# Lecture 07

# Ensemble Methods

STAT 479: Machine Learning, Fall 2018

Sebastian Raschka

http://stat.wisc.edu/~sraschka/teaching/stat479-fs2018/

# Announcements First

**[DataSci] Machine Learning Series**
Scheduled: Thursday Oct 10, 6:30 PM
Location: Genetics-Biotech Center 1441

UW DATA SCIENCE CLUB

Machine Learning Series :

## Soft-Biometric Attributes Prediction from Face Images with PyTorch

Sebastian Raschka

Thursday Oct 10, 6:30pm
eGenetics-Biotech 1441

Soft-biometric characteristics include a person's age, gender, race, and health status. As many Deep Learning-centric applications are developed in recent years, the automatic extraction of soft biometric attributes can happen without the user's agreement, thereby raising several privacy concerns. This talk will introduce how to extract soft-biometric attributes from facial images, as well as how to conceal soft-biometric information for enhancing privacy.

Don't worry if you do not have programming experience with Python! Dr. Rashka will also give a tutorial introducing PyTorch and how we can use it to train a simple gender classifier and ordinal regression model for estimating the apparent age from face images.

Best,

Lareina Liu
UW Data Science Club
.Data

We are here
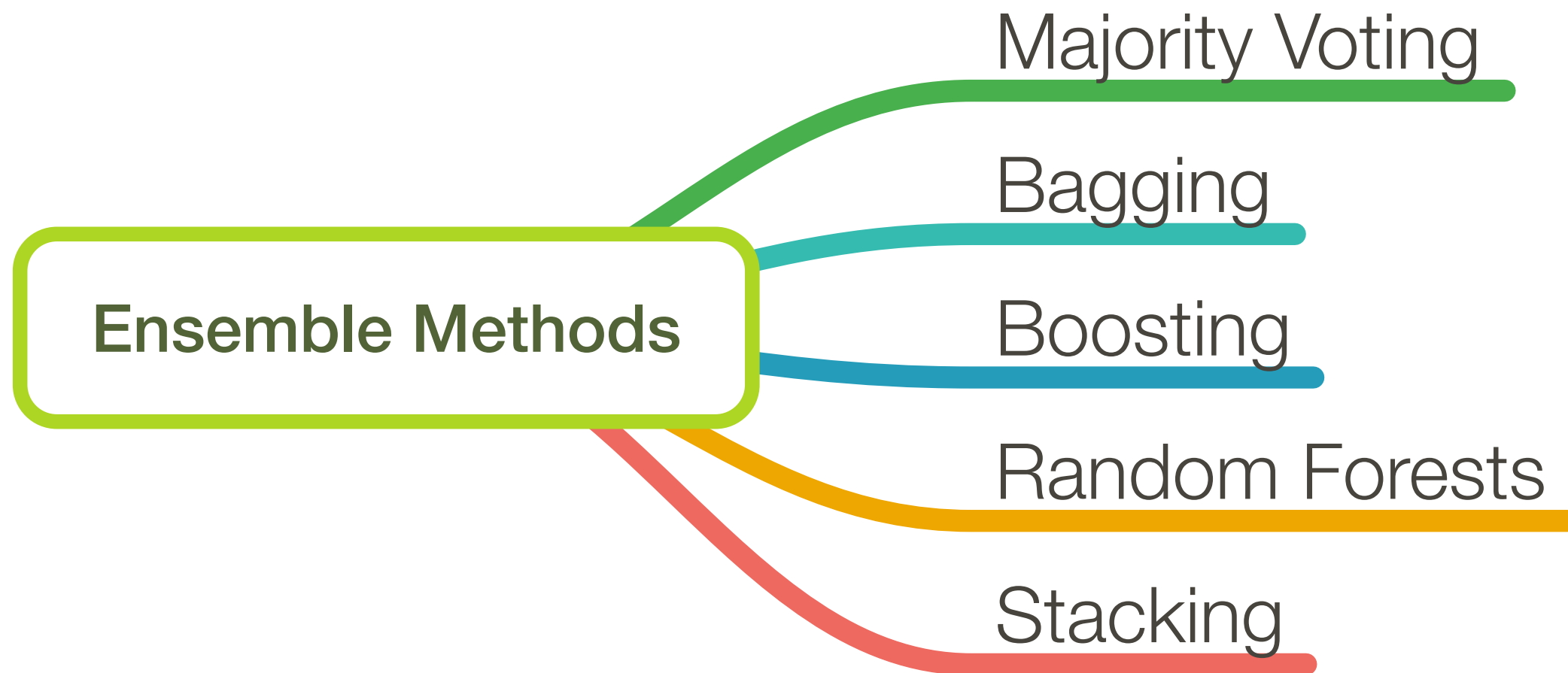
Genetics-Biotech Center 1441

# Example Exam Question

(6 points) Does the (computational) time complexity of a $k$-Nearest Neighbor classifier grow linearly, quadratically, or exponentially with the number of samples in the training dataset? Explain your answer in 1-2 sentences.

# Example Exam Question

(6 points) Can you represent the following boolean function with a decision tree? If you answer "no," explain why in 1-2 sentences. Otherwise, draw a decision tree that separates the data records perfectly.
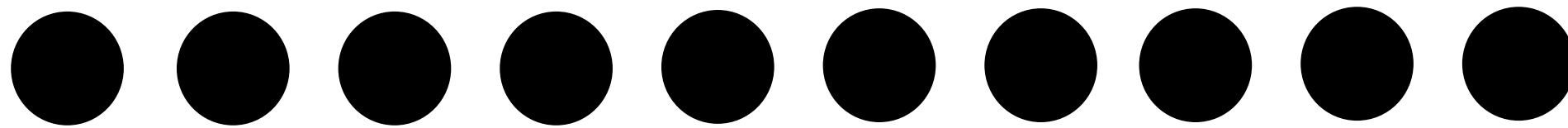
| $x_1$ | $x_2$ | $f(x_1, x_2)$ |
|-------|-------|---------------|
| 1 | 1 | 0 |
| 0 | 0 | 0 |
| 1 | 0 | 1 |
| 0 | 1 | 0 |

# Lecture Overview



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

# Majority Voting

Unanimity

Majority

Plurality

# Majority Vote Classifier



Training set

Classification models

$h_1$  $h_2$  $\ldots$  $h_n$

New data

Predictions

$\hat{y}_1$  $\hat{y}_2$  $\ldots$  $\hat{y}_n$

Voting

Final prediction

$\hat{y}_f$  $\hat{y}_f = mode\{h_1(\mathbf{x}), h_2(\mathbf{x}), \ldots h_n(\mathbf{x})\}$

where $h_i(\mathbf{x}) = \hat{y}_i$

# Why Majority Vote?

- assume $n$ independent classifiers with a base error rate $\epsilon$

- here, independent means that the errors are uncorrelated

- assume a binary classification task

- assume the error rate is better than random guessing (i.e., lower than 0.5 for binary classification)

$$\forall \epsilon_i \in \{\epsilon_1, \epsilon_2, \ldots, \epsilon_n\}, \epsilon_i < 0.5$$

# Why Majority Vote?

- assume $n$ independent classifiers with a base error rate $\epsilon$
- here, independent means that the errors are uncorrelated
- assume a binary classification task
- assume the error rate is better than random guessing (i.e., lower than 0.5 for binary classification)

$$\forall \epsilon_i \in \{\epsilon_1, \epsilon_2, \ldots, \epsilon_n\}, \epsilon_i < 0.5$$
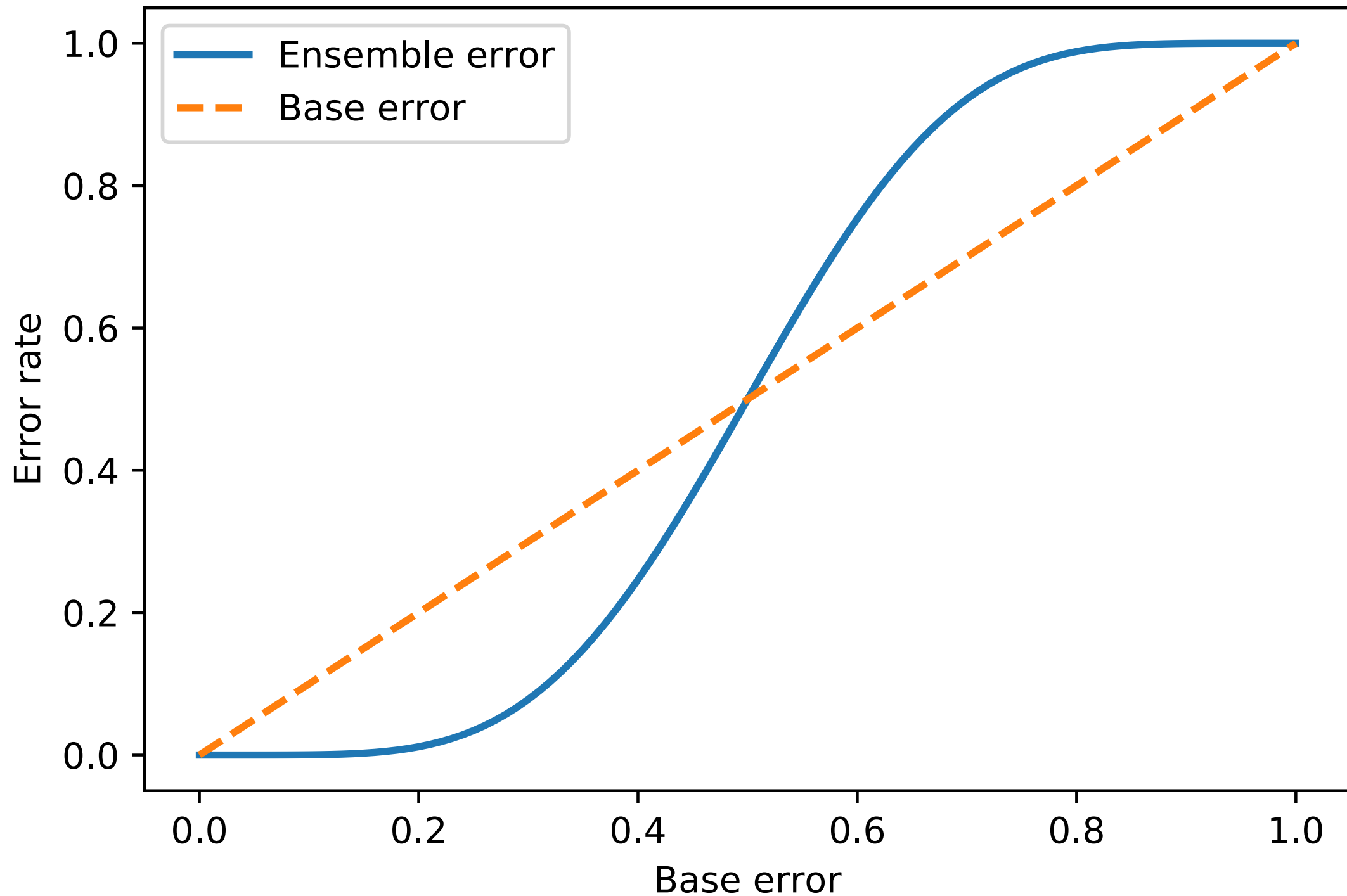
The probability that we make a wrong prediction via the ensemble if $k$ classifiers predict the same class label

$$P(k) = \binom{n}{k} \epsilon^k (1 - \epsilon)^{n-k} \qquad k > \lceil n/2 \rceil$$

(Probability mass func. of a binomial distr.)

# Why Majority Vote?

The probability that we make a wrong prediction via the ensemble if $k$ classifiers predict the same class label

$$P(k) = \binom{n}{k} \epsilon^k (1 - \epsilon)^{n-k} \qquad k > \lceil n/2 \rceil$$

Ensemble error:

$$\epsilon_{ens} = \sum_{k}^{n} \binom{n}{k} \epsilon^k (1 - \epsilon)^{n-k}$$

$$\epsilon_{ens} = \sum_{k=6}^{11} \binom{11}{k} 0.25^k (1 - 0.25)^{11-k} = 0.034$$

$$\epsilon_{ens} = \sum_{k}^{n} \binom{n}{k} \epsilon^k (1-\epsilon)^{n-k}$$

# "Soft" Voting

$$\hat{y} = \arg\max_{j} \sum_{i=1}^{n} w_i p_{i,j}$$

$p_{i,j}$ : predicted class membership probability of the $i$th classifier for class label $j$

$w_j$ : optional weighting parameter, default $w_i = 1/n, \forall w_i \in \{w_1, \ldots, w_n\}$

# "Soft" Voting

$$\hat{y} = \arg \max_{j} \sum_{i=1}^{n} w_i p_{i,j}$$

$p_{i,j}:$ predicted class membership probability of the $i$th classifier for class label $j$

$w_j:$ optional weighting parameter, default $w_i = 1/n, \forall w_i \in \{w_1, \ldots, w_n\}$

# "Soft" Voting

$$\hat{y} = \arg\max_{j} \sum_{i=1}^{n} w_i p_{i,j}$$

Binary classification example

$$j \in \{0,1\} \qquad h_i(i \in \{1,2,3\})$$

$$h_1(\mathbf{x}) \to [0.9, 0.1]$$
$$h_2(\mathbf{x}) \to [0.8, 0.2]$$
$$h_3(\mathbf{x}) \to [0.4, 0.6]$$

# "Soft" Voting

$$\hat{y} = \arg\max_j \sum_{i=1}^{n} w_i p_{i,j}$$

Binary classification example

$$j \in \{0,1\} \qquad h_i(i \in \{1,2,3\})$$

$$h_1(\mathbf{x}) \rightarrow [0.9, 0.1]$$
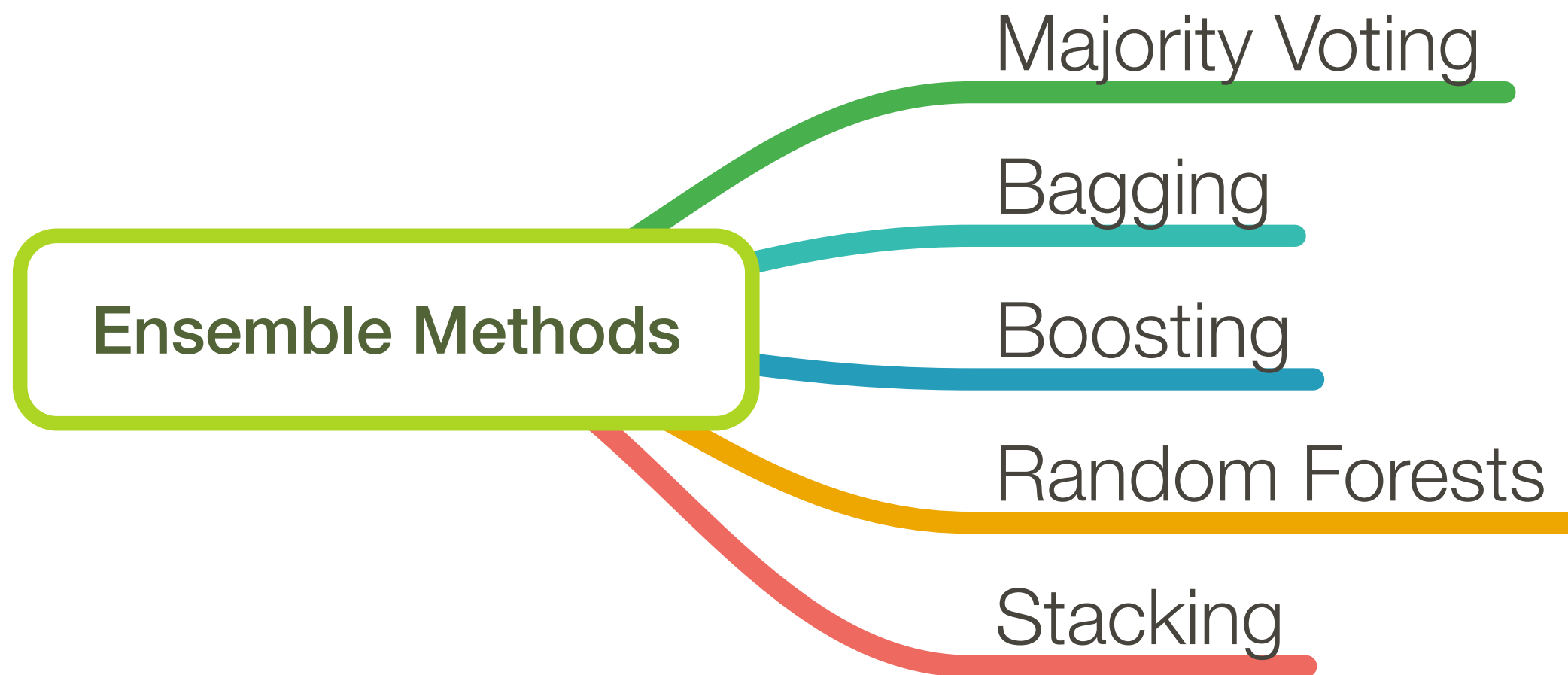
$$h_2(\mathbf{x}) \rightarrow [0.8, 0.2]$$

$$h_3(\mathbf{x}) \rightarrow [0.4, 0.6]$$

$$p(j = 0 \,|\, \mathbf{x}) = 0.2 \cdot 0.9 + 0.2 \cdot 0.8 + 0.6 \cdot 0.4 = 0.58$$

$$p(j = 1 \,|\, \mathbf{x}) = 0.2 \cdot 0.1 + 0.2 \cdot 0.2 + 0.6 \cdot 0.6 = 0.42$$

$$\hat{y} = \arg\max_j \left\{ p(j = 0 \,|\, \mathbf{x}), p(j = 1 \,|\, \mathbf{x}) \right\}$$

