

# Rules Rules Rules! Cubist Regression Models

Boston R User Group

Max Kuhn

Pfizer R&D

# Tree-based Regression Models

Classification and Regression Trees (CART) are a framework for machine learning models.

A CART searches through each predictor to find a value of a single variable that best splits the data into two groups.

- typically, the best split minimizes the RMSE of the outcome in the resulting data subsets.

For the two resulting groups, the process is repeated until a hierarchical structure (a tree) is created.

- in effect, trees partition the  $X$  space into rectangular sections that assign a single value to samples within the rectangle.

To demonstrate, we'll walk through the first two iterations of this process.

# Example Data

The data used to illustrate the models are sale prices of homes in Sacramento CA.

The original data were obtained from the website for the SpatialKey software. From their website:

*The Sacramento real estate transactions file is a list of 985 real estate transactions in the Sacramento area reported over a five-day period, as reported by the Sacramento Bee.*

Google was used to fill in missing/incorrect data.

# Example Data

```
> library(caret)
> data(Sacramento)
> str(Sacramento, vec.len = 1)
```

```
'data.frame': 932 obs. of 9 variables:
 $ city      : Factor w/ 37 levels "ANTELOPE","AUBURN",...: 34 34 ...
 $ zip       : Factor w/ 68 levels "z95603","z95608",...: 64 52 ...
 $ beds      : int  2 3 ...
 $ baths     : num  1 1 ...
 $ sqft      : int  836 1167 ...
 $ type      : Factor w/ 3 levels "Condo","Multi_Family",...: 3 3 ...
 $ price     : int  59222 68212 ...
 $ latitude  : num  38.6 ...
 $ longitude : num  -121 ...
```

# Example Data

A random split was used to create a test set with 20% of the data. The data are:

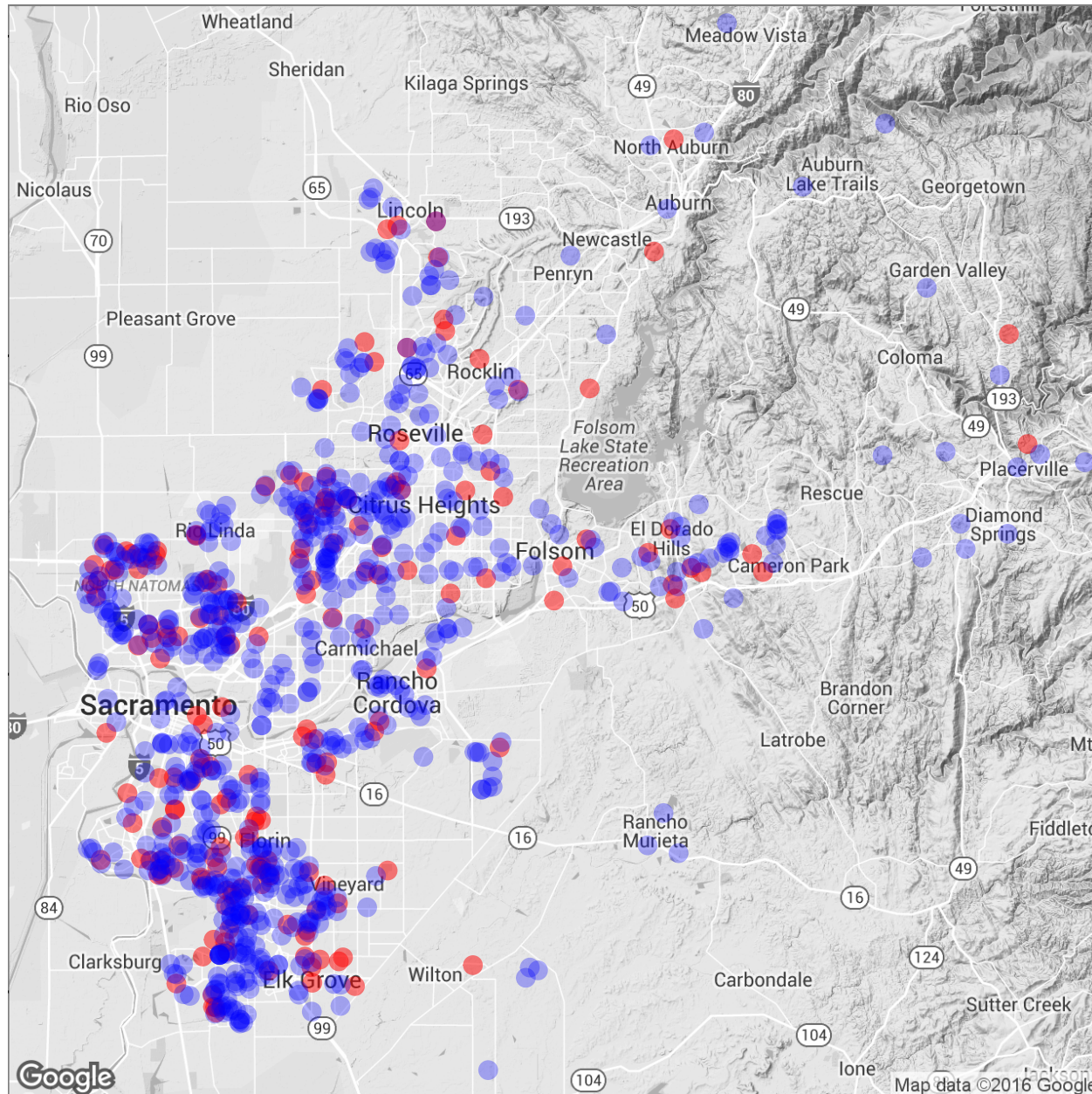
```
> set.seed(955)
> in_train <- createDataPartition(log10(Sacramento$price), p = .8, list = FALSE)
>
> training <- Sacramento[ in_train,]
> testing  <- Sacramento[-in_train,]
> nrow(training)

[1] 747

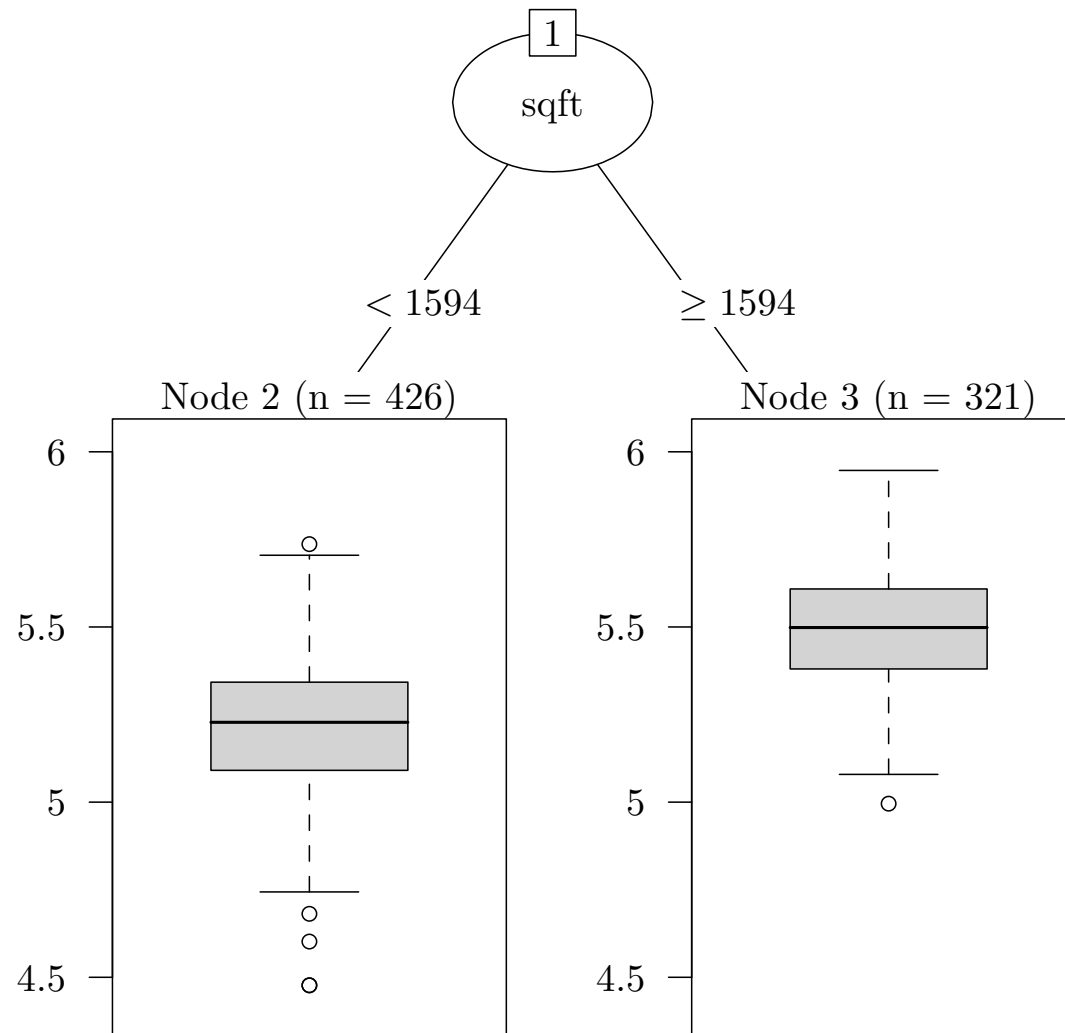
> nrow(testing)

[1] 185
```

