

CS109 – Data Science

Joe Blitzstein, Hanspeter Pfister, Verena Kaynig-Fittkau

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

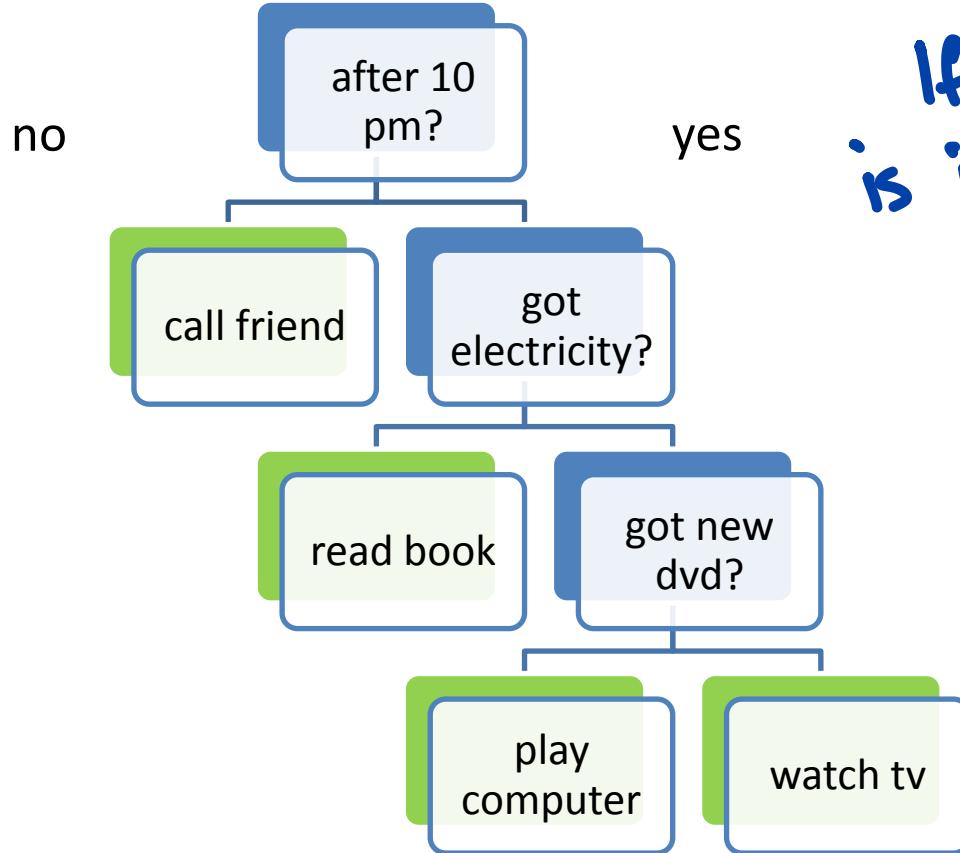
- “Elements of Statistical Learning”
- <http://statweb.stanford.edu/~tibs/ElemStatLearn/>
- “Pattern Recognition and Machine Learning”
- <http://research.microsoft.com/en-us/um/people/cmbishop/PRML/>

Next Topics

- Tree classifier
- Bagging
- Random Forest



Decision Tree



If small, it
is interpretable.
Simple.
Intuitive.
Is not the
core with
Support Vector
machines
(SVM)

Decision Trees

- Fast training
- Fast prediction
- Easy to understand
- Easy to interpret

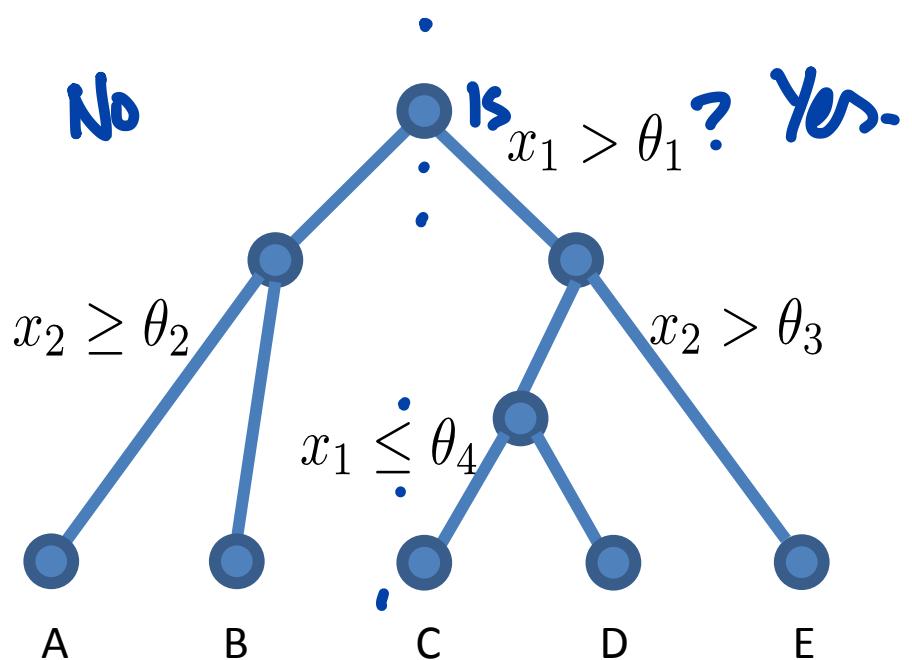
leads to

Not predictive.

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

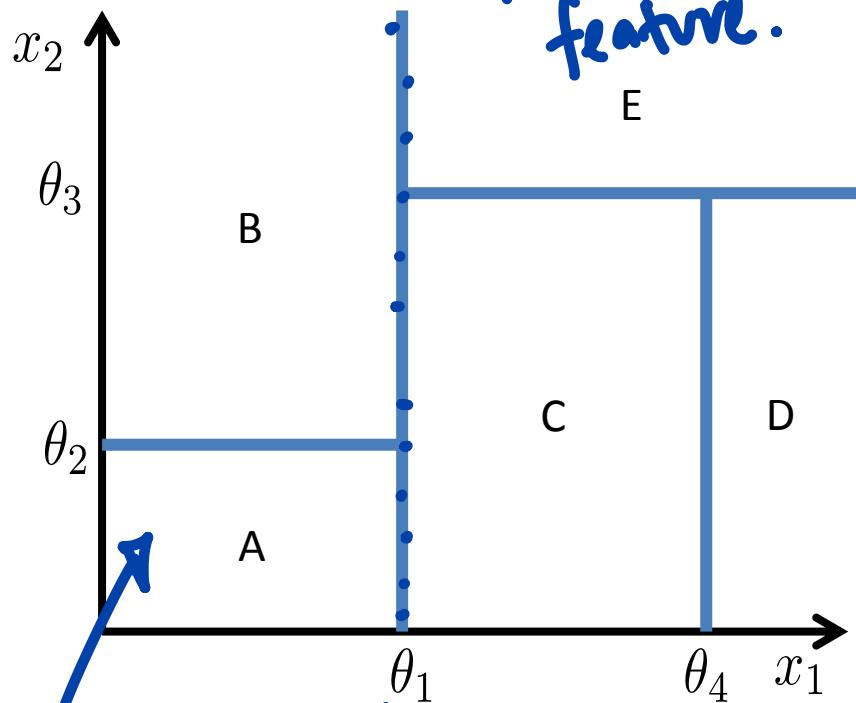
Decision Tree - Idea

↳ Single decision rule.



• Uses thresholds.

Cells /
classes



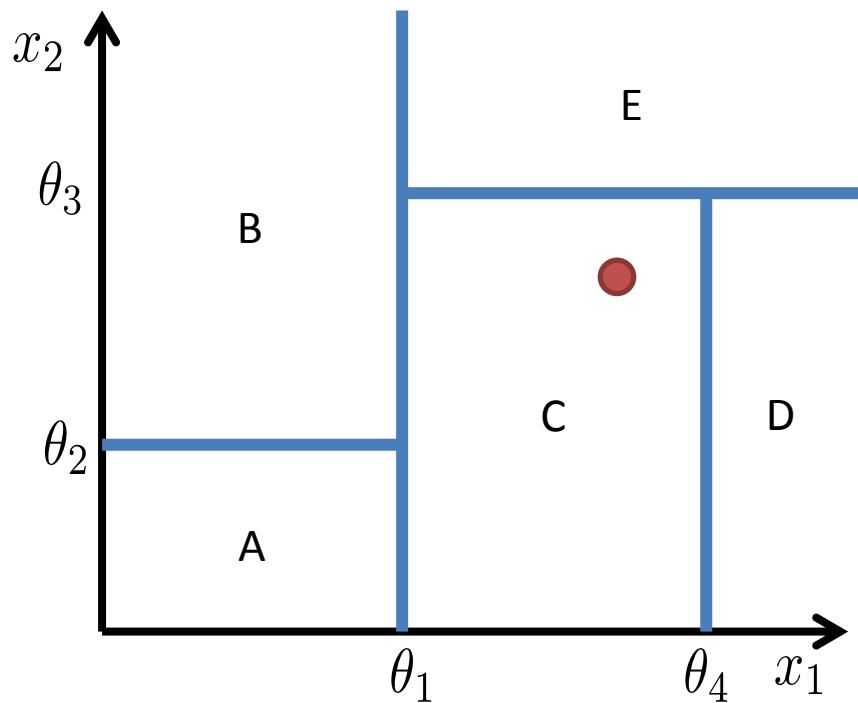
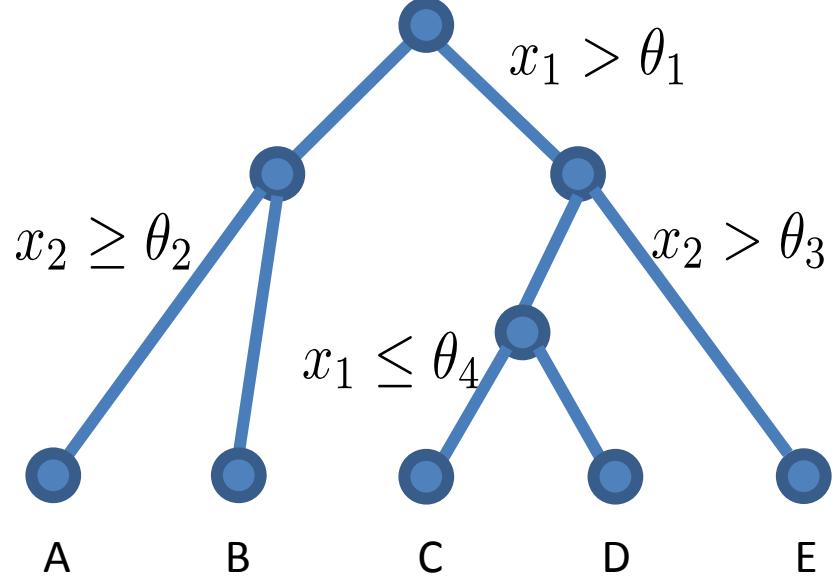
4 splits in half.

straight line
from 1 feature.

