

ModStatSAP 19/11/2023

# **Peut-on améliorer les prédictions d'une épiphytie en combinant les résultats de plusieurs méthodes de machine learning ? Leçons d'une expérience numérique**

**César Martinez**<sup>1</sup>, Edith Gabriel<sup>1</sup>, Ghislain Géniaux<sup>2</sup>, Dorian Chauvin<sup>1</sup>, Samuel Soubeyrand<sup>1</sup>

<sup>1</sup> INRAE Avignon BioSP

<sup>2</sup> INRAE Avignon Ecodev

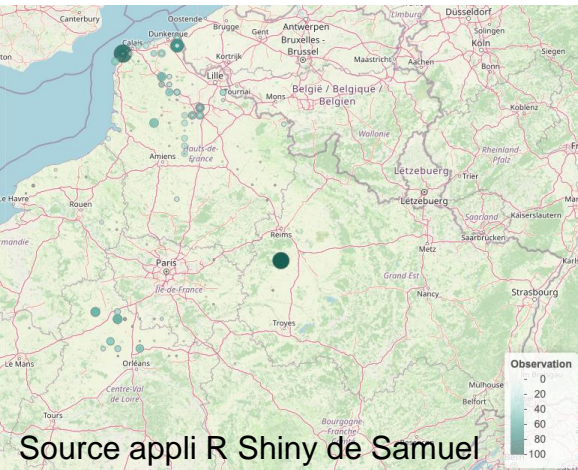
# PLAN

- 1) Contexte
- 2) Expérience numérique
- 3) Conclusion

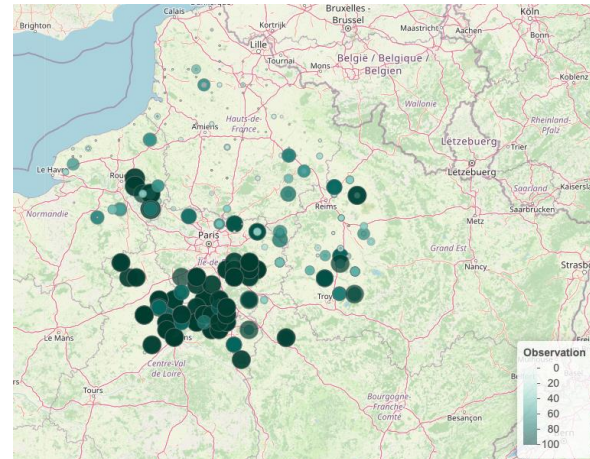
# Jaunisse de la betterave et données



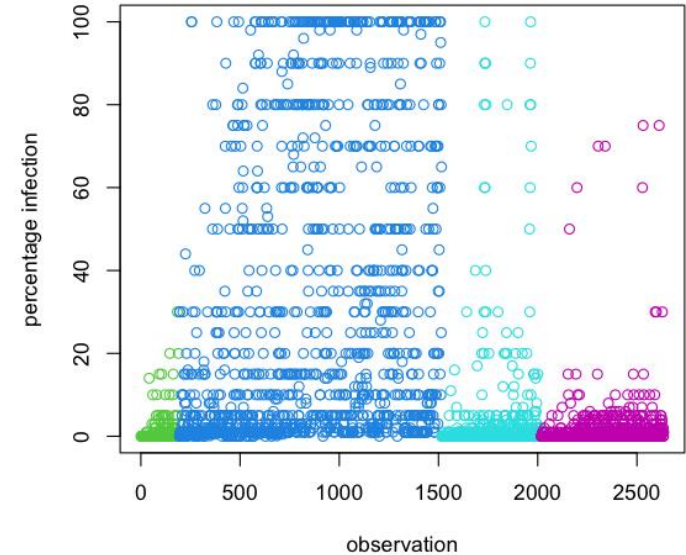
$Y_{obs} = \%$   
parcelle jaune



2019



2020



Observations 2019, 2020,  
2021, 2022

# Dans SEPIM (travail de Dorian)

## Bases de données

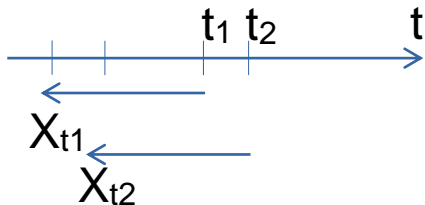


- Variables agronomiques
- Variables environnementales

Variables fixes



Variables « glissantes » ex :



**Modèles** de machine learning  
(prédire la jaunisse)

Earth  
Rf  
Cubist  
XgBoost

$M_1(X), M_2(X), \dots, M_k(X)$

**Agrégation ?**

$Y = f(M_1(X), M_2(X), \dots, M_k(X))$  ?

