

Learning Ensembles

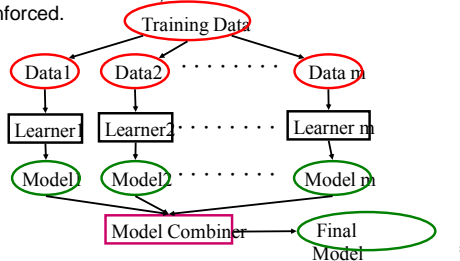
Tao Yang, 290N UCSB, 2013

Outlines

- Learning Assembles
- Random Forest
- Adaboost

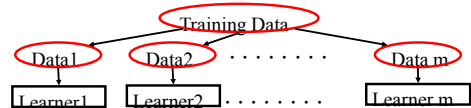
Learning Ensembles

- Learn multiple classifiers separately
- Combine decisions (e.g. using weighted voting)
- When combining multiple decisions, random errors cancel each other out, correct decisions are reinforced.



Homogenous Ensembles

- Use a single, arbitrary learning algorithm but manipulate training data to make it learn multiple models.
 - $Data_1 \neq Data_2 \neq \dots \neq Data_m$
 - $Learner_1 = Learner_2 = \dots = Learner_m$
- Methods for changing training data:
 - Bagging: Resample training data
 - Boosting: Reweight training data
 - DECORATE: Add additional artificial training data



Bagging

- Create ensembles by repeatedly randomly resampling the training data (Breiman, 1996).
- Given a training set of size n , create m sample sets
 - Each *bootstrap sample set* will on average contain 63.2% of the unique training examples, the rest are replicates.
- Combine the m resulting models using majority vote.
- Decreases error by decreasing the variance in the results due to *unstable learners*, algorithms (like decision trees) whose output can change dramatically when the training data is slightly changed.

5

Random Forests

- **Introduce two sources of randomness: “Bagging” and “Random input vectors”**
 - Each tree is grown using a bootstrap sample of training data
 - At each node, best split is chosen from random sample of m variables instead of all variables M .
- m is held constant during the forest growing
- Each tree is grown to the largest extent possible
- Bagging using decision trees is a special case of random forests when $m=M$

Random Forests

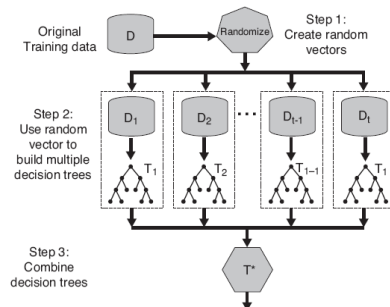


Figure 5.40. Random forests.

Random Forest Algorithm

- Good accuracy without over-fitting
- Fast algorithm (can be faster than growing/pruning a single tree); easily parallelized
- Handle high dimensional data without much problem

Boosting: AdaBoost

Yoav Freund 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):119–139, August 1997.

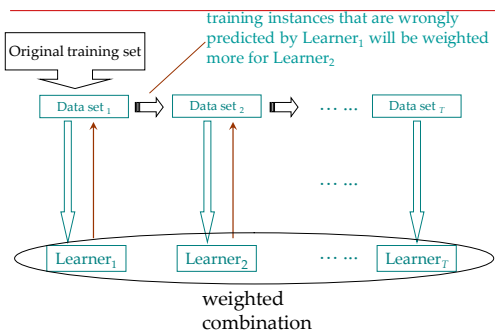
- Simple with theoretical foundation

Adaboost - Adaptive Boosting

- Use training set re-weighting
 - Each training sample uses a weight to determine the probability of being selected for a training set.
- AdaBoost is an algorithm for constructing a “strong” classifier as linear combination of “simple” “weak” classifier
$$f(x) = \sum_{t=1}^T \alpha_t h_t(x)$$
- Final classification based on weighted sum of weak classifiers

10

AdaBoost: An Easy Flow

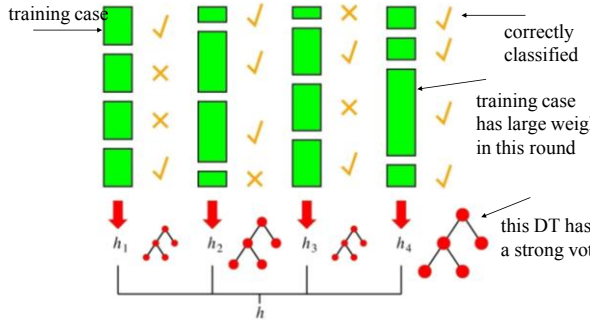


Adaboost Terminology

- $h_t(x)$... “weak” or basis classifier
- $H(x) = \text{sign}(f(x))$... “strong” or final classifier
- Weak Classifier: < 50% error over any distribution
- Strong Classifier: thresholded linear combination of weak classifier outputs

12

And in a Picture



AdaBoost.M1

- Given **training set** $X = \{(x_i, y_i)\}$,
 - $y_i \in \{-1, +1\}$ **correct label**
 - Initialize **distribution** $D_t(i) = 1/m$; (**weight of training cases**)
- for $t = 1, \dots, T$:
 - Find a **weak classifier** ("rule of the form" $h_t: X \rightarrow \{-1, +1\}$) with small error ϵ_t on D_t ;
 - Update distribution D_t on $\{1, \dots, m\}$

$$D_{t+1}(i) = \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$
 - $y_i * h_t(x_i) > 0$, if correct
 - $y_i * h_t(x_i) < 0$, if wrong
- Output $H(x) = \text{sign}(\sum_{t=1}^T \alpha_t h_t(x))$
 - where Z_t is a normalization factor (chosen so that D_{t+1} will be a distribution).

