

Combining Classifiers: Boosting the Margin and Mixtures of Experts

Tuesday 12 November 2002

William H. Hsu

Department of Computing and Information Sciences, KSU

http://www.kddresearch.org http://www.cis.ksu.edu/~bhsu

Readings: "Bagging, Boosting, and *C4.5*", Quinlan Section 5, "*MLC++* Utilities 2.0", Kohavi and Sommerfield



CIS 732: Machine Learning and Pattern Recognition

Lecture Outline

- Readings: Section 5, MLC++ 2.0 Manual [Kohavi and Sommerfield, 1996]
- Paper Review: "Bagging, Boosting, and C4.5", J. R. Quinlan
- Boosting the Margin
 - Filtering: feed examples to trained inducers, use them as "sieve" for consensus
 - <u>Resampling</u>: aka subsampling (S[i] of fixed size m' resampled from D)
 - <u>Reweighting</u>: fixed size S[i] containing weighted examples for inducer
- Mixture Model, aka Mixture of Experts (ME)
- <u>Hierarchical Mixtures of Experts (HME)</u>
- Committee Machines
 - Static structures: *ignore input signal*
 - Ensemble averaging (single-pass: weighted majority, bagging, stacking)
 - Boosting the margin (some single-pass, some multi-pass)
 - Dynamic structures (multi-pass): *use input signal to improve classifiers*
 - Mixture of experts: training in combiner inducer (aka gating network)
 - Hierarchical mixtures of experts: hierarchy of inducers, combiners

Quick Review: Ensemble Averaging

- Intuitive Idea
 - Combine experts (aka prediction algorithms, classifiers) using combiner function
 - Combiner may be weight vector (WM), vote (bagging), trained inducer (stacking)
- <u>Weighted Majority (WM)</u>
 - Weights each algorithm in proportion to its training set accuracy
 - Use this weight in performance element (and on test set predictions)
 - Mistake bound for WM
- <u>Bootstrap Aggregating (Bagging)</u>
 - Voting system for collection of algorithms
 - Training set for each member: sampled with replacement
 - Works for <u>unstable inducers</u> (search for *h* sensitive to perturbation in *D*)
- Stacked Generalization (aka Stacking)
 - Hierarchical system for combining inducers (ANNs or other inducers)
 - Training sets for "leaves": sampled with replacement; combiner: validation set
- Single-Pass: Train Classification and Combiner Inducers Serially
- Static Structures: Ignore Input Signal



Boosting: Idea

- Intuitive Idea
 - Another type of static committee machine: can be used to improve any inducer
 - Learn set of classifiers from *D*, but reweight examples to *emphasize misclassified*
 - − Final classifier ← weighted combination of classifiers
- Different from Ensemble Averaging
 - WM: all inducers trained on same D
 - Bagging, stacking: training/validation partitions, i.i.d. *subsamples S*[*i*] of *D*
 - Boosting: data sampled according to different distributions
- Problem Definition
 - <u>Given</u>: collection of multiple inducers, large data set or example stream
 - <u>Return</u>: combined predictor (trained committee machine)
- Solution Approaches
 - <u>Filtering</u>: use weak inducers in <u>cascade</u> to filter examples for downstream ones
 - <u>Resampling</u>: reuse data from D by subsampling (don't need huge or "infinite" D)
 - <u>Reweighting</u>: reuse $x \in D$, but measure error over weighted x



Boosting: Procedure

- Algorithm Combiner-AdaBoost (D, L, k)
 - $m \leftarrow D.size$
 - FOR $i \leftarrow 1$ TO m DO
 - Distribution[i] $\leftarrow 1 / m$
 - FOR $i \leftarrow 1$ TO k DO
 - $P[j] \leftarrow L[j]$. Train-Inducer (Distribution, D) // assume L[j] identical; $h_j \equiv P[j]$
 - $Error[i] \leftarrow Count-Errors(P[i], Sample-According-To (Distribution, D))$
 - $\beta[i] \leftarrow Error[i] / (1 Error[i])$
 - FOR $i \leftarrow 1$ TO m DO // update distribution on D $Distribution[i] \leftarrow Distribution[i] * ((P[i](D[i]) = D[i], target) ? \beta[i] : 1)$
 - Distribution.Renormalize ()
 - RETURN (*Make-Predictor* (P, D, β))
- Function *Make-Predictor* (P, D, β) ٠
 - // Combiner(x) = $\operatorname{argmax}_{v \in V} \sum_{j:P[j](x) = v} \lg (1/\beta[j])$
 - RETURN (fn $x \Rightarrow$ *Predict-Argmax-Correct* (*P*, *D*, *x*, fn $\beta \Rightarrow \lg (1/\beta)$))

