))` ↳ checks if 'ms' is subset of 'ys'
`↳ op: True`

`print(ms.issuperset(ys))` ↳ checks if 'ms' is superset of 'ys'
`↳ op: False`

Learn Python Part 2

`a = False` ↳ op: `c3`
`b = True` ↳ done

`if(a and b):`
`... print('c1')`

`elif(a):` ↳ # we can use many
`... print('c2')` ↳ elif's in our code!

`else:`
`... print('c3')`
`print('done')`

If I do: ↳ `a = 'hello'`
`b = 5`

`#(if, elif, stdt!)` ↳ op: `c1`
`done`

because our 'if' basically converts conditions to boolean val

`{... bool('hello') → True, bool(5) → True}`

(similar we have) ↳ (Falsey value) ↳ (Truthy value)

→ Shortcut for if else: "Ternary Operator"

↳ syntax: `condition_if_true if condition else condition_if_false`

ex: `a = False`
`d = "allow" if a else "not allow"`
`print(d)` ↳ op: not allow

Now we can also use 'short-circuiting' in 'if-else'!

ex: `if (a and b)`
`... ↳ T/F`

`False`

`if a = False, then it won't check 'b'` ↳ and
`both`

ex: `if (a or b)`
`... ↳ F/T`

`True`

`if a = True, then it won't check 'b'` ↳ OR
`either (this) or (that)`

→ ASCII of A=65, B=66, a=97, b=98 ; ($a \neq b$ means $a \neq b$)

→ '==' checks equality of value.

ex: `print(True == 1)`
`" ("1" == 1)"`
`" ([] == 1)"`
`" (20 == 20.0)"`
`" ([1, 2, 3] == [1, 2, 3])"`

↳ op: True
`False`
`False`
`True`
`True`

↳ `print([] == [])`
`" ([2, 3] == [2, 3])"`

↳ op: False
`False`

↳ checks if location in memory where value is stored is same or not.
`(these 2 lists are stored in diff. locations)`

→ For loops
for ; in 'Ankit Verma':
 print(i)

can be anything → list / tuple / set
for i in [1, 2, 3, 4, 5]
 print(i)
 print(i)
 print(i)
 print(i)
 ↓
(this points to end = 5)

Nested loops, for i in (1, 2, 3):
 for j in ['a', 'b', 'c']:
 print(i, j)

Iterable simply means 'it is an object/collection that can be iterated over'.
Ex: user = [
 {'name': 'Golem',
 ('one by one check item')
 }]

(list, dictionary, tuple,
set, string
(prints keys))

```
for i in user:  
    print(i)
```

```

• for i in user:
    for i in user.values():
        for i in user.keys():
            print(i)           print(i)           print(i)
            name             Golam             name
            age              500               age

• for i in user.items():
    print(i)
    { ('name', 'Golam'), ('age', 500) }
    { op in tuple form

    (if we don't want commas)

```

```
for key, val in user.items():
    print(key, val)
```

→ print(range(100)) → range(0, 100)

```
for i in range(5):
    print(i)
```

for λ in range(2):

```
print(list(crange(10)))
```

```
→for i,char in enumerate('Hello'): print(i,char)
```

• Enumerate is very useful if we need the index of counter of the item that you're looping through.

```

→ i=0
while i<50:
    print(i)
    i+=1
else: # break
    print("done")

```

0
1
2
3
4a
done

* If 'break' is present, -then
'else' won't work

- * When 'while' loop done, then 'else' will run.

```
→ for item in [1,2,3]:  
    continue  
    print(item)
```

- * "Continue" makes the current iteration stops & moves to next iteration, so restarting for while loop.

→ for item in [1, 2, 3]: {
 pass
 print(item) } → ↗ 1
 2
 3

"pass" basically tells python
to move to next line of code

④ picture = $\begin{bmatrix} 0, 0, 0, 1, 0, 0, 0 \\ 0, 0, 1, 1, 0, 0 \\ 0, 1, 1, 1, 1, 0 \\ 1, 1, 1, 1, 1, 1 \\ 0, 0, 0, 1, 0, 0, 0 \end{bmatrix}$

```

for i in picture:
    for j in i:
        if (j == 0):
            print(' ', end=' ')
        else:
            print('*', end=' ')
    print()

```

The default value of `>end = \n` so everytime code runs, it goes to next

Here we changed from '`\n`' to `"\t"` after end of inner for loop
we print(`\n`) `(end = \n)`!

Q) Check for duplicates in the list : [a', 'b', 'c', 'b', 'd', 'm', 'n', 'o', 'n']
A) some_list = ['a', 'b', 'c', 'b', 'd', 'm', 'n', 'o', 'n']

duplicates = []

```
for i in some_list:
```

```
if some_list.count(i) > 1:
```

if i not in duplicates:) if this line not present
duplicates.append(i) (key word)

Print (duplicated)

Functions

```
def <func.name>: { ex: def hi():
    print('Hello')
}
hi() ~> Hello
```

But these func. can only be used after the declaration line.

Now we can also pass values!
parameters → variables that we receive on def hi(name, emoji): print(f'Hello {name}{emoji}')
Hi('ankit', '@') arguments → actual values we provide to func.

→ def hello(name='Chimp', emoji='2'):
print(f'Hello {name}{emoji}')
hello('Dan', '@') → positional arguments
hello(name='Bi', emoji='+') → keyword argument
hello('Tommy') > default
hello()

→ Returns: def sum(num1, num2):
num1 + num2
print(sum(10,5))

We can also do func. inside func.

```
def sum(num1, num2):
    def add(n1, n2):
        return n1 + n2
    return add(num1, num2)
print(sum(10,5))
```

Methods: (we use dot notation)
ex: "hello".capitalize
"count"
"find"

Docstrings: def test(a):
"""
Info: this func. tests if prints param a
print(a)
test('!!!!')
"""

Now doc strings are helpful, because if you hover over func., you will come to know it's use written in:

We can also get info on user-based func. by help(test) (or) print(test.__doc__)

→ Keyword

ex: def super(*args):
print(args)
print(*args)
return sum(args)
print(super(1,2,3,4,5))
→ (stores in tuple form)

*args → let us grab positional elements

→ def super(*args, **kwargs):
total = 0
print(kwargs)
stores key values
for i in kwargs.values():
 total += i
return sum(args) + total
print(super(1,2,3,4, num1=5, num2=6))
→ (stores in dictionary form)

**kwargs → keyword args allow us to grab any no. of keyword arguments & get a dictionary.

In python, scope follows: Local
↓
Parent Local
↓
Global

If we want to use global func. instead we "global" keyword:
ex: total = 0
def count():
 global total
 total += 1
 return total
count()
print(count())

non-local:

ex: def outer():
 x = "local"
 def inner():
 nonlocal x
 x = "nonlocal"
 print("inner:", x)
 inner()
 print("outer:", x)
outer()

nonlocal helps us refer to the parent variable! So in the next line, we are making changes to parent variable 'x'.

Once func. is called, it's destroyed by garbage collector of Python, so we can't access it anymore!

ex: li = [1, 2, 3]

```
def mult_by_2(item):
    return item*2
print(list(map(mult_by_2, li)))
print(li)
```

(if 'list' not mentioned, then we will get address of resultant map) (no need to mention) ('[]')

eg: [2, 4, 6] > (total, sum = total)

[1, 2, 3]

"map" → we give it some data & it will act upon it. It iterates automatically.

ex: li = [1, 2, 3]

```
def onlyodd(item):
    return item%2 != 0
print(list(filter(onlyodd, li)))
```

(basically filters) (it's going to try to receive a boolean value) (then it checks whether it needs to be filtered or not)

if filtered → doesn't get odd in final list

eg: [1, 3]

ex: li = [1, 2, 3]

li2 = [(10, 20), (30)]

```
print(list(zip(li, li2)))
```

(it takes the first ele., zips it), (if 'li2' has only 2 ele.)

so then goes to next ele.

eg: [(1, 10), (2, 20)]

ex: "reduce" → for this we need to import:

```
from functools import reduce
li = [1, 2, 3]
def accumulator(acc, item):
    print(acc, item)
    return acc + item
print(reduce(accumulator, li, 0))
```

(initial value of 'acc') (doesn't support op in this format)

syntax: reduce(function, sequence, initial) (reduces to single value) (iterable)

eg: 01
22
33
6

List Comprehensions: are used to make list in an easier way

Prev: → my_list = []
for char in 'hello':
 my_list.append(char)
print(my_list)

eg: [h, e, l, l, o]

But using list compre.: → [parameter for parameter in iterable]

ex: li = [char for char in 'hello']

print(li)

[h, e, l, l, o]

(we can name anything)

ex: li2 = [num for num in range(0, 100)]

print(li2)

[0, 1, 2, 3, ..., 98, 99]

If we want numbers to be multiplied by 2.

ex: li3 = [num*2 for num in range(0, 100)]

print(li3)

[0, 2, 4, 6, ..., 196, 198]

If we want $(\text{num})^2$ → with only even no.s

ex: li4 = [num**2 for num in range(0, 100) if num%2 == 0]

print(li4)

[0, 4, 16, ..., 9604]

Set Comprehension: Just change '[]' to '{ }'

ex: li = {c for c in 'tikkiankit'} {eg: {i, t, a, k, n}}

print(li)

Dictionary Comprehension: di = { 'a': 1, 'b': 2 } (grabs whole dictionary)
my_di = { k: v**2 for k, v in di.items() } (key (value))

eg: { 'a': 1, 'b': 2 }

We can also add 'if' at end;

my_di2 = { k: v**2 for k, v in di.items() if v%2 == 0 } {eg: { 'b': 4 }}

print(my_di2)

Suppose we want a dict. of num & ($\text{num} \times 2$)

ex: di = {num: num*2 for num in [1, 2, 3]} {eg: {1: 2, 2: 4, 3: 6}}

print(di)

④ Find duplicates using comprehensions: [a, b, c, b, d, m, n, n]

ex: li = ['a', 'b', 'c', 'b', 'd', 'm', 'n', 'n']
li2 = [char for char in li if li.count(char) > 1]
print(li2)

li3 = set([char for char in li if li.count(char) > 1])
print(li3)

or we change li2 brackets from "[]" to "[]"
Same app.

