

# Winnowing with Gradient Descent

**Ehsan Amid**

*Google Brain, Mountain View, CA 94043*

EAMID@GOOGLE.COM

**Manfred K. Warmuth**

*Google Brain, Mountain View, CA 94043*

MANFRED@GOOGLE.COM

**Editors:** Jacob Abernethy and Shivani Agarwal

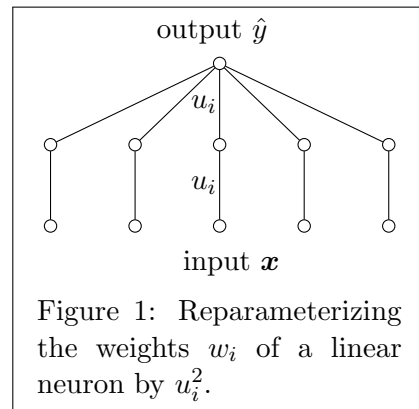
## Abstract

The performance of multiplicative updates is typically logarithmic in the number of features when the targets are sparse. Strikingly, we show that the same property can also be achieved with gradient descent updates. We obtain this result by rewriting the non-negative weights  $w_i$  of multiplicative updates by  $u_i^2$  and then performing a gradient descent step w.r.t. the new  $u_i$  parameters. We apply this method to the Winnow update, the Hedge update, and the unnormalized and normalized exponentiated gradient (EG) updates for linear regression. When the original weights  $w_i$  are scaled to sum to one (as done for Hedge and normalized EG), then in the corresponding reparameterized update, the  $u_i$  parameters are now divided by  $\|\mathbf{u}\|_2$  after the gradient descent step. We show that these reparameterizations closely track the original multiplicative updates by proving in each case the same online regret bounds (albeit in some cases, with slightly different constants).

As a side, our work exhibits a simple two-layer linear neural network that, when trained with gradient descent, can experimentally solve a certain sparse linear problem (known as the Hadamard problem) with exponentially fewer examples than any kernel method.

## 1. Introduction

Multiplicative updates started with the Winnow algorithm for learning disjunctions and linear threshold functions (Littlestone, 1988). These algorithms update their weights by multiplicative factors and often learn sparse targets with a logarithmic dependence on the number of features. For example, Winnow makes at most  $k \log n$  mistakes on the sequence of  $n$ -dimensional boolean vectors labeled consistently with a  $k$ -literal disjunction. Another paradigmatic case is the so called “expert setting” where the learner is to perform as well as a single feature/expert. For sequences of  $n$ -dimensional feature vectors with one consistent feature, the total loss of the update is typically  $\mathcal{O}(\log n)$ . The situation is repeated for linear regression; when the weight vector realizing the labels is sparse, then the normalized and unnormalized EG algorithms (Kivinen and Warmuth, 1997) incur loss at most  $\mathcal{O}(\log n)$ . So far we focused on the noise-free case. In the noisy case (when the solution has non-zero loss), additional square root terms appear in the regret bounds of these algorithms.



We first observe that in continuous-time, multiplicative updates can be exactly reparameterized as continuous-time gradient descent (GD). This is done by replacing the linear weights  $w_i$  by  $u_i^2$  and applying continuous-time GD w.r.t. the new parameter  $u_i$  (See Figure 1). We show how the discretization of the continuous multiplicative updates results in the usual Winnow and EG updates. More importantly, the discretization of the reparameterized continuous multiplicative updates results in discrete-time GD variants of these updates. However, the equivalence between the multiplicative updates and the GD reparameterizations does not hold after discretization. Nevertheless, we show that the GD reparameterizations closely track the multiplicative originals by proving the same regret bounds for the new updates, as were known for the originals. In each case, the constants we obtain for the second order terms are the same or slightly different in either direction.

For the discrete GD updates, lower bounds that grow linearly in the number of features  $n$  are known in some settings (Kivinen et al., 1997; Vishwanathan and Warmuth, 2005). In these lower bounds, the instances are essentially the rows of an  $n \times n$  Hadamard matrix and the target, in the simplest case, is a single column of this matrix. When averaged over targets, the lower bound for linear regression holds even if the input vectors  $\mathbf{x}$  are replaced by a feature map  $\phi(\mathbf{x})$  of any kernel (Vishwanathan and Warmuth, 2005). It was conjectured in (Derezinski and Warmuth, 2014) that the linear lower bound holds for any neural network, as long as it is trained with GD. This conjecture is now contradicted by the linear network in Figure 1. If this network is trained with online GD (i.e. backpropagation) w.r.t. the squared Euclidean loss, then the update coincides with our discrete-time GD reparameterization of the EGU update. Experimentally, GD on this network solves the Hadamard problem as efficiently as EG and EGU (on a single neuron). Interestingly, this simple two-layer linear network, when trained with GD, can realize the  $\log n$  dependence on the number of examples. Also surprisingly, when all the missing  $n^2 - n$  connections in the bottom layer of Figure 1 are reinstated and initialized to zero, then experimentally the loss of the now fully-connected two-layer linear network again decays only linearly when trained with GD.

**Previous work.** Online regret bounds for multiplicative updates have a long history (see e.g. (Littlestone, 1988; Vovk, 1990; Kivinen and Warmuth, 1997)). A key insight of Gunasekar et al. (2017) (followed by Woodworth et al. (2019) and Vaskevicius et al. (2019)) was that if the weights  $w_i$  of a linear predictor are replaced by  $u_i^2$ , then the continuous-time gradient descent update on the  $u_i$ 's imposes an implicit regularization on the  $w_i$ 's that makes the reparameterized weight vector  $\mathbf{w}$  converge to the minimum  $L_1$ -norm solution.<sup>1</sup> The EG updates are known to connect with the  $L_1$ -norm regularization. Indeed, the relative entropy (the divergence motivating the EG updates) is strongly convex w.r.t. the  $L_1$ -norm (Shalev-Shwartz et al., 2012). We show that the continuous-time unnormalized EG algorithm is in fact equivalent to the continuous-time GD on the  $u_i$ 's. Curiously enough, this reparameterization method has been used by game theorists to convert problems defined on the unit simplex to the unit sphere (Akin, 1979; Sandholm, 2010). In particular, the replicator dynamics of evolutionary game theory corresponds to EG, and GD equivalents have been investigated for this dynamic. In the continuous case, it is easy to see the

1. Note that Gunasekar et al. (2017) analyzed the much richer matrix context. For the sake of simplicity, we focus on the diagonal (i.e. vector) case.

equivalence between the EG updates and the corresponding reparameterizations as GD. However, the main contribution of this paper is showing that essentially the same regret bounds hold for the discretized GD reparameterizations of the multiplicative updates.

Note that obtaining an online regret bound is one of the most stringent learning criterion since it must hold for worst-case sequences of examples. Standard conversions to randomized settings also exist (see e.g. (Kivinen and Warmuth, 1997)). Furthermore, EGU is a special case of mirror descent where the link function is the component-wise logarithm. A more general framework for reparameterizing mirror descent updates was developed in a recent paper (Amid and Warmuth, 2020).

**Outline.** In the next section, we show the equivalence between the continuous-time exponentiated gradient updates and their reparameterized gradient descent versions. We also motivate the discretizations of the continuous updates. In the following three chapters, we reprove the regret bounds for the discrete gradient descent variants of Winnow, the Hedge algorithm, and the EG and EGU algorithms for linear regression. Lower bounds and some implications of neural network training are briefly discussed in Section 6. We conclude the paper with a number of open problems in Section 7.

**Notation.** We use  $\odot$  and  $\oslash$  to denote element-wise vector multiplication and division, respectively. The dot symbol denotes the time derivative, i.e.  $\dot{\mathbf{v}}(t) := \frac{\partial}{\partial t} \mathbf{v}(t)$ .

## 2. Reparameterizing the Continuous-time Exponentiated Gradient

In the following, we discuss the continuous-time EGU and EG updates<sup>2</sup> and derive their reparameterized forms. Although the reparameterization method presented here is more versatile and applies to a wider class of continuous-time mirror descent updates (Amid and Warmuth, 2020), we only focus on these two main cases in this paper.

The continuous-time EGU update can be seen as the solution of the following ordinary differential equation (ODE) defined on  $\mathbb{R}_{\geq 0}^n$ :

$$\dot{\log \mathbf{w}}(t) = -\eta \nabla_{\mathbf{w}} L(\mathbf{w}(t)), \text{ with initial condition } \mathbf{w}(t=0) = \mathbf{w}^0 \in \mathbb{R}_{\geq 0}^n. \quad (1)$$

Here,  $\nabla_{\mathbf{w}} L(\mathbf{w}(t))$  denotes the gradient of the loss  $L(\mathbf{w})$  w.r.t.  $\mathbf{w}$  evaluated at  $\mathbf{w}(t)$  and  $\eta$  is a positive learning rate. Typically, the loss  $L(\cdot)$  also depends on a given example  $\mathbf{x}(t)$  and possibly a label  $y(t)$ . Although the dynamic of the continuous-time EGU (1) is fundamentally different than a dynamic based on gradient descent, the following theorem shows that the same update can be realized using gradient descent via a simple reparameterization. Namely, we substitute  $\mathbf{w}(t) = q(\mathbf{u}(t)) := \mathbf{u}(t) \odot \mathbf{u}(t)$  where the new parameter  $\mathbf{u}(t) \in \mathbb{R}^n$  is updated via gradient descent on the composite loss  $L \circ q(\cdot)$ .<sup>3</sup>

**Theorem 1 (Reparameterized continuous-time EGU)** *The solution of the ODE (1) is equal to  $\mathbf{w}(t) = q(\mathbf{u}(t))$ , for all  $t \geq 0$ , where  $q(\mathbf{u}) := \mathbf{u} \odot \mathbf{u}$  and the new parameter  $\mathbf{u}(t)$  is the solution of the following ODE defined on  $\mathbb{R}^n$ :*

$$\dot{\mathbf{u}}(t) = -\eta/4 \nabla_{\mathbf{u}} L \circ q(\mathbf{u}(t)), \text{ s.t. } \mathbf{u}(t=0) = \mathbf{u}^0 \in \mathbb{R}^n \text{ and } \mathbf{u}^0 \odot \mathbf{u}^0 = \mathbf{w}^0. \quad (2)$$

- 
2. Recall that for any loss  $L(\mathbf{w})$ , where  $\mathbf{w} \in \mathbb{R}_{\geq 0}^n$ , the unnormalized exponentiated gradient update (EGU) resets  $\mathbf{w}$  to  $\mathbf{w} \odot \exp(-\eta \nabla L(\mathbf{w})) \in \mathbb{R}_{\geq 0}^n$ . The normalized version (EG) divides by the sum of the updated weights, thus projecting the weight vector onto the unit simplex  $\Delta^{n-1}$ .
  3. Even if  $L(\mathbf{w})$  is convex in  $\mathbf{w}$ , the composite loss  $L \circ q(\mathbf{u})$  might not be convex in the new parameter  $\mathbf{u}$ .

