



# Tree-based models

# 1. Predictive modeling

2

Predictive modeling is required for :

- Estimating the function behind the data
- Taking decisions based on science, not arbitrarily

Viettel Digital Talent 2023

## 2. Choosing a predictive model

3

Predictive models approximate a mapping function  $F$  that maps input feature  $\mathbf{X}$  to the target feature  $Y$ :  $F(\mathbf{X}) + \varepsilon$ . The form of the true function  $F$  being unknown, models make assumptions regarding this functional form and estimate the function.

There are 2 types of models:

- **Parametric models:** Models that simplify the function  $F$  to a known functional form<sup>1</sup>
- **Non-parametric models:** Models that are free to learn any functional form

There are 2 types of predictive tasks:

- **Regression task:** Task of approximating a mapping function  $F$  when the target feature is quantitative
- **Classification task:** Task of approximating a mapping function  $F$  when the target feature is qualitative<sup>2</sup>

<sup>1</sup> Usually, the assumed functional form is a linear combination of the input features and, as such, those predictive models refer to linear models

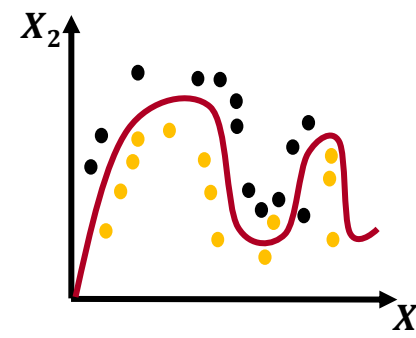
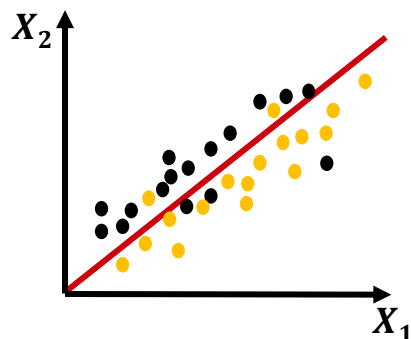
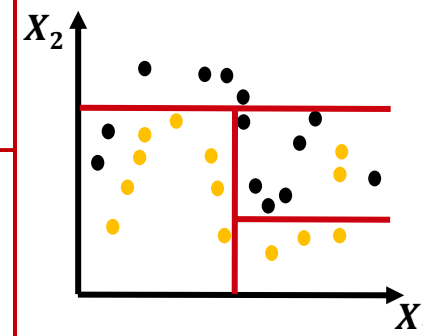
<sup>2</sup> Either binary or multi-class

## 2. Choosing a predictive model

4

There are 2 types of models:

	Parametric models	Non-parametric models
Benefits	Less data required to estimate the function Fast to learn from data as a fixed number of parameters (small) needs to be estimated Simple to understand Robust to overfitting	Flexible to fit a large number of functional forms Robust as no/few assumptions are made about the function Good fit as many forms are involved
Limitations	Constrained to fit a specific functional form Not robust as strong assumptions are made about the function Poor fit if the function is not linear	More data required to estimate the function Slow to learn from data as the number of parameters grow with the size of data Difficult to understand Prone to overfitting <sup>1</sup>



<sup>1</sup> Non-parametric models estimate the function while getting the closest possible to the data points and are then more prone to overfitting



## 2. Choosing a predictive model

The choice of the predictive model will then depend on:

- The type of the model
- The type of the predictive task
- The strengths and weaknesses of the model

Model	Type	Task
Linear regression	Parametric	Regression
Logistic regression	Parametric	Classification
Poisson regression	Parametric	Regression
Linear discriminant analysis	Parametric	Classification
Penalized regression (Ridge, Lasso, Elastic Net)	Parametric	Both
Decision tree	Non-parametric	Both
Support vector machine	Both	Both
K-nearest neighbors	Non-parametric	Both
Tree-based models	Non-parametric	Both

	Linear/Logistic/ Poisson regression	Penalized regression	Linear SVC	Non-linear SVC	CART tree	Bagging	Random Forests	Boosting
<b>Robustness to overfitting</b>	No penalization	Penalization	- Support vectors only - Penalization	- Support vectors only - Penalization	Instable method <sup>2</sup>	Aggregation	Aggregation	Aggregation
<b>Performance</b>	- Dataset-dependent - Linear relation <sup>1</sup>	- Dataset-dependent - Linear relation <sup>1</sup>	- Margin maximization - Linear relation <sup>1</sup>	Margin maximization	Local optimum	Variance optimization	Variance optimization	Bias optimization
<b>Interpretability</b>	Coefficient analysis	Coefficient analysis	- Support vectors - Black box	- Support vectors - Black box	- Binary decisions - Tree diagram	- Multiple trees - Black box	- Multiple trees - Black box	- Multiple trees - Black box
<b>Implementation</b>	No hyper-parameter tuning	Few hyper- parameters	Few hyper- parameters	Few hyper- parameters	Some hyper- parameters	Many hyper- parameters	Many hyper- parameters	Many hyper- parameters
<b>Computational cost</b>	Simple model	Simple model	Scalar product	Scalar product	Greedy search	Many iterations	Many iterations	Many iterations

<sup>1</sup> GLM models and linear SVC might be under-performing if the relationship with the target feature is not linear

<sup>2</sup> Even though CART tree might be penalized (trimmed), the strategy still relies on a greedy search which makes the method unstable and prone to overfitting

# 2. Choosing a predictive model

6

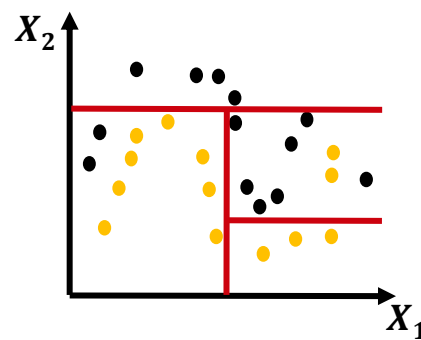
## 1. CART Decision tree<sup>1</sup>

Mapping: No functional form

Loss function at leaf  $l$ <sup>3</sup>:

- Misclassification or 0-1 loss function:  $Loss(l) = 1 - \max_k(\hat{p}_{lk})$
- Gini index:  $Loss(l) = 1 - \sum_{k=1}^K (\hat{p}_{lk})^2$
- Cross entropy:  $Loss(l) = - \sum_{k=1}^K \hat{p}_{lk} \log(\hat{p}_{lk})$

Assumption: None



The idea is to find recursively the input feature and split point (i.e. best binary decision) that minimizes an impurity measure of a sample of records forming a region

Type	Non-Parametric
Target feature	Qualitative <sup>2</sup>
Distribution	No assumption
Task	Classification <sup>2</sup>
Loss	Misclassification error or Gini index or Cross entropy
Method	Node purity maximization through greedy search

