

On-line Learning - Methods and Open Problems

Manfred K. Warmuth
UC Santa Cruz, USA

Dagstuhl, Germany
July 26, 2001

Early contributors:

Nick Littlestone

Volodya Vovk

Help with this tutorial:

Claudio Gentile

High-level Themes

- On-line versus Off-line
- Geometric versus information theoretic
Kernels partial kernels

Why on-line

- Simple algs
- Data too large
- Superiour bounds
- Data inherently on-line
- The game **MAFIA**

Overview

- Differential motivation of updates
- Motivation with Bregman divergences
- Bregman divergences as loss functions
- Pythagorean Theorem
- Rotation invariance and kernel trick
- Proving relative loss bounds
- Bregman divergences and the exponential family

On-line Linear Regression

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

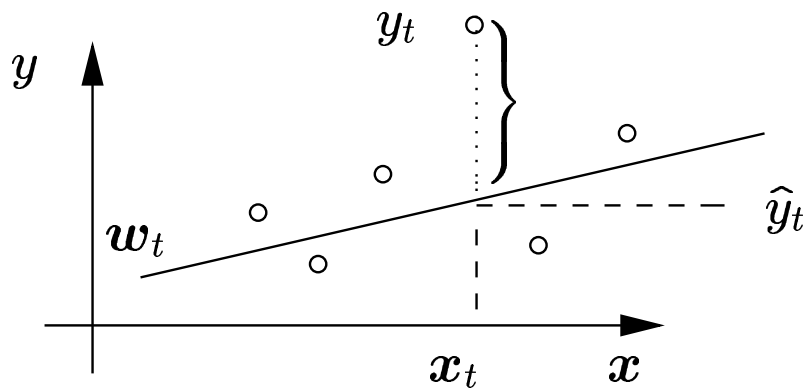
Get instance $\mathbf{x}_t \in \mathbf{R}^n$

Predict $\hat{y}_t = \mathbf{w}_t \cdot \mathbf{x}_t$

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

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

Update \mathbf{w}_t to \mathbf{w}_{t+1}



Comparison class is a set of **linear** predictors

What if no comparator consistent?

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

- $L_A(S)$ be the total loss of alg. A
- $L_u(S)$ be the total loss of linear weight vector u

Want bounds of the form:

$$\forall S : L_A(S) \leq \min_u (L_u(S) + \text{additional})$$

Bounds loss of algorithm **relative to** loss of best linear predictor u

Examples of Updates

Gradient descent

($w \in \mathbb{R}^n$)

$$w_{t+1} = w_t - \eta \nabla L_t(w_t)$$

$$= w_t - \eta (w_t \cdot x_t - y_t) x_t \quad [\text{WH}]$$

Exponentiated Gradient Algorithm

[KW]

(w is probability vector)

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

Continuous Updates

Gradient Descent

$$\mathbf{w} \in \mathbf{R}^n$$

$$\dot{\mathbf{w}}_t = -\eta \nabla_{\mathbf{w}} L_t(\mathbf{w}_t)$$

Unnormalized Exponentiated Gradient Alg.

$$\mathbf{w} \geq \mathbf{0}$$

$$\log(\dot{\mathbf{w}}_t) = -\eta \nabla_{\mathbf{w}} L_t(\mathbf{w}_t)$$

[WJ]

Characterization of algs.

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

[WJ,MW]

$F(\mathbf{w})$ convex

$$\underbrace{\dot{f}(\mathbf{w}_t)}_{\theta_t} = -\eta \nabla_{\mathbf{w}} L_t(\mathbf{w}_t)$$

Alg.	$f(\mathbf{w})$	Domain of \mathbf{w}
GD	$f(\mathbf{w}) = \mathbf{w}$	$\mathbf{w} \in \mathbf{R}^n$
EGU	$f(\mathbf{w}) = \log \mathbf{w}$	$\mathbf{w} \in [0, \infty)^n$
EG	$f(\mathbf{w}) = \ln \frac{\mathbf{w}}{1 - \ \mathbf{w}\ _1}$	$\mathbf{w} \in [0, 1]^{n-1}, \ \mathbf{w}\ _1 \leq 1$
BEG	$f(\mathbf{w}) = \ln \frac{\mathbf{w}}{1 - \mathbf{w}}$	$\mathbf{w} \in [0, 1]^n$

f is barrier function

Discretization

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

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

We use $h = 1$

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

Conjecture: **Forward Euler** better:

Replace $\nabla_{\mathbf{w}} L_t(\mathbf{w}_t)$ by $\nabla_{\mathbf{w}} L_t(\mathbf{w}_{t+h})$

Alternate Motivation

[KW]

GD

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

EG

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

Families of update algorithms

parameter divergence	name of family	update algs.
$\ w - w_t\ _2^2$	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	-	-

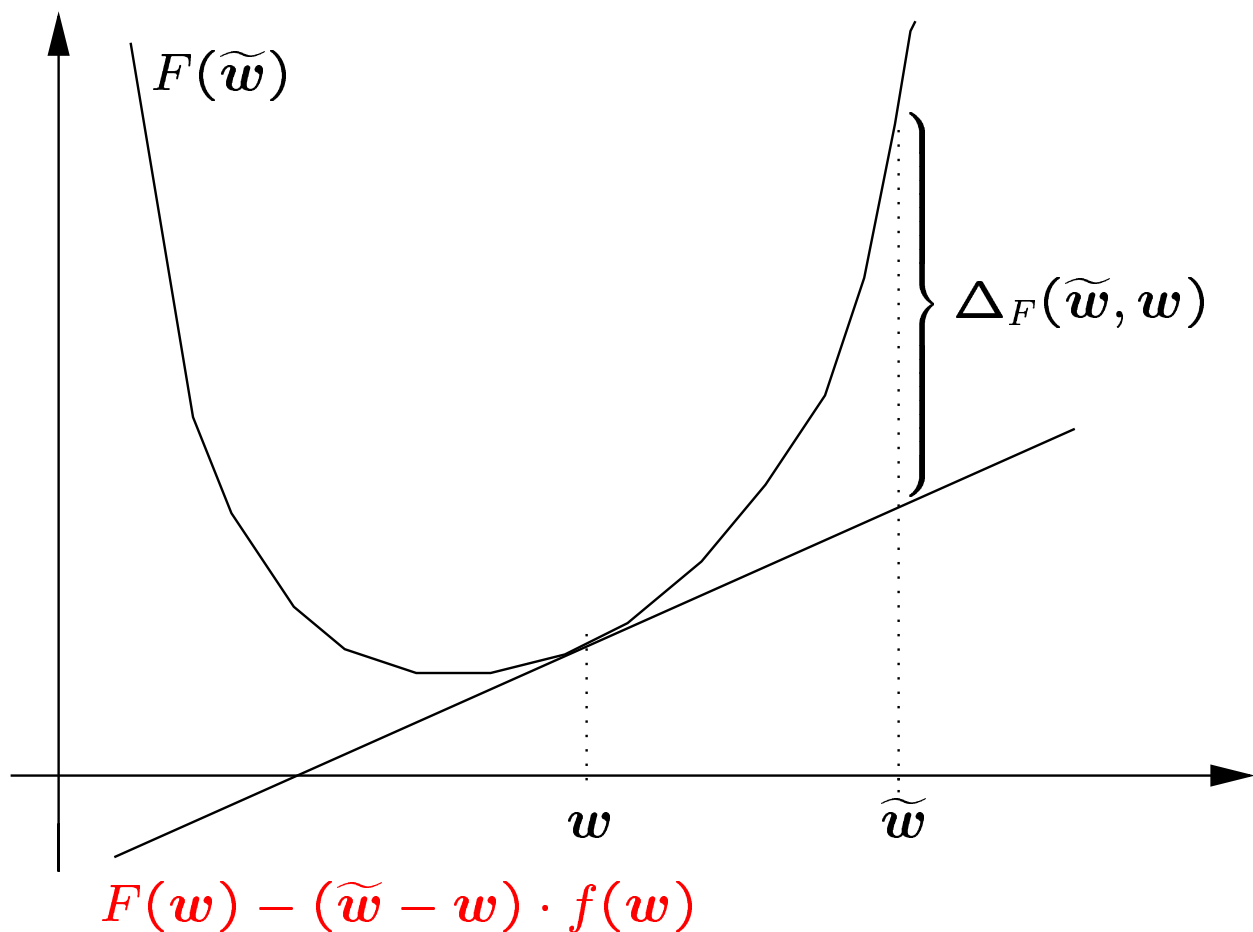
Members of different families exhibit different behavior

