

07 - Trees: Demo

R Code for analyzing CART regression trees

SYS 6018 | Fall 2020

07-trees_demo.pdf

Contents

1	Trees Intro	2
1.1	Required R Packages	2
1.2	Baseball Salary Data	2
2	Regression Tree	2
2.1	Build Tree	2
2.2	Evaluate Tree	4
2.3	Regression Tree example with 2 dimensions only	6
3	Details of Splitting (for Regression Trees)	10
3.1	First Split	10
3.2	Second Split	15

1 Trees Intro

1.1 Required R Packages

We will be using the R packages of:

- `rpart` for classification and regression trees (CART)
- `rpart.plot` for `prp()` which allows more plotting control for trees
- `randomForest` for `randomForest()` function
- `ISLR` for Hitters baseball data
- `tidyverse` for data manipulation and visualization

```
1 library(ISLR)
2 library(rpart)
3 library(rpart.plot)
4 library(randomForest)
5 library(tidyverse)
```

1.2 Baseball Salary Data

The goal is to build models to predict the (log) salary of baseball players

```
1 head(bball)
2 #>   AtBat Hits HmRun Runs RBI Walks Years CAtBat CHits CHmRun CRuns CRBI CWalks
3 #> 1   315   81    7   24  38   39   14   3449   835    69   321   414   375
4 #> 2   479  130   18   66  72   76    3   1624   457    63   224   266   263
5 #> 4   321   87   10   39  42   30    2    396   101    12    48    46    33
6 #> 5   594  169    4   74  51   35   11   4408  1133    19   501   336   194
7 #> 6   185   37    1   23   8   21    2    214    42     1    30    9    24
8 #> 7   298   73    0   24  24    7    3    509   108     0    41   37    12
9 #>   League Division PutOuts Assists Errors      Y NewLeague
10 #> 1      N         W     632     43    10 6.163          N
11 #> 2      A         W     880     82    14 6.174          A
12 #> 4      N         E     805     40     4 4.516          N
13 #> 5      A         W     282    421    25 6.620          A
14 #> 6      N         E      76    127     7 4.248          A
15 #> 7      A         W     121    283     9 4.605          A
```

2 Regression Tree

2.1 Build Tree

```
1 #####
2 #-- Regression Trees in R
3 # trees are in many packages: rpart, tree, party, ...
4 # there are also many packages to display tree results
5 #
6 # Formulas: you don't need to specify interactions as the tree does this
7 # naturally.
8 #####
9 #-- Build Tree
10 library(rpart)
11 tree = rpart(Y~., data=bball)
12 summary(tree, cp=1)
13 #> Call:
14 #> rpart(formula = Y ~ ., data = bball)
```

```

15 #> n= 200
16 #>
17 #>      CP nsplit rel error xerror  xstd
18 #> 1 0.61187      0    1.0000 1.0034 0.07223
19 #> 2 0.08077      1    0.3881 0.4361 0.04334
20 #> 3 0.05616      2    0.3074 0.3754 0.04487
21 #> 4 0.05037      3    0.2512 0.3534 0.04483
22 #> 5 0.01840      4    0.2008 0.2500 0.02677
23 #> 6 0.01381      5    0.1824 0.2553 0.02637
24 #> 7 0.01000      6    0.1686 0.2433 0.02594
25 #>
26 #> Variable importance
27 #> CAtBat CHits CRuns CRBI CWalks CHmRun Hits AtBat Runs Walks RBI
28 #> 17 17 16 15 14 11 2 2 2 2 1
29 #> HmRun
30 #> 1
31 #>
32 #> Node number 1: 200 observations
33 #> mean=5.898, MSE=0.8053

```

```

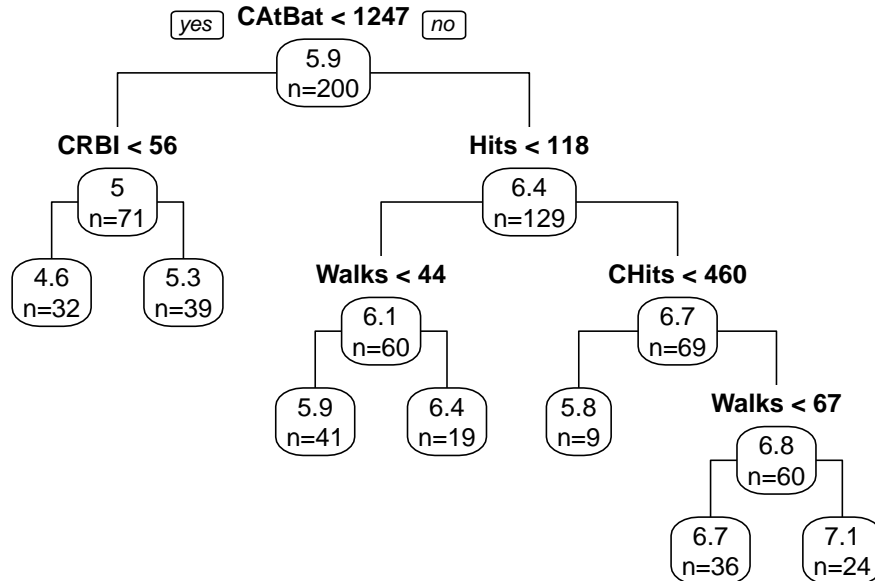
1 length(unique(tree$where)) # number of leaf nodes
2 #> [1] 7

```

```

1
2 #-- Plot Tree
3 library(rpart.plot) # for prp() which allows more plotting control
4 prp(tree, type=1, extra=1, branch=1)

