#### Game Theory and Boosting

#### Manfred K. Warmuth

University of California, Santa Cruz

#### Game Theory Class, UCSC, March 3, 2009

Joint work and help from Gunnar Rätsch, Vishy Vishwanathan, Karen Glocer





3 What is Boosting?



#### Outline



2 A Machine learning problem

3 What is Boosting?

4 LPBoost and entropy regularized version

Zero-sum games

# <u>R</u>ock, <u>P</u>aper, <u>S</u>cissors game

	column player						
			R	Ρ	S		
			$\alpha_1$	$\alpha_2$	$\alpha_3$		
	R	$d_1$	0	1	-1		
row player	Ρ	$d_2$	-1	0	1		
	S	<i>d</i> <sub>3</sub>	1	-1	0		
		рау	off m	atrix			

payoff = 
$$\mathbf{d}^T M \boldsymbol{\alpha}$$
  
=  $\sum_{i,j} d_i M_{i,j} \alpha_j$ 

Zero-sum games

#### Two-player Zero Sum Game

- Gains of row player = losses of column player row player minimizes, column player maximizes
- Single row is pure strategy of row player and d is mixed strategy
- Single column is pure strategy of column player and *α* is mixed strategy

#### **Optimum Strategy**

			R	Ρ	S
			α <sub>1</sub> .33	α <sub>2</sub> .33	$lpha_{3}$ .33
R	$d_1$	.33	0	1	-1
Ρ	$d_2$	.33	-1	0	1
S	<b>d</b> 3	.33	1	-1	0

• Min-max theorem:

[Van Neumann 1928]

 $\min_{\mathbf{d}} \max_{\alpha} \mathbf{d}^{\mathsf{T}} M \alpha = \max_{\alpha} \min_{\mathbf{d}} \mathbf{d}^{\mathsf{T}} M \alpha$ = value of the game (0 in example)

#### Pure strategies

- e<sub>i</sub> pure strategy of row player
  - e<sub>j</sub> pure strategies of column player
- Inner strategy can be pure

$$\min_{\mathbf{d}} \max_{\alpha} \mathbf{d}^{\mathsf{T}} M \alpha = \min_{\mathbf{d}} \max_{j} \mathbf{d}^{\mathsf{T}} M \mathbf{e}_{j}$$
$$\max_{\alpha} \min_{\mathbf{d}} \mathbf{d}^{\mathsf{T}} M \alpha = \max_{\alpha} \min_{j} \mathbf{e}_{j}^{\mathsf{T}} M \alpha$$

All equal value of game

Zero-sum games

#### New column added



Value of game increases from 0 to .11

#### Row added



Value of game decreases from 0 to -.11

M.K.Warmuth et.al. (UCSC)

Game Theory and Boosting

Game Theory Class, UCSC, March 3, 2009 / 45

#### Incremental games

Column adding game (Boosting)

- Column player has large pool of columns available
- In each iteration one is added

#### 

iteration 1 iteration 2 iteration 3

- $\bullet\,$  In each iteration solve optimization problem to update d
- Column player always picks column which has largest edge wrt the current **d**

[FS

Zero-sum games

#### Desired Properties of Algorithm

- Want an algorithm that makes the value increase as quickly as possible
- Final game matrix should have value not too much smaller than optimum
- Number of columns needs should be as small as possible

#### Outline



2 A Machine learning problem

#### 3 What is Boosting?

4 LPBoost and entropy regularized version



- 11 apples (examples)
- labeled +1 if natural and -1 if artificial
- want to classify the apples
- given weak hypotheses decision stumps along the two features
- examples = rows
- weak hypotheses = possible columns



- 11 apples (examples)
- labeled +1 if natural and -1 if artificial
- want to classify the apples
- given weak hypotheses decision stumps along the two features
- examples = rows
- weak hypotheses = possible columns



- 11 apples (examples)
- labeled +1 if natural and -1 if artificial
- want to classify the apples
- given weak hypotheses decision stumps along the two features
- examples = rows
- weak hypotheses = possible columns