Decision Tree - Idea

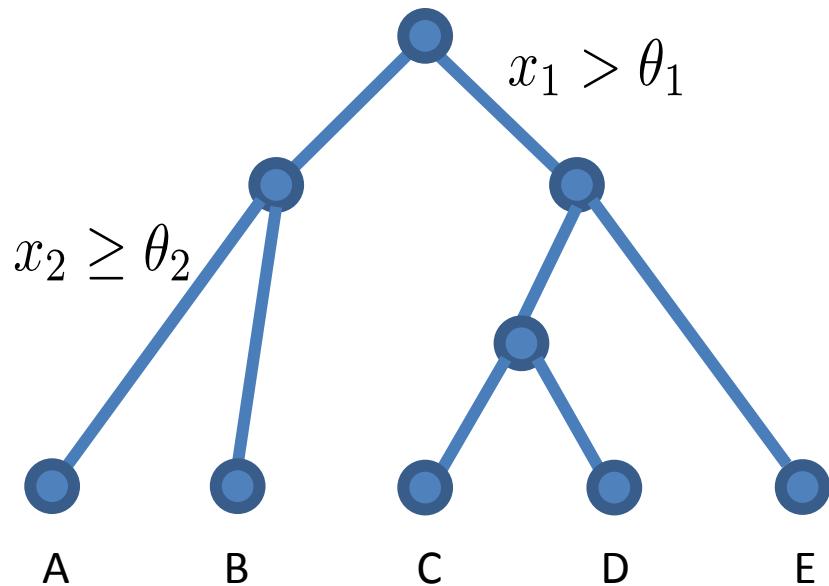
- What is the benefit of using only one feature at a time? *Fast performance. Only need to test @ 1 feature. Binary search w/ lots of splits. Also interpretable. Sequential process w/ categories.*
- What is the drawback?
Straight lines, would need to have a lot of splits → big tree
↓
SVM requires # communion

Decision Tree - Prediction

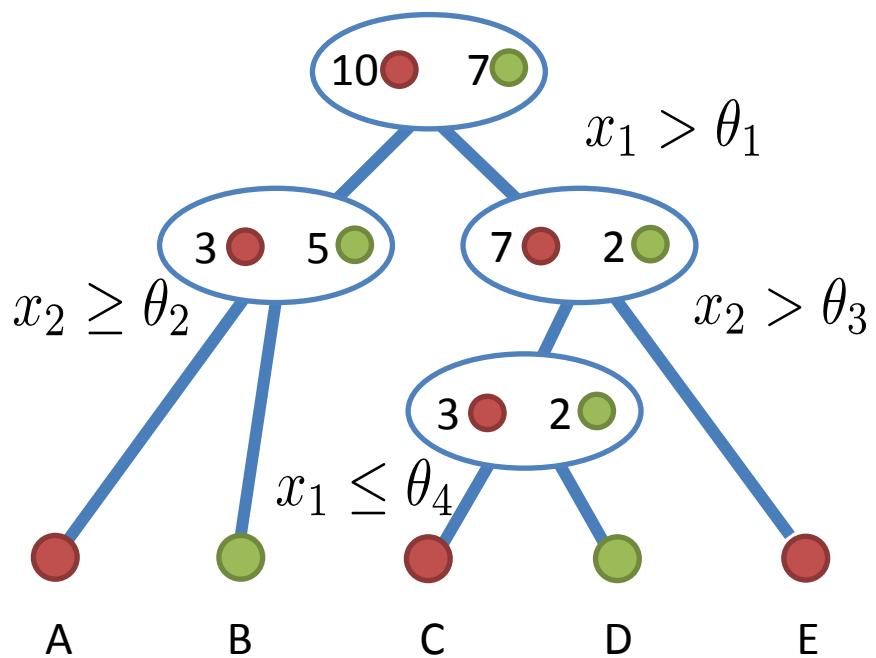


Decision Tree -Training

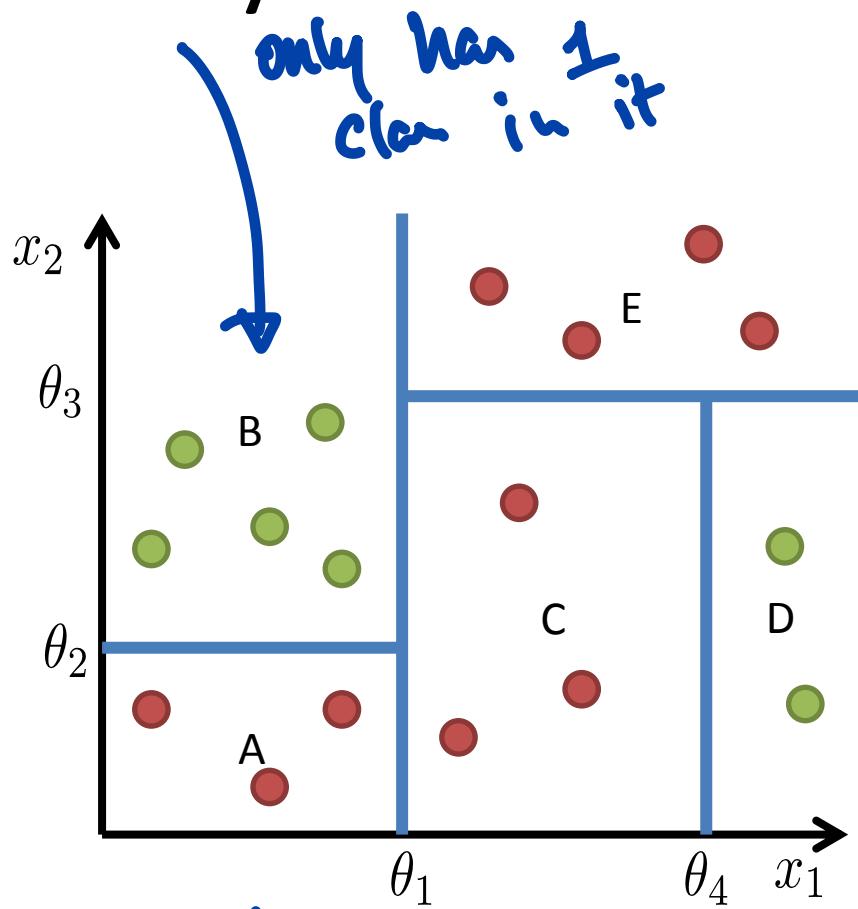
- Learn the tree structure:
 - which feature to query x_1, x_2
 - which threshold to choose θ_1, θ_2



