

Prédiction du Niveau de la Pollution Atmosphérique des Villes

Antoine Ghidini / Laetitia Kamwag / Sebila Doubaeva



99% DE POPULATION

Respirent un air qui dépasse les limites fixées



IMPACT SUR LA SANTE

Cardiopathies, cancers du poumon, de système cardi-vasculaire



7 MILLIONS

Meurent chaque année suite à la pollution de l'air

Quelles sont nos motivations?

Amener le projet au stade de l'utilité pour tous : non seulement apprendre, mais aussi partager nos résultats !



ALGORITHMES CONCRETS



DONNEE REELLES



APPLICATION RESULTANTE

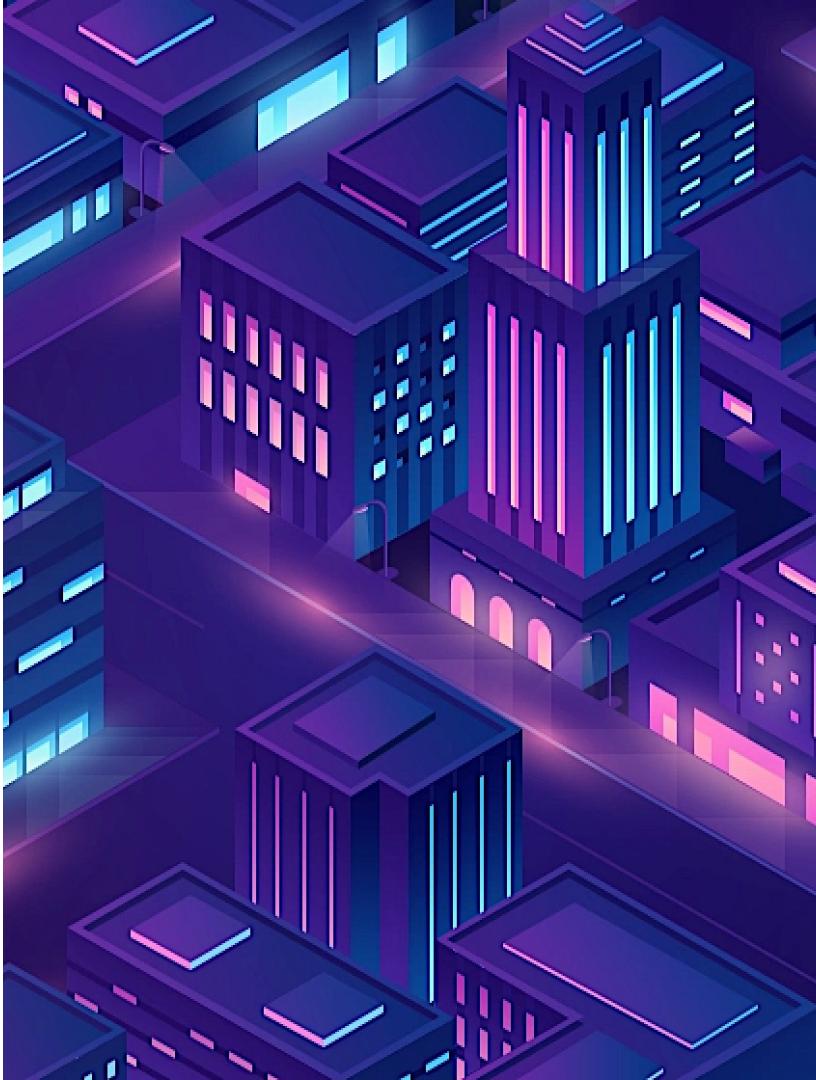




Machine Learning

- Numpy
- Matplotlib, Seaborn
- Pandas, Skit Learn

- Jupyter Notebook
- GitHub
- Dash & Heroku



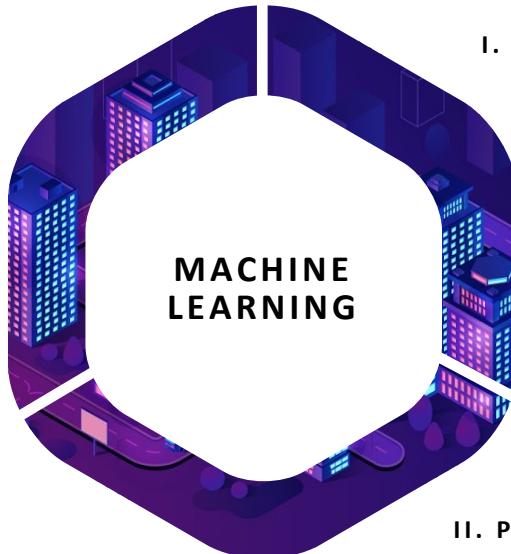


Méthodologie

Les étapes à suivre dans
un cycle de projet en ML

III. LEARNING

- Decision Tree Regressor
- Random Forest Regressor
- Linear Regression



I. EXPLORATORY DATA ANALYSIS

- Comprendre les données
- Identifier les erreurs évidentes
- Détecter les anomalies
- Trouver des relations

