
Learning threshold neurons via the “edge of stability”

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Existing analyses of neural network training often operate under the unrealistic as-
2 sumption of an extremely small learning rate. This lies in stark contrast to practical
3 wisdom and empirical studies, such as the work of J. Cohen et al. (ICLR 2021),
4 which exhibit startling new phenomena (the “edge of stability” or “unstable conver-
5 gence”) and potential benefits for generalization in the large learning rate regime.
6 Despite a flurry of recent works on this topic, however, the latter effect is still
7 poorly understood. In this paper, we take a step towards understanding genuinely
8 non-convex training dynamics with large learning rates by performing a detailed
9 analysis of gradient descent for simplified models of two-layer neural networks.
10 For these models, we provably establish the edge of stability phenomenon and
11 discover a sharp phase transition for the step size below which the neural network
12 fails to learn “threshold-like” neurons (i.e., neurons with a non-zero first-layer bias).
13 This elucidates one possible mechanism by which the edge of stability can in fact
14 lead to better generalization, as threshold neurons are basic building blocks with
15 useful inductive bias for many tasks.

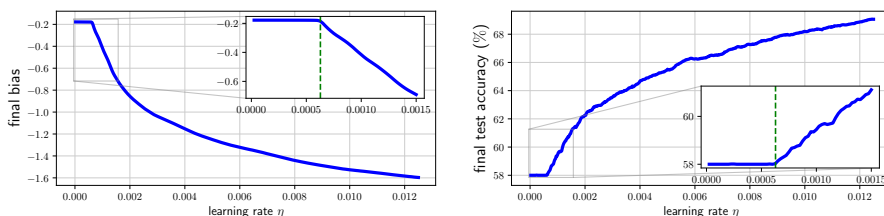


Figure 1: **Large step sizes are necessary to learn the “threshold neuron” of a ReLU network (2) for a simple binary classification task (1).** We choose $d = 200$, $n = 300$, $\lambda = 3$, and run gradient descent with the logistic loss. The weights are initialized as $a^-, a^+ \sim \mathcal{N}(0, 1/(2d))$ and $b = 0$. For each learning rate η , we set the iteration number such that the total time elapsed (iteration $\times \eta$) is 10. The vertical dashed lines indicate **our theoretical prediction** of the phase transition phenomenon (precise threshold at $\eta = 8\pi/d^2$).

16 1 Introduction

17 How much do we understand about the training dynamics of neural networks? We begin with a
18 simple and canonical learning task which indicates that the answer is still “far too little”.

19 **Motivating example:** Consider a binary classification task of labeled pairs $(\mathbf{x}^{(i)}, y^{(i)}) \in \mathbb{R}^d \times \{\pm 1\}$
20 where each covariate $\mathbf{x}^{(i)}$ consists of a 1-sparse vector (in an unknown basis) corrupted by additive
21 Gaussian noise, and the label $y^{(i)}$ is the sign of the non-zero coordinate of the 1-sparse vector. Due to

22 rotational symmetry, we can take the unknown basis to be the standard one and write

$$\mathbf{x}^{(i)} = \lambda y^{(i)} \mathbf{e}_{j(i)} + \boldsymbol{\xi}^{(i)} \in \mathbb{R}^d, \quad (1)$$

23 where $y^{(i)} \in \{\pm 1\}$ is a random label, $j(i) \in [d]$ is a random index, $\boldsymbol{\xi}^{(i)}$ is Gaussian noise, and
 24 $\lambda > 1$ is the unknown signal strength. In fact, (1) is a special case of the well-studied sparse coding
 25 model [OF97; VG00; OF04; Yan+09; AL22; KR18]. We ask the following fundamental question:

26 *How do neural networks learn to solve the sparse coding problem (1)?*

27 In spite of the simplicity of the setting, a full resolution to this question requires a thorough under-
 28 standing of surprisingly rich dynamics which **lies out of reach of existing theory**. To illustrate this
 29 point, consider an extreme simplification in which the basis $\mathbf{e}_1, \dots, \mathbf{e}_d$ is known in advance, for
 30 which it is natural to parametrize a two-layer ReLU network as

$$f(\mathbf{x}; a^-, a^+, b) = a^- \sum_{i=1}^d \text{ReLU}(-\mathbf{x}[i] + b) + a^+ \sum_{i=1}^d \text{ReLU}(\mathbf{x}[i] + b). \quad (2)$$

31 The parametrization (2) respects the latent data structure (1) well: a good network has a negative bias
 32 b to threshold out the noise, and has $a^- < 0$ and $a^+ > 0$ to output correct labels. We are particularly
 33 interested in understanding the mechanism by which the bias b becomes negative, thereby allowing
 34 the non-linear ReLU activation to act as a threshold function; we refer to this as the problem of
 35 learning “threshold neurons”. More broadly, such threshold neurons are of interest as they constitute
 36 basic building blocks for producing neural networks with useful inductive bias.

