

Model Improvements

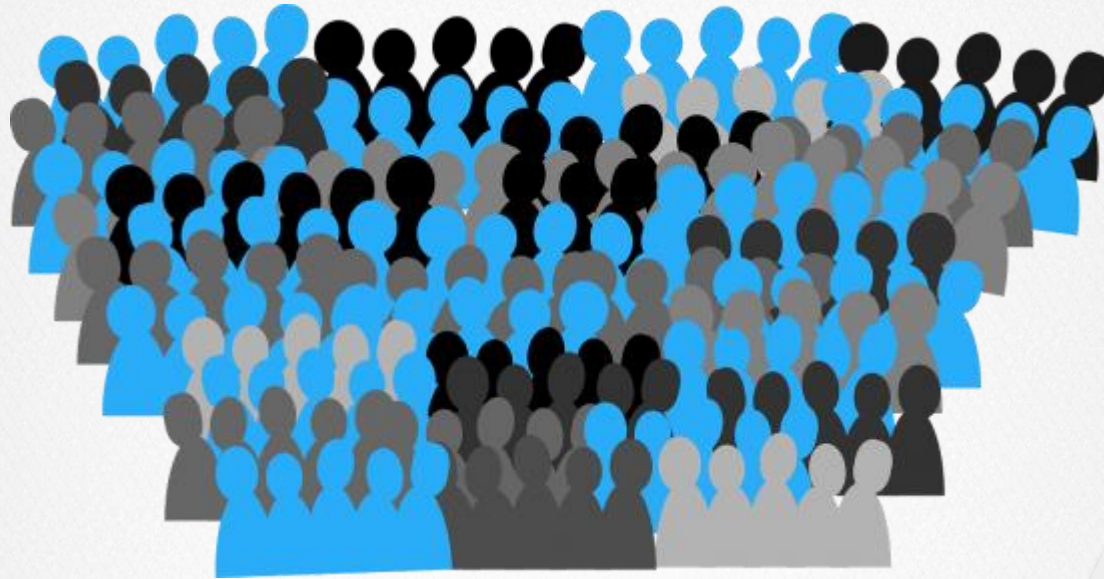
Practical Machine Learning (with R)
UC Berkeley

