

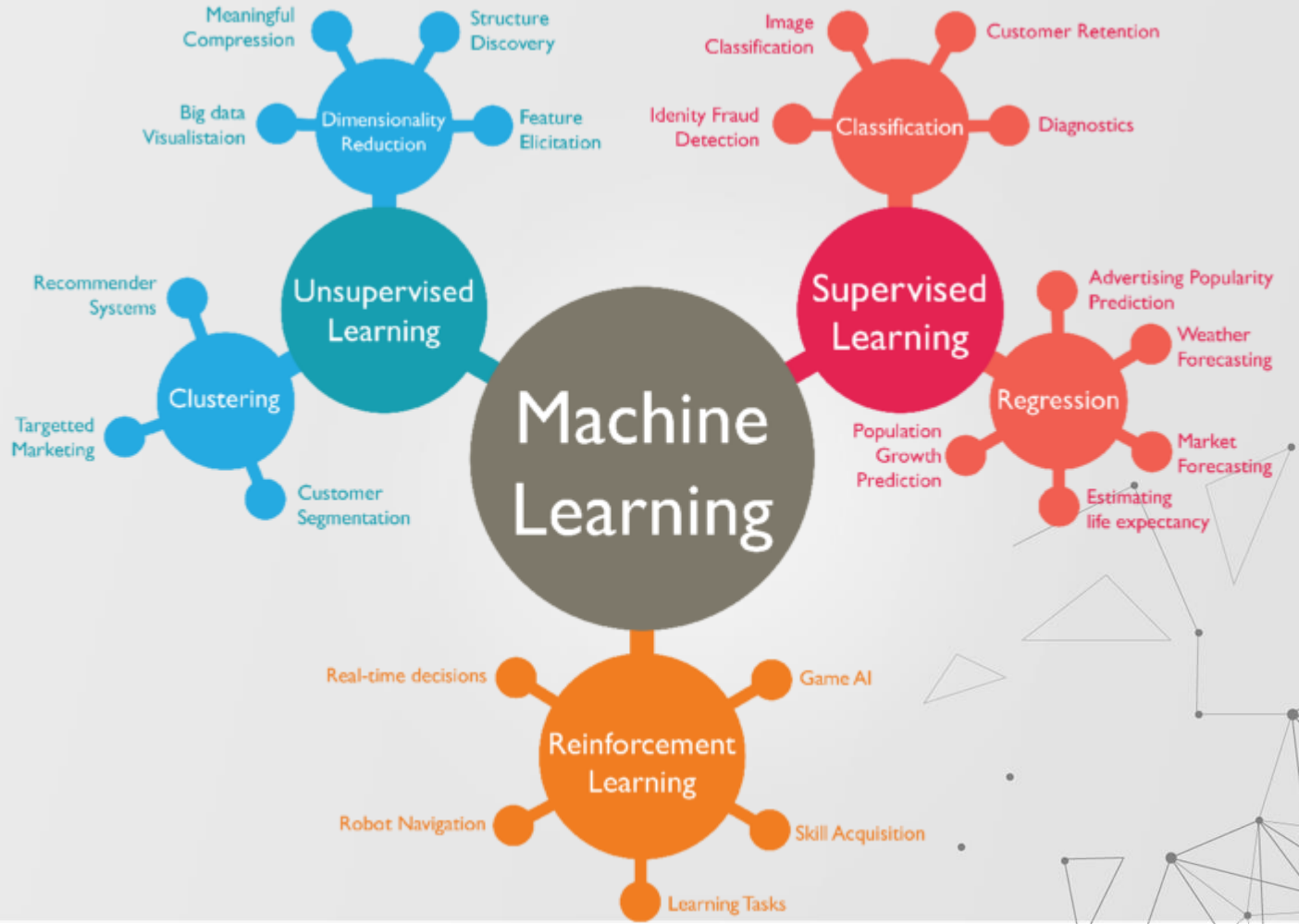
Machine Learning Pipelines

Fei GAO
Jan 2025



01

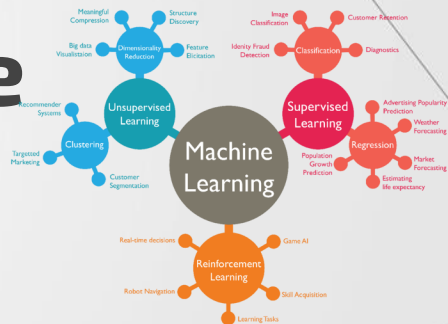
Machine Learning predictive data analysis



Données historique



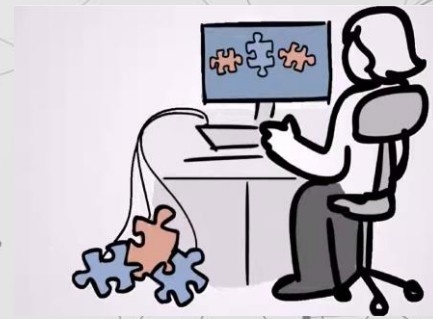
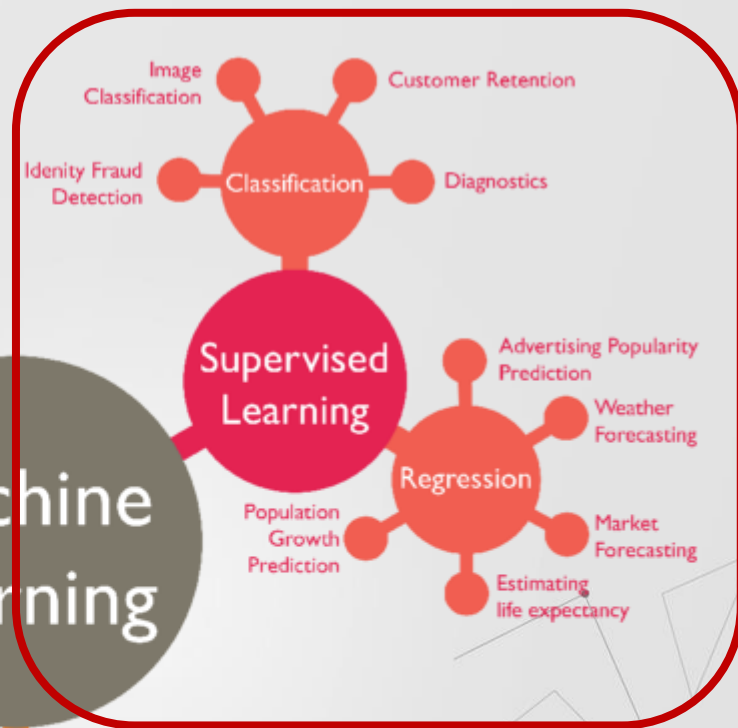
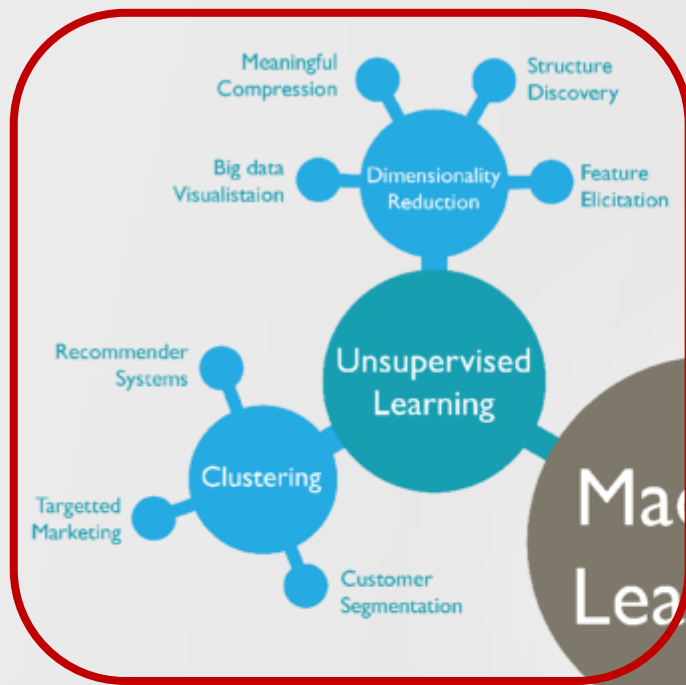
Mémorisation /
apprentissage



Nouvelles données



Généralisation /
Prédiction



APPRENTISSAGE SUPERVISÉ

Features + Label

$(X_1, X_2, X_3, X_4, \dots) (Y)$

APPRENTISSAGE NON SUPERVISÉ

Features + ~~Label~~

$(X_1, X_2, X_3, X_4, \dots) (\text{PAS d}'Y)$



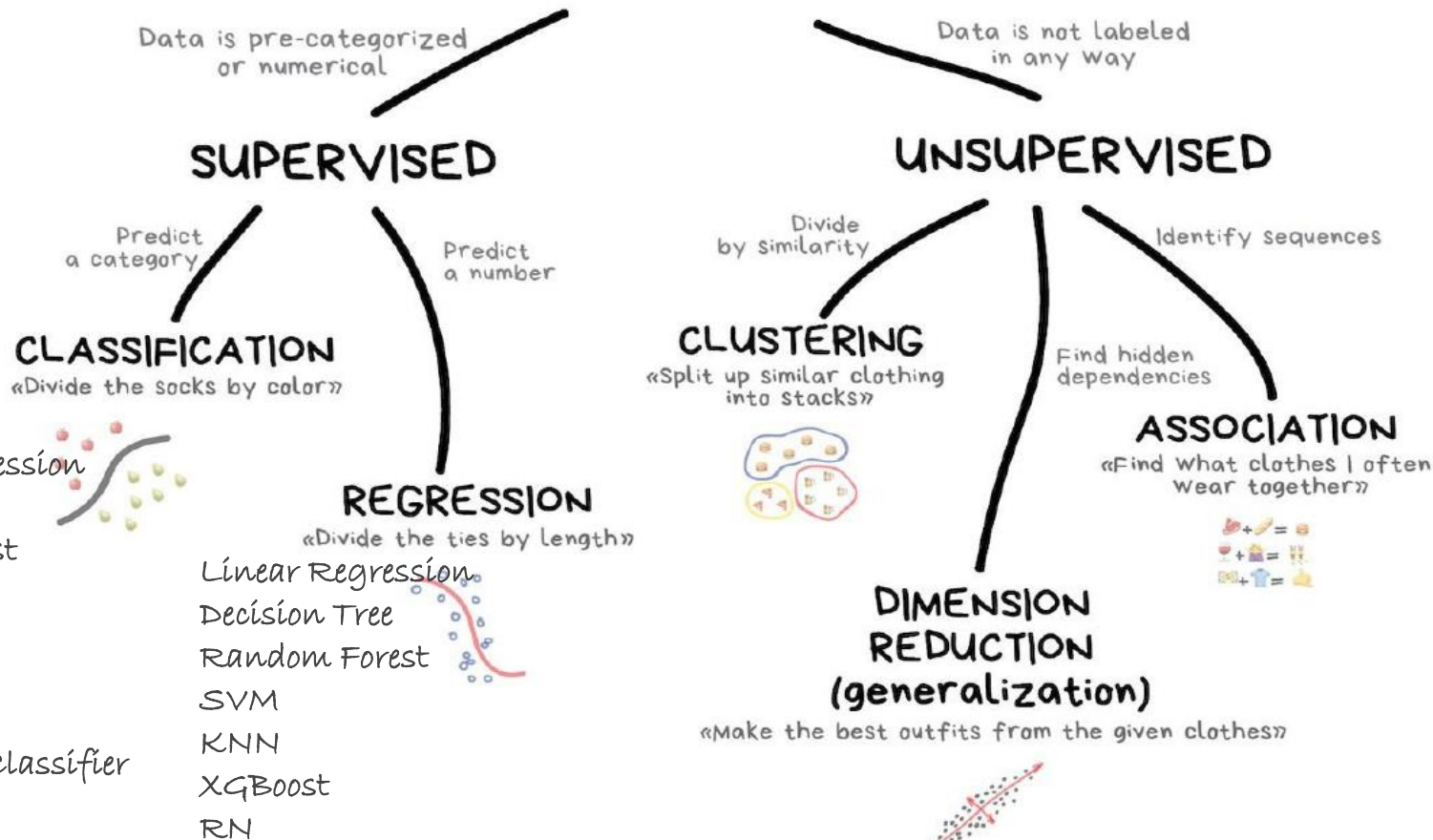
Problématique de l'analyse prédictive supervisée

