Regularisation of CART trees by summation of p-values

Nils Engler^{*}, Mathias Lindholm[†], Filip Lindskog[‡] and Taariq Nazar[§]

May 27, 2025

Abstract

The standard procedure to decide on the complexity of a CART regression tree is to use cross-validation with the aim of obtaining a predictor that generalises well to unseen data. The randomness in the selection of folds implies that the selected CART tree is not a deterministic function of the data. We propose a deterministic in-sample method that can be used for stopping the growing of a CART tree based on node-wise statistical tests. This testing procedure is derived using a connection to change point detection, where the null hypothesis corresponds to that there is no signal. The suggested *p*-value based procedure allows us to consider covariate vectors of arbitrary dimension and allows us to bound the *p*-value of an entire tree from above. Further, we show that the test detects a not-too-weak signal with a high probability, given a not-too-small sample size.

We illustrate our methodology and the asymptotic results on both simulated and real world data. Additionally, we illustrate how our pvalue based method can be used as an automatic deterministic early stopping procedure for tree-based boosting. The boosting iterations stop when the tree to be added consists only of a root node.

Keywords: Regression trees, CART, *p*-value, stopping criterion, multiple testing, max statistics

^{*}nils.engler@math.su.se, Department of Mathematics, Stockholm University, Sweden [†]lindholm@math.su.se, Department of Mathematics, Stockholm University, Sweden [‡]lindskog@math.su.se, Department of Mathematics, Stockholm University, Sweden [§]taariq.nazar@math.su.se, Department of Mathematics, Stockholm University, Sweden

1 Introduction

When using binary-split regression trees in practice an important question is how to decide on the complexity of the constructed tree expressed in terms of, e.g., the number of binary splits in the tree, given data. Many applications focus on predictive modeling, where the objective is to construct a tree that generalises well to unseen data. The standard approach to decide on the tree complexity is then to use hold-out data and apply cross-validation techniques, see e.g. [Hastie et al., 2009]. When constructing a tree by sequentially deciding on continuing to split, adding new leaves to the tree in each step, cross-validation corresponds to a method for so-called "early stopping". When using a cross-validation-based early stopping rule, the constructed tree obviously depends on the hold-out-data for the different steps of the procedure. In particular, a randomised selection of hold-out data will inevitably result in the constructed tree being a random function of the data. This is not always desirable. In the present paper a deterministic in-sample early stopping rule is introduced, which is based on *p*-values for whether to accept a binary split or not.

In order to explain the suggested tree-growing method, let T_m denote a greedily grown optimal L^2 CART regression tree (L^2 refers to using a squarederror loss function) with m leaves (suppressing the dependence on covariates), see e.g. [Breiman et al., 1984]. Input to the tree-growing method is a given sequence of nested regression trees T_{m_1}, T_{m_2}, \ldots , where $1 =: m_1 < m_2 < \ldots$, i.e. the first tree is simply a root node, each tree is a subtree of the next tree in the sequence, and no tree appears more than once. Note that T_{m_j} and $T_{m_{j+1}}$ may differ by more than one leaf, i.e. $m_{j+1} - m_j \geq 1$. The tree-growing process starts from the root node T_{m_1} by testing whether increasing the treecomplexity from T_{m_1} to T_{m_2} corresponds to a significant improvement in terms of the L^2 loss. If this is the case, the tree-growing process continues to test whether the tree-complexity should be increased from T_{m_2} and T_{m_3} ; otherwise the tree-growing process stops. If $m_{j+1} - m_j > 1$, all added splits are tested. The tree-growing process is

- (i) based on *p*-values so hypotheses and significance levels need to be specified,
- (ii) an iterative procedure, possibly resulting in a large number of tests.

Concerning (i): The null hypothesis, H_0 , is that there is no signal in data. The alternative hypothesis, H_A , is that there is a sufficiently strong signal making a binary split appropriate. The significance level of the test can be seen as a subjectively chosen hyper-parameter, depending on the modeler's view on the Type I-error. Concerning (ii): We cannot perfectly adjust for multiple testing, but it is possible to use Bonferroni arguments to bound the Type I-error from above. By doing so the tree-growing process is stopped once the *sum* of the *p*-values is greater than the subjectively chosen overall significance level for testing the significance of the entire tree. If $m_{j+1} - m_j >$ 1, then more than one *p*-value is added is added to sum. Since the *p*-value based stopping rule relies on a Bonferroni bound, this tree-growing procedure will be conservative, tending to avoid fitting too large trees to the data.

Relating to the previous paragraph it is important to recall that the treegrowing process is based on a given sequence of nested greedily-grown L^2 CART regression trees, and it is whether these binary splits provide significant loss improvements or not that is being tested. In order to compute a *p*-value for such a split it is crucial to account for that the split was found to be optimal in a step of the greedy recursive partitioning process that generated the tree. This is done by representing the tree-growing process as a certain change-point detection problem, building on results and constructions from [Yao and Davis, 1986]. The usefulness of these results for change-point detection when analysing regression trees was noted in [Shih and Tsai, 2004]. It is important to stress that the *p*-values used are defined with respect to loss improvements and not with respect to potential errors in the estimators for the mean values within a leaf. In the latter problem one needs to adjust for selective inference and this is discussed in a CART-tree context in [Neufeld et al., 2022]. By focusing on the loss improvement and properly taking into account that the tested splits are locally optimal (as described above), selective inference will not be an issue here. Moreover, since the treegrowing process is based on a given sequence of nested CART-trees, we do not address variable selection issues. For more on CART-trees and variable selection, see [Shih and Tsai, 2004].

The *p*-values for loss improvements for a single locally optimally chosen binary split can be calculated exactly for small sample sizes n, but in practice large values for the sample size require approximations. In the current paper an asymptotic approximation is used, which is based on results from [Yao and Davis, 1986] for a single covariate. A contribution of the current paper is to show that for covariate vectors of arbitrary dimension, the accuracy of the *p*-value approximation for a single binary split does not deteriorate substantially if we increase the dimension of the covariate vector. The *p*-value approximation for an entire tree, accounting for multiple testing issues, results in

(a) a conservative stopping rule, given that the null hypothesis H_0 of no signal is true, i.e. the tree-growing process will not be stopped too late,

due to that we are using a Bonferroni upper bound,

(b) that a not-too-weak signal should be detected with a high probability, given a sufficient sample size, i.e. given that the alternative hypothesis H_A is true, the signal will be detected as the sample size tends to infinity.

So far we have focused on deterministic *p*-value-based early stopping when constructing a single greedily grown optimal L^2 CART tree. In practice, however, trees are commonly used as so-called "weak learners" in boosting. The use of *p*-value based early stopping in tree-based L^2 boosting is considered in Section 4. This is similar to the so-called ABT-machine introduced in [Huyghe et al., 2024], which uses another deterministic (not based on e.g. cross-validation) stopping rule based on a sequence of nested trees obtained from so-called cost-complexity pruning, see [Breiman et al., 1984].

Although we focus only on CART trees, one may, of course, consider other types of regression trees and inference based procedures to construct trees. For more on this, see e.g. [Hothorn et al., 2006].

Our main contribution. Given an arbitrary sequence of nested L^2 CART trees, grown by greedy optimal recursive partitioning, we provide an easy-to-use deterministic stopping rule for deciding on the regression tree with suitable complexity. We allow for covariate vectors of arbitrary dimension and the stopping rule is formulated in terms of an easily computable upper bound for the *p*-value corresponding to testing the hypothesis of no signal. Because of the upper bound, the stopping rule is conservative. However, we provide a theoretical guarantee that if there exists signal, then we will detect the existence of this signal if the sample size is sufficiently large. In particular, it is unlikely that we will stop the tree-growing process too early. The asymptotic theoretical guarantee is confirmed by numerical experiments.

