

Ensemble Methods (Part I) Bagging and Random Forests

Schwartz

August 1, 2017

Objectives

- ▶ What is *Bagging*?

Objectives

- ▶ What is *Bagging*?
 - ▶ How do you do it?

Objectives

- ▶ What is *Bagging*?
 - ▶ How do you do it?
 - ▶ Why is it good?

Objectives

- ▶ What is *Bagging*?
 - ▶ How do you do it?
 - ▶ Why is it good?
- ▶ What are *Random Forests*?

Objectives

- ▶ What is *Bagging*?
 - ▶ How do you do it?
 - ▶ Why is it good?
- ▶ What are *Random Forests*?
 - ▶ How do you make them?

Objectives

- ▶ What is *Bagging*?
 - ▶ How do you do it?
 - ▶ Why is it good?
- ▶ What are *Random Forests*?
 - ▶ How do you make them?
 - ▶ Why might they be an improvement over bagging?

Objectives

- ▶ What is *Bagging*?
 - ▶ How do you do it?
 - ▶ Why is it good?
- ▶ What are *Random Forests*?
 - ▶ How do you make them?
 - ▶ Why might they be an improvement over bagging?
- ▶ What is “out of bag error” (OOB) relative to K -folds?

Objectives

- ▶ What is *Bagging*?
 - ▶ How do you do it?
 - ▶ Why is it good?
- ▶ What are *Random Forests*?
 - ▶ How do you make them?
 - ▶ Why might they be an improvement over bagging?
- ▶ What is “out of bag error” (OOB) relative to K -folds?
- ▶ How can one interpret variables in tree ensembles?

Single Predictors/Estimators

- ▶ A single tree is pretty great!

Single Predictors/Estimators

- ▶ A single tree is pretty great!
- ▶ Well, it's not that good, actually

Single Predictors/Estimators

- ▶ A single tree is pretty great!
- ▶ Well, it's not that good, actually
What could be wrong with it?

Single Predictors/Estimators

- ▶ A single tree is pretty great!
- ▶ Well, it's not that good, actually
What could be wrong with it?

Biased?

High variance?

Single Predictors/Estimators

- ▶ A single tree is pretty great!
- ▶ Well, it's not that good, actually

What could be wrong with it?

Biased?

High variance?

How so with trees?

Single Predictors/Estimators

- ▶ A single tree is pretty great!
- ▶ Well, it's not that good, actually
What could be wrong with it?

Biased?

High variance?

How so with trees?

Do we think a prediction from a single tree is good?

Ensemble Challenge

Suppose you are trying to predict the election...

You have 5 expert opinions, each with a 70% chance of being right, and each expert pick is *independent* of the other expert picks.

How could you leverage expert picks to improve accuracy and how often would you be right?

Ensemble Challenge

Suppose you are trying to predict the election...

You have 5 expert opinions, each with a 70% chance of being right, and each expert pick is *independent* of the other expert picks.

How could you leverage expert picks to improve accuracy and how often would you be right?

How many independent predictors would you need to have a 99.9% of being correct?

Ensemble Challenge

Suppose you are trying to predict the election...

You have 5 expert opinions, each with a 70% chance of being right, and each expert pick is *independent* of the other expert picks.

How could you leverage expert picks to improve accuracy and how often would you be right?

How many independent predictors would you need to have a 99.9% of being correct?

```
1 - stats.binom.cdf(26,n=53,p=.7) < .5  
< 1 - stats.binom.cdf(27,n=55,p=.7)
```

Averaging Multiple Predictors/Estimators

- ▶ Let $\hat{f}^{(j)}(\mathbf{x}_0)$ be the j^{th} predictor of outcome $f(\mathbf{x}_0)$ constructed from a training set with features $\mathbf{x}^{(j)}$ and outcomes $\mathbf{Y}^{(j)}$
- ▶ If each $\hat{f}^{(j)}(\mathbf{x}_0)$ is an unbiased predictor of $f(\mathbf{x}_0)$ then

$$E \left[\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0) \right] = f(\mathbf{x}_0)$$

Averaging Multiple Predictors/Estimators

- ▶ Let $\hat{f}^{(j)}(\mathbf{x}_0)$ be the j^{th} predictor of outcome $f(\mathbf{x}_0)$ constructed from a training set with features $\mathbf{x}^{(j)}$ and outcomes $\mathbf{Y}^{(j)}$
- ▶ If each $\hat{f}^{(j)}(\mathbf{x}_0)$ is an unbiased predictor of $f(\mathbf{x}_0)$ then

$$E \left[\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0) \right] = f(\mathbf{x}_0)$$

$\hat{f}^{(j)}(\mathbf{x}_0)$ can be an unbiased predictor by being highly complex

Averaging Multiple Predictors/Estimators

- ▶ Let $\hat{f}^{(j)}(\mathbf{x}_0)$ be the j^{th} predictor of outcome $f(\mathbf{x}_0)$ constructed from a training set with features $\mathbf{x}^{(j)}$ and outcomes $\mathbf{Y}^{(j)}$
- ▶ If each $\hat{f}^{(j)}(\mathbf{x}_0)$ is an unbiased predictor of $f(\mathbf{x}_0)$ then

$$E \left[\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0) \right] = f(\mathbf{x}_0)$$

$\hat{f}^{(j)}(\mathbf{x}_0)$ can be an unbiased predictor by being highly complex
E.g., like an overfit tree!

Averaging Multiple Predictors/Estimators

- ▶ Let $\hat{f}^{(j)}(\mathbf{x}_0)$ be the j^{th} predictor of outcome $f(\mathbf{x}_0)$ constructed from a training set with features $\mathbf{x}^{(j)}$ and outcomes $\mathbf{Y}^{(j)}$
- ▶ If each $\hat{f}^{(j)}(\mathbf{x}_0)$ is an unbiased predictor of $f(\mathbf{x}_0)$ then

$$E \left[\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0) \right] = f(\mathbf{x}_0)$$

$\hat{f}^{(j)}(\mathbf{x}_0)$ can be an unbiased predictor by being highly complex
E.g., like an overfit tree!

- ▶ And if $\text{Var} \left[\hat{f}^{(j)}(\mathbf{x}_0) \right] = \sigma^2$, then $\text{Var} \left[\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0) \right] = ?$

Averaging Multiple Predictors/Estimators

- ▶ Let $\hat{f}^{(j)}(\mathbf{x}_0)$ be the j^{th} predictor of outcome $f(\mathbf{x}_0)$ constructed from a training set with features $\mathbf{x}^{(j)}$ and outcomes $\mathbf{Y}^{(j)}$
- ▶ If each $\hat{f}^{(j)}(\mathbf{x}_0)$ is an unbiased predictor of $f(\mathbf{x}_0)$ then

$$E \left[\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0) \right] = f(\mathbf{x}_0)$$

$\hat{f}^{(j)}(\mathbf{x}_0)$ can be an unbiased predictor by being highly complex
E.g., like an overfit tree!

- ▶ And if $\text{Var} \left[\hat{f}^{(j)}(\mathbf{x}_0) \right] = \sigma^2$, then $\text{Var} \left[\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0) \right] = \frac{\sigma^2}{m}$?

Averaging Multiple Predictors/Estimators

- ▶ Let $\hat{f}^{(j)}(\mathbf{x}_0)$ be the j^{th} predictor of outcome $f(\mathbf{x}_0)$ constructed from a training set with features $\mathbf{x}^{(j)}$ and outcomes $\mathbf{Y}^{(j)}$
- ▶ If each $\hat{f}^{(j)}(\mathbf{x}_0)$ is an unbiased predictor of $f(\mathbf{x}_0)$ then

$$\mathbb{E} \left[\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0) \right] = f(\mathbf{x}_0)$$

$\hat{f}^{(j)}(\mathbf{x}_0)$ can be an unbiased predictor by being highly complex
E.g., like an overfit tree!

- ▶ And if $\text{Var} \left[\hat{f}^{(j)}(\mathbf{x}_0) \right] = \sigma^2$, then $\text{Var} \left[\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0) \right] = \frac{\sigma^2}{m}$?

