

Machine Learning Basics II

Agenda

- Why should we combine classifiers ?
- Voting Classifier
- Ensemble learning
- Ensemble: Parallel Learners
- Bagging
- Out of bag errors
- Random Forest
- Variable Importance
- Ensemble: Boosted Learners
- AdaBoost
- XGBoost
- TPOT: Genetic Evolution

Why should we combine classifiers ?

Voting Classifier

How to combine the classifiers?

- (weighted) Majority voting
 - Class label output
 - Select the class most voted for

$$\sum_{t=1}^T d_{t,J} = \max_{j=1}^C \sum_{t=1}^T d_{t,j}.$$

- Mean rule
 - Continuous output
 - Support for class w_j is average of classifier output

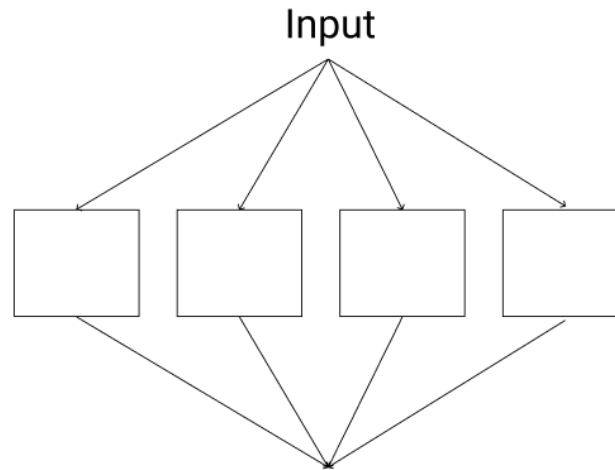
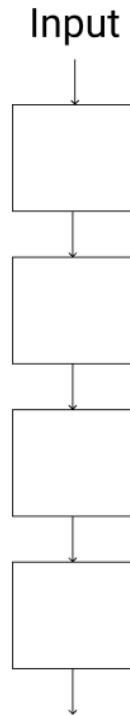
$$\sum_{t=1}^T w_t d_{t,J} = \max_{j=1}^C \sum_{t=1}^T w_t d_{t,j}$$

- Product rule
 - Continuous output
 - Product of classifier output

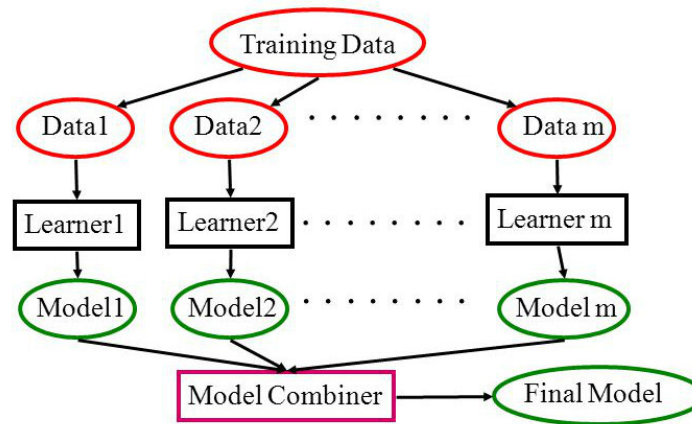
$$\mu_j(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T d_{t,j}(\mathbf{x})$$

