

# Reparameterizing Mirror Descent as Gradient Descent

Ehsan Amid & Manfred K. Warmuth  
Google Research, Brain Team  
{eamid, manfred}@google.com

# Mirror descent

$$\mathbf{w}_{s+1} = f^{-1}(f(\mathbf{w}_s) - h \nabla L(\mathbf{w}_s))$$

(where  $f$  is (coordinate-wise) strictly monotonic link function)

Gradient Descent (GD):

$$\mathbf{w}_{s+1} = \mathbf{w}_s - h \nabla L(\mathbf{w}_s) \quad (f(\mathbf{w}) = \mathbf{w})$$

Unnormalized Exponentiated Gradient Descent (EGU): [\[KW97\]](#)

$$\mathbf{w}_{s+1} = \mathbf{w}_s \odot \exp(-h(\nabla L(\mathbf{w}_s))) \quad (f(\mathbf{w}) = \log \mathbf{w})$$

(with  $w_i \geq 0$ )

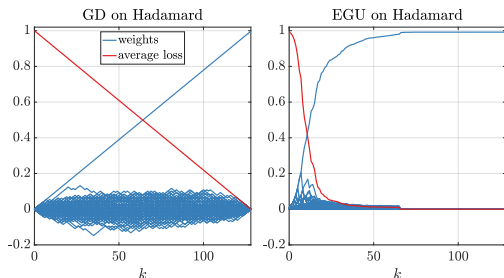
# Major differences between the two families

GD: backprop, kernel methods

EGU: Winnow, expert algorithms, Boosting, Bayes

Setup: **128x128 Hadamard matrix**

**Permuted** rows are instances, labels are any fixed column



x-axis:  $k = 1..128$

y-axis: **all 128 weights** **Loss when trained on examples 1..k**

Upshot: After half examples, GD has average loss =  $1/2$

EG family converges in  $\log(n)$  many examples

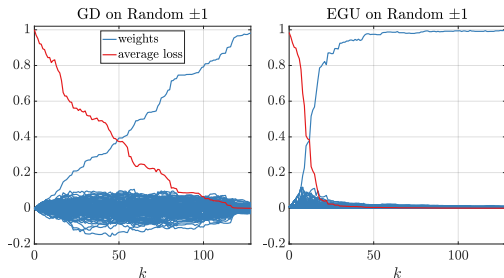
# Major differences between the two families

GD: backprop, kernel methods

EG: Winnow, expert algorithms, Boosting, Bayes

Setup: **128x128 random  $\pm 1$  matrix**

Rows are instances, labels are the first column, square loss



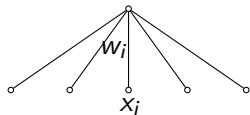
x-axis:  $k = 1..128$

y-axis: all 128 weights Loss when trained on examples  $1..k$

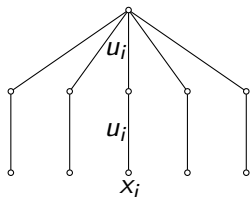
Upshot: After half examples, GD has average loss  $\approx 1/2$

EG family converges in  $\log(n)$  many examples

# Surprises



When linear neuron is trained with GD, then lower bound for linear decrease of avg. loss [WV05]



Reparameterize weights  $w_i$  by  $u_i^2$  [Akin79,GWBNS17]

Continuous GD on  $u_i$  simulates continuous EGU on  $w_i$

Discretizations learn Hadamard with Back-prop with essentially  $\mathcal{O}(\log n)$  examples

Experimentally indistinguishable from discrete EGU

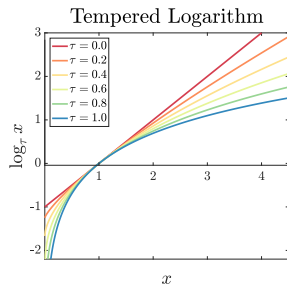
# Main focus here: continuous MD

$$\dot{f}(\mathbf{w}) = -\nabla L(\mathbf{w})$$

Main examples:

GD ( $f(\mathbf{w}) = \mathbf{w}$ ) and EGU ( $f(\mathbf{w}) = \log(\mathbf{w})$ )

Between  $f(\mathbf{w}) = \log \mathbf{w}$  and  $f(\mathbf{w}) = \mathbf{w}$ :  
 $\log_{\tau}(\mathbf{w}) := \frac{1}{1-\tau}(\mathbf{w}^{1-\tau} - 1)$   
(for  $\tau \in [0, 1]$ )



**Main Theorem:** For the reparameterization function  $\mathbf{w} = q(\mathbf{u})$  with the property that  $\text{range}(q) = \text{dom}(f)$ , the two updates

$$\dot{\mathbf{f}}(\mathbf{w}) = -\nabla L(\mathbf{w}) \quad \text{and} \quad \dot{\mathbf{u}} = -\nabla L \circ q(\mathbf{u}),$$

coincide if that  $\mathbf{w}(0) = q(\mathbf{u}(0))$ ,  $\text{range}(q) \subseteq \text{dom}(F)$ , and we have

$$(\mathbf{J}_f(\mathbf{w}))^{-1} = \mathbf{J}_q(\mathbf{u}) (\mathbf{J}_q(\mathbf{u}))^\top$$