**Proof** By the chain rule,  $\dot{\log} \mathbf{w}(t) = \dot{\mathbf{w}}(t) \odot \mathbf{w}(t)$ . Thus the ODE (1) can be written as

$$\dot{\mathbf{w}}(t) = -\eta \mathbf{w}(t) \odot \nabla_{\mathbf{w}} L(\mathbf{w}(t)). \quad (3)$$

Similarly, 
$$\dot{\mathbf{w}}(t) = \frac{\partial}{\partial t} (\mathbf{u}(t) \odot \mathbf{u}(t)) = 2 \mathbf{u}(t) \odot \dot{\mathbf{u}}(t). \quad (4)$$

Plugging (2) into (4) and using  $\nabla_{\mathbf{u}} L \circ q(\mathbf{u}(t)) = 2 \mathbf{u}(t) \odot \nabla_{\mathbf{w}} L(\mathbf{w}(t))$  and  $\mathbf{w}(t) = q(\mathbf{u}(t))$  yields (3) and concludes the proof.  $\blacksquare$

The discrete-time EGU can be derived as the finite difference approximation of (1) with a step-size of one, that is,  $\log \mathbf{w}^{t+1} - \log \mathbf{w}^t = -\eta \nabla_{\mathbf{w}} L(\mathbf{w}^t)$ . Alternatively, this is also the solution to the following minimization problem:

$$\mathbf{w}^{t+1} = \operatorname{argmin}_{\tilde{\mathbf{w}} \in \mathbb{R}_{\geq 0}^n} \left\{ 1/\eta D_{\text{RE}}(\tilde{\mathbf{w}}, \mathbf{w}^t) + \hat{L}(\tilde{\mathbf{w}} | \mathbf{w}^t) \right\}, \quad (5)$$

where  $D_{\text{RE}}(\tilde{\mathbf{w}}, \mathbf{w}) = \sum_{i=1}^n \tilde{w}_i \log \frac{\tilde{w}_i}{w_i} - \tilde{w}_i + w_i$  denotes the relative entropy divergence and  $\hat{L}(\tilde{\mathbf{w}} | \mathbf{w}^t) = L(\tilde{\mathbf{w}}) + (\tilde{\mathbf{w}} - \mathbf{w}^t) \cdot \nabla_{\mathbf{w}} L(\mathbf{w}^t)$  is the Taylor approximation<sup>4</sup> of the loss  $L(\tilde{\mathbf{w}})$  at  $\mathbf{w}^t$ . Similarly, discretizing (2) yields the discrete-time reparameterized EGU update

$$\mathbf{u}^{t+1} - \mathbf{u}^t = -\eta/4 \nabla_{\mathbf{u}} L \circ q(\mathbf{u}^t) = -\eta/2 \mathbf{u}^t \odot \nabla_{\mathbf{w}} L(\mathbf{w}^t), \quad (6)$$

with  $\mathbf{w}^t = \mathbf{u}^t \odot \mathbf{u}^t$ . This can be seen as the solution to the minimization problem

$$\mathbf{u}^{t+1} = \operatorname{argmin}_{\tilde{\mathbf{u}} \in \mathbb{R}^n} \left\{ 2/\eta \|\tilde{\mathbf{u}} - \mathbf{u}^t\|_2^2 + \widehat{L \circ q}(\tilde{\mathbf{u}} | \mathbf{u}^t) \right\}, \quad (7)$$

where  $\widehat{L \circ q}(\tilde{\mathbf{u}} | \mathbf{u}^t)$  is the first-order Taylor series approximation of  $L \circ q(\tilde{\mathbf{u}})$  at  $\mathbf{u}^t$ :

$$\widehat{L \circ q}(\tilde{\mathbf{u}} | \mathbf{u}^t) = L \circ q(\mathbf{u}^t) + (\tilde{\mathbf{u}} - \mathbf{u}^t) \cdot \nabla_{\mathbf{u}} L \circ q(\mathbf{u}^t). \quad (8)$$

A similar continuous-time dynamic can be constructed for the normalized EG update and its equivalent reparameterized form.<sup>5</sup> This involves applying projected gradient updates to maintain the constraint  $\|\mathbf{w}(t)\|_1 = 1$  (respectively,  $\|\mathbf{u}(t)\|_2^2 = 1$ ). Here, we state the result for the discretized reparameterized form and refer the reader to (Amid and Warmuth, 2020) for further details. Adding the Lagrange multiplier  $\lambda$  to (7) to enforce the constraint  $(\mathbf{u}^{t+1} \odot \mathbf{u}^{t+1}) \cdot \mathbf{1} = 1$ , we have

$$\mathbf{u}^{t+1} = \operatorname{argmin}_{\tilde{\mathbf{u}} \in \mathbb{R}^n} \left\{ 1/\eta \|\tilde{\mathbf{u}} - \mathbf{u}^t\|_2^2 + \widehat{L \circ q}(\tilde{\mathbf{u}} | \mathbf{u}^t) + \lambda \left( (\tilde{\mathbf{u}} \odot \tilde{\mathbf{u}}) \cdot \mathbf{1} - 1 \right) \right\}, \quad (9)$$

which results in the reparameterized EG update

$$\mathbf{u}^{t+1} = \frac{\mathbf{u}^t - \eta/2 \nabla_{\mathbf{u}} L \circ q(\mathbf{u}^t)}{\|\mathbf{u}^t - \eta/2 \nabla_{\mathbf{u}} L \circ q(\mathbf{u}^t)\|_2}, \quad \text{thus } \mathbf{u}^{t+1} \odot \mathbf{u}^{t+1} \in \Delta^{n-1}. \quad (10)$$

4. Note that a backward Euler approximation of (1) yields the *prox* or *implicit form* of the update (Rockafellar, 1976; Kivinen et al., 2006), which corresponds to using  $L(\tilde{\mathbf{w}})$  instead of  $\hat{L}(\tilde{\mathbf{w}} | \mathbf{w}^t)$  in (5).

5. In the remainder of the paper, we use  $\mathbf{w} = q(\mathbf{u}) = \mathbf{u} \odot \mathbf{u}$  as the reparameterization. The constants can be absorbed into the learning rates.

Clearly, after the discretization, the EGU and EG updates and their reparameterized forms are no longer equivalent. The key question that naturally arises is the following: *How well does the discretized reparameterized EGU update (and its normalized form) approximate the original EGU (respectively, the original EG) update?* In the following, we address this question by comparing the worst-case regrets of the original (un)normalized EG and its reparameterizations on three problems, namely, the Winnow algorithm for binary classification, the Hedge algorithm for the expert setting, and EGU and EG for linear regression. The following regret bounds for the reparameterizations are proven by bounding the progress towards a comparator. Curiously enough, the progress is measured i.t.o. the relative entropy divergence even though the reparameterized updates are motivated by regularizing with the squared Euclidean distance.

### 3. The Reparameterization of the Winnow algorithm

The Winnow algorithm learns a linear threshold function for the task of binary classification. It is a special case of EGU when the loss is the hinge loss. This loss on example  $(\mathbf{x}^t, y^t)$  (for  $y^t \in \{\pm 1\}$ ) is defined as  $L_H(\mathbf{w} | \mathbf{x}^t, y^t) = [-y^t(\mathbf{w} \cdot \mathbf{x}^t - \theta)]_+$ , where  $\theta \in \mathbb{R}$  is a threshold. GD w.r.t. the same loss results in the Perceptron update. (See (Kivinen et al., 1997) for a comparative study.) Below we give the Winnow algorithm and its reparameterization as GD. Both algorithms update their weights only when a mistake occurs because when the prediction is correct, the gradient of the hinge loss is zero (Gentile and Warmuth, 1998).

---

#### Algorithm 1 Winnow Algorithm

---

**Parameters** initial weight  $w^1 > 0$ , learning rate  $\eta > 0$ , threshold  $\theta > 0$   
**Initialize**  $\mathbf{w}^1 = w^1 \mathbf{1}_n$   
**for**  $t = 1$  to  $T$  **do**  
  Receive instance  $\mathbf{x}^t \in [0, 1]^n$   
  Predict  $\hat{y}^t = \begin{cases} +1 & \text{if } \mathbf{w}^t \cdot \mathbf{x}^t \geq \theta \\ -1 & \text{otherwise} \end{cases}$   
  Receive label  $y^t$  and update:  
 $\mathbf{w}^{t+1} = \begin{cases} \mathbf{w}^t & \text{if } \hat{y}^t = y^t \\ \mathbf{w}^t \odot \exp(\eta y^t \mathbf{x}^t) & \text{otherwise} \end{cases}$

---



---

#### Algorithm 2 Reparameterized Winnow

---

**Parameters** initial weight  $u^1 \in \mathbb{R}$ , learning rate  $\eta > 0$ , threshold  $\theta > 0$   
**Initialize**  $\mathbf{u}^1 = u^1 \mathbf{1}_n$   
**for**  $t = 1$  to  $T$  **do**  
  Receive instance  $\mathbf{x}^t \in [0, 1]^n$   
  Predict  $\hat{y}^t = \begin{cases} +1 & \text{if } (\mathbf{u}^t \odot \mathbf{u}^t) \cdot \mathbf{x}^t \geq \theta \\ -1 & \text{otherwise} \end{cases}$   
  Receive label  $y^t$  and update:  
 $\mathbf{u}^{t+1} = \begin{cases} \mathbf{u}^t & \text{if } \hat{y}^t = y^t \\ \mathbf{u}^t + \eta y^t (\mathbf{u}^t \odot \mathbf{x}^t) & \text{otherwise} \end{cases}$

---

**Theorem 2 (Winnow bound (Littlestone, 1988; Warmuth, 2007))** *Given any sequence of examples  $(\mathbf{x}^t, y^t)$  such that  $\mathbf{x}^t \in [0, 1]^n$ , the labels  $y^t$  are  $\pm 1$ , and there is a weight  $\mathbf{r}$  with  $k$  ones and  $n - k$  zeros such that*

$$\mathbf{r} \cdot \mathbf{x}^t = \begin{cases} \geq \frac{1}{2} & \text{if } y^t = +1 \\ 0 & \text{otherwise,} \end{cases}$$

*then the Winnow algorithm makes at most  $7.18 k \log \frac{n}{k}$  mistakes on this sequence, when  $\eta \approx 1.28$ ,  $\theta = 0.19$ , and  $w^1 = n/k$ .*

Note for the sake of simplicity we only address the case when there exists a consistent disjunction. In the more general case, there are the additional terms in the mistake bound that involve the number of attribute errors w.r.t. the best disjunction.

