

Inductive Bias of Gradient Descent for Weight Normalized Smooth Homogeneous Neural Nets

Depen Morwani

Harish G. Ramaswamy

*Department of Computer Science and Engineering, and RBCDSAI,
Indian Institute of Technology Madras, India.*

DEPENMORWANI@GMAIL.COM

HARIGURU@CSE.IITM.AC.IN

Editors: Sanjoy Dasgupta and Nika Haghtalab

Abstract

We analyze the inductive bias of gradient descent for weight normalized smooth homogeneous neural nets, when trained on exponential or cross-entropy loss. We analyse both standard weight normalization (SWN) and exponential weight normalization (EWN), and show that the gradient flow path with EWN is equivalent to gradient flow on standard networks with an adaptive learning rate. We extend these results to gradient descent, and establish asymptotic relations between weights and gradients for both SWN and EWN. We also show that EWN causes weights to be updated in a way that prefers asymptotic relative sparsity. For EWN, we provide a finite-time convergence rate of the loss with gradient flow and a tight asymptotic convergence rate with gradient descent. We demonstrate our results for SWN and EWN on synthetic data sets. Experimental results on simple datasets support our claim on sparse EWN solutions, even with SGD. This demonstrates its potential applications in learning neural networks amenable to pruning.

Keywords: Deep Learning Theory, Inductive Bias, Gradient Descent, Weight Normalization

1. Introduction

The prevailing hypothesis for explaining the generalization ability of deep neural nets, despite their ability to fit even random labels (Zhang et al., 2017), is that the optimisation/training algorithms such as gradient descent have a ‘bias’ towards ‘simple’ solutions. This property is often called inductive bias, and has been an active research area over the past few years.

It has been shown that gradient descent does indeed seem to prefer ‘simpler’ solutions over more ‘complex’ solutions, where the notion of complexity is often problem/architecture specific. The predominant line of work typically shows that gradient descent prefers a least norm solution in some variant of the L_2 -norm. This is satisfying, as gradient descent over the parameters abides by the rules of L_2 geometry, i.e. the weight vector moves along direction of steepest descent, with length measured using the Euclidean norm. However, there is nothing special about the Euclidean norm in the parameter space, and hence several other notions of ‘length’ and ‘steepness’ are equally valid. In recent years, several alternative parameterizations of the weight vector, such as Batch normalization and Weight normalization, have seen immense success and these do not seem to respect L_2 geometry in the ‘weight space’. We pose the question of inductive bias of gradient descent for some of these parameterizations, and demonstrate interesting inductive biases. In particular, it can still be argued that gradient descent with these reparameterizations prefers simpler solutions, but the notion of complexity is different.

Our Contributions. The main contributions of the paper are as follows.

- We establish that the gradient flow path with exponential weight normalization is equal to the gradient flow path of an unnormalized network using an adaptive neuron dependent learning rate. This provides a crisp description of the difference between exponential weight normalized networks and unnormalized networks.
- While most of the previous works on the inductive bias of non-linear deep learning architectures work under the assumption of directional convergence of weights and gradients, we show that gradient convergence implies weight convergence, even for gradient descent, for both standard and exponentially weight normalized network (which is not homogeneous in its parameters)¹.
- We establish the asymptotic relations between weights and gradients for gradient descent on standard weight normalized and exponentially weight normalized networks and show that exponential weight normalization is likely to lead to asymptotic sparsity in weights. We demonstrate the relative sparsity of exponential weight normalization on MNIST dataset, by showing that it leads to networks with better pruning efficacy.
- We establish finite-time convergence rates for gradient flow and tight asymptotic convergence rates for gradient descent on exponentially weight normalized networks.

2. Related Work

The literature most closely related to this paper can be broadly classified into two categories - the inductive biases established for neural networks, and the theoretical studies of normalization methods in deep learning.

2.1. Inductive Bias

[Soudry et al. \(2018\)](#) showed that gradient descent(GD) on the logistic loss with linearly separable data converges to the L_2 maximum margin solution for almost all datasets. These results were extended to loss functions with super-polynomial tails in [Nacson et al. \(2019b\)](#). [Nacson et al. \(2019c\)](#) extended these results to hold for stochastic gradient descent(SGD) and [Gunasekar et al. \(2018a\)](#) extended the results for other optimization geometries. [Ji and Telgarsky \(2019b\)](#) provided tight convergence bounds in terms of dataset size as well as training time. [Ji and Telgarsky \(2019a\)](#) provide similar results when the data is not linearly separable.

[Ji and Telgarsky \(2019c\)](#) showed that for deep linear nets, under certain conditions on the initialization, for almost all linearly separable datasets, the network, in function space, converges to the maximum margin solution. [Gunasekar et al. \(2018b\)](#) established that for linear convolutional nets, under certain assumptions regarding convergence of gradients etc, the function converges to a KKT point of the maximum margin problem in fourier space. [Nacson et al. \(2019a\)](#) shows that for smooth homogeneous nets, the network converges to a KKT point of the maximum margin problem in parameter space. [Lyu and Li \(2020\)](#) established these results with weaker assumptions and also provide asymptotic convergence rates for the loss. [Chizat and Bach \(2020\)](#) explore the inductive bias for a 2-layer infinitely wide ReLU neural net in function space and show that the function learnt is a

1. The assumptions about weights and gradients converging in direction have been recently shown to hold for gradient flow on homogeneous neural nets without normalization under some regularity conditions related to o-minimality of the architecture([Ji and Telgarsky, 2020](#))

max-margin classifier for variation norm. [Moroshko et al. \(2020\)](#) established the inductive bias for linear diagonal networks and showed that the network transitions between max-margin and L_1 -norm margin, depending on the relation between the initialization scale and training accuracy.

2.2. Normalization

[Salimans and Kingma \(2016\)](#) introduced weight normalization and demonstrated that it replicates the convergence speedup of BatchNorm. Similarly, other normalization techniques have been proposed as well ([Ba et al., 2016](#)) ([Qiao et al., 2020](#)) ([Li et al., 2019](#)), but only a few have been theoretically explored. [Santurkar et al. \(2018\)](#) demonstrated that batch normalization makes the loss surface smoother and L_2 normalization in batchnorm can even be replaced by L_1 normalizations. [Kohler et al. \(2019\)](#) showed that for GD, batchnorm speeds up convergence in the case of GLM by splitting the optimization problem into learning the direction and the norm. [Cai et al. \(2019\)](#) analyzed GD on BN for squared loss and showed that it converges for a wide range of η . [Bjorck et al. \(2018\)](#) showed that the primary reason BN allows networks to achieve higher accuracy is by enabling higher learning rates. [Arora et al. \(2019\)](#) showed that in case of GD or SGD with batchnorm, η for scale-invariant parameters does not affect the convergence rate towards stationary points. [Du et al. \(2018\)](#) showed that for GD over one-hidden-layer weight normalized CNN, with a constant probability over initialization, iterates converge to global minima. [Qiao et al. \(2019\)](#) compared different normalization techniques from the perspective of whether they lead to points, where neurons are consistently deactivated. [Wu et al. \(2020\)](#) established the inductive bias of gradient flow with weight normalization for overparameterized least squares and showed that for a wider range of initializations as compared to normal parameterization, it converges to the minimum norm solution. [Dukler et al. \(2020\)](#) analyzed weight normalization for multilayer ReLU net in the infinite width regime and showed that it may speedup convergence. Some other papers ([Luo et al., 2019](#); [Roburin et al., 2020](#)) also provide other perspectives to think about normalization techniques.

3. Problem Setup

We use a standard view of neural networks as a collection of nodes/neurons grouped by layers. Each node u is associated with a weight vector w_u , that represents the incoming weight vector for that node. In case of CNNs, weights can be shared across different nodes. w represents all the parameters of the network arranged in form of a vector (In general, for any vector v associated with the entire network, v_u represent its components corresponding to the node u). The training dataset consists of $(x_i; y_i)$ pairs with a total of m points in the dataset. The function represented by the neural network is denoted by $f(w; \cdot)$. The loss for a single data point is given by $\ell(y_i; f(w; x_i))$ and the loss vector is represented by $\ell(w)$ and is given by $\ell(w) = \sum_{i=1}^m \ell(y_i; f(w; x_i))$. We sometimes abbreviate $\ell(w(t))$ as L when the context is clear.

In standard weight normalisation (SWN), each weight vector w_u is reparameterized as $\frac{v_u}{\|v_u\|_k}$. This was proposed by [Salimans and Kingma \(2016\)](#), as a substitute for Batch Normalization and has been practically used in multiple papers such as [Sokolic et al. \(2017\)](#), [Dauphin et al. \(2017\)](#), [Kim et al. \(2018\)](#) and [Hieber et al. \(2018\)](#). The corresponding update equations for gradient descent

are given by

$$w_u(t+1) = w_u(t) - \eta \frac{v_u(t) \langle w_u(t), v_u(t) \rangle}{\|v_u(t)\|^2} r_{w_u} L \quad (1)$$

$$v_u(t+1) = v_u(t) - \eta \frac{v_u(t)}{\|v_u(t)\|^2} \left(\langle w_u(t), v_u(t) \rangle - \frac{v_u(t) \langle w_u(t), v_u(t) \rangle}{\|v_u(t)\|^2} \right) r_{w_u} L \quad (2)$$

In exponential weight normalization (EWN), each weight vector is reparameterized as $w_u = \frac{v_u}{\|v_u\|}$. This was mentioned in [Salimans and Kingma \(2016\)](#), but to the best of our knowledge, has not been widely used. The corresponding update equations for gradient descent with learning rate η are given by

$$w_u(t+1) = w_u(t) - \eta e^{-\langle w_u(t), v_u(t) \rangle} \frac{v_u(t) \langle w_u(t), v_u(t) \rangle}{\|v_u(t)\|^2} r_{w_u} L \quad (3)$$

$$v_u(t+1) = v_u(t) - \eta \frac{e^{-\langle w_u(t), v_u(t) \rangle}}{\|v_u(t)\|^2} \left(\langle w_u(t), v_u(t) \rangle - \frac{v_u(t) \langle w_u(t), v_u(t) \rangle}{\|v_u(t)\|^2} \right) r_{w_u} L \quad (4)$$

The update equations for gradient flow are the continuous counterparts for the same. In gradient flow, for both SWN and EWN, we set $\|v_u(0)\| = 1$, to simplify the update equations.

4. Inductive Bias of Weight Normalization

In this section, we state our main results for weight normalized smooth homogeneous models on exponential loss $\ell(y_i; (w; x_i)) = e^{-y_i \langle w; x_i \rangle}$. The results for cross-entropy loss and proofs have been deferred to the appendix due to space constraints. First, we state the main proposition that helps in establishing these results for EWN.

Theorem 1 The gradient flow path with learning rate $\eta(t)$ for EWN and SWN are given as follows:

$$\text{EWN: } \frac{dw_u(t)}{dt} = -\eta(t) \frac{v_u(t) \langle w_u(t), v_u(t) \rangle}{\|v_u(t)\|^2} r_{w_u} L \quad (5)$$

$$\text{SWN: } \frac{dw_u(t)}{dt} = -\eta(t) \left(\frac{v_u(t) \langle w_u(t), v_u(t) \rangle}{\|v_u(t)\|^2} r_{w_u} L + \frac{1 - \langle w_u(t), v_u(t) \rangle^2}{\|v_u(t)\|^2} (w_u(t) \langle w_u(t), v_u(t) \rangle) r_{w_u} L \right) \quad (6)$$

Thus, the gradient flow path of EWN can be replicated by an adaptive (neuron dependent) learning rate given by $\eta(t) \frac{v_u(t) \langle w_u(t), v_u(t) \rangle}{\|v_u(t)\|^2}$ on the unnormalized network (Unnorm).

4.1. Assumptions

The assumptions in the paper can be broadly divided into loss function/architecture based assumptions and trajectory based assumptions. The loss functions/architecture based assumptions are shared across both gradient flow and gradient descent.

Loss function/Architecture based assumptions

- $\ell(y_i; (w; x_i)) = e^{-y_i \langle w; x_i \rangle}$

- $\ell(\cdot; x)$ is a C^1 function (i.e. continuously differentiable), for any $x \in \mathbb{R}^d$

$$\exists (\epsilon; \delta) = \epsilon^L (\delta; \delta), \text{ for some } \epsilon > 0 \text{ and } L > 0$$

Gradient Descent. For gradient descent, we make the following trajectory based assumptions

(A1) There exists a time t_0 such that $L(w(t_0)) < 1$.

$$(A2) \lim_{t \rightarrow \infty} \frac{r_{w(t)} L(w(t))}{k_{w(t)}} := \epsilon.$$

The first trajectory assumption is simply a separability assumption and means that the network is able to correctly classify the dataset at some point during the training process. This is not a completely unreasonable assumption, given recent papers demonstrating neural networks with sufficient overparameterization can fit even random labels (Zhang et al., 2017; Jacot et al., 2018). The second assumption has been used in multiple previous works (Gunasekar et al., 2018b; Chizat and Bach, 2020; Nacson et al., 2019a), and is standard in the literature related to the inductive bias of non-linear deep learning architectures. Moreover, we remove one of the assumptions related to directional convergence of weights used in these works, and instead show that it is implied by the directional convergence of gradients.

Gradient Descent For gradient descent, we require the learning rate to not grow too fast, and a slightly stronger assumption on loss.

$$(B1) \lim_{t \rightarrow \infty} L(w(t)) = 0$$

$$(B2) \lim_{t \rightarrow \infty} \frac{r_{w(t)} L(w(t))}{k_{w(t)}} := \epsilon$$

$$(B3) \lim_{t \rightarrow \infty} (t) k_{w_u(t)} k_{w_u} L(w(t)) = 0 \text{ for all } u \text{ in the network.}$$

The assumption (B3) is mild, as the norm of the gradient of the exponential loss goes down exponentially fast as compared to norm of the weights. We demonstrate that these assumptions hold for multiple datasets including MNIST in Appendix P.

4.2. Asymptotic relations between weights and gradients

This section contains the main theorems that establish asymptotic relations between weights and gradients for SWN and EWN. First, we will state a common proposition for both SWN and EWN.

Proposition 2 Under assumption (A1) for gradient descent, for both SWN and EWN, $\lim_{t \rightarrow \infty} L(w(t)) = 0$.

Although the above proposition was established for homogeneous nets by Lyu and Li (2020) we extend it for the non-homogeneous parameterization of EWN. Now, we provide one of our main theorem that establishes gradient convergence implies weight convergence.

Theorem 3 Consider a node u in the network with $k_{w_u} > 0$ and $\lim_{t \rightarrow \infty} k_{w_u(t)} = 1$. Under assumptions (A1), (A2) for gradient descent and (B1)-(B3) for gradient descent, for both SWN and EWN

$$(i) \lim_{t \rightarrow \infty} \frac{w_u(t)}{k_{w_u(t)}} := \mathbf{w}_u \text{ exists.}$$

$$(ii) \mathbf{w}_u = \epsilon_u \text{ for some } \epsilon > 0.$$

The above theorem relaxes one of the assumptions regarding weight convergence used in many of the previous works, by showing that even for non-homogeneous parameterization under gradient descent, gradient convergence implies weight convergence. Moreover, it also shows, that weights and gradients eventually get aligned opposite to each other.

2. Homogeneous networks in the space are also homogeneous in the space. Therefore results regarding convergence rates and monotonic margin hold from Lyu and Li (2020). However, the results for convergence to a KKT point of the max margin problem do not hold. For details, refer Appendix L

(a) (b) (c)

Figure 1: Demonstration of Results for EWN in Lin-Sep experiment: (a) Evolution of $\|w_u\|$ - norm of the incoming weights for neuron (b) Cosine between weights and gradients for neurons 5, 7 and 8. (c) Weight and gradient norms for weights 5, 7 and 8.

Now, we provide the main theorem that distinguishes SWN and EWN.

Theorem 4 Consider two nodes u and v in the network with $\|g_u\|, \|g_v\| > 0$; $\lim_{t \rightarrow \infty} \|w_u(t)\| = 1$ and $\lim_{t \rightarrow \infty} \|w_v(t)\| = 1$. Let $\frac{\|g_u\|}{\|g_v\|}$ be denoted by c . Under assumptions (A1), (A2) for gradient flow and (B1)-(B3) for gradient descent,

- (i) for SWN, $\lim_{t \rightarrow \infty} \frac{\|w_u(t)\|}{\|w_v(t)\|} = c$
- (ii) for EWN, $\lim_{t \rightarrow \infty} \frac{\|w_u(t)\|}{\|w_v(t)\|}$ is either 0, 1 or $\frac{1}{c}$

Thus, if $\frac{\|w_u(t)\|}{\|w_v(t)\|}$ converges to a finite non-zero value for EWN, then $\|w_u(t)\| \|r_{w_u} L(w(t))\| = \|w_v(t)\| \|r_{w_v} L(w(t))\| = k_1(t)$ asymptotically. While, for SWN $\frac{\|w_u(t)\|}{\|w_v(t)\|} = \frac{\|w_v(t)\|}{\|w_u(t)\|} = k_2(t)$ asymptotically, where $k_1(t)$ and $k_2(t)$ are independent of u and v . The exact conditions under which $\frac{\|w_u(t)\|}{\|w_v(t)\|}$ tends to 0 or 1 for EWN are provided in Proposition 5.

We demonstrate Theorem 3 and Theorem 4 for EWN on a linearly separable dataset in Figure 1. In this experiment, a 2-layered neural network, with 8 neurons in the hidden layer and a ReLU-squared activation function, is trained on a linearly separable dataset. The learning rate schedule used was $\frac{1}{\sqrt{0.97}}$ and the network was trained till a loss of 10^{-300} . As can be seen in Figure 1, for weights 5, 7 and 8, whose norms keep on growing, weights and gradients eventually become oppositely aligned, and their norms are inversely proportional to each other. The results for SWN, along with results on other datasets including MNIST have been deferred to Appendix N.

4.3. Sparsity Inductive Bias for Exponential Weight Normalisation

The inverse relation between $\|w_u(t)\|$ and $\|r_{w_u} L(w(t))\|$ in the EWN trajectory results in an interesting inductive bias that favours movement along sparse directions.

Proposition 5 Consider two nodes u and v in the network such that $\|g_v\|, \|g_u\| > 0$ and $\|w_u(t)\|, \|w_v(t)\| \neq 1$. Let $\frac{\|g_u\|}{\|g_v\|}$ be denoted by c (note that c and c will be different for SWN and EWN trajectory). Consider any ϵ such that $0 < \epsilon < c$ and $0 < \delta < \frac{1}{2}$. Then, the following holds:

(a) Network Architecture (b) Weight trajectories

Figure 2: (a) Network architecture for the Simple-Traj experiment. (b) Trajectories of the two weights for EWN and Unnorm, starting from 5 different initialization points.

(i) There exists a time t_1 , such that for all $t > t_1$ both SWN and EWN trajectories have the following properties:

$$(a) \frac{kr_{w_u} L(w(t))k}{kr_{w_v} L(w(t))k} \in [c^-; c^+] \quad (b) \frac{kw_u(t)k}{kw_v(t)k} > \frac{r_{w_u} L(w(t))}{r_{w_v} L(w(t))} \cos(\theta)$$

$$(c) \frac{kw_v(t)k}{kw_u(t)k} > \frac{r_{w_v} L(w(t))}{r_{w_u} L(w(t))} \cos(\theta).$$

(ii) for SWN, $\lim_{t \rightarrow \infty} \frac{kw_u(t)k}{kw_v(t)k} = c$

(iii) for EWN, if at some time $t_2 > t_1$,

$$(a) \frac{kw_u(t_2)k}{kw_v(t_2)k} > \frac{1}{(c^-) \cos(\theta)} \Rightarrow \lim_{t \rightarrow \infty} \frac{kw_u(t)k}{kw_v(t)k} = 1$$

$$(b) \frac{kw_u(t_2)k}{kw_v(t_2)k} < \frac{\cos(\theta)}{c^+} \Rightarrow \lim_{t \rightarrow \infty} \frac{kw_u(t)k}{kw_v(t)k} = 0$$

The above proposition shows that the limit property of the weights in Theorem 4, makes non-sparse w an unstable convergent direction for EWN. But that is not the case for SWN. We demonstrate the relative sparsity between EWN, SWN and Unnorm through two toy experiments Simple-Traj and XOR

In the Simple-Traj experiment, we illustrate the notion of asymptotic relative sparsity occurring in the EWN parameterization. In this dataset, we have a single data point $(0, 1)$ that is labelled positive and train a network with linear activations. The architecture is shown in Figure 2, where weights in blue and red are frozen to values 0 and 1 respectively. Thus, there are effectively only two scalar parameters w_1 and w_2 . The network is trained till a loss value of 10^{-50} starting from 5 different initialization points. The weight trajectories in Figure 2 shows that EWN prefers to converge either along the x or y axis, and hence has an asymptotic relative sparsity property. We provide a theoretical proof for the same in the general d-dimensional case.

