

# CS109 – Data Science

Joe Blitzstein, Hanspeter Pfister, Verena Kaynig-Fittkau

[vkaynig@seas.harvard.edu](mailto:vkaynig@seas.harvard.edu)

staff@cs109.org

# Announcements

- HW2 is due today!
- Please execute your notebooks, but without test output.
- Help with lecture material

# Books

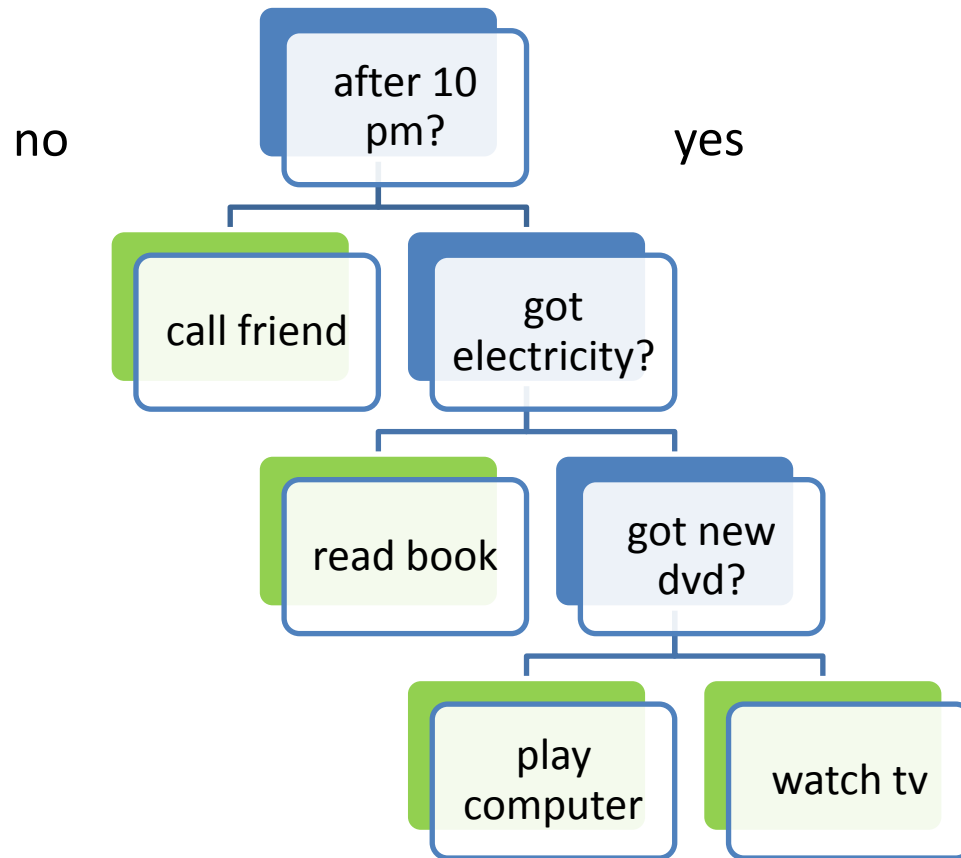
- “Elements of Statistical Learning”
- <http://statweb.stanford.edu/~tibs/ElemStatLearn/>
- “Pattern Recognition and Machine Learning”
- <http://research.microsoft.com/en-us/um/people/cmbishop/PRML/>

# Next Topics

- Tree classifier
- Bagging
- Random Forest



# Decision Tree



# Decision Trees

SVM weight vector harder to interpret  
than decision tree

- Fast training
- Fast prediction      Going down splits is reasonably fast.
- Easy to understand
- Easy to interpret

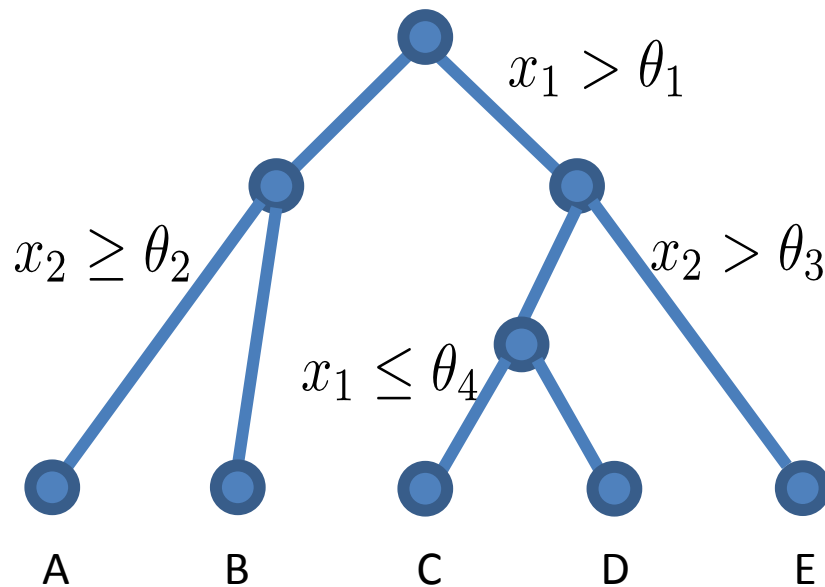
<http://en.akinator.com/personnages/jeu>

data vector  $x$ , with two features:  $x_1, x_2$

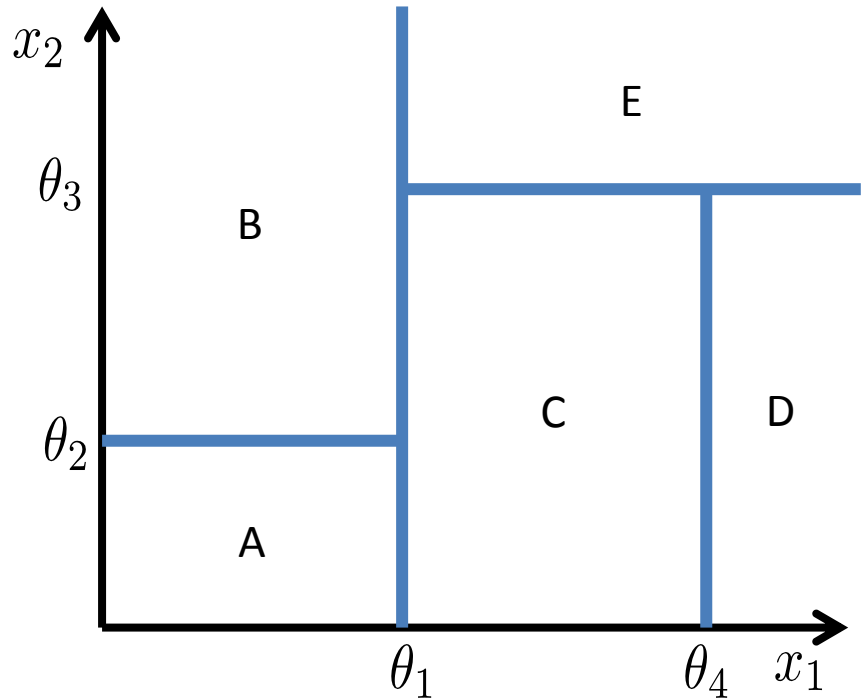
only query  $x_1$  with threshold  $\theta_1 \rightarrow$  split data into two sets

# Decision Tree - Idea

second split with  $\theta_2$  for  $x_2$



Finally, assign classes into tree:  
 $\rightarrow$  Use chart w/ cells.



# Decision Tree - Idea

- What is the benefit of using only one feature at a time?

Scaling/normalization not important:  
do one feature at a time -> tree don't care.

Performance: super fast

Very intuitive.

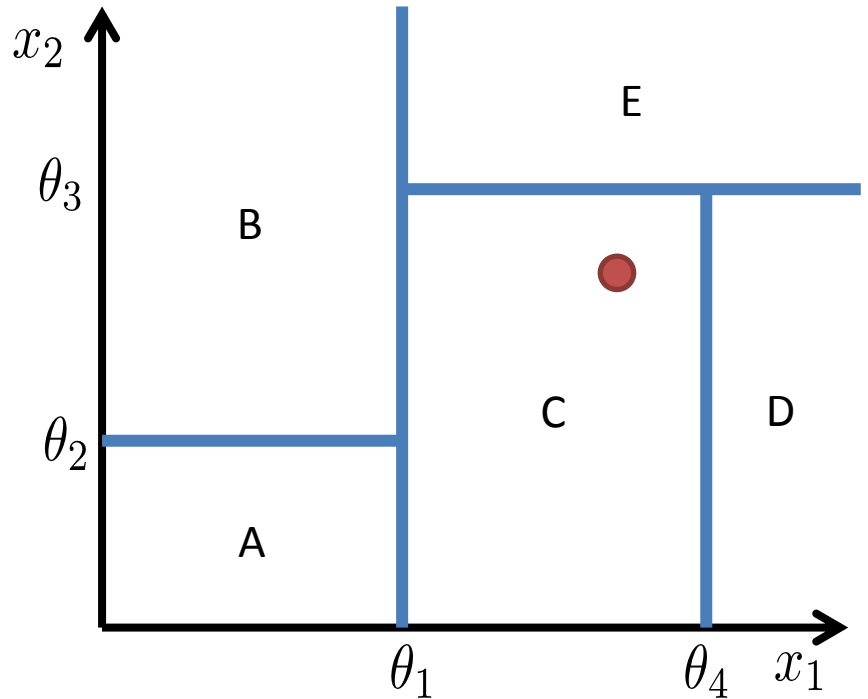
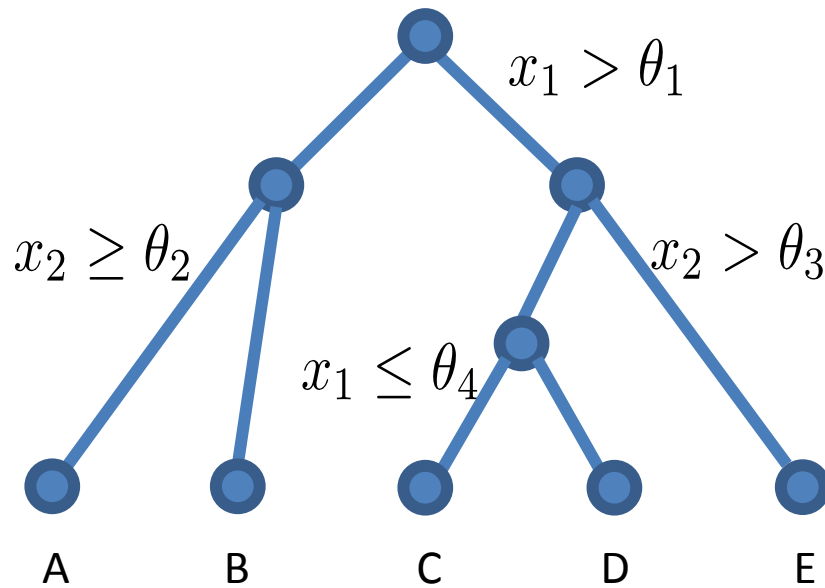
Categories: can work with nonnumerical data (without conversion).

