

# COVID STEROID 2: HTE Analysis

## Analysis

First the required R packages are loaded, as well as the data set. We also source a couple of functions to allow for running BART models with multiple MCMC chains (without needing multiple cores).

```
rm(list = ls())
library(BART)
library(caret)
library(rpart)
library(rpart.plot)
source("chainfunctions.R")

# Load data from appropriate directory (edit as needed)
dat <- read.csv2("~/Downloads/synth_covid.csv")
```

Next we do a small amount of data cleaning/preparation.

```
# Clean up data variables types and remove the small amount of missing data
dat$resp_sup <- as.factor(dat$resp_sup)
dat$dead90 <- ifelse(dat$dead90 == TRUE, 1, 0)
dat <- dat[complete.cases(dat), ]

# Standardize continuous covariates
dat$age <- (dat$age - mean(dat$age)) / sd(dat$age)
dat$BL9_Weight <- (dat$BL9_Weight - mean(dat$BL9_Weight)) / sd(dat$BL9_Weight)

# Make datasets under each counterfactual
dat1 <- dat0 <- dat
dat1$allocation <- TRUE
dat0$allocation <- FALSE
```

Then we run a BART analysis focused on the binary mortality outcome.

```
# Fit BART models under default hyperparameters, get predictions under each trt
set.seed(60622)
bartmod1 <- lbart.chained(x.train = dat[, c(1, 4:13)], y.train = dat$dead90,
                        x.test = dat1[, c(1, 4:13)], nchains = 4)
bartmod0 <- lbart.chained(x.train = dat[, c(1, 4:13)], y.train = dat$dead90,
                        x.test = dat0[, c(1, 4:13)], nchains = 4)

# Collapse predictions across chains for certain calculations
bartmod1$yhat.train.collapse <- apply(bartmod1$yhat.train, 2, rbind)
bartmod1$yhat.test.collapse <- apply(bartmod1$yhat.test, 2, rbind)
bartmod0$yhat.train.collapse <- apply(bartmod0$yhat.train, 2, rbind)
bartmod0$yhat.test.collapse <- apply(bartmod0$yhat.test, 2, rbind)
```

Then conditional average treatment effects are estimated using the predictions under each counterfactual.

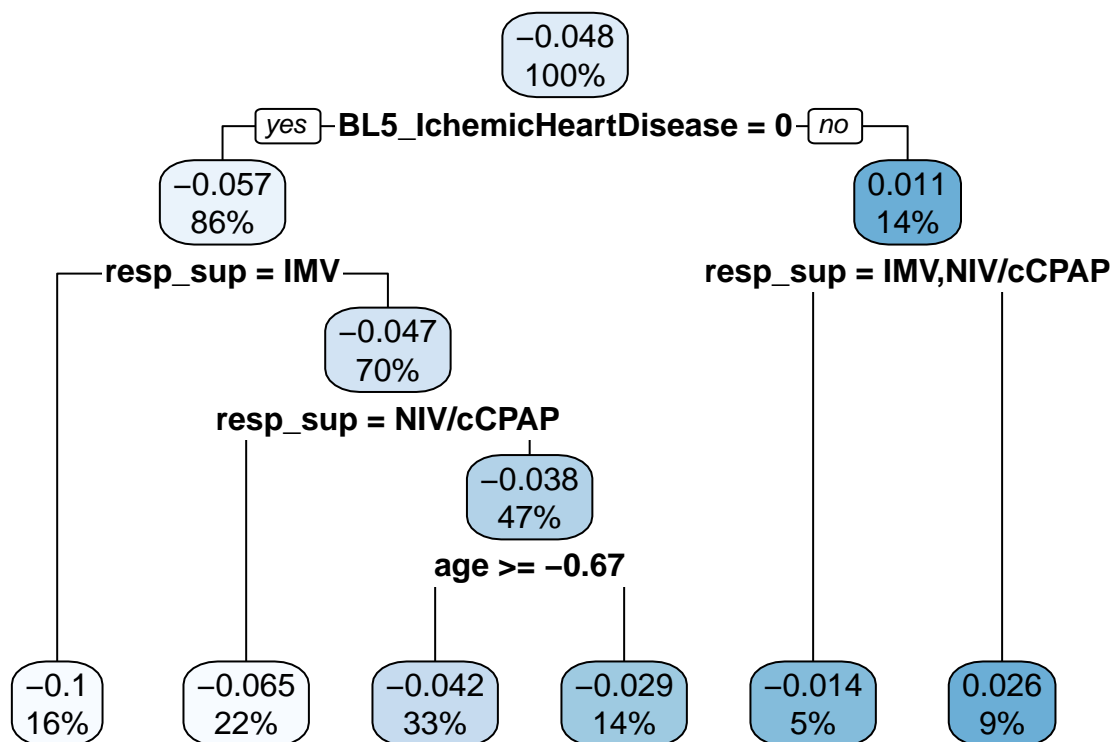
```
dat$cate <-  
  exp(colMeans(bartmod1$yhat.test.collapse)) /  
    (1 + exp(colMeans(bartmod1$yhat.test.collapse))) -  
  exp(colMeans(bartmod0$yhat.test.collapse)) /  
    (1 + exp(colMeans(bartmod0$yhat.test.collapse)))
```

This full process is then repeated for the continuous outcome (days alive without life support by day 90).

```
# Fit BART models under default hyperparameters, get predictions under each trt  
set.seed(60622)  
bartmod1_c <- wbart.chained(x.train = dat[, c(1, 4:13)],  
                           y.train = dat$dawols90,  
                           x.test = dat1[, c(1, 4:13)], nchains = 4)  
bartmod0_c <- wbart.chained(x.train = dat[, c(1, 4:13)],  
                           y.train = dat$dawols90,  
                           x.test = dat0[, c(1, 4:13)], nchains = 4)  
  
# Collapse predictions across chains for certain calculations  
bartmod1_c$yhat.train.collapse <- apply(bartmod1_c$yhat.train, 2, rbind)  
bartmod1_c$yhat.test.collapse <- apply(bartmod1_c$yhat.test, 2, rbind)  
bartmod0_c$yhat.train.collapse <- apply(bartmod0_c$yhat.train, 2, rbind)  
bartmod0_c$yhat.test.collapse <- apply(bartmod0_c$yhat.test, 2, rbind)  
  
# Estimate CATEs  
dat$cate_c <- colMeans(bartmod1_c$yhat.test.collapse) -  
  colMeans(bartmod0_c$yhat.test.collapse)
```

Finally, the “fit-the-fit” approach is used to find subgroups exhibiting heterogeneity of treatment effect, starting with the binary outcome. In particular, a CART model is fit with the CATE for 90-day mortality as the outcome and the covariates as possible predictors. The model is first fit under default CART hyperparameter settings.

```
# CART model for 90 day mortality with default CART hyperparameter and  
# all covariates considered  
cartmod <- rpart(cate ~ ., data = dat[, c(4:14)], method = "anova")  
rpart.plot(cartmod)
```



Now we prune the tree for interpretability using the stepwise approach of Hu et al. In particular, covariates are added to the CART model sequentially according to greatest increase in model  $R^2$  until an increase of less than 1% occurs.

```

# Pruned tree using stepwise approach of Hu et al.
r2 <- c()
covar_include <- c()
covar_cols <- 4:13
currentr2 <- 0
maxcovars <- 4

for (numcovars in 1:maxcovars) {

  for (covar in covar_cols) {
    cartmod <- rpart(cate ~ ., data = dat[, c(covar, covar_include, 14)],
                     method = "anova")
    r2[which(covar_cols == covar)] <-
      1 - printcp(cartmod)[dim(printcp(cartmod))[1], 3]
  }

  if ((max(r2) - currentr2) / currentr2 > 0.01) {
    covar_include <- c(covar_include, covar_cols[which.max(r2)])
    covar_cols <- covar_cols[covar_cols != covar_cols[which.max(r2)]]
    currentr2 <- max(r2)
    r2 <- c()
  } else {

```

```

    break
  }
}

```

Print out the number of covariates selected and the resulting CART model output.

```
print("Number of covariates:")
```

```
## [1] "Number of covariates:"
```

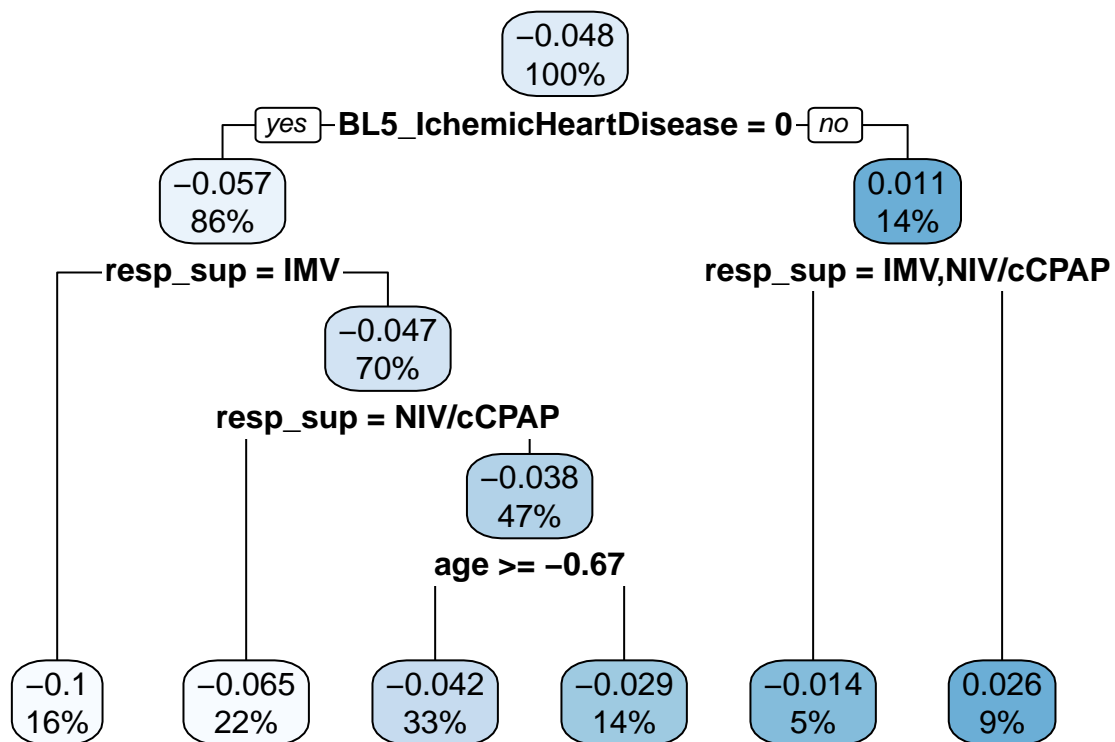
```
print(numcovars - 1)
```

```
## [1] 3
```

```

cartmod <- rpart(cate ~ ., data = dat[, c(covar_include, 14)],
                 method = "anova")
rpart.plot(cartmod)

