
Characterizing Implicit Bias in Terms of Optimization Geometry

Suriya Gunasekar¹ Jason Lee² Daniel Soudry³ Nathan Srebro¹

Abstract

We study the implicit bias of generic optimization methods, including mirror descent, natural gradient descent, and steepest descent with respect to different potentials and norms, when optimizing under determined linear regression or separable linear classification problems. We explore the question of whether the specific global minimum (among the many possible global minima) reached by optimization can be characterized in terms of the potential or norm of the optimization geometry, and independently of hyperparameter choices such as step size and momentum.

1. Introduction

Implicit bias from the optimization algorithm plays a crucial role in learning deep neural networks, as it introduces effective capacity control not directly specified in the objective (Neyshabur et al., 2015b;a; Zhang et al., 2017; Keskar et al., 2016; Wilson et al., 2017; Neyshabur et al., 2017). In over-parameterized models where the training objective has many global minima, optimizing using a specific algorithm, such as gradient descent, *implicitly biases* the solutions to some special global minima. These special global minima in turn specify the properties of the learned model, including its generalization performance. In deep neural networks especially, characterizing the specific global minima reached by local search methods such as stochastic gradient descent (SGD) is essential for understanding what the inductive bias of the learned model is and why such large capacity networks often show remarkably good generalization performance even in the absence of explicit regularization (Zhang et al., 2017) or early stopping (Hoffer et al., 2017).

The implicit bias depends on the choice of optimization algorithm, and changing the optimization algorithm, or even

changing associated hyperparameter can change the implicit bias. For example, Wilson et al. (2017) showed that on standard deep learning architectures, common variants of SGD methods for different choices of momentum and adaptive gradient updates (AdaGrad and Adam) exhibit different biases and thus have different generalization performance; Keskar et al. (2016), Hoffer et al. (2017) and Smith (2018) study how the size of the mini-batches used in SGD influences generalization; and Neyshabur et al. (2015a) compare the bias of path-SGD (steepest descent with respect to a scale invariant path-norm) to standard SGD.

To rigorously understand learning and generalization in deep neural network training, it is therefore important to understand what the implicit biases are for different algorithms. Can we explicitly relate between the choice of algorithm and the implicit bias? Can we precisely characterize which global minima different optimization algorithms converge to? How does this depend on the loss functions? What other choices including initialization, step size, momentum, stochasticity, and adaptivity, does the implicit bias depend on? In this paper, we provide answers to some of these questions for simple linear regression and classification models.

We already have an understanding of the implicit bias of gradient descent for linear models. For underdetermined least squares objective, gradient descent can be shown to converge to the minimum Euclidean norm solution. Recently, Soudry et al. (2017) studied gradient descent for linear logistic regression. The logistic loss is fundamentally different from the squared loss in that the loss function has no attainable global minima. Gradient descent iterates therefore diverge (the norm goes to infinity), but Soudry et al. showed that they diverge in the direction of the hard margin support vector machine solution, and therefore the decision boundary converges to this max margin solution.

Can we extend such characterization to other optimization methods that work under different (non-Euclidean) geometries, such as mirror descent with respect to some potential, natural gradient descent with respect to a Riemannian metric, and steepest descent with respect to a generic norm? Can we relate the implicit bias to these geometries?

As we shall see, the answer depends on whether the loss function is similar to a squared loss or to a logistic loss. This difference is captured by two family of losses: (a) loss

¹TTI Chicago, USA ²USC Los Angeles, USA ³Technion, Israel. Correspondence to: Suriya Gunasekar <suriya@ttic.edu>, Jason Lee <jasonlee@marshall.usc.edu>, Daniel Soudry <daniel.soudry@gmail.com>, Nathan Srebro <nati@ttic.edu>.

functions that have a unique finite root, like the squared loss and (b) strictly monotone loss functions where the infimum is unattainable, like the logistic loss. For losses with a unique finite root, we study the *limit point* of the optimization iterates, $w_\infty = \lim_{t \rightarrow \infty} w_{(t)}$. For monotone losses, we study the *limit direction* $\bar{w}_\infty = \lim_{t \rightarrow \infty} \frac{w_{(t)}}{\|w_{(t)}\|}$.

In Section 2 we study linear models with loss functions that have unique finite roots. We obtain a robust characterization of the limit point for mirror descent, and discuss how it is independent of step size and momentum. For natural gradient descent, we show that the step size does play a role, but get a characterization for infinitesimal step size. For steepest descent, we show that not only does step size affects the limit point, but even with infinitesimal step size the expected characterization does not hold. The situation is fundamentally different for strictly monotone losses such as the logistic loss (Section 3) where we do get a precise characterization of the implicit bias of the limit direction for generic steepest descent. We also study the adaptive method AdaGrad and optimization over a matrix factorization. Recent studies considered the bias of such methods for least squares problems (Wilson et al., 2017; Gunasekar et al., 2017), and here we study these algorithms for monotone loss functions, obtaining a more robust characterization for matrix factorization problems, while concluding that the implicit bias of AdaGrad depends on initial conditions including step sizes even for strict monotone losses.

2. Losses with a Unique Finite Root

We first consider learning linear models using losses with a unique finite root, such as the squared loss, where loss function $\ell(f(x), y)$ between a predictor $f(x)$ and label y is minimized at a unique and finite value of $f(x)$.

Property 1 (Losses with a unique finite root). *For any y and sequence \hat{y}_t , $\ell(\hat{y}_t, y) \rightarrow \inf_{\hat{y}} \ell(\hat{y}, y)$ if and only if $\hat{y}_t \rightarrow y$. Here we assumed without loss of generality that $\inf_{\hat{y}} \ell(\hat{y}, y) = 0$ and the root of $\ell(\hat{y}, y)$ is at $\hat{y} = y$.*

Denote the training dataset $\{(x_n, y_n) : n = 1, 2, \dots, N\}$ with features $x_n \in \mathbb{R}^d$ and labels $y_n \in \mathbb{R}$. The empirical loss (or risk) minimizer of a linear model $f(x) = \langle w, x \rangle$ with parameters $w \in \mathbb{R}^d$ is given by,

$$\min_w \mathcal{L}(w) := \sum_{n=1}^N \ell(\langle w, x_n \rangle; y_n). \quad (1)$$

We are particularly interested in the case when $N < d$ and the observations are realizable, i.e., $\min_w \mathcal{L}(w) = 0$. In this case $\mathcal{L}(w)$ is under constrained and we have multiple global minima denoted by $\mathcal{G} = \{w : \mathcal{L}(w) = 0\} = \{w : \forall n, \langle w, x_n \rangle = y_n\}$. Note that the set \mathcal{G} is the same for any loss ℓ with unique finite root (Property 1), including, e.g., the Huber loss, the truncated squared loss. The question we

want to answer here is which specific global minima $w \in \mathcal{G}$ do different optimization algorithms reach when minimizing the empirical loss objective $\mathcal{L}(w)$.

2.1. Gradient descent

Consider gradient descent updates with step size sequence $\{\eta_t\}$ for minimizing $\mathcal{L}(w)$, $w_{(t+1)} = w_{(t)} - \eta_t \nabla \mathcal{L}(w_{(t)})$. If $w_{(t)}$ minimizes the empirical loss in eq. (1), then the iterates converge to the unique global minimum that is closest to initialization $w_{(0)}$ in ℓ_2 distance, $\operatorname{argmin}_{w \in \mathcal{G}} \|w - w_{(0)}\|_2$. This can be easily seen as at any w , the gradients $\nabla \mathcal{L}(w) = \sum_n \ell'(\langle w, x_n \rangle, y_n) x_n$ are always constrained to a fixed subspace spanned by the data $\{x_n\}_n$, and thus the iterates $w_{(t)}$ are confined to the low dimensional affine manifold $w_{(0)} + \operatorname{span}(\{x_n\}_n)$. Within this N -dimensional manifold, there is a unique global minimizer w that satisfies the N linear constraints $\mathcal{G} = \{w : \langle w, x_n \rangle = y_n, \forall n \in [N]\}$.

It is also evident that this bias also holds for updates with instance-wise stochastic gradients, where in place of the gradient $\nabla \mathcal{L}(w_{(t)})$ over the entire dataset, we use stochastic gradients computed from a random subset of instances $S_t \subseteq [N]$ as defined below:

$$\tilde{\nabla} \mathcal{L}(w_{(t)}) = \sum_{n \in S_t \subseteq [n]} \nabla_W \ell(\langle w_{(t)}, x_n \rangle; y_n). \quad (2)$$

Moreover, when initialized with $w_{(0)} = 0$, the implicit bias characterization also extends to the following generic momentum and acceleration based updates:

$$w_{(t+1)} = w_{(t)} + \beta_t \Delta w_{(t-1)} - \eta_t \nabla \mathcal{L}(w_{(t)}) + \gamma_t \Delta w_{(t-1)}, \quad (3)$$

where $\Delta w_{(t-1)} = w_{(t)} - w_{(t-1)}$. This includes Nesterov's acceleration ($\beta_t = \gamma_t$) (Nesterov, 1983) and Polyak's heavy ball momentum ($\gamma_t = 0$) (Polyak, 1964).

For losses with a unique finite root, the implicit bias of gradient descent therefore depends only on the initialization and not on the step size or momentum or mini-batch size. Can we get such succinct characterization for other optimization algorithms? That is, characterize the bias in terms of the optimization geometry and initialization, but independent of choices of step sizes, momentum, and stochasticity.

2.2. Mirror descent

Mirror descent (MD) (Beck & Teboulle, 2003; Nemirovskii & Yudin, 1983) was introduced as a generalization of gradient descent for optimization over geometries beyond the Euclidean geometry of gradient descent. In particular, mirror descent updates are defined for any strongly convex and differentiable potential ψ as,

$$w_{(t+1)} = \operatorname{argmin}_{w \in \mathcal{W}} \eta_t \langle w, \nabla \mathcal{L}(w_{(t)}) \rangle + D_\psi(w, w_{(t)}), \quad (4)$$

where $D_\psi(w, w') = \psi(w) - \psi(w') - \langle \nabla \psi(w'), w - w' \rangle$ is the *Bregman divergence* (Bregman, 1967) w.r.t. ψ , and \mathcal{W} is some constraint set for parameters w .

We look at unconstrained optimization where $\mathcal{W} = \mathbb{R}^d$ and the update in eq. (4) can be equivalently written as:

$$\nabla\psi(w_{(t+1)}) = \nabla\psi(w_{(t)}) - \eta_t \nabla\mathcal{L}(w_{(t)}). \quad (5)$$

For a strongly convex potential ψ , the link function $\nabla\psi$ is invertible and hence, the above updates are uniquely defined. $\nabla\psi(w)$ is often referred as the *dual* variable corresponding to the *primal* variable w . Examples of strongly convex potentials ψ for mirror descent include, the squared ℓ_2 norm $\psi(w) = \|w\|_2^2$, which leads to gradient descent; the entropy potential $\psi(w) = \sum_i w[i] \log w[i] - w[i]$; the spectral entropy for matrix valued w , wherein $\psi(w)$ is the entropy potential on the singular values of w ; general quadratic potentials $\psi(w) = \|w\|_D^2 = w^\top D w$ for any positive definite matrix D ; and the squared ℓ_p norm for $p \in (1, 2]$.

From the updates given in eq. (5), we can see that, rather than the primal iterates $w_{(t)}$, it is the dual iterates $\nabla\psi(w_{(t)})$ that are constrained to the low dimensional data manifold $\nabla\psi(w_{(0)}) + \text{span}(\{x_n\}_{n \in [N]})$. The arguments for gradient descent can now be generalized to get the following result.

Theorem 1. *For any loss ℓ with a unique finite root (Property 1), any initialization $w_{(0)}$, any step size sequence $\{\eta_t\}$, and any strongly convex potential ψ , consider the corresponding mirror descent iterates $w_{(t)}$ from eq. (5). If the limit point of the iterates $w_\infty = \lim_{t \rightarrow \infty} w_{(t)}$ is a global minimizer of \mathcal{L} , i.e., $\mathcal{L}(w_\infty) = 0$, then w_∞ is given by*

$$w_\infty = \underset{w: \forall n, \langle w, x_n \rangle = y_n}{\text{argmin}} D_\psi(w, w_{(0)}). \quad (6)$$

In particular, if we start at $w_{(0)} = \underset{w}{\text{argmin}} \psi(w)$, then we get to $w_\infty = \underset{w \in \mathcal{G}}{\text{argmin}} \psi(w)$, where recall that $\mathcal{G} = \{w : \forall n, \langle w, x_n \rangle = y_n\}$ is the set of global minima for $\mathcal{L}(w)$.

Let us now consider momentum for mirror descent. There are two generalizations of gradient descent momentum in eq. (3): adding momentum either to primal variables $w_{(t)}$, or to dual variables $\nabla\psi(w_{(t)})$,

$$\text{Dual momentum: } \nabla\psi(w_{(t+1)}) = \nabla\psi(w_{(t)}) + \beta_t \Delta z_{(t-1)} - \eta_t \nabla\mathcal{L}(w_{(t)} + \gamma_t \Delta w_{(t-1)}) \quad (7)$$

$$\text{Primal momentum: } \nabla\psi(w_{(t+1)}) = \nabla\psi(w_{(t)} + \beta_t \Delta w_{(t-1)}) - \eta_t \nabla\mathcal{L}(w_{(t)} + \gamma_t \Delta w_{(t-1)}) \quad (8)$$

where $\Delta z_{(-1)} = \Delta w_{(-1)} = 0$, and for $t \geq 1$, $\Delta z_{(t-1)} = \nabla\psi(w_{(t)}) - \nabla\psi(w_{(t-1)})$ and $\Delta w_{(t-1)} = w_{(t)} - w_{(t-1)}$ are the momentum terms in the primal and dual space, respectively; and $\{\beta_t \geq 0, \gamma_t \geq 0\}$ are the momentum parameters.

If we initialize to $w_{(0)} = \underset{w}{\text{argmin}} \psi(w)$, then even with dual momentum, $\nabla\psi(w_{(t)})$ continues to remain in the data manifold, leading to the following extension of Theorem 1.

Theorem 1a. *Under the conditions in Theorem 1, if initialized at $w_{(0)} = \underset{w}{\text{argmin}} \psi(w)$, then the mirror descent*

updates with dual momentum also converges to (6), i.e., for all $\{\eta_t\}_t, \{\beta_t\}_t, \{\gamma_t\}_t$, if $w_{(t)}$ from eq. (7) converges to $w_\infty \in \mathcal{G}$, then $w_\infty = \underset{w \in \mathcal{G}}{\text{argmin}} \psi(w)$.

Theorem 1–1a also hold when stochastic gradients defined in eq. (2) are used in place of $\nabla\mathcal{L}(w_{(t)})$ in the mirror descent updates in eq. (4). Furthermore, the exact analysis also extends when updates in eq. (4) are carried out for constrained optimization with affine equality constraints, i.e., $\mathcal{L}(w)$ in eq. (1) is minimized over $\mathcal{W} = \{w : Gw = h\}$ for some G and h , as long as there exists a feasible solution $w \in \mathcal{W}$ with $\mathcal{L}(w) = 0$. For example, the results will also extend for exponentiated gradient descent (Kivinen & Warmuth, 1997), which is MD w.r.t $\psi(w) = \sum_i w[i] \log w[i] - w[i]$ under the explicit simplex constraint $\mathcal{W} = \{w : \sum_i w[i] = 1\}$.

For quadratic potentials $\psi(w) = \|w\|_D^2$, the primal momentum in eq. (8) is equivalent to the dual momentum in eq. (7). For general potentials ψ , the dual of iterates $\nabla\psi(w_{(t)})$ from the primal momentum can fall off the data manifold and the additional components influence the final solution. Thus, the specific global minimum that the iterates $w_{(t)}$ converge to can depend on the values of momentum parameters $\{\beta_t, \gamma_t\}$ and step sizes $\{\eta_t\}$.

Example 2. *Consider optimizing $\mathcal{L}(w)$ with dataset $\{(x_1 = [1, 2], y_1 = 1)\}$ and squared loss $\ell(u, y) = (u - y)^2$ using primal momentum updates from eq. (8) for MD w.r.t the entropy potential $\psi(w) = \sum_i w[i] \log w[i] - w[i]$ and initialization $w_{(0)} = \underset{w}{\text{argmin}} \psi(w)$. Figure 1a shows how different choices of momentum $\{\beta_t, \gamma_t\}$ change the limit point w_∞ . Additionally, we show the following:*

Proposition 2a. *In Example 2, consider the case where primal momentum is used only in the first step, but $\gamma_t = 0$ and $\beta_t = 0$ for all $t > 1$. For any $\beta_1 > 0$, there exists $\{\eta_t\}_t$, such that $w_{(t)}$ from (8) converges to a global minimum, but not to $\underset{w \in \mathcal{G}}{\text{argmin}} \psi(w)$.*

2.3. Natural gradient descent

Natural gradient descent (NGD) was introduced by Amari (1998) as a modification of gradient descent, wherein the updates are chosen to be the steepest descent direction w.r.t a Riemannian metric tensor H that maps w to a positive definite local metric $H(w)$. The updates are given by,

$$w_{(t+1)} = w_{(t)} - \eta_t H(w_{(t)})^{-1} \nabla\mathcal{L}(w_{(t)}) \quad (9)$$

In many instances, the metric tensor H is specified by the Hessian $\nabla^2\psi$ of a strongly convex potential ψ . For example, when the metric over the Riemannian manifold is the KL divergence between distributions P_w and $P_{w'}$ parameterized by w , the metric tensor is given by $H(w) = \nabla^2\psi(P_w)$ where the potential ψ is the entropy potential over P_w .