# Overview



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

# Bagging

## (Bootstrap Aggregating)

Breiman, L. (1996). Bagging predictors. *Machine learning*, *24*(2), 123-140.
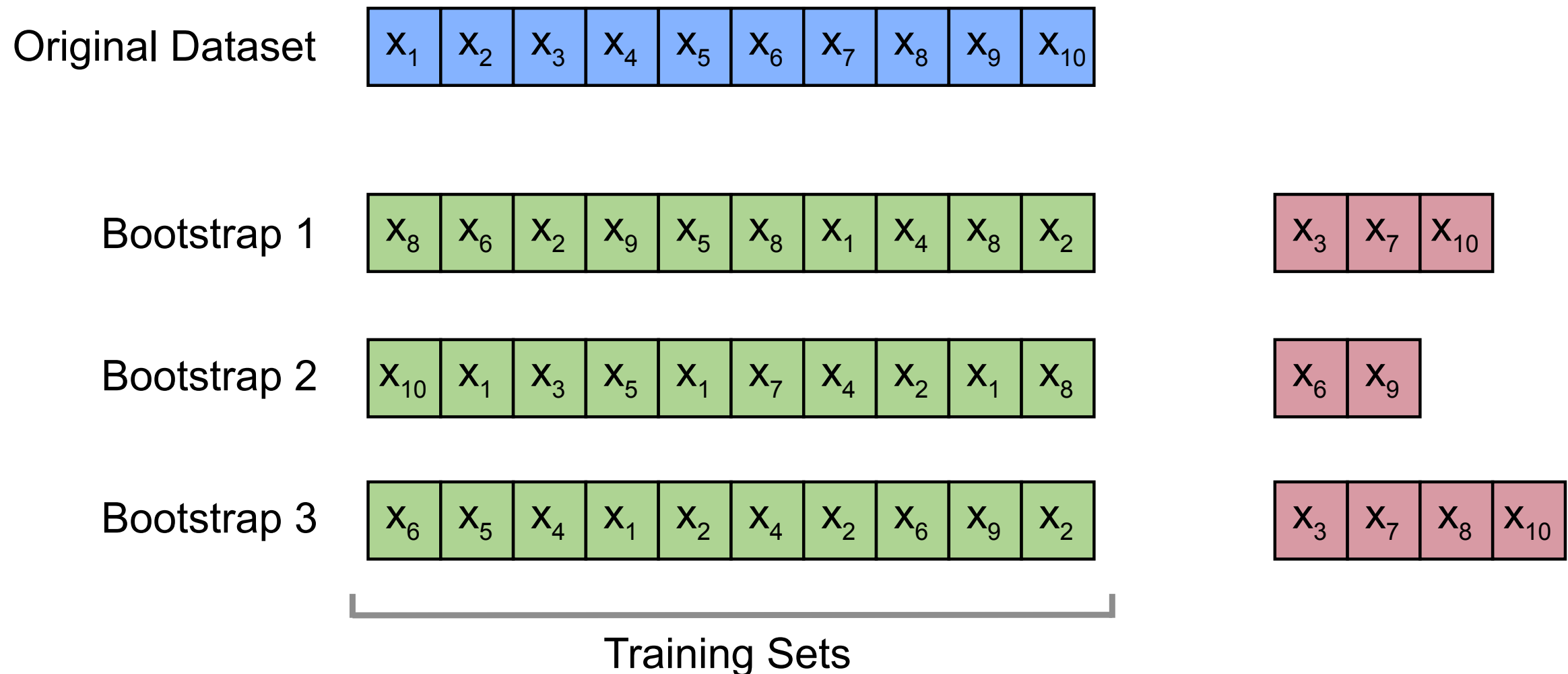
# Bagging
## (Bootstrap Aggregating)

---

**Algorithm 1** Bagging

---

1: Let $n$ be the number of bootstrap samples

2:

3: **for** i=1 to $n$ **do**

4:        Draw bootstrap sample of size $m$, $\mathcal{D}_i$

5:        Train base classifier $h_i$ on $\mathcal{D}_i$

6: $\hat{y} = mode\{h_1(\mathbf{x}), ..., h_n(\mathbf{x})\}$

---

# Bootstrap Sampling

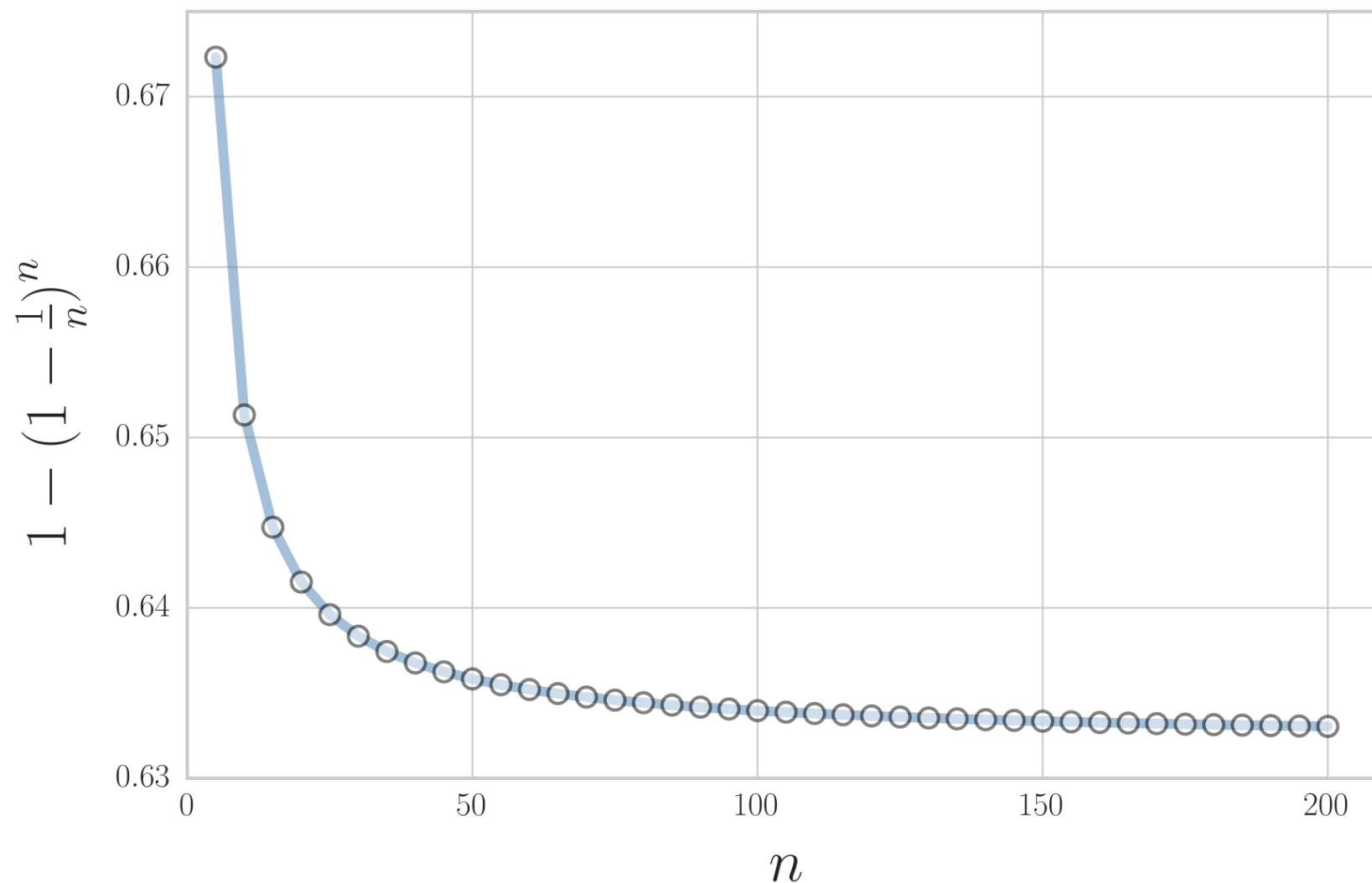**Original Dataset** $\quad$ | $X_1$ | $X_2$ | $X_3$ | $X_4$ | $X_5$ | $X_6$ | $X_7$ | $X_8$ | $X_9$ | $X_{10}$ |

**Bootstrap 1** $\quad$ | $X_8$ | $X_6$ | $X_2$ | $X_9$ | $X_5$ | $X_8$ | $X_1$ | $X_4$ | $X_8$ | $X_2$ | $\quad$ | $X_3$ | $X_7$ | $X_{10}$ |

**Bootstrap 2** $\quad$ | $X_{10}$ | $X_1$ | $X_3$ | $X_5$ | $X_1$ | $X_7$ | $X_4$ | $X_2$ | $X_1$ | $X_8$ | $\quad$ | $X_6$ | $X_9$ |

**Bootstrap 3** $\quad$ | $X_6$ | $X_5$ | $X_4$ | $X_1$ | $X_2$ | $X_4$ | $X_2$ | $X_6$ | $X_9$ | $X_2$ | $\quad$ | $X_3$ | $X_7$ | $X_8$ | $X_{10}$ |

**Training Sets**

# Bootstrap Sampling

$$P(\textbf{not chosen}) = \left( 1 - \frac{1}{n} \right)^{n},$$

$$\frac{1}{e} \approx 0.368, \quad n \rightarrow \infty \, .$$

$$P(\textbf{not chosen}) = \left(1 - \frac{1}{n}\right)^{n},$$
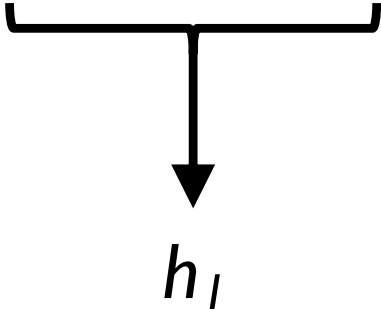
$$\frac{1}{e} \approx 0.368, \quad n \to \infty.$$

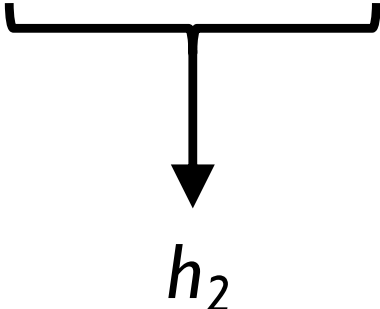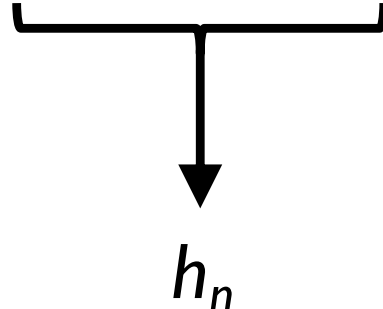$$P(\textbf{chosen}) = 1 - \left(1 - \frac{1}{n}\right)^{n} \approx 0.632$$
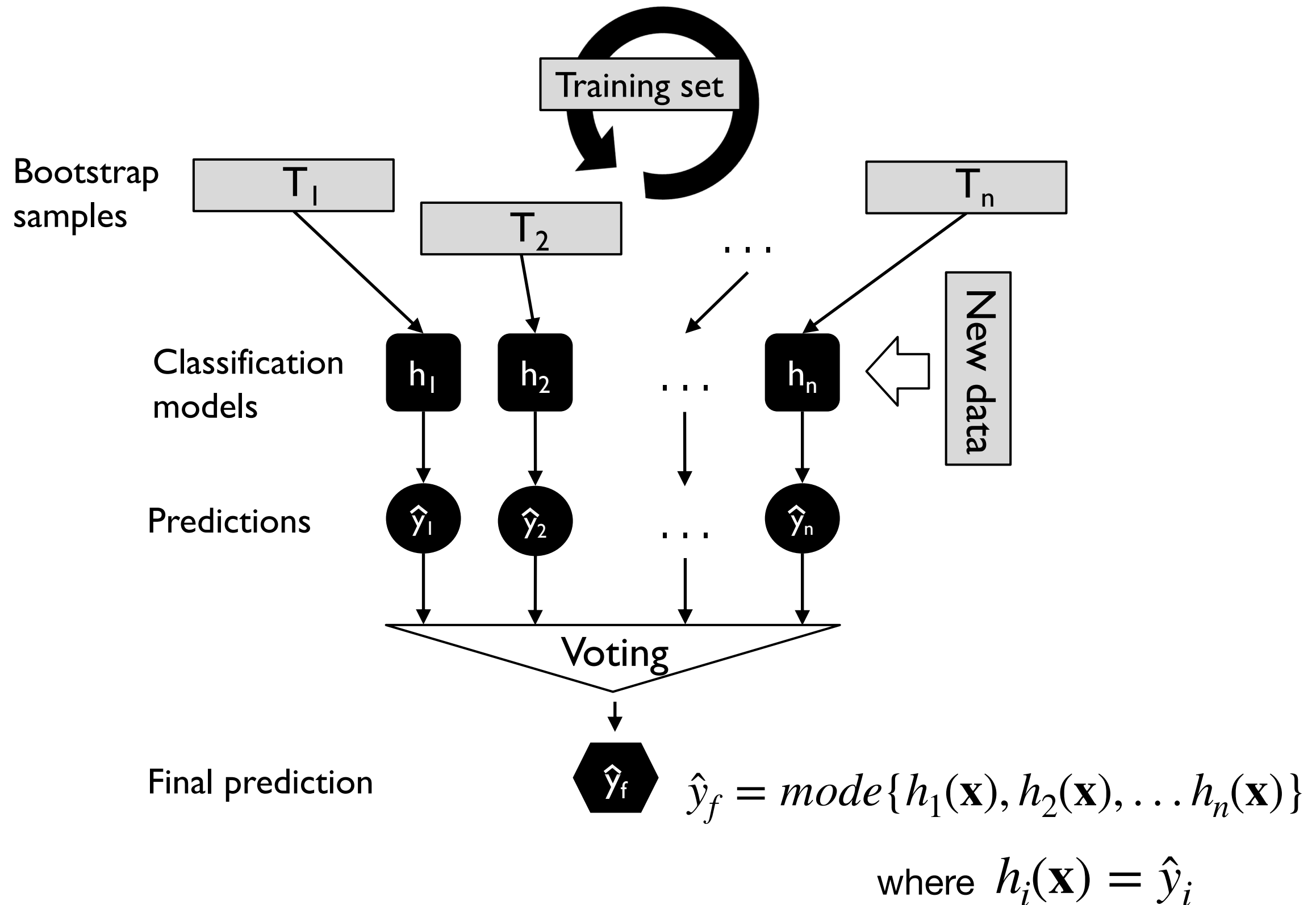
# Bootstrap Sampling

| Training example indices | Bagging round 1 | Bagging round 2 | ... |
|---|---|---|---|
| 1 | 2 | 7 | ... |
| 2 | 2 | 3 | ... |
| 3 | 1 | 2 | ... |
| 4 | 3 | 1 | ... |
| 5 | 7 | 1 | ... |
| 6 | 2 | 7 | ... |
| 7 | 4 | 7 | ... |

