

Optimal Parallelization of Boosting



Arthur C.W. da Cunha*



Mikael Møller Høgsgaard*



Kasper Green Larsen

Aarhus University - Denmark

NeurIPS
2024

*Looking for jobs

- Domain: \mathcal{X} (and a distribution \mathcal{D} over it)

- Domain: \mathcal{X} (and a distribution \mathcal{D} over it)

\mathcal{X}



- Domain: \mathcal{X} (and a distribution \mathcal{D} over it)
- Label space: \mathcal{Y} , with $|\mathcal{Y}| = 2$

\mathcal{X}



- Domain: \mathcal{X} (and a distribution \mathcal{D} over it)
- Label space: \mathcal{Y} , with $|\mathcal{Y}| = 2$

\mathcal{X}



\mathcal{Y}

“CAT”

“DOG”

- Domain: \mathcal{X} (and a distribution \mathcal{D} over it)
- Label space: \mathcal{Y} , with $|\mathcal{Y}| = 2$
- Target concept: $c: \mathcal{X} \rightarrow \mathcal{Y}$

\mathcal{X}

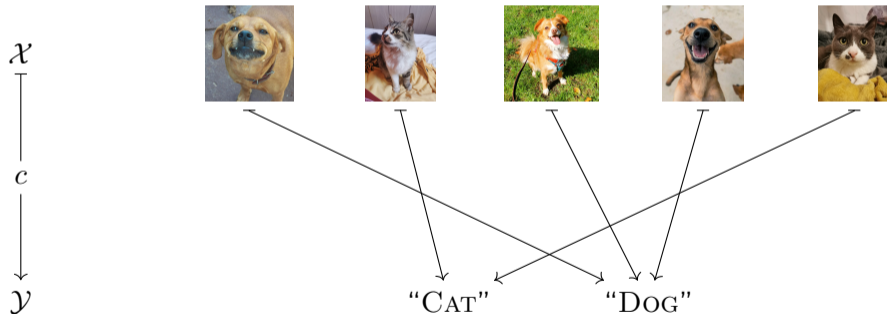


\mathcal{Y}

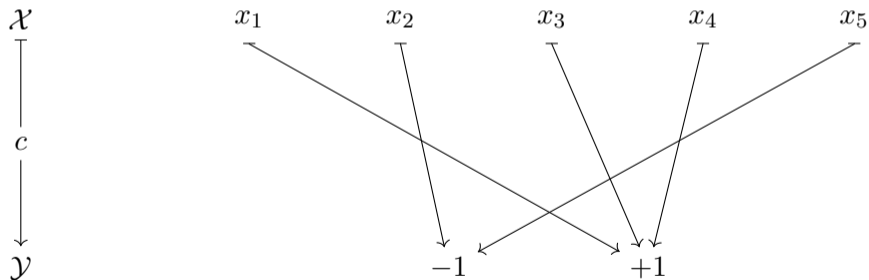
“CAT”

“DOG”

- Domain: \mathcal{X} (and a distribution \mathcal{D} over it)
- Label space: \mathcal{Y} , with $|\mathcal{Y}| = 2$
- Target concept: $c: \mathcal{X} \rightarrow \mathcal{Y}$



- Domain: \mathcal{X} (and a distribution \mathcal{D} over it)
- Label space: \mathcal{Y} , with $|\mathcal{Y}| = 2$
- Target concept: $c: \mathcal{X} \rightarrow \mathcal{Y}$

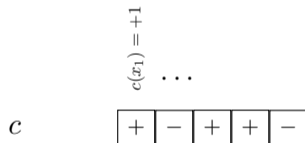


- Domain: \mathcal{X} (and a distribution \mathcal{D} over it)
- Label space: \mathcal{Y} , with $|\mathcal{Y}| = 2$
- Target concept: $c: \mathcal{X} \rightarrow \mathcal{Y}$

c

+	-	+	+	-
---	---	---	---	---

- Domain: \mathcal{X} (and a distribution \mathcal{D} over it)
- Label space: \mathcal{Y} , with $|\mathcal{Y}| = 2$
- Target concept: $c: \mathcal{X} \rightarrow \mathcal{Y}$



c

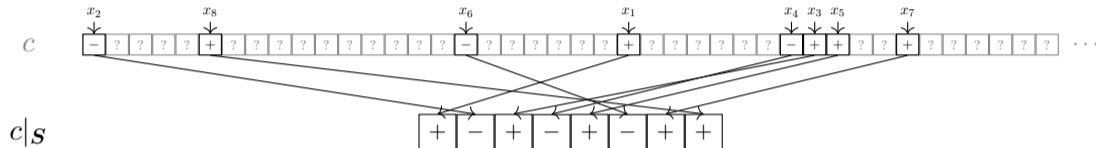
-	-	+	-	+	+	-	-	+	+	+	+	-	+	-	+	-	+	-	-	-	+	-	+	-	+	-	+	-	+	+	+	-	+	+	+	+	-	-	-	-
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

 ...

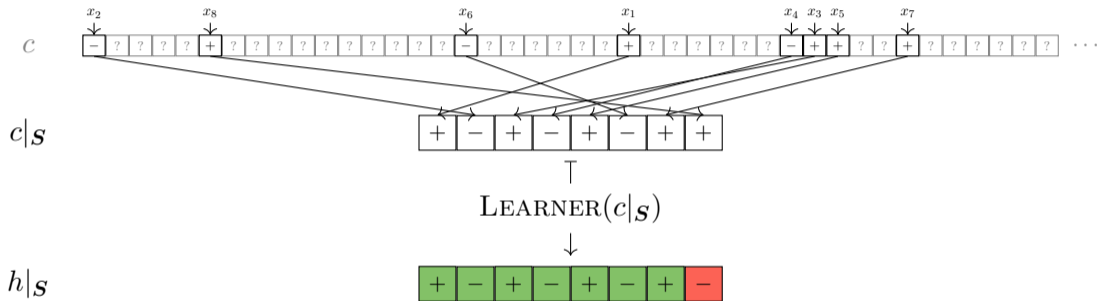
- We only see a sample $S \sim \mathcal{D}^m$ ($m = 8$)



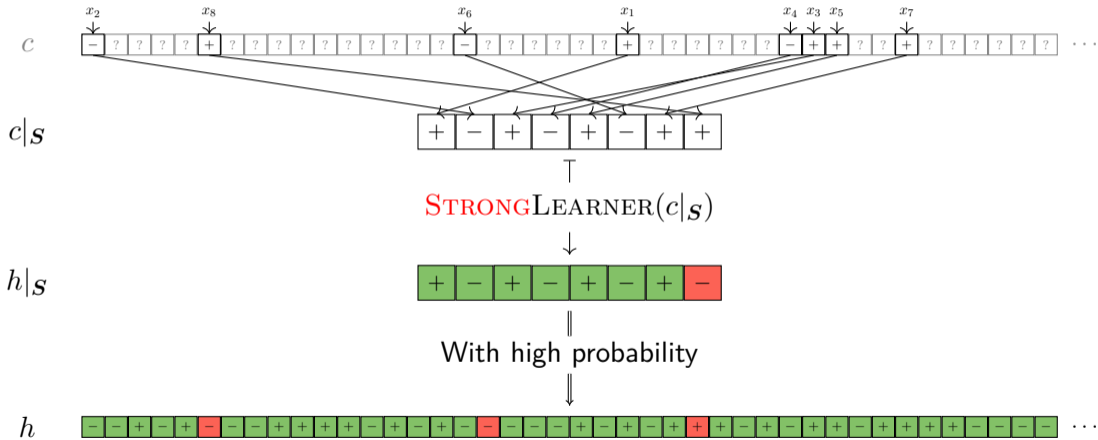
- We only see a sample $S \sim \mathcal{D}^m$ ($m = 8$)



- We only see a sample $S \sim \mathcal{D}^m$ ($m = 8$)



- We only see a sample $S \sim \mathcal{D}^m$ ($m = 8$)



Algorithm STRONGLEARNER such that

Algorithm STRONGLEARNER such that

For all Precision $\varepsilon \in (0, 1)$.

Algorithm STRONGLEARNER such that

For all Precision $\varepsilon \in (0, 1)$.

Given Sufficiently large $m = m(\varepsilon)$ sample $\mathcal{S} \sim \mathcal{D}^m$.

Algorithm STRONGLEARNER such that

For all Precision $\varepsilon \in (0, 1)$.

Given Sufficiently large $m = m(\varepsilon)$ sample $\mathcal{S} \sim \mathcal{D}^m$.

Given Examples of a target concept on sample $c|_{\mathcal{S}}$.

Algorithm STRONGLEARNER such that

For all Precision $\varepsilon \in (0, 1)$.

Given Sufficiently large $m = m(\varepsilon)$ sample $\mathcal{S} \sim \mathcal{D}^m$.

Given Examples of a target concept on sample $c|_{\mathcal{S}}$.

Satisfies With high probability returns a classifier $h = \text{STRONGLEARNER}(c|_{\mathcal{S}})$ that approximates c **well on the entire domain**.

Algorithm STRONGLEARNER such that

For all Precision $\varepsilon \in (0, 1)$.

Given Sufficiently large $m = m(\varepsilon)$ sample $\mathcal{S} \sim \mathcal{D}^m$.

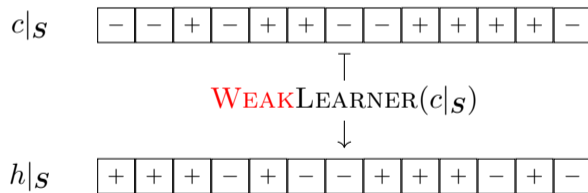
Given Examples of a target concept on sample $c|_{\mathcal{S}}$.

