

CSC 4020 Fundamentals of Machine Learning: Boosting

Baoyuan Wu

March 31

Boosting

Recall that an *ensemble* is a set of predictors whose individual decisions are combined in some way to classify new examples.

Bagging: Train classifiers independently on random subsets of the training data.

Boosting: Train classifiers sequentially, each time focusing on training data points that were previously misclassified.

Let us start with the concept of **weak learner/classifier** (or base classifiers).

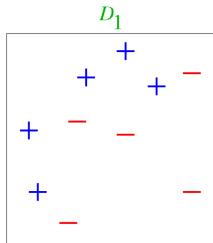
Weak Learner/Classifier

(Informal) Weak learner is a learning algorithm that outputs a hypothesis (e.g., a classifier) that performs slightly better than chance, e.g., it predicts the correct label with probability 0.6.

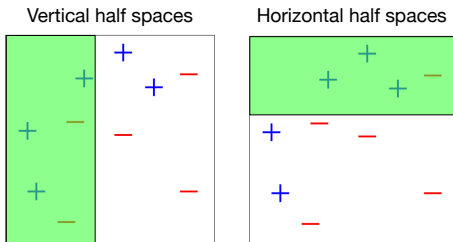
We are interested in weak learners that are *computationally* efficient.

- ▶ Decision trees
- ▶ Even simpler: **Decision Stump**: A decision tree with only a single split

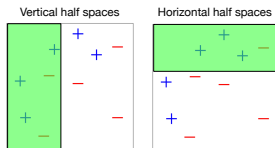
Weak Classifiers



These weak classifiers, which are decision stumps, consist of the set of horizontal and vertical half spaces.



Weak Classifiers



A *single* weak classifier is not capable of making the training error very small. It only performs slightly better than chance, i.e., the error of classifier h according to the given weights $\mathbf{w} = (w_1, \dots, w_N)$ (with $\sum_{i=1}^N w_i = 1$ and $w_i \geq 0$)

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

is at most $\frac{1}{2} - \gamma$ for some $\gamma > 0$.

Can we combine a set of weak classifiers in order to make a better ensemble of classifiers?

Boosting: Train classifiers sequentially, each time focusing on training data points that were previously misclassified.

AdaBoost (Adaptive Boosting)

Key steps of AdaBoost:

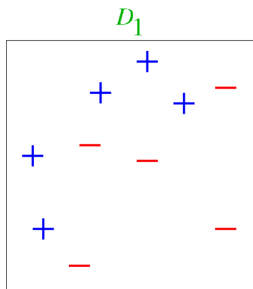
1. At each iteration we re-weight the training samples by **assigning larger weights** to samples (i.e., data points) that were **classified incorrectly**.
2. We train a new weak classifier based on the **re-weighted samples**.
3. We **add this weak classifier to the ensemble** of classifiers. This is our new classifier.
4. Weight each weak classifier in the ensemble with some weights.
5. We repeat the process many times.

The weak learner needs to minimize weighted error.

AdaBoost reduces **bias** by making each classifier focus on previous mistakes.

AdaBoost Example

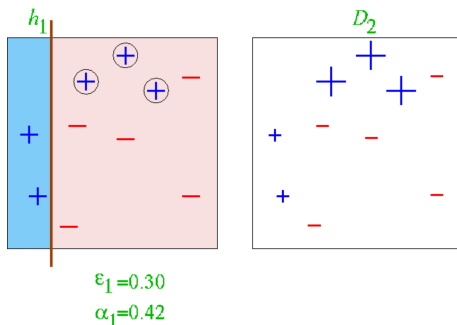
Training data



AdaBoost Example

Round 1

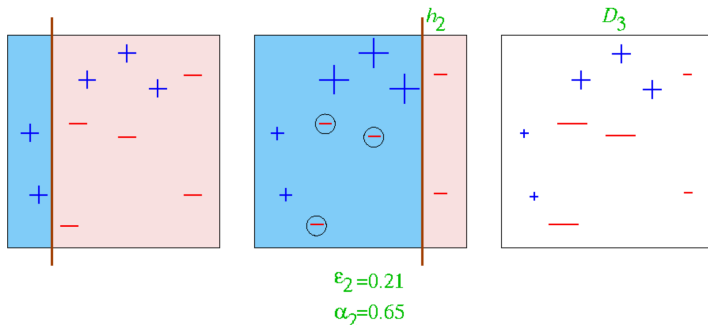
ϵ : Training error, α : Weighting of the current tree.