- Syntax for type of variable in Py: `print(type(x))`
- Syn. to return 1st char of string: `x="Hello"[0]`
- Method to replace parts of string: `replace()`
- List is ordered, changeable & allows duplicate members.
- Multiline comment: `''''''`
- Print length of any string

```
x = "Hello World"
print(len(x))
```
- Get char. from index 2 to index 4
`txt = "Hello World"` `x=txt[2:5]`
- Return String without whitespaces
`x=txt.replace(" ", "")`
- Replace char 'H' with 'J'
`txt = "Hello World"` `txt=txt.replace("H", "J")`
- Insert "lemon" as 2nd item
`fruits = ["apple", "banana"]` `fruits.insert(1, "lemon")`
- Use (-)ve indexing & print last item
`print(fruits[-1])`
- No. of items in a list: `print(len(<list_name>))`
- Add multiple items to a set:
`fruits = {"apple", "mango"}` `more_fruits = {"kiwi", "grapes"}` `fruits.update(more_fruits)`
`print(fruits)` `{'apple', 'grapes', 'mango', 'kiwi'}`
(in any order)
- Print value of "model" in car dictionary, change year to 2020
`car = {"brand": "ford", "model": "Mustang", "year": 1964}`
`print(car.get("model"))`
`car["year"] = 2020`
- Create class → "Hi"
`class Hi:`
`x=5`
- `p1 = Hi()` # Object "P1" of class Hi
`print(p1.x)` # Using object to print value of x
- If we don't know the no. of arguments that will be passed into my func, we can use '*' prefix in func. definition.
`def my_function(*kids):`
- If we don't know the no. of keyword arguments that will be passed into my func, we can use '**' prefix in func. definition.
`def my_function(**hi):`
- Syntax to define "init" func. to a class : or class Person:
`def __init__(self, name, age)`
`self.name = name`
`self.age = age`
- Create class Student that inherits from class Person :
`class Student(Person):`
- Print all variables & func.'s names of "mymodule" module;
`import mymodule`
`print(dir(mymodule))`
- Syntax of importing only person1 dictionary of "mymodule" module;
`from mymodule import person1`
- Execute 'printname' method of obj x : class Person:
`def printname(self):`
`print(self.firstname)`
`class Student(Person):`
`pass`
`x = Student("Mike")`
`x.printname()`
- Modules are nothing but python files. We generally communicate b/w diff files.
`main.py`
`import utility` # this is how we access other files
`utility.mul(4, 5)`
`def mul(num1, num2):`
`return num1 * num2`
- Packages are folders in which .py files are there.
`shopping`
`shopping.cart.py` `def buy(item):`
`__init__.py` `cart = []`
`(It indicates that this is a Python package)` `cart.append(item)`
`return cart`
- Using in main.py
`import shopping.shopping.cart`
`print(shopping.shopping.cart.buy('apple'))`
`['apple']`
- Now we can import in many different ways.
 - from shopping.shopping.cart import buy # importing a func from a package
 - from shopping import shopping.cart # importing module from a package
 - print(buy('apple'))
 - print(shopping.cart.buy('apple'))
 - if there is a func. defined in a package & it's a keyword with diff use, then it gets overwritten if we call it. (like max)
 - Keyword → func. name(x)
 - user-defined use

Data Science

Environment

Setup

A collection of tools/packages is called environment.

We use "Conda" to create the environment & then use it to install different packages. It also allows us to share.

.ipynb = python notebook
 * runs in order of execution of cells

Example Heart Disease Project

import pandas as pd
 (becomes its shorthand)

df = pd.read_csv("heart-disease.csv")
dataframe
 ↳ (comma separated values)
 ↳ (press "tab" for autofill)

df.head() → shows first 5 rows of table
 ↳ if we write any no. here, it will show those many rows of table

import matplotlib.pyplot as plt ↳ (to make graphs)

df.target.value_counts().plot(xkind="bar")
 ↳ counts no. of occurrences of diff. values in that column
 ↳ means create a plot
 ↳ means "bar graph"

if we want to upload an image
(markdown mode) (type image name) ↳ now shift+enter

![] (6-step ml-framework.png)

ls { it works in the same way as
if we did 'ls' in command prompt. It shows us all files

In order to close the notebook → go to Jupyter notebook

(press "ctrl+c" to exit)
 ↳ you can see all updates which you did

Environment

Setup

Pandas: Data Analysis

Data Science (DS)

Data Analysis (DA)

Machine Learning (ML)

Data Modelling

1) Problem Definition

2) Data

3) Evaluation

4) Features

5) Modelling

6) Experiment

DA → pandas, matplotlib, numpy

ML → tensorflow, Pytorch, scikit learn, CatBoost, dmlc XGBoost

DS → jupyter, anaconda, miniconda

In this section we will learn:
→ Most useful funcs. → pandas Datatypes
→ Importing & exporting data → Describing data
→ Viewing & selecting data → Manipulating data

How to open Jupyter Notebook:

→ Open "Jupyter" → click on its link

• Series = 1-D
 ↳ takes "List" as input
 ↳ 1 → to enter heading

• Dataframe = 2-D
 ↳ takes "dictionary" as input
 ↳ column name

ex: check1 = pd.Series(["Hi", "Bye"])
check1
 ↳ o Hi
 ↳ 1 Bye
 ↳ dtype: object

• Importing data:
car_sales = pd.read_csv("car-sales.csv")
car_sales
 ↳ Keyword
 ↳ file name
 ↳ prints the csv data

--- ≈ CHANGES MADE
• Exporting data frame:
car_sales.to_csv("exported-car-sales.csv", index=False)
new_sales = pd.read_csv("exported-car-sales.csv")
new_sales
 ↳ file name to be read

* We can also read by using (---), github link! car("data/car.csv")
 ↳ file location if in some directory

DESCRIBING DATA

• Attributes → no brackets

 ↳ it's just some meta info which is stored in our file header

Functions → brackets there

 ↳ performs some steps

so we have "car_sales.csv" file!

④ (Attribute)

- car_sales.dtypes: returns columns & datatype of values under it
- car_sales.columns: returns column name
- car_sales.index: returns start, stop, step

⑤ (Function)

- car_sales.describe(): return info like: mean, min of numeric column
- car_sales.info(): it's basically ((index) + (dtypes))
- car_sales.mean(): returns mean value of numeric columns
- car_sales.sum(): adds the numeric values
↳ appends the string values
- car_sales["Doors"].sum(): returns 'sum' of specific column named "Doors".

- len(car_sales): returns length of our data

- car_prices = pd.Series([3000, 1500, 111256])

car_prices.mean()

↳ Ans: 38583.336

⑥ SELECTING AND VIEWING DATA

- car_sales.head(): returns first 5 values

↳ if we write any no. here, then it returns those many values from top/bottom.

- car_sales.tail(): returns last 5 values

"loc" → refers to index no.
"iloc" → refers to position

→ animals = pd.Series(["cat", "dog", "bird", "panda", "snake"], index=[0, 3, 9, 8, 7], dtype="object")

If we don't want custom index, no need to mention it

→ animals.loc[3]

↳ 3 dog

→ animals.iloc[9] {→ animals.loc[3]}
↳ 9 bird

→ animals.loc[7] {→ animals.loc[3]}
↳ 7 panda

→ animals.loc[:3] → 0 1 2

0 cat
1 dog
2 bird

→ car_sales.loc[:3]

↳ calling head with val. u.

If we want a particular column!

→ car_sales["Make"] → ① OR car_sales.Make → ②

Column name

Column name

In ① → if we have spaces in column, it will run ✓
② → if we have spaces in column, it will not run ✗

→ car_sales[car_sales["Make"] == "Toyota"] → boolean indexing

↳ ↳ returns all values whose column name = Toyota

→ car_sales[car_sales["Odometer (km)"] > 100000] ↳ selects rows with over 100000 on odometer

→ pd.crosstab(car_sales.Make, car_sales.Doors)

↳ ↳ makes a dataframe of 'Make' & 'Doors'

→ car_sales.groupby(["Make"]).mean() # grouping → similar to crosstab but focusing on more columns at a time

↳ ↳ print mean of diff. columns (using "Make" as reference)
(only numeric)

→ car_sales["Odometer (km)"].plot()

↳ ↳ plots line graph of odometer!

→ car_sales[["Odometer (km)"]].hist()

↳ ↳ plots histogram

Now we have "Price" column in our car-sales.csv! It's dtype = "object"
↳ ↳ (\$4500.00)

Category

Now to convert string to integer!

→ Syntax → dataframe["amount"] = dataframe["amount"].str.replace("\$", "", regex=False)

→ car_sales["price"] = car_sales["Price"].str.replace("\$", "", regex=True).astype(int)

car_sales

↳ ↳ changes price values → (450000)

access car_sales["Price"] → we want to keep it as it is but access the string & replace the dollar sign, commas & dots with nothing
& empty string & then change it from string to integer

Now to ignore decimal values → car_sales["Price"] = car_sales["Price"].apply(lambda x: int(x[0:3]))

↳ ↳ 450000

MANIPULATING DATA

`car_sales["Make"].str.lower()` : returns lowered alphabets/words of column "Make"

But if we print `"car_sales"`, our lowered words won't be visible;
change(so to change our original table, we can reassign it)

`car_sales["Make"] = car_sales["Make"].str.lower()`

Now our original table, make column changed to lowered words.

