

Projet Intégrateur CMI - 2021/2022

Étude de modèles et construction d'une
interface R-shiny sur des données COVID-19.

Pierre LAGUE, Diamondra RAKOTONDRAZAKA, Jean ROBIN, Damien TANNEAU, Titouan QUINTIN, Victor HAGUET, Ewen DENIAU

8 juin 2022



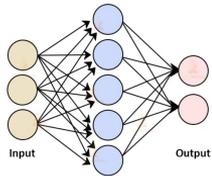
Sommaire

- Rédaction de l'état de l'art
- Préparation des données
- Développement des modèles
- Conception et implémentation de l'application web R Shiny
- Conclusion sur le projet

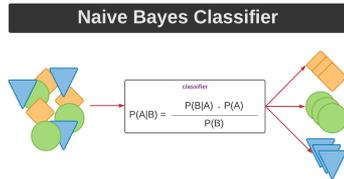
Rédaction de l'état de l'art

- Constitue une base de travail
- Permet une connaissance plus précise du modèle étudié/implémenté
- Mise à jour des pratiques et des standards du monde professionnel

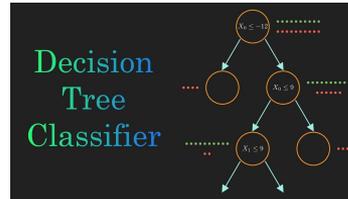
BUT -> Travailler sur des modèles/algorithmes standards et se les approprier.



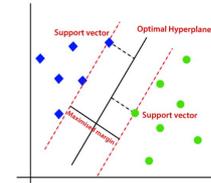
Réseau de neurones



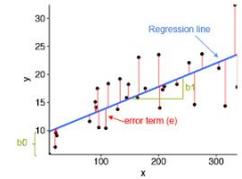
Classifieur Naïf Bayes



Arbre de décision



Support Vecteur Machine



Régression Linéaire et ses dérivées



Régression pénalisée

Pénalités

Ridge $\lambda \times (variable_1^2 + \dots + variable_x^2)$

Lasso $\lambda \times (|variable_1| + \dots + |variable_x|)$

Elastic-Net $\lambda_1 \times (variable_1^2 + \dots + variable_x^2) + \lambda_2 \times (|variable_1| + \dots + |variable_x|)$

Préparation des données - Description



Données
Covid-19

COLUMN COUNT 67 <small>2022-06-06 09:48</small>	FILE COUNT 1 <small>2022-06-06 09:48</small>	SIZE 45.28 MB <small>2022-06-06 09:48</small>
RECORD COUNT 166 554 <small>2022-06-06 09:48</small>		

Description du dataset

- Données manquantes
- Colonnes inutiles

iso_code	continent	location	date	total_cases	new_cases	new_cases_smoothed	total_deaths
missing	missing	missing	missing	missing	missing	missing	missing
Country	Text	Country	Date (unparsed)	Decimal	Decimal	Decimal	Decimal
AFG	Asia	Afghanistan	2020-02-24	5.0	5.0		
AFG	Asia	Afghanistan	2020-02-25	5.0	0.0		
AFG	Asia	Afghanistan	2020-02-26	5.0	0.0		
AFG	Asia	Afghanistan	2020-02-27	5.0	0.0		
AFG	Asia	Afghanistan	2020-02-28	5.0	0.0		
AFG	Asia	Afghanistan	2020-02-29	5.0	0.0		
AFG	Asia	Afghanistan	2020-03-01	5.0	0.0		0.714
AFG	Asia	Afghanistan	2020-03-02	5.0	0.0		0.0
AFG	Asia	Afghanistan	2020-03-03	5.0	0.0		0.0
AFG	Asia	Afghanistan	2020-03-04	5.0	0.0		0.0
AFG	Asia	Afghanistan	2020-03-05	5.0	0.0		0.0
AFG	Asia	Afghanistan	2020-03-06	5.0	0.0		0.0
AFG	Asia	Afghanistan	2020-03-07	8.0	3.0		0.429
AFG	Asia	Afghanistan	2020-03-08	8.0	0.0		0.429
AFG	Asia	Afghanistan	2020-03-09	8.0	0.0		0.429
AFG	Asia	Afghanistan	2020-03-10	8.0	0.0		0.429
AFG	Asia	Afghanistan	2020-03-11	11.0	3.0		0.857
AFG	Asia	Afghanistan	2020-03-12	11.0	0.0		0.857
AFG	Asia	Afghanistan	2020-03-13	11.0	0.0		0.857
AFG	Asia	Afghanistan	2020-03-14	14.0	3.0		0.857
AFG	Asia	Afghanistan	2020-03-15	20.0	6.0		1.714
AFG	Asia	Afghanistan	2020-03-16	25.0	5.0		2.429
AFG	Asia	Afghanistan	2020-03-17	26.0	1.0		2.571
AFG	Asia	Afghanistan	2020-03-18	26.0	0.0		2.143
AFG	Asia	Afghanistan	2020-03-19	26.0	0.0		2.143