```



```

1
2 # rpart() functions can also plot (just not as good):
3 # plot(tree, uniform=TRUE)
4 # text(tree, use.n=TRUE, xpd=TRUE)

```

2.2 Evaluate Tree

```

1  #- mean squared error function
2  mse <- function(yhat, y){
3    yhat = as.matrix(yhat)
4    apply(yhat, 2, function(f) mean((f-y)^2))
5  }
6
7
8  mse(predict(tree), bball$Y)          # training error
9  #> [1] 0.1358

1  mse(predict(tree, X.test), Y.test)   # testing error
2  #> [1] 0.4931

```

Build a more complex tree

```

1  #-- More complex tree
2  # see ?rpart.control() for details
3  # xval: number of cross-validations
4  # minsplit: min obs to still allow a split
5  # cp: complexity parameter
6
7  tree2 = rpart(Y~., data=bball, xval=0, minsplit=5, cp=0.005)
8  summary(tree2, cp=1)
9  #> Call:
10 #> rpart(formula = Y ~ ., data = bball, xval = 0, minsplit = 5,
11 #>      cp = 0.005)
12 #>      n= 200
13 #>
14 #>      CP nsplit rel error
15 #> 1  0.611866      0  1.00000
16 #> 2  0.080767      1  0.38813
17 #> 3  0.056162      2  0.30737
18 #> 4  0.050368      3  0.25121
19 #> 5  0.018404      4  0.20084
20 #> 6  0.013809      5  0.18243
21 #> 7  0.008264      6  0.16862
22 #> 8  0.007883      7  0.16036
23 #> 9  0.007740      8  0.15248
24 #> 10 0.007267      9  0.14474
25 #> 11 0.007129     10  0.13747
26 #> 12 0.006491     11  0.13034
27 #> 13 0.005916     12  0.12385
28 #> 14 0.005816     14  0.11202
29 #> 15 0.005332     15  0.10620
30 #> 16 0.005078     16  0.10087
31 #> 17 0.005000     18  0.09071
32 #>
33 #> Variable importance
34 #> CAtBat CHits CRuns CRBI CWalks CHmRun Hits AtBat Runs RBI Walks
35 #>      17      17      16      15      13      10       2       2       2       2       2
36 #> HmRun  Years
37 #>       1       1
38 #>
39 #> Node number 1: 200 observations
40 #>      mean=5.898, MSE=0.8053

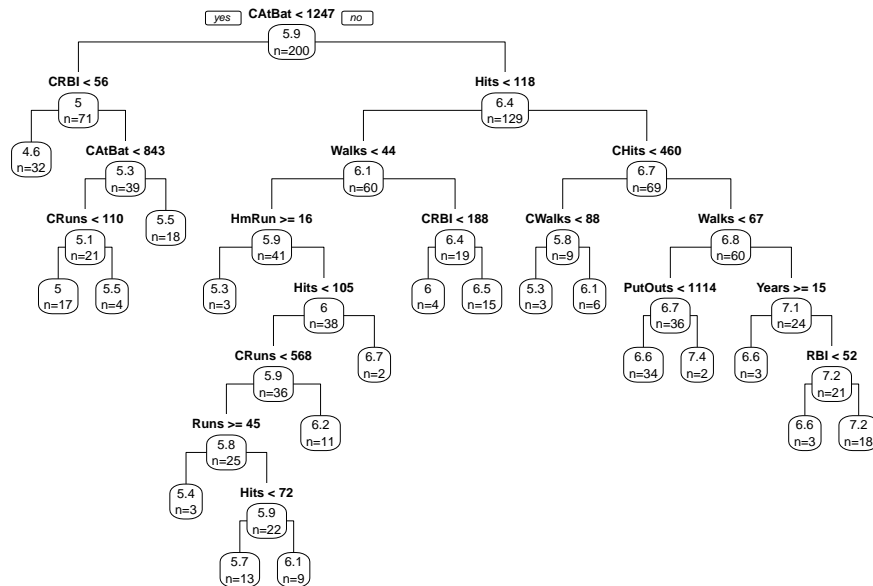
1  length(unique(tree2$where))
2  #> [1] 19

```

```

1
2 prp(tree2, type=1, extra=1, branch=1)

```



```

1
2 mse(predict(tree2), bball$Y)           # training error
3 #> [1] 0.07305

1 mse(predict(tree2, X.test), Y.test)    # testing error
2 #> [1] 0.552

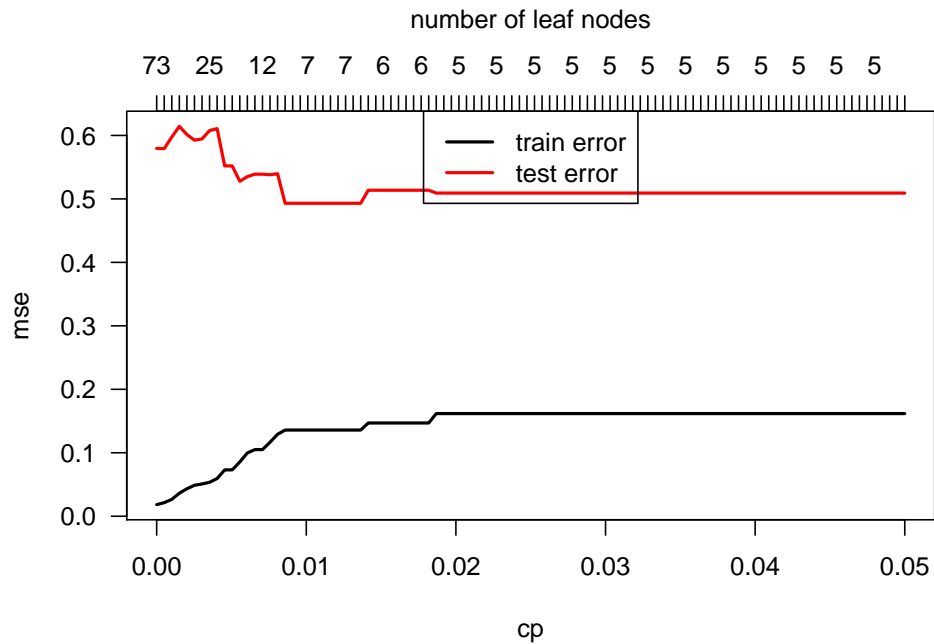
```

