

2021

# A brief introduction to machine learning and deep learning

David Makowski

INRAE

<https://www6.inrae.fr/mia-paris/Equipes/Membres/David-Makowski>

# Outline

- Definition & main principles
- Several extensions of linear regression
- Trees and forests
- Deep learning

# Outline

- Definition & main principles
- Several extensions of linear regression
- Trees and forests
- Deep learning

Artificial intelligence  
Machine learning

# Artificial intelligence

## Machine learning

### Supervised learning

Objective: « Learning a function that maps an input to an output based on examples of input-output pairs »

# Statistical Modeling: The Two Cultures (Breiman, 2001)

$$y = f(x) + e$$

Modelling approach 1: Try to find the true  $f(x)$

Modelling approach 2: Predict  $y$  from  $x$  as accurately as possible

# Statistical Modeling: The Two Cultures (Breiman, 2001)

$$y = f(x) + e$$

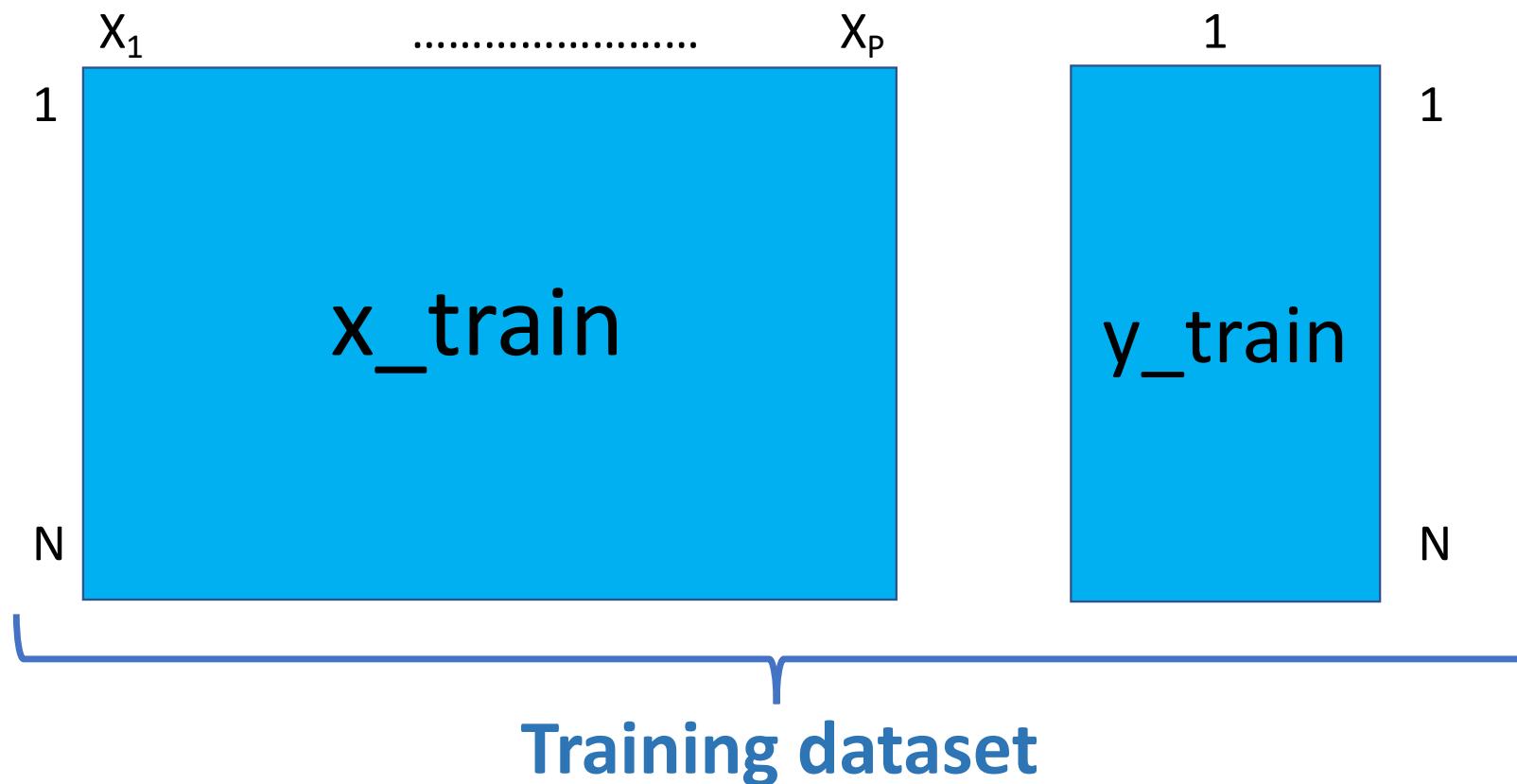
Modelling approach 1: Try to find the true  $f(x)$

**Modelling approach 2: Predict  $y$  from  $x$  as accurately as possible**

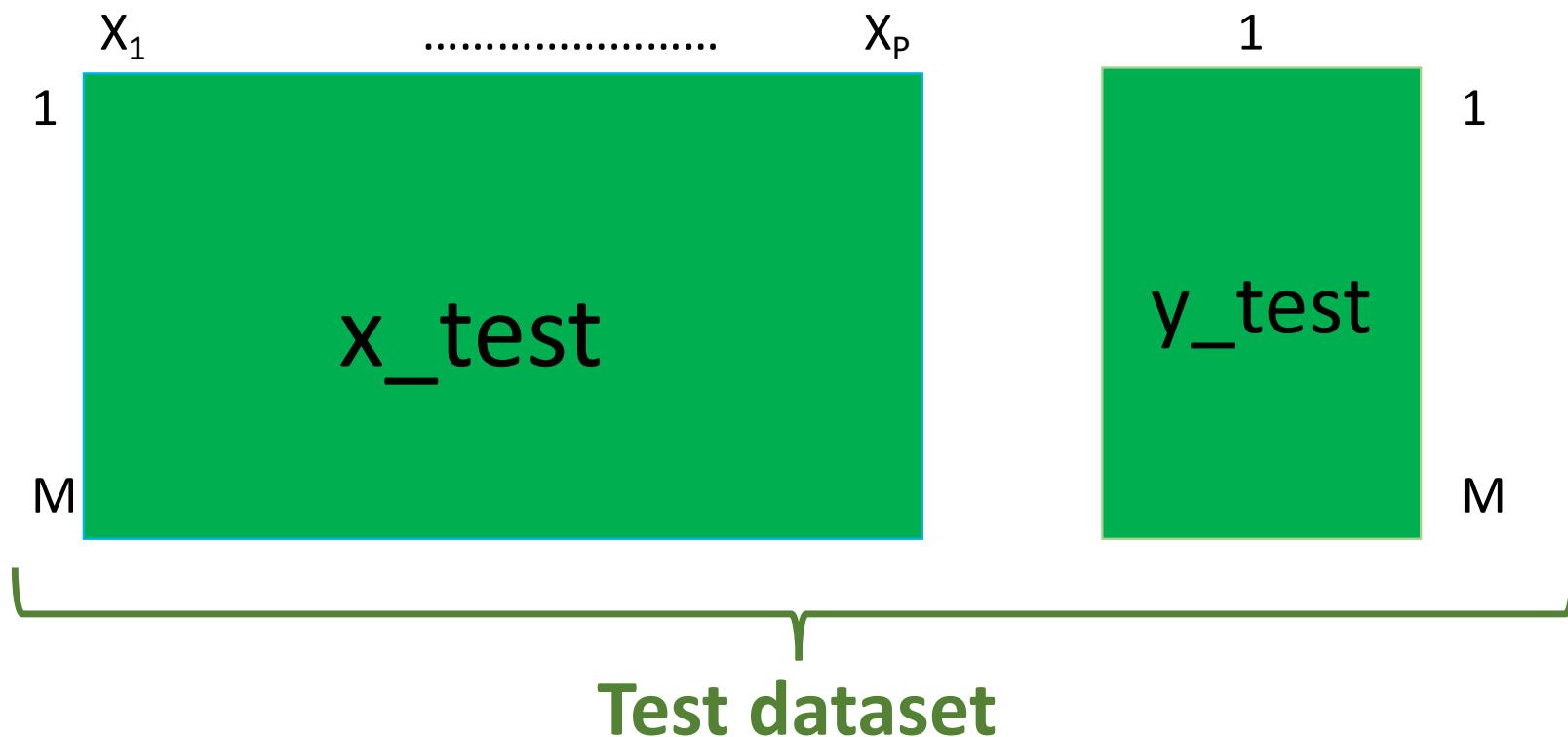
# Two main steps

- Step 1: Training
- Step 2: Test

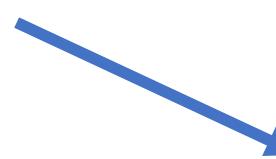
**Step 1:** Train an algorithm predicting  $Y$  as a function of  $X_1, \dots, X_p$  using a **training dataset**



**Step 2:** Assess the predictive capability of the trained algorithm using a **test dataset**

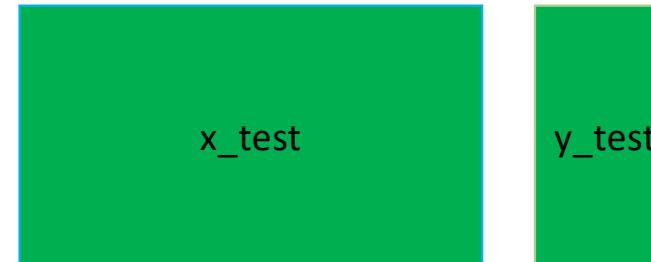


## Training dataset



$f(X_1, \dots, X_p)$   
*Trained algorithm*

## Training dataset

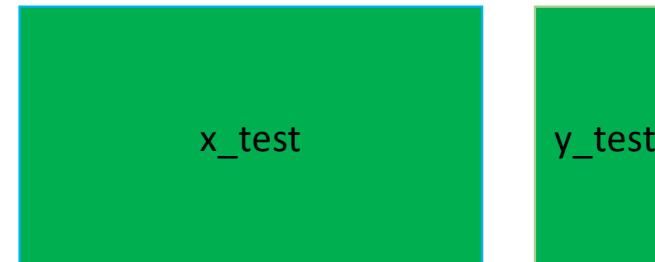
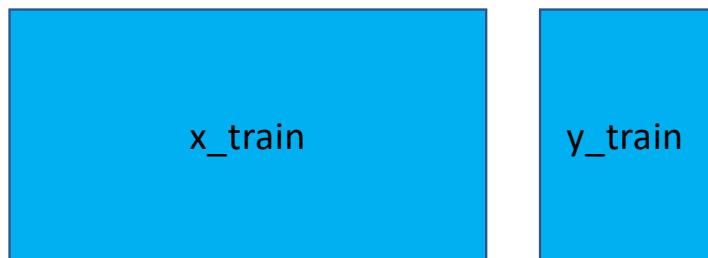


