

Statistical Methods for Data Science

Semester 1 2024

DATA7202

Tree Models

Slava Vaisman

The University of Queensland

r.vaisman@uq.edu.au

- 1 Decision trees
 - Regional Prediction Functions
 - Termination Criterion
 - Additional Considerations
 - Controlling the Tree Shape
 - Advantages Decision Trees
- 2 Bootstrap Aggregation and Random Forests
- 3 Boosting

Introduction

- Statistical learning methods based on decision trees have gained tremendous popularity due to their simplicity, intuitive representation, and predictive accuracy.
- In this unit, we give an introduction to the construction and use of such trees.
- We also discuss two key ensemble methods, namely bootstrap aggregation and boosting, which can further improve the efficiency of decision trees and other learning methods.

Decision trees

- Tree-based methods provide a simple, intuitive, and powerful mechanism for both regression and classification.
- The main idea is to divide a (potentially complicated) feature space \mathcal{X} into smaller regions and fit a simple prediction function to each region.
- For example, in a regression setting, one could take the mean of the training responses associated with the training features that fall in that specific region.
- In the classification setting, a commonly used prediction function takes the majority vote among the corresponding response variables.

We start with a simple classification example.

Decision trees example (1)

The left panel of Figure 1 shows a training set of 15 two-dimensional points (features) falling into two classes (red and blue). How should the new feature vector (black point) be classified?

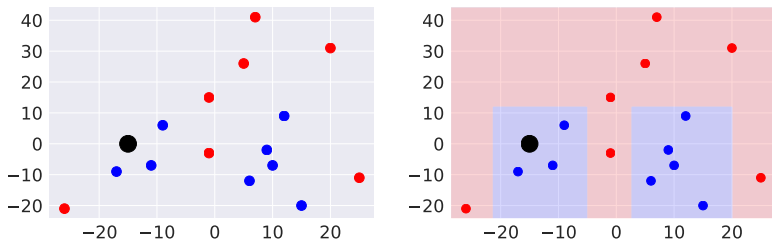
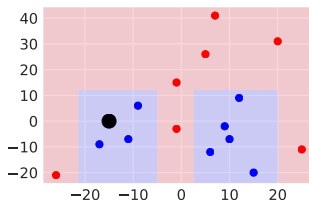
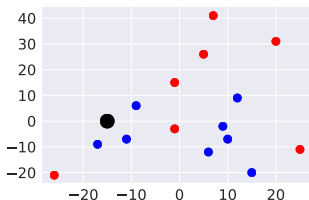


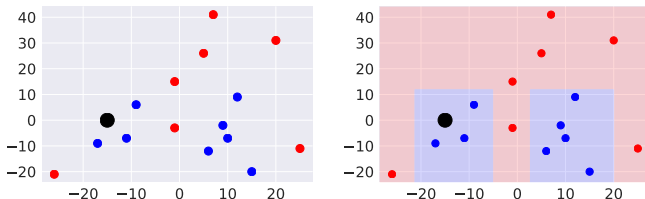
Figure 1: Left: training data and a new feature. Right: a partition of the feature space.

Decision trees example (2)

- It is not possible to linearly separate the training set, but we can partition the feature space $X = \mathbb{R}^2$ into rectangular regions and assign a class (color) to each region, as shown in the right panel of Figure 1.
- Points in these regions are classified accordingly as blue or red.
- The partition thus defines a classifier (prediction function) g (in this unit we will use g instead of h), that assigns to each feature vector \mathbf{x} a class “red” or “blue”.



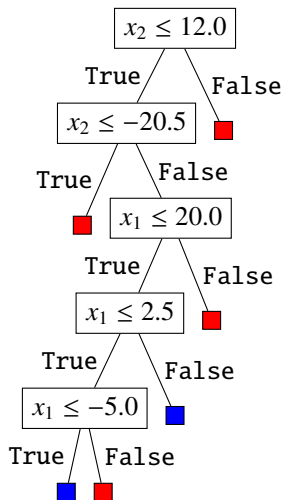
Decision trees example (3)



- For example, for $\mathbf{x} = [-15, 0]^T$ (solid black point), $g(\mathbf{x}) = \text{"blue"}$, since it belongs to a blue region of the feature space.
- Both the classification procedure and the partitioning of the feature space can be conveniently represented by a binary decision tree. This is a tree where each node v corresponds to a region (subset) R_v of the feature space X — the root node corresponding to the feature space itself.

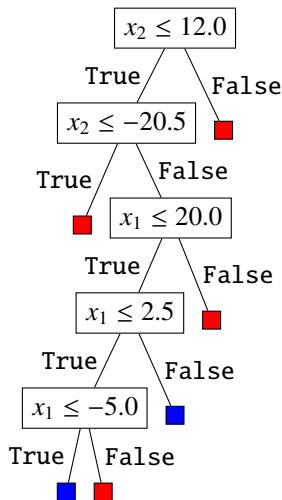
Decision trees example (4)

- Each internal node v contains a logical condition that divides R_v into two disjoint subregions.
- The leaf nodes (the terminal nodes of the tree) are not subdivided, and their corresponding regions form a partition of X , as they are disjoint and their union is X .
- Associated with each leaf node w is also a regional prediction function g^w on R_w .
- The partitioning of the 2D space, was obtained from the decision tree shown on the right.



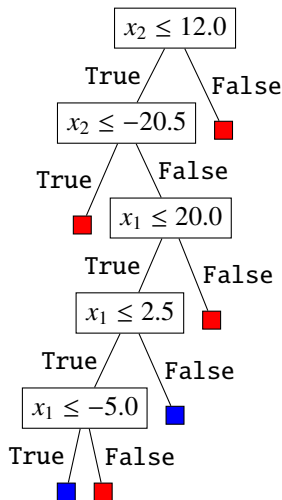
Decision trees example (5)

- As an illustration of the decision procedure, consider again the input $\mathbf{x} = [x_1, x_2]^T = [-15, 0]^T$.
- The classification process starts from the tree root, which contains the condition $x_2 \leq 12.0$.
- As the second component of \mathbf{x} is 0, the root condition is satisfied and we proceed to the left child, which contains the condition $x_2 \leq -20.5$.



Decision trees example (6)

- The next step is similar. As $0 > -20.5$, the condition is not satisfied and we proceed to the right child.
- Such an evaluation of logical conditions along the tree path will eventually bring us to a leaf node and its associated region.
- In this case the process terminates in a leaf that corresponds to the left blue region in the right-hand panel of the Figure.



The general setting (1)

- More generally, a binary tree T will partition the feature space \mathcal{X} into as many regions as there are leaf nodes.
- Denote the set of leaf nodes by \mathcal{W} . The overall prediction function g that corresponds to the tree can then be written as

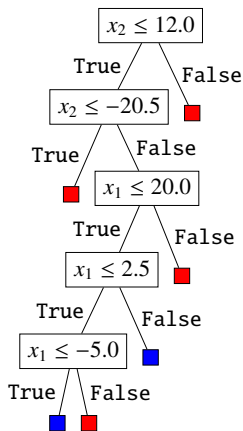
$$g(\mathbf{x}) = \sum_{w \in \mathcal{W}} g^w(\mathbf{x}) \mathbb{1}\{\mathbf{x} \in R_w\}, \quad (1)$$

where $\mathbb{1}$ denotes the indicator function.

- The representation (1) is very general and depends on
 - ① how the regions $\{R_w\}$ are constructed via the logical conditions in the decision tree, and
 - ② how the regional prediction functions of the leaf nodes are defined.
- Simple logical conditions of the form $x_j \leq \xi$ split a Euclidean feature space into rectangles aligned with the axes.

The general setting (2)

For example, the below Figure partitions the feature space into six rectangles: two blue and four red rectangles.