$h_1$      $h_2$      $h_n$

# Bagging Classifier



$$\hat{y}_f = mode\{h_1(\mathbf{x}), h_2(\mathbf{x}), \dots h_n(\mathbf{x})\}$$

where $h_i(\mathbf{x}) = \hat{y}_i$
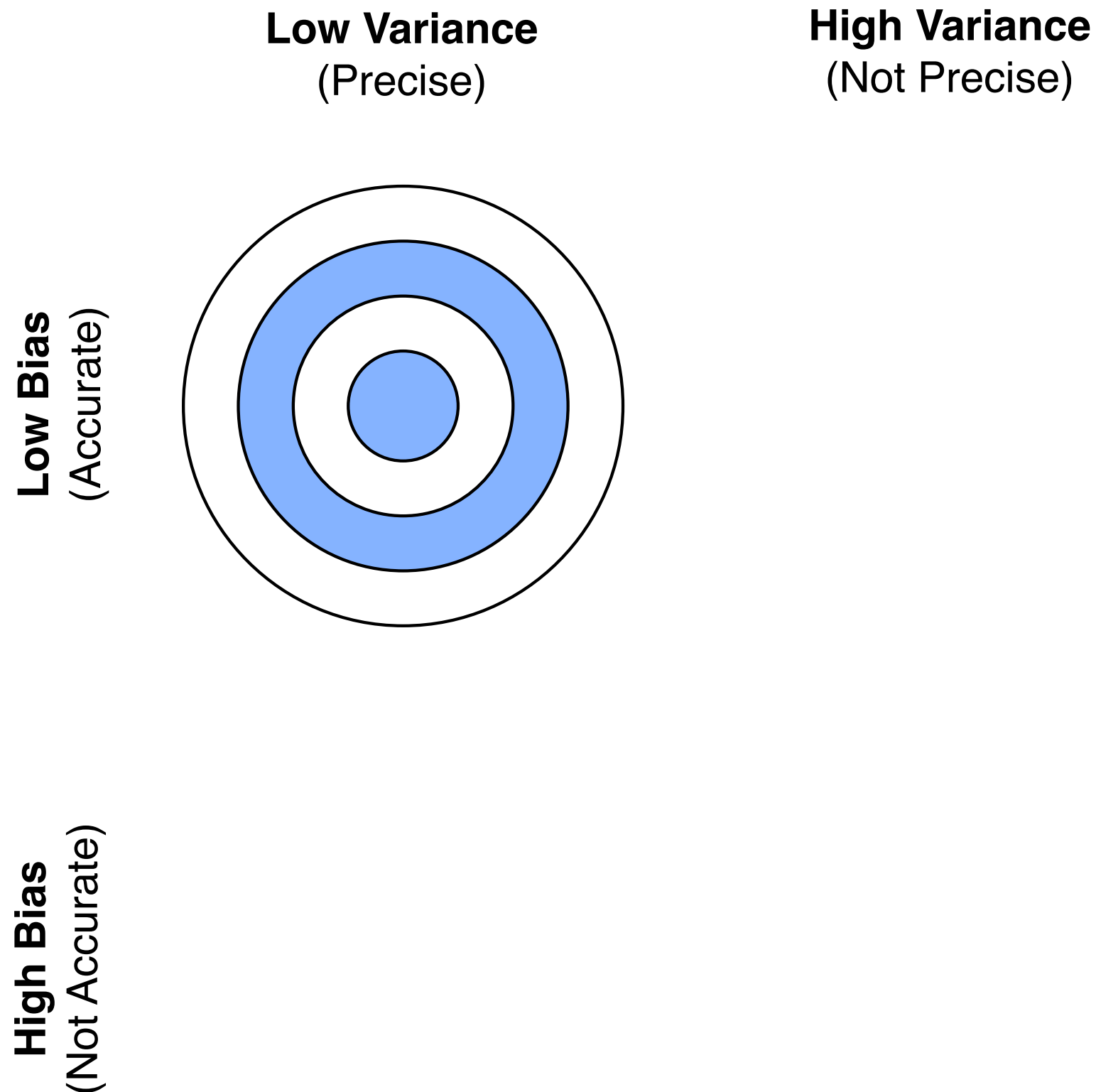
# Bias-Variance Decomposition
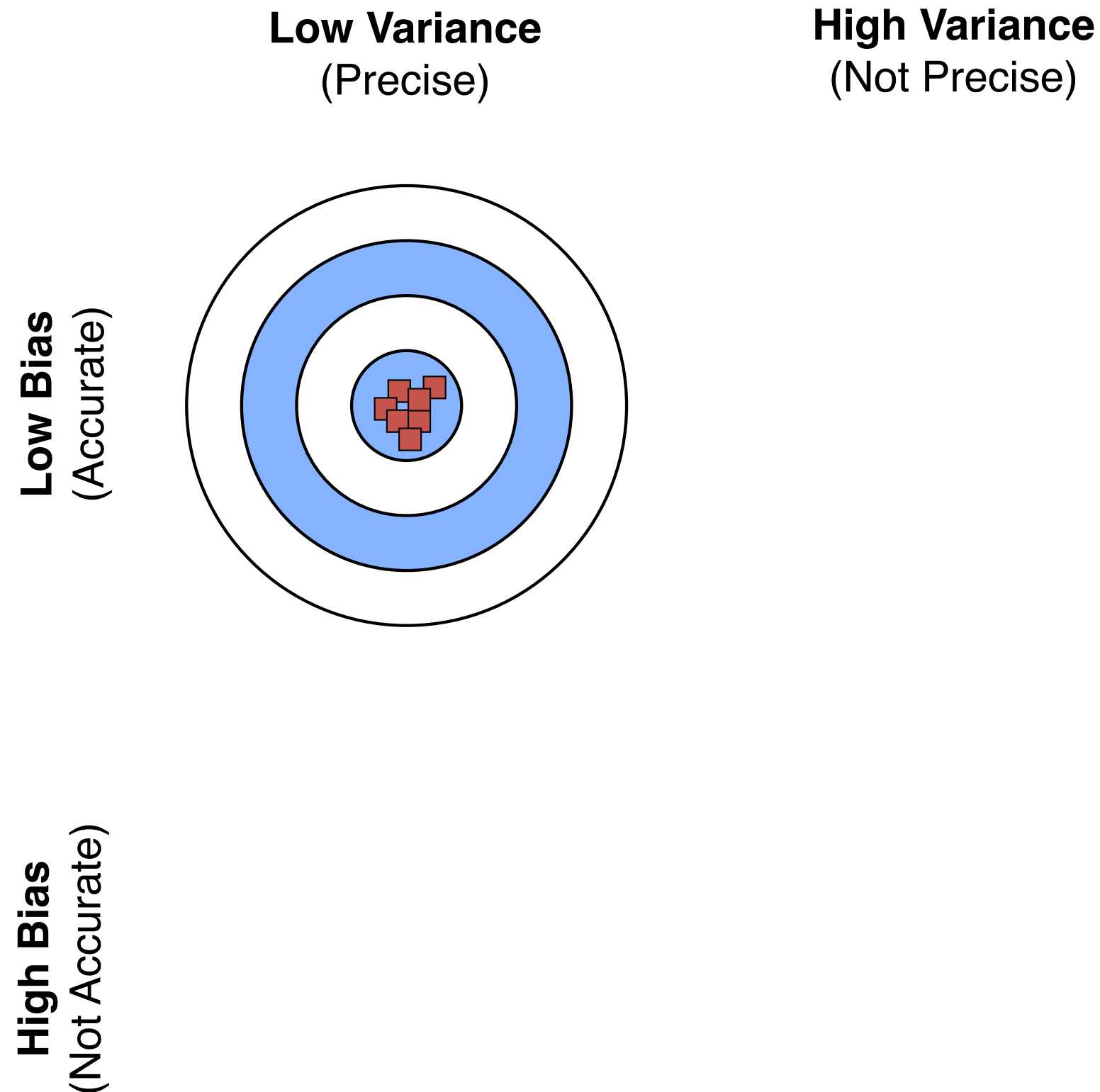
Loss = Bias + Variance + Noise

(more technical details in next lecture on model evaluation)

# Bias-Variance Intuition



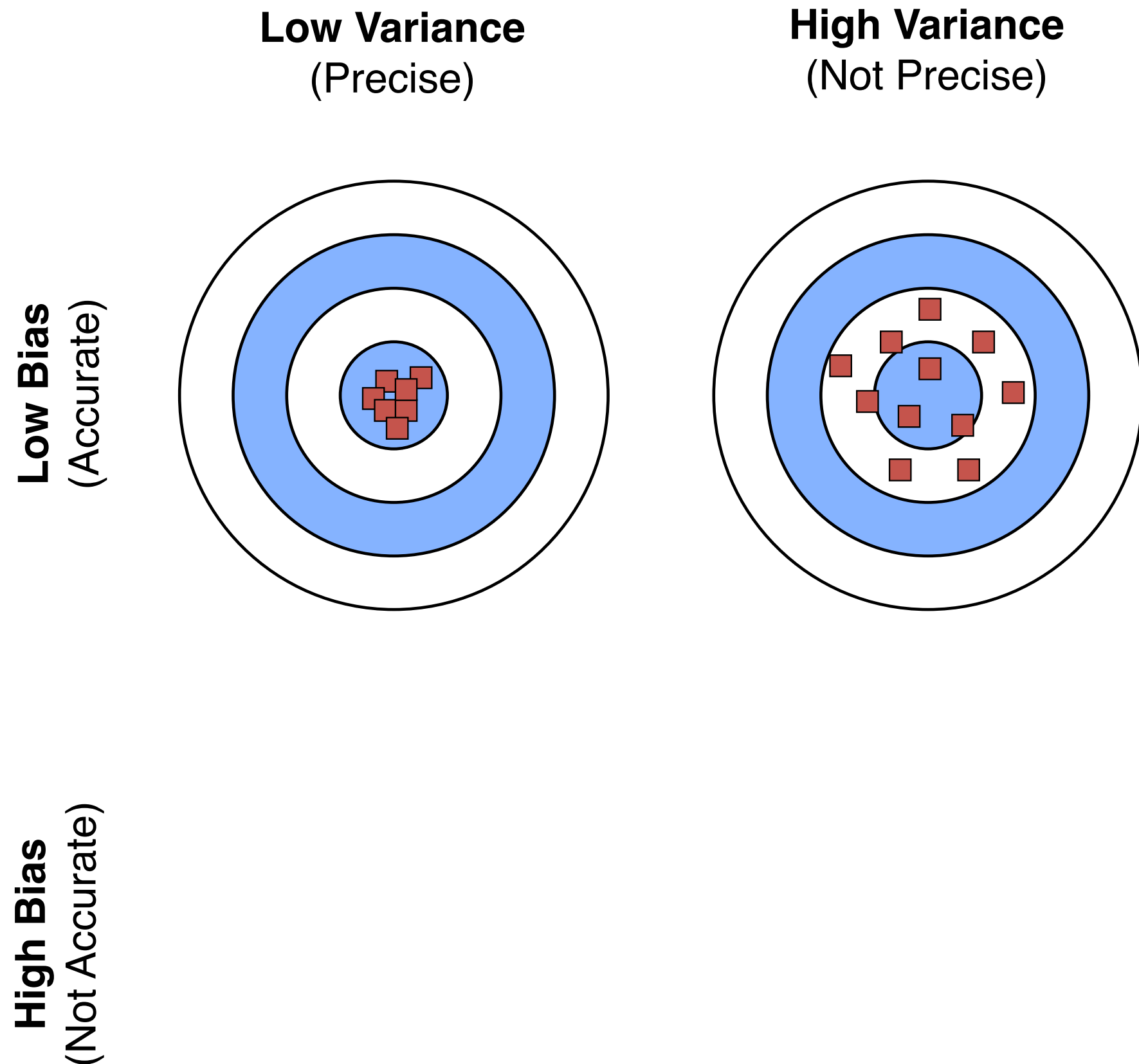**Low Variance**
(Precise)
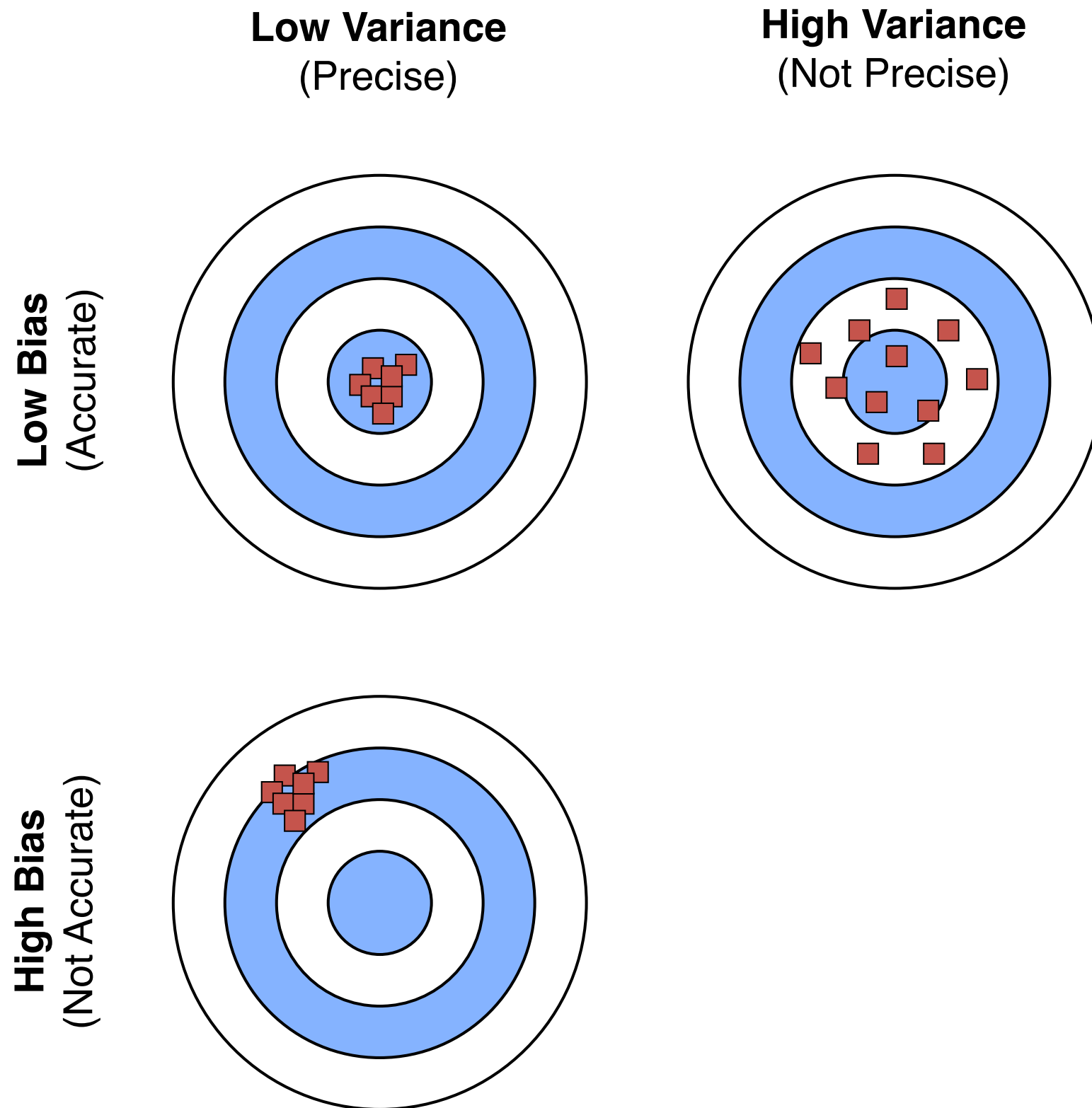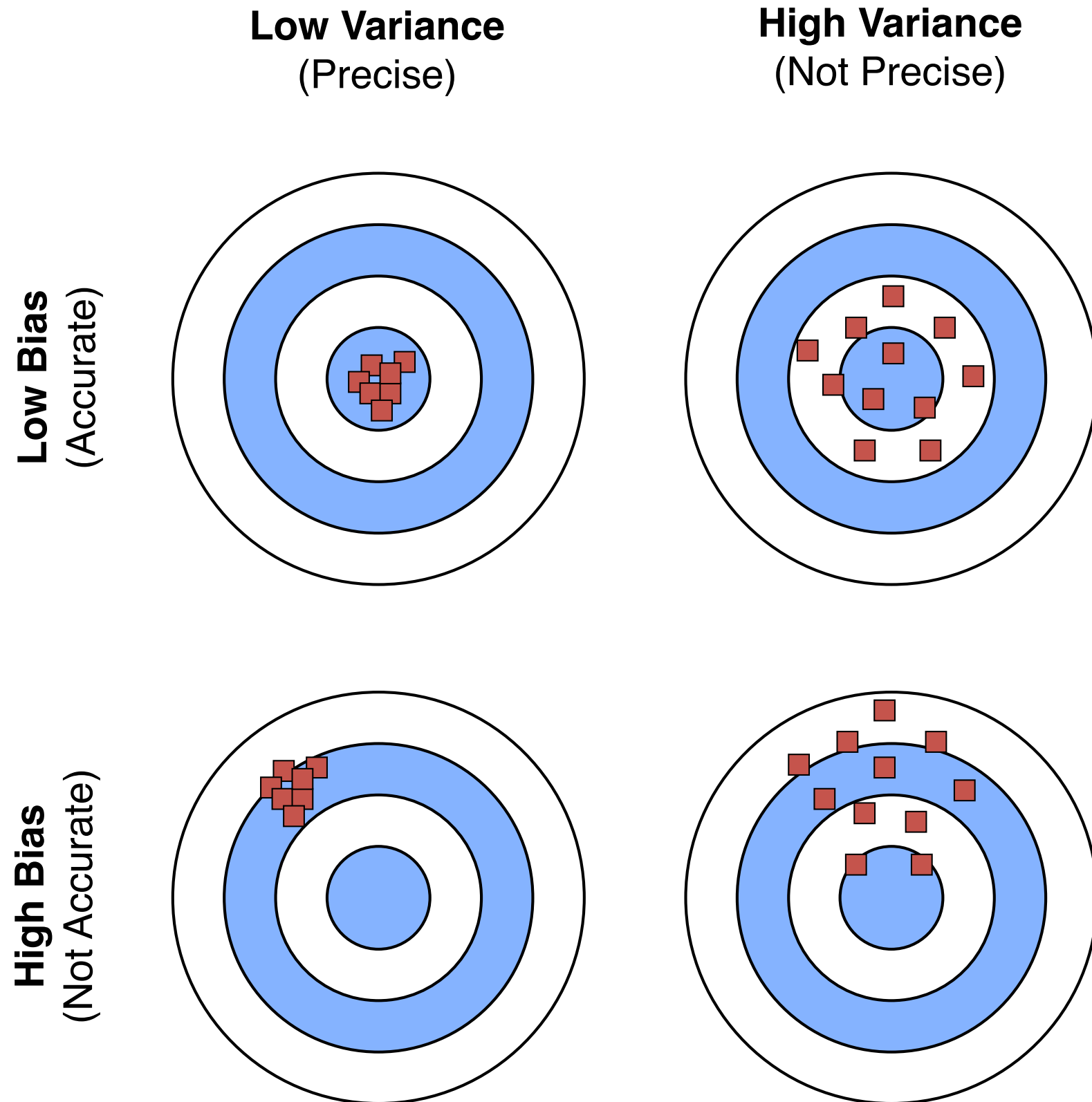
**High Variance**
(Not Precise)

