

[9] Big Data: Trees

Matt Taddy, University of Chicago Booth School of Business

`faculty.chicagobooth.edu/matt.taddy/teaching`

Decision Trees: Using tree-logic to make predictions.

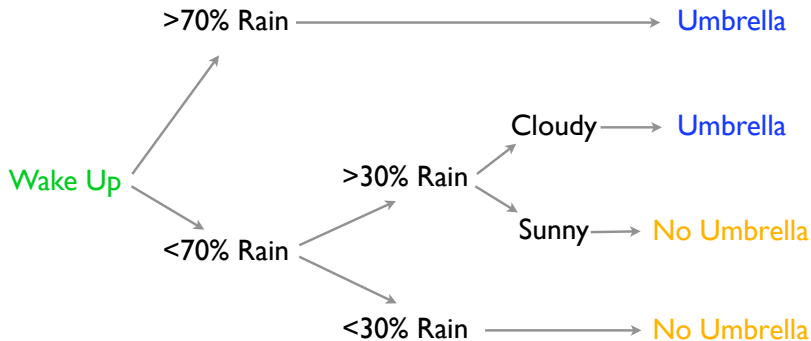
CART: Classification and Regression Trees.

Random Forests: averaging over many possible trees.

Simple examples: NBC, prostate cancer, motorcycle crashes.

Larger example on house prices in california.

What is a Decision Tree?



Tree-logic uses a series of steps to come to a conclusion.
The trick is to have mini-decisions combine for good choices.
Each decision is a node, and the final prediction is a **leaf node**

Decision Trees are a Regression Model

You have inputs \mathbf{x} (forecast, current conditions)
and an output of interest y (need for an umbrella).

Based on previous data, the goal is to specify branches of choices that lead to good predictions in new scenarios.

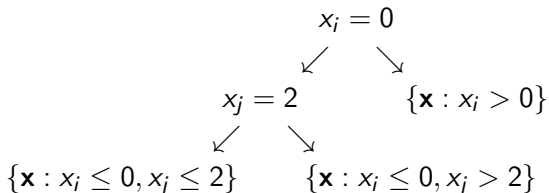
In other words, you want to estimate a **Tree Model**.

Instead of linear coefficients, we need to find 'decision nodes':
split-rules defined via thresholds on some dimension of \mathbf{x} .

Nodes have a parent-child structure: every node except the root has a parent, and every node except the leaves has two children.

Decision trees are like a game of mousetrap

You drop your \mathbf{x} covariates in at the top, and each decision node bounces you either left or right. Finally, you end up in a **leaf node** which contains the data subset defined by these decisions (splits).



The **prediction rule** at each leaf (a class probability or predicted \hat{y}) is the average of the sample y values that ended up in that leaf.

Estimation of Decision Trees

As usual, we'll maximize data likelihood (minimize deviance).

But what are the observation probabilities in a tree model?

Two types of likelihood: classification and regression trees.

A given covariate \mathbf{x} dictates your path through tree nodes, leading to a **leaf node** at the end.

Classification trees have **class probabilities** at the leaves.

Probability I'll be in heavy rain is 0.9 (so take an umbrella).

Regression trees have a **mean response** at the leaves.

The expected amount of rain is 2in (so take an umbrella).

Tree deviance is the same as in linear models

Regression Deviance: $\sum_{i=1}^n (y_i - \hat{y}_i)^2$

Classification Deviance: $-\sum_{i=1}^n \log(\hat{p}_{y_i})$

It is also common to use Gini Deviance, $\sum_{i=1}^n \hat{p}_{y_i}(1 - \hat{p}_{y_i})$

Instead of being based on $\mathbf{x}'\beta$, predicted \hat{p} and \hat{y} are functions of \mathbf{x} passed through the decision nodes.

We need a way to estimate the sequence of decisions.

- How many are they? What is the order?

There is a *huge* set of possible tree configurations.