# Méthodes d'agrégation

(doivent fonctionner avec des modèles de ML)



Moyenne

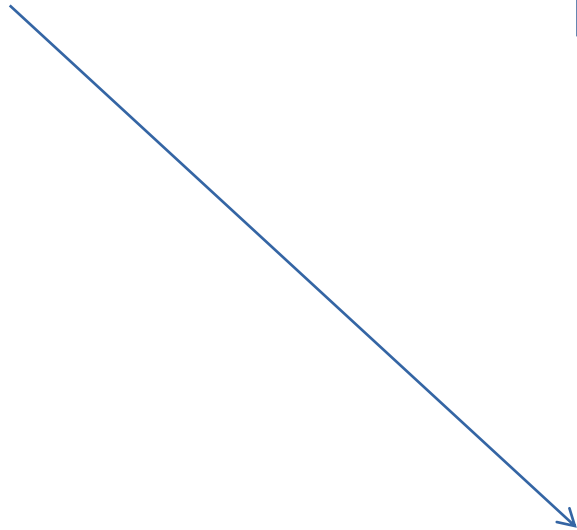
Médiane

Régression  
 $Y_{out} =$   
 $EARTH(M_1(X), \dots, M_k(X))$



Stacking

BMA-EM

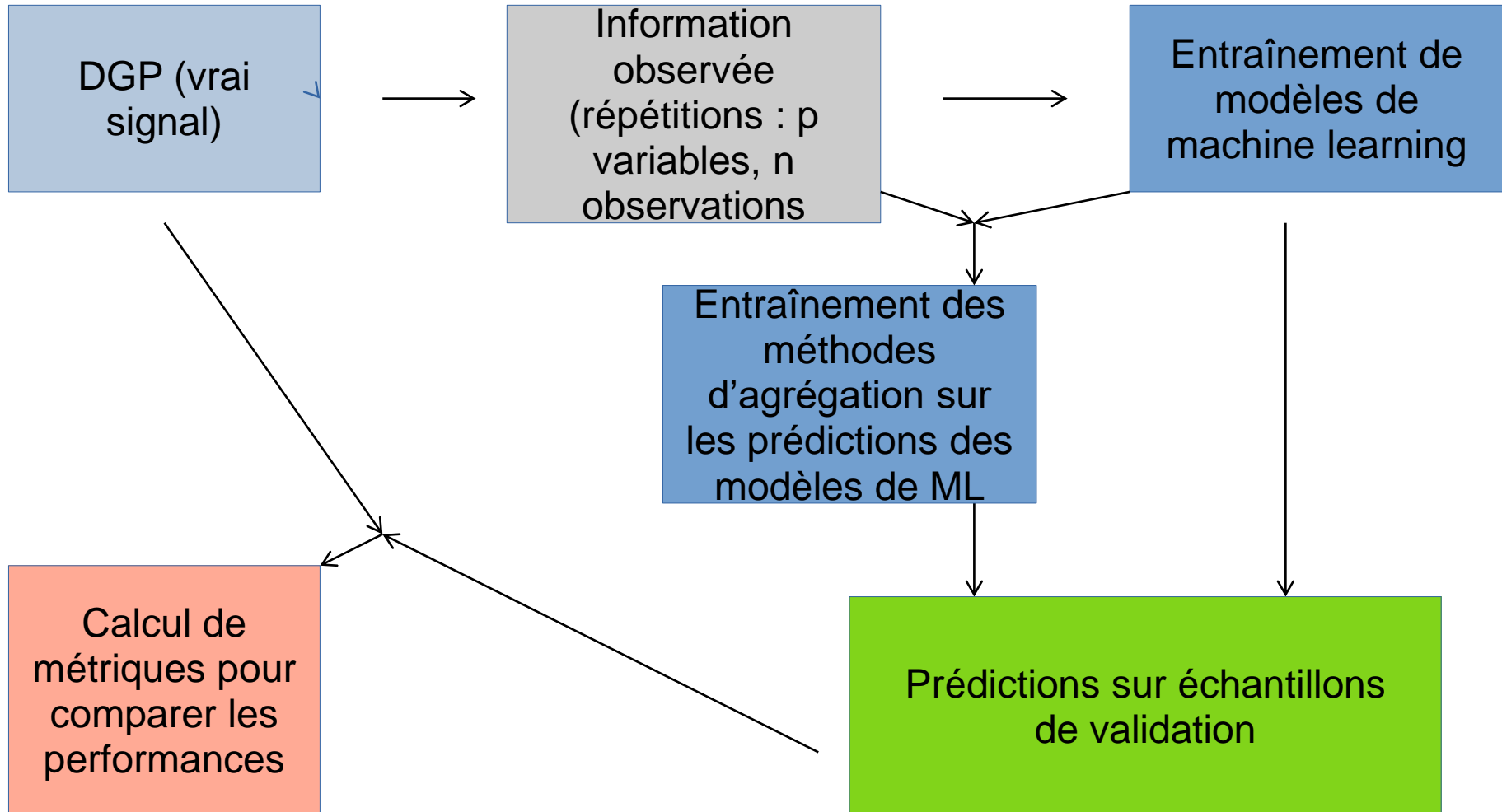


MGWR

Clustering

