

# Combining Models

Jiachun Jin  
jinjch@shanghaitech.edu.cn

# Outline

- ▶ Committees
- ▶ Boosting

# Committees

- ▶ [Definition] A combinations of models. Practically we train  $M$  different models, make predictions using the **average** of the predictions made by each model.

Predictive model:  $y_m(\mathbf{x}), \forall m \in [M]$

Committee predictive model:  $y_{\text{COM}}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M y_m(\mathbf{x})$

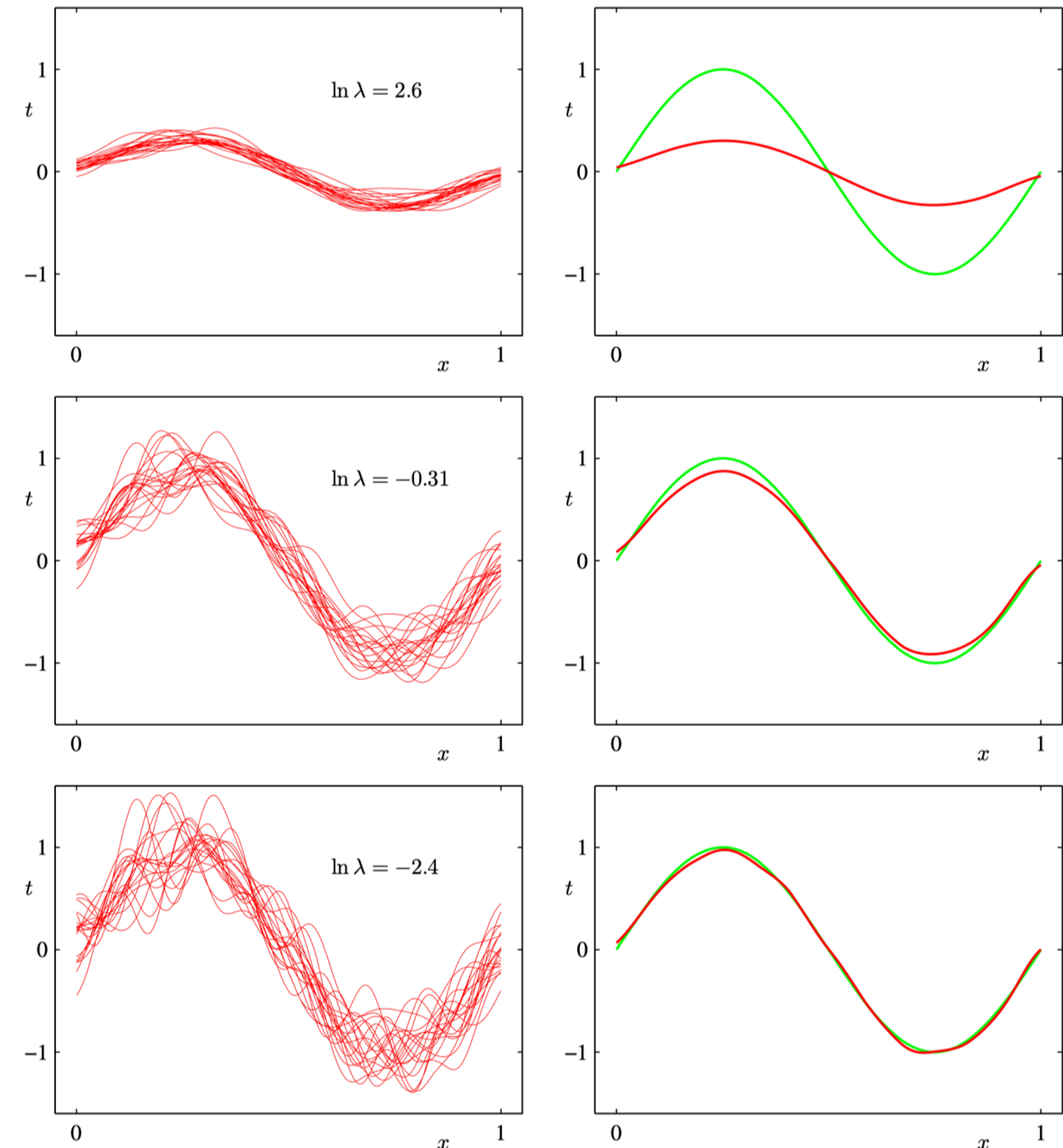
- ▶ An example

# Committees

- ▶ Use polynomial curve to fit sinusoidal function