Connection to mirror descent When $H(w) = \nabla\psi^2(w)$ for a strongly convex potential ψ , as the step size η goes to

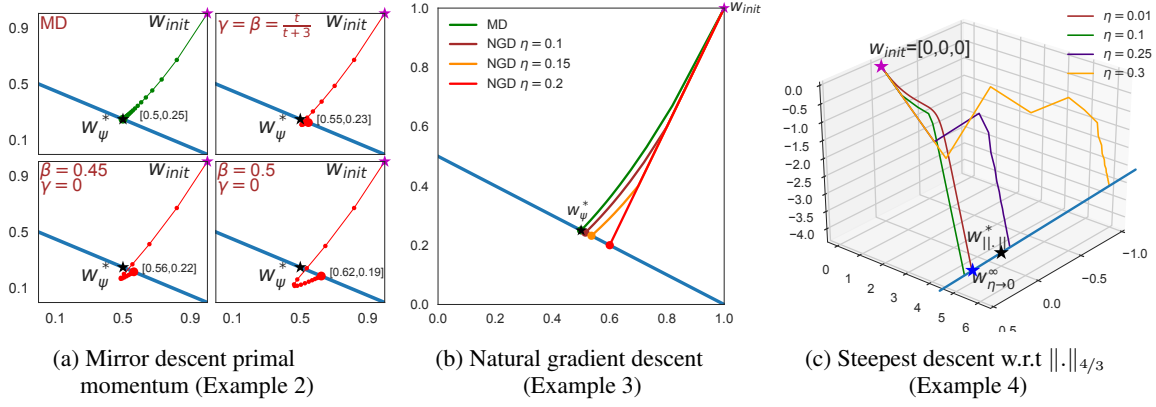


Figure 1: Dependence of implicit bias on step size and momentum: In (a)–(c), the blue line denotes the set \mathcal{G} of global minima for the respective examples. In (a) and (b) ψ is the entropy potential and all algorithms are initialized with $w_{(0)} = [1, 1]$ such that $\psi(w_{(0)}) = \operatorname{argmin}_w \psi(w)$ and $w_\psi^* = \operatorname{argmin}_{\psi \in \mathcal{G}} \psi(w)$ denotes the minimum potential global minima we expect to converge to. (a) **Mirror descent with primal momentum (Example 2)**: the global minimum that (8) converges to depends on the momentum parameters—the sub-plots contain the trajectories of (8) for different choices of β and γ ; (b) **Natural gradient descent (Example 3)**: for different choices of step size η , (9) converges to different global minima. Here, η was chosen to be small enough to ensure $w_{(t)} \in \operatorname{dom}(\psi)$. (c) **Steepest descent w.r.t. $\|\cdot\|_{4/3}$ (Example 4)**: the global minimum to which (11) converges depends on η . Here $w_{(0)} = [0, 0, 0]$, the minimum norm global minima is denoted as $w_{\|\cdot\|}^* = \operatorname{argmin}_{\psi \in \mathcal{G}} \|w\|_{4/3}$, and $w_{\eta \rightarrow 0}^\infty$ is the solution of infinitesimal SD with $\eta \rightarrow 0$.

zero, the iterates $w_{(t)}$ from natural gradient descent in eq. (9) and mirror descent w.r.t ψ in eq. (4) converge to each other and the common dynamics in the limit is given by:

$$\begin{aligned} \frac{d\nabla\psi(w_{(t)})}{dt} &= -\nabla\mathcal{L}(w_{(t)}) \\ \implies \frac{dw_{(t)}}{dt} &= -\nabla^2\psi(w_{(t)})^{-1}\nabla\mathcal{L}(w_{(t)}) \end{aligned} \quad (10)$$

So as the step sizes are made infinitesimal, the limit point of natural gradient descent $w_\infty = \lim_{t \rightarrow \infty} w_{(t)}$ is also the limit point of mirror descent and hence will be biased towards solutions with minimum divergence to the initialization, i.e., as $\eta \rightarrow 0$, $w_\infty \rightarrow \operatorname{argmin}_{w \in \mathcal{G}} D_\psi(w, w_{(0)})$.

For general step sizes $\{\eta_t\}$, if the potential ψ is quadratic, $\psi(w) = 1/2\|w\|_D^2$ for some positive definite D , we get linear link functions $\nabla\psi(w) = Dw$ and constant metric tensors $\nabla^2\psi(w) = H(w) = D$, and the mirror descent (5) and natural gradient descent (9) updates are exactly equal for all values of $\{\eta_t\}$. Otherwise the updates differ in that (9) is only an approximation of the mirror descent update $\nabla\psi^{-1}(\nabla\psi(w_{(t)}) - \eta_t\nabla\mathcal{L}(w_{(t)}))$.

For natural gradient descent with finite step size and non-quadratic potentials ψ , the characterization in eq. (6) generally does not hold. We can see this as, for any initialization $w_{(0)}$, a finite $\eta_1 > 0$ will easily lead to $w_{(1)}$ for which the dual variable $\nabla\psi(w_{(1)})$ is no longer in the data manifold $\operatorname{span}(\{x_n\}) + \nabla\psi(w_{(0)})$, and hence will converge to a different global minimum dependent on the step sizes $\{\eta_t\}$.

Example 3. Consider optimizing $\mathcal{L}(w)$ with squared loss over dataset $\{(x_1 = [1, 2], y_1 = 1)\}$ using the natu-

ral gradient descent w.r.t. the metric tensor given by, $H(w) = \nabla^2\psi(w)$, where $\psi(w) = \sum_i w[i] \log w[i] - w[i]$, and initialization $w_{(0)} = [1, 1]$. Figure 1b shows that NGD with different step sizes η converges to different global minima. For a simple analytical example: take one finite step $\eta_1 > 0$ and then follow the continuous time path in eq. (10).

Proposition 3a. For almost all $\eta_1 > 0$, $\lim_{t \rightarrow \infty} w_{(t)} = \operatorname{argmin}_{w \in \mathcal{G}} D_\psi(w, w_{(1)}) \neq \operatorname{argmin}_{w \in \mathcal{G}} D_\psi(w, w_{(0)})$.

2.4. Steepest Descent

Gradient descent is also a special case of steepest descent (SD) w.r.t a generic norm $\|\cdot\|$ (Boyd & Vandenberghe, 2004) with updates given by,

$$\begin{aligned} w_{(t+1)} &= w_{(t)} + \eta_t \Delta w_{(t)}, \\ \text{where } \Delta w_{(t)} &= \operatorname{argmin}_v \langle \nabla\mathcal{L}(w_{(t)}), v \rangle + \frac{1}{2}\|v\|^2. \end{aligned} \quad (11)$$

The optimality of $\Delta w_{(t)}$ in eq. (11) requires $-\nabla\mathcal{L}(w_{(t)}) \in \partial\|\Delta w_{(t)}\|^2$, which is equivalent to,

$$\langle \Delta w_{(t)}, -\nabla\mathcal{L}(w_{(t)}) \rangle = \|\Delta w_{(t)}\|^2 = \|\nabla\mathcal{L}(w_{(t)})\|_*^2. \quad (12)$$

Examples of steepest descent include gradient descent, which is steepest descent w.r.t ℓ_2 norm and coordinate descent, which is steepest descent w.r.t ℓ_1 norm. In general, the update $\Delta w_{(t)}$ in eq. (11) is not uniquely defined and there could be multiple direction $\Delta w_{(t)}$ that minimize eq. (11). In such cases any minimizer of eq. (11) is a valid steepest descent update and satisfies eq. (12).

Generalizing gradient descent, we might expect the limit point w_∞ of steepest descent w.r.t an arbitrary norm $\|\cdot\|$.

to be the solution closest to initialization in corresponding norm, $\operatorname{argmin}_{w \in \mathcal{G}} \|w - w_{(0)}\|$. This is indeed the case for quadratic norms $\|v\|_D = \sqrt{v^\top D v}$ when eq. 11 is equivalent to mirror descent with $\psi(w) = 1/2 \|w\|_D^2$. Unfortunately, this does not hold for general norms.

Example 4. Consider minimizing $\mathcal{L}(w)$ with dataset $\{(x_1 = [1, 1, 1], y_1 = 1), (x_1 = [1, 2, 0], y_1 = 10)\}$ and loss $\ell(u, y) = (u - y)^2$ using steepest descent updates w.r.t. the $\ell_{4/3}$ norm. The empirical results for this problem in Figure 1c clearly show that even for ℓ_p norms where the $\|\cdot\|_p^2$ is smooth and strongly convex, the corresponding steepest descent converges to a global minimum that depends on the step size. Further, even in the continuous step size limit of $\eta \rightarrow 0$, $w_{(t)}$ does not converge to $\operatorname{argmin}_{w \in \mathcal{G}} \|w - w_{(0)}\|$.

Coordinate descent Steepest descent w.r.t. the ℓ_1 norm can be written as coordinate descent, with updates:

$$\Delta w_{(t+1)} \in \operatorname{conv} \left\{ -\eta_t \frac{\partial \mathcal{L}(w)}{\partial w_{[j_t]}} e_{j_t} : j_t = \operatorname{argmax}_j \left| \frac{\partial \mathcal{L}(w)}{\partial w_{[j]}} \right| \right\},$$

where $\operatorname{conv}(S)$ denotes the convex hull of the set S , and $\{e_j\}$ are the standard basis. Thus, when multiple partial derivatives are maximal, we can choose any convex combination of the maximizing directions, leading to many possible coordinate descent optimization paths.

The connection between optimization paths of coordinate descent and the ℓ_1 regularization path given by, $\hat{w}(\lambda) = \operatorname{argmin}_w \mathcal{L}(w) + \lambda \|w\|_1$, has been studied by Efron et al. (2004). The specific coordinate descent path where updates are along the average of all optimal coordinates and the step sizes are infinitesimal is equivalent to forward stage-wise selection, a.k.a. ϵ -boosting (Friedman, 2001). When the ℓ_1 regularization path $\hat{w}(\lambda)$ is monotone in each of the coordinates, it is identical to this stage-wise selection path, i.e., to a coordinate descent optimization path (and also to the related LARS path) (Efron et al., 2004). In this case, at the limit of $\lambda \rightarrow 0$ and $t \rightarrow \infty$, the optimization and regularization paths, both converge to the minimum ℓ_1 norm solution. However, when the regularization path $\hat{w}(\lambda)$ is not monotone, which can and does happen, the optimization and regularization paths diverge, and forward stage-wise can converge to solutions with sub-optimal ℓ_1 norm. This matches our understanding that steepest descent w.r.t. a norm $\|\cdot\|$, in this case the ℓ_1 norm, might converge to a solution that is *not* the minimum $\|\cdot\|$ norm solution.

2.5. Summary for losses with a unique finite root

For losses with a unique finite root, we characterized the implicit bias of generic mirror descent algorithm in terms of the potential function and initialization. This characterization extends for momentum in the dual space as well as to natural gradient descent in the limit of infinitesimal step

size. We also saw that the characterization breaks for mirror descent with primal momentum and natural gradient descent with finite step sizes. Moreover, for steepest descent with general norms, we were unable to get a useful characterization even in the infinitesimal step size limit. In the following section, we will see that for strictly monotone losses, we can get a characterization also for steepest descent.

3. Strictly Monotone Losses

We now turn to strictly monotone loss functions ℓ where the behavior of the implicit bias is fundamentally different, and as are the situations when the implicit bias can be characterized. Such losses are common in classification problems where $y = \{-1, 1\}$ and $\ell(f(x), y)$ is typically a continuous surrogate of the 0-1 loss. Examples of such losses include logistic loss, exponential loss, and probit loss.

Property 2 (Strict monotone losses). $\ell(\hat{y}, y)$ is bounded from below, and $\forall y$, $\ell(\hat{y}, y)$ is strictly monotonically decreasing in \hat{y} . Without loss of generality, $\forall y$, $\inf_{\hat{y}} \ell(\hat{y}, y) = 0$ and $\ell(\hat{y}, y) \xrightarrow{\hat{y} \rightarrow \infty} 0$.

We look at classification models that fit the training data $\{x_n, y_n\}$ with linear decision boundaries $u(x) = \langle w, x \rangle$ with decision rule given by $\hat{y}(x) = \operatorname{sign}(u(x))$. In many instances of the proofs, we also assume without loss of generality that $y_n = 1$ for all n , since for linear models, the sign of y_n can equivalently be absorbed into x_n .

When the training data $\{x_n, y_n\}$ is not linearly separable, the empirical objective $\mathcal{L}(w)$ in eq. (1) can have a finite global minimum. However, if the data set $\{x_n, y_n\}_{n=1}^N$ is linearly separable, the empirical loss $\mathcal{L}(w)$ in eq. (1) is again ill-posed, and moreover $\mathcal{L}(w)$ does not have any finite minimizer, i.e., $\mathcal{L}(w) \rightarrow 0$ only as $\|w\| \rightarrow \infty$. Then for iterates $w_{(t)}$ from any algorithm, if $\mathcal{L}(w_{(t)}) \rightarrow 0$, then the iterates diverge to infinity rather than converge and hence, we cannot talk about $\lim_{t \rightarrow \infty} w_{(t)}$. Instead, we look at the limit direction $\bar{w}_\infty = \lim_{t \rightarrow \infty} \frac{w_{(t)}}{\|w_{(t)}\|}$ whenever the limit exists. When the limit exists, we refer to this as convergence in direction. For classification problems, the limit direction is all we care about as it entirely specifies the separating hyperplane or the decision rule, and hence the generalization properties with respect to 0-1 error.

We focus on the exponential loss $\ell(u, y) = \exp(-uy)$. However, our results can be extended to loss functions with tight exponential tails, including logistic and sigmoid losses, along the lines of Soudry et al. (2017) and Telgarsky (2013).

3.1. Gradient descent

Soudry et al. (2017) showed that for almost all linearly separable datasets, gradient descent with *any initialization and any bounded step size* converges in direction to maximum

margin separator with unit ℓ_2 norm, i.e., the hard margin support vector machine classifier,

$$\bar{w}_\infty = \lim_{t \rightarrow \infty} \frac{w(t)}{\|w(t)\|_2} = w_{\|\cdot\|_2}^* := \operatorname{argmax}_{\|w(t)\|_2 \leq 1} \min_n y_n \langle w, x_n \rangle.$$

This characterization of the implicit bias is independent of both the step size as well as the initialization. We already see a fundamental difference from the implicit bias of gradient descent for losses with a unique finite root (Section 2.1) where the characterization depended on the initialization.

Can we similarly characterize the implicit bias of different algorithms establishing $w(t)$ converges in direction and calculating \bar{w}_∞ ? Can we do this even when we *could not* characterize the limit point $w_\infty = \lim_{t \rightarrow \infty} w(t)$ for losses with unique finite roots? As we will see in the following section, we can indeed answer these questions for steepest descent w.r.t arbitrary norms.

3.2. Steepest Descent

Recall that for the squared loss, the limit point of steepest descent could depend strongly on the step size, and we were unable obtain a useful characterization even for infinitesimal step size. In contrast, the following theorem provides a crisp characterization of the limit direction of steepest descent as a maximum margin solution, independent of step size (as long as it is small enough) and initialization. Let $\|\cdot\|_*$ denote the dual norm of $\|\cdot\|$.

Theorem 5. *For any separable dataset $\{x_n, y_n\}_{n=1}^N$ and any norm $\|\cdot\|$, consider the steepest descent updates from eq. (12) for minimizing $\mathcal{L}(w)$ in eq. (1) with the exponential loss $\ell(u, y) = \exp(-uy)$. For all initializations $w(0)$, and all bounded step sizes satisfying $\eta_t \leq \max\{\eta_+, \frac{1}{B^2 \mathcal{L}(w(t))}\}$, where $B := \max_n \|x_n\|_*$ and $\eta_+ < \infty$ is some maximum step size bound, the iterates $w(t)$ satisfy the following,*

$$\lim_{t \rightarrow \infty} \min_n \frac{y_n \langle w(t), x_n \rangle}{\|w(t)\|} = \max_{w: \|w\| \leq 1} \min_n y_n \langle w, x_n \rangle.$$

In particular, if there is a unique maximum- $\|\cdot\|$ margin solution $w_{\|\cdot\|}^ = \operatorname{argmax}_w \min_n \frac{y_n \langle w, x_n \rangle}{\|w\|}$, then the limit direction converges to it: $\bar{w}_\infty = \lim_{t \rightarrow \infty} \frac{w(t)}{\|w(t)\|} = w_{\|\cdot\|}^*$.*

A special case of Theorem 5 is for steepest descent w.r.t. the ℓ_1 norm, which as we already saw corresponds to coordinate descent. More specifically, coordinate descent on the exponential loss can be thought of as an alternative presentation of AdaBoost (Schapire & Freund, 2012), where each coordinate represents the output of one ‘‘weak learner’’. Indeed, initially mysterious generalization properties of boosting have been understood in terms of implicit ℓ_1 regularization (Schapire & Freund, 2012), and later on AdaBoost with

small enough step size was shown to converge in direction precisely to the maximum ℓ_1 margin solution (Zhang et al., 2005; Shalev-Shwartz & Singer, 2010; Telgarsky, 2013), just as guaranteed by Theorem 5. In fact, Telgarsky (2013) generalized the result to a richer variety of exponential tailed loss functions including logistic loss, and a broad class of non-constant step size rules. Interestingly, coordinate descent with exact line search might result in step sizes that are too large leading the iterates to converge to a different direction that is not a max- ℓ_1 -margin direction (Rudin et al., 2004), hence the maximum step size bound in Theorem 5.

