

# Data Analysis and Machine Learning: From Decision Trees to Forests and all that

Morten Hjorth-Jensen<sup>1,2</sup>

<sup>1</sup>Department of Physics, University of Oslo

<sup>2</sup>Department of Physics and Astronomy and National Superconducting Cyclotron Laboratory, Michigan State University

May 26, 2020

## Decision trees, overarching aims

We start here with the most basic algorithm, the so-called decision tree. With this basic algorithm we can in turn build more complex networks, spanning from homogeneous and heterogeneous forests (bagging, random forests and more) to one of the most popular supervised algorithms nowadays, the extreme gradient boosting, or just XGBoost. But let us start with the simplest possible ingredient.

Decision trees are supervised learning algorithms used for both, classification and regression tasks.

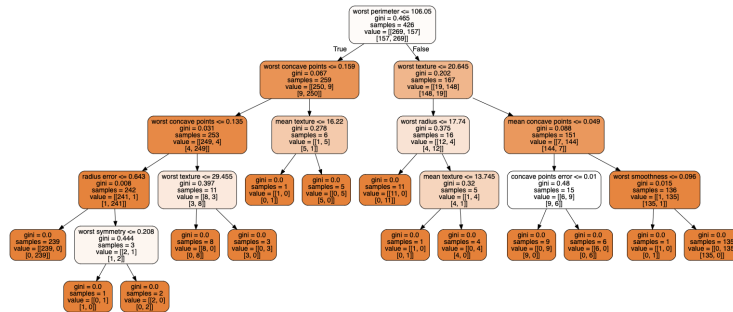
The main idea of decision trees is to find those descriptive features which contain the most **information** regarding the target feature and then split the dataset along the values of these features such that the target feature values for the resulting underlying datasets are as pure as possible.

The descriptive features which reproduce best the target/output features are normally said to be the most informative ones. The process of finding the **most informative** feature is done until we accomplish a stopping criteria where we then finally end up in so called **leaf nodes**.

A decision tree is typically divided into a **root node**, the **interior nodes**, and the final **leaf nodes** or just **leaves**. These entities are then connected by so-called **branches**.

The leaf nodes contain the predictions we will make for new query instances presented to our trained model. This is possible since the model has learned the underlying structure of the training data and hence can, given some assumptions, make predictions about the target feature value (class) of unseen query instances.

## A typical Decision Tree with its pertinent Jargon, Classification Problem



This tree was produced using the Wisconsin cancer data (discussed here as well, see code examples below) using **Scikit-Learn**'s decision tree classifier. Here we have used the so-called **gini** index (see below) to split the various branches.

## General Features

The overarching approach to decision trees is a top-down approach.

- A leaf provides the classification of a given instance.
- A node specifies a test of some attribute of the instance.
- A branch corresponds to a possible values of an attribute.
- An instance is classified by starting at the root node of the tree, testing the attribute specified by this node, then moving down the tree branch corresponding to the value of the attribute in the given example.

This process is then repeated for the subtree rooted at the new node.

## How do we set it up?

In simplified terms, the process of training a decision tree and predicting the target features of query instances is as follows:

1. Present a dataset containing of a number of training instances characterized by a number of descriptive features and a target feature
2. Train the decision tree model by continuously splitting the target feature along the values of the descriptive features using a measure of information gain during the training process
3. Grow the tree until we accomplish a stopping criteria create leaf nodes which represent the *predictions* we want to make for new query instances

4. Show query instances to the tree and run down the tree until we arrive at leaf nodes

Then we are essentially done!

## Decision trees and Regression

### Building a tree, regression

There are mainly two steps

1. We split the predictor space (the set of possible values  $x_1, x_2, \dots, x_p$ ) into  $J$  distinct and non-overlapping regions,  $R_1, R_2, \dots, R_J$ .
2. For every observation that falls into the region  $R_j$ , we make the same prediction, which is simply the mean of the response values for the training observations in  $R_j$ .

How do we construct the regions  $R_1, \dots, R_J$ ? In theory, the regions could have any shape. However, we choose to divide the predictor space into high-dimensional rectangles, or boxes, for simplicity and for ease of interpretation of the resulting predictive model. The goal is to find boxes  $R_1, \dots, R_J$  that minimize the MSE, given by

$$\sum_{j=1}^J \sum_{i \in R_j} (y_i - \bar{y}_{R_j})^2,$$

where  $\bar{y}_{R_j}$  is the mean response for the training observations within box  $j$ .

### A top-down approach, recursive binary splitting

Unfortunately, it is computationally infeasible to consider every possible partition of the feature space into  $J$  boxes. The common strategy is to take a top-down approach

The approach is top-down because it begins at the top of the tree (all observations belong to a single region) and then successively splits the predictor space; each split is indicated via two new branches further down on the tree. It is greedy because at each step of the tree-building process, the best split is made at that particular step, rather than looking ahead and picking a split that will lead to a better tree in some future step.