- What is the drawback?

Compared to SVN or KNN boundary, everything is straight (not diagonal)  
ie, we chose orthogonal relative to feature axis: would need more splits to approximate diagonals (tree gets deeper).



# Decision Tree - Prediction

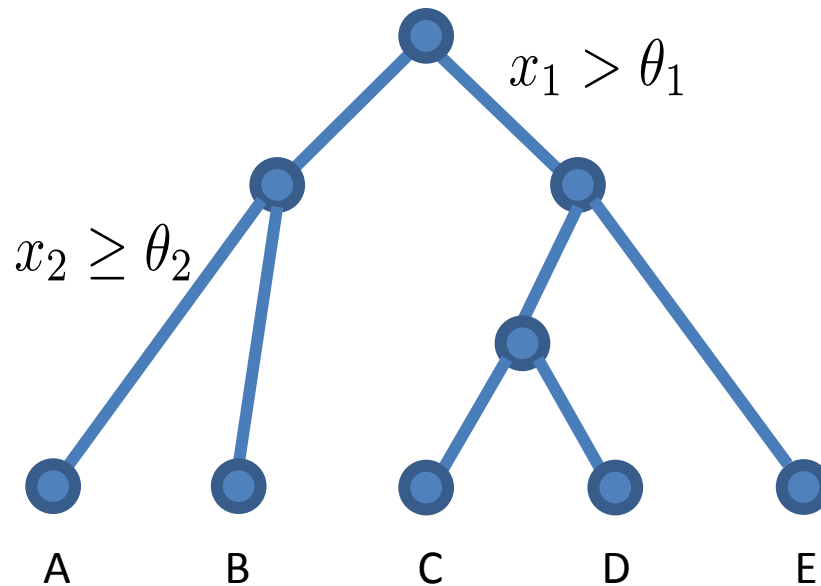


Add new point to cells: can see label/class and not the category

Decision tree: natural multiclass, just assign diff classes to leaves (the labeling category)

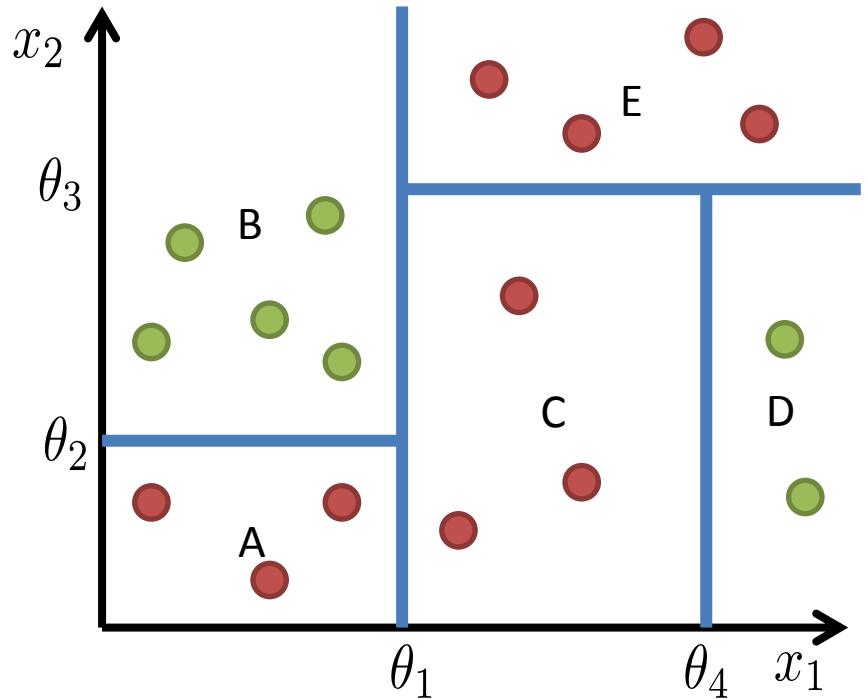
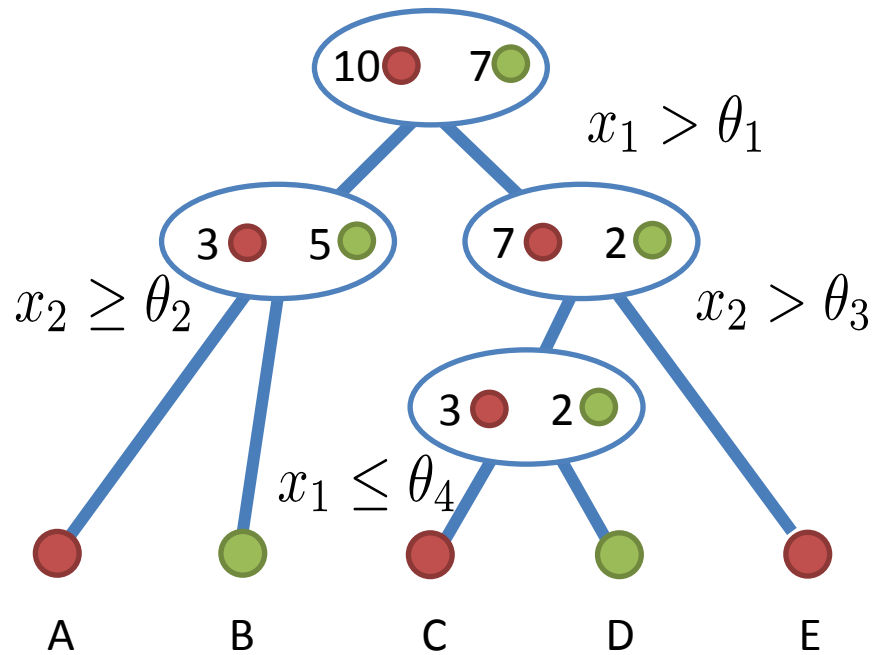
# Decision Tree -Training

- Learn the tree structure:
  - which feature to query
  - which threshold to choose



Criteria to decide threshold:

# Node Purity



Come up with splits that give cells that are pure and only have one cell in them

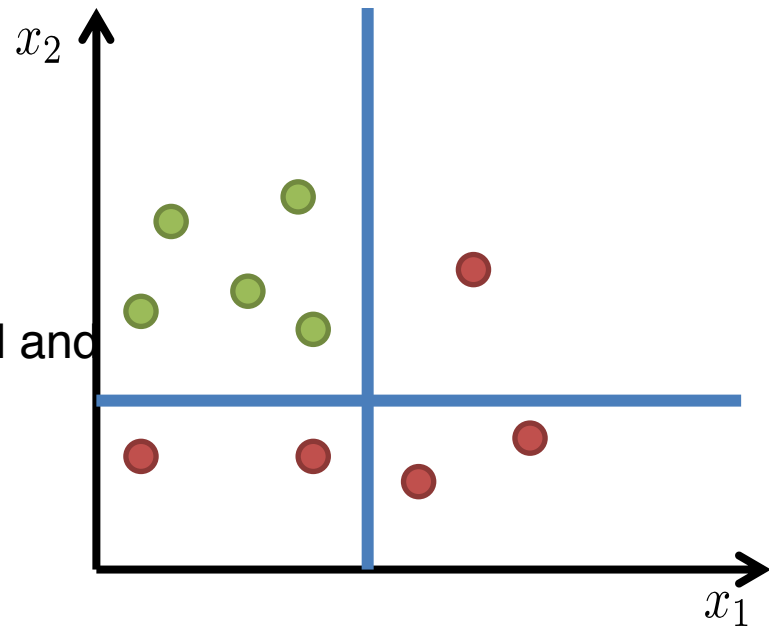
# Gini Impurity

- Expected error
- if you randomly choose a sample
- and predict the class of the entire node based on it.

begin: no split 50/50 chance

split vertical: we have one pure and  
a node that is 5/7 green

but, if do horizontal split: one pure red and  
other is 5/6 green (less impurity)



# Gini Impurity

Example:

4 **red**, 3 **green**, 3 **blue** data points

- Class probabilities:
  - **red**:  $4/10$       **green**:  $3/10$       **blue**:  $3/10$
- misclassification:
  - **red**:  $4/10 * (3/10 + 3/10)$

Picking  
red

Making an  
error

# Gini Impurity

- misclassification:

— red:

$$4/10 * (3/10 + 3/10) = 0.24$$

— green and blue:

$$3/10 * (4/10 + 3/10) = 0.21$$

- gini impurity: **0.24** + **0.21** + **0.21** = 0.66

sum of probs of mistakes of sorting into  
classes over all classes.