ANN

# Mise en place d'une expérience numérique

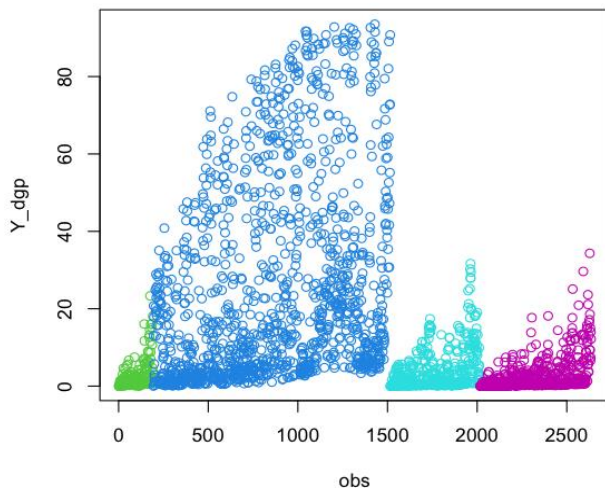


# DGP, échantillonnage

$$Y\_DGP = f(X(\text{SEPIM}))$$

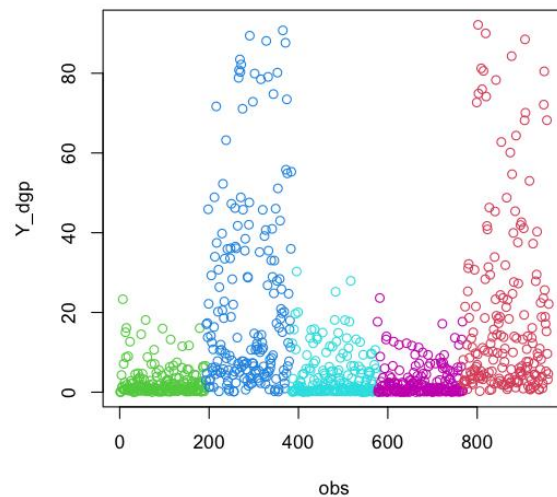
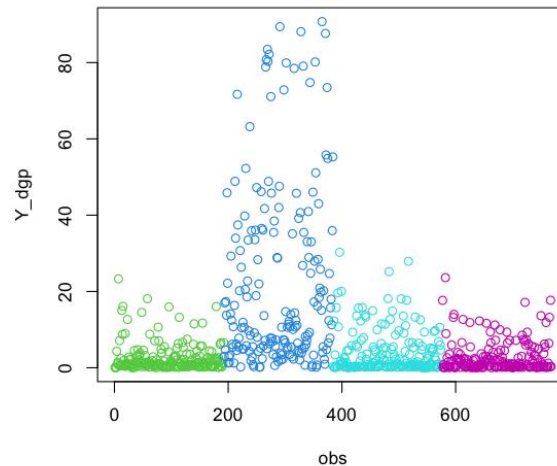
→Modèle de régression  
logistique

