# Data Mining and Machine Learning

Madhavan Mukund

Lecture 20, Jan–Apr 2020
https://www.cmi.ac.in/~madhavan/courses/dmml2020jan/

# Limitations of classification models

Recall

- Bias : Expressiveness of model limits classification

- Variance: Variation in model based on sample of training data

Overcoming limitations

- Bagging is an effective way to overcome high variance

  - ▶ Ensemble models

    - ⋆ Sequence of models based on independent bootstrap samples
    - ⋆ Use voting to get an overall classifier

- How can we cope with high bias?

# Dealing with bias

- A biased model always makes mistakes
  - Build an ensemble of models to average out mistakes

- Mistakes should be compensated across models in the ensemble
  - How to build a sequence of models, each biased a different way?
  - Again, we assume we have only one set of training data

# Boosting

- Build a sequence of weak classifiers $M_1$, $M_2$, ..., $M_n$ on inputs $D_1$, $D_2$, ..., $D_n$
  - A weak classifier is any classifier that has error rate strictly below 50%

- Each $D_i$ is a weighted variant of original training data $D$
  - Initially all weights equal, $D_1$
  - Going from $D_i$ to $D_{i+1}$ : increase weights where $M_i$ makes mistakes on $D_i$
  - $M_{i+1}$ will compensate for errors of $M_i$

- Also, each model $M_i$ gets a weight $\alpha_i$ based on its accuracy on $D_i$

- Ensemble output
  - Individual classification outcomes are $\{-1, +1\}$
  - Unknown input $x$: ensemble outcome is weighted sum $\sum_{i=1}^{n} \alpha_i M_i(x)$
  - Check if weighted sum is negative/positive

# The boosting algorithm — Adaboost

- Initially, all data items have equal weight

**AdaBoost**($D$, $Y$, BaseLeaner, $k$)

1. Initialize $D_1(w_i) \leftarrow 1/n$ for all $i$;
2. **for** $t = 1$ to $k$ **do**
3.     $f_t \leftarrow$ BaseLearner($D_t$);
4.     $e_t \leftarrow \displaystyle\sum_{i:\, f_t(D_t(\mathbf{x}_i)) \neq y_i} D_t(w_i)$;
5.     **if** $e_t > \frac{1}{2}$ **then**
6.       $k \leftarrow k - 1$;
7.       exit-loop
8.     **else**
9.       $\beta_t \leftarrow e_t / (1 - e_t)$;
10.       $D_{t+1}(w_i) \leftarrow D_t(w_i) \times \begin{cases} \beta_t & \text{if } f_t(D_t(\mathbf{x}_i)) = y_i \\ 1 & \text{otherwise} \end{cases}$;
11.       $D_{t+1}(w_i) \leftarrow \dfrac{D_{t+1}(w_i)}{\sum_{i=1}^{n} D_{t+1}(w_i)}$
12.     **endif**
13. **endfor**
14. $f_{final}(\mathbf{x}) \leftarrow \underset{y \in Y}{\operatorname{argmax}} \displaystyle\sum_{t:\, f_t(\mathbf{x}) = y} \log \frac{1}{\beta_t}$

# The boosting algorithm — Adaboost

- Initially, all data items have equal weight

- Build a new model and compute its weighted error

**AdaBoost**($D$, $Y$, BaseLeaner, $k$)
1. Initialize $D_1(w_i) \leftarrow 1/n$ for all $i$;
2. **for** $t = 1$ to $k$ **do**
3. $\quad$ $f_t \leftarrow$ BaseLearner($D_t$);
4. $\quad$ $e_t \leftarrow \sum_{i:f_t(D_t(\mathbf{x}_i)) \neq y_i} D_t(w_i)$;
5. $\quad$ **if** $e_t > \frac{1}{2}$ **then**
6. $\quad\quad$ $k \leftarrow k - 1$;
7. $\quad\quad$ exit-loop
8. $\quad$ **else**
9. $\quad\quad$ $\beta_t \leftarrow e_t / (1 - e_t)$;
10. $\quad\quad$ $D_{t+1}(w_i) \leftarrow D_t(w_i) \times \begin{cases} \beta_t & \text{if } f_t(D_t(\mathbf{x}_i)) = y_i \\ 1 & \text{otherwise} \end{cases}$;
11. $\quad\quad$ $D_{t+1}(w_i) \leftarrow \dfrac{D_{t+1}(w_i)}{\sum_{i=1}^{n} D_{t+1}(w_i)}$
12. $\quad$ **endif**
13. **endfor**
14. $f_{final}(\mathbf{x}) \leftarrow \underset{y \in Y}{\mathrm{argmax}} \sum_{t:f_t(\mathbf{x})=y} \log \dfrac{1}{\beta_t}$