$$\mu_j(\mathbf{x}) = \frac{1}{T} \prod_{t=1}^T d_{t,j}(\mathbf{x})$$

Ensemble learning



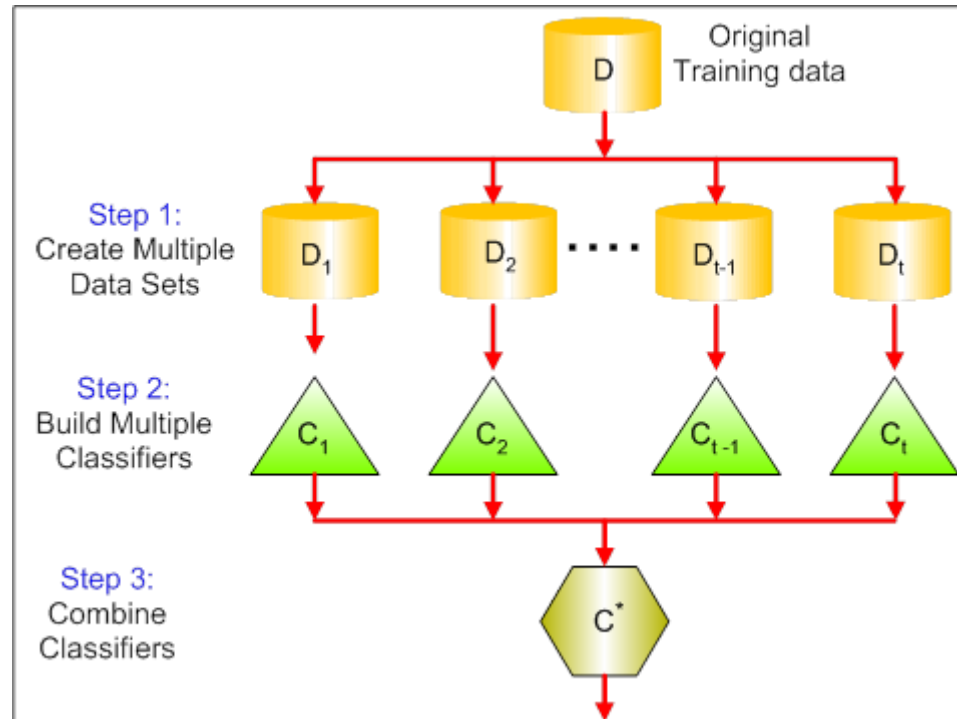
Ensemble learning: Parallel Learning



Source: Ray Mooney

Carla P. Gomes
CS4700


Bagging



Out of Bag Errors

The out-of-bag estimate is as accurate as using a test set of the same size as the training set.

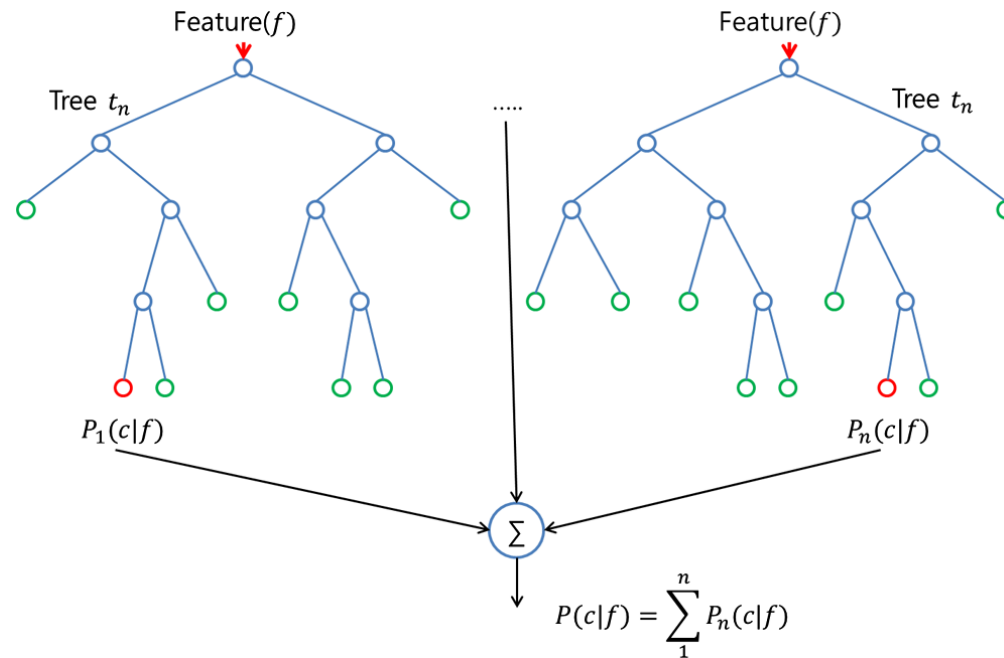
10 Total elements



Trained	0,1,3	2,6,8,9	1,5,7,9	3,7,8
	clf 1	clf 2	clf 3	clf 4
Out of Bag Test Elements	2,4,5,6, 7,8,9	1,3,4,5, 7	2,3,4,6,8	1,2,4,5, 6,9
Accuracy	.6	.8	.9	.9