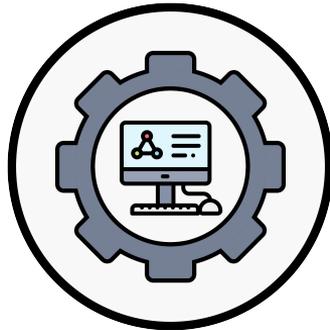
Échantillon du dataset

Préparation des données - Engineering



Données
Covid-19

$$x_{\text{norm}} = \frac{x - \min(x)}{\max(x) - \min(x)}$$

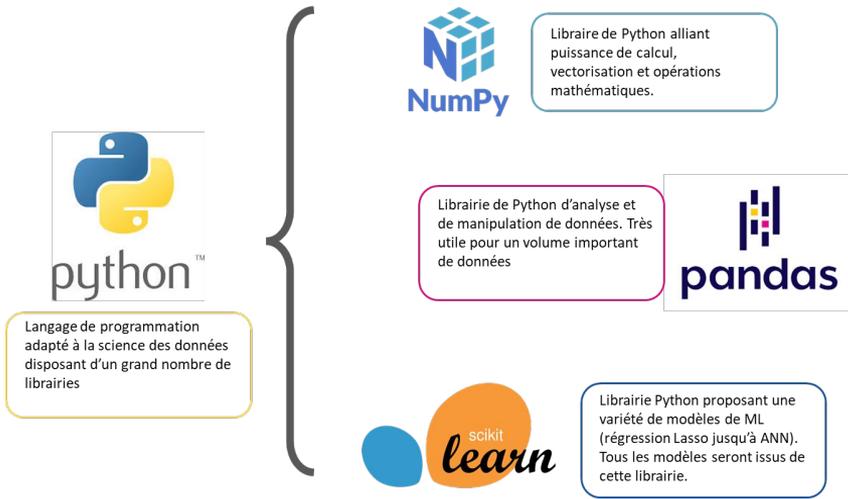


COLUMN COUNT	FILE COUNT	SIZE	RECORD COUNT
28	1	8.79 MB	134 996

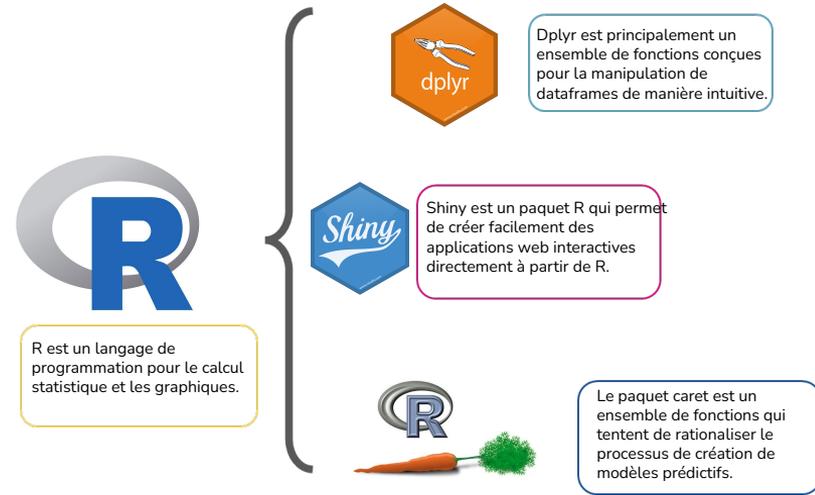
Détails du dataset final

- Conversion des données en types numériques
- Remplissage des valeurs manquantes (si NA < 5% insérer la moyenne des données).
- Parsing de la date

