

The Many Faces of Optimal Weak-to-Strong Learning

Mikael Møller Høgsgaard – Aarhus University

Kasper Green Larsen – Aarhus University

Markus Engelund Mathiasen – Aarhus University

NeurIPS 2024

Weak-to-Strong Learning

Binary classification

Weak Learner \mathcal{W} :

For any distribution \mathcal{D}

Samples: $S \sim \mathcal{D}^{m_0}$

Let $h_S \leftarrow \mathcal{W}(S)$

Then with prob. $1 - \delta_0$:

$$\text{er}_{\mathcal{D}}(h) \leq \frac{1}{2} - \gamma$$

Weak-to-Strong Learning

Binary classification

Weak Learner \mathcal{W} :

For any distribution \mathcal{D}

Samples: $S \sim \mathcal{D}^{m_0}$

Let $h_S \leftarrow \mathcal{W}(S)$

Then with prob. $1 - \delta_0$:

$$\text{er}_{\mathcal{D}}(h) \leq \frac{1}{2} - \gamma$$

Strong Learner \mathcal{A} :

For any distribution \mathcal{D} and parameters $0 < \varepsilon, \delta < 1$

Samples: $S \sim \mathcal{D}^{m(\varepsilon, \delta)}$

$m(\varepsilon, \delta)$ is sample-complexity of \mathcal{A} .

Let $h_S \leftarrow \mathcal{A}(S)$

Then with prob. $1 - \delta$:

$$\text{er}_{\mathcal{D}}(h) \leq \varepsilon$$

Weak-to-Strong Learning

Binary classification

Weak Learner \mathcal{W} :

For any distribution \mathcal{D}

Samples: $S \sim \mathcal{D}^{m_0}$

Let $h_S \leftarrow \mathcal{W}(S)$

Then with prob. $1 - \delta_0$:

$$\text{er}_{\mathcal{D}}(h) \leq \frac{1}{2} - \gamma$$



Strong Learner \mathcal{A} :

For any distribution \mathcal{D} and parameters $0 < \varepsilon, \delta < 1$

Samples: $S \sim \mathcal{D}^{m(\varepsilon, \delta)}$

$m(\varepsilon, \delta)$ is sample-complexity of \mathcal{A} .

Let $h_S \leftarrow \mathcal{A}(S)$

Then with prob. $1 - \delta$:

$$\text{er}_{\mathcal{D}}(h) \leq \varepsilon$$

Majority of 5

Algorithm 1: MAJORITY-OF-5(S, \mathcal{W})

Input: Training set $S = (x_1, y_1), \dots, (x_m, y_m)$. Weak learner \mathcal{W} .

Result: Hypothesis $g : \mathcal{X} \rightarrow \{-1, 1\}$.

1 Partition S into 5 disjoint pieces S_1, \dots, S_5 of size $m/5$.

2 **for** $t = 1, \dots, 5$ **do**

3 | Run AdaBoost on S_t with \mathcal{W} to obtain $f_t : \mathcal{X} \rightarrow \{-1, 1\}$.

4 $g \leftarrow \text{sign}(\sum_t f_t)$.

5 **return** g

S

Majority of 5

Algorithm 1: MAJORITY-OF-5(S, \mathcal{W})

Input: Training set $S = (x_1, y_1), \dots, (x_m, y_m)$. Weak learner \mathcal{W} .

Result: Hypothesis $g : \mathcal{X} \rightarrow \{-1, 1\}$.

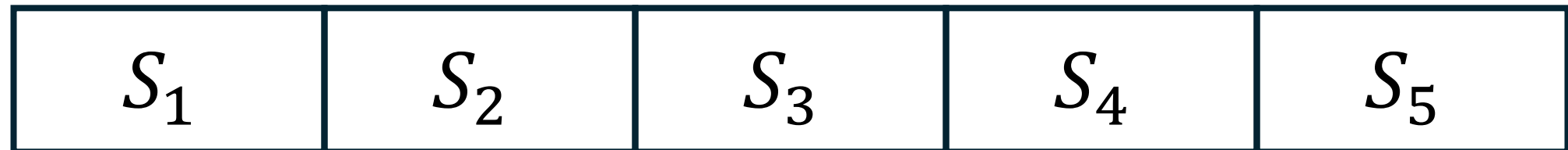
1 Partition S into 5 disjoint pieces S_1, \dots, S_5 of size $m/5$.

2 **for** $t = 1, \dots, 5$ **do**

3 | Run AdaBoost on S_t with \mathcal{W} to obtain $f_t : \mathcal{X} \rightarrow \{-1, 1\}$.

4 $g \leftarrow \text{sign}(\sum_t f_t)$.

5 **return** g



Majority of 5

Algorithm 1: MAJORITY-OF-5(S, \mathcal{W})

Input: Training set $S = (x_1, y_1), \dots, (x_m, y_m)$. Weak learner \mathcal{W} .

Result: Hypothesis $g : \mathcal{X} \rightarrow \{-1, 1\}$.