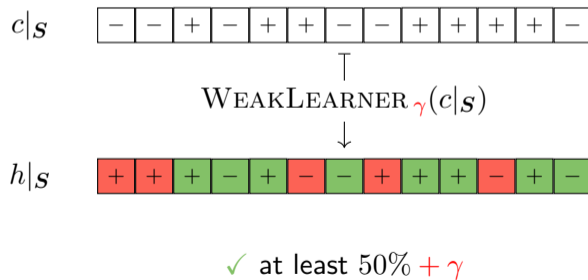
Satisfies With high probability returns a classifier $h = \text{STRONGLEARNER}(c|_{\mathcal{S}})$ that approximates c well on the entire domain. More precisely,

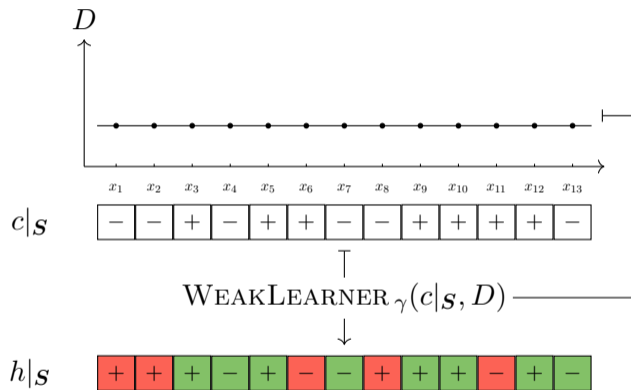
$$\text{err}_{\mathcal{D}}(h) := \Pr_{\mathbf{x} \sim \mathcal{D}} [h(\mathbf{x}) \neq c(\mathbf{x})] < \varepsilon.$$

$c|s$

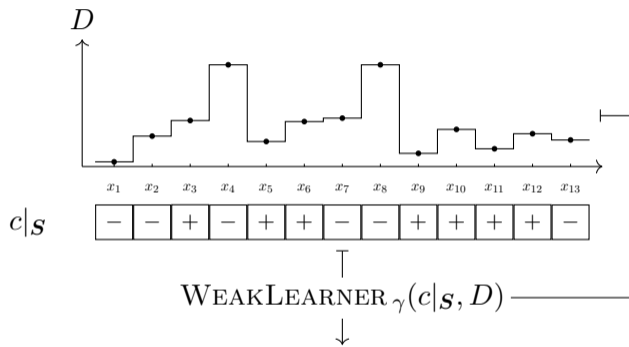
-	-	+	-	+	+	-	-	+	+	+	+	-
---	---	---	---	---	---	---	---	---	---	---	---	---

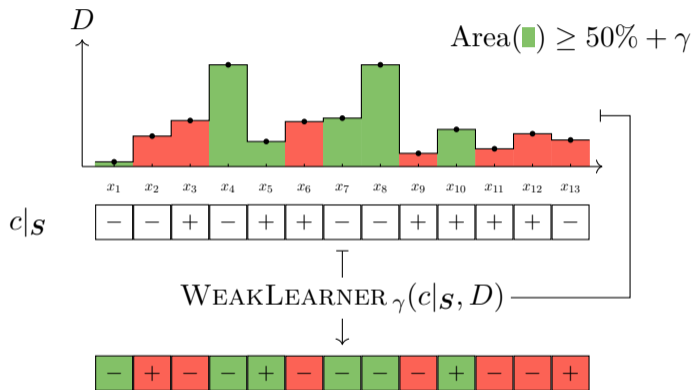






✓ at least $50\% + \gamma$





Algorithm $\text{WEAKLEARNER}_\gamma$ such that

Algorithm \mathcal{W}_γ such that

Algorithm \mathcal{W}_γ such that

Given Set $S \subseteq \mathcal{X}$.

Algorithm \mathcal{W}_γ such that

Given Set $S \subseteq \mathcal{X}$.

Given Examples of the target concept $c|_S$.

Algorithm \mathcal{W}_γ such that

Given Set $S \subseteq \mathcal{X}$.

Given Examples of the target concept $c|_S$.

Given Any distribution (weighing) D over S .

Algorithm \mathcal{W}_γ such that

Given Set $S \subseteq \mathcal{X}$.

Given Examples of the target concept $c|_S$.

Given Any distribution (weighing) D over S .

Satisfies Returns a classifier $h = \mathcal{W}_\gamma(c|_S, D)$ that approximates c a bit better than chance on the training data.

Algorithm \mathcal{W}_γ such that

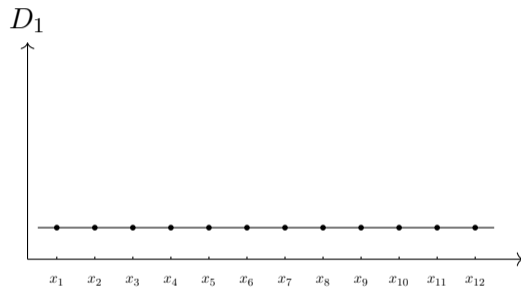
Given Set $S \subseteq \mathcal{X}$.

Given Examples of the target concept $c|_S$.

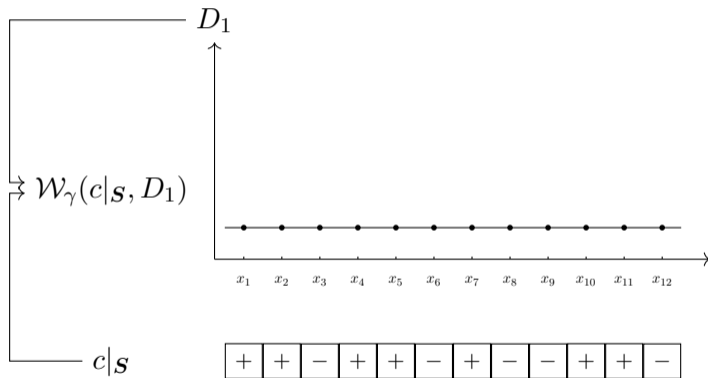
Given Any distribution (weighing) D over S .

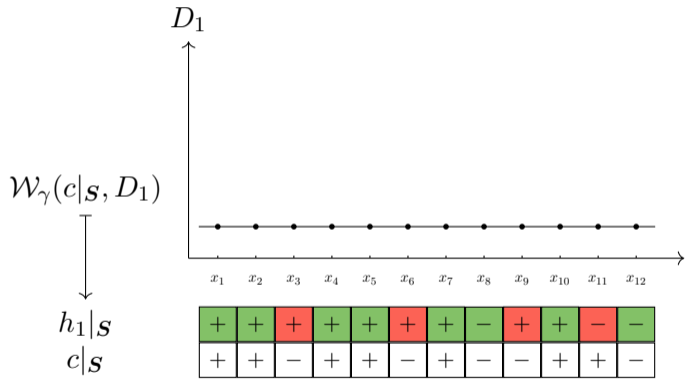
Satisfies Returns a classifier $h = \mathcal{W}_\gamma(c|_S, D)$ that approximates c a bit better than chance on the training data. More precisely,

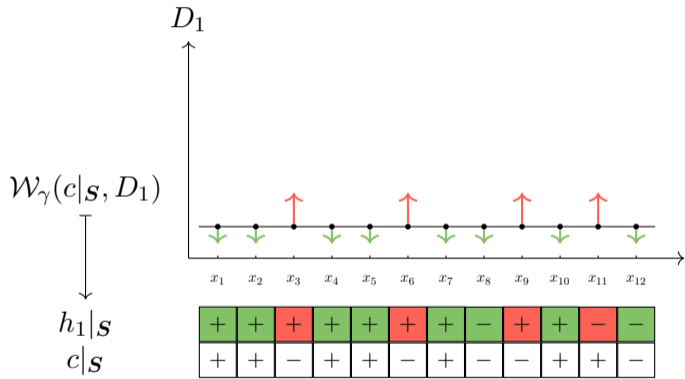
$$\text{err}_D(h) := \Pr_{\mathbf{x} \sim D} [h(\mathbf{x}) \neq c(\mathbf{x})] < \frac{1}{2} - \gamma.$$

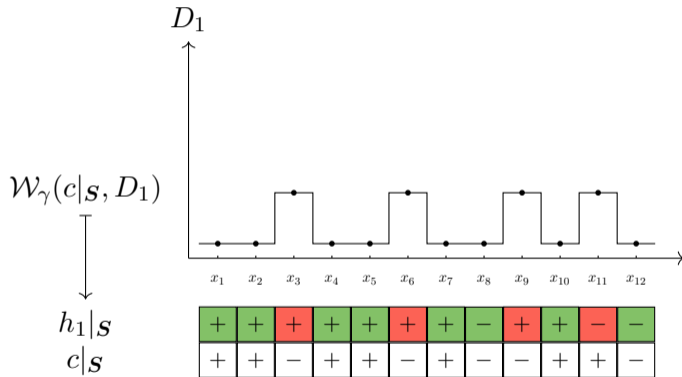
 $c|_S$

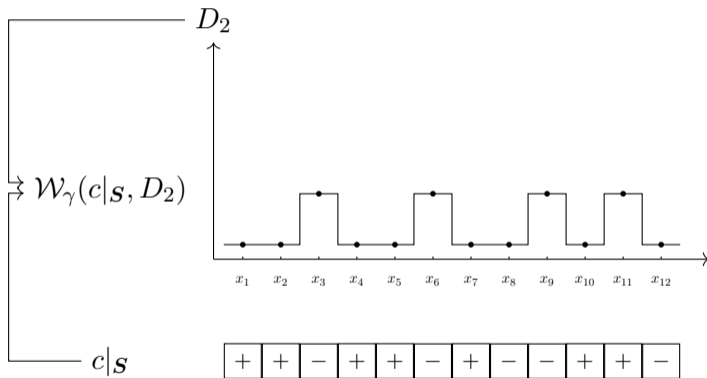
+	+	-	+	+	-	+	-	-	+	+	-
---	---	---	---	---	---	---	---	---	---	---	---

