

Introduction to Machine Learning

Ranga Raju Vatsavai, Ph.D.

Chancellors Faculty Excellence Associate Professor in Geospatial Analytics
Department of Computer Science, North Carolina State University (NCSU)

Feb. 25-27, 2019

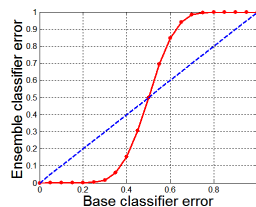
Ensemble Methods

- Construct a set of classifiers from the training data
- Predict class label of test records by combining the predictions made by multiple classifiers

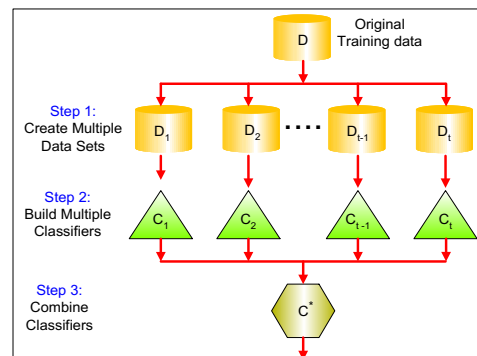
Why Ensemble Methods work?

- Suppose there are 25 binary classifiers
 - Each classifier has error rate, $\epsilon = 0.35$
 - Assume errors made by classifiers are uncorrelated
 - Probability that the ensemble classifier makes a wrong prediction:

$$P(X \geq 13) = \sum_{i=13}^{25} \binom{25}{i} \epsilon^i (1-\epsilon)^{25-i} = 0.06$$



General Approach



Types of Ensemble Methods

- Bayesian ensemble
 - Example: Mixture of Gaussian
- Manipulate data distribution
 - Example: Resampling method
- Manipulate input features
 - Example: Feature subset selection
- Manipulate class labels
 - Example: error-correcting output coding
- Introduce randomness into learning algorithm
 - Example: Random forests

Bagging

- Sampling with replacement

| Original Data | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|-------------------|---|---|----|----|---|---|----|----|---|----|
| Bagging (Round 1) | 7 | 8 | 10 | 8 | 2 | 5 | 10 | 10 | 5 | 9 |
| Bagging (Round 2) | 1 | 4 | 9 | 1 | 2 | 3 | 2 | 7 | 3 | 2 |
| Bagging (Round 3) | 1 | 8 | 5 | 10 | 5 | 5 | 9 | 6 | 3 | 7 |

- Build classifier on each bootstrap sample
- Each sample has probability $1 - (1 - 1/n)^n$ of being selected
- On average, each bootstrap data contains ~63% of original data

Bagging Algorithm

Algorithm 5.6 Bagging Algorithm

```

1: Let  $k$  be the number of bootstrap samples.
2: for  $i = 1$  to  $k$  do
3:   Create a bootstrap sample of size  $n$ ,  $D_i$ .
4:   Train a base classifier  $C_i$  on the bootstrap sample  $D_i$ .
5: end for
6:  $C^*(x) = \arg \max_y \sum_i \delta(C_i(x) = y)$ ,  $\{\delta(\cdot) = 1$  if its argument is true, and 0 otherwise. $\}$ 

```

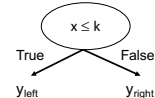
Bagging Example

- Consider 1-dimensional data set:

Original Data:

| x | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | 1 |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|
| y | 1 | 1 | 1 | -1 | -1 | -1 | -1 | 1 | 1 | 1 |

- Classifier is a decision stump
 - Decision rule: $x \leq k$ versus $x > k$
 - Split point k is chosen based on entropy



NC STATE UNIVERSITY

Bagging Example

Bagging Round 1:

| | | | | | | | | | | |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| x | 0.1 | 0.2 | 0.2 | 0.3 | 0.4 | 0.4 | 0.5 | 0.6 | 0.9 | 0.9 |
| y | 1 | 1 | 1 | 1 | -1 | -1 | -1 | -1 | 1 | 1 |