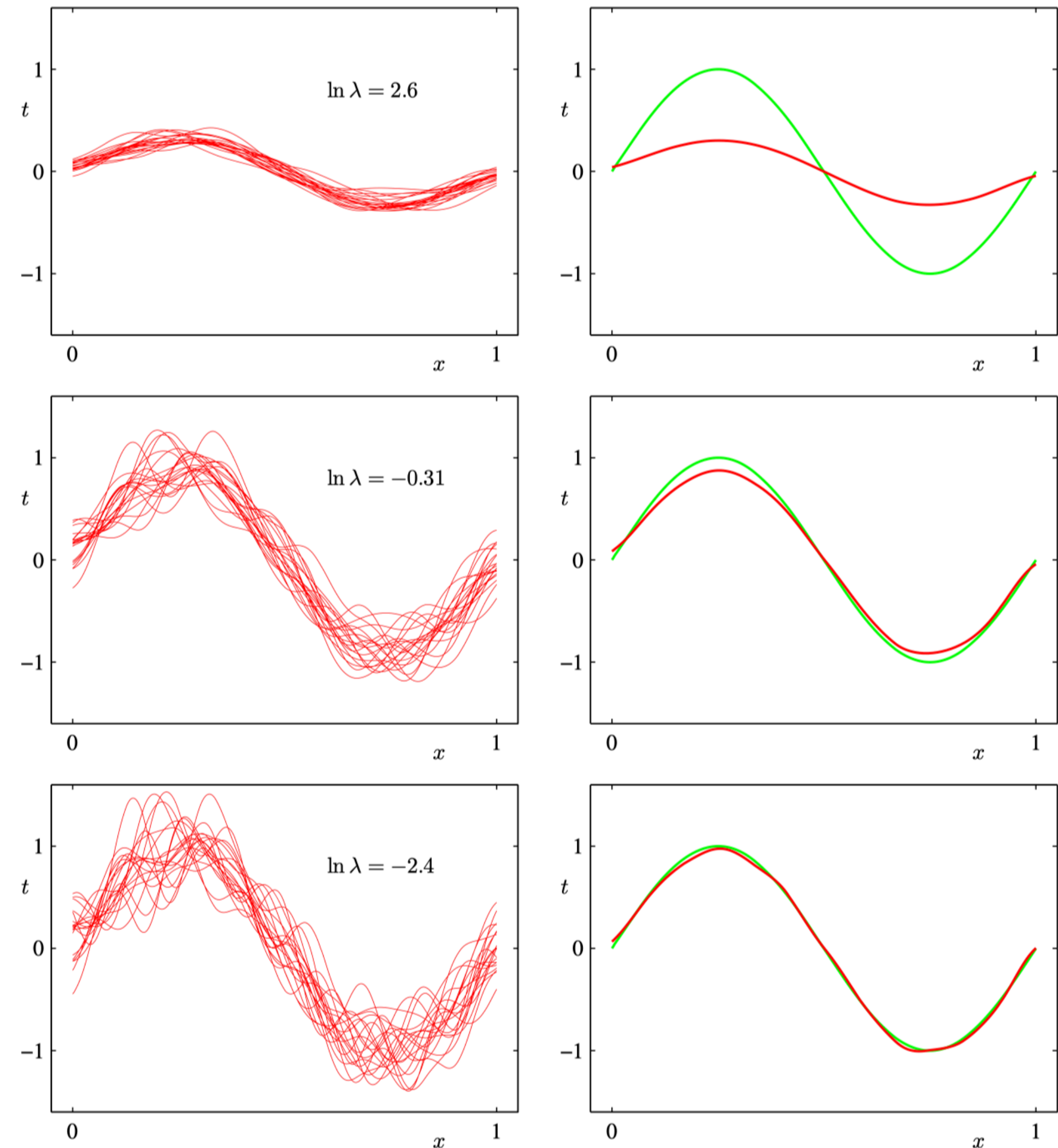
$$\frac{1}{2} \sum_{n=1}^N \left\{ y_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) \right\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- ▶  $L = 100$  data sets, each having  $N = 25$  data points, and  $M = 25$  including the bias term
- ▶ Green line: the sinusoidal function from which the data sets were generated
- ▶ Red line: average of the 100 fits
- ▶ Observation: the contribution arising from the **variance term tended to cancel**, leading to improved predictions
- ▶ Idea: averaged a set of low-bias models (corresponding to higher order polynomials)



