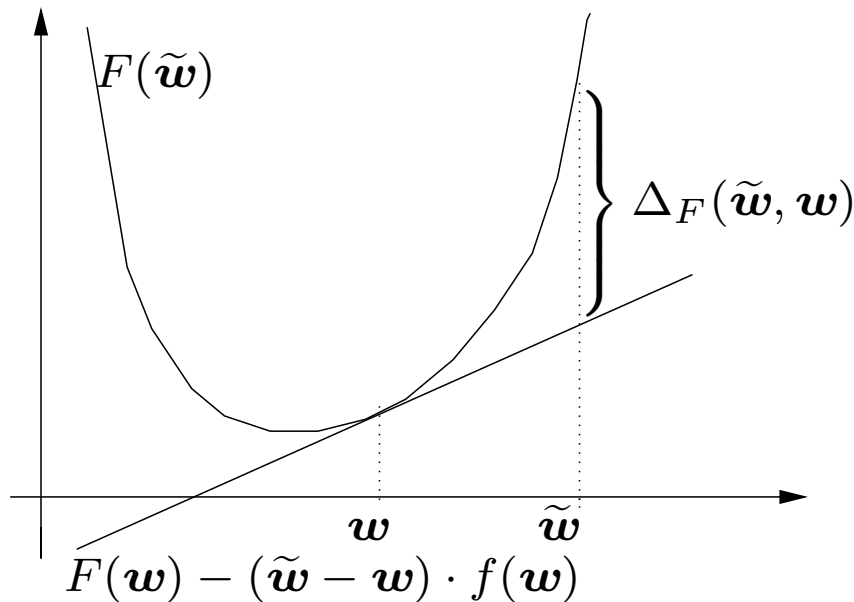


Bregman Divergences [Br,CL,Cs]

For **any** differentiable convex function F

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

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



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

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

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

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

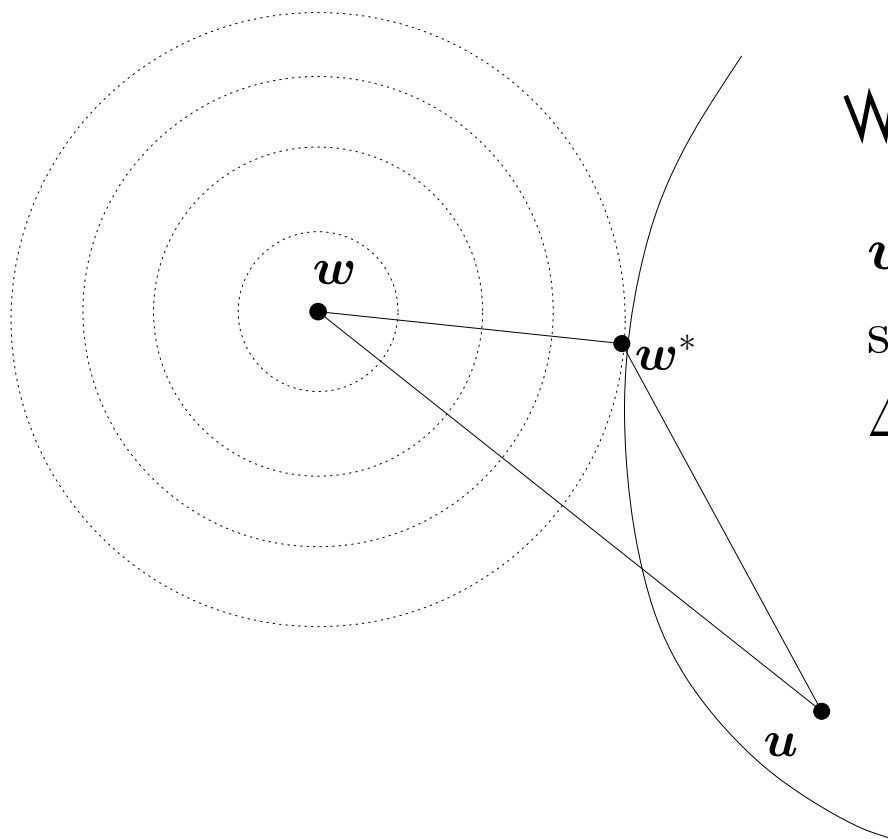
$$7. \Delta_F(\mathbf{w}_1, \mathbf{w}_2) + \Delta_F(\mathbf{w}_2, \mathbf{w}_3)$$