New file \rightarrow car_sales_missing (has few blocks with no values)

(shown with 'NaN' in pandas)

If we want to fill those empty values, one way is reassigning

2nd way:
`car_sales_missing["Odometer"].fillna(car_sales_missing["Odometer"].mean(), inplace=True)`
helps us make changes in copy table without reassigning by default = False
returns the table with 'NaN' values now filled with mean of Odometer

`car_sales_missing.dropna()` : returns only those rows which doesn't have 'NaN' values

If we have to reaccess our prev. data after 'inplacing', we have to import again.

`car_sales_missing_dropped = car_sales_missing.dropna()
car_sales_missing_dropped`

returns fully filled rows!

Now exporting it!

`car_sales_missing_dropped.to_csv("car-sales-missing-dropped.csv")`

Making column from series

`seats_column = pd.Series([5,5,5,5,5])`

`car_sales["Seats"] = seats_column`

fill first '5' values with '5'
rest all 'NaN'

`car_sales.fillna(5, inplace=True)` \rightarrow fills all values by 5

Making column from Python List

`fuel_economy = [1.5, 9.2, 5.0, 9.6, 8.7, 4.5, 6.0, 1.2, 2.0, 3.0]`

`car_sales["Fuel per 100km"] = fuel_economy`

list size = datafram size
New column name

Making column from another column

`car_sales["Total fuel(L)"] = car_sales["Odometer (km)"] / 100 * car_sales["Fuel per 100km"]`

(column name)

Creating column from a single value

`car_sales["Num of wheels"] = 4` \rightarrow creates column "Num of wheels" with all its value = 4

Now if we want to remove a column,

`car_sales.drop("Total fuel used", axis=1, inplace=True)`

column name
which is to be removed

If we want to randomize our sample input:

`car_sales_shuffle = car_sales.sample(frac=1)` '1' means shuffle 100% of car_sales_shuffle

{randomizes rows}
{prints table} if we want to make changes in "car_sales" only, we have to do "inplace"!

If we want to re-order/reset our index;

`car_sales_shuffle.reset_index(drop=True, inplace=True)`

for permanent changes
if we didn't mention this, then it will give new index for those random indices, so to avoid it, we write!

1 mile = 1.6 km

Now if we want to change km to miles in Odometer:

`car_sales["Odometer (km)"] = car_sales["Odometer (km)"].apply(lambda x: x/1.6)`

{prints table with odometer values divided by 1.6}

So it says: apply \rightarrow reset $x = x/1.6$

apply() \rightarrow lets you apply some kind of func. whether it be numpy or lambda to a certain column

lambda : is a keyword in python which is short for anonymous func.

Now if want to rename a column,

`car_sales = car_sales.rename(columns = {"Diameter (cm)": "Diameter (Mile)"})`

Imp point:

Previously - when we changed 'Price'! (rewriting in better way)

(Price → \$4000.00)

#remove punctuation from price column

`car_sales["Price"] = car_sales["Price"].str.replace("[\$,\.,\,]", "")`

#checking changes in Price column

`car_sales["Price"]` → 40000

#removing the 2 extra zeroes at the end

`car_sales["Price"] = car_sales["Price"].str[:-2]`

#checking changes

`car_sales["Price"]`

→ 4000

→ changing datatype of Price column to integers

`car_sales["Price"] = car_sales["Price"].astype(int)`

NumPy

Numerical Python

In this section, we will learn:

- Most useful func.'s
- NumPy datatypes & attributes (ndarray)
- Creating arrays
- Viewing arrays
- Manipulating & comparing arrays
- Sorting arrays
- Use cases

DATATYPES AND ATTRIBUTES

NumPy's main datatype is "ndarray" ~ N-dimensional array

`a1 = np.array([1, 2, 3])`, `a2 = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])`

a2 → float becomes 16-bit datatype by default
→ returns above array

`a3 = np.array([[[1, 2, 3], [4, 5, 6], [7, 8, 9]], [[10, 11, 12], [13, 14, 15], [16, 17, 18]]])`

→ a1.shape, a2.shape, a3.shape → ((3,), (2, 3), (2, 3, 3))
→ keeping common will also print our ans. in command

→ a1.ndim, a2.ndim, a3.ndim → (1, 2, 3) → no. of dimensions

→ a1.dtype, a2.dtype : {int, float} → returns datatype of arrays

→ a1.size, a2.size, a3.size : → (3, 6, 18) → returns no. of elements in array

→ type(a1), type(a2), type(a3) : {`numpy.ndarray`}

.Creating a dataframe from NumPy array.

`import pandas as pd`

`df = pd.DataFrame(a2)`

`df`

0	1	2
2	3	4
1	5	6

.Anatomy of NumPy array:

Data	<code>a1</code>	<code>a2</code>	<code>a3</code>
NumPy	<code>[[1, 2, 3]]</code>	<code>[[[1, 2, 3], [4, 5, 6], [7, 8, 9]]]</code>	<code>[[[1, 2, 3], [4, 5, 6], [7, 8, 9]]]</code>
Details	<ul style="list-style-type: none"> ↳ Name: array, vector ↳ 1D ↳ Shape = (1, 3) 	<ul style="list-style-type: none"> ↳ Name: array, matrix ↳ > 1D ↳ Shape = (2, 3) 	<ul style="list-style-type: none"> ↳ Name: array, matrix ↳ > 1D ↳ Shape = (2, 3, 3)

CREATING ARRAYS

→ np.{func_name}

→ (press 'Shift + tab' → opens docstring which gives us info. about that func.)

→ ones = `np.ones((2, 2))`
ones → returns new array of given shape & type, filled with ones

→ zeros = `np.zeros((2, 3))`
zeros → returns new array of given shape & type, filled with zeros

→ np.arange(0, 10, 2)
→ returns evenly spaced values with given interval
→ e.g. ([0, 2, 4, 6, 8]) → (not saved in anything)

Now making random arrays!

→ random_array = `np.random.rand(3, 5)`

random_array → (returns random integers from low to high)

→ 3x5 matrix with random values between 0 to 10

$\rightarrow \text{random_array} = \text{np.random.random}(\text{size}=(3, 3))$
 $\rightarrow \text{np.random.random}((3, 3))$
 ↳ returns random floats in half open interval $[0, 1)$
 $\{ \text{if random } 5 \times 3 \text{ matrix with random values b/w } [0.0 \text{ & } 1.0]$
 $\rightarrow \text{np.random.rand}(5, 3) \rightarrow \text{np.random.randn}(5, 3)$
 $\{ \text{creates an array of normally distributed numbers size } 5 \times 3$
 Now whatever the above op is generated is not random but pseudo random!
 $\rightarrow \text{np.random.seed(Seed=0)}$ ↳ any no.
 $\{ \text{for diff. seed, diff. random no.s will be saved}$
 $\text{random_array_4} = \text{np.random.randint}(0, \text{size}=(2, 3))$
 random_array_4
 $\{ \text{shape } ((5, 0, 3), [3, 7, 9])$
 $\} \{ \text{Now no matter how many times we run this code, output will be same unless until seed value is changed}$
 So those random no.s generated are basically having seed = (random value)
 You can run the code in any device, output will be same!

VIEWING ARRAYS AND MATRICES

If we want unique no.s in an array:

$\rightarrow \text{np.unique}(\text{random_array_4}) \rightarrow \text{np.array}([0, 3, 5, 7, 9])$

We can also use slicing concept to view arrays!

$\{ \text{a3}[:, :, 2] \# \text{means 1st matrix} \rightarrow 2 \text{ rows + 2 column elements}$
 $\{ \text{shape } [[1, 2], [4, 5]]$

$\{ \text{a3}[:, :, :1] \# \text{means both matrix} \rightarrow \text{all rows + 1st column elements}$
 $\{ \text{shape } [[1], [4], [7], [10], [13], [16]]$

Furthest element on the right (or furthest number of the shape gets displayed on the innermost section.)

MANIPULATING AND COMPARING ARRAYS

Arithmetic

- $\rightarrow \text{a1} + \text{ones}$ ↳ $\text{a1} + \text{np.ones}(2, 3, 1)$
- $\rightarrow \text{a1} - \text{ones}$ ↳ $\text{a1} - \text{np.ones}(2, 3, 1)$
- $\rightarrow \text{a1} * \text{ones}$ ↳ $\text{a1} * \text{np.ones}(2, 3, 1)$
- $\rightarrow \text{a1} * \text{a2}$ ↳ $\text{a1} * \text{a2}$
- $\rightarrow \text{a2} / \text{a1}$ ↳ floor division - removes decimals & rounds down
- $\rightarrow \text{a2} ** 2$ ↳ power 2

Above are simple python operations, but we can do numpy operations too

$\rightarrow \text{np.add(a1, ones)} \rightarrow \text{np.square(a2)} = \text{a2} ** 2$

$\rightarrow \text{np.exp(a1)} \rightarrow \text{np.log(a1)}$

$\{ \text{exp}, \text{exponential}$
 $\{ \text{of elements}$
 $\{ \text{log of elements}$

Aggregation = performing same operation on number of things!

