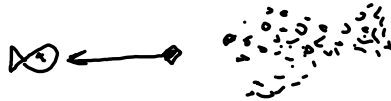


(CART alg.)
training

DECISION
TREES



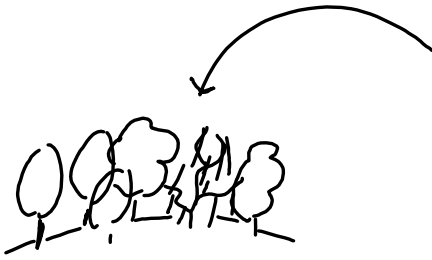
improvement
ENSEMBLE LEARNING



sequential
BOOSTING



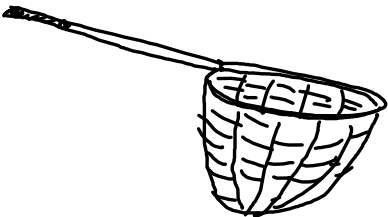
total error



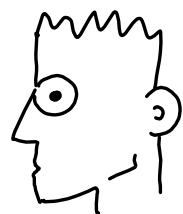
RANDOM FORESTS



non sequential
BOOSTING



BART model

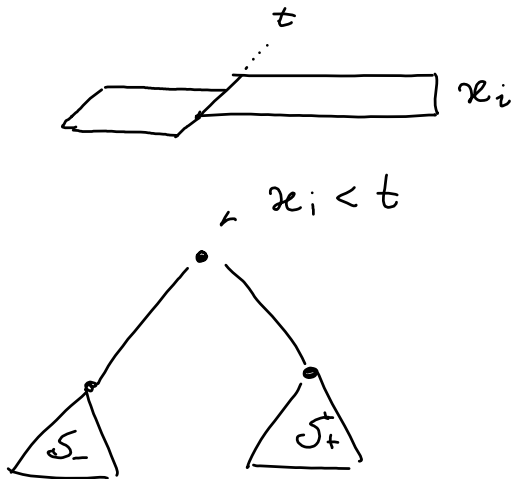


DECISION TREES

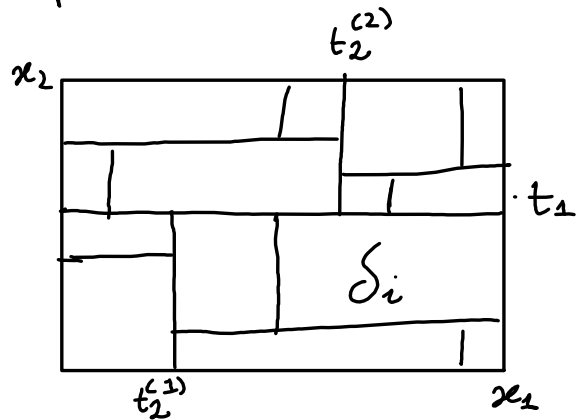
x_1, \dots, x_m features

$y = \{C_1, \dots, C_m\}$. class

Recursive Definition



Each split corresponds to a subdivision of the feature space:

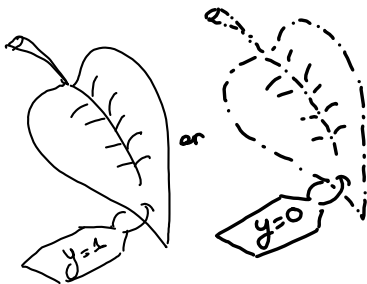


At each level we try to find best pair of (x_i, t) to divide the data in order to:

- minimize RSS (REGRESSION)
- maximize IG (CLASSIFICATION)

Prediction formula: In every hyperplane S_i of the feature space the same formula for prediction is used.

Final leaves



$$f_S := \begin{cases} \frac{1}{|S|} \sum_{i \in S} y_i = \bar{y}_S & \text{in Regression} \\ \max_k P(y \in C_k | x \in S) & \text{in Classif.} \end{cases}$$

LIKE: infected tree.

DECISION TREES (REGRESSION)

$$f_S = \bar{y}_S = \frac{1}{|S|} \sum_{y \in S} y \quad (\text{prediction of } S \text{ by p. in feature space})$$

GOAL FUNCTION:

$$\min_J \underbrace{\sum_{j=1}^J \sum_{i \in S_j} (y_i - \bar{y}_{S_j})^2}_{RSS} \quad J = \# \text{ of Subsets } S \text{ we find.}$$

The goal is finding the best set of S_1, \dots, S_J to divide the feature space, so that the predictions we make minimize the goal function.

SOLUTION:

Top-down building of the tree.

At each level we seek the best pair (X_j, t) such that splitting the feature space into:

$$\left[\begin{array}{c} X | X_j < t \\ \hline t \\ \hline X | X_j > t \end{array} \right] X_j$$

Leads to biggest reduction in RSS.

DECISION TREES (CLASSIFICATION)

Mathematical form of the optimization problem.

the threshold t is found as the value that, each split, maximizes the information gain:

$$IG_t(S) = H(S) - \sum_{i=1}^2 \frac{|S_i|}{|S|} H(S_i)$$

This means that at every split we seek the t that divides the training set in the most balanced way:

$$\min IG_t(S) = \max \sum_{i=1}^2 \frac{|S_i|}{|S|} H(S_i)$$

where $S_r := \{ \text{records in the root} \}$

Eventually the IG we can achieve by splitting the node will be very small.

