# Overview

Ensemble Methods

- Majority Voting
- Bagging
- Boosting
- Random Forests
- Stacking

# Boosting

# General Boosting

**Training Sample** $\longrightarrow$ $h_1(\mathbf{x})$

**Weighted Training Sample** $\longrightarrow$ $h_2(\mathbf{x})$

**Weighted Training Sample** $\longrightarrow$ $h_m(\mathbf{x})$

$$h_m(\mathbf{x}) = sign\left( \sum_{j=1}^{m} w_j \, h_j(\mathbf{x}) \right) \qquad \text{for} \quad h(\mathbf{x}) \in \{-1,1\}$$

$$\text{or} \quad h_m(\mathbf{x}) = \arg\max_i \left( \sum_{j=1}^{m} w_j \, \mathbf{1}[h_j(\mathbf{x}) = i] \right) \text{for} \quad h(\mathbf{x}) = i, \qquad i \in \{1,...,n\}$$

# General Boosting

- Initialize a weight vector with uniform weights

- Loop:
  - Apply weak learner* to weighted training examples (instead of orig. training set, may draw bootstrap samples with weighted probability)

  - Increase weight for misclassified examples

- (Weighted) majority voting on trained classifiers

* a learner slightly better than random guessing

# AdaBoost

**Algorithm 1** AdaBoost

1: Initialize $k$: the number of AdaBoost rounds
2: Initialize $\mathcal{D}$: the training dataset, $\mathcal{D} = \{\langle \mathbf{x}^{[1]}, y^{[1]} \rangle, ..., \mathbf{x}^{[n]}, y^{[n]} \rangle\}$
3: Initialize $w_1(i) = 1/n, \quad i = 1, ..., n, \ \mathbf{w}_1 \in \mathbb{R}^n$

4:

5: **for** r=1 to $k$ **do**
6:      For all $i : \mathbf{w}_r(i) := w_r(i) / \sum_i w_r(i)$     [normalize weights]
7:      $h_r := FitWeakLearner(\mathcal{D}, \mathbf{w}_r)$
8:      $\epsilon_r := \sum_i w_r(i) \, \mathbf{1}(h_r(i) \neq y_i)$     [compute error]
9:      if $\epsilon_r > 1/2$ then stop
10:      $\alpha_r := \frac{1}{2} \log[(1 - \epsilon_r)/\epsilon_r]$     [small if error is large and vice versa]

11:      $w_{r+1}(i) := w_r(i) \times \begin{cases} e^{-\alpha_r} & \text{if } h_r(\mathbf{x}^{[i]}) = y^{[i]} \\ e^{\alpha_r} & \text{if } h_r(\mathbf{x}^{[i]}) \neq y^{[i]} \end{cases}$

12: Predict: $h_{ens}(\mathbf{x}) = \arg\max_j \sum_r^k \alpha_r \mathbf{1}[h_r(\mathbf{x}) = j]$
13:

# AdaBoost

0/1 loss

$$\mathbf{1}\left(h_r(i) \neq y_i\right) = \begin{cases} 0 & \textbf{if } h_r(i) = y_i \\ 1 & \textbf{if } h_r(i) \neq y_i \end{cases}$$

---

**Algorithm 1** AdaBoost

---

1: Initialize $k$: the number of AdaBoost rounds

2: Initialize $\mathcal{D}$: the training dataset, $\mathcal{D} = \{\langle \mathbf{x}^{[1]}, y^{[1]} \rangle, ..., \mathbf{x}^{[n]}, y^{[n]} \rangle\}$

3: Initialize $w_1(i) = 1/n, \quad i = 1, ..., n, \ \mathbf{w}_1 \in \mathbb{R}^n$

4:

5: **for** r=1 to $k$ **do**

6:      For all $i : \mathbf{w}_r(i) := w_r(i) / \sum_i w_r(i)$     [normalize weights]

7:      $h_r := FitWeakLearner(\mathcal{D}, \mathbf{w}_r)$

8:      $\epsilon_r := \sum_i w_r(i) \, \mathbf{1}(h_r(i) \neq y_i)$     [compute error]

9:      if $\epsilon_r > 1/2$ then stop

10:      $\alpha_r := \frac{1}{2} \log[(1 - \epsilon_r)/\epsilon_r]$     [small if error is large and vice versa]

11:      $w_{r+1}(i) := w_r(i) \times \begin{cases} \mathrm{e}^{-\alpha_r} & \text{if } h_r(\mathbf{x}^{[i]}) = y^{[i]} \\ \mathrm{e}^{\alpha_r} & \text{if } h_r(\mathbf{x}^{[i]}) \neq y^{[i]} \end{cases}$

