# Random Forests

Last time: $V_1(x) =$ Win Probability if you have a
1st down and 10 out
game-state $x = ($yardline,
game seconds remaining, pointspread,
timeouts, recent 2nd half
kickoff, ___ )

We introduced Decision Trees as a machine
learning model to capture interacting relationships
between variables.

Problem: decision trees are unstable and
prone to overfitting.

Unstable: if you slightly perturb the training dataset,
you can get a very different model

Overfitting: memorizing the random noise in the
training dataset rather than the "true"
underlying trend (signal)

What to do about this? Ask Leo Breiman.

Unstable methods can have their accuracy improved by perturbing and combining, that is, generate multiple versions of the predictor by perturbing the training set or construction method, then combine these multiple versions into a single predictor. ~~For instance,~~ ~~(1995)~~

- small changes in the training set or in the construction of a decision tree can lead to larger changes in the aggregated predictor

---

## <u>Bagging</u> = Bootstrap - Aggregating

Training set $T$ consists of $N$ instance $n=1, ..., N$ (e.g. $N$ football plays)

put equal probabilities $p(n) = \frac{1}{N}$ on each instance.

Sample with Replacement (Bootstrap) $N$ times from the training set $T$, forming a "new" Resampled training set $T^{(B)}$.

Then fit a decision tree $\hat{f}^{(CB)}$ on $T^{(B)}$.

Lets do this, say, 1000 times, so we have 1000 bootstrapped decision trees $\{\hat{f}^{(1)}, ..., \hat{f}^{(1000)}\}$

Aggregate $\hat{f} = \frac{1}{1000} \sum_{B=1}^{1000} \hat{f}^{(B)}$.

# Why might bagging be a good idea?

1. First, pretend you have 1000 **independently** **drawn** datasets of size $N$, $\{T^{(1)}, \ldots, T^{(1000)}\}$.

Fit 1000 decision trees $\{\hat{f}^{(1)}, \ldots, \hat{f}^{(1000)}\}$ where $\hat{f}^{(B)}$ fit on $T^{(B)}$,

$$\hat{f} = \frac{1}{1000} \sum_{B=1}^{1000} \hat{f}^{(B)} \quad \text{and}$$

$$\mathbb{E}\hat{f} = \frac{1}{1000} \sum_{B=1}^{1000} \mathbb{E}\hat{f}^{(B)} \quad \text{by linearity of expectation}$$

$$= \frac{1}{1000} \sum_{B=1}^{1000} \mathbb{E}\hat{f}^{(1)} \quad \text{because } \mathbb{E}\hat{f}^{(B)} = \mathbb{E}\hat{f}^{(1)} \text{ for all } B$$

$$= \mathbb{E}\hat{f}^{(1)}.$$

The aggregated predictor has the same expected value as 1 individual decision tree.

$$\text{Var}(\hat{f}) = \text{Var}\left(\frac{1}{1000} \sum_{B=1}^{1000} \hat{f}^{(B)}\right)$$

$$= \frac{1}{1000^2} \sum_{B=1}^{1000} \text{Var}(\hat{f}^{(B)}) \quad \text{by independence}$$

$$= \frac{1}{1000^2} \sum_{B=1}^{1000} \sigma^2$$

$$= \frac{\sigma^2}{1000}$$

$$\mu = \mathbb{E}\hat{f}^{(B)}, \quad \sigma^2 = \text{Var}(\hat{f}^{(B)}), \quad \text{then}$$

$$\mathbb{E}\hat{f} = \mu \qquad \text{Var}(\hat{f}) = \frac{\sigma^2}{1000}.$$

Aggregating independent decision trees has same expected value but lower variance.

Assumed 1000 independently drawn training sets from some "true" function.

\* In reality we have only **one** independently drawn dataset of size N, not 1000.

Bootstrapping: created 1000 different training datasets by Resampling N Rows with Replacement.

**Idea behind Bagging** one bootstrapped dataset is a reasonably close approximation to one independently drawn dataset from the "true" data source.

Why is bootstrapping a reasonable approximation to independently drawing a new training set?

## Have

"tRue" underlying data generating pRocess

$\downarrow$

1 observed training dataset

$\downarrow$

1 decision tree $\hat{f}$

## Something better

"tRue" underlying data generating pRocess

Observed training Set 1 | Observed training set 2 | ... | observed trainy Set 1000

we do not have this

1000 independently drawn trainy sets

decision tree $\hat{f}^{(1)}$ | $\hat{f}^{(2)}$ | ... | $\hat{f}^{(1000)}$

$\downarrow$

aggregated decision tree $\hat{f}$

## Approximate by bootstrapping

"tRue" underlying data generating pRocess $\rightarrow$ 1 observed training dataset

bootstrapped training set 1 | boot training set 2 | ... | boot training set 1000

$\hat{f}^{(1)}$ | $\hat{f}^{(2)}$ | ... | $\hat{f}^{(1000)}$

$\rightarrow \hat{f} \leftarrow$

bootstrapping = Sample N rows with replacement

Why is bootstrapped $\hat{f}^{(B)}$ reasonably similar to an independently drawn training set?