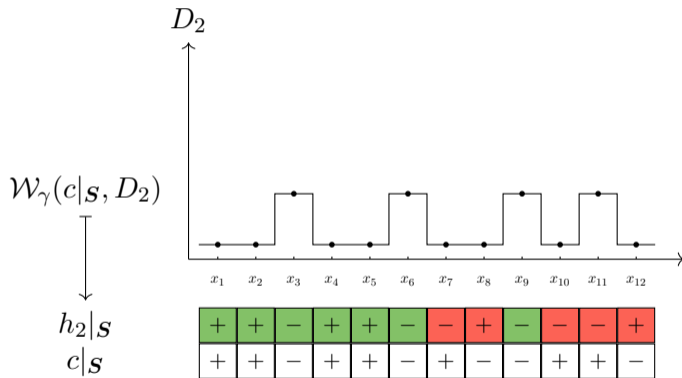


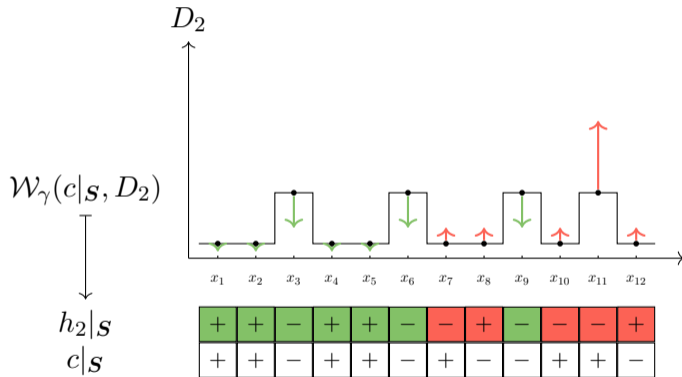


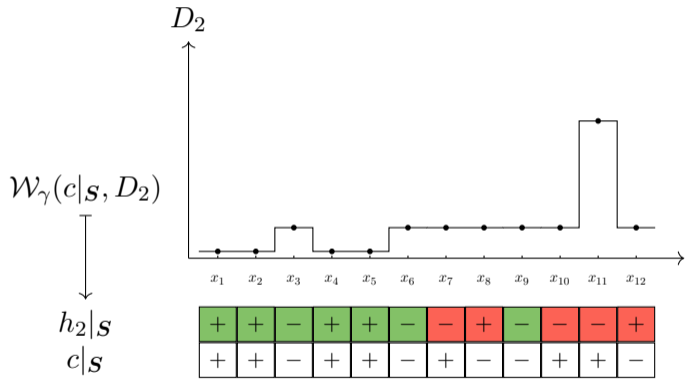


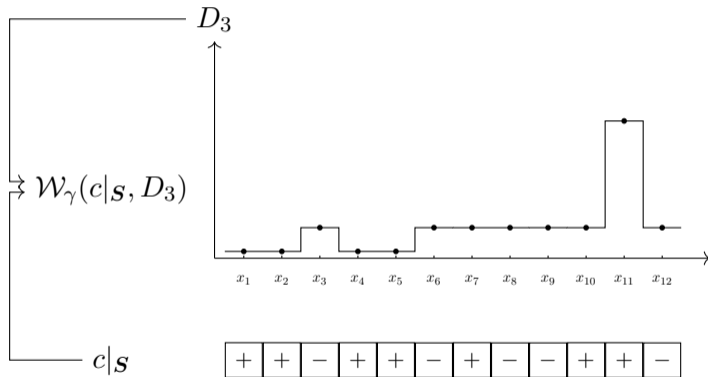


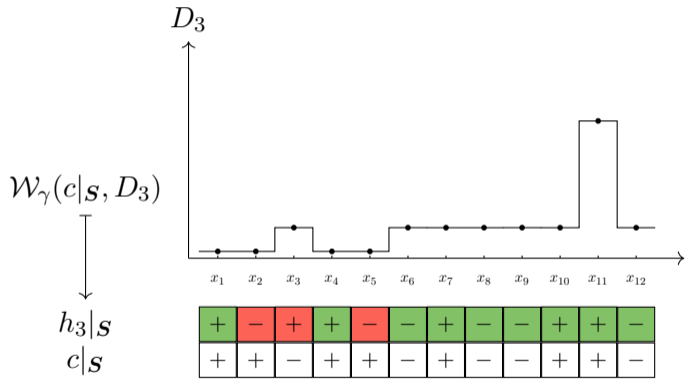


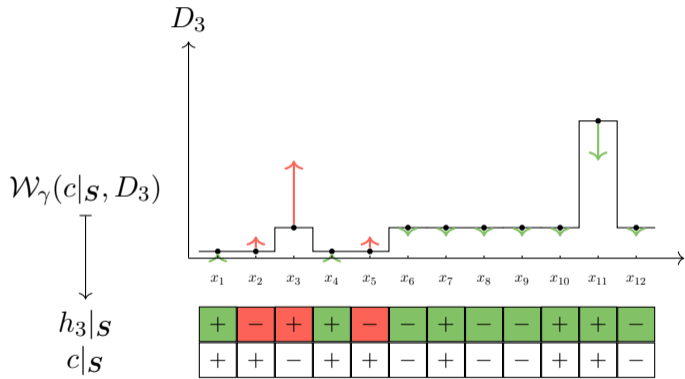


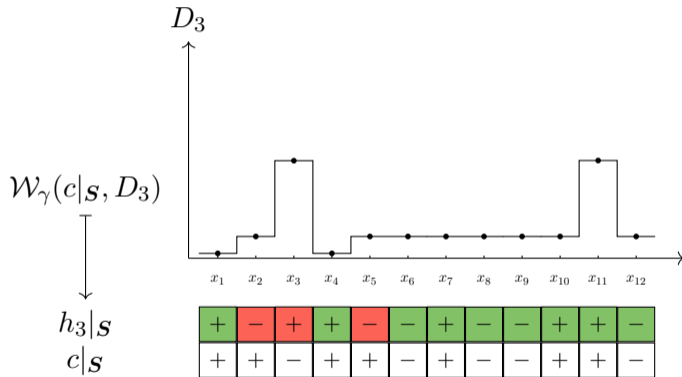


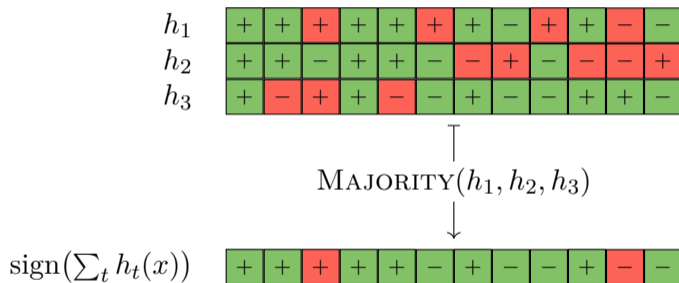












- State-of-the-art in practice:

- State-of-the-art in practice:
 - Gradient boosters: e.g., XGBOOST (Chen and Guestrin, 2016) and LIGHTGBM (Ke et al., 2017);

- State-of-the-art in practice:
 - Gradient boosters: e.g., XGBOOST (Chen and Guestrin, 2016) and LIGHTGBM (Ke et al., 2017);
 - Base learners: low to medium depth decision trees;

- State-of-the-art in practice:
 - Gradient boosters: e.g., XGBOOST (Chen and Guestrin, 2016) and LIGHTGBM (Ke et al., 2017);
 - Base learners: low to medium depth decision trees;
 - Often win kaggle™ competitions with small datasets and/or tabular data.

- State-of-the-art in practice:
 - Gradient boosters: e.g., XGBOOST (Chen and Guestrin, 2016) and LIGHTGBM (Ke et al., 2017);
 - Base learners: low to medium depth decision trees;
 - Often win kaggle™ competitions with small datasets and/or tabular data.
- Drawbacks:

- State-of-the-art in practice:
 - Gradient boosters: e.g., XGBOOST (Chen and Guestrin, 2016) and LIGHTGBM (Ke et al., 2017);
 - Base learners: low to medium depth decision trees;
 - Often win kaggle™ competitions with small datasets and/or tabular data.
- Drawbacks:
 - Achieving the best performance often takes 1000s of iterations.

- State-of-the-art in practice:
 - Gradient boosters: e.g., XGBOOST (Chen and Guestrin, 2016) and LIGHTGBM (Ke et al., 2017);
 - Base learners: low to medium depth decision trees;
 - Often win kaggle™ competitions with small datasets and/or tabular data.
- Drawbacks:
 - Achieving the best performance often takes 1000s of iterations.
 - Sequential nature: even with many computers available, it's not obvious how to speed it up.

- State-of-the-art in practice:
 - Gradient boosters: e.g., XGBOOST (Chen and Guestrin, 2016) and LIGHTGBM (Ke et al., 2017);
 - Base learners: low to medium depth decision trees;
 - Often win kaggle™ competitions with small datasets and/or tabular data.
- Drawbacks:
 - Achieving the best performance often takes 1000s of iterations.
 - Sequential nature: even with many computers available, it's not obvious how to speed it up.
 - Infeasible for large datasets or “expensive” base learners.