Bregman Divergences

[Br, CL, Cs]

For any differentiable convex function F

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



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

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

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

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

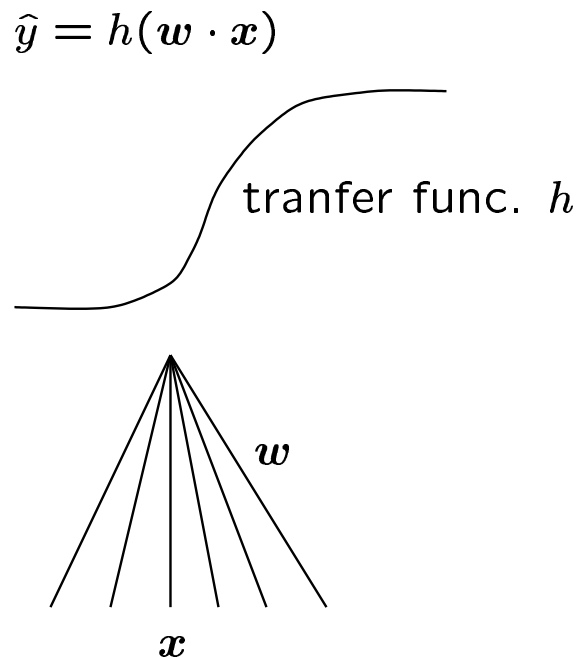
(Unnormalized) Relative Entropy

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

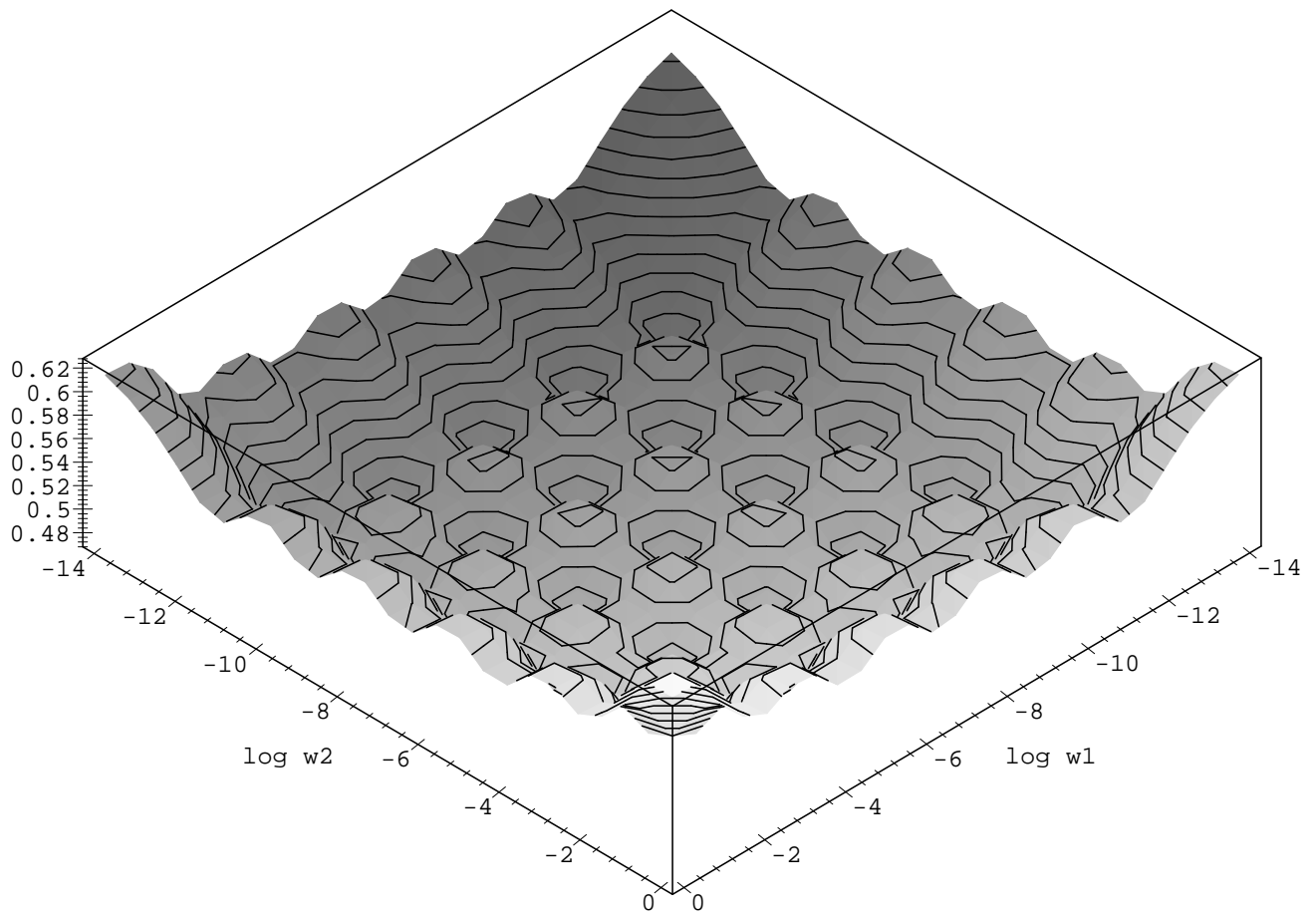
$$f(\mathbf{w}) = \ln \mathbf{w}$$

