Online Learning and Bregman Divergences

Tutorial at Machine Learning Summer School Taipei, Taiwan, July 2006

Manfred K. Warmuth

University of California at Santa Cruz, USA http://www.cse.ucsc.edu/~manfred

Help with the tutorial: Gunnar Rätsch



Content of this tutorial

- P I: Introduction to Online Learning
 - The Learning setting
 - Predicting as good as the best expert
 - Predicting as good as the best linear combination of experts
- P II: Bregman divergences and Loss bounds
 - Introduction to Bregman divergences
 - Relative loss bounds for the linear case
 - Nonlinear case & matching losses
 - Duality and relation to exponential families
 - On-line algorithms motivated by game theory
- P III: on-line to batch conversion, applications
 - Simple conversions
 - Caching and the disk spin down problem
 - $-\,$ Other applications and conclusion

Goal: How can we prove relative loss bounds?

 $\mathbf{2}$

Predicting almost as good as the best expert predic true E_1 E_2 E_3 E_n ... tion label loss 1 1 0 . . . 0 0 1 1 day 1 dav 20 0 1 0 1 1 1 . . . 0 1 0 day 3 1 1 . . . 1 1 day t $x_{t,1}$ $x_{t,2}$ $x_{t,3}$ $|y_t - \hat{y}_t|$. . . $x_{t,n}$ \hat{y}_t y_t Master Algorithm For t = 1 To T Do Get instance $x_t \in \{0, 1\}^n$ Predict $\hat{y}_t \in \{0, 1\}$ $y_t \in \{0, 1\}$ Get label Incur loss $|y_t - \hat{y}_t|$

More general online settings

Protocol:

F	or	t	=	1	То	T	Do	

Get instance	$oldsymbol{x}_t \in \mathbf{R}^n$
Predict	$\hat{y}_t \in \mathcal{Y}$
Get label	$y_t \in \mathcal{Y}$
Incur loss	$L(y_t, \hat{y}_t, \boldsymbol{x}_t)$

Problem instances:

- Classification: $\hat{y}_t, y_t \in \{0, 1\}$, e.g. $L(y, \hat{y}, \boldsymbol{x}) = |y \hat{y}|$
- Regression: $\hat{y}_t, y_t \in \mathbf{R}$, e.g. $L(y, \hat{y}, \boldsymbol{x}) = (y \hat{y})^2$
- Density estimation: no label, e.g. $L(y, \hat{y}, \boldsymbol{x}) = -\log P(\boldsymbol{x}|\theta)$

 $\mathbf{5}$

Goal: small total loss $\sum_{t} L(y_t, \hat{y}_t, \boldsymbol{x}_t)$

		A r	A run of the Halving Algorithm										
E_1	E_2	E_3	E_4	E_5	E_6	E_7	E_8	${}^{majo}_{rity}$	$_{label}^{true}$	loss			
1	1	0	0	1	1	0	0	1	0	1			
x	x	0	1	х	х	1	1	1	1	0			
x	x	х	1	х	х	0	0	0	1	1			
x	x	х	\uparrow	х	х	х	х						
		со	nsiste	ent									
⁷ or a Ialvi	ny seo ng Al	quenc gorith	e witł 1m ma	n a <mark>co</mark> akes a	nsiste it mos	ent exp st ≤ 1	$\operatorname{pert}_{\operatorname{Og}_2 n}$	mistako	es				
(Good	d bou	nd, b	ut no	t opti	mal)	<u> </u>	5 <u>57</u> 77	mistak					



where a, b are constants

Bounds loss of algorithm relative to loss of best expert





Relative Loss bound for Weighted Majority

Solving for M:

$$M \leq \frac{\ln \frac{1}{\beta}}{\ln \frac{2}{1+\beta}} M_i + \frac{1}{\ln \frac{2}{1+\beta}} \ln n$$
$$M \leq \frac{2.63}{i} \min_i M_i + \frac{2.63}{2.63} \ln n$$

For all sequences, loss of master algorithm is comparable to loss of best expert

 \Rightarrow Relative loss bounds

[Fr]

Other Loss Functions

absolute loss	$L(y, \hat{y})$	=	$ y-\hat{y} $
square loss	$L(y, \hat{y})$	=	$(y-\hat{y})^2$
entropic loss	$L(y, \hat{y})$	=	$y\ln\frac{y}{\hat{y}} + (1-y)\ln\frac{1-y}{1-\hat{y}},$
			$y, \hat{y} \in [0, 1]$
entropic loss \pm	$L(y, \hat{y})$	=	$\frac{1+y}{2}\ln\frac{1+y}{1+\hat{y}} + \frac{1-y}{2}\ln\frac{1-y}{1-\hat{y}},$
			$y, \hat{y} \in [-1, +1]$
hellinger loss	$L(y, \hat{y})$	=	$\frac{1}{2} \left(\sqrt{1-y} - \sqrt{1-\hat{y}} \right)^2 + \frac{1}{2} \left(\sqrt{y} - \sqrt{\hat{y}} \right)^2$
			$y, \hat{y} \in [0, 1]$

How does it work with other loss functions?One weight per expert:[V] $w_{t,i} = \beta^{L_{t,i}} = e^{-\eta L_{t,i}}$

[KW]

where $L_{t,i}$ is total loss of E_i before trial tand η is a positive learning rate

Master predicts with the weighted average (WA)

$$v_{t,i} = \frac{w_{t,i}}{\sum_{i=1}^{n} w_{t,i}} \text{ normalized weights}$$
$$\hat{y}_t = \sum_{i=1}^{n} v_{t,i} x_{t,i} = \boldsymbol{v}_t \cdot \boldsymbol{x}_t$$

where $x_{t,i}$ is the prediction of E_i in trial t

14

Wo	rst-case :	regre	t boun	\mathbf{ds}		
	b	WA	fancy			
	entropic	1	1			
	square	2	1/2			
	hellinger	1	.71			
T I I				c	·	
• Improved constat [V]	nts of $b = 1$	$/\eta$ whe	en Mastei	uses fa	ancier pi	red.
• For the discrete	loss and the	e absol	ute loss:	a > 1		

Bounds for other Loss Functions

13

 $\forall \text{ sequences } S \text{ of examples } \langle (\boldsymbol{x}_t, y_t) \rangle_{1 \leq t \leq T}$ $\text{where } \boldsymbol{x}_t \in [0, 1]^n \text{ and } y_t \in [0, 1]$ $L_{WA}(S) \leq \min_i \underbrace{1}_a L_i(S) + \underbrace{1/\eta}_b \ln(n)$ $\underbrace{L_{WA}(S) - \min_i L_i(S)}_{\text{regret}} \leq b \ln(n)$



Usefulness:

- Easy to combine many pretty good experts (algorithms) so that Master is guaranteed to be almost as good as the best
- Bounds logarithmic in number of experts (multiplicative updates)