**Low Bias**
(Accurate)

**High Bias**
(Not Accurate)

# Bias-Variance Intuition

**Low Variance**
(Precise)

**High Variance**
(Not Precise)

**Low Bias**
(Accurate)

**High Bias**
(Not Accurate)

# Bias-Variance Intuition

# Bias-Variance Intuition



**Low Variance**
(Precise)

**High Variance**
(Not Precise)

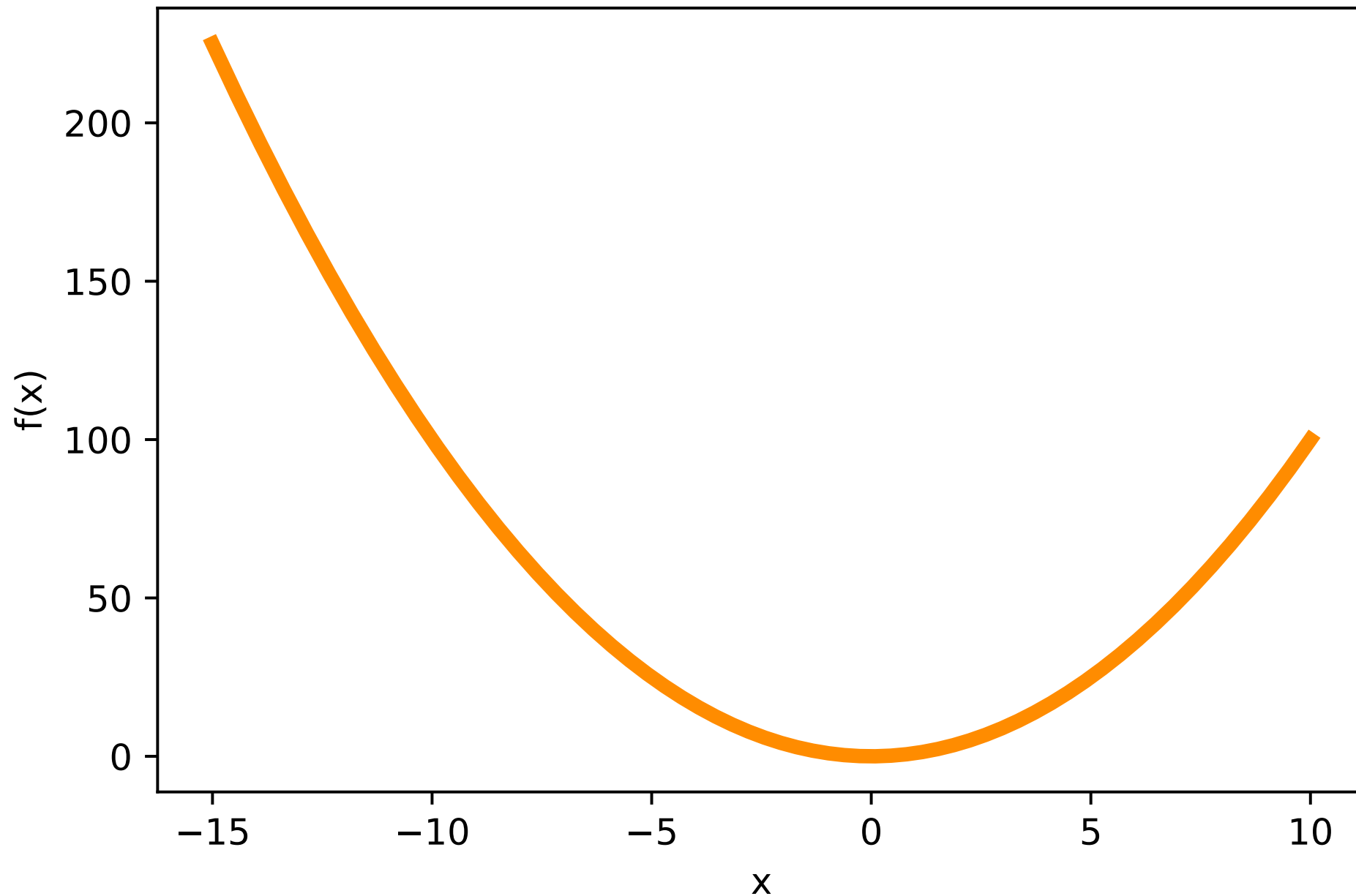**Low Bias**
(Accurate)

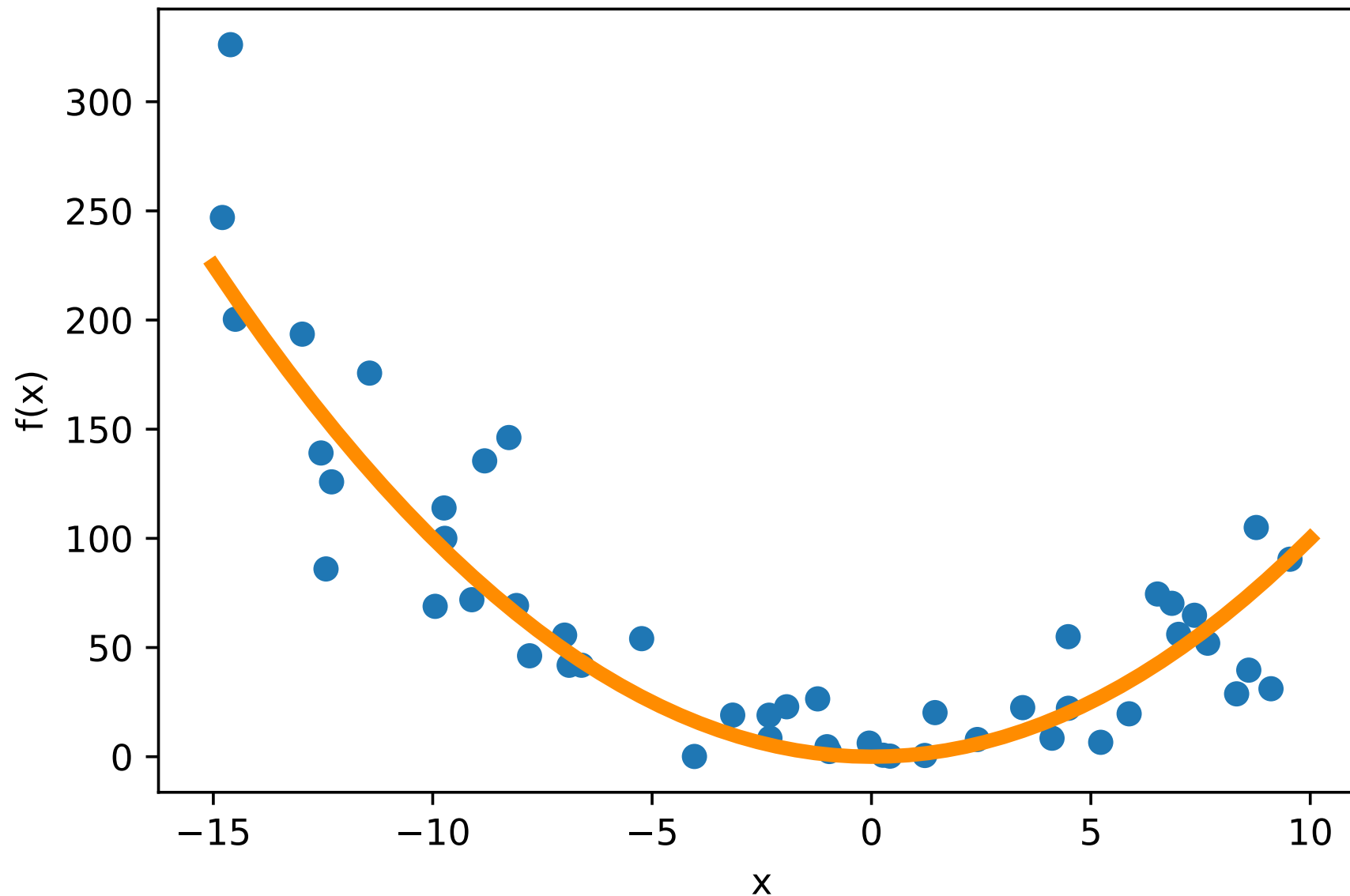**High Bias**
(Not Accurate)

# Bias-Variance Intuition

# Bias and Variance Example



where f(x) is some true (target) function

# Bias and Variance Example



where f(x) is some true (target) function

the blue dots are a training dataset;
here, I added some random Gaussian noise

# Bias and Variance Example



where f(x) is some true (target) function

the blue dots are a training dataset;
here, I added some random Gaussian noise

here, suppose I fit a simple linear model  (linear regression)
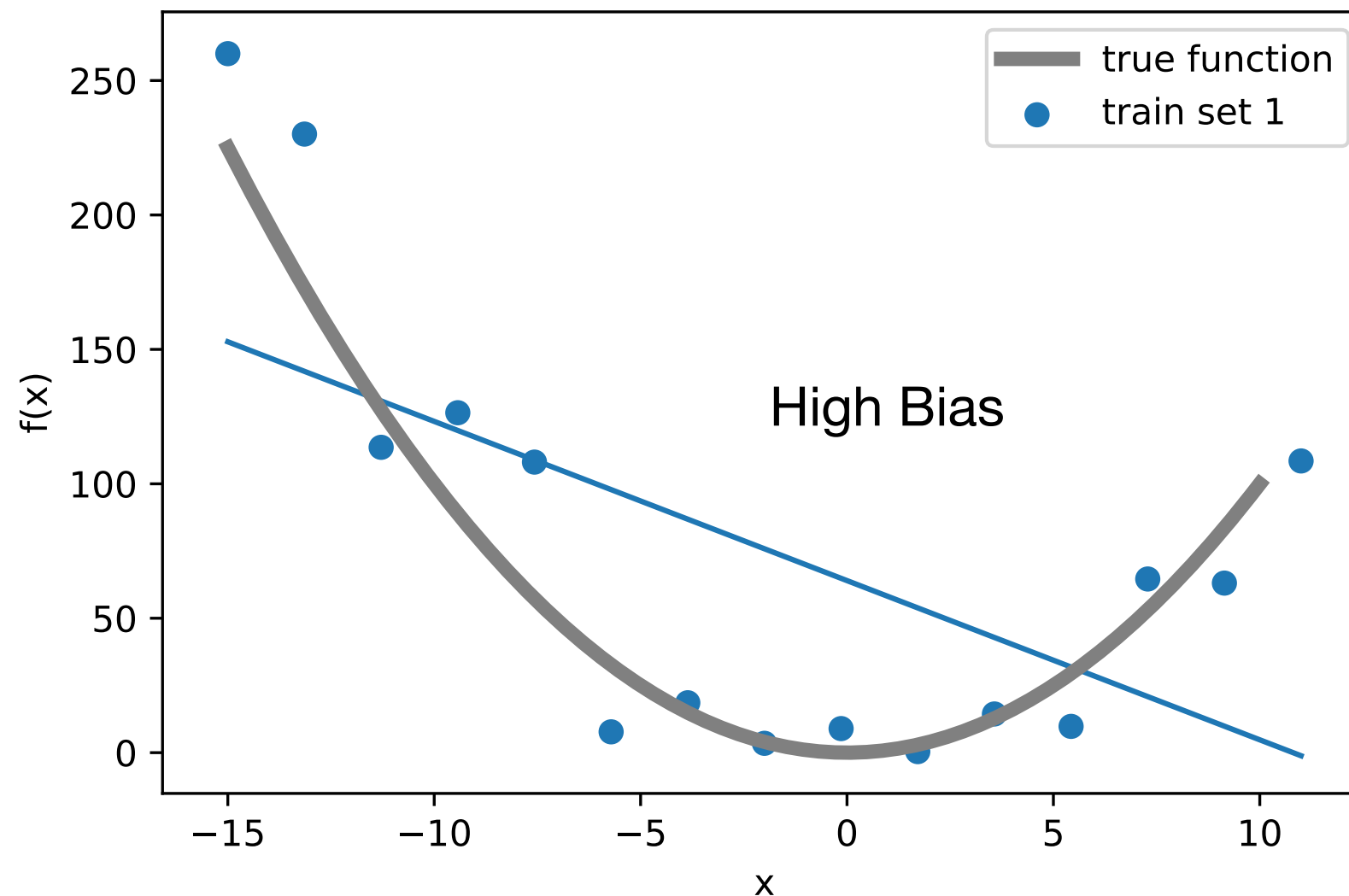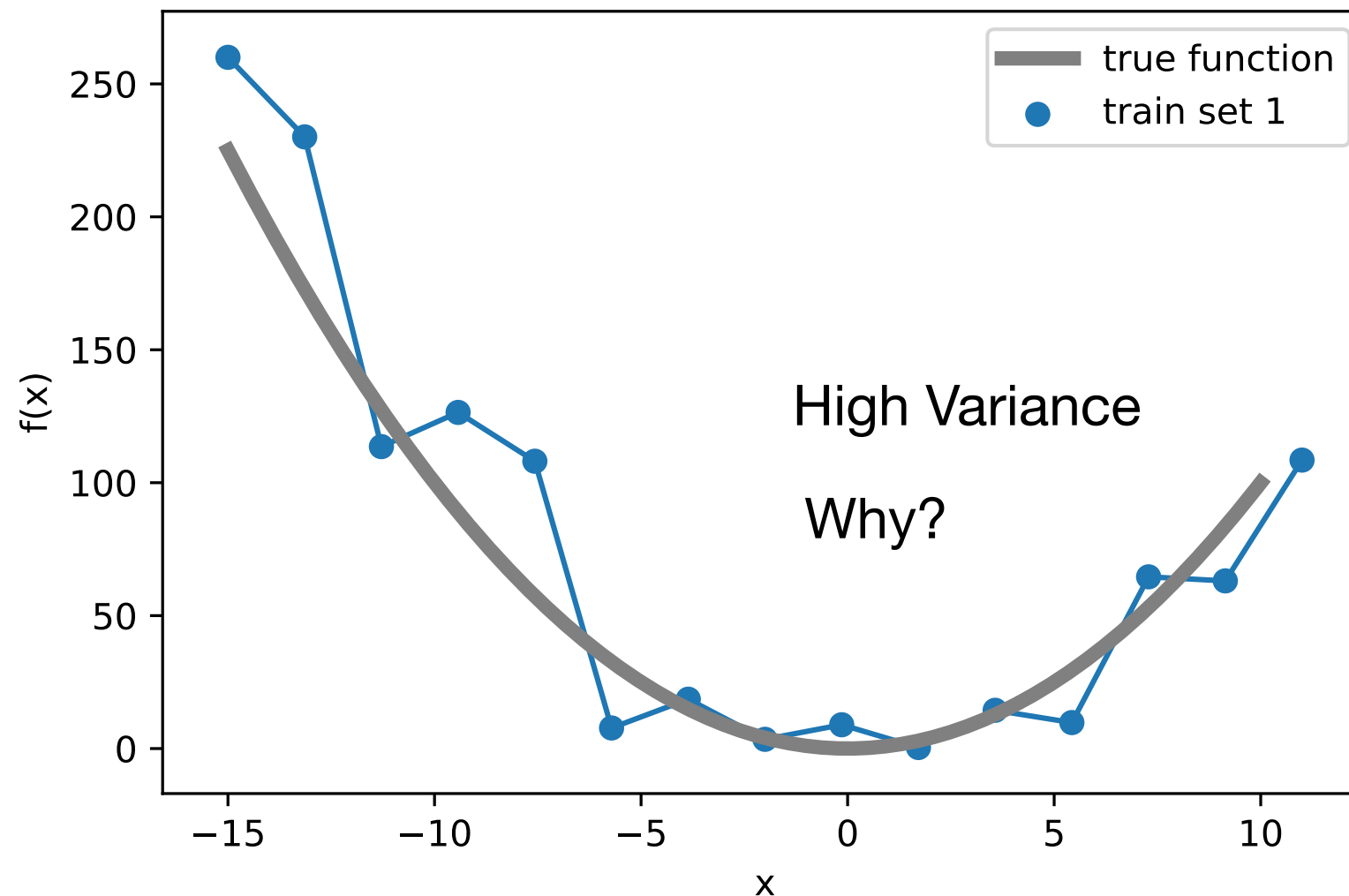or a decision tree stump
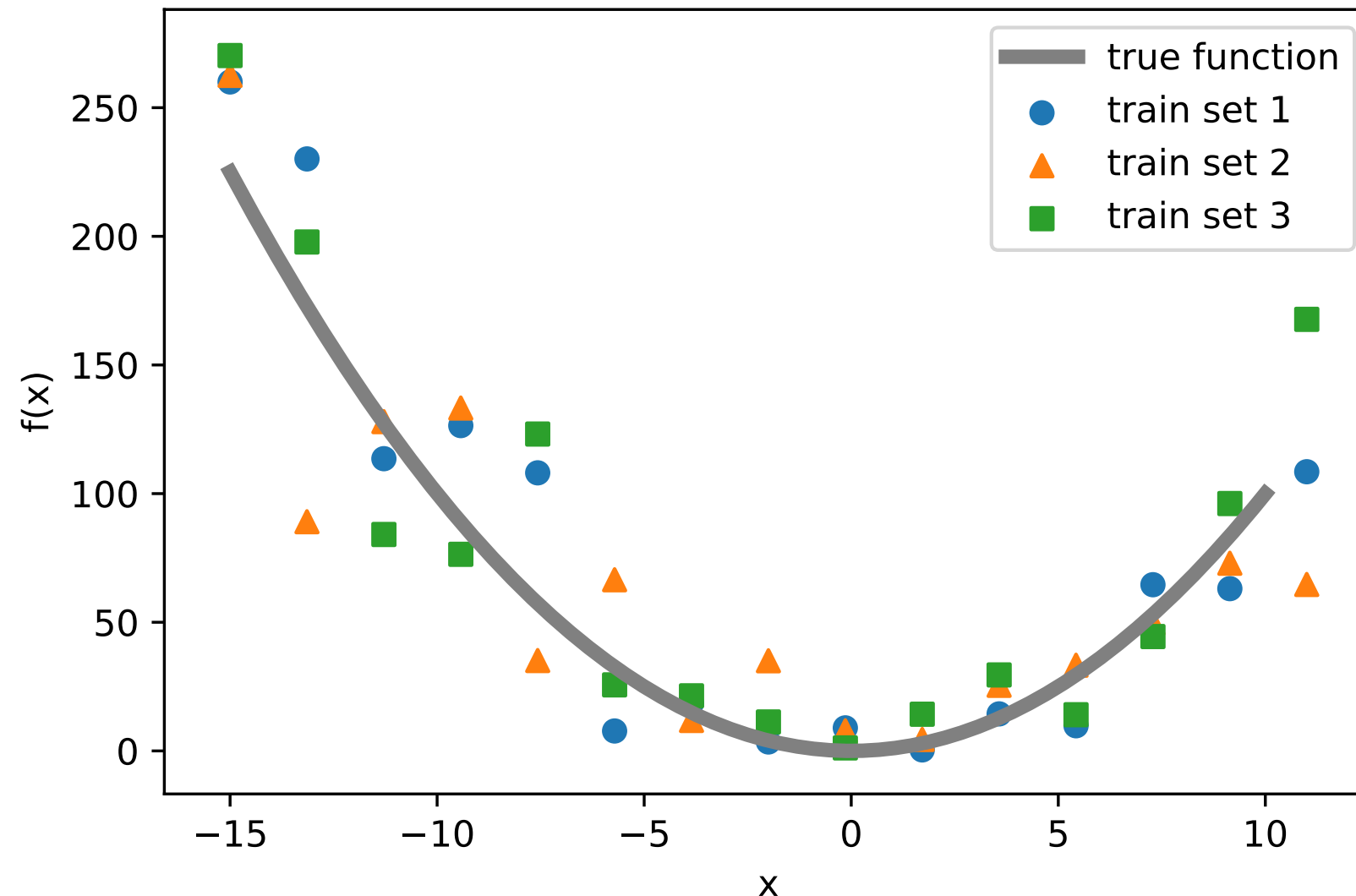
# Bias and Variance Example