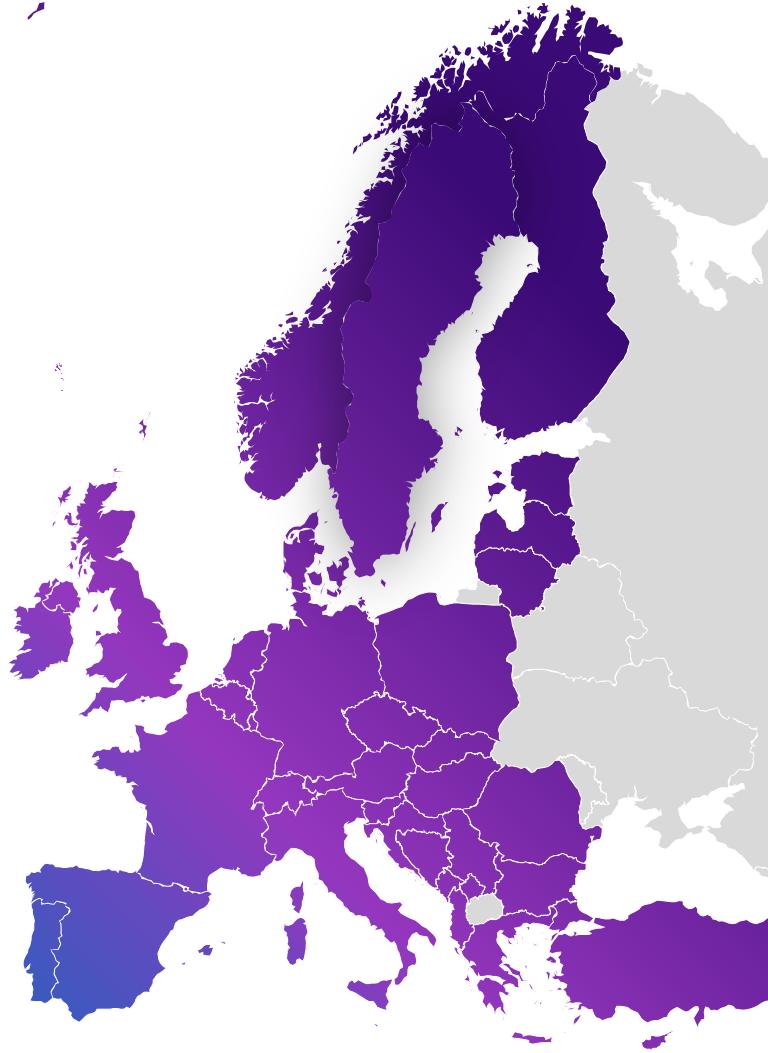
II. PREPROCESSING

- Amputation de données
- Validation et nettoyage de données brutes
- Encodage & Fenêtre glissante



European Environment Agency

- Convention sur la pollution atmosphérique à longue distance (LRTAP Convention)
- Inventaire de la Directive sur les Plafond d'Emission Nationaux (NEC)
- Portail Européen des Emissions Industrielles : Directive sur les Emissions Industrielles et du Registre Européen des Rejets et Transferts de Polluants (E-PRTR)





DataSets

- 2 732 184 (72%) données d'émissions manquantes pour dataset de LRTAPC
- 95 données d'émissions manquantes pour dataset de E-PRTR
- 224 villes manquantes
- 19 nom d'installations

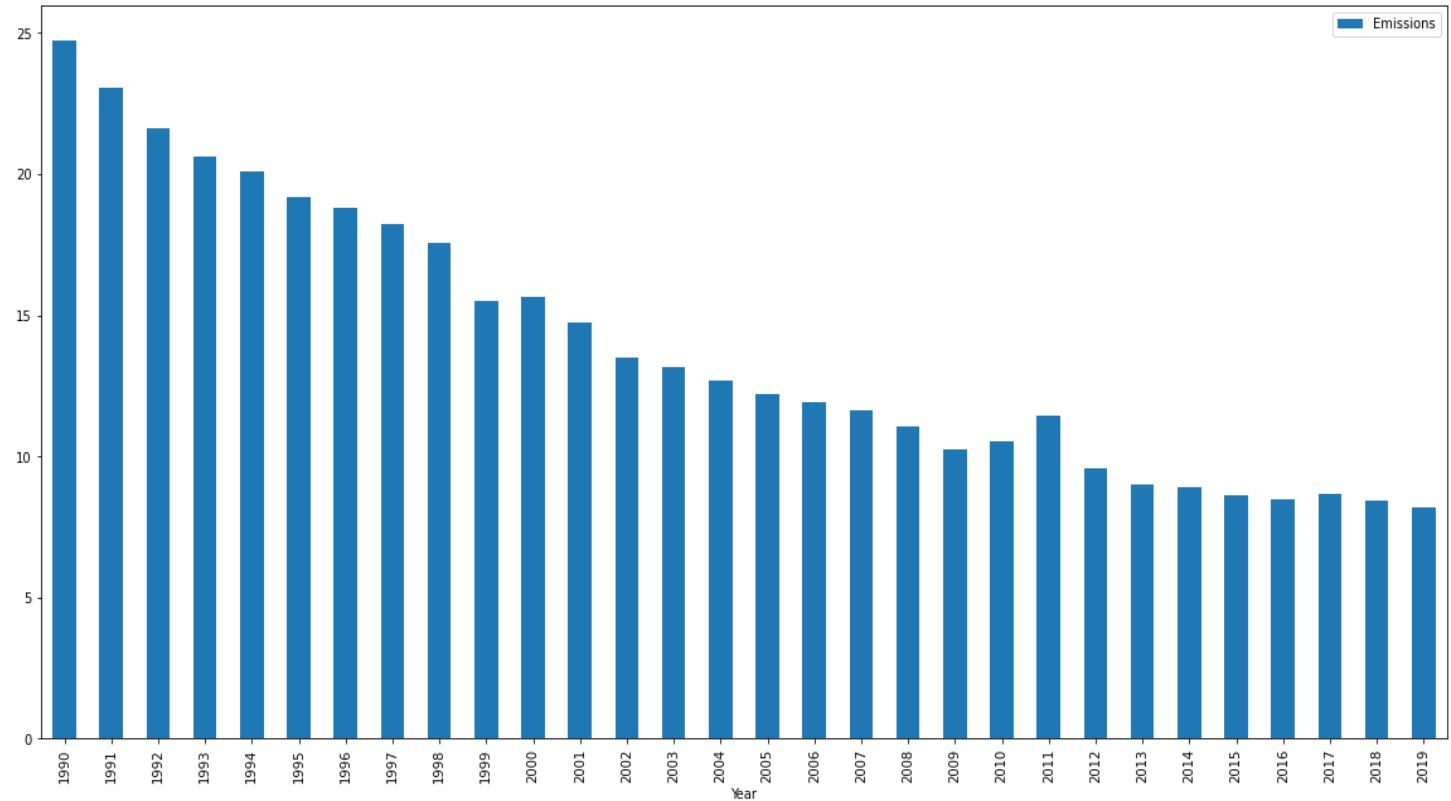
RangeIndex : 3794700 entries		
#	Column	Dtype
0	Country_Code	object
1	Country	object
2	Pollutant_name	object
3	Format_name	object
4	Sector_code	object
5	Year	int64
6	Emissions	float64
7	Unit	object
8	Notation	object
9	VersionId	int64
10	Parent_sector_code	object
11	Sector_name	object

RangeIndex : 288566 entries		
#	Column	Dtype
0	countryName	object
1	EPRTRSectorCode	object
2	eprtrSectorName	object
3	EPRTRAnnexIMainActivityCode	object
4	EPRTRAnnexIMainActivityLabel	object
5	FacilityInspireID	object
6	facilityName	object
7	Longitude	float64
8	Latitude	float64
9	City	object
10	targetRelease	object
11	pollutant	object
12	emissions	float64
13	reportingYear	int64