### Making a tree

In order to implement the recursive binary splitting we start by selecting the predictor  $x_j$  and a cutpoint  $s$  that splits the predictor space into two regions  $R_1$  and  $R_2$

$$\{X | x_j < s\},$$

and

$$\{X|x_j \geq s\},$$

so that we obtain the lowest MSE, that is

$$\sum_{i:x_i \in R_j} (y_i - \bar{y}_{R_1})^2 + \sum_{i:x_i \in R_2} (y_i - \bar{y}_{R_2})^2,$$

which we want to minimize by considering all predictors  $x_1, x_2, \dots, x_p$ . We consider also all possible values of  $s$  for each predictor. These values could be determined by randomly assigned numbers or by starting at the midpoint and then proceed till we find an optimal value.

For any  $j$  and  $s$ , we define the pair of half-planes where  $\bar{y}_{R_1}$  is the mean response for the training observations in  $R_1(j, s)$ , and  $\bar{y}_{R_2}$  is the mean response for the training observations in  $R_2(j, s)$ .

Finding the values of  $j$  and  $s$  that minimize the above equation can be done quite quickly, especially when the number of features  $p$  is not too large.

Next, we repeat the process, looking for the best predictor and best cutpoint in order to split the data further so as to minimize the MSE within each of the resulting regions. However, this time, instead of splitting the entire predictor space, we split one of the two previously identified regions. We now have three regions. Again, we look to split one of these three regions further, so as to minimize the MSE. The process continues until a stopping criterion is reached; for instance, we may continue until no region contains more than five observations.

## Pruning the tree

The above procedure is rather straightforward, but leads often to overfitting and unnecessarily large and complicated trees. The basic idea is to grow a large tree  $T_0$  and then prune it back in order to obtain a subtree. A smaller tree with fewer splits (fewer regions) can lead to smaller variance and better interpretation at the cost of a little more bias.

The so-called Cost complexity pruning algorithm gives us a way to do just this. Rather than considering every possible subtree, we consider a sequence of trees indexed by a nonnegative tuning parameter  $\alpha$ .

## Cost complexity pruning

For each value of  $\alpha$  there corresponds a subtree  $T \in T_0$  such that

$$\sum_{m=1}^{\bar{T}} \sum_{i:x_i \in R_m} (y_i - \bar{y}_{R_m})^2 + \alpha \bar{T},$$

is as small as possible. Here  $\bar{T}$  is the number of terminal nodes of the tree  $T$ ,  $R_m$  is the rectangle (i.e. the subset of predictor space) corresponding to the  $m$ -th terminal node.

The tuning parameter  $\alpha$  controls a trade-off between the subtree's complexity and its fit to the training data. When  $\alpha = 0$ , then the subtree  $T$  will simply equal  $T_0$ , because then the above equation just measures the training error. However, as  $\alpha$  increases, there is a price to pay for having a tree with many terminal nodes. The above equation will tend to be minimized for a smaller subtree.

It turns out that as we increase  $\alpha$  from zero branches get pruned from the tree in a nested and predictable fashion, so obtaining the whole sequence of subtrees as a function of  $\alpha$  is easy. We can select a value of  $\alpha$  using a validation set or using cross-validation. We then return to the full data set and obtain the subtree corresponding to  $\alpha$ .

## Schematic Regression Procedure

### Building a Regression Tree.

1. Use recursive binary splitting to grow a large tree on the training data, stopping only when each terminal node has fewer than some minimum number of observations.
2. Apply cost complexity pruning to the large tree in order to obtain a sequence of best subtrees, as a function of  $\alpha$ .
3. Use for example  $K$ -fold cross-validation to choose  $\alpha$ . Divide the training observations into  $K$  folds. For each  $k = 1, 2, \dots, K$  we:
  - repeat steps 1 and 2 on all but the  $k$ -th fold of the training data.
  - Then we evaluate the mean squared prediction error on the data in the left-out  $k$ -th fold, as a function of  $\alpha$ .
  - Finally we average the results for each value of  $\alpha$ , and pick  $\alpha$  to minimize the average error.
4. Return the subtree from Step 2 that corresponds to the chosen value of  $\alpha$ .

## A Classification Tree

A classification tree is very similar to a regression tree, except that it is used to predict a qualitative response rather than a quantitative one. Recall that for a regression tree, the predicted response for an observation is given by the mean response of the training observations that belong to the same terminal node. In contrast, for a classification tree, we predict that each observation belongs to the most commonly occurring class of training observations in the region to which it belongs. In interpreting the results of a classification tree, we are often interested not only in the class prediction corresponding to a particular terminal node region, but also in the class proportions among the training observations that fall into that region.

## Growing a classification tree

The task of growing a classification tree is quite similar to the task of growing a regression tree. Just as in the regression setting, we use recursive binary splitting to grow a classification tree. However, in the classification setting, the MSE cannot be used as a criterion for making the binary splits. A natural alternative to MSE is the **classification error rate**. Since we plan to assign an observation in a given region to the most commonly occurring error rate class of training observations in that region, the classification error rate is simply the fraction of the training observations in that region that do not belong to the most common class.