## Test dataset

$f(X_1, \dots, X_p)$   
*Trained algorithm*

Predicted  
y values

## Training dataset

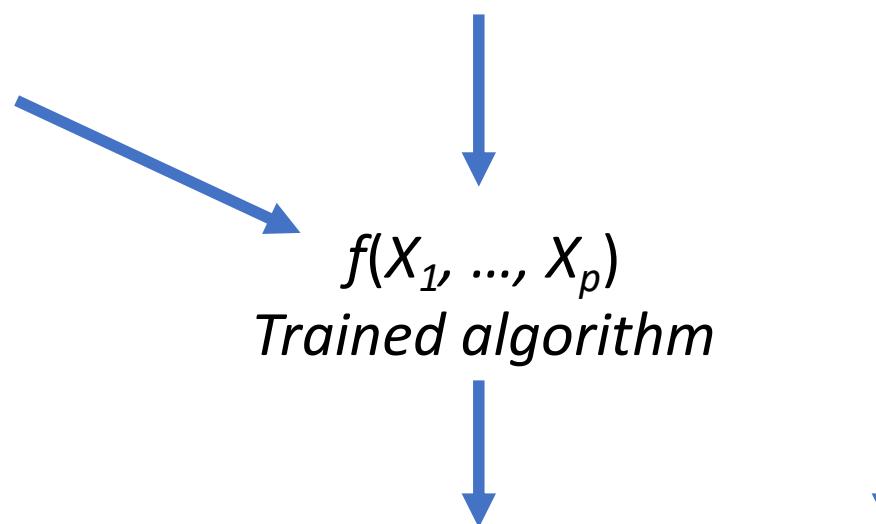


$f(X_1, \dots, X_p)$   
*Trained algorithm*

Predicted  
y values

y\_test

Comparison



kaggle

Q Search

Competitions Datasets Notebooks

# Competitions

-  Home
-  Compete
-  Data
-  Notebooks
-  Discuss
-  Courses
-  More

## Genentech

A Member of the Roche Group

### Flu Forecasting

Predict when, where and how strong the flu will be

\$125,000 · 50 teams · 6 years ago

Overview Data Discussion Leaderboard Rules

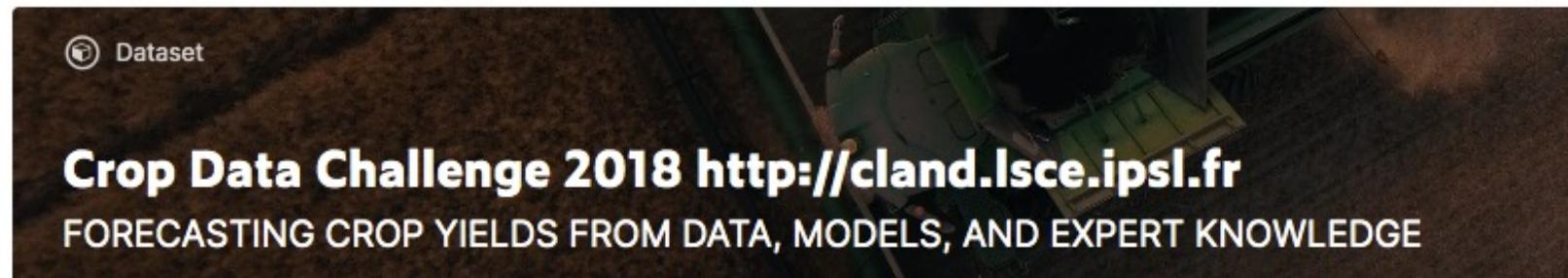
« The objective of this competition is to build an algorithm that helps predict the occurrence, peak and severity of influenza in a given season ».

 In the money    Gold    Silver    Bronze

#	△pub	Team Name	Notebook	Team Members	Score ⓘ
1	—	Alfonso Nieto-Castanon			0.47415
2	—	J.A. Guerrero (Datrik Intelligen...			0.47567
3	—	Zhanpeng Fang			0.47573
4	—	Tim Salimans			0.47650
5	—	Victor			0.47708
6	—	Nitai Dean			0.48110
7	—	BenPlus			0.48665

RMSE

-  Home
-  Compete
-  Data
-  Notebooks

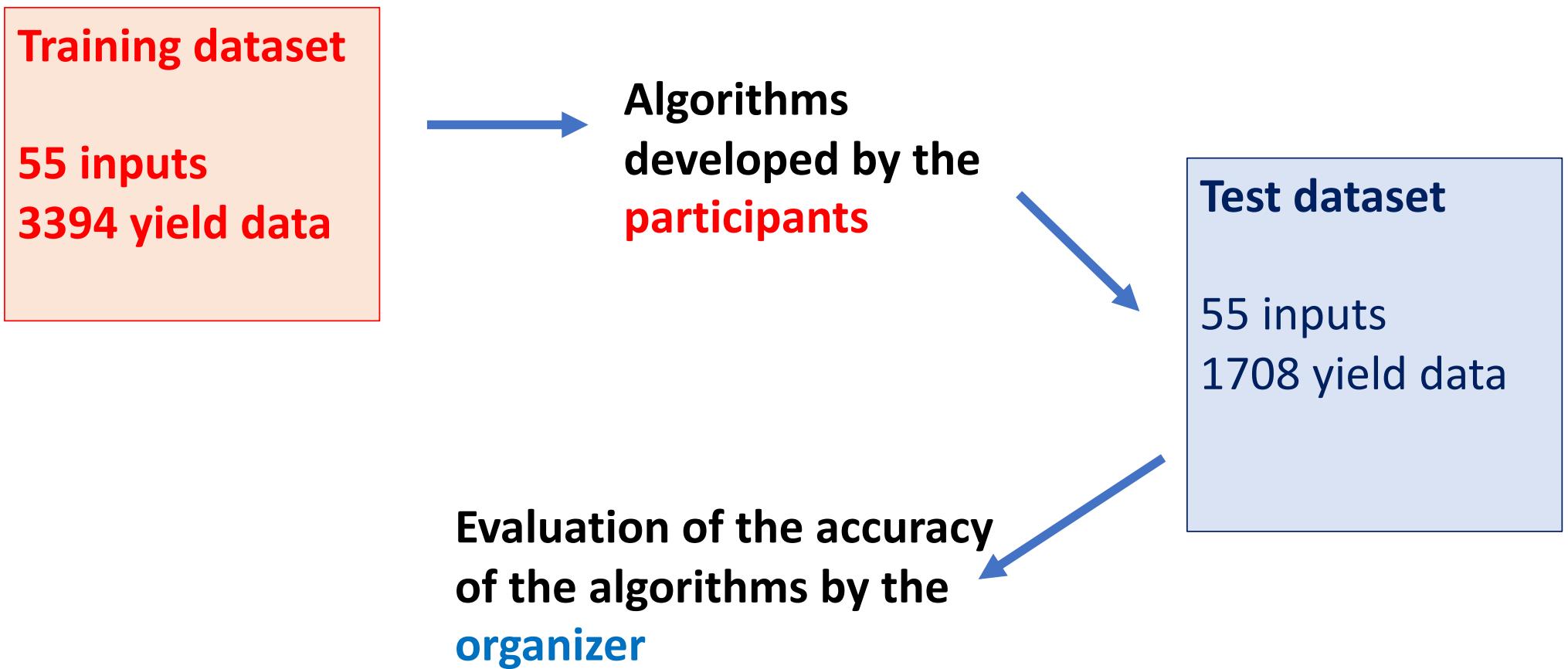


## Data (5 MB)

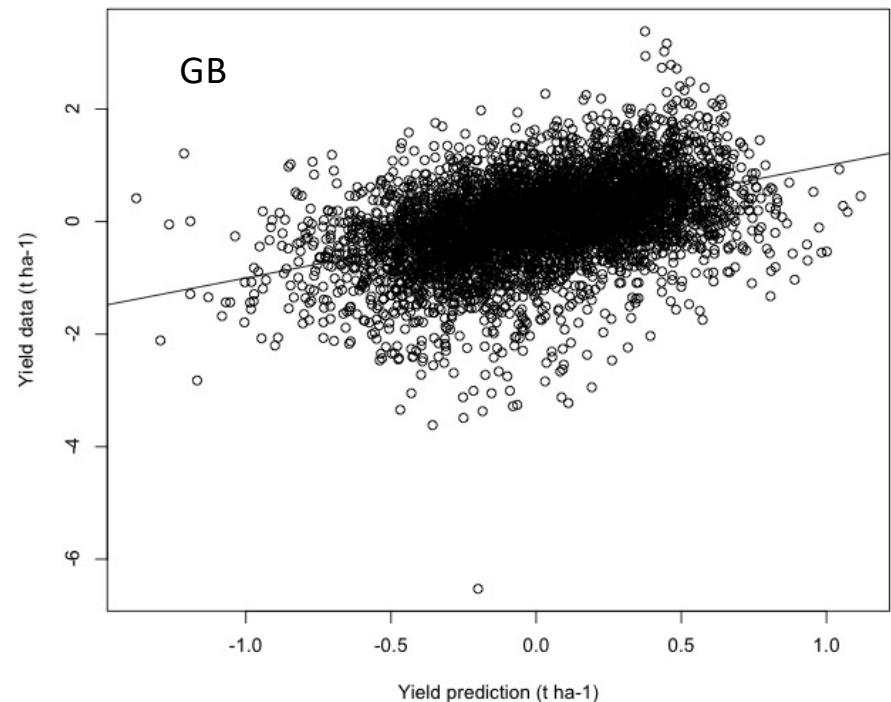
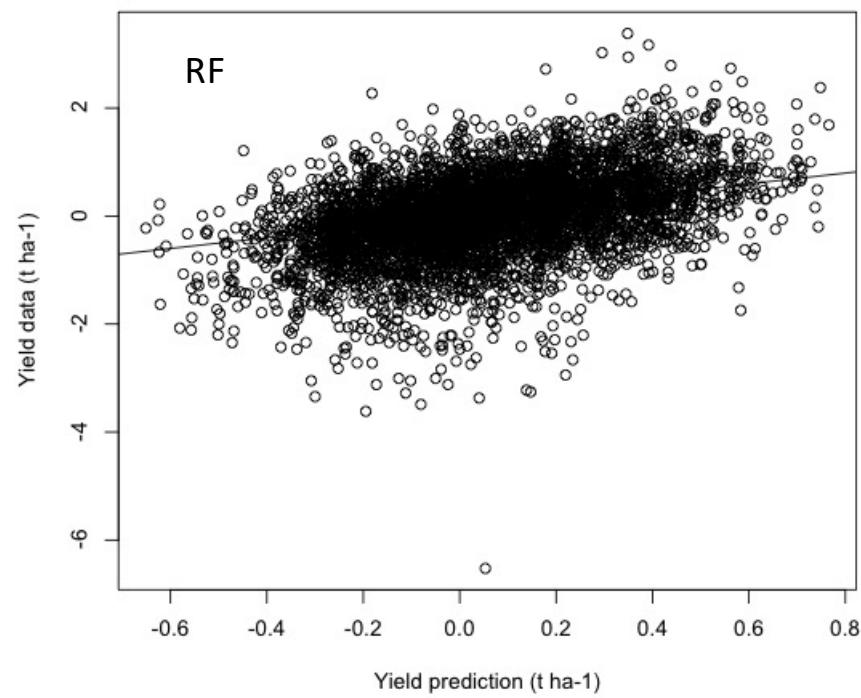
### Data Sources

-  TestDataSet\_Ma... 57 columns
-  TestDataSet\_W... 92 columns
-  TrainingDataSet... 58 columns
-  TrainingDataSet... 93 columns