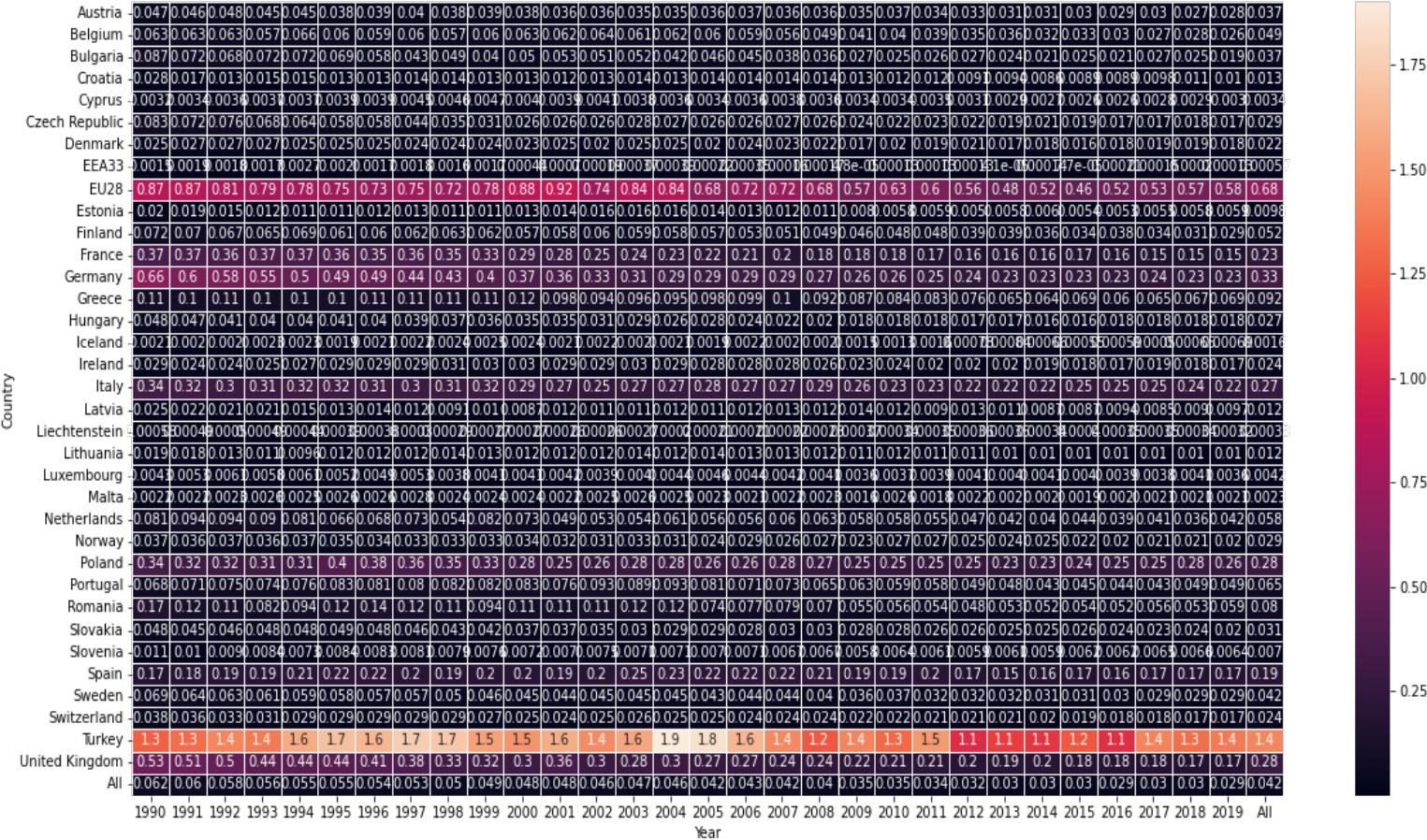


EDA - Premier Dataset

Graphique d'évolution du niveau des émission de dataset de LRTAP Convention

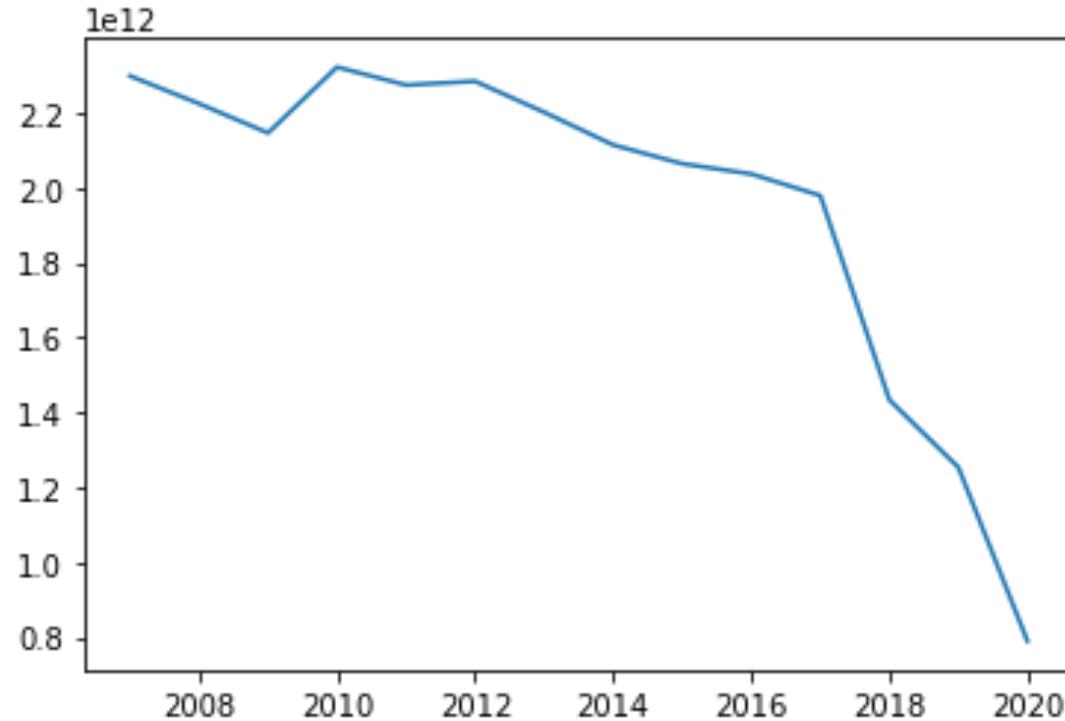


Heatmap avec Seaborn de dataset de LRTAP Convention

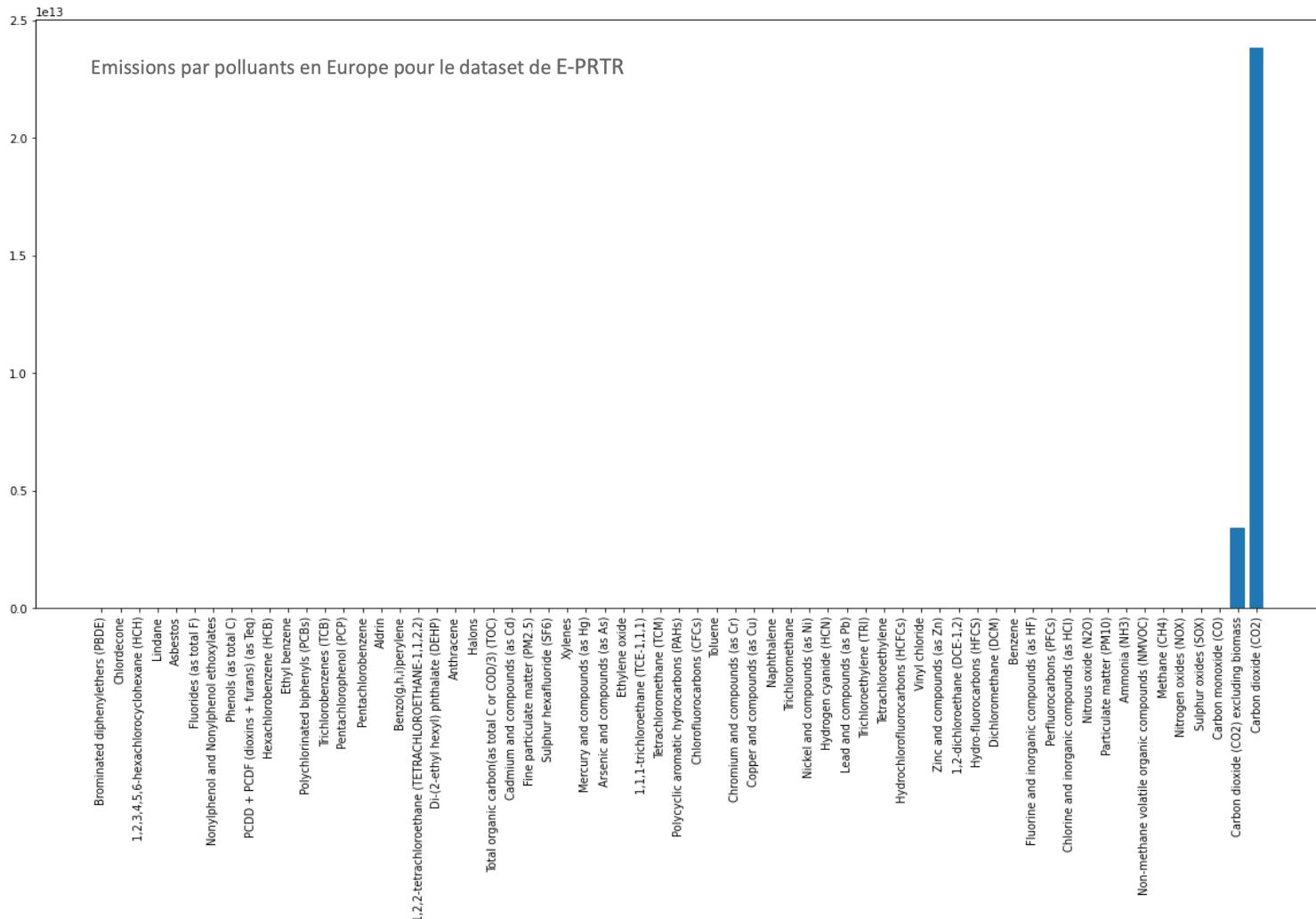




EDA - Second Dataset

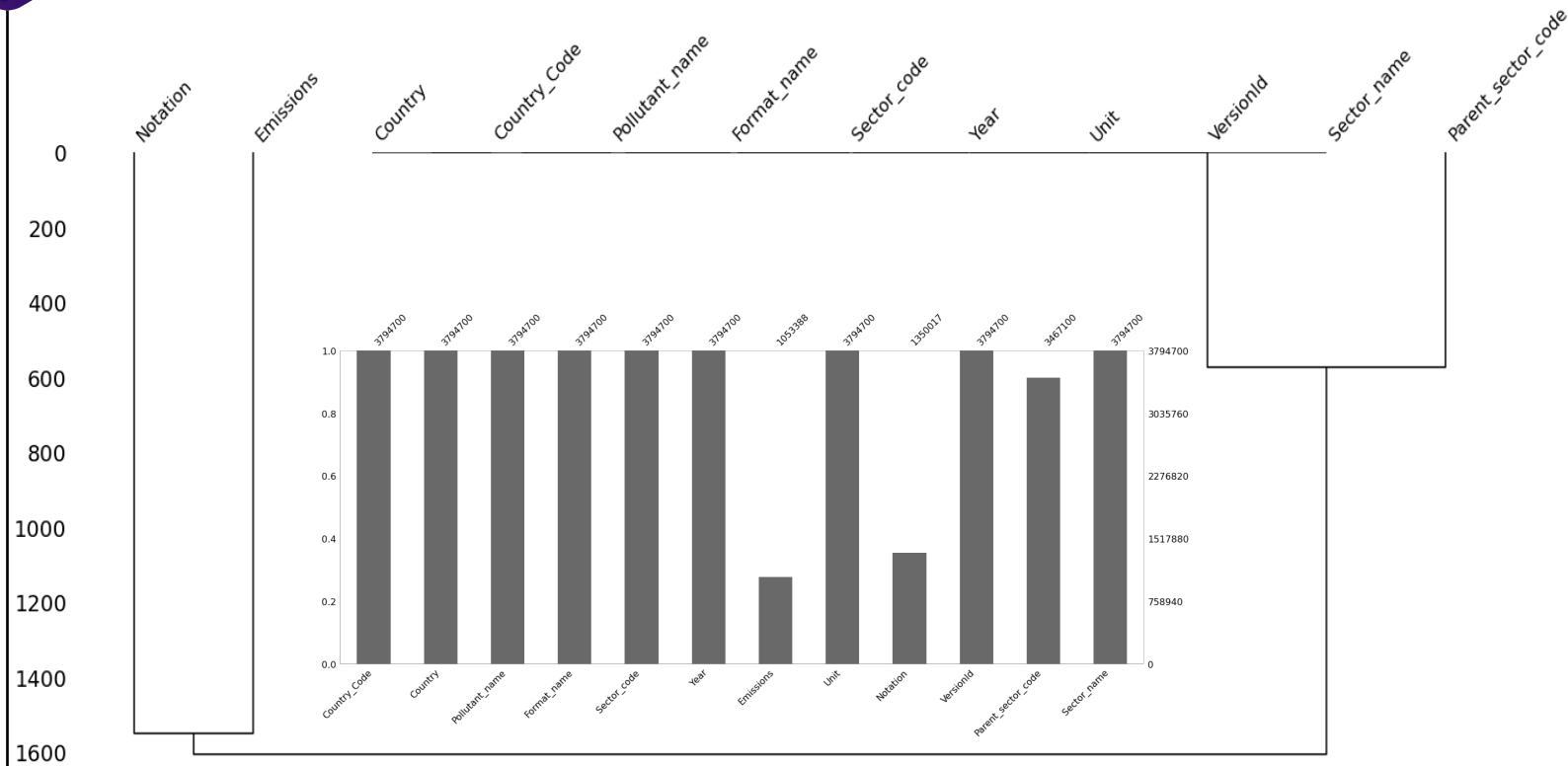


Evolution des émissions en Europe entre 2007 et 2020 pour le dataset de E-PRTR





Data Amputation



Le dendrogramme du regroupement hiérarchique des colonnes ainsi que barplot de nullité par colonne