Node Purity

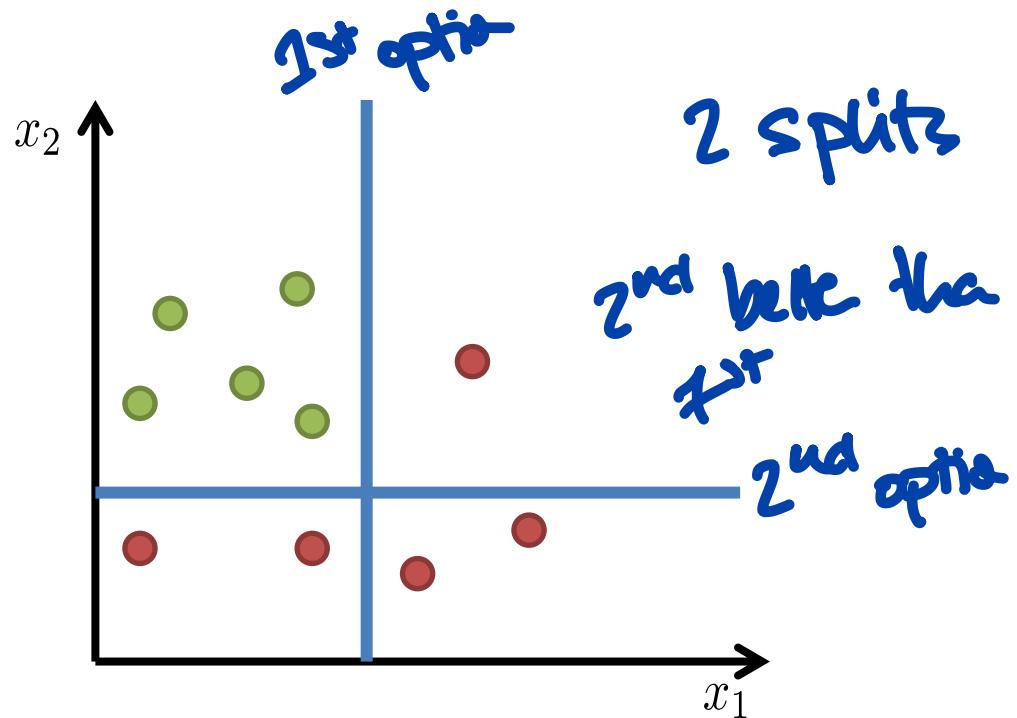


Extreme would be a cell/cluster
for each node



Gini Impurity *(default in sklearn)*

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



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

Gini Impurity

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

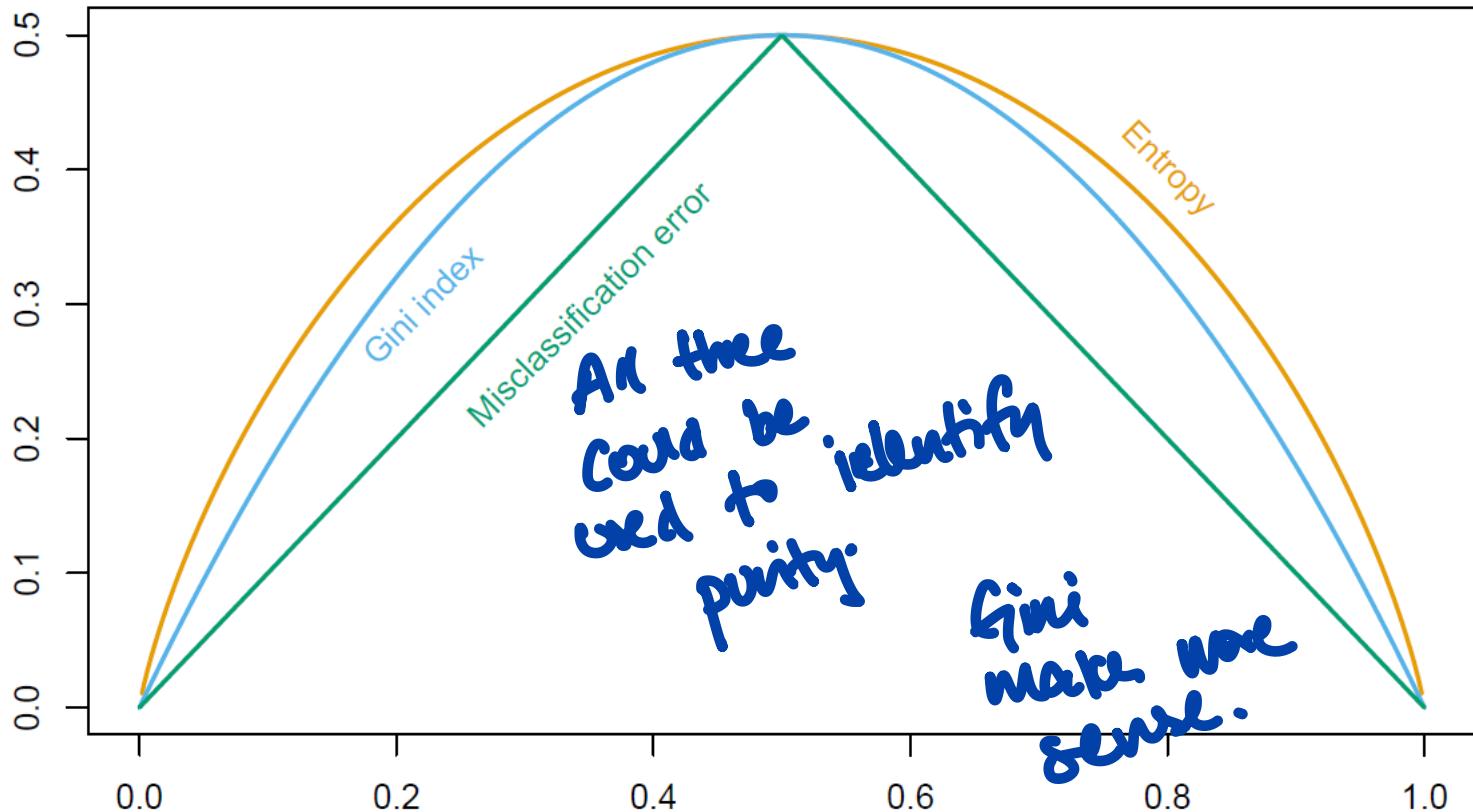
→ red, green, blue from above.

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

true
class

wrong
prediction

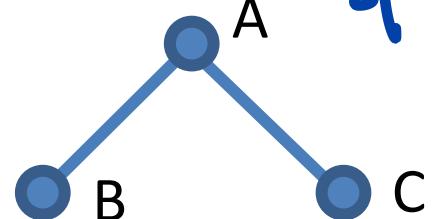
Gini Impurity



Node Purity Gain

Want to judge the split not just
in node after n split.

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



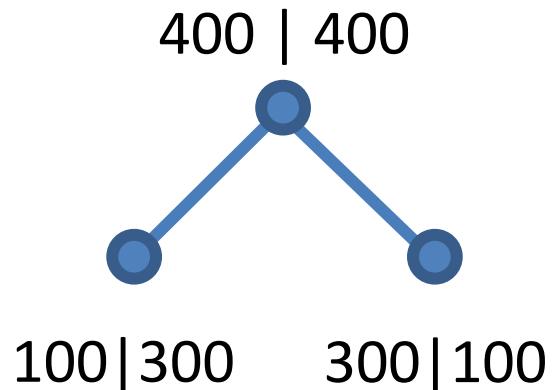
$$\Delta I_G = I_G(A) - \frac{N(B)}{N(A)} I_G(B) - \frac{N(C)}{N(A)} I_G(C)$$

Misclassification

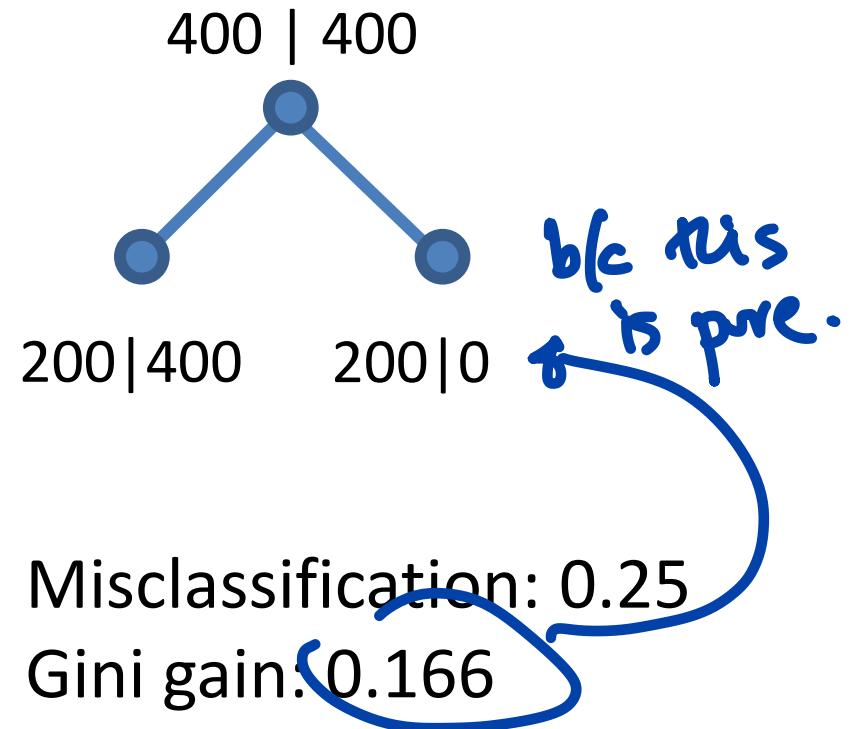
- $\frac{1}{N} \sum_i^N \mathbf{1}(\hat{y}_i \neq y_i)$
Is the predicted label \neq same
as the actual label?
- not differentiable

Comparison Gini vs Misclassification

- Binary problem: 400 samples per class



Misclassification: 0.25
Gini gain: 0.125



Misclassification: 0.25
Gini gain: 0.166

Pseudocode

- Check if already finished *-if leaf node is pure.*
- 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
 - Take the best next step, doesn't look ahead.
- Does this produce an optimal tree?
- What does optimal tree mean?
 - ↳ Not easy to define → purity? interpretable?