Développement des modèles - Implémentation



Tous les modèles + comparaison



Application R-Shiny et quelques modèles



Développement des modèles - Entraînements

```
k-Nearest Neighbors
101247 samples
  7 predictor

Pre-processing: centered (7), scaled (7)
Resampling: Bootstrapped (25 reps)
Summary of sample sizes: 101247, 101247, 101247, 101247, 101247, 101247, ...
Resampling results across tuning parameters:

 k  RMSE      Rsquared  MAE
 5  0.01461566  0.9436960  0.00229603
 7  0.01530359  0.9383102  0.00261205
 9  0.01596484  0.9328988  0.00288967

RMSE was used to select the optimal model using the smallest value.
The final value used for the model was k = 5.
```

Résultats de l'optimisation de l'algorithme KNN sur les données du Covid-19



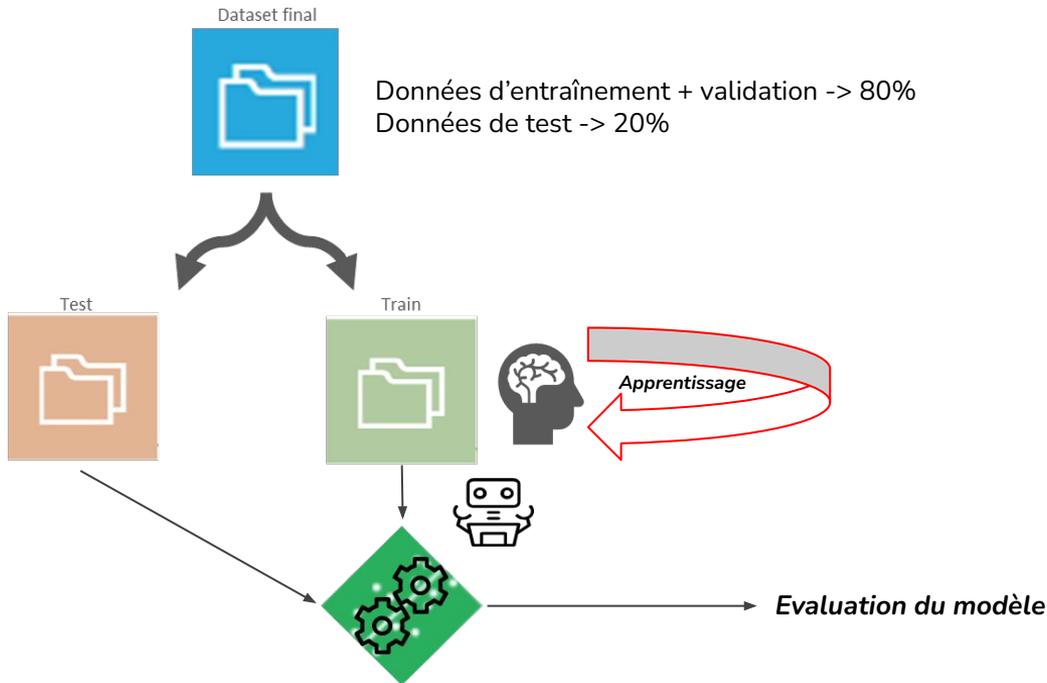
Développement des modèles - Entraînements

- Chaque modèle est entraîné avec ses hyperparamètres par défaut
- Chaque modèle est entraîné sur la même base de données train (même seed)
- Chaque modèle est entraîné sur “total_deaths” puis sur “total_cases”

<input type="checkbox"/>	TOTAL_DEATHS		
<input type="checkbox"/>	Random forest (total_deaths)	🏆 0.999	☆
<input type="checkbox"/>	Ordinary Least Squares (total_deaths)	0.869	☆
<input type="checkbox"/>	Ridge (L2) regression (total_deaths)	0.822	☆
<input type="checkbox"/>	Lasso (L1) regression (total_deaths)	0.823	☆
<input type="checkbox"/>	Decision Tree (total_deaths)	0.997	☆
<input type="checkbox"/>	SVM (total_deaths)	-0.074	☆
<input type="checkbox"/>	K Nearest Neighbors (grid) (total_deaths)	0.996	☆
<input type="checkbox"/>	Artificial Neural Network (total_deaths)	0.319	☆

