# Machine learning Introduction and use in yield and yield gap decomposition

# Machine learning

- Learning pattern from data

https://observablehq.com/@yizhe-ang/interactive-visualization-of-linear-regression

**Traditional Programming**

Say computer to fit a linear regression or frontier or boundary line based on some assumption about what should be the property of parameters, How residuals should be distributed ?

Data → Computer → Output

Program →

**Machine Learning**

You provide your management and yield to the computer and ask computer how these set of management leads to this yield. How management interacts to affect yield ?

Data → Computer → Program

Output →

By default, the model are black box

Slide credit: Pedro Domingos

There exists some other technique for model explanation !!

# Machine learning not new for you !

**Amazon:**

Recommended for you, bought together, recently viewed, etc. ML clusters the things based on search criteria

**YouTube:** If you regularly sees news in morning, When you open you tube in morning, mostly news comes.
If you search for flights in google, you may get an advertisement related to flight price while watching a video

(Video recommendation come based on recent searches ,clicks, likes and dislikes, watch time, and shares)

**Movie recommendation from Netflix:**

Browsing history and ratings issued, movie type and popularity, seasonal trends, and item-item similarity

Establishing relationship between two things based on data, maybe we don't know earlier how they are related

# Yield decomposition vs Yield gap decomposition

Explainable Machine learning says how management directly and interacting with other variable affects yield ?

For example, ML says the farms with low irrigation and late sowing has less yield.

Does not it mean, management of these factors can improve yield and hence close yield gap ?

| N | P | Irrigation | Sowing | Weed | Soil | | | Yield |
|---|---|---|---|---|---|---|---|---|
| 100 | 50 | 5 | 05-Nov | Most | Heavy | Farmer 1 | | 2500 |
| 150 | 55 | 5 | 06-Nov | Most | Heavy | Farmer 2 | | 3500 |

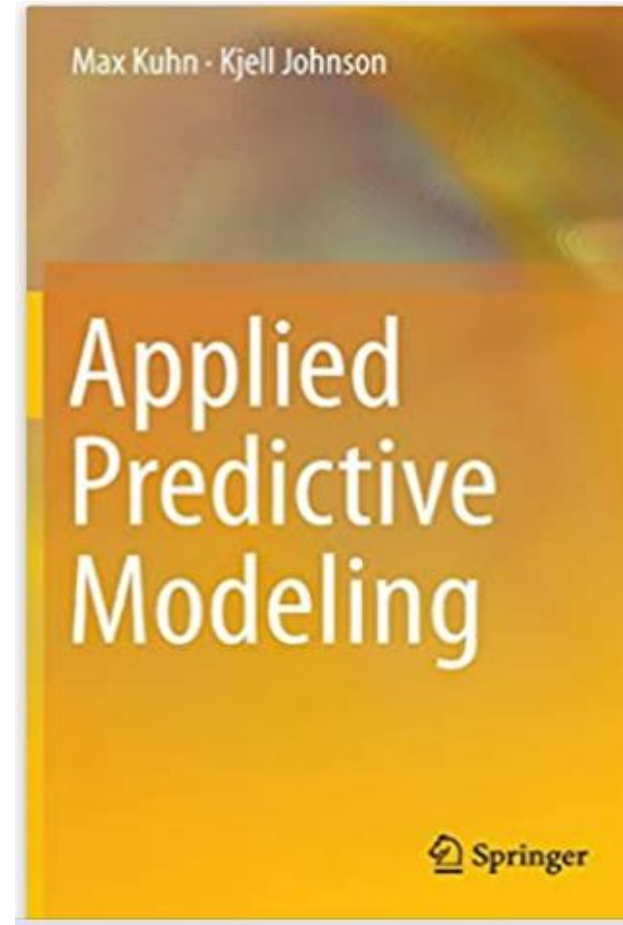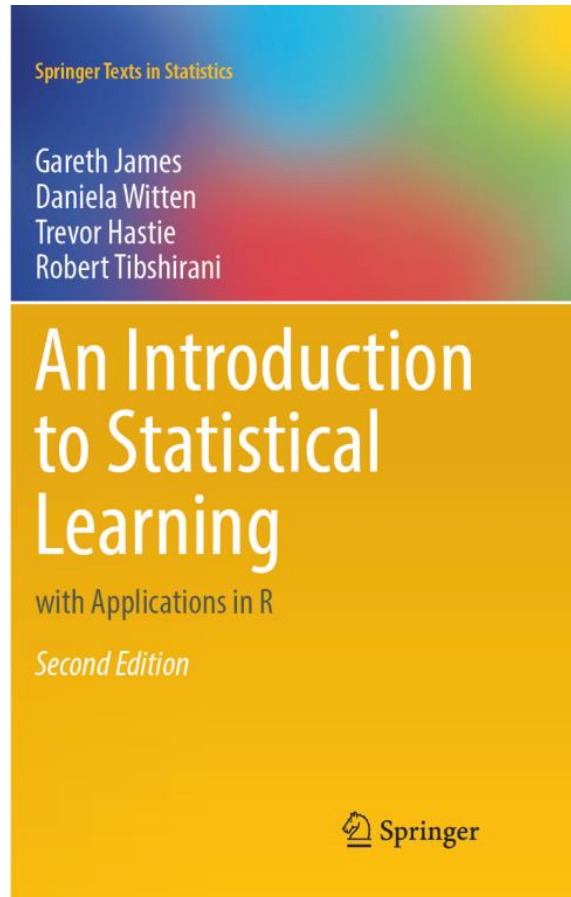Explainable ML says Farmer 1 attained 1000 kg less yield because of 50 kg less N → Part of yield gap

But it does not say anything about, what will be the yield if farmer 1 uses 50 kg additional N with all other existing management. → Predictive ML and Yield gap