Decision trees for classification and regression

In a classification setting:

- the regional prediction function g^w corresponding to a leaf node w takes values in the set of possible class labels.
- In most cases, as in our tree example, it is taken to be constant on the corresponding region R_w .

In a regression setting:

- g^w is real-valued and also usually takes only one value.
- That is, every feature vector in R_w leads to the same predicted value. Of course, different regions will usually have different predicted values.

Constructing a tree (1)

- Constructing a tree with a training set $\tau = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ amounts to minimizing the training loss

$$\ell_{\tau}(g) = \frac{1}{n} \sum_{i=1}^n \text{Loss}(y_i, g(\mathbf{x}_i)) \quad (2)$$

for some loss function.

- With g of the form $g(\mathbf{x}) = \sum_{w \in \mathcal{W}} g^w(\mathbf{x}) \mathbb{1}\{\mathbf{x} \in R_w\}$, we can write

$$\begin{aligned} \ell_{\tau}(g) &= \frac{1}{n} \sum_{i=1}^n \text{Loss}(y_i, g(\mathbf{x}_i)) = \frac{1}{n} \sum_{i=1}^n \sum_{w \in \mathcal{W}} \mathbb{1}\{\mathbf{x}_i \in R_w\} \text{Loss}(y_i, g(\mathbf{x}_i)) \\ &= \sum_{w \in \mathcal{W}} \underbrace{\frac{1}{n} \sum_{i=1}^n \mathbb{1}\{\mathbf{x}_i \in R_w\} \text{Loss}(y_i, g^w(\mathbf{x}_i))}_{(*)}, \end{aligned}$$

where $(*)$ is the contribution by the regional prediction function g^w to the overall training loss.

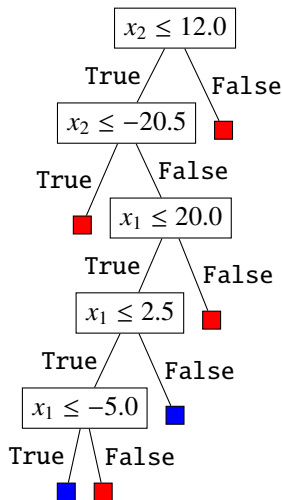
Constructing a tree (2)

- In the case where all $\{\mathbf{x}_i\}$ are different, finding a decision tree T that gives a zero squared-error or zero-one training loss is easy (assignment?).
- but such an “overfitted” tree will have poor predictive behavior, expressed in terms of the generalization risk.
- Instead we consider a restricted class of decision trees and aim to minimize the training loss within that class.
- It is common to use a top-down greedy approach, which can only achieve an approximate minimization of the training loss.

Specifically, we start from the root, and, for every split, consider a rule that minimizes the training loss.

Regional Prediction Functions

- In general, there is no restriction on how to choose the prediction function g^w for a leaf node $v = w$ in Line 2 of the tree construction Algorithm.
- In principle we can train any model from the data; e.g., via linear regression. However, in practice very simple prediction functions are used.



Regional Prediction Functions — classification problems

- In the *classification* setting with class labels $0, \dots, c - 1$, the regional prediction function g^w for leaf node w is usually chosen to be *constant* and equal to the most common class label of the training data in the associated region R_w (ties can be broken randomly).
- More precisely, let n_w be the number of feature vectors in region R_w and let

$$p_z^w = \frac{1}{n_w} \sum_{\{(\mathbf{x}, y) \in \tau : \mathbf{x} \in R_w\}} \mathbb{1}_{\{y=z\}},$$

be the proportion of feature vectors in R_w that have class label $z = 0, \dots, c - 1$. The regional prediction function for node w is chosen to be the constant

$$g^w(\mathbf{x}) = \operatorname{argmax}_{z \in \{0, \dots, c-1\}} p_z^w. \quad (3)$$

Regional Prediction Functions — regression problems

- In the *regression* setting, g^w is usually chosen as the mean response in the region; that is,

$$g^w(\mathbf{x}) = \bar{y}_{R_w} := \frac{1}{n_w} \sum_{\{(\mathbf{x}, y) \in \tau : \mathbf{x} \in R_w\}} y, \quad (4)$$

where n_w is again the number of feature vectors in R_w .

- It is not difficult to show that $g^v(\mathbf{x}) = \bar{y}_{R_w}$ minimizes the squared-error loss with respect to all constant functions, in the region R_w (assignment?).

The termination criterion

- When building a tree, one can define various types of termination conditions.
- For example, we might stop when the number of data points in the tree node (the size of the input σ set in the tree construction algorithm is less than or equal to some predefined number.
- Or we might choose the maximal depth of the tree in advance.
- Another possibility is to stop when there is no significant advantage, in terms of training loss, to split regions.

Ultimately, the quality of a tree is determined by its predictive performance (generalization risk) and the termination condition should aim to strike a balance between minimizing the approximation error and minimizing the statistical error.

Example — Fixed Tree Depth (1)

To illustrate how the tree depth impacts on the generalization risk, consider Figure 2, which shows the typical behavior of the cross-validation loss as a function of the tree depth.

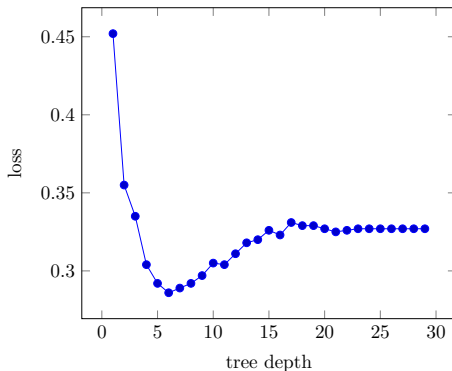


Figure 2: The ten-fold cross-validation loss as a function of the maximal tree depth for a classification problem. The optimal maximal tree depth is here 6.

Example — Fixed Tree Depth (2)

- Recall that the cross-validation loss is an estimate of the expected generalization risk. Complicated (deep) trees tend to overfit the training data by producing many divisions of the feature space.
- As we have seen, this overfitting problem is typical of all learning methods.
- To conclude, increasing the maximal depth does not necessarily result in better performance.

Example — Fixed Tree Depth (3) (TreeDepthCV.py)

```
import numpy as np
from sklearn.datasets import make_blobs
from sklearn.model_selection import cross_val_score
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import zero_one_loss
import matplotlib.pyplot as plt

def ZeroOneScore(clf, X, y):
    y_pred = clf.predict(X)
    return zero_one_loss(y, y_pred)

# Construct the training set
X, y = make_blobs(n_samples=5000, n_features=10, centers=3,
                  random_state=10, cluster_std=10)

# construct a decision tree classifier
clf = DecisionTreeClassifier(random_state=0)
```

Example — Fixed Tree Depth (4)

```
# Cross-validation loss as a function of tree depth (1 to 30)
xdepthlist = []
cvlist = []
tree_depth = range(1,30)
for d in tree_depth:
    xdepthlist.append(d)
    clf.max_depth=d
    cv = np.mean(cross_val_score(clf, X, y, cv=10, scoring=ZeroOneScore))
    cvlist.append(cv)

plt.xlabel('tree depth', fontsize=18, color='black')
plt.ylabel('loss', fontsize=18, color='black')
plt.plot(xdepthlist, cvlist, '-*' , linewidth=0.5)
```

Additional Considerations: Binary Versus Non-Binary Trees

- While it is possible to split a tree node into more than two groups (multiway splits), it generally produces inferior results compared to the simple binary split.
- The major reason is that multiway splits can lead to too many nodes near the tree root that have only a few data points, thus leaving insufficient data for later splits.
- As multiway splits can be represented as several binary splits, the latter is preferred.

