

# Winnowing with Gradient Descent

Ehsan Amid

**Manfred K. Warmuth**

Google Brain & UCSC

COLT 2020

# Two main families of updates

## **Additive updates:**

GD: stochastic gradient descent, backprop, Newton's update, kernel methods

## **Multiplicative updates:**

EG: expert algorithms, Boosting, Bayes

EGU: Winnow

Performance of GD linear in  $n$  for sparse targets

Performance of EG & EGU grows as  $\log n$  for sparse targets

Here we will reparameterize EG & EGU as GD:

Reparameterized forms act like original EG & EGU

Winnowing with GD!

# Two main families of updates

## **Additive updates:**

GD: stochastic gradient descent, backprop, Newton's update, kernel methods

## **Multiplicative updates:**

EG: expert algorithms, Boosting, Bayes

EGU: Winnow

Performance of GD linear in  $n$  for sparse targets

Performance of EG & EGU grows as  $\log n$  for sparse targets

Here we will reparameterize EG & EGU as GD:

Reparameterized forms act like original EG & EGU

**Winnowing with GD!**

# Paradigmatic sparse linear problem

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

Hadamard matrix or random  $\pm$

After receiving example  $(\mathbf{x}_t, y_t)$   
and incurring loss  $(\mathbf{x}_t^\top \mathbf{w}_t - y_t)^2$  update:

**multiplicative**, EGU:  $w_{t+1,i} = w_{t,i} \exp(-2\eta(\mathbf{x}_t^\top \mathbf{w}_t - y_t)x_{t,i})$

**additive**, GD:  $w_{t+1,i} = w_{t,i} - \underbrace{2\eta(\mathbf{x}_t^\top \mathbf{w}_t - y_t)x_{t,i}}_{\text{gradient}}$

Special cases of mirror descent (MD):

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

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

## Paradigmatic sparse linear problem

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

Hadamard matrix or random  $\pm$

After receiving example  $(\mathbf{x}_t, y_t)$   
and incurring loss  $(\mathbf{x}_t^\top \mathbf{w}_t - y_t)^2$  update:

**multiplicative**, EGU:  $w_{t+1,i} = w_{t,i} \exp(-2\eta(\mathbf{x}_t^\top \mathbf{w}_t - y_t)x_{t,i})$

**additive**, GD:  $w_{t+1,i} = w_{t,i} - \underbrace{2\eta(\mathbf{x}_t^\top \mathbf{w}_t - y_t)x_{t,i}}_{\text{gradient}}$

Special cases of mirror descent (MD):

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

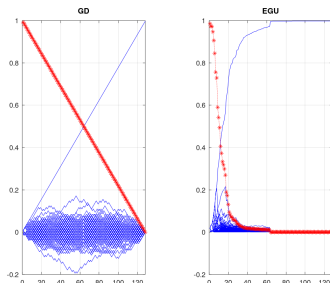
with  $f(\mathbf{w}) = \log \mathbf{w}$  or  $f(\mathbf{w}) = \mathbf{w}$

# Major differences between the two families

Paradigmatic setup:

**128x128 Hadamard matrix**

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



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

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

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

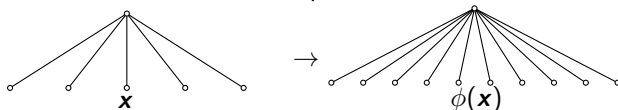
EG family converges in  $\log n$  many examples

# Hardness for GD Hadamard

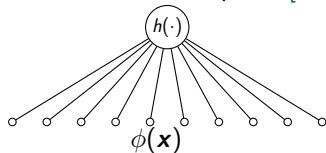
- ▶ Linear decay of loss remains for GD even if

- ▶ linear neuron with kernel inputs

[WV05]



- ▶ neuron with any transfer function  $h$  and kernel inputs [DW14]