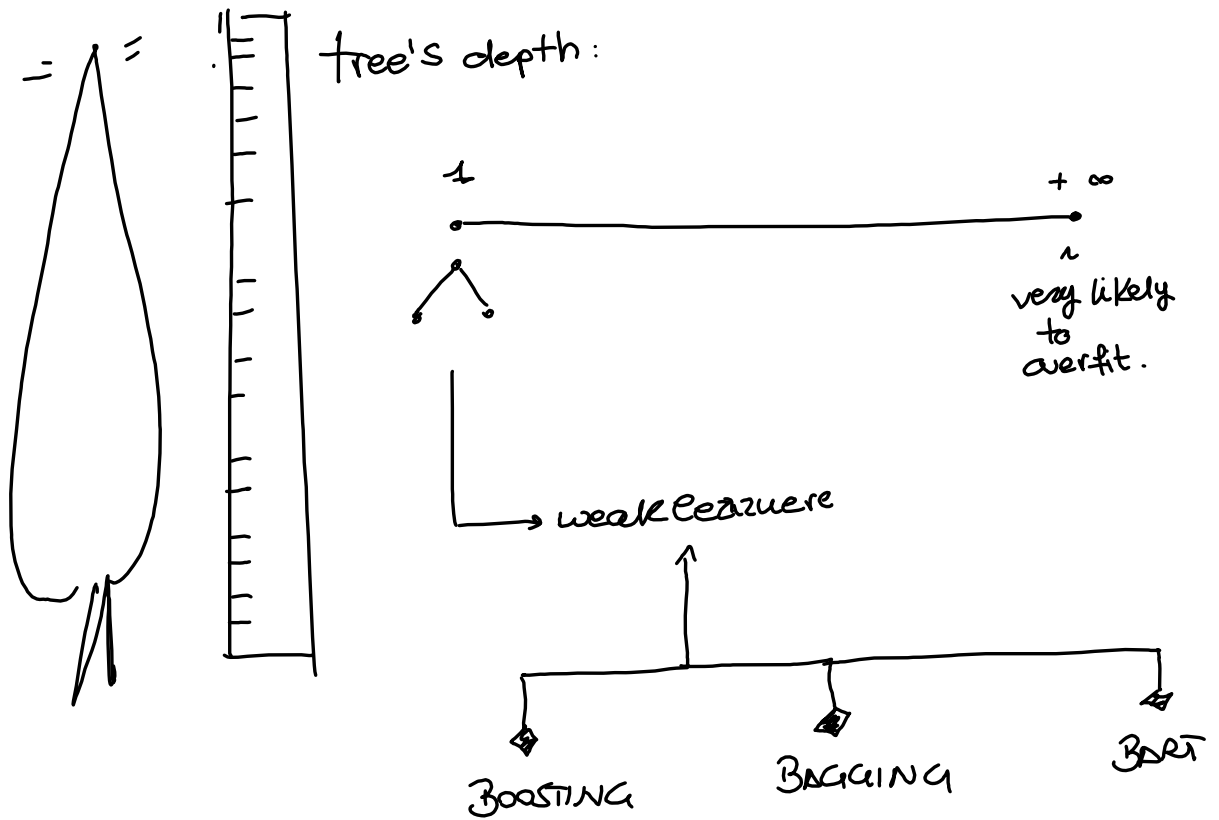
$$|IG| < \epsilon$$

This means that the leaf we arrived at during the splitting process cannot be divided into 2 balanced subsets of records. So the leaf is already very imbalanced: its records belong mostly to the same class. So:

$$H(S_r) = - \sum_{i=1}^m \phi_i \ln \phi_i \approx -1 \ln 1 \approx 0$$

DECISION TREES

Hyperparameters

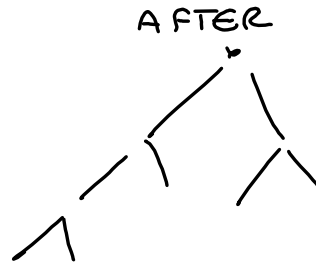
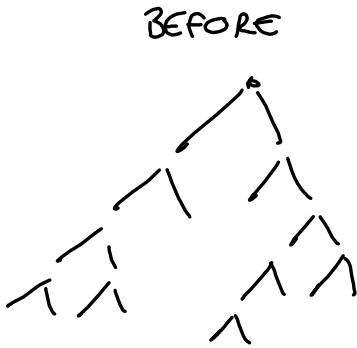


NOTE

For how they work, decision trees do NOT demand any kind of feature scaling. In fact features (X) are analyzed singularly and there's no direct comparison between them of any kind. This means that relative scales don't matter.

DECISION TREES (PRUNING)

PROBLEM: decision trees tend to overfit the data. They become too complex and the feat. space (X) subd. too specific.



Pruning is the operation of undoing the splits (rejoining 2 subspaces of the feature space) that grant the less advancement wrt the optimization problem.



