#### **Ensemble Learning**

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

# **Ensemble Learning**

Consider a set of classifiers  $h_1, ..., h_L$ 

**Idea:** construct a classifier  $H(\mathbf{x})$  that combines the individual decisions of  $h_1, ..., h_L$ 

- e.g., could have the member classifiers vote, or
- e.g., could use different members for different regions of the instance space

## **Ensemble Learning**

Consider a set of classifiers  $h_1, ..., h_L$ 

**Idea:** construct a classifier  $H(\mathbf{x})$  that combines the individual decisions of  $h_1, ..., h_L$ 

- e.g., could have the member classifiers vote, or
- e.g., could use different members for different regions of the instance space

Successful ensembles require diversity

- Classifiers should make different mistakes
- Can have different types of base learners

## **Combining Classifiers: Averaging**



• Final hypothesis is a simple vote of the members

#### Combining Classifiers: Weighted Average



 Coefficients of individual members are trained using a validation set

# **Combining Classifiers: Gating**



- Coefficients of individual members depend on input
- Train gating function via validation set

# **Combining Classifiers: Stacking**



- Predictions of 1<sup>st</sup> layer used as input to 2<sup>nd</sup> layer
- Train 2<sup>nd</sup> layer on validation set

## How to Achieve Diversity

Cause of the Mistake	<b>Diversification Strategy</b>
Pattern was difficult	We will come back to this
Overfitting	Vary the training sets
Some features are noisy	Vary the set of input features

# Manipulating the Training Data

#### **Bootstrap replication:**

- Given n training examples, construct a new training set by sampling *n* instances with replacement
- Excludes ~35% of the training instances

#### **Bootstrap aggregating (Bagging):**

- Create bootstrap replicates of training set
- Train a classifier (e.g., a decision tree) for each replicate
- Estimate classifier performance using out-of-bootstrap data
- Average output of all classifiers

#### Boosting: (in just a minute...)

# Manipulating the Features

#### **Random Forests**

- Construct decision trees on bootstrap replicas
  - Restrict the node decisions to a small subset of features picked randomly for each node
- Do not prune the trees
  - Estimate tree performance on out-of-bootstrap data
- Average the output of all trees (or choose mode decision)

#### Adaptive Boosting (AdaBoost)

[Freund & Schapire, 1997]

- A meta-learning algorithm with great theoretical and empirical performance
- Turns a base learner (i.e., a "weak hypothesis") into a high performance classifier
- Creates an ensemble of weak hypotheses by repeatedly emphasizing mispredicted instances

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:

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

- 7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution 8: end for
- 9: **Return** the hypothesis

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



• Size of point represents the instance's weight

1: Initialize a vector of n uniform weights  $\mathbf{w}_1$ 2: for t = 1, ..., 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:

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

- 7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution 8: end for
- 9: **Return** the hypothesis

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



Initialize a vector of n uniform weights w<sub>1</sub>
for t = 1,...,T
Train model h<sub>t</sub> on X, y with weights w<sub>t</sub>

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

7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution 8: end for

9: **Return** the hypothesis

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



1: Initialize a vector of n uniform weights  $\mathbf{w}_1$ 2: for t = 1, ..., T3: 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:  $w_{t+1,i} = w_{t,i} \exp\left(-\beta_t y_i h_t(\mathbf{x}_i)\right)$ 7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution

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

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



- $\beta_{\rm t}$  measures the importance of  $h_{\rm t}$
- If  $\epsilon_t \leq 0.5$ , then  $\beta_t \geq 0$  (can trivially guarantee)

1: Initialize a vector of n uniform weights  $\mathbf{w}_1$ 2: for t = 1, ..., T3: 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:  $w_{t+1,i} = w_{t,i} \exp\left(-\beta_t y_i h_t(\mathbf{x}_i)\right)$ 7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution

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

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



- $\beta_{\rm t}$  measures the importance of  $h_{\rm t}$
- If  $\epsilon_t \leq 0.5$ , then  $\beta_t \geq 0$  ( $\beta_t$  grows as  $\epsilon_t$  gets smaller)

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



- Weights of correct predictions are multiplied by  $\ e^{-eta_t} \leq 1$
- Weights of incorrect predictions are multiplied by  $\,e^{eta_t}\geq 1$

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



- Weights of correct predictions are multiplied by  $\ e^{-eta_t} \leq 1$
- Weights of incorrect predictions are multiplied by  $e^{eta_t} \geq 1$

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



Disclaimer: Note that resized points in the illustration above are not necessarily to scale with  $\beta_t$ 

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:

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

- 7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution 8: end for
- 9: **Return** the hypothesis

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



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:

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

- 7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution 8: end for
- 9: **Return** the hypothesis

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



#### Training a Model with Weighted Instances

• For algorithms like logistic regression, can simply incorporate weights w into the cost function