# Gini Impurity

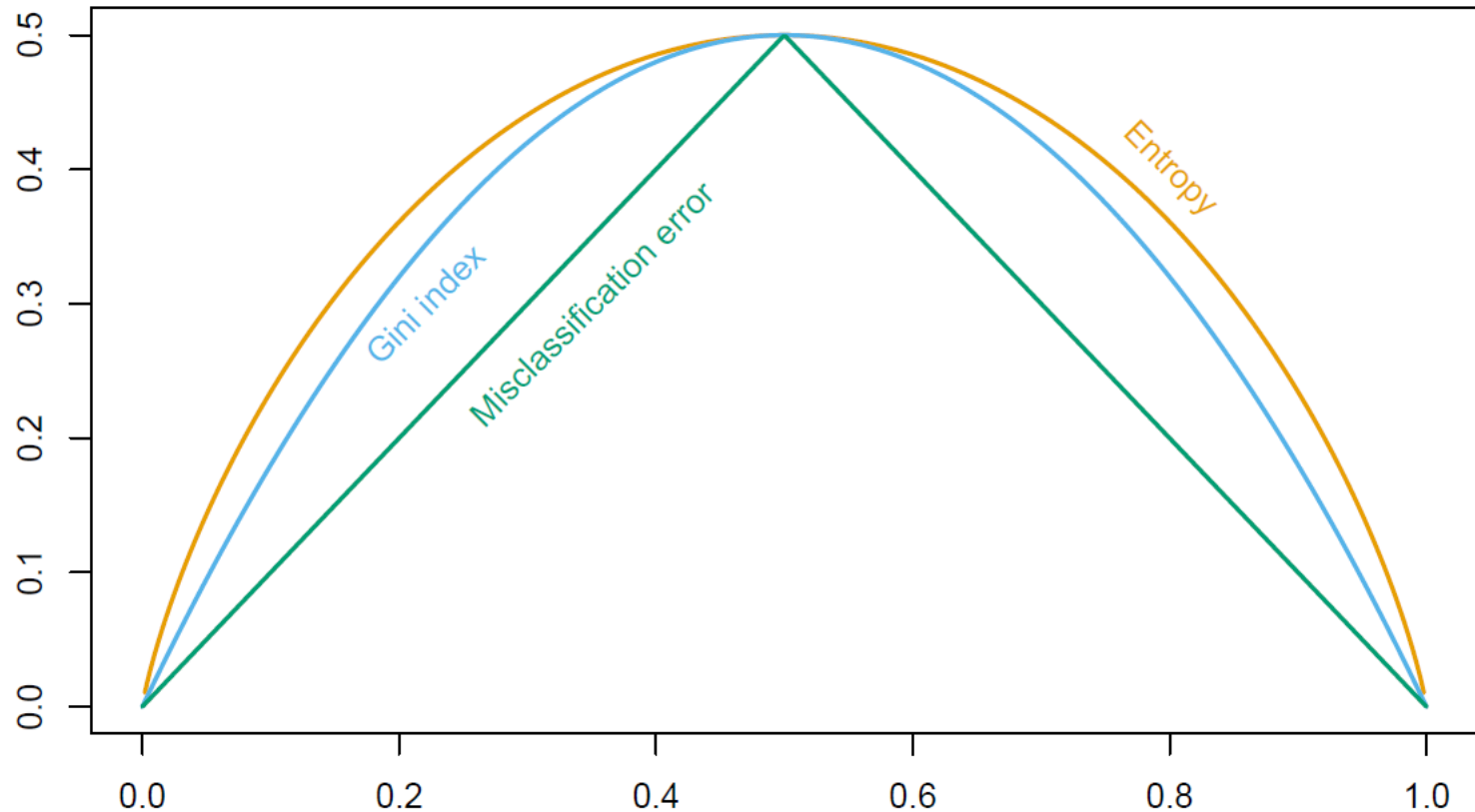
- Number of classes:  $C$
- Number of data points:  $N$
- Number of data points of class  $i$ :  $N_i$

$$I_G = \sum_{i=1}^C \frac{N_i}{N} \left( 1 - \frac{N_i}{N} \right)$$

true  
class

wrong  
prediction

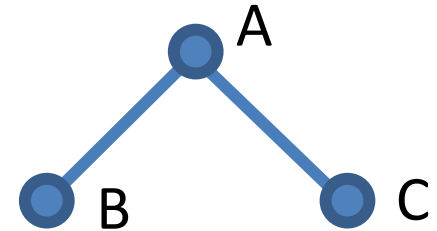
# Gini Impurity





# Node Purity Gain

- Compare:
  - Gini impurity of parent node
  - Gini impurity of child nodes



$$\Delta I_G = I_G(A) - \frac{N(B)}{N(A)} I_G(B) - \frac{N(C)}{N(A)} I_G(C)$$

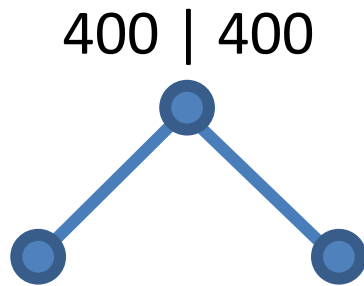
want a lot of pts into a pure node.

# Misclassification

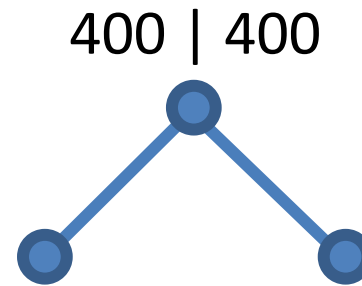
- $\frac{1}{N} \sum_i^N \mathbf{1}(\hat{y}_i \neq y_i)$
- not differentiable

# Comparison Gini vs Misclassification

- Binary problem: 400 samples per class



100 | 300      300 | 100  
clear majority in both, neither pure



200 | 400      200 | 0  
one pure node, but worse other node.

same misclassification: same number of mistakes (doesn't account for distribution)

Misclassification: 0.25

Gini gain: 0.125

Misclassification: 0.25

Gini gain: 0.166

Gini gain tells us more about distribution into higher quality nodes.

# Pseudocode

- Check if already finished
- For each feature  $x_i$ 
  - Calculate the gain from splitting on  $x_i$
  - Let  $x_{best}$  be the feature with highest gain
- Create a decision *node* that splits on  $x_{best}$
- Repeat on the sub-nodes

Depth,

- Does this produce an optimal tree?
- What does optimal tree mean?

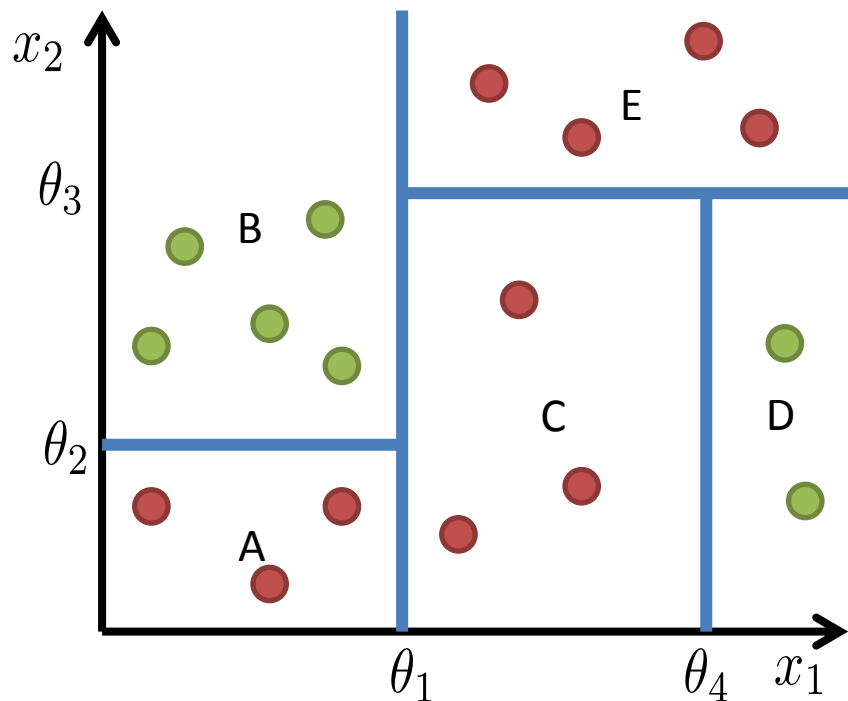
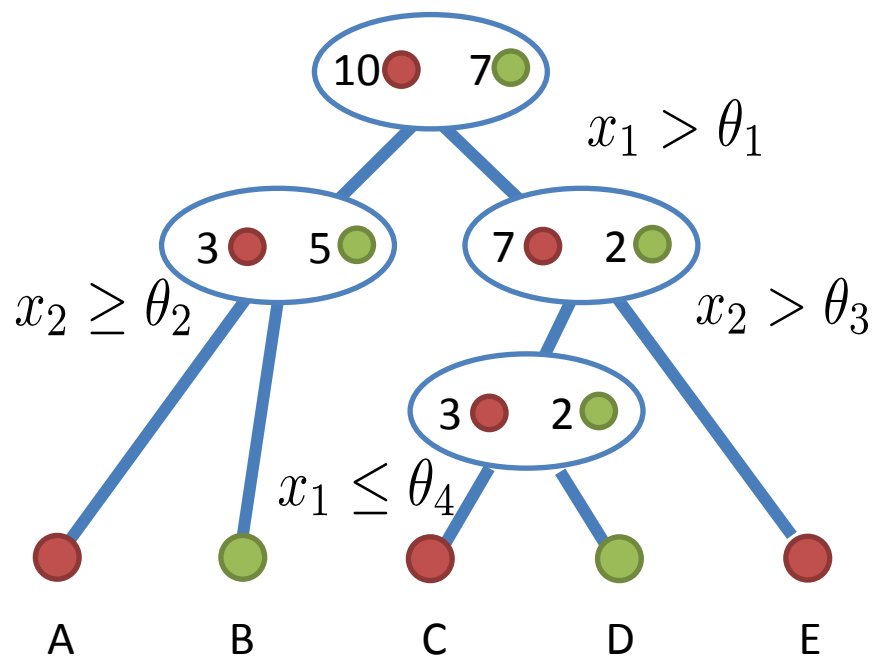
# When to Stop

- node contains only one class
- node contains less than x data points prevent small cells:  
don't want one cell  
for each data pt.
- max depth is reached fast tree.
- node purity is sufficient don't need completely pure:  
base on degree of index.
- you start to overfit => cross-validation

majority vote by leaf.

or you can do a probability ratio to predict based on purity of parent.