Organisation of the paper. The remainder of the paper is structured as follows. Section 2 introduces L^2 CART trees and sequences of nested such trees. Section 2.1 presents and motivates the suggested stopping rule. Section 2.2 describes that the stopping rule naturally leads to considering a change-point-detection problem and presents theoretical results that guarantee statistical soundness of our approach for large sample size. Section 3 compares, for a single split, our approach to well-established regularisation techniques. Section 4 provides a range of numerical illustrations, both in order to clarify the finite-sample performance of our approach and also to illustrate useful applications for tree-based boosting without cross-validation. The proofs of the main results are found in the appendix.

2 Regression trees

The Classification and Regression Tree (CART) method was introduced in the 1980s and uses a greedy approach to build a piecewise constant predictor based on binary splits of the covariate space, one covariate at a time, see e.g. [Breiman et al., 1984]. If we let x be a d-dimensional covariate vector with $x \in \mathbb{X} \subseteq \mathbb{R}^d$, a regression tree with m leaves can be expressed as

$$x \mapsto T_m(x) := \sum_{k=1}^m \zeta_k \mathbb{1}_{\{x \in \mathbb{A}_k\}},\tag{1}$$

where $\zeta_k \in \mathbb{R}$, where $\mathbb{A}_k \subset \mathbb{X}, \bigcup_{k=1}^m \mathbb{A}_k = \mathbb{X}$, and where $\mathbb{1}_{\{x \in \mathbb{A}_k\}}$ is the indicator such that $\mathbb{1}_{\{x \in \mathbb{A}_k\}} = 1$ if $x \in \mathbb{A}_k$, and 0 otherwise. For binary split regression trees, having *m* leaves corresponds to having made m - 1 binary splits.

The construction of a CART tree is based on recursive greedy binary splitting. A split is decided by, for each covariate dimension j, considering the best threshold value ξ for the given covariate dimension, and finally choosing to split based on the best covariate dimension and the associated best threshold value. Splitting the covariate space X based on the jth covariate dimension and threshold value ξ corresponds to the two regions

$$\mathbb{R}_{\text{left}}(j,\xi) = \{ x \in \mathbb{X} : x_j \le \xi \}, \quad \mathbb{R}_{\text{right}}(j,\xi) = \{ x \in \mathbb{X} : x_j > \xi \}.$$

The CART algorithm estimates a regression tree by recursively minimising the empirical risk based on the observed data $(Y^{(1)}, X^{(1)}), \ldots, (Y^{(n)}, X^{(n)})$ that are independent copies of (Y, X), where Y is a real-valued response variable and X is a X-valued covariate vector. When using the L^2 loss and considering a split w.r.t. covariate j, this means that we want to minimise

$$\sum_{i:X^{(i)}\in\mathbb{R}_{left}(j,\xi)} (Y^{(i)} - \overline{Y}_{left}(j,\xi))^2 + \sum_{i:X^{(i)}\in\mathbb{R}_{right}(j,\xi)} (Y^{(i)} - \overline{Y}_{right}(j,\xi))^2, \quad (2)$$

where $\overline{Y}_{\text{left}}(j,\xi)$ is the average of all $Y^{(i)}$ for which $X^{(i)} \in \mathbb{R}_{\text{left}}(j,\xi)$, and similarly for $\overline{Y}_{\text{right}}(j,\xi)$. A regression tree with a single binary split w.r.t. covariate j and threshold value ξ is therefore

$$T_2(x) = \overline{Y}_{\text{left}}(j,\xi) \mathbb{1}_{\{x \in \mathbb{R}_{\text{left}}(j,\xi)\}} + \overline{Y}_{\text{right}}(j,\xi) \mathbb{1}_{\{x \in \mathbb{R}_{\text{right}}(j,\xi)\}}.$$

In order to ease notation, it is convenient to fix a covariate dimension index j and considered the the ordered pairs $(Y^{(1)}, X^{(1)}), \ldots, (Y^{(n)}, X^{(n)})$ of (Y, X), where we assume ordered covariate values $X_j^{(1)} \leq \cdots \leq X_j^{(n)}$ and that the response variables appear in the order corresponding to the size of the covariate values. Hence, $(Y^{(1)}, X^{(1)})$ satisfies $X_j^{(1)} = \min_i X_j^{(i)}$, etc. A different choice of index j would therefore imply a particular permutation of the n response-covariate pairs. By suppressing the dependence on j, this allows us to introduce

$$S_{\leq r} := \sum_{i=1}^{r} (Y^{(i)} - \overline{Y}_{\leq r})^2, \quad S_{>r} := \sum_{i=r+1}^{n} (Y^{(i)} - \overline{Y}_{>r})^2, \quad S := S_{\leq n}, \quad (3)$$

where

$$\overline{Y}_{\leq r} := \frac{1}{r} \sum_{i=1}^{r} Y^{(i)}, \quad \overline{Y}_{>r} := \frac{1}{n-r} \sum_{i=r+1}^{n} Y^{(i)}$$

That is, minimisation of (2) is equivalent to minimising $S_{\leq r} + S_{>r}$ with respect to r, or alternatively we can consider maximising the relative L^2 loss improvement, given by

$$\frac{S - (S_{\le r} + S_{>r})}{S}.$$
 (4)

Further, note that unless we build balanced trees with a pre-specified number of splits we need to add a stopping criterion to the tree-growing process. The perhaps most natural choice is to consider a threshold value, ϑ , say, such that the recursive splitting only continues if the optimal r, denoted r^* , for the optimally chosen covariate dimension $j^* \in \{1, \ldots, d\}$ satisfies

$$\frac{S - (S_{\le r^*} + S_{>r^*})}{S} > \vartheta.$$
(5)

This means that the threshold parameter ϑ functions as a hyper-parameter. In particular, if we let T_m denote a recursively grown L^2 optimal CART-tree with m leaves created using the threshold parameter ϑ , then for any subtree T_m of $T_{m'}$, m < m', the corresponding threshold parameters satisfy $\vartheta > \vartheta'$. Threshold parameters $\vartheta_1 > \vartheta_2 > \ldots > \vartheta_{\tau}$ generate a sequence of nested trees $T_{m_1}, T_{m_2}, \ldots, T_{m_{\tau}}$ with $m_1 \leq m_2 \leq \ldots \leq m_{\tau}$. In applications we will consider sequences $\vartheta_1 > \vartheta_2 > \ldots$ such that $1 = m_1 < m_2 < \ldots$ Note that such a decreasing sequence of threshold parameters will not necessarily result in a sequence of nested trees that only increases by one split at a time.

One procedure to construct a sequence of nested trees is to first pick $\vartheta = 0$ and build a maximal CART-tree, which is pruned from the leaves to the root. One such procedure is the cost-complexity pruning introduced in [Breiman et al., 1984], which likely will lead to a sequence of nested trees where more than one leaf is added in each iteration. For more on this, see Section 3.1.

The threshold parameter ϑ controls the complexity of the tree that is constructed using recursive binary splitting, but it is not clear how to choose ϑ . One option is to base the choice of ϑ on out-of-sample validation techniques, such as cross-validation. The drawback with this is that the tree construction then becomes random: given a fixed dataset repeated application of the procedure may generate different regression trees. We do not want a procedure for constructing regression trees to have this feature. The focus of the current paper is to start from a sequence of nested greedy binary split regression trees, from shallow to deep, and use a particular stopping criterion to decide when to stop the greedy binary splitting in the tree-growing process. The stopping criterion is based entirely on the data used for building the regression trees and is a deterministic mapping from the data to the elements in the sequence of regression trees.