Now, fit a set of tree for range of cp values.

```

1 cp = seq(.05,0,length=100) # cp is like a penalty on the tree size
2 for(i in 1:length(cp)) {
3   if(i == 1){train.error = test.error = nleafs = numeric(length(cp))}
4   tree.fit = rpart(Y~.,data=bball, xval=0, minsplit=5, cp=cp[i])
5   train.error[i] = mse(predict(tree.fit), bball$Y)           # training error
6   test.error[i] = mse(predict(tree.fit, X.test), Y.test)    # testing error
7   nleafs[i] = length(unique(tree.fit$where))
8 }
9
10 plot(range(cp), range(train.error, test.error), typ='n', xlab="cp", ylab="mse", las=1)
11 lines(cp, train.error, col="black", lwd=2)
12 lines(cp, test.error, col="red", lwd=2)
13 legend("top", c('train error', 'test error'), col=c("black", "red"), lwd=2)
14 axis(3, at=cp, labels=nleafs)
15 mtext("number of leaf nodes", 3, line=2.5)

```



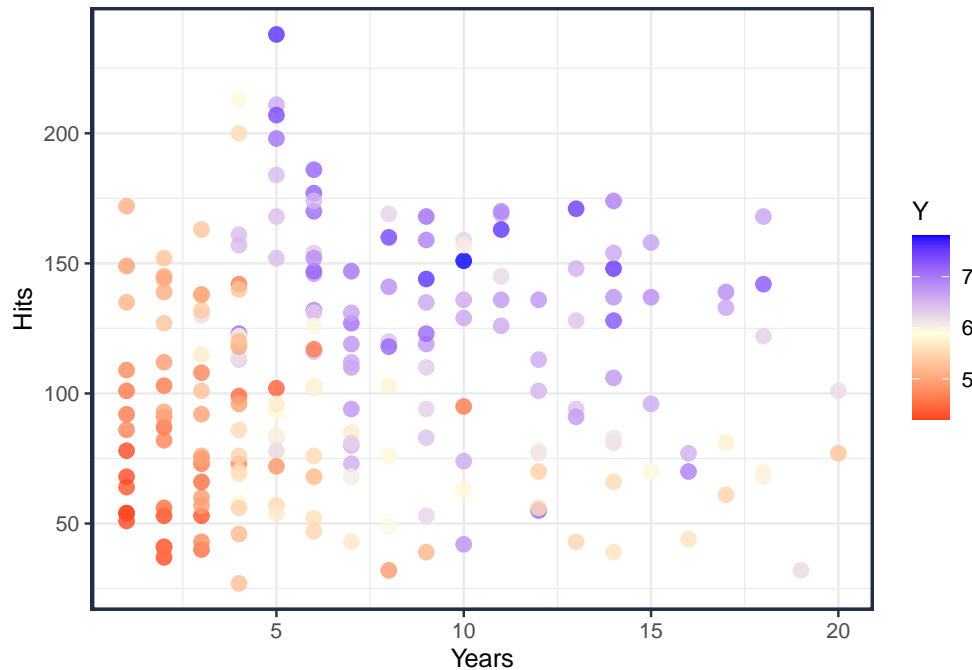
2.3 Regression Tree example with 2 dimensions only

Consider the two variables `Years` and `Hits` and their relationship to `Y`.

```

1 #####
2 #-- Regression Tree Examples for 2D
3 #####
4 library(ggplot2)
5
6 #-- 2D plot (using only Years and Hits)
7 p2D = ggplot(bball) + #scale_size_area(max_size=5) +
8     scale_color_gradient2(midpoint=mean(bball$Y), mid="lightyellow", low="red", high="blue")
9 p2D +
10     geom_point(aes(x=Years, y=Hits, color=Y), alpha=.8, size=3)

