

Ensemble Learning

Songpan Yang

School of Mathematics Sun Yat-sen University

Dec, 3, 2020

Table of Contents

- 1 Ensemble Learning
- 2 Random Forest
- 3 Incremental Learning

Table of Contents

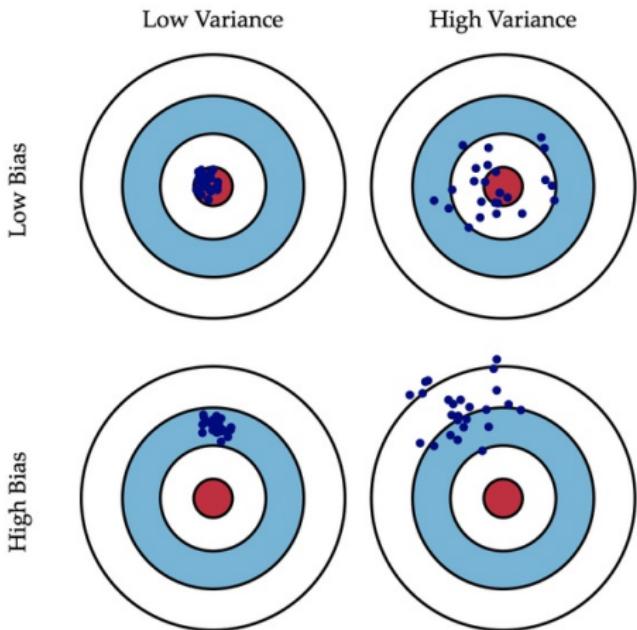
- 1 Ensemble Learning
- 2 Random Forest
- 3 Incremental Learning

Why do we study Ensemble Learning?

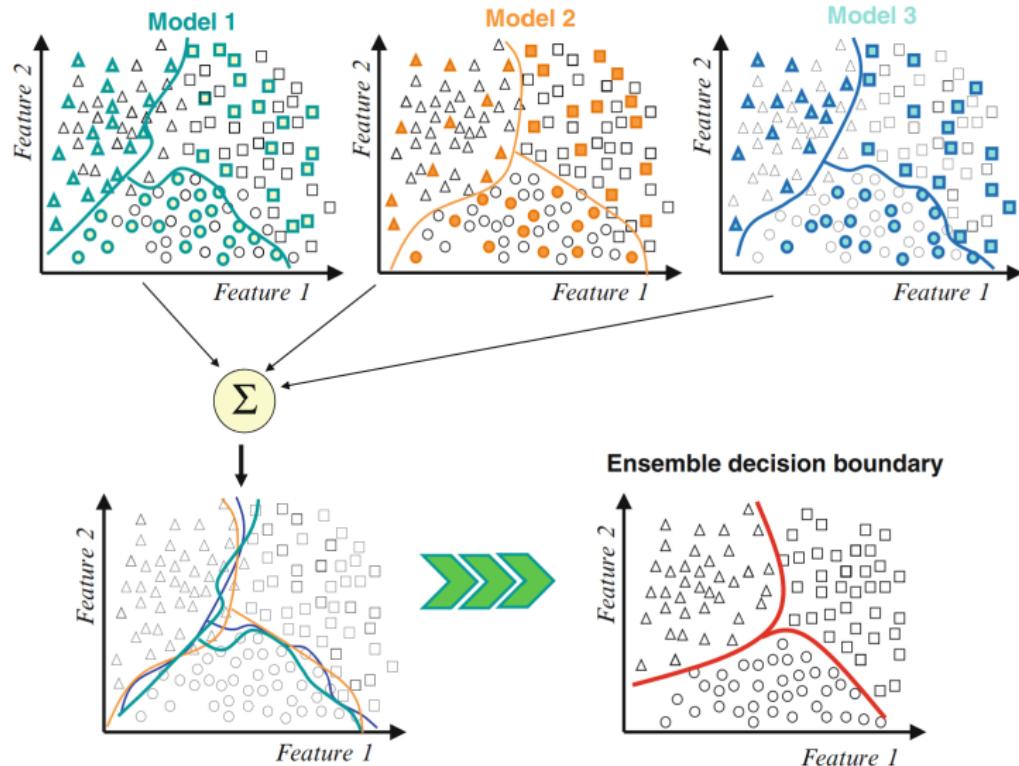
- We often consult others before making a decision
- An expert (or oracle) VS no expert
- The accuracy of each decision (classifier) has a nonzero variability.
- Two components of any classification error: bias and variance
- The goal of ensemble systems:
 1. create several classifiers with relatively fixed (or similar) bias
 2. combining their outputs, to reduce the variance

Bias-variance tradeoff

- The true model (An expert or oracle):
 $y = f_{true}(x)$
- Suppose we have a model (classifier):
 $y = f(x) + \epsilon$
- $E(y - f(x)) = \sigma_\epsilon^2 + \text{Bias}^2(f) + \text{Var}(f)$



An example



The process of Ensemble Learning

1. Data sampling/selection
2. Training Member Classifiers: bagging, boosting, stacking
3. Combining Ensemble Members

Data sampling/selection

- Making different errors on any given sample
- Diversity in ensembles.
- The most common approach:
using different subsets of the training data
 1. bootstrapped replicas of the training data
 2. sampling from a distribution that favors previously misclassified samples
 3. use different subsets of the available features to train each classifier
- Other less common approaches:
 1. using different parameters of the base classifier
 2. using different base classifiers as the ensemble members

Popular Ensemble-Based Algorithms

- Three popular algorithms: Bagging, Boosting, Stacking
- Denote the classifier f_i , the weight of classifier w_i , the correlation coefficient of each classifier ρ
- The Full model be the linear combination of basic classifiers of Bagging and Boosting

$$\begin{aligned} E(F) &= E\left(\sum_i^m w_i f_i\right) = \sum_i^m w_i E(f_i) = w \sum_i^m E(f_i) \\ Var(F) &= Var\left(\sum_i^m w_i f_i\right) = Cor\left(\sum_i^m w_i f_i, \sum_i^m w_i f_i\right) \\ &= \sum_i^m w_i^2 Var(f_i) + \sum_i^m \sum_{j \neq i}^m 2\rho w_i w_j \sqrt{Var(f_i) Var(f_j)} \\ &= m^2 w^2 \sigma^2 \rho + mw^2 \sigma^2 (1 - \rho) \end{aligned}$$

bagging and boosting

- For bagging:

$$E(F) = w \sum_i^m E(f_i) = \frac{1}{m} m\mu$$

$$Var(F) = \sigma^2\rho + \frac{\sigma^2(1-\rho)}{m}$$

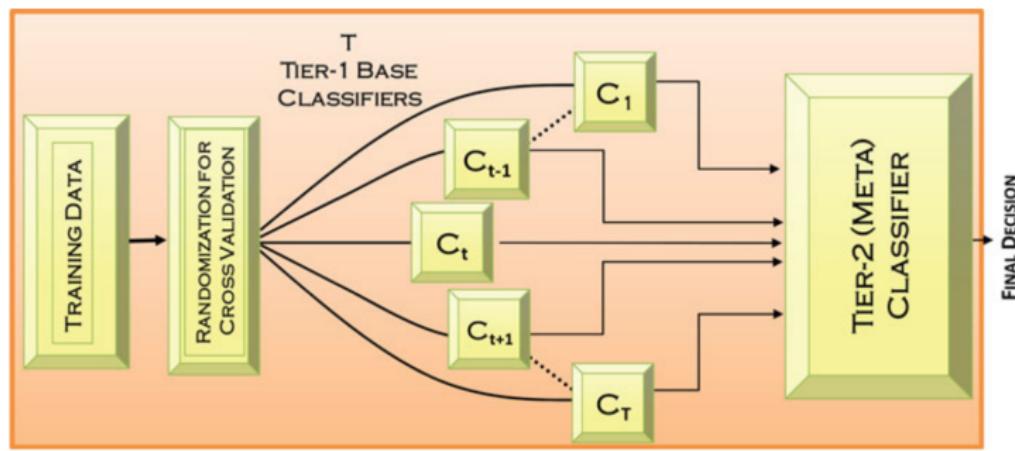
- For boosting:

$$E(F) = w \sum_i^m E(f_i)$$

$$Var(F) = m^2 w^2 \sigma^2$$

Stacking

- A separate classifier, trained on the outputs of the ensemble members
- Suppose we first train T classifiers base on a cross-validation partition of the training data
- The outputs of these classifiers on their pseudo-training blocks constitute the training data



Combining Ensemble Members

- Combine the individual classifiers: depends on the type of classifiers
- Combining Class Labels:

$$d_{t,c} \in \{0, 1\}, t = 1, \dots, T, c = 1, \dots, C$$

- Majority voting: unanimous voting, simple majority, plurality voting

$$\sum_{t=1}^T d_{t,c^*} = \max_c \sum_{t=1}^T d_{t,c}$$

- Weighted Majority Voting:

$$\sum_{t=1}^T w_t d_{t,c^*} = \max_c \sum_{t=1}^T w_t d_{t,c}$$

- Borda Count: rank order the classifier outputs

Combining Ensemble Members

- How do we assign the weights?
- A plausible and commonly used strategy is to use the performance of a classifier on a separate validation (or even training) dataset, as an estimate of that classifier's generalization performance (AdaBoost)

Combining Continuous Outputs

- An estimate of the posterior probability for that class

$$P(\omega_c | \mathbf{x}) \approx \tilde{g}_c(\mathbf{x}) = \frac{e^{g_c(\mathbf{x})}}{\sum_{i=1}^C e^{g_i(\mathbf{x})}} \Rightarrow \sum_{i=1}^C \tilde{g}_i(\mathbf{x}) = 1$$

- Kuncheva's decision profile matrix $DP(\mathbf{x})$, let $d_{t,c} \in [0, 1]$ be the $P(\omega_c | \mathbf{x})$ of the t classifier

Support from all classifiers $h_1 \dots h_T$
for class ω_c – one of the C classes.

$$DP(\mathbf{x}) = \begin{bmatrix} d_{1,1}(\mathbf{x}) & \cdots & d_{1,c}(\mathbf{x}) & \cdots & d_{1,C}(\mathbf{x}) \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ d_{t,1}(\mathbf{x}) & \cdots & d_{t,c}(\mathbf{x}) & \cdots & d_{t,C}(\mathbf{x}) \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ d_{T,1}(\mathbf{x}) & \cdots & d_{T,c}(\mathbf{x}) & \cdots & d_{T,C}(\mathbf{x}) \end{bmatrix}$$

Support given by classifier h_t
to each of the classes

Combining Continuous Outputs

- Algebraic Combiners

$$\mu_c(\mathbf{x}) = F [d_{1,c}(\mathbf{x}), \dots, d_{T,c}(\mathbf{x})]$$

- F is one of the following combination functions

- mean rules: $\mu_c(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T d_{t,c}(\mathbf{x})$
- weighted Average: $\mu_c(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T w_t d_{t,c}(\mathbf{x})$
- Trimmed mean: Sometimes classifiers may erroneously give unusually low or high support to a particular class
- Minimum/Maximum/Median Rule:

$$\mu_c(\mathbf{x}) = \min_{t=1, \dots, T} \{d_{t,c}(\mathbf{x})\}$$

$$\mu_c(\mathbf{x}) = \max_{t=1, \dots, T} \{d_{t,c}(\mathbf{x})\}$$

$$\mu_c(\mathbf{x}) = \text{median}_{t=1, \dots, T} \{d_{t,c}(\mathbf{x})\}$$

Combining Continuous Outputs

5. Product Rule: $\mu_c(\mathbf{x}) = \frac{1}{T} \prod_{t=1}^T d_{t,c}(\mathbf{x})$
6. Generalized Mean: $\mu_c(\mathbf{x}) = \left(\frac{1}{T} \sum_{t=1}^T (d_{t,c}(\mathbf{x}))^\alpha \right)^{1/\alpha}$
7. Decision Template: $DT_c = 1/N_c \sum_{X_c \in \omega_c} DP(X_c)$

$$\mu_c(x) = S(DP(x), DT_c), c = 1, \dots, C$$

7. S is usually a squared Euclidean distance

$$\mu_c(\mathbf{x}) = 1 - \frac{1}{T \times C} \sum_{t=1}^T \sum_{i=1}^C (DT_c(t, i) - d_{t,i}(\mathbf{x}))^2$$

Table of Contents

- 1 Ensemble Learning
- 2 Random Forest
- 3 Incremental Learning

Review: Classification and Regression Trees

- Grow a binary tree
- At each node, “split” the data into two “daughter” nodes
- Splits are chosen using a splitting criterion
- For regression the predicted value at a node is the average response variable for all observations in the node
- For classification the predicted class is the most common class in the node (majority vote)

Splitting criterion

- For classification: ID3, C4.5, CART
- For regression: CART
- ID3:

$$H(D) = - \sum_{k=1}^K \frac{|C_k|}{|D|} \log_2 \frac{|C_k|}{|D|}$$

$$H(D | A) = \sum_{i=1}^n \frac{|D_i|}{|D|} H(D_i)$$

$$\text{Gain}(D, A) = H(D) - H(D | A)$$

Splitting criterion

- C4.5:

$$\text{Gain}_{\text{ratio}}(D, A) = \frac{\text{Gain}(D, A)}{H_A(D)}$$
$$H_A(D) = - \sum_{i=1}^n \frac{|D_i|}{|D|} \log_2 \frac{|D_i|}{|D|}$$