- Y : label
- X_1, X_2, \dots : features
- Nous cherchons $f()$: une fonction qui essaie d'établir la relation $Y = f(X_1, X_2, \dots)$ **pour faire de la prédiction**
- $f()$ doit être aussi précise que possible pour pouvoir prédire sur les nouvelles données



Les principales familles

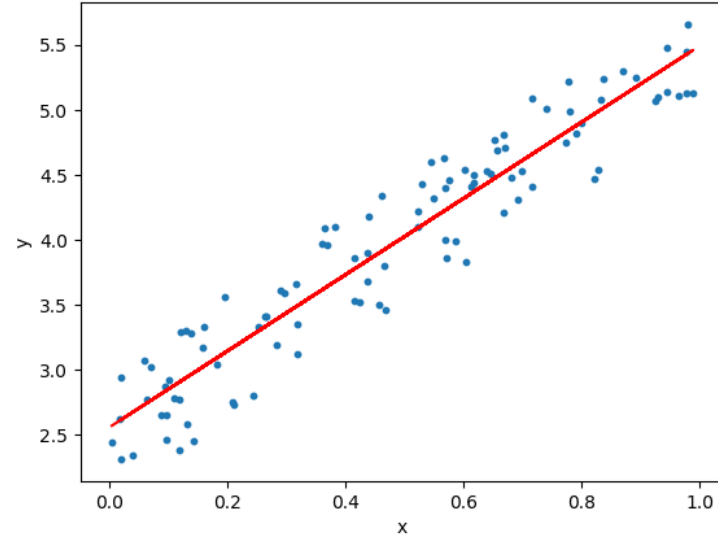
CLASSICAL MACHINE LEARNING



Usuel MACHINE LEARNING ALGORITHMS

Linear Regression

Linear Regression tends to establish a relationship between a dependent variable(Y) and one or more independent variable(X) by finding the best fit of the straight line.

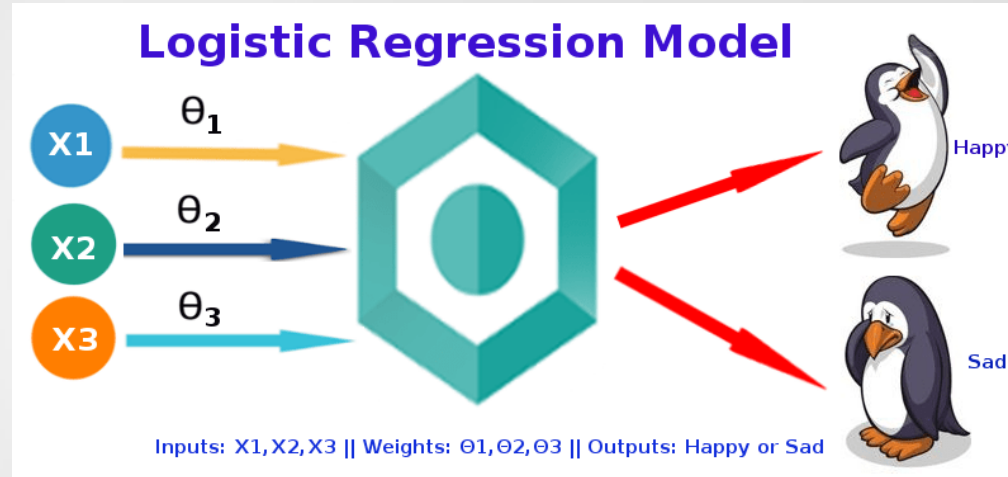


Regressor problems

Usuel MACHINE LEARNING ALGORITHMS

Logistic Regression

The logistic regression technique involves the dependent variable, which can be represented in the binary (0 or 1, true or false, yes or no) values or the probability of a successful or fail event.



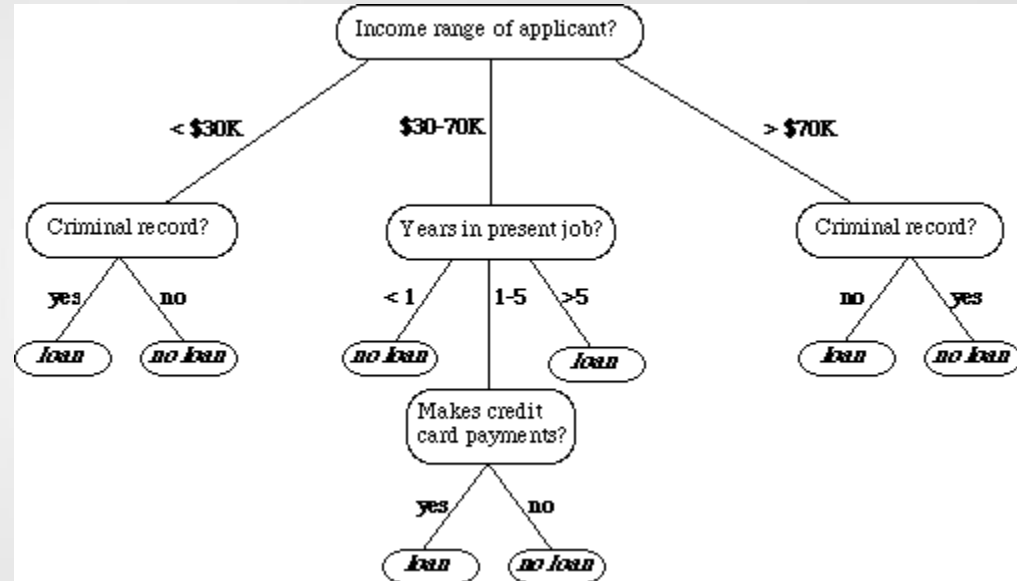
Classification problems

Usuel MACHINE LEARNING ALGORITHMS

Decision Tree

The decision tree works on an if-then statement.

Decision tree tries to solve a problem by using tree representation

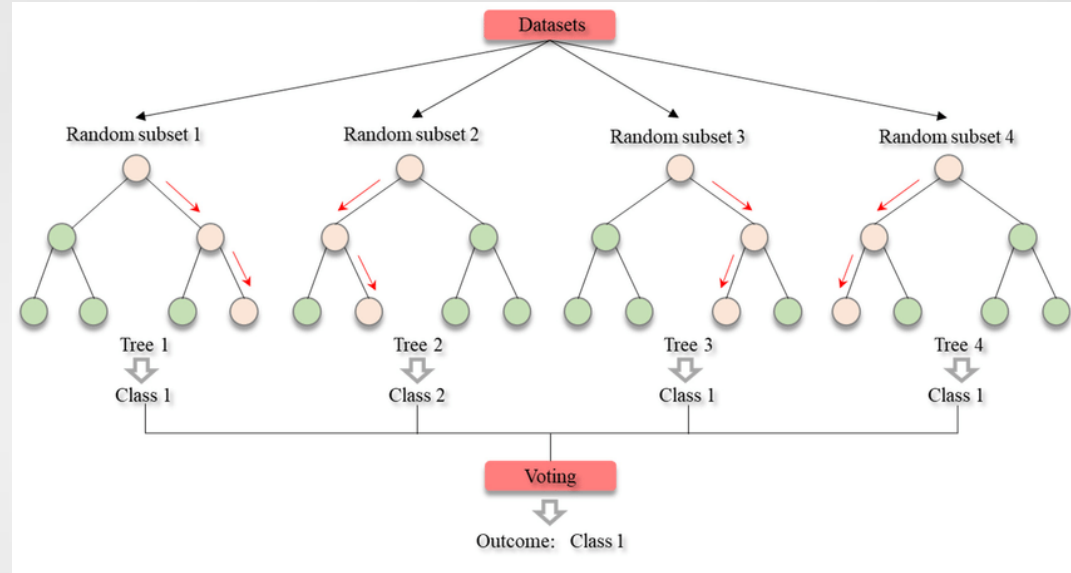


Classification as well as Regressor problems

Usuel MACHINE LEARNING ALGORITHMS

Random Forest

Random Forest is an ensemble machine learning algorithm that follows the bootstrapping technique. Random forest randomly selects a set of features that are used to decide the best split at each node of the decision tree.

