

Seminararbeit

Conditional tress

Christoph Molnar

Supervisor:
Stephanie MÃ¶st

1. June 2012
Department of Statistics
University of Munich



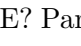
Contents

1. Introduction	1
2. Motivation	3
3. Algorithm	5
3.1. The test statistic	6
3.2. Stop criteria (1) and variable selection (2)	7
3.3. Splitting criteria (3)	7
3.4. Repeat(4)	9
4. Continuous Regression: Example bodyfat	11
5. Classification: Example glaucoma	15
6. Other scales	17
A. Computational details	19
B. Digital Appendix	21
Bibliography	23
List of Figures	25

1. Introduction

Recursive partitioning is a powerful yet simple tool in predictive and explanatory statistics. Models build by partitioning take on the form of decision trees, which makes the approach easy to understand for everyone without understanding the algorithm behind. The model partitiones the data to reconstruct the relationship

$$Y = f(X)$$

, where Y is called the response variable, which depends on a function f of the covariates matrix X .  Partitioning can be done with many different approaches and therefore the landscape of algorithms is very vivid. The differences of trees algorithms his the way trees are grown. They can be divided into those which can do regression, those which can do classification and those which can do both. Another characteristic is the number of split per partition step. There are binary splits, which divides the partition into two new partitions and multiway splits which yield more than two partitions. The variety gets big in the philosophy of how to determine which variable to take for the next step and where to split it. The point where the tree is not grown any more differs for the algorithms. Somewhere between all of those algorithms is the conditional trees framework.

2. Motivation

Recursive partitioning suffers from different problems, some of which are already solved by some approaches:

- Overfitting
- High variance
- Variable selection bias
- Heuristic approach, lack of statistical model behind
- Restriction on possible measurement scales of Y and X

Most approaches suffers at least of one of those problems. CART (Classification And Regression Trees) is a famous and widely used example of partitioning algorithms. Let us take a closer look at the problems with CART as example and how the conditional trees approach solves them.

If a tree is allowed to grow full length, pathological split could happen and if the covariate space is large enough we would end up with a tree, which contains only one observation in each terminal node. This tree would very likely be **overfitting** on the training data and deliver very bad results on new data. Approaches to avoid this problem are techniques called early stopping and pruning. Early stopping forces trees to stop growing when some criterion is not fulfilled. This criterion could be a minimum number of observations in a node. Pruning let's the tree grow at first and prunes the leafs back afterwards. The CART algorithm uses both early stopping and pruning to avoid overfitting.

As the tree strongly depends on the first splits, different variables at for the first split can yield two structurally different trees. Therefore trees (also CART) are sensitive to variance in the data and resulting trees themselves have a **high variance**.

Exhaustive search procedures as used by the CART algorithm tend to choose variables with more possible split points (**variable selection bias**). This is a problem of multiple comparison. Covariates with many possible splits are searched more often for the best split.

The next split is just a heuristic, as the algorithm only searches for the next best split (like CART). Conditional trees algorithm measures the association and uses the covariate with the strongest association with the response variable. The algorithm is embedded in a well-defined framework of hypothesis.

In the family of partitioning algorithm, the CART algorithm is one of the more powerful ones, as it can do regression as well as classification. Many other algorithm are restricted to one of the both tasks. Though, CART still lacks support for other scales of X and Y . Examples are: ordinal regression and censored data, just to name two.

To achieve the above goals, the conditional tree algorithm embeds all decisions into statistical hypothesis tests. The tests are based on conditional inference, i.e. permutation tests. The test statistic used is based on the work of [STRASSER AND WEBER].

3. Algorithm

All decisions are embedded into hypothesis tests. The conditional trees algorithm uses permutation tests to test the hypothesis of independence between a covariate and the response. This will be described further in the single steps of the algorithm.

The algorithm:

Permutation test related steps are written in red.