- 11 apples (examples)
- labeled +1 if natural and -1 if artificial
- want to classify the apples
- given weak hypotheses decision stumps along the two features
- examples = rows
- weak hypotheses = possible columns



- 11 apples (examples)
- labeled +1 if natural and -1 if artificial
- want to classify the apples
- given weak hypotheses decision stumps along the two features
- examples = rows
- weak hypotheses = possible columns

#### Examples and Hypotheses

M.K.Warmuth

Example	s Labels	h <sup>1</sup> : redness	
é	-1	-1	
Ć.	-1	-1	
é	-1	-1	
é	-1	1	mistake
۲	1	1	
	1	1	
9	1	1	
۲	1	-1	mistake
et.al. (UCSC)	Game The	ory and Boosting	/ 45

#### Boosting: 1st Iteration



First hypothesis: • error rate<sub>t</sub> =  $\frac{2}{11}$ =  $\sum_{i=1}^{n} d_i^t \mathbf{I}(h^t(\mathbf{x}_i) \neq y_i)$ • edge<sub>t</sub> =  $\frac{9}{22}$ =  $\sum_{i=1}^{n} d_i^t y_i h^t(\mathbf{x}_i)$ =  $1 - 2\epsilon_t$ 

#### Boosting: 1st Iteration



First hypothesis: • error rate  $t = \frac{2}{11}$  $=\sum_{i=1}^{n} d_i^t \mathbf{I}(\bar{h}^{\bar{t}}(\mathbf{x}_i) \neq y_i)$ •  $edge_t = \frac{9}{22}$  $=\sum_{i=1}^{n} \bar{d}_{i}^{t} y_{i} h^{t}(\mathbf{x}_{i})$  $= 1 - 2\epsilon_t$ Edge 0.5 1 Error -0.5 Rate

#### Connection to column adding game?

- Rows are the examples (fixed)
- Columns the weak hypotheses



- Column sum: edge of weak hypothesis
- Row sum: margin of example
- Value of game as large as possible

#### Game Matrix M

Example	s Labels	<i>h</i> <sup>1</sup> : redness	$M_{1,*}$	
Ś	-1	-1	1	
Ć	-1	-1	1	
é	-1	-1	1	
é	-1	1	-1	
۲	1	1	1	
۱	1	1	1	
<b>ö</b>	1	1	1	
<b>(</b>	- 1	-1	_1	
I. (UCSC)	⊥ Game Theory	and Boosting	▲ Game Theory Clas	ss, UCSC, March 3

M.K.Warmuth et.a

2009 / 45

## Margins and Edges

- The margin of example  $x_i$  at iteration t is summing *i*th row =  $\sum_{j=1}^{t} \alpha_j y_i h^j(x_i)$
- The edge of hypothesis  $h^j$  at iteration t is summing jth column =  $\sum_{i=1}^{n} d_i^t y_i h^j(x_i)$
- Example:

$$\begin{array}{cccc} h^1 & h^2 & h^3 \\ \alpha_1 & \alpha_2 & \alpha_3 \\ .2 & .4 & .4 \end{array}$$
 margin

#### Update Distribution



# $\begin{array}{l} \mbox{Misclassified examples} \\ \Rightarrow \mbox{Increased weights} \end{array}$

# fter update:Minimum edge small

#### Update Distribution



# $\begin{array}{l} \mbox{Misclassified examples} \\ \Rightarrow \mbox{Increased weights} \end{array}$

#### After update:

• Minimum edge small

#### Before 2nd Iteration



Hard examples

• High weight

#### Boosting: 2nd Hypothesis



Pick hypotheses with high edge

#### Update Distribution



After update: edge of all chosen hypotheses is small

#### Boosting: 3nd Hypothesis



### Boosting: 4th Hypothesis



#### All Hypotheses



Decision: 
$$f_{\alpha}(\mathbf{x}) = \sum_{t=1}^{T} \alpha_t h^t(\mathbf{x}) > 0$$
?