When building a classification tree, either the Gini index or the entropy are typically used to evaluate the quality of a particular split, since these two approaches are more sensitive to node purity than is the classification error rate.

### Classification tree, how to split nodes

If our targets are the outcome of a classification process that takes for example  $k = 1, 2, \dots, K$  values, the only thing we need to think of is to set up the splitting criteria for each node.

We define a PDF  $p_{mk}$  that represents the number of observations of a class  $k$  in a region  $R_m$  with  $N_m$  observations. We represent this likelihood function in terms of the proportion  $I(y_i = k)$  of observations of this class in the region  $R_m$  as

$$p_{mk} = \frac{1}{N_m} \sum_{x_i \in R_m} I(y_i = k).$$

We let  $p_{mk}$  represent the majority class of observations in region  $m$ . The three most common ways of splitting a node are given by

- Misclassification error

$$p_{mk} = \frac{1}{N_m} \sum_{x_i \in R_m} I(y_i \neq k) = 1 - p_{mk}.$$

- Gini index  $g$

$$g = \sum_{k=1}^K p_{mk}(1 - p_{mk}).$$

- Information entropy or just entropy  $s$

$$s = - \sum_{k=1}^K p_{mk} \log p_{mk}.$$

## Visualizing the Tree, Classification

## Visualizing the Tree, The Moons

## Algorithms for Setting up Decision Trees

Two algorithms stand out in the set up of decision trees:

1. The CART (Classification And Regression Tree) algorithm for both classification and regression
2. The ID3 algorithm based on the computation of the information gain for classification

We discuss both algorithms with applications here. The popular library **Scikit-Learn** uses the CART algorithm. For classification problems you can use either the **gini** index or the **entropy** to split a tree in two branches.

## The CART algorithm for Classification

For classification, the CART algorithm splits the data set in two subsets using a single feature  $k$  and a threshold  $t_k$ . This could be for example a threshold set by a number below a certain circumference of a malign tumor.

How do we find these two quantities? We search for the pair  $(k, t_k)$  that produces the purest subset using for example the **gini** factor  $G$ . The cost function it tries to minimize is then

$$C(k, t_k) = \frac{m_{\text{left}}}{m} G_{\text{left}} + \frac{m_{\text{right}}}{m} G_{\text{right}},$$

where  $G_{\text{left/right}}$  measures the impurity of the left/right subset and  $m_{\text{left/right}}$  is the number of instances in the left/right subset

Once it has successfully split the training set in two, it splits the subsets using the same logic, then the subsubsets and so on, recursively. It stops recursing once it reaches the maximum depth (defined by the *max\_depth* hyperparameter), or if it cannot find a split that will reduce impurity. A few other hyperparameters control additional stopping conditions such as the *min\_samples\_split*, *min\_samples\_leaf*, *min\_weight\_fraction\_leaf*, and *max\_leaf\_nodes*.

## The CART algorithm for Regression

The CART algorithm for regression works is similar to the one for classification except that instead of trying to split the training set in a way that minimizes say the **gini** or **entropy** impurity, it now tries to split the training set in a way that minimizes our well-known mean-squared error (MSE). The cost function is now

$$C(k, t_k) = \frac{m_{\text{left}}}{m} \text{MSE}_{\text{left}} + \frac{m_{\text{right}}}{m} \text{MSE}_{\text{right}}.$$

Here the MSE for a specific node is defined as

$$\text{MSE}_{\text{node}} = \frac{1}{m_{\text{node}}} \sum_{i \in \text{node}} (\bar{y}_{\text{node}} - y_i)^2,$$

with

$$\bar{y}_{\text{node}} = \frac{1}{m_{\text{node}}} \sum_{i \in \text{node}} y_i,$$

the mean value of all observations in a specific node.

Without any regularization, the regression task for decision trees, just like for classification tasks, is prone to overfitting.

## Computing the Gini index

The example we will look at is a classical one in many Machine Learning applications. Based on various meteorological features, we have several so-called attributes which decide whether we at the end will do some outdoor activity like skiing, going for a bike ride etc etc. The table here contains the features **outlook**, **temperature**, **humidity** and **wind**. The target or output is whether we ride (True=1) or whether we do something else that day (False=0). The attributes for each feature are then sunny, overcast and rain for the outlook, hot, cold and mild for temperature, high and normal for humidity and weak and strong for wind.

The table here summarizes the various attributes and

Day	Outlook	Temperature	Humidity	Wind	Ride
1	Sunny	Hot	High	Weak	0
2	Sunny	Hot	High	Strong	1
3	Overcast	Hot	High	Weak	1
4	Rain	Mild	High	Weak	1
5	Rain	Cool	Normal	Weak	1
6	Rain	Cool	Normal	Strong	0
7	Overcast	Cool	Normal	Strong	1
8	Sunny	Mild	High	Weak	0
9	Sunny	Cool	Normal	Weak	1
10	Rain	Mild	Normal	Weak	1
11	Sunny	Mild	Normal	Strong	1
12	Overcast	Mild	High	Strong	1
13	Overcast	Hot	Normal	Weak	1
14	Rain	Mild	High	Strong	0



## Simple Python Code to read in Data and perform Classification

### Computing the Gini Factor

The above functions (gini, entropy and misclassification error) are important components of the so-called CART algorithm. We will discuss this algorithm below after we have discussed the information gain algorithm ID3.

