# Randomization methods: Bagging

## -Statistical Machine Learning-

Lecturer: Darren Homrighausen, PhD

# NOTATION

REMINDER: For either classification or regression, we produce predictions for a given covariate vector $X$

That is, we form

$$\hat{Y} = \hat{f}(X) \qquad \text{or} \qquad \hat{Y} = \hat{g}(X)$$

where

- $\hat{f}$ or $\hat{g}$ is some procedure formed with the training data

  (EXAMPLES: $\hat{\beta}$ formed by least squares)

- The prediction $\hat{Y}$ formed at a desired covariate vector $X$

  (EXAMPLE: $\hat{Y} = X^\top \hat{\beta}$ formed by least squares)

# BAGGING

Many methods (trees included) tend to be designed to have lower bias but high variance

This means that if we split the training data into two parts at random and fit a decision tree to each part, the results could be quite different

A low variance estimator would yield similar results if applied repeatedly to distinct data sets
(consider $\hat{f}(X) = 0$ for all $X$)

Bagging, also known as Bootstrap AGgregation, is a general purpose procedure for reducing variance.

We'll use it specifically in the context of trees, but it can be applied more broadly.

# Bagging: The main idea

Suppose we have $n$ uncorrelated observations $Z_1, \ldots, Z_n$, each with variance $\sigma^2$.

What is the variance of

$$\overline{Z} = \frac{1}{n} \sum_{i=1}^{n} Z_i?$$

# Bagging: The main idea

Suppose we have $n$ uncorrelated observations $Z_1, \ldots, Z_n$, each with variance $\sigma^2$.

What is the variance of

$$\overline{Z} = \frac{1}{n} \sum_{i=1}^{n} Z_i?$$

Answer: $\sigma^2/n$.

More generally, if we have $B$ separate (uncorrelated) training sets, we could form $B$ separate model fits,

$$\hat{f}^1(X), \ldots, \hat{f}^B(X)$$

Then average them:

$$\hat{f}_B(X) = \frac{1}{B} \sum_{b=1}^{B} \hat{f}^b(X)$$

# Bagging: The bootstrap part

Of course, this isn't practical as having access to many training sets is unlikely.

We therefore turn to the bootstrap to simulate having many training sets.

The bootstrap is a widely applicable statistical tool that can be used to quantify uncertainty without Gaussian approximations.

Let's look at an example.

# Bootstrap detour

# BOOTSTRAP DETOUR

Suppose we are looking to invest in two financial instruments, $X$ and $Y$. The return on these investments is random, but we still want to allocate our money in a risk minimizing way.

That is, for some $\alpha \in (0, 1)$, we want to minimize

$$\text{Var}(\alpha X + (1 - \alpha)Y)$$

The minimizing $\alpha$ is:

$$\alpha_* = \frac{\sigma_Y^2 - \sigma_{XY}^2}{\sigma_X^2 + \sigma_Y^2 - 2\sigma_{XY}^2}$$

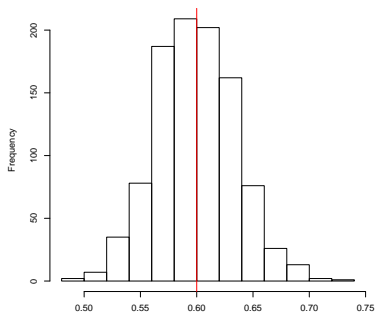(Here, $\sigma_{XY}^2$ is the covariance between $X$ and $Y$)

which we can estimate via

$$\hat{\alpha} = \frac{\hat{\sigma}_Y^2 - \hat{\sigma}_{XY}^2}{\hat{\sigma}_X^2 + \hat{\sigma}_Y^2 - 2\hat{\sigma}_{XY}^2}$$
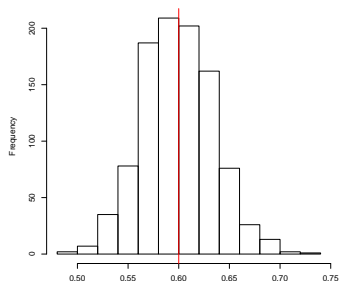
# Bootstrap detour

Now that we have an estimator of $\alpha$, it would be nice to have an estimator of its variability. In this case, computing a standard error is difficult.