12: Predict: $h_{ens}(\mathbf{x}) = \arg\max_j \sum_r^k \alpha_r \mathbf{1}[h_r(\mathbf{x}) = j]$
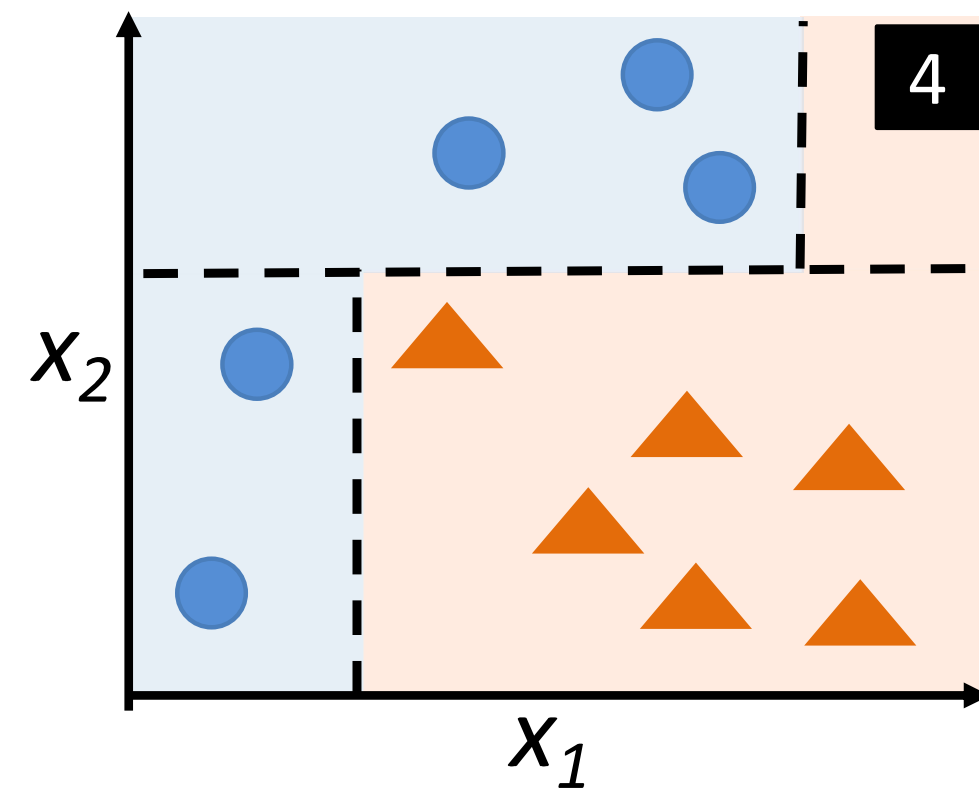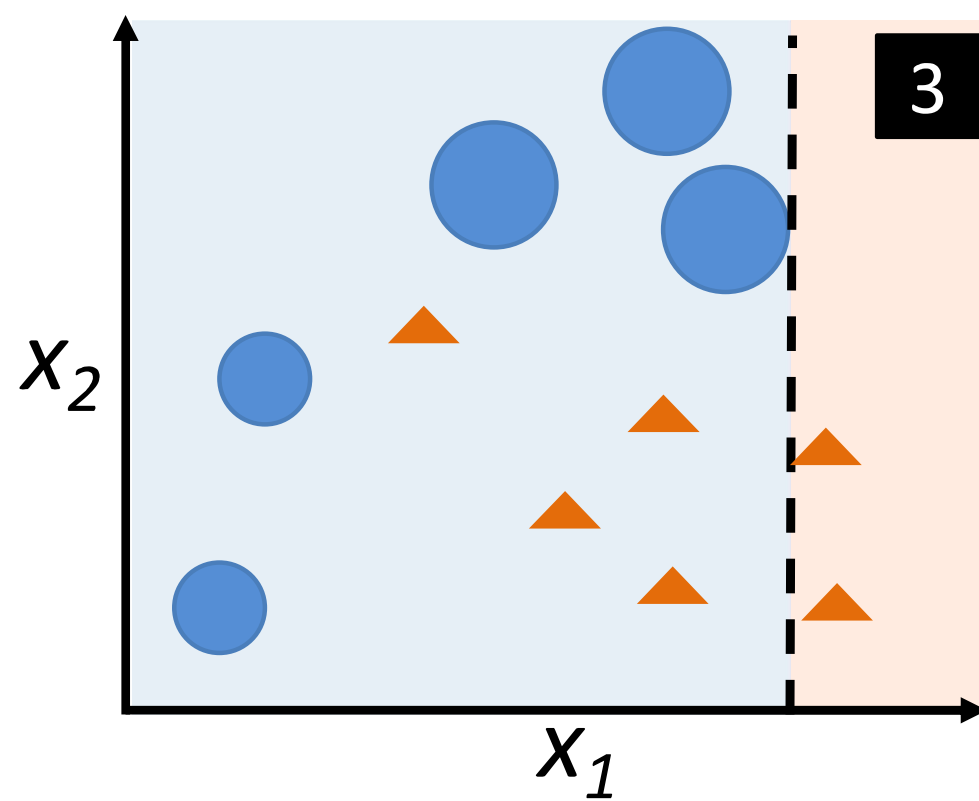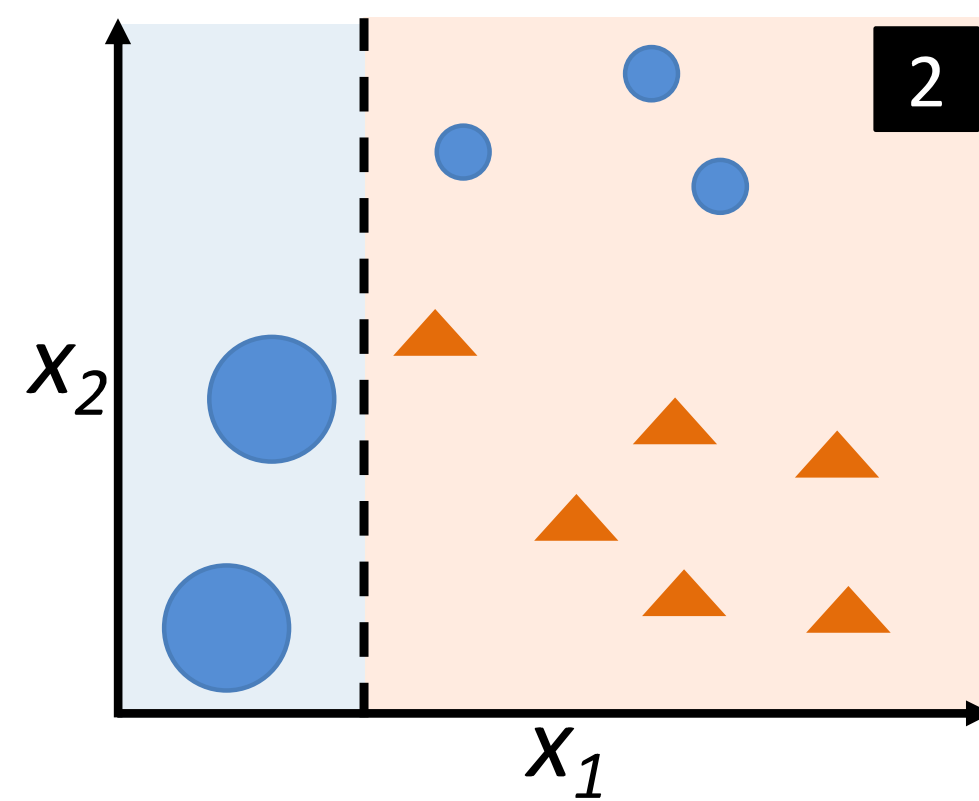
13:

---
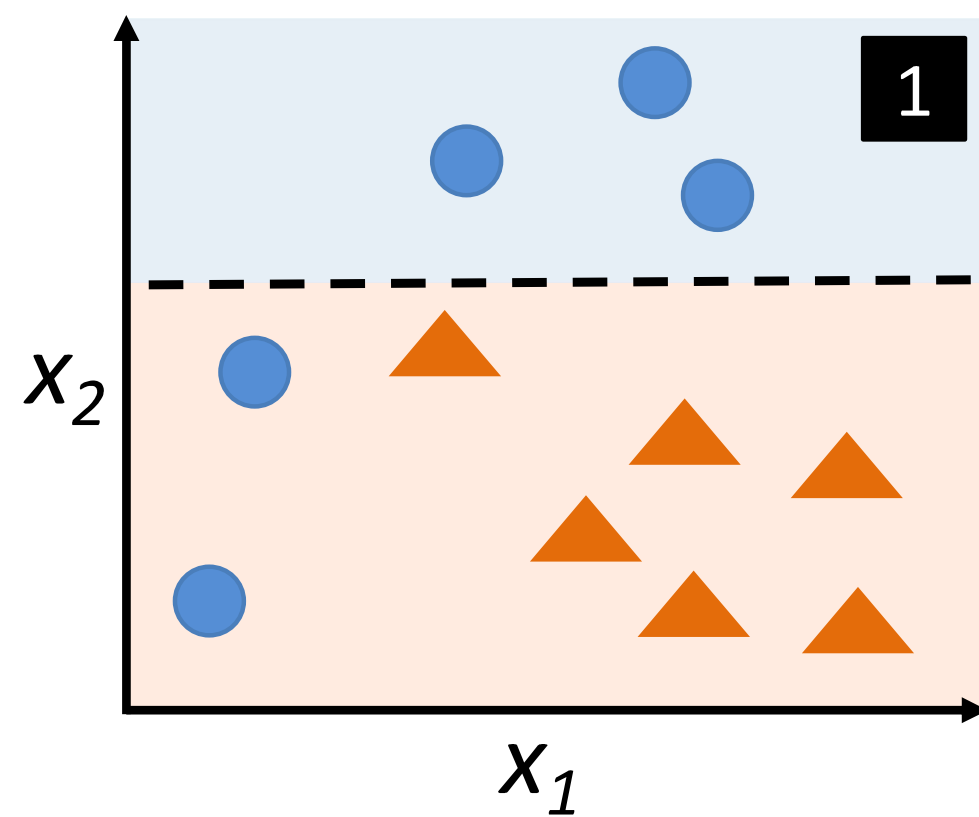
**Assumes binary classification problem**

# AdaBoost

---

**Algorithm 1** AdaBoost

---

1: Initialize $k$: the number of AdaBoost rounds
2: Initialize $\mathcal{D}$: the training dataset, $\mathcal{D} = \{\langle \mathbf{x}^{[1]}, y^{[1]} \rangle, ..., \mathbf{x}^{[n]}, y^{[n]} \rangle\}$
3: Initialize $w_1(i) = 1/n, \quad i = 1, ..., n, \ \mathbf{w}_1 \in \mathbb{R}^n$
4:
5: **for** r=1 to $k$ **do**
6:     For all $i : \mathbf{w}_r(i) := w_r(i)/\sum_i w_r(i)$    [normalize weights]
7:     $h_r := FitWeakLearner(\mathcal{D}, \mathbf{w}_r)$
8:     $\epsilon_r := \sum_i w_r(i) \mathbf{1}(h_r(i) \neq y_i)$    [compute error]
9:     if $\epsilon_r > 1/2$ then stop
10:    $\alpha_r := \frac{1}{2} \log[(1 - \epsilon_r)/\epsilon_r]$    [small if error is large and vice versa]

11:    $w_{r+1}(i) := w_r(i) \times \begin{cases} \mathrm{e}^{-\alpha_r} & \text{if } h_r(\mathbf{x}^{[i]}) = y^{[i]} \\ \mathrm{e}^{\alpha_r} & \text{if } h_r(\mathbf{x}^{[i]}) \neq y^{[i]} \end{cases}$

12: Predict: $h_{ens}(\mathbf{x}) = \arg\max_j \sum_r^k \alpha_r \mathbf{1}[h_r(\mathbf{x}) = j]$
13:

---

Estimator weight

Sample weight

# Gradient Boosting

# Gradient Boosting

Gradient boosting is somewhat similar to AdaBoost:
- trees are fit sequentially to improve error of previous trees
- boost weak learners to a strong learner

The way how the trees are fit sequentially differs in AdaBoost and Gradient Boosting, though ...