2 Big Ideas



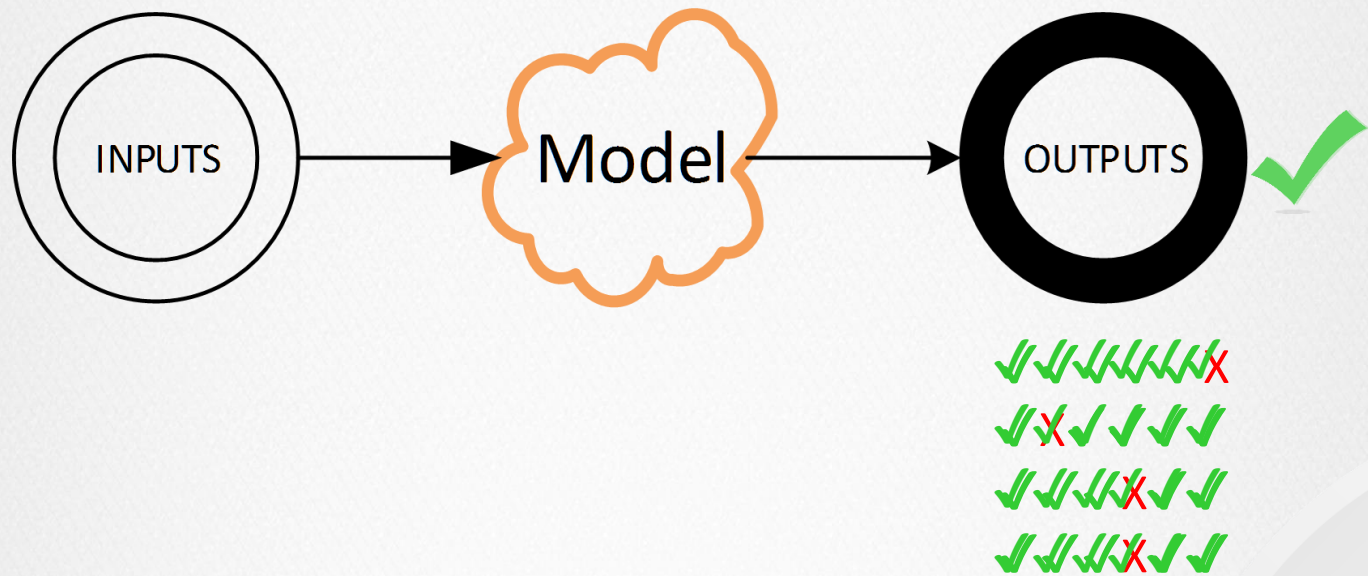
Idea 1



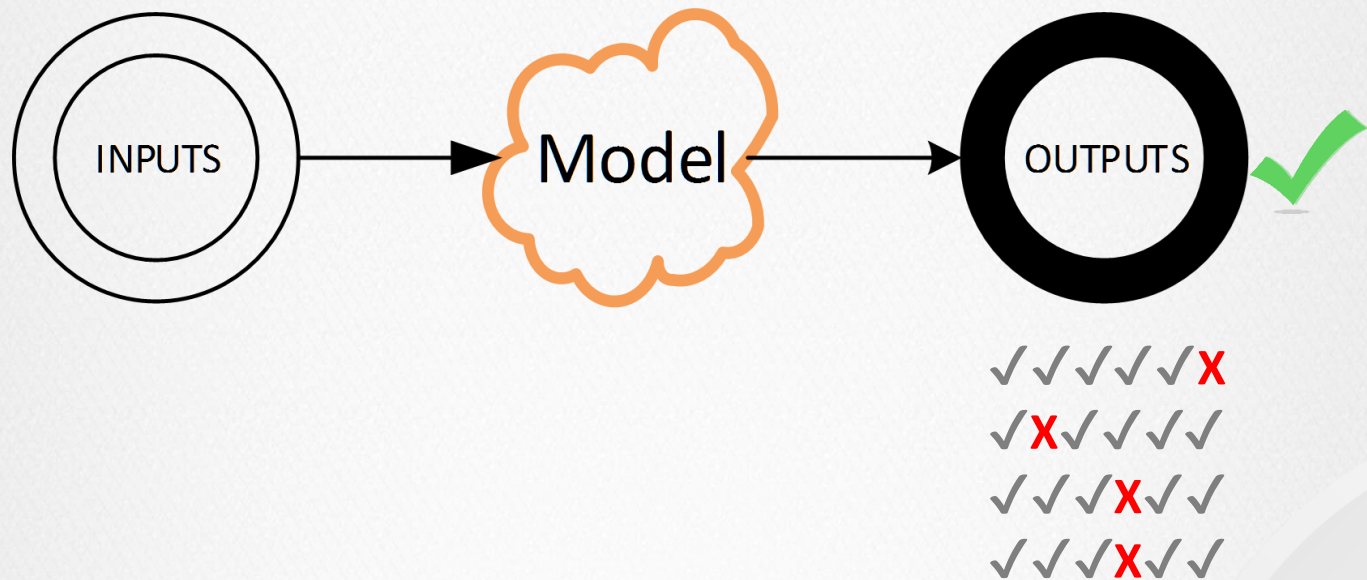


Wisdom of the Crowds

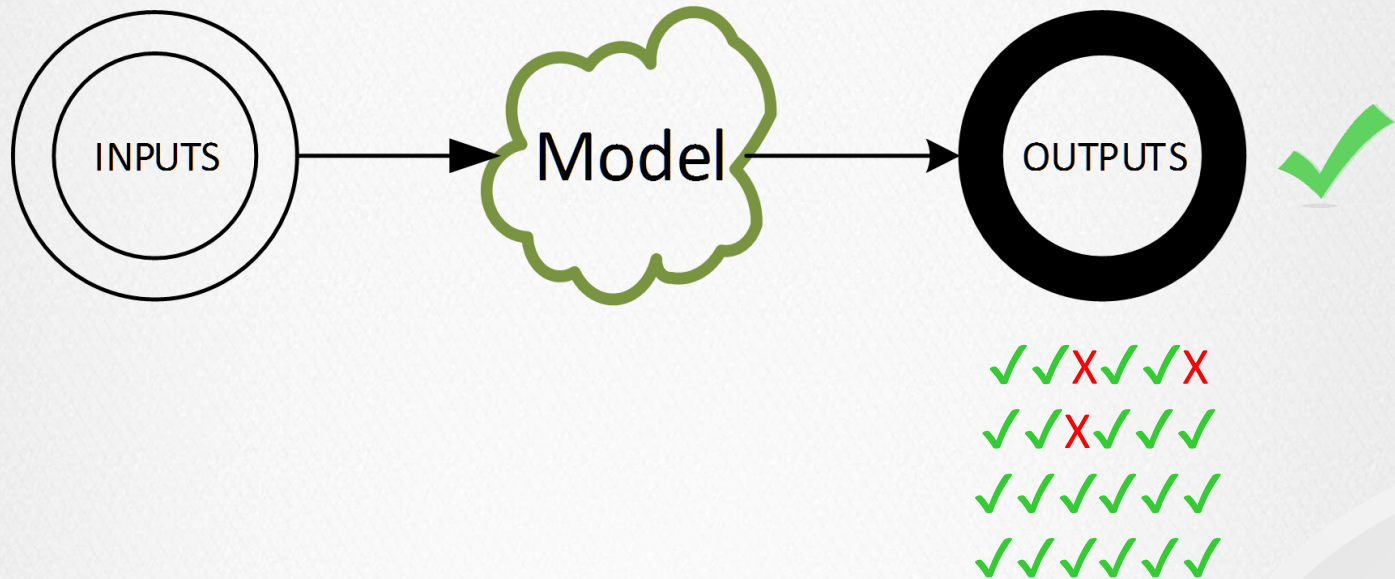




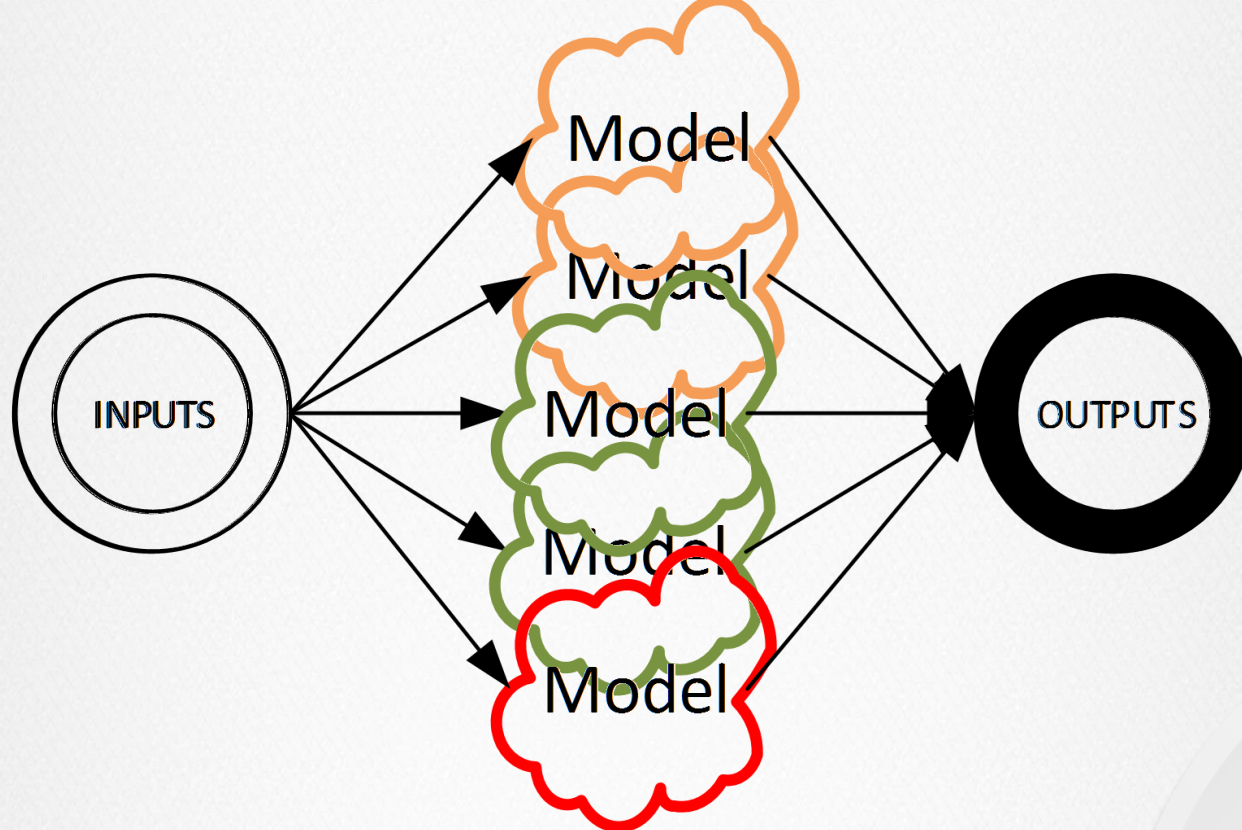
A model fits data well
(generally) ...



but may provide **poor**
estimates for certain
observations.