Class of **parallel** Boosting algorithms considered

1
2
3
4
5

Class of **parallel** Boosting algorithms considered

Input : training examples $c|_S$, γ -weak learner \mathcal{W}_γ

1
2
3
4
5

Class of **parallel** Boosting algorithms considered

Input : training examples $c|_S$, γ -weak learner \mathcal{W}_γ

1 **for** $p \leftarrow 1$ **to** P **do**

2 |

3 |

4 |

5 |

Class of **parallel** Boosting algorithms considered

Input : training examples $c|_S$, γ -weak learner \mathcal{W}_γ

1 **for** $p \leftarrow 1$ **to** P **do**

2 | **parallel for** $t \leftarrow 1$ **to** T **do**

3 | | $h_{p,t} \leftarrow$ Query \mathcal{W}_γ with some distribution $D_{p,t}$

4 |

5

Class of **parallel** Boosting algorithms considered

Input : training examples $c|_S$, γ -weak learner \mathcal{W}_γ

```
1 for  $p \leftarrow 1$  to  $P$  do
2   parallel for  $t \leftarrow 1$  to  $T$  do
3      $h_{p,t} \leftarrow$  Query  $\mathcal{W}_\gamma$  with some distribution  $D_{p,t}$ 
4     ... // Do something with hypotheses found so far
5
```

Class of **parallel** Boosting algorithms considered

Input : training examples $c|_S$, γ -weak learner \mathcal{W}_γ

1 **for** $p \leftarrow 1$ **to** P **do**

2 | **parallel for** $t \leftarrow 1$ **to** T **do**

3 | | $h_{p,t} \leftarrow$ Query \mathcal{W}_γ with some distribution $D_{p,t}$

4 | | ... // Do something with hypotheses found so far

5 **return** Classifier H with generalization error **not much worse than ADABOOST's** ($\tilde{O}\left(\frac{d}{\gamma^2 \cdot |S|}\right)$), with $d = \text{VC}(\text{"base classifiers"})$

Class of **parallel** Boosting algorithms considered

Input : training examples $c|_S$, γ -weak learner \mathcal{W}_γ

1 **for** $p \leftarrow 1$ **to** P **do**

2 | **parallel for** $t \leftarrow 1$ **to** T **do**

3 | | $h_{p,t} \leftarrow$ Query \mathcal{W}_γ with some distribution $D_{p,t}$

4 | | ... // Do something with hypotheses found so far

5 **return** Classifier H with generalization error **not much worse than ADABOOST's** ($\tilde{O}\left(\frac{d}{\gamma^2 \cdot |S|}\right)$), with $d = \text{VC}(\text{"base classifiers"})$

- E.g., ADABOOST: $P = \Theta\left(\frac{\ln|S|}{\gamma^2}\right)$ and $T = 1$ (no parallelism).

Class of **parallel** Boosting algorithms considered

Input : training examples $c|_S$, γ -weak learner \mathcal{W}_γ

1 **for** $p \leftarrow 1$ **to** P **do**

2 | **parallel for** $t \leftarrow 1$ **to** T **do**

3 | | $h_{p,t} \leftarrow$ Query \mathcal{W}_γ with some distribution $D_{p,t}$

4 | | ... // Do something with hypotheses found so far

5 **return** Classifier H with generalization error **not much worse than ADABOOST's** ($\tilde{O}\left(\frac{d}{\gamma^2 \cdot |S|}\right)$), with $d = \text{VC}(\text{"base classifiers"})$

- E.g., ADABOOST: $P = \Theta\left(\frac{\ln|S|}{\gamma^2}\right)$ and $T = 1$ (no parallelism).
- Karbasi and Larsen (2024):

Class of **parallel** Boosting algorithms considered

Input : training examples $c|_S$, γ -weak learner \mathcal{W}_γ

1 **for** $p \leftarrow 1$ **to** P **do**

2 | **parallel for** $t \leftarrow 1$ **to** T **do**

3 | | $h_{p,t} \leftarrow$ Query \mathcal{W}_γ with some distribution $D_{p,t}$

4 | | ... // Do something with hypotheses found so far

5 **return** Classifier H with generalization error **not much worse than ADABOOST's** ($\tilde{O}\left(\frac{d}{\gamma^2 \cdot |S|}\right)$), with $d = \text{VC}(\text{"base classifiers"})$

- E.g., ADABOOST: $P = \Theta\left(\frac{\ln|S|}{\gamma^2}\right)$ and $T = 1$ (no parallelism).
- Karbasi and Larsen (2024):
 - Boosting algorithm with $P = 1$ and $T = \exp\left(O\left(\frac{d \ln m}{\gamma^2}\right)\right)$.

Class of **parallel** Boosting algorithms considered

Input : training examples $c|_S$, γ -weak learner \mathcal{W}_γ

1 **for** $p \leftarrow 1$ **to** P **do**

2 | **parallel for** $t \leftarrow 1$ **to** T **do**

3 | | $h_{p,t} \leftarrow$ Query \mathcal{W}_γ with some distribution $D_{p,t}$

4 | | ... // Do something with hypotheses found so far

5 **return** Classifier H with generalization error **not much worse than ADABOOST's** $(\tilde{O}(\frac{d}{\gamma^2 \cdot |S|}))$, with $d = \text{VC}(\text{"base classifiers"})$

- E.g., ADABOOST: $P = \Theta(\frac{\ln|S|}{\gamma^2})$ and $T = 1$ (no parallelism).
- Karbasi and Larsen (2024):
 - Boosting algorithm with $P = 1$ and $T = \exp(O(\frac{d \ln m}{\gamma^2}))$.
 - Lower bound: any **significant parallelization of Boosting requires exponential total work.**

Class of **parallel** Boosting algorithms considered

Input : training examples $c|_S$, γ -weak learner \mathcal{W}_γ

1 **for** $p \leftarrow 1$ **to** P **do**

2 | **parallel for** $t \leftarrow 1$ **to** T **do**

3 | | $h_{p,t} \leftarrow$ Query \mathcal{W}_γ with some distribution $D_{p,t}$

4 | | ... // Do something with hypotheses found so far

5 **return** Classifier H with generalization error **not much worse than ADABOOST's** $(\tilde{O}\left(\frac{d}{\gamma^2 \cdot |S|}\right))$, with $d = \text{VC}(\text{"base classifiers"})$

- E.g., ADABOOST: $P = \Theta\left(\frac{\ln|S|}{\gamma^2}\right)$ and $T = 1$ (no parallelism).
- Karbasi and Larsen (2024):
 - Boosting algorithm with $P = 1$ and $T = \exp\left(O\left(\frac{d \ln m}{\gamma^2}\right)\right)$.
 - Lower bound: any **significant parallelization of Boosting requires exponential total work.**
- Lyu, Wu and Yang (2024):

Class of **parallel** Boosting algorithms considered

Input : training examples $c|_S$, γ -weak learner \mathcal{W}_γ

1 **for** $p \leftarrow 1$ **to** P **do**

2 | **parallel for** $t \leftarrow 1$ **to** T **do**

3 | | $h_{p,t} \leftarrow$ Query \mathcal{W}_γ with some distribution $D_{p,t}$

4 | | ... // Do something with hypotheses found so far

5 **return** Classifier H with generalization error **not much worse than ADABOOST's** ($\tilde{O}\left(\frac{d}{\gamma^2 \cdot |S|}\right)$), with $d = \text{VC}(\text{"base classifiers"})$)

- E.g., ADABOOST: $P = \Theta\left(\frac{\ln|S|}{\gamma^2}\right)$ and $T = 1$ (no parallelism).
- Karbasi and Larsen (2024):
 - Boosting algorithm with $P = 1$ and $T = \exp\left(O\left(\frac{d \ln m}{\gamma^2}\right)\right)$.
 - Lower bound: any **significant parallelization of Boosting requires exponential total work.**
- Lyu, Wu and Yang (2024):
 - $P = O\left(\frac{\ln m}{\gamma^2 R}\right)$ and $T = \exp\left(O(dR^2)\right) \cdot \ln \frac{1}{\gamma}$.

Class of **parallel** Boosting algorithms considered

Input : training examples $c|_S$, γ -weak learner \mathcal{W}_γ

1 **for** $p \leftarrow 1$ **to** P **do**

2 | **parallel for** $t \leftarrow 1$ **to** T **do**

3 | | $h_{p,t} \leftarrow$ Query \mathcal{W}_γ with some distribution $D_{p,t}$

4 | | ... // Do something with hypotheses found so far

5 **return** Classifier H with generalization error **not much worse than ADABOOST's** ($\tilde{O}\left(\frac{d}{\gamma^2 \cdot |S|}\right)$), with $d = \text{VC}(\text{"base classifiers"})$

- E.g., ADABOOST: $P = \Theta\left(\frac{\ln|S|}{\gamma^2}\right)$ and $T = 1$ (no parallelism).
- Karbasi and Larsen (2024):
 - Boosting algorithm with $P = 1$ and $T = \exp\left(O\left(\frac{d \ln m}{\gamma^2}\right)\right)$.
 - Lower bound: any **significant parallelization of Boosting requires exponential total work**.
- Lyu, Wu and Yang (2024):
 - $P = O\left(\frac{\ln m}{\gamma^2 R}\right)$ and $T = \exp\left(O(dR^2)\right) \cdot \ln \frac{1}{\gamma}$.
 - Improved lower bound.

Class of **parallel** Boosting algorithms considered

Input : training examples $c|_S$, γ -weak learner \mathcal{W}_γ

1 **for** $p \leftarrow 1$ **to** P **do**

2 | **parallel for** $t \leftarrow 1$ **to** T **do**

3 | | $h_{p,t} \leftarrow$ Query \mathcal{W}_γ with some distribution $D_{p,t}$

4 | | ... // Do something with hypotheses found so far

5 **return** Classifier H with generalization error **not much worse than ADABOOST's** ($\tilde{O}\left(\frac{d}{\gamma^2 \cdot |S|}\right)$), with $d = \text{VC}(\text{"base classifiers"})$

- E.g., ADABOOST: $P = \Theta\left(\frac{\ln|S|}{\gamma^2}\right)$ and $T = 1$ (no parallelism).
- Karbasi and Larsen (2024):
 - Boosting algorithm with $P = 1$ and $T = \exp\left(O\left(\frac{d \ln m}{\gamma^2}\right)\right)$.
 - Lower bound: any **significant parallelization of Boosting requires exponential total work.**
- Lyu, Wu and Yang (2024):
 - $P = O\left(\frac{\ln m}{\gamma^2 R}\right)$ and $T = \exp\left(O(dR^2)\right) \cdot \ln \frac{1}{\gamma}$.
 - Improved lower bound.
- **A gap remains**

Class of **parallel** Boosting algorithms considered

Input : training examples $c|_S$, γ -weak learner \mathcal{W}_γ

1 **for** $p \leftarrow 1$ **to** P **do**

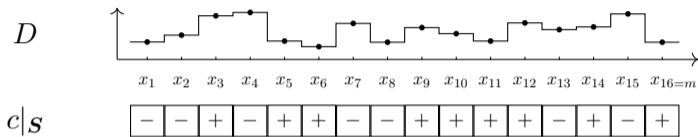
2 | **parallel for** $t \leftarrow 1$ **to** T **do**