- CART(classification):

$$\text{Gini}(D) = \sum_{k=1}^K \frac{|C_k|}{|D|} \left(1 - \frac{|C_k|}{|D|}\right)$$
$$= 1 - \sum_{k=1}^K \left(\frac{|C_k|}{|D|}\right)^2$$

Splitting criterion

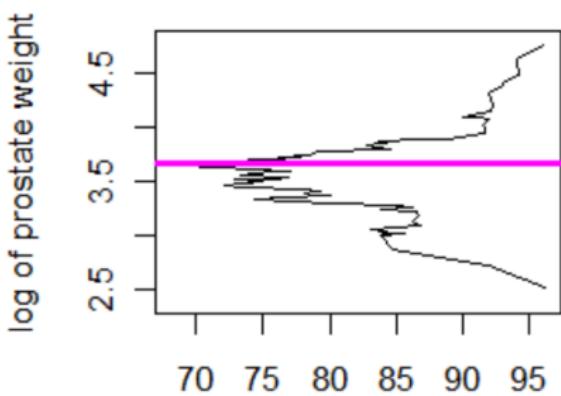
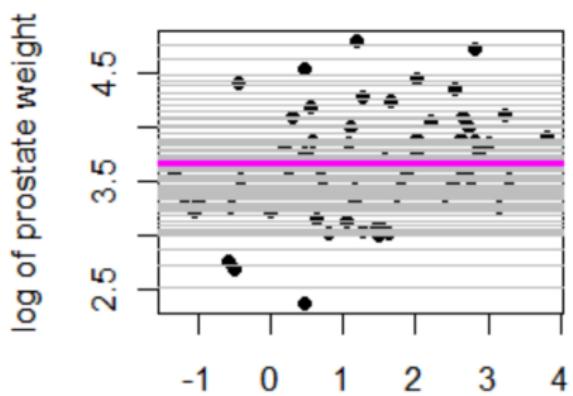
- CART(classification):

$$\text{Gini}(D \mid A) = \sum_{i=1}^n \frac{|D_i|}{|D|} \text{Gini}(D_i)$$

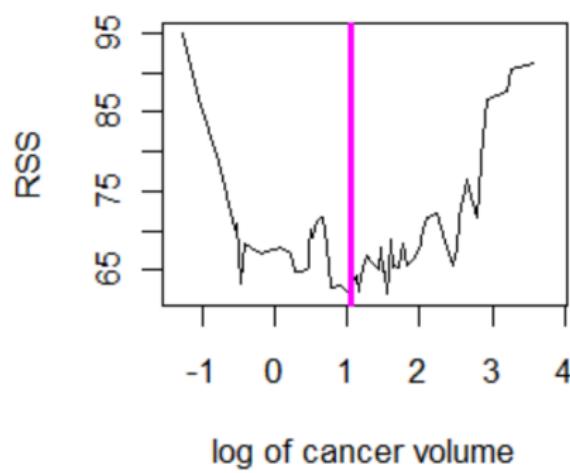
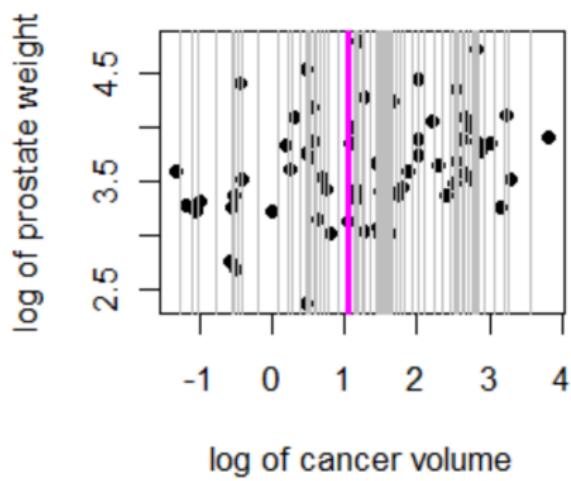
- CART(regression):

$$\min_{a,s} \left[\min_{c_1} \sum_{x_i \in D_1} (y_i - c_1)^2 + \min_{c_2} \sum_{x_i \in D_2} (y_i - c_2)^2 \right]$$

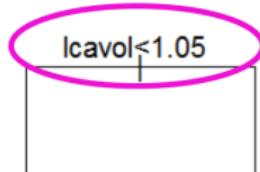
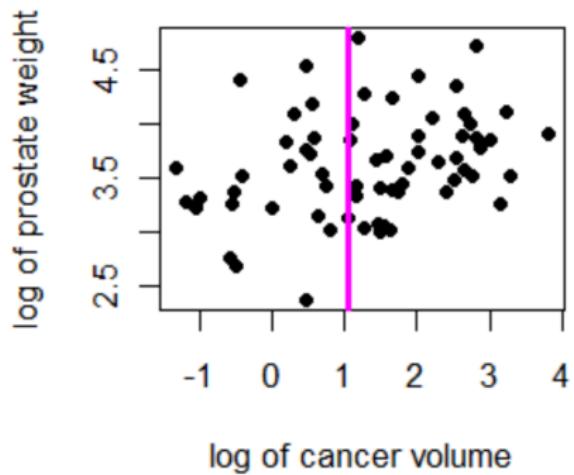
Choosing the best horizontal split



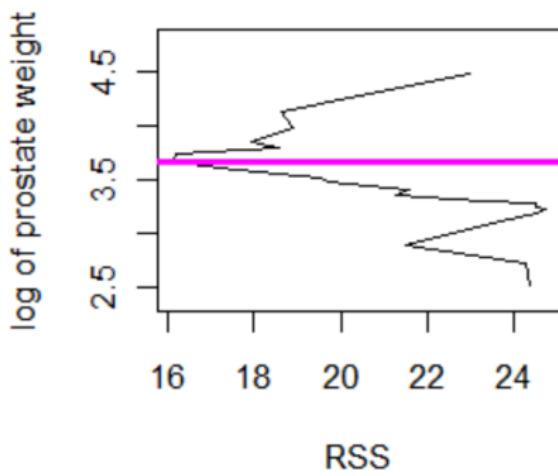
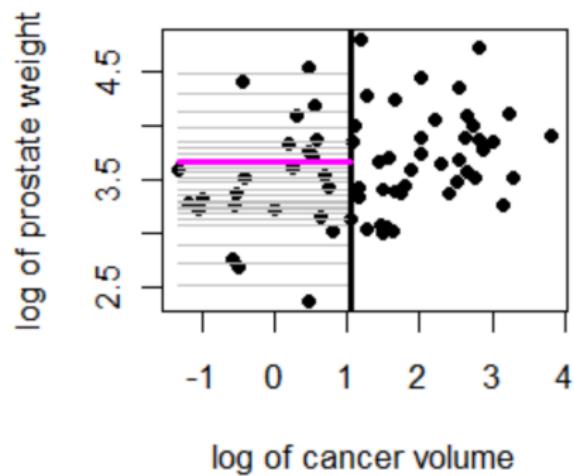
Choosing the best vertical split