DEAD BRANCH:

- PRO:
- reducing overfitting tendency.
- not ↓ RSS (regression).
- not ↑ T1C (classif.)

DECISION TREES in Sklearn

The algorithm that is used to train the model in Python Sklearn is called Cart. It works both for classification and regression trees, by changing the cost function.

At every split considered we try to find the best pair (k, t_k) that minimizes a problem tailored cost function.

CLASSIFICATION (cost: Gini impurity ; $\varphi_i(x) = \max_j P_i(y \in C_j | x \in S_i)$)

$$\min_{k, t} J(k, t_k) = \min_{k, t} \frac{m_{\text{left}}}{m} G_1 + \frac{m_{\text{right}}}{m} G_2$$

REGRESSION (cost: MSE ; $\varphi_i(x) = \frac{1}{|S_i|} \sum_{y \in S_i} y$)

$$\min_{k, t} J(k, t_k) = \min_{k, t} \frac{m_{\text{left}}}{m} \text{MSE}_{\text{left}} + \frac{m_{\text{right}}}{m} \text{MSE}_{\text{right}}$$

where $\left\{ \begin{array}{l} G_1, G_2 \text{ measures the impurity of the left, right subset} \\ m_{\text{left/right}} \text{ is the no of instances in the left/right subset.} \\ k \text{ is a feature (from } x) \text{ and } t_k \text{ the splitting threshold.} \end{array} \right.$

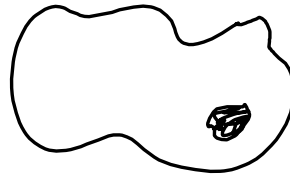
① DECISION TREES

Give impurity or entropy?

Let's assume we have: m different classes to choose from: C_1, \dots, C_m

Which is the best cost function for training a classification tree?

$$G_{S_i} = 1 - \sum_{k=1}^m p_{i,k}^2$$



$$H(S_i) = - \sum_{\substack{k=1 \\ p_{i,k} > 0}}^m p_{i,k} \log_2(p_{i,k})$$

empirical distr.



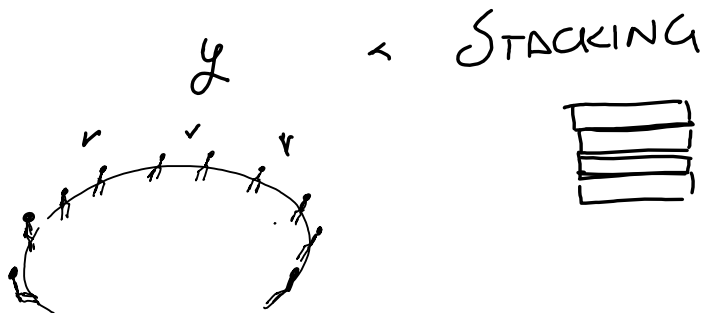
$$\text{where: } p_{i,k} = \frac{\# (x, y) \in S_i \text{ so that } y = C_k}{\# (x, y) \in S_i}$$

ANSWER: there's no big difference in the result. Give impurity is a little more computationally efficient.

$$\begin{aligned} G_{S_i} &= 0 \\ H(S_i) &= 0 \end{aligned} \iff \text{All instances in } S_i \text{ have the same label.}$$

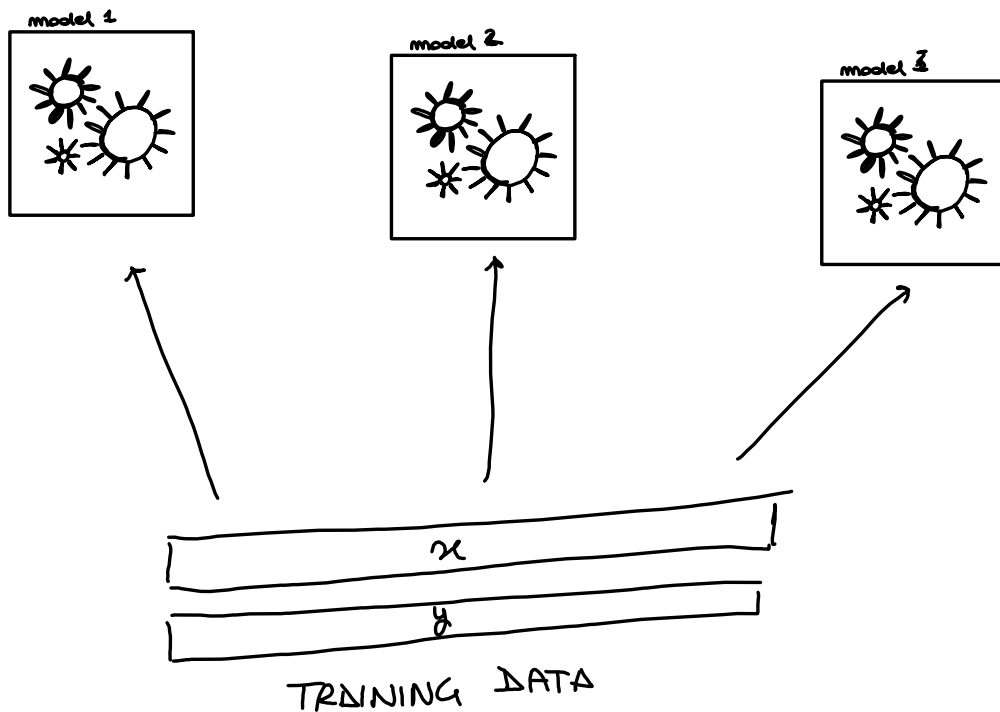
ENSEMBLE LEARNING

Class of techniques that uses the wisdom of the crowd effect

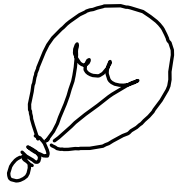


(ex: mean, mode, majority vote...)