When to Stop

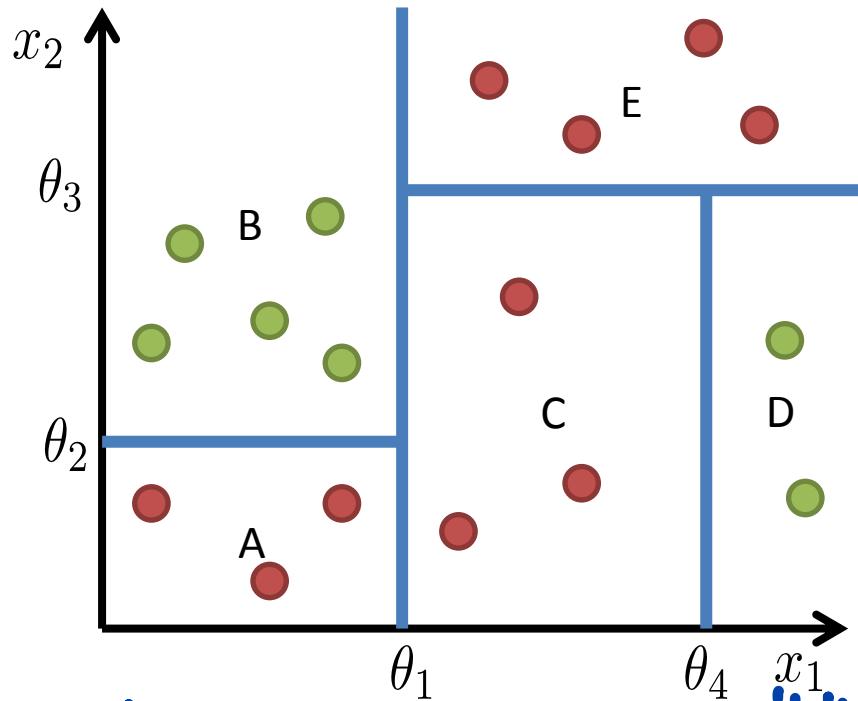
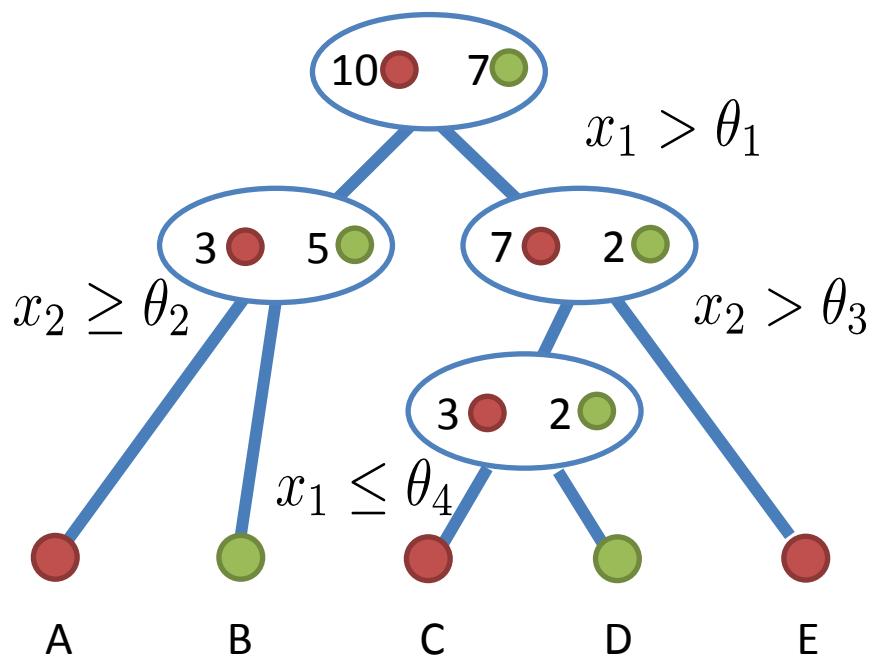
- node contains only one class
- node contains less than x data points
- max depth is reached \rightarrow limits time to complete.
- node purity is sufficient \rightarrow not completely pure
- you start to overfit \Rightarrow cross-validation

\hookrightarrow Computationally intense but w/ current computer can do.

\hookrightarrow DT then to overfit.

Tree Pruning

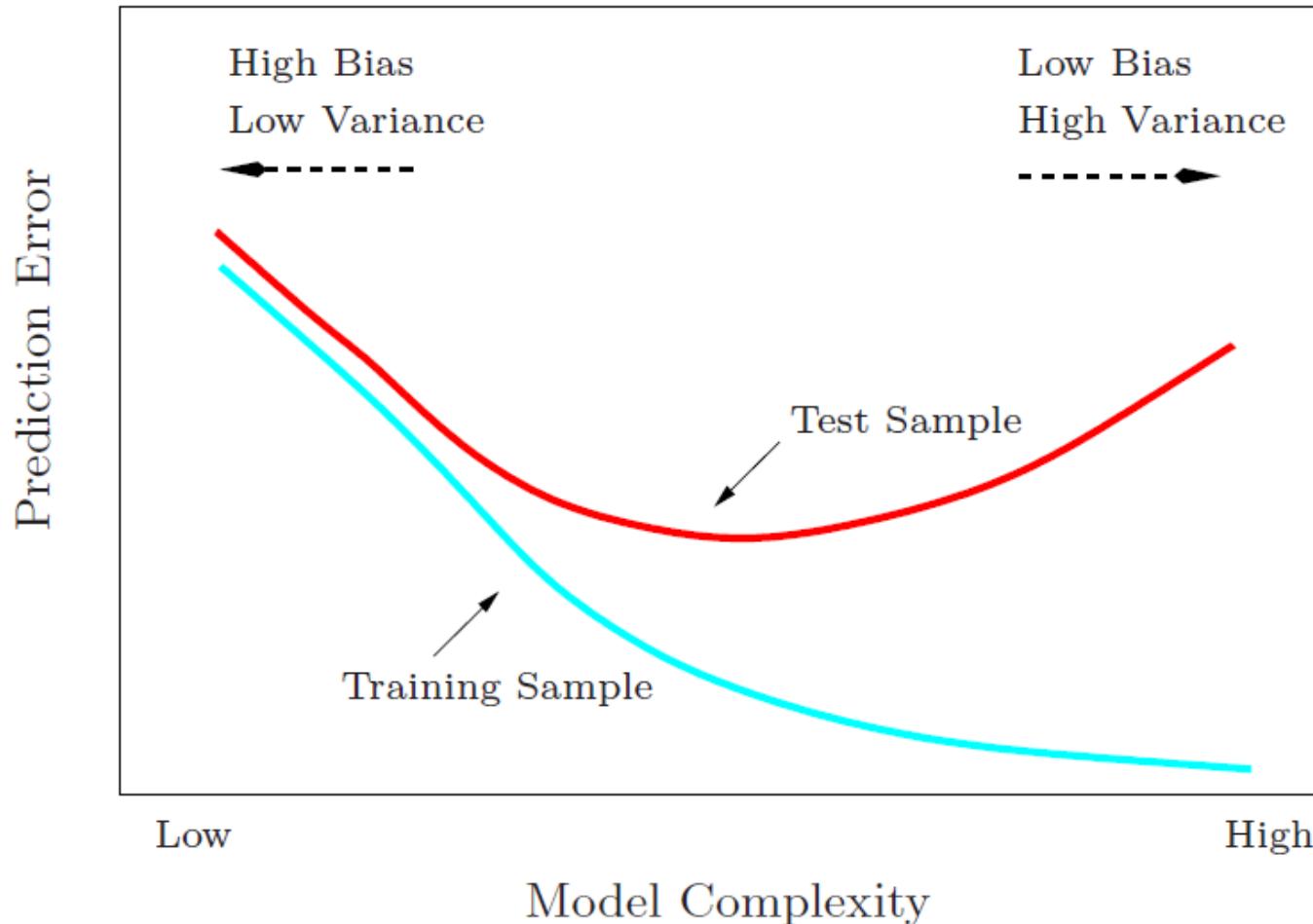
Grow the whole tree then remove leaves.



If impure, take the majority. Or use probability ratio?

How do you make a prediction for the merged cell?

Pruning and Complexity



Decision Trees - Disadvantages

- Sensitive to small changes in the data
 - Overfitting
 - Only axis aligned splits
- | _____
- \Rightarrow high variance.
leads to overfitting.

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	▲	○

▷ Needs #s.

→ b/c 1 feature @ a time.

has to care about dot product.
↳ means computation expensive.

if small.

▷ doesn't generalize well.

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



A large, semi-transparent purple sphere with concentric rings of light emanating from its center. A smaller, glowing purple cube sits at the bottom of the sphere. The background is dark, making the purple glow stand out.

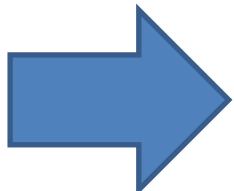
Netflix Prize

Netflix Prize

- Take home messages:
 - Early gains are easy. (80%).
 - Gains through an ensemble approach.
(avg. of 800 models)

Ensemble Methods

- A single decision tree does not perform well
- But, it is super fast (*training + testing*)
- What if we learn multiple trees?



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

} Can't just
run it
100 times.