# Gradient Boosting -- Conceptual Overview

- **Step 1:** Construct a base tree (just the root node)

- **Step 2:** Build next tree based on errors of the previous tree

- **Step 3:** Combine tree from step 1 with trees from step 2. Go back to step 2.

**Algorithm 1** Gradient Boosting

1: Initialize $T$: the number of trees for gradient boosting rounds
2: Initialize $\mathcal{D}$: the training dataset, $\{\langle \mathbf{x}^{(i)}, y^{(i)} \rangle\}_{i=1}^n$
3: Choose $L\big(y^{(i)}, h(\mathbf{x}^{(i)})\big)$, a differentiable loss function
4: **Step 1**: Initialize model $h_0(\mathbf{x}) = \underset{\hat{y}}{\arg\min} \sum_{i=1}^n L\big(y^{(i)}, \hat{y}\big)$ [root node]

5: **Step 2**:
6: **for** t=1 to $T$ **do**

7:      **A.** Compute pseudo residual $r_{i,t} = -\left[ \dfrac{\partial L(y^{(i)}, h(\mathbf{x}^{(i)}))}{\partial h(\mathbf{x}^{(i)})} \right]_{h(\mathbf{x})=h_{t-1}(\mathbf{x})}$ , for $i=1$ to $n$

8:      **B.** Fit tree to $r_{i_t}$ values, and create terminal nodes $R_{j,t}$ for $j=1,...,J_t$.

9:      **C.**

10:      **for** j=1 to $J_t$ **do**

11:          $\hat{y}_{j,t} = \underset{\hat{y}}{\arg\min} \sum_{\mathbf{x}^{(i)} \in R_{i,j}} L\big(y^{(i)}, h_{t-1}(\mathbf{x}^{(i)}) + \hat{y}\big)$

12:      **D.** Update $h_t(\mathbf{x}) = h_{t-1}(\mathbf{x}) + \alpha \sum_{j=1}^{J_t} \hat{y}_{j,t}\, \mathbb{I}\big(\mathbf{x} \in R_{j,t}\big)$

13: **Step 3**: Return $h_t(\mathbf{x})$

# Gradient Boosting -- Conceptual Overview
## --> A Regression-based Example

In million US Dollars

| x1# Rooms | x2=City | x3=Age | y=Price |
|:---:|:---:|:---:|:---:|
| 5 | Boston | 30 | 1.5 |
| 10 | Madison | 20 | 0.5 |
| 6 | Lansing | 20 | 0.25 |
| 5 | Waunakee | 10 | 0.1 |

- **<u>Step 1:</u>** Construct a base tree (just the root node)

$$\hat{y}_1 = \frac{1}{n} \sum_{i=1}^{n} y^{(i)} = 0.5875$$

# Gradient Boosting -- Conceptual Overview --> A Regression-based Example

- **Step 2:** Build next tree based on errors of the previous tree

First, compute (pseudo) residuals: $r_1 = y_1 - \hat{y}_1$

In million US Dollars

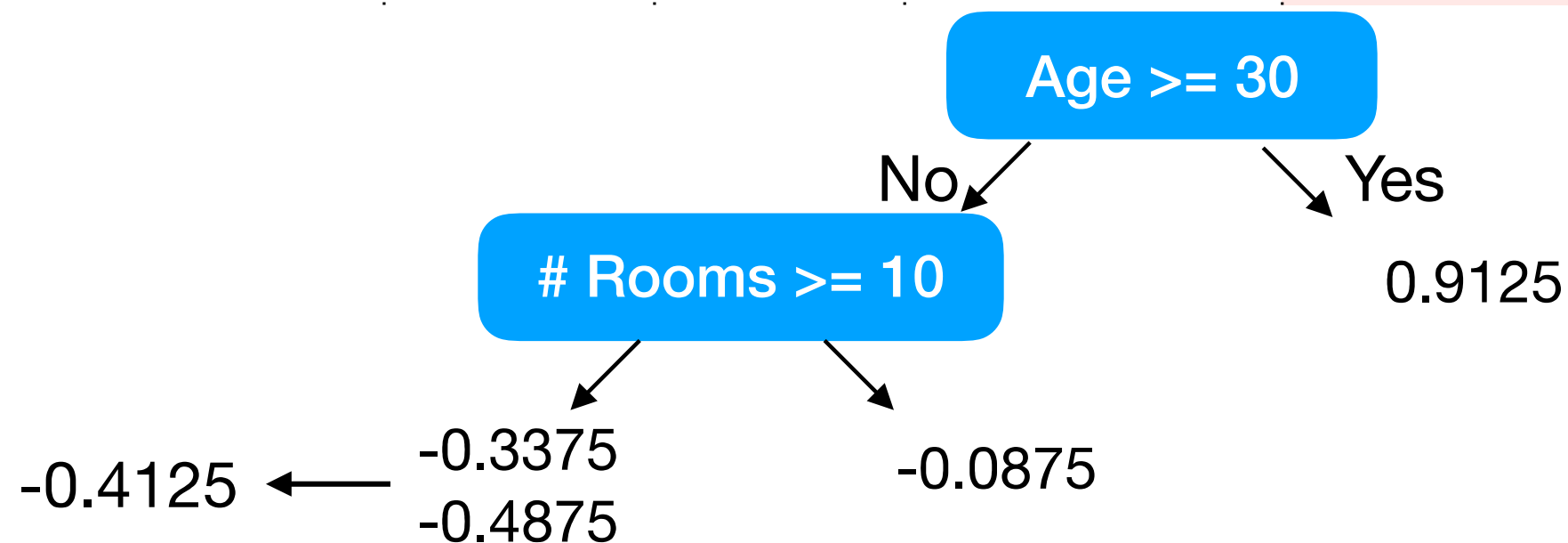| x1# | x2=City | x3=Age | y=Price | r1=Res |
|---|---|---|---|---|
| 5 | Boston | 30 | 1.5 | 1.5 - 0.5875 = 0.9125 |
| 10 | Madison | 20 | 0.5 | 0.5 - 0.5875 = -0.0875 |
| 6 | Lansing | 20 | 0.25 | 0.25 - 0.5875 = -0.3375 |
| 5 | Waunake | 10 | 0.1 | 0.1 - 0.5875 = -0.4875 |

# Gradient Boosting -- Conceptual Overview
## --> A Regression-based Example

- **<u>Step 2:</u>** Build next tree based on errors of the previous tree

Then, create a tree based on $x_1, \ldots, x_m$ to fit the residuals

| x1# | x2=City | x3=Age | y=Price | r1=Residual |
|-----|---------|--------|---------|-------------|
| 5 | Boston | 30 | 1.5 | 1.5 - 0.5875 = 0.9125 |
| 10 | Madison | 20 | 0.5 | 0.5 - 0.5875 = -0.0875 |
| 6 | Lansing | 20 | 0.25 | 0.25 - 0.5875 = -0.3375 |
| 5 | Waunake | 10 | 0.1 | 0.1 - 0.5875 = -0.4875 |

