

Back to the expert setting?

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

- Form of bounds: for all sequence $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$

$$L_A \leq \min_i (L_i + c \log n)$$

where c is constant

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

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
- P III: Extensions, interpretations, applications
 - Online to Batch Conversions
 - Prior information on the weight vector
 - Some applications

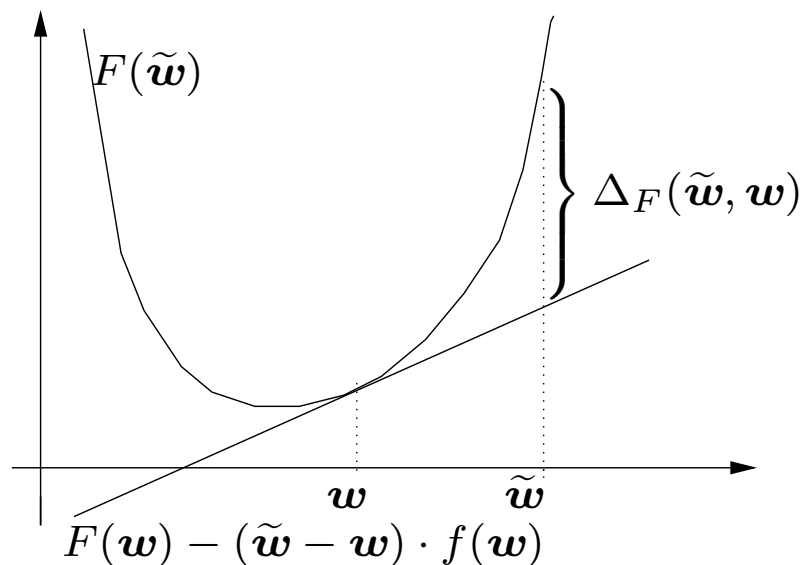
Goal: How can we prove relative loss bounds?

Bregman Divergences [Br,CL,Cs]

For **any** differentiable convex function F

$$\Delta_F(\tilde{\mathbf{w}}, \mathbf{w}) = F(\tilde{\mathbf{w}}) - F(\mathbf{w}) - (\tilde{\mathbf{w}} - \mathbf{w}) \cdot \underbrace{\nabla_{\mathbf{w}} F(\mathbf{w})}_{f(\mathbf{w})}$$

$$= F(\tilde{\mathbf{w}}) - \text{supporting hyperplane through } (\mathbf{w}, F(\mathbf{w}))$$



Bregman Divergences: Simple Properties

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

Bregman Divergences: More Properties

$$6. \nabla_{\tilde{\mathbf{w}}} \Delta_F(\tilde{\mathbf{w}}, \mathbf{w})$$

$$= \nabla F(\tilde{\mathbf{w}}) - \nabla_{\tilde{\mathbf{w}}}(\tilde{\mathbf{w}} \nabla_{\mathbf{w}} F(\mathbf{w}))$$

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

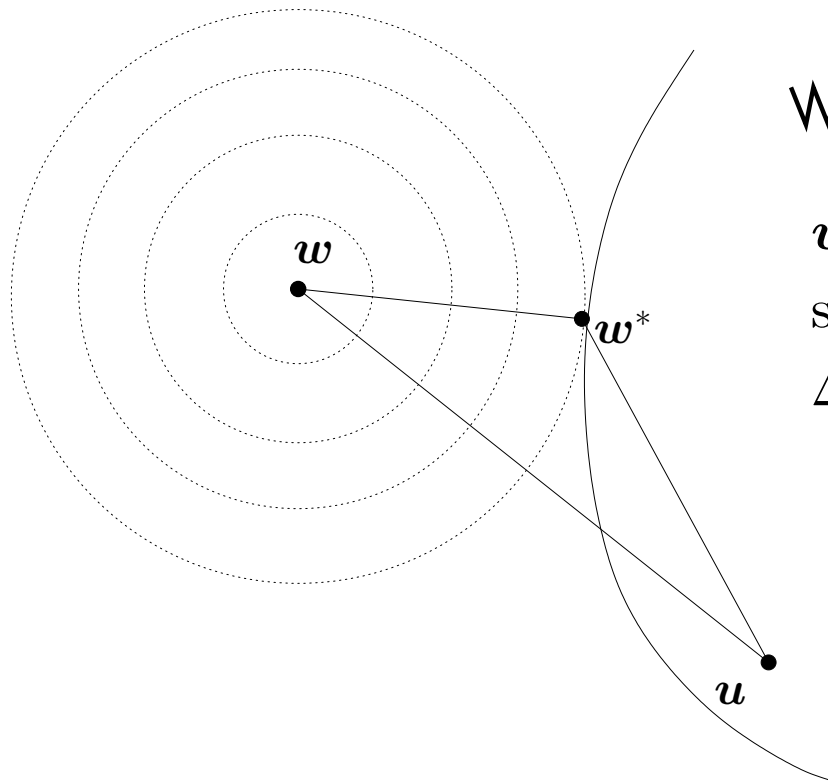
$$7. \Delta_F(\mathbf{w}_1, \mathbf{w}_2) + \Delta_F(\mathbf{w}_2, \mathbf{w}_3)$$

$$= F(\mathbf{w}_1) - F(\mathbf{w}_2) - (\mathbf{w}_1 - \mathbf{w}_2)f(\mathbf{w}_2)$$

$$F(\mathbf{w}_2) - F(\mathbf{w}_3) - (\mathbf{w}_2 - \mathbf{w}_3)f(\mathbf{w}_3)$$

$$= \Delta_F(\mathbf{w}_1, \mathbf{w}_3) + (\mathbf{w}_1 - \mathbf{w}_2) \cdot (f(\mathbf{w}_3) - f(\mathbf{w}_2))$$

A Pythagorean Theorem [Br,Cs,A,HW]



W

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

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

Theorem:

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

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*: $\frac{1}{p} + \frac{1}{q} = 1$)

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

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

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

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

General Motivation of Updates [KW]

Trade-off between two term:

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

$\Delta_F(\mathbf{w}, \mathbf{w}_t)$ 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