$\rightarrow \text{lis} = [1, 2, 3]$
 $\text{type(lis)} \rightarrow \text{list}$
 $\text{sum(lis)} \rightarrow \text{int}$

$\{ \text{a1} = \text{np.array}([1, 2, 3])$
 $\text{sum(a1)} \rightarrow \text{int}$
 $\text{np.sum(a1)} \rightarrow \text{int}$

Use Python's method (sum()) on Python datatypes so NumPy's methods on NumPy arrays (np.sum())

{Magic func. → has % symbol}

Ex:
 $\# \text{Creating a massive array}$
 $\text{massive_array} = \text{np.random.random}(100000)$
 $\text{massive_array}[:, :10]$ ↳ Shows us first 10 values (random values b/w 0 & 1)

%timeit sum(massive_array) ↳ Python's sum()
%timeit np.sum(massive_array) ↳ NumPy's np.sum()

$\{ \text{avg } 7.79 \text{ ms} \pm 986 \text{ ys per loop (mean of 7 runs, 100 loops each)}$
 $\{ \text{avg } 411.2 \mu\text{s} \pm 3.98 \mu\text{s per loop (mean of 7 runs, 1000 loops each)}$

∴ We can see "np.sum()" is a lot faster than "sum()"

$\rightarrow \text{a2}$
 $\{ \text{shape } [[1, 2, 3, 3], [4, 5, 6, 5]]$
 $\rightarrow \text{np.mean(a2)} : \text{np: } 3.6333$
 $\rightarrow \text{np.max(a2)} : \text{np: } 6.5$
 $\rightarrow \text{np.min(a2)} : \text{np: } 1.0$

$\rightarrow \text{np.std(a2)} : \text{np: } 1.82269$
 $\rightarrow \text{np.var(a2)} : \text{np: } 3.32222$
 $\rightarrow \text{np.sqrt(np.var(a2))} : \text{np: } 1.82269$

Standard deviation (std) = measure of how spread out of a group of numbers is mean

Variance (var) = measure of the average degree to which each number is different to the mean

High variance: wide range of numbers
Low variance: lower range of numbers

Standard deviation = square root of variance = $\sqrt{\text{Variance}}$ (sqrt)

Understanding std & var:

high_var_array = np.array([3, 100, 200, 300, 4000, 5000])
low_var_array = np.array([2, 4, 6, 8, 10])

→ np.var(high_var_array), np.var(low_var_array) : op: (4296133.492..., 2.0)

→ np.std(high_var_array), np.std(low_var_array) : op: (2072.71162..., 2.87842)

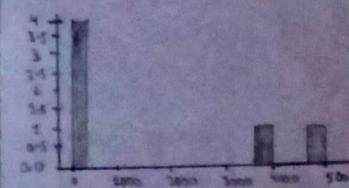
→ np.mean(high_var_array), np.mean(low_var_array) : op: (1600.1666..., 6.0)

! matplotlib inline

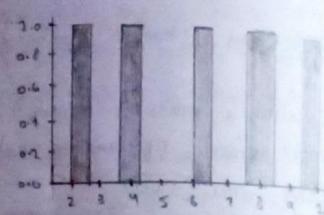
import matplotlib.pyplot as plt

→ plt.hist(high_var_array)

plt.show()



→ plt.hist(low_var_array)
plt.show()



Reshaping And Transposing:

→ a2 * a3 : op: error

→ a2.shape, a3.shape : op: ((2,3),(2,3,3))

→ a2_new = a2.reshape(2,3,1)

↳ reshapes the dimension

Now we can multiply!

→ a2_new * a3 : op: ✓

transpose : switches the axis

→ a2.T : op: ((2,2,1), [2, 5, 3], [1, 3, 6, 5]))

• General Broadcasting Rules:

- 2 dimensions are compatible when:
 - ① They're equal (or)
 - ② One of them is 1

Element - wise:

$$\begin{matrix} A & B \\ C & D \end{matrix} \times \begin{matrix} E & F \\ G & H \end{matrix} = \begin{matrix} A+E & B+F \\ C+G & D+H \end{matrix}$$

Dot - product:

$$\begin{matrix} A & B & C \\ D & E & F \\ G & H & I \end{matrix} \cdot \text{dot} \cdot \begin{matrix} I & K \\ L & M \\ N & O \end{matrix} = \begin{matrix} A*I + B*L + C*N & A*K + B*M + C*O \\ D*I + E*L + F*N & D*K + E*M + F*O \\ G*I + H*L + I*N & G*K + H*M + I*O \end{matrix}$$

Let mat1, mat2 = 2 matrices

→ mat1.shape, mat2.shape = ((5,3), (5,3))

→ mat1 * mat2 : op: ✓ → element wise multiplication

→ np.dot(mat1, mat2) : op: error → as shapes not aligned
↳ dot product

This is where we use "Transpose"!

→ mat3 = np.dot(mat1, mat2.T)

→ mat3 : op: (5, 5)

→ mat3 : op: ✓

Nut Butter Sales (example → dot product) → refer jupyter for full code, here only imp. ones are there!

↳ Weekly_sales = pd.DataFrame(sales_amount, index=['Mon', 'Tue', 'Wed', 'Thu', 'Fri'], columns=['Almond', 'Peanut', 'Cashew'])

In NumPy, we can use same comparison operators which are in Python

→ a1 → [1, 2, 3] & a2 → [[1, 2, 3], [4, 5, 6, 5]]

→ a1 > a2 : op: Array([[False, False, False], [False, False, False]])

→ Similarly: bool_arr = a1 >= a2

→ type(bool_arr), bool_arr.dtype
↳ (numpy.ndarray, dtype('bool'))

→ a1 >= a2 : op: [[True, True, True], [False, False, False]]

• Sorting arrays:

random_array = np.random.randint(10, size=(3, 5))
random_array

↳ op: [[7, 8, 2, 5, 9], [8, 9, 4, 3, 0], [3, 5, 0, 2, 1]]

→ np.argsort(random_array)

↳ op: [[1, 5, 7, 8, 9], [0, 3, 4, 8, 9], [0, 2, 3, 3, 5]]

→ a1 = [1, 2, 3]

→ np.argsort(a1)

↳ op: [0, 1, 2]

→ np.argmin(a1)

↳ op: 0 (index of minimum number)

→ np.argmax(a1)

↳ op: 2 (index of maximum number)

(sorts the indexes)

→ np.argsort(random_array)

↳ op: [[2, 3, 0, 2, 4], [4, 3, 2, 0, 1], [2, 3, 0, 4, 1]]
(columns → 12 R)

→ np.argmax(random_array, axis=0)

↳ op: [1, 2, 0, 0]

→ np.argmax(random_array, axis=1)

↳ op: [4, 1, 1]

(rows → 72 B)

• Practical example:

```
from IPython.display import Image, display
# Specify path to image
image_path = "numpy-images/panda.png"
# Display the image
display(Image(image_path))
```



display image

Turn image into numpy array

```
from matplotlib.image import imread
panda = imread("numpy-images/panda.png")
print(type(panda))
# <class 'numpy.ndarray'>
```

→ panda.size, panda.shape, panda.ndim # if we zoom in, we see pixels
 → (24465000, (2330, 3500, 3), 3)
 → np.linspace(1, 100, 23) → creates an array with 23 evenly spaced no.'s between 1 to 100
 Start stop total evenly spaced no's

Matplotlib : Plotting And Data Visualization

In this section, we will learn:

- Matplotlib workflow
- Importing matplotlib & the 2 ways of plotting
- Plotting data from Numpy arrays
- Plotting data from pandas Dataframes
- Customizing plots
- Saving and sharing plots

A matplot lib workflow:

- a) Create Data
- b) Create plot (figure)
- c) Plot data (axes on figure)
- d) Save/Share plot

• Creating empty graphs:

→ plt.plot() → plt.plot();

→ x = [1, 2, 3, 4]
 y = [11, 22, 33, 44]
 plt.plot(x, y);

→ fig = plt.figure()

ax = fig.add_axes([1, 1, 1, 1])
 ax.plot(x, y)
 plt.show()

→ plt.plot()
 plt.show()

creates figure
 → fig = plt.figure()
 ax = fig.add_subplot()
 plt.show()

adds some axes

• Recommended method to create a figure:

→ fig, ax = plt.subplots() (we can also replace this with a list)
 ax.plot(x, y) → adding data
 type(fig), type(ax)

→ fig (matplotlib.figure.Figure), matplotlib.axes._subplots.AxesSubplot

• Example Workflow

→ Step 0: Import matplotlib & get it ready for plotting in Jupyter

→ %matplotlib inline
 import matplotlib.pyplot as plt

Step 1: Prepare data

x = [1, 2, 3, 4]
 y = [11, 22, 33, 44]

Step 2: Setup plot width height → figure

fig, ax = plt.subplots(figsize=(5, 5))

Step 3: Plot data

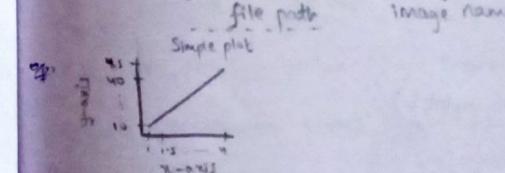
ax.plot(x, y) → we are adding data to our axes, above we just created them

Step 4: Customize plot

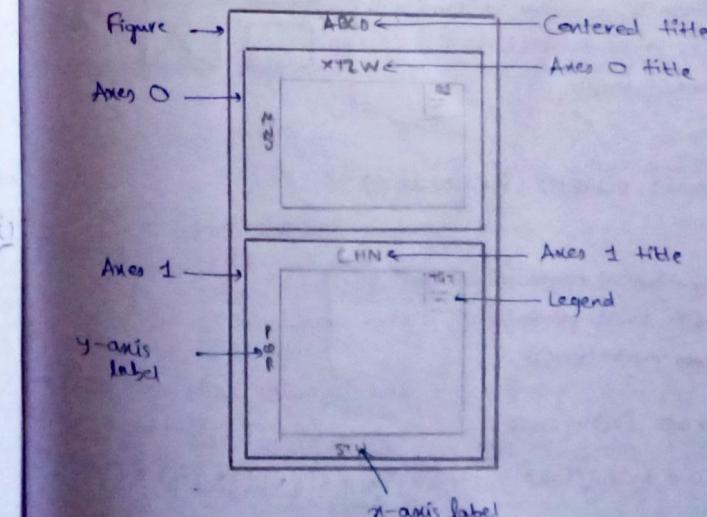
ax.set(title="Simple Plot", xlabel="x-axis", ylabel="y-axis")

Step 5: Save & show

fig.savefig("File Images/sample-plot.png")

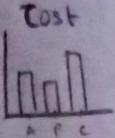


• Anatomy of a Matplotlib plot



Making figures with NumPy arrays

- Creating some data
- $x = np.linspace(0, 10, 100)$
- $fig, ax = plt.subplots()$
- $ax.plot(x, x**2)$
- (Line plot by default)
- Making plot from dictionary:
- $prices = {"A": 10, "B": 8, "C": 12}$
- $fig, ax = plt.subplots()$
- $ax.bar(prices.keys(), prices.values())$
- $ax.set(title="Cost", ylabel="Price")$
- (not necessary)



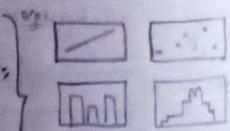
If we want horizontal bar graph, we have to pass "dict" as list!

