

Proving Relative Loss Bounds for On-Line Learning Algorithms Using Bregman Divergences

Claudio Gentile
Universita' di Milano, Italy

Manfred K. Warmuth
UC Santa Cruz, USA

July 13, 2000

Early contributors:

Nick Littlestone
Volodya Vovk
Tom Cover

On-Line Learning

experts					prediction	<i>true label</i>	loss
	E_1	E_2	E_3	E_n			
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 $

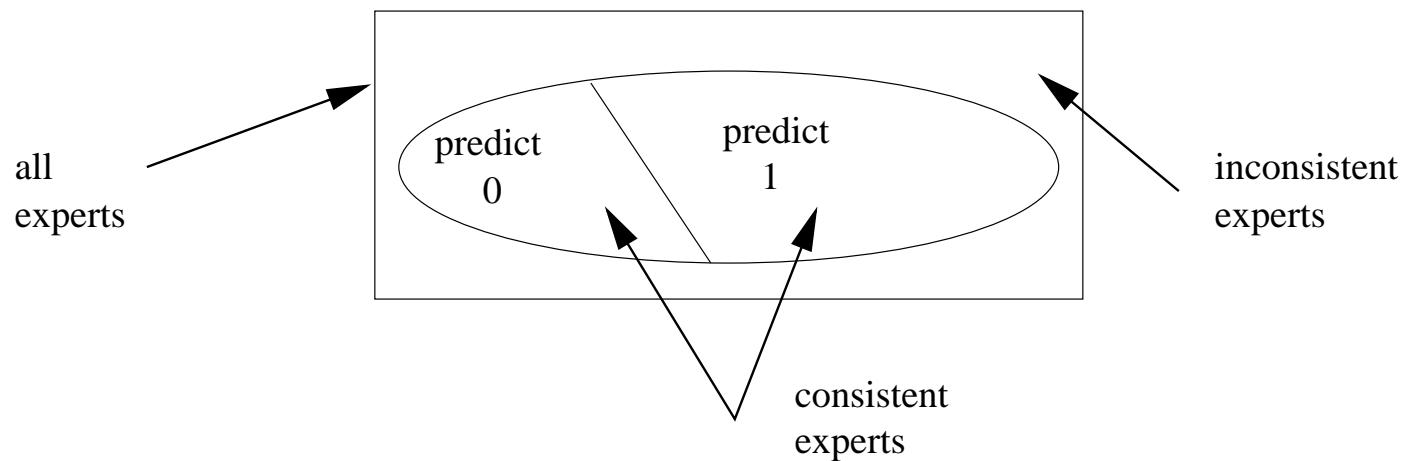
Protocol of the Master Algorithm

For $t = 1$ To T Do

- Get instance $x_t \in \{0, 1\}^n$
- Predict $\hat{y}_t \in \{0, 1\}$
- Get label $y_t \in \{0, 1\}$
- Incur loss $|y_t - \hat{y}_t|$

Halving Algorithm

[BF]



- Predicts with majority
- If mistake is made then number of **consistent** experts is (at least) **halved**

A run of the Halving Algorithm

E_1	E_2	E_3	E_4	E_5	E_6	E_7	E_8	<i>majo rity</i>	<i>true label</i>	loss
1	1	0	0	1	1	0	0	1	0	1
x	x	0	1	x	x	1	1	1	1	0
x	x	x	1	x	x	0	0	0	1	1
x	x	x	↑	x	x	x	x			
consistent										

For any sequence with a **consistent** expert
 HA makes $\leq \log_2 n$ mistakes

What if no expert is consistent?

Sequence of examples $S = (x_1, y_1), \dots, (x_T, y_T)$

- $L_A(S)$ be the total loss of alg. A
- $L_i(S)$ be the total loss of i -th expert E_i

Want bounds of the form:

$$\forall S : L_A(S) \leq a \min_i L_i(S) + b \log(n)$$

where a, b are constants

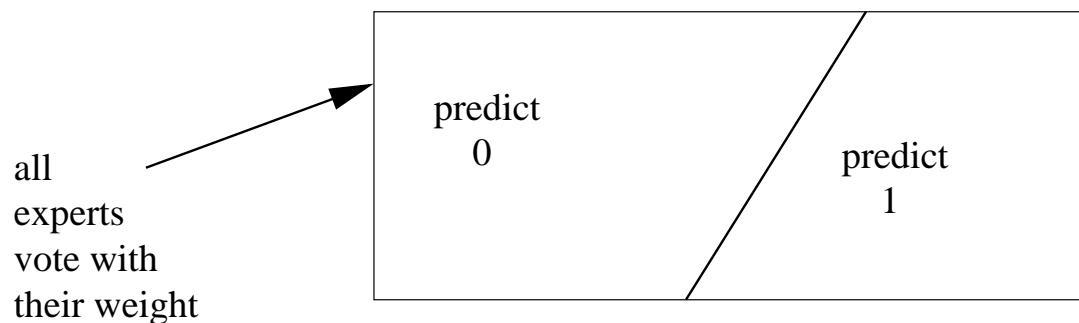
Bounds loss of algorithm **relative to** loss of best expert

Can't wipe out experts!

One weight per expert

Weighted Majority Algorithm

[LW]



- Predicts with larger side
- Weights of wrong experts are multiplied by $\beta \in [0, 1)$

Number of mistakes of the WM algorithm

$$\begin{aligned} M_{t,i} &= \# \text{ of mistakes of } E_i \text{ before trial } t \\ w_{t,i} &= \beta^{M_{t,i}} \text{ weight of } E_i \text{ at beginning of trial } t \\ W_t &= \sum_{i=1}^n w_{t,i} \text{ total weight at trial } t \end{aligned}$$

$$\text{Minority} \leq \frac{1}{2}W_t$$

$$\text{Majority} \geq \frac{1}{2}W_t$$

If no mistake then

minority multiplied by β :

$$W_{t+1} \leq 1 W_t$$

If mistake then

majority multiplied by β :

$$\begin{aligned} W_{t+1} &\leq 1 \frac{\frac{1}{2}W_t}{\text{minority}} + \beta \frac{\frac{1}{2}W_t}{\text{majority}} \\ &= \frac{1+\beta}{2} W_t \end{aligned}$$

Hence

$$\begin{aligned} \underset{\substack{\text{total final} \\ \text{weight}}}{W_{T+1}} &\leq \left(\frac{1+\beta}{2}\right)^M W_1 \\ W_{T+1} &= \sum_{j=1}^n w_{T+1,j} = \sum_{j=1}^n \beta^{M_j} \geq \beta^{M_i} \end{aligned}$$

We got:

$$\left(\frac{1+\beta}{2}\right)^M \underbrace{W_1}_n \geq \beta^{M_i}$$

Solving for M :

$$\begin{aligned} M &\leq \frac{\ln \frac{1}{\beta}}{\ln \frac{2}{1+\beta}} M_i + \frac{1}{\ln \frac{2}{1+\beta}} \ln n \\ M &\leq \underbrace{\frac{2.63}{a} \min_i M_i}_{\beta = 1/e} + \underbrace{\frac{2.63}{b} \ln n}_{\beta = 1/e} \end{aligned}$$

For all sequences, loss of master alg.
is comparable to loss of best expert
Relative loss bounds [Fr]

Other Loss Functions

$$\text{absolute loss } L(y, \hat{y}) = |y - \hat{y}|$$

$$\text{square loss } L(y, \hat{y}) = (y - \hat{y})^2$$

$$\text{entropic loss } L(y, \hat{y}) = y \ln \frac{y}{\hat{y}} + (1 - y) \ln \frac{1-y}{1-\hat{y}}$$
$$y, \hat{y} \in [0, 1]$$

One weight per expert:

[V]

$$w_{t,i} = \beta^{L_{t,i}} = e^{-\eta} L_{t,i}$$

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

Master predicts with the weighted average [KW]

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

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

\forall sequences S of examples $\langle(x_t, y_t)\rangle_{1 \leq t \leq T}$
 where $x_t \in [0, 1]^n$ and $y_t \in [0, 1]$

$$L_{WA}(S) \leq \min_i \underbrace{\frac{1}{a} L_i(S)}_a + \underbrace{\frac{1}{\eta} \ln(n)}_b$$

$1/\eta$	dot pred	fancy
entropic	1	1
square	2	$1/2$
hellinger	1	.71

- Slightly improved constants of $1/\eta$ when Master uses fancier prediction [V]
- For the discrete loss and the absolute loss $a > 1$

Potential: $\frac{1}{\eta} \ln W_t$

$$\begin{aligned}\text{Key inequality: } L(y, \mathbf{v}_t \cdot \mathbf{x}_t) &\leq \frac{1}{\eta} \ln W_t - \frac{1}{\eta} \ln W_{t+1} \\ &= -\frac{1}{\eta} \ln \frac{W_{t+1}}{W_t}\end{aligned}$$

Telescoping:

$$\begin{aligned}L_{\text{WA}}(S) &\leq -\frac{1}{\eta} \ln \frac{W_{T+1}}{W_1} \\ &= -\frac{1}{\eta} \ln \sum_{j=1}^n \frac{1}{n} e^{-\eta L_j(S)} \\ &\leq -\frac{1}{\eta} \ln \frac{1}{n} e^{-\eta L_i(S)} \\ &= -\frac{1}{\eta} \ln \frac{1}{n} e^{-\eta L_i(S)} \\ &= L_t(S) + \frac{1}{\eta} \ln n\end{aligned}$$

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)

Questions:

- How to obtain algs. 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?

A more general setting

Instance	Prediction of alg A	Label	Loss of alg A
x_1	\hat{y}_1	y_1	$L(y_1, \hat{y}_1)$
\vdots	\vdots	\vdots	\vdots
x_t	\hat{y}_t	y_t	$L(y_t, \hat{y}_t)$
\vdots	\vdots	\vdots	\vdots
x_T	\hat{y}_T	y_T	$L(y_T, \hat{y}_T)$
Total Loss			$L_A(S)$

Sequence of examples $S = (x_1, y_1), \dots, (x_T, y_T)$

Comparison class $\{u\}$

Relative loss

$$L_A(S) - \inf_{\{u\}} \text{Loss } u(S)$$

Goal: Bound relative loss
for arbitrary sequence S

Learning Disjunctions of Experts

variables/experts				<i>true label</i>	$E_1 \vee E_3$	$E_3 \vee E_4$
E_1	E_2	E_3	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
$x_{t,1}$	$x_{t,2}$	$x_{t,3}$	$x_{t,4}$		↑	↑
					3	2
					mistakes	

$$E_1 \vee E_3 \quad \text{becomes} \quad \mathbf{u} = (1, 0, 1, 0)$$

$$E_1 \vee E_3 \text{ is one on } \mathbf{x}_t \in \{0, 1\}^n \quad \text{iff} \quad \mathbf{u} \cdot \mathbf{x}_t \geq 1$$

Do as well as best

k out of n literal (monotone) disjunction

Each disjunction is an expert

Keep one weight per disjunction: $\binom{n}{k}$ weights

$$\begin{aligned} \text{\# of mistakes} \\ \text{of WM} \end{aligned} \leq 2.63 M + 2.63 k \ln \frac{n}{k}$$

M is # of mistakes of best

Time (and space) exponential in k

Efficient algs: one weight per **literal**

The Perceptron Algorithm

In trial t : Get instance $x_t \in \{0, 1\}^n$
If $w_t \cdot x_t \geq 1/2$ then $\hat{y}_t = 1$
else $\hat{y}_t = 0$
Get label $y_t \in \{0, 1\}$
If mistake then
 $w_{t+1} = w_t - \eta (\hat{y}_t - y_t) x_t$

Perc. Conv. Th. ($\eta = \frac{1}{2n}$)

$$\# \text{ of mistakes} \leq 4 A + 4 k n$$

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 rot. inv. algs:

[KWA]

$$\#\text{mistakes} = \Omega(n)$$

The Winnow Algorithm

[L]

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

If $w_t \cdot x_t \geq \theta$ 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 ($e^{-\eta} = 1/3$, $\theta = \frac{3 \ln 3}{8}$) [AW]

$$\# \text{ of mistakes} \leq 4 A + 3.6 k \ln \frac{n}{k}$$

Not rotation invariant!

k -term DNF via Feature Expansion [KW]

$$\begin{array}{c}
 & & & 1 \\
 & & & x_1 \\
 & & & x_2 \\
 & & & \vdots \\
 & & & x_n \\
 x = & \begin{matrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{matrix} & \Rightarrow & \Phi(x) = & \begin{matrix} x_1 x_2 \\ \vdots \\ x_1 x_2 \dots x_n \end{matrix} \\
 & n \text{ inputs} & & & 2^n \text{ features}
 \end{array}$$

k -term DNF in input space
is k -literal disjunction in feature space

$$\underbrace{\Phi(x)}_{2^n} \cdot \underbrace{\Phi(y)}_{2^n} = \underbrace{\prod_{i=1}^n (1 + x_i y_i)}_{O(n) \text{ time}} = K(\underbrace{\underline{x}_n}_{\text{ }}, \underbrace{\underline{y}_n}_{\text{ }})$$

(Simple ANOVA kernel)

Perceptron: $w_t = \sum_{q \text{ mistake}} \alpha_q \Phi(x_i)$

Prediction:

$$\begin{aligned}
 w_t \cdot \Phi(x) &= \left(\sum_{q \text{ mistake}} \alpha_q \Phi(x_q) \right) \cdot \Phi(x) \\
 &= \sum_{q \text{ mistake}} \alpha_q \Phi(x_q) \cdot \Phi(x) \\
 &= \underbrace{\sum_{q \text{ mistake}} \alpha_q K(x_q, x)}_{\text{time: } O(n \cdot \# \text{ mistakes})}
 \end{aligned}$$

Mistake bound: $O(k 2^n)$

Winnow: $w_{t,i} = \exp\left(-\eta \sum_{q \text{ mistake}} \alpha_q \Phi(x_q)_i\right)$

log of weights is

linear comb of past examples

Mistake bound: $O(k \ln 2^n) = O(k n)$

prediction time: $\Omega(2^n \# \text{ mistakes})$

No kernel trick with purely mult. updates!