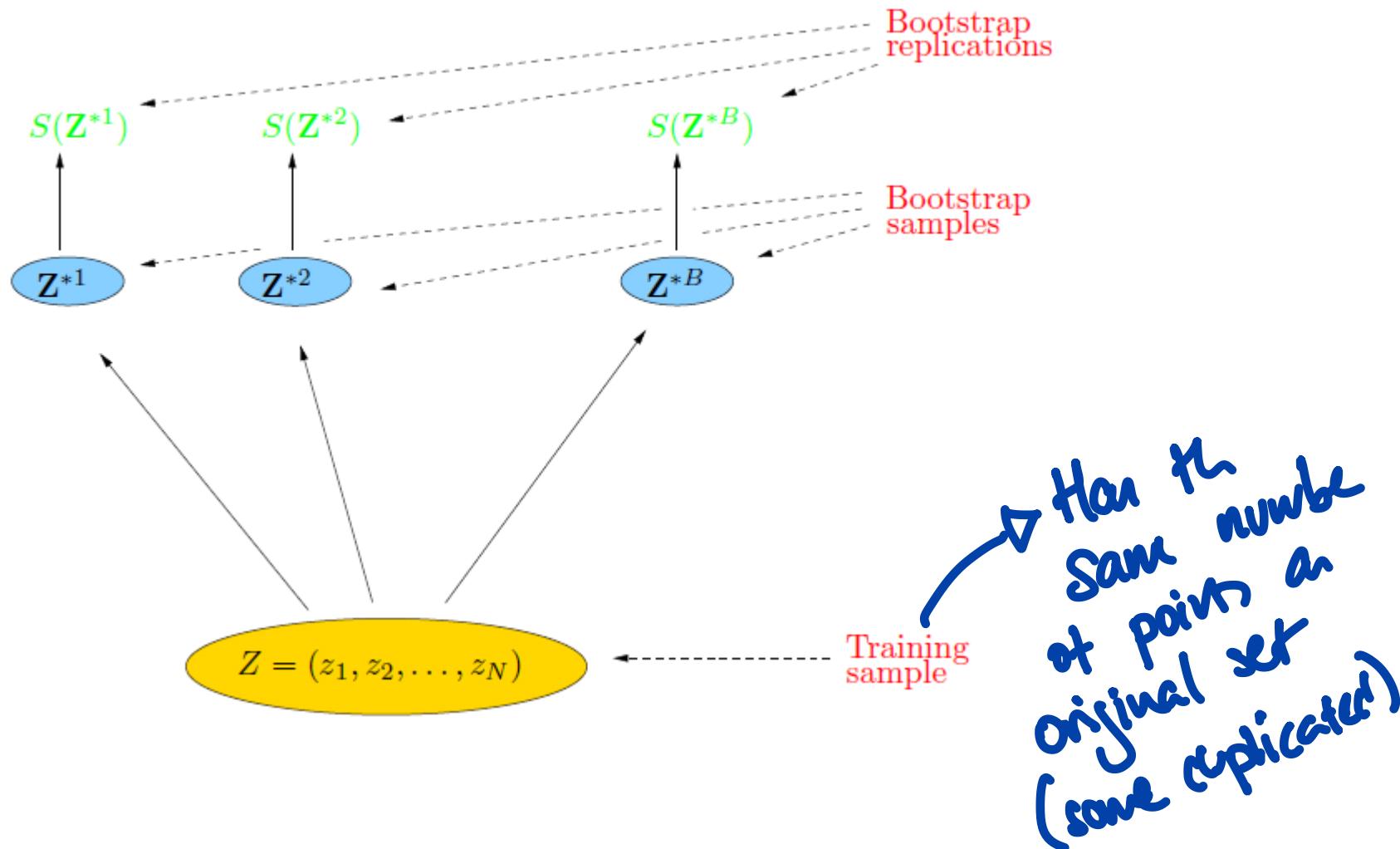
Bootstrap



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?

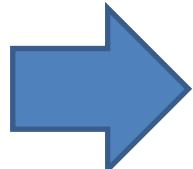
No. B/c same data points.

↳ Bootstrap sets will have overlap.

↳ Okay for DT

Bootstrap vs Cross-validation

- Bootstrap has overlap in data sets



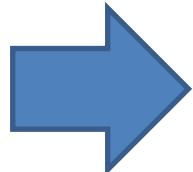
Do not use simple bootstrap to generate train and test data from same data set.

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

Probability of choosing n

Bootstrap vs Cross-validation

- Bootstrap has overlap in data sets



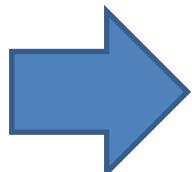
Do not use simple bootstrap to generate train and test data from same data set.

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

Probability of not choosing n

Bootstrap vs Cross-validation

- Bootstrap has overlap in data sets



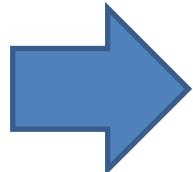
Do not use simple bootstrap to generate train and test data from same data set.

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

Probability of not choosing n in N draws

Bootstrap vs Cross-validation

- Bootstrap has overlap in data sets



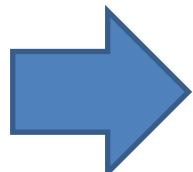
Do not use simple bootstrap to generate train and test data from same data set.

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

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.