- $fig, ax = plt.subplots()$
- $ax.barh(list(prices.keys()), list(prices.values()))$
- $x = np.random.randint(1000) \rightarrow$ picks no. that lies in standard normal distribution curve
- $fig, ax = plt.subplots()$
- $ax.hist(x)$

Now, there are 2 ways to create subplots;

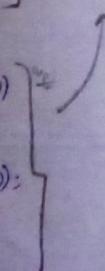
Subplot option 1

- $fig, ((ax1, ax2), (ax3, ax4)) = plt.subplots(nrows=2, ncols=2, figsize=(10, 5))$
- # Plot to each diff. axis
- $ax1.plot(x/2)$
- $ax2.scatter(np.random.random(10), np.random.random(10))$
- $ax3.bar(prices.keys(), prices.values())$
- $ax4.hist(np.random.randint(1000))$



Subplot option 2

- $fig, ax = plt.subplots(nrows=2, ncols=2, figsize=(10, 5))$
- # Plot to each diff. axis
- $ax[0, 0].plot(x/2)$
- $ax[0, 1].scatter(np.random.random(10), np.random.random(10))$
- $ax[1, 0].bar(prices.keys(), prices.values())$
- $ax[1, 1].hist(np.random.randint(1000))$



Plotting from pandas Dataframe

- $car_sales["Price"] = car_sales["Price"].str.replace("[\$, \u20ac]", "")$
- car_sales
- (e.g. \$4500.00 → 4500.00)
- \$500 → 500.00
- replaces "\$", "\u20ac", "," with "" (empty string)

car_sales["Price"] = car_sales["Price"].str[-2:] removes last 2 digits

car_sales

400000 → 4000
550000 → 5000

car_sales["Sale Date"] = pd.date_range("1/3/2023", periods=len(car_sales))

car_sales

Sale Date
2023-01-01
2023-01-02

we convert to int.
else it will convert string.

car_sales["Total Sales"] = car_sales["Price"].astype(int).cumsum()

car_sales

Total Sales
4000
9000
4000 + 5000 + 5000

Cumulative Sales

it adds current value to prev. value

`# type(car_sales["Price"][0]) : str`

Scatter plot

car_sales["Price"] = car_sales["Price"].astype(int)

car_sales.plot(x="Odometer (km)", y="Price", kind="scatter")

Line plot

car_sales.plot(x="Sale Date", y="Total Sales")

Bar graph

$x = np.random.rand(10, 4) \neq$ 10 rows, 4 columns

df = pd.DataFrame(x, columns=["a", "b", "c", "d"])

df.plot.bar(), # or df.plot(kind="bar");

car_sales.plot(x="Make", y="Odometer (km)", kind="bar");

Hist

car_sales["Odometer (km)"].plot.hist(bins=20); # plot(kind="hist");

default - 10
(need to mention)
basically - the no. of columns
our histogram is divided into

Now trying for another dataset heart_disease → has 14 columns

heart_disease.plot.hist(figsize=(10, 30), subplots=True)

Which one should we use? → when plotting something quickly, okay to use pyplot method (.plot())

When plotting something more advanced, use the OO method
(object oriented)

Creating a sub frame

over_50 = heart_disease[heart_disease["age"] > 50]

Pyplot method

over50.plot(kind="scatter", x="age", y="chol", c="target");

gives us the "scale"
Kinda like the scatterplot

```

# If plt + 00
→ fig, ax = plt.subplots(figsize=(10,6))
over_50 = plot(kind='scatter', x='age', y='chol', c='target', ax=ax);
over_50.set_xlim([45, 100]); # changes 'x' range to 100
low    high

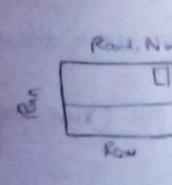
if we don't write → R: (50-75)
if we write → R: (50-100)

# 00 from scratch
# subplot of chol, age, thalach
→ fig, (ax0, ax1) = plt.subplots(nrows=2, ncols=1, figsize=(10,10), sharex=True)
# add data to ax0
scatter = ax0.scatter(x=over_50['age'], y=over_50['chol'], c=over_50['target'])
# customize ax0
ax0.set(title='Heart Level', ylabel='Cholesterol');
# add a legend to ax0
ax0.legend(scatter.legend_elements(), title='Target')
# add a meanline
ax0.axhline(y=over_50['chol'].mean(), linestyle='--') ← go to 'scatter' variable & find all legend elements
# add a horizontal line ← y-value
# add data to ax1
scatter = ax1.scatter(x=over_50['age'], y=over_50['thalach'], c=over_50['target'])
# customize ax1
ax1.set(title='H and MH Rate', xlabel='Age', ylabel='MH Rate');
# add a legend to ax1
ax1.legend(scatter.legend_elements(), title='Target')
# add a meanline
ax1.axhline(y=over_50['thalach'].mean(), linestyle='--') ← go to 'scatter' variable & find all legend elements
# add title to figure
fig.suptitle("HD Analysis", fontsize=16, fontweight="bold")

```

→ if we don't need in subplot
→ fig, ax = plt.subplots(figsize=(10,6))
| # not all same. replace ax0->ax

- Customizing Matplotlib plots so getting stylish:
→ plt.style.available → shows different styles available
→ plt.style.use('dark_background') → updates the plot internally
car_sales[['Price']].plot() → graph plotted with 'dark_background' style
→ x=np.random.randint(10,40) # Create some data
df = pd.DataFrame(x, columns=['a', 'b', 'c', 'd'])
ax = df.plot(kind='bar')
ax.set(title='Rand. Num', xlabel='None', ylabel='Rand.') ← make legends visible



Customizing our full-fledged heart-disease plot!
→ plt.style.use('dark_background')

Now a small change: if we want to change 'plotted line/dot' colour
| exits
| scatter = ax.scatter(x=over_50['age'], y=over_50['chol'], c=over_50['target'], cmap='winter') ← colour map
| (cmap)
| (winter)
| (cmap)

Now to save image/graph, we have 2 options:
| Right click → 'Save image as'
| in code
| fig.savefig("heart-disease-analysis-plot.png") ← gets saved where your jupyter notebook is saved!

Resetting the figure.

→ fig, ax = plt.subplots();

Scikit-Learn: Creating Machine Learning Models

In this section, we will learn:

- An end-to-end Scikit-Learn workflow
- Getting data ready
- choosing a ML model (to be used with ML models)
- fitting a model to the data (learning patterns)
- Making predictions with a model (using patterns)
- Evaluating model predictions
- Improving model predictions
- Saving & loading models

Q. An end-to-end Scikit-Learn workflow:

	age	sex	cp	chol	fbs	thal	target
0	57	1	3	233	1	3	1
1	57	0	0	234	2	3	0

from sklearn.ensemble import RandomForestClassifier → choose right model
clf = RandomForestClassifier() # clf - classifier
clf.get_params() → shows its parameters

from sklearn.model_selection import train_test_split → fitting model
x_train, x_test, y_train, y_test = train_test_split(x,y, test_size=0.2) ← splits x,y to training & test data with size 0.2 means 80% - 20%

clf.fit(x_train, y_train) ← says classification model, random forest finds patterns in training data

$y_{pred} = \text{clf.predict}(x_{\text{test}})$ → makes a prediction
 $\text{clf.score}(x_{\text{train}}, y_{\text{train}})$ → evaluating model → if model found patterns
 → systems mean accuracy on given test data
 $\text{clf.score}(x_{\text{test}}, y_{\text{test}})$ → shows accuracy on test dataset
 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
 $\text{print(classification_report(y_{\text{test}}, y_{\text{pred}}))}$ → shows classification metrics
 compare test labels to predictions
 $\text{accuracy_score}(y_{\text{test}}, y_{\text{pred}})$

run 13 times with
 $\text{np.random.seed(23)}$
 $\text{for i in range(10, 140, 10):}$ diff. n.est. parameter is giving diff. accuracy each time
 $\quad \text{clf} = \text{RandomForestClassifier(n_estimators=i)} \cdot \text{fit}(x_{\text{train}}, y_{\text{train}})$
 $\quad \text{print(f'Accuracy with {i} estimators: {clf.score(x_{\text{test}}, y_{\text{test}}) * 100 : .2f} %')}$
 import pickle; save model to load it writing binary
 $\text{pickle.dump(Clf, open('rfm_1.pkl', 'wb'))}$ (upto 2 decimal places in python shorthand)
 Saved file name
 $\text{loaded_model} = \text{pickle.load(open('rfm_1.pkl', 'rb'))}$ read binary
 $\text{loaded_model.score(x_{\text{test}}, y_{\text{test}})}$ → (scores) as per last n_estimator value i.e.
 $= 13$

1. Getting our data ready to be used with machine learning:

3 main things: → Split data into features (x) & labels (y)
 → Filling (aka imputing) or disregarding missing values
 → Converting non-numerical → numerical (aka feature encoding)

