A Theory of Multiclass Boosting

Indraneel Mukherjee*, R. E. Schapire Princeton University

Wrigley Field prepared for college football game

Dublin warned over ECB liquidity

Newest senators Coons and Manchin sworn in



Politics

Business

Sports

Wrigley Field prepared for college football game

Dublin warned over ECB liquidity

Newest senators Coons and Manchin sworn in

Business











• Boost *simplest* weak classifiers

- Boost *simplest* weak classifiers
 - Use right weak learning condition (WLC)

- Boost *simplest* weak classifiers
 - Use right weak learning condition (WLC)
- Important for generalization error:

- Boost *simplest* weak classifiers
 - Use right weak learning condition (WLC)
- Important for generalization error:
 - Simple weak classifier may imply less overfitting

- Boost *simplest* weak classifiers
 - Use right weak learning condition (WLC)
- Important for generalization error:
 - Simple weak classifier may imply less overfitting
 - Too simple could lead to underfitting

- Boost *simplest* weak classifiers
 - Use right weak learning condition (WLC)
- Important for generalization error:
 - Simple weak classifier may imply less overfitting
 - Too simple could lead to underfitting
- Theory known for binary, not for multiclass

• Existing frameworks inadequate for multiclass

- Existing frameworks inadequate for multiclass
 - Most resulting WLC's are too weak or too strong

- Existing frameworks inadequate for multiclass
 - Most resulting WLC's are too weak or too strong
- Introduce new framework for multiclass boosting

- Existing frameworks inadequate for multiclass
 - Most resulting WLC's are too weak or too strong
- Introduce new framework for multiclass boosting
 - Captures the *minimal WLC*

- Existing frameworks inadequate for multiclass
 - Most resulting WLC's are too weak or too strong
- Introduce new framework for multiclass boosting
 - Captures the *minimal WLC*
- Boosting algorithm using the minimal WLC

- Existing frameworks inadequate for multiclass
 - Most resulting WLC's are too weak or too strong
- Introduce new framework for multiclass boosting
 - Captures the *minimal WLC*
- Boosting algorithm using the minimal WLC
 - Provably drives down error efficiently

- Existing frameworks inadequate for multiclass
 - Most resulting WLC's are too weak or too strong
- Introduce new framework for multiclass boosting
 - Captures the *minimal WLC*
- Boosting algorithm using the minimal WLC
 - Provably drives down error efficiently
 - Experiments to complement the theory

Binary boosting

Binary boosting

Input: (x_1, y_1) , ..., (x_m, y_m)

Booster

4 = {weak classifiers}









Final model: (weighted) majority{h₁, ... , h_T}

Binary boostingMore weight on
misclassified
examplesInput:
$$(x_1,y_1)$$
, ..., (x_m, y_m)
 d_1 , ..., d_m $300ster$ d_1 , ..., d_m
 $\mathcal{A} = \{weak classifiers\}$
 $h \in \mathcal{A}$. $h: \{Example\} \Rightarrow \{Label\}$ $b \in \mathcal{A}$. $h: \{Example\} \Rightarrow \{Label\}$ Condition : $\widehat{err}_d(h) \leq \frac{1}{2} - \gamma$

Final model: (weighted) majority{ $h_1, ..., h_T$ }

$$\begin{array}{l} \textbf{Binary boosting} \\ \textbf{More weight on} \\ \textbf{misclassified} \\ \textbf{examples} \\ \textbf{More weight on} \\ \textbf{misclassified} \\ \textbf{miscl$$

• Required tasks easy. Only better than random

- Required tasks easy. Only better than random
- Sufficient. # satisfies binary WLC => # is boostable

- Required tasks easy. Only better than random
- Sufficient. # satisfies binary WLC => # is boostable
 - Boostable space: contains perfect combination

- Required tasks easy. Only better than random
- Sufficient. # satisfies binary WLC => # is boostable
 - Boostable space: contains perfect combination
- Necessary. Boostable space satisfies binary WLC

- Required tasks easy. Only better than random
- Sufficient. # satisfies binary WLC => # is boostable
 - Boostable space: contains perfect combination
- Necessary. Boostable space satisfies binary WLC
- Effective. Allows efficient boosting algorithm

Extending to Multiclass




SAMME [Zhu, Zou, Rosset, Hastie '09]



SAMME [Zhu, Zou, Rosset, Hastie '09]