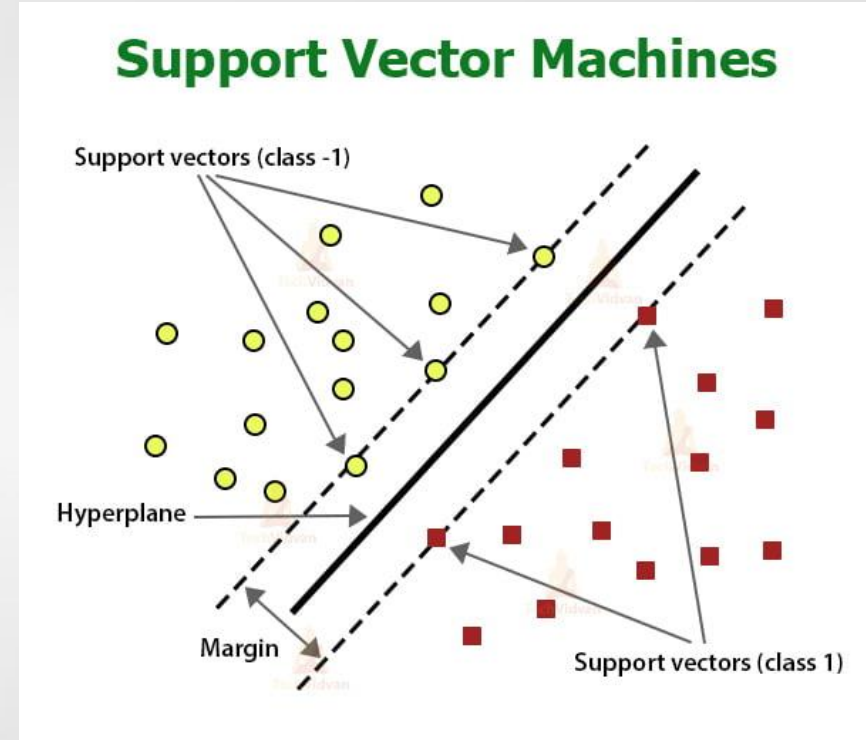


Classification as well as Regressor problems

Usuel MACHINE LEARNING ALGORITHMS

Support Vector Machine

SVM tries to find a line/hyperplane (in multidimensional space) that separates these two classes. Then it classifies the new point depending on whether it lies on the positive or negative side of the hyperplane depending on the classes to predict.



Classification as well as Regressor problems

Usuel MACHINE LEARNING ALGORITHMS

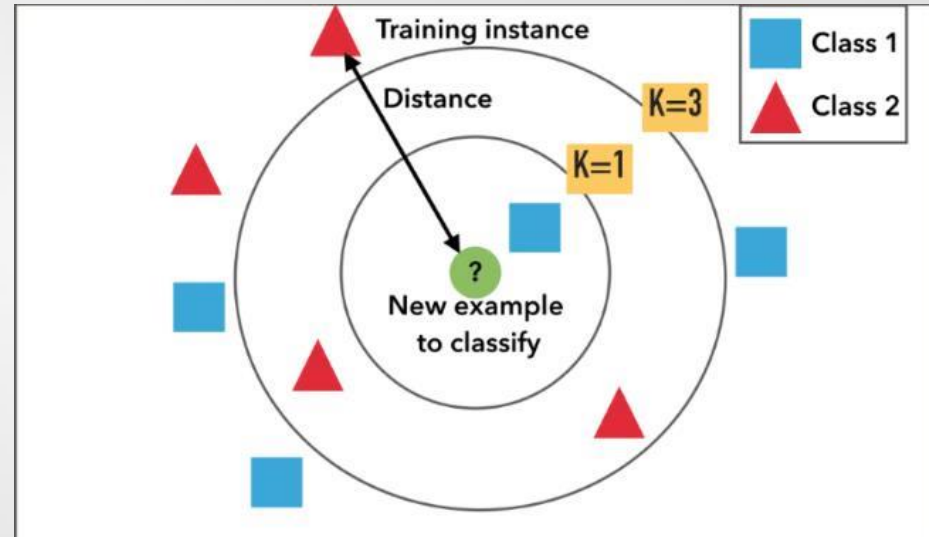
K Nearest Neighbor

In k-NN classification, the output is a class membership. An object is classified by a plurality vote of its neighbors, with the object being assigned to the class most common among its k nearest neighbors.

It does not create a generalized model during the time of training.

Testing is very costly.

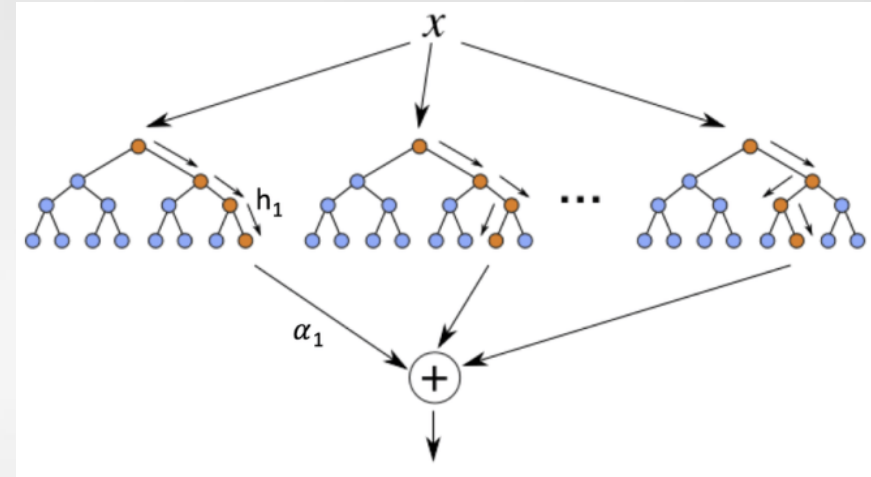
Classification as well as Regressor problems



Usuel MACHINE LEARNING ALGORITHMS

Gradient Boosting

XGBoost is a decision-tree-based ensemble Machine Learning algorithm that uses a gradient boosting framework. The approach consists in using a gradient descent to minimize the empirical risk



Classification as well as Regressor problems

Usuel MACHINE LEARNING ALGORITHMS

Naive Bayes Classifier

It is a classification technique based on Bayes' Theorem with an assumption of independence among predictors. In simple terms, a Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature.

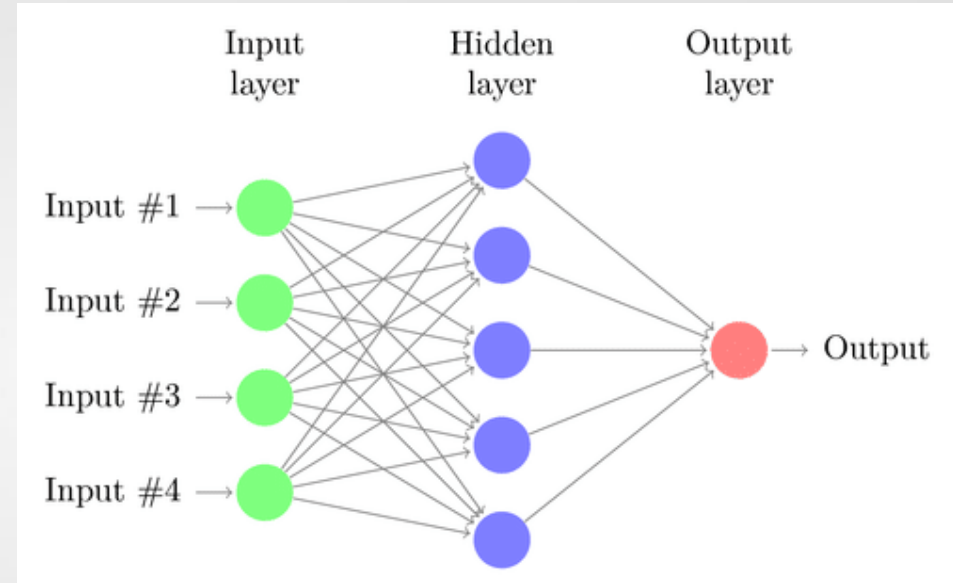
$$\begin{aligned}
 \text{Posterior} \quad P(A|B) &= \frac{P(B|A)P(A)}{P(B)} \\
 &\quad \text{Likelihood} \quad \text{Prior} \\
 &\quad \text{Normalizing constant} \\
 P(B) &= \sum_Y P(B|A)P(A)
 \end{aligned}$$

Classification problems

Usuel MACHINE LEARNING ALGORITHMS

Neural Network

As neural suggests, they are braininspired systems which are intended to replicate the way that we humans learn. NNs consist of input and output layers, as well as a hidden layer consisting of units that transform the input.



Classification as well as Regressor problems

Pourquoi existe-t-il autant d'algorithmes ?

Le « théorème » du « No Free Lunch »

Il n'existe pas d'algorithme qui soit le meilleur quelque soit le problème d'apprentissage...



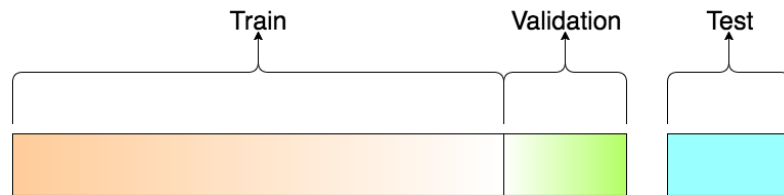
Schéma typique (standard) de l'analyse prédictive supervisée

Y : variable cible

X1, X2, ... : variables explicatives

f(.) une fonction qui essaie d'établir la relation $Y = f(X1, X2, \dots)$

f(.) doit être « aussi précise que possible »...



Construction de la fonction f(.) à partir des données d'apprentissage

Training set

A small thumbnail image of a data table with multiple columns and rows, representing the training set.

$$Y = f(X1, X2, \dots) + \epsilon$$

Application du modèle (prédiction) sur l'ensemble de test

Validation set

A small thumbnail image of a data table with multiple columns and rows, representing the validation set.

$$(Y, \hat{Y})$$

Y : valeurs observées

\hat{Y} : valeurs prédites par f(.)