→ heart_disease.drop("target", axis=1) → removes "target" column

→ car_sales → Make Colour Odometer(km) Doors Price

0	Toyota	Red	150000	4	4000
1	Honda	Blue	35000	3	50000

$x = \text{car_sales.drop("Price", axis=1)}$
 $y = \text{car_sales["Price"]}$

Now we have to change them into numbers

from sklearn.preprocessing import OneHotEncoder
 from sklearn.compose import ColumnTransformer
 categorical_features = ["Make", "Colour", "Doors"]
 one_hot = OneHotEncoder() → handling it, a process used to turn categories → numbers
 transformer = ColumnTransformer([{"one-hot": one_hot, are hot, categorical features}], remainder = "passthrough") → accepts tuples with a name

transformed_x = transformer.fit_transform(x) → applies transformation in 'cat_features'
 takes one hot features to dataframes x
 remainder = passthrough (don't do anything)
 categorical features & rest all i.e.

pd.DataFrame(transformed_x) →

0	1	2	3	4
0	0	0	0	1
1	0	0	1	0
2	0	1	0	0

another way to do:

dummies = pd.get_dummies(car_sales[["Make", "Colour", "Doors"]])

dummies →

	Door	Make_Toyota	Make_Honda	Colour_Brown	Colour_Blue
0	4	1	0	0	0
1	3	0	1	0	0
2	3	0	0	1	0

from sklearn.ensemble import RandomForestRegressor → regression can predict non-i

model = RandomForestRegressor()

np.random.seed(26)

x_train, x_test, y_train, y_test = train_test_split(transformed_x, y, test_size=0.2)

model.fit(x_train, y_train)

model.score(x_test, y_test) → returns co-eff. of determination of the prediction

→ (B)

(★) (aka R²: statistical measure of how well model predicts)

If missing data is there → imputation

disregarding

car_sales_missing = pd.read_csv("car-sales-missing-data.csv")

car_sales_missing.isna().sum() → shows how many missing rows there in each col
 is non

car_sales_missing["Doors"].value_counts() → checking how many doors of each type

• fill missing data with Pandas:

car_sales_missing["Make"].fillna("missing", inplace=True) → fills missing values

"Colour". → "

"Odometer(km)". → (car_sales_missing["Odometer"].mean(), inplace=True)

"Pclass" → (4, inplace=True)

→ car_sales_missing.dropna(inplace=True) → removes empty rows (here only rows with "Price" empty will be removed)

car_sales_missing.isna().sum() → shows no. of rows empty for each column

①

• fill missing values with Scikit-Learn

from sklearn.impute import SimpleImputer as Num, as "mean"
 from sklearn.compose import ColumnTransformer

cat_imputer = SimpleImputer(strategy="constant", fill_value="missing")

door_imputer = SimpleImputer(" ", fill_value=4)

num_imputer = SimpleImputer(strategy="mean")

cat_features = ["Make", "Colour"] defining columns

door_features = ["Doors"]

num_features = ["Odometer"] takes list of multiple diff. transformer

imputer = ColumnTransformer([{"cat_imputer": cat_imputer, cat_features},

("door_imputer": door_imputer, door_features),

("num_imputer": num_imputer, num_features)])

(something that fills the missing data)

`.predict_proba()`: returns possibilities of a classification model
`clf.predict_proba(x_test[5])` → e.g. $\begin{bmatrix} \text{prob. of '0'} \\ \text{prob. of '1'} \end{bmatrix}$
 (predicting on the same data)
 $\rightarrow [0.18, 0.82]$
`clf.predict(x_test[5])` → $\begin{bmatrix} 0, 1, 1, 0, 1 \end{bmatrix}$

from sklearn.ensemble import RandomForestRegressor

(D) `RandomForestClassifier(2)`
`x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)`

model = RandomForestRegressor() creating model instance

model.fit(x_train, y_train)

`y_preds = model.predict(x_test)` → getting output similar to `y-test`)

from sklearn.metrics import mean_absolute_error → comparing prediction to the truth

`mean_absolute_error(y_test, y_preds)`
 $\rightarrow \text{Avg. } 0.333636$
 → on avg. each one of our models prediction is 0.33 diff. to the actual test value

4. Evaluating a machine learning model

3 ways to evaluate Scikit-learn models/estimators → estimator's built-in ✓ score() method
 problem specific- metric functions → the scoring parameter

Score() method:

(E) `//heart_disease`
`clf = RandomForestClassifier()`
 (F) `max value = 1.0 & min = 0.0`
`clf.score(x_train, y_train)`
`clf.score(x_test, y-test)`

→ returns co-eff. of determination
 → in simple words, returns mean accuracy
 → portion of variation in the dependent variable that is predictable from the independent variable

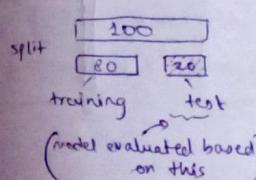
(G) `model = RandomForestRegressor()`

`model.score(x_test, y-test)` → returning R² value!

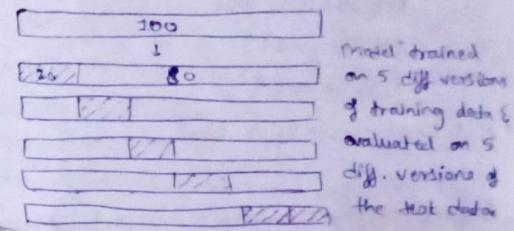
- from sklearn.model_selection import cross_val_score
 from sklearn.ensemble import RandomForestClassifier
 (H) `clf = RandomForestClassifier()`
 (I) `cross_val_score(clf, x, y, cv=5)`
 → by default 'scoring=None'
 If we want, we can import our own scoring methods!
 → shows array of 5 values (diff. mean acc.)
`clf_single_score = clf.score(x-test, y-test)`
`clf_cross_val_score = np.mean(cross_val_score(clf, x, y, cv=5))`
`clf_single_score, clf_cross_val_score` → comparing both
 $\rightarrow 0.80, 0.82$

cross-validation: (cv = we can write any no.)

Normal train & test split



5-fold cross validation



Classification models evaluation metrics:

(J)

`cross_val_score = cross_val_score(clf, x, y, cv=5)`

`print(f"Heart_d cross-val. acc: {np.mean(cross_val_score)*100:.2f} %")`
 ↑ up to 2 decimal

Area under ROC curve:

ROC curves are comparison of a model's true positive rate (tpr) versus model's false positive rate (fpr)

- (K) True (+) = model predicts 1 when truth is 1
- (L) False (+) = model " 1 " truth " 0 "
- (M) True (-) = model " 0 " truth " 0 "
- (N) False (-) = model " 0 " truth " 1 "

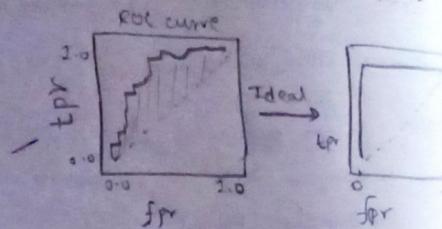
from sklearn.metrics import roc_curve

`clf.fit(x_train, y_train)`

make predictions with probabilities

`y_probs = clf.predict_proba(x_test)`
`y_probs[:5]`
 $\rightarrow \begin{bmatrix} \text{prob. of '0'} \\ \text{prob. of '1'} \end{bmatrix} \rightarrow \begin{bmatrix} 0.45, 0.55 \\ 0.95, 0.05 \end{bmatrix}$

$y_probs_positive = y_prob[:, 1] \rightarrow$ slicing of column 1
 fpr, tpr, thresholds = roc_curve(y_test, y_probs_positive)
 fpr → shows fpr
 # (Creating a func. to plot ROC curve)
 import matplotlib.pyplot as plt
 def plot_roc_curve(fpr, tpr):
 '''
 Plots a ROC curve given fpr & tpr of a model.
 '''
 plt.plot(fpr, tpr, color="orange", label="ROC")
 plt.plot([0, 1], [0, 1], color="darkblue", linestyle="--", label="Random")
 plt.xlabel("FPR")
 plt.ylabel("TPR")
 plt.title("ROC curve")
 plt.legend()
 plt.show()



from sklearn.metrics import roc_auc_score
 roc_auc_score(y_test, y_probs_positive) → takes test labels & prob.
 ↪ 0.92 → gives area under curve wrt ideal fig!
 roc_auc_score(y_test, y_test) ↪ perfect area score

ROC & AUC curves are evaluation metrics for binary classification models (a model which predicts one thing or another, such as heart disease or not)

ROC curves compares tpr vs fpr at diff. classification thresholds

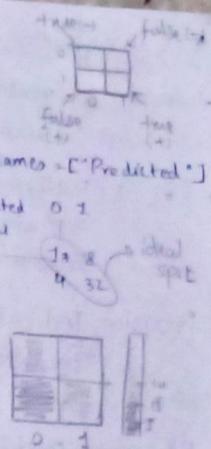
The AUC metric tells you how well your model is at choosing between classes.

Confusion Matrix: quick way to compare the labels a model predicts & the actual labels it was supposed to predict.

In essence, giving you an idea where the model is getting confused.

from sklearn.metrics import confusion_matrix
 y_pred = y.Predict(x_test)
 confusion_matrix(y_test, y_pred)
 pd.crosstab(y_test, y_pred, rownames=["Actual"], colnames=["Predicted"])

Making our matrix more visual with
 sns.heatmap() → plots rectangular data
 import seaborn as sns as a color-enabled matrix
 sns.set(font_scale=1.5) → setting font scale
 generate a confusion matrix
 conf_mat = confusion_matrix(y_test, y_pred)
 plotting it
 sns.heatmap(conf_mat),



Now creating matrix using scikit-learn → 2 methods!

from sklearn.metrics import ConfusionMatrixDisplay
 ConfusionMatrixDisplay.from_estimator(estimator=clf, X=x, y=y) → Plot confusion matrix given an estimator & data
 ConfusionMatrixDisplay.from_predictions(y_true=y_test, y_pred=y_pred) → Plot confusion matrix given true & predicted labels

Classification Report

It's like report bank of a no. of different parameters evaluating our classification models

from sklearn.metrics import classification_report
 print(classification_report(y_test, y_pred))

	Precision	Recall	f1-score	Support
0	0.82	0.66	0.74	25
1	--	--	--	36
accuracy	--	--	0.8	61
macro avg	--	--	--	--
weighted avg	--	--	--	--

Now precision & recall becomes important when there aren't much balanced classes

no. of test cases 0-0 & 1-1, here 25-36 so balanced
 not like 1 in 60! we need to identify which one patient is '1'!