```

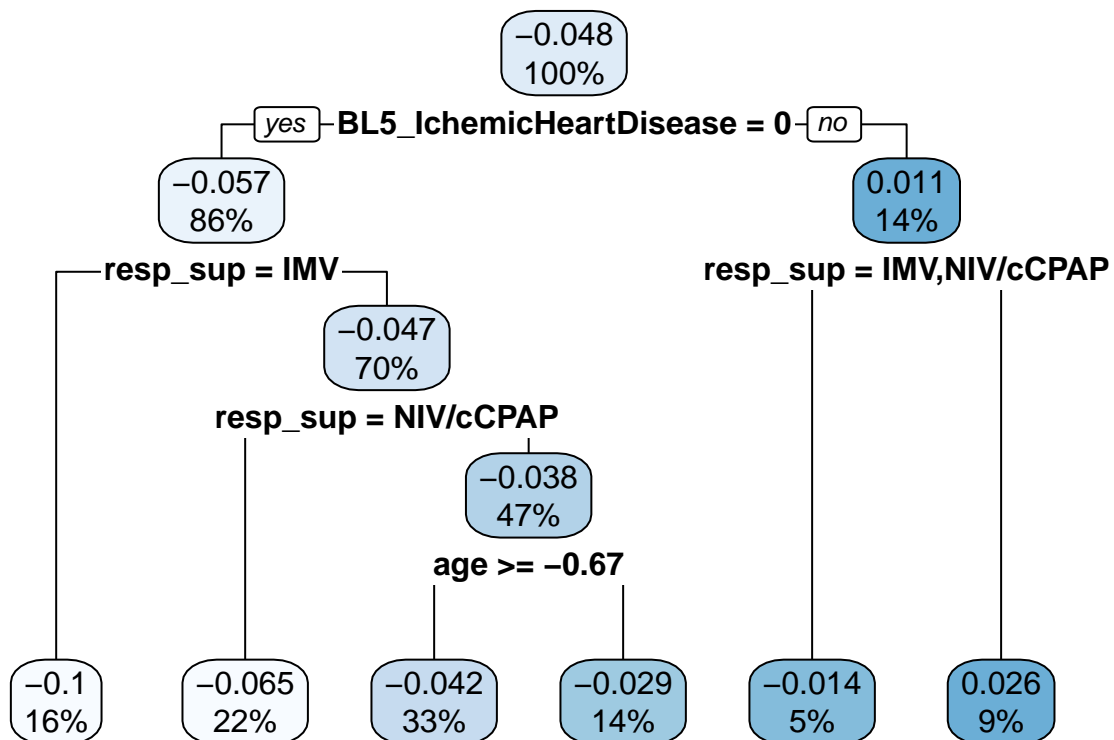


Finally we explore further pruning:

```
printcp(cartmod)
```

```
##
## Regression tree:
## rpart(formula = cate ~ ., data = dat[, c(covar_include, 14)],
##       method = "anova")
##
## Variables actually used in tree construction:
## [1] age          BL5_IchemicHeartDisease resp_sup
##
## Root node error: 1.3132/968 = 0.0013566
##
## n= 968
##
##      CP nsplit rel error  xerror    xstd
## 1 0.417395     0  1.00000 1.00124 0.0452141
## 2 0.292665     1  0.58260 0.58564 0.0236997
## 3 0.081458     2  0.28994 0.29301 0.0126259
## 4 0.039191     3  0.20848 0.21131 0.0109342
## 5 0.012666     4  0.16929 0.17161 0.0088641
## 6 0.010000     5  0.15662 0.16257 0.0093332
```

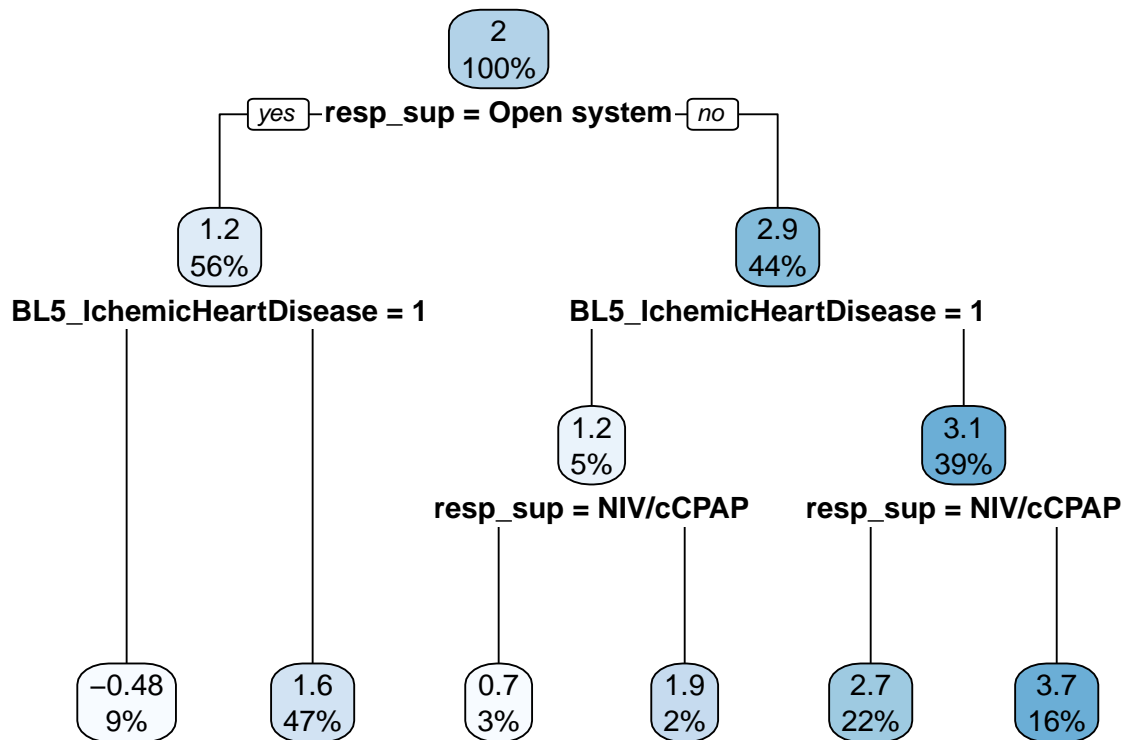
```
prunedtree <-
  prune(cartmod,
        cp = cartmod$sctable[which.min(cartmod$sctable[, "xerror"]), "CP"])
rpart.plot(prunedtree)
```



Next the same fit-the-fit approach is used to summarize the results for the continuous outcome, starting with

a CART model using all covariates and default hyperparameter.

```
# CART model for 90 day mortality with default CART hyperparameter and
# all covariates considered
cartmod <- rpart(cate_c ~ ., data = dat[, c(4:13, 15)], method = "anova")
rpart.plot(cartmod)
```



Now we prune the tree for interpretability using the stepwise approach of Hu et al.

```
# Pruned tree using stepwise approach of Hu et al.
r2 <- c()
covar_include <- c()
covar_cols <- 4:13
currentr2 <- 0
maxcovars <- 4

for (numcovars in 1:maxcovars) {

  for (covar in covar_cols) {
    cartmod <- rpart(cate_c ~ ., data = dat[, c(covar, covar_include, 15)],
                      method = "anova")
    r2[which(covar_cols == covar)] <-
      1 - printcp(cartmod)[dim(printcp(cartmod))[1], 3]
  }

  if ((max(r2) - currentr2) / currentr2 > 0.01) {
```

```

    covar_include <- c(covar_include, covar_cols[which.max(r2)])
    covar_cols <- covar_cols[covar_cols != covar_cols[which.max(r2)]]
    currentnr2 <- max(r2)
    r2 <- c()
  } else {
    break
  }
}

```

Print out the number of covariates selected and the resulting CART model output.

```
print("Number of covariates:")
```

```
## [1] "Number of covariates:"
```

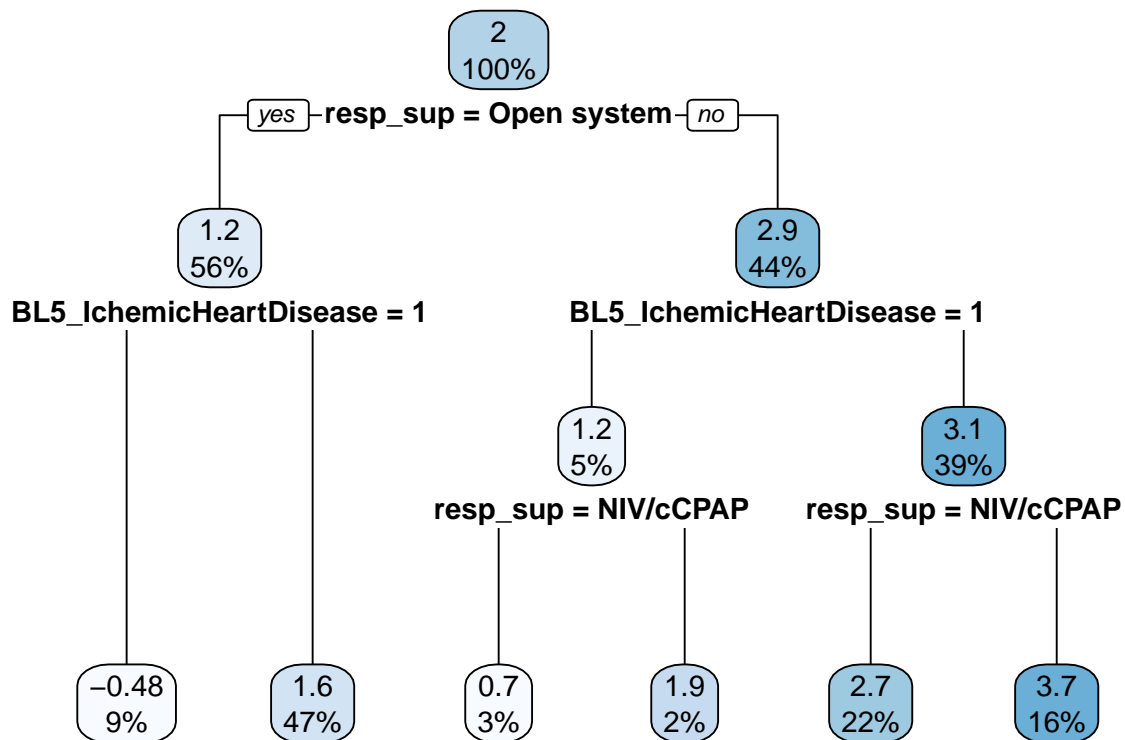
```
print(numcovars - 1)
```

```
## [1] 2
```

```

cartmod <- rpart(cate_c ~ ., data = dat[, c(covar_include, 15)],
                 method = "anova")
rpart.plot(cartmod)

```



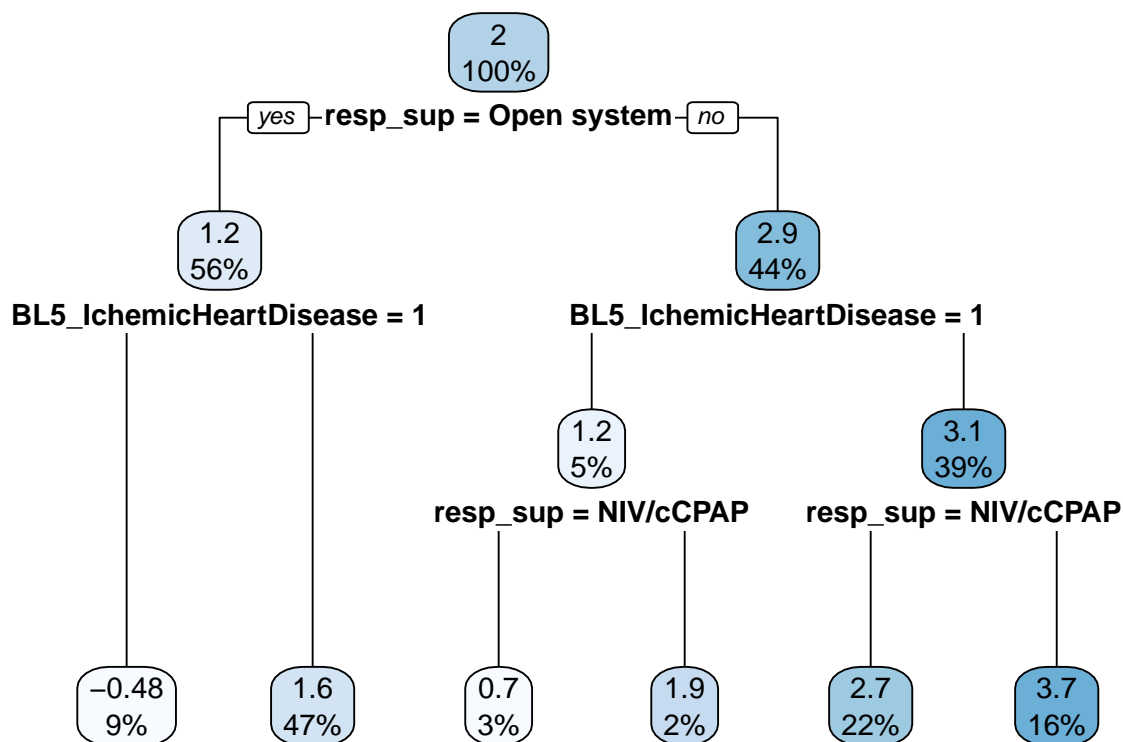
Finally we explore further pruning:

```
printcp(cartmod)
```

```
##
## Regression tree:
## rpart(formula = cate_c ~ ., data = dat[, c(covar_include, 15)],
##       method = "anova")
##
## Variables actually used in tree construction:
## [1] BL5_IchemicHeartDisease resp_sup
##
## Root node error: 1301.2/968 = 1.3442
##
## n= 968
##
##      CP nsplit rel error   xerror   xstd
## 1 0.507328     0 1.000000 1.003428 0.0455471
## 2 0.233092     1 0.492672 0.494884 0.0292256
## 3 0.135542     2 0.259580 0.260543 0.0223696
## 4 0.067051     3 0.124038 0.125918 0.0058326
## 5 0.013916     4 0.056986 0.058687 0.0039371
## 6 0.010000     5 0.043071 0.047032 0.0030947
```

```
prunedtree <-
  prune(cartmod,
        cp = cartmod$cptable[which.min(cartmod$cptable[, "xerror"]), "CP"])
rpart.plot(prunedtree)
```





Further (non-automated) pruning may need to be considered here for trees that remain too large to easily interpret.