Mesures de **performances** par confrontation entre Y et \hat{Y} : matrice de confusion + mesures

Age	Sexe	Numéro de la carte	Montant	Statut
1	M	123456789	100	positive
2	F	987654321	200	positive
3	M	567890123	300	positive
4	F	456789012	400	positive
5	M	345678901	500	positive
6	F	234567890	600	positive
7	M	123456789	700	positive
8	F	987654321	800	positive
9	M	567890123	900	positive
10	F	456789012	1000	positive
11	M	345678901	1100	positive
12	F	234567890	1200	positive
13	M	123456789	1300	positive
14	F	987654321	1400	positive
15	M	567890123	1500	positive
16	F	456789012	1600	positive
17	M	345678901	1700	positive
18	F	234567890	1800	positive
19	M	123456789	1900	positive
20	F	987654321	2000	positive
21	M	567890123	2100	positive
22	F	456789012	2200	positive
23	M	345678901	2300	positive
24	F	234567890	2400	positive
25	M	123456789	2500	positive
26	F	987654321	2600	positive
27	M	567890123	2700	positive
28	F	456789012	2800	positive
29	M	345678901	2900	positive
30	F	234567890	3000	positive
31	M	123456789	3100	positive
32	F	987654321	3200	positive
33	M	567890123	3300	positive
34	F	456789012	3400	positive
35	M	345678901	3500	positive
36	F	234567890	3600	positive
37	M	123456789	3700	positive
38	F	987654321	3800	positive
39	M	567890123	3900	positive
40	F	456789012	4000	positive
41	M	345678901	4100	positive
42	F	234567890	4200	positive
43	M	123456789	4300	positive
44	F	987654321	4400	positive
45	M	567890123	4500	positive
46	F	456789012	4600	positive
47	M	345678901	4700	positive
48	F	234567890	4800	positive
49	M	123456789	4900	positive
50	F	987654321	5000	positive
51	M	567890123	5100	positive
52	F	456789012	5200	positive
53	M	345678901	5300	positive
54	F	234567890	5400	positive
55	M	123456789	5500	positive
56	F	987654321	5600	positive
57	M	567890123	5700	positive
58	F	456789012	5800	positive
59	M	345678901	5900	positive
60	F	234567890	6000	positive
61	M	123456789	6100	positive
62	F	987654321	6200	positive
63	M	567890123	6300	positive
64	F	456789012	6400	positive
65	M	345678901	6500	positive
66	F	234567890	6600	positive
67	M	123456789	6700	positive
68	F	987654321	6800	positive
69	M	567890123	6900	positive
70	F	456789012	7000	positive
71	M	345678901	7100	positive
72	F	234567890	7200	positive
73	M	123456789	7300	positive
74	F	987654321	7400	positive
75	M	567890123	7500	positive
76	F	456789012	7600	positive
77	M	345678901	7700	positive
78	F	234567890	7800	positive
79	M	123456789	7900	positive
80	F	987654321	8000	positive
81	M	567890123	8100	positive
82	F	456789012	8200	positive
83	M	345678901	8300	positive
84	F	234567890	8400	positive
85	M	123456789	8500	positive
86	F	987654321	8600	positive
87	M	567890123	8700	positive
88	F	456789012	8800	positive
89	M	345678901	8900	positive
90	F	234567890	9000	positive
91	M	123456789	9100	positive
92	F	987654321	9200	positive
93	M	567890123	9300	positive
94	F	456789012	9400	positive
95	M	345678901	9500	positive
96	F	234567890	9600	positive
97	M	123456789	9700	positive
98	F	987654321	9800	positive
99	M	567890123	9900	positive
100	F	456789012	10000	positive

Metrics Regression

1

MAE : Mean Average Error

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

2

RMSE : Root Mean Squared Error

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

3

Root Mean Squared Logarithmic Error

$$RMSLE = \sqrt{\frac{1}{n} \sum (\log(1 + \text{prédiction}) - \log(1 + \text{cible}))^2}$$

4

Autres mesures

Mean Squared Error, Weigthed Mean Average Error, ...

$$MSE = \frac{1}{n} \sum_{i=1}^n (\text{prédiction}_i - \text{cible}_i)^2$$

Metrics Classification

1. Accuracy, Precision, Recall...

$$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{FP} + \text{FN} + \text{TN})$$

Precision = $\text{TP} / (\text{TP} + \text{FP})$: is the proportion of all positive cases that were correctly classified => **FP intolerable**

Sensitivity/Recall = $(\text{TP}) / (\text{TP} + \text{FN})$ => **FN intolerable**

Specificity = $\text{TN} / (\text{TN} + \text{FP})$: is the proportion of all negative cases that were correctly classified

False positive rate (FPR) = $1 - \text{specificity} = 1 - (\text{TN} / (\text{TN} + \text{FP}))$: is the proportion of false positives among

2

$$F_1 = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}}$$

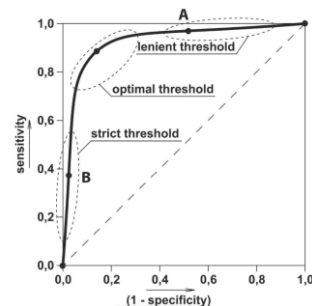
Plus d'information :

<https://towardsdatascience.com/the-5-classification-evaluation-metrics-you-must-know-aa97784ff226>

3. AUC

AUC ROC indicates how well the probabilities from the positive classes are separated from the negative classes => Allowing threshold selection

		Actual	
		Positive	Negative
Predicted	Positive	True Positive	False Positive
	Negative	False Negative	True Negative



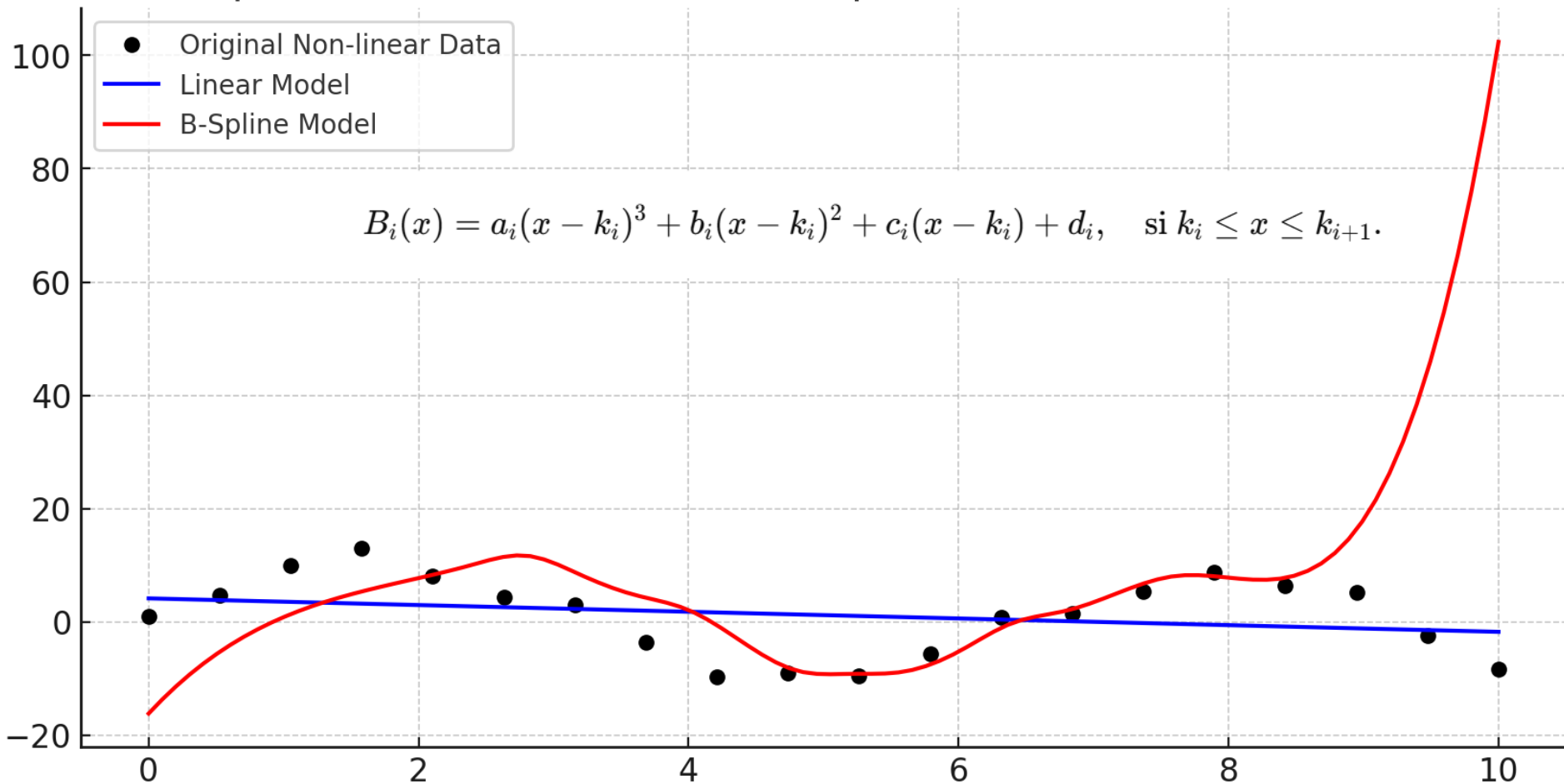
Interaction

- Une promotion (variable 1) augmente les ventes de 10 unités.
- La publicité (variable 2) augmente les ventes de 15 unités.
- Cependant, si vous combinez une promotion **et** une publicité, les ventes augmentent de 40 unités.
- Ici, l'effet combiné (40) est supérieur à la somme des effets individuels (10 + 15 = 25).

Modèle : $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \epsilon$

Si $\beta_3 > 0$, cela signifie qu'il existe une interaction amplificatrice entre x_1 et x_2 .

Comparison: Linear Model vs. B-Spline Model on Non-linear Data



Méthode	Quand l'utiliser
Interactions	Quand il existe une relation forte entre deux variables (effets combinés) ou à tester avec une hypothèse.
Polynomial Features	Quand la relation non linéaire est globale et simple (par exemple, courbe quadratique ou cubique).
Splines	Quand la relation non linéaire est complexe et locale (varie dans différentes parties des données).