- Classification Model Evaluation Metrics/Techniques:
- Accuracy: good measure to start if classes are balanced. Perfect accuracy is 1.0.
 - Recall: indicates the proportion of actual positives which were correctly classified. Model which produces no false negatives has a recall of 1.0.
 - Precision: indicates the proportion of true identifications (Model predicted class 1) which were actually correct. Model with no false positives has precision of 1.
 - F1 score: combination of precision & recall. (max value = 1.0)
 - Support: no. of samples each metric was calculated on
 - Macro avg: the avg. precision, recall & f1-score between classes. If you have class imbalances - pay attention to this metric.
 - Weighted avg: the weighted avg. precision, recall & f1-score between classes. Weighted means each metric is calculated w.r.t how many samples there are in each class.
 - Confusion Matrix: compares predicted values with true values in a tabular way. If 100% correct, all values in the matrix will be top left to bottom right (diagonal line).
 - Cross-Validation: splits your dataset into multiple parts & train & test your model on each part then evaluate performance as an avg.
 - Classification Report: returns some of the main classification metrics such as precision, recall & f1-score.
 - ROC curve: aka "receiver operating characteristic" is a plot of true positive rate vs false-positive rate.
 - Area Under Curve (AUC) score: The area underneath the ROC curve. A perfect model achieves an AUC score of 1.0.
- What classification metric should you use?
- Accuracy is a good measure to start if all classes are balanced.
 - Precision & recall become more imp. when classes are imbalanced.
 - If false-positive predictions are worse than false-negatives, aim for high precision.
 - If false-negative predictions are worse than false-positives, aim for higher recall.
 - F1 score is a combination of precision & recall.
 - A confusion matrix is always a good way to visualize how a classification model is going.
- Regression Model Evaluation Metrics:
- R^2 / Co-eff. of determination
 - Mean Absolute Error (MAE)
 - Mean Squared Error (MSE)
 - R^2 : it compares my models prediction to the mean of the target. Values can range from -ve oo (a very poor model) to 1. For ex. if all my model does is predict the mean of the targets, its R^2 value = 0. If my model predicts a range of nos., R^2 = 1
 - from sklearn.metrics import r2_score
 $r2_score(y_true = y_test, y_pred = y_test) \rightarrow 1.0$
 - Mean Absolute Error (MAE): its avg of the absolute differences between predictions and actual values. It gives you an idea of how wrong your model predictions are!
 - from sklearn.metrics import mean_absolute_error
 $y_preds = model.predict(x_test)$
 $mae = mean_absolute_error(y_test, y_preds)$
 $mae \rightarrow 0.33$
 - MAE turns -ive to +ive
 - df = pd.DataFrame(data={"actual values": y_test, "predicted values": y_pred})
 $df["differences"] = df["predicted values"] - df["actual values"]$
 $df.head(5)$

	actual values	predicted values	differences
df["differences"].mean()	0.02	0.700	-0.678
np.abs(df["differences"]).mean()	0.33	0.81	0.48
			getting diff from mae & MAE didn't change to 0.476

Mean Squared Error (MSE): MSE is the mean of the square of the errors b/w actual & predicted values

from sklearn.metrics import mean_squared_error

y_preds = Model.predict(x-test)

mse = mean_absolute_error(y-test, y-preds)

mse ≈ 0.37

df[["squared differences"]] = np.square(df[["differences"]])

↳ now here diff. is in decimals but actually represents real world values so you square large no. \rightarrow Large error!

★ Regression Model Evaluation Metrics / Techniques:

- R²/co-eff. of determination - MAE - MSE

★ Which regression metric should you use?

- R² is similar to accuracy. It gives you a quick indication of how well your model might be doing. Generally, the closer your R² value is to 1.0, the better the model.

But it doesn't really tell exactly how wrong your model is in terms of how far off each prediction is.

- MAE gives a better indication of how far off each of your model's prediction are on average.

- MSE amplifies larger differences.

• Scoring() parameter:

from sklearn.model_selection import cross_val_score

from sklearn.ensemble import RandomForestClassifier

np.random.seed(32)

(E)

clf = RandomForestClassifier()

cv_acc = cross_val_score(clf, x, y, cv=5, scoring="accuracy")

cv_acc, np.mean(cv_acc) * 100 : .2f

↳ returns array with 5 values

((Cross-Validated accuracy))

cv_acc = cross_val_score(clf, x, y, cv=5, scoring="accuracy")

cv_acc, np.mean(cv_acc) * 100 : .2f

↳ returns 5 values ↳ cross-validated accuracy ↳ ap of above & this one is same!

cv_precision = cross_val_score(clf, x, y, cv=5, scoring="precision")

cv_precision, np.mean(cv_precision) * 100 : .2f

↳ returns 5 values ↳ cross-validated precision

cv_recall = cross_val_score(clf, x, y, cv=5, scoring="recall")

cv_recall, np.mean(cv_recall) * 100 : .2f

↳ returns 5 values ↳ cross-validated recall

Now using it for regression problem

from -- import --

" "

model = RandomForestRegressor()

cv_r2 = cross_val_score(model, x, y, cv=5, scoring=None)

np.mean(cv_r2) \rightarrow 0.653 Mean squared error

cv_mse = cross_val_score(model, x, y, cv=5, scoring="neg_mean_squared_error")

np.mean(cv_mse) \rightarrow -0.430 Mean absolute error

cv_mae = cross_val_score(model, x, y, cv=5, scoring="neg_mean_absolute_error")

np.mean(cv_mae) \rightarrow -0.468

• Diff. evaluation metrics as Scikit-learn functions:

- from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

" " ensemble -> RandomForestClassifier

" " Model_selection -> train_test_split

(F)

clf.fit(x-train, y-train) making predictions

y_preds = clf.predict(x-test) evaluating model using evaluate function

Print("Classifier metrics on test set")

Print(f"Accuracy : {accuracy_score(y-test, y-preds)}") \rightarrow 0.885 ..

Print(f"Precision : {precision_score(y-test, y-preds)}") \rightarrow 0.842 ..

Print(f"Recall : {recall_score(y-test, y-preds)}") \rightarrow 0.969 ..

Print(f" F1 : {f1_score(y-test, y-preds)}") \rightarrow 0.902 ..

```

from sklearn.metrics import r2_score, mean_absolute_error, mean_squared_error
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import train_test_split
np.random.seed(42)

```

(1)

```

x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
model = RandomForestRegressor()
model.fit(x_train, y_train)
y_preds = model.predict(x_test)
print("Regressor metrics on the test set")
print(f"R2 score: {r2_score(y_test, y_preds)}") → 0.822
print(f"MAE: {mean_absolute_error(y_test, y_preds)}") → 0.321
print(f"MSE: {mean_squared_error(y_test, y_preds)}") → 0.236

```

5.0 Improving a model

Hyperparameters vs Parameters

model finds these patterns in data

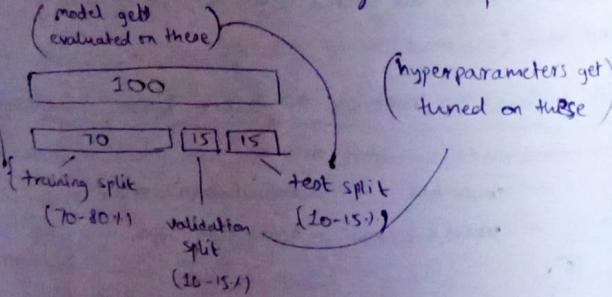
Settings on a model you can adjust to (potentially) improve its ability to find patterns

3 ways to adjust hyperparameters → by hand

randomly with RandomSearchCV

exhaustively with GridSearchCV

(model gets trained on these)



from — import RandomForestClassifier

clf = RandomForestClassifier()

clf.get_params()

get parameters for this estimator

We will try to adjust max_depth, max_features, min_samples_leaf, n_estimators, min_samples_split

Tuning hyperparameters by hand

def evaluate_preds(y_true, y_preds):

...

Performs evaluation comparison on y_true labels vs y_preds labels on a classification.

...

accuracy = accuracy_score(y_true, y_preds)

precision = precision_score(y_true, y_preds)

recall = recall_score(y_true, y_preds)

f1 = f1_score(y_true, y_preds)

metric_dict = { "accuracy": round(accuracy, 2),

"precision": round(precision, 2),

"recall": round(recall, 2),

"f1": round(f1, 2) }

print(f"Acc: {accuracy * 100:.2f}%")

print(f"Precision: {precision:.2f}")

print(f"Recall: {recall:.2f}")

print(f"F1 score: {f1:.2f}")

return metric_dict

from sklearn.ensemble import RandomForestClassifier

np.random.seed(42)

heart_disease = pd.read_csv("heart-disease.csv") shuffling data

heart_disease_shuffled = heart_disease.sample(frac=1) means 100% of data

x = heart_disease_shuffled.drop("target", axis=1)

y = heart_disease_shuffled["target"] split into 'x' & 'y'

train_split = round(0.7 * len(heart_disease_shuffled)) → 70% data split into train, validation, test sets

valid_split = round(train_split + 0.15 * len(heart_disease_shuffled)) → 15% of data

x_train, y_train = x[:train_split], y[:train_split] → 15% of data

x_valid, y_valid = x[train_split:valid_split], y[train_split:valid_split]

x_test, y_test = x[valid_split:], y[valid_split:]

clf = RandomForestClassifier() instantiate classifier with baseline hyperparameters

clf.fit(x_train, y_train) → make baseline predictions

y_preds = clf.predict(x_valid)

evaluate on validation sets

baseline_metrics = evaluate_preds(y_valid, y_preds)

→ op: Acc: 75.56%

Precision: 0.83

Recall: 0.74

F1 score: 0.78

Training it on diff dataset
from sklearn.ensemble import RandomForestClassifier

clf_2 = RandomForestClassifier(max_depth=10)

clf_2.fit(x_train, y_train)

y_preds_2 = clf_2.predict(x_valid)

clf_2.metrics = evaluate_preds(y_valid, y_preds_2)

② Acc: 82.2%

Precision: 0.89

Recall: 0.81

F1 score: 0.85

You notice that searching each parameter & writing values for it, this method not preferable.