Additional Considerations: Data preprocessing

- Sometimes, it can be beneficial to preprocess the data prior to the tree construction.
- For example, PCA can be used with a view to identify the most important dimensions, which in turn will lead to simpler and possibly more informative splitting rules in the internal nodes.

Additional Considerations: Alternative Splitting Rules (1)

- We restricted our attention to splitting rules of the type

$$s(\mathbf{x}) = \mathbb{1}\{x_j \leq \xi\},$$

where $j \in \{1, \dots, p\}$ and $\xi \in \mathbb{R}$.

- These types of rules may not always result in a simple partition of the feature space, as illustrated by the binary data in the Figure on the next slide.
- In this case, the feature space could have been partitioned into just two regions, separated by a straight line.

Additional Considerations: Alternative Splitting Rules (2)

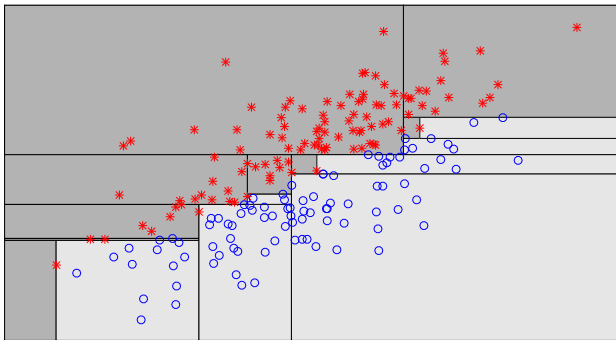


Figure 3: The two groups of points can here be separated by a straight line. Instead, the classification tree divides up the space into many rectangles, leading to an unnecessarily complicated classification procedure.

Additional Considerations: Alternative Splitting Rules (3)

- In this case many classification methods such as linear discriminant analysis will work very well, whereas the classification tree is rather elaborate, dividing the feature set into too many regions.
- An obvious remedy is to use splitting rules of the form

$$s(\mathbf{x}) = \mathbb{1}\{\mathbf{a}^\top \mathbf{x} \leq \xi\}.$$

- In some cases, such as the one just discussed, it may be useful to use a splitting rule that involves several variables, as opposed to a single one.
- The decision regarding the split type clearly depends on the problem domain. For example, for logical (binary) variables our domain knowledge may indicate that a different behavior is expected when both x_i and x_j ($i \neq j$) are True.
- In this case, we will naturally introduce a decision rule of the form:

$$s(\mathbf{x}) = \mathbb{1}\{x_i = \text{True and } x_j = \text{True}\}.$$

Additional Considerations: Categorical Variables

- When an explanatory variable is categorical with labels (levels) say $\{1, \dots, k\}$, the splitting rule is generally defined via a partition of the label set $\{1, \dots, k\}$ into two subsets.
- Specifically, let L and R be a partition of $\{1, \dots, k\}$. Then, the splitting rule is defined via

$$s(\mathbf{x}) = \mathbb{1}\{x_j \in L\}.$$

- For the general supervised learning case, finding the optimal partition in the sense of minimal loss requires one to consider 2^k subsets of $\{1, \dots, k\}$.
- Consequently, finding a good splitting rule for categorical variables can be challenging when the number of labels p is large.

Controlling the Tree Shape (1)

- Eventually, we are interested in getting the right-size tree. Namely, a tree that shows good generalization properties.
- Shallow trees tend to underfit and deep trees tend to overfit the data.
- Basically, a shallow tree does not produce a sufficient number of splits and a deep tree will produce many partitions and thus many leaf nodes.
- If we grow the tree to a sufficient depth, each training sample will occupy a separate leaf and we will observe a zero loss with respect to the training data.
- The above phenomenon is illustrated in Figure 4 (next slide), which presents the cross-validation loss and the training loss as a function of the tree depth.

Controlling the Tree Shape (2)

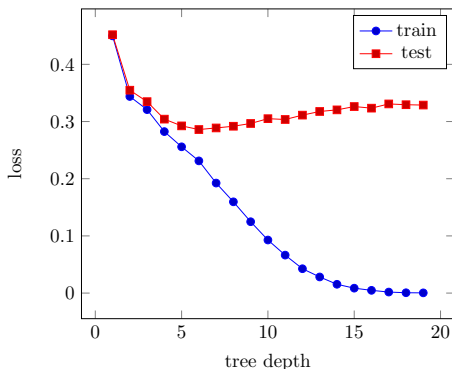


Figure 4: The cross-validation and the training loss as a function of the tree depth for a binary classification problem.

Controlling the Tree Shape (3)

- To address these issues, a so-called *pruning* routine can be employed. The idea is as follows.
- We first grow a very deep tree and then prune (remove nodes) it upwards until we reach the root node.
- Consequently, the pruning process causes the number of tree nodes to decrease.
- While the tree is being pruned, the generalization risk gradually decreases up to the point where it starts increasing again, at which point the pruning is stopped.
- This decreasing/increasing behavior is due to the bias–variance tradeoff.
- To start with, let v and v' be tree nodes. We say that v' is a *descendant* of v if there is a path down the tree, which leads from v to v' .
- If such a path exists, we also say that v is an *ancestor* of v' . Consider the tree in Figure 5 (see next slide).

Controlling the Tree Shape (4)

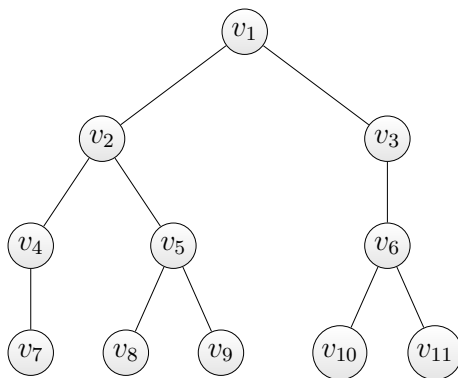


Figure 5: The node v_9 is a descendant of v_2 , and v_2 is an ancestor of $\{v_4, v_5, v_7, v_8, v_9\}$, but v_6 is not a descendant of v_2 .

Controlling the Tree Shape (5)

To formally define pruning, we will require the following Definition 1. An example of pruning is demonstrated in Figure 6.

Definition (Branches and Pruning)

- 1 A tree branch T_v of the tree T is a sub-tree of T rooted at node $v \in T$.
- 2 The pruning of branch T_v from a tree T is performed via deletion of the entire branch T_v from T except the branch's root node v . The resulting pruned tree is denoted by $T - T_v$.
- 3 A sub-tree $T - T_v$ is called a pruned sub-tree of T . We indicate this with the notation $T - T_v \prec T$ or $T \succ T - T_v$.

Controlling the Tree Shape (6)

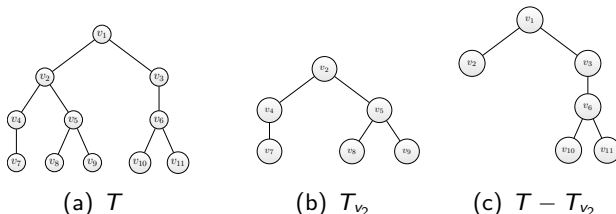


Figure 6: The pruned tree $T - T_{v_2}$ in (c) is the result of pruning the T_{v_2} branch in (b) from the original tree T in (a).

A basic decision tree pruning procedure is summarized in Algorithm 1 (see next slide).

Controlling the Tree Shape (7)

Algorithm 1: Decision Tree Pruning

Input: Training set τ .