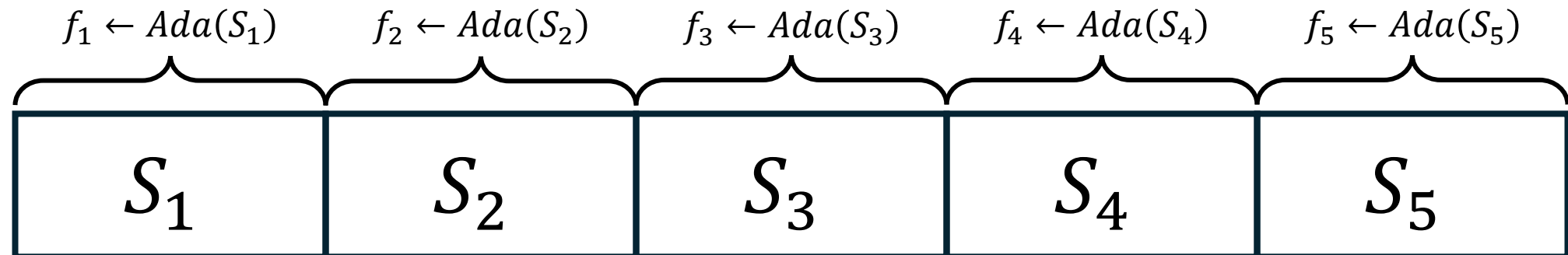
1 Partition S into 5 disjoint pieces S_1, \dots, S_5 of size $m/5$.

2 **for** $t = 1, \dots, 5$ **do**

3 | Run AdaBoost on S_t with \mathcal{W} to obtain $f_t : \mathcal{X} \rightarrow \{-1, 1\}$.

4 $g \leftarrow \text{sign}(\sum_t f_t)$.

5 **return** g



Majority of 5

Algorithm 1: MAJORITY-OF-5(S, \mathcal{W})

Input: Training set $S = (x_1, y_1), \dots, (x_m, y_m)$. Weak learner \mathcal{W} .

Result: Hypothesis $g : \mathcal{X} \rightarrow \{-1, 1\}$.

1 Partition S into 5 disjoint pieces S_1, \dots, S_5 of size $m/5$.

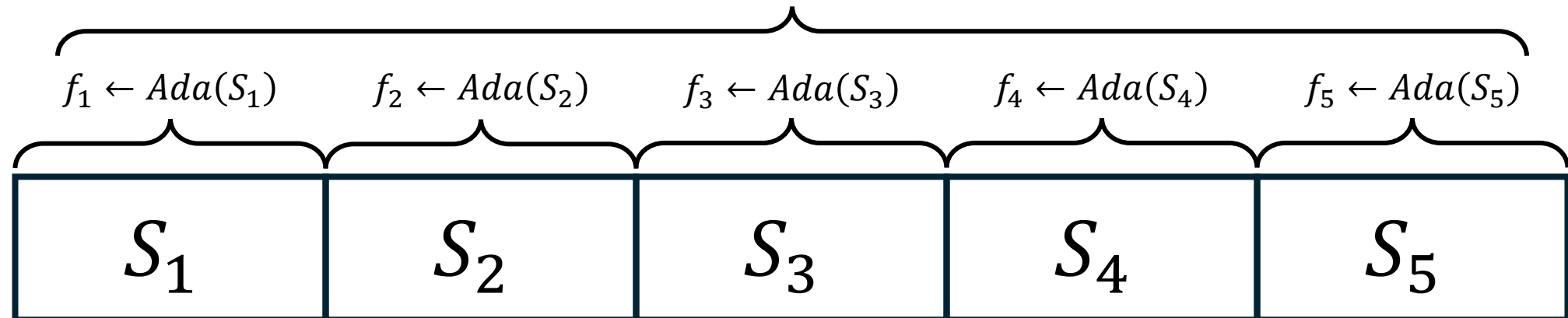
2 **for** $t = 1, \dots, 5$ **do**

3 | Run AdaBoost on S_t with \mathcal{W} to obtain $f_t : \mathcal{X} \rightarrow \{-1, 1\}$.

4 $g \leftarrow \text{sign}(\sum_t f_t)$.

5 **return** g

$$g \leftarrow \text{Maj}(f_1, \dots, f_5)$$



Majority of 5 - Guarantee

Theorem 1

$$\mathbb{E}_{S \sim \mathcal{D}^m} [\text{er}_{\mathcal{D}}(g)] = O\left(\frac{d}{\gamma^2 m}\right)$$

Algorithm 1: MAJORITY-OF-5(S, \mathcal{W})

Input: Training set $S = (x_1, y_1), \dots, (x_m, y_m)$. Weak learner \mathcal{W} .

Result: Hypothesis $g : \mathcal{X} \rightarrow \{-1, 1\}$.

- 1 Partition S into 5 disjoint pieces S_1, \dots, S_5 of size $m/5$.
 - 2 **for** $t = 1, \dots, 5$ **do**
 - 3 | Run AdaBoost on S_t with \mathcal{W} to obtain $f_t : \mathcal{X} \rightarrow \{-1, 1\}$.
 - 4 $g \leftarrow \text{sign}(\sum_t f_t)$.
 - 5 **return** g
-

Sample-Optimal Weak-to-Strong Learners

Algorithm	Error	Invocations of Weak Learner
AdaBoost	$O\left(\frac{d \ln^2(m)}{\gamma^2 m}\right)$	$O\left(\frac{\ln(m)}{\gamma^2}\right)$
LarsenRitzert ¹	$O\left(\frac{d}{\gamma^2 m}\right)$	$O\left(\frac{m^{0.8}}{\gamma^2}\right)$
Bagged AdaBoost ²	$O\left(\frac{d}{\gamma^2 m}\right)$	$O\left(\frac{\ln^2(m)}{\gamma^2}\right)$
Majority of 5	$O\left(\frac{d}{\gamma^2 m}\right)$	$O\left(\frac{\ln(m)}{\gamma^2}\right)$

1) Green Larsen, K., & Ritzert, M. (2022). Optimal weak to strong learning. *Advances in Neural Information Processing Systems*, 35, 32830-32841.

2) Larsen, K. G. (2023, July). Bagging is an optimal PAC learner. In *The Thirty Sixth Annual Conference on Learning Theory* (pp. 450-468). PMLR.

Open Questions

- Is Majority of 5 optimal in high probability setting?
- Is Majority of 3 sufficient?