

ggRandomForests: Exploring Random Forest Survival

John Ehrlinger

Microsoft

Abstract

Random forest (Leo Breiman 2001a) (RF) is a non-parametric statistical method requiring no distributional assumptions on covariate relation to the response. RF is a robust, nonlinear technique that optimizes predictive accuracy by fitting an ensemble of trees to stabilize model estimates. Random survival forests (RSF) (Ishwaran and Kogalur 2007; Ishwaran et al. 2008) are an extension of Breimans RF techniques allowing efficient non-parametric analysis of time to event data. The **randomForestSRC** package (Ishwaran and Kogalur 2014) is a unified treatment of Breimans random forest for survival, regression and classification problems. Predictive accuracy makes RF an attractive alternative to parametric models, though complexity and interpretability of the forest hinder wider application of the method. We introduce the **ggRandomForests** package, tools for visually understand random forest models grown in R (R Core Team 2014) with the **randomForestSRC** package. The **ggRandomForests** package is structured to extract intermediate data objects from **randomForestSRC** objects and generate figures using the **ggplot2** (Wickham 2009) graphics package. This document is structured as a tutorial for building random forest for survival with the **randomForestSRC** package and using the **ggRandomForests** package for investigating how the forest is constructed. We analyse the Primary Biliary Cirrhosis of the liver data from a clinical trial at the Mayo Clinic (Fleming and Harrington 1991). We demonstrate random forest variable selection using Variable Importance (VIMP) (Leo Breiman 2001a) and Minimal Depth (Ishwaran et al. 2010), a property derived from the construction of each tree within the forest. We will also demonstrate the use of variable dependence and partial dependence plots (Friedman 2000) to aid in the interpretation of RSF results. We then examine variable interactions between covariates using conditional variable dependence plots. Our aim is to demonstrate the strength of using Random Forest methods for both prediction and information retrieval, specifically in time to event data settings.

Keywords: random forest, survival, vimp, minimal depth, R, **randomForestSRC**, **ggRandomForests**, **randomForest**.

1. Introduction

Random forest (Leo Breiman 2001a) (RF) is a non-parametric statistical method which requires no distributional assumptions on covariate relation to the response. RF is a robust, nonlinear technique that optimizes predictive accuracy by fitting an ensemble of trees to stabilize model estimates. Random Survival Forest (RSF) (Ishwaran and Kogalur 2007; Ishwaran et al. 2008) is an extension of Breiman's RF techniques to survival settings, allowing efficient non-parametric analysis of time to event data. The **randomForestSRC** package (<http://CRAN.R-project.org/package=randomForestSRC>) (Ishwaran and Kogalur 2014) is

a unified treatment of Breiman’s random forest for survival, regression and classification problems.

Predictive accuracy make RF an attractive alternative to parametric models, though complexity and interpretability of the forest hinder wider application of the method. We introduce the **ggRandomForests** package (<http://CRAN.R-project.org/package=ggRandomForests>) for visually exploring random forest models. The **ggRandomForests** package is structured to extract intermediate data objects from **randomForestSRC** objects and generate figures using the **ggplot2** graphics package (<http://CRAN.R-project.org/package=ggplot2>) (Wickham 2009).

Many of the figures created by the **ggRandomForests** package are also available directly from within the **randomForestSRC** package. However **ggRandomForests** offers the following advantages:

- Separation of data and figures: **ggRandomForests** contains functions that operate on either the **rfsrc** forest object directly, or on the output from **randomForestSRC** post processing functions (i.e., `plot.variable`, `var.select`) to generate intermediate **ggRandomForests** data objects. **ggRandomForests** functions are provide to further process these objects and plot results using the **ggplot2** graphics package. Alternatively, users can use these data objects for their own custom plotting or analysis operations.
- Each data object/figure is a single, self contained unit. This allows simple modification and manipulation of the data or **ggplot** objects to meet users specific needs and requirements.
- We chose to use the **ggplot2** package for our figures for flexibility in modifying the output. Each **ggRandomForests** plot function returns either a single **ggplot** object, or a **list** of **ggplot** objects, allowing the use of additional **ggplot2** functions to modify and customize the final figures.

This document is structured as a tutorial for using the **randomForestSRC** package for building and post-processing random survival forest models and using the **ggRandomForests** package for understanding how the forest is constructed. In this tutorial, we will build a random survival forest for the primary biliary cirrhosis (PBC) of the liver data set (Fleming and Harrington 1991), available in the **randomForestSRC** package.

In Section ?? we introduce the **pbc** data set and summarize the proportional hazards analysis of this data from Chapter 4 of (Fleming and Harrington 1991). In Section ??, we describe how to grow a random survival forest with the **randomForestSRC** package. Random forest is not a parsimonious method, but uses all variables available in the data set to construct the response predictor. We demonstrate random forest variable selection techniques (Section ??) using Variable Importance (VIMP) (Leo Breiman 2001a) in Section ?? and Minimal Depth (Ishwaran et al. 2010) in Section ?. We then compare both methods with variables used in the (Fleming and Harrington 1991) model.

Once we have an idea of which variables we are most interested in, we use dependence plots (Friedman 2000) (Section ??) to understand how these variables are related to the response. Variable dependence (Section ??) plots give us an idea of the overall trend of a variable/response relation, while partial dependence plots (Section ??) show us the risk adjusted relation by averaging out the effects of other variables. Dependence plots often show

strongly non-linear variable/response relations that are not easily obtained through parametric modeling.

We then graphically examine forest variable interactions with the use of variable and partial dependence conditioning plots (coplots) (Chambers 1992; Cleveland 1993) (Section ??) and close with concluding remarks in Section ??.

2. Data summary: primary biliary cirrhosis (PBC) data set

The *primary biliary cirrhosis* of the liver (PBC) study consists of 424 PBC patients referred to Mayo Clinic between 1974 and 1984 who met eligibility criteria for a randomized placebo controlled trial of the drug D-penicillamine (DPCA). The data is described in (?, Chapter 0.2) and a partial likelihood model (Cox proportional hazards) is developed in Chapter 4.4. The `pb` data set, included in the `randomForestSRC` package, contains 418 observations, of which 312 patients participated in the randomized trial (?, Appendix D).

```
R> data("pb", package = "randomForestSRC")
```

For this analysis, we modify some of the data for better formatting of our results. Since the data contains about 12 years of follow up, we prefer using `years` instead of `days` to describe survival. We also convert the `age` variable to years, and the `treatment` variable to a factor containing levels of `c("DPCA", "placebo")`. The variable names, type and description are given in Table ??.

2.1. Exploratory data analysis

It is good practice to view your data before beginning analysis. Exploratory Data Analysis (EDA) (Tukey 1977) will help you to understand the data, and find outliers, missing values and other data anomalies within each variable before getting deep into the analysis. To this end, we use `ggplot2` figures with the `facet_wrap` function to create two sets of panel plots, one of histograms for categorical variables (Figure 1), and another of scatter plots for continuous variables (Figure 2). Variables are plotted along a continuous variable on the X-axis to separate the individual observations.

In categorical EDA plots (Figure 1), we are looking for patterns of missing data (white portion of bars). We often use surgical date for our X-axis variable to look for possible periods of low enrollment. There is not a comparable variable available in the `pb` data set, so instead we used follow up time (`years`). Another reasonable choice may have been to use the patient `age` variable for the X-axis. The important quality of the selected variable is to spread the observations out to aid in finding data anomalies.

In continuous data EDA plots (Figure 2), we are looking for missingness (rug marks) and extreme or non-physical values. For survival settings, we color and shape the points as red 'x's to indicate events, and blue circles to indicate censored observation.

Extreme value examples are evident in a few of the variables in Figure 2. We are typically looking for values that are outside of the biological range. This is often caused by measurements recorded in differing units, which can sometimes be corrected algorithmically. Since we can not ask the original investigator to clarify these values in this particular study, we will continue without modifying the data.

Table 1: ‘pbc’ data set variable dictionary.

Variable name	Description	Type
years	Time (years)	numeric
status	Event (F = censor, T = death)	logical
treatment	Treatment (DPCA, Placebo)	factor
age	Age (years)	numeric
sex	Female = T	logical
ascites	Presence of Ascites	logical
hepatom	Presence of Hepatomegaly	logical
spiders	Presence of Spiders	logical
edema	Edema (0, 0.5, 1)	factor
bili	Serum Bilirubin (mg/dl)	numeric
chol	Serum Cholesterol (mg/dl)	integer
albumin	Albumin (gm/dl)	numeric
copper	Urine Copper (ug/day)	integer
alk	Alkaline Phosphatase (U/liter)	numeric
sgot	SGOT (U/ml)	numeric
trig	Triglycerides (mg/dl)	integer
platelet	Platelets per cubic ml/1000	integer
prothrombin	Prothrombin time (sec)	numeric
stage	Histologic Stage	factor

Table 2: Missing value counts in ‘pbc’ data set and pbc clinical trial observations (‘pbc.trial’).

	pbc	pbc.trial
treatment	106	0
ascites	106	0
hepatom	106	0
spiders	106	0
chol	134	28
copper	108	2
alk	106	0
sgot	106	0
trig	136	30
platelet	11	4
prothrombin	2	0
stage	6	0

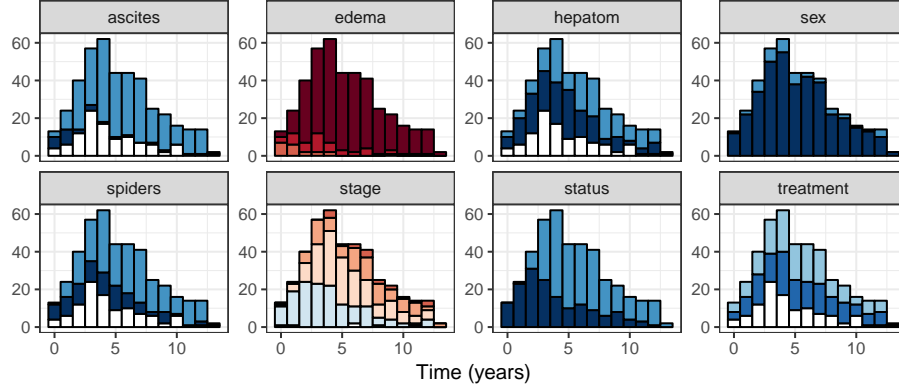


Figure 1: EDA plots for categorical variables (logicals and factors). Bars indicate number of patients within 1 year of followup interval for each categorical variable. Colors correspond to class membership within each variable. Missing values are included in white.

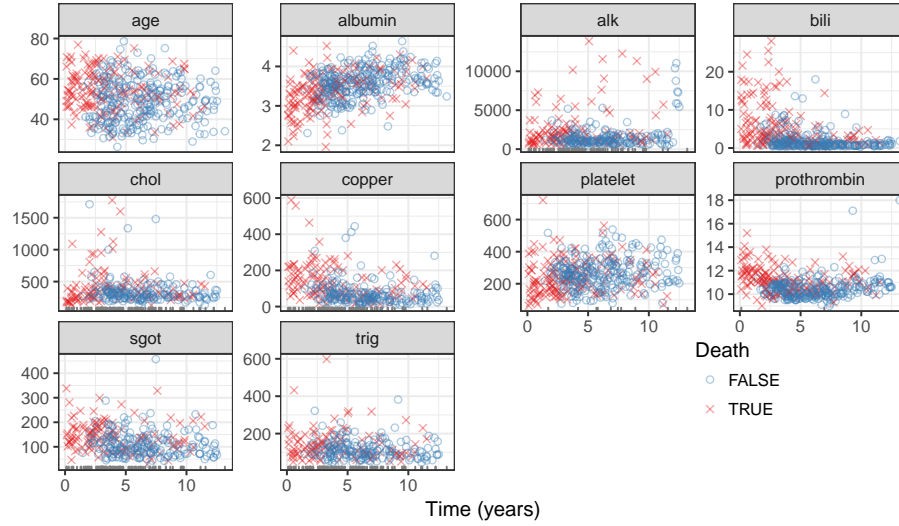


Figure 2: EDA plots for continuous variables. Symbols indicate observations with variable value on Y-axis against follow up time in years. Symbols are colored and shaped according to the death event ('status' variable). Missing values are indicated by rug marks along the X-axis

Both EDA figures indicate the `pub` data set contains quite a bit of missing data. Table ?? shows the number of missing values in each variable of the `pub` data set. Of the 19 variables in the data, 12 have missing values. The `pub` column details variables with missing data in the full `pub` data set, though there are patients that were not randomized into the trial. If we restrict the data to the trial only, most of the missing values are also removed, leaving only 4 variables with missing values. Therefore, we will focus on the 312 observations from the clinical trial for the remainder of this document. We will discuss how `randomForestSRC` handles missing values in Section ??.

2.2. PBC Model Summary

We conclude the data set investigation with a summary of (Fleming and Harrington 1991) model results from Chapter 4.4. We start by generating Kaplan–Meier (KM) survival estimates comparing the treatment groups of DPCA and placebo. We use the **ggRandomForests** `gg_survival` function to generate these estimates from the data set as follows.

```
R> # Create the trial and test data sets.
R> pbc.trial <- pbc %>% filter(!is.na(treatment))
R> pbc.test <- pbc %>% filter(is.na(treatment))
R>
R> # Create the gg_survival object
R> gg_dta <- gg_survival(interval = "years",
R+                      censor = "status",
R+                      by = "treatment",
R+                      data = pbc.trial,
R+                      conf.int = 0.95)
```

The code block reduces the `pbc` data set to the `pbc.trial` which only include observations from the clinical trial. The remaining observations are stored in the `pbc.test` data set for later use. The **ggRandomForests** package is designed to use a two step process in figure generation. The first step is data generation, where we store a `gg_survival` data object in the `gg_dta` object. The `gg_survival` function uses the `data` set, follow up `interval`, `censor` indicator and an optional grouping argument (`by`). By default `gg_survival` also calculates 95% confidence band, which we can control with the `conf.int` argument.

