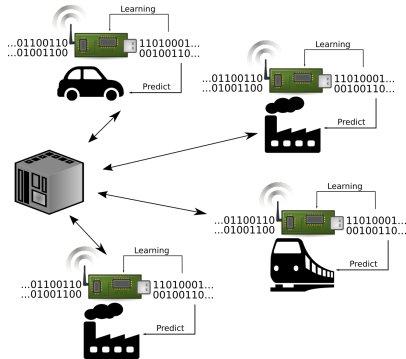


Shrub Ensembles for Online Classification

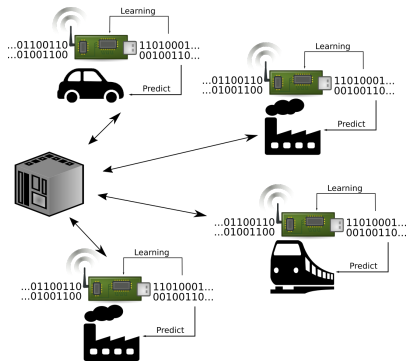
Sebastian Buschjäger, Sibylle Hess and Katharina Morik

Thirty-Sixth AAAI Conference on Artificial Intelligence 2022

Artificial Intelligence Group@TU Dortmund University  - Data Mining Group@ TU Eindhoven **TU/e** - Collaborative Research Center 876 



In 2018 Roughly 22 Billion IoT Devices world wide
By 2025 \approx 38 Billion IoT Devices world wide¹



In 2018 Roughly 22 Billion IoT Devices world wide
By 2025 \approx 38 Billion IoT Devices world wide¹

Clear benefits

- + Privacy and independence of data analysis
- + Reduced communication infrastructure costs
- + Faster response times

¹<https://findstack.com/internet-of-things-statistics/>
Shrub Ensembles for Online Classification by Sebastian Buschjäger, Sibylle Hess and Katharina Morik

Computational efficiency

The algorithm must process examples at least as fast as new examples arrive.

Computational efficiency

The algorithm must process examples at least as fast as new examples arrive.

Memory efficiency

The algorithm has only a limited budget of memory and fails if more memory is required.

Computational efficiency

The algorithm must process examples at least as fast as new examples arrive.

Memory efficiency

The algorithm has only a limited budget of memory and fails if more memory is required.

Evolving data streams

The algorithm must adapt to changes in the new data distribution and preserve its performance.

Online Decision Tree Ensemble Learning

Online Decision Tree Ensemble Learning



Gradient-based Learners

- + constant memory consumption
- + joint optimization of all trees
- backpropagation through entire tree

Online Decision Tree Ensemble Learning



Gradient-based Learners

- + constant memory consumption
- + joint optimization of all trees
- backpropagation through entire tree

Incremental Learners

- + fast
- + proven in practice
- nodes are not removed

Online Decision Tree Ensemble Learning



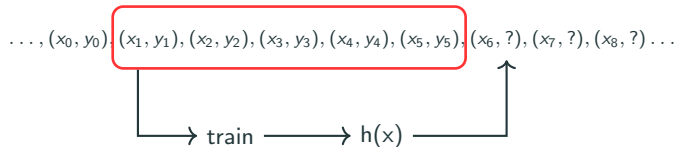
Gradient-based Learners

- + constant memory consumption
- + joint optimization of all trees
- backpropagation through entire tree

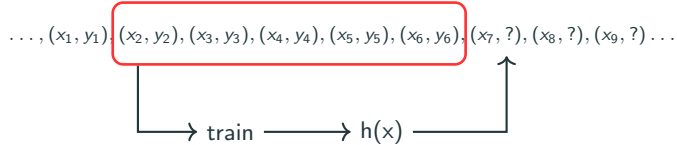
Incremental Learners

- + fast
- + proven in practice
- nodes are not removed

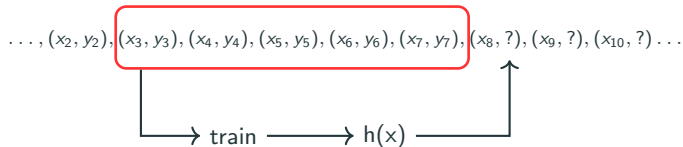
Can we design simple and small online ensembles?