Proof continued

$$L_{\text{WA}}(S) \leq -\frac{1}{\eta} \ln \frac{W_{T+1}}{W_1}$$

= $-\frac{1}{\eta} \ln \sum_{j=1}^n \frac{1}{n} e^{-\eta L_{T+1,j}(S)}$
 $\leq -\frac{1}{\eta} \ln \frac{1}{n} e^{-\eta L_i(S)}$ $i = \underset{j}{\operatorname{argmin}} L_{T+1,j}$
= $-\frac{1}{\eta} \ln \frac{1}{n} e^{-\eta L_{T+1,i}(S)}$
= $L_{T+1,i}(S) + \frac{1}{\eta} \ln n$

18

Questions:

Telescoping:

- How to obtain algorithms that do well compared to best linear combination or best thresholded linear combination of experts?
- How to motivate the updates?
- What are good measures of progress?
- What are good loss functions?
- Methods for proving relative loss bounds?

va	riables	s/expe	rts			
E_1	E_2	E_3	E_4	$\left egin{array}{c} true \\ label \end{array} ight $	$E_1 \vee E_3$	$E_3 \vee E_4$
1	1	0	0	0	1	0
1	0	1	0	1	1	1
0	1	1	1	0	1	1
0	1	0	0	1	0	0
$x_{t,1}$	$x_{t,2}$	$x_{t,3}$	$x_{t,4}$		\uparrow	Ť
					3	2
					mist	akes

21

The Perceptron Algorithm

In trial t: Get instance
$$\boldsymbol{x}_t \in \{0,1\}^n$$

If $\boldsymbol{w}_t \cdot \boldsymbol{x}_t \ge \theta$ then $\hat{y}_t = 1$
else $\hat{y}_t = 0$
Get label $y_t \in \{0,1\}$
If mistake then
 $\boldsymbol{w}_{t+1} = \boldsymbol{w}_t - \eta (\hat{y}_t - y_t) \boldsymbol{x}_t$

Rotation invariant if $\boldsymbol{w}_1 = (0, \ldots, 0)$

Weighted Majority on k-literal DisjunctionsOne expert per disjunction $\binom{n}{k}$ weightsDo as well as best k out of n literal (monotone) disjunction# of mistakes
of WM $\leq 2.63 M + 2.63 k \ln \frac{n}{k}$ M is # of mistakes of bestTime (and space) exponential in kEfficient algorithm have only one weight per literal
instead of one weight per disjunction

22

k-literal Disjunctions with Perceptron

Perceptron Convergence Theorem $(\boldsymbol{w}_1 = (0, \dots, 0), \theta = \frac{1}{2}, \eta = \frac{1}{2n})$

of mistakes $\leq 4 A + 4 k n$

where A is # of attribute errors of best disjunction of size k, i.e., the minimum # of attributes that need to be flipped to make the disjunction consistent

$A \leq kM$

Lower bound for rotation invariant algorithms:

[KWA]

#mistakes = $\Omega(n)$

The Winnow Algorithm [L]

In trial t: Get instance $x_t \in \{0,1\}^n$ If $w_t \cdot x_t \ge \theta$ then $\hat{y}_t = 1$ else $\hat{y}_t = 0$ Get label $y_t \in \{0,1\}$ If mistake then $w_{t+1,i} = w_{t,i} e^{-\eta (\hat{y}_t - y_t) x_{t,i}}$ Mistake bound $(w_1 := \frac{k}{n}(1, \dots, 1), \theta = \frac{3 \ln 3}{8}n, e^{-\eta} = \frac{1}{3})$ [AW] # of mistakes $\le 4 A + 3.6 k \ln \frac{n}{k}$ Not rotation invariant! 25

For $t = 1, \ldots, T$ do $oldsymbol{x}_t \in \mathbf{R}^n$ Get instance Predict $\hat{y}_t = \boldsymbol{w}_t \cdot \boldsymbol{x}_t$ Get label $y_t \in \mathbf{R}$ $L_t(\boldsymbol{w}_t) = (y_t - \hat{y}_t)^2$ Incur loss Update \boldsymbol{w}_t to \boldsymbol{w}_{t+1} y0 \hat{y}_t 0 x_t \boldsymbol{x} Assume comparison class $\{u\}$ is a set of linear predictors $oldsymbol{u}$: $oldsymbol{x} ightarrow oldsymbol{u} \cdot oldsymbol{x}$

So far

- Learning relative to best expert for various loss functions
- Learning relative to best disjunction
- Perceptron versus Winnow and expansion into feature space

26





Motivation of Updates [KW]

Gradient descent

$$\begin{split} \boldsymbol{w}_{t+1} &= \operatorname*{argmin}_{\boldsymbol{w}} \left(||\boldsymbol{w} - \boldsymbol{w}_t||_2^2 / 2 + \boldsymbol{\eta} (y_t - \boldsymbol{w} \cdot \boldsymbol{x}_t)^2 / 2 \right) \\ &= \boldsymbol{w}_t - \boldsymbol{\eta} (\underbrace{\boldsymbol{w}_{t+1} \cdot \boldsymbol{x}_t}_{\approx \boldsymbol{w}_t \cdot \boldsymbol{x}_t} - y_t) \, \boldsymbol{x}_t \end{split}$$

Exponentiated Gradient Algorithm

$$\begin{aligned} \boldsymbol{w}_{t+1} &= \arg \min_{\boldsymbol{w}} \left(\sum_{i=1}^{n} w_i \ln \frac{w_i}{w_{t,i}} + \eta (y_t - \boldsymbol{w} \cdot \boldsymbol{x}_t)^2 / 2 \right) \\ &= w_{t,i} \exp \left[-\eta \left(\underbrace{\boldsymbol{w}_{t+1} \cdot \boldsymbol{x}_t}_{\approx \boldsymbol{w}_t \cdot \boldsymbol{x}_t} - y_t \right) x_{t,i} \right] / \text{normalization} \end{aligned}$$

More examples of Updates (cont)

p-norm Algorithms $(- \mathbf{p}^n)$

 $(\boldsymbol{w} \in \mathbf{R}^n)$

$$\boldsymbol{w}_{t+1} = f^{-1} \big(f(\boldsymbol{w}_t) - \boldsymbol{\eta} \, \nabla L_t(\boldsymbol{w}_t) \big)$$

[GLS,GL]

where

$$f(\boldsymbol{w}) = \nabla \frac{1}{2} ||\boldsymbol{w}||_q^2 = \nabla \frac{1}{2} \left(\sum_i |w_i|^q \right)^{2/q}$$

and q dual to p (i.e., $\frac{1}{p} + \frac{1}{q} = 1$)

- p = 2 becomes gradient descent
- $p = O(\log n)$ becomes EG-like algs
- 2 interpolates between the two extremes

30

Ľ	amilies of u	pdate algorithms
narameter	name of	undate
"divergence"	family	algorithms
$\ m{w} - m{w}_t\ _2^2$	Gradient	Widrow Hoff (LMS)
	Descent	Linear Least Squares. Backpropagation Perceptron Algorithms kernel based algorithms,
$\sum_{i=1}^{n} w_i \ln \frac{w_i}{w_{t,i}}$	Exponentiated Gradient Algorithm	expert algs / Bayes update Normalized Winnow "AdaBoost"

parameter "divergence"	name of family	update algorithms
$\sum_{i=1}^{n} w_i \ln \frac{w_i}{w_{t,i}}$	Unnormalized	Winnow
$+w_{t,i}-w_i$	Exp. Grad. Alg.	
$\sum_{i=1}^{n} w_i \ln \frac{w_i}{w_{t,i}}$	Binary	
$+(1-w_i)\ln\frac{1-w_i}{1-w_{t,i}}$	Exp. Grad. Alg.	
any	-	-
"Bregman divergence"		
Members of different fa	amilies exhibit differ	ent behavior
$ _2^2 $ V	ersus	entropic regulariza
	33	