# The boosting algorithm — Adaboost

- Initially, all data items have equal weight

- Build a new model and compute its weighted error

- Discard if error rate is above 50%

**AdaBoost**($D$, $Y$, BaseLeaner, $k$)

1.    Initialize $D_1(w_i) \leftarrow 1/n$ for all $i$;
2.    **for** $t = 1$ to $k$ **do**
3.      $f_t \leftarrow$ BaseLearner($D_t$);
4.      $e_t \leftarrow \displaystyle\sum_{i: f_t(D_t(\mathbf{x}_i)) \neq y_i} D_t(w_i)$ ;
5.      **if** $e_t > \frac{1}{2}$ **then**
6.       $k \leftarrow k - 1$;
7.       exit-loop
8.      **else**
9.       $\beta_t \leftarrow e_t / (1 - e_t)$;
10.     $D_{t+1}(w_i) \leftarrow D_t(w_i) \times \begin{cases} \beta_t & \text{if } f_t(D_t(\mathbf{x}_i)) = y_i \\ 1 & \text{otherwise} \end{cases}$ ;
11.     $D_{t+1}(w_i) \leftarrow \dfrac{D_{t+1}(w_i)}{\sum_{i=1}^{n} D_{t+1}(w_i)}$
12.     **endif**
13.  **endfor**
14. $f_{final}(\mathbf{x}) \leftarrow \displaystyle\operatorname*{argmax}_{y \in Y} \sum_{t: f_t(\mathbf{x}) = y} \log \frac{1}{\beta_t}$

# The boosting algorithm — Adaboost

- Initially, all data items have equal weight

- Build a new model and compute its weighted error

- Discard if error rate is above 50%

- Damping factor — reduce weight of correct inputs

**AdaBoost**($D$, $Y$, BaseLeaner, $k$)
1.    Initialize $D_1(w_i) \leftarrow 1/n$ for all $i$;
2.    **for** $t = 1$ to $k$ **do**
3.      $f_t \leftarrow$ BaseLearner($D_t$);
4.      $e_t \leftarrow \sum_{i: f_t(D_t(\mathbf{x}_i)) \neq y_i} D_t(w_i)$;
5.      **if** $e_t > \frac{1}{2}$ **then**
6.        $k \leftarrow k - 1$;
7.        exit-loop
8.      **else**
9.        $\boxed{\beta_t \leftarrow e_t / (1 - e_t);}$
10.     $D_{t+1}(w_i) \leftarrow D_t(w_i) \times \begin{cases} \beta_t & \text{if } f_t(D_t(\mathbf{x}_i)) = y_i \\ 1 & \text{otherwise} \end{cases}$;
11.     $D_{t+1}(w_i) \leftarrow \dfrac{D_{t+1}(w_i)}{\sum_{i=1}^{n} D_{t+1}(w_i)}$
12.    **endif**
13. **endfor**
14. $f_{final}(\mathbf{x}) \leftarrow \underset{y \in Y}{\operatorname{argmax}} \sum_{t: f_t(\mathbf{x}) = y} \log \dfrac{1}{\beta_t}$

# The boosting algorithm — Adaboost

- Initially, all data items have equal weight

- Build a new model and compute its weighted error

- Discard if error rate is above 50%

- Damping factor — reduce weight of correct inputs

- Reweight data items and normalize

**AdaBoost**($D$, $Y$, BaseLeaner, $k$)
1. Initialize $D_1(w_i) \leftarrow 1/n$ for all $i$;
2. **for** $t = 1$ to $k$ **do**
3.     $f_t \leftarrow$ BaseLearner($D_t$);
4.     $e_t \leftarrow \displaystyle\sum_{i: f_t(D_t(\mathbf{x}_i)) \neq y_i} D_t(w_i)$;
5.     **if** $e_t > \frac{1}{2}$ **then**
6.       $k \leftarrow k - 1$;
7.       exit-loop
8.     **else**
9.       $\beta_t \leftarrow e_t / (1 - e_t)$;
10.     $D_{t+1}(w_i) \leftarrow D_t(w_i) \times \begin{cases} \beta_t & \text{if } f_t(D_t(\mathbf{x}_i)) = y_i \\ 1 & \text{otherwise} \end{cases}$;
11.     $D_{t+1}(w_i) \leftarrow \dfrac{D_{t+1}(w_i)}{\sum_{i=1}^{n} D_{t+1}(w_i)}$
12.     **endif**
13. **endfor**
14. $f_{final}(\mathbf{x}) \leftarrow \displaystyle\operatorname*{argmax}_{y \in Y} \sum_{t: f_t(\mathbf{x}) = y} \log \frac{1}{\beta_t}$