Suppose for a moment that we can simulate a large number of draws (say 1000) of the data, which has actual value $\alpha = 0.6$. Then we could get estimates $\hat{\alpha}_1, \ldots, \hat{\alpha}_{1000}$:



This is the sampling distribution of $\hat{\alpha}$

# Bootstrap detour



The mean of all of these is:

$$\overline{\alpha} = \frac{1}{1000} \sum_{r=1}^{1000} \hat{\alpha}_r = 0.599,$$

which is very close to 0.6 (red line), and the standard error is

$$\sqrt{\frac{1}{1000-1} \sum_{r=1}^{1000} (\hat{\alpha}_r - \overline{\alpha})^2} = 0.035.$$

# Bootstrap detour

The standard error of 0.035 gives a very good idea of the accuracy of $\hat{\alpha}$ for a single sample. Roughly speaking, for a new random sample, we expect $\hat{\alpha} \in (\alpha - 2 * 0.035, \alpha + 2 * 0.035)$.

In practice, of course, we cannot use this procedure as it relies on being able to draw a large number of (independent) samples from the same distribution as our data.

This is where the bootstrap comes in.

We instead draw a large number of samples directly from our observed data. This sampling is done with replacement, which means that the same data point can be drawn multiple times.

# Bootstrap detour: Small example

Suppose we have data $\mathcal{D} = (4.3, 3, 7.2, 6.9, 5.5)$.

Then we can draw bootstrap samples, which might look like:

$$\mathcal{D}_1^* = (7.2, 4.3, 7.2, 5.5, 6.9)$$
$$\mathcal{D}_2^* = (6.9, 4.3, 3.0, 4.3, 6.9)$$
$$\vdots$$
$$\mathcal{D}_B^* = (4.3, 3.0, 3.0, 5.5, 6.9)$$

It turns out each of these $\mathcal{D}_b^*$ have very similar properties as $\mathcal{D}$
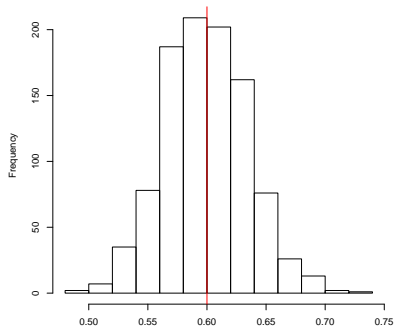
# BOOTSTRAP DETOUR: SMALL EXAMPLE

Now, we form the bootstrap mean:

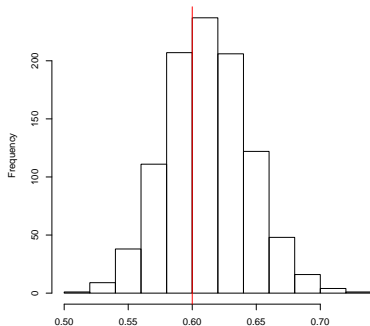$$\text{mean}_B = \frac{1}{B} \sum_{b=1}^{B} \hat{\alpha}_b^*$$

The bootstrap estimator of the standard error is:

$$\text{SE}_B = \sqrt{\frac{1}{B} \sum_{b=1}^{B} (\hat{\alpha}_b^* - \text{mean}_B)^2}$$

# Bootstrap detour



Sampling distribution of $\hat{\alpha}$
(impossible to form)

Bootstrap distribution of $\hat{\alpha}$
(possible to form)

# BOOTSTRAP: END DETOUR

## SUMMARY:

Suppose we have data $\mathcal{D} = (Z_1, \ldots, Z_n)$ and we want to get an idea of the sampling distribution of some statistic $\hat{f}$ trained on $\mathcal{D}$.