where f(x) is some true (target) function

the blue dots are a training dataset;
here, I added some random Gaussian noise

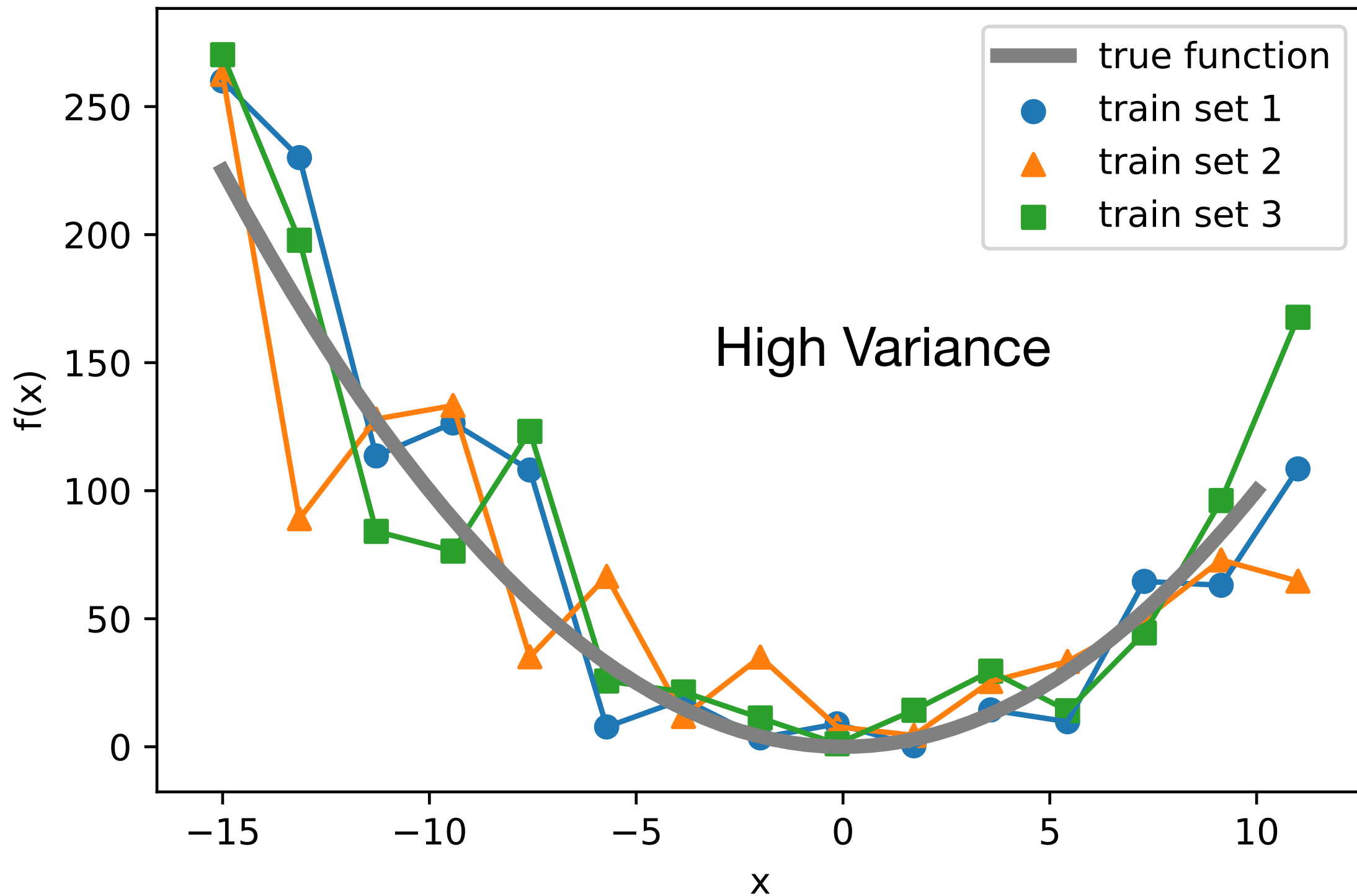here, suppose I fit an unpruned decision tree

# Bias and Variance Example



where f(x) is some true (target) function
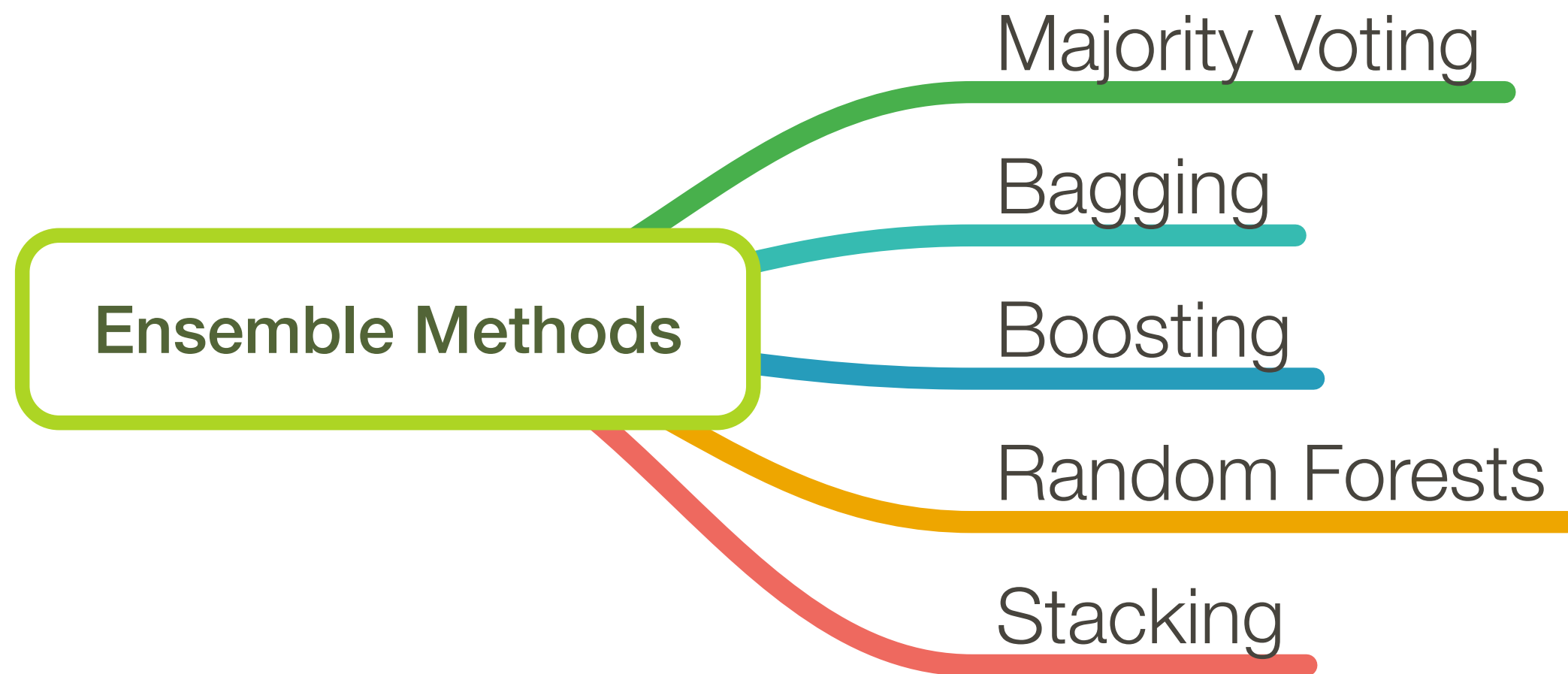
suppose we have multiple training sets

# Bias and Variance Example

# So, why does bagging work/what does it do?

# Overview



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

# Boosting

# Adaptive Boosting

e.g., AdaBoost

Freund, Y., & Schapire, R. E. (1997). A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, *55*(1), 119-139.

# Gradient Boosting

e.g., LightGBM, XGBoost, scikit-learn's GradientBoostingClassifier

Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of statistics*, 1189-1232.

Chen, T., & Guestrin, C. (2016). 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.

# Adaptive Boosting

e.g., AdaBoost

Freund, Y., & Schapire, R. E. (1997). A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, *55*(1), 119-139.

Differ mainly in terms of how
- weights are updated
- classifiers are combined

# Gradient Boosting
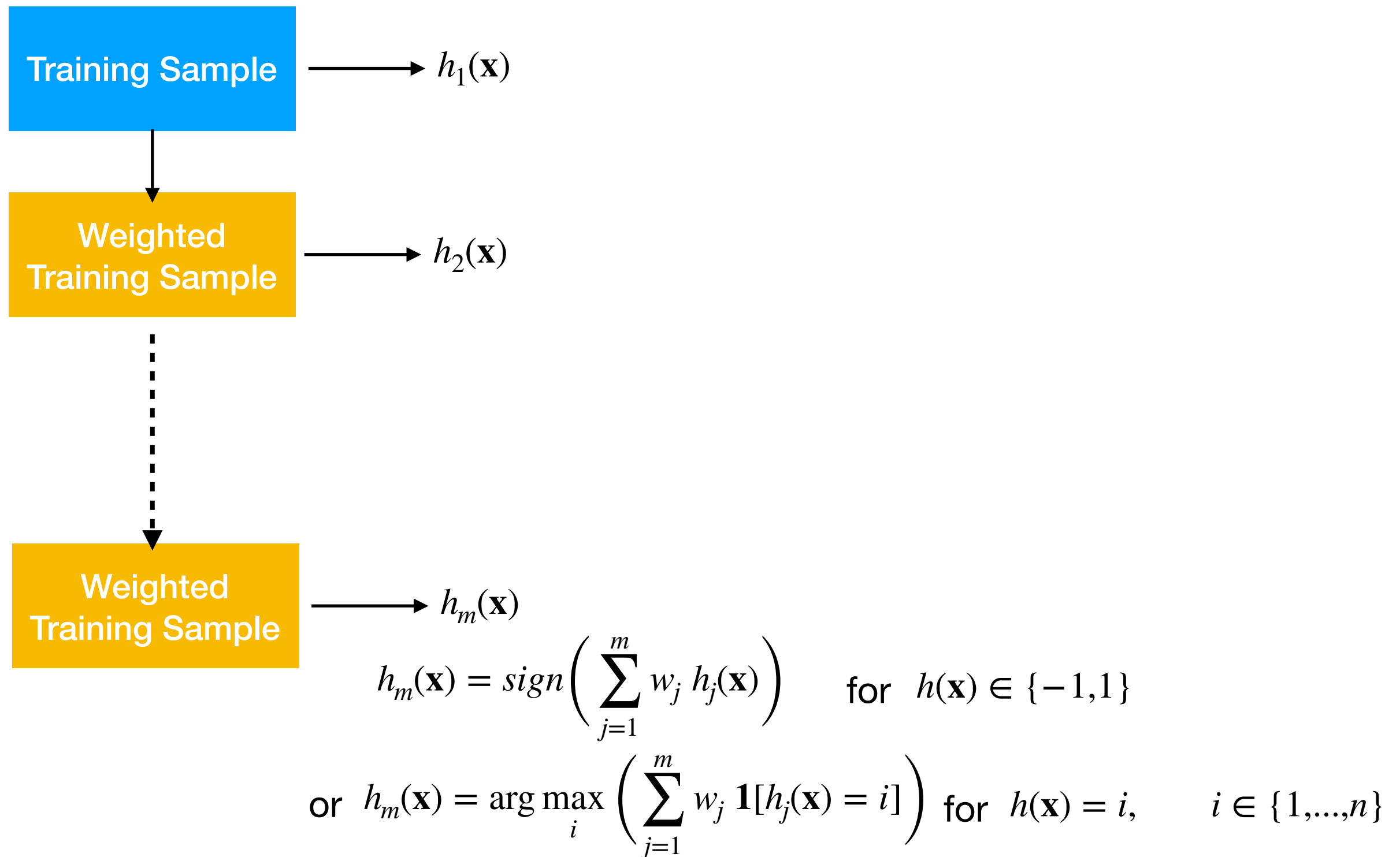
e.g., LightGBM, XGBoost, scikit-learn's GradientBoostingClassifier

Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of statistics*, 1189-1232.

Chen, T., & Guestrin, C. (2016). 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.