```



Let's fit a tree with the two predictors

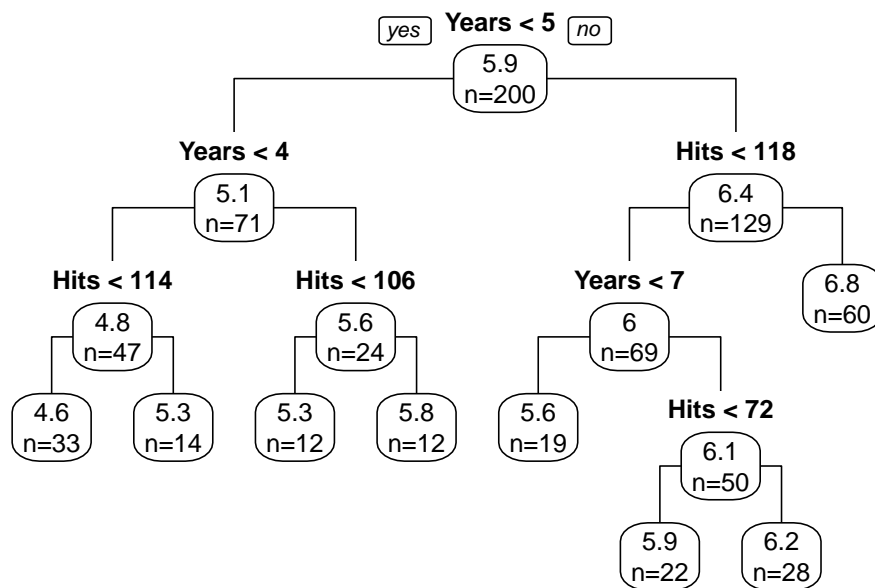
```

1  #-- Fit tree to only Years and Hits
2  tree3 = rpart(Y~Years+Hits, data=bball)
3  summary(tree3,cp=1)
4  #> Call:
5  #> rpart(formula = Y ~ Years + Hits, data = bball)
6  #>   n= 200
7  #>
8  #>      CP nsplit rel error xerror  xstd
9  #> 1 0.47952    0   1.0000 1.0110 0.07293
10 #> 2 0.15162    1   0.5205 0.5722 0.06131
11 #> 3 0.05761    2   0.3689 0.4242 0.05369
12 #> 4 0.02587    3   0.3113 0.3831 0.04738
13 #> 5 0.01892    4   0.2854 0.3593 0.04643
14 #> 6 0.01019    5   0.2665 0.3427 0.04869
15 #> 7 0.01011    6   0.2563 0.3539 0.05074
16 #> 8 0.01000    7   0.2462 0.3539 0.05074
17 #>
18 #> Variable importance
19 #> Years Hits
20 #>   73   27
21 #>
22 #> Node number 1: 200 observations
23 #>   mean=5.898, MSE=0.8053

1  length(unique(tree3$where))           # number of leaf nodes
2  #> [1] 8

1  prp(tree3, type=1, extra=1, branch=1)

```



```

1 mse(predict(tree3), bball$Y)           # training error
2 #> [1] 0.1982

```

```

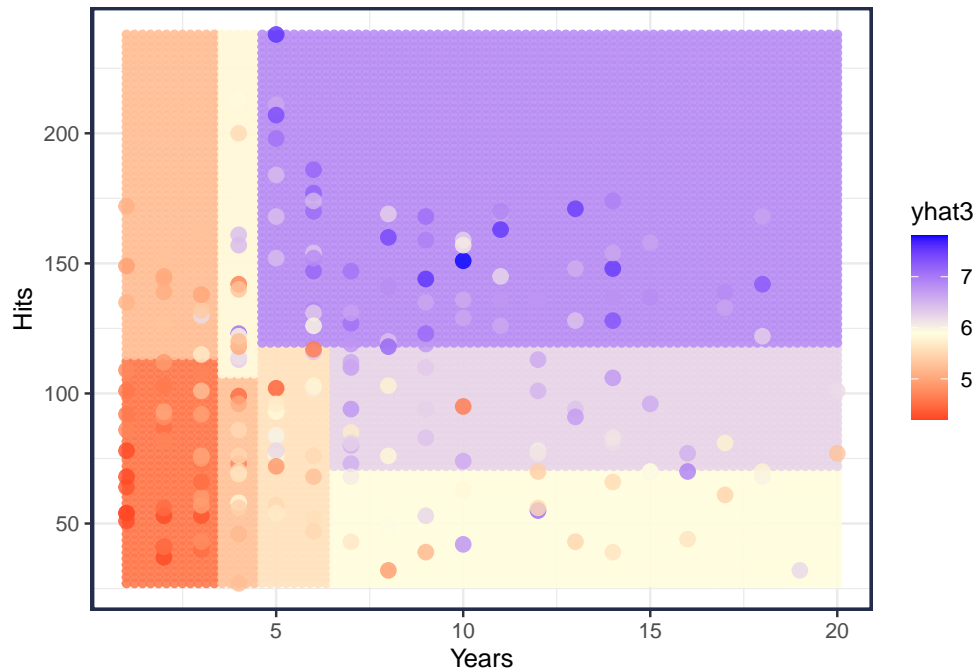
1 mse(predict(tree3,X.test),Y.test)      # testing error
2 #> [1] 0.5179

```

```

1
2
3 #-- Plot Results
4 grid = expand.grid(Years = seq(min(bball$Years),max(bball$Years),length=90),
5                   Hits = seq(min(bball$Hits),max(bball$Hits),length=90))
6 grid$yhat3 = predict(tree3,newdata = grid)
7
8 p2D + geom_point(data=grid,aes(x=Years, y=Hits, color=yhat3),alpha=.9) +
9   geom_point(aes(x=Years, y=Hits, color=Y), alpha=.8, size=3)

```

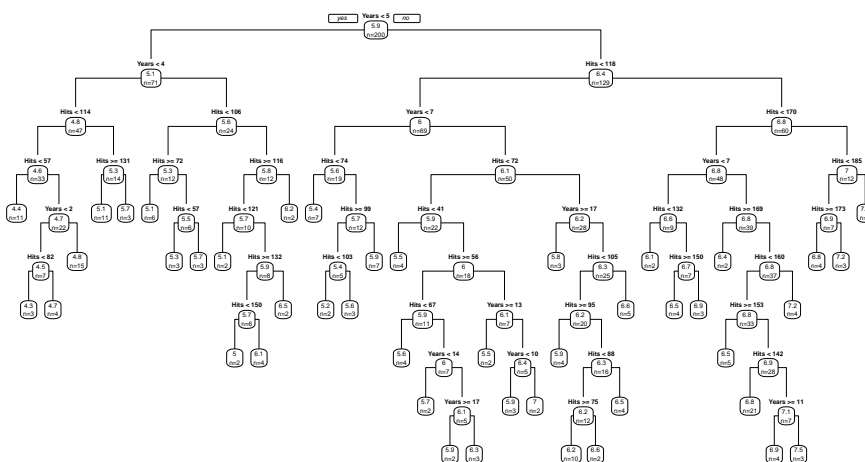



This shows the leaf regions (in 2D).