<sup>1</sup> Only the CART decision tree is considered here, which is a decision tree that splits a region into 2 sub-regions thanks to binary decisions

<sup>2</sup> Only the classification task, in a binary set-up, is considered here but the target feature could be quantitative with a resulting task being regression

<sup>3</sup> The loss at node level refers as the node impurity

# 2. Choosing a predictive model

7

## 1. CART Decision tree

Algorithm:

1. From a dataset  $D$  composed of  $N$  records
2. For each feature  $X_j$  and splitting point  $S_j$  from  $X_j$ , calculate the purity gain<sup>1</sup> from splitting the region into 2 sub-regions: 
$$\text{gain} = \text{loss}_{R_{\text{root}}} - \frac{\sum_{i=1}^N \mathbb{I}(X_i \in R_{\text{left}})}{\sum_{i=1}^N \mathbb{I}(X_i \in R_{\text{root}})} \text{loss}_{R_{\text{left}}} - \frac{\sum_{i=1}^N \mathbb{I}(X_i \in R_{\text{right}})}{\sum_{i=1}^N \mathbb{I}(X_i \in R_{\text{root}})} \text{loss}_{R_{\text{right}}}$$
3. Choose  $(X_j, S_j)$  in order to maximize the purity gain
4. Split the region accordingly to  $(X_j, S_j)$
5. Repeat steps 2-4 until a stopping criterion<sup>2</sup> is met
6. Prune tree<sup>3</sup>:  $\arg \min_{\alpha} \left[ \sum_{l=1}^L \left[ \sum_{i=1}^N \mathbb{I}(X_i \in R_l) \right] \text{loss}_{R_l} \right] + \alpha L$

<sup>1</sup> Which corresponds to the impurity decrease from splitting the region

<sup>2</sup> Usually, a stopping criteria is intended to prevent overfitting (number of records in a leaf, depth of a tree,...) but is different from pruning

<sup>3</sup> The higher the  $m$ , the smaller the tree

## 2. Choosing a predictive model

### 2. Bootstrap Aggregating or Bagging

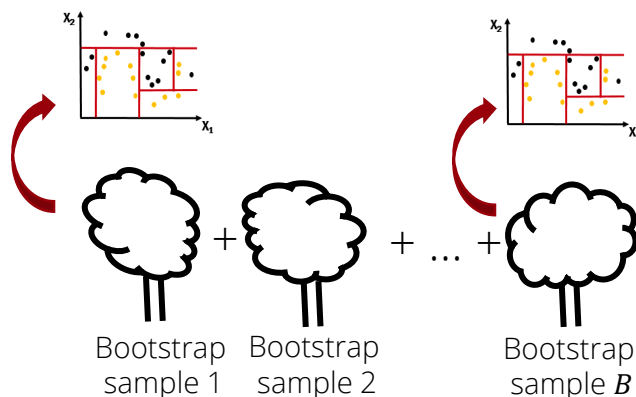
Mapping: No functional form

Loss function at leaf  $l$ :

- Misclassification or 0-1 loss function:  $Loss(l) = 1 - \max_k(\hat{p}_{lk})$
- Gini index:  $Loss(l) = 1 - \sum_{k=1}^K (\hat{p}_{lk})^2$
- Cross entropy:  $Loss(l) = - \sum_{k=1}^K \hat{p}_{lk} \log(\hat{p}_{lk})$

Assumption: None

Type	Non-Parametric
Target feature	Qualitative <sup>1</sup>
Distribution	No assumption
Task	Classification <sup>1</sup>
Loss	Misclassification error or Gini index or Cross entropy
Method	Variance reduction through simultaneous creation of different experts (bootstrap sampling) and aggregation



The idea is to generate simultaneously  $B$  bootstrap samples<sup>2</sup> from the dataset, train their own  $B$  CART trees and aggregating their votes through a majority voting system

<sup>1</sup> Again, only the classification task, in a binary set-up, is considered

<sup>2</sup> See Annex 3 about Bootstrapping for further explanation

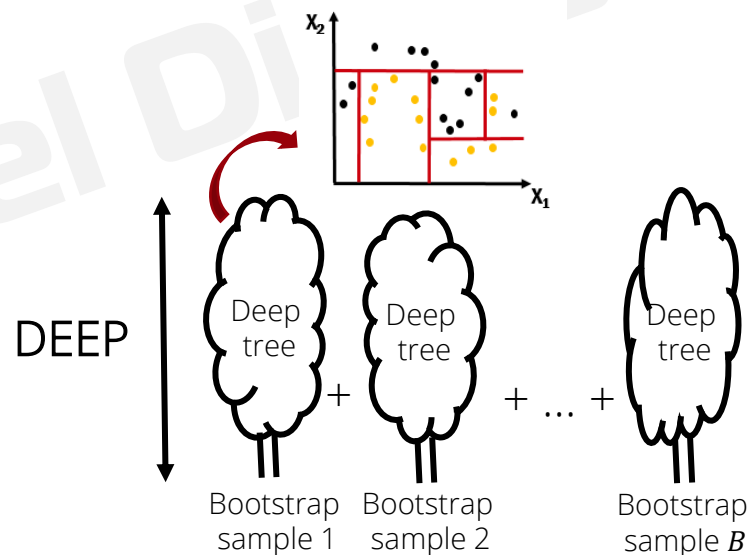


## 2. Choosing a predictive model

9

### 2. Bootstrap Aggregating or Bagging

- In case of the squared loss (regression task): it decreases the test error of the aggregated model whatever the tree's structure
- In case of the 0-1 loss (classification task): it decreases the test error of the aggregated model only if the trees are rather unbiased<sup>1</sup>



As Bagging decreases the **variance** while leaving the **bias** unchanged, it makes sense to train quite unbiased **trees** in order to act upon the variance

<sup>1</sup> Trees are then usually grown without trimming, at least deep enough

# 2. Choosing a predictive model

10

## 2. Bootstrap Aggregating or Bagging

Algorithm :

1. From a dataset  $D$  composed of  $N$  records
2. For each iteration  $b$  ( $b = 1, \dots, B$ ):
  - I. Create a bootstrap sample  $D_b$  by randomly drawing  $N$  records with replacement from  $D$
  - II. Train a model  $\hat{F}_b$  from  $D_b$
3. Let each of the  $B$  models classify the records and assign the class  $k$  based on a majority voting:

$$\hat{F}(X) = \arg \min_k \sum_{b=1}^B \mathbb{I}(\hat{F}_b(X) \neq k)$$

## 2. Choosing a predictive model

### 3. Random Forests

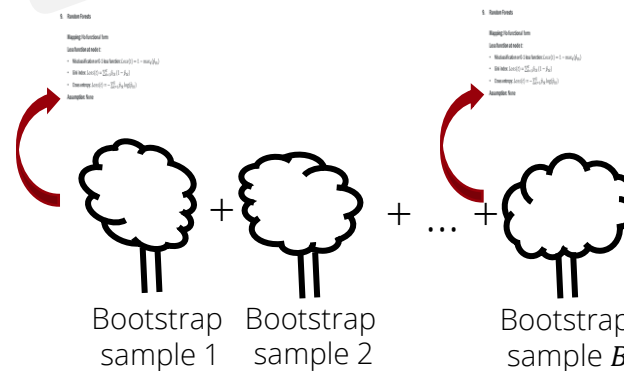
Mapping: No functional form

Loss function at leaf  $l$ :

