



LUND UNIVERSITY

Asymptotically Optimal Regression Trees

Mohlin, Erik

2018

Document Version:
Other version

[Link to publication](#)

Citation for published version (APA):

Mohlin, E. (2018). *Asymptotically Optimal Regression Trees*. (Working Papers ; No. 2018:12).

Total number of authors:

1

General rights

Unless other specific re-use rights are stated the following general rights apply:

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

Read more about Creative commons licenses: <https://creativecommons.org/licenses/>

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

LUND UNIVERSITY

PO Box 117
221 00 Lund
+46 46-222 00 00

Working Paper 2018:12

Department of Economics
School of Economics and Management

Asymptotically Optimal Regression Trees

Erik Mohlin

May 2018



LUND
UNIVERSITY

Asymptotically Optimal Regression Trees*

Erik Mohlin[†]
Lund University

This version May 22, 2018. First version May 5, 2014.

Abstract

Regression trees are evaluated with respect to mean square error (MSE), mean integrated square error (MISE), and integrated squared error (ISE), as the size of the training sample goes to infinity. The asymptotically MSE- and MISE minimizing (locally adaptive) regression trees are characterized. Under an optimal tree, MSE is $O(n^{-2/3})$. The estimator is shown to be asymptotically normally distributed. An estimator for ISE is also proposed, which may be used as a complement to cross-validation in the pruning of trees.

Keywords: Piece-Wise Linear Regression; Partitioning Estimators; Non-Parametric Regression; Categorization; Partition; Prediction Trees; Decision Trees; Regression Trees; Regressogram; Mean Squared Error.

JEL codes: C14; C38.

1 Introduction

Regression trees are an important and widely used tool in machine learning. They are easy to interpret and relatively fast to construct. Moreover they are resistant to the inclusion of irrelevant predictor variables, and are able to handle non-smooth regression surfaces (Hastie et al. 2009, chapter 10). The main drawback is that, for any particular data set, there is usually a more accurate method available. This paper contributes to the theoretical understanding of regression trees, and their predictive accuracy, by analysing mean square error (MSE), mean integrated square error (MISE), and integrated squared error (ISE), as the size of the training sample goes to infinity. The asymptotically MSE-

*This paper has benefited from comments by Liang Chen, Toru Kitagawa, and Bent Nielsen. I am grateful for financial support from the Swedish Research Council (Grant 2015-01751), and the Knut and Alice Wallenberg Foundation (Wallenberg Academy Fellowship 2016-0156).

[†]Address: Department of Economics, Lund University, Tycho Brahes väg 1, 220 07 Lund, Sweden.
E-mail: erik.mohlin@nek.lu.se

and MISE-minimising regression trees are characterized. Asymptotic normality is proved. Moreover, an estimator for ISE is proposed.

A regression tree is a tool for generating piece-wise constant regressions, and as such it can be seen as a particular method for constructing regressograms, introduced by Tukey (1947), and later developed and studied under the name of partitioning estimators (see Györfi et al. (2002) and references therein). Regression trees were introduced by Morgan and Sonquist (1963), and further developed by Breiman et al. (1984), Quinlan (1992), and Loh (2002), among others. Loh (2011) is a recent survey. By a successive procedure of binary splits, a set $\mathcal{X} \subseteq \mathbb{R}^d$ is split into a number of cells in the shape of hyper-rectangles. Together the cells partition the set \mathcal{X} . In each step of the splitting procedure, one dimension $\mathcal{X}_j \subseteq \mathbb{R}$ and one splitting point $z \in \mathcal{X}_j$ is used to divide \mathcal{X}_j into two halves. Thus the sequence of splits can be represented as a binary tree τ . The prediction for each cell is equal to the sample mean of the data in that cell. The criterion for evaluating a tree is usually based on the residual sum of squares $R(\tau)$, with the addition of a complexity cost α per split k . First one grows a large tree based on minimisation of $R(\tau)$, and then one reduces the number of cells/splits by pruning. Typically pruning is based on minimisation of $R_\alpha(\tau) = R(\tau) + \alpha k$, where the value of α is determined by some cross-validation procedure. In contrast the current paper studies the bias-variance trade-off analytically. The estimator of ISE that is proposed below may be used as a complement to cross-validation in the pruning of trees.

Decreasing the size of cells has two effects: The within-cell differences between objects tend to decrease, but the number of training observations in each cell tends to decrease, thereby making inference less reliable. It follows that as the size of the training set n is increased, the optimal number of cells is also increased, but at a slower rate, $O(n^{-1/3})$. Asymptotically the width of the cell to which x belongs, is (i) increasing in the variance of y conditional on x , (ii) decreasing in the derivative of the mean of y conditional on x , and (iii) decreasing in the marginal density at x . Under an optimal tree, MSE is $O(n^{2/3})$.

In the literature on regression trees the main focus has been on developing algorithms rather than deriving analytical results. Consequently there are few asymptotic results and no results on adaptive optimal cells. Similarly, Györfi et al. (2002) collect many results on consistency of partitioning estimators, but it does not seem to exist any results regarding locally adaptive optimal partitionings. The most closely related results are due to Cattaneo and Farrell (2013). They analyse the asymptotic mean square error of partitionings but do not allow the cells be locally adaptive.¹

My results on optimal regression trees can be viewed as lying in between those previously obtained for, on the one hand, asymptotically optimal locally adaptive kernels for non-parametric regression (e.g. Fan and Gijbels 1992), and on the other hand, asymp-

¹When their theorem 3 is restricted to the case of estimation of the regression line (conditional mean) by means of a constant fit, then it makes essentially the same statement about the mean square error as my theorem 1 when all cells are restricted to have the same shape and volume.

totically optimal locally adaptive histograms for non-parametric density estimation (e.g. Kogure 1987). The results partly generalize Mohlin (2014b).

2 Model

2.1 Data

Consider a sample, or training set, $\{(X^s, Y^s)\}_{s=1}^n$ of size n . Each realised observation is a $d + 1$ -dimensional attribute vector $(x^s, y^s) = (x_1^s, x_2^s, \dots, x_d^s, y^s) \in \mathcal{X} \times \mathcal{Y}$, where $\mathcal{Y} = \mathbb{R}$, and $\mathcal{X} = \times_{j=1}^d \mathcal{X}_j$ is a compact subset of \mathbb{R}^n , such that $\mathcal{X}_j = [a_j, b_j] \subseteq \mathbb{R}$ for all dimensions j . Observations are i.i.d. draws from a probability distribution with joint density $f(x, y)$ and marginal densities $f_X(x)$ and $f_Y(y)$.² The density of Y conditional on $X = x$ is $f_Y(y|x) = f(x, y) / f_X(x)$, assuming $f_X(x) > 0$ for all $x \in \mathcal{X}$. For all $X \in \mathcal{X}$, write

$$Y = m(X) + \varepsilon(X),$$

where $m(X) = \mathbb{E}[Y|X]$ is the conditional mean, and $\varepsilon(X) = Y - \mathbb{E}[Y|X]$ is a noise term with mean zero. Assume $|\mathbb{E}[Y|X = x]| < \infty$ so that the variance $\sigma^2(x) = \text{Var}(Y|X = x)$ is well defined.

The set of samples, or training sets, of size n is denoted $\mathcal{T}(n)$.

2.2 Trees

A tree τ induces a set of cells (or categories) $\mathcal{C} = \{\mathcal{C}^1, \dots, \mathcal{C}^k\}$ that partitions \mathcal{X} . Each cell is a hyper-rectangle i.e. $\mathcal{C}^i = \times_{j=1}^d \mathcal{C}_j^i$, where component \mathcal{C}_j^i is an interval (open, closed or half-open) of length h_j^i on \mathbb{R} , such that

$$\inf_{x_j} \mathcal{C}_j^i = a_j^i, \quad \sup_{x_j} \mathcal{C}_j^i = a_j^i + h_j^i = b_j^i.$$

We identify a tree τ by the set of cells (or categories) \mathcal{C} that it induces. The probability that an object belongs to cell/category i is $p^i = \int_{x \in \mathcal{C}^i} f_X(x) dx$, which is required to be strictly positive. The conditional marginal density of y in \mathcal{C}^i is $f_Y(y|x \in \mathcal{C}^i) = \int_{x \in \mathcal{C}^i} f(x, y) dx / p^i$, and the conditional marginal density of x in \mathcal{C}^i is $f_X(x|x \in \mathcal{C}^i) = f_X(x) / p^i$. Denote the within-category variance $\text{Var}(Y^i) = \text{Var}(Y|X \in \mathcal{C}^i)$, and the within-category mean $\mu^i = \mathbb{E}[m(X)|X \in \mathcal{C}^i]$.