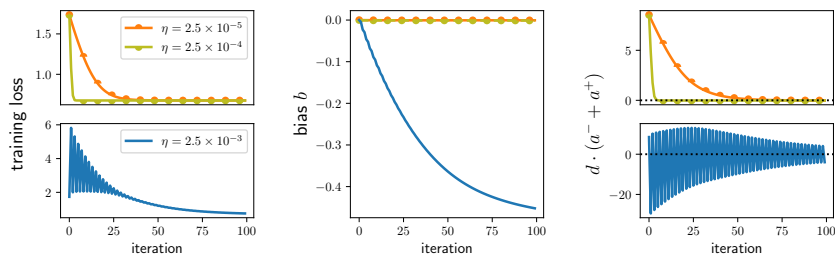


Figure 2: **Large learning rates lead to unexpected phenomena: non-monotonic loss and wild oscillations of weights.** We choose the same setting as Figure 1. With a small learning rate ($\eta = 2.5 \cdot 10^{-5}$), the bias does not decrease noticeably, and the same is true even when we increase the learning rate by ten times ($\eta = 2.5 \cdot 10^{-4}$). When we increase the learning rate by another ten times ($\eta = 2.5 \cdot 10^{-3}$), we finally see a noticeable decrease in the bias, but with this we observe unexpected behavior: *the loss decreases non-monotonically and the sum of second-layer weights $d \cdot (a^- + a^+)$ oscillates wildly.*

37 We train the parameters a^-, a^+, b using gradient descent with step size $\eta > 0$ on the logistic loss
 38 $\sum_{i=1}^n \ell_{\text{logi}}(y^{(i)} f(\mathbf{x}^{(i)}; a^-, a^+, b))$, where $\ell_{\text{logi}}(z) := \log(1 + \exp(-z))$, and we report the results
 39 in Figures 1 and 2. The experiments reveal a compelling picture of the optimization dynamics.

- 40 ■ **Large learning rates are necessary, both for generalization and for learning threshold**
 41 **neurons.** Figure 1 shows that the bias decreases and the test accuracy increases as we increase
 42 η ; note that we plot the results after a fixed *time* (iteration $\times \eta$), so the observed results are not
 43 simply because larger learning rates track the continuous-time gradient flow for a longer time.
- 44 ■ **Large learning rates lead to unexpected phenomena: non-monotonic loss and wild oscillations**
 45 **of $a^- + a^+$.** Figure 2 shows that large learning rates also induce stark phenomena, such as
 46 non-monotonic loss and large weight fluctuations, which lie firmly outside the explanatory power
 47 of existing analytic techniques based on principles from convex optimization.
- 48 ■ **There is a phase transition between small and large learning rates.** In Figure 1, we zoom in
 49 on learning rates around $\eta \approx 0.0006$ and observe *sharp* phase transition phenomena.

50 We have presented these observations in the context of the simple ReLU network (2), but we
 51 emphasize that **these findings are indicative of behaviors observed in practical neural network**
 52 **training settings.** In Figure 3, we display results for a two-layer ReLU network trained on the full
 53 sparse coding model (1) with unknown basis, as well as a deep neural network trained on CIFAR-10.

54 In each case, we again observe non-monotonic loss coupled with steadily decreasing bias parameters.
 55 For these richer models, the transition from small to large learning rates is oddly reminiscent of well-
 56 known separations between the “lazy training” or “NTK” regime [JGH18] and the more expressive
 57 “feature learning” regime. For further experimental results, see §A.

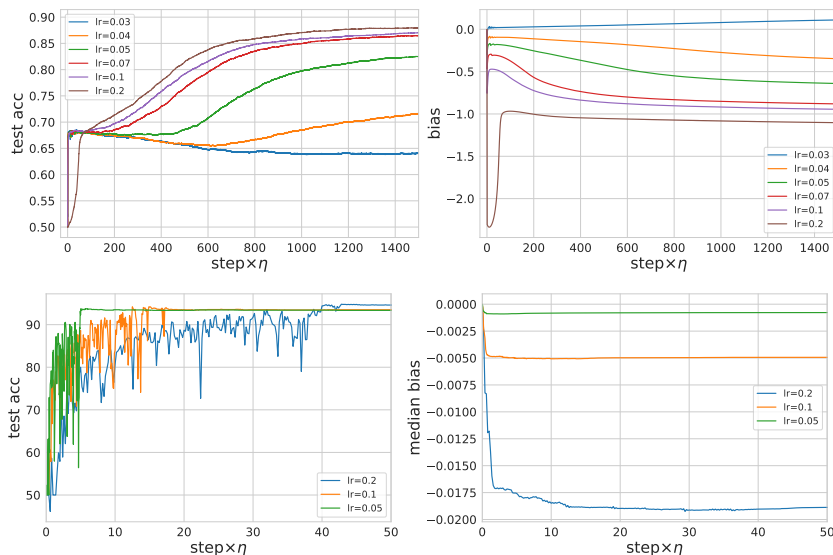


Figure 3: (Top) Results for training an over-parametrized two-layer neural network $f(\mathbf{x}; \mathbf{a}, \mathbf{W}, b) = \sum_{i=1}^m a_i \text{ReLU}(\mathbf{w}_i^\top \mathbf{x} + b)$ with $m \gg d$ for the **full sparse coding model (1)**; in this setting, the basis vectors are unknown, and the neural network learn them through additional parameters $\mathbf{W} = (\mathbf{w}_i)_{i=1}^m$. Also, we use m different weights $\mathbf{a} = (a_i)_{i=1}^m$ for the second layer. (Bottom) Full-batch gradient descent dynamics of **ResNet-18 on (binary) CIFAR-10 with various learning rates**. Details are deferred to §A.