"True" underlying data generating process → Observed training set

"True" underlying data generating process → Observed training dataset → bootstrapped training set
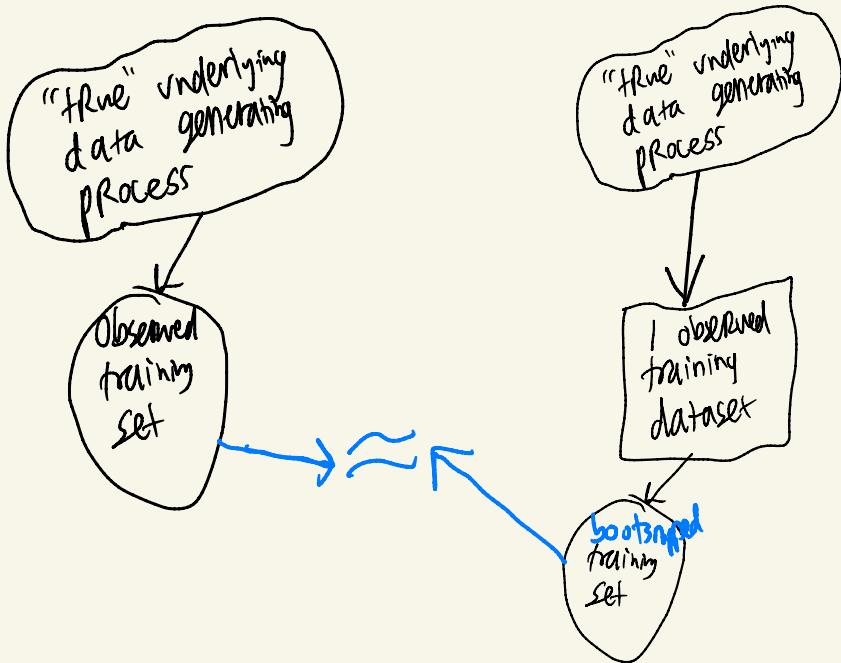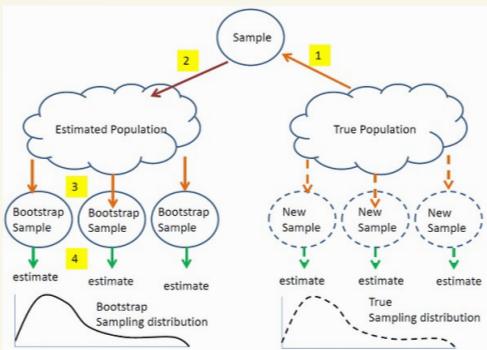
Observed training set ≈ bootstrapped training set

Sampling from the "true" underlying data generating process is similar to sampling from the observed training set because the observed training set itself is sampled from the "true" underlying data generating process.

# Random Forest

Input: observed training data $T = (X, y)$

Hyperparameters: $n, m, A = 1000$

for $B = 1, \ldots, A$:

1. Bootstrap (Row subsampling):
   Sample $n$ of $N$ rows of $T$
   with replacement

2. Column Subsampling:
   Sample $m$ of the columns (features)
   of $X$

   These 2 things yield a new training set $T^{(B)}$
   with $n$ Rows and $m$ Columns.

3. fit a decision tree $\hat{f}^{(B)}$ on $T^{(B)}$.

Aggregate: $\hat{f} = \dfrac{1}{1000} \displaystyle\sum_{B=1}^{1000} \hat{f}^{(B)}$.

\* Random forest is <u>much</u> better than a single decision tree because it <u>reduces variance.</u>

## Bias-Variance Tradeoff

out-of sample prediction error $=$ Bias $+$ Variance $+$ Irreducible error

↓↓ lower for R.F. than dec. tree

$\mathbb{E}(\hat{f}-f)^2$

Same for 1 dec. tree and 1 Random Forest

R.F. has lower variance than 1 decision tree

Same no matter what model you use

\* one decision tree $\hat{f}(T)$ is more unstable and sensitive to the Random idiosyncrasies of the traing set than 1000 aggregated decision trees $\hat{f}(T) = \frac{1}{1000} \sum_{}^{1000} \hat{f}^{(B)}(T)$.