


Ensemble Models

The idea



- The idea: construct a set of classifiers and combine them in certain way to classify new samples.
- The key to winning competitions 
- To make it success: more accurate than each individual classifier
 - Need “accurate” (i.e. **better than random guessing**) and “diverse” (i.e. **make errors on different data samples**) classifiers
- The question: why does the idea of ensemble work?

- The three fundamental reasons

- Statistical

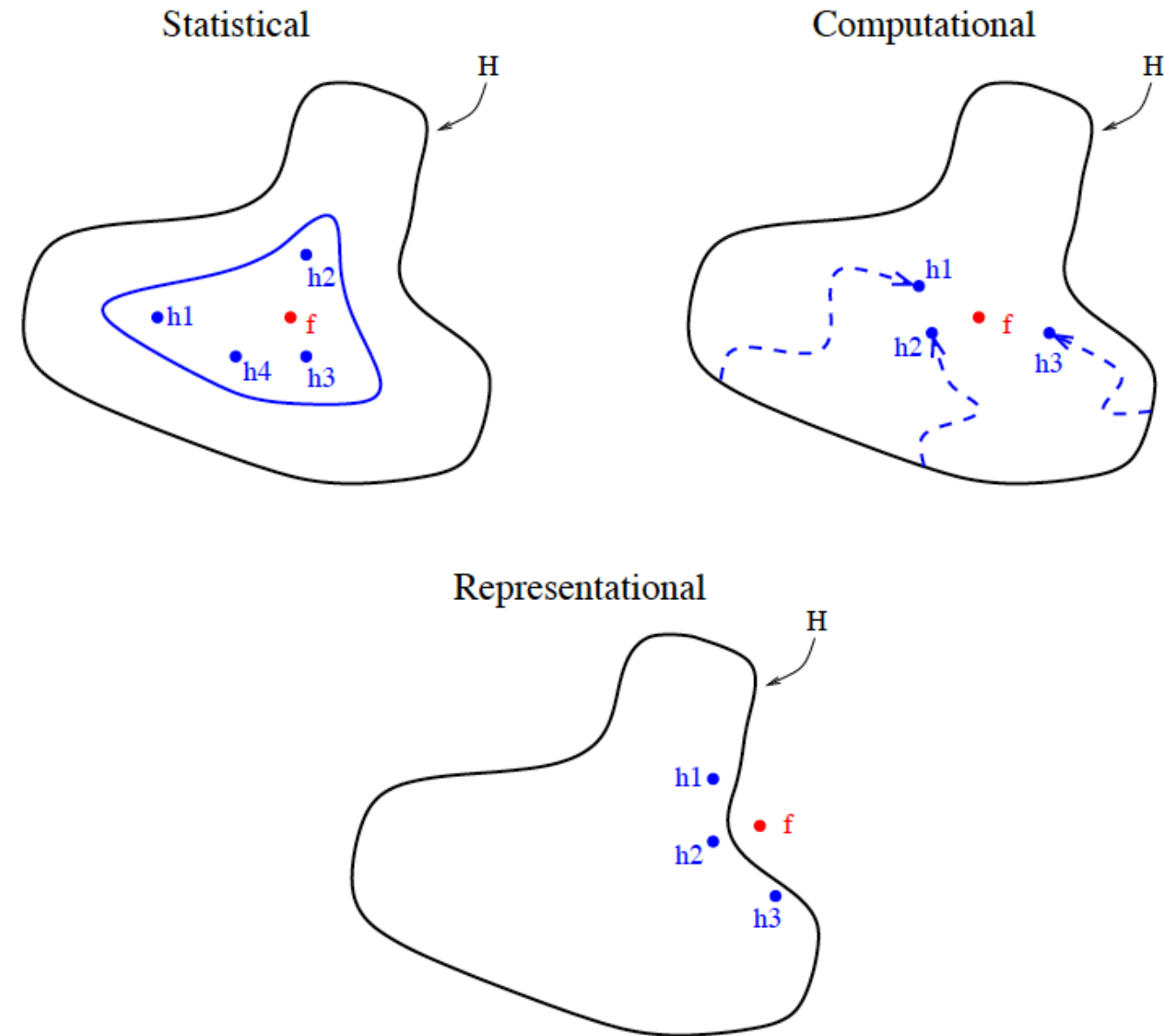
- Average out the votings to reduce the variance or bias

- Computational

- Combine the solutions obtained from local search

- Representational

- Possible to expand the space of the representation functions



Several Representatives

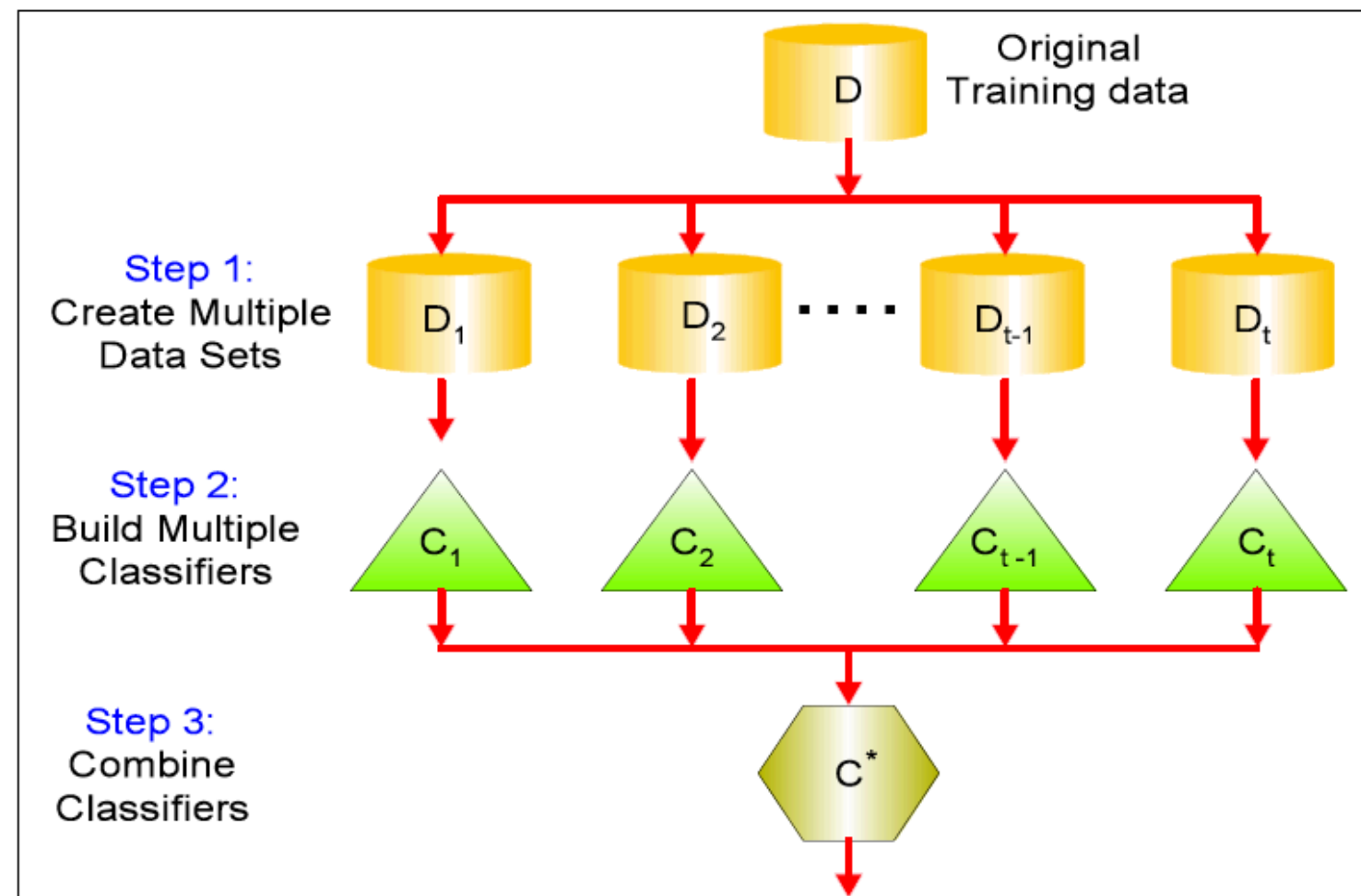


- General steps
 - 1. Producing a distribution of simple models on the subsets of the original data
 - 2. Combining the distribution into an “aggregated” model
- Three strategies
 - Bagging: decreasing variance
 - Boosting: decreasing bias
 - Stacking: improving the prediction power