In the example here we have converted all our attributes into numerical values 0, 1, 2 etc.

### Entropy and the ID3 algorithm

ID3, learns decision trees by constructing them topdown, beginning with the question **which attribute should be tested at the root of the tree?**

1. Each instance attribute is evaluated using a statistical test to determine how well it alone classifies the training examples.
2. The best attribute is selected and used as the test at the root node of the tree.
3. A descendant of the root node is then created for each possible value of this attribute.
4. Training examples are sorted to the appropriate descendant node.
5. The entire process is then repeated using the training examples associated with each descendant node to select the best attribute to test at that point in the tree.
6. This forms a greedy search for an acceptable decision tree, in which the algorithm never backtracks to reconsider earlier choices.

The ID3 algorithm selects, which attribute to test at each node in the tree.

We would like to select the attribute that is most useful for classifying examples.

What is a good quantitative measure of the worth of an attribute?

Information gain measures how well a given attribute separates the training examples according to their target classification.

The ID3 algorithm uses this information gain measure to select among the candidate attributes at each step while growing the tree.

## Cancer Data again now with Decision Trees and other Methods

### Pros and cons of trees, pros

- White box, easy to interpret model. Some people believe that decision trees more closely mirror human decision-making than do the regression

and classification approaches discussed earlier (think of support vector machines)

- Trees are very easy to explain to people. In fact, they are even easier to explain than linear regression!
- No feature normalization needed
- Tree models can handle both continuous and categorical data (Classification and Regression Trees)
- Can model nonlinear relationships
- Can model interactions between the different descriptive features
- Trees can be displayed graphically, and are easily interpreted even by a non-expert (especially if they are small)

## Disadvantages

- Unfortunately, trees generally do not have the same level of predictive accuracy as some of the other regression and classification approaches
- If continuous features are used the tree may become quite large and hence less interpretable
- Decision trees are prone to overfit the training data and hence do not well generalize the data if no stopping criteria or improvements like pruning, boosting or bagging are implemented
- Small changes in the data may lead to a completely different tree. This issue can be addressed by using ensemble methods like bagging, boosting or random forests
- Unbalanced datasets where some target feature values occur much more frequently than others may lead to biased trees since the frequently occurring feature values are preferred over the less frequently occurring ones.
- If the number of features is relatively large (high dimensional) and the number of instances is relatively low, the tree might overfit the data
- Features with many levels may be preferred over features with less levels since for them it is *more easy* to split the dataset such that the sub datasets only contain pure target feature values. This issue can be addressed by preferring for instance the information gain ratio as splitting criteria over information gain

However, by aggregating many decision trees, using methods like bagging, random forests, and boosting, the predictive performance of trees can be substantially improved.

## Ensemble Methods: From a Single Tree to Many Trees and Extreme Boosting, Meet the Jungle of Methods

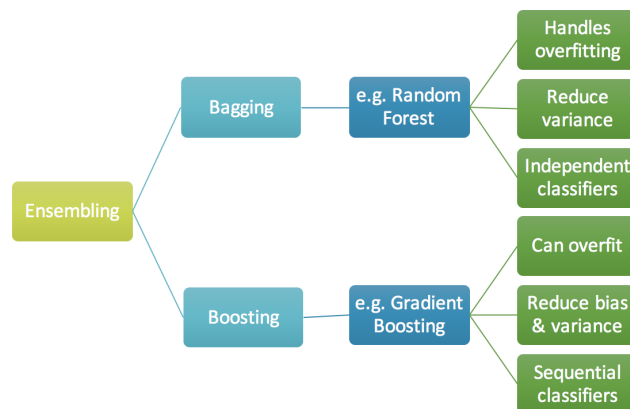
As stated above and seen in many of the examples discussed here about a single decision tree, we often end up overfitting our training data. This normally means that we have a high variance. Can we reduce the variance of a statistical learning method?

This leads us to a set of different methods that can combine different machine learning algorithms or just use one of them to construct forests and jungles of trees, homogeneous ones or heterogenous ones. These methods are recognized by different names which we will try to explain here. These are

1. Voting classifiers
2. Bagging and Pasting
3. Random forests
4. Boosting methods, from adaptive to Extreme Gradient Boosting (XGBoost)

We discuss these methods here.

### An Overview of Ensemble Methods



### Bagging

The **plain** decision trees suffer from high variance. This means that if we split the training data into two parts at random, and fit a decision tree to both halves, the results that we get could be quite different. In contrast, a procedure with low variance will yield similar results if applied repeatedly to distinct data sets;

linear regression tends to have low variance, if the ratio of  $n$  to  $p$  is moderately large.

**Bootstrap aggregation**, or just **bagging**, is a general-purpose procedure for reducing the variance of a statistical learning method.