Data Cleaning I



- Energy production and distribution
- Energy use in industry
- Non road transport
- Commercial, institutional. and households
- Industrial processes and product use
- National total for the entire territory
- Agriculture
- Waste
- Other

- CO
- NH₃
- NMVOC
- Nox
- PM10
- PM2.5
- Sox
- TSP



Data Cleaning II

RangeIndex : 288566 entries		
#	Column	Dtype
0	countryName	object
1	EPRTRAnnexIMainActivityCode	object
2	Longitude	float64
3	Latitude	float64
4	pollutant	object
5	emissions	float64
6	reportingYear	int64

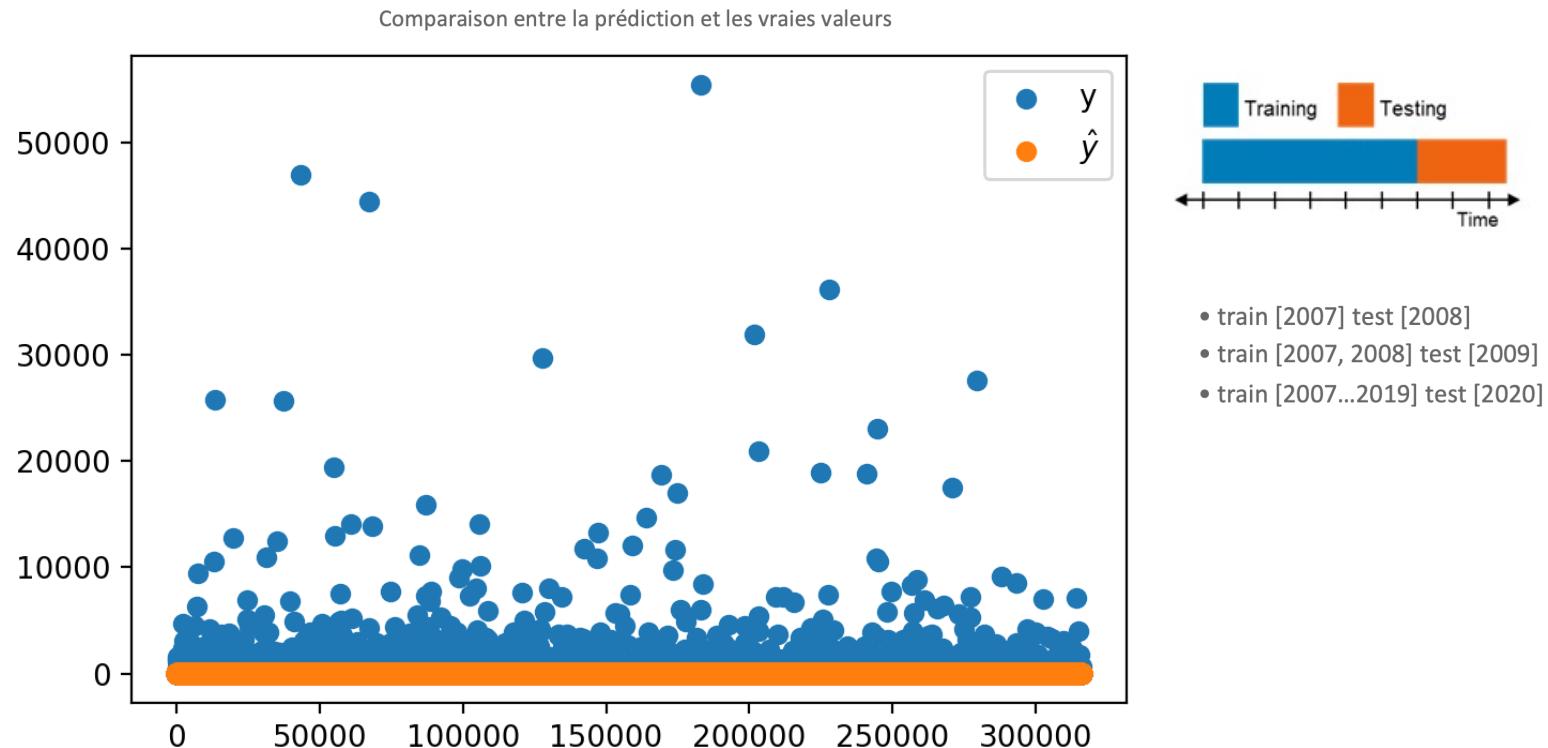


RangeIndex : 288566 entries		
#	Column	Dtype
0	countryName	object
1	EPRTRSectorCode	object
2	eprtrSectorName	object
3	EPRTRAnnexIMainActivityCode	object
4	EPRTRAnnexIMainActivityLabel	object
5	FacilityInspireID	object
6	facilityName	object
7	Longitude	float64
8	Latitude	float64
9	City	object
10	targetRelease	object
11	pollutant	object
12	emissions	float64
13	reportingYear	int64