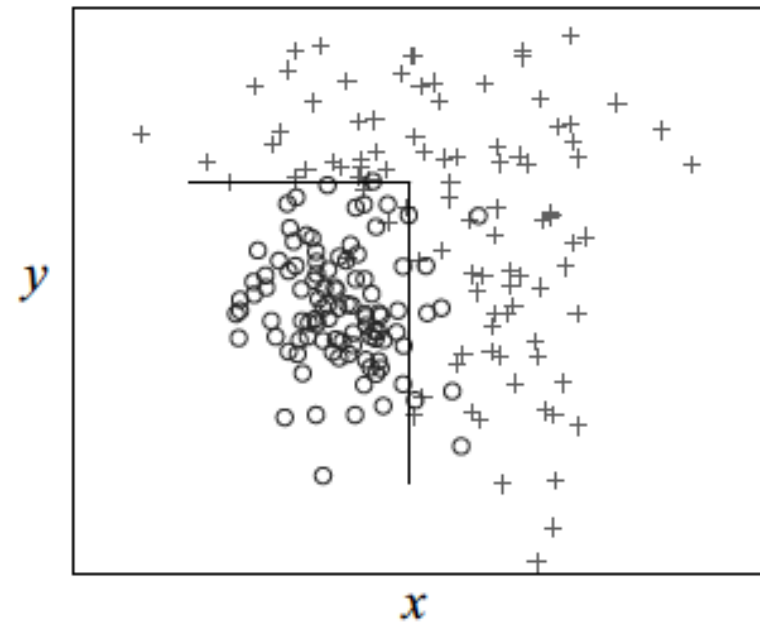
Bagging



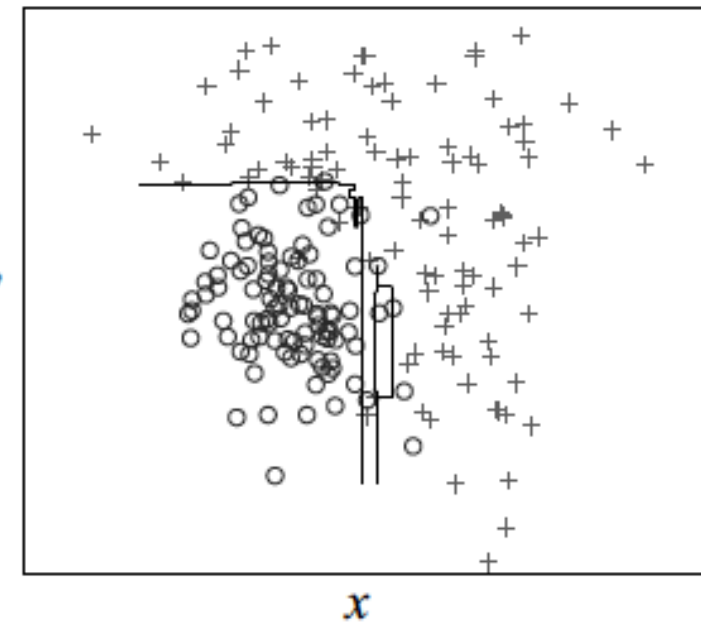
- Abbreviated for **bootstrap aggregating**
- Generating m new training sets by bootstrapping (i.e. sampling with replacement)
- Training each base classifier
- Combining them by averaging the output



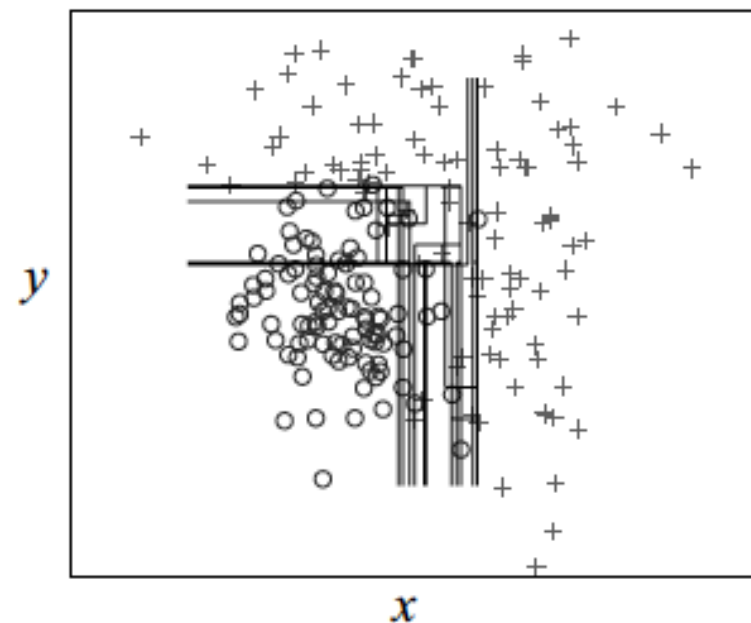
Bagging 10 decision trees



(a)



(b)



(c)

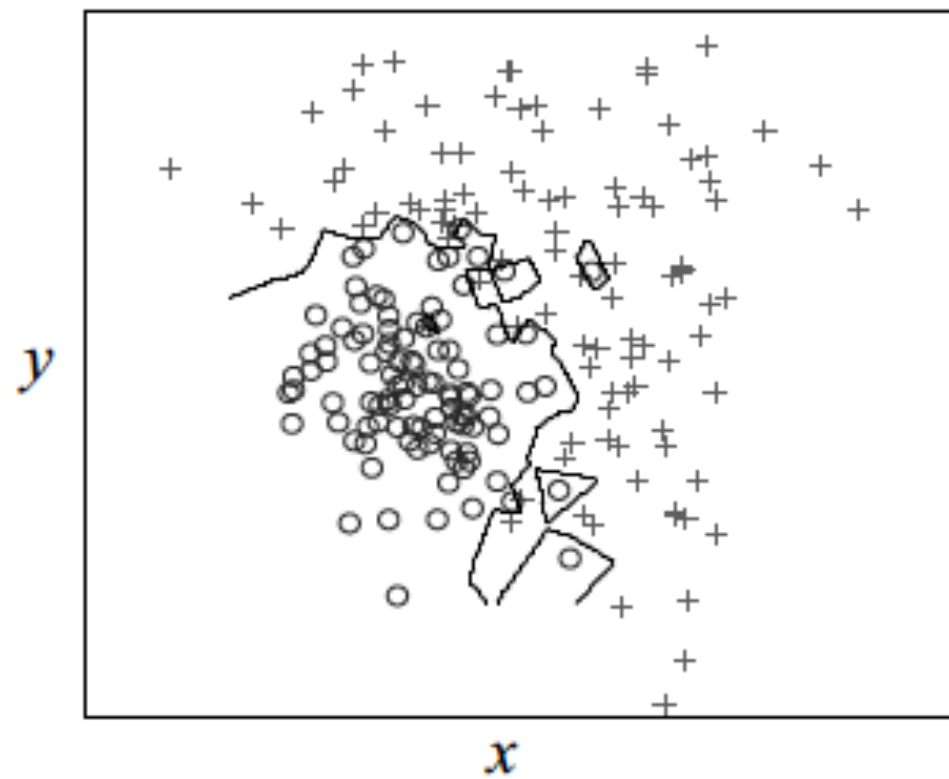
- Out-of-bag (oob) examples for estimating generalization error

$$H^{oob}(\mathbf{x}) = \arg \max_{y \in \mathcal{Y}} \sum_{t=1}^T \mathbb{I}(h_t(\mathbf{x}) = y) \cdot \mathbb{I}(\mathbf{x} \notin D_t)$$

$$err^{oob} = \frac{1}{|D|} \sum_{(\mathbf{x}, y) \in D} \mathbb{I}(H^{oob}(\mathbf{x}) \neq y)$$

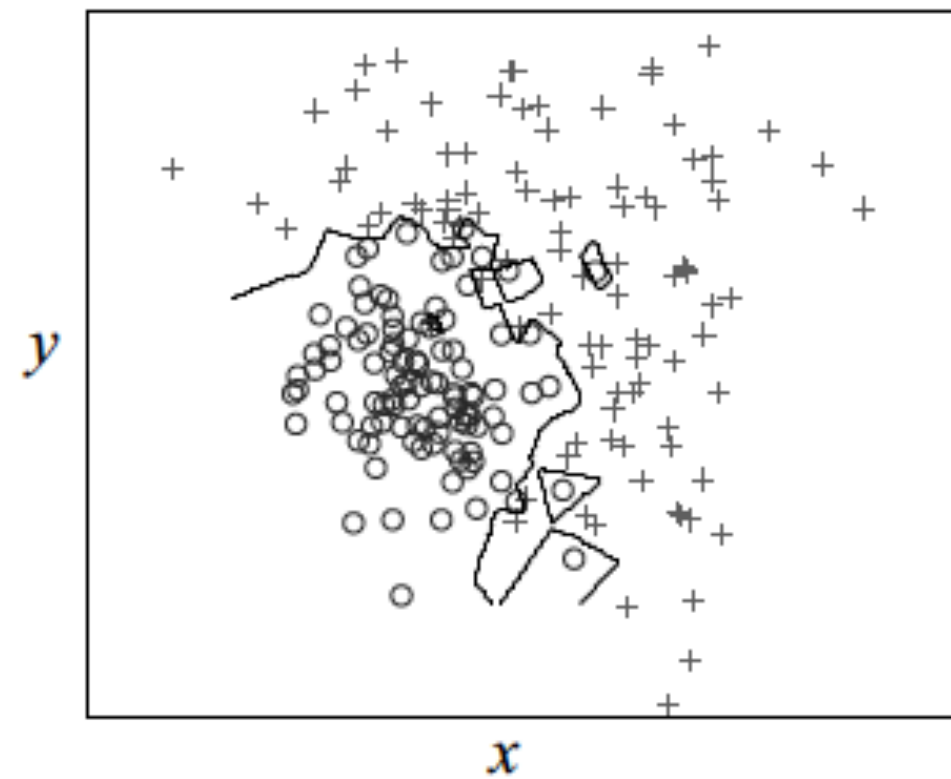
- Bagging requires diverse/independent learners, and so **stable learners do not help**.
- KNN is stable.
- Decision trees are unstable learners w.r.t. bootstrap subsets, particularly unpruned trees. 😎

