Boosting

Sargur N. Srihari srihari@buffalo.edu

Topics in Boosting

- What is boosting
- Boosting vs Bagging
- Adaboost algorithm

Boosting

- Incrementally adding models to the ensemble
 - After a weak learner is added, the data are reweighted:
 - examples that are misclassified gain weight and examples that are classified correctly lose weight
- Has been applied to ensembles of neural networks, by incrementally adding neural networks to the ensemble
- Also interpreting a neural network as an ensemble, incrementally adding hidden units to the network

AdaBoost (Adaptive Boosting)

- Most widely used form of boosting is the AdaBoost algorithm
- Boosting can give good results even if base classifiers have performance, only slightly better than random
 - Hence base learners are called weak learners
- Boosting can be extended to regression

- The base classifiers are trained in sequence
- Each base classifier is trained using a weighted form of the data set in which the weighting coefficient associated with each data point depends on the performance of the previous classifiers
 - Misclassified points are given greater weight when used to train the next classifier
- Once all classifiers have been trained, their predictions are combined through a weighted majority weighting scheme

Boosting framework



Each base classifier $y_m(\mathbf{x})$ is trained on a weighted form of the training set.

Weights $w^{(m)}{}_n$ depend on the performance of the previous base classifier $y_{m-1}(\mathbf{x})$

Once all base classifiers are trained, they are combined to give the final classifier $Y_M(\mathbf{x})$

Adaboost (informal)

- Consider 2-class problem
 - Data set $X = \{x_1, .., x_N\}$
 - Binary target variables $t_1, .., t_N, t_n \in \{-1, 1\}$
 - Each data point has weight w_n , initially set to 1/N
 - We have a procedure available for training a base classifier using weighted data to give $y(\mathbf{x}) \in \{-1,1\}$
- At each stage, train a new classifier
 - using a data set in which the weights are adjusted giving higher weight to those misclassified
- When desired no of base classifiers are trained they are combined to form a committee
 - With different weights to base classifiers

Machine Learning

Adaboost Algorithm

1. Initialize $w^{(1)}_{n} = 1/N$ for n=1,..N

2. For m=1,..M

- a) Fit a classifier $y_m(\mathbf{x})$ by minimizing weighted error $J_m = \sum_{n=1}^{N} w_n^{(m)} I(y_m(x_n) \neq t_n)$ where $I(y_m(x_n) \neq t_n)$ is the indicator function and equals 1 when $I(y_m(x_n) \neq t_n)$ and 0 otherwise ϵ_m ; weighted error
- b) Evaluate the quantities $\epsilon_m = \frac{J_m}{\sum_{m=1}^N w_m^{(m)}}$

and then use these to evaluate $\alpha_m = \ln \left\{ \frac{1 - \epsilon_m}{\epsilon_m} \right\}$

c) Update the data weighting coefficients

 $w_n^{(m+1)} = w_n^{(m)} \exp\{\alpha_m I(y_m(x_n) \neq t_n)\}$

3. Make predictions using final model $Y_M(x) = sign(\sum_{m=1}^{M} \alpha_m y_m(x))$ ϵ_m : weighted error rate of classifier

 $\alpha_{\rm m}$: weighting coeffts give greater weight to more accurate classifiers

$$w_n^{(m+1)}$$
: increase
weight of
misclassified data
with exponential erro

Notes on Adaboost algorithm

- First base classifier y₁(x) is trained using weighting coefficients w⁽¹⁾_n that are all set equal Usual procedure for training a single classifier
- From step 2(c): $w_n^{(m+1)} = w_n^{(m)} \exp\{\alpha_m I(y_m(x_n) \neq t_n)\}$ in subsequent iterations $w_n^{(m)}$ are increased for data points that are misclassified by previous classifiers and data points that continue to be misclassified receive even greater weight

Data set for illustrating Adaboost



True posterior probabilities are shown on right

Base Classifier

- Threshold on one of the input variables
- Equivalent to a decision tree with decision stumps (decision tree with a single node)

Srihari

Illustration of Boosting



Dashed line: boundary of latest learner Each data point is a circle whose radius indicates weight Green line: combined decision boundary of ensemble

Minimizing exponential errors

- Boosting was originally motivated using statistical learning theory
 - Leading to upper bounds on generalization error
 - Bounds turn out to be too loose
 - Actual performance is much better
- Instead, a simple interpretation of boosting in terms of sequential minimization of an exponential error function

Exponential error function

Consider the exponential error function

$$E = \sum_{n=1}^{N} \exp[-t_n f_m(x_n)]$$

- where $f_m(x)$ is a classifier defined in terms of a linear combination of base classifiers of the form

$$f_m(x) = \frac{1}{2} \sum_{i=1}^m \alpha_i y_i(x)$$

- and $t_n \in \{-1,1\}$ are the training set values

- Our goal is to minimize E wrt both the weighting coefficients and the parameters of the base classifiers $y_{l}(x)$ $^{\rm 14}$

Error Functions for Boosting

- Exponential error function minimized by AdaBoost differs from previous ML methods
- Consider expected error

 $E_{x,t}[\exp\{-t\ y(x)\}] = \sum_t \int \exp\{-t\ y(x)\}p(t|x)p(x)dx$

- If we perform variational minimization wrt all possible functions y(x) we obtain

$$y(x) = \frac{1}{2} \ln \left\{ \frac{p(t=1|x)}{p(t=-1|x)} \right\}$$

- Which is half the log odds
- Thus AdaBoost is seeking the best approximation to the log odds ratio