Processus général d'un problème de data science

1

Définition de la cible et des variables explicatives

- Identification ou construction de la cible
- Identification des variables explicatives

2

Construction du jeu de données

- Filtres de la population concernée
- Exploration et visualisation

3

Préparation des données

- Feature Engineering : création de variables explicatives **metier et DS**
- Traitement des valeurs manquantes
- Encodage numérique des variables catégoriques **et La mise à l'échelle**
- Séparation des données en échantillons

4

Entraînement du modèle

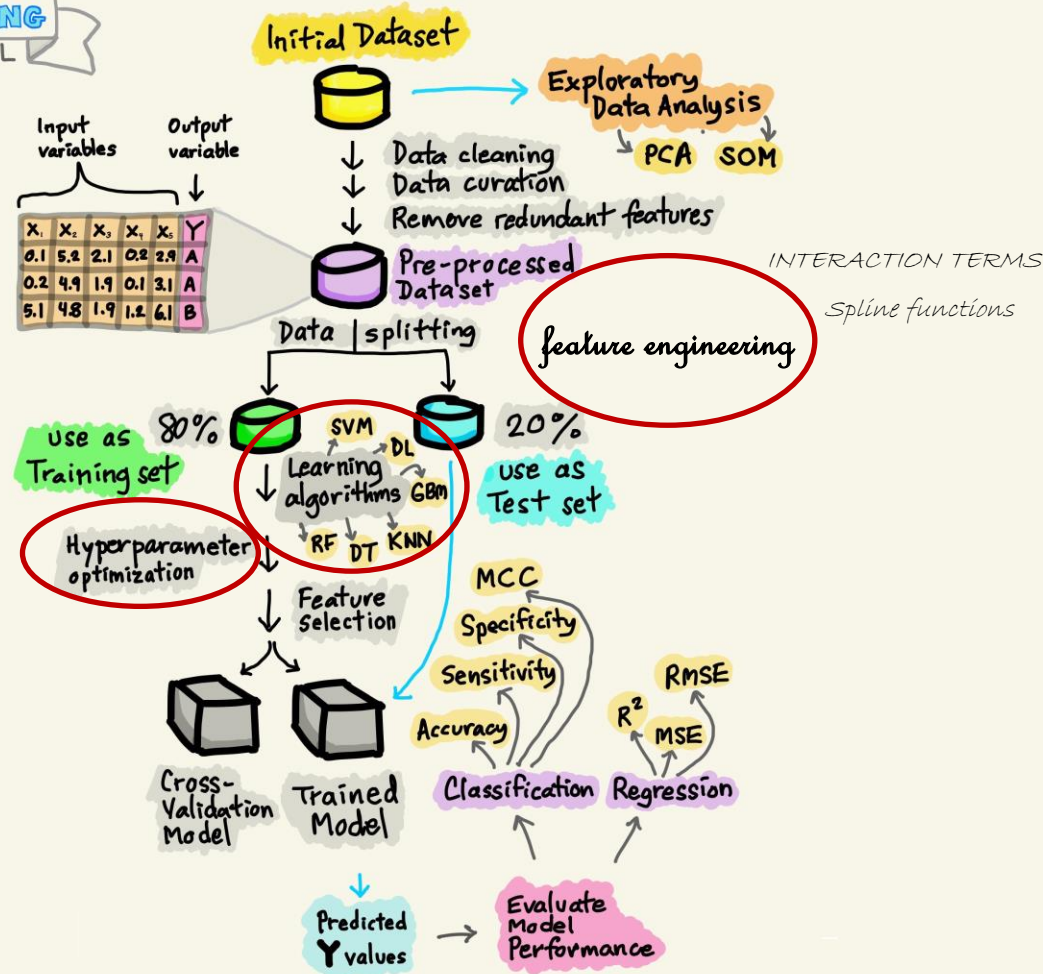
- Sélection du meilleur jeu de paramètres
- Entraînement du modèle



Résultats et performance du modèle

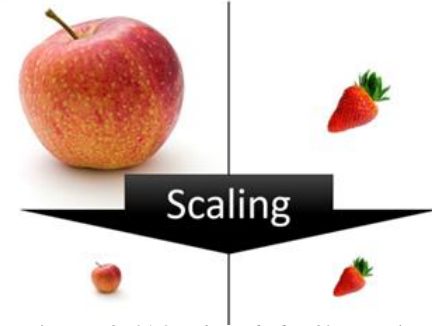
- Application du modèle sur l'échantillon de test
- Etude des caractéristiques clés du modèle

BUILDING THE MACHINE LEARNING MODEL



La mise à l'échelle

Name	Weight	Price
Orange	15	1
Apple	18	3
Banana	12	2
Grape	10	5



La mise à l'échelle

Sr No.	Algorithms	Feature Scaling
1.	Linear/Non-Linear Regressions	Yes
2.	Logistic Regression	Yes
3.	KNN	Yes
4.	SVM	Yes
5.	Neural Networks	Yes
6.	K-means clustering	Yes
7.	CART	No
8.	Random Forests	No
9.	Gradient Boosted Decision Trees	No
10.	Naïve Bayes	NO
11.	PCA	Yes
12.	SVD	Yes
13.	Factorization Machines	Yes

Fig: Feature Scaling requires based on Machine Learning

Schéma typique (standard) de l'analyse prédictive supervisée

Y : variable cible

X1, X2, ... : variables explicatives

f(.) une fonction qui essaie d'établir la relation $Y = f(X1, X2, \dots)$

f(.) doit être « aussi précise que possible »...



Ensemble
d'apprentissage
Training set

Construction de la fonction f(.) à
partir des données d'apprentissage

A small thumbnail image of a dataset table with multiple columns and rows of data.

$$Y = f(X1, X2, \dots) + \epsilon$$

Application du modèle (prédiction)
sur l'ensemble de test

Test set

A small thumbnail image of a dataset table, similar to the one in the training set section.

$$(Y, \hat{Y})$$

Mesures de **performances**
par confrontation entre Y
et \hat{Y} : matrice de
confusion + mesures

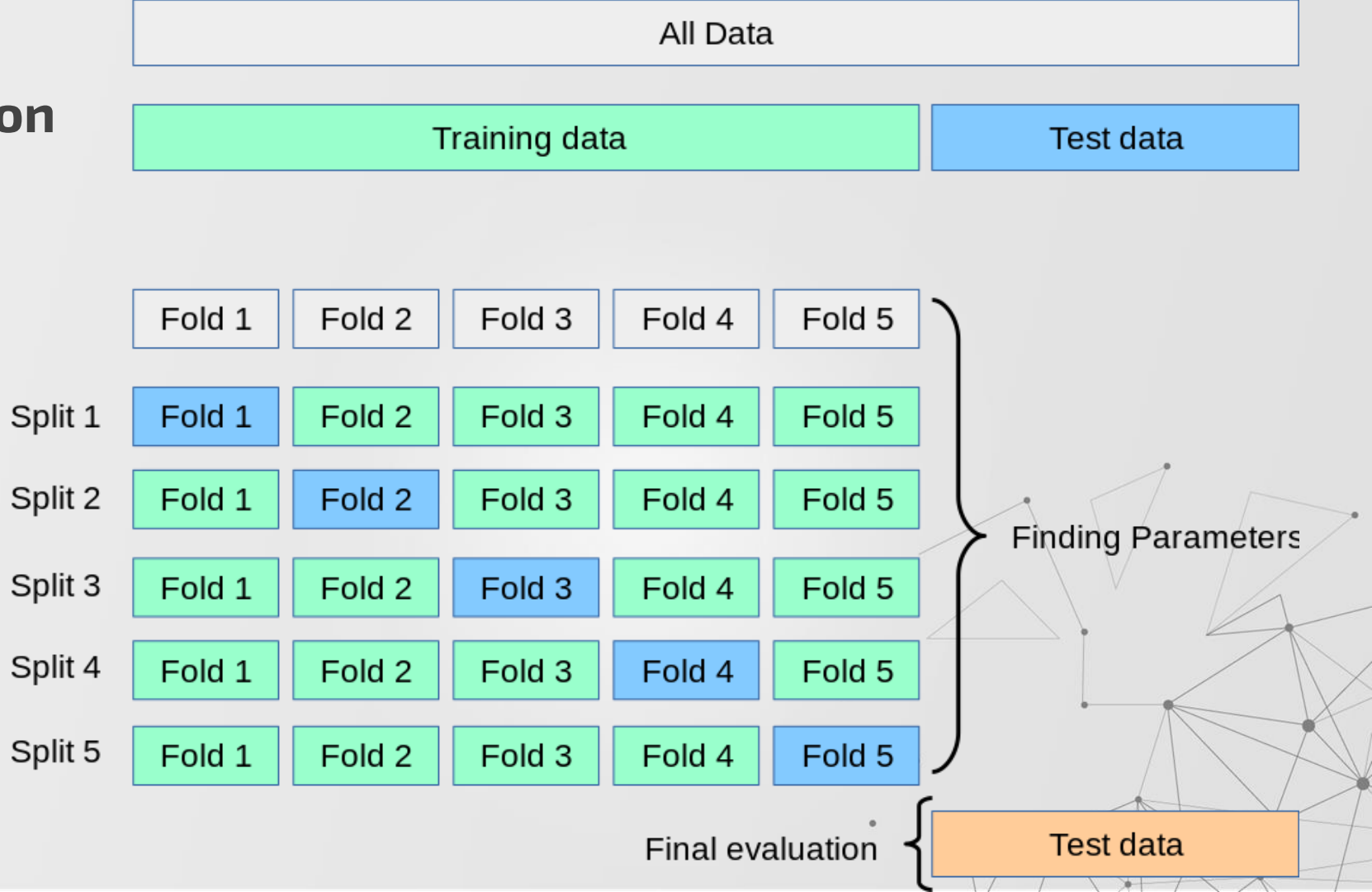
Y : valeurs observées
 \hat{Y} : valeurs prédites par f(.)

5	F	29.11.1	1.621	50	145	0	positive
6	M	29.78.8	0.205	25	89	0	negative
7	M	33.3	0.272	32	105	0	positive
8	M	22.28	0.107	25	89	94	negative
9	M	34.43	0.206	31	117	100	positive
10	F	2.28	0.205	30	110	0	negative
11	M	40.50	0.206	33	107	0	positive
12	M	0	0.202	30	120	0	negative
13	M	0	0.206	31	100	0	negative
14	M	0	0.202	30	120	0	negative
15	M	0	0.206	31	100	0	negative
16	M	0.17	0.405	57	170	0	negative
17	M	22.28	0.206	30	100	100	positive
18	M	10.70	0.507	51	100	170	positive
19	M	0	0.206	32	100	0	positive
20	M	47.40	0.505	56	110	200	positive
21	M	2.28	0.206	31	107	0	positive
22	M	30.43	0.105	30	100	0	negative
23	M	20.50	0.206	32	120	0	positive
24	M	40.50	0.706	27	150	200	negative
25	M	1.20	0.206	30	90	0	positive
26	M	0.50	0.405	40	100	0	positive
27	M	30.50	0.206	31	100	0	positive
28	M	20.50	0.706	31	140	140	positive
29	M	20.50	0.206	40	120	120	positive
30	M	0.50	0.207	40	140	0	positive
31	M	10.70	0.407	27	90	140	negative
32	M	10.70	0.206	32	140	140	negative
33	M	0.50	0.207	30	110	0	negative

Ensemble
de données
(dataset)