- Misclassification or 0-1 loss function:  $Loss(l) = 1 - \max_k(\hat{p}_{lk})$
- Gini index:  $Loss(l) = 1 - \sum_{k=1}^K (\hat{p}_{lk})^2$
- Cross entropy:  $Loss(l) = - \sum_{k=1}^K \hat{p}_{lk} \log(\hat{p}_{lk})$

Assumption: None

Type	Non-Parametric
Target feature	Qualitative <sup>1</sup>
Distribution	No assumption
Task	Classification <sup>1</sup>
Loss	Misclassification error or Gini index or Cross entropy
Method	Variance reduction <sup>2</sup> through bagging and random subset of features



The idea is similar to Bagging<sup>3</sup> except that an input feature is chosen for binary decision from a random sample of  $v$  features without replacement

<sup>1</sup> Again, only the classification task, in a binary set-up, is considered

<sup>2</sup> Variance decreases even more thanks to extra tree decorrelation (random sampling of records + random subset of features)

<sup>3</sup> See supra for explanation about Bagging

## 2. Choosing a predictive model

12

### 3. Random Forests

Algorithm :

1. From a dataset  $D$  composed of  $N$  records
2. For each iteration  $b$  ( $b = 1, \dots, B$ ):
  - I. Create a bootstrap sample  $D_b$  by randomly drawing  $N$  records with replacement from  $D$
  - II. Train a model  $\hat{F}_b$  from  $D_b$  while:
    - A. Choosing randomly  $v$  features without replacement from the  $p$  features
    - B. Choosing  $X_j$  from  $v$  and  $S_j$  in order to maximize the gain and splitting the region accordingly
3. Let each of the  $B$  models classify the records and assign the class  $k$  based on a majority voting:  
$$\hat{F}(X) = \arg \min_k \sum_{b=1}^B \mathbb{I}(\hat{F}_b(X) \neq k)$$

## 2. Choosing a predictive model

### 4. Adaptive Boosting or AdaBoost

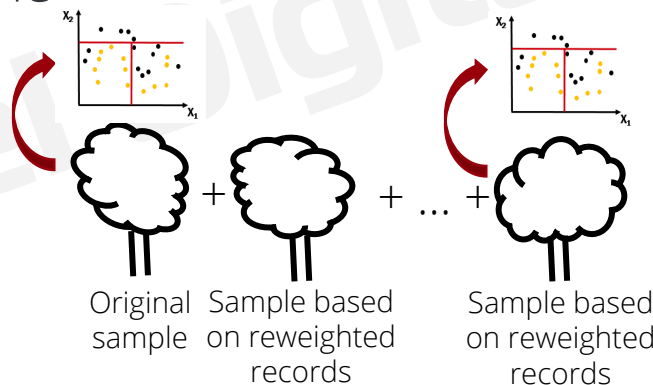
Mapping: No functional form

Loss function:

- Exponential loss:  $Loss = e^{-Y\hat{Y}}$

Assumption: None

Type	Non-Parametric
Target feature	Qualitative <sup>1</sup>
Distribution	No assumption
Task	Classification <sup>1-2</sup>
Loss	Exponential loss
Method <sup>3</sup>	Bias reduction (local margin maximization) through sequential creation of corrective models and weighted aggregation



The idea is to improve a weak tree<sup>4</sup> over a sequence of  $B$  iterations of resampled data<sup>5</sup> and aggregating their votes through a weighted majority voting system

<sup>1</sup> Again, only the classification task, in a binary set-up, is considered

<sup>2</sup> Adaboost assumes the recoding of the binary target to  $\{-1, +1\}$ , which is important for the concept of "margin"

<sup>3</sup> While SVC maximizes the minimum margin for all records at once, Adaboost maximizes the margin locally for each record

<sup>4</sup> A weak tree is a tree that performs slightly better than a random choice in a binary set-up ( $\epsilon \leq 50\%$ )

<sup>5</sup> Last iteration's erroneous predictions get their weights increase and get more chance to get resampled at the next iteration in order to be further modeled for correction

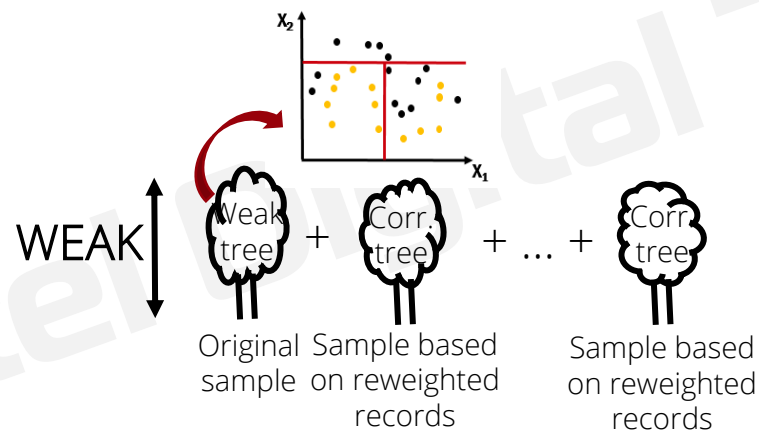


## 2. Choosing a predictive model

14

### 4. Adaptive Boosting or AdaBoost

It decreases the test error of the aggregated model only if the trees are rather biased<sup>1</sup>



As Adaboost decreases the **bias** mostly even though variance might decrease through aggregation, it makes sense to **train quite biased trees** in order to act upon the bias

<sup>1</sup> Trees are then usually grown simple (depth of 2 to 6 levels)

## 2. Choosing a predictive model

15

### 4. Adaptive Boosting or AdaBoost

Weights:

- First iteration:  $w_{1,i} = 1/N$  → Records are given equal weight (uniform weighting)
- Next iterations:  $w_{b+1,i} = w_{b,i} \exp(c_b \mathbb{I}(\hat{F}_b(X_i) \neq Y_i))$  → Record's weights are updated with a constant<sup>1</sup>

Constants<sup>2-3</sup>, which are functions of the test error:

- Breiman's constant:  $c_b = \ln\left(\frac{1-\varepsilon_b}{\varepsilon_b}\right)$
- Freund & Schapire's constant:  $c_b = \frac{1}{2} \ln\left(\frac{1-\varepsilon_b}{\varepsilon_b}\right)$

$$\text{Margin}^4: m(X_i) = \frac{\sum_{b=1}^B c_b \mathbb{I}(\hat{F}_b(X_i) = Y_i)}{\sum_{b=1}^B c_b} - \frac{\sum_{b=1}^B c_b \mathbb{I}(\hat{F}_b(X_i) \neq Y_i)}{\sum_{b=1}^B c_b}$$