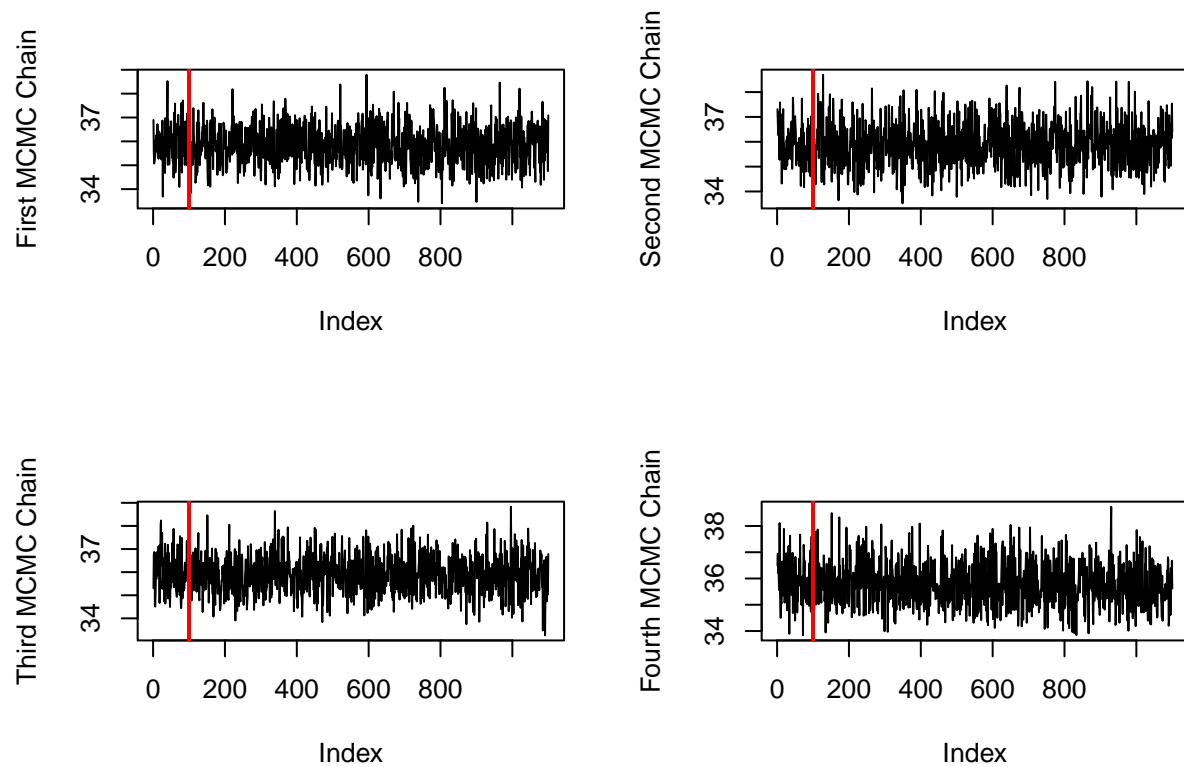
## Diagnostics

Next we run standard diagnostics for each of the models. We begin with diagnostics for the continuous outcome models, which are simpler.

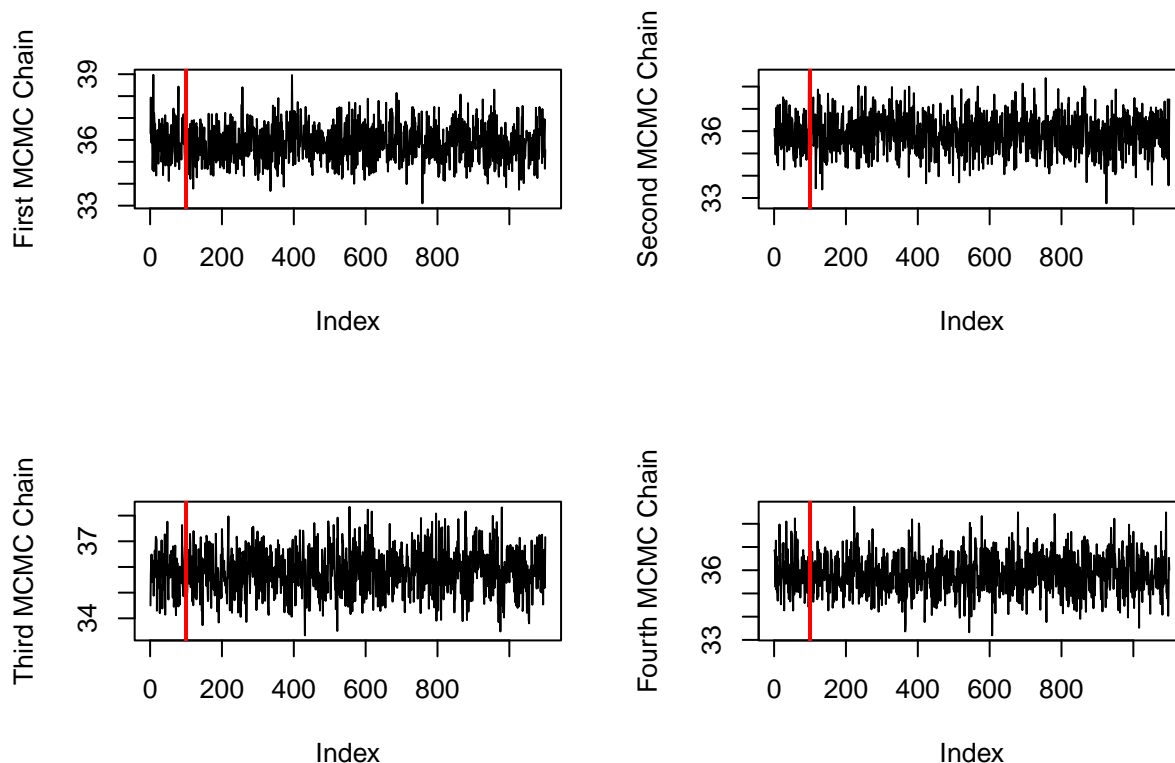
```

# MCMC chains for parameter in models of continuous outcome
par(mfrow = c(2, 2))
plot(bartmod0_c$sigma[, 1], type = "l", ylab = "First MCMC Chain")
abline(v = 100, lwd = 2, col = "red")
plot(bartmod0_c$sigma[, 2], type = "l", ylab = "Second MCMC Chain")
abline(v = 100, lwd = 2, col = "red")
plot(bartmod0_c$sigma[, 3], type = "l", ylab = "Third MCMC Chain")
abline(v = 100, lwd = 2, col = "red")
plot(bartmod0_c$sigma[, 4], type = "l", ylab = "Fourth MCMC Chain")
abline(v = 100, lwd = 2, col = "red")

```



```
plot(bartmod1_c$sigma[, 1], type = "l", ylab = "First MCMC Chain")
abline(v = 100, lwd = 2, col = "red")
plot(bartmod1_c$sigma[, 2], type = "l", ylab = "Second MCMC Chain")
abline(v = 100, lwd = 2, col = "red")
plot(bartmod1_c$sigma[, 3], type = "l", ylab = "Third MCMC Chain")
abline(v = 100, lwd = 2, col = "red")
plot(bartmod1_c$sigma[, 4], type = "l", ylab = "Fourth MCMC Chain")
abline(v = 100, lwd = 2, col = "red")
```



One should check that each chain has converged after burn-in (designated by the vertical red lines). In general, the chains should be converging to approximately the same value.

Next consider diagnostics for the models with the binary mortality outcome as described in Sparapani (2021). First consider the autocorrelation of the estimated response surface from BART from 10 randomly selected subjects. This may start somewhat correlated for nearby observations, but should reduce to 0 correlation for observations further apart.

```
# First for bartmod0, one panel for each chain
par(mfrow = c(2, 2))

auto.corr <- acf(bartmod0$yhat.train[ , sample(1:dim(dat)[1], 10), 1],
                 plot = FALSE)
max.lag <- max(auto.corr$lag[ , 1, 1])

j <- seq(-0.5, 0.4, length.out = 10)
for (h in 1:10) {
  if (h == 1) {
    plot(1:max.lag + j[h], auto.corr$acf[1 + (1:max.lag), h, h],
         type = 'h', xlim = c(0, max.lag + 1), ylim = c(-1, 1),
         ylab = 'acf', xlab = 'lag')
  } else {
    lines(1:max.lag + j[h], auto.corr$acf[1 + (1:max.lag), h, h],
          type = 'h', col = h)
  }
}
```

```

auto.corr <- acf(bartmod0$yhat.train[ , sample(1:dim(dat)[1], 10), 2],
                plot = FALSE)
max.lag <- max(auto.corr$lag[ , 1, 1])

j <- seq(-0.5, 0.4, length.out = 10)
for (h in 1:10) {
  if (h == 1) {
    plot(1:max.lag + j[h], auto.corr$acf[1 + (1:max.lag), h, h],
         type = 'h', xlim = c(0, max.lag + 1), ylim = c(-1, 1),
         ylab = 'acf', xlab = 'lag')
  } else {
    lines(1:max.lag + j[h], auto.corr$acf[1 + (1:max.lag), h, h],
         type = 'h', col = h)
  }
}

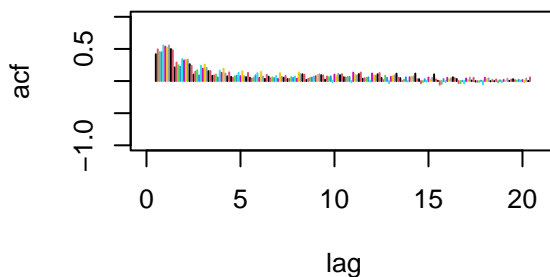
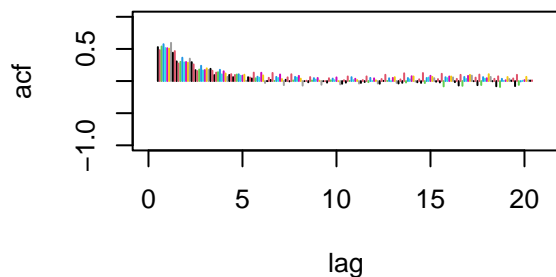
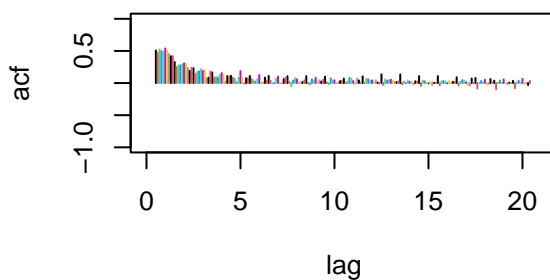
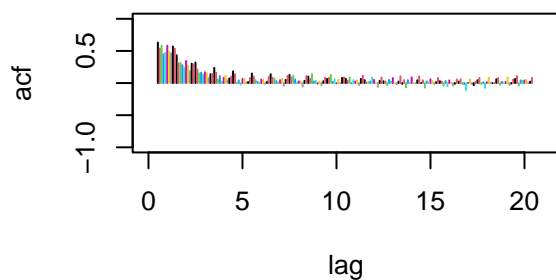
auto.corr <- acf(bartmod0$yhat.train[ , sample(1:dim(dat)[1], 10), 3],
                plot = FALSE)
max.lag <- max(auto.corr$lag[ , 1, 1])

j <- seq(-0.5, 0.4, length.out = 10)
for (h in 1:10) {
  if (h == 1) {
    plot(1:max.lag + j[h], auto.corr$acf[1 + (1:max.lag), h, h],
         type = 'h', xlim = c(0, max.lag + 1), ylim = c(-1, 1),
         ylab = 'acf', xlab = 'lag')
  } else {
    lines(1:max.lag + j[h], auto.corr$acf[1 + (1:max.lag), h, h],
         type = 'h', col = h)
  }
}

auto.corr <- acf(bartmod0$yhat.train[ , sample(1:dim(dat)[1], 10), 4],
                plot = FALSE)
max.lag <- max(auto.corr$lag[ , 1, 1])

j <- seq(-0.5, 0.4, length.out = 10)
for (h in 1:10) {
  if (h == 1) {
    plot(1:max.lag + j[h], auto.corr$acf[1 + (1:max.lag), h, h],
         type = 'h', xlim = c(0, max.lag + 1), ylim = c(-1, 1),
         ylab = 'acf', xlab = 'lag')
  } else {
    lines(1:max.lag + j[h], auto.corr$acf[1 + (1:max.lag), h, h],
         type = 'h', col = h)
  }
}