Conjecture: Hadamard problem remains hard  
for any neural net trained with GD

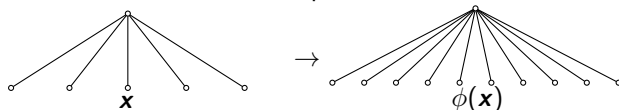
[DW14]

# Hardness for GD Hadamard

- ▶ Linear decay of loss remains for GD even if

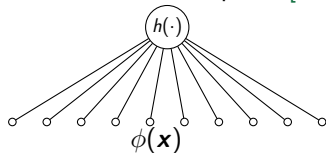
- ▶ linear neuron with kernel inputs

[WV05]



- ▶ neuron with any transfer function  $h$  and kernel inputs

[DW14]



Conjecture: Hadamard problem remains hard  
for any neural net trained with GD

[DW14]



- ▶ Parameter vector  $\mathbf{w}(t)$  continuous function of time
- ▶ Continuous update

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

- ▶ Examples are still discrete

$$(\mathbf{x}_s, y_s) \text{ for time } t \in [s, s + 1)$$

Again two main updates:

$$\text{GD} \quad \dot{\mathbf{w}}(t) = -\eta \nabla L(\mathbf{w}(t))$$

$$\text{EGU} \quad \dot{\log}(\mathbf{w}(t)) = -\eta \nabla L(\mathbf{w}(t))$$

We motivate updates in the continuous domain  
and then “discretize” these updates

# Two stunning surprises

- I) Continuous EGU can be simulated with continuous GD

Here: discretized versions of continuous GD simulation solves the Hadamard problem efficiently

Conjecture about GD training of neural nets is false  
Neural nets trained w. GD more powerful than kernel methods

- II) The structure of the network determines regularization when training with GD

# Two stunning surprises

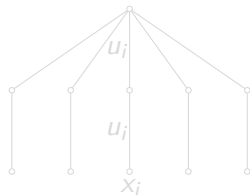
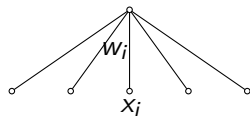
- I) Continuous EGU can be simulated with continuous GD

Here: discretized versions of continuous GD simulation solves the Hadamard problem efficiently

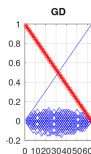
Conjecture about GD training of neural nets is false  
Neural nets trained w. GD more powerful than kernel methods

- II) The structure of the network determines regularization when training with GD

# I) Pictorially



When linear neuron is trained with GD,  
then linear decrease of loss



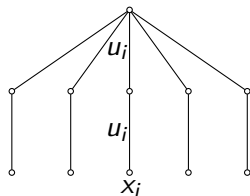
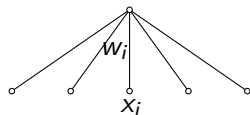
Reparameterize weights  $w_i$  by  $u_i^2$

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

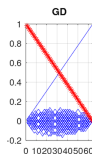
$$\dot{\mathbf{u}} = -2\eta (\mathbf{u} \odot \mathbf{u} \cdot \mathbf{x} - y) \mathbf{u} \odot \mathbf{x} \text{ simulates}$$

$$\dot{\mathbf{w}} = -2\eta (\mathbf{w} \cdot \mathbf{x} - y) \mathbf{x}$$

# I) Pictorially



When linear neuron is trained with GD,  
then linear decrease of loss



Reparameterize weights  $w_i$  by  $u_i^2$

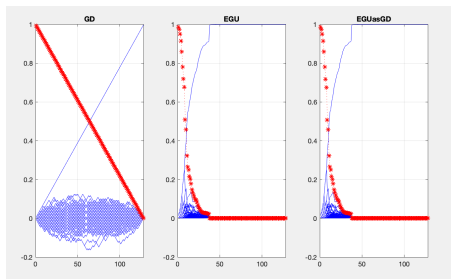
Continuous GD on  $u_i$  exactly  
simulates EGU on  $w_i$