**Theorem 3 (Reparameterized Winnow bound)** *Given that the assumptions of Theorem 2 hold, for any sequence of examples  $(\mathbf{x}^t, y^t)$  the reparameterized Winnow algorithm makes at most  $5.66 k \log \frac{n}{k}$  mistakes on this sequence, when  $\eta \approx 0.85$ ,  $\theta = 0.18$ , and  $\mathbf{u}^1 = \sqrt{n/k}$ .*

**Proof** We lower bound the per trial progress  $D_{\text{RE}}(\mathbf{r}, \mathbf{u}^t \odot \mathbf{u}^t) - D_{\text{RE}}(\mathbf{r}, \mathbf{u}^{t+1} \odot \mathbf{u}^{t+1})$  towards any comparator  $\mathbf{r}$  which satisfies the constraints in Theorem 2. Note that if no mistake occurs,  $\mathbf{u}^{t+1} = \mathbf{u}^t$ . Otherwise,  $\mathbf{u}^{t+1}$  is updated to  $\mathbf{u}^t + \eta y^t (\mathbf{u}^t \odot \mathbf{x}^t) = \mathbf{u}^t \odot (\mathbf{1} + \eta y^t \mathbf{x}^t)$ . Assuming  $\eta \leq 1$ , and using the facts that  $\log(1 + \eta y^t x_i^t) \geq x_i^t \log(1 + \eta y^t)$  and  $(x_i^t)^2 \leq x_i^t$  for  $x_i^t \in [0, 1]$ , we have

$$\begin{aligned} D_{\text{RE}}(\mathbf{r}, \mathbf{u}^t \odot \mathbf{u}^t) - D_{\text{RE}}(\mathbf{r}, \mathbf{u}^{t+1} \odot \mathbf{u}^{t+1}) &= 2\mathbf{r} \cdot \log(\mathbf{1} + \eta y^t \mathbf{x}^t) - (\mathbf{u}^t \odot \mathbf{u}^t) \cdot (2\eta y^t \mathbf{x}^t + \eta^2 \mathbf{x}^t \odot \mathbf{x}^t) \\ &\geq 2\mathbf{r} \cdot \mathbf{x}^t \log(1 + \eta y^t) - (\mathbf{u}^t \odot \mathbf{u}^t) \cdot \mathbf{x}^t (2\eta y^t + \eta^2). \end{aligned}$$

A mistake occurs if  $\hat{y}^t \neq y^t$ . If  $y^t = -1$ , then  $\mathbf{r} \cdot \mathbf{x}^t = 0$  by the assumption and  $(\mathbf{u}^t \odot \mathbf{u}^t) \cdot \mathbf{x}^t \geq \theta$ . Thus, the lower bound on the progress becomes

$$(\mathbf{u}^t \odot \mathbf{u}^t) \cdot \mathbf{x}^t (2\eta - \eta^2) \geq \theta (2\eta - \eta^2).$$

If  $y^t = +1$ , then we have  $\mathbf{r} \cdot \mathbf{x}^t \geq \frac{1}{2}$  by the assumption and also  $(\mathbf{u}^t \odot \mathbf{u}^t) \cdot \mathbf{x}^t \leq \theta$ . Thus in this case the progress is lower bounded by

$$\log(1 + \eta) - (\mathbf{u}^t \odot \mathbf{u}^t) \cdot \mathbf{x}^t (2\eta + \eta^2) \geq \log(1 + \eta) - \theta (2\eta + \eta^2).$$

Setting the two values to be equal, we obtain  $\theta = \frac{\log(1+\eta)}{4\eta}$ . With this choice of  $\theta$ , the progress per mistake is always lower bounded by  $1/4 (2 - \eta) \log(1 + \eta)$  which is maximized for  $\eta \approx 0.85$ . This choice of  $\eta$  yields  $\theta \approx 0.18$  and the progress is  $\approx 0.18$ . Summing over all trials and denoting the number of mistakes by  $M$ , we have

$$\underbrace{D_{\text{RE}}(\mathbf{r}, \mathbf{u}^1 \odot \mathbf{u}^1)}_{k \log \frac{n}{k}} - \underbrace{D_{\text{RE}}(\mathbf{r}, \mathbf{u}^{T+1} \odot \mathbf{u}^{T+1})}_{\geq 0} \geq 0.18 M, \text{ and thus } M \leq 5.56 k \log \frac{n}{k}. \quad \blacksquare$$

#### 4. The Reparameterization of the Hedge Algorithm

An important algorithm in the online expert setting is the Randomized Weighted Majority algorithm (Littlestone and Warmuth, 1994). Here, we only discuss the simplified version known as the Hedge algorithm (Freund and Schapire, 1997). This algorithm maintains a non-negative probability vector  $\mathbf{w}^t \in \Delta^{n-1}$  such that  $\sum_i w_i^t = 1$ . At trial  $t$ , the algorithm draws an expert/feature  $i$  with probability  $w_i^t$  and upon receiving the loss vector  $\boldsymbol{\ell}^t \in [0, 1]^n$ , it incurs the expected loss  $\mathbf{w}^t \cdot \boldsymbol{\ell}^t$ . The weights are then updated by a multiplicative exponentiated gradient term and re-normalized (see Algorithm 3) to assure that the non-negative weights sum to one.

The Hedge update and its reparameterization are a special case of EG and reparameterized EG when the losses are the dot loss  $\tilde{\mathbf{w}} \cdot \ell^t$  and  $\tilde{\mathbf{u}} \odot \tilde{\mathbf{u}} \cdot \ell^t$ , respectively. The Hedge update at round  $t$  is motivated by minimizing the relative entropy to the current weight  $\mathbf{w}^t$  plus the dot loss:

$$\mathbf{w}^{t+1} = \operatorname{argmin}_{\tilde{\mathbf{w}} \in \Delta^{n-1}} \left\{ \frac{1}{\eta} D_{\text{RE}}(\tilde{\mathbf{w}}, \mathbf{w}^t) + \tilde{\mathbf{w}} \cdot \ell^t \right\}. \quad (11)$$

---

**Algorithm 3** Hedge Algorithm

**Parameters** initial probability vector  $\mathbf{w}^1 \in \Delta^{n-1}$ , learning rate  $\eta > 0$   
**for**  $t = 1$  to  $T$  **do**  
 Draw expert  $i$  with probability  $w_i^t$   
 Incur loss  $\ell_i^t$  & expected loss  $\mathbf{w}^t \cdot \ell^t$   
 Update:

$$\mathbf{w}^{t+1} = \frac{\mathbf{w}^t \odot \exp(-\eta \ell^t)}{\sum_i w_i^t \exp(-\eta \ell_i^t)}$$


---

---

**Algorithm 4** Reparameterized Hedge Alg.

**Parameters** initial weight vector  $\mathbf{u}^1 \in \mathbb{R}^n$  s.t.  $\|\mathbf{u}^1\|_2 = 1$ , learning rate  $\eta > 0$   
**for**  $t = 1$  to  $T$  **do**  
 Draw expert  $i$  with probability  $(u_i^t)^2$   
 Incur loss  $\ell_i^t$  & expected loss  $(\mathbf{u}^t \odot \mathbf{u}^t) \cdot \ell^t$   
 Update:

$$\mathbf{u}^{t+1} = \frac{\mathbf{u}^t - \eta \mathbf{u}^t \odot \ell^t}{\|\mathbf{u}^t - \eta \mathbf{u}^t \odot \ell^t\|_2}$$


---

**Theorem 4 (Hedge bound (Freund and Schapire, 1997))** *For any sequence of loss vectors  $\{\ell^t\}_{t=1}^T \in [0, 1]^n$  such that  $\min_{i \in [n]} \sum_{t=1}^T \ell_i^t \leq L$ , any comparator  $\mathbf{r} \in \Delta^{n-1}$ , and any start vector  $\mathbf{w}^1 \in \Delta^{n-1}$  such that  $D_{\text{RE}}(\mathbf{r}, \mathbf{w}^1) \leq D \leq \log n$ , the total expected loss of the Hedge algorithm with start vector  $\mathbf{w}^1$  and learning rate  $\eta = \log(1 + \sqrt{2D/L})$  is bounded as*

$$\sum_t \mathbf{w}^t \cdot \ell^t \leq \sum_t \mathbf{r} \cdot \ell^t + \sqrt{2LD} + D_{\text{RE}}(\mathbf{r}, \mathbf{w}^1).$$

In the reparameterized Hedge update,  $\mathbf{w}$  is replaced with  $\mathbf{u} \odot \mathbf{u}$  where  $\mathbf{u} \in \mathbb{R}^n$  and  $\|\mathbf{u}\|_2 = 1$ . That is, the reparameterized weight  $\mathbf{u}$  lies on the unit sphere and the squared weight  $\mathbf{u} \odot \mathbf{u}$  corresponds to an  $n$ -dimensional probability vector. The update is motivated by minimizing the squared Euclidean distance as the inertia term plus the expected loss,<sup>6</sup>

$$\mathbf{u}^{t+1} = \operatorname{argmin}_{\tilde{\mathbf{u}} \text{ s.t. } \|\tilde{\mathbf{u}}\|_2=1} \left\{ \frac{1}{\eta} \|\tilde{\mathbf{u}} - \mathbf{u}^t\|_2^2 + (\tilde{\mathbf{u}} \odot \tilde{\mathbf{u}}) \cdot \ell^t \right\}. \quad (12)$$

Using a Lagrange multiplier to enforce the constraint that  $\sum_i (u_i^{t+1})^2 = 1$  and solving for  $\mathbf{u}^{t+1}$  yields the reparameterized Hedge update of Algorithm 4.