1. Development of Model
2. Feature engineering and fine tuning
3. Model interpretation for yield explanation...
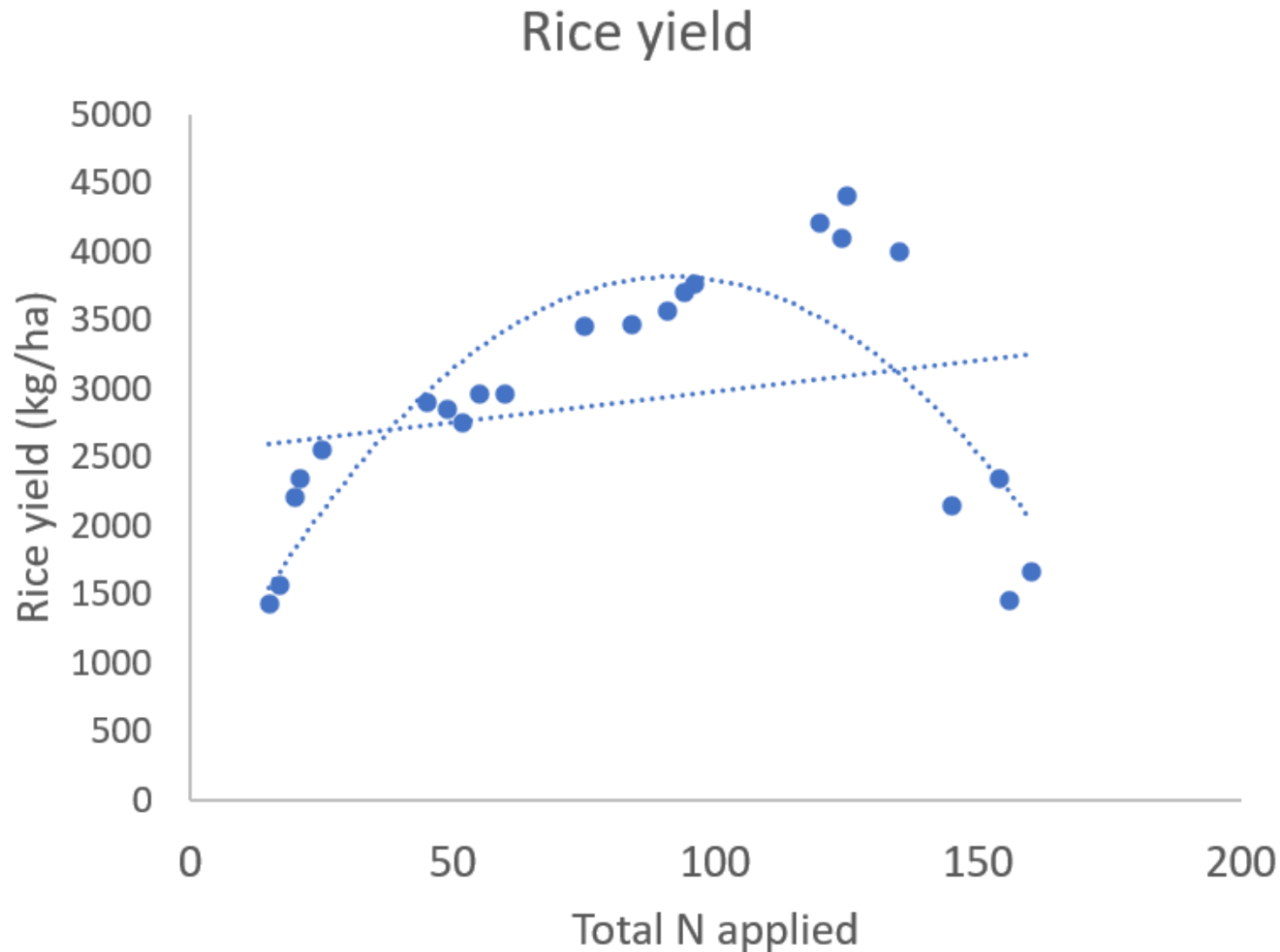
# Step 1: Model comparison

## Why we need this ?

- Based on data structure various model may fit differently.

- If underlying structure is linear and small data set linear regression, Lasso or Ridge regression may perform better.

- If too much complex, nonlinear interactions tree-based models works best.

- if too-too much complex ensemble model may be required (mixture of two model one overcoming the others limitation).

Useful in many agronomic data not limited to yield gap only…. (Clustering of farms, Resource use efficiency, homogenous of weed occurrence, …)

