



Machine Learning

Bertrand Borel, Loïc Plessis, Pauline Sanchez, Victor Lannurien.



TABLE DES MATIÈRES

- Définition du ML
- Apprentissage Supervisé
- Dataset (features, target)
- Modèle (régression linéaire simple, polynomiale, multiple)
- Fonction de coût (l'erreur quadratique moyenne)
- Gradient en ML
- Algorithme de minimisation (Descente de gradient)
- Coefficient de détermination



Définition du machine learning :

- Science moderne permettant de découvrir des patterns et d'effectuer des prédictions à partir de données
- Basé sur des statistiques, du forage de données, de la reconnaissance de patterns et des analyses prédictives
- Très efficace dans les situations où les insights doivent être découvertes à partir de larges ensembles de données diverses et changeantes (Big Data)
- Branche de l'intelligence artificielle englobant de nombreuses méthodes. Ces méthodes sont des algorithmes
- Un système ML ne suit pas d'instructions, mais apprend à partir de l'expérience.
- Ses performances s'améliorent au fil de son "entraînement"



Apprentissage supervisé

- Modèle d'apprentissage le plus populaire en ML
- Consiste à superviser l'apprentissage de la machine en lui montrant des exemples (des données) de la tâche qu'elle doit réaliser
- Fonctionne en 4 étapes :
 - Importer un Dataset (x, y) qui contient les exemples
 - Développer un **Modèle** aux paramètres aléatoires
 - Développer une **Fonction Coût** qui mesure les erreurs entre le modèle et le Dataset
 - Développer un **Algorithme d'apprentissage** pour trouver les **paramètres** du modèle qui **minimisent** la **Fonction Coût**

Dataset (features, target)

Les Features

	A	B	C	D	E	F	G	H	I	J	K	L
1	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
2	1	0	3	Braund, Mr.	male	22	1	0	A/5 21171	7.25		S
3	2	1	1	Cumings, Mr.	female	38	1	0	PC 17599	71.2833	C85	C
4	3	1	3	Heikkinen, M.	female	26	0	0	STON/O2. 31	7.925		S
5	4	1	1	Futrelle, Mrs.	female	35	1	0	113803	53.1	C123	S
6	5	0	3	Allen, Mr. W.	male	35	0	0	373450	8.05		S
7	6	0	3	Moran, Mr. J.	male		0	0	330877	8.4583		Q



Target

Prenons un exemple : on cherche à déterminer si une personne, dans la base de données d'un magasin de chaussures, va cliquer sur le lien contenu dans le mail promotionnel qui lui a été envoyé. Pour cela, nous observons dans un premier temps l'historique, dont voici un extrait :

Premier email promotionnel	Âge	A cliqué
OUI	52	OUI
NON	19	NON
NON	37	NON
OUI	25	OUI



Modèle (régression linéaire simple, polynomiale, multiple)

Régression linéaire:

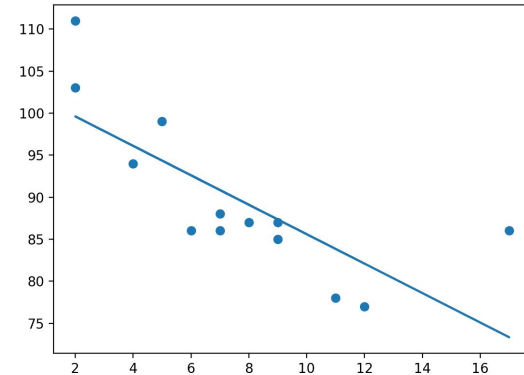
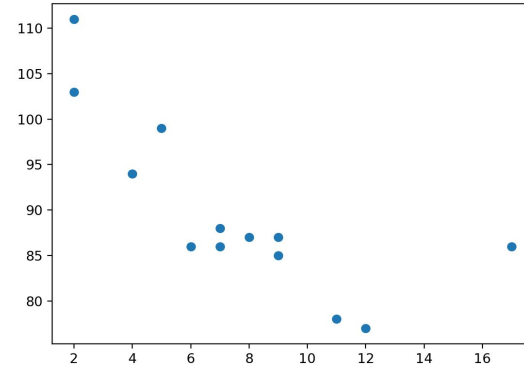
Il existe 3 types de régression linéaire:

- la régression linéaire simple
- la régression linéaire polynomiale
- la régression linéaire multiple