# French maize yield prediction (départements)



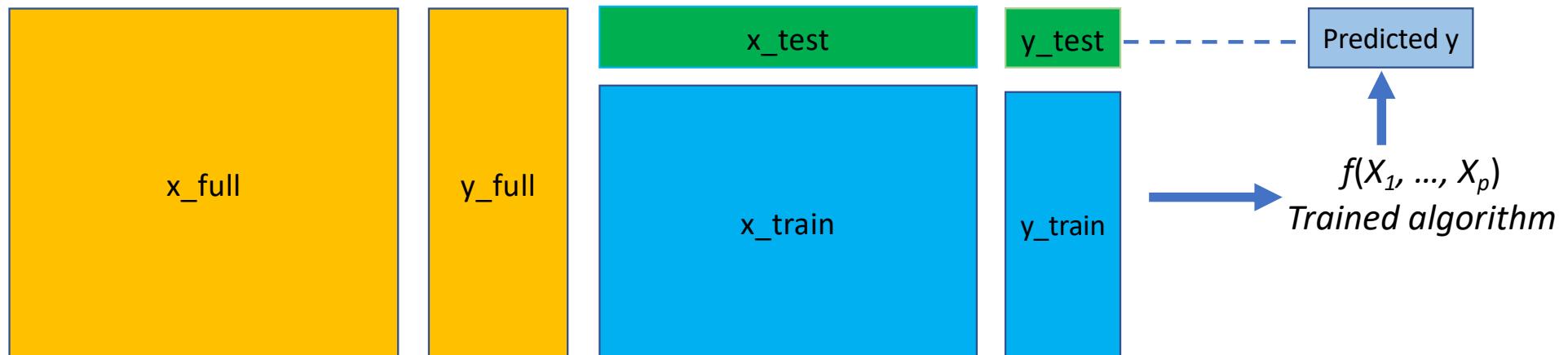
Method	RMSEP (maize yield)
Random Forest (RF)	0.71 t/ha
Gradient boosting (GB)	0.70 t/ha



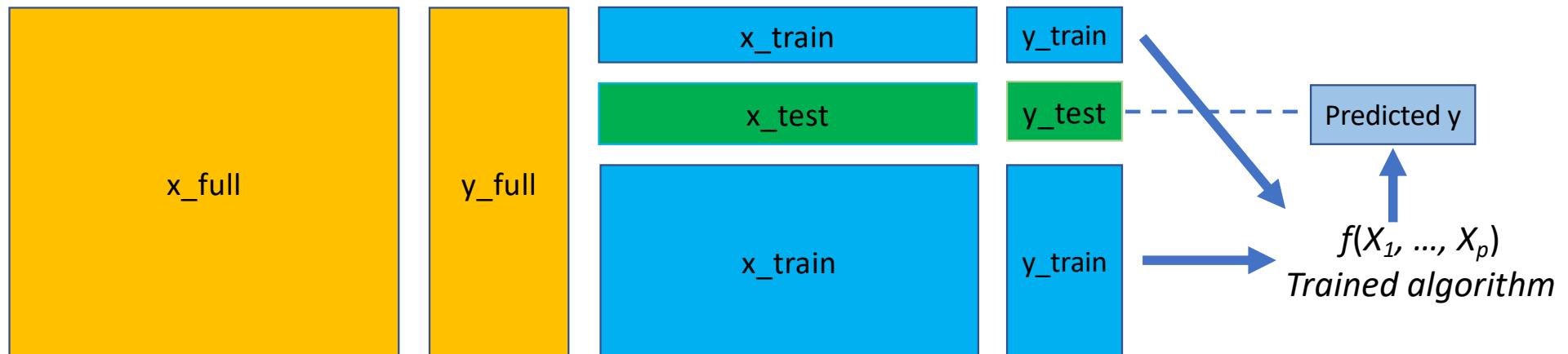
Model testing should be taken seriously to avoid risk of overfitting

Cross-validation is used when no independent test dataset is available

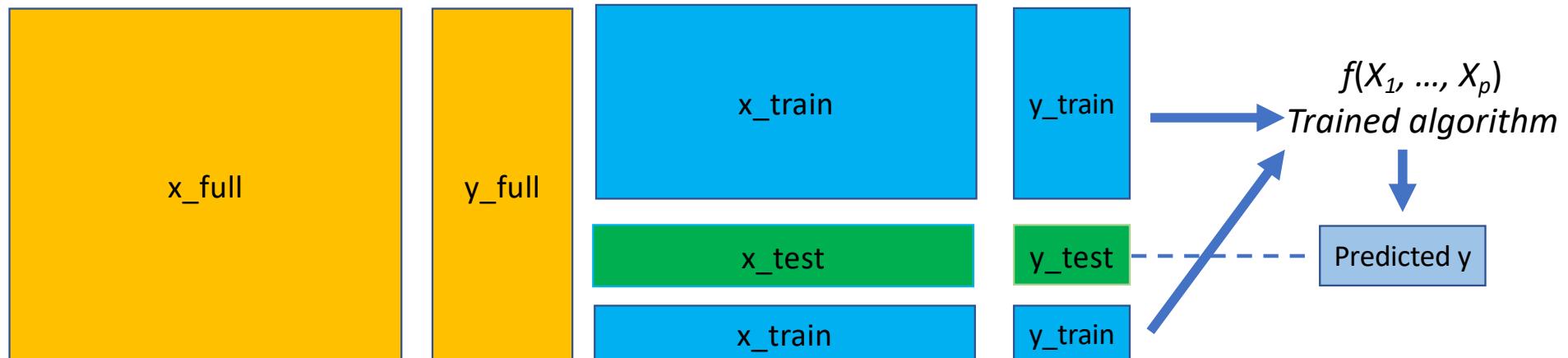
## Full dataset



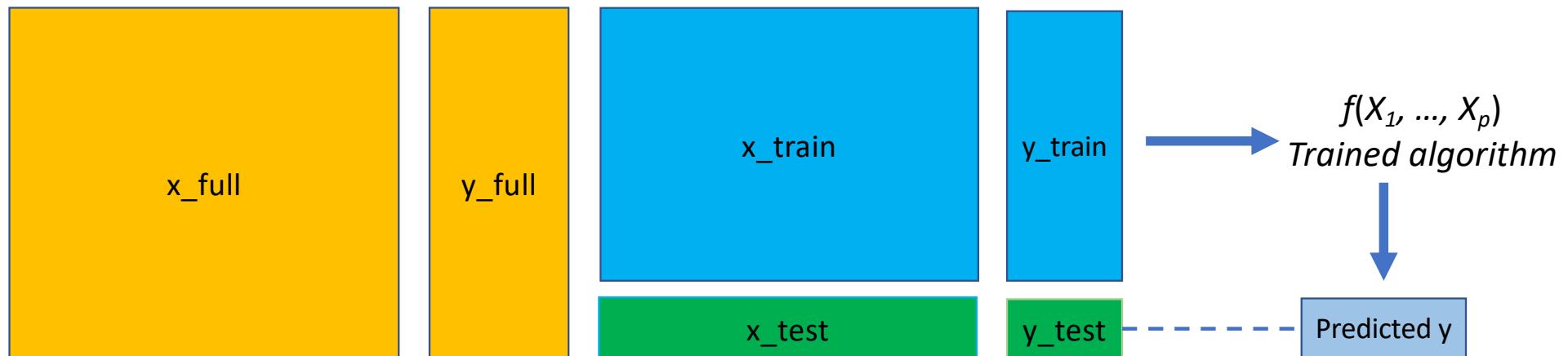
## Full dataset



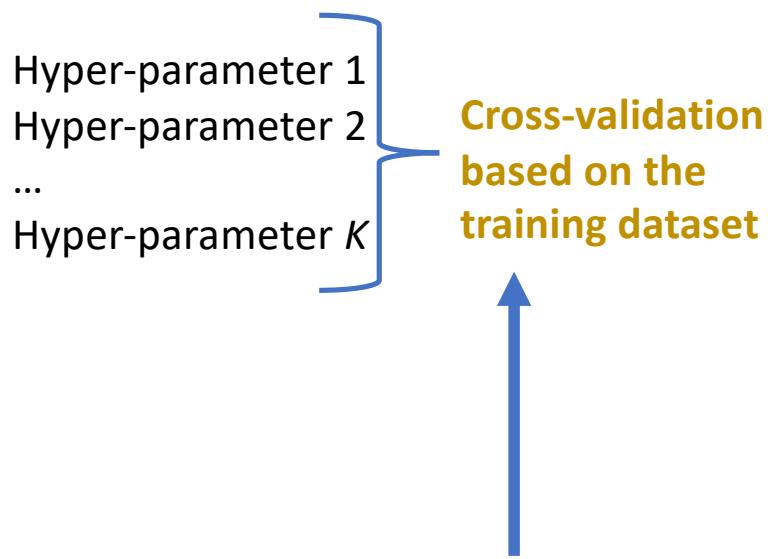
## Full dataset



## Full dataset



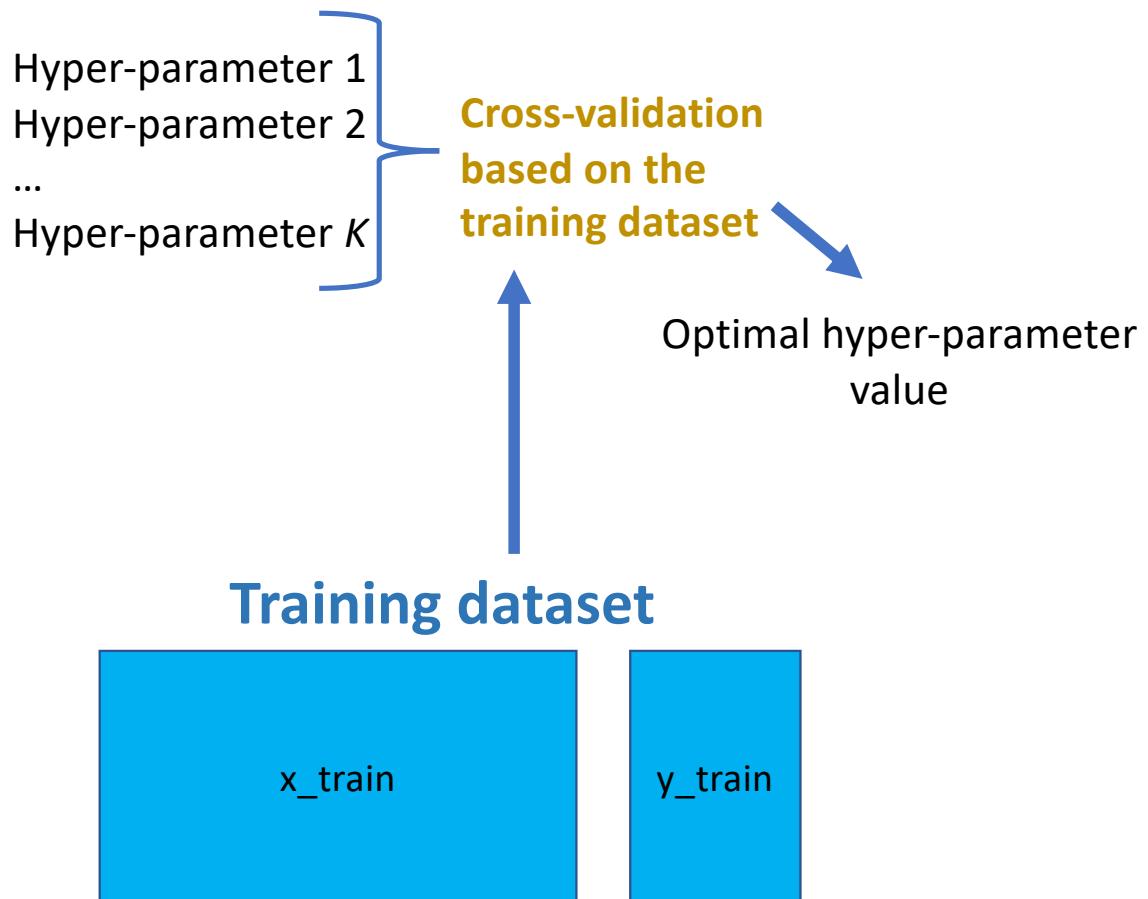
Cross-validation and test dataset can be combined together