# Linear and quadratic regression
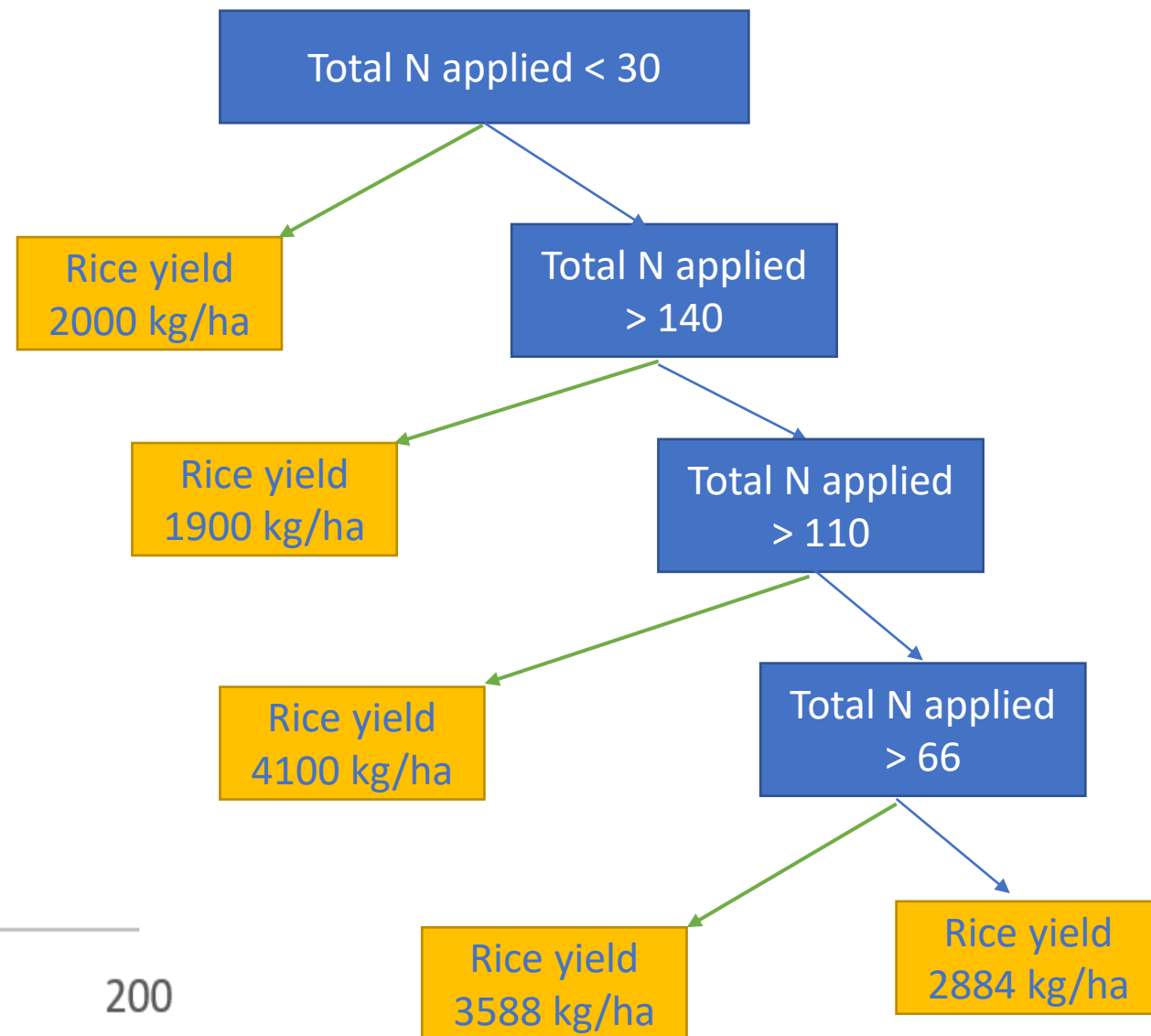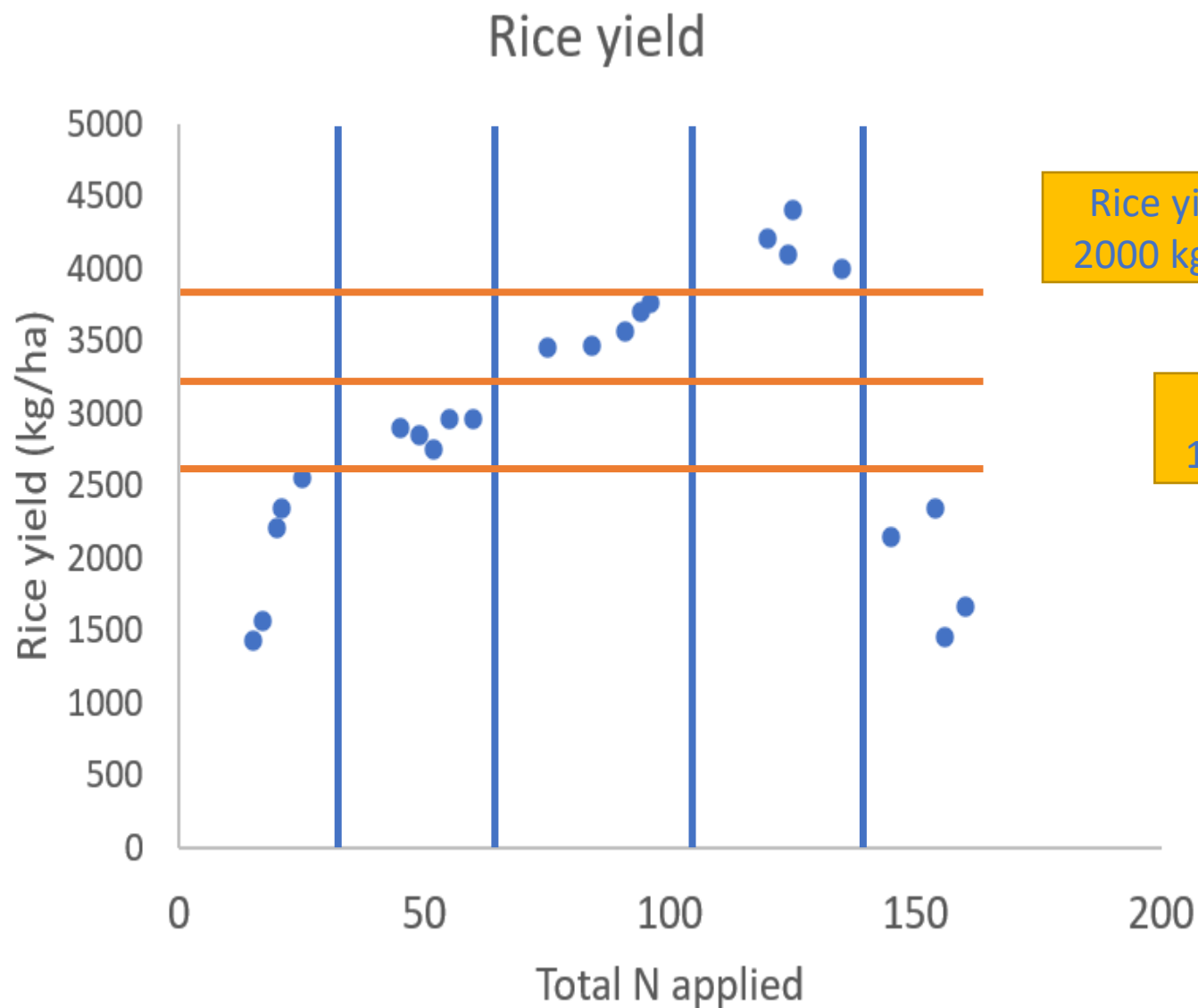
## Rice yield



