

# QU'EST CE QU'UN DATASET

**X** : Features / Colonnes / Inputs

**Y** : Target / Variables / Labels

<b>X</b>					<b>Y</b>
22	rouge	France	32.6	1/10/20	Yes
30	vert	Allemagne	7.2	9/09/20	No
42	bleu	Ø	44.4	10/09/19	Yes

\* Les lignes représentent des "exemples"

\* Chaque ligne de **X** est associée à une ligne de **Y**

Les Data proviennent  
de différentes sources



Scraper  
www.

.CSV

<XML/>

Se DÉBROUILLER  
pour tout RASSEMBLER  
DANS UN FICHIER .CSV



MAINTENANT IL FAUT  
LES LABELLISER  
(si ce n'est pas déjà fait)

ET \*TADAA\*...

UN DATASET

# SUPERVISED LEARNING

## The Mission

1

À partir d'un  
DATASET  
LABELLISÉ...



ISOLER LES  
Features des  
TARGETS...



Pour obtenir  
X et Y

2

À partir de  
X et Y



CRÉER UN MODÈLE  
f



Tel que  
•  $f(X)$  = prédictions  
• prédictions  $\approx Y$

The Mission

L'objectif est d'obtenir un **modèle F**, qui lorsqu'on lui présente des **données inédites** provenant de l'environnement de **production**, fassent des **prédictions** qui se révèlent **correctes** suffisamment souvent pour être utilisable en production

ON APELLE ÇA  
LA GÉNÉRALISATION  
↗ à connaître

# Les grandes étapes du PREPROCESSING

Preprocessing: Les étapes préalables au traitement (=processing)  
par le modèle de Machine Learning

DONNÉES  
ENTRANTES

PREPROCESSING

MACHINE LEARNING

SORTIE DE LA  
PIPELINE

- Data Cleaning
- Data Transform
- Feature Engineering
- Dimensionality Reduction

Toutes ces étapes  
ne sont pas  
"obligatoires" ...

Il est même rare  
de toutes les utiliser  
en même temps...

EN FAIT...  
Ce qu'il faut RETENIR  
C'EST QUE...

LE PREPROCESSING A 2 "GRANDS" OBJECTIFS:

- \* Rendre les données entrantes compatibles avec nos modèles de Machine Learning
- \* Transformer les données pour améliorer les performances prédictives de nos Modèles  
(et seulement cette étape là est obligatoire)

# The Scikit-Learn A.P.I.

## TRANSFORMER VS. ESTIMATOR

Utilisés pour  
le PREPROCESSING

TRANSFORMER

.fit()

.transform()

CALIBRATION

APPLICATION

Lorsqu'on .fit() un TRANSFORMER, on enregistre des statistiques du dataset... ce qui nous permet de mémoriser une transformation pour pouvoir la réutiliser plus tard

Permet d'appliquer une transformation enregistrée préalablement à un nouvel ensemble de données

Souvent se sont des  
ALGORITHMES de  
MACHINE LEARNING

ESTIMATOR

.fit()

.predict()

Permet d'appliquer la fonction mathématique entraînée préalablement (le modèle) à de nouvelles données.  
On appelle ça faire des PRÉDICTIONS

Lorsqu'on .fit() un ESTIMATOR on applique l'algorithme associé afin de trouver les paramètres qui s'ajustent le mieux au dataset présenté  
On appelle ça l'entraînement

BREAK INTO.AI

COIN VOCABULAIRE

LA DIFFÉRENCE  
ENTRE ALGORITHME  
ET MODÈLES  
DE MACHINE LEARNING

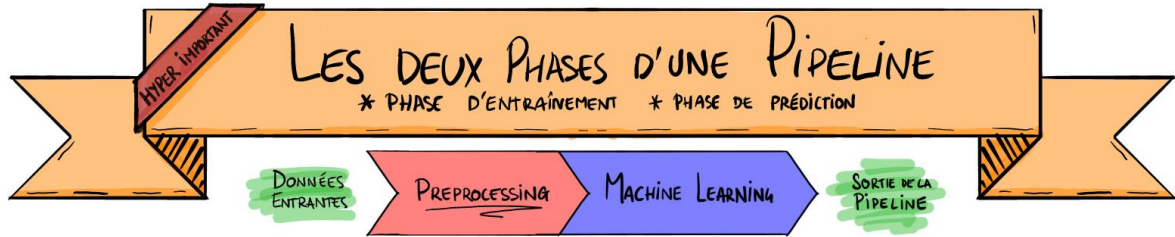
UN ALGORITHME DE M.L.  
EST UNE RECETTE DE CUISINE (procédure)  
POUR CRÉER UN MODÈLE DE ML.