Exemple de file d'entraînement sous Dataluku - R^2

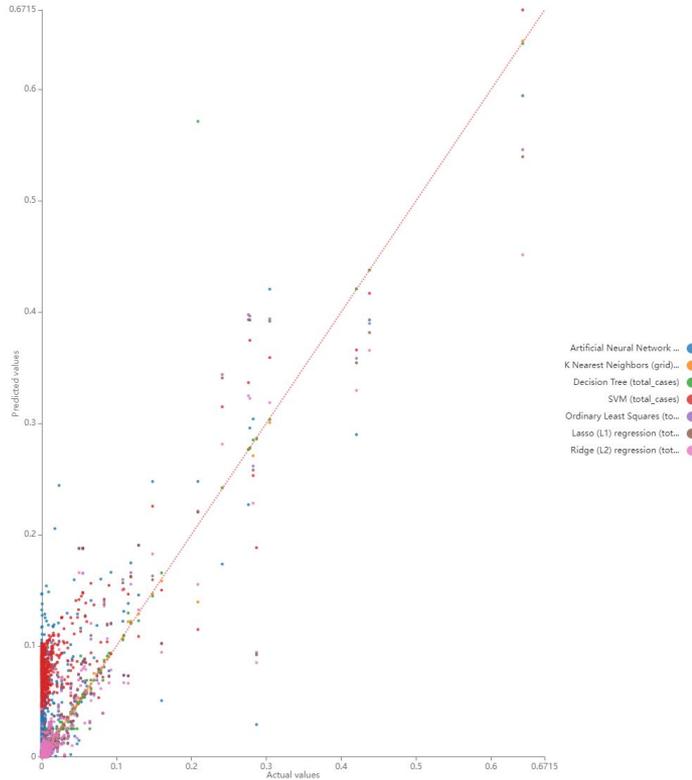
Développement des modèles - Evaluation



Evaluation du modèle :

- R^2
 - MSE
 - RMSE
 - Correlation
- + Allure de la courbe
“*actual vs predicted*”

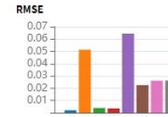
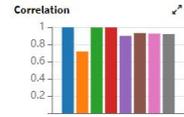
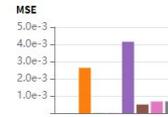
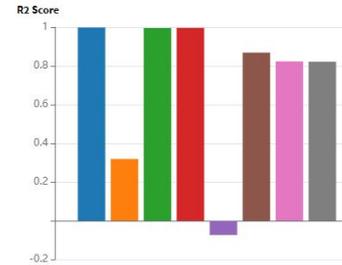
Développement des modèles - Evaluation



Scatter Plot des test des modèles

Prediction type: Regression

Metric: MSE, RMSE, R2 Score, Co

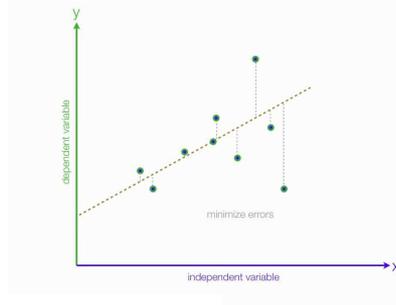


Model Evaluation		evaluation	evaluationDataset	model	metrics			
Name	date	dataset-name	name	R2 Score	MSE	RMSE	Correlation	
Random forest (total_deaths)	2022-06-04T09:...	train	Random forest (...)	0.999	0.000	0.002	0.999	
Artificial Neural Network (total_deaths)	2022-06-04T08:...	train	Artificial Neural ...	0.319	0.003	0.051	0.718	
K Nearest Neighbors (grid) (total_deaths)	2022-06-04T08:...	train	K Nearest Neigh...	0.996	0.000	0.004	0.998	
Decision Tree (total_deaths)	2022-06-04T08:...	train	Decision Tree (t...	0.997	0.000	0.003	0.998	
SVM (total_deaths)	2022-06-04T08:...	train	SVM (total_deaths)	-0.074	0.004	0.064	0.900	
Ordinary Least Squares (total_deaths)	2022-06-04T08:...	train	Ordinary Least ...	0.869	0.001	0.023	0.933	
Lasso (L1) regression (total_deaths)	2022-06-04T08:...	train	Lasso (L1) regre...	0.823	0.001	0.026	0.926	
Ridge (L2) regression (total_deaths)	2022-06-04T08:...	train	Ridge (L2) regre...	0.822	0.001	0.026	0.921	