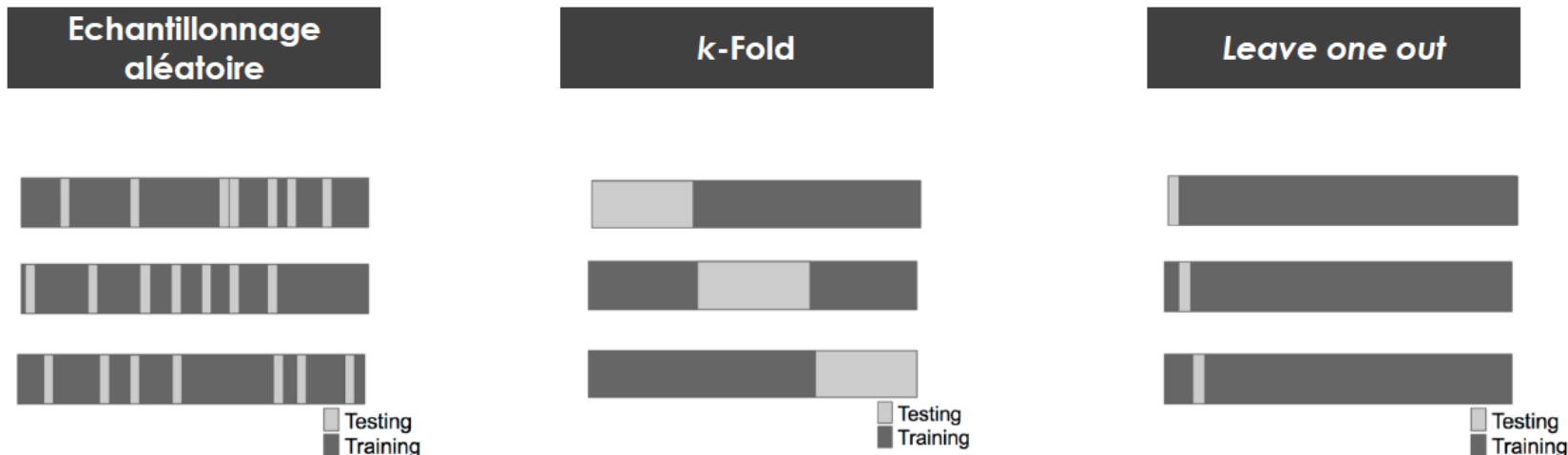
Ensemble de test

Cross validation



Model cross-validation

Différentes méthodes pour s'assurer que le modèle « apprend » bien, qu'il « généralise »

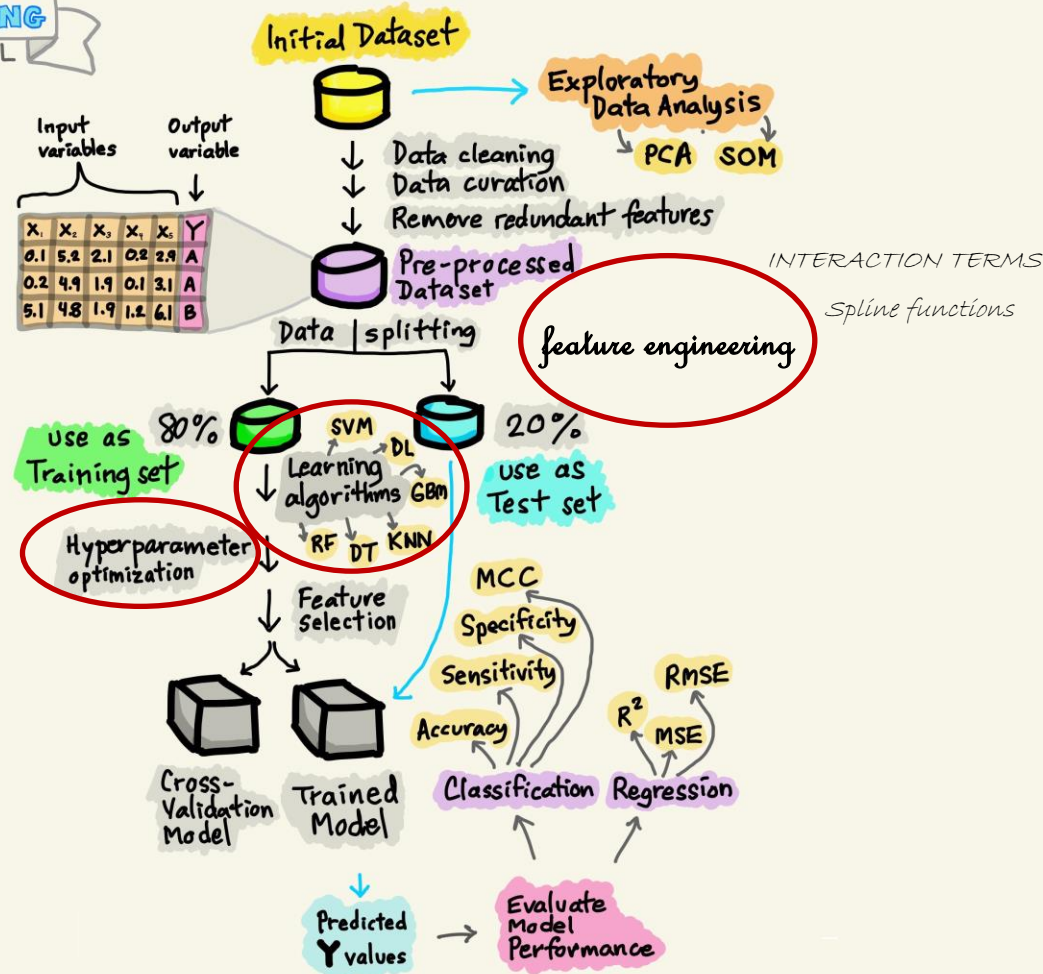


Idée de la validation croisée pour comparer la performance des modèles, indépendamment des « individus » utilisés pour apprendre (risque si l'on utilise 1 seul jeu d'apprentissage).

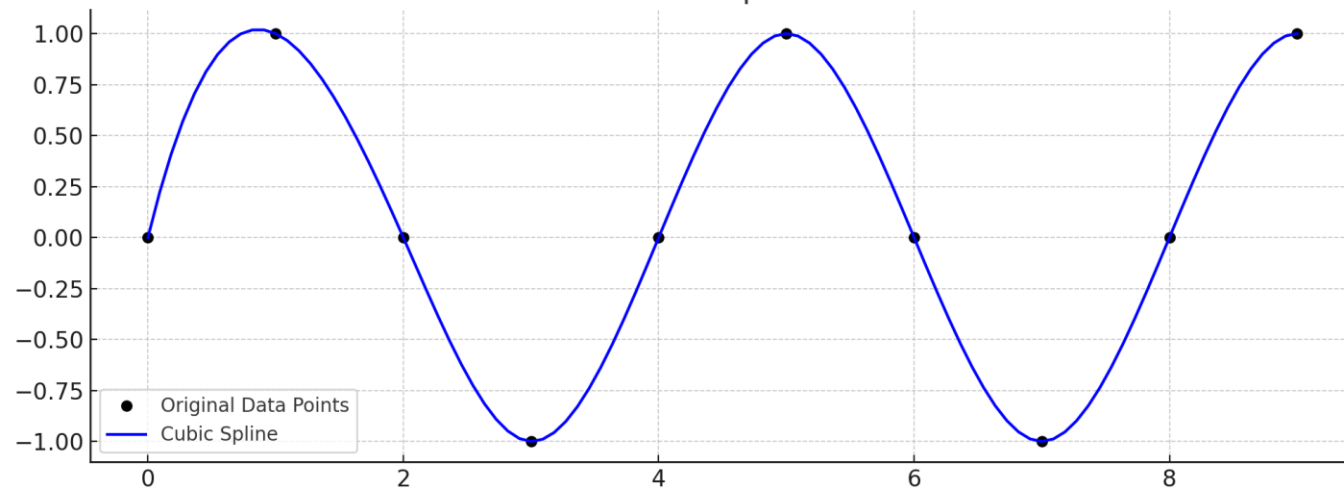
Habituellement, on utilise entre 5 ou 10 « folds ».



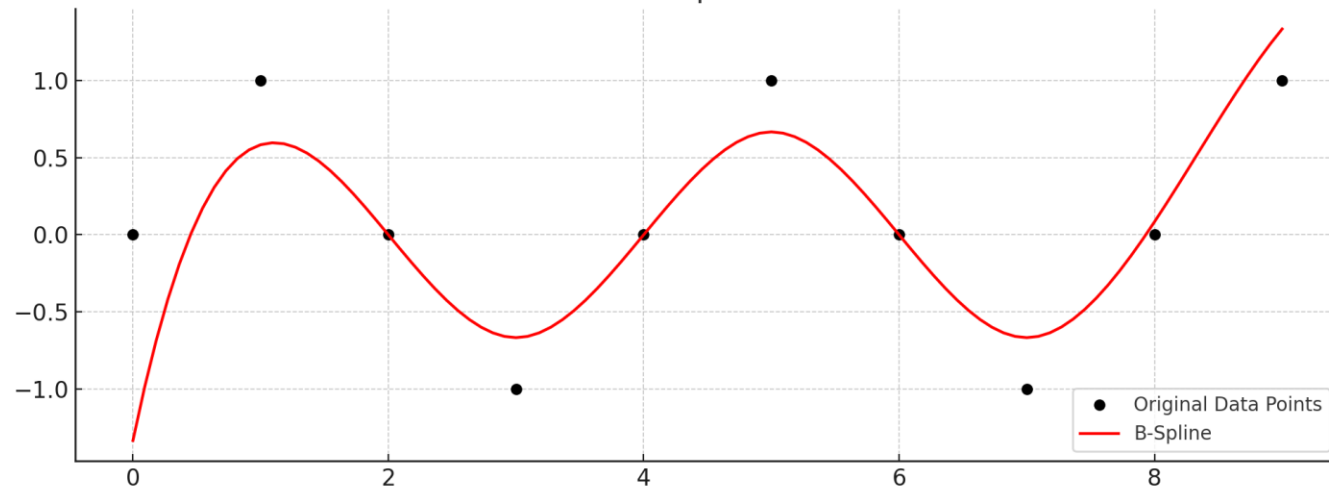
BUILDING THE MACHINE LEARNING MODEL



Cubic Spline



B-Spline



B-spline de degré 3 (cubic B-spline)

Si vous avez des nœuds $t = [1, 2, 3, 4, 5]$, une cubic B-spline est définie par la formule récursive :

$$B_{i,3}(x) = \frac{x - t_i}{t_{i+3} - t_i} B_{i,2}(x) + \frac{t_{i+4} - x}{t_{i+4} - t_{i+1}} B_{i+1,2}(x),$$

où $B_{i,2}(x)$ est une B-spline de degré 2.