1. We get coefficient which can be directly interpreted
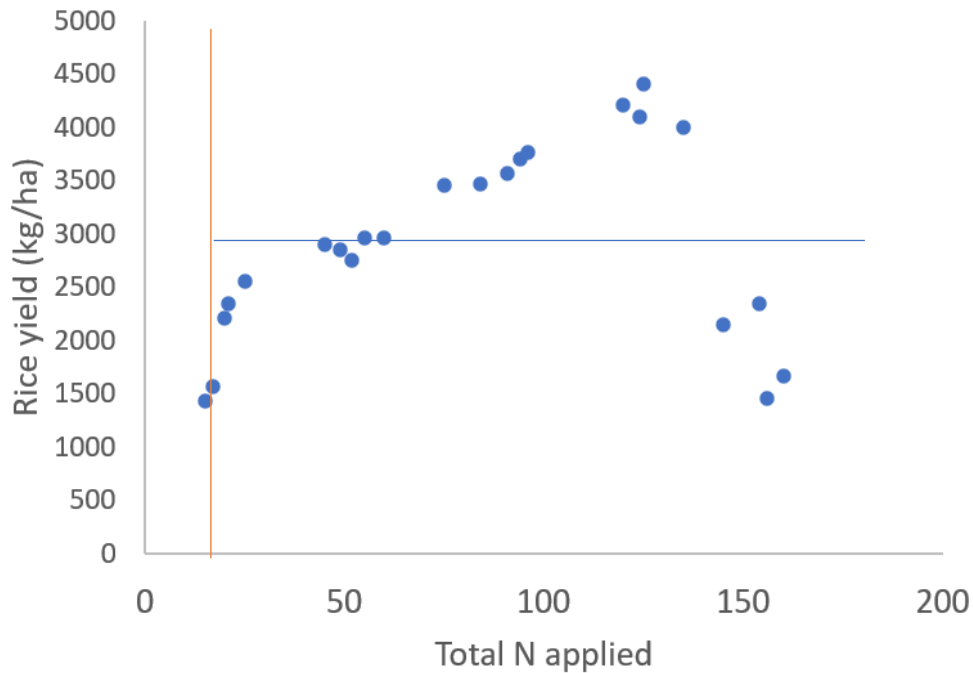2. No iteration required
3. Easy to calculate and interpret
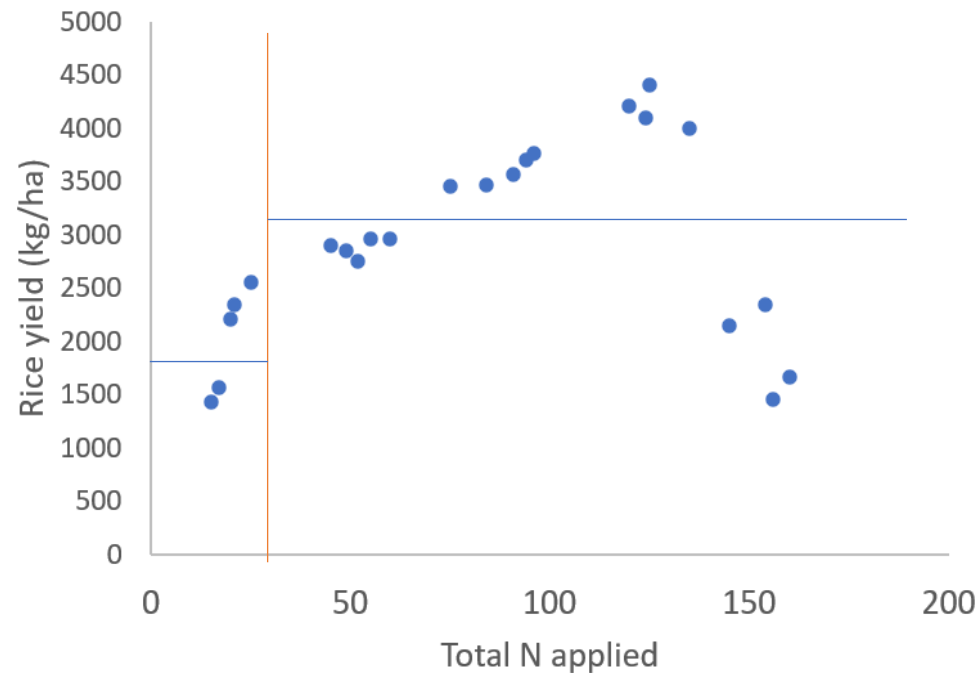
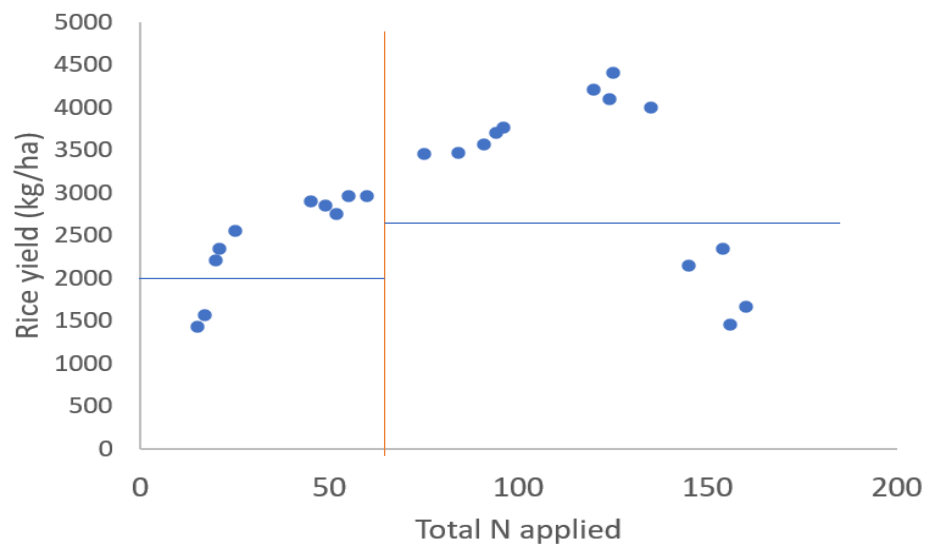1. Not good for nonlinear and complex interactions

# Regression tree intuition

How regression tree choose variable and threshold

Sowing days after December 15 ≤ −1.5 → 4.5 t ha$^{-1}$, $n$ = 107; 25%

Environment: Fultala, Sadar, Kaligonj, Dumuria, Ujipur → 4.2 t ha$^{-1}$, $n$ = 142; 34%; RE=0.09

Sowing days after December 15 ≥ −1.5 → 3.5 t ha$^{-1}$, $n$ = 35; 8%

4.0 t ha$^{-1}$, $n$ = 162; 38%; RE=0.14

Environment: Kalapara → 2.7 t ha$^{-1}$, $n$ = 20; 5%

Sowing days after December 15 ≤ 2.5

Sowing days after December 15 < 10.5 → 3.5 t ha$^{-1}$, $n$ = 30; 7%

3.1 t ha$^{-1}$, $n$ = 60; 14%; RE=0.10

Sowing days after December 15 ≥ 10.5 → 2.7 t ha$^{-1}$, $n$ = 30; 7%

3.3 t ha$^{-1}$, $n$ = 279; 66%; RE=0.21

Sowing days after December 15 ≥ 2.5

Environment: Fultala, Sadar, Kaligonj

2.4 t ha$^{-1}$, $n$ = 117; 28%; RE=0.16

N rate ≥ 34 kg ha$^{-1}$

Environment: Dumuria, Kalapara, Ujipur

Sowing days after December 15 < 13 → 1.9 t ha$^{-1}$, $n$ = 34; 8%

1.6 t ha$^{-1}$, $n$ = 56; 66%; RE=0.11

Sowing days after December 15 ≥ 13 → 1.1 t ha$^{-1}$, $n$ = 23; 5%

2.6 t ha$^{-1}$, $n$ = 422; 100%; RE=.52

N rate < 34 kg ha$^{-1}$

Environment: Fultala, Sadar, Kaligonj → 1.6 t ha$^{-1}$, $n$ = 72; 17%

1.3 t ha$^{-1}$, $n$ = 143; 34%; RE=0.29

Environment: Dumuria, Kalapara, Ujipur → 0.98 t ha$^{-1}$, $n$ = 71; 17%

(Krupnik et al., 2015)