2.1 The stopping rule

Our approach relies on that all binary splits in the sequence of nested regression trees have been chosen in a greedy optimal manner. That is, if we consider an arbitrary binary split in the sequence of nested trees, the reduction in squared error loss is given by the statistic

$$U_{\max} := \max_{1 \le j \le d} U_j, \quad U_j := \max_{1 \le r \le n-1} \frac{S - (S_{\le r} + S_{>r})}{S}, \tag{6}$$

where the sums $S_{\leq r}$ and $S_{>r}$ depend on j because of the implicit ordering of the terms as outlined above, see (3). Given any sample size n and any observed value u_{obs} for the test statistic U_{max} we easily compute, under the null hypothesis of no signal, an upper bound $p_{obs} \geq \mathbb{P}_{\mathcal{N}}(U_{max} > u_{obs})$, where the subscript \mathcal{N} emphasizes the null hypothesis. Therefore, for a regression tree T_m resulting from m - 1 binary splits, it holds that

$$\mathbb{P}_{\mathcal{N}}\left(\bigcup_{k=1}^{m-1}\left\{U_{\max,k} > u_{\mathrm{obs},k}\right\}\right) \leq \sum_{k=1}^{m-1} \mathbb{P}_{\mathcal{N}}(U_{\max,k} > u_{\mathrm{obs},k}) \leq \sum_{k=1}^{m-1} p_{\mathrm{obs},k}.$$

Note that the summation is over all m-1 splits (or internal nodes) of the tree with m leaves. We emphasize that, for every binary split k, $u_{\text{obs},k}$ is observed and $p_{\text{obs},k}$ is easily computed from $u_{\text{obs},k}$. If for a pre-chosen tolerance $\delta \in (0, 1)$ close to zero,

$$\sum_{k=1}^{m-1} p_{\text{obs},k} \le \delta,\tag{7}$$

then we conclude that the event $\cup_{k=1}^{m-1} \{U_{\max,k} > u_{\text{obs},k}\}$ is very unlikely and we reject the null hypothesis of no signal. Consequently, we proceed by considering the next, larger, regression tree $T_{m'}$, m' > m, in the sequence of nested regression trees. If, when considering the regression tree $T_{m'}$ we find that

$$\sum_{k=1}^{m'-1} p_{\text{obs},k} > \delta, \tag{8}$$

then the procedure stops and the previous regression tree T_m is selected as the optimal regression tree.

Since we consider an upper bound for the probability (under the null hypothesis) of the event $\bigcup_{k=1}^{m-1} \{U_{\max,k} > u_{obs,k}\}$ and since we consider upper bounds $p_{obs,k}$ for the probabilities of the events $U_{\max,k} > u_{obs,k}$, we are more likely to stop – observe that (8) holds – compared to a hypothetical situation where the probability of the event $\bigcup_{k=1}^{m-1} \{U_{\max,k} > u_{obs,k}\}$ could be computed and were found to exceed the tolerance level δ . Hence, our stopping criterion is conservative. We therefore have to be concerned with the possibility of a too conservative stopping criterion. However, it is shown in Proposition 1 below that under an alternative hypothesis of a sufficiently strong signal, the computable upper bound p_{obs} for the true p-value is very small. Hence, our stopping criterion is not too conservative.

2.2 Change point detection for a single binary split

The question of whether a candidate binary split should be rejected or not can be phrased as a change-point-detection problem. This observation has been made already in [Shih and Tsai, 2004], where the aim was to target inference based variable selection. The idea here is to make inference on squared-errorloss reduction, where a significant loss reduction translates into not rejecting a split, hence continuing the tree-growing process. This approach builds on the analysis of change-point detection from [Yao and Davis, 1986] that uses a scaled version of (6) according to

$$U_{j}^{(n)} := \max_{1 \le r \le n-1} \frac{S - (S_{\le r} + S_{>r})}{S/n}, \quad U_{\max}^{(n)} := \max_{1 \le j \le d} U_{j}^{(n)},$$

where the dependence of $U_j^{(n)}$ on j is implicit in the order of $Y^{(1)}, \ldots, Y^{(n)}$ which determines the sums of squares $S_{\leq r}$ and $S_{>r}$, as before. That is, the optimal candidate change point w.r.t. covariate dimension j is expressed in terms of the statistic $U_j^{(n)}$, which, hence, is identical to a candidate split point. The test for rejecting a candidate split is based on the null hypothesis saying that observing X gives no information about Y. The null hypothesis corresponds to a simple model \mathcal{N} for (Y, X).

Definition 1 (Null hypothesis, H_0). For model \mathcal{N} , Y and X are independent and Y is normally distributed: there exist $\mu \in \mathbb{R}$ and $\sigma^2 \in (0, \infty)$ such that

$$\mathbb{P}_{\mathcal{N}}(Y \in \cdot \mid X) = \mathbb{P}_{\mathcal{N}}(Y \in \cdot) = N(\mu, \sigma^2).$$

When considering a nested sequence of binary regression trees, $U_{\max}^{(n)}$ is the random variable whose outcome is the observed test statistic for a single candidate binary split. Under the null hypothesis, the common distribution of the statistics $U_1^{(n)}, \ldots, U_d^{(n)}$ does not depend on μ and σ . Hence, under the null hypothesis, the distribution of $U_{\max}^{(n)}$ does not depend on μ and σ . Clearly,

$$\mathbb{P}_{\mathcal{N}}\left(U_{\max}^{(n)} > u\right) = \mathbb{P}_{\mathcal{N}}\left(\cup_{j=1}^{d} \left\{U_{j}^{(n)} > u\right\}\right) \le d\mathbb{P}_{\mathcal{N}}\left(U_{j}^{(n)} > u\right)$$
(9)

which does not depend on j since the probability is evaluated under the null hypothesis. We approximate the tail probability $\mathbb{P}_{\mathcal{N}}(U_j^{(n)} > u)$ by $p_n(u)$, where

$$p_n(u) := 1 - \Phi\left(u^{1/2} - \frac{\ln_3(n) + \ln(2)}{(2\ln_2(n))^{1/2}}\right)^{2\ln(n/2)},\tag{10}$$

where $\ln_k(n)$ corresponds to the k times iterated logarithm, e.g. $\ln_2(n) = \ln(\ln(n))$. The approximation $p_n(u)$ from (10) corresponds to Eq. (2.5) on p. 345 in [Yao and Davis, 1986]. The true p-value is the function $u \mapsto \mathbb{P}_{\mathcal{N}}(U_{\max}^{(n)} > u)$ evaluated at the observed value for $U_{\max}^{(n)}$. The true p-value is approximated from above by

$$P_{\max}^{(n)} := dp_n(U_{\max}^{(n)}).$$
(11)

We emphasise that given an observation $u_{\text{obs},k}$ of $U_{\max}^{(n)}$, $p_{\text{obs},k}$ is the observed outcome of $P_{\max}^{(n)}$.

If the true signal is not too weak, which means that the conditional expectation of Y given X should fluctuate sufficiently in size, then for any significance level we want to reject the null hypothesis in a setting with sufficiently large sample size n. In order to make the meaning of this statement precise, and in order to verify it, we must consider the alternative hypothesis as a sequence of hypotheses indexed by the sample size n. The alternative hypothesis corresponds to a sequence of models $\mathcal{A} = (\mathcal{A}^{(n)})$.