And we can also use more complex trees:

```
1  #-- Fit more complex tree to only Years and Hits
2  tree4 = rpart(Y~Years+Hits,data=bball,xval=0,minsplit=5,cp=0.001)
3  length(unique(tree4$where))           # number of leaf nodes
4  #> [1] 44
```

```
1  prp(tree4, type=1, extra=1, branch=1)
```



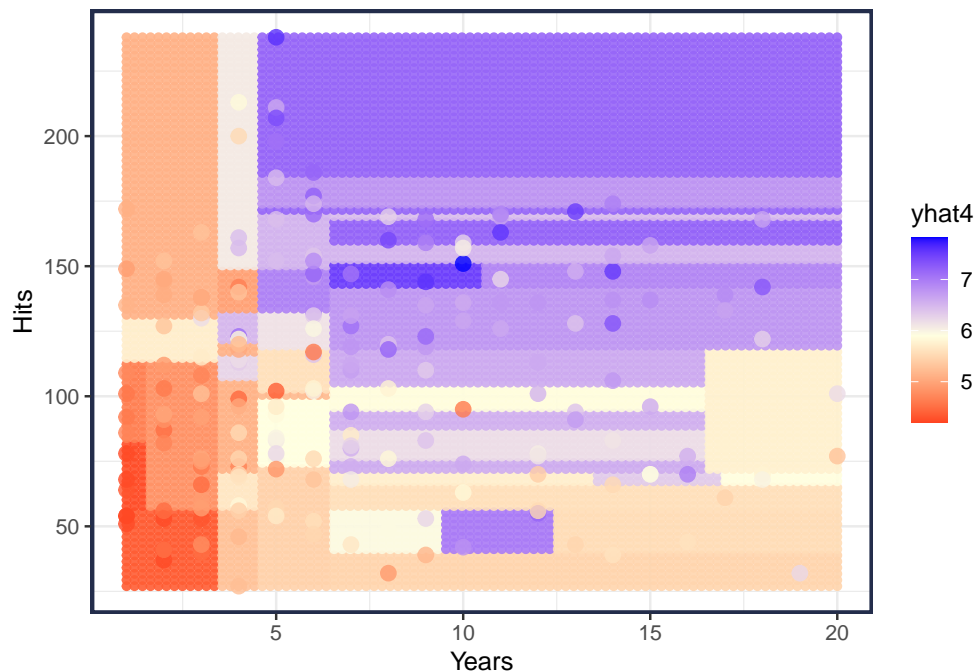
```

1 mse(predict(tree4), bball$Y)           # training error
2 #> [1] 0.09876

1 mse(predict(tree4,X.test), Y.test)     # testing error
2 #> [1] 0.6959

1
2
3 #-- Plot Results
4 grid$yhat4 = predict(tree4,newdata = grid)
5
6 p2D + geom_point(data=grid,aes(x=Years,y=Hits,color=yhat4),alpha=.9) +
7   geom_point(aes(x=Years,y=Hits,color=Y),alpha=.8, size=3)

```



3 Details of Splitting (for Regression Trees)

Consider only two dimensions, hits and years on which to make first split.

See `trees.R` for the `split_info()` and `split_metrics()` functions

3.1 First Split

3.1.1 Split on Years

```

1 ## Split by Years
2 years = split_info(x=bball$Years, y=bball$Y)
3 head(years)
4 #> # A tibble: 6 x 9
5 #>   split.pt  n.L  n.R est.L est.R SSE.L SSE.R  SSE  gain
6 #>   <dbl> <int> <int> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
7 #> 1     1.5    13   187  4.57  5.99  1.37 135.  137.   24.5
8 #> 2     2.5    29   171  4.68  6.11  3.63 107.  110.   50.7