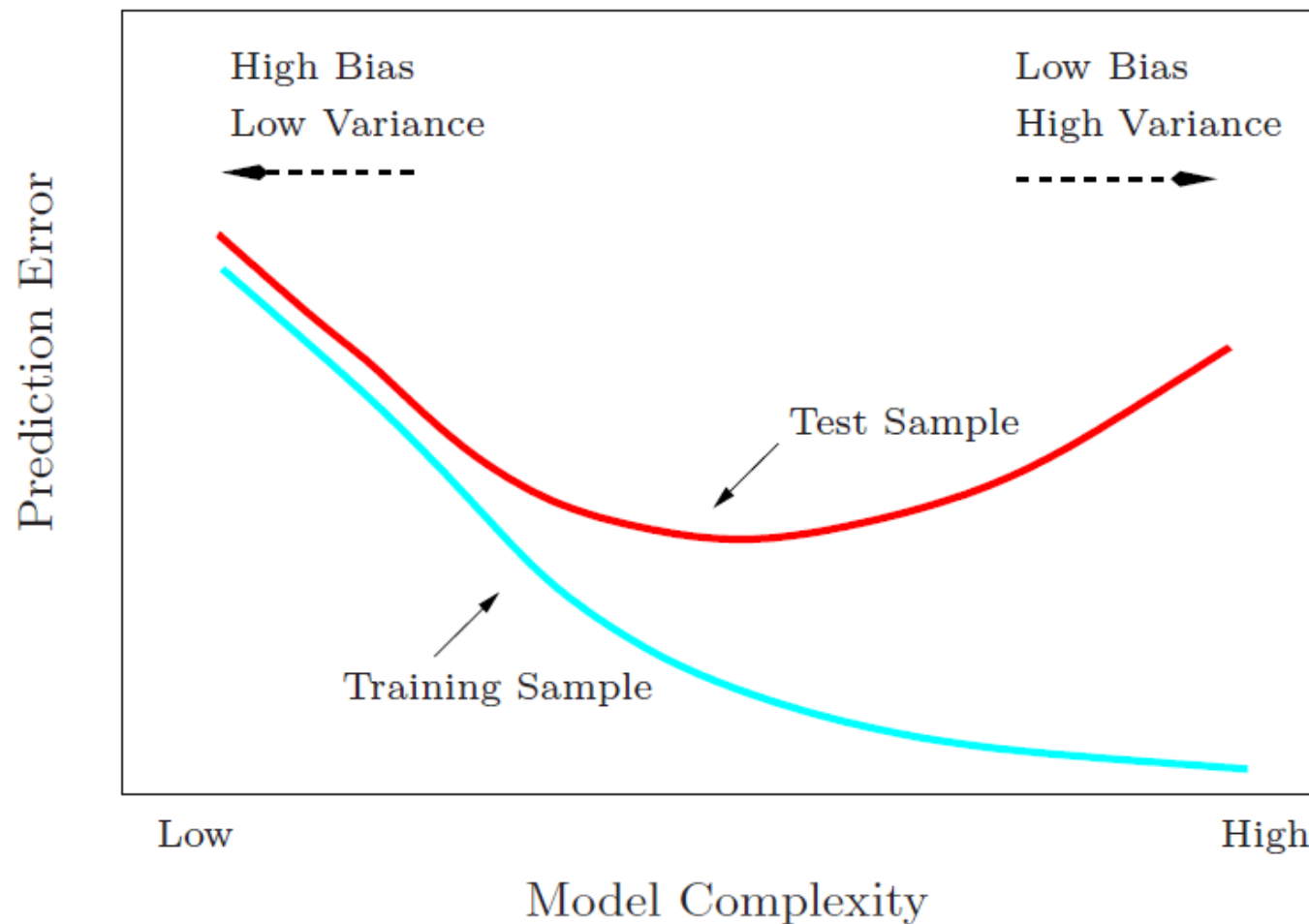
# Tree Pruning



How do you make a prediction for the merged cell?

Tree prone to overfitting without pruning.

# Pruning and Complexity



# Decision Trees - Disadvantages

- Sensitive to small changes in the data  
high variance: sensitive to changes.
- Overfitting
- Only axis aligned splits



# Decision Trees vs SVM

Characteristic	SVM	Trees
Natural handling of data of “mixed” type	▼	▲
Handling of missing values	▼	▲
Robustness to outliers in input space	▼	▲
Insensitive to monotone transformations of inputs	▼	▲
Computational scalability (large $N$ )	▼	▲
Ability to deal with irrelevant inputs	▼	▲
Ability to extract linear combinations of features	▲	▼
Interpretability	▼	◆
Predictive power	▲	▼

rescaling impt for SVM, not trees.

tree scales up well with size of data:  
very simple decisions made.

Diagonal in feature space:  
easy w/ SM by linear combination.

Tree doesn't generalize well.

Use many trees, aggregate to better result.

# Wisdom of Crowds

The collective knowledge of a **diverse and independent** body of people typically exceeds the knowledge of any single individual, and can be harnessed by voting.

James Surowiecki





**Netflix Prize**

# Netflix Prize

- Take home messages:

Early gains very easy, then hit diminishing returns

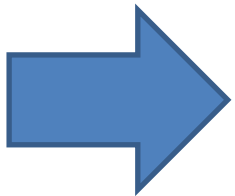
Plateaued, got gains with an ensemble technique

Napoleon Dynamite: had one big error that hurt the % gain.

if at descretion (not in a competition) may want to change metric to avoid outlier issue

# Ensemble Methods

- A single decision tree does not perform well
- But, it is super fast
- What if we learn multiple trees?



We need to make sure they  
do not all just learn the  
same. (colinearity)

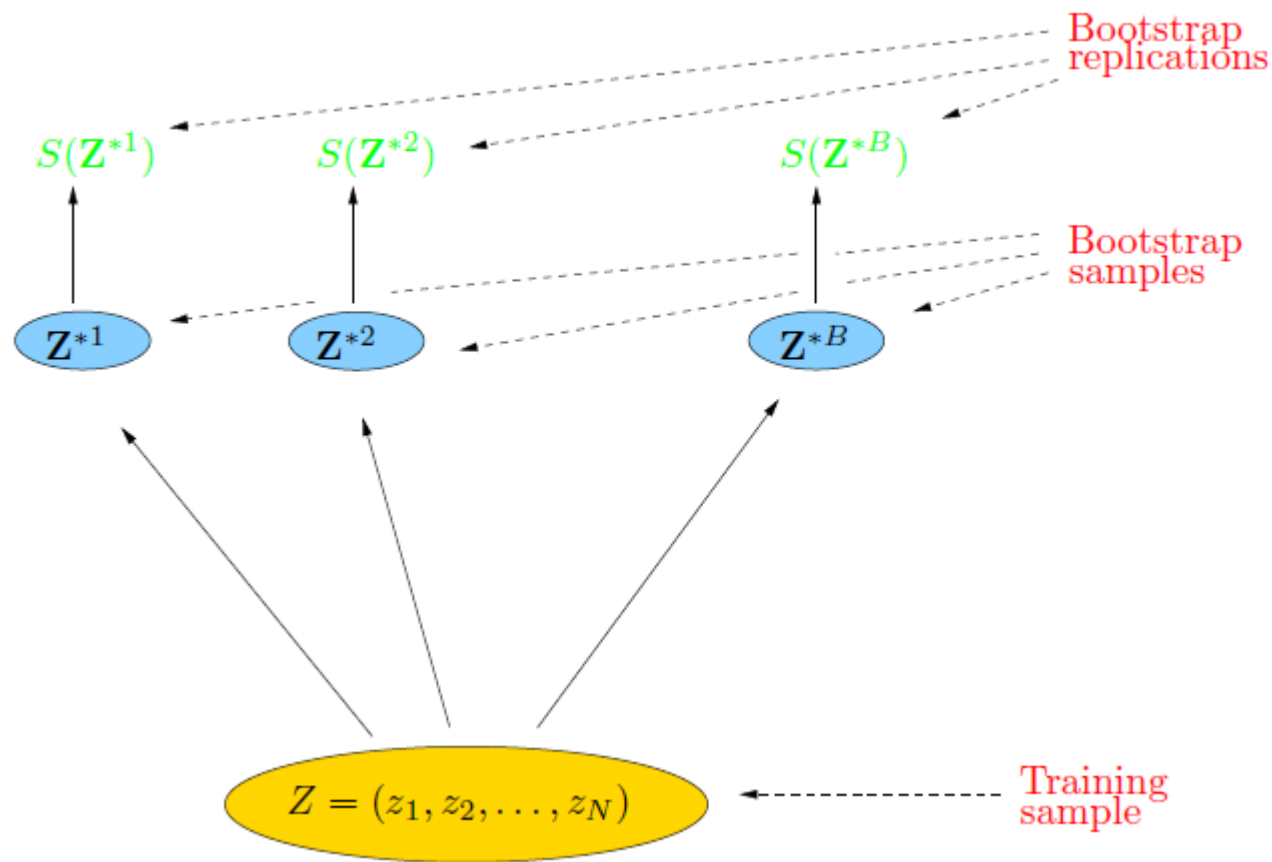
# Bootstrap



# Bootstrap

- Resampling method from statistics
- Useful to get error bars on estimates
- Take  $N$  data points
- Draw  $N$  times with replacement
- Get estimate from each bootstrapped sample

# Bootstrap



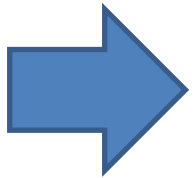


# Bootstrap

- I can generate more data!
- Can I do cross validation on this?  
No? —> same data points, the overlap is too strong.