# The boosting algorithm — Adaboost

- Initially, all data items have equal weight

- Build a new model and compute its weighted error

- Discard if error rate is above 50%

- Damping factor — reduce weight of correct inputs

- Reweight data items and normalize

- Verdict: weighted sum of individual scores — weights derived from error rate

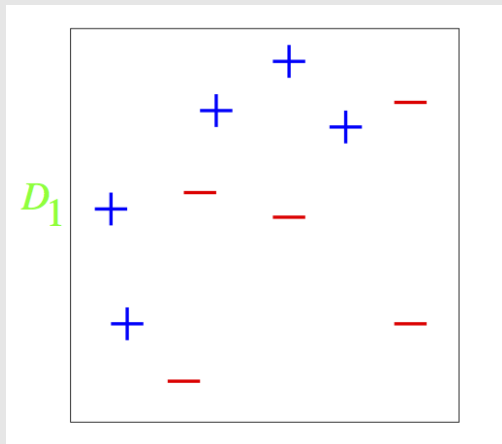**AdaBoost**($D$, $Y$, BaseLeaner, $k$)
1.    Initialize $D_1(w_i) \leftarrow 1/n$ for all $i$;
2.    **for** $t = 1$ to $k$ **do**
3.      $f_t \leftarrow \text{BaseLearner}(D_t)$;
4.      $e_t \leftarrow \sum\limits_{i: f_t(D_t(\mathbf{x}_i)) \neq y_i} D_t(w_i)$;
5.      **if** $e_t > \frac{1}{2}$ **then**
6.        $k \leftarrow k - 1$;
7.        exit-loop
8.      **else**
9.        $\beta_t \leftarrow e_t / (1 - e_t)$;
10.      $D_{t+1}(w_i) \leftarrow D_t(w_i) \times \begin{cases} \beta_t & \text{if } f_t(D_t(\mathbf{x}_i)) = y_i \\ 1 & \text{otherwise} \end{cases}$;
11.      $D_{t+1}(w_i) \leftarrow \dfrac{D_{t+1}(w_i)}{\sum_{i=1}^{n} D_{t+1}(w_i)}$
12.      **endif**
13. **endfor**
14. $f_{final}(\mathbf{x}) \leftarrow \underset{y \in Y}{\operatorname{argmax}} \sum\limits_{t: f_t(\mathbf{x}) = y} \log \dfrac{1}{\beta_t}$

# The boosting algorithm — Adaboost

- Each $M_i$ could be a different type of model

- Can we pick best $n$ out of $N$ weak classifiers?

- Initially all data items have equal weight, select $M_1$ as model with lowest error rate among $N$ candidates

- Inductively, assume we have selected $M_1, \ldots M_j$, with model weights $\alpha_1, \ldots, \alpha_j$, and dataset is updated with new weights as $D_{j+1}$
  - Pick model with lowest error rate on $D_{j+1}$ as $M_{j+1}$
  - Calculate $\alpha_{j+1}$ based on error rate of $M_{j+1}$
  - Reweight all training data based on error rate of $M_{j+1}$

- Note that same model $M$ may be picked in multiple iterations, assigned different weights $\alpha$

# Boosting: An example

- Weak classifiers are horizontal and vertical lines

- Initial training data has equal weights

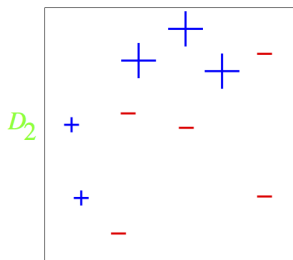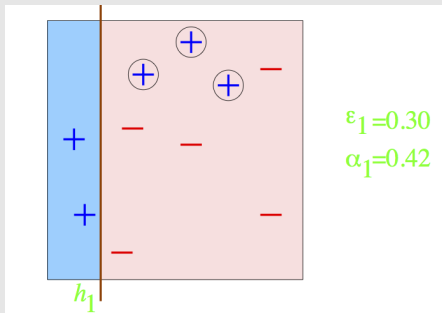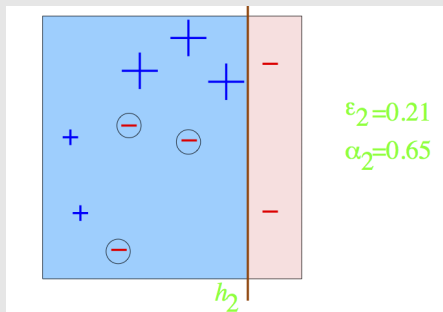# Boosting: An example