# Features of Classification and regression tree

1. Interpretation is easy based on some rule (Structure of tree for model explanation)

2. Variable selection is based on some criteria (minimization of sum of square of error)

3. We can grow tree up to any complexity depth (Hyperparameter; cp)

4. Sensitive to data, variable, and outlier

5. Solution require Iterative running

6. Also, we can take classification problem based on Ginni's Purity index

7. rpart, rpart.plot

# Random forest; bagging method

## **Step 1:** Create a "bootstrapped" dataset

| Original Dataset | $X_1$ | $X_2$ | $X_3$ | $X_4$ | $X_5$ | $X_6$ | $X_7$ | $X_8$ | $X_9$ | $X_{10}$ |
|---|---|---|---|---|---|---|---|---|---|---|

| Bootstrap 1 | $X_8$ | $X_6$ | $X_2$ | $X_9$ | $X_5$ | $X_8$ | $X_1$ | $X_4$ | $X_8$ | $X_2$ |
|---|---|---|---|---|---|---|---|---|---|---|

| Bootstrap 2 | $X_{10}$ | $X_1$ | $X_3$ | $X_5$ | $X_1$ | $X_7$ | $X_4$ | $X_2$ | $X_1$ | $X_8$ |
|---|---|---|---|---|---|---|---|---|---|---|

| Bootstrap 3 | $X_6$ | $X_5$ | $X_4$ | $X_1$ | $X_2$ | $X_4$ | $X_2$ | $X_6$ | $X_9$ | $X_2$ |
|---|---|---|---|---|---|---|---|---|---|---|

**Step 2:** Create a decision tree using the bootstrapped dataset, but only use a random subset of variables (or columns) at each step.
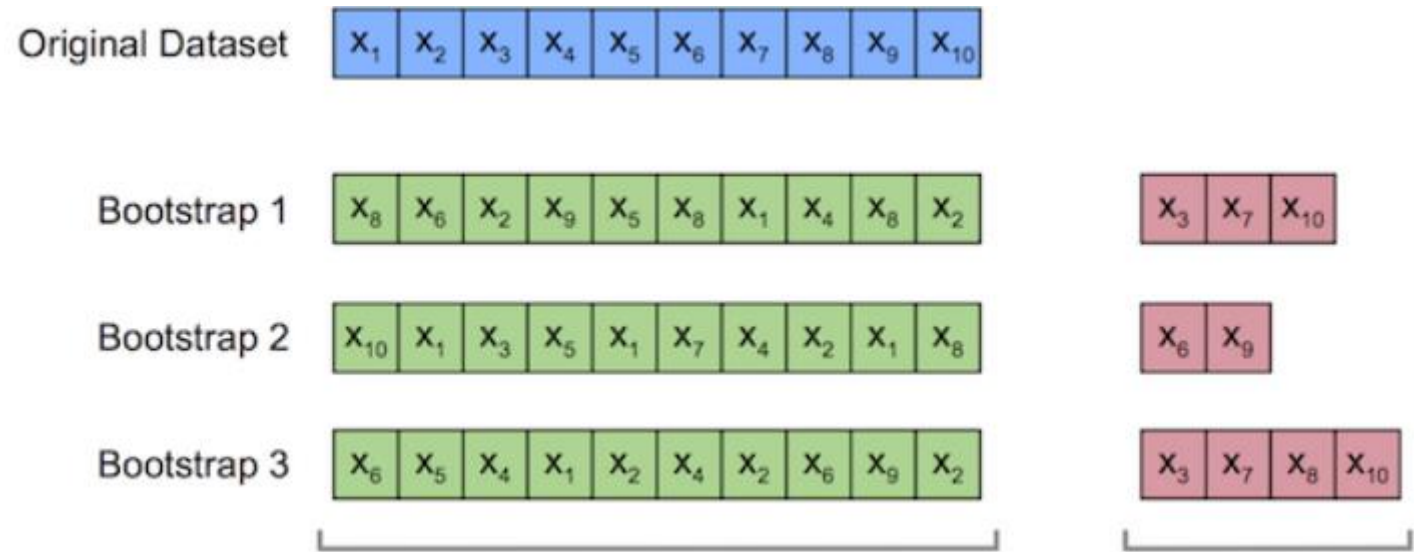
1) Using a bootstrapped dataset (Now less sensitive to outlier)

2) Only considering a random subset of variable at each step (mtry or max.nodesize)

   (Not over rely on a variable)

3) Repeat ideally for 100, 500, or 1000 times (num.tree)

4) We get wide variety of tree (More effective than individual decision tree)

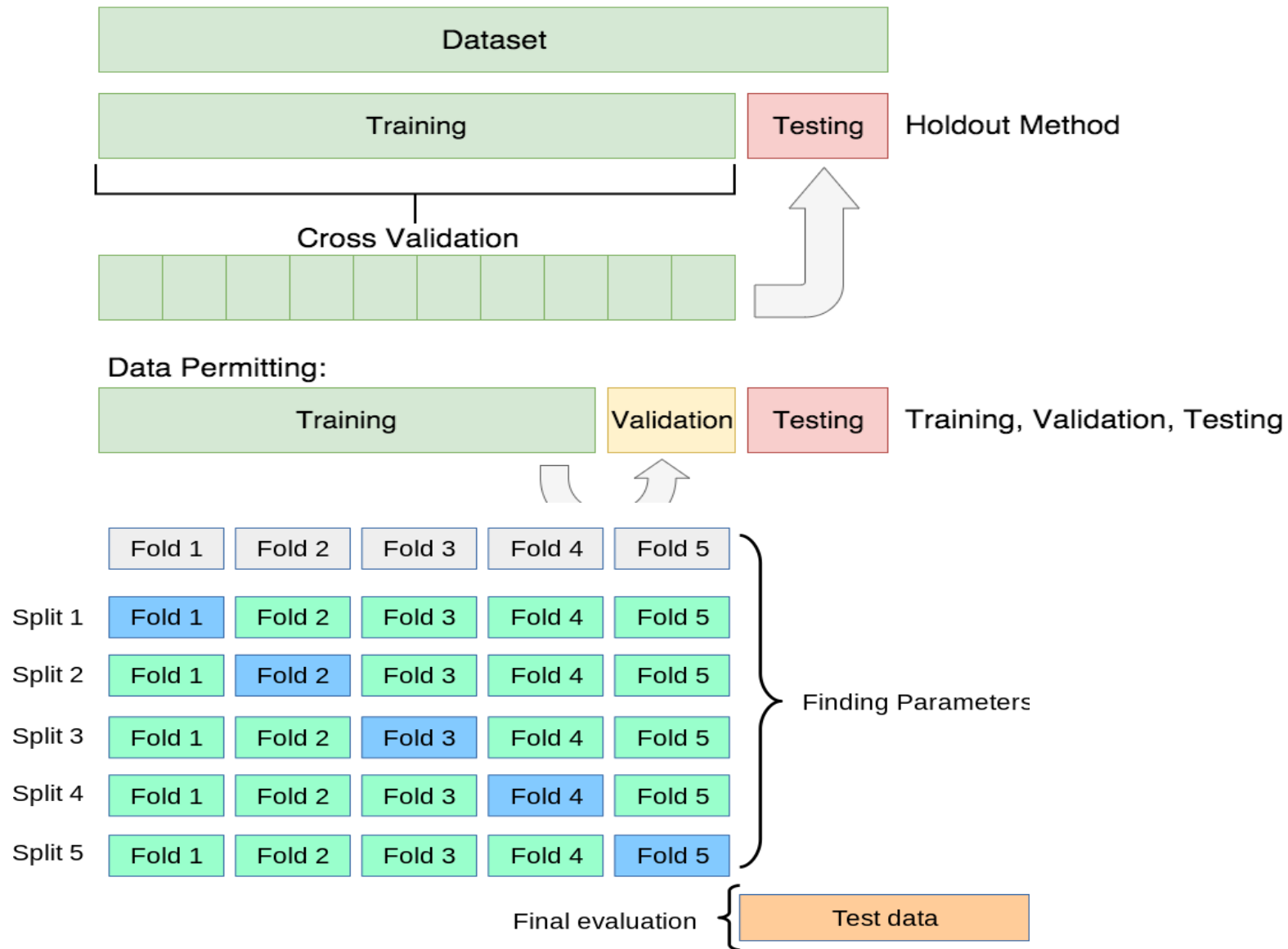# Bootstrap the data + Aggregating the decision = Bagging