Modèle : régression linéaire simple

```
1 import matplotlib.pyplot as plt
2
3 x = [5,7,8,7,2,17,2,9,4,11,12,9,6]
4 y = [99,86,87,88,111,86,103,87,94,78,85,86]
5
6 plt.scatter(x, y)
7 plt.show()
```

```
1 import matplotlib.pyplot as plt
2 from scipy import stats
3
4 x = [5,7,8,7,2,17,2,9,4,11,12,9,6]
5 y = [99,86,87,88,111,86,103,87,94,78,85,86]
6
7 slope, intercept, r, p, std_err = stats.linregress(x, y)
8
9 def myfunc(x):
10     return slope * x + intercept
11
12 mymodel = list(map(myfunc, x))
13
14 plt.scatter(x, y)
15 plt.plot(x, mymodel)
16 plt.show()
```



Modèle : régression linéaire simple

```
1  from scipy import stats
2
3  x = [5,7,8,7,2,17,2,9,4,11,12,9,6]
4  y = [99,86,87,88,111,86,103,87,94,78,77,85,86]
5
6  slope, intercept, r, p, std_err = stats.linregress(x, y)
7
8  print(r)
```

-0.758591524376155

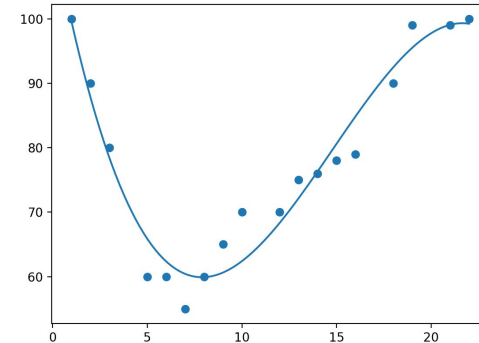
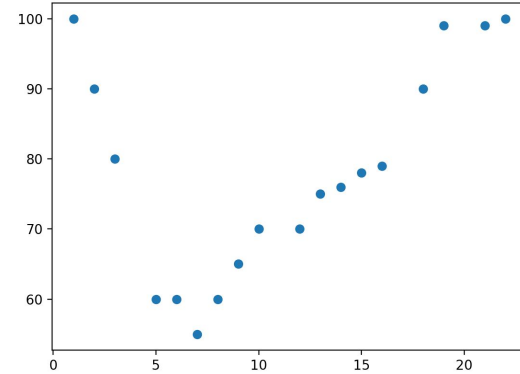
```
1  from scipy import stats
2
3  x = [5,7,8,7,2,17,2,9,4,11,12,9,6]
4  y = [99,86,87,88,111,86,103,87,94,78,77,85,86]
5
6  slope, intercept, r, p, std_err = stats.linregress(x, y)
7
8  def myfunc(x):
9      return slope * x + intercept
10
11  speed = myfunc(10)
12
13  print(speed)
```

85.59308314937454

Modèle : régression linéaire polynomiale

```
1 import matplotlib.pyplot as plt
2
3 x = [1,2,3,5,6,7,8,9,10,12,13,14,15,16,18,19,21,22]
4 y = [100,90,80,60,60,55,60,65,70,70,75,76,78,79,90,99,99,100]
5
6 plt.scatter(x, y)
7 plt.show()
```

```
1 import numpy
2 import matplotlib.pyplot as plt
3
4 x = [1,2,3,5,6,7,8,9,10,12,13,14,15,16,18,19,21,22]
5 y = [100,90,80,60,60,55,60,65,70,70,75,76,78,79,90,99,99,100]
6
7 mymodel = numpy.poly1d(numpy.polyfit(x, y, 3))
8
9 myline = numpy.linspace(1, 22, 100)
10
11 plt.scatter(x, y)
12 plt.plot(myline, mymodel(myline))
13 plt.show()
```



Modèle : régression linéaire polynomiale

```
1 import numpy
2 from sklearn.metrics import r2_score
3
4 x = [1,2,3,5,6,7,8,9,10,12,13,14,15,16,18,19,21,22]
5 y = [100,90,80,60,60,55,60,65,70,70,75,76,78,79,90,99,99,100]
6
7 mymodel = numpy.poly1d(numpy.polyfit(x, y, 3))
8
9 print(r2_score(y, mymodel(x)))
```

0.9432150416451027

```
1 import numpy
2 from sklearn.metrics import r2_score
3
4 x = [1,2,3,5,6,7,8,9,10,12,13,14,15,16,18,19,21,22]
5 y = [100,90,80,60,60,55,60,65,70,70,75,76,78,79,90,99,99,100]
6
7 mymodel = numpy.poly1d(numpy.polyfit(x, y, 3))
8
9 speed = mymodel(17)
10 print(speed)
```