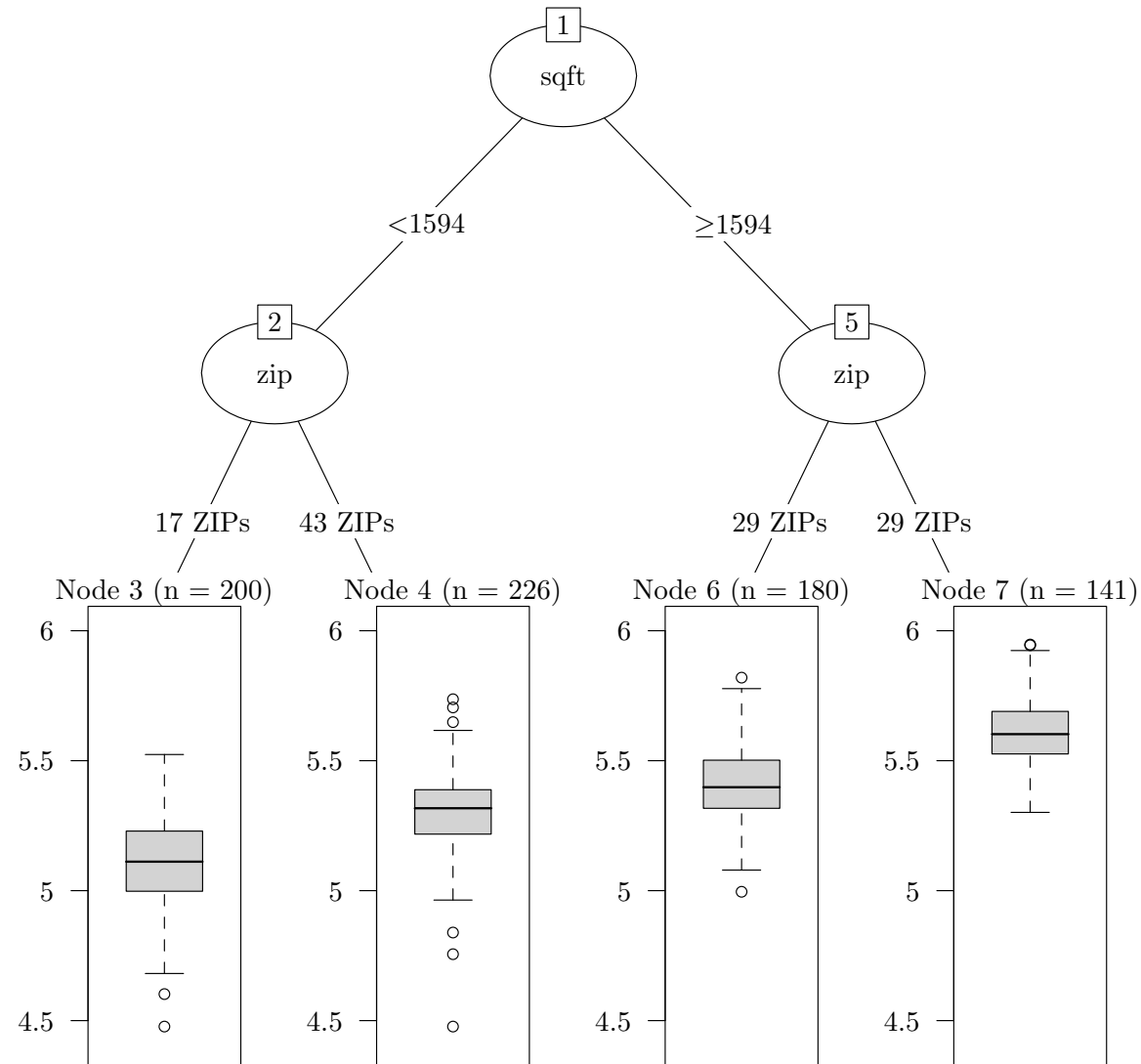
# Training in Blue, Testing in Red



# First Split of a CART Tree

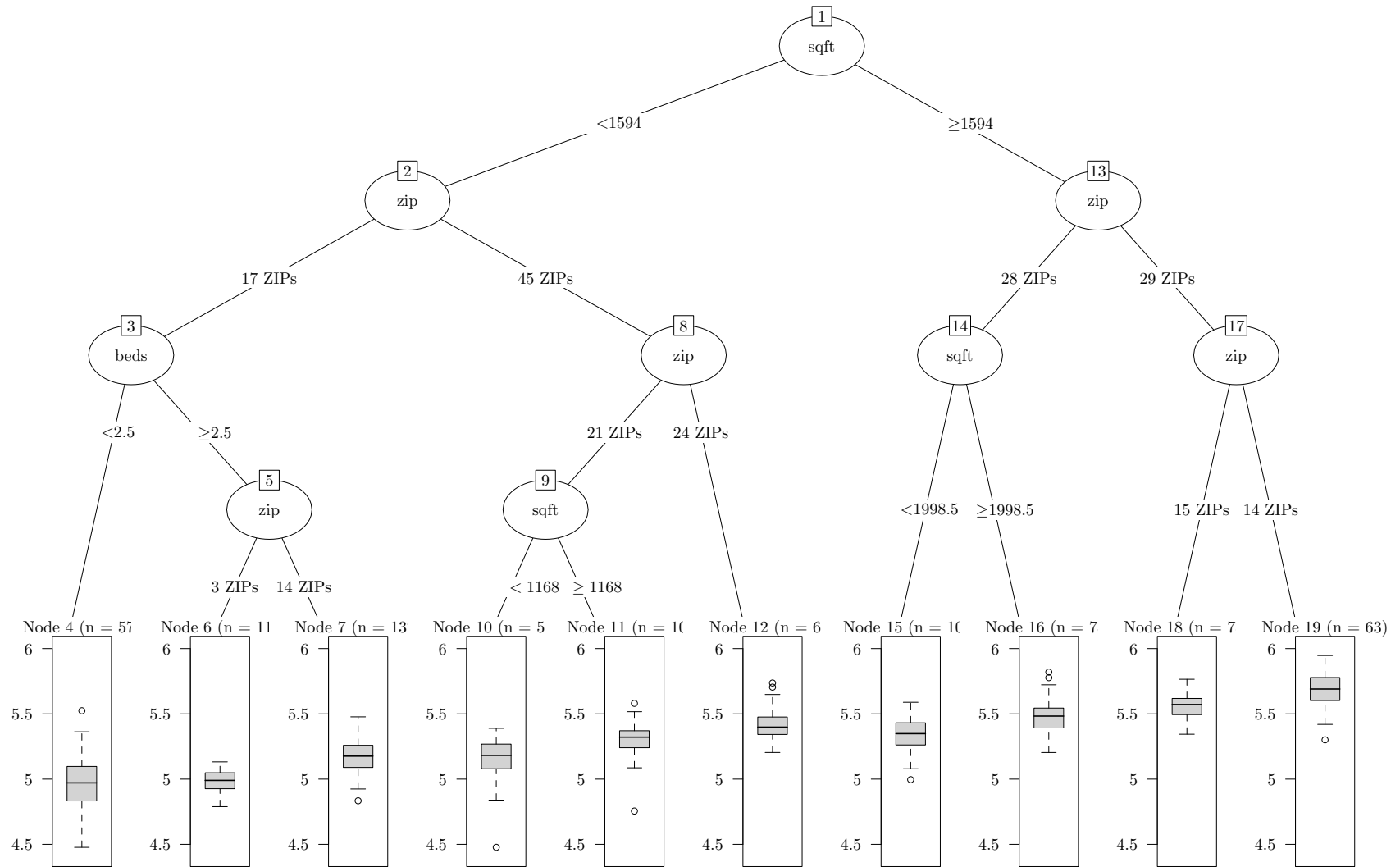


# Second Split





# Full Tree



# The Good and Bad of Trees

Trees can be computed very quickly and have simple interpretations.

Also, they have built-in feature selection; if a predictor was not used in any split, the model is completely independent of that data.

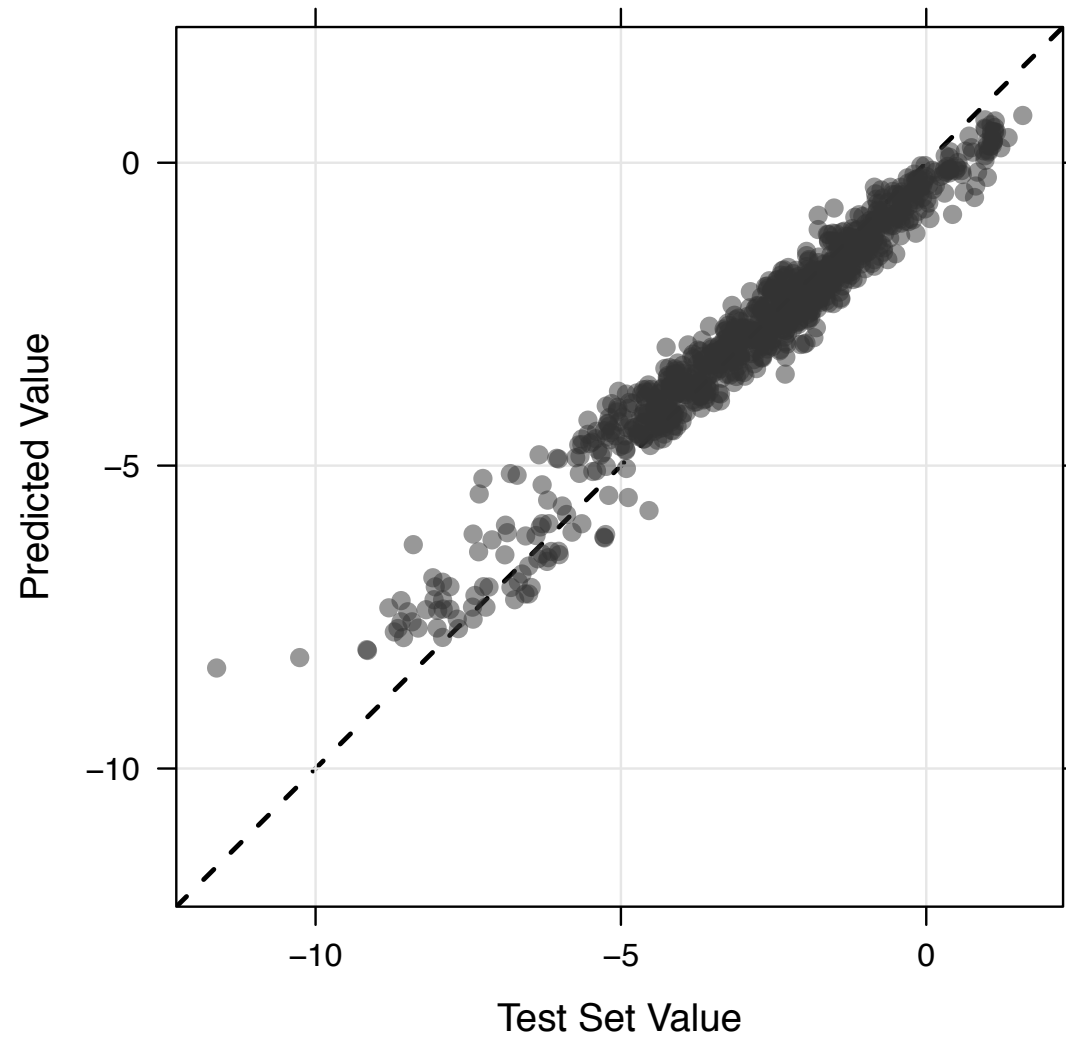
Unfortunately, trees do not usually have optimal performance when compared to other methods.

Also, small changes in the data can drastically affect the structure of a tree.

This last point has been exploited to improve the performance of trees via ensemble methods where many trees are fit and predictions are aggregated across the trees. Examples are bagging, boosting and random forests.

Trees may not fit the data well in the extremes of the outcome range.

# Poor Fits in the Tails



# Model Trees

The *model tree* approach described in Quinlan (1992) called M5, which is similar to regression trees except:

- the splitting criterion is different,
- the terminal nodes predict the outcome using a linear model (as opposed to the simple average), and
- when a sample is predicted, it is often a combination of the predictions from different models along the same path through the tree.

The main implementation of this technique is a “rational reconstruction” of this model called M5’, which is described by Wang and Witten (1997) and is included in the [Weka](#) software package.