Then we do the following: Fix a large number $B$

($B$ could be, say, 1000)

Then for each $b = 1, \ldots, B$

1. Form a new bootstrap draw from $\mathcal{D}$, call it $\mathcal{D}^*$
2. Compute $\hat{f}_b^*$ from $\mathcal{D}^*$

Now, we can estimate the distribution of $\hat{f}$ trained on $\mathcal{D}$ by looking at the distribution of the $B$ draws, $\hat{f}_b^*$

End detour

# Bagging: The bootstrap part

Now, instead of having $B$ separate training sets, we train on $B$ bootstrap draws:
$$\hat{f}_1^*(X), \ldots, \hat{f}_B^*(X)$$

and then average them:

$$\hat{f}_{\text{bag}}(X) = \frac{1}{B} \sum_{b=1}^{B} \hat{f}_b^*(X)$$

This process is known as Bagging

# Bagging trees

# BAGGING TREES

The procedure for trees is the following

1. Choose a large number $B$.
2. For each $b = 1, \ldots, B$, grow an unpruned tree on the $b^{th}$ bootstrap draw from the data.
3. Average all these trees together.

Each tree, since it is unpruned, will have (low/high) variance and (low/high) bias

# BAGGING TREES

The procedure for trees is the following

1. Choose a large number $B$.
2. For each $b = 1, \ldots, B$, grow an unpruned tree on the $b^{th}$ bootstrap draw from the data.
3. Average all these trees together.

Each tree, since it is unpruned, will have (low/high) variance and (low/high) bias

Therefore averaging many trees results in an estimator that has lower variance and still low bias.

# Additional tree bagging topics

# BAGGING TREES

Now that we are growing a large number ($B$) of random trees, we can't directly look at the dendrogram

We no longer have that nice diagram that shows the segmentation of the predictor space (More accurately, we have $B$ of them)

However, we do get some helpful information instead:

- Mean decrease variable importance
- Permutation variable importance
- Out-of-Bag error estimation (OOB)
  (Each time a tree is grown, we can get its prediction error on the unused observations. We average this over all bootstrap samples)
- Proximity plot

Observation: At every split of a node, the loss function decreases

Hence, adding up the amount of decrease for each covariate over all trees gives an indication of feature importance

Intuitively an important covariate is one that if split upon, it leads to a large drop in loss function
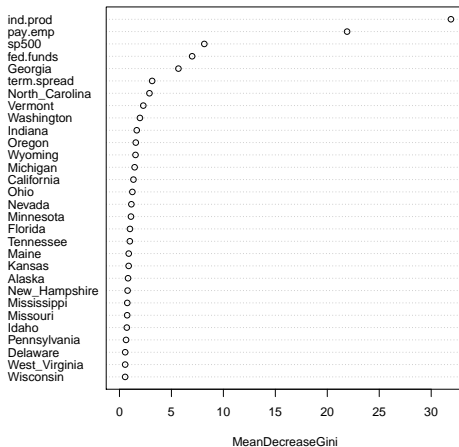
# Mean decrease variable importance

To recover some information, we can do the following:

1. For each of the $B$ trees and each of the $p$ features, we record the amount that the Gini index (or cross-entropy) is reduced by the addition of that feature
2. Report the average reduction over all $B$ trees

This gives us an indication of the importance of a feature

# Mean decrease variable importance



MeanDecreaseGini

# OUT-OF-BAG SAMPLES (OOB)

One can show that, on average, drawing $n$ samples from $n$ observations with replacement results in about $2/3$ of the observations being selected.

The remaining one-third of the observations not used are referred to as out-of-bag (OOB)

# Out-of-Bag samples (OOB)

We can think of it as a for-free cross-validation

The observations that aren't included serve as test data

This provides a free estimate of prediction risk for each tree

We can therefore get an overall estimate of prediction risk by averaging these estimates over all bootstrapped trees

# PERMUTATION VARIABLE IMPORTANCE

Consider the $b^{th}$ bootstrap sample

1. The OOB prediction accuracy is recorded for the $b^{th}$ tree
2. The $j^{th}$ feature is randomly permuted in the OOB sample
   (ie: If $\mathcal{D}_{OOB,b} = (Z_{t_1}, \ldots, Z_{t_{n/3}})$ then permute $X_t^j$ to form $\tilde{X}_t^j$ and hence
   $\tilde{Z}_{t_i} = (Y_{t_i}, X_{t_i}^1, \ldots, \tilde{X}_{t_i}^j, \ldots, X_{t_i}^p).$)
3. The prediction error is recomputed and the change in prediction error is recorded
   (That is, get the OOB error rate for $\tilde{\mathcal{D}}_{OOB,b} = (\tilde{Z}_{t_1}, \ldots, \tilde{Z}_{t_{n/3}})$)