M.K.Warmu

#### Apple Classification Problem in Matrix Form

			$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_4$	margin	
			-	-	-	-		
Ś	$d_1$	1/11	1	-1	-1	1	-	
Ć	$d_2$	1/11	1	1	-1	1	-	
Ś	$d_3$	1/11	1	1	-1	1	-	
<b>É</b>	$d_4$	1/11	1	1	1	-1	-	
é	$d_5$	1/11	1	1	-1	1	-	
é	$d_6$	1/11	-1	1	1	-1	-	
٥	$d_7$	1/11	1	-1	1	1	-	
٥	$d_8$	1/11	1	1	1	1	-	
۱	$d_9$	1/11	1	1	1	1	-	
۲	$d_{10}$	1/11	1	1	1	1	-	
۲	$d_{11}$	1/11	-1	1	1	1	-	
	edge		.64	.64	.27	.64		
	value	-1				Game	Theory Class LICS	SC.
th et.al.	(UCSC)	Ga	me Theory	Ganle	Theory Class, OCC	, C,		

/ 45

March 3, 2009

M.K.Warmuth

#### Apple Classification Problem in Matrix Form

			$\alpha_1$	$\alpha_2$	$\alpha_{3}$	$lpha_{4}$	margin
			1	0	0	0	
Ś	$d_1$	0	1	-1	-1	1	1
Ć	$d_2$	0	1	1	-1	1	1
Ś	$d_3$	0	1	1	-1	1	1
ć	$d_4$	0	1	1	1	-1	1
<b>é</b>	$d_5$	0	1	1	-1	1	1
é	$d_6$	.5	-1	1	1	-1	-1
0	$d_7$	0	1	-1	1	1	1
٥	$d_8$	0	1	1	1	1	1
۵	$d_9$	0	1	1	1	1	1
۲	$d_{10}$	0	1	1	1	1	1
۲	$d_{11}$	.5	-1	1	1	1	-1
	edge		-1	1	1	0	
	value	-1	-1			c	Same Theory Class LICSC M
et.al. (UCSC)		Ga	ame Theo	ory and B	oosting		

, March 3, 2009 / 45

#### Apple Classification Problem in Matrix Form

			$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_4$	margin
			.5	.5	0	0	
Ś	$d_1$	.25	1	-1	-1	1	0
Ć	$d_2$	0	1	1	-1	1	1
Ś	$d_3$	0	1	1	-1	1	1
<u></u>	$d_4$	0	1	1	1	-1	1
<b>é</b>	$d_5$	0	1	1	-1	1	1
é	$d_6$	.25	-1	1	1	-1	0
Ø	$d_7$	.25	1	-1	1	1	0
ø	$d_8$	0	1	1	1	1	1
۲	$d_9$	0	1	1	1	1	1
۲	$d_{10}$	0	1	1	1	1	1
۲	$d_{11}$	.25	-1	1	1	1	0
	edge		0	0	.5	.5	
	value	-1	-1	0		62	me Theory Class, UCSC March
M.K.Warmuth et.al. (UCSC)		Ga	me Theor	y and Bo	osting	00	

3. 2009

45

M.K.Warmut

#### Apple Classification Problem in Matrix Form

		$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_4$	margin	
		.5	.18	.32	0		
🗯 d_1	.5	1	-1	-1	1	0	
ć d <sub>2</sub>	0	1	1	-1	1	.36	
🗯 d <sub>3</sub>	0	1	1	-1	1	.36	
🔹 d <sub>4</sub>	0	1	1	1	-1	1	
<b>≤</b> d₅	0	1	1	-1	1	.36	
<ul> <li><i>d</i><sub>6</sub></li> </ul>	.25	-1	1	1	-1	0	
) d	0	1	-1	1	1	.64	
ه d <sub>8</sub>	0	1	1	1	1	1	
🗎 d <sub>9</sub>	0	1	1	1	1	1	
d <sub>10</sub>	0	1	1	1	1	1	
ĕ d <sub>11</sub>	.25	-1	1	1	1	0	
edg	e	0	0	0	.5		
valu	<b>e</b> -1	-1	0	0	Gan	ne Theory Class, UCSC, M	arch
h et.al. (UCSC)	G	ame Theo	ory and Bo	osting		, , , , , , , , , , , , , , , , , , , ,	