So even if $\hat{f}^{(j)}$ is a high variance predictor

Averaging Multiple Predictors/Estimators

- ▶ Let $\hat{f}^{(j)}(\mathbf{x}_0)$ be the j^{th} predictor of outcome $f(\mathbf{x}_0)$ constructed from a training set with features $\mathbf{x}^{(j)}$ and outcomes $\mathbf{Y}^{(j)}$
- ▶ If each $\hat{f}^{(j)}(\mathbf{x}_0)$ is an unbiased predictor of $f(\mathbf{x}_0)$ then

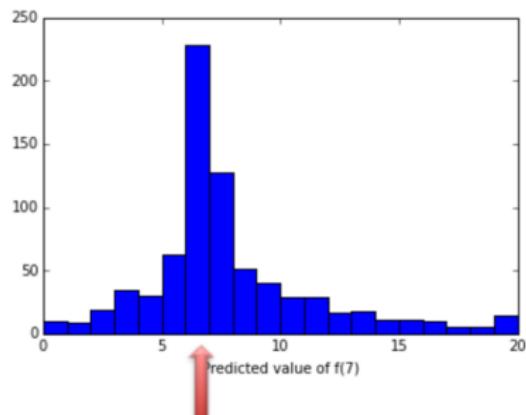
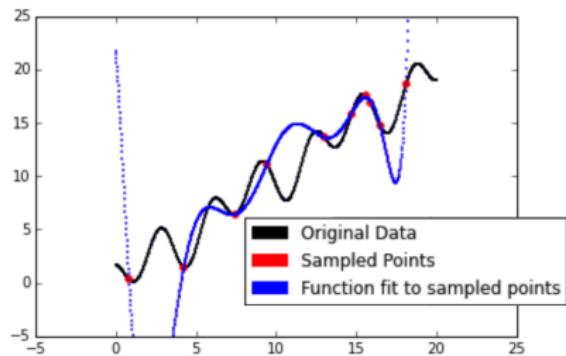
$$E \left[\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0) \right] = f(\mathbf{x}_0)$$

$\hat{f}^{(j)}(\mathbf{x}_0)$ can be an unbiased predictor by being highly complex
E.g., like an overfit tree!

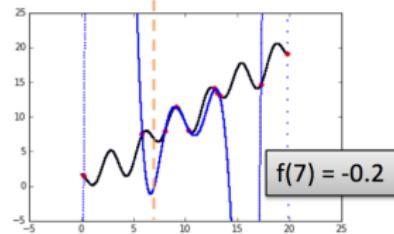
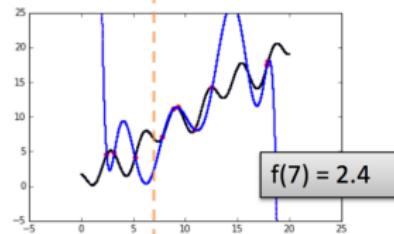
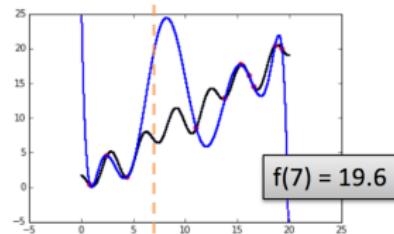
- ▶ And if $\text{Var} \left[\hat{f}^{(j)}(\mathbf{x}_0) \right] = \sigma^2$, then $\text{Var} \left[\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0) \right] = \frac{\sigma^2}{m}$?

So even if $\hat{f}^{(j)}$ is a high variance predictor the variance of the averaged predictor $\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}$ decreases with m

What is happening?



Actual value: $f(7) = 6.8$



Bootstrapping Aggregation (Bagging)

- ▶ Let $\hat{f}^{(j)}(\mathbf{x}_0)$ be a predictor of outcome $f(\mathbf{x}_0)$ constructed from a training set with features $\mathbf{x}^{(j)}$ and outcomes $\mathbf{Y}^{(j)}$

Bootstrapping Aggregation (Bagging)

- ▶ Let $\hat{f}^{(j)}(\mathbf{x}_0)$ be a predictor of outcome $f(\mathbf{x}_0)$ constructed from a training set with features $\mathbf{x}^{(j)}$ and outcomes $\mathbf{Y}^{(j)}$
- ▶ *Bagging* is simply predicting $f(\mathbf{x}_0)$ with
$$\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0)$$

an average of predictors made from bootstrapped samples

Bootstrapping Aggregation (Bagging)

- ▶ Let $\hat{f}^{(j)}(\mathbf{x}_0)$ be a predictor of outcome $f(\mathbf{x}_0)$ constructed from a training set with features $\mathbf{x}^{(j)}$ and outcomes $\mathbf{Y}^{(j)}$
- ▶ *Bagging* is simply predicting $f(\mathbf{x}_0)$ with
$$\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0)$$

an average of predictors made from bootstrapped samples
- ▶ $\mathbf{x}^{(j)}$ and $\mathbf{Y}^{(j)}$ is a *bootstrapped* sample from \mathbf{X} and \mathbf{Y}
- ▶ $\hat{f}^{(j)}$ is the model fit on that bootstrapped training set

Bootstrapping Aggregation (Bagging)

- ▶ Let $\hat{f}^{(j)}(\mathbf{x}_0)$ be a predictor of outcome $f(\mathbf{x}_0)$ constructed from a training set with features $\mathbf{x}^{(j)}$ and outcomes $\mathbf{Y}^{(j)}$
 - ▶ *Bagging* is simply predicting $f(\mathbf{x}_0)$ with
$$\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0)$$

an average of predictors made from bootstrapped samples
 - ▶ $\mathbf{x}^{(j)}$ and $\mathbf{Y}^{(j)}$ is a *bootstrapped* sample from \mathbf{X} and \mathbf{Y}
 - ▶ $\hat{f}^{(j)}$ is the model fit on that bootstrapped training set
1. Bootstrapped trees provides unbiased, high variance predictors*

Bootstrapping Aggregation (Bagging)

- ▶ Let $\hat{f}^{(j)}(\mathbf{x}_0)$ be a predictor of outcome $f(\mathbf{x}_0)$ constructed from a training set with features $\mathbf{x}^{(j)}$ and outcomes $\mathbf{Y}^{(j)}$
- ▶ *Bagging* is simply predicting $f(\mathbf{x}_0)$ with
$$\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0)$$

an average of predictors made from bootstrapped samples
- ▶ $\mathbf{x}^{(j)}$ and $\mathbf{Y}^{(j)}$ is a *bootstrapped* sample from \mathbf{X} and \mathbf{Y}
- ▶ $\hat{f}^{(j)}$ is the model fit on that bootstrapped training set

1. Bootstrapped trees provides unbiased, high variance predictors*
2. Averaged estimators are lower variance than single predictors

Bootstrapping Aggregation (Bagging)

- ▶ Let $\hat{f}^{(j)}(\mathbf{x}_0)$ be a predictor of outcome $f(\mathbf{x}_0)$ constructed from a training set with features $\mathbf{x}^{(j)}$ and outcomes $\mathbf{Y}^{(j)}$
- ▶ *Bagging* is simply predicting $f(\mathbf{x}_0)$ with
$$\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0)$$

an average of predictors made from bootstrapped samples
- ▶ $\mathbf{x}^{(j)}$ and $\mathbf{Y}^{(j)}$ is a *bootstrapped* sample from \mathbf{X} and \mathbf{Y}
- ▶ $\hat{f}^{(j)}$ is the model fit on that bootstrapped training set

1. Bootstrapped trees provides unbiased, high variance predictors*
2. Averaged estimators are lower variance than single predictors
3. Averaged predictors are still unbiased

Bootstrapping Aggregation (Bagging)

- ▶ Let $\hat{f}^{(j)}(\mathbf{x}_0)$ be a predictor of outcome $f(\mathbf{x}_0)$ constructed from a training set with features $\mathbf{x}^{(j)}$ and outcomes $\mathbf{Y}^{(j)}$
- ▶ *Bagging* is simply predicting $f(\mathbf{x}_0)$ with $\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0)$
an average of predictors made from bootstrapped samples
- ▶ $\mathbf{x}^{(j)}$ and $\mathbf{Y}^{(j)}$ is a *bootstrapped sample* from \mathbf{X} and \mathbf{Y}
- ▶ $\hat{f}^{(j)}$ is the model fit on that bootstrapped training set