Theorem 5 is a generalization of the result of Telgarsky to steepest descent with respect to other norms, and our proof follows the same strategy as Telgarsky. We first prove a generalization of the duality result of Shalev-Shwartz & Singer (2010): if there is a unit norm linear separator that achieves margin γ , then $\|\nabla \mathcal{L}(w)\|_* \geq \gamma \mathcal{L}(w)$ for all w . By using this lower bound on the dual norm of the gradient, we are able to show that the loss decreases faster than the increase in the norm of the iterates, establishing convergence in a margin maximizing direction.

In relating the optimization path to the regularization path, it is also relevant to relate Theorem 5 to the result by Rosset et al. (2004) that for monotone loss functions and ℓ_p norms, the ℓ_p regularization path $\hat{w}(c) = \operatorname{argmin}_{w: \|w\|_p \leq c} \mathcal{L}(w(t))$ also converges in direction to the maximum margin separator, i.e., $\lim_{c \rightarrow \infty} \hat{w}(c) = w_{\|\cdot\|_p}^*$. Although the optimization path and regularization path are not the same, they both converge to the same max-margin separator in the limits of $c \rightarrow \infty$ and $t \rightarrow \infty$, for the regularization path and steepest descent optimization path, respectively.

3.3. Adaptive Gradient Descent (AdaGrad)

Adaptive gradient methods, such as AdaGrad (Duchi et al., 2011) or Adam (Kingma & Adam, 2015) are very popular for neural network training. We look at the implicit bias of AdaGrad in this section. To examine this, we focus on the basic (diagonal) AdaGrad,

$$w_{(t+1)} = w_{(t)} - \eta \mathbf{G}_{(t)}^{-1/2} \nabla \mathcal{L}(w_{(t)}), \quad (13)$$

where $\mathbf{G}_{(t)}$ is a diagonal matrix such that,

$$\forall i: \mathbf{G}_{(t)}[i, i] = \sum_{u=0}^t (\nabla \mathcal{L}(w_{(u)})[i])^2. \quad (14)$$

AdaGrad updates described above corresponds to a pre-conditioned gradient descent, except the pre-conditioning matrix $\mathbf{G}_{(t)}$ adapts across iterations. It was observed by Wilson et al. (2017) that for neural networks with squared loss, adaptive methods tend to degrade generalization performance in comparison to non-adaptive methods (e.g., SGD with momentum), even when both methods are used to train

the network until convergence to a global minimum of training loss. This suggests that adaptivity does indeed affect the implicit bias. For squared loss, by inspection the updates in eq. (13), we do not expect to get a characterization of the limit point w_∞ that is independent of the step sizes.

However, we might hope that, like for steepest descent, the situation might be different for strictly monotone losses, where the asymptotic behavior could potentially nullify the initial conditions. Examining the updates in eq. (13), we can hypothesize that the robustness to initialization and initial updates depend on whether the matrices $\mathbf{G}_{(t)}$ diverge or converge: if $\mathbf{G}_{(t)}$ diverges, then we expect the asymptotic effects to dominate, while if it converges, then the limit direction will indeed depend on the initial conditions.

Unfortunately, the following proposition shows that, the components of $\mathbf{G}_{(t)}$ matrix are bounded, and hence even for strict monotone losses, the initial conditions $w_{(0)}$, $\mathbf{G}_{(0)}$ and step size η will have a non-vanishing contribution to the asymptotic behavior of $\mathbf{G}_{(t)}$ and hence to the limit direction $\bar{w}_\infty = \lim_{t \rightarrow \infty} \frac{w_{(t)}}{\|w_{(t)}\|}$, whenever it exists. In other words, the implicit bias of AdaGrad does indeed depend on the initial conditions, including initialization and step size.

Proposition 6. *For any training data $\{x_n, y_n\}$, consider the AdaGrad iterates $w_{(t)}$ from eq. (13) for minimizing $\mathcal{L}(w)$ with exponential loss $\ell(u, y) = \exp(-uy)$. For any fixed and bounded step size $\eta < \infty$, and any initialization of $w_{(0)}$ and $\mathbf{G}_{(0)}$, such that $\frac{\eta}{2}\mathcal{L}(w_{(0)}) < 1$, and $\|\mathbf{G}_{(0)}^{-1/4} x_n\|_2 \leq 1, \forall i, \forall t : \mathbf{G}_{(t)}[i, i] < \infty$.*

4. Gradient descent on the factorized parameterization

Consider the empirical risk minimization in eq. (1) for matrix valued $X_n \in \mathbb{R}^{d \times d}$, $W \in \mathbb{R}^{d \times d}$

$$\min_W \mathcal{L}(W) = \ell(\langle W, X_n \rangle; y_n). \quad (15)$$

This is the exact same setting as eq. (1) obtained by arranging w and x_n as matrices. We can now study another class of optimization algorithms for learning linear models based on matrix factorization where we reparameterize W as $W = UV^\top$ with unconstrained $U \in \mathbb{R}^{d \times d}$ and $V \in \mathbb{R}^{d \times d}$ to get the following equivalent objective,

$$\min_{U, V} \mathcal{L}(UV^\top) = \sum_{n=1}^N \ell(\langle UV^\top, X_n \rangle; y_n) \quad (16)$$

Note that although non-convex eq. (16) is equivalent to eq. (15) with the exact same set of global minima over $W = UV^\top$. Gunasekar et al. (2017) studied this problem for squared loss $\ell(u, y) = (u - y)^2$ and noted that gradient descent on the factorization yields radically different

implicit bias compared to gradient descent on W . In particular, gradient descent on U, V is observed to be biased towards low nuclear norm solutions, which in turns ensures generalization (Srebro et al., 2005) and low rank matrix recovery (Recht et al., 2010; Candes & Recht, 2009). Since the matrix factorization objective in eq. (16) can be viewed as a two-layer neural network with linear activation, understanding the implicit bias then could provide insights into characterizing the implicit bias in more complex neural networks with non-linear activation.

Gunasekar et al. (2017) noted that, the optimization problem over non-p.s.d. factorization $W = UV^\top$ in eq. (16) is a special case of the optimizing \mathcal{L} over p.s.d. matrices $\tilde{W} \succcurlyeq 0$ parameterized using unconstrained symmetric factorization $W = UU^\top$ with $U \in \mathbb{R}^{d \times d}$:

$$\min_{U \in \mathbb{R}^{d \times d}} \bar{\mathcal{L}}(U) = \mathcal{L}(UU^\top) = \sum_{n=1}^N \ell(\langle UU^\top, X_n \rangle; y_n) \quad (17)$$

In particular, both the objective as well as gradient descent updates of eq. (16) can be derived as an instance of the problem in eq. (17) over a larger p.s.d. matrix $\tilde{W} = \tilde{U}\tilde{U}^\top = \begin{bmatrix} A_1 & W \\ W^\top & A_2 \end{bmatrix}$, and data matrices $\tilde{X}_n = \begin{bmatrix} 0 & X_n \\ X_n^\top & 0 \end{bmatrix}$.

Henceforth, we will also consider the symmetric matrix factorization in (17). Let $U_{(0)} \in \mathbb{R}^{d \times d}$ be any full rank initialization, gradient descent updates in U are given by,

$$U_{(t+1)} = U_{(t)} - \eta_t \nabla \bar{\mathcal{L}}(U_{(t)}), \quad (18)$$

with corresponding updates in $W_{(t)} = U_{(t)}U_{(t)}^\top$ given by,

$$W_{(t+1)} = W_{(t)} - \eta_t [\nabla \mathcal{L}(W_{(t)})W_{(t)} + W_{(t)}\nabla \mathcal{L}(W_{(t)})] + \eta_t^2 \nabla \mathcal{L}(W_{(t)})W_{(t)}\nabla \mathcal{L}(W_{(t)}) \quad (19)$$

Losses with a unique finite root For squared loss, Gunasekar et al. (2017) showed that the implicit bias of iterates in eq. (19) crucially depended on both the initialization $U_{(0)}$ as well as the step size η . Gunasekar et al. conjectured, and provided theoretical and empirical evidence that gradient descent on the factorization converges to the minimum nuclear norm global minimum, but only if the initialization is infinitesimally close to zero and the step-sizes are infinitesimally small. Li et al. (2017), later proved the conjecture under additional assumption that the measurements X_n satisfy certain *restricted isometry property (RIP)*.

In the case of squared loss, it is evident that for finite step sizes and finite initialization, the implicit bias towards the minimum nuclear norm global minima is not exact. In practice, not only do we need $\eta > 0$, but we also cannot initialize very close to zero since zero is a saddle point for eq. (17). The natural question motivated by the results in Section 3 is: for strictly monotone losses, can we get a

characterization of the implicit bias of gradient descent for the factorized objective in eq. (17) that is more robust to initialization and step size?

Strict monotone losses In the following theorem, we again see that the characterization of the implicit bias of gradient descent for factorized objective is more robust in the case of strict monotone losses.

Theorem 7. *For almost all linearly separable datasets $\{X_n, y_n\}_{n=1}^N$, consider the gradient descent iterates $U_{(t)}$ in eq. (18) for minimizing $\tilde{\mathcal{L}}(U)$ with the exponential loss $\ell(u, y) = \exp(-uy)$ and the corresponding sequence of linear predictors $W_{(t)}$ in eq. (19). For any full rank initialization $U_{(0)}$ and any sufficiently small step size sequences η_t such that η_t is smaller than the local Lipschitz at $W_{(t)}$, if $W_{(t)}$ converges to a global minimum i.e., $\mathcal{L}(W_{(t)}) \rightarrow 0$, and additionally the incremental updates $W_{(t+1)} - W_{(t)}$ and the gradients $\nabla \mathcal{L}(W_{(t)})$ converge in direction, then the limit direction $\bar{W}_\infty = \lim_{t \rightarrow \infty} \frac{W_{(t)}}{\|W_{(t)}\|_*}$ exists, and is given by the maximum margin separator with unit nuclear norm $\|\cdot\|_*$,*

$$\bar{W}_\infty = \operatorname{argmax}_{W \succ 0} \min_n y_n \langle W, X_n \rangle \text{ s.t., } \|W\|_* \leq 1.$$

Here we note that convergence of $\frac{W_{(t)}}{\|W_{(t)}\|}$ is necessary for the characterization of implicit bias to be relevant, but in Theorem 7 we require stronger conditions that the incremental updates $W_{(t+1)} - W_{(t)}$ and the gradients $\nabla \mathcal{L}(W_{(t)})$ converge in direction, which might not hold in general. Relaxing this condition is of interest for future work.

Key property Let us look at exponential loss when $W_{(t)}$ converges in direction to, say \bar{W}_∞ as $W_{(t)} = \bar{W}_\infty g(t) + \rho(t)$ for some scalar $g(t) \rightarrow \infty$ and $\frac{\rho(t)}{g(t)} \rightarrow 0$. Consequently, the gradients $\nabla \mathcal{L}(W_{(t)}) = \sum_n e^{-g(t)y_n \langle W_\infty, X_n \rangle} e^{-y_n \langle \rho(t), X_n \rangle} y_n X_n$ will asymptotically be dominated by linear combinations of examples X_n that have the smallest distance to the decision boundary, i.e., the support vectors of W_∞ . This behavior can be used to show optimality of W_∞ to the maximum margin solution subject to nuclear norm constraint in Theorem 7.

This idea formalized in the following lemma, which is of interest beyond the results in this paper.

Lemma 8. *For almost all linearly separable datasets $\{x_n, y_n\}_n$, consider any sequence $w_{(t)}$ that minimizes $\mathcal{L}(w)$ in eq. (1) with exponential loss, i.e., $\mathcal{L}(w_{(t)}) \rightarrow 0$. If $\frac{w_{(t)}}{\|w_{(t)}\|}$ converges, then for every accumulation point z_∞ of $\left\{ \frac{-\nabla \mathcal{L}(w_{(t)})}{\|\nabla \mathcal{L}(w_{(t)})\|} \right\}_t$, $\exists \{\alpha_n \geq 0\}_{n \in S}$ s.t., $z_\infty = \sum_{n \in S} \alpha_n y_n x_n$, where $\bar{w}_\infty = \lim_{t \rightarrow \infty} \frac{w_{(t)}}{\|w_{(t)}\|}$ and $S = \{n : y_n \langle \bar{w}_\infty, x_n \rangle = \min_n y_n \langle \bar{w}_\infty, x_n \rangle\}$ are the indices of the data points with smallest margin to \bar{w}_∞ .*

5. Summary

We studied the implicit bias of different optimization algorithms for two families of losses—losses with a unique finite root, and strict monotone losses, where the biases are fundamentally different. In the case of losses with a unique finite root, we have a simple characterization of the limit point $w_\infty = \lim_{t \rightarrow \infty} w_{(t)}$ for mirror descent, but for this family of losses, such a succinct characterization does not extend to steepest descent with respect to general norms. On the other hand, for strict monotone losses, we noticed that the initial updates of the algorithm, including initialization and initial step sizes are nullified when we analyze the asymptotic limit direction $\bar{w}_\infty = \lim_{t \rightarrow \infty} \frac{w_{(t)}}{\|w_{(t)}\|}$. In particular, we show that for steepest descent, the limit direction is a maximum margin separator within the unit ball of the corresponding norm. We also looked at other optimization algorithms for strictly monotone losses. For matrix factorization, we again get a more robust characterization of the implicit bias as the maximum margin separator with unit nuclear norm. This again, in contrast to squared loss Gunasekar et al. (2017), is independent of the initialization and step size. However, for AdaGrad, we show that even for strict monotone losses, the limit direction \bar{w}_∞ could depend on the initial conditions.

In our results, we characterize the implicit bias for linear models as minimum norm (potential) or maximum margin solutions. These are indeed very special among all the solutions that fit the training data, and in particular, their generalization performance can in turn be understood from standard analyses (Bartlett & Mendelson, 2003).

For more complicated non-linear models, especially neural networks, further work is required in order to get a more complete understanding of the implicit bias. The preliminary result for matrix factorization provides us tools to attempt extensions to multi-layer linear models, and eventually to non-linear networks. Even for linear models, the question of what is the implicit bias is when $\mathcal{L}(w)$ is optimized with explicitly constraints $w \in \mathcal{W}$ is an open problem. We believe similar characterizations can be obtained when there are multiple feasible solutions with $\mathcal{L}(w) = 0$. We also believe, the results for single outputs considered in this paper can also be extended for multi-output loss functions.

Finally, we would like a more fine grained analysis connecting the iterates $w_{(t)}$ along the optimization path of various algorithms to the estimates along regularization path, $\hat{w}(c) = \operatorname{argmin}_{\mathcal{R}(w) \leq c} \mathcal{L}(w)$, where an explicit regularization is added to the optimization objective. In particular, what we show in this paper is that the optimization path and the regularization path meet in their limit points, $t \rightarrow \infty$ and $c \rightarrow \infty$, respectively. It would be desirable to further understand the relations between the entire optimization and regularization paths, which will help us understand the non-asymptotic effects from early stopping.

Acknowledgments The authors are grateful to M.S. Nacson, Y. Carmon, and the anonymous ICML reviewers for helpful comments on the manuscript. The research was supported in part by NSF IIS award 1302662. The work of DS was supported by the Taub Foundation.