INTUITION: If a feature is highly important, then the OOB prediction error should increase substantially after permuting the OOB values for that feature

# PROXIMITY PLOT

For the $b^{th}$ tree, we can examine which OOB observations are assigned to the same terminal node

Form an $n \times n$ matrix $P$ and increment $P[i, i'] \leftarrow P[i, i'] + 1$ if $Z_i$ and $Z_{i'}$ are assigned to the same terminal node

Now, use some sort of dimension reduction technique to visualize the data in 2-3 dimensions

(Multidimensional scaling is most commonly used (between observation distances are preserved))

The idea is that even if the data have combinations of qualitative/quantitative variables and/or have high dimension, we can view their similarity through the forward operator of the bagged estimator

# Random forest

# Random Forest

Random Forest is a small extension of Bagging, in which the bootstrap trees are decorrelated

The idea is, we draw a bootstrap sample and start to build a tree.

- At each split, we randomly select $m$ of the possible $p$ features as candidates for the split.
- A new sample of size $m$ of the features is taken at each split.

Usually, we use about $m = \sqrt{p}$

(this would be 7 out of 56 features for GDP data)

In other words, at each split, we aren't even allowed to consider the majority of possible features!

# Random Forest

Suppose there is 1 really strong feature and many mediocre ones.

- Then each tree will have this one feature in it,
- Therefore, each tree will look very similar (i.e. highly correlated).
- Averaging highly correlated things leads to much less variance reduction than if they were uncorrelated.

If we don't allow some trees/splits to use this important feature, each of the trees will be much less similar and hence much less correlated.

Bagging is Random Forest when $m = p$, that is, when we can consider all the features at each split.

# Random forest

An average of $B$ i.i.d random variables has variance

$$\frac{\sigma^2}{B}$$

An average of $B$ random variables has variance

$$\rho\sigma^2 + \frac{(1-\rho)\sigma^2}{B}$$

for correlation $\rho$

As $B \to \infty$, the second term goes to zero, but the first term remains

Hence, correlation of the trees limits the benefit of averaging

# Sensitivity and specificity

SENSITIVITY: The proportion of times we label recession, given that recession is the correct answer.

SPECIFICITY: The proportion of times we label no recession, given that no recession is the correct answer.

We can think of this in terms of hypothesis testing. If

$$H_0 : \text{ no recession,}$$

then

SENSITIVITY: $P(\text{reject } H_0 | H_0 \text{ is false})$, [1 - $P$(Type II error)]

SPECIFICITY: $P(\text{accept } H_0 | H_0 \text{ is true})$, [1 - $P$(Type I error)]

# Confusion matrix

We can report our results in a matrix:

|  |  | Truth | |
|---|---|---|---|
|  |  | Recession | No Recession |
| Our | Recession | (A) | (B) |
| Predictions | No Recession | (C) | (D) |

The total number of each combination is recorded in the table.

The overall miss-classification rate is

$$\frac{(B) + (C)}{(A) + (B) + (C) + (D)} = \frac{(B) + (C)}{\text{total observations}}$$

What is the sensitivity/specificity?

# Confusion matrix

We can report our results in a matrix:

|  |  | Truth | |
|---|---|---|---|
|  |  | Recession | No Recession |
| Our | Recession | (A) | (B) |
| Predictions | No Recession | (C) | (D) |

The total number of each combination is recorded in the table.

The overall miss-classification rate is

$$\frac{(B) + (C)}{(A) + (B) + (C) + (D)} = \frac{(B) + (C)}{\text{total observations}}$$

What is the sensitivity/specificity?