58 We currently do not have right tools to understand these phenomena. First of all, a drastic change
 59 in behavior between the small and the large learning rates cannot be captured through well-studied
 60 regimes, such as the “neural tangent kernel” (NTK) regime [JGH18; ALS19; Aro+19; COB19;
 61 Du+19; OS20] or the mean-field regime [CB18; MMM19; Chi22; NWS22; RV22]. In addition,
 62 understanding why a large learning rate is required to learn the bias is beyond the scope of prior
 63 theoretical works on the sparse coding model [Aro+15; Kar+21]. Our inability to explain these
 64 findings points to a serious gap in our grasp of neural network training dynamics and calls for a
 65 detailed theoretical study.

66 1.1 Main scope of this work

67 In this work, we do not aim to understand the sparse coding problem (1) in its full generality. Instead,
 68 we pursue the more modest goal of shedding light on the following question.

69 **Q.** What is the role of a large step size in learning the bias for the ReLU network (2)?

70 As discussed above, the dynamics of the simple ReLU network (2) is a microcosm of emergent
 71 phenomena beyond the convex optimization regime. In fact, there is a recent growing body of
 72 work [Coh+21; ALP22; AZS22; CB22; LLA22; Ma+22; WLL22; DNL23; Zhu+23] on training
 73 with large learning rates, which largely aims at explaining a striking empirical observation called the
 74 “**edge of stability (EoS)**” phenomenon.

75 The edge of stability (EoS) phenomenon is a set of distinctive behaviors observed recently by
 76 [Coh+21] when training neural networks with gradient descent (GD). Here we briefly summarize the
 77 salient features of the EoS and defer a discussion of prior work to §1.3. Recall that if we use GD to
 78 optimize an L -smooth loss function with step size η , then the well-known descent lemma from convex
 79 optimization ensures monotonic decrease in the loss so long as $L < 2/\eta$. In contrast, when $L > 2/\eta$,
 80 it is easy to see on simple convex quadratic examples that GD can be unstable (or divergent). The

81 main observation of [Coh+21] is that when training neural networks¹ with constant step size $\eta > 0$,
 82 the largest eigenvalue of the Hessian at the current iterate (dubbed the “sharpness”) initially increases
 83 during training (“progressive sharpening”) and saturates near or above $2/\eta$ (“EoS”).

84 A surprising message of the present work is that **the answer to our main question is intimately**
 85 **related to the EoS.** Indeed, Figure 11 shows that the GD iterates of our motivating example exhibit
 86 the EoS during the initial phase of training when the bias decreases rapidly.

87 Consequently, we first set out to thoroughly understand the workings of the EoS phenomena through
 88 a simple example. Specifically, we consider a single-neuron linear neural network in dimension 1,
 89 corresponding to the loss

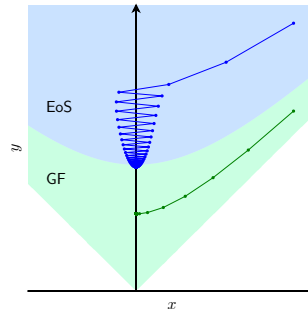
$$\mathbb{R}^2 \ni (x, y) \mapsto \ell(xy), \quad \text{where } \ell \text{ is convex, even, and Lipschitz.} \quad (3)$$

90 Although toy models have appeared in works on the EoS (see §1.3), our example is simpler than all
 91 prior models, and we provably establish the EoS for (3) with transparent proofs.

92 We then use the newfound insights gleaned from the analysis of (3) to answer our main question.
 93 To the best of our knowledge, we provide the first explanation of the mechanism by which a large
 94 learning rate can be *necessary* for learning threshold neurons.

95 1.2 Our contributions

96 **Explaining the EoS with a single-neuron example.** Al-
 97 though the EoS has been studied in various settings (see
 98 §1.3 for a discussion), these works either do not rigorously
 99 establish the EoS phenomenon, or they operate under com-
 100 plex settings with opaque assumptions. Here, we study
 101 a simple two-dimensional loss function, $(x, y) \mapsto \ell(xy)$,
 102 where ℓ is convex, even, and Lipschitz. Some examples in-
 103 clude² $\ell(s) = \frac{1}{2} \log(1 + \exp(-s)) + \frac{1}{2} \log(1 + \exp(+s))$
 104 and $\ell(s) = \sqrt{1 + s^2}$. Surprisingly, GD on this loss already
 105 exhibits rich behavior (Figure 4).



106 En route to this result, we rigorously establish the quasi-
 107 static dynamics formulated in [Ma+22].

Figure 4: Illustration of two different regimes (the “gradient flow” regime and the “EoS” regime) of the GD dynamics.

108 The elementary nature of our example leads to transparent
 109 arguments, and consequently our analysis isolates generalizable principles for “bouncing” dynamics.
 110 To demonstrate this, we use our insights to study our main question of learning threshold neurons.

111 **Learning threshold neurons with the mean model.** The connection between the single-neuron
 112 example and the ReLU network (2) can already be anticipated via a comparison of the dynamics: (i)
 113 for the single neuron example, x oscillates wildly while y decreases (Figure 4); (ii) for the ReLU
 114 network (2), the sum of weights ($a^- + a^+$) oscillates while b decreases (Figure 2).