Another ***model*** may fit these observations well.



So why not use 2 models?
Or three?
Or twenty?

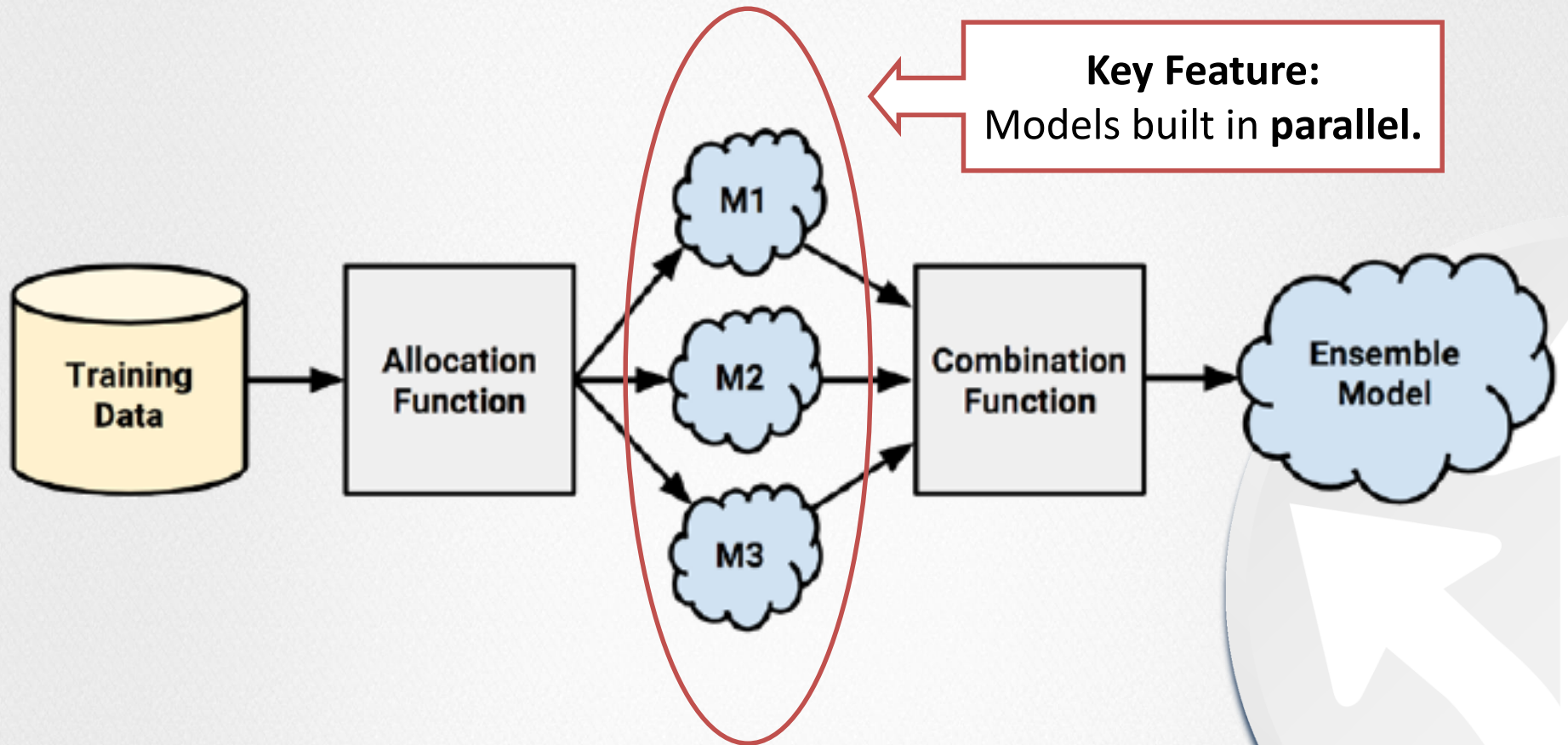
WISDOM OF THE CROWDS

Combine results from multiple models (**ensembles**) to get:

- Better predictions
- Lowers variance for the same model



Ensembles



Source: *Machine Learning with R*

Allocation

⇒ Easiest thing to do?

- **All** data to each model
- **Split** data randomly between models

Are other strategies?



Combination Function

- ➔ Easiest thing to do?
 - **Average** model results

Other strategies?