Given a **parent** set of data $\{\mathbf{x}_i, y_i\}_{i=1}^n$, the **optimal split** is that location x_{ij} on some dimension j on some observation i , so that the **child** sets

$$\text{left: } \{\mathbf{x}_k, y_k : x_{kj} \leq x_{ij}\} \quad \text{and} \quad \text{right: } \{\mathbf{x}_k, y_k : x_{kj} > x_{ij}\}$$

are as homogeneous in response y as possible.

For example, we will minimize the sum of squared errors

$$\sum_{k \in \text{left}} (y_k - \bar{y}_{\text{left}})^2 + \sum_{k \in \text{right}} (y_k - \bar{y}_{\text{right}})^2$$

for *regression trees*, or gini impurity for *classification trees* (e.g., the sum across children 'c' of $n_c \bar{y}_c (1 - \bar{y}_c)$ if $y \in \{0, 1\}$).

We estimate decision trees by being recursive and greedy

CART grows the tree through a sequence of splits:

- ▶ Given any set (node) of data, you can find the **optimal split** (the error minimizing split) and divide into two child sets.
- ▶ We then look at each child set, and again find the optimal split to divide it into two homogeneous subsets.
- ▶ The children become parents, and we look again for the optimal split on their new children (the grandchildren!).

You stop splitting and growing when the size of the leaf nodes hits some minimum threshold (e.g., say no less than 10 obsv per leaf). Often there are also minimum deviance improvement thresholds.

Use the `tree` library for CART in R

The syntax is essentially the same as for `glm`:

```
mytree = tree(y ~x1 + x2 + x3 ..., data=mydata)
```

There are only a few other possible arguments, all of which dictate possible types of new children

- ▶ `mincut` is the minimum size for a new child.
- ▶ `mindev` is the minimum (proportion) deviance improvement for proceeding with a new split.

These are important: you may want to make them smaller than their defaults: `mincut=5`, `mindev=0.01`.

As usual, you can `print`, `summarize`, and `plot` the tree.

Recall our NBC data: viewer demographic % by show.

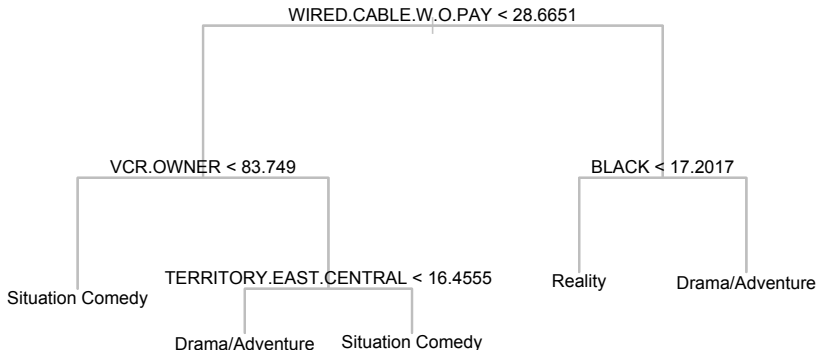
Consider a classification tree to predict genre from demographics.

```
> genretree
node), split, n, deviance, yval, (yprob)
  * denotes terminal node

1) root 40 75.800 Drama/Adventure ( 0.475 0.425 0.100 )
 2) CABLE.W.O.PAY < 28.6651 22 33.420 Drama/Adventure (0.73 0.09 0.18)
 4) VCR.OWNER < 83.749 5 6.730 Situation Comedy (0.00 0.40 0.60) *
 5) VCR.OWNER > 83.749 17 7.606 Drama/Adventure (0.941 0.000 0.059)
   10) TERRITORY.EAST.CENTRAL < 16.4555 16 0.000 Drama/Adventure(1 0 0) *
   11) TERRITORY.EAST.CENTRAL > 16.4555 1 0.000 Situation Comedy(0 0 1) *
 3) CABLE.W.O.PAY > 28.6651 18 16.220 Reality ( 0.16667 0.83333 0 )
   6) BLACK < 17.2017 15 0.000 Reality ( 0 1 0 ) *
   7) BLACK > 17.2017 3 0.000 Drama/Adventure ( 1 0 0 ) *
```

Output from `tree` shows a series of decision nodes and the proportion in each genre at these nodes, down to the leaves.