88.87331269697987



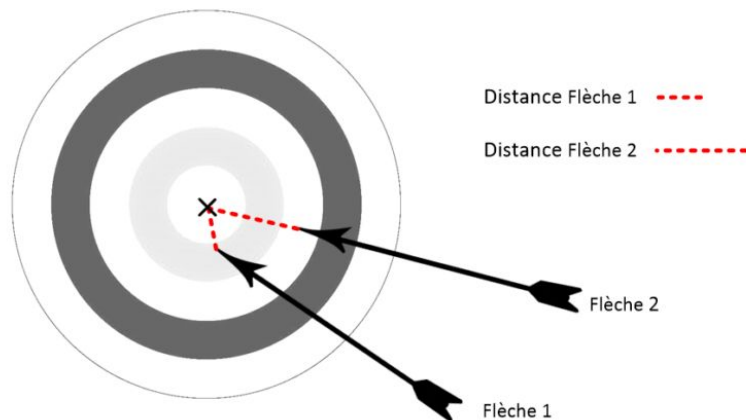
Modèle : régression linéaire multiple

Exactement la même chose que pour les deux autres régressions précédentes mais basée sur 2 variables ou plus.

Fonction de coût (l'erreur quadratique moyenne)

Pour savoir quel modèle est le meilleur parmi 2 candidats, il faut les évaluer. Pour cela, on mesure l'erreur entre un modèle et le Dataset, et on appelle ça la **Fonction Coût**.

Dans le cas d'une régression, on peut par exemple mesurer l'erreur entre la prédiction du modèle et la valeur qui est associée à ce dans notre Dataset. C'est similaire à l'idée de mesurer la distance entre votre flèche () et le centre de la cible, qui n'est autre que le point () qu'elle est sensée atteindre.



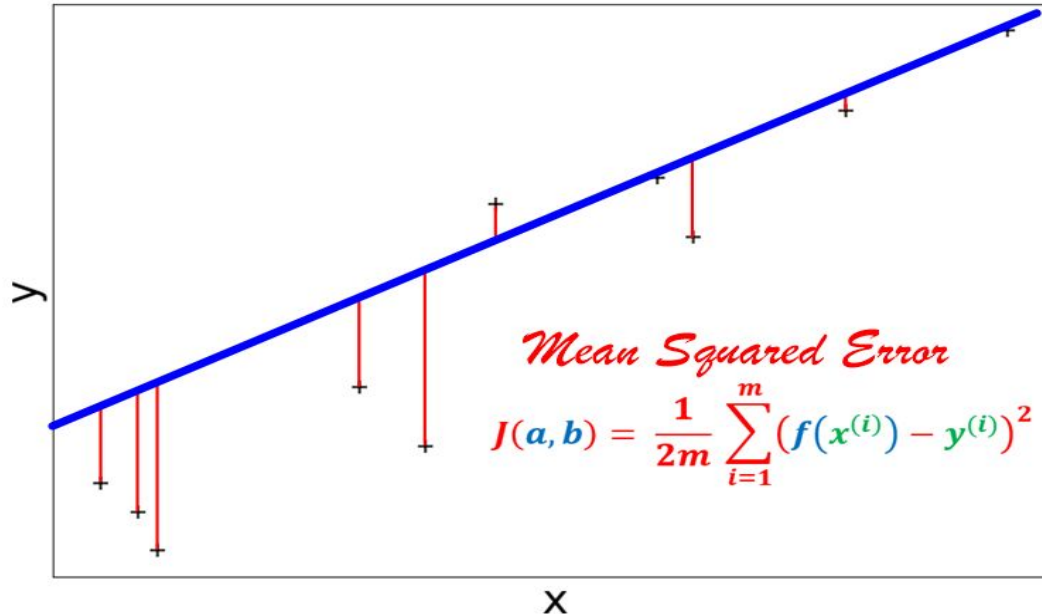
L'erreur quadratique moyenne



$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n \underbrace{(y_i - \hat{y}_i)^2}_{\substack{\text{predicted value} \quad \text{actual value}}}$$

test set

l'Erreur Quadratique Moyenne ou autrement appelé anglais la fonction **Mean Squared Error**.