- Weak classifiers are horizontal and vertical lines

- Initial training data has equal weights

- First separator: vertical line



$\varepsilon_1 = 0.30$
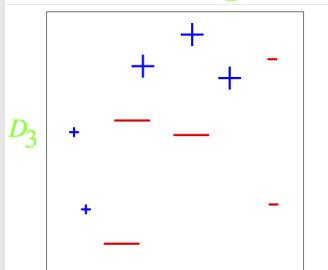$\alpha_1 = 0.42$

$h_1$

# Boosting: An example

- Weak classifiers are horizontal and vertical lines

- Initial training data has equal weights

- First separator: vertical line
  - Increase weight of misclassified inputs



$\varepsilon_1 = 0.30$
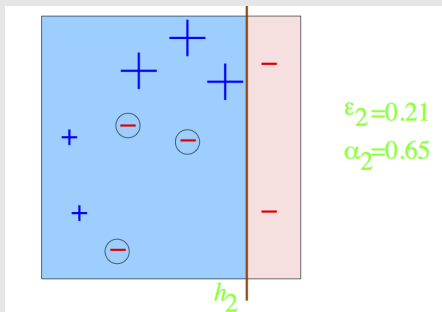$\alpha_1 = 0.42$

$h_1$

$D_2$

# Boosting: An example

- Weak classifiers are horizontal and vertical lines

- Initial training data has equal weights

- First separator: vertical line
  - Increase weight of misclassified inputs

- Second separator: vertical line
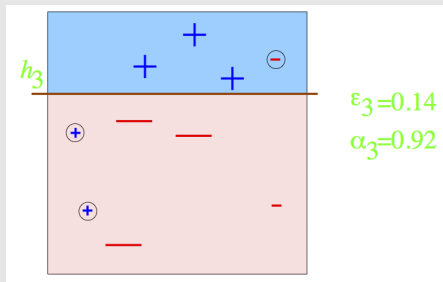


$\varepsilon_2=0.21$
$\alpha_2=0.65$

$h_2$

# Boosting: An example

- Weak classifiers are horizontal and vertical lines

- Initial training data has equal weights

- First separator: vertical line
  - Increase weight of misclassified inputs

- Second separator: vertical line
  - Increase weight of misclassified inputs



$\varepsilon_2 = 0.21$
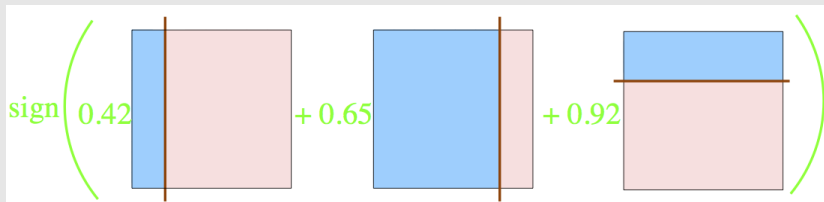
$\alpha_2 = 0.65$

$h_2$

$D_3$

# Boosting: An example

- Weak classifiers are horizontal and vertical lines

- Initial training data has equal weights

- First separator: vertical line
  - Increase weight of misclassified inputs

- Second separator: vertical line
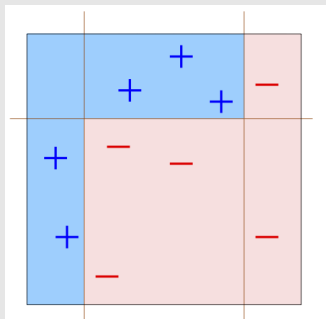  - Increase weight of misclassified inputs

- Third separator: horizontal line



$\varepsilon_3 = 0.14$
$\alpha_3 = 0.92$

# Boosting: An example

- Final classifier is weighted sum of three weak classifiers



$$\text{sign}\left( 0.42 \quad + 0.65 \quad + 0.92 \right)$$

- Pictorially

Given a sample $S$ of $n$ labeled examples $\mathbf{x}_1, \ldots, \mathbf{x}_n$, initialize each example $\mathbf{x}_i$ to have a weight $w_i = 1$. Let $\mathbf{w} = (w_1, \ldots, w_n)$.

For $t = 1, 2, \ldots, t_0$ do

Call the weak learner on the weighted sample $(S, \mathbf{w})$, receiving hypothesis $h_t$.

Multiply the weight of each example that was misclassified by $h_t$ by $\alpha = \frac{\frac{1}{2} + \gamma}{\frac{1}{2} - \gamma}$. Leave the other weights as they are.