# Internal check on performance on Out-of-bag or Out-of-boot data set



| Original Dataset | $x_1$ $x_2$ $x_3$ $x_4$ $x_5$ $x_6$ $x_7$ $x_8$ $x_9$ $x_{10}$ | |
| --- | --- | --- |
| Bootstrap 1 | $x_8$ $x_6$ $x_2$ $x_9$ $x_5$ $x_8$ $x_1$ $x_4$ $x_8$ $x_2$ | $x_3$ $x_7$ $x_{10}$ |
| Bootstrap 2 | $x_{10}$ $x_1$ $x_3$ $x_5$ $x_1$ $x_7$ $x_4$ $x_2$ $x_1$ $x_8$ | $x_6$ $x_9$ |
| Bootstrap 3 | $x_6$ $x_5$ $x_4$ $x_1$ $x_2$ $x_4$ $x_2$ $x_6$ $x_9$ $x_2$ | $x_3$ $x_7$ $x_8$ $x_{10}$ |

- Performance of random forest is checked based on correct prediction of out-of-bag sample.
- **Out-of-bag error**

Cross-validation: Techniques to check performance on unseen data set while building model

# Permutation based feature importance

# Feature selection methods: Recursive feature elimination



## Algorithm 1: Recursive feature elimination

1.1 Tune/train the model on the training set using all predictors

1.2 Calculate model performance

1.3 Calculate variable importance or rankings

1.4 **for** *Each subset size* $S_i$, $i = 1 \ldots S$ **do**

1.5     Keep the $S_i$ most important variables

1.6     [Optional] Pre–process the data

1.7     Tune/train the model on the training set using $S_i$ predictors

1.8     Calculate model performance

1.9     [Optional] Recalculate the rankings for each predictor

1.10 **end**

| | Round 1 | Round 2 | Round 3 | Round 4 | Round 5 | Round 6 | Round 7 |
|---|---|---|---|---|---|---|---|
| | V2 | V2 | V2 | V2 | V2 | V2 | V2 |
| | V1 | | V1 | V1 | V1 | V1 | V1 |
| | V3 | | | V3 | V3 | V3 | V3 |
| | V5 | | | | V5 | V5 | V5 |
| | V4 | | | | | V4 | V4 |
| | V6 | | | | | | V6 |
| Number of variable | | 1 | 2 | 3 | 4 | 5 | 6 |
| RMSE | | | | | | | |
| R2 | | | | | | | |

(https://www.youtube.com/watch?v=SljoN0cO95Q)
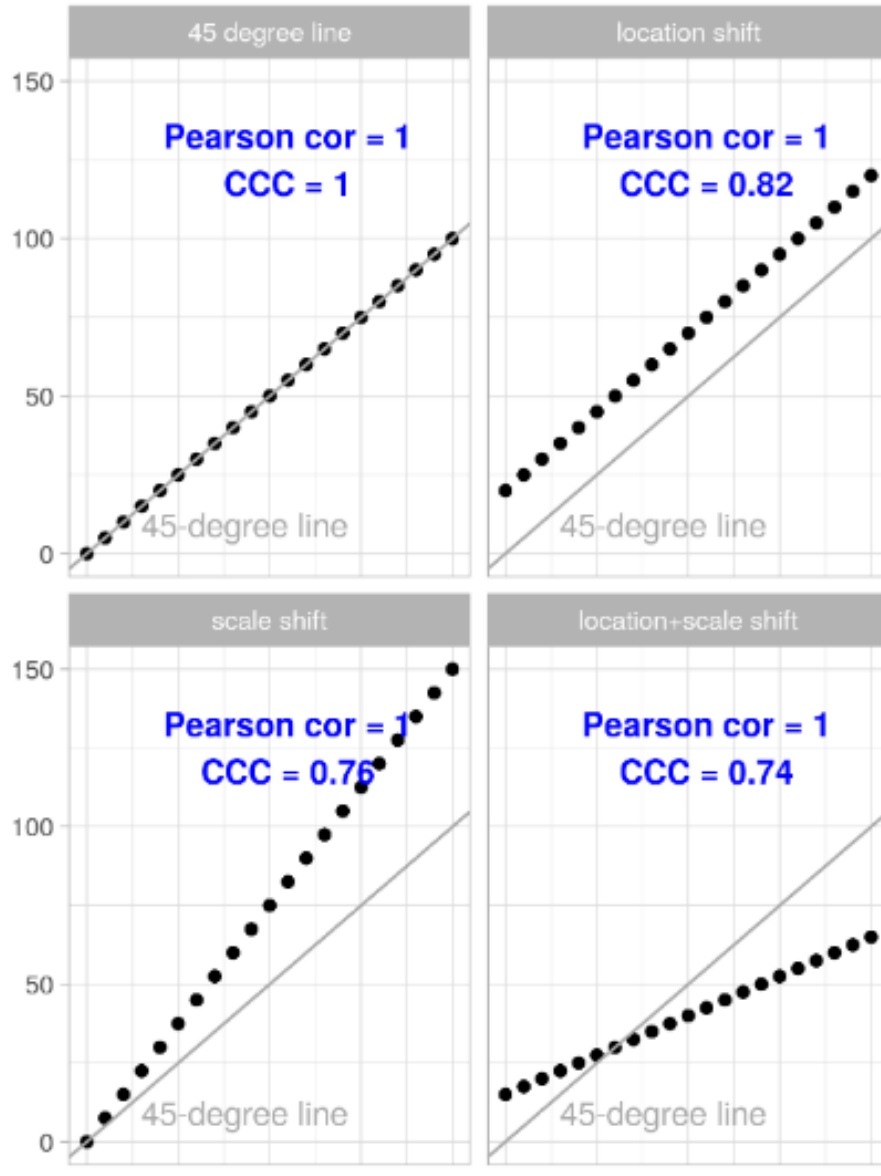
# There can be many random forest..