Proposition 6 Consider a linear model over \mathbb{R}^d given by $f(x) = w^T x$, where each w_i is further reparameterized as $w_i = \alpha_i \cdot e_i$. Consider a dataset consisting of a single data point $(0, 1)$, that is labelled

(a) EWN (b) SWN (c) Unnorm

Figure 3: (a), (b) and (c) demonstrate the evolution of weight norms for each neuron in the XOR experiment. EWN weights grow sparsely when compared to Unnorm and SWN

(a) EWN (b) SWN (c) Unnorm

Figure 4: Variation of convergence rate of train loss with number of layers for multilayer linear nets on a linearly separable dataset

as $+1$. According to the initialization of w , define a relation R on $\{1, \dots, d\}$, given by $i \sim j$ if $w_i(0)z_i = w_j(0)z_j$. Then, R is an equivalence relation on $\{1, \dots, d\}$. Let these equivalent sets be denoted by I_1, I_2, \dots, I_k . Define a total order on these sets given by $I_a > I_b$ if $\exists i \in I_a, j \in I_b$ such that $w_i(0)z_i > w_j(0)z_j$. Let the maximum set according to this order be denoted by I_{max} . Then, for gradient flow on exponential loss, the following holds

- (i) For any $i \in I_{max}$, $\lim_{t \rightarrow \infty} w_i(t) = 1$
- (ii) For $i, j \in I_{max}$, $\frac{w_i(t)}{w_j(t)} = \frac{x_j}{x_i}$.
- (iii) For any $i \notin I_{max}$, $\lim_{t \rightarrow \infty} w_i(t) = \frac{1}{w_i(0)} \frac{x_i}{w_j(0)x_j}$, where j is any element in I_{max} .

Thus, if w is initialized from a continuous distribution, then with probability 1, the cardinality of I_{max} is 1 and hence $\|w(t)\|_k$ will approach a sparse vector.

In the XOR experiment, we train a 2-layer ReLU network with 20 hidden neurons on XOR dataset, till a loss value of 10^{-50} . The second layer is fixed to the values 1 or -1 randomly. For attaining 100% accuracy on this dataset with this architecture, at least 4 hidden units are needed. As

can be seen in Figure 3, EWN asymptotically uses exactly 4 neurons out of 20, while Unnorm and SWN use almost all the 20 neurons.

5. Convergence Rates

In this section, we provide convergence rate of loss for EWN.

Gradient Flow: We provide a finite-time convergence rate of loss for gradient flow in case of EWN.

Theorem 7 For Exponential Weight Normalization, under assumption (A1), the following hold for $t > t_0$ in case of gradient flow

- (i) $\|w(t)\|_k$ grows with $\text{as } O((\log t)^{\frac{1}{L}})$
- (ii) $L(t)$ goes down with $\text{as } O(\frac{1}{t})$

Gradient Descent: For establishing convergence rates for gradient descent, we are going to make an additional assumption that the overall weight vector converges in direction, i.e., $\frac{w(t)}{\|w(t)\|_k}$ exists (B4). Although we have already shown this is indeed true for nodes with $\|w_i(t)\|_k > 0$, we need this assumption to take into account the nodes with $\|w_i(t)\|_k = 0$. Under this assumption, w can be represented as $w = g(t)v + r(t)$, where $\lim_{t \rightarrow \infty} \frac{\|r(t)\|_k}{g(t)} = 0$. Let $d : \mathbb{N} \rightarrow \mathbb{R}$, given by $d(t) = \sum_{i=0}^{t-1} \eta$ denote total step size. Let $\gamma = \min_i y_i (v; x_i)$ be the normalized margin at convergence.

The asymptotic convergence rate of loss for SWN and Unnorm have already been established in Lyu and Li (2020) as $\frac{1}{d(t)(\log d(t))^{\frac{2}{L}}}$. For EWN, the corresponding theorem is provided below

Theorem 8 For Exponential Weight Normalization, under Assumptions (B1)-(B4), $\lim_{t \rightarrow \infty} \frac{\|r(t)\|_k}{g(t)} = 0$, $\lim_{t \rightarrow \infty} \frac{\|r(t+1)\|_k}{g(t+1)} - \frac{\|r(t)\|_k}{g(t)} = 0$, the following hold

- (i) $\|w(t)\|_k$ asymptotically grows with $\text{as } O((\log d(t))^{\frac{1}{L}})$
- (ii) $L(w(t))$ asymptotically goes down with $\text{as } O(\frac{1}{d(t)(\log d(t))^2})$.

Although the additional assumption $\lim_{t \rightarrow \infty} \frac{\|r(t+1)\|_k}{g(t+1)} - \frac{\|r(t)\|_k}{g(t)} = 0$ is not standard, we empirically demonstrate that, for EWN, the convergence rate is almost independent of the number of layers. Moreover, the learning rate assumption used still covers the constant case, that is generally used in practice.

For multilayer linear nets, the variation of convergence rate with number of layers for a linearly separable dataset is illustrated in Figure 4. All of these networks were explicitly initialized to represent the same point in function space. It can be seen that EWN, SWN and unnormalized networks all converge faster with more layers, but the effect is much less pronounced for EWN.

6. Pruning Experiments

As EWN leads to asymptotically sparse solutions, it is likely that a sufficiently trained EWN network would be comparatively robust to pruning. In this section, we compare the pruning efficacy of

(a) $L = e^{-10}$ (b) $L = e^{-100}$ (c) $L = e^{-300}$

Figure 5: Variation of test accuracy vs percentage of neurons pruned in first layer at different loss values for MNIST experiment

EWN, SWN and Unnorm on MNIST (LeCun et al., 2010) dataset. We use a 2-layer ReLU network with 1024 neurons in the hidden layer. In case of EWN and SWN, only the first layer is weight normalized as only this layer needs to be pruned.

Pruning Strategy: The natural pruning strategy of removing neurons with least $\|w_u\|_k$ gives inordinate importance to the initialisation and the initial optimization epochs. In order to minimize the effect of initialization and initial movement of weights, we prune according to the weight norm increase from a reference point. For example, when pruning at a loss value $L = e^{-300}$, we consider 4 reference points: 0, weight at initialization, weight when $L = e^{-10}$ and $L = e^{-100}$. We then choose the pruning strategy that gives maximum testing accuracy for a given level of pruning. Similarly, for pruning at a loss value $L = e^{-100}$, we consider three reference points: 0, weight at initialization and weight when $L = e^{-10}$, and for pruning at a loss value $L = e^{-10}$, we consider 0, weight at initialization and weight when $L = e^{-5}$. More detailed description of the pruning strategy is provided in Appendix Q.

The pruning graphs for MNIST at different loss values averaged across multiple seeds are shown in Figure 5. It can be seen that when the loss levels are sufficiently low, the EWN network becomes better adapted for pruning, significantly outperforming SWN and the unnormalized network in terms of test accuracy for a given level of pruning. Further details along with convergence rate plots are provided in Appendix O.

6.1. Combining EWN with other sparsity regularizers

EWN by itself has a sparsifying effect, and in addition it can be combined with other sparsity regularizers. e.g. the neural net can be trained using any existing sparsity regularizer with standard parameterization for the initial few epochs, and use EWN (without the regularizer) for the later phase.

We conducted experiment on MNIST dataset, by initially training the network till convergence with $\ell_{2,1}$ group sparsity regularizer ($\sum_{u \in \text{nodes of network}} \|w_u\|_k$), and later on using EWN till a loss value of e^{-30} . In this case, the pruning strategy is based on the ℓ_2 norms of the weights as the sparsity regularizer already induces an asymmetry among different neurons. The pruning results are

(a) (b)

Figure 6: (a) Variation of test accuracy vs percentage of neurons pruned on MNIST dataset after training the network initially with $\ell_{2,1}$ regularizer and later on continuing training with either EWN, SWN or unnormalized net $\ell_{1,1} = e^{-30}$ (b) Zoomed in view on EWN

shown in Figure 6. As can be seen, training further with EWN improves the pruning efficacy of the network.

7. Proof Sketch

In this section, we provide the proof for part (iii)a of Proposition 5 for gradient flow, i.e, for EWN, $\frac{\|w_u(t_2)\|}{\|w_v(t_2)\|} > \frac{1}{c \cos(\cdot)} \Rightarrow \lim_{t \rightarrow \infty} \frac{\|w_u(t)\|}{\|w_v(t)\|} = 1$. Remaining proofs are given in appendix.

The update equations for u and v_u in case of gradient flow for EWN are given by

$$\frac{du(t)}{dt} = -e^{-u(t)} \frac{v_u(t) \langle r_{w_u} L \rangle}{\|v_u(t)\|} \quad (7)$$

$$\frac{dv_u(t)}{dt} = -e^{-u(t)} \left(1 - \frac{v_u(t) \langle v_u(t) \rangle}{\|v_u(t)\|^2} \right) \langle r_{w_u} L \rangle \quad (8)$$

Using Equation (7) (along with the fact that $\|w_u(t)\|$ does not change with time and $\|w_u(0)\| = 1$),

$$\frac{d\|w_u(t)\|}{dt} = \frac{de^{-u(t)}}{dt} = -e^{-u(t)} \langle r_{w_u} L \rangle \quad (9)$$

Using Equation (9) and part 1 of Proposition 5, we can say for t_1 ,

$$\begin{aligned} \frac{d \frac{\|w_u(t)\|}{\|w_v(t)\|}}{dt} &= \frac{\|w_v(t)\| \frac{d\|w_u(t)\|}{dt} - \|w_u(t)\| \frac{d\|w_v(t)\|}{dt}}{\|w_v(t)\|^2} \\ &= \frac{e^{-u(t)} \langle r_{w_u} L \rangle (\|w_u(t)\| \cos(\cdot) - \|w_v(t)\| \langle r_{w_u} L \rangle)}{\|w_v(t)\|^2} \\ &= \frac{e^{-u(t)} \langle r_{w_u} L \rangle}{\|w_v(t)\|} \cos(\cdot) - \frac{1}{c} \end{aligned} \quad (10)$$

In this case, using Equation (10), we can see $\frac{d \frac{\|w_u(t)\|_k}{\|w_v(t)\|_k}}{dt} > 0$ at t_2 . Thus, $\frac{\|w_u(t)\|_k}{\|w_v(t)\|_k}$ always remains greater than $\frac{1}{\cos(\theta)}$ and keeps on increasing. Let's denote $\frac{\|w_u(t_2)\|_k}{\|w_v(t_2)\|_k}$ by β . Then, for $t > t_2$,

$$\frac{d \frac{\|w_u(t)\|_k}{\|w_v(t)\|_k}}{dt} = \cos(\theta) \frac{1}{c} \int_{w_u} L(w(t)) dw$$

As $\beta \rightarrow 1$, $\int_{t_2}^{R_1} \int_{w_u} L(w(t)) dw dt = 1$ using Equation (7). Thus, integrating both the sides of the equation above from t_2 to 1 , we get

$$\int_{t_2}^1 \frac{d \frac{\|w_u(t)\|_k}{\|w_v(t)\|_k}}{dt} dt = 1$$

Thus $\lim_{t \rightarrow 1} \frac{\|w_u(t)\|_k}{\|w_v(t)\|_k} = 1$. A similar proof works for the part (iii)b of Proposition 5 as well.

8. Conclusion

In this paper, we analyze the inductive bias of weight normalization for smooth homogeneous neural nets and show that exponential weight normalization is likely to lead to asymptotically sparse solutions and has a faster convergence rate than unnormalized or standard weight normalized networks.

The smooth homogeneity assumptions made in the paper are satisfied by any positive power of ReLU greater than 1. The primary issue with ReLU is that the assumption of gradients converging in direction does not make sense for the non-smooth case. However, as our experiments demonstrate, the implication of relative sparsity holds for ReLU activation as well. Therefore, extending the results in the paper for the non-smooth case is a promising research direction.

Although the trajectory based assumptions have been shown to hold for gradient flow on unnormalized nets under certain regularity conditions, establishing similar conditions for weight normalized networks remains an open question. Moreover, extending the directional convergence results from gradient flow to gradient descent is also an interesting research direction.

References

- Sanjeev Arora, Zhiyuan Li, and Kaifeng Lyu. Theoretical analysis of auto rate-tuning by batch normalization. In International Conference on Learning Representations, 2019. URL <https://openreview.net/forum?id=rkxQ-nA9FX>.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. arXiv preprint arXiv:1607.06450, 2016.
- Nils Bjorck, Carla P Gomes, Bart Selman, and Kilian Q Weinberger. Understanding batch normalization. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, Advances in Neural Information Processing Systems, pages 7694–7705. Curran Associates, Inc., 2018. URL <http://papers.nips.cc/paper/7996-understanding-batch-normalization.pdf>.

- Yongqiang Cai, Qianxiao Li, and Zuwei Shen. A quantitative analysis of the effect of batch normalization on gradient descent. volume 97 of Proceedings of Machine Learning Research pages 882–890, Long Beach, California, USA, 09–15 Jun 2019. PMLR. URL <http://proceedings.mlr.press/v97/cai19a.html>
- Lénác Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss. volume 125 of Proceedings of Machine Learning Research pages 1305–1338. PMLR, 09–12 Jul 2020. URL <http://proceedings.mlr.press/v125/chizat20a.html>
- Yann N. Dauphin, Angela Fan, Michael Auli, and David Grangier. Language modeling with gated convolutional networks. volume 70 of Proceedings of Machine Learning Research pages 933–941, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR. URL <http://proceedings.mlr.press/v70/dauphin17a.html>
- Simon Du, Jason Lee, Yuandong Tian, Aarti Singh, and Barnabas Poczos. Gradient descent learns one-hidden-layer CNN: Don't be afraid of spurious local minima. volume 80 of Proceedings of Machine Learning Research pages 1339–1348, Stockholm, Sweden, Stockholm Sweden, 10–15 Jul 2018. PMLR. URL <http://proceedings.mlr.press/v80/du18b.html>
- Yonatan Dukler, Quanquan Gu, and Guido Marafioti. Optimization theory for relu neural networks trained with normalization layers. International Conference on Machine Learning, 2020.
- Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias in terms of optimization geometry. volume 80 of Proceedings of Machine Learning Research pages 1832–1841, Stockholm, Sweden, Stockholm Sweden, 10–15 Jul 2018a. PMLR. URL <http://proceedings.mlr.press/v80/gunasekar18a.html>
- Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent on linear convolutional networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, Advances in Neural Information Processing Systems, pages 9461–9471. Curran Associates, Inc., 2018b. URL <http://papers.nips.cc/paper/8156-implicit-bias-of-gradient-descent-on-linear-convolutional-networks.pdf>
- Felix Hieber, Tobias Domhan, Michael Denkowski, David Vilar, Artem Sokolov, Ann Clifton, and Matt Post. Sockeye: A toolkit for neural machine translation, 2018.
- Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, Advances in Neural Information Processing Systems, 31 pages 8571–8580. Curran Associates, Inc., 2018. URL <http://papers.nips.cc/paper/8076-neural-tangent-kernel-convergence-and-generalization-in-neural-networks.pdf>
- Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. volume 99 of Proceedings of Machine Learning Research pages 1772–1798, Phoenix, USA, 25–28 Jun 2019a. PMLR. URL <http://proceedings.mlr.press/v99/ji19a.html>

- Ziwei Ji and Matus Telgarsky. A refined primal-dual analysis of the implicit bias, 2019b.
- Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. In International Conference on Learning Representations, 2019c. URL <https://openreview.net/forum?id=HJflg30qKX>.
- Ziwei Ji and Matus Telgarsky. Directional convergence and alignment in deep learning. preprint arXiv:2006.06657, 2020.
- Jin-Hwa Kim, Jaehyun Jun, and Byoung-Tak Zhang. Bilinear attention networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, Advances in Neural Information Processing Systems, volume 31, pages 1564–1574. Curran Associates, Inc., 2018. URL <https://proceedings.neurips.cc/paper/2018/file/96ea64f3a1aa2fd00c72faacf0cb8ac9-Paper.pdf>.
- Jonas Kohler, Hadi Daneshmand, Aurelien Lucchi, Thomas Hofmann, Ming Zhou, and Klaus Neymeyr. Exponential convergence rates for batch normalization: The power of length-direction decoupling in non-convex optimization. volume 89 of Proceedings of Machine Learning Research, pages 806–815. PMLR, 16–18 Apr 2019. URL <http://proceedings.mlr.press/v89/kohler19a.html>.
- Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. AT&T Labs [Online]. Available: <http://yann.lecun.com/exdb/mnist>, 2010.
- Boyi Li, Felix Wu, Kilian Q Weinberger, and Serge Belongie. Positional normalization. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'AchBuc, E. Fox, and R. Garnett, editors, Advances in Neural Information Processing Systems, pages 1622–1634. Curran Associates, Inc., 2019. URL <http://papers.nips.cc/paper/8440-positional-normalization.pdf>.
- Ping Luo, Xinjiang Wang, Wenqi Shao, and Zhanglin Peng. Towards understanding regularization in batch normalization. In International Conference on Learning Representations, 2019. URL <https://openreview.net/forum?id=HJILKjR9FQ>.
- Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks. In International Conference on Learning Representations, 2020. URL <https://openreview.net/forum?id=SJeLIgBKPS>.
- Edward Moroshko, Suriya Gunasekar, Blake Woodworth, Jason D. Lee, Nathan Srebro, and Daniel Soudry. Implicit bias in deep linear classification: Initialization scale vs training accuracy, 2020.
- M. Muresan. A Concrete Approach to Classical Analysis. CMS Books in Mathematics. Springer New York, 2015. ISBN 9780387789330. URL https://books.google.co.in/books?id=N8rBgtlu_qgC.
- Mor Shpigel Nacson, Suriya Gunasekar, Jason Lee, Nathan Srebro, and Daniel Soudry. Lexicographic and depth-sensitive margins in homogeneous and non-homogeneous deep models. volume 97 of Proceedings of Machine Learning Research, pages 4683–4692, Long Beach, California, USA, 09–15 Jun 2019a. PMLR. URL <http://proceedings.mlr.press/v97/nacson19a.html>.

- Mor Shpigel Nacson, Jason Lee, Suriya Gunasekar, Pedro Henrique Pamplona Savarese, Nathan Srebro, and Daniel Soudry. Convergence of gradient descent on separable data. volume 89 of Proceedings of Machine Learning Research, pages 3420–3428. PMLR, 16–18 Apr 2019b. URL <http://proceedings.mlr.press/v89/nacson19b.html>
- Mor Shpigel Nacson, Nathan Srebro, and Daniel Soudry. Stochastic gradient descent on separable data: Exact convergence with a fixed learning rate. volume 89 of Proceedings of Machine Learning Research, pages 3051–3059. PMLR, 16–18 Apr 2019c. URL <http://proceedings.mlr.press/v89/nacson19a.html>
- Siyuan Qiao, Huiyu Wang, Chenxi Liu, Wei Shen, and Alan Yuille. Rethinking normalization and elimination singularity in neural networks. arXiv preprint arXiv:1911.09738, 2019.
- Siyuan Qiao, Huiyu Wang, Chenxi Liu, Wei Shen, and Alan Yuille. Micro-batch training with batch-channel normalization and weight standardization, 2020.
- Simon Roburin, Yann de Mont-Marin, Andrei Bursuc, Renaud Marlet, Patrick P. and Mathieu Aubry. Spherical perspective on learning with batch normalization. arXiv preprint arXiv:2006.13382, 2020.
- Tim Salimans and Durk P Kingma. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, Advances in Neural Information Processing Systems, 29 pages 901–909. Curran Associates, Inc., 2016. URL <http://papers.nips.cc/paper/6114-weight-normalization-a-simple-reparameterization-to-accelerate-training-of-deep-neural-networks.pdf>
- Shibani Santurkar, Dimitris Tsipras, Andrew Ilyas, and Aleksander Madry. How does batch normalization help optimization? In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, Advances in Neural Information Processing Systems, 31 pages 2483–2493. Curran Associates, Inc., 2018. URL <http://papers.nips.cc/paper/7515-how-does-batch-normalization-help-optimization.pdf>
- Jure Sokolic, Raja Giryes, Guillermo Sapiro, and Miguel R. D. Rodrigues. Robust large margin deep neural networks. IEEE Transactions on Signal Processing, 65(16):4265–4280, Aug 2017. ISSN 1941-0476. doi: 10.1109/tsp.2017.2708039. URL <http://dx.doi.org/10.1109/TSP.2017.2708039>
- Daniel Soudry, Elad Hoffer, and Nathan Srebro. The implicit bias of gradient descent on separable data. In International Conference on Learning Representations, 2018. URL <https://openreview.net/forum?id=r1q7n9gAb>
- Xiaoxia Wu, Edgar Dobriban, Tongzheng Ren, Shanshan Wu, Zhiyuan Li, Suriya Gunasekar, Rachel Ward, and Qiang Liu. Implicit regularization and convergence for weight normalization. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, Advances in Neural Information Processing Systems, volume 33, pages 2835–2847. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/1de7d2b90d554be9f0db1c338e80197d-Paper.pdf>

Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. International Conference on Learning Representation, 2017. URL <https://openreview.net/forum?id=Sy8gdB9xx>.