Trees are easiest to understand in a **dendrogram**



Shows the sequence of internal splits, ending in leaf-node decisions.
Here, the decision is the genre of highest proportion in each leaf.
To get the dendrogram, do `plot(mytree)` then `text(mytree)`.

Consider predicting **engagement** from **ratings** and **genre**.

Split on genre by turning it into a group of numeric variables:

```
x <- model.matrix(PE ~ Genre + GRP, data=nbc)[-1]
names(x) <- c("reality","comedy","GRP")
```

We have a reference factor level (drama)

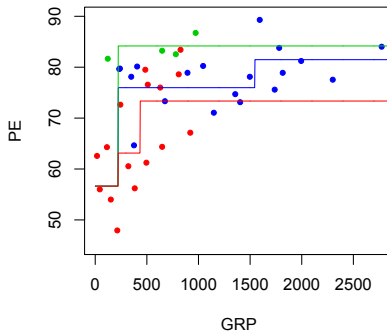
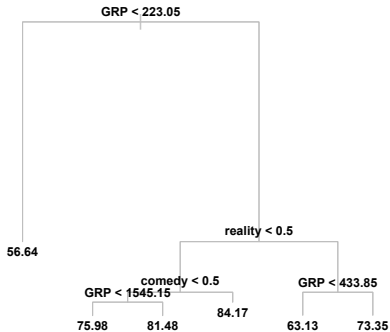
A regression tree:

```
nbctree <- tree(PE ~ ., data=x, mincut=1)
```

Instead of genre, leaf predictions are expected engagement.

`mincut=1` allows for leaves containing a single show,
with expected engagement that single show's PE.

An NBC Show Engagement Tree



Green is comedy, blue is drama, red is reality

Nonlinear: PE increases with GRP, but in jumps

Follow how the tree translates into changing $\mathbb{E}[\text{PE}]$

Trees provide **Automatic Interaction Detection**

For example, different genres are more/less dependent on GRP.

AID was an original motivation for building decision trees.

Older algorithms have it in their name: CHAID, US-AID, ...

This is pretty powerful technology: nonlinearity and interaction without having to specify it in advance.

Moreover, nonconstant variance is no problem.

Methods with these characteristics are called **nonparametric**.

No assumed parametric model (eg, $y = x\beta + \varepsilon$, $\varepsilon \sim N(0, \sigma^2)$).

Pruning your tree for cross validation

Biggest challenge with such flexible models is avoiding overfit. For CART, the usual solution is to rely on cross validation.

The basic constraints (`mincut`, `mindev`) lead to a full tree fit. Prune this tree by removing split rules from the bottom up:

At each step, remove the split that contributes least to deviance reduction, thus reversing CART's growth process.

Pruning yields candidate trees, and we use CV to choose. Each prune step produces a candidate tree model, and we can compare their out-of-sample prediction performance.

Example: Prostate Cancer Prognosis

After tumor detection, there are many treatment options.

- ▶ Various chemo + radiation, surgical removal.

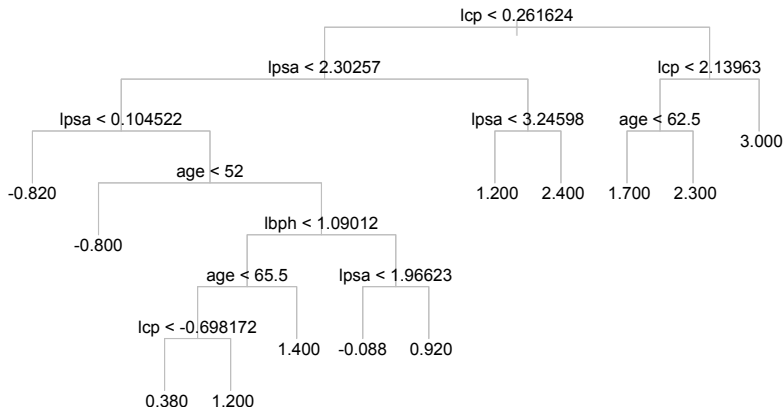
Biopsy information is available to help in deciding treatment

- ▶ **Gleason Score**: microscopic pattern classes.
- ▶ **Prostate Specific Antigen**: protein production.
- ▶ **Capsular Penetration**: reach of cancer into gland lining.
- ▶ **Benign Prostatic Hyperplasia Amount**: size of prostate.