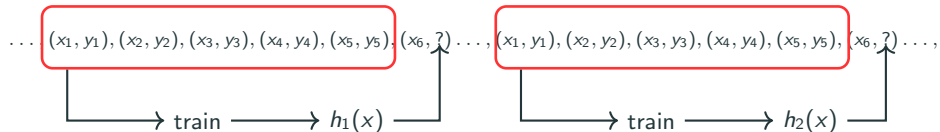
- + We can use a batch algorithm to train $h(x)$
- + We quickly adapt to concept drift
- No long-term learning possible



- + We can use a batch algorithm to train $h(x)$
- + We quickly adapt to concept drift
- No long-term learning possible

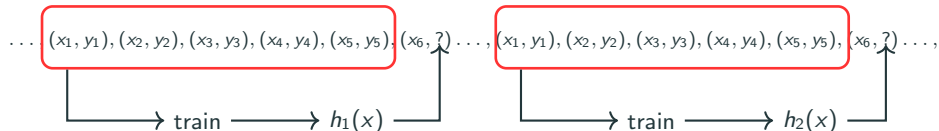


- + We can use a batch algorithm to train $h(x)$
- + We quickly adapt to concept drift
- No long-term learning possible



Often Local Patterns repeat over long running processes

Hence The two classifiers h_1 and h_2 are redundant. We should use h_1



Often Local Patterns repeat over long running processes

Hence The two classifiers h_1 and h_2 are redundant. We should use h_1

Idea Learn an ensemble of local patterns from sliding windows

$$f(x) = \sum_{i=1}^M w_i h_i(x)$$

Formally Let $\mathcal{H} = \{h_1, h_2, \dots, h_K\}$ be the set of trees learned on local patterns:

$$\arg \min_{w \in \mathbb{R}^K} \sum_{t=1}^T \ell(f_{S[0:t-1]}(x_t), y_t) \quad \text{s.t.} \quad \|w\|_0 \leq M, w_i \geq 0, \sum_{i=1}^K w_i = 1$$

with

- $f_{S[0:t-1]}: \mathbb{R}^d \rightarrow \mathbb{R}^C$ is the model at time t
- $\ell: \mathbb{R}^C \times \mathcal{Y} \rightarrow \mathbb{R}_+$ is the loss

Formally Let $\mathcal{H} = \{h_1, h_2, \dots, h_K\}$ be the set of trees learned on local patterns:

$$\arg \min_{w \in \mathbb{R}^K} \sum_{t=1}^T \ell(f_{S[0:t-1]}(x_t), y_t) \quad \text{s.t.} \quad \|w\|_0 \leq M, w_i \geq 0, \sum_{i=1}^K w_i = 1$$

with

- $f_{S[0:t-1]}: \mathbb{R}^d \rightarrow \mathbb{R}^C$ is the model at time t
- $\ell: \mathbb{R}^C \times \mathcal{Y} \rightarrow \mathbb{R}_+$ is the loss
- $M \geq 1$ is the maximum number of ensemble members

Formally Let $\mathcal{H} = \{h_1, h_2, \dots, h_K\}$ be the set of trees learned on local patterns:

$$\arg \min_{w \in \mathbb{R}^K} \sum_{t=1}^T \ell(f_{S[0:t-1]}(x_t), y_t) \quad \text{s.t.} \quad \|w\|_0 \leq M, w_i \geq 0, \sum_{i=1}^K w_i = 1$$

with

- $f_{S[0:t-1]}: \mathbb{R}^d \rightarrow \mathbb{R}^C$ is the model at time t
- $\ell: \mathbb{R}^C \times \mathcal{Y} \rightarrow \mathbb{R}_+$ is the loss
- $M \geq 1$ is the maximum number of ensemble members
- $\|w\|_0 = \sum_{i=1}^K \mathbb{1}\{w_i \neq 0\}$ is the 0-norm