# General Boosting



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

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

# 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

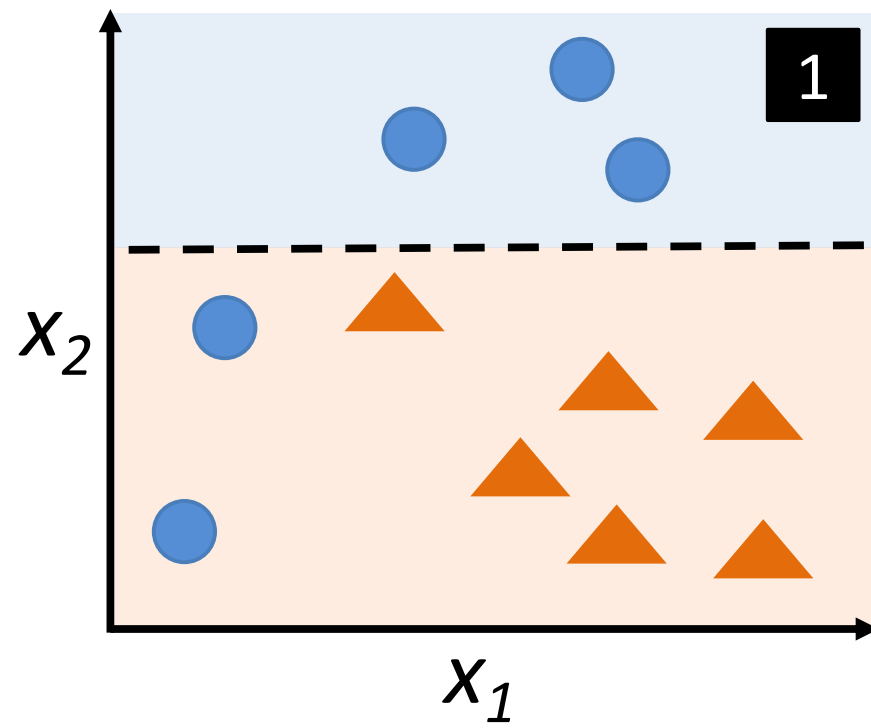* a learner slightly better than random guessing
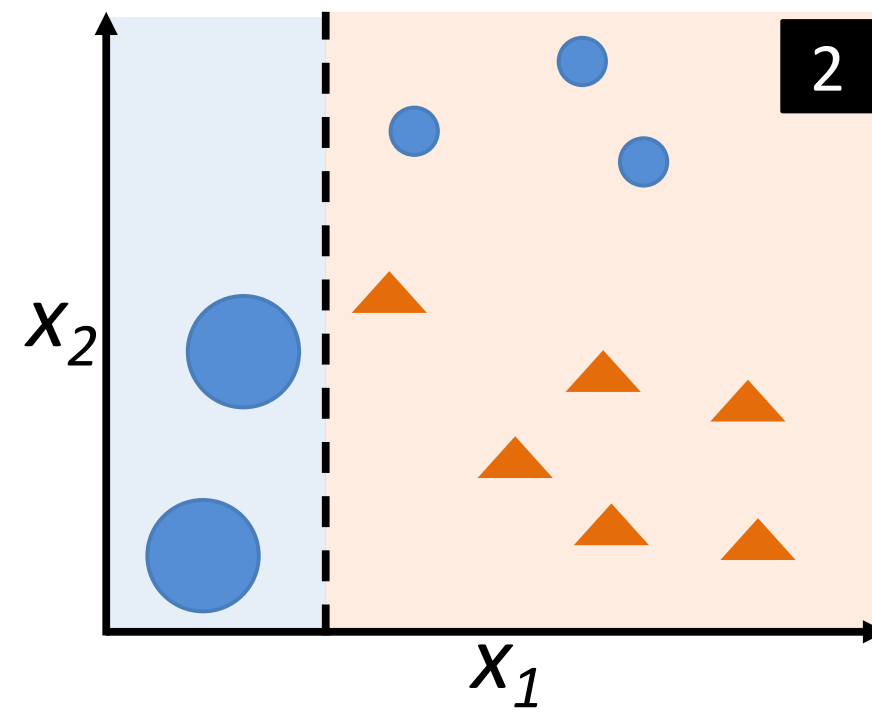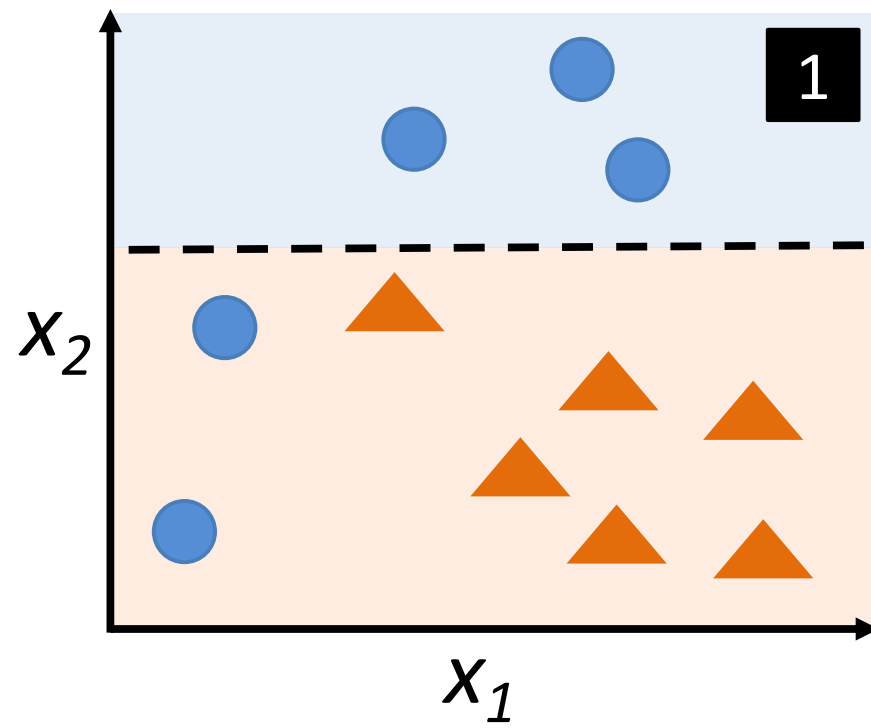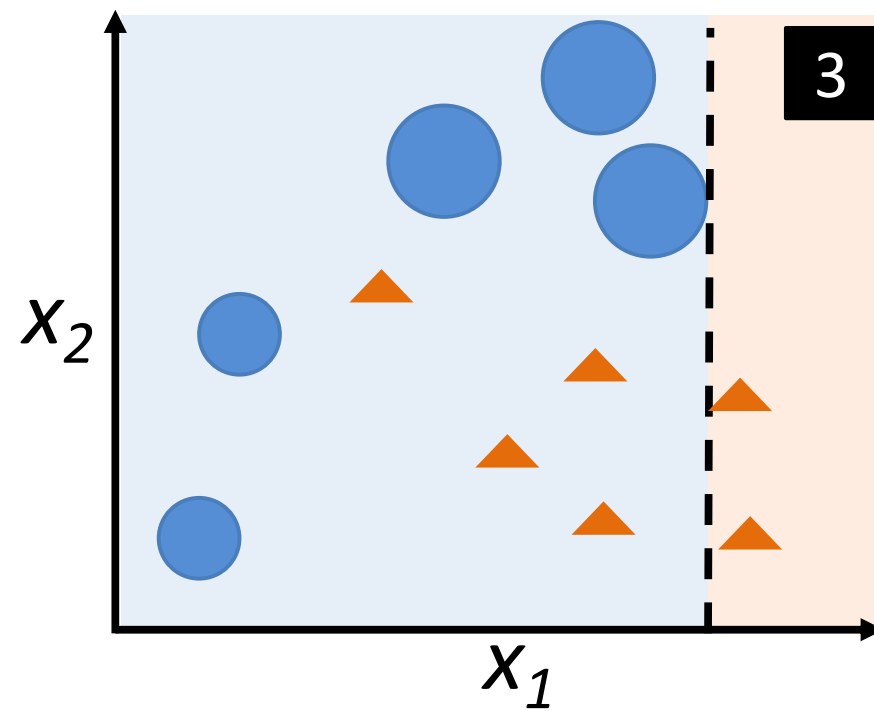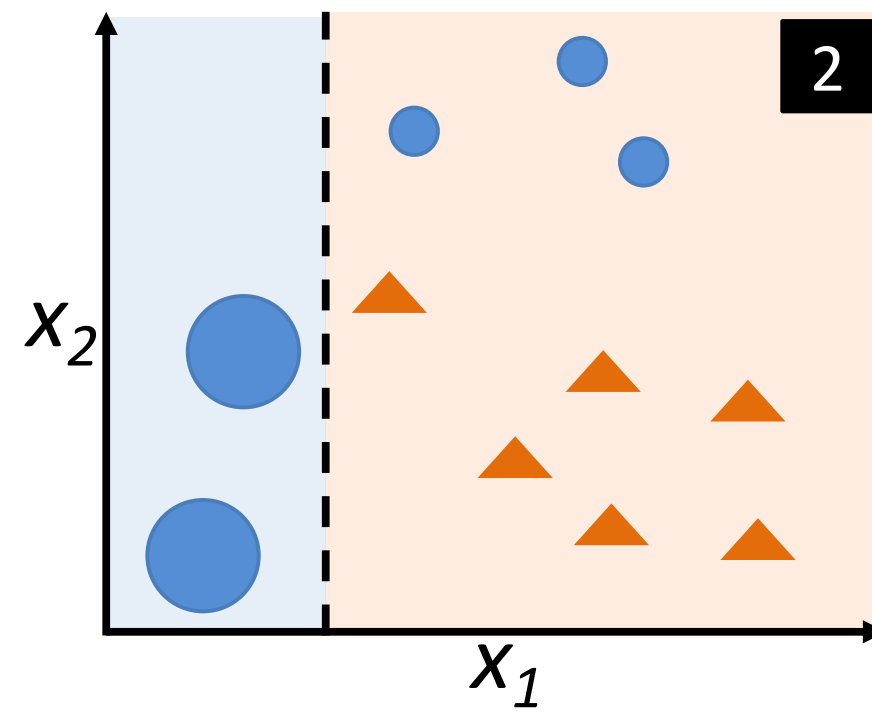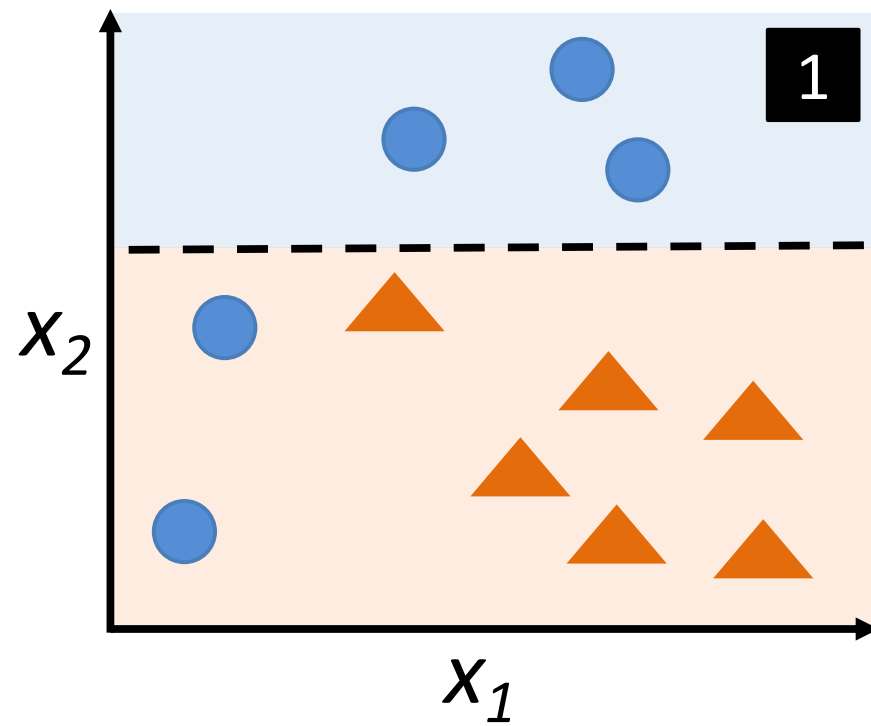
# AdaBoost

---

**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: $\quad$ For all $i : \mathbf{w}_r(i) := w_r(i)/\sum_i w_r(i)$ $\quad$ [normalize weights]
7: $\quad$ $h_r := FitWeakLearner(\mathcal{D}, \mathbf{w}_r)$
8: $\quad$ $\epsilon_r := \sum_i w_r(i) \mathbf{1}(h_r(i) \neq y_i)$ $\quad$ [compute error]
9: $\quad$ if $\epsilon_r > 1/2$ then stop
10: $\quad$ $\alpha_r := \frac{1}{2} \log[(1 - \epsilon_r)/\epsilon_r]$ $\quad$ [small if error is large and vice versa]
11: $\quad$ $w_{r+1}(i) := w_r(i) \times \begin{cases} \mathrm{e}^{-\alpha_r} & \text{if } h_r(\mathbf{x}^{[i]}) = y^{[i]} \\ \mathrm{e}^{\alpha_r} & \text{if } h_r(\mathbf{x}^{[i]}) \neq y^{[i]} \end{cases}$

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

---

# Decision Tree Stumps

# 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

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

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

- **<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: $r_1 = y - \hat{y}_1$

| x1# | x2=City | x3=Age | y=Price | r=Res |
|:---:|:---:|:---:|:---:|:---:|
| 5 | Boston | 30 | $1.5 x 10^6 | 1.5 - 0.5875 = 0.9125 |
| 10 | Madison | 20 | $0.5 x 10^6 | 0.5 - 0.5875 = -0.0875 |
| 6 | Lansing | 20 | $0.25 x 10^6 | 0.25 - 0.5875 = -0.3375 |
| 5 | Waunake | 10 | $0.1 x 10^6 | 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 x1, ..., x3 to fit the residuals

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

Age >= 30

No          Yes

# Rooms >= 10          0.9125

-0.4125 ←   -0.3375          -0.0875
            -0.4875

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

- **Step 3:** Combine tree from step 1 with trees from step 2

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

-0.3375
-0.4875

→ -0.4125

-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 x 10^6 | 1.5 - 0.5875 = 0.9125 |
| 10 | Madison | 20 | $0.5 x 10^6 | 0.5 - 0.5875 = -0.0875 |
| 6 | Lansing | 20 | $0.25 x 10^6 | 0.25 - 0.5875 = -0.3375 |
| 5 | Waunakee | 10 | $0.1 x 10^6 | 0.1 - 0.5875 = -0.4875 |

E.g., predict Lansing →

**Age >= 30**

No    Yes

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

**# Rooms >= 10**

0.9125

-0.4125 ← -0.3375
-0.4875

-0.0875

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

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

**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

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

**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

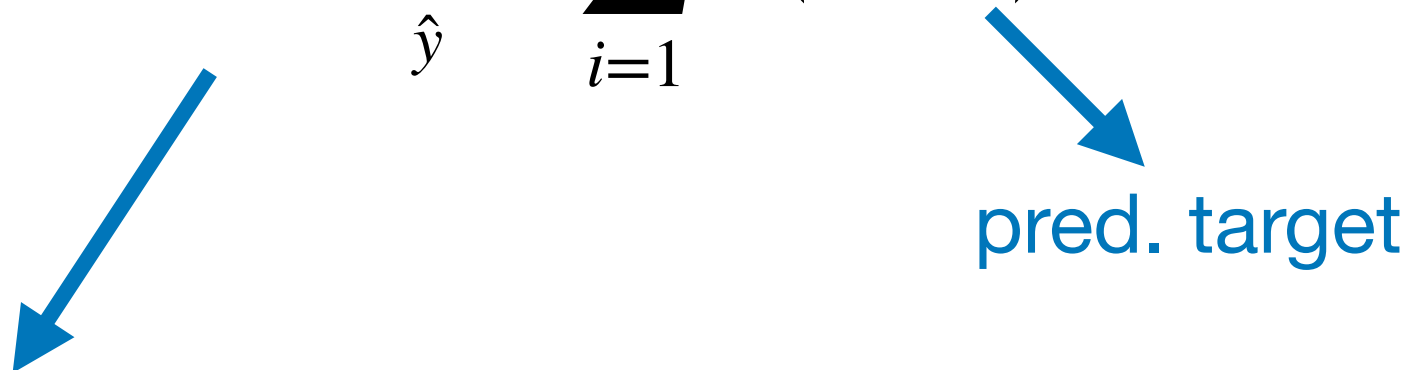$$SSE' = \frac{1}{2}\big(y^{(i)} - h(\mathbf{x}^{(i)})\big)^2$$

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

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

[neg. residual]

# Gradient Boosting -- Algorithm Overview <span style="color:#1b82c4">Discussion</span>

**<u>Step 1:</u>** Initialize model $h_0(\mathbf{x}) = \underset{\hat{y}}{\arg\min} \sum_{i=1}^{n} L\left(y^{(i)}, \hat{y}\right)$

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

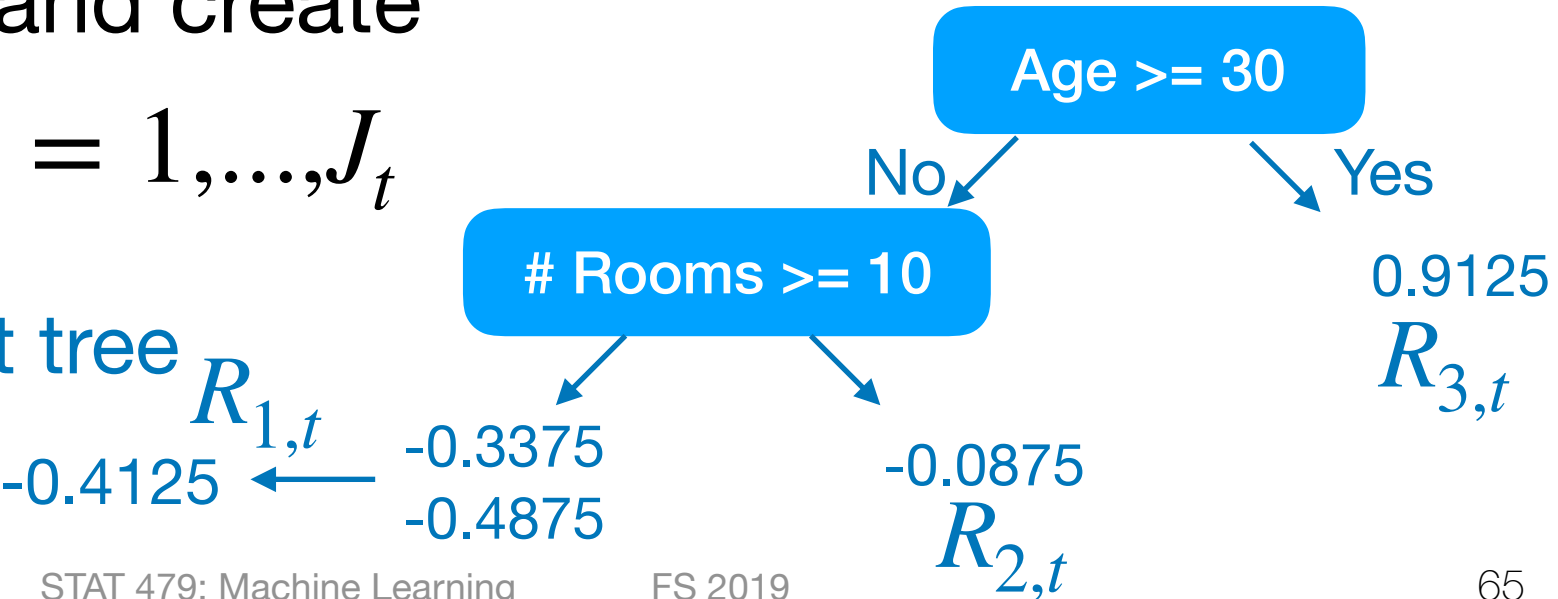# Gradient Boosting -- Algorithm Overview **Discussion**

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

0.9125
$R_{3,t}$

$R_{1,t}$

-0.3375
-0.4875

-0.4125

-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

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

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

**<u>Step 2:</u>** 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

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

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

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.,

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

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

...