$$\dot{\mathbf{u}} = -2\eta (\mathbf{u} \odot \mathbf{u} \cdot \mathbf{x} - y) \mathbf{u} \odot \mathbf{x} \text{ simulates}$$

$$\dot{\log(\mathbf{w})} = -2\eta (\mathbf{w} \cdot \mathbf{x} - y) \mathbf{x}$$

# I) Simulations

Discretization  $\mathbf{u}_{t+1} = \mathbf{u}_t - 2\eta (\mathbf{u}_t \odot \mathbf{u}_t \cdot \mathbf{x}_t - y_t) \mathbf{u}_t \odot \mathbf{x}_t$  tracks  
 $\mathbf{w}_{t+1} = \mathbf{w}_t \odot \exp(-2\eta (\mathbf{w}_t \cdot \mathbf{x}_t - y_t) \mathbf{x}_t)$



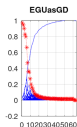
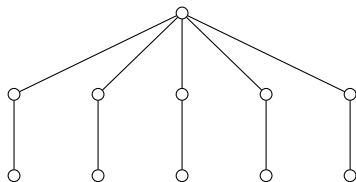
Simulation visually identical but slightly different numerically

Same regret bounds

Upshot: 2-layer neural net trained w. GD cracks Hadamard

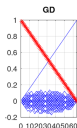
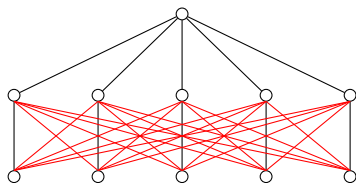
# Not just a matter of initialization

Case A



When trained with GD: approximates EGU and cracks Hadamard

Case B



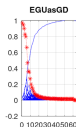
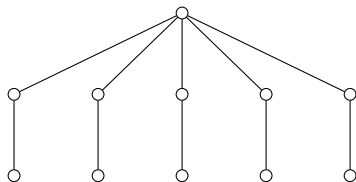
Red weights initialized to zero

Linear loss on Hadamard when trained with GD

Also true if all bottom weights initialized to zero

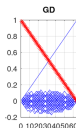
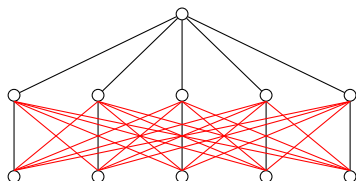
# Not just a matter of initialization

Case A



When trained with GD: approximates EGU and cracks Hadamard

Case B



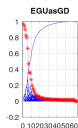
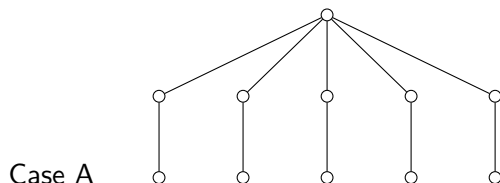
Red weights initialized to zero

Linear loss on Hadamard when trained with GD

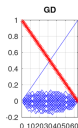
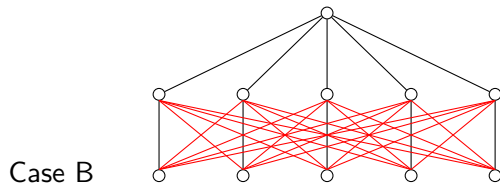
Also true if all bottom weights initialized to zero



## II) Structure determines regularization



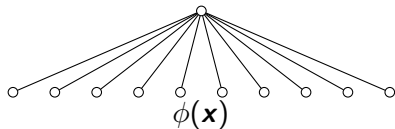
In continuous case, converges to **smallest  $L_1$  norm** solution  
In discrete case, same regret bounds as for EGU



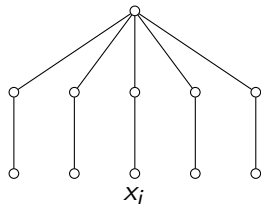
→ **smallest  $L_2$  norm** solution when bottom weights initialized to 0  
More complicated for other initializations, but experimentally satisfies linear lower bound