- // Resampling Algorithm
- // initialization
- // subsampling distribution

// Invariant: *Distribution* is a pdf



Boosting: Properties

- Boosting in General
 - Empirically shown to be effective
 - Theory still under development
 - Many variants of boosting, active research (see: references; current ICML, COLT)
- Boosting by Filtering
 - Turns weak inducers into strong inducer (committee machine)
 - Memory-efficient compared to other boosting methods
 - **<u>Property</u>**: improvement of weak classifiers (trained inducers) guaranteed
 - Suppose 3 experts (subhypotheses) each have error rate $\varepsilon < 0.5$ on D[i]
 - Error rate of committee machine $\leq g(\varepsilon) = 3\varepsilon^2 2\varepsilon^3$
- Boosting by Resampling (*AdaBoost*): Forces *Error*_D toward *Error*
- References
 - Filtering: [Schapire, 1990] MLJ, 5:197-227
 - Resampling: [Freund and Schapire, 1996] ICML 1996, p. 148-156
 - Reweighting: [Freund, 1995]
 - Survey and overview: [Quinlan, 1996; Haykin, 1999]



Mixture Models: Idea

- Intuitive Idea
 - Integrate knowledge from multiple experts (or data from multiple sensors)
 - Collection of inducers organized into committee machine (e.g., modular ANN)
 - **Dynamic structure**: *take input signal into account*
 - References
 - [Bishop, 1995] (Sections 2.7, 9.7)
 - [Haykin, 1999] (Section 7.6)
- Problem Definition
 - <u>Given</u>: collection of inducers ("experts") *L*, data set *D*
 - <u>Perform</u>: supervised learning using inducers and self-organization of experts
 - <u>Return</u>: committee machine with trained <u>gating network</u> (combiner inducer)
- Solution Approach
 - Let combiner inducer be generalized linear model (e.g., threshold gate)
 - Activation functions: linear combination, vote, "smoothed" vote (softmax)



Mixture Models: Procedure

- Algorithm Combiner-Mixture-Model (D, L, Activation, k)
 - $m \leftarrow D.size$
 - FOR $j \leftarrow 1$ TO k DO

w[*j*] ← 1

- UNTIL the termination condition is met, DO
 - FOR $j \leftarrow 1$ TO k DO
 - $P[j] \leftarrow L[j]. Update-Inducer(D)$
 - FOR *i* ← 1 TO *m* DO

Sum[i] ← 0 FOR $j \leftarrow 1$ TO k DO Sum[i] += P[j](D[i]) Net[i] ← Compute-Activation (Sum[i]) // compute $g_j \equiv Net[i][j]$ FOR $j \leftarrow 1$ TO k DO $w[j] \leftarrow Update-Weights (w[j], Net[i], D[i])$

- RETURN (Make-Predictor (P, w))
- Update-Weights: Single Training Step for Mixing Coefficients



// single training step for L[j]

// initialization

Mixture Models: Properties

- Unspecified Functions
 - Update-Inducer
 - Single training step for each expert module
 - e.g., ANN: one backprop cycle, aka epoch
 - Compute-Activation
 - Depends on ME <u>architecture</u>
 - Idea: smoothing of "winner-take-all" ("hard" max)
 - Softmax activation function (Gaussian mixture model)

$$\boldsymbol{g}_{l} = \frac{\boldsymbol{e}^{\vec{\boldsymbol{w}}_{l} \cdot \vec{\boldsymbol{x}}}}{\sum_{i=1}^{k} \boldsymbol{e}^{\vec{\boldsymbol{w}}_{j} \cdot \vec{\boldsymbol{x}}}}$$