45

M.K.Warmi

#### Apple Classification Problem in Matrix Form

			$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_4$	margin	
			.25	.44	.06	.25		
Ś	$d_1$	.5	1	-1	-1	1	0	
Ć	$d_2$	0	1	1	-1	1	.88	
ć	$d_3$	0	1	1	-1	1	.88	
Ć.	$d_4$	0	1	1	1	-1	.5	
é	$d_5$	0	1	1	-1	1	.88	
é	$d_6$	.5	-1	1	1	-1	0	
٥	$d_7$	0	1	-1	1	1	.12	
Ö	$d_8$	0	1	1	1	1	1	
۵	$d_9$	0	1	1	1	1	1	
۲	$d_{10}$	0	1	1	1	1	1	
۲	$d_{11}$	0	-1	1	1	1	.5	
	edge		0	0	0	0		
	value	-1	-1	0	0	<b>0</b> <sub>Gan</sub>	ne Theory Class, U(	SC. March 3, 2009
uth et.al. (	UCSC)		Game The	ory and B	oosting	Gui		/ 45

#### Outline



2 A Machine learning problem

#### 3 What is Boosting?

4 LPBoost and entropy regularized version

What is Boosting?

#### Boosting = greedy method for increasing margin

Converges to optimum marging w.r.t. all hypotheses



#### Want small number of iterations

M.K.Warmuth et.al. (UCSC)

Game Theory Class, UCSC, March 3, 2009 / 45

#### Assumption on next weak hypothesis

For current weighting of examples, oracle returns hypothesis of edge  $\geq g$ 

Goal

- For given  $\epsilon,$  produce convex combination of weak hypotheses with margin  $\geq g-\epsilon$
- Number of iterations  $O(\frac{\ln n}{\epsilon^2})$

#### Outline



2 A Machine learning problem

#### 3 What is Boosting?

4 LPBoost and entropy regularized version

#### LPBoost

# [GS98,RSS+00,DBST02]

Choose distribution that minimizes the maximum edge via LP

$$= \min_{\substack{\sum_n d_n=1, d_i \ge 0 \ q=1,2,...,t \\ \sum_n d_n=1, d_i \ge 0 \\ M_{*,q} \cdot \mathbf{d} \le c}} \max_{\substack{M_{*,q} \cdot \mathbf{d} \le c}} M_{*,q} \cdot \mathbf{d}$$

- Good practical boosting algorithm
- All weight is put on examples with minimum margin
- Brittle: iteration bound can be linear in number of examples *n* on malign artificial data sets [WGR07]

		$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_4$	$\alpha_5$	margin
		0	0	0	0	0	
$d_1$	.125	+1	95	93	91	99	—
$d_2$	.125	+1	95	93	91	99	—
$d_3$	.125	+1	95	93	91	99	—
$d_4$	.125	+1	95	93	91	99	_
$d_5$	.125	98	+1	93	91	+.99	—
$d_6$	.125	97	96	+1	91	+.99	_
d7	.125	97	95	94	+1	+.99	—
$d_8$	.125	97	95	93	92	+.99	_
edge		.0137	7075	6900	6725	.0000	
value	-1						

		$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_4$	$\alpha_5$	margin
		1	0	0	0	0	
$d_1$	0	+1	95	93	91	99	1
$d_2$	0	+1	95	93	91	99	1
$d_3$	0	+1	95	93	91	99	1
$d_4$	0	+1	95	93	91	99	1
$d_5$	1	98	+1	93	91	+.99	98
$d_6$	0	97	96	+1	91	+.99	97
$d_7$	0	97	95	94	+1	+.99	97
$d_8$	0	97	95	93	92	+.99	97
edge		98	1	93	91	.99	
value	-1	98					