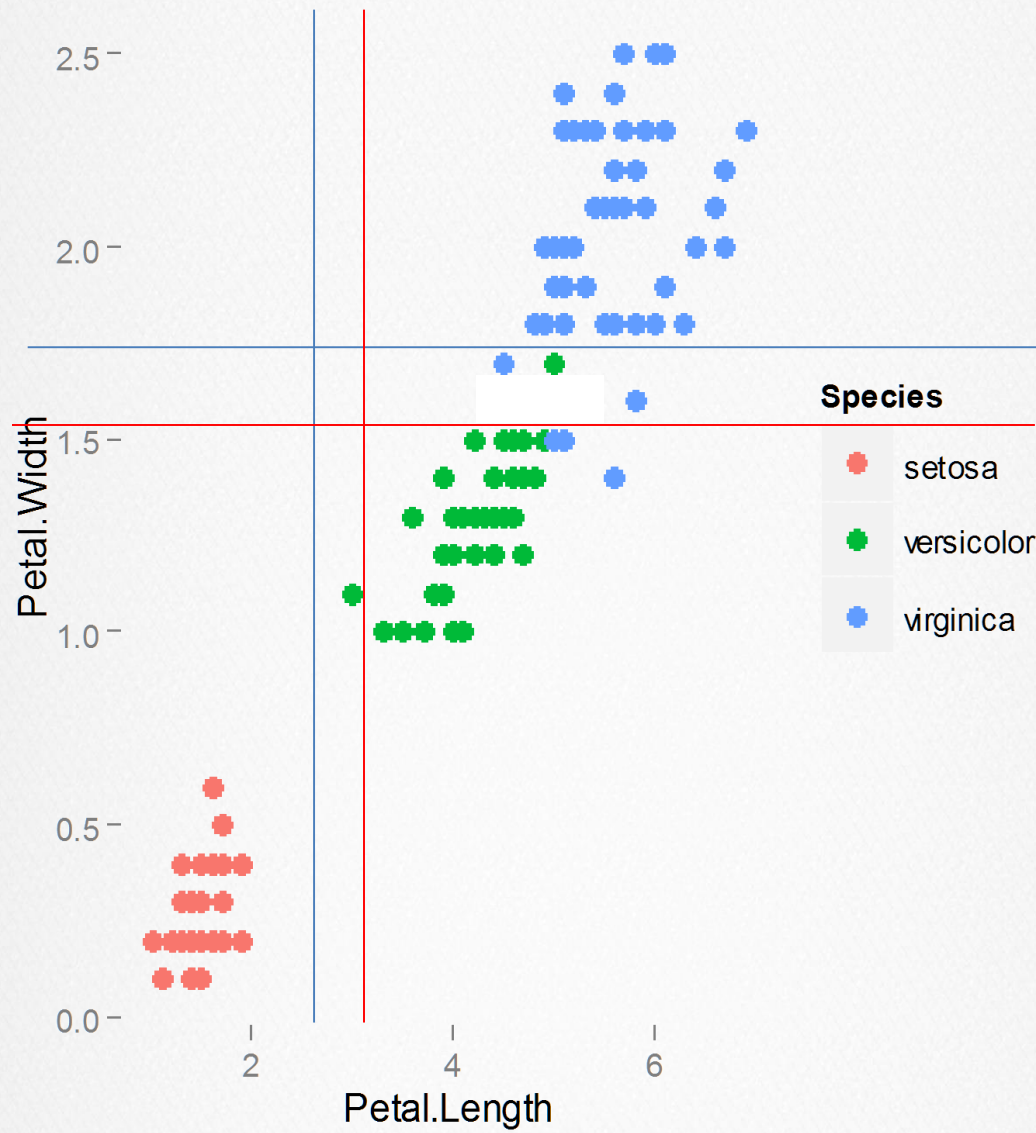
BAGGING MODELS

➔ Brieman:

"Bagging is a general approach that uses bootstrapping in conjunction with any regression (or classification) model to construct an ensemble."

```
1 for  $i = 1$  to  $m$  do
2   |   Generate a bootstrap sample of the original data
3   |   Train an unpruned tree model on this sample
4 end
```

$$\hat{y} = \frac{\sum_i \hat{y}_i}{m}$$



BAGGING NOTES

➔ Advantages

- Lowers variance
- Increases stability
- Has less effect on lower variance models (e.g. linear models)
- More effective with weak learners

➔ Disadvantages

- Computational cost → but parallelizable
- Reduces Interpretability



RANDOM FORESTS



RANDOM FOREST

- ➔ Bagged trees
- ➔ Consider subset of predictors at each split

```
1 Select the number of models to build,  $m$ 
2 for  $i = 1$  to  $m$  do
3     Generate a bootstrap sample of the original data
4     Train a tree model on this sample
5     for each split do
6         Randomly select  $k$  ( $< P$ ) of the original predictors
7         Select the best predictor among the  $k$  predictors and
           partition the data
8     end
9     Use typical tree model stopping criteria to determine when a
       tree is complete (but do not prune)
10 end
```


TUNING PARAMETER

m_{try} : number of predictors to use at each split

- **regression** 1/3rd of number predictors
 - **classification** $\sqrt{\text{number of predictors}}$
- ➔ Kuhn: “Starting with five values of k that are somewhat evenly spaced across the range from 2 to P ”.

ADVANTAGES

- ➔ No overfitting
- ➔ More trees better (limited by computation time/power only)
- ➔ In caret, parameters are considered independently
- ➔ Because each learner is selected independently of all previous learners, Random Forests is robust to a noisy response
- ➔ Computationally efficient -- each tree built on subset of predictors at each split.
- ➔ Use any tree variants as "base learner": CART, ctree, etc



DISADVANTAGES

- ➔ Makes Decision Trees harder to interpret
- ➔ More time/work to tune the model
 - Number of trees
 - m_{try}



Idea 2





VS



***Greed* is
Bad**

***Patience* is
Good**



GREED IS BAD; PATIENCE IS GOOD

Optimizations/search methods generally wish to arrive at a solution **efficiently**:

- ⇒ Take as few steps as possible ...
- ⇒ go as far as possible in each step



Greedy

GREED IS BAD; PATIENCE IS GOOD

Patience is better:

- ⇒ Slowly approach your solution
- ⇒ Reconsidered at each step

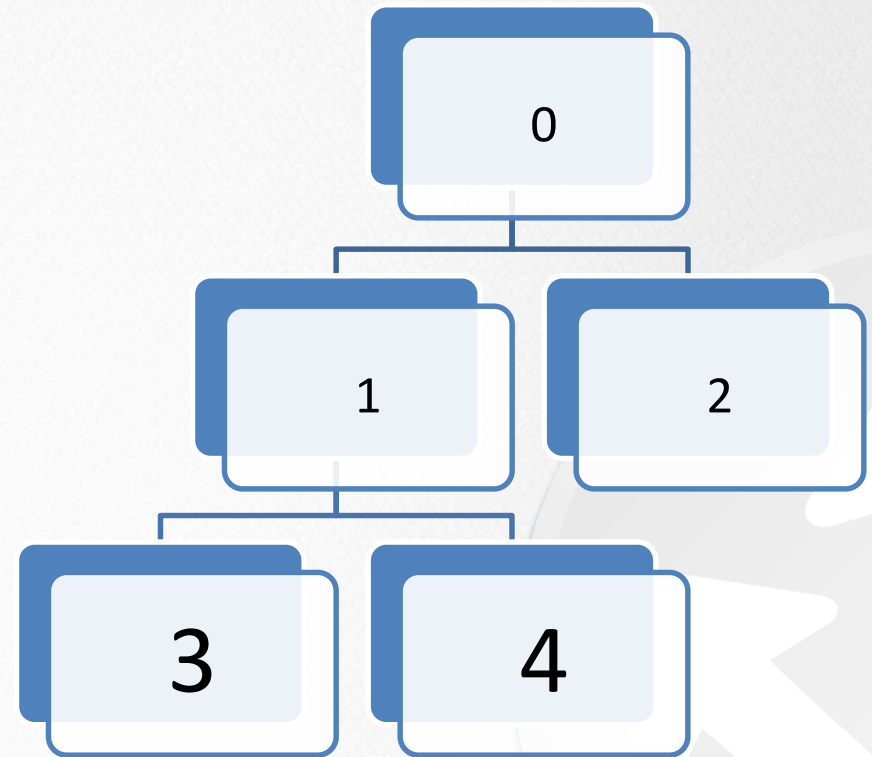


CART

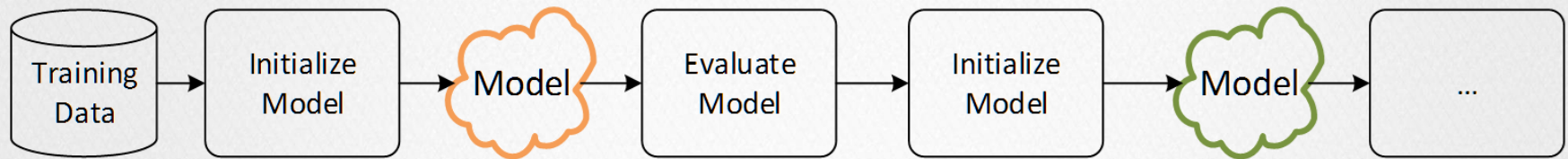
Does CART use patient strategy?

Is a single binary split a model?.

Does the model reconsiders/resets at each iteration?



GREED VS PATIENCE



Process is:

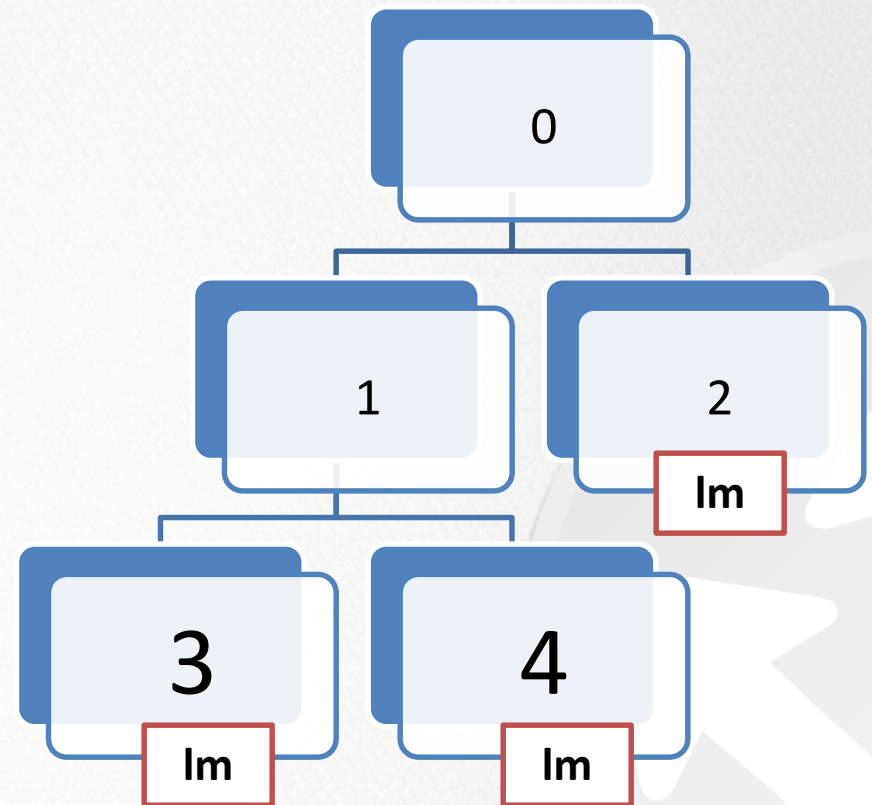
- Sequential
- Recursive

Implication for our algorithm



Tree Enhancement: M5

- Having one value represent the entirety of the node leaves information in the node.
- Function in the node is a simple average
- Use something better
 - **M5** put linear models in nodes of trees

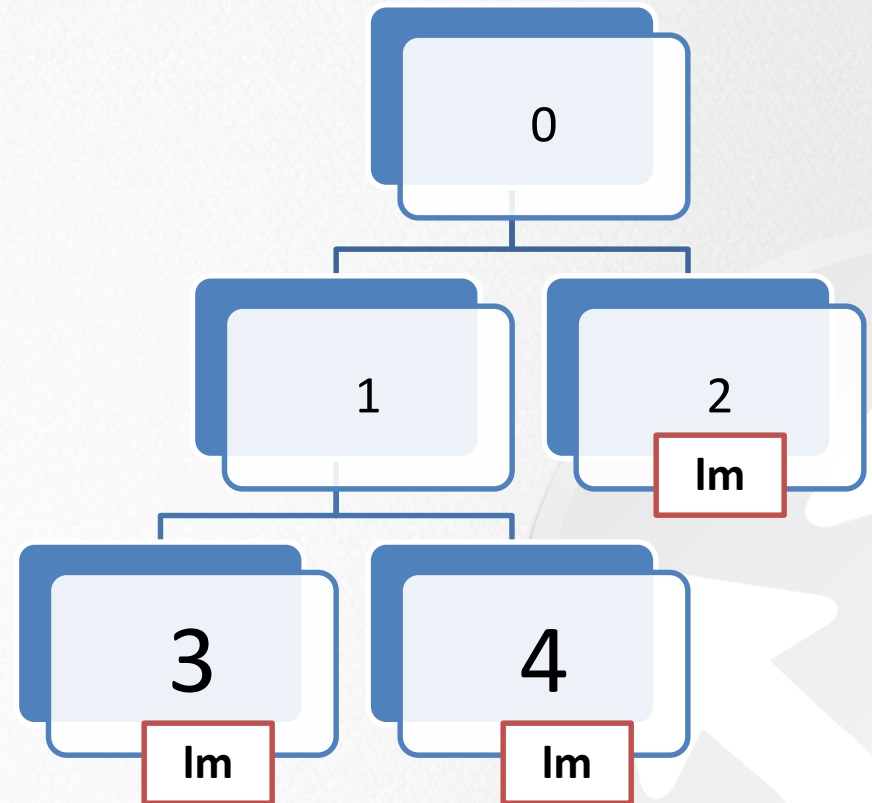
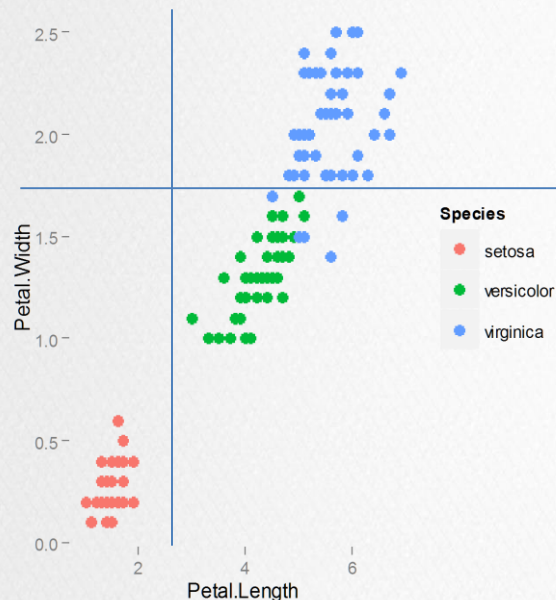