```



```
# Then for bartmod1
auto.corr <- acf(bartmod1$yhat.train[ , sample(1:dim(dat)[1], 10), 1],
                 plot = FALSE)
max.lag <- max(auto.corr$lag[ , 1, 1])

j <- seq(-0.5, 0.4, length.out = 10)
for (h in 1:10) {
  if (h == 1) {
    plot(1:max.lag + j[h], auto.corr$acf[1 + (1:max.lag), h, h],
         type = 'h', xlim = c(0, max.lag + 1), ylim = c(-1, 1),
         ylab = 'acf', xlab = 'lag')
  } else {
    lines(1:max.lag + j[h], auto.corr$acf[1 + (1:max.lag), h, h],
         type = 'h', col = h)
  }
}

auto.corr <- acf(bartmod1$yhat.train[ , sample(1:dim(dat)[1], 10), 2],
                 plot = FALSE)
max.lag <- max(auto.corr$lag[ , 1, 1])

j <- seq(-0.5, 0.4, length.out = 10)
for (h in 1:10) {
  if (h == 1) {
    plot(1:max.lag + j[h], auto.corr$acf[1 + (1:max.lag), h, h],
         type = 'h', xlim = c(0, max.lag + 1), ylim = c(-1, 1),
```

```

        ylab = 'acf', xlab = 'lag')
    } else {
        lines(1:max.lag + j[h], auto.corr$acf[1 + (1:max.lag), h, h],
              type = 'h', col = h)
    }
}

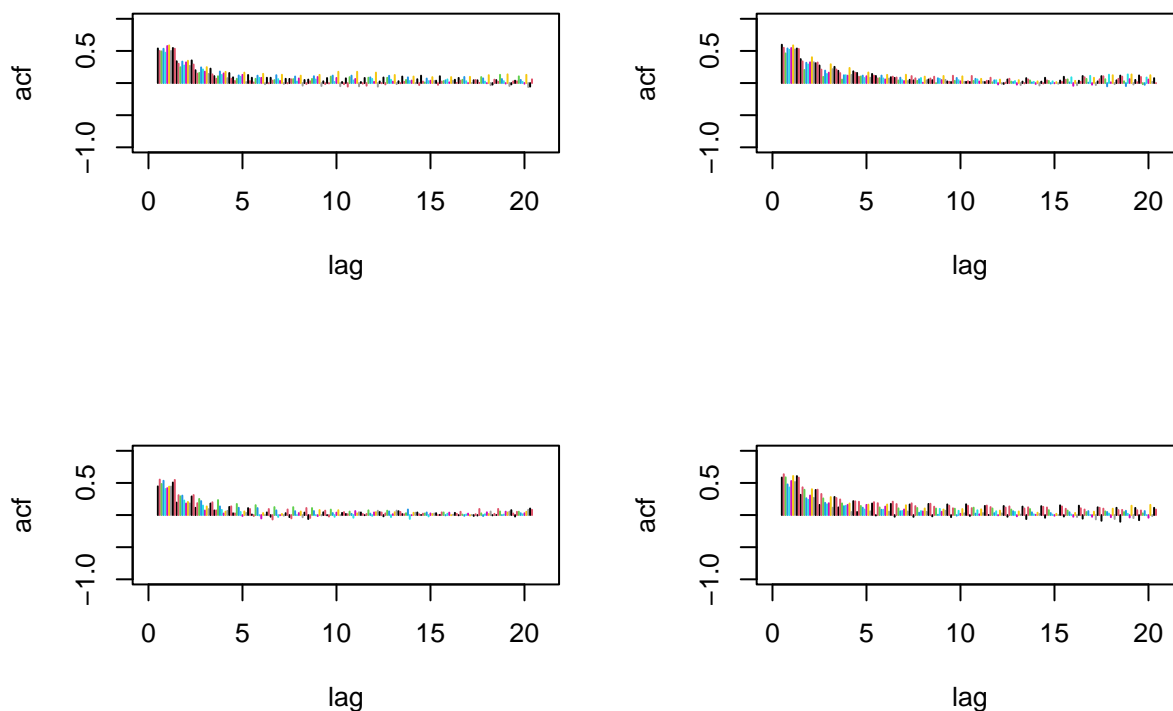
auto.corr <- acf(bartmod1$yhat.train[ , sample(1:dim(dat)[1], 10), 3],
                 plot = FALSE)
max.lag <- max(auto.corr$lag[ , 1, 1])

j <- seq(-0.5, 0.4, length.out = 10)
for (h in 1:10) {
    if (h == 1) {
        plot(1:max.lag + j[h], auto.corr$acf[1 + (1:max.lag), h, h],
              type = 'h', xlim = c(0, max.lag + 1), ylim = c(-1, 1),
              ylab = 'acf', xlab = 'lag')
    } else {
        lines(1:max.lag + j[h], auto.corr$acf[1 + (1:max.lag), h, h],
              type = 'h', col = h)
    }
}

auto.corr <- acf(bartmod1$yhat.train[ , sample(1:dim(dat)[1], 10), 4],
                 plot = FALSE)
max.lag <- max(auto.corr$lag[ , 1, 1])

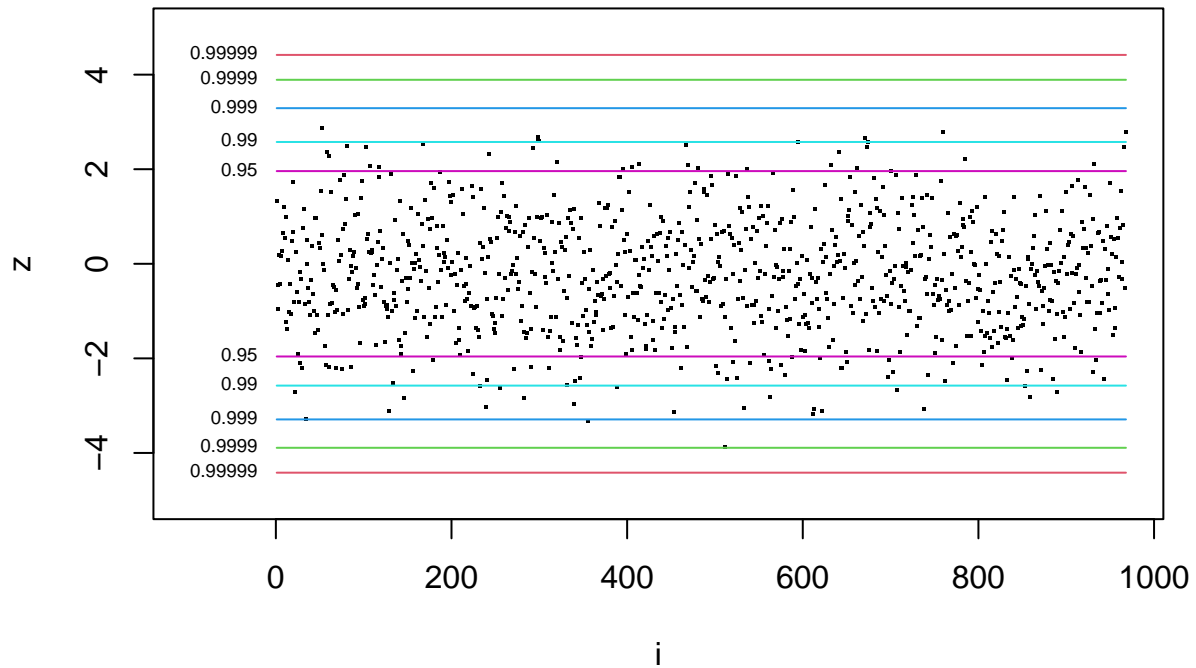
j <- seq(-0.5, 0.4, length.out = 10)
for (h in 1:10) {
    if (h == 1) {
        plot(1:max.lag + j[h], auto.corr$acf[1 + (1:max.lag), h, h],
              type = 'h', xlim = c(0, max.lag + 1), ylim = c(-1, 1),
              ylab = 'acf', xlab = 'lag')
    } else {
        lines(1:max.lag + j[h], auto.corr$acf[1 + (1:max.lag), h, h],
              type = 'h', col = h)
    }
}

```



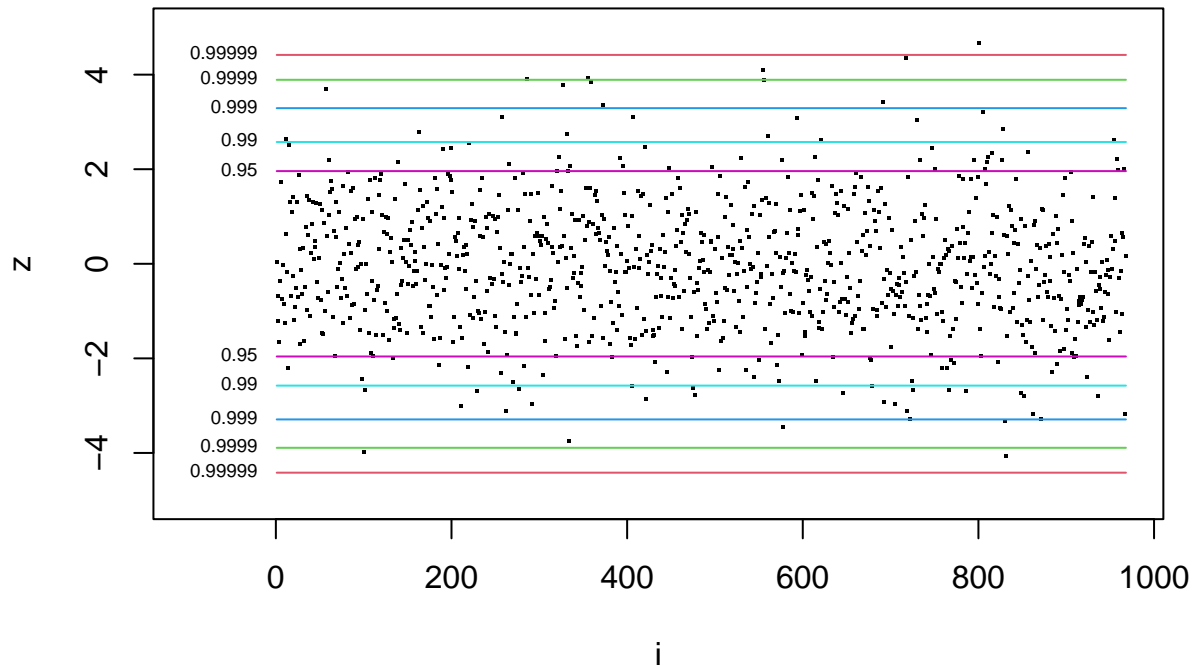
Next, we plot the Geweke Z statistics for each individual, which should be approximately distributed as a standard Normal.

```
# First for bartmod0
geweke <- gewekediag(bartmod0$yhat.train.collapse)
n <- dim(dat)[1]
j <- -10^(log10(n) - 1)
plot(geweke$z, pch = '.', cex = 2, ylab = 'z', xlab = 'i',
     xlim=c(j, n), ylim=c(-5, 5))
lines(1:n, rep(-1.96, n), type='l', col=6)
lines(1:n, rep(+1.96, n), type='l', col=6)
lines(1:n, rep(-2.576, n), type='l', col=5)
lines(1:n, rep(+2.576, n), type='l', col=5)
lines(1:n, rep(-3.291, n), type='l', col=4)
lines(1:n, rep(+3.291, n), type='l', col=4)
lines(1:n, rep(-3.891, n), type='l', col=3)
lines(1:n, rep(+3.891, n), type='l', col=3)
lines(1:n, rep(-4.417, n), type='l', col=2)
lines(1:n, rep(+4.417, n), type='l', col=2)
text(c(1, 1), c(-1.96, 1.96), pos=2, cex=0.6, labels='0.95')
text(c(1, 1), c(-2.576, 2.576), pos=2, cex=0.6, labels='0.99')
text(c(1, 1), c(-3.291, 3.291), pos=2, cex=0.6, labels='0.999')
text(c(1, 1), c(-3.891, 3.891), pos=2, cex=0.6, labels='0.9999')
text(c(1, 1), c(-4.417, 4.417), pos=2, cex=0.6, labels='0.99999')
```



```
# Then for bartmod1
geweke <- gewekediag(bartmod1$yhat.train.collapse)
plot(geweke$z, pch = '.', cex = 2, ylab = 'z', xlab = 'i',
      xlim=c(j, n), ylim=c(-5, 5))
lines(1:n, rep(-1.96, n), type='l', col=6)
lines(1:n, rep(+1.96, n), type='l', col=6)
lines(1:n, rep(-2.576, n), type='l', col=5)
lines(1:n, rep(+2.576, n), type='l', col=5)
lines(1:n, rep(-3.291, n), type='l', col=4)
lines(1:n, rep(+3.291, n), type='l', col=4)
lines(1:n, rep(-3.891, n), type='l', col=3)
lines(1:n, rep(+3.891, n), type='l', col=3)
lines(1:n, rep(-4.417, n), type='l', col=2)
lines(1:n, rep(+4.417, n), type='l', col=2)
text(c(1, 1), c(-1.96, 1.96), pos=2, cex=0.6, labels='0.95')
text(c(1, 1), c(-2.576, 2.576), pos=2, cex=0.6, labels='0.99')
text(c(1, 1), c(-3.291, 3.291), pos=2, cex=0.6, labels='0.999')
text(c(1, 1), c(-3.891, 3.891), pos=2, cex=0.6, labels='0.9999')
text(c(1, 1), c(-4.417, 4.417), pos=2, cex=0.6, labels='0.99999')
```





If several points lie beyond the dark blue line or further, consider using more thinning when fitting the models.