Age >= 30

No     Yes

0.9125

\# Rooms >= 10

-0.4125 ← -0.3375
-0.4875

-0.0875

# Gradient Boosting -- Conceptual Overview
## --> A Regression-based Example

- **<u>Step 3:</u>** Combine tree from step 1 with trees from step 2

| x1# | x2=City | x3=Age | y=Price | r=Res |
|---|---|---|---|---|
| 5 | Boston | 30 | 1.5 | 1.5 - 0.5875 = 0.9125 |
| 10 | Madison | 20 | 0.5 | 0.5 - 0.5875 = -0.0875 |
| 6 | Lansing | 20 | 0.25 | 0.25 - 0.5875 = -0.3375 |
| 5 | Waunake | 10 | 0.1 | 0.1 - 0.5875 = -0.4875 |

$$\hat{y}_1 = \frac{1}{n}\sum_{i=1}^{n} y^{(i)} = 0.5875 \quad +$$

Age >= 30

No    Yes

# Rooms >= 10     0.9125

-0.4125 ←   -0.3375
       -0.4875     -0.0875

# Gradient Boosting -- Conceptual Overview
## --> A Regression-based Example

- **<u>Step 3:</u>** Combine tree from step 1 with trees from step 2

| x1# | x2=City | x3=Age | y=Price | r=Res |
|---|---|---|---|---|
| 5 | Boston | 30 | 1.5 | 1.5 - 0.5875 = 0.9125 |
| 10 | Madison | 20 | 0.5 | 0.5 - 0.5875 = -0.0875 |
| 6 | Lansing | 20 | 0.25 | 0.25 - 0.5875 = -0.3375 |
| 5 | Waunakee | 10 | 0.1 | 0.1 - 0.5875 = -0.4875 |

E.g., predict Lansing →

**Age >= 30**

No ↙        Yes ↘

**# Rooms >= 10**        0.9125

↙        ↘

-0.3375
-0.4125 ←        -0.0875
-0.4875

$$\hat{y}_1 = \frac{1}{n}\sum_{i=1}^{n} y^{(i)} = 0.5875 \quad \textbf{+}$$

E.g., predict Lansing

$$0.5875 + \alpha \times (-0.4125)$$

where $\alpha$ learning rate between 0 and 1 (if $\alpha = 1$, low bias but high variance)

# Gradient Boosting -- Algorithm Overview

**Step 0:**   Input data  $\{\langle \mathbf{x}^{(i)}, y^{(i)} \rangle\}_{i=1}^{n}$

Differentiable Loss function  $L\left(y^{(i)}, h(\mathbf{x}^{(i)})\right)$

**Step 1:**   Initialize model  $h_0(\mathbf{x}) = \underset{\hat{y}}{\operatorname{argmin}} \sum_{i=1}^{n} L\left(y^{(i)}, \hat{y}\right)$

**Step 2:**   for  $t = 1$   to  $T$

**A.** Compute pseudo residual  $r_{i,t} = -\left[\dfrac{\partial L(y^{(i)}, h(\mathbf{x}^{(i)}))}{\partial h(\mathbf{x}^{(i)})}\right]_{h(\mathbf{x}) = h_{t-1}(\mathbf{x})}$

for  $i = 1$  to  $n$

**B.** Fit tree to $r_{i,t}$ values, and create terminal nodes $R_{j,t}$ for $j = 1,...,J_t$

■ ■ ■

# Gradient Boosting -- Algorithm Overview

**Step 2:**    for $t = 1$   to   $T$

**A.** Compute pseudo residual $r_{i,t} = - \left[ \dfrac{\partial L(y^{(i)}, h(\mathbf{x}^{(i)}))}{\partial h(\mathbf{x}^{(i)})} \right]_{h(\mathbf{x})=h_{t-1}(\mathbf{x})}$

for $i = 1$ to $n$

**B.** Fit tree to $r_{i,t}$ values, and create terminal nodes $R_{j,t}$ for $j = 1,...,J_t$

**C.** for $j = 1,...,J_t$, compute

$$\hat{y}_{j,t} = \operatorname*{argmin}_{\hat{y}} \sum_{\mathbf{x}^{(i)} \in R_{i,j}} L\left(y^{(i)}, h_{t-1}(\mathbf{x}^{(i)}) + \hat{y}\right)$$

**D.** Update $h_t(\mathbf{x}) = h_{t-1}(\mathbf{x}) + \alpha \sum_{j=1}^{J_t} \hat{y}_{j,t} \, \mathbb{I}\left(\mathbf{x} \in R_{j,t}\right)$

**Step 3:**   Return $h_t(\mathbf{x})$

# Gradient Boosting -- Algorithm Overview Discussion

**Step 0:**  Input data  $\{\langle \mathbf{x}^{(i)}, y^{(i)} \rangle\}_{i=1}^n$

Differentiable Loss function  $L\big(y^{(i)}, h(\mathbf{x}^{(i)})\big)$

E.g., Sum-squared error in regression

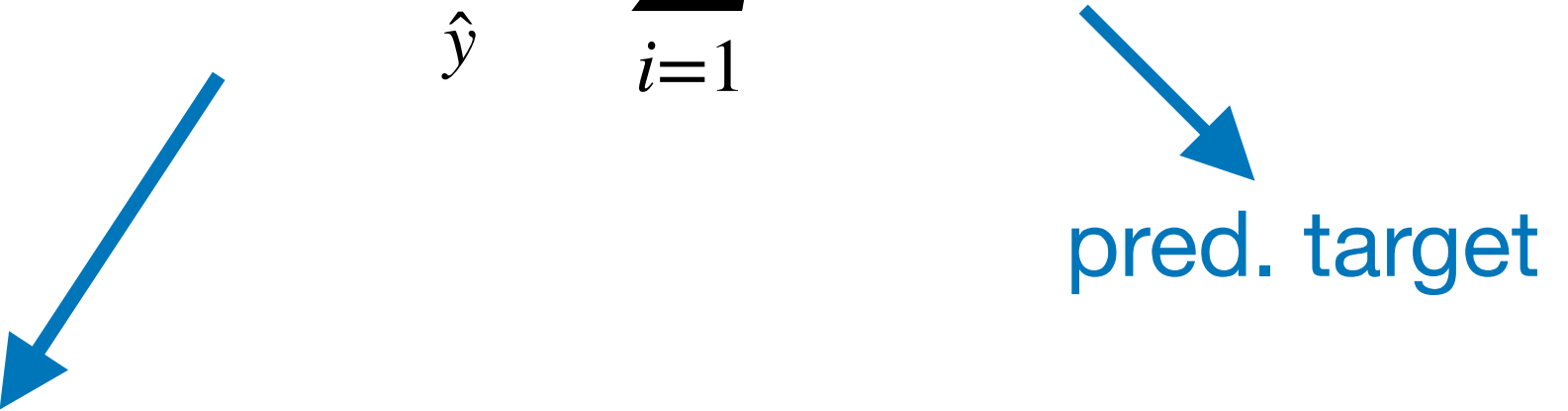$$SSE' = \frac{1}{2}\big(y^{(i)} - h(\mathbf{x}^{(i)})\big)^2$$

$$\frac{\partial}{\partial h(\mathbf{x}^{(i)})} \frac{1}{2}\big(y^{(i)} - h(\mathbf{x}^{(i)})\big)^2 \quad \text{[chain rule]}$$

$$= 2 \times \frac{1}{2}\big(y^{(i)} - h(\mathbf{x}^{(i)})\big) \times (0 - 1) = -\big(y^{(i)} - h(\mathbf{x}^{(i)})\big)$$

[neg. residual]

# Gradient Boosting -- Algorithm Overview Discussion

**Step 1:** Initialize model $h_0(\mathbf{x}) = \underset{\hat{y}}{\text{argmin}} \sum_{i=1}^{n} L\left(y^{(i)}, \hat{y}\right)$

pred. target

turns out to be the average (in regression)

$$\frac{1}{n} \sum_{i=1}^{n} y^{(i)}$$

# Gradient Boosting -- Algorithm Overview Discussion

Loop to make *T* trees (e.g., *T=100*)