- Suppression des colonnes
- Ajout des villes



Fenêtre Glissante





Encodage

country
Austria
Belgium
Bulgaria
Croatia



country
0
1
2
3

- Encodage avec OrdinalEncoder() de Sklearn
- Encodage avec One Hot Coder

color
red
red
yellow
green
yellow



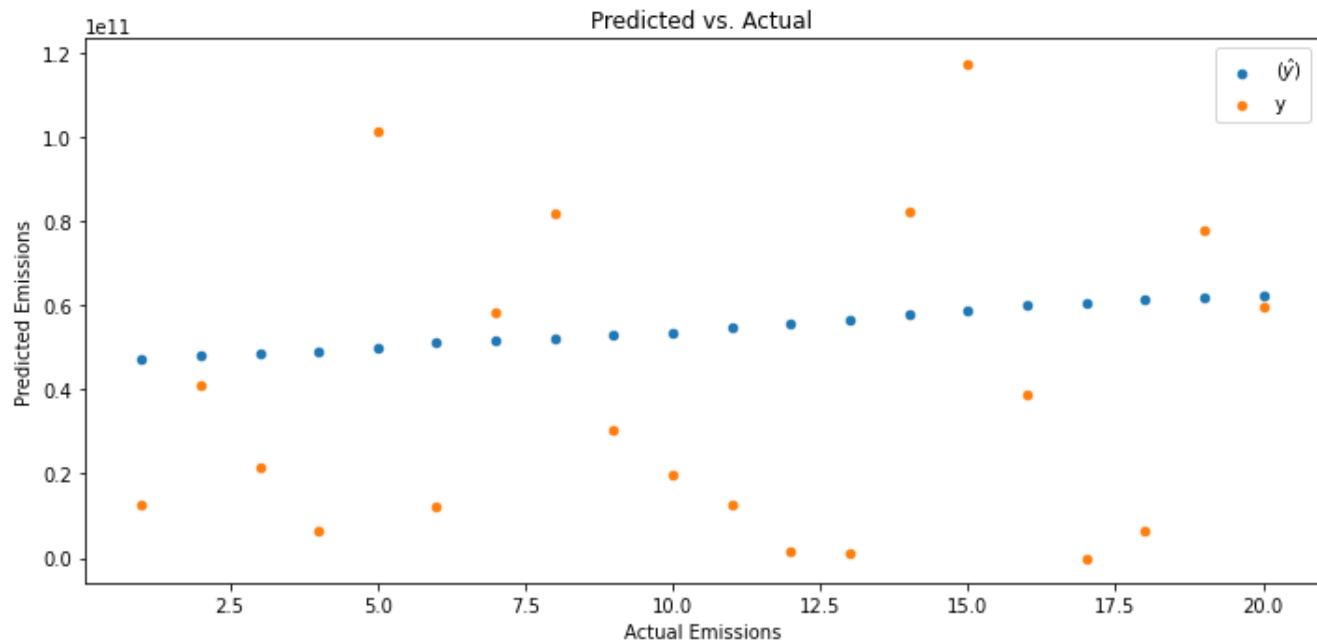
red	yellow	green
1	0	0
1	0	0
0	1	0
0	0	1
0	1	0