$$\mathbf{w} = \left(\frac{1}{10}, \dots, \frac{1}{10} \right) \Rightarrow \text{Train a classifier (using } \mathbf{w} \text{)} \Rightarrow \text{err}_1 = \frac{\sum_{i=1}^{10} w_i \mathbb{I}\{h_1(\mathbf{x}^{(i)}) \neq t^{(i)}\}}{\sum_{i=1}^N w_i} = \frac{3}{10}$$
$$\Rightarrow \alpha_1 = \frac{1}{2} \log \frac{1 - \text{err}_1}{\text{err}_1} = \frac{1}{2} \log \left(\frac{1}{0.3} - 1 \right) \approx 0.42 \Rightarrow H(\mathbf{x}) = \text{sign}(\alpha_1 h_1(\mathbf{x}))$$

AdaBoost Example

Round 2

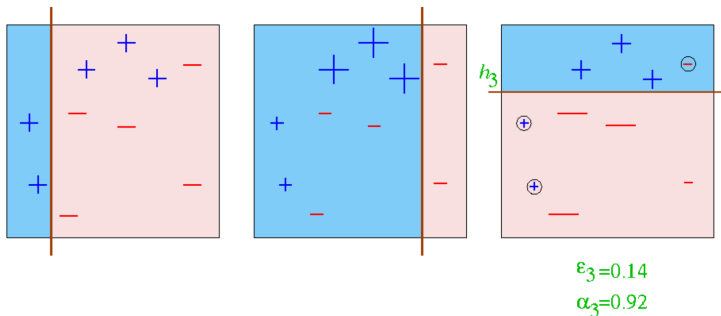


$$\mathbf{w} = \text{updated weights} \Rightarrow \text{Train a classifier (using } \mathbf{w} \text{)} \Rightarrow \text{err}_2 = \frac{\sum_{i=1}^{10} w_i \mathbb{I}\{h_2(\mathbf{x}^{(i)}) \neq t^{(i)}\}}{\sum_{i=1}^N w_i} = 0.21$$

$$\Rightarrow \alpha_2 = \frac{1}{2} \log \frac{1 - \text{err}_3}{\text{err}_3} = \frac{1}{2} \log \left(\frac{1}{0.21} - 1 \right) \approx 0.66 \Rightarrow H(\mathbf{x}) = \text{sign}(\alpha_1 h_1(\mathbf{x}) + \alpha_2 h_2(\mathbf{x}))$$

AdaBoost Example

Round 3

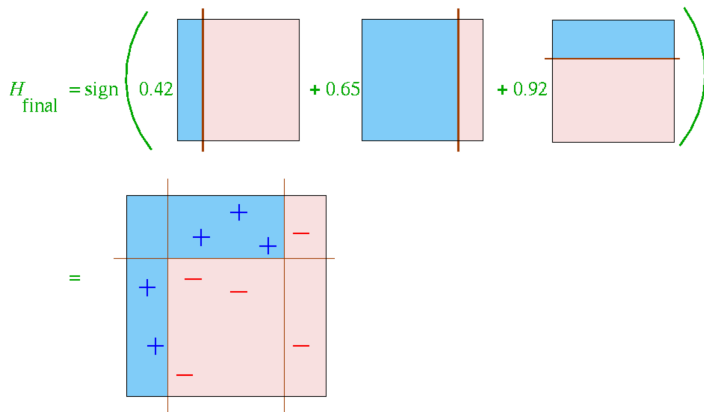


$$\mathbf{w} = \text{updated weights} \Rightarrow \text{Train a classifier (using } \mathbf{w} \text{)} \Rightarrow \text{err}_3 = \frac{\sum_{i=1}^{10} w_i \mathbb{I}\{h_3(\mathbf{x}^{(i)}) \neq t^{(i)}\}}{\sum_{i=1}^N w_i} = 0.14$$