# Committees

- ▶ Problem: we have only a single data set
- ▶ Solution: use bootstrap datasets
  - ▶ Suppose our original data set consists of  $N$  data points  
 $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
  - ▶ We can create a new data set  $\mathbf{X}_B$  by drawing  $N$  points at random from  $\mathbf{X}$ , with replacement, some points in  $\mathbf{X}$  may be replicated in  $\mathbf{X}_B$
  - ▶ repeat  $M$  times



# Committees

- ▶ Predictive model:  $y_m(\mathbf{x}), \forall m \in [M]$
- ▶ Committee prediction:  $y_{\text{COM}}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M y_m(\mathbf{x})$
- ▶ True regression function:  $h(\mathbf{x})$
- ▶ Output of each model:  $y_m(\mathbf{x}) = h(\mathbf{x}) + \epsilon_m(\mathbf{x})$
- ▶ Average error made by the models acting individually:  
$$E_{\text{AV}} = \frac{1}{M} \sum_{m=1}^M \mathbb{E}_{\mathbf{x}} [\epsilon_m(\mathbf{x})^2]$$
- ▶ EPE: trade-off between bias and variance:

$$\mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2] = \underbrace{\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2}_{(\text{bias})^2} + \underbrace{\mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2]}_{\text{variance}}$$

- ▶ The expected error from the committee  $y_{\text{COM}}$ :

$$E_{\text{COM}} = \mathbb{E}_{\mathbf{x}} \left[ \left\{ \frac{1}{M} \sum_{m=1}^M y_m(\mathbf{x}) - h(\mathbf{x}) \right\}^2 \right] = \mathbb{E}_{\mathbf{x}} \left[ \left\{ \frac{1}{M} \sum_{m=1}^M \epsilon_m(\mathbf{x}) \right\}^2 \right]$$

- ▶ If we assume:

$$\mathbb{E}_{\mathbf{x}} [\epsilon_m(\mathbf{x})] = 0$$

$$\mathbb{E}_{\mathbf{x}} [\epsilon_m(\mathbf{x}) \epsilon_l(\mathbf{x})] = 0, \quad m \neq l$$

errors due to the individual models are uncorrelated (not true in practice)

- ▶ Then  $E_{\text{COM}} = \frac{1}{M} E_{\text{AV}}$

# Outline

- ▶ Committees
- ▶ Boosting



# Boosting

## AdaBoost

1. Initialize the data weighting coefficients  $\{w_n\}$  by setting  $w_n^{(1)} = 1/N$  for  $n = 1, \dots, N$ .
2. For  $m = 1, \dots, M$ :

(a) Fit a classifier  $y_m(\mathbf{x})$  to the training data by minimizing the weighted error function

$$J_m = \sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n) \quad (14.15)$$

where  $I(y_m(\mathbf{x}_n) \neq t_n)$  is the indicator function and equals 1 when  $y_m(\mathbf{x}_n) \neq t_n$  and 0 otherwise.

(b) Evaluate the quantities

$$\epsilon_m = \frac{\sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}} \quad (14.16)$$

and then use these to evaluate

$$\alpha_m = \ln \left\{ \frac{1 - \epsilon_m}{\epsilon_m} \right\}. \quad (14.17)$$

(c) Update the data weighting coefficients

$$w_n^{(m+1)} = w_n^{(m)} \exp \{ \alpha_m I(y_m(\mathbf{x}_n) \neq t_n) \} \quad (14.18)$$

3. Make predictions using the final model, which is given by

$$Y_M(\mathbf{x}) = \text{sign} \left( \sum_{m=1}^M \alpha_m y_m(\mathbf{x}) \right). \quad (14.19)$$

$$t_n \in \{-1, 1\} \quad y(\mathbf{x}) \in \{-1, 1\}$$

the exponential error function:  $E = \sum_{n=1}^N \exp \{ -t_n f_m(\mathbf{x}_n) \}$

linear combination of base classifiers:  $f_m(\mathbf{x}) = \frac{1}{2} \sum_{l=1}^m \alpha_l y_l(\mathbf{x})$

Q: How do they come up with this idea?