<sup>1</sup> Records with erroneous predictions will get their weights updated upwards, while others will get their weights unchanged

<sup>2</sup> The lower  $\varepsilon_b$ , the higher  $c_b$ . Given weak trees ( $\varepsilon \leq 50\%$ ), the constants will be positive

<sup>3</sup> Constants, which are functions of the error, are used to weight the vote of each tree in the ensemble

<sup>4</sup> Correctly predicted records have positive margin while incorrectly predicted records have negative margin (link with SVC and Adaboost)

## 2. Choosing a predictive model

16

### 4. Adaptive Boosting or AdaBoost

Algorithm :

1. From a dataset  $D$  composed of  $N$  records
2. Set initial weights  $w_{1,i} = 1/N \quad \forall i$
3. For each iteration  $b$  ( $b = 1, \dots, B$ ):
  - I. Create a sample  $D_b$  by randomly drawing  $N$  records with replacement from  $D$  based on their weights
  - II. Train a model  $\hat{F}_b$  from  $D_b$
  - III. Calculate the error  $\varepsilon_b = \sum_{i=1}^N w_{b,i} \mathbb{I}(\hat{F}_b(X_i) \neq Y_i)$  and constant  $c_b = \ln\left(\frac{1-\varepsilon_b}{\varepsilon_b}\right)$  or  $\frac{1}{2} \ln\left(\frac{1-\varepsilon_b}{\varepsilon_b}\right)$
  - IV. If  $\varepsilon_b > 50\%$ , the model  $\hat{F}_b$  is not weak, reinitialize weights to  $1/N$  and get to 3.1.
  - V. Update the weights  $w_{b+1,i} = w_{b,i} \exp\left(c_b \mathbb{I}(\hat{F}_b(X_i) \neq Y_i)\right)$  and normalize them
4. Let each of the  $B$  models classify the records and assign the class  $k$  based on a weighted majority voting:  
$$\hat{F}(X) = \arg \min_k \sum_{b=1}^B c_b \mathbb{I}(\hat{F}_b(X) \neq k)$$

# 2. Choosing a predictive model

17

## 5. Gradient Boosting

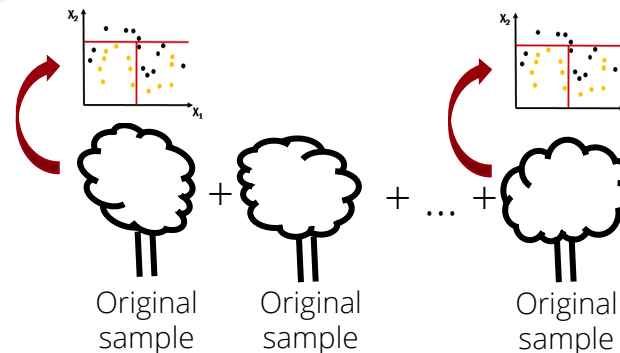
Mapping: No functional form

Loss function:

- Log-loss:  $Loss = -Y \ln(\hat{Y}) - (1 - Y) \ln(1 - \hat{Y})$
- Logistic loss:  $Loss = \ln(1 + e^{-Y\hat{Y}})$
- Exponential loss:  $Loss = e^{-Y\hat{Y}}$

Assumption: None

Type	Non-Parametric
Target feature	Qualitative <sup>1</sup>
Distribution	No assumption
Task	Classification <sup>1-2</sup>
Loss	Log-loss or logistic loss or exponential loss
Method	Bias reduction through sequential creation of corrective models and aggregation <sup>5</sup>



The idea is to improve a weak tree<sup>3</sup> over a sequence of  $B$  iterations while focusing on the residuals<sup>4</sup> and aggregating<sup>5</sup> predictions

<sup>1</sup> Again, only the classification task, in a binary set-up, is considered

<sup>2</sup> Gradient Boosting assumes the binary target within the space  $\{0, +1\}$

<sup>3</sup> Again, weak trees are used, due to the boosting process

<sup>4</sup> As we boost the residuals, which are obtained from the first derivative of the loss function, we boost the gradients

<sup>5</sup> Aggregation with Gradient Boosting is not weighted (i.e. each tree has the same weight) since no resampling of data is involved, unlike in Adaboost

# 2. Choosing a predictive model

18

## 5. Gradient Boosting

Algorithm :

1. From a dataset  $D$  composed of  $N$  records
2. Initialize a first model with a constant value<sup>1-2</sup>:  $\hat{F}_0(X) = \arg \min_{\gamma} \sum_{i=1}^N L(Y_i, \gamma)$
3. For each iteration  $b$  ( $b = 1, \dots, B$ ):
  - I. Calculate the pseudo residuals<sup>3-4</sup>  $r_{ib} = - \left[ \frac{\partial L(Y_i, \hat{F}_{b-1}(X_i))}{\partial \hat{F}_{b-1}(X_i)} \right] \quad i = 1, \dots, N$
  - II. Train a regression tree to the pseudo residuals  $r_{ib}$  to minimize  $\sum_{i=1}^N L(r_{ib}, \hat{r}_{ib})$  and create the leaves  $R_{jb}, j = 1, \dots, J_b$

Given the loss is the log-loss:

<sup>1</sup>  $L(Y, \hat{Y}) = -Y \ln(\hat{p}) - (1 - Y) \ln(1 - \hat{p})$  in terms of probability and  $L(Y, \hat{Y}) = -Y \ln(\hat{p}/1 - \hat{p}) + \ln(1 + e^{\ln(\hat{p}/1 - \hat{p})})$  in terms of odds

<sup>2</sup>  $\gamma$  corresponds to the log-odds:  $\gamma = \ln(\hat{p}/1 - \hat{p})$

<sup>3</sup> First derivative regards to the log-odds  $\frac{\partial L(Y_i, \hat{F}_{b-1}(X_i))}{\partial \hat{F}_{b-1}(X_i)} = g_{ib} = -(Y_i - \hat{p}_{ib})$  and refers to gradient, second derivative  $\frac{\partial^2 L(Y_i, \hat{F}_{b-1}(X_i))}{\partial^2 \hat{F}_{b-1}(X_i)} = h_{ib} = \hat{p}_{ib}(1 - \hat{p}_{ib})$  and refers to the hessian

<sup>4</sup> Pseudo residuals can be expressed as:  $r_{ib} = -g_{ib} = Y_i - \hat{p}_{ib}$  in terms of probability or  $r_{ib} = Y_i - \frac{e^{\ln(\hat{p}_{ib}/1 - \hat{p}_{ib})}}{1 + e^{\ln(\hat{p}_{ib}/1 - \hat{p}_{ib})}}$  in terms of odds, with  $\hat{p}_{ib}$  relating to  $b - 1$