In the figure generation step, we use the **ggRandomForests** plot routine `plot.gg_survival` as shown in the following code block. The `plot.gg_survival` function uses the `gg_dta` data object to plot the survival estimate curves for each group and corresponding confidence interval ribbons. We have used additional **ggplot2** commands to modify the axis and legend labels (`labs`), the legend location (`theme`) and control the plot range of the y-axis (`coord_cartesian`) for this figure.

```
R> plot(gg_dta) +
R+   labs(y = "Survival Probability", x = "Observation Time (years)",
R+        color = "Treatment", fill = "Treatment") +
R+   theme(legend.position = c(0.2, 0.2)) +
R+   coord_cartesian(y = c(0, 1.01))
```

The `gg_survival` plot of Figure 3 is analogous to (Fleming and Harrington 1991) Figure 0.2.3 and Figure 4.4.1, showing there is little difference between the treatment and control groups.

The `gg_survival` function generates a variety of time-to-event estimates, including the cumulative hazard. The follow code block creates a cumulative hazard plot (Figure 0.2.1) in Figure 4 using the same data object generated by the original `gg_survival` function call. The red DPCA line is directly comparable to Figure 0.2.1, we've add the cumulative hazard estimates for the placebo population in blue.

```
R> plot(gg_dta, type = "cum_haz") +
R+   labs(y = "Cumulative Hazard", x = "Observation Time (years)",
```

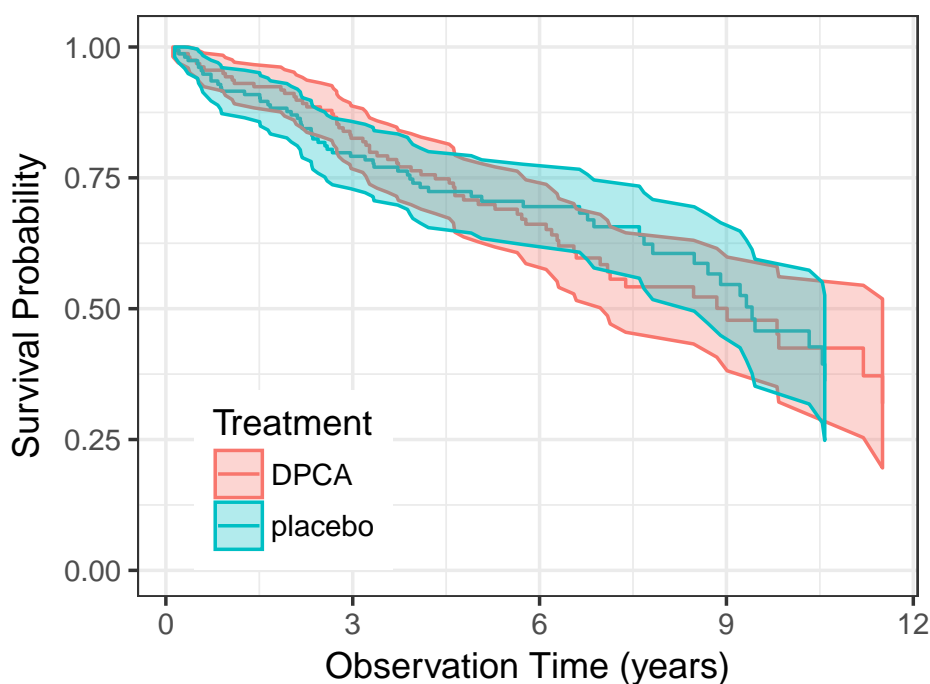


Figure 3: Kaplan–Meier survival estimates comparing the DPCA treatment (red) with placebo (blue) groups for the `pbc.trial` data set. Median survival with shaded 95% confidence band.

```
R>      color = "Treatment", fill = "Treatment") +
R> theme(legend.position = c(0.2, 0.8)) +
R> coord_cartesian(ylim = c(-0.02, 1.22))
```

In Figure 3, we demonstrated grouping on the categorical variable (`treatment`). To demonstrate plotting grouped survival on a continuous variable, we examine KM estimates of survival within stratified groups of bilirubin measures. The groupings are obtained directly from (Fleming and Harrington 1991) Figure 4.4.2, where they presented univariate model results of predicting survival on a function of bilirubin.

We set up the `bili` groups on a temporary data set (`pbc.bili`) using the `cut` function with intervals matching the reference figure. For this example we combine the data generation and plot steps into a single line of code. The `error` argument of the `plot.gg_survival` function is used to control display of the confidence bands. We suppress the intervals for this figure with `error = "none"` and again modify the plot display with `ggplot2` commands to generate Figure 5.

```
R> pbc.bili <- pbc.trial
R> pbc.bili$bili_grp <- cut(pbc.bili$bili, breaks = c(0, 0.8, 1.3, 3.4, 29))
R>
R> plot(gg_survival(interval = "years", censor = "status", by = "bili_grp",
R>      data = pbc.bili), error = "none") +
R>   labs(y = "Survival Probability", x = "Observation Time (years)",
R>      color = "Bilirubin")
```

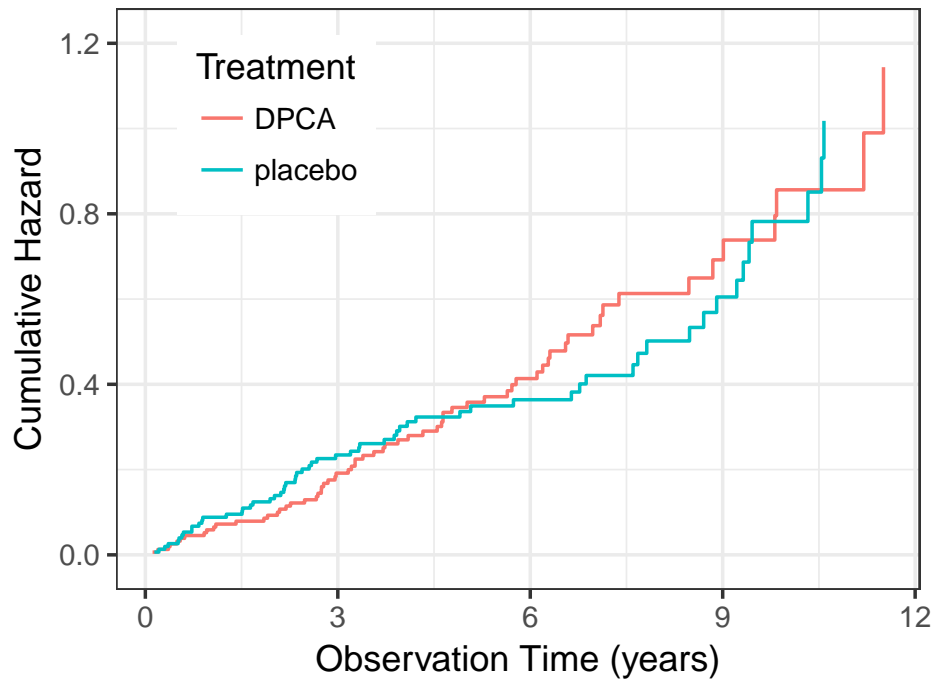


Figure 4: Kaplan–Meier cumulative hazard estimates comparing the DPCA treatment (red) with placebo (blue) groups for the pbc data set.

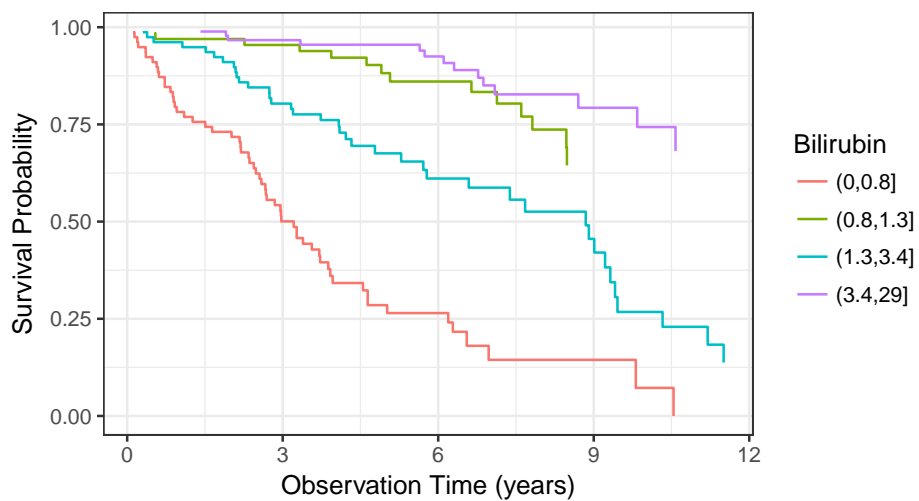


Figure 5: Kaplan–Meier survival estimates comparing different groups of Bilirubin measures (bili) for the pbc data set. Groups defined in Chapter 4 of [Fleming:1991].

In Chapter 4, (Fleming and Harrington 1991) use partial likelihood methods to build a linear model with log transformations on some variables. We summarize the final, biologically reasonable model in Table ?? for later comparison with our random forest results.

Table 3: ‘pbc’ proportional hazards model summary of 312 randomized cases in ‘pbc.trial’ data set. (Table 4.4.3c [fleming:1991])

	Coef.	Std. Err.	Z stat.
Age	0.033	0.009	3.84
log(Albumin)	-3.055	0.724	-4.22
log(Bilirubin)	0.879	0.099	8.90
Edema	0.785	0.299	2.62
log(Prothrombin Time)	3.016	1.024	2.95

3. Random survival forest

A Random Forest (Leo Breiman 2001a) is grown by *bagging* (L Breiman 1996a) a collection of *classification and regression trees* (CART) (Breiman et al. 1984). The method uses a set of B *bootstrap* (Efron and Tibshirani 1994) samples, growing an independent tree model on each sub-sample of the population. Each tree is grown by recursively partitioning the population based on optimization of a *split rule* over the p -dimensional covariate space. At each split, a subset of $m \leq p$ candidate variables are tested for the split rule optimization, dividing each node into two daughter nodes. Each daughter node is then split again until the process reaches the *stopping criteria* of either *node purity* or *node member size*, which defines the set of *terminal (unsplit) nodes* for the tree. In regression trees, node impurity is measured by mean squared error, whereas in classification problems, the Gini index is used (Friedman 2000).

Random forest sorts each training set observation into one unique terminal node per tree. Tree estimates for each observation are constructed at each terminal node, among the terminal node members. The Random Forest estimate for each observation is then calculated by aggregating, averaging (regression) or votes (classification), the terminal node results across the collection of B trees.

Random Survival Forest (Ishwaran 2007; Ishwaran et al. 2008) (RSF) are an extension of Random Forest to analyze right censored, time to event data. A forest of survival trees is grown using a log-rank splitting rule to select the optimal candidate variables. Survival estimate for each observation are constructed with a Kaplan–Meier (KM) estimator within each terminal node, at each event time.

Random Survival Forests adaptively discover nonlinear effects and interactions and are fully nonparametric. Averaging over many trees enables RSF to approximate complex survival functions, including non-proportional hazards, while maintaining low prediction error. (Ishwaran and Kogalur 2010) showed that RSF is uniformly consistent and that survival forests have a uniform approximating property in finite-sample settings, a property not possessed by individual survival trees.

The **randomForestSRC** `rfsrc` function call grows the forest, determining the type of forest by the response supplied in the `formula` argument. In the following code block, we grow a random forest for survival, by passing a survival (`Surv`) object to the forest. The forest uses all remaining variables in the `pbc.trial` data set to generate the RSF survival model.

```
R> rfsrc_pbc <- rfsrc(Surv(years, status) ~ ., data = pbc.trial,
```

```
R> nsplit = 10, na.action = "na.impute",
R> tree.err = TRUE, importance = TRUE)
```

The `print.rfsrc` function returns information on how the random forest was grown. Here the `family = "surv"` forest has `ntree = 1000` trees (the default `ntree` argument). The forest selected from $\text{ceil}(\sqrt{p} = 17) = 5$ randomly selected candidate variables for splitting at each node, stopping when a terminal node contained three or fewer observations. For continuous variables, we used a random logrank split rule, which randomly selects from `nsplit = 10` split point values, instead of optimizing over all possible values.

3.1. Generalization error

One advantage of random forest is a built in generalization error estimate. Each bootstrap sample selects approximately 63.2% of the population on average. The remaining 36.8% of observations, the Out-of-Bag (L Breiman 1996b) (OOB) sample, can be used as a hold out test set for each tree. An OOB prediction error estimate can be calculated for each observation by predicting the response over the set of trees which were not trained with that particular observation. Out-of-Bag prediction error estimates have been shown to be nearly identical to n -fold cross validation estimates (Hastie, Tibshirani, and Friedman 2009). This feature of random forest allows us to obtain both model fit and validation in one pass of the algorithm. The `gg_error` function operates on the random forest (`rfsrc_pbc`) object to extract the error estimates as a function of the number of trees in the forest. The following code block first creates a `gg_error` data object, then uses the `plot.gg_error` function to create a `ggplot` object for display in a single line of code.

```
R> plot(gg_error(rfsrc_pbc))
```

The `gg_error` plot of Figure 6 demonstrates that it does not take a large number of trees to stabilize the forest prediction error estimate. However, to ensure that each variable has enough of a chance to be included in the forest prediction process, we do want to create a rather large random forest of trees.