UN MODÈLE DE M.L. EST  
UNE FONCTION MATHÉMATIQUE  
POUR FAIRE DES PRÉDICTIONS

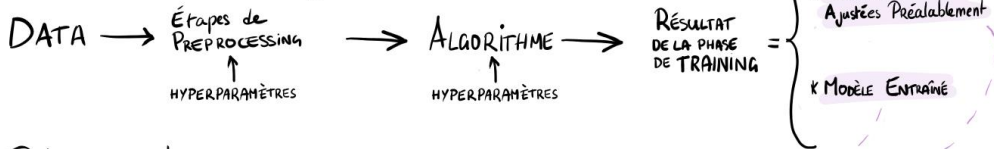
Lorsqu'on choisit un ALGO  
on choisit en plus des  
"réglages"... Ces "réglages"  
choisis avant l'entraînement  
sont appelés des **HYPER-PARAMÈTRES**

Lorsqu'un ALGO entraîne un  
MODÈLE, il cherche l'ensemble  
de **PARAMÈTRES** du modèle qui est  
le mieux ajusté au dataset

Les **HYPER-PARAMÈTRES** et les **PARAMÈTRES**  
ONT UNE INFLUENCE SUR LES  
PRÉDICTIONS DU MODÈLE



## Phase de Training



NOTRE MISSION ÉTAIT DE CRÉER UN MODÈLE CAPABLE DE PRÉDIRE AUTOMATIQUÉMENT. CE QUE NOUS AVONS ACCOMPLI

## Phase de Prédiction



"TOUT ÇA POUR ÇA?"

## 1 TRAINING SUR TOUT LE DATASET

PROBLÈME: Aucun moyen de vérifier les résultats

## 2 SPLIT TRAIN / TEST

STRATIFIÉ

PROBLÈME: GROS FACTEUR CHANCE...  
(changer la seed peut faire changer la mesure de performances)

## 3 K-FOLD CROSS VALIDATION



PERMET D'OBTENIR UN SCORE FIABLE  
POUR UNE PIPELINE



Après avoir validé une pipeline, il faut la  
ré-entraîner sur la totalité du DATASET

PROBLÈME: PEUT SE RÉVÉLER TRÈS CÔUTEUX EN  
TEMPS / PUISSANCE DE CALCUL

# BREAK INTO AI

## NOTRE MISSION RECHERCHE DE LA MEILLEURE PIPELINE

- 1) Créez un schéma de crossvalidation
- 2) TESTEZ UNE GRANDE VARIÉTÉ DE PIPELINE  
( $\hookrightarrow \neq$  ALGOS,  $\neq$  PREPROCESSING,  $\neq$  HYPER-PARAMÈTRES)
- 3) SÉLECTIONNER LA PIPELINE avec le meilleur score de C.V.
- 4) RÉ-ENTRAÎNEZ LA PIPELINE SUR LA TOTALITÉ DES DONNÉES  
DISPONIBLES POUR L'ENTRAÎNEMENT



Le petit schéma  
de la

# K-FOLD CROSS VALIDATION

ici  $K=5$

ON DÉCOUPE  
le Train-Set  
en  
 $K$  morceaux  
STRATIFIÉS

CE QUI NOUS  
DONNE  
 $K$  Splits  
Différents

A	B	C	D	E
A	B	C	D	E
A	B	C	D	E
A	B	C	D	E
A	B	C	D	E

ON ENTRAÎNE LA PIPELINE  
 $K$ -Fois (1 fois par Split)

TRAIN SETS

$A + B + C + D$   
 $A + B + C + E$   
 $A + B + D + E$   
 $A + C + D + E$   
 $B + C + D + E$

VALID SETS

E  
D  
C  
B  
A

Score de  
Cross-Validation

ce qui nous donne

ON FAIT  
LA

MOYENNE

DES  
SCORES  
SUR CHAQUE

VALIDATION SET