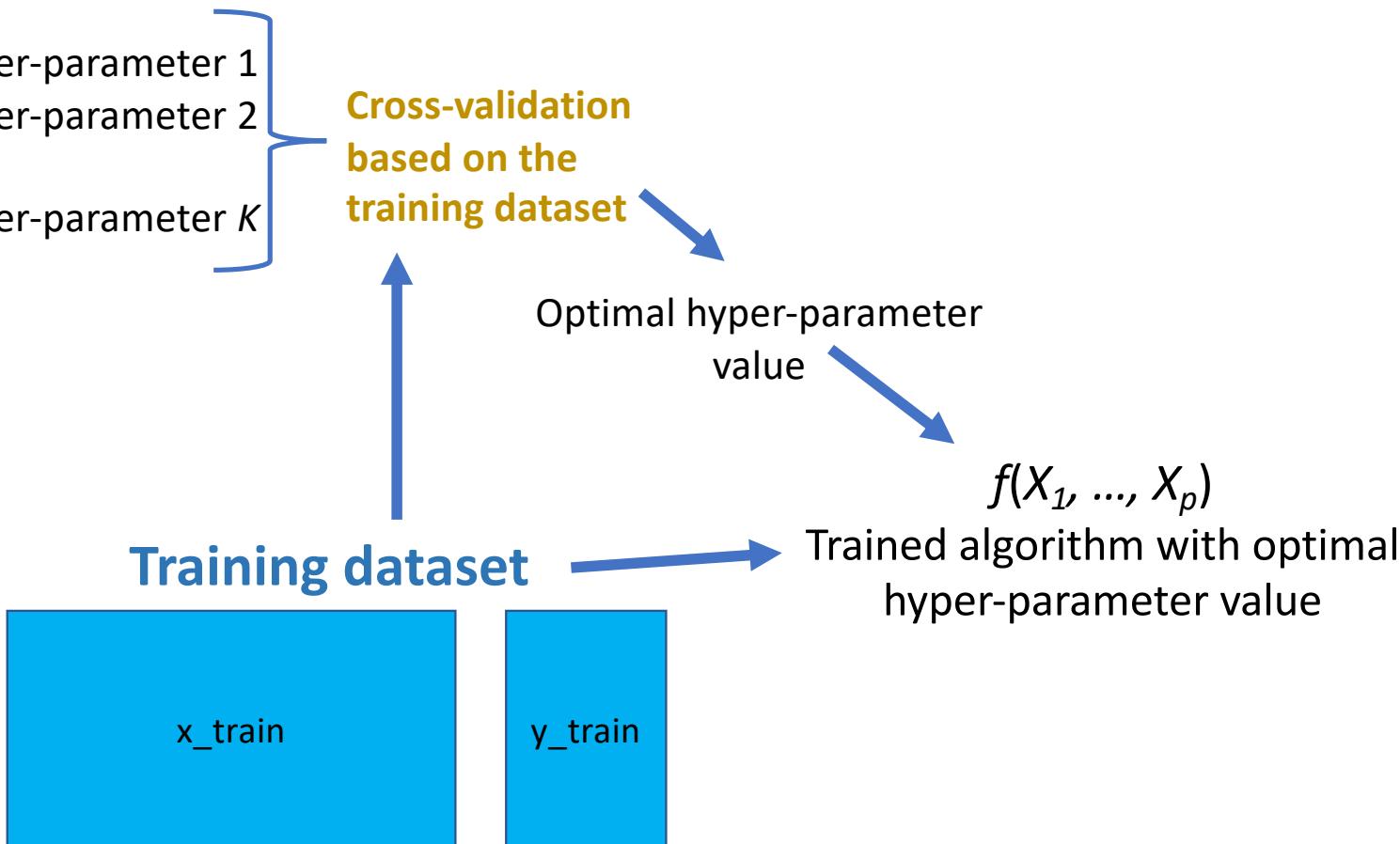


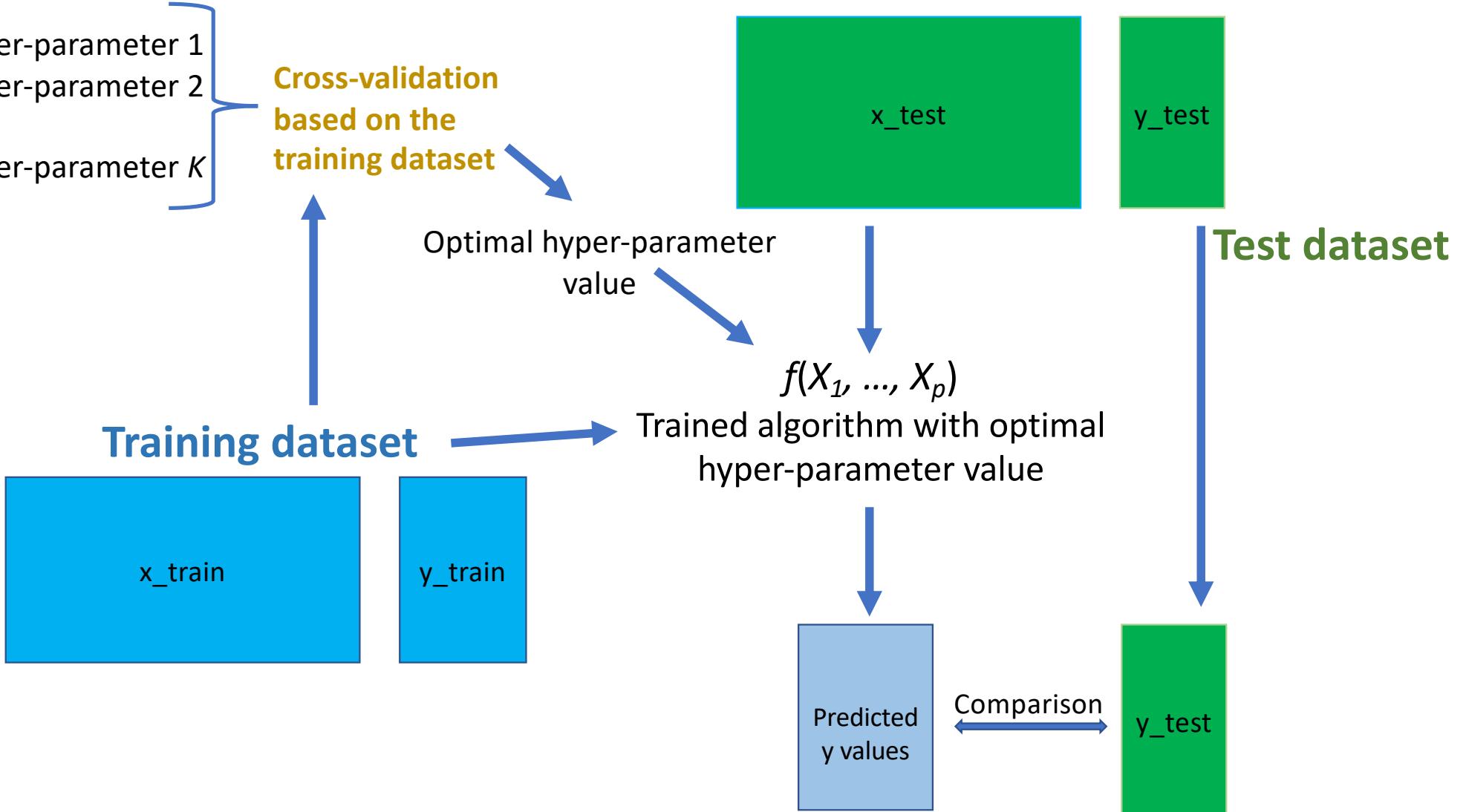
Training dataset

x\_train

y\_train







# Why machine learning is powerful?

Very flexible methods

+

Computational power

+

Large datasets

Increased chance to  
obtain accurate  
predictions

# Why machine learning is powerful?

Prediction error =  $g(\text{Bias}, \text{Variance})$

# Why machine learning is powerful?

Prediction error =  $g(\text{Bias}, \text{Variance})$

**ML is able to find a good balance  
between bias and variance**

<b>Several « ML tricks »</b>	<b>Principle</b>	<b>Effect</b>
Regularization	Add information to prevent overfitting and simplify the model	Reduce variance at the cost of a small increase of bias
Bagging	Bootstrap aggregation: average together multiple models fitted to resampled dataset	Reduce variance
Boosting	Fit a sequence of weak models to weighted versions of the data (more weight given to poorly predicted data at earlier rounds).	Reduce bias

Numerous methods available

- Regressions (standard, PLS, LASSO, Elastic net...)
- SVM
- Tree and random forest
- Gradient boosting
- Neural network
- Deep neural network
- Deep learning
- Bayesian classification

## Numerous methods available

- Regressions (standard, PLS, LASSO, Elastic net...)
- SVM
- Tree and random forest
- Gradient boosting
- Neural network
- Deep neural network
- Deep learning
- Bayesian classification

**Relatively easy to run these methods with specialized packages (with R or Python)**

# Are machine learning models « black boxes »?

This is less true than before.