Discretization

$$\frac{\boldsymbol{f}(\boldsymbol{w}_{t+\boldsymbol{h}}) - \boldsymbol{f}(w_t)}{\boldsymbol{h}} = -\eta \nabla L_t(\boldsymbol{w}_t)$$
$$w_{t+\boldsymbol{h}} = \boldsymbol{f}^{-1} \left(\boldsymbol{f}(w_t) - \eta \ \boldsymbol{h} \nabla \boldsymbol{w} L_t(\boldsymbol{w}_t) \right)$$

We use h = 1

$$\boldsymbol{w}_{t+1} = \boldsymbol{f}^{-1} \left(\boldsymbol{f}(\boldsymbol{w}_t) - \eta_t \nabla L_t(\boldsymbol{w}_t) \right)$$

Conjecture: Forward Euler better: Replace $\nabla_{\boldsymbol{w}} L_t(\boldsymbol{w}_t)$ by $\nabla_{\boldsymbol{w}} L_t(\boldsymbol{w}_{t+\boldsymbol{h}})$

Alternate motivation: continuous updates [WJ]

 Continuous time
$$\Rightarrow$$
 Ordinary differential equations

 Gradient Descent

 $w \in \mathbb{R}^n$
 $\dot{w}_t = -\eta \nabla_w L_t(w_t)$

 Unnormalized Exponentiated Gradient Alg.

 $w \ge 0$
 $\dot{o}(w_t) = -\eta \nabla_w L_t(w_t)$

 34

 Image: Second Sec

• We now give a problem that favors the multiplicative updates

Linear versus

logarithmic dependence in n







The Kernel Trick[BGV92]If w linear combination of expanded instances, then $\hat{y} = \sum_{t} \alpha_t \phi(x_t) \cdot \phi(x) = \sum_{t} \alpha_t \phi(x_t) \cdot \phi(x)$ $\hat{y} = \underbrace{\sum_{t} \alpha_t \phi(x_t) \cdot \phi(x) = \sum_{t} \alpha_t \phi(x_t) \cdot \phi(x)}_{w}$ Kernel function $\mathbf{K}(x_t, x)$ often efficient to compute $\phi(\underbrace{(x_1, \dots, x_n)}_n) = \underbrace{(1, \dots, x_i, \dots, x_i x_j \dots, x_i x_j x_k \dots)}_{2^n \text{ products}}$ Kernel magic $K(x, z) = \phi(x) \cdot \phi(z) = \sum_{I \subseteq 1..n} \prod_{i \in I} x_i \prod_{i \in I} z_i = \prod_{i=1}^n (1 + x_i z_i) \bigoplus_{O(n) \text{ time}} 0$



Linear or non-linear ?

- We give a problem for which kernel algorithms behave like linear algorithms
- Embeddings don't help

:-(



Main Result Theorem • No matter how the instances are embedded • No matter what k training instances chosen by the learner • No matter what linear combination used For one of the targets average square loss on all n instances is $1 - \frac{k}{n}$ 46 Additional Constraints $w_i \ge 0$ and $\sum_{i=1}^n w_i = 1$ -1 +1 -1 +1 -1 +1 -1 +1Problem matrix -1 -1 +1 +1 -1 -1 +1 +1-1 -1 -1 -1 +1 +1 +1 +1• For above k instances, labeled by one of the 2^k columns, only consistent weight vector is unit identifying that column • Weight space can have rank 2^k











53

Questions

Gave problem that cannot be learned by kernel base algs

- What is the optimal kernal for a given problem?
- What is the hard problem for multiplicative updates?
- Can multiplicative updates be kernalized? Some cases are given by

[TW]

• Can entropy regularization be replaced by $|| ||_2^2$ regularization plus non-negativity and total weight constraints?

Which matrix?

Kernel matrix	dot products of instances							
Problem matrix	instances as rows - targets as columns							
• If eigen-spectrum then kernel not us	of kernel matrix has heavy tail seful							
- Picked wrong	kernel							
 Problem too hard 								
• If svd-spectrum of problem matrix has heavy tail then problem not learnable								
We showed:								
• Hadamard proble	m matrix has heavy tail							
• Adding random fe	eatures makes tail of kernel matrix heavy							
	54							
Co	ntent of this tutorial							

- P I: Introduction to Online Learning
 - The Learning setting
 - Predicting as good as the best expert
 - Predicting as good as the best linear combination of experts
- P II: Bregman divergences and Loss bounds
 - Introduction to Bregman divergences
 - $-\,$ Relative loss bounds for the linear case
 - Nonlinear case & matching losses
 - Duality and relation to exponential families
 - On-line algorithms motivated by game theory
- P III: on-line to batch conversion, applications
 - Simple conversions
 - $-\,$ Caching and the disk spin down problem
 - Other applications and conclusion
- Goal: How can we prove relative loss bounds?





Examples-3 Burg entropy $F(\boldsymbol{w}) = \sum_{i} -\ln w_{i}$ $f(\boldsymbol{w}) = -\frac{1}{\boldsymbol{w}}$ $\Delta_{F}(\widetilde{\boldsymbol{w}}, \boldsymbol{w}) = \sum_{i} \left(-\ln \frac{\widetilde{w_{i}}}{w_{i}} + \frac{\widetilde{w_{i}}}{w_{i}} \right) - n$

Examples-2 [GLS,GL] *p*-norm Algs (*q* is dual to *p*: $\frac{1}{p} + \frac{1}{q} = 1$) $F(\boldsymbol{w}) = \frac{1}{2}||\boldsymbol{w}||_{q}^{2}$ $f(\boldsymbol{w}) = \nabla \frac{1}{2}||\boldsymbol{w}||_{q}^{2}$ $\Delta_{F}(\boldsymbol{\tilde{w}}, \boldsymbol{w}) = \frac{1}{2}||\boldsymbol{\tilde{w}}||_{q}^{2} + \frac{1}{2}||\boldsymbol{w}||_{q}^{2} - \boldsymbol{\tilde{w}} \cdot f(\boldsymbol{w})$ When p = q = 2 this reduces to squared Euclidean distance (Widrow-Hoff).