Bagging 1-NN



(a)

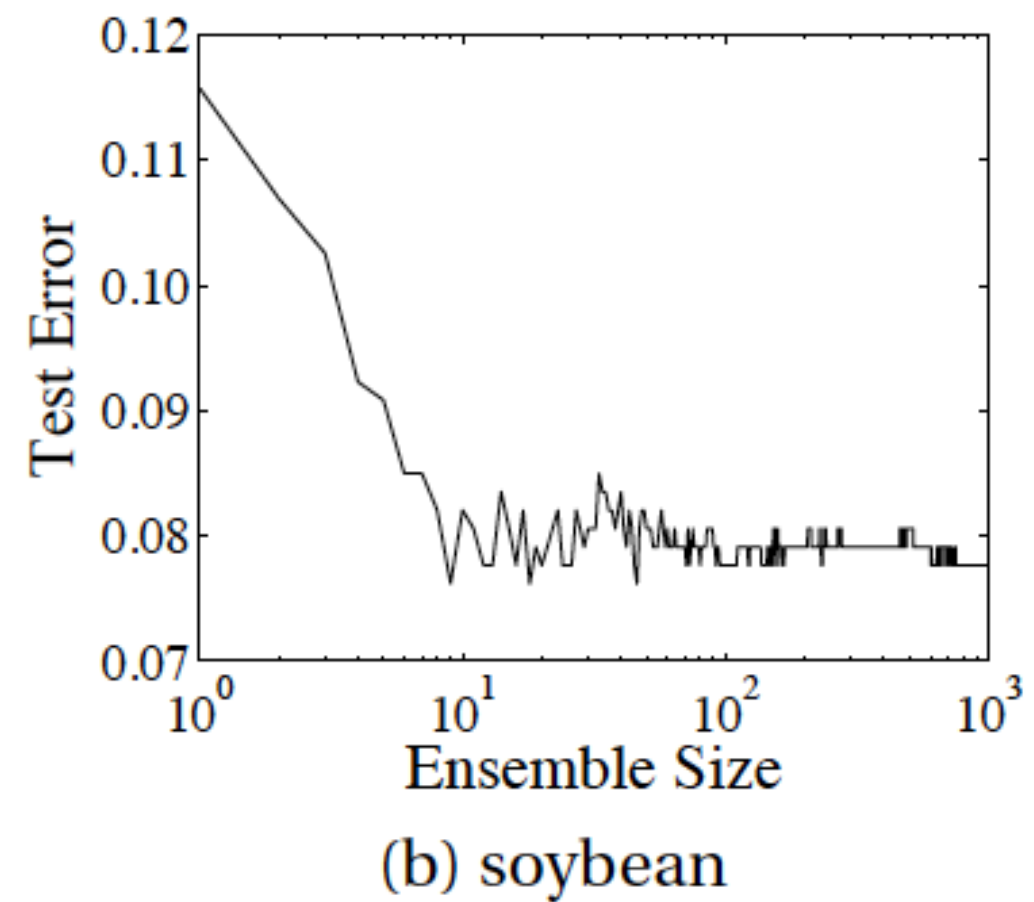
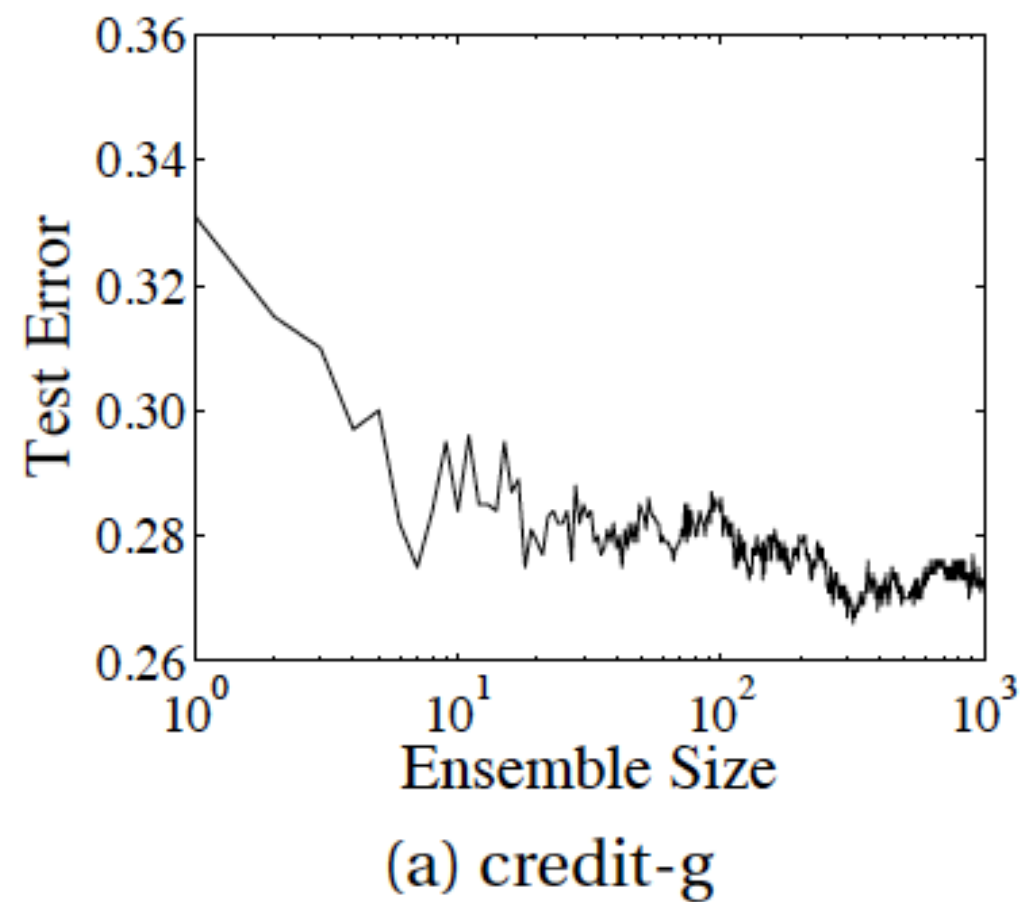
1-NN