1. Bootstrapped trees provides unbiased, high variance predictors*
2. Averaged estimators are lower variance than single predictors
3. Averaged predictors are still unbiased
4. Bootstrapping explores the “data set” / “predictors” space: then uses the average of predictors we might see under sampling

Bootstrapping Aggregation (Bagging)

- ▶ Let $\hat{f}^{(j)}(\mathbf{x}_0)$ be a predictor of outcome $f(\mathbf{x}_0)$ constructed from a training set with features $\mathbf{x}^{(j)}$ and outcomes $\mathbf{Y}^{(j)}$
- ▶ *Bagging* is simply predicting $f(\mathbf{x}_0)$ with $\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0)$
an average of predictors made from bootstrapped samples
- ▶ $\mathbf{x}^{(j)}$ and $\mathbf{Y}^{(j)}$ is a *bootstrapped sample* from \mathbf{X} and \mathbf{Y}
- ▶ $\hat{f}^{(j)}$ is the model fit on that bootstrapped training set

1. Bootstrapped trees provides unbiased, high variance predictors*
2. Averaged estimators are lower variance than single predictors
3. Averaged predictors are still unbiased
4. Bootstrapping explores the “data set” / “predictors” space: then uses the average of predictors we might see under sampling

* *Bagging is a general ensemble procedure often used with trees*

“Out of Bag” (OOB) testing error

- ▶ Let $\mathbf{x}_0^{(j)}$ and $\mathbf{Y}_0^{(j)}$ be the features and outcomes not in $\mathbf{x}^{(j)}$

“Out of Bag” (OOB) testing error

- ▶ Let $\mathbf{x}_0^{(j)}$ and $\mathbf{Y}_0^{(j)}$ be the features and outcomes not in $\mathbf{x}^{(j)}$
- ▶ OOB prediction error is calculated from

$$\mathbf{Y}_0^{(j)} \text{ vs } \hat{f}^{(j)}(\mathbf{x}_0^{(j)})$$

i.e., actual $\mathbf{Y}_0^{(j)}$ versus prediction $\hat{\mathbf{Y}}_0^{(j)}$ from $\hat{f}^{(j)}(\mathbf{x}_0^{(j)})$

“Out of Bag” (OOB) testing error

- ▶ Let $\mathbf{x}_0^{(j)}$ and $\mathbf{Y}_0^{(j)}$ be the features and outcomes not in $\mathbf{x}^{(j)}$
- ▶ OOB prediction error is calculated from

$$\mathbf{Y}_0^{(j)} \text{ vs } \hat{f}^{(j)}(\mathbf{x}_0^{(j)})$$

i.e., actual $\mathbf{Y}_0^{(j)}$ versus prediction $\hat{\mathbf{Y}}_0^{(j)}$ from $\hat{f}^{(j)}(\mathbf{x}_0^{(j)})$

(j)	Y_1	Y_2	Y_3	Y_4	Y_5	Y_6	Y_7	Y_8	Y_9	\dots
(1)	■	■	■	□	□	■	■	□	■	
(2)	□	□	■	■	■	□	■	■	■	
(3)	■	□	■	■	■	□	■	□	■	
:	Bootstrapped samples of size n leave out $\sim \frac{1}{3}$ of the data									

“Out of Bag” (OOB) testing error

- ▶ Let $\mathbf{x}_0^{(j)}$ and $\mathbf{Y}_0^{(j)}$ be the features and outcomes not in $\mathbf{x}^{(j)}$
- ▶ OOB prediction error is calculated from

$$\mathbf{Y}_0^{(j)} \text{ vs } \hat{f}^{(j)}(\mathbf{x}_0^{(j)})$$

i.e., actual $\mathbf{Y}_0^{(j)}$ versus prediction $\hat{\mathbf{Y}}_0^{(j)}$ from $\hat{f}^{(j)}(\mathbf{x}_0^{(j)})$

(j)	Y_1	Y_2	Y_3	Y_4	Y_5	Y_6	Y_7	Y_8	Y_9	\dots
(1)	■	■	■	□	□	■	■	□	■	
(2)	□	□	■	■	■	□	■	■	■	
(3)	■	□	■	■	■	□	■	□	■	
:	Bootstrapped samples of size n leave out $\sim \frac{1}{3}$ of the data									

- ▶ When the number of bootstrapped samples $j = 1, \dots, m$ is large, this closely approximates LOO test error estimation.

“Out of Bag” (OOB) testing error

- ▶ Let $\mathbf{x}_0^{(j)}$ and $\mathbf{Y}_0^{(j)}$ be the features and outcomes not in $\mathbf{x}^{(j)}$
- ▶ OOB prediction error is calculated from

$$\mathbf{Y}_0^{(j)} \text{ vs } \hat{f}^{(j)}\left(\mathbf{x}_0^{(j)}\right)$$

i.e., actual $\mathbf{Y}_0^{(j)}$ versus prediction $\hat{\mathbf{Y}}_0^{(j)}$ from $\hat{f}^{(j)}\left(\mathbf{x}_0^{(j)}\right)$

(j)	Y_1	Y_2	Y_3	Y_4	Y_5	Y_6	Y_7	Y_8	Y_9	...
(1)	■	■	■	□	□	■	■	□	■	
(2)	□	□	■	■	■	□	■	■	■	
(3)	■	□	■	■	■	□	■	□	■	
:	Bootstrapped samples of size n leave out $\sim \frac{1}{3}$ of the data									

- ▶ When the number of bootstrapped samples $j = 1, \dots, m$ is large, this closely approximates LOO test error estimation.
- ▶ Parameter tuning can be done “for free” when model fitting

But wait...

$$\text{Var} \left[\sum_{j=1}^m \frac{1}{m} \hat{f}^{(j)}(\mathbf{x}_0) \right] \neq \frac{\sigma^2}{m}$$

But wait...

$$\text{Var} \left[\sum_{j=1}^m \frac{1}{m} \hat{f}^{(j)}(\mathbf{x}_0) \right] \neq \frac{\sigma^2}{m}$$

- ▶ Why?

But wait...

$$\text{Var} \left[\sum_{j=1}^m \frac{1}{m} \hat{f}^{(j)}(\mathbf{x}_0) \right] \neq \frac{\sigma^2}{m}$$

- ▶ Why?

$$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}[X, Y]$$

i.e.,

$$\begin{aligned} \text{Var} \left[\hat{f}^{(j)}(\mathbf{x}_0) + \hat{f}^{(k)}(\mathbf{x}_0) \right] &= \text{Var} \left[\hat{f}^{(j)}(\mathbf{x}_0) \right] + \text{Var} \left[\hat{f}^{(k)}(\mathbf{x}_0) \right] \\ &\quad + 2\text{Cov} \left[\hat{f}^{(j)}(\mathbf{x}_0), \hat{f}^{(k)}(\mathbf{x}_0) \right] \end{aligned}$$

But wait...

$$\text{Var} \left[\sum_{j=1}^m \frac{1}{m} \hat{f}^{(j)}(\mathbf{x}_0) \right] \neq \frac{\sigma^2}{m}$$

- ▶ Why?

$$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}[X, Y]$$

i.e.,

$$\begin{aligned} \text{Var} \left[\hat{f}^{(j)}(\mathbf{x}_0) + \hat{f}^{(k)}(\mathbf{x}_0) \right] &= \text{Var} \left[\hat{f}^{(j)}(\mathbf{x}_0) \right] + \text{Var} \left[\hat{f}^{(k)}(\mathbf{x}_0) \right] \\ &\quad + 2\text{Cov} \left[\hat{f}^{(j)}(\mathbf{x}_0), \hat{f}^{(k)}(\mathbf{x}_0) \right] \end{aligned}$$

$\hat{f}^{(j)}$ and $\hat{f}^{(k)}$ are correlated: they're predicting the same thing from bootstrapped samples $\mathbf{x}^{(j)}/\mathbf{Y}^{(j)}$ & $\mathbf{x}^{(k)}/\mathbf{Y}^{(k)}$ from \mathbf{x}/\mathbf{Y} !

