

Assembler la diversité des modèles classiques et «deep learning» pour développer un pipeline de calibration SPIR performant et générique.

Cornet D, Beurier G, Cormier F, Meghar K, Davrieux F, Arnau G (Cirad)

Desfontaines L, Marie-Magdeleine C (INRA)

Rencontres HelioSPIR, Session 5 - Place aux « jeunes » !

15 octobre 2019

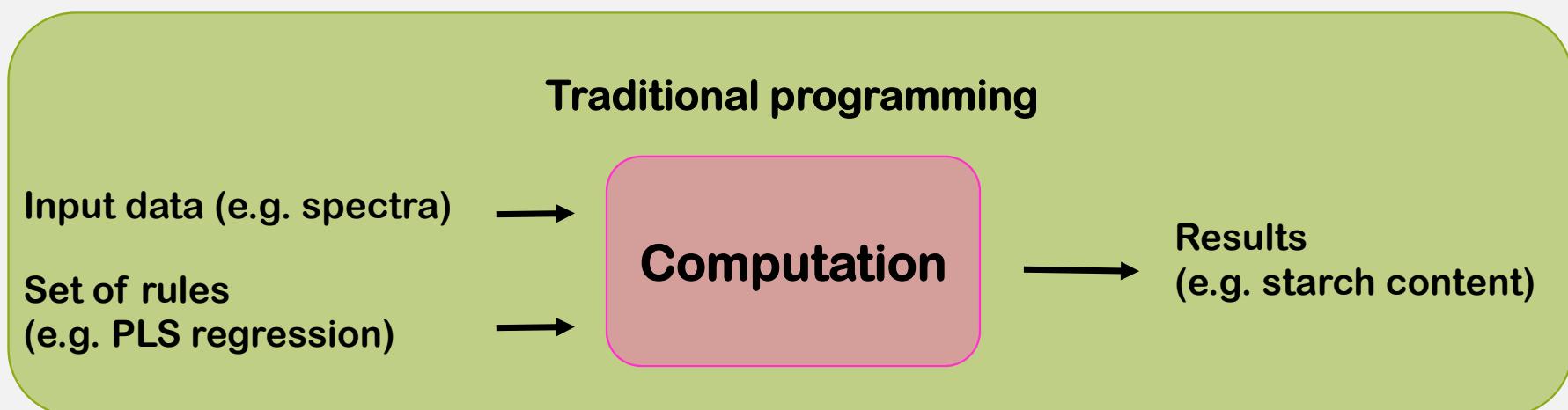
Contenu

- Justification des choix méthodologique
 - Apprentissage automatisé
 - Pipeline générique
 - Assemblage de modèle
- Assemblage de modèles
- Analysis pipeline under development
- Perspectives

Pourquoi l'apprentissage automatisé ?

- calibration SPIR = développement du meilleur modèle **prédictif**
 - Pas/peu d'intérêt de comprendre pourquoi et comment cela fonctionne
 - > 1000 variables explicatives et taille échantillon << variables explicatives
 - Relations non linéaires

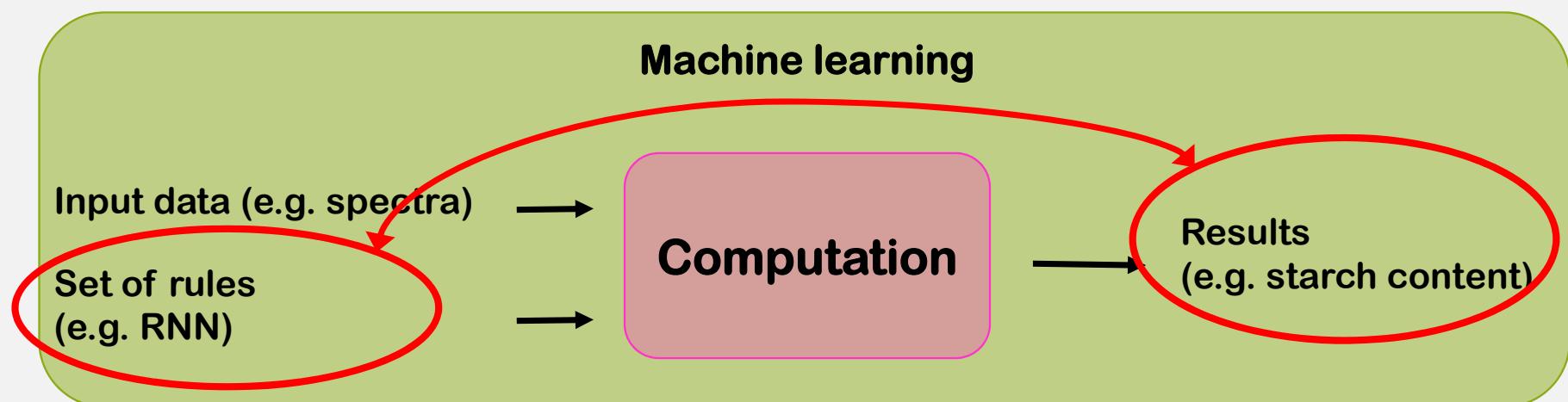
⇒ l'apprentissage automatisé se révèle plus performant