M5 Tree Enhancement (cont.)

→ Greed is bad

- linear models are built on the residuals of the tree model.

- Models are recursive



BOOSTING



BOOSTING

- ⇒ Single models work;
 - Multiple models work better

- ⇒ Idea is simple:

- **Fit first** model: $\hat{y}_1 \sim f_1(x)$

- **Fit** errors/residuals:
$$\begin{aligned}\hat{y}_2 &= f_2(y - \hat{y}_1) \\ &= f_2(y - f_1(x)) \\ &= f_2(x)\end{aligned}$$

- **Iterate:** $\hat{y}_i = (y - \hat{y}_{i-1}) \sim f_i(x)$

- **Predict:** $\hat{y} \sim \sum_i f_i(x)$

Boosting

Fit successive series of models each subsequent model to the prior models residuals.

BOOSTING NOTES

- ➔ Additive models
- ➔ Works best with “weak learners”
 - i.e. ungreedy, low bias, low variance
 - ~~Any~~ Most models with a tuning parameter can be a weak learner
 - Trees are excellent weak learners
 - Weak → “restricted depth”
- ➔ Residuals or errors define a gradient
- ➔ Interpreted as forward step-wise regression with exponential loss



BOOSTING EXAMPLES



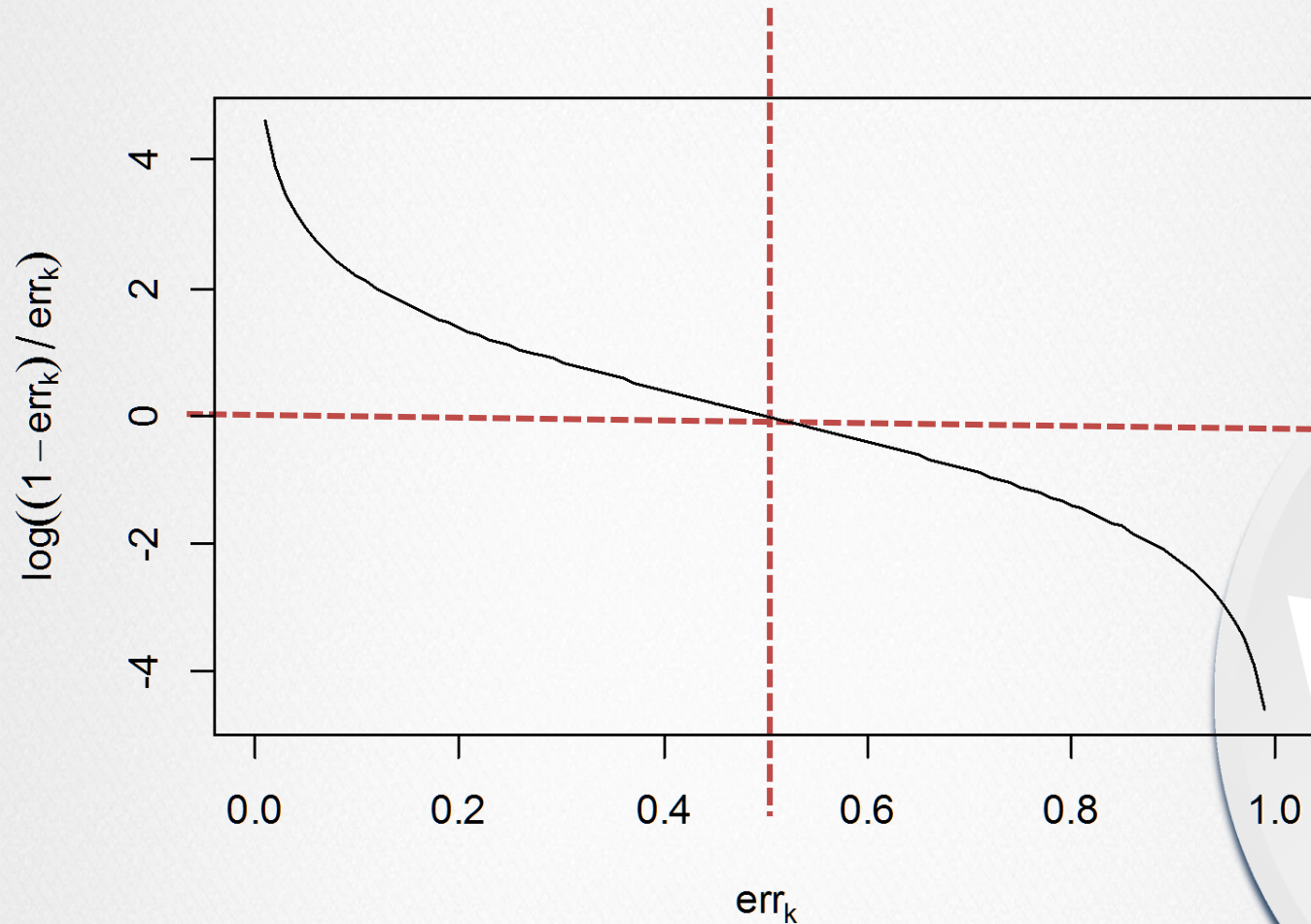
ADABOOST (SHAPIORE/FREUND)