References

- Amari, S. I. Natural gradient works efficiently in learning. *Neural computation*, 1998.
- Bartlett, P. L. and Mendelson, S. Rademacher and Gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 2003.
- Beck, A. and Teboulle, M. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 2003.
- Boyd, S. and Vandenberghe, L. *Convex optimization*. Cambridge university press, 2004.
- Bregman, L. M. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR computational mathematics and mathematical physics*, 1967.
- Candes, E. J. and Recht, B. Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 2009.
- Duchi, J., Hazan, E., and Singer, Y. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 2011.
- Efron, B., Hastie, T., Johnstone, I., and Tibshirani, R. Least angle regression. *The Annals of statistics*, 2004.
- Friedman, J. H. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, 2001.
- Gunasekar, S., Woodworth, B. E., Bhojanapalli, S., Neyshabur, B., and Srebro, N. Implicit regularization in matrix factorization. In *Advances in Neural Information Processing Systems*, pp. 6152–6160, 2017.
- Hoffer, E., Hubara, I., and Soudry, D. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. In *NIPS*, pp. 1–13, may 2017. URL <http://arxiv.org/abs/1705.08741>.
- Keskar, N. S., Mudigere, D., Nocedal, J., Smelyanskiy, M., and Tang, P. T. P. On large-batch training for deep learning: Generalization gap and sharp minima. In *International Conference on Learning Representations*, 2016.
- Kingma, D. and Adam, J. B. Adam: A method for stochastic optimisation. In *International Conference for Learning Representations*, volume 6, 2015.
- Kivinen, J. and Warmuth, M. K. Exponentiated gradient versus gradient descent for linear predictors. *Information and Computation*, 1997.
- Li, Y., Ma, T., and Zhang, H. Algorithmic regularization in over-parameterized matrix recovery. *arXiv preprint arXiv:1712.09203*, 2017.
- Muresan, M. and Muresan, M. *A concrete approach to classical analysis*, volume 14. Springer, 2009.
- Nemirovskii, A. and Yudin, D. *Problem complexity and method efficiency in optimization*. Wiley, 1983.
- Nesterov, Y. A method of solving a convex programming problem with convergence rate $o(1/k^2)$. In *Soviet Mathematics Doklady*, 1983.
- Neyshabur, B., Salakhutdinov, R. R., and Srebro, N. Pathsgd: Path-normalized optimization in deep neural networks. In *Advances in Neural Information Processing Systems*, pp. 2422–2430, 2015a.
- Neyshabur, B., Tomioka, R., and Srebro, N. In search of the real inductive bias: On the role of implicit regularization in deep learning. In *International Conference on Learning Representations*, 2015b.
- Neyshabur, B., Tomioka, R., Salakhutdinov, R., and Srebro, N. Geometry of optimization and implicit regularization in deep learning. *arXiv preprint arXiv:1705.03071*, 2017.
- Polyak, B. T. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 1964.
- Recht, B., Fazel, M., and Parrilo, P. A. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM review*, 2010.
- Ross, K. A. *Elementary analysis*. Springer, 1980.
- Rosset, S., Zhu, J., and Hastie, T. Boosting as a regularized path to a maximum margin classifier. *Journal of Machine Learning Research*, 2004.
- Rudin, C., Daubechies, I., and Schapire, R. E. The dynamics of adaboost: Cyclic behavior and convergence of margins. *Journal of Machine Learning Research*, 5(Dec):1557–1595, 2004.
- Schapire, R. E. and Freund, Y. *Boosting: Foundations and algorithms*. MIT press, 2012.
- Shalev-Shwartz, S. and Singer, Y. On the equivalence of weak learnability and linear separability: New relaxations and efficient boosting algorithms. *Machine learning*, 80(2-3):141–163, 2010.

- Smith, Kindermans, L. Don't Decay the Learning Rate, Increase the Batch Size. In *ICLR*, 2018.
- Soudry, D., Hoffer, E., and Srebro, N. The implicit bias of gradient descent on separable data. *arXiv preprint arXiv:1710.10345*, 2017.
- Srebro, N., Alon, N., and Jaakkola, T. S. Generalization error bounds for collaborative prediction with low-rank matrices. In *Advances In Neural Information Processing Systems*, pp. 1321–1328, 2005.
- Telgarsky, M. Margins, shrinkage and boosting. In *Proceedings of the 30th International Conference on International Conference on Machine Learning-Volume 28*, pp. II–307. JMLR. org, 2013.
- Wilson, A. C., Roelofs, R., Stern, M., Srebro, N., and Recht, B. The marginal value of adaptive gradient methods in machine learning. In *Advances in Neural Information Processing Systems*, pp. 4151–4161, 2017.
- Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations*, 2017.
- Zhang, T., Yu, B., et al. Boosting with early stopping: Convergence and consistency. *The Annals of Statistics*, 33(4):1538–1579, 2005.

A. Losses with a unique finite root

Let $\mathcal{P}_{\mathcal{X}} = \text{span}(\{x_n : n \in [N]\}) = \{\sum_n \nu_n x_n : \nu_n \in \mathbb{R}\}$ and $\ell'(u, y)$ be the derivative of ℓ w.r.t u . We have,

$$\forall w \in \mathbb{R}^d, \nabla \mathcal{L}(w) = \sum_{n=1}^N \ell'(\langle w, x_n \rangle; y_n) x_n \in \mathcal{P}_{\mathcal{X}}. \quad (20)$$

A.1. Proof of Theorem 1-1a

For a strongly convex potential ψ , denote the global optimum with minimum Bregman divergence $D_{\psi}(\cdot, w_{(0)})$ to the initialization $w_{(0)}$ as

$$w_{\psi}^* = \underset{w}{\operatorname{argmin}} D_{\psi}(w, w_{(0)}) \text{ s.t., } \forall n, \langle w, x_n \rangle = y_n, \quad (21)$$

where recall that $D_{\psi}(w, w_{(0)}) = \psi(w) - \psi(w_{(0)}) - \langle \nabla \psi(w_{(0)}), w - w_{(0)} \rangle$.

The KKT optimality conditions for (21) are as follows,

$$\begin{aligned} \text{Stationarity:} \quad & \nabla \psi(w_{\psi}^*) - \nabla \psi(w_{(0)}) \in \mathcal{P}_{\mathcal{X}}, \text{ or } \exists \{\nu_n\}_{n=1}^N \text{ s.t., } \nabla \psi(w_{\psi}^*) - \nabla \psi(w_{(0)}) = \sum_{n=1}^N \nu_n x_n \\ \text{Primal feasibility:} \quad & \forall n, \langle w_{\psi}^*, x_n \rangle = y_n, \text{ or } w_{\psi}^* \in \mathcal{G} \end{aligned} \quad (22)$$

Recall Theorem 1–1a from Section 2.2.

Theorem 1. *For any loss ℓ with a unique finite root (Property 1), any initialization $w_{(0)}$, any step size sequence $\{\eta_t\}$, and any strongly convex potential ψ , consider the corresponding mirror descent iterates $w_{(t)}$ from eq. (5). If the limit point of the iterates $w_{\infty} = \lim_{t \rightarrow \infty} w_{(t)}$ is a global minimizer of \mathcal{L} , i.e., $\mathcal{L}(w_{\infty}) = 0$, then w_{∞} is given by*

$$w_{\infty} = \underset{w: \forall n, \langle w, x_n \rangle = y_n}{\operatorname{argmin}} D_{\psi}(w, w_{(0)}).$$

Theorem 1a. *Under the conditions in Theorem 1, if initialized at $w_{(0)} = \underset{w}{\operatorname{argmin}} \psi(w)$, then the mirror descent updates with dual momentum also converges to (6), i.e., for all $\{\eta_t\}_t, \{\beta_t\}_t, \{\gamma_t\}_t$, if $w_{(t)}$ from eq. (7) converges to $w_{\infty} \in \mathcal{G}$, then $w_{\infty} = \underset{w \in \mathcal{G}}{\operatorname{argmin}} \psi(w)$.*

Proof. (a) **Generic mirror descent:** Recall the updates of mirror descent: $\nabla \psi(w_{(t+1)}) - \nabla \psi(w_{(t)}) = -\eta_t \nabla \mathcal{L}(w_{(t)})$. Using a telescoping sum, we have,

$$\forall t, \nabla \psi(w_{(t)}) - \nabla \psi(w_{(0)}) = \sum_{t' < t} \nabla \psi(w_{(t'+1)}) - \nabla \psi(w_{(t')}) = \sum_{t' < t} -\eta_{t'} \nabla \mathcal{L}(w_{(t')}) \in \mathcal{P}_{\mathcal{X}}, \quad (23)$$

where the last inclusion follows as $\forall t', -\eta_{t'} \nabla \mathcal{L}(w_{(t')}) \in \mathcal{P}_{\mathcal{X}}$ from (20).

Thus, for all t , $w_{(t)}$ from mirror descent updates in eq. (5) always satisfies the stationarity condition of eq. (22). Additionally, if $w_{(t)}$ converges to a global minimum, then $w_{\infty} = \lim_{t \rightarrow \infty} w_{(t)} \in \mathcal{G} = \{w : \forall n, \langle w, x_n \rangle = y_n\}$ also satisfies the primal feasibility condition in eq. (22). Combining the above arguments, we have that if $\mathcal{L}(w_{\infty}) = 0$, then $w_{\infty} = \underset{w \in \mathcal{G}}{\operatorname{argmin}} D_{\psi}(w, w_{(0)})$.

(b) **Dual momentum:** For any $\tilde{\beta}_{t'}, \tilde{\gamma}_{t'} \in \mathbb{R}$ and $\tilde{w}_{(t')} \in \mathbb{R}^d$, consider a general update of the form

$$\nabla \psi(w_{(t+1)}) = \sum_{t' \leq t} \tilde{\beta}_{t'} \nabla \psi(w_{(t')}) + \tilde{\gamma}_{t'} \nabla \mathcal{L}(\tilde{w}_{(t')}). \quad (24)$$

Claim: If $\nabla \psi(w_{(0)}) = 0$, then for all updates of the form (24) satisfies $\nabla \psi(w_{(t)}) \in \mathcal{P}_{\mathcal{X}}$ —this can be easily proved by induction: (a) for $t = 0$, $\nabla \psi(w_{(0)}) = 0 \in \mathcal{P}_{\mathcal{X}}$; (b) let $\forall t' \leq t$, $\nabla \psi(w_{(t')}) \in \mathcal{P}_{\mathcal{X}}$, (c) then using the inductive assumption and eq. (20), we have $\nabla \psi(w_{(t+1)}) = \sum_{t' \leq t} \tilde{\beta}_{t'} \nabla \psi(w_{(t')}) + \tilde{\gamma}_{t'} \nabla \mathcal{L}(\tilde{w}_{(t')}) \in \mathcal{P}_{\mathcal{X}}$.

Dual momentum in eq. (7) is a special case of eq. (24) with appropriate choice of $\tilde{\beta}_{t'}, \tilde{\gamma}_{t'} \in \mathbb{R}$, and $\tilde{w}_{t'} \in \mathbb{R}^d$.

□

A.2. Proofs of propositions in Section 2

A.2.1. PRIMAL MOMENTUM AND NATURAL GRADIENT DESCENT

Recall the optimization problem in Examples 2–3: $\{(x_1 = [1, 2], y_1 = 1)\}$, and $\ell(u, y) = (u - y)^2$. We have $\mathcal{P}_{\mathcal{X}} = \text{span}(x_1) = \{z : 2z[1] - z[2] = 0\}$.

For entropy potential $\psi(w) = \sum_i w[i] \log w[i] - w[i]$, we have $\nabla\psi(w) = \log w$ (where the log is taken elementwise), and initialization $w_{(0)} = [1, 1]$ satisfies $\nabla\psi(w_{(0)}) = 0$.

1. Proof of Proposition 2a: we use primal momentum with $\beta_1 > 0$ only in the first step, and $\forall t \geq 2, \beta_t = \gamma_t = 0$.

First we note that, for $t > 2$, the updates follow the path of standard MD initialized at $\nabla\psi(w_{(2)})$ for a convex loss function. This implies the following:

- for appropriate choice of $\{\eta_t\}_{t \geq 2}$ (given by convergence analysis of mirror descent for convex functions), we can get $w_\infty = \lim_{t \rightarrow \infty} w_{(t)} \in \mathcal{G}$, and
- from eq. (23), w_∞ satisfies $\nabla\psi(w_\infty) - \nabla\psi(w_{(2)}) \in \mathcal{P}_{\mathcal{X}} \Rightarrow \nabla\psi(w_\infty) \in \nabla\psi(w_{(2)}) + \mathcal{P}_{\mathcal{X}}$.

Since w_∞ satisfies primal feasibility, from stationarity condition in eq. (22), we have

$$w_\infty = w_\psi^* = \underset{w \in \mathcal{G}}{\text{argmin}} \psi(w) \text{ if and only if } \nabla\psi(w_{(2)}) \in \mathcal{P}_{\mathcal{X}}.$$

We show that this is not the case for any $\beta_1 > 0$ and any $\gamma_1 \geq 0$.

Recall that $\Delta w_{(-1)} = 0$, $\nabla\psi(w_{(0)}) = 0$ and $\nabla\psi(w) = \log w$. Working through the steps in eq. (8), for scalars $r_0 = \eta_0(y_1 - \langle w_{(0)}, x_1 \rangle)$ and $\tilde{r}_1 = \eta_1(y_1 - \langle w_{(1)} + \gamma_1 \Delta w_{(0)}, x_1 \rangle)$, and any $\beta_1 > 0$, we have:

- $\nabla\psi(w_{(1)}) = r_0 x_1 \implies w_{(1)} = \exp(r_0 x_1)$, and
- $\nabla\psi(w_{(2)}) = \nabla\psi((1 + \beta_1)w_{(1)}) + \tilde{r}_1 x_1 = \log(1 + \beta_1) + r_0 x_1 + \tilde{r}_1 x_1 \in \log(1 + \beta_1) + \mathcal{P}_{\mathcal{X}} \notin \mathcal{P}_{\mathcal{X}}$. \square

2. Proof of Proposition 3a: The arguments are similar to the proof of Proposition 2a. In Example 3 we again use a finite $\eta_1 > 0$ to get $w_{(1)}$ and then follow the NGD using infinitesimal η initialized at $w_{(1)}$.

We know that for infinitesimal step size, the NGD path starting at $w_{(1)}$ follows the corresponding infinitesimal MD path on a convex problem and hence from eq. (23), the NGD updates for this example converges to a global minimum $w_\infty = \lim_{t \rightarrow \infty} w_{(t)} \in \mathcal{G}$, that satisfies $\nabla\psi(w_\infty) - \nabla\psi(w_{(1)}) \in \mathcal{P}_{\mathcal{X}} \Rightarrow \nabla\psi(w_\infty) \in \nabla\psi(w_{(1)}) + \mathcal{P}_{\mathcal{X}}$.

From stationarity condition in (22), $w_\infty = w_\psi^* = \underset{w \in \mathcal{G}}{\text{argmin}} \psi(w)$ if and only if $\nabla\psi(w_{(1)}) \in \mathcal{P}_{\mathcal{X}}$.

For natural gradient descent, $w_{(1)} = w_{(0)} - \eta_1 \nabla^2 \psi(w_{(0)})^{-1} \nabla \mathcal{L}(w_{(0)}) = [1 + \eta_1 r_0, 1 + 2\eta_1 r_0]$, where $r_0 = \eta_0(y_1 - \langle w_{(0)}, x_1 \rangle)$. We then have $\nabla\psi(w_{(1)}) \in \mathcal{P}_{\mathcal{X}} \Leftrightarrow 2\nabla\psi(w_{(1)})[1] - \nabla\psi(w_{(1)})[2] = 0 \Leftrightarrow 2\log(w_{(1)}[1]) - \log(w_{(1)}[2]) = 0 \Leftrightarrow \log\left(1 + \frac{\eta_1^2 r_0^2}{1 + 2\eta_1 r_0}\right) = 0$.

For any η_1 such that $\frac{\eta_1^2 r_0^2}{1 + 2\eta_1 r_0} \neq 0$, we get a contradiction. \square

B. Steepest descent for strictly monotone losses

We prove Theorem 5 in this section.

Theorem 5. For any separable dataset $\{x_n, y_n\}_{n=1}^N$ and any norm $\|\cdot\|$, consider the steepest descent updates from eq. (12) for minimizing $\mathcal{L}(w)$ in eq. (1) with the exponential loss $\ell(u, y) = \exp(-uy)$. For all initializations $w_{(0)}$, and all bounded step sizes satisfying $\eta_t \leq \max\{\eta_+, \frac{1}{B^2 \mathcal{L}(w_{(t)})}\}$, where $B := \max_n \|x_n\|_*$ and $\eta_+ < \infty$ is some maximum step size bound, the iterates $w_{(t)}$ satisfy the following,

$$\lim_{t \rightarrow \infty} \min_n \frac{y_n \langle w_{(t)}, x_n \rangle}{\|w_{(t)}\|} = \max_{w: \|w\| \leq 1} \min_n y_n \langle w, x_n \rangle.$$

In particular, if there is a unique maximum- $\|\cdot\|$ margin solution $w_{\|\cdot\|}^* = \underset{w}{\text{argmax}} \min_n \frac{y_n \langle w, x_n \rangle}{\|w\|}$, then the limit direction converges to it: $\bar{w}_\infty = \lim_{t \rightarrow \infty} \frac{w_{(t)}}{\|w_{(t)}\|} = w_{\|\cdot\|}^*$.

The following lemma is a standard result in convex analysis.

Lemma 9 (Fenchel Duality). *Let $A \in \mathbb{R}^{m \times n}$, and $f : \mathbb{R}^m \rightarrow \mathbb{R}$, $g : \mathbb{R}^n \rightarrow \mathbb{R}$ be two closed convex functions and f^* , g^* be their Fenchel conjugate functions, respectively. Then,*

$$\max_{w \in \mathbb{R}^n} -f^*(Aw) - g^*(-w) \leq \min_{r \in \mathbb{R}^n} f(r) + g(A^\top r). \quad (25)$$

Let $X \in \mathbb{R}^{N \times d}$ be the data matrix with x_n along the rows of X . Without loss of generality $y_n = 1$, as for linear models y_n can be absorbed into x_n . Let e_n denote the n^{th} standard basis in \mathbb{R}^N .

We define the $\|\cdot\|$ -maximum margin as,

$$\gamma = \max_w \min_{n \in [N]} \frac{\langle w, x_n \rangle}{\|w\|} = \max_{\|w\|=1} \min_{n \in [N]} e_n^\top X w. \quad (26)$$

Our primary technical novelty is the following duality lemma that generalizes similar result in Telgarsky (2013) for ℓ_1 norm to general norms: we want to show that $\|\nabla \mathcal{L}(w)\|_* \geq \gamma \mathcal{L}(w)$ for all w , where recall that $\|\cdot\|_*$ is the dual norm of $\|\cdot\|$.

Define $r_n(w) = \exp(-w^\top x_n)$ and let $r(w) = [r_n(w)]_{n=1}^N \in \mathbb{R}^N$. For succinctness, we often refer $r(w)$ without the dependence on w as r . Note that $\mathcal{L}(w) = \|r\|_1$ and $\nabla \mathcal{L}(w) = X^\top r$.