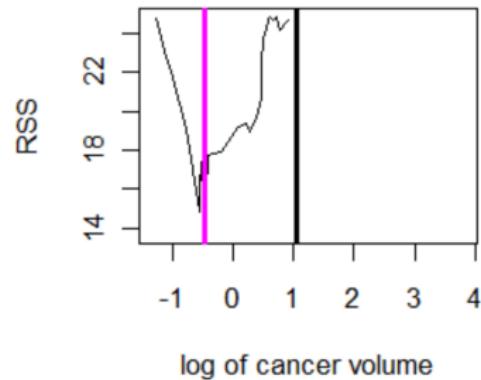
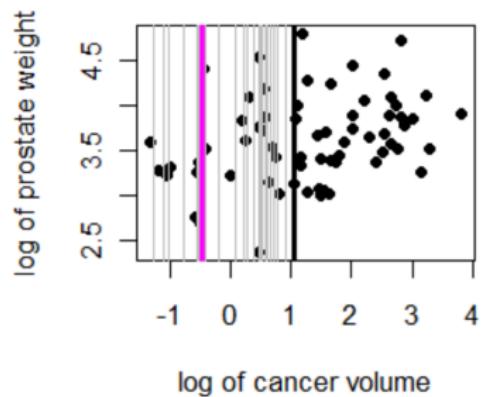
Regression tree (prostate cancer)



Choosing the best split in the left node

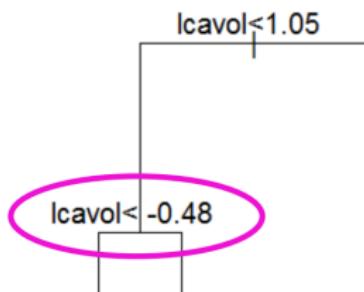
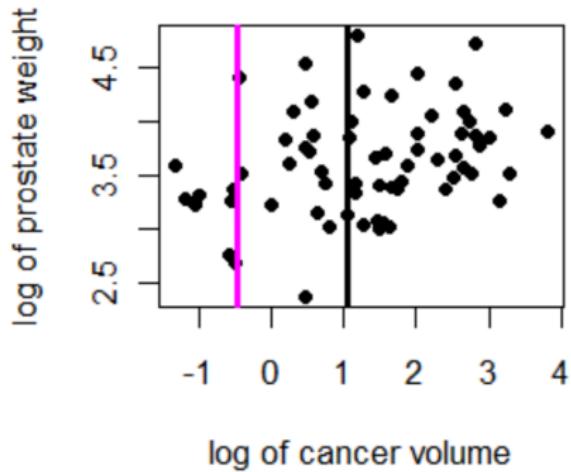


Choosing the best split in the left node

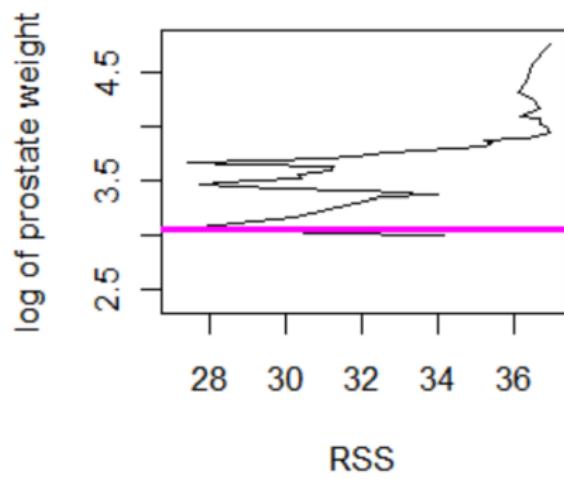
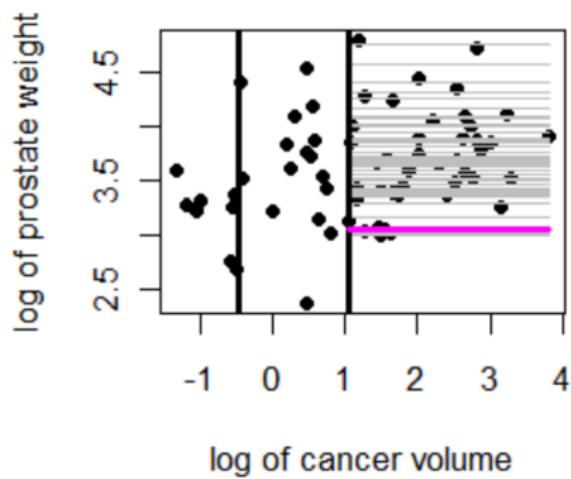


Best vertical split is at -0.48 with $\text{RSS} = 13.61$.

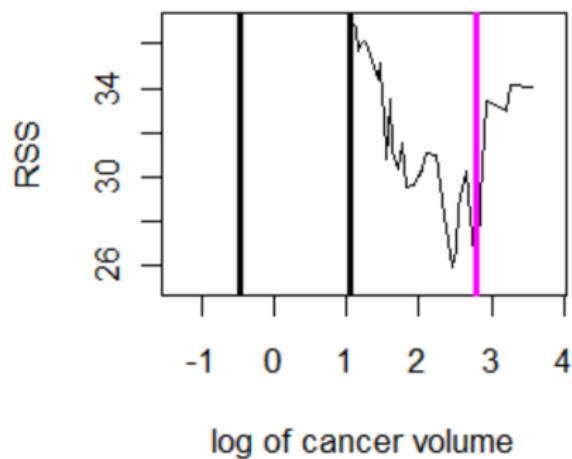
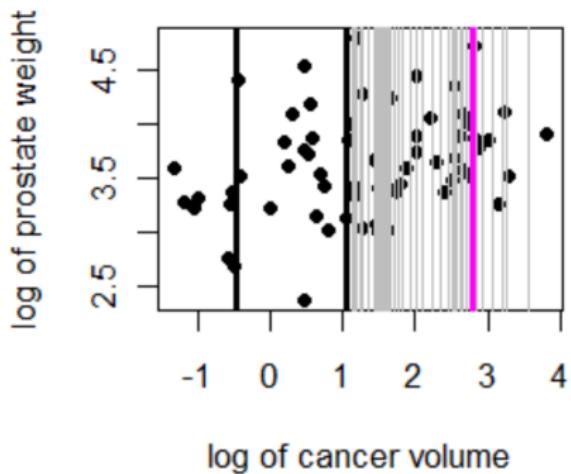
Regression tree (prostate cancer)