## 2. Choosing a predictive model

19

### 5. Gradient Boosting

Algorithm :

- III. Calculate the leaf output<sup>1-2</sup>  $\gamma_{jb} = \arg \min_{\gamma} \sum_{X_i \in R_{jb}} L(Y_i, \hat{F}_{b-1}(X_i) + \gamma)$ , for each leaf  $j = 1, \dots, J_b$
- IV. Update the predictions<sup>3-4</sup>  $\hat{F}_b(X) = \hat{F}_{b-1}(X) + \varphi \sum_{j=1}^{J_b} \gamma_{jb} \mathbb{I}(X \in R_{jb})$
- 4. Let the last model predict the records<sup>5</sup>:  $\hat{F}(X) = \hat{F}_B(X)$

<sup>1</sup> One leaf may indeed have several records and thus several residuals, so we need to aggregate the results leaf-wise

<sup>2</sup> Given the loss is the log-loss, leaf output corresponds to:  $\gamma_{jb} = \frac{-\sum_{X_i \in R_{jb}} g_{ib}}{\sum_{X_i \in R_{jb}} h_{ib}} = \frac{-\sum_{X_i \in R_{jb}} g_{ib}}{\text{cover}} = \frac{\sum_{X_i \in R_{jb}} r_{ib}}{\sum_{X_i \in R_{jb}} [\hat{p}_{ib}(1-\hat{p}_{ib})]}$  with  $\hat{p}_{ib}$  relating to  $b-1$

<sup>3</sup>  $\varphi$  is the learning rate

<sup>4</sup> The binary function refers to whether a record ends up in multiple leaves

<sup>5</sup> The final output is nothing more than the sum over all previous models' updated predictions (simple aggregation): since the first prediction is a log-odds prediction, the final prediction will also be and can be turned to a probability prediction with  $p = \frac{e^{\ln(p/1-p)}}{1+e^{\ln(p/1-p)}}$

## 2. Choosing a predictive model

20

### 6. Extreme Gradient Boosting or XGBoost

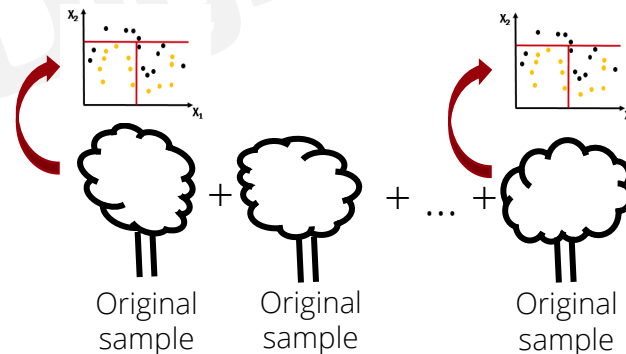
Mapping: No functional form

Loss function:

- Log-loss:  $Loss = -Y \ln(\hat{Y}) - (1 - Y) \ln(1 - \hat{Y})$

Assumption: None

Type	Non-Parametric
Target feature	Qualitative <sup>1</sup>
Distribution	No assumption
Task	Classification <sup>1-2</sup>
Loss	Log-loss
Method	Bias reduction through sequential creation of corrective models and aggregation



The idea is similar to Gradient Boosting except that the weak trees<sup>3</sup> are regularized<sup>4</sup>

<sup>1</sup> Again, only the classification task, in a binary set-up, is considered

<sup>2</sup> Extreme Gradient Boosting assumes the binary target within the space  $\{0, +1\}$

<sup>3</sup> Again, weak trees are used, due to the boosting process

<sup>4</sup> The regularization is similar to Ridge Regression

## 2. Choosing a predictive model

21

### 6. Extreme Gradient Boosting or XGBoost

Algorithm :

1. From a dataset  $D$  composed of  $N$  records
2. Initialize a first model with a constant value<sup>1</sup>:  $\hat{F}_0(X) = 0,5$
3. For each iteration  $b$  ( $b = 1, \dots, B$ ):
  - I. Calculate the residuals<sup>2</sup>  $r_{ib} = -g_{ib} = Y_i - \hat{p}_{ib} \quad i = 1, \dots, N$
  - II. Train a regression tree to the residuals  $r_{ib}$  to minimize<sup>3</sup>  $\sum_{i=1}^N L(r_{ib}, \hat{r}_{ib}) + \frac{1}{2} \lambda \gamma_b^2 + \psi T_b$  and create the leaves  $R_{jb}, j = 1, \dots, J_b$ :
    - A. Calculate the leaf output<sup>2-4</sup>  $\gamma_{jb} = \frac{-\sum_{X_i \in R_{jb}} g_{ib}}{[\sum_{X_i \in R_{jb}} h_{ib}] + \lambda} = \frac{\sum_{X_i \in R_{jb}} r_{ib}}{[\sum_{X_i \in R_{jb}} \hat{p}_{ib}(1 - \hat{p}_{ib})] + \lambda}$  for each leaf  $j = 1, \dots, J_b$

<sup>1</sup> XGBoost chooses 50% as initial model

<sup>2</sup> Expressed in probability, given the loss is the log-loss, with  $\hat{p}_{ib}$  relating to  $b - 1$

<sup>3</sup>  $\lambda$  is the shrinkage parameter,  $\gamma$  is the leaf output,  $\psi$  is the tree pruning penalty,  $T_b$  is the number of leaves in the tree

<sup>4</sup>  $\sum_{X_i \in R_{jb}} h_{ib}$  is the cover and should be set with 0 in the classification task (~min\_child\_weight)

## 2. Choosing a predictive model

22

### 6. Extreme Gradient Boosting or XGBoost

---

Algorithm :

- B. Calculating the similarity score<sup>1-2</sup>  $score_{jb} = \frac{(\sum_{X_i \in R_{jb}} r_{ib})^2}{[\sum_{X_i \in R_{jb}} \hat{p}_{ib}(1-\hat{p}_{ib})] + \lambda'}$  for each leaf  $j = 1, \dots, J_b$
  - C. Calculating the split gain  $gain = score_{R_{left}} + score_{R_{right}} - score_{R_{root}}$
  - D. Choose  $(X_k, S_k)$  in order to maximize the gain
  - E. Prune branches if<sup>3</sup>  $gain - \psi < 0$  and get the final tree
  - III. Update the predictions<sup>4-5</sup>  $\hat{F}_b(X) = \hat{F}_{b-1}(X) + \varphi \sum_j^{J_b} \gamma_{jb} \mathbb{I}(X \in R_{jb})$
  - 4. Let the last model predict the records<sup>6</sup>:  $\hat{F}(X) = \hat{F}_B(X)$
- 

<sup>1</sup> Expressed in probability, given the loss is the log-loss, with  $\hat{p}_{ib}$  relating to  $b - 1$