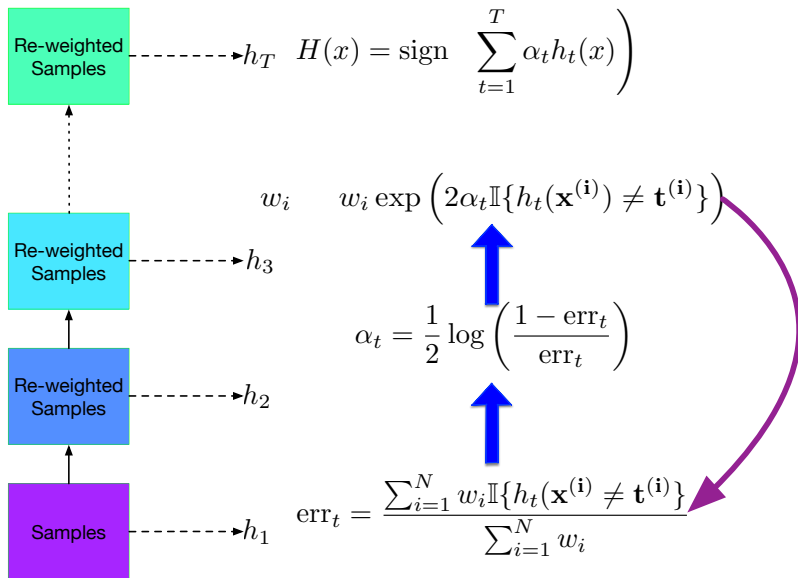
$$\Rightarrow \alpha_3 = \frac{1}{2} \log \frac{1 - \text{err}_3}{\text{err}_3} = \frac{1}{2} \log \left(\frac{1}{0.14} - 1 \right) \approx 0.91 \Rightarrow H(\mathbf{x}) = \text{sign}(\alpha_1 h_1(\mathbf{x}) + \alpha_2 h_2(\mathbf{x}) + \alpha_3 h_3(\mathbf{x}))$$

AdaBoost Example

Final classifier



AdaBoost Algorithm



AdaBoost Algorithm

Input: Data $\mathcal{D}_N = \{\mathbf{x}^{(i)}, t^{(i)}\}_{i=1}^N$, weak classifier *WeakLearn* (a classification procedure that return a classifier from base hypothesis space \mathcal{H} with $h : \mathbf{x} \rightarrow \{-1, +1\}$ for $h \in \mathcal{H}$), number of iterations T

Output: Classifier $H(\mathbf{x})$

Initialize sample weights: $w_i = \frac{1}{N}$ for $i = 1, \dots, N$

For $t = 1, \dots, T$

- Fit a classifier to data using weighted samples ($h_t \leftarrow \text{WeakLearn}(\mathcal{D}_N, \mathbf{w})$), e.g.,

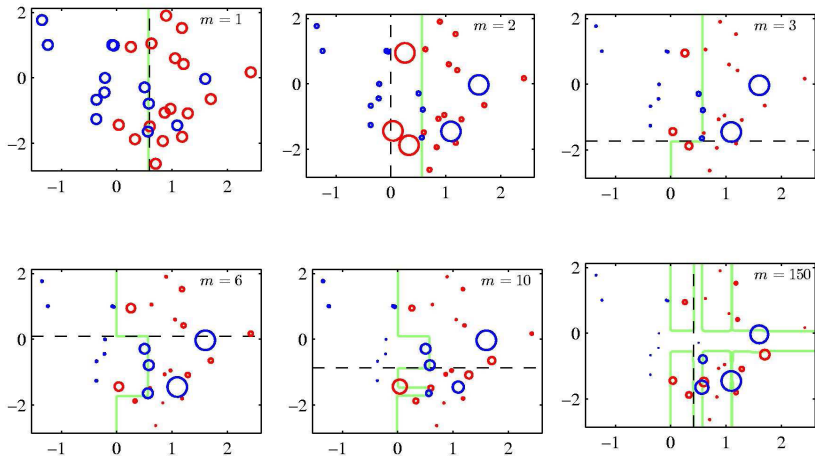
$$h_t \leftarrow \operatorname{argmin}_{h \in \mathcal{H}} \sum_{i=1}^N w_i \mathbb{I}\{h(\mathbf{x}^{(i)}) \neq t^{(i)}\}$$