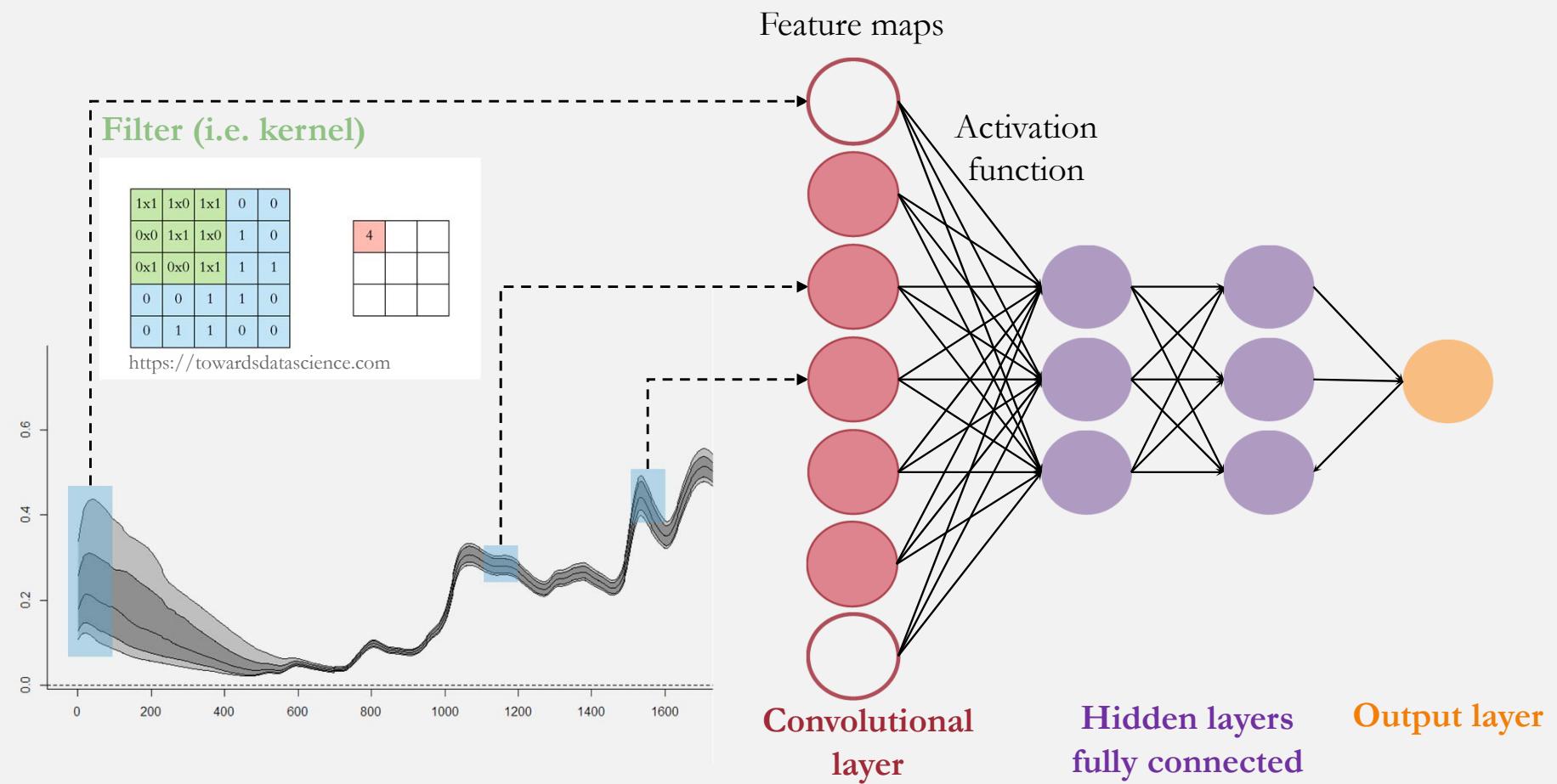
Pourquoi l'apprentissage automatisé ?

- calibration SPIR = développement du meilleur modèle **prédictif**
 - Pas/peu d'intérêt de comprendre pourquoi et comment cela fonctionne
 - > 1000 variables explicatives et taille échantillon << variables explicatives
 - Relations non linéaires