1. Stop criterion

- Test global null hypothesis H_0 of independence between Y and all X_j with $H_0 = \cap_{j=1}^m H_0^j$ and $H_0^j : D(\mathbf{Y}|X_j) = D(\mathbf{Y})$ (permutation tests for each X_j)
- If H_0 not rejected (no significance for all X_j) \Rightarrow Stop

2. Variable selection

Select covariate X_{j^*} with strongest association (smallest p-value)

3. Best split point search

Search best split for X_{j^*} (max. test statistic c) and partition data

4. Repeat

Repeat steps 1.), 2.) and 3.) for both of the new partitions

The algorithm starts with the whole data set and tests if it should be splitted. If the answer is positive the variable with the strongest association with the response is chosen and the data set will be split into two partitions. The steps will be repeated within both of the new partitions. Covariates chosen for a split can be chosen again later (only in the case of a bivariate covariate it doesn't make sense). First if in all partitions the global null hypothesis of independence cannot be rejected (Stop criterion) the tree does not grow any further and the algorithm stops.

The next Sections describe in detail how the single steps work and especially how permutation tests are applied.

3.1. The test statistic

All decisions of the algorithm are embeded in hypothesis tests. These are done with permutation tests (conditional inference).

PERMUTATION TESTS EXPLANATION?

Strasser and Weber [LINK] have formulated a very general test statistic, which can be used to do a permutation test if a response Y and a covariate X are independent.

$$\mathbf{T}_j(L_n, w) = \text{vec} \left(\sum_{i=1}^n w_i g_j(X_{ij}) h(Y_i, (Y_1, \dots, Y_n))^T \right) \in \mathbb{R}^{p_j q}$$

It may look difficult in the first place, but it can be broken down:

- $\text{vec}()$ The core of the test statistic can be a matrix. In this case $\text{vec}()$ - Operator vectorizes the matrix
- \sum The test statistic is a sum over all observations
- w I lied: Not all observations, because observations which are not in the current partition will get the weight $w = 0$ and otherwise $w = 1$. This ensures us, that only the data in the current node is in focus.
- g_j A transformation of the j -th covariate X_j . Transformation depends on scale of the covariate
- h Influence function. Transformation of the response Y .

The test statistic has an expectation and variance:

$$\begin{aligned} \mu_j &= \mathbb{E}(\mathbf{T}_j(\mathcal{L}_n, \mathbf{w}) | S(\mathcal{L}_n, \mathbf{w})) = \text{vec} \left(\left(\sum_{i=1}^n w_i g_j(X_{ji}) \right) \mathbb{E}(h | S(\mathcal{L}_n, \mathbf{w}))^T \right) \\ \Sigma_j &= \mathbb{V}(\mathbf{T}_j(\mathcal{L}_n, \mathbf{w}) | S(\mathcal{L}_n, \mathbf{w})) \\ &= \frac{\mathbf{w} \cdot}{\mathbf{w} \cdot - 1} \mathbb{V}(h | S(\mathcal{L}_n, \mathbf{w})) \otimes \left(\sum_i w_i g_j(X_{ji}) \otimes w_i g_j(X_{ji})^T \right) \\ &\quad - \frac{1}{\mathbf{w} \cdot - 1} \mathbb{V}(h | S(\mathcal{L}_n, \mathbf{w})) \otimes \left(\sum_i w_i g_j(X_{ji}) \right) \otimes \left(\sum_i w_i g_j(X_{ji}) \right)^T \\ \mathbf{w} \cdot &= \sum_{i=1}^n w_i \end{aligned}$$

$$\begin{aligned}\mathbb{E}(h|S(\mathcal{L}_n, \mathbf{w})) &= \mathbf{w}^{-1} \sum_i w_i h(\mathbf{Y}_i, (\mathbf{Y}_1, \dots, \mathbf{Y}_n)) \in \mathbb{R}^q \\ \mathbb{V}(h|S(\mathcal{L}_n, \mathbf{w})) &= \mathbf{w}^{-1} \sum_i w_i (h(\mathbf{Y}_i, (\mathbf{Y}_1, \dots, \mathbf{Y}_n)) - \mathbb{E}(h|S(\mathcal{L}_n, \mathbf{w}))) \\ &\quad (h(\mathbf{Y}_i, (\mathbf{Y}_1, \dots, \mathbf{Y}_n)) - \mathbb{E}(h|S(\mathcal{L}_n, \mathbf{w})))^T\end{aligned}$$

Thus we can standardize the linear tests statistic: $c(\mathbf{t}, \mu, \Sigma) = \max_{k=1, \dots, pq} \left| \frac{(\mathbf{t} - \mu)_k}{\sqrt{(\Sigma)_{kk}}} \right|$

CONVERGENCE BLA FROM STRASSER PAPER??