Output: Sequence of decision trees $T^0 \succ T^1 \succ \dots$

- 1 Build a large decision tree T^0 via the recursive tree construction algorithm. [A possible termination criterion for that algorithm is to have some small predetermined number of data points at each terminal node of T^0 .]
- 2 $T' \leftarrow T^0$
- 3 $k \leftarrow 0$
- 4 **while** T' has more than one node **do**
 - 5 $k \leftarrow k + 1$
 - 6 Choose $v \in T'$.
 - 7 Prune the branch rooted at v from T' .
 - 8 $T^k \leftarrow T' - T_v$ and $T' \leftarrow T^k$.
- 9 **end**
- 0 **return** T^0, T^1, \dots, T^k

Controlling the Tree Shape (8)

- Let T^0 be the initial (deep) tree and let T^k be the tree obtained after the k -th pruning operation, for $k = 1, \dots, K$.
- As soon as the sequence of trees $T^0 \succ T^1 \succ \dots \succ T^K$ is available, one can choose the best tree of $\{T^k\}_{k=1}^K$ according to the smallest generalization risk.
- Specifically, we can split the data into training and validation sets. In this case, Algorithm 1 is executed using the training set and the generalization risks of $\{T^k\}_{k=1}^K$ are estimated via the validation set.
- While Algorithm 1 and the corresponding best tree selection process look appealing, there is still an important question to consider; namely, how to choose the node v and the corresponding branch T_v in Line 6 of the algorithm.
- In order to overcome this problem, Breiman proposed a method called *cost complexity pruning*, which we discuss next.

Controlling the Tree Shape (9)

- Let $T \prec T^0$ be a tree obtained via pruning of a tree T^0 .
- Denote the set of leaf (terminal) nodes of T by \mathcal{W} . The number of leaves $|\mathcal{W}|$ is a measure for the complexity of the tree; recall that $|\mathcal{W}|$ is the number of regions $\{R_w\}$ in the partition of \mathcal{X} .
- Corresponding to each tree T is a prediction function g , as in (1).
- In cost-complexity pruning the objective is to find a prediction function g (or, equivalently, tree T) that minimizes the training loss $\ell_T(g)$ while taking into account the complexity of the tree.
- The idea is to regularize the training loss, by adding a penalty term for the complexity of the tree. This leads to the following definition.

Controlling the Tree Shape (10)

Definition (Cost-Complexity Measure)

Let $\tau = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ be a data set and $\gamma \geq 0$ be a real number. For a given tree T , the cost-complexity measure $C_\tau(\gamma, T)$ is defined as:

$$\begin{aligned} C_\tau(\gamma, T) &:= \frac{1}{n} \sum_{w \in \mathcal{W}} \left(\sum_{i=1}^n \mathbb{1}\{\mathbf{x}_i \in R_w\} \text{Loss}(y_i, g^w(\mathbf{x}_i)) \right) + \gamma |\mathcal{W}| \\ &= \ell_\tau(g) + \gamma |\mathcal{W}|, \end{aligned}$$

where $\ell_\tau(g)$ is the training loss.

Controlling the Tree Shape (11)

$$C_{\tau}(\gamma, T) := \ell_{\tau}(g) + \gamma |W|,$$

- Small values of γ result in a small penalty for the tree complexity $|W|$, and thus large trees (that fit the entire *training* data well) will minimize the measure $C_{\tau}(\gamma, T)$. In particular, for $\gamma = 0$, $T = T^0$ will be the minimizer of $C_{\tau}(\gamma, T)$.
- On the other hand, large values of γ will prefer smaller trees or, more precisely, trees with fewer leaves. For sufficiently large γ , the solution T will collapse to a single (root) node.
- It can be shown that, for every value of γ , there exists a smallest minimizing sub-tree with respect to the cost-complexity measure. In practice, a suitable γ is selected via observing the performance of the learner on the validation set or by cross-validation.

Advantages of Decision Trees (1)

- The tree structure can handle both categorical and numerical features in a natural and straightforward way. Specifically, there is no need to pre-process categorical features, say via the introduction of dummy variables.
- The final tree obtained after the training phase can be compactly stored for the purpose of making predictions for new feature vectors. The prediction process only involves a single tree traversal from the tree root to a leaf.
- The hierarchical nature of decision trees allows for an efficient encoding of the feature's conditional information. Specifically, after an internal split of a feature x_j via the standard splitting rule we will only consider such subsets of data that were constructed based on this split, thus implicitly exploiting the corresponding conditional information from the initial split of x_j .

Advantages of Decision Trees (2)

- The tree structure can be easily understood and interpreted by domain experts with little statistical knowledge, since it is essentially a logical decision flow diagram.
- Decision trees are invariant under monotone transformations of the data. To see this, consider the (optimal) splitting rule $s(\mathbf{x}) = \mathbb{1}\{x_3 \leq 2\}$, where x_3 is a positive feature. Suppose that x_3 is transformed to $x'_3 = x_3^2$. Now, the optimal splitting rule will take the form $s(\mathbf{x}) = \mathbb{1}\{x'_3 \leq 4\}$.

Advantages of Decision Trees (3)

- In the classification setting, it is common to report not only the predicted value of a feature vector, but also the respective class probabilities. Decision trees handle this task without any additional effort. Specifically, consider a new feature vector. During the estimation process, we will perform a tree traversal and the point will end up in a certain leaf w . The probability of this feature vector lying in class z can be estimated as the proportion of training points in w that are in class z .
- As each training point is treated equally in the construction of a tree, their structure of the tree will be relatively robust to outliers. In a way, trees exhibit a similar kind of robustness as the sample median does for real-valued data.

Limitations of Decision Trees

- Despite the fact that the decision trees are extremely interpretable, the predictive accuracy is generally inferior to other established statistical learning methods.
- In addition, decision trees, and in particular very deep trees that were not subject to pruning, are heavily reliant on their training set.
- A small change in the training set can result in a dramatic change of the resulting decision tree.
- Their inferior predictive accuracy, however, is a direct consequence of the bias–variance tradeoff. Specifically, a decision tree model generally exhibits a high variance.
- To overcome the above limitations, several promising approaches such as *bagging*, *random forest*, and *boosting* were developed.

Bootstrap Aggregation (1)

- The major idea of the *bootstrap aggregation* or *bagging* method is to combine prediction functions learned from multiple data sets, with a view to improving overall prediction accuracy.
- Bagging is especially beneficial when dealing with predictors that tend to overfit the data, such as in decision trees, where the (unpruned) tree structure is very sensitive to small changes in the training set.
- To start with, consider an idealized setting for a regression tree, where we have access to B iid copies¹ $\mathcal{T}_1, \dots, \mathcal{T}_B$ of a training set \mathcal{T} .
- Then, we can train B separate regression models (B different decision trees) using these sets, giving learners $g_{\mathcal{T}_1}, \dots, g_{\mathcal{T}_B}$, and take their average:

$$g_{\text{avg}}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B g_{\mathcal{T}_b}(\mathbf{x}). \quad (5)$$

¹In this section \mathcal{T}_k means the k -th training set, not a training set of size k .

Bootstrap Aggregation (2)

- By the law of large numbers, as $B \rightarrow \infty$, the average prediction function converges to the expected prediction function $g^\dagger := \mathbb{E}g_{\mathcal{T}}$.
- The following result shows that using g^\dagger as a prediction function (if it were known) would result in an expected squared-error generalization risk that is less than or equal to the expected generalization risk for a general prediction function $g_{\mathcal{T}}$.
- It thus suggests that taking an average of prediction functions may lead to a better expected squared-error generalization risk.

Theorem

Expected Squared-Error Generalization Riskbootstrap Let \mathcal{T} be a random training set and let \mathbf{X}, Y be a random feature vector and response that are independent of \mathcal{T} . Then,

$$\mathbb{E}\left(Y - g_{\mathcal{T}}(\mathbf{X})\right)^2 \geq \mathbb{E}\left(Y - g^\dagger(\mathbf{X})\right)^2.$$

Bootstrap Aggregation (3)

Proof.