→Prédictions in sample



192  
observations  
par an

Ajout d'un  
deuxième  
tirage de  
192 dans  
2020





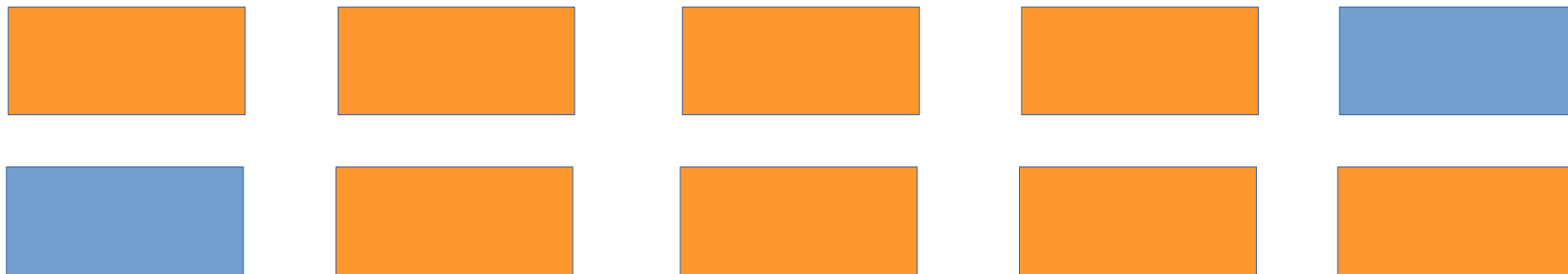


# Cross-validation

Données sur 5 années qui permettront de définir les feuillets

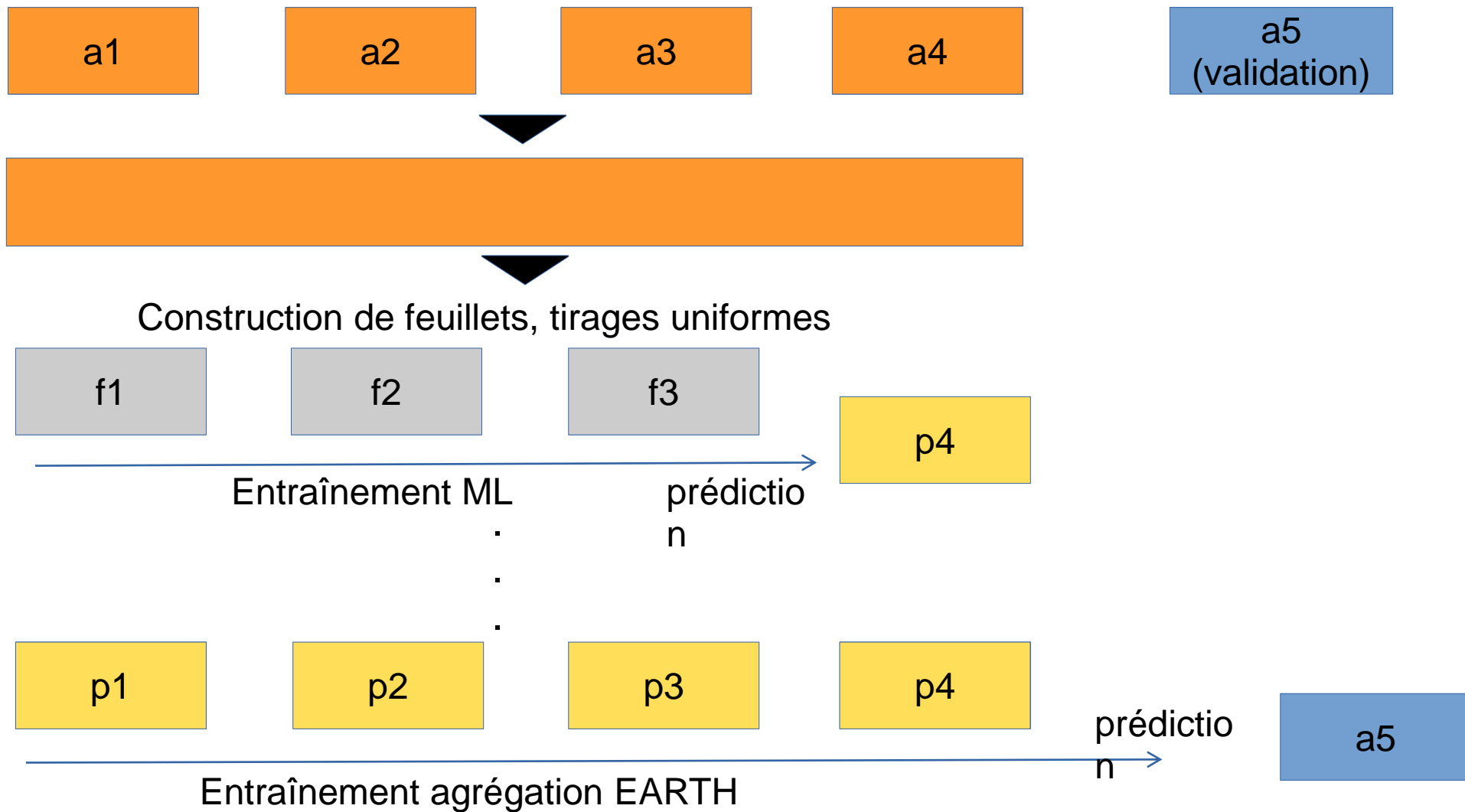
4 années pour entraîner les modèles

1 année de validation



▪  
▪  
▪