⇒ l'apprentissage automatisé se révèle plus performant

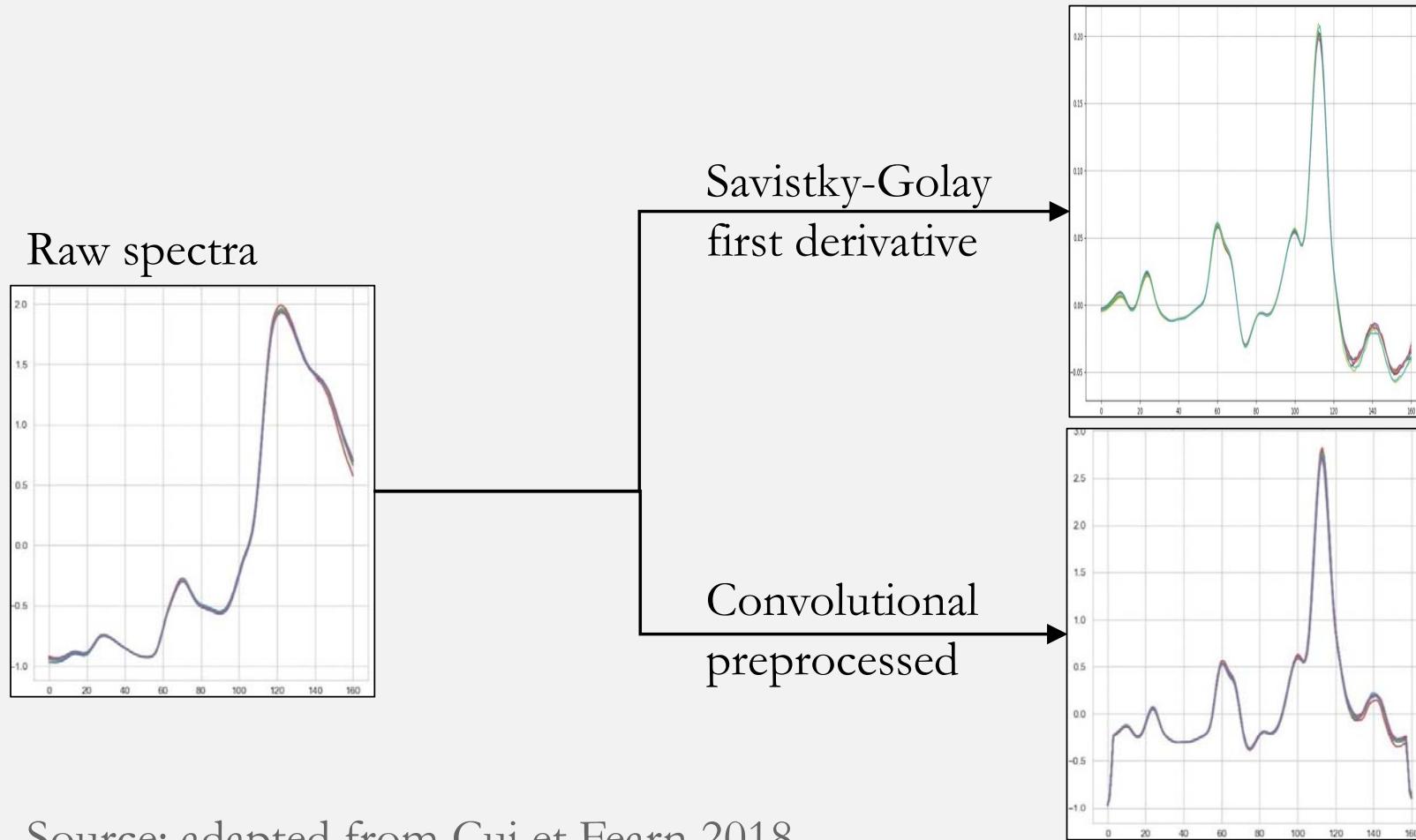


Exemple d'application Convolutional network



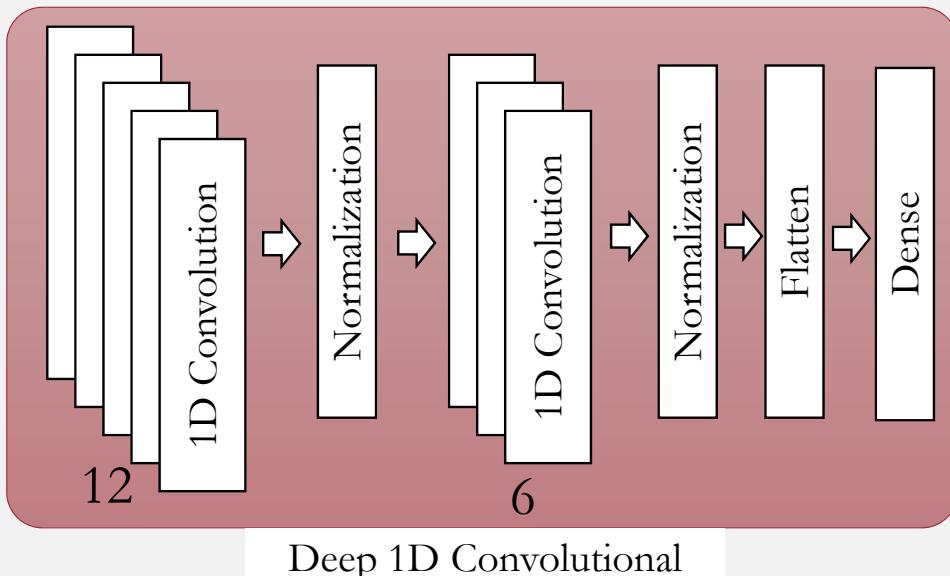
Exemple d'application

Convolutional network

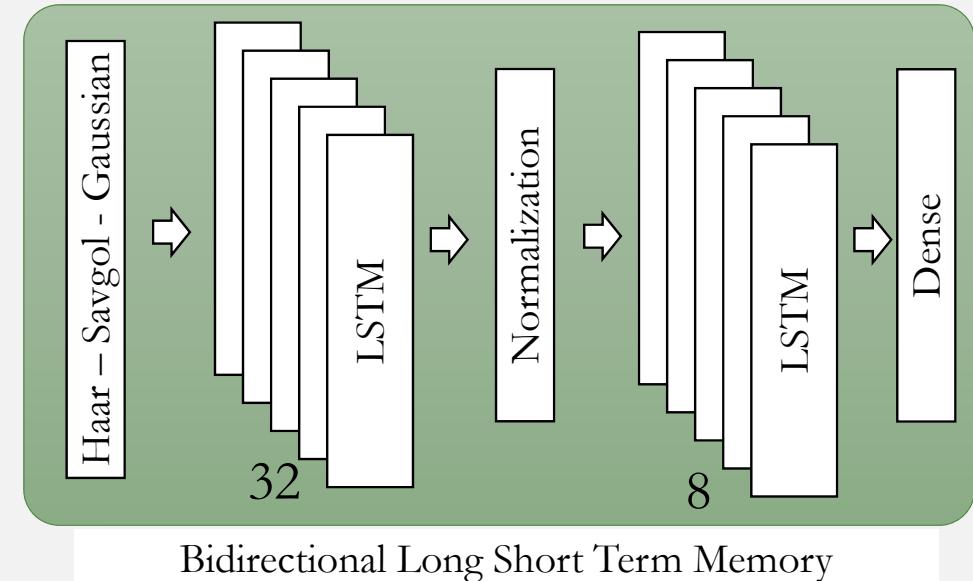


Source: adapted from Cui et Fearn 2018

Exemple d'application : texture et caisson du manioc



KFold – 5 Fold 4 Repeats
1500 epochs – **No Hyperparametrization**
40 % Dropout on convolutions



KFold – 5 Fold 4 Repeats
800 epochs – **No Hyperparametrization**
50 % Dropout on input

Cooking Time
Validation Set MAE 19.49 % \pm 2.77%
Global Set MAE 10.92 % \pm 0.82%

Texture
Validation Set MAE 13.39 % \pm 2.91%
Global Set MAE 7.67 % \pm 1.74%

Cooking Time
Validation Set MAE 20.32 % \pm 1.85%
Global Set MAE 20.74 % \pm 2.26%

Texture
Validation Set MAE 10.31 % \pm 0.60%
Global Set MAE 11.93 % \pm 2.14%

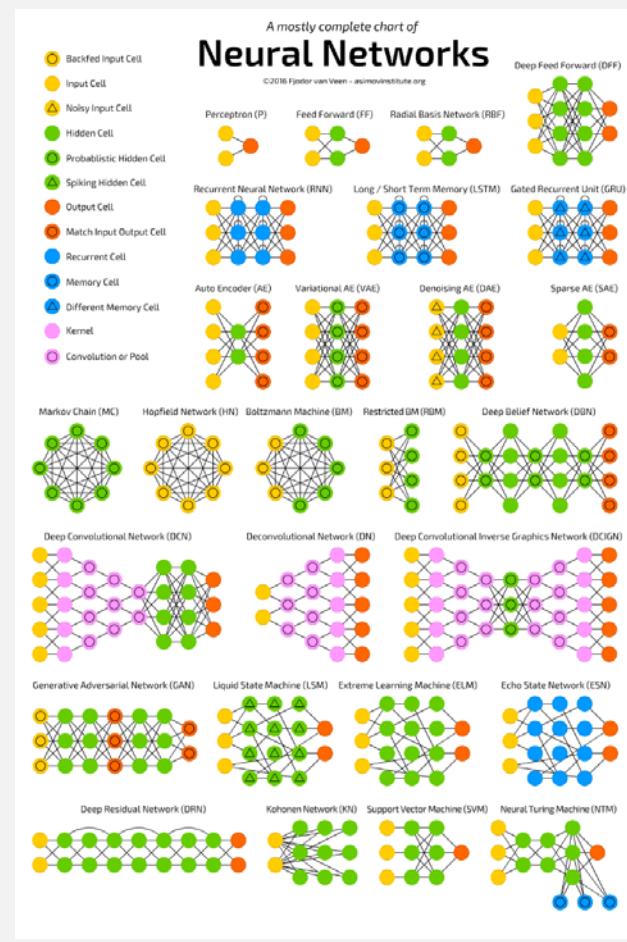
Etude préliminaire avec plusieurs machines

KFold – 5 Fold 4 Repeats
 1500 epochs – No Hyperparametrization
 40 % Dropout on convolutions

Training Set	CNN Spectro 1	CNN Spectro 2	BiLSTM Spectro 1	BiLSTM Spectro 2
Spectro 1 & Spectro 2	7.57% \pm 0.19%	13.74% \pm 0.45%	7.95% \pm 0.49%	9.95% \pm 0.66%
Spectro 1 only	7.16% \pm 1.05%	25.81% \pm 1.51%	6.30% \pm 0.73%	26.68% \pm 2.76%
Spectro 2 only	17.40% \pm 0.44 %	2.93% \pm 0.20%	13.28% \pm 0.37%	4.44% \pm 0.12%
Spectro 1 then Spectro 2	7.28% \pm 0.40%	4.45% \pm 0.39%	6.64% \pm 1.27%	3.87% \pm 0.36%
Spectrums Auto-Encoding then Spectro	12.28% \pm 1.48%	16.57% \pm 1.03%		

Pourquoi un pipeline générique ?