Appendix A. Proof of Theorem 1

Theorem The gradient flow path with learning rate $\eta(t)$ for EWN and SWN are given as follows:

$$\begin{aligned} \text{EWN: } \frac{dw_u(t)}{dt} &= -\eta(t)kw_u(t)k^2r_{w_u}L \\ \text{SWN: } \frac{dw_u(t)}{dt} &= -\eta(t)kw_u(t)k^2r_{w_u}L + \frac{1}{kw_u(t)k^2} \langle w_u(t), r_{w_u}L \rangle w_u(t) \end{aligned}$$

The proof for the two parts will be provided in different subsections, where the corresponding part will be restated for ease of the reader.

A.1. Exponential Weight Normalization

Theorem The gradient flow path with learning rate $\eta(t)$ for EWN is given by:

$$\frac{dw_u(t)}{dt} = -\eta(t)kw_u(t)k^2r_{w_u}L$$

Proof In case of EWN, weights are reparameterized as $w_u = e^{-\eta \int_0^t \langle v_u, r_{w_u}L \rangle} \frac{v_u}{\|v_u\|}$. Then

$$\begin{aligned} r_{w_u}L &= e^{-\eta \int_0^t \langle v_u, r_{w_u}L \rangle} \frac{v_u \langle v_u, r_{w_u}L \rangle}{\|v_u\|} \\ r_{v_u}L &= \frac{e^{-\eta \int_0^t \langle v_u, r_{w_u}L \rangle}}{\|v_u\|} \left(\langle v_u, r_{w_u}L \rangle \frac{v_u \langle v_u, r_{w_u}L \rangle}{\|v_u\|^2} - r_{w_u}L \right) \end{aligned}$$

Now, in case of gradient flow with learning rate $\eta(t)$, we can say

$$\begin{aligned} \frac{d}{dt} \langle v_u, r_{w_u}L \rangle &= -\eta(t) \langle v_u, r_{w_u}L \rangle = -\eta(t) e^{-\eta \int_0^t \langle v_u, r_{w_u}L \rangle} \frac{\langle v_u, r_{w_u}L \rangle \langle v_u, r_{w_u}L \rangle}{\|v_u\|} \\ \frac{d\|v_u\|}{dt} &= -\eta(t) \langle v_u, r_{v_u}L \rangle = -\eta(t) \frac{e^{-\eta \int_0^t \langle v_u, r_{w_u}L \rangle}}{\|v_u\|} \left(\langle v_u, r_{w_u}L \rangle \frac{\langle v_u, r_{w_u}L \rangle \langle v_u, r_{w_u}L \rangle}{\|v_u\|^2} - \langle v_u, r_{w_u}L \rangle \right) \end{aligned}$$

Now, using these equations, we can say

$$\frac{d\|v_u\|^2}{dt} = 2 \langle v_u, r_{v_u}L \rangle \langle v_u, r_{w_u}L \rangle = 0$$

Thus, $\|v_u\|$ does not change with time. As we assume $\|v_u(0)\|$ to be 1, therefore for any t , $\|v_u(t)\| = 1$. Using this simplification, we can write

$$\begin{aligned} \frac{dw_u(t)}{dt} &= \frac{d(e^{-\eta \int_0^t \langle v_u, r_{w_u}L \rangle} v_u(t))}{dt} \\ &= e^{-\eta \int_0^t \langle v_u, r_{w_u}L \rangle} \left(-\eta(t) e^{-\eta \int_0^t \langle v_u, r_{w_u}L \rangle} \langle v_u, r_{w_u}L \rangle v_u(t) + e^{-\eta \int_0^t \langle v_u, r_{w_u}L \rangle} \frac{dv_u(t)}{dt} \right) \\ &= -\eta(t) e^{-2\eta \int_0^t \langle v_u, r_{w_u}L \rangle} r_{w_u}L \end{aligned}$$

Thus, the gradient flow path with exponential weight normalization can be replicated by an adaptive learning rate given by $\eta(t) = \frac{\eta_0}{\|w_u(t)\|^2}$. ■

A.2. Standard Weight Normalization

Theorem The gradient flow path with learning rate $\eta(t)$ for SWN is given by:

$$\frac{dw_u(t)}{dt} = -\eta(t) \|w_u(t)\|^2 r_{w_u} L + \frac{1}{\|w_u(t)\|^2} (w_u(t)^\top r_{w_u} L) w_u(t)$$

Proof In case of SWN, weights are reparameterized as $w_u = \frac{v_u}{\|v_u\|}$. Then

$$r_{w_u} L = \frac{v_u^\top r_{w_u} L}{\|v_u\|}$$

$$r_{v_u} L = \frac{v_u^\top r_{w_u} L}{\|v_u\|^2} \left(\|v_u\| + \frac{v_u^\top v_u^\top r_{w_u} L}{\|v_u\|^2} r_{w_u} L \right)$$

Now, in case of gradient flow with learning rate $\eta(t)$, we can say

$$\frac{dv_u(t)}{dt} = -\eta(t) r_{v_u} L = -\eta(t) \frac{v_u(t)^\top r_{w_u} L}{\|v_u(t)\|^2}$$

$$\frac{dv_u(t)}{dt} = -\eta(t) r_{v_u} L = -\eta(t) \frac{v_u(t)}{\|v_u(t)\|} \left(\|v_u(t)\| + \frac{v_u(t)^\top v_u(t)^\top r_{w_u} L}{\|v_u(t)\|^2} r_{w_u} L \right)$$

Now, similar to EWN, $\|v_u(t)\|$ does not change with time. Using the fact that $\|v_u(t)\| = 1$ for all t , we can say

$$\begin{aligned} \frac{dw_u(t)}{dt} &= \frac{d(v_u(t)/\|v_u(t)\|)}{dt} \\ &= v_u(t) \left(-\eta(t) \frac{v_u(t)^\top r_{w_u} L}{\|v_u(t)\|^2} \right) - \eta(t) (v_u(t)^\top r_{w_u} L) v_u(t) \\ &= -\eta(t) \|v_u(t)\|^2 r_{w_u} L + (1 - v_u(t)^\top v_u(t)^\top r_{w_u} L) v_u(t) \\ &= -\eta(t) \|v_u(t)\|^2 r_{w_u} L + \frac{1 - v_u(t)^\top v_u(t)^\top r_{w_u} L}{\|v_u(t)\|^2} (w_u(t)^\top r_{w_u} L) w_u(t) \end{aligned}$$

Replacing $v_u(t)$ by $\|w_u(t)\| w_u(t)$ gives the required expression. ■

Appendix B. Proof of Proposition 2

Proposition Under assumption (A1) for gradient flow, for both SWN and EWN, $\nabla_{w(t)} L(w(t)) = 0$.

The proof for SWN, as it is homogeneous in its parameters, was provided by [Lyu and Li \(2020\)](#). We provide the proof for EWN.

First of all, for exponential loss

$$\frac{dL(t)}{dw} = \sum_i e^{y_i (w(t); x_i)} y_i r_w (w(t); x_i) \quad (11)$$

Now, using Theorem 1,

$$\frac{dL(t)}{dt} = \frac{dL(t)}{dw} \frac{dw(t)}{dt} = \sum_u k w_u(t) k^2 \frac{dL(t)}{dw_u}^2$$

Let k be the total number of neurons in the network. Then using the elementary inequality $\sum_{i=1}^n a_i^2 \geq \frac{1}{n} (\sum_{i=1}^n a_i)^2$, we get

$$\frac{dL(t)}{dt} \geq \frac{1}{k} \sum_u k w_u(t) k \frac{dL(t)}{dw_u}^2$$

Again using the fact that $w(t) \frac{dL(t)}{dw} \geq \sum_u k w_u(t) k \frac{dL(t)}{dw_u}$, we get

$$\frac{dL(t)}{dt} \geq \frac{1}{k} w(t) \frac{dL(t)}{dw}^2 \quad (12)$$

Taking the dot product with w on both sides of Equation (11) and using $r_w (w; x_i) = L (w; x_i)$ (Euler's homogeneity theorem), we get

$$w(t) \frac{dL(t)}{dw} = L \sum_i e^{y_i (w(t); x_i)} y_i (w(t); x_i)$$

Now, using the fact, that at time t_0 , $L(t_0) < 1$, which means $\min_i y_i (w(t_0); x_i) = \epsilon > 0$. Also, as we know, for gradient flow, the loss cannot go up, therefore, for any time $t \geq t_0$, $\min_i y_i (w(t); x_i) > \epsilon > 0$. Using this, we can say, for any $t \geq t_0$,

$$w(t) \frac{dL(t)}{dw} \geq \epsilon L(t)$$

Substituting this in Equation (12), we get

$$\frac{dL(t)}{dt} \leq \frac{\epsilon^2}{k} L(t)^2$$

Integrating this equation from t_0 to t , we get

$$\frac{1}{L(t)} \leq \frac{1}{L(t_0)} + \frac{\epsilon^2}{k} (t - t_0) \quad (13)$$

Clearly as t tends to ∞ , RHS tends to ∞ and thus L tends to 0.

Appendix C. Proof of Theorem 3

Theorem Consider a node u in the network with $\|g_u\| > 0$ and $\lim_{t \rightarrow \infty} \|w_u(t)\| = 1$. Under assumptions (A1), (A2) for gradient flow and (B1)-(B3) for gradient descent, for both SWN and EWN

- (i) $\lim_{t \rightarrow \infty} \frac{w_u(t)}{\|w_u(t)\|} := v_u$ exists. (ii) $v_u = \frac{g_u}{\|g_u\|}$ for some $\alpha > 0$.

The proof for different cases will be split into different subsections and corresponding theorem will be stated there for ease of the reader. The proof will depend on the Stolz Cesaro theorem (stated in Appendix K), Integral Form of Stolz-Cesaro Theorem (stated and proved in Appendix K) and following lemmas that have been proved in Appendix J.

Lemma 9 Consider sequence a satisfying the following properties

1. $a_k > 0$
2. $\sum_{k=0}^{\infty} a_k = 1$
3. $\lim_{k \rightarrow \infty} a_k = 0$

Then $\sum_{k=0}^{\infty} \frac{a_k}{\sum_{j=0}^k a_j^2} = 1$

Lemma 10 Consider two sequences a and b satisfying the following properties

1. $a_k > 0$; $\sum_{k=0}^{\infty} a_k = 1$ and $\lim_{k \rightarrow \infty} a_k = 0$
2. $b_0 > 0$, b is increasing and $b_{k+1}^2 = b_k^2 + \frac{a_k}{b_k}$

Then $\sum_{k=0}^{\infty} \frac{a_k}{b_k} = 1$.

C.1. Exponential Weight Normalization

In this section, we will use $\|w_u(t)\|$ and $\|v_u(t)\|$ interchangeably.

C.1.1. GRADIENT FLOW

Theorem Consider a node u in the network with $\|g_u\| > 0$ and $\lim_{t \rightarrow \infty} \|w_u(t)\| = 1$. Under assumptions (A1), (A2) for gradient flow, for EWN

- (i) $\lim_{t \rightarrow \infty} \frac{w_u(t)}{\|w_u(t)\|} := v_u$ exists. (ii) $v_u = \frac{g_u}{\|g_u\|}$ for some $\alpha > 0$.

Update Equations:

$$\frac{d w_u(t)}{dt} = -\|w_u(t)\|^{-1} (v_u(t) \cdot \nabla_{w_u} L(w(t))) w_u(t) \quad (14)$$

$$\frac{d v_u(t)}{dt} = -\|w_u(t)\|^{-1} (I - v_u(t)v_u(t)^\top) \nabla_{w_u} L(w(t)) \quad (15)$$

Proof As $\|g_u\| > 0$, therefore $\nabla_{w_u} L(w(t))$ converges in direction. Therefore, for every $\epsilon > 0$, there exists a time $t_1(\epsilon)$, such that for $t > t_1$, $\frac{\| \nabla_{w_u} L(w(t)) \|}{\| \nabla_{w_u} L(w(t_1)) \|} > \frac{\|g_u\|}{\|g_u\|} \cos(\epsilon)$.

Now, Let's assume that $\mathbf{w}_u(t)$ does not converge in the direction \mathbf{g}_u . Then, there must exist a $\delta > 0$, such that for this δ , there exists a time $t_2 > t_1$ satisfying $\langle \mathbf{w}_u(t_2), \mathbf{g}_u \rangle < \cos(\delta)$, where $\delta > 0$.

Now, we are going to show that for any $\delta > 0$, there exists a time $t_3 > t_2$ such that $\langle \mathbf{w}_u(t_3), \mathbf{g}_u \rangle > \cos(\delta)$. Let's say for a given δ , no such t_3 exists. Then, taking dot product with $\frac{\mathbf{g}_u}{\|\mathbf{g}_u\|}$ on both sides of Equation (15), we can say

$$\frac{\mathbf{g}_u}{\|\mathbf{g}_u\|} \cdot \frac{d\mathbf{w}_u(t)}{dt} = -\eta e^{-\eta t} \|\mathbf{w}_u(t)\| \frac{\mathbf{g}_u}{\|\mathbf{g}_u\|} \cdot (\|\mathbf{w}_u(t)\| \langle \mathbf{w}_u(t), \mathbf{g}_u \rangle) \frac{\mathbf{g}_u}{\|\mathbf{g}_u\|} \cdot \frac{\mathbf{w}_u(t)}{\|\mathbf{w}_u(t)\|}$$

Now, as $\langle \mathbf{w}_u(t_2), \mathbf{g}_u \rangle < \cos(\delta)$ and $\langle \mathbf{w}_u(t_2), \mathbf{g}_u \rangle > \cos(\delta)$, we can say

$$\frac{\mathbf{g}_u}{\|\mathbf{g}_u\|} \cdot \frac{d\mathbf{w}_u(t)}{dt} \leq -\eta e^{-\eta t} \|\mathbf{w}_u(t)\| (\cos(\delta) - \cos(\delta)) \tag{16}$$

Now, using the fact that $\|\mathbf{w}_u(t)\| \geq 1$ and using Equation (14), we can say

$$\int_{t_2}^{\infty} \eta e^{-\eta t} dt = 1$$

Using this fact and integrating the Equation (16) on both the sides from t_2 to ∞ , we get a contradiction as vectors on LHS have a finite norm while RHS tends to $-\infty$. Thus, for every δ between 0 and $\pi/2$, there must exist t_3 , such that $\langle \mathbf{w}_u(t_3), \mathbf{g}_u \rangle > \cos(\delta)$.

Now, we are going to show for all t_3 , $\langle \mathbf{w}_u(t) > \frac{\mathbf{g}_u}{\|\mathbf{g}_u\|} \cdot \mathbf{w}_u(t) > \cos(\delta)$. Now, consider any $t > t_3$ such that $\langle \mathbf{w}_u(t), \mathbf{g}_u \rangle < \cos(\delta)$. Using similar argument as in Equation (16), we can say, if for any $t > t_3$, $\langle \mathbf{w}_u(t), \mathbf{g}_u \rangle < \cos(\delta)$, then

$$\frac{\mathbf{g}_u}{\|\mathbf{g}_u\|} \cdot \frac{d\mathbf{w}_u(t)}{dt} \leq -\eta e^{-\eta t} \|\mathbf{w}_u(t)\| (\cos(\delta) - \cos(\delta)) \tag{17}$$

This means that the dot product between $\frac{\mathbf{g}_u}{\|\mathbf{g}_u\|}$ and $\mathbf{w}_u(t)$ goes up, whenever $\langle \mathbf{w}_u(t), \mathbf{g}_u \rangle < \cos(\delta)$. Therefore, its not possible that $\langle \mathbf{w}_u(t), \mathbf{g}_u \rangle < \cos(\delta)$ for any $t > t_3$. As δ can be arbitrarily chosen between 0 and $\pi/2$, $\mathbf{w}_u(t)$ converges in the direction \mathbf{g}_u . ■

C.1.2. GRADIENT DESCENT

Theorem Consider a node u in the network with $\|\mathbf{g}_u\| > 0$ and $\lim_{t \rightarrow \infty} \|\mathbf{w}_u(t)\| = 1$. Under assumptions (B1)-(B3) for gradient descent, for EWN

- (i) $\lim_{t \rightarrow \infty} \frac{\mathbf{w}_u(t)}{\|\mathbf{w}_u(t)\|} := \mathbf{w}_u$ exists.
- (ii) $\mathbf{w}_u = \frac{\mathbf{g}_u}{\|\mathbf{g}_u\|}$ for some $\eta > 0$.

Update Equations:

$$\mathbf{w}_u(t+1) = \mathbf{w}_u(t) - \eta e^{-\eta t} \frac{\langle \mathbf{w}_u(t), \mathbf{g}_u \rangle \mathbf{g}_u}{\|\mathbf{w}_u(t)\|} \tag{18}$$

$$v_u(t+1) = v_u(t) \left(t \right) \frac{e^{-u(t)}}{kv_u(t)k} \left| \frac{v_u(t)v_u(t)^>}{kv_u(t)k^2} r_{w_u} L(w(t)) \right. \quad (19)$$

Proof

As $\frac{g_u}{kg_u k} > 0$, therefore $r_{w_u} L(w(t))$ converges in direction. Therefore, for every ϵ satisfying $0 < \epsilon < 2$, there exists a time $t_1(\epsilon)$, such that for $t > t_1(\epsilon)$, $\frac{r_{w_u} L(w(t))}{kr_{w_u} L(w(t))k} > \frac{g_u}{kg_u k} \cos(\epsilon)$. Now, Let's assume that $v_u(t)$ does not converge in the direction of g_u . Then, there must exist a ϵ satisfying $0 < \epsilon < 2$, such that for this ϵ , there exists a time $t_2 > t_1(\epsilon)$ satisfying $v_u(t_2)^> \frac{g_u}{kg_u k} = \cos(\epsilon)$, where $\epsilon > \epsilon_0$.

Now, we are going to show that for any ϵ satisfying $\epsilon < \epsilon_0$, there exists a time $t_3 > t_2$ such that $\frac{v_u(t_3)^>}{kv_u(t_3)k} > \frac{g_u}{kg_u k} > \cos(\epsilon)$. Let's say for a given ϵ , no such t_3 exists. Then, taking dot product with $\frac{g_u}{kg_u k}$ on both sides of Equation (19), we can say

$$\frac{v_u(t+1)^> g_u}{kg_u k} = \frac{v_u(t)^> g_u}{kg_u k} + \left(t \right) \frac{e^{-u(t)}}{kv_u(t)k} kr_{w_u} L(w(t))k \frac{g_u}{kg_u k} > \left| \frac{v_u(t)v_u(t)^>}{kv_u(t)k^2} \frac{r_{w_u} L(w(t))}{kr_{w_u} L(w(t))k} \right.$$

Now, as $\frac{g_u}{kg_u k} > \frac{r_{w_u} L(w(t))}{kr_{w_u} L(w(t))k} \cos(\epsilon)$ and $\frac{g_u}{kg_u k} > \frac{v_u(t)}{kv_u(t)k} \cos(\epsilon)$, we can say

$$\frac{v_u(t+1)^> g_u}{kg_u k} = \frac{v_u(t)^> g_u}{kg_u k} + (\cos(\epsilon) - \cos(\epsilon)) \left(t \right) \frac{e^{-u(t)}}{kv_u(t)k} kr_{w_u} L(w(t))k \quad (20)$$

However, in this case $kv_u(t)k$ doesn't stay constant and thus increase in dot product doesn't directly correspond to an increase in angle. Now, using Equation (19), we can say

$$kv_u(t+1)k^2 = kv_u(t)k^2 + \left(t \right) \frac{e^{-u(t)}}{kv_u(t)k} kr_{w_u} L(w(t))k \quad (21)$$

Using the above two equations, we can say, for time t_2 ,

$$\frac{v_u(t+1)^> g_u}{kv_u(t+1)k kg_u k} = \frac{\frac{v_u(t)^> g_u}{kg_u k} + (\cos(\epsilon) - \cos(\epsilon)) \left(t \right) \frac{e^{-u(t)}}{kv_u(t)k} kr_{w_u} L(w(t))k}{kv_u(t)k^2 + \left(t \right) \frac{e^{-u(t)}}{kv_u(t)k} kr_{w_u} L(w(t))k^2}$$

Unrolling the equation above, we get

$$\frac{v_u(t+1)^> g_u}{kv_u(t+1)k kg_u k} = \frac{\frac{v_u(t_2)^> g_u}{kg_u k} + \prod_{k=t_2}^t (\cos(\epsilon) - \cos(\epsilon)) \left(k \right) \frac{e^{-u(k)}}{kv_u(k)k} kr_{w_u} L(w(k))k}{kv_u(t_2)k^2 + \prod_{k=t_2}^t \left(k \right) \frac{e^{-u(k)}}{kv_u(k)k} kr_{w_u} L(w(k))k^2} \quad (22)$$

Now, as $v_u(t) \neq 1$, therefore, using Equation (18), we can say

$$\prod_{k=t_2}^t \left(k \right) \frac{e^{-u(k)}}{kv_u(k)k} kr_{w_u} L(w(k))k = 1$$

Now, using this identity, along with the Assumption (A5), Equation (21) and Lemma 10, we can say

$$\sum_{k=t_2}^{\infty} \frac{e^{-u(k)}}{kv_u(k)} \leq \frac{1}{w_u} L(w(k)) \leq 1$$

Using this along with Equation (22) and Lemma 9, we can say

$$\lim_{t \rightarrow \infty} \frac{v_u(t+1) \cdot \theta_u}{kv_u(t+1) \cdot k\theta_u} = 1$$

However, this is not possible as the vectors on LHS have bounded norm. This contradicts. Thus there must exist t_3 such that $\frac{v_u(t_3)}{kv_u(t_3)} \cdot \theta_u > \cos(\alpha)$.