# Model Tree Structure

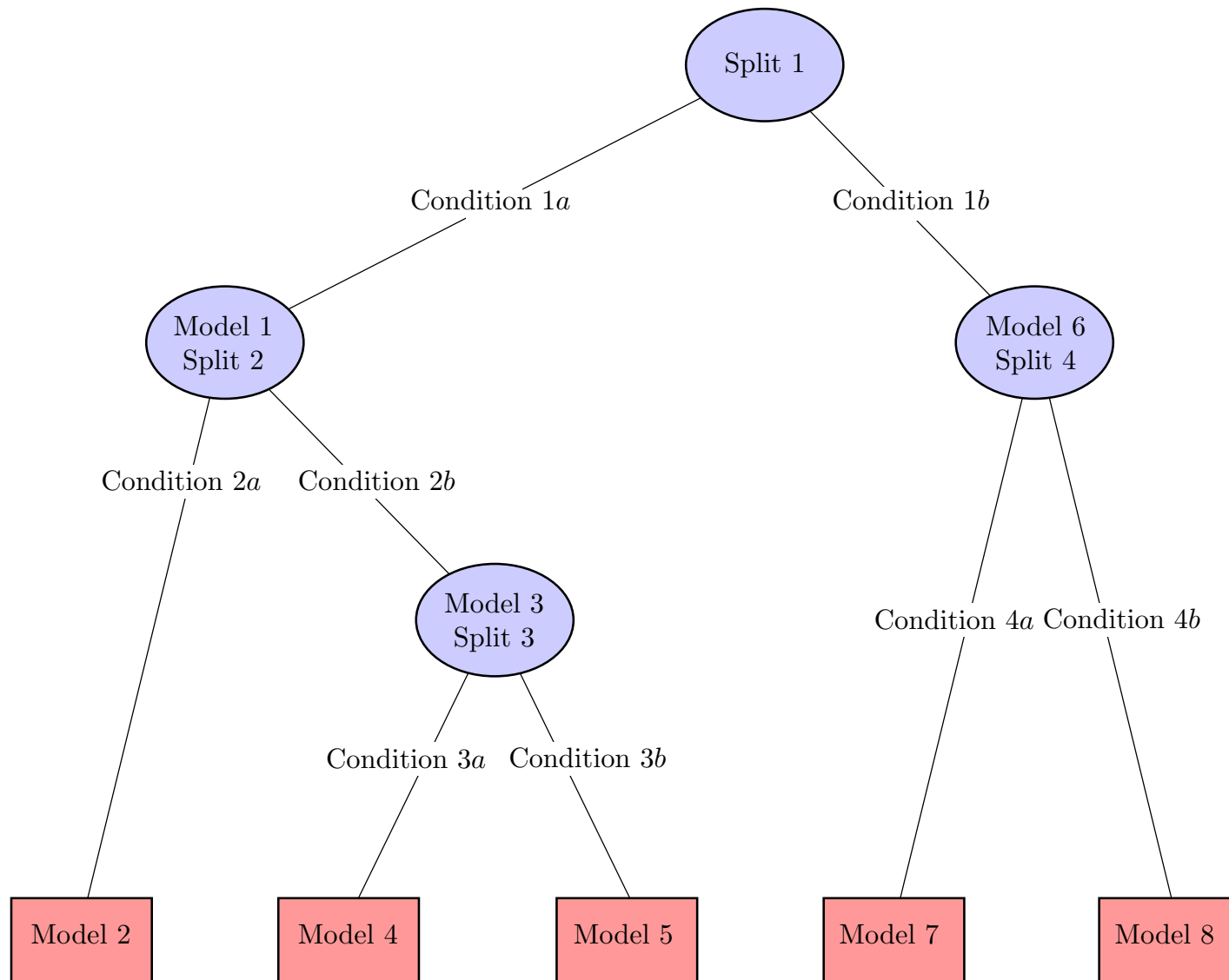
When model trees make a split of the data, they fit a linear model to the current subset using all the predictors involved in the splits along the path.

This process proceeds until there are not enough samples to split and/or fit the model.

A *pruning* stage is later used to simplify the model.

**Note:** Many of the models here are fit with and without encoding categorical predictors as dummy variables. Tree- and rule-based models usually do not require dummy variables.

# Model Tree Structure



# Model Tree Predictions

When a sample is predicted, all of the linear models along the path are combined using:

$$\hat{y}_{par} = \frac{n_{kid} \hat{y}_{kid} + c \hat{y}_{par}}{n_{kid} + c}$$

$\hat{y}_{kid}$  is the prediction from the child node

$n_{kid}$  is the number of training set data points in the child node

$\hat{y}_{par}$  is the prediction from the parent node

$c$  is a constant with a default value of 15.

For the example data, the unpruned model had 81 paths through the tree and the pruned version used 2 paths.

# From Trees to Rules

Tree-based models consist of one or more nested **if-then** statements for the predictors that partition the data.

Within these partitions, a model is used to predict the outcome.

For example, a very simple tree could be defined as:

```
if >= 1.7 then  
| if X2 >= 202.1 then Y = 1.3  
| else Y = 5.6  
else Y = 2.5
```



# From Trees to Rules

Notice that the *if-then* statements generated by a tree define a unique route to one terminal node for any sample.

A *rule* is a set of *if-then* conditions (possibly created by a tree) that have been collapsed into independent conditions.

For the example above, there would be three rules:

```
if X1 >= 1.7 & X2 >= 202.1 then Y = 1.3
if X1 >= 1.7 & X2 < 202.1 then Y = 5.6
if X1 < 1.7 then Y = 2.5
```

Rules can be simplified or pruned in a way that samples are covered by multiple rules, eg.

```
if X1 >= 1.7
```

# Rule-Based Models

One path to a terminal node in an unpruned model is

```
sqft <= 1594 &  
zip not in {z95631, z95833, z95758, z95670, 45 others} &  
beds <= 2.5 &  
latitude > 38.543 &  
latitude > 38.615 &  
latitude > 38.637 &  
latitude <= 38.688 &  
latitude <= 38.673
```

We can convert our model tree to a rule-based model. Many conditions can be simplified

## “Separate and Conquer” Approach to Rules

First, an initial model tree is created and only the rule with the largest coverage is saved from this model.

The samples covered by the rule are removed from the training set and another model tree is created with the remaining data.

Again, only the rule with the maximum coverage is retained.

This process repeats until all the training set data has been covered by at least one rule.

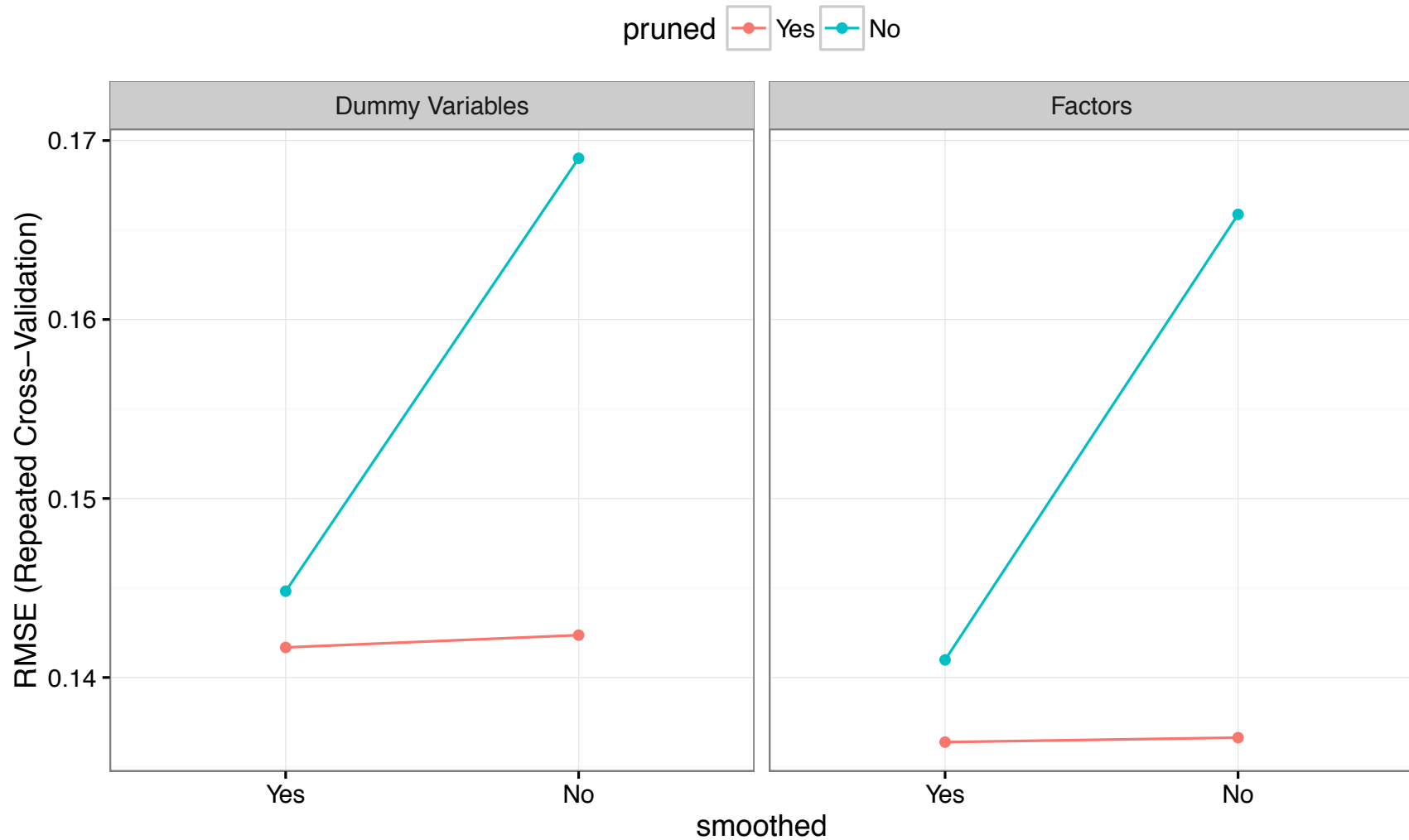
A new sample is predicted by determining which rule(s) it falls under then applies the linear model associated with the largest coverage.

For our data, the unpruned model has 81 and can be reduced shown to two rules based on  $\text{sqft} \leq 1594$ .

# An Example of a Terminal Node Model

```
log10(price) =  
  - 6.3388  
  - 0.0032 * city in {GALT, POLLOCK_PINES, ..., GRANITE_BAY}  
  + 0.0209 * zip in {z95820, z95822, z95626, ..., z95746}  
  + 0.0015 * zip in {z95673, z95832, z95621, ..., z95746}  
  + 0.0098 * zip in {z95631, z95833, z95758, ..., z95746}  
  + 0.0091 * zip in {z95818, z95608, z95662, ..., z95746}  
  + 0.0033 * zip in {z95814, z95765, z95667, ..., z95746}  
  + 0.0005 * beds  
  + 0.0001 * sqft  
  + 0.3097 * latitude  
  + 0.0056 * longitude
```

# Effect of Smoothing and Pruning Results



# Model Trees in R

```
> library(RWeka)
> model_tree <- M5P(log10(price) ~ ., data = training,
+                 ## Make the minimum number of instances per
+                 ## leaf higher than the default of 4
+                 control = Weka_control(M = 15))
>
> model_tree_unpruned <- M5P(log10(price) ~ ., data = training,
+                            control = Weka_control(M = 15, N = TRUE))
```

Note that the formula method is used but factors are *not* converted to dummy variables.

# Tuning Model Trees in R

```
> ctrl <- trainControl(method = "repeatedcv", repeats = 5)
>
> mt_grid <- expand.grid(rules = "Yes",
+                       pruned = c("No", "Yes"),
+                       smoothed = c("No", "Yes"))
>
> ## will use dummy variables:
> set.seed(139)
> mt_tune_dv <- train(log10(price) ~ ., data = training,
+                   method = "M5",
+                   tuneGrid = mt_grid,
+                   trControl = ctrl)
> ## will not:
> set.seed(139)
> mt_tune <- train(x = training[, -7], y = log10(training$price),
+                method = "M5",
+                tuneGrid = mt_grid,
+                trControl = ctrl)
```

Setting the seed prior to each call ensures that the same resamples are used.

# Cubist

Some specific differences between Cubist and the previously described approaches for model trees and their rule-based variants are:

- the specific techniques used for linear model smoothing, creating rules and pruning are different,
- an optional boosting-like procedure called *committees* can be used, and
- the predictions generated by the model rules can be adjusted using nearby points from the training set data.

We are indebted to the work of Chris Keefer, who extensively studied the Cubist source code to figure out the details.



# Cubist

Cubist does not use the Separate and Conquer approach to creating rules from trees.

A single tree is created then “flattened” into a set of rules.

The pruning and smoothing procedures are similar to those implemented in M5, but ...