## Sensitivity analysis

We conduct two simple sensitivity analyses, one assuming that all missing mortality data correspond to an alive status and one assuming that all missing mortality data correspond to a deceased status. In the former scenario, missing days without life support will be set to a random sample with replacement from the observed days without life support among those known to be alive at day 90; in the latter scenario it will be set to a random sample with replacement from the observed days without life support among those not alive at day 90.

Note that these analyses are unlikely to output the exact same trees as the original analysis, but they should hopefully output trees which largely tell the same HTE story as the original analysis.

### Sensitivity analysis 1: Impute missing outcomes as alive

First the memory is cleared and the data set is reloaded.

```
rm(list = ls())
library(BART)
library(caret)
library(rpart)
library(rpart.plot)
source("chainfunctions.R")
```

```
# Load data from appropriate directory (edit as needed)
dat <- read.csv2("~/Downloads/synth_covid.csv")
```

Next we do a small amount of data cleaning/preparation.

```
# Clean up data variables types and impute the small amount of missing data
dat$resp_sup <- as.factor(dat$resp_sup)
dat$dead90 <- ifelse(dat$dead90 == TRUE, 1, 0)
dat$dawols90[is.na(dat$dawols90)] <-
  sample(dat$dawols90[!is.na(dat$dawols90) & dat$dead90 == 0],
    sum(is.na(dat$dawols90)), replace = TRUE)
dat$dead90[is.na(dat$dead90)] <- 0

# Standardize continuous covariates
dat$age <- (dat$age - mean(dat$age)) / sd(dat$age)
dat$BL9_Weight <- (dat$BL9_Weight - mean(dat$BL9_Weight)) / sd(dat$BL9_Weight)

# Make datasets under each counterfactual
dat1 <- dat0 <- dat
dat1$allocation <- TRUE
dat0$allocation <- FALSE
```

Then we run a BART analysis focused on the binary mortality outcome.

```
# Fit BART models under default hyperparameters, get predictions under each trt
set.seed(60622)
bartmod1 <- lbart.chained(x.train = dat[, c(1, 4:13)], y.train = dat$dead90,
  x.test = dat1[, c(1, 4:13)], nchains = 4)
bartmod0 <- lbart.chained(x.train = dat[, c(1, 4:13)], y.train = dat$dead90,
  x.test = dat0[, c(1, 4:13)], nchains = 4)

# Collapse predictions across chains for certain calculations
bartmod1$yhat.train.collapse <- apply(bartmod1$yhat.train, 2, rbind)
bartmod1$yhat.test.collapse <- apply(bartmod1$yhat.test, 2, rbind)
bartmod0$yhat.train.collapse <- apply(bartmod0$yhat.train, 2, rbind)
bartmod0$yhat.test.collapse <- apply(bartmod0$yhat.test, 2, rbind)
```

Then conditional average treatment effects are estimated using predictions under each counterfactual.

```
dat$cate <-
  exp(colMeans(bartmod1$yhat.test.collapse)) /
    (1 + exp(colMeans(bartmod1$yhat.test.collapse))) -
  exp(colMeans(bartmod0$yhat.test.collapse)) /
    (1 + exp(colMeans(bartmod0$yhat.test.collapse)))
```

This full process is then repeated for the continuous outcome (days alive without life support by day 90).

```
# Fit BART models under default hyperparameters, get predictions under each trt
set.seed(60622)
bartmod1_c <- wbart.chained(x.train = dat[, c(1, 4:13)],
  y.train = dat$dawols90,
```

```

x.test = dat1[, c(1, 4:13)], nchains = 4)
bartmod0_c <- wbart.chained(x.train = dat[, c(1, 4:13)],
y.train = dat$dawols90,
x.test = dat0[, c(1, 4:13)], nchains = 4)

# Collapse predictions across chains for certain calculations
bartmod1_c$yhat.train.collapse <- apply(bartmod1_c$yhat.train, 2, rbind)
bartmod1_c$yhat.test.collapse <- apply(bartmod1_c$yhat.test, 2, rbind)
bartmod0_c$yhat.train.collapse <- apply(bartmod0_c$yhat.train, 2, rbind)
bartmod0_c$yhat.test.collapse <- apply(bartmod0_c$yhat.test, 2, rbind)

# Estimate CATEs
dat$cate_c <- colMeans(bartmod1_c$yhat.test.collapse) -
colMeans(bartmod0_c$yhat.test.collapse)

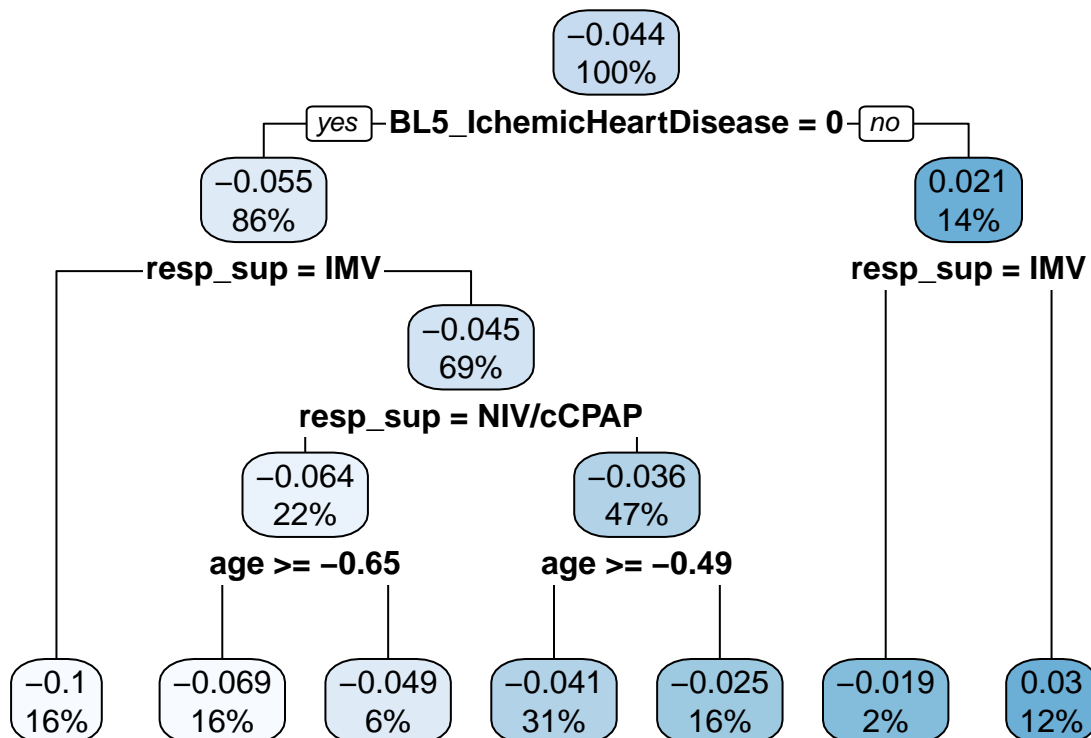
```

Finally, the “fit-the-fit” approach is used to find subgroups exhibiting heterogeneity of treatment effect, starting with the binary outcome. In particular, a CART model is fit with the CATE for 90-day mortality as the outcome and the covariates as possible predictors. The model is first fit under default CART hyperparameter settings.

```

# CART model for 90 day mortality with default CART hyperparameter and
# all covariates considered
cartmod <- rpart(cate ~ ., data = dat[, c(4:14)], method = "anova")
rpart.plot(cartmod)

```



Now we prune the tree for interpretability using the stepwise approach of Hu et al. In particular, covariates are added to the CART model sequentially according to greatest increase in model  $R^2$  until an increase of less than 1% occurs.

```
# Pruned tree using stepwise approach of Hu et al.
r2 <- c()
covar_include <- c()
covar_cols <- 4:13
currentr2 <- 0
maxcovars <- 4

for (numcovars in 1:maxcovars) {

  for (covar in covar_cols) {
    cartmod <- rpart(cate ~ ., data = dat[, c(covar, covar_include, 14)],
                     method = "anova")
    r2[which(covar_cols == covar)] <-
      1 - printcp(cartmod)[dim(printcp(cartmod))[1], 3]
  }

  if ((max(r2) - currentr2) / currentr2 > 0.01) {
    covar_include <- c(covar_include, covar_cols[which.max(r2)])
    covar_cols <- covar_cols[covar_cols != covar_cols[which.max(r2)]]
    currentr2 <- max(r2)
    r2 <- c()
  } else {
    break
  }
}
```

Print out the number of covariates selected and the resulting CART model output.

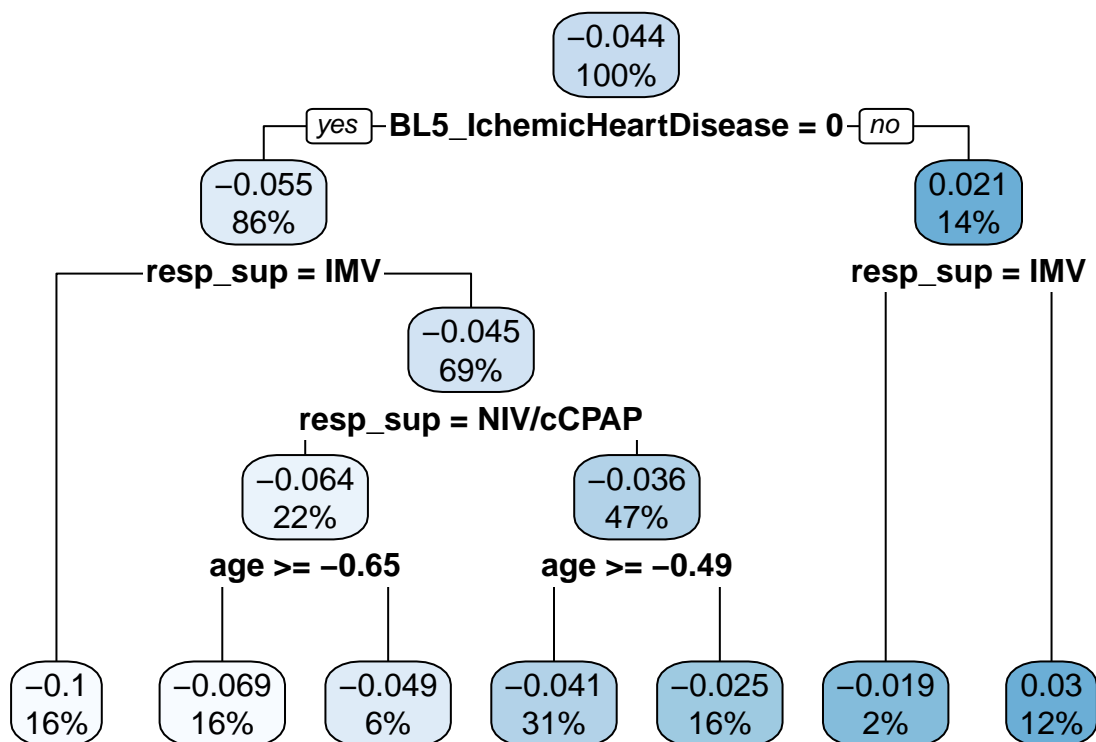
```
print("Number of covariates:")
```

```
## [1] "Number of covariates:"
```

```
print(numcovars - 1)
```

```
## [1] 3
```

```
cartmod <- rpart(cate ~ ., data = dat[, c(covar_include, 14)],
                 method = "anova")
rpart.plot(cartmod)
```