So far

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

Rest of tutorial

- Motivation of updates with Bregman divergences
- Bregman divergences as loss functions
- Pythagorean Theorem
- Proving relative loss bounds
- Conversions to batch model
- Bregman divergences and the exponential family
- Comparator shifts with time

On-line Linear Regression

For $t = 1, \dots, T$ do

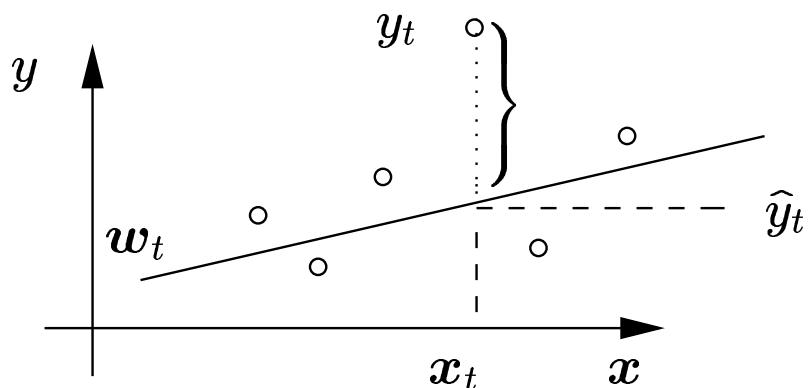
Get instance $x_t \in \mathbf{R}^n$

Predict $\hat{y}_t = w_t \cdot x_t$

Get label $y_t \in \mathbf{R}$

Incur loss $L_t(w_t) = (y_t - \hat{y}_t)^2$

Update w_t to w_{t+1}



Assume comparison class $\{u\}$ is a set of linear predictors

$$u : x \rightarrow u \cdot x$$

Examples of Updates

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

$$\begin{aligned}\mathbf{w}_{t+1} &= \mathbf{w}_t - \eta \nabla L_t(\mathbf{w}_t) \\ &= \mathbf{w}_t - \eta (\mathbf{w}_t \cdot \mathbf{x}_t - y_t) \mathbf{x}_t \quad [\text{WH}]\end{aligned}$$

Exponentiated Gradient Algorithm [KW]
 $(\mathbf{w}$ is probability vector)

$$\mathbf{w}_{t+1,i} = w_{t,i} \exp \left[-\eta \frac{\partial L_t(\mathbf{w}_t)}{\partial w_{t,i}} \right] / \text{normaliz.}$$

More examples of Updates

Unnormalized Exponentiated Gradient Alg. [KW]
 $(\mathbf{w} \geq 0)$

$$\mathbf{w}_{t+1,i} = w_{t,i} \exp \left[-\eta \frac{\partial L_t(\mathbf{w}_t)}{\partial w_{t,i}} \right]$$

Binary Exponentiated Gradient Algorithm [By]
 $(\mathbf{w} \in [0, 1]^n)$

$$\mathbf{w}_{t+1,i} = \frac{w_{t,i} \exp \left[-\eta \frac{\partial l_t(\mathbf{w}_t)}{\partial w_{t,i}} \right]}{1 - w_{t,i} + w_{t,i} \exp \left[-\eta \frac{\partial L_t(\mathbf{w}_t)}{\partial w_{t,i}} \right]}$$

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

[GLS, GL]

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

where

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

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

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

Motivation of Updates

[KW]

Gradient descent

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

Exponentiated Gradient Algorithm

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

Families of update algorithms

parameter divergence	name of family	update algs.
$\ \mathbf{w} - \mathbf{w}_t\ _2^2$	Grad. Desc.	Widrow Hoff (LMS) Lin. Least Squ. Backprop. Perceptron Alg. kernel based algs.,...
$\sum_{i=1}^n w_i \ln \frac{w_i}{w_{t,i}}$	Exponentiated Gradient Alg.	expert algs Normalized Winnow “AdaBoost”

Families of update algorithms (cont)

$\sum_{i=1}^n w_i \ln \frac{w_i}{w_{t,i}}$ + $w_{t,i} - w_i$	Unnormalized Exp. Grad. Alg.	Winnow
$\sum_{i=1}^n w_i \ln \frac{w_i}{w_{t,i}}$ + $(1 - w_i) \ln \frac{1-w_i}{1-w_{t,i}}$	Binary Exp. Grad. Alg.	
any Bregman divergence	-	-

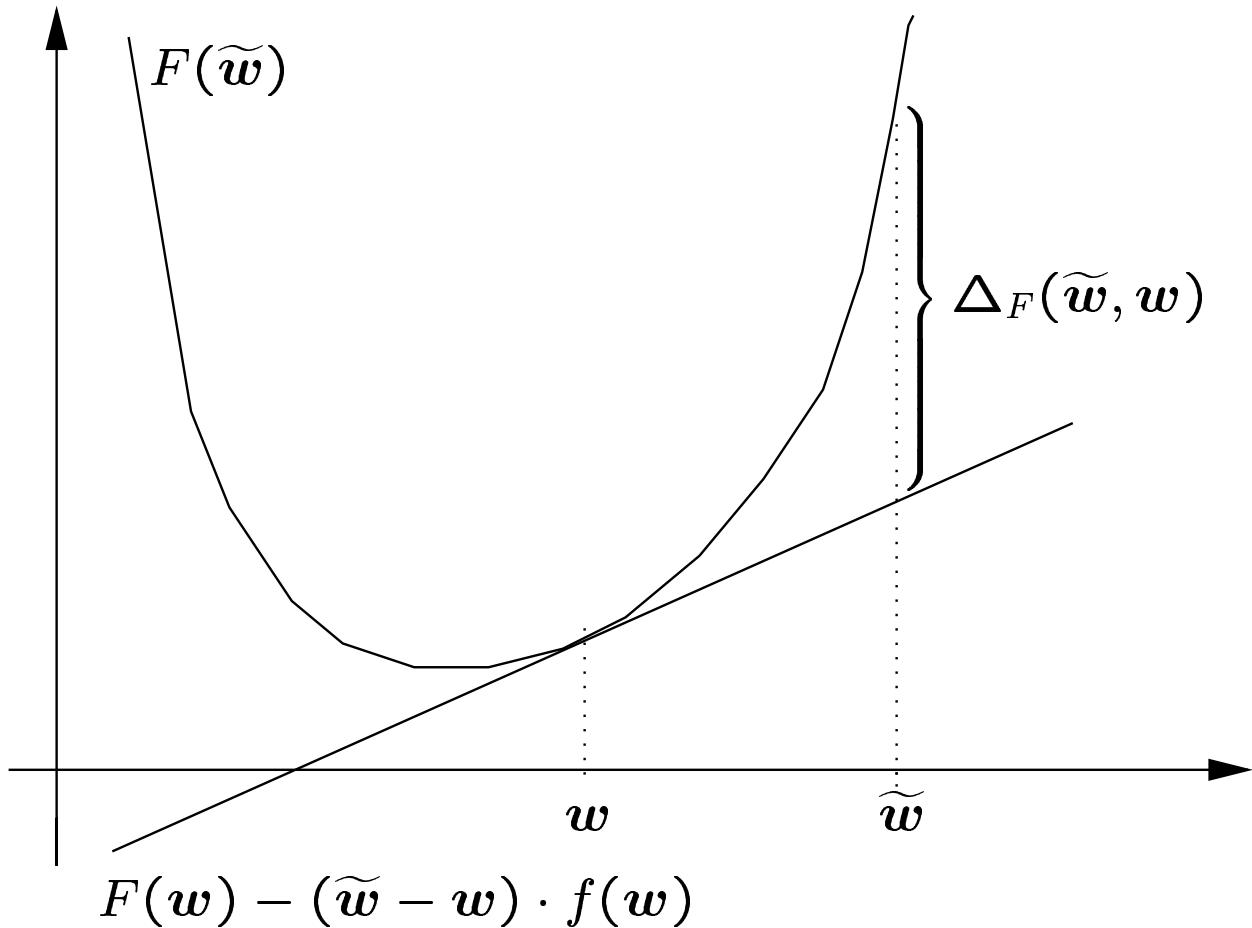
Members of different families exhibit different behavior

Bregman Divergences

[Br, CL, Cs]

For **any** differentiable convex function F

$$\begin{aligned}\Delta_F(\tilde{\mathbf{w}}, \mathbf{w}) &= F(\tilde{\mathbf{w}}) - F(\mathbf{w}) - (\tilde{\mathbf{w}} - \mathbf{w}) \cdot \underbrace{\nabla_{\mathbf{w}} F(\mathbf{w})}_{f(\mathbf{w})} \\ &= F(\tilde{\mathbf{w}}) - \text{supporting hyperplane} \\ &\quad \text{through } (\mathbf{w}, F(\mathbf{w}))\end{aligned}$$



Bregman Divergences, Simple Properties [AW]

1. $\Delta_F(\tilde{w}, w)$ is convex in \tilde{w}
2. $\Delta_F(\tilde{w}, w) \geq 0$
If F convex equality holds iff $\tilde{w} = w$
3. $\nabla_{\tilde{w}} \Delta_F(\tilde{w}, w) = f(\tilde{w}) - f(w)$
4. Usually not symmetric: $\Delta_F(\tilde{w}, w) \neq \Delta_F(w, \tilde{w})$
5. Linearity (for $a \geq 0$):
$$\Delta_{F+aH}(\tilde{w}, w) = \Delta_F(\tilde{w}, w) + a \Delta_H(\tilde{w}, w)$$
6. Unaffected by linear terms ($a \in \mathbf{R}$, $b \in \mathbf{R}^n$):
$$\Delta_{H+a\tilde{w}+b}(\tilde{w}, w) = \Delta_H(\tilde{w}, w)$$
7.
$$\begin{aligned} \Delta_F(w_1, w_2) + \Delta_F(w_2, w_3) \\ = \Delta_F(w_1, w_3) + (w_1 - w_2) \cdot (f(w_3) - f(w_2)) \end{aligned}$$

Examples

Squared Euclidean Distance

$$F(\mathbf{w}) = \|\mathbf{w}\|_2^2 / 2$$

$$f(\mathbf{w}) = \mathbf{w}$$

$$\begin{aligned} \Delta_F(\tilde{\mathbf{w}}, \mathbf{w}) &= \|\tilde{\mathbf{w}}\|_2^2 / 2 - \|\mathbf{w}\|_2^2 / 2 - (\tilde{\mathbf{w}} - \mathbf{w}) \cdot \mathbf{w} \\ &= \|\tilde{\mathbf{w}} - \mathbf{w}\|_2^2 / 2 \end{aligned}$$

(Unnormalized) Relative Entropy

$$F(\mathbf{w}) = \sum_i (w_i \ln w_i - w_i)$$

$$f(\mathbf{w}) = \ln \mathbf{w}$$

$$\Delta_F(\tilde{\mathbf{w}}, \mathbf{w}) = \sum_i \left(\tilde{w}_i \ln \frac{\tilde{w}_i}{w_i} + w_i - \tilde{w}_i \right)$$

Examples (cont)

[GLS, GL]

p-norm Algs (*q* is dual to *p*)

$$F(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_q^2$$

$$f(\mathbf{w}) = \nabla \frac{1}{2} \|\mathbf{w}\|_q^2$$

$$\Delta_F(\tilde{\mathbf{w}}, \mathbf{w}) = \frac{1}{2} \|\tilde{\mathbf{w}}\|_q^2 + \frac{1}{2} \|\mathbf{w}\|_q^2 - \tilde{\mathbf{w}} \cdot f(\mathbf{w})$$

When $p = q = 2$ this reduces to squared Euclidean distance (Widrow-Hoff).