Now, we are going to show that there exists t_4 , such that for all $t > t_4$, $\frac{v_u(t)}{kv_u(t)} \cdot \theta_u > \cos(\alpha)$. Consider a δ such that $\cos(\alpha) - \delta < \cos(\alpha)$. Now, if at any time t , $\frac{v_u(t)}{kv_u(t)} \cdot \theta_u < \cos(\alpha)$, then, similar to Equation (20), we can say

$$\frac{v_u(t+1) \cdot \theta_u}{k\theta_u} \leq \frac{v_u(t) \cdot \theta_u}{k\theta_u} + (\cos(\alpha) - \cos(\alpha - \delta)) \cdot \frac{e^{-u(t)}}{kv_u(t)} \leq \frac{v_u(t) \cdot \theta_u}{k\theta_u} + \delta \cdot \frac{e^{-u(t)}}{kv_u(t)}$$

Using the upper bound $\frac{v_u(t)}{kv_u(t)}$ from Equation (21), we can say

$$\frac{v_u(t+1) \cdot \theta_u}{kv_u(t+1) \cdot k\theta_u} \leq \frac{\frac{v_u(t) \cdot \theta_u}{k\theta_u} + (\cos(\alpha) - \cos(\alpha - \delta)) \cdot \frac{e^{-u(t)}}{kv_u(t)}}{\frac{kv_u(t)k^2 + \frac{e^{-u(t)}}{kv_u(t)}}{k^2}} \quad (23)$$

Let $\frac{e^{-u(t)}}{kv_u(t)}$ be denoted by $\rho(t)$. Then, the above equation can be rewritten as

$$\frac{v_u(t+1) \cdot \theta_u}{kv_u(t+1) \cdot k\theta_u} \leq \frac{v_u(t) \cdot \theta_u}{kv_u(t) \cdot k\theta_u} \rho + \frac{kv_u(t)k}{kv_u(t)k^2 + \rho(t)^2} + (\cos(\alpha) - \cos(\alpha - \delta)) \rho \frac{\rho(t)}{kv_u(t)k^2 + \rho(t)^2}$$

Now, we are going to show that for a small enough δ , RHS is greater than $\frac{v_u(t) \cdot \theta_u}{kv_u(t) \cdot k\theta_u}$.

$$\begin{aligned} & \frac{v_u(t) \cdot \theta_u}{kv_u(t) \cdot k\theta_u} \rho + \frac{kv_u(t)k}{kv_u(t)k^2 + \rho(t)^2} + (\cos(\alpha) - \cos(\alpha - \delta)) \rho \frac{\rho(t)}{kv_u(t)k^2 + \rho(t)^2} > \frac{v_u(t) \cdot \theta_u}{kv_u(t) \cdot k\theta_u} \\ \Rightarrow & (\cos(\alpha) - \cos(\alpha - \delta)) \rho \frac{\rho(t)}{kv_u(t)k^2 + \rho(t)^2} > \frac{v_u(t) \cdot \theta_u}{kv_u(t) \cdot k\theta_u} \rho - \frac{kv_u(t)k}{kv_u(t)k^2 + \rho(t)^2} \\ \Rightarrow & (\cos(\alpha) - \cos(\alpha - \delta)) > \frac{v_u(t) \cdot \theta_u}{kv_u(t) \cdot k\theta_u} \frac{\rho}{\rho(t)} - \frac{kv_u(t)k}{kv_u(t)k^2 + \rho(t)^2} \end{aligned}$$

Clearly as $\rho(t) \rightarrow 0$, the RHS tends to 0, therefore the equation is satisfied. Thus for a small enough δ , RHS of Equation (23) is greater than $\frac{v_u(t) \cdot \theta_u}{kv_u(t) \cdot k\theta_u}$. As $kv_u(t)k$ keeps on increasing and

by Assumption (B3) $\lim_{t \rightarrow \infty} \|w_u(t) - w_u\| = 0$, we can say there exists a time t_5 such that for any $t > t_5$, $\frac{v_u(t) \cdot g_u}{\|v_u(t)\| \|g_u\|} > \cos(\theta)$ goes up whenever $\frac{v_u(t)}{\|v_u(t)\|} > \frac{g_u}{\|g_u\|} < \cos(\theta)$.

Also, by using Equation (19) and Assumption (B3), we can say, that there exists a time t_6 such that for $t > t_6$, $\frac{v_u(t)}{\|v_u(t)\|} > \frac{g_u}{\|g_u\|} > \cos(\theta) \Rightarrow \frac{v_u(t+1)}{\|v_u(t+1)\|} > \frac{g_u}{\|g_u\|} > \cos(\theta)$, as the RHS of Equation (19) goes to 0 norm in limit. Now, determine $t_4 > \max(t_5, t_6)$ such that $\frac{v_u(t_4)}{\|v_u(t_4)\|} > \frac{g_u}{\|g_u\|} > \cos(\theta)$ (must exist from previous arguments). Then, as the dot product always goes up when between $\cos(\theta)$ and $\cos(\theta)$, and can't go in a single step from being greater than $\cos(\theta)$ to less than $\cos(\theta)$, therefore, for every $t > t_4$, $\frac{v_u(t)}{\|v_u(t)\|} > \frac{g_u}{\|g_u\|} > \cos(\theta)$.

Now as the above argument holds for any θ between 0 and $\pi/2$, and for any $\epsilon > 0$, we can say that $w_u(t)$ converges in direction of g_u . ■

C.2. Standard Weight Normalization

C.2.1. GRADIENT FLOW

Theorem Consider a node u in the network with $\|g_u\| > 0$ and $\lim_{t \rightarrow \infty} \|w_u(t)\| = 1$. Under assumptions (A1), (A2) for gradient flow, for SWN

- (i) $\lim_{t \rightarrow \infty} \frac{w_u(t)}{\|w_u(t)\|} := w_u$ exists.
- (ii) $w_u = \frac{g_u}{\|g_u\|}$ for some $\epsilon > 0$.

Update Equations:

$$\frac{dw_u(t)}{dt} = -\frac{v_u(t) \cdot r_{w_u} L(w(t))}{\|v_u(t)\|} \quad (24)$$

$$\frac{dv_u(t)}{dt} = -\frac{v_u(t)}{\|v_u(t)\|} \left(\frac{v_u(t) \cdot v_u(t)}{\|v_u(t)\|^2} \cdot r_{w_u} L(w(t)) \right) \quad (25)$$

Proof The proof will be given for $\|w_u\| = 1$. The one for $\|w_u\| \neq 1$ can be handled similarly.

As $\|g_u\| > 0$, therefore $r_{w_u} L(w(t))$ converges in direction. Therefore, for every ϵ satisfying $0 < \epsilon < 2$, there exists a time $t_1(\epsilon)$, such that for $t > t_1(\epsilon)$, $\frac{r_{w_u} L(w(t))}{\|r_{w_u} L(w(t))\|} > \frac{g_u}{\|g_u\|} \cos(\epsilon)$. Now, Let's assume that $w_u(t)$ does not converge in the direction of g_u . Then, there must exist a δ satisfying $0 < \delta < 2$, such that for this δ , there exists a time $t_2 > t_1(\epsilon)$ satisfying $v_u(t_2) \cdot \frac{g_u}{\|g_u\|} = \cos(\delta)$, where $\delta > \epsilon$.

Now, we are going to show that for any ϵ satisfying $0 < \epsilon < \delta$, there exists a time $t_3 > t_2$ such that $v_u(t_3) \cdot \frac{g_u}{\|g_u\|} > \cos(\epsilon)$. Let's say for a given ϵ , no such t_3 exists. Then, taking dot product with $\frac{g_u}{\|g_u\|}$ on both sides of Equation (25), we can say

$$\frac{g_u}{\|g_u\|} \cdot \frac{dv_u(t)}{dt} = -\frac{v_u(t) \cdot r_{w_u} L(w(t))}{\|v_u(t)\|} \left(\frac{g_u}{\|g_u\|} \cdot \left(\frac{v_u(t) \cdot v_u(t)}{\|v_u(t)\|^2} \right) \cdot \frac{r_{w_u} L(w(t))}{\|r_{w_u} L(w(t))\|} \right)$$

Now, as $\frac{g_u}{\|g_u\|} \cdot \frac{r_{w_u} L(w(t))}{\|r_{w_u} L(w(t))\|} > \cos(\epsilon)$ and $\frac{g_u}{\|g_u\|} \cdot v_u(t) > \cos(\delta)$, we can say

$$\frac{g_u}{\|g_u\|} \cdot \frac{dv_u(t)}{dt} < -\frac{v_u(t) \cdot r_{w_u} L(w(t))}{\|v_u(t)\|} (\cos(\epsilon) - \cos(\delta)) \quad (26)$$

Now, using the fact that $\|w_u\| \leq 1$ and using Equation (24), we can say

$$\int_{t=t_2}^1 \|w_u(t)\|_{w_u} L(w(t)) dt = 1$$

Using this fact and integrating the Equation (26) on both the sides from t_2 to 1, we get a contradiction as vectors on LHS have a finite norm while RHS tends to infinity. Thus, for every ϵ between $\cos(\theta)$ and $\frac{\|g_u\|}{\|g_u\|}$, there must exist t_3 , such that $v_u(t_3) > \frac{\|g_u\|}{\|g_u\|} > \cos(\theta)$.

Now, we are going to show for all $t_3 < t < 1$, $v_u(t) > \frac{\|g_u\|}{\|g_u\|} > \cos(\theta)$. Now, consider any t_4 such that $\cos(\theta) < t_4 < 1$. Using similar argument as in Equation (26), we can say, if for any $t_3 < t_4$, $v_u(t_4) > \frac{\|g_u\|}{\|g_u\|} < \cos(\theta)$, then

$$\frac{\|g_u\|}{\|g_u\|} > \frac{dv_u(t_4)}{dt} \quad (t_4) \quad v_u(t_4) \|w_u(t_4)\| L(w(t_4)) (\cos(\theta) - \cos(\theta)) \quad (27)$$

This means that the dot product between $\frac{g_u}{\|g_u\|}$ and $v_u(t)$ goes up, whenever $\frac{g_u}{\|g_u\|} \cdot v_u(t) < \cos(\theta)$. Therefore, its not possible that $v_u(t) > \frac{\|g_u\|}{\|g_u\|} < \cos(\theta)$ for any $t > t_3$. As ϵ can be arbitrarily chosen between $\cos(\theta)$ and $\frac{\|g_u\|}{\|g_u\|}$, $w_u(t)$ converges in the direction of g_u . ■

C.2.2. GRADIENT DESCENT

Theorem Consider a node u in the network with $\|g_u\| > 0$ and $\lim_{t \rightarrow 1} \|w_u(t)\| = 1$. Under assumptions (B1)-(B3) for gradient descent, for SWN

- (i) $\lim_{t \rightarrow 1} \frac{w_u(t)}{\|w_u(t)\|} := w_u$ exists.
- (ii) $w_u = \frac{g_u}{\|g_u\|}$ for some $\theta > 0$.

Update Equations:

$$w_u(t+1) = w_u(t) - \eta \frac{v_u(t) \cdot r_{w_u} L(w(t))}{\|v_u(t)\|} \quad (28)$$

$$v_u(t+1) = v_u(t) - \eta \left(\frac{v_u(t)}{\|v_u(t)\|} \cdot \frac{v_u(t) \cdot v_u(t)}{\|v_u(t)\|^2} - r_{w_u} L(w(t)) \right) \quad (29)$$

Proof The proof will be given for $\|g_u\| \leq 1$. The one for $\|g_u\| > 1$ can be handled similarly.

As $\|g_u\| > 0$, therefore $r_{w_u} L(w(t))$ converges in direction. Therefore, for every ϵ satisfying $0 < \epsilon < 2$, there exists a time $t_1(\epsilon)$, such that for $t > t_1(\epsilon)$, $\frac{r_{w_u} L(w(t))}{\|r_{w_u} L(w(t))\|} > \frac{\|g_u\|}{\|g_u\|} > \cos(\theta)$. Now, Let's assume that $w_u(t)$ does not converge in the direction of g_u . Then, there must exist a δ satisfying $0 < \delta < 2$, such that for this δ , there exists a time $t_2 > t_1(\epsilon)$ satisfying $v_u(t_2) \cdot \frac{g_u}{\|g_u\|} = \cos(\theta) - \delta$, where $\delta > 0$.

Now, we are going to show that for any ϵ satisfying $0 < \epsilon < 2$, there exists a time $t_3 > t_2$ such that $\frac{v_u(t_3)}{\|v_u(t_3)\|} \cdot \frac{g_u}{\|g_u\|} > \cos(\theta)$. Let's say for a given ϵ , no such t_3 exists. Then, taking

dot product with $\frac{\mathbf{g}_u}{\|\mathbf{g}_u\|}$ on both sides of Equation (29), we can say

$$\frac{\mathbf{v}_u(t+1) \cdot \mathbf{g}_u}{\|\mathbf{g}_u\|} = \frac{\mathbf{v}_u(t) \cdot \mathbf{g}_u}{\|\mathbf{g}_u\|} + \left(\frac{\mathbf{v}_u(t)}{\|\mathbf{v}_u(t)\|} \right)^T \mathbf{r}_{w_u} \mathbf{L}(w(t)) \frac{\mathbf{g}_u}{\|\mathbf{g}_u\|} \quad | \quad \frac{\mathbf{v}_u(t) \cdot \mathbf{v}_u(t)}{\|\mathbf{v}_u(t)\|^2} \frac{\mathbf{r}_{w_u} \mathbf{L}(w(t))}{\|\mathbf{r}_{w_u} \mathbf{L}(w(t))\|}$$

Now, as $\frac{\mathbf{g}_u}{\|\mathbf{g}_u\|} \cdot \frac{\mathbf{r}_{w_u} \mathbf{L}(w(t))}{\|\mathbf{r}_{w_u} \mathbf{L}(w(t))\|} = \cos(\theta)$ and $\frac{\mathbf{g}_u}{\|\mathbf{g}_u\|} \cdot \frac{\mathbf{v}_u(t)}{\|\mathbf{v}_u(t)\|} = \cos(\phi)$, we can say

$$\frac{\mathbf{v}_u(t+1) \cdot \mathbf{g}_u}{\|\mathbf{g}_u\|} = \frac{\mathbf{v}_u(t) \cdot \mathbf{g}_u}{\|\mathbf{g}_u\|} + (\cos(\theta) - \cos(\phi)) \left(\frac{\mathbf{v}_u(t)}{\|\mathbf{v}_u(t)\|} \right)^T \mathbf{r}_{w_u} \mathbf{L}(w(t)) \frac{\mathbf{g}_u}{\|\mathbf{g}_u\|} \quad (30)$$

However, in this case, $\|\mathbf{v}_u(t)\|$ doesn't stay constant and thus increase in dot product doesn't directly correspond to an increase in angle. Now, using Equation (29), we can say

$$\|\mathbf{v}_u(t+1)\|^2 = \|\mathbf{v}_u(t)\|^2 + \left(\frac{\mathbf{v}_u(t)}{\|\mathbf{v}_u(t)\|} \right)^T \mathbf{r}_{w_u} \mathbf{L}(w(t)) \frac{\mathbf{v}_u(t)}{\|\mathbf{v}_u(t)\|} \quad (31)$$

Using the above two equations, we can say, for time t_2 ,

$$\frac{\mathbf{v}_u(t+1) \cdot \mathbf{g}_u}{\|\mathbf{v}_u(t+1)\| \|\mathbf{g}_u\|} = \frac{\frac{\mathbf{v}_u(t) \cdot \mathbf{g}_u}{\|\mathbf{g}_u\|} + \left(\frac{\mathbf{v}_u(t)}{\|\mathbf{v}_u(t)\|} \right)^T \mathbf{r}_{w_u} \mathbf{L}(w(t)) \frac{\mathbf{g}_u}{\|\mathbf{g}_u\|}}{\sqrt{\|\mathbf{v}_u(t)\|^2 + \left(\frac{\mathbf{v}_u(t)}{\|\mathbf{v}_u(t)\|} \right)^T \mathbf{r}_{w_u} \mathbf{L}(w(t)) \frac{\mathbf{v}_u(t)}{\|\mathbf{v}_u(t)\|}}}$$

Unrolling the equation above, we get

$$\frac{\mathbf{v}_u(t+1) \cdot \mathbf{g}_u}{\|\mathbf{v}_u(t+1)\| \|\mathbf{g}_u\|} = \frac{\frac{\mathbf{v}_u(t_2) \cdot \mathbf{g}_u}{\|\mathbf{g}_u\|} + \sum_{k=t_2}^t \left(\frac{\mathbf{v}_u(k)}{\|\mathbf{v}_u(k)\|} \right)^T \mathbf{r}_{w_u} \mathbf{L}(w(k)) \frac{\mathbf{g}_u}{\|\mathbf{g}_u\|}}{\sqrt{\|\mathbf{v}_u(t_2)\|^2 + \sum_{k=t_2}^t \left(\frac{\mathbf{v}_u(k)}{\|\mathbf{v}_u(k)\|} \right)^T \mathbf{r}_{w_u} \mathbf{L}(w(k)) \frac{\mathbf{v}_u(k)}{\|\mathbf{v}_u(k)\|}}} \quad (32)$$

Now, as $\mathbf{v}_u(t) \cdot \mathbf{1} = 1$, therefore, using Equation (29), we can say

$$\sum_{k=t_2}^t \left(\frac{\mathbf{v}_u(k)}{\|\mathbf{v}_u(k)\|} \right)^T \mathbf{r}_{w_u} \mathbf{L}(w(k)) \frac{\mathbf{1}}{\|\mathbf{1}\|} = 1$$

Now, using this identity, along with the assumption (A5), Equation (31) and Lemma 10, we can say

$$\sum_{k=t_2}^t \left(\frac{\mathbf{v}_u(k)}{\|\mathbf{v}_u(k)\|} \right)^T \mathbf{r}_{w_u} \mathbf{L}(w(k)) \frac{\mathbf{1}}{\|\mathbf{1}\|} = 1$$

Using this along with Equation (32) and Lemma 9, we can say

$$\lim_{t \rightarrow \infty} \frac{\mathbf{v}_u(t+1) \cdot \mathbf{g}_u}{\|\mathbf{v}_u(t+1)\| \|\mathbf{g}_u\|} = 1$$

However, this is not possible as the vectors on LHS have bounded norm. This contradicts. Thus there must exist t_3 such that $\frac{v_u(t_3)}{kv_u(t_3)k} > \frac{g_u}{kg_u k} > \cos(\theta)$.