## 2-layer linear neural net GD can beat any kernel

For Hadamard problem



Any kernel has linear decaying loss on average



EGUasGD has exponentially decaying loss

# Two ways for obtaining discrete updates

1. As discretizations of continuous updates
2. Regularizing with Bregman divergences

For a strictly-convex function  $F(\mathbf{w})$ , the Bregman divergence is

$$\begin{aligned}\Delta_F(\mathbf{w}, \mathbf{w}_s) &= F(\mathbf{w}) - F(\mathbf{w}_s) - f(\mathbf{w}_s)^\top (\mathbf{w} - \mathbf{w}_s) \\ &= \Delta_{F^*}(\underbrace{f(\mathbf{w}_s)}_{\mathbf{w}_s^*}, \underbrace{f(\mathbf{w})}_{\mathbf{w}^*})\end{aligned}\quad (\text{duality})$$

$F(\mathbf{w})$  convex,  $\nabla F(\mathbf{w}) =: f(\mathbf{w}) = \mathbf{w}^*$  is the gradient  
 $f(\mathbf{w})$  called the link function

# Two ways for obtaining discrete updates

1. As discretizations of continuous updates
2. Regularizing with Bregman divergences

For a strictly-convex function  $F(\mathbf{w})$ , the Bregman divergence is

$$\begin{aligned}\Delta_F(\mathbf{w}, \mathbf{w}_s) &= F(\mathbf{w}) - F(\mathbf{w}_s) - f(\mathbf{w}_s)^\top (\mathbf{w} - \mathbf{w}_s) \\ &= \Delta_{F^*}(\underbrace{f(\mathbf{w}_s)}_{\mathbf{w}_s^*}, \underbrace{f(\mathbf{w})}_{\mathbf{w}^*})\end{aligned}\quad (\text{duality})$$

$F(\mathbf{w})$  convex,  $\nabla F(\mathbf{w}) =: f(\mathbf{w}) = \mathbf{w}^*$  is the gradient  
 $f(\mathbf{w})$  called the link function

$$\mathbf{w}_{s+1} = \operatorname{argmin}_{\tilde{\mathbf{w}}} \Delta_F(\tilde{\mathbf{w}}, \mathbf{w}_s) + \eta L(\tilde{\mathbf{w}})$$

Setting derivative at  $\mathbf{w}_{s+1}$  to zero

$$f(\mathbf{w}_{s+1}) - f(\mathbf{w}_s) + \eta \nabla L(\mathbf{w}_{s+1}) = \mathbf{0}$$

Implicit/Prox MD update

[R76,NY83]

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

Explicit MD update

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

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

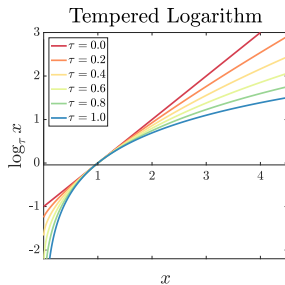
Main examples:

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

$$\log_{\tau}(\mathbf{w}) := \frac{1}{1-\tau}(\mathbf{w}^{1-\tau} - 1)$$

$\tau$  is temperature

(we use  $\tau \in [0, 1]$ )



[N02]

# Motivation with Bregman momentum

$$\mathbf{w}(t) = \operatorname{argmin}_{\tilde{\mathbf{w}}(t)} \underbrace{\dot{\Delta}_F(\tilde{\mathbf{w}}(t), \mathbf{w}_s)}_{\text{Bregman momentum}} + \eta L(\tilde{\mathbf{w}}(t))$$