3.2. Training Set Prediction

The `gg_rfsrc` function extracts the OOB prediction estimates from the random forest. This code block executes the data extraction and plotting in one line, since we are not interested in holding the prediction estimates for later reuse. Each of the `ggRandomForests` plot commands return `ggplot` objects, which we can also store for modification or reuse later in the analysis (`ggRFsrc` object). Note that we again use additional `ggplot2` commands to modify the display of the plot object.

```
R> ggRFsrc <- plot(gg_rfsrc(rfsrc_pbc), alpha = 0.2) +
R>   scale_color_manual(values = strCol) +
R>   theme(legend.position = "none") +
R>   labs(y = "Survival Probability", x = "Time (years)") +
R>   coord_cartesian(ylim = c(-0.01, 1.01))
R> show(ggRFsrc)
```

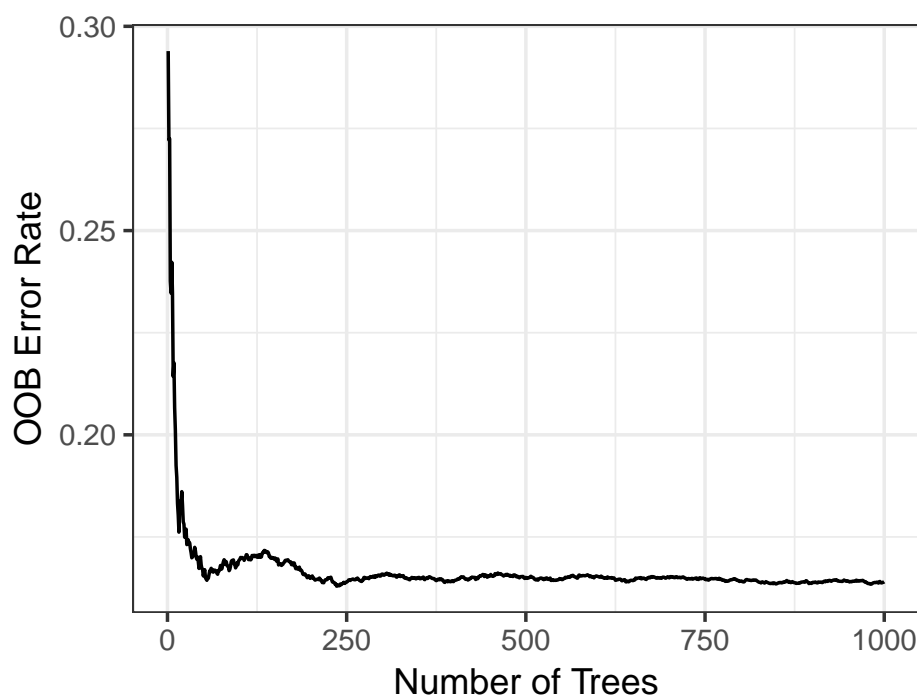


Figure 6: Random forest OOB prediction error estimates as a function of the number of trees in the forest.

The `gg_rfsrc` plot of Figure 7 shows the predicted survival from our RSF model. Each line represents a single patient in the training data set, where censored patients are colored blue, and patients who have experienced the event (death) are colored in red. We extend all predicted survival curves to the longest follow up time (12 years), regardless of the actual length of a patient's follow up time.

Interpretation of general survival properties from Figure 7 is difficult because of the number of curves displayed. To get more interpretable results, it is preferable to plot a summary of the survival results. The following code block compares the predicted survival between treatment groups, as we did in Figure 3.

```
R> plot(gg_rfsrc(rfsrc_pbc, by = "treatment")) +
R+   theme(legend.position = c(0.2, 0.2)) +
R+   labs(y = "Survival Probability", x = "Time (years)") +
R+   coord_cartesian(ylim = c(-0.01, 1.01))
```

The `gg_rfsrc` plot of Figure 8 shows the median survival with a 95% shaded confidence band for the DPCA group in red, and the placebo group in blue. When calling `gg_rfsrc` with either a `by` argument or a `conf.int` argument, the function calculates a bootstrap confidence interval around the median survival line. By default, the function will calculate the `conf.int=0.95` confidence interval, with the number of `bs.samples` equal to the number of observations.

3.3. Random forest imputation

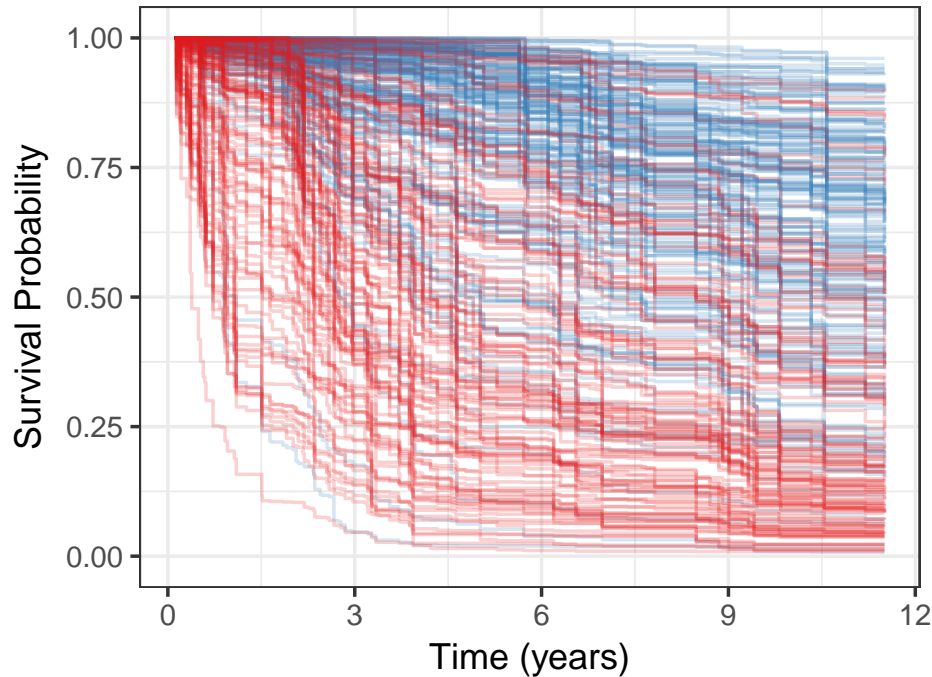


Figure 7: Random forest OOB predicted survival. Blue curves correspond to censored observations, red curves correspond to observations experiencing death events.

There are two modeling issues when dealing with missing data values: How does the algorithm build a model when values are missing from the training data?, and How does the algorithm predict a response when values are missing from the test data?. The standard procedure for linear models is to either remove or impute the missing data values before modelling. Removing the missingness is done by either removing the variable with missing values (column wise) or removing the observations (row wise). Removal is a simple solution, but may bias results when either observations or variables are scarce.

The **randomForestSRC** package imputes missing values using *adaptive tree imputation* (Ishwaran et al. 2008). Rather than impute missing values before growing the forest, the algorithm takes a ‘just--in--time’ approach. At each node split, the set of candidate variables is checked for missing values. Missing values are then imputed by randomly drawing values from non-missing data within the node. The split-statistic is then calculated on observations that were not missing values. The imputed values are used to sort observations into the subsequent daughter nodes and then discarded before the next split occurs. The process is repeated until the stopping criteria is reached and all observations are sorted into terminal nodes.

A final imputation step can be used to fill in missing values from within the terminal nodes. This step uses a process similar to the previous imputation but uses the OOB non-missing terminal node data for the random draws. These values are aggregated (averaging for continuous variables, voting for categorical variables) over the `ntree` trees in the forest to estimate an imputed data set. By default, the missing values are not filled into the training data, but are available within the forest object for later use if desired.

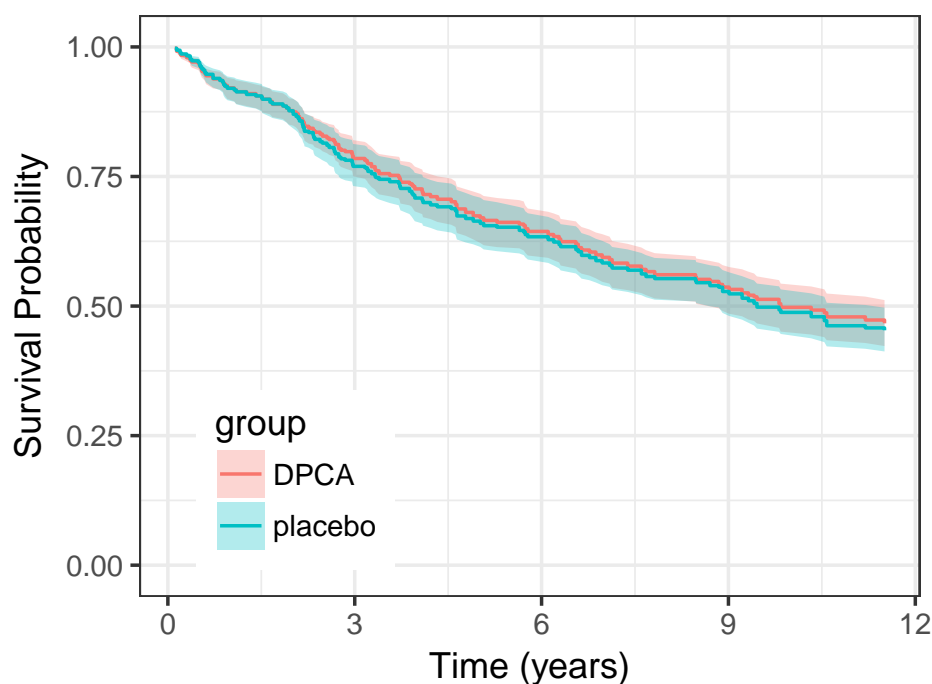


Figure 8: Random forest predicted survival stratified by treatment groups. DPCA group in red, placebo in blue with shaded 95% confidence bands.

Adaptive tree imputation still requires the missing at random assumptions (Rubin 1976). At each imputation step, the random forest assumes that similar observations are grouped together within each node. The random draws used to fill in missing data do not bias the split rule, but only sort observations similar in non-missing data into like nodes. An additional feature of this approach is the ability of predicting on test set observations with missing values.

3.4. Test set predictions

The strength of adaptive tree imputation becomes clear when doing prediction on test set observations. If we want to predict survival for patients that did not participate in the trial using the model we created in Section ??, we need to somehow account for the missing values detailed in Table ??.

The `predict.rfsrc` call takes the forest object (`rfsrc_pbc`), and the test data set (`pbc_test`) and returns a predicted survival using the same forest imputation method for missing values within the test data set (`na.action="na.impute"`).

```
R> rfsrc_pbc_test <- predict(rfsrc_pbc, newdata = pbc.test,
R+                          na.action = "na.impute",
R+                          importance = TRUE)
```

The forest summary indicates there are 106 test set observations with 36 deaths and the predicted error rate is 19.1%. We plot the predicted survival just as we did the training set

estimates.

```
R> plot(gg_rfsrc(rfsrc_pbc_test), alpha=.2) +
R+   scale_color_manual(values = strCol) +
R+   theme(legend.position = "none") +
R+   labs(y = "Survival Probability", x = "Time (years)") +
R+   coord_cartesian(ylim = c(-0.01, 1.01))
```

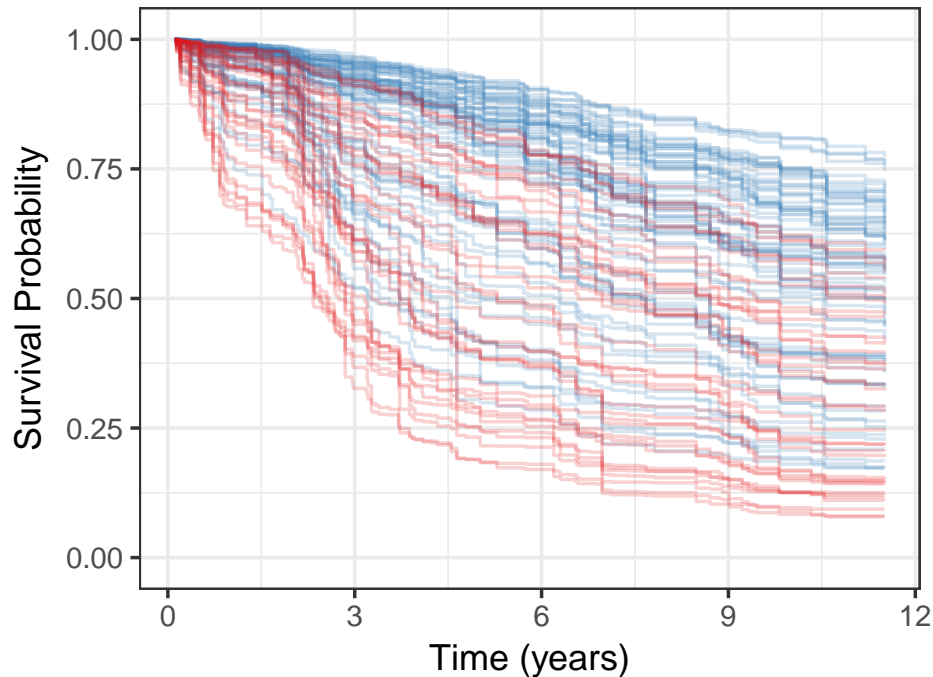


Figure 9: Random forest survival estimates for patients in the pbc.test data set. Blue curves correspond to censored patients, red curves correspond to patients experiencing a death event.

The `gg_rfsrc` plot of Figure 9 shows the test set predictions, similar to the training set predictions in Figure 7, though with fewer patients the survival curves do not cover the same area of the figure. It is important to note that because Figure 7 is constructed with OOB estimates, the survival results are comparable as estimates from unseen observations in Figure 9.

4. Variable selection

Random forest is not a parsimonious method, but uses all variables available in the data set to construct the response predictor. Also, unlike parametric models, random forest does not require the explicit specification of the functional form of covariates to the response. Therefore there is no explicit p -value/significance test for variable selection with a random forest model. Instead, RF ascertains which variables contribute to the prediction through the split rule optimization, optimally choosing variables which separate observations.

The typical goal of a random forest analysis is to build a *prediction* model, in contrast to extracting *information* regarding the underlying process (Leo Breiman 2001b). There is not usually much care given in how variables are included into the training data set. Since the goal is prediction, investigators often include the “kitchen sink” if it can help.

In contrast, in survival settings we are typically also interested in how we can possibly improve the the outcome of interest. To achieve this, for understandable inference, it is important to avoid both duplication and transformations of variables whenever possible when building our data sets. Duplication of variables, including multiple measures of a similar covariate, can reduce or mask the importance of the covariate. Transformations can also mask importance as well as make interpretation of the inference results difficult to impossible.

In this Section, We explore two separate approaches to investigate the RF variable selection process. Variable Importance (Section ??), a property related to variable misspecification, and Minimal Depth (Section ??), a property derived from the construction of the trees within the forest.