How good is bagging?

- ▶ If $\text{Cor}[\hat{f}^{(j)}, \hat{f}^{(k)}] = \rho$ for all j and k

How good is bagging?

- ▶ If $\text{Cor}\left[\hat{f}^{(j)}, \hat{f}^{(k)}\right] = \rho$ for all j and k
then

$$\text{Var} \left[\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0) \right] = \rho\sigma^2 + (1 - \rho)\frac{\sigma^2}{m}$$

How good is bagging?

- ▶ If $\text{Cor}[\hat{f}^{(j)}, \hat{f}^{(k)}] = \rho$ for all j and k

then

$$\text{Var} \left[\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0) \right] = \rho\sigma^2 + (1 - \rho)\frac{\sigma^2}{m}$$

- ▶ Only “uncorrelated parts” of $\hat{f}^{(j)}$ and $\hat{f}^{(k)}$ get “CLT effect”

How good is bagging?

- ▶ If $\text{Cor}[\hat{f}^{(j)}, \hat{f}^{(k)}] = \rho$ for all j and k

then

$$\text{Var} \left[\frac{1}{m} \sum_{j=1}^m \hat{f}^{(j)}(\mathbf{x}_0) \right] = \rho\sigma^2 + (1 - \rho)\frac{\sigma^2}{m}$$

- ▶ Only “uncorrelated parts” of $\hat{f}^{(j)}$ and $\hat{f}^{(k)}$ get “CLT effect”
- ▶ So is there any way to get ρ close to 0?

Decorrelating trees

We might try to decorrelate any class of predictors $\hat{f}^{(j)}$

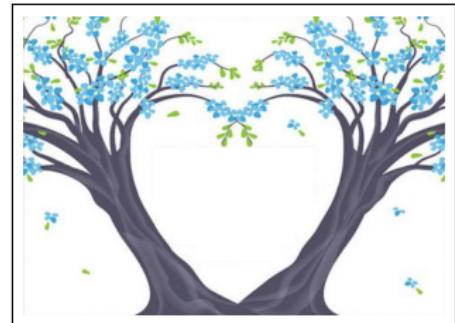
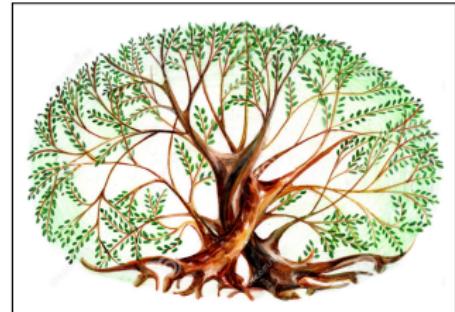
But let's focus on



Decorrelating trees

We might try to decorrelate any class of predictors $\hat{f}^{(j)}$

But let's focus on



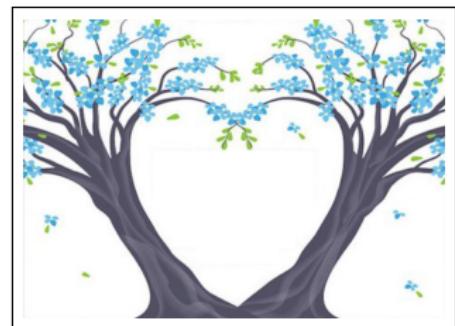
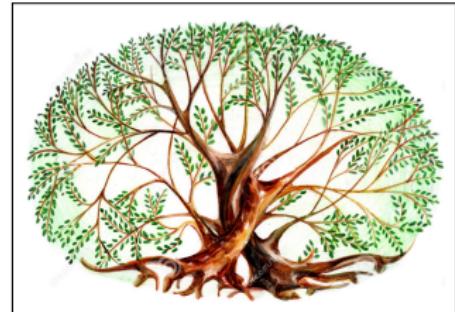
Decorrelating trees

We might try to decorrelate any class of predictors $\hat{f}^{(j)}$

But let's focus on



Trees are “correlated” because



Decorrelating trees

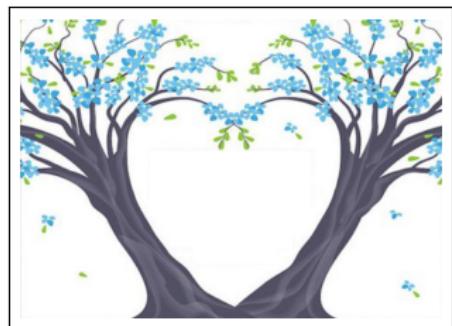
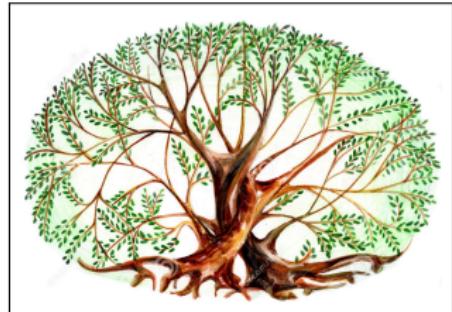
We might try to decorrelate any class of predictors $\hat{f}^{(j)}$

But let's focus on



Trees are “correlated” because

1. bootstrap samples are about the same [nothing to do there]



Decorrelating trees

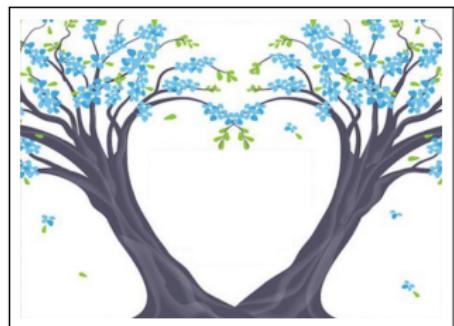
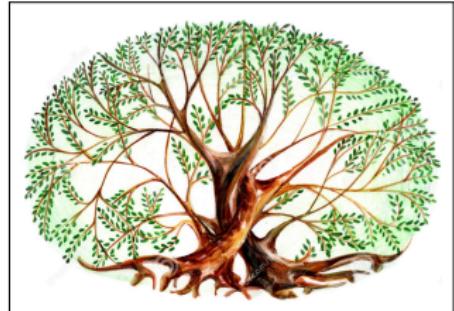
We might try to decorrelate any class of predictors $\hat{f}^{(j)}$

But let's focus on



Trees are “correlated” because

1. bootstrap samples are about the same [nothing to do there]
2. influential features tend to be same [this we could influence...]



Decorrelating trees

We might try to decorrelate any class of predictors $f^{(j)}$

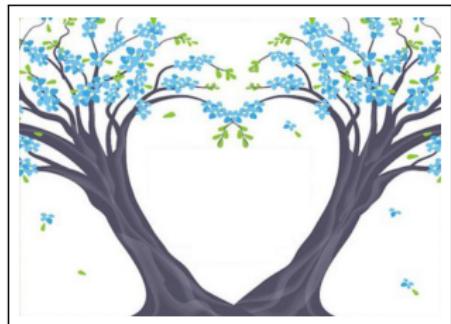
But let's focus on



Trees are “correlated” because

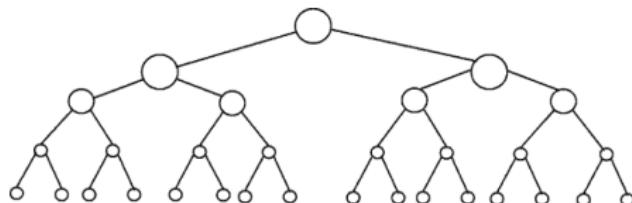
1. bootstrap samples are about the same [nothing to do there]
2. influential features tend to be same [this we could influence...]

How can we make constructed tree not exactly the same?