Cubic spline (non-B-spline)

Pour les mêmes nœuds, une cubic spline s'écrit explicitement comme :

$$f(x) = a_i(x - t_i)^3 + b_i(x - t_i)^2 + c_i(x - t_i) + d_i, \quad \text{si } t_i \leq x \leq t_{i+1}.$$

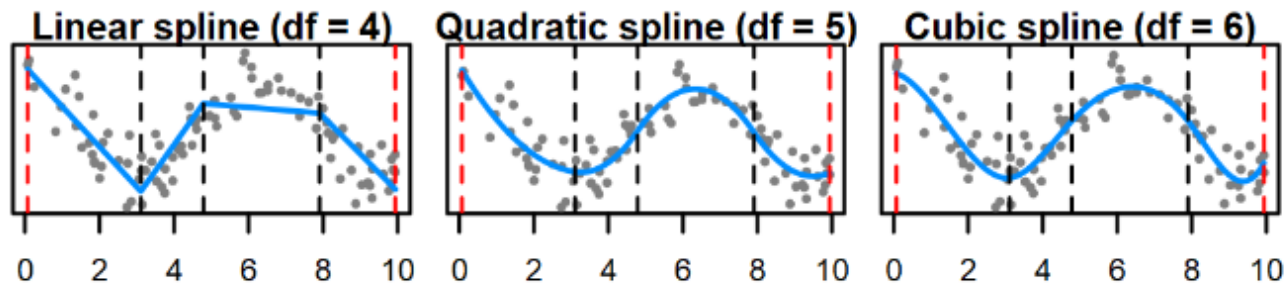
Les coefficients a_i, b_i, c_i, d_i sont calculés en résolvant un système linéaire basé sur les conditions de continuité.

$$B_i(x) = a_i(x - k_i)^3 + b_i(x - k_i)^2 + c_i(x - k_i) + d_i, \quad \text{si } k_i \leq x \leq k_{i+1}.$$

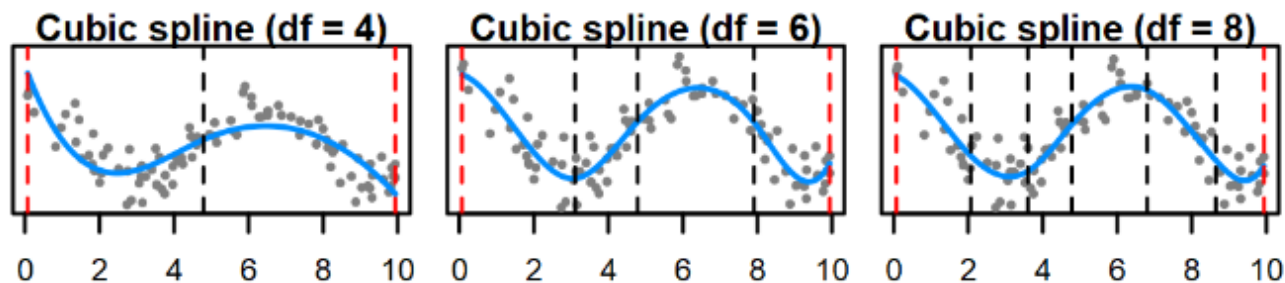
For B-splines: $\text{df} = k + \text{degree}$ if you specify the knots or $k = \text{df} - \text{degree}$ if you specify the degrees of freedom and the degree. For natural (restricted) cubic splines: $\text{df} = k - 1$ if you specify the knots or $k = \text{df} + 1$ if you specify the degrees of freedom.

As an example: A cubic spline ($\text{degree} = 3$) with 4 (internal) knots will have $\text{df} = 4 + 3 = 7$ degrees of freedom. Or: A cubic spline ($\text{degree} = 3$) with 5 degrees of freedom will have $k = 5 - 3 = 2$ knots.

Let's see some illustrations. In the scatterplots below you see some artificial data together with the spline fits of different degrees but the same amount of knots ($k = 3$). The knots are indicated by dashed vertical lines (Boundary knots by red dashed lines) and are placed at the 25th, 50th and 75th percentile of x . The first plot shows a linear spline (degree = 1), the second one a quadratic spline (degree = 2) and the third is a cubic spline with degree = 3.



In the next plot, you see three cubic splines with different degrees of freedom. As before, the knots are shown as dashed vertical lines. With increasing degrees of freedom, the number of knots gets larger (from 1 to 3 to 5). The spline gets wigglier although the difference is only really noticeable between the first and second plot.



spline function procedures

A cubic spline function, with three knots (τ_1, τ_2, τ_3) will have 7 degrees of freedom. Using representation given in Eq. [2](#), the function can be written as:

$$f(X) = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \beta_4 (X - \tau_1)^3 + \beta_5 (X - \tau_2)^3 + \beta_6 (X - \tau_3)^3$$



Detecting Interaction Effects



REF :

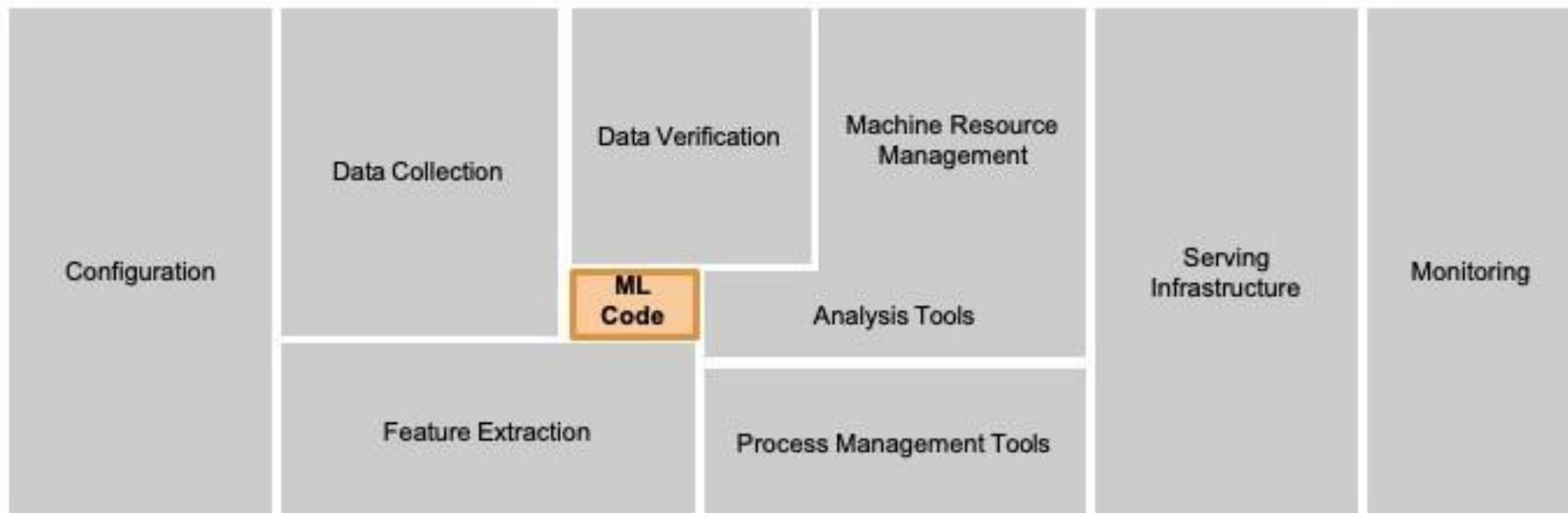
<https://bookdown.org/max/FES/detecting-interaction-effects.html>



Excercise



The Requirements Surrounding ML Infrastructure



Expressions régulières (1/1)

- **ion** : recherche les mots qui contiennent la chaîne "ion", dans n'importe quelle position
- **ion\$** : mots se terminant par "ion" (\$=fin de mot)
- **^anti** : recherche tous les mots commençant par "anti" (^=début de mot)
- **^maison\$** : recherche exactement le mot "maison"
- **p.r** : recherche les mots qui contiennent un "p", suivi d'une lettre quelconque, puis d'un "r" (le point correspond à n'importe quel caractère)
- **^p...r\$** : mots commençant par "p", suivi de trois lettres quelconques, et finissant par "r" (le symbole . dans une regex correspond à n'importe quel caractère)
- **^p.*r\$** : mots commençant par "p" et finissant par "r" (*= répétitions – 0 ou plusieurs fois – du caractère précédent, ici '.', donc n'importe quel caractère)
- **oid|ion|ein** : recherche les mots qui contiennent (au moins) une des trois chaînes "iod", "ion" ou "ein" (| = ou).

Source : http://www.lexique.org/?page_id=101
<https://buzut.net/la-puissance-des-regex/>

Expressions régulières (2/2)

- `[A-Za-z]` : n'importe quoi comme caractère, majuscule ou minuscule
- `+` : 1 ou plus
- `()` : contenu à extraire
- `\w` : un caractère alphanumérique ou un tiret de soulignement (tiret de 8).

Source : http://www.lexique.org/?page_id=101
<https://buzut.net/la-puissance-des-regex/>