- Possible Modifications
 - Batch (as opposed to online) updates: lift Update-Weights out of outer FOR loop
 - Classification learning (versus concept learning): multiple y_i values
 - Arrange gating networks (combiner inducers) in *hierarchy* (HME)



<u>Generalized Linear Models (GLIMs)</u>

• Recall: Perceptron (Linear Threshold Gate) Model



- Generalization of LTG Model [McCullagh and Nelder, 1989]
 - Model parameters: connection weights as for LTG
 - Representational power: depends on transfer (activation) function
- Activation Function
 - Type of mixture model depends (in part) on this definition
 - e.g., o(x) could be softmax $(x \cdot w)$ [Bridle, 1990]
 - *NB*: *softmax* is computed across *j* = 1, 2, ..., *k* (cf. "hard" max)
 - Defines (multinomial) pdf over experts [Jordan and Jacobs, 1995]



<u>Hierarchical Mixture of Experts (HME):</u> Idea



Department of Computing and Information Sciences

<u>Hierarchical Mixture of Experts (HME):</u> Procedure

- Algorithm Combiner-HME (D, L, Activation, Level, k, Classes)
 - $m \leftarrow D.size$
 - FOR *j* ← 1 TO *k* DO *w*[*j*] ← 1

// initialization

- UNTIL the termination condition is met DO
 - IF Level > 1 THEN
 - FOR $j \leftarrow 1$ TO k DO

 $P[j] \leftarrow Combiner-HME(D, L[j], Activation, Level - 1, k, Classes)$

• ELSE

_

FOR $j \leftarrow 1$ TO k DO $P[j] \leftarrow L[j]$. Update-Inducer (D)

• FOR *i* ← 1 TO *m* DO

```
Sum[i] \leftarrow 0
FOR j \leftarrow 1 TO k DO
Sum[i] += P[j](D[i])
Net[i] \leftarrow Compute-Activation (Sum[i])

FOR l \leftarrow 1 TO Classes DO w[l] \leftarrow Update-Weights (w[l], Net[i], D[i])

RETURN (Make-Predictor (P, w))
```



<u>Hierarchical Mixture of Experts (HME):</u> Properties

- Advantages
 - Benefits of ME: base case is single level of expert and gating networks
 - More combiner inducers \Rightarrow more capability to <u>decompose</u> complex problems
- Views of HME
 - Expresses <u>divide-and-conquer</u> strategy
 - Problem is distributed across subtrees "on the fly" by combiner inducers
 - Duality: data fusion ⇔ problem redistribution
 - Recursive decomposition: until good fit found to "local" structure of D
 - Implements soft decision tree
 - Mixture of experts: 1-level decision tree (decision stump)
 - Information preservation compared to traditional (hard) decision tree
 - Dynamics of HME improves on greedy (high-commitment) strategy of decision tree induction



Training Methods for <u>H</u>ierarchical <u>M</u>ixture of <u>Experts</u> (HME)

- Stochastic Gradient Ascent
 - Maximize <u>log-likelihood function</u> $L(\Theta) = \lg P(D | \Theta)$
 - Compute

$$\frac{\partial \boldsymbol{L}}{\partial \boldsymbol{w}_{ij}}, \frac{\partial \boldsymbol{L}}{\partial \boldsymbol{a}_{j}}, \frac{\partial \boldsymbol{L}}{\partial \boldsymbol{a}_{ij}}$$

- Finds MAP values
 - Expert network (leaf) weights w_{ij}
 - Gating network (interior node) weights at lower level (a_{ij}) , upper level (a_j)
- <u>Expectation-Maximization (EM)</u> Algorithm
 - Recall definition
 - Goal: maximize incomplete-data log-likelihood function $L(\Theta) = \lg P(D \mid \Theta)$
 - Estimation step: calculate E[unobserved variables $| \Theta]$, assuming current Θ
 - <u>Maximization step</u>: update Θ to maximize $E[Ig P(D | \Theta)], D \equiv all variables$
 - Using EM: estimate with gating networks, then adjust $\Theta = \{w_{ij}, a_{ij}, a_{j}\}$