<sup>2</sup> The higher the  $\lambda$  (which appears with one residual in the leaf or if residuals are similar), the lower the similarity score and the lower the gain, the more likely the tree to be pruned

<sup>3</sup> The higher the pruning penalty  $\psi$ , the more likely the tree to be pruned (warning:  $\psi = 0$  doesn't turn off pruning if  $\lambda$  is involved)

<sup>4</sup>  $\varphi$  is again the learning rate

<sup>5</sup> Since the first prediction is a probability ( $\hat{F}_0(X) = 0.5$ ), we need to convert this probability to log-odds before updating:  $\ln(p/1-p)$

<sup>6</sup> The final prediction, which is again the sum over all previous models' updated predictions (simple aggregation), will be a log-odds prediction and can be then turned to a probability prediction with  $p = \frac{e^{\ln(p/1-p)}}{1 + e^{\ln(p/1-p)}}$

## 2. Choosing a predictive model

23

Recap: Key differences about tree-based ensemble

### Dependence

- Bagging: Independent trees
- Boosting: Sequential trees (~dependence)

### Bagging

- Bagging: Randomization through record sampling
- Random Forests: Randomization through record sampling and feature sampling

### Boosting

- Adaboost: Boosting based on sampling of erroneous records, weighted aggregation
- Gradient Boosting: Boosting based on shrinking pseudo residuals, aggregation
- Extreme Gradient Boosting: Gradient Boosting with regularization



# 3. Model's hyper-parameters

24

## 1. CART Decision tree

Hyper-parameter	Space	Explanation
Loss	[misclassification, gini, entropy]	Loss function used for calculating the loss over all records
Depth	Positive integer	Maximum depth of the tree: the deeper the tree, the more overfitting
Split samples	Positive integer or float	Minimum number of records for splitting a node: the larger the number, the smaller the tree and the less overfitting
Leave samples	Positive integer or float	Minimum number of records to be part a leaf: the larger the number, the smaller the tree and the less overfitting
Leave number	Positive integer	Maximum number of leaves: the smaller the number, the smaller the tree and the less overfitting
Impurity decrease	Positive float	Minimum impurity decrease for splitting a node: the larger the decrease requested, the smaller the tree and the less overfitting
$\alpha^1$	Positive float	Amount of penalization: the larger the amount, the smaller the tree and the less overfitting

<sup>1</sup> CART decision tree also exists with a penalized form which accounts for the size of tree

# 3. Model's hyper-parameters

25

## 2. Bootstrap Aggregating or Bagging

Hyper-parameter	Space	Explanation
All CART tree's hyper-parameters are valid for Bagging		
Tree number	Positive integer	Number of trees: the more trees, the larger the reduction in variance and the less overfitting <sup>1</sup>

<sup>1</sup> In practice, however, it is not always the case so it is advised to perform hyper-parameter tuning

# 3. Model's hyper-parameters

26

## 3. Random Forests

Hyper-parameter	Space	Explanation
All CART tree's hyper-parameters are valid for Random Forests		
Tree number	Positive integer	Number of trees: the more trees, the larger the reduction in variance and the less overfitting <sup>1</sup>
Feature number	Positive integer or float	Maximum number of features for splitting a node: the smaller the number, the more tree decorrelation

<sup>1</sup> In practice, however, it is not always the case so it is advised to perform hyper-parameter tuning

# 3. Model's hyper-parameters

27

## 4. Adaptive Boosting or AdaBoost

Hyper-parameter	Space	Explanation
All CART tree's hyper-parameters are valid for Adaboost		
Loss <sup>1</sup>	[exponential]	Loss function used for calculating the loss over all records
Tree number	Positive integer	Number of trees: the more trees, the larger the reduction in bias but the more overfitting <sup>2</sup>

<sup>1</sup> All CART decision tree's hyper-parameters are valid for Adaboost except for the loss

<sup>2</sup> As the number of trees increases, the algorithm keeps trying to maximize the margin locally and continuously correct the errors which eventually leads to overfitting

# 3. Model's hyper-parameters

28

## 5. Gradient Boosting

Hyper-parameter	Space	Explanation
All CART tree's hyper-parameters are valid for Gradient Boosting		
Loss <sup>1</sup>	[exponential, logistic, log]	Loss function used for calculating the loss over all records
Tree number	Positive integer	Number of trees: the more trees, the larger the reduction in bias but the more overfitting <sup>2</sup>
Learning rate	Positive float	Step size of the coefficients at each iteration while moving to the minimum of the loss function: the lower the rate, the slower the convergence but the more chance to converge to the global optimum

<sup>1</sup> All CART decision tree's hyper-parameters are valid for Gradient Boosting except for the loss

<sup>2</sup> As the number of trees increases, the algorithm keeps trying to boost the errors which eventually leads to overfitting



# 3. Model's hyper-parameters

## 6. Extreme Gradient Boosting or XGBoost

Hyper-parameter	Space	Explanation
All CART tree's hyper-parameters are valid for Extreme Gradient Boosting		
Loss <sup>1</sup>	[log]	Loss function used for calculating the loss over all records
Tree number	Positive integer	Number of trees: the more trees, the larger the reduction in bias but the more overfitting <sup>2</sup>
Learning rate	Positive float	Step size of the coefficients at each iteration while moving to the minimum of the loss function: the lower the rate, the slower the convergence but the more chance to converge to the global optimum
l1/l2 penalty	[0, +∞[	Amount of penalization: the larger the amount, the more the model is regularized, the less overfitting

<sup>1</sup> All CART decision tree's hyper-parameters are valid for XGBoost except for the loss

<sup>2</sup> As the number of trees increases, the algorithm keeps trying to boost the errors which eventually leads to overfitting



# Annexes

# Annex 1 – Loss functions

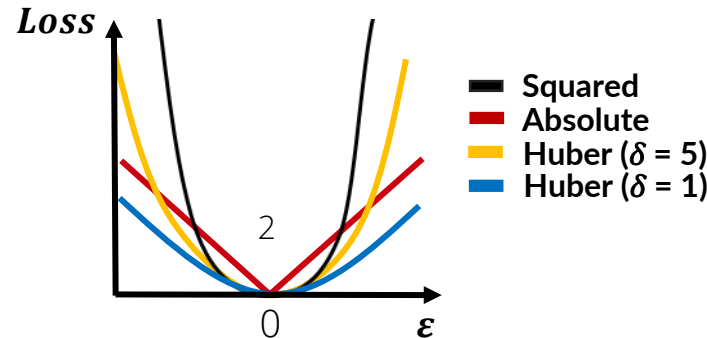
31

Loss functions differ depending on the nature of the predictive task.