```

```

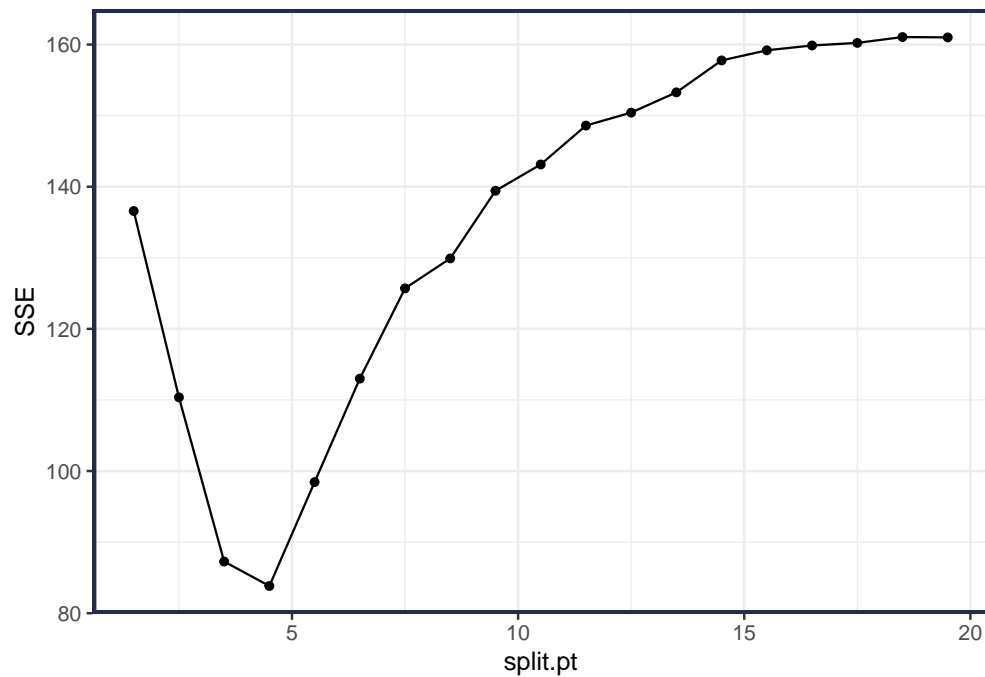
9  #> 3      3.5    47   153  4.80  6.23  8.29  79.0  87.3  73.8
10 #> 4      4.5    71   129  5.06  6.36 26.5  57.3  83.8  77.2
11 #> 5      5.5    88   112  5.27  6.39 52.7  45.7  98.4  62.6
12 #> 6      6.5   108    92  5.45  6.43 78.3  34.7 113.   48.1

```

```

1
2 ggplot(years,aes(x=split.pt,y=SSE)) + geom_line() + geom_point()

```



```

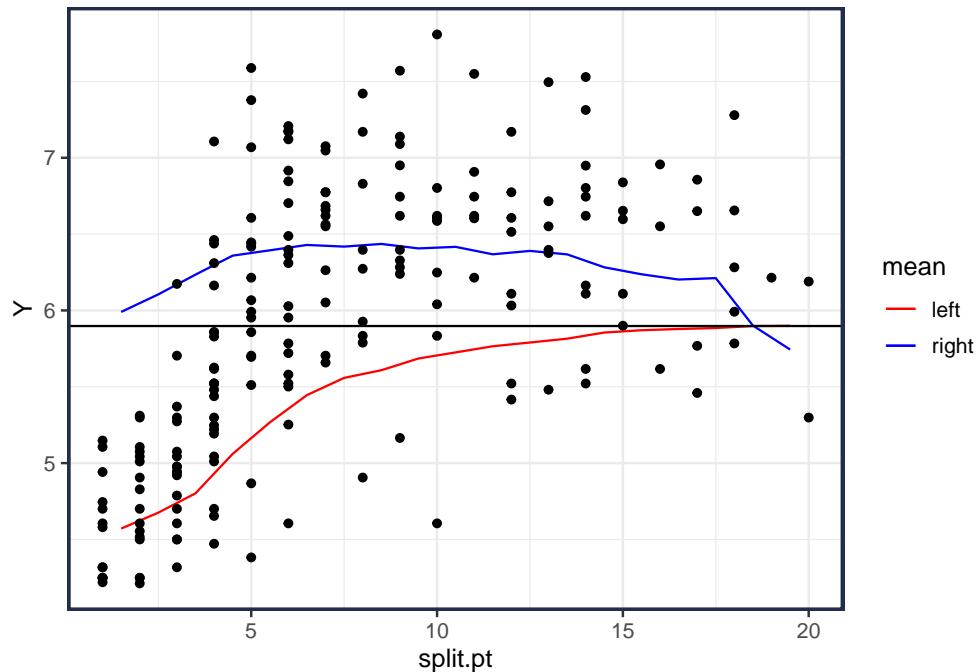
1
2 filter(years, min_rank(SSE) == 1) # optimal split point for Years
3 #> # A tibble: 1 x 9
4 #>   split.pt  n.L  n.R est.L est.R SSE.L SSE.R  SSE  gain
5 #>   <dbl> <int> <int> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
6 #> 1     4.5    71   129  5.06  6.36  26.5  57.3  83.8  77.2

```

```

1
2 ggplot(years,aes(x=split.pt)) +
3   geom_line(aes(y=est.L,color="left")) + # mean left of split pt
4   geom_line(aes(y=est.R,color="right")) + # mean right of split pt
5   geom_hline(yintercept=mean(bball$Y)) + # overall mean
6   scale_color_manual("mean",values=c('left'='red','right'='blue')) +
7   geom_point(data=bball,aes(x=Years,y=Y)) + # add points
8   labs(y="Y")