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

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

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

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



\mathcal{W}

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

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

Theorem:

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

Examples

Squared Euclidean Distance

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

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

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

(Unnormalized) Relative Entropy

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

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

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

Examples-2 [GLS, GL]

p-norm Algs (*q* is dual to *p*: $\frac{1}{p} + \frac{1}{q} = 1$)

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

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

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

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

Examples-3

Burg entropy

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

$$f(\boldsymbol{w}) = -\frac{1}{\boldsymbol{w}}$$

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

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

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

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

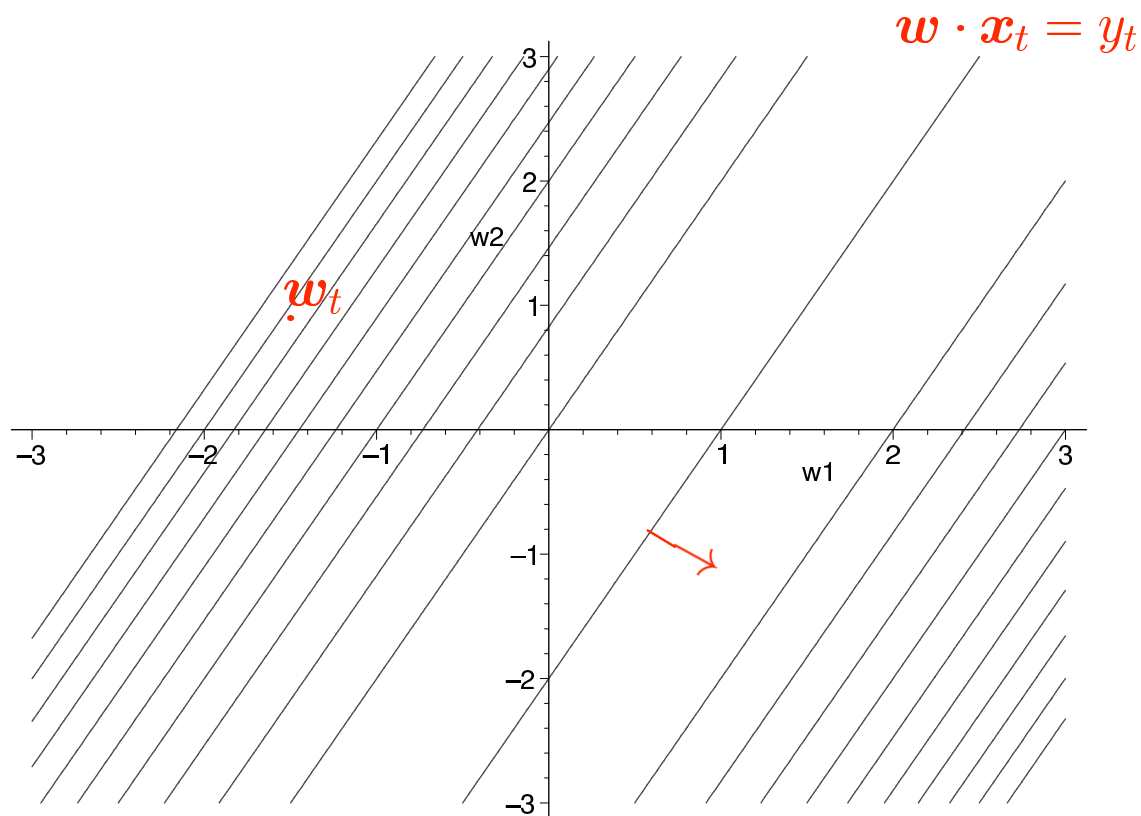
Quadratic Loss

$$L_t(\mathbf{w}) = \frac{1}{2}(\mathbf{y}_t - \mathbf{w} \cdot \mathbf{x}_t)^2$$