Mesures MSE, RMSE, R² et Correlation des chaque modèles

Développement des modèles - Conclusion

Phase d'analyse



Les modèles linéaires sont plus adaptés à notre problème de régression



Il nous faut beaucoup de données préparées et sans valeurs manquantes



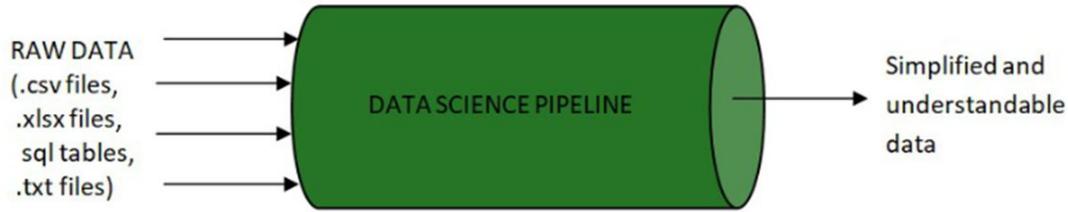
Les variables explicatives doivent être choisies de manière à expliquer la variable cible sur plusieurs dimensions (temps, espace)



Conception et implémentation de l'application web R Shiny

Shiny

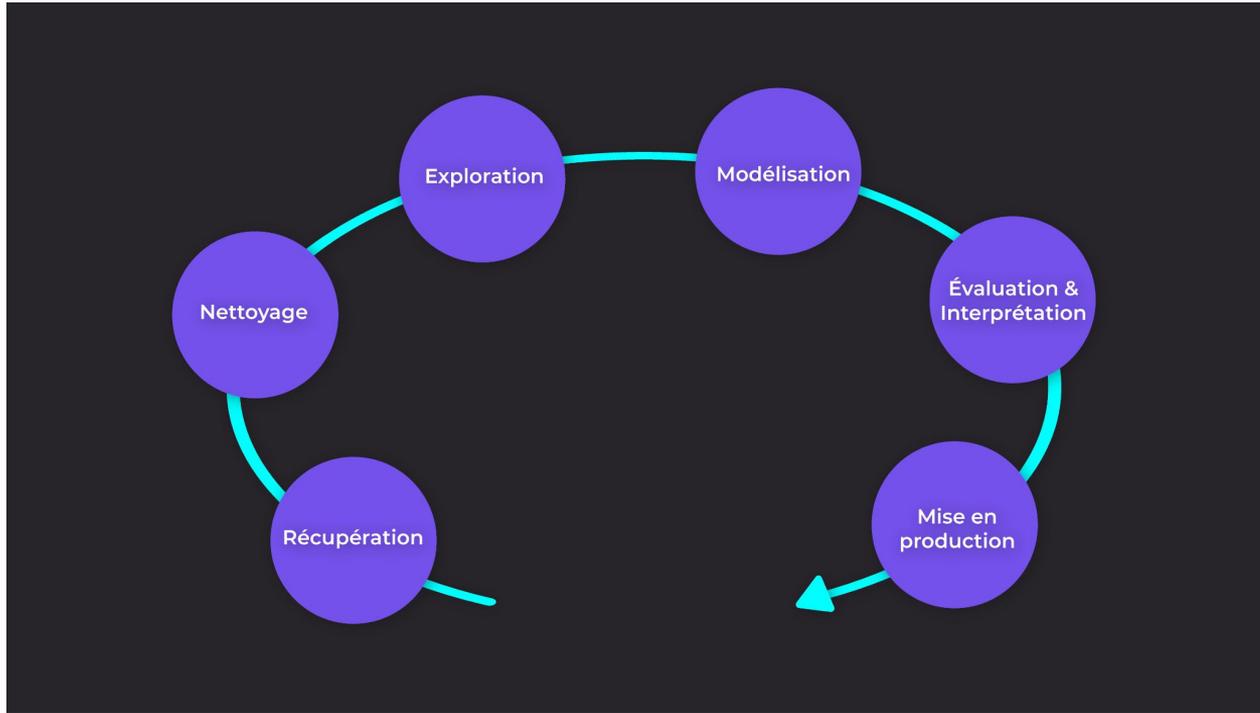
from  Studio



Représentation simplifiée d'un data pipeline

Source : <https://www.geeksforgeeks.org/whats-data-science-pipeline/>

- Obtention des données
- Nettoyage des données
- Analyse exploratoire des données
- Modélisation
- Interprétation



