# Ensemble Learning (continued)

These slides were assembled by Byron Boots based on the slides assembled by Eric Eaton, with grateful acknowledgement of the many others who made their course materials freely available online. Feel free to reuse or adapt these slides for your own academic purposes, provided that you include proper attribution.

Robot Image Credit: Viktoriya Sukhanova © 123RF.com

t = T

- 1: Initialize a vector of n uniform weights  $\mathbf{w}_1$ 2: for  $t = 1, \ldots, T$
- 3: Train model  $h_t$  on X, y with weights  $\mathbf{w}_t$
- 4: Compute the weighted training error of  $h_t$
- 5: Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1 \epsilon_t}{\epsilon_t} \right)$
- 6: Update all instance weights:  $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i))$

Normalize  $\mathbf{w}_{t+1}$  to be a distribution

8: end for

7:

9: **Return** the hypothesis

$$H(\mathbf{x}) = \operatorname{sign}\left(\sum_{t=1}^{T} \beta_t h_t(\mathbf{x})\right)$$



- Final model is a weighted combination of members
  - Each member weighted by its importance

[Freund & Schapire, 1997]

**INPUT:** training data  $X, y = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , the number of iterations T

2: for 
$$t = 1, ..., T$$

- 3: Train model  $h_t$  on X, y with instance weights  $\mathbf{w}_t$
- 4: Compute the weighted training error rate of  $h_t$ :

$$\epsilon_t = \sum_{i: y_i \neq h_t(\mathbf{x}_i)} w_{t,i}$$

- 5: Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1 \epsilon_t}{\epsilon_t} \right)$
- 6: Update all instance weights:

$$w_{t+1,i} = w_{t,i} \exp\left(-\beta_t y_i h_t(\mathbf{x}_i)\right) \quad \forall i = 1, \dots, n$$

7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution:

$$w_{t+1,i} = \frac{w_{t+1,i}}{\sum_{j=1}^{n} w_{t+1,j}} \quad \forall i = 1, \dots, n$$

- 8: end for
- 9: **Return** the hypothesis

$$H(\mathbf{x}) = \operatorname{sign}\left(\sum_{t=1}^{T} \beta_t h_t(\mathbf{x})\right)$$

**INPUT:** training data  $X, y = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , the number of iterations T1: Initialize a vector of n uniform weights  $\mathbf{w}_1 = \begin{bmatrix} \frac{1}{n}, \dots, \frac{1}{n} \end{bmatrix}$ 2: **for**  $t = 1, \dots, T$ 3: Train model  $h_t$  on X, y with instance weights  $\mathbf{w}_t$ 4: Compute the weighted training error rate of  $h_t$ :  $\epsilon_t = \sum w_{t,i}$  **W**<sub>t</sub>

$$-\sum_{i:y_i\neq h_t(\mathbf{x}_i)} w_t$$

- 5: Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1 \epsilon_t}{\epsilon_t} \right)$
- 6: Update all instance weights:  $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i)) \quad \forall i = 1, ..., n$
- 7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution:

$$w_{t+1,i} = \frac{w_{t+1,i}}{\sum_{j=1}^{n} w_{t+1,j}} \quad \forall i = 1, \dots, n$$

- 8: end for
- 9: **Return** the hypothesis

$$H(\mathbf{x}) = \operatorname{sign}\left(\sum_{t=1}^{T} \beta_t h_t(\mathbf{x})\right)$$

w<sub>t</sub> is a vector of weights over the instances at iteration t

All points start with equal weight

**INPUT:** training data  $X, y = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , the number of iterations T1: Initialize a vector of n uniform weights  $\mathbf{w}_1 = \begin{bmatrix} \frac{1}{n}, \dots, \frac{1}{n} \end{bmatrix}$ 

2: for 
$$t = 1, ..., T$$

- 3: Train model  $h_t$  on X, y with instance weights  $\mathbf{w}_t$
- 4: Compute the weighted training error rate of  $h_t$ :

$$\epsilon_t = \sum_{i: y_i \neq h_t(\mathbf{x}_i)} w_{t,i}$$

5: Choose 
$$\beta_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$$

6: Update all instance weights:  $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i))$ 

We weight instances differently when learning the model, either in the cost function or by bootstrap replication.