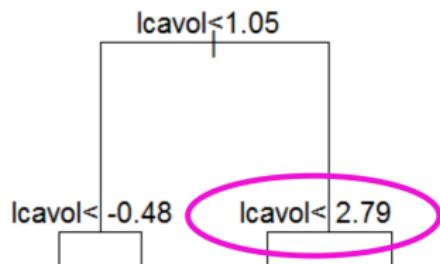
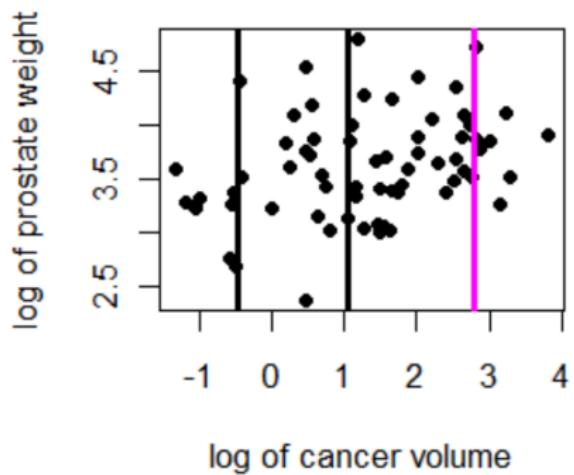
Choosing the best split in the right node



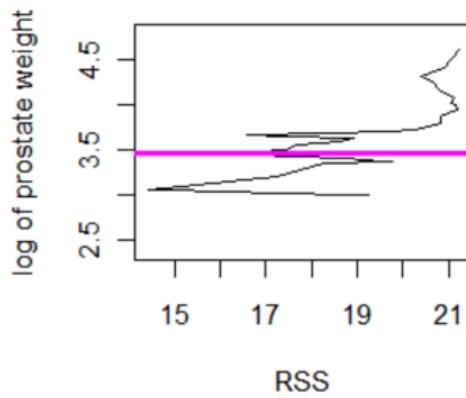
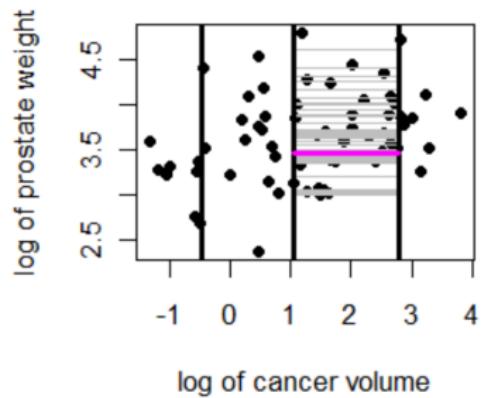
Choosing the best split in the right node



Regression tree (prostate cancer)

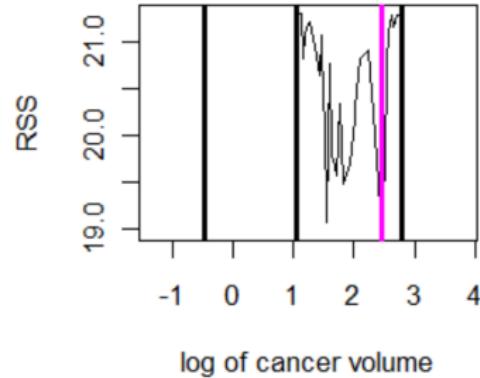
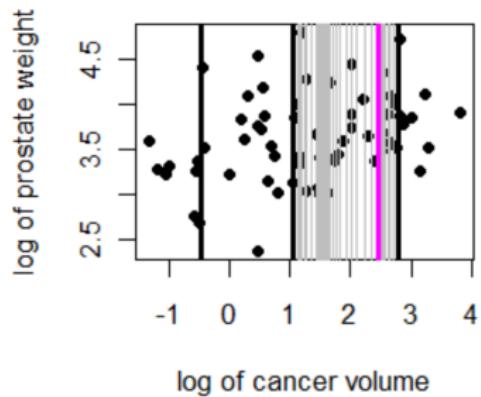


Choosing the best split in the third node



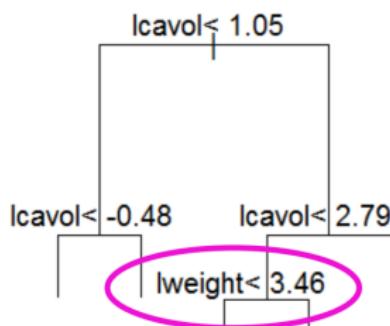
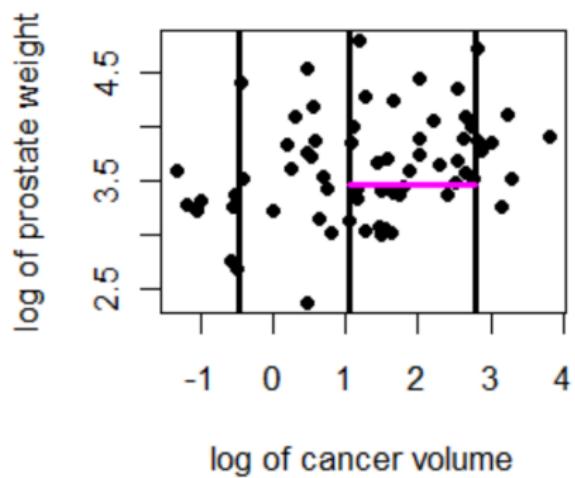
Best horizontal split is at 3.07 with RSS = 14.42, but this is too close to the edge. Use 3.46 with RSS = 16.14.

Choosing the best split in the third node

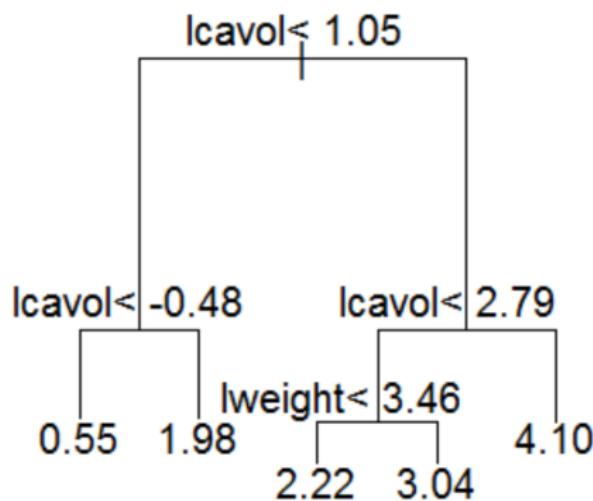
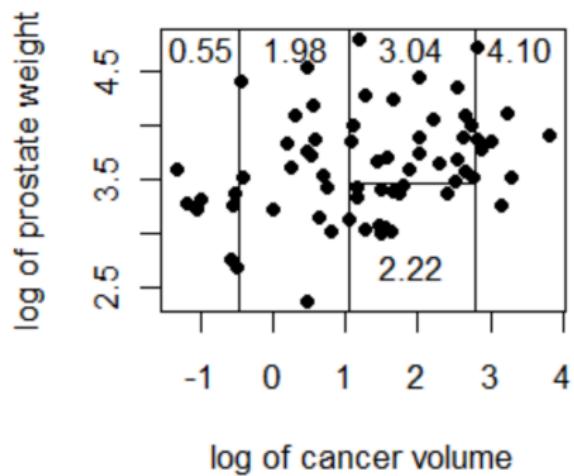


Best vertical split is at 2.46 with RSS = 18.97.