# Smoothing Models in Cubist

Cubist has a different formula for combining models up the tree:

$$\hat{y}_{par} = a \times \hat{y}_{kid} + (1 - a) \times \hat{y}_{par}$$

where

$$a = \frac{Var(\hat{y}_{par}) - b}{Var(\hat{y}_{par}) + Var(\hat{y}_{kid}) - 2b}$$

$$b = \frac{S_{11} - \frac{1}{n}S_1S_2}{n - 1}$$

$$S_1 = \sum_{i=1}^n (y_i - \hat{y}_{i_{par}})$$

$$S_2 = \sum_{i=1}^n (y_i - \hat{y}_{i_{kid}})$$

$$S_{12} = \sum_{i=1}^n (y_i - \hat{y}_{i_{kid}})(y_i - \hat{y}_{i_{par}})$$

# Cubist in R

```
> library(Cubist)
> cb <- cubist(x = training[, -7], y = log10(training$price))
> ## To see the rules + models
> summary(cb)
```

# Cubist Base–Model Results

A basic cubist model resulted in 8 rules. For example:

Rule 1: [58 cases, mean 4.980517, range 4.477121 to 5.523746, est err 0.152796]

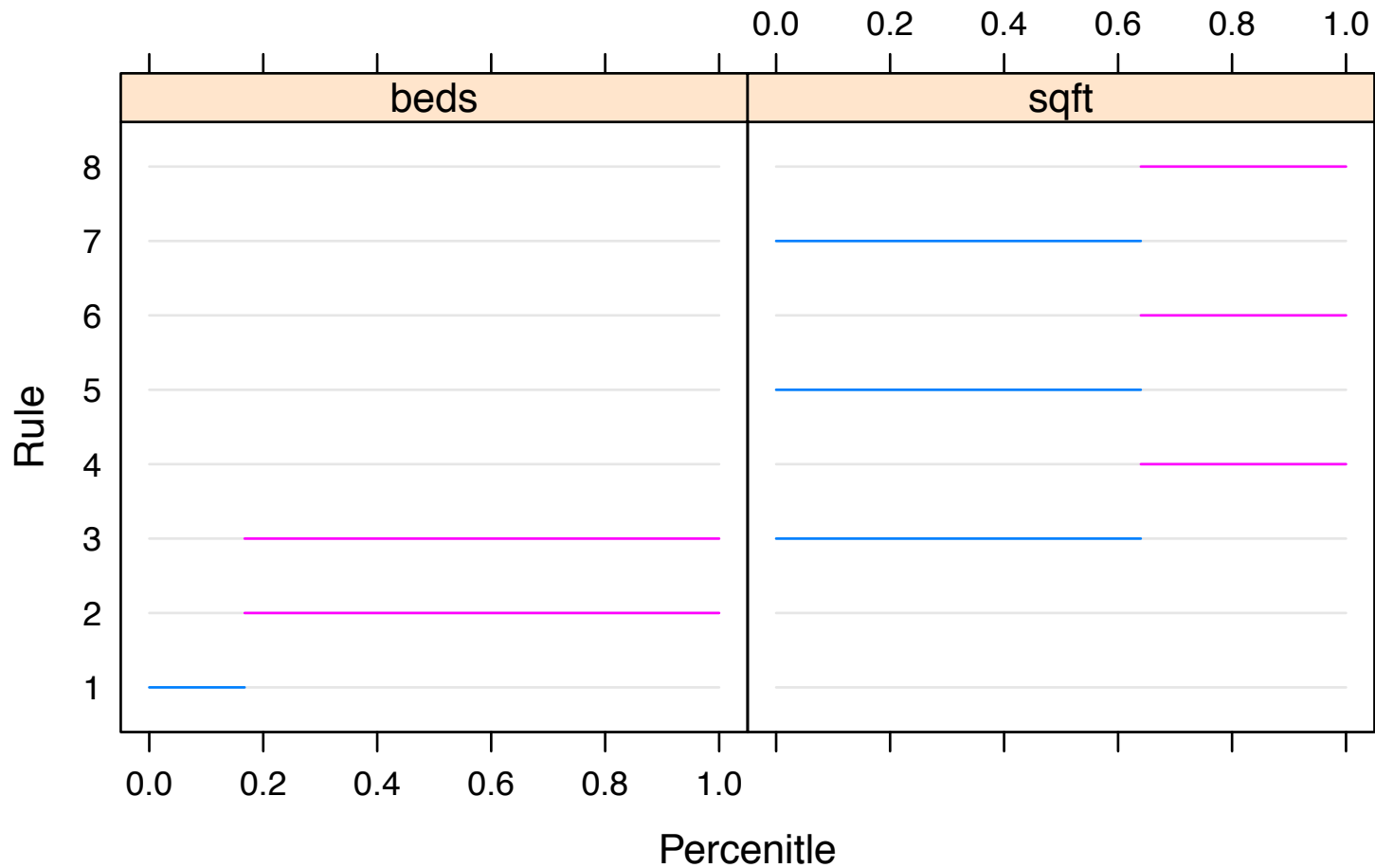
```
if zip in {z95621, z95626, z95660, z95673, z95683, z9581, z95817, z95820,  
          z95822, z95823, z95824, z95826, z95827, z95828, z95832, z95838,  
          z95841, z95842, z95843} and  
  beds <= 2 then  
outcome = 7.944631 + 0.323 beds + 4e-05 sqft + 0.03 longitude
```

Rule 2: [126 cases, mean 5.200466, range 4.788875 to 5.662758, est err 0.090147]

```
if zip in {z95626, z95660, z95683, z95815, z95823, z95824, z95827, z95832,  
          z95838, z95841} and  
  beds > 2 then  
outcome = 8.524561 - 0.056 beds + 0.000342 sqft + 0.03 longitude + 0.003 baths
```

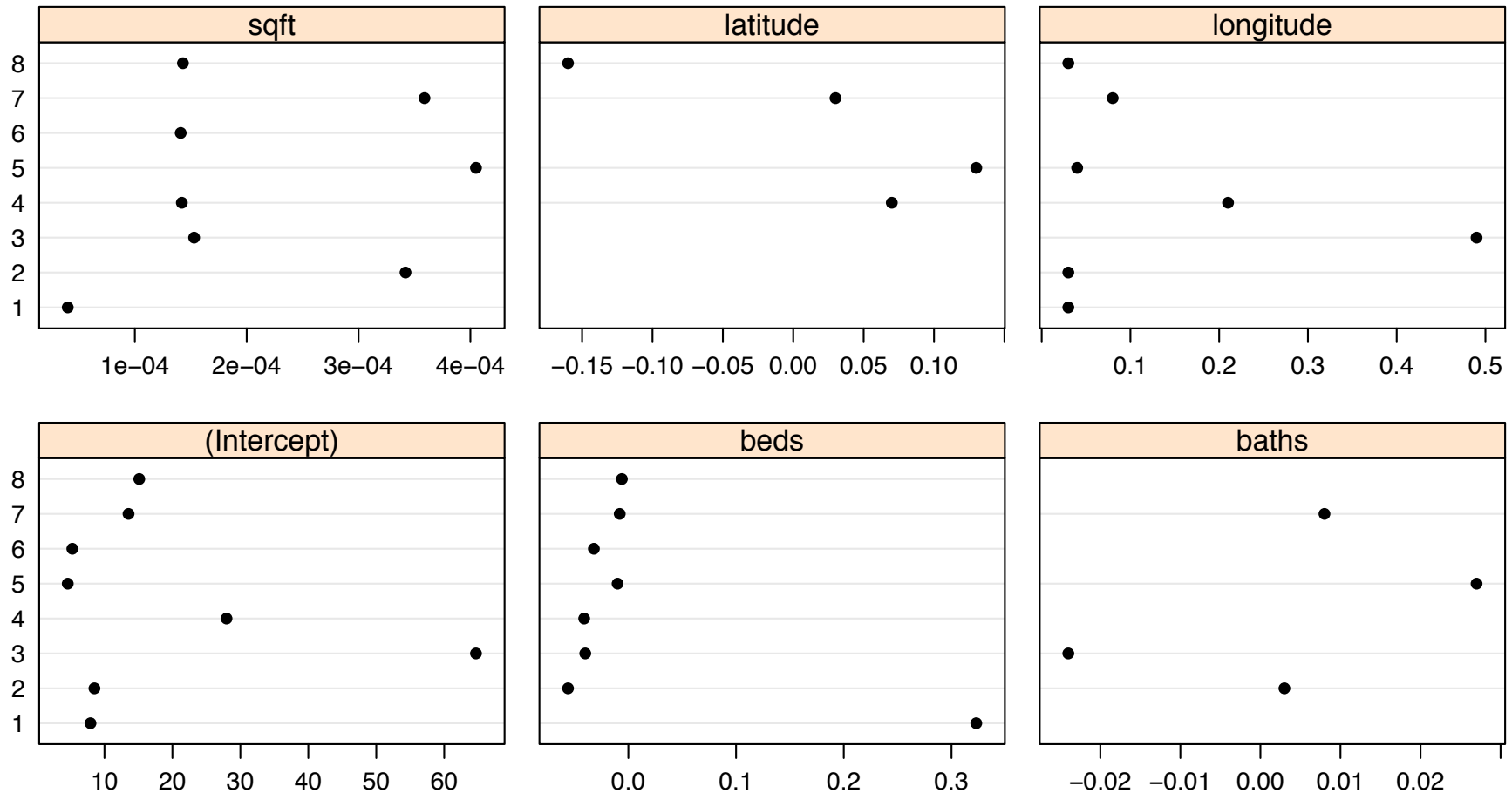
# Plotting the Splits

```
> dotplot(cb)
```



# Plotting the Slopes

```
> dotplot(cb, what = "coefs")
```



# Cubist Committees

Model committees can be created by generating a sequence of rule-based models (similar to boosting).

The training set outcome is adjusted based on the prior model fit and then builds a new set of rules using this pseudo-response.