(Sensitivity is (A)/[(A) + (C)], Specificity is (D)/[(B) + (D)])

# TREE RESULTS: CONFUSION MATRICES

| | | | Truth | | |
|---|---|---|---|---|---|
| | | | Growth | Recession | Mis-Class |
| Our Preds | NULL | Growth | 111 | 26 | |
| | | Recession | 0 | 0 | 18.9% |
| | TREE | Growth | 99 | 3 | |
| | | Recession | 12 | 23 | 10.9% |
| | RANDOM FOREST | Growth | 102 | 5 | |
| | | Recession | 9 | 21 | 10.2% |
| | BAGGING | Growth | 104 | 3 | |
| | | Recession | 7 | 23 | 7.3% |

# Tree results: Sensitivity & specificity

|  | Sensitivity | Specificity |
|---|---|---|
| NULL | 0.000 | 1.000 |
| TREE | 0.884 | 0.891 |
| RANDOM FOREST | 0.807 | 0.918 |
| BAGGING | 0.884 | 0.936 |

# OUT-OF-BAG ERROR ESTIMATION FOR BAGGING

| | | Truth | | |
|---|---|---|---|---|
| | | Growth | Recession | Miss-Class |
| OOB BAGGING | Growth | 400 | 10 | |
| | Recession | 21 | 46 | 6.5% |
| | | | | |
| TEST BAGGING | Growth | 104 | 3 | |
| | Recession | 7 | 23 | 7.3% |

# Random Forest in R

```
require(randomForest)
out.rf    = randomForest(X,Y,importance=T,mtry=p)
class.rf = predict(out.rf,X_0)
```

Notes:

- The importance statement tells it to produce the variable importance measures
- the mtry = p tells randomForest to consider all the covariates at each split

  (This particular choice corresponds to bagging)

- randomForest also supports formulae

  ```
  out.rf    = randomForest(Y~.,data=X)
  ```

  However, it can take much longer to run

# RANDOM FOREST IN R

```
> out.rf

Call:
randomForest(formula = Y~.,data = X, import = T, mtry = p)
               Type of random forest: classification
                     Number of trees: 500
No. of variables tried at each split: 56


         OOB estimate of  error rate: 7.33%
Confusion matrix:
    0  1 class.error
0 508 13  0.02495202
1  32 61  0.34408602
```

# Random Forest in R

```
#Permutation variable importance
> head(importance(out.rf,type=1))
          MeanDecreaseAccuracy
Alabama              3.7277511
Alaska               1.7941463
Arizona              2.9659623
Arkansas             0.8341577
California           7.2973572
#Mean decrease variable importance
> head(importance(out.rf,type=2))
          MeanDecreaseGini
Alabama          0.4551073
Alaska           1.6440170
Arizona          0.7025527
Arkansas         0.3503138
California       1.4616203
```

# ADDITIONAL RANDOM FOREST TOPICS

CLAIM: Random forest cannot overfit. This is and isn't true.

Write

$$\hat{f}_{rf}^B = \frac{1}{B} \sum_{b=1}^{B} T(x; \Theta_b)$$

where $\Theta_b$ characterizes the $b^{th}$ tree

(That is, the split variables, cutpoints of each node, terminal node values)

Increasing $B$ does not cause Random forest to overfit, rather removes the Monte-Carlo-like approximation error

$$\hat{f}_{rf}(x) = \mathbb{E}_\Theta T(x, \Theta) = \lim_{B \to \infty} \hat{f}_{rf}^B$$

However, this limit can overfit the data, the average of fully grown trees can result in too complex of a model

(Note that Segal (2004) shows that a small benefit can be derived by stopping each tree short, but thus induce another tuning parameter)

# ADDITIONAL RANDOM FOREST TOPICS

Things I'd like to cover but may not

(AKA possible short presentation topics)

- Variance and decorrelation effect (showing precisely how random forest may work/not work)
- Introduction of noise covariates improves performance
- Using subsampling instead of bootstrap to generate trees

  (This is an idea I had while writing this up. I don't know if this exists, but it seems to work in many cases the bootstrap doesn't and is easier to do theory (just use Hoeffding's inequalty for $U$-statistics))

- Adaptive nearest neighbors