Out of Bag Error Rate = Mean of Out of Bag Accuracies
= 0.8

Random Forest



Random Forest: Parameters

The main parameters to consider tuning are:

- `n_estimators` : The number of trees in the forest.
- `criterion` : The function to measure the quality of a split. Supported criteria are `"gini"` for the Gini impurity and `"entropy"` for the information gain. Note: this parameter is tree-specific.
- `max_features` : The number of features to consider when looking for the best split
- `max_depth` : The maximum depth of the tree.

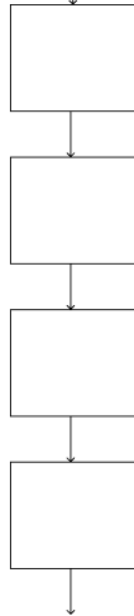
Variable Importance

The importance of each variable can be found by using 2 techniques:

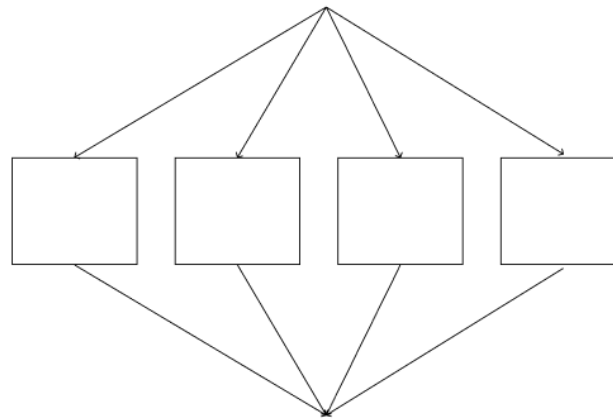
- Gini Importance: It is the average of the entropy gain in each tree (bias towards features with lots of categories).
- Permutation Importance: It is the effect on the OOB if the values of the feature are randomly permuted.

Ensemble: Boosted Learners

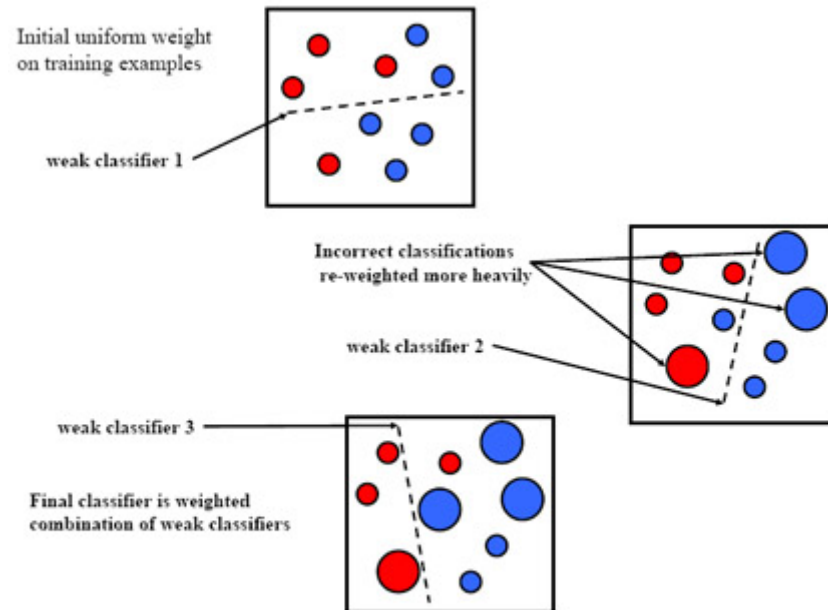
Input



Input



Adaboost



$$H(x) = \text{sign}(\alpha_1 h_1(x) + \alpha_2 h_2(x) + \alpha_3 h_3(x))$$

Adaboost (II)