Now, we are going to show that there exists $t_4 > t_3$, such that for all $t > t_4$, $\frac{v_u(t)}{kv_u(t)k} > \frac{g_u}{kg_u k} > \cos(\theta)$. Consider a t such that $\frac{v_u(t)}{kv_u(t)k} < \frac{g_u}{kg_u k} < \cos(\theta)$. Now, if at any time t , $\frac{v_u(t)}{kv_u(t)k} < \frac{g_u}{kg_u k} < \cos(\theta)$, then, similar to Equation (30), we can say

$$\frac{v_u(t+1) \cdot g_u}{kg_u k} = \frac{v_u(t) \cdot g_u}{kg_u k} + (\cos(\theta) - \cos(\theta)) \cdot \left(\frac{v_u(t)}{kv_u(t)k} \right) \cdot kr_{w_u} L(w(t))k$$

Using the upper bound $kr_{w_u} L(w(t))k$ from Equation (31), we can say

$$\frac{v_u(t+1) \cdot g_u}{kv_u(t+1)k kg_u k} \leq \frac{\frac{v_u(t) \cdot g_u}{kg_u k} + (\cos(\theta) - \cos(\theta)) \cdot \left(\frac{v_u(t)}{kv_u(t)k} \right) \cdot kr_{w_u} L(w(t))k}{kv_u(t)k^2 + \left(\frac{v_u(t)}{kv_u(t)k} \right) \cdot kr_{w_u} L(w(t))k} \quad (33)$$

Let $\left(\frac{v_u(t)}{kv_u(t)k} \right) \cdot kr_{w_u} L(w(t))k$ be denoted by (t) . Then, the above equation can be rewritten as

$$\frac{v_u(t+1) \cdot g_u}{kv_u(t+1)k kg_u k} \leq \frac{v_u(t) \cdot g_u}{kv_u(t)k kg_u k} \cdot \frac{kv_u(t)k}{kv_u(t)k^2 + (t)^2} + (\cos(\theta) - \cos(\theta)) \cdot \frac{(t)}{kv_u(t)k^2 + (t)^2}$$

Now, we are going to show that for a small enough (t) , RHS is greater than $\frac{v_u(t) \cdot g_u}{kv_u(t)k kg_u k}$.

$$\begin{aligned} & \frac{v_u(t) \cdot g_u}{kv_u(t)k kg_u k} \cdot \frac{kv_u(t)k}{kv_u(t)k^2 + (t)^2} + (\cos(\theta) - \cos(\theta)) \cdot \frac{(t)}{kv_u(t)k^2 + (t)^2} > \frac{v_u(t) \cdot g_u}{kv_u(t)k kg_u k} \\ \Rightarrow & (\cos(\theta) - \cos(\theta)) \cdot \frac{(t)}{kv_u(t)k^2 + (t)^2} > \frac{v_u(t) \cdot g_u}{kv_u(t)k kg_u k} \cdot \left(1 - \frac{kv_u(t)k}{kv_u(t)k^2 + (t)^2} \right) \\ \Rightarrow & (\cos(\theta) - \cos(\theta)) > \frac{v_u(t) \cdot g_u}{kv_u(t)k kg_u k} \cdot \frac{kv_u(t)k}{(t)} \end{aligned}$$

Clearly as $(t) \rightarrow 0$, the RHS tends to 0, therefore the equation is satisfied. Thus for a small enough (t) , RHS of Equation (33) is greater than $\frac{v_u(t) \cdot g_u}{kv_u(t)k kg_u k}$. As $kv_u(t)k$ keeps on increasing and by Assumption (A5) $\lim_{t \rightarrow \infty} \left(\frac{v_u(t)}{kv_u(t)k} \right) \cdot kr_{w_u} L(w(t))k = 0$, we can say there exists a time t_5 such that for any $t > t_5$, $\frac{v_u(t) \cdot g_u}{kv_u(t)k kg_u k}$ goes up whenever $\frac{v_u(t)}{kv_u(t)k} > \frac{g_u}{kg_u k} < \cos(\theta)$.

Also, by using Equation (29) and Assumption (A5), we can say, that there exists a time t_6 such that for $t > t_6$, $\frac{v_u(t)}{kv_u(t)k} > \frac{g_u}{kg_u k} > \cos(\theta) \Rightarrow \frac{v_u(t+1)}{kv_u(t+1)k} > \frac{g_u}{kg_u k} > \cos(\theta)$, as the RHS of Equation (29) goes to 0 norm in limit. Now, define $t_4 = \max(t_5, t_6)$ such that $\frac{v_u(t_4)}{kv_u(t_4)k} > \frac{g_u}{kg_u k} > \cos(\theta)$ (must exist from previous arguments). Then, as the dot product always goes up when between $\cos(\theta)$ and $\cos(\theta)$, and can't go in a single step from being greater than $\cos(\theta)$ to less than $\cos(\theta)$, therefore, for every $t > t_4$, $\frac{v_u(t)}{kv_u(t)k} > \frac{g_u}{kg_u k} > \cos(\theta)$.

Now as the above argument holds for any ϵ between ϵ_1 and ϵ_2 , and for any $\delta > 0$, we can say that $w_u(t)$ converges in direction of g_u . ■

Appendix D. Proof of Theorem 4

Theorem Consider two nodes u and v in the network with $\|g_u\|, \|g_v\| > 0$; $\lim_{t \rightarrow \infty} \|w_u(t)\| = 1$ and $\lim_{t \rightarrow \infty} \|w_v(t)\| = 1$. Let $\frac{\|g_u\|}{\|g_v\|}$ be denoted by c . Under assumptions (A1), (A2) for gradient flow and (B1)-(B3) for gradient descent,

- (i) for SWN, $\lim_{t \rightarrow \infty} \frac{\|w_u(t)\|}{\|w_v(t)\|} = c$
- (ii) for EWN, $\lim_{t \rightarrow \infty} \frac{\|w_u(t)\|}{\|w_v(t)\|}$ is either 0, 1 or $\frac{1}{c}$

Proof for different cases will be split into different subsections and the corresponding case will be restated there for ease of the reader.

D.1. Exponential Weight Normalization

D.1.1. GRADIENT FLOW

Theorem Consider two nodes u and v in the network with $\|g_u\|, \|g_v\| > 0$; $\lim_{t \rightarrow \infty} \|w_u(t)\| = 1$ and $\lim_{t \rightarrow \infty} \|w_v(t)\| = 1$. Let $\frac{\|g_u\|}{\|g_v\|}$ be denoted by c . Under assumptions (A1), (A2) for gradient flow, for EWN, $\lim_{t \rightarrow \infty} \frac{\|w_u(t)\|}{\|w_v(t)\|}$ is either 0, 1 or $\frac{1}{c}$

Proof Using Theorem 3 and the fact that $\lim_{t \rightarrow \infty} \frac{\|w_u(t)\|}{\|w_v(t)\|} = c$, for any $0 < \epsilon < c$ and $0 < \delta < 2$, there exists a time t_1 , such that for $t > t_1$, the following hold

- (i) $\frac{\|w_u(t)\|}{\|w_v(t)\|} \geq 2 [c - \epsilon; c + \epsilon]$ (ii) $\frac{\|w_u(t)\|}{\|w_v(t)\|} > \frac{r_{w_u} L(w(t))}{r_{w_v} L(w(t))} \cos(\delta)$
- (iii) $\frac{\|w_v(t)\|}{\|w_u(t)\|} > \frac{r_{w_v} L(w(t))}{r_{w_u} L(w(t))} \cos(\delta)$.

Now, we will provide the proof of part (iii) of Proposition 5, i.e, for EWN, if at some time $t_2 > t_1$,

- (a) $\frac{\|w_u(t_2)\|}{\|w_v(t_2)\|} > \frac{1}{(c - \epsilon) \cos(\delta)} \Rightarrow \lim_{t \rightarrow \infty} \frac{\|w_u(t)\|}{\|w_v(t)\|} = 1$
- (b) $\frac{\|w_u(t_2)\|}{\|w_v(t_2)\|} < \frac{\cos(\delta)}{c + \epsilon} \Rightarrow \lim_{t \rightarrow \infty} \frac{\|w_u(t)\|}{\|w_v(t)\|} = 0$
- (a) $\frac{\|w_u(t_2)\|}{\|w_v(t_2)\|} > \frac{1}{(c - \epsilon) \cos(\delta)} \Rightarrow \lim_{t \rightarrow \infty} \frac{\|w_u(t)\|}{\|w_v(t)\|} = 1$
Using Equation (14),

$$\frac{d\|w_u(t)\|}{dt} = \frac{d\|w_v(t)\|}{dt} = \frac{r_{w_u} L(w(t)) - r_{w_v} L(w(t))}{\|w_u(t)\|^2} \quad (34)$$

Using the equation above, we can say for $t > t_1$,

$$\begin{aligned} \frac{d \frac{\|w_u(t)\|}{\|w_v(t)\|}}{dt} &= \frac{\|w_v(t)\| \frac{d\|w_u(t)\|}{dt} - \|w_u(t)\| \frac{d\|w_v(t)\|}{dt}}{\|w_v(t)\|^2} \\ &= \frac{\|w_u(t)\|}{\|w_v(t)\|} (\|w_u(t)\| \kappa_{w_u} L(w(t)) - \|w_v(t)\| \kappa_{w_v} L(w(t))) \cos(\theta) - \frac{1}{c} \\ &= (\|w_u(t)\| \kappa_{w_u} L(w(t)) - \|w_v(t)\| \kappa_{w_v} L(w(t))) \frac{\|w_u(t)\|}{\|w_v(t)\|} \cos(\theta) - \frac{1}{c} \end{aligned} \quad (35)$$

In this case, using Equation (35), we can see $\frac{d \frac{\|w_u(t)\|}{\|w_v(t)\|}}{dt} > 0$ at t_2 . Thus, $\frac{\|w_u(t)\|}{\|w_v(t)\|}$ always remains greater than $\frac{1}{(c) \cos(\theta)}$ and keeps on increasing. Let's denote $\frac{\|w_u(t_2)\|}{\|w_v(t_2)\|}$ by R_1 . Then, for $t > t_2$,

$$\frac{d \frac{\|w_u(t)\|}{\|w_v(t)\|}}{dt} = \cos(\theta) \frac{1}{c} (\|w_u(t)\| \kappa_{w_u} L(w(t)) - \|w_v(t)\| \kappa_{w_v} L(w(t)))$$

As $\|w_u\| \neq 1$, therefore using Equation (14), we can say $\int_{t_2}^{R_1} (\|w_u(t)\| \kappa_{w_u} L(w(t)) - \|w_v(t)\| \kappa_{w_v} L(w(t))) dt \neq 1$. Thus, integrating both the sides of the equation above from t_2 to R_1 , we get

$$\int_{t_2}^{R_1} \frac{d \frac{\|w_u(t)\|}{\|w_v(t)\|}}{dt} dt = 1$$

Thus $\lim_{t \rightarrow 1} \frac{\|w_u(t)\|}{\|w_v(t)\|} = 1$.

(b) $\frac{\|w_u(t_2)\|}{\|w_v(t_2)\|} < \frac{\cos(\theta)}{c+}$ $\Rightarrow \lim_{t \rightarrow 1} \frac{\|w_u(t)\|}{\|w_v(t)\|} = 0$.

Using Equation (34), we can say for $t > t_1$,

$$\begin{aligned} \frac{d \frac{\|w_u(t)\|}{\|w_v(t)\|}}{dt} &= \frac{\|w_v(t)\| \frac{d\|w_u(t)\|}{dt} - \|w_u(t)\| \frac{d\|w_v(t)\|}{dt}}{\|w_v(t)\|^2} \\ &= \frac{\|w_u(t)\|}{\|w_v(t)\|} (\|w_u(t)\| \kappa_{w_u} L(w(t)) - \|w_v(t)\| \kappa_{w_v} L(w(t))) \cos(\theta) \\ &= (\|w_u(t)\| \kappa_{w_u} L(w(t)) - \|w_v(t)\| \kappa_{w_v} L(w(t))) \frac{\|w_u(t)\|}{\|w_v(t)\|} \cos(\theta) \\ &= (\|w_u(t)\| \kappa_{w_u} L(w(t)) - \|w_v(t)\| \kappa_{w_v} L(w(t))) \frac{\|w_u(t)\|}{\|w_v(t)\|} (c+ \theta) \cos(\theta) \end{aligned} \quad (36)$$

In this case, using Equation (36), we can see $\frac{d \frac{\|w_u(t)\|}{\|w_v(t)\|}}{dt} < 0$ at t_2 . Thus, $\frac{\|w_u(t)\|}{\|w_v(t)\|}$ always remains smaller than $\frac{\cos(\theta)}{c+}$ and keeps on decreasing. Now, let's say $\lim_{t \rightarrow 1} \frac{\|w_u(t)\|}{\|w_v(t)\|} > 0$. This means that $\frac{\|w_u(t)\|}{\|w_v(t)\|} > \frac{\cos(\theta)}{c+}$, for some $\theta > 0$. Also, let's denote $\frac{\|w_u(t_2)\|}{\|w_v(t_2)\|}$ by R_2 . Then we can say

$$\frac{d \frac{\|w_u(t)\|}{\|w_v(t)\|}}{dt} = (\cos(\theta) - (c+ \theta)) (\|w_u(t)\| \kappa_{w_u} L(w(t)) - \|w_v(t)\| \kappa_{w_v} L(w(t)))$$

As $\nu \rightarrow 1$, therefore using Equation (14), we can say $\int_{t_2}^{t_1} \frac{dw_u(t)}{dw_v(t)} dt = 1$. Thus, integrating both the sides of the equation above from t_2 to t_1 , we get

$$\int_{t_2}^{t_1} \frac{dw_u(t)}{dw_v(t)} dt = 1$$

This is not possible as $\frac{dw_u(t)}{dw_v(t)}$ is lower bounded by 0. Thus $\lim_{t \rightarrow 1} \frac{dw_u(t)}{dw_v(t)} = 0$.

Now, as α and β tend to 0, the length of the interval of stability $[\frac{\cos(\alpha)}{c+}; \frac{1}{(c-)\cos(\alpha)}]$ shrinks to zero, around the point $\frac{1}{c}$. Thus, either $\frac{dw_u(t)}{dw_v(t)}$ moves out of the interval of stability and converges to either 0 or 1, or it always remains within the interval of stability and converges to $\frac{1}{c}$. ■

D.1.2. GRADIENT DESCENT

Theorem Consider two nodes u and v in the network with $k_u, k_v > 0; \lim_{t \rightarrow 1} kw_u(t)k = 1$ and $\lim_{t \rightarrow 1} kw_v(t)k = 1$. Let $\frac{k_u}{k_v}$ be denoted by γ . Under assumptions (B1)-(B3) for gradient descent, for EWN $\lim_{t \rightarrow 1} \frac{kw_u(t)k}{kw_v(t)k}$ is either 0; $\frac{1}{c}$ or 1

Proof Using Theorem 3 and the fact that $\lim_{t \rightarrow 1} \frac{kr_{w_u}L(w(t))k}{kr_{w_v}L(w(t))k} = c$, for any $0 < \epsilon < c$ and $0 < \delta < 2$, there exists a time t_1 , such that for $t > t_1$, the following hold

- (i) $\frac{kr_{w_u}L(w(t))k}{kr_{w_v}L(w(t))k} \geq 2 [c - \epsilon; c + \epsilon]$ (ii) $\frac{w_u(t)}{kw_u(t)k} > \frac{r_{w_u}L(w(t))}{kr_{w_u}L(w(t))k} \cos(\delta)$
- (iii) $\frac{w_v(t)}{kw_v(t)k} > \frac{r_{w_v}L(w(t))}{kr_{w_v}L(w(t))k} \cos(\delta)$.

Now, we will provide the proof of part (iii) of Proposition 5, i.e, for EWN, if at some time $t_2 > t_1$,

(a) $\frac{kw_u(t_2)k}{kw_v(t_2)k} > \frac{1}{(c-)\cos(\delta)} \Rightarrow \lim_{t \rightarrow 1} \frac{kw_u(t)k}{kw_v(t)k} = 1$

(b) $\frac{kw_u(t_2)k}{kw_v(t_2)k} < \frac{\cos(\delta)}{c+} \Rightarrow \lim_{t \rightarrow 1} \frac{kw_u(t)k}{kw_v(t)k} = 0$

(a) $\frac{kw_u(t_2)k}{kw_v(t_2)k} > \frac{1}{(c-)\cos(\delta)} \Rightarrow \lim_{t \rightarrow 1} \frac{kw_u(t)k}{kw_v(t)k} = 1$

Using Equation (18) and part 1 of the Proposition, we can say

$$\begin{aligned} \frac{kw_u(t_2+1)k}{kw_v(t_2+1)k} &= \frac{kw_u(t_2)k + (\delta) \cos(\delta) kw_u(t_2)k^2 kr_{w_u}L(w(t_2))k}{kw_v(t_2)k + (\delta) kw_v(t_2)k^2 kr_{w_v}L(w(t_2))k} \\ &= \frac{kw_u(t_2)k}{kw_v(t_2)k} \frac{1 + \cos(\delta) (\delta) kw_u(t_2)k kr_{w_u}L(w(t_2))k}{1 + (\delta) kw_v(t_2)k kr_{w_v}L(w(t_2))k} \\ &= \frac{kw_u(t_2)k}{kw_v(t_2)k} \end{aligned}$$

Thus, $\frac{kw_u(t)k}{kw_v(t)k}$ keeps on increasing for $t > t_2$. It can either diverge to infinity or converge to a finite value. If it converges to a finite value, then by Stolz Cesaro theorem,

$$\lim_{t \rightarrow \infty} \frac{kw_u(t)k}{kw_v(t)k} = \lim_{t \rightarrow \infty} \frac{kw_u(t)k^2kr_{w_u}L(w(t))k}{kw_v(t)k^2kr_{w_v}L(w(t))k}$$

However, this is not possible as $\frac{kw_u(t)k}{kw_v(t)k} > \frac{1}{c}$ for every $t > t_2$. Thus, $\frac{kw_u(t)k}{kw_v(t)k}$ diverges to infinity.

(b) $\frac{kw_u(t_2)k}{kw_v(t_2)k} < \frac{\cos(\cdot)}{c}$ $\Rightarrow \lim_{t \rightarrow \infty} \frac{kw_u(t)k}{kw_v(t)k} = 0$

Using Equation (18) and part 1 of the Proposition, we can say

$$\begin{aligned} \frac{kw_u(t_2+1)k}{kw_v(t_2+1)k} &= \frac{kw_u(t_2)k + (t_2)kw_u(t_2)k^2kr_{w_u}L(w(t_2))k}{kw_v(t_2)k + (t_2)\cos(\cdot)kw_v(t_2)k^2kr_{w_v}L(w(t_2))k} \\ &= \frac{kw_u(t_2)k}{kw_v(t_2)k} \frac{1 + (t_2)kw_u(t_2)kkr_{w_u}L(w(t_2))k}{1 + (t_2)\cos(\cdot)kw_v(t_2)kkr_{w_v}L(w(t_2))k} \\ &= \frac{kw_u(t_2)k}{kw_v(t_2)k} \end{aligned}$$

Thus, $\frac{kw_u(t)k}{kw_v(t)k}$ keeps on decreasing for $t > t_2$. As it is always greater than zero, it must converge. Therefore, by Stolz Cesaro Theorem,

$$\lim_{t \rightarrow \infty} \frac{kw_u(t)k}{kw_v(t)k} = \lim_{t \rightarrow \infty} \frac{kw_u(t)k^2kr_{w_u}L(w(t))k}{kw_v(t)k^2kr_{w_v}L(w(t))k}$$

For $\frac{kw_u(t)k}{kw_v(t)k} < \frac{1}{c}$, this can only be satisfied when $\lim_{t \rightarrow \infty} \frac{kw_u(t)k}{kw_v(t)k} = 0$.

Now, as ϵ and δ tend to 0, the length of the interval of stability $[\frac{\cos(\cdot)}{c+\epsilon}; \frac{1}{(c-\epsilon)\cos(\cdot)}]$ shrinks to zero, around the point $\frac{1}{c}$. Thus, either $\frac{kw_u(t)k}{kw_v(t)k}$ moves out of the interval of stability and converges to either 0 or 1, or it always remains within the interval of stability and converges to $\frac{1}{c}$. ■

D.2. Standard Weight Normalization

D.2.1. GRADIENT FLOW

Theorem Consider two nodes u and v in the network with $k_{g_u}k < k_{g_v}k < 0$; $\lim_{t \rightarrow \infty} kw_u(t)k = 1$ and $\lim_{t \rightarrow \infty} kw_v(t)k = 1$. Let $\frac{k_{g_u}k}{k_{g_v}k}$ be denoted by c . Under assumptions (A1), (A2) for gradient flow, for SWN $\lim_{t \rightarrow \infty} \frac{kw_u(t)k}{kw_v(t)k} = c$

Proof From Theorem 3, we can say, for both u and v , weights and gradients converge in opposite directions.

