eleven | ENPC - Département IMI Refresher on Supervised Learning

To the attention of IMI Department

February 13th, 2024







Supervised Learning



Artificial intelligence can be organized around three families of machine learning models

| | *x* Approach | Model | Goal | |
|---------------------|---------------------------------|-------------------------|---|--|
| Machine learning | Supervised learning | Regression | Learn a function that maps an input to an output based on example input-output pairs (labeled data) | |
| | | Classification | | |
| | Unsupervised learning | Clustering | Identify and uncover previously undetected patterns | |
| | | Collaborative filtering | with no preexisting labeled data | |
| | Not necessary for the hackathon | | | |
| | Reinforcement learning | Q-table based | Find the optimal solution based on: An environment | |
| | | Deep learning based | A set of rules A playing agent | |



Regression models aim at modeling continuous variables whereas classification models aim at modeling categorical variables



Example of the difference between regression and classification



Regression

What is the temperature going to be tomorrow?

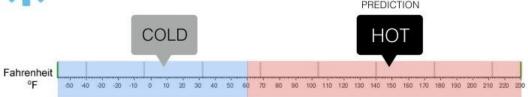


The output of the model is the temperature which is a continuous variable



Classification

Will it be Cold or Hot tomorrow?



The output of the model is a category: cold or hot





Understand

- Understanding the behavior or a given process, e.g.,
 - The relation between temperature and a building infrastructure
 - The relation between the intensity of a press and the quality of the animal food
 - ..
- If the goal is to understand a behavior, then the value is in the estimation of the parameters (the causal parameters)



Predict

- Predicting the output of a process given a new occurrence, e.g.,
 - What will the energy price be in one week?
 - What is this task's time to completion?
 - •



Optimize

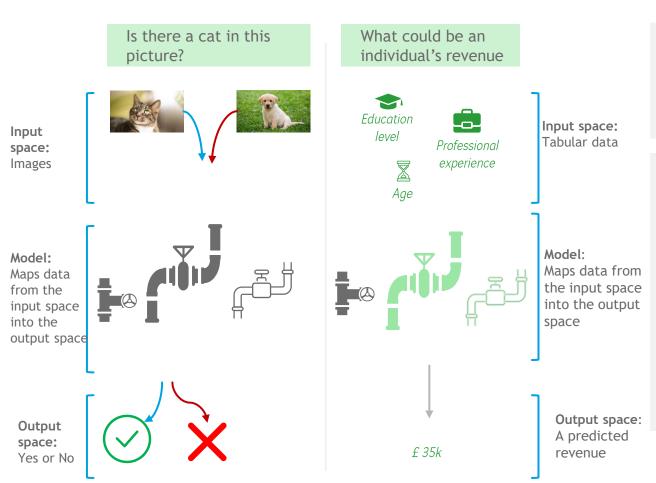
- Optimizing a process given a certain amount of information and constraints, e.g.,
 - Generating an efficient sprinkler network based on the building's blueprint and the involved regulatory principles
 - •

- If the goal is to predict an output or optimize a process, then the value is in the goodness of the output estimation
- It would be enough to have a « good » estimation of the parameter to be able to make accurate predictions or efficient optimizations



A machine learning model refers to both a group of functions to solve a task, and the chosen function within that group

lllustration of machine learning models



Machine Learning tasks consist in making a statistical model learn a mapping function from the inputs to the desired outputs from the data available

Formally, a model (e.g., Linear Regression) is a class of functions to consider for that task (e.g., all linear functions)

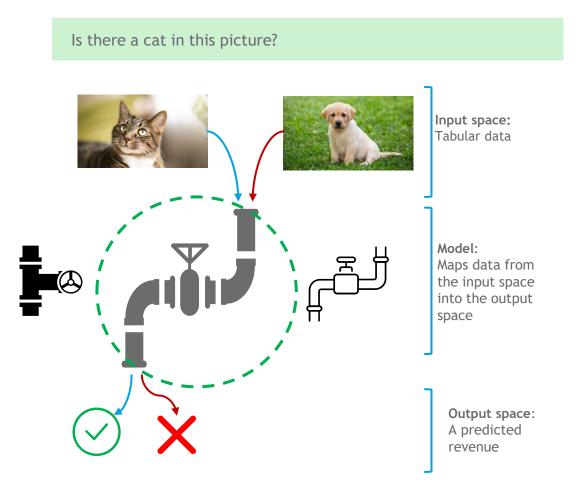
$$S = \{f_{\theta}, \theta \in \Theta\}$$

In other words, a group of mappings to choose from:



A machine learning model refers to both a group of functions to solve a task, and the chosen function within that group

Illustration of machine learning models



Training a model amounts to finding the optimal function with that class:

Find best $s \in S$

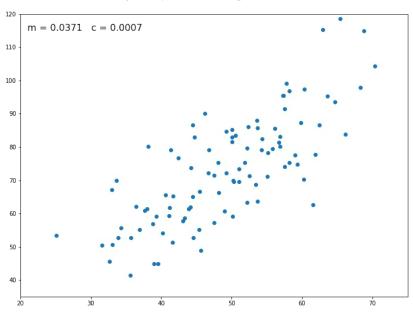
By extension, we call « model » both the ensemble of considered mappings (before training) and the one selected after the training process

- The goal of ML model training is to find the configuration of parameters' values of the model that most suit the observed data
- Once the model is "trained" it can then make proper inferences about a previously unknown observation (test set, real-life data)

Training or fitting a model amounts to finding the parameters making the model closer to the dataset with regards to a metric

Illustration of machine learning models

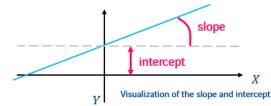
Example of linear regression



i Linear regression example

Our objective is to find the line that "best represents" (fits) the available observations (the point cloud).

That line is defined by two parameters: the slope θ_1 and the intercept θ_2



- We must find the intercept and slope (the parameters) that best fit the points
- We need a metric (the Loss) in order to tell apart the different parameters. We choose the distance between each point and the line.
- We find the best parameters by minimizing this distance.

Supervised learning algorithms aim to understand the fundamental relationship between explanatory variables X and a variable to be predicted Y

TWO MAIN PARADIGMS

Estimation of the relationship Y = f(X)

• The problem is to estimate the general function based on the available data and assuming an a priori form of the function

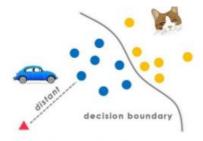
Approximation of the function dependent on the dataset

$$Y \approx \hat{f}_{\tau}(X_1, \dots, X_n)$$

Target variable or events to be predicted

Explanatory variables or Features

- The key question: is the approximation good enough to generalise to new observations, beyond the available data?
- Illustration:



Estimation of the probability distribution P(X,Y)

• The problem is to estimate the joint distribution based on the available data (e.g. Bayesian learning)

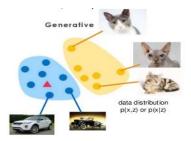
Approximation of the distribution dependent on the data set

$$\hat{P}_{\tau}(Y, X_1, \dots, X_n)$$

Target variable or events to be predicted

Explanatory variables or Features

- The key question: is the approximation good enough to generalise to new observations, beyond the available data?
- Illustration:



Regression attempts to estimate the mapping function from the input variables to numerical or continuous output variables

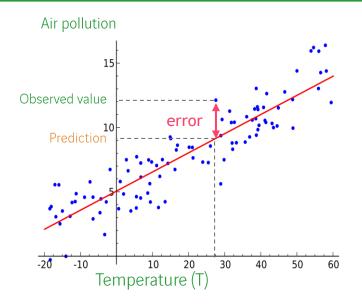
Description: Linear Regression

- Find the relation between input variables and continuous output variable
- Example: Understanding the relation between air pollution and temperature
- Class of model selected: linear regression

$$\hat{y} = f_{\theta}(x) = \theta_0 + \theta_1 x_1$$

Parameter $\theta = (\theta_0, \theta_1)$

Illustration





Metric

The error is the difference between model's prediction and the real value

$$error = (y - \hat{y})$$

A classical metric for regression is the Root Mean Squared Error (RMSE):

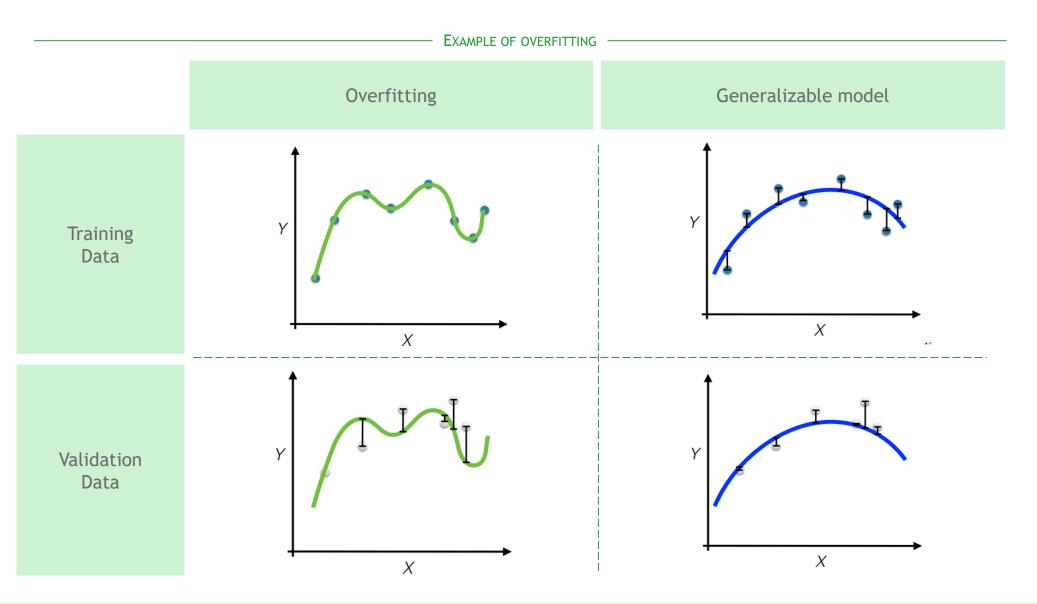
$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^{N} (Y_i - \hat{Y}_i)^2}$$
* ** Yi the observed output values ** \hat{Y}_i represents the **