Another influential variable is the patient's **age**.

The goal is to predict tumor log-volume (size, spread).

Full tree fit to 97 prostate cancer patients



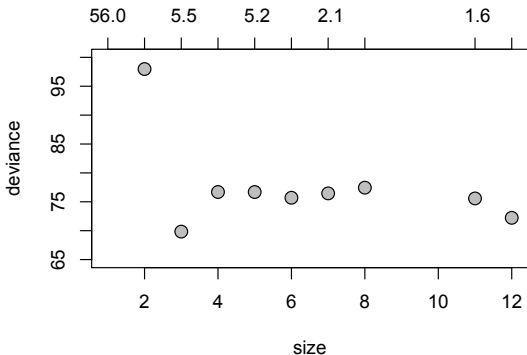
Leaf node labels are expected tumor $\log(\text{volume})$.

Do we need all the splits? Is the tree just fitting noise?

Cross-Validated Tree Pruning

`cv.tree` does cross-validation across pruning levels.

```
cvpst <- cv.tree(pstree, K=90) # K is nfolds
```



The output can be plotted, and it holds out-of-sample deviance for each tree size (the number of leaf nodes).

Pruning the Prostate Cancer Tree

Out-of-sample deviance can be used to choose tree size.

```
> cvpst$size
[1] 12 11  8  7  6  5  4  3  2  1
> cvpst$dev
[1] 72  75  77  76  76  77  77  70  97 160
```

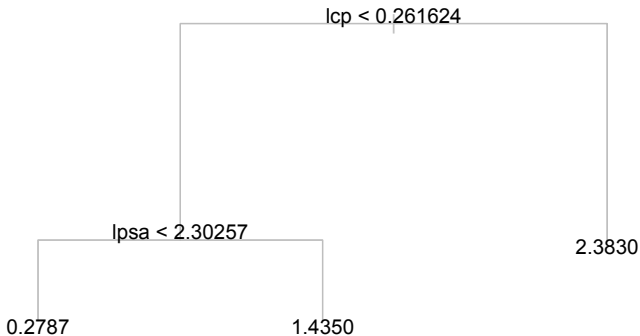
Since size 3 has lowest CV deviance, it is 'best'.

To fit this tree, use the `prune.tree` function:

```
pstcut <- prune.tree(pstree, best=3)
```

`pstcut` is then itself a new tree object.

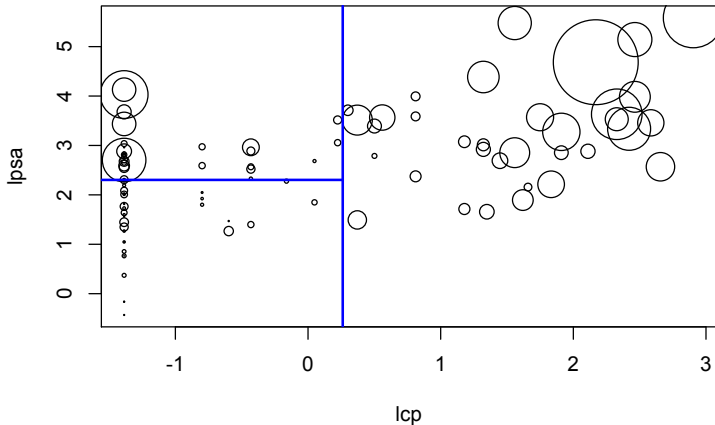
The Pruned Treatment Tree



CV chooses PSA and penetration as deciding variables.

Note the interaction: penetration effect depends on PSA.