Definition 2 (Alternative hypothesis, H_A). For the sequence of models $(\mathcal{A}^{(n)})$ there exist $j \in \{1, \ldots, d\}, \xi \in \mathbb{R}, t_0 \in (0, 1), \sigma^2 \in (0, \infty)$ and $\mu_l, \mu_r \in \mathbb{R}, \mu_l \neq \mu_r$, such that for all n

$$\mathbb{P}_{\mathcal{A}^{(n)}}(X_j \le \xi) = t_0, \\ \mathbb{P}_{\mathcal{A}^{(n)}}(Y \in \cdot \mid X_j = x) = N(\mu_l \mathbb{1}_{\{x \le \xi\}} + \mu_r \mathbb{1}_{\{x \ge \xi\}}, \sigma^2),$$

where $|\mu_r - \mu_l| = \sigma \theta_n > 0$. The sequence θ_n satisfies

$$\theta_n = \frac{(2\ln_2(n))^{1/2} + \eta_n}{n^{1/2}(t_0(1-t_0))^{1/2}}$$
(12)

for some increasing sequence η_n with $\lim_{n\to\infty} \eta_n = \infty$ and $\limsup_{n\to\infty} \theta_n < \infty$.

The requirement under the alternative hypothesis of a shift in mean of size $\sigma \theta_n$ says that the amplitude of the signal is allowed to decrease towards zero with n, but not too fast. We could consider $\theta_n = n^{-r}$ for some r < 1/2. We may also consider a constant signal amplitude θ . However, that situation is not very interesting since such a signal should eventually be easily detectable as the sample size n becomes very large. The expression for θ_n in (12) comes from [Yao and Davis, 1986] (Eq. (3.2) on p. 347) and corresponds to an at least slightly stronger signal compared to what was considered in [Yao and Davis, 1986] ($\eta_n \to \infty$ instead of $\eta_n = \eta + o(1)$).

We want to show that under the alternative hypothesis we will reject the null hypothesis with a probability tending to one. The null hypothesis is not rejected at significance level $\varepsilon > 0$ if $P_{\text{max}}^{(n)} > \varepsilon$. We want to show that under the alternative hypothesis, the probability of falsely not rejecting the null hypothesis is very small. More precisely, we show the following:

Proposition 1. $\lim_{n\to\infty} \mathbb{P}_{\mathcal{A}^{(n)}}(P_{\max}^{(n)} > \varepsilon) = 0$ for every $\varepsilon > 0$.

The proof of Proposition 1 is given in the Appendix. To conclude, using the p-value approximation (11) results in

- (i) a conservative stopping rule, given that the null hypothesis H_0 of no signal is true, i.e. the tree-growing process will not be stopped too early, due to that we are using a Bonferroni upper bound,
- (*ii*) that a not too weak signal should be detected with a high probability, given a sufficient sample size, i.e. given that the alternative hypothesis H_A is true, the signal will be detected as the sample size tends to infinity.

3 Relation to classical regularisation techniques

The focus of this section is on a single binary split. Let $T_2(X)$ denote an optimal binary split CART tree with a single split, and let $T_1(X)$ denote the root tree of $T_2(X)$. Based on the notation in Section 2 the split is accepted at significance level ε if

$$U_{\max}^{(n)} = \frac{S - (S_{\le r^*} + S_{>r^*})}{S/n} > u_{\varepsilon},$$
(13)

where "*" indicates that we consider the optimal split, and where u_{ε} is the solution to

$$dp_n(u_\varepsilon) = \varepsilon, \tag{14}$$

where $p_n(u)$ is from (10). An equivalent rephrasing of (13) is

$$MSE_1 - MSE_2 - u_{\varepsilon}\widehat{\sigma}^2 > 0, \qquad (15)$$

where $MSE_1 := S, MSE_2 := S_{\leq r^*} + S_{>r^*}$, together with $\hat{\sigma}^2 := S/n$. A natural question, which is partially answered below, is how u_{ε} depends on n for a fixed significance level ε .

Proposition 2. u_{ε} solving (14) satisfies $u_{\varepsilon} = o(\ln_2(n))$ as $n \to \infty$.

The proof of Proposition 2 is given in the Appendix.

Based on (15) it is seen that u_{ε} can be thought of as a regularisation term (or penalty), and from Proposition 2 it is seen that this term behaves almost like a constant. We will continue with a short comparison with other techniques that can be used to decide on accepting a split or not.

3.1 Cost-complexity pruning

cost-complexity pruning was introduced in [Breiman et al., 1984] and is described in terms of the so-called "cost" w.r.t. a split tolerance ϑ , denoted by $R_{\vartheta}(T)$, defined as

$$R_{\vartheta}(T) := R(T) + \vartheta |T|, \tag{16}$$

where, in our sitting, we have $R(T) = \sum_{i=1}^{n} (Y^{(i)} - T(X^{(i)}))^2$ (other loss functions may be considered). The parameter ϑ is also referred to as the "cost-complexity" parameter. Note that the critical ϑ value needed in order

to accept $T_2(X)$ in favour of $T_1(X)$ is the threshold value ϑ for which the so-called "gain" $R_{\vartheta}(T_1) - R_{\vartheta}(T_2)$ is 0, which gives

$$R_{\vartheta}(T_1) - R_{\vartheta}(T_2) = R(T_1) - R(T_2) + \vartheta(|T_1| - |T_2|) = MSE_1 - MSE_2 - \vartheta = 0,$$

or equivalently, the split is accepted if $MSE_1 - MSE_2 - \vartheta > 0$. The choice of ϑ used in applications is typically based on out-of-sample performance using, e.g., cross-validation; also recall the discussion in relation to (5) above. Using the specific choice $\vartheta := u_{\varepsilon} \hat{\sigma}^2$ is equivalent to using the *p*-value based penalty from (15). Note that this equivalence only applies to the situation concerning whether one should accept a single split or not, whereas, as mentioned above, the cost-complexity pruning is a procedure that evaluates entire subtrees.

3.2 Covariance penalty and information criteria

Another alternative is to assess a candidate split based on its predictive performance using the mean squared error of prediction (MSEP), conditioning on the observed covariate values. When working with linear Gaussian models this corresponds to using Mallows' C_p , where p corresponds to the number of regression parameters, see e.g. [Mallows, 1973], which is an example of an estimate of the prediction error using covariance based penalty, see e.g. [Efron, 2004]. The C_p statistic can then be expressed as

$$C_p := \frac{1}{n} \big(\operatorname{MSE}_p + 2p\widehat{\sigma}^2 \big),$$

which is the formulation used in [Hastie et al., 2009, Ch. 7.5, Eq. (7.26)]. Consequently, since a binary single-split L^2 regression tree with predetermined split point can be interpreted as fitting a Gaussian model with a single binary covariate, C_p can in this situation be used to evaluate predictive performance. By considering the C_p improvement when going from no split, i.e. p = 1, to one split, p = 2, corresponds to $C_1 - C_2 > 0$, which is equivalent to

$$MSE_1 - MSE_2 - 2\hat{\sigma}^2 > 0.$$

Thus, using Mallows' C_p , targeting the predictive performance of the estimator, will be asymptotically too liberal compared to the *p*-value based stopping rule. This, however, should not be too surprising, since the above application of the C_p statistic does *not* take into account that the candidate split point has been chosen by minimising an L^2 loss.