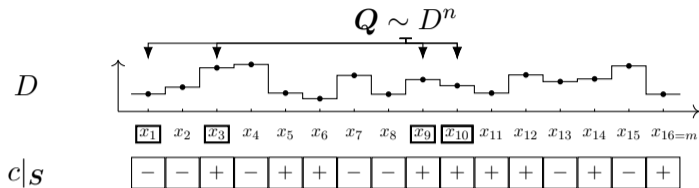
3 | | $h_{p,t} \leftarrow$ Query \mathcal{W}_γ with some distribution $D_{p,t}$

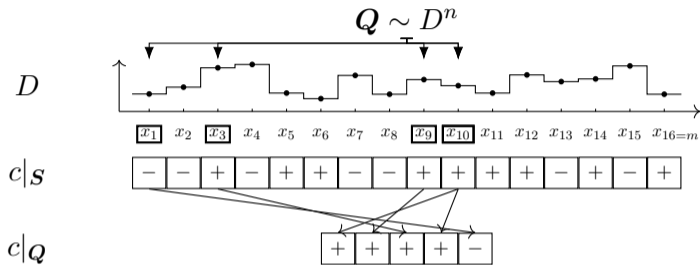
4 | | ... // Do something with hypotheses found so far

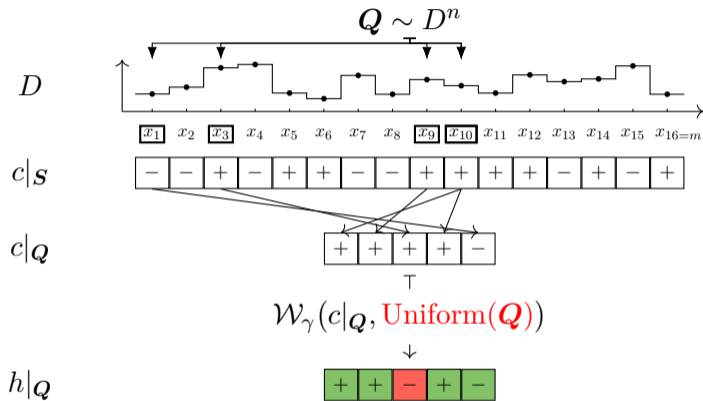
5 **return** Classifier H with generalization error **not much worse than ADABOOST's** ($\tilde{O}\left(\frac{d}{\gamma^2 \cdot |S|}\right)$), with $d = \text{VC}(\text{"base classifiers"})$

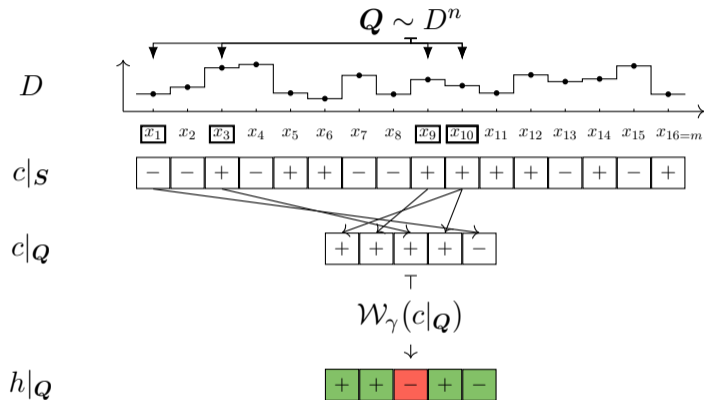
- E.g., ADABOOST: $P = \Theta\left(\frac{\ln|S|}{\gamma^2}\right)$ and $T = 1$ (no parallelism).
- Karbasi and Larsen (2024):
 - Boosting algorithm with $P = 1$ and $T = \exp\left(O\left(\frac{d \ln m}{\gamma^2}\right)\right)$.
 - Lower bound: any **significant parallelization of Boosting requires exponential total work.**
- Lyu, Wu and Yang (2024):
 - $P = O\left(\frac{\ln m}{\gamma^2 R}\right)$ and $T = \exp\left(O(dR^2)\right) \cdot \ln \frac{1}{\gamma}$.
 - Improved lower bound.
- A gap remains: **this work closes it.**

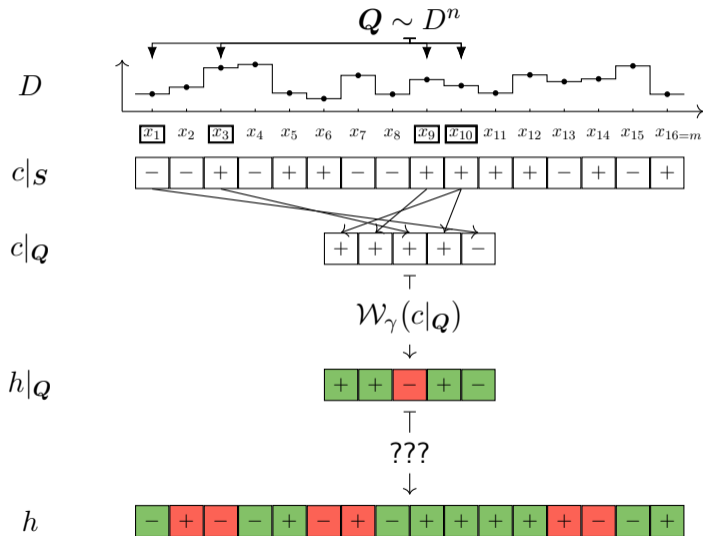


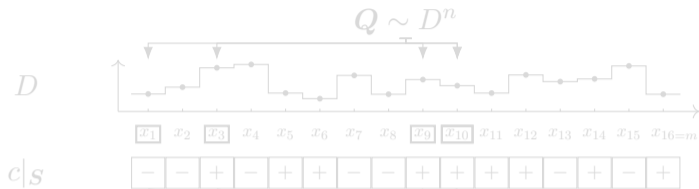




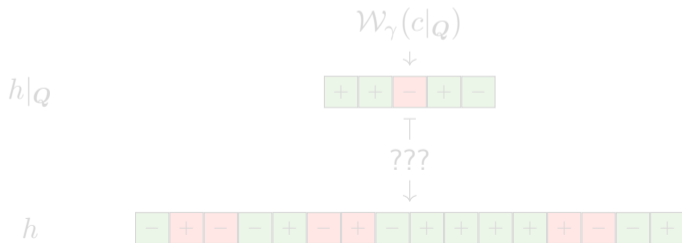


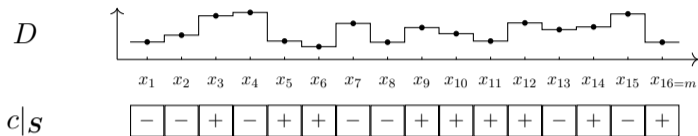


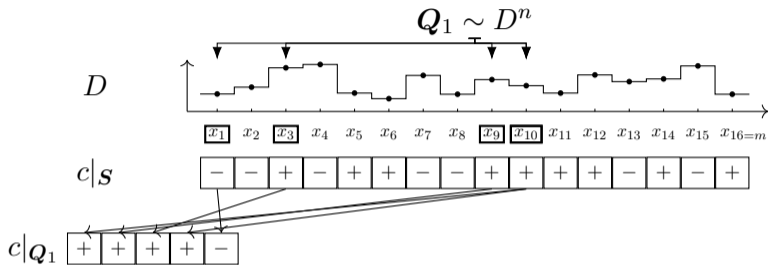


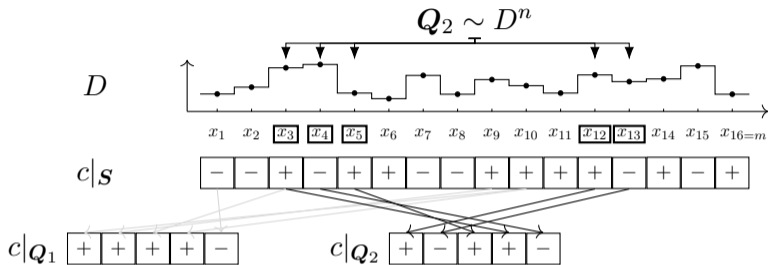


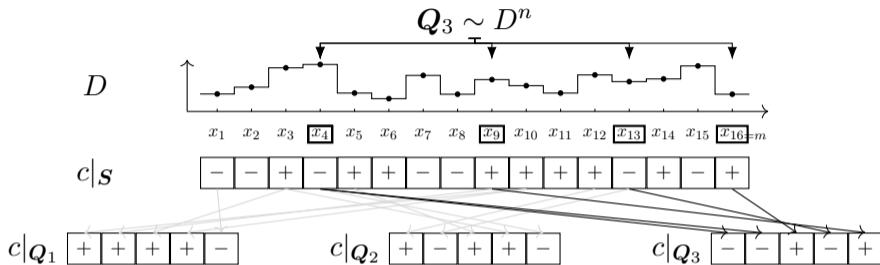
We'll need to understand how performance on a (sub)sample **generalises to the population (sample)**