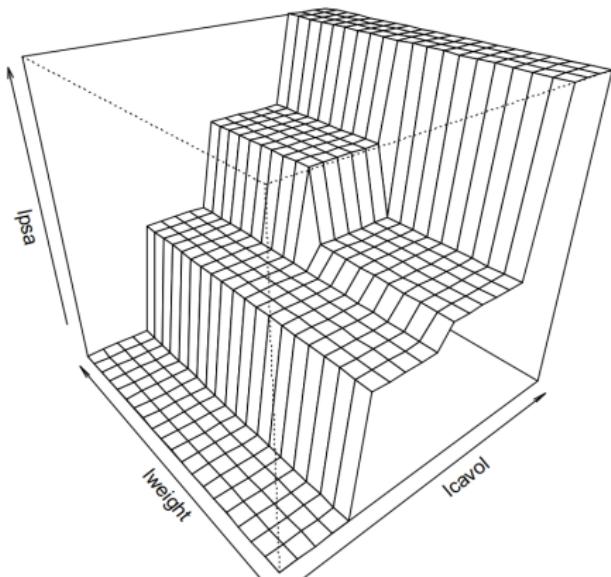
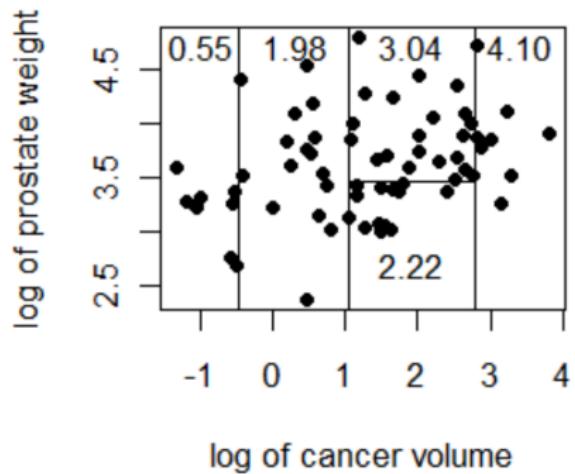
Regression tree (prostate cancer)



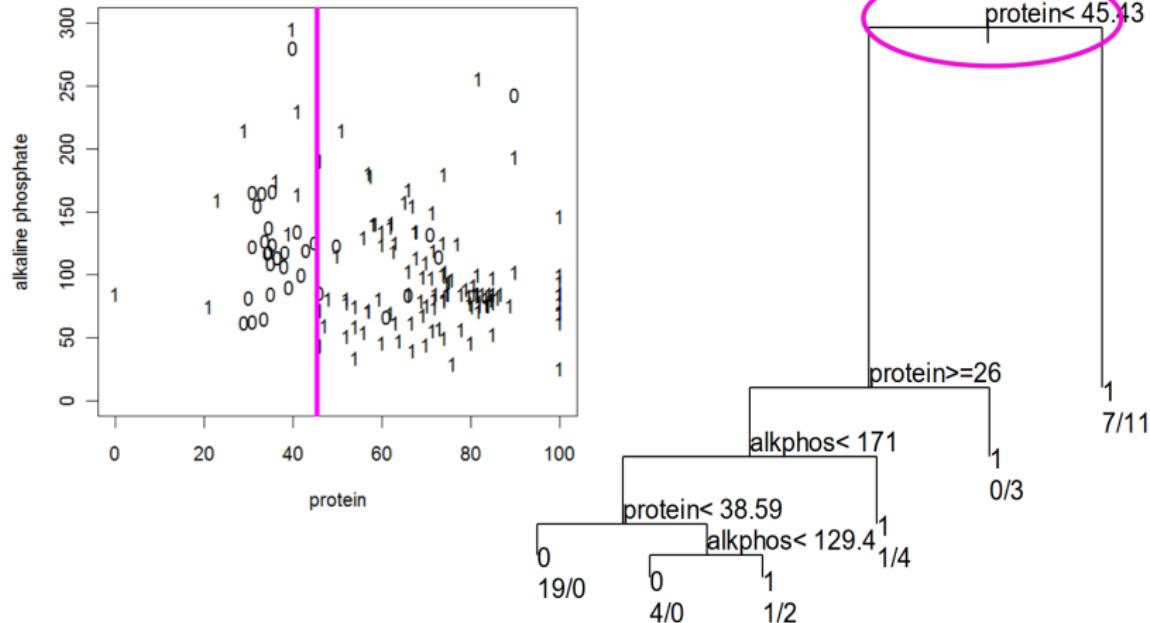
Regression tree (prostate cancer)



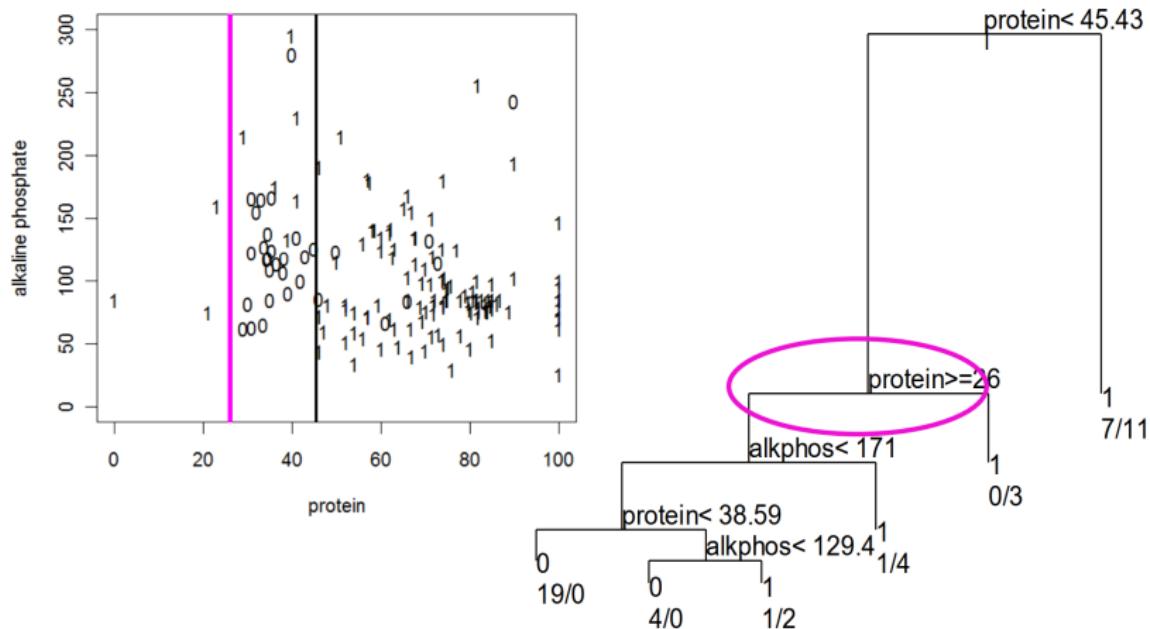
Regression tree (prostate cancer)