## More bagging

Bagging typically results in improved accuracy over prediction using a single tree. Unfortunately, however, it can be difficult to interpret the resulting model. Recall that one of the advantages of decision trees is the attractive and easily interpreted diagram that results.

However, when we bag a large number of trees, it is no longer possible to represent the resulting statistical learning procedure using a single tree, and it is no longer clear which variables are most important to the procedure. Thus, bagging improves prediction accuracy at the expense of interpretability. Although the collection of bagged trees is much more difficult to interpret than a single tree, one can obtain an overall summary of the importance of each predictor using the MSE (for bagging regression trees) or the Gini index (for bagging classification trees). In the case of bagging regression trees, we can record the total amount that the MSE is decreased due to splits over a given predictor, averaged over all  $B$  possible trees. A large value indicates an important predictor. Similarly, in the context of bagging classification trees, we can add up the total amount that the Gini index is decreased by splits over a given predictor, averaged over all  $B$  trees.

## Making your own Bootstrap: Changing the Level of the Decision Tree

Let us bring up our good old bootstrap example from the linear regression lectures. We change the linear regression algorithm with a decision tree with different depths and perform a bootstrap aggregate (in this case we perform as many bootstraps as data points  $n$ ).

## Random forests

Random forests provide an improvement over bagged trees by way of a small tweak that decorrelates the trees.

As in bagging, we build a number of decision trees on bootstrapped training samples. But when building these decision trees, each time a split in a tree is considered, a random sample of  $m$  predictors is chosen as split candidates from the full set of  $p$  predictors. The split is allowed to use only one of those  $m$  predictors.

A fresh sample of  $m$  predictors is taken at each split, and typically we choose

$$m \approx \sqrt{p}.$$

In building a random forest, at each split in the tree, the algorithm is not even allowed to consider a majority of the available predictors.

The reason for this is rather clever. Suppose that there is one very strong predictor in the data set, along with a number of other moderately strong predictors. Then in the collection of bagged variable importance random forest trees, most or all of the trees will use this strong predictor in the top split. Consequently, all of the bagged trees will look quite similar to each other. Hence the predictions from the bagged trees will be highly correlated. Unfortunately, averaging many highly correlated quantities does not lead to as large of a reduction in variance as averaging many uncorrelated quantities. In particular, this means that bagging will not lead to a substantial reduction in variance over a single tree in this setting.

## Random Forest Algorithm

The algorithm described here can be applied to both classification and regression problems.

We will grow a forest of say  $B$  trees.

1. For  $b = 1 : B$ 
  - Draw a bootstrap sample of from the training data organized in our  $\mathbf{X}$  matrix.
  - We grow then a random forest tree  $T_b$  based on the bootstrapped data by repeating the steps outlined till we reach the maximum node size is reached
    - (a) we select  $m \leq p$  variables at random from the  $p$  predictors/features
    - (b) pick the best split point among the  $m$  features using either the CART algorithm or the ID3 for classification and create a new node
    - (c) split the node into daughter nodes
2. Output then the ensemble of trees  $\{T_b\}_1^B$  and make predictions for either a regression type of problem or a classification type of problem.

## Random Forests Compared with other Methods on the Cancer Data

### Compare Bagging on Trees with Random Forests

### Boosting, a Bird's Eye View

The basic idea is to combine weak classifiers in order to create a good classifier. With a weak classifier we often intend a classifier which produces results which are only slightly better than we would get by random guesses.

This is done by applying in an iterative way a weak (or a standard classifier like decision trees) to modify the data. In each iteration we emphasize those observations which are misclassified by weighting them with a factor.

## What is boosting? Additive Modelling/Iterative Fitting

Boosting is a way of fitting an additive expansion in a set of elementary basis functions like for example some simple polynomials. Assume for example that we have a function

$$f_M(x) = \sum_{i=1}^M \beta_m b(x; \gamma_m),$$

where  $\beta_m$  are the expansion parameters to be determined in a minimization process and  $b(x; \gamma_m)$  are some simple functions of the multivariable parameter  $x$  which is characterized by the parameters  $\gamma_m$ .

As an example, consider the Sigmoid function we used in logistic regression. In that case, we can translate the function  $b(x; \gamma_m)$  into the Sigmoid function

$$\sigma(t) = \frac{1}{1 + \exp(-t)},$$

where  $t = \gamma_0 + \gamma_1 x$  and the parameters  $\gamma_0$  and  $\gamma_1$  were determined by the Logistic Regression fitting algorithm.

As another example, consider the cost function we defined for linear regression

$$C(\mathbf{y}, \mathbf{f}) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - f(x_i))^2.$$

In this case the function  $f(x)$  was replaced by the design matrix  $\mathbf{X}$  and the unknown linear regression parameters  $\beta$ , that is  $\mathbf{f} = \mathbf{X}\beta$ . In linear regression we can simply invert a matrix and obtain the parameters  $\beta$  by

$$\beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

In iterative fitting or additive modeling, we minimize the cost function with respect to the parameters  $\beta_m$  and  $\gamma_m$ .