End

Output the classifier MAJ$(h_1, \ldots, h_{t_0})$ which takes the majority vote of the hypotheses returned by the weak learner. Assume $t_0$ is odd so there is no tie.

Given a sample $S$ of $n$ labeled examples $\mathbf{x}_1, \ldots, \mathbf{x}_n$, initialize each example $\mathbf{x}_i$ to have a weight $w_i = 1$. Let $\mathbf{w} = (w_1, \ldots, w_n)$.

For $t = 1, 2, \ldots, t_0$ do

Call the weak learner on the weighted sample $(S, \mathbf{w})$, receiving hypothesis $h_t$.

Multiply the weight of each example that was misclassified by $h_t$ by $\alpha = \frac{\frac{1}{2} + \gamma}{\frac{1}{2} - \gamma}$. Leave the other weights as they are.

# Theoretical analysis — simplified Adaboost

$\gamma$-weak classifier

> *For any training dataset $D$, for any assignment of non-negative real weights $w_i$ to $x_i \in D$, classifier correctly labels subset with weight at least $(\frac{1}{2} + \gamma) \sum_{i=1}^{n} w_i$*

### Theorem

Let $A$ be a $\gamma$-weak classification algorithm. Let $D$ be a training dataset of size $n$. Then $t_0 = O(\frac{1}{\gamma^2} \ln n)$ is sufficient for the boosted classifier $MAJ(h_1, h_2, \ldots, h_{t_0})$ to have training error zero.

For any training dataset $D$, a $\gamma$-weak classifier can be boosted to make correct predictions on all of $D$.

# Theoretical analysis — simplified Adaboost

## Proof

- Let $m$ be the number of examples misclassified by final classifier
  - Majority classifier — each such item misclassified at least $t_0/2$ times
  - Each item has weight at least $\alpha^{t_0/2}$
  - Total weight of misclassified items at least $m\alpha^{t_0/2}$

- Iteration $t+1$, only items misclassifed by $h_t$ are reweighted
  - $\gamma$-weak classifier — total weight of misclassified items at most $\frac{1}{2} - \gamma$

- $weight(t)$ – total weight at time $t$
  - Recall that $\alpha = \frac{\frac{1}{2}+\gamma}{\frac{1}{2}-\gamma}$

$$weight(t+1) \leq \left( \underbrace{\alpha\left(\frac{1}{2} - \gamma\right)}_{\text{misclassified}} + \underbrace{\left(\frac{1}{2} + \gamma\right)}_{\text{correct}} \right) weight(t) \leq (1+2\gamma)weight(t)$$

# Theoretical analysis — simplified Adaboost

**Proof**

- $weight(0) = n$, so after $t_0$ iterations, $weight(t_0) \leq n(1 + 2\gamma)^{t_0}$

- Total weight of misclassified items at least $m\alpha^{t_0/2}$,

  so $m\alpha^{t_0/2} \leq n(1 + 2\gamma)^{t_0}$

- Rewrite $\alpha = \frac{\frac{1}{2} + \gamma}{\frac{1}{2} - \gamma}$ as $\alpha = \frac{1 + 2\gamma}{1 - 2\gamma}$

- $m\alpha^{t_0/2} \leq n(1 + 2\gamma)^{t_0} \Rightarrow m \leq n(1 + 2\gamma)^{t_0}\dfrac{(1 - 2\gamma)^{t_0/2}}{(1 + 2\gamma)^{t_0/2}}$

  $\Rightarrow m \leq n(1 - 2\gamma)^{t_0/2}(1 + 2\gamma)^{t_0/2} \Rightarrow m \leq n(1 - 4\gamma^2)^{t_0/2} \Rightarrow$
  $m \leq n(1 - 4\gamma^2)^{t_0/2} \Rightarrow m \leq n(1 - 4\gamma^2)^{t_0/2}$

- Since $1 - x \leq e^{-x}$, $m \leq n\left(e^{-4\gamma^2}\right)^{t_0/2} \Rightarrow m \leq ne^{-2t_0\gamma^2}$

- If $t_0 > \frac{\ln n}{2\gamma^2}$, $m < 1$, so zero training error

# Summary

- Boosting provably improves the performance of a biased classifier

- Adjust weights of incorrectly classified inputs to iteratively produce new classifiers that compensate for errors

- Can also do this to select best $n$ of $N$ diverse classifiers
  - Combining expert advice

- Variation: Combining sleeping experts
  - Each expert (classifier) need not classify all inputs
  - For each input, some experts may be "sleeping"