Random Forests

- ▶ Tree is constructed by recursive best splits



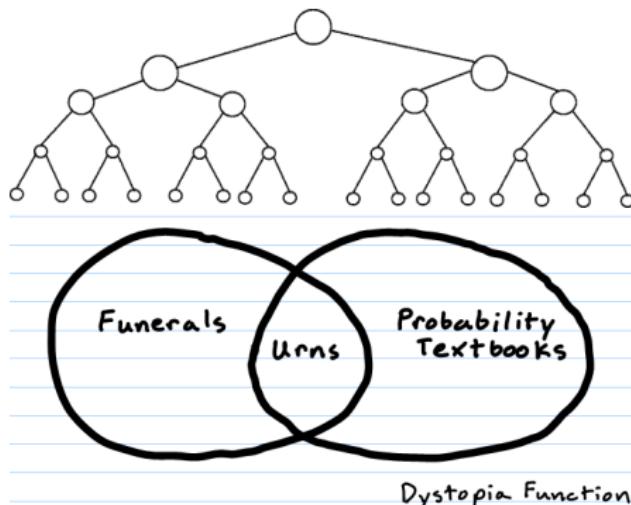
FEATURES

♣	♠	□	■
◀	▷	▽	○
◇	♡	★	△
⊗	⊕	◎	⊗

Features sampled with
replacement for possible
use **at each node**

Random Forests

- ▶ Tree is constructed by recursive best splits



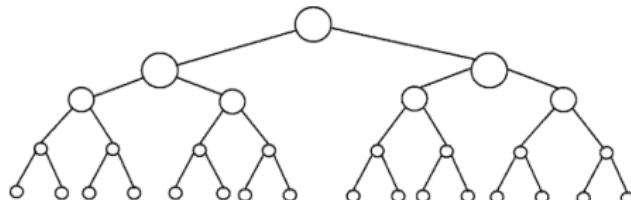
FEATURES

♣	♠	□	■
△	▷	▽	○
◊	♡	★	△
⊗	⊕	◎	⊗

Features sampled with replacement for possible use **at each node**

Random Forests

- ▶ Tree is constructed by recursive best splits
- ▶ We only use a random subset of features **at each split**



♠ □ ⊖ ⊕

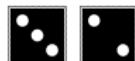
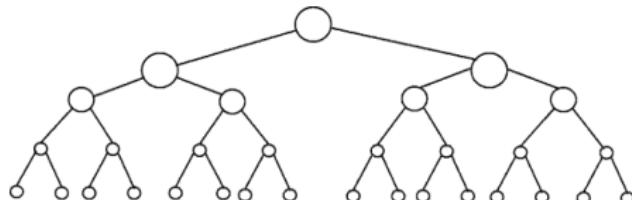
FEATURES

♣	♠	□	■
◀	▷	▽	○
◇	♡	★	△
⊖	⊕	⊙	⊗

Features sampled with
replacement for possible
use at each node

Random Forests

- ▶ Tree is constructed by recursive best splits
- ▶ We only use a random subset of features **at each split**



♣ ♦ ⊗ ⊙

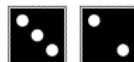
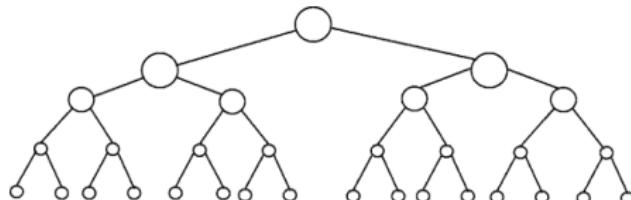
FEATURES

♣	♠	□	■
◀	▷	▽	○
◇	♡	★	△
⊗	⊕	⊙	⊗

Features sampled with replacement for possible use **at each node**

Random Forests

- ▶ Tree is constructed by recursive best splits
- ▶ We only use a random subset of features **at each split**



♣ ♦ ⊗ ⊙

FEATURES

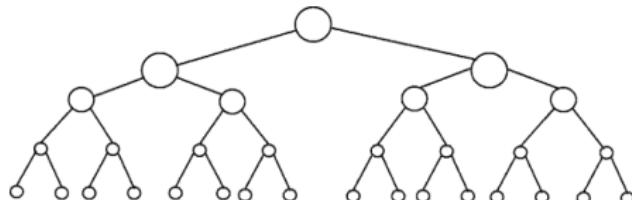
♣	♠	□	■
◀	▷	▽	○
◇	♡	★	△
⊗	⊕	⊙	⊗

Features sampled with replacement for possible use **at each node**

- ▶ Temporary exclusion – features includable again at each split

Random Forests

- ▶ Tree is constructed by recursive best splits
- ▶ We only use a random subset of features **at each split**



♣ ♦ ⊗ ⊙

FEATURES

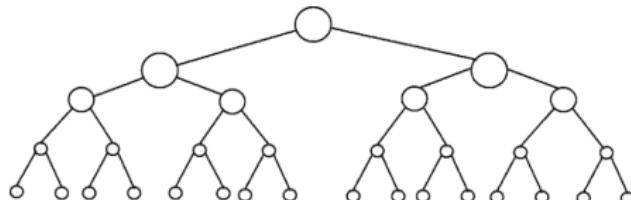
♣	♠	□	■
◀	▷	▽	○
◇	♡	★	△
⊗	⊕	⊙	⊗

Features sampled with replacement for possible use **at each node**

- ▶ Temporary exclusion – features includable again at each split
- ▶ Typically select \sqrt{p} (classification) & $\frac{p}{3}$ (regression) features

Random Forests

- ▶ Tree is constructed by recursive best splits
- ▶ We only use a random subset of features **at each split**



♣ ♦ ⊗ ⊙

FEATURES

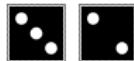
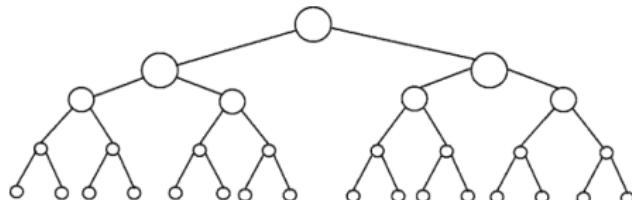
♣	♠	□	■
◀	▷	▽	○
◊	♡	★	△
⊗	⊕	⊙	⊗

Features sampled with replacement for possible use **at each node**

- ▶ Temporary exclusion – features includable again at each split
 - ▶ Typically select \sqrt{p} (classification) & $\frac{p}{3}$ (regression) features
- Number features considered is a tuning parameter for what?

Random Forests

- ▶ Tree is constructed by recursive best splits
- ▶ We only use a random subset of features **at each split**



♣ ♦ ⊗ ⊙

FEATURES

♣	♦	□	■
◀	▷	▽	○
◊	♡	★	△
⊗	⊕	⊙	⊗

Features sampled with replacement for possible use **at each node**

- ▶ Temporary exclusion – features includable again at each split
- ▶ Typically select \sqrt{p} (classification) & $\frac{p}{3}$ (regression) features
Number features considered is a tuning parameter for what?
- ▶ More features means stronger but also more correlated trees...