```



3.1.2 Split on Hits

```

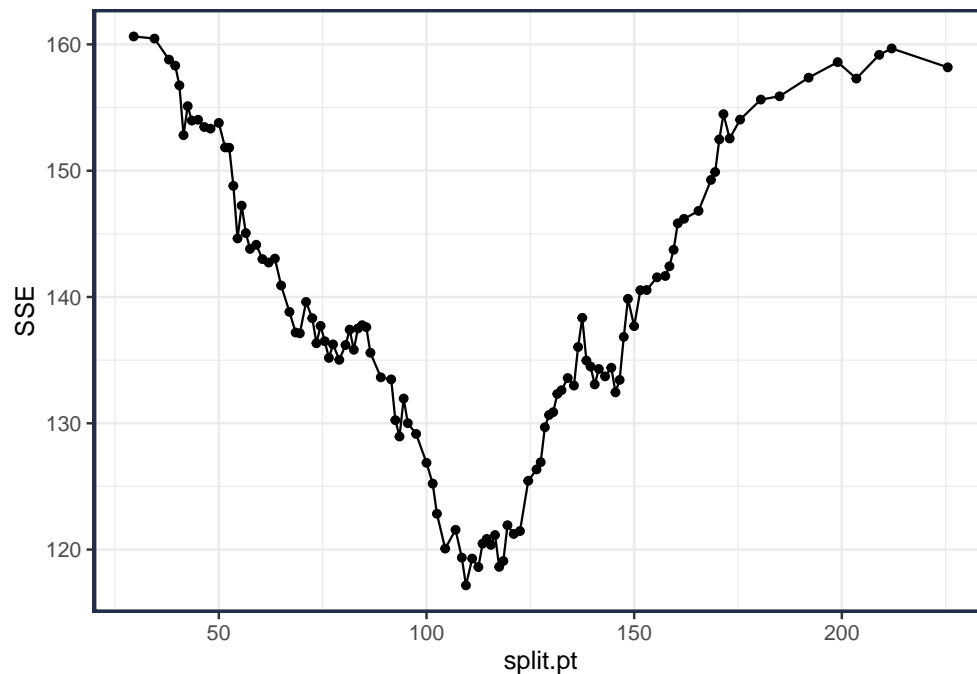
1  ## Split by Hits
2  hits = split_info(x=bball$Hits,y=bball$Y)
3  head(hits)
4  #> # A tibble: 6 x 9
5  #>   split.pt  n.L  n.R est.L est.R SSE.L SSE.R  SSE  gain
6  #>   <dbl> <int> <int> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
7  #> 1    29.5     1   199  5.25  5.90  0      161.  161.  0.426
8  #> 2    34.5     3   197  5.46  5.90  0.922  160.  160.  0.596
9  #> 3     38      4   196  5.15  5.91  2.02   157.  159.  2.26
10 #> 4    39.5     6   194  5.23  5.92  2.19   156.  158.  2.74
11 #> 5    40.5     7   193  5.13  5.93  2.65   154.  157.  4.30
12 #> 6    41.5     9   191  4.96  5.94  3.57   149.  153.  8.25

```

```

1
2  ggplot(hits,aes(x=split.pt,y=SSE)) + geom_line() + geom_point()

```



```

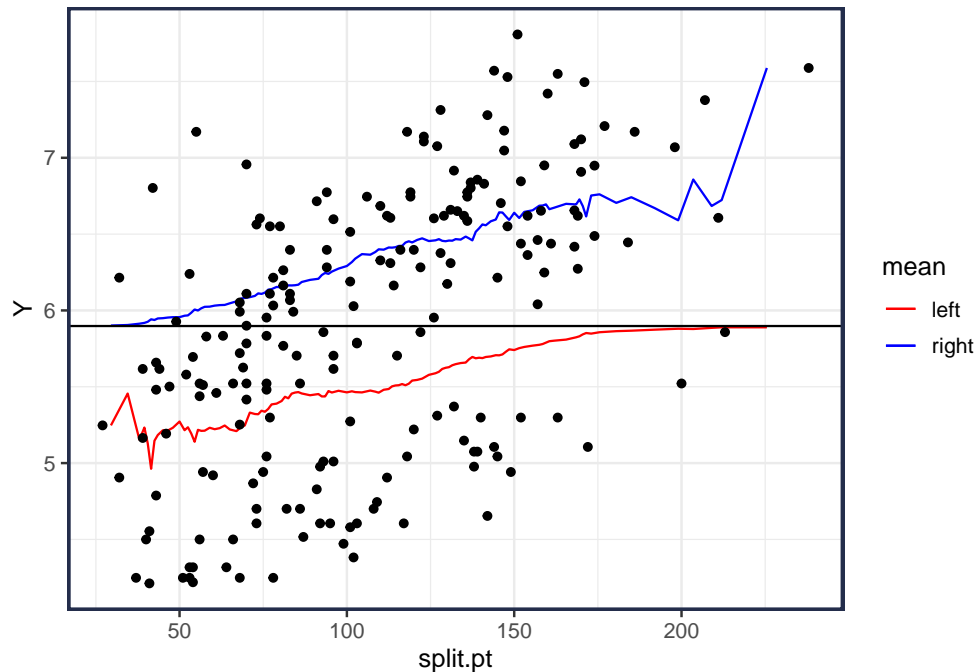
1
2 filter(hits, min_rank(SSE) == 1)      # optimal split point for Hits
3 #> # A tibble: 1 x 9
4 #>   split.pt  n.L  n.R est.L est.R SSE.L SSE.R  SSE  gain
5 #>   <dbl> <int> <int> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
6 #> 1    110.   107    93  5.46  6.40  61.3  55.9  117.  43.9

```

```

1
2 ggplot(hits, aes(x=split.pt)) +
3   geom_line(aes(y=est.L, color="left")) +           # mean left of split pt
4   geom_line(aes(y=est.R, color="right")) +          # mean right of split pt
5   geom_hline(yintercept=mean(bball$Y)) +            # overall mean
6   scale_color_manual("mean", values=c('left'='red', 'right'='blue')) +
7   geom_point(data=bball, aes(x=Hits, y=Y)) +        # add points
8   labs(y = "Y")