A: Motivated by finding a local optimal of  $E$



# Boosting

(a) Fit a classifier  $y_m(\mathbf{x})$  to the training data by minimizing the weighted error function

$$J_m = \sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n) \quad (14.15)$$

where  $I(y_m(\mathbf{x}_n) \neq t_n)$  is the indicator function and equals 1 when  $y_m(\mathbf{x}_n) \neq t_n$  and 0 otherwise.

the exponential error function:  $E = \sum_{n=1}^N \exp \{ -t_n f_m(\mathbf{x}_n) \}$   
 linear combination of base classifiers:  $f_m(\mathbf{x}) = \frac{1}{2} \sum_{l=1}^m \alpha_l y_l(\mathbf{x})$

$$t_n \in \{-1, 1\} \quad y(\mathbf{x}) \in \{-1, 1\}$$

Goal: minimize  $E$  with respect to both the weighting coefficients  $\alpha_l$  and the parameters of the base classifiers  $y_l(\mathbf{x})$

Suppose:  $y_1(\mathbf{x}), \dots, y_{m-1}(\mathbf{x}), \alpha_1, \dots, \alpha_{m-1}$  are fixed, minimize only with respect to  $y_m(\mathbf{x})$  and  $\alpha_m$

$$\begin{aligned} E &= \sum_{n=1}^N \exp \left\{ -t_n f_{m-1}(\mathbf{x}_n) - \frac{1}{2} t_n \alpha_m y_m(\mathbf{x}_n) \right\} \\ &= \sum_{n=1}^N w_n^{(m)} \exp \left\{ -\frac{1}{2} t_n \alpha_m y_m(\mathbf{x}_n) \right\} \\ &= e^{-\alpha_m/2} \sum_{n \in \mathcal{T}_m} w_n^{(m)} + e^{\alpha_m/2} \sum_{n \in \mathcal{M}_m} w_n^{(m)} \\ &= (e^{\alpha_m/2} - e^{-\alpha_m/2}) \sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n) + e^{-\alpha_m/2} \sum_{n=1}^N w_n^{(m)} \end{aligned}$$

minimize  $E$  w.r.t  $y_m(\mathbf{x})$ :

$$\hat{y}_m = \arg \min_{y_m} \sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)$$

# Boosting

## AdaBoost

1. Initialize the data weighting coefficients  $\{w_n\}$  by setting  $w_n^{(1)} = 1/N$  for  $n = 1, \dots, N$ .
2. For  $m = 1, \dots, M$ :

- (a) Fit a classifier  $y_m(\mathbf{x})$  to the training data by minimizing the weighted error function

$$J_m = \sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n) \quad (14.15)$$

where  $I(y_m(\mathbf{x}_n) \neq t_n)$  is the indicator function and equals 1 when  $y_m(\mathbf{x}_n) \neq t_n$  and 0 otherwise.

- (b) Evaluate the quantities

$$\epsilon_m = \frac{\sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}} \quad (14.16)$$

and then use these to evaluate

$$\alpha_m = \ln \left\{ \frac{1 - \epsilon_m}{\epsilon_m} \right\}. \quad (14.17)$$

- (c) Update the data weighting coefficients

$$w_n^{(m+1)} = w_n^{(m)} \exp \{ \alpha_m I(y_m(\mathbf{x}_n) \neq t_n) \} \quad (14.18)$$

3. Make predictions using the final model, which is given by

$$Y_M(\mathbf{x}) = \text{sign} \left( \sum_{m=1}^M \alpha_m y_m(\mathbf{x}) \right). \quad (14.19)$$

$$t_n \in \{-1, 1\} \quad y(\mathbf{x}) \in \{-1, 1\}$$

the exponential error function:  $E = \sum_{n=1}^N \exp \{ -t_n f_m(\mathbf{x}_n) \}$

linear combination of base classifiers:  $f_m(\mathbf{x}) = \frac{1}{2} \sum_{l=1}^m \alpha_l y_l(\mathbf{x})$

Q: How do they come up with this idea?

A: Motivated by finding a local optimal of  $E$

# Boosting

(c) Update the data weighting coefficients

$$w_n^{(m+1)} = w_n^{(m)} \exp \{ \alpha_m I(y_m(\mathbf{x}_n) \neq t_n) \} \quad (14.18)$$

the exponential error function:  $E = \sum_{n=1}^N \exp \{ -t_n f_m(\mathbf{x}_n) \}$   
 linear combination of base classifiers:  $f_m(\mathbf{x}) = \frac{1}{2} \sum_{l=1}^m \alpha_l y_l(\mathbf{x})$   
 $t_n \in \{-1, 1\} \quad y(\mathbf{x}) \in \{-1, 1\}$

Goal: minimize  $E$  with respect to both the weighting coefficients  $\alpha_l$  and the parameters of the base classifiers  $y_l(\mathbf{x})$