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

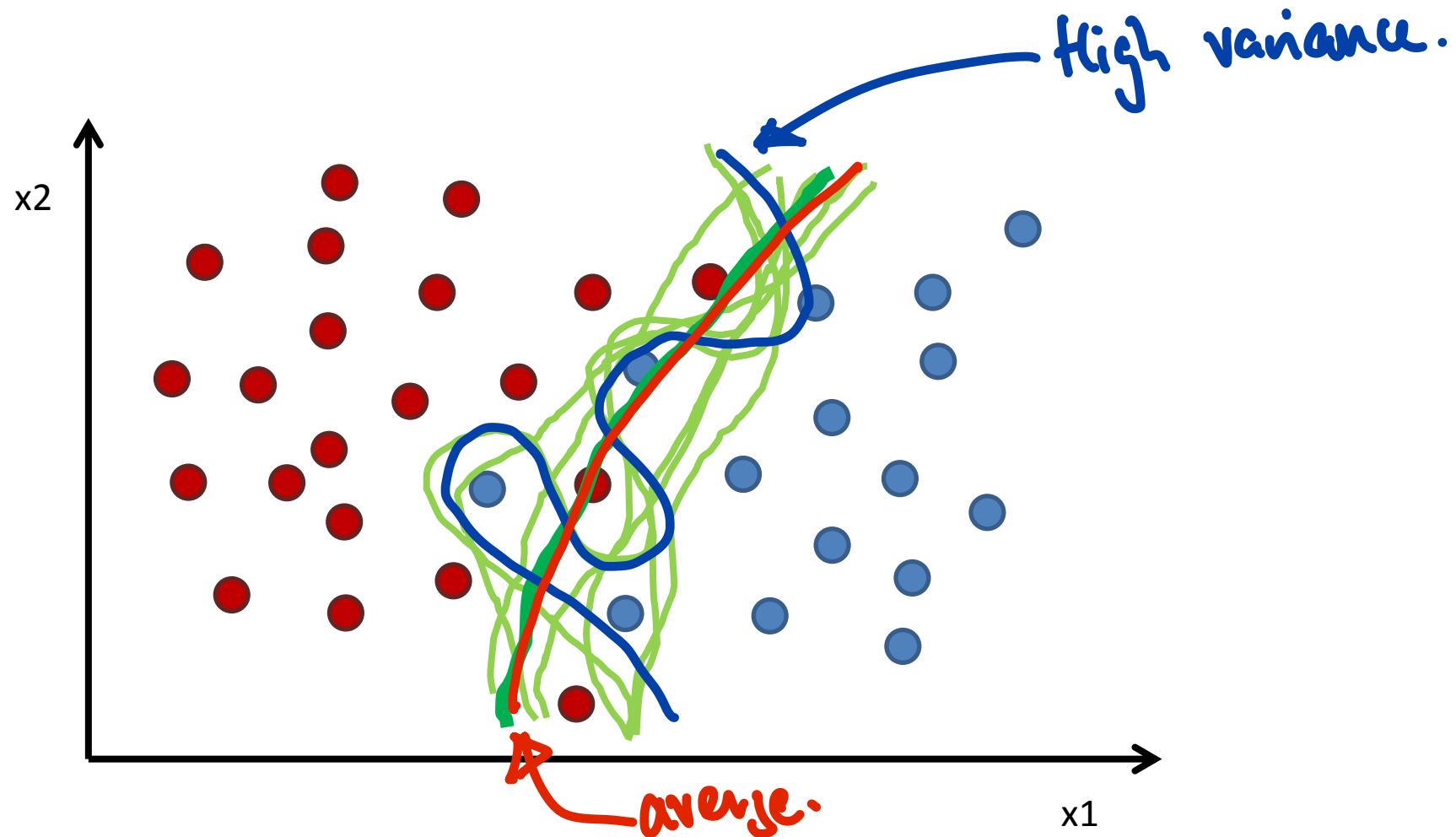
$$\approx 0.632$$

This number is important later

Bagging

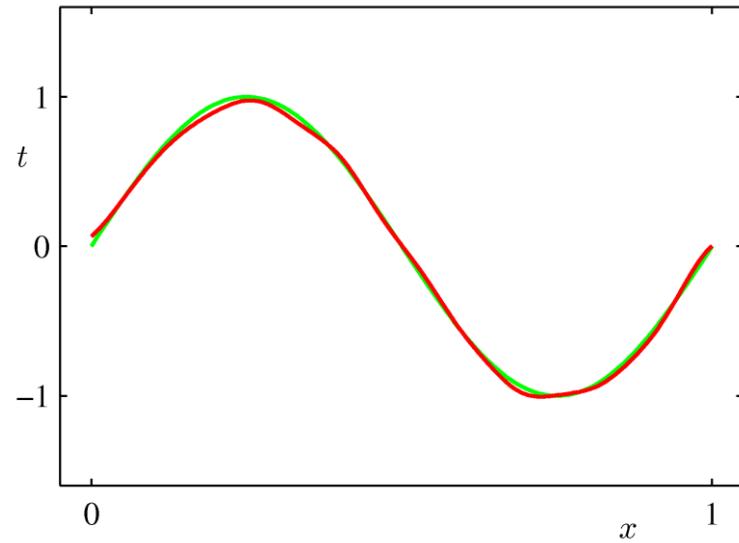
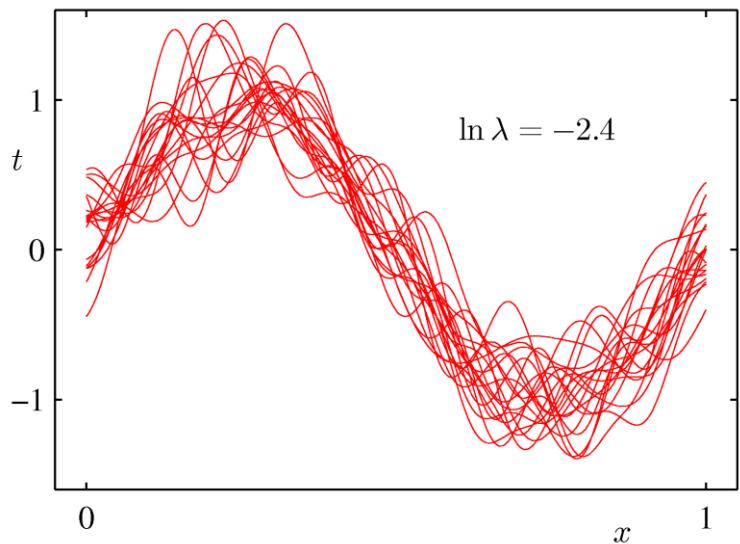
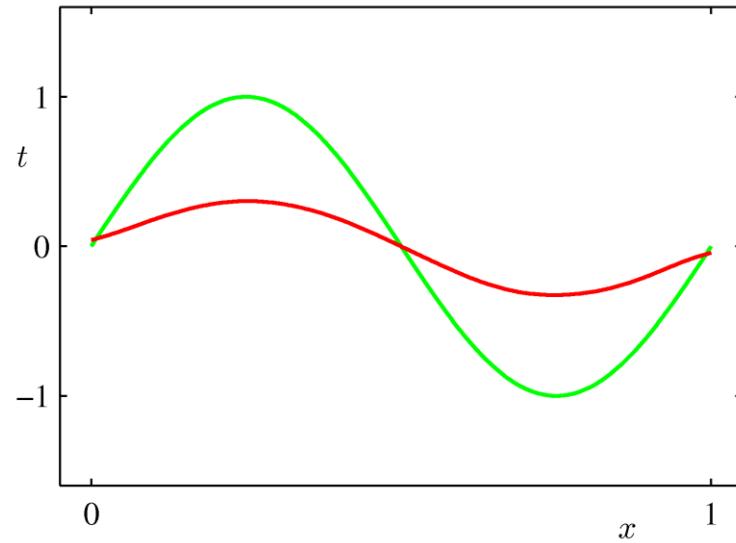
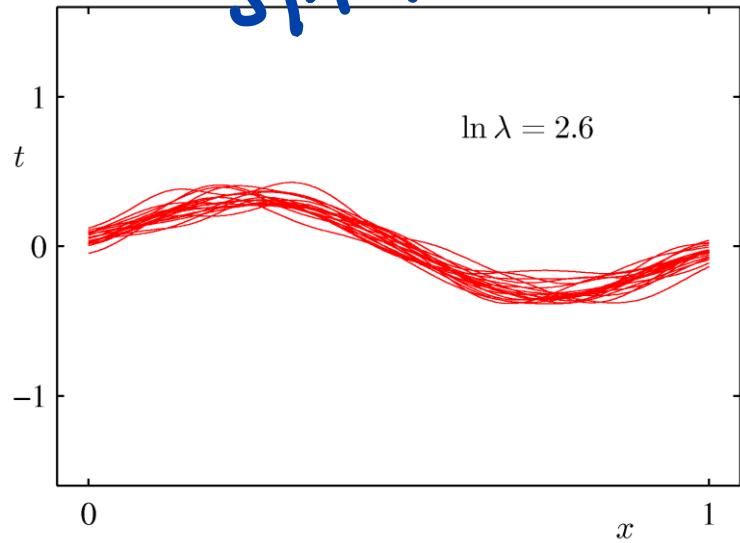
- Bootstrap aggregating
- Sample with replacement from your data set
- Learn a classifier for each bootstrap sample
- Average the results

Bagging Example



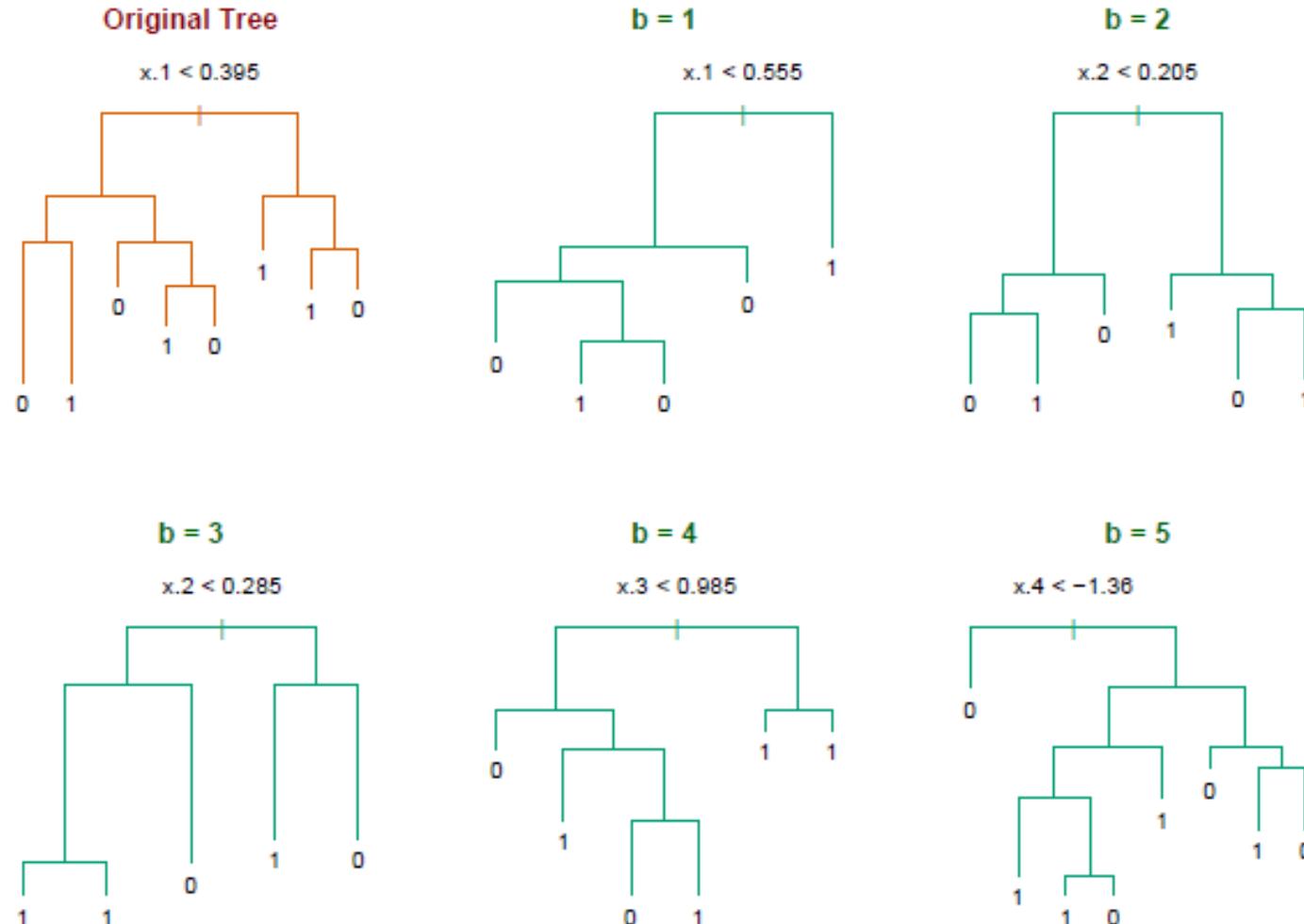
Bias-Variance Trade-off

↳ Bagging reduces variance w/o increasing bias too much.