Regression task

Loss functions include:

- Squared loss function:  $Loss = (Y - \hat{Y})^2$
- Absolute loss function:  $Loss = |Y - \hat{Y}|$
- Huber loss function<sup>1</sup>:  $Loss = \begin{cases} \frac{1}{2}(Y - \hat{Y})^2 & \text{if } |Y - \hat{Y}| \leq \delta \\ \delta|Y - \hat{Y}| - \frac{1}{2}\delta^2 & \text{otherwise} \end{cases}$



<sup>1</sup>  $\delta$  is the tuning parameter used to determine which data point is considered as an outlier

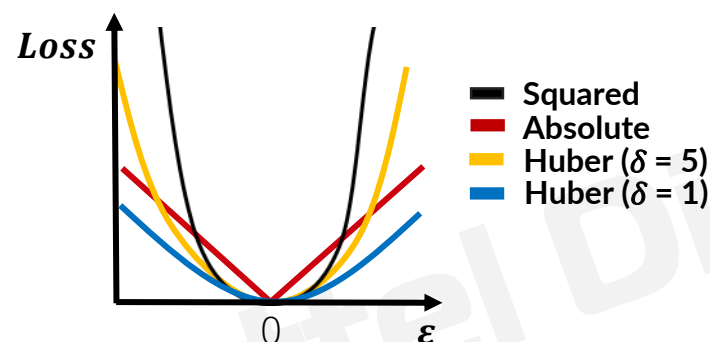
<sup>2</sup> Derivability and convexity are nice properties for a function to have because respectively that's easier at optimizing and the minimum found is the global one

# Annex 1 – Loss functions

32

Loss functions differ depending on the nature of the predictive task.

Regression task



Loss function	Optimization <sup>1-2-3</sup>	Comments
Squared loss	Continuous, differentiable and convex	Penalizing more heavily large errors and thus more sensitive to outliers
Absolute loss	Continuous but not differentiable at $x=0$ , convex	Penalizing less heavily large errors and thus more robust to outliers
Huber loss	Continuous, differentiable and convex	Combination of squared loss (when $\delta = +\infty$ ) and absolute loss (when $\delta = 0$ ) to determine the level of sensitiveness to outliers

<sup>1</sup> A function is continuous if we can ensure arbitrarily small changes by triggering small changes in the input

<sup>2</sup> A real-valued function is differentiable if it is differentiable at every point in its domain – A function that is differentiable at  $x$  is by default continuous at  $x$  but the contrary doesn't hold true

<sup>3</sup> A real-value function is convex if the line segment between any 2 points lies above the graph between these 2 points

# Annex 1 – Loss functions

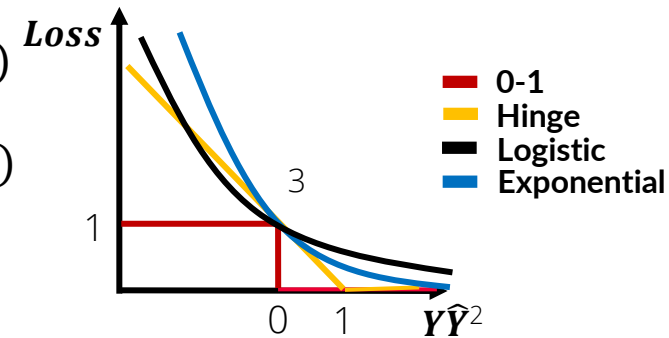
33

Loss functions differ depending on the nature of the predictive task.

Classification task<sup>1</sup>

Loss functions include:

- Misclassification or 0-1 loss function<sup>1</sup>:  $Loss = \mathbb{I}(Y \neq \hat{Y})$   $\hat{Y} \rightarrow class$
- Hinge loss function<sup>2</sup>:  $Loss = \max(0; 1 - Y\hat{Y})$   $\hat{Y} \rightarrow X\hat{\beta} (< or > 0)$
- Squared Hinge loss function<sup>2</sup>:  $Loss = (\max(0; 1 - Y\hat{Y}))^2$   $\hat{Y} \rightarrow X\hat{\beta} (< or > 0)$
- Log-loss function:  $Loss = -Y \ln(\hat{Y}) - (1 - Y) \ln(1 - \hat{Y})$   $\hat{Y} \rightarrow \hat{p}$
- Logistic or Binomial Deviance loss function:  $Loss = \ln(1 + e^{\hat{Y}\hat{Y}})$   $\hat{Y} \rightarrow \hat{p}$
- Exponential loss function (Adaboost)<sup>2</sup>:  $Loss = e^{-Y\hat{Y}}$   $\hat{Y} \rightarrow weighted\ sum\ of\ -1/1 (< or > 0)$



<sup>1</sup> Only the binary case is considered here with the target feature coded to  $\{-1, +1\}$ , which is important for the concept of “margin”

<sup>2</sup>  $Y\hat{Y}$  refers to the term “margin”: Correctly predicted records have positive margin while incorrectly predicted records have negative margin (link with SVC and Adaboost)

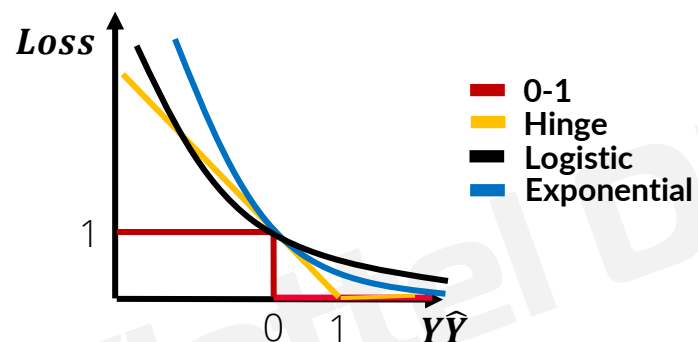
<sup>3</sup> Derivability and convexity are nice properties for a function to have because respectively that's easier at optimizing and the minimum found is the global one

# Annex 1 – Loss functions

34

Loss functions differ depending on the nature of the predictive task.

**Classification task**



Loss function	Optimization <sup>1-2-3</sup>	Comments
<b>0-1 loss</b>	Not continuous and not differentiable at $x=0$ , not convex	- Giving unit penalty for negative margins and no penalty for positive margins - Robust to outliers
<b>Hinge loss</b>	Continuous but not differentiable at $x=1$ , convex	- Penalizing continuously small margin and increasingly negative margin while giving no penalty if enough margin
<b>Logistic loss</b>	Continuous, differentiable and convex	- Penalizing continuously decreasingly negative margin more heavily than rewarding increasingly positive margin - Penalizing less heavily negative margin and thus more robust to outliers
<b>Exponential loss</b>	Continuous, differentiable and convex	- Penalizing continuously decreasingly negative margin more heavily than rewarding increasingly positive margin - Penalizing more heavily negative margin and thus more sensitive to outliers

<sup>1</sup> A function is continuous if we can ensure arbitrarily small changes by triggering small changes in the input

<sup>2</sup> A real-valued function is differentiable if it is differentiable at every point in its domain – A function that is differentiable at  $x$  is by default continuous at  $x$  but the contrary doesn't hold true