Extending to Multiclass
Input:
$$(x_1,y_1), ..., (x_m, y_m)$$

 $d_1, ..., d_m$
3000ster
 $h \in \mathcal{H}. h: \{Example\} \Rightarrow \{Multiclass Label\}$
 $\widehat{\operatorname{err}_d(h) \leq 1 - \frac{1}{k} - \gamma}$
SAMME [Zhu, Zou, Rosset, Hastie '09]
 $\widehat{\operatorname{err}_d(h) \leq \frac{1}{2} - \gamma}$
AdaBoost.MI [Freund, Schapire '96]

Extending to Multiclass
Input:
$$(x_1,y_1), ..., (x_m, y_m)$$

 $d_1, ..., d_m$
Booster
 $h \in \mathcal{H}. h: \{Example\} \Rightarrow \{Multiclass Label\}$
 $\widehat{\operatorname{err}_d(h) \leq 1 - \frac{1}{k} - \gamma}$
SAMME [Zhu, Zou, Rosset, Hastie '09]
 $\widehat{\operatorname{err}_d(h) \leq \frac{1}{2} - \gamma}$
AdaBoost.MI [Freund, Schapire '96]





•One-against-all, all-pairs, ECOC. E.g.



One-against-all, all-pairs, ECOC. E.g.
One-against-all: (AdaBoost.MH)[Schapire & Singer '99]



One-against-all, all-pairs, ECOC. E.g.
One-against-all: (AdaBoost.MH)[Schapire & Singer '99]
All-pairs: (AdaBoost.MR)[Freund & Schapire '96, Schapire & Singer '99]



One-against-all, all-pairs, ECOC. E.g.
One-against-all: (AdaBoost.MH)[Schapire & Singer '99]
All-pairs: (AdaBoost.MR)[Freund & Schapire '96, Schapire & Singer '99]
Practical, but poorly understood



- One-against-all, all-pairs, ECOC. E.g.
 One-against-all: (AdaBoost.MH)[Schapire & Singer '99]
 All-pairs: (AdaBoost.MR)[Freund & Schapire '96, Schapire & Singer '99]
 Practical, but poorly understood
- •Sometimes too strong



- One-against-all, all-pairs, ECOC. E.g.
 One-against-all: (AdaBoost.MH)[Schapire & Singer '99]
 All-pairs: (AdaBoost.MR)[Freund & Schapire '96, Schapire & Singer '99]
- •Practical, but poorly understood
- Sometimes too strong
 e.g. One-against-all (AdaBoost.MH)

• Booster sends cost matrix C, not distribution

- Booster sends cost matrix C, not distribution
 - C(i, ℓ): cost of predicting label ℓ on example i

• Cost(C, h) =
$$\sum_{i} C(i, h(x_i))$$

- Booster sends cost matrix C, not distribution
 - C(i, ℓ): cost of predicting label ℓ on example i

• Cost(C, h) =
$$\sum_{i} C(i, h(x_i))$$

• Perform as well as fixed baseline predictor B

- Booster sends cost matrix C, not distribution
 - C(i, ℓ): cost of predicting label ℓ on example i

• Cost(C, h) =
$$\sum_{i} C(i, h(x_i))$$

- Perform as well as fixed baseline predictor B
 - B(i, ℓ): prob. with which B predicts ℓ on i
 - Cost(C, B) = $\sum_{i} \mathbb{E}[C(i, B(x_i))] = \sum_{i} \sum_{\ell} C(i, I) B(i, \ell)$

- Booster sends cost matrix C, not distribution
 - C(i, ℓ): cost of predicting label ℓ on example i

• Cost(C, h) =
$$\sum_{i} C(i, h(x_i))$$

- Perform as well as fixed baseline predictor B
 - B(i, ℓ): prob. with which B predicts ℓ on i
 - Cost(C, B) = $\sum_{i} \mathbb{E}[C(i, B(x_i))] = \sum_{i} \sum_{\ell} C(i, I) B(i, \ell)$
 - Restriction: Cost $(C, h) \leq Cost (C, B)$

Parameter: Fixed baseline B

Parameter: Fixed baseline B



10

Parameter: Fixed baseline B



$$Cost(C, h) \leq Cost(C, B)$$



• Edge-over-random baseline Q

- Edge-over-random baseline Q
 - B(i, correct) \geq B(i, wrong) + 2γ
 - $B(i, \cdot)$ is a distribution