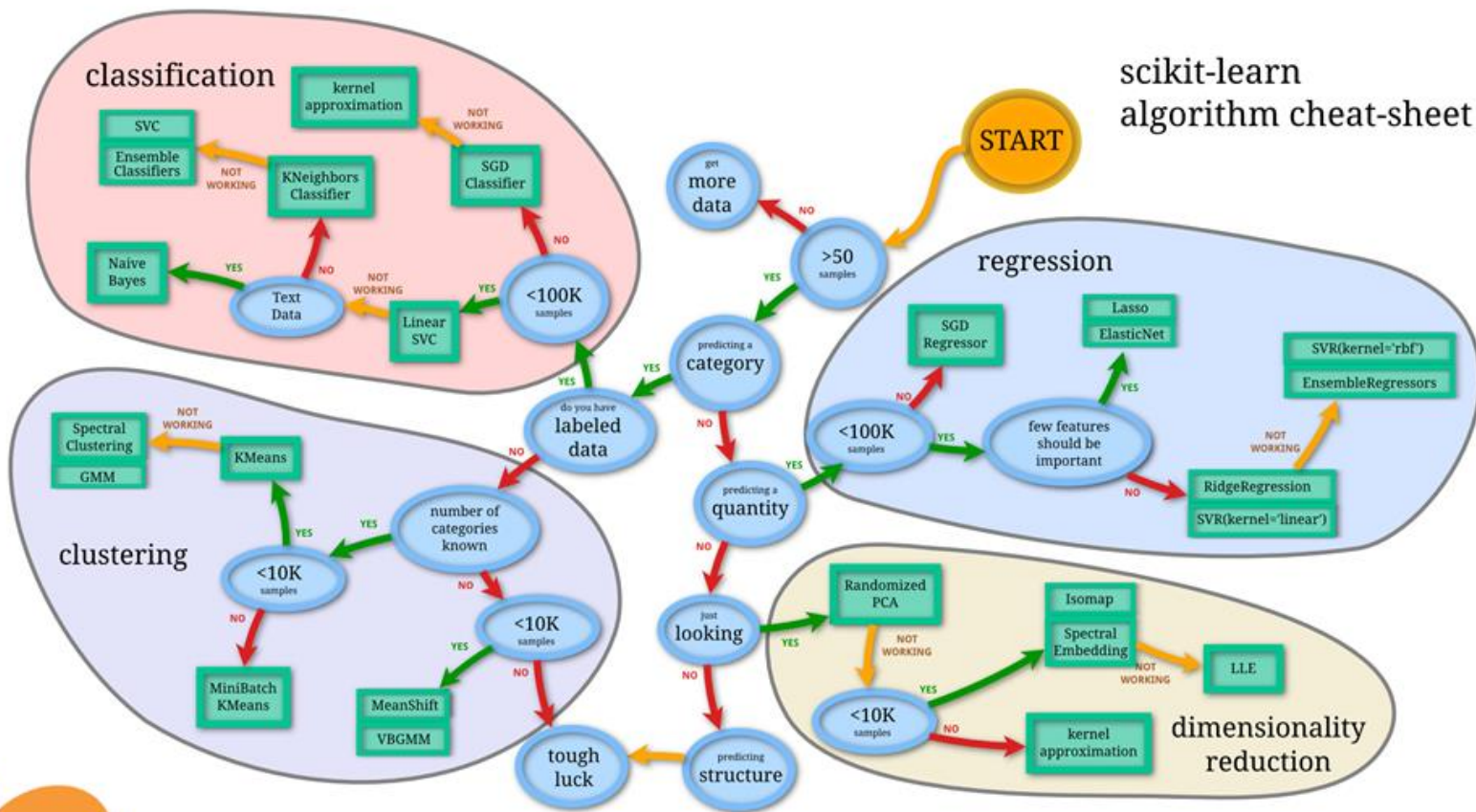
Python vs. R: What's the Difference?

**What is the difference between
Anaconda Prompt and Anaconda
Powershell Prompt?**



Choisir son algorithme

scikit-learn
algorithm cheat-sheet



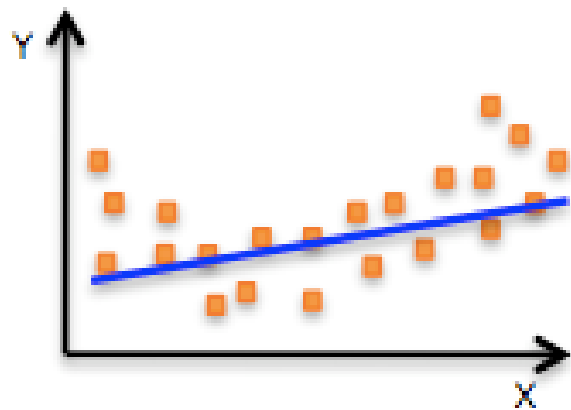
Back



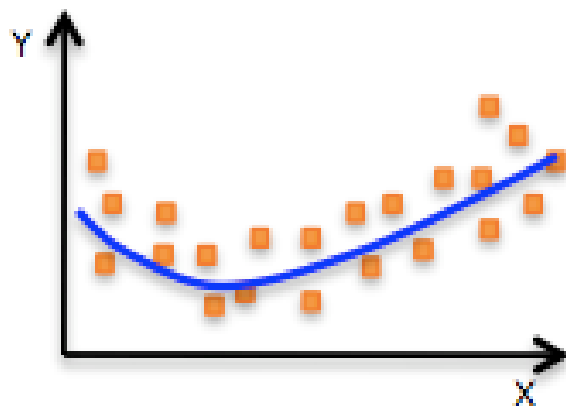
THANKS

Does anyone have any questions?

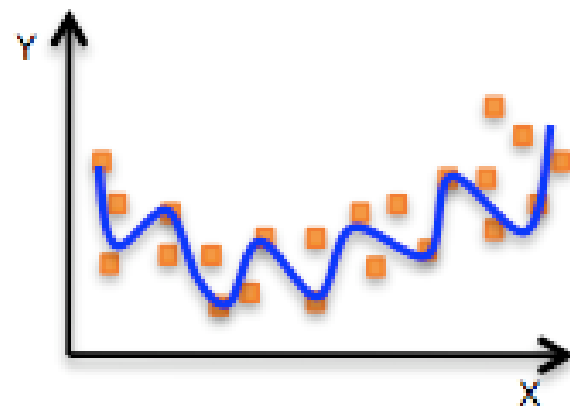
CREDITS: Fei GAO (EHESP)



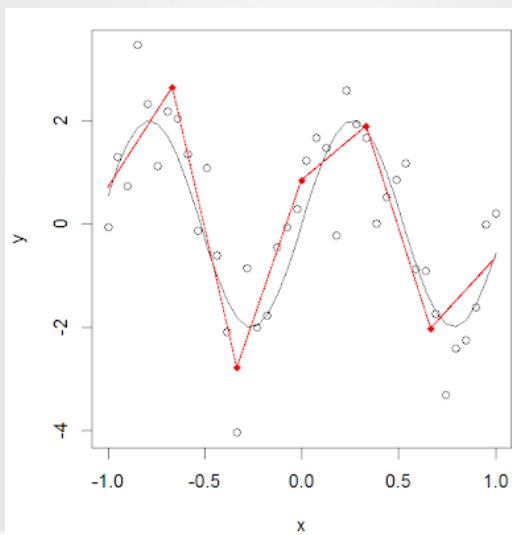
Underfitting

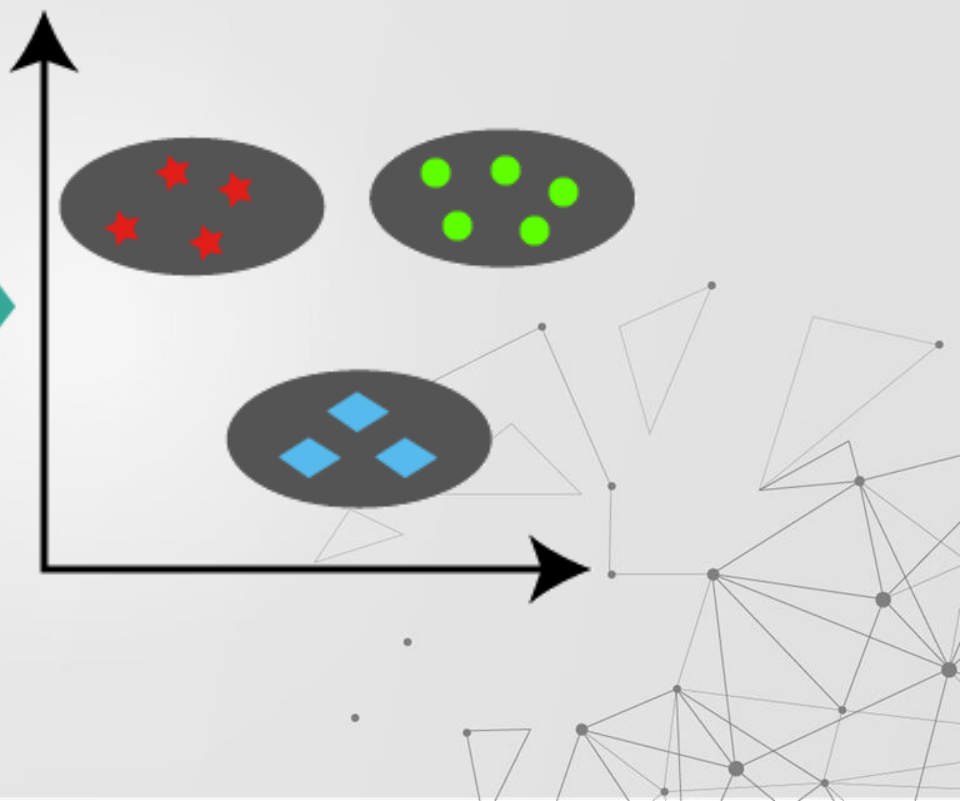
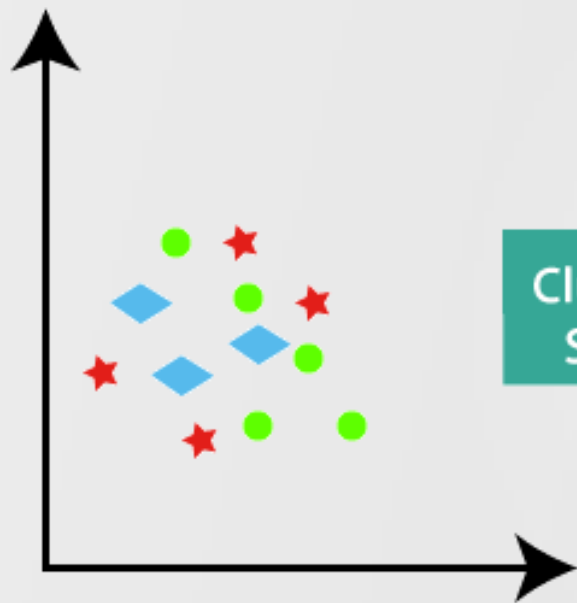


Just right!



overfitting





Créez votre compte Databricks

1/2

Prénom

Fei

Nom

GAO

E-mail professionnel

constance.fe.gao.pro@gmail.com

Entreprise

hdh

Intitulé de poste

DS

Numéro de téléphone (facultatif)

0661213981

Pays

France

- ☐ Oui, je souhaite recevoir les communications marketing concernant les services, les événements et les produits open source de Databricks. Je suis conscient(e) que je peux modifier [mes préférences](#) à tout moment.

Continuer

Sélectionnez un fournisseur de cloud

2/2



Amazon Web Services



Microsoft Azure



Google Cloud Platform

Continuer

En cliquant sur « Essayer », vous acceptez la [Politique de confidentialité](#) et les [Conditions de service](#).

Vous n'avez pas de compte cloud ?

Community Edition est un environnement Databricks limité destiné à un usage personnel et à la formation.

[Essayer Community Edition →](#)

En cliquant sur « Essayer Community Edition », vous acceptez la [Politique de confidentialité](#) et les [Conditions de service](#) undefined.