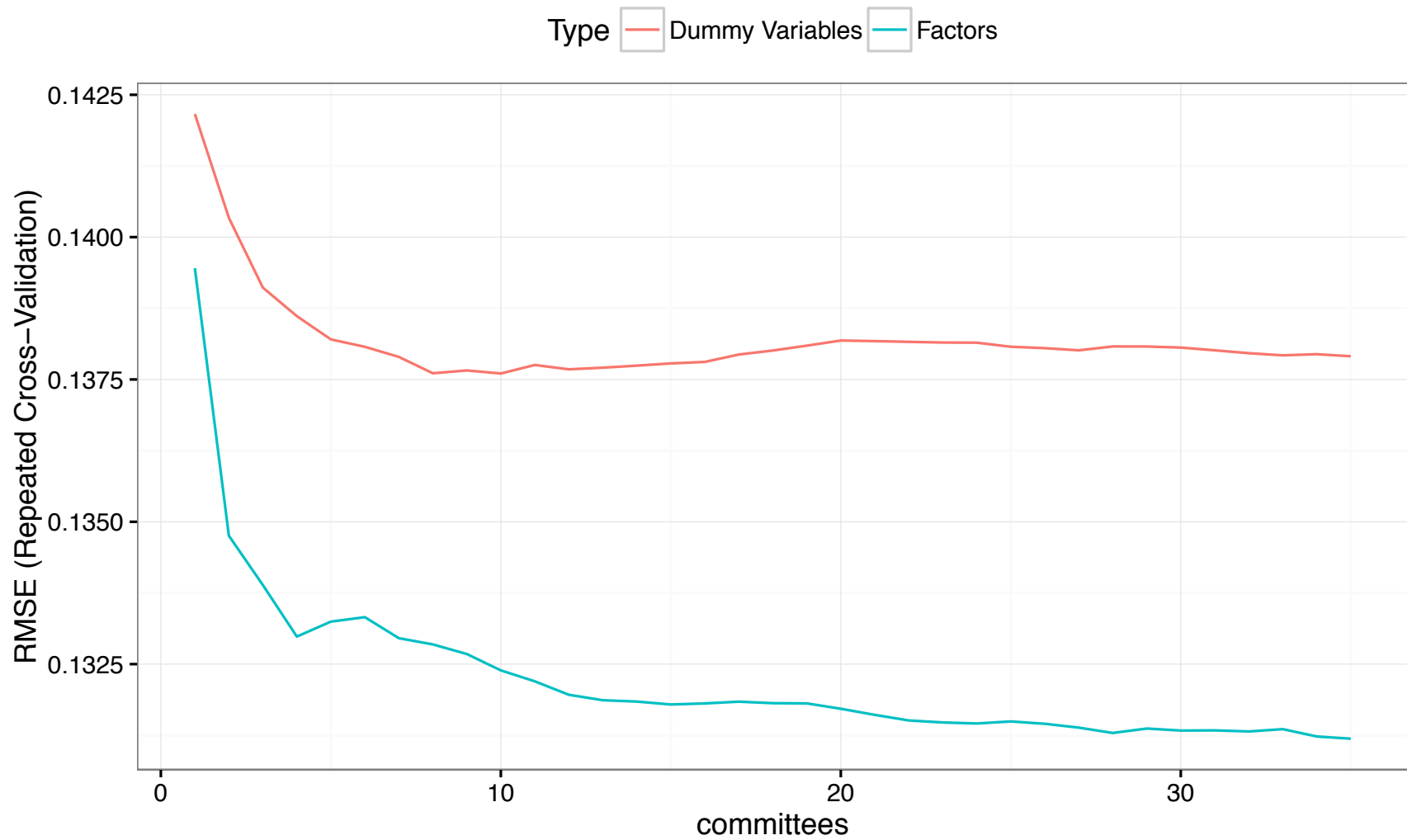
Specifically, the  $k^{th}$  committee model uses an adjusted response:

$$y_{i(k)} = 2y_{i(k-1)} - \hat{y}_{i(k-1)}$$

Once the full set of committee models are created, new samples are predicted using each model and the final rule-based prediction is the simple average of the individual model predictions.

```
> cb <- cubist(x = training[, -7], y = log10(training$price), committees = 17)
```

# Committee Results





# Neighbor–Based Adjustments

Cubist has the ability to adjust the model prediction using samples from the training set (Quinlan 1993).

When predicting a new sample, the  $K$  most similar neighbors are determined from the training set.

$$\hat{y} = \frac{1}{K} \sum_{\ell=1}^K w_{\ell} [(t_{\ell} - \hat{t}_{\ell}) + \hat{y}]$$

$t_{\ell}$  is the observed outcome for a training set neighbor,

$\hat{t}_{\ell}$  is the model prediction of that neighbor and

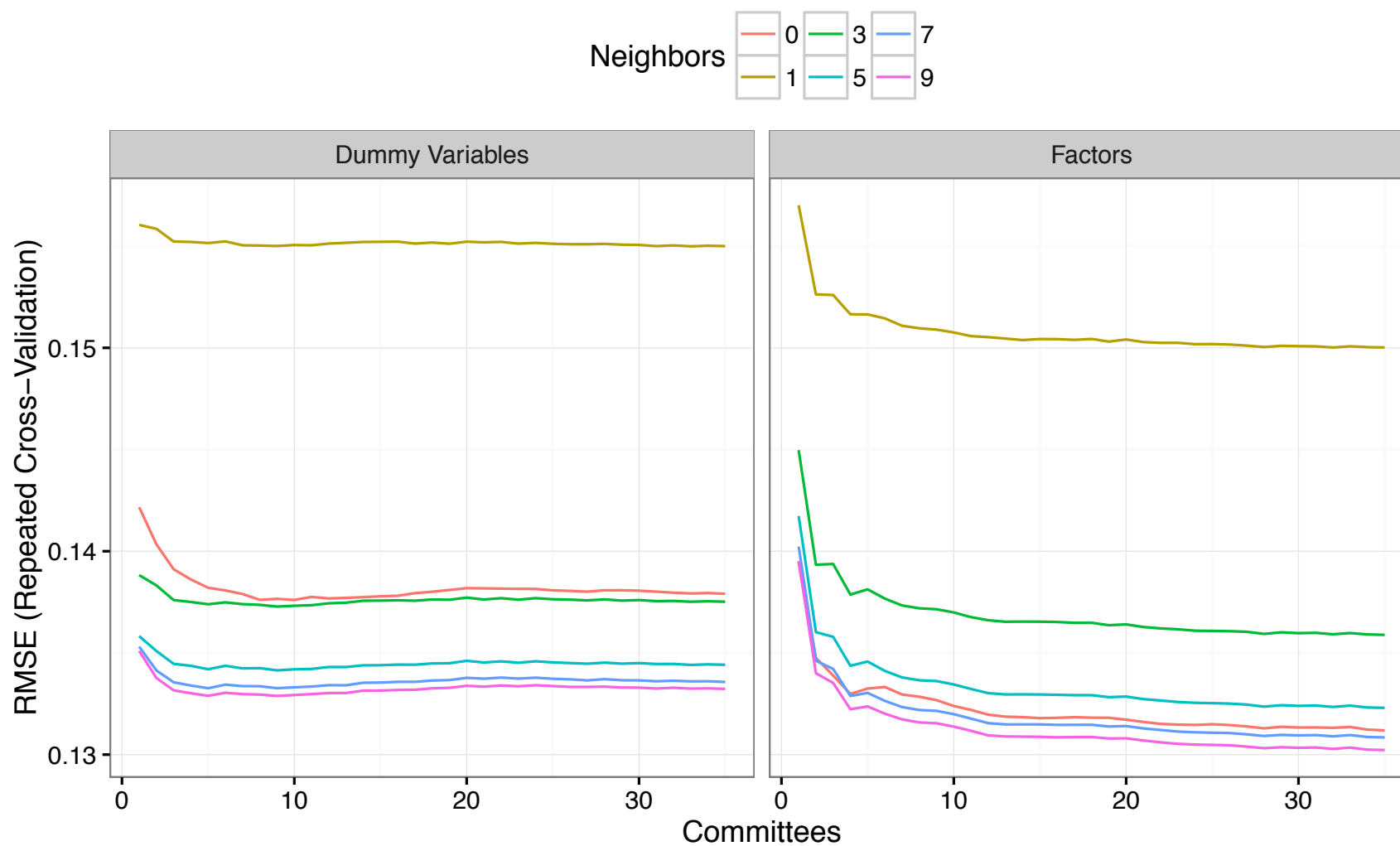
$w_{\ell}$  is a weight calculated using the distance of the training set neighbors to the new sample.

```
> predict(cb, newdata = testing, neighbors = 4)
```

# Tuning Model Trees in R

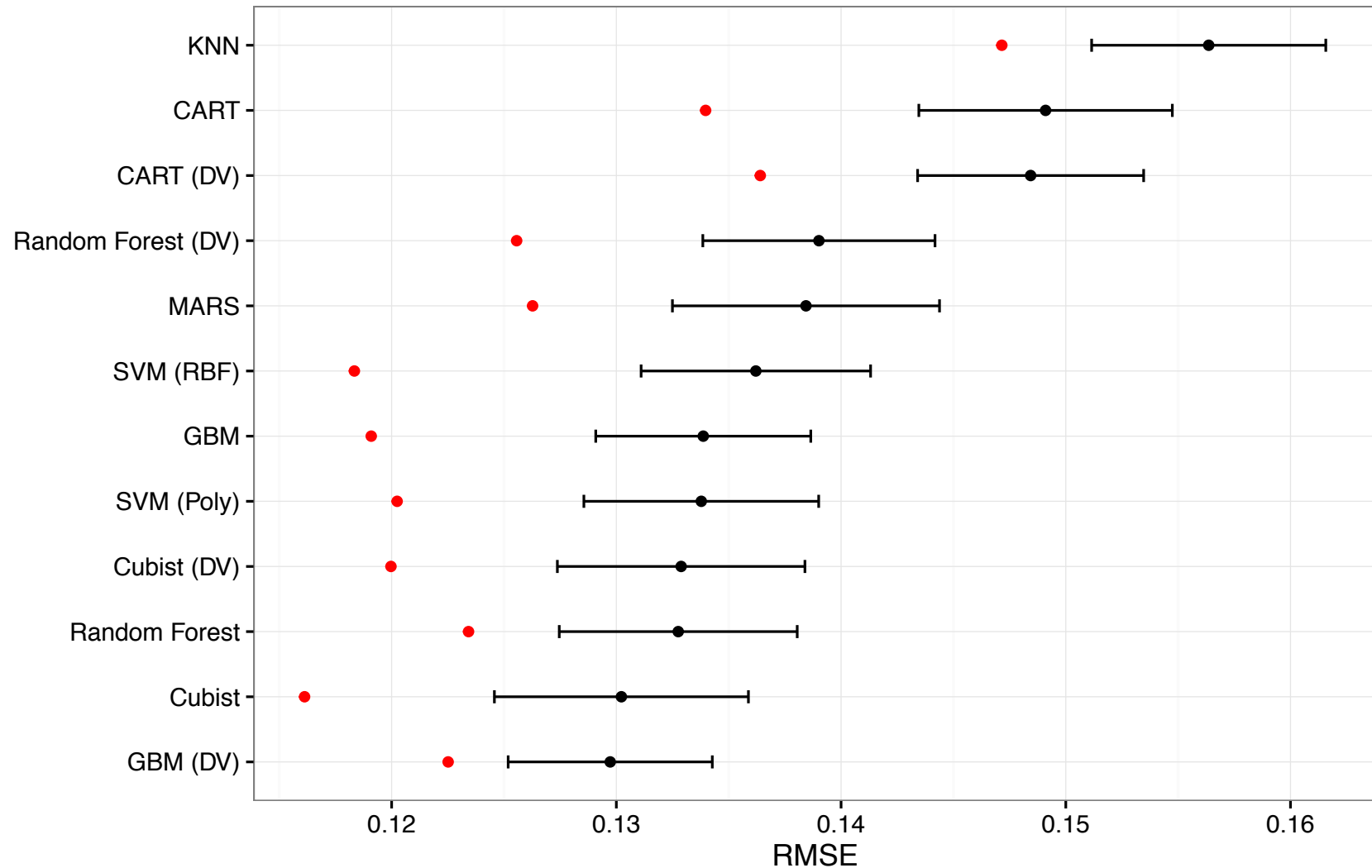
```
> cb_grid <- expand.grid(committees = c(1:35), neighbors = c(0, 1, 3, 5, 7, 9))
> set.seed(139)
> cb_tune_dv <- train(log10(price) ~ ., data = training,
+                   method = "cubist",
+                   tuneGrid = cb_grid,
+                   trControl = ctrl)
> set.seed(139)
> cb_tune <- train(x = training[, -7], y = log10(training$price),
+                method = "cubist",
+                tuneGrid = cb_grid,
+                trControl = ctrl)
> ggplot(cb_tune) ## to see the profiles
```