115 We study this example in §2 and delineate a transition from the “gradient flow” regime to the “EoS
 116 regime”, depending on the step size η and the initialization. Moreover, in the EoS regime, we
 117 rigorously establish asymptotics for the limiting sharpness which depend on the higher-order behavior
 118 of ℓ . In particular, for the two losses mentioned above, the limiting sharpness is $2/\eta + O(\eta)$, whereas
 119 for losses ℓ which are exactly quadratic near the origin the limiting sharpness is $2/\eta + O(1)$.

120 In fact, this connection can be made formal by considering an approximation for the GD dynamics
 121 for the ReLU network (2). It turns out (see §3.1 for details) that during the initial phase of training,
 122 the dynamics of $A_t := d(a_t^- + a_t^+)$ and b_t due to the ReLU network are well-approximated by the
 123 “rescaled” GD dynamics on the loss $(A, b) \mapsto \ell_{\text{sym}}(Ag(b))$, where the step size for the A -dynamics
 124 is multiplied by $2d^2$, $g(b) := \mathbb{E}_{z \sim \mathcal{N}(0,1)} \text{ReLU}(z + b)$ is the “smoothed” ReLU, and ℓ_{sym} is the
 125 symmetrized logistic loss; see §3.1 and Figure 10. We refer to these dynamics as the **mean model**.

¹The phenomenon in [Coh+21] is most clearly observed for tanh activations, although the appendix of [Coh+21] contains thorough experimental results for various neural network architectures.

²Suppose that we have a single-layer linear neural network $f(x; a, b) = abx$, and that the data is drawn according to $x = 1$, $y \sim \text{unif}(\{\pm 1\})$. Then, the population loss under the logistic loss is $(a, b) \mapsto \ell_{\text{sym}}(ab)$ with $\ell_{\text{sym}}(s) = \frac{1}{2} \log(1 + \exp(-s)) + \frac{1}{2} \log(1 + \exp(+s))$.

126 The mean model bears a great resemblance to the single-neuron example $(x, y) \mapsto \ell(xy)$, and hence
127 we can leverage the techniques developed for the latter in order to study the former.

128 Our main result for the mean model precisely explains the phase transition in [Figure 1](#). For any $\delta > 0$,

- 129 • if $\eta \leq (8 - \delta)\pi/d^2$, then **the mean model fails to learn threshold neurons**: the limiting bias
130 satisfies $|b_\infty| = O_\delta(1/d^2)$.
- 131 • if $\eta \geq (8 + \delta)\pi/d^2$, then **the mean model enters the EoS and learns threshold neurons**: the
132 limiting bias satisfies $b_\infty \leq -\Omega_\delta(1)$.

133 1.3 Related work

134 **Edge of stability:** Our work is motivated by the extensive empirical study of [\[Coh+21\]](#), which
135 identified the EoS phenomenon. Subsequently, there has been a flurry of works aiming at developing
136 a theoretical understanding of the EoS, which we briefly summarize here.

137 *Properties of the loss landscape:* The works [\[AZS22; Ma+22\]](#) study the properties of the loss
138 landscape that lead to the EoS. Namely, [\[AZS22\]](#) argue that the existence of forward-invariant subsets
139 near the minimizers allows GD to convergence even in the unstable regime. They also explore various
140 characteristics of EoS in terms of loss and iterates. Also, [\[Ma+22\]](#) empirically show that the loss
141 landscape of neural networks exhibits subquadratic growth locally around the minimizers. They prove
142 that for a one-dimensional loss, subquadratic growth implies that GD finds a 2-periodic trajectory.

143 *Limiting dynamics:* Other works characterize the limiting dynamics of the EoS in various regimes.
144 [\[ALP22; LLA22\]](#) show that (normalized) GD tracks a “sharpness reduction flow” near the manifold
145 of minimizers. The recent work of [\[DNL23\]](#) obtains a different predicted dynamics based on self-
146 stabilization of the GD trajectory. Also, [\[Ma+22\]](#) describes a quasi-static heuristic for understanding
147 the overall trajectory of GD when one component of the iterate is oscillating.

148 *Simple models and beyond:* Closely related to our own approach, there are prior works which
149 carefully study simple models. [\[CB22\]](#) prove global convergence of GD for the two-dimensional
150 function $(x, y) \mapsto (xy - 1)^2$ and a single-neuron student-teacher setting; note that unlike our results,
151 they do not study the limiting sharpness. [\[WLL22\]](#) study progressive sharpening for a neural
152 network model. Also, the recent and concurrent work of [\[Zhu+23\]](#) studies the two-dimensional loss
153 $(x, y) \mapsto (x^2y^2 - 1)^2$; to our knowledge, their work is the first to asymptotically and rigorously
154 show that the limiting sharpness of GD is $2/\eta$ in a simple setting, at least when initialized locally. In
155 comparison, in [§2](#), we perform a global analysis of the limiting sharpness of GD for $(x, y) \mapsto \ell(xy)$
156 for a class of convex, even, and Lipschitz losses ℓ , and in doing so we clearly delineate the “gradient
157 flow regime” from the “EoS regime”.