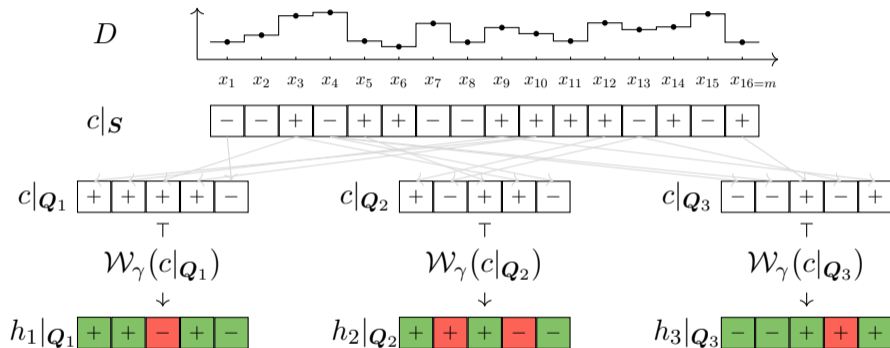


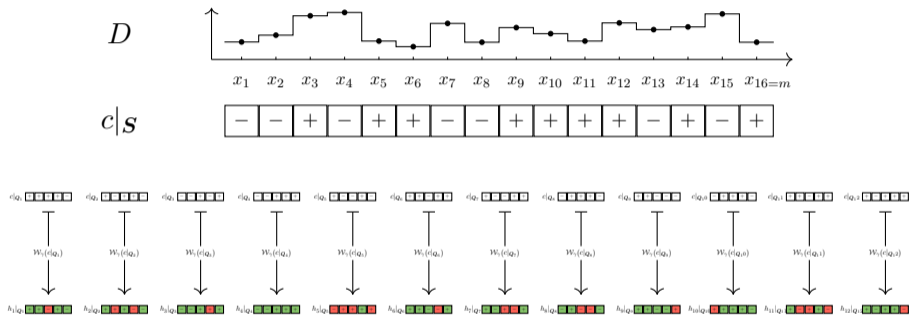


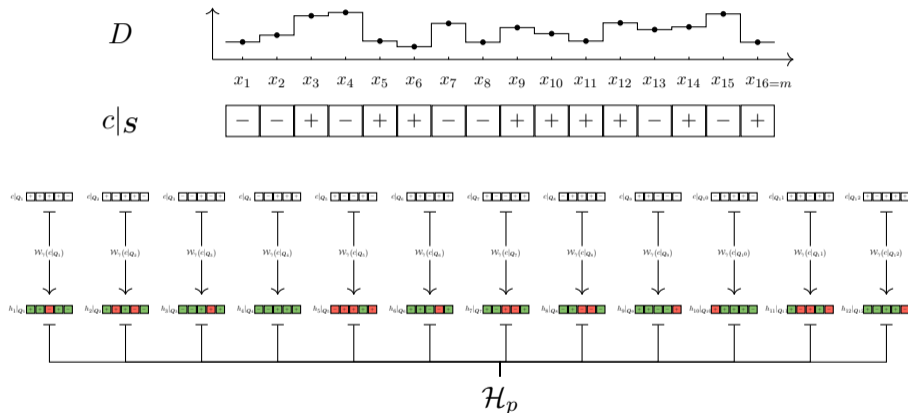












Bag of hypotheses used to perform **multiple** boosting steps

Input : Training data $c|_{\mathcal{S}}$, γ -weak learner \mathcal{W}_γ

Input : Training data $c|_{\mathcal{S}}$, γ -weak learner \mathcal{W}_γ

- 1 $\mathbf{D}_1 \leftarrow$ Uniform distribution over the m examples

Input : Training data $c|_{\mathcal{S}}$, γ -weak learner \mathcal{W}_γ

1 $D_1 \leftarrow$ Uniform distribution over the m examples

2 **for** $p \leftarrow 1$ **to** P **do**

|

Input : Training data $c|_{\mathcal{S}}$, γ -weak learner \mathcal{W}_γ

1 $D_1 \leftarrow$ Uniform distribution over the m examples

2 **for** $p \leftarrow 1$ **to** P **do**

 // Bagging step

3 $\mathcal{H}_p \leftarrow \emptyset$

 // Bag of weak hypotheses

Input : Training data $c|_{\mathcal{S}}$, γ -weak learner \mathcal{W}_γ

1 $D_1 \leftarrow$ Uniform distribution over the m examples

2 **for** $p \leftarrow 1$ **to** P **do**

 // Bagging step

3 $\mathcal{H}_p \leftarrow \emptyset$

 // Bag of weak hypotheses

4 **parallel for** $t \leftarrow 1$ **to** T **do**

 |

Input : Training data $c|_{\mathcal{S}}, \gamma$ -weak learner \mathcal{W}_γ

1 $D_1 \leftarrow$ Uniform distribution over the m examples

2 **for** $p \leftarrow 1$ **to** P **do**

 // Bagging step

3 $\mathcal{H}_p \leftarrow \emptyset$

 // Bag of weak hypotheses

4 **parallel for** $t \leftarrow 1$ **to** T **do**

5 $\mathbf{h} \leftarrow$ Query \mathcal{W}_γ on subsample following the current distribution ($D_{(p-1)R+1}$)

6 Add \mathbf{h} to \mathcal{H}_p

Input : Training data $c|_{\mathcal{S}}, \gamma$ -weak learner \mathcal{W}_γ

1 $D_1 \leftarrow$ Uniform distribution over the m examples

2 **for** $p \leftarrow 1$ **to** P **do**

 // Bagging step

3 $\mathcal{H}_p \leftarrow \emptyset$

 // Bag of weak hypotheses

4 **parallel for** $t \leftarrow 1$ **to** T **do**

5 $h \leftarrow$ Query \mathcal{W}_γ on subsample following the current distribution ($D_{(p-1)R+1}$)

6 Add h to \mathcal{H}_p

 // Boosting steps

Input : Training data $c|_{\mathcal{S}}$, γ -weak learner \mathcal{W}_γ

1 $D_1 \leftarrow$ Uniform distribution over the m examples

2 **for** $p \leftarrow 1$ **to** P **do**

 // Bagging step

3 $\mathcal{H}_p \leftarrow \emptyset$ // Bag of weak hypotheses

4 **parallel for** $t \leftarrow 1$ **to** T **do**

5 $h \leftarrow$ Query \mathcal{W}_γ on subsample following the current distribution ($D_{(p-1)R+1}$)

6 Add h to \mathcal{H}_p

 // Boosting steps

7 **for** $r \leftarrow 1$ **to** R **do**

8 $h_{(p-1)R+r} \leftarrow$ Simulate $\frac{\gamma}{2}$ -weak learner: search \mathcal{H}_p for h s.t. $\text{err}_{D_{(p-1)R+r}}(h) \leq \frac{1}{2} - \frac{\gamma}{2}$

Input : Training data $c|_S, \gamma$ -weak learner \mathcal{W}_γ

```

1  $D_1 \leftarrow$  Uniform distribution over the  $m$  examples
2 for  $p \leftarrow 1$  to  $P$  do
   // Bagging step
3  $\mathcal{H}_p \leftarrow \emptyset$  // Bag of weak hypotheses
4 parallel for  $t \leftarrow 1$  to  $T$  do
5    $h \leftarrow$  Query  $\mathcal{W}_\gamma$  on subsample following the current distribution ( $D_{(p-1)R+1}$ )
6   Add  $h$  to  $\mathcal{H}_p$ 
   // Boosting steps
7 for  $r \leftarrow 1$  to  $R$  do
8    $h_{(p-1)R+r} \leftarrow$  Simulate  $\frac{\gamma}{2}$ -weak learner: search  $\mathcal{H}_p$  for  $h$  s.t.  $\text{err}_{D_{(p-1)R+r}}(h) \leq \frac{1}{2} - \frac{\gamma}{2}$ 
9   ... // Omitted details

```


Input : Training data $c|_S, \gamma$ -weak learner \mathcal{W}_γ

```

1  $D_1 \leftarrow$  Uniform distribution over the  $m$  examples
2 for  $p \leftarrow 1$  to  $P$  do
   // Bagging step
3  $\mathcal{H}_p \leftarrow \emptyset$  // Bag of weak hypotheses
4 parallel for  $t \leftarrow 1$  to  $T$  do
5    $h \leftarrow$  Query  $\mathcal{W}_\gamma$  on subsample following the current distribution ( $D_{(p-1)R+1}$ )
6   Add  $h$  to  $\mathcal{H}_p$ 
   // Boosting steps
7 for  $r \leftarrow 1$  to  $R$  do
8    $h_{(p-1)R+r} \leftarrow$  Simulate  $\frac{\gamma}{2}$ -weak learner: search  $\mathcal{H}_p$  for  $h$  s.t.  $\text{err}_{D_{(p-1)R+r}}(h) \leq \frac{1}{2} - \frac{\gamma}{2}$ 
9   ... // Omitted details
10   $D_{(p-1)R+r+1} \leftarrow$  Usual "ADABOOST update" of  $D_{(p-1)R+r}$ 

```

Input : Training data $c|_S, \gamma$ -weak learner \mathcal{W}_γ

```

1  $D_1 \leftarrow$  Uniform distribution over the  $m$  examples
2 for  $p \leftarrow 1$  to  $P$  do
   // Bagging step
3  $\mathcal{H}_p \leftarrow \emptyset$  // Bag of weak hypotheses
4 parallel for  $t \leftarrow 1$  to  $T$  do
5    $h \leftarrow$  Query  $\mathcal{W}_\gamma$  on subsample following the current distribution ( $D_{(p-1)R+1}$ )
6   Add  $h$  to  $\mathcal{H}_p$ 
   // Boosting steps
7   for  $r \leftarrow 1$  to  $R$  do
8      $h_{(p-1)R+r} \leftarrow$  Simulate  $\frac{\gamma}{2}$ -weak learner: search  $\mathcal{H}_p$  for  $h$  s.t.  $\text{err}_{D_{(p-1)R+r}}(h) \leq \frac{1}{2} - \frac{\gamma}{2}$ 
9     ... // Omitted details
10     $D_{(p-1)R+r+1} \leftarrow$  Usual "ADABOOST update" of  $D_{(p-1)R+r}$ 
11 return Majority aggregation of  $\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_{PR}$ 

```

Input : Training data $c|_S, \gamma$ -weak learner \mathcal{W}_γ

```

1  $D_1 \leftarrow$  Uniform distribution over the  $m$  examples
2 for  $p \leftarrow 1$  to  $P$  do
   | // Bagging step
3    $\mathcal{H}_p \leftarrow \emptyset$  // Bag of weak hypotheses
4   parallel for  $t \leftarrow 1$  to  $T$  do
5     |  $\mathbf{h} \leftarrow$  Query  $\mathcal{W}_\gamma$  on subsample following the current distribution ( $D_{(p-1)R+1}$ )
6     | Add  $\mathbf{h}$  to  $\mathcal{H}_p$ 
   | // Boosting steps
7   for  $r \leftarrow 1$  to  $R$  do
8     |  $\mathbf{h}_{(p-1)R+r} \leftarrow$  Simulate  $\frac{\gamma}{2}$ -weak learner: search  $\mathcal{H}_p$  for  $h$  s.t.  $\text{err}_{D_{(p-1)R+r}}(h) \leq \frac{1}{2} - \frac{\gamma}{2}$ 
9     | ... // Omitted details
10    |  $D_{(p-1)R+r+1} \leftarrow$  Usual "ADABOOST update" of  $D_{(p-1)R+r}$ 
11 return Majority aggregation of  $\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_{PR}$ 

```

Main challenge

- Bag \mathcal{H}_p : hypotheses trained on samples from $\mathcal{D}_{(p-1)R+1}$.