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

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

$$y_t = 1$$



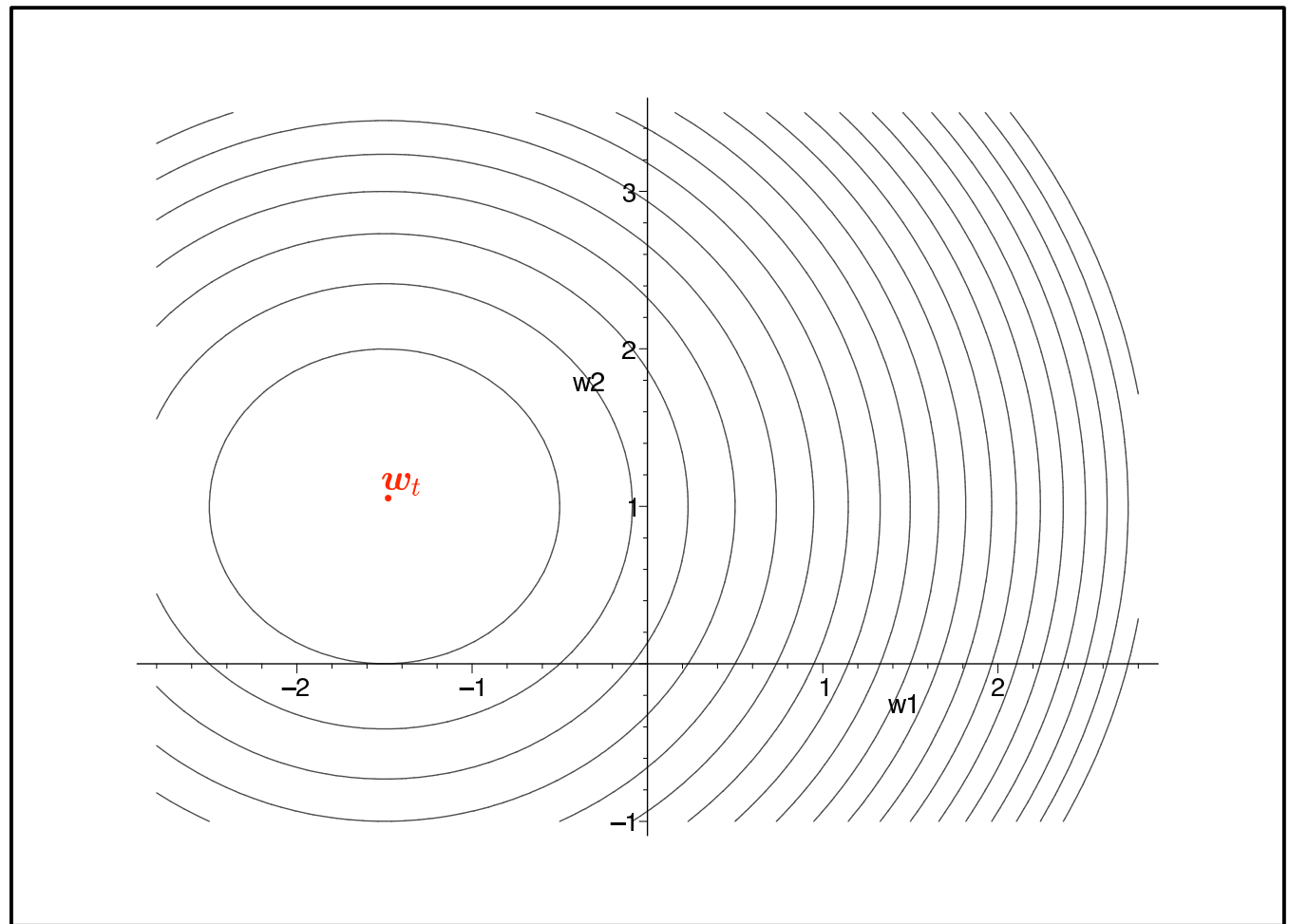
Divergence: Euclidean Distance Squared

$$\Delta_F(\mathbf{w}, \mathbf{w}_t) = \frac{1}{2} \|\mathbf{w} - \mathbf{w}_t\|_2^2$$