**Step 2:** for $t = 1$ to $T$

**A.** Compute pseudo residual $r_{i,t} = - \left[ \dfrac{\partial L(y^{(i)}, h(\mathbf{x}^{(i)}))}{\partial h(\mathbf{x}^{(i)})} \right]_{h(\mathbf{x}) = h_{t-1}(\mathbf{x})}$

for $i = 1$ to $n$

pseudo residual of the *t*-th tree and *i*-th example

Derivative of the loss function

# Gradient Boosting -- Algorithm Overview **Discussion**

Loop to make *T* trees (e.g., *T=100*)
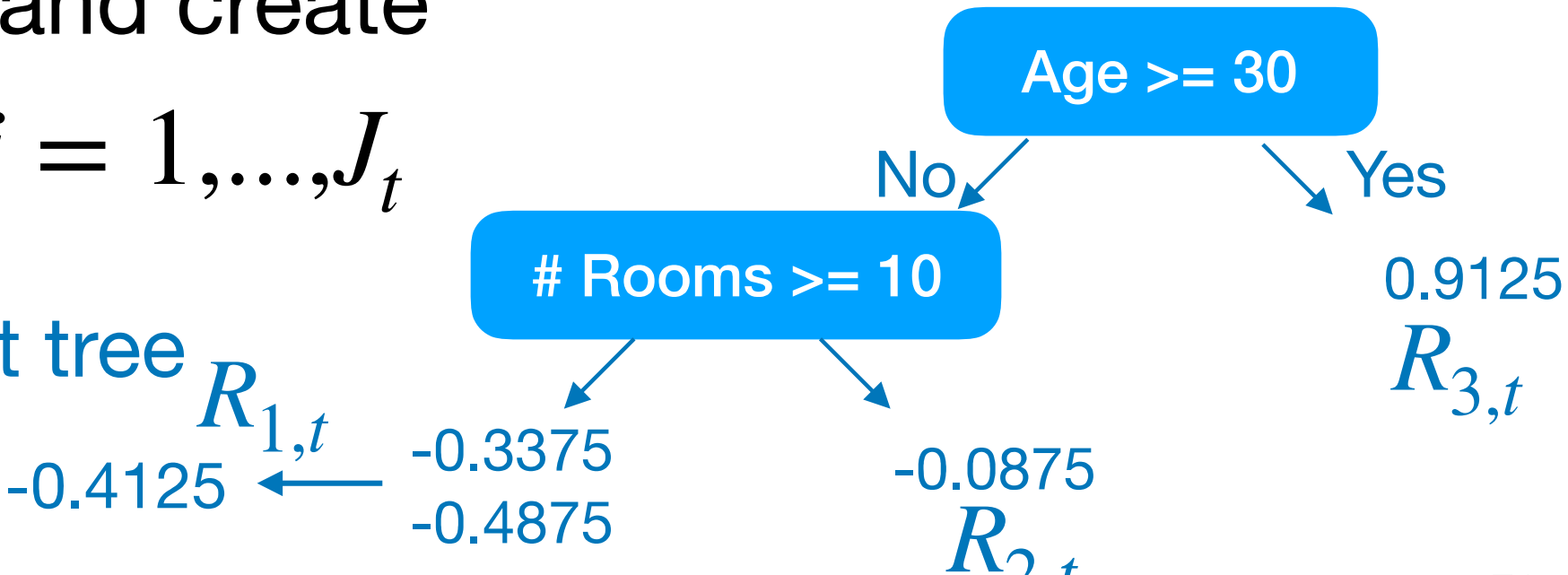
**Step 2:** for $t = 1$ to $T$

**A.** Compute pseudo residual $r_{i,t} = - \left[ \dfrac{\partial L(y^{(i)}, h(\mathbf{x}^{(i)}))}{\partial h(\mathbf{x}^{(i)})} \right]_{h(\mathbf{x})=h_{t-1}(\mathbf{x})}$

for $i = 1$ to $n$

pseudo residual of the *t*-th tree and *i*-th example

Derivative of the loss function

**B.** Fit tree to $r_{i,t}$ values, and create terminal nodes $R_{j,t}$ for $j = 1,...,J_t$

Use features in dataset to fit tree

Age >= 30

No     Yes

# Rooms >= 10

$R_{1,t}$

0.9125
$R_{3,t}$

-0.4125 ←   -0.3375   -0.4875     -0.0875

$R_{2,t}$

# Gradient Boosting -- Algorithm Overview **Discussion**

**Step 2:** for $t = 1$ to $T$

**A.** Compute pseudo residual $r_{i,t} = -\left[ \dfrac{\partial L(y^{(i)}, h(\mathbf{x}^{(i)}))}{\partial h(\mathbf{x}^{(i)})} \right]_{h(\mathbf{x})=h_{t-1}(\mathbf{x})}$

for $i = 1$ to $n$

**B.** Fit tree to $r_{i,t}$ values, and create terminal nodes $R_{j,t}$ for $j = 1,...,J_t$

**C.** for $j = 1,...,J_t$, compute

$$\hat{y}_{j,t} = \underset{\hat{y}}{\arg\min} \sum_{\mathbf{x}^{(i)} \in R_{i,j}} L\left(y^{(i)}, h_{t-1}(\mathbf{x}^{(i)}) + \hat{y}\right)$$

Compute the residual for each leaf node

Only consider examples at that leaf node

Like step 1 but add previous prediction

# Gradient Boosting -- Algorithm Overview Discussion

**Step 2:**    for $t = 1$   to   $T$

**A.** Compute pseudo residual $r_{i,t} = -\left[ \dfrac{\partial L(y^{(i)}, h(\mathbf{x}^{(i)}))}{\partial h(\mathbf{x}^{(i)})} \right]_{h(\mathbf{x})=h_{t-1}(\mathbf{x})}$

                   for $i = 1$ to $n$

**B.** Fit tree to $r_{i,t}$ values, and create terminal nodes $R_{j,t}$ for $j = 1,\ldots,J_t$

**C.** for $j = 1,\ldots,J_t$, compute

$$\hat{y}_{j,t} = \underset{\hat{y}}{\arg\min} \sum_{\mathbf{x}^{(i)} \in R_{i,j}} L\big(y^{(i)}, h_{t-1}(\mathbf{x}^{(i)}) + \hat{y}\big)$$

**D.** Update $h_t(\mathbf{x}) = h_{t-1}(\mathbf{x}) + \alpha \displaystyle\sum_{j=1}^{J_t} \hat{y}_{j,t}\, \mathbb{I}\big(\mathbf{x} \in R_{j,t}\big)$

learning rate between 0 and 1 (usually 0.1)

Summation just in case examples end up in multiple nodes

# Gradient Boosting -- Algorithm Overview Discussion

For prediction, combine all $T$ trees, e.g.,

$$h_0(\mathbf{x}) = \underset{\hat{y}}{\text{argmin}} \sum_{i=1}^{n} L\big(y^{(i)}, \hat{y}\big)$$

$$+\alpha\, \hat{y}_{j,t=1} = \underset{\hat{y}}{\text{argmin}} \sum_{\mathbf{x}^{(i)} \in R_{i,j}} L\big(y^{(i)}, h_{(t=1)-1}(\mathbf{x}^{(i)}) + \hat{y}\big)$$

$$\ldots$$

$$+\alpha\, \hat{y}_{j,T} = \underset{\hat{y}}{\text{argmin}} \sum_{\mathbf{x}^{(i)} \in R_{i,j}} L\big(y^{(i)}, h_{T-1}(\mathbf{x}^{(i)}) + \hat{y}\big)$$

For prediction, combine all $T$ trees, e.g.,

$$h_0(\mathbf{x}) = \underset{\hat{y}}{\text{argmin}} \sum_{i=1}^{n} L\left(y^{(i)}, \hat{y}\right)$$

$+\alpha\ \hat{y}_{j,t=1}$

The idea is that we decrease the pseudo residuals by a small amount at each step

$\dots$

$+\alpha\ \hat{y}_{j,T}$

# XGBoost

Summary and Main Points:

- scalable implementation of gradient boosting

- Improvements include: regularized loss, sparsity-aware algorithm, weighted quantile sketch for approximate tree learning, caching of access patterns, data compression, sharding

- Decision trees based on CART

- Regularization term for penalizing model (tree) complexity

- Uses second order approximation for optimizing the objective

- Options for column-based and row-based subsampling

- Single-machine version of XGBoost supports the exact greedy algorithm

Chen, T., & Guestrin, C. (2016, August). Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining* (pp. 785-794). ACM.