For a *p*-parameter Gaussian model the C_p statistic coincides with the Akaike information criterion (AIC), see e.g. [Hastie et al., 2009, Ch. 7.5,

Eq. (7.29)]. For a *p*-parameter Gaussian model, the Bayesian information criterion (BIC) considers the quantity

$$\operatorname{BIC}_p := \frac{n}{\sigma^2} (\operatorname{MSE}_p + \ln(n)p\sigma^2),$$

see e.g. [Hastie et al., 2009, Ch. 7.7, Eq. (7.36)], as the basis for model selection. In practice σ^2 is replaced by a suitable estimator, $\hat{\sigma}^2$, see, e.g., the discussion in the paragraph following [Hastie et al., 2009, Ch. 7.7, Eq. (7.36)]. Hence, it follows that accepting a split based on BIC-improvement in a single split corresponds to

$$\operatorname{BIC}_1 - \operatorname{BIC}_2 > 0,$$

which is equivalent to

$$MSE_1 - MSE_2 - \ln(n)\widehat{\sigma}^2$$
.

Thus, using BIC as a stopping criterion is more conservative than the p-value based stopping criterion, despite not taking into account that the split point is given as a result of an optimisation procedure.

4 Numerical illustrations

4.1 The *p*-value approximation for a single split

In this section we investigate the error from applying the two approximations in (9) and (10). Both together provide the *p*-value approximation used to test for signal. Since we do not have access to the true distribution of U_{max} under H_0 , we compute its empirical distribution from 10,000 realisations in order to compare to the approximations.

Figure 1 shows the approximated and true cdfs for varying sample size and covariate dependence. Here, the covariate dimension is set to d = 10. Table 1 compares the approximated and true critical quantile values at a 0.95-level for varying sample size, covariate dimension and covariate dependence. Note that for d = 1, varying dependence is not an issue so that the entries of the first two tables are identical.

As was noted in [Yao and Davis, 1986] [Remark 2.3], the approximation (10) yields satisfactory results even for small sample sizes $20 \le n \le 50$. This is confirmed by the first row of Table 1. The second row of Figure 1 as well as the middle part of Table 1 show that a strong positive pairwise correlation of $\rho = 0.8$ between covariates does not substantially affect the upper tail of the



Figure 1: Blue curves: empirical cdf of U_{\max} given H_0 computed from 10,000 realisations. Orange curves: Approximation $1-dp_n(u)$. Left column: n = 50, right column: n = 1000. Top row: independent standard normal covariates, bottom row: dependent normal covariates with common pairwise correlation $\rho = 0.8$ and unit variance. The points of intersection with the dashed blue line illustrate empirical and approximate 0.95-quantile of U_{\max} .

	n = 50	n = 1000	n = 50	n = 1000	n = 50	n = 1000
d = 1	8.55	10.78	8.55	10.78	9.12	11.09
d = 2	9.79	12.10	9.62	12.00	10.67	12.68
d = 10	12.46	15.51	11.94	14.84	14.23	16.31

Table 1: Left table: 0.95-level quantiles based on the empirical cdf of U_{max} given H_0 computed from 10,000 realisations for independent standard normal covariates. Middle table: The analogous quantiles for dependent normal covariates with a common pairwise correlation of $\rho = 0.8$ and standard variances. Right table: Quantile approximation corresponding to (10).

distribution of U_{max} under H_0 and that the quantile approximations provide good upper bounds.

We now turn to assuming that the alternative hypothesis H_A according to Definition 2 holds. In order to illustrate Proposition 1, we pick $\varepsilon = 0.05$, $\sigma^2 = 1$, j = 1, $\xi = 0$, $t_0 = 1/2$, $\mu_l = 0$ and $\mu_r = n^{-1/5}$. Note that the step size is chosen to decrease slowly enough towards zero in order to fulfil the assumptions of H_A in Definition 2.

In Figure 2, we plot the fraction of correct signal detections from 1000 realisations of the event $\{U_{\max}^{(n)} > u_{\varepsilon}\}$, where u_{ε} is given in (14). We run the simulations for an increasing number of data points n. Figure 2 confirms the findings of Proposition 1 that the probability of detecting a slowly decreasing signal converges to one as n tends to infinity.

It can be noted that the upper tail of $U_{\max}^{(n)}$ is not affected much by introducing dependence between the covariates, as the orange and blue curves in the right plot of Figure 2 differ little.

4.2 Simulated examples from Neufeldt et al.

In this section we fix a simple tree and then generate residuals around its level values in order to illustrate the detection performance of our method. We consider the following example as proposed by [Neufeld et al., 2022, section 5]. Consider independent standard normal covariates and a regression function given by

$$\mu(x) = b \Big(\mathbb{1}_{\{x_1 \le 0\}} \Big(1 + a \mathbb{1}_{\{x_2 > 0\}} + \mathbb{1}_{\{x_2 x_3 > 0\}} \Big) \Big), \tag{17}$$

for $x \in \mathbb{R}^{10}$ and parameters $a, b \in \mathbb{R}$ determining the step size between the level values (signal strength). The step size between siblings at level two is



Figure 2: Blue curves: Fraction of correct signal detections according to $\{U_{\max}^{(n)} > u_{\varepsilon}\}$ for an increasing number of data points n and independent standard normal covariates. Orange curve: The analogous fraction based on dependent multivariate normal covariates with common pairwise correlation $\rho = 0.8$ and unit variance. Green curve: The signal strength $|\mu_r - \mu_l| = n^{-1/5}$. The blue dashed line shows the 0.95-level. The left and right plots correspond to d = 1 and d = 10 covariates, respectively.

ab while the step size between siblings at level three is *b*. An illustration of the tree corresponding to (17) is given in Figure 3. We generate 500 iid covariate vectors X_1, \ldots, X_{500} of $N(0, I_{10})$ and corresponding response variables Y_1, \ldots, Y_{500} , where, given X_i, Y_i is drawn from $N(\mu(X_i), 1)$.



Figure 3: Regression tree corresponding to (17) with a = 1 adopted from [Neufeld et al., 2022, section 5]. Each left child answers the inequality with "true".

Using the python package sklearn.tree.DecisionTreeRegressor, we grow a full CART tree of maximal depth 4 with a minimal number of data points per

leaf set to 20. For each tree in the nested sequence of cost-complexity-pruned subtrees (from the root to the fully grown CART tree), we compute the insample error (MSE) and out-of-sample error (MSEP), where the latter is done using independently generated test data of the same size n = 500 which was neither used to fit the CART tree, nor to compute *p*-values, but serves only as a data set for pure out-of-sample testing.



Figure 4: Left plot: MSEP (blue) and MSE (orange) for each tree in the nested sequence of cost-complexity-pruned subtrees. Right plot: cumulative *p*-value for each tree in the nested sequence of cost-complexity-pruned subtrees. The *x*-axis depicts the number of leaves of the subtree considered. The dashed blue line marks our method's output tree, i.e. the largest subtree whose cumulative *p*-value lies below $\delta = 0.05$. The signal strength parameters are a = b = 1.

In the example of Figure 4, the proposed method detects the correct complexity of μ which is given by 5 leaves and which minimises MSEP. The cumulative *p*-values of all smaller subtrees are very close to zero (0, 0.0001, 0.0003), while jumping sharply to 1.08 after the first "unnecessary" split (cf. Figure 6). The results in this example are hence not sensitive to the choice of the tolerance parameter δ . Note that individual *p*-values may exceed one due to the approximation (11). Comparing Figure 3 with the upper tree of Figure 6, we note that also the split points and mean values are accurate.

We repeat the simulation for a decreased signal parameter b = 0.5, while keeping a = 1, $\sigma^2 = 1$ and n = 500. As can be observed in Figure 5 and the bottom tree of Figure 6, the method stops after already one split not capable of detecting the weak signal in the lower part of the tree. However, it regularises well in the sense that MSEPs are close to minimal. Even though the sample size n = 500 is chosen rather small, the results of Figures 4 and 5 do not vary much between runs with different random seeds for the training and validation data generation.