Derivation of the optimum curve  $\mathbf{w}(t)$ :

$$\begin{aligned} & \frac{\partial}{\partial \tilde{\mathbf{w}}(t)} \left( \frac{\partial}{\partial t} \left( F(\tilde{\mathbf{w}}(t)) - f(\mathbf{w}_s)^\top \tilde{\mathbf{w}}(t) \right) + \eta L(\tilde{\mathbf{w}}(t)) \right) \quad (\text{differentiate}) \\ &= \frac{\partial}{\partial \tilde{\mathbf{w}}(t)} \left( (f(\tilde{\mathbf{w}}(t)) - f(\mathbf{w}_s))^\top \dot{\tilde{\mathbf{w}}}(t) + \eta \nabla L(\tilde{\mathbf{w}}(t)) \right) \\ &= (\mathbf{J}f(\tilde{\mathbf{w}}) \dot{\tilde{\mathbf{w}}}(t) + \underbrace{\left( \frac{\partial \dot{\tilde{\mathbf{w}}}(t)}{\partial \tilde{\mathbf{w}}(t)} \right)^\top}_{\mathbf{0}} (f(\tilde{\mathbf{w}}(t)) - f(\mathbf{w}_s)) + \eta \nabla L(\tilde{\mathbf{w}}(t))) \end{aligned}$$

(By calculus of variations,  $\tilde{\mathbf{w}}(t)$  and  $\dot{\tilde{\mathbf{w}}}(t)$  are independent variables)

$$= \dot{f}(\tilde{\mathbf{w}}(t)) + \eta \nabla L(\tilde{\mathbf{w}}(t)) \stackrel{\tilde{\mathbf{w}}(t) = \mathbf{w}(t)}{=} \mathbf{0}$$

## Reparameterizing cont. MD w. link $f$ i.t.o. link $g$

**Theorem** For the reparameterization function  $\mathbf{w} = q(\mathbf{u})$  with the property that  $\text{range}(q) = \text{dom}(f)$ ,  $\dot{\mathbf{g}}(\mathbf{u}) = -\eta \nabla L \circ q(\mathbf{u})$  simulates  $\dot{\mathbf{f}}(\mathbf{w}) = -\eta \nabla L(\mathbf{w})$  if

$$(\mathbf{J}f(\mathbf{w}))^{-1} = \mathbf{J}q(\mathbf{u}) (\mathbf{J}g(\mathbf{u}))^{-1} (\mathbf{J}q(\mathbf{u}))^\top$$

and  $q(\mathbf{u}(0)) = \mathbf{w}(0)$

For reparameterization as GD use  $g = id$

[details in a companion paper under review]



# Our main example: EGU as GD

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

$$\dot{\log}(\mathbf{w}) = -\eta \nabla L(\mathbf{w}) \text{ equals } \dot{\mathbf{u}} = -\eta \underbrace{\nabla L \circ q(\mathbf{u})}_{\nabla_{\mathbf{u}} L(1/4 \mathbf{u} \odot \mathbf{u})} = -\eta 1/2 \mathbf{u} \odot \nabla L(\mathbf{w})$$

Link

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

Reparameterization

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

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

$$(\mathbf{J}f(\mathbf{w}))^{-1} = \text{diag}\left(\frac{1}{\mathbf{w} \odot \mathbf{w}}\right)^{-1} = \text{diag}(\mathbf{w})^2$$

$$\mathbf{J}q(\mathbf{u})(\mathbf{J}q(\mathbf{u}))^\top = \text{diag}(\exp(\mathbf{u})) \text{diag}(\exp(\mathbf{u}))^\top = \text{diag}(\mathbf{w})^2$$

Conclusion

$$\left(-\frac{\dot{1}}{\mathbf{w}}\right) = -\eta \nabla L(\mathbf{w}) \text{ equals } \dot{\mathbf{u}} = -\eta \underbrace{\nabla L \circ q(\mathbf{u})}_{\nabla_{\mathbf{u}} L(\exp(\mathbf{u}))} = -\eta \exp(\mathbf{u}) \odot \nabla L(\mathbf{w})$$