## Iterative Fitting, Regression and Squared-error Cost Function

The way we proceed is as follows (here we specialize to the squared-error cost function)

1. Establish a cost function, here  $\mathcal{C}(\mathbf{y}, \mathbf{f}) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - f_M(x_i))^2$  with  $f_M(x) = \sum_{i=1}^M \beta_m b(x; \gamma_m)$ .
2. Initialize with a guess  $f_0(x)$ . It could be one or even zero or some random numbers.

3. For  $m = 1 : M$

- (a) minimize  $\sum_{i=0}^{n-1} (y_i - f_{m-1}(x_i) - \beta b(x; \gamma))^2$  wrt  $\gamma$  and  $\beta$
- (b) This gives the optimal values  $\beta_m$  and  $\gamma_m$
- (c) Determine then the new values  $f_m(x) = f_{m-1}(x) + \beta_m b(x; \gamma_m)$

We could use any of the algorithms we have discussed till now. If we use trees,  $\gamma$  parameterizes the split variables and split points at the internal nodes, and the predictions at the terminal nodes.

### Squared-Error Example and Iterative Fitting

To better understand what happens, let us develop the steps for the iterative fitting using the above squared error function.

For simplicity we assume also that our functions  $b(x; \gamma) = 1 + \gamma x$ .

This means that for every iteration  $m$ , we need to optimize

$$(\beta_m, \gamma_m) = \operatorname{argmin}_{\beta, \gamma} \sum_{i=0}^{n-1} (y_i - f_{m-1}(x_i) - \beta b(x; \gamma))^2 = \sum_{i=0}^{n-1} (y_i - f_{m-1}(x_i) - \beta(1 + \gamma x_i))^2.$$

We start our iteration by simply setting  $f_0(x) = 0$ . Taking the derivatives with respect to  $\beta$  and  $\gamma$  we obtain

$$\frac{\partial \mathcal{C}}{\partial \beta} = -2 \sum_i (1 + \gamma x_i)(y_i - \beta(1 + \gamma x_i)) = 0,$$

and

$$\frac{\partial \mathcal{C}}{\partial \gamma} = -2 \sum_i \beta x_i (y_i - \beta(1 + \gamma x_i)) = 0.$$

We can then rewrite these equations as (defining  $\mathbf{w} = \mathbf{e} + \gamma \mathbf{x}$  with  $\mathbf{e}$  being the unit vector)

$$\gamma \mathbf{w}^T (\mathbf{y} - \beta \gamma \mathbf{w}) = 0,$$

which gives us  $\beta = \mathbf{w}^T \mathbf{y} / (\mathbf{w}^T \mathbf{w})$ . Similarly we have

$$\beta \gamma \mathbf{x}^T (\mathbf{y} - \beta(1 + \gamma \mathbf{x})) = 0,$$

which leads to  $\gamma = (\mathbf{x}^T \mathbf{y} - \beta \mathbf{x}^T \mathbf{e}) / (\beta \mathbf{x}^T \mathbf{x})$ . Inserting for  $\beta$  gives us an equation for  $\gamma$ . This is a non-linear equation in the unknown  $\gamma$  and has to be solved numerically.

The solution to these two equations gives us in turn  $\beta_1$  and  $\gamma_1$  leading to the new expression for  $f_1(x)$  as  $f_1(x) = \beta_1(1 + \gamma_1 x)$ . Doing this  $M$  times results in our final estimate for the function  $f$ .

## Iterative Fitting, Classification and AdaBoost

Let us consider a binary classification problem with two outcomes  $y_i \in \{-1, 1\}$  and  $i = 0, 1, 2, \dots, n-1$  as our set of observations. We define a classification function  $G(x)$  which produces a prediction taking one or the other of the two values  $\{-1, 1\}$ .

The error rate of the training sample is then

$$\overline{\text{err}} = \frac{1}{n} \sum_{i=0}^{n-1} I(y_i \neq G(x_i)).$$

The iterative procedure starts with defining a weak classifier whose error rate is barely better than random guessing. The iterative procedure in boosting is to sequentially apply a weak classification algorithm to repeatedly modified versions of the data producing a sequence of weak classifiers  $G_m(x)$ .

Here we will express our function  $f(x)$  in terms of  $G(x)$ . That is

$$f_M(x) = \sum_{i=1}^M \beta_m b(x; \gamma_m),$$

will be a function of

$$G_M(x) = \text{sign} \sum_{i=1}^M \alpha_m G_m(x).$$

## Adaptive Boosting, AdaBoost

In our iterative procedure we define thus

$$f_m(x) = f_{m-1}(x) + \beta_m G_m(x).$$

The simplest possible cost function which leads (also simple from a computational point of view) to the AdaBoost algorithm is the exponential cost/loss function defined as

$$C(\mathbf{y}, \mathbf{f}) = \sum_{i=0}^{n-1} \exp(-y_i(f_{m-1}(x_i) + \beta G(x_i))).$$

We optimize  $\beta$  and  $G$  for each value of  $m = 1 : M$  as we did in the regression case. This is normally done in two steps. Let us however first rewrite the cost function as

$$C(\mathbf{y}, \mathbf{f}) = \sum_{i=0}^{n-1} w_i^m \exp(-y_i \beta G(x_i)),$$

where we have defined  $w_i^m = \exp(-y_i f_{m-1}(x_i))$ .