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

Bregman Divergences Lead to Good Loss Functions

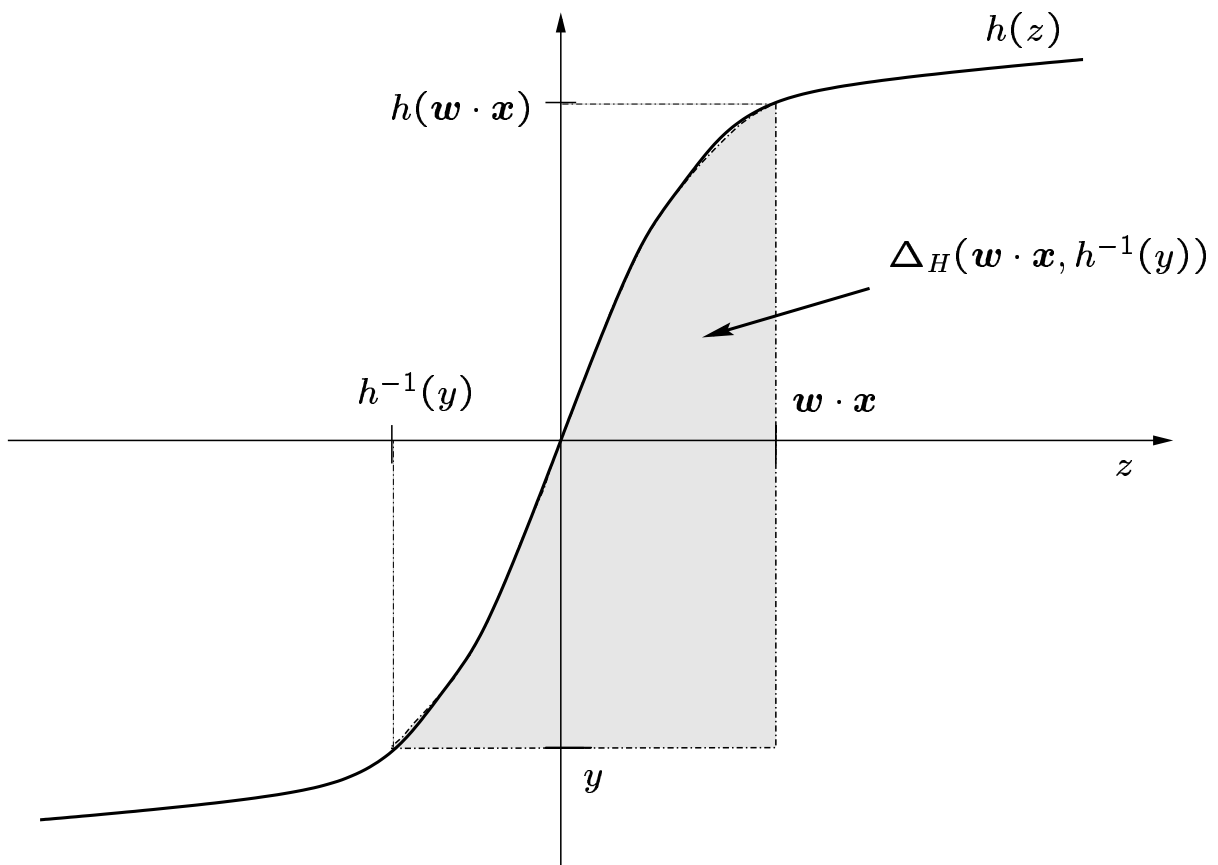


- Sigmoid function $h(z) = \frac{1}{1+e^{-z}}$
- For a set of examples $(x_1, y_1), \dots, (x_T, y_T)$
total loss $\sum_{t=1}^T (h(\mathbf{w} \cdot \mathbf{x}) - y_t)^2 / 2$
can have **exponentially many minima**
in weight space [Bu, AHW]



Want loss that is convex in w

Bregman Divergences Lead to Good Loss Functions (cont)



$$(h = \nabla H)$$

$$\int_{h^{-1}(y)}^{w \cdot x} (h(z) - y) dz$$