Classification tree (hepatitis)



Classification tree (hepatitis)



The idea of Random Forest

- The essential idea in bagging is to average many noisy but approximately unbiased models
- If B trees are iid with σ^2 , then average variance is $\frac{1}{B}\sigma^2$
- If B trees are simply i.d. , the variance of the average is $\rho\sigma^2 + \frac{1-\rho}{B}\sigma^2$
- As B increases, the second term disappears, but the first remains
- Random Forest builds a large collection of de-correlated trees, averages trees

Random Forest Algorithm

Algorithm 15.1 *Random Forest for Regression or Classification.*

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 :

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

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

The difference of RF and CART

- RF inherit many of the advantages of CART:
 1. Applicable to both regression and classification problems
 2. No formal distributional assumptions (non-parametric)
 3. Can handle highly non-linear interactions and classification boundaries
 4. Handles missing values
- The improvement of RF:
 1. Accuracy
 2. Instability

Details of Random Forest

- Why random forests do not overfit as more trees are added?
- Define the margin function:

$$mg(X, Y) = \frac{1}{B} \sum_{b=1}^B I(T_b(X) = Y) - \max_{j \neq Y} \frac{1}{B} \sum_{b=1}^B I(T_b(X) = j)$$

- The generalization error:

$$PE^* = P_{X,Y}(mg(X, Y) < 0)$$

- As the number of trees increases, for almost surely converges to

$$P_{X,Y} \left(P_\Theta(T(X, \Theta) = Y) - \max_{j \neq Y} P_\Theta(T(X, \Theta) = j) < 0 \right)$$

Details of Random Forest

- The inventors make the following recommendations:
 1. For classification, the default value for m is $\lfloor \sqrt{p} \rfloor$ and the minimum node size is one.
 2. For regression, the default value for m is $\lfloor p/3 \rfloor$ and the minimum node size is five.
- Out of Bag Samples: the original sample is not chosen in bootstraps process
- Suppose we have N samples, the probability of a sample being OOB is:

$$\left(1 - \frac{1}{N}\right)^N \approx e^{-1} \approx 0.368$$

Details of Random Forest

- An oob error estimate is almost identical to that obtained by N-fold crossvalidation
- Denote S_b be the index set of observations in b^{th} tree
- $C_i = \{b_i : i \in S_b\}$ is the trees set using the i^{th} sample
- For classification:

$$\hat{C}_{rf}^B(x) = \text{majorityvote}\{\hat{C}_b(x)\}_1^B$$

$$Err_{oob} = \frac{1}{N} \sum_{i=1}^N I(y_i \neq \hat{C}_{rf}^B(x_i))$$

$$Err_{cv} = \frac{1}{N} \sum_{i=1}^N I(y_i \neq F_{cv}^{(-i)}(x_i))$$

Details of Random Forest: OOB estimate

- Out-of bag estimates for strength and correlation
- We have:

$$P_{\Theta}(T(x, \Theta) = y) - \max_{j \neq y} P_{\Theta}(T(x, \Theta) = j)$$

- $Q(x, j)$ is an estimate for $P_{\Theta}(T(x, \Theta) = j)$

$$Q(x, j) = \sum_k I(T(x, \Theta_k) = j; (y, x) \notin T_{k,B}) / \sum_k I((y, x) \notin T_{k,B})$$

- The estimate strength is s

$$s = Q(x, y) - \max_{j \neq y} Q(x, j)$$

Details of Random Forest: OOB estimate

- The estimate correlation is $\bar{\rho}$

$$\bar{\rho} = \text{var}(mr) / (E_{\Theta} sd(\Theta))^2$$

- where

$$\text{var}(mr) = E_{X,Y} \left[P_{\Theta}(T(x, \Theta) = y) - \max_{j \neq y} P_{\Theta}(T(x, \Theta) = j) \right]^2 - s^2$$

$$sd(\Theta) = \left[p_1 + p_2 + (p_1 - p_2)^2 \right]^{1/2}$$

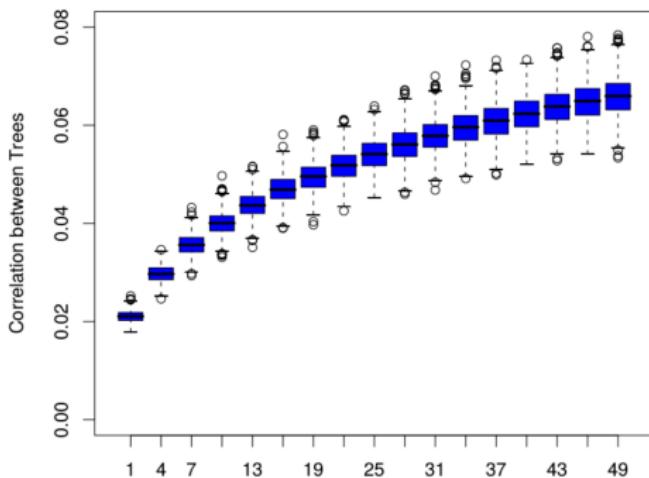
$$p_1 = E_{X,Y}(T(X, \Theta) = Y)$$

$$p_2 = E_{X,Y}(T(X, \Theta) = \hat{j}(X, Y))$$

$$\hat{j}(X, Y) = \arg \max_{j \neq Y} P_{\Theta}(T(X, \Theta) = j)$$

Details of Random Forest

- $Y = \frac{1}{\sqrt{50}} \sum_{j=1}^{50} X_j + \epsilon$
- pairs of tree predictions at x for different training sets Z are likely to be less similar if they do not use the same splitting variables



Details of Random Forest: variable importance

- Four methods to measure variable importance by Breiman
- Method 1:
 1. randomly permute all values of the m^{th} variable
 2. computing a new internal error rate
 3. The amount by which this new error exceeds the original test set error is defined as the importance of the m^{th} variable
- Method 2:
 1. for the n^{th} sample, calculate the margin
 2. the average lowering of the margin across all cases when the m^{th} variable is randomly permuted as in method 1