SHORT ABOUT PERMUTATION TESTS

3.2. Stop criteria (1) and variable selection (2)

In every partition the first step of the algorithm asks: “Split at all?”. This is formulated in a proper Null Hypothesis, the global null hypothesis. Note that global means here global in this partition and not the whole tree.

$H_0 :=$ The response Y is independent from all covariates X_j , $j \in 1, \dots, m$. The hypothesis is a joint hypothesis:

$$H_0 = \cap_{j=1}^m H_0^j \text{ and } H_0^j : D(\mathbf{Y}|X_j) = D(\mathbf{Y})$$

Each hypothesis H_0^j is tested seperatly. Most simple approach would be to look at the m resulting p-values and to rejected the global null hypothesis of independence if one the p-values exceeds the predetermined significance level α (e.g. $\alpha = 0.05$). Though multiple comparison has to be considered and any multiple testing procedure can be used at this part of the algorithm to determine t he if H_0 can be rejected.

RESULT P-VALUE

AGAIN HYPOTHESIS, REJECTED - NOT REJECTED

VARIABLE SELECTION

3.3. Splitting criteria (3)

Steps 1) and 2) of the algorithm are completed. The covariate with the strongest association is chosen for the next partition step. Every covariate (which is not binary) has more than one possible split. To determine where to split, a criterion which measures the goodness of

the split has to be applied. The CART algorithm uses Gini for classification and sum of squares for regression. Both of the criteria could be used by the Conditional Tree algorithm as well, but the approach is different. Because of the different types of possible regression-/classification - models (categorical, ordinal, numeric, censored, ...) a more general approach is suitable. Again the test statistic framework from Strasser and Weber [CITE] can be used. A special linear test statistic, of the same kind, which is used for the stop criterion and variable selection, can be used. The formula is:

$$T_j^A(L_n, w) = vec \left(\sum_{i=1}^n w_i I(X_{ji} \in A) \cdot h(Y_i, (Y_1, \dots, Y_n))^T \right)$$

The difference to the test statistic for the association test is the transformation of X . We only look at the different partitions of X . Therefore the scale of X is not of any interest anymore, but the partition which emerges by a certain split point. An appropriate function to capture only the difference in the partition, the transformation of X is the indicator function. It is defined as: $I(X_{ji} \in A) = \begin{cases} 1, & X_{ij} \in A \\ 0 & X_{ij} \notin A \end{cases}$, where A is one possible partition. This results in the statistic

$$c_{max}(\mathbf{t}, \mu, \Sigma) = max_k \left| \frac{(\mathbf{t}^A - \mu)_k}{\sqrt{(\Sigma)_{kk}}} \right|$$

For regression and the identity function for the influence function h , the test statistic is the following:

$$c_{max}(\mathbf{t}, \mu, \Sigma) = \max_k \left| \frac{(\mathbf{t}^A - \mu)_k}{\sqrt{(\Sigma)_{kk}}} \right| \quad (3.1)$$

$$= \left| \frac{\sum_{i=1}^n w_i I(X_{ji} \in A) \cdot Y_i - \sum_{i=1}^n w_i I(X_{ji} \in A) \cdot n_{node}^{-1} \sum_{i=1}^n w_i Y_i}{\sqrt{\frac{n_{node}}{n_{node}-1} \frac{1}{n_{node}} \sum_{i=1}^n w_i (Y_i - \bar{Y}_{node})^2 \cdot \sum_{i=1}^n w_i I(X_{ji} \in A)^2 - \frac{1}{n_{node}-1} \frac{1}{n_{node}} \sum_{i=1}^n w_i (Y_i - \bar{Y}_{node})^2}} \right| \quad (3.2)$$

$$= \left| \frac{\sum_{i: x_{ij} \in A} Y_i - n_A \bar{Y}_{node}}{\sqrt{\frac{1}{n_{node}-1} \sum_{i=1}^n (Y_i - \bar{Y}_{node})^2 n_A (1 - \frac{n_A}{n_{node}})}} \right| \quad (3.3)$$

$$= n_A \left| \frac{\bar{Y}_A - \bar{Y}_{node}}{\sqrt{\frac{1}{n_{node}-1} \sum_{i \in node} (Y_i - \bar{Y}_{node})^2 \cdot n_{node} (\frac{n_A}{n_{node}}) (1 - (\frac{n_A}{n_{node}}))}} \right| \quad (3.4)$$