We have

$$\begin{aligned}\mathbb{E} \left[\left(Y - g_{\mathcal{T}}(\mathbf{X}) \right)^2 \mid \mathbf{X}, Y \right] &\geq \left(\mathbb{E}[Y \mid \mathbf{X}, Y] - \mathbb{E}[g_{\mathcal{T}}(\mathbf{X}) \mid \mathbf{X}, Y] \right)^2 \\ &= \left(Y - g^{\dagger}(\mathbf{X}) \right)^2,\end{aligned}$$

where the inequality follows from $\mathbb{E}U^2 \geq (\mathbb{E}U)^2$ for any (conditional) expectation. Consequently, by the tower property,

$$\mathbb{E} \left(Y - g_{\mathcal{T}}(\mathbf{X}) \right)^2 = \mathbb{E} \left[\mathbb{E} \left[\left(Y - g_{\mathcal{T}}(\mathbf{X}) \right)^2 \mid \mathbf{X}, Y \right] \right] \geq \mathbb{E} \left(Y - g^{\dagger}(\mathbf{X}) \right)^2.$$



Bootstrap Aggregation (4)

- Unfortunately, multiple independent data sets are rarely available. But we can substitute them by bootstrapped ones.
- Specifically, instead of the $\mathcal{T}_1, \dots, \mathcal{T}_B$ sets, we can obtain random training sets $\mathcal{T}_1^*, \dots, \mathcal{T}_B^*$ by resampling them from a *single* (fixed) training set τ , and use them to train B separate models.
- By model averaging we obtain the bootstrapped aggregated estimator or *bagged estimator of the form*:

$$g_{\text{bag}}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B g_{\mathcal{T}_b^*}(\mathbf{x}). \quad (6)$$

The Bootstrap Sampling Algorithm

Algorithm 2: Bootstrap Aggregation Sampling

Input: Training set $\tau = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and resample size B .

Output: Bootstrapped data sets.

```
1 for  $b = 1$  to  $B$  do
2    $\mathcal{T}_b^* \leftarrow \emptyset$ 
3   for  $i = 1$  to  $n$  do
4     Draw  $U \sim \mathcal{U}(0, 1)$ ;
5      $I \leftarrow \lceil nU \rceil$            // select random index;
6      $\mathcal{T}_b^* \leftarrow \mathcal{T}_b^* \cup \{(\mathbf{x}_I, y_I)\}$ .
7   end
8 end
9 return  $\mathcal{T}_b^*, b = 1, \dots, B$ .
```

Bootstrap Aggregation for Classification Problems

- Note that the bagging idea was discussed in the context of regression problems.
- However, the bagging idea can be readily extended to handle classification settings as well.
- For example, g_{bag} can take the majority vote among $\{g_{\mathbf{T}_b^*}\}, b = 1, \dots, B$; that is, to accept the most frequent class among B predictors.
- While bagging can be applied for any statistical model (such as decision trees, neural networks, linear regression, K -nearest neighbors, and so on), it is most effective for predictors that are sensitive to small changes in the training set.

Where can we use the Bootstrap Aggregation? (1)

- The reason becomes clear when we decompose the expected generalization risk as

$$\mathbb{E} \ell(g_{\mathcal{T}}) = \ell^* + \underbrace{\mathbb{E} (\mathbb{E}[g_{\mathcal{T}}(\mathbf{X}) \mid \mathbf{X}] - g^*(\mathbf{X}))^2}_{\text{expected squared bias}} + \underbrace{\mathbb{E}[\text{Var}[g_{\mathcal{T}}(\mathbf{X}) \mid \mathbf{X}]]}_{\text{expected variance}}, \quad (7)$$

- Compare this with the same decomposition for the average prediction function g_{bag} in $g_{\text{avg}}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B g_{\mathcal{T}_b}(\mathbf{x})$.
- As $\mathbb{E} g_{\text{bag}}(\mathbf{x}) = \mathbb{E} g_{\mathcal{T}}(\mathbf{x})$, we see that any possible improvement in the generalization risk must be due to the expected variance term.
- **Averaging and bagging are thus only useful for predictors with a large expected variance, relative to the other two terms.**

Where can we use the Bootstrap Aggregation? (2)

- Examples of such “unstable” predictors include decision trees, neural networks, and subset selection in linear regression.
- On the other hand, “stable” predictors are insensitive to small data changes, an example being the K -nearest neighbors method.
- Note that for independent training sets $\mathcal{T}_1, \dots, \mathcal{T}_B$ a reduction of the variance by a factor B is achieved:

$$\text{Var } g_{\text{bag}}(\mathbf{x}) = B^{-1} \text{Var } g_{\mathcal{T}}(\mathbf{x}).$$

- Again, it depends on the squared bias and irreducible loss how significant this reduction is for the generalization risk.

Limitations of Bagging

- It is important to remember that g_{bag} is not exactly equal to g_{avg} , which in turn is not exactly g^\dagger .
- Specifically, g_{bag} is constructed from the bootstrap approximation of the sampling pdf f .
- As a consequence, for stable predictors, it can happen that g_{bag} will perform worse than g_T .
- In addition to the deterioration of the bagging performance for stable procedures, it can also happen that g_T has already achieved a near optimal predictive accuracy given the available training data.
- In this case, bagging will not introduce a significant improvement.

Out-of-Bag Loss Estimation (1)

- The bagging process provides an opportunity to estimate the generalization risk of the bagged model without an additional test set.
- Specifically, recall that we obtain the $\mathcal{T}_1^*, \dots, \mathcal{T}_B^*$ sets from a single training set τ by sampling via the bootstrap sampling algorithm, and use them to train B separate models.
- It can be shown that, for large sample sizes, on average about a third (more precisely, a fraction $e^{-1} \approx 0.37$) of the original sample points are not included in bootstrapped set \mathcal{T}_b^* for $1 \leq b \leq B$.
- Therefore, these samples can be used for the loss estimation. These samples are called *out-of-bag* (OOB) observations.

Out-of-Bag Loss Estimation (2)

- For each sample from the original data set, we calculate the OOB loss using predictors that were trained without this particular sample.
- Tibshirani observe that the OOB loss is almost identical to the n -fold cross-validation loss.
- In addition, the OOB loss can be used to determine the number of trees required.
- Specifically, we can train predictors until the OOB loss stops changing. Namely, decision trees are added until the OOB loss stabilizes.

The Out-of-Bag Algorithm

Algorithm 3: Out-of-Bag Loss Estimation

Input: The original data set $\tau = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, the bootstrapped data sets $\{\mathcal{T}_1^*, \dots, \mathcal{T}_B^*\}$, and the trained predictors $\{g_{\mathcal{T}_1^*}, \dots, g_{\mathcal{T}_B^*}\}$.