Suppose:  $y_1(\mathbf{x}), \dots, y_{m-1}(\mathbf{x}), \alpha_1, \dots, \alpha_{m-1}$  are fixed, minimize only with respect to  $y_m(\mathbf{x})$  and  $\alpha_m$

$$\begin{aligned} E &= \sum_{n=1}^N \exp \left\{ -t_n f_{m-1}(\mathbf{x}_n) - \frac{1}{2} t_n \alpha_m y_m(\mathbf{x}_n) \right\} \\ &= \sum_{n=1}^N w_n^{(m)} \exp \left\{ -\frac{1}{2} t_n \alpha_m y_m(\mathbf{x}_n) \right\} \end{aligned}$$

$$\begin{aligned} w_n^{(m+1)} &= \exp \{ -t_n f_m(\mathbf{x}_n) \} \\ &= \exp \left\{ -t_n \left( f_{m-1}(\mathbf{x}_n) + \frac{1}{2} \alpha_m y_m(\mathbf{x}_n) \right) \right\} \\ &= \exp \left\{ -t_n f_{m-1}(\mathbf{x}_n) - \frac{1}{2} t_n \alpha_m y_m(\mathbf{x}_n) \right\} \\ &= w_n^{(m)} \exp \left\{ -\frac{1}{2} t_n \alpha_m y_m(\mathbf{x}_n) \right\} \\ &\quad t_n y_m(\mathbf{x}_n) = 1 - 2\mathbb{I}(y_m(\mathbf{x}_n) \neq t_n) \\ \Rightarrow \\ w_n^{(m+1)} &= w_n^{(m)} \exp \left\{ -\frac{1}{2} \alpha_m (1 - 2\mathbb{I}(y_m(\mathbf{x}_n) \neq t_n)) \right\} \\ &= w_n^{(m)} \exp \left\{ -\frac{1}{2} \alpha_m \right\} \exp \left\{ \alpha_m \mathbb{I}(y_m(\mathbf{x}_n) \neq t_n) \right\} \\ &\quad \text{same for all } n \in [N] \end{aligned}$$

# Boosting

## AdaBoost

1. Initialize the data weighting coefficients  $\{w_n\}$  by setting  $w_n^{(1)} = 1/N$  for  $n = 1, \dots, N$ .
2. For  $m = 1, \dots, M$ :

- (a) Fit a classifier  $y_m(\mathbf{x})$  to the training data by minimizing the weighted error function

$$J_m = \sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n) \quad (14.15)$$

where  $I(y_m(\mathbf{x}_n) \neq t_n)$  is the indicator function and equals 1 when  $y_m(\mathbf{x}_n) \neq t_n$  and 0 otherwise.

- (b) Evaluate the quantities

$$\epsilon_m = \frac{\sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}} \quad (14.16)$$

and then use these to evaluate

$$\alpha_m = \ln \left\{ \frac{1 - \epsilon_m}{\epsilon_m} \right\}. \quad (14.17)$$

- (c) Update the data weighting coefficients

$$w_n^{(m+1)} = w_n^{(m)} \exp \{ \alpha_m I(y_m(\mathbf{x}_n) \neq t_n) \} \quad (14.18)$$

3. Make predictions using the final model, which is given by

$$Y_M(\mathbf{x}) = \text{sign} \left( \sum_{m=1}^M \alpha_m y_m(\mathbf{x}) \right). \quad (14.19)$$

$$t_n \in \{-1, 1\} \quad y(\mathbf{x}) \in \{-1, 1\}$$

the exponential error function:  $E = \sum_{n=1}^N \exp \{ -t_n f_m(\mathbf{x}_n) \}$

linear combination of base classifiers:  $f_m(\mathbf{x}) = \frac{1}{2} \sum_{l=1}^m \alpha_l y_l(\mathbf{x})$

Q: How do they come up with this idea?

A: Motivated by finding a local optimal of  $E$



# Boosting

(b) Evaluate the quantities

$$\epsilon_m = \frac{\sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}} \quad (14.16)$$

and then use these to evaluate

$$\alpha_m = \ln \left\{ \frac{1 - \epsilon_m}{\epsilon_m} \right\}. \quad (14.17)$$

the exponential error function:  $E = \sum_{n=1}^N \exp \{ -t_n f_m(\mathbf{x}_n) \}$

linear combination of base classifiers:  $f_m(\mathbf{x}) = \frac{1}{2} \sum_{l=1}^m \alpha_l y_l(\mathbf{x})$

$$t_n \in \{-1, 1\} \quad y(\mathbf{x}) \in \{-1, 1\}$$

Goal: minimize  $E$  with respect to both the weighting coefficients  $\alpha_l$  and the parameters of the base classifiers  $y_l(\mathbf{x})$