7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution:

$$w_{t+1,i} = \frac{w_{t+1,i}}{\sum_{j=1}^{n} w_{t+1,j}} \quad \forall i = 1, \dots, n$$

- 8: end for
- 9: **Return** the hypothesis

$$H(\mathbf{x}) = \operatorname{sign}\left(\sum_{t=1}^{T} \beta_t h_t(\mathbf{x})\right)$$

#### **Base Learner Requirements**

- AdaBoost works best with "weak" learners
  - Should not be complex
  - Typically high bias classifiers
  - Works even when weak learner has an error rate just slightly under 0.5 (i.e., just slightly better than random)
    - Can prove training error goes to 0 in O(log n) iterations
- Examples:
  - Decision stumps (1 level decision trees)
  - Depth-limited decision trees
  - Linear classifiers

**INPUT:** training data  $X, y = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ,

the number of iterations T

1: Initialize a vector of n uniform weights  $\mathbf{w}_1 = \begin{bmatrix} \frac{1}{n}, \dots, \frac{1}{n} \end{bmatrix}$ 2: for t - 1 T

- 2: for t = 1, ..., T
- 3: Train model  $h_t$  on X, y with instance weights  $\mathbf{w}_t$
- 4: Compute the weighted training error rate of  $h_t$ :

 $\epsilon_t = \sum_{i: y_i \neq h_t(\mathbf{x}_i)} w_{t,i}$ 

5: Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$ 

6: Update all instance weights:

$$w_{t+1,i} = w_{t,i} \exp\left(-\beta_t y_i h_t(\mathbf{x}_i)\right) \quad \forall i = 1, \dots, n$$

7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution:

$$w_{t+1,i} = \frac{w_{t+1,i}}{\sum_{j=1}^{n} w_{t+1,j}} \quad \forall i = 1, \dots, n$$

- 8: end for
- 9: **Return** the hypothesis

$$H(\mathbf{x}) = \operatorname{sign}\left(\sum_{t=1}^{T} \beta_t h_t(\mathbf{x})\right)$$

**INPUT:** training data  $X, y = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , the number of iterations T1: Initialize a vector of n uniform weights  $\mathbf{w}_1 = \left|\frac{1}{n}, \ldots, \frac{1}{n}\right|$ 2: for t = 1, ..., TTrain model  $h_t$  on X, y with instance weights  $\mathbf{w}_t$ 3: Compute the weighted training error rate of  $h_t$ : 4:  $\epsilon_t = \sum w_{t,i}$  $i: y_i \neq h_t(\mathbf{x}_i)$ Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$ 5:Update all instance weights 6:  $\beta_{t}$  measures the importance of  $h_{t}$  $w_{t+1,i} = w_{t,i} \exp\left(-\beta_t y_i h_t\right)$ If  $\epsilon_t \leq 0.5$ , then  $\beta_t \geq 0$ 7: Normalize  $\mathbf{w}_{t+1}$  to be a dist  $\circ$  Trivial, otherwise flip h<sub>t</sub>'s predictions  $w_{t+1,i} = \frac{w_{t+1,i}}{\sum_{i=1}^{n} w_{t+1,i}} \quad \forall \bullet \quad \beta_t \text{ grows as } h_t \text{'s error shrinks}$ 8: end for 9: **Return** the hypothesis  $H(\mathbf{x}) = \operatorname{sign}\left(\sum_{t=1}^{T} \beta_t h_t(\mathbf{x})\right)$ 



$$H(\mathbf{x}) = \operatorname{sign}\left(\sum_{t=1}^{I} \beta_t h_t(\mathbf{x})\right)$$