- N is the number of observations
- \hat{Y}_i represents the predicted values

The goal of regression algorithms is to minimize the RMSE

How can we make sure that the parameters found are also optimal for new samples?

Various sources of error are then possible: measurement errors, errors in estimating the function f, errors related to the choice of the regularity of the function



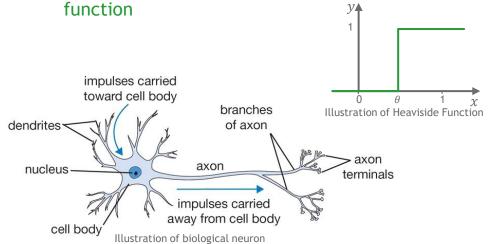
Artificial neural networks are based on biological neurons because the human brain processes information efficiently



Biological neurons

- Basis constituent of the nervous system
- Information is propagated by electric signals through dendrites and axons
- A neuron is activated when the input signal is higher than a threshold

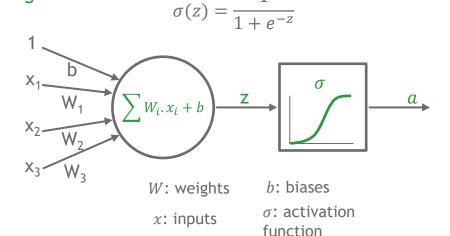
• A neuron is equivalent to a Heaviside step





Artificial neuron

- A simplified model of a biological neuron
- An artificial neuron is a linear function, on which an activation function is added
- The step function is not adapted due to its sensitivity to noise
- The common activation function used is the sigmoid function:

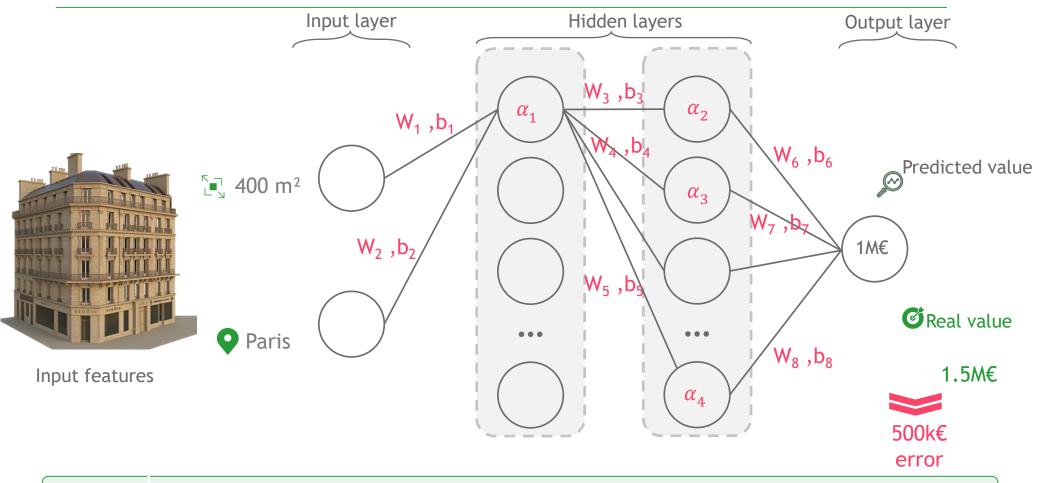




- Each neuron processes a bit of information and passes it to its children
- Overall, the network processes raw information into general concepts