158 **Effect of learning rate on learning:** Recently, several works have sought to understand how
159 the choice of learning rate affects the learning process, in terms of the properties of the resulting
160 minima [\[Jas+18; WME18; MMS21; Nac+22\]](#) and the behavior of optimization dynamics [\[Xin+18;](#)
161 [Jas+19; Jas+20; Lew+20; Jas+21\]](#).

162 [\[LWM19\]](#) demonstrate for a synthetic data distribution and a two-layer ReLU network model that
163 choosing a larger step size for SGD helps with generalization. Subsequent works have shown similar
164 phenomena for regression [\[Nak20; Wu+21; Ba+22\]](#), kernel ridge regression [\[BMR22\]](#), and linear
165 diagonal networks [\[Nac+22\]](#). However, the large step sizes considered in these work still fall under
166 the scope of descent lemma, and most prior works do not theoretically investigate the effect of large
167 step size in the EoS regime. A notable exception is the work of [\[Wan+22\]](#), which studies the impact
168 of learning rates greater than $2/\text{smoothness}$ for a matrix factorization problem. Also, the recent work
169 of [\[And+22\]](#) seeks to explain the generalization benefit of SGD in the large step size regime by
170 relying on a heuristic SDE model for the case of linear diagonal networks. Despite this similarity,
171 their main scope is quite different from ours, as we *(i)* focus on GD instead of SGD and *(ii)* establish a
172 direct and detailed analysis of the GD dynamics for a model of the motivating sparse coding example.

173 2 Single-neuron linear network

174 In this section, we analyze the single-neuron linear network model $(x, y) \mapsto f(x, y) := \ell(xy)$.

175 **2.1 Basic properties and assumptions**

176 **Basic properties.** If ℓ is minimized at 0, then the *global minimizers* of f are the x - and y -axes. The
 177 GD iterates x_t, y_t , for step size $\eta > 0$ and iteration $t \geq 0$ can be written as

$$x_{t+1} = x_t - \eta \ell'(x_t y_t) y_t, \quad y_{t+1} = y_t - \eta \ell'(x_t y_t) x_t.$$

178 **Assumptions.** From here onward, we assume $\eta < 1$ and the following conditions on $\ell : \mathbb{R} \rightarrow \mathbb{R}$.

179 (A1) ℓ is convex, even, 1-Lipschitz, and of class \mathcal{C}^2 near the origin with $\ell''(0) = 1$.

180 (A2) There exist constants $\beta > 1$ and $c > 0$ with the following property: for all $s \neq 0$,

$$\ell'(s)/s \leq 1 - c |s|^\beta \mathbb{1}\{|s| \leq c\}.$$

181 We allow $\beta = +\infty$, in which case we simply require that $\frac{\ell'(s)}{s} \leq 1$ for all $s \neq 0$.

182 Assumption (A2) imposes decay of $s \mapsto \ell'(s)/s$ locally away from the origin in order to obtain more
 183 fine-grained results on the limiting sharpness in [Theorem 2](#). As we show in [Lemma 5](#) below, when ℓ
 184 is smooth and has a strictly negative fourth derivative at the origin, then Assumption (A2) holds with
 185 $\beta = 2$. See [Example 1](#) for some simple examples of losses satisfying our assumptions.

186 **2.2 Two different regimes for GD depending on the step size**

187 Before stating rigorous results, in this section we begin by giving an intuitive understanding of the
 188 GD dynamics. It turns out that for a given initialization (x_0, y_0) , there are two different regimes for
 189 the GD dynamics depending on the step size η . Namely, there exists a threshold on the step size such
 190 that (i) below the threshold, GD remains close to the gradient flow for all time, and (ii) above the
 191 threshold, GD enters the edge of stability and diverges away from the gradient flow.

192 First, recall that the GD dynamics are symmetric in x, y and that the lines $y = \pm x$ are invariant.
 193 Hence, we may assume without loss of generality that

$$y_0 > x_0 > 0, \quad y_t > |x_t| \text{ for all } t \geq 1, \quad \text{and GD converges to } (0, y_\infty) \text{ for } y_\infty > 0.$$

194 From the expression (8) for the Hessian of f and our normalization $\ell''(0) = 1$, it follows that **the**
 195 **sharpness reached by GD in this example is precisely y_∞^2 .**

196 Initially, in both regimes, the GD dynamics tracks the continuous-time gradient flow. Our first
 197 observation is that the gradient flow admits a conserved quantity, thereby allowing us to predict the
 198 dynamics in this initial phase.

199 **Lemma 1** (conserved quantity). *Along the gradient flow for f , the quantity $y^2 - x^2$ is conserved.*

200 *Proof.* Differentiating $y_t^2 - x_t^2$ with respect to t gives $2y_t (-\ell'(x_t y_t) x_t) - 2x_t (-\ell(x_t y_t) y_t) = 0$. \square

201 [Lemma 1](#) implies that the gradient flow converges to $(0, y_\infty^{\text{GF}}) = (0, \sqrt{y_0^2 - x_0^2})$. For GD with
 202 step size $\eta > 0$, the quantity $y^2 - x^2$ is no longer conserved, but we show in [Lemma 6](#) that it is
 203 *approximately* conserved until the GD iterate lies close to the y -axis. Hence, GD initialized at (x_0, y_0)
 204 also reaches the y -axis approximately at the point $(x_{t_0}, y_{t_0}) \approx (0, \sqrt{y_0^2 - x_0^2})$.