# EGU as GD: The squaring trick

Link

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

Reparameterization

$$\mathbf{w} = q(\mathbf{u}) := 1/4 \mathbf{u} \odot \mathbf{u}$$

$$\mathbf{u} = 2\sqrt{\mathbf{w}}$$

$$(\mathbf{J}_f(\mathbf{w}))^{-1} = (\text{diag}(\mathbf{w})^{-1})^{-1} = \text{diag}(\mathbf{w})$$

$$\mathbf{J}_q(\mathbf{u})(\mathbf{J}_q(\mathbf{u}))^\top = 1/2 \text{diag}(\mathbf{u}) (1/2 \text{diag}(\mathbf{u}))^\top = \text{diag}(\mathbf{w})$$

Conclusion

$$\dot{\log}(\mathbf{w}) = -\nabla L(\mathbf{w}) \quad \text{equals} \quad \dot{\mathbf{u}} = -\underbrace{\nabla L \circ q}_{\nabla_{\mathbf{u}} L(1/4 \mathbf{u} \odot \mathbf{u})}(\mathbf{u})$$



Link

$$f(\mathbf{w}) = -\mathbf{1} \oslash \mathbf{w}$$

Reparameterization

$$\mathbf{w} = q(\mathbf{u}) := \exp(\mathbf{u})$$

$$\mathbf{u} = \log(\mathbf{w})$$

$$\begin{aligned} (\mathbf{J}_f(\mathbf{w}))^{-1} &= \text{diag}(\mathbf{1} \oslash (\mathbf{w} \odot \mathbf{w}))^{-1} = \text{diag}(\mathbf{w})^2 \\ \mathbf{J}_q(\mathbf{u})(\mathbf{J}_q(\mathbf{u}))^\top &= \text{diag}(\exp(\mathbf{u})) \text{diag}(\exp(\mathbf{u}))^\top = \text{diag}(\mathbf{w})^2 \end{aligned}$$

Conclusion

$$(-\mathbf{1} \oslash \mathbf{w}) = -\nabla L(\mathbf{w}) \quad \text{equals} \quad \dot{\mathbf{u}} = \underbrace{-\nabla L \circ q(\mathbf{u})}_{\nabla_{\mathbf{u}} L(\exp(\mathbf{u}))}$$

$$\log_{\tau} \mathbf{w} = \frac{1}{1-\tau} (\mathbf{w}^{1-\tau} - 1) \text{ as GD}$$

Link

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

Reparameterization

$$\mathbf{w} = q(\mathbf{u}) := \left( \frac{2-\tau}{2} \right)^{\frac{2}{2-\tau}} \mathbf{u}^{\frac{2}{2-\tau}}$$

$$\mathbf{u} = \frac{2}{2-\tau} \mathbf{w}^{\frac{2-\tau}{2}}$$

$$(\mathbf{J}_{\log_{\tau}}(\mathbf{w}))^{-1} = (\text{diag}(\mathbf{w})^{-\tau})^{-1} = \text{diag}(\mathbf{w})^{\tau}$$

$$\mathbf{J}_q(\mathbf{u})(\mathbf{J}_q(\mathbf{u}))^{\top} = \left( \left( \frac{2-\tau}{2} \right)^{\frac{\tau}{2-\tau}} \text{diag}(\mathbf{u})^{\frac{\tau}{2-\tau}} \right)^2 = \text{diag}(\mathbf{w})^{\tau}$$

Conclusion

$$\dot{\log}_{\tau}(\mathbf{w}) = -\nabla L(\mathbf{w}) \quad \text{equals} \quad \dot{\mathbf{u}} = - \underbrace{\nabla L \circ q(\mathbf{u})}_{\nabla_{\mathbf{u}} L \left( \left( \frac{2-\tau}{2} \right)^{\frac{2}{2-\tau}} \mathbf{u}^{\frac{2}{2-\tau}} \right)}$$

# Open problems

- ▶ World of continuous updates more succinct

$$\begin{aligned} \text{Euler discr.: } \quad & \frac{f(\mathbf{w}(t+h)) - f(\mathbf{w}(t))}{h} = -\nabla L(\mathbf{w}(t)) \\ & \iff \mathbf{w}(t+h) = f^{-1}(f(\mathbf{w}(t)) - h \nabla L(\mathbf{w}(t))) \end{aligned}$$

- ▶ Under what conditions does the discrete MD track continuous MD

Discretization of reparameterized EGU as GD tracks discrete EGU well enough so that the same regret bounds hold [AW20]

Discretization of reparameterized EGU as GD sample efficiently learns Hadamard problem