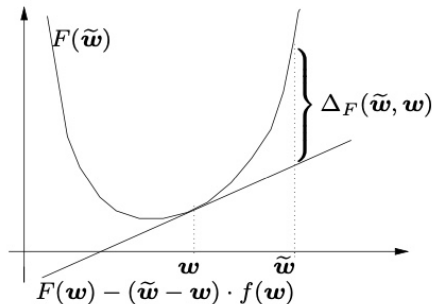


Bregman Divergences [Br,CL,Cs]

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

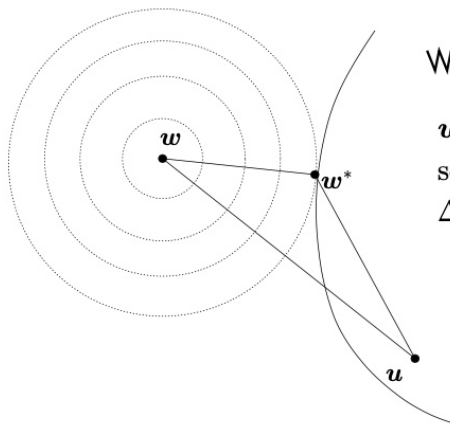
1. $\Delta_F(\tilde{\mathbf{w}}, \mathbf{w})$ is convex in $\tilde{\mathbf{w}}$
2. $\Delta_F(\tilde{\mathbf{w}}, \mathbf{w}) \geq 0$
If F convex equality holds iff $\tilde{\mathbf{w}} = \mathbf{w}$
3. Usually not symmetric: $\Delta_F(\tilde{\mathbf{w}}, \mathbf{w}) \neq \Delta_F(\mathbf{w}, \tilde{\mathbf{w}})$
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

$$\begin{aligned} 6. \quad \nabla_{\tilde{\mathbf{w}}} \Delta_F(\tilde{\mathbf{w}}, \mathbf{w}) &= \nabla F(\tilde{\mathbf{w}}) - \nabla_{\tilde{\mathbf{w}}}(\tilde{\mathbf{w}} \nabla_{\mathbf{w}} F(\mathbf{w})) \\ &= f(\tilde{\mathbf{w}}) - f(\mathbf{w}) \end{aligned}$$

$$\begin{aligned} 7. \quad \Delta_F(\mathbf{w}_1, \mathbf{w}_2) + \Delta_F(\mathbf{w}_2, \mathbf{w}_3) &= F(\mathbf{w}_1) - F(\mathbf{w}_2) - (\mathbf{w}_1 - \mathbf{w}_2)f(\mathbf{w}_2) \\ &\quad F(\mathbf{w}_2) - F(\mathbf{w}_3) - (\mathbf{w}_2 - \mathbf{w}_3)f(\mathbf{w}_3) \\ &= \Delta_F(\mathbf{w}_1, \mathbf{w}_3) + (\mathbf{w}_1 - \mathbf{w}_2) \cdot (f(\mathbf{w}_3) - f(\mathbf{w}_2)) \end{aligned}$$

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

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

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

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

(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-2 [GLS, GL]

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

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

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

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

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

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

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

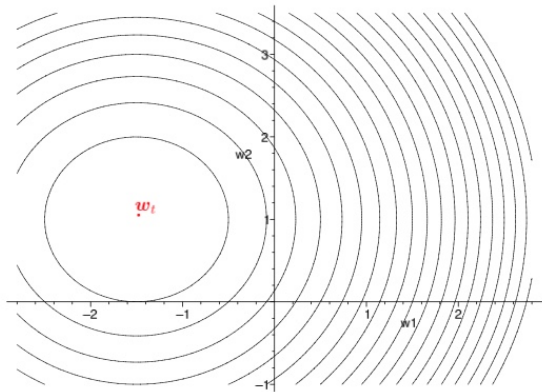
Divergence: Euclidean Distance Squared

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

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

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

$$y_t = 1$$



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]

How to prove relative loss bounds?

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

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

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

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

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]

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]

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]

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!