Finally we explore further pruning:

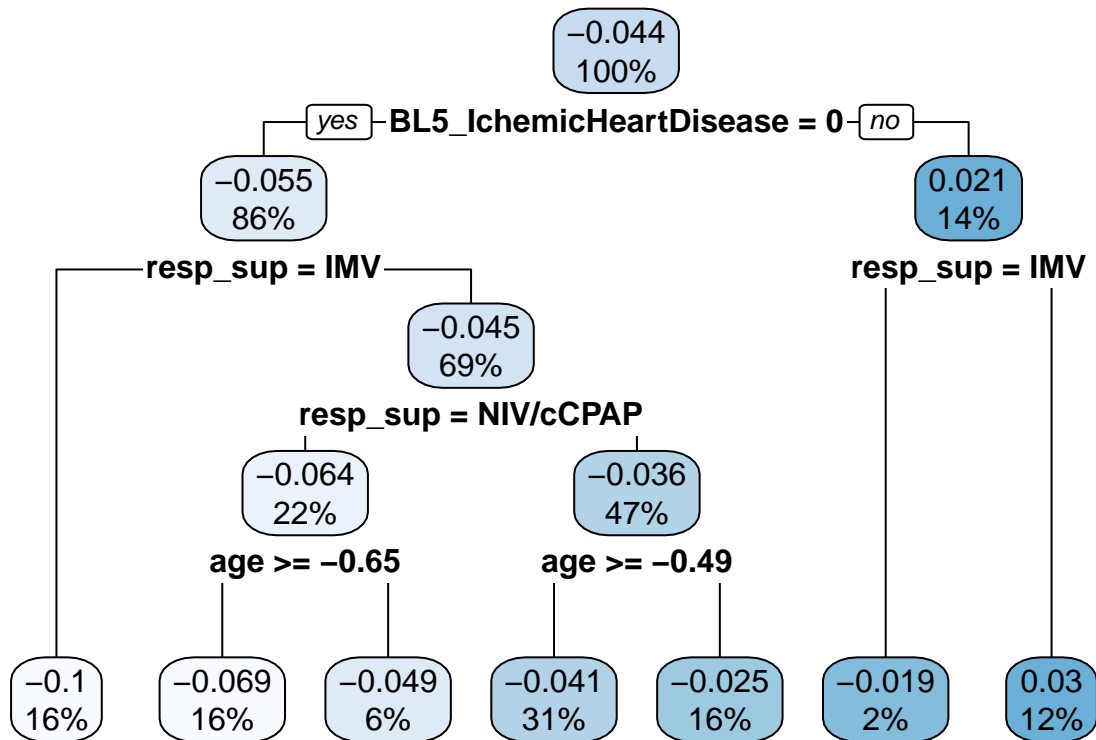
```
printcp(cartmod)
```

```
##
## Regression tree:
## rpart(formula = cate ~ ., data = dat[, c(covar_include, 14)],
##       method = "anova")
##
## Variables actually used in tree construction:
## [1] age          BL5_IchemicHeartDisease resp_sup
##
## Root node error: 1.5163/982 = 0.0015441
##
## n= 982
##
##      CP nsplit rel error  xerror    xstd
## 1 0.472054      0  1.00000 1.00061 0.0453561
## 2 0.259996      1  0.52795 0.52967 0.0213342
## 3 0.075276      2  0.26795 0.26929 0.0113978
## 4 0.031268      3  0.19267 0.19433 0.0094964
## 5 0.018728      4  0.16141 0.16259 0.0076075
## 6 0.011530      5  0.14268 0.14503 0.0075145
## 7 0.010000      6  0.13115 0.14095 0.0078628
```

```

prunedtree <-
  prune(cartmod,
        cp = cartmod$cptable[which.min(cartmod$cptable[, "xerror"]), "CP"])
rpart.plot(prunedtree)

```

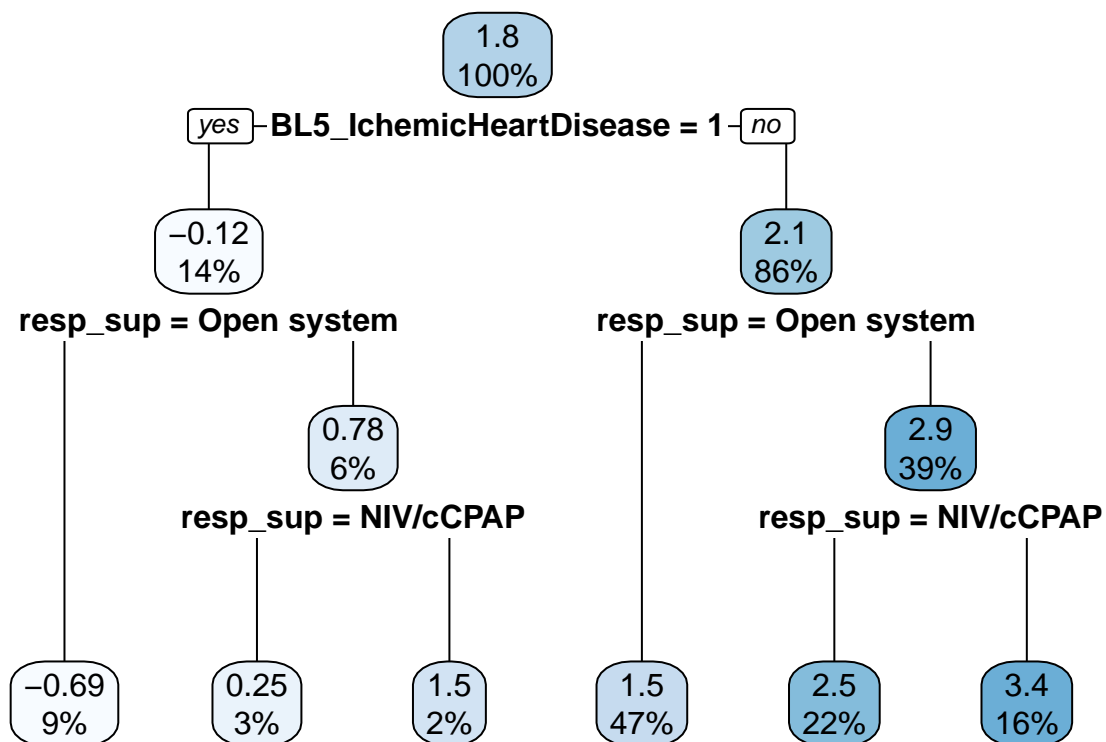


Next the same fit-the-fit approach is used to summarize the results for the continuous outcome, starting with a CART model using all covariates and default hyperparameter.

```

# CART model for 90 day mortality with default CART hyperparameter and
# all covariates considered
cartmod <- rpart(cate_c ~ ., data = dat[, c(4:13, 15)], method = "anova")
rpart.plot(cartmod)

```



Now we prune the tree for interpretability using the stepwise approach of Hu et al.

```

# Pruned tree using stepwise approach of Hu et al.
r2 <- c()
covar_include <- c()
covar_cols <- 4:13
currentr2 <- 0
maxcovars <- 4

for (numcovars in 1:maxcovars) {

  for (covar in covar_cols) {
    cartmod <- rpart(cate_c ~ ., data = dat[, c(covar, covar_include, 15)],
                     method = "anova")
    r2[which(covar_cols == covar)] <-
      1 - printcp(cartmod)[dim(printcp(cartmod))[1], 3]
  }

  if ((max(r2) - currentr2) / currentr2 > 0.01) {
    covar_include <- c(covar_include, covar_cols[which.max(r2)])
    covar_cols <- covar_cols[covar_cols != covar_cols[which.max(r2)]]
    currentr2 <- max(r2)
    r2 <- c()
  } else {
    break
  }
}

```

```
}
```

Print out the number of covariates selected and the resulting CART model output.

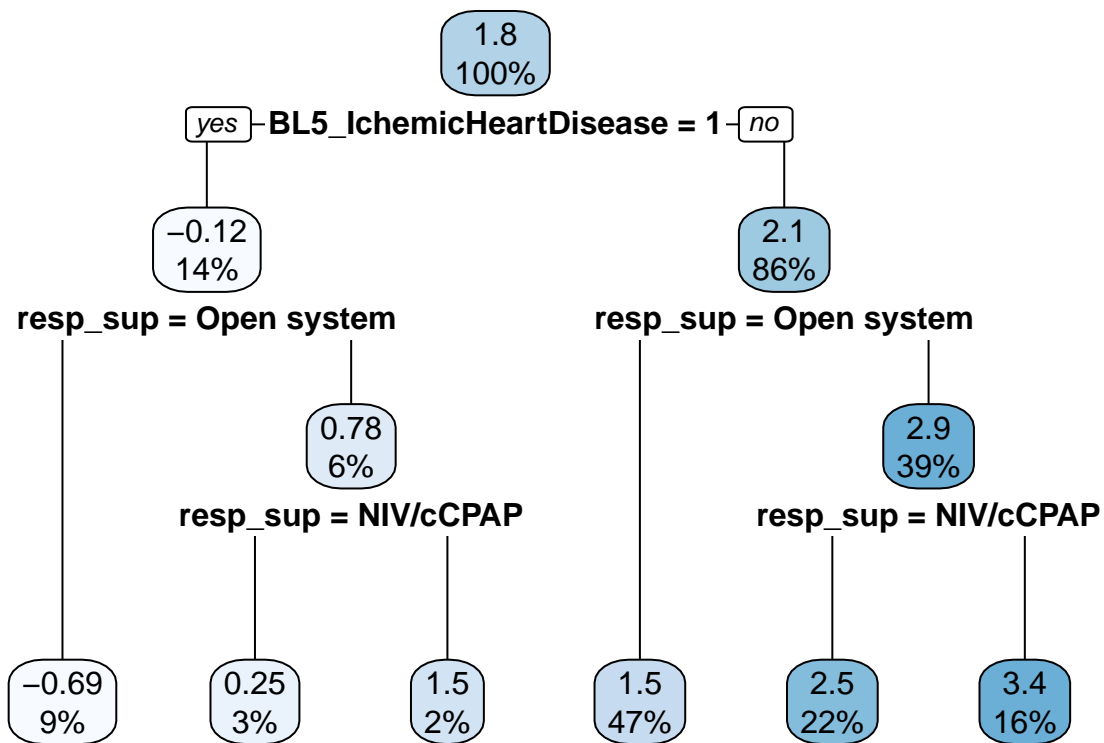
```
print("Number of covariates:")
```

```
## [1] "Number of covariates:"
```

```
print(numcovars - 1)
```

```
## [1] 2
```

```
cartmod <- rpart(cate_c ~ ., data = dat[, c(covar_include, 15)],  
                 method = "anova")  
rpart.plot(cartmod)
```