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

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

$$y_t = 1$$



Divergence + η Loss

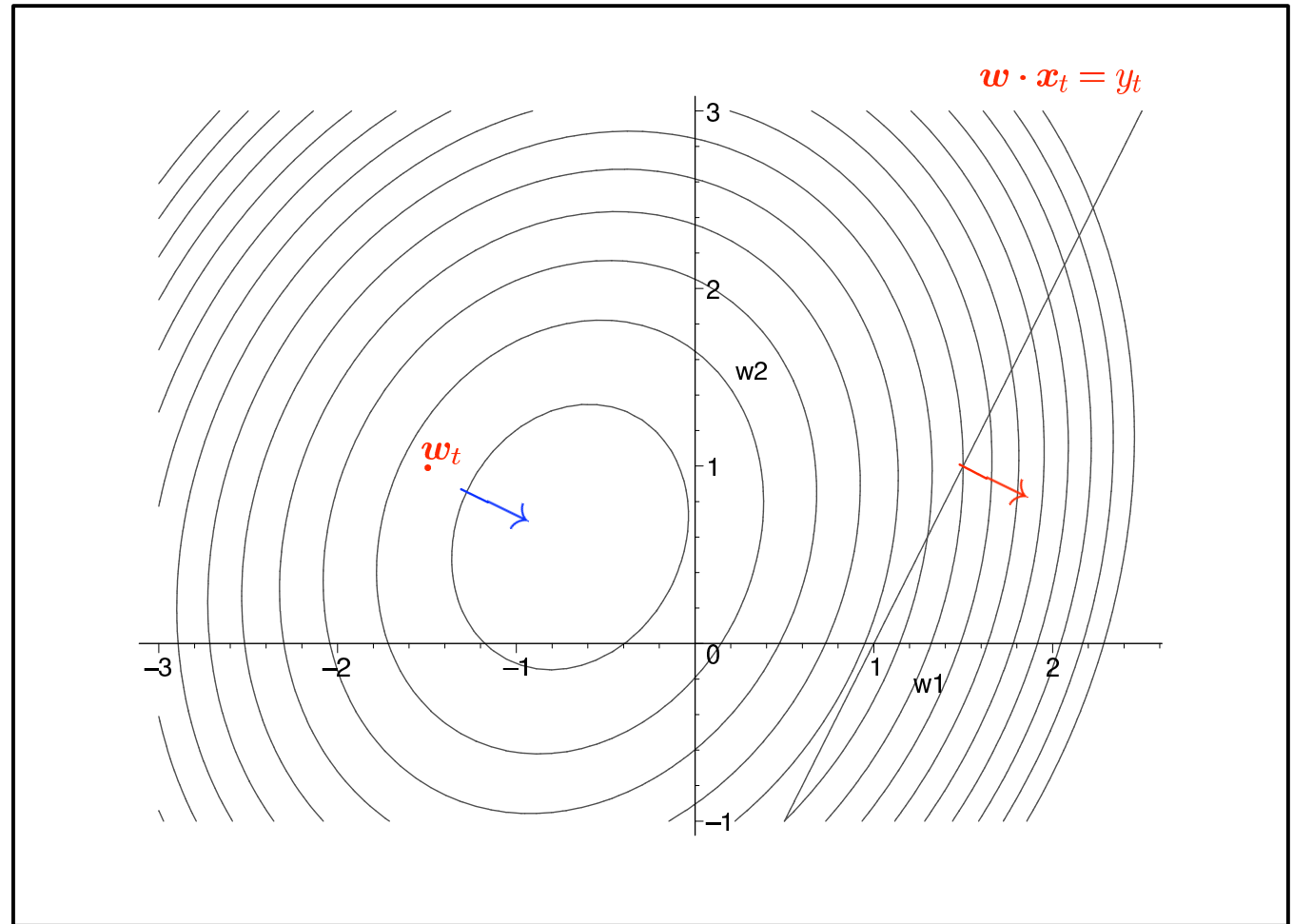
$$\frac{1}{2} \|w - w_t\|_2^2 + \eta \frac{1}{2} (y_t - w \cdot x_t)^2$$