$$= H(w \cdot x) - H(h^{-1}(y)) - (w \cdot x - h^{-1}(y)) y$$

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

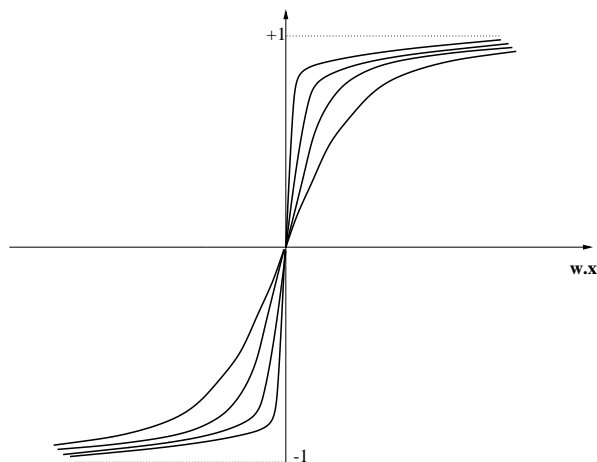
Use $\Delta_H(\mathbf{w} \cdot \mathbf{x}, h^{-1}(y))$ as loss of \mathbf{w} on (\mathbf{x}, y)

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
$\frac{e^z}{1+e^z}$	$\ln(1 + e^z)$	logistic loss
$\text{sign}(z)$	$ z $	$\max\{0, -y\mathbf{w} \cdot \mathbf{x}\}$ hinge loss

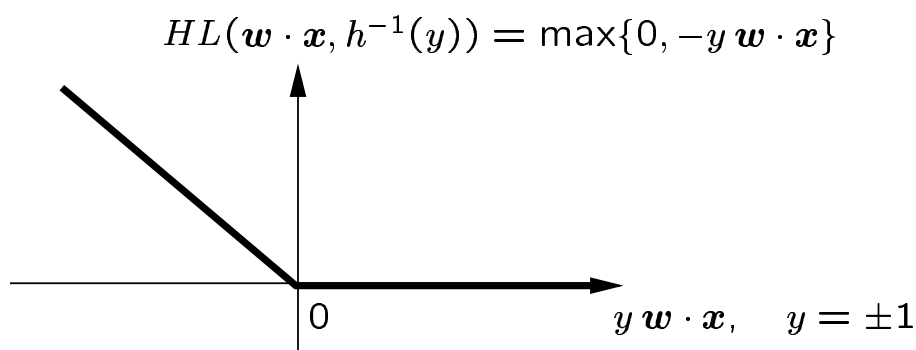
For transfer function $h(z) = \text{sign}(z)$