Binary Classification

- 1 Let one class be represented with a value of $+1$ and the other with a value of -1
- 2 Let each sample have the same starting weight ($1/n$)
- 3 **for** $k = 1$ **to** K **do**
 - 4 Fit a weak classifier using the weighted samples and compute the k th model's misclassification error (err_k)
 - 5 Compute the k th stage value as $\ln((1 - err_k) / err_k)$.
 - 6 Update the sample weights giving more weight to incorrectly predicted samples and less weight to correctly predicted samples
- 7 **end**
- 8 Compute the boosted classifier's prediction for each sample by multiplying the k th stage value by the k th model prediction and adding these quantities across k . If this sum is positive, then classify the sample in the $+1$ class, otherwise the -1 class.



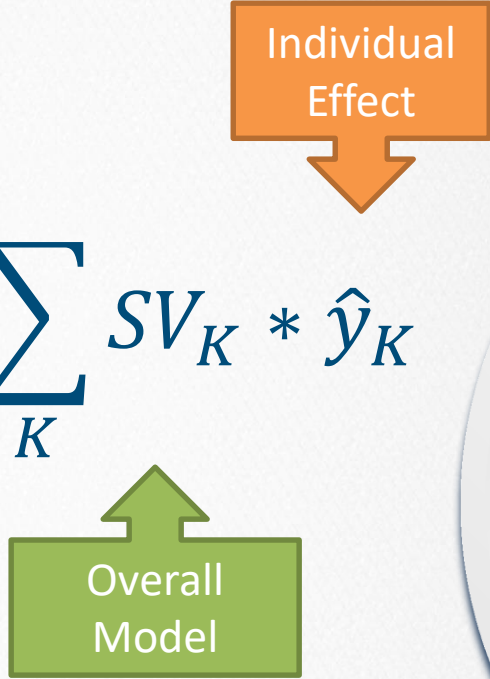
Adaboost Stage Value



<----- Model Quality

ADABOOST MODEL COMBINATION

→ K models

$$PV \mid Class = \sum_K SV_K * \hat{y}_K$$


The diagram illustrates the Adaboost model combination process. It features the equation $PV \mid Class = \sum_K SV_K * \hat{y}_K$ in the center. Above the equation, an orange box labeled "Individual Effect" has a downward-pointing arrow directed at the \hat{y}_K term. Below the equation, a green box labeled "Overall Model" has an upward-pointing arrow directed at the summation symbol \sum_K . To the right of the equation, there is a large, faint, light-gray circular graphic containing a white arrow that curves around itself in a clockwise direction.

SIMPLE GRADIENT BOOSTING (REGRESSION)

- 1 Select tree depth, D , and number of iterations, K .
- 2 Compute the average response, \bar{y} , and use this as the initial predicted value for each sample
- 3 **for** $k = 1$ **to** K **do**
 - 4 Compute the residual, the difference between the observed value and the *current* predicted value, for each sample
 - 5 Fit a regression tree of depth, D , using the residuals as the response
 - 6 Predict each sample using the regression tree fit in the previous step
 - 7 Update the predicted value of each sample by adding the previous iteration's predicted value to the predicted value generated in the previous step
- 8 **end**

Naïve
Guess

SIMPLE GRADIENT BOOSTING SCORING

Models are Additive

$$\hat{y} \sim \sum_i f_i(x)$$

Simple Sum of the Models
Models are Additive



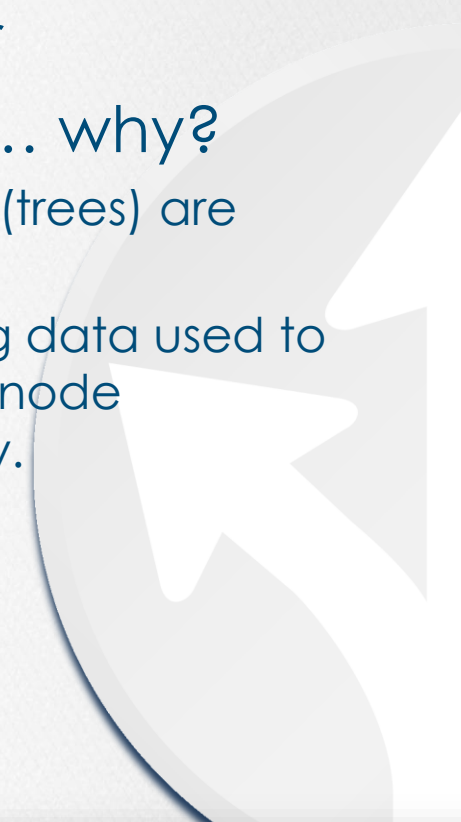
Simple Gradient Boosting

Benefits

- Better predictive performance
- Lower variance

Limitations

- ➔ Less interpretability than base learner
- ➔ Can overfit ... why?
 - base learners (trees) are greedy
 - model training data used to split and eval node simultaneously.



Stochastic Gradient Boosting Machines



- ➔ Proposed by Jerome Freidman (Stanford) circa 1999

M. Kuhn:

“now widely accepted as the boosting algorithm of choice among practitioners”



Simple Gradient Boosting + 2 Innovations

(to fix overfitting problem)