# Stacking

# Stacking Algorithm

Wolpert, David H. "Stacked generalization." Neural networks 5.2 (1992): 241-259.
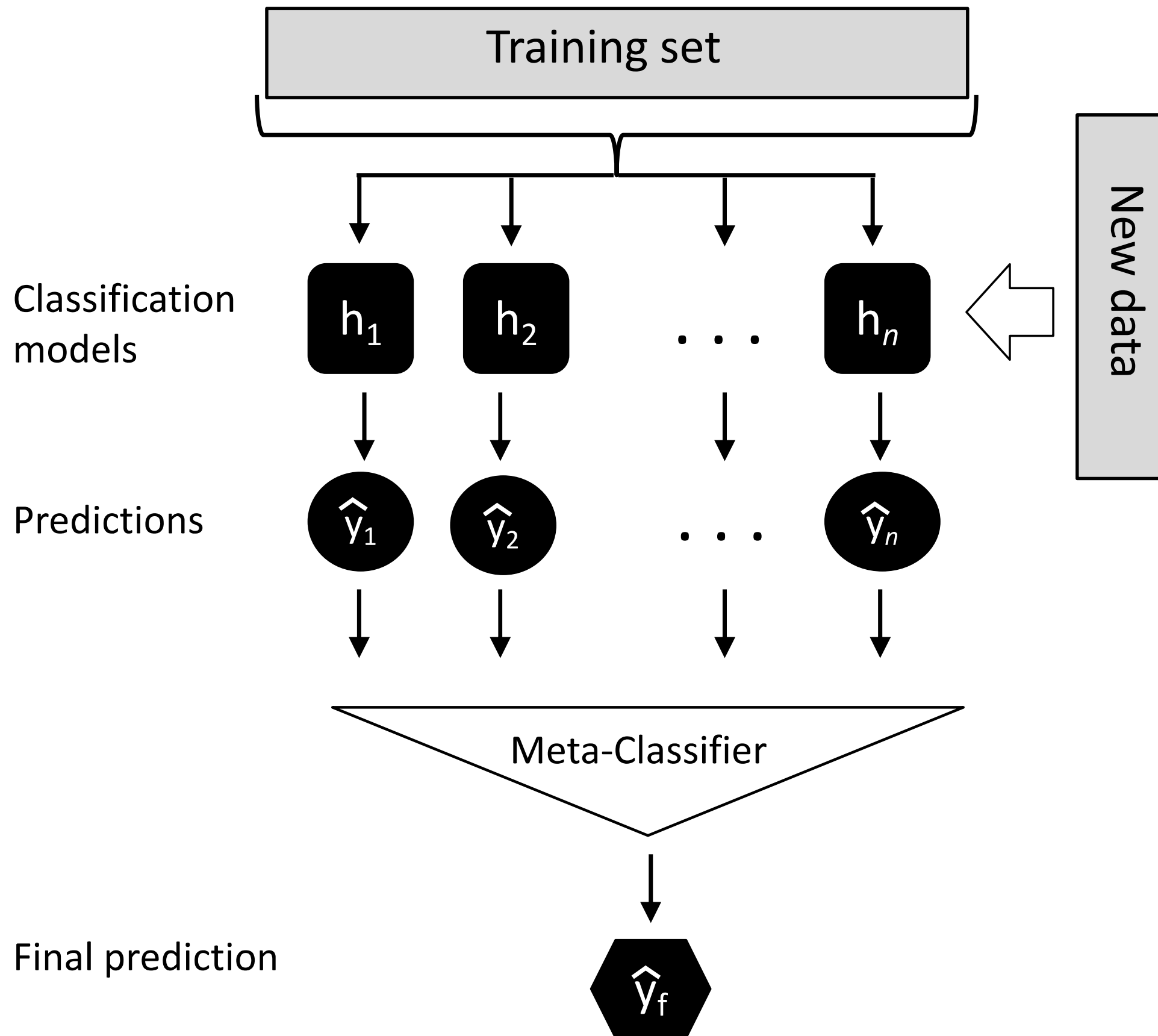
**Algorithm 19.7 Stacking**

**Input:** Training data $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^{m}$ $(\mathbf{x}_i \in \mathbb{R}^n, y_i \in \mathcal{Y})$
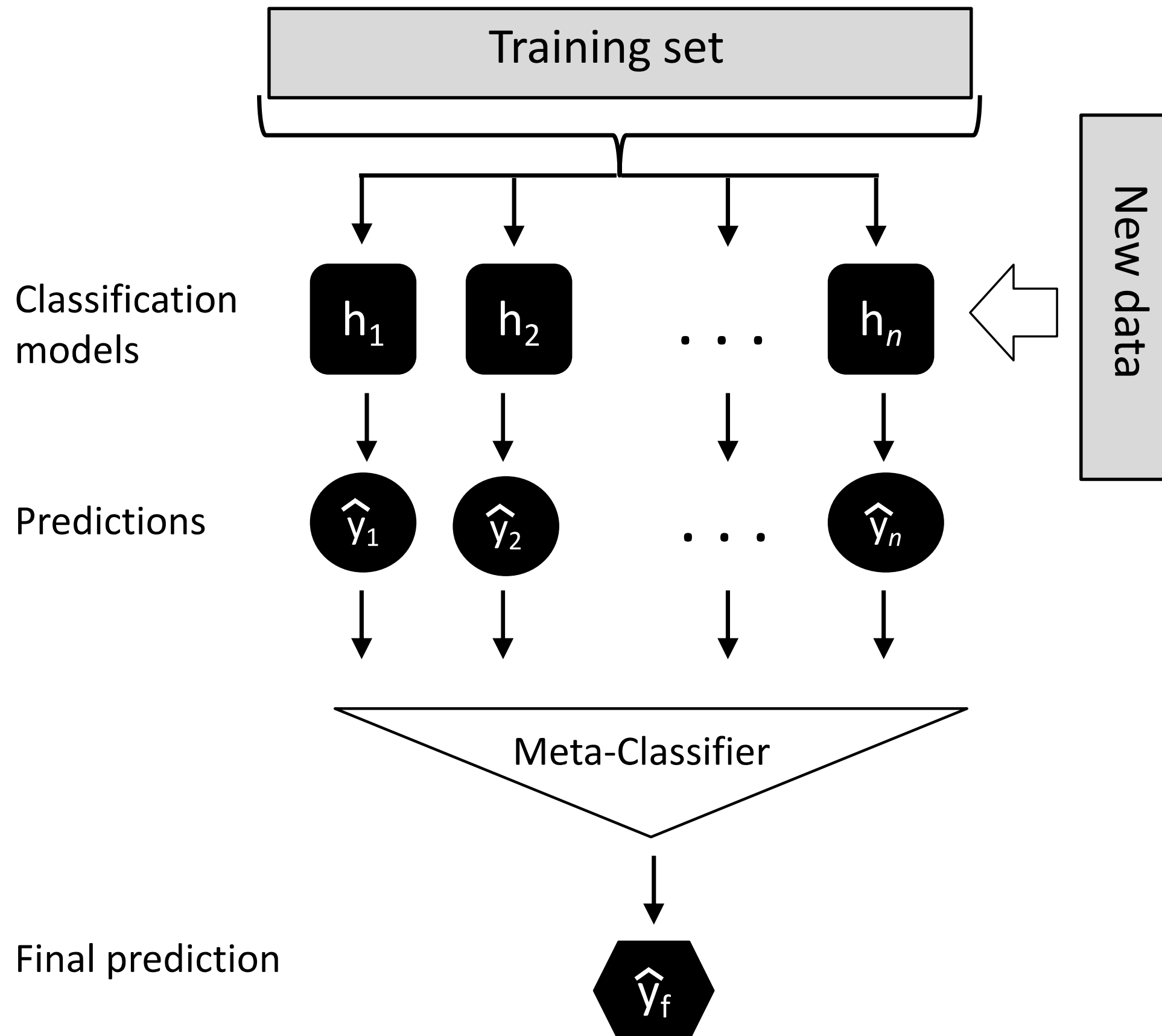**Output:** An ensemble classifier $H$