- Start same weight for all points: $\alpha_i = 1/N$

- For $t = 1, \dots, T$

- Learn $f_t(\mathbf{x})$ with data weights α_i

- Compute coefficient \hat{w}_t

- Recompute weights α_i

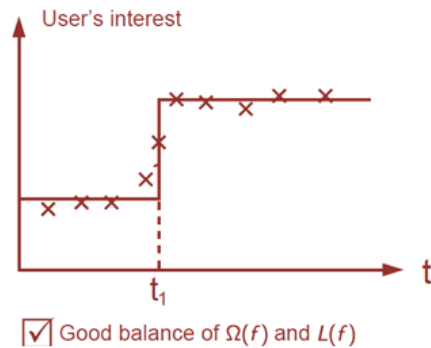
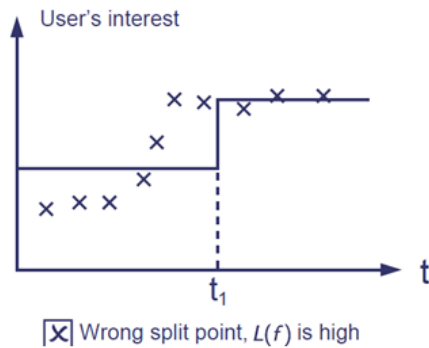
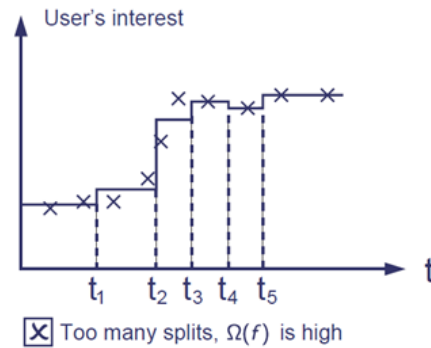
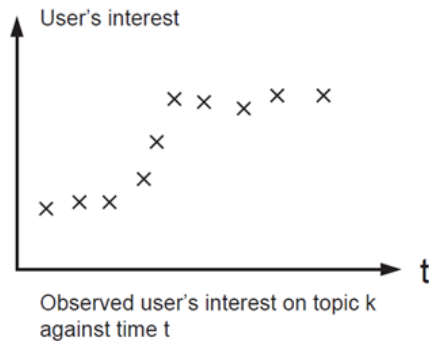
$$\hat{w}_t = \frac{1}{2} \ln \left(\frac{1 - \text{weighted_error}(f_t)}{\text{weighted_error}(f_t)} \right)$$

- Final model predicts by:

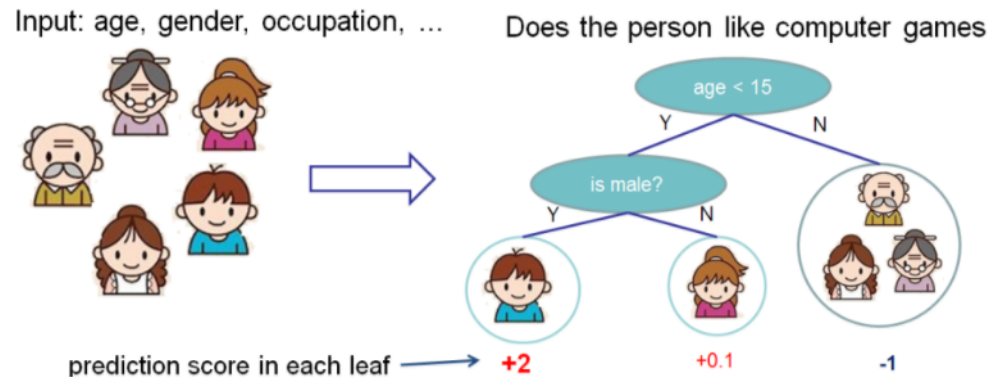
$$\hat{y} = \text{sign} \left(\sum_{t=1}^T \hat{w}_t f_t(\mathbf{x}) \right)$$

$$\alpha_i \leftarrow \begin{cases} \alpha_i e^{-\hat{w}_t}, & \text{if } f_t(\mathbf{x}_i) = y_i \\ \alpha_i e^{\hat{w}_t}, & \text{if } f_t(\mathbf{x}_i) \neq y_i \end{cases}$$