- Compute weighted error $\text{err}_t = \frac{\sum_{i=1}^N w_i \mathbb{I}\{h_t(\mathbf{x}^{(i)}) \neq t^{(i)}\}}{\sum_{i=1}^N w_i}$
- Compute classifier coefficient $\alpha_t = \frac{1}{2} \log \frac{1 - \text{err}_t}{\text{err}_t}$
- Update data weights

$$w_i \leftarrow w_i \exp \left(-\alpha_t t^{(i)} h_t(\mathbf{x}^{(i)}) \right) \left[\equiv w_i \exp \left(2\alpha_t \mathbb{I}\{h_t(\mathbf{x}^{(i)}) \neq t^{(i)}\} \right) \right]$$

Return $H(\mathbf{x}) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(\mathbf{x}) \right)$

AdaBoost Example



Each figure shows the number m of base learners trained so far, the decision of the most recent learner (dashed black), and the boundary of the ensemble (green)

AdaBoost Minimizes the Training Error

Theorem

Assume that at each iteration of AdaBoost the WeakLearn returns a hypothesis with error $\text{err}_t \leq \frac{1}{2} - \gamma$ for all $t = 1, \dots, T$ with $\gamma > 0$. The training error of the output hypothesis $H(\mathbf{x}) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(\mathbf{x}) \right)$ is at most

$$L_N(H) = \frac{1}{N} \sum_{i=1}^N \mathbb{I}\{H(\mathbf{x}^{(i)}) \neq t^{(i)}\} \leq \exp(-2\gamma^2 T).$$

This is under the simplifying assumption that each weak learner is γ -better than a random predictor.

Analyzing the convergence of AdaBoost is generally difficult.

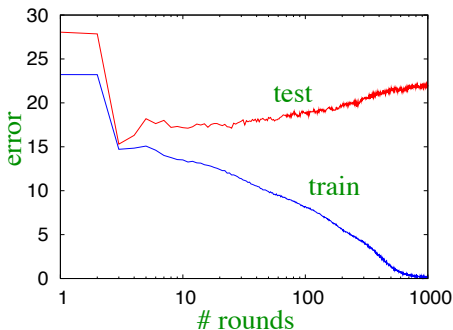
Generalization Error of AdaBoost

AdaBoost's training error (loss) converges to zero. What about the test error of H ?

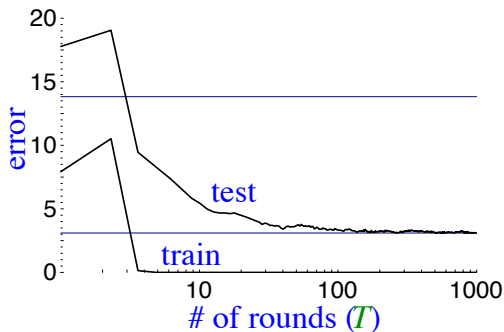
As we add more weak classifiers, the overall classifier H becomes more “complex”.

We expect more complex classifiers overfit.

If one runs AdaBoost long enough, it can in fact overfit.



Generalization Error of AdaBoost



How does that happen?

Summary

Boosting reduces bias by generating an ensemble of weak classifiers.

Each classifier is trained to reduce errors of previous ensemble.

It is quite resilient to overfitting, though it can overfit.

We will later provide a loss minimization viewpoint to AdaBoost. It allows us to derive other boosting algorithms for regression, ranking, etc.

Ensembles Recap

Ensembles combine classifiers to improve performance

Boosting

- ▶ Reduces bias
- ▶ Increases variance (large ensemble can cause overfitting)
- ▶ Sequential
- ▶ High dependency between ensemble elements

Bagging

- ▶ Reduces variance (large ensemble can't cause overfitting)
- ▶ Bias is not changed (much)
- ▶ Parallel
- ▶ Want to minimize correlation between ensemble elements.