1: Step 1: Learn first-level classifiers
2: **for** $t \leftarrow 1$ to $T$ **do**
3:     Learn a base classifier $h_t$ based on $\mathcal{D}$
4: **end for**
5: Step 2: Construct new data sets from $\mathcal{D}$
6: **for** $i \leftarrow 1$ to $m$ **do**
7:     Construct a new data set that contains $\{\mathbf{x}_i', y_i\}$, where $\mathbf{x}_i' = \{h_1(\mathbf{x}_i), h_2(\mathbf{x}_i), \dots, h_T(\mathbf{x}_i)\}$
8: **end for**
9: Step 3: Learn a second-level classifier
10: Learn a new classifier $h'$ based on the newly constructed data set
11: **return** $H(\mathbf{x}) = h'(h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_T(\mathbf{x}))$

Tang, J., S. Alelyani, and H. Liu. "Data Classification: Algorithms and Applications." Data Mining and Knowledge Discovery Series, CRC Press (2015): pp. 498-500.
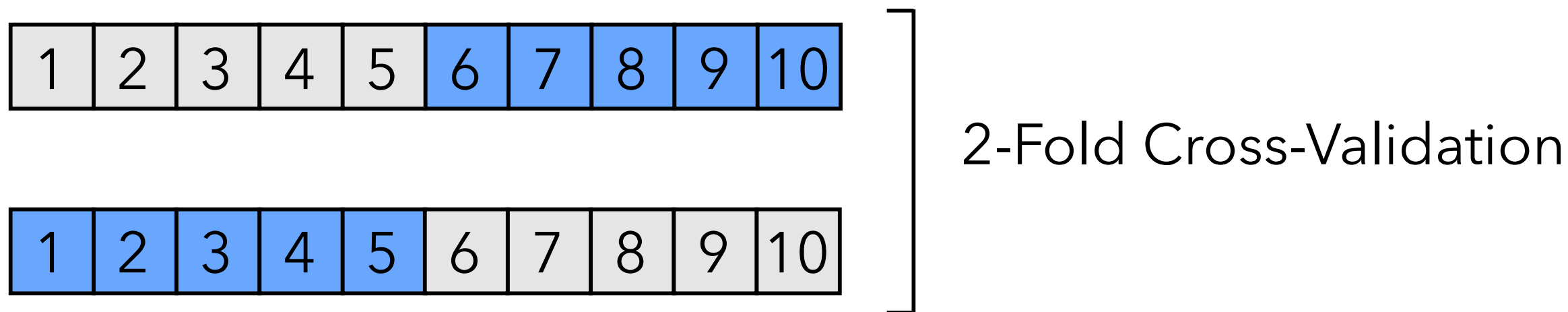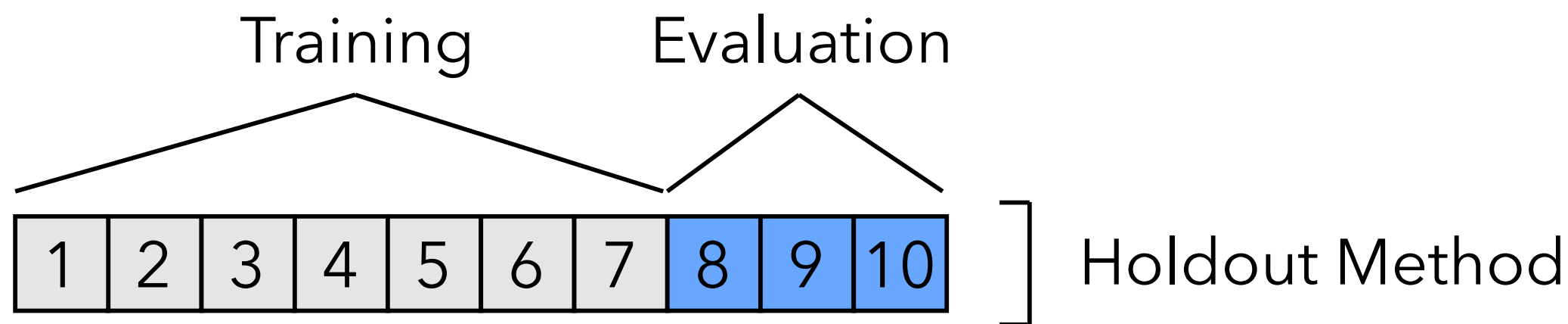
# Stacking Algorithm

# Cross-Validation

# *k*-fold Cross-Validation



Validation Fold

Training Fold

K Iterations (K-Folds)

1st — Performance$_1$

2nd — Performance$_2$

3rd — Performance$_3$

4th — Performance$_4$

5th — Performance$_5$

$$\text{Performance} = \frac{1}{5}\sum_{i=1}^{5}\text{Performance}_i$$

A

Validation Fold

Training Fold

1st    Performance$_1$

2nd    Performance$_2$

3rd    Performance$_3$

4th    Performance$_4$

5th    Performance$_5$

K Iterations (K-Folds)

