

Lecture 1: Ensemble Learning: Bagging

Professor: Dr Ernest Fokoué

Scribe: Miss Sailee Rumao

I. INTRODUCTION

When different models which are qualitatively very different and have almost similar model selection scores (eg. BIC), **Model selection** is not necessarily the most ideal way for classification because when you choose one, you are making strong rejection against other variables that you do not choose.

This is one of the motivation for **Model Aggregation**.

We also know that,

$$V(\bar{x}) < V(x_i) \quad (\forall i = 1, 2, \dots, n) \quad (1)$$

This is another motivation for Model Aggregation. In such a case, ensemble (Aggregated Model) will have very low variance which results in low prediction error based on the **Bias-Variance Decomposition**.

II. ENSEMBLE LEARNING

It is a method used for increasing model accuracy by aggregating k learned models C_1, C_2, \dots, C_k (classifiers or predictors) with an aim of creating an improved composite model C^* .

Here, C : Classifier.

Consider $\mathbb{D} = \{Z_1, Z_2, \dots, Z_n\}$ where $Z_i = (X_i, Y_i) \in \mathcal{X} \times \mathcal{Y}$ where are two random variables.

Assume f : function that models relationship between \mathcal{X} and \mathcal{Y}

$f_1(\cdot), f_2(\cdot), \dots, f_m(\cdot)$: M different estimators of f

α_m : Weights that measure relative importance of each estimator $\hat{f}_m(\cdot)$ $m=1,2,3,\dots,M$

Our goal is to find an estimator of f with the smallest prediction error.

Mathematically,

$$\hat{f}^{(agg)}(\cdot) = \sum_{m=1}^M \alpha_m \hat{f}_m(\cdot) \quad i = 1, 2, \dots, M \quad (2)$$

Each of the estimator $\hat{f}_m(\cdot)$ is called **Base Learner**.

Note: If the base estimators $\hat{f}_m(\cdot)$ are strongly correlated then the performance of the ensemble is likely to decrease.

The two techniques/Methods of Ensemble learning are:

1. Bagging
2. Boosting

II.1 Bagging(Leo Breiman (1996))

Bagging is also known as Bootstrap aggregation.

The steps for bagging are as follows:

1. Drawing $\mathbb{D}^{(m)} = (\mathbf{x}_i^{(m)}, \mathbf{y}_i^{(m)})$ with replacement from \mathbb{D} . This is called bootstrap sampling.
2. Build base learners $\hat{f}_m(\cdot)$ using $\mathbb{D}^{(m)}$ and the chosen learning machine.
3. Then if all base learners have equal weights i.e $\hat{\alpha}_m = \frac{1}{M}$ (in case of bagging),

$$\hat{f}^{(agg)}(\cdot) = \frac{1}{M} \sum_{m=1}^M \hat{f}_m(\cdot) \quad (3)$$

Note:

1. Many $\mathbf{z}_i = (\mathbf{x}_i, \mathbf{y}_i)$ from \mathbb{D} are replaced many times in $\mathbb{D}^{(m)}$
2. $|\mathbb{D}^{(m)}| = n$ (which is the size of \mathbb{D})
3. In Regression Model avergaing is obvious and is given as,

$$\hat{f}^{(agg)}(X_{new}) \stackrel{Reg}{=} \frac{1}{M} \sum_{m=1}^M \hat{f}^{(m)}(X_{new}) \quad (4)$$

In classification we use the majority rule (for the class labels obtained) to make a decision which is given as follows:

$$\hat{f}^{(agg)}(X_{new}) \stackrel{Classif}{=} \underset{g=1,2,3...G}{\operatorname{argmax}} \left\{ \frac{1}{M} \sum_{m=1}^M \mathbb{1}(\hat{f}^{(m)}(X_{new}) = g) \right\} \quad (5)$$

4. Bias of Machine depends on the bias of the base learners.

Out of Bag Error

Consider $\mathbf{z} \in \mathbb{D}$ be a random pair in \mathbb{D} (original dataset of size n) and
Consider bootstrapped sample $\mathbb{D}^{(m)}$ of size n , then,

$$\begin{aligned} Pr(\mathbf{z} \in \mathbb{D}^{(m)}) &= \text{Proportion of observations from } \mathbb{D} \text{ present in } \mathbb{D}^{(m)} = 1 - (1 - \frac{1}{n})^n \dots\dots (*) \\ Pr(\mathbf{z} \notin \mathbb{D}^{(m)}) &= \text{Proportion of observations from } \mathbb{D} \text{ present in } \mathbb{D}^{(m)} = (1 - \frac{1}{n})^n = Pr[O_n] \dots\dots (**) \end{aligned}$$

As $n \rightarrow \infty$, $Pr[O_n] \rightarrow e^{-1} = 0.37$

This implies that approximately one third of the training set is not used to build m th bootstrapped base learner. This proportion is used to calculate the out of bag error.(used in random forest to estimate variable importance.)