Formally Let $\mathcal{H} = \{h_1, h_2, \dots, h_K\}$ be the set of trees learned on local patterns:

$$\arg \min_{w \in \mathbb{R}^K} \sum_{t=1}^T \ell(f_{S[0:t-1]}(x_t), y_t) \quad \text{s.t.} \quad \|w\|_0 \leq M, w_i \geq 0, \sum_{i=1}^K w_i = 1$$

with

- $f_{S[0:t-1]}: \mathbb{R}^d \rightarrow \mathbb{R}^C$ is the model at time t
- $\ell: \mathbb{R}^C \times \mathcal{Y} \rightarrow \mathbb{R}_+$ is the loss
- $M \geq 1$ is the maximum number of ensemble members
- $\|w\|_0 = \sum_{i=1}^K \mathbb{1}\{w_i \neq 0\}$ is the 0-norm

Formally Let $\mathcal{H} = \{h_1, h_2, \dots, h_K\}$ be the set of trees learned on local patterns:

$$\arg \min_{w \in \mathbb{R}^K} \sum_{t=1}^T \ell(f_{S[0:t-1]}(x_t), y_t) \quad \text{s.t.} \quad \|w\|_0 \leq M, w_i \geq 0, \sum_{i=1}^K w_i = 1$$

with

- $f_{S[0:t-1]}: \mathbb{R}^d \rightarrow \mathbb{R}^C$ is the model at time t
- $\ell: \mathbb{R}^C \times \mathcal{Y} \rightarrow \mathbb{R}_+$ is the loss
- $M \geq 1$ is the maximum number of ensemble members
- $\|w\|_0 = \sum_{i=1}^K \mathbb{1}\{w_i \neq 0\}$ is the 0-norm

For this talk Focus on the MSE loss

$$\ell(f_{S[0:t-1]}(x_t), y_t) = \frac{1}{C} \|f_{S[0:t-1]}(x_t) - y_t\|^2$$

Proximal Gradient Descent

$$w \leftarrow \mathcal{P}(w - \alpha \nabla_w L_{\mathcal{B}}(w)),$$

where

- \mathcal{B} is the current window with $|\mathcal{B}| = B$ examples

Proximal Gradient Descent

$$w \leftarrow \mathcal{P}(w - \alpha \nabla_w L_{\mathcal{B}}(w)),$$

where

- \mathcal{B} is the current window with $|\mathcal{B}| = B$ examples
- $\nabla_w L_{\mathcal{B}}(w)$ is the gradient on the current window

Proximal Gradient Descent

$$w \leftarrow \mathcal{P}(w - \alpha \nabla_w L_{\mathcal{B}}(w)),$$

where

- \mathcal{B} is the current window with $|\mathcal{B}| = B$ examples
- $\nabla_w L_{\mathcal{B}}(w)$ is the gradient on the current window
- $\alpha \in \mathbb{R}_+$ is the step-size

Proximal Gradient Descent

$$w \leftarrow \mathcal{P}(w - \alpha \nabla_w L_{\mathcal{B}}(w)),$$

where

- \mathcal{B} is the current window with $|\mathcal{B}| = B$ examples
- $\nabla_w L_{\mathcal{B}}(w)$ is the gradient on the current window
- $\alpha \in \mathbb{R}_+$ is the step-size
- $\mathcal{P}(w)$ is the prox-operator for the feasible set $\Delta = \left\{ w \in \mathbb{R}_+^K \mid \sum_{i=1}^K w_i = 1, \|w\|_0 = M \right\}$

Proximal Gradient Descent

$$w \leftarrow \mathcal{P}(w - \alpha \nabla_w L_{\mathcal{B}}(w)),$$

where

- \mathcal{B} is the current window with $|\mathcal{B}| = B$ examples
- $\nabla_w L_{\mathcal{B}}(w)$ is the gradient on the current window
- $\alpha \in \mathbb{R}_+$ is the step-size
- $\mathcal{P}(w)$ is the prox-operator for the feasible set $\Delta = \left\{ w \in \mathbb{R}_+^K \mid \sum_{i=1}^K w_i = 1, \|w\|_0 = M \right\}$