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

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

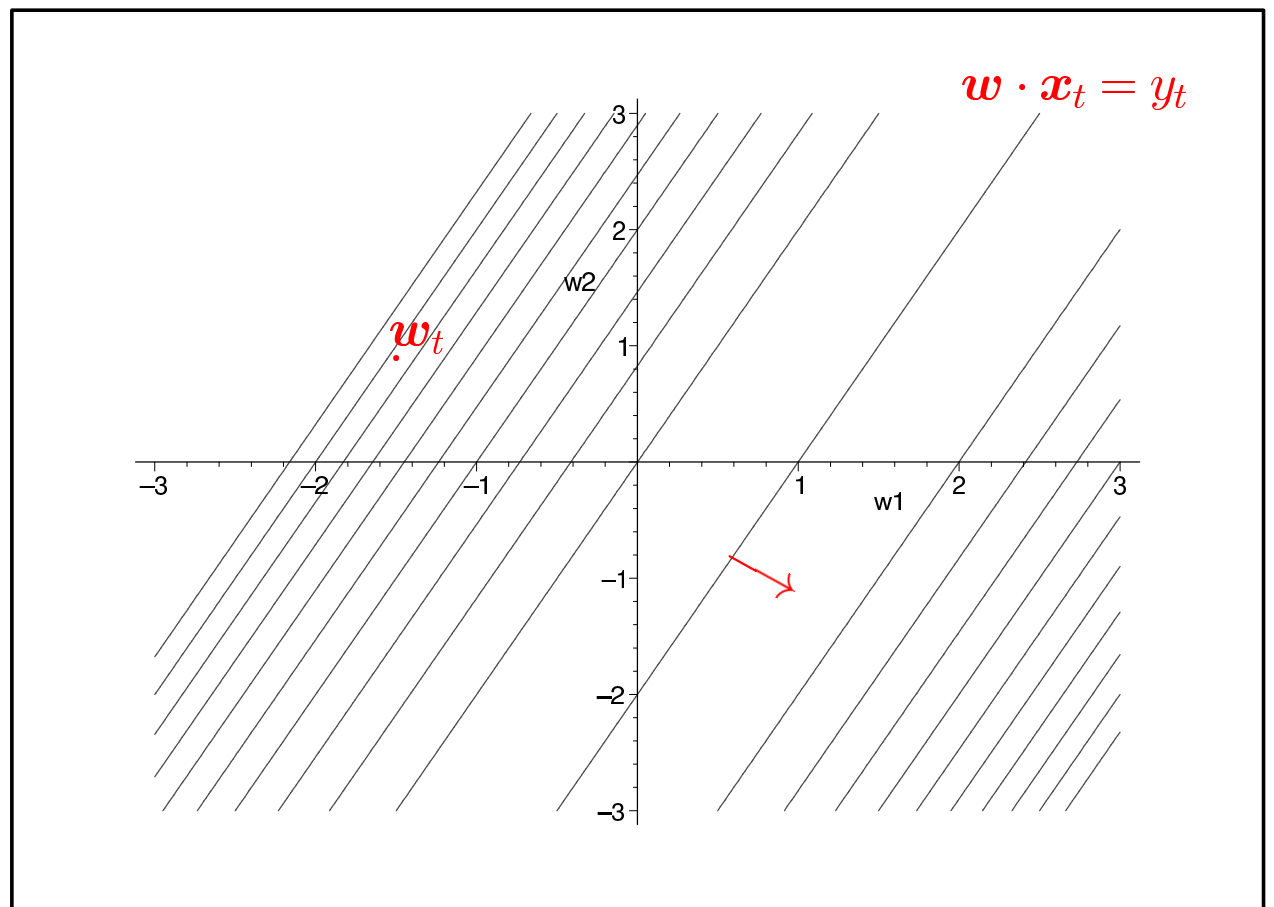
Quadratic Loss

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

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

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

$$y_t = 1$$



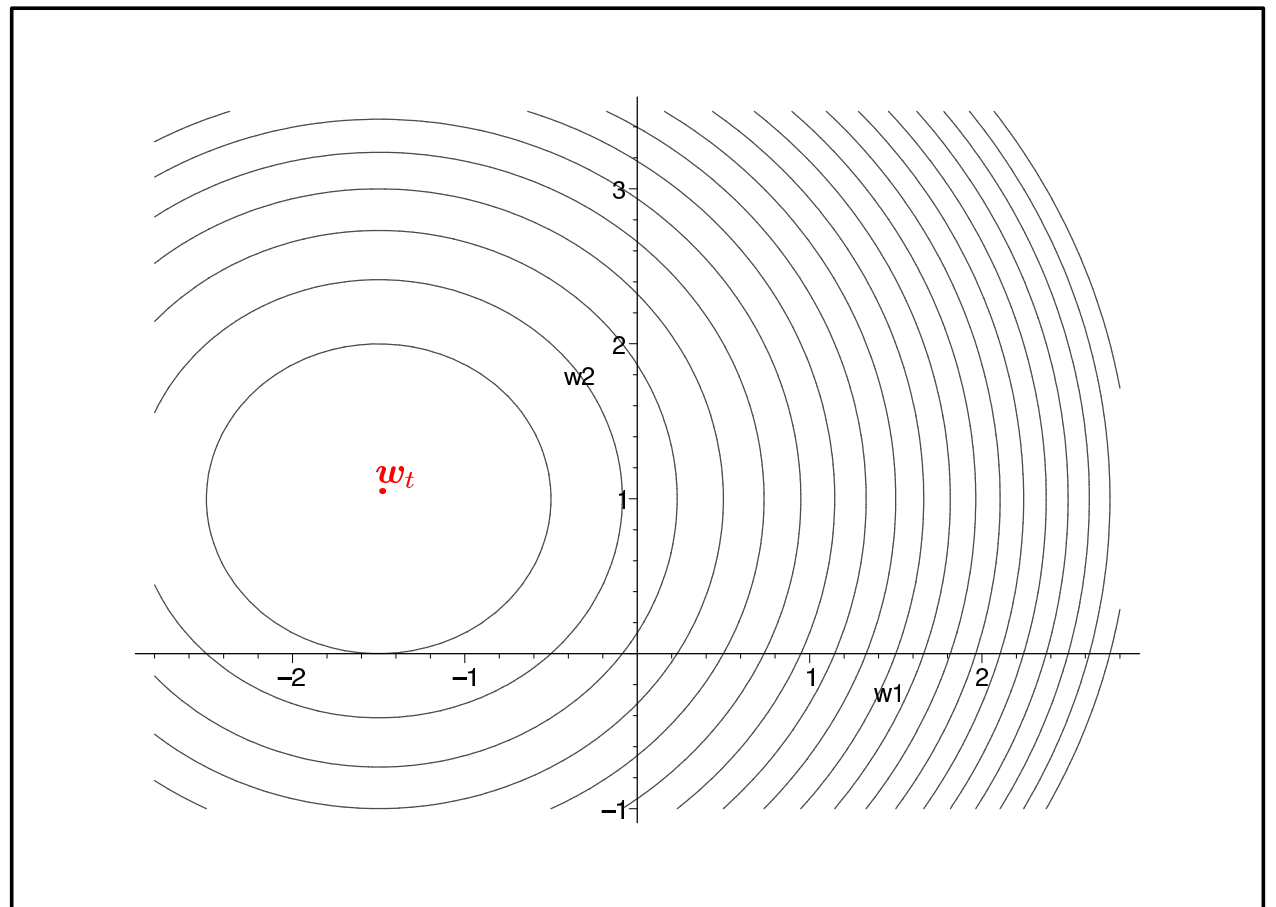
Divergence: Euclidean Distance Squared

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