Etapes d'un projet de machine learning

Source : <https://openclassrooms.com/fr>

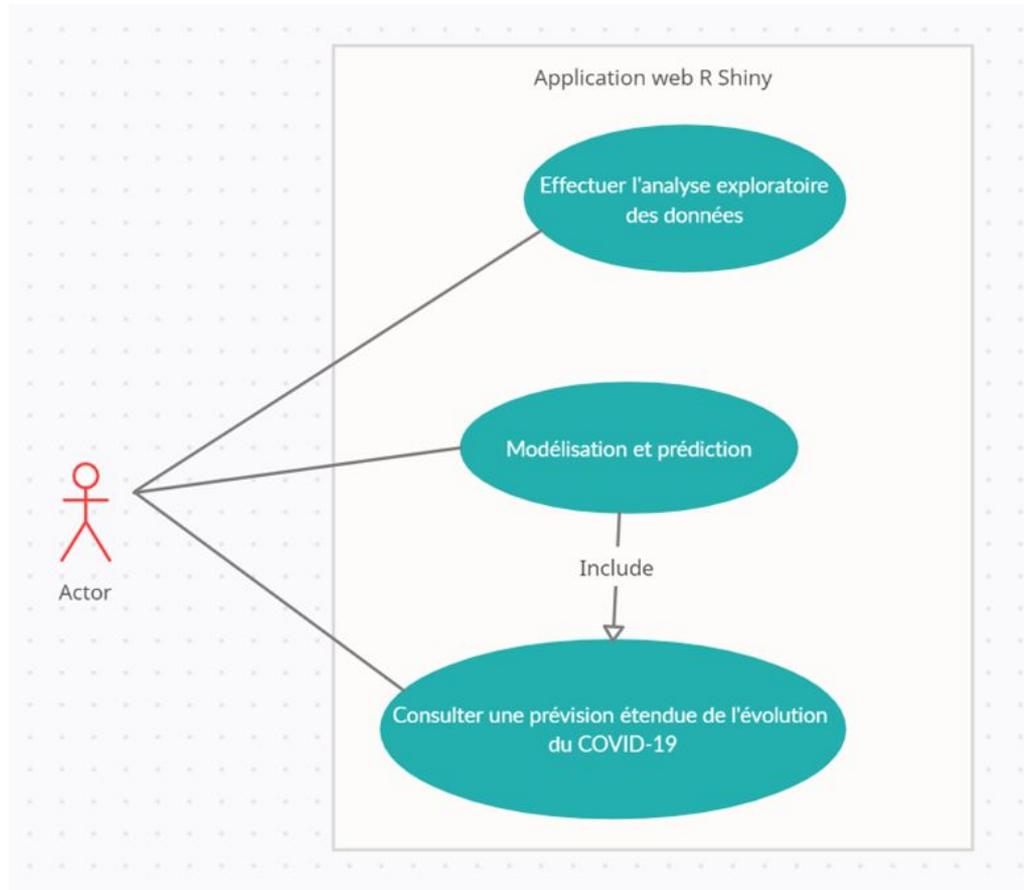
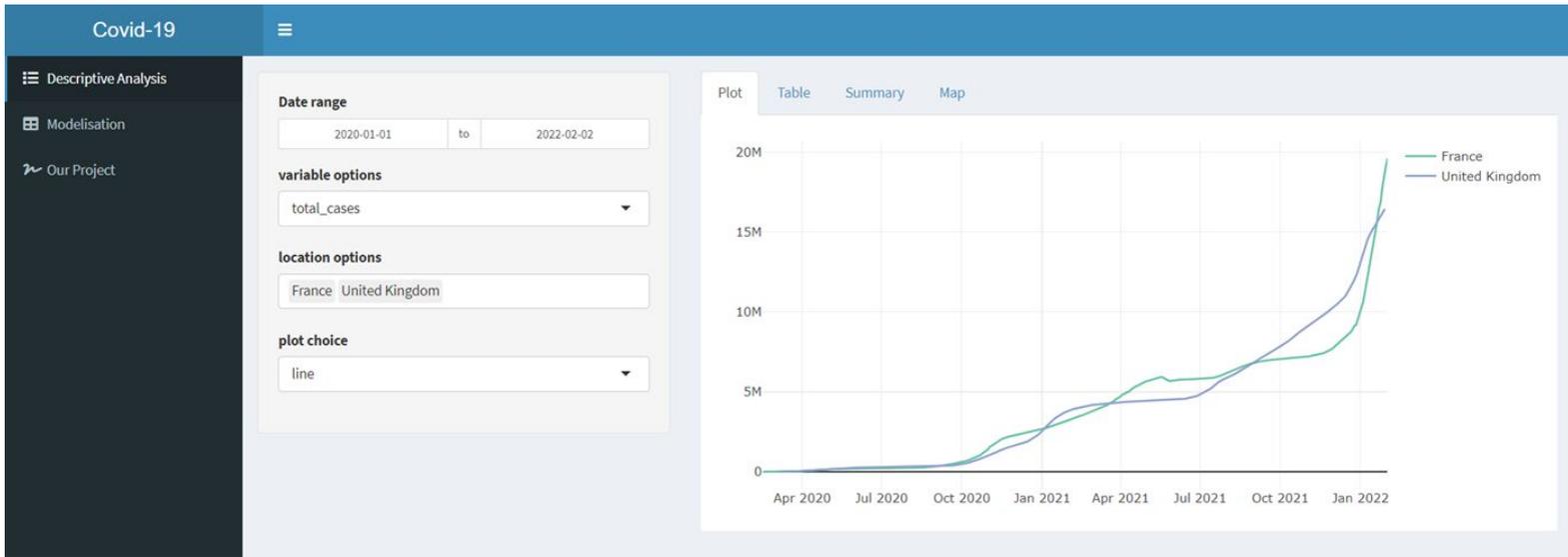