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

Matching loss is **hinge loss**

[GW]



Convex in 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]

$$w_{t+1} = \operatorname{argmin}_w \left(\underbrace{\Delta_F(w, w_t)}_{\text{parameter divergence}} + \eta \underbrace{\Delta_H(w \cdot x_t, h^{-1}(y_t))}_{\text{matching loss}} \right)$$

Both divergences are convex in w

$$w_{t+1} = f^{-1} (f(w_t) - \eta(h(w_t \cdot x_t) - y_t)x_t)$$

Generalization of the “delta”-rule

Projections

$$w_{t+1} = \operatorname{argmin}_w (\Delta_F(w, w_t) + \eta(w \cdot x_t - y_t)^2)$$

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

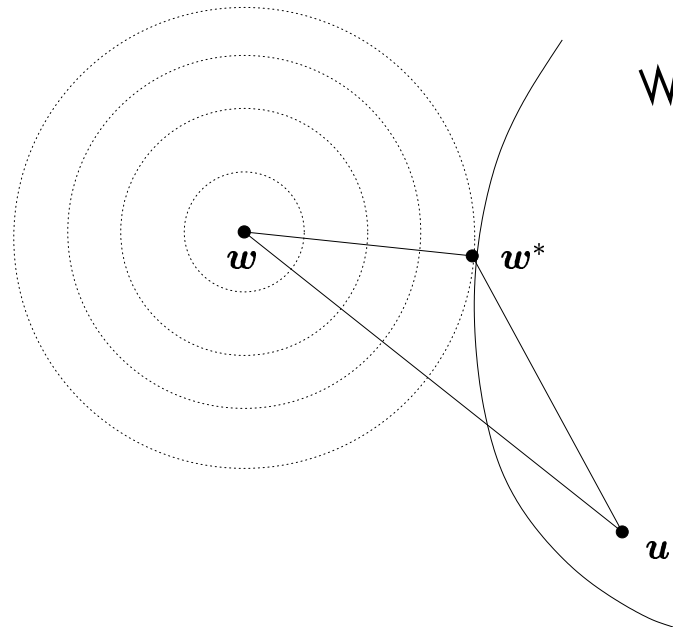
$$w_{t+1} = \operatorname{argmin}_{\{w : w \cdot x_t = y_t\}} \Delta_F(w, w_t)$$

The **AdaBoost** update of the probability vector w_t on the examples is a projection w.r.t. divergence $\Delta_F(w, w_t) = \sum_i w_i \ln \frac{w_i}{w_{t,i}}$

[La, KW]

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:

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

Kernel trick

[BGV]

- Prediction determined by $w \cdot x$
- w is linear combination of past examples

$$\begin{array}{ccc} & x \rightarrow \Phi(x) & \\ \underbrace{\left(\sum_t \alpha_t \Phi(x_t) \right)}_w \cdot \Phi(x) & = & \sum_t \alpha_t \underbrace{\Phi(x_t) \cdot \Phi(x)}_{K(x_t, x)} \end{array}$$

When trick applicable?

A linear predictions algorithm is **rotation invariant** when rotating the instances does not change the predictions

Rotation invariance $\Rightarrow w$ lin. comb. of ex.

[KWA]

When applicable (cont.)?

Bregman update:

$$\begin{aligned} f(\mathbf{w}_T) &= f(\mathbf{w}_1) + \eta \sum_t \nabla_{\mathbf{w}} L(y_t, h(\mathbf{w}_t \cdot \mathbf{x}_t)) \\ &= \underbrace{f(\mathbf{w}_1)}_0 + \text{lin. comb. of ex.} \end{aligned}$$

geometric	information theoretic
rotation invariant	no
$f = id$	$f = \log$
kernels	???