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

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

$$y_t = 1$$



Loss + η Divergence

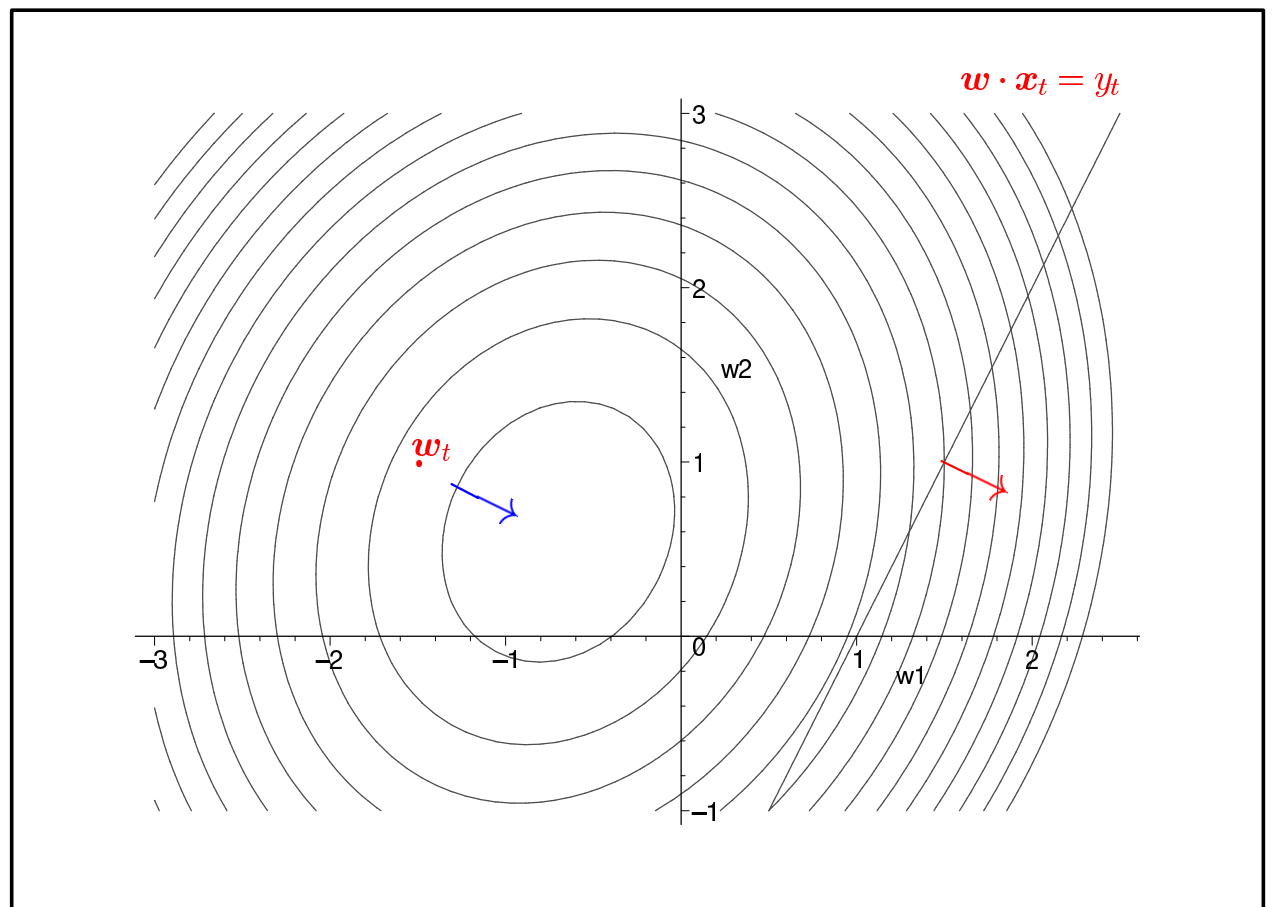
$$\Delta_F(\mathbf{w}, \mathbf{w}_t) = \|\mathbf{w} - \mathbf{w}_t\|_2^2 / 2$$

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

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

$$y_t = 1$$

$$\eta = 0.2$$



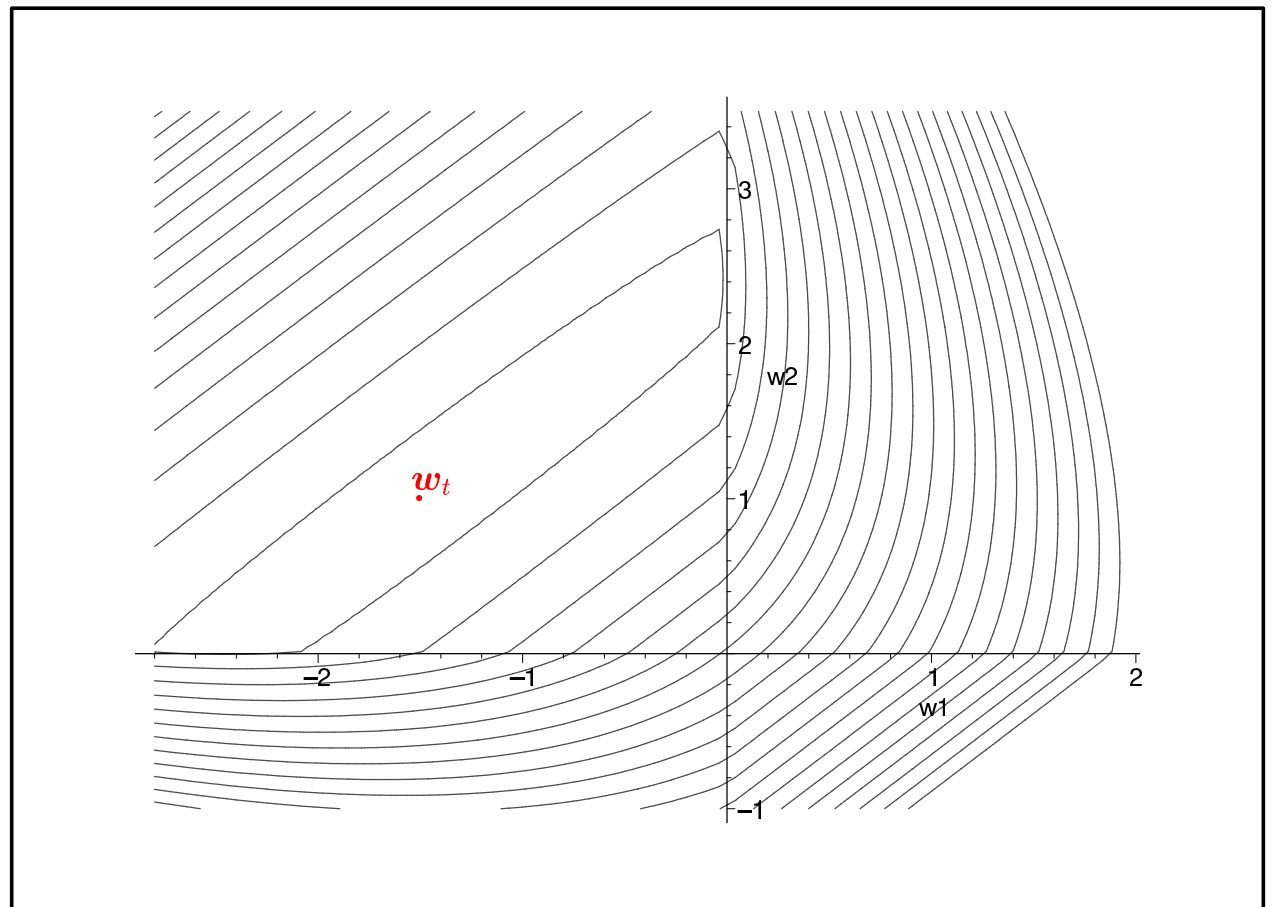
Divergence: 10-norm algorithm divergence