## Building up AdaBoost

First, for any  $\beta > 0$ , we optimize  $G$  by setting

$$G_m(x) = \text{sign} \sum_{i=0}^{n-1} w_i^m I(y_i \neq G(x_i)),$$

which is the classifier that minimizes the weighted error rate in predicting  $y$ .

We can do this by rewriting

$$\exp(-\beta) \sum_{y_i=G(x_i)} w_i^m + \exp(\beta) \sum_{y_i \neq G(x_i)} w_i^m,$$

which can be rewritten as

$$(\exp(\beta) - \exp(-\beta)) \sum_{i=0}^{n-1} w_i^m I(y_i \neq G(x_i)) + \exp(-\beta) \sum_{i=0}^{n-1} w_i^m = 0,$$

which leads to

$$\beta_m = \frac{1}{2} \log \frac{1 - \overline{\text{err}}}{\overline{\text{err}}},$$

where we have redefined the error as

$$\overline{\text{err}}_m = \frac{1}{n} \frac{\sum_{i=0}^{n-1} w_i^m I(y_i \neq G(x_i))}{\sum_{i=0}^{n-1} w_i^m},$$

which leads to an update of

$$f_m(x) = f_{m-1}(x) + \beta_m G_m(x).$$

This leads to the new weights

$$w_i^{m+1} = w_i^m \exp(-y_i \beta_m G_m(x_i))$$

## Adaptive boosting: AdaBoost, Basic Algorithm

The algorithm here is rather straightforward. Assume that our weak classifier is a decision tree and we consider a binary set of outputs with  $y_i \in \{-1, 1\}$  and  $i = 0, 1, 2, \dots, n-1$  as our set of observations. Our design matrix is given in terms of the feature/predictor vectors  $\mathbf{X} = [\mathbf{x}_0 \mathbf{x}_1 \dots \mathbf{x}_{p-1}]$ . Finally, we define also a classifier determined by our data via a function  $G(x)$ . This function tells us how well we are able to classify our outputs/targets  $\mathbf{y}$ .

We have already defined the misclassification error  $\text{err}$  as

$$\text{err} = \frac{1}{n} \sum_{i=0}^{n-1} I(y_i \neq G(x_i)),$$

where the function  $I()$  is one if we misclassify and zero if we classify correctly.

## Basic Steps of AdaBoost

With the above definitions we are now ready to set up the algorithm for AdaBoost. The basic idea is to set up weights which will be used to scale the correctly classified and the misclassified cases.

1. We start by initializing all weights to  $w_i = 1/n$ , with  $i = 0, 1, 2, \dots, n-1$ . It is easy to see that we must have  $\sum_{i=0}^{n-1} w_i = 1$ .

2. We rewrite the misclassification error as

$$\overline{\text{err}}_m = \frac{\sum_{i=0}^{n-1} w_i^m I(y_i \neq G(x_i))}{\sum_{i=0}^{n-1} w_i},$$

1. Then we start looping over all attempts at classifying, namely we start an iterative process for  $m = 1 : M$ , where  $M$  is the final number of classifications. Our given classifier could for example be a plain decision tree.

- (a) Fit then a given classifier to the training set using the weights  $w_i$ .
- (b) Compute then err and figure out which events are classified properly and which are classified wrongly.
- (c) Define a quantity  $\alpha_m = \log(1 - \overline{\text{err}}_m) / \overline{\text{err}}_m$
- (d) Set the new weights to  $w_i = w_i \times \exp(\alpha_m I(y_i \neq G(x_i)))$ .

2. Compute the new classifier  $G(x) = \sum_{i=0}^{n-1} \alpha_m I(y_i \neq G(x_i))$ .

For the iterations with  $m \leq 2$  the weights are modified individually at each steps. The observations which were misclassified at iteration  $m-1$  have a weight which is larger than those which were classified properly. As this proceeds, the observations which were difficult to classify correctly are given a larger influence. Each new classification step  $m$  is then forced to concentrate on those observations that are missed in the previous iterations.

## AdaBoost Examples

Using **Scikit-Learn** it is easy to apply the adaptive boosting algorithm, as done here.

## AdaBoost for Regression

Here we present [Drucker's AdaBoost](#) tailored for regression.

In bagging, each training example is equally likely to be picked. In boosting, the probability of a particular example being in the training set of a particular machine depends on the performance of the prior machines on that example. The following is a modification of Adaboost by Drucker.