$$\text{Performance} = \frac{1}{5} \sum_{i=1}^{5} \text{Performance}_i$$

B

Training Fold Data

Training Fold Labels

Hyperparameter Values

**Learning Algorithm**
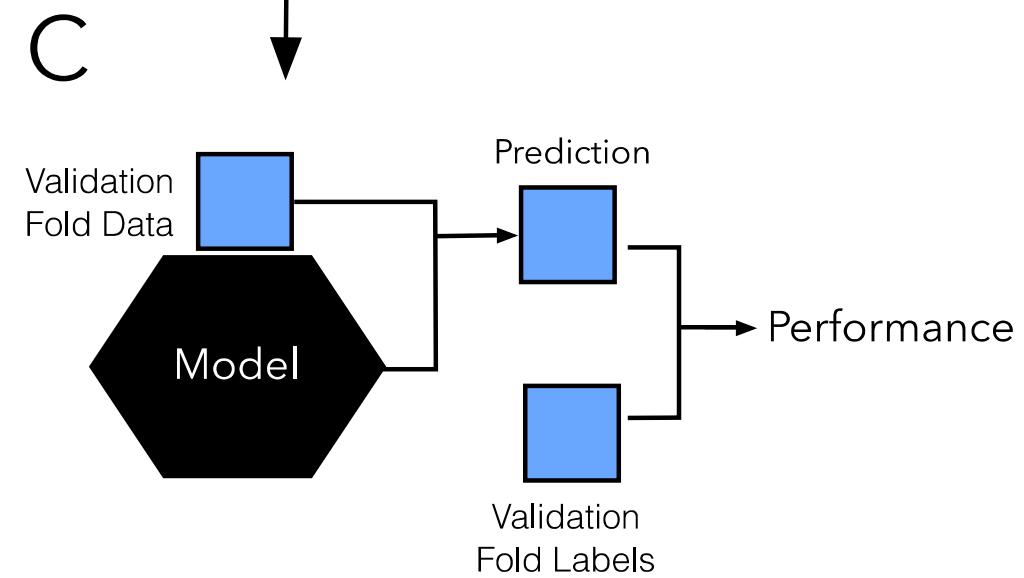
Model

C

Validation Fold Data

Prediction

Model
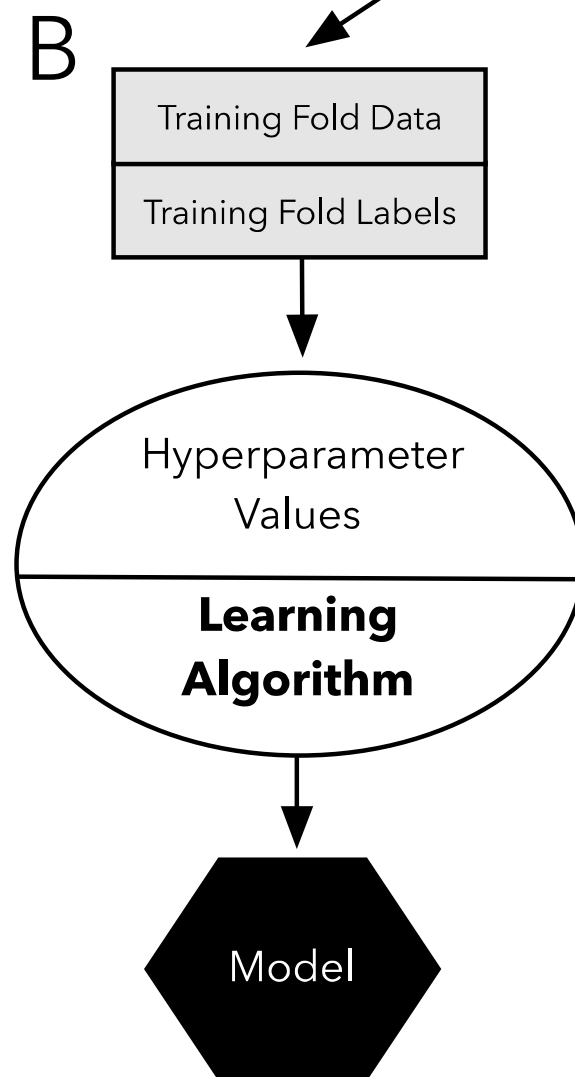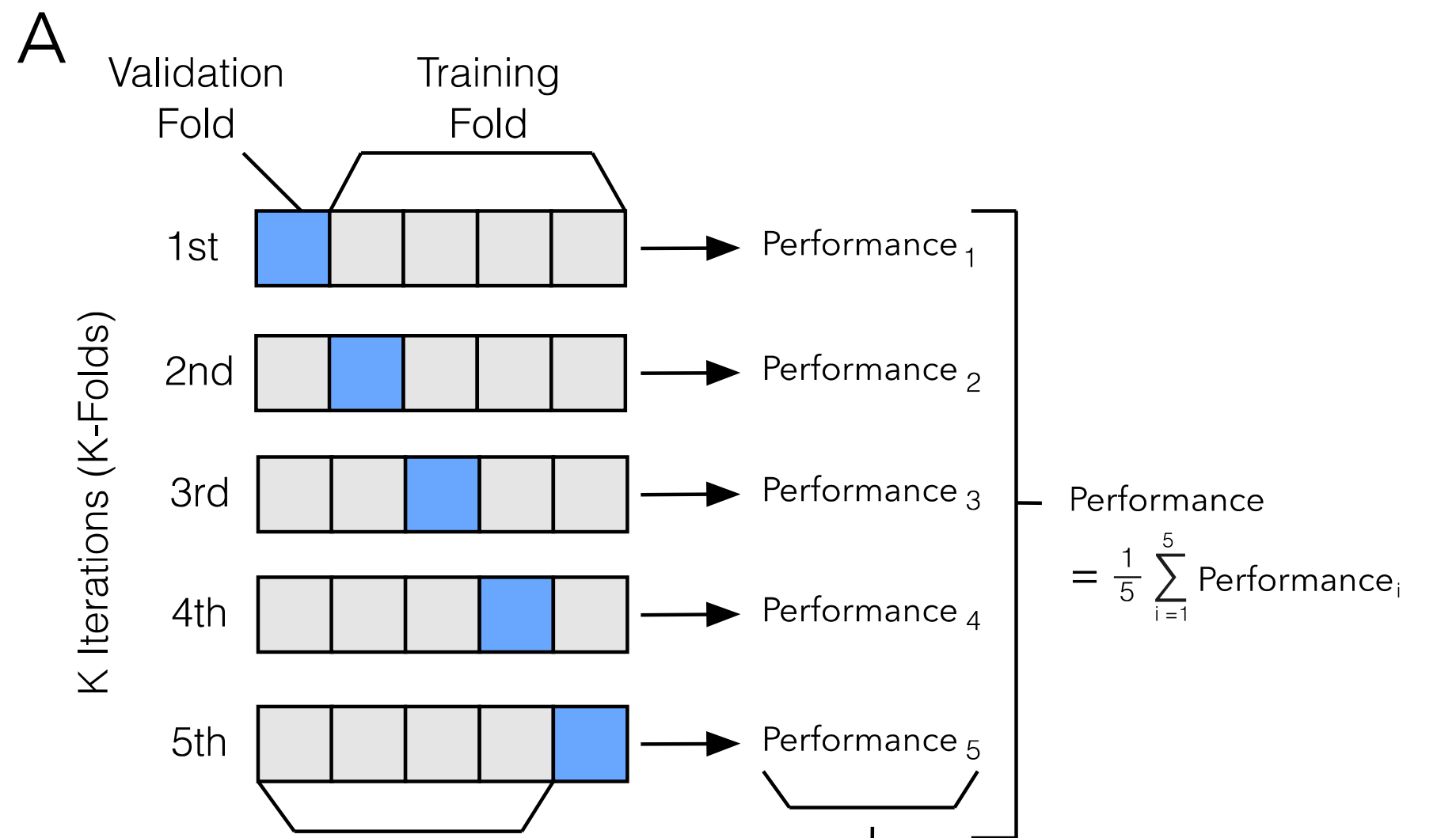
Validation Fold Labels

Performance

# Stacking Algorithm with Cross-Validation

Wolpert, David H. "Stacked generalization." Neural networks 5.2 (1992): 241-259.

---

**Algorithm 19.8 Stacking with $K$-fold Cross Validation**

---

**Input:** Training data $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^m$ ($\mathbf{x}_i \in \mathbb{R}^n$, $y_i \in \mathcal{Y}$)

**Output:** An ensemble classifier $H$

---

1: Step 1: Adopt cross validation approach in preparing a training set for second-level classifier
2: Randomly split $\mathcal{D}$ into $K$ equal-size subsets: $\mathcal{D} = \{\mathcal{D}_1, \mathcal{D}_2, \ldots, \mathcal{D}_K\}$
3: **for** $k \leftarrow 1$ to $K$ **do**
4:     Step 1.1: Learn first-level classifiers
5:     **for** $t \leftarrow 1$ to $T$ **do**
6:         Learn a classifier $h_{kt}$ from $\mathcal{D} \setminus \mathcal{D}_k$
7:     **end for**
8:     Step 1.2: Construct a training set for second-level classifier
9:     **for** $\mathbf{x}_i \in \mathcal{D}_k$ **do**
10:         Get a record $\{\mathbf{x}_i', y_i\}$, where $\mathbf{x}_i' = \{h_{k1}(\mathbf{x}_i), h_{k2}(\mathbf{x}_i), \ldots, h_{kT}(\mathbf{x}_i)\}$
11:     **end for**
12: **end for**
13: Step 2: Learn a second-level classifier
14: Learn a new classifier $h'$ from the collection of $\{\mathbf{x}_i', y_i\}$
15: Step 3: Re-learn first-level classifiers
16: **for** $t \leftarrow 1$ to $T$ **do**
17:     Learn a classifier $h_t$ based on $\mathcal{D}$
18: **end for**
19: **return** $H(\mathbf{x}) = h'(h_1(\mathbf{x}), h_2(\mathbf{x}), \ldots, h_T(\mathbf{x}))$

---

Tang, J., S. Alelyani, and H. Liu. "Data Classification: Algorithms and Applications." Data Mining and Knowledge Discovery Series, CRC Press (2015): pp. 498-500.

# Stacking Algorithm with Cross-Validation

Training set

Training folds   Validation fold

$h_1$   $h_2$   ...   $h_n$   Base Classifiers

$\hat{y}_1$   $\hat{y}_2$   ...   $\hat{y}_n$   Level-1 predictions in $k$-th iteration

Repeat $k$ times

All level-1 predictions

Train

Meta-Classifier

$\hat{y}_f$   Final prediction