Finally we explore further pruning:

```
printcp(cartmod)
```

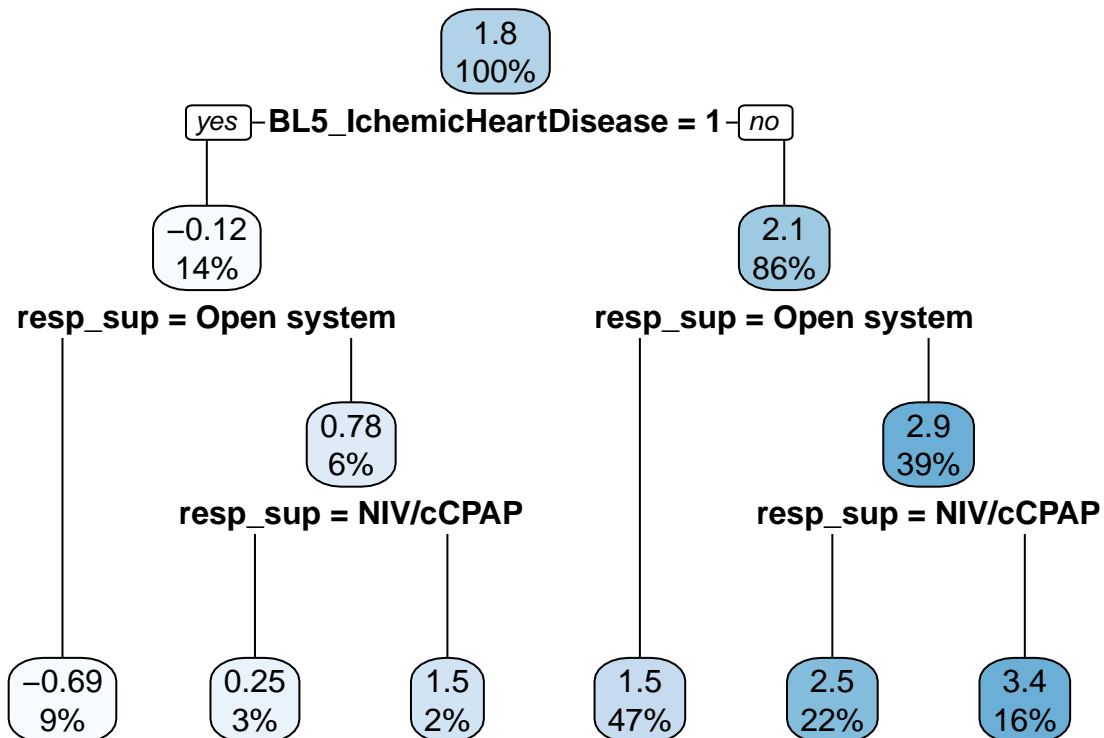
```
##
```

```
## Regression tree:
```



```
## rpart(formula = cate_c ~ ., data = dat[, c(covar_include, 15)],
##       method = "anova")
##
## Variables actually used in tree construction:
## [1] BL5_IchemicHeartDisease resp_sup
##
## Root node error: 1244.1/982 = 1.2669
##
## n= 982
##
##      CP nsplit rel error   xerror   xstd
## 1 0.495222    0 1.000000 1.001577 0.0485097
## 2 0.322680    1 0.504778 0.508665 0.0179710
## 3 0.060805    2 0.182098 0.185450 0.0122326
## 4 0.058695    3 0.121293 0.124356 0.0109316
## 5 0.016008    4 0.062598 0.072427 0.0051450
## 6 0.010000    5 0.046590 0.047438 0.0025556
```

```
prunedtree <-
  prune(cartmod,
        cp = cartmod$cpstable[which.min(cartmod$cpstable[, "xerror"]), "CP"])
rpart.plot(prunedtree)
```



## Sensitivity analysis 2: Impute missing outcomes as deceased

First the memory is cleared and the data set is reloaded.

```
rm(list = ls())
library(BART)
library(caret)
library(rpart)
library(rpart.plot)
source("chainfunctions.R")

# Load data from appropriate directory (edit as needed)
dat <- read.csv2("~/Downloads/synth_covid.csv")
```

Next we do a small amount of data cleaning/preparation.

```
# Clean up data variables types and impute the small amount of missing data
dat$resp_sup <- as.factor(dat$resp_sup)
dat$dead90 <- ifelse(dat$dead90 == TRUE, 1, 0)
dat$dawols90[is.na(dat$dawols90)] <-
  sample(dat$dawols90[!is.na(dat$dawols90) & dat$dead90 == 1],
    sum(is.na(dat$dawols90)), replace = TRUE)
dat$dead90[is.na(dat$dead90)] <- 1

# Standardize continuous covariates
dat$age <- (dat$age - mean(dat$age)) / sd(dat$age)
dat$BL9_Weight <- (dat$BL9_Weight - mean(dat$BL9_Weight)) / sd(dat$BL9_Weight)

# Make datasets under each counterfactual
dat1 <- dat0 <- dat
dat1$allocation <- TRUE
dat0$allocation <- FALSE
```

Then we run a BART analysis focused on the binary mortality outcome.

```
# Fit BART models under default hyperparameters, get predictions under each trt
set.seed(60622)
bartmod1 <- lbart.chained(x.train = dat[, c(1, 4:13)], y.train = dat$dead90,
  x.test = dat1[, c(1, 4:13)], nchains = 4)
bartmod0 <- lbart.chained(x.train = dat[, c(1, 4:13)], y.train = dat$dead90,
  x.test = dat0[, c(1, 4:13)], nchains = 4)

# Collapse predictions across chains for certain calculations
bartmod1$yhat.train.collapse <- apply(bartmod1$yhat.train, 2, rbind)
bartmod1$yhat.test.collapse <- apply(bartmod1$yhat.test, 2, rbind)
bartmod0$yhat.train.collapse <- apply(bartmod0$yhat.train, 2, rbind)
bartmod0$yhat.test.collapse <- apply(bartmod0$yhat.test, 2, rbind)
```

Then conditional average treatment effects are estimated using predictions under each counterfactual.

```
dat$cate <-
  exp(colMeans(bartmod1$yhat.test.collapse)) /
  (1 + exp(colMeans(bartmod1$yhat.test.collapse))) -
```

```
exp(colMeans(bartmod0$yhat.test.collapse)) /
  (1 + exp(colMeans(bartmod0$yhat.test.collapse)))
```

This full process is then repeated for the continuous outcome (days alive without life support by day 90).

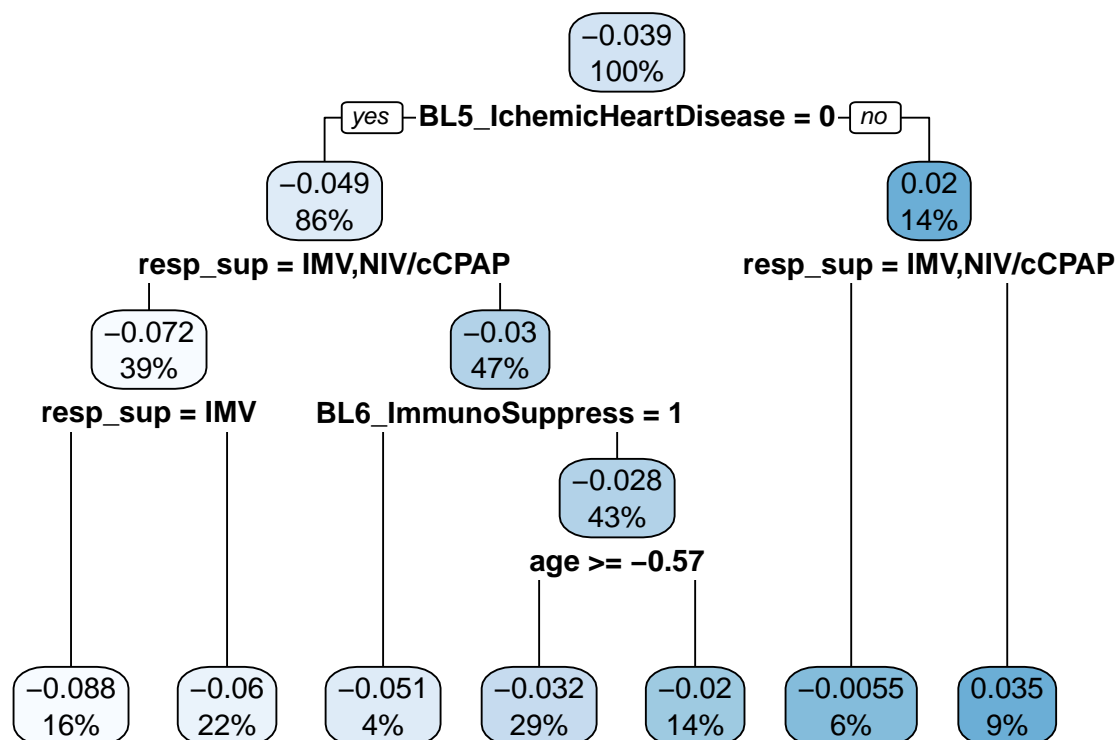
```
# Fit BART models under default hyperparameters, get predictions under each trt
set.seed(60622)
bartmod1_c <- wbart.chained(x.train = dat[, c(1, 4:13)],
                           y.train = dat$dawols90,
                           x.test = dat1[, c(1, 4:13)], nchains = 4)
bartmod0_c <- wbart.chained(x.train = dat[, c(1, 4:13)],
                           y.train = dat$dawols90,
                           x.test = dat0[, c(1, 4:13)], nchains = 4)

# Collapse predictions across chains for certain calculations
bartmod1_c$yhat.train.collapse <- apply(bartmod1_c$yhat.train, 2, rbind)
bartmod1_c$yhat.test.collapse <- apply(bartmod1_c$yhat.test, 2, rbind)
bartmod0_c$yhat.train.collapse <- apply(bartmod0_c$yhat.train, 2, rbind)
bartmod0_c$yhat.test.collapse <- apply(bartmod0_c$yhat.test, 2, rbind)

# Estimate CATEs
dat$cate_c <- colMeans(bartmod1_c$yhat.test.collapse) -
  colMeans(bartmod0_c$yhat.test.collapse)
```

Finally, the “fit-the-fit” approach is used to find subgroups exhibiting heterogeneity of treatment effect, starting with the binary outcome. In particular, a CART model is fit with the CATE for 90-day mortality as the outcome and the covariates as possible predictors. The model is first fit under default CART hyperparameter settings.

```
# CART model for 90 day mortality with default CART hyperparameter and
# all covariates considered
cartmod <- rpart(cate ~ ., data = dat[, c(4:14)], method = "anova")
rpart.plot(cartmod)
```



Now we prune the tree for interpretability using the stepwise approach of Hu et al. In particular, covariates are added to the CART model sequentially according to greatest increase in model  $R^2$  until an increase of less than 1% occurs.

```

# Pruned tree using stepwise approach of Hu et al.
r2 <- c()
covar_include <- c()
covar_cols <- 4:13
currentr2 <- 0
maxcovars <- 4

for (numcovars in 1:maxcovars) {

  for (covar in covar_cols) {
    cartmod <- rpart(cate ~ ., data = dat[, c(covar, covar_include, 14)],
                     method = "anova")
    r2[which(covar_cols == covar)] <-
      1 - printcp(cartmod)[dim(printcp(cartmod))[1], 3]
  }

  if ((max(r2) - currentr2) / currentr2 > 0.01) {
    covar_include <- c(covar_include, covar_cols[which.max(r2)])
    covar_cols <- covar_cols[covar_cols != covar_cols[which.max(r2)]]
    currentr2 <- max(r2)
    r2 <- c()
  } else {

```

```

    break
  }
}

```

Print out the number of covariates selected and the resulting CART model output.

```
print("Number of covariates:")
```

```
## [1] "Number of covariates:"
```

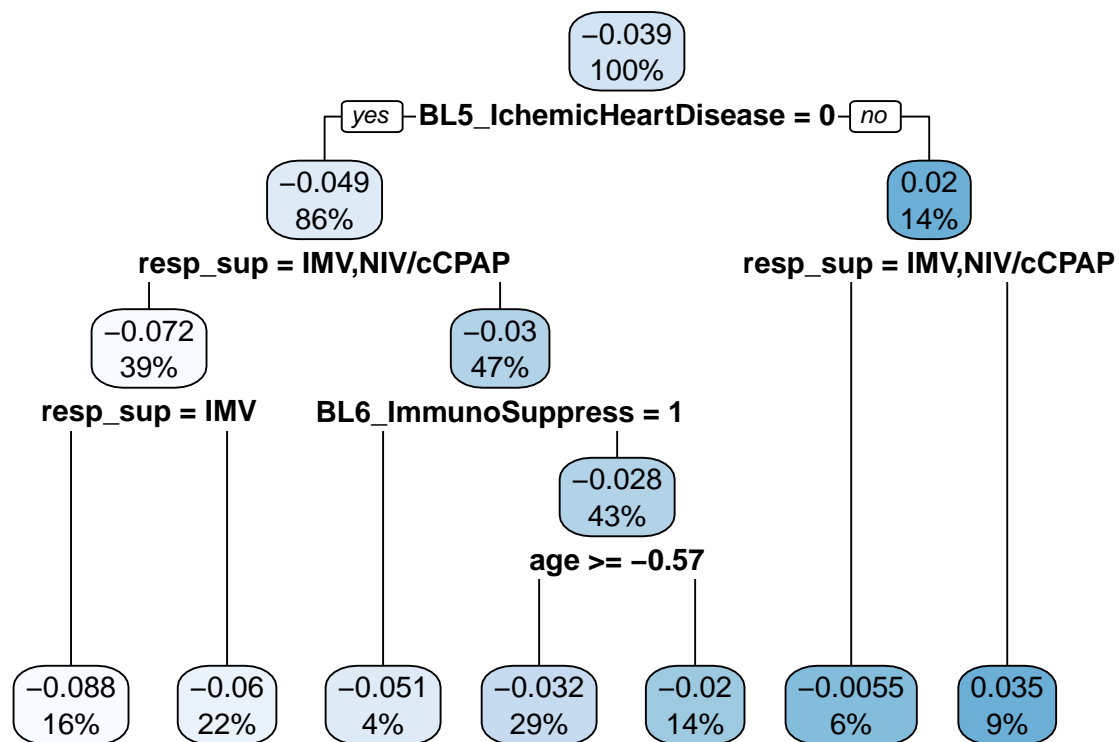
```
print(numcovars - 1)
```

```
## [1] 3
```

```

cartmod <- rpart(cate ~ ., data = dat[, c(covar_include, 14)],
                 method = "anova")
rpart.plot(cartmod)