62

Examples-4: div. between density matrices

Umegaki Divergence

$$F(\mathbf{W}) = \operatorname{tr} (\mathbf{W} \ln \mathbf{W} - \mathbf{W})$$
$$f(\mathbf{W}) = \ln \mathbf{W}$$

$$\Delta_F(\widetilde{\mathbf{W}}, \mathbf{W}) = \operatorname{tr} \left(\widetilde{\mathbf{W}}(\ln \widetilde{\mathbf{W}} - \ln \mathbf{W}) + \mathbf{W} - \widetilde{\mathbf{W}} \right)$$

LogDet Divergence

$$F(\mathbf{W}) = -\ln |\mathbf{W}|$$

$$f(\mathbf{W}) = \mathbf{W}^{-1}$$

$$\Delta_F(\widetilde{\mathbf{W}}, \mathbf{W}) = -\ln \frac{|\widetilde{\mathbf{W}}|}{|\mathbf{W}|} + \operatorname{tr}(\widetilde{\mathbf{W}}\mathbf{W}^{-1}) - n$$

General Motivation of Updates [KW]

Trade-off between two term:

$$\boldsymbol{w}_{t+1} = \operatorname*{argmin}_{\boldsymbol{w}}(\underbrace{\Delta_F(\boldsymbol{w}, \boldsymbol{w}_t)}_{weight \ domain} + \eta_t \underbrace{L_t(\boldsymbol{w})}_{label \ domain})$$

 $\Delta_F(\boldsymbol{w}, \boldsymbol{w}_t)$ is "regularization term" and serves as measure of progress in the analysis.

When loss L is convex (in \boldsymbol{w})

$$\nabla \boldsymbol{w}(\Delta_F(\boldsymbol{w}, \boldsymbol{w}_t) + \boldsymbol{\eta_t} L_t(\boldsymbol{w})) = 0$$

iff

$$f(\boldsymbol{w}) - f(\boldsymbol{w}_t) + \boldsymbol{\eta}_t \underbrace{\nabla L_t(\boldsymbol{w})}_{\approx \nabla L_t(\boldsymbol{w}_t)} = 0$$

$$\Rightarrow \quad \boldsymbol{w}_{t+1} = f^{-1} \left(f(\boldsymbol{w}_t) - \boldsymbol{\eta}_t \nabla L_t(\boldsymbol{w}_t) \right)$$

65









$$L_t(\boldsymbol{u}) \xrightarrow{\geq} L_t(\boldsymbol{w}_t) + (\boldsymbol{u} - \boldsymbol{w}_t) \cdot \underbrace{\nabla \boldsymbol{w} L_t(\boldsymbol{w}_t)}_{\text{update}}$$

$$= L_t(\boldsymbol{w}_t) - \frac{1}{\eta} \underbrace{(\boldsymbol{u} - \boldsymbol{w}_t) \cdot (f(\boldsymbol{w}_{t+1}) - f(\boldsymbol{w}_t))}_{\text{prop. 7 of } \Delta_F}$$

$$= L_t(\boldsymbol{w}_t) + \frac{1}{\eta} \left(\Delta_F(\boldsymbol{u}, \boldsymbol{w}_{t+1}) - \Delta_F(\boldsymbol{u}, \boldsymbol{w}_t) - \Delta_F(\boldsymbol{w}_t, \boldsymbol{w}_{t+1})\right)$$





Second step: Relate $\Delta_F(\boldsymbol{w}_t, \boldsymbol{w}_{t+1})$ to loss $L_t(\boldsymbol{w}_t)$

Loss & divergence are dependent

Get
$$\Delta_F(\boldsymbol{w}_t, \boldsymbol{w}_{t+1}) \leq \text{const.} L_t(\boldsymbol{w}_t)$$

Then solve for
$$\sum_{t} L_t(\boldsymbol{w}_t)$$

Yield bounds of the form

$$\sum_{t} L_t(\boldsymbol{w}_t) \leq a \sum_{t} L_t(\boldsymbol{u}) + b \,\Delta_F(\boldsymbol{u}, \boldsymbol{w}_1)$$

a, b constants, a > 1.

Regret bounds (a = 1):

time changing η , subtler analysis

[AG]



Bounds for Linear Regression with Square Loss

Gradient Descent

$$\sum_{t} L_t(\boldsymbol{w}_t) \le (1+c) \sum_{t} L_t(\boldsymbol{u}) + \frac{1+c}{c} X_2^2 U_2^2$$

 $||\boldsymbol{x}_t||_2 \leq X_2, ||\boldsymbol{u}||_2 \leq U_2, c > 0, \eta = f(c, X_2)$ Scaled Exponentiated Gradient

$$\sum_{t} L_t(\boldsymbol{w}_t) \le (1+c) \sum_{t} L_t(\boldsymbol{u}) + \frac{1+c}{c} \ln n X_{\infty}^2 U_1^2$$

 $||\boldsymbol{x}_t||_{\infty} \leq X_{\infty}, ||\boldsymbol{u}||_1 \leq U_1, c > 0, \eta = f(c, X_{\infty})$ *p*-norm Algorithm

$$\sum_{t} L_t(\boldsymbol{w}_t) \le (1+c) \sum_{t} L_t(\boldsymbol{u}) + \frac{1+c}{c} (p-1) X_p^2 U_q^2$$
$$||\boldsymbol{x}_t||_p \le X_p, ||\boldsymbol{u}||_q \le U_q, c > 0, \eta = f(c, X_\infty)$$





77

Idea behind the matching loss

If transfer function and loss match, then

$$\nabla \boldsymbol{w} \Delta_H(\boldsymbol{w} \cdot \boldsymbol{x}, h^{-1}(\boldsymbol{y})) = h(\boldsymbol{w} \cdot \boldsymbol{x}) - \boldsymbol{y}$$

Then update has simple form:

$$f(\boldsymbol{w}_{t+1}) = f(\boldsymbol{w}_t) - \eta_t (\boldsymbol{h}(\boldsymbol{w}_t \cdot \boldsymbol{x}) - \boldsymbol{y}_t) \boldsymbol{x}_t$$

This can be exploited in proofs

But not absolutely necessary One only needs convexity of $L(h(\boldsymbol{w} \cdot \boldsymbol{x}), y)$ in \boldsymbol{w}