$x \leq 0.35 \Rightarrow y = 1$
 $x > 0.35 \Rightarrow y = -1$

NC STATE UNIVERSITY

Bagging Example

Bagging Round 1:

| | | | | | | | | | | |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| x | 0.1 | 0.2 | 0.2 | 0.3 | 0.4 | 0.4 | 0.5 | 0.6 | 0.9 | 0.9 |
| y | 1 | 1 | 1 | 1 | -1 | -1 | -1 | -1 | 1 | 1 |

$x \leq 0.35 \Rightarrow y = 1$
 $x > 0.35 \Rightarrow y = -1$

Bagging Round 2:

| | | | | | | | | | | |
|---|-----|-----|-----|-----|-----|-----|-----|---|---|---|
| x | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.5 | 0.9 | 1 | 1 | 1 |
| y | 1 | 1 | 1 | -1 | -1 | -1 | 1 | 1 | 1 | 1 |

$x \leq 0.7 \Rightarrow y = 1$
 $x > 0.7 \Rightarrow y = -1$

Bagging Round 3:

| | | | | | | | | | | |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| x | 0.1 | 0.2 | 0.3 | 0.4 | 0.4 | 0.5 | 0.7 | 0.7 | 0.8 | 0.9 |
| y | 1 | 1 | 1 | -1 | -1 | -1 | -1 | -1 | 1 | 1 |

$x \leq 0.35 \Rightarrow y = 1$
 $x > 0.35 \Rightarrow y = -1$

Bagging Round 4:

| | | | | | | | | | | |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| x | 0.1 | 0.1 | 0.2 | 0.4 | 0.4 | 0.5 | 0.5 | 0.7 | 0.8 | 0.9 |
| y | 1 | 1 | 1 | -1 | -1 | -1 | -1 | -1 | 1 | 1 |

$x \leq 0.3 \Rightarrow y = 1$
 $x > 0.3 \Rightarrow y = -1$

Bagging Round 5:

| | | | | | | | | | | |
|---|-----|-----|-----|-----|-----|-----|-----|---|---|---|
| x | 0.1 | 0.1 | 0.2 | 0.5 | 0.6 | 0.6 | 0.6 | 1 | 1 | 1 |
| y | 1 | 1 | 1 | -1 | -1 | -1 | -1 | 1 | 1 | 1 |

$x \leq 0.35 \Rightarrow y = 1$
 $x > 0.35 \Rightarrow y = -1$

NC STATE UNIVERSITY

Bagging Example

Bagging Round 6:

| | | | | | | | | | | |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|
| x | 0.2 | 0.4 | 0.5 | 0.6 | 0.7 | 0.7 | 0.7 | 0.8 | 0.9 | 1 |
| y | 1 | -1 | -1 | -1 | -1 | -1 | -1 | 1 | 1 | 1 |

$x \leq 0.75 \Rightarrow y = -1$
 $x > 0.75 \Rightarrow y = 1$

Bagging Round 7:

| | | | | | | | | | | |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|
| x | 0.1 | 0.4 | 0.4 | 0.6 | 0.7 | 0.8 | 0.9 | 0.9 | 0.9 | 1 |
| y | 1 | -1 | -1 | -1 | -1 | 1 | 1 | 1 | 1 | 1 |

$x \leq 0.75 \Rightarrow y = -1$
 $x > 0.75 \Rightarrow y = 1$

Bagging Round 8:

| | | | | | | | | | | |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|
| x | 0.1 | 0.2 | 0.5 | 0.5 | 0.5 | 0.7 | 0.7 | 0.8 | 0.9 | 1 |
| y | 1 | 1 | -1 | -1 | -1 | -1 | -1 | 1 | 1 | 1 |

$x \leq 0.75 \Rightarrow y = -1$
 $x > 0.75 \Rightarrow y = 1$

Bagging Round 9:

| | | | | | | | | | | |
|---|-----|-----|-----|-----|-----|-----|-----|-----|---|---|
| x | 0.1 | 0.3 | 0.4 | 0.4 | 0.6 | 0.7 | 0.7 | 0.8 | 1 | 1 |
| y | 1 | 1 | -1 | -1 | -1 | -1 | -1 | 1 | 1 | 1 |