$$w_t = (-3/2, 1)$$

$$x_t = (1, -0.5)$$

$$y_t = 1$$

$$\eta = 0.2$$



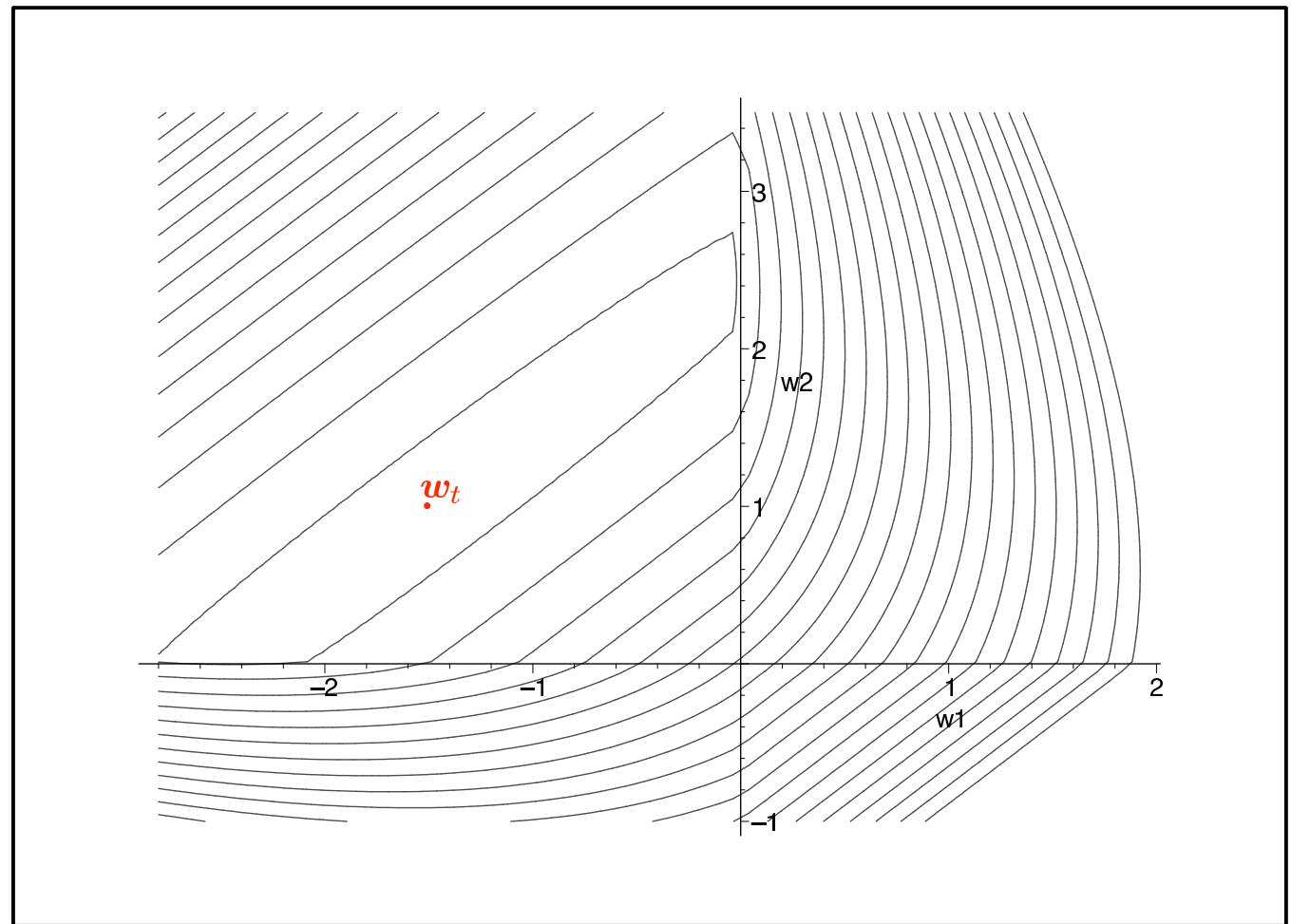
Divergence: 10-norm algorithm divergence

$$\Delta_F(\mathbf{w}, \mathbf{w}_t) \text{ where } F(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_{10}^2$$