**Theorem 5 (Reparameterized Hedge bound)** *For any sequence of loss vectors  $\{\ell^t\}_{t=1}^T \in [0, 1]^n$  such that  $\min_{i \in [n]} \sum_{t=1}^T \ell_i^t \leq L$ , any comparator  $\mathbf{r} \in \Delta^{n-1}$  and any start vector  $\mathbf{u}^1 \in \mathbb{R}^n$  such that  $\|\mathbf{u}^1\|_2 = 1$  and  $D_{\text{RE}}(\mathbf{r}, \mathbf{u}^1 \odot \mathbf{u}^1) \leq D \leq \log n$ , the total expected loss of reparameterized Hedge with learning rate  $\eta = (1 + \sqrt{L/D})^{-1}$  is bounded as*

$$\sum_t (\mathbf{u}^t \odot \mathbf{u}^t) \cdot \ell^t \leq \sum_t \mathbf{r} \cdot \ell^t + 2\sqrt{LD} + D_{\text{RE}}(\mathbf{r}, \mathbf{u}^1 \odot \mathbf{u}^1).$$

---

6. A similar first-order Taylor approximation to (8) is required to obtain the explicit update.



**Proof** We lower bound the progress of the algorithm towards an arbitrary comparator  $\mathbf{r} \in \Delta^{n-1}$ . Assuming  $\eta < 1$ , the progress can be written as

$$\begin{aligned} D_{\text{RE}}(\mathbf{r}, \mathbf{u}^t \odot \mathbf{u}^t) - D_{\text{RE}}(\mathbf{r}, \mathbf{u}^{t+1} \odot \mathbf{u}^{t+1}) & \\ &= 2\mathbf{r} \cdot \log(\mathbf{1} - \eta \boldsymbol{\ell}^t) - \log((\mathbf{u}^t \odot \mathbf{u}^t) \cdot (\mathbf{1} - 2\eta \boldsymbol{\ell}^t + \eta^2 \boldsymbol{\ell}^t \odot \boldsymbol{\ell}^t)) \\ &\geq 2\mathbf{r} \cdot \boldsymbol{\ell}^t \log(1 - \eta) - \log(1 - (2\eta - \eta^2)(\mathbf{u}^t \odot \mathbf{u}^t) \cdot \boldsymbol{\ell}^t). \end{aligned}$$

Using  $\log(1 - x) \geq -x/(1 - x)$  and  $-\log(1 - x) \geq x$  for  $0 \leq x < 1$  yields

$$D_{\text{RE}}(\mathbf{r}, \mathbf{u}^t \odot \mathbf{u}^t) - D_{\text{RE}}(\mathbf{r}, \mathbf{u}^{t+1} \odot \mathbf{u}^{t+1}) \geq -\frac{2\eta}{1 - \eta} \mathbf{r} \cdot \boldsymbol{\ell}^t + (2\eta - \eta^2)(\mathbf{u}^t \odot \mathbf{u}^t) \cdot \boldsymbol{\ell}^t.$$

Summing over all trials and re-arranging the terms results in the following bound on the loss of the algorithm

$$\sum_t (\mathbf{u}^t \odot \mathbf{u}^t) \cdot \boldsymbol{\ell}^t \leq \frac{\frac{2\eta}{1 - \eta} \sum_t \mathbf{r} \cdot \boldsymbol{\ell}^t + D_{\text{RE}}(\mathbf{r}, \mathbf{u}^1 \odot \mathbf{u}^1) - D_{\text{RE}}(\mathbf{r}, \mathbf{u}^{T+1} \odot \mathbf{u}^{T+1})}{2\eta - \eta^2}.$$

Setting  $\eta$  to  $(1 + \sqrt{L/D})^{-1}$  and substituting the values of  $L$  and  $D$  yields the bound.  $\blacksquare$

Note that the regret bound of Theorem 5 for the reparameterized Hedge has an additional  $\sqrt{2}$  factor before the square root term in its regret bound. By plotting the progress, we can show that this additional factor disappears if you use the alternate learning rate  $\eta = \sqrt{D/(D + 2L)}$ . Based on this evidence, we conjecture that with the alternate tuning, the reparameterized Hedge has the same regret bound as the original.

Another approximation of the Hedge update has been analysed in (Cesa-Bianchi et al., 2007), called the *Prod* update. It replaces the exponential factors  $\exp(-\eta \ell_i)$  used in Hedge by their Taylor expansions  $(1 - \eta \ell_i)$  and normalizes multiplicatively.<sup>7</sup> Our reparameterized GD update is subtly different: it is a GD update w.r.t. the square roots of the weights as the parameters.

## 5. Reparameterizations of EGU and EG for Linear Regression

The regret bounds for linear regression using GD, EG, and EGU have been analyzed extensively in (Kivinen and Warmuth, 1997). Here, we derive similar bounds for the reparameterized EG and EGU updates.

We first recall the original EGU algorithm for linear regression which maintains a weight vector  $\mathbf{w}^t \in \mathbb{R}_{\geq 0}^n$ . Upon receiving input  $\mathbf{x}^t \in \mathbb{R}^n$  at round  $t$ , the algorithm predicts with  $\hat{y}^t = \mathbf{w}^t \cdot \mathbf{x}^t$ . It then receives the response  $y^t$ , incurs loss  $(y^t - \hat{y}^t)^2$ , and updates as:

$$\mathbf{w}^{t+1} = \mathbf{w}^t \odot \exp(-2\eta(\hat{y}^t - y^t)\mathbf{x}^t).$$

Kivinen and Warmuth (1997) analyze a slight variant of EGU: given a sequence of trials  $(\mathbf{x}^t, y^t)$  for which,  $y^t \in [0, Y]$  for all  $t$  and some  $Y > 0$ , the variant uses the following clipped prediction in its update:  $\hat{y}^t = \mathbf{w}^t \cdot \mathbf{x}^t$  if  $\hat{y}^t \leq Y$ , and  $\hat{y}^t = Y$  otherwise.

7. The same Taylor expansion is used in the *Approximated EG* update of (Kivinen and Warmuth, 1997), but that update keeps the weight sum as one by subtracting a term. It is motivated by the  $\chi^2$ -divergence.



**Theorem 6 (Linear regression with EGU (Kivinen and Warmuth, 1997))** *Let  $\{(\mathbf{x}^t, y^t)\}_{t=1}^T$  be any sequence such that  $\mathbf{x}^t \in [0, X]^n$  and  $y^t \in [0, Y]$  for some constants  $X, Y > 0$ . Then for any comparator  $\mathbf{r} \in \mathbb{R}_{\geq 0}^n$ , EGU with learning rate  $\eta = 1/(3XY)$  and arbitrary start point  $\mathbf{w}^1 \in \mathbb{R}_{\geq 0}^n$  satisfies the total loss bound*

$$\sum_{t=1}^T (y^t - \mathbf{w}^t \cdot \mathbf{x}^t)^2 \leq 3 \left( \sum_{t=1}^T (y^t - \mathbf{r} \cdot \mathbf{x}^t)^2 + XY D_{\text{RE}}(\mathbf{r}, \mathbf{w}^1) \right). \quad (13)$$

Furthermore, let  $L$  and  $D$  be constants such that  $\sum_t (y^t - \mathbf{r} \cdot \mathbf{x}^t)^2 \leq L$  and  $D_{\text{RE}}(\mathbf{r}, \mathbf{w}^1) \leq D$  and let

$$\eta = \frac{\sqrt{D}}{\sqrt{2LXY} + 2XY\sqrt{D}}. \quad (14)$$

Then we have

$$\sum_{t=1}^T (y^t - \mathbf{w}^t \cdot \mathbf{x}^t)^2 \leq \sum_{t=1}^T (y^t - \mathbf{r} \cdot \mathbf{x}^t)^2 + 2\sqrt{2LXYD} + 2XY D_{\text{RE}}(\mathbf{r}, \mathbf{w}^1). \quad (15)$$

Reparameterized EGU uses  $\mathbf{u}^t \odot \mathbf{u}^t$  as its weights,<sup>8</sup> where  $\mathbf{u}^1 \in \mathbb{R}^n$ , and updates as

$$\mathbf{u}^{t+1} = \mathbf{u}^t - \eta (\hat{y}^t - y^t) \mathbf{u}^t \odot \mathbf{x}^t, \quad (16)$$

where  $\hat{y}^t$  is again the clipped prediction:  $\min(\mathbf{w}^t \cdot \mathbf{x}^t, Y)$ .

**Theorem 7 (Reparameterized EGU linear regression)** *Let  $\{(\mathbf{x}^t, y^t)\}_{t=1}^T$  be any sequence such that  $\mathbf{x}^t \in [0, X]^n$  and  $y^t \in [0, Y]$  for some constants  $X, Y > 0$ . Then for any comparator  $\mathbf{r} \in \mathbb{R}_{\geq 0}^n$ , reparameterized EGU with learning rate  $\eta = 1/(3XY)$  and arbitrary start point  $\mathbf{w}^1 \in \mathbb{R}_{\geq 0}^n$  satisfies the same total loss bound (13) as the original EGU when replacing  $\mathbf{w}^t$  with  $\mathbf{u}^t \odot \mathbf{u}^t$ . Moreover, for the same constants  $L$  and  $D$  as in Theorem 6 and by setting  $\eta$  to (14), the reparameterized EGU achieves the same bound as in (15).*

**Proof sketch** We lower bound the progress  $D_{\text{RE}}(\mathbf{r}, \mathbf{u}^t \odot \mathbf{u}^t) - D_{\text{RE}}(\mathbf{r}, \mathbf{u}^{t+1} \odot \mathbf{u}^{t+1})$  of the algorithm towards any comparator  $\mathbf{r} \in \mathbb{R}_{\geq 0}^n$  by  $a(y^t - \hat{y})^2 - b(y^t - \mathbf{r} \cdot \mathbf{x}^t)^2$ , for some constants  $a, b \geq 0$ . We can lower bound the progress as

$$\begin{aligned} & D_{\text{RE}}(\mathbf{r}, \mathbf{u}^t \odot \mathbf{u}^t) - D_{\text{RE}}(\mathbf{r}, \mathbf{u}^{t+1} \odot \mathbf{u}^{t+1}) \\ & \geq \frac{2\mathbf{r} \cdot \mathbf{x}^t \log(1 - \eta(\hat{y}^t - y^t)X)}{X} + (2\eta(\hat{y}^t - y^t) - \eta^2(\hat{y}^t - y^t)^2 X) (\mathbf{u}^t \odot \mathbf{u}^t) \cdot \mathbf{x}^t. \end{aligned}$$

Denoting by  $s := \mathbf{r} \cdot \mathbf{x}^t$  and  $p := (\mathbf{u}^t \odot \mathbf{u}^t) \cdot \mathbf{x}^t$ , it suffices to show that  $G(p, \hat{y}, y, s) \leq 0$  where (omitting the superscript  $t$ )

$$G(p, \hat{y}, y, s) = -\frac{2s \log(1 - \eta(\hat{y} - y)X)}{X} - (2\eta(\hat{y} - y) - \eta^2(\hat{y} - y)^2 X) p - a(y - \hat{y})^2 + b(y - s)^2.$$

It suffices to show the result for  $p = \hat{y}$ . The function  $G(\hat{y}, \hat{y}, y, s)$  is maximized for  $s = y - \log(1 - \eta(\hat{y} - y)X)/(Xb)$ . Setting  $\eta = b/(1 + 2XYb)$ , the inequality holds for any

8. Note that the squared loss  $(\mathbf{w} \cdot \mathbf{x} - y)^2$  is convex in  $\mathbf{w}$ , but  $(\mathbf{u} \odot \mathbf{u} \cdot \mathbf{x} - y)^2$  is not convex in  $\mathbf{u}$ .