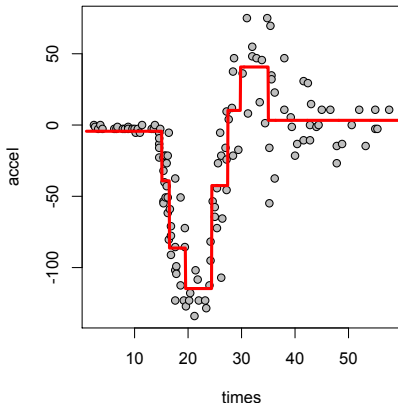
Prostate Cancer Prognosis Tree



With only 2 relevant inputs, we can plot the data and tree fit. Points proportional to tumor size, leaf partitions are in blue.

They automatically learn non-linear response functions and will discover interactions between variables.

x is time from impact, y is acceleration on the helmet.



Unfortunately, it is tough to avoid overfit with CART:

Deep tree structure is so unstable that optimal depth is not easily chosen via cross validation, and there's no theory to fall back on.

Instead, we can average over a **bootstrapped** sample of trees:

- ▶ repeatedly re-sample the data, **with-replacement**, to get a 'jittered' dataset of n observations.
- ▶ for each resample, **fit a CART tree**.
- ▶ when you want to predict y for some \mathbf{x} , take the **average** prediction from this forest of trees.

Real structure that persists across datasets shows up in the average. Noisy useless signals will average out to have no effect.

This is a Random Forest

Random Forests

- Sample B subsets of the data + variables:
e.g., observations 1, 5, 20, ...
- Fit a tree to each subset, to get B fitted trees is \mathcal{T}_b .
- Average prediction across trees:
 - for regression average $\mathbb{E}[y|\mathbf{x}] = \frac{1}{B} \sum_{b=1}^B \mathcal{T}_b(\mathbf{x})$.
 - for classification let $\{\mathcal{T}_b(\mathbf{x})\}_{b=1}^B$ vote on \hat{y} .

The observation resample is usually *with-replacement*, so that this is taking the *average of bootstrapped trees* (i.e., 'bagging')

Understanding **Random Forests**

Recall how CART is used in practice.

- ▶ Split to lower deviance until leaves hit minimum size.
- ▶ Create a set of candidate trees by pruning back from this.
- ▶ Choose the best among those trees by cross validation.

Random Forests avoid the need for CV.

Each tree ' b ' is not overly complicated because you only work with a limited set of variables.

Your predictions are not 'optimized to noise' because they are averages of trees fit to many different subsets.

RFs are a great go-to model for nonparametric prediction.

Model Averaging

This technique of 'Model Averaging' is central to many advanced nonparametric learning algorithms.

ensemble learning, mixture of experts, Bayesian averages, ...

Recall lasso as a stabilized version of stepwise regression (if you jitter the data your estimates stay pretty constant).

Model averaging is a way to take arbitrary *unstable* methods, and make them stable. **This makes training easier.**

Probability of rain on a new day is the average $P(\text{rain})$ across some trees that split on forecast, others on sky.

We don't get tied to one way of deciding about umbrellas.

Random Forests in R

R has the `randomForest` package,
which works essentially the same as `tree`

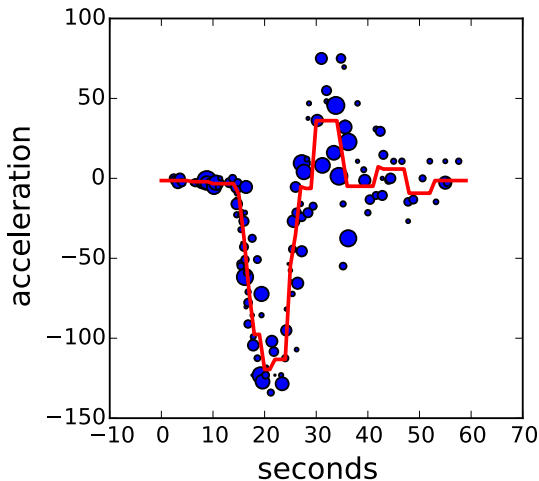
```
spamrf <- randomForest(spam ~ ., data=xemail)
```

For big datasets, use `x = x`, `y = y` like in `gamlr`.