Diagramme de cas d'utilisation de l'application



Application Shiny

<https://q1ov6nn1.shinyapps.io/app-covid19/>



Covid-19

- Descriptive Analysis
- Modelisation
- Our Project

Date range: 2020-01-01 to 2022-02-02

variable options: total_cases

location options: France, United Kingdom

plot choice: line

Plot | Table | Summary | Map

location	mean	sd	min	median	max
France	4214134.32	3750394.54	12.00	3900344.00	19609576.00
United Kingdom	4097492.67	4047010.63	373.00	3976063.00	16447073.00



Covid-19

- Descriptive Analysis
- Modelisation
- Our Project

Date range
2020-01-13 to 2022-02-02

variable options
total_cases

location options
France United Kingdom

plot choice
line

Plot Table Summary **Map**

Country	Approximate Case Count
United States	15
France	10
United Kingdom	10
Germany	10
Italy	10
Spain	10
China	10
India	10
Japan	10
South Korea	10
Other countries	5

Modeling and predicting the evolution of COVID-19

target variable

total_cases

model

LR

location

France

population

67422000

people_vaccinated

1240

total_tests

34791440

icu_patients

2694

reproduction_rate

1,05

total_deaths_per_million

936,208

date

2021-01-01

2021

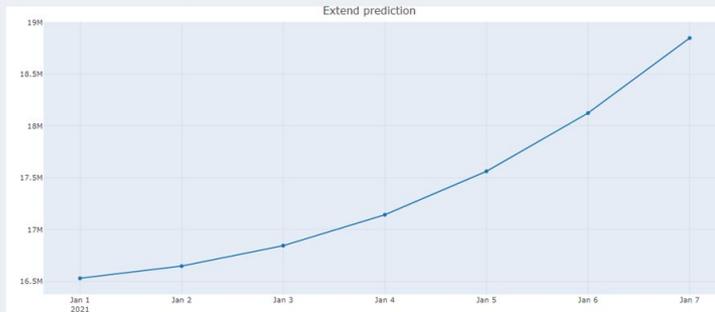


THE PREDICTED VALUE IS:

1683046

Extend prediction of total_cases

1 week



Fit of the model



LR RMSE : 2867498.94512285



Merci de votre attention