- Démocratisation/vulgarisation de la SPIR apporte de + en + d'utilisateurs naïfs
- La majorité des publications s'attachent à démontrer la supériorité d'une méthode
 - Applique 1 stratégie de modélisation
 - Utilise 1 combinaison spécifique de prétraitement
 - Optimise la calibration pour 1 analyte
- L'utilisateur
 - Dispose souvent de plusieurs analytes (e.g. sucre, amidon, protéine)
 - Se heurte à un mur grandissant de possibilités de combinaisons de prétraitements/modèles
- Si la combinaison optimale diffère d'une étude à l'autre, la manière de l'identifier peut être généralisé



Pourquoi un pipeline générique ?

Construction : les prétratements

- Noise correction (additive and multiplicative effects)
 - **Multi scatter correction** (remove some physical effects like particle size or Shining Light Reflection)
 - **Savitsky-Golay smoothing and derivative**
- Enhance contrast
 - Discret Wavelet Transform: **Haar transform**
 - **SD feature selection**
- Baseline correction
 - **Detrend**
 - **Continuum removal** (albedo normalization)
 - **Derivatives**

⇒ 616 combinations of 7 pretreatments

Pourquoi un pipeline générique ?

Construction : les modèles

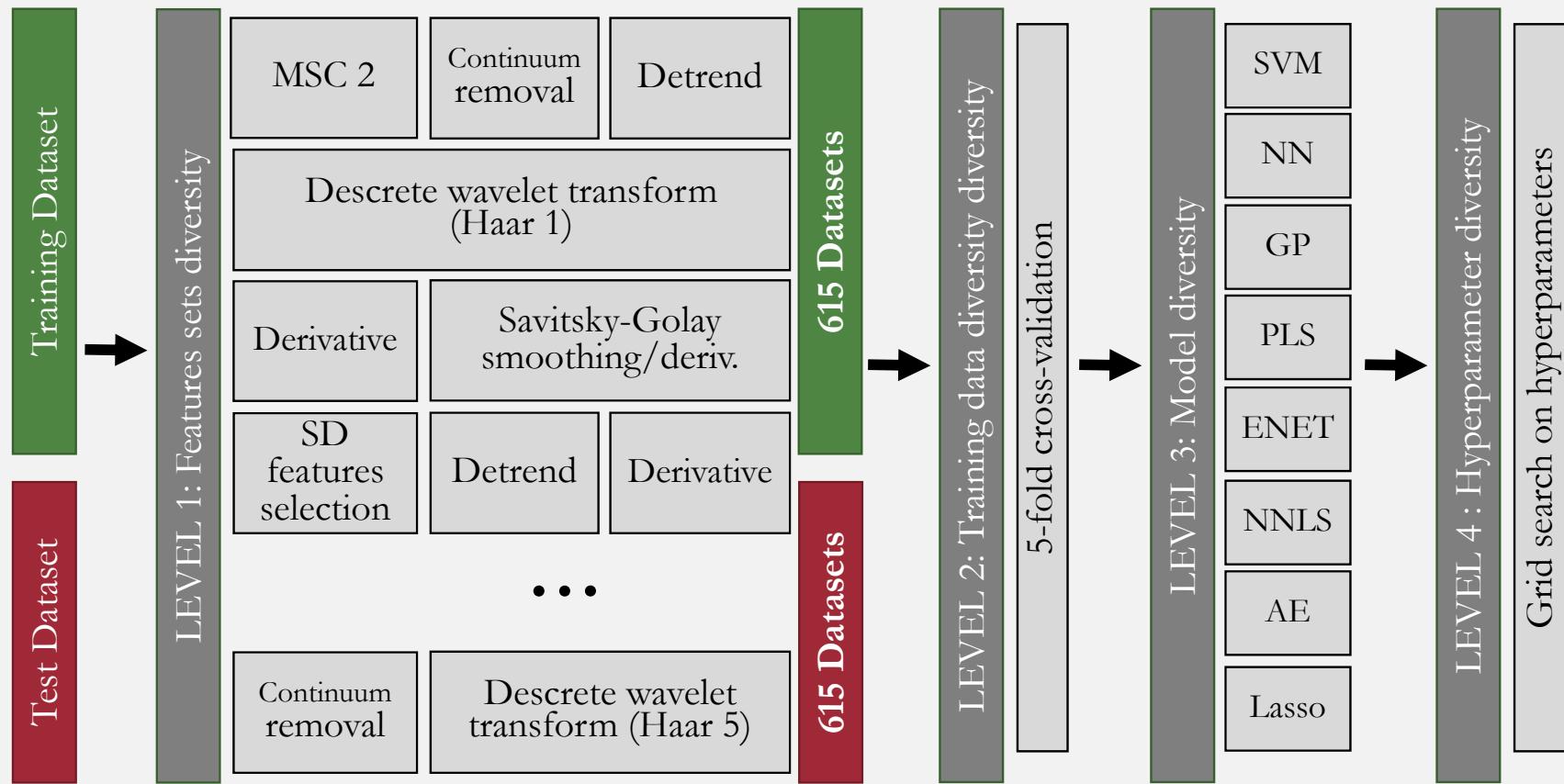
- Traditional programming
 - Partial least squares regression (PLS)
 - Principal component regression (PCR)
 - Non negative least squares (NNLS)
- Machine learning
 - Regularization
 - Lasso regression (least absolute shrinkage and selection operator, LASS)
 - Elastic net regularization (ENET)
 - Bayesian
 - Gaussian process (GP)
 - Support vector machine
 - Linear support vector regression (SVMI)
 - Radial support vector regression (SVMr)
 - Neural networks
 - Neural network (NNET)
 - LSTM, RNN...

Hyperparameter tuning:

- Grid selection
- Random selection

Pipeline d'analyse en cours de construction

Stacking avec R, python, keras et tensorflow



Pourquoi l'assemblage de modèles ?