4.1. Variable Importance

Variable importance (VIMP) was originally defined in CART using a measure involving surrogate variables (see Chapter 5 of (Breiman et al. 1984)). The most popular VIMP method uses a prediction error approach involving “noising-up” each variable in turn. VIMP for a variable x_v is the difference between prediction error when x_v is randomly permuted, compared to prediction error under the observed values (Leo Breiman 2001a; Liaw and Wiener 2002; Ishwaran 2007; Ishwaran et al. 2008).

Since VIMP is the difference in OOB prediction error before and after permutation, a large VIMP value indicates that misspecification detracts from the predictive accuracy in the forest. VIMP close to zero indicates the variable contributes nothing to predictive accuracy, and negative values indicate the predictive accuracy *improves* when the variable is misspecified. In the later case, we assume noise is more informative than the true variable. As such, we ignore variables with negative and near zero values of VIMP, relying on large positive values to indicate that the predictive power of the forest is dependent on those variables.

The `gg_vimp` function extracts VIMP measures for each of the variables used to grow the forest. The `plot.gg_vimp` function shows the variables, in VIMP rank order, labeled with the named vector in the `lbls` argument.

```
R> plot(gg_vimp(rfsrc_pbc), lbls = st.labs) +
R+   theme(legend.position = c(0.8, 0.2)) +
R+   labs(fill = "VIMP > 0")
```

The `gg_vimp` plot of Figure 10 details VIMP ranking for the `pbctrial` baseline variables, from the largest (Serum Bilirubin) at the top, to smallest (Treatment (DPCA, Placebo)) at the bottom. VIMP measures are shown using bars to compare the scale of the error increase under permutation and colored by the sign of the measure (red for negative values). Note that four of the five highest ranking variables by VIMP match those selected by the (Fleming and Harrington 1991) model listed in Table ??, with urine copper (2) ranking higher than age (8). We will return to this in Section ??.

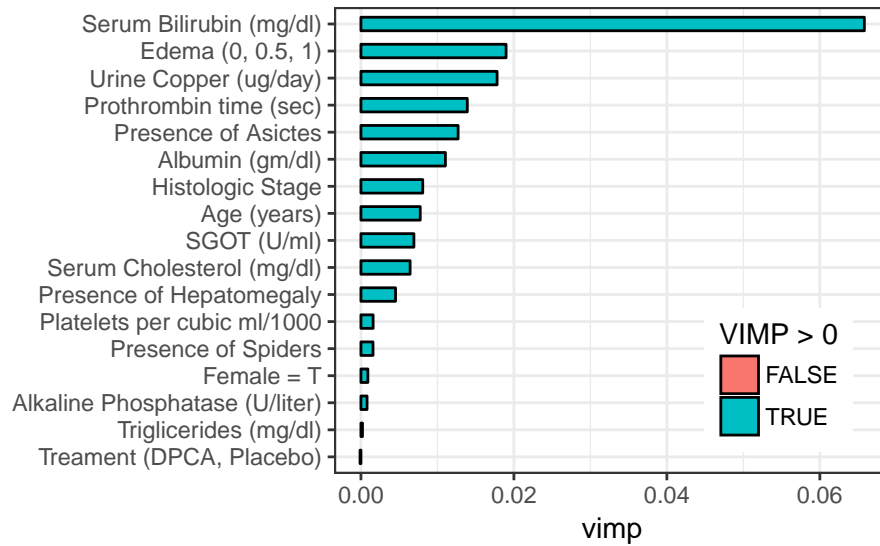


Figure 10: Random forest Variable Importance (VIMP). Blue bars indicates positive VIMP, red indicates negative VIMP. Importance is relative to positive length of bars.

4.2. Minimal Depth

In VIMP, prognostic risk factors are determined by testing the forest prediction under alternative data settings, ranking the most important variables according to their impact on predictive ability of the forest. An alternative method uses inspection of the forest construction to rank variables. *Minimal depth* (Ishwaran et al. 2010; Ishwaran et al. 2011) assumes that variables with high impact on the prediction are those that most frequently split nodes nearest to the root node, where they partition the largest samples of the population.

Within each tree, node levels are numbered based on their relative distance to the root of the tree (with the root at 0). Minimal depth measures important risk factors by averaging the depth of the first split for each variable over all trees within the forest. The assumption in the metric is that smaller minimal depth values indicate the variable separates large groups of observations, and therefore has a large impact on the forest prediction.

In general, to select variables according to VIMP, we examine the VIMP values, looking for some point along the ranking where there is a large difference in VIMP measures. Given minimal depth is a quantitative property of the forest construction, we also derive an analytic threshold for evidence of variable impact. A simple optimistic threshold rule uses the mean of the minimal depth distribution, classifying variables with minimal depth lower than this threshold as important in forest prediction.

The **randomForestSRC** `var.select` function uses the minimal depth methodology for variable selection, returning an object with both minimal depth and vimp measures. The **ggRandomForests** `gg_minimal_depth` function is analogous to the `gg_vimp` function. Variables are ranked from most important at the top (minimal depth measure), to least at the bottom (maximal minimal depth).

```
R> varsel_pbc <- var.select(rfsrc_pbc)
```


minimal depth variable selection ...

```
-----
family           : surv
var. selection   : Minimal Depth
conservativeness : medium
x-weighting used? : TRUE
dimension        : 17
sample size      : 312
ntree             : 1000
nsplit           : 10
mtry             : 5
nodesize         : 3
refitted forest  : FALSE
model size       : 13
depth threshold  : 6.7621
PE (true OOB)    : 16.402
```

Top variables:

	depth	vimp
bili	1.840	0.066
albumin	2.426	0.011
copper	2.743	0.018
prothrombin	2.837	0.014
edema	3.179	0.019
chol	3.258	0.006
platelet	3.403	0.002
age	3.661	0.008
sgot	3.723	0.007
alk	3.915	0.001
trig	4.493	0.000
stage	5.233	0.008
ascites	5.463	0.013

```
-----
```

```
R> gg_md <- gg_minimal_depth(varsel_pbc, lbls = st.labs)
R> # print(gg_md)
```

The `gg_minimal_depth` summary mostly reproduces the output from the `var.select` function from the **randomForestSRC** package. We report the minimal depth threshold (`threshold` 6.762) and the number of variables with depth below that threshold (`model size` 13). We also list a table of the top (13) selected variables, in minimal depth rank order with the associated VIMP measures. The minimal depth numbers indicate that `bili` tends to split between the first and second node level, and the next three variables (`albumin`, `copper`, `prothrombin`) split between the second and third levels on average.

```
R> plot(gg_md, lbls = st.labs)
```

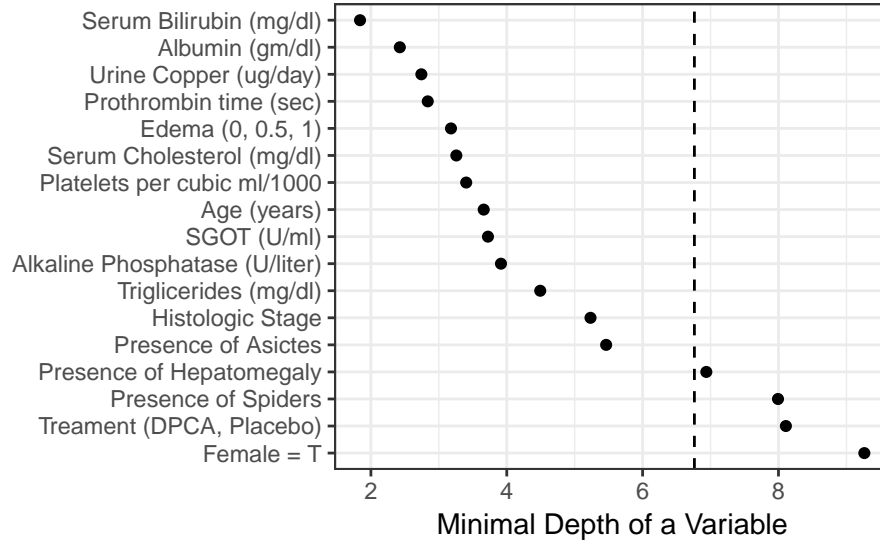


Figure 11: Minimal Depth variable selection. Low minimal depth indicates important variables. The dashed line is the threshold of maximum value for variable selection.

The `gg_minimal_depth` plot of Figure 11 is similar to the `gg_vimp` plot in Figure 10, ranking variables from most important at the top (minimal depth measure), to least at the bottom (maximal minimal depth). The vertical dashed line indicates the minimal depth threshold where smaller minimal depth values indicate higher importance and larger values indicate lower importance.

4.3. Variable selection comparison

Since the VIMP and Minimal Depth measures use different criteria, we expect the variable ranking to be somewhat different. We use `gg_minimal_vimp` function to compare rankings between minimal depth and VIMP in Figure 12.

```
R> plot(gg_minimal_vimp(gg_md), lbls = st.labs) +  
R+   theme(legend.position=c(0.8, 0.2))
```

The points along the red dashed line indicate where the measures are in agreement. Points above the red dashed line are ranked higher by VIMP than by minimal depth, indicating the variables are more sensitive to misspecification. Those below the line have a higher minimal depth ranking, indicating they are better at dividing large portions of the population. The further the points are from the line, the more the discrepancy between measures.

We examine the ranking of the different variable selection methods further in Table ?? . We can use the Z statistic from Table ?? to rank variables selected in the (Fleming and Harrington 1991) model to compare with variables selected by minimal depth and VIMP. The table is constructed by taking the top ranked minimal depth variables (below the selection threshold) and matching the VIMP ranking and (Fleming and Harrington 1991) model transforms. We

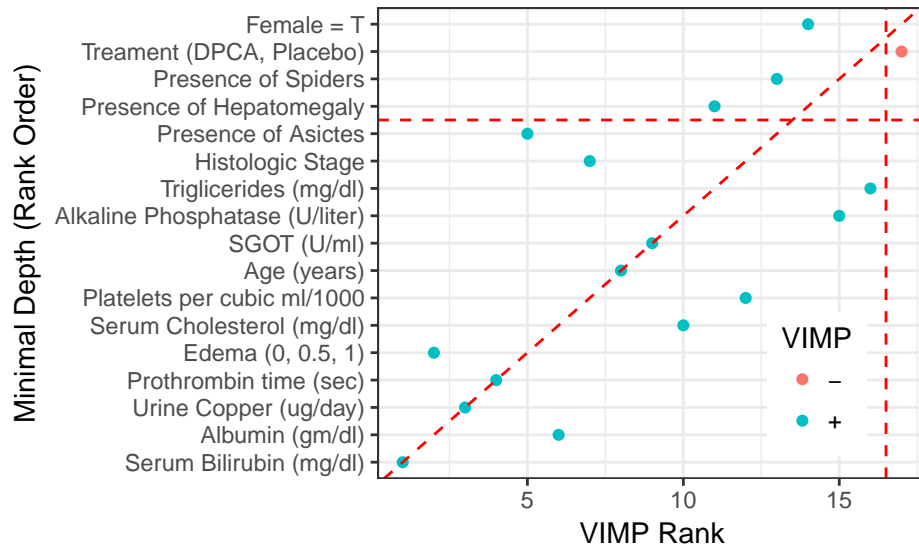


Figure 12: Comparing Minimal Depth and Vimp rankings. Points on the red dashed line are ranked equivalently, points above have higher VIMP ranking, those below have higher minimal depth ranking.

see all three methods indicate a strong relation of serum bilirubin to survival, and overall, the minimal depth and VIMP rankings agree reasonably well with the (Fleming and Harrington 1991) model.

The minimal depth selection process reduced the number of variables of interest from ~17 to 13, which is still a rather large subset of interest. An obvious selection set is to examine the five variables selected by (Fleming and Harrington 1991). Combining the Minimal Depth and (Fleming and Harrington 1991) model, there may be evidence to keep the top 7 variables. Though minimal depth does not indicate the **edema** variable is very interesting, VIMP ranking does agree with the proportional hazards model, indicating we might not want to remove the **edema** variable. Both minimal depth and VIMP suggest including **copper**, a measure associated with liver disease.

Regarding the **chol** variable, recall missing data summary of Table ?? . In in the trial data set, there were 28 observations missing **chol** values. The forest imputation randomly sorts observations with missing values into daughter nodes when using the **chol** variable, which is also how **randomForestSRC** calculates VIMP. We therefore expect low values for VIMP when a variable has a reasonable number of missing values.

Restricting our remaining analysis to the five (Fleming and Harrington 1991) variables, plus the **copper** retains the biological sense of these analysis. We will now examine how these six variables are related to survival using variable dependence methods to determine the direction of the effect and verify that the log transforms used by (Fleming and Harrington 1991) are appropriate.

5. Variable dependence

Table 4: Comparison of variable selection criteria. Minimal depth ranking, VIMP ranking and [fleming:1991] (FH) proportional hazards model ranked according to ‘abs(Z stat)’ from Table ??.

Variable	FH	Min depth	VIMP
bili	1	1	1
albumin	2	2	6
copper	NA	3	3
prothrombin	4	4	4
edema	5	5	2
chol	NA	6	10
platelet	NA	7	12
age	3	8	8
sgot	NA	9	9
alk	NA	10	15
trig	NA	11	16
stage	NA	12	7
ascites	NA	13	5

As random forest is not parsimonious, we have used minimal depth and VIMP to reduce the number of variables to a manageable subset. Once we have an idea of which variables contribute most to the predictive accuracy of the forest, we would like to know how the response depends on these variables.

Although often characterized as a *black box* method, the forest predictor is a function of the predictor variables $\hat{f}_{RF} = f(x)$. We use graphical methods to examine the forest predicted response dependency on covariates. We again have two options, variable dependence plots (Section ??) are quick and easy to generate, and partial dependence plots (Section ??) are more computationally intensive but give us a risk adjusted look at variable dependence.

5.1. Variable Dependence

Variable dependence plots show the predicted response relative to a covariate of interest, with each training set observation represented by a point on the plot. Interpretation of variable dependence plots can only be in general terms, as point predictions are a function of all covariates in that particular observation.