A neural network learn how to imitate a function from a training dataset

© Computation process: real estate price prediction example

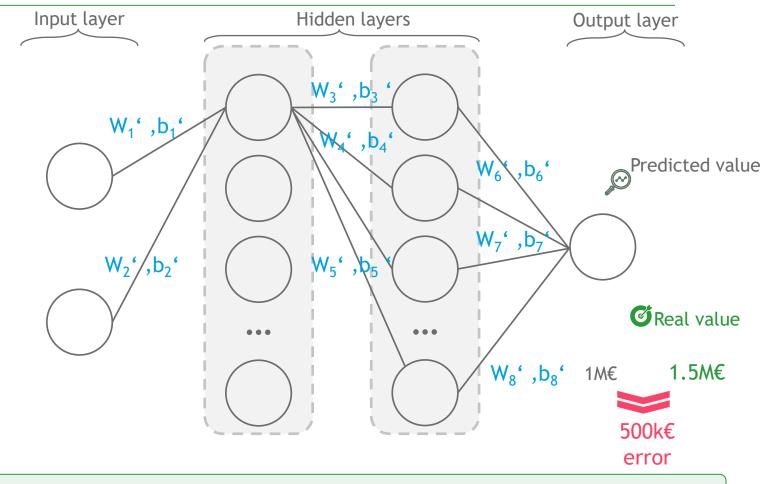




- This process is called the forward propagation and allows to compute the loss function of the model
- How can we find the model's optimal parameters for this task?

A neural network learn how to imitate a function from a training dataset

Computation process: real estate price prediction example





- This process is called the backward propagation and allows to update the neural network parameters
- Forward and backward propagation are done iteratively until having a high-performance model, by doing a gradient descent to minimize loss

The gradient descent can be controlled by a hyperparameter of the network: the learning rate



Overview of the learning rate

- The learning rate (α) controls the speed at which the model learn
- It controls the update of the weights
- If the learning rate is high, the weights will be strongly modified at each epoch
- A technique to speed up the gradient descent is the learning rate decay, which consists to slowly reduce the learning rate over time: allows large weight changes at the start of the learning process, and more fine-tuning towards the end



Illustration of the learning rate

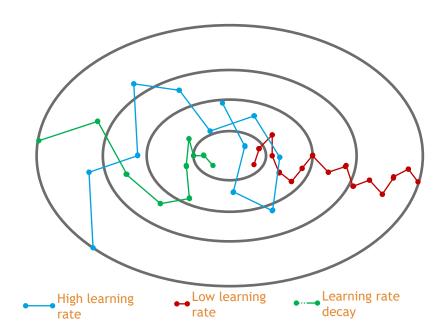


Illustration of the impact of the learning rate on the gradient descent



- The learning rate is equivalent to the length of the step, and it may be the most important hyperparameter to configure our model, and is generally in [10^{-6} , 1]
- Adaptive Learning Rates is a learning rate schedule adapting the LR to the performance of the model. Typically, once the performance of the model plateaus, LR can be decreased by a factor of 2



On Python, using an open-source machine learning framework like PyTorch can accelerate the path from research prototyping to production deployment



Deep Learning computation



- PyTorch is defined as an open-source machine learning library for Python
 It is designed for Deep Learning and used for applications such as NLP, Computer Vison



Main Features

- Easy Interface:
 - PyTorch offers easy to use API; hence it is considered to be very simple to operate and runs on Python
- Python usage:
 - Pytorch can leverage all the services and functionalities offered by the Python environment
- Computational graphs:
 - PyTorch provides a platform which offers dynamic computational graphs. Thus, a user can change them during runtime



Key advantages

- Easy to debug
- It includes many layers as Torch
- It includes lot of loss functions
- It can be considered as NumPy extension to **GPUs**
- It allows building networks whose structure is dependent on computation itself