General Motivation of Updates

[KW]

Trade-off between two divergences:

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w}} \left(\underbrace{\Delta_F(\mathbf{w}, \mathbf{w}_t)}_{\text{weight domain}} + \eta_t \underbrace{L_t(\mathbf{w})}_{\text{label domain}} \right)$$

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

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

$$\nabla_{\mathbf{w}} (\Delta_F(\mathbf{w}, \mathbf{w}_t) + \eta_t L_t(\mathbf{w})) = 0$$

iff

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

\Rightarrow

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

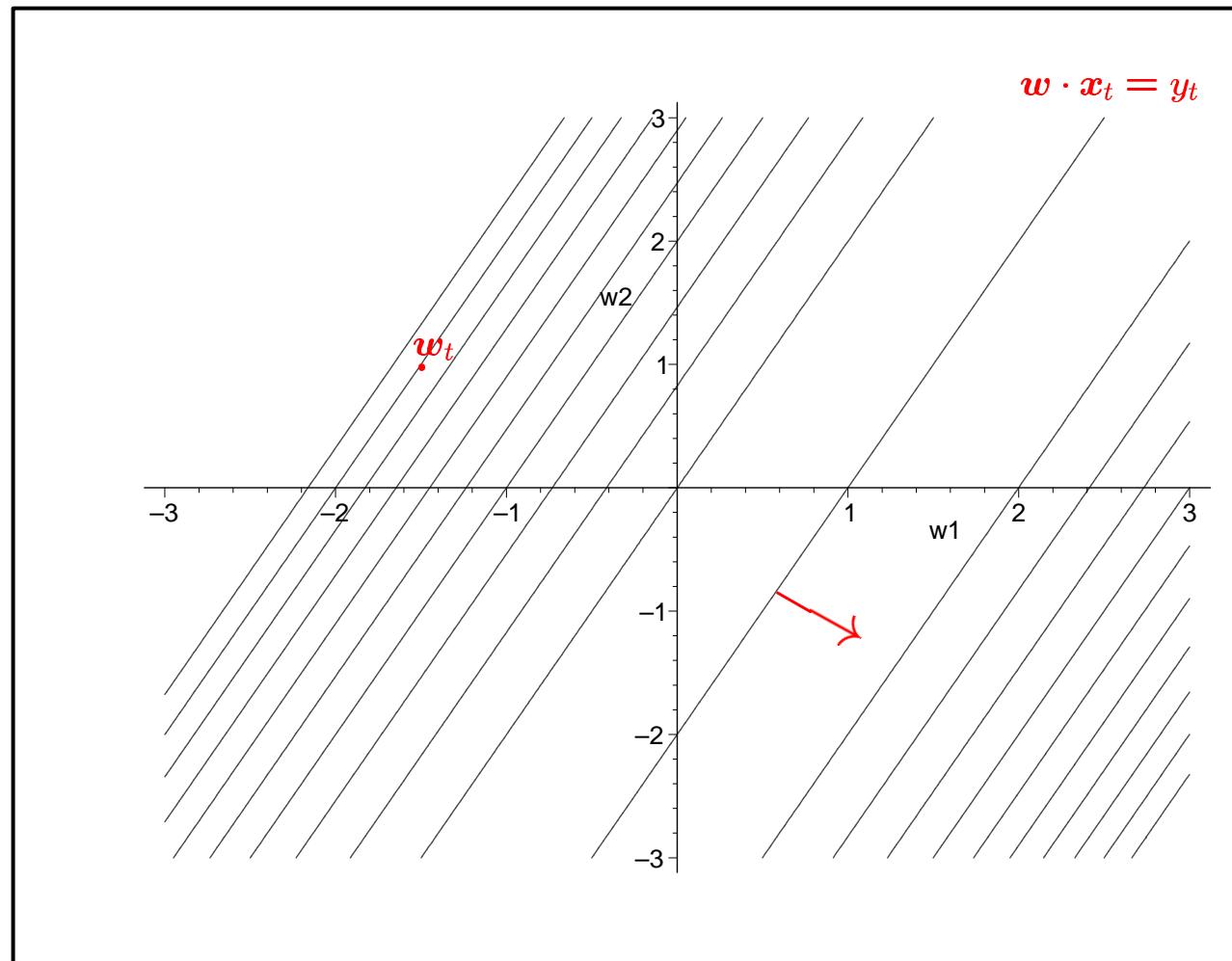
Square Loss

$$L_t(\mathbf{w}) = (y_t - \mathbf{w} \cdot \mathbf{x}_t)^2 / 2$$

$$\mathbf{w}_t = (-3/2, 1)$$

$$\mathbf{x}_t = (1, -0.5)$$

$$y_t = 1$$



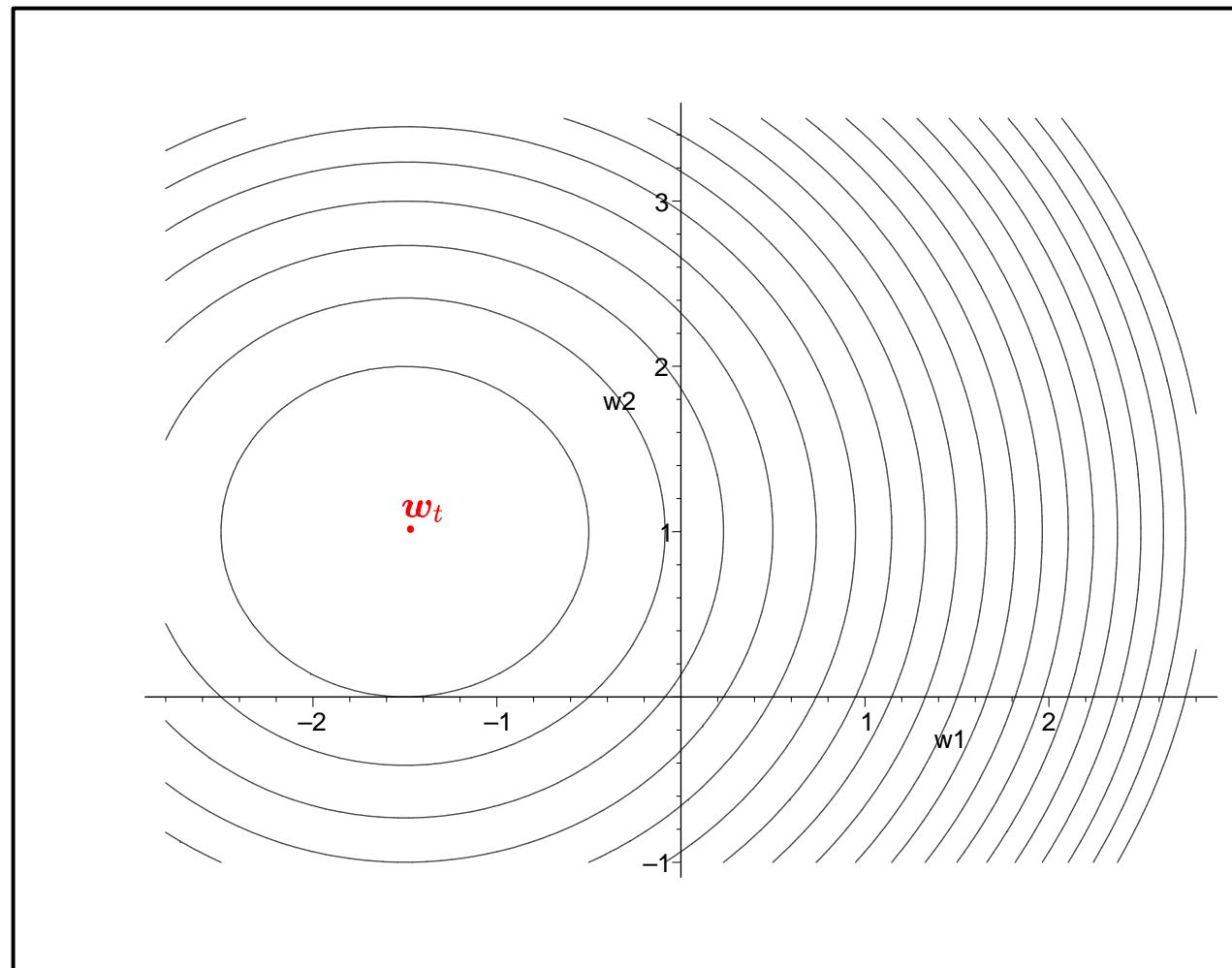
Divergence: Euclidean Distance Squared

$$\Delta_F(\mathbf{w}, \mathbf{w}_t) = \|\mathbf{w} - \mathbf{w}_t\|_2^2 / 2$$

$$\mathbf{w}_t = (-3/2, 1)$$

$$\mathbf{x}_t = (1, -0.5)$$

$$y_t = 1$$



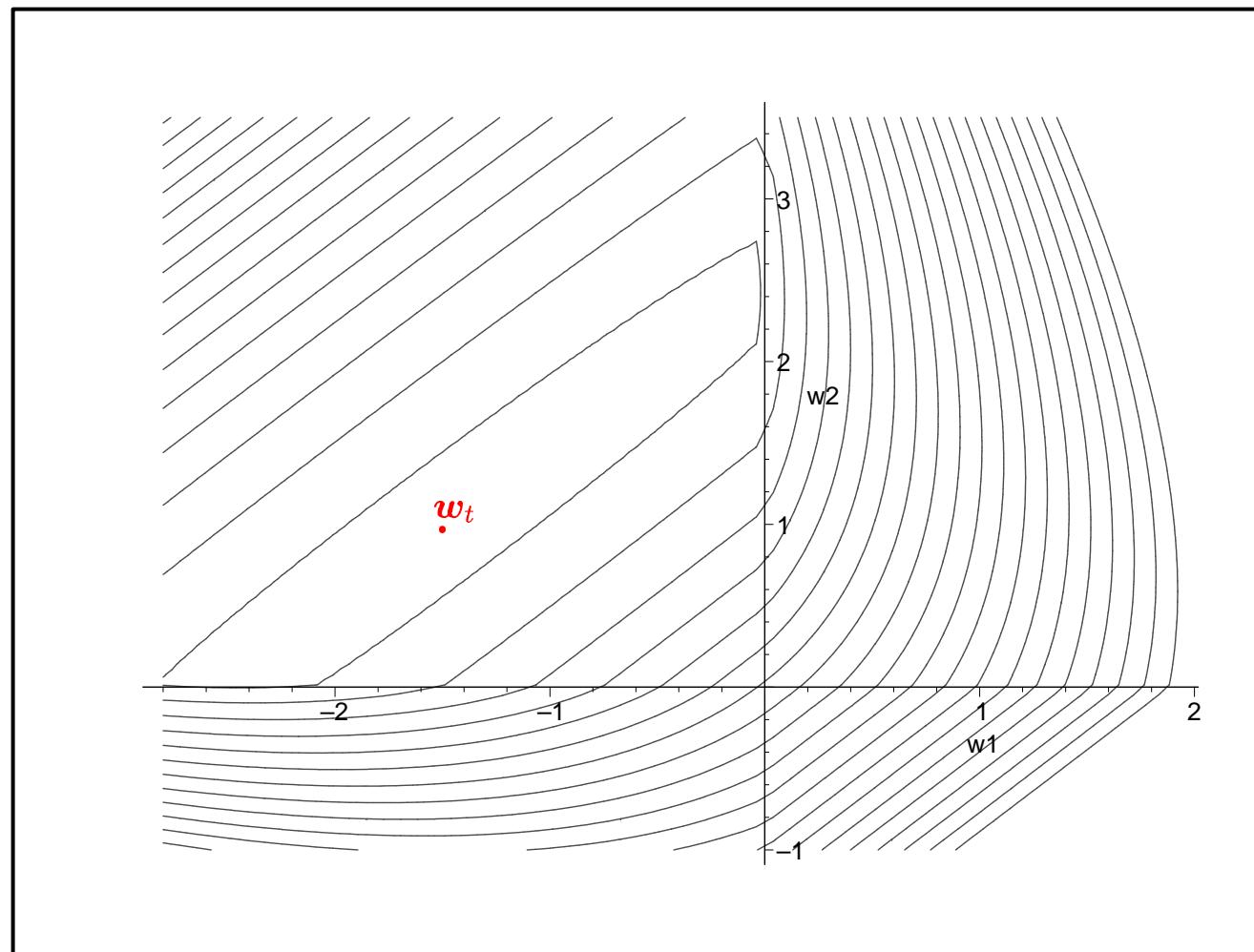
Divergence: ℓ_0 -norm alg. divergence

$\Delta_F(\mathbf{w}, \mathbf{w}_t)$ is
 ℓ_0 -norm alg.
divergence