Consider a time t_1 , such that for any $t > t_1$,

- $r_{w_u}L(w(t))$ and $w_u(t)$ at most make an angle with each other
- $r_{w_v}L(w(t))$ and $w_v(t)$ at most make an angle with each other

Then, using Equation (24), we can say for any $t \geq t_2$,

$$\begin{aligned} \|w_u(t)\| &\leq \|w_u(t_2)\| + \cos(\alpha) \int_{t_2}^t \int_{k=t_2}^k (k)^{r_{w_u}} L(w(k)) dk \\ \|w_v(t)\| &\leq \|w_v(t_2)\| + \cos(\beta) \int_{t_2}^t \int_{k=t_2}^k (k)^{r_{w_v}} L(w(k)) dk \end{aligned}$$

Using the above equations, we can say, for $t \geq t_2$,

$$\begin{aligned} \frac{\|w_u(t)\|}{\|w_v(t)\|} &\leq \frac{\|w_u(t_2)\| + \cos(\alpha) \int_{t_2}^t \int_{k=t_2}^k (k)^{r_{w_u}} L(w(k)) dk}{\|w_v(t_2)\| + \cos(\beta) \int_{t_2}^t \int_{k=t_2}^k (k)^{r_{w_v}} L(w(k)) dk} \\ \frac{\|w_u(t)\|}{\|w_v(t)\|} &\geq \frac{\|w_u(t_2)\| + \cos(\alpha) \int_{t_2}^t \int_{k=t_2}^k (k)^{r_{w_u}} L(w(k)) dk}{\|w_v(t_2)\| + \cos(\beta) \int_{t_2}^t \int_{k=t_2}^k (k)^{r_{w_v}} L(w(k)) dk} \end{aligned}$$

We know that both integrals diverge as $t \rightarrow \infty$ and $\|w_u(t)\|, \|w_v(t)\| \rightarrow 1$, $\lim_{t \rightarrow \infty} \frac{\|w_u(t)\|}{\|w_v(t)\|} = c$ and $\lim_{t \rightarrow \infty} \frac{\int_{t_2}^t \int_{k=t_2}^k (k)^{r_{w_u}} L(w(k)) dk}{\int_{t_2}^t \int_{k=t_2}^k (k)^{r_{w_v}} L(w(k)) dk}$ exist. Taking limit $t \rightarrow \infty$ on both the equations and using the Integral form of Stolz-Cesaro theorem, we get

$$\begin{aligned} \liminf_{t \rightarrow \infty} \frac{\|w_u(t)\|}{\|w_v(t)\|} &= \frac{\cos(\alpha) \int_{t_2}^t \int_{k=t_2}^k (k)^{r_{w_u}} L(w(k)) dk}{\cos(\beta) \int_{t_2}^t \int_{k=t_2}^k (k)^{r_{w_v}} L(w(k)) dk} \\ \limsup_{t \rightarrow \infty} \frac{\|w_u(t)\|}{\|w_v(t)\|} &= \frac{\int_{t_2}^t \int_{k=t_2}^k (k)^{r_{w_u}} L(w(k)) dk}{\cos(\beta) \int_{t_2}^t \int_{k=t_2}^k (k)^{r_{w_v}} L(w(k)) dk} \end{aligned}$$

As this holds for any $\epsilon > 0$, therefore

$$\lim_{t \rightarrow \infty} \frac{\|w_u(t)\|}{\|w_v(t)\|} = c$$

■

D.2.2. GRADIENT DESCENT

Theorem Consider two nodes u and v in the network with $\|g_u\|, \|g_v\| > 0$; $\lim_{t \rightarrow \infty} \|w_u(t)\| = 1$ and $\lim_{t \rightarrow \infty} \|w_v(t)\| = 1$. Let $\frac{\|g_u\|}{\|g_v\|}$ be denoted by c . Under assumptions (B1) - (B3) for gradient descent, for SWIN, $\lim_{t \rightarrow \infty} \frac{\|w_u(t)\|}{\|w_v(t)\|} = c$

Proof For ease of notation, we will denote the two nodes as u and v . From Theorem 3, we can say, for both u and v , weights and gradients converge in opposite directions. Now from Equation (28), we can say

$$w_u(t) = w_u(0) - \int_{k=0}^{t-1} \frac{v_u(k) \cdot r_{w_u} L(w(k))}{\|v_u(k)\|} dk$$

$$s(t) = s(0) \prod_{k=0}^{t-1} \frac{v_s(k) r_{w_s} L(w(k))}{k v_s(k)}$$

Now, $s(t)$ either diverges to ∞ or 0 . In both the cases, it is a strictly monotonic sequence for large enough t . Also $\lim_{t \rightarrow \infty} \frac{s(t+1)}{s(t)}$ exists. Therefore, using Stolz-Cesaro Theorem, we can say

$$\lim_{t \rightarrow \infty} \frac{kw_u(t)}{kw_s(t)} = c$$

■

Appendix E. Proof of Proposition 5

Proposition Consider two nodes u and v in the network such that $k_{uv} > 0$ and $kw_u(t), kw_v(t) \leq 1$. Let $\frac{k_{uv}}{k_{vv}}$ be denoted by c . Consider any ϵ such that $0 < \epsilon < c$ and $0 < \delta < \frac{\epsilon}{2}$. Then, the following holds:

(i) There exists a time t_1 , such that for all $t > t_1$ both SWN and EWN trajectories have the following properties:

$$(a) \frac{kr_{w_u} L(w(t))}{kr_{w_v} L(w(t))} \in [c - \delta, c + \delta] \quad (b) \frac{kw_u(t)}{kw_v(t)} > \frac{r_{w_u} L(w(t))}{r_{w_v} L(w(t))} \cos(\delta)$$

$$(c) \frac{kw_v(t)}{kw_u(t)} > \frac{r_{w_v} L(w(t))}{r_{w_u} L(w(t))} \cos(\delta).$$

(ii) for SWN, $\lim_{t \rightarrow \infty} \frac{kw_u(t)}{kw_v(t)} = c$

(iii) for EWN, if at some time $t_2 > t_1$,

$$(a) \frac{kw_u(t_2)}{kw_v(t_2)} > \frac{1}{(c - \delta) \cos(\delta)} \Rightarrow \lim_{t \rightarrow \infty} \frac{kw_u(t)}{kw_v(t)} = 1$$

$$(b) \frac{kw_u(t_2)}{kw_v(t_2)} < \frac{\cos(\delta)}{c + \delta} \Rightarrow \lim_{t \rightarrow \infty} \frac{kw_u(t)}{kw_v(t)} = 0$$

Proof (i) It follows from Theorem 3 and the fact that $\lim_{t \rightarrow \infty} \frac{kr_{w_u} L(w(t))}{kr_{w_v} L(w(t))} = c$.

(ii) Proof provided in Appendix D.2

(iii) Proof provided in Appendix D.1

■

Appendix F. Proof of Proposition 6

Proposition Consider a linear model over \mathbb{R}^d given by $f(x) = w^T x$, where each w_i is further reparameterized as $w_i = \sum_{j=1}^d a_{ij} z_j$. Consider a dataset consisting of a single data point 0 , that is labelled as $+1$. According to the initialization of w , define a relation R on $\{1, \dots, d\}$, given by $i R j$ if $w_i(0) z_i = w_j(0) z_j$. Then, R is an equivalence relation on $\{1, \dots, d\}$. Let these equivalent sets be denoted by I_1, I_2, \dots, I_k . Define a total order on these sets given by $I_a > I_b$ if $\exists i \in I_a, j \in I_b$ such that $w_i(0) z_i > w_j(0) z_j$. Let the maximum set according to this order be denoted by I_{max} . Then, for gradient flow on exponential loss, the following holds

(i) For any $i \in I_{max}$, $\lim_{t \rightarrow \infty} w_i(t) = 1$

(ii) For $i, j \in I$, $\frac{w_i(t)}{w_j(t)} = \frac{z_j}{z_i}$.

(iii) For any $i \in I$, $\lim_{t \rightarrow \infty} w_i(t) = \frac{1}{w_i(0)} \frac{z_i}{w_j(0)z_j}$, where j is any element in I .

Proof In this case, the loss is given by

$$L(w) = e^{-w^T z}$$

As each w_i is further exponentially reparameterized, therefore, using Theorem 1,

$$\frac{dw_i}{dt} = w_i^2 z_i e^{-w^T z} \tag{37}$$

Using this equation, we can say, for any pair of indices

$$\frac{d\left(\frac{z_i}{w_j} \frac{z_j}{w_i}\right)}{dt} = 0$$

Thus, for any pair of indices i, j

$$\frac{z_i}{w_j(t)} \frac{z_j}{w_i(t)} = \frac{z_i}{w_j(0)} \frac{z_j}{w_i(0)} \tag{38}$$

Thus, we can say for every equivalent set I_a , then as $\frac{z_j}{w_i(0)} \frac{z_i}{w_j(0)} = 0$, therefore, at any time t ,

$$\frac{w_i(t)}{w_j(t)} = \frac{z_j}{z_i}$$

Now, clearly the gradient flow stops only when $w^T z = 1$. This means at least one of the w_i must tend to 1 (Notice from Equation (37) that w_i always goes up). For two equivalent sets I_a and I_b , such that $I_a > I_b$, we can say, for $i \in I_a$ and $j \in I_b$, the RHS of Equation (38) is positive. Thus, it is not possible for w_j to go to 1, otherwise the quantity on the LHS will be negative when $w_i = 1$. Using this argument, we can say, only for $i \in I$, w_i tend to 1. Then, for $i \in I$, considering $j \in I$ and using Equation (38),

$$\lim_{t \rightarrow \infty} w_i(t) = \frac{1}{w_i(0)} \frac{z_i}{w_j(0)z_j}$$

■

Appendix G. Proof of Theorem 7

Theorem For Exponential Weight Normalization, under assumption (A1), the following hold for $t > t_0$ in case of gradient flow

- (i) $\|w(t)\|$ grows with $\text{as } o((\log t)^{\frac{1}{L}})$
- (ii) $L(t)$ goes down with $\text{as } O\left(\frac{1}{t}\right)$

Proof Using Equation (13), we can say

$$\frac{1}{L(t)} = \frac{1}{L(t_0)} + \frac{L^2}{k}(t - t_0)$$

Thus $L(t)$ goes down with t as $\frac{1}{t}$. Now, we know, $L(w) = \sum_i y_i (w; x_i)$. Clearly, as $(w; x_i)$ is smooth, therefore $L(w)$ attains a minima over the compact set $\|w\| = 1$. Using homogeneity of $L(w)$ and the fact that $L(w)$ goes down as $\frac{1}{t}$, we can say $\|w(t)\|$ grows with t as $(\log t)^{\frac{1}{c}}$. ■

Appendix H. Proof of Theorem 8

Theorem For Exponential Weight Normalization, under Assumptions (B1)-(B4), $0 < c < 1$ and $\lim_{t \rightarrow \infty} \frac{kr(t+1) - r(t)k}{g(t+1) - g(t)} = 0$, the following hold

- (i) $\|w(t)\|$ asymptotically grows with as $(\log d(t))^{\frac{1}{c}}$
- (ii) $L(w(t))$ asymptotically goes down with as $\frac{1}{d(t)(\log d(t))^2}$.

First, we will establish rates for gradient flow to elucidate the proof technique and then go to the case of gradient descent.

H.1. Gradient Flow

Although the asymptotic convergence rates for smooth homogeneous neural nets have been established in [Lyu and Li \(2020\)](#), the proof technique becomes easier to understand for smooth homogeneous nets, without weight normalization. In this case, we will use an assumption (A3) that $\lim_{t \rightarrow \infty} \frac{w(t)}{\|w(t)\|}$ exists.

H.1.1. UNNORMALIZED NETWORK

Theorem For Unnorm, under Assumptions (A1)-(A3) for gradient flow, $0 < c < 1$ and $\lim_{t \rightarrow \infty} \frac{k \frac{dr(t)}{dt} k}{g^0(t)} = 0$, the following hold

1. $\|w(t)\|$ asymptotically grows at $(\log t)^{\frac{1}{c}}$
2. $L(w(t))$ asymptotically goes down at the rate of $\frac{1}{t(\log t)^2}$.

Proof Consider $w = g(t)w + r(t)$, where $\lim_{t \rightarrow \infty} \frac{kr(t)k}{g(t)} = 0$ and $r(t) > 0 = 0$. Now, we make an additional assumption that $\lim_{t \rightarrow \infty} \frac{k \frac{dr(t)}{dt} k}{g^0(t)} = 0$. This basically avoids any oscillations in $r(t)$ for target, where it can have a higher derivative, but the value may be bounded. Now, we know

$$\frac{dw(t)}{dt} = \sum_{i=1}^n y_i e^{-y_i (w(t); x_i)} y_i r_w (w(t); x_i) \tag{39}$$

Now, we know $\frac{dw(t)}{dt} \neq 0$ for any finite t , otherwise w won't change and L can't converge to 0. Thus, for all t , we can say

$$\frac{\sum_{i=1}^m e^{-y_i(w(t);x_i)} y_i r_w(w; x_i)}{\sum_{i=1}^m e^{-y_i(w(t);x_i)} y_i r_w(w; x_i)} = 1$$

Taking limit $t \rightarrow \infty$ on both the sides, we get

$$\lim_{t \rightarrow \infty} \frac{\sum_{i=1}^m e^{-y_i(w(t);x_i)} y_i r_w(w; x_i)}{\sum_{i=1}^m e^{-y_i(w(t);x_i)} y_i r_w(w; x_i)} = 1$$

Now, we know

$$\frac{dw(t)}{dt} = g^0(t)w + \frac{dr(t)}{dt}$$

$$\sum_{i=1}^m e^{-y_i(w(t);x_i)} y_i r_w(w; x_i) = \sum_{i=1}^m e^{-y_i g^0(t)w + \frac{r(t)}{g(t)}; x_i} (y_i g(t))^{L-1} r_w(w; x_i) + \frac{r(t)}{g(t)}; x_i$$

Thus, we can say

$$\sum_{i=1}^m e^{-y_i(w(t);x_i)} y_i r_w(w; x_i) = L(w(t))g(t)^{L-1} \sum_{i=1}^m \frac{e^{-y_i(w(t);x_i)}}{L(w(t))} y_i r_w(w; x_i) + \frac{r(t)}{g(t)}; x_i$$

Now, for large enough t , $\sum_{i=1}^m \frac{e^{-y_i(w(t);x_i)}}{L(w(t))} y_i r_w(w; x_i) + \frac{r(t)}{g(t)}; x_i$ is bounded as, using Euler's homogeneous theorem, we can say

$$\lim_{t \rightarrow \infty} \sum_{i=1}^m \frac{e^{-y_i(w(t);x_i)}}{L(w(t))} y_i r_w(w; x_i) + \frac{r(t)}{g(t)}; x_i = \lim_{t \rightarrow \infty} L \sum_{i=1}^m \frac{e^{-y_i(w(t);x_i)}}{L(w(t))} (w; x_i)$$

Thus, it's a convex combination of positive defined terms and hence bounded. Thus, we can say

$$k_1 \lim_{t \rightarrow \infty} \frac{\sum_{i=1}^m e^{-y_i(w(t);x_i)} y_i r_w(w; x_i)}{L(w(t))g(t)^{L-1}} \leq k_2 \tag{40}$$

where k_1 and k_2 are some constants. Also, by the assumption

$$\lim_{t \rightarrow \infty} \frac{\sum_{i=1}^m e^{-y_i(w(t);x_i)} y_i r_w(w; x_i)}{g^0(t)} = 1$$

Thus, we can say

$$k_1 \lim_{t \rightarrow \infty} \frac{g^0(t)}{L(w(t))g(t)^{L-1}} \leq k_2$$

Now, as for large enough t , for all i that satisfy $(w; x_i) = 0$ and any ϵ satisfying $0 < \epsilon < 1$, we can say

$$\sum_{i=1}^m \frac{r(t)}{g(t)}; x_i \leq \epsilon$$

therefore, for large enough t we can say

$$c_1 e^{-g(t)^L} < L(w(t)) < c_2 e^{-g(t)^L}$$

where c_1 and c_2 are some constants. Using the above equations, we can say

$$\lim_{t \rightarrow \infty} \frac{g^0(t)}{e^{-g(t)^L} g(t)^{L-1}} = k_1 c_1 \quad (41)$$

$$\lim_{t \rightarrow \infty} \frac{g^0(t)}{e^{-g(t)^L} g(t)^{L-1}} = k_2 c_2 \quad (42)$$

Now, in Equation 41, multiplying numerator and denominator by $g(t)^{L-1}$ and denoting $g(t)^{L-1}$ by $h_1(t)$, we get

$$\lim_{t \rightarrow \infty} \frac{h_1^0(t)}{e^{-h_1(t)} h_1(t)^{\frac{2}{L}}}$$

where $\frac{2}{L}$ is some constant. This leads us to conclude

$$\lim_{t \rightarrow \infty} \frac{h_1(t)}{\log(t)} = 1 \Rightarrow \lim_{t \rightarrow \infty} \frac{g(t)^L}{\log(t)} = 1$$

Similarly, in Equation 42, multiplying numerator and denominator by $g(t)^{L-1}$ and denoting $g(t)^{L-1}$ by $h_2(t)$, we get

$$\lim_{t \rightarrow \infty} \frac{h_2^0(t)}{e^{-h_2(t)} h_2(t)^{\frac{2}{L}}}$$

where $\frac{2}{L}$ is some constant. This leads us to conclude

$$\lim_{t \rightarrow \infty} \frac{h_2(t)}{\log(t)} = 1 \Rightarrow \lim_{t \rightarrow \infty} \frac{g(t)^L}{\log(t)} = 1$$

As $\frac{2}{L}$ can be chosen to be arbitrarily small, we can conclude

$$\lim_{t \rightarrow \infty} \frac{g(t)^L}{\log(t)} = 1$$

Substituting this in Equation 40, we get that loss asymptotically goes down at $\frac{1}{t(\log t)^2}$. ■

H.1.2. EXPONENTIAL WEIGHT NORMALIZATION

Theorem For EWN, under Assumptions (A1)-(A3) for gradient flow $\lim_{t \rightarrow \infty} \frac{d}{dt} \|k\| = 0$, the following hold

1. $\|k(t)\|$ asymptotically grows at $\log(t)^{\frac{1}{L}}$
2. $L(w(t))$ asymptotically goes down at the rate of $\frac{1}{t(\log t)^2}$.

Proof Consider $w = g(t)w + r(t)$, where $\lim_{t \rightarrow \infty} \frac{kr(t)k}{g(t)} = 0$ and $r(t) > w = 0$. Now, we make an additional assumption that $\lim_{t \rightarrow \infty} \frac{k \frac{dr(t)}{dt} k}{g^0(t)} = 0$.

In this case,

$$\frac{dL(w(t))}{dw} = \sum_{i=1}^n e^{-y_i(w(t); x_i)} y_i r_w(w(t); x_i)$$

However, in this case, for a node

$$\frac{dw_u(t)}{dt} = k w_u(t) k^2 \frac{dL(w(t))}{dw_u}$$

Consider a vector $a(t)$ of equal dimension as w , and its components corresponding to a node given by $a_u(t) = k w_u(t) k^2 \frac{dL(w(t))}{dw_u}$. Now as we know w converges in direction w , therefore, using the update equation above, we can say

$$\lim_{t \rightarrow \infty} \frac{k \frac{dw(t)}{dt} k}{g(t)^2 k a(t) k} = 1$$

Using the update equation for $\frac{dL(w(t))}{dw}$, which is the same as Equation 39, and using the same arguments as for Unnorm in Appendix H.1.1, we can say

$$k_1 \lim_{t \rightarrow \infty} \frac{k \frac{dL(w(t))}{dw} k}{L(w(t)) g(t)^{L-1}} \leq k_2 \tag{43}$$

where k_1 and k_2 are some constants. Now, using the expression $a(t)$, we can say

$$k_3 \lim_{t \rightarrow \infty} \frac{k a(t) k}{L(w(t)) g(t)^{L-1}} \leq k_4$$

where k_3 and k_4 are some constants. Now, from the assumption, we know

$$\lim_{t \rightarrow \infty} \frac{k \frac{dw(t)}{dt} k}{g^0(t)} = 1$$

Using the equations above, we can say

$$k_3 \frac{g^0(t)}{L(w(t)) g(t)^{L+1}} \leq k_4 \tag{44}$$

Using similar reasoning as for Unnorm in Appendix H.1.1, we can say, for a large enough

$$c_1 e^{-g(t)^{L+1}} \leq L(w(t)) \leq c_2 e^{-g(t)^L}$$

where c_1 and c_2 are some constants. Substituting this in equation above, we get

$$\lim_{t \rightarrow \infty} \frac{g^0(t)}{e^{-g(t)^{L+1}} g(t)^{L+1}} \leq k_3 c_1 \tag{45}$$

$$\lim_{t \rightarrow \infty} \frac{g^0(t)}{e^{-g(t)^L} g(t)^{L+1}} \leq k_4 c_2 \tag{46}$$

Now, using similar arguments as for Unnorm in Appendix H.1.1, we can say

$$\lim_{t \rightarrow \infty} \frac{g(t)^L}{\log(t)} = 1$$

Substituting this in Equation 44, we get that loss asymptotically goes down at $\frac{1}{t(\log t)^2}$. ■

H.2. Gradient Descent

Theorem For Exponential Weight Normalization, under Assumptions (B1)-(B4), $0 < c < 1$ and $\lim_{t \rightarrow \infty} \frac{kr(t+1)}{g(t+1)} = 0$, the following hold

- (i) $\|w(t)\|$ asymptotically grows with as $(\log d(t))^{\frac{1}{L}}$
- (ii) $L(w(t))$ asymptotically goes down with as $\frac{1}{d(t)(\log d(t))^2}$.

Proof Consider $w = g(t)v + r(t)$, where $\lim_{t \rightarrow \infty} \frac{kr(t)}{g(t)} = 0$ and $r(t) \perp v = 0$. Now, we make additional assumptions that $\lim_{t \rightarrow \infty} \frac{kr(t+1)}{g(t+1)} = 0$.