$$= n_A \left| \frac{\bar{Y}_A - \bar{Y}_{node}}{\sqrt{Var(Y) \cdot Var(Z)}} \right| \quad (3.5)$$

$$(3.6)$$

with $Z \sim B(n_{node}, \pi = \frac{n_A}{n_{node}})$ Can be interpreted as the probability that z observations would be assigned to A if the process of assigning would be random with probability $\frac{n_A}{n_{node}}$. Is maximal for $n_A = \frac{n_{node}}{2}$. The closer $\frac{n_A}{n_{node}}$ to 0.5 the bigger is c (assuming Y_A stays the same). Thus the test statistic favors bigger partitions. $n_{node} := \sum_{i=1}^n w_i$ n_{node} : Number of observation in node \bar{Y}_{node} : Mean of Y in node n_A : Number of observations in partition A \bar{Y}_A : Mean of Y in A

Additional stopping criteria like stopping when the resulting partitions would become too small can be implemented by restricting the searched split points.

3.4. Repeat(4)

[Is repeated; choice of alpha;]

4. Continuous Regression: Example bodyfat

The first example is a continuous regression model, where both the response and the covariates are measured on a numeric scale. The model is illustrated with the *bodyfat* available in the mboost LINK package.

The data set

The data set contains observations of 71 healthy women. The measurements contain body fat, which is measured by DXA (Dual-energy X-ray absorptiometry), a method to determine the amount of body fat. Other variables in the data set are anthropometric measurements like the breadth of the knee, the waist circumference, the hip circumference etc.. The objective is to predict the body fat with the anthropometric measurements, because the DXA method is more expensive and not always available.

The test statistic

Bodyfat measured by DXA as well as body measurements are numeric variables. Thus one possible choice for the influence function h and the transformation function g_j , $\forall j \in 1, \dots, m$ is the identity function, which means the variables will not be transformed at all. Thus:

$$h = \mathbf{Y}_i$$

$$g = \mathbf{X}_i$$

This yields following not-standardized test statistic:

$$\mathbf{T}_j(\mathcal{L}_n, \mathbf{w}) = \sum_{i=1}^n w_i \mathbf{X}_{ij} \mathbf{Y}_i = \sum_{i:node} \mathbf{X}_{ij} \mathbf{Y}_i$$

The next step is to standardize the test statistic:

$$c_{max}(\mathbf{t}, \mu, \Sigma) = \max_{k=1, \dots, pq} \left| \frac{(t - \mu)_k}{\sqrt{(\Sigma)_{kk}}} \right| = \left| \frac{t - \mu}{\sqrt{\Sigma}} \right|$$

With LINK TO MU AND SIGMA FORMULA:

$$\mu_j = \sum_{i=1}^n X_{ij} \mathbb{E}(h|S) = n \cdot \bar{X}_j \bar{Y} \quad (4.1)$$

$$\Sigma = \frac{n_{node}}{n_{node} - 1} \mathbb{V}(h) \cdot \sum_{i=1}^n X_{ij}^2 - \frac{1}{n_{node} - 1} \mathbb{V}(h) n_{node}^2 \bar{X}_j^2 \quad (4.2)$$

$$= \frac{1}{n_{node}} \sum_{i=1}^n (Y_i - \bar{Y})^2 \sum_{i=1}^n X_{ij}^2 - \frac{1}{n_{node} - 1} \cdot \frac{1}{n_{node}} \sum_{i=1}^n (Y_i - \bar{Y})^2 n_{node}^2 \bar{X}_j^2 \quad (4.3)$$

$$= \frac{1}{n_{node} - 1} \sum_{i=1}^n (Y_i - \bar{Y})^2 \left(\sum_{i=1}^n X_{ij}^2 - n_{node} \bar{X}_j^2 \right) \quad (4.4)$$

$$= \frac{1}{n_{node} - 1} \left(\sum_{i=1}^n (Y_i - \bar{Y})^2 \right) \left(\sum_{i=1}^n (X_{ij} - \bar{X}_j)^2 \right) \quad (4.5)$$

Therefore:

$$c \propto \left| \frac{\sum_{i:node} X_{ij} Y_i - n_{node} \bar{X}_j \bar{Y}}{\sqrt{\left(\sum_{i:node} (Y_i - \bar{Y})^2 \right) \left(\sum_{i:node} (X_{ij} - \bar{X}_j)^2 \right)}} \right|$$

The linear test statistic is proportional to the pearson correlation coefficient. This means, that the permutation test tests if the correlation between Y and any X_j is different than 0. Thus by choosing the identity function for h and g_j the null hypothesis of independence between Y and X_j is formulated as “The correlation between Y and X_j is zero”.