Other useful libraries



Pandas: handle dataframes and structured data easily



Scikit-Learn: train and evaluate shallow MI models



Jupyter: iterate quickly on your code, keeping your variables across time



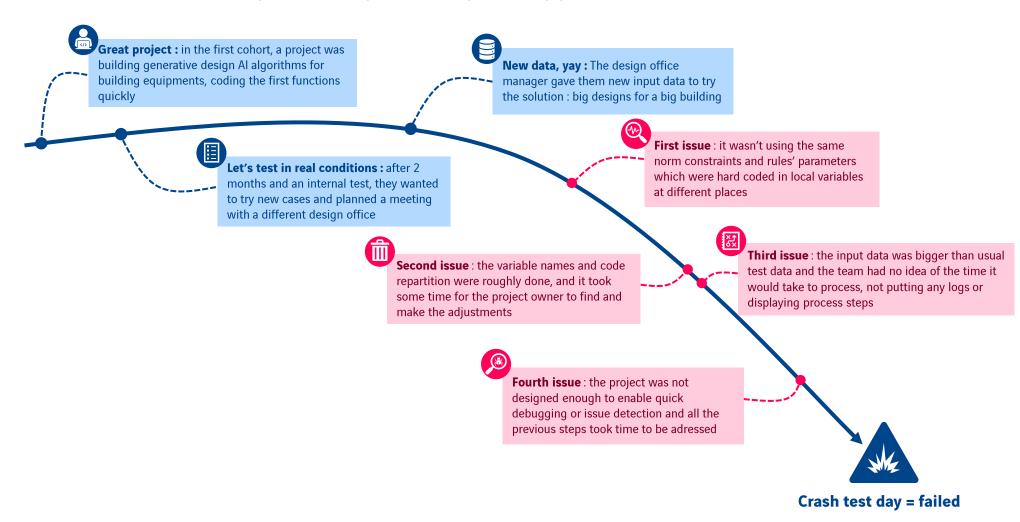
GitHub: share your code across the team

Some coding best practices



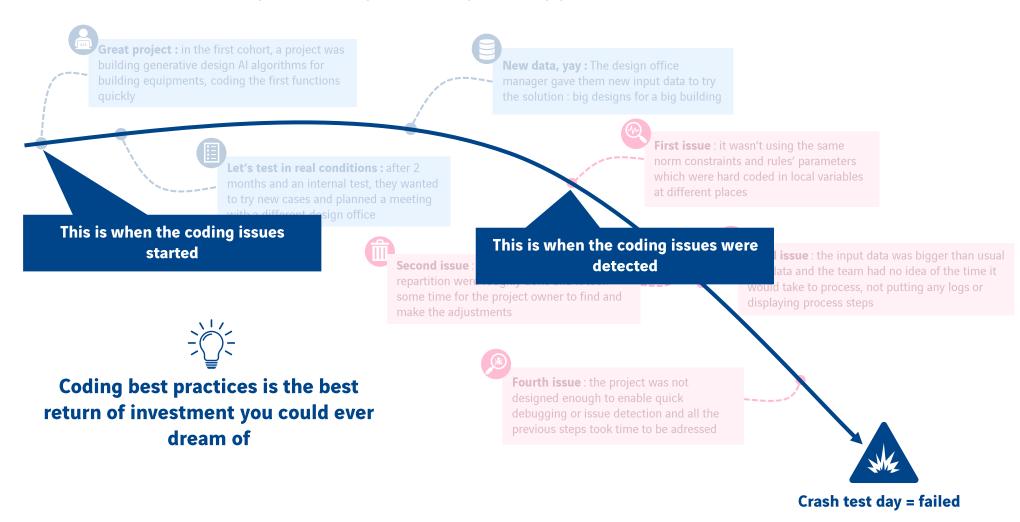
>

WHY IS IT IMPORTANT TO WRITE CLEAN CODE?



>

WHY IS IT IMPORTANT TO WRITE CLEAN CODE?







PEP8 IS A SET OF CONVENTIONS PROVIDING NAMING GUIDELINES FOR PYTHON

Type Convention Examples

| Function | snake_case: use lowercase words, separate words by underscores Start with a verb | function, write_name, delete_date, |
|----------|--|--|
| Variable | snake_case: use lowercase words, separate words by underscores Use a noun | • first_name, last_name, month, |
| Class | CamelCase: start each word with a capital letter, without separators Use a noun | MyClass, Player, |
| Method | snake_case: Use lowercase words, separate words by underscores Use verb | class_method, set_player, get_age, |
| Constant | Use uppercase words, separate words by underscores Use a noun | CONSTANT, INITIAL_DATE, |
| Module | snake_case: use lowercase words, separate words by underscores Be short, use a noun | module.py, feaure_engineering.py, |
| Package | Use short lowercase words Do not separate words by underscores | • mypackage, |





OTHER GUIDELINES TO NAME VARIABLES

- Use descriptive names to make it clear what the object represents and use pronounceable names
- Use English language
- + Avoid starting a name with Python specifics words like list, from, dict, ...
- Never use I, O, or I single letter names as these can be mistaken for 1 and 0
- + And much more to be found here: Pep8



PEP8 ALSO PROVIDES A SET OF RULES FOR CODE LAYOUT

Convention **Type Examples**

Line length

Limit all lines to a maximum of 79 characters

Indentation

- To keep lines under 79 characters, it can be necessary to break the codes with line continuation
- Use 4 spaces per indentation level
- Align the indented block either with the opening delimiter (1) or use a hanging incident (2) where every line is indented except the first one
- Use '\' to separate an expression into two lines

Closing brace

 Line-up closing brace with the first nonwhitespace character of the previous line or with the first character of the first line

White space

- Add a single whitespace before and after assignment operators (=, +=,...), comparisons (==, !=,...) and Booleans (and, or, ...)
- Use a whitespace after a ',' or ';'
- Do not use a whitespace after '(' or '[' or '}'

```
def function(arg one, arg two,
             arg_three, arg_four):
    return arg_one
```

```
def function(
                                  var = function(
        arg_one, arg_two,
                                      arg one, arg two,
        arg_three, arg_four):
                                     arg_three, arg_four)
    return arg one
```