Variable dependence is straight forward to calculate, involving only the getting the predicted response for each observation. In survival settings, we must account for the additional dimension of time. We plot the response at specific time points of interest, for example survival at 1 or 3 years.

```
R> ggRFsrc + geom_vline(aes(xintercept = 1), linetype = "dashed") +
R+   geom_vline(aes(xintercept = 3), linetype = "dashed") +
R+   coord_cartesian(xlim = c(0, 5))
```

The `gg_rfsrc` of Figure 13 identical to Figure 7 (stored in the `ggRFsrc` variable) with the addition of a vertical dashed line at the 1 and 3 year survival time. A variable dependence plot

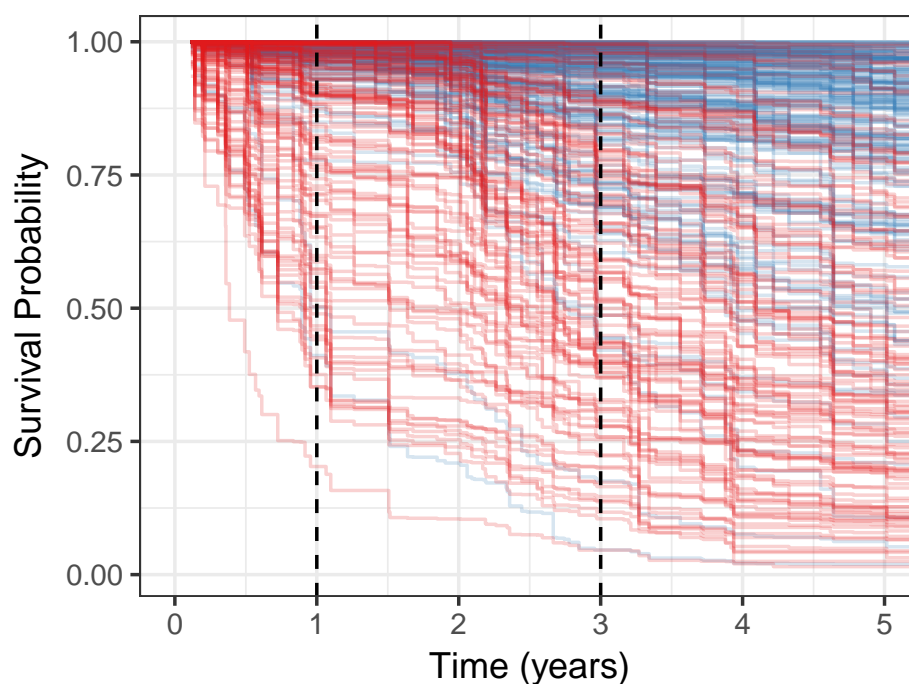


Figure 13: Random forest predicted survival (Figure 7) with vertical dashed lines indicate the 1 and 3 year survival estimates.

is generated from the predicted response value of each survival curve at the intersecting time line plotted against covariate value for that observation. This can be visualized as taking a slice of the predicted response at each time line, and spreading the resulting points out along the variable of interest.

The `gg_variable` function extracts the training set variables and the predicted OOB response from `rfsrc` and `predict` objects. In the following code block, we store the `gg_variable` data object for later use (`gg_v`), as all remaining variable dependence plots can be constructed from this object.

```
R> gg_v <- gg_variable(rfsrc_pbc, time = c(1, 3),
R+                   time.labels = c("1 Year", "3 Years"))
R>
R> plot(gg_v, xvar = "bili", alpha = 0.4) + #, se=FALSE
R+   labs(y = "Survival", x = st.labs["bili"]) +
R+   theme(legend.position = "none") +
R+   scale_color_manual(values = strCol, labels = event.labels) +
R+   scale_shape_manual(values = event.marks, labels = event.labels) +
R+   coord_cartesian(ylim = c(-0.01, 1.01))
```

The `gg_variable` plot of Figure 14 shows variable dependence for the Serum Bilirubin (`bili`) variable. Again censored cases are shown as blue circles, events are indicated by the red `x` symbols. Each predicted point is dependent on the full combination of all other covariates, not only on the covariate displayed in the dependence plot. The smooth loess line (Cleveland

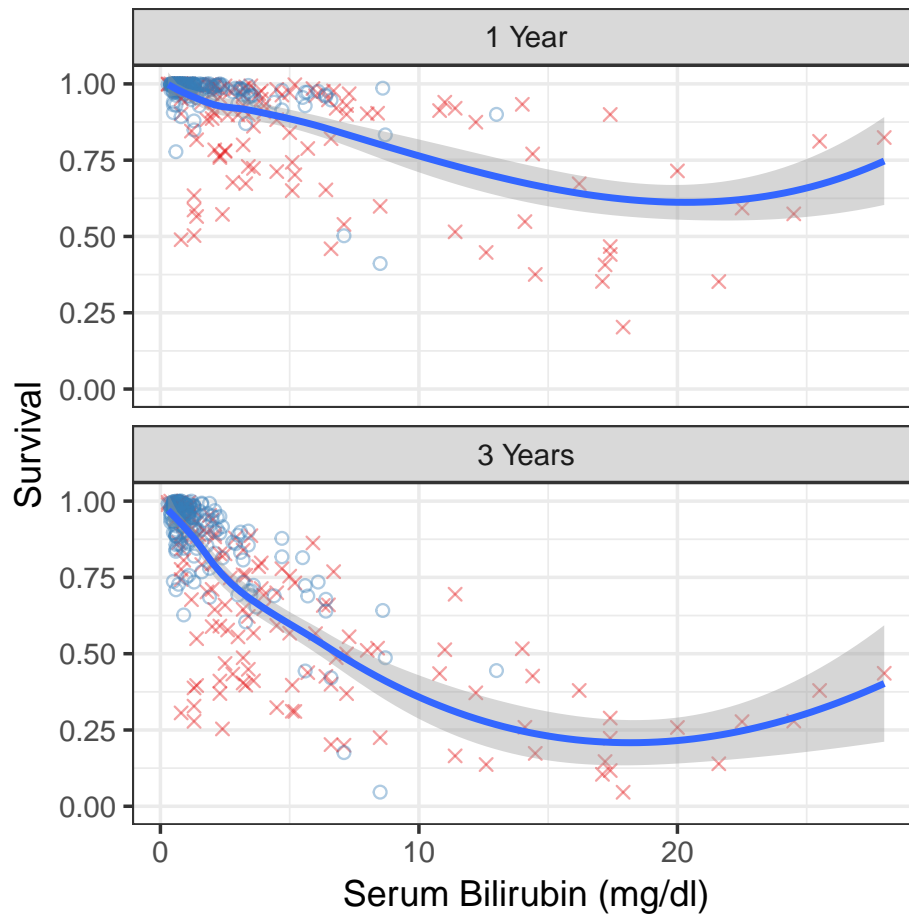


Figure 14: Variable dependence of survival at 1 and 3 years on bili variable. Individual cases are marked with blue circles (alive or censored) and red x (dead). Loess smooth curve with shaded 95% confidence band indicates decreasing survival with increasing bilirubin.

1981; Cleveland and Devlin 1988) indicates the trend of the prediction over the change in the variable.

Examination of Figure 14 indicates most of the cases are grouped in the lower end of bili values. We also see that most of the higher values experienced an event. The “normal” range of Bilirubin is from 0.3 to 1.9 mg/dL, indicating the distribution from our population is well outside the normal range. These values make biological sense considering Bilirubin is a pigment created in the liver, the organ effected by the PBC disease. The figure also shows that the risk of death increases as time progresses. The risk at 3 years is much greater than that at 1 year for patients with high Bilirubin values compared to those with values closer to the normal range.

The `plot.gg_variable` function call operates on the `gg_variable` object controlled by the list of variables of interest in the `xvar` argument. By default, the `plot.gg_variable` function returns a list of `ggplot` objects, one figure for each variable named in `xvar`. The remaining arguments are passed to internal `ggplot2` functions controlling the display of the figure. The `se` argument is passed to the internal call to `geom_smooth` for fitting smooth lines to the data.

The `alpha` argument lightens the coloring points in the `geom_point` call, making it easier to see point over plotting. We also demonstrate modification of the plot labels using the `labs` function and point attributes with the `scale_` functions.

An additional `plot.gg_variable` argument (`panel = TRUE`) can be used to combine multiple variable dependence plots into a single figure. In the following code block, we plot the remaining continuous variables of interest found in Section ??.

```
R> xvar <- c("bili", "albumin", "copper", "prothrombin", "age")
R> xvar.cat <- c("edema")
R>
R> plot(gg_v, xvar = xvar[-1], panel = TRUE, alpha = 0.4) +
R+   labs(y = "Survival") +
R+   theme(legend.position = "none") +
R+   scale_color_manual(values = strCol, labels = event.labels) +
R+   scale_shape_manual(values = event.marks, labels = event.labels) +
R+   coord_cartesian(ylim = c(-0.05, 1.05))
```

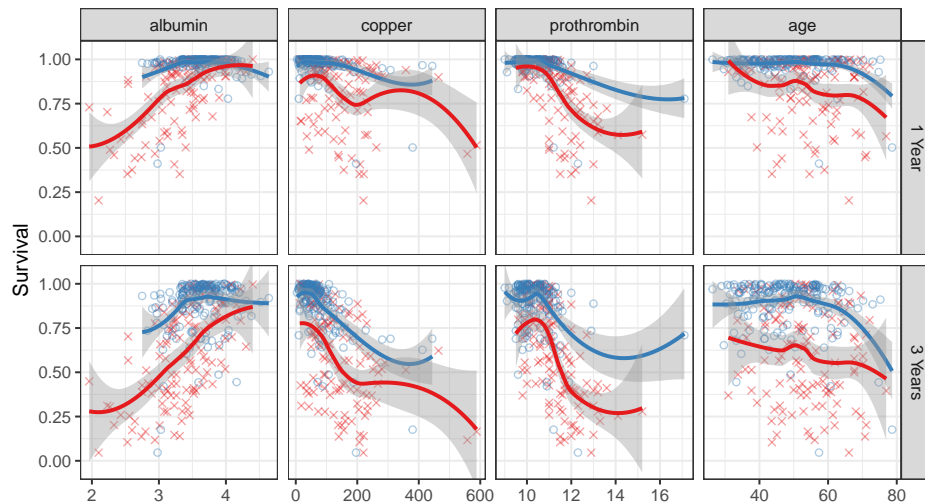


Figure 15: Variable dependence of predicted survival at 1 and 3 years on continuous variables of interest. Individual cases are marked with blue circles for censored cases and red x for death events. Loess smooth curve indicates the survival trend with increasing values.

The `gg_variable` plot in Figure 15 displays a panel of the remaining continuous variable dependence plots. The panels are sorted in the order of variables in the `xvar` argument and include a smooth loess line (Cleveland 1981; Cleveland and Devlin 1988) to indicate the trend of the prediction dependence over the covariate values. The `se=FALSE` argument turns off the loess confidence band, and the `span=1` argument controls the degree of smoothing.

The figures indicate that survival increases with `albumin` level, and decreases with `bili`, `copper`, `prothrombin` and `age`. Note the extreme value of `prothrombin` (> 16) influences the loess curve more than other points, which would make it a candidate for further investigation. We expect survival at 3 years to be lower than at 1 year. However, comparing the two time plots for each variable does indicate a difference in response relation for `bili`, `copper` and

prothrombine. The added risk for high levels of these variables at 3 years indicates a non-proportional hazards response. The similarity between the time curves for **albumin** and **age** indicates the effect of these variables is constant over the disease progression.

There is not a convenient method to panel scatter plots and boxplots together, so we recommend creating panel plots for each variable type separately. We plot the categorical variable (**edema**) in Figure 16 separately from the continuous variables in Figure 15.

```
R> plot(gg_v, xvar = xvar.cat, alpha = 0.4) + labs(y = "Survival") +
R+   theme(legend.position = "none") +
R+   scale_color_manual(values = strCol, labels = event.labels) +
R+   scale_shape_manual(values = event.marks, labels = event.labels) +
R+   coord_cartesian(ylim = c(-0.01, 1.02))
```

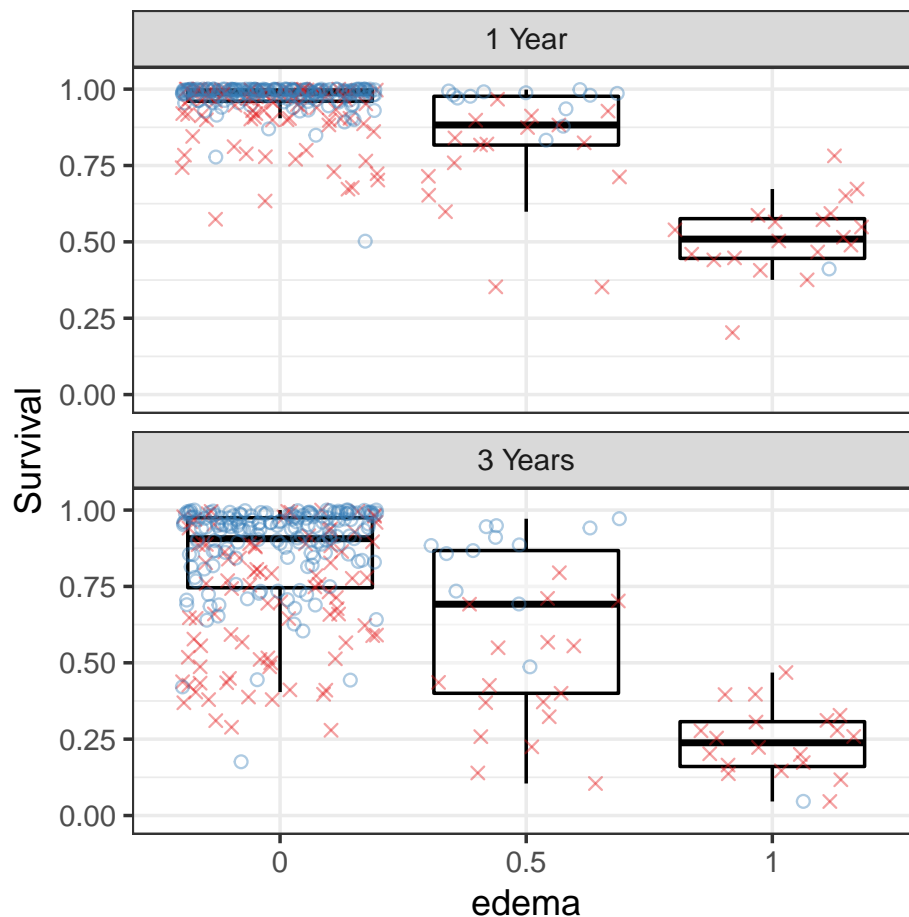


Figure 16: Variable dependence of survival 1 and 3 years on edema categorical variable. Symbols with blue circles indicate censored cases and red x indicate death events. Boxplots indicate distribution of predicted survival for all observations within each edema group.