The relative size of cells is constrained by some (small) number $\rho \in (0, 1)$ such that, for all i and j , if $p^i \geq p^j$ then $p^j \geq \rho p^i$. For any finite number of cells this simply means

²More formally observations are drawn i.i.d. according to an absolutely continuous cumulative distribution function $F : \mathcal{V} \rightarrow [0, 1]$, with a bounded probability density function, $f : \mathcal{V} \rightarrow \mathbb{R}_+$. The marginal densities are $f_X(x) = \int_{y \in \mathcal{Y}} f(x, y) dy$ and $f_Y(y) = \int_{x \in \mathcal{X}} f(x, y) dx$.

that all cells have positive probability. When the number of cells goes to infinity, the constraint implies that no cell becomes relatively infinitely larger than another cell. The set of categorizations satisfying the assumptions above, called feasible categorizations, is denoted Ψ . For a given training set t the set of feasible categorizations with non-empty cells is $\Psi(t)$.

2.3 Prediction Error and Optimality

Given a training set t , each cell \mathcal{C}^i is associated with a unique (point) prediction \hat{y}^i . Let $\hat{y}(x)$ be the prediction for the cell to which x belongs. The *squared error* (SE) associated with an object (x, y) is

$$SE(\mathcal{C}, t)(x, y) = (y - \hat{y}(x))^2. \quad (1)$$

In the statistical learning literature this is also known as the prediction error. Taking expectation over objects (X, Y) we obtain the *integrated squared error* (ISE), which is also known as generalization error, or test error,

$$ISE(\mathcal{C}, t) = \mathbb{E}[(Y - \hat{y}(X))^2] = \int \mathbb{E}[(Y(x) - \hat{y}(x))^2] f_X(x) dx. \quad (2)$$

Fixing $X = x$ and taking expectation over Y one may define a point-wise version of ISE , $ISE(\mathcal{C}, t)(x) = \mathbb{E}[(Y - \hat{y}(X))^2 | X = x]$.

Ex ante the training set is a random variable T , and hence the resulting prediction is also a random variable \hat{Y} . Fixing $X = x$ and taking the expectation of $SE(\mathcal{C}, t)(x)$ over \hat{Y} and Y yields the *mean square error* (MSE)

$$MSE(\mathcal{C}, n)(x) = \mathbb{E}\left[\left(Y(x) - \hat{Y}(x)\right)^2\right] = \mathbb{E}\left[\left(Y - \hat{Y}(X)\right)^2 | X = x\right]. \quad (3)$$

Now taking expectation over X yields the *mean integrated squared error* ($MISE$)

$$MISE(\mathcal{C}, n) = \mathbb{E}\left[\left(Y - \hat{Y}(X)\right)^2\right] = \int \mathbb{E}\left[\left(Y(x) - \hat{Y}(x)\right)^2\right] f_X(x) dx. \quad (4)$$

In the statistical learning literature (MSE) and ($MISE$) are also known as the expected prediction error or expected test error.

Remark 1 *Some authors (e.g. Härdle 1990) focus on conditional square error $(m(x) - \hat{y}(x))$, and hence on conditional MSE and conditional MISE, thereby ignoring the irreducible error $\sigma^2(x)$. Results would essentially be unaffected by using $(m(x) - \hat{y}(x))^2$ instead of $(y - \hat{y}(x))^2$ as the basis of evaluation. In particular the solution derived in theorems 1 and 2 would be the same.*

2.4 Prediction

The prediction rule remains to be specified. Fixing the training set one may look for predictions $\{\hat{y}^i\}_{i=1}^k$ that minimize $ISE(\mathcal{C}, t)$, for a given categorization \mathcal{C} . Thus, for each cell i the optimal prediction \hat{y}^i solves

$$\min_{\hat{y}^i} \mathbb{E} \left[(Y - \hat{y}^i)^2 | X \in \mathcal{C}^i \right]. \quad (5)$$

It is basic result of statistical decision theory that mean square error is minimized when the prediction is equal to the conditional mean (see e.g. theorem 2.1 of Li and Racine 2007). In the current set-up this means that the solution to (5) is $\hat{y}^i = \mathbb{E}[Y | X \in \mathcal{C}^i] = \mu^i$. The true mean μ^i is not observed, but the sample mean \bar{y}^i is an unbiased estimator of the true mean. For this reason the within-cell mean \bar{y}^i will be used as prediction, as long as the cell is non-empty. The choice of predictor for empty cells will turn out to be unimportant for the results. For simplicity I assume that the prediction for an empty cell is the mean across all observations, \bar{y} . The training set for cell i is $t^i = t \cap \mathcal{C}^i$. Let $n^i = |t^i|$, so that $\sum_{i=1}^k n^i = n$. The prediction for cell i is

$$\hat{y}^i = \begin{cases} \bar{y}^i = \frac{1}{n^i} \sum_{y^s \in t^i} y^s & \text{if } n^i > 0 \\ \bar{y} = \frac{1}{n} \sum_{s=1}^n y^s & \text{if } n^i = 0 \end{cases}. \quad (6)$$

Note that ex ante the number of observations in cell i is a random variable N^i .

3 Results

3.1 Preliminary Results

The first two results provide expressions for ISE , MSE and $MISE$.

Lemma 1 *The integrated square error for a categorization $\mathcal{C} \in \Psi(t)$, conditional on a training set t , is³*

$$ISE(\mathcal{C}, t) = \sum_{i=1}^k p^i \left(Var(Y^i) + (\bar{y}^i - \mu^i)^2 \right).$$

Lemma 2 *The mean squared error of categorization $\mathcal{C} \in \Psi$, at $x \in \mathcal{C}^i$, is*

$$MSE(\mathcal{C}, n)(x) = \sigma^2(x) + (Bias(\bar{Y}^i))^2 + Var(\bar{Y}^i),$$

where

$$(Bias(\bar{Y}^i))^2 = (m(x) - \mu^i)^2,$$

³Point-wise one may also calculate, $ISE(\mathcal{C}, t)(x) = \sigma^2(x) + (m(x) - \bar{y}_{it})^2$.

and

$$\text{Var}(\bar{Y}^i) = \text{Var}(Y^i) \sum_{r=1}^n \Pr(N^i = r) \frac{1}{r} + \Pr(N^i = 0) \mathbb{E} \left[(\bar{Y}_n - \mu^i)^2 | N^i = 0 \right],$$

with $N^i \sim \text{Bin}(n, p^i)$.

Lemma 3 τ The mean integrated squared error is

$$\begin{aligned} \text{MISE}(\mathcal{C}, n) &= \sum_{i=1}^k p^i \left(\mathbb{E}[\sigma^2(X) | X \in \mathcal{C}^i] + \text{Var}(m(X) | X \in \mathcal{C}^i) + \text{Var}(\bar{Y}^i) \right) \\ &= \sum_{i=1}^k p^i \left(\text{Var}(Y^i) + \text{Var}(\bar{Y}^i) \right). \end{aligned}$$

The expression for $\text{MSE}(\mathcal{C}, n)(x)$ reveals the fundamental bias-variance trade-off. In addition the prediction error is affected by *irreducible noise* $\sigma^2(x)$.

Remark 2 It is easy to see that ISE , MSE , and MISE are continuous in the distribution f . Mohlin (2014a) defines a metric on partitions which can be employed to show that prediction errors ISE , MSE , and MISE are also continuous in the decision variable \mathcal{C} .