Now let $\hat{f}^{(m)}(\cdot)$ be the base learner from $\mathbb{D}^{(m)}$.
Then the out of bag error of $\hat{f}^{(m)}(\cdot)$ is ,

$$err_{OOB}(\hat{f}^{(m)}(\cdot)) = \frac{\sum_{i=1}^n \mathbb{1}(\mathbf{z}_i \notin \mathbb{D}^{(m)}) \ell(\mathbf{y}_i, \hat{f}^{(m)}(\mathbf{x}_i))}{\sum_{i=1}^n \mathbb{1}(\mathbf{z}_i \notin \mathbb{D}^{(m)})} \quad (6)$$

where,

$\mathbb{1}(\cdot)$: Indicator function

$\ell(\cdot, \cdot)$: Loss function

Thus, the OOB error for the ensemble is,

$$err_{OOB}(\hat{f}^{(agg)}(\cdot)) = \frac{1}{M} \sum_{m=1}^M err_{OOB}(\hat{f}^{(m)}(\cdot)) \quad (7)$$

Note:

1. Out of bag error is a good estimate of test error.
2. $\lceil e^{-1} \rceil =$ size of the out of bag sample(OOB sample) \approx Number of observations from \mathbb{D} not in $\mathbb{D}^{(m)}$

Properties of Bagged Ensembles:

1. If $\hat{f}^{(m)}(\cdot)$ is an unbiased estimator of f^* i.e $\mathbb{E}[\hat{f}^{(m)}(\cdot)] = f^*$ then $\hat{f}^{(agg)}(\cdot)$ will also be unbiased i.e $\mathbb{E}[\hat{f}^{(agg)}(\cdot)] = f^*$

Expectation of final predictor is,

$$\mathbb{E}[\hat{f}^{(agg)}(\cdot)] \stackrel{bagging}{=} \mathbb{E}\left[\frac{1}{M} \sum_{m=1}^M \hat{f}^{(m)}(\cdot)\right] = \frac{1}{M} \sum_{m=1}^M \mathbb{E}[\hat{f}^{(m)}(\cdot)] \quad (8)$$

Thus the above statement.

2. Assume,

$$\mathbb{V}(\hat{f}^{(m)}(\cdot)) = \sigma^2 \quad \forall m = (1, 2, \dots, M)$$

$$\text{Cov}(\hat{f}^{(l)}(\cdot), \hat{f}^{(m)}(\cdot)) = \tau$$

Then,

$$\mathbb{V}(\hat{f}^{(agg)}(\cdot)) = \mathbb{V}\left(\frac{1}{M} \sum_{m=1}^M \hat{f}^{(m)}(\cdot)\right) = \frac{1}{M^2} \mathbb{V}\left(\sum_{m=1}^M \hat{f}^{(m)}(\cdot)\right) + \frac{1}{M^2} \sum_{l=1}^M \sum_{m=1}^M \text{Cov}(\hat{f}^{(l)}(\cdot), \hat{f}^{(m)}(\cdot)) \quad (9)$$

Now, if $\tau = 0$

then ,

$$\mathbb{V}(\hat{f}^{(agg)}(\cdot)) = \frac{1}{M^2} \mathbb{V}\left(\sum_{m=1}^M \hat{f}^{(m)}(\cdot)\right) = \frac{\sigma^2}{M} \quad (10)$$

and $\lim_{M \rightarrow \infty} (\mathbb{V}(\hat{f}^{(agg)}(\cdot))) = 0$

i.e As M tends to infinity , variance of the aggregate Machine decreases without any decrease in its bias resulting in bias-variance tradeoff and therefore leading to an optimal predictive performance.

In other words, if the individual machines are unbiased and they are near orthogonal , the aggregate machine will become a consistent estimator of the true function.

Note: Bagging always uses the whole p-dimensional vector X_i . This can be bad for $n \ll p$ (curse of dimensionality) and for variable selection.

Random Forest

Random Forest overcomes the drawback of Curse of Dimensionality in Bagging.

Bagging is performed in Random forest in the following steps:

1.

1. Pick a subset of d variables from the original p -variables
($\mathbf{x}_{j1}, \mathbf{x}_{j2} \dots, \mathbf{x}_{jd}$) : Subset of variables (where $d < p$)
Here \mathbf{x}_{jk} are chosen without replacement
 $\mathbb{D}^{(m)}$ is an $n \times d$ matrix

2. Choose a bootstrap sample

3. Construct base learner

2. Device for Variable Importance

1. Take out of bag sample of size $S = \lceil e^{-1}n \rceil$

2. After $\hat{f}^{(m)}(\cdot)$ is built, calculate out of bag error i.e $\hat{\mathbb{E}}_{OOB}$ (raw)

3. For $k=1$ to d ,

- (a) Perform permutation of column j_k in OOB sample (shuffle the entries in column j_k)

- (b) Compute $\hat{\mathbb{E}}_{OOB}^{(k)}$

- (c) Compute $\hat{\mathbb{E}}_{OOB} - \hat{\mathbb{E}}_{OOB}^{(k)}$

This Mechanism allows Random forest to measure variable importance.

Advantages of Random Forest:

1. Achieves a more powerful and flexible machine.
2. Gives a measure of variable importance.

II.2 Boosting

The importance of boosting lies in combining various base learners assumed to be weak (descent machines) in the sense that each of them is just slightly better than random guessing.

Boosting uses complicated α .

We'll look at the working of Adaptive boosting which is a type of Boosting in the next lecture