# Entraînement de la méthode de régression pour l'agrégation



# Mesure des performances

Mean absolute error (MAE) par rapport au DGP, sur l'ensemble des données

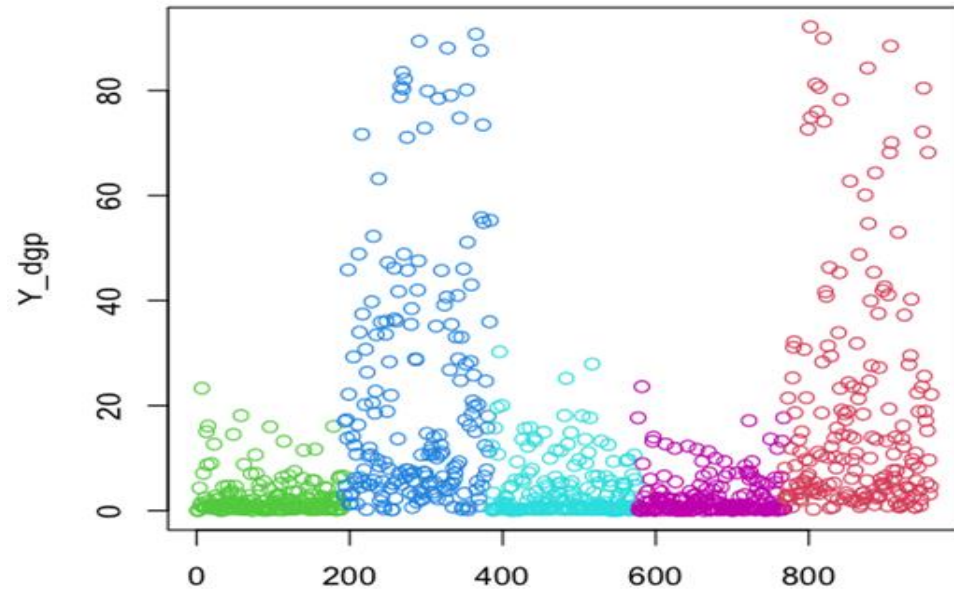
Pour des infections > 20 %

$$\text{TPR} = \text{TP}/\text{P}$$

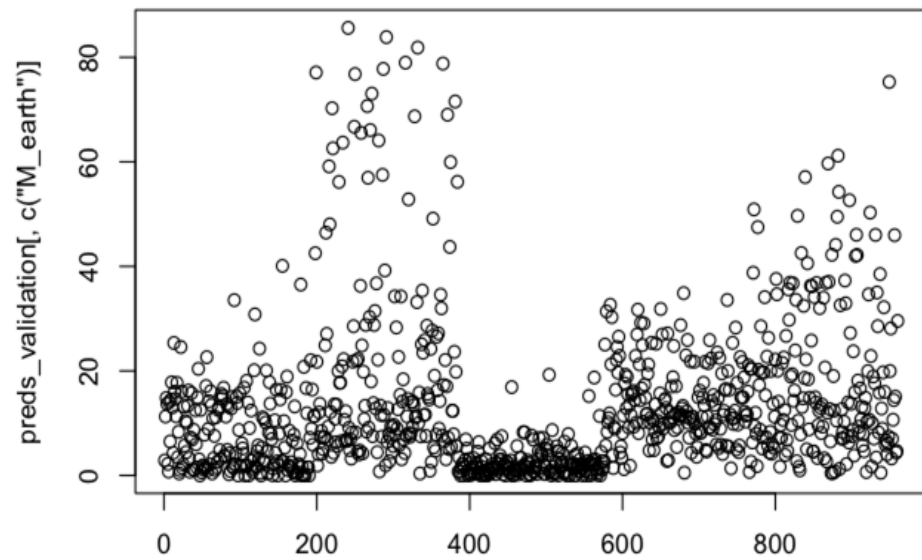
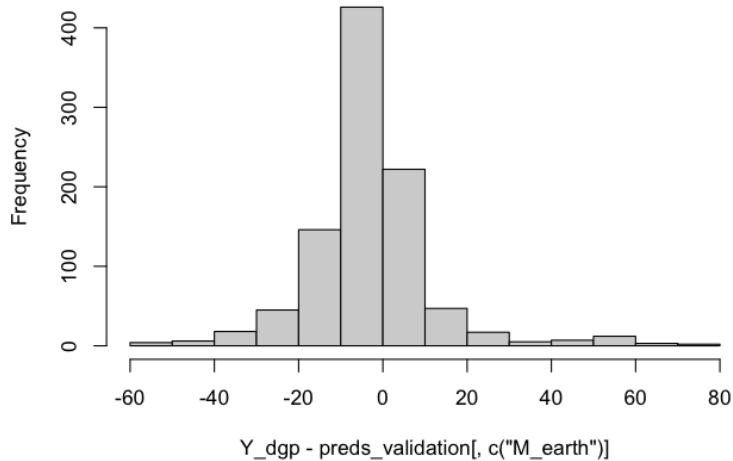
TPR pour des infections < 3 %