3.2 MSE - and MISE -Optimal Trees

Intuitively, as n increases one may obtain a reasonable approximation of $m(x)$ by increasing the number of cells k , though at a slower rate than n . The following lemma formalizes this intuition:

Lemma 4 For any $\varepsilon > 0$ there are $n_0 > 0$ and $\delta > 0$, such that if $n > n_0$ and

$$\text{MISE}(\mathcal{C}, n) - \inf_{\mathcal{C}' \in \Psi} \text{MISE}(\mathcal{C}', n) < \delta,$$

then \mathcal{C} satisfies $1/k < \varepsilon$, and $k/n < \varepsilon$.

The lemma says that minimization of $\text{MISE}(\mathcal{C}, n)$ implies that if the number of observations goes to infinity ($n \rightarrow \infty$) then it is optimal to let the number of cells go to infinity too ($k \rightarrow \infty$), but at a slower rate ($k/n \rightarrow 0$). This is the basis for the asymptotic results that will now be derived.

3.2.1 Fixed Design

Consider first a fixed design model where data is taken such that $n^i = np^i$. (See e.g. Jennen-Steinmetz and Gasser (1988) for a similar fixed design assumption.) If f is three times differentiable, Taylor approximations may be used to derive the asymptotically optimal width of cells as $n \rightarrow \infty$.

As we let $n \rightarrow \infty$ there exist corresponding sequences of optimal trees (for each n there is at least one optimal tree τ^n). For a given sequence of optimal trees $\{\tau^n\}_{n=1}^\infty$ we examine the sequence of cells $\{\mathcal{C}^{i(n)}(x)\}_{n=1}^\infty$ that contain an arbitrary point x .

Theorem 1 *Let $\{\tau^n\}_{n=1}^\infty$ be a sequence of optimal trees, and let $\{\mathcal{C}^{i(n)}(x)\}_{n=1}^\infty$ be the sequence of cells that contain x . Suppressing the dependence on n , let the volume of $\mathcal{C}^{i(n)}(x)$ be denoted $h = \times_{l=1}^d h_l$. Suppose that f is three times differentiable at x , and $f_X(x) > 0$. Assume a fixed design model such that $n^i = np^i$. If $n \rightarrow \infty$, $k \rightarrow \infty$, and $k/n \rightarrow 0$ (implying $p^i \rightarrow 0$ and $np^i \rightarrow \infty$), then asymptotically,*

$$MSE(\mathcal{C}, n)(x) = \underbrace{\sigma^2(x) + h^2 \left(\sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \right)^2}_{Bias[\bar{Y}(x)]^2} + \underbrace{\frac{\sigma^2(x)}{f_X(x)nh}}_{Var[\bar{Y}(x)]} + O(n^{-1}), \quad (7)$$

where $\delta_j \in [0, 1]$ is such that $\left| \mathbb{E}[s_j | s_j \in \mathcal{C}_j^{i(n)}] - x_j \right| = \delta_j h_j$. If $\sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \neq 0$ then the asymptotic $MSE(\mathcal{C}, n)(x)$ is minimized by

$$h = \left(\frac{\sigma^2(x)}{2nf_X(x) \left(\sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \right)^2} \right)^{\frac{1}{3}}, \quad (8)$$

which yields

$$MSE(\mathcal{C}', n)(x) = \sigma^2(x) + O\left(n^{-\frac{2}{3}}\right). \quad (9)$$

The size of cells should decrease (and the number of cells should increase) at a rate $O(n^{-d/3})$. Asymptotically, the size of the MSE -minimising cell is decreasing in the density $f_X(x)$, and the curvature $\sum_j \frac{\partial m(x)}{\partial x_j}$ of the conditional mean. It is increasing in the variance $\sigma^2(x)$.

Remark 3 *If f is not continuously differentiable but Lipschitz continuous we may still derive an upper bound on the mean square error, and an expression for the associated width h^* . Restrict attention to $d = 1$ and a fixed design model with $n^i = np^i$. Assume*

that $\sigma^2(x)$, $m(x)$, and $f(x)$ are Lipschitz, so that for any $x, \alpha, \beta, \gamma \in \mathcal{C}^i$ it holds that

$$\begin{aligned} |\sigma^2(x) - \sigma^2(\alpha)| &< \lambda |x - \alpha|, \\ |m(x) - m(\beta)| &< \eta |x - \beta|, \\ |f_X(x) - f(\gamma)| &< \phi |x - \gamma|. \end{aligned}$$

In this case we find that

$$MSE(C, n)(x) < \sigma^2(x) + \eta^2 h^2 + \frac{\sigma^2(x)}{f(x)hn} + O(n^{-1}).$$

The right hand side is minimised by

$$h = \left(\frac{\sigma^2(x)}{2n\eta^2 f(x)} \right)^{\frac{1}{3}}.$$

Evaluating $MSE(C, n)(x)$ with this solution we find

$$MSE(C, n)(x) = \sigma^2(x) + O\left(n^{-\frac{2}{3}}\right).$$

The results are is similar to those obtained with Taylor approximations. Note that η is a measure of the roughness of the conditional mean, just like m' .

3.2.2 Random Design

We now verify that the speed of convergence is the same under a random design model as under the fixed design model.

Theorem 2 Let $\{\tau^n\}_{n=1}^\infty$ be a sequence of optimal trees, and let $\{\mathcal{C}^{i(n)}(x)\}_{n=1}^\infty$ be the sequence of cells that contain x . Supressing the dependence on n , let the volume of $\mathcal{C}^{i(n)}(x)$ be denoted $h = \times_{l=1}^d h_l$. Suppose that f is three times differentiable at x , and $f_X(x) > 0$. Assume a random design model. If $n \rightarrow \infty$, $k \rightarrow \infty$, and $k/n \rightarrow 0$ (implying $p^i \rightarrow 0$ and $np^i \rightarrow \infty$), then asymptotically,

$$MSE(\mathcal{C}, n)(x) = \sigma^2(x) + \underbrace{h^2 \left(\sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \right)^2}_{Bias[\bar{Y}(x)]^2} + O(h^3) + \underbrace{\frac{\sigma^2(x)}{f_X(x)} O((nh)^{-1})}_{Var[\bar{Y}(x)]} + O(n^{-1}), \quad (10)$$

where $\delta_j \in [0, 1]$ is such that $\left| \mathbb{E} \left[s_j | s_j \in \mathcal{C}_j^{i(n)} \right] - x_j \right| = \delta_j h_j$. If $\sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \neq 0$ then the

asymptotic $MSE(\mathcal{C}, n)(x)$ is minimized by

$$h = O\left(n^{-\frac{1}{3}}\right), \quad (11)$$

which yields

$$MSE(\mathcal{C}', n)(x) = \sigma^2(x) + O\left(n^{-\frac{2}{3}}\right). \quad (12)$$

3.2.3 Asymptotic Normality

The regression tree estimator is asymptotically normally distributed. This holds regardless of whether a fixed or random design is assumed.