$a \leq b/(1 + 2XYb)$ . Setting  $b = 1/(XY)$ , the bound is achieved for  $\eta = a = 1/(3XY)$ .  $\blacksquare$

Algorithm 5 EGU Linear Regression	Algorithm 6 Reparam. EGU Lin. Reg.
<b>Parameters</b> initial weight $\mathbf{w}^1 \in \mathbb{R}_{\geq 0}^n$ , learning rate $\eta > 0$ <b>for</b> $t = 1$ to $T$ <b>do</b> Receive instance $\mathbf{x}^t \in [0, X]^n$  Predict $\hat{y}^t = \begin{cases} \mathbf{w}^t \cdot \mathbf{x}^t & \text{if } \mathbf{w}^t \cdot \mathbf{x}^t \leq Y \\ Y & \text{otherwise} \end{cases}$  Receive label $y^t$ and update: $\mathbf{w}^{t+1} = \mathbf{w}^t \odot \exp(-2\eta(\hat{y}^t - y^t)\mathbf{x}^t)$	<b>Parameters</b> initial weight $\mathbf{u}^1 \in \mathbb{R}^n$ , learn- ing rate $\eta > 0$ <b>for</b> $t = 1$ to $T$ <b>do</b> Receive instance $\mathbf{x}^t \in [0, X]^n$  Predict $\hat{y}^t = \begin{cases} \mathbf{u}^t \odot \mathbf{u}^t \cdot \mathbf{x}^t & \text{if } \mathbf{u}^t \odot \mathbf{u}^t \cdot \mathbf{x}^t \leq Y \\ Y & \text{otherwise} \end{cases}$  Receive label $y^t$ and update: $\mathbf{u}^{t+1} = \mathbf{u}^t - \eta(\hat{y}^t - y^t)\mathbf{u}^t \odot \mathbf{x}^t$

Alternatively, the (normalized) EG algorithm maintains a probability vector  $\mathbf{w}^t \in \Delta^{n-1}$  as its weight vector. The update for EG is the same as EGU, except for a multiplicative normalization:

$$\mathbf{w}^{t+1} = \frac{\mathbf{w}^t \odot \exp(-2\eta(\hat{y}^t - y^t)\mathbf{x}^t)}{\sum_i w_i^{t+1} \exp(-2\eta(\hat{y}^t - y^t)x_i^t)}.$$

The following theorem from (Kivinen and Warmuth, 1997) expresses the worst-case bound of running EG for linear regression.

**Theorem 8 (EG linear regression (Kivinen and Warmuth, 1997))** *Let  $\{(\mathbf{x}^t, y^t)\}_{t=1}^T$  be any sequences such that  $\max_i x_i^t - \min_i x_i^t \leq R$  for some  $R \geq 0$ . Then for any comparator  $\mathbf{r} \in \Delta^{n-1}$ , the EG algorithm with learning rate  $\eta = 2/(3R^2)$  and arbitrary start point  $\mathbf{w}^1 \in \Delta^{n-1}$  satisfies the total loss bound*

$$\sum_{t=1}^T (y^t - \mathbf{w}^t \cdot \mathbf{x}^t)^2 \leq \frac{3}{2} \left( \sum_{t=1}^T (y^t - \mathbf{r} \cdot \mathbf{x}^t)^2 + R^2 D_{\text{RE}}(\mathbf{r}, \mathbf{w}^1) \right). \quad (17)$$

Furthermore, let  $L$  and  $D$  be constants such that  $\sum_t (y^t - \mathbf{r} \cdot \mathbf{x}^t)^2 \leq L$  and  $D_{\text{RE}}(\mathbf{r}, \mathbf{w}^1) \leq D$  and let

$$\eta = \frac{2\sqrt{D}}{R\sqrt{2L} + R^2\sqrt{D}}. \quad (18)$$

Then we have

$$\sum_{t=1}^T (y^t - \mathbf{w}^t \cdot \mathbf{x}^t)^2 \leq \sum_{t=1}^T (y^t - \mathbf{r} \cdot \mathbf{x}^t)^2 + R\sqrt{2LD} + \frac{R^2}{2} D_{\text{RE}}(\mathbf{r}, \mathbf{w}^1). \quad (19)$$

Using (10), the reparameterized EG update for linear update can be written as

$$\mathbf{u}^{t+1} = \frac{\mathbf{u}^t - \eta(\hat{y}^t - y^t)\mathbf{u}^t \odot \mathbf{x}^t}{\|\mathbf{u}^t - \eta(\hat{y}^t - y^t)\mathbf{u}^t \odot \mathbf{x}^t\|_2}. \quad (20)$$

The algorithm for EG on linear regression is similar to Algorithm 6, except the initial weight  $\mathbf{u}^1$  satisfies  $\|\mathbf{u}^1\|_2 = 1$  and the update is replaced with (20). In the following, we

show the regret for the reparameterized EG update for linear regression. For simplicity, we first consider the case where the input  $\mathbf{x} \in [0, X]^n$  for some  $X > 0$ .

**Theorem 9 (Reparameterized EG linear regression)** *Let  $\{(\mathbf{x}^t, y^t)\}_{t=1}^T$  be any sequences such that  $\mathbf{x}^t \in [0, X]^n$  for some  $X \geq 0$ . Then for any comparator  $\mathbf{r} \in \Delta^{n-1}$ , the Reparameterized EG algorithm with learning rate  $\eta = 1/(3X^2)$  and arbitrary start point  $\mathbf{u}^1 \in \mathbf{R}^n$  such that  $\|\mathbf{u}^1\|_2 = 1$  satisfies the total loss bound*

$$\sum_{t=1}^T (y^t - (\mathbf{u}^t \odot \mathbf{u}^t) \cdot \mathbf{x}^t)^2 \leq 3 \left( \sum_t (y^t - \mathbf{r} \cdot \mathbf{x}^t)^2 + X^2 D_{\text{RE}}(\mathbf{r}, \mathbf{u}^1 \odot \mathbf{u}^1) \right). \quad (21)$$

Furthermore, let  $L$  and  $D$  be constants such that  $\sum_t (y^t - \mathbf{r} \cdot \mathbf{x}^t)^2 \leq L$  and  $D_{\text{RE}}(\mathbf{r}, \mathbf{u}^1 \odot \mathbf{u}^1) \leq D$  and let

$$\eta = \frac{2\sqrt{D}}{X\sqrt{2L} + X^2\sqrt{D}}. \quad (22)$$

Then we have

$$\sum_{t=1}^T (y^t - (\mathbf{u}^t \odot \mathbf{u}^t) \cdot \mathbf{x}^t)^2 \leq \sum_{t=1}^T (y^t - \mathbf{r} \cdot \mathbf{x}^t)^2 + 2X\sqrt{2LD} + 2X^2 D_{\text{RE}}(\mathbf{r}, \mathbf{u}^1). \quad (23)$$

The proof is given in the appendix. In the following, we claim a bound for the reparameterized EG for the more general case where  $\mathbf{x}^t \in [-X, X]^n$  for all  $t$ . We provide a proof sketch for the claim in the appendix. Our partial proof shows strong evidence for the existence of such regret bound. However, finding an analytical proof by simplifying the final form is left for future work.

**Claim 1 (More general bound for reparameterized EG linear regression)** *For the same setting as in Theorem 9 but a more general case where  $\mathbf{x}^t \in [-X, X]^n$  for all  $t$ , the reparameterized EG algorithm achieves the same bound (23) in which  $X$  is replaced by  $2X$ .*

## 6. Lower Bounds and Simulations

Linear lower bounds for GD have been studied in the previous work (Kivinen et al., 1997; Vishwanathan and Warmuth, 2005) for the hinge and squared loss. These lower bounds are all based on the ‘‘Hadamard problem’’: the instances are the  $n$  rows of the  $n$ -dimensional Hadamard matrix in random order and the target is one of the columns.<sup>9</sup> Thus target is a unit weight vector selecting the right column. For example, for linear regression, as online GD passes over the  $n$  examples, it brings up the weight of the right column slowly and the average squared loss on all  $n$  instances decreases linearly (Figure 3(a)). Multiplicative updates such as EGU bring up the target weight dramatically faster and

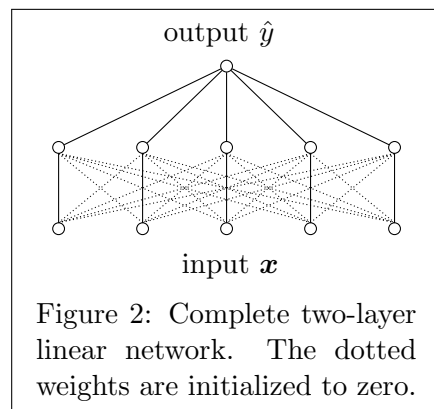


Figure 2: Complete two-layer linear network. The dotted weights are initialized to zero.