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

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

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

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

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

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

$$\ldots$$

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

# Gradient Boosting -- Classification

Replace "average" of $h_0$ by *log(odds)*

probability of an event

$$odds = \frac{p}{1-p}$$

$$p = \frac{odds}{1 + odds} = \frac{e^{log(odds)}}{1 + e^{log(odds)}}$$

# Gradient Boosting -- Classification

Replace "average" of $h_0$ by *log(odds)*

Pseudo residual becomes $(y - p)$

Prediction (*log(odds)* & *residual*) are computed via the following transform:

$$\frac{\sum_i residual^{(i)}}{\sum_i p_{t-1}^{(i)} \times (1 - p_{t-1}^{(i)})}$$

# Gradient Boosting -- Classification

Replace "average" of $h_0$ by *log(odds)*

Pseudo residual becomes $(y - p)$

Prediction (*log(odds)* & *residual*) are computed via the following transform:
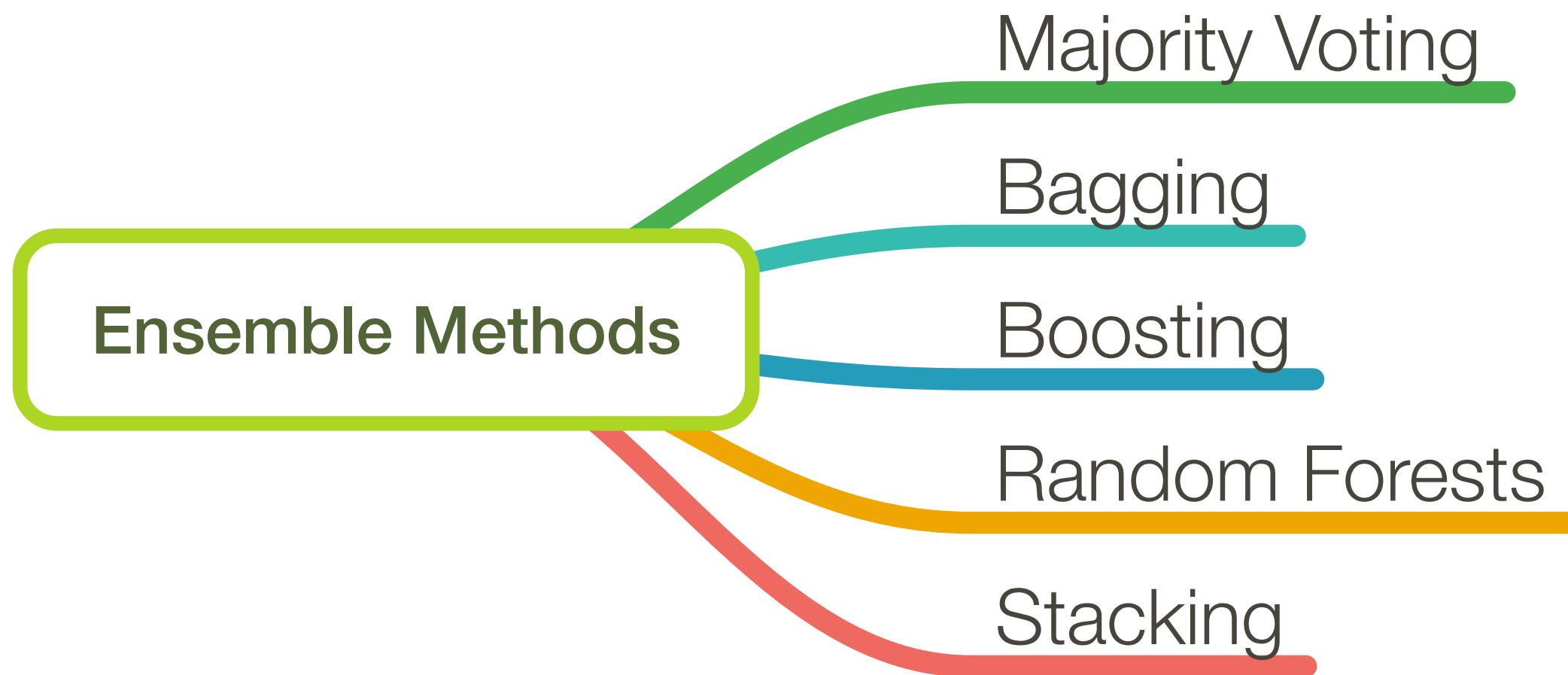
$$\frac{\sum_i residual^{(i)}}{\sum_i p_{t-1}^{(i)} \times (1 - p_{t-1}^{(i)})}$$

Loss function becomes neg. log likelihood

$$L(y^{(i)}, h(\mathbf{x})) = -\sum_{i=1}^{n} y^{(i)} \times \log(p) + (1 - y^{(i)}) \times \log(1 - p)$$

# Overview

**Ensemble Methods**

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

# Random Forests

# Random Forests

= Bagging w. trees + random feature subsets

# Random Feature Subset for each Tree or Node?

**Tin Kam Ho** used the **"random subspace method,"** where each tree got a random subset of features.

"Our method relies on an autonomous, pseudo-random procedure to select a small number of dimensions from a given feature space …"

- Ho, Tin Kam. "The random subspace method for constructing decision forests." IEEE transactions on pattern analysis and machine intelligence 20.8 (1998): 832-844.

**"Trademark" random forest:**

*"… random forest with random features is formed by selecting at random, at each node, a small group of input variables to split on."*

- **Breiman**, Leo. "Random Forests" Machine learning 45.1 (2001): 5-32.

# Random Feature Subset for each Tree or Node?

Tin Kam Ho used the "random subspace method," where each tree got a random subset of features.

"Our method relies on an autonomous, pseudo-random procedure to select a small number of dimensions from a given feature space …"

- Ho, Tin Kam. "The random subspace method for constructing decision forests." IEEE transactions on pattern analysis and machine intelligence 20.8 (1998): 832-844.

"Trademark" random forest:

*"… random forest with random feature* [...] *andom, at each node, a small group of input variab* [...]

**num features** $= \log_2 m + 1$

where $m$ is the number of input features

- Breiman, Leo. "Random Forests" Ma [...] 2.

In contrast to the original publication
[Breiman, "Random Forests", Machine Learning, 45(1), 5-32, 2001]
the scikit-learn implementation combines classifiers by averaging their
probabilistic prediction, instead of letting each classifier vote for a single
class.

"Soft Voting"

# Will discuss Random Forests and feature importance in *Feature Selection* lecture

# (Loose) Upper Bound for the Generalization Error

Breiman, "Random Forests", Machine Learning, 45(1), 5-32, 2001

$$\textbf{PE} \leq \frac{\bar{\rho} \cdot (1 - s^2)}{s^2}$$

$\bar{\rho}$ : Average correlation among trees

$s$ : "Strength" of the ensemble

# Extremely Randomized Trees (ExtraTrees)

Geurts, P., Ernst, D., & Wehenkel, L. (2006). Extremely randomized trees. *Machine learning*, *63*(1), 3-42.
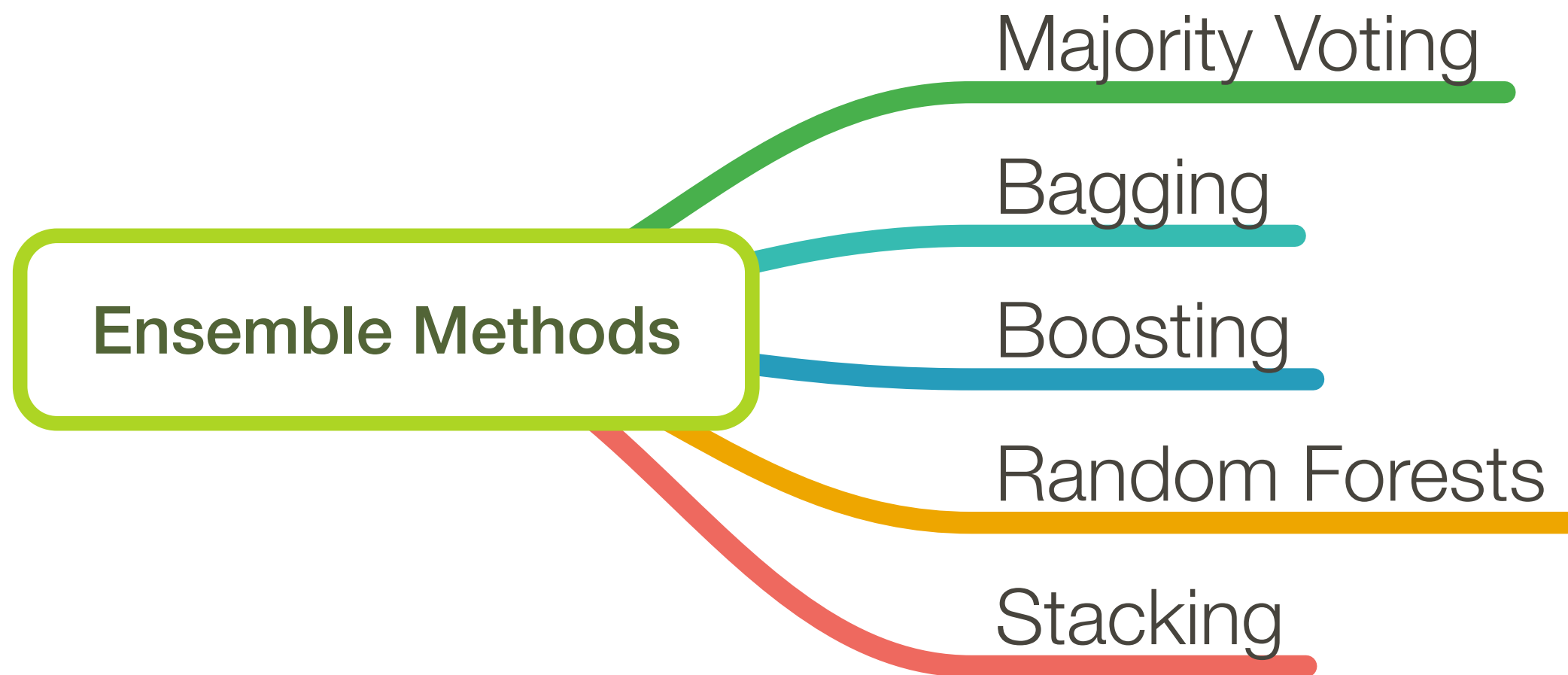
Random Forest random components:

1) _____

2) _____

ExtraTrees algorithm adds one more random component

3) _____

# Overview

Ensemble Methods

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

# Stacking

# Stacking Algorithm

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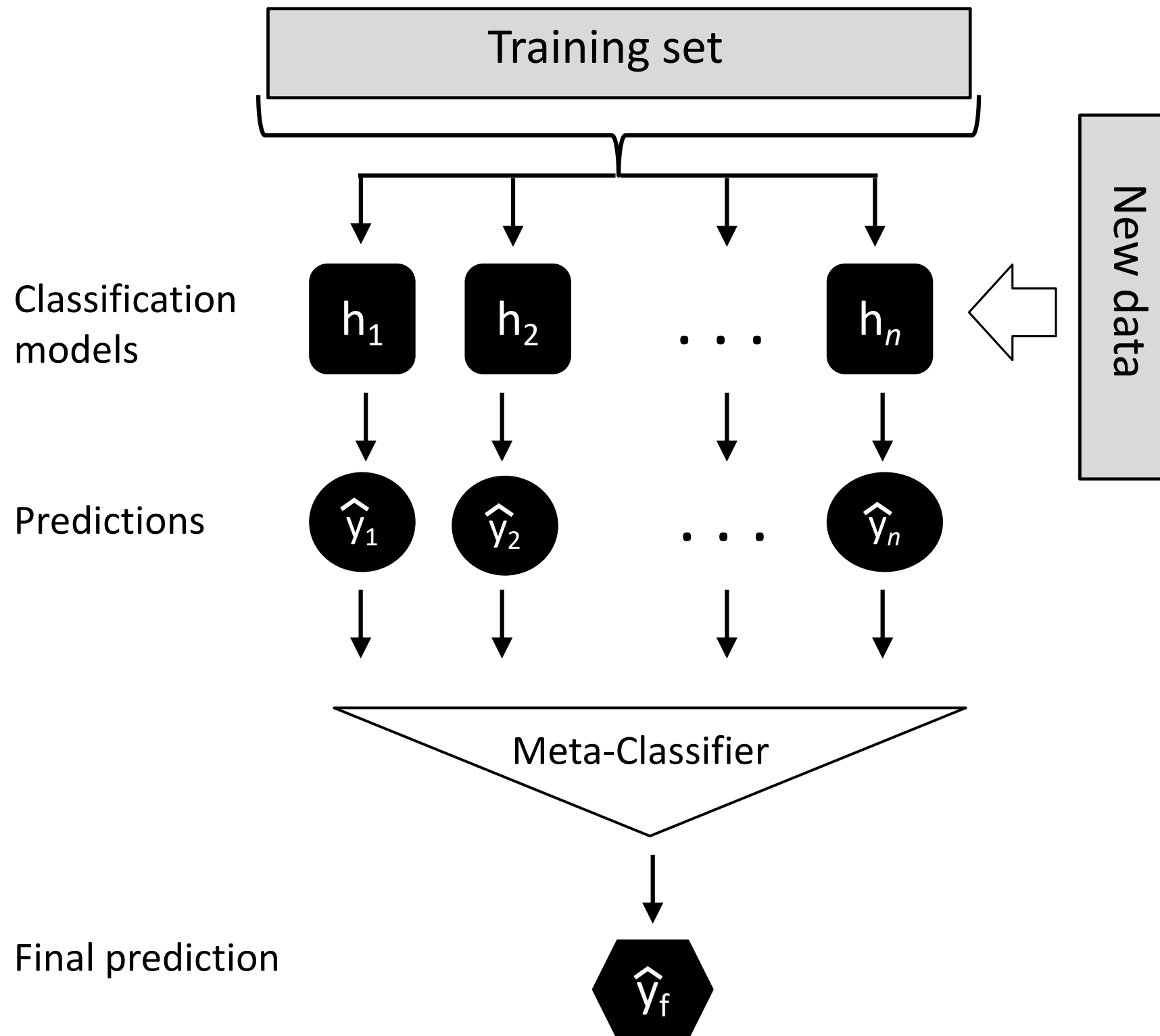
---

**Algorithm 19.7 Stacking**

---

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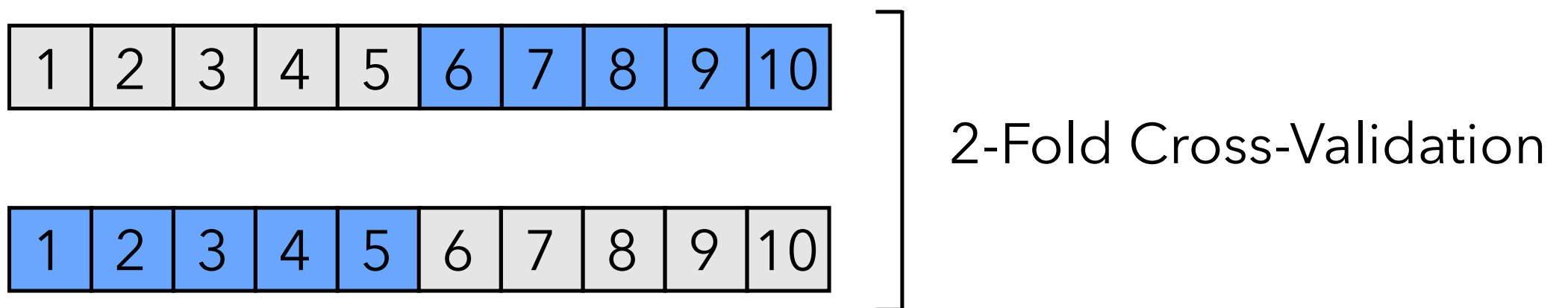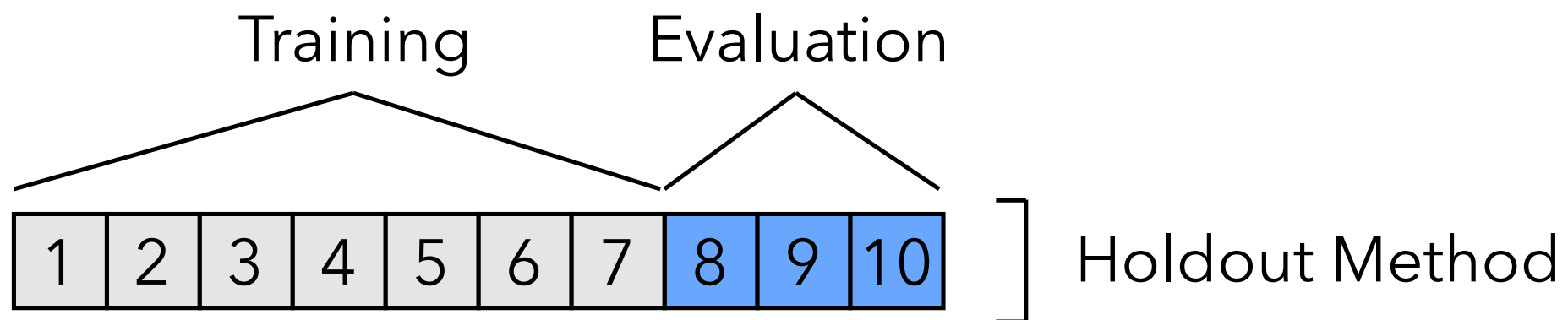**Output:** An ensemble classifier $H$