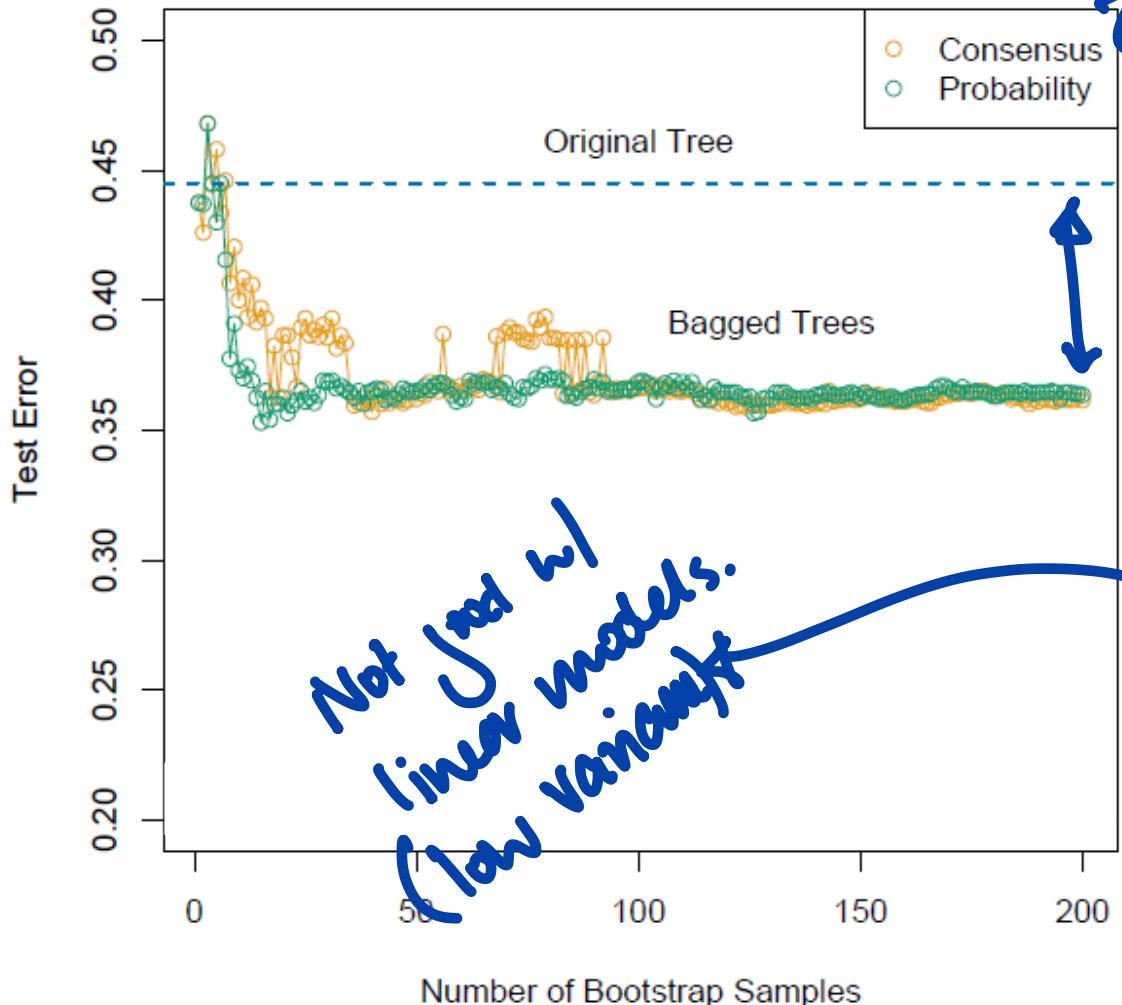


Bagging Decision Trees

↳ randomness is introduced.



Bagging Decision Trees



Two option.
Consensus is
majority
approach.

Reduction due
to bagging

Bagging can
be done w/
oth approach
that have high
variance + are
fast.

Bagging

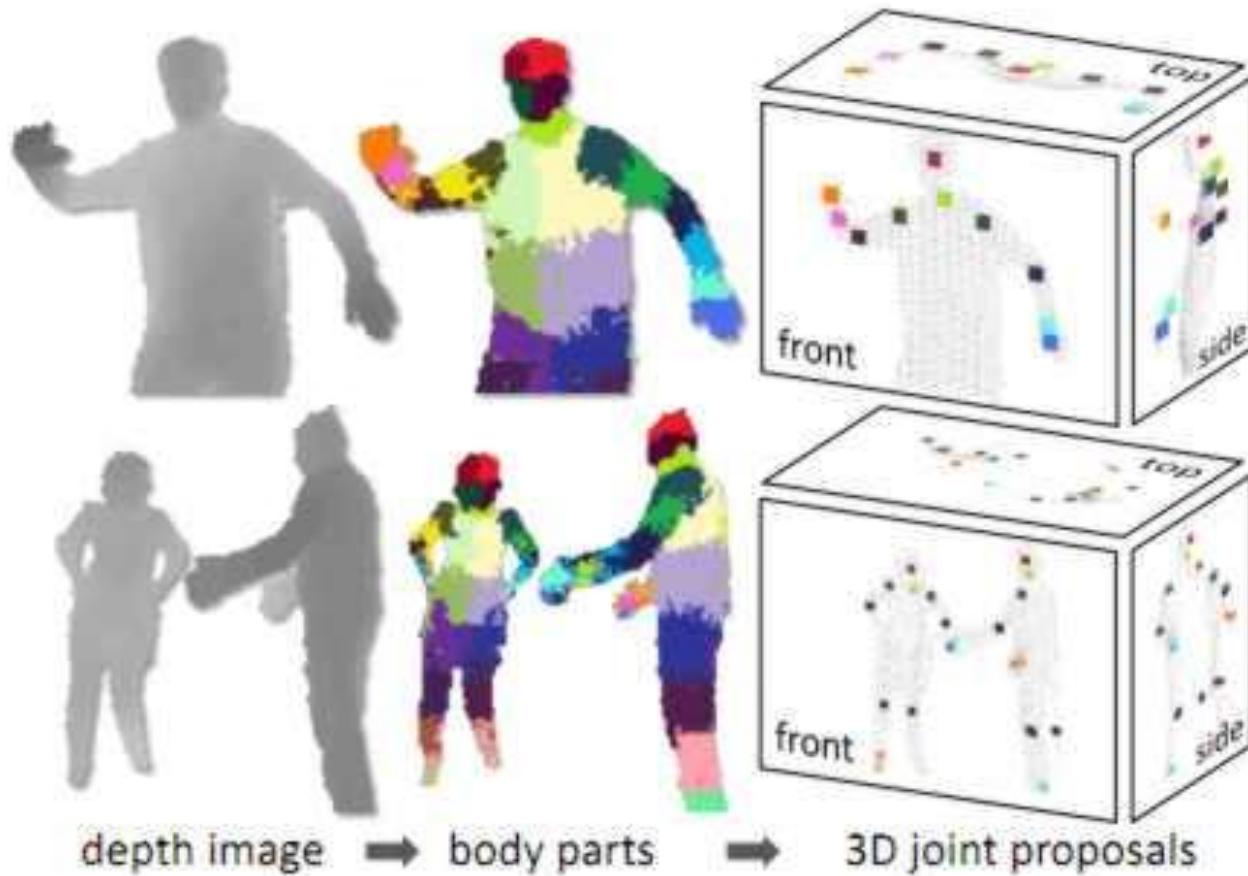
- Reduces overfitting (variance)
- Normally uses one type of classifier
- Decision trees are popular
- Not helping with linear models
- Easy to parallelize