9. Similar behavior is observed if a random  $\pm 1$  matrix is used.

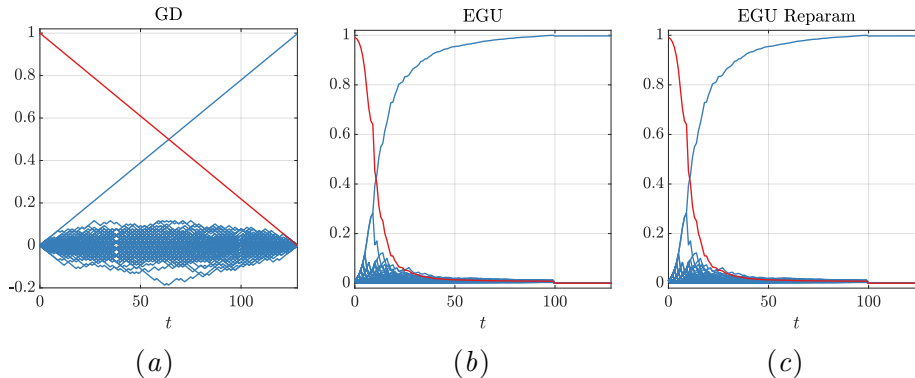


Figure 3: GD, EGU, and Reparameterized EGU on the online Hadamard problem ( $n=128$ ): at round  $t$ , we train until consistency on the past  $t$  examples. For  $1 \leq t \leq 128$ , all 128 weights are shown in blue and the average loss over all 128 examples is shown in red.

the average squared loss decays essentially after  $\log n$  examples (Figure 3(b)). Surprisingly, GD on the  $u_i$ 's of the sparse linear network of Figure 1 has visually identical weights<sup>10</sup> (Figure 3(c)) and average loss trajectories as EGU (Figure 3(b)). Furthermore, if we run GD on the fully-connected two-layer network of Figure 2 with the dotted weights initialized to zero, then the combined linear weights of both layers again behave as when a simple linear neuron is trained with GD (Appendix D). This behavior seems to be related to the observation that GD focuses on keeping the highest weight as small as possible and uses the additional weights to overfit.

## 7. Open Problems

There are a number of technical open problems associated with the current paper. A complete analytical proof is still needed for the regret of the reparameterized EG when the instance domain is two-sided, i.e.  $\mathbf{x}^t \in [-X, X]^n$ . Note also that regret bounds for the two-sided domain do not exist for the original and reparameterized EGU. Experimentally however, the two updates show the typical  $\log n$  behavior of multiplicative updates. For the lower bound discussion, a next step would be to show that for the Hadamard problem, GD on a two-layer linear network with any initialization has the typical linear decaying loss. Also in this paper, we were able to repeat the regret bounds for the Reparameterized Hedge algorithm (i.e. when the loss is the dot loss and the loss components lie in  $[0, 1]$ ). However, so far we were not able to repeat  $\mathcal{O}(\log n)$  regret bounds (Vovk, 1990; Haussler et al., 1998) for a reparameterized version of the multiplicative expert algorithm when the loss is the dot loss and the loss components are the squared loss (a.k.a. the Brier scores). We conjecture that  $\mathcal{O}(\log n)$  regret bounds are not possible for the reparameterized version of the expert algorithm when the loss components are log losses. The reason is the unboundedness of the log loss. Lastly, there are many open problems regarding neural networks. For example, are reparameterized multiplicative updates useful for training large networks and how the optimization routines such as AdaGrad and Adam interact with the new reparameterized updates?

**Acknowledgement** We thank Zakaria Mhammedi for valuable feedback. Also E. Amid was partially supported by the NSF grant IIS 1546459 and a gift grant from the Intel Corp.

<sup>10</sup>. The trajectories are very close, but not numerically equal.

## References

- Ethan Akin. *The geometry of population genetics*, volume 31 of *Lecture Notes in Biomathematics*. Springer-Verlag, Berlin-New York, 1979.
- Ehsan Amid and Manfred K. Warmuth. Reparameterizing mirror descent as gradient descent. *arXiv preprint arXiv:2002.10487*, 2020.
- Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit acceleration by overparameterization. In *International Conference on Machine Learning (ICML)*, 2018.
- Nicolo Cesa-Bianchi, Yishay Mansour, and Gilles Stoltz. Improved second-order bounds for prediction with expert advice. *Machine Learning*, 66(2-3):321–352, 2007.
- Michal Dereziński and Manfred K. Warmuth. The limits of squared Euclidean distance regularization. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 2807–2815, 2014.
- Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997.
- C. Gentile and M. K. Warmuth. Hinge loss and average margin. In *Advances in Neural Information Processing Systems (NeurIPS)*, 1998.
- Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Implicit regularization in matrix factorization. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 6151–6159, 2017.
- David Haussler, Jyrki Kivinen, and Manfred K Warmuth. Sequential prediction of individual sequences under general loss functions. *IEEE Transactions on Information Theory*, 44(5):1906–1925, 1998.
- Jyrki Kivinen and Manfred K. Warmuth. Exponentiated gradient versus gradient descent for linear predictors. *Information and Computation*, 132(1):1–63, 1997.
- Jyrki Kivinen, Manfred K Warmuth, and Peter Auer. The Perceptron algorithm versus Winnow: linear versus logarithmic mistake bounds when few input variables are relevant. *Artificial Intelligence*, 97(1-2):325–343, 1997.
- Jyrki Kivinen, Manfred K Warmuth, and Babak Hassibi. The p-norm generalization of the LMS algorithm for adaptive filtering. *IEEE Transactions on Signal Processing*, 54(5):1782–1793, 2006.
- N Littlestone and MK Warmuth. The weighted majority algorithm. *Information and Computation*, 108(2):212–261, 1994.
- Nick Littlestone. Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm. *Machine learning*, 2(4):285–318, 1988.

- R Tyrrell Rockafellar. Monotone operators and the proximal point algorithm. *SIAM journal on control and optimization*, 14(5):877–898, 1976.
- William H Sandholm. *Population games and evolutionary dynamics*. MIT Press, 2010.
- Shai Shalev-Shwartz et al. Online learning and online convex optimization. *Foundations and Trends® in Machine Learning*, 4(2):107–194, 2012.
- Tomas Vaskevicius, Varun Kanade, and Patrick Rebeschini. Implicit regularization for optimal sparse recovery. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 2968–2979, 2019.
- S.V.N. Vishwanathan and M.K. Warmuth. Leaving the span. In *Proceedings of the 18th Annual Conference on Learning Theory (COLT)*, 2005.
- Volodimir G Vovk. Aggregating strategies. In *Proceedings of the third annual workshop on Computational learning theory*, pages 371–386, 1990.
- Manfred K. Warmuth. Winnowing subspaces. In *Proceedings of the 24th International Conference on Machine Learning (ICML)*, pages 999–1006, 2007.
- Blake Woodworth, Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Kernel and deep regimes in overparametrized models. *arXiv preprint arXiv:1906.05827*, 2019.

## Appendix A. Proof of Theorem 7

**Proof** We first establish a lower bound of the form,

$$a(y^t - \hat{y}^t)^2 - b(y^t - \mathbf{r} \cdot \mathbf{x}^t)^2 \leq D_{\text{RE}}(\mathbf{r}, \mathbf{u}^t \odot \mathbf{u}^t) - D_{\text{RE}}(\mathbf{r}, \mathbf{u}^{t+1} \odot \mathbf{u}^{t+1}), \quad (24)$$

on the progress of the algorithm towards the comparator  $\mathbf{r}$ , for some constants  $0 \leq a, b$ . Assuming that  $\eta \leq 1/(2XY)$ , we have  $(1 - \eta(\hat{y}^t - y^t)x_i^t) \geq 0$  for all  $i$ . Thus, we can lower bound the progress as

$$\begin{aligned} & D_{\text{RE}}(\mathbf{r}, \mathbf{u}^t \odot \mathbf{u}^t) - D_{\text{RE}}(\mathbf{r}, \mathbf{u}^{t+1} \odot \mathbf{u}^{t+1}) \\ &= 2\mathbf{r} \cdot \log(\mathbf{1} - \eta(\hat{y}^t - y^t)\mathbf{x}^t) - \log(\mathbf{u}^t \cdot (\mathbf{1} - 2\eta(\hat{y}^t - y^t)\mathbf{x}^t + \eta(\hat{y}^t - y^t)^2\mathbf{x}^t \odot \mathbf{x}^t - \mathbf{1})). \end{aligned}$$

Applying the inequalities  $\log(1 - \eta(\hat{y}^t - y^t)x_i^t) \geq \frac{x_i^t \log(1 - \eta(\hat{y}^t - y^t)X)}{X}$  and  $\log(1 + x) \leq x$  for  $x > -1$  and using the fact that  $(x_i^t)^2 \leq x_i^t X$ , we have

$$\begin{aligned} & D_{\text{RE}}(\mathbf{r}, \mathbf{u}^t \odot \mathbf{u}^t) - D_{\text{RE}}(\mathbf{r}, \mathbf{u}^{t+1} \odot \mathbf{u}^{t+1}) \\ & \geq \frac{2\mathbf{r} \cdot \mathbf{x}^t \log(1 - \eta(\hat{y}^t - y^t)X)}{X} + (2\eta(\hat{y}^t - y^t) - \eta^2(\hat{y}^t - y^t)^2 X) (\mathbf{u}^t \odot \mathbf{u}^t) \cdot \mathbf{x}^t. \end{aligned}$$

Denoting by  $s := \mathbf{r} \cdot \mathbf{x}^t$  and  $p := (\mathbf{u}^t \odot \mathbf{u}^t) \cdot \mathbf{x}^t$ , it suffices to show that  $G(p, \hat{y}, y, s) \leq 0$  where (omitting the superscript  $t$ )

$$G(p, \hat{y}, y, s) = -\frac{2s \log(1 - \eta(\hat{y} - y)X)}{X} - (2\eta(\hat{y} - y) - \eta^2(\hat{y} - y)^2 X) p + a(y - \hat{y})^2 - b(y - s)^2.$$

Recall that the prediction  $\hat{y}^t$  of the reparameterized EGU is given by  $\hat{y}^t = p = (\mathbf{u}^t \odot \mathbf{u}^t) \cdot \mathbf{x}^t$  if  $p \leq Y$  holds; otherwise  $\hat{y}^t = Y$ . Thus, we need to show  $G(p, \hat{y}, y, s) \leq 0$  for two cases: when  $\hat{y} = p$ , and for  $0 \leq \hat{y} = Y < p$ . Recall that by the assumption  $0 \leq y \leq Y$  and  $\eta \leq 1/(2XY)$ . Therefore,  $G(p, \hat{y}, y, s)$  is non-increasing in  $p$  for  $\hat{y} \geq y$ . Hence, the condition  $G(p, \hat{y}, y, s) \leq 0$  for  $\hat{y} = Y < p$  is satisfied if  $G(Y, Y, y, s) \leq 0$  holds. Thus, it suffices to show the result for  $0 \leq \hat{y} = p \leq Y$ .