Linear Regression

- $R^2 = -0.174$
- RMSE = 385972242

Prediction des émissions pour 2020 pour le CO2

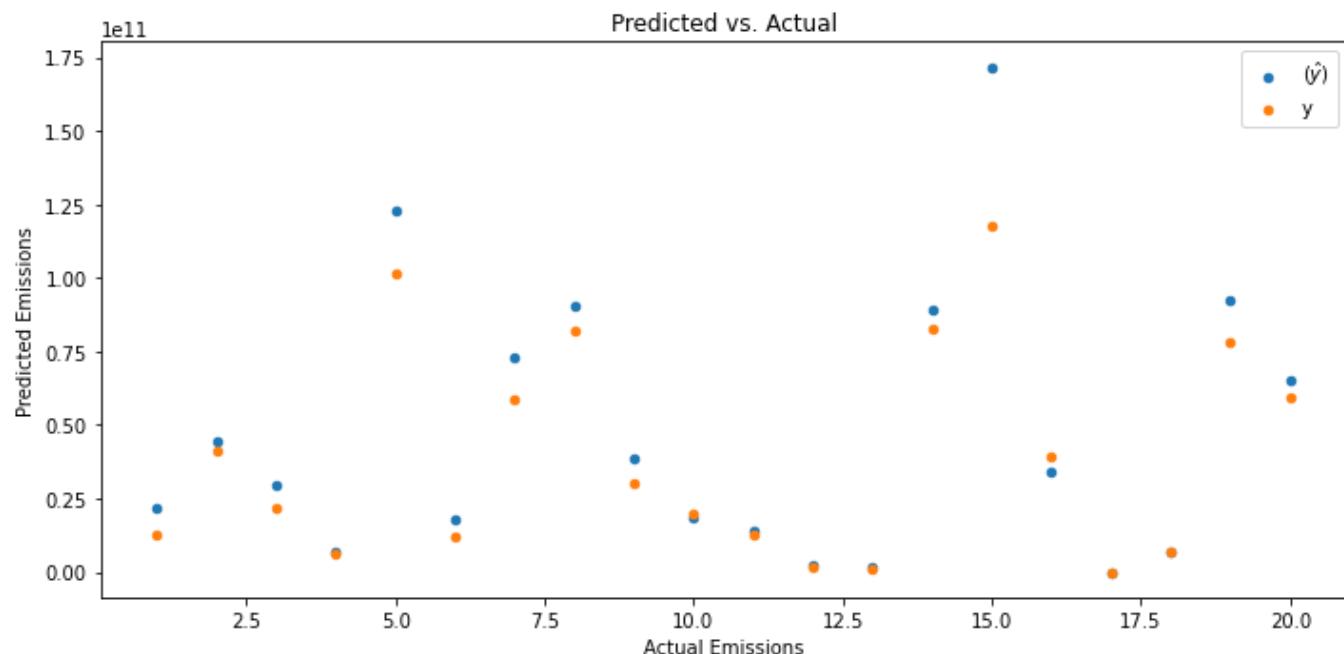




Decision Tree Regressor

- $R^2 = 0.8338$
- RMSE = 14540912642

Prediction des émissions pour 2020 pour le CO2 avec random_state=1

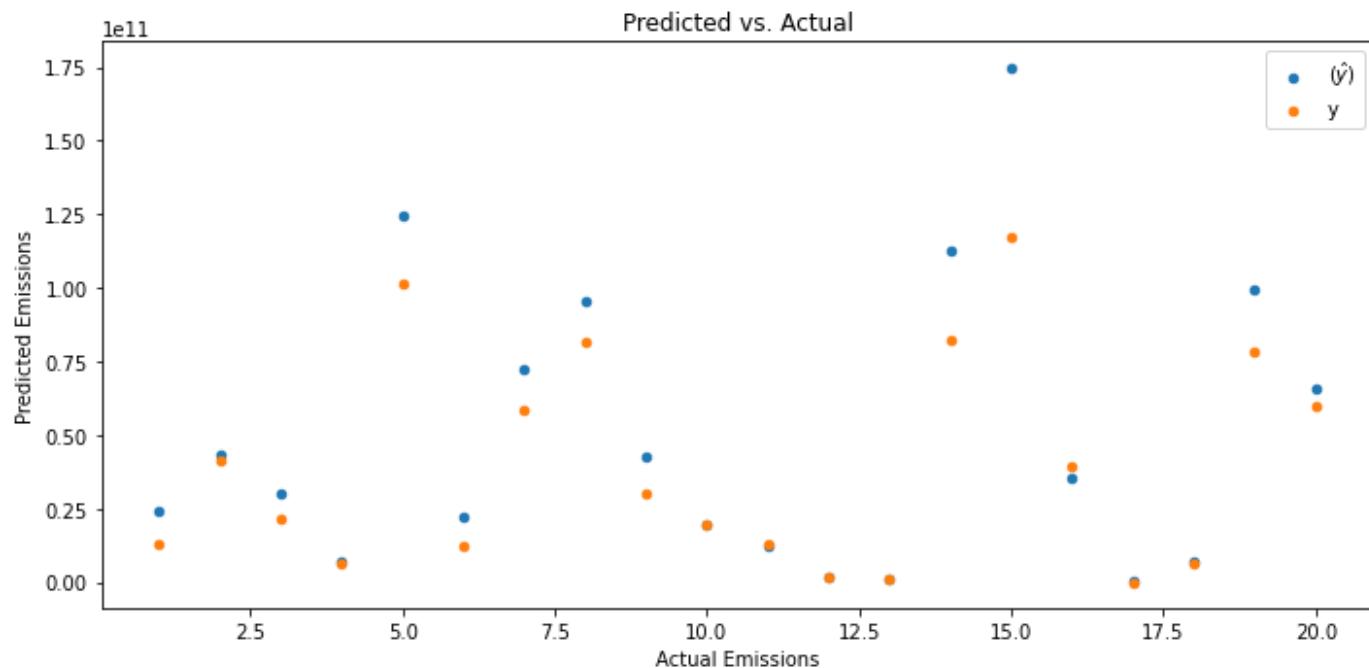




Random Forest Regressor

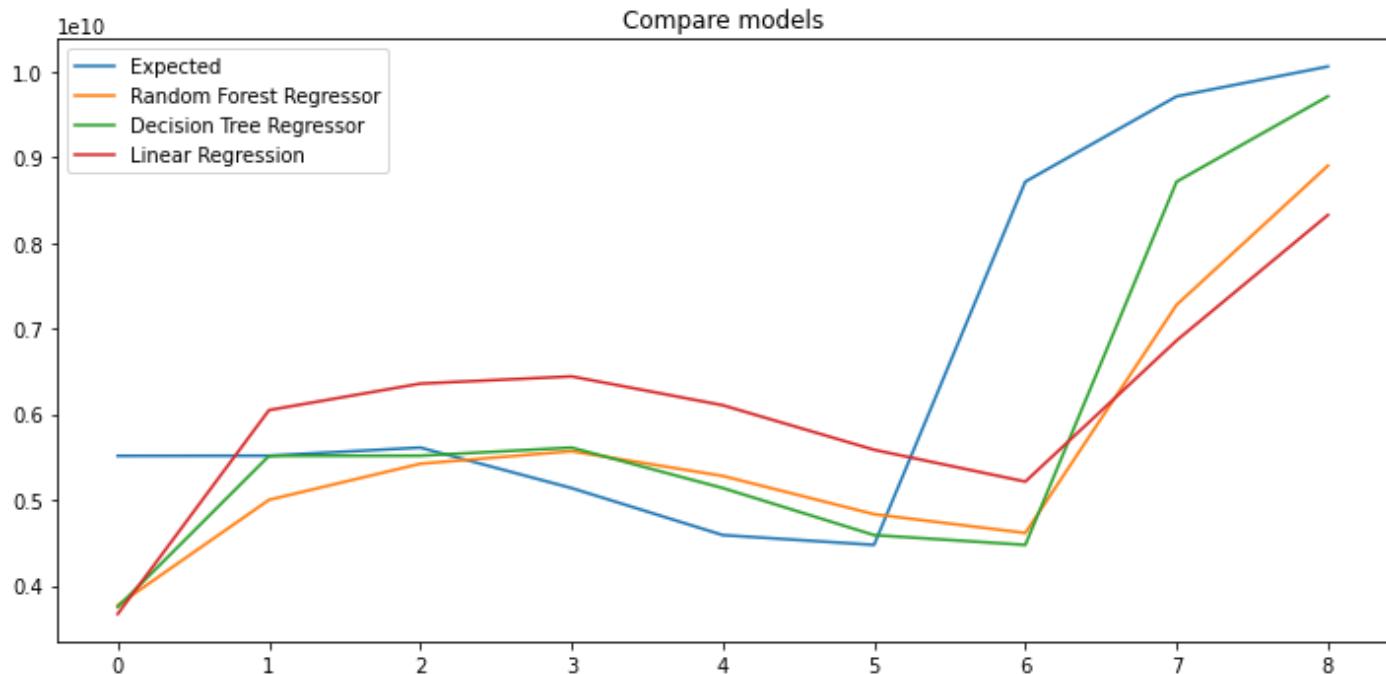
- $R^2 = 0.7613$
- RMSE = 17429321368

Prediction des émissions pour 2020 pour le CO2 avec n_estimators=1000, n_jobs=-1, random_state=1