**INPUT:** training data  $X, y = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , the number of iterations T1: Initialize a vector of n uniform weights  $\mathbf{w}_1 = \left\lfloor \frac{1}{n}, \ldots, \frac{1}{n} \right\rfloor$ 2: for t = 1, ..., TTrain model  $h_t$  on X, y with instance weights  $\mathbf{w}_t$ 3: Compute the weighted training error rate of  $h_t$ : 4:  $\epsilon_t = \sum w_{t,i}$  $i: y_i \neq h_t(\mathbf{x}_i)$ Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$ Make  $w_{t+1}$  sum to 1 5:6: Update all instance weights:  $w_{t+1,i} = w_{t,i} \exp\left(-\beta_t y_i h_t(\mathbf{x}_i)\right) \quad \forall i = 1, \dots, n$ Normalize  $\mathbf{w}_{t+1}$  to be a distribution: 7:  $w_{t+1,i} = \frac{w_{t+1,i}}{\sum_{i=1}^{n} w_{t+1,i}} \quad \forall i = 1, \dots, n$ 8: end for 9: **Return** the hypothesis  $H(\mathbf{x}) = \operatorname{sign}\left(\sum_{t=1}^{T} \beta_t h_t(\mathbf{x})\right)$ 

**INPUT:** training data  $X, y = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , the number of iterations T

1: Initialize a vector of n uniform weights  $\mathbf{w}_1 = \begin{bmatrix} \frac{1}{n}, \dots, \frac{1}{n} \end{bmatrix}$ 2: for  $t = 1, \dots, T$ 

- 3: Train model  $h_t$  on X, y with instance weights  $\mathbf{w}_t$
- 4: Compute the weighted training error rate of  $h_t$ :

$$\epsilon_t = \sum_{i: y_i \neq h_t(\mathbf{x}_i)} w_{t,i}$$

- 5: Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1 \epsilon_t}{\epsilon_t} \right)$
- 6: Update all instance weights:  $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i)) \quad \forall i = 1,$
- 7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution:

$$w_{t+1,i} = \frac{w_{t+1,i}}{\sum_{j=1}^{n} w_{t+1,j}} \quad \forall i = 1, \dots, n$$

- 8: end for
- 9: **Return** the hypothesis

$$H(\mathbf{x}) = \operatorname{sign}\left(\sum_{t=1}^{T} \beta_t h_t(\mathbf{x})\right)$$

Member classifiers with less error are given more weight in the final ensemble hypothesis

Final prediction is a weighted combination of each member's prediction

### Dynamic Behavior of AdaBoost

- If a point is repeatedly misclassified...
  - Each time, its weight is increased
  - Eventually it will be emphasized enough to generate a hypothesis that correctly predicts it
- Successive member hypotheses focus on the hardest parts of the instance space
  - Instances with highest weight are often outliers

### AdaBoost and Overfitting

- VC Theory predicts that AdaBoost will overfit as the number of weak classifiers T grows large
  - Hypothesis keeps growing more complex
- In practice, AdaBoost often does <u>not</u> overfit, performing significantly better than VC theory suggests
- AdaBoost does not explicitly regularize the model, and yet generalizes well

## Explaining Why AdaBoost Works



- Empirically, boosting resists overfitting
- Note that it continues to drive down the test error even AFTER the training error reaches zero
- Boosting maximizing confidence

### AdaBoost in Practice

Strengths:

- Fast and simple to program
- No parameters to tune (besides T)
- No assumptions on weak learner

#### When boosting can fail:

- Given insufficient data
- Overly complex weak hypotheses
- Can be susceptible to noise
- When there are a large number of outliers

## **Boosted Decision Trees**

- Boosted decision trees are one of the best "off-the-shelf" classifiers

   i.e., no parameter tuning
- Limit member hypothesis complexity by limiting tree depth
- Boosting methods are typically used with trees in practice



"AdaBoost with trees is the best off-the-shelf classifier in the world" -Breiman, 1996 (Also, see results by Caruana & Niculescu-Mizil, ICML 2006)