Theorem 3 *Suppose that f is three times differentiable at $x \in \mathcal{C}^i$, a point in the interior of the support of X , and $f_X(x) > 0$, $\sigma^2(x) > 0$. If $n \rightarrow \infty$, $k \rightarrow \infty$, and $k/n \rightarrow 0$, then asymptotically,*

$$\sqrt{n^i} \left[\bar{Y}(x) - m(x) - h^2 \left(\sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \right)^2 \right] \xrightarrow{d} N(0, \sigma^2(x)).$$

3.3 ISE-Optimal Trees

Deriving the optimal trees by computing MSE and $MISE$ relies on knowledge of the underlying distribution f . In practice this information is not available to the analyst, so the optimal tree is usually determined via cross-validation. An alternative approach would be to find an estimator of MSE and $MISE$. I would like to suggest the following estimator:

Definition 1 *Let $\hat{\Psi}(t)$ denote the set of feasible trees in which all cells have at least two elements ($n^i \geq 2$) given the training set t . The sample integrated square error (SISE) for a tree $\mathcal{C} \in \hat{\Psi}(t)$, conditional on a training set t , is*

$$SISE(\mathcal{C}, t) = \sum_{i=1}^k \frac{n^i}{n} \left(1 + \frac{1}{n^i} \right) s_i^2, \quad s_i^2 = \frac{1}{n^i - 1} \sum_{y^s \in t^i} (y^s - \bar{y}^i)^2.$$

Note that this can be written

$$SISE(\mathcal{C}, t) = \sum_{i=1}^k \frac{1}{n} \frac{n^i + 1}{n^i - 1} \sum_{y^s \in t^i} (y^s - \bar{y}^i)^2.$$

In the absence of the factor $(n^i + 1)/(n^i - 1)$ this expression would simply be the average residual sum of squares (RSS), which is also known as training error in the statistical

learning literature. The factor $(n^i + 1) / (n^i - 1)$ penalises cells with few members. The motivation for this particular choice of adjustment of RSS comes from the following observations:

Theorem 4 *For a given tree τ , inducing categorization \mathcal{C} , and training set \tilde{t} with an allocation of observations to cells $\{\tilde{n}^1, \tilde{n}^2, \dots, \tilde{n}^k\}$, such that $\tilde{n}^i \geq 2$ for all i , let $\mathcal{T}(\mathcal{C}, \tilde{t})$ be the set of training sets t such that $\tilde{n}^i = n^i$ for each cell i in \mathcal{C} .*

(a) *If expectation is taken over $\mathcal{T}(\mathcal{C}, \tilde{t})$, then*

$$\mathbb{E}[ISE(\mathcal{C}, t)] = \sum_{i=1}^k p^i \left(1 + \frac{1}{n^i}\right) \text{Var}(Y^i),$$

and

$$\mathbb{E}[SISE(\mathcal{C}, t)] = \sum_{i=1}^k \frac{n^i}{n} \left(1 + \frac{1}{n^i}\right) \text{Var}(Y^i).$$

(b) *For any sequence of training sets from $\mathcal{T}(\mathcal{C}, \tilde{t})$*

$$P \lim_{n \rightarrow \infty} ISE(\mathcal{C}, t) = P \lim_{n \rightarrow \infty} SISE(\mathcal{C}, t) = \sum_{i=1}^k p^i \text{Var}(Y^i).$$

Part (a) of the theorem implies that if the actual fraction of objects in each cell, n^i/n , is equal to the probability of receiving an object in the corresponding cell p^i (as it was under the fixed design assumption above), then $ISE(\mathcal{C}, t)$ and $SISE(\mathcal{C}, t)$ have the same expected value on $\mathcal{T}(\mathcal{C}, t)$ (the set of training sets such that $\tilde{n}^i = n^i$ for each cell i in \mathcal{C}). Part (b) confirms that for large enough n it is highly probable that $n^i/n = p^i$.

The estimator $SISE(\mathcal{C}, t)$ could be employed in the pruning of regression trees. (Note that $SISE(\mathcal{C}, t)$ is well-defined also when the input variables are categorical rather than numerical.) This would not require any use of cross-validation, since the cost-complexity trade-off is decided on analytically. Indeed a $SISE(\mathcal{C}, t)$ -based pruning process might be seen as a complement to traditional cross-validation methods. It should be straightforward to modify the tree algorithm in the statistical software R since the command `prune.tree` has an optional `method` argument whose default `method="deviance"` minimises average RSS.

4 Discussion

Non-parametric kernel regression is another common method for estimating the conditional mean $m(x)$. To facilitate comparison with the results above restrict attention to

the one-dimensional case, $d = 1$. In this case (11) becomes

$$h_C^*(x) = \left(\frac{\sigma^2(x)}{2nf_X(x)(\delta m'(x))^2} \right)^{1/3}. \quad (13)$$

Fan and Gijbels (1992) derive the following expression for the locally adaptive asymptotically optimal kernel bandwidth (where optimal is understood in the sense of minimizing conditional MISE),

$$h_K^*(x) = q_K \left(\frac{\sigma^2(x)}{nf_X(x)(m''(x))^2} \right)^{1/5}, \quad (14)$$

where q_K is a constant which is independent of x . This is very similar to the expression (13) for the optimal cell width $h_C^*(x)$. One difference is that the curvature of the conditional mean enters through the second derivative $m''(x)$ here, compared to the first derivative $m'(x)$ above. The reason is that the kernel is symmetric around x , whereas the cell induced by a regression tree is not, except for a measure zero set of points. Another difference is that speed at which the bandwidth vanishes is slower than the speed at which the optimal cell width vanishes.

The kernels due to Priestley and Chao (1972), and Gasser and Müller (1979), lead to expressions that are similar to (14), see e.g. Härdle (1990) and Brockmann et al. (1993). The kernel of Nadaraya (1964) and Watson (1964) (being of first order) induces a bias term that involves both the first and second order derivatives of m at x .

Within density estimation histograms is the closest related approach. Kogure (1987) (see also Scott 1992) derives the following expression for the locally adaptive asymptotically optimal bin width,

$$h_H^*(x) = q_H \left(\frac{f_X(x)}{n(f'_X(x))^2} \right)^{1/3}, \quad (15)$$

where q_H is a constant which is independent of x . In contrast to (13) and (14) this expression is independent of the variance $\sigma^2(x)$, and the optimal width at x is increasing in density $f_X(x)$, rather than decreasing. Interestingly, the width of the bins optimally vanishes at the rate $O(n^{-1/3})$, exactly as in (13).

In conclusion, the optimality result in theorem 2 lies between the results derived for kernels in regression analysis and histograms in density estimation. This reflects the fact that while the prediction *task* of regression trees and is essentially the same as that of kernel regression, the *tool* (producing an estimate for each cell in a partition) is more similar to histograms.

5 Binary Outcomes / Classification

The framework developed above can be amended to handle the task of estimating the probability of binary outcomes, with the aid of a tree. The binary outcome space could signify membership in one of two classes. However, the question that is being asked is different from the one posed by classification trees (e.g. Ripley (1996)). In that literature the prediction is class membership, not the probability of a particular class membership. Moreover the framework of this paper uses loss functions based on squared error, something that is rarely the criterion for evaluation of classification trees.

5.1 Model

5.1.1 Data and Trees

Instead of a continuous outcome set, there is a binary outcome set $\mathcal{Y} = \{0, 1\}$. Furthermore, instead of making a point prediction of what value Y will take (in the set \mathcal{Y}), the task will be to estimate the probability that Y takes either of the values in $\mathcal{Y} = \{0, 1\}$. The density $f(x, y)$, the marginal densities $f_X(x)$ and $f_Y(y)$, and the conditional density $f_Y(y|x)$ are all defined as above. As before, let $m(x) = \mathbb{E}[Y|X = x]$. Note that $\mathbb{E}[Y|X = x] = f_Y(1|x)$. Thus the task is to estimate

$$m(x) = f_Y(1|x) = \mathbb{E}[Y = 1|X = x] \in [0, 1].$$

One may still use the decomposition $Y = m(X) + \varepsilon(X)$.

A tree is defined as above, along with the probability p^i , the conditional marginal densities $f_Y(y|x \in \mathcal{C}^i)$ and $f_X(x|x \in \mathcal{C}^i)$. Feasible categorizations are also defined as before.