```



3.1.3 Find best variable to split on

```

1 ## No splits
2 sum((bball$Y - mean(bball$Y))^2) # SSE if no splits are made
3 #> [1] 161.1

1 # (nrow(bball)-1)*var(bball$Y)
2
3
4
5 ## Results (see function split_metrics at top of file)
6 # splitting on Years gives the best reduction in SSE, so we would split on
7 # Years (at a value of 4.5).
8 sum((bball$Y - mean(bball$Y))^2) # no split
9 #> [1] 161.1

1 filter(years, min_rank(SSE) == 1) # split on years
2 #> # A tibble: 1 x 9
3 #>   split.pt  n.L  n.R est.L est.R SSE.L SSE.R  SSE  gain
4 #>   <dbl> <int> <int> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
5 #> 1     4.5    71   129  5.06  6.36  26.5  57.3  83.8  77.2

1 filter(hits, min_rank(SSE) == 1) # split on hits
2 #> # A tibble: 1 x 9
3 #>   split.pt  n.L  n.R est.L est.R SSE.L SSE.R  SSE  gain
4 #>   <dbl> <int> <int> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
5 #> 1    110.   107    93  5.46  6.40  61.3  55.9  117.  43.9

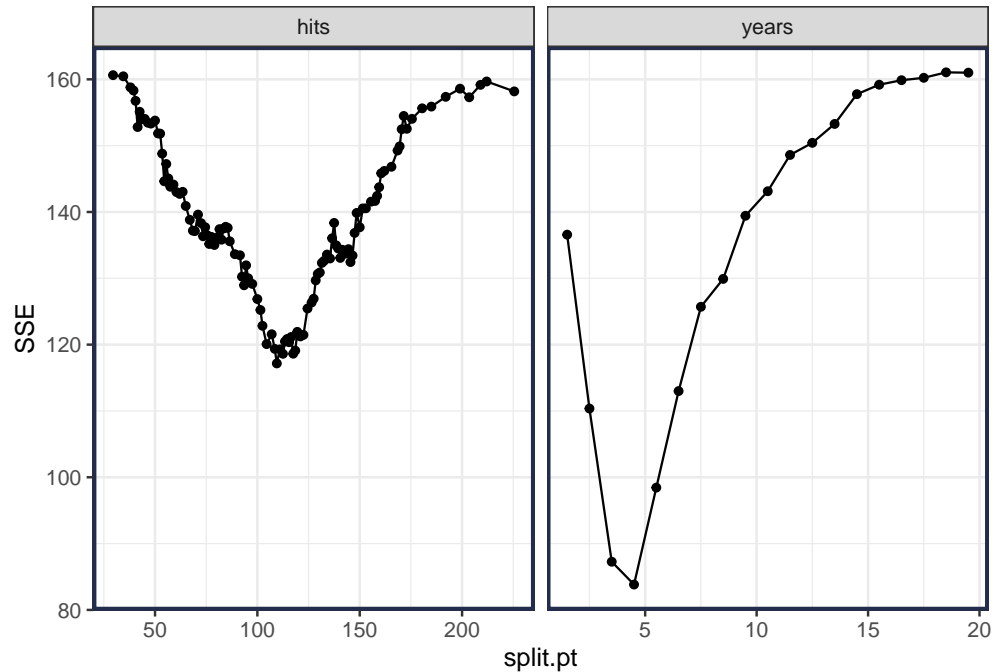
1 split_metrics(bball$Years, bball$Y, 4.5)
2 #> # A tibble: 2 x 3
3 #>   region  SSE    n
4 #>   <chr> <dbl> <int>
5 #> 1 LEFT    26.5    71
6 #> 2 RIGHT   57.3   129

```

```

1
2  ## Comparison of splitting on both variables
3  bind_rows(hits=hits, years=years, .id="split.var") %>%
4    ggplot(aes(x=split.pt, y=SSE)) + geom_line() + geom_point() +
5    facet_wrap(~split.var, scales="free_x")

```



3.2 Second Split

```

1  #-- 2nd Split
2  # now we have to compare 4 possibilities. We can split on Years or Hits, but
3  # use data that has Years < 4.5 or Years > 4.5
4
5  left = (bball$Years<=4.5) # split point from previous step
6  years2.L = split_info(x=bball$Years[left],y=bball$Y[left])
7  years2.R = split_info(x=bball$Years[!left],y=bball$Y[!left])
8  hits2.L = split_info(x=bball$Hits[left],y=bball$Y[left])
9  hits2.R = split_info(x=bball$Hits[!left],y=bball$Y[!left])
10
11 #-- Find best region to split on
12 max(years2.L$gain,na.rm=TRUE)
13 #> [1] 9.278
14
15 max(years2.R$gain,na.rm=TRUE)
16 #> [1] 1.576
17
18 max(hits2.L$gain,na.rm=TRUE)
19 #> [1] 8.726
20
21 max(hits2.R$gain,na.rm=TRUE)
22 #> [1] 24.42
23
24 hits2.R[which.max(hits2.R$gain),]
25 #> # A tibble: 1 x 9
26 #>   split.pt  n.L  n.R est.L est.R SSE.L SSE.R  SSE  gain

```

```

5  #>      <dbl> <int> <int> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
6  #> 1      118.    69    60  5.95  6.83  22.6  10.3  32.9  24.4
1
2  # 2nd split on Hits <= 117.5 in region 2.
3
4
5
6  #-- Summary of Splits
7  # Rule 1: Years < 4.5
8  # Rule 2: Years >= 4.5 & Hits < 117.5
9  # ...
10 prp(tree3, type=1, extra=1, branch=1)

```