We can now restate $\|\nabla \mathcal{L}(w)\|_* \geq \gamma \mathcal{L}(w)$ as $\frac{\|X^\top r\|_*}{\|r\|_1} \geq \gamma$. In the following lemma, we show this holds for any $r_n \geq 0$, and since norms are homogeneous, this is equivalent to $\min_{r \in \Delta_{N-1}} \|X^\top r\|_* \geq \gamma$, where $\Delta_{N-1} = \{v \in \mathbb{R}^N : v \geq 0, \|v\|_1 = 1\}$ is the N -dimensional probability simplex.

Lemma 10. *For any norm $\|\cdot\|$, the following duality holds:*

$$\min_{r \in \Delta_{N-1}} \|X^\top r\|_* \geq \max_{\|w\|=1} \min_{n \in [N]} e_n^\top X w = \gamma. \quad (27)$$

This implies, for exponential loss $\ell(u, y) = \exp(-uy)$, the following holds

$$\forall w, \|\nabla \mathcal{L}(w)\|_* \geq \gamma \mathcal{L}(w). \quad (28)$$

Proof. Let $\mathbf{1}_E$ denote the indicator function which takes value 0 if E is satisfied and ∞ otherwise.

Define $f(r) = \mathbf{1}_{r \in \Delta_{N-1}}$ and $g(z) = \|z\|_*$, so that

$$\min_{r \in \Delta_{N-1}} \|X^\top r\|_* = \min_{r \in \mathbb{R}^N} f(r) + g(X^\top r). \quad (29)$$

The conjugates are $f^*(w) = \max_{z \in \Delta_{N-1}} \langle w, z \rangle = \max_n w_n$, and $g^*(w) = \mathbf{1}_{\|w\| \leq 1}$. The LHS of Lemma 9 is

$$\begin{aligned} \max_w (-f^*(Xw) - g^*(-w)) &= \max_w (-\max_n e_n^\top Xw - \mathbf{1}_{\|w\| \leq 1}) \\ &= \max_{\|w\| \leq 1} \min_n e_n^\top X(-w) \stackrel{(a)}{=} \max_{\|w\| \leq 1} \min_n e_n^\top Xw \stackrel{(b)}{=} \gamma, \end{aligned} \quad (30)$$

where (a) follows from central symmetry of $\{w : \|w\| \leq 1\}$, and (b) from definition of maximum $\|\cdot\|$ -margin in eq. (26).

Using weak duality (Lemma 9) on eqs. (29) and (30), we have $\forall r, \|X^\top r\|_* \geq \gamma \|r\|_1$. Finally, recalling that for exponential loss $r_n(w) = \exp(-w^\top x_n)$, $\mathcal{L}(w) = \|r(w)\|_1$ and $\nabla \mathcal{L}(w) = X^\top r(w)$, we have $\forall w, \|\nabla \mathcal{L}(w)\|_* \geq \gamma \mathcal{L}(w)$. \square

Recall the steepest descent updates in eqs. (11) and (12) :

$$\begin{aligned} w_{(t+1)} &= w_{(t)} + \eta_t \Delta w_{(t)}, \text{ where } \Delta w_{(t)} \text{ satisfies} \\ \langle \Delta w_{(t)}, -\nabla \mathcal{L}(w_{(t)}) \rangle &= \|\Delta w_{(t)}\|^2 = \|\nabla \mathcal{L}(w_{(t)})\|_*^2. \end{aligned} \quad (31)$$

Lemma 11. *For exponential loss $\ell(u; y) = \exp(-uy)$, consider the steepest descent iterates $w_{(t)}$ for minimizing $\mathcal{L}(w_{(t)})$, with any initialization $w_{(0)}$ and any finite step size η_t that leads to a strictly decreasing sequence $\mathcal{L}(w_{(t)})$ and satisfies $0 < \eta_t \leq \max\{\eta_+, \frac{1}{B^2 \mathcal{L}(w_{(t)})}\}$, where $B = \max_n \|x_n\|_*$. Then the following holds:*

- (A) $\sum_{t=0}^{\infty} \eta_t \|\nabla \mathcal{L}(w_{(t)})\|_{\star}^2 \leq \infty$, and hence $\|\nabla \mathcal{L}(w_{(t)})\|_{\star} \rightarrow 0$.
 (B) Iterates $w_{(t)}$ converge to a global minima $\mathcal{L}(w_{(t)}) \rightarrow 0$, and hence $\forall n \langle w_{(t)}, x_n \rangle \rightarrow \infty$.
 (C) $\sum_{t=0}^{\infty} \eta_t \|\nabla \mathcal{L}(w_{(t)})\|_{\star} = \infty$.

Proof. 1. **Proof of (A):** We have that $\|x_n\|_{\star} \leq B$ for all n . Recall that $r_n(w) = \exp(-\langle w, x_n \rangle)$, $\mathcal{L}(w) = \sum_n r_n(w)$, and $\nabla \mathcal{L}(w) = \sum_n r_n(w) x_n$. Thus, for all v , we have

$$v^{\top} \nabla^2 \mathcal{L}(w) v = \sum_n r_n(w) (x_n^{\top} v)^2 \leq \sum_n r_n(w) \|x_n\|_{\star}^2 \|v\|^2 \leq \mathcal{L}(w) B^2 \|v\|^2. \quad (32)$$

Using Taylor's reminder theorem for the convex loss \mathcal{L} , we have

$$\begin{aligned} \mathcal{L}(w_{(t+1)}) &\leq \mathcal{L}(w_{(t)}) + \eta_t \langle \nabla \mathcal{L}(w_{(t)}), \Delta w_{(t)} \rangle + \sup_{\beta \in (0,1)} \frac{\eta_t^2}{2} \Delta w_{(t)}^{\top} \nabla^2 \mathcal{L}(w_{(t)} + \beta \eta_t \Delta w_{(t)}) \Delta w_{(t)} \\ &\stackrel{(a)}{\leq} \mathcal{L}(w_{(t)}) - \eta_t \|\nabla \mathcal{L}(w_{(t)})\|_{\star}^2 + \frac{\eta_t^2 B^2}{2} \sup_{\beta \in (0,1)} \mathcal{L}(w_{(t)} + \beta \eta_t \Delta w_{(t)}) \|\Delta w_{(t)}\|^2 \\ &\stackrel{(b)}{\leq} \mathcal{L}(w_{(t)}) - \eta_t \|\nabla \mathcal{L}(w_{(t)})\|_{\star}^2 + \frac{\eta_t^2 B^2}{2} \mathcal{L}(w_{(t)}) \|\Delta w_{(t)}\|^2 \\ &\stackrel{(c)}{\leq} \mathcal{L}(w_{(t)}) - \frac{\eta_t}{2} \|\nabla \mathcal{L}(w_{(t)})\|_{\star}^2, \end{aligned} \quad (33)$$

where (a) follows from eq. (32) and from the condition on update direction in eq. (31); (b) follows as $\eta_t \Delta w_{(t)}$ is a descent step and along with convexity of $\mathcal{L}(w)$ we have $\sup_{\beta \in (0,1)} \mathcal{L}(w_{(t)} + \beta \eta_t \Delta w_{(t)}) \leq \mathcal{L}(w_{(t)})$; and (c) follows as $\eta_t \leq \frac{1}{B^2 \mathcal{L}(w_{(t)})}$ from the assumption and also using $\|\Delta w_{(t)}\| = \|\nabla \mathcal{L}(w_{(t)})\|_{\star}$ from eq. 31.

Thus, $\mathcal{L}(w_{(t)}) - \mathcal{L}(w_{(t+1)}) \geq \frac{\eta_t}{2} \|\nabla \mathcal{L}(w_{(t)})\|_{\star}^2$, which implies

$$\forall t, \sum_{u=0}^t \eta_u \|\nabla \mathcal{L}(w_{(u)})\|_{\star}^2 \leq 2 \sum_{u=0}^t (\mathcal{L}(w_{(u)}) - \mathcal{L}(w_{(u+1)})) = 2(\mathcal{L}(w_{(0)}) - \mathcal{L}(w_{(t+1)})) < \infty. \quad (34)$$

where the final inequality follows as $\mathcal{L}(w_{(0)}) < \infty$ and $\mathcal{L}(w_{(t)}) \geq 0 \forall t$.

In the continuous time limit of $\eta \rightarrow 0$, (A) is equivalently expressed as $\int_0^t \|\nabla \mathcal{L}(w_{(t)})\|_{\star}^2 < \infty$. Thus, we have $\lim_{t \rightarrow \infty} \|\nabla \mathcal{L}(w_{(t)})\|_{\star} = 0$ —both for any finite $\eta_t > 0$ as well as in the continuous time limit of $\eta \rightarrow 0$.

2. **Proof of (B) and (C) :** Consider any $v \in \mathbb{R}^d$ that linearly separates the data, i.e., $\forall n, \langle v, x_n \rangle > 0$ (such a v always exists for any linearly separable data), then $\forall t < \infty$, $v^{\top} \nabla \mathcal{L}(w_{(t)}) = \sum_{n \in [N]} \exp(-\langle w_{(t)}, x_n \rangle) x_n^{\top} v > 0$.

Since $\lim_{t \rightarrow \infty} v^{\top} \nabla \mathcal{L}(w_{(t)}) = 0$, it must be that $\forall n$, $\lim_{t \rightarrow \infty} \exp(-\langle w_{(t)}, x_n \rangle) = 0$, and $\lim_{t \rightarrow \infty} \|w_{(t)}\| = \infty$.

Finally, using triangle inequality, we have

$$\infty = \lim_{t \rightarrow \infty} \|w_{(t)}\| \leq \|w_{(0)}\| + \eta \sum_{t=0}^{\infty} \|\Delta w_{(t)}\| = \|w_{(0)}\| + \sum_{t=0}^{\infty} \eta_t \|\nabla \mathcal{L}(w_{(t)})\|_{\star}, \quad (35)$$

where we used $\|\Delta w_{(t)}\| = \|\nabla \mathcal{L}(w_{(t)})\|_{\star}$ from (31). This gives us $\sum_{t=0}^{\infty} \eta_t \|\nabla \mathcal{L}(w_{(t)})\|_{\star} = \infty$ in (C). \square

We next show that under the conditions of Theorem 5, $\mathcal{L}(w_{(t)})$ forms a decreasing sequence, and hence satisfies the assumption in Lemma 11.

Lemma 12. *If step sizes η_t satisfy $\eta_t = \frac{c_t}{B^2 \mathcal{L}(w_{(t)})}$ for $c_t \leq \sqrt{2}$, then $\mathcal{L}(w_{(t+1)}) \leq \mathcal{L}(w_{(t)})$.*

Proof. From the Taylor expansion of $\mathcal{L}(w)$ in eq. (33), we have

$$\begin{aligned} \mathcal{L}(w_{(t+1)}) &\leq \mathcal{L}(w_{(t)}) - \eta_t \|\nabla \mathcal{L}(w_{(t)})\|_{\star}^2 + \frac{\eta_t^2 B^2}{2} \|\nabla \mathcal{L}(w_{(t)})\|_{\star}^2 \sup_{\beta \in (0,1)} \mathcal{L}(w_{(t)} + \beta \eta_t \Delta w_{(t)}) \\ &\stackrel{(a)}{\leq} \mathcal{L}(w_{(t)}) - \eta_t \|\nabla \mathcal{L}(w_{(t)})\|_{\star}^2 \left(1 - \frac{\eta_t B^2}{2} \max(\mathcal{L}(w_{(t)}), \mathcal{L}(w_{(t+1)})) \right) \end{aligned} \quad (36)$$

We want to show that $\mathcal{L}(w_{(t+1)}) \leq \mathcal{L}(w_{(t)})$. Lets assume contrary that $\mathcal{L}(w_{(t+1)}) \geq \mathcal{L}(w_{(t)})$.

Denote $\gamma_t = \|\nabla \mathcal{L}(w_{(t)})\|_*$. From eq. (36), we have

$$\mathcal{L}(w_{(t)}) \stackrel{(a)}{\leq} \mathcal{L}(w_{(t+1)}) \leq \mathcal{L}(w_{(t)}) - \eta_t \gamma_t^2 \left(1 - \frac{\eta_t B^2 \mathcal{L}(w_{(t+1)})}{2}\right) \stackrel{(b)}{\implies} \left(1 - \frac{\eta_t B^2 \mathcal{L}(w_{(t+1)})}{2}\right) \leq 0, \quad (37)$$

where (a) follows from contrary assumption, and (b) follows as $\eta_t \gamma_t^2 \geq 0$

Further, let $\eta_t = \frac{c_t}{B^2 \mathcal{L}(w_t)}$ for some $0 < c_t \leq \sqrt{2}$. Following up from eq. (37), we have

$$\begin{aligned} \mathcal{L}(w_{(t+1)}) &\leq \mathcal{L}(w_{(t)}) + \eta_t \gamma_t^2 \left(\frac{\eta_t B^2 \mathcal{L}(w_{(t+1)})}{2} - 1\right) \\ &\stackrel{(a)}{\leq} \mathcal{L}(w_{(t)}) - c_t \mathcal{L}(w_t) + \frac{c_t^2 \mathcal{L}(w_{(t+1)})}{2} \\ \implies \mathcal{L}(w_{(t+1)}) &\leq \frac{1 - c_t}{1 - 0.5c_t^2} \mathcal{L}(w_{(t)}) \stackrel{(b)}{\leq} \mathcal{L}(w_{(t)}). \end{aligned}$$

where in (a) we used $\gamma_t = \|\nabla \mathcal{L}(w_t)\|_* = \|\sum_n \exp(w_t^T x_n) x_n\|_* \leq B \mathcal{L}(w_t)$ from triangle inequality as well as $\left(\frac{\eta_t B^2 \mathcal{L}(w_{(t+1)})}{2} - 1\right) \geq 0$ from eq. (37), and (b) follows as for $0 < c_t \leq \sqrt{2}$, $\frac{1-c_t}{1-0.5c_t^2} \leq 1$. This shows $\mathcal{L}(w_{(t+1)}) \leq \mathcal{L}(w_t)$ which is a contradiction. \square

B.1. Remaining steps in the proof of Theorem 5

The steepest descent updates in eq. (31) can be equivalently written as:

$$\begin{aligned} w_{(t+1)} &= w_{(t)} - \eta_t \gamma_t p_{(t)}, \text{ where} \\ \gamma_t &\triangleq \|\nabla \mathcal{L}(w_{(t)})\|_*, \text{ and } p_{(t)} = \frac{\Delta w_{(t)}}{\|\nabla \mathcal{L}(w_{(t)})\|_*}, \text{ which satisfies} \\ \langle p_{(t)}, \nabla \mathcal{L}(w_{(t)}) \rangle &= \|\nabla \mathcal{L}(w_{(t)})\|_*, \|p_{(t)}\| = 1. \end{aligned} \quad (38)$$

From eq. 33, using $\gamma_t = \|\nabla \mathcal{L}(w_{(t)})\|_* = \|\Delta w_{(t)}\|$, we have that

$$\begin{aligned} \mathcal{L}(w_{(t+1)}) &\leq \mathcal{L}(w_{(t)}) - \eta_t \gamma_t^2 + \frac{\eta_t^2 B^2 \mathcal{L}(w_{(t)}) \gamma_t^2}{2} = \mathcal{L}(w_{(t)}) \left[1 - \frac{\eta_t \gamma_t^2}{\mathcal{L}(w_{(t)})} + \frac{\eta_t^2 B^2 \gamma_t^2}{2}\right] \\ &\stackrel{(a)}{\leq} \mathcal{L}(w_{(t)}) \exp\left(-\frac{\eta_t \gamma_t^2}{\mathcal{L}(w_{(t)})} + \frac{\eta_t^2 B^2 \gamma_t^2}{2}\right) \\ &\stackrel{(b)}{\leq} \mathcal{L}(w_{(0)}) \exp\left(-\sum_{u \leq t} \frac{\eta_u \gamma_u^2}{\mathcal{L}(w_{(u)})} + \sum_{u \leq t} \frac{\eta_u^2 B^2 \gamma_u^2}{2}\right), \end{aligned} \quad (39)$$

where we get (a) by using $(1+x) \leq \exp(x)$, and (b) using recursion.