1. Based on number of random subset of variable selected (mtry) or Number of data points at terminal node (nodesize)

2. Based on number of trees in the random forest (num.tree)

Aim is to select

(The) random forest with minimum out-of-bag error and (The) random forest which performs well on un-seen data (Test data)

## Hyperparameter tuning

# Model evaluation statistics



Internal  Only using the calibration/training information and model diagnostics

External  Using external information, not used in model calibration/training

Cross-validation  Simulating external assessment with the same dataset used for model building

$$\text{ME} = \frac{1}{n}\sum_{i=1}^{n}(y_i - \hat{y}_i) \qquad \text{RMSE} = \left[\frac{1}{n}\sum_{i=1}^{n}(\hat{y}_i - y_i)^2\right]^{1/2}$$

The CCC includes all sources of deviation from a perfect model:
- location shift (bias) $(\mu_1 - \mu_2)/\sqrt{\sigma_1\sigma_2}$
- scale shift (slope not 1) $\sigma_1/\sigma_2$
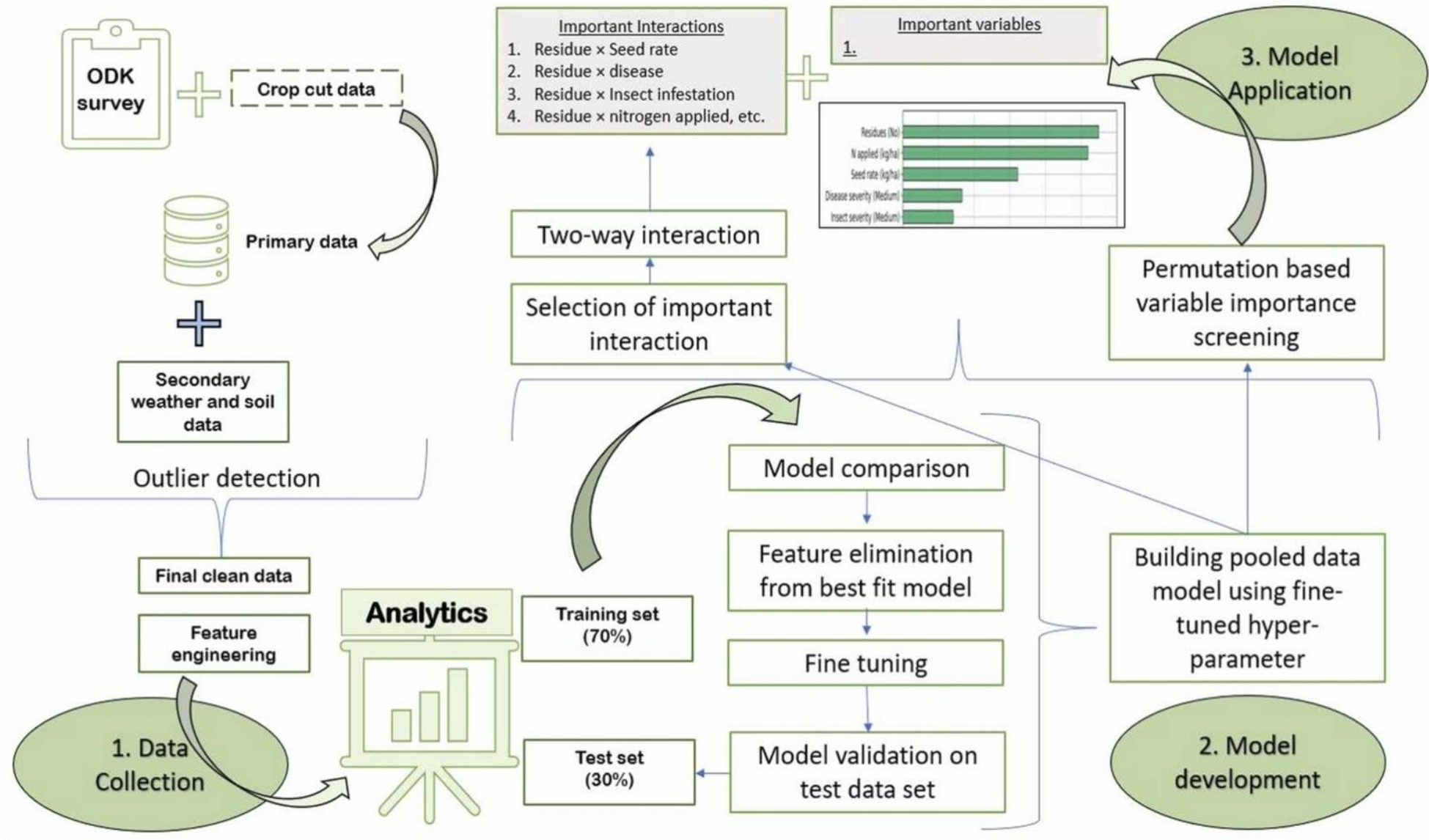- lack of correlation (spread) $1 - \rho_{1,2}$

(David Rosetire)

Numerous model (and variants) exists in literature and are increasing day by day too…
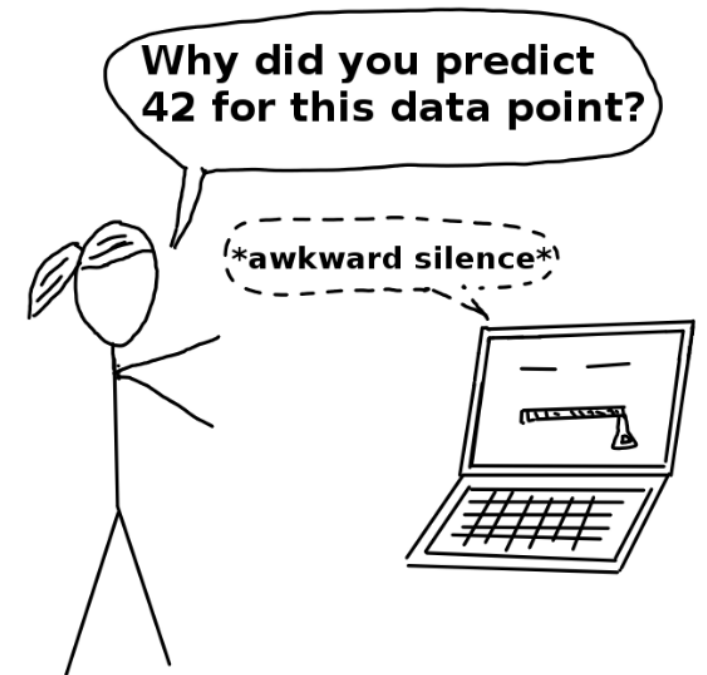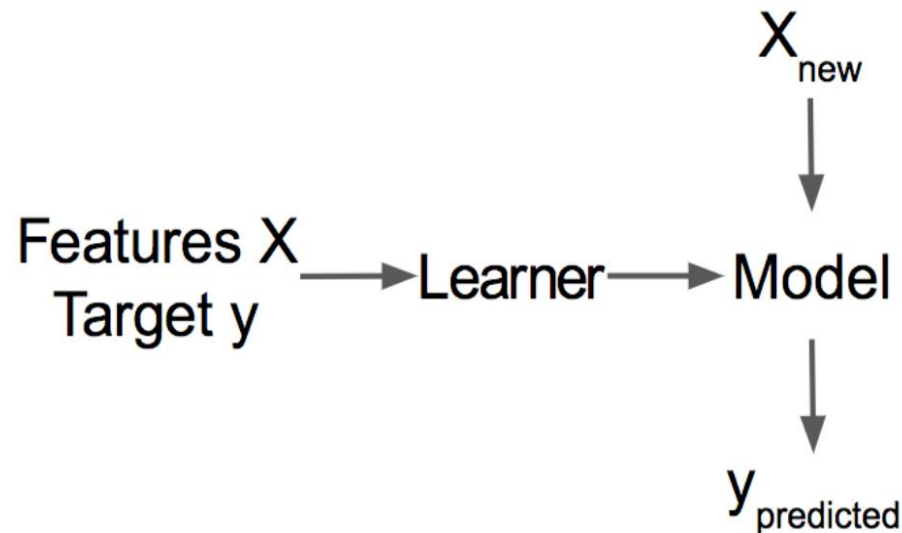
Show 238 entries