How do we prove relative loss bounds?

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

Divergence: $\Delta_F(\mathbf{u}, \mathbf{w})$

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

$L_t(\mathbf{u})$

convexity

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

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

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

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

Summing over t

[WJ,KW]

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

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

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

[AW]

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

Gradient Descent

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

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

Scaled Exponentiated Gradient

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

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

p -norm alg

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

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

Hadamard example

[KW]

$$\begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ -1 \\ -1 \end{pmatrix}$$

H unit labels
instances target

Rotation invariant alg.	$O(n)$ total loss
Information theoretic alg.	$O(\log n)$ total loss

Rotating

$$\begin{pmatrix} n & 0 & 0 & 0 \\ 0 & n & 0 & 0 \\ 0 & 0 & n & 0 \\ 0 & 0 & 0 & n \end{pmatrix} \begin{pmatrix} \frac{1}{n} \\ \frac{1}{n} \\ -\frac{1}{n} \\ -\frac{1}{n} \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ -1 \\ -1 \end{pmatrix}$$

Product of norms explains behavior

GD	$\underbrace{\ x\ _2^2}_{\sqrt{n}} \underbrace{\ w\ _2^2}_1 = n$
EG	$\underbrace{\ x\ _\infty^2}_1 \underbrace{\ w\ _1^2}_1 \log n = \log n$

$$F(w) = \frac{\|w\|_2^2}{2} \text{ only rotation invariant norm}$$

Where do Bregman divergences come from?

- Exponential family of distributions
- Inherent duality

Exponential Family of Distributions

- Parametric density functions

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

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

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

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

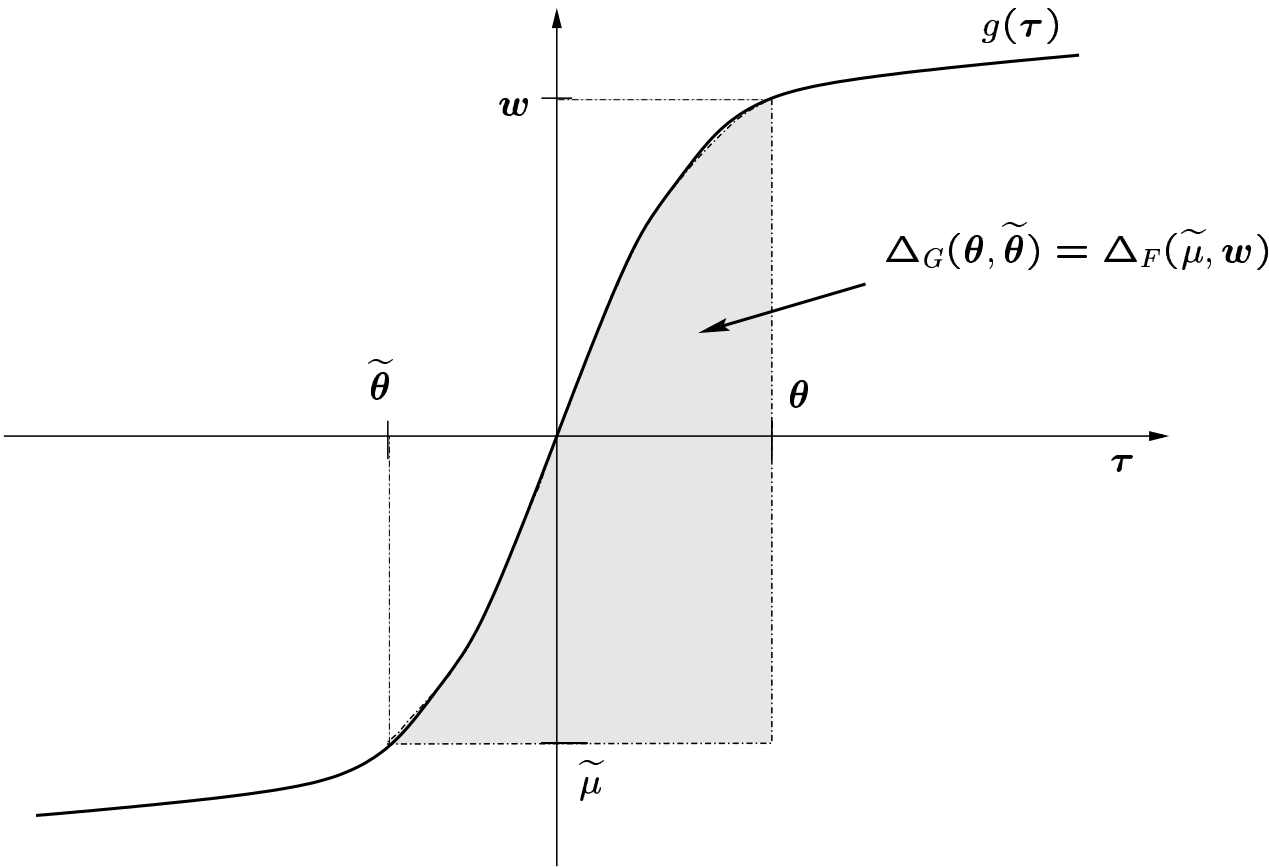
Bregman Divergences as Relative Entropies between Exponential Distributions