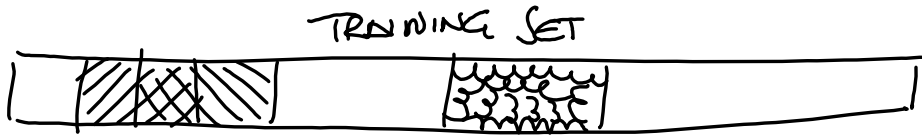
ENSEMBLE PREDICTION



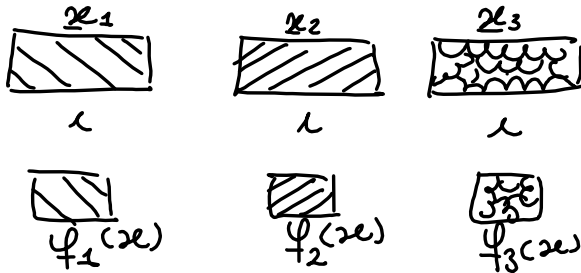
BAGGING ..



model that are good will tend to agree on the same prediction, while bad model tend to disagree on different predictions



B bootstrap samples from the training data



train a weak model on each one of them



prediction happens through rotation of weak model



$$y \leftarrow f(x) = \frac{1}{B} \sum_{i=1}^B f_i(x)$$

BAGGING

theory.

Bagging is a method used for reducing the variance of a statistical learning method.

$$f(x) = \frac{1}{B} \sum_{i=1}^B f_i(x)$$

The model prediction is computed as average output of B other models, trained on Bootstrap sets.

In fact Bagging stands for Bootstrap Aggregation.

Since we output an avg., from the CLT:

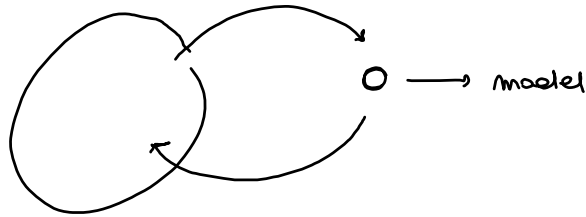
$$f(x) \sim \mathcal{N}(\bar{f}(x) = y; \frac{V(f(x))}{B})$$

As B increases the final model variance decreases. This is why we operate bagging with high variance learning methods.

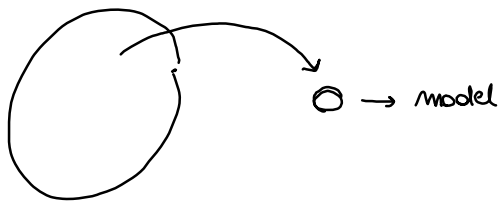
Bagging improves prediction accuracy at the expenses of readability.

BAGGING and PASTING and Out of bag evaluation (OOB)

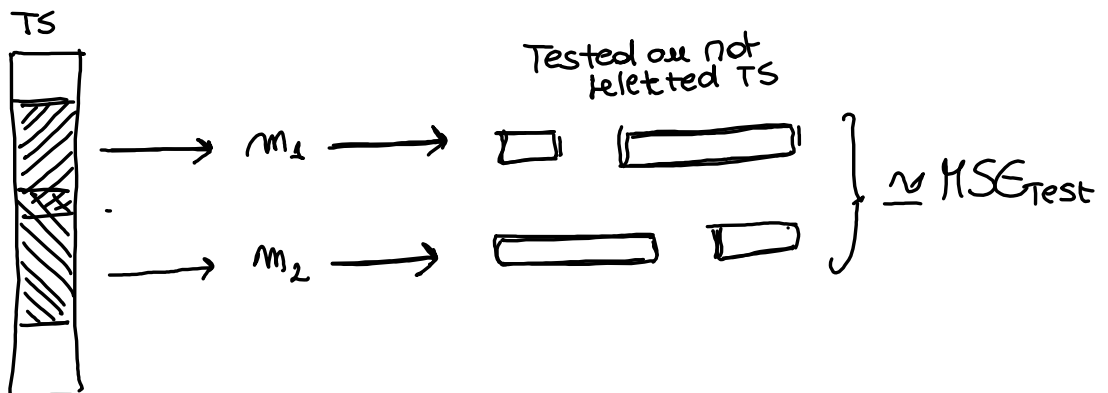
BAGGING: performs sampling with replacement



