

CS725 : Ensemble Models : Boosting and Bagging

Ensemble model or learner : One which combines models for improved predictions compared to each model in the ensemble

Boosting is a very popular ensembling technique

Weak learner is slightly better than chance

Main idea behind boosting : Can we combine weak learners to improve performance and come up with a strong learner

E.g., Decision stump, Perceptron learner, etc

Many forms of boosting \Rightarrow Adaboost, gradient boosting, XG boost, etc

Template of a boosting algorithm:

- ① Initialize equal weights (that sum to 1) for all the N training points
- ① Train a weak learner on the weighted training data
- ② Compute a weighted error based on the weak learner from ① on the training set
- ③ Compute an importance estimate corresponding to the weak learner
- ④ Recompute weights for all the training examples \uparrow so as to add more wt to misclassified examples using the importance estimate in ③
- ⑤ Repeat steps ② to ④ until stopping criterion is met

ADABOOST (Freund & Schapire, 1995)

Given $\{x_i, y_i\}_{i=1}^n$, $y_i \in \{-1, 1\}$, initial weights $D_1(i) = \frac{1}{n} \forall i$ $D_t(i)$ is normalized to be a prob distribution

for round $t = 1, \dots, T \rightarrow$ hyperparameter

Step 1: Train a weak learner on the trainset weighted by $D_t(i)$

Call it $h_t(x_i) \in \{-1, 1\}$

Step 2: Compute the weighted error with using $h_t(x_i)$

$$\epsilon_t = \sum_{i=1}^n D_t(i) \mathbb{1}_{[h_t(x_i) \neq y_i]}$$

INDICATOR FUNCTION

Step 3: Compute an importance estimate α_t associated with the weak learner $h_t(x)$

$$\alpha_t = \frac{1}{2} \log \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

for rounds $t = 1 \dots T$

Step 4: Recompute the weights of the training examples

$$D_{t+1}(x_i) \propto \begin{cases} D_t(x_i) \exp(-\alpha_t) & \text{if } h_t(x_i) = y_i \\ D_t(x_i) \exp(\alpha_t) & \text{if } h_t(x_i) \neq y_i \end{cases}$$

upweighting
misclassified
examples

$$D_{t+1}(x_i) = \frac{D_t(x_i) \exp(-\alpha_t y_i h_t(x_i))}{\sum_{j=1}^n D_t(x_j) \exp(-\alpha_t y_j h_t(x_j))}$$

downweighting
examples that
 $h_t(x)$ gets correct

Final classifier after boosting: $H(x) = \text{Sign} \left(\underbrace{\sum_{t=1}^T \alpha_t h_t(x)}_{\text{Ensemble classifier weighted by importance estimates } \alpha_t} \right)$
for a test instance x

Ensemble classifier weighted
by importance estimates α_t

Adaboost minimizes the following exponential loss: $L_{\text{boost}} = \frac{1}{n} \sum_i \exp(-y_i H(x_i))$

If I want to add a new weak predictor h to the final classifier i.e., $H \leftarrow H + \alpha h$

the optimal α which minimizes L_{boost} is $\alpha_t = \frac{1}{2} \log \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$

Bagging (Bootstrap Aggregating)

Consider D_1, D_2, \dots, D_m independent datasets sampled i.i.d. from an underlying distribution P .

— Say we train a predictor to classify an x based on $D_i \Rightarrow h(x; D_i)$

An average predictor would just be $\frac{1}{m} \sum_{i=1}^m h(x; D_i)$

What can we say about the following for the average predictor?

① BIAS ?

$$\text{Bias : } \mathbb{E}_{D_1, \dots, D_m} \left[\frac{1}{m} \sum_i h(x; D_i) \right]$$

$$= \frac{1}{m} \sum_i \mathbb{E}_{D_i} [h(x; D_i)]$$

$$= \mathbb{E}_D [h(x; D)]$$

UNCHANGED

② VARIANCE?

$$\text{Var} \left[\frac{1}{m} \sum_{i=1}^m h(x; D_i) \right]$$

$$= \frac{1}{m^2} \sum_{i=1}^m \text{Var} [h(x; D_i)]$$

$$= \frac{1}{m} \text{Var} [h(x; D)]$$

REDUCES by a factor of m !

Catch is that we typically do not have access to independent datasets D_1, \dots, D_m . Bootstrap Aggregating \Rightarrow BAGGING

Solution: Bootstrap! Sample at random with replacement m times from a training set D of size n

x_8 x_7 x_1
 x_6 x_2
 x_5 x_3
 x_4

D_1

D_2

D_3

x_8 x_7 x_7 x_1 x_2 x_4 x_2 x_5
BOOTSTRAP SAMPLE

"

"

train a model h_1
for test pt x

train a model h_2
for test pt x

train a model h_3
for test pt x

for regression

majority voting for classification

$$\frac{1}{m} \sum_{i=1}^m h_i(x)$$

The bootstrap samples are not independent of each other. But the overall variance of the bagged classifier reduces.

Key point: In an ensemble, we want decorrelated predictions for the ensemble to generalize well

Very popular bagging technique is the RANDOM FOREST classifier

Random forests \Rightarrow Decision trees + Bagging + an additional tweak

① Bootstrap from the original dataset D , $|D|=n$, and train

DTs one for each bootstrap sample

② Tweak is each DT only uses a subset of features at each split

[of size $k \approx \sqrt{d}$]

where d is the size of the original feature space

Consider a bootstrap sample. Any training instance that is not part of the bootstrap sample is called an OOB (out-of-bag) sample.

OOB samples can be used as a bagged estimate of a test error without using an explicit test set. OOB error estimate

- How?
- ① For every i^{th} training instance, predict labels from every bootstrap model for which this instance was OOB
 - ② Compute the mean (or maj vote) across these predictions & compute its error
 - ③ Average errors across all the training instances