Suppose:  $y_1(\mathbf{x}), \dots, y_{m-1}(\mathbf{x}), \alpha_1, \dots, \alpha_{m-1}$  are fixed, minimize only with respect to  $y_m(\mathbf{x})$  and  $\alpha_m$

minimize  $E$  w.r.t  $\alpha_m$ :

$$\begin{aligned} E &= \sum_{n=1}^N \exp \left\{ -t_n f_{m-1}(\mathbf{x}_n) - \frac{1}{2} t_n \alpha_m y_m(\mathbf{x}_n) \right\} \\ &= \sum_{n=1}^N w_n^{(m)} \exp \left\{ -\frac{1}{2} t_n \alpha_m y_m(\mathbf{x}_n) \right\} \\ &= e^{-\alpha_m/2} \sum_{n \in \mathcal{T}_m} w_n^{(m)} + e^{\alpha_m/2} \sum_{n \in \mathcal{M}_m} w_n^{(m)} \\ &= (e^{\alpha_m/2} - e^{-\alpha_m/2}) \sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n) + e^{-\alpha_m/2} \sum_{n=1}^N w_n^{(m)} \end{aligned}$$

# Boosting

(b) Evaluate the quantities

$$\epsilon_m = \frac{\sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}} \quad (14.16)$$

and then use these to evaluate

$$\alpha_m = \ln \left\{ \frac{1 - \epsilon_m}{\epsilon_m} \right\}. \quad (14.17)$$

We want to minimize  $E$  w.r.t  $\alpha_m$ , where

$$E = (e^{\alpha_m/2} - e^{-\alpha_m/2}) \sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n) + e^{-\alpha_m/2} \sum_{n=1}^N w_n^{(m)}$$

$$\frac{\partial E}{\partial \alpha_m} = \frac{1}{2} \left( (e^{\alpha_m/2} + e^{-\alpha_m/2}) \sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n) - e^{-\alpha_m/2} \sum_{n=1}^N w_n^{(m)} \right)$$

set it to zero and rearrange:

$$\frac{\sum_n w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)}{\sum_n w_n^{(m)}} = \frac{e^{-\alpha_m/2}}{e^{\alpha_m/2} + e^{-\alpha_m/2}} = \frac{1}{e^{\alpha_m} + 1} = \epsilon_m$$

$$\Rightarrow e^{\alpha_m} = \frac{1 - \epsilon_m}{\epsilon_m}$$

$$\Rightarrow \alpha_m = \log \frac{1 - \epsilon_m}{\epsilon_m}$$

# Boosting

## AdaBoost

1. Initialize the data weighting coefficients  $\{w_n\}$  by setting  $w_n^{(1)} = 1/N$  for  $n = 1, \dots, N$ .
2. For  $m = 1, \dots, M$ :

- (a) Fit a classifier  $y_m(\mathbf{x})$  to the training data by minimizing the weighted error function

$$J_m = \sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n) \quad (14.15)$$

where  $I(y_m(\mathbf{x}_n) \neq t_n)$  is the indicator function and equals 1 when  $y_m(\mathbf{x}_n) \neq t_n$  and 0 otherwise.

- (b) Evaluate the quantities

$$\epsilon_m = \frac{\sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}} \quad (14.16)$$

and then use these to evaluate

$$\alpha_m = \ln \left\{ \frac{1 - \epsilon_m}{\epsilon_m} \right\}. \quad (14.17)$$

- (c) Update the data weighting coefficients

$$w_n^{(m+1)} = w_n^{(m)} \exp \{ \alpha_m I(y_m(\mathbf{x}_n) \neq t_n) \} \quad (14.18)$$

3. Make predictions using the final model, which is given by

$$Y_M(\mathbf{x}) = \text{sign} \left( \sum_{m=1}^M \alpha_m y_m(\mathbf{x}) \right). \quad (14.19)$$

$$t_n \in \{-1, 1\} \quad y(\mathbf{x}) \in \{-1, 1\}$$

This is trivial

the exponential error function:  $E = \sum_{n=1}^N \exp \{ -t_n f_m(\mathbf{x}_n) \}$   
linear combination of base classifiers:  $f_m(\mathbf{x}) = \frac{1}{2} \sum_{l=1}^m \alpha_l y_l(\mathbf{x})$

Q: How do they come up with this idea?

A: Motivated by finding a local optimal of  $E$



# Reference

- ▶ PRML chapter 14