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

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

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

$$y_t = 1$$



Loss + η Divergence

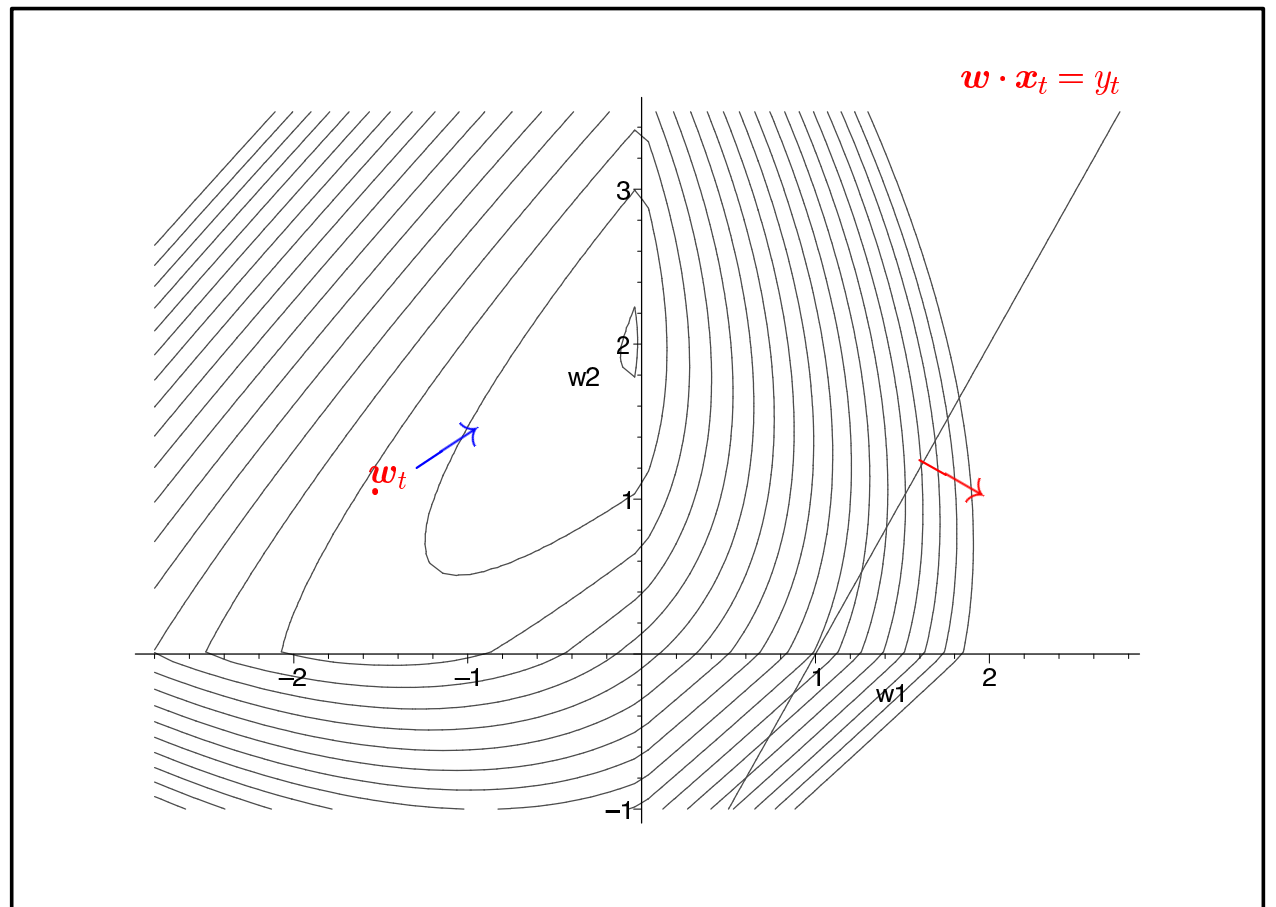
$$\Delta_F(\mathbf{w}, \mathbf{w}_t) = \|\mathbf{w} - \mathbf{w}_t\|_2^2 / 2$$

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

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

$$y_t = 1$$

$$\eta = 0.2$$



How to prove relative loss bounds?

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

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

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

$$\begin{aligned} L_t(\mathbf{u}) &\stackrel{\text{convexity}}{\geq} L_t(\mathbf{w}_t) + (\mathbf{u} - \mathbf{w}_t) \cdot \underbrace{\nabla_{\mathbf{w}} L_t(\mathbf{w}_t)}_{\text{update}} \\ &= L_t(\mathbf{w}_t) - \frac{1}{\eta} \underbrace{(\mathbf{u} - \mathbf{w}_t) \cdot (f(\mathbf{w}_{t+1}) - f(\mathbf{w}_t))}_{\text{prop. 7 of } \Delta_F} \\ &= L_t(\mathbf{w}_t) + \frac{1}{\eta} (\Delta_F(\mathbf{u}, \mathbf{w}_{t+1}) - \Delta_F(\mathbf{u}, \mathbf{w}_t) - \Delta_F(\mathbf{w}_t, \mathbf{w}_{t+1})) \end{aligned}$$

First step: Teleskopung

Summing over t

[WJ,KW]

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

Any convex loss and any Bregman divergence!

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

Loss & divergence are 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]

Bounds for 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

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

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

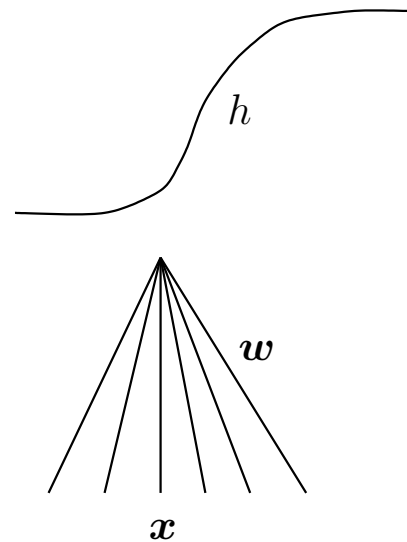
p -norm Algorithm

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

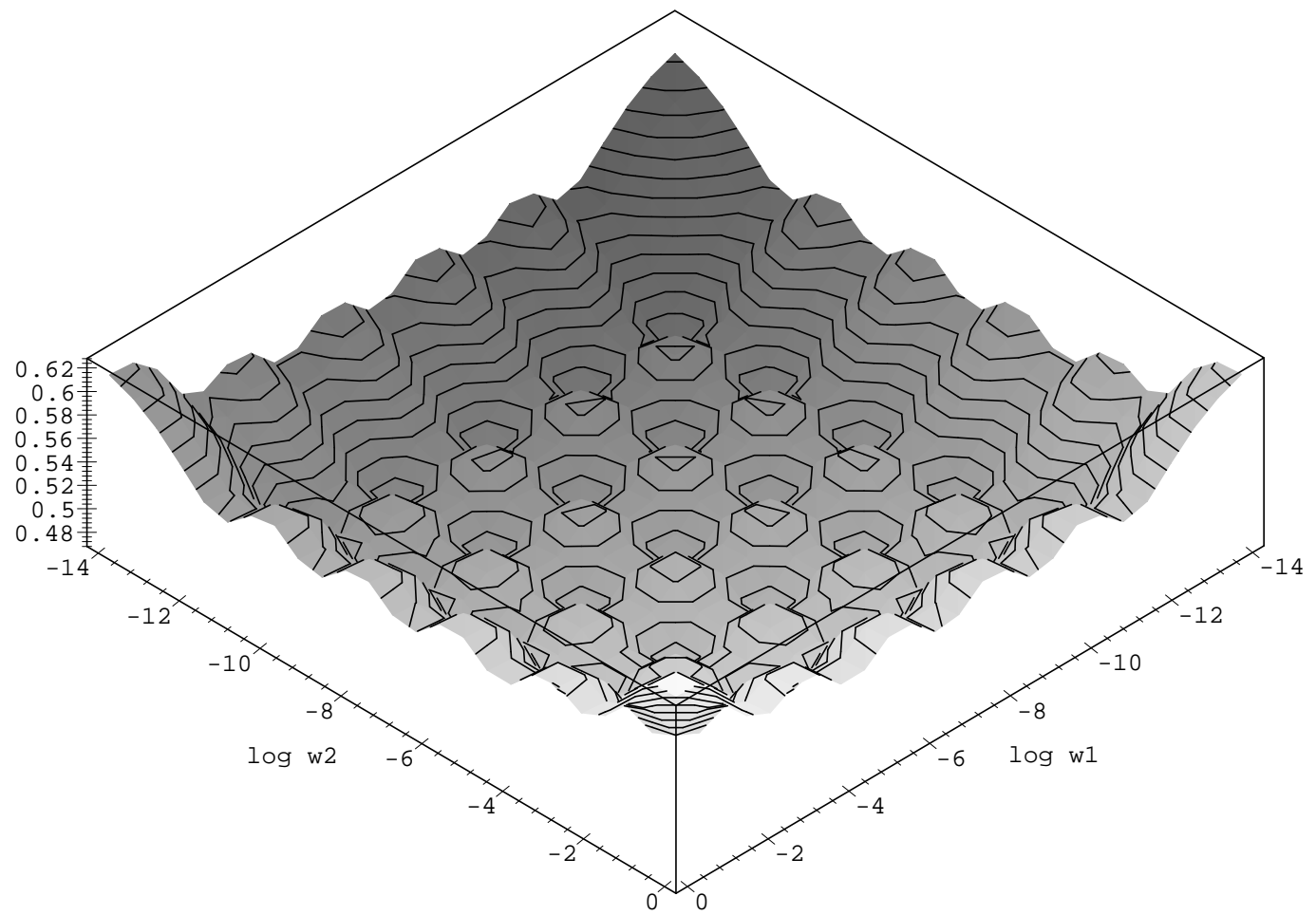
$$\|\mathbf{x}_t\|_p \leq X_p, \|\mathbf{u}\|_q \leq U_q, c > 0$$