5.1.2 Prediction Error

Given a training set t , each cell \mathcal{C}^i is associated with a unique (point) estimate \hat{m}^i . Let $\hat{m}(x)$ be the estimate for the cell to which x belongs. Using $\hat{m}(x)$ instead of $\hat{y}(x)$, and $m(x)$ instead of $y(x)$ (and $Y(x)$), in equations (1), (2), (3), and (4), redefine squared error (SE), integrated squared error (ISE), mean square error (MSE), and mean integrated squared error (MISE), for the present set-up.⁴ For example

$$ISE'(\mathcal{C}, t) = \mathbb{E}[m(X) - \hat{m}(X)]^2 = \int \mathbb{E}[(m(x) - \hat{m}(x))^2] f_X(x) dx.$$

⁴This brings the definitions closer to the ones of Härdle (1990), mentioned above.

and

$$MSE'(\mathcal{C}, n)(x) = \mathbb{E} \left[\left(m(x) - \hat{M}(x) \right)^2 \right] = \mathbb{E} \left[\left(m(X) - \hat{M}(X) \right)^2 | X = x \right].$$

Here \hat{M} reflects the fact that the estimator is itself a random variable, just as the use of $\hat{Y}(X)$ in equation (3)

5.1.3 Prediction

By the same line of reasoning as before, look for a prediction \hat{y}^i , for each cell i , that minimize

$$\min_{\hat{y}^i} \mathbb{E} \left[\left(m(X) - \hat{m}^i \right)^2 | X \in \mathcal{C}^i \right],$$

and for the same reasons as before the solution is $\hat{m}^i = \mathbb{E} [m(X) | X \in \mathcal{C}^i] = \mu^i$. Again, the true mean μ^i is not observed, but the sample mean \bar{y}^i is an unbiased estimator of it. Consequently the estimator of \hat{m}^i should be equal to the predictor \hat{y}^i used above:

$$\hat{m}^i = \begin{cases} \bar{y}^i = \frac{1}{n^i} \sum_{y^s \in t_i} y^s & \text{if } n^i > 0 \\ \bar{y} = \frac{1}{n} \sum_{s=1}^n y^s & \text{if } n^i = 0 \end{cases}.$$

5.2 Results

The results are almost identical to the case of a continuous outcome variable. The only difference is that the irreducible noise vanishes. This is due to the fact that Y has been replaced with $m(X)$ in the definitions of prediction error. I only review the results for (the revised version of) mean square error. Corresponding to lemma 2 the following holds:

Lemma 5 *The mean squared error of categorization $\mathcal{C} \in \Psi$, at $x \in \mathcal{C}^i$, is*

$$MSE'(\mathcal{C}, n)(x) = \left(Bias(\bar{Y}^i) \right)^2 + Var(\bar{Y}^i),$$

where $\left(Bias(\bar{Y}^i) \right)^2$ and $Var(\bar{Y}^i)$ are the same as in lemma 2.

In the same way theorem 1 is recovered, the only difference being that the irreducible noise is taken away from equation (7).

Theorem 5 *Let $\{\tau^n\}_{n=1}^\infty$ be a sequence of optimal trees, and let $\{\mathcal{C}^{i(n)}(x)\}_{n=1}^\infty$ be the sequence of cells that contain x . Suppressing the dependence on n , let the volume of $\mathcal{C}^{i(n)}(x)$ be denoted $h = \times_{l=1}^d h_l$. Suppose that f is three times differentiable at x , and $f_X(x) > 0$. Assume a fixed design model such that $n^i = np^i$. If $n \rightarrow \infty$, $k \rightarrow \infty$, and $k/n \rightarrow 0$*

(implying $p^i \rightarrow 0$ and $np^i \rightarrow \infty$), then asymptotically,

$$MSE(\mathcal{C}, n)(x) = \underbrace{h^2 \left(\sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \right)^2}_{Bias[\bar{Y}(x)]^2} + O(h^3) + \underbrace{\frac{\sigma^2(x)}{f_X(x)nh}}_{Var[\bar{Y}(x)]} + O(n^{-1}),$$

where $\delta_j \in [0, 1]$ is such that $|\mathbb{E}[s_j | s_j \in \mathcal{C}_j^i] - x_j| = \delta_j h_j$. If $\sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \neq 0$ then the asymptotic $MSE(\mathcal{C}, n)(x)$ is minimized by (8), which yields (9).

6 Conclusion

Hopefully the theoretical investigations of this paper will somehow prove useful in the further development and application of regression trees and related machine learning methods.

References

- Breiman, L., Friedman, J., Olshen, R. and Stone, C. (1984), *Classification and Regression Trees*, Wadsworth, New York.
- Brockmann, M., Gasser, T. and Herrmann, E. (1993), ‘Locally adaptive bandwidth choice for kernel regression estimators’, *Journal of the American Statistical Association* **88**(424), 1302–1309.
- Cattaneo, M. D. and Farrell, M. H. (2013), ‘Optimal convergence rates, bahadur representation, and asymptotic normality of partitioning estimators’, *Journal of Econometrics* **174**(2), 127–143.
- Chung, F. and Lu, L. (2002), ‘Connected components in random graphs with given expected degree sequences’, *Annals of Combinatorics* **6**(2), 125–145.
- Fan, J. and Gijbels, I. (1992), ‘Variable bandwidth and local linear regression smoothers’, *The Annals of Statistics* **20**, 2008–2036.
- Gasser, T. and Müller, H.-G. (1979), Kernel estimation of regression functions, in T. Gasser and M. Rosenblatt, eds, ‘Smoothing Techniques for Curve Estimation’, Springer Verlag, New York, pp. 23–68.
- Györfi, L., Krzyżak, A., Kohler, M. and Walk, H. (2002), *A distribution-free theory of nonparametric regression*, Springer.

- Hastie, T., Tibshirani, R. and Friedman, J. (2009), *The Elements of Statistical Learning Theory*, Springer, New York.
- Härdle, W. (1990), *Applied Nonparametric Regression*, Cambridge University Press, Cambridge, England.
- Jennen-Steinmetz, C. and Gasser, T. (1988), ‘A unifying approach to nonparametric regression estimation’, *Journal of the American Statistical Association* **83**(404), 1084–1089.
- Kogure, A. (1987), ‘Asymptotically optimal cells for a histogram’, *Annals of Statistics* **15**, 1023–1030.
- Li, Q. and Racine, J. S. (2007), *Nonparametric Econometrics: Theory and Practice*, Princeton University Press.
- Loh, W.-Y. (2002), ‘Regression trees with unbiased variable selection and interaction detection’, *Statistica Sinica* **12**(2), 361–386.
- Loh, W.-Y. (2011), ‘Classification and regression trees’, *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* **1**(1), 14–23.
- Mohlin, E. (2014a), A metric for measurable partitions. Mimeo.
- Mohlin, E. (2014b), ‘Optimal categorization’, *Journal of Economic Theory* **152**, 356–381.
- Morgan, J. M. and Sonquist, J. A. (1963), ‘Problems in the analysis of survey data, and a proposal’, *Journal of the American Statistical Association* **58**, 415–434.
- Nadaraya, E. (1964), ‘On estimating regression’, *Theory of Probability and Its Applications* **9**(1), 141–142.
- Priestley, M. B. and Chao, M. T. (1972), ‘Non-parametric function fitting’, *Journal of the Royal Statistical Society. Series B* **34**, 385–392.
- Quinlan, J. R. (1992), Learning with continuous classes, in ‘Proceedings of the 5th Australian joint Conference on Artificial Intelligence’, Vol. 92, Singapore, pp. 343–348.
- Ripley, B. D. (1996), *Pattern Recognition and Neural Networks*, Cambridge University Press, Cambridge.
- Scott, D. W. (1992), *Multivariate Density Estimation*, John Wiley, New York.
- Tukey, J. W. (1947), ‘Non-parametric estimation ii. statistically equivalent blocks and tolerance regions—the continuous case’, *The Annals of Mathematical Statistics* **18**(4), 529–539.

Watson, G. S. (1964), ‘Smooth regression analysis’, *Sankhyā: The Indian Journal of Statistics, Series A* pp. 359–372.