Vizualisation tools:

- Importance ranking
- Partial dependence plots (PDP)
- Accumulated Local Effects (ALE) Plot

# Example 1

RESEARCH ARTICLE

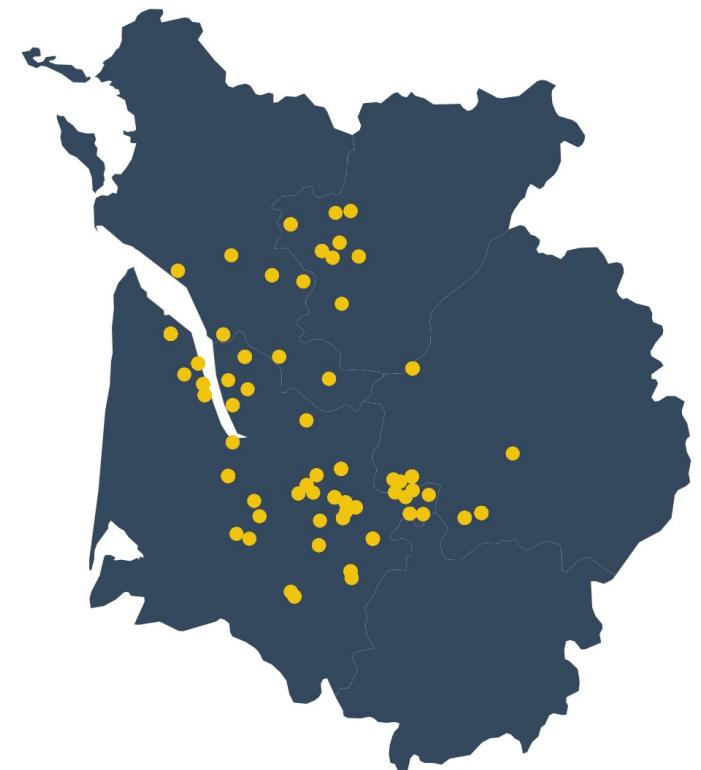
## Forecasting severe grape downy mildew attacks using machine learning

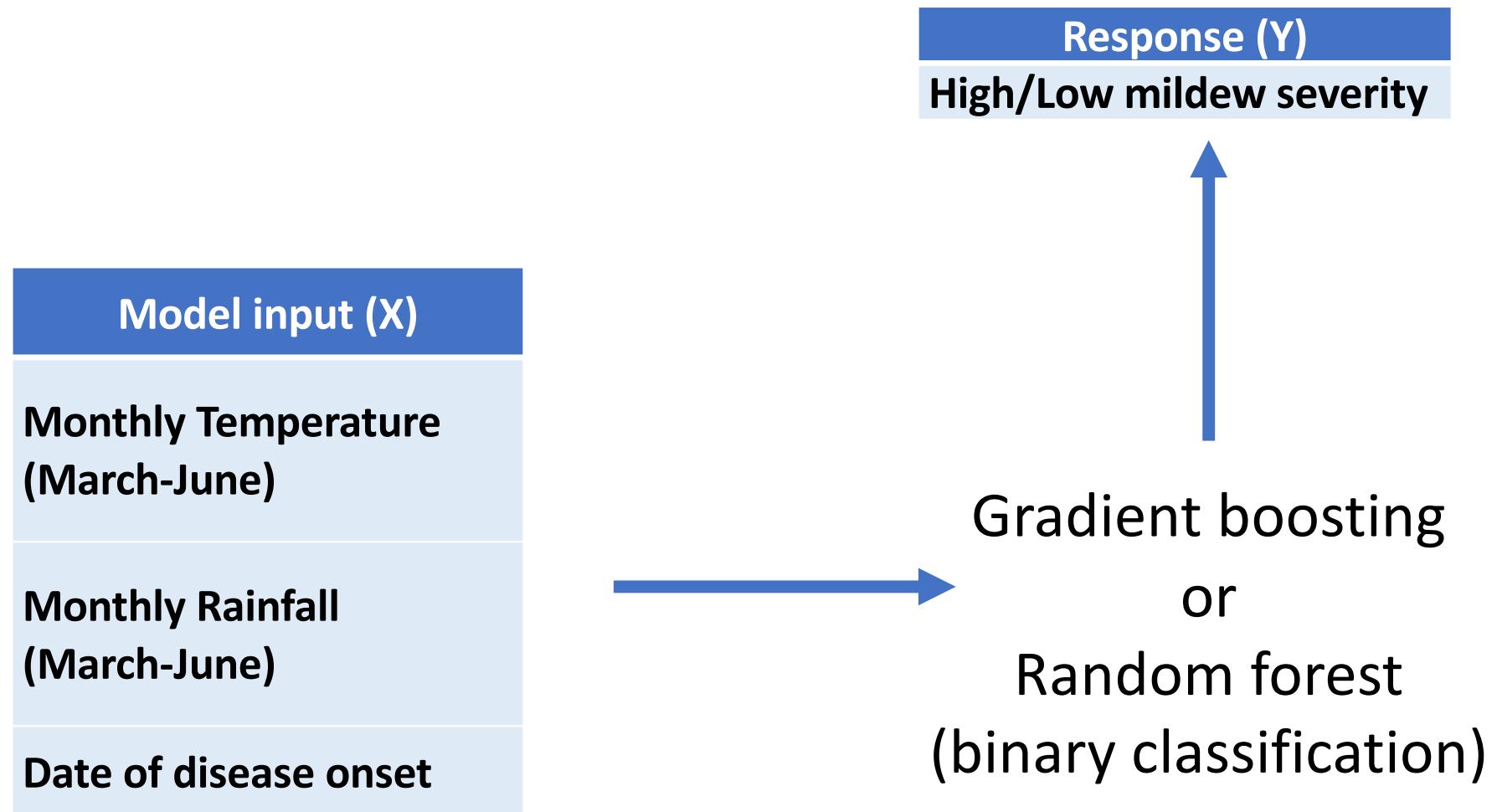
Mathilde Chen<sup>1\*</sup>, François Brun<sup>2</sup>, Marc Raynal<sup>3</sup>, David Makowski<sup>4,5</sup>

<https://doi.org/10.1371/journal.pone.0230254>

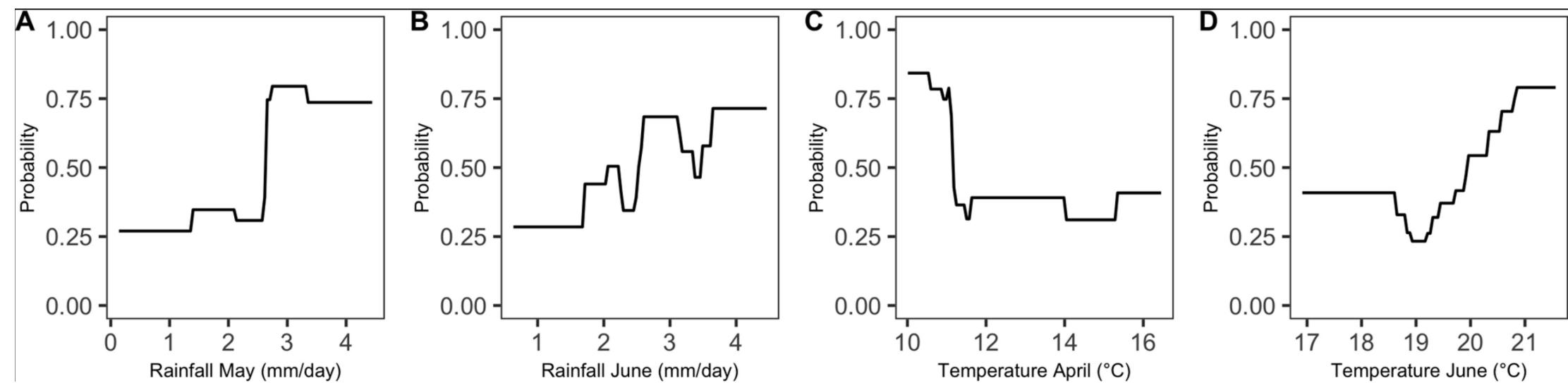
# Forecasting severe grape downy mildew attacks using machine learning

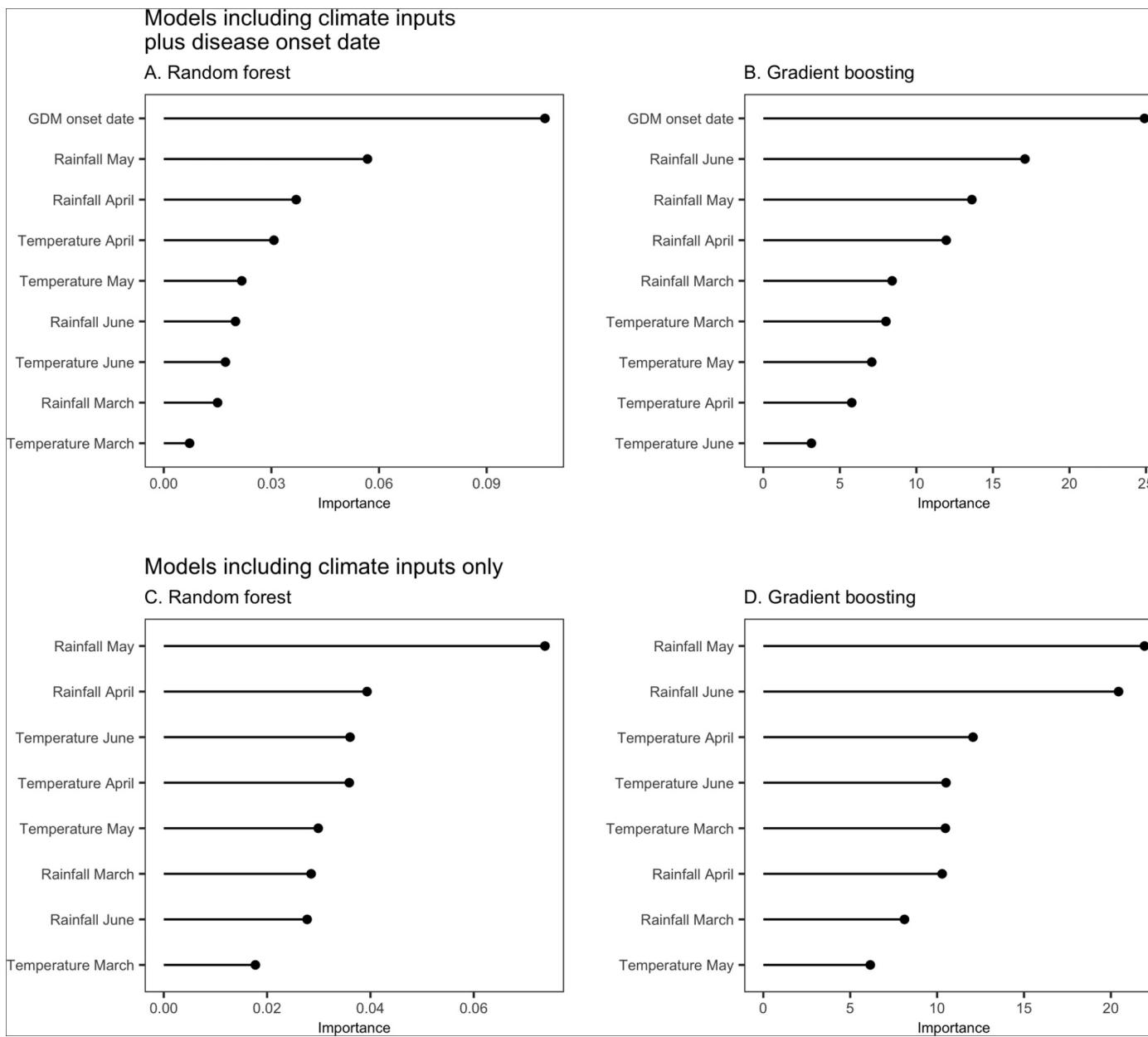
Incidence data collected on  
266 untreated plots between 2010 and  
2017





## Probability of high severity on leaves as a function of climate inputs





« Results show that the use of this forecasting tool would reduce the number of treatments against grape downy mildew in the Bordeaux vineyards compared to current practices **by at least 50%**. »

# Example 2: Prediction of root biomass

<https://doi.org/10.5194/essd-2021-25>  
Preprint. Discussion started: 29 March 2021  
© Author(s) 2021. CC BY 4.0 License.