Output: Out-of-bag loss for the averaged model.

```
1 for  $i = 1$  to  $n$  do
2    $C_i \leftarrow \emptyset$ ; // Indices of predictors not depending on  $(\mathbf{x}_i, y_i)$ 
3   for  $b = 1$  to  $B$  do
4     if  $(\mathbf{x}_i, y_i) \notin \mathcal{T}_b^*$  then  $C_i \leftarrow C_i \cup \{b\}$ ;
5   end
6    $Y_i' \leftarrow |C_i|^{-1} \sum_{b \in C_i} g_{\mathcal{T}_b^*}(\mathbf{x}_i)$ 
7    $L_i \leftarrow \text{Loss}(y_i, Y_i')$ 
8 end
9  $L_{\text{OOB}} \leftarrow \frac{1}{n} \sum_{i=1}^n L_i$ 
10 return  $L_{\text{OOB}}$ .
```

Bagging for a Regression Tree Example (1)

BaggingExample.py

Consider a basic bagging example for a regression tree, in which we compare the decision tree estimator with the corresponding bagged estimator. We use the R^2 metric (coefficient of determination) for comparison.

```
import numpy as np
from sklearn.datasets import make_friedman1
from sklearn.tree import DecisionTreeRegressor
from sklearn.model_selection import train_test_split
from sklearn.metrics import r2_score

# create regression problem
n_points = 1000 # points
x, y = make_friedman1(n_samples=n_points, n_features=15,
                      noise=1.0, random_state=100)

# split to train/test set
x_train, x_test, y_train, y_test = \
    train_test_split(x, y, test_size=0.33, random_state=100)
```

Bagging for a Regression Tree Example (2)

Simple regression tree (single tree) code.

```
# training
regTree = DecisionTreeRegressor(random_state=100)
regTree.fit(x_train,y_train)
# test
yhat = regTree.predict(x_test)

print("DecisionTreeRegressor R^2 score = ",r2_score(y_test, yhat))
```

Bagging for a Regression Tree Example (3)

Constructing bagged estimator with 500 bootstrapped datasets.

```
# Bagging construction
n_estimators=500
bag = np.empty((n_estimators), dtype=object)
bootstrap_ds_arr = np.empty((n_estimators), dtype=object)
for i in range(n_estimators):
    # sample bootstrapped data set
    ids = np.random.choice(range(0,len(x_test)),size=len(x_test),
                           replace=True)
    x_boot = x_train[ids]
    y_boot = y_train[ids]
    bootstrap_ds_arr[i] = np.unique(ids)

    bag[i] = DecisionTreeRegressor()
    bag[i].fit(x_boot,y_boot)
```

Bagging for a Regression Tree Example (4)

Constructing prediction and OOB estimator.

```
# bagging prediction
yhatbag = np.zeros(len(y_test))
for i in range(n_estimators):
    yhatbag = yhatbag + bag[i].predict(x_test)

yhatbag = yhatbag/n_estimators

# out of bag loss estimation
oob_pred_arr = np.zeros(len(x_train))
for i in range(len(x_train)):
    x = x_train[i].reshape(1, -1)
    C = []
    for b in range(n_estimators):
        if(np.isin(i, bootstrap_ds_arr[b])==False):
            C.append(b)
    for pred in bag[C]:
        oob_pred_arr[i] = oob_pred_arr[i] + (pred.predict(x)/len(C))
}
L_oob = r2_score(y_train, oob_pred_arr)
```

Bagging for a Regression Tree Example (5)

```
DecisionTreeRegressor R^2 score = 0.575438224929718  
Bagging R^2 score = 0.7612121189201985  
Bagging OOB R^2 score = 0.7758253149069059
```

The results are as follows.

- The decision tree bagging improves the test-set R^2 score by about 32% (from 0.575 to 0.761).
- Moreover, the OOB score (0.776) is very close to the true generalization risk (0.761) of the bagged estimator.

Random Forests (1)

- We discussed the intuition behind the prediction averaging procedure.
- Specifically, for some feature vector \mathbf{x} let $Z_b = g_{\mathbb{T}_b}(\mathbf{x})$, $b = 1, 2, \dots, B$ be iid prediction values, obtained from independent training sets $\mathbb{T}_1, \dots, \mathbb{T}_B$.
- Suppose that $\text{Var } Z_b = \sigma^2$ for all $b = 1, \dots, B$. Then the variance of the average prediction value \bar{Z}_B is equal to σ^2/B .
- However, if bootstrapped data sets $\{\mathbb{T}_b^*\}$ are used instead, the corresponding random variables $\{Z_b\}$ will be *correlated*.
- In particular, $Z_b = g_{\mathbb{T}_b^*}(\mathbf{x})$ for $b = 1, \dots, B$ are identically distributed (but not independent) with some positive pairwise correlation ρ . It then holds that

$$\text{Var } \bar{Z}_B = \rho \sigma^2 + \sigma^2 \frac{(1 - \rho)}{B}. \quad (8)$$

Random Forests (2)

- While the second term of

$$\text{Var } \bar{Z}_B = \rho \sigma^2 + \sigma^2 \frac{(1 - \rho)}{B}$$

goes to zero as the number of observation B increases, the first term remains constant.

- This issue is particularly relevant for bagging with decision trees.
- For example, consider a situation in which there exists a feature that provides a very good split of the data.
- Such a feature will be selected and split for every $\{g_{T_b}^*\}_{b=1}^B$ at the root level and we will consequently end up with highly correlated predictions.
- In such a situation, prediction averaging will not introduce the desired improvement in the performance of the bagged predictor.

Random Forests (3)

- The major idea of random forests is to perform bagging in combination with a “decorrelation” of the trees by including only a subset of features during the tree construction.
- For each bootstrapped training set T_b^* we build a decision tree, where for each split in this tree, we will use a randomly selected subset of $m \leq p$ features for the splitting rules.
- This simple but powerful idea will decorrelate the trees, since strong predictors will have a smaller chance to be considered at the root levels.
- Consequentially, we can expect to improve the predictive performance of the bagged estimator. The resulting predictor (random forest) construction is summarized in the following Algorithm.

Random Forests Construction Algorithm

Algorithm 4: Random Forest Construction

Input: Training set $\tau = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, the number of trees in the forest B , and the number $m \leq p$ of features to be included, where p is the total number of features in \mathbf{x} .

Output: Ensemble of trees.

- 1 Generate bootstrapped training sets $\{\mathcal{T}_1^*, \dots, \mathcal{T}_B^*\}$ via Algorithm 2.
 - 2 **for** $b = 1$ **to** B **do**
 - 3 Randomly select m out of p features, without replacement for each split in the tree.
 - 4 Using these, train a decision tree $g_{\mathcal{T}_b^*}$ via the tree construction algorithm.
 - 5 **end**
 - 6 **return** $\{g_{\mathcal{T}_b^*}\}_{b=1}^B$.
-

Random Forests Example: BaggingExampleRF.py

- For regression problems, the output of Algorithm 4 is combined to yield the random forest prediction function:

$$g_{\text{RF}}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B g_{\mathbf{T}_b^*}(\mathbf{x}).$$

- In the classification setting, we take instead the majority vote from the $\{g_{\mathbf{T}_b^*}\}$.

Generally, a random forest will outperform bagging estimators. We next continue the example in which we compared the decision tree estimator with the corresponding bagged estimator. It can be seen that the random forest's R^2 score is outperforming that of the bagged estimator. Specifically, we obtained:

RF R^2 score = 0.8106589580845707

RF OOB R^2 score = 0.8260541058404149

```
from sklearn.datasets import make_friedman1
from sklearn.model_selection import train_test_split
from sklearn.metrics import r2_score
from sklearn.ensemble import RandomForestRegressor

# create regression problem
n_points = 1000 # points
x, y = make_friedman1(n_samples=n_points, n_features=15,
                      noise=1.0, random_state=100)
# split to train/test set
x_train, x_test, y_train, y_test = \
    train_test_split(x, y, test_size=0.33, random_state=100)
rf = RandomForestRegressor(n_estimators=500, oob_score = True, max_features=8,
                          random_state=100)
rf.fit(x_train,y_train)
yhatrf = rf.predict(x_test)

print("RF R^2 score = ", r2_score(y_test, yhatrf),
      "\nRF OOB R^2 score = ", rf.oob_score_)
```

Random Forests — The Optimal Number of Subset Features m

- The default values for m are $\lfloor p/3 \rfloor$ and $\lfloor \sqrt{p} \rfloor$ for regression and classification setting, respectively.
- However, the standard practice is to treat m as a hyperparameter that requires tuning, depending on the specific problem at hand.