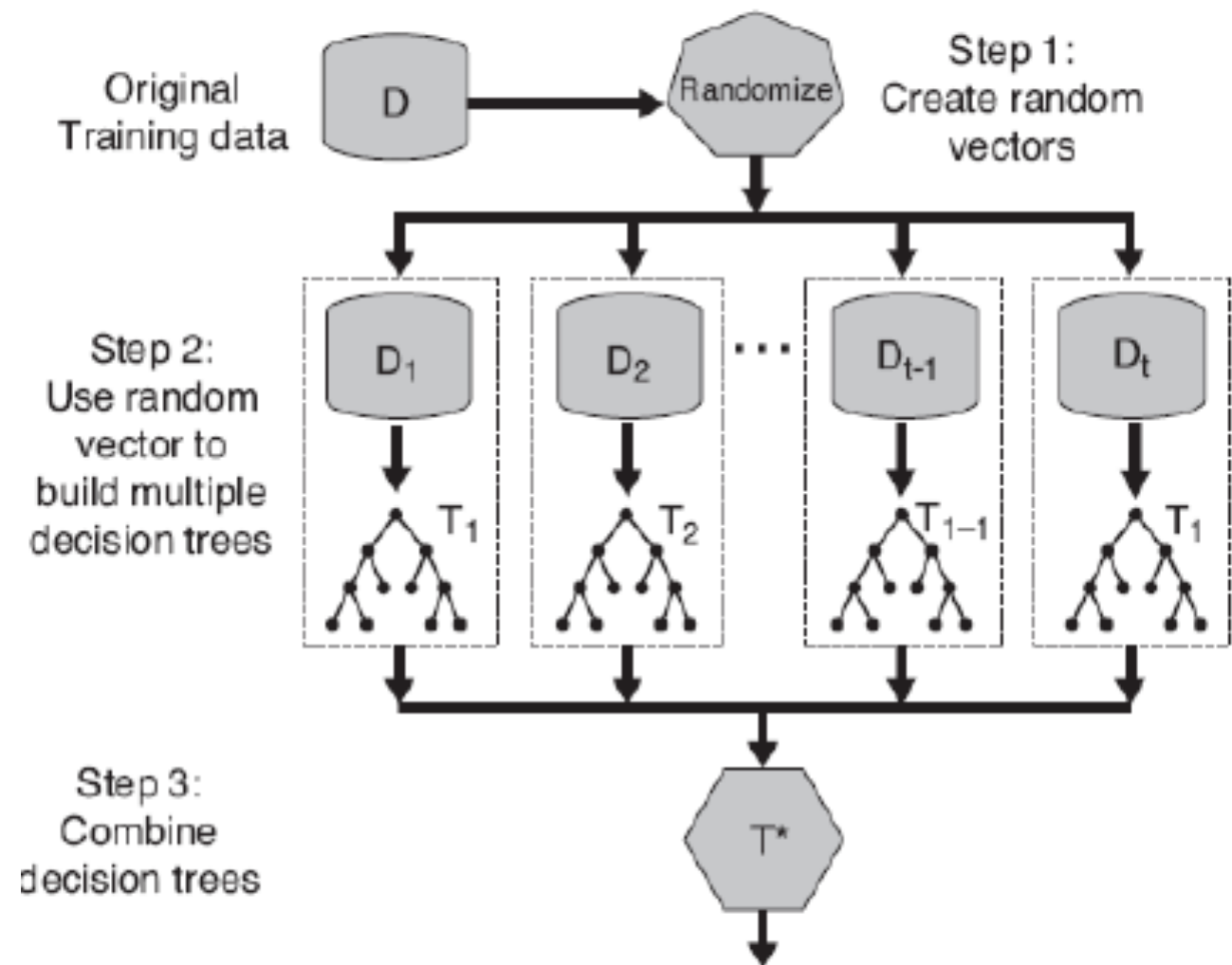
(b)

Bagging 1-NN

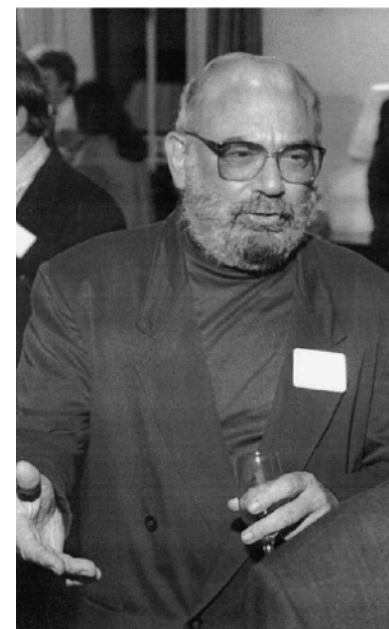
- With increasing ensemble size, the generalization converges.
- Why?



- An example: **random forest**
- Ensemble method designed for decision trees.
- Two sources of randomness.
 - Bagging: each tree is grown using a bootstrap sample of data
 - Random input vectors: at each node, best split is chosen from random m attributes instead of all attributes



Leo Breiman, 1928 - 2005



1954: PhD Berkeley (mathematics)
 1960 -1967: UCLA (mathematics)
 1969 -1982: Consultant
 1982 - 1993 Berkeley (statistics)
 1984 "Classification & Regression Trees"
 (with Friedman, Olshen, Stone)
 1996 "Bagging"
 2001 "Random Forests"

- Methods for growing the trees

Fix a $m \leq p$, at each node

- Method 1

- Choose d attributes randomly, compute their information gain, and choose the attribute with the largest gain to split.

- Method 2

- Compute a linear combination of L random attributes using weights generated from $[-1, 1]$ randomly.

- Method 3

- Compute the information gain of all D attributes. Select the top d attributes by information gain. Randomly select one of d attributes as the splitting node.

Random Forest based on method 1

1. For $b = 1$ to B :

- (a) Draw a **bootstrap sample** \mathbf{Z}^* of size N from the training data.
- (b) Grow a random-forest tree T_b to the bootstrapped data, by recursively repeating the following steps for each terminal node of the tree, until the minimum node size n_{min} is reached.
 - i. Select **m variables at random** from the p variables.
 - ii. Pick the best variable/split-point among the m .
 - iii. Split the node into two daughter nodes.

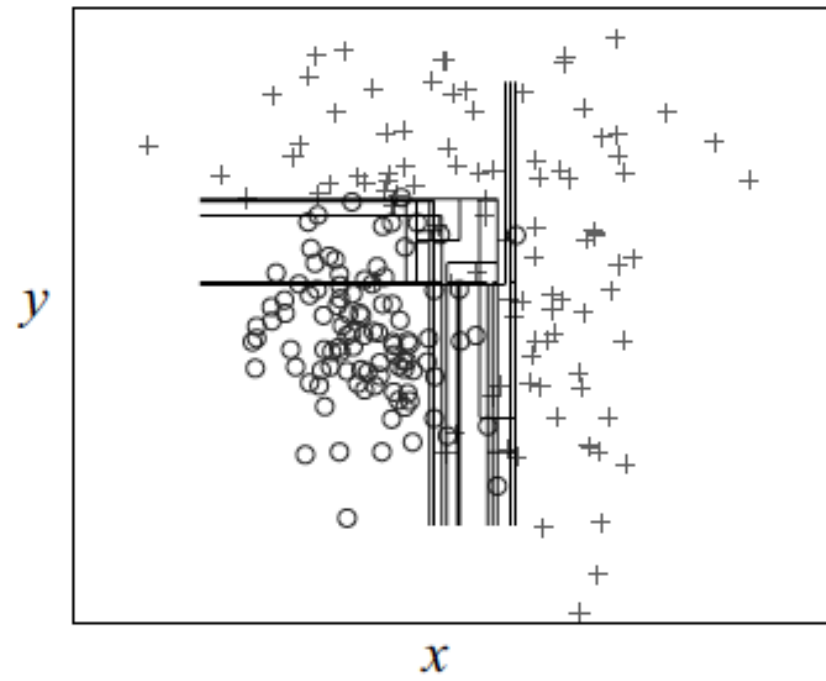
2. Output the ensemble of trees $\{T_b\}_1^B$.

To make a prediction at a new point x :

Regression: $\hat{f}_{\text{rf}}^B(x) = \frac{1}{B} \sum_{b=1}^B T_b(x)$.

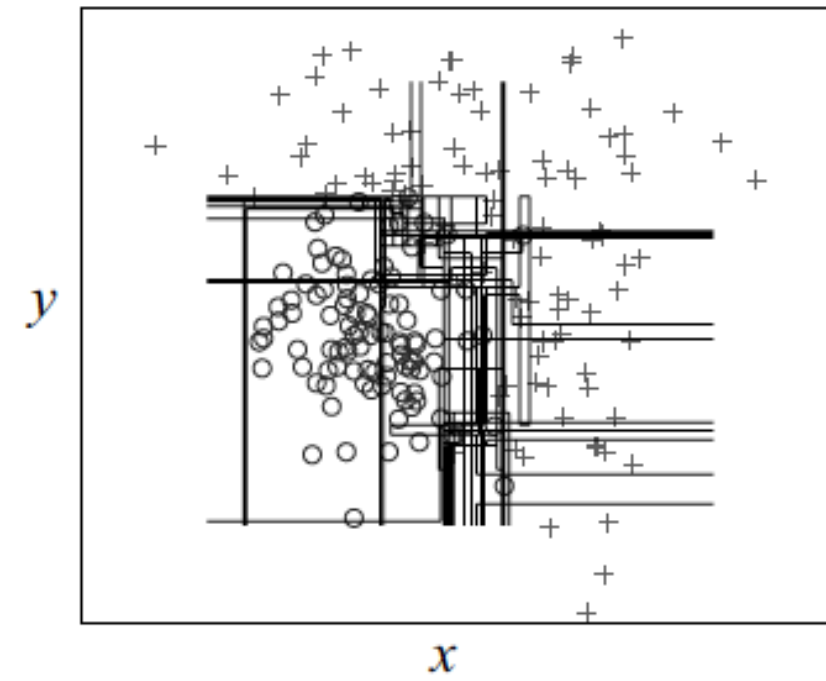
Classification: Let $\hat{C}_b(x)$ be the class prediction of the b th random-forest tree. Then $\hat{C}_{\text{rf}}^B(x) = \text{majority vote } \{\hat{C}_b(x)\}_1^B$.