Unfortunately, you lose the interpretability of a single tree.

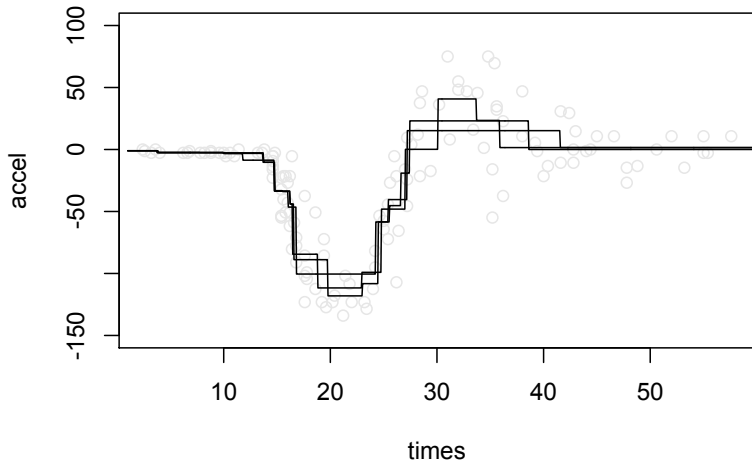
However, if you set `importance=TRUE`, Random Forest will evaluate each \mathcal{T}_b 's performance on the *left-out sample* (recall each tree is fit on a sub-sample). This yields nice OOS stats.

They can be slow (due to many tree fits) but they can also be fit in parallel or on distributed data...



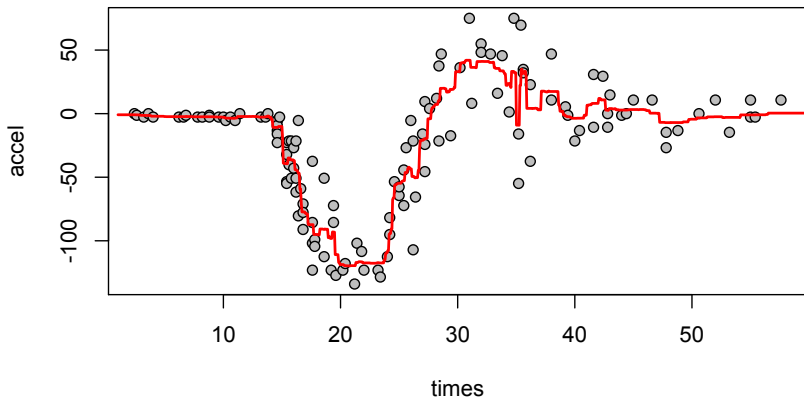
Fitting CART to sampled-with-replacement data is equivalent to randomly **weighting** your observations in the deviance calculations.
(size proportional to weight in this picture).

Random Trees for the Motorcycle Data



If you fit to random subsets of the data,
you get a slightly different tree each time.

Model Averaging with **Random Forests**



Averaging many trees yields a single response surface.

Still looks like a bit of overfit to me, which remains a danger.

A larger example: **California Housing Data**

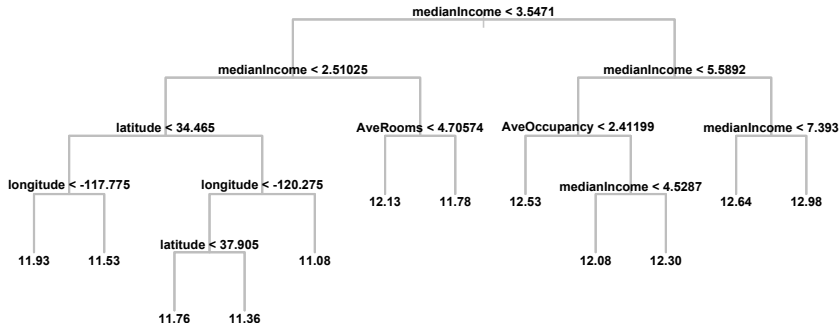
Median home values in census tracts, along with

- ▶ Latitude and Longitude of tract centers.
- ▶ Population totals and median income.
- ▶ Average room/bedroom numbers, home age.