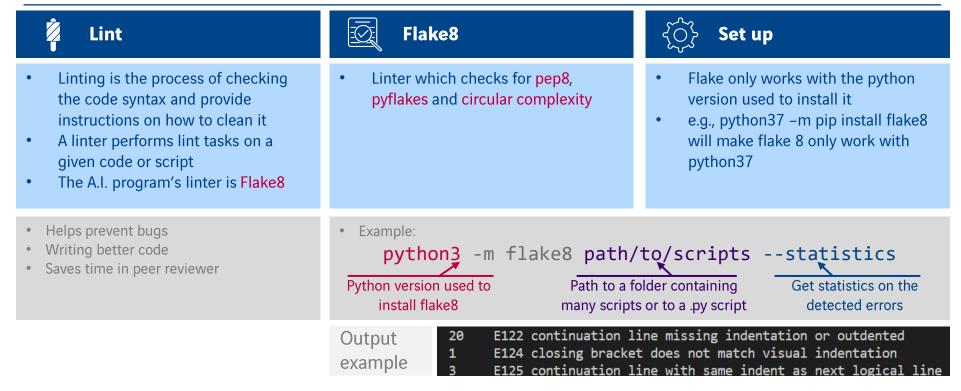
```
my list = [
   1, 2, 3,
                                1, 2, 3,
   4, 5, 6,
                                  4, 5, 6,
```

```
print(x, y)
def double(x):
x = 5
                                      list[3]
V = 6
```



FOR THE AI PROGRAM WE WILL USE A TOOL TO CHECK THE CODE SYNTAX AND LAYOUT CALLED FLAKES

** Quick introduction to linting and flake8





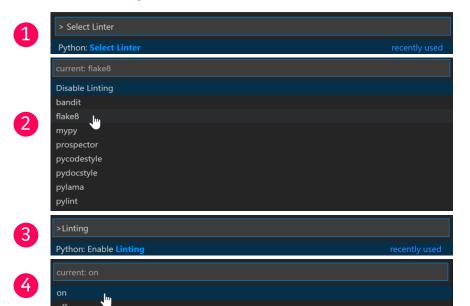


TUTORIAL: ONCE INSTALLED, FLAKE8 CAN BE DIRECTLY USED ON VSC

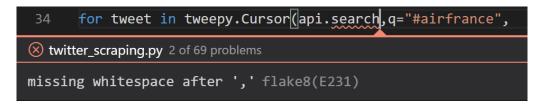


Setup on VSC

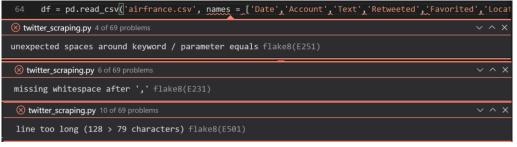
- Make sure Flake8 is installed
- Open the command palette with CTRL + SHIFT + P
- Open the category "Python: Select Linter" and choose Flake8
- Enable linting







★ 1 line is responsible for 10 linting problems !!



✓ The code is now PEP8 compliant





CODE DOCUMENTATION IS ESSENTIAL TO UNDERSTAND CODE YOU WROTE IN THE PAST OR TO SHARE YOUR CODE WITH OTHERS



Code documentation

- Start a comment with "# "
- Explain WHY decisions are taken and not HOW
- Include information in the variable as much as possible
- DRY: Don't Repeat Yourself
- Use codetags (see PEP350) for tasks like TODO, TODOC, ...
- # A list of cities to keep list = ['Paris', 'London']
- </> cities_to_keep = ['Paris', 'London']



Potential risks and inefficiencies

- Over-commenting code
- Lack of comments
- Update the code but not the comments
- Misleading/controversy comment
- Obvious comments
- Comments just because you feel obliged
- Funny comments



- Clean code is a good substitute to many comments
- Comments should be mostly used to explain main ideas or specific details



Python variables

Variables are the corner stones of any python code

Primitive types

m Container types

| str | "hello | |
|------|---------------------|----------------|
| set | world" {1, 5, 9} | set(|
| list | [1, 5, 9] | list(|
| dict | {1:"one",2:"t | wo"}dict() |

→ "concat"

🖶 Print

Concatenation

"con" +

| "cat" 1 + 2 → 3 |
|--|
| 1 + 2.5 → 3. |
| 1 + "cat" × |
| $[2, 3] + [4, \rightarrow [2, 3, 4,$ |
| 5] 5] 5] {2, 3} + {4, × |
| 5} {2, → {2, 3, 3}.add(4) 4} {1:"one"} + × {2:"two"} |

E Container indexing - for lists, strings...

```
negative index-5 -4 -3 -2 -1
            0 1 2
positive
index lst = [10, 20, 30, 40]
positiv561ce
negative slice5
Individual access to items via
```