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 $\mathbf{x}_i' = \{h_1(\mathbf{x}_i), h_2(\mathbf{x}_i), \ldots, h_T(\mathbf{x}_i)\}$
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** $H(\mathbf{x}) = h'(h_1(\mathbf{x}), h_2(\mathbf{x}), \ldots, h_T(\mathbf{x}))$

---

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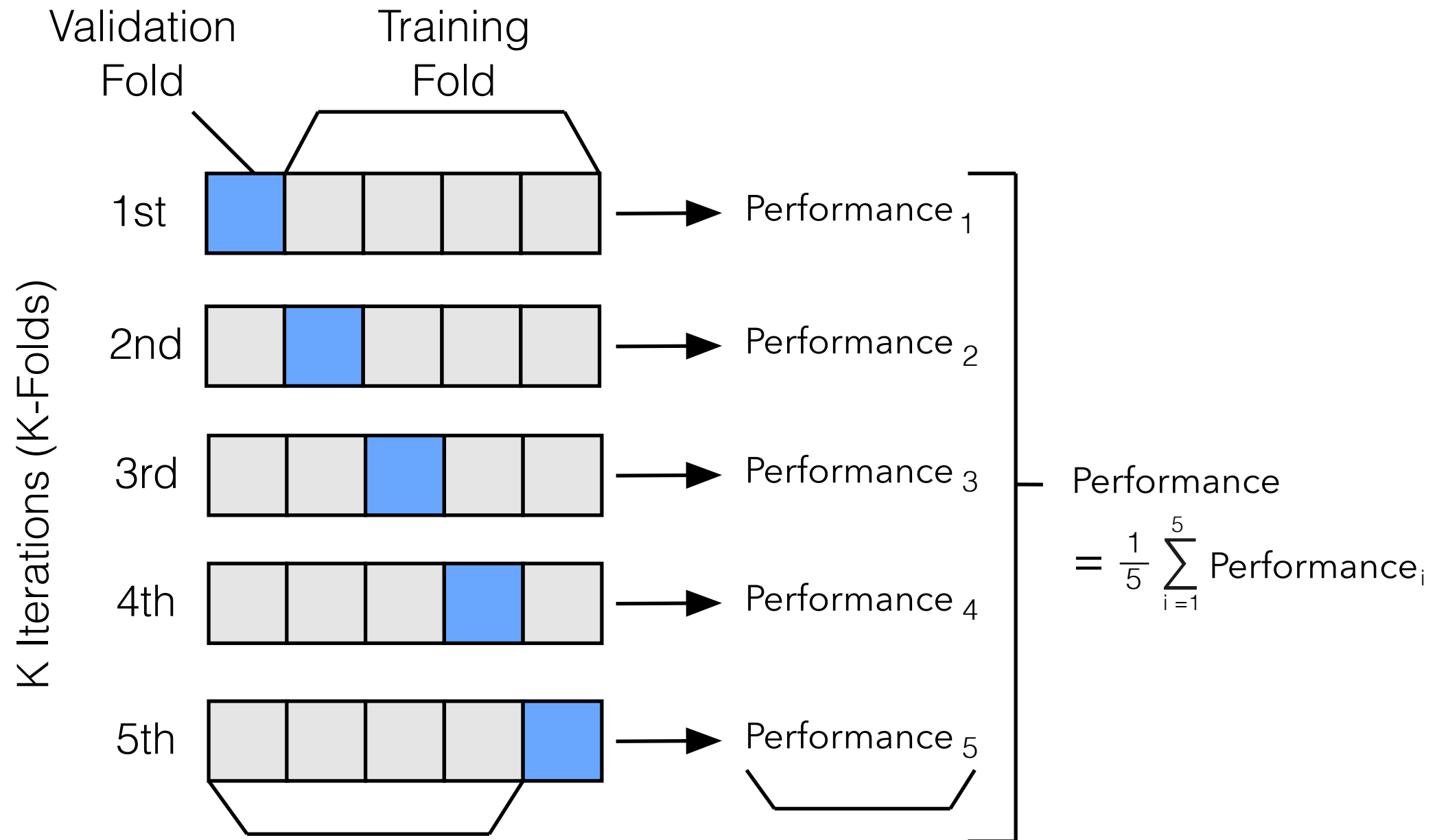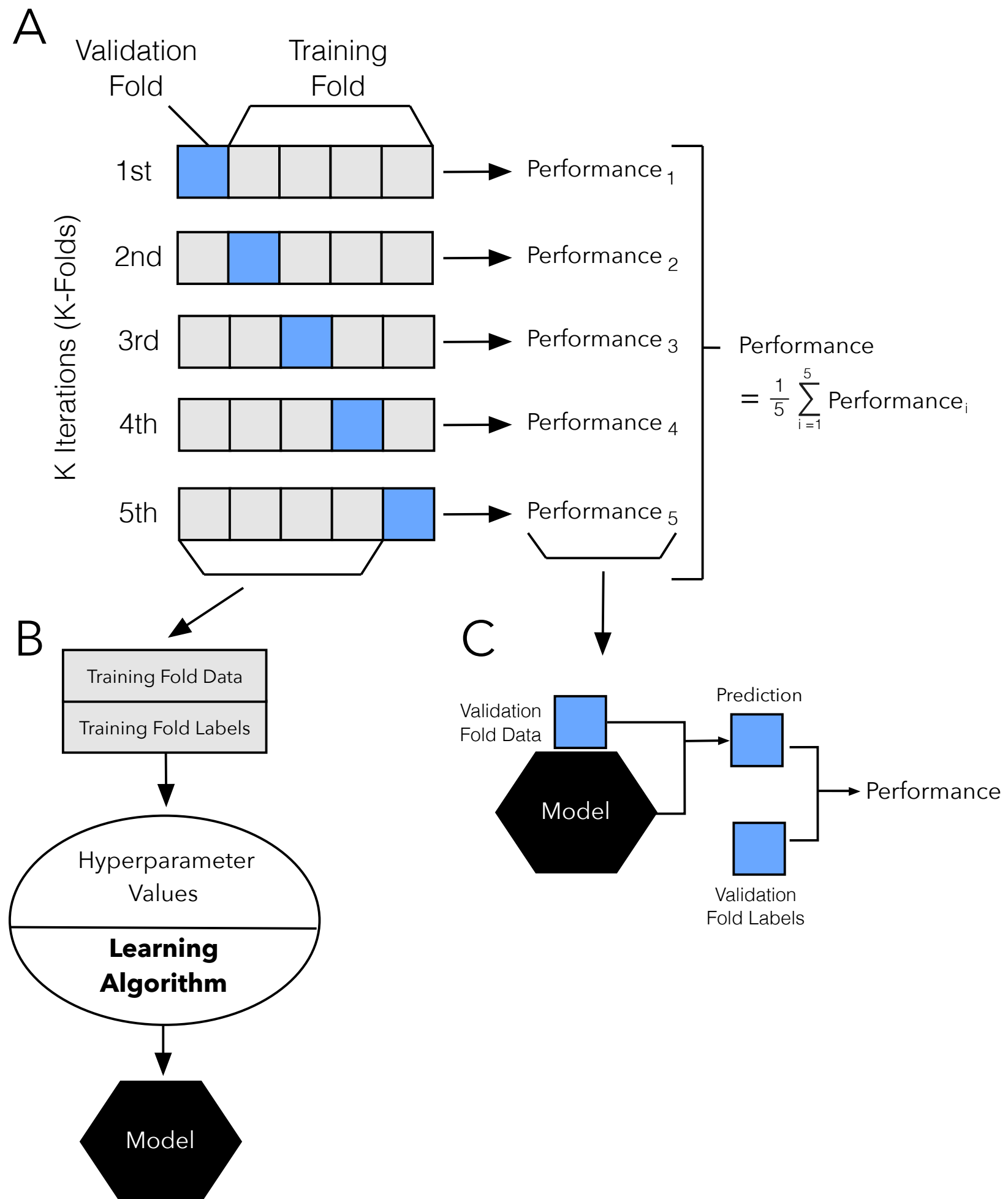
# Stacking Algorithm

# Cross-Validation

# *k*-fold Cross-Validation



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

**A**

Validation Fold | Training Fold

K Iterations (K-Folds)

1st → Performance$_1$

2nd → Performance$_2$

3rd → Performance$_3$

4th → Performance$_4$

5th → Performance$_5$

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

**B**

Training Fold Data

Training Fold Labels

Hyperparameter Values

**Learning Algorithm**

Model

**C**

Validation Fold Data

Model

Prediction

Validation Fold Labels

Performance

# Stacking Algorithm with Cross-Validation

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

---

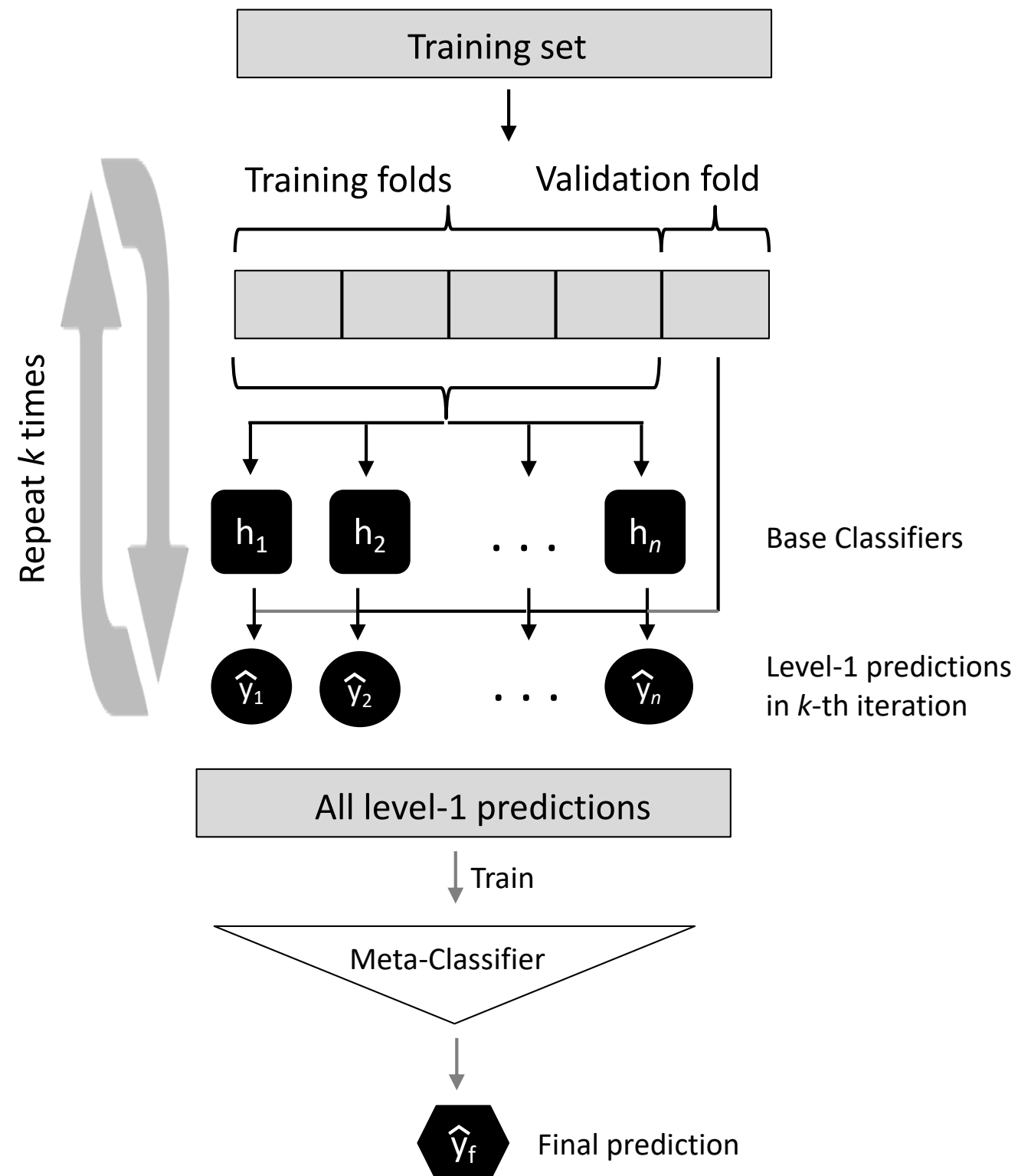**Algorithm 19.8 Stacking with $K$-fold Cross Validation**

---

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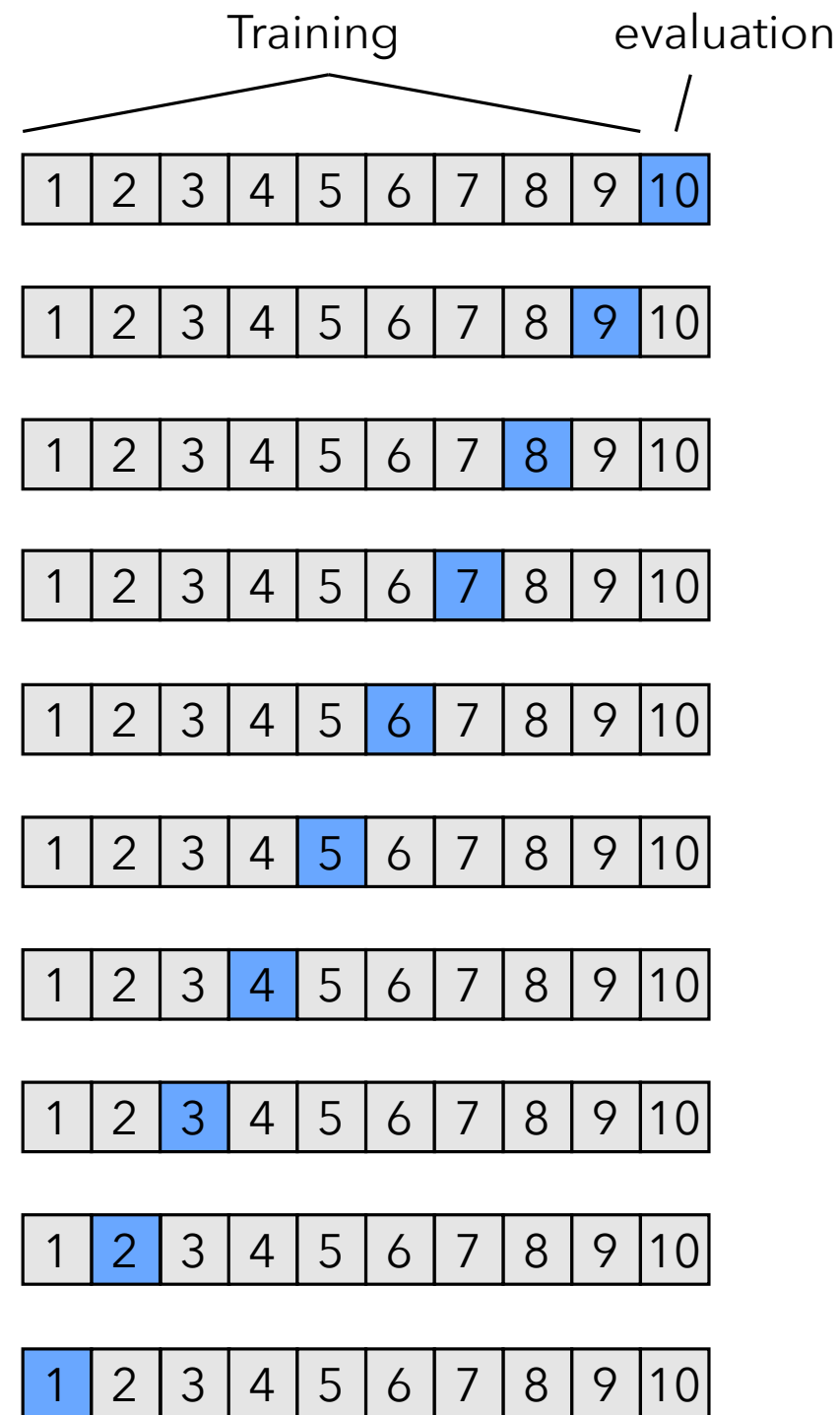**Output:** An ensemble classifier $H$

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

# Leave-One-Out CV

# Demos

http://rasbt.github.io/mlxtend/user_guide/classifier/EnsembleVoteClassifier/

http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.VotingClassifier.html

http://scikit-learn.org/stable/auto_examples/ensemble/
plot_bias_variance.html#sphx-glr-auto-examples-ensemble-plot-bias-variance-py

http://scikit-learn.org/stable/auto_examples/ensemble/
plot_adaboost_hastie_10_2.html#sphx-glr-auto-examples-ensemble-plot-adaboost-
hastie-10-2-py

https://scikit-learn.org/stable/modules/generated/
sklearn.ensemble.GradientBoostingClassifier.html

http://rasbt.github.io/mlxtend/user_guide/classifier/StackingClassifier/

http://rasbt.github.io/mlxtend/user_guide/classifier/StackingCVClassifier/

# Reading Assignments

Python Machine Learning, 2nd Ed., Ch07