Each training sample has a weight, which determines the probability of being selected for training the component classifier

Weighted voting for the final decision-making

case last-ssis

Reweighting

Effect on the training set

Reweighting formula:

$$D_{t+1}(i) = \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

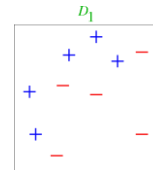
$$\exp(-\alpha_t y_i h_t(x_i)) \begin{cases} < 1, & y_i = h_t(x_i) \\ > 1, & y_i \neq h_t(x_i) \end{cases}$$

$y * h(x) = 1$

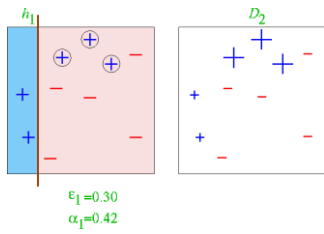
$y * h(x) = -1$

⇒ Increase (decrease) weight of wrongly (correctly) classified examples

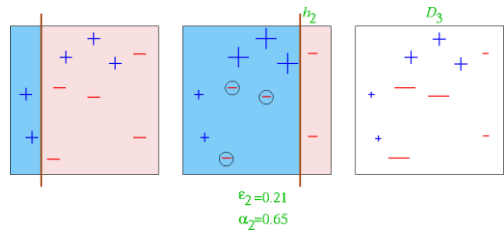
Toy Example



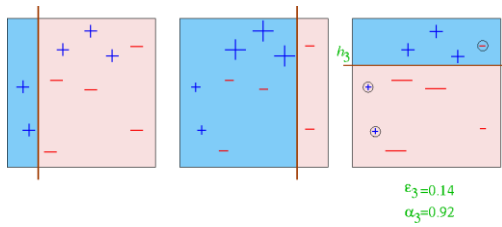
Round 1



Round 2



Round 3



Final Combination

$$H_{\text{final}} = \text{sign} \left(0.42 \left(\begin{array}{|c|} \hline \text{blue} \\ \hline \end{array} \right) + 0.65 \left(\begin{array}{|c|} \hline \text{blue} \\ \hline \end{array} \right) + 0.92 \left(\begin{array}{|c|} \hline \text{blue} \\ \hline \end{array} \right) \right)$$

$$= \begin{array}{|c|} \hline \text{blue} \\ \hline \end{array}$$

Pros and cons of AdaBoost

Advantages

- Very simple to implement
- Does feature selection resulting in relatively simple classifier
- Fairly good generalization

Disadvantages

- Suboptimal solution
- Sensitive to noisy data and outliers

21

References

- Duda, Hart, ect - *Pattern Classification*
- Freund - "An adaptive version of the boost by majority algorithm"
- Freund - "Experiments with a new boosting algorithm"
- Freund, Schapire - "A decision-theoretic generalization of on-line learning and an application to boosting"
- Friedman, Hastie, etc - "Additive Logistic Regression: A Statistical View of Boosting"
- Jin, Liu, etc (CMU) - "A New Boosting Algorithm Using Input-Dependent Regularizer"
- Li, Zhang, etc - "Floatboost Learning for Classification"
- Opitz, Maclin - "Popular Ensemble Methods: An Empirical Study"
- Ratsch, Warmuth - "Efficient Margin Maximization with Boosting"
- Schapire, Freund, etc - "Boosting the Margin: A New Explanation for the Effectiveness of Voting Methods"
- Schapire, Singer - "Improved Boosting Algorithms Using Confidence-Weighted Predictions"
- Schapire - "The Boosting Approach to Machine Learning: An overview"
- Zhang, Li, etc - "Multi-view Face Detection with Floatboost"

22

AdaBoost: Training Error Analysis

- Suppose $f(x) = \sum_{t=1}^T \alpha_t h_t(x)$ Equivalent $H(x) = \text{sign}(f(x))$
if $H(x_i) \neq y_i$ then $y_i f(x_i) \leq 0$ implying that $\exp(-y_i f(x_i)) \geq 1$. Thus,

- Therefore, training error is: $\frac{1}{m} \sum_{i: H(x_i) \neq y_i} 1 \leq \frac{1}{m} \sum_{i=1}^m \exp(-y_i f(x_i))$.

- As: $\frac{1}{m} \sum_{i: H(x_i) \neq y_i} 1 \leq \frac{1}{m} \sum_{i=1}^m \exp(-y_i f(x_i))$
Considering $D_{r+1}(i) = \frac{\exp(-\sum_{t=1}^r \alpha_t y_t h_t(x_i))}{\sum_{j=1}^m \exp(-\sum_{t=1}^r \alpha_t y_t h_t(x_j))} = \prod_{t=1}^r Z_t$
Finally: $\frac{1}{m} \sum_{i: H(x_i) \neq y_i} 1 \leq \prod_{t=1}^r Z_t$

$\{i: H(x_i) \neq y_i\}$ is a vector which i -th element is $[H(x_i) \neq y_i]$.
 $|\{i: H(x_i) \neq y_i\}|$ is the sum of all the element in the vector

AdaBoost: How to choose α_t

- According to $\frac{1}{m} \sum_{i: H(x_i) \neq y_i} 1 \leq \prod_{t=1}^r Z_t$
This equation is obvious if we treat u_t as a binary-valued variable.
Let $u_t = y_t r_t$
By setting $\alpha_t = u_t$ easily get
Actually AdaBoost can just minimize the training error.
Minimize the error bound could be done.
By setting $\alpha_t = u_t$ easily get
Let $\epsilon = \sum_{i: H(x_i) \neq y_i} D_t(i)$, we have $\epsilon = \frac{1-r}{2}$
- The right term is minimized
 $\alpha = \frac{1}{2} \ln \left(\frac{1+r}{1-r} \right)$
Then the training error of H is at most $\prod_{t=1}^r \sqrt{1-r_t^2}$.