 $\begin{array}{c}
|st[index]| \\
|st[0]| \longrightarrow 10 \text{ (first one)}
\end{array}$ $[st[-1] \rightarrow 50 \text{ (last one)}$ lst[1] → 20

Modify with assignment Ist[4]

```
Access to sub-sequence via <a href="Istart slice">Ist</a>[start slice : end slice
: step]
[st[]: -1] → [10, 20, 30, 40]
   lst[1:-1] \rightarrow [20, 30, 40]
                                           Index from
   lst[::- → [50, 40, 30, 20,
                                           (hQe from 0 to 4)
   lst[::- → [50, 30,
                                           Items count
   ist[ 1: → [20, 30]
                                           len(lst)→ 5
   [st[:] → [10, 20, 30, 40,
```

Modify with assignment [st[1:3] = [15, 25]

Conversion

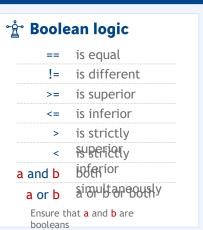
Additional remarks

Python is a dynamically typed language, which means the type of variable don't have to be declared

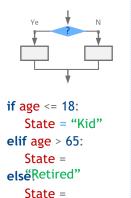


Loops and conditions

Loops and conditions lead to more complex program



Can go with several *elif*, and only one final *else*. Only the bloc of the first true condition is executed

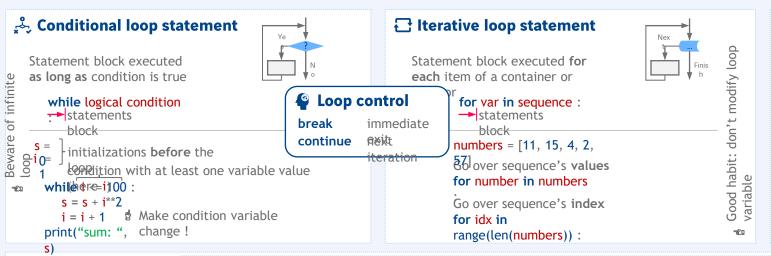


"Active"

Integer sequence

```
range([start,] end [, step])

start default 0, end not included in sequence, step
default 1
range(5→0 1 2 3 4 range(5, 12, → 5 8 11
)
range(3, → 3 4 5 6 range(20, 5, - → 20 15 10
8) 7 5)
range(len(seq → sequence of index of values in
) ) seq
```



Title Iterative loop tricks

Go simultaneously over sequence's **index** and **values**:

for idx, val in
enumerate(lst) :

Go simultaneously over two sequences' values:

for val1, val2 in zip(lst1, lst2):

. × @ & . *

Additional remarks

"For" loops are used for iterating through a sequence like a list, a tuple, a string or a dictionary

Appendix

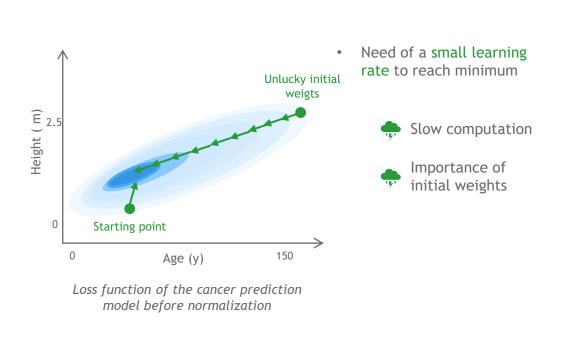


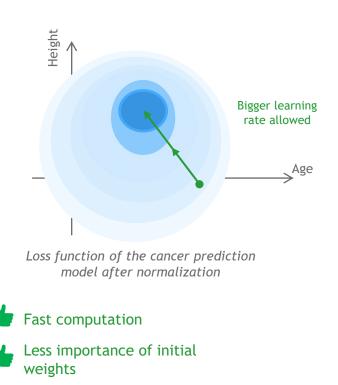
Normalisation throughout the neural network helps ease the optimization process



Why should layers be normalized?

Use case of cancer prediction based on Height and Age







- Use BatchNorm1d to normalize a layer with PyTorch
- Weight initialization can be "random" by default. Otherwise, it can depend on the type of activation function used: for "Sigmoid", a xavier initialization method is usually used



How to prevent Neural Network from overfitting?