Random Forest

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



Random Forest – Fun Fact





hand_tracking_kinect.mp4

<http://research.microsoft.com/en-us/projects/handpose/>

Random Forest

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

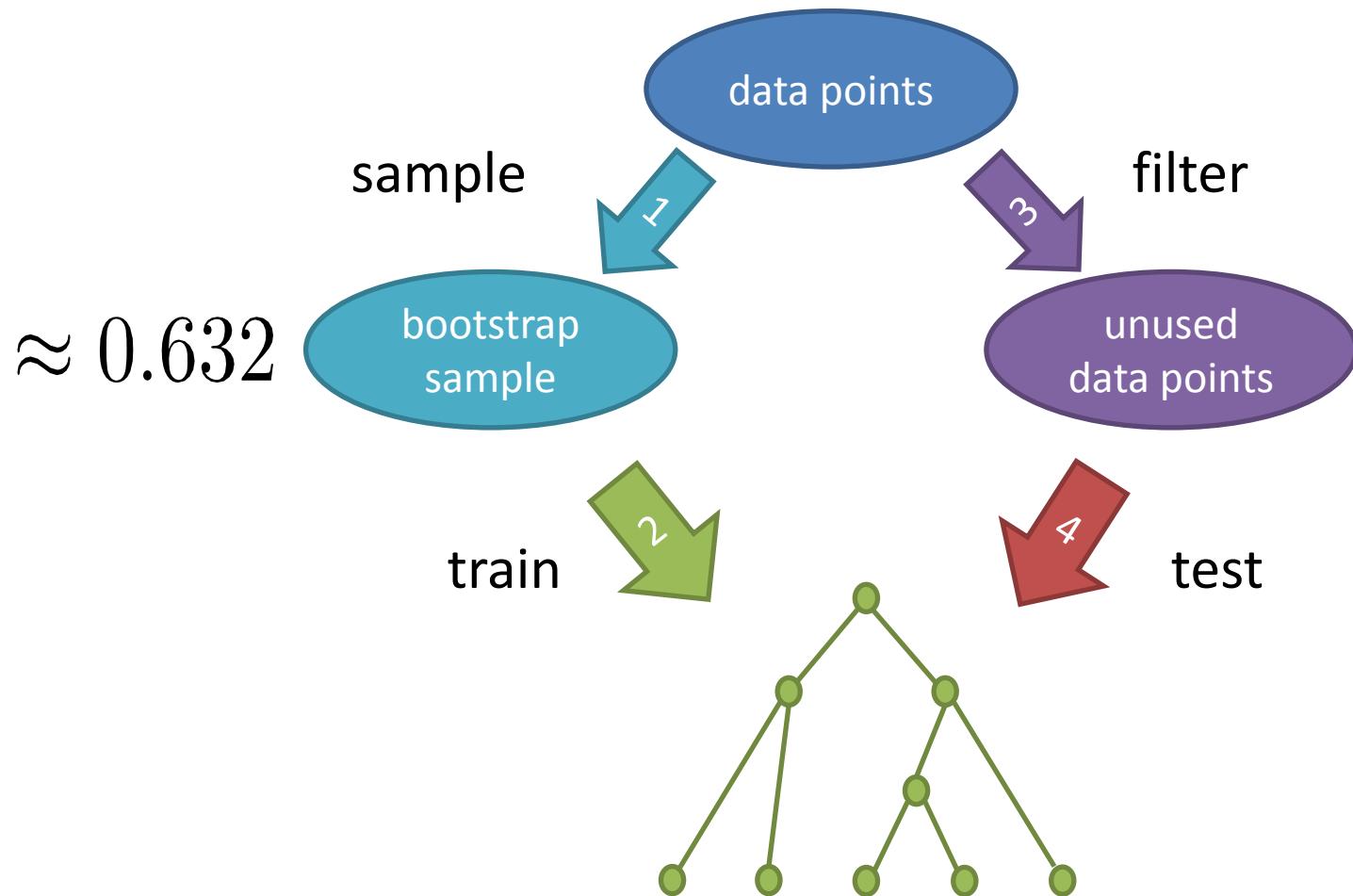
Random Forest Error Rate

- 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



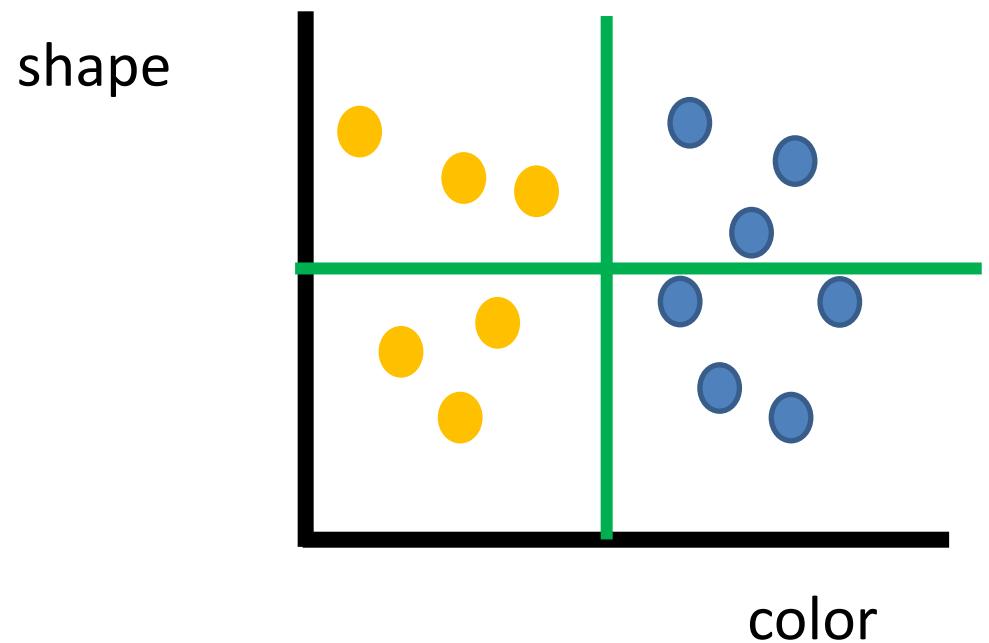
Out of Bag Error

- 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

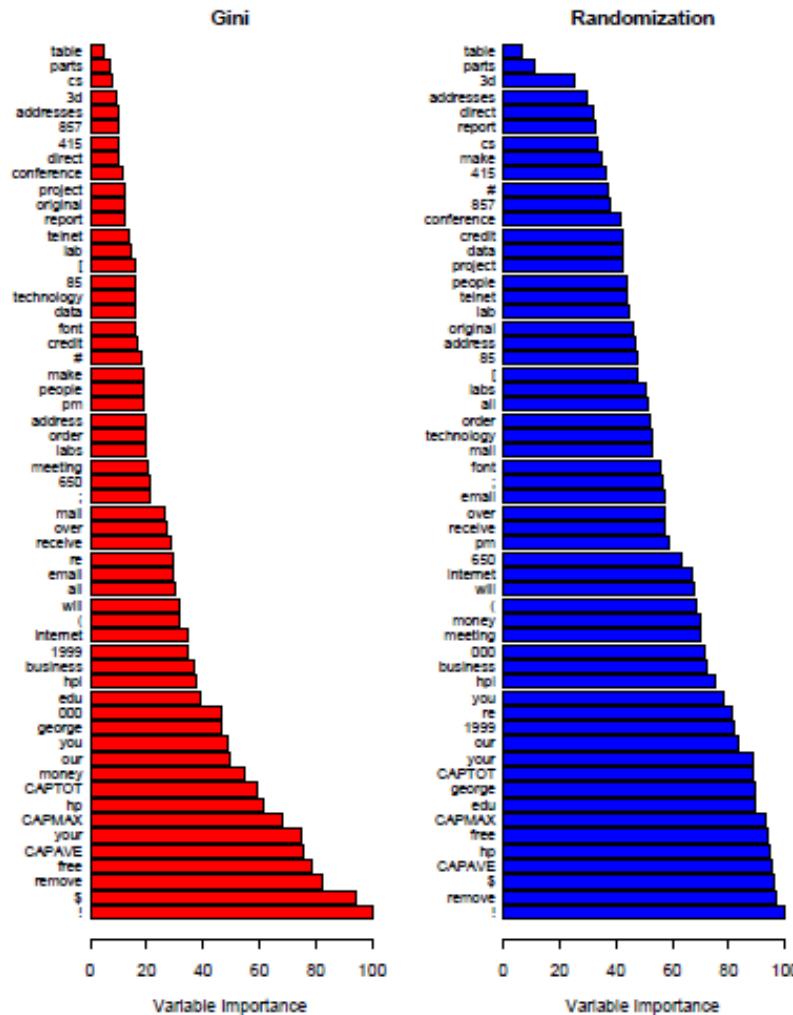
Variable Importance - 1



Variable Importance - 2

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

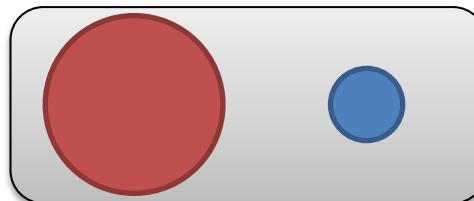
Example: Spam classification



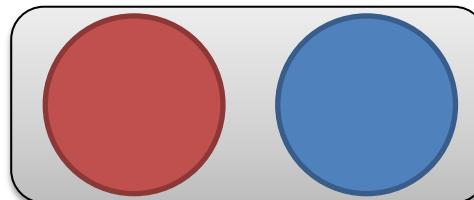
Randomization tends to spread out the variable importance more uniformly.

Unbalanced Classes

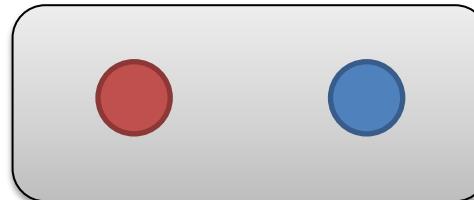
- The Problem:



- Oversample:

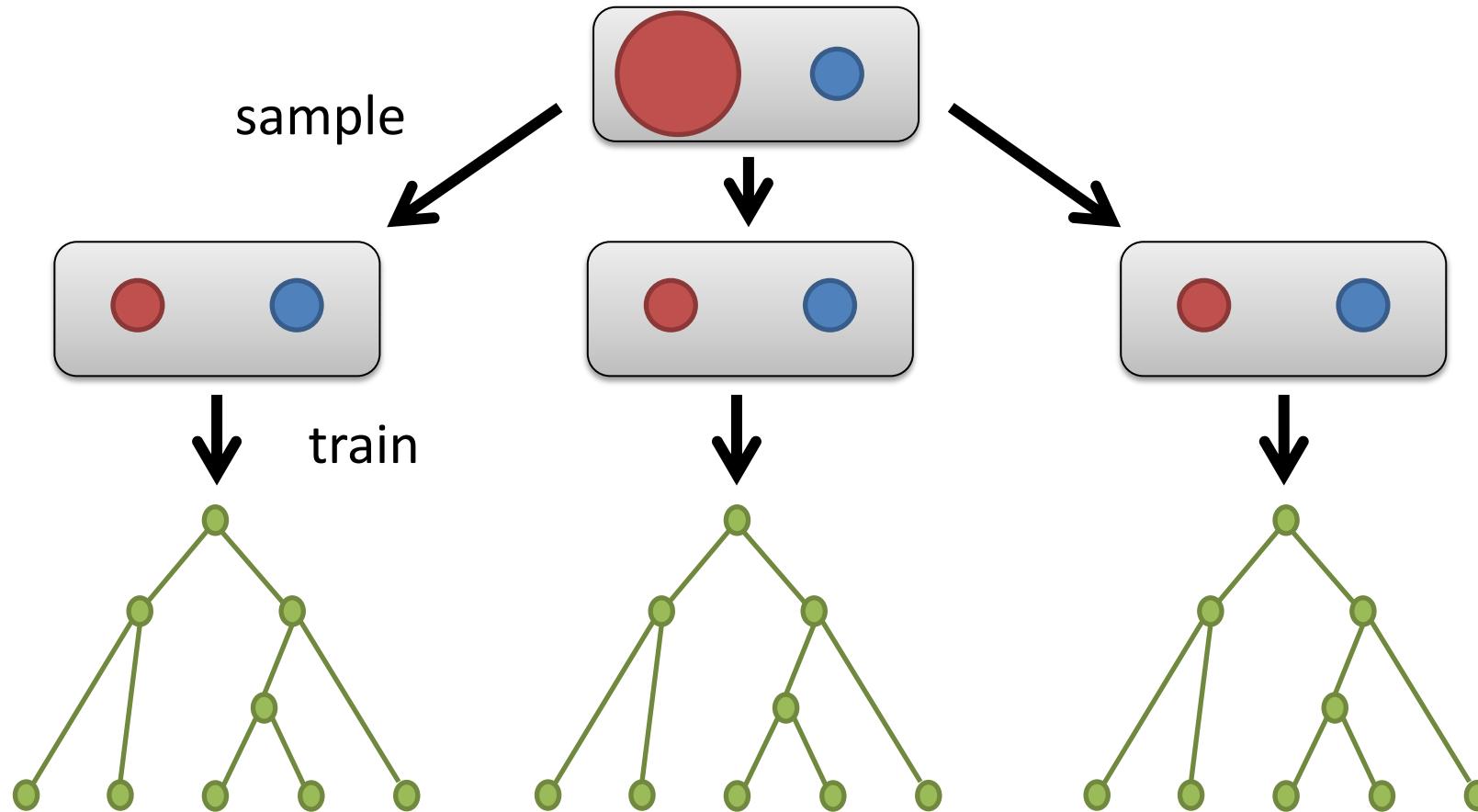


- Subsample:



- 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