Step 1: Lower bound the unnormalized margin: From eq. (39), we have,

$$\max_{n \in [N]} \exp(-\langle w_{(t+1)}, x_n \rangle) \leq \mathcal{L}(w_{(t+1)}) \leq \mathcal{L}(w_{(0)}) \exp\left(-\sum_{u \leq t} \frac{\eta_u \gamma_u^2}{\mathcal{L}(w_{(u)})} + \sum_{u \leq t} \frac{\eta_u^2 B^2 \gamma_u^2}{2}\right). \quad (40)$$

By applying $-\log$,

$$\min_{n \in [N]} \langle w_{(t+1)}, x_n \rangle \geq \sum_{u \leq t} \frac{\eta_u \gamma_u^2}{\mathcal{L}(w_{(u)})} - \sum_{u \leq t} \frac{\eta_u^2 B^2 \gamma_u^2}{2} - \log \mathcal{L}(w_{(0)}). \quad (41)$$

Step 2: Upper bound $\|w_{(t+1)}\|$: Using $\|\Delta w_{(u)}\| = \|\nabla \mathcal{L}(w_{(u)})\|_* = \gamma_u$, we have,

$$\|w_{(t+1)}\| \leq \|w_{(0)}\| + \sum_{u \leq t} \eta_u \|\Delta w_{(u)}\| \leq \|w_{(0)}\| + \sum_{u \leq t} \eta_u \gamma_u. \quad (42)$$

Combining eqs. (41) and (42), $\forall n \in [N]$, we have that

$$\frac{\langle w_{(t+1)}, x_n \rangle}{\|w_{(t+1)}\|} \geq \frac{\sum_{u \leq t} \frac{\eta_u \gamma_u^2}{\mathcal{L}(w_{(u)})}}{\sum_{u \leq t} \eta_u \gamma_u + \|w_{(0)}\|} - \frac{\sum_{u \leq t} \frac{\eta_u^2 B^2 \gamma_u^2}{2}}{\|w_{(t+1)}\|} - \frac{\log \mathcal{L}(w_{(0)})}{\|w_{(t+1)}\|}. \quad (43)$$

$$:= (I) + (II) + (III). \quad (44)$$

We look at the three terms separately,

(I) From the duality Lemma 10, we have $\gamma_u = \|\nabla \mathcal{L}(w_{(u)})\|_* \geq \gamma \mathcal{L}(w_{(u)})$. Hence, $\sum_{u \leq t} \frac{\eta_u \gamma_u^2}{\mathcal{L}(w_{(u)})} \geq \gamma \sum_{u \leq t} \eta_u \gamma_u$ and further using $\sum_{u \leq t} \eta_u \gamma_u \rightarrow \infty$ from Lemma 11, we have

$$\frac{\sum_{u \leq t} \frac{\eta_u \gamma_u^2}{\mathcal{L}(w_{(u)})}}{\sum_{u \leq t} \eta_u \gamma_u + \|w_{(0)}\|} \geq \gamma \frac{\sum_{u \leq t} \eta_u \gamma_u}{\sum_{u \leq t} \eta_u \gamma_u + \|w_{(0)}\|} \rightarrow \gamma$$

(II) For any bounded $\eta \leq \eta_+$, $\sum_{u \leq t} \frac{\eta_u^2 B^2 \gamma_u^2}{2} \leq \frac{\eta_+ B^2}{2} \sum_{u \leq t} \eta_u \gamma_u^2 < \infty$ (from Lemma 11).

Along with using $\|w_{(t)}\| \rightarrow \infty$ from Lemma 11, we get $\frac{\sum_{u \leq t} \frac{\eta_u^2 B^2 \gamma_u^2}{2}}{\|w_{(t+1)}\|} \rightarrow 0$.

(III) $\frac{\log \mathcal{L}(w_{(0)})}{\|w_{(t+1)}\|} \rightarrow 0$ since $\|w_{(t)}\| \rightarrow \infty$.

Combining the above in (44), we get $\lim_{t \rightarrow \infty} \frac{w_{(t+1)}^\top x_n}{\|w_{(t+1)}\|} \geq \gamma := \max_w \frac{w^\top x_n}{\|w\|}$ \square

C. Adagrad

Lemma 13. Let $\mathcal{L}(w) = \sum_{n=1}^N \exp(-w^\top x_n)$, $\|\cdot\|_t$ be some $w_{(t)}$ -dependent norm, and $\|\cdot\|_{t,*}$ be its dual, and assume that and $\forall t : \|x_n\|_{t,*} \leq 1$. We examine the steepest descent (SD) Sequence:

$$w_{(t+1)} = w_{(t)} - \eta \beta_t p_{(t)} \quad (45)$$

where

$$p_{(t)}^\top \nabla \mathcal{L}(w_{(t)}) = \|\nabla \mathcal{L}(w_{(t)})\|_{t,*} \triangleq \beta_t ; \|p_{(t)}\|_t = 1.$$

Then, for any $w_{(0)}$ such that $\frac{\eta}{2} \mathcal{L}(w_{(0)}) < 1$, we have that $\sum_{u=0}^{\infty} \beta_u^2 < \infty$ and therefore $\lim_{t \rightarrow \infty} \|\beta_t\| = 0$.

Lemma 14. Let $\mathcal{L}(w) = \sum_{n=1}^N \exp(-w^\top x_n)$. We examine the AdaGrad

$$w_{(t+1)} = w_{(t)} - \eta \mathbf{G}_{(t)}^{-1/2} \nabla \mathcal{L}(w_{(t)}) \quad (46)$$

where $\mathbf{G}_{(t)}$ is a diagonal matrix such that

$$\forall i : \mathbf{G}_{(t),ii} = \sum_{u=0}^t (\nabla \mathcal{L}(w_{(u)}))_i^2.$$

Then, for any $w_{(0)}$ such that $\frac{\eta}{2} \mathcal{L}(w_{(0)}) < 1$, and if $\left\| \mathbf{G}_{(0)}^{-1/4} x_n \right\|_2 \leq 1$, $\exists C < \infty$ such that

$$\forall i, \forall t : G_{(t),ii} < C.$$

Proof. First, we note that AdaGrad is a special case of the steepest descent algorithm as in Lemma 13 with respect to the norm $\|v\|_t = \|\mathbf{G}_{(t)}^{1/2}v\|_2$. Here the dual norm $\|v\|_{t,*} = \|\mathbf{G}_{(t)}^{-1/2}v\|_2$. Since $\mathbf{G}_{(t),ii}^{-1}$ is monotonically decreasing for all t , $\|\mathbf{G}_{(t)}^{-1/2}x_n\|_2 \leq \|\mathbf{G}_{(0)}^{-1/2}x_n\|_2 \leq 1$, and so we can apply Lemma 13. This implies that

$$\begin{aligned} \infty &> \sum_{t=0}^{\infty} \left\| \mathbf{G}_{(t)}^{-1/2} \nabla \mathcal{L}(w_{(t)}) \right\|_2^2 \\ &= \sum_{i=1}^d \sum_{t=0}^{\infty} (\nabla \mathcal{L}(w_{(t)}))_i^2 \left[\sum_{u=0}^t (\nabla \mathcal{L}(w_{(u)}))_i^2 \right]^{-1/2} \\ &\geq \sum_{i=1}^d \sum_{t=0}^{\infty} (\nabla \mathcal{L}(w_{(t)}))_i^2 \left[\sum_{u=0}^{\infty} (\nabla \mathcal{L}(w_{(u)}))_i^2 \right]^{-1/2} \\ &= \sum_{i=1}^d \sqrt{\sum_{t=0}^{\infty} (\nabla \mathcal{L}(w_{(t)}))_i^2} \end{aligned}$$

This implies that

$$\forall i : \sum_{t=0}^{\infty} (\nabla \mathcal{L}(w_{(t)}))_i^2 < \infty,$$

□

D. Gradient descent on factorized parameterization

We first prove the Lemma 8 and Lemma D.2 that hold for any general linear model (1) with exponential loss $\ell(u, y) = \exp(-uy)$. These results are not specific to matrix factorization setup in Section 4.

D.1. Convergence of $-\nabla \mathcal{L}(w_{(t)})$ in direction

Recall Lemma 8:

Lemma 8. *For almost all linearly separable datasets $\{x_n, y_n\}_n$, consider any sequence $w_{(t)}$ that minimizes $\mathcal{L}(w)$ in eq. (1) with exponential loss, i.e., $\mathcal{L}(w_{(t)}) \rightarrow 0$. If $\frac{w_{(t)}}{\|w_{(t)}\|}$ converges, then for every accumulation point z_{∞} of $\left\{ \frac{-\nabla \mathcal{L}(w_{(t)})}{\|\nabla \mathcal{L}(w_{(t)})\|} \right\}_t$, $\exists \{\alpha_n \geq 0\}_{n \in S}$ s.t., $z_{\infty} = \sum_{n \in S} \alpha_n y_n x_n$, where $\bar{w}_{\infty} = \lim_{t \rightarrow \infty} \frac{w_{(t)}}{\|w_{(t)}\|}$ and $S = \{n : y_n \langle \bar{w}_{\infty}, x_n \rangle = \min_n y_n \langle \bar{w}_{\infty}, x_n \rangle\}$ are the indices of the data points with smallest margin to \bar{w}_{∞} .*

Here for almost all $\{x_n\}$ means that with probability 1 over $\{x_n\}$ are drawn from a distribution that is absolutely continuous w.r.t the d dimensional Lebesgue measure.

Proof. Without loss of generality assume $\forall n, y_n = 1$ —as the sign of y can be absorbed into x , i.e., $x_n \leftarrow y_n x_n$. Let $X \in \mathbb{R}^{N \times d}$ denote the data matrix with $x_n \in \mathbb{R}^d$ along the rows of X . Also, for any $J \subseteq [N]$, $X_J \in \mathbb{R}^{|J| \times d}$ denotes the submatrix of X with only the rows corresponding to indices in J .

We have that $\lim_{t \rightarrow \infty} \mathcal{L}(w_{(t)}) = 0$ for strictly monotone loss over separable data, this implies asymptotically $w_{(t)}$ is in $\mathcal{G} = \{w : Xw > 0, \|w\| \rightarrow \infty\}$. Also, since $w_{(t)}$ converges in direction to \bar{w}_{∞} , we can write $w_{(t)} = g(t)\bar{w}_{\infty} + \rho(t)$ for a scalar $g(t) = \|w_{(t)}\| \rightarrow \infty$ and $\frac{\rho(t)}{g(t)}$.

From these conditions, we also have $\forall n, X\bar{w}_{\infty} > 0$. We introduce some additional notation:

- Let $\gamma = \min_n \langle x_n, \bar{w}_{\infty} \rangle = \min_n e_n^{\top} X\bar{w}_{\infty} > 0$ be the margin of \bar{w}_{∞} , where e_n are standard basis in \mathbb{R}^N .
- Denote by $S := \{n : \langle x_n, \bar{w}_{\infty} \rangle = \gamma\}$ the indices of support vectors of \bar{w}_{∞} .
- Denote the second smallest margin of \bar{w}_{∞} as $\bar{\gamma} := \min_{n \in S^c} \langle x_n, \bar{w}_{\infty} \rangle > \gamma$.
- Denote the margin for datapoint n as $\bar{\gamma}_n := \langle x_n, \bar{w}_{\infty} \rangle \geq \gamma$.

- Define $\alpha_n(t) := \exp(-\langle \rho(t), x_n \rangle)$, $\alpha \in \mathbb{R}^N$ be a vector of α_n stacked. For any $J \subset [N]$, similar to X_J , let $\alpha_J \in \mathbb{R}^{|J|}$ be a sub-vector with components corresponding to the indices in J
- $B = \max_n \|x_n\|_2$,

Since $\rho(t)/g(t) \rightarrow 0$, $\forall \epsilon_1, \epsilon_2, \exists t_{\epsilon_1}, t_{\epsilon_2}$ such that

$$\begin{aligned} \forall t > t_{\epsilon_1}, \quad \min_n -\langle \rho(t), x_n \rangle &\geq -\|\rho(t)\|_2 B \geq -\epsilon_1 \gamma g(t), \\ \forall t > t_{\epsilon_2}, \quad \max_n \langle \rho(t), x_n \rangle &\leq \|\rho(t)\|_2 B \leq \epsilon_2 \bar{\gamma} g(t) \end{aligned} \quad (47)$$

The following claim is useful:

Claim 1. For almost $\{x_n\}$ we have, $|S| < d$ and $\sigma_{|S|}(X_S) > 0$, where $\sigma_k(A)$ is the k^{th} singular value of A .

Proof. Since, $S = \{n : \langle \bar{w}_\infty, x_n \rangle = \gamma\}$, we have $X_S \bar{w}_\infty = \gamma 1_S \in \mathbb{R}^{|S|}$. For any fixed subset J if $|J| > d$, then with probability 1,

$$\mathbb{R}^{|J|} \ni 1_J \notin \text{colspan}(X_J), \text{ for almost all } X_J \in \mathbb{R}^{|J| \times d} \quad (48)$$

This is because if X is random and from a continuous distribution, and the rank deficient column span of X_J will miss any fixed vector v (that is independent of X) with probability 1. Since $1_S \in \text{span}(X_S)$, then for almost all X , $|S| \leq d$ and $\sigma_{|S|}(X_S) > 0$. \square

Exponential loss: For exponential loss, the gradients are given by

$$\begin{aligned} -\nabla \mathcal{L}(w(t)) &= \sum_{n \in S} \exp(-\gamma g(t)) \exp(-\rho(t)^\top x_n) x_n + \sum_{n \in S^c} \exp(-\bar{\gamma}_n g(t)) \exp(-\rho(t)^\top x_n) x_n \\ &= I(t) + II(t), \end{aligned} \quad (49)$$

where $I(t) = \sum_{n \in S} \exp(-\gamma g(t)) \exp(-\rho(t)^\top x_n) x_n$ and $II(t) = \sum_{n \in S^c} \exp(-\bar{\gamma}_n g(t)) \exp(-\rho(t)^\top x_n) x_n$.

We will show that $\lim_{t \rightarrow \infty} \frac{\|II(t)\|}{\|I(t)\|} = 0$.

Recall that $\alpha(t) = [\alpha_n(t)]_n$ is defined as $\alpha_n(t) = \exp(-\langle \rho(t), x_n \rangle)$ and $\alpha_S(t) \in \mathbb{R}^{|S|}$ is a subvector restricted to indices in S . The following are true for any $\epsilon_1, \epsilon_2 > 0$.

Step 1: Lower bound on $I(t)$: For large enough $t > t_{\epsilon_1}$, we have

$$\begin{aligned} \|I\|_2 &= \exp(-\gamma g(t)) \|X_S \alpha_S(t)\|_2 \geq \exp(-\gamma g(t)) \sigma_{|S|}(X_S) \|\alpha_S(t)\|_2 \geq \exp(-\gamma g(t)) \sigma_{|S|}(X_S) \min_{n \in S} \alpha_n(t) \\ &\stackrel{(a)}{\geq} \sigma_{|S|}(X_S) \exp(-(1 + \epsilon_1)\gamma g(t)) := C_1 \exp(-(1 + \epsilon_1)\gamma g(t)), \end{aligned} \quad (50)$$

where (a) follows from the definition of $\alpha_n = \exp(-\langle \rho(t), x_n \rangle)$ and (47), and $C_1 > 0$ is a constant independent of t .

Step 2: Upper bound on $II(t)$: Again, for large enough $t > t_{\epsilon_2}$, we have

$$\begin{aligned} \|II(t)\|_2 &= \sum_{n \in S^c} \exp(-\bar{\gamma}_n g(t)) \exp(-\rho(t)^\top x_n) x_n \leq N \max_n \exp(-\bar{\gamma}_n g(t)) \alpha_n \|x_n\|_2 \stackrel{(a)}{\leq} \exp(-\bar{\gamma} g(t)) B N \max_n \alpha_n \\ &\stackrel{(a)}{\leq} B N \exp(-(1 - \epsilon_2)\bar{\gamma} g(t)) := C_2 \exp(-(1 - \epsilon_2)\bar{\gamma} g(t)), \end{aligned} \quad (51)$$

where (a) uses $\forall n \notin S, \bar{\gamma}_n \geq \bar{\gamma}$ and (b) follows from the definition of $\alpha_n = \exp(-\langle \rho(t), x_n \rangle)$ and (47), and $C_2 > 0$ is a constant independent of t .

Remaining steps in the proof: By combining (50) and (51) using $\epsilon_1 = (\bar{\gamma} - \gamma)/4\gamma$ and $\epsilon_2 = (\bar{\gamma} - \gamma)/4\bar{\gamma}$ and an appropriate constant $C > 0$, we have for any norm $\|\cdot\|$

$$\frac{\|II(t)\|}{\|I(t)\|} \leq C \exp\left(-\frac{1}{2}(\bar{\gamma} - \gamma)g(t)\right) \stackrel{(a)}{\rightarrow} 0, \quad (52)$$

where (a) follows from $\bar{\gamma} > \gamma$ and $g(t) = \|w_{(t)}\| \rightarrow \infty$

$$\text{Finally, } -\frac{\nabla \mathcal{L}(w_{(t)})}{\|\nabla \mathcal{L}(w_{(t)})\|} = \frac{I(t)}{\|I(t)+II(t)\|} + \frac{II(t)}{\|I(t)+II(t)\|}.$$

Since $\left\| \frac{II(t)}{\|I(t)+II(t)\|} \right\| \leq \frac{\|II(t)\|/\|I(t)\|}{1-\|II(t)\|/\|I(t)\|} \xrightarrow{t \rightarrow \infty} 0$, and $I(t) \propto \sum_{n \in S} \alpha_n(t) x_n$ for $\alpha_n(t) > 0$, then every limit point of $-\frac{\nabla \mathcal{L}(w_{(t)})}{\|\nabla \mathcal{L}(w_{(t)})\|} \rightarrow \sum_{n \in S} \alpha_n x_n$ for some $\alpha_n > 0$. Recall that in the beginning of the proof we made a change of variable that $x_n \leftarrow y_n x_n$. Reversing this change of variable finishes the proof for exponential loss. \square

D.2. Convergence of $w_{(t)}$ in direction

Lemma 15 ($w_{(t)}$ converges in direction if $\Delta w_{(t)}$ converges in direction). *Assume $\|w_{(t)}\| \rightarrow \infty$ and $\|\Delta w_{(t)}\| > 0, \forall t < \infty$. If $\frac{\Delta w_{(t)}}{\|\Delta w_{(t)}\|} \rightarrow \bar{w}_\infty$, then $\frac{w_{(t)}}{\|w_{(t)}\|} \rightarrow \bar{w}_\infty$ under (a) any bounded discrete step-size update of $w_{(t+1)} = w_{(t)} + \eta_t \Delta w_{(t)}$ for $0 < \eta_t < \infty$, or (b) the continuous time dynamics of $\eta_t \rightarrow 0$ with $\frac{dw_{(t)}}{dt} = \Delta w_{(t)}$.*

Proof. Discrete Update: Let $0 < \eta_- \leq \eta_t \leq \eta_+ < \infty$. Since $\frac{\Delta w_{(t)}}{\|\Delta w_{(t)}\|} \rightarrow \bar{w}_\infty$, we can write $\Delta w_{(t)} = \bar{w}_\infty h(t) + \xi(t)$ where $h(t) = \|\Delta w_{(t)}\|$ and $\frac{\xi(t)}{h(t)} \rightarrow 0$.