The `gg_variable` plot of Figure 16 for categorical variable dependence displays boxplots to examine the distribution of predicted values within each level of the variable. The points

are plotted with a jitter to see the censored and event markers more clearly. The boxes are shown with horizontal bars indicating the median, 75th (top) and 25th (bottom) percentiles. Whiskers extend to 1.5 times the interquartile range. Points plotted beyond the whiskers are considered outliers.

When using categorical variables with linear models, we use boolean dummy variables to indicate class membership. In the case of `edema`, we would probably create two logical variables for `edema = 0.5` (complex Edema presence indicator) and `edema = 1.0` (Edema with diuretics) contrasted with the `edema = 0` variable (no Edema). Random Forest can use factor variables directly, separating the populations into homogeneous groups of `edema` at nodes that split on that variable. Figure 16 indicates similar survival response distribution between 1 and 3 year when `edema = 1.0`. The distribution of predicted survival does seem to spread out more than for the other values, again indicating a possible non-proportional hazards response.

5.2. Partial Dependence

Partial dependence plots are a risk adjusted alternative to variable dependence. Partial plots are generated by integrating out the effects of variables beside the covariate of interest. The figures are constructed by selecting points evenly spaced along the distribution of the variable of interest. For each of these points ($X = x$), we calculate the average RF prediction over all remaining covariates in the training set by

$$\tilde{f}(x) = \frac{1}{n} \sum_{i=1}^n \hat{f}(x, x_{i,o}),$$

where \hat{f} is the predicted response from the random forest and $x_{i,o}$ is the value for all other covariates other than $X = x$ for observation i (Friedman 2000).

Generating partial dependence data is effectively averaging the response for a series of nomograms constructed for each observation by varying the variable of interest. The operation is computationally intensive, especially when there are a large number of observations. The default parameters for the `plot.variable` function generate partial dependence estimates at `pts = 25` points along the variable of interest. For each point of interest, the `plot.variable` function averages the `n` response predictions. This process is repeated for each of the variables of interest.

For time to event data, we also have to deal with the additional time dimension, as with variable dependence. The following code block uses the `mclapply` function from the `parallel` package to run the `plot.variable` function for three time points (`time=1, 3` and `5` years) in parallel. For RSF models, we calculate a risk adjusted survival estimates (`surv.type="surv"`), suppressing the internal base graphs (`show.plots = FALSE`) and store the point estimates in the `partial_pbc` list.

```
R> xvar <- c(xvar, xvar.cat)
R>
R> time_index <- c(which(rfsrc_pbc$time.interest > 1)[1]-1,
R+               which(rfsrc_pbc$time.interest > 3)[1]-1,
R+               which(rfsrc_pbc$time.interest > 5)[1]-1)
R> partial_pbc <- mclapply(rfsrc_pbc$time.interest[time_index],
```

```

R+           function(tm){
R+               plot.variable(rfsrc_pbc, surv.type = "surv",
R+                           time = tm, xvar.names = xvar,
R+                           partial = TRUE ,
R+                           show.plots = FALSE)
R+           })

```

Because partial dependence data is collapsed onto the risk adjusted response, we can show multiple time curves on a single panel. The following code block converts the `plot.variable` output into a list of `gg_partial` objects, and then combines these data objects, with descriptive labels, along each variable of interest using the `combine.gg_partial` function.

```

R> gg_dta <- mclapply(partial_pbc, gg_partial)
R> pbc_ggpart <- combine.gg_partial(gg_dta[[1]], gg_dta[[2]],
R+                               lbls = c("1 Year", "3 Years"))

```

We then segregate the continuous and categorical variables, and generate a panel plot of all continuous variables in the `gg_partial` plot of Figure 17. The panels are ordered by minimal depth ranking. Since all variables are plotted on the same Y-axis scale, those that are strongly related to survival make other variables look flatter. The figures also confirm the strong non-linear contribution of these variables. Non-proportional hazard response is also evident in at least the `bili` and `copper` variables by noting the divergence of curves as time progresses.

Categorical partial dependence is displayed as boxplots, similar to categorical variable dependence. Risk adjustment greatly reduces the spread of the response as expected, and may also move the mean response compared to the unadjusted results. The categorical `gg_partial` plot of Figure 18 indicates that, adjusting for other variables, survival decreases with rising `edema` values. We also note that the risk adjusted distribution does spread out as we move further out in time.

```

R> ggplot(pbc_ggpart[["edema"]], aes(y=yhat, x=edema, col=group))+
R+   geom_boxplot(notch = TRUE,
R+               outlier.shape = NA) + # panel=TRUE,
R+   labs(x = "Edema", y = "Survival (%)", color="Time", shape="Time") +
R+   theme(legend.position = c(0.2, 0.2))

```

Partial dependence is an extrapolation operation. By averaging over a series of nomograms, the algorithm constructs observations for all values of the variable of interest, regardless of the relation with other variables. In contrast, variable dependence only uses observations from within the training set. A simple example would be for a model including BMI, weight and height. When examining partial dependence of BMI, the algorithm only manipulates BMI values, height or weight values. The averaging operation is then confounded in two directions. First, dependence on height and weight is shared with BMI, making it difficult to see the true response dependence. Second, partial dependence is calculated over nomograms that can not physically occur. For simple variable combinations, like BMI, it is not difficult to recognize this and modify the independent variable list to avoid these issues. However, care must be taken when interpreting more complex biological variables.

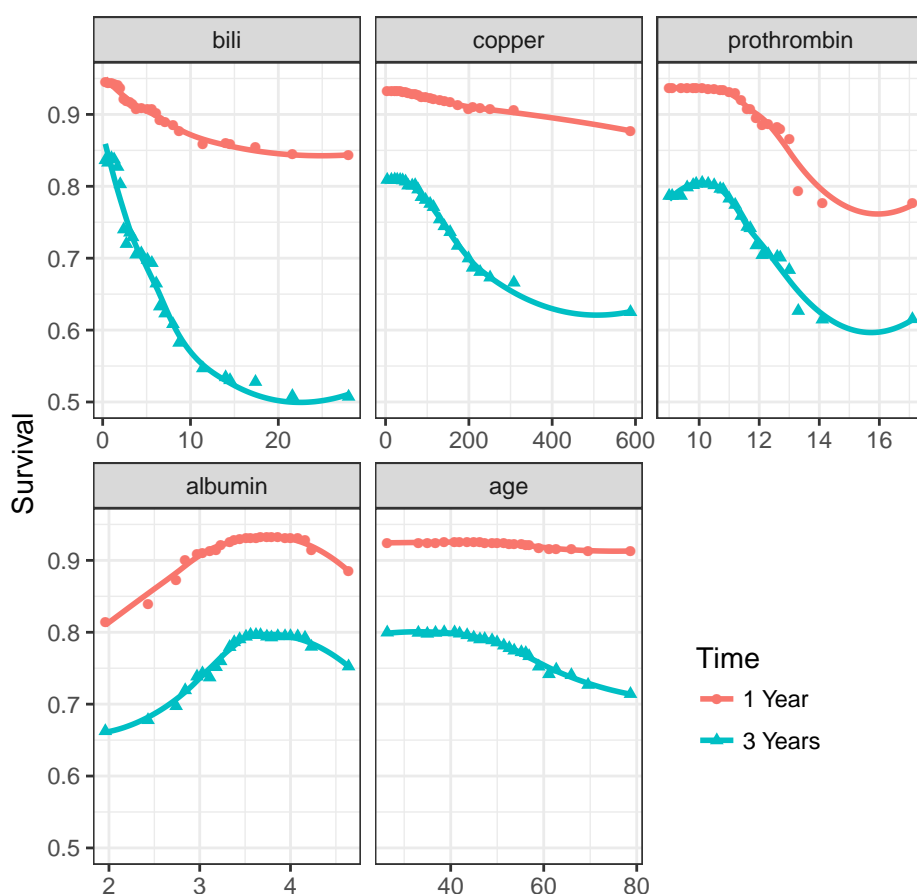


Figure 17: Partial dependence of predicted survival at 1 year (red circle) and 3 years (blue triangle) as a function continuous variables of interest. Symbols are partial dependence point estimates with loess smooth line to indicate trends.

5.3. Partial dependence as a function of time

In the previous section, we calculated risk adjusted (partial) dependence at two time points (1 and 3 years). The selection of these points can be driven by biological times of interest (i.e., 1 year and 5 year survival in cancer studies) or by investigating time points of interest from a `gg_rfsrc` prediction plot. We typically restrict generating `gg_partial` plots to the variables of interest at two or three time points of interest due to computational constraints.

It is instructive to see a more detailed map of the risk adjusted response to get a feel for interpreting partial and variable dependence plots. In Figure 17, we can visualize the two curves as extending into the plane of the page along a time axis. Filling in more partial dependence curves, it is possible to create a partial dependence surface.

For this exercise, we will generate a series of 50 `gg_partial` plot curves for the `bili` variable. To fill the surface in, we also increased the number of points along the distribution of `bili` to `npts=50` to create a grid of 50×50 risk adjusted estimates of survival along time in one dimension and the `bili` variable in the second.

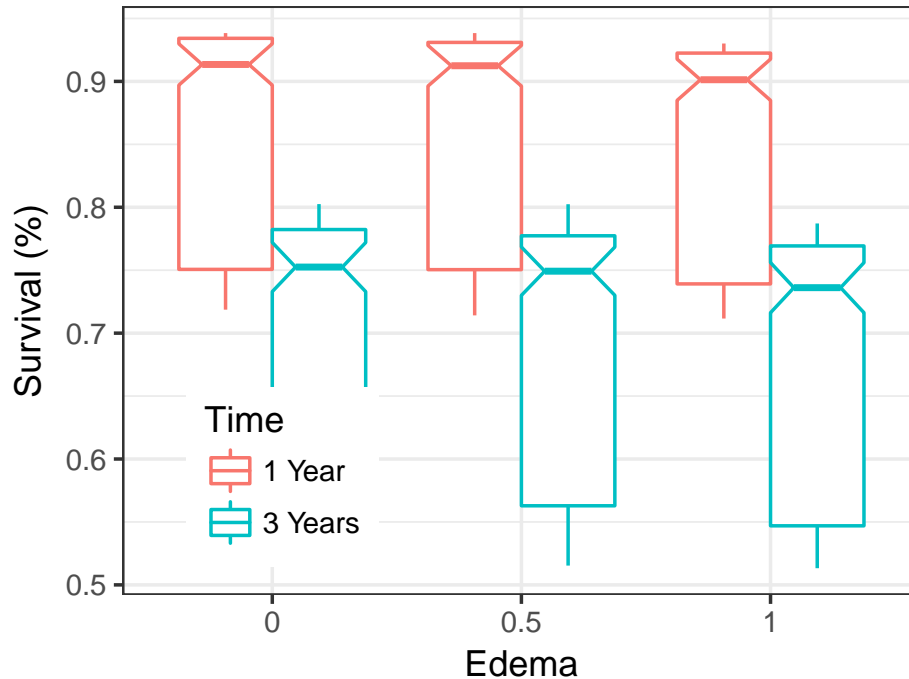


Figure 18: Partial dependence plot of predicted survival at 1 year (red) and 3 years (blue) as a function of edema groups (categorical variable). Boxplots indicate distribution within each group.

The `gg_partial` surface of Figure ?? was constructed using the `surf3D` function from the `plot3D` package (?, <http://CRAN.R-project.org/package=plot3D>). Source code for generating this figure is shown in Appendix ??.

The figure shows partial dependence of survival (Z-axis) as a function of `bili` over a five year follow up time period. Lines perpendicular to the Bilirubin axis are distributed along the `bili` variable. Lines parallel to the Bilirubin axis are taken at 50 training set event times, the first event after $t = 0$ at the back to last event before $t = 5$ years at the front. The distribution of the time lines is also evenly selected using the same procedure as selecting points for partial dependence curves.

The 2500 estimated partial dependence points are joined together with a simple straight line interpolation to create the surface, colored according to the survival estimates (yellow close to 1, red for lower values) to aid the visualization of 3 dimensions on a 2 dimensional page. The blue lines in Figure ?? correspond to the 1 and 3 year partial dependence, as shown in the `bili` panel of Figure 17.

Viewed as a surface, we see how the partial dependence changes with time. For low values of `bili`, survival decreases at a constant rate. For higher values, the rate seems constant until somewhere near 2 years, where it increases rapidly before slowing again as we approach the 5 year point.

6. Variable Interactions

We could stop with the results that our RF analysis has found these six variables to be important in predicting survival. Where the survival response is decreasing with increasing `bili`, `copper`, `prothrombin`, `age` and `edema` and increasing with increasing `albumin`. These results agree with the sign of the `?` model coefficients shown in Table ???. The `gg_partial` plot in Figure 17 supports the `log` transform of `bili`, `albumin` and `prothrombin` and suggest a similar transform for including the `copper` variable in a proportional hazards model. The `age` variable does seem to have a more linear response than the other continuous variables, and using dummy variables for `edema` would preclude the need for a transformation.

Using minimal depth, it is also possible to calculate measures of pairwise interactions among variables. Recall that minimal depth measure is defined by averaging the tree depth of variable i relative to the root node. To detect interactions, this calculation can be modified to measure the minimal depth of a variable j with respect to the maximal subtree for variable i (??).