For fixed  $\hat{y}$  and  $y$ , the function  $G(\hat{y}, \hat{y}, y, s)$  is maximized for

$$s = y - A/(Xb), \quad \text{where} \quad A := \log(1 - \eta(\hat{y} - y)X).$$

Plugging in this value, we have  $H(\hat{y}, y) := G(\hat{y}, \hat{y}, y, y - A/(Xb))$  where

$$H(\hat{y}, y) = -2\eta(\hat{y} - y)\hat{y} + \eta^2(\hat{y} - y)^2 \hat{y} X + a(\hat{y} - y)^2 + A^2/(bX^2) - 2Ay/X.$$

In order to obtain the bound (24), we show that  $H(\hat{y}, y) \leq 0$  holds for the choice of  $\eta = b/(1 + 2XYb)$  and  $a = b/(1 + 2XYb)$ . Substituting these values for  $\eta$  and  $a$ , we have  $H(y, y) = \frac{\partial H(\hat{y}, y)}{\partial \hat{y}} \Big|_{\hat{y}=y} = 0$  and Furthermore,

$$\frac{\partial^2 H(\hat{y}, y)}{\partial \hat{y}^2} \Big|_{\hat{y}=y} = -\frac{4b^2 X(Y - y)}{(1 + 2XYb)^2} \leq 0 \quad \text{for all } 0 \leq y \leq Y.$$



Thus, it suffices to show that  $y = \hat{y}$  is the only maximum of the function. For this, we first write

$$\begin{aligned} \frac{\partial^2 H(\hat{y}, y)}{\partial y^2} = & -\frac{1}{\left(\frac{1}{2} + X\left(Y + \frac{y}{2} - \frac{\hat{y}}{2}\right)b\right)^2 \left(\frac{1}{2} + XYb\right)^2} \times \left( \left(\frac{1}{2} + XYb\right)^2 \log\left(1 + \frac{X(y - \hat{y})b}{1 + 2XYb}\right)b \right. \\ & + \left. \left(Y^3 - \left(\frac{\hat{y} + y}{2}\right)Y^2 + \left(\frac{\hat{y}^2 - y^2}{4}\right)Y + \frac{\hat{y}(\hat{y} - y)}{8}\right)b^4 X^3 \right. \\ & \left. + \left(Y^2 - \left(\frac{\hat{y} + y}{2}\right)Y + \frac{\hat{y}^2 - y^2}{8}\right)b^3 X^2 + \left(\frac{Y}{4} + \frac{\hat{y} - y}{8}\right)b^2 X \right). \end{aligned}$$

Using the inequality  $\log(1 + x) \leq x$  for  $x > -1$ , we can upper bound the log-term as

$$\log\left(1 + \frac{X(y - \hat{y})b}{1 + 2XYb}\right)b \leq \frac{X(y - \hat{y})b^2}{1 + 2XYb},$$

and write the new function as  $Q(\hat{y}, y)$  such that  $\frac{\partial^2 H(\hat{y}, y)}{\partial y^2} \leq Q(\hat{y}, y)$ . It is trivial to check that the derivative

$$\frac{\partial Q(\hat{y}, y)}{\partial y} = \frac{3\left(\left(Y + \frac{y}{6} - \frac{\hat{y}}{2}\right)\left(Y + \frac{y}{2} - \frac{\hat{y}}{2}\right)bX + \frac{Y}{2} + \frac{y}{3} - \frac{\hat{y}}{2}\right)X^2 b^3}{2\left(\frac{1}{2} + X\left(Y + \frac{y}{2} - \frac{\hat{y}}{2}\right)b\right)^4} \geq 0,$$

for any  $0 \leq y \leq Y$  and  $0 \leq \hat{y} \leq Y$ . Hence, it remains to check that  $Q(\hat{y}, Y) \leq 0$  holds:

$$\begin{aligned} Q(\hat{y}, Y) = & -\frac{1}{9\left(\frac{1}{3} + X\left(Y - \frac{\hat{y}}{3}\right)b\right)^3 \left(\frac{1}{2} + XYb\right)^2} \times \\ & (Y - \hat{y})Xb^2 \left(\frac{1}{3} + \left(Y - \frac{\hat{y}}{3}\right)\left(Y - \hat{y}\right)\left(Y - \frac{\hat{y}}{2}\right)b^3 X^3 + \frac{5\left(Y^2 - \frac{3}{5}Y\hat{y} + \frac{2}{15}\hat{y}^2\right)b^2 X^2}{2} + \frac{5\left(Y - \frac{\hat{y}}{5}\right)bX}{3}\right) \leq 0, \end{aligned}$$

which holds for any  $0 \leq \hat{y} \leq Y$ . This implies  $\frac{\partial^2 H(\hat{y}, y)}{\partial y^2} \leq Q(\hat{y}, y) \leq 0$  for any  $0 \leq \hat{y} \leq Y$  and  $0 \leq y \leq Y$ .

The remainder of the proof follows similarly to (Kivinen and Warmuth, 1997). Specifically, by setting  $b = c/(XY)$ , we obtain

$$\sum_{t=1}^T (y^t - (\mathbf{u}^t \odot \mathbf{u}^t) \cdot \mathbf{x}^t)^2 \leq (1 + 2c) \sum_{t=1}^T (y^t - \mathbf{r} \cdot \mathbf{x}^t)^2 + \left(2 + \frac{1}{c}\right) XY D_{\text{RE}}(\mathbf{r}, \mathbf{u}^1 \odot \mathbf{u}^1).$$

Setting  $c = 1$ , the bound in (13) is achieved for  $\eta = a = 1/(3XY)$ . Using the values  $L$  and  $D$  and tuning for  $c$  achieves (15) for the choice of  $\eta$  as in (14). ■

## Appendix B. Proof of Theorem 9

Similar to the proof of Theorem 7, we establish a lower bound of the form

$$a(y^t - \hat{y}^t)^2 - b(y^t - \mathbf{r} \cdot \mathbf{x}^t)^2 \leq D_{\text{RE}}(\mathbf{r}, \mathbf{u}^t \odot \mathbf{u}^t) - D_{\text{RE}}(\mathbf{r}, \mathbf{u}^{t+1} \odot \mathbf{u}^{t+1}), \quad (25)$$

on the progress of the algorithm towards the comparator  $\mathbf{r} \in \Delta^{n-1}$ , and for some constants  $a, b \geq 0$ . Assuming that  $\eta \leq 1/(2X^2)$ , we have  $(1 - \eta(\hat{y}^t - y^t)x_i^t) \geq 0$  for all  $i$ . Thus, we can lower bound the progress as

$$\begin{aligned} & D_{\text{RE}}(\mathbf{r}, \mathbf{u}^t \odot \mathbf{u}^t) - D_{\text{RE}}(\mathbf{r}, \mathbf{u}^{t+1} \odot \mathbf{u}^{t+1}) \\ &= 2\mathbf{r} \cdot \log(\mathbf{1} - \eta(\hat{y}^t - y^t)\mathbf{x}^t) - \log(\mathbf{w}^t \cdot (\mathbf{1} - 2\eta(\hat{y}^t - y^t)\mathbf{x}^t + \eta(\hat{y}^t - y^t)^2\mathbf{x}^t \odot \mathbf{x}^t)). \end{aligned}$$

Applying the inequalities  $\log(1 - \eta(\hat{y}^t - y^t)x_i^t) \geq \frac{x_i^t \log(1 - \eta(\hat{y}^t - y^t)X)}{X}$  and  $-\log(1 - x) \geq x$  for  $0 \leq x \leq 1$  and using the fact that  $(x_i^t)^2 \leq x_i^t X$ , we have

$$\begin{aligned} & D_{\text{RE}}(\mathbf{r}, \mathbf{u}^t \odot \mathbf{u}^t) - D_{\text{RE}}(\mathbf{r}, \mathbf{u}^{t+1} \odot \mathbf{u}^{t+1}) \\ & \geq \frac{2\mathbf{r} \cdot \mathbf{x}^t \log(1 - \eta(\hat{y}^t - y^t)X)}{X} + (2\eta(\hat{y}^t - y^t) - \eta^2(\hat{y}^t - y^t)^2 X) \hat{y}^t. \end{aligned}$$

Denoting by  $s := \mathbf{r} \cdot \mathbf{x}^t$ , it suffices to show that  $G(\hat{y}, y, s) \leq 0$  where (omitting the superscript  $t$ )

$$G(\hat{y}, y, s) = -\frac{2r \log(1 - \eta(\hat{y} - y)X)}{X} - (2\eta(\hat{y} - y) - \eta^2(\hat{y} - y)^2 X) \hat{y} + a(y - \hat{y})^2 - b(y - s)^2.$$

For fixed  $\hat{y}$  and  $y$ , the function  $G(\hat{y}, y, s)$  is maximized for

$$s = y - A/(Xb), \quad \text{where } A := \log(1 - \eta(\hat{y} - y)X).$$

Plugging in this value, we have  $H(\hat{y}, y) := G(\hat{y}, y, y - A/(Xb))$  where

$$H(\hat{y}, y) = -2\eta(\hat{y} - y)\hat{y} + \eta^2(\hat{y} - y)^2 \hat{y} X + a(\hat{y} - y)^2 + A^2/(bX^2) - 2Ay/X.$$