Define

$$g(t) := \sum_{u < t} \eta_u h(u), \text{ and } \rho(t) := w_{(t)} - \bar{w}_\infty g(t) = \sum_{u < t} \eta_u \xi(u) + w_{(0)}. \quad (53)$$

In order to prove the lemma, we need to show that $g(t) \rightarrow \infty$ and $\frac{\rho(t)}{g(t)} \rightarrow 0$. We have the following,

1. As $\|w_{(t)}\| \rightarrow \infty$, we have $\|w_{(t)}\| \leq \|w_{(0)}\| + \sum_{t' < t} \eta_{t'} \|\Delta w_{(t')}\| = \|w_{(0)}\| + g(t) \xrightarrow{t \in \infty} \infty$.
2. Also, $g(t) = \sum_{t' < t} h(t')$ is a strictly monotonically increasing as $\forall t < \infty, h(t) > 0$.
3. Finally,

$$\lim_{t \rightarrow \infty} \frac{\rho(t+1) - \rho(t)}{g(t+1) - g(t)} = \lim_{t \rightarrow \infty} \frac{\xi(t)}{\eta_t h(t)} \stackrel{(a)}{=} 0, \quad (54)$$

where in (a) we use that $\eta_t \in [\eta_-, \eta_+]$ and thus $0 = \lim_{t \rightarrow \infty} \frac{\xi(t)}{\eta_+ h(t)} \leq \lim_{t \rightarrow \infty} \frac{\xi(t)}{\eta_t h(t)} \leq \lim_{t \rightarrow \infty} \frac{\xi(t)}{\eta_- h(t)} = 0$.

Summarizing, we have $g(t)$ strictly monotone and unbounded and $\lim_{t \rightarrow \infty} \frac{\rho(t+1) - \rho(t)}{g(t+1) - g(t)} = 0$. Thus, by using the Stolz-Cesaro (Theorem 21), we have $\lim_{t \rightarrow \infty} \frac{\rho(t)}{g(t)} = 0$.

Continuous time dynamics Here we have $\dot{w}_{(t)} := \frac{dw_{(t)}}{dt} = \Delta w_{(t)} = \bar{w}_\infty h(t) + \xi(t)$.

Define $g(t) := \int_{u=0}^t h(u) du$ and $\rho(t) := w_{(t)} - \bar{w}_\infty g(t) = \int_{u=0}^t \xi(u) + w_{(0)}$ with $\dot{g}(t) := \frac{dg(t)}{dt} = h(t)$ and $\dot{\rho}(t) := \frac{d\rho(t)}{dt} = \xi(t)$ and $\lim_{t \rightarrow \infty} \frac{\dot{\rho}(t)}{\dot{g}(t)} = 0$. In order to prove the lemma, we need to show that $g(t) \rightarrow \infty$ and $\frac{\rho(t)}{g(t)} \rightarrow 0$.

Since $\|w_{(t)}\| \rightarrow \infty$, $\|w_{(t)}\| \leq \int_{t'=0}^t \|\Delta w_{(t')}\| dt = g(t) \xrightarrow{t \in \infty} \infty$.

Thus, we have $g(t) \rightarrow \infty, \forall t < \infty, \dot{g}(t) = h(t) > 0$, and $\lim_{t \rightarrow \infty} \frac{\dot{\rho}(t)}{\dot{g}(t)} = 0$. Thus, using L'Hopitals Rule (Theorem 20), we have $\lim_{t \rightarrow \infty} \frac{\rho(t)}{g(t)} = 0$. \square

D.3. Proof of Theorem 7

Recall Theorem 7:

Theorem 7. For almost all linearly separable datasets $\{X_n, y_n\}_{n=1}^N$, consider the gradient descent iterates $U_{(t)}$ in eq. (18) for minimizing $\tilde{\mathcal{L}}(U)$ with the exponential loss $\ell(u, y) = \exp(-uy)$ and the corresponding sequence of linear predictors $W_{(t)}$ in eq. (19). For any full rank initialization $U_{(0)}$ and any sufficiently small step size sequences η_t such that η_t is smaller than the local Lipschitz at $W_{(t)}$, if $W_{(t)}$ converges to a global minimum i.e., $\mathcal{L}(W_{(t)}) \rightarrow 0$, and additionally the incremental updates $W_{(t+1)} - W_{(t)}$ and the gradients $\nabla \mathcal{L}(W_{(t)})$ converge in direction, then the limit direction $\bar{W}_\infty = \lim_{t \rightarrow \infty} \frac{W_{(t)}}{\|W_{(t)}\|_*}$ exists, and is given by the maximum margin separator with unit nuclear norm $\|\cdot\|_*$,

$$\bar{W}_\infty = \operatorname{argmax}_{W \succcurlyeq 0} \min_n y_n \langle W, X_n \rangle \text{ s.t., } \|W\|_* \leq 1.$$

Proof. From the assumption that $\frac{\Delta W_{(t)}}{\|\Delta W_{(t)}\|_*}$ converges, let $\bar{W}_\infty = \lim_{t \rightarrow \infty} \frac{\Delta W_{(t)}}{\|\Delta W_{(t)}\|_*}$. Lemma D.2 shows the first part of the theorem that $W_{(t)}$ normalized by the nuclear norm converges $\frac{W_{(t)}}{\|W_{(t)}\|_*} \xrightarrow{t \rightarrow \infty} \bar{W}_\infty$.

Also, since $W_{(t)}$ minimizes a strictly monotone loss, we have that $\|W_{(t)}\|_* \rightarrow \infty$ and $\forall n, y_n \langle \bar{W}_\infty, X_n \rangle > 0$.

Let $\gamma = \min_n y_n \langle \bar{W}_\infty, X_n \rangle$, in order to show (b), we equivalently show that $\bar{\bar{W}}_\infty := \bar{W}_\infty / \gamma$ is the solution to the following nuclear norm constrained maximum margin solution:

$$W^* = \operatorname{argmin}_{W \succcurlyeq 0} \|W\|_* \text{ s.t., } \forall n, y_n \langle W, X_n \rangle \geq 1. \quad (55)$$

The KKT optimality conditions of (55) is given by

Stationarity:
$$W^* = \sum_n \alpha_n X_n X_n^T W^*, \quad (56)$$

Complementary slackness:
$$\alpha_n = 0, \forall i \notin S := \{i \in [m] : y_n \langle W^*, X_n \rangle = 1\}, \quad (57)$$

Dual feasibility:
$$\alpha \geq 0, \text{ and } I - \sum_n \alpha_n X_n X_n^T \succeq 0, \quad (58)$$

Primal feasibility:
$$y_n \langle W^*, X_n \rangle \geq 1, W^* \succeq 0. \quad (59)$$

Primal feasibility We already get primal feasibility for $\bar{\bar{W}}_\infty := \bar{W}_\infty / \gamma$ as it has unit margin by the scaling, and from (19), $\forall t W_{(t)} = U_{(t)} U_{(t)}^\top \succcurlyeq 0$, and hence it must converge in direction to a p.s.d. matrix $\bar{W}_\infty \succcurlyeq 0$.

Dual feasibility and complementary slackness: Let $S = \{n : y_n \langle \bar{W}_\infty, X_n \rangle = \gamma\} = \{n : y_n \langle \bar{\bar{W}}_\infty, X_n \rangle = 1\}$, and $\lambda_{\max}(\cdot)$ denote the maximum eigenvalue of a symmetric matrix Z .

For a p.s.d. \bar{W}_∞ , we have $\langle -\nabla \mathcal{L}(W_{(t)}), \bar{W}_\infty \rangle = \sum_n \exp(-y_n \langle W_{(t)}, X_n \rangle) \langle \bar{W}_\infty, X_n \rangle > 0$ as $\langle \bar{W}_\infty, X_n \rangle \geq \gamma > 0$. This implies, $\lambda_{\max}(-\nabla \mathcal{L}(W_{(t)})) > 0$. From the assumption in Theorem 7, we have that $\frac{-\nabla \mathcal{L}(W_{(t)})}{\|-\nabla \mathcal{L}(W_{(t)})\|}$ converges, thus denote $Z_\infty := \lim_{t \rightarrow \infty} \frac{-\nabla \mathcal{L}(W_{(t)})}{\lambda_{\max}(-\nabla \mathcal{L}(W_{(t)}))}$.

Using Lemma 8, we have $Z_\infty = \sum_{n \in S} \alpha_n X_n$ for some $\alpha_n \geq 0$ (with $\alpha_n = 0, \forall n \notin S$). We propose $\{\alpha_n\}$ as our candidate dual certificate for $\bar{\bar{W}}_\infty$ which satisfies complementary slackness by definition. Further, as $\lambda_{\max}(Z_\infty) = 1$, we also have $I - Z_\infty = I - \sum_n \alpha_n X_n X_n^T \succcurlyeq 0$.

Stationarity: This is the main condition to verify: $\bar{\bar{W}}_\infty = Z_\infty \bar{\bar{W}}_\infty$, or equivalently $\bar{W}_\infty = Z_\infty \bar{W}_\infty$.

From Lemma D.2 and Lemma 8), we have the following

1. From assumption that $\frac{\Delta W_{(t)}}{\|\Delta W_{(t)}\|_*} \rightarrow \bar{W}_\infty$, we define the following (note that unlike eq. (53), we have absorbed η into definition of $\Delta W_{(t)}$ here):

$$\Delta W_{(t)} = W_{(t+1)} - W_{(t)} = \bar{W}_\infty h(t) + \xi(t) \text{ s.t., } h(t) = \|\Delta W_{(t)}\|_*, \frac{\xi(t)}{h(t)} \rightarrow 0, \|\bar{W}_\infty\|_* = 1, \quad (60)$$

2. From the construction in eq. (53) in proof of Lemma D.2:

$$W(t) = \bar{W}_\infty g(t) + \rho(t), \text{ where } g(t) = \sum_u h(u), \rho(t) = \sum_u \xi(u) \text{ and } \frac{\rho(t)}{g(t)} \rightarrow 0. \quad (61)$$

3. Since $\mathcal{L}(W(t)) \rightarrow 0$, $\nabla \mathcal{L}(W(t)) \rightarrow 0$. Thus, using Lemma 8 we have $Z_\infty = \sum_{n \in S} \alpha_n$, such that

$$-\nabla \mathcal{L}(W(t)) = Z_\infty p(t) + \zeta(t), \text{ where } \frac{\zeta(t)}{p(t)} \rightarrow 0 \text{ and } p(t) = \lambda_{\max}(-\nabla \mathcal{L}(W(t))) \rightarrow 0. \quad (62)$$

4. $\nabla \mathcal{L}(W(t)) W(t) \nabla \mathcal{L}(W(t)) = p(t) g(t) \delta_1(t)$ where $\delta_1(t) := p(t) Z_\infty \bar{W}_\infty Z_\infty + Z_\infty \bar{W}_\infty \frac{\zeta(t)}{p(t)} + Z_\infty \frac{\rho(t)}{g(t)} Z_\infty + Z_\infty \frac{\rho(t)}{g(t)} \frac{\zeta(t)}{p(t)} + \frac{\zeta(t)}{p(t)} \bar{W}_\infty Z_\infty + \frac{\zeta(t)}{p(t)} \bar{W}_\infty \frac{\zeta(t)}{p(t)} + \frac{\zeta(t)}{p(t)} \frac{\rho(t)}{g(t)} Z_\infty + \frac{\zeta(t)}{p(t)} \frac{\rho(t)}{g(t)} \frac{\zeta(t)}{p(t)} \rightarrow 0$.

Using $W(t)$ and $-\nabla \mathcal{L}(W(t))$ from (61) and (62), respectively, for the updates $\Delta W(t)$ from eq. (19), we have

$$\begin{aligned} \Delta W(t) &= -\eta \nabla \mathcal{L}(W(t)) W(t) - \eta W(t) \nabla \mathcal{L}(W(t)) + \eta^2 \nabla \mathcal{L}(W(t)) W(t) \nabla \mathcal{L}(W(t)) \\ &= \eta p(t) g(t) \left[Z_\infty \bar{W}_\infty + \bar{W}_\infty Z_\infty + Z_\infty \frac{\rho(t)}{g(t)} + \frac{\rho(t)}{g(t)} Z_\infty + \frac{\zeta(t)}{p(t)} \bar{W}_\infty + \bar{W}_\infty \frac{\zeta(t)}{p(t)} + \eta \delta_1(t) \right] \\ &\stackrel{(a)}{=} \eta p(t) g(t) [Z_\infty \bar{W}_\infty + \bar{W}_\infty Z_\infty + \delta(t)] \end{aligned} \quad (63)$$

where in (a) we set $\delta(t) = Z_\infty \frac{\rho(t)}{g(t)} + \frac{\rho(t)}{g(t)} Z_\infty + \frac{\zeta(t)}{p(t)} \bar{W}_\infty + \bar{W}_\infty \frac{\zeta(t)}{p(t)} + \eta \delta_1(t) \rightarrow 0$ as $\frac{\rho(t)}{g(t)}, \frac{\zeta(t)}{p(t)}, \delta_1(t) \rightarrow 0$.

In eq. (63), using $\frac{\Delta W(t)}{\|\Delta W(t)\|_*} \rightarrow \bar{W}_\infty$ with $W_\infty \succcurlyeq 0$ and $h(t) = \|\Delta W(t)\|_*$, we have,

$$1 = \|\bar{W}_\infty\|_* = \langle \bar{W}_\infty, I \rangle = \lim_{t \rightarrow \infty} \left\langle \frac{\Delta W(t)}{\|\Delta W(t)\|_*}, I \right\rangle = \lim_{t \rightarrow \infty} \frac{2\eta p(t) g(t)}{h(t)} \langle Z_\infty W_\infty, I \rangle. \quad (64)$$

$$\Rightarrow \lim_{t \rightarrow \infty} \frac{2\eta p(t) g(t)}{h(t)} = \frac{1}{\langle Z_\infty, \bar{W}_\infty \rangle} := D \stackrel{(a)}{\geq} 1, \quad (65)$$

$$\Rightarrow \bar{W}_\infty = \lim_{t \rightarrow \infty} \frac{\eta p(t) g(t)}{\Delta g(t)} [Z_\infty \bar{W}_\infty + \bar{W}_\infty Z_\infty + \delta(t)] = \frac{D}{2} (Z_\infty \bar{W}_\infty + \bar{W}_\infty Z_\infty) \quad (66)$$

where constant D defined above is independent of t , and in (a) we use $\langle \bar{W}_\infty, I - Z_\infty \rangle \geq 0$ as $\bar{W}_\infty, I - Z_\infty \succcurlyeq 0$, and hence $\langle \bar{W}_\infty, Z_\infty \rangle \leq \langle \bar{W}_\infty, I \rangle = 1$.

Claim 2. If $D = 1$, then the stationarity condition in (56) holds.

Proof. If $D = 1$, from eq. 64, $\langle \bar{W}_\infty, I \rangle = \langle \bar{W}_\infty, Z_\infty \rangle \implies \text{trace}(\bar{W}_\infty (I - Z_\infty)) = 0 \implies \bar{W}_\infty = \bar{W}_\infty Z_\infty$, where the last implication follows as both \bar{W}_∞ and $(I - Z_\infty)$ are p.s.d. \square

Showing $D = 1$ Let $Z_\infty = \sum_i \lambda_i z_i z_i^\top$ be the eigenvalue decomposition of Z_∞ with $\lambda_1 = \lambda_{\max}(Z_\infty) = 1$. Let $\mu_i(t) = \langle W(t), z_i z_i^\top \rangle \geq 0$, $\bar{\mu}_i(t) = \frac{\mu_i(t)}{g(t)}$, and $\bar{\mu}_i^\infty = \lim_{t \rightarrow \infty} \bar{\mu}_i(t) = \langle \bar{W}_\infty, z_i z_i^\top \rangle$.

For any W , we have $\langle Z_\infty W, z_i z_i^\top \rangle = \lambda_i \langle W, z_i z_i^\top \rangle$ from the eigenvalue decomposition of Z_∞ .

Claim 3. For all i , $\bar{\mu}_i^\infty > 0 \implies \lambda_i = 1/D$.

Proof. From (66) we have, $\mu_i^\infty = \langle \bar{W}_\infty, z_i z_i^\top \rangle = D \langle Z_\infty \bar{W}_\infty, z_i z_i^\top \rangle = D \lambda_i \bar{\mu}_i^\infty \implies D = \lambda_i^{-1}$. \square

In particular, from above proposition, if $\exists i : \lambda_i = 1$ and $z_i^\top \bar{W}_\infty z_i > 0$, then $D = 1$ and thus, $\bar{W}_\infty = Z_\infty \bar{W}_\infty$.