- Diversité des combinaisons de prétraitements/modèles/analytes
- À un seul modèle
 - Ne bénéficie pas de toute l'information générée pour le sélectionner
 - Risque de s'ajuster au bruit lorsque très/trop spécifique
- Dans le monde singulier des compétitions de données (kaggle), l'immense majorité des équipes gagnantes utilise une combinaison de modèle plutôt qu'un modèle optimal
- L'assemblage de modèles de base au sein d'un métamodèle s'appelle l'assemblage de modèle
- Utilisé pour la première fois en 1992 par Wolpert mais il faut attendre 2007 pour avoir une démonstration mathématique de son avantage (van de Laar *et al.*)

The screenshot shows a competition page for 'Quora Question Pairs' on Kaggle. At the top, it displays the title 'Quora Question Pairs', a brief description 'Can you identify question pairs that have the same intent?', and a '\$25,000 Prize Money'. Below this, the section 'WINNING SOLUTION' is prominently displayed in large blue letters, followed by the text 'KAGGLE QUORA'. A horizontal line separates this from the 'SUMMARY' section. The 'SUMMARY' section contains a numbered list of five items: 1. Introduction, 2. Deep Learning approach, 3. Graphical approach, 4. Ensembling and stacking, and 5. Conclusion. The fourth item, 'Ensembling and stacking', is highlighted with a yellow oval and a blue circle around it.

Maximilien BAUDRY
Featured Prediction Competition
Quora Question Pairs
Can you identify question pairs that have the same intent?
\$25,000 Prize Money