Nonlinear Regression

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

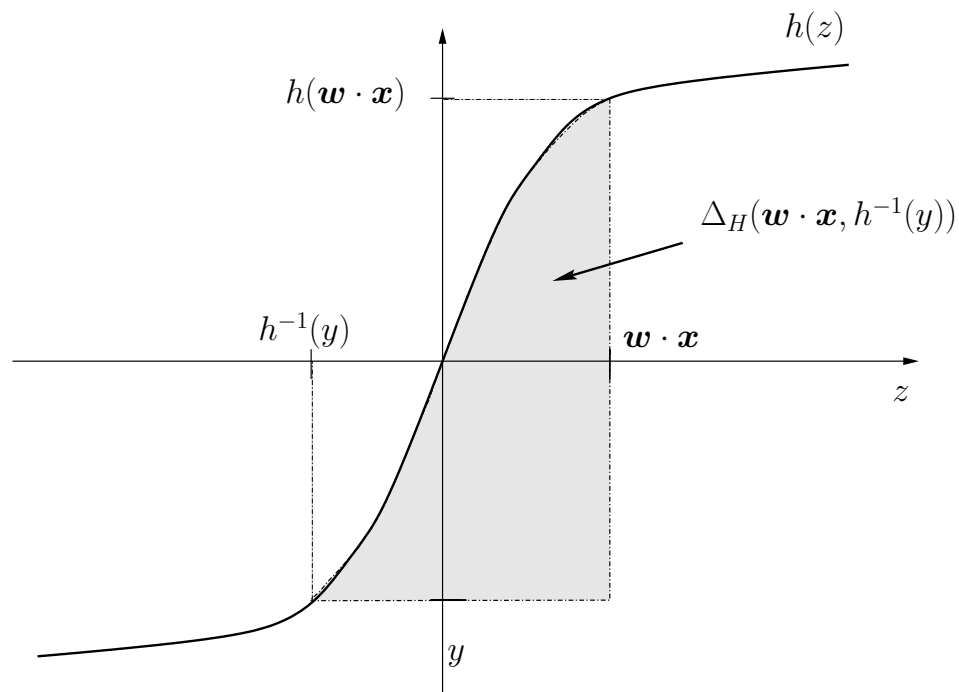


- Sigmoid function $h(z) = \frac{1}{1+e^{-z}}$
- For a set of examples $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$
total loss $\sum_{t=1}^T (h(\mathbf{w} \cdot \mathbf{x}_t) - y_t)^2 / 2$
can have **exponentially many minima**
in weight space



Want loss that is convex in w

Bregman Div. Lead to Good Loss Function



$$(h = \nabla H)$$

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

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

Idea behind the matching loss

If transfer function and loss match, then

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

Then update has simple form:

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

This can be exploited in proofs

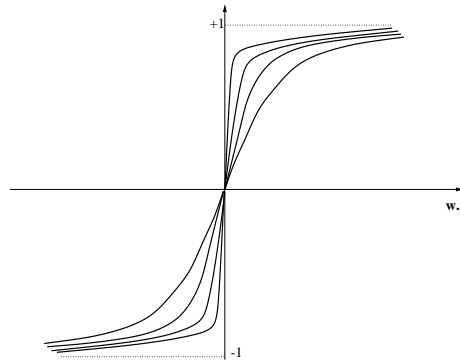
But not absolutely necessary

One only needs convexity of $L(h(\mathbf{w} \cdot \mathbf{x}), y)$ in \mathbf{w}

[Ce]

Sigmoid in the Limit

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

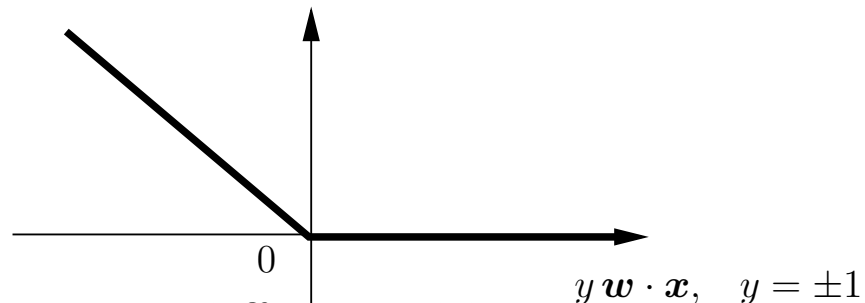


$$H(z) = |z|$$

Matching loss is **hinge loss**

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

[GW]



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 algorithms for ± 1 -classification case are
gradient-based algorithms with hinge loss

Perceptron

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

Normalized Winnow

w_{t+1}

$$= \operatorname{argmin}_{\mathbf{w}} \left(\sum_{i=1}^n w_i \ln \frac{w_i}{w_{t,i}} + \eta HL(\mathbf{w} \cdot \mathbf{x}_t, g^{-1}(y_t)) \right)$$

$$= w_{t,i} e^{-\eta (\operatorname{sign}(\mathbf{w} \cdot \mathbf{x}_t) - y_t) x_{t,i}} / \text{normalization}$$

$$\approx w_{t,i} e^{-\eta \underbrace{(\operatorname{sign}(\mathbf{w}_t \cdot \mathbf{x}_t) - y_t)}_{\hat{y}_t} x_{t,i}} / \text{normalization}$$

Trade-off between two divergences [KW]

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

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

Duality

Special case:

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

General:

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

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

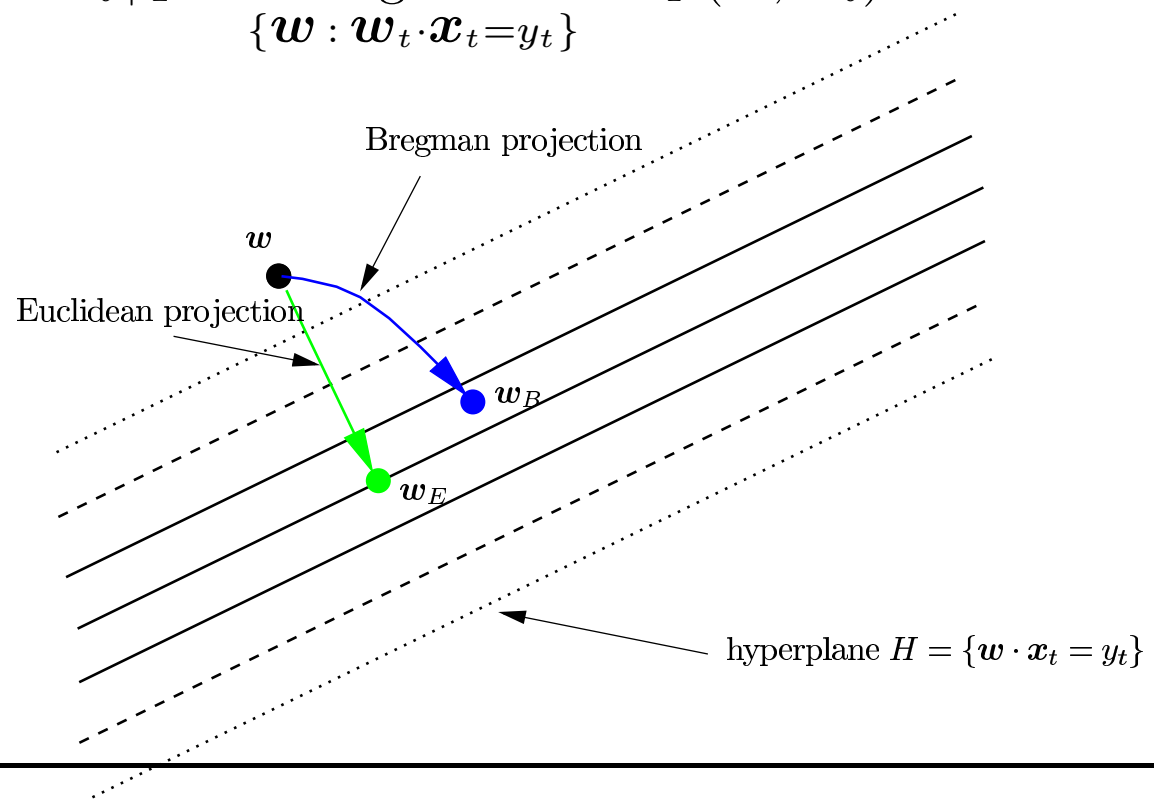
$$\mathcal{F}(\mathbf{x}) = \sup_{\mathbf{y}} \mathbf{x} \cdot \mathbf{y} - F(\mathbf{y}) = \mathbf{x} \cdot (\nabla F)^{-1}(\mathbf{x}) - F((\nabla F)^{-1}(\mathbf{x}))$$

Projections onto Hyperplanes

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

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

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



Relation to Boosting

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

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

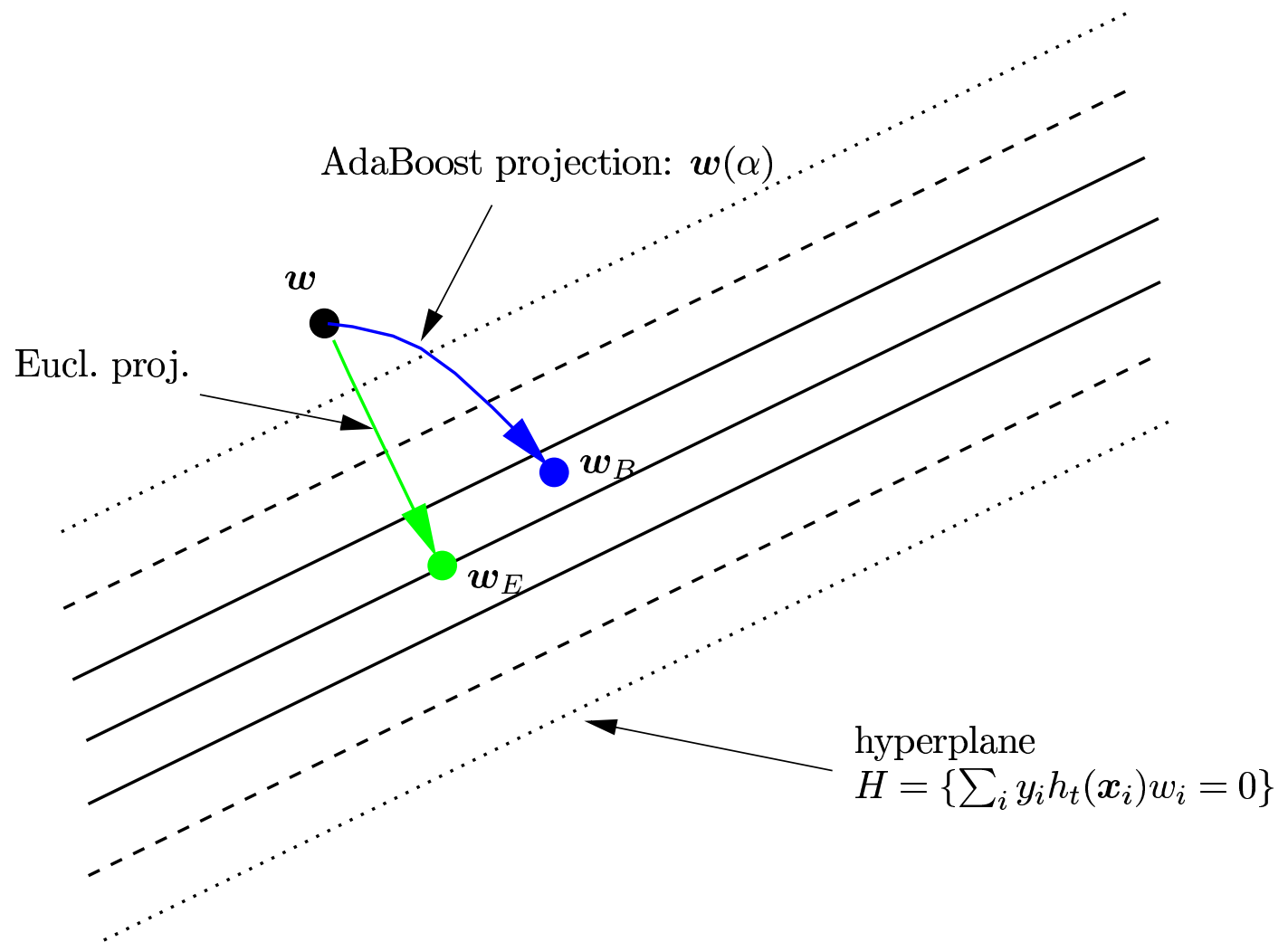
Is a projection w.r.t. divergence

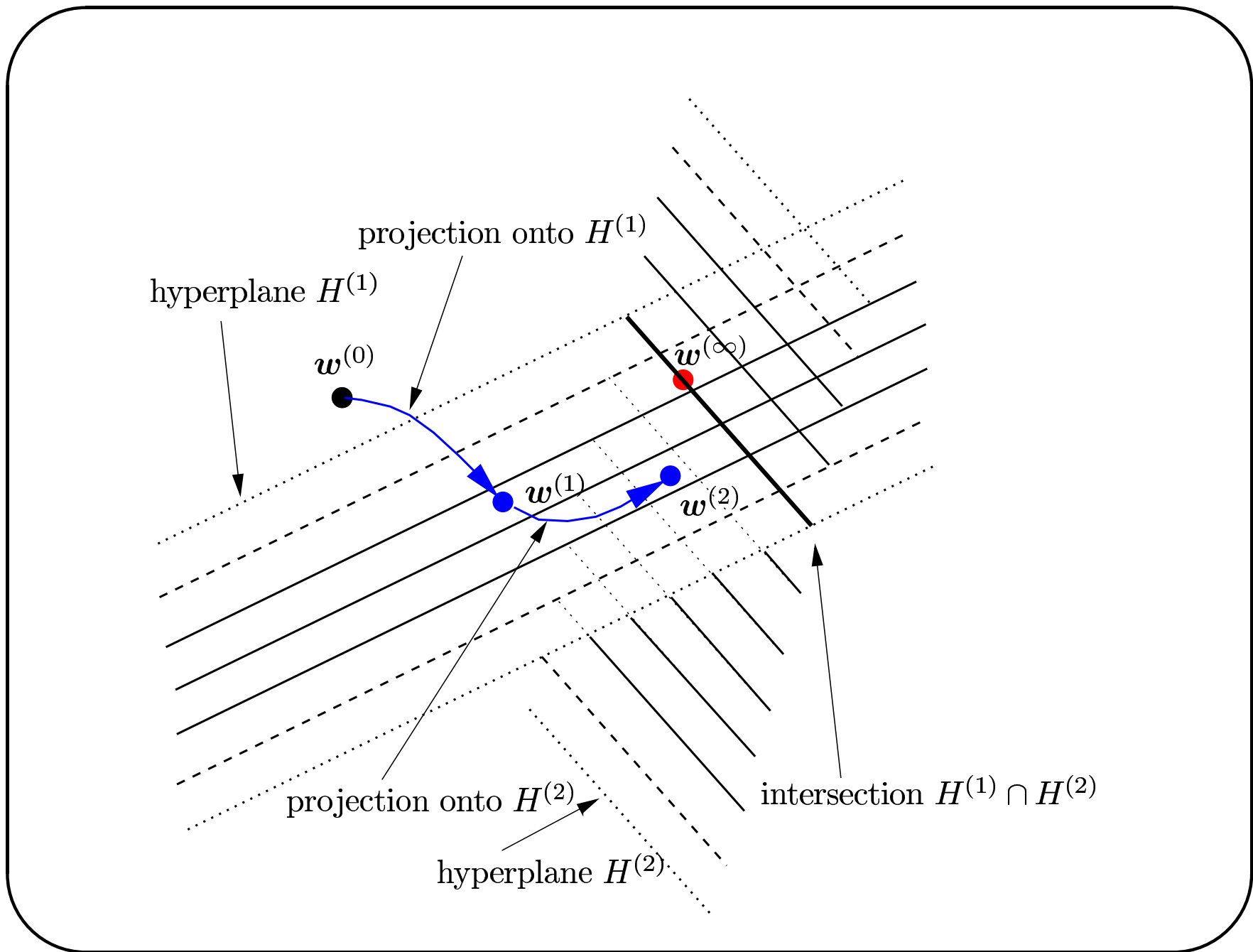
[CKW,La,KW,CSS]