Figure 5: Analogue of Figure 4 with b = 0.5 instead of b = 1.

Moreover, from Figure 2 in Section 4.2 we can observe that a larger number of data points of around n = 2500 would ensure (with a 95 percent probability) the detection of an even lower signal 0.21 < b = 0.5 in each split of the tree. We conclude that n = 500 is insufficient in this example with b = 0.5.

4.2.1 Illustrating the randomness of tree construction using crossvalidation

Above we mention the drawback of training trees using cross-validation which is that the resulting tree depends of the randomness inherent in the crossvalidation procedure. In this section we illustrate this fact for CART-trees. We generate data according to the model from [Neufeld et al., 2022], as presented in Section 4.2, with parameters a = 1, b = 1 and $\sigma^2 = 1$. Here we consider sample size n = 1000 (rather than n = 500 considered in Section 4.2). We split the data into a 80% training set and a 20% test set. The CARTtree is trained using 5-fold cross-validation on the training set, which entails optimally choosing a cost-complexity parameter ϑ . An optimal CART-tree is trained on the complete training set using the cost-complexity parameter ϑ . Finally, the trained model is evaluated on the test set. This procedure is repeated 500 times, allowing us to estimate RMSE values empirically. It turns out that throughout the 500 iterations of the procedure, only two distinct trees are selected by the cross-validation procedure: either a tree with two leaves or a tree made up of only the root node. Since cross-validation results in non-deterministic ϑ , we realise two distinct ϑ values corresponding to two distinct trees in the sequence cost-complexity pruned trees.



Figure 6: Regularised output trees for b = 1 (top) and b = 0.5 (bottom). First row of each node: split point selected by CART. Second row: mean value. Third row: node *p*-value. Fourth row: cumulative *p*-value of the smallest subtree the node appears in as a non-leaf. Nodes shaded red violate the condition that the cumulative *p*-value lies below 0.05.

We evaluate our model on the same dataset with identical CART-tree parameters and significance levels $\delta = 0.1, 0.05, 0.01$ and find that our method attains an even lower RMSE for all three choices of δ . The results can be seen in 2. Further, we can see the shape of the estimated trees in Figure 7.

		-
cost-complexity parameter ϑ	RMSE	number of leaves
0.000	1.045	7
0.008	1.012	5
0.072	1.046	4
0.075	1.062	3
0.113	1.144	2
1.016	1.537	1

Table 2: Evaluation of trees in the sequence of cost-complexity pruned trees.

4.3 An application to L²-boosting

In this section, we illustrate how our proposed method performs when it is used as a weak learner in a standard L^2 -boosting setting applied to the datasets California Housing and beMTPL16 from [Dutang and Charpentier, 2024].

Throughout these illustrations we compare the L^2 -boosting version of our method to the Gradient Boosting Machine (GBM) with identical configurations. For both methods, we split the data into a 80% training set and a 20% test set. We train the models on the same training set and evaluate them on the same test set. We fix the max depth of the weak learners to 3, i.e. a tree with at most 8 leaves can be added in a single iteration, the minimum samples per leaf is set to 20 and we set the learning rate for the boosting procedures to 0.1.

In each boosting iteration, we use the residuals from the previous iteration as the working response. In each boosting iteration, we determine a nested sequence of trees (as described above) and the weak learner is selected as the maximally split tree that satisfies the criterion $\sum p_j < \delta$ for the chosen significance level δ . We stop the boosting procedure when the candidate weak learner is the root node, i.e., no statistically significant split can be made. Note that the complexity of the weak learner for our method is dynamic, determined by the criterion $\sum p_j < \delta$.

The California Housing dataset consists of n = 20640 data points, and the number of covariates is d = 8. The beMTPL16 dataset consists of n = 70791 data points, and the number of covariates is d = 6. In Figure 8 we see how our method compares to the GBM when applied to the two datasets for varying



Figure 7: Regression trees corresponding to different cost-complexity parameters related to the test of cross-validation randomness; Panels (a) and (b) show the trees obtained using CV, Panel (c) shows the tree obtained using the *p*-value method. Leaf values correspond to mean values. All RMSE values can be found in Table 2.

levels of δ . The value $\delta = \infty$ gives a boosted-trees procedure similar to the ABT-machine from [Huyghe et al., 2024] in the case of L^2 -boosting. It should be noted that the GBM stopping criterion implies, for the California housing dataset, that it is trained for approximately 2500 iterations before stopping. One could consider tuning the shrinkage parameter in order to adjust the number of boosting steps, but this has not been investigated further in the present paper. It can be seen from Figure 8 that the number of iterations for the *p*-value based method is not necessarily monotone in δ . However, this is not contradictory since different values of δ will result in that the trees added in each iteration may have a rather different tree complexity. We find that the *p*-value based stopping criterion for the weak learner in L^2 -boosting generates promising results and should be investigated further, including comparisons with, e.g., the ABT-machine from [Huyghe et al., 2024].



Figure 8: RMSE on test data (y-axis) as a function of the number of boosting iterations (x-axis). The left plot corresponds to the California Housing dataset, the right plot to the beMTPL16 dataset. The blue curve corresponds to our method using $\delta = \infty$, the orange curve $\delta = 0.10$, the green curve $\delta = 0.05$, the red curve $\delta = 0.01$ and the purple curve corresponds to the GBM. The vertical dashed lines corresponds to where the iterations stop and the horizontal dashed lines corresponds to the lowest RMSE achieved for the respective methods.

References

- [Breiman et al., 1984] Breiman, L., Friedman, J. H., Olshen, R. A., and Stone, C. J. (1984). *Classification and regression trees*. Wadsworth, Belmont, Calif.
- [Dutang and Charpentier, 2024] Dutang, C. and Charpentier, A. (2024). CASdatasets: Insurance datasets. R package version 1.2-0, DOI 10.57745/P0KHAG.
- [Efron, 2004] Efron, B. (2004). The estimation of prediction error: covariance penalties and cross-validation. *Journal of the American Statistical Association*, 99(467):619–632.
- [Hastie et al., 2009] Hastie, T., Tibshirani, R., and Friedman, J. (2009). The Elements of Statistical Learning [Elektronisk resurs] Data Mining, Inference, and Prediction. Springer New York, New York, NY, second. edition.
- [Hothorn et al., 2006] Hothorn, T., Hornik, K., and Zeileis, A. (2006). Unbiased recursive partitioning: A conditional inference framework. *Journal* of Computational and Graphical statistics, 15(3):651–674.

- [Huyghe et al., 2024] Huyghe, J., Trufin, J., and Denuit, M. (2024). Boosting cost-complexity pruned trees on tweedie responses: the abt machine for insurance ratemaking. *Scandinavian Actuarial Journal*, 2024(5):417–439.
- [Mallows, 1973] Mallows, C. (1973). Some comments on c_p . Technometrics, 15(4):661-675.
- [Neufeld et al., 2022] Neufeld, A. C., Gao, L. L., and Witten, D. M. (2022). Tree-values: selective inference for regression trees. *The Journal of Machine Learning Research*, 23(1):13759–13801.
- [Shih and Tsai, 2004] Shih, Y.-S. and Tsai, H.-W. (2004). Variable selection bias in regression trees with constant fits. *Computational statistics & data analysis*, 45(3):595–607.
- [Yao and Davis, 1986] Yao, Y.-C. and Davis, R. A. (1986). The asymptotic behavior of the likelihood ratio statistic for testing a shift in mean in a sequence of independent normal variates. *Sankhyā: The Indian Journal of Statistics, Series A*, pages 339–353.