# Bootstrap vs Cross-validation

- Bootstrap has overlap in data sets



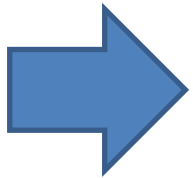
Do not use simple bootstrap to generate train and test data from same data set.

$$p(n \in Z^{*i}) = \frac{1}{N}$$

Probability of choosing n

# Bootstrap vs Cross-validation

- Bootstrap has overlap in data sets



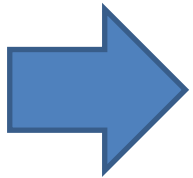
Do not use simple bootstrap to generate train and test data from same data set.

$$p(n \in Z^{*i}) = 1 - \frac{1}{N}$$

Probability of not choosing n

# Bootstrap vs Cross-validation

- Bootstrap has overlap in data sets



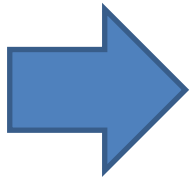
Do not use simple bootstrap to generate train and test data from same data set.

$$p(n \in Z^{*i}) = \left(1 - \frac{1}{N}\right)^N$$

Probability of not choosing  $n$  in  $N$  draws

# Bootstrap vs Cross-validation

- Bootstrap has overlap in data sets



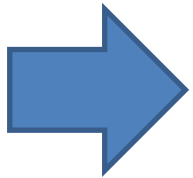
Do not use simple bootstrap to generate train and test data from same data set.

$$p(n \in Z^{*i}) = 1 - \left(1 - \frac{1}{N}\right)^N$$

Probability of (not not) choosing  $n$  in  $N$  draws

# Bootstrap vs Cross-validation

- Bootstrap has overlap in data sets



Do not use simple bootstrap to generate train and test data from same data set.

$$p(n \in Z^{*i}) = 1 - e^{-1}$$

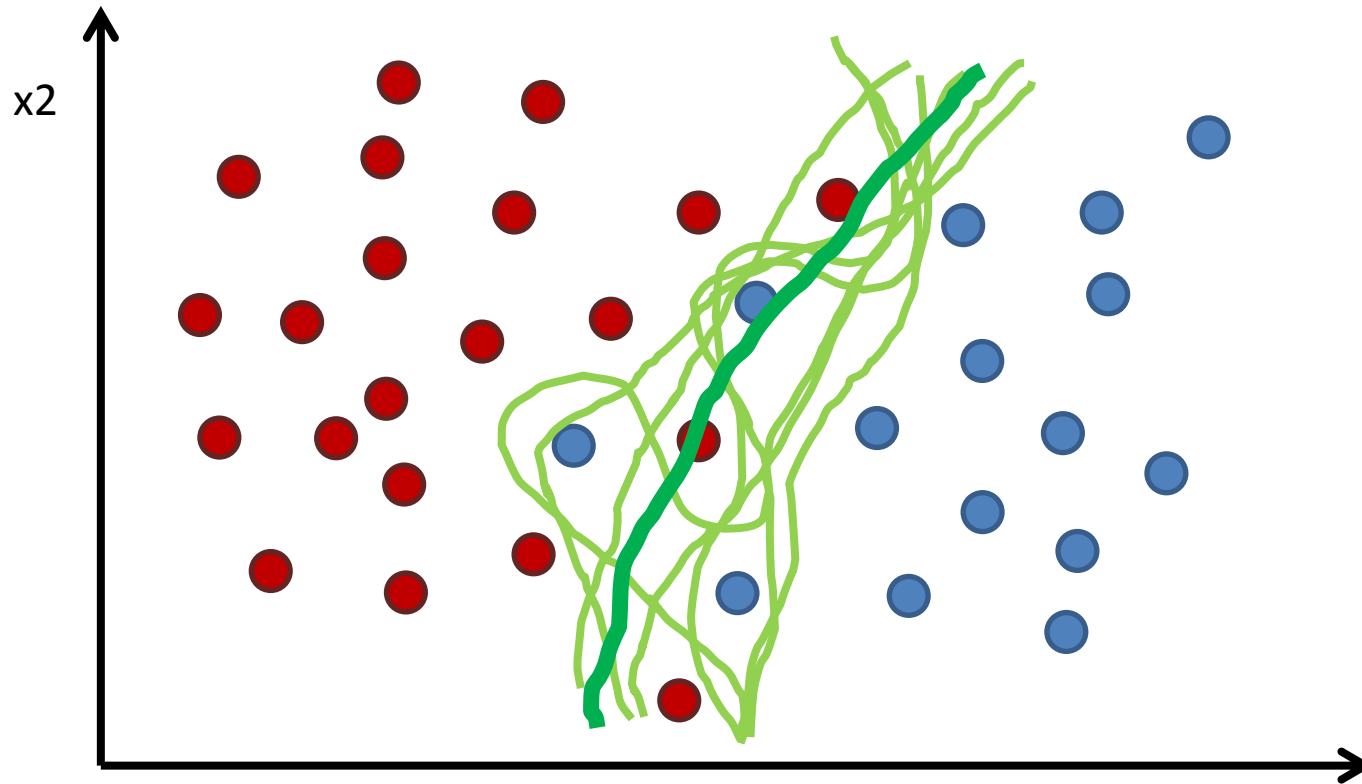
$$\approx 0.632$$

This number is important later

# Bagging

- Bootstrap aggregating
- Sample with replacement from your data set
- Learn a classifier for each bootstrap sample
- Average the results

# Bagging Example

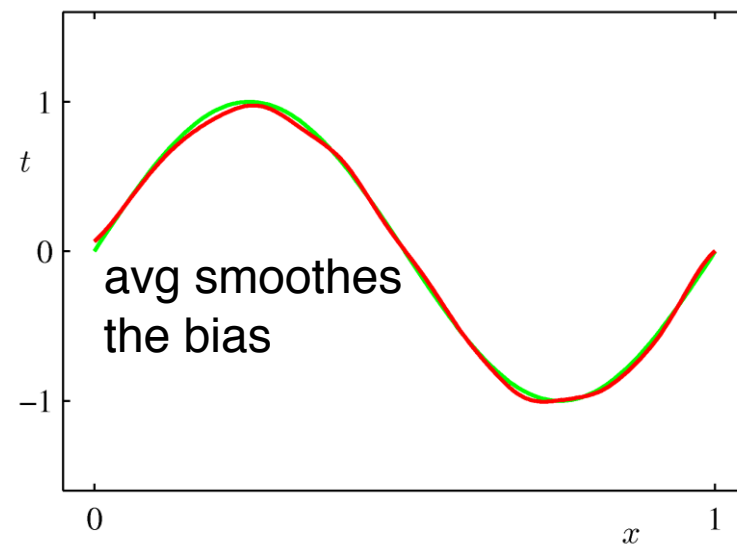
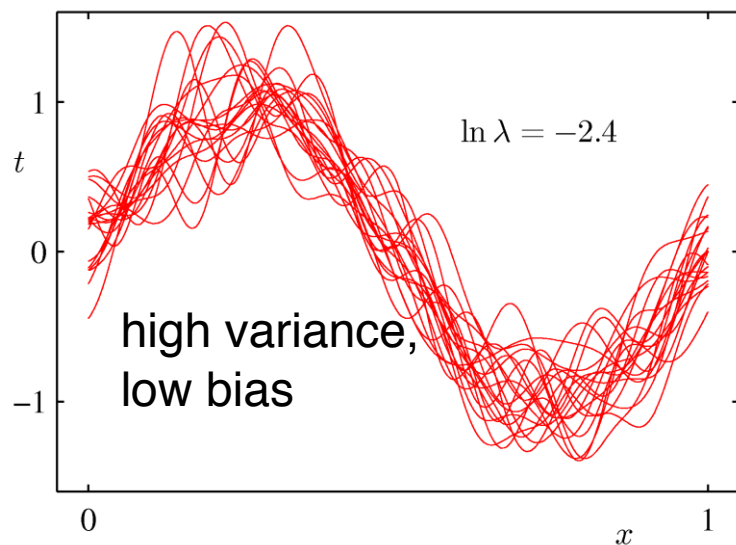
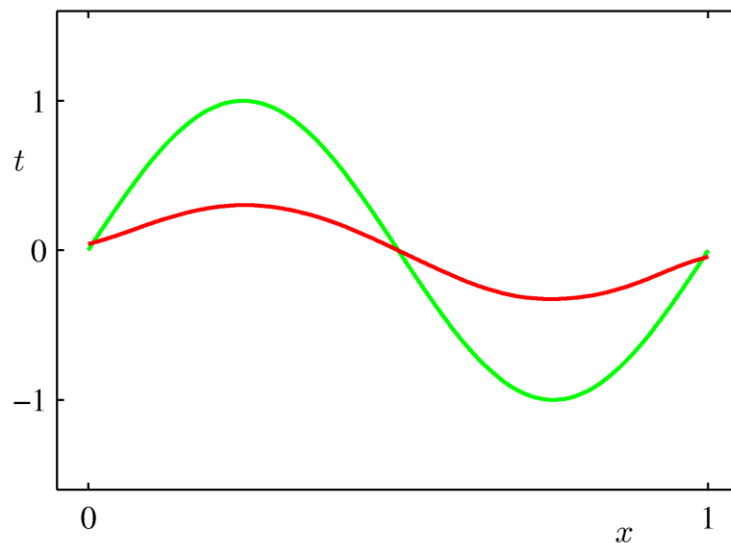
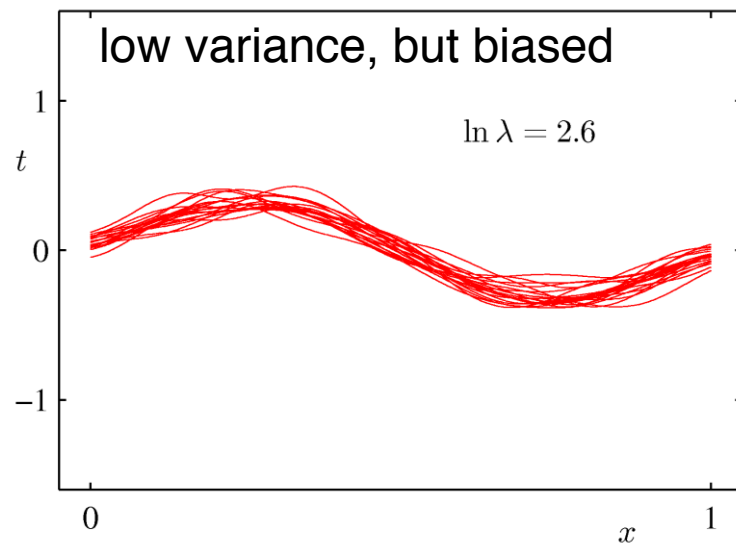