$\Delta_{\mu}(\boldsymbol{w}\cdot\boldsymbol{x})$	$h^{-1}(u)$) as le	oss of \boldsymbol{w} on $(\boldsymbol{x}, \boldsymbol{y})$	
d matchin	$\log \log h$ for h	[AH	W,HKW]
thing loss i	is convex in ²	w	
ansfer f.	H(z)	match. loss	
h(z)		$d_H(oldsymbol{w}\cdotoldsymbol{x},h^{-1}(y)$	
~	$\frac{1}{2} \gamma^2$	$rac{1}{2}(oldsymbol{w}\cdotoldsymbol{x}-y)^2$	
2	$\overline{2}^{2}$	square loss	
		$\ln(1+e^{\boldsymbol{W}\cdot\boldsymbol{X}})-y\boldsymbol{w}\cdot\boldsymbol{x}$	
$\frac{e^z}{1+e^z}$	$\ln(1+e^z)$	$+y\ln y + (1-y)\ln(1-y)$	
		logistic loss	
ign(z)	2	$\max\{0, -y\boldsymbol{w}\cdot\boldsymbol{x}\}$	
1811(~)	~	hinge loss	
	$\Delta_H(\boldsymbol{w} \cdot \boldsymbol{x}, d \text{ matching loss})$ $\Delta_H($	$\begin{array}{c c} \Delta_{H}(\boldsymbol{w} \cdot \boldsymbol{x}, h^{-1}(\boldsymbol{y})) \text{ as left} \\ \text{d matching loss for } h \\ \text{thing loss is convex in ansfer f.} & H(z) \\ \hline \boldsymbol{x} & \frac{1}{2}z^{2} \\ \hline \frac{e^{z}}{1+e^{z}} & \ln(1+e^{z}) \\ \hline \text{ign}(z) & z \end{array}$	$\begin{array}{c c} \Delta_{H}(\boldsymbol{w}\cdot\boldsymbol{x},h^{-1}(y)) \text{ as loss of } \boldsymbol{w} \text{ on } (\boldsymbol{x},y) \\ \text{d matching loss for } h & [AH] \\ \text{thing loss is convex in } \boldsymbol{w} \\ \text{ansfer f.} & H(z) & \text{match. loss} \\ h(z) & d_{H}(\boldsymbol{w}\cdot\boldsymbol{x},h^{-1}(y)) \\ \hline z & \frac{1}{2}z^{2} & \frac{1}{2}(\boldsymbol{w}\cdot\boldsymbol{x}-y)^{2} \\ & \text{square loss} \\ \hline \frac{e^{z}}{1+e^{z}} & \ln(1+e^{z}) & \ln(1+e^{\boldsymbol{w}\cdot\boldsymbol{x}}) - y\boldsymbol{w}\cdot\boldsymbol{x} \\ & +y\ln y + (1-y)\ln(1-y) \\ & \log \text{istic loss} \\ \hline \text{ign}(z) & z & \max\{0,-y\boldsymbol{w}\cdot\boldsymbol{x}\} \\ & \text{hinge loss} \end{array}$

78



[Ce]



Duality

Special case:

$$\min_{\boldsymbol{w}} \Delta_F(\boldsymbol{w}, \boldsymbol{w}_t) + \Delta_H(\boldsymbol{x}_t \cdot \boldsymbol{w}, h^{-1}(\boldsymbol{y}_t)) = -\min_{\alpha} \Delta_{\mathcal{G}}(\alpha + y_t, h(0)) + \Delta_{\mathcal{F}}(f(\boldsymbol{w}_t) - \alpha \boldsymbol{x}_t, f(\mathbf{0})) + \text{const.}$$

General:

$$\min_{\boldsymbol{w}} \Delta_F(\boldsymbol{w} + \boldsymbol{\mu}, \underbrace{f^{-1}(\boldsymbol{\phi})}_{\boldsymbol{w}_t}) + \Delta_H(X\boldsymbol{w} + \boldsymbol{\nu}, h^{-1}(\boldsymbol{y}))$$
$$= -\min_{\boldsymbol{\alpha}} \Delta_{\mathcal{G}}(\boldsymbol{\alpha} + \boldsymbol{y}, h(\boldsymbol{\nu})) + \Delta_{\mathcal{F}}(\boldsymbol{\phi} - X^{\top}\boldsymbol{\alpha}, f(\boldsymbol{\mu})) + \text{const.}$$

where F and \mathcal{F} are convex conjugate functions:

$$\mathcal{F}(\boldsymbol{x}) = \sup_{\boldsymbol{y}} \boldsymbol{x} \cdot \boldsymbol{y} - F(\boldsymbol{y}) = \boldsymbol{x} \cdot (\nabla F)^{-1}(\boldsymbol{x}) - F((\nabla F)^{-1}(\boldsymbol{x}))$$

85

Relation to Boosting

The AdaBoost update of the probability vector \boldsymbol{w}_t :

$$w_i^{(t+1)} = w_i^{(t)} \exp(-\alpha_t y_i h_t(\boldsymbol{x}_i))$$

Is a projection w.r.t. divergence

[CKW,La,KW,CSS]

$$\Delta_F(\boldsymbol{w}, \boldsymbol{w}_t) = \sum_i w_i \ln \frac{w_i}{w_{t,i}}$$

Such that the weighted training error of h_t w.r.t. $\boldsymbol{w}^{(t+1)}$ is $\frac{1}{2}$ ("diversification" of Boosting mentioned in Ron Meir's talk)

Projections onto Hyperplanes

$$oldsymbol{w}_{t+1} = \operatorname*{argmin}_{oldsymbol{w}} (\Delta_F(oldsymbol{w}, oldsymbol{w}_t) + oldsymbol{\eta} (oldsymbol{w} \cdot oldsymbol{x}_t - y_t)^2)$$

When η is large then w_{t+1} is projection of w_t onto plane $w \cdot x_t = y_t$







Exponential Family of Distributions

• Parametric density functions

$$P_G(\boldsymbol{x}|\boldsymbol{\theta}) = e^{\boldsymbol{\theta} \cdot \boldsymbol{x} - \boldsymbol{G}(\boldsymbol{\theta})} P_0(\boldsymbol{x})$$

- $\boldsymbol{\theta}$ and \boldsymbol{x} vectors in R^d
- Cumulant function $G(\theta)$ assures normalization

$$G(\boldsymbol{ heta}) = \ln \int e^{\boldsymbol{ heta} \cdot \boldsymbol{x}} P_0(\boldsymbol{x}) \ d\boldsymbol{x}$$

- $G(\boldsymbol{\theta})$ is convex function on convex set $\boldsymbol{\Theta} \subseteq R^d$
- *G* characterizes members of the family
- $\boldsymbol{\theta}$ is *natural* parameter



• Expectation parameter

$$\boldsymbol{\mu} = \int_{\boldsymbol{x}} \boldsymbol{x} P_G(\boldsymbol{x}|\boldsymbol{\theta}) d\boldsymbol{x} = E_{\boldsymbol{\theta}}(\boldsymbol{x}) = g(\boldsymbol{\theta})$$

where $g(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} G(\boldsymbol{\theta})$

• Second convex function $F(\boldsymbol{\mu})$ on space $g(\boldsymbol{\Theta})$

$$F(\boldsymbol{\mu}) = \boldsymbol{\theta} \cdot \boldsymbol{\mu} - G(\boldsymbol{\theta})$$

- $G(\boldsymbol{\theta})$ and $F(\boldsymbol{\mu})$ are *convex conjugate* functions
- Let $f(\boldsymbol{\mu}) = \nabla \boldsymbol{\mu} F(\boldsymbol{\mu})$
- $f(\boldsymbol{\mu}) = g^{-1}(\boldsymbol{\mu})$



Bernoulli

Examples x_t are coin flips in $\{0, 1\}$

 $P(x|\mu) = \mu^{x} (1-\mu)^{1-x}$

 μ is the probability (expectation) of 1 Natural parameter: $\theta = \ln \frac{\mu}{1-\mu}$

$$P(x|\theta) = \exp\left(\theta x - \ln(1 + e^{\theta})\right)$$

Cumulant function: $G(\theta) = \ln(1 + e^{\theta})$

Parameter transformations:

$$\mu = g(\theta) = \frac{e^{\theta}}{1 + e^{\theta}}$$
 and $\theta = f(\mu) = \ln \frac{\mu}{1 - \mu}$