Open Access  
Earth System  
Science  
Data  
Discussions

## A global map of root biomass across the world's forests

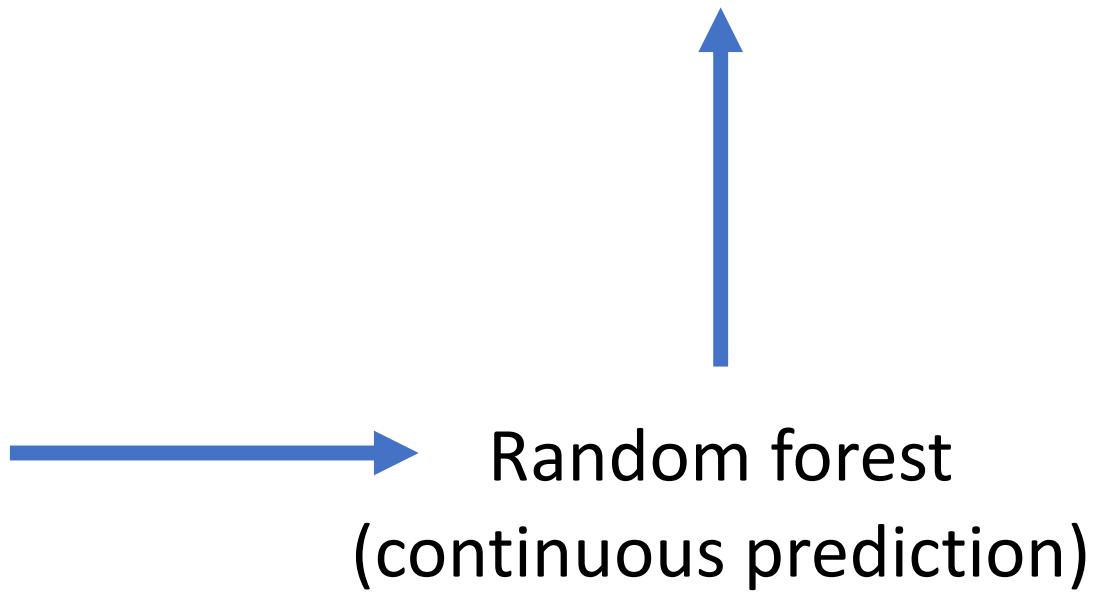
Yuanyuan Huang<sup>1,2</sup>, Phillippe Ciais<sup>1</sup>, Maurizio Santoro<sup>3</sup>, David Makowski<sup>4,5</sup>, Jerome Chave<sup>6</sup>, Dmitry Schepaschenko<sup>7,8,9</sup>, Rose Z. Abramoff<sup>1</sup>, Daniel S. Goll<sup>1</sup>, Hui Yang<sup>1</sup>, Ye Chen<sup>10</sup>, Wei Wei<sup>11</sup>, Shilong Piao<sup>12,13,14</sup>

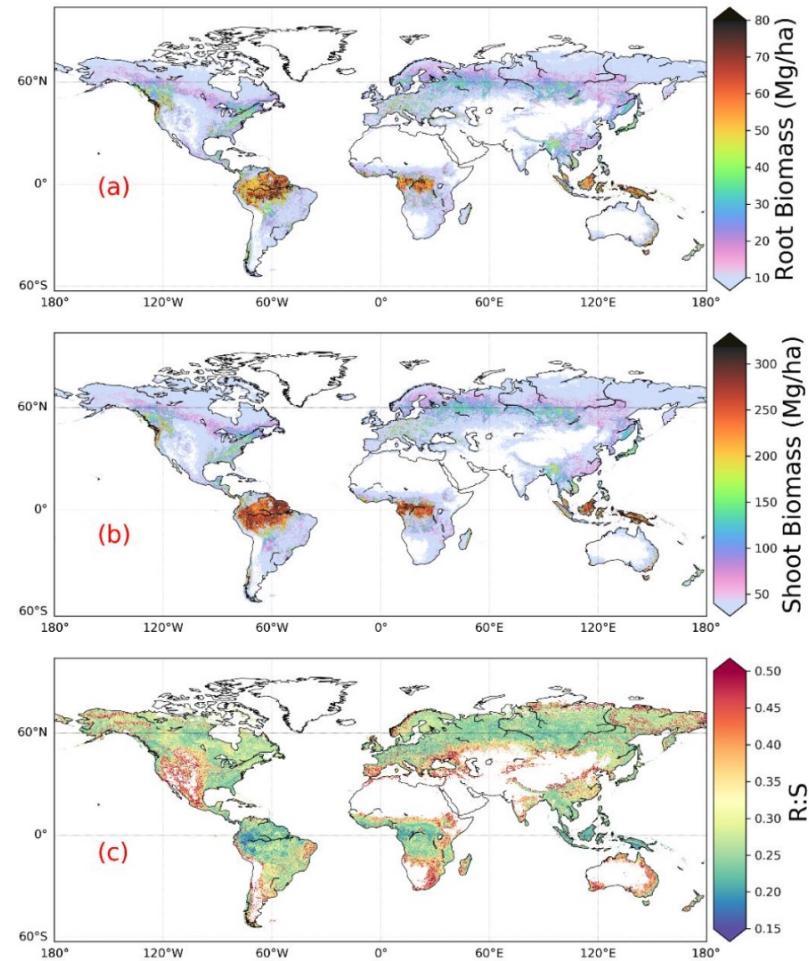
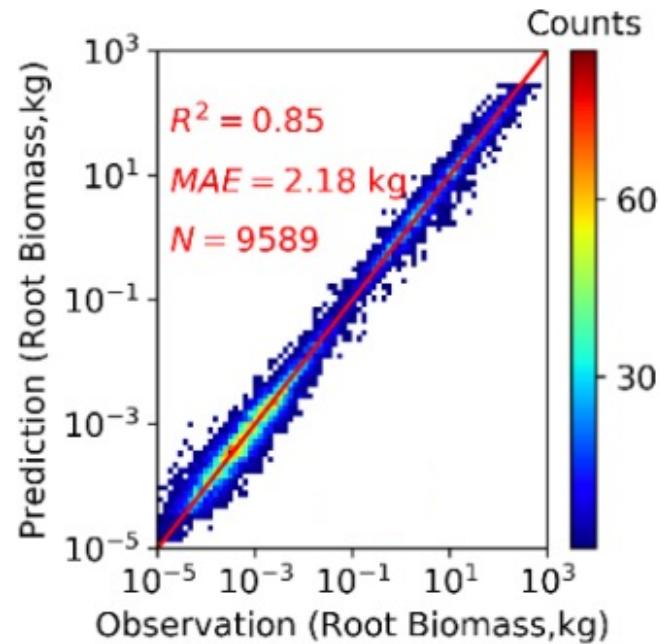
Dataset:

10,307 in-situ measurements of the biomass of roots and shoots for individual woody plants, covering 465 species across 10 biomes.

47 model input (X)
Shoot biomass
Height
Age
Species
Soil bulk density
Soil organic C
pH
Sand content
Clay content
Total N
...

Response (Y)
Root biomass





Global maps of **forest root biomass** generated through a machine learning model (a), shoot biomass from GlobBiomass-AGB(Santoro, 2018b) (b) and Root:Shoot ratio (c).

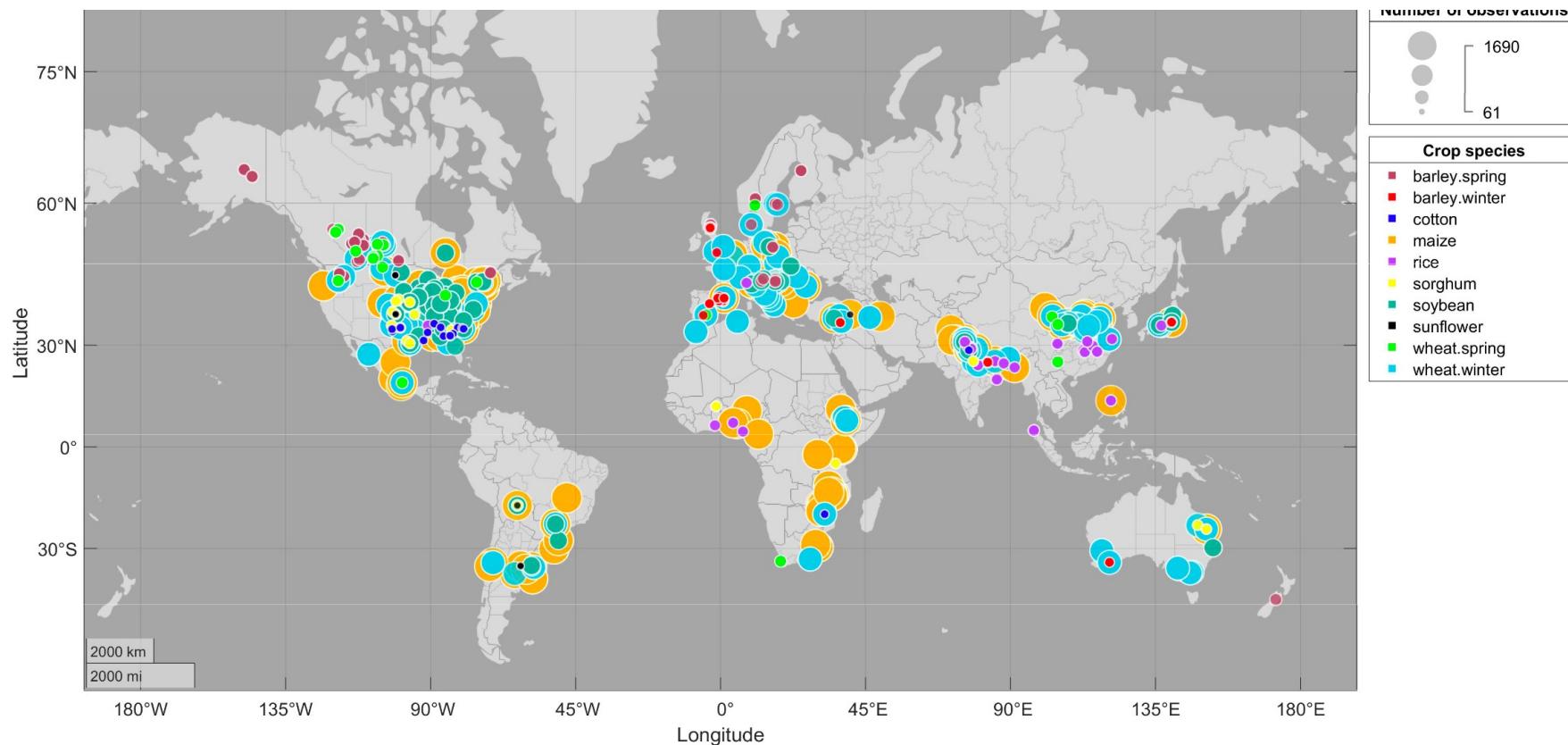
Example 3: Map the probability of yield increase of converting “conventional tillage system” to “conservation agriculture” at the global scale