smooth out variance with bootstrap and averaging over many values.

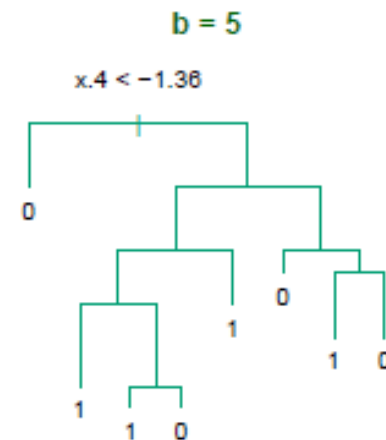
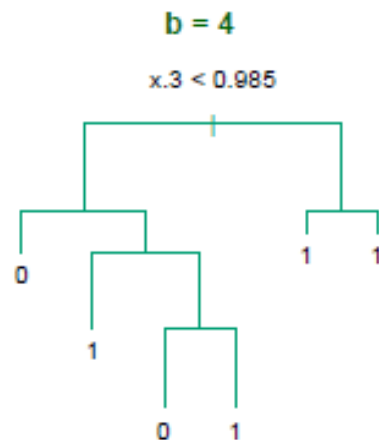
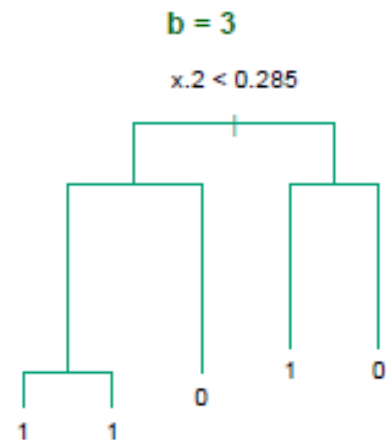
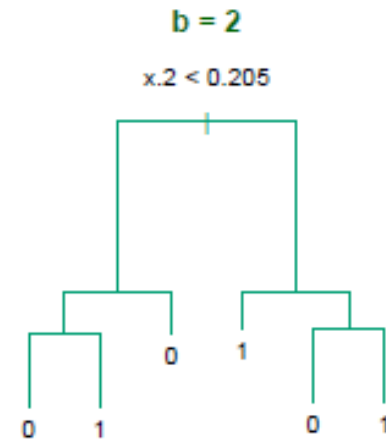
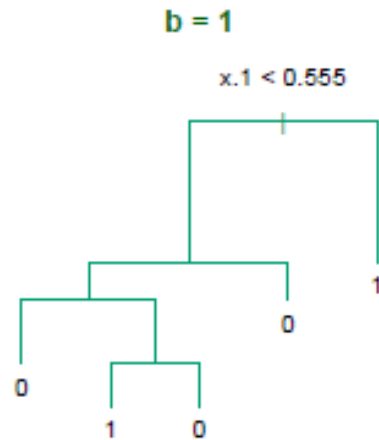
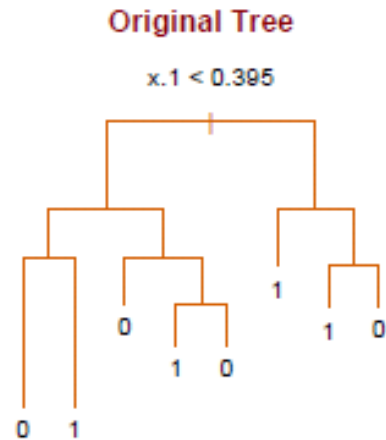


reduce variance w/o introducing much bias

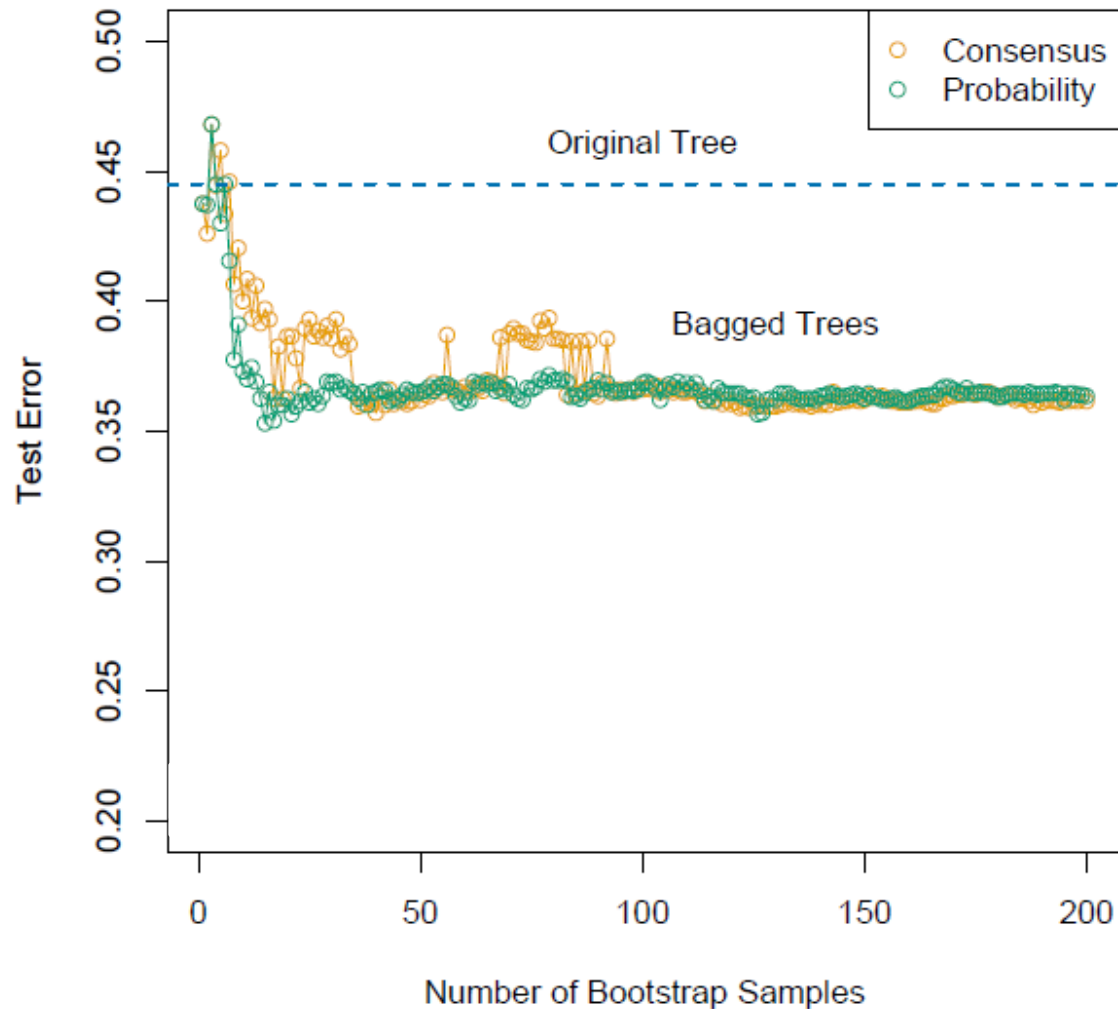
# Bias-Variance Trade-off



# Bagging Decision Trees



# Bagging Decision Trees



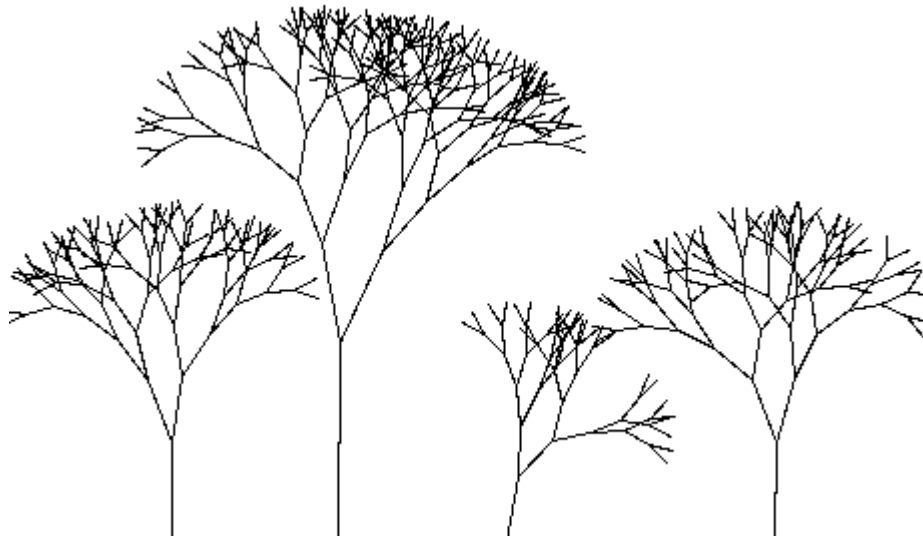
# Bagging

- Reduces overfitting (variance)
  - Normally uses one type of classifier
  - Decision trees are popular
  - Not helping with linear models
  - Easy to parallelize
- average linear models: you still get a line  
so you can't reduce variance this way
- linear model already reduces variance through  
its restrictions and assumptions

Each tree looks even more different.