base classifiers of bagging

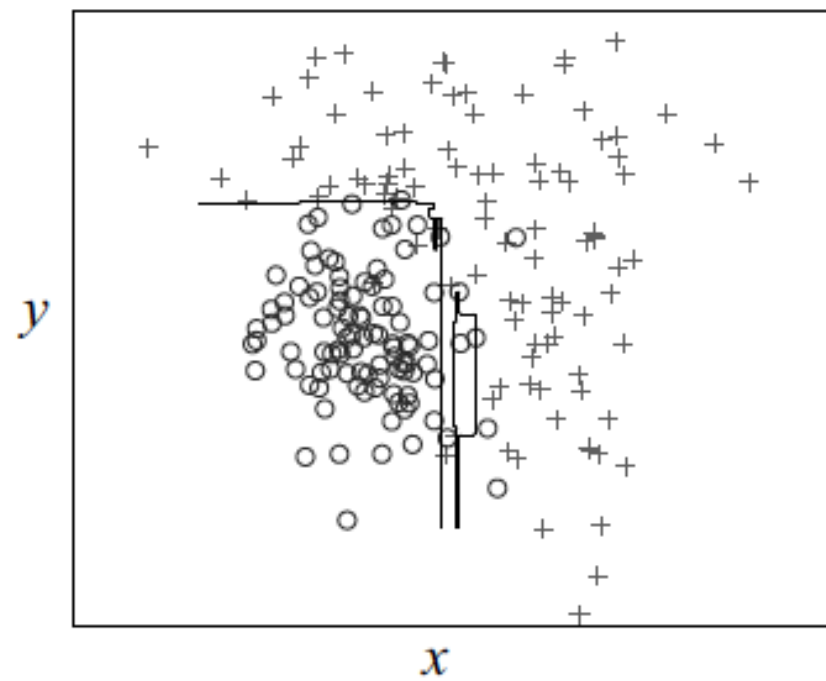


(a)

base classifiers of RF

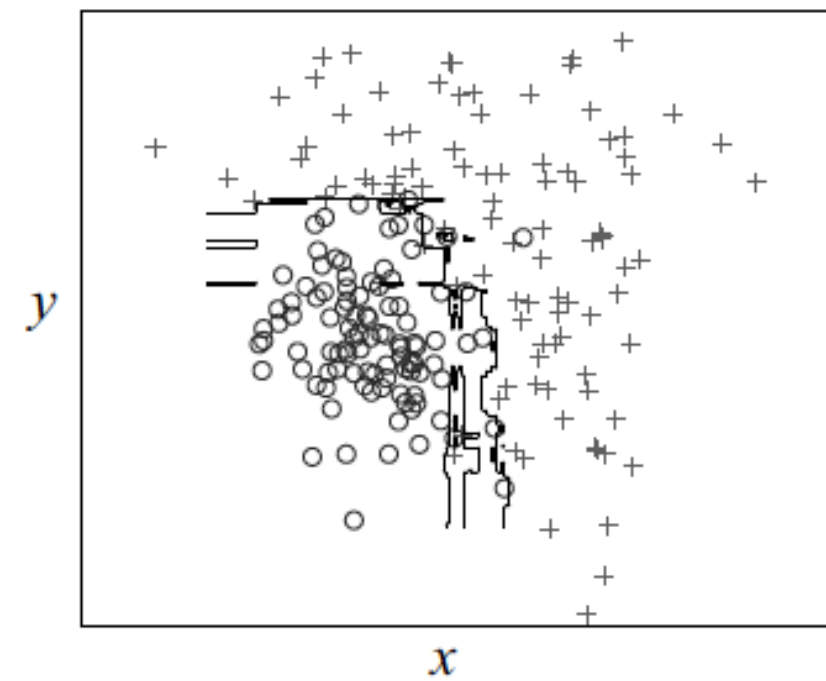


(b)



(c)