• Hyperparameter tuning with RandomizedSearchCV:

from sklearn.model_selection import RandomizedSearchCV

grid = {"n_estimators": [10, 100, 200, 500, 1000, 1200],
"max_depth": [None, 5, 10, 20, 30],
"max_features": ["auto", "sqrt"],
"min_samples_split": [2, 4, 6],
"min_samples_leaf": [1, 2, 4]}

} in form of dictionary

np.random.seed(23)

③

x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)

clf = RandomForestClassifier(n_jobs=1) *dictate how much parallel process to run during training & prediction*

rs_clf = RandomizedSearchCV(estimator=clf,

param_distributions=grid,

n_iter=10, *-> no. of models to try*

cv=5, *(RSCV -> chooses any random value among above hyp-param & does this 10 times)*

verbose=2)

rs_clf.fit(x_train, y_train)

↳ fitting 5 folds for each of 10 candidates, totalling 50 fits

rs_clf.best_params_

↳ checking best esp among above 50

rs_y_preds = rs_clf.predict(x_test) *→ make predictions with best hyperparameters*
↳ last value called of rs_clf is chosen

evaluate the predictions

rs_metrics = evaluate_preds(y_test, rs_y_preds)

④ rs: Acc: 83.61%

Precision: 0.89

Recall: 0.81

F1 score: 0.85

• Hyperparameter tuning with GridSearchCV

In this, we will go through each combination of custom hyp-param given

grid_2 = {"n_estimators": [100, 200, 500],

"max_depth": [None], "max_features": ["auto", "sqrt"]},

"min_samples_split": [2], "min_samples_leaf": [1, 2]} *↳ total combo & with grid through = 3x2x2x1x2 = 24*

from sklearn.model_selection import GridSearchCV
np.random.seed(23)

⑤

⑥

gs_clf = GridSearchCV(estimator=clf,

param_grid=grid_2,

cv=5,

verbose=2) *↳ helps in running smoothly*

gs_clf.fit(x_train, y_train) *→ fitting 5 folds for each of 12 candidates, totalling 60 fits*

gs_clf.best_params_ *→ best hyp. param*

gs_y_preds = gs_clf.predict(x_test) *→ making predictions*

gs_metrics = evaluate_preds(y_test, gs_y_preds) *→ op: Acc: 83.61%*

Precision: 0.86

Recall: 0.86

F1 score: 0.86

• Comparing different models:

Note: Data was split differently for diff. models, so the comparisons aren't fully correct

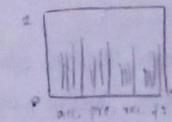
compare_metrics = pd.DataFrame({ "baseline": baseline_metrics,

"clf_2": clf_2.metrics,

"random search": rs.metrics,

"grid search": gs.metrics})

compare_metrics.plot.bar(figsize=(5, 4)) *→ plots graph*



acc, pre, rec, f1

6. Saving and loading a trained machine learning models

2 ways to save & load → Python's *pickle* module

or *joblib* module

• import pickle
// save existing file to model saved filename write binary
pickle.dump(gs_clf, open("gs_random_forest_model_1.pkl", "wb")) *→ pickle*

loaded_pickle_model = pickle.load(open("gs_random_forest_model_1.pkl", "rb")) *→ read binary*

make some predictions

pickle_y_preds = loaded_pickle_model.predict(x_test)

evaluate_preds(y_test, pickle_y_preds)

```

from joblib import dump, load
# save model to file
dump(gs_clf, "gs_random_forest_model_1.joblib")
# import a saved joblib file
loaded_joblib_model = load(filename="gs_random_forest_model_1.joblib")
# make & evaluate joblib predictions
joblib_y_preds = loaded_joblib_model.predict(x_test)
evaluate_preds(y_test, joblib_y_preds)

```

If model is large, it's better to use joblib to save & load.

7. Putting it all together

data = pd.read_csv("car-sales-extended-missing-data.csv")

	Make	Colour	Odometer (km)	Doors	Price
1	BMW	Red	120000	5	30000
2	BMW	Red	120000	5	30000

Things to remember

- All data should be numerical
- there should be no missing values
- Manipulate the test set the same as the training set
- Never test on data you've trained on
- Tune hyperparameters on validation set OR use cross-validation
- One best performance metric doesn't mean the best model.

• Getting data ready

```

import pandas as pd
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import OneHotEncoder

```

Modelling

```

from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import train_test_split, GridSearchCV

```

Setup random seed

```

import numpy as np
np.random.seed(42)

```

Import data & drop rows with missing labels

```

data = pd.read_csv("car-sales-extended-missing-data.csv")
data.dropna(subset=["Price"], inplace=True) # drop rows with missing
                                             "Price" value

```

Define different features & transformer pipeline

categorical_features = ["Make", "Colour"]

categorical_transformer = Pipeline(steps=[

("imputer", SimpleImputer(strategy="constant", fill_value="missing")),
 ("onehot", OneHotEncoder(handle_unknown="ignore"))
)
 value 'missing' & then applies one-hot encoding

door_feature = ["Doors"]

door_transformer = Pipeline(steps=[

("imputer", SimpleImputer(strategy="constant", fill_value=4))
)

numeric_features = ["Odometer (km)"]

numeric_transformer = Pipeline(steps=[

("imputer", SimpleImputer(strategy="mean"))
)

)

Setup preprocessing steps (fill missing values, then convert to nos.)

preprocessor = ColumnTransformer(transformers=[

("cat", categorical_transformer, categorical_features),
 ("door", door_transformer, door_feature),
 ("num", numeric_transformer, numeric_features)
)
 "Transformer for Door Feature"
 "(door)" - applies the door_transformer to column door_feature

)

Creating a preprocessing and modelling pipeline

```

model = Pipeline(steps=[("preprocessor", preprocessor),
                       ("model", RandomForestRegressor()))

```

Split data

```

x = data.drop("Price", axis=1)

```

```

y = data["Price"]

```

```

x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)

```

Fit & score our model

```

model.fit(x_train, y_train)

```

```

model.score(x_test, y_test)

```

We can also use 'GridSearchCV' or 'RandomizedSearchCV' with our Pipeline

```

from sklearn.model_selection import GridSearchCV

```

pipe_grid = Pipeline([

("preprocessor_num_imputer_strategy": ["mean", "median"]),
 "means in 'preprocessor', go to 'num', go to

"model_n_estimators": [100, 1000],
 "model_max_depth": [None, 5],
 "model_max_features": ["auto"],
 "model_min_samples_split": [2, 4]
)
 "model_max_depth": [None, 5], "model_max_features": ["auto"], "model_min_samples_split": [2, 4]
 "mean" & "median"

"model_max_depth": [None, 5], "model_max_features": ["auto"], "model_min_samples_split": [2, 4]
)
 "mean" & "median"

```

grid = GridSearchCV(model, pipe_grid, cv=5, verbose=2)
grid.fit(x_train, y_train)
grid.score(x_test, y_test)

```

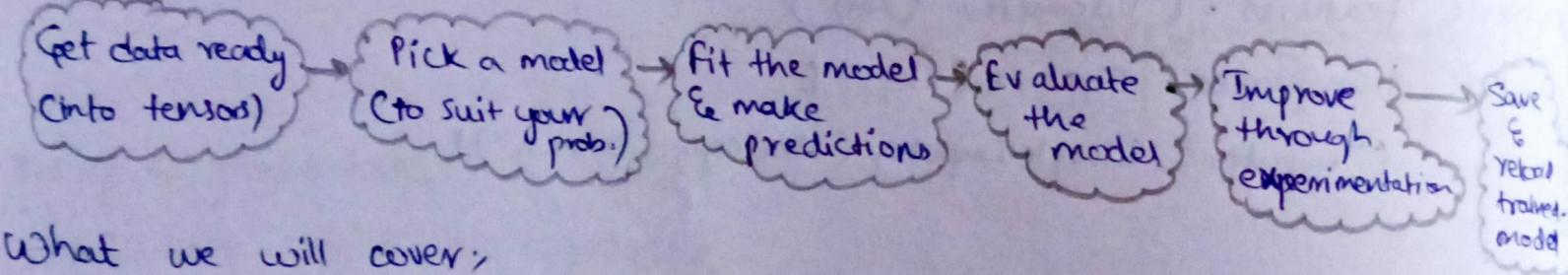
Why use 'Pipeline'?

- ↳ code organization
- reusability
- consistent application
- easy hyperparameter tuning.

Note: If Scikit-Learn version > 1.2 ; then "plot_roc_curve" will show error
so, use "RocCurveDisplay" syntax!

Neural Networks : Deep Learning, Transfer Learning & Tensorflow

A tensorflow workflow:



What we will cover;

- An end-to-end multi-class classification workflow with TensorFlow
- Preprocessing image data (getting it into Tensors)
- Choosing a deep learning model → fitting a model to the data.
(learning patterns)
- Making predictions with a model (using patterns)
- Evaluating Model predictions → Saving & loading models
- Using a trained model to make predictions on custom data