$$\mathbf{w}_t = (-3/2, 1)$$

$$\mathbf{x}_t = (1, -0.5)$$

$$y_t = 1$$



Loss + η Divergence

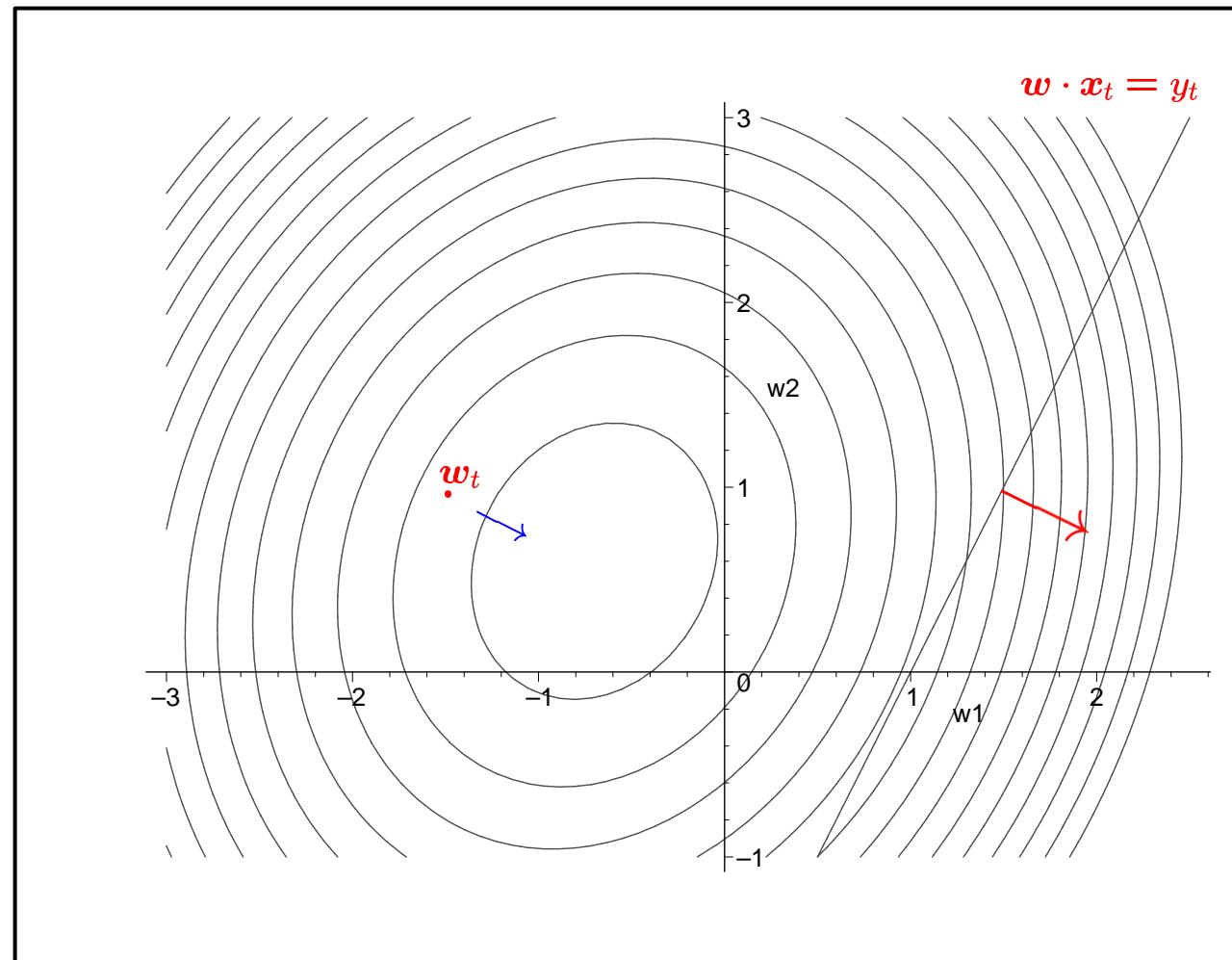
$$\Delta_F(\mathbf{w}, \mathbf{w}_t) = \|\mathbf{w} - \mathbf{w}_t\|_2^2 / 2$$

$$\mathbf{w}_t = (-3/2, 1)$$

$$\mathbf{x}_t = (1, -0.5)$$

$$y_t = 1$$

$$\eta = 0.2$$



Loss + η Divergence

$\Delta_F(\mathbf{w}, \mathbf{w}_t)$
is 10 -norm alg.
divergence

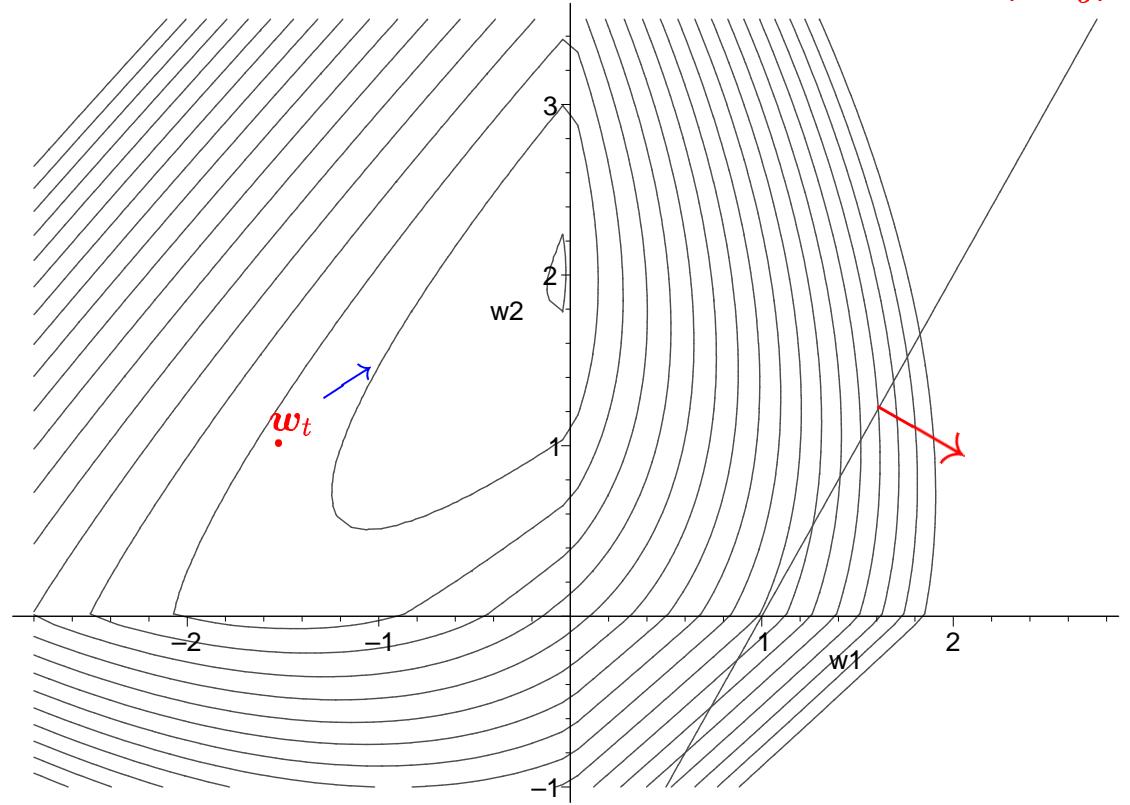
$$\mathbf{w}_t = (-3/2, 1)$$

$$\mathbf{x}_t = (1, -0.5)$$

$$y_t = 1$$

$$\eta = 0.2$$

$$\mathbf{w} \cdot \mathbf{x}_t = y_t$$



Characterization of algs.

i.t.o. link function $f = \nabla F$

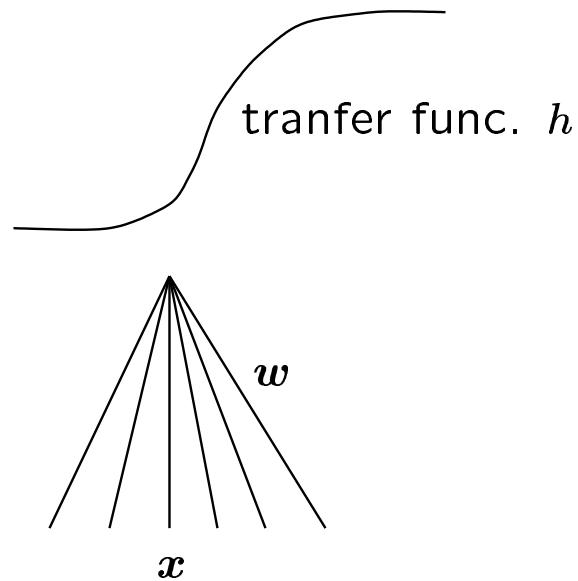
[WJ]

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

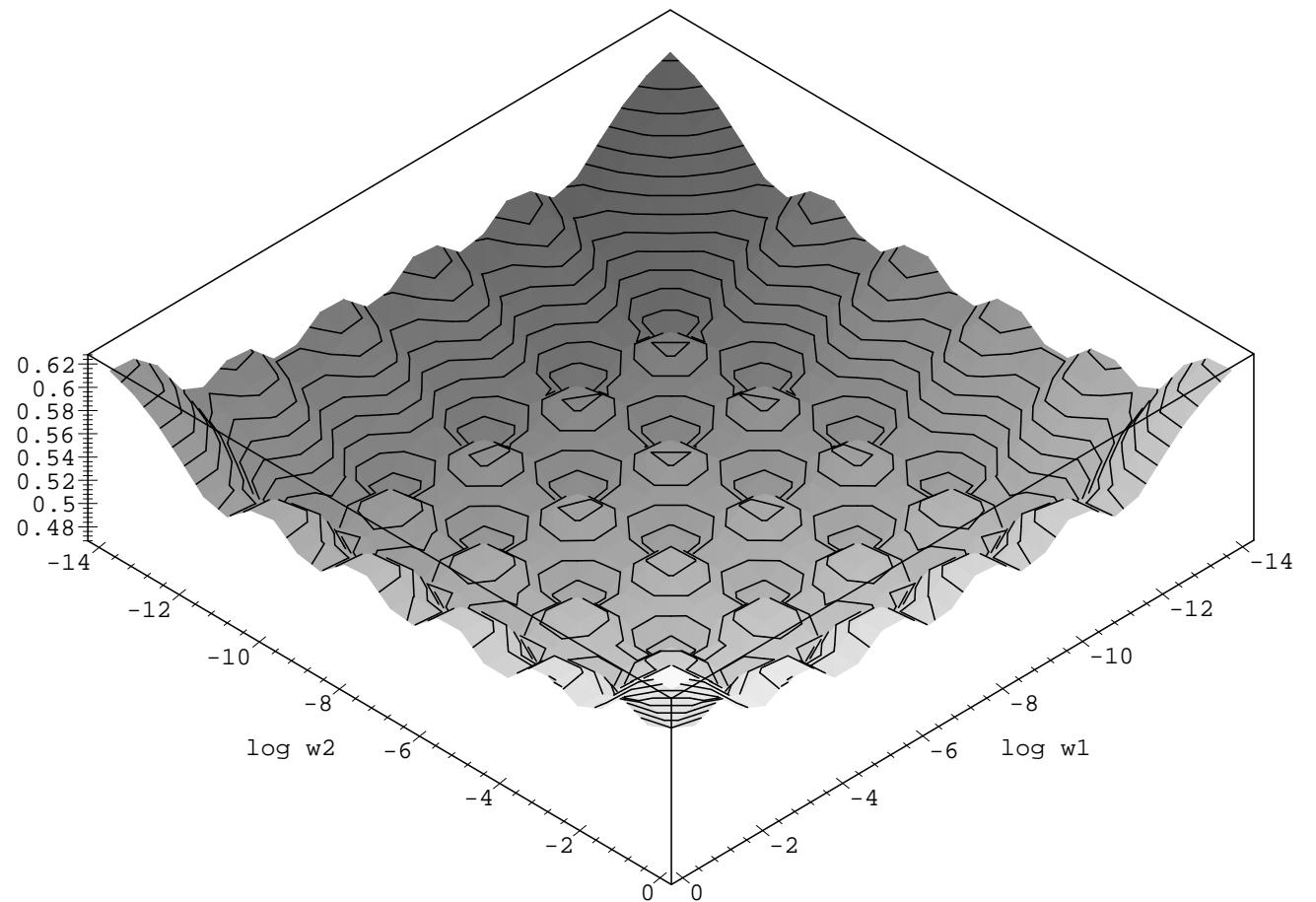
Alg.	$f(\mathbf{w})$	Domain of \mathbf{w}
GD	$f(\mathbf{w}) = \mathbf{w}$	$\mathbf{w} \in \mathbf{R}^n$
BEG	$f(\mathbf{w}) = \ln \frac{\mathbf{w}}{1-\mathbf{w}}$	$\mathbf{w} \in [0, 1]^n$
EG	$f(\mathbf{w}) = \ln \frac{\mathbf{w}}{1- \mathbf{w} _1}$	$\mathbf{w} \in [0, 1]^{n-1}, \mathbf{w} _1 \leq 1$
p -norm	$f(\mathbf{w}) = \nabla \frac{1}{2} \mathbf{w} _q^2$ where $\frac{1}{p} + \frac{1}{q} = 1$	$\mathbf{w} \in \mathbf{R}^n$