Main challenge

- Bag \mathcal{H}_1 : hypotheses trained on samples from D_1 .

Main challenge

- Bag \mathcal{H}_1 : hypotheses trained on samples from D_1 .
- Easy to find good $\mathbf{h} \in \mathcal{H}_1$ for the first step.

Main challenge

- Bag \mathcal{H}_1 : hypotheses trained on samples from D_1 .
- Easy to find good $h \in \mathcal{H}_1$ for the first step.
- Then we change the distribution but not the bag.

Main challenge

- Bag \mathcal{H}_1 : hypotheses trained on samples from D_1 .
- Easy to find good $\mathbf{h} \in \mathcal{H}_1$ for the first step.
- Then we change the distribution but not the bag.

Solution

- Track the divergence from starting point: $\text{KL}(D_r \parallel D_1)$.

Main challenge

- Bag \mathcal{H}_1 : hypotheses trained on samples from D_1 .
- Easy to find good $h \in \mathcal{H}_1$ for the first step.
- Then we change the distribution but not the bag.

Solution

- Track the divergence from starting point: $KL(D_r \parallel D_1)$.
- Low divergence:

Main challenge

- Bag \mathcal{H}_1 : hypotheses trained on samples from D_1 .
- Easy to find good $h \in \mathcal{H}_1$ for the first step.
- Then we change the distribution but not the bag.

Solution

- Track the divergence from starting point: $KL(D_r \parallel D_1)$.
- Low divergence:
 - Lemma: likely to find good $h \in \mathcal{H}_1$.

Main challenge

- Bag \mathcal{H}_1 : hypotheses trained on samples from D_1 .
- Easy to find good $h \in \mathcal{H}_1$ for the first step.
- Then we change the distribution but not the bag.

Solution

- Track the divergence from starting point: $KL(D_r \parallel D_1)$.
- Low divergence:
 - Lemma: likely to find good $h \in \mathcal{H}_1$.
- Large divergence:

Main challenge

- Bag \mathcal{H}_1 : hypotheses trained on samples from D_1 .
- Easy to find good $h \in \mathcal{H}_1$ for the first step.
- Then we change the distribution but not the bag.

Solution

- Track the divergence from starting point: $KL(D_r \parallel D_1)$.
- Low divergence:
 - Lemma: likely to find good $h \in \mathcal{H}_1$.
- Large divergence:
 - Hard to find good $h \in \mathcal{H}_1$.

Main challenge

- Bag \mathcal{H}_1 : hypotheses trained on samples from D_1 .
- Easy to find good $h \in \mathcal{H}_1$ for the first step.
- Then we change the distribution but not the bag.

Solution

- Track the divergence from starting point: $KL(D_r \parallel D_1)$.
- Low divergence:
 - Lemma: likely to find good $h \in \mathcal{H}_1$.
- Large divergence:
 - Hard to find good $h \in \mathcal{H}_1$.
 - Perhaps some $h \in \mathcal{H}_1$ is so bad that $-h$ is good.

Main challenge

- Bag \mathcal{H}_1 : hypotheses trained on samples from D_1 .
- Easy to find good $\mathbf{h} \in \mathcal{H}_1$ for the first step.
- Then we change the distribution but not the bag.

Solution

- Track the divergence from starting point: $\text{KL}(D_r \parallel D_1)$.
- Low divergence:
 - Lemma: likely to find good $\mathbf{h} \in \mathcal{H}_1$.
- Large divergence:
 - Hard to find good $\mathbf{h} \in \mathcal{H}_1$.
 - Perhaps some $\mathbf{h} \in \mathcal{H}_1$ is so bad that $-\mathbf{h}$ is good.
 - Otherwise already made enough progress: early stopping.

Recall:

- $T :=$ Number of parallel calls to \mathcal{W}_γ .
- $R :=$ Number of steps of “simulated” $\gamma/2$ -weak learner.
- $P :=$ Number of iterations of those.
- $d :=$ VC(“base classifiers”).

Results

Recall:

- $T :=$ Number of parallel calls to \mathcal{W}_γ .
- $R :=$ Number of steps of “simulated” $\gamma/2$ -weak learner.
- $P :=$ Number of iterations of those.
- $d := \text{VC}(\text{“base classifiers”})$.

Results

- Given $R \in \mathbb{N}$,

Recall:

- $T :=$ Number of parallel calls to \mathcal{W}_γ .
- $R :=$ Number of steps of “simulated” $\gamma/2$ -weak learner.
- $P :=$ Number of iterations of those.
- $d :=$ VC(“base classifiers”).

Results

- Given $R \in \mathbb{N}$,
 - With high probability, the algorithm described performs well (generalization error no worse than ADABOOST's).

Recall:

- $T :=$ Number of parallel calls to \mathcal{W}_γ .
- $R :=$ Number of steps of “simulated” $\gamma/2$ -weak learner.
- $P :=$ Number of iterations of those.
- $d :=$ VC(“base classifiers”).

Results

- Given $R \in \mathbb{N}$,
 - With high probability, the algorithm described performs well (generalization error no worse than ADABOOST's).
 - Satisfies
 - $P = O\left(\frac{\ln|S|}{\gamma^2 \cdot R}\right)$,
 - $T = e^{O(d \cdot R)}$.

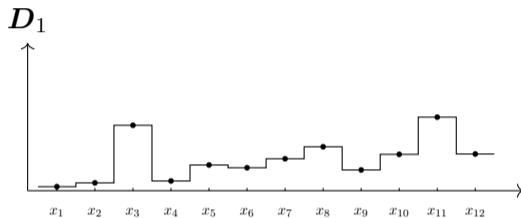
Recall:

- $T :=$ Number of parallel calls to \mathcal{W}_γ .
- $R :=$ Number of steps of “simulated” $\gamma/2$ -weak learner.
- $P :=$ Number of iterations of those.
- $d :=$ VC(“base classifiers”).

Results

- Given $R \in \mathbb{N}$,
 - With high probability, the algorithm described performs well (generalization error no worse than ADABOOST's).
 - Satisfies
 - $P = O\left(\frac{\ln|S|}{\gamma^2 \cdot R}\right)$,
 - $T = e^{O(d \cdot R)}$.
- Matching lower bounds (up to logarithmic factors) for all values of R .

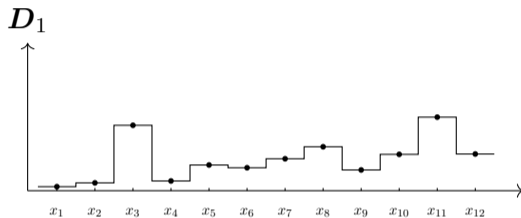
We'll be presenting this work's poster in 20 minutes from now (at West Ballroom A-D).
Come chat with us!

 $c|s$

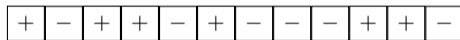
+	-	+	+	-	+	-	-	-	+	+	-
---	---	---	---	---	---	---	---	---	---	---	---

Goal Simulate a $\gamma/2$ -weak learner.

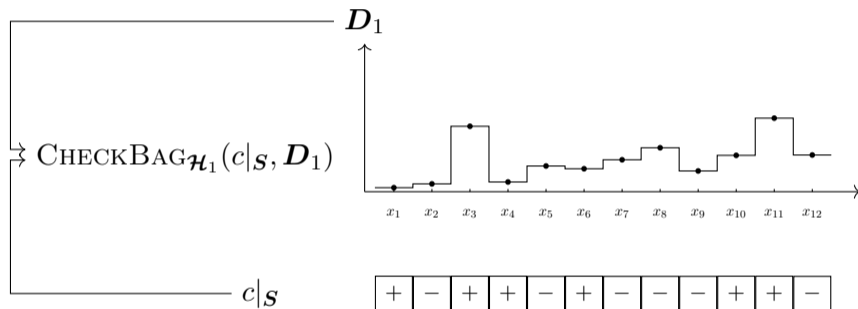
CHECKBAG \mathcal{H}_1



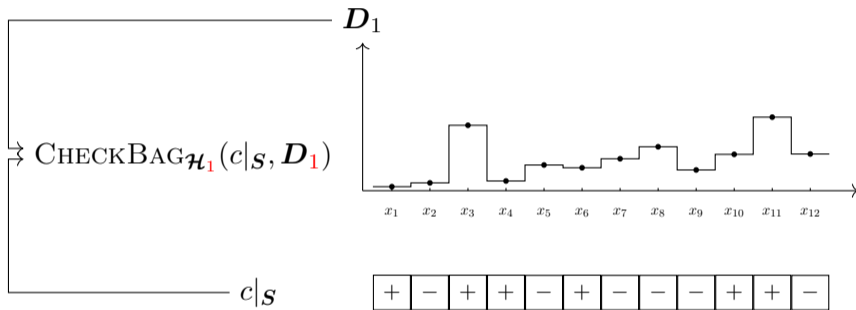
$c|_S$



Goal Simulate a $\gamma/2$ -weak learner.