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

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

Reparameterization

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

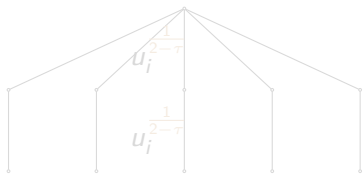
$$\mathbf{u} = \frac{2}{2-\tau} \mathbf{w}^{\frac{2-\tau}{2}}$$

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

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

Conclusion

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



$\tau = 1$ : EGU     $\tau = 0$ : GD

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

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

Reparameterization

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

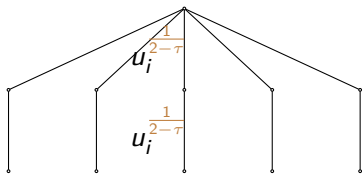
$$\mathbf{u} = \frac{2}{2-\tau} \mathbf{w}^{\frac{2-\tau}{2}}$$

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

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

Conclusion

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



$$\tau = 1: \text{ EGU} \quad \tau = 0: \text{ GD}$$

# Discrete multiplicative updates for dot loss $\sum_i w_i \ell_i$

EGU	$\tilde{w}_i = w_i \exp(-\eta \ell_i)$
Approx. EGU/PRODU	$\tilde{w}_i = w_i(1 - \eta \ell_i)$
EGUasGD	$\tilde{u}_i = u_i(1 - \eta \ell_i)$ $(\tilde{u}_i^2 = u_i^2(1 - \eta \ell_i)^2)$
EG/HEDGE	$\tilde{w}_i = \frac{w_i \exp(-\eta \ell_i)}{\sum_j w_j \exp(-\eta \ell_j)}$
Approx. EG	$\tilde{w}_i = w_i(1 - \eta \ell_i + \eta \sum_j w_j \ell_j)$
PROD	$\tilde{w}_i = \frac{w_i(1 - \eta \ell_i)}{\sum_j w_j(1 - \eta \ell_j)}$
EGasGD	$\tilde{u}_i = \frac{u_i(1 - \eta \ell_i)}{\ \sum_j u_j^2(1 - \eta \ell_j)^2\ _2^2}$ $(\tilde{u}_i^2 = \frac{u_i^2(1 - \eta \ell_i)^2}{\sum_j u_j^2(1 - \eta \ell_j)^2})$

# Regret bounds

total online loss of update

$$\leq \text{total online loss of best comparator} + \text{norms} \sqrt{\text{loss of best}}$$

update	regret bound
EGUasGD, hinge loss	as Winnow
EGUasGD, linear regression	as EGU but only one-sided case
EGasGD, linear regression	as EG
EGasGD, dot loss	as Hedge

All proofs done with relative entropy as a measure of progress

# Open problems

- ▶ Need 2-sided regret bound for linear regression EGU and EGUasGD  
Or prove linear lower bound when Hadamard is replaced by 0/1 matrix
- ▶ **Is there any natural problem in which GD beats EGU $^{\pm}$ ?**
- ▶ Revisit vanishing gradient issue, batch normalization, dropout, learning rate heuristics for neural nets where all linear activations are replaced by the sparse network
- ▶ Large scale simulations
  - Do multiplicative updates lead to sparse solutions?
- ▶ **Revised open question about limited power of GD:  
Does any GD trained neural net with complete input neurons satisfy the linear lower bound for the Hadamard problem?**

Thank you!

# Open problems

- ▶ Need 2-sided regret bound for linear regression EGU and EGUasGD  
Or prove linear lower bound when Hadamard is replaced by 0/1 matrix
- ▶ **Is there any natural problem in which GD beats EGU $^{\pm}$ ?**
- ▶ Revisit vanishing gradient issue, batch normalization, dropout, learning rate heuristics for neural nets where all linear activations are replaced by the sparse network
- ▶ Large scale simulations
  - Do multiplicative updates lead to sparse solutions?
- ▶ **Revised open question about limited power of GD:  
Does any GD trained neural net with complete input neurons satisfy the linear lower bound for the Hadamard problem?**

Thank you!