Our paper / Kyrillidis et al. 2013 Details for prox-operator

Idea If a tree is helpful to the ensemble, then it should have a large weight

Algorithm 1: Shrub Ensembles.

```
1:  $w \leftarrow (0)$ ;  $\mathcal{B} \leftarrow []$ ;  $\mathcal{H} \leftarrow []$    ▷ Init.
2: for next item  $(x, y)$  do
3:   if  $|\mathcal{B}| = B$  then                         ▷ Update batch
4:      $\mathcal{B}.\text{pop\_first}()$ 
5:      $\mathcal{B}.\text{append}((x, y))$ 
6:      $h_{\text{new}} \leftarrow \text{train}(\mathcal{B})$              ▷ Add new classifier
7:      $\mathcal{H}.\text{append}(h_{\text{new}})$ 
8:      $w \leftarrow (w_1, \dots, w_M, 0)$            ▷ Initialize weight
9:      $w \leftarrow w - \alpha \nabla_w L_{\mathcal{B}}(w)$      ▷ Gradient step
10:     $w, \mathcal{H} \leftarrow \text{sorted}(w, \mathcal{H})$      ▷ Sort decreasing order
11:     $w \leftarrow \mathcal{P}(w)$                        ▷ Project on feasible set
12:     $w, \mathcal{H} \leftarrow \text{prune}(w, \mathcal{H})$     ▷ Remove zero weights
```

Runtime $\mathcal{O}(dB^2 \log B + \log M)$ per example with d features

Runtime $\mathcal{O}(dB^2 \log B + \log M)$ per example with d features

CART \longleftarrow \uparrow

Runtime $\mathcal{O}(dB^2 \log B + \log M)$ per example with d features

CART \longleftarrow $dB^2 \log B$ \longleftarrow $\log M$ prox operator

Runtime $\mathcal{O}(dB^2 \log B + \log M)$ per example with d features

CART \longleftarrow $dB^2 \log B$ \longleftarrow $\log M$ prox operator

Memory $\mathcal{O}(dB + 2 \cdot B \cdot (M + 1))$ per example with d features

Runtime $\mathcal{O}(dB^2 \log B + \log M)$ per example with d features

CART \longleftarrow $dB^2 \log B$ \longleftarrow $\log M$ prox operator

Memory $\mathcal{O}(dB + 2 \cdot B \cdot (M + 1))$ per example with d features

window \longleftarrow dB

Runtime $\mathcal{O}(dB^2 \log B + \log M)$ per example with d features

CART \longleftarrow $dB^2 \log B$ \longleftarrow prox operator $\longleftarrow \log M$

Memory $\mathcal{O}(dB + 2 \cdot B \cdot (M + 1))$ per example with d features

window $\longleftarrow dB$ \longleftarrow trees \rightarrow shrubs $\longleftarrow 2 \cdot B \cdot (M + 1)$

Runtime $\mathcal{O}(dB^2 \log B + \log M)$ per example with d features

CART \longleftarrow $dB^2 \log B$ \longleftarrow prox operator \longleftarrow $\log M$

Memory $\mathcal{O}(dB + 2 \cdot B \cdot (M + 1))$ per example with d features

window \longleftarrow dB \longleftarrow trees \rightarrow shrubs \longleftarrow $2 \cdot B \cdot (M + 1)$

Behaviour Let $m \leq M$ be the number of models in the ensemble and let $\forall j = 1, \dots, m: h_j(x_B) \neq y_B$. If SE trains fully-grown trees such that $\forall i = 1, \dots, B: h(x_i) = y_i$ it holds for $\alpha > \frac{BC}{4m}$ that:

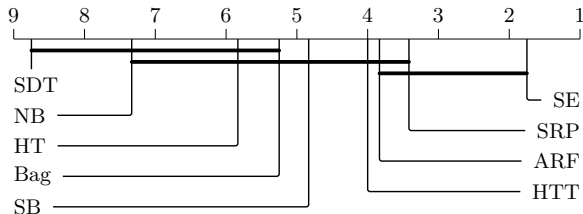
- (1) If $m < M$, then h is added to the ensemble
- (2) If $m = M$, then h replaces the tree with the smallest weight from the ensemble

Goal Compare Accuracy-Memory Trade-off of different configurations

- 1) Plot Pareto Front of best performing configurations against accuracy
- 2) Compute Area Under the Pareto Front to measure accuracy-memory trade-off
- 3) Rank each method according to its trade-off and plot a CD-diagram