7 Appendix

7.1 Proofs: Preliminaries

Proof of Lemma 1.

$$\begin{aligned}
ISE(\mathcal{C}, t) &= \mathbb{E}[(Y|X - \bar{y}(X))^2] \\
&= \sum_i p^i \mathbb{E}[(Y|X - \bar{y}^i)^2 | X \in \mathcal{C}^i] \\
&= \sum_i p^i \mathbb{E}[(Y|X - \mu^i)^2 + (\bar{y}^i - \mu^i)^2 - 2(Y|X - \mu^i)(\bar{y}^i - \mu^i) | X \in \mathcal{C}^i] \\
&= \sum_i p^i (Var(Y^i) + (\bar{y}^i - \mu^i)^2)
\end{aligned}$$

■

Proof of Lemma 2. Suppose $x \in \mathcal{C}^i$. Start by noting

$$\begin{aligned}
MSE(\mathcal{C}, n)(x) &= \mathbb{E}[(Y(x) - \bar{Y}^i)^2] \\
&= \mathbb{E}[(Y(x) - m(x))^2 + (m(x) - \bar{Y}^i)^2 + 2(Y(x) - m(x))(m(x) - \bar{Y}^i)] \\
&= \sigma^2(x) + \mathbb{E}[(m(x) - \bar{Y}^i)^2]
\end{aligned}$$

Using $m(x) - \bar{Y}^i = m(x) - \mu^i - (\bar{Y}^i - \mu^i)$ we have

$$\mathbb{E}[(m(x) - \bar{Y}^i)^2] = Var(\bar{Y}^i) + (m(x) - \mu^i)^2.$$

Furthermore

$$Var(\bar{Y}^i) = \sum_{r=1}^n \Pr(N^i = r) \mathbb{E}[(\bar{Y}^i - \mu^i)^2 | N^i = r] + \Pr(N^i = 0) \mathbb{E}[(\bar{Y}_n - \mu^i)^2 | N^i = 0].$$

Note that if $r \geq 1$ then $E[\bar{Y}^i | N^i = r] = \mu^i$, so

$$\mathbb{E}[(\bar{Y}^i - \mu^i)^2 | N^i = r] = Var(\bar{Y}^i | N^i = r) = \frac{1}{r} Var(Y^i).$$

It is evident that the number of objects in a cell, N^i , has a binomial distribution with parameters p^i and n . ■

Proof of Lemma 3. It is straightforward to verify that

$$\begin{aligned} \text{Var}(Y^i) &= \int_{x \in \mathcal{C}^i} \sigma^2(x) f_X(x|x \in \mathcal{C}^i) dx + \int_{x \in \mathcal{C}^i} (m(x) - \mu^i)^2 f_X(x|x \in \mathcal{C}^i) dx \\ &= \mathbb{E}[\sigma^2(X) | X \in \mathcal{C}^i] + \text{Var}(m(X) | X \in \mathcal{C}^i). \end{aligned} \quad (16)$$

Using this we have

$$\begin{aligned} \text{MISE}(\mathcal{C}, n) &= \sum_i p^i \int_{x \in \mathcal{C}^i} \left(\sigma^2(x) + \text{Var}(\bar{Y}^i) + (m(x) - \mu^i)^2 \right) f_X(x|x \in \mathcal{C}^i) dx \\ &= \sum_i p^i \left(\mathbb{E}[\sigma^2(X) | X \in \mathcal{C}^i] + \text{Var}(\bar{Y}^i) + \text{Var}(m(X) | X \in \mathcal{C}^i) \right). \end{aligned}$$

■

7.2 Proofs: *MSE*- and *MISE* Optimality

Lemma 4 was proved in Mohlin (2014b).

7.2.1 Fixed design

Proof of Theorem 1. By lemma 4 an asymptotically optimal tree must satisfy $k \rightarrow \infty$ and $k/n \rightarrow 0$ as $n \rightarrow \infty$. This allows us to use lemma 2, to write, for $x \in \mathcal{C}^i$,

$$\text{MSE}(\mathcal{C}, n)(x) = \sigma^2(x) + \underbrace{(m(x) - \mu^i)^2}_{(\text{Bias}[\bar{Y}(x)])^2} + \underbrace{\text{Var}(Y^i) \frac{1}{np^i}}_{\text{Var}[\bar{Y}(x)]}. \quad (17)$$

Also note

$$\text{Var}(Y^i) = \int_{s \in \mathcal{C}^i} \left[\sigma^2(s) + (m(s) - \mu^i)^2 \right] f_X(s|s \in \mathcal{C}^i) ds. \quad (18)$$

A Taylor approximation yields,

$$\begin{aligned} p^i &= \int_{a_d^i}^{b_d^i} \dots \int_{a_1^i}^{b_1^i} [f_X(s)] ds_1 \dots ds_d \\ &= \int_{a_d^i}^{b_d^i} \dots \int_{a_1^i}^{b_1^i} [f_X(x) + D(f_X(x))(s - x) + \dots] ds_1 \dots ds_d. \end{aligned}$$

Note that if $x_j \neq a_j^i + h_j/2$, then

$$\int_{a_j^i}^{b_j^i} (s_j - x_j) ds_j = \left(a_j^i + \frac{h_j}{2} - x_j \right) h_j^i = \left(\frac{h_j}{2} + O(h_j) \right) h_j = O(h_j^2),$$

and if $x_j = a_j^i + h_j/2$, then this is zero, so if there is at least one j such that $x_j \neq a_j^i + h_j/2$, then

$$\begin{aligned} \int_{a_d^i}^{b_d^i} \dots \int_{a_1^i}^{b_1^i} [D(f_X(x))(s-x)] ds_1 \dots ds_d &= \int_{a_d^i}^{b_d^i} \dots \int_{a_1^i}^{b_1^i} \left[\sum_j \frac{\partial f_X(x)}{\partial x_j} (s_j - x_j) \right] ds_1 \dots ds_d \\ &= O\left(\times_{l=1}^d h_l\right) O(h_j), \end{aligned}$$

and if $x_j = a_j^i + h_j/2$, for all j , then this is zero. Moreover, it can be verified that

$$\int_{a_l^i}^{b_l^i} \frac{1}{2} (s-x) D^2(f_X(x))(s-x) ds_l = \frac{1}{2} \sum_{i \neq l} \sum_{j \neq l} \frac{\partial^2 f_X(x)}{\partial x_j \partial x_i} (s_i - x_i) (s_j - x_j) h_l + O(h_l^3),$$

so

$$\begin{aligned} &\int_{a_d^i}^{b_d^i} \dots \int_{a_1^i}^{b_1^i} \left[\frac{1}{2} (s-x) D^2(f_X(x))(s-x) \right] ds_1 \dots ds_d \\ &= \left(\int_{x_1-h_1/2}^{x_1+h_1/2} \frac{\partial^2 f_X(x)}{\partial^2 x_1} (s_1 - x_1)^2 ds_1 \right) h_2 h_3 \dots h_d + \dots \\ &+ h_1 h_3 \dots h_{n-1} \left(\int_{x_d-h_2/2}^{x_d+h_2/2} \frac{\partial^2 f_X(x)}{\partial^2 x_d} (s_d - x_d)^2 ds_d \right) \\ &= O\left(\times_{l=1}^d h_l \cdot \sum_j h_j^2\right). \end{aligned}$$

Thus if there is some j with $x_j \neq a_j^i + h_j/2$, then

$$p^i = \left(\times_{l=1}^d h_l\right) [f_X(x)] + O\left(\times_{l=1}^d h_l \cdot h_j\right),$$

and if not then

$$p^i = \left(\times_{l=1}^d h_l\right) [f_X(x)] + \left(\times_{l=1}^d h_l \cdot \sum_j h_j^2\right).$$

Consequently, we can assume

$$\frac{1}{np^i} = \frac{1}{f_X(x) \cdot n \cdot (\times_{l=1}^d h_l)} + O\left((\times_{l=1}^d h_l \cdot h_j \cdot n)^{-1}\right). \quad (19)$$

Another Taylor approximation yields

$$\mu^i = \int_{s \in \mathcal{C}^i} \left[m(x) + D(m(x))(s-x) + \frac{1}{2}(s-x)D^2(m(x))(s-x) + \dots \right] f_X(s|s \in \mathcal{C}^i) ds.$$