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

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

$$w_i(\alpha) = w_i^{(t)} \exp(-\alpha y_i h_t(\mathbf{x}_i)) \quad \alpha = \operatorname{argmin}_{\alpha'} \sum_i w_i(\alpha')$$





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

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

- $\boldsymbol{\theta}$ 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

- 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 $g(\boldsymbol{\Theta})$

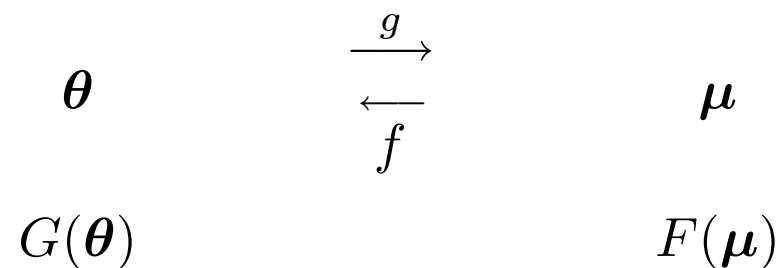
$$F(\boldsymbol{\mu}) = \boldsymbol{\theta} \cdot \boldsymbol{\mu} - G(\boldsymbol{\theta})$$

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

Primal & Dual Parameters

natural
parameter

expectation
parameter



- $\boldsymbol{\theta}$ and $\boldsymbol{\mu}$ are dual parameters
- Parameter transformations
 $g(\boldsymbol{\theta}) = \boldsymbol{\mu}$ and $f(\boldsymbol{\mu}) = \boldsymbol{\theta}$

[A,BN]

Gaussian (unit variance)

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

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: $L_t(\boldsymbol{\theta}) = \frac{1}{2}(\boldsymbol{\theta}_t - \mathbf{x}_t)^2$

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: $\theta = \ln \mu$

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

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 Div. as Rel. Ent. between Distributions

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

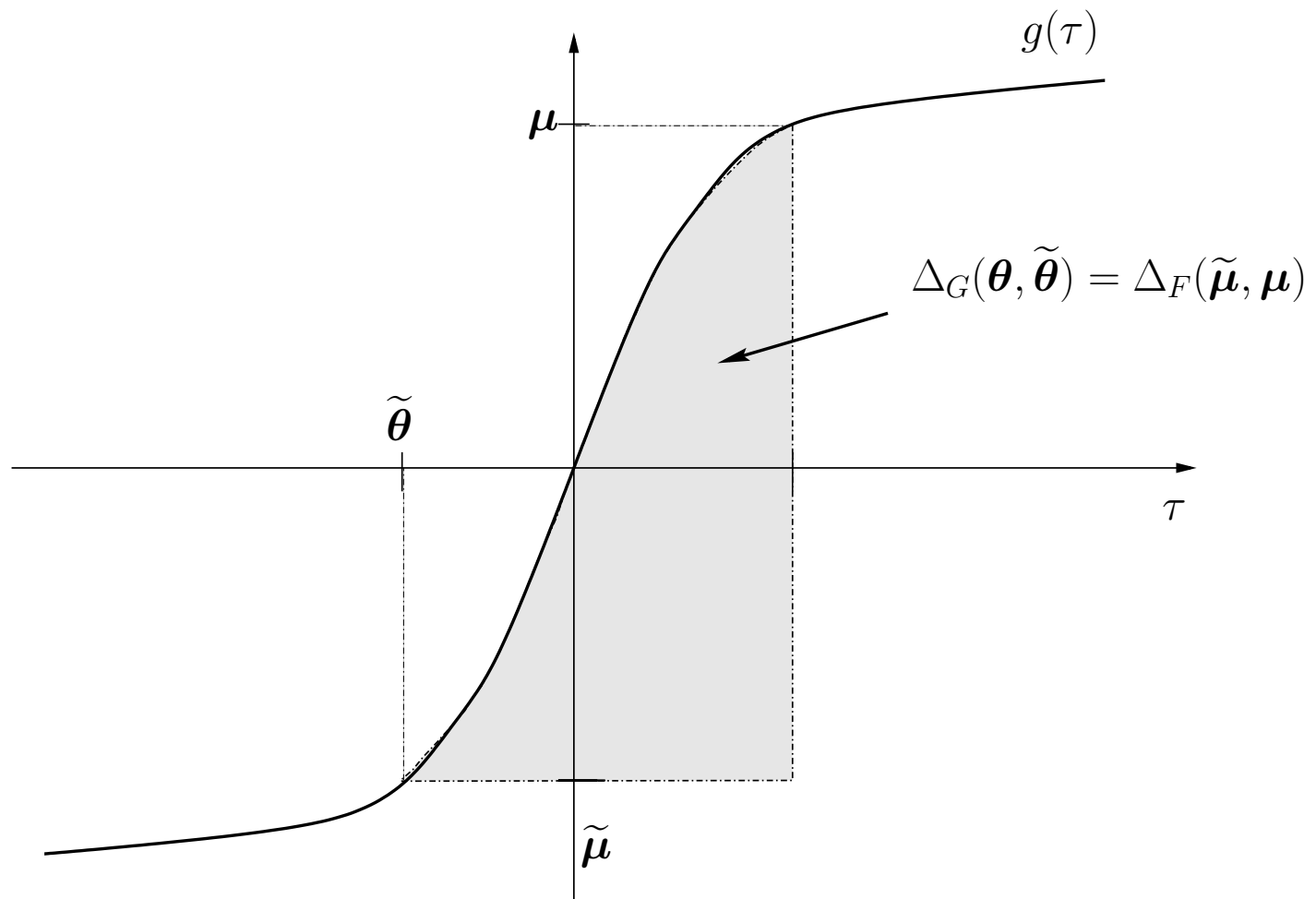
$$\Delta_G(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}) = \int_{\mathbf{x}} P_G(\mathbf{x}|\boldsymbol{\theta}) \ln \frac{P_G(\mathbf{x}|\boldsymbol{\theta})}{P_G(\mathbf{x}|\tilde{\boldsymbol{\theta}})} d\mathbf{x}$$