PASTING: performs sampling without replacement

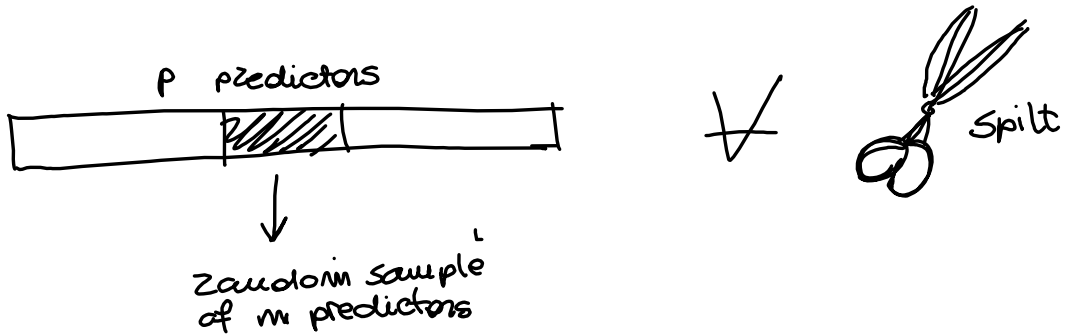


Out of BAG evaluation: since every model is trained on a random subset of the training data, the training instances not used for training can be used as test instances:



RANDOM FORESTS

We apply a bagging procedure to a decision tree. Each time a split in a tree is considered, a random sample of m predictors is chosen as split candidates.



Avoiding collinearity
between bagged trees

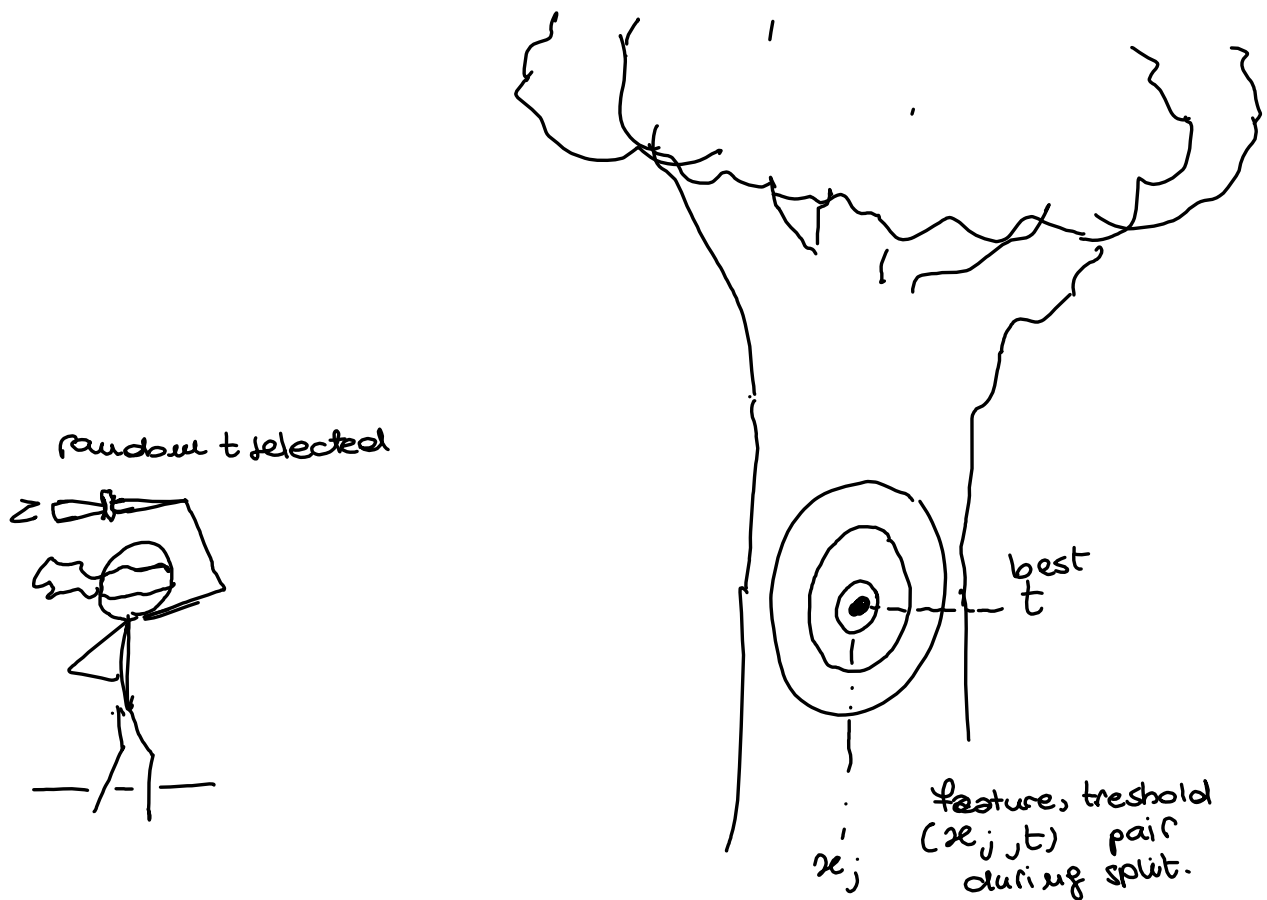


Having correlated trees doesn't lead to a reduction in variance as having many un-correlated trees. For this reason RANDOM FORESTS represent a quick improvement of bagged trees

NOTE: usually $m \approx \sqrt{p}$

EXTREMELY RANDOMIZED TREES

Like random forests, but also the threshold t at each split during training is randomly selected

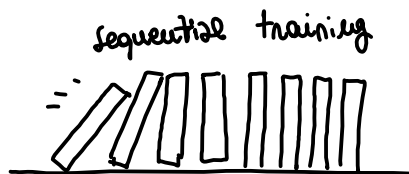


PRO: • training is faster than Random Forests.