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

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

$$y_t = 1$$



Loss + η Divergence

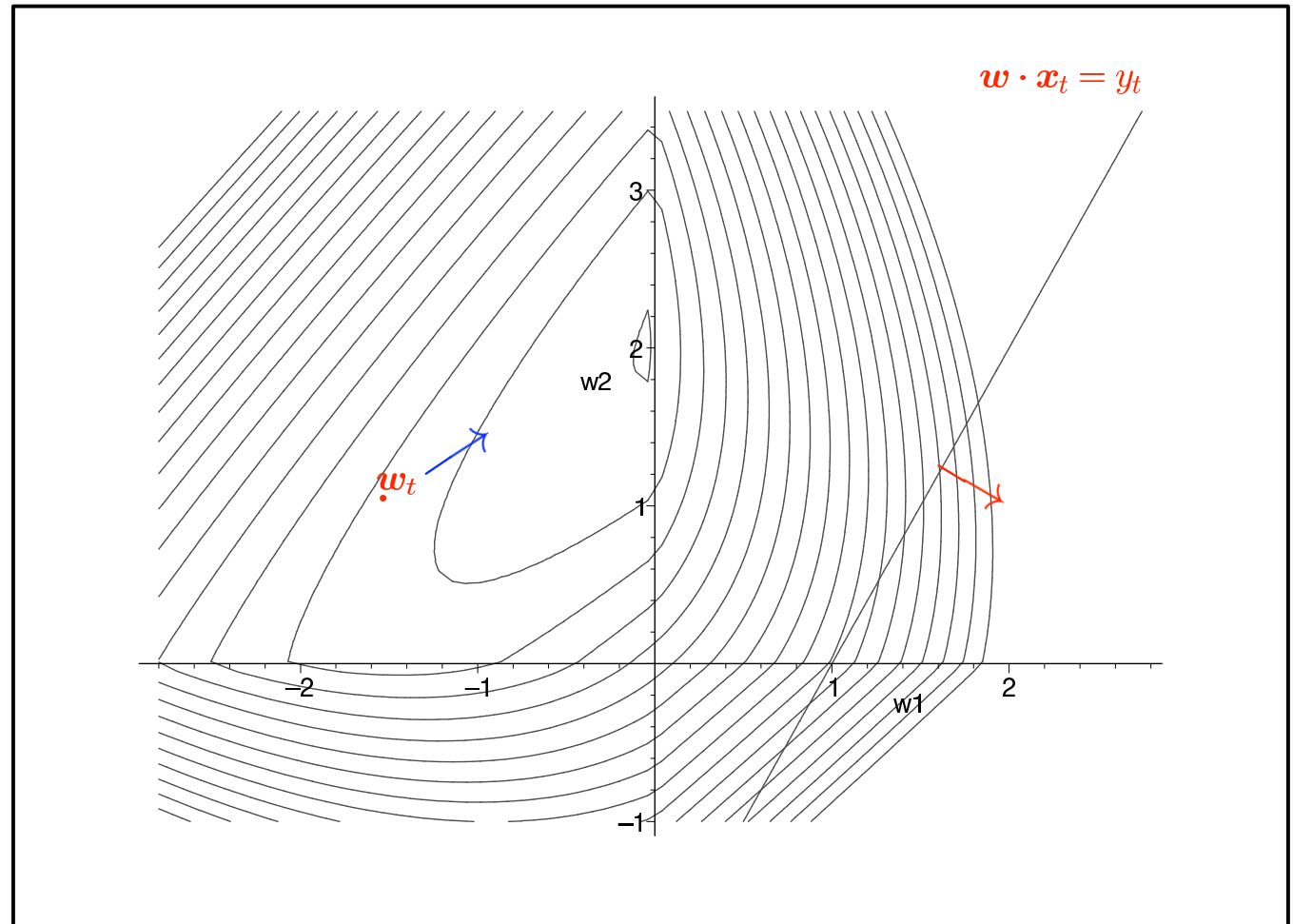
$$\Delta_F(\mathbf{w}, \mathbf{w}_t) + \eta \frac{1}{2} (y_t - \mathbf{w} \cdot \mathbf{x}_t)^2, \text{ where } F(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_{10}^2$$