Bregman Divergences Lead to Good Loss Functions

$$\hat{y} = h(\mathbf{w} \cdot \mathbf{x})$$

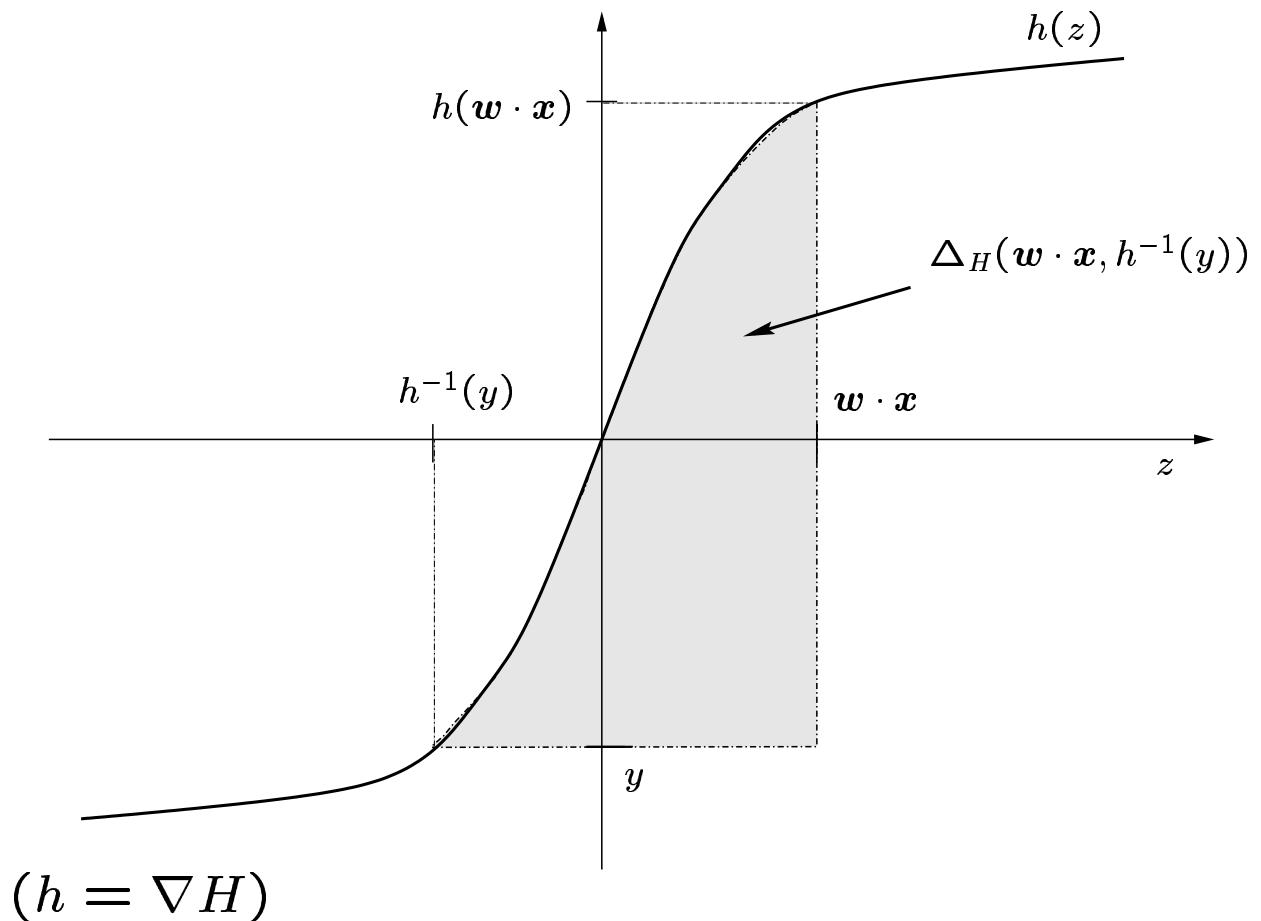


- Sigmoid function $h(z) = \frac{1}{1+e^{-z}}$
- For a set of examples $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$ total loss $\sum_{t=1}^T (h(\mathbf{w} \cdot \mathbf{x}) - y_t)^2 / 2$ can have exponentially many minima in weight space [Bu,AHW]



Want loss that is convex in w

Bregman Divergences Lead to Good Loss Functions (cont)



$$\begin{aligned}
 & \int_{h^{-1}(y)}^{w \cdot x} (h(z) - y) dz \\
 &= H(w \cdot x) - H(h^{-1}(y)) - (w \cdot x - h^{-1}(y)) y \\
 &= \Delta_H(w \cdot x, h^{-1}(y))
 \end{aligned}$$

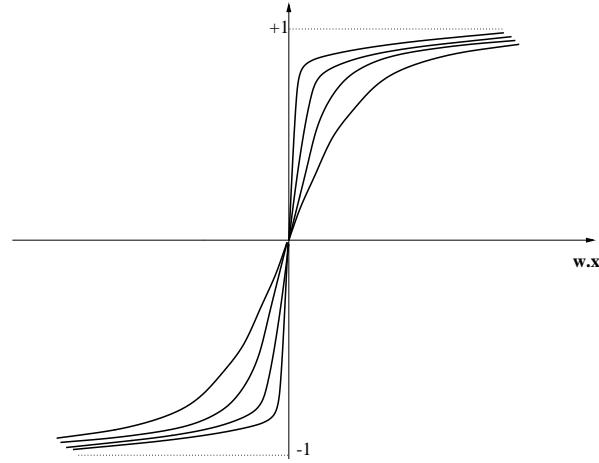
Use $\Delta_H(\mathbf{w} \cdot \mathbf{x}, h^{-1}(y))$ as loss of \mathbf{w} on (\mathbf{x}, y)

Called **matching loss** for h [AHW,HKW]

Matching loss is **convex** in \mathbf{w}

transfer f. $h(z)$	$H(z)$	match. loss $d_H(\mathbf{w} \cdot \mathbf{x}, h^{-1}(y))$
z	$\frac{1}{2}z^2$	$\frac{1}{2}(\mathbf{w} \cdot \mathbf{x} - y)^2$ square loss
$\frac{e^z}{1+e^z}$	$\ln(1 + e^z)$	$\ln(1 + e^{\mathbf{w} \cdot \mathbf{x}}) - y\mathbf{w} \cdot \mathbf{x}$ $+ y \ln y + (1 - y) \ln(1 - y)$ logistic loss
$\text{sign}(z)$	$ z $	$\max\{0, -y\mathbf{w} \cdot \mathbf{x}\}$ hinge loss

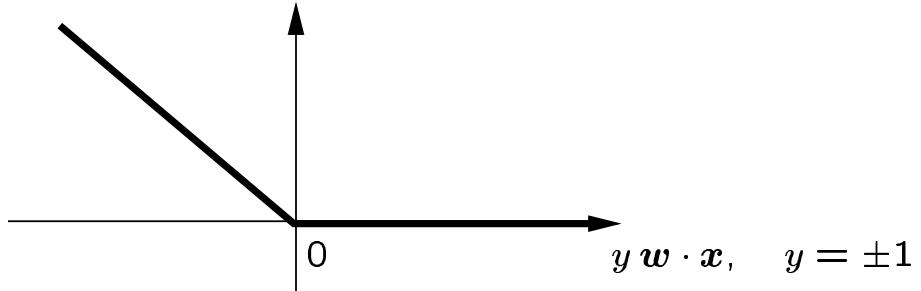
For transfer function $h(z) = \text{sign}(z)$



$$H(z) = |z|$$

Matching loss is **hinge loss** [GW]

$$HL(\mathbf{w} \cdot \mathbf{x}, h^{-1}(y)) = \max\{0, -y \mathbf{w} \cdot \mathbf{x}\}$$



Convex in \mathbf{w} but not differentiable

Motivation of linear threshold algs

Gradient descent
with
Hinge Loss

Perceptron

Expon. gradient
with
Hinge Loss

Normalized
Winnow

Known linear threshold algs for ± 1 -class are
gradient-based algs with hinge loss

Trade-off between two divergences [KW]

$$\begin{aligned} \mathbf{w}_{t+1} &= \operatorname{argmin}_{\mathbf{w}} \left(\underbrace{\Delta_F(\mathbf{w}, \mathbf{w}_t)}_{\text{parameter divergence}} + \underbrace{\eta_t \Delta_H(\mathbf{w} \cdot \mathbf{x}_t, h^{-1}(y_t))}_{\text{matching loss}} \right) \end{aligned}$$

Both divergences are convex in \mathbf{w}

$$\mathbf{w}_{t+1} = f^{-1} (f(\mathbf{w}_t) - \eta_t (h(\mathbf{w}_t \cdot \mathbf{x}_t) - y_t) \mathbf{x}_t)$$

Generalization of the “delta”-rule

Projections

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

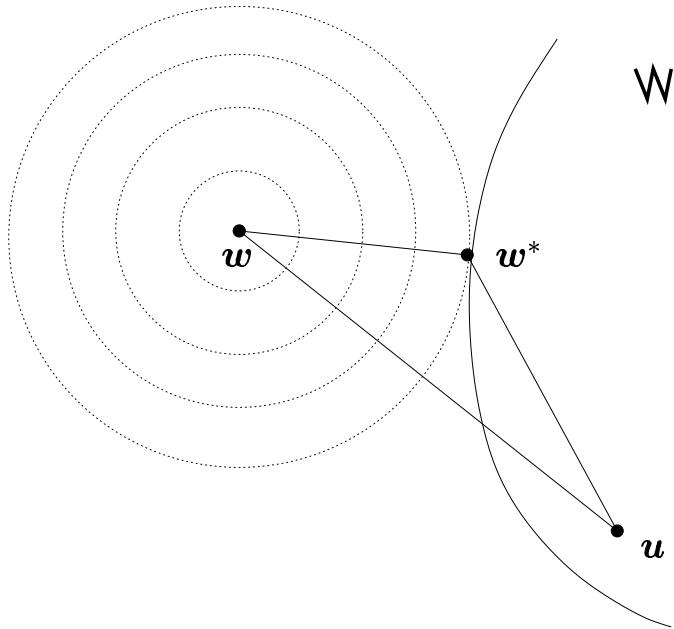
When η is large then \mathbf{w}_{t+1} is projection of \mathbf{w}_t onto plane $\mathbf{w} \cdot \mathbf{x}_t = y_t$

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\{\mathbf{w} : \mathbf{w}_t \cdot \mathbf{x}_t = y_t\}} \Delta_F(\mathbf{w}, \mathbf{w}_t)$$

The AdaBoost update of the probability vector \mathbf{w}_t on the examples is a projection w.r.t. divergence $\Delta_F(\mathbf{w}, \mathbf{w}_t) = \sum_i w_i \ln \frac{w_i}{w_{t,i}}$ [CKW, La, KW, CSS]

A Pythagorean Theorem

[Br,Cs,A,HW]



w^* is projection of w onto convex set \mathcal{W} w.r.t.
Bregman divergence Δ_F :

$$w^* = \operatorname{argmin}_{u \in \mathcal{W}} \Delta_F(u, w)$$

Th:

$$\Delta_F(u, w) \geq \Delta_F(u, w^*) + \Delta_F(w^*, w)$$

How do we prove relative loss bounds?