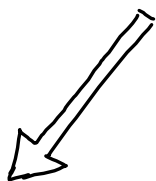
There's no way of knowing if ERT will fit the data better than Random Forests beforehand.

BOOSTING

overview

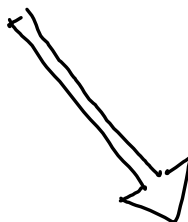


BOOSTING



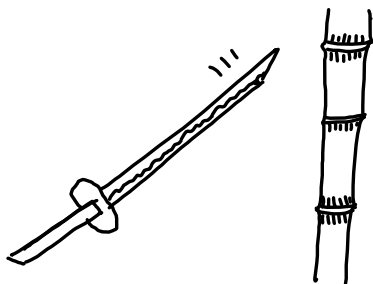
training
on residuals

GRADIENT BOOSTING



+ on
weighted
training set

ADA = BOOST



$$f(x) = \alpha^T \underline{f(x)}$$

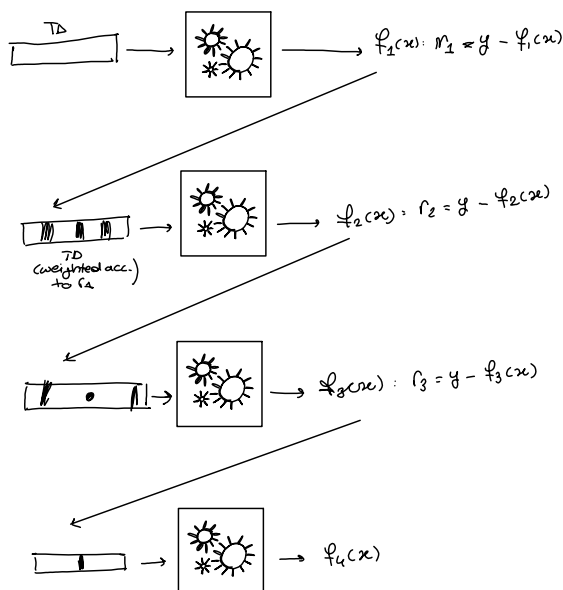


$$f(x) = h(\alpha_j)$$

model weights.

BOOSTING: Gradient Boosting

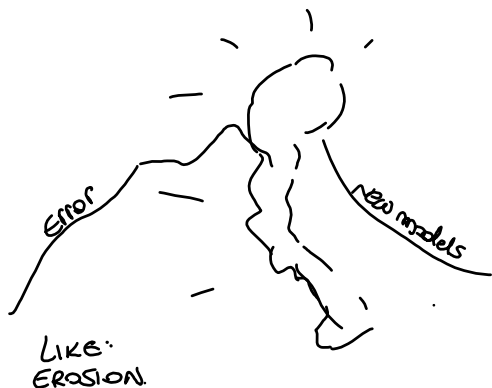
Create sequential models to predict the errors the other ones make.



boosted model:
$$f(x) = \sum_{i=1}^n \alpha f_i(x) = \alpha^T \underline{f}(x)$$

Each new model is trained having as labels the error of the previous ones:

$$y^{(k)} = y - \sum_{i=1}^{k-1} f_i(x)$$



NOTE: Similarity to Grad. Desc.
At each step the overall error is implicitly reduced by training a new mod. on the past error.

STEP UPD.) $x_{k+1} = x_k - \alpha \nabla f(x_k)$

RESID. UPD.) $y_{k+1} = y_k - \alpha f_k(x)$

For this reason, Boosting is also referred as Gradient Boosting.

BOOSTING: Gradient Boosting Hyperparameters

1) $B :=$ no of models to train



Chosen through cross validation

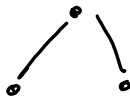
2) $\alpha :=$ learning rate (step of GD)

Typical values: $\alpha = 10^{-2}, 10^{-3}$

3) $d :=$ no of splits in each tree

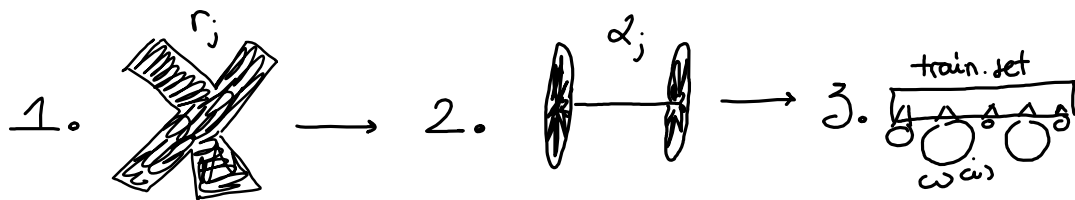
This controls the complexity of the boosted ensemble.