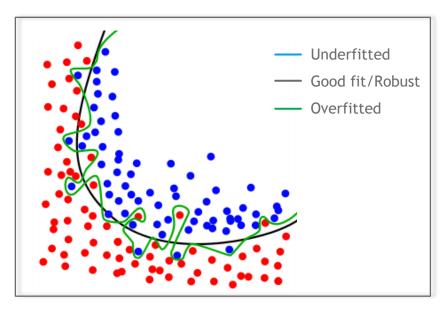
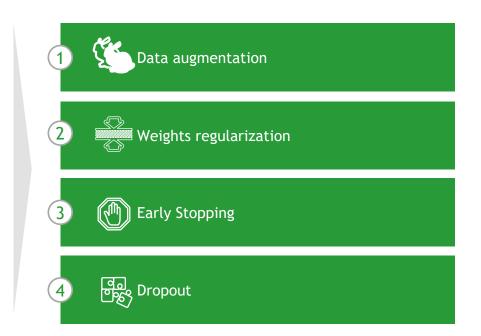
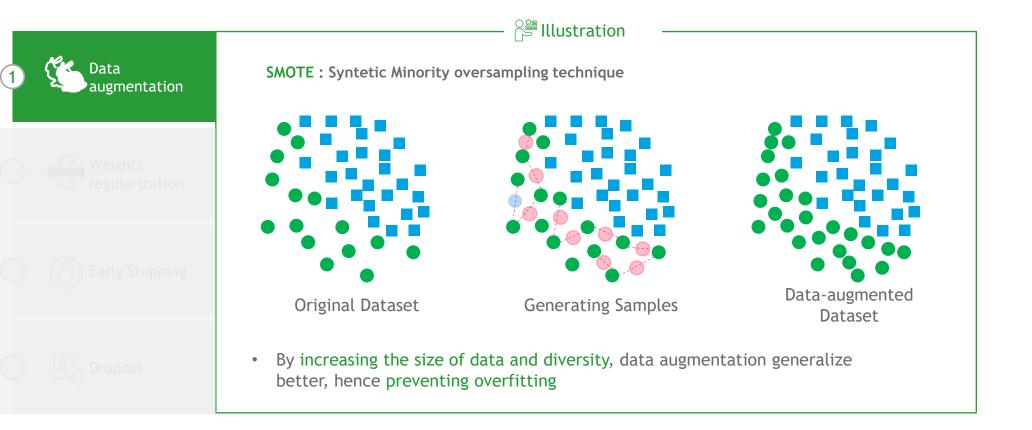


Illustration of a regression problem and overfitting possibilities





How to prevent Neural Network from overfitting?







How to prevent Neural Network from overfitting?











Lasso regularization

New objective function:

$$\frac{1}{N} \sum_{i=1}^{N} (Y_i - \hat{Y}_i)^2 + \lambda_1 \sum_{j=0}^{n} |\beta_j|$$

Ridge regularization

New objective function:

$$\frac{1}{N} \sum_{i=1}^{N} (Y_i - \hat{Y}_i)^2 + \lambda_2 \sum_{j=0}^{n} \beta_j^2$$

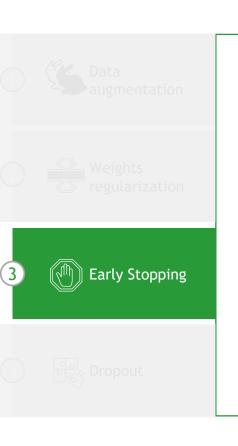
ElasticNet data regularization is the combination of both Lasso and Ridge regularization



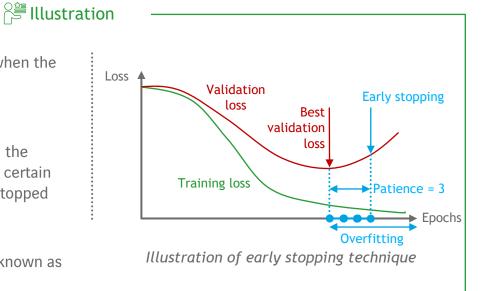
torch.optim.SGD(net.parameters(), lr=0.2, weight_decay=1e-2)



How to prevent Neural Network from overfitting?

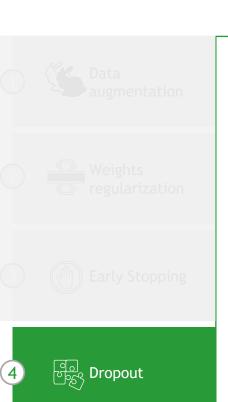


- Early stopping stops the training when the validation loss is increasing
- If the validation loss is higher than the current lowest validation lost for a certain number of epochs the training is stopped
- The number of epochs we wait is known as the patience





How to prevent Neural Network from overfitting?



Illustration

- Not training certain weights at a given update step, leading to greater generalization over the whole network
- The classic formulation of this is dropout where neuron activations are "dropped out" (set to zero) with probability p

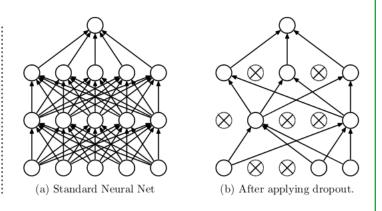


Illustration of dropout technique



torch.nn.Dropout(p=0.5)