205 At this point, GD either approximately converges to the gradient flow solution $(0, \sqrt{y_0^2 - x_0^2})$ or
 206 diverges away from it, depending on whether or not $y_{t_0}^2 > 2/\eta$. To see this, for $|x_{t_0} y_{t_0}| \ll 1$, we can
 207 Taylor expand ℓ' near zero to obtain the approximate dynamics for x (recalling $\ell''(0) = 1$),

$$x_{t_0+1} \approx x_{t_0} - \eta x_{t_0} y_{t_0}^2 = (1 - \eta y_{t_0}^2) x_{t_0}. \quad (4)$$

208 From (4), we deduce the following conclusions.

209 (i) If $y_{t_0}^2 < 2/\eta$, then $|1 - \eta y_{t_0}^2| < 1$. Since y_t is decreasing, it implies that $|1 - \eta y_t^2| < 1$ for all
 210 $t \geq t_0$, and so $|x_t|$ converges to zero exponentially fast.

211 (ii) On the other hand, if $y_{t_0}^2 > 2/\eta$, then $|1 - \eta y_{t_0}^2| > 1$, i.e., the magnitude of x_{t_0} increases in the
 212 next iteration, and hence GD cannot stabilize. In fact, in the approximate dynamics, x_{t_0+1} has
 213 the opposite sign as x_{t_0} , i.e., x_{t_0} jumps across the y -axis. One can show that the ‘‘bouncing’’ of
 214 the x variable continues until y_t^2 has decreased past $2/\eta$, at which point we are in the previous
 215 case and GD approximately converges to $(0, 2/\eta)$.

216 This reasoning, combined with the expression for the Hessian of f , shows that

$$\begin{aligned} \text{sharpness}(0, y_\infty) &:= \lambda_{\max}(\nabla^2 f(0, y_\infty)) \approx \min\{y_0^2 - x_0^2, 2/\eta\} \\ &= \min\{\text{gradient flow sharpness, EoS prediction}\}. \end{aligned}$$

217 Accordingly, we refer to the case $y_0^2 - x_0^2 < 2/\eta$ as the **gradient flow regime**, and the case
218 $y_0^2 - x_0^2 > 2/\eta$ as the **EoS regime**.

219 See Figure 4 for illustrations of these two regimes. In the subsequent sections, we aim to make the
220 above reasoning rigorous. For example, instead of the approximate dynamics (4), we consider the
221 original GD dynamics and justify the Taylor approximation. Also, in the EoS regime, rather than
222 loosely asserting that $|x_t| \searrow 0$ exponentially fast and hence the dynamics stabilizes “quickly” once
223 $y_t^2 < 2/\eta$, we track precisely how long this convergence takes so that we can bound the gap between
224 the limiting sharpness and the prediction $2/\eta$.

225 2.3 Results

226 **Gradient flow regime.** Our first rigorous result is that when $y_0^2 - x_0^2 = \frac{2-\delta}{\eta}$ for some constant
227 $\delta \in (0, 2)$, then the limiting sharpness of GD with step size η is $y_0^2 - x_0^2 + O(1) = \frac{2-\delta}{\eta} + O(1)$,
228 which is precisely the sharpness attained by the gradient flow up to a controlled error term.

229 In fact, our theorem is slightly more general, as it covers initializations in which δ can scale mildly
230 with η . The precise statement is as follows.

231 **Theorem 1** (gradient flow regime; see §B.2). *Suppose we run GD with step size $\eta > 0$ on the*
232 *objective f , where $f(x, y) := \ell(xy)$, and ℓ satisfies Assumptions (A1) and (A2). Let $(\tilde{x}, \tilde{y}) \in \mathbb{R}^2$*
233 *satisfy $\tilde{y} > \tilde{x} > 0$ with $\tilde{y}^2 - \tilde{x}^2 = 1$. Suppose we initialize GD at $(x_0, y_0) := (\frac{2-\delta}{\eta})^{1/2}(\tilde{x}, \tilde{y})$, where*
234 *$\delta \in (0, 2)$ and $\eta \lesssim \delta^{1/2} \wedge (2 - \delta)$. Then, GD converges to $(0, y_\infty)$ satisfying*

$$\frac{2-\delta}{\eta} - O(2-\delta) - O\left(\frac{\eta}{\min\{\delta, 2-\delta\}}\right) \leq \lambda_{\max}(\nabla^2 f(0, y_\infty)) \leq \frac{2-\delta}{\eta} + O\left(\frac{\eta}{2-\delta}\right),$$

235 where the implied constants depend on \tilde{x}, \tilde{y} , and ℓ , but not on δ, η .

236 The proof of Theorem 1 is based on a two-stage analysis. In the first stage, we use Lemma 6 on
237 the approximate conservation of $y^2 - x^2$ along GD in order to show that GD lands near the y -axis
238 with $y_{t_0}^2 \approx \frac{2-\delta}{\eta}$. In the second stage, we use the assumptions on ℓ in order to control the rate of
239 convergence of $|x_t|$ to 0, which is subsequently used to control the final deviation of y_∞^2 from $\frac{2-\delta}{\eta}$.

240 **EoS regime.** Our next result states that when $y_0^2 - x_0^2 > 2/\eta$, then the limiting sharpness of GD is
241 close to the EoS prediction of $2/\eta$, up to an error term which depends on the exponent β in (A2).