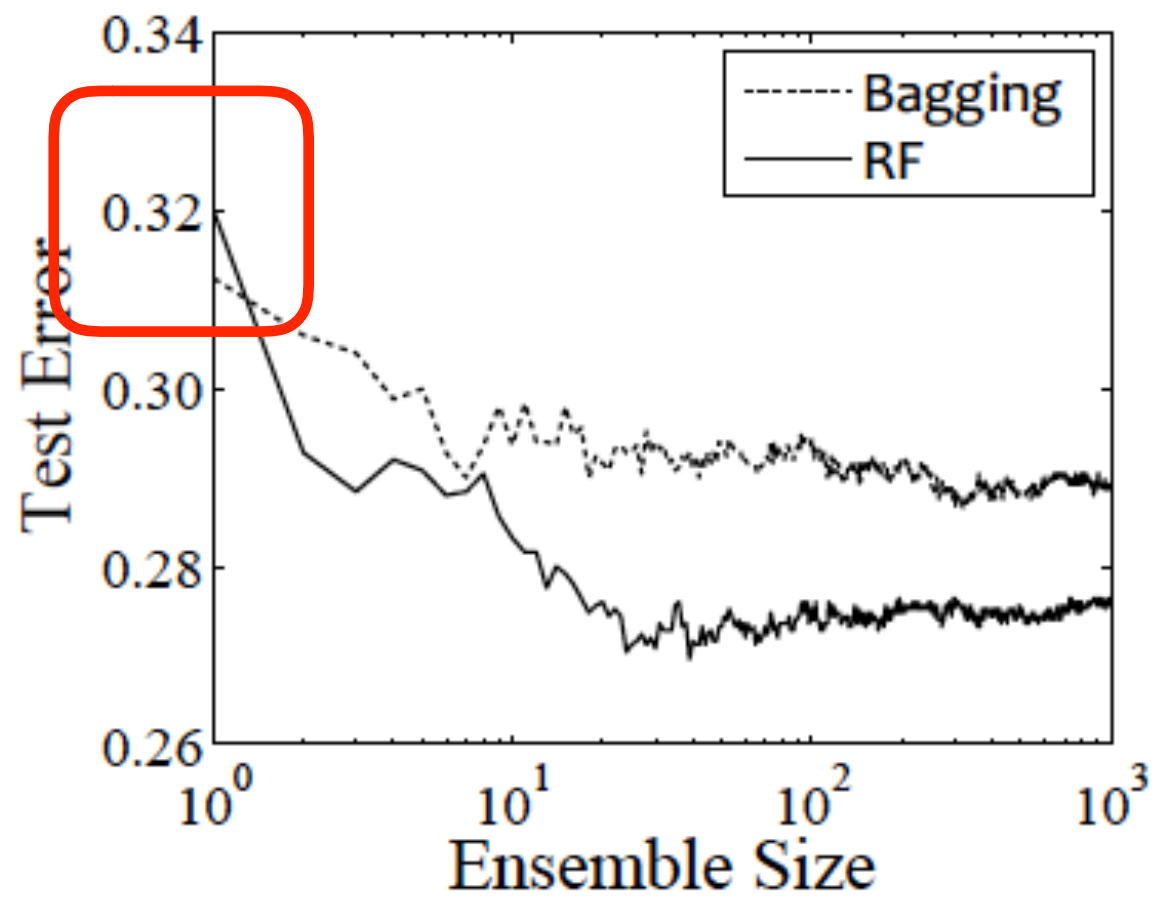
bagging



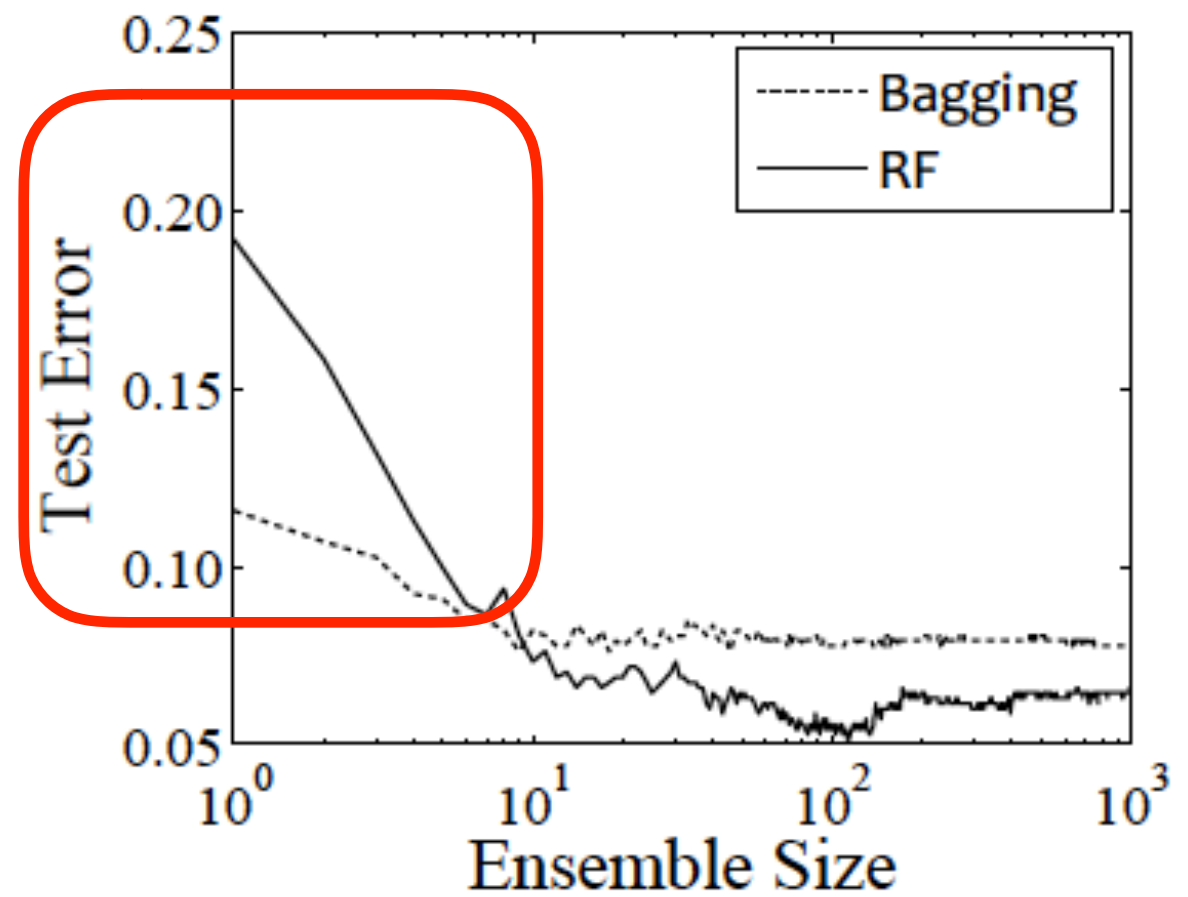
(d)

RF

Convergence property



(a) credit-g



(b) soybean

Some Remarks

- Bagging works through reducing the prediction variance.
- Some simple analysis

$$P(h_i(x) \neq f(x)) = \epsilon$$

H makes errors when at least half of the base classifiers make error

$$H(x) = \text{sign} \left(\sum_{i=1}^T h_i(x) \right)$$

Assuming independence

$$P(H(x) \neq f(x)) = \sum_{k=0}^{\lfloor T/2 \rfloor} \binom{T}{k} (1-\epsilon)^k \epsilon^{T-k} \leq \exp \left(-\frac{1}{2} T (2\epsilon - 1)^2 \right)$$

Hoeffding inequality

- Bagging is amenable to parallel processing.

**Bagging can reduce variance,
can we decrease bias?**

Boosting

- A family of algorithms that **turn weak learners to strong ones**.
- The idea
 - **Sequentially** train a set of classifiers and combine them
 - Late learners focus more on mistakes of the earlier learners

Input: Sample distribution \mathcal{D} ;
Base learning algorithm \mathcal{L} ;
Number of learning rounds T .

Process:

1. $\mathcal{D}_1 = \mathcal{D}$. % Initialize distribution
2. **for** $t = 1, \dots, T$:
3. $h_t = \mathcal{L}(\mathcal{D}_t)$; % Train a weak learner from distribution \mathcal{D}_t
4. $\epsilon_t = P_{\mathbf{x} \sim \mathcal{D}_t}(h_t(\mathbf{x}) \neq f(\mathbf{x}))$; % Evaluate the error of h_t
5. $\mathcal{D}_{t+1} = \text{Adjust_Distribution}(\mathcal{D}_t, \epsilon_t)$
6. **end**

Output: $H(\mathbf{x}) = \text{Combine_Outputs}(\{h_1(\mathbf{x}), \dots, h_t(\mathbf{x})\})$

- An instantiation of boosting: **AdaBoost** (Freund and Schapire 1997, Friedman et.al 2000)

Input: Data set $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$;
 Base learning algorithm \mathcal{L} ;
 Number of learning rounds T .

Process:

1. $\mathcal{D}_1(x) = 1/m$. % Initialize the weight distribution
2. **for** $t = 1, \dots, T$:
3. $h_t = \mathcal{L}(D, \mathcal{D}_t)$; % Train a classifier h_t from D under distribution \mathcal{D}_t
4. $\epsilon_t = P_{\mathbf{x} \sim \mathcal{D}_t}(h_t(\mathbf{x}) \neq f(\mathbf{x}))$; % Evaluate the error of h_t
5. **if** $\epsilon_t > 0.5$ **then break**
6. $\alpha_t = \frac{1}{2} \ln \left(\frac{1-\epsilon_t}{\epsilon_t} \right)$; % Determine the weight of h_t
7. $\mathcal{D}_{t+1}(x) = \frac{\mathcal{D}_t(x)}{Z_t} \times \begin{cases} \exp(-\alpha_t) & \text{if } h_t(x) = f(x) \\ \exp(\alpha_t) & \text{if } h_t(x) \neq f(x) \end{cases}$ ↗ **Re-weighting**
 $= \frac{\mathcal{D}_t(x) \exp(-\alpha_t f(x) h_t(x))}{Z_t}$ % Update the distribution, where
 % Z_t is a normalization factor which
 % enables \mathcal{D}_{t+1} to be a distribution
8. **end**

Output: $H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$

- Minimizing exponential loss

$$\ell_{\text{exp}}(h \mid \mathcal{D}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})h(\mathbf{x})}] \quad H(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x})$$

When the H is optimal, the following should hold.

$$\begin{aligned} \frac{\partial e^{-f(\mathbf{x})H(\mathbf{x})}}{\partial H(\mathbf{x})} &= -f(\mathbf{x})e^{-f(\mathbf{x})H(\mathbf{x})} \\ &= -e^{-H(\mathbf{x})}P(f(\mathbf{x}) = 1 \mid \mathbf{x}) + e^{H(\mathbf{x})}P(f(\mathbf{x}) = -1 \mid \mathbf{x}) \\ &= 0. \end{aligned}$$



$$H(\mathbf{x}) = \frac{1}{2} \ln \frac{P(f(\mathbf{x}) = 1 \mid \mathbf{x})}{P(f(\mathbf{x}) = -1 \mid \mathbf{x})}$$

$$\begin{aligned}
 \text{sign}(H(\mathbf{x})) &= \text{sign} \left(\frac{1}{2} \ln \frac{P(f(\mathbf{x}) = 1 \mid \mathbf{x})}{P(f(\mathbf{x}) = -1 \mid \mathbf{x})} \right) \\
 &= \begin{cases} 1, & P(f(\mathbf{x}) = 1 \mid \mathbf{x}) > P(f(\mathbf{x}) = -1 \mid \mathbf{x}) \\ -1, & P(f(\mathbf{x}) = 1 \mid \mathbf{x}) < P(f(\mathbf{x}) = -1 \mid \mathbf{x}) \end{cases} \\
 &= \arg \max_{y \in \{-1, 1\}} P(f(\mathbf{x}) = y \mid \mathbf{x}), \quad \text{Bayesian error rate}
 \end{aligned}$$

Construct from individual learner:

$$\begin{aligned}
 \ell_{\text{exp}}(\alpha_t h_t \mid \mathcal{D}_t) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [e^{-f(\mathbf{x})\alpha_t h_t(\mathbf{x})}] \\
 &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [e^{-\alpha_t} \mathbb{I}(f(\mathbf{x}) = h_t(\mathbf{x})) + e^{\alpha_t} \mathbb{I}(f(\mathbf{x}) \neq h_t(\mathbf{x}))] \\
 &= e^{-\alpha_t} P_{\mathbf{x} \sim \mathcal{D}_t}(f(\mathbf{x}) = h_t(\mathbf{x})) + e^{\alpha_t} P_{\mathbf{x} \sim \mathcal{D}_t}(f(\mathbf{x}) \neq h_t(\mathbf{x})) \\
 &= e^{-\alpha_t} (1 - \epsilon_t) + e^{\alpha_t} \epsilon_t,
 \end{aligned}$$

Optimal weight:

$$\frac{\partial \ell_{\text{exp}}(\alpha_t h_t \mid \mathcal{D}_t)}{\partial \alpha_t} = -e^{-\alpha_t} (1 - \epsilon_t) + e^{\alpha_t} \epsilon_t = 0 \longrightarrow \alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

Minimizing the loss:

$$\begin{aligned}\ell_{\text{exp}}(H_{t-1} + h_t \mid \mathcal{D}) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})(H_{t-1}(\mathbf{x}) + h_t(\mathbf{x}))}] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} e^{-f(\mathbf{x})h_t(\mathbf{x})}]\end{aligned}$$