#1

Apply Bootstrap

← Like RF

At each iteration, bootstrap training data with **bagging fraction**, BF .*

Additional Benefits:

lower variance, improve predication accuracy, reduce computation cost

*Typical value of BF is 0.5

#2

Introduce Patience/Limit Greediness

regularization/shrinkage/learning
rate/decay parameter λ



LEARNING RATE, λ

Use Modified Residuals

$$y - \hat{y}$$

“usual” residuals


$$y - \lambda \hat{y}$$

“shrunk” residual

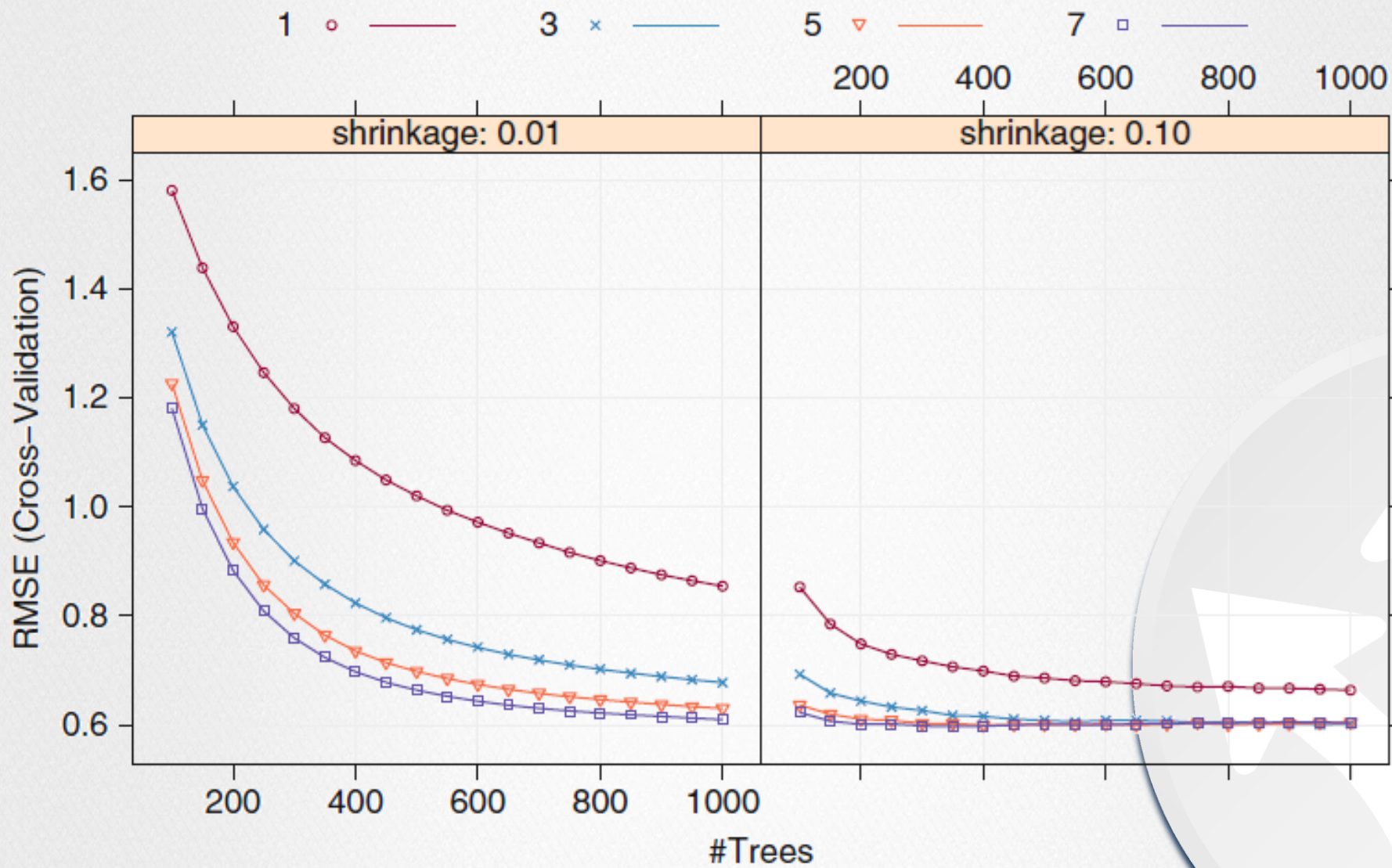
$$0 < \lambda \leq 1$$

For Example

y	\hat{y}	$y - \hat{y}$	$y - 0.1\hat{y}$
10	9	1	9.1
9	10	-1	8



Learning Rate
leaves more
opportunity
for subsequent
models



APM Fig 8.20/

STOCHASTIC GRADIENT BOOSTING SCORING

Models are additive

$$\hat{y} \sim \lambda \sum_i f_i(x)$$

Simple sum of the models, but reduced
by the learning rate



LEARNING RATE, λ

What are the effects of a small (~ 0.01) learning rate vs. large rate?

- more opportunity left at each step
- more steps for similar performance
- increased computational costs
 - computational time, storage size
 $\lambda \sim 1/\text{computational time} \sim 1/\text{storage size}$



STOCHASTIC GRADIENT BOOSTING NOTES

- ⇒ Can be applied to ***any*** learner
- ⇒ Yields lower variance, better performing
- ⇒ Reduced interpretability



Simple Gradient Boosting – Comparison To Random Forest

Similarities

Differences



APPENDIX



Random Forests

Strengths	Weaknesses
<ul style="list-style-type: none">• An all-purpose model that performs well on most problems• Can handle noisy or missing data as well as categorical or continuous features• Selects only the most important features• Can be used on data with an extremely large number of features or examples	<ul style="list-style-type: none">• Unlike a decision tree, the model is not easily interpretable• May require some work to tune the model to the data

GRADIENT BOOSTING (CLASSIFICATION)

Naïve Guess

- 1 Initialized all predictions to the sample log-odds: $f_i^{(0)} = \log \frac{\hat{p}}{1-\hat{p}}$.
- 2 **for** *iteration* $j = 1 \dots M$ **do**
- 3 Cor
- 4 Rar **!**
- 5 Tra This is Algorithm 14.3 from APM. The caption says simple the residuals as the gradient boosting but is the recipe for stochastic gradient boosting.
- 6 Compute the terminal node estimates or the Pearson residuals:
$$r_i = \frac{1/n \sum_i^n (y_i - \hat{p}_i)}{1/n \sum_i^n \hat{p}_i (1 - \hat{p}_i)}$$
- 7 Update the current model using $f_i = f_i + \lambda f_i^{(j)}$
- 8 **end**

STOCHASTIC GRADIENT BOOSTING

- ⇒ Simple Gradient Boosting can overfit (greedy)
 - Apply “regularization/shrinkage”
 - Use λ (“Learning Rate”)
Rather than add the entirety of the residuals, add a fraction of the residuals at each iteration.

$$\hat{y} \sim \lambda \sum_i f_i(x) \quad 0 < \lambda \leq 1$$

- Small values for λ (~ 0.01) work best
 - $\lambda \sim 1/\text{computational time} \sim 1/\text{storage size}$
- ⇒ Use bagging, as well
 - Bagging Fraction: a sample of data in each loop iteration