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

$$\begin{aligned}
 \Delta_G(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}) &= \int_{\mathbf{x}} P_G(\mathbf{x}|\boldsymbol{\theta}) \ln \frac{P_G(\mathbf{x}|\boldsymbol{\theta})}{P_G(\mathbf{x}|\tilde{\boldsymbol{\theta}})} d\mathbf{x} \\
 &= G(\tilde{\boldsymbol{\theta}}) - G(\boldsymbol{\theta}) - (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}) \cdot \underbrace{g(\boldsymbol{\theta})}_{\mathbf{w}} \\
 \underbrace{F(\mathbf{w})}_{= \boldsymbol{\theta} \cdot \mathbf{w} - G(\boldsymbol{\theta})} &= F(\mathbf{w}) - F(\tilde{\mathbf{w}}) - (\mathbf{w} - \tilde{\mathbf{w}}) \cdot \underbrace{f(\tilde{\mathbf{w}})}_{\tilde{\boldsymbol{\theta}}} \\
 &= \Delta_F(\mathbf{w}, \tilde{\mathbf{w}})
 \end{aligned}$$

[A, BN, AW]

Area unchanged When Slide Flipped



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

Max likelihood $\theta_t = \frac{\sum_{q=1}^{t-1} \mathbf{x}_q}{t-1}$

Forward Alg. $\theta_t = \frac{\sum_{q=1}^{t-1} \mathbf{x}_q}{t-1+1}$

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

Minimax Alg. $\theta_t = \frac{\sum_{q=1}^{t-1} \mathbf{x}_q}{t + \ln T - \ln(t + O(\ln T))}$

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]

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

For Gaussian and linear regression

Last-step Minimax is same as Forward Alg.

For Bernoulli Last-step Minimax slightly better than Laplace Estimator

Open problem solved

[P]

Relative entropy is **regularizer** and **barrier**

Conjecture:

$$\underset{\substack{\mathbf{w} \\ |\mathbf{w}|_1 = 1 \\ \text{constraints}}}{\operatorname{argmin}} \sum_i w_i \log \frac{w_i}{1/n}$$

and

$$\underset{\substack{\mathbf{w} \\ |\mathbf{w}|_1 = 1 \\ \text{constraints} \\ w_i \geq 0}}{\operatorname{argmin}} \sum_i w_i^2$$

behave similar!

Open problems

- When does shrinkage help?
- Kernel's for other link functions
- Better conversions to off-line bounds [CG]
- Good on-line updates for learning rates?
- Sample size $\geq \frac{d}{\epsilon} \log \frac{1}{\epsilon} + \frac{\log \frac{1}{\delta}}{\epsilon}$
produces (ϵ, δ) -good hypotheses

When can the $\log \frac{1}{\epsilon}$ factor be dropped?

Is it true when class **intersection closed** and alg. uses smallest consistent concept?

- VC dim. d
 \Rightarrow **compression scheme** of size d