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

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

$$= \Delta_F(\boldsymbol{\mu}, \tilde{\boldsymbol{\mu}})$$

Area unchanged When Slide Flipped



$$\Delta_G(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) = G(\boldsymbol{\theta}) - G(\tilde{\boldsymbol{\theta}}) - (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}) \cdot g(\tilde{\boldsymbol{\theta}})$$

$$= \int_{\tilde{\boldsymbol{\theta}}}^{\boldsymbol{\theta}} (g(\boldsymbol{\tau}) - g(\tilde{\boldsymbol{\theta}})) \cdot d\boldsymbol{\tau}$$

$$\stackrel{\text{flip}}{=} \int_{\boldsymbol{\mu}}^{\tilde{\boldsymbol{\mu}}} (f(\boldsymbol{\sigma}) - f(\boldsymbol{\mu})) \cdot d\boldsymbol{\sigma}$$

$$= F(\tilde{\boldsymbol{\mu}}) - F(\boldsymbol{\mu}) - (\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu}) \cdot f(\boldsymbol{\mu})$$

$$= \Delta_F(\tilde{\boldsymbol{\mu}}, \boldsymbol{\mu})$$

Dual divergence for Bernoulli

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

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

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

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

Binary relative entropy

Dual divergence for Poisson

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

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

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

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

Unnormalized relative entropy

Dual matching loss for sigmoid transfer func.

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

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

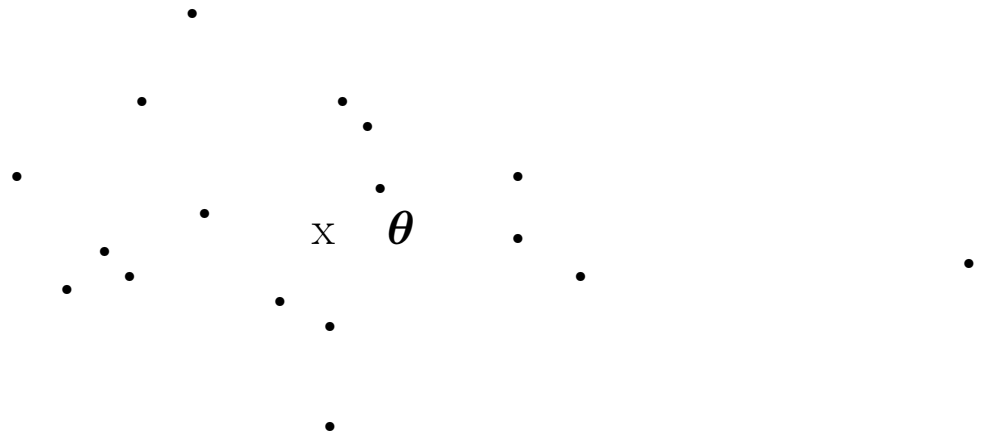
K dual to H and $k = h^{-1}$

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

By duality **logistic loss** is same as **entropic loss**

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

Example: Gaussian density estimation



Off-line versus on-line

- Loss on example \mathbf{x}_t

$$L_t(\boldsymbol{\theta}) = -\ln P(\mathbf{x}_t|\boldsymbol{\theta}) = \frac{1}{2}(\mathbf{x}_t - \boldsymbol{\theta})^2$$

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

$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$$

- Setup for choosing best parameter setting

$$\boldsymbol{\theta}_B = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} (\eta_B^{-1} \Delta_G(\boldsymbol{\theta}, \boldsymbol{\theta}_1) + L_{1..T}(\boldsymbol{\theta}))$$

divergence total
to initial loss

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

On-line Algorithm [AW]

- In trial t , the first t examples

$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t\}$$

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

$$\boldsymbol{\theta}_{t+1} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \left(\eta_1^{-1} \Delta_G(\boldsymbol{\theta}, \boldsymbol{\theta}_1) + L_{1..t}(\boldsymbol{\theta}) \right)$$

divergence loss
to initial so far

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

Alternate Motivation of Same On-Line Update

$$\boldsymbol{\theta}_{t+1} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \left(\eta_t^{-1} \Delta_G(\boldsymbol{\theta}, \boldsymbol{\theta}_t) + L_t(\boldsymbol{\theta}) \right)$$

divergence to last current loss

where $\eta_t = \frac{1}{\eta_1^{-1} + t - 1}$

Parameter Updates

Off-line: $\boldsymbol{\mu}_B = \frac{\eta_B^{-1} \boldsymbol{\mu}_1 + \sum_{t=1}^T \boldsymbol{x}_t}{\eta_B^{-1} + T}$

On-Line in trial t : $\boldsymbol{\mu}_{t+1} = \frac{\eta_1^{-1} \boldsymbol{\mu}_1 + \sum_{q=1}^t \boldsymbol{x}_q}{\eta_1^{-1} + t} = \boldsymbol{\mu}_t - \eta_{t+1}(\boldsymbol{\mu}_t - \boldsymbol{x}_t)$

$$\boldsymbol{\theta}_{t+1} = g^{-1} \left(g(\boldsymbol{\theta}_t) - \eta_{t+1}(\boldsymbol{\mu}_t - \boldsymbol{x}_t) \right)$$

- 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 $\eta_B^{-1}(\eta_B^{-1} + T)^{-1}$

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

Key Lemma [AW]

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

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

Main Theorem

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

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

Proven by simply summing the Key Lemma

Bounds for the Forward Algorithm

$$\sum_{t=1}^T L_t(\boldsymbol{\theta}_t) - \inf_{\boldsymbol{\theta}} L_{1..T}(\boldsymbol{\theta}) \stackrel{\text{Gauss}}{=} \sum_{t=1}^T \eta_t \mathbf{x}_t^2 / 2 - \sum_{t=1}^{T-1} \eta_t \boldsymbol{\mu}_{t+1}^2 / 2 \quad [\text{AW}]$$

$$\leq \frac{X^2}{2} \ln\left(1 + \frac{T}{\eta_1^{-1} - 1}\right)$$

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

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

$$X^2 = \max_{t=1}^T \mathbf{x}_t^2, \quad Y = \max_{t=1}^T y_t, \quad \mathbf{w}_t = \left(a\mathbf{I} + \sum_{q=1}^t \mathbf{x}_q \mathbf{x}_q' \right)^{-1} \sum_{q=1}^{t-1} \mathbf{x}_q y_q$$

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

Learner against adversary

$$\inf_{\boldsymbol{\theta}_1} \sup_{\mathbf{x}_1} \inf_{\boldsymbol{\theta}_2} \sup_{\mathbf{x}_2} \inf_{\boldsymbol{\theta}_3} \sup_{\mathbf{x}_3} \dots \inf_{\boldsymbol{\theta}_T} \sup_{\mathbf{x}_T}$$

$$\sum_{t=1}^T \frac{1}{2} (\boldsymbol{\theta}_t - \mathbf{x}_t)^2 \quad - \quad \inf_{\boldsymbol{\theta}} \left(\sum_{t=1}^T \frac{1}{2} (\boldsymbol{\theta} - \mathbf{x}_t)^2 \right)$$

total loss of

total loss of

on-line

off-line

algorithm

algorithm

Instances must be bounded: $\|\mathbf{x}_t\|_2 \leq X$

Minimax algorithm usually **intractable**

Gaussian and Bernoulli are exceptions

[TW,Sh]

Gaussian

Forward Alg. $\boldsymbol{\theta}_t = \frac{\sum_{q=1}^{t-1} \mathbf{x}_q}{t}$

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

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

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]

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

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)