Methods for Combining Classifiers: Committee Machines

- Framework
 - Think of collection of trained inducers as committee of experts
 - Each produces predictions given input (s(D_{test}), i.e., new x)
 - Objective: combine predictions by vote (subsampled D_{train}), learned weighting function, or more complex combiner inducer (trained using D_{train} or $D_{validation}$)
- Types of Committee Machines
 - Static structures: based only on *y* coming out of local inducers
 - Single-pass, same data or independent subsamples: WM, bagging, stacking
 - Cascade training: AdaBoost
 - Iterative reweighting: boosting by reweighting
 - Dynamic structures: take x into account
 - Mixture models (mixture of experts aka ME): one combiner (gating) level
 - <u>Hierarchical Mixtures of Experts (HME)</u>: multiple combiner (gating) levels
 - <u>Specialist-Moderator (SM)</u> networks: partitions of x given to combiners



Comparison of Committee Machines

	Aggregating Mixtures			Partitioning Mixtures	
	Stacking	<u>Bagging</u>	SM Networks	<u>Boosting</u>	<u>HME</u>
Sampling Method	Round-robin (cross- validation)	Random, with replacement	Attribute partitioning/ clustering	Least squares (proportionate)	Linear gating (proportionate)
Splitting of Data	Length-wise	Length-wise	Length-wise	Width-wise	Width -wise
Guaranteed improvement of weak classifiers?	Νο	Νο	No	Yes	No
Hierarchical?	Yes	No, but can be extended	Yes	No	Yes
Training	Single bottom- up pass	N/A	Single bottom- up pass	Multiple passes	Multiple top- down passes
Wrapper or mixture?	Both	Wrapper	Mixture, can be both	Wrapper	Mixture, can be both



Terminology

- <u>Committee Machines</u> aka Combiners
- <u>Static Structures</u>
 - Ensemble averaging
 - Single-pass, separately trained inducers, common input
 - Individual outputs combined to get scalar output (e.g., linear combination)
 - **Boosting the margin: separately trained inducers**, *different input distributions*
 - <u>Filtering</u>: feed examples to trained inducers (<u>weak classifiers</u>), pass on to next classifier *iff* conflict encountered (<u>consensus model</u>)
 - <u>Resampling</u>: *aka* subsampling (*S*[*i*] of fixed size *m*' *resampled* from *D*)
 - <u>Reweighting</u>: fixed size *S*[*i*] containing *weighted examples* for inducer
- Dynamic Structures
 - Mixture of experts: training in combiner inducer (*aka* <u>gating network</u>)
 - Hierarchical mixtures of experts: hierarchy of inducers, combiners
- <u>Mixture Model</u>, *aka* <u>Mixture of Experts (ME)</u>
 - <u>Expert</u> (classification), <u>gating</u> (combiner) inducers (<u>modules</u>, "networks")
 - <u>Hierarchical Mixtures of Experts (HME)</u>: multiple combiner (gating) levels



Summary Points

- Committee Machines aka Combiners
- Static Structures (Single-Pass)
 - Ensemble averaging
 - For improving <u>weak</u> (especially <u>unstable</u>) classifiers
 - e.g., <u>weighted majority</u>, bagging, stacking
 - Boosting the margin
 - Improve performance of any inducer: weight examples to emphasize errors
 - Variants: filtering (*aka* consensus), resampling (*aka* subsampling), reweighting
- Dynamic Structures (Multi-Pass)
 - Mixture of experts: training in combiner inducer (aka gating network)
 - Hierarchical mixtures of experts: hierarchy of inducers, combiners
- Mixture Model (*aka* Mixture of Experts)
 - Estimation of mixture coefficients (i.e., weights)
 - Hierarchical Mixtures of Experts (HME): multiple combiner (gating) levels
- Next Week: Intro to GAs, GP (9.1-9.4, Mitchell; 1, 6.1-6.5, Goldberg)