The next step is to calculate the test statistic (the pearson correlation coefficient multiplied with a constant) for the observation in the current partition (where $w \neq 0$). The response of the observations will be permuted and the test statistic calculated again. This will be done often enough to approximate the distribution of the test statistic for the sample. If the correlation coefficient of the original data is very extreme compared to the permuted test statistics, the p-value will be very low.

The procedure of calculating the test statistic for the original data and the permutations is done for every covariate X_j , $j \in 1, \dots, m$ separately.

EXAMPLE FIRST SPLIT? R-CODE

For regression and the identity function for the influence function h , the test statistic is

the following:

$$c_{max}(\mathbf{t}, \mu, \Sigma) = \max_k \left| \frac{(\mathbf{t}^A - \mu)_k}{\sqrt{(\Sigma)_{kk}}} \right| \quad (4.6)$$

$$= \left| \frac{\sum_{i=1}^n w_i I(X_{ji} \in A) \cdot Y_i - \sum_{i=1}^n w_i I(X_{ji} \in A) \cdot n_{node}^{-1} \sum_{i=1}^n w_i Y_i}{\sqrt{\frac{n_{node}}{n_{node}-1} \frac{1}{n_{node}} \sum_{i=1}^n w_i (Y_i - \bar{Y}_{node})^2 \cdot \sum_{i=1}^n w_i I(X_{ji} \in A)^2 - \frac{1}{n_{node}-1} \frac{1}{n_{node}} \sum_{i=1}^n w_i (Y_i - \bar{Y}_{node})^2}} \right| \quad (4.7)$$

$$= \left| \frac{\sum_{i: x_{ij} \in A} Y_i - n_A \bar{Y}_{node}}{\sqrt{\frac{1}{n_{node}-1} \sum_{i=1}^n (Y_i - \bar{Y}_{node})^2 n_A (1 - \frac{n_A}{n_{node}})}} \right| \quad (4.8)$$

$$= n_A \left| \frac{\bar{Y}_A - \bar{Y}_{node}}{\sqrt{\frac{1}{n_{node}-1} \sum_{i \in node} (Y_i - \bar{Y}_{node})^2 \cdot n_{node} (\frac{n_A}{n_{node}}) (1 - (\frac{n_A}{n_{node}}))}} \right| \quad (4.9)$$

$$= n_A \left| \frac{\bar{Y}_A - \bar{Y}_{node}}{\sqrt{Var(Y) \cdot Var(Z)}} \right| \quad (4.10)$$

$$(4.11)$$

with $Z \sim B(n_{node}, \pi = \frac{n_A}{n_{node}})$ Can be interpreted as the probability that z observations would be assigned to A if the process of assigning would be random with probability $\frac{n_A}{n_{node}}$. Is maximal for $n_A = \frac{n_{node}}{2}$. The closer $\frac{n_A}{n_{node}}$ to 0.5 the bigger is c (assuming Y_A stays the same). Thus the test statistic favors bigger partitions. $n_{node} := \sum_{i=1}^n w_i$ n_{node} : Number of observation in node \bar{Y}_{node} : Mean of Y in node n_A : Number of observations in partition A \bar{Y}_A : Mean of Y in A

Additional stopping criteria like stopping when the resulting partitions would become to small can be implemented by restricting the searched split points.

5. Classification: Example glaucoma

6. Other scales

A. Computational details

B. Digital Appendix

Bibliography

- J. Friedman, T. Hastie, and R. Tibshirani. *The elements of statistical learning*, volume 1. Springer Series in Statistics, 2001.
- T. Hothorn, K. Hornik, and A. Zeileis. Unbiased recursive partitioning: A conditional inference framework. *Journal of Computational and Graphical Statistics*, 15(3):651–674, 2006.
- H. Strasser and C. Weber. On the asymptotic theory of permutation statistics. 1999.

List of Figures