Gradient en ML

Repose sur la notion de dérivée

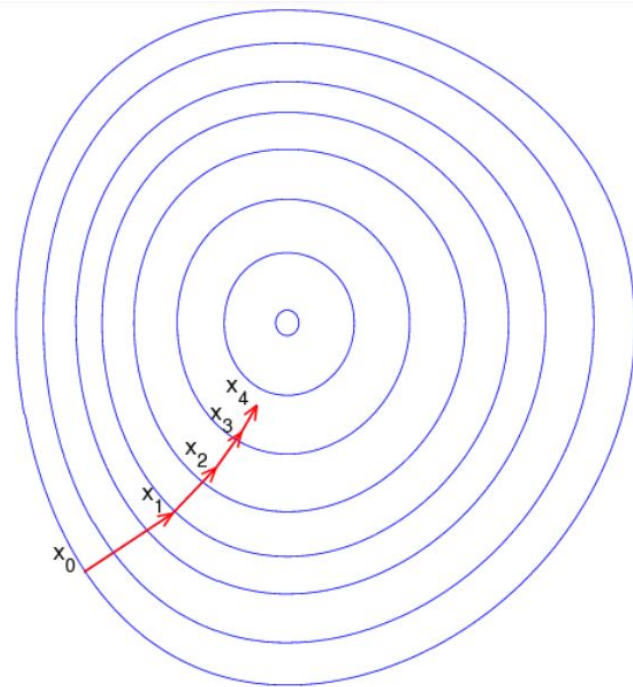
Un gradient est un dérivé d'une fonction qui

a plus d'une variable d'entrée.

En lien avec le calcul vectoriel.

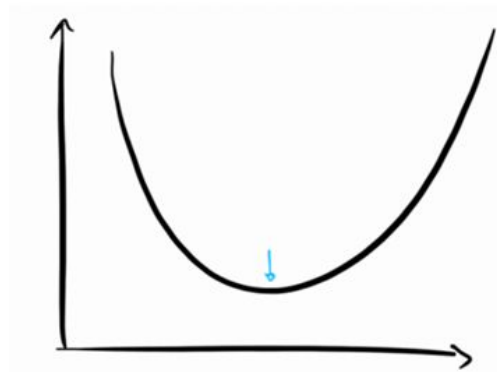
Permet de calculer la pente locale de la fonction

= progression “pas-à-pas”

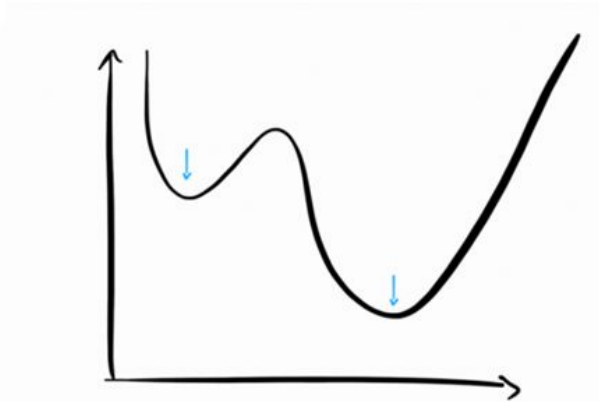


Algorithme de minimisation (Descente de gradient)

Algorithme d'optimisation permettant de trouver le minimum d'une fonction convexe.



Fonction Convexe



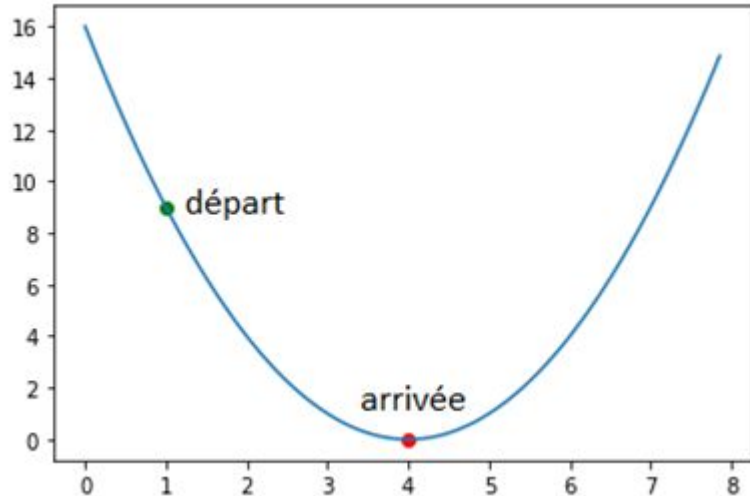
Fonction Non-Convexe

Fonction convexe = un minimum global

Fonction non-convexe = plusieurs minimums locaux

Sert pour la minimisation de la fonction coût : la machine apprend et trouve le meilleur modèle

L'algorithme permet de trouver la valeur idéale entre le départ et l'arrivée.



Calcul d'une dérivée (pour la courbe)

valeur négative en pente / valeur
positive en montée

Importance de l'hyper-paramètre
Alpha



Coefficient de détermination (linéaire de Pearson)

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

- Noté R^2 , c'est une mesure de la qualité de la prédiction d'une régression linéaire
- C'est le carré de la corrélation de Pearson entre les vraies valeurs et les valeurs prédites
- Permet d'indiquer à quel point les valeurs prédites sont corrélées aux vraies valeurs



Merci de votre attention