<https://www.nature.com/articles/s41598-021-82375-1.pdf>

# Locations of the experiments included in the dataset

Each experiment includes yield data for

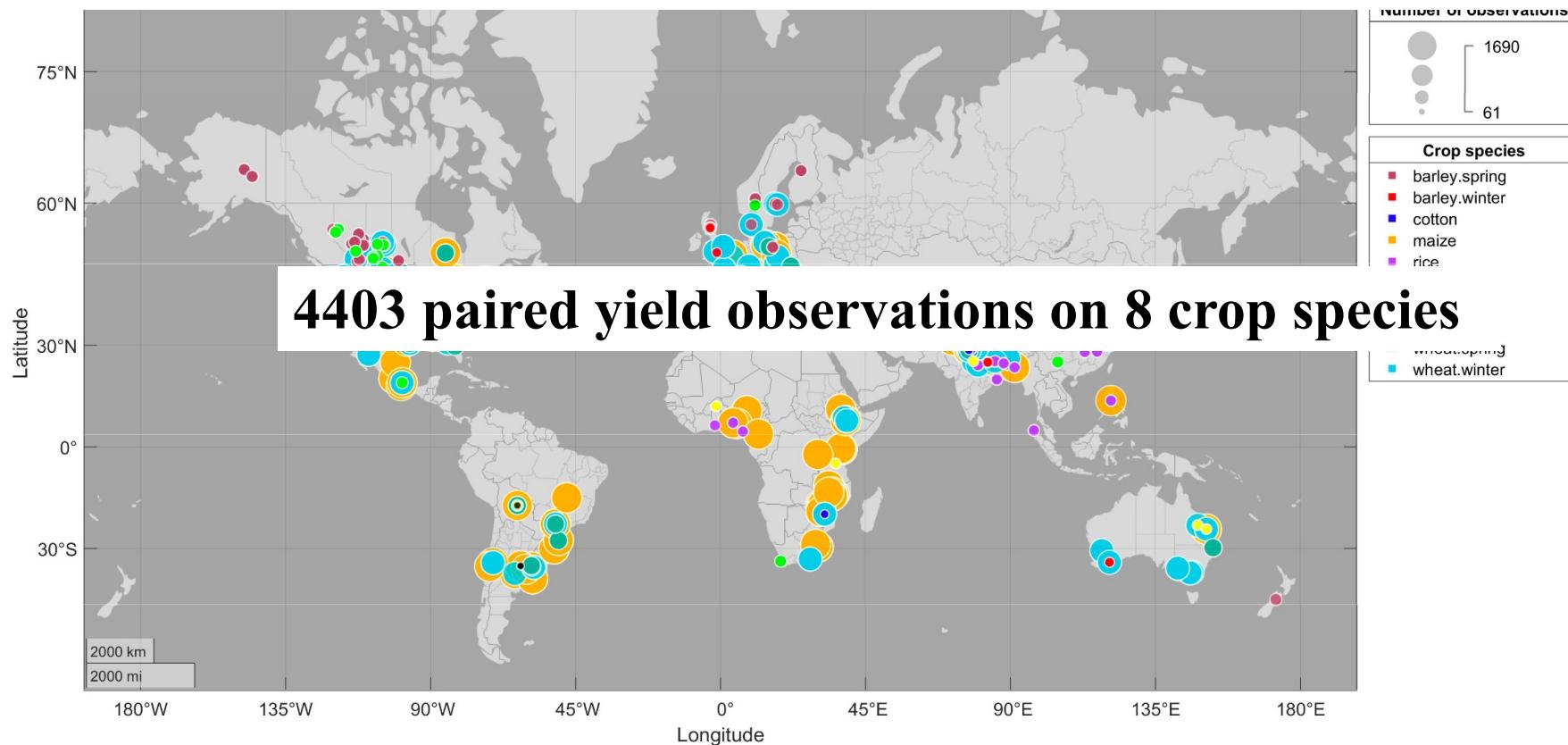
- a conservation agriculture system
- a conventional tillage system



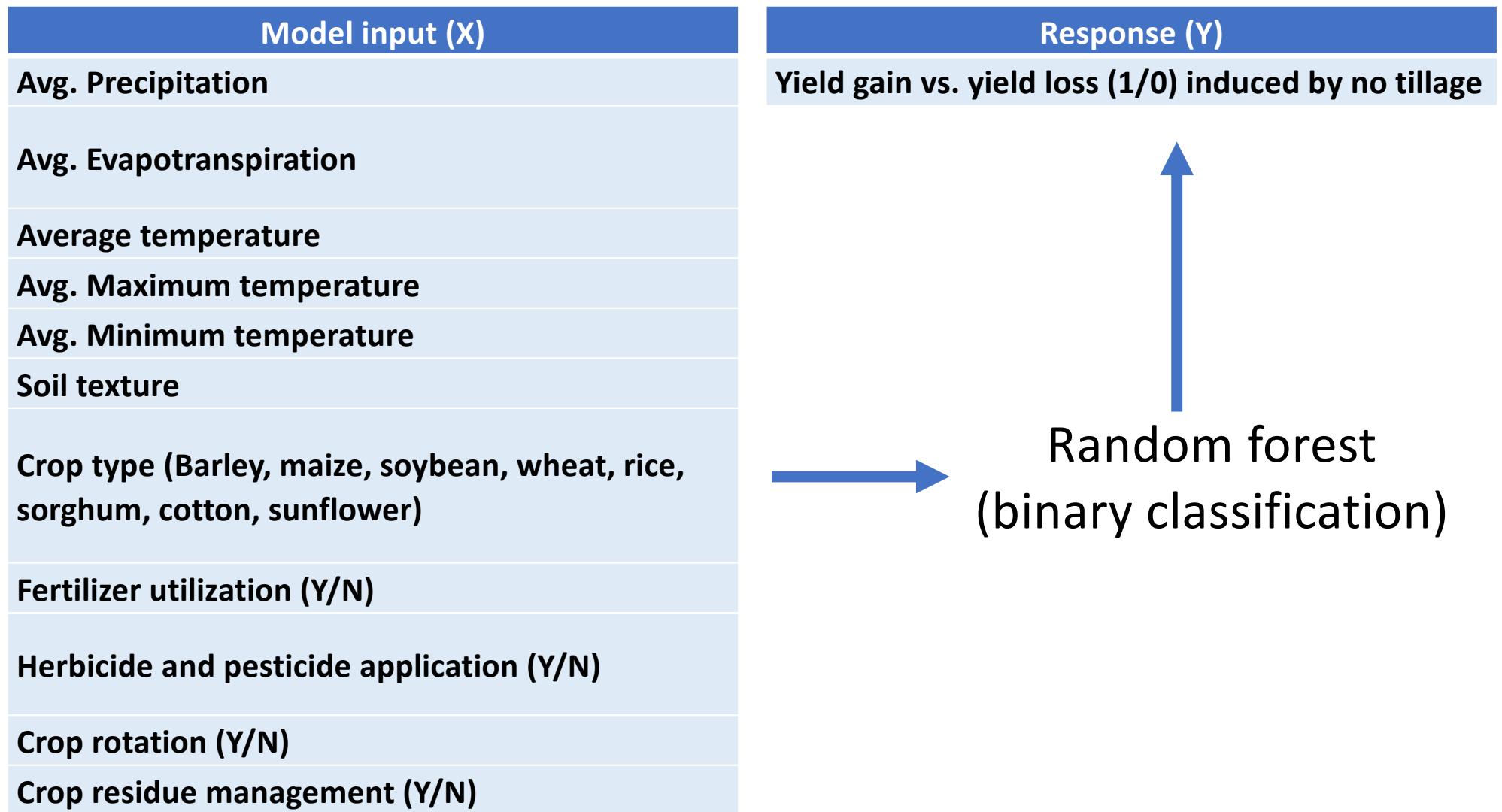
# Locations of the experiments included in the dataset

Each experiment includes yield data for

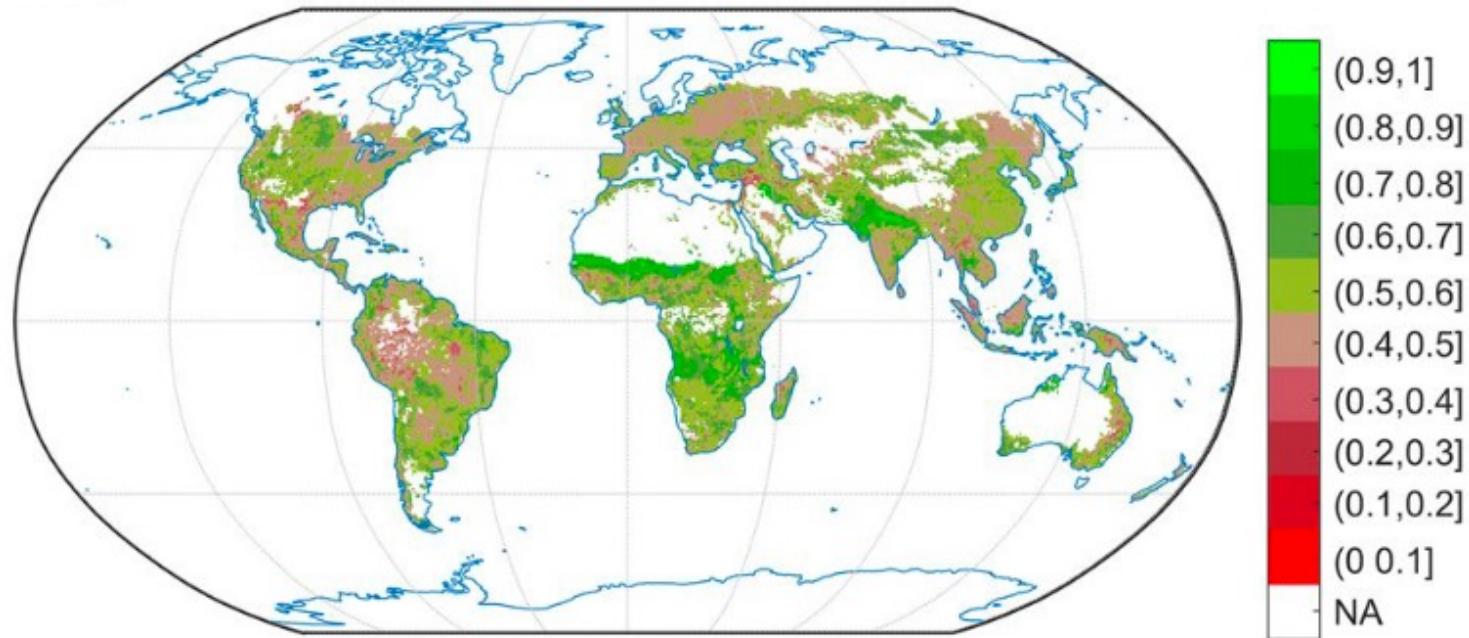
- a conservation agriculture system
- a conventional tillage system