Note that

$$\int_{a_j^i}^{b_j^i} s_j f_X(s_j|s_j \in \mathcal{C}_j^i) ds_j = \mathbb{E}[s_j|s_j \in \mathcal{C}_j^i],$$

so

$$\int_{a_d^i}^{b_d^i} \dots \int_{a_1^i}^{b_1^i} \left[\frac{\partial m(x)}{\partial x_j} (s_j - x_j) f_X(s|s \in \mathcal{C}^i) \right] ds_1 \dots ds_d = \frac{\partial m(x)}{\partial x_j} \delta_j (\times_{l=1}^n h_l).$$

and

$$\begin{aligned} & \int_{a_d^i}^{b_d^i} \dots \int_{a_1^i}^{b_1^i} [D(m(x))(s-x) f_X(s|s \in \mathcal{C}^i)] ds_1 \dots ds_d \\ &= \int_{a_d^i}^{b_d^i} \dots \int_{a_1^i}^{b_1^i} \left[\sum_j \frac{\partial m(x)}{\partial x_j} (s_j - x_j) f_X(s|s \in \mathcal{C}^i) \right] ds_1 \dots ds_d \\ &= (\times_{l=1}^n h_l) \sum_j \frac{\partial m(x)}{\partial x_j} \delta_j. \end{aligned}$$

Furthermore, note that

$$\int_{a_d^i}^{b_d^i} \dots \int_{a_1^i}^{b_1^i} \left[\frac{1}{2} (s-x) D^2(m(x))(s-x) f_X(s|s \in \mathcal{C}^i) \right] ds_1 \dots ds_d = O\left((\times_{l=1}^d h_l)^2\right).$$

Thus,

$$\mu^i = m(x) + (\times_{l=1}^d h_l) \sum_j \frac{\partial m(x)}{\partial x_j} \delta_j + O\left((\times_{l=1}^d h_l)^2\right),$$

and consequently

$$(m(x) - \mu^i)^2 = \left((\times_{l=1}^d h_l) \sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \right)^2 + O\left((\times_{l=1}^d h_l)^3\right). \quad (20)$$

Next use this to find

$$\begin{aligned}
& \int_{s \in \mathcal{C}^i} (m(s) - \mu^i)^2 f_X(s|s \in \mathcal{C}^i) ds \\
&= (\times_{l=1}^d h_l)^2 \int_{a_d^i}^{b_d^i} \dots \int_{a_1^i}^{b_1^i} \left(\sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \right)^2 f_X(s|s \in \mathcal{C}^i) ds_1 \dots ds_d \\
&+ O\left((\times_{l=1}^d h_l)^3\right)
\end{aligned} \tag{21}$$

Finally, a Taylor approximation yields

$$\int_{s \in \mathcal{C}^i} \sigma^2(s) f_X(s|s \in \mathcal{C}^i) ds = \sigma^2(x) + \sum_j \frac{\partial \sigma^2(x)}{\partial x_j} (\times_{l=1}^d h_l) \delta_j + O\left((\times_{l=1}^d h_l)^3\right). \tag{22}$$

Thus, using (20) and (21), equation (18) becomes,

$$Var(Y^i) = \sigma^2(x) + O\left(\times_{l=1}^d h_l\right).$$

Using (18)-(22) in (17) gives (7). This can be simplified to

$$MSE(\mathcal{C}, n)(x) = \sigma^2(x) + h^2 \left(\sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \right)^2 + \frac{\sigma^2(x)}{f_X(x) nh}. \tag{23}$$

Maximise w.r.t. h . The first order condition is

$$\frac{\partial MSE(\mathcal{C}, n)(x)}{\partial h} = 2h \left(\sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \right)^2 - \frac{\sigma^2(x)}{f_X(x) nh^2} = 0,$$

or (8), and the second order condition is satisfied,

$$\frac{\partial^2 MSE(\mathcal{C}, n)(x)}{\partial h^2} = \frac{\sigma^2(x)}{nf_X(x)} \frac{1}{h^3} > 0.$$

Evaluating (23) at (8) yields (9). ■

The proof of the claims made in remark 3 are omitted, but available upon request.

7.2.2 Random Design

The proof of theorem 2 relies on lemmas 6 and 7 below.

Lemma 6 *If $n \rightarrow \infty$, $k \rightarrow \infty$, and $k/n \rightarrow 0$, then $\sum_{r=1}^n \Pr(N^i = r)/r = O(1/n p^i)$.*

Proof of Lemma 6. Note that $k/n \rightarrow 0$ implies $np^i \rightarrow \infty$. A Chernoff inequality for the binomial distribution (lemma 2.1 of Chung and Lu 2002) states that, for any real number λ ,

$$\Pr(N \leq tp - \lambda) \leq e^{\frac{-\lambda^2}{2tp^i}},$$

so for $\lambda = \frac{1}{2}np^i$,

$$\Pr\left(N \leq \frac{1}{2}tp\right) \leq e^{\frac{-np^i}{8}}.$$

Let $r^* = \lfloor \frac{1}{2}np^i \rfloor$. Thus

$$\begin{aligned} \sum_{r=1}^n \Pr(N^i = r) \frac{1}{r} &= \sum_{r=1}^{r^*-1} \Pr(N^i = r) \frac{1}{r} + \sum_{r=r^*}^n \Pr(N^i = r) \frac{1}{r} \\ &\leq \sum_{r=1}^{r^*-1} \Pr(N^i = r) + \sum_{r=r^*}^n \Pr(N^i = r) \frac{1}{r^*} \\ &\leq \Pr(N \leq r^* - 1) + \frac{1}{r^*} \\ &\leq e^{\frac{-np^i}{8}} + \frac{1}{\lfloor \frac{1}{2}np^i \rfloor} \\ &\leq e^{\frac{-np^i}{8}} + \frac{1}{\frac{1}{2}np^i - 1}. \end{aligned}$$

Since $np^i/e^{np^i} \rightarrow 0$, as $np^i \rightarrow \infty$, we have

$$O\left(e^{\frac{-np^i}{8}} + \frac{1}{\frac{1}{2}np^i - 1}\right) = O\left(\frac{1}{\frac{1}{2}np^i - 1}\right) = O\left(\frac{1}{np^i}\right).$$

■

Lemma 7 *If $n \rightarrow \infty$, $k \rightarrow \infty$, and $k/n \rightarrow 0$, then, for all i ,*

$$\Pr(N^i = 0) \mathbb{E}\left[(\bar{Y}_n - \mu^i)^2 | N^i = 0\right] = O\left(e^{-\frac{n}{k}}\right).$$

Proof of Lemma 7. Using $\mathbb{E}[\bar{Y}_n | N^i = 0] = \mu_{-i}$, we have

$$\begin{aligned} \mathbb{E}\left[(\bar{Y}_n - \mu^i)^2 | N^i = 0\right] &= \mathbb{E}\left[(\bar{Y}_n - \mu_{-i})^2 | N^i = 0\right] + (\mu_{-i} - \mu^i)^2 \\ &= \text{Var}(\bar{Y}_n | N^i = 0) + (\mu_{-i} - \mu^i)^2 \\ &= \frac{1}{n} \text{Var}(Y | X \notin \mathcal{C}^i) + (\mu_{-i} - \mu^i)^2. \end{aligned} \tag{24}$$

Next consider $\Pr(N^i = 0) = (1 - p^i)^n$. Let $p_{\max} = \max_i p^i$ and $p_{\min} = \min_i p^i$. Note that $p_{\min} > 1 - (k - 1)p_{\max}$. Since $p_{\min}/\rho \geq p_{\max}$ we have $p_{\min} > 1 - (k - 1)p_{\min}/\rho$, or equivalently $p_{\min} > \rho/(\rho + (k - 1))$. Hence

$$(1 - p^i)^n \leq (1 - p_{\min})^n \leq \left(\left(1 - \frac{\rho}{\rho + k - 1} \right)^{(\rho + k - 1)} \right)^{\frac{n}{\rho + k - 1}}.$$

As $\rho + k - 1 \rightarrow \infty$ this approaches $e^{-\rho \frac{n}{\rho + k - 1}}$. The desired result follows. ■

Proof of Theorem 2. By lemma 4 an asymptotically optimal tree must satisfy $k \rightarrow \infty$ and $k/n \rightarrow 0$ as $n \rightarrow \infty$. This allows us to use lemma 6 and lemma 7 together with lemma 2, to write, for $x \in \mathcal{C}^i$,

$$MSE(\mathcal{C}, n)(x) = \sigma^2(x) + \underbrace{(m(x) - \mu^i)^2}_{Bias[\bar{Y}(x)]} + \underbrace{Var(Y^i) O\left(\frac{1}{p^i n}\right)}_{Var[\bar{Y}(x)]} + O\left(e^{-\frac{n}{k}}\right). \quad (25)$$