# Results with Neighbor Correction

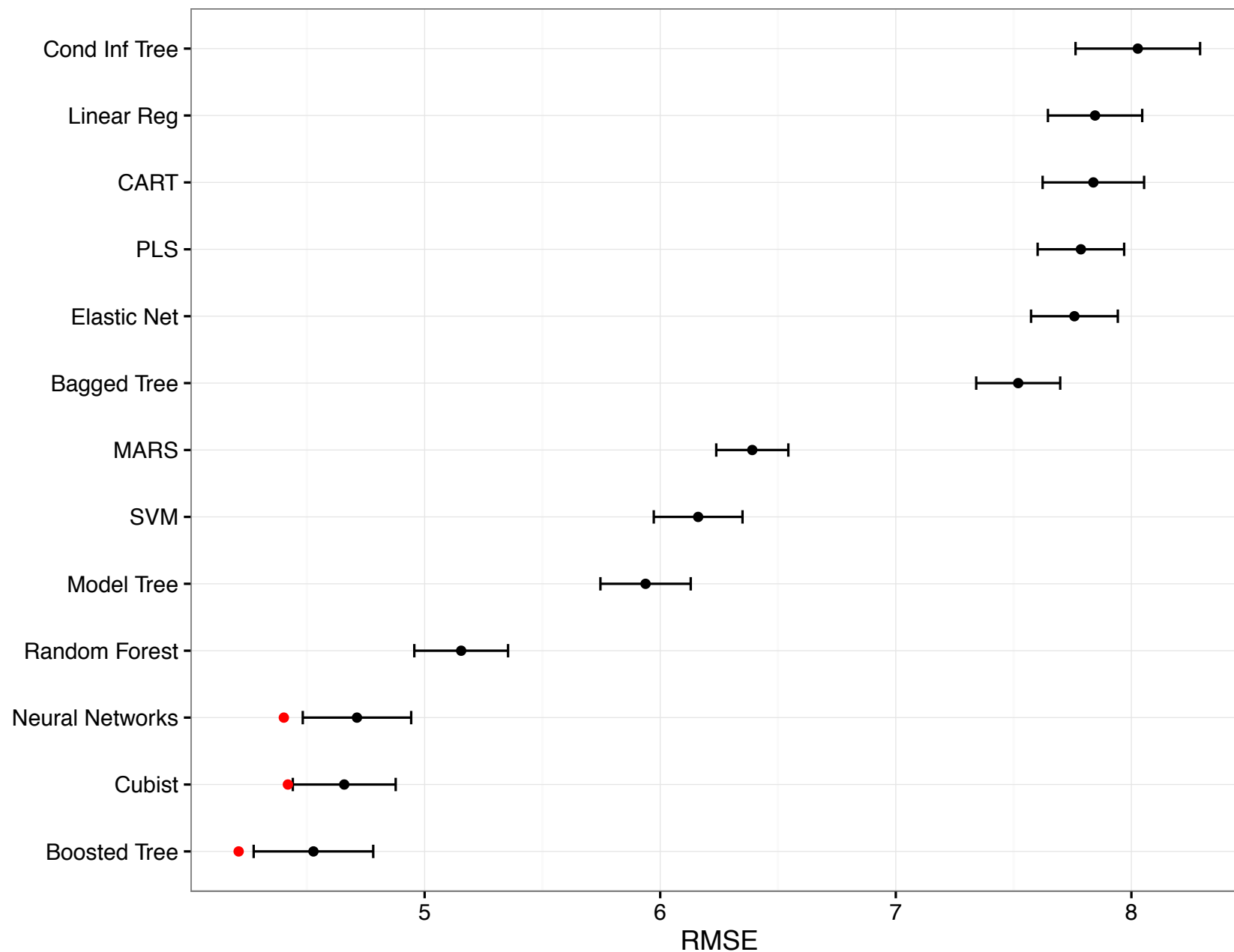


# Comparisons with Other Models

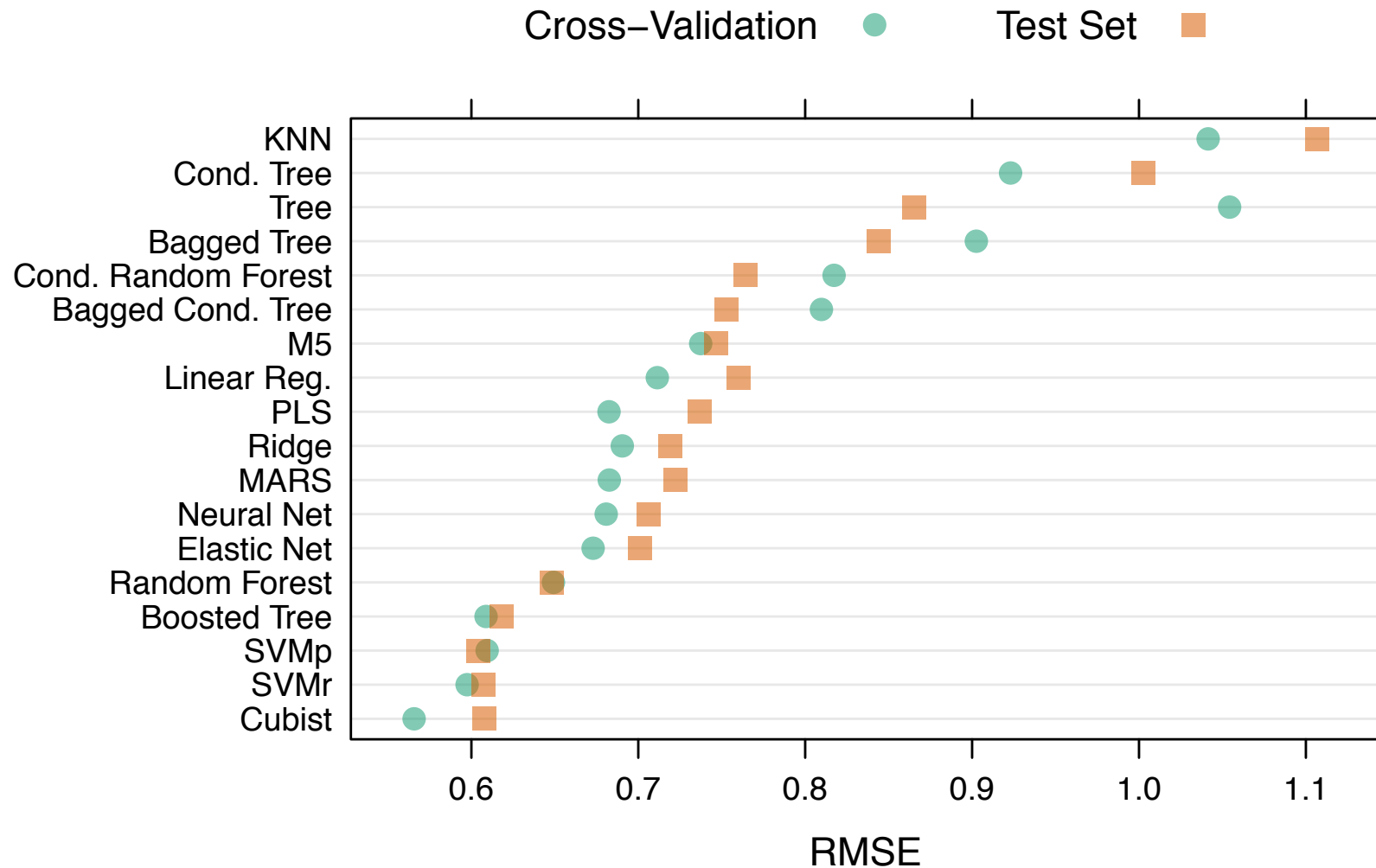
Test results in red



# Results with *APM's* Concrete Data Analysis

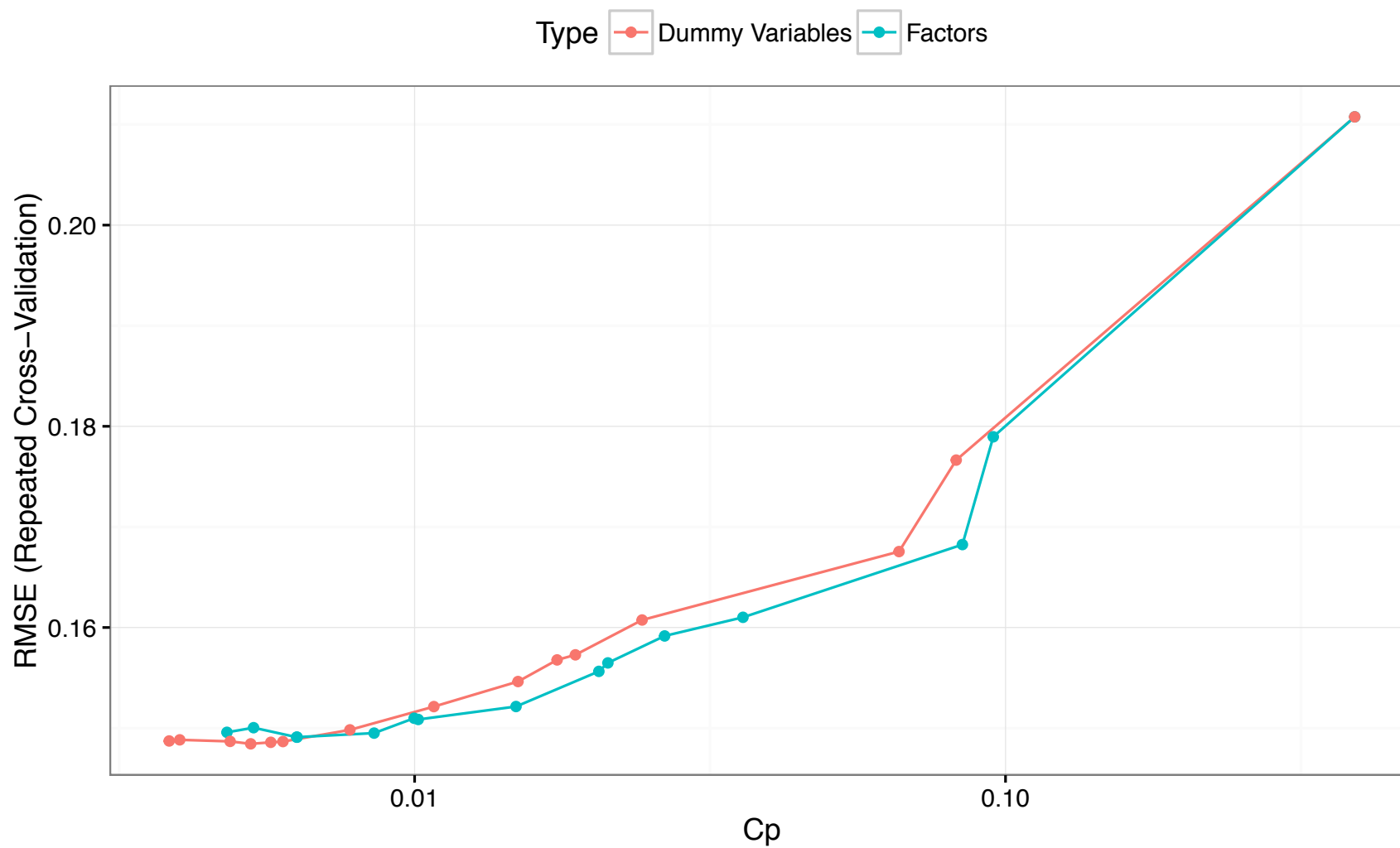


# Results with *APM's* Solubilty Data Analysis



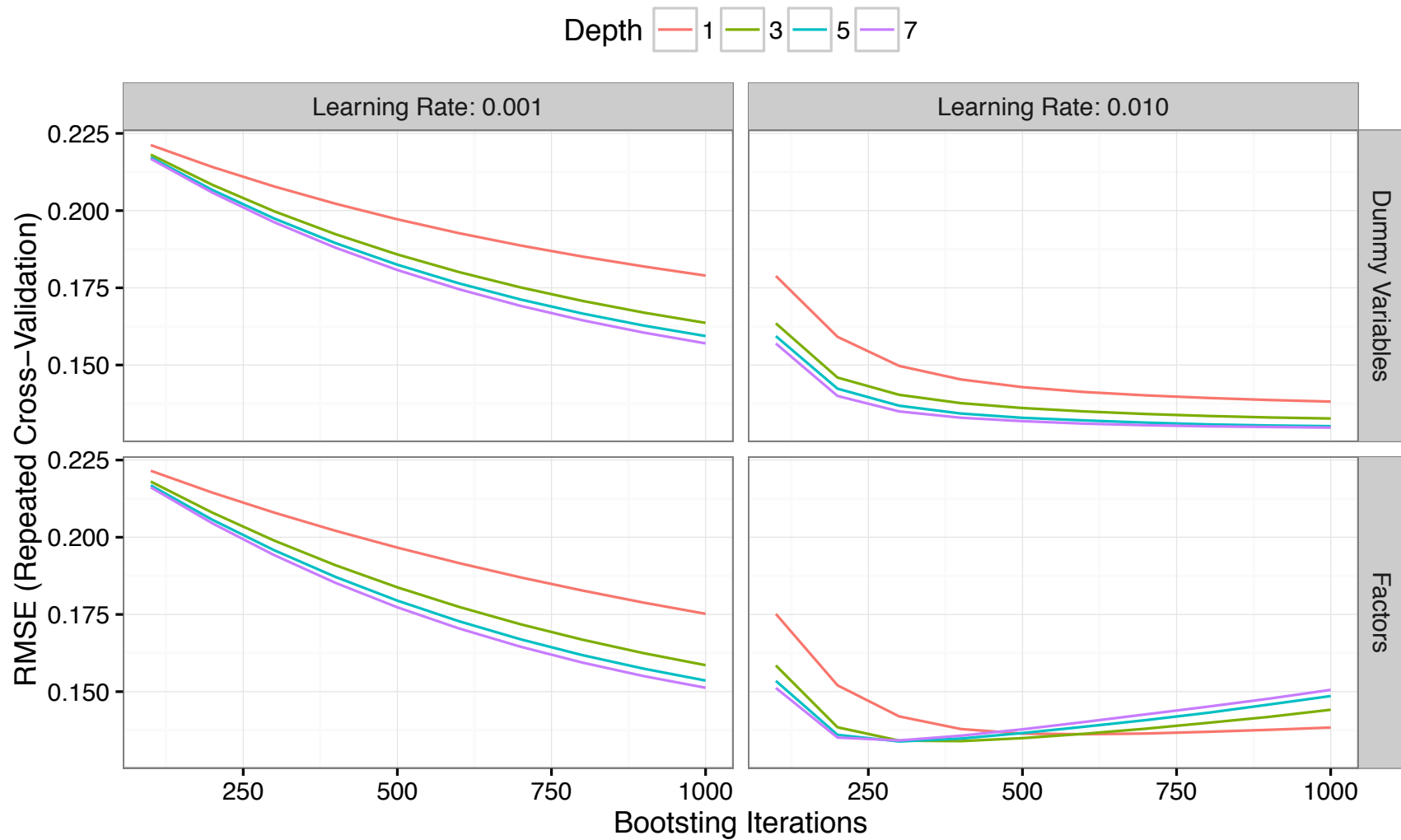
# Backup Slides

# CART Profiles

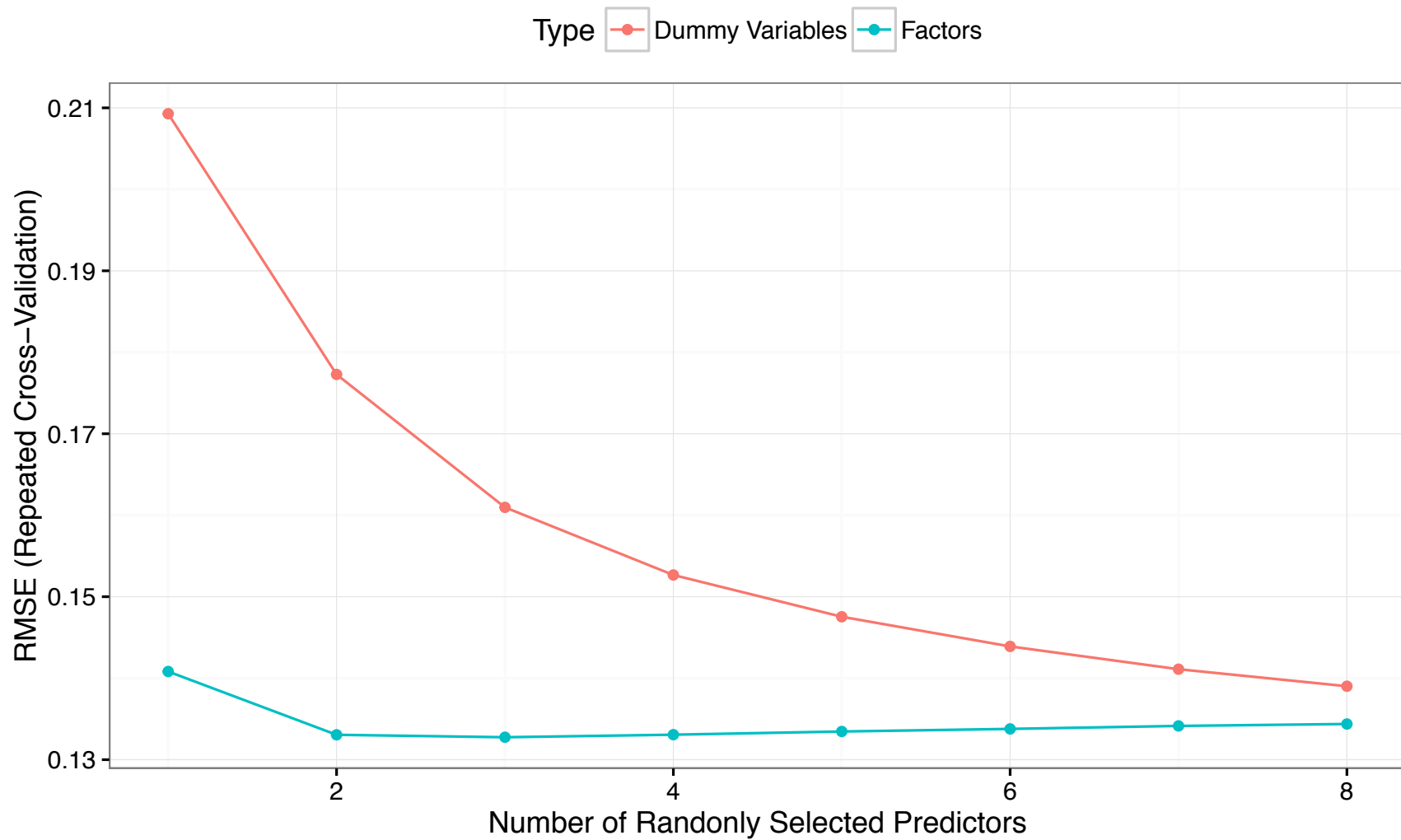




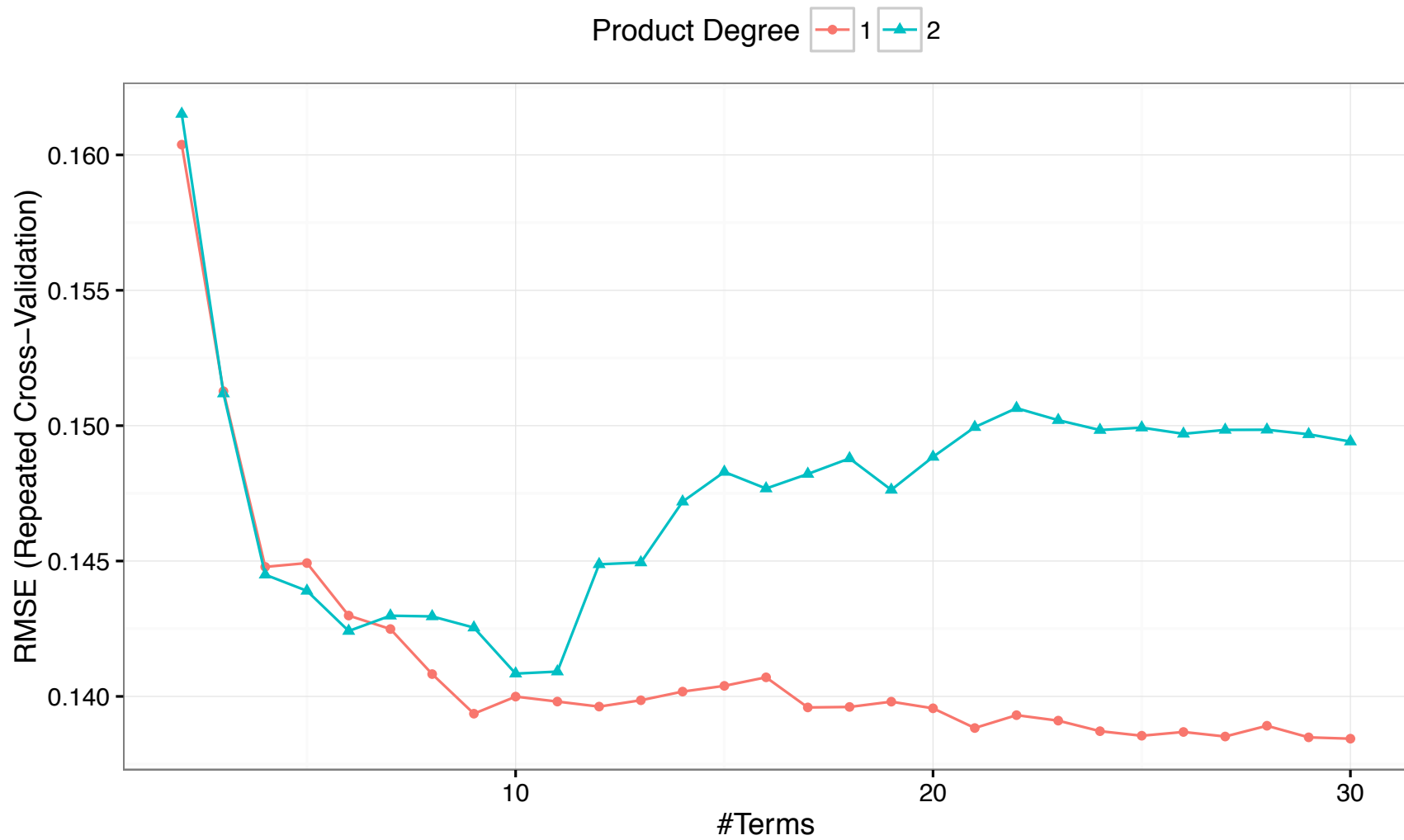