Often: $d = 1$



BOOSTING

AdaBoost basic operations. Models are sequentially giving each time more importance to « hard » training instance



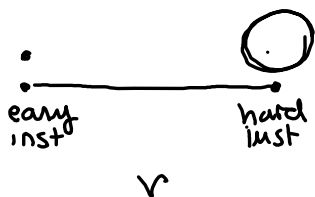
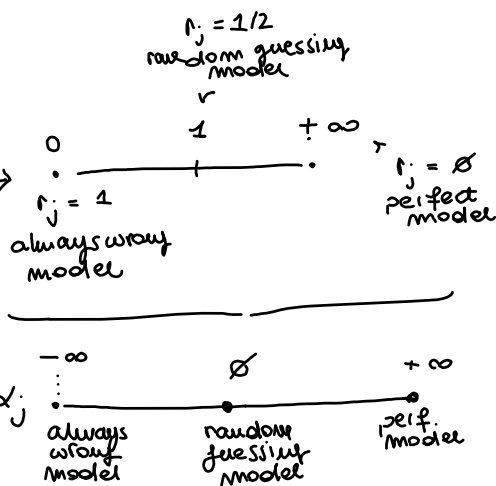
\forall model sequentially trained :

1. weighted error rate of the j -th model

$$r_j := \frac{\sum_{i=1}^m \omega_i \mathbb{1}_{\hat{y}_j^{(i)} \neq y^{(i)}}}{\sum_{i=1}^m \omega_i}$$

2. generating weight of the j -th model

$$\alpha_j := \eta \log \frac{1 - r_j}{r_j}$$



3. weight updating

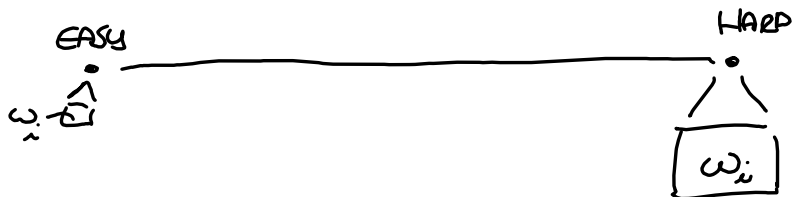
$$\omega^{(i)} = \begin{cases} \omega^{(i)} & \text{correct prediction: } \hat{y}_j^{(i)} = y^{(i)} \\ \omega^{(i)} e^{\alpha} & \text{wrong pred.: } \hat{y}_j^{(i)} \neq y^{(i)} \end{cases}$$

BOOSTING: AdaBoost

example algorithm: each new model is trained on data that are weighted by how hard they find for the past models to get right

Training data are weighted:

$$\text{NOTE: } y_i = \begin{cases} 1 \\ -1 \end{cases}$$



AdaBoost:

FOR $t = 1, \dots, T$:

- choose weak classifier:

- find weak learner $h_t(x)$ that minimize the weighted sum of errors: $\epsilon_t = \sum_i w_{i,t} \epsilon_{i,t} \leq 0.5$

- choose $\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$

- ensemble:

$$F_t(x) = F_{t-1}(x) + \alpha_t h_t(x) = h_1(x) + h_2(x) + \dots$$

- update weights:

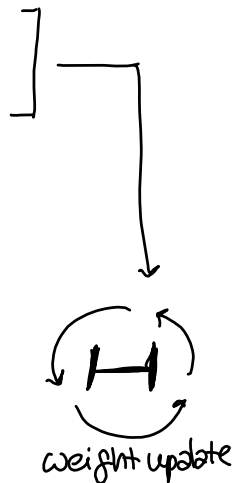
$$w_{i,t+1} = (w_{i,t}) e^{-y_i \alpha_t h_t(x_i)} \quad \forall i$$

- Renormalize weights.

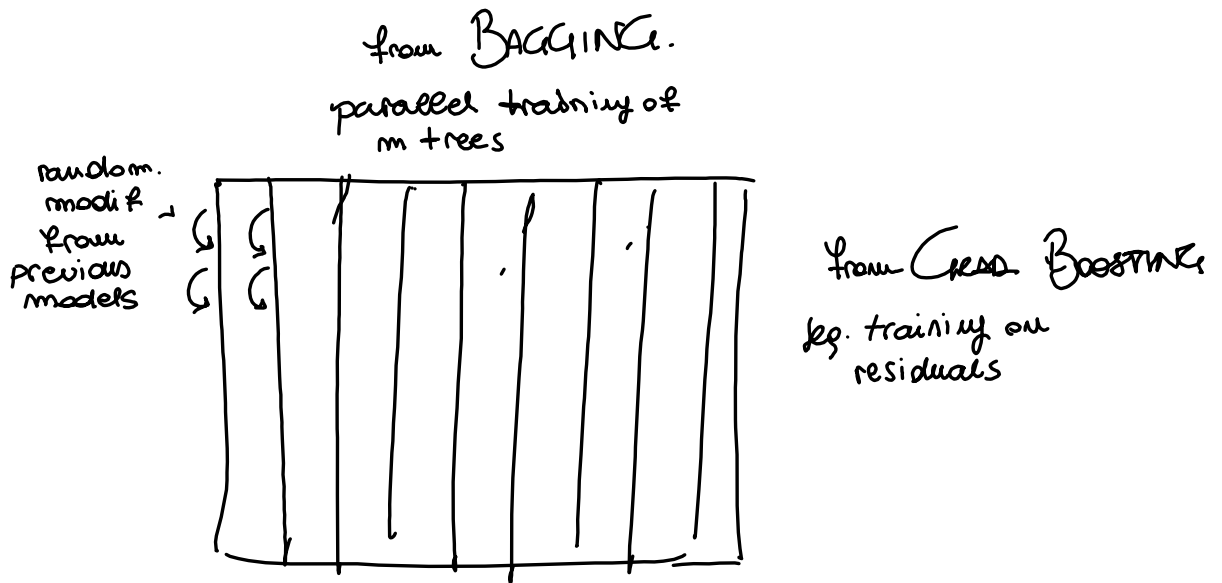
normalize weights.