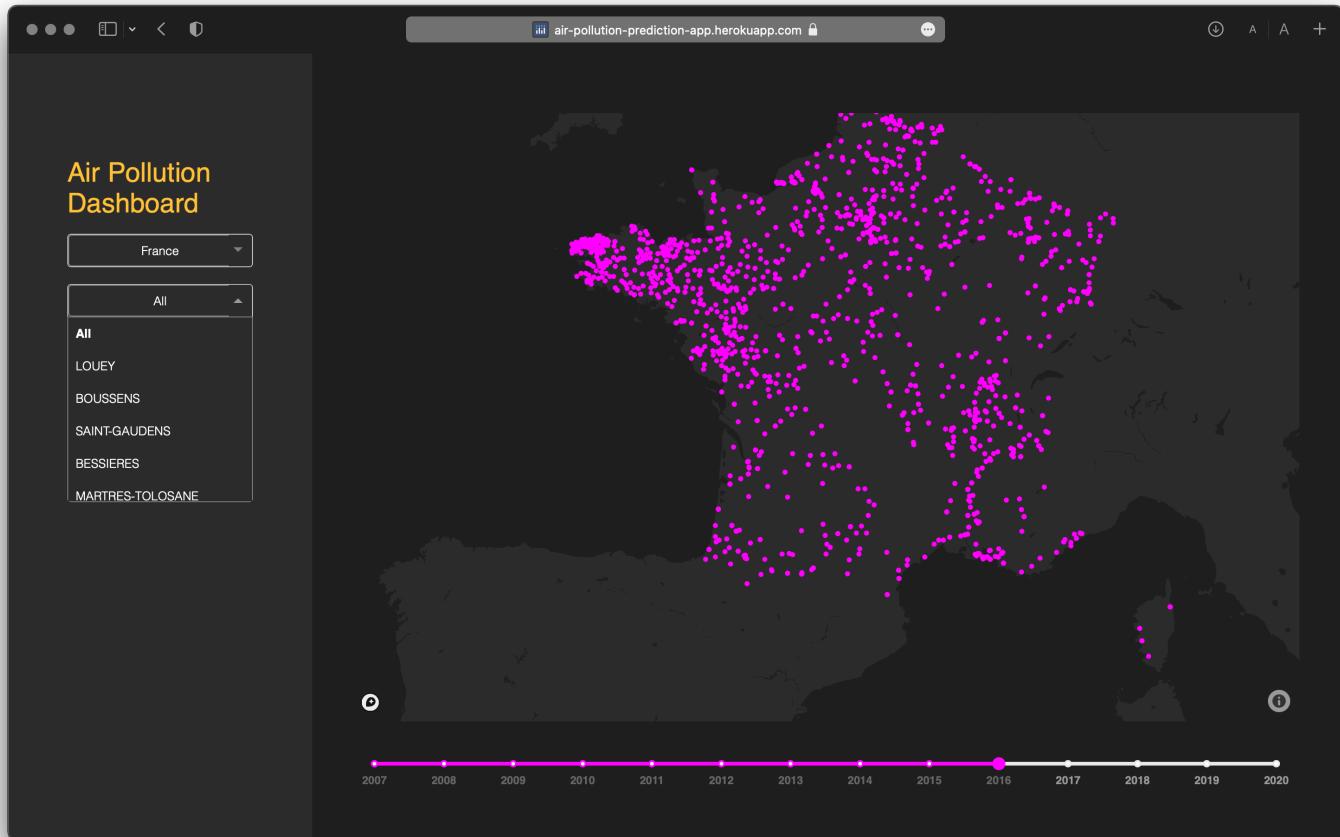


Comparer trois modèles





Dashboard





Analyse et Interpretation

- Decision Tree Regressor est plus rapide
- apprend et s'adapte vite
- lien émissions – secteur code vue en amputation

```
Results for DecisionTreeRegressor
R2 for train: [2007] test: [2008] : 0.997
R2 for train: [2007, 2008] test: [2009] : 0.966
R2 for train: [2007, 2008, 2009] test: [2010] : 0.983
R2 for train: [2007, 2008, 2009, 2010] test: [2011] : 0.997
R2 for train: [2007, 2008, 2009, 2010, 2011] test: [2012] : 0.995
R2 for train: [2007, 2008, 2009, 2010, 2011, 2012] test: [2013] : 0.994
R2 for train: [2007, 2008, 2009, 2010, 2011, 2012, 2013] test: [2014] : 0.989
R2 for train: [2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014] test: [2015] : 0.995
R2 for train: [2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015] test: [2016] : 0.992
R2 for train: [2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016] test: [2017] : 0.987
R2 for train: [2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017] test: [2018] : 0.988
R2 for train: [2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018] test: [2019] : 0.930
R2 for train: [2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019] test: [2020] : 0.834
```

- combine plusieurs arbres
- apprentissage lent
- adaptation lente aux changements

```
Results for RandomForestRegressor
R2 for train: [2007] test: [2008] : 0.826
R2 for train: [2007, 2008] test: [2009] : 0.918
R2 for train: [2007, 2008, 2009] test: [2010] : 0.952
R2 for train: [2007, 2008, 2009, 2010] test: [2011] : 0.990
R2 for train: [2007, 2008, 2009, 2010, 2011] test: [2012] : 0.990
R2 for train: [2007, 2008, 2009, 2010, 2011, 2012] test: [2013] : 0.995
R2 for train: [2007, 2008, 2009, 2010, 2011, 2012, 2013] test: [2014] : 0.987
R2 for train: [2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014] test: [2015] : 0.991
R2 for train: [2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015] test: [2016] : 0.987
R2 for train: [2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016] test: [2017] : 0.985
R2 for train: [2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017] test: [2018] : 0.987
R2 for train: [2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018] test: [2019] : 0.921
R2 for train: [2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019] test: [2020] : 0.761
```



Conclusion et Perspectives

ACCES AUX RESULTATS
Mise à jour des analyses, amélioration de l'application



RECHERCHE DE CAUSALITE

Utilisation des statistiques pour prédiction avancée



AJOUTER DES DONNEES

Remplissage des données, recherche des nouvelles





Merci
pour votre attention!