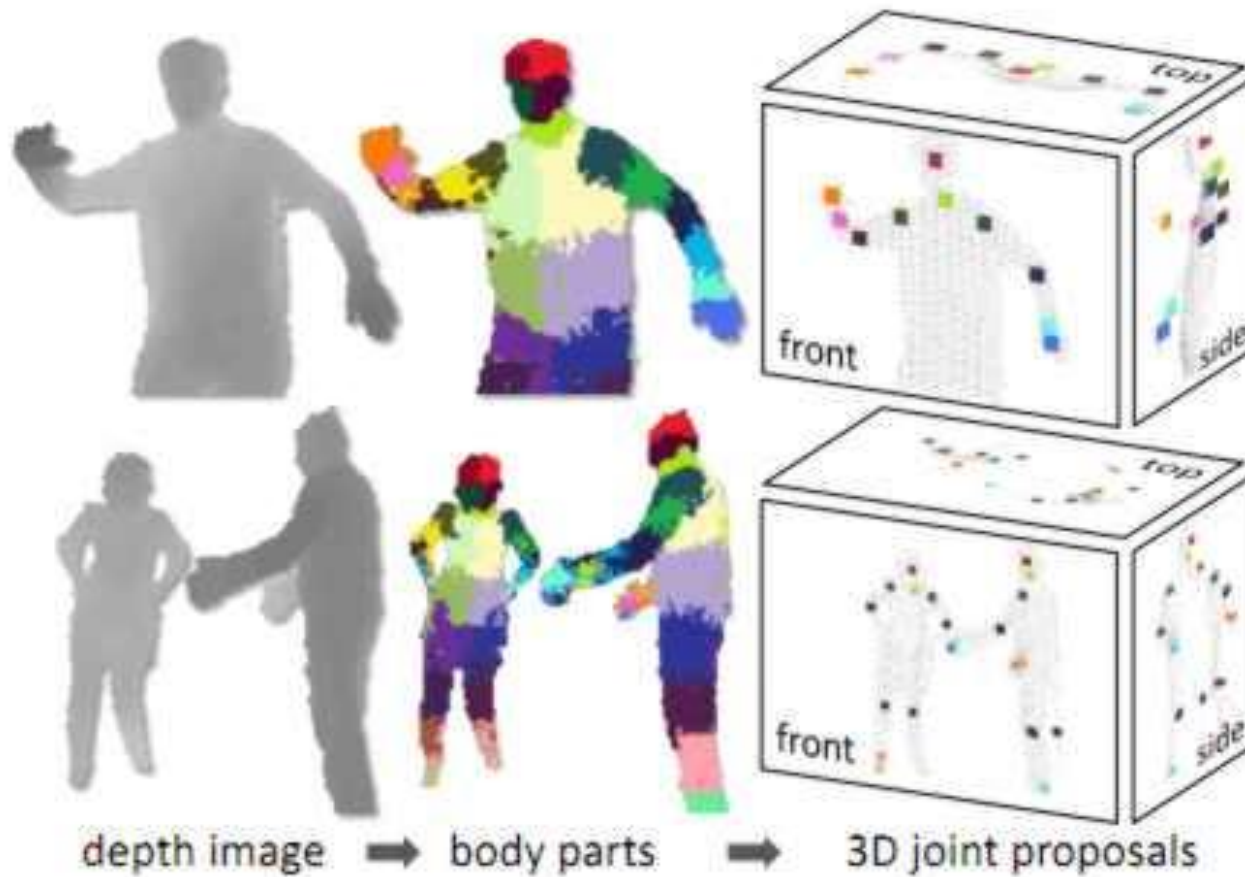
# Random Forest

- Builds upon the idea of bagging
- Each tree build from bootstrap sample
- Node splits calculated from **random feature subsets**



fairly generalizable detection of body parts working in real-time with arbitrary background.

# Random Forest – Fun Fact





hand\_tracking\_kinect.mp4

# Random Forest

- All trees are fully grown
- No pruning
- Two parameters
  - Number of trees
  - Number of features



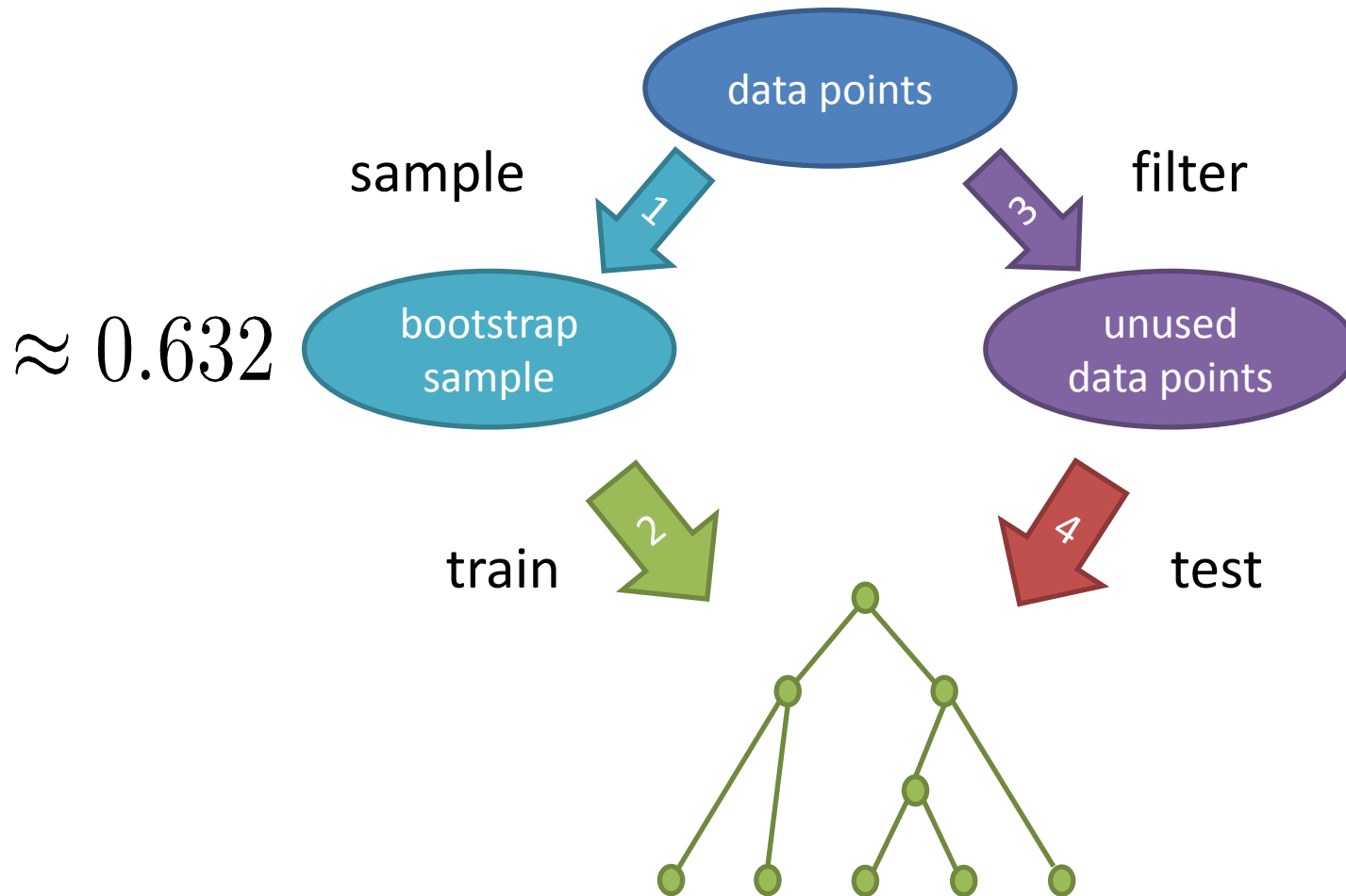
# Random Forest Error Rate

- Error depends on:
  - Correlation between trees (higher is worse)
  - Strength of single trees (higher is better)
- Increasing number of features for each split:
  - Increases correlation
  - Increases strength of single trees

# Out of Bag Error

- Each tree is trained on a bootstrapped sample
- About  $1/3$  of data points not used for training
- Predict unseen points with each tree
- Measure error