Note that the procedure of bagging decision trees is a special case of a random forest construction. Consequently, the OOB loss is readily available for random forests.

Random Forests — inference

Recall the advertisement example — what is the the most promising domain of advertisement (radio/tv/internet).

- While the advantage of bagging in the sense of enhanced accuracy is clear, we should also consider its negative aspects and, in particular, the loss of interpretability.
- Specifically a random forest consists of many trees, thus making the prediction process both hard to visualize and interpret.
- For example, given a random forest, it is not easy to determine a subset of features that are essential for accurate prediction.
- The feature importance measure intends to address this issue. The idea is as follows.

Random Forests — feature importance (1)

- Each *internal* node of a decision tree induces a certain decrease in the training loss.
- Let us denote this decrease in the training loss by $\Delta_{\text{Loss}}(v)$, where v is not a leaf node of T .
- In addition, recall that for splitting rules of the type $\mathbb{1}\{x_j \leq \xi\}$ ($1 \leq j \leq p$), each node v is associated with a feature x_j that determines the split.
- Using the above definitions, we can define the *feature importance* of x_j as

$$\mathcal{I}_T(x_j) = \sum_{v \text{ internal} \in T} \Delta_{\text{Loss}}(v) \mathbb{1}\{x_j \text{ is associated with } v\}, \quad 1 \leq j \leq p.$$

Random Forests — feature importance (2)

- While

$$\mathcal{I}_T(x_j) = \sum_{v \text{ internal} \in T} \Delta_{\text{Loss}}(v) \mathbb{1}\{x_j \text{ is associated with } v\}, \quad 1 \leq j \leq p.$$

is defined for a single tree, it can be readily extended to random forests.

- Specifically, the feature importance in that case will be averaged over all trees of the forest; that is, for a forest consisting of B trees $\{T_1, \dots, T_B\}$, the feature importance measure is:

$$\mathcal{I}_{\text{RF}}(x_j) = \frac{1}{B} \sum_{b=1}^B \mathcal{I}_{T_b}(x_j), \quad 1 \leq j \leq p.$$

Random Forests — feature importance example (1)

VarImportance.py

- We consider a classification problem with 15 features.
- The data is specifically designed to contain only 5 informative features out of 15.
- In the code below, we apply the random forest procedure and calculate the corresponding feature importance measures.

```
import numpy as np
from sklearn.datasets import make_classification
from sklearn.ensemble import RandomForestClassifier
import matplotlib.pyplot as plt, pylab

n_points = 1000 # create regression data with 1000 data points
x, y = make_classification(n_samples=n_points, n_features=15,
    n_informative=5, n_redundant=0, n_repeated=0, random_state=100,
    shuffle=False)
```


Random Forests — feature importance example (2)

```
rf = RandomForestClassifier(n_estimators=200, max_features="log2")
rf.fit(x,y)

importances = rf.feature_importances_
indices = np.argsort(importances)[::-1]

for f in range(15):
    print("Feature %d (%f)" % (indices[f]+1, importances[indices[f]]))

std = np.std([rf.feature_importances_ for tree in rf.estimators_],
             axis=0)
f = plt.figure()
plt.bar(range(x.shape[1]), importances[indices],
        color="b", yerr=std[indices], align="center")
plt.xticks(range(x.shape[1]), indices+1)
plt.xlim([-1, x.shape[1]])
pylab.xlabel("feature index")
pylab.ylabel("importance")
plt.show()
```

Random Forests — feature importance example

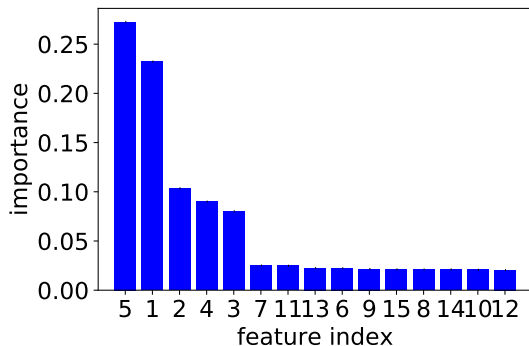


Figure 7: Importance measure for the 15-feature data set with only 5 informative features x_1, x_2, x_3, x_4 , and x_5 . Clearly, it is hard to visualize and understand the prediction process based on 200 trees. However, Figure 7 shows that the features x_1, x_2, x_3, x_4 , and x_5 were correctly identified as being important.

Boosting (1)

- Boosting is a powerful idea that aims to improve the accuracy of any learning algorithm, especially when involving *weak learners* — simple prediction functions that exhibit performance slightly better than random guessing. Shallow decision trees typically yield weak learners.
- Originally, boosting was developed for binary classification tasks, but it can be readily extended to handle general classification and regression problems.
- The boosting approach has some similarity with the bagging method in the sense that boosting uses an ensemble of prediction functions. Despite this similarity, there exists a fundamental difference between these methods.
- Specifically, while bagging involves the fitting of prediction functions to bootstrapped data, the predicting functions in boosting are learned *sequentially*.
- That is, each learner uses information from previous learners.

Boosting (2)

- The idea is to start with a simple model (weak learner) g_0 for the data $\tau = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and then to improve or “boost” this learner to a learner

$$g_1 := g_0 + h_1.$$

- Here, the function h_1 is found by minimizing the training loss for $g_0 + h_1$ over all functions h in some class of functions \mathcal{H} .
- For example, \mathcal{H} could be the set of prediction functions that can be obtained via a decision tree of maximal depth 2.

Given a loss function Loss , the function h_1 is thus obtained as the solution to the optimization problem

$$h_1 = \operatorname{argmin}_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \text{Loss}(y_i, g_0(\mathbf{x}_i) + h(\mathbf{x}_i)). \quad (9)$$

Boosting (3)

- This process can be repeated for g_1 to obtain $g_2 = g_1 + h_2$, and so on, yielding the boosted prediction function

$$g_B(\mathbf{x}) = g_0 + \sum_{b=1}^B h_b(\mathbf{x}). \quad (10)$$

- Instead of using the updating step

$$g_b = g_{b-1} + h_b,$$

one prefers to use the smooth updating step

$$g_b = g_{b-1} + \gamma h_b,$$

for some suitably chosen step-size parameter γ . As we shall see shortly, this helps reduce overfitting.

Boosting (4)

- Boosting can be used for regression and classification problems.
- We start with a simple regression setting, using the squared-error loss; thus, $\text{Loss}(y, \hat{y}) = (y - \hat{y})^2$.
- In this case, it is common to start with $g_0(\mathbf{x}) = n^{-1} \sum_{i=1}^n y_i$, and each h_b for $b = 1, \dots, B$ is chosen as a learner for the data set τ_b of residuals corresponding to g_{b-1} .
- That is, $\tau_b := \left\{ \left(\mathbf{x}_i, e_i^{(b)} \right) \right\}_{i=1}^n$, with

$$e_i^{(b)} := y_i - g_{b-1}(\mathbf{x}_i). \quad (11)$$

- This leads to the following boosting procedure for regression with squared-error loss.

Boosting — Construction Algorithm

Algorithm 5: Regression Boosting with Squared-Error Loss

Input: Training set $\tau = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, the number of boosting rounds B , and a shrinkage step-size parameter γ .