Assume the contrary that

$$\lim_{t \rightarrow \infty} \bar{\mu}_1(t) \rightarrow 0 \text{ and } \exists k \text{ s.t., } \lambda_k < 1 \text{ and } \lim_{t \rightarrow \infty} \bar{\mu}_k(t) = 1/D > 0 \implies \lim_{t \rightarrow \infty} \frac{\bar{\mu}_1(t)}{\bar{\mu}_k(t)} = 0. \quad (67)$$

We will show that this is not possible:

Step I: Expression for $\mu_i(t+1)$: From (63), we have $\Delta W_{(t)} = \eta p(t)[Z_\infty W_{(t)} + W_{(t)} Z_\infty + g(t)\bar{\delta}(t)]$ for $\bar{\delta}(t) = (p(t)g(t))^{-1}(\zeta(t)W_{(t)} + W_{(t)}\zeta(t)) \rightarrow 0$. Thus,

$$\Delta\mu_i(t) := \mu_i(t+1) - \mu_i(t) = \langle \Delta W_{(t)}, z_i z_i^\top \rangle = 2\eta p(t)\lambda_i \langle W_{(t)}, z_i z_i^\top \rangle + \eta p(t)g(t) \langle \bar{\delta}(t), z_i z_i^\top \rangle.$$

Defining $\bar{\delta}_i(t) = 1/2 \langle \bar{\delta}(t), z_i z_i^\top \rangle \rightarrow 0$ and using $\mu_i(t) = g(t)\bar{\mu}_i(t)$, we have $\forall i$,

$$\mu_i(t+1) = g(t) \left[(1 + 2\eta p(t)\lambda_i)\bar{\mu}_i + 2\eta p(t)\bar{\delta}_i(t) \right].$$

Step II: Bound on $\frac{\mu_1(t+1)}{\mu_k(t+1)}$: Defining $\kappa(t) := \frac{\bar{\mu}_1(t)}{\bar{\mu}_k(t)} = \frac{\mu_1(t)}{\mu_k(t)}$, we have the following:

$$\begin{aligned} \kappa(t+1) &= \frac{\mu_1(t+1)}{\mu_k(t+1)} = \frac{(1 + 2\eta p(t)\lambda_1)\bar{\mu}_1}{(1 + 2\eta p(t)\lambda_k)\bar{\mu}_k + 2\eta p(t)\bar{\delta}_k(t)} + \frac{2\eta p(t)\bar{\delta}_1(t)}{(1 + 2\eta p(t)\lambda_k)\bar{\mu}_k + 2\eta p(t)\bar{\delta}_k(t)} \\ &\stackrel{(a)}{\geq} \left(\frac{1 + 2\eta p(t)\lambda_1}{1 + 2\eta p(t) \left(\lambda_k + \frac{\bar{\delta}_k(t)}{\bar{\mu}_k(t)} \right)} \right) \frac{\bar{\mu}_1(t)}{\bar{\mu}_k(t)} - \frac{|2\eta p(t)\bar{\delta}_1(t)|}{(1 + 2\eta p(t)\lambda_k)\bar{\mu}_k + 2\eta p(t)\bar{\delta}_k(t)} \\ &:= \tau(t)\kappa(t) - \tilde{\delta}(t), \end{aligned} \tag{68}$$

where (a) follows by dividing the numerator and denominator by $\bar{\mu}_k(t) > 0$.

Step III: Show $\kappa(t) \rightarrow \infty$: The following propositions are proved in Section D.4

Proposition 16. $\exists \bar{\epsilon} > 0, t_0$ s.t., $\forall t > t_0, \tau(t) \geq 1 + 2\eta p(t)\bar{\epsilon}$. In particular, $\bar{\epsilon} = \frac{\lambda_1 - \lambda_k}{2(\lambda_1 + \lambda_k)} > 0$.

Proposition 17. For any t_0 , we have $\frac{\sum_{u=t_0}^t \bar{\delta}(u)}{\sum_{u=t_0}^t 2\eta p(u)} \rightarrow 0$, and further, $\sum_{u=t_0}^t 2\eta p(u) \rightarrow \infty$.

Thus, extending eq. 68, we have $\forall t > t_0$ and $\bar{\epsilon} = \frac{\lambda_1 - \lambda_k}{2(\lambda_1 + \lambda_k)} > 0$,

$$\begin{aligned} \kappa(t+1) &= \tau(t)\kappa(t) - \tilde{\delta}(t) \stackrel{(a)}{\geq} \prod_{u=t_0}^t (1 + 2\eta p(u)\bar{\epsilon}) \kappa(t_0) - \sum_{u=t_0}^t \tilde{\delta}(t) \\ &\stackrel{(b)}{\geq} \left(1 + \sum_{u=t_0}^t 2\eta p(u)\bar{\epsilon} \right) \kappa(t_0) - \sum_{u=t_0}^t \tilde{\delta}(t) \\ &= \kappa(t_0) + \left(\sum_{u=t_0}^t 2\eta p(u) \right) \left[\bar{\epsilon}\kappa(t_0) - \frac{\sum_{u=t_0}^t \tilde{\delta}(u)}{\sum_{u=t_0}^t 2\eta p(u)} \right], \end{aligned} \tag{69}$$

where (a) follows from iterating over t and Proposition 16; and (b) by expanding the product and ignoring the higher order positive terms as $2\eta p(t) > 0$.

Since we start with a full rank $U_{(0)}$ and hence a full rank $W_{(0)}$, and $\{X_n\}$ are in general position for almost all datasets, with small enough η (smaller than the inverse of local Lipschitz), for any finite t_0 , the iterates $U_{(t)} = U_{(t-1)} - \eta \nabla \mathcal{L}(W_{(t-1)})U_{(t-1)}$ remains full rank and hence $\kappa(t_0) > 0$, and thus, $\kappa(t_0)\bar{\epsilon} > 0$.

Also, from Proposition 17, we have $\frac{\sum_{u=t_0}^t \bar{\delta}(u)}{\sum_{u=t_0}^t 2\eta p(u)} \rightarrow 0$ and hence for large enough t , $\bar{\epsilon}\kappa(t_0) - \frac{\sum_{u=t_0}^t \bar{\delta}(u)}{\sum_{u=t_0}^t 2\eta p(u)} > 0$. Combining this with $\sum_{u=t_0}^t 2\eta p(u) \rightarrow \infty$, we have $\kappa(t) \rightarrow \infty$, thus contradicting 67. \square

D.4. Proof of Proposition 16 and Proposition 17

Proposition 16. $\exists \bar{\epsilon} > 0, t_0$ s.t., $\forall t > t_0, \tau(t) \geq 1 + 2\eta p(t)\bar{\epsilon}$. In particular, $\bar{\epsilon} = \frac{\lambda_1 - \lambda_k}{2(\lambda_1 + \lambda_k)} > 0$.

Proof. We have $p(t) \rightarrow 0$ (as $\nabla \mathcal{L}(W_{(t)}) \rightarrow 0$), and $\bar{\delta}_k(t) \rightarrow 0$, $\bar{\mu}_k(t) \rightarrow \bar{\mu}_k^\infty > 0 \Rightarrow \frac{\bar{\delta}_k(t)}{\bar{\mu}_k(t)} \rightarrow 0$.

Thus, $\forall \epsilon > 0$, $\exists t_0$ such that $\forall t > t_0$, $2\eta p(t) \leq 1$ and $\frac{\bar{\delta}_k(t)}{\lambda_1 \bar{\mu}_k(t)} < \epsilon$.

Define $\bar{\epsilon} = \frac{1-\lambda_k}{2(1+\lambda_k)} > 0$ since $1 = \lambda_1 > \lambda_k$. We then pick $\epsilon = \frac{1-\lambda_k - \bar{\epsilon}(1+\lambda_k)}{1+\bar{\epsilon}} = \frac{1-\lambda_k}{2(1+\bar{\epsilon})} > 0$.

For $t > t_0$, $2\eta p(t) \leq 1$ and $\frac{\bar{\delta}_k(t)}{\bar{\mu}_k(t)} \leq \epsilon$, where we have

$$0 < \epsilon = \frac{1 - \lambda_k - \bar{\epsilon}(1 + \lambda_k)}{1 + \bar{\epsilon}} \leq \frac{1 - \lambda_k - \bar{\epsilon}(1 + 2\eta p(t)\lambda_k)}{1 + 2\eta p(t)\bar{\epsilon}} = \frac{1 - \bar{\epsilon}}{1 + 2\eta p(t)\bar{\epsilon}} - \lambda_k \quad (70)$$

Now using $\frac{\bar{\delta}_k(t)}{\lambda_1 \bar{\mu}_k(t)} \leq \epsilon$ and $\lambda_1 = 1$, we have the following for $t \geq t_0$

$$\begin{aligned} \tau(t) &= \frac{1 + 2\eta p(t)\lambda_1}{1 + 2\eta p(t) \left(\lambda_k + \frac{\bar{\delta}_k(t)}{\bar{\mu}_k(t)} \right)} \geq \frac{1 + 2\eta p(t)}{1 + 2\eta p(t) (\lambda_k + \epsilon)} \\ &\stackrel{(a)}{\geq} \frac{1 + 2\eta p(t)}{1 + 2\eta p(t) \left(\frac{1 - \bar{\epsilon}}{1 + 2\eta p(t)\bar{\epsilon}} \right)} = 1 + 2\eta p(t)\bar{\epsilon}, \end{aligned} \quad (71)$$

where (a) follows from eq. (70) □

Proposition 17. For any t_0 , we have $\frac{\sum_{u=t_0}^t \bar{\delta}(u)}{\sum_{u=t_0}^t 2\eta p(u)} \rightarrow 0$.

Proof. **First show $\sum_{u=t_0}^t 2\eta p(u) \rightarrow \infty$:** Recall that $g(t) = \sum_{u < t} h(u)$, where $h(u) = \|\Delta W_{(t)}\|_* > 0$ (from eq. (61)-(60)). We then have

$$0 < \log g(t+1) - \log g(t) = \log \frac{g(t+1)}{g(t)} \stackrel{(a)}{\leq} \frac{g(t+1)}{g(t)} - 1 = \frac{g(t+1) - g(t)}{g(t)} = \frac{h(t)}{g(t)},$$

where in (a) we used $\log(x) \leq x - 1$. Summing over t , we have $\sum_{u=t_0}^t \frac{h(u)}{g(u)} \geq \log \frac{g(t+1)}{g(t_0)} \rightarrow \infty$.

Also, recall from eq. (65) that $\frac{2\eta p(t)}{h(t)/g(t)} \rightarrow D$.

Now using $a_t = \sum_{u=t_0}^t 2\eta p(u)$ and monotonically increasing divergent sequence $b_t = \sum_{u=t_0}^t \frac{h(u)}{g(u)} \rightarrow \infty$ in Stolz-Cesaro theorem (Theorem 21), we get

$$\lim_{t \rightarrow \infty} \frac{\sum_{u=t_0}^t 2\eta p(u)}{\sum_{u=t_0}^t \frac{h(u)}{g(u)}} = \lim_{t \rightarrow \infty} \frac{a_t}{b_t} = \lim_{t \rightarrow \infty} \frac{a_t - a_{t-1}}{b_t - b_{t-1}} = \lim_{t \rightarrow \infty} \frac{2\eta p(t)g(t)}{h(t)} = D. \quad (72)$$

Hence, $\lim_{t \rightarrow \infty} \sum_{u=t_0}^t 2\eta p(u) = D \lim_{t \rightarrow \infty} \sum_{u=t_0}^t \frac{\Delta g(u)}{g(u)} = \infty$.

Bound $\sum_{u=t_0}^t \tilde{\delta}(t)$: We have from definition $0 < \tilde{\delta}(t) = \frac{|2\eta p(t)\bar{\delta}_1(t)|}{(1+2\eta p(t)\lambda_k)\bar{\mu}_k + 2\eta p(t)\bar{\delta}_k(t)} \leq 2\eta p(t) \frac{|\bar{\delta}_1(t)|}{\bar{\mu}_k(t)} \rightarrow 0$

Since $\bar{\delta}_1(t) \rightarrow 0$ and $\bar{\mu}_k(t) \rightarrow \bar{\mu}_k^\infty > 0$, we also have $\frac{\tilde{\delta}(t)}{2\eta p(t)} \rightarrow 0$.

Again using Stolz-Cesaro theorem with $c_t = \sum_{u=t_0}^t \tilde{\delta}(u)$ and $d_t = \sum_{u=t_0}^t 2\eta p(u) \rightarrow \infty$, we have

$$\lim_{t \rightarrow \infty} \frac{\sum_{u=t_0}^t \tilde{\delta}(u)}{\sum_{u=t_0}^t 2\eta p(u)} = \lim_{t \rightarrow \infty} \frac{c_t - c_{t-1}}{d_t - d_{t-1}} = \lim_{t \rightarrow \infty} \frac{\tilde{\delta}(t)}{2\eta p(t)} = 0. \quad (73)$$

□

E. Preliminaries

Lemma 18 (Sub-differentials of norms). *For a generic norm $\|v\|$ for $v \in \mathcal{V}$, recall the dual norm $\|y\|_* = \sup_{\|v\| \leq 1} \langle y, v \rangle$. The sub-differential of a norm $\|\cdot\|$ at v is defined as $\partial\|v\| = \{y : \forall \Delta \in \mathcal{V}, \|v + \Delta\| \geq \|v\| + \langle y, \Delta \rangle\}$.*

We have the following results on the properties on the sub-differentials are readily established:

1. $\partial\|v\| = \{y : \|y\|_* = 1, \text{ and } \langle y, v \rangle = \|v\|\}$
2. $y \in \partial\|v\|^2$ if and only if $v \in \partial\|v\|^2$
3. if there exists $v_1, v_2 \in \mathcal{V}$ and $g \in \mathcal{V}^*$ such that $g \in \partial\|v_1\|$ and $g \in \partial\|v_2\|$, then for all $\alpha, \beta > 0$, $g \in \partial\|\alpha v_1 + \beta v_2\|$.

Proof. 1. It can be easily verified that $\{y : \|y\|_* = 1, \text{ and } \langle y, v \rangle = \|v\|\} \subseteq \partial\|v\|$. Conversely, $\forall y \in \partial\|v\|$, from the definition, we have $\forall \Delta, \|v\| + \|\Delta\| \geq \|v + \Delta\| \geq \|v\| + \langle y, \Delta \rangle \implies \|y\|_* = \sup_{\|\Delta\| \neq 0} \langle \frac{\Delta}{\|\Delta\|}, y \rangle \leq 1$. Using $\|y\|_* \leq 1$ along with $\Delta = -v$, we have $\langle y, v \rangle \geq \|v\| = \sup_{\|y\|_* \leq 1} \langle v, y \rangle \implies \langle y, v \rangle = \sup_{\|y\|_* \leq 1} \langle v, y \rangle \|v\|$, which by homogeneity of norms implies $\|y\|_* = 1$.

2. From above result, $y \in \partial\frac{1}{2}\|v\|^2 \Leftrightarrow \|y\|_* = \|v\|$ and $\langle y, v \rangle = \|v\|^2 = \|y\|_*^2 \Leftrightarrow v \in \partial\|y\|_*^2$.

3. $g \in \partial\|v_1\| \cap \partial\|v_2\|$ implies $\|g\|_* = 1$, $\|v_1\| = \langle g, v_1 \rangle$, and $\|v_2\| = \langle g, v_2 \rangle$. Using triangle inequality, $\|\alpha v_1 + \beta v_2\| \leq \alpha\|v_1\| + \beta\|v_2\| = \langle g, \alpha v_1 + \beta v_2 \rangle \leq \sup_{\|y\|_* \leq 1} \langle y, \alpha v_1 + \beta v_2 \rangle = \|\alpha v_1 + \beta v_2\| \implies \|\alpha v_1 + \beta v_2\| = \langle g, \alpha v_1 + \beta v_2 \rangle$. □

Lemma 19 (Limit points of a compact sets). *If $\{a_t\}_{t=1}^\infty$ is a sequence contained in a compact set $a_t \in C$, then there exists at least one limit point of $\{a_t\}$ in C . That is, $\exists a^\infty \in C$ and a subsequence $\{a_{t_k}\}_{k=1}^\infty$, such that $\lim_{k \rightarrow \infty} a_{t_k} = a^\infty$.*

Theorem 20 (L-Hopital's Rule, proof in Theorem 30.2 of Ross (1980)). *Let $s \in \mathbb{R} \cup \{-\infty, \infty\}$, and $f(x)$ and $g(x)$ be continuous and differentiable functions such that $\lim_{x \rightarrow s} \frac{f'(x)}{g'(x)} = L$ exists. If either (a) $\lim_{x \rightarrow s} f(x) = \lim_{x \rightarrow s} g(x) = 0$, or (b) $\lim_{x \rightarrow s} |g(x)| = \infty$, then $\lim_{x \rightarrow s} \frac{f(x)}{g(x)}$ exists and is equal to L .*

Theorem 21 (Stolz–Cesaro theorem, proof in Theorem 1.22 of Muresan & Muresan (2009)). *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 (i.e., monotonic increasing with $b_k \rightarrow \infty$ or monotonic decreasing with $b_k \rightarrow -\infty$). 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 .*