Consider a node u in the network that has $\|w_u\| > 0$. The update equations for $v_u(t)$ and $\theta_u(t)$ are given by

$$\begin{aligned} v_u(t+1) &= v_u(t) + \theta_u(t) \frac{v_u(t) \cdot r_{w_u} L(w(t))}{\|v_u(t)\|} \\ \theta_u(t+1) &= \theta_u(t) e^{-\theta_u(t)} \left(1 + \frac{v_u(t) \cdot v_u(t)}{\|v_u(t)\|^2} \cdot r_{w_u} L(w(t)) \right) \end{aligned}$$

Now, we will first estimate $\frac{\|v_u(t+1)\|}{\|v_u(t)\|} = e^{-\theta_u(t)} \frac{\|v_u(t)\|}{\|v_u(t)\|} = e^{-\theta_u(t)}$. Let $\theta_u(t)$ denote $\theta_u(t)$ and $\phi_u(t)$ denote the angle between $v_u(t)$ and $r_{w_u} L(w(t))$. We know $\lim_{t \rightarrow \infty} \theta_u(t) = 0$ and $\lim_{t \rightarrow \infty} \phi_u(t) = 0$. Now, rewriting update equations in terms of these symbols, we get

$$e^{-\theta_u(t+1)} = e^{-\theta_u(t)} e^{-\theta_u(t) \cos(\phi_u(t))}$$

$$v_u(t+1) = v_u(t) + \theta_u(t) \sin(\phi_u(t)) \frac{r_{w_u} L(w(t))_{\perp v_u(t)}}{\|r_{w_u} L(w(t))_{\perp v_u(t)}\|}$$

where $r_{w_u} L(w(t))_{\perp v_u(t)}$ denotes the component of $r_{w_u} L(w(t))$ perpendicular to $v_u(t)$. Now, using these equations we can say

$$\begin{aligned} \frac{\|v_u(t+1)\|}{\|v_u(t)\|} &= e^{-\theta_u(t)} e^{-\theta_u(t) \cos(\phi_u(t))} \frac{\|v_u(t)\|}{\|v_u(t)\|} + \frac{\theta_u(t) \sin(\phi_u(t))}{\|v_u(t)\|} \frac{\|r_{w_u} L(w(t))_{\perp v_u(t)}\|}{\|r_{w_u} L(w(t))_{\perp v_u(t)}\|} \\ &= e^{-\theta_u(t)} e^{-\theta_u(t) \cos(\phi_u(t))} + \frac{\theta_u(t) \sin(\phi_u(t))}{\|v_u(t)\|} \frac{\|r_{w_u} L(w(t))_{\perp v_u(t)}\|}{\|r_{w_u} L(w(t))_{\perp v_u(t)}\|} \end{aligned} \quad (47)$$

Now as $\lim_{t \rightarrow \infty} \frac{\|v_u(t)\|}{\|v_u(t+1)\|} = 1$, therefore we can say

$$\lim_{t \rightarrow \infty} \frac{e^{-u(t)\cos(u(t))} \frac{\|v_u(t)\|}{\|v_u(t+1)\|}}{u(t)\cos(u(t))} = 1$$

Now, as $\|v_u(t)\|$ keeps on increasing during the gradient descent trajectory, therefore we can say $\frac{1}{\|v_u(t)\| \|v_u(t+1)\|} = k$, where $k > 0$ is some constant. Now dividing both sides of Equation (47) by $e^{-u(t)} u(t) \cos(u(t))$ and analyzing the coefficient of the second term on RHS, we get

$$\lim_{t \rightarrow \infty} \frac{e^{-u(t)\cos(u(t))} \sin(u(t))}{\|v_u(t)\| \|v_u(t+1)\| \cos(u(t))} = 0$$

Taking norm on both sides of Equation (47), using Pythagoras theorem and the limits established above, we can say

$$\lim_{t \rightarrow \infty} \frac{k e^{-u(t+1)} \frac{\|v_u(t+1)\|}{\|v_u(t+1)\|} e^{-u(t)} \frac{\|v_u(t)\|}{\|v_u(t)\|} k}{e^{-u(t)} u(t)} = 1$$

Now, we also know

$$\lim_{t \rightarrow \infty} \frac{k e^{-u(t+1)} \frac{\|v_u(t+1)\|}{\|v_u(t+1)\|} e^{-u(t)} \frac{\|v_u(t)\|}{\|v_u(t)\|} k}{g(t+1) g(t)} = k \|v_u\|$$

Now, using equations above and Equation (43), we can say

$$k_1 \lim_{t \rightarrow \infty} \frac{g(t+1) g(t)}{(t)^L (w(t)) g(t)^{L+1}} = k_2$$

where k_1 and k_2 are some constants. Using similar reasoning as for U_{norm} in Appendix H.1.1, we can say

$$c_1 e^{-g(t)^L (t+1)^L (w(t))} = c_2 e^{-g(t)^L (t)^L}$$

where c_1 and c_2 are some constants. Substituting in the equation above, we get

$$\lim_{t \rightarrow \infty} \frac{g(t+1) g(t)}{(t) e^{-g(t)^L (t+1)^L} g(t)^{L+1}} = k_1 c_1$$

$$\lim_{t \rightarrow \infty} \frac{g(t+1) g(t)}{(t) e^{-g(t)^L (t)^L} g(t)^{L+1}} = k_2 c_2$$

These equations govern the rate of $g(t)$ for any (t) that satisfies assumption (A4). Now, to obtain a better closed form, we will need the new assumption (A5), i.e., $(t) = O(\log \frac{1}{\epsilon}^c)$, where $c < 1$.

Now, define a map $d: \mathbb{N} \rightarrow \mathbb{R}$, given by $d(t) = \prod_{s=0}^{t-1} (s)$ and a real analytic function $f(t)$ satisfying $f(d(t)) = g(t)$ for all $t \in \mathbb{N}$ and $\lim_{t \rightarrow \infty} \frac{f(d(t+1))}{f(d(t))} = 1$. We will later verify that the $f(t)$ obtained indeed satisfies this for the given $f(t)$. Substituting in the equations above,

$$k_1 \lim_{t \rightarrow \infty} \frac{f(d(t))}{L(w(t)) f(d(t))^{L+1}} = k_2 \tag{48}$$

$$\lim_{t \rightarrow \infty} \frac{f^Q(d(t))}{e^{-f(d(t))^L} f(d(t))^{L+1}} = k_1 C_1$$

$$\lim_{t \rightarrow \infty} \frac{f^Q(d(t))}{e^{-f(d(t))^L} f(d(t))^{L+1}} = k_2 C_2$$

As $t \rightarrow \infty$, $d(t) \rightarrow \infty$, therefore

$$\lim_{t \rightarrow \infty} \frac{f^Q(t)}{e^{-f(t)^L} f(t)^{L+1}} = k_1 C_1$$

$$\lim_{t \rightarrow \infty} \frac{f^Q(t)}{e^{-f(t)^L} f(t)^{L+1}} = k_2 C_2$$

Now, using similar arguments as in Appendix H.1.1,

$$\lim_{t \rightarrow \infty} \frac{f(t)^L}{\log(t)} = 1$$

Substituting in the Equation 48, we get that the loss goes down at $\frac{1}{d(t)(\log d(t))^2}$.

We also verify that $\lim_{t \rightarrow \infty} \frac{f(d(t+1)) - f(d(t))}{(t)f^Q(d(t))} = 1$ for $(t) = O(\log \frac{1}{L^c})$, where $c < 1$. This can be easily verified by using mean value theorem, and simply verifying $\frac{(t)f^Q(d(t))}{f^Q(d(t))} = 0$. Obtaining the expressions for $f^Q(d(t))$ and $f^Q(d(t))$, we get

$$\lim_{t \rightarrow \infty} \frac{(t)f^Q(d(t))}{f^Q(d(t))} = \lim_{t \rightarrow \infty} \frac{(t) d(t) (\frac{1}{L} - 1) \log d(t)}{d(t) \log(d(t))}$$

As loss goes down at $\frac{1}{d(t)(\log d(t))^2}$, therefore if $(t) = O(\log \frac{1}{L^c})$ for $c < 1$, the above limit tends to 0 as $(t) \rightarrow \infty$. ■

Appendix I. Cross-Entropy Loss

In this section, we will provide the corresponding assumptions and theorems, along with their proofs, for cross-entropy loss.

I.1. Notations

Let k denote the total number of classes. As $(w; x_i)$ is a multidimensional function for multi-class classification, let's denote the j^{th} component of the output by $y_j(w; x_i)$. Also, denote the asymptotic normalized margin for j^{th} class corresponding to i^{th} data point $(w; y_i)$ by γ_{ij} , i.e., $\gamma_{ij} = y_i(w; x_i) - y_j(w; x_i)$. Margin for a data point is defined as $\gamma_i = \min_j y_i - y_j$. The margin for the entire network is defined as $\gamma = \min_i \gamma_i$.

I.2. Assumptions

The assumptions can be broadly divided into loss function/architecture based assumptions and trajectory based assumptions. The loss functions/architecture based assumptions are shared across both gradient flow and gradient descent.

Loss function/Architecture based assumptions

$$1 \quad \sigma(y_i; (\mathbf{w}; \mathbf{x}_i)) = \log \left(1 + \prod_{j \in \mathcal{Y}_i} e^{-(y_i(\mathbf{w}; \mathbf{x}_i) - j(\mathbf{w}; \mathbf{x}_i))} \right)$$

2 $\sigma(\cdot; \mathbf{x})$ is a C^1 function, for a fixed \mathbf{x}

3 $\sigma(\mathbf{w}; \mathbf{x}) = L(\mathbf{w}; \mathbf{x})$, for some $\epsilon > 0$ and $L > 0$

Gradient flow. For gradient flow, we make the following trajectory based assumptions

(A1) There exists a time t_0 such that $L(\mathbf{w}(t_0)) < \log 2$.

(A2) $\lim_{t \rightarrow \infty} \frac{r_{\mathbf{w}} L(\mathbf{w}(t))}{k r_{\mathbf{w}} L(\mathbf{w}(t)) k} := \epsilon$.

Gradient Descent For gradient descent, we require the learning rate to not grow too fast, and a slightly stronger assumption on loss.

(B1) $\lim_{t \rightarrow \infty} L(\mathbf{w}(t)) = 0$

(B2) $\lim_{t \rightarrow \infty} \frac{r_{\mathbf{w}} L(\mathbf{w}(t))}{k r_{\mathbf{w}} L(\mathbf{w}(t)) k} := \epsilon$

(B3) $\lim_{t \rightarrow \infty} (t) k w_u(t) k r_{\mathbf{w}_u} L(\mathbf{w}(t)) k = 0$ for all u in the network.

The assumption (B3) is mild, as the norm of the gradient of cross-entropy loss goes down exponentially fast as compared to norm of the weights.

1.3. Asymptotic relations between weights and gradients

This section contains the main theorems that establish asymptotic relations between weights and gradients for SWN and EWN. First, we will state a common proposition for both SWN and EWN.

Proposition 11 Under assumption (A1) for gradient flow, for both SWN and EWN, $\lim_{t \rightarrow \infty} L(\mathbf{w}(t)) = 0$.

Proof First of all, for cross-entropy loss

$$\frac{dL(t)}{dw} = \sum_i \frac{\prod_{j \in \mathcal{Y}_i} e^{-(y_i(\mathbf{w}; \mathbf{x}_i) - j(\mathbf{w}; \mathbf{x}_i))} (r_{\mathbf{w}} y_i(\mathbf{w}; \mathbf{x}_i) - r_{\mathbf{w}} j(\mathbf{w}; \mathbf{x}_i))}{1 + \prod_{j \in \mathcal{Y}_i} e^{-(y_i(\mathbf{w}; \mathbf{x}_i) - j(\mathbf{w}; \mathbf{x}_i))}} \quad (49)$$

Now, using Theorem 1,

$$\frac{dL(t)}{dt} = \frac{dL(t)}{dw} \frac{dw(t)}{dt} = \sum_u k w_u(t) k^2 \frac{dL(t)}{dw_u}^2$$

Let k be the total number of neurons in the network. Then using the elementary inequality $\sum_{i=1}^n a_i^2 \leq n \max_{i=1, \dots, n} a_i^2$, we get

$$\frac{dL(t)}{dt} \leq \frac{1}{k} \sum_u k w_u(t) k \frac{dL(t)}{dw_u}^2$$

Again using the fact that $w(t) > \frac{dL(t)}{dw} \sum_u k w_u(t) k \frac{dL(t)}{dw_u}$, we get

$$\frac{dL(t)}{dt} \leq \frac{1}{k} w(t) \frac{dL(t)}{dw}^2 \quad (50)$$

Taking the dot product with w on both sides of Equation (49) and using $\nabla_w L(w; x_i) = L(w; x_i) w$ (Euler's homogeneity theorem), we get

$$w(t) \frac{dL(t)}{dw} = L \sum_i \frac{\prod_{j \in y_i} e^{-(y_j(w; x_i) - j(w; x_i))} (y_j(w; x_i) - j(w; x_i))}{1 + \prod_{j \in y_i} e^{-(y_j(w; x_i) - j(w; x_i))}}$$

Now, using the fact, that at time t_0 , $L(t_0) < \log 2$, which means $\min_i \min_{j \in y_i} (y_j(w; x_i) - j(w; x_i)) = \delta > 0$. Also, as we know, for gradient flow, the loss cannot go up, therefore, for any time $t > t_0$, $\min_i \min_{j \in y_i} (y_j(w; x_i) - j(w; x_i)) > \delta > 0$. Using this, we can say, for any $t > t_0$,

$$w(t) \frac{dL(t)}{dw} \geq L \sum_i \frac{\prod_{j \in y_i} e^{-(y_j(w; x_i) - j(w; x_i))}}{1 + \prod_{j \in y_i} e^{-(y_j(w; x_i) - j(w; x_i))}}$$

Using the fact that $\ln(1 + t) > \frac{t}{1+t}$ for $t \in (0; 1)$, therefore,

$$w(t) \frac{dL(t)}{dw} \geq L L(t)$$

Substituting this in Equation (50), we get

$$\frac{dL(t)}{dt} \leq \frac{L^2}{k} L(t)^2$$

Integrating this equation from t_0 to t , we get

$$\frac{1}{L(t)} \leq \frac{1}{L(t_0)} + \frac{L^2}{k} (t - t_0) \tag{51}$$

Clearly as t tends to ∞ , RHS tends to ∞ and thus L tends to 0. ■

Now, we provide one of our main theorem that establishes gradient convergence implies weight convergence.

Theorem 12 Consider a node u in the network with $k_{u,k} > 0$ and $\lim_{t \rightarrow \infty} \|kw_u(t)\| = 1$. Under assumptions (A1), (A2) for gradient flow and (B1)-(B3) for gradient descent, for both SWN and EWN

- (i) $\lim_{t \rightarrow \infty} \frac{kw_u(t)}{\|kw_u(t)\|} := w_u$ exists.
- (ii) $w_u = g_u$ for some $c > 0$.

Proof Same as in Appendix C. ■

Now, we provide the main theorem that distinguishes SWN and EWN.

Theorem 13 Consider two nodes u and v in the network with $k_{u,k} > k_{v,k} > 0$; $\lim_{t \rightarrow \infty} \|kw_u(t)\| = 1$ and $\lim_{t \rightarrow \infty} \|kw_v(t)\| = 1$. Let $\frac{k_{u,k}}{k_{v,k}}$ be denoted by c . Under assumptions (A1), (A2) for gradient flow and (B1)-(B3) for gradient descent,

- (i) for SWN, $\lim_{t \rightarrow \infty} \frac{\|kw_u(t)\|}{\|kw_v(t)\|} = c$
- (ii) for EWN, $\lim_{t \rightarrow \infty} \frac{\|kw_u(t)\|}{\|kw_v(t)\|}$ is either 0, 1 or $\frac{1}{c}$

Proof Same as in Appendix D. ■

I.4. Sparsity Inductive Bias for Exponential Weight Normalisation

The inverse relation between $\|w_u(t)\|_k$ and $\|r_{w_u} L(w(t))\|_k$ in the EWN trajectory results in an interesting inductive bias that favours movement along sparse directions.

Proposition 14 Consider two nodes u and v in the network such that $\|g_u\|_k, \|g_v\|_k > 0$ and $\|w_u(t)\|_k, \|w_v(t)\|_k \leq 1$. Let $\frac{\|g_u\|_k}{\|g_v\|_k}$ be denoted by c . Consider any θ such that $0 < \theta < c$ and $0 < \alpha < 2$. Then, the following holds:

(i) There exists a time t_1 , such that for all $t > t_1$ both SWN and EWN trajectories have the following properties:

(a) $\frac{\|r_{w_u} L(w(t))\|_k}{\|r_{w_v} L(w(t))\|_k} \geq [c - \alpha; c + \alpha]$ (b) $\frac{\|w_u(t)\|_k}{\|w_v(t)\|_k} > \frac{\|r_{w_u} L(w(t))\|_k}{\|r_{w_v} L(w(t))\|_k} \cos(\theta)$

(c) $\frac{\|w_v(t)\|_k}{\|w_u(t)\|_k} > \frac{\|r_{w_v} L(w(t))\|_k}{\|r_{w_u} L(w(t))\|_k} \cos(\theta)$.

(ii) for SWN, $\lim_{t \rightarrow \infty} \frac{\|w_u(t)\|_k}{\|w_v(t)\|_k} = c$

(iii) for EWN, if at some time $t_2 > t_1$,

(a) $\frac{\|w_u(t_2)\|_k}{\|w_v(t_2)\|_k} > \frac{1}{(c - \alpha) \cos(\theta)} \Rightarrow \lim_{t \rightarrow \infty} \frac{\|w_u(t)\|_k}{\|w_v(t)\|_k} = 1$

(b) $\frac{\|w_u(t_2)\|_k}{\|w_v(t_2)\|_k} < \frac{\cos(\theta)}{c + \alpha} \Rightarrow \lim_{t \rightarrow \infty} \frac{\|w_u(t)\|_k}{\|w_v(t)\|_k} = 0$

Proof The proof follows from Appendix E. ■

I.5. Convergence rates

In this section, we provide convergence rate of loss for EWN.

Gradient Flow: We provide a finite-time convergence rate of loss for gradient flow in case of EWN.

Theorem 15 For Exponential Weight Normalization, under assumption (A1), the following hold for $t > t_0$ in case of gradient flow

(i) $\|w(t)\|_k$ grows with $\text{as } O((\log t)^{\frac{1}{L}})$ (ii) $L(t)$ goes down with $\text{as } O\left(\frac{1}{t}\right)$

Proof Follow from Equation (51) ■

Gradient Descent:

Theorem 16 For Exponential Weight Normalization, under Assumptions (B1)-(B4), $L(t) = O\left(\log \frac{1}{L}\right)^c$ for $c < 1$ and $\lim_{t \rightarrow \infty} \frac{\|r_{w(t+1)}\|_k}{\|g(t+1)\|_k} = 0$, the following hold

(i) $\|w(t)\|_k$ asymptotically grows with $\text{as } (\log d(t))^{\frac{1}{L}}$

(ii) $L(w(t))$ asymptotically goes down with $\text{as } \frac{1}{d(t)(\log d(t))^2}$.

Proof The proof follows Appendix H.2, the only difference is in the gradient update. Let w be represented as $w = g(t)w + r(t)$, where $\lim_{t \rightarrow \infty} \frac{kr(t)k}{g(t)} = 0$. Using Equation (49), we can say

$$k_1 \lim_{t \rightarrow \infty} \frac{k \frac{dL(w(t))}{dw} k}{L(w(t))g(t)^L} \leq k_2 \quad (52)$$

where k_1 and k_2 are some constants. As the order remains the same as in the proof for exponential loss, the proof follows from Appendix H.2. ■

Appendix J. Lemma Proofs

Lemma Consider sequence a satisfying the following properties

1. $a_k > 0$
2. $\prod_{k=0}^{\infty} a_k = 1$
3. $\lim_{k \rightarrow \infty} a_k = 0$

Then $\prod_{k=0}^{\infty} \frac{a_k}{\sum_{j=0}^k a_j^2} = 1$

Proof If $\prod_{k=0}^{\infty} a_k^2$ is bounded, then the statement is obvious. Let's consider the case where $\prod_{k=0}^{\infty} a_k^2$ diverges. As $\lim_{k \rightarrow \infty} a_k = 0$, there must be an index k_1 , such that for $k > k_1$, $a_k < \frac{1}{2}$. Now, as $\prod_{k=0}^{\infty} a_k^2$ diverges, therefore, there must be an index $k_2 > k_1$, such that for any $k > k_2$,

$$\begin{aligned} \prod_{j=k_1}^k \frac{a_j}{\sum_{l=0}^j a_l^2} &= \prod_{j=k_1}^k \frac{a_j}{\sum_{l=k_1}^j a_l^2} \\ &= \prod_{j=k_1}^k \frac{a_j}{\sum_{l=k_1}^j a_l} \\ &= \prod_{j=k_1}^k \frac{a_j}{\sum_{l=k_1}^j a_l} \\ &= \prod_{j=k_1}^k \frac{a_j}{\sum_{l=k_1}^j a_l} \end{aligned}$$

As $\prod_{k=0}^{\infty} a_k$ diverges, therefore $\prod_{k=0}^{\infty} \frac{a_k}{\sum_{j=0}^k a_j^2}$ diverges as well. ■

Lemma Consider two sequences a and b satisfying the following properties

1. $a_k > 0$; $\prod_{k=0}^{\infty} a_k = 1$ and $\lim_{k \rightarrow \infty} a_k = 0$

2. $b_0 > 0$, b is increasing and $b_{k+1}^2 = b_k^2 + \frac{a_k}{b_k}^2$

Then $\sum_{k=0}^{\infty} \frac{a_k}{b_k} = 1$.