<sup>3</sup> A real-value function is convex if the line segment between any 2 points lies above the graph between these 2 points



# Annex 2 – Bias and variance trade-off

35

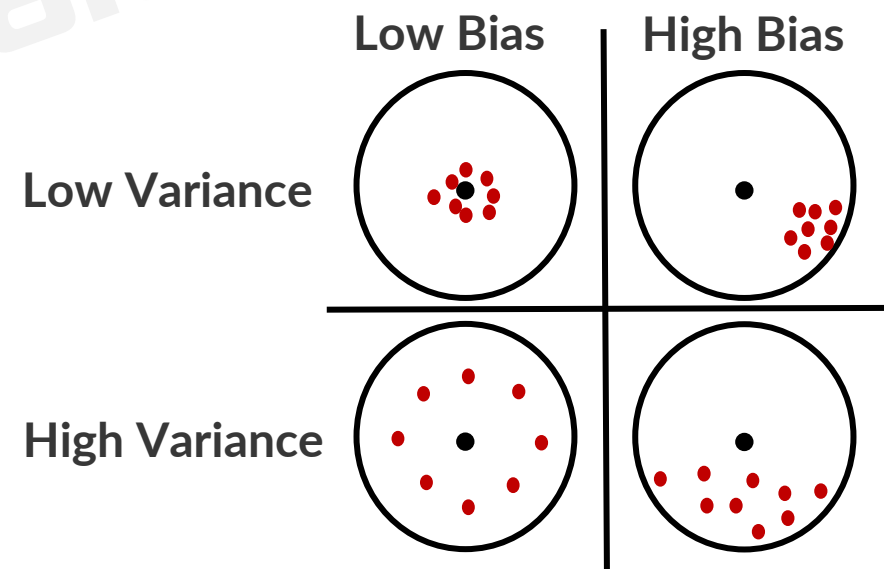
While approximating the true mapping function  $F$  with  $\hat{F}$ , predictive models make errors.

The generalization error in the case of the quadratic loss is<sup>1-2-3</sup>:

$$ERR(\hat{F}) = E \left[ [y - \hat{F}(X)]^2 \right] = Bias(\hat{F})^2 + Var(\hat{F}) + irreducible\ error$$

With:

- $Bias(\hat{F}) = E[\hat{F}(X) - F(X)] = E[\hat{F}(X)] - F(X)$
- $Var(\hat{F}) = E \left[ [\hat{F}(X) - E[\hat{F}(X)]]^2 \right]$



<sup>1</sup> Bias is the error by how much the average prediction differs from the reality: it refers to erroneous assumptions or decisions about the data which affect the representation of results

<sup>2</sup> Variance is the error by how much the prediction fluctuates around the average prediction

<sup>3</sup> Error which cannot be reduced whatever the quality of the model, it is a measure of the amount of noise in the data

# Annex 2 – Bias and variance trade-off

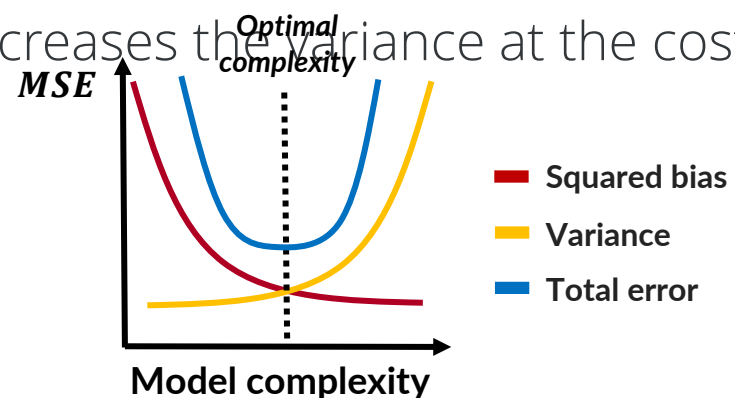
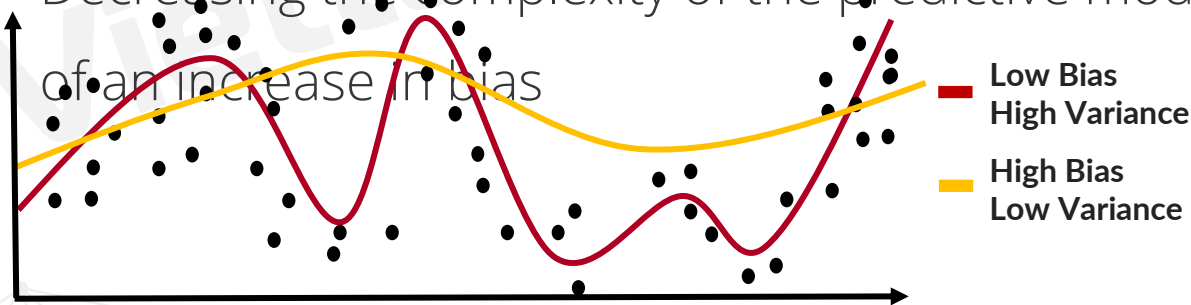
36

If bias is high,  $\hat{F}$  is unable to capture the underlying pattern from the training data **UNDERFITTING**

If variance is high,  $\hat{F}$  captures too much the noise from the training data **OVERFITTING**

Bias and variance are 2 quantities that compete against each other and are part of a trade-off<sup>1</sup>:

- Increasing the complexity of the predictive model decreases the bias at the cost of an increase in variance
- Decreasing the complexity of the predictive model decreases the variance at the cost



<sup>1</sup> Model will be chosen and its complexity will be tuned with hyper-parameters in order to minimize the generalization loss

# Annex 3 – Bootstrapping

37

Bootstrapping is a resampling method used to estimate the sampling distribution<sup>1</sup> of almost any statistic.

Given  $\theta$  the unknown parameter of interest,  $T$  the related sample's statistic,  $D$  a sample of  $N$  independent and identically distributed records<sup>2</sup>,  $T(D)$  is an estimator of  $\theta$  and depends on  $D$ . The procedure consists of :

- **Bootstrap sampling:** Drawing randomly  $B$  samples  $D_b$  of size  $N$  with replacement (~phantom samples) from  $D$
- **Bootstrap's statistic:** Calculating  $T$  for each bootstrap sample  $D_b$ ,  $T(D_b)$
- **Bootstrap distribution:** Collecting all the statistic's values across the bootstrap samples and obtaining the bootstrap distribution of  $T$  which approximates the sampling distribution of  $T$

<sup>1</sup> The population parameter being unknown, we want to estimate it from a sample's statistic; as any statistic of interest fluctuate across samples, we want to know the magnitude of those fluctuations and, as a result, want to know the sampling distribution of that statistic and its standard deviation (= standard error)

<sup>2</sup> Independent records: records are unrelated to each other – Identically distributed records: records come from the same probability distribution with same parameters