Loss: $L_t(\mathbf{w}) = L((x_t, y_t), \mathbf{w})$ convex in \mathbf{w}

Divergence: $\Delta_F(\mathbf{u}, \mathbf{w})$
 $= F(\mathbf{u}) - F(\mathbf{w}) - (\mathbf{u} - \mathbf{w}) \cdot f(\mathbf{w})$

Update: $f(\mathbf{w}_{t+1}) - f(\mathbf{w}_t) = -\eta \nabla_{\mathbf{w}} L_t(\mathbf{w}_t)$

$$L_t(\mathbf{u})$$

convexity

$$\overbrace{\geq}^{\text{convexity}} L_t(\mathbf{w}_t) + (\mathbf{u} - \mathbf{w}_t) \cdot \underbrace{\nabla_{\mathbf{w}} L_t(\mathbf{w}_t)}_{\text{update}}$$

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

$$= L_t(\mathbf{w}_t)$$

$$+ \frac{1}{\eta} (\Delta_F(\mathbf{u}, \mathbf{w}_{t+1}) - \Delta_F(\mathbf{u}, \mathbf{w}_t) - \Delta_F(\mathbf{w}_t, \mathbf{w}_{t+1}))$$

Summing over t

[WJ, KW]

$$\begin{aligned}
 \sum_t L_t(\mathbf{w}_t) &\leq \sum_t L_t(\mathbf{u}) \\
 &+ \frac{1}{\eta} \sum_t \left(\Delta_F(\mathbf{u}, \mathbf{w}_t) - \Delta_F(\mathbf{u}, \mathbf{w}_{t+1}) \right. \\
 &\quad \left. + \Delta_F(\mathbf{w}_t, \mathbf{w}_{t+1}) \right) \\
 &\leq \sum_t L_t(\mathbf{u}) \\
 &+ \frac{1}{\eta} \left(\Delta_F(\mathbf{u}, \mathbf{w}_1) - \underbrace{\Delta_F(\mathbf{u}, \mathbf{w}_{T+1})}_{\geq 0} \right) \\
 &+ \frac{1}{\eta} \sum_t \Delta_F(\mathbf{w}_t, \mathbf{w}_{t+1})
 \end{aligned}$$

$$\begin{aligned}
 \sum_t L_t(\mathbf{w}_t) &\leq \sum_t L_t(\mathbf{u}) + \frac{1}{\eta} \Delta_F(\mathbf{u}, \mathbf{w}_1) \\
 &+ \frac{1}{\eta} \sum_t \Delta_F(\mathbf{w}_t, \mathbf{w}_{t+1})
 \end{aligned}$$

Any convex loss and any Bregman divergence!

Key step:

Relate $\Delta_F(\mathbf{w}_t, \mathbf{w}_{t+1})$ to loss $L_t(\mathbf{w}_t)$

Loss & divergence dependent

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

Then solve for $\sum_t L_t(\mathbf{w}_t)$

Yield bounds of the form

$$\sum_t L_t(\mathbf{w}_t) \leq a \sum_t L_t(\mathbf{u}) + b \Delta_F(\mathbf{u}, \mathbf{w}_1)$$

a, b constants, $a > 1$.

Regret bounds ($a = 1$):

time changing η , subtler analysis

[AG]

Some Bounds [KW, GLS, GL] (Linear Regression with Square Loss)

Gradient Descent

$$\sum_t L_t(\mathbf{w}_t) \leq (1 + c) \sum_t L_t(\mathbf{u}) + \frac{1 + c}{c} X_2^2 U_2^2$$

$$||\mathbf{x}_t||_2 \leq X_2, \quad ||\mathbf{u}||_2 \leq U_2, \quad c > 0$$

Scaled Exponentiated Gradient

$$\sum_t L_t(\mathbf{w}_t) \leq (1 + c) \sum_t L_t(\mathbf{u}) + \frac{1 + c}{c} \ln n X_\infty^2 U_1^2$$

$$||\mathbf{x}_t||_\infty \leq X_\infty, \quad ||\mathbf{u}||_1 \leq U_1, \quad c > 0$$

p -norm alg

$$\sum_t L_t(\mathbf{w}_t) \leq (1 + c) \sum_t L_t(\mathbf{u}) + \frac{1 + c}{c} (p - 1) X_p^2 U_q^2$$

$$||\mathbf{x}_t||_p \leq X_p, \quad ||\mathbf{u}||_q \leq U_q, \quad c > 0$$

Generalization Bounds

Loss bounds for on-line algs

⇒

Bounds on expected loss in i.i.d. case

- On-line to batch [L]
- Leave-one out [HW,CB+,KW]

On-line to batch

Simplest case: 0-1 Loss, binary y_t

$$S = (x_1, y_1), \dots, (x_T, y_T) \sim D^T$$

Alg A updates only if mistake

If A has mistake bound M then \exists alg A' :

$$\Pr_{S \sim D^T} \left(\text{err}_D(A'(S)) \leq \epsilon \right) \geq 1 - \delta$$

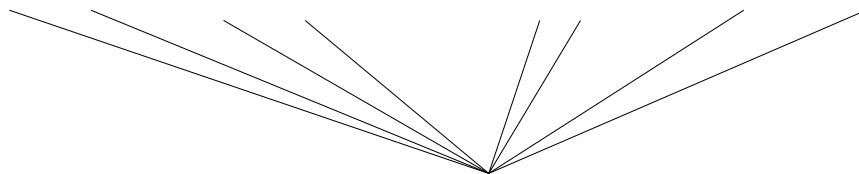
$$\text{err}_D(A'(S)) = \Pr_{(\mathbf{x}, \mathbf{y}) \sim D} \left(A'(S)(\mathbf{x}) \neq \mathbf{y} \right)$$

$$T = O \left(\frac{M}{\epsilon} + \frac{1}{\epsilon} \ln \frac{1}{\delta} \right)$$

Run A

X	X	X	X				X X		X	X	
---	---	---	---	--	--	--	-----	--	---	---	--

 $\#X \leq M+1$



Test hypotheses and pick best

(One must be "good")

(Simplest) Leave One Out

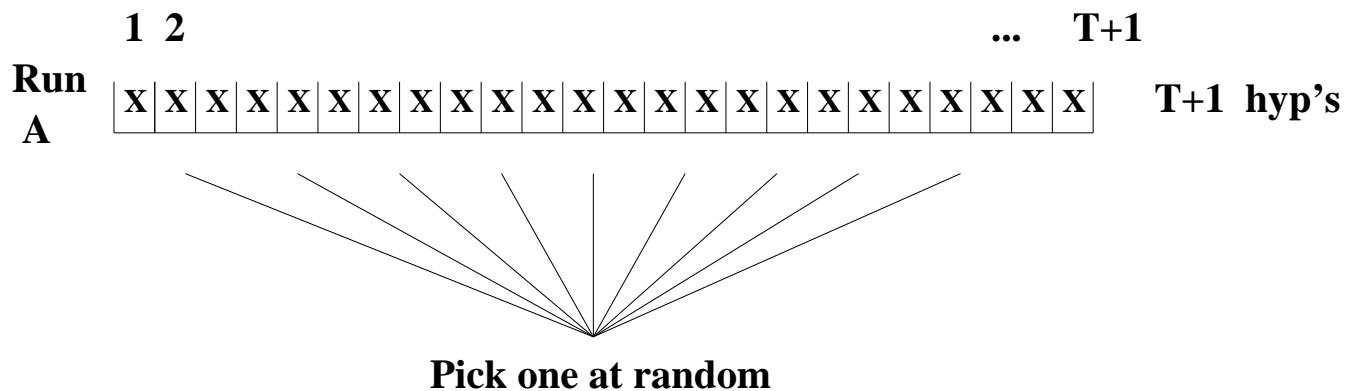
[HW]

Loss $L : \mathbf{R}^2 \rightarrow \mathbf{R}$

$$S = (x_1, y_1), \dots, (x_T, y_T) \sim D^T$$

Given alg A , want to bound

$$E_{(\mathbf{x}, \mathbf{y}) \sim D}[L(\mathbf{y}, A(\mathbf{x}))]$$



call it h_i and predict
on new instance x as $h_i(x)$

(Simplest) Leave One Out (cont)

$$Loo_A(S)(x) = h_i(x) \sim \frac{1}{T+1}$$

Then:

$$\begin{aligned} & E [L(\textcolor{red}{y}, Loo_A(\textcolor{blue}{S})(\textcolor{red}{x}))] \\ & \leq \frac{E_{S \sim D^{T+1}}[\text{cum. loss of } A]}{T+1} \\ & \leq \frac{\text{worst-case loss bound}}{T+1} \end{aligned}$$

Applied to the Perceptron Alg.

[FS]

Where do Bregman divergences come from?

- Exponential family of distributions
- Inherent duality

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

primal param. dual param.

$$\begin{array}{ccc} \mathbf{w}_t & \xrightarrow{f} & f(\mathbf{w}_t) \\ \mathbf{w}_{t+1} & \xleftarrow{f^{-1}} & -\eta \nabla L_t(\mathbf{w}_t) \end{array}$$

Exponential Family of Distributions

- Parametric density functions

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

- θ and x vectors in R^d

- Cumulant function $G(\theta)$
assures normalization

$$G(\theta) = \ln \int e^{\theta \cdot x} P_0(x) dx$$

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

- Expectation parameter

$$\mu = \int_x x P_G(x|\theta) dx = E_{\theta}(x) = g(\theta)$$

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

- Second convex function $F(\mu)$ on space $g(\Theta)$

$$F(\mu) = \theta \cdot \mu - G(\theta)$$

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

Summary

natural
par.

expectation
par.

$$\begin{array}{ccc} \theta & \xrightarrow{g} & \mu \\ G(\theta) & \xleftarrow{f} & F(\mu) \end{array}$$

- θ and μ are dual parameters
- Parameter transformations
 $g(\theta) = \mu$ and $f(\mu) = \theta$

[A, BN]

Gaussian (unit variance)

$$\begin{aligned} P(x|\theta) &\sim e^{-\frac{1}{2}(\theta-x)^2} \\ &= e^{\theta \cdot x - \frac{1}{2}\theta^2} e^{\frac{1}{2}x^2} \end{aligned}$$

Cumulant function: $G(\theta) = \frac{1}{2}\theta^2$

Parameter transformations:

$$g(\theta) = \theta = \mu \quad \text{and} \quad f(\mu) = \mu = \theta$$

Dual convex function: $F(\mu) = \theta \cdot \mu - G(\theta)$
 $= \frac{1}{2}\mu^2$

Square loss: $L_t(\theta) = \frac{1}{2}(\theta_t - x_t)^2$

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:

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

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

$$\begin{aligned} \text{Log loss: } L_t(\theta) &= -x_t \theta + \ln(1 + e^\theta) \\ &= -x_t \ln \mu - (1 - x_t) \ln(1 - \mu) \end{aligned}$$

Poisson

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

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

μ is expectation of x

Natural parameter: $\theta = \ln \mu$

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

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

Parameter transformations:

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

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

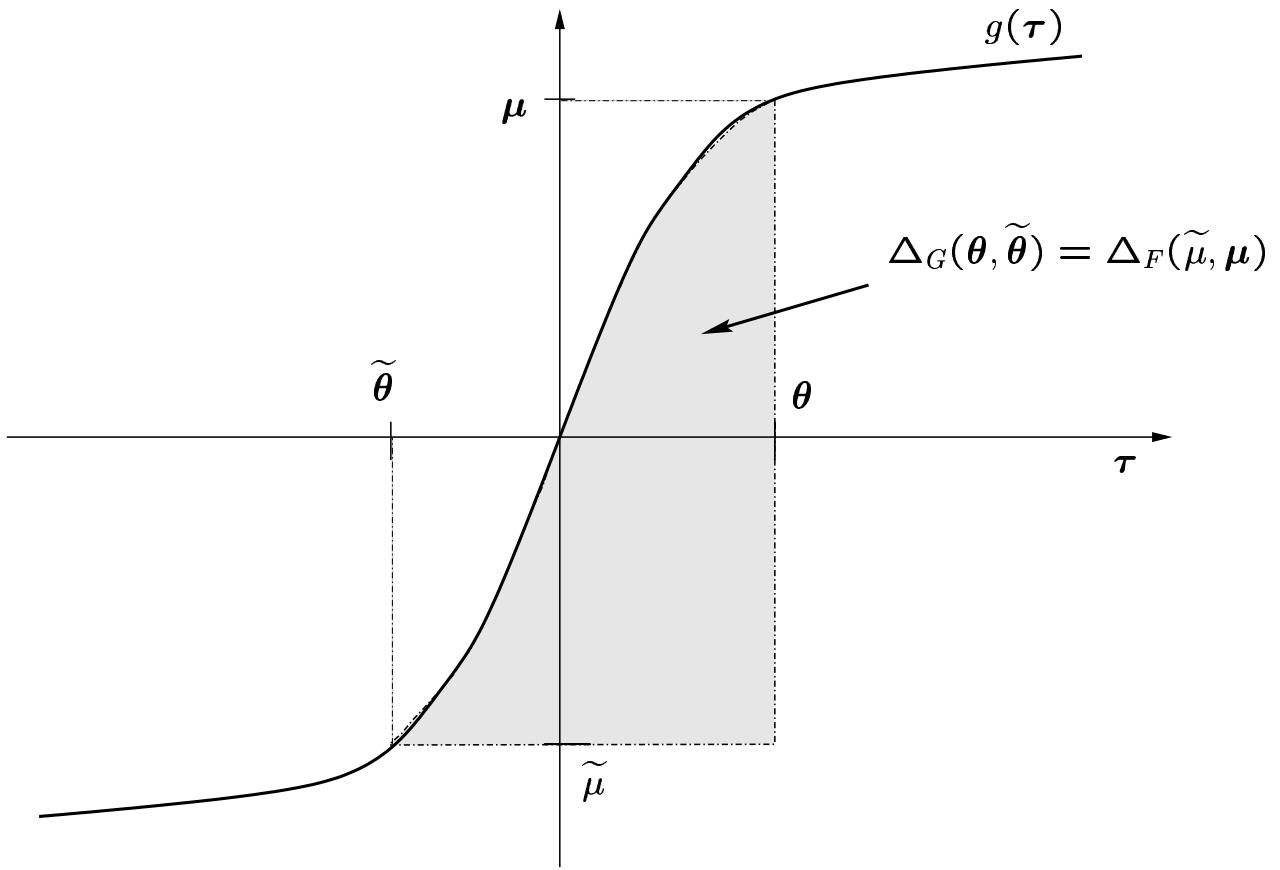
$$\begin{aligned} \text{Loss: } L_t(\theta) &= -x_t\theta + e^\theta + \ln x_t! \\ &= -x_t \ln \mu + \mu + \ln x_t! \end{aligned}$$

Bregman Divergences as Relative Entropies between Exponential Distributions

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

$$\begin{aligned}
 \Delta_G(\tilde{\theta}, \theta) &= \int_x P_G(x|\theta) \ln \frac{P_G(x|\theta)}{P_G(x|\tilde{\theta})} dx \\
 &= G(\tilde{\theta}) - G(\theta) - (\tilde{\theta} - \theta) \cdot \mu \\
 F(\mu) &\stackrel{\theta \cdot \mu - G(\theta)}{=} F(\mu) - F(\tilde{\mu}) - (\mu - \tilde{\mu}) \cdot \tilde{\theta} \\
 &= \Delta_F(\mu, \tilde{\mu}) \\
 &\quad [\text{A, BN, AW}]
 \end{aligned}$$

Area unchanged When Slide Flipped



$$\begin{aligned}
\Delta_G(\theta, \tilde{\theta}) &= G(\theta) - G(\tilde{\theta}) - (\theta - \tilde{\theta}) \cdot g(\tilde{\theta}) \\
&= \int_{\tilde{\theta}}^{\theta} (g(\tau) - g(\tilde{\theta})) \cdot d\tau \\
&\stackrel{\text{flip}}{=} \int_{\mu}^{\tilde{\mu}} (f(\sigma) - f(\mu)) \cdot d\sigma \\
&= F(\tilde{\mu}) - F(\mu) - (\tilde{\mu} - \mu) \cdot f(\mu) \\
&= \Delta_F(\tilde{\mu}, \mu)
\end{aligned}$$

Dual divergence for Bernoulli

$$G(\theta) = \ln(1 + e^\theta) \quad F(\mu) = \mu \ln \mu + (1 - \mu) \ln(1 - \mu)$$

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

$$\Delta_G(\tilde{\theta}, \theta) = \ln(1 + e^{\tilde{\theta}}) - \ln(1 + e^\theta) - (\tilde{\theta} - \theta) \frac{e^\theta}{1 + e^\theta}$$

$$\Delta_F(\mu, \tilde{\mu}) = \mu \ln \frac{\mu}{\tilde{\mu}} + (1 - \mu) \ln \frac{1 - \mu}{1 - \tilde{\mu}}$$

Binary relative entropy

Dual divergence for Poisson

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

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

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

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

Unnormalized relative entropy

Dual matching loss for sigmoid tranfer func.

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

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

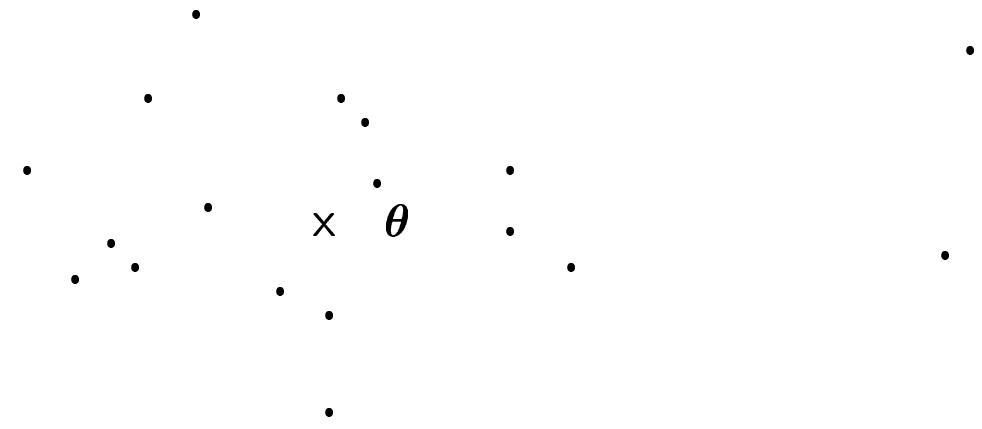
K dual to H and $k = h^{-1}$

$$\begin{aligned}\Delta_H(\mathbf{w} \cdot \mathbf{x}, h^{-1}(y)) \\ = \ln(1 + e^{\mathbf{w} \cdot \mathbf{x}}) - y\mathbf{w} \cdot \mathbf{x} + y \ln y + (1 - y) \ln(1 - y)\end{aligned}$$

By duality logistic loss is same as entropic loss

$$\begin{aligned}\Delta_K(y, h(\mathbf{w} \cdot \mathbf{x})) \\ = y \ln \frac{y}{h(\mathbf{w} \cdot \mathbf{x})} + (1 - y) \ln \frac{1 - y}{1 - h(\mathbf{w} \cdot \mathbf{x})}\end{aligned}$$

Gaussian density estimation (Fixed variance)



Off-line versus on-line

- Loss on example x_t

$$L_t(\theta) = -\ln P(x_t|\theta) = \frac{1}{2}(x_t - \theta)^2$$

Derivation of Updates

- Want to bound

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

- Off-line algorithm has all T examples

$$\{x_1, x_2, \dots, x_T\}$$

- Setup for choosing best parameter setting

$$\theta_B = \operatorname{argmin}_{\theta} (\eta_B^{-1} \Delta_G(\theta, \theta_1) + L_{1..T}(\theta))$$

divergence total
 to initial loss

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

- Equivalent to Bayesian MAP

where $\eta_B^{-1} \Delta_G(\theta, \theta_1)$ is the log of the conjugate prior

and $L_{1..T}(\theta)$ is the log of the joint likelihood

- Alternate:

$\eta_B^{-1} \Delta_G(\theta, \theta_1)$ loss on initial set of examples

On-line Algorithm

[AW]

- In trial t , the first t examples

$$\{x_1, x_2, \dots, x_t\}$$

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

$$\theta_{t+1} = \operatorname{argmin}_{\theta} (\eta_1^{-1} \Delta_G(\theta, \theta_1) + L_{1..t}(\theta))$$

divergence loss
to initial so far

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

Alternate Motivation of Same On-Line Update

$$\theta_{t+1} = \operatorname{argmin}_{\theta} (\eta_t^{-1} \Delta_G(\theta, \theta_t) + L_t(\theta))$$

divergence current
to last 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

$$\begin{aligned}\mu_{t+1} &= \frac{\eta_1^{-1} \mu_1 + \sum_{q=1}^t x_q}{\eta_1^{-1} + t} \\ &= \mu_t - \eta_{t+1}(\mu_t - x_t)\end{aligned}$$

$$\theta_{t+1} = g^{-1} (g(\theta_t) - \eta_{t+1}(\mu_t - x_t))$$

- 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]

	Off-line	Forward on-line
Gauss $\mu_1 = 0, \eta_B^{-1} = 0$	$\mu_B = \frac{\sum_{t=1}^T x_t}{T}$	$\mu_t = \frac{\sum_{q=1}^{t-1} x_q}{t}$
Bernoulli $\mu_1 = \frac{1}{2}, \eta_B^{-1} = 0$	$\mu_B = \frac{\sum_{t=1}^T x_t}{T}$	$\mu_t = \frac{\frac{1}{2} + \sum_{q=1}^{t-1} x_q}{t}$

Key Lemma

[AW]

For any example x_t and any $\theta \in \Theta$

$$\begin{aligned} & L_t(\theta_t) - L_t(\theta) \\ & \text{loss of} \quad \quad \quad \text{loss of} \\ & \text{algorithm} \quad \quad \quad \text{comparator } \theta \\ = & \eta_t^{-1} \Delta_G(\theta, \theta_t) - \eta_{t+1}^{-1} \Delta_G(\theta, \theta_{t+1}) \\ & \text{divergence} \quad \quad \quad \text{divergence} \\ & \text{to last par.} \quad \quad \quad \text{to updated par.} \\ + & \eta_{t+1}^{-1} \Delta_G(\theta_t, \theta_{t+1}) \\ & \text{cost of} \\ & \text{update} \end{aligned}$$

Main Theorem

For any sequence of examples and any $\theta \in \Theta$

$$\begin{aligned} & \sum_{t=1}^T L_t(\theta_t) - L_{1..T}(\theta) \\ & \text{total loss of} \quad \quad \quad \text{total loss of} \\ & \text{algorithm} \quad \quad \quad \text{comparator } \theta \\ \\ & = \eta_1^{-1} \Delta_G(\theta, \theta_1) - \eta_{T+1}^{-1} \Delta_G(\theta, \theta_{T+1}) \\ & \quad \quad \quad \text{divergence} \quad \quad \quad \text{divergence} \\ & \quad \quad \quad \text{to initial par.} \quad \quad \quad \text{to last par.} \\ \\ & + \sum_{t=1}^T \eta_{t+1}^{-1} \Delta_G(\theta_t, \theta_{t+1}) \\ & \quad \quad \quad \text{cost of all} \\ & \quad \quad \quad \text{updates} \end{aligned}$$

Proven by simply summing the Key Lemma

Bounds for the Forward Algorithm

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

$$\begin{aligned} \mathsf{G}_{auss} &= \sum_{t=1}^T \eta_t x_t^2 / 2 - \sum_{t=1}^{T-1} \eta_t \mu_{t+1}^2 / 2 & [AW] \\ &\leq \frac{X^2}{2} \ln\left(1 + \frac{T}{\eta_1^{-1} - 1}\right) \end{aligned}$$

where $X^2 = \max_{t=1}^T x_t^2$

$$\mathsf{B}_{Bernoulli} \leq \frac{1}{2} \ln(T+1) + 1 \quad [Fr, XB, AW]$$

$$\mathsf{lin.regr.} \leq \frac{1}{2} Y^2 n \ln\left(1 + \frac{TX^2}{a}\right) \quad [V, Fo, AW]$$

where $Y = \max_{t=1}^T y_t$

and $\mathbf{w}_t = \left(a\mathbf{I} + \sum_{q=1}^{\textcolor{red}{t}} \mathbf{x}_q \mathbf{x}_q'\right)^{-1} \sum_{q=1}^{\textcolor{red}{t-1}} \mathbf{x}_q y_q$