Proof As we know b is increasing and $b_{k+1}^2 = b_k^2 + (\frac{a_k}{b_k})^2$, we get

$$b_k \sum_{j=0}^{k-1} \frac{a_j}{b_j} \leq \sum_{j=0}^{k-1} \frac{a_j}{b_0 + \frac{1}{b_0^2} \sum_{l=0}^{j-1} a_l^2}$$

Using this, we can say

$$\sum_{j=0}^k \frac{a_j}{b_j} \leq \sum_{j=0}^k \frac{a_j}{b_0 + \frac{1}{b_0^2} \sum_{l=0}^{j-1} a_l^2} \leq \sum_{j=0}^k \frac{a_j}{b_0 + \frac{1}{b_0^2} \sum_{l=0}^{k-1} a_l^2}$$

Now, if $\sum_{k=0}^{\infty} a_k^2$ does not diverge to infinity, then b remains bounded using the bound above and then it's trivial to establish that $\sum_{k=0}^{\infty} \frac{a_k}{b_k}$ diverges. In case $\sum_{k=0}^{\infty} a_k^2$ diverges to infinity, then there must be an index k_1 such that for any $k > k_1$, we can say $\sum_{j=0}^k a_j^2 \geq b_0^4$. So, for $k > k_1$, we can say

$$\sum_{j=0}^k \frac{a_j}{b_j} \leq \sum_{j=0}^k \frac{b_0}{2} \frac{a_j}{\sum_{l=0}^{k-1} a_l^2}$$

Now, as we have assumed a_k tends to zero, so there must be an index k_2 such that for any $k > k_2$, $a_k < \frac{1}{2}$. Also, as we have assumed $\sum_{j=0}^{\infty} a_j^2$ diverges, therefore there must be an index $k_3 > k_2$, such that for $k > k_3$, $\sum_{j=k_2}^k a_j^2 \geq \sum_{j=0}^{k_2} a_j^2$. Using these things and that $a_j < \frac{1}{2}$, then $a_j^2 < a_j$, we can say for $k > k_3$,

$$\sum_{j=k_3}^k \frac{a_j}{b_j} \leq \sum_{j=k_3}^k \frac{b_0}{2} \frac{a_j}{\sum_{l=k_3}^{k-1} a_l} \leq \frac{b_0}{2} \sum_{j=k_3}^k \frac{1}{a_j}$$

Now, as $\sum_{k=0}^{\infty} a_k$ diverges, thus $\sum_{k=0}^{\infty} \frac{a_k}{b_k}$ diverges as well. ■

Appendix K. Integral Form of Stolz-Cesaro Theorem

We first state the Stolz-Cesaro Theorem.

Theorem (Muresan, 2015) Assume that $\{a_k\}_{k=1}^{\infty}$ and $\{b_k\}_{k=1}^{\infty}$ are two sequences of real numbers such that $\{b_k\}_{k=1}^{\infty}$ is strictly monotonic and diverging. Additionally, if $\lim_{k \rightarrow \infty} \frac{a_{k+1} - a_k}{b_{k+1} - b_k} = L$ exists, then $\lim_{k \rightarrow \infty} \frac{a_k}{b_k}$ exists and is equal to L .

Now, we state and prove the Integral Form of Stolz-Cesaro Theorem.

Theorem Consider two functions $f(t)$ and $g(t)$ greater than zero satisfying $\int_a^b f(t) dt < 1$ and $\int_a^b g(t) dt < 1$ for every a, b . For any time t , it's known that $\int_t^1 f(t) dt = 1 - \int_t^1 g(t) dt = 1$. If $\lim_{t \rightarrow 1} \frac{f(t)}{g(t)}$ exist and is equal to L , then $\lim_{t \rightarrow 1} \frac{\int_c^t f(t) dt}{\int_c^t g(t) dt}$ exists for any c and is equal to L .

Proof Case 1: $L = 0$ or 1 :

We will prove for $L = 1$. The case for 0 can be handled similarly. For any $\epsilon > 0$, there must exist a time $t_1 > c$, such that $\frac{f(t)}{g(t)} > M - \epsilon$, for $t > t_1$. Thus we can say for $t > t_1$,

$$\int_c^t f(t) dt > \int_c^{t_1} f(t) dt + M \int_{t_1}^t g(t) dt$$

Adding $M \int_c^{t_1} g(t) dt$ on both the sides, we get

$$\int_c^t f(t) dt + M \int_c^{t_1} g(t) dt > \int_c^{t_1} f(t) dt + M \int_c^{t_1} g(t) dt$$

Dividing both sides by $\int_c^t g(t) dt$ and taking $\limsup_{t \rightarrow \infty}$ (using also the fact that $\int_a^b f(t) dt < 1$ and $\int_a^b g(t) dt < 1$ for every finite a, b), we get

$$\limsup_{t \rightarrow \infty} \frac{\int_c^t f(t) dt}{\int_c^t g(t) dt} > M - \epsilon$$

Similarly the equation holds for \liminf as well. Thus, both \liminf and \limsup are greater than $M - \epsilon$ for any M . Hence $\lim_{t \rightarrow \infty} \frac{\int_c^t f(t) dt}{\int_c^t g(t) dt} = 1$.

Case 2: L is finite

In this case, there must exist some time $t_1 > c$, such that $L - \epsilon < \frac{f(t)}{g(t)} < L + \epsilon$. Thus, we can say for $t > t_1$,

$$\int_c^{t_1} f(t) dt + (L - \epsilon) \int_{t_1}^t g(t) dt < \int_c^t f(t) dt < \int_c^{t_1} f(t) dt + (L + \epsilon) \int_{t_1}^t g(t) dt$$

Taking the left inequality, adding $(L - \epsilon) \int_c^{t_1} g(t) dt$ on both the sides, dividing both the sides by $\int_c^t g(t) dt$ and taking $\liminf_{t \rightarrow \infty}$, we get

$$L - \epsilon \leq \liminf_{t \rightarrow \infty} \frac{\int_c^t f(t) dt}{\int_c^t g(t) dt}$$

Similarly, taking the right inequality, adding $(L + \epsilon) \int_c^{t_1} g(t) dt$ on both the sides, dividing both the sides by $\int_c^t g(t) dt$ and taking $\limsup_{t \rightarrow \infty}$, we get

$$\limsup_{t \rightarrow \infty} \frac{\int_c^t f(t) dt}{\int_c^t g(t) dt} \leq L + \epsilon$$

Using the two inequalities, we get, for any $\epsilon > 0$,

$$\limsup_{t \rightarrow \infty} \frac{\int_c^t f(t) dt}{\int_c^t g(t) dt} \leq \liminf_{t \rightarrow \infty} \frac{\int_c^t f(t) dt}{\int_c^t g(t) dt} + 2\epsilon$$

Thus, $\lim_{t \rightarrow \infty} \frac{\int_c^t f(t) dt}{\int_c^t g(t) dt}$ exists and is equal to L . ■

Appendix L. Standard Weight Normalization is not Locally Lipschitz in its parameters

In this section, we will denote \mathbf{w}_u by \mathbf{w}_k so as to be consistent with the notation in [Lyu and Li \(2020\)](#). SWN (in its parameters and \mathbf{v}) is also a homogeneous network. Therefore, results from [Lyu and Li \(2020\)](#) should directly apply to the case of SWN as well. However, a crucial point to be noted is that it is not even locally Lipschitz around $\|\mathbf{w}_u\|_k = 0$. Therefore, the assumptions from [Lyu and Li \(2020\)](#) do not hold.

However, during gradient descent or gradient flow, if started from a $\|\mathbf{w}_u\|_k > 0$, for all u , then during the entire trajectory $\|\mathbf{w}_u\|_k$ cannot go down. Therefore, the network is still locally Lipschitz along the trajectory it takes. Examining the proofs from [Lyu and Li \(2020\)](#), it's clear that the proof regarding monotonicity of margin and convergence rates are just dependent on the path that gradient descent/flow takes and thus the proofs hold.

However, the result regarding the limit points of \mathbf{w}_k do not hold. One of the crucial theorems the proof relies on is stated below

Theorem Let $\{x_k\}_{k \in \mathbb{N}} \subset \mathbb{R}^d : k \in \mathbb{N}$ be a sequence of feasible points of an optimization problem (P), $f_k > 0 : k \in \mathbb{N}$ and $\lambda_k > 0 : k \in \mathbb{N}$ be two sequences. x_k is an (λ_k, f_k) -KKT point for every k and $\lambda_k \rightarrow 0; f_k \rightarrow 0$. If $x_k \rightarrow x^*$ and MFCQ holds at x^* , then x^* is a KKT point of (P)

The above statement requires MFCQ to be satisfied, which was shown in [Lyu and Li \(2020\)](#) assuming local Lipschitzness/smoothness. However, in this case, for gradient flow, $\|\mathbf{w}_u\|_k$ does not grow, while $\|\mathbf{v}_u\|_1 \rightarrow 1$, therefore the convergent point of \mathbf{w}_k will always have the component corresponding to \mathbf{v}_u as 0. Thus, the network is not locally Lipschitz and the proof that MFCQ holds is violated. Similarly, for gradient descent as well, it can't be said that it has a non-zero component in \mathbf{v}_k . Thus, the proof does not hold.

Appendix M. Experiment Details

In all the experiments, techniques for handling numerical underflow were used as described in [Lyu and Li \(2020\)](#). However, the learning rate they used was $\frac{1}{L^c}$, but in our case, we generally modify it to be $O\left(\frac{1}{L^c}\right)$, where $c < 1$.

M.1. Lin-Sep

The learning rate used was $\frac{k(t)}{L^{0.97}}$, so that it speeds up at the beginning of training, but slows down as loss approaches 0. The constant $k(t)$ was initialized at 0.01, and was increased by a factor of 1:1 every time loss went down and decreased by a factor of 1:1 every time loss went up after a gradient step. Its value was capped at 1 for EWN and SWN.

M.2. Simple-Traj

The learning rate used was $\frac{k(t)}{L^{0.9}}$, so that it speeds up at the beginning of training, but slows down as loss approaches 0. The constant $k(t)$ was initialized at 0.01, and was increased by a factor of 1:1 every time loss went down and decreased by a factor of 1:1 every time loss went up after a gradient step. Its value was capped at 1 for EWN and Unnorm.

(a) (b) (c)

Figure 7: Demonstration of Results for SWN in Lin-Sep experiment: (a) Evolution of $\|w_{\cup k}\|$ (b) Cosine between weights and gradients for weights 5, 6 and 8. (c) Weight and gradient norms for weights 5, 6 and 8.

(a) (b) (c)

Figure 8: Demonstration of Results for EWN in XOR experiment with ReLU-square activation: (a) Evolution of $\|w_{\cup k}\|$ (b) Cosine between weights and gradients for weights 0, 1, 13 and 17. (c) Weight and gradient norms for weights 0, 1, 13 and 17.

M.3. XOR

The learning rate used was $\frac{k(t)}{0.93}$ for SWN and Unnorm, while $\frac{k(t)}{0.8}$ for EWN, so that it speeds up at the beginning of training, but slows down as loss approaches 0. The constant $k(t)$ was initialized at 0.01, and was increased by a factor of 1.1 every time loss went down and decreased by a factor of 0.9 every time loss went up after a gradient step. Its value was capped at 1.0 for EWN and Unnorm and 0.01 for SWN.

M.4. Convergence rate experiment

For all SWN, EWN and Unnorm, the learning rate was constant at 0.001 and they were trained for 5000 steps. All the networks were explicitly initialized to the same point in function space.

(a) (b) (c)

Figure 9: Demonstration of Results for SWN in XOR experiment with ReLU-square activation:
 (a) Evolution of $\|w_u\|$ (b) Cosine between weights and gradients for weights 0, 1, 13 and 17. (c) Weight and gradient norms for weights 0, 1, 13 and 17.

(a) (b) (c)

Figure 10: Demonstration of Results for EWN on MNIST dataset with 2-class classification:
 (a) Evolution of $\|w_u\|$ (b) Cosine between weights and gradients for weights 96 and 105. (c) Weight and gradient norms for weights 96 and 105.

M.5. Pruning Experiments

The learning rate used was $\frac{\eta}{L}$. The constant η was initialized at 0.01, and was increased by a factor of 1.1 every time loss went down and decreased by a factor of 1.05 every time loss went up after an epoch.

Appendix N. Demonstration of Theorems on various datasets

In this section, we demonstrate Theorem 3 and 4 on various datasets: Lin-Sep, XOR and MNIST (2-class classification).

N.1. Lin-Sep

We demonstrate Theorem 3 and Theorem 4 for EWN and SWN on a linearly separable dataset (Lin-Sep) in Figure 1 and 7 respectively. As can be seen in Figure 7, for weights 5, 6 and 8,

(a) $L = e^{-10}$ (b) $L = e^{-100}$ (c) $L = e^{-300}$

Figure 11: Convergence rate of EWN, SWN and Unnorm on the MNIST dataset for seed - 2356 at different loss values

(a) $L = e^{-10}$ (b) $L = e^{-100}$ (c) $L = e^{-300}$

Figure 12: Convergence rate of EWN, SWN and Unnorm on the MNIST dataset for seed - 3576 at different loss values

whose norms keep on growing, weights and gradients eventually become oppositely aligned, and their norms are directly proportional to each other.

N.2. XOR

In this experiment, we train a 2-layer network with 20 hidden neurons and ReLU-square activation on XOR dataset, till a loss value of 10^{-40} . We demonstrate Theorems 3 and 4 for EWN and SWN in Figure 8 and Figure 9 respectively.

As can be seen in Figure 8, for weights 0, 1, 13 and 17, whose norms keep on growing, weights and gradients eventually become oppositely aligned, and their norms are inversely proportional to each other.

Similarly, in Figure 9, for weights 0, 1, 13 and 17, whose norms keep on growing, weights and gradients eventually become oppositely aligned, and their norms are directly proportional to each other.

N.3. MNIST

In this experiment, we train a 2-layer network with 128 hidden neurons and ReLU-square activation on MNIST dataset with 2 classes. Since even after considering just 2 classes, this dataset is huge, therefore training takes longer. Therefore, we only consider exponential weight normalized network in this case.

We demonstrate Theorems 3 and 4 for MNIST dataset for EWN in Figure 10. As can be seen, for weights 96 and 105, whose norms keep on growing, weights and gradients eventually become oppositely aligned, and their norms are inversely proportional to each other.

Appendix O. Convergence rate plots for pruning experiment

In this section, we provide convergence rate plots for the pruning experiments.

The convergence rate for 2 different seeds at various loss levels are shown in Figure 11 and 12. As can be seen, initially the convergence rate of all the normalizations are comparable. But at extremely low loss values, Unnorm becomes slightly faster as compared to SWN or EWN. Note that the results regarding asymptotic convergence rate do not apply in this case, as we are training at extremely high learning rates of $\frac{1}{L}$.

Appendix P. Verification of Assumptions for various datasets

In this section, we verify the assumptions on three datasets: Lin-Sep, XOR and MNIST (2-class classification). Note that in Figures 13, 14, 15, 16 and 17, plots demonstrating the components of unit gradient vector, each line corresponds to a single parameter of the network, while in the plots demonstrating the evolution of $\|w_u(t)\|$ or $\|w_u L(w(t))\|$, each line corresponds to a neuron of the network.

P.1. Lin-sep

We verify assumptions (B1)-(B3) for the Lin-Sep experiment for EWN and SWN in Figure 13 and Figure 14 respectively. As can be seen, for both the cases, the components of unit gradient vector become constant as training proceeds. Another thing to note, is that even for an aggressive learning rate schedule of the form $\frac{1}{L}$, $\|w_u(t)\|$ or $\|w_u L(w(t))\|$ still goes down to 0.

P.2. XOR

In this experiment, we train a 2-layer network with 20 hidden neurons and ReLU-square activation on XOR dataset, till a loss value of 10^{-40} . We verify assumptions (B1)-(B3) for the XOR experiment for EWN and SWN in Figure 15 and Figure 16 respectively. As can be seen, for both the cases, the components of unit gradient vector become constant as training proceeds. Another thing to note, is that even for an aggressive learning rate schedule of the form $\frac{1}{L}$, $\|w_u(t)\|$ or $\|w_u L(w(t))\|$ still goes down to 0.

P.3. MNIST

In this experiment, we train a 2-layer network with 128 hidden neurons and ReLU-square activation on MNIST dataset with 2 classes. Since even after considering just 2 classes, this dataset is

(a) (b)

Figure 13: Verification of Assumptions for EWN in Lin-Sep experiment: (a) Evolution of $\frac{r_w L(t)}{kr_w L(t)k}$ (b) Evolution of $(t)kw_u(t)kkr_w L(t)k$

(a) (b)

Figure 14: Verification of Assumptions for SWN in Lin-Sep experiment: (a) Evolution of $\frac{r_w L(t)}{kr_w L(t)k}$ (b) Evolution of $(t)kw_u(t)kkr_w L(t)k$

huge, therefore training takes longer. Therefore, we only consider exponential weight normalized network in this case.

We verify assumptions (B1)-(B3) for MNIST dataset for EWN in Figure 17. As can be seen, the components of unit gradient vector become constant as training proceeds. Another thing to note, is that even for an aggressive learning rate schedule of the form $(t)kw_u(t)kkr_w L(w(t))k$ still goes down to 0.

Appendix Q. Pruning algorithm

We will explain the pruning algorithm used in Figure 5b. The same algorithm is used for SWN, EWN as well as the unnormalized net.

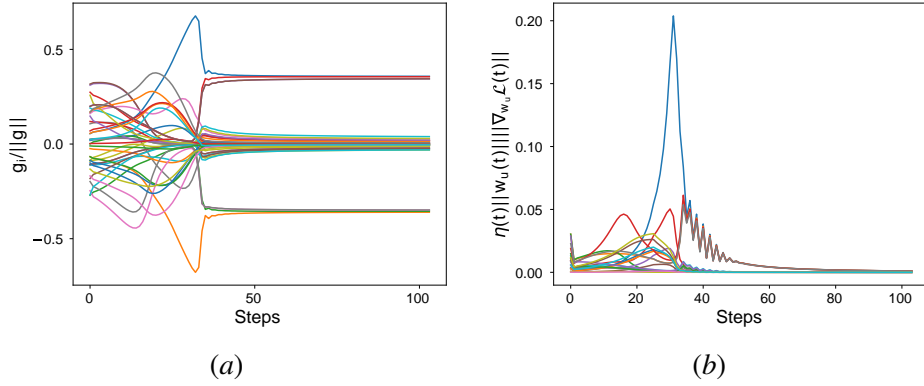


Figure 15: **Verification of Assumptions for EWN in XOR experiment with ReLU-square activation:** (a) Evolution of $\frac{r_{w^L(t)}}{kr_{w^L(t)}k}$ (b) Evolution of $\eta(t)kw_u(t)kkr_{w^L(t)}k$

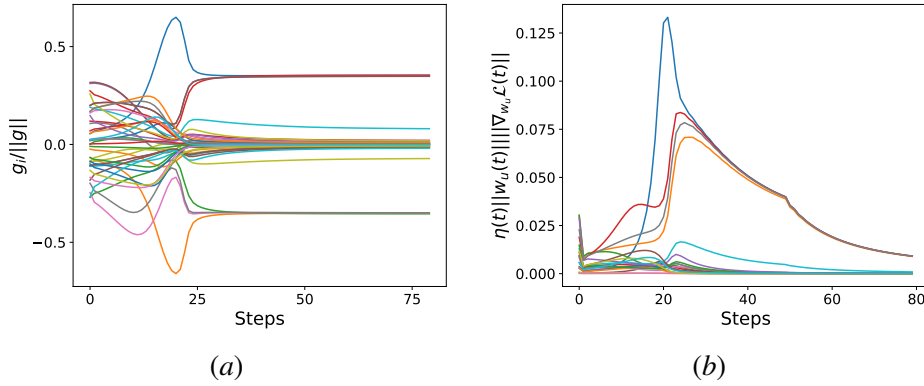


Figure 16: **Verification of Assumptions for SWN in XOR experiment with ReLU-square activation:** (a) Evolution of $\frac{r_{w^L(t)}}{kr_{w^L(t)}k}$ (b) Evolution of $\eta(t)kw_u(t)kkr_{w^L(t)}k$

Let t_1 and t_2 denote the optimization iteration indices when log-loss has a value of 10 and 100 respectively. Consider three pruning strategies available at the time instant given by t_2 :

- Prune weights on the basis of $kw_u(t_2)k$
- Prune weights on the basis of $kw_u(t_2) - w_u(0)k$
- Prune weights on the basis of $kw_u(t_2) - w_u(t_1)k$

For a given level of pruning, we try all the 3 strategies, and then pick the one with the best test performance.

Variants of option (a) are the most prevalent pruning algorithms. The non-standard options of (b) and (c) represent the intuition of pruning neurons whose weight has not moved in the recent