Search:

| Model | *method* Value | Type | Libraries | Tuning Parameters |
|---|---|---|---|---|
| AdaBoost Classification Trees | adaboost | Classification | fastAdaboost | nIter, method |
| AdaBoost.M1 | AdaBoost.M1 | Classification | adabag, plyr | mfinal, maxdepth, coeflearn |
| Adaptive Mixture Discriminant Analysis | amdai | Classification | adaptDA | model |
| Adaptive-Network-Based Fuzzy Inference | ANFIS | Regression | frbs | num.labels, max.iter |

Snapshot from caret github page on 13/05/2023

# What a typical ML workflows should be ?



(Nayak, H.S., Silva J V., Parihar C.M., et al., 2022)

# RF model developed! How to interpret ??

1. We don't have any coefficient alike in linear regression for interpreting or publishing

2. We don't have a regression tree like structure (We have 500s of trees)

3. How to say which variable are important ? And their effect ?

# Model interpretation

Linear regression
Classification and regression tree
…

Intrinsically interpretable models
            Vs
Post hoc (and model-agnostic) interpretation models.

Gradient boosting
Random forest
…

**Local or global?**

Which variable is more important for describing yield at population level ?

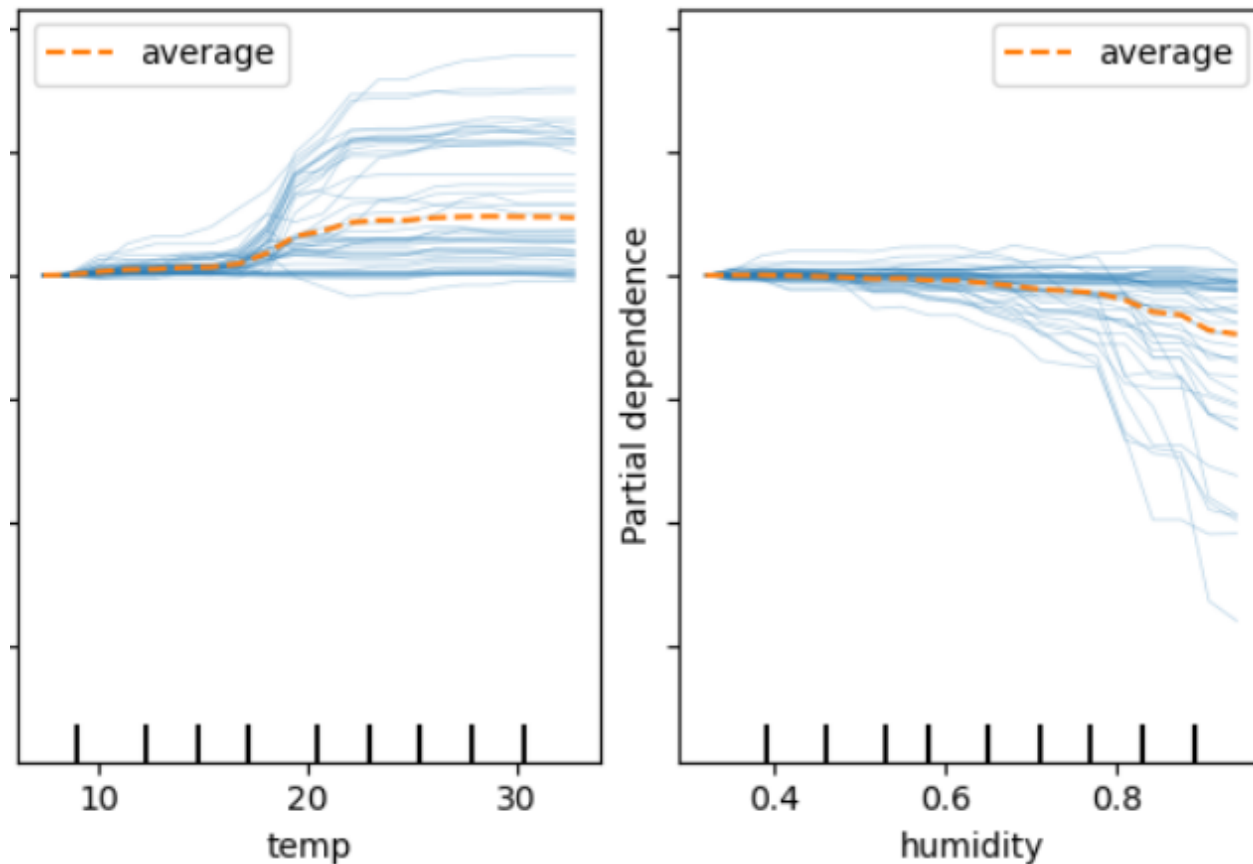What is the effect of the variable (say N) on yield ?

With which factor N is interacting to affect yield ? (Important two-way interactions)

At farm level what is the effect of any production practices on rice yield ?

# Effect of production practices on yield or (any outcome of interest)

Partial dependence plots (Single variable or two-way interactions)



ICE and PDP representations

$$\hat{f}_S(x_S) = E_{X_C}\left[\hat{f}(x_S, X_C)\right]$$

By marginalizing the machine learning model output over the distribution of the features in set C, PDP shows the relationship between the features in set S we are interested in and the predicted outcome

# Feature interaction

We will look at the exercise..

# Thank You