Taylor expansion:

$$\begin{aligned}\ell_{\text{exp}}(H_{t-1} + h_t \mid \mathcal{D}) &\approx \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} \left(1 - f(\mathbf{x})h_t(\mathbf{x}) + \frac{f(\mathbf{x})^2 h_t(\mathbf{x})^2}{2} \right) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} \left(1 - f(\mathbf{x})h_t(\mathbf{x}) + \frac{1}{2} \right) \right]\end{aligned}$$

$$\begin{aligned}h_t(\mathbf{x}) &= \arg \min_h \ell_{\text{exp}}(H_{t-1} + h \mid \mathcal{D}) \\ &= \arg \min_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} \left(1 - f(\mathbf{x})h(\mathbf{x}) + \frac{1}{2} \right) \right] \\ &= \arg \max_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} f(\mathbf{x})h(\mathbf{x}) \right] \\ &= \arg \max_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\frac{e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]} f(\mathbf{x})h(\mathbf{x}) \right],\end{aligned}$$

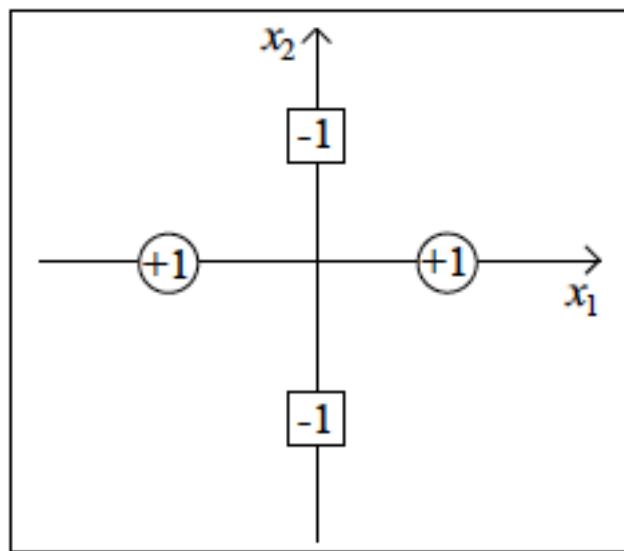
$$\begin{aligned}
 h_t(\mathbf{x}) &= \arg \max_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\frac{e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]} f(\mathbf{x})h(\mathbf{x}) \right] \\
 &= \arg \max_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [f(\mathbf{x})h(\mathbf{x})] .
 \end{aligned}$$

$$f(\mathbf{x})h_t(\mathbf{x}) = 1 - 2\mathbb{I}(f(\mathbf{x}) \neq h_t(\mathbf{x}))$$

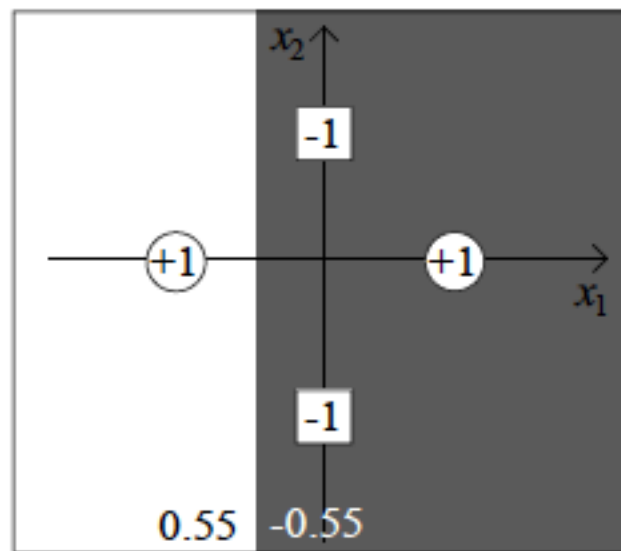
$$h_t(\mathbf{x}) = \arg \min_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [\mathbb{I}(f(\mathbf{x}) \neq h(\mathbf{x}))]$$

$$\begin{aligned}
 \mathcal{D}_{t+1}(\mathbf{x}) &= \frac{\mathcal{D}(\mathbf{x})e^{-f(\mathbf{x})H_t(\mathbf{x})}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_t(\mathbf{x})}]} \\
 &= \frac{\mathcal{D}(\mathbf{x})e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}e^{-f(\mathbf{x})\alpha_t h_t(\mathbf{x})}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_t(\mathbf{x})}]} \\
 &= \mathcal{D}_t(\mathbf{x}) \cdot e^{-f(\mathbf{x})\alpha_t h_t(\mathbf{x})} \frac{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_t(\mathbf{x})}]}
 \end{aligned}$$

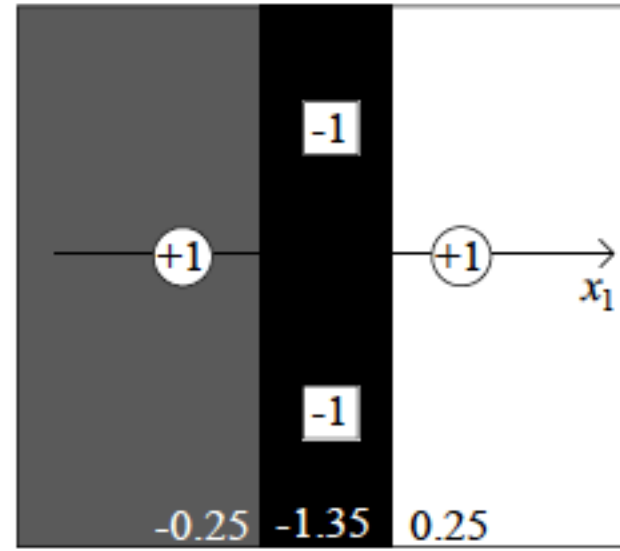
- A simple example:



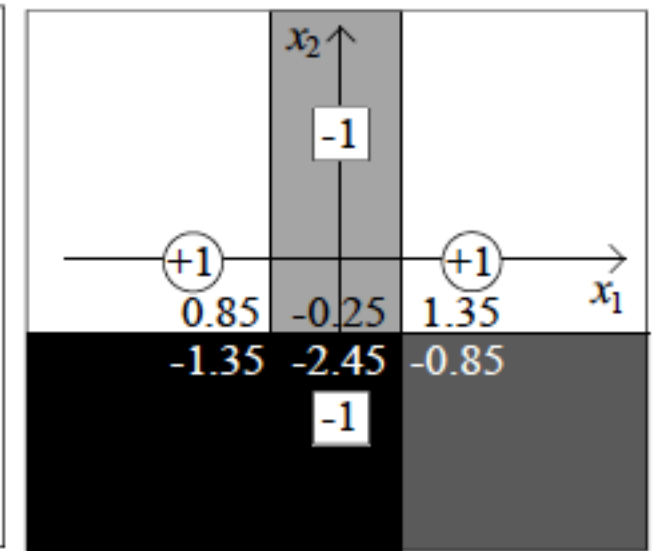
(a) The XOR data



(b) 1st round



(c) 2nd round



(d) 3rd round

8 base learners

$$h_1(x) = \begin{cases} +1, & \text{if } (x_1 > -0.5) \\ -1, & \text{otherwise} \end{cases} \quad h_2(x) = \begin{cases} -1, & \text{if } (x_1 > -0.5) \\ +1, & \text{otherwise} \end{cases}$$

$$h_3(x) = \begin{cases} +1, & \text{if } (x_1 > +0.5) \\ -1, & \text{otherwise} \end{cases} \quad h_4(x) = \begin{cases} -1, & \text{if } (x_1 > +0.5) \\ +1, & \text{otherwise} \end{cases}$$

$$h_5(x) = \begin{cases} +1, & \text{if } (x_2 > -0.5) \\ -1, & \text{otherwise} \end{cases} \quad h_6(x) = \begin{cases} -1, & \text{if } (x_2 > -0.5) \\ +1, & \text{otherwise} \end{cases}$$

$$h_7(x) = \begin{cases} +1, & \text{if } (x_2 > +0.5) \\ -1, & \text{otherwise} \end{cases} \quad h_8(x) = \begin{cases} -1, & \text{if } (x_2 > +0.5) \\ +1, & \text{otherwise} \end{cases}$$

where x_1 and x_2 are the values of x at the first and the second dimension, respectively.