The goal is to predict $\log(\text{MedVal})$ for census tracts.

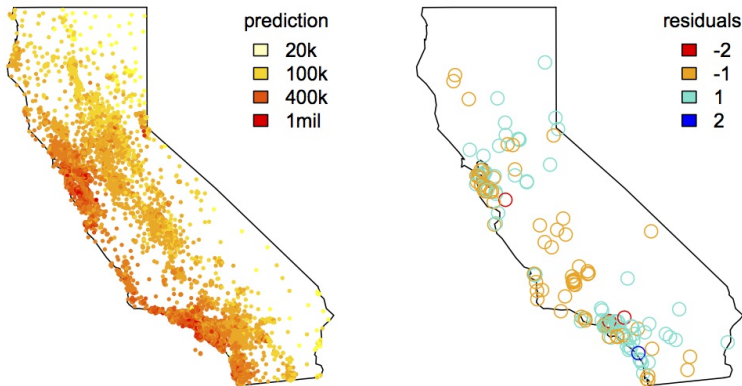
Difficult regression: Covariate effects change with location.
How they change is probably not linear.

CART Dendrogram for CA housing



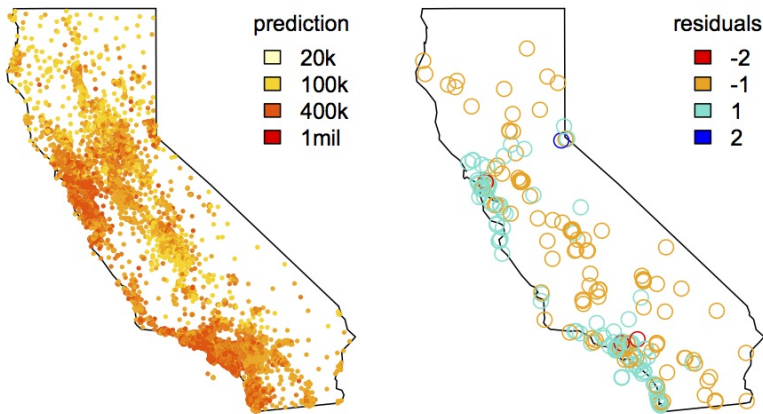
Income is dominant, with location important for low income.
Cross Validation favors the most complicated tree: 12 leaves.

LASSO fit for CA housing data



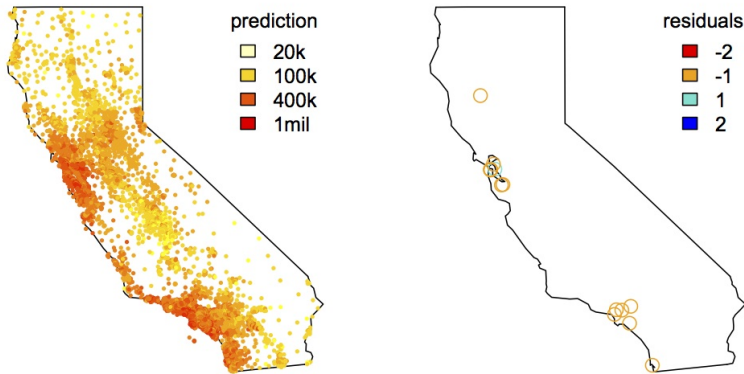
Looks like over-estimates in the Bay, under-estimates in OC.