Dual function: $F(\mu) = \mu \ln \mu + (1 - \mu) \ln(1 - \mu)$

Log loss:
$$L_t(\theta) = -x_t \theta + ln(1+e^{\theta})$$

= $-x_t \ln \mu - (1-x_t) \ln(1-\mu)$

 $\begin{array}{l} \textbf{Gaussian (unit variance)} \\ P(\boldsymbol{x}|\boldsymbol{\theta}) &\sim e^{-\frac{1}{2}(\boldsymbol{\theta}-\boldsymbol{x})^2} \\ &= e^{\boldsymbol{\theta}\cdot\boldsymbol{x}-\frac{1}{2}\boldsymbol{\theta}^2}e^{\frac{1}{2}\boldsymbol{x}^2} \end{array}$ Cumulant function: $G(\boldsymbol{\theta}) = \frac{1}{2}\boldsymbol{\theta}^2$ Parameter transformations: $g(\boldsymbol{\theta}) = \boldsymbol{\theta} = \boldsymbol{\mu} \quad \text{and} \quad f(\boldsymbol{\mu}) = \boldsymbol{\mu} = \boldsymbol{\theta}$ Dual convex function: $F(\boldsymbol{\mu}) = \boldsymbol{\theta} \cdot \boldsymbol{\mu} - G(\boldsymbol{\theta}) \\ &= \frac{1}{2}\boldsymbol{\mu}^2$ Square loss: $L_t(\boldsymbol{\theta}) = \frac{1}{2}(\boldsymbol{\theta}_t - \boldsymbol{x}_t)^2$

94

Poisson

Examples x_t are natural numbers in $\{0, 1, \ldots\}$

$$P(x|\mu) = \frac{e^{-\mu}\mu^2}{x!}$$

 μ is expectation of xNatural parameter: $\theta = \ln \mu$

 $P(x|\theta) = \exp\left(\theta x - e^{\theta}\right) \frac{1}{x!}$

Cumulant function: $G(\theta) = e^{\theta}$

Parameter transformations:

$$\mu = g(\theta) = e^{\theta}$$
 and $\theta = f(\mu) = \ln \mu$

Dual function: $F(\mu) = \mu \ln \mu - \mu$ Loss: $L_t(\theta) = -x_t \theta + e^{\theta} + \ln x_t!$

$$= -x_t \ln \mu + \mu + \ln x_t!$$

Bregman Div. as Rel. Ent. between Distributions

Let $P(x|\theta)$ and $P(x|\tilde{\theta})$ denote two distributions with cumulant function G

$$\begin{split} \Delta_{G}(\widetilde{\boldsymbol{\theta}},\boldsymbol{\theta}) &= \int_{\boldsymbol{x}} P_{G}(\boldsymbol{x}|\boldsymbol{\theta}) \ln \frac{P_{G}(\boldsymbol{x}|\boldsymbol{\theta})}{P_{G}(\boldsymbol{x}|\widetilde{\boldsymbol{\theta}})} d\boldsymbol{x} \\ &= \int_{\boldsymbol{x}} P_{G}(\boldsymbol{x}|\boldsymbol{\theta}) (\boldsymbol{\theta} \cdot \boldsymbol{x} - G(\boldsymbol{\theta}) - \boldsymbol{\theta} \cdot \boldsymbol{x} + G(\widetilde{\boldsymbol{\theta}})) d\boldsymbol{x} \\ &= G(\widetilde{\boldsymbol{\theta}}) - G(\boldsymbol{\theta}) - (\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}) \cdot (\int_{\boldsymbol{x}} P_{G}(\boldsymbol{x}|\boldsymbol{\theta}) \boldsymbol{x} d\boldsymbol{x}) \\ &= G(\widetilde{\boldsymbol{\theta}}) - G(\boldsymbol{\theta}) - G(\boldsymbol{\theta}) - (\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}) \cdot \boldsymbol{\mu} \\ F(\boldsymbol{\mu}) = \boldsymbol{\theta} \cdot \boldsymbol{\mu}^{-G}(\boldsymbol{\theta}) \quad F(\boldsymbol{\mu}) - F(\widetilde{\boldsymbol{\mu}}) - (\boldsymbol{\mu} - \widetilde{\boldsymbol{\mu}}) \cdot \widetilde{\boldsymbol{\theta}} \\ &= \Delta_{F}(\boldsymbol{\mu}, \widetilde{\boldsymbol{\mu}}) \quad [A, BN, AW] \end{split}$$

$$egin{aligned} &\Delta_G(heta, \widetilde{ heta}) &= & G(heta) - G(\widetilde{ heta}) - (heta - \widetilde{ heta}) \cdot g(\widetilde{ heta}) \ &= & \int_{\widetilde{ heta}}^{eta} (g(au) - g(\widetilde{ heta})) \cdot d au \ &= & \int_{\mu}^{\widetilde{eta}} (f(\sigma) - f(\mu)) \cdot d\sigma \ &= & F(\widetilde{\mu}) - F(\mu) - (\widetilde{\mu} - \mu) \cdot f(\mu) \ &= & \Delta_F(\widetilde{\mu}, \mu) \end{aligned}$$





for BEG

Dual divergence for Poisson

$$G(\theta) = e^{\theta}$$
 $F(\mu) = \mu \ln \mu - \mu$

$$g(\theta) = e^{\theta} = \mu \quad f(\mu) = \ln \mu = \theta$$

$$\Delta_G(\widetilde{\theta}, \theta) = e^{\widetilde{\theta}} - e^{\theta} - (\widetilde{\theta} - \theta)e^{\theta}$$

$$\Delta_F(\mu,\widetilde{\mu}) = \mu \ln \frac{\mu}{\widetilde{\mu}} + \widetilde{\mu} - \mu$$

Unnormalized relative entropy

Sum of unnormalized relative entropies is parameter for UEG (e.g. Winnow)

101



Dual matching loss for sigmoid transfer func.

$$H(z) = \ln(1 + e^z) K(r) = r \ln r + (1 - r) \ln(1 - r)$$

$$h(z) = \frac{e^z}{1 + e^z} = r k(r) = \ln \frac{r}{1 - r} = z$$

K dual to H and $k = h^{-1}$

$$\Delta_H(\boldsymbol{w}\cdot\boldsymbol{x},h^{-1}(y))$$

= ln(1 + e^{\boldsymbol{w}\cdot\boldsymbol{x}}) - y\boldsymbol{w}\cdot\boldsymbol{x} + y\ln y + (1-y)\ln(1-y)

By duality logistic loss is same as entropic loss

$$\Delta_K(y, h(\boldsymbol{w} \cdot \boldsymbol{x}))$$

= $y \ln \frac{y}{h(\boldsymbol{w} \cdot \boldsymbol{x})} + (1 - y) \ln \frac{1 - y}{1 - h(\boldsymbol{w} \cdot \boldsymbol{x})}$

Matching loss for logistic transfer function

102

Derivation of Updates

• Want to bound

$$\sum_{t=1}^{T} L_t(\boldsymbol{\theta}_t) - \inf_{\boldsymbol{\theta}} L_{1..T}(\boldsymbol{\theta})$$