Rethinking AdaBoost

- Boosting is an additive model constructed sequentially.
- The exponential loss is an surrogate loss, an upper bound of 0-1 loss.

$$P(f(\mathbf{x}) = 1 \mid \mathbf{x}) = \frac{e^{H(\mathbf{x})}}{e^{H(\mathbf{x})} + e^{-H(\mathbf{x})}}$$

Log loss: $\ell_{\log}(h \mid \mathcal{D}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\ln \left(1 + e^{-2f(\mathbf{x})h(\mathbf{x})} \right) \right]$

- Optimization log loss with gradient descent with base regression models: LogitBoost
- Other variants: L2Boost, or more generally Gradient Boosting (functional gradient descent)
 - XGBoost Library



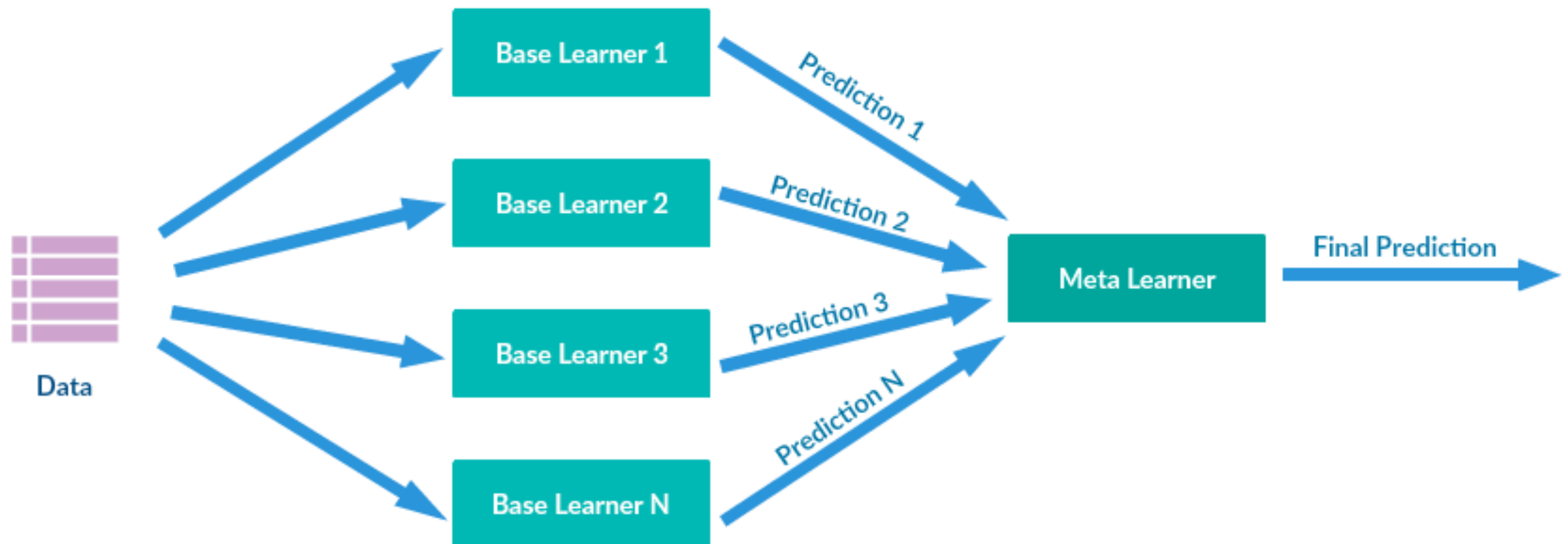
Process:

- Output:** $H(x) = \text{sign} \left(\sum_{t=1}^T h_t(x) \right)$

LogitBoost

Stacking

- Meta-learning methods: learning to combine (Wolpert 1992, Smyth and Wolpert 1998)
- A generalization of many ensemble methods
- The idea: the output of the base learners (the first-level learners) are used as the input features of the meta-learners (the second-level learners).



Input: Data set $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$;
First-level learning algorithms $\mathcal{L}_1, \dots, \mathcal{L}_T$;
Second-level learning algorithm \mathcal{L} .

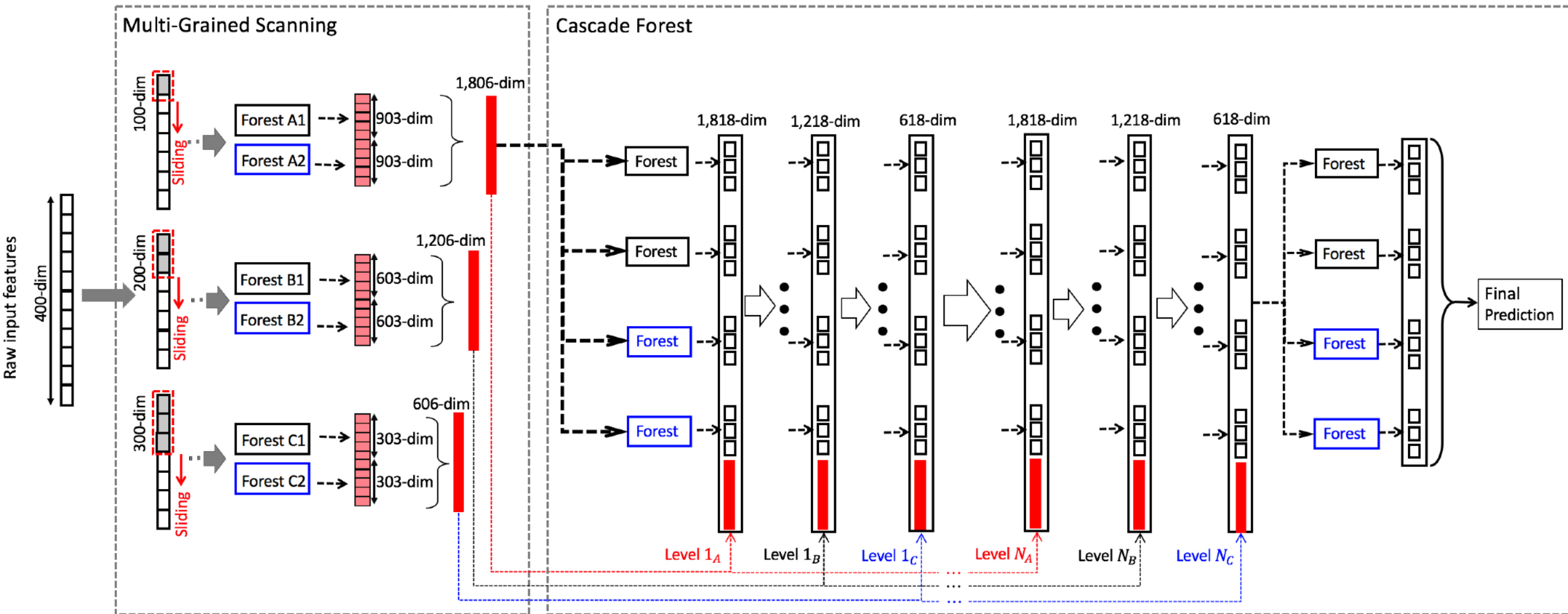
Process:


```
1.  for  $t = 1, \dots, T$ :    % Train a first-level learner by applying the
2.       $h_t = \mathcal{L}_t(D)$ ;    % first-level learning algorithm  $\mathcal{L}_t$ 
3.  end
4.   $D' = \emptyset$ ;          % Generate a new data set
5.  for  $i = 1, \dots, m$ :
6.      for  $t = 1, \dots, T$ :
7.           $z_{it} = h_t(x_i)$ ;
8.      end
9.       $D' = D' \cup ((z_{i1}, \dots, z_{iT}), y_i)$ ;
10. end
11.  $h' = \mathcal{L}(D')$ ;      % Train the second-level learner  $h'$  by applying
                          % the second-level learning algorithm  $\mathcal{L}$  to the
                          % new data set  $D'$ .
```

Output: $H(x) = h'(h_1(x), \dots, h_T(x))$

A general procedure for stacking


- A recent example: Deep Forest (Zhou and Feng 2017)



- 
- A decorative blue line graphic that starts with a jagged, wave-like pattern on the left and then transitions into a straight horizontal line extending across the top of the slide.
- More topics of Boosting (for further study)
 - How to extend to multi-class cases?
 - Noise sensitivity (why sensitive?)
 - Generalization bound

Summary



- Ensemble models: combining different learners in various ways
 - Bagging: variance reduction, parallel
 - Random forest
 - Boosting: bias reduction, sequentially
 - AdaBoost
 - Stacking: learn how to combine
- Key: how to construct diverse base learners by introducing various types of randomness
- Typically an important trick to win competitions 
- Generalization error analysis

Exercises



- Derive how LogitBoost works.
- Readings
 - Generalization error analysis of AdaBoost
 - Freund, Yoav, and Robert E. Schapire. "A decision-theoretic generalization of on-line learning and an application to boosting." *Journal of computer and system sciences* 55.1 (1997): 119-139.
 - Schapire, Robert E., et al. "Boosting the margin: A new explanation for the effectiveness of voting methods." *The annals of statistics* 26.5 (1998): 1651-1686.