CART fit for CA housing data



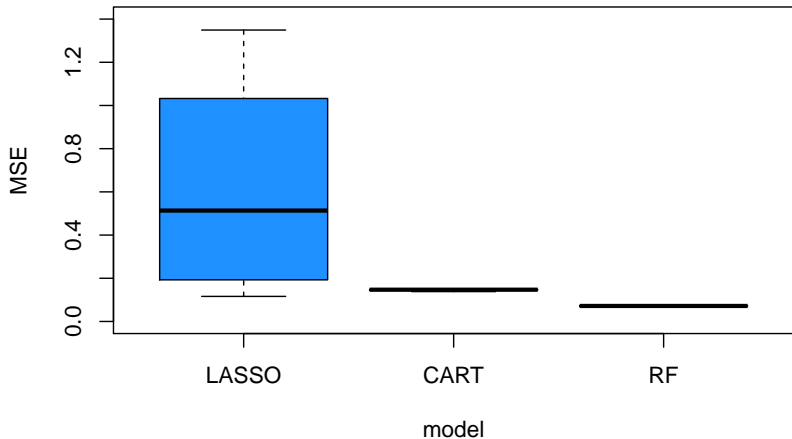
Under-estimating the coast, over-estimating the central valley?

randomForest fit for CA housing data



No big residuals! (although still missing the LA and SF effects)
Overfit? From out-of-sample prediction it appears not.

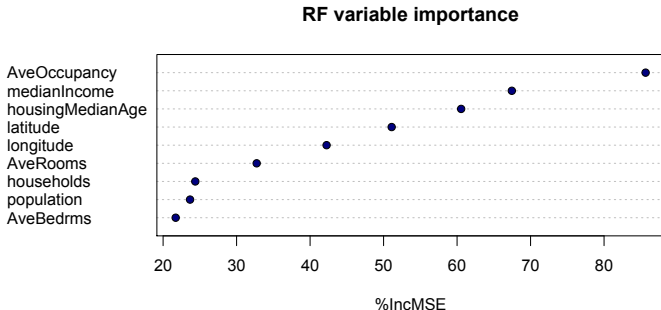
CA housing: out-of-sample prediction



Trees outperform LASSO: gain from nonlinear interaction.
RF is better still than CART: benefits of model averaging.

Although you don't have a nice single tree to interpret, `randomForest` provides **OOS** variable importance plots.

You need to run `randomForest` with `importance=TRUE`. Otherwise it doesn't store the necessary information.



The x-axis here is the % amount that removing splits on that variable would increase the MSE.
For classification it plots increase in % misclassified.

Roundup on **Tree-based learning**

We've seen two techniques for building tree models.

- ▶ **CART**: recursive partitions, pruned back by CV.
- ▶ **randomForest**: average many simple CART trees.

There are many other tree-based algorithms.

- ▶ **Boosted Trees**: repeatedly fit simple trees to residuals.
Fast, but it is tough to avoid over-fit (requires full CV).
- ▶ **Bayes Additive Regression Trees**: mix many simple trees.
Robust prediction, but suffers with non-constant variance.
- ▶ **Dynamic Trees**: grow sequential 'particle' trees
Good online, but fit depends on data ordering

Trees are poor in high dimension, but fitting them to low dimension factors (principle components) is a good option.

Roundup on **Nonlinear Regression** and **Classification**

Many other *nonparametric learning* algorithms

- ▶ Neural Networks (and deep learning):
many recursive logistic regressions.
- ▶ Support Vector Machines:
Project to HD, then classify.
- ▶ Gaussian Processes, splines, wavelets, etc:
Use sums of curvy functions in regression.

Some of these are great, but all take a ton of tuning.

Nothing's better out-of-the-box in low dimension than trees.

But: when the (simpler) linear model fits, it will do better.

This is most often the case in very high dimension.

Homework due next week

Use your project data for this homework!

If you don't have it yet, get it now.

Build and interpret both a single tree and a random forest.

Compare the result to other techniques we've learned.

CART: fit, prune, + plot. Concentrate on interpretation.

RF: plot variable importance and predictive performance.

[+] *Use Random Forests in treatment effects estimation.*

Look back to Lecture 4. Replace lasso estimation for $\hat{d}(\mathbf{x})$ with a random forest. Compare to the original two-lassos algorithm.