Random Forest *Tuning*

$g(p)$: Number of features considered at each split

Random Forest *Tuning*

$g(p)$: Number of features considered at each split

m : Number of trees (i.e., number of bootstrapped samples)

Random Forest Tuning

$g(p)$: Number of features considered at each split

m : Number of trees (i.e., number of bootstrapped samples)

n_b : Sample size of each bootstrapped sample

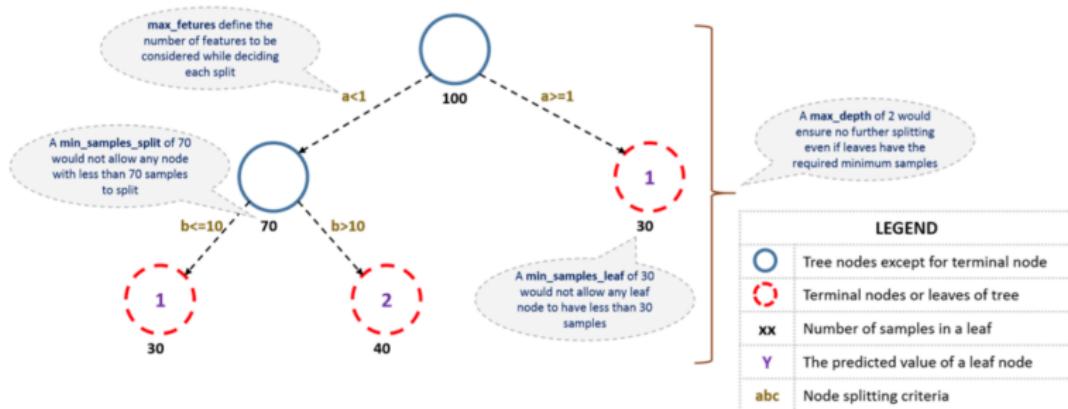
Random Forest Tuning

$g(p)$: Number of features considered at each split

m : Number of trees (i.e., number of bootstrapped samples)

n_b : Sample size of each bootstrapped sample

χ : Tree characteristics



Random Forest Tuning

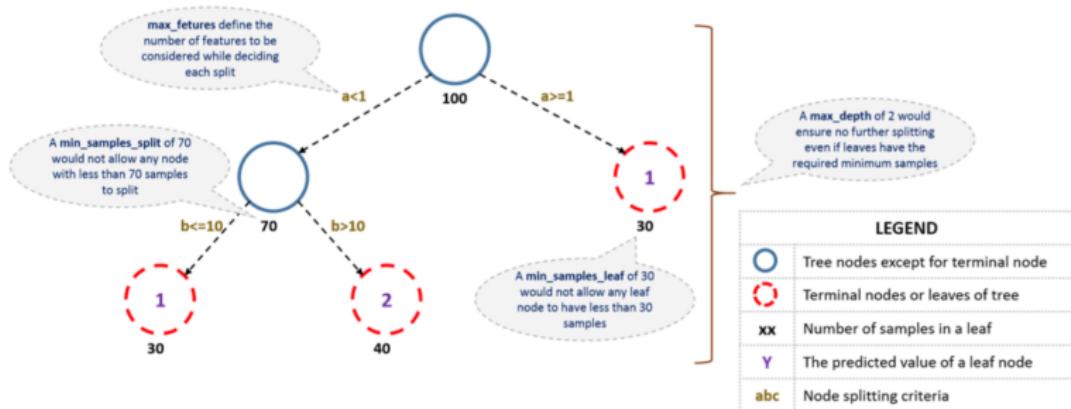
$g(p)$: Number of features considered at each split

m : Number of trees (i.e., number of bootstrapped samples)

n_b : Sample size of each bootstrapped sample

χ : Tree characteristics

$$\Rightarrow \sum(Y_i - \hat{Y}_i)^2 + \alpha|T|$$



Random Forest Tuning

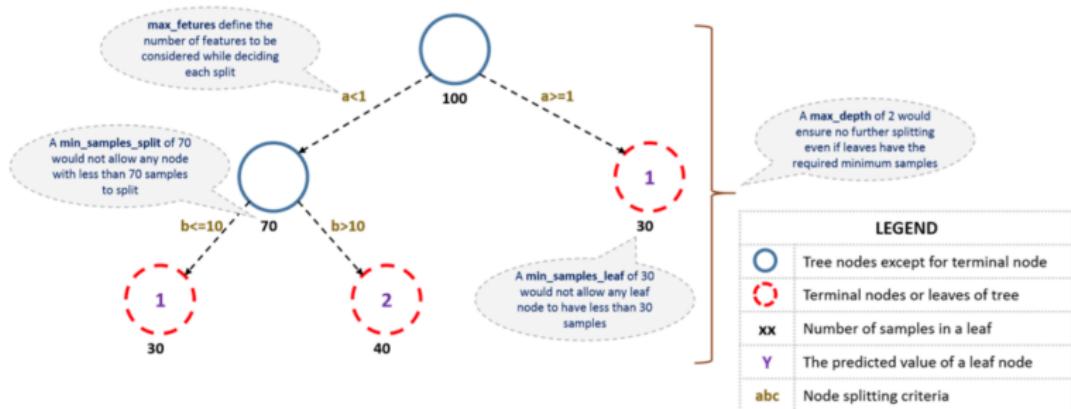
$g(p)$: Number of features considered at each split

m : Number of trees (i.e., number of bootstrapped samples)

n_b : Sample size of each bootstrapped sample

χ : Tree characteristics

$$\Rightarrow \sum(Y_i - \hat{Y}_i)^2 + \alpha|T|$$



- ▶ Random forests are often robust to tree characteristics:

Random Forest Tuning

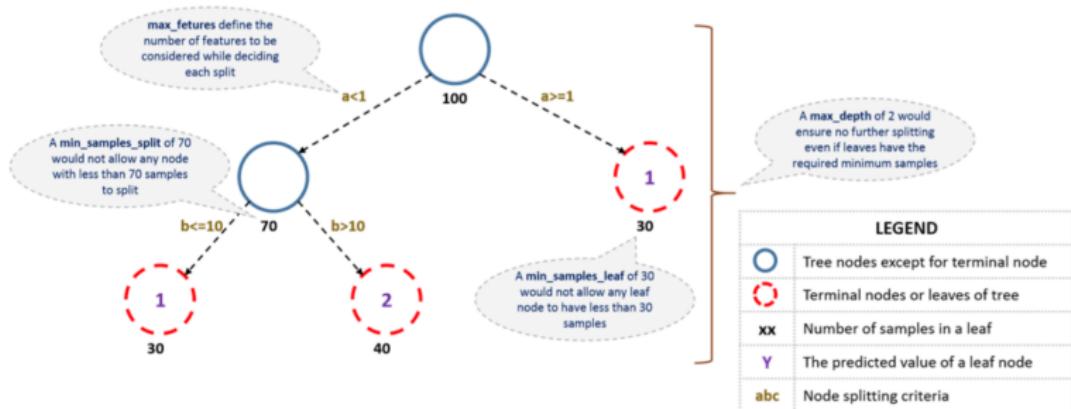
$g(p)$: Number of features considered at each split

m : Number of trees (i.e., number of bootstrapped samples)

n_b : Sample size of each bootstrapped sample

χ : Tree characteristics

$$\Rightarrow \sum(Y_i - \hat{Y}_i)^2 + \alpha|T|$$



- ▶ Random forests are often robust to tree characteristics:
very large numbers of “bushy” trees typically used &

Random Forest Tuning

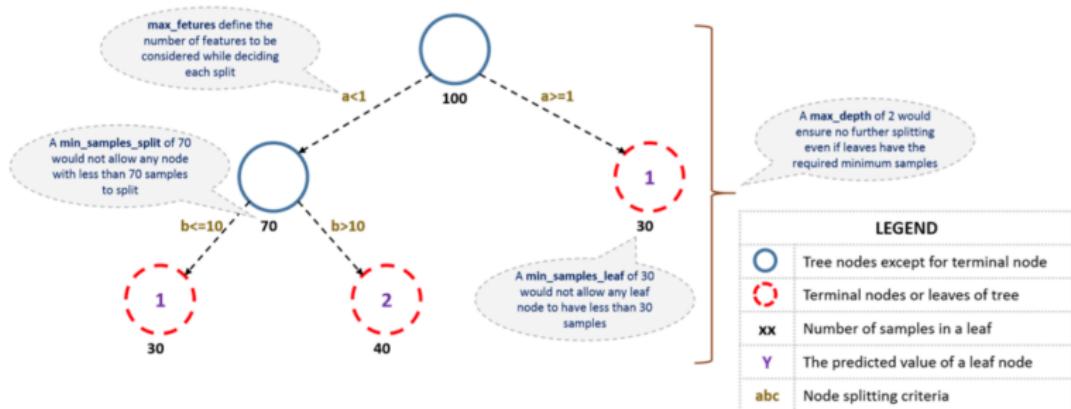
$g(p)$: Number of features considered at each split

m : Number of trees (i.e., number of bootstrapped samples)

n_b : Sample size of each bootstrapped sample

χ : Tree characteristics

$$\Rightarrow \sum(Y_i - \hat{Y}_i)^2 + \alpha|T|$$



- ▶ Random forests are often robust to tree characteristics:
very large numbers of “bushy” trees typically used & approach state of the art performance out of the box