The `randomForestSRC::find.interaction` function traverses the forest, calculating all pairwise minimal depth interactions, and returns a $p \times p$ matrix of interaction measures. The diagonal terms are normalized to the root node, and off diagonal terms are normalized measures of pairwise variable interaction.

```
R> ggint <- gg_interaction(rfsrc_pbc)
```

```

Method: maxsubtree
No. of variables: 17
Variables sorted by minimal depth?: TRUE

      bili albumin copper prothrombin edema chol platelet age sgot
bili    0.12   0.26   0.29         0.29 0.51 0.29      0.29 0.25 0.31
albumin  0.32   0.16   0.35         0.38 0.64 0.36      0.38 0.35 0.38
copper   0.34   0.37   0.18         0.39 0.66 0.41      0.40 0.36 0.41
prothrombin 0.37   0.40   0.39         0.19 0.67 0.43      0.44 0.41 0.45
edema    0.44   0.48   0.48         0.49 0.21 0.50      0.49 0.49 0.51
chol     0.41   0.44   0.44         0.45 0.77 0.22      0.48 0.43 0.46
platelet 0.52   0.51   0.53         0.56 0.81 0.55      0.23 0.51 0.53
age      0.41   0.42   0.44         0.48 0.81 0.46      0.46 0.24 0.45
sgot     0.47   0.47   0.50         0.53 0.79 0.51      0.51 0.47 0.25
alk      0.56   0.56   0.56         0.59 0.85 0.58      0.60 0.55 0.57
trig     0.56   0.56   0.57         0.60 0.85 0.60      0.58 0.54 0.58
stage    0.59   0.62   0.61         0.63 0.84 0.63      0.63 0.59 0.63
ascites  0.54   0.56   0.58         0.59 0.75 0.59      0.60 0.58 0.61
hepatom  0.70   0.71   0.72         0.73 0.87 0.71      0.74 0.70 0.70
spiders  0.81   0.82   0.81         0.83 0.94 0.83      0.83 0.81 0.82
treatment 0.87   0.88   0.88         0.87 0.97 0.90      0.88 0.87 0.89
sex      0.84   0.85   0.86         0.88 0.96 0.88      0.85 0.86 0.87

      alk trig stage ascites hepatom spiders treatment sex
bili    0.31 0.34  0.45    0.67   0.63    0.63      0.59 0.67
albumin 0.39 0.41  0.56    0.78   0.69    0.69      0.67 0.74
copper  0.41 0.43  0.56    0.79   0.68    0.69      0.68 0.77

```

prothrombin	0.42	0.48	0.60	0.81	0.73	0.74	0.71	0.79
edema	0.49	0.52	0.64	0.81	0.73	0.75	0.74	0.81
chol	0.46	0.50	0.63	0.86	0.76	0.77	0.73	0.82
platelet	0.54	0.58	0.70	0.92	0.81	0.81	0.78	0.86
age	0.46	0.50	0.62	0.90	0.77	0.76	0.71	0.81
sgot	0.52	0.55	0.67	0.89	0.80	0.80	0.77	0.84
alk	0.26	0.61	0.75	0.93	0.85	0.84	0.81	0.89
trig	0.59	0.30	0.73	0.91	0.84	0.85	0.80	0.88
stage	0.61	0.66	0.35	0.92	0.85	0.84	0.81	0.88
ascites	0.60	0.63	0.71	0.37	0.79	0.80	0.78	0.83
hepatom	0.71	0.74	0.80	0.93	0.46	0.87	0.86	0.89
spiders	0.83	0.85	0.88	0.96	0.93	0.53	0.91	0.95
treatment	0.88	0.90	0.92	0.99	0.95	0.96	0.54	0.97
sex	0.85	0.88	0.92	0.98	0.95	0.94	0.94	0.61

The `gg_interaction` function wraps the `find_interaction` matrix for use with the **ggRandomForests** plot and print functions. The `xvar` argument is used to restrict the variables of interest and the `panel = TRUE` argument displays the results in a single figure.

```
R> plot(ggint, xvar = xvar)
```

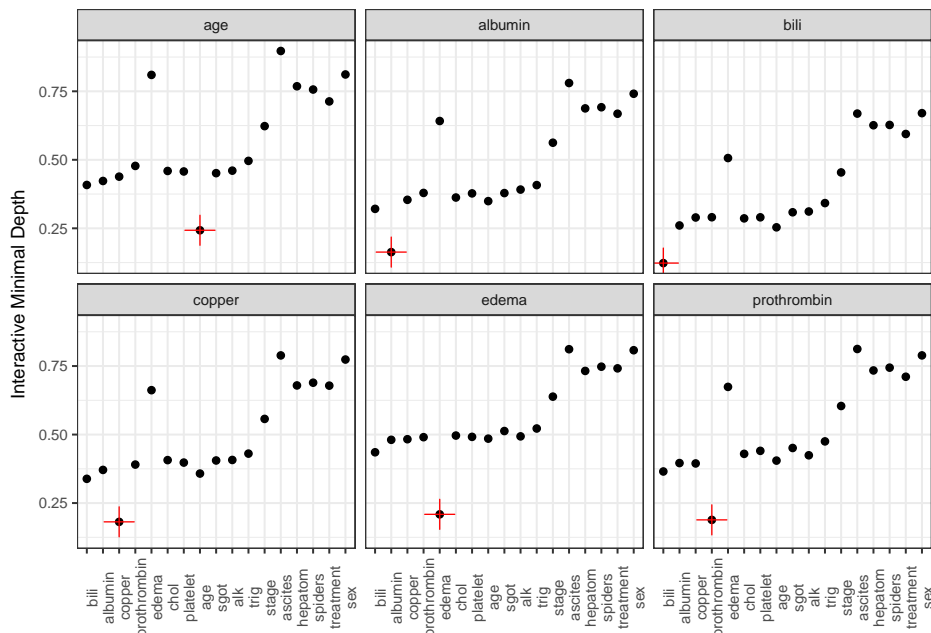


Figure 19: Minimal depth variable interaction plot for six variables of interest. Higher values indicate lower interactivity with target variable marked in red.

The `gg_interaction` plots in Figure 19 show interactions for the target variable (shown with the red cross) with interaction scores for all remaining variables. We expect the covariate with lowest minimal depth (**bili**) to be associated with almost all other variables, as it typically

splits close to the root node, so viewed alone it may not be as informative as looking at a collection of interactive depth plots. Scanning across the panels, we see each successive target depth increasing, as expected. We also see the interactive variables increasing with increasing target depth.

7. Conditional dependence plots

Conditioning plots (coplots) (??) are a powerful visualization tool to efficiently study how a response depends on two or more variables (?). The method allows us to view data by grouping observations on some conditional membership. The simplest example involves a categorical variable, where we plot our data conditional on class membership, for instance on groups of the `edema` variable. We can view a coplot as a stratified variable dependence plot, indicating trends in the RF prediction results within panels of group membership.

Interactions with categorical data can be generated directly from variable dependence plots. Recall the variable dependence for bilirubin shown in Figure 14. We recreated the `gg_variable` plot in Figure 20, modified by adding a linear smooth as we intend on segregating the data along conditional class membership.

```
R> # Get variable dependence at 1 year
R> ggvar <- gg_variable(rfsrc_pbc, time = 1)
R>
R> # For labeling coplot membership
R> ggvar$edema <- paste("edema = ", ggvar$edema, sep = "")
R>
R> # Plot with linear smooth (method argument)
R> var_dep <- plot(ggvar, xvar = "bili",
R+               alpha = 0.5) +
R+ #   geom_smooth(method = "glm", se = FALSE) +
R+   labs(y = "Survival",
R+        x = st.labs["bili"]) +
R+   theme(legend.position = "none") +
R+   scale_color_manual(values = strCol, labels = event.labels) +
R+   scale_shape_manual(values = event.marks, labels = event.labels) +
R+   coord_cartesian(y = c(-.01, 1.01))
R>
R> var_dep
```

We can view the conditional dependence of survival against bilirubin, conditional on `edema` group membership (categorical variable) in Figure 21 by reusing the saved `ggplot` object (`var_dep`) and adding a call to the `facet_grid` function.

```
R> var_dep + facet_grid(~edema)
```

Comparing Figure 20 with conditional panels of Figure 21, we see the overall response is similar to the `edema=0` response. The survival for `edema=0.5` is slightly lower, though the slope of the smooth indicates a similar relation to `bili`. The `edema=1` panel shows that the

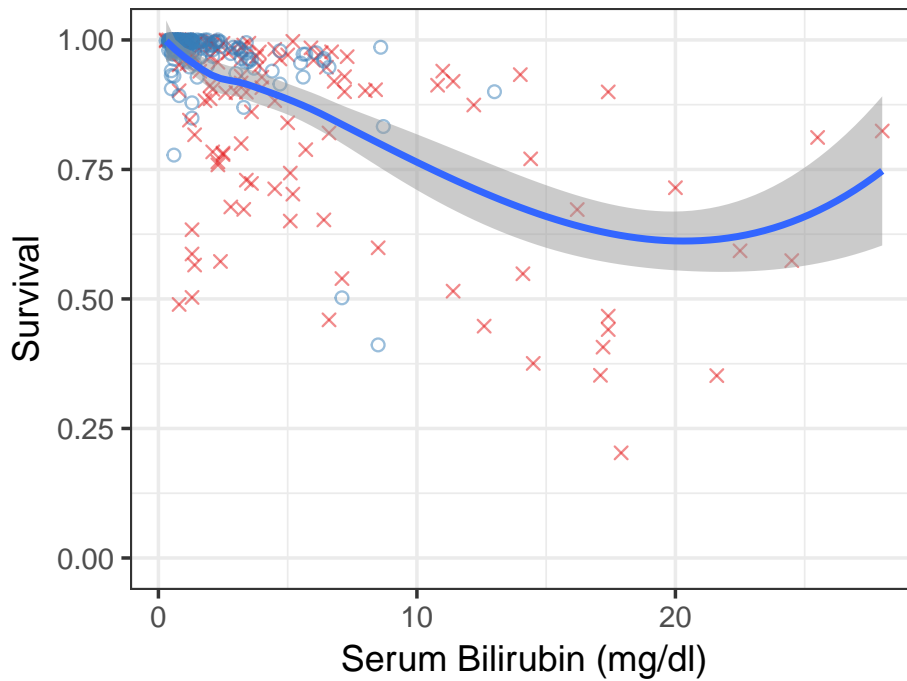


Figure 20: Variable dependence of survival at 1 year against bili variable. Reproduction of top panel of Figure 14 with a linear smooth to indicate trend.

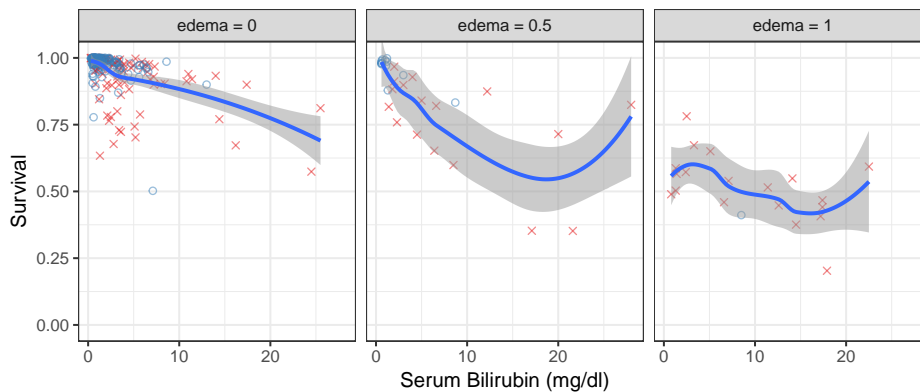


Figure 21: Variable dependence coplot of survival at 1 year against bili, conditional on edema group membership. Linear smooth indicates trend of variable dependence.

survival for this (smaller) group of patients is worse, but still follows the trend of decreasing with increasing `bili`.

Conditional membership within a continuous variable requires stratification at some level. We can sometimes make these stratification along some feature of the variable, for instance a variable with integer values, or 5 or 10 year age group cohorts. However with our variables of interest, there are no logical stratification indications. Therefore we arbitrarily stratify our variables into 6 groups of roughly equal population size using the `quantile_cuts` function.

We pass the break points located by `quantile_cuts` to the `cut` function to create grouping intervals, which we can then add to the `gg_variable` object before plotting with the `plot.gg_variable` function. This time we use the `facet_wrap` function to generate the panels grouping interval, which automatically sorts the six panels into two rows of three panels each.

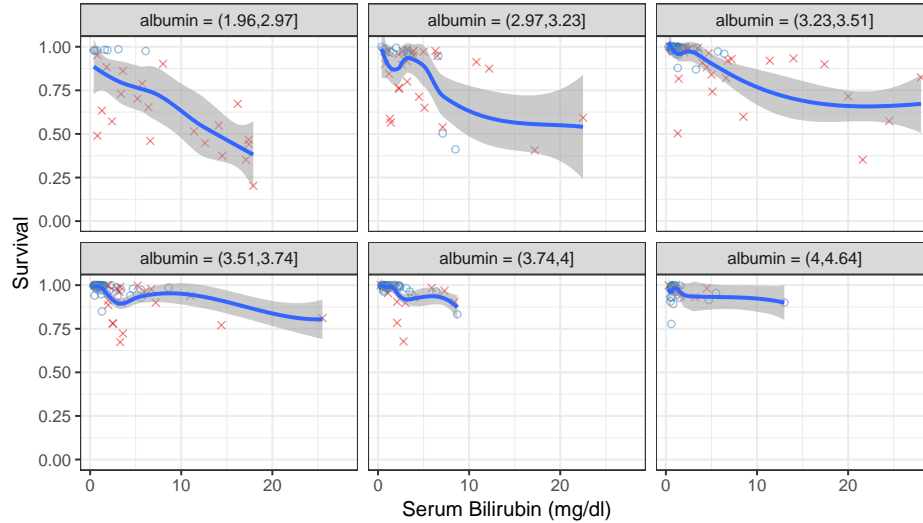


Figure 22: Variable dependence coplot of survival at 1 year against bili, conditional on albumin interval group membership.

The `gg_variable` coplot of Figure 22 indicates that the effect of `bili` decreases conditional on membership within increasing `albumin` groups. To get a better feel for how the response depends on both these variables together, it is instructive to look at the compliment coplot of `albumin` conditional on membership in `bili` groups. We repeat the previous coplot process, predicted survival as a function of the `albumin` variable, conditional on membership within 6 groups `bili` intervals. As the code to create the coplot of Figure 23 is nearly identical to the code for creating Figure 22, we include the source code for this figure in Appendix ??.