# Expérience :

5 échantillons : 2019,  
2020<sub>1</sub>, 2021, 2022, 2020<sub>2</sub>

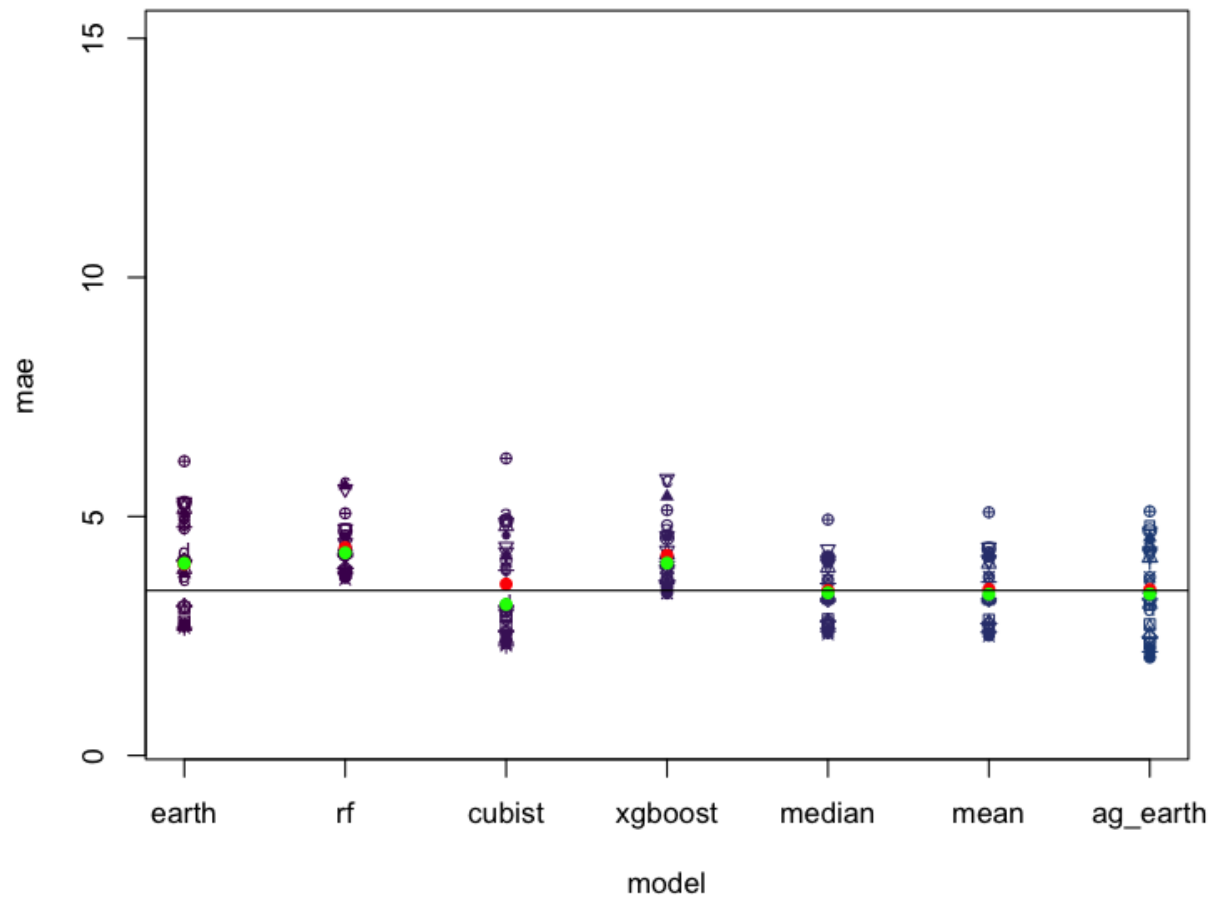


Histogram of Y\_dgp - preds\_validation[, c("M\_earth")]