• Off-line algorithm has all T examples

$$\{oldsymbol{x}_1,oldsymbol{x}_2,\ldots,oldsymbol{x}_T\}$$

• Setup for choosing best parameter setting

Here $\eta_B^{-1} > 0$ is a tradeoff parameter

On-line Algorithm [AW]

• In trial t, the first t examples

$$\{oldsymbol{x}_1,oldsymbol{x}_2,\ldots,oldsymbol{x}_t\}$$

have been presented

- Motivation for on-line parameter update: do as well as best off-line algorithm up to trial t
- At end of trial t algorithm minimizes

 $\begin{aligned} \boldsymbol{\theta}_{t+1} &= \operatorname*{argmin}_{\boldsymbol{\theta}} \left(\eta_1^{-1} \ \Delta_G(\boldsymbol{\theta}, \boldsymbol{\theta}_1) \ + \ L_{1..t}(\boldsymbol{\theta}) \right) \\ \boldsymbol{\theta} & \operatorname{divergence} & \operatorname{loss} \\ & \operatorname{to initial} & \operatorname{so far} \end{aligned}$

Tradeoff parameter $\eta_1^{-1} \ge 0$

105

Alternate Motivation of Same On-Line Update $\theta_{t+1} = \underset{\theta}{\operatorname{argmin}} (\eta_t^{-1} \quad \Delta_G(\theta, \theta_t) + L_t(\theta))$ divergence to last current loss where $\eta_t = \frac{1}{\eta_1^{-1} + t - 1}$ Parameter Updates Off-line: $\mu_B = \frac{\eta_B^{-1} \mu_1 + \sum_{t=1}^T x_t}{\eta_B^{-1} + T}$ On-Line in trial t: $\mu_{t+1} = \frac{\eta_1^{-1} \mu_1 + \sum_{q=1}^t x_q}{\eta_1^{-1} + t} = \mu_t - \eta_{t+1}(\mu_t - x_t)$ $\theta_{t+1} = g^{-1} (g(\theta_t) - \eta_{t+1}(\mu_t - x_t))$

106

- On-line algorithm has freedom to use a tradeoff parameter η_1^{-1} that could be different from the off-line parameter η_B^{-1}
- Two choices for η_1^{-1}
 - Case $\eta_1^{-1} = \eta_B^{-1}$:

Incremental Off-Line Algorithm

Case $\eta_1^{-1} = \eta_B^{-1} + 1$: Forward Algorithm [V] Shrinkage Towards Initial

$$\boldsymbol{\mu}_B = \overline{\boldsymbol{x}_T} - \eta_B^{-1} (\eta_B^{-1} + T)^{-1} (\overline{\boldsymbol{x}_T} - \boldsymbol{\mu}_1)$$

where $\overline{\boldsymbol{x}_T} = \frac{\sum_{t=1}^T \boldsymbol{x}_t}{T}$ Shrinkage factor $\eta_B^{-1}(\eta_B^{-1} + T)^{-1}$







Why Bregman divergences?

- No need to check whether there is an underlying exponential family
- More general than exponential families
- As parameter divergence and matching loss
- Used in motivation and analysis of updates
- When $\eta \to \infty$, updates morph into Bregman projection
- Generalized Pythagorean Theorem for Bregman projections

General setup of on-line learning

- We hide some information from the learner
- The relative loss bound quantifies the price for hiding the information

114

• So far the future examples are hidden Off-line algorithm knows all examples On-line algorithm knows past examples

113

Minimax Algorithm for T Trials

Learner against adversary

 $\inf_{\boldsymbol{\mu}_1} \sup_{\boldsymbol{x}_1} \inf_{\boldsymbol{\mu}_2} \sup_{\boldsymbol{x}_2} \inf_{\boldsymbol{\mu}_3} \sup_{\boldsymbol{x}_3} \ldots \inf_{\boldsymbol{\mu}_T} \sup_{\boldsymbol{x}_T}$

$\begin{array}{rcl} \sum_{t=1}^{T} \frac{1}{2} (\boldsymbol{\mu}_t - \boldsymbol{x}_t)^2 & - & \inf \boldsymbol{\mu} \sum_{t=1}^{T} \frac{1}{2} (\boldsymbol{\mu} - \boldsymbol{x}_t)^2 \\ \text{total loss of} & \text{total loss of} \\ \text{on-line} & & \text{off-line} \\ \text{algorithm} & & \text{algorithm} \\ \text{Instances must be bounded:} & ||\boldsymbol{x}_t||_2 \leq X \\ \text{Minimax algorithm usually intractable} \end{array}$

Gaussian and Bernoulli are exceptions



[TW,Sh]

Last-step Minimax
Assumes that current trial is last trial [Fo,TW]

$$\mu_{t} = \underset{\mu}{\operatorname{arginf}} \sup_{x_{t}} \sum_{q=1}^{t} L_{q}(\mu_{q}) - \inf_{\mu} L_{1..t}(\mu)$$

$$= \underset{\mu}{\operatorname{arginf}} \sup_{x_{t}} L_{t}(\mu_{t}) - \inf_{\mu} L_{1..t}(\mu)$$
For Gaussian and linear regression
Last-step Minimax is same as Forward Alg.
117
IIT
Synopsis of methods
Game theoretic
• Slightly better bounds

• Harder to find

Bregman divergences

• Closer to Bayes and standard convex optimization

Last-step Minimax: Bernoulli

 Forward alg:

$$\mu_t = \frac{s + \frac{1}{2}}{t - 1 + 1}$$
, where $s = \sum_{q=1}^{t-1} x_q$

 Last-step:

 $\mu_t = \frac{(s + 1)^{(s + 1)}(t - s - 1)^{t - s - 1}}{s^s(t - s)^{t - s} + (s + 1)^{s + 1}(t - s - 1)^{t - s - 1}}$

 Worst-case regret bounds: $\ln(T + 1) + \mathbf{c}$

 Forward: $\mathbf{c} = \frac{\ln \pi}{2}$

 Last step: $\mathbf{c} = \frac{1}{2}$

 118

 Content of this tutorial

 • P I: Introduction to Online Learning

 - The Learning setting

- Predicting as good as the best expert
- Predicting as good as the best linear combination of experts
- P II: Bregman divergences and Loss bounds
 - Introduction to Bregman divergences
 - $-\,$ Relative loss bounds for the linear case
 - Nonlinear case & matching losses
 - Duality and relation to exponential families
 - On-line algorithms motivated by game theory
- P III: on-line to batch conversion, applications
 - Simple conversions
 - Caching and the disk spin down problem
 - Other applications and conclusion
- Goal: How can we prove relative loss bounds?