Random Forest *Computational Notes*

- ▶ Trees can be computationally expensive to fit

Random Forest *Computational Notes*

- ▶ Trees can be computationally expensive to fit
 - ▶ But they can be fit in parallel

Random Forest *Computational Notes*

- ▶ Trees can be computationally expensive to fit
 - ▶ But they can be fit in parallel
 - ▶ At some point $\rho\sigma^2 + (1 - \rho)\frac{\sigma^2}{m}$ reaches an asymptote

Random Forest *Computational Notes*

- ▶ Trees can be computationally expensive to fit
 - ▶ But they can be fit in parallel
 - ▶ At some point $\rho\sigma^2 + (1 - \rho)\frac{\sigma^2}{m}$ reaches an asymptote
- ▶ Cross-validation can be computationally demanding

Random Forest *Computational Notes*

- ▶ Trees can be computationally expensive to fit
 - ▶ But they can be fit in parallel
 - ▶ At some point $\rho\sigma^2 + (1 - \rho)\frac{\sigma^2}{m}$ reaches an asymptote
- ▶ Cross-validation can be computationally demanding
 - ▶ Parameters can be tuned using OOB

You have built-in test samples so you don't need K-folds – you score a model right when it's fit at a specific tuning parameter!

Random Forest *Computational Notes*

- ▶ Trees can be computationally expensive to fit
 - ▶ But they can be fit in parallel
 - ▶ At some point $\rho\sigma^2 + (1 - \rho)\frac{\sigma^2}{m}$ reaches an asymptote
- ▶ Cross-validation can be computationally demanding
 - ▶ Parameters can be tuned using OOB

You have built-in test samples so you don't need K-folds – you score a model right when it's fit at a specific tuning parameter!

- ▶ Cross-validation is still needed for methodology comparisons...
- ▶ OOB samples don't appear in other model fitting contexts – so you still need a common “out of sample” set to compare on...

Interpreting Models

Change X_i and see what happens

It's just really not that hard...

Interpreting Models

Change X_i and see what happens

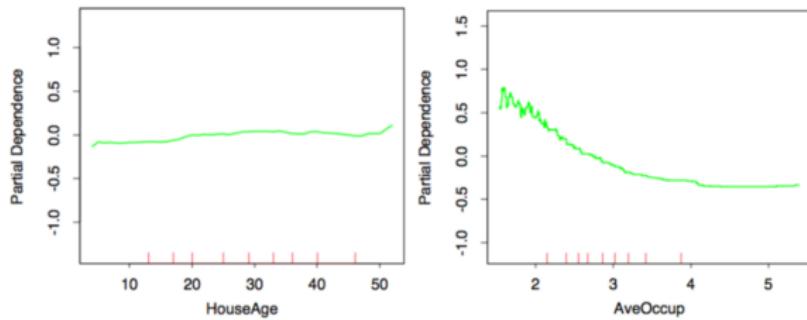
It's just really not that hard...

That's all you're doing in a linear model:

$$\beta_0 + \beta_1 X_1 + \beta_1 X_1 + \beta_2 X_2 + \dots$$

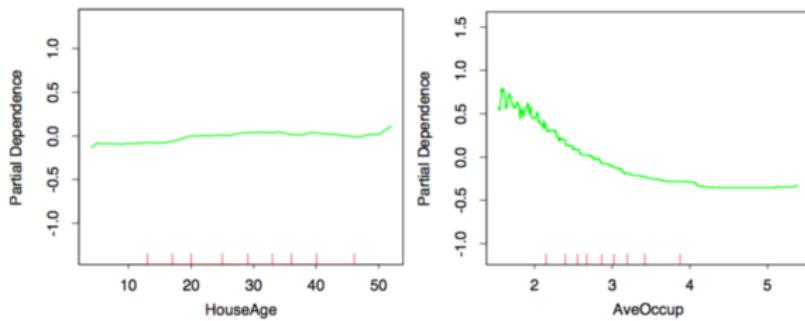
Partial Dependency Plots

1. Set feature $V = v$ for a single values v



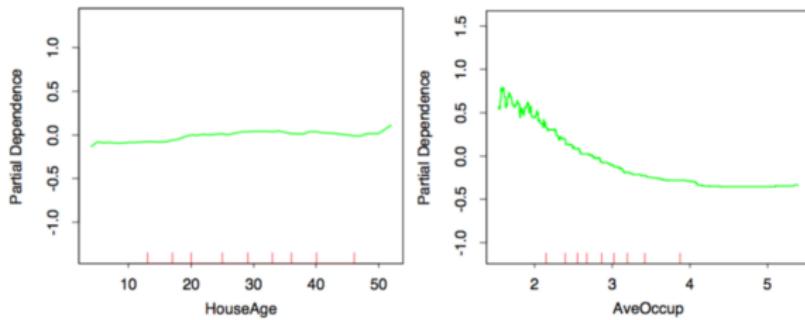
Partial Dependency Plots

1. Set feature $V = v$ for a single values v
2. Record average predictions $\frac{1}{n} \sum_{i=1}^n \hat{Y}_i^{(v)}$ over “synthetic” data



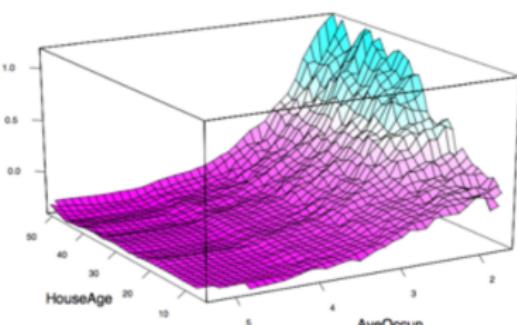
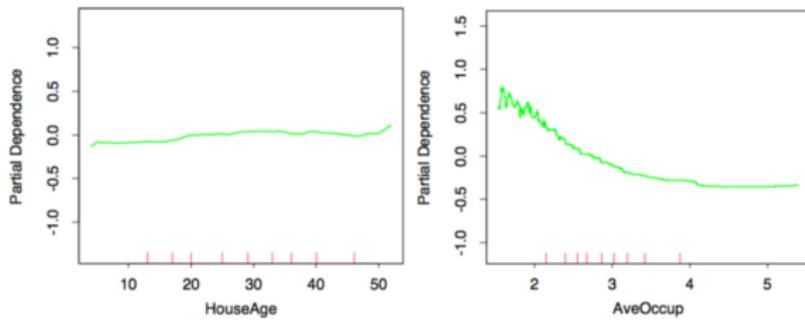
Partial Dependency Plots

1. Set feature $V = v$ for a single values v
2. Record average predictions $\frac{1}{n} \sum_{i=1}^n \hat{Y}_i^{(v)}$ over “synthetic” data
3. Plot $\frac{1}{n} \sum_{i=1}^n \hat{Y}_i^{(v)}$ at v



Partial Dependency Plots

1. Set feature $V = v$ for a single values v
2. Record average predictions $\frac{1}{n} \sum_{i=1}^n \hat{Y}_i^{(v)}$ over “synthetic” data
3. Plot $\frac{1}{n} \sum_{i=1}^n \hat{Y}_i^{(v)}$ at v



Variable Importance

But there are some other *Very Cool* model interpretation options available in tree based contexts as well...