242 **Theorem 2** (EoS; see §B.4). *Suppose we run GD on f with step size $\eta > 0$, where $f(x, y) := \ell(xy)$,*
243 *and ℓ satisfies Assumptions (A1) and (A2). Let $(\tilde{x}, \tilde{y}) \in \mathbb{R}^2$ satisfy $\tilde{y} > \tilde{x} > 0$ with $\tilde{y}^2 - \tilde{x}^2 = 1$.*
244 *Suppose we initialize GD at $(x_0, y_0) := (\frac{2+\delta}{\eta})^{1/2}(\tilde{x}, \tilde{y})$, where $\delta > 0$ is a constant. Also, assume*
245 *that for all $t \geq 1$ such that $y_t^2 > 2/\eta$, we have $x_t \neq 0$. Then, GD converges to $(0, y_\infty)$ satisfying*

$$2/\eta - O(\eta^{1/(\beta-1)}) \leq \lambda_{\max}(\nabla^2 f(0, y_\infty)) \leq 2/\eta,$$

246 where the implied constants depend on $\tilde{x}, \tilde{y}, \delta \wedge 1$, and ℓ , but not on η .

247 Remarks on the assumptions.

- 248 1. **Choice of initialization.** The initialization in our results is such that both y_0 and $y_0 - x_0$ are
249 on the same scale, i.e., $y_0, y_0 - x_0 = \Theta(1/\sqrt{\eta})$. This rules out extreme initializations such as
250 $y_0 \approx x_0$, which are problematic because they lie too close to the invariant line $y = x$. Since our
251 aim in this work is not to explore every edge case, we focus on this setting for simplicity.
- 252 2. **Assumption that $x_t \neq 0$ in Theorem 2.** We imposed the additional assumption that the iterates
253 of GD do not exactly hit the y -axis before crossing $y^2 = 2/\eta$. This is necessary because if $x_t = 0$
254 for some iteration t , then $(x_{t'}, y_{t'}) = (x_t, y_t)$ for all $t' > t$, and hence the limiting sharpness may
255 not be close to $2/\eta$. This assumption holds generically, e.g., if we perturb each iterate of GD with
256 a vanishing amount of noise from a continuous distribution, and we conjecture that for any step
257 size $\eta > 0$, the assumption holds for all but a measure zero set of initializations.

258 When $\beta = +\infty$, which is the case for the Huber loss in Example 1, the limiting sharpness is
 259 $2/\eta + O(1)$. When $\beta = 2$, which is the case for the logistic and square root losses in Example 1, the
 260 limiting sharpness is $2/\eta + O(\eta)$. Numerical experiments show that **our error bound of $O(\eta^{1/(\beta-1)})$**
 261 **is sharp**; see Figure 9 below.

262 We make a few remark about the proof. As we outline the proof in §B.3, it turns out in order to bound
 263 the gap $2/\eta - y_\infty^2$, the proof requires a control of the size $|x_{\mathbf{t}} y_{\mathbf{t}}|$, where \mathbf{t} is the first iteration such
 264 that $y_{\mathbf{t}}^2$ crosses $2/\eta$. However, controlling the size of $|x_{\mathbf{t}} y_{\mathbf{t}}|$ is surprisingly delicate as it requires a
 265 fine-grained understanding of the bouncing phase. The insight that guides the proof is the observation
 266 that during the bouncing phase, the GD iterates lie close to a certain envelope (Figure 12).

267 As a by-product of our analysis, we obtain a rigorous version of the quasi-static principle from
 268 which can more accurately track the sharpness gap and convergence rate (see §B.5). The results
 269 of Theorem 1, Theorem 2, and Theorem 5 are displayed pictorially as Figure 12.

270 3 Understanding the bias evolution of the ReLU network

271 In this section, we use the insights from §2 to answer our main question, namely understanding
 272 the role of a large step size in learning threshold neurons for the ReLU network (2). Based on the
 273 observed dynamics (Figure 2), we can make our question more concrete as follows.

274 **(Refined) Q.** What is the role of a large step size during the “initial phase” of training in which (i)
 275 the bias b rapidly decreases and (ii) the sum of weights $a^- + a^+$ oscillates?

276 3.1 Approximating the initial phase of GD with the “mean model”

277 Deferring details to §C, the GD dynamics for the ReLU net-
 278 work (2) in the initial phase are well-approximated by

$$\text{GD dynamics on } (a^-, a^+, b) \mapsto \ell_{\text{sym}}(d(a^- + a^+)g(b)),$$

279 where $\ell_{\text{sym}}(s) := \frac{1}{2}(\log(1 + \exp(-s)) + \log(1 + \exp(+s)))$
 280 and $g(b) := \mathbb{E}_{z \sim \mathcal{N}(0,1)} \text{ReLU}(z + b)$ is the ‘smoothed’ ReLU. Figure 5: The ‘smoothed’ ReLU $g(b)$
 281 The GD dynamics can be compactly written in terms of the parameter $A_t := d(a_t^- + a_t^+)$.

$$A_{t+1} = A_t - 2d^2 \eta \ell'_{\text{sym}}(A_t g(b_t)) g(b_t), \quad b_{t+1} = b_t - \eta \ell'_{\text{sym}}(A_t g(b_t)) A_t g'(b_t). \quad (5)$$

282 We call these dynamics **the mean model**. Figure 10 shows that the mean model closely captures the
 283 GD dynamics for the ReLU network (2), and we henceforth focus on analyzing the mean model.