General Setup

- We hide some information from the learner
- The relative loss bound quantifies the price for hiding the information
- So far the future examples are hidden
Off-line algorithm knows **all** examples
On-line algorithm knows **past** examples

Minimax Algorithm for T Trials

Gaussian

[TW]

Learner against adversary

$$\inf_{\theta_1} \sup_{x_1} \inf_{\theta_2} \sup_{x_2} \inf_{\theta_3} \sup_{x_3} \dots \inf_{\theta_T} \sup_{x_T}$$

$$\sum_{t=1}^T \frac{1}{2}(\theta_t - x_t)^2 \quad - \quad \inf_{\theta} \left(\sum_{t=1}^T \frac{1}{2}(\theta - x_t)^2 \right)$$

total loss of
on-line
algorithm total loss of
 off-line
 algorithm

Instances must be bounded: $\|x_t\|_2 \leq X$

Minimax algorithm usually intractable

Bernoulli is another exception

[Sh]

Gaussian

Forward Alg.

$$\theta_t = \frac{\sum_{q=1}^{t-1} x_q}{t}$$

Bound

$$\frac{1}{2}X^2(1 + \ln T)$$

Minimax Alg.

$$\theta_t = \frac{\sum_{q=1}^{t-1} x_q}{t + \ln T - \ln(t + O(\ln T))}$$

Bound

$$\frac{1}{2}X^2(\ln T - \ln \ln T) + o(1)$$

Minimax alg. needs to know T

Last-step Minimax

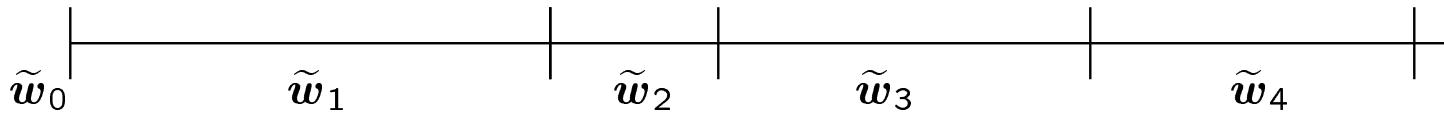
Assumes that current trial is last trial [Fo,TW]

$$\begin{aligned}\theta_t &= \underset{\theta}{\operatorname{arginf}} \sup_{x_t} \sum_{q=1}^t L_q(\theta_q) - \inf_{\theta} L_{1..t}(\theta) \\ &= \underset{\theta}{\operatorname{arginf}} \sup_{x_t} L_t(\theta_t) - \inf_{\theta} L_{1..t}(\theta)\end{aligned}$$

For Gaussian and linear regression
Last-step Minimax is same as Forward Alg.

For Bernoulli Last-step Minimax slightly better
than Forward Alg (Laplace Estimator)

Comparator shifts with time



On-line examples and on-line comparator

$$\sum_{t=1}^T L_t(\mathbf{w}_t) - \inf_{\widetilde{\mathbf{w}}_t} \sum_{t=1}^T (L_t(\widetilde{\mathbf{w}}_t) + \Delta(\widetilde{\mathbf{w}}_{t-1}, \widetilde{\mathbf{w}}_t))$$

total loss of shifting off-line comparator

on-line algorithm

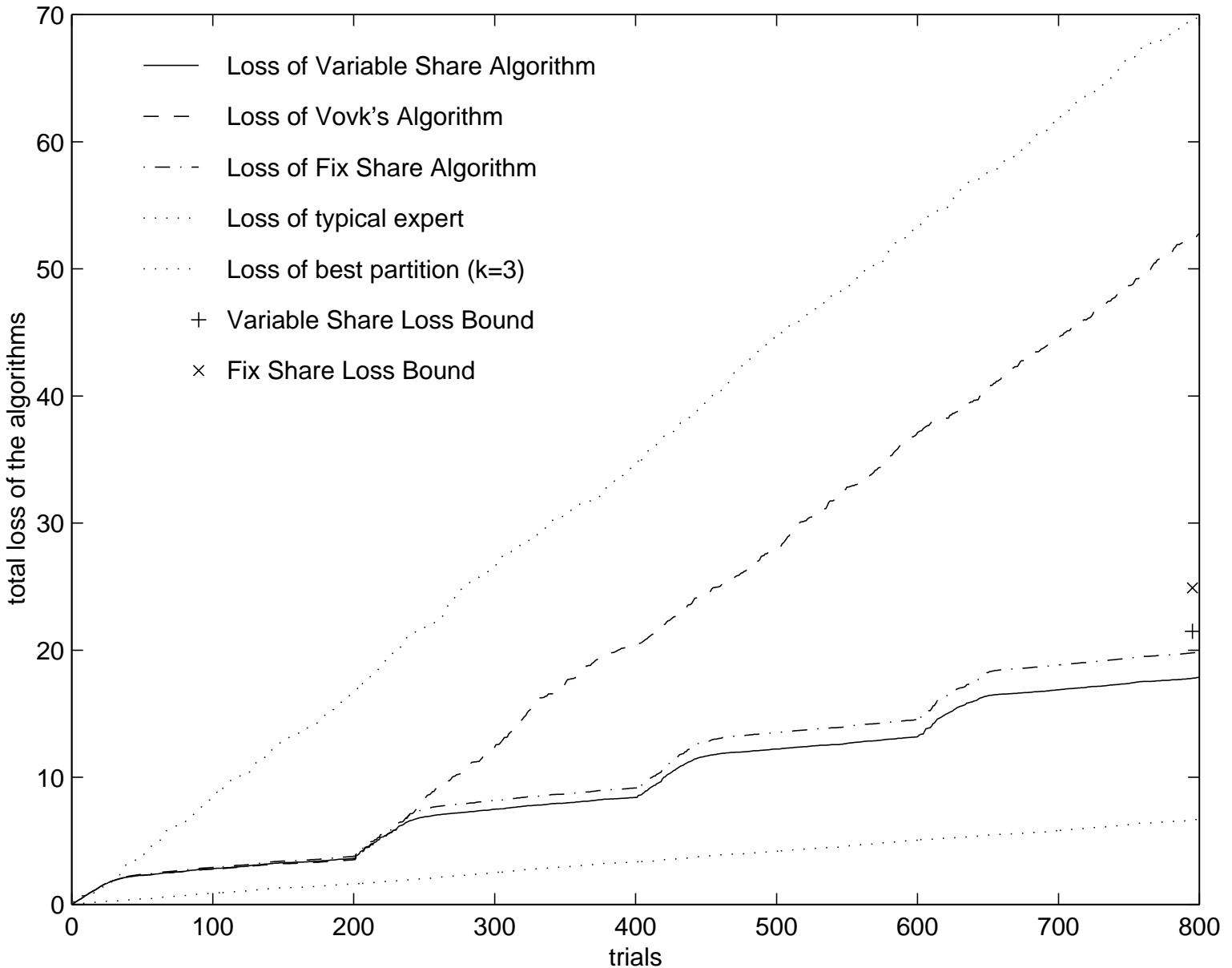
Modifications to the Expert Alg. [HW]

Predict $\hat{y}_t = v_t \cdot x_t$,
where $v_{t,i} = \frac{w_{t,i}}{\sum_{i=1}^n w_{t,i}}$

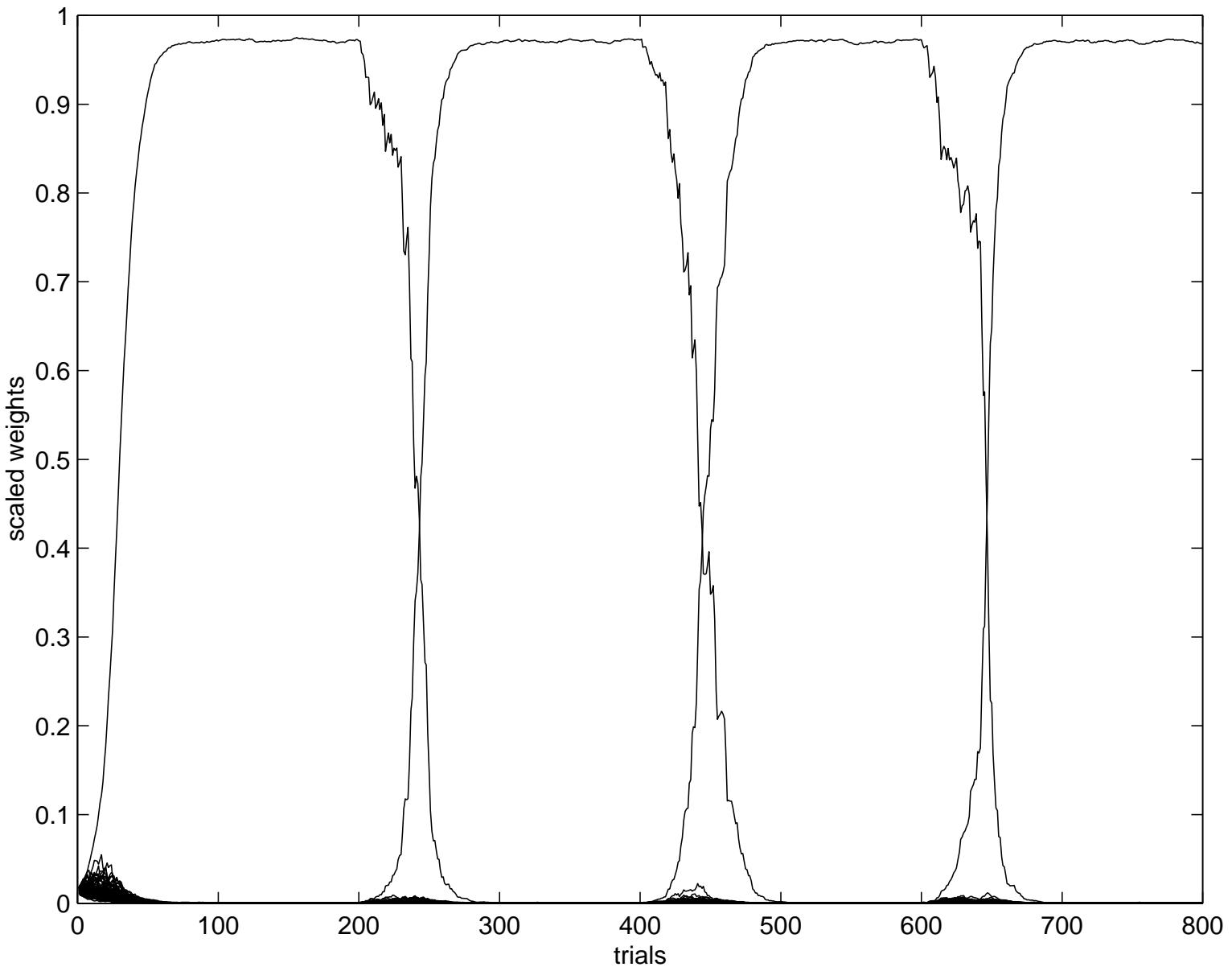
Loss Update $w_{t,i} := w_{t,i} e^{-\eta L_{y_t, x_{t,i}}}$

Share Updates ($\alpha \in [0, 1)$)

- Static-expert: Blank
- Fixed-share:
Each expert sends $\frac{\alpha}{n-1}$ of its weight to the other $n - 1$ experts
- Variable-share: Replace $\frac{\alpha}{n-1}$ by
$$\frac{1}{n-1}(1 - (1 - \alpha)^{L(y_t, x_{t,i})})$$



Loss of the share algorithms
versus Static Expert Algorithm



Relative weights of the Fixed Share Algorithm

Shifting bounds

- The Static Expert bounds

$$L_{\text{Alg}}(S) \leq \min_i L_i(S) + O(\log n)$$

become [HW]

$$L_{\text{Alg}}(S) \leq \min_P L_i(S) + O(\text{size}(P) \log n)$$

where $\text{size}(P)$ is # of shifts in partition P

- For shifting disjunctions [AW]



$$L_{\text{Alg}}(S) \leq O(\min_{\tau} A_{\tau}(S) + \text{size}(\tau) \log n)$$

where $\text{size}(\tau)$ is # of literals in τ
and $A_{\tau}(S)$ is # of attrib. errors w.r.t. τ

Applications

- Calendar managing
Many features (sleeping experts) [BI,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 [BB]

Future

- Apply clean setup for density estimation to regression and classification problems
- Other frameworks for deriving on-line algorithms such as Last-Step Minimax Alg.
- Shifting [H]
- More applications for multiplicative updates