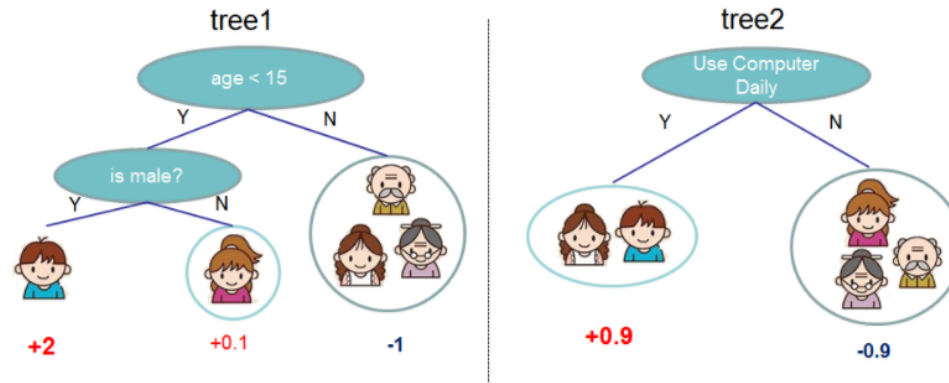
XGboost: Complexity



XGboost: Does a person like video games ?



XGboost: Trees



XGboost: Optimization Function

$$Obj(\Theta) = L(\Theta) + \Omega(\Theta)$$

Training Loss measures how well model fit on training data

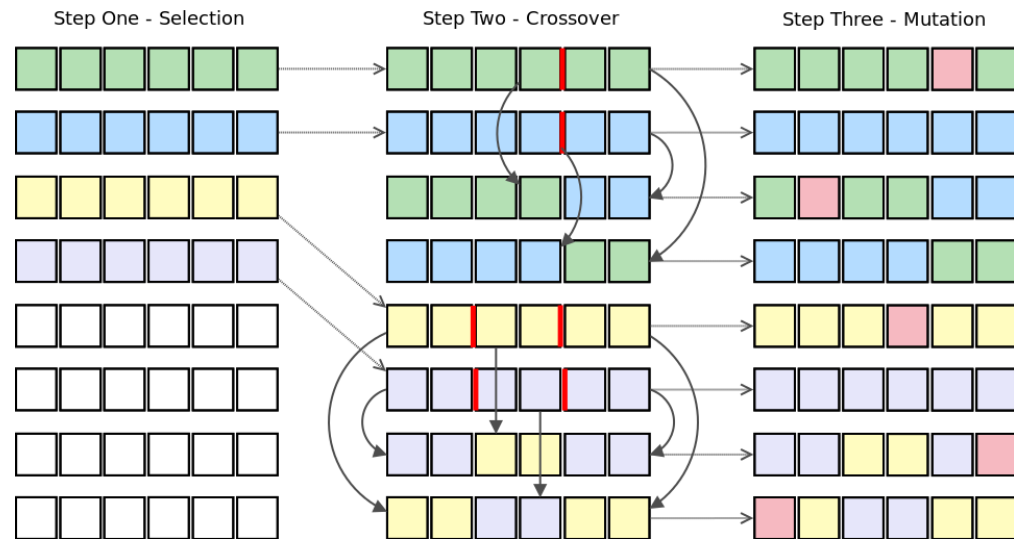
Regularization, measures complexity of model

XGboost: Parameters

The parameters to consider tuning are:

- The number and size of trees (n_estimators and max_depth).
- The learning rate and number of trees (learning_rate and n_estimators).
- The row and column subsampling rates (subsample, colsample_bytree and colsample_bylevel).

TPOT: Genetic Algorithm



TPOT