Variable Importance (v1.0): “share of information gain”

1. Split are made on the feature minimizing RSS or Gini/Entropy

Variable Importance (v1.0): “share of information gain”

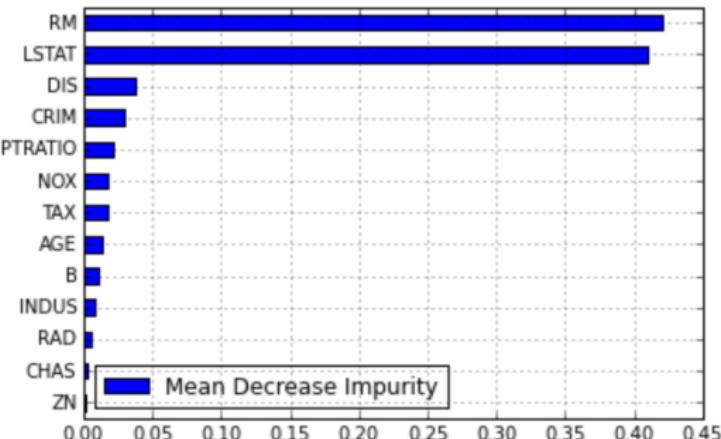
1. Split are made on the feature minimizing RSS or Gini/Entropy
2. For split s , cumulate the explanatory contribution $\xi_V^{(s)}$ for V

Variable Importance (v1.0): “share of information gain”

1. Split are made on the feature minimizing RSS or Gini/Entropy
2. For split s , cumulate the explanatory contribution $\xi_V^{(s)}$ for V
3. Average the scores $\xi_V^{(s)}$ for each feature V **across all trees**

Variable Importance (v1.0): “share of information gain”

1. Split are made on the feature minimizing RSS or Gini/Entropy
2. For split s , cumulate the explanatory contribution $\xi_V^{(s)}$ for V
3. Average the scores $\xi_V^{(s)}$ for each feature V **across all trees**



In Boston, the most relevant associations with neighborhood home values are (1) number of rooms and (2) proportion of low income households in the neighborhood

Variable Importance (v2.0): “error caused by losing V”

1. Get OOB accuracy using all features, including original V

Variable Importance (v2.0): “error caused by losing V ”

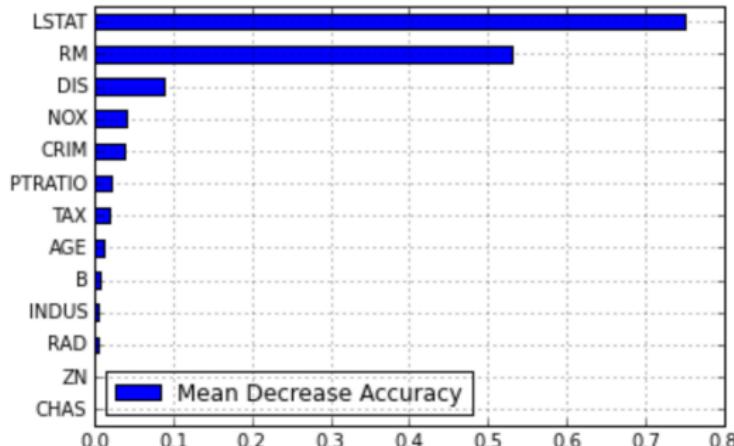
1. Get OOB accuracy using all features, including original V
2. Get OOB accuracy using all features & V^* , a permuted V

Variable Importance (v2.0): “error caused by losing V ”

1. Get OOB accuracy using all features, including original V
2. Get OOB accuracy using all features & V^* , a permuted V
3. Compare the differences between 1 and 2 for all features

Variable Importance (v2.0): “error caused by losing V ”

1. Get OOB accuracy using all features, including original V
2. Get OOB accuracy using all features & V^* , a permuted V
3. Compare the differences between 1 and 2 for all features



This approach suggests prediction is more sensitive to
(1) low income proportion rather than (2) room number

Variable Importance (v3.0): “what influences predictions”

1. Follow test samples down the tree through each split s on V

Variable Importance (v3.0): “what influences predictions”

1. Follow test samples down the tree through each split s on V
2. Record the change in prediction $\epsilon_V^{(s)}$ due to V split at s

Variable Importance (v3.0): “what influences predictions”

1. Follow test samples down the tree through each split s on V
2. Record the change in prediction $\epsilon_V^{(s)}$ due to V split at s
3. Sum up all the changes $\epsilon_V^{(s)}$ due to feature V

Variable Importance (v3.0): “what influences predictions”

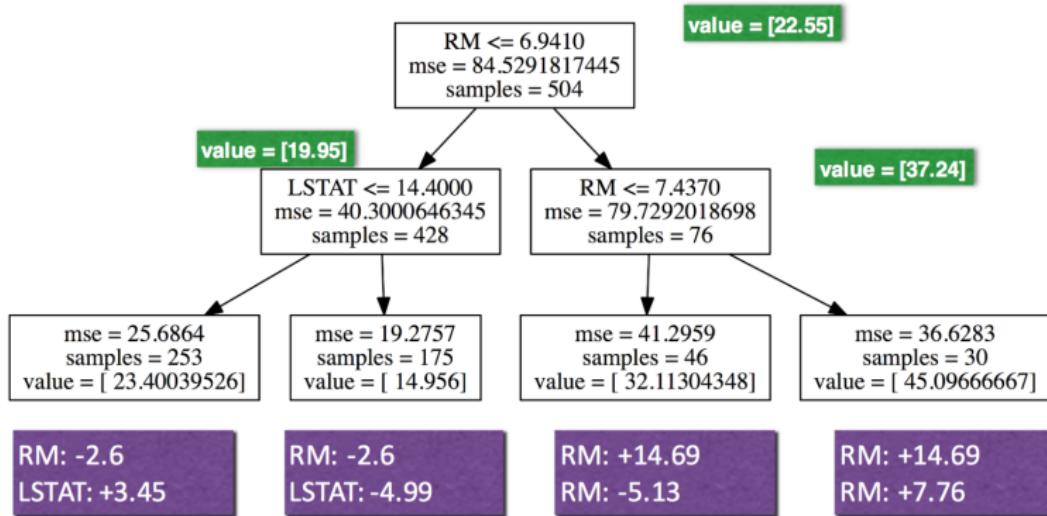
1. Follow test samples down the tree through each split s on V
2. Record the change in prediction $\epsilon_V^{(s)}$ due to V split at s
3. Sum up all the changes $\epsilon_V^{(s)}$ due to feature V

These values are then averaged across trees

Variable Importance (v3.0): “what influences predictions”

1. Follow test samples down the tree through each split s on V
2. Record the change in prediction $\epsilon_V^{(s)}$ due to V split at s
3. Sum up all the changes $\epsilon_V^{(s)}$ due to feature V

These values are then averaged across trees



Variable Importance (v3.0): “what influences predictions”

1. Follow test samples down the tree through each split s on V
2. Record the change in prediction $\epsilon_V^{(s)}$ due to V split at s
3. Sum up all the changes $\epsilon_V^{(s)}$ due to feature V

These values are then averaged by tree

	DIS	INDUS	LSTAT	NOX	PTRATIO	RAD	RM	TAX	ZN
Prediction 1	6.11	0.10	2.67	-0.02	-0.19	0.06	-2.63	-0.20	0.03
Prediction 2	6.22	0.16	2.56	-0.01	-0.19	0.06	-3.15	-0.16	0.03
Prediction 3	-0.70	0.04	7.42	-0.11	0.42	-0.02	1.10	-0.14	-0.05
Prediction 4	-0.53	0.25	3.50	0.16	1.46	0.13	2.21	-0.24	0.11
Prediction 5	-0.68	0.15	7.86	0.03	0.85	0.01	-1.14	-0.16	0.18
Prediction 6	0.18	-0.26	8.62	-0.19	-0.02	-0.07	-1.83	-0.34	-0.05

Distance from city-center Fraction of lower-class residents Average number of rooms

Features effects can be characterized on individual samples (!)

Variable Importance (v4.0): “what the tree predicts with”

1. Calculate proportion of samples visiting feature V in each tree

Variable Importance (v4.0): “what the tree predicts with”

1. Calculate proportion of samples visiting feature V in each tree
(conditions higher in the tree typically visited by more data)

Variable Importance (v4.0): “what the tree predicts with”

1. Calculate proportion of samples visiting feature V in each tree
(conditions higher in the tree typically visited by more data)
2. Average the “proportion of samples visited” across all trees

Fine

So don't tell me machine learning models, *particularly tree-based methods* are uninterpretable “Black-Boxes”. I'll kick your butt.

Fine

So don't tell me machine learning models, *particularly tree-based methods* are uninterpretable “Black-Boxes”. I'll kick your butt.

Now:

Tree < Bagging < Random Forests < _____ ?