WINNING SOLUTION

KAGGLE QUORA

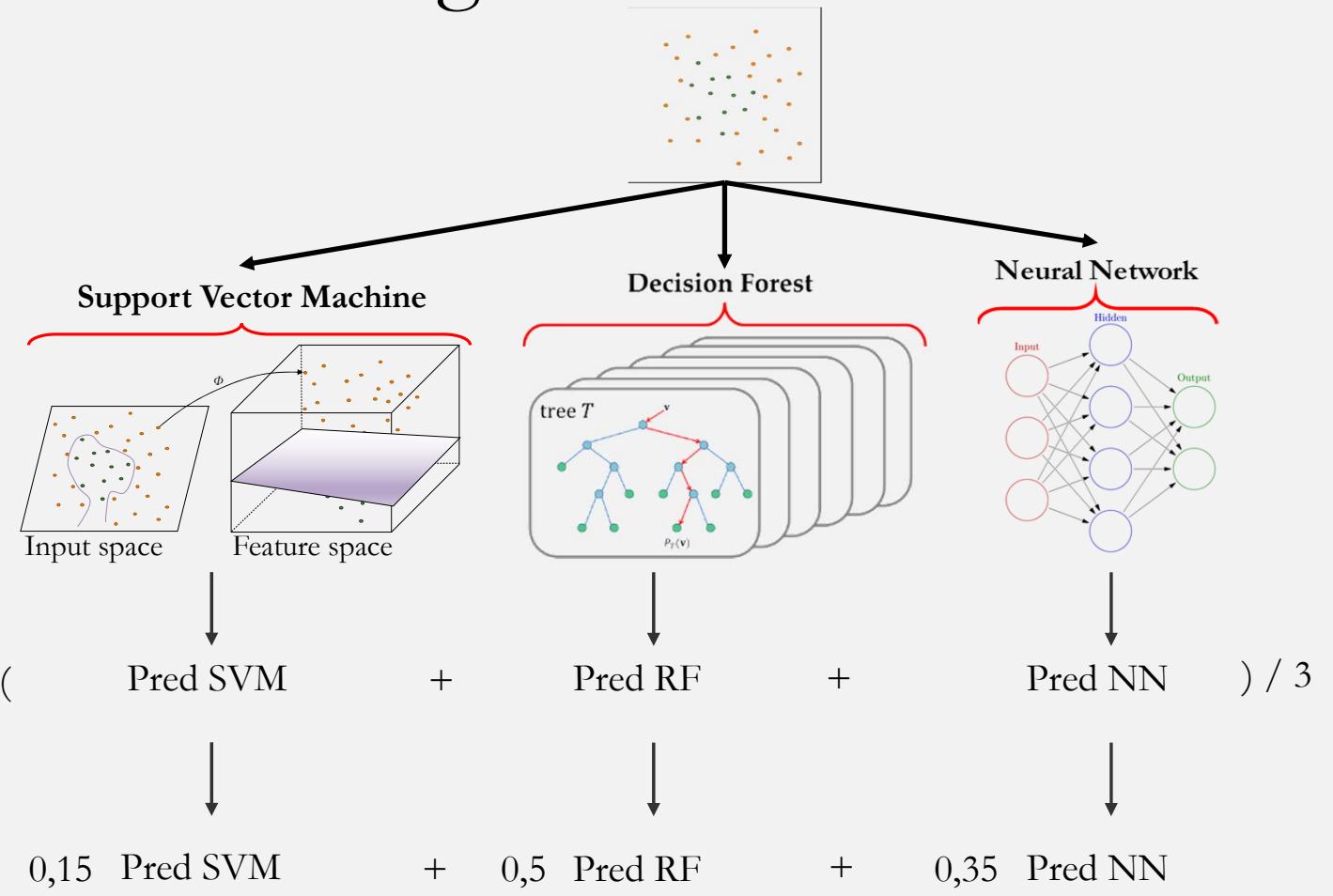
KAGGLE QUORA - WINNING SOLUTION

SUMMARY

1. Introduction
2. Deep Learning approach
3. Graphical approach
4. **Ensembling and stacking**
5. Conclusion

Techniques d'assemblage

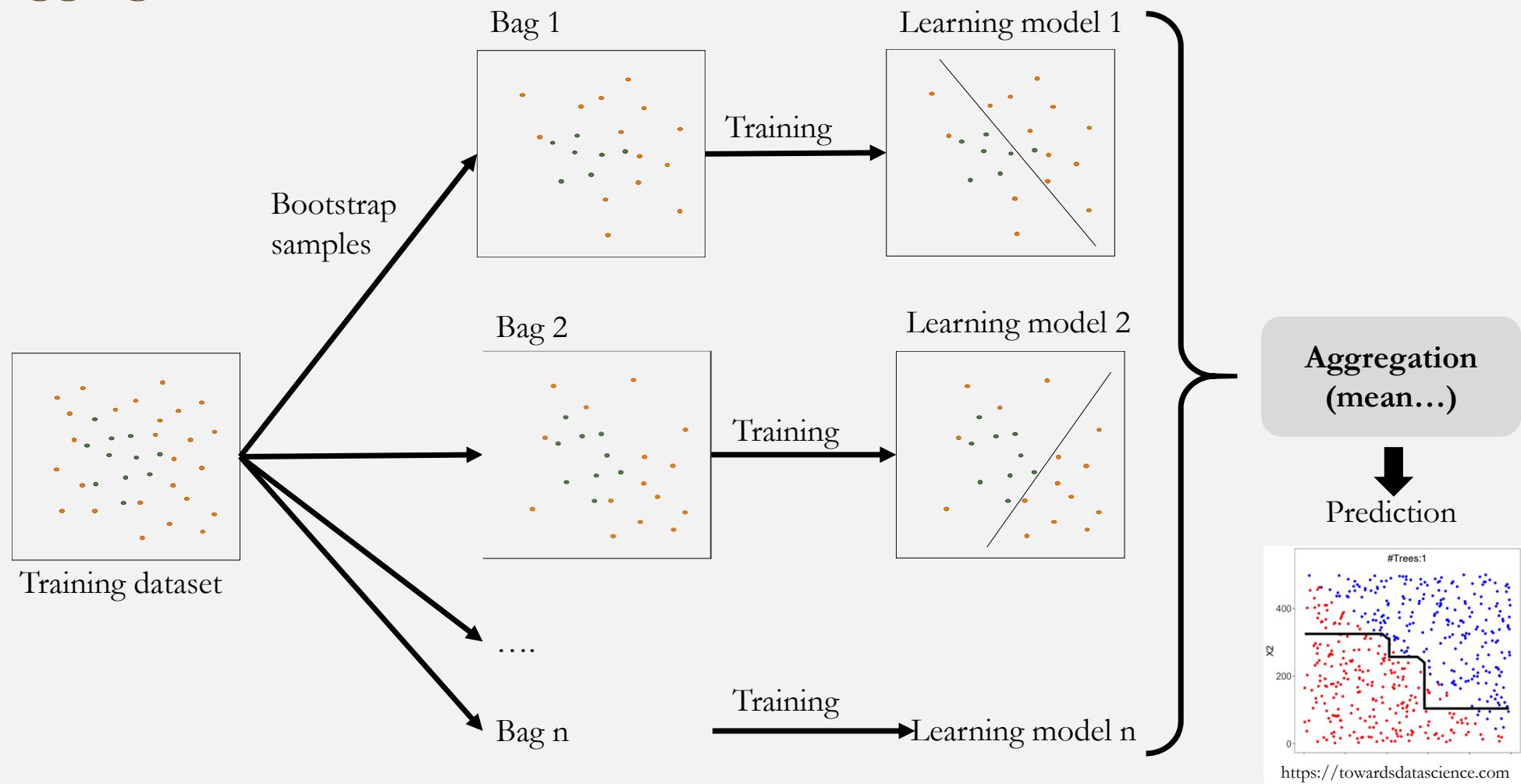
Moyenne



Techniques d'assemblage

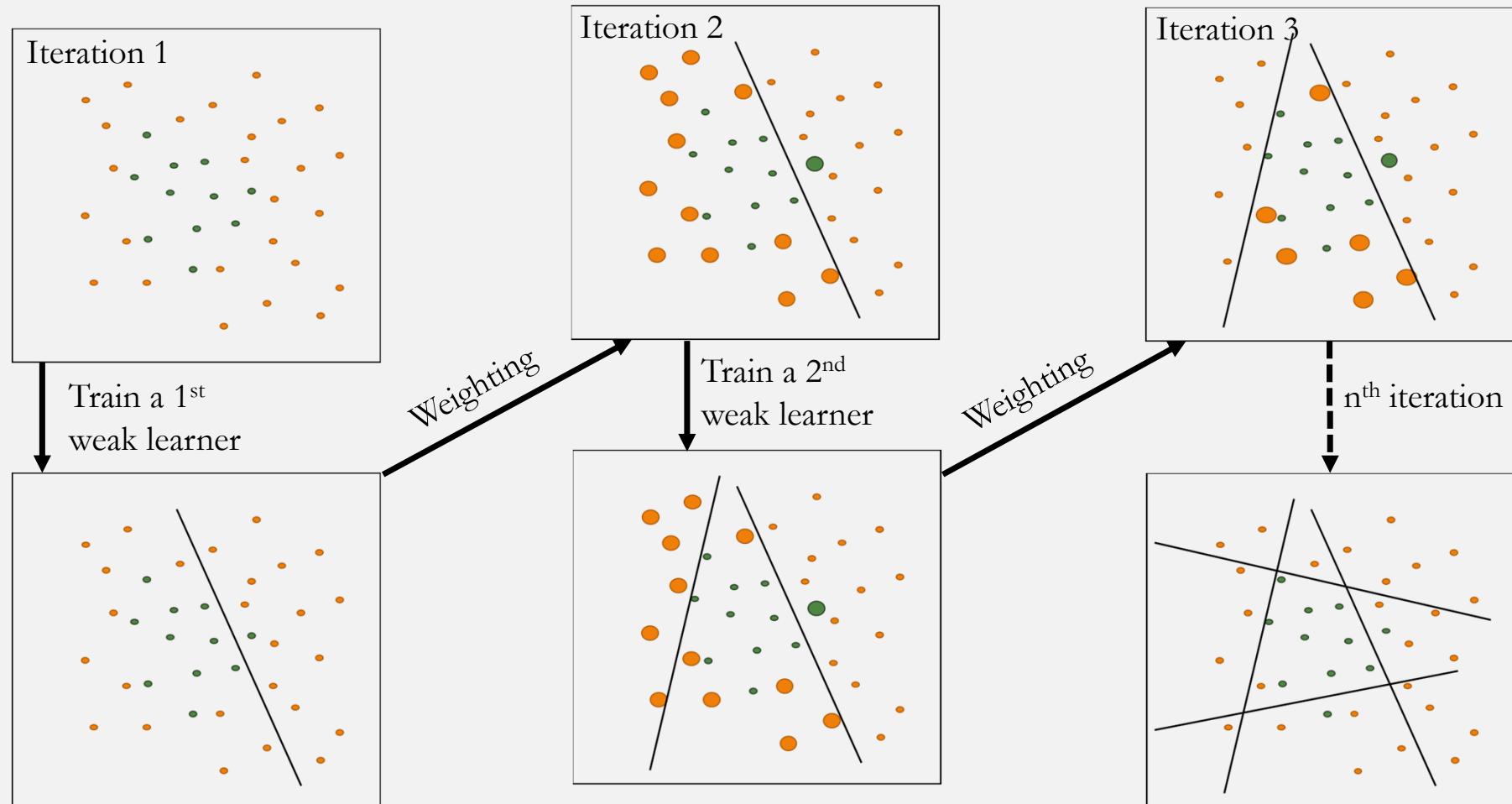
Bagging

Bagging (bootstrap aggregation, e.g. randomforest)