- Edge-over-random baseline Q
 - B(i, correct) \geq B(i, wrong) + 2γ
 - $B(i, \cdot)$ is a distribution
- Many choices for B (only one for binary)

- Edge-over-random baseline Q
 - B(i, correct) \geq B(i, wrong) + 2γ
 - $B(i, \cdot)$ is a distribution
- Many choices for B (only one for binary)
- Condition with such baseline:

- Edge-over-random baseline Q
 - B(i, correct) \geq B(i, wrong) + 2γ
 - $B(i, \cdot)$ is a distribution
- Many choices for B (only one for binary)
- Condition with such baseline:
 Edge-over-random WLC

• Required tasks easy. Only beat random

- Required tasks easy. Only beat random
- Sufficient. Satisfying EOR implies boostability

- Required tasks easy. Only beat random
- Sufficient. Satisfying EOR implies boostability
- Effective. Allows efficient boosting

- Required tasks easy. Only beat random
- Sufficient. Satisfying EOR implies boostability
- Effective. Allows efficient boosting
- Not Necessary. For any EOR (B), there is some boostable space # that does not satisfy it.

- Required tasks easy. Only beat random
- Sufficient. Satisfying EOR implies boostability
- Effective. Allows efficient boosting
- Not Necessary. For any EOR (B), there is some boostable space # that does not satisfy it.
- Relaxed necessity. For any boostable space #, there is some EOR (B) that # satisfies

- Required tasks easy. Only beat random
- Sufficient. Satisfying EOR implies boostability
- Effective. Allows efficient boosting
- Not Necessary. For any EOR (B), there is some boostable space # that does not satisfy it.
- Relaxed necessity. For any boostable space #, there is some EOR (B) that # satisfies
- Combine to form single minimal WLC

- Required tasks easy. Only beat random
- Sufficient. Satisfying EOR implies boostability
- Effective. Allows efficient boosting
- Not Necessary. For any EOR (B), there is some boostable space # that does not satisfy it.
- Relaxed necessity. For any boostable space #, there is some EOR (B) that # satisfies
- Combine to form single minimal WLC
 - Necessary and sufficient for boostability
- Optimally efficient algorithm for any fixed EOR
 - Like Boost-by-majority [Freund '95]

- Optimally efficient algorithm for any fixed EOR
 - Like Boost-by-majority [Freund '95]
 - Non-adaptive. Requires knowledge of γ

- Like Boost-by-majority [Freund '95]
- Non-adaptive. Requires knowledge of γ
- Adaptive algorithm assuming the minimal WLC

- Like Boost-by-majority [Freund '95]
- Non-adaptive. Requires knowledge of γ
- Adaptive algorithm assuming the minimal WLC
 - Based on multiplicative updates, like AdaBoost

- Like Boost-by-majority [Freund '95]
- Non-adaptive. Requires knowledge of γ
- Adaptive algorithm assuming the minimal WLC
 - Based on multiplicative updates, like AdaBoost
 - Not optimal, but still provably very efficient

• In each round t:

- In each round t:
 - Create cost matrix C_t

- In each round t:
 - Create cost matrix C_t
 - Receive weak classifier h_t with edge δ_t

- In each round t:
 - Create cost matrix C_t
 - Receive weak classifier h_t with edge δ_t
 - Compute weight α_t and update $f_t = f_{t-1} + \alpha_t h_t$

- In each round t:
 - Create cost matrix C_t
 - Receive weak classifier h_t with edge δ_t
 - Compute weight α_t and update $f_t = f_{t-1} + \alpha_t h_t$

Weight

$$\alpha_t = \ln\left\{\frac{1+\delta_t}{1-\delta_t}\right\}$$

Cost Matrix

$$C_{t+1}(i,l) = \begin{cases} e^{f_t(i,l) - f_t(i,y_i)} & \text{if } l \neq y_i \\ -\sum_{l' \neq y_i} e^{f_t(i,l') - f_t(i,y_i)} & \text{if } l = y_i \end{cases}$$

Experiments

- Ran adaptive algorithm using minimal WLC
- Compared with AdaBoost.MI, AdaBoost.MH
- Tested on benchmark datasets
- Weak classifiers: bounded size decision trees









• What happens with multi-label / confidence rated weak classifiers?

- What happens with multi-label / confidence rated weak classifiers?
- Consistency of the algorithms.

- What happens with multi-label / confidence rated weak classifiers?
- Consistency of the algorithms.
- Extensions to ranking.

Thank you