A Proofs

A.1 Proof of Proposition 1

Before starting the proof of Proposition 1 we note the following:

Remark 3. The distribution of the observed test statistic $U_{\max}^{(n)}$ does not depend on σ under the alternative hypothesis. Under the alternative hypothesis, for any $r \in \{1, \ldots, n\}$ and $b \in \{1, \ldots, r\}$ such that $X_j^{(i)} < \xi$ for $i \leq b$ and $X_j^{(i)} \geq \xi$ for i > b, we may write

$$Y^{(i)} = Z^{(i)} + \begin{cases} \mu_l, & i = 1, \dots, b, \\ \mu_r, & i = b + 1, \dots, n. \end{cases}$$

where $Z^{(1)}, \ldots, Z^{(n)}$ are independent and $N(0, \sigma^2)$ distributed. Then $\overline{Y}_{\leq r} = \overline{Z}_{\leq r} + (b\mu_l + (r-b)\mu_r)/r$ and

$$Y^{(i)} - \overline{Y}_{\leq r} = Z^{(i)} - \overline{Z}_{\leq r} + \begin{cases} (\mu_l - \mu_r)(r-b)/r, & i = 1, \dots, b, \\ (\mu_r - \mu_l)b/r, & i = b+1, \dots, r. \end{cases}$$

Hence, $Y^{(i)} - \overline{Y}_{\leq r}$ equals σ times a random variable whose distribution does not depend on σ . This also holds for $Y^{(i)} - \overline{Y}_{>r}$. We conclude that the distribution of $U_j^{(n)}$ does not depend on σ under the alternative hypothesis. Remark 4. By construction

$$U_{j}^{(n)}\frac{S/n}{\sigma^{2}} = \max_{1 \le r \le n-1} \frac{1}{\sigma^{2}} \left(S - S_{\le r} - S_{>r} \right).$$
(18)

Under the alternative hypothesis, by [Yao and Davis, 1986] p. 347,

$$\max_{1 \le r \le n-1} \frac{1}{\sigma^2} \left(S - S_{\le r} - S_{>r} \right) \stackrel{d}{=} \max_{1 \le nt \le n-1} \frac{(W_0(t) - f_n(t))^2}{t(1-t)},$$

where W_0 is a standard Brownian bridge and

$$f_n(t) = \begin{cases} n^{1/2} \theta_n t (1 - [nt_0]/n), & \text{if } nt \le [nt_0], \\ n^{1/2} \theta_n (1 - t) [nt_0]/n, & \text{if } nt > [nt_0]. \end{cases}$$

Proof of Proposition 1. Since p_n is a decreasing function we know that $P_{\max}^{(n)} \leq dp_n(U_j^{(n)})$ for every j, in particular for j for which there is signal with amplitude $\sigma \theta_n$ according to the model $\mathcal{A}^{(n)}$. Hence,

$$\mathbb{P}_{\mathcal{A}^{(n)}}(P_{\max}^{(n)} > \varepsilon) \le \mathbb{P}_{\mathcal{A}^{(n)}}(p_n(U_j^{(n)}) > \varepsilon/d)$$

= 1 - $\mathbb{P}_{\mathcal{A}^{(n)}}(U_j^{(n)} > p_n^{-1}(\varepsilon/d))$

Let T_n^2 denote the quantity $U_j^{(n)}(S/n)/\sigma^2$ in (18). Then

$$\mathbb{P}_{\mathcal{A}^{(n)}}(U_j^{(n)} > p_n^{-1}(\varepsilon/d)) = \mathbb{P}_{\mathcal{A}^{(n)}}\left(T_n^2\left(\frac{c_n^2}{p_n^{-1}(\varepsilon/d)}\frac{\sigma^2}{S/n}\right) > c_n^2\right)$$

for any positive sequence (c_n^2) . We consider the choice of sequence

$$c_n^2 = \left(\frac{2^{-1}\ln_3(n) - \ln(2^{-1}\pi^{1/2}\ln((1-\alpha)^{-1}))}{(2\ln_2(n))^{1/2}} + (2\ln_2(n))^{1/2}\right)^2$$
(19)

in order to relate the tail probability $\mathbb{P}_{\mathcal{A}^{(n)}}(U_j^{(n)} > p_n^{-1}(\varepsilon/d))$ to the tail probability $\mathbb{P}_{\mathcal{A}^{(n)}}(T_n^2 > c_n^2)$ studied by [Yao and Davis, 1986]. By Lemma 5,

$$\liminf_{n \to \infty} \mathbb{P}_{\mathcal{A}^{(n)}}(U_j^{(n)} > p_n^{-1}(\varepsilon/d)) \ge \liminf_{n \to \infty} \mathbb{P}_{\mathcal{A}^{(n)}}(T_n^2 > c_n^2).$$

For any $\eta \in \mathbb{R}$, by Lemma 6,

$$\liminf_{n \to \infty} \mathbb{P}_{\mathcal{A}^{(n)}}(T_n^2 > c_n^2) \ge \alpha + \Phi(\eta)(1 - \alpha).$$

Hence, for any $\eta \in \mathbb{R}$,

$$\limsup_{n \to \infty} \mathbb{P}_{\mathcal{A}^{(n)}}(P_{\max}^{(n)} > \varepsilon) \leq \limsup_{n \to \infty} \left(1 - \mathbb{P}_{\mathcal{A}^{(n)}}(T_n^2 > c_n^2) \right)$$
$$\leq 1 - \alpha - \Phi(\eta)(1 - \alpha).$$

Since we may choose η arbitrarily large, the proof is complete.

Lemma 5. $\liminf_{n\to\infty} \mathbb{P}_{\mathcal{A}^{(n)}}(U_j^{(n)} > p_n^{-1}(\varepsilon/d)) \ge \liminf_{n\to\infty} \mathbb{P}_{\mathcal{A}^{(n)}}(T_n^2 > c_n^2)$ *Proof.* Let

$$F_n := \frac{c_n^2}{p_n^{-1}(\varepsilon/d)} \frac{\sigma^2}{S/n}$$

and note that

$$\mathbb{P}_{\mathcal{A}^{(n)}}(U_{j}^{(n)} > p_{n}^{-1}(\varepsilon/d)) = \mathbb{P}_{\mathcal{A}^{(n)}}(T_{n}^{2}F_{n} > c_{n}^{2})$$

$$\geq \mathbb{P}_{\mathcal{A}^{(n)}}(T_{n}^{2}F_{n} > c_{n}^{2} \mid F_{n} < 1)\mathbb{P}_{\mathcal{A}^{(n)}}(F_{n} < 1)$$

$$+ \mathbb{P}_{\mathcal{A}^{(n)}}(T_{n}^{2} > c_{n}^{2}).$$

We will show that $\lim_{n\to\infty} \mathbb{P}_{\mathcal{A}^{(n)}}(F_n < 1) = 0$ from which the conclusion follows. By Lemma 7,

$$\lim_{n \to \infty} p_n^{-1}(\varepsilon/d)/c_n^2 = 0.$$
(20)

Under $\mathcal{A}^{(n)}$ there exist independent $Z^{(i)} \sim N(0, \sigma^2)$ and $r \in \{1, \ldots, n\}$ such that $Y^{(i)} = Z^{(i)}$ for $i \leq r$, and $Y^{(i)} = Z^{(i)} + \sigma \theta_n$ for i > r. Therefore,

$$S = \sum_{i=1}^{r} (Y^{(i)} - \overline{Y}_{\leq n})^2 + \sum_{i=r+1}^{n} (Y^{(i)} - \overline{Y}_{\leq n})^2$$