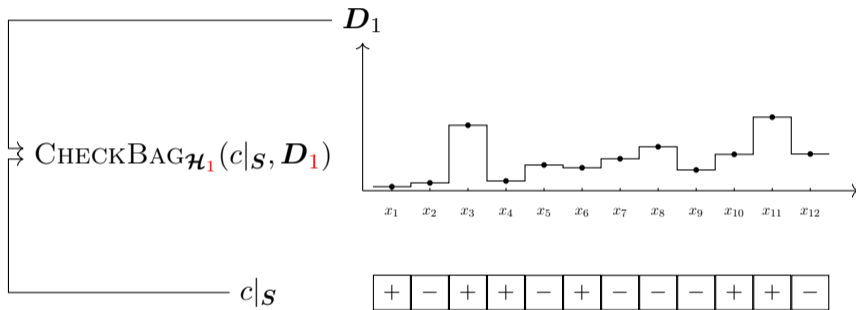


Goal Simulate a $\gamma/2$ -weak learner.



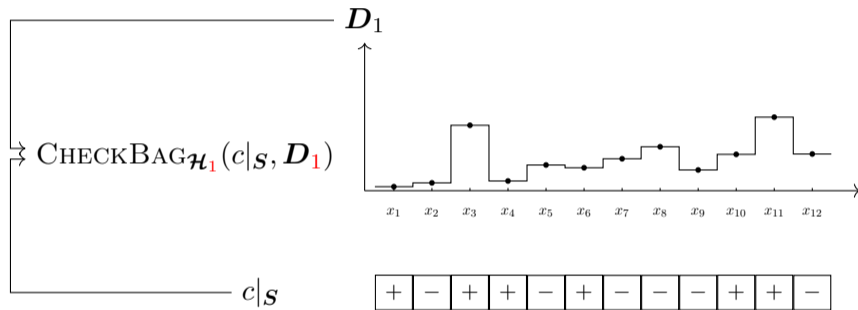
- \mathcal{H}_1 contains weak-hypotheses for subsamples following D_1 .

Goal Simulate a $\gamma/2$ -weak learner.



- \mathcal{H}_1 contains weak-hypotheses for subsamples following D_1 .
- Classical LT: large enough ($O(d/\gamma^2)$) samples are likely to be (γ -)representative

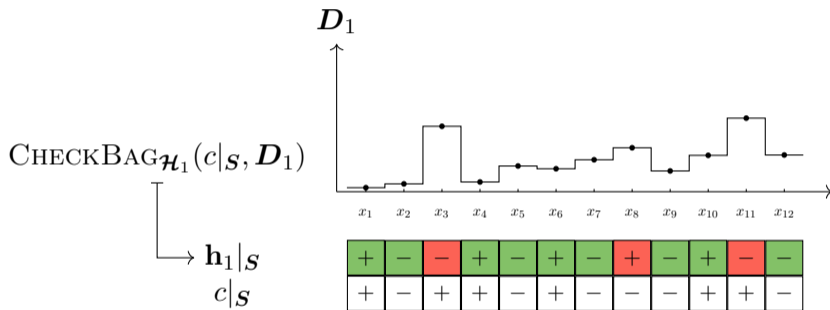
Goal Simulate a $\gamma/2$ -weak learner.



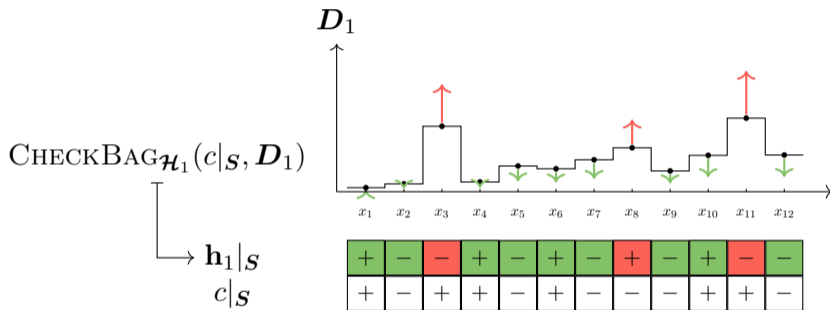
- \mathcal{H}_1 contains weak-hypotheses for subsamples following D_1 .
- Classical LT: large enough ($O(d/\gamma^2)$) samples are likely to be (γ -)representative:

$$|\text{err}_{Q \sim D_1^n}(h) - \text{err}_{D_1}(h)| < \gamma/2 \quad (\text{with high probability})$$

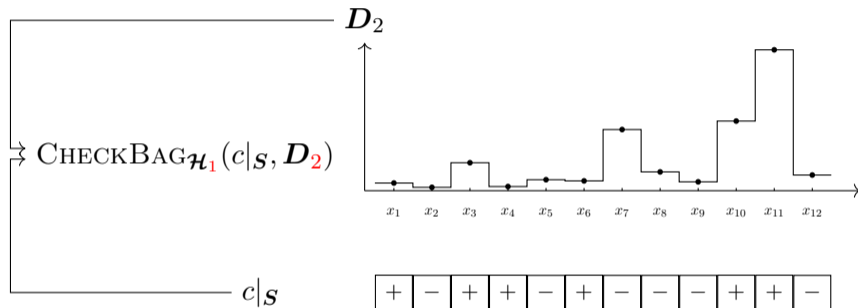
Goal Simulate a $\gamma/2$ -weak learner.



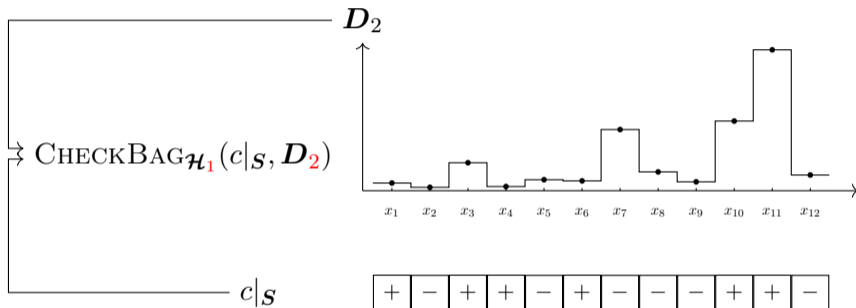
Goal Simulate a $\gamma/2$ -weak learner.



Goal Simulate a $\gamma/2$ -weak learner.

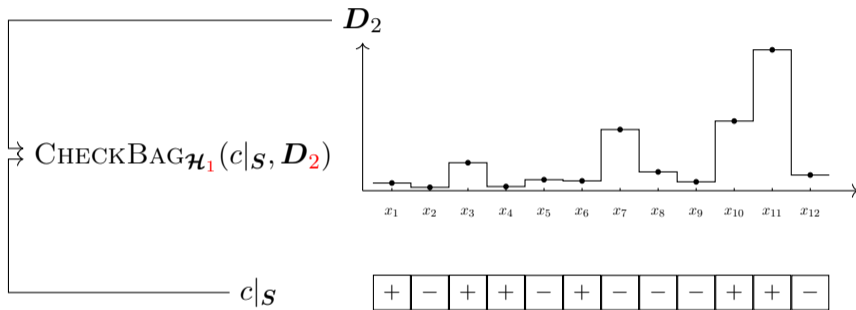


Goal Simulate a $\gamma/2$ -weak learner.



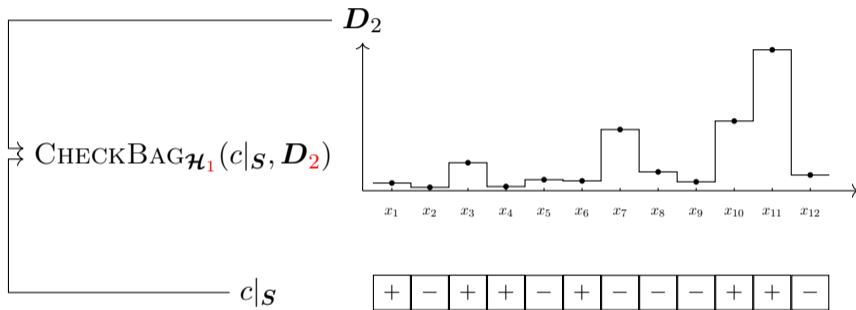
- \mathcal{H}_1 contains weak-hypotheses for subsamples following D_1 .

Goal Simulate a $\gamma/2$ -weak learner.



- \mathcal{H}_1 contains weak-hypotheses for subsamples following D_1 .
- Does performance on subsamples from D_1 generalise to performance under D_2 ?

Goal Simulate a $\gamma/2$ -weak learner.



- \mathcal{H}_1 contains weak-hypotheses for subsamples following D_1 .
- Does performance on subsamples from D_1 generalise to performance under D_2, D_3, \dots ?

$$|\text{err}_{Q \sim D_1^n}(h) - \text{err}_{D_r}(h)| < ??? \quad (\text{with high probability})$$

-  Chen, Tianqi and Carlos Guestrin (2016). 'XGBoost: A Scalable Tree Boosting System.'. In: *KDD*. ACM, pp. 785–794. ISBN: 978-1-4503-4232-2.
-  Karbasi, Amin and Kasper Green Larsen (2024). 'The Impossibility of Parallelizing Boosting'. In: *Proceedings of The 35th International Conference on Algorithmic Learning Theory*. Ed. by Claire Vernade and Daniel Hsu. Vol. 237. Proceedings of Machine Learning Research. PMLR, pp. 635–653. URL: <https://proceedings.mlr.press/v237/karbasi24a.html>.
-  Ke, Guolin et al. (2017). 'LightGBM: A Highly Efficient Gradient Boosting Decision Tree'. In: *NIPS*.
-  Lyu, Xin, Hongxun Wu and Junzhao Yang (2024). 'The Cost of Parallelizing Boosting'. In: *Proceedings of the 2024 ACM-SIAM Symposium on Discrete Algorithms, SODA 2024, Alexandria, VA, USA, January 7-10, 2024*. Ed. by David P. Woodruff. SIAM, pp. 3140–3155. DOI: 10.1137/1.9781611977912.112. URL: <https://doi.org/10.1137/1.9781611977912.112>.