		$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_4$	$\alpha_5$	margin
		0	1	0	0	0	
$d_1$	0	+1	95	93	91	99	95
$d_2$	0	+1	95	93	91	99	95
$d_3$	0	+1	95	93	91	99	95
$d_4$	0	+1	95	93	91	99	95
$d_5$	0	98	+1	93	91	+.99	1
$d_6$	1	97	96	+1	91	+.99	96
$d_7$	0	97	95	94	+1	+.99	95
$d_8$	0	97	95	93	92	+.99	95
edge		97	96	1	91	.99	
value	-1	98	96				

		$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_4$	$\alpha_5$	margin
		0	0	1	0	0	
$d_1$	0	+1	95	93	91	99	93
$d_2$	0	+1	95	93	91	99	93
$d_3$	0	+1	95	93	91	99	93
$d_4$	0	+1	95	93	91	99	93
$d_5$	0	98	+1	93	91	+.99	93
$d_6$	0	97	96	+1	91	+.99	1
$d_7$	1	97	95	94	+1	+.99	94
$d_8$	0	97	95	93	92	+.99	93
edge		97	95	94	1	.99	
value	-1	98	96	94			

		$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_4$	$\alpha_5$	margin
		0	0	0	1	0	
$d_1$	0	+1	95	93	91	99	91
$d_2$	0	+1	95	93	91	99	91
$d_3$	0	+1	95	93	91	99	91
$d_4$	0	+1	95	93	91	99	91
$d_5$	0	98	+1	93	91	+.99	91
$d_6$	0	97	96	+1	91	+.99	91
$d_7$	0	97	95	94	+1	+.99	1
$d_8$	1	97	95	93	92	+.99	92
edge		97	95	94	92	.99	
value	-1	98	96	94	92		

		$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_4$	$\alpha_{5}$	margin
		.5	.0026	0	0	.4975	
$d_1$	0.4974	+1	95	93	91	99	.0051
$d_2$	0	+1	95	93	91	99	.0051
$d_3$	0	+1	95	93	91	99	.0051
$d_4$	0	+1	95	93	91	99	.0051
$d_5$	0	98	+1	93	91	+.99	.0051
$d_6$	.4898	97	96	+1	91	+.99	.0051
d7	0	97	95	94	+1	+.99	.0051
$d_8$	.0127	97	95	93	92	+.99	.0051
edge		.0051	.0051	.9055	.9100	.0051	
value	-1	98	96	94	92	.0051	

# Entropy Regularized LPBoost

$$\min_{\sum_n d_n=1} \max_{q=1,2,\dots,t} M_{*,q} \cdot \mathbf{d} + \frac{1}{\eta} \Delta(\mathbf{d}, \mathbf{d}^0),$$

where regularizer  $\Delta(\mathbf{d}, \mathbf{d}^0)$  is relative entropy  $\sum_i d_i \ln \frac{d_i}{d_i^0}$ See visualization in part 2 of Lecture 3:

http://www.soe.ucsc.edu/classes/cmps290c/Spring07/

$$\mathbf{d}_n = \frac{\exp^{-\eta \text{ margin of example } n}}{Z} \qquad \text{"soft min"}$$

- Within  $\epsilon$  of maximum margin in  $O(\frac{\log n}{\epsilon^2})$  iterations
- Above form of weights first appeared in  $\nu$ -Arc algorithm [RSS+00]

M.K.Warmuth et.al. (UCSC)

# The effect of entropy regularization

#### Different distribution on the examples



LPBoost: lots of zeros sometimes  $\Omega(n)$  iterations



ERLPBoost: smoother distribution always  $O(\frac{\log n}{\epsilon^2})$  iterations

#### Conclusion

- Machine learning problems often modeled as games against nature or adversary
- Often end up with zero-sum games
- Its all about efficiency: Even against worst adversary we only need so much resources