Start by selecting a set of training data  $n$  and assign to each entry a weight  $w_i = 1$  for  $i = 1, 2, \dots, n$ . As we have done earlier, we could pick say 80% of the data set for training. The algorithm runs as follows:

1. We define the probability that the training sample  $i$  is in the set by  $p_i = w_i / \sum_i w_i$ . We pick  $n$  samples (with replacement) to form our training set. We pick a number uniformly in the range  $[0, \sum_i w_i]$ .
2. We choose then a regression machine (for example plain linear regression or a simple decision tree). A given regression machine makes then a hypothesis.
3. Using every member of the training set with the chosen regression machine we obtain then a prediction  $\tilde{y}_i$ .
4. We calculate then the loss function  $L_i$  for each training sample. We can use various types of loss function as long as we have a value

$$L_i \in [0, 1].$$

## Gradient boosting: Basics with Steepest Descent

Gradient boosting is again a similar technique to Adaptive boosting, it combines so-called weak classifiers or regressors into a strong method via a series of iterations.

In order to understand the method, let us illustrate its basics by bringing back the essential steps in linear regression, where our cost function was the least squares function.

## The Squared-Error again! Steepest Descent

We start again with our cost function  $\mathcal{C}(\mathbf{y} \mathbf{m} \mathbf{f}) = \sum_{i=0}^{n-1} \mathcal{L}(y_i, f(x_i))$  where we want to minimize. This means that for every iteration, we need to optimize

$$(\hat{\mathbf{f}}) = \operatorname{argmin}_{\mathbf{f}} \sum_{i=0}^{n-1} (y_i - f(x_i))^2.$$

We define a real function  $h_m(x)$  that defines our final function  $f_M(x)$  as

$$f_M(x) = \sum_{m=0}^M h_m(x).$$

In the steepest decent approach we approximate  $h_m(x) = -\rho_m g_m(x)$ , where  $\rho_m$  is a scalar and  $g_m(x)$  the gradient defined as

$$g_m(x_i) = \left[ \frac{\partial \mathcal{L}(y_i, f(x_i))}{\partial f(x_i)} \right]_{f(x_i)=f_{m-1}(x_i)}.$$

With the new gradient we can update  $f_m(x) = f_{m-1}(x) - \rho_m g_m(x)$ . Using the above squared-error function we see that the gradient is  $g_m(x_i) = -2(y_i - f(x_i))$ .

Choosing  $f_0(x) = 0$  we obtain  $g_m(x) = -2y_i$  and inserting this into the minimization problem for the cost function we have

$$(\rho_1) = \operatorname{argmin}_{\rho} \sum_{i=0}^{n-1} (y_i + 2\rho y_i)^2.$$

### Steepest Descent Example

Optimizing with respect to  $\rho$  we obtain (taking the derivative) that  $\rho_1 = -1/2$ . We have then that

$$f_1(x) = f_0(x) - \rho_1 g_1(x) = -y_i.$$

We can then proceed and compute

$$g_2(x_i) = \left[ \frac{\partial \mathcal{L}(y_i, f(x_i))}{\partial f(x_i)} \right]_{f(x_i)=f_1(x_i)=y_i} = -4y_i,$$

and find a new value for  $\rho_2 = -1/2$  and continue till we have reached  $m = M$ . We can modify the steepest descent method, or steepest boosting, by introducing what is called **gradient boosting**.

### Gradient Boosting, algorithm

Suppose we have a cost function  $C(f) = \sum_{i=0}^{n-1} L(y_i, f(x_i))$  where  $y_i$  is our target and  $f(x_i)$  the function which is meant to model  $y_i$ . The above cost function could be our standard squared-error function

$$C(\mathbf{y}, \mathbf{f}) = \sum_{i=0}^{n-1} (y_i - f(x_i))^2.$$

The way we proceed in an iterative fashion is to

1. Initialize our estimate  $f_0(x)$ .
2. For  $m = 1 : M$ , we
  - (a) compute the negative gradient vector  $\mathbf{u}_m = -\partial C(\mathbf{y}, \mathbf{f}) / \partial \mathbf{f}(x)$  at  $f(x) = f_{m-1}(x)$ ;
  - (b) fit the so-called base-learner to the negative gradient  $h_m(u_m, x)$ ;
  - (c) update the estimate  $f_m(x) = f_{m-1}(x) + \nu h_m(u_m, x)$ ;
3. The final estimate is then  $f_M(x) = \sum_{m=1}^M \nu h_m(u_m, x)$ .

## Gradient Boosting, Examples of Regression

### Gradient Boosting, Classification Example

#### XGBoost: Extreme Gradient Boosting

[XGBoost](#) or Extreme Gradient Boosting, is an optimized distributed gradient boosting library designed to be highly efficient, flexible and portable. It implements machine learning algorithms under the Gradient Boosting framework. XGBoost provides a parallel tree boosting that solve many data science problems in a fast and accurate way. See the [article by Chen and Guestrin](#).

The authors design and build a highly scalable end-to-end tree boosting system. It has a theoretically justified weighted quantile sketch for efficient proposal calculation. It introduces a novel sparsity-aware algorithm for parallel tree learning and an effective cache-aware block structure for out-of-core tree learning.

It is now the algorithm which wins essentially all ML competitions!!!

### Regression Case

#### Xgboost on the Cancer Data

As you will see from the confusion matrix below, XGBoots does an excellent job on the Wisconsin cancer data and outperforms essentially all agorithms we have discussed till now.