# Boosted Tree Profiles



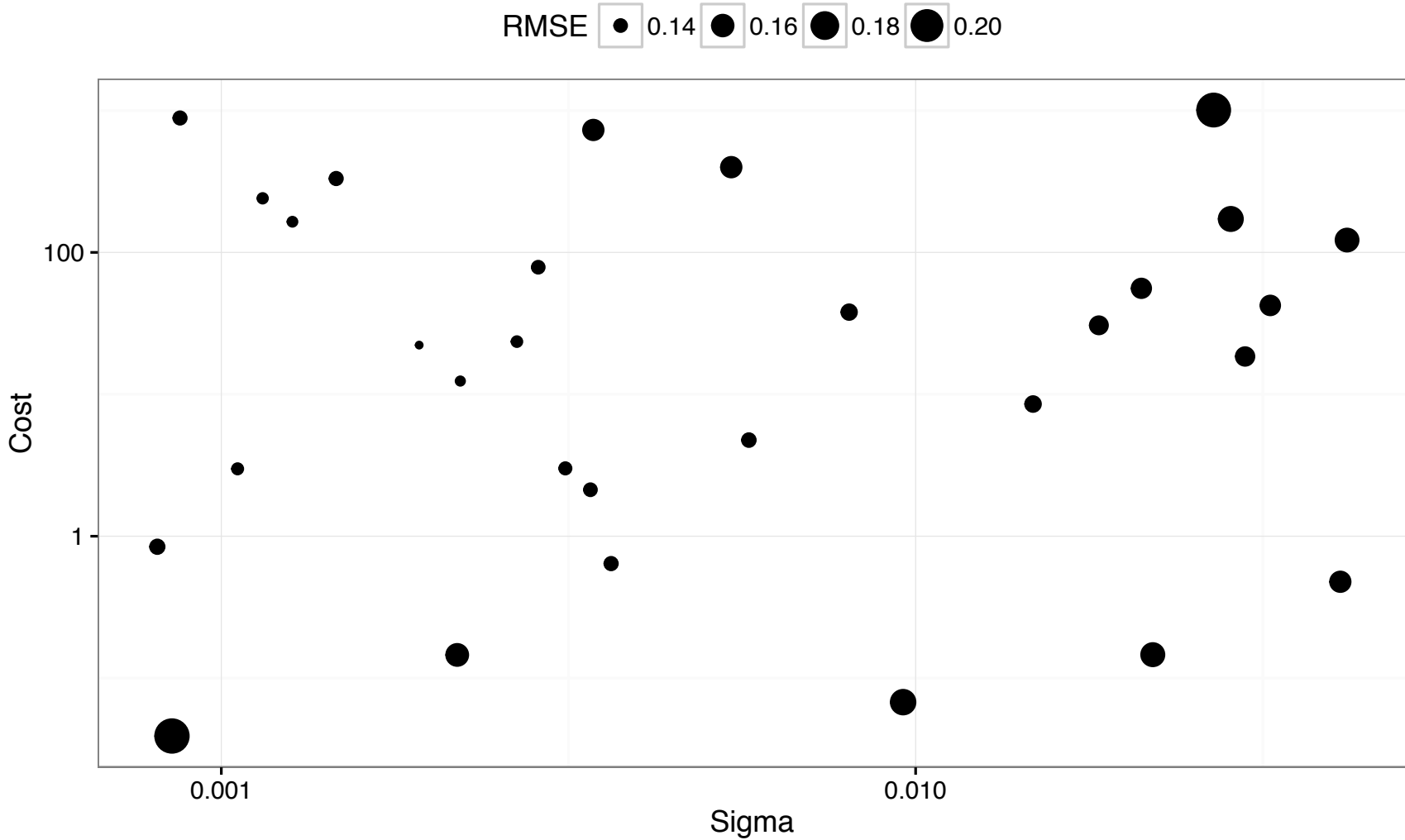
# Random Forest Profiles



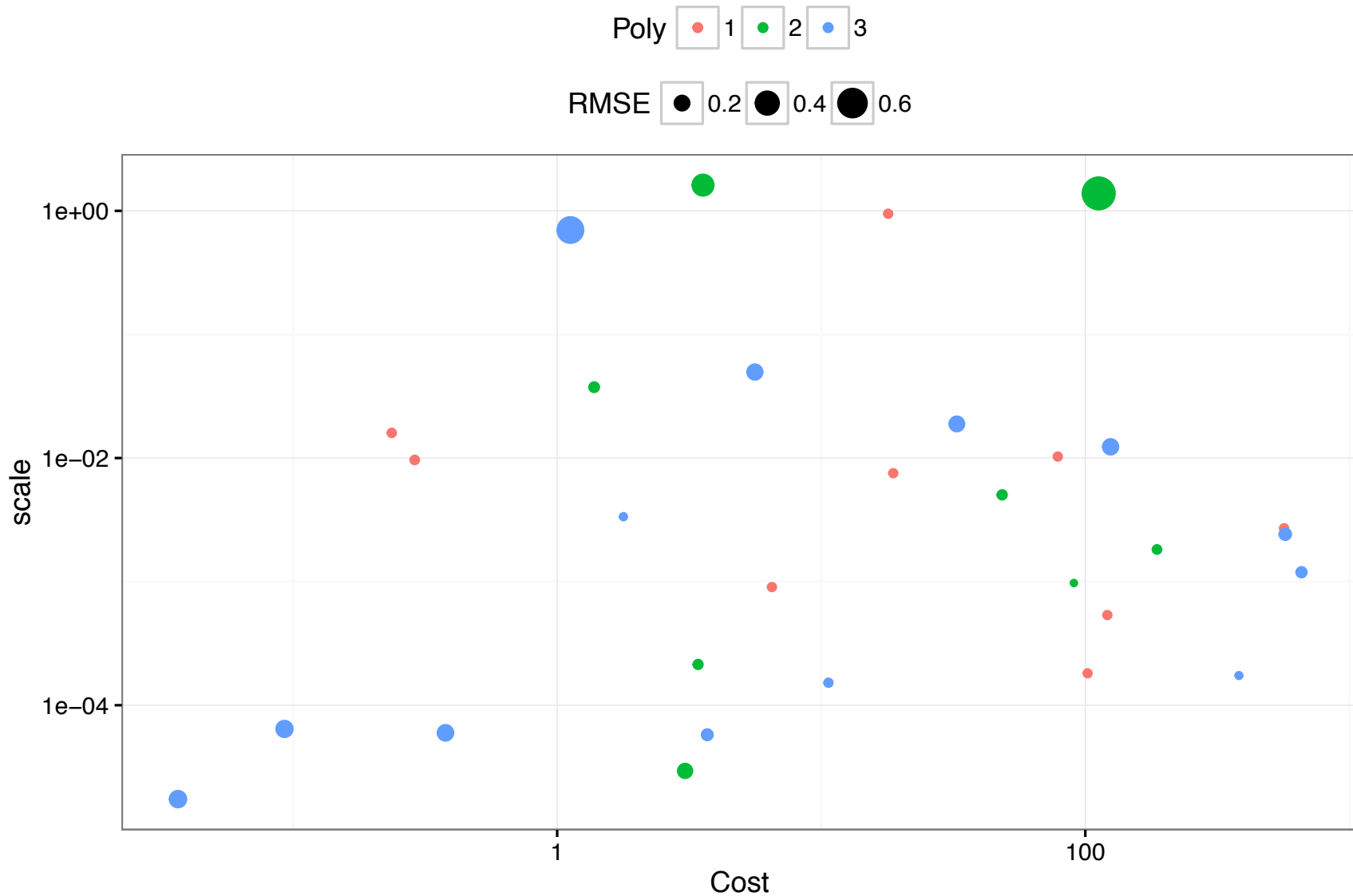
# MARS Profiles



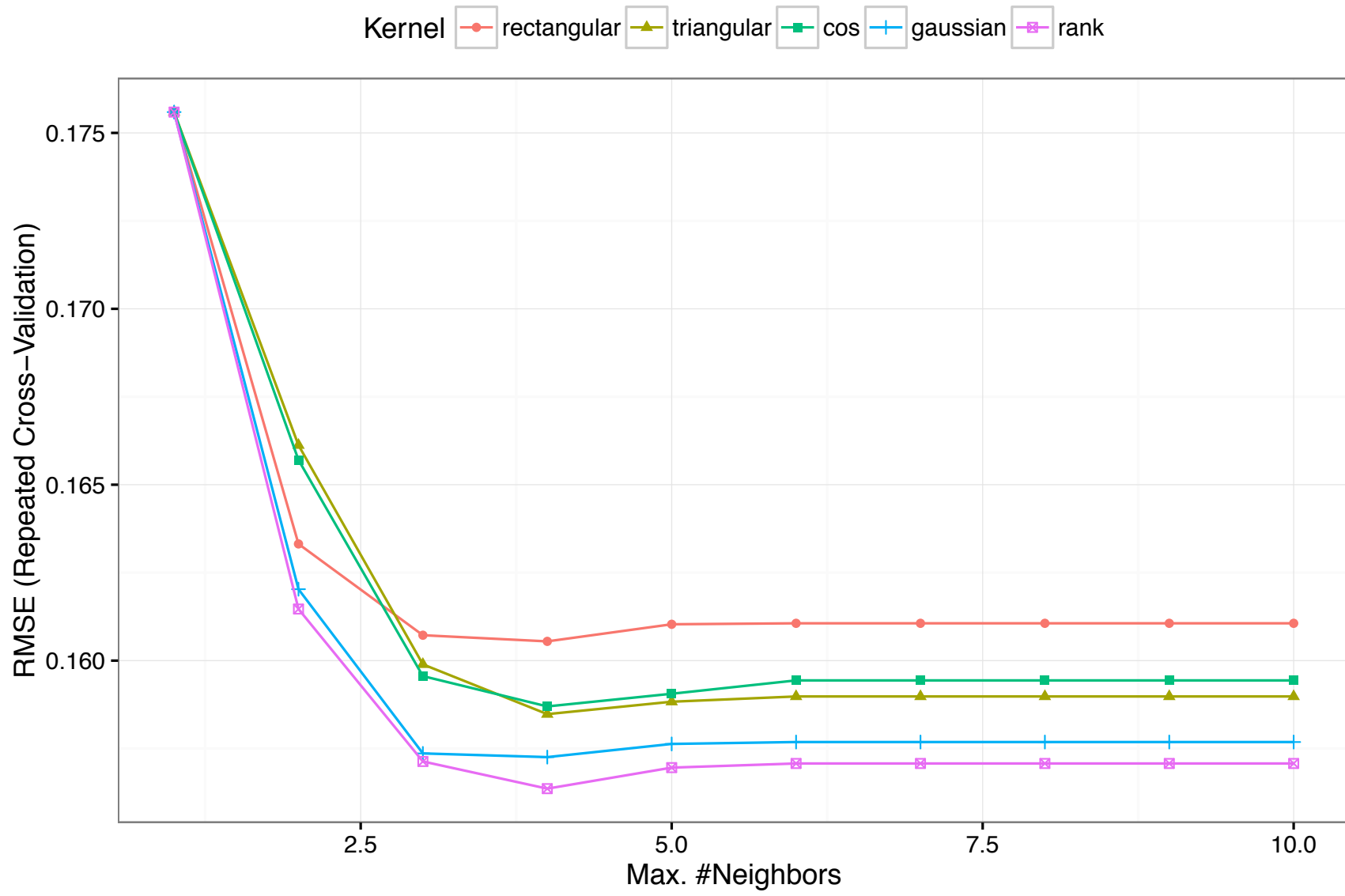
# SVM (RBF) Profiles using Random Search



# SVM (Poly) Profiles using Random Search



# KNN Profiles



# Thanks

Kirk Mettler for the invitation to speak tonight

Chris Keefer for his work with Cubist

Steve Weston, Chris Keefer and Nathan Coulter for adapting the Cubist C code to R.

# References

Quinlan R (1992). "Learning with Continuous Classes." Proceedings of the 5th Australian Joint Conference On Artificial Intelligence, pp. 343-348.

Quinlan R (1993). "Combining InstanceBased and ModelBased Learning." Proceedings of the Tenth International Conference on Machine Learning, pp. 236-243.

Wang Y, Witten I (1997). "Inducing Model Trees for Continuous Classes." Proceedings of the Ninth European Conference on Machine Learning, pp. 128-137.