```

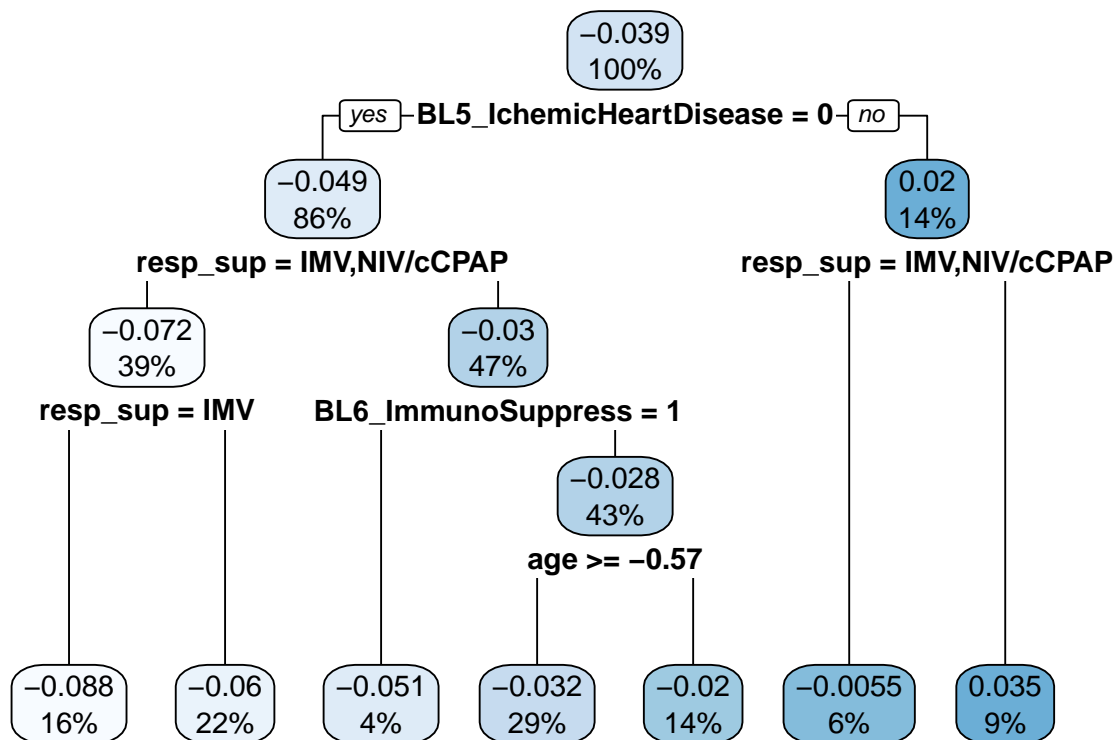


Finally we explore further pruning:

```
printcp(cartmod)
```

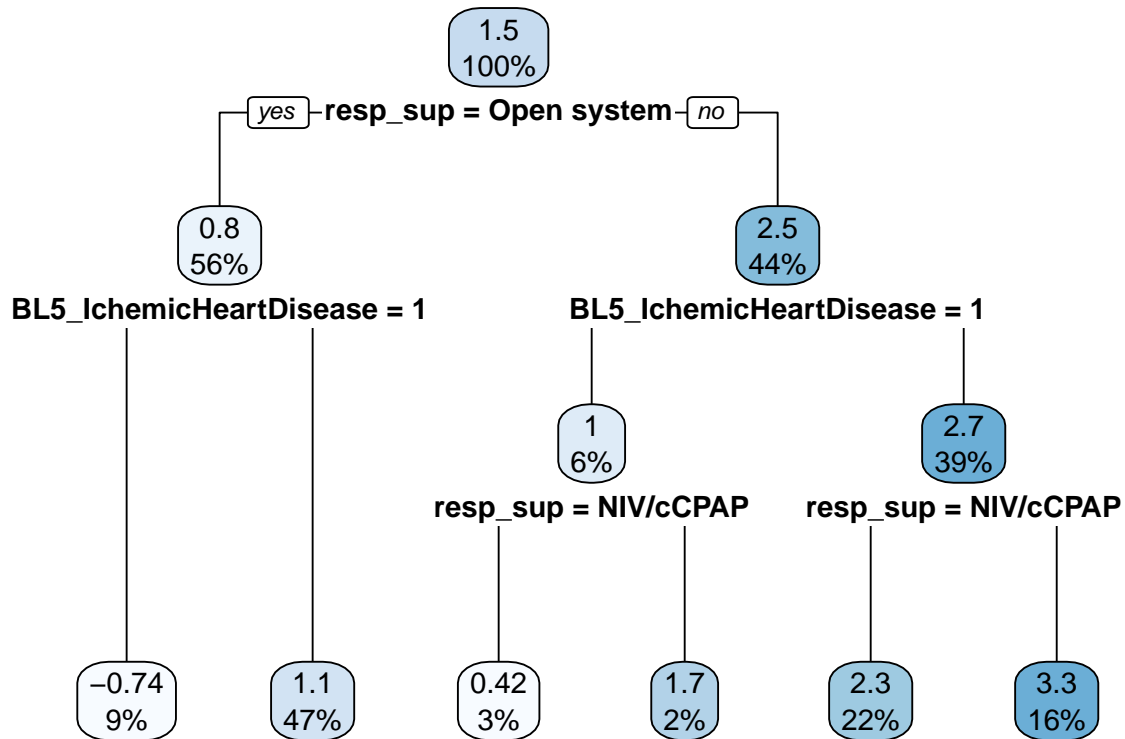
```
##
## Regression tree:
## rpart(formula = cate ~ ., data = dat[, c(covar_include, 14)],
##       method = "anova")
##
## Variables actually used in tree construction:
## [1] age          BL5_IchemicHeartDisease BL6_ImmunoSuppress
## [4] resp_sup
##
## Root node error: 1.2489/982 = 0.0012718
##
## n= 982
##
##      CP nsplit rel error  xerror   xstd
## 1 0.459150    0  1.00000 1.00331 0.0466035
## 2 0.288499    1  0.54085 0.54239 0.0189474
## 3 0.055730    2  0.25235 0.25361 0.0111441
## 4 0.045092    3  0.19662 0.19892 0.0101958
## 5 0.014793    4  0.15153 0.15288 0.0077625
## 6 0.010603    5  0.13674 0.13839 0.0073919
## 7 0.010000    6  0.12613 0.13484 0.0076368
```

```
prunedtree <-
  prune(cartmod,
        cp = cartmod$sctable[which.min(cartmod$sctable[, "xerror"]), "CP"])
rpart.plot(prunedtree)
```



Next the same fit-the-fit approach is used to summarize the results for the continuous outcome, starting with a CART model using all covariates and default hyperparameter.

```
# CART model for 90 day mortality with default CART hyperparameter and
# all covariates considered
cartmod <- rpart(cate_c ~ ., data = dat[, c(4:13, 15)], method = "anova")
rpart.plot(cartmod)
```



Now we prune the tree for interpretability using the stepwise approach of Hu et al.

```
# Pruned tree using stepwise approach of Hu et al.
r2 <- c()
covar_include <- c()
covar_cols <- 4:13
currentr2 <- 0
maxcovars <- 4

for (numcovars in 1:maxcovars) {

  for (covar in covar_cols) {
    cartmod <- rpart(cate_c ~ ., data = dat[, c(covar, covar_include, 15)],
                      method = "anova")
    r2[which(covar_cols == covar)] <-
      1 - printcp(cartmod)[dim(printcp(cartmod))[1], 3]
  }
}
```

```

if ((max(r2) - currentr2) / currentr2 > 0.01) {
  covar_include <- c(covar_include, covar_cols[which.max(r2)])
  covar_cols <- covar_cols[covar_cols != covar_cols[which.max(r2)]]
  currentr2 <- max(r2)
  r2 <- c()
} else {
  break
}
}

```

Print out the number of covariates selected and the resulting CART model output.

```
print("Number of covariates:")
```

```
## [1] "Number of covariates:"
```

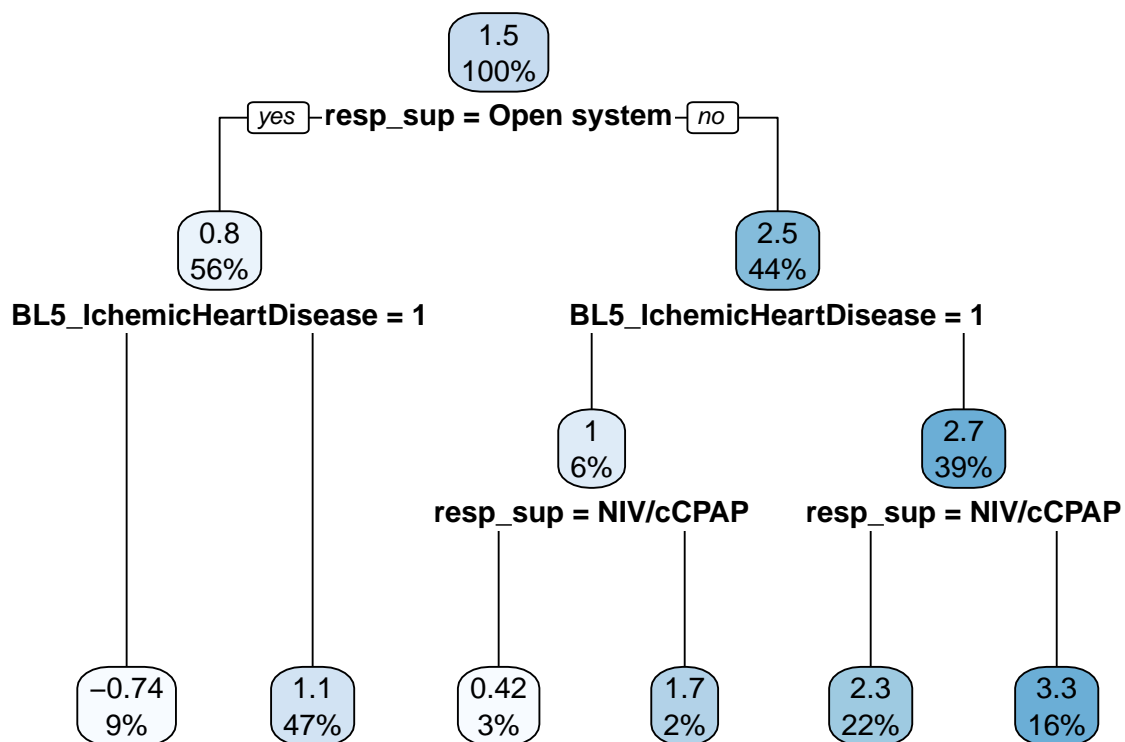
```
print(numcovars - 1)
```

```
## [1] 2
```

```

cartmod <- rpart(cate_c ~ ., data = dat[, c(covar_include, 15)],
  method = "anova")
rpart.plot(cartmod)

```





Finally we explore further pruning:

```
printcp(cartmod)
```

```
##
## Regression tree:
## rpart(formula = cate_c ~ ., data = dat[, c(covar_include, 15)],
##       method = "anova")
##
## Variables actually used in tree construction:
## [1] BL5_IchemicHeartDisease resp_sup
##
## Root node error: 1257.8/982 = 1.2809
##
## n= 982
##
##      CP nsplit rel error   xerror   xstd
## 1 0.554046      0 1.000000 1.001602 0.0423371
## 2 0.193368      1 0.445954 0.447664 0.0263096
## 3 0.111685      2 0.252586 0.253709 0.0211178
## 4 0.069989      3 0.140901 0.141787 0.0063733
## 5 0.018810      4 0.070912 0.071802 0.0049041
## 6 0.010000      5 0.052103 0.052839 0.0027216
```

```
prunedtree <-
  prune(cartmod,
        cp = cartmod$cptable[which.min(cartmod$cptable[, "xerror"]), "CP"])
rpart.plot(prunedtree)
```