The `gg_variable` coplot of Figure 23 indicates the probability of survival increases with increasing `albumin` and increases within groups of increasing `bili`.

Typically, conditional plots for continuous variables include overlapping intervals along the grouped variable (?). We chose to use mutually exclusive continuous variable intervals for the following reasons:

- **Simplicity** - We can create the coplot figures directly from the `gg_variable` object by adding a conditional group column directly to the object.
- **Interpretability** - We find it easier to interpret and compare the panels if each observation is only in a single panel.
- **Clarity** - We prefer using more space for the data portion of the figures than typically displayed in the `coplot` function which requires the bar plot to present the overlapping segments.

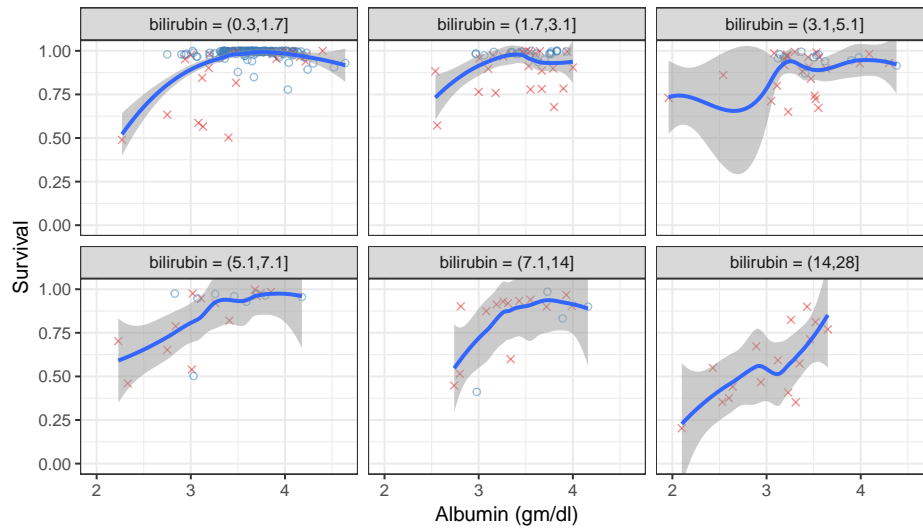


Figure 23: Variable dependence coplot of survival at 1 year against albumin, conditional on bili interval group membership.

It is still possible to augment the `gg_variable` to include overlapping conditional membership with continuous variables by duplicating rows of the training set data within the `rfsrc$xvar` object, and then setting the conditional group membership as described. The `plot.gg_variable` function recipe above could be used to generate the panel plot, with panels ordered according to the factor levels of the grouping variable. We leave this as an exercise for the reader.

7.1. Partial dependence coplots

By characterizing conditional plots as stratified variable dependence plots, the next logical step would be to generate an analogous conditional partial dependence plot. The process is similar to variable dependence coplots, first determine conditional group membership, then calculate the partial dependence estimates on each subgroup using the `plot.variable` function with a `subset` argument for each grouped interval. The `ggRandomForests` `gg_partial_coplot` function is a wrapper for generating conditional partial dependence data objects. Given a random forest (`rfsrc`) object and a `groups` vector for conditioning the training data set observations, `gg_partial_coplot` calls the `plot.variable` function the training set observations conditional on `groups` membership. The function returns a `gg_partial_coplot` object, a subclass of the `gg_partial` object, which can be plotted with the `plot.gg_partial` function.

The following code block will generate the data object for creating partial dependence coplot of 1 year survival as a function of `bili` conditional on membership within the 6 groups of `albumin` intervals that we examined in the Figure 22.

```
\begin{Schunk} \begin{Sinput} R> partial_coplot_pbc <- gg_partial_coplot(rfsrc_pbc, xvar
= "bili", R+ groups = ggvar$albumin_grp, R+ surv_type = "surv", R+ time = 1, R+
show.plots = FALSE) \end{Sinput} \end{Schunk}
```

```
R> ggplot(partial_coplot_pbc, aes(x=bili, y=yhat, col=group, shape=group)) +
R+   geom_smooth(se = FALSE) +
```

```
R+ labs(x = st.labs["bili"], y = "Survival at 1 year (%)",
R+       color = "albumin", shape = "albumin")
```

The `gg_partial_coplot` of Figure ?? shows point estimates of the risk adjusted survival as a function of `bili` conditional on group membership defined by `albumin` intervals. The figure is slightly different than the `gg_partial` plot of Figure 17 as each set of partial dependence estimates is calculated over a subset of the training data. We again connect the point estimates with a Loess curve.

For completeness, we construct the compliment coplot view of one year survival as a function of `albumin` conditional on `bili` interval group membership in Figure ?. We list the source code for this figure in Appendix ?.

8. Partial plot surfaces

Just as in partial dependence, we can view the partial coplot curves as slices along a surface that could extend along an axis into the page. This visualization is made a bit difficult by our choice to select groups of similar population size, as the curves are not evenly spaced along the grouping variables. So, similar to the partial dependence surface we created along time in Section ?, we can examine the relation of these two variables using a partial dependence surface.

A difficulty with conditional dependence for this exercise is the reduction of the sample sizes for calculating a coplot surface. So instead, we calculate the full partial dependence surface by generating 50 `albumin` values spaced evenly along the data distribution. For each value of `albumin`, we calculate the partial dependence on `bili` at `npts = 50` points with the `plot.variable` function. We generate the surface again using the `surf3D` function.

The partial dependence surface of Figure ? shows partial dependence of 1 year survival on the Z-axis against values of Bilirubin and Albumin. We again use linear interpolation between the 2500 estimates, and color the surface by the response. Here blue corresponds to lower and yellow to higher risk adjusted survival. The blue lines are placed at the cut points between groups of `albumin` and `bili` used in the partial coplots of Figures ? and ? respectively.

To construct the partial coplot for groups of `albumin` in Figure ?, we arbitrarily segmented the training set into 6 groups of equal membership size. The segments between blue lines parallel to the Bilirubin axis indicate where on the surface these observations are located. Similarly, the blues lines perpendicular to the Bilirubin axis segment observations into the 6 groups of `bili` intervals. Figure ? indicates the arbitrary grouping for groups of `bili` in Figure ?.

The figure indicates that partial dependence of higher `albumin` levels are similar, which results in the over plotting seen in Figure ?. The distribution is sparser at lower `albumin` levels, creating the larger area in lowest `albumin` values, where the partial dependence changes the most.

9. Conclusion

In this vignette, we have demonstrated the use of Random Survival Forest methods with the `ggRandomForests` (<http://CRAN.R-project.org/package=ggRandomForests>) package.

We have shown how to grow a random forest model and determine which variables contribute to the forest prediction accuracy using both VIMP and Minimal Depth measures. We outlined how to investigate variable associations with the response variable using variable dependence and the risk adjusted partial dependence plots. We've also explored variable interactions by using pairwise minimal depth interactions and directly viewed these interactions using variable dependence coplots and partial dependence coplots. Along the way, we've demonstrated the use of additional commands from the **ggplot2** package (<http://CRAN.R-project.org/package=ggplot2>) package for modifying and customizing plots from **ggRandomForests** functions.

9.1. Computational details

This document is a package vignette for the **ggRandomForests** package for “Visually Exploring Random Forests” (<http://CRAN.R-project.org/package=ggRandomForests>). The **ggRandomForests** package is designed for use with the **randomForestSRC** package (<http://CRAN.R-project.org/package=randomForestSRC>) for growing survival, regression and classification random forest models and uses the **ggplot2** package (<http://CRAN.R-project.org/package=ggplot2>) for plotting diagnostic and variable association results. **ggRandomForests** is structured to extract data objects from **randomForestSRC** objects and provides functions for printing and plotting these objects.

The vignette is a tutorial for using the **ggRandomForests** package with the **randomForestSRC** package for building and post-processing random survival forests. In this tutorial, we explore a random forest for survival model constructed for the primary biliary cirrhosis (PBC) of the liver data set ([?](#)), available in the **randomForestSRC** package. We grow a random survival forest and demonstrate how **ggRandomForests** can be used when determining how the survival response depends on predictive variables within the model. The tutorial demonstrates the design and usage of many of **ggRandomForests** functions and features and also how to modify and customize the resulting **ggplot** graphic objects along the way.

The vignette is written in \LaTeX using the **knitr** package (<http://CRAN.R-project.org/package=knitr>), which facilitates weaving R ([?](#)) code, results and figures into document text.

This vignette is available within the **ggRandomForests** package on the Comprehensive R Archive Network (CRAN) (<http://cran.r-project.org>). Once the package has been installed, the vignette can be viewed directly from within R with the following command:

```
R> vignette("randomForestSRC-Survival", package = "ggRandomForests")
```

A development version of the **ggRandomForests** package is also available on GitHub (<https://github.com>). We invite comments, feature requests and bug reports for this package at <https://github.com/ehrlinger/ggRandomForests>.

9.2. Acknowledgement

This work was supported in part by the National Institutes of Health grant R01-HL103552-01A1.

A. Source Code

Throughout this document, we have listed all R source code to create the figures included here with a few exceptions. For completeness, we include the missing code blocks in this appendix. The code blocks are included here in order of appearance in the document.

A.1. Partial Dependence in Time Dimension

The surface plot of ?? demonstrates how partial dependence curves relate to the survival curves. This code block is the R source code for creating Figure ??.

A.2. Bilirubin Coplot

In Section ??, we generate variable dependence coplots for the `bili` variable conditional on grouping on intervals of the `albumin` variable, and the complimentary `albumin` variable conditional on grouping on intervals of the `bili` variable. We include the source code for Figure 22 in the document. Since the code is nearly identical for the later case, we include the source code for generating Figure 23 here.

A.3. Bilirubin Partial Coplot

Similar to variable dependence coplots, In Section ??, we compare the partial dependence coplots for the same `albumin` and `bili` variable groupings. Again, the source code for Figure ?? is nearly identical to the source code for generating Figure ???. We include the partial dependence coplot source code for the `albumin` variable conditional on grouping on intervals of the `bili` variable.

A.4. Partial Dependence in Multiple Variable Dimensions

In Section ??, we generate a partial dependence surface of one year survival dependence on both `bili` and `albumin` variables. We include the Source code for generating Figure ??? here.

References

- R Core Team. 2014. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <http://www.R-project.org/>.
- Breiman, L. 1996a. “Bagging Predictors.” *Machine Learning* 26: 123–40.
- . 1996b. “Out-Of-Bag Estimation.” Statistics Department, University of California, Berkeley, CA. 94708. <https://www.stat.berkeley.edu/~breiman/OOBestimation.pdf>.
- Breiman, L, Jerome H Friedman, R Olshen, and C Stone. 1984. *Classification and Regression Trees*. Monterey, CA: Wadsworth; Brooks.
- Breiman, Leo. 2001a. “Random Forests.” *Machine Learning* 45 (1). Kluwer Academic Publishers, Boston: 5–32.
- . 2001b. “Statistical Modeling: The Two Cultures.” *Statistical Science* 16 (3): 199–231.
- Chambers, J. M. 1992. *Statistical Models in S*. Wadsworth & Brooks/Cole.
- Cleveland, William S. 1981. “LOWESS: A Program for Smoothing Scatterplots by Robust

- Locally Weighted Regression.” *The American Statistician* 35 (1): 54.
- . 1993. *Visualizing Data*. Summit Press.
- Cleveland, William S., and Susan J. Devlin. 1988. “Locally-Weighted Regression: An Approach to Regression Analysis by Local Fitting.” *Journal of the American Statistical Association* 83 (403): 596–610.
- Efron, Bradley, and Robert Tibshirani. 1994. *An Introduction to the Bootstrap*. Chapman & Hall/CRC.
- Fleming, Thomas R., and David P. Harrington. 1991. *Counting Processes and Survival Analysis*. John Wiley & Sons, New York.
- Friedman, Jerome H. 2000. “Greedy Function Approximation: A Gradient Boosting Machine.” *Annals of Statistics* 29: 1189–1232.
- Hastie, Trevor, Robert Tibshirani, and Jerome H. Friedman. 2009. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Second. New York: Springer-Verlag.
- Ishwaran, Hemant. 2007. “Variable Importance in Binary Regression Trees and Forests.” *Electronic Journal of Statistics* 1: 519–37.
- Ishwaran, Hemant, and Udaya B. Kogalur. 2007. “Random Survival Forests for R.” *R News* 7 (2): 25–31.
- . 2010. “Consistency of Random Survival Forests.” *Statistics and Probability Letters* 80: 1056–64.
- . 2014. “Random Forests for Survival, Regression and Classification (RF-SRC), R package version 1.6.” <http://CRAN.R-project.org/package=randomForestSRC>.
- Ishwaran, Hemant, Udaya B. Kogalur, Eugene H. Blackstone, and Michael S. Lauer. 2008. “Random Survival Forests.” *The Annals of Applied Statistics* 2 (3): 841–60.
- Ishwaran, Hemant, Udaya B. Kogalur, Xi Chen, and Andy J. Minn. 2011. “Random Survival Forests for High-Dimensional Data.” *Statist. Anal. Data Mining* 4: 115–32.
- Ishwaran, Hemant, Udaya B. Kogalur, Eiran Z. Gorodeski, Andy J. Minn, and Michael S. Lauer. 2010. “High-Dimensional Variable Selection for Survival Data.” *J. Amer. Statist. Assoc.* 105: 205–17.
- Liaw, Andy, and Matthew Wiener. 2002. “Classification and Regression by RandomForest.” *R News* 2 (3): 18–22.
- Rubin, D.B. 1976. “Inference and Missing Data.” *Biometrika* 63: 581–92.
- Tukey, John W. 1977. *Exploratory Data Analysis*. Pearson.
- Wickham, Hadley. 2009. *Ggplot2: Elegant Graphics for Data Analysis*. New York: Springer-Verlag.

Affiliation:

John Ehrlinger

Microsoft

First line Second line

E-mail: john.ehrlinger@gmail.com

URL: <https://github.com/ehrlinger/ggRandomForests>