Simple conversions

Worst case loss bounds for on-line algs

are converted to

algorithms with good performance bounds in the i.i.d. case

• Expected loss bounds [HW,CB+,KW]

121

• Tail bounds

[CCG]

Expected loss bounds [HW]

Loss function $L : \mathbf{R}^2 \to \mathbf{R}$ $S = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T) \sim D^T$ Instantaneous loss of hypothesis h w.r.t. distribution DInstLoss $(h, D) = E_{e \sim D} L(h, e)$

122

$$\begin{split} \mathbf{E}_{S \sim D^{T}}(\text{TotLoss}(A, S)) \\ &= \mathbf{E}_{(e_{1}, \dots, e_{T}) \sim D^{T}}\left(\sum_{t=1}^{T} L(A((e_{1}, \dots, e_{t-1})), e_{t})\right) \\ &= \sum_{t=1}^{T} \mathbf{E}_{(e_{1}, \dots, e_{t-1}) \sim D^{t-1}}\left(\mathbf{E}_{e \sim D}\left(L(A((e_{1}, \dots, e_{t-1})), e))\right) \\ &= \sum_{t=1}^{T} \mathbf{E}_{(e_{1}, \dots, e_{t-1}) \sim D^{t-1}}\left(\text{InstLoss}(A(e_{1}, \dots, e_{t-1}), D))\right) \\ \end{split}$$
So expected total loss is total instaneous loss



- Choose h_i uniformly at random from the T + 1 hypotheses
 On new instance x predict with h_i(x)
 Instanteneous loss of this algorithm is expected total loss of original algorithm over T + 1
- Applied to the Perceptron Algorithm [FS]





0.5

0.4

0.3

0.2 0.1

0

205000

210000

215000



220000

225000

230000

gd

235000





Better Master Policy

- Use Expert Framework from On-line learning
- Maintain one weight w_i for each base policy / expert
- w_i is estimate of current relative performance of policy i
- Weights updated after each request:
 - Loss update punishes policies quickly that score misses

141

Share update [LW94,HW98,BW01]
 Keeps weights of poor policies from becoming too small
 Helps recovery

Fixed Share to Uniform Past

Loss Update:

$$w'_{t,i} = \frac{w_{t,i}\beta^{\mathrm{miss}_{t,i}}}{\mathrm{normaliz}}, \quad \beta \in (0,1)$$

Share Update:

$$\mathbf{w}_{t+1} = (1 - \alpha) \mathbf{w}'_t + \alpha \mathbf{r}_{t-1},$$

where $\mathbf{r}_{t-1} = \sum_{q=1}^{t-1} \mathbf{w}'_q / (t-1)$

• Prevents weights that did well in past from becoming too small Helps when these weights need to recover

142





On-line Learning											
		\exp	erts								
	E_1	E_2	E_3	E_n	tion	label	loss				
day 1	1	1	0	0	0	1	1				
day 2	1	0	1	0	1	0	1				
day 3	0	1	1	1	1	1	0				
day t	$x_{t,1}$	$x_{t,2}$	$x_{t,3}$	$x_{t,n}$	\hat{y}_t	y_t	$(y_t - \hat{y}_t)^2$				
• Cho	ose co	mparis	on cla	ss of p	redictor	s (expe	rts)				
• Mas	ter Alg	gorithr	n com	bines p	oredictic	ons of ex	perts				
• X4 V	ector (of expe	ert's pr	redictio	ns						
	00001	on pe	itos pi	carotro							
				1	45						
				Ţ	40						
WI	not k	ind c	of not	rform	2220	<u></u>	o ovport?				
VV I	lat K	mu c	n per	IOIII	lance	call w	e expect :				
T	- bo	tha ta	tal log	n of al	rorithm	Δ					
- L ₁₇	r, <u>A</u> be		11		gor rennin	A F					
- L_{17}	r, <mark>i</mark> be t	the tot	al loss	of <i>i</i> -th	ı expert	E_i					
• Form	n of b	ounde									
• 1'011		Junus				/	λ.				
		$\forall S$:	: L_{1}	$T, A \leq$	min	$L_{1T,i}$	$+ c \log n$				

where c is constant

• Bounds the loss of the algorithm relative to the loss of best expert



• The weights are updated according to the Loss Update

$$w_{t+1,i} := rac{w_{t,i} \ e^{-\eta \ L_{t,i}}}{ ext{normaliz.}}, \ e^{-\eta} = eta$$

where $L_{t,i}$ is loss of expert *i* in trial *t*

- \rightarrow Weighted Majority Algorithm [LW89]
- \rightarrow Generalized by Vovk [Vovk90]





























Master Policy Protocol

- Process request on virtual caches
- Apply Loss and Share Updates
- Process request on real cache
 - based on combined weightings of all caches
- Refetch objects into real cache (if desired)

Virtual Cache Rankings

• Priorities induce ranks over virtually cached objects:

object	<i>o</i> ₁₂	07	02	022	03	06	02	<i>o</i> ₁₅	09
priority	31.2	30.2	24.1	17.1	9.3	8	4.1	2.5	1.2
rank	9	8	7	6	5	4	3	2	1

• o_9 first discarded

• o_{12} last

Master Rank

• Master priority *P* constructed from weights and ranks of virtual caches

$$P_o = \begin{cases} \sum_{n:o \in \mathrm{VC}_n} w_n r_{n,o} & \text{if } \exists n:o \in \mathrm{VC}_n \\ 0 & \text{if } \forall n:o \notin \mathrm{VC}_n \end{cases}$$

166

• R is corresponding master rank

165



Managing Real Cache: Instantaneous Rollover
Keep RealCache = Ideal Cache

i.e. Refetch all o ∈ IdealCache – RealCache

Too much refetching











Refetches as % of Total Requests

n

Summary

- Demand Rollover is already as good or better than BestFixed
- Small amounts of refetching always beats Best Fixed
 - 15-22% fewer misses than BestFixed
 - 45-70% fewer misses than LRU
- Can be as good as BestRefetching
 - always less I/O's than LRU
 - can result in less I/O than BestFixed

177

Conclusion

- Operating Systems have many parameter tweaking problems suitable for on-line learning
- Previous work using same updates:
 - Tuning time-out for spinning down disk of a PC [HLSS00]
 - Load balancing between processors [BB97]
 - Tracking with GPS

benchmark data sets

178

Too expensive?

- Not for web caching and filesystem's caching
- Not clear for paging
- Implement in Linux kernel

Two approaches Use existing caching strategies as experts Use set of fine-grained experts from which all existing caching policies are built Machine Learners will get interested if there are realistic

Application: Disk Spin Down [HLSS]

Problem of adapt. spinning down hard disks in mobile computers

Common approach

- Fixed time-out (e.g. 2 min)
- Does not exploit changing usage patterns

181







Idea

- Use about 20 experts with different time-outs
- Apply shifting expert algorithm with mixing to decaying past
- Efficient but proofs don't apply because on unusual loss function

Does it work?

Comparators:

- Best fixed idle time chosen in hindsight
- Optimal algorithm: Spin down if cost of next idle period > spin down cost

Performance

- Better than best fixed
- Close to optimal
- Parameters easy to tune and algorithm very stable over a large variety of data
- Better than other algorithms that provable have good competitve ratios

Other Applications

• Calendar managing	
Many features (sleeping experts)	[Bl,FSSW]
• Text categorization	[LSCP]
One attribute per word in text	
• Spelling correction	[Ro]
• Portfolio prediction	[Co,CO,HSSW,BK]
• Boosting	[Sc, Fr, SS]
• Load Balancing based on shifting expert algorithms [BE	

185

186