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

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

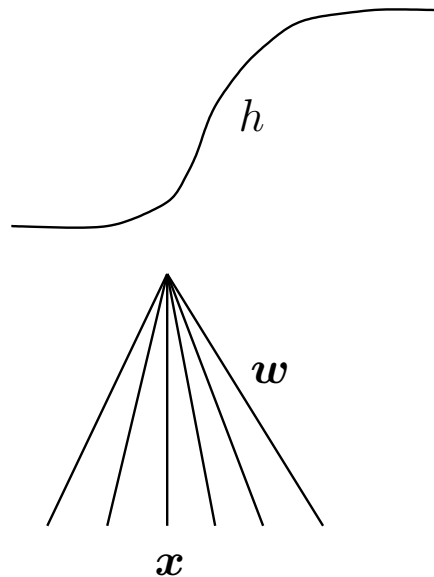
$$y_t = 1$$

$$\eta = 0.2$$



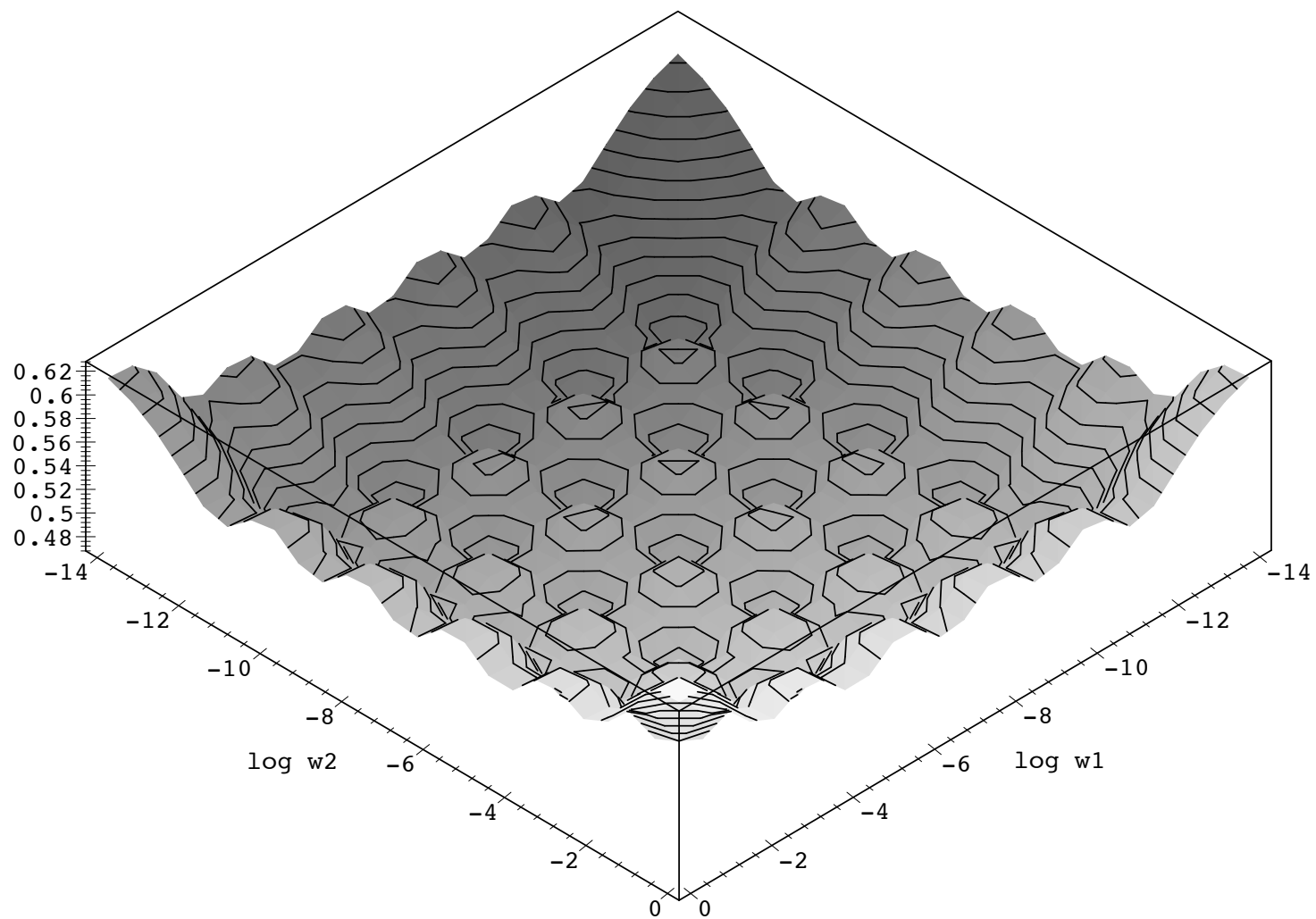
Nonlinear Regression

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



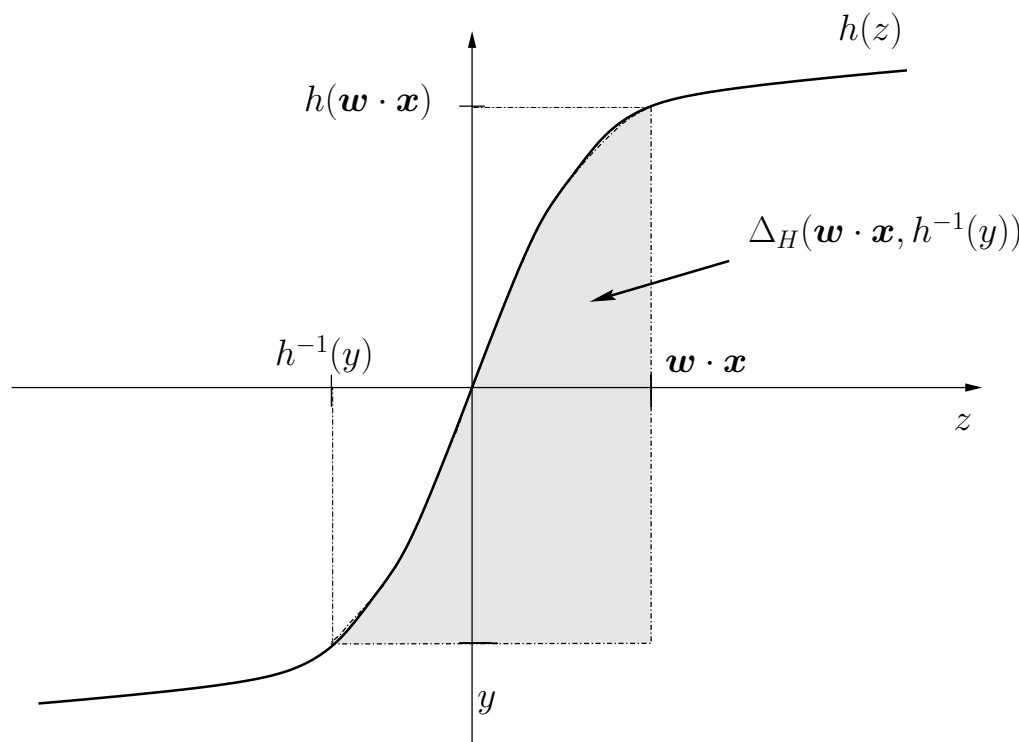
- Sigmoid function $h(z) = \frac{1}{1+e^{-z}}$
- For a set of examples $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$
total loss $\sum_{t=1}^T h(\mathbf{w} \cdot \mathbf{x}) - y_t)^2 / 2$
can have **exponentially many minima**
in weight space

[Bu,AHW]



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

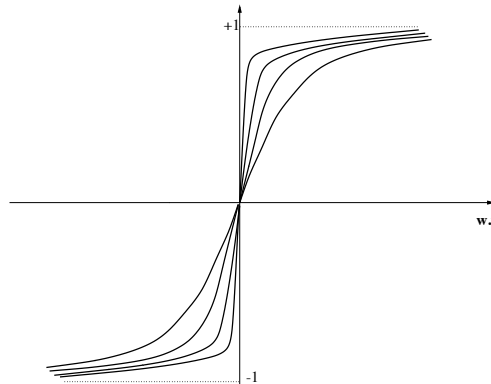
But not absolutely necessary

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

[Ce]

Sigmoid in the Limit

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

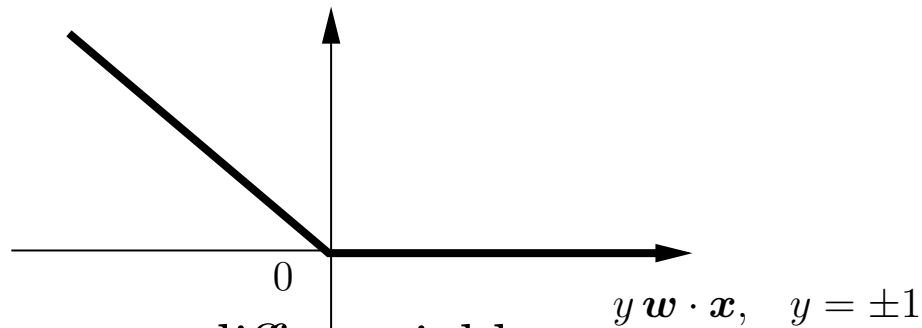


$$H(z) = |z|$$

Matching loss is **hinge loss**

[GW]

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



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

Perceptron

$$\mathbf{w}_{t+1}$$

$$= \underset{\mathbf{w}}{\operatorname{argmin}} \left(\|\mathbf{w} - \mathbf{w}_t\|^2 / 2 + \eta HL(\mathbf{w} \cdot \mathbf{x}_t, g^{-1}(y_t)) \right)$$

$$= \mathbf{w}_t - \eta (\operatorname{sign}(\mathbf{w}_{t+1} \cdot \mathbf{x}_t) - y_t) \mathbf{x}_t$$

$$\approx \mathbf{w}_t - \eta \underbrace{(\operatorname{sign}(\mathbf{w}_t \cdot \mathbf{x}_t))}_{\hat{y}_t} - y_t) \mathbf{x}_t$$

Normalized Winnow

$$\mathbf{w}_{t+1}$$

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

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

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

Trade-off between two divergences [KW]

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

Both divergences are convex in \mathbf{w}

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

Generalization of the “delta”-rule