284 The main advantage of the representation (5) is that it makes apparent the connection to the single-
 285 neuron example that we studied in §2. More specifically, (5) can be interpreted as the “rescaled”
 286 GD dynamics on the objective $(A, b) \mapsto \ell_{\text{sym}}(Ag(b))$, where the step size for the A -dynamics is
 287 multiplied by $2d^2$. Due to this resemblance, we can apply the techniques from §2.

288 3.2 Two different regimes for the mean model

289 Throughout the section, we use the shorthand $\ell := \ell_{\text{sym}}$, and focus on initializing with $a_0^\pm = \Theta(1/d)$,
 290 $a^- + a^+ \neq 0$, and $b_0 = 0$. This implies $A_0 = \Theta(1)$. We also note the following fact for later use.

291 **Lemma 2** (formula for the smoothed ReLU; see §D.1). *The smoothed ReLU function g can be
 292 expressed in terms of the PDF φ and the CDF Φ of the standard Gaussian distribution as $g(b) =$
 293 $\varphi(b) + b\Phi(b)$. In particular, $g' = \Phi$.*

294 Note also that b_t is monotonically decreasing. This is because $\ell'(A_t g(b_t)) A_t g'(b_t) \geq 0$ since ℓ' is
 295 an odd function and $g(b), g'(b) > 0$ for any $b \in \mathbb{R}$.

296 Following the strategy of §2.2, we begin with the continuous-time dynamics of the mean model:

$$\dot{A} = -2d^2 \ell'(Ag(b)) g(b), \quad \dot{b} = -\ell'(Ag(b)) Ag'(b). \quad (6)$$

297 **Lemma 3** (conserved quantity; see §D.1). *Let $\kappa : \mathbb{R} \rightarrow \mathbb{R}$ be defined as $\kappa(b) := \int_0^b g/g'$. Along the
 298 gradient flow (6), the quantity $\frac{1}{2}A^2 - 2d^2\kappa(b)$ is conserved.*

299 Based on [Lemma 3](#), if we initialize the continuous-time dynamics (6) at $(A_0, 0)$ and if $A_t \rightarrow 0$, then
 300 the limiting value of the bias b_∞^{GF} satisfies $\kappa(b_\infty^{\text{GF}}) = -\frac{1}{4d^2} A_0^2$, which implies that $b_\infty^{\text{GF}} = -\Theta(\frac{1}{d^2})$;
 301 indeed, this holds since $\kappa'(0) = g(0)/g'(0) > 0$, so there exist constants $c_0, c_1 > 0$ such that
 302 $c_0 b \leq \kappa(b) \leq c_1 b$ for all $-1 \leq b \leq 0$. Since the mean model (5) tracks the continuous-time
 303 dynamics (6) until it reaches the b -axis, the mean model initialized at $(A_0, 0)$ also approximately
 304 reaches $(A_{t_0}, b_{t_0}) \approx (0, -\Theta(\frac{1}{d^2})) \approx (0, 0)$ in high dimension $d \gg 1$. In other words, **the continuous-**
 305 **time dynamics (6) fails to learn threshold neurons.**

306 Once the mean model reaches the b -axis, we again identify two different regimes depending on the
 307 step size. A Taylor expansion of ℓ' around the origin yields the following approximate dynamics
 308 (here $\ell''(0) = 1/4$): $A_{t_0+1} \approx A_{t_0} - \frac{\eta d^2}{2} A_{t_0} g(b_{t_0})^2 = A_{t_0} (1 - \frac{\eta d^2}{2} g(b_{t_0})^2)$. We conclude that the
 309 condition which now dictates whether we have bouncing or convergence is $\frac{1}{2} d^2 g(b_{t_0})^2 > 2/\eta$.

- 310 (i) **Gradient flow regime:** If $2/\eta > d^2 g(0)^2/2 = d^2/(4\pi)$ (since $g(0)^2 = 1/(2\pi)$), i.e., the
 311 step size η is *below* the threshold $8\pi/d^2$, then the final bias of the mean model b_∞^{MM} satisfies
 312 $b_\infty^{\text{MM}} \approx b_\infty^{\text{GF}} \approx 0$. In other words, **the mean model fails to learn threshold neurons.**
 313 (ii) **EoS regime:** If $2/\eta < d^2/(4\pi)$, i.e., the step size η is *above* the threshold $8\pi/d^2$, then
 314 $\frac{1}{2} d^2 g^2(b_\infty^{\text{MM}}) < 2/\eta$, i.e., $b_\infty^{\text{MM}} < g^{-1}(2/\sqrt{\eta d^2})$. For instance, if $\eta = \frac{10\pi}{d^2}$, then $b_\infty^{\text{MM}} < -0.087$.
 315 In other words, **the mean model successfully learns threshold neurons.**

316 3.3 Results for the mean model

317 **Theorem 3** (mean model, gradient flow regime; see §D). *Consider the mean model (5) initialized at*
 318 *$(A_0, 0)$, with step size $\eta = \frac{(8-\delta)\pi}{d^2}$ for some $\delta > 0$. Let $\gamma := \frac{1}