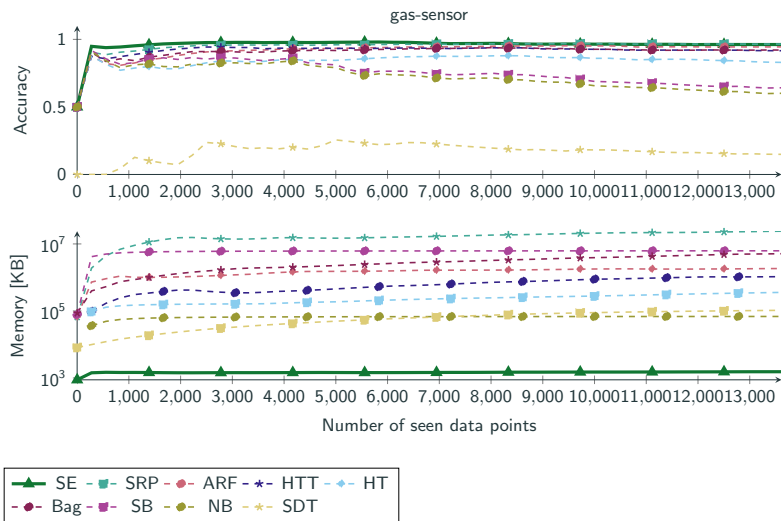
Goal Compare Accuracy-Memory Trade-off of different configurations

- 1) Plot Pareto Front of best performing configurations against accuracy
- 2) Compute Area Under the Pareto Front to measure accuracy-memory trade-off
- 3) Rank each method according to its trade-off and plot a CD-diagram

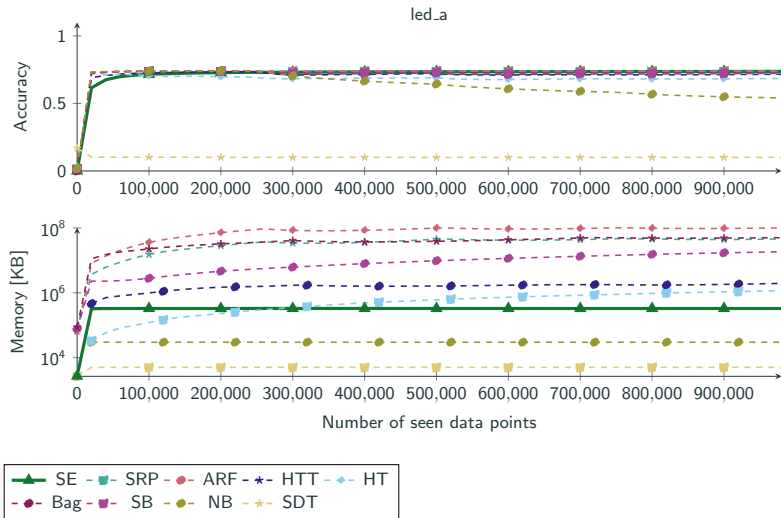


- SE: Shrub Ensembles
- SRP: Streaming Random Patches
- ARP: Adaptive Random Forest
- HTT: HoeffdingAnyTree
- SB: (Online) Smooth Boost
- Bag: (Online) Bagging
- HT: HoeffdingTrees
- NB: (Online) NaiveBayes
- SDT: Soft Decision Trees

Experimental Insights: Qualitative Analysis (1)



Experimental Insights: Qualitative Analysis (2)



Edge learning becomes more important every year

- Computational efficiency
- Memory efficiency
- Evolving data streams

Edge learning becomes more important every year

- Computational efficiency
- Memory efficiency
- Evolving data streams

Current approaches Unbounded memory or costly gradients

Edge learning becomes more important every year

- Computational efficiency
- Memory efficiency
- Evolving data streams

Current approaches Unbounded memory or costly gradients

Shrub Ensembles Bounded memory and simple gradients

- Train small trees (\rightarrow Shrubs) on sliding window
- Maintain tree weights via proximal gradient descent
- Prune unimportant trees

Results Better accuracy-memory trade-off than existing methods!