It is easy to verify that  $H(y, y) = \frac{\partial H(\hat{y}, y)}{\partial y} \Big|_{\hat{y}=y} = 0$ . Moreover, notice that

$$\frac{\partial^2 H(\hat{y}, y)}{\partial y^2} \Big|_{\hat{y}=y} = \frac{(4Xy + 2)\eta^2 - 4b\eta + 2ab}{b},$$

which is minimized at  $y = X$  for  $\eta = b/(1 + 2X^2b)$ . Plugging in this value, we have  $\frac{\partial^2 H(\hat{y}, y)}{\partial y^2} \leq 0$  for  $a \leq b/(1 + 2X^2b)$ . Now there remains to show that  $\hat{y} = y$  is the only maximum of the function  $H(\hat{y}, y)$ . This is easily verified by plugging in  $a = b/(1 + 2X^2b)$  and observing

$$\begin{aligned} \frac{\partial^2 H(\hat{y}, y)}{\partial y^2} &= -\frac{b}{\left(\frac{1}{2} + X^2b\right)^2 \left(\frac{1}{2} + \left(X + \frac{y}{2} - \frac{\hat{y}}{2}\right)Xb\right)^2} \times \\ & \quad \left(\frac{1}{2} \left(\frac{1}{2} + X^2b\right)^2 \log\left(1 + \frac{(y - yh)Xb}{1 + 2X^2b}\right)\right. \\ & \quad \left.+ Xb \left(\left(X^3 - \left(\frac{y + yh}{2}\right)X^2 + \left(\frac{\hat{y}^2 - y^2}{4}\right)X - \frac{\hat{y}(y - \hat{y})^2}{8}\right)X^2b^2\right.\right. \\ & \quad \left.\left.+ X^2b \left(X^2 - \left(\frac{y + yh}{2}\right)X + \left(\frac{\hat{y}^2 - y^2}{8}\right)\right) + \frac{X^2}{4} - \frac{y + \hat{y}}{8}X\right) \leq 0, \end{aligned}$$

for  $0 \leq y, \hat{y} \leq X$  and

$$\frac{\partial^2 H(\hat{y}, y)}{\partial y^2} \Big|_{\hat{y}=y} = -\frac{4b^2 X(X-y)}{(1+2X^2b)^2}.$$

Finally, setting  $b = c/(X^2)$  for some  $c > 0$ , we have

$$\sum_{t=1}^T (y^t - (\mathbf{u}^t \odot \mathbf{u}^t) \cdot \mathbf{x}^t)^2 \leq (1+2c) \sum_{t=1}^T (y^t - \mathbf{r} \cdot \mathbf{x}^t)^2 + \left(2 + \frac{1}{c}\right) XY D_{\text{RE}}(\mathbf{r}, \mathbf{u}^1 \odot \mathbf{u}^1).$$

Setting  $c = 1$  establishes the first bound. Similarly, optimizing for  $c$  yields

$$c = \frac{X\sqrt{D}}{\sqrt{2L}}.$$

For this choice of  $c$  and  $\eta$  as in (22), we obtain the second bound.

### Appendix C. Proof Sketch of Claim 1

We can lower bound the progress of the algorithm as (assuming  $\eta \leq 1/(2X^2)$ ),

$$\begin{aligned} & D_{\text{RE}}(\mathbf{r}, \mathbf{u}^t \odot \mathbf{u}^t) - D_{\text{RE}}(\mathbf{r}, \mathbf{u}^{t+1} \odot \mathbf{u}^{t+1}) \\ &= 2\mathbf{r} \cdot \log(\mathbf{1} - \eta(\hat{y}^t - y^t)\mathbf{x}^t) - \log(\mathbf{w}^t \cdot (\mathbf{1} - 2\eta(\hat{y}^t - y^t)\mathbf{x}^t + \eta(\hat{y}^t - y^t)^2\mathbf{x}^t \odot \mathbf{x}^t)) \\ &\geq \frac{(s+X)}{X} \log\left(\frac{1 - \eta(\hat{y}^t - y^t)X}{1 + \eta(\hat{y}^t - y^t)X}\right) + 2 \log(1 + \eta(\hat{y}^t - y^t)X) \\ &\quad - \log(1 - 2\eta(\hat{y}^t - y^t)\hat{y}^t + \eta^2(\hat{y}^t - y^t)^2 X^2), \end{aligned}$$

where  $s = \mathbf{r} \cdot \mathbf{x}^t$ . For fixed  $y$  and  $\hat{y}$ , the lower bound can be maximized for (omitting  $t$ )

$$s = y - \frac{1}{Xb} \log\left(\frac{1 - \eta(\hat{y} - y)X}{1 + \eta(\hat{y} - y)X}\right).$$

Thus, plugging back for  $s$  and introducing the new variable  $\delta = y - \hat{y}$ , it suffices to show that the function  $G(a, b, \eta, \delta, \hat{y}) \leq 0$  for the choices of  $a, b$ , and  $\eta$  in the claim, and for all  $-X \leq \hat{y} \leq X$ , where

$$\begin{aligned} G(a, b, \eta, \delta, \hat{y}) &:= \log(1 + 2\eta\delta\hat{y} + \eta^2\delta^2X^2) + \frac{1}{4X^2b} \log\left(\frac{1 + \eta\delta X}{1 - \eta\delta X}\right)^2 \\ &\quad - \frac{1}{X}(X + \delta + \hat{y}) \log\left(\frac{1 + \eta\delta X}{1 - \eta\delta X}\right) - 2 \log(1 - \eta\delta X) + a\delta^2. \end{aligned}$$

For a fixed  $\delta$ , this function is maximized at

$$\hat{y}_{\text{opt}} = \frac{X}{\log\left(\frac{1+\eta\delta X}{1-\eta\delta X}\right)} - \frac{1}{2}\eta\delta X^2 \quad \text{since} \quad \frac{\partial^2}{\partial \hat{y}^2} G(a, b, \eta, \delta, \hat{y}) \Big|_{\hat{y}=\hat{y}_{\text{opt}}} = -\frac{\log\left(\frac{1+\eta\delta X}{1-\eta\delta X}\right)^2}{X}.$$

Substituting this value for  $\hat{y}$ , we need to show  $H(a, b, \eta) := G(a, b, \eta, \delta, \hat{y}_{\text{opt}}) \leq 0$  where

$$H(a, b, \eta) = \log\left(\frac{1 + \eta\delta X}{1 - \eta\delta X}\right) \left(\frac{1}{4X^2b} \log\left(\frac{1 + \eta\delta X}{1 - \eta\delta X}\right) + \frac{\eta\delta X}{2} - \frac{\delta}{X} + \frac{1}{2\eta\delta X}\right) \\ + \log\left(\frac{\eta\delta X}{\log\left(\frac{1 + \eta\delta X}{1 - \eta\delta X}\right)}\right) + a\delta^2 - \log(1 - \eta^2\delta^2 X^2) + \log(2) - 1.$$

This function can be simplified further by substituting  $a = \eta = b/(1 + 2X^2b)$ . Also, we can introduce two new variables by defining  $b = c/X^2$  for some  $c > 0$  and  $\delta = pX$  for  $-2 \leq p \leq 2$ . As final step of the proof, there remains to show that the function

$$K(c, p) := \log\left(\frac{1 + (2 + p)c}{1 + (2 - p)c}\right) \left(\frac{1}{4c} \log\left(\frac{1 + (2 + p)c}{1 + (2 - p)c}\right) + \frac{pc}{2(2c + 1)} - p + \frac{2c + 1}{2pc}\right) \\ + \log\left(\frac{pc}{(2c + 1) \log\left(\frac{1 + (2 + p)c}{1 + (2 - p)c}\right)}\right) + \frac{p^2c}{2c + 1} - \log((1 + 2c)^2 - p^2) + \log(2) - 1 \leq 0,$$

for all values of  $c > 0$  and  $-2 \leq p \leq 2$ .

#### Appendix D. Behavior of GD and Reparameterized EGU with Different Initializations

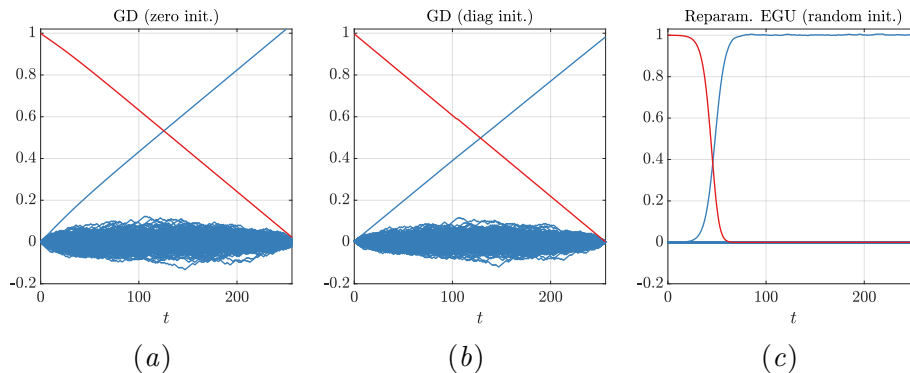


Figure 4: Results of two-layer linear networks on the online Hadamard problem ( $n=128$ ). We perform one pass over the examples. The product of the weights of the first and the second layer (128 weights in total) are shown in blue and the average loss over all 128 examples is shown in red. Results of the two-layer network of Figure 2 using GD where the first layer weights are (a) initialized to zero, (b) initialized uniformly on both diagonals, (c) results of the sparse network of Figure 1 where the weights are initialized randomly.

Here, we discuss two interesting observations on the two-layer networks in Figure 1 and Figure 2. We extend the experiments on learning a column of the Hadamard matrix to two-layer networks. We consider the online Hadamard problem ( $n=128$ ) where we perform one pass over the examples. The results are shown in Figure 4. The figures stress two points about training two-layer linear networks with GD. First, if all the weights in the first layer are initialized to zero (Figure 4(a)), then the product of the weight matrix of the

first layer times the weight vector of the second layer behaves qualitatively the same as the weights of a single linear neuron trained with GD (compare to Figure 3(a)). The same is true if the dotted weights in Figure 2 are initialized to zero and the remaining weights (solid diagonal weights of the first layer and the weights of the second layer) are all set to uniform (Figure 4(b)). This is essentially the same initialization used for the reparameterized EGU algorithm, except in this case, the dotted connections are not removed.

Second, the initialization of the sparse linear network (Figure 1) is rather robust experimentally. Equal initializations or random initializations of both diagonals all make the combined weight behave qualitatively like EGU on a single neuron (Figure 4(c) compared to Figure 3(b)). In (Arora et al., 2018) it was observed that multiplying the weight vector  $\mathbf{w}$  of a linear neuron by a scalar weight  $w'$  and training both with GD behaves similarly to a momentum update. We experimentally show that even when the weights of the sparse network are initialized randomly, then training it with GD shows a drastically different dynamic than training a single neuron with GD. Notice that in the sparse network, each weight on the bottom layer is now multiplied by a separate scalar weight (on the top layer), instead of one common weight as in (Arora et al., 2018).

Finally, note that if we initialize any consecutive weights (along the path from the input to the output) of the sparse linear network in Figure 1 to equal values, then these weights will remain equal throughout the training using GD (since the two edges will have equal gradients at any time). Thus, keeping both weights equal during training does not require any additional constraints on the network.