 $\boldsymbol{n}$ 

 Essentially, weigh the cost of misclassification differently for each instance

$$J_{\text{reg}}(\boldsymbol{\theta}) = -\sum_{i=1}^{n} w_i \left[ y_i \log h_{\boldsymbol{\theta}}(\mathbf{x}_i) + (1 - y_i) \log \left(1 - h_{\boldsymbol{\theta}}(\mathbf{x}_i)\right) \right] + \lambda \|\boldsymbol{\theta}_{[1:d]}\|_2^2$$

- For algorithms that don't directly support instance weights (e.g., ID3 decision trees, etc.), use weighted bootstrap sampling
  - Form training set by resampling instances with replacement according to w

1: Initialize a vector of n uniform weights  $\mathbf{w}_1$ 2: for t = 1, ..., T3: 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:  $w_{t+1,i} = w_{t,i} \exp\left(-\beta_t y_i h_t(\mathbf{x}_i)\right)$ 7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution

8: end for

9: **Return** the hypothesis

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



1: Initialize a vector of n uniform weights  $\mathbf{w}_1$ 2: for t = 1, ..., T3: 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:  $w_{t+1,i} = w_{t,i} \exp \left(-\beta_t y_i h_t(\mathbf{x}_i)\right)$ 7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution

8: end for

9: **Return** the hypothesis

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

$$t = 2$$

- $\beta_{\rm t}$  measures the importance of  $h_{\rm t}$
- If  $\epsilon_t \leq 0.5$ , then  $\beta_t \geq 0$  ( $\beta_t$  grows as  $\epsilon_t$  gets smaller)



- Weights of correct predictions are multiplied by  $e^{-eta_t} \leq 1$
- Weights of incorrect predictions are multiplied by  $e^{eta_t} \geq 1$



- Weights of correct predictions are multiplied by  $e^{-eta_t} \leq 1$
- Weights of incorrect predictions are multiplied by  $e^{\beta_t} \geq 1$

1: Initialize a vector of n uniform weights  $\mathbf{w}_1$ 2: for t = 1, ..., T3: 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:  $w_{t+1,i} = w_{t,i} \exp\left(-\beta_t y_i h_t(\mathbf{x}_i)\right)$ 

- 7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution 8: end for
- 9: **Return** the hypothesis

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



1: Initialize a vector of n uniform weights  $\mathbf{w}_1$ 2: for t = 1, ..., T3: 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:  $w_{t+1,i} = w_{t,i} \exp\left(-\beta_t y_i h_t(\mathbf{x}_i)\right)$ 7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution

8: end for

9: **Return** the hypothesis

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



1: Initialize a vector of n uniform weights  $\mathbf{w}_1$ 2: for t = 1, ..., T3: 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:  $w_{t+1,i} = w_{t,i} \exp\left(-\beta_t y_i h_t(\mathbf{x}_i)\right)$ 7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution

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

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

$$t = 3$$

- $\beta_{\rm t}$  measures the importance of  $h_{\rm t}$
- If  $\epsilon_t \leq 0.5$ , then  $\beta_t \geq 0$  ( $\beta_t$  grows as  $\epsilon_t$  gets smaller)



- Weights of correct predictions are multiplied by  $e^{-\beta_t} \leq 1$
- Weights of incorrect predictions are multiplied by  $e^{\beta_t} \geq 1$



- Weights of correct predictions are multiplied by  $e^{-eta_t} \leq 1$
- Weights of incorrect predictions are multiplied by  $e^{\beta_t} \geq 1$

1: Initialize a vector of n uniform weights  $\mathbf{w}_1$ 2: for t = 1, ..., T3: 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:  $w_{t+1,i} = w_{t,i} \exp\left(-\beta_t y_i h_t(\mathbf{x}_i)\right)$ 

- 7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution 8: end for
- 9: **Return** the hypothesis

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



1: Initialize a vector of n uniform weights  $\mathbf{w}_1$ 2: for t = 1, ..., T3: 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:  $w_{t+1,i} = w_{t,i} \exp \left(-\beta_t y_i h_t(\mathbf{x}_i)\right)$ 

- 7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution 8: end for
- 9: **Return** the hypothesis

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



1: Initialize a vector of n uniform weights  $\mathbf{w}_1$ 2: for t = 1, ..., T3: 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:  $w_{t+1,i} = w_{t,i} \exp \left(-\beta_t y_i h_t(\mathbf{x}_i)\right)$ 7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution 8: end for

9: **Return** the hypothesis

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



t = T

1: Initialize a vector of n uniform weights  $\mathbf{w}_1$ 

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

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

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

9: **Return** the hypothesis

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



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:  $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\mathbf{x}_i))$
- 7: Normalize  $\mathbf{w}_{t+1}$  to be a distribution 8: end for
- 9: **Return** the hypothesis

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



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