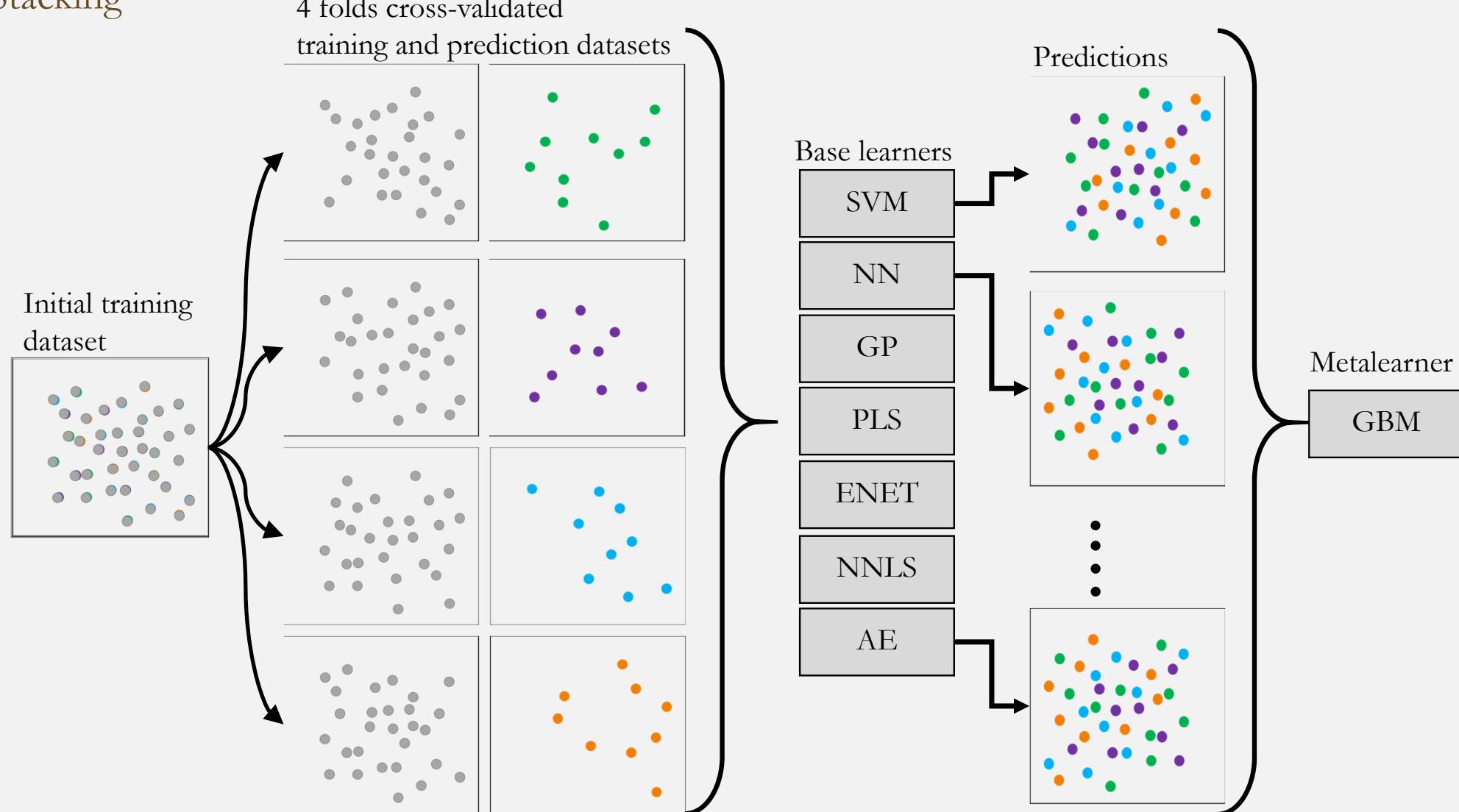
Techniques d'assemblage

Boosting



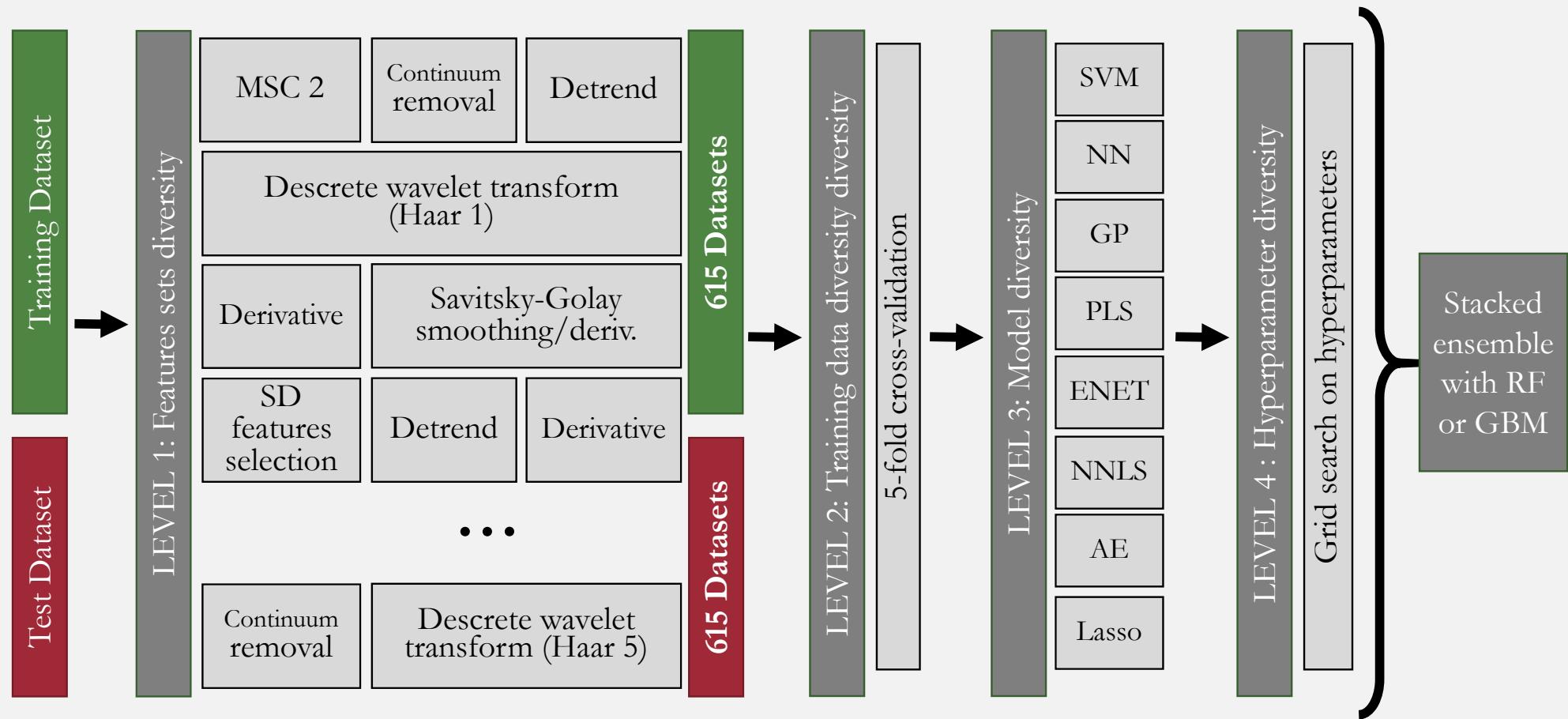
Technique d'assemblage

Stacking



Pipeline d'analyse en cours de construction

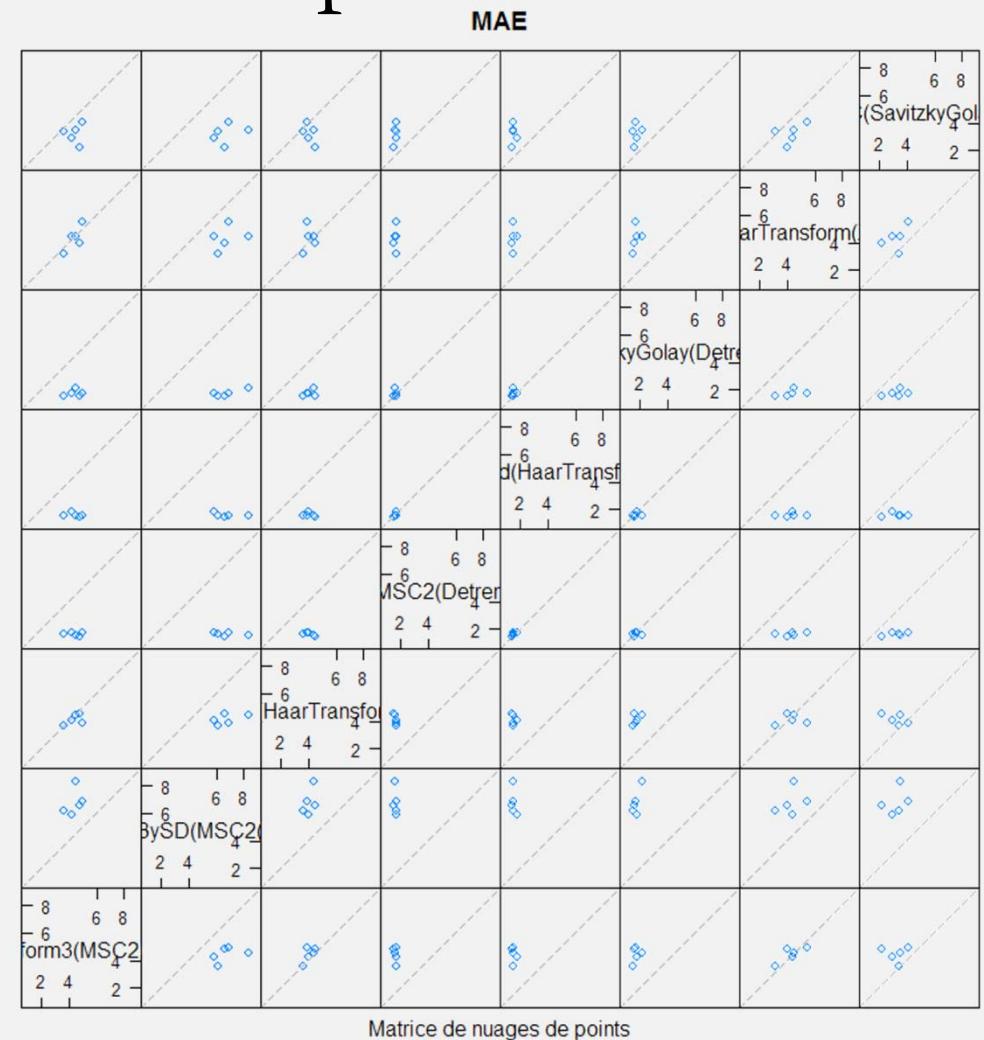
Stacking avec R, python, keras et tensorflow



Analysis pipeline under development

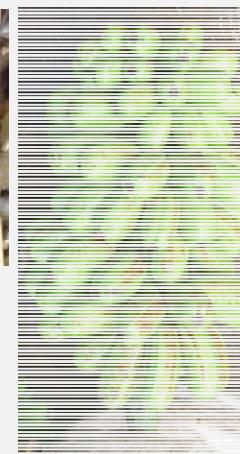
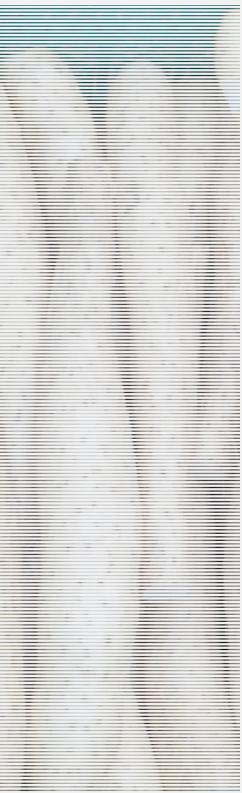
Base model choice

- Conditions
 - Performant
 - ⇒ Strong learners ($MAE < 1$)
 - ⇒ Steady learners (low variation between cross validated samples MAE)
 - Diverse
 - ⇒ Correlation: Pearson < 0.9
 - ⇒ Distribution: Kolmogorov-Smirnov > 0.1
 - 6150 models → 8 base learners



Perspectives

- Finaliser la construction du pipeline
- Jouer avec les résultats d'analyse (performance relative, analyse de sensibilité....)
- Tester la générericité du pipeline (produits, des tailles de jeux de données, specromètres...) => projet RTBfoods
- Tester la plus-value de travailler sur un vecteur d'analytes plutôt qu'un seul analyte à la fois (meilleure gestion des outliers, ...)
- Tester la nouvelle generation d'algorithmes d'apprentissage profond (e.g. RNN, CNN, LSTM, attention based NN)



Merci