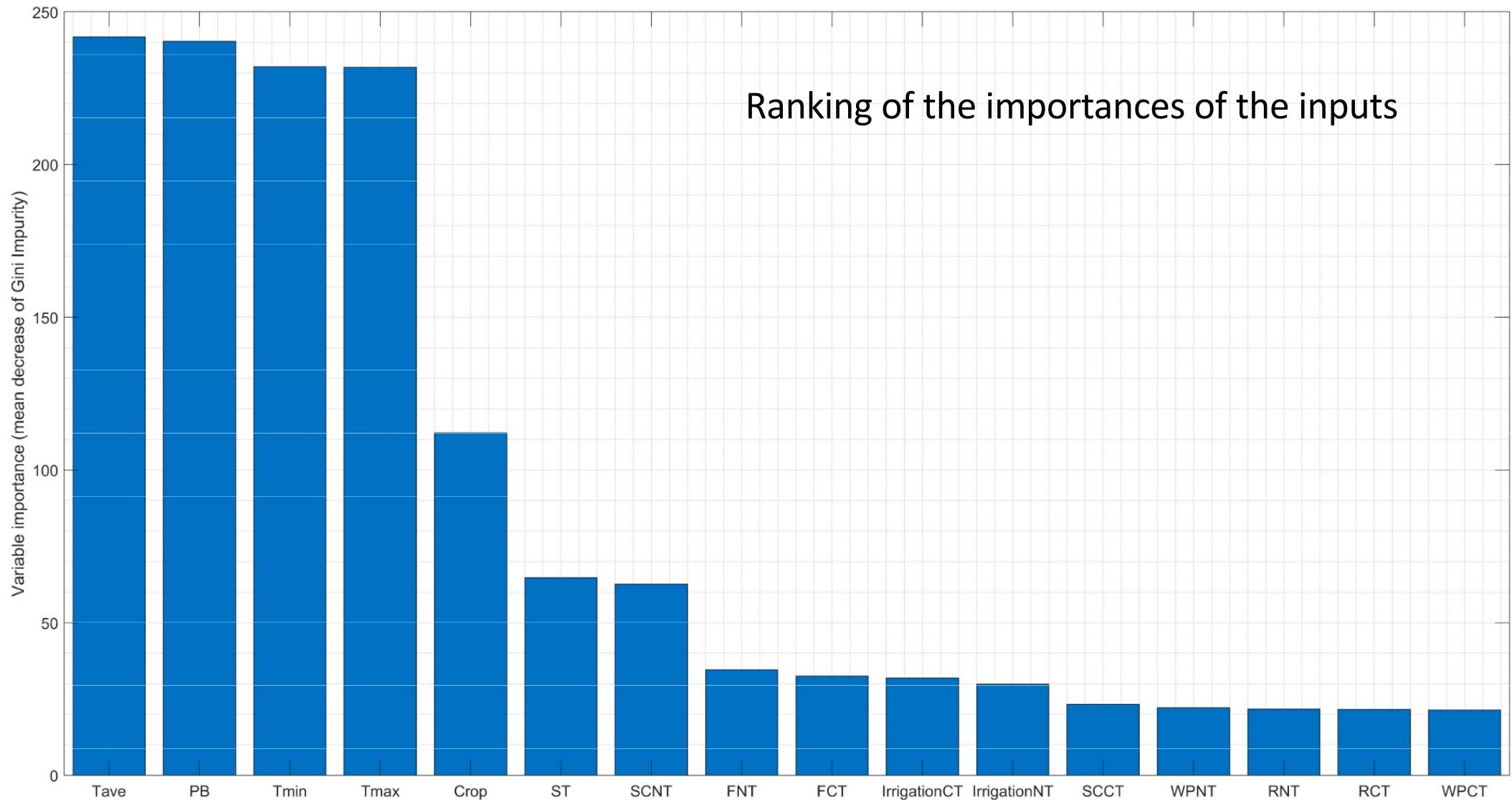
Model input (X)	Response (Y)
Avg. Precipitation	Yield gain vs. yield loss (1/0) induced by no tillage
Avg. Evapotranspiration	
Average temperature	
Avg. Maximum temperature	
Avg. Minimum temperature	
Soil texture	
Crop type (Barley, maize, soybean, wheat, rice, sorghum, cotton, sunflower)	
Fertilizer utilization (Y/N)	
Herbicide and pesticide application (Y/N)	
Crop rotation (Y/N)	
Crop residue management (Y/N)	



## probability of yield increase for maize with Conservation agriculture vs. Tillage



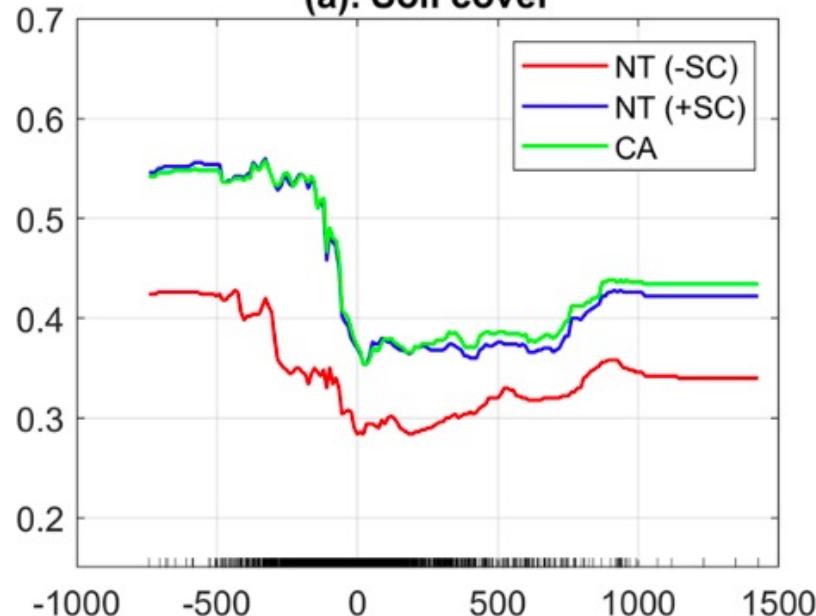
<https://www.nature.com/articles/s41598-021-82375-1.pdf>



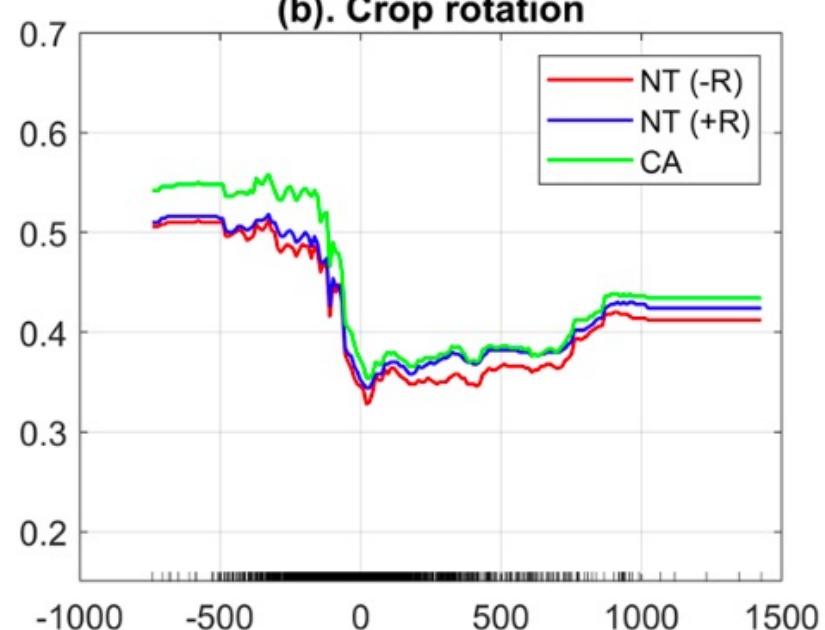
Probability of yield gain with CA or NT vs Tillage system

## 1D-partial dependence plot

(a). Soil cover



(b). Crop rotation



NT: No tillage system

R: Rotation

SC: Soil cover

CA: Conservation agriculture (NT+R+SC)

# Main challenges in machine learning projects

- Choose a relevant question (Which Y? Which X?)
- Find reliable data
- Calibrate the hyper-parameters
- Assess prediction accuracy without bias
- Optimize computation time
- Vizualisation of output responses

# Start simple

Start with two simple methods:

- Penalized linear regression (ex: LASSO)
- Random forest

# Some trends

- Visualization tools (to open « the black boxes »)
- Image and text analyses (text mining, deep learning)
- Packages to streamline the development of predictive models (keras, caret, H2O...)
- Including expert knowledge in machine learning

# Machine learning to emulate complexe models



Contents lists available at [ScienceDirect](#)

Journal of Environmental Management

journal homepage: [www.elsevier.com/locate/jenvman](http://www.elsevier.com/locate/jenvman)

Research article

Meta-modeling methods for estimating ammonia volatilization from nitrogen fertilizer and manure applications

Maharavo Marie Julie Ramanantenasoa<sup>a,b</sup>, Sophie Génermont<sup>a,\*</sup>, Jean-Marc Gilliot<sup>a</sup>,  
Carole Bedos<sup>a</sup>, David Makowski<sup>c,d</sup>

ACCEPTED MANUSCRIPT • OPEN ACCESS

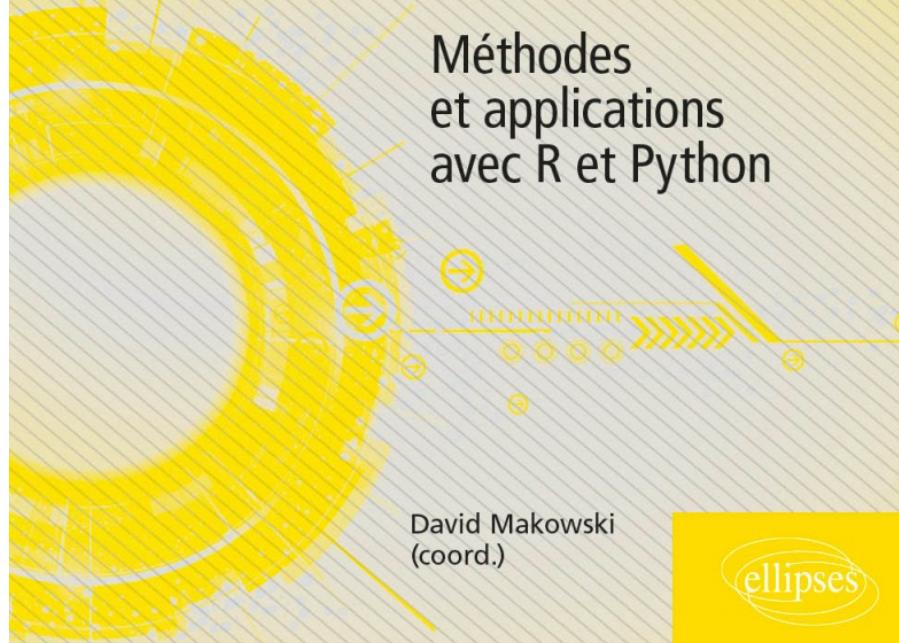
Maize yield and nitrate loss prediction with machine learning algorithms

Mohsen Shahhosseini<sup>1</sup>, Rafael A Martinez-Feria<sup>2</sup> , Guiping Hu<sup>3</sup> and Sotirios Archontoulis<sup>1</sup>  
Accepted Manuscript online 29 October 2019 • © 2019 The Author(s). Published by IOP Publishing Ltd

Formations & Techniques

# DATA SCIENCE POUR L'AGRICULTURE ET L'ENVIRONNEMENT

Méthodes  
et applications  
avec R et Python



David Makowski  
(coord.)

ellipses