# Out of Bag Error



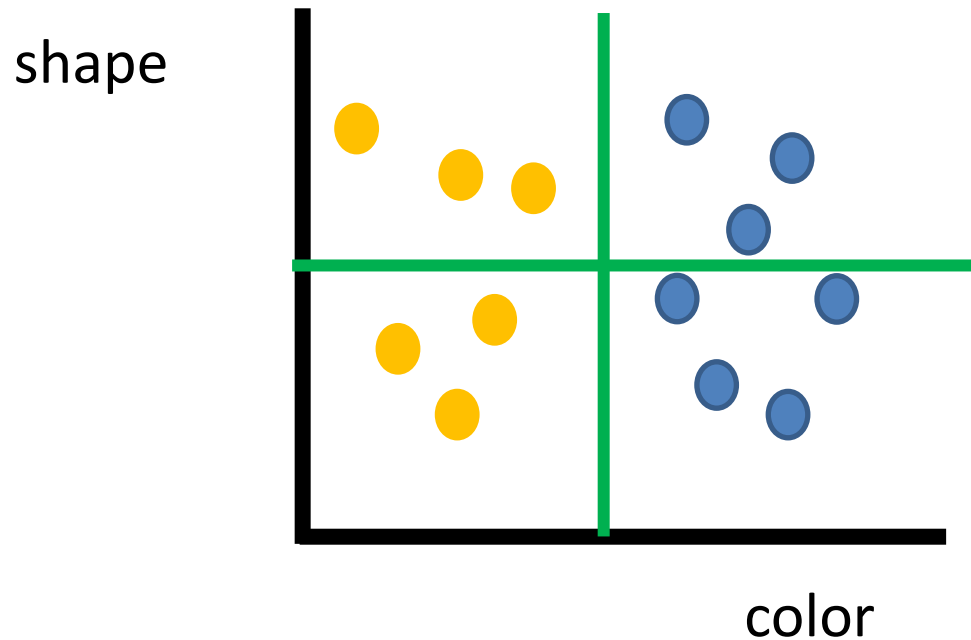
# Out of Bag Error

- Very similar to cross-validation
- Measured during training
- Can be too optimistic

# Variable Importance - 1

- Again use out of bag samples
- Predict class for these samples
- Randomly permute values of one feature
- Predict classes again
- Measure decrease in accuracy

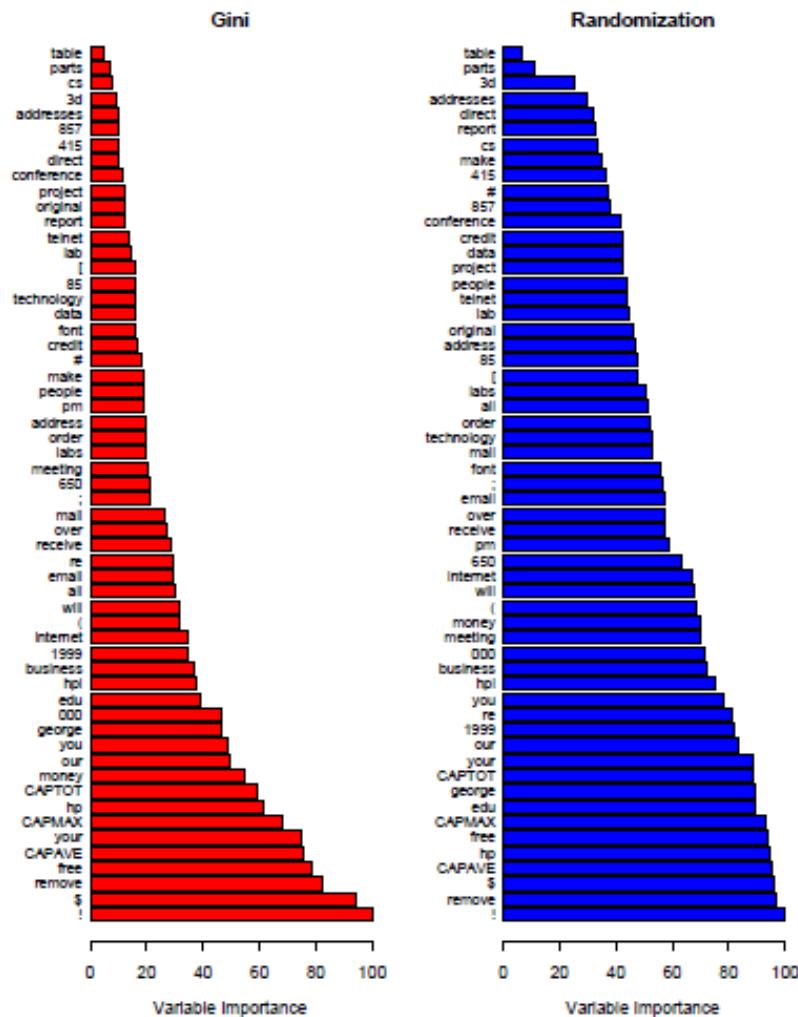
# Variable Importance - 1



# Variable Importance - 2

- Measure split criterion improvement
- Record improvements for each feature
- Accumulate over whole ensemble

# Example: Spam classification

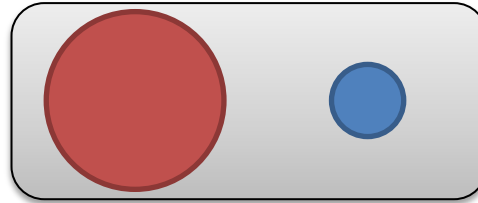


Randomization tends to spread out the variable importance more uniformly.

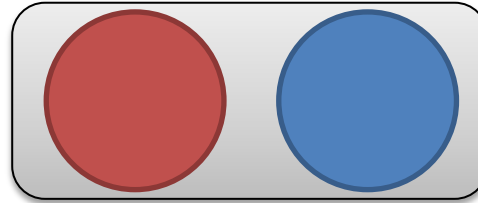


# Unbalanced Classes

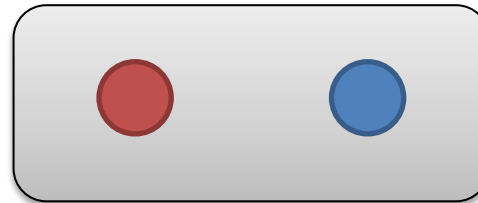
- The Problem:



- Oversample:

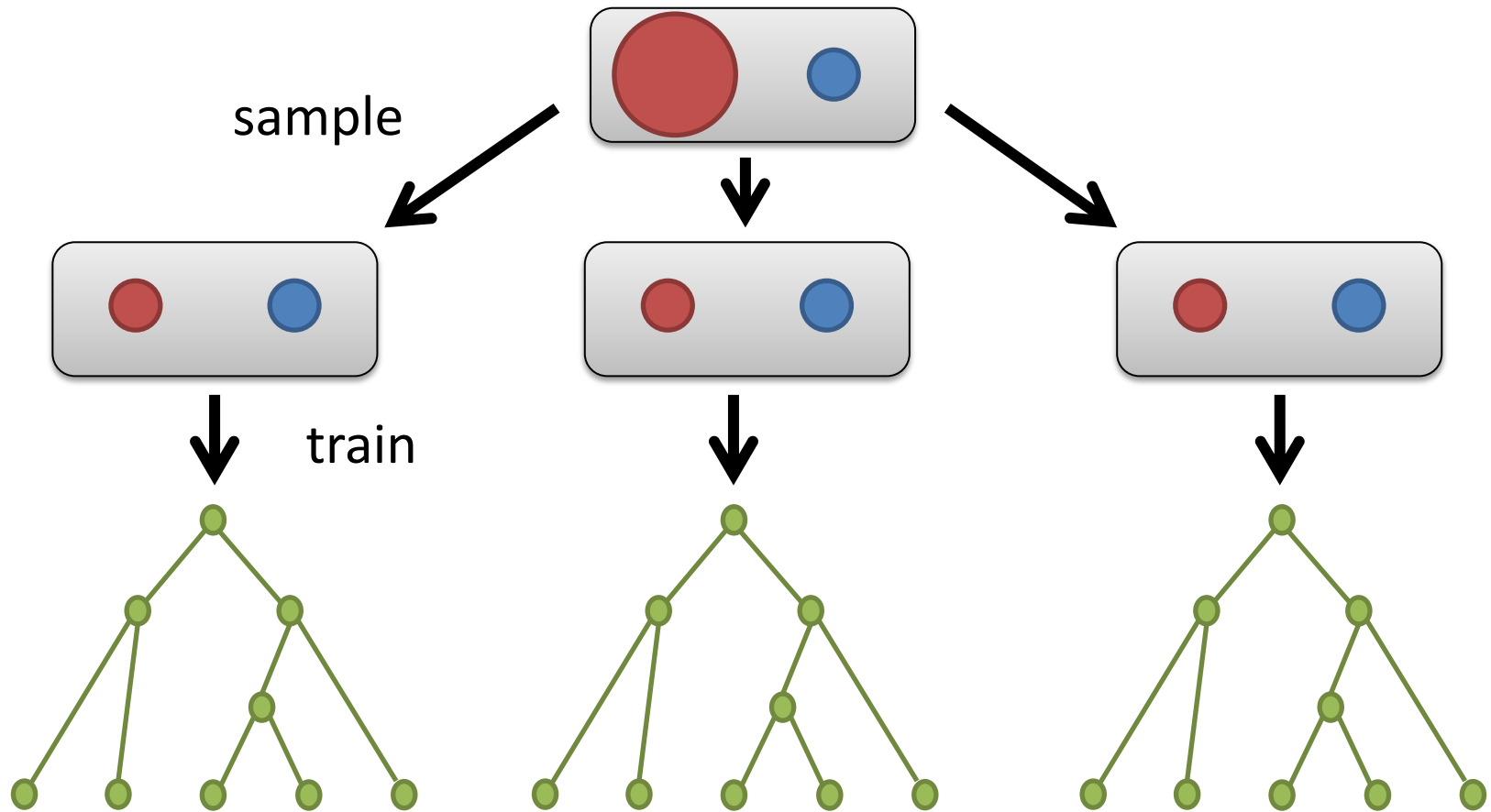


- Subsample:



- Subsample for each tree!

# Random Forest Subsampling



# Random Forest

- Similar to Bagging
- Easy to parallelize
- Packaged with some neat functions:
  - Out of bag error
  - Feature importance measure
  - Proximity estimation