$x \leq 0.75 \Rightarrow y = -1$
 $x > 0.75 \Rightarrow y = 1$

Bagging Round 10:

| | | | | | | | | | | |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| x | 0.1 | 0.1 | 0.1 | 0.1 | 0.3 | 0.3 | 0.8 | 0.8 | 0.9 | 0.9 |
| y | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

$x \leq 0.05 \Rightarrow y = 1$
 $x > 0.05 \Rightarrow y = -1$

NC STATE UNIVERSITY

Bagging Example

- Summary of Training sets:

| Round | Split Point | Left Class | Right Class |
|-------|-------------|------------|-------------|
| 1 | 0.35 | 1 | -1 |
| 2 | 0.7 | 1 | -1 |
| 3 | 0.35 | 1 | -1 |
| 4 | 0.3 | 1 | -1 |
| 5 | 0.35 | 1 | -1 |
| 6 | 0.75 | -1 | 1 |
| 7 | 0.75 | -1 | 1 |
| 8 | 0.75 | -1 | 1 |
| 9 | 0.75 | -1 | 1 |
| 10 | 0.05 | 1 | 1 |

Bagging Example

- Assume test set is the same as the original data
- Use majority vote to determine class of ensemble classifier

| Round | x=0.1 | x=0.2 | x=0.3 | x=0.4 | x=0.5 | x=0.6 | x=0.7 | x=0.8 | x=0.9 | x=1.0 |
|-----------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1 | 1 | 1 | 1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 3 | 1 | 1 | 1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| 4 | 1 | 1 | 1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| 5 | 1 | 1 | 1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
| 6 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | 1 | 1 | 1 |
| 7 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | 1 | 1 | 1 |
| 8 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | 1 | 1 | 1 |
| 9 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | 1 | 1 | 1 |
| 10 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| Sum | 2 | 2 | 2 | -6 | -6 | -6 | -6 | 2 | 2 | 2 |
| Predicted Class | 1 | 1 | 1 | -1 | -1 | -1 | -1 | 1 | 1 | 1 |

Boosting

- An iterative procedure to adaptively change distribution of training data by focusing more on previously misclassified records
 - Initially, all "n" records are assigned equal weights
 - Unlike bagging, weights may change at the end of each boosting round

Boosting

- Records that are wrongly classified will have their weights increased
- Records that are classified correctly will have their weights decreased

| Original Data | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|--------------------|---|---|---|----|---|---|---|----|---|----|
| Boosting (Round 1) | 7 | 3 | 2 | 8 | 7 | 9 | 4 | 10 | 6 | 3 |
| Boosting (Round 2) | 5 | 4 | 9 | 4 | 2 | 5 | 1 | 7 | 4 | 2 |
| Boosting (Round 3) | 4 | 4 | 8 | 10 | 4 | 5 | 4 | 6 | 3 | 4 |

- Example 4 is hard to classify
- Its weight is increased, therefore it is more likely to be chosen again in subsequent rounds

Random Forests

- Class of ensemble methods specifically designed for decision tree classifiers
- Bagging using decision trees is a special case of random forests
- RF method combines Breiman's "bagging" idea with the random feature selection.

Random Forests

Each tree is constructed using the following algorithm:

Given

- Let the number of training cases be N , and the number of variables in the classifier be M .
- m be the number of input variables used to determine the decision at a node of the tree; m should be much less than M .

Generate training data (bootstrap)

- Choose a training set for each tree by choosing n times with replacement from all N available training cases (i.e. take a bootstrap sample). Use the rest of the cases to estimate the error of the tree, by predicting their classes.
- For each node of the tree, randomly choose m variables on which to base the decision at that node. Calculate the best split based on these m variables in the training set.
- Each tree is fully grown and not pruned

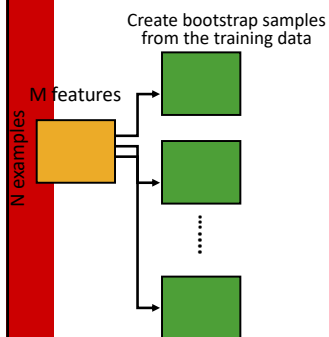
Prediction

- Each tree is used to predict the label. Voting scheme is used to determine the final class label.

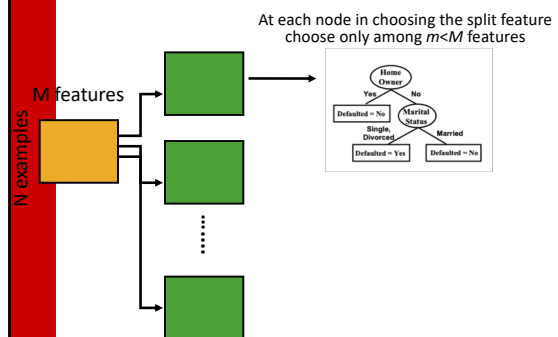
Random Forests

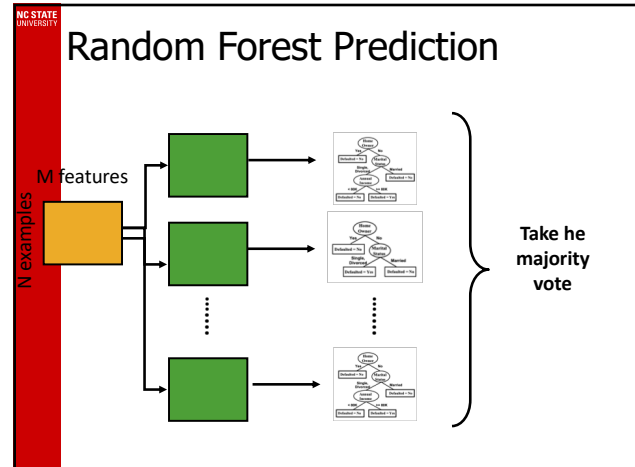
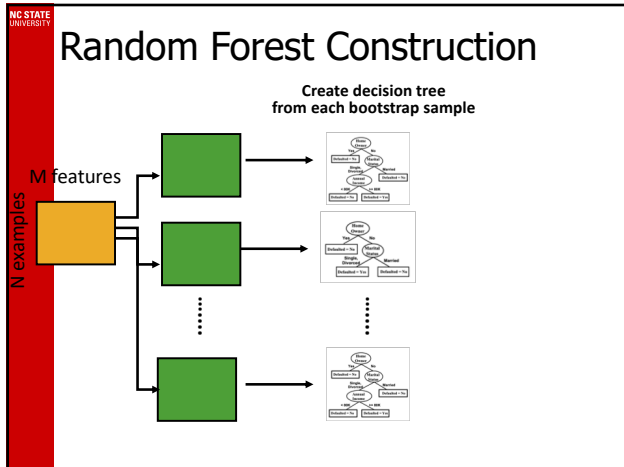
- Class of ensemble methods specifically designed for decision tree classifiers
- Bagging using decision trees is a special case of random forests
- RF method combines Breiman's "bagging" idea with the random feature selection.

Random Forest Construction



Random Forest Construction





NC STATE UNIVERSITY

Random Forests

- Advantages
 - One accurate and widely used algorithm
 - Can handle large number of features without need for feature selection first
 - Provides estimates of which variables are important
 - Handles missing data well
- Disadvantages
 - May lead to overfitting on some datasets
 - May be biased towards attributes that result in large number splits

NC STATE UNIVERSITY

Practice

- On one or more of UCI ML Repository (<http://archive.ics.uci.edu/ml/>)
 - Construct decision tree
 - Construct random forest
 - Compare and contrast both
- Acknowledgements
 - Slides adopted from UMN and CMU classes