=
$$\sum_{i=1}^{n} (Z^{(i)} - \overline{Z}_{\leq n})^2 + \sigma^2 \theta_n^2 \left(r \left(\frac{n-r}{n} \right)^2 + (n-r) \left(\frac{r}{n} \right)^2 \right)$$

+
$$2 \sum_{i=1}^{r} (Z^{(i)} - \overline{Z}_{\leq n}) \sigma \theta_n \frac{n-r}{n} + 2 \sum_{i=r+1}^{n} (Z^{(i)} - \overline{Z}_{\leq n}) \sigma \theta_n \frac{r}{n}$$

Therefore, by Hölder's inequality applied to the sum of the last two terms above,

$$\frac{S}{n} \leq \frac{1}{n} \sum_{i=1}^{n} (Z^{(i)} - \overline{Z}_{\leq n})^2 + \sigma^2 \theta_n^2 + 2\left(\frac{1}{n} \sum_{i=1}^{n} (Z^{(i)} - \overline{Z}_{\leq n})^2\right)^{1/2} \sigma \theta_n$$
$$= \left(\left(\frac{1}{n} \sum_{i=1}^{n} (Z^{(i)} - \overline{Z}_{\leq n})^2\right)^{1/2} + \sigma \theta_n\right)^2.$$

Since the first term inside the square converges in probability to σ and since the second term is bounded we conclude that $\lim_{n\to\infty} \mathbb{P}_{\mathcal{A}^{(n)}}(F_n < 1) = 0$. The proof is complete.

Lemma 6. For every $\eta \in \mathbb{R}$, $\liminf_{n \to \infty} \mathbb{P}_{\mathcal{A}^{(n)}}(T_n^2 > c_n^2) \ge \alpha + \Phi(\eta)(1-\alpha)$.

Proof. Fix $\eta \in \mathbb{R}$. From the expression for the tail probability on page 350 in [Yao and Davis, 1986] we see that for each n,

$$\mathbb{P}_{\mathcal{A}^{(n)}}(T_n^2 > c_n^2) \ge \mathbb{P}(B_{n,1} \cap B_{n,2} \cup A_n(\theta_n)).$$

The events $B_{n,1}, B_{n,2}$ are independent of θ_n and given by

$$B_{n,1} = \left\{ \max_{t \in D_{n,1}} \frac{|W(t)|}{t^{1/2}} > c_n \right\}, \quad B_{n,2} = \left\{ \max_{t \in D_{n,2}} \frac{|W(t) - W(1)|}{(1-t)^{1/2}} > c_n \right\},$$

where W is standard Brownian motion and $D_{n,1}, D_{n,2}$ are index sets. The event $A_n(\theta_n)$ is increasing in θ_n and given by the expression on p. 350 in [Yao and Davis, 1986] (there with θ instead of θ_n). Writing $\theta_n = \theta(n, \eta_n)$ for θ_n in (12), note that $\theta(n, \eta_n) \ge \theta(n, \eta)$ for n sufficiently large since $\eta_n \to \infty$ as $n \to \infty$. Hence, for n sufficiently large,

$$\mathbb{P}_{\mathcal{A}^{(n)}}(T_n^2 > c_n^2) \ge \mathbb{P}(B_{n,1} \cap B_{n,2} \cup A_n(\theta(n,\eta)))$$

and the right-hand side converges to $\alpha + \Phi(\eta)(1-\alpha)$ as concluded on p. 350 in [Yao and Davis, 1986]. The proof is complete.

Lemma 7. $\lim_{n\to\infty} p_n^{-1}(\varepsilon/d)/c_n^2 = 0$

Proof. We have, from the definition of p_n ,

$$p_n^{-1}(\varepsilon/d) = \left(\frac{\ln_3(n) + \ln(2)}{(2\ln_2(n))^{1/2}} + \Phi^{-1}\left((1 - \varepsilon/d)^{1/(2\ln(n/2))}\right)\right)^2, \quad (21)$$

where the first term vanishes asymptotically and the second term tends to ∞ as $n \to \infty$. Similarly, in (19) the first term vanishes asymptotically and the second term tends to ∞ as $n \to \infty$. Hence, it is sufficient to compare the two terms that are not vanishing asymptotically and show that

$$\lim_{n \to \infty} \frac{\Phi^{-1}(x_n)}{\Phi^{-1}(y_n)} = 0, \quad x_n := (1 - \varepsilon/d)^{1/(2\ln(n/2))}, \quad y_n := \Phi((2\ln_2(n))^{1/2}).$$

By l'Hospital's rule, the convergence follows if we verify that

$$\lim_{n \to \infty} \frac{\phi(\Phi^{-1}(y_n))}{\phi(\Phi^{-1}(x_n))} = 0.$$

Note that $\phi(\Phi^{-1}(y_n)) = (\sqrt{2\pi}\ln(n))^{-1} \to 0$ as $n \to \infty$. The Mill's ratio bound $(1 - \Phi(z))/\phi(z) < 1/z$ for z > 0 yields, with $z = \Phi^{-1}(x_n)$,

$$\phi(\Phi^{-1}(x_n)) > \Phi^{-1}(x_n)(1-x_n), \quad x_n > 1/2.$$

Hence,

$$\frac{\phi(\Phi^{-1}(y_n))}{\phi(\Phi^{-1}(x_n))} \le \frac{1}{\sqrt{2\pi}\ln(n)\Phi^{-1}(x_n)(1-x_n)}.$$

We claim that $\ln(n)(1-x_n)$ converges to a positive limit as $n \to \infty$. Since $\Phi^{-1}(x_n) \to \infty$ as $n \to \infty$ verifying this claim will prove the statement of the lemma. Note that

$$\left(1-\varepsilon/d\right)^{1/(2\ln(n/2))} = \exp\left(\frac{\ln\left(1-\varepsilon/d\right)}{2\ln(n/2)}\right)$$

and hence

$$1 + \frac{\ln(1 - \varepsilon/d)}{2\ln(n/2)} < (1 - \varepsilon/d)^{1/(2\ln(n/2))} < 1 + \frac{\ln(1 - \varepsilon/d)}{2\ln(n/2)} + \frac{1}{2} \left(\frac{\ln(1 - \varepsilon/d)}{2\ln(n/2)}\right)^2.$$

Hence, with $\gamma := -\ln(1 - \varepsilon/d)$,

$$\frac{\gamma}{2}\frac{\ln(n)}{\ln(n/2)} > \ln(n)(1-x_n) > \frac{\gamma}{2}\frac{\ln(n)}{\ln(n/2)} - \frac{\gamma^2}{8}\frac{\ln(n)}{\ln(n/2)^2}$$

which shows that $\lim_{n\to\infty} \ln(n)(1-x_n) = \gamma/2$. The proof is complete. \Box

A.2 Proof of Proposition 2

Proof. Note that (14) is equivalent to $u_{\varepsilon} = p_n^{-1}(\varepsilon/d)$. Lemma 7 says that $\lim_{n\to\infty} p_n^{-1}(\varepsilon/d)/c_n^2 = 0$. Note that c_n^2 given by (19) takes the form

$$c_n^2 = (d_n + (2\ln_2(n))^{1/2})^2,$$

where $\lim_{n\to\infty} d_n = 0$. The inequality $(a+b)^2 \le 2(a^2+b^2)$ gives

$$c_n^2 = (d_n + (2\ln_2(n))^{1/2})^2 \le 2(d_n^2 + 2\ln_2(n)).$$

Hence, $\lim_{n\to\infty} u_{\varepsilon}/\ln_2(n) = 0$ which completes the proof.