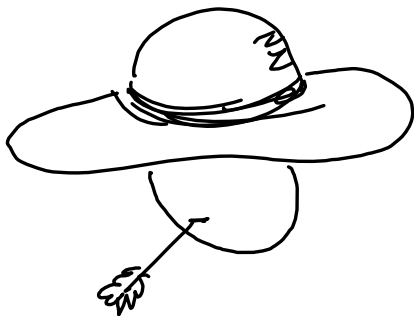
$$-y_i \alpha_t h_t(x_i) = \begin{cases} > 0 & \text{WRONG PRED. } h_t(x_i) \neq y_i \rightarrow \uparrow w_i \\ < 0 & \text{RIGHT PRED. } h_t(x_i) = y_i \rightarrow \downarrow w_i \end{cases}$$



BART: visualization and influences.



BART



BART (For regression)

Bayesian Additive Regression trees

K = # regression trees

B = # iterations

$f_k^{(b)}(x)$ = prediction at x of the k -th regression tree at iteration b

n = # train. records

$$1. f_1^{(1)}(x) = \dots = f_K^{(1)}(x) = \frac{1}{nK} \sum_{i=1}^n y_i \quad \left. \vphantom{\sum_{i=1}^n y_i} \right] \text{like bagging}$$

$$2. \varphi^{(1)}(x) = \frac{1}{n} \sum_{i=1}^n f_i^{(1)}(x)$$

3. for $b = 2, \dots, B$

2) for $k = 1, \dots, K$:

I) for $i = 1, \dots, n$ compute partial residual $\left. \vphantom{\sum_{k' \neq k} f_k^{(b)}(x)} \right] \text{like boosting}$

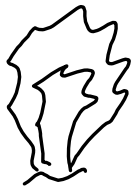
$$r_i = y_i - \sum_{k' \neq k} f_{k'}^{(b)}(x)$$

II) Fit a new tree $\varphi_k^{(b)}(x)$ by rand. perturbing $\varphi_k^{(b-1)}(x)$.

$$b) \text{ compute } f^{(b)}(x) = \sum_{k=1}^K \varphi_k^{(b)}(x)$$

4. Compute the mean after L burn-in samples $\left. \vphantom{\sum_{b=L+1}^B \varphi^{(b)}(x)} \right] \text{like}$

$$\varphi(x) = \frac{1}{B-L} \sum_{b=L+1}^B \varphi^{(b)}(x)$$



In words:

1, 2: initialization of K regression trees. $\varphi^{(1)}(x)$ is computed as avg output of the K trees

3: B iterations. At each iteration everyone of the K trees is updated from its previous version, by a random modification

4: the model is an avg. of the not-burnt prediction models.

BART (for regression)

key insights

3. At this step we do not fit a new tree to the current partial residual error but a modification of the previous one

- AVOID overfitting: improving the older tree prevents from fitting «hard» the data at each iteration

4. The burn-in period (the first L iterations) corresponds to the number of model we do not consider in the final version model. Generally the first models found $(f^{(1)}(x), \dots, f^{(L)}(x))$ tend to not provide very good results.

4. As in bagging the final model is an average of the best models we have found:

$$\begin{aligned} f(x) &= \frac{1}{B-L} \sum_{i=L+1}^B f^{(i)}(x) \\ &= \frac{1}{B-L} \sum_{i=L+1}^B \sum_{k=1}^K f_k^{(i)}(x) \end{aligned}$$

Having 2 avg in the final model, maximally leverages the variance reduction effect of the mean.

As CLT states:

$$\bar{X} \sim \mathcal{N}(E[\tilde{X}], \frac{V[\tilde{X}]}{n})$$

BART

hyperparameters setting

For how it is thought: the BART:

- prevents overfitting
- is good out of the box

Standard value for its hyperparameters are:

$K = 200$ (no of trees)

$B = 1000$ (no of iterations)

$L = 100$ (burn-in period)

STACKING.

Instead of using a voting aggregation function like the mean or majority, we can train a model to aggregate votes.