Using (18)-(22) (from the proof of theorem 1) in (25) gives

$$\begin{aligned} MSE(\mathcal{C}, n)(x) &= \sigma^2(x) + h^2 \left(\sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \right)^2 + O(h^3) \\ &\quad + \frac{\sigma^2(x)}{f_X(x)} O((n \cdot h)^{-1}) + O(n^{-1}) + O\left(e^{-\frac{n}{k}}\right). \end{aligned} \quad (26)$$

From above, we know that asymptotically $p^i = O(h)$, and $k^{-1} = O(h)$ so that $O\left(e^{-\frac{n}{k}}\right) = O\left(e^{-np^i}\right)$. Since $np^i \cdot e^{-np^i} \rightarrow 0$ this means that (26) can be simplified to (10) or

$$MSE(\mathcal{C}, n)(x) = \sigma^2(x) + h^2 \left(\sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \right)^2 + \frac{\sigma^2(x)}{f_X(x)} O((nh)^{-1}).$$

Maximise w.r.t. h . The first order condition is

$$2h \left(\sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \right)^2 = \frac{\sigma^2(x)}{f_X(x)} O(n^{-1}h^{-2}),$$

and the second order condition is satisfied, so optimally (11). Using this in (10) yields (12). Having established that $h = O(n^{-2/3})$ it can be shown that $e^{-\frac{n}{k}}$ vanishes faster than h^3 and n^{-1} . ■

7.2.3 Asymptotic Normality

Proof of Theorem 3. Since

$$\mathbb{E} [\bar{Y}(x)] - m(x) = \text{Bias}(\bar{Y}(x)) = h^2 \left(\sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \right)^2 + O(h^3),$$

we have, for $x \in \mathcal{C}^i$,

$$\begin{aligned} & \sqrt{\frac{n^i}{\text{Var}(Y^i)}} \left[\bar{Y}(x) - m(x) - h^2 \left(\sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \right)^2 \right] \\ &= \sqrt{\frac{n^i}{\text{Var}(Y^i)}} [\bar{Y}(x) - \mathbb{E}[\bar{Y}(x)]] + \sqrt{\frac{n^i}{\text{Var}(Y^i)}} \left[\mathbb{E}[\bar{Y}(x)] - m(x) - h^2 \left(\sum_j \frac{\partial m(x)}{\partial x_j} \delta_j \right)^2 \right] \\ &= \sqrt{\frac{n^i}{\text{Var}(Y^i)}} [\bar{Y}(x) - \mathbb{E}[\bar{Y}(x)]] + O(\sqrt{n^i} h^3) \\ &= \sqrt{\frac{n^i}{\text{Var}(Y^i)}} \sum_{Y^s \in t^i} \left[\frac{Y^s}{n^i} - \mathbb{E} \left[\frac{Y^s}{n^i} \right] \right] + o(1) \\ &= \sqrt{\frac{n^i}{\text{Var}(Y^i)}} \left[\sum_{Y^s \in t^i} \frac{Y^s}{n^i} - \mathbb{E} \left[\sum_{Y^s \in t^i} \frac{Y^s}{n^i} \right] \right] + o(1). \end{aligned}$$

We now use the Liapunov Central Limit Theorem (e.g. Lemma A.5 in Li and Racine 2007). Define $Z_{n,s} \equiv Y^s/n^i$ and $S_n = \sum_{Y^s \in t^i} Z_{n,s}$ so that $\mathbb{E}[Z_{n,s}] = \mu^i/n^i$ and $\text{Var}(Z_{n,s}) = \text{Var}(Y^s)/(n^i)^2 = \text{Var}(Y^i)/(n^i)^2$, and $\text{Var}(S_n) = \sum_{Y^s \in t^i} \text{Var}(Z_{n,s}) = \text{Var}(Y^i)/n^i$. Since Y^s is bounded we have $\mathbb{E}|Z_{n,s}|^{2+q} < \infty$ for some $q > 0$. From (16) we see that

$$\text{Var}(Y^i) = \mathbb{E}[\sigma^2(X) | X \in \mathcal{C}^i] + o(1).$$

Finally note that

$$\lim_{n^i \rightarrow \infty} \sum_{Y^s \in t^i} \mathbb{E}|Z_{n,s} - \mathbb{E}[Z_{n,s}]|^{2+q} = 0,$$

for some $q > 0$. It follows that the Liapunov Central Limit Theorem implies

$$\sqrt{\frac{n^i}{\text{Var}(Y^i)}} \left[\sum_{Y^s \in t^i} \frac{Y^s}{n^i} - \mathbb{E} \left[\sum_{Y^s \in t^i} \frac{Y^s}{n^i} \right] \right] \xrightarrow{d} N(0, 1),$$

or

$$\sqrt{n^i} \left[\sum_{Y^s \in t^i} \frac{Y^s}{n^i} - \mathbb{E} \left[\sum_{Y^s \in t^i} \frac{Y^s}{n^i} \right] \right] \xrightarrow{d} N(0, \sigma^2(x)).$$

■

7.3 Proofs: *ISE*-Optimality

Proof of Theorem 4. (a) Follows directly from the facts that

$$\mathbb{E}_{T \in \mathcal{T}(\mathcal{C}, n)} [S^{i2}] = \mathbb{E}_{T \in \mathcal{T}(\mathcal{C}, n)} \left[\frac{1}{n^i - 1} \sum_{s \in T^i} (Y^s - \bar{Y}^i)^2 \right] = \text{Var}(Y^i),$$

and

$$\mathbb{E}_{T \in \mathcal{T}(\mathcal{C}, t)} [(\bar{Y}^i - \mu^i)^2] = \frac{1}{n^i} \text{Var}(Y^i).$$

(b) First we prove the result concerning *ISE*(\mathcal{C}, t). Consider a cell $\mathcal{C}_i \in \mathcal{C}$. Suppose that $n^i \geq 1$. It can be verified that

$$(\bar{Y}^i - \mu^i)^2 = \left(\frac{1}{n^i} \sum_{s \in t^i} Y^s \right)^2 + \mu^i - 2\mu^i \frac{1}{n^i} \sum_{s \in t^i} Y^s.$$

Let $n^i \rightarrow \infty$. Since, for each cell i , $\{Y_s\}$ is an i.i.d. sequence with $\mathbb{E}[Y_s] = \mu^i$ we can use Kinchine's law of large numbers and Slutsky's lemma to conclude that $P \lim_{n^i \rightarrow \infty} (\bar{Y}^i - \mu^i)^2 = \mu^i + \mu^i - 2\mu^i \mu^i = 0$. In other words, for any $\varepsilon > 0$ and $\delta \in (0, 1)$, there is an \bar{n} such that if $n^i > \bar{n}$ then $\Pr((\bar{Y}^i - \mu^i)^2 < \varepsilon) > \delta^{1/2}$. Moreover, for any \bar{n} there is a n such that if $n > \bar{n}$ then $\Pr(N^i > \bar{n}) > \delta^{1/2}$. This implies that, for any $\varepsilon > 0$ and $\delta \in (0, 1)$, there is a n such that if $n > n$ then $\Pr((\bar{Y}^i - \mu^i)^2 < \varepsilon) > \delta$. Thus, for any $\varepsilon > 0$ and $\delta \in (0, 1)$, there is a n such that if $n > n$ then

$$\Pr \left(\sum_{i=1}^k p^i (\bar{Y}^i - \mu^i)^2 < \varepsilon \right) > \delta.$$

Since *ISE*(\mathcal{C}, t) ≥ 0 , the desired result follows.

The proof for *SISE*(\mathcal{C}, t) is similar to the proof for *ISE*(\mathcal{C}, t), using the standard result $P \lim_{n^i \rightarrow \infty} s_i^2 = \text{Var}(Y^i)$. ■

7.4 Proofs: Binary Outcomes

The proof of lemma 5 is the same as second half of the proof of lemma 2, and the proof of theorem 5 is essentially the same as the proof of theorem 7.