Output: Boosted prediction function.

```
1 Set  $g_0(\mathbf{x}) \leftarrow n^{-1} \sum_{i=1}^n y_i$ .
2 for  $b = 1$  to  $B$  do
3   Set  $e_i^{(b)} \leftarrow y_i - g_{b-1}(\mathbf{x}_i)$  for  $i = 1, \dots, n$ , and let
      $\tau_b \leftarrow \left\{ \left( \mathbf{x}_i, e_i^{(b)} \right) \right\}_{i=1}^n$ .
4   Fit a prediction function  $h_b$  on the training data  $\tau_b$ .
5   Set  $g_b(\mathbf{x}) \leftarrow g_{b-1}(\mathbf{x}) + \gamma h_b(\mathbf{x})$ .
6 end
7 return  $g_B$ .
```

Boosting step size (1)

- The *step-size parameter* γ introduced in Algorithm 5 controls the speed of the fitting process.
- Specifically, for small values of γ , boosting takes smaller steps towards the training loss minimization.
- The step-size γ is of great practical importance, since it helps the boosting algorithm to avoid overfitting.
- This phenomenon is demonstrated in the Figure on the next slide.

Boosting step size (2)

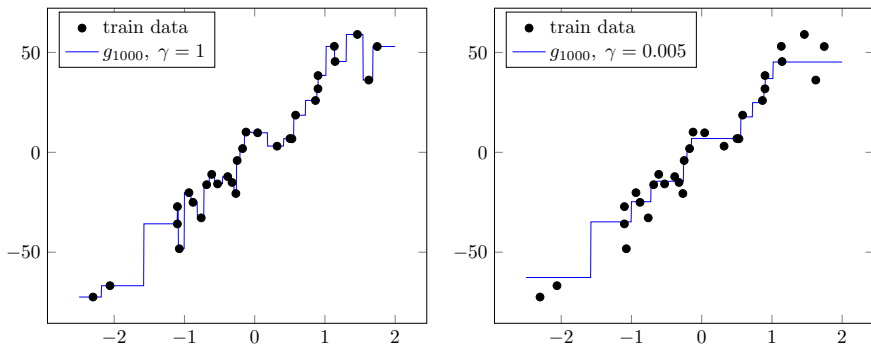


Figure 8: The left and the right panels show the fitted boosting regression model g_{1000} with $\gamma = 1.0$ and $\gamma = 0.005$, respectively. Note the overfitting on the left.

A very basic implementation of 1 which reproduces Figure 8 is provided below.

Boosting — basic example (1) RegressionBoosting.py

```
# create data set
x,y = make_regression(n_samples=sz, n_features=1, n_informative=1,noise=10.0)

# boosting algorithm
BoostingRounds = 1000
alphas = [1, 0.005]

for alpha in alphas:
    g_0, g_boost = TrainBoost(alpha,BoostingRounds,x,y)
    yhat = Predict(g_0, g_boost, alpha, x)

    # plot
    tmpX = np.reshape(np.linspace(-2.5,2,1000),(1000,1))
    yhatX = Predict(g_0, g_boost, alpha, tmpX)
    f = plt.figure()
    plt.plot(x,y,'*')
    plt.plot(tmpX,yhatX)
    plt.show()
```

Boosting — basic example (2)

```
def TrainBoost(alpha, BoostingRounds, x, y):  
    g_0 = np.mean(y)  
    residuals = y - alpha * g_0  
    # list of basic regressor  
    g_boost = []  
  
    for i in range(BoostingRounds):  
        h_i = DecisionTreeRegressor(max_depth=1)  
        h_i.fit(x, residuals)  
        residuals = residuals - alpha * h_i.predict(x)  
        g_boost.append(h_i)  
    return g_0, g_boost  
  
def Predict(g_0, g_boost, alpha, x):  
    yhat = alpha * g_0 * np.ones(len(x))  
    for j in range(len(g_boost)):  
        yhat = yhat + alpha * g_boost[j].predict(x)  
  
    return yhat
```

Gradient Boosting (1)

- The parameter γ can be viewed as a step size made in the direction of the negative gradient of the squared-error training loss.
- To see this, note that the negative gradient

$$-\left. \frac{\partial \text{Loss}(y_i, z)}{\partial z} \right|_{z=g_{b-1}(\mathbf{x}_i)} = -\left. \frac{\partial (y_i - z)^2}{\partial z} \right|_{z=g_{b-1}(\mathbf{x}_i)} = 2(y_i - g_{b-1}(\mathbf{x}_i))$$

is two times the residual $e_i^{(b)}$ given in

$$e_i^{(b)} := y_i - g_{b-1}(\mathbf{x}_i).$$

This was used in the Boosting Algorithm to fit the prediction function h_b .

Gradient Boosting (2)

- In fact, one of the major advances in the theory of boosting was the recognition that one can use a similar gradient descent method for any differentiable loss function.
- The resulting algorithm is called *gradient boosting*.
- The general gradient boosting algorithm is summarized on the next slide.
- The main idea is to mimic a gradient descent algorithm in the following sense. At each stage of the boosting procedure, we calculate a negative gradient on n training points $\mathbf{x}_1, \dots, \mathbf{x}_n$ (Lines 3–4).
- Then, we fit a simple model (such as a shallow decision tree) to approximate the gradient (Line 6) for any feature \mathbf{x} .
- Finally, similar to the gradient descent method, we make a γ -sized step in the direction of the negative gradient (Line 7).

Algorithm 6: Gradient Boosting

Input: Training set $\tau = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, the number of boosting rounds B , a differentiable loss function $\text{Loss}(y, \hat{y})$, and a gradient step-size parameter γ .

Output: Gradient boosted prediction function.

1 Set $g_0(\mathbf{x}) \leftarrow 0$.

2 **for** $b = 1$ **to** B **do**

3 **for** $i = 1$ **to** n **do**

4 Evaluate the negative gradient of the loss at (\mathbf{x}_i, y_i) via

$$r_i^{(b)} \leftarrow - \left. \frac{\partial \text{Loss}(y_i, z)}{\partial z} \right|_{z=g_{b-1}(\mathbf{x}_i)} \quad i = 1, \dots, n.$$

5 **end**

6 Approximate the negative gradient by solving

$$h_b = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \left(r_i^{(b)} - h(\mathbf{x}_i) \right)^2. \quad (12)$$

7 Set $g_b(\mathbf{x}) \leftarrow g_{b-1}(\mathbf{x}) + \gamma h_b(\mathbf{x})$.

8 **end**

9 **return** g_B

Gradient Boosting Example GradientBoostingRegression.py

- We continue with the basic bagging and random forest examples for a regression tree, where we compared the standard decision tree estimator with the corresponding bagging and random forest estimators.
- Now, we use the gradient boosting estimator as implemented in `sklearn`.
- We use $\gamma = 0.1$ and perform $B = 100$ boosting rounds.
- As a prediction function h_b for $b = 1, \dots, B$ we use small regression trees of depth at most 3. Note that such individual trees do not usually give good performance; that is, they are weak prediction functions.
- We can see that the resulting boosting prediction function gives the R^2 score equal to 0.899, which is better than R^2 scores of simple decision tree (0.5754), the bagged tree (0.761), and the random forest (0.8106).

```

import numpy as np
from sklearn.datasets import make_friedman1
from sklearn.tree import DecisionTreeRegressor
from sklearn.model_selection import train_test_split
from sklearn.metrics import r2_score
# create regression problem
n_points = 1000 # points
x, y = make_friedman1(n_samples=n_points, n_features=15,
                      noise=1.0, random_state=100)
# split to train/test set
x_train, x_test, y_train, y_test = \
    train_test_split(x, y, test_size=0.33, random_state=100)

# boosting sklearn
from sklearn.ensemble import GradientBoostingRegressor
breg = GradientBoostingRegressor(learning_rate=0.1,
                                n_estimators=100, max_depth=3, random_state=100)
breg.fit(x_train,y_train)
yhat = breg.predict(x_test)
print("Gradient Boosting R^2 score = ",r2_score(y_test, yhat))
-----
>>Gradient Boosting R^2 score = 0.8993055635639531

```