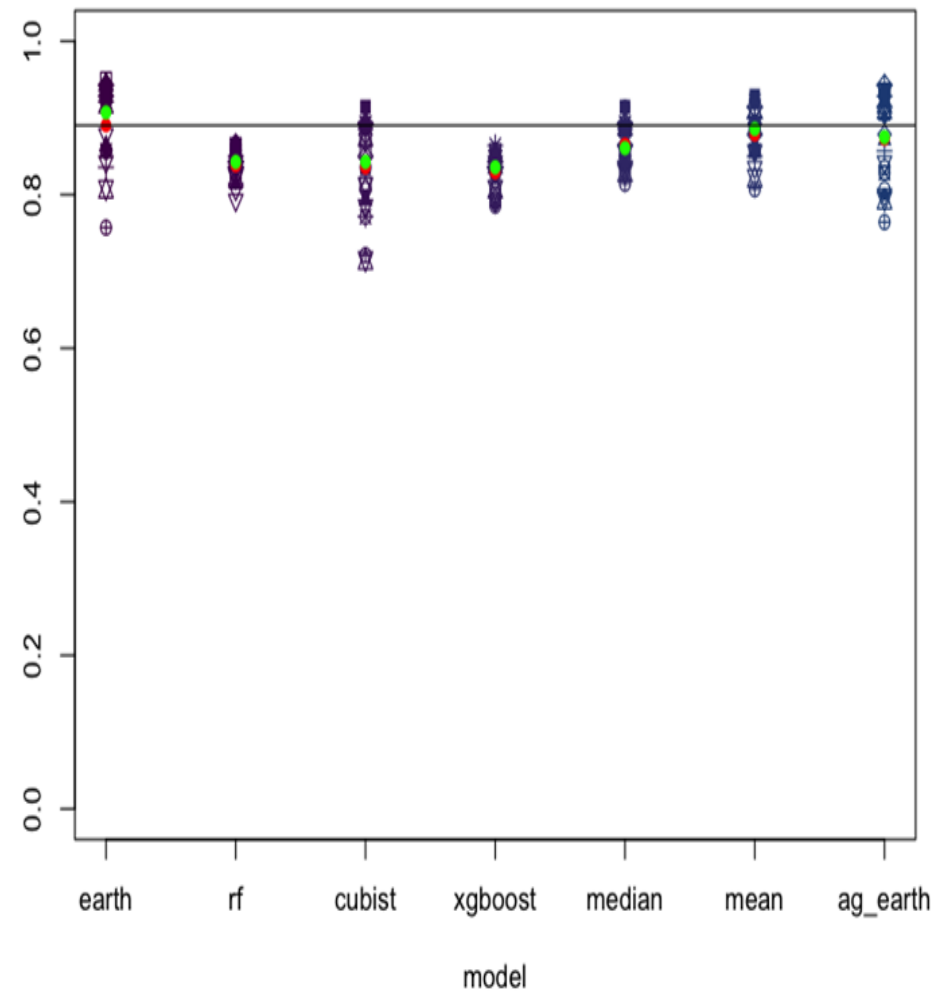


# Expérience : mae

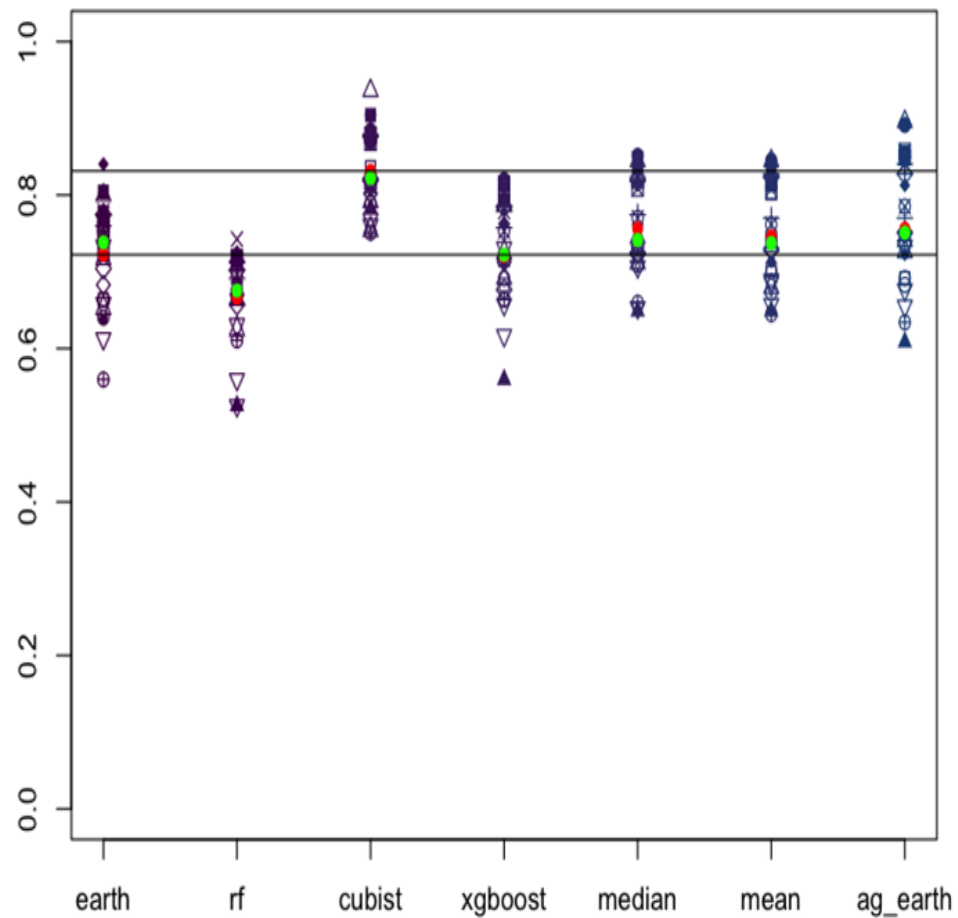
5 years, 30 repetitions, dgp logistic

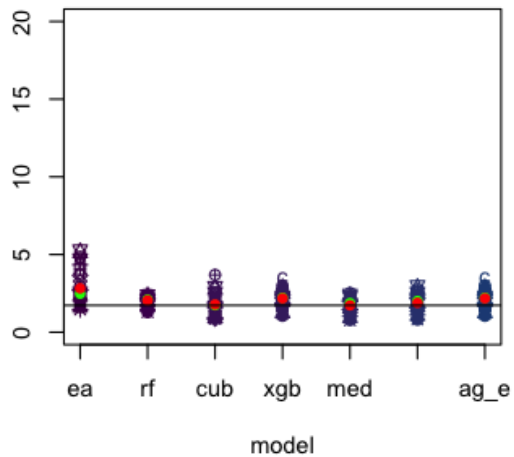
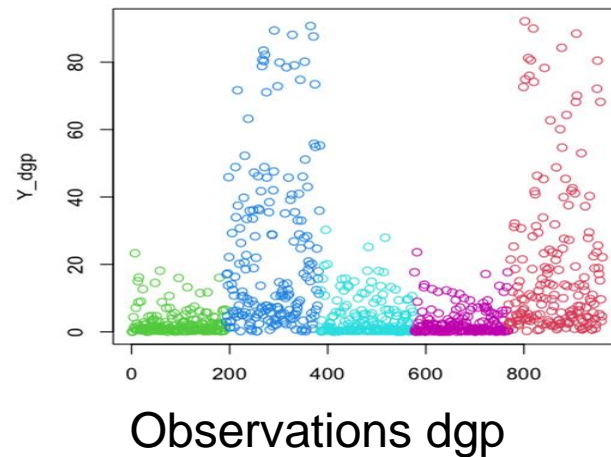
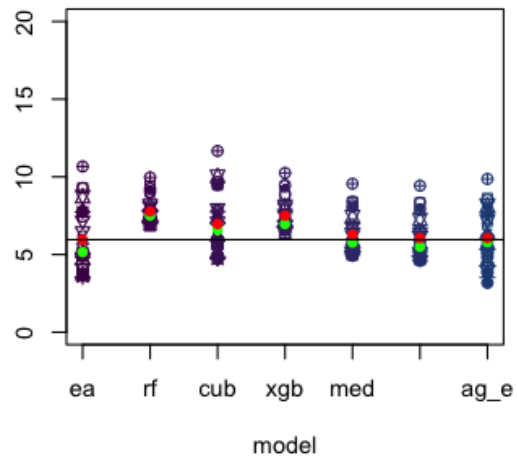
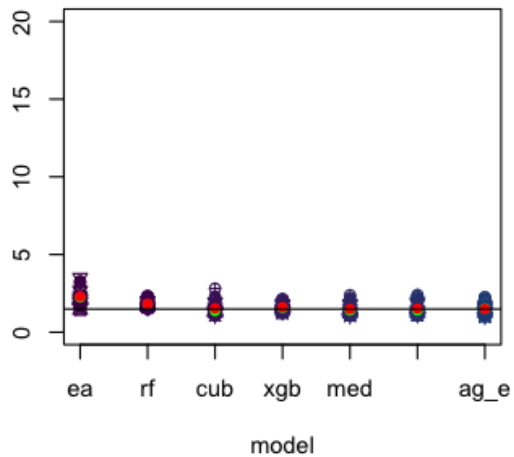
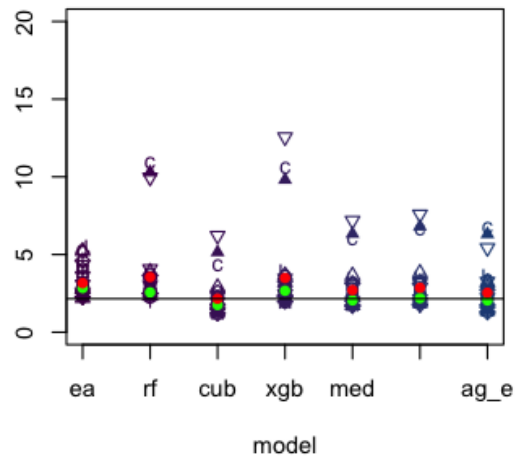
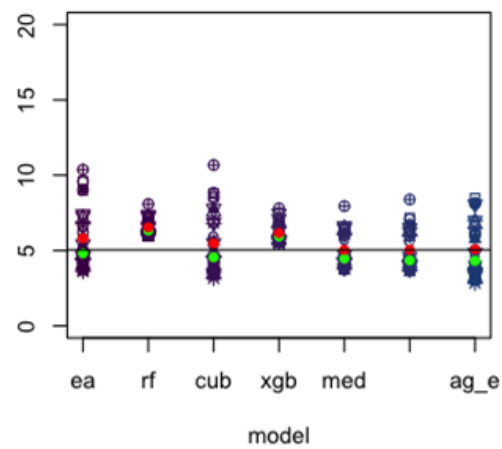


>20 % bien identifiés



<3 % bien identifiés



**mae 1****mae 2****mae 3****mae 4****mae 5**

# Conclusion

## Expérience préliminaire

- 1) Tenir compte des effets d'échantillonnage
- 2) Choisir des métriques adaptées à la question
- 3) Agrégation de modèles intéressante dans le cas étudié

## Pour poursuivre :

### **Finesse des méthodes d'agrégation :**

- comparaison des méthodes
- prendre en compte des variations spatiales

### **Enrichir le jeu de modèles :**

Couplage des modèles statistiques avec des modèles « physiques »

- date d'arrivée des pucerons
- modèle de développement de la betterave
- modèle de propagation épidémique (spatio-temporel)



MERCI POUR VOTRE ATTENTION !