Details of Random Forest: variable importance

- Method 3:
 1. for the n^{th} sample, calculate the margin
 2. the count of how many margins are lowered minus the number of margins raised
- Method 4:
 1. the gini criterion

Table of Contents

- 1 Ensemble Learning
- 2 Random Forest
- 3 Incremental Learning

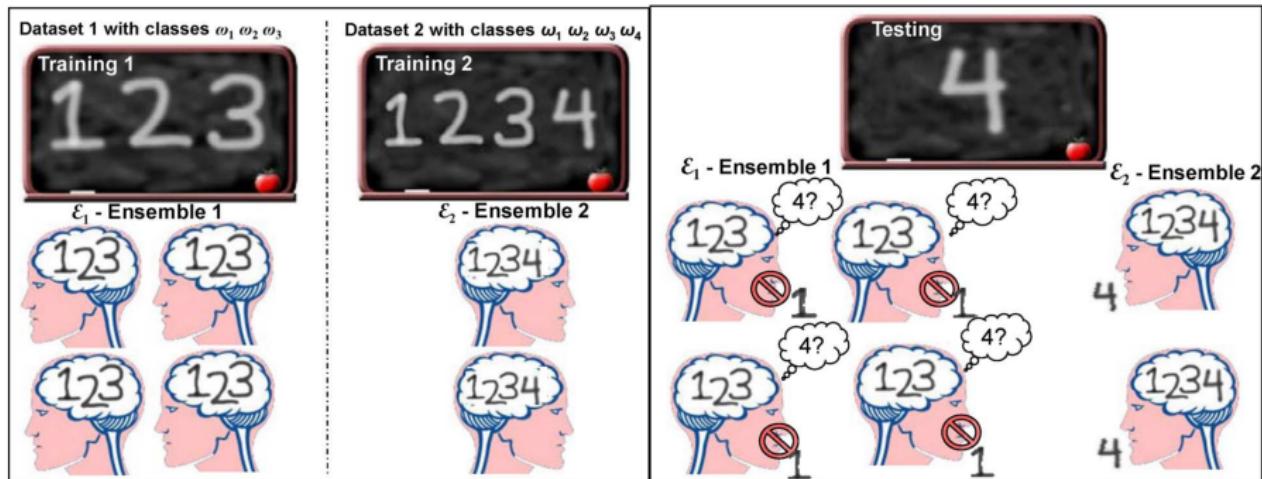
Incremental Learning

- stability–plasticity dilemma
- data often become available in batches over a period of time
- Goal: learn the novel information while retaining as much of the previous knowledge as possible

The idea of Incremental Learning

- discard the existing classifier, and retrain a new one using all data accumulated thus far
- Online learning algorithms: Winston's seminal work, Littlestone's Winnow
 - restricted forms of classifiers
 - Boolean or linear threshold functions,
- train an ensemble of classifiers for each data set
- combine their outputs using an appropriate combination rule
- To generate a diverse set of classifiers

- the weighted majority voting



Incremental Learning

Input: For each dataset $k=1, 2, \dots, K$

- Training data $S_k = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_i, y_i), \dots, (\mathbf{x}_{m_k}, y_{m_k})\}$, $\mathbf{x}_i \in X; y_i \in Y_k \subset \{\omega_1, \dots, \omega_C\}$
- Weak learning algorithm **BaseClassifier**.
- Integer T_k , specifying the number of **BaseClassifiers** to create using \mathcal{D}_k .

Do for $k=1, 2, \dots, K$

 Initialize the instance weights to be uniform when new data arrive $w_t^k(i) = 1/m_k, i = 1, \dots, m_k$

If $k \neq 1$, Set $t=0$ and **Go to** step 5 of the inner **Do loop** to adjust initialization weights

Do for $t=1, 2, \dots, T_k$:

1. Set $D_t^k = w_t^k / \sum_{i=1}^{m_k} w_t^k(i)$ so that D_t^k is a distribution.
2. Call **BaseClassifier**, providing it with $TR_t^k \subset S_k$, drawn according to D_t^k .
3. Obtain a hypothesis $h_t^k: X \rightarrow Y_k$, and calculate its error $\varepsilon_t^k = \sum_{i: h_t^k(\mathbf{x}_i) \neq y_i} D_t^k(i)$. If $\varepsilon_t^k > \frac{1}{2}$,
 discard h_t^k and go to step 2. Otherwise, compute normalized error $\beta_t^k = \varepsilon_t^k / (1 - \varepsilon_t^k)$.
4. Let $CL_t^k = Y_{k(t)} \subset \{\omega_1, \dots, \omega_C\}$ be the class labels used in training h_t^k for dataset \mathcal{D}_k .
5. Call **DW-CAV** to obtain the composite hypothesis H_t^k .
6. Compute the error of the composite hypothesis $E_t^k = \sum_{i: H_t^k(\mathbf{x}_i) \neq y_i} D_t^k(i)$
7. Set $B_t^k = E_t^k / (1 - E_t^k)$, and update the weights: $w_{t+1}^k(i) = w_t^k(i) \times \begin{cases} B_t^k, & \text{if } H_t^k(\mathbf{x}_i) = y_i \\ 1, & \text{otherwise} \end{cases}$

End

End

Call **DW-CAV** to obtain the final hypothesis, H_{final}

Incremental Learning

Inputs:

- Instance \mathbf{x}_i to be classified
- All classifiers h_t^k generated thus far
- Normalized error values, β_t^k .
- Class labels, CL_i^k used in training h_t^k .

Initialize classifier voting weights $W_t^k = \log(1/\beta_t^k)$

Calculate for each $\omega_c \in \{\omega_1, \dots, \omega_C\}$

- Normalization factor $Z_c = \sum_k \sum_{t: c \in CL_t^k} W_t^k$
$$Z_c = \sum_k \sum_{t: h_t^k(\mathbf{x}_i) = \omega_c} W_t^k$$
- Class-specific confidence $P_c(i) = \frac{W_t^k}{Z_c}$

If $P_k(i) = P_l(i) = 1, k \neq l$ such that

$\mathcal{E}_k \cap \mathcal{E}_l = \emptyset \rightarrow \text{Set } P_k(i) = P_l(i) = 0$

(where \mathcal{E}_k is the set of classifiers that have seen class ω_k)

Update voting weights for instance \mathbf{x}_i

$$W_t^k(i) = W_t^k \cdot \prod_{c: \omega_c \notin CL_t^k} (1 - P_c(i))$$

Compute final (or current composite) hypothesis

$$H_{final}(\mathbf{x}_i) = \arg \max_{\omega_c \in \Omega} \sum_k \sum_{t: h_t^k(\mathbf{x}_i) = \omega_c} W_t^k(i)$$