

7. Boosting

Like bagging and random forests, boosting involves creating many models.

Unlike bagging and random forests, boosting creates these models sequentially, and there is no resampling involved.

Boosting **learns slowly** by repeatedly fitting new models based on the residuals of earlier models.

For regression trees, the algorithm works as follows:

1. Set $\hat{\mathbf{f}}(\mathbf{x}) = \mathbf{0}$ and $r_i = y_i$ for all i .
2. For $m = 1, \dots, M$, repeat:
 - i. Fit a tree $\hat{\mathbf{f}}^m(\mathbf{x})$ with d splits to the training data (\mathbf{x}, \mathbf{r}) .
 - ii. Update $\hat{\mathbf{f}}$ by adding a shrunk version of $\hat{\mathbf{f}}^m(\mathbf{x})$ to it:

$$\hat{\mathbf{f}}(\mathbf{x}) \leftarrow \hat{\mathbf{f}}(\mathbf{x}) + \lambda \hat{\mathbf{f}}^m(\mathbf{x}). \quad (1)$$

iii. Update the residuals:

$$r_i \leftarrow r_i - \lambda \hat{f}^m(\mathbf{x}_i). \quad (2)$$

3. The boosted model is

$$\hat{\mathbf{f}}(\mathbf{x}) = \sum_{m=1}^M \lambda \hat{\mathbf{f}}^m(\mathbf{x}). \quad (3)$$

Recall that the initial value of $\hat{\mathbf{f}}$ was $\mathbf{0}$. Thus all of the explanatory power comes from the $\hat{\mathbf{f}}^m(\mathbf{x})$.

There are three tuning parameters:

1. The number of (in this case) trees. There is a risk of overfitting if M gets too large, so we need to use cross-validation.
2. The **shrinkage parameter** λ . Typical values are 0.01 and 0.001. When λ is very small, M needs to be large.
3. The number of splits d , called the **interaction depth**. This tends to be small, perhaps just $d = 1$.

When $d = 1$, every split is a stump, and (3) becomes an additive model.

It seems odd that λ is not set to $1/M$. Do we have to rescale λ if M becomes larger?

See ISLR Figure 8.11.

For squared error loss (which may not be a good thing to use; see ESL Section 10.6), the objective function is

$$\sum_{i=1}^N L(y_i, f(\mathbf{x}_i)) = \sum_{i=1}^N (y_i - f(\mathbf{x}_i))^2. \quad (4)$$

Since we build up the $f(\mathbf{x}_i)$ slowly in (1), we see that

$$\begin{aligned} L(y_i, \hat{f}(\mathbf{x}_i)) &\leftarrow L(y_i, \hat{f}(\mathbf{x}_i) + \lambda \hat{f}^m(\mathbf{x}_i)) \\ &= (y_i - \hat{f}(\mathbf{x}_i) - \lambda \hat{f}^m(\mathbf{x}_i))^2 \\ &= (r_i - \lambda \hat{f}^m(\mathbf{x}_i))^2, \end{aligned} \quad (5)$$

where r_i is simply the i^{th} residual for the current model, before we have added the m^{th} term to it.

At each step in the boosting algorithm, we add λ times the term $\hat{f}^m(\mathbf{x}_i)$ that best fits the function

$$\sum_{i=1}^N (r_i - \hat{f}^m(\mathbf{x}_i))^2, \quad (6)$$

which just depends on the current residuals and on $\hat{f}^m(\mathbf{x}_i)$.

The discussion above assumes that we use a tree to obtain $\hat{f}^m(\mathbf{x}_i)$, but many other models can also be boosted.

It does not make sense to boost a linear regression model, because the residuals are orthogonal to all the predictors.

7.1. AdaBoost.M1 for two-way classification

We code the output as $\{-1, 1\}$ instead of $\{0, 1\}$.

The error rate on the training sample is

$$\overline{\text{err}} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i \neq G(\mathbf{x}_i)), \quad (7)$$

where $G(\mathbf{x}_i)$ denotes the classifier. It also takes values $\{-1, 1\}$.

A **weak classifier** is only slightly better than random guessing.

Boosting sequentially applies the weak classifier to repeatedly modified versions of the data (in the regression case above, these were residuals). Eventually, we have M of these, denoted $G_m(\mathbf{x})$ for $m = 1, \dots, M$.

We then combine these as follows to produce the final prediction

$$G(\mathbf{x}) = \text{sign} \left(\sum_{m=1}^M \alpha_m G_m(\mathbf{x}) \right), \quad (8)$$

where the weights α_m have to be determined.

Instead of using residuals from successive steps, classification uses weights that change as the algorithm proceeds.

Observations that were misclassified get more weight, and observations that were correctly classified get less weight.

The `AdaBoost.M1` algorithm works as follows:

1. Initialize the observations weights to $1/N$ for all i .
2. For $m = 1, \dots, M$:
 - i. Fit a classifier $G_m(\mathbf{x})$ to the training data using weights w_i .
 - ii. Compute

$$\text{err}_m = \frac{\sum_{i=1}^N w_i \mathbb{I}(y_i \neq G_m(\mathbf{x}_i))}{\sum_{i=1}^N w_i}.$$

- iii. Compute $\alpha_m = \log((1 - \text{err}_m)/\text{err}_m)$.
- iv. Set $w_i \leftarrow w_i \exp(\alpha_m \mathbb{I}(y_i \neq G_m(\mathbf{x}_i)))$ for $i = 1, \dots, N$.

3. The final classifier is given in (8): $G(\mathbf{x}) = \text{sign}\left(\sum_{m=1}^M \alpha_m G_m(\mathbf{x})\right)$.

The algorithm just described is called “Discrete AdaBoost,” because the classifier always reports -1 or 1 .

A modified version called “Real AdaBoost” was proposed in Friedman, Hastie, and Tibshirani (AMS, 2000). This seems to be a key paper. Real AdaBoost yields real numbers in the $[0, 1]$ interval as outputs instead of $\{-1, 1\}$.

In the above algorithm, step 1. is unchanged. Step 2. becomes

2. For $m = 1, \dots, M$:

- i. Fit a classifier $p_m(\mathbf{x})$ to the training data using weights w_i .
- ii. Compute $f_m = \frac{1}{2} \log(p_m(\mathbf{x})/(1 - p_m(\mathbf{x})))$.
- iii. Set $w_i \leftarrow w_i \exp(-y_i f_m(\mathbf{x}_i))$ for all i , and renormalize so that $\sum_{i=1}^N w_i = 1$.

3. The final classifier is the sign of $\sum_{m=1}^M f_m(\mathbf{x})$.

Observe that the f_m here are the logs of odds ratios.

ESL present a simulated example in which $p = 10$ and each input X_j is $N(0, 1)$.

The output is 1 if

$$\sum_{j=1}^{10} X_j^2 > \chi_{10}^2(0.5) = 9.34182. \quad (9)$$

Thus, on average, half of the outputs will be 1 and half will be -1 . But, since we observe the realized random variables, we should be able to make predictions.

There are 2000 training cases and 10,000 test ones.

The weak classifier is a tree with two terminal nodes, a “stump.”

The error rate of the initial classifier is 45.8%. After boosting with $M = 400$, it is only 5.8%. A single large tree has an error rate of 24.7%.

See ESL, Figure 10.2.

7.2. Boosting and additive models

Because boosting is additive—see (3) and (8)—it can be thought of as a form of basis expansion.

Basis expansions take the form

$$f(\mathbf{x}) = \sum_{m=1}^M \beta_m b(\mathbf{x}, \gamma_m), \quad (10)$$

where the β_m are expansion coefficients, and the $b(\cdot)$ are (usually) simple functions of \mathbf{x} and the coefficient vectors γ_m .

For neural nets, the $b(\cdot)$ are sigmoid functions of linear combinations of \mathbf{x} .

For MARS, the $b(\cdot)$ are splines, and the γ_m parametrize the variables and the values for the knots.

For trees, the γ_m parametrize the split variables and split points.

In general, we estimate this sort of model by minimizing a loss function with respect to the β_m and the γ_m :

$$\sum_{i=1}^N L \left(y_i, \sum_{m=1}^M \beta_m b(\mathbf{x}, \gamma_m) \right) \quad (11)$$

This is hard if we have to minimize with respect to all parameters at once, but it can be easy if we can minimize sequentially with respect to the parameters of one basis function at a time.

The idea of **forward stagewise additive modeling** is to minimize (11) by adding additional basis functions *without* modifying previous ones.

This is exactly what boosting does; see (3) and the algorithm that includes it.

7.3. Other issues in boosting

What does AdaBoost minimize?

According to ESL, it minimizes

$$L(y, f(\mathbf{x})) = \exp(-y f(\mathbf{x})); \quad (12)$$

see Section 10.4.

They show that this is equivalent to minimizing the deviance

$$\ell(y, f(\mathbf{x})) = \log(1 + \exp(-2y f(\mathbf{x}))), \quad (13)$$

which is what a logit model would minimize, except for the factor of 2 that arises because of coding the output as $\{-1, 1\}$ instead of $\{0, 1\}$.

There is a long discussion of robust and non-robust loss functions in Section 10.6.

There is a detailed discussion of boosting trees in Section 10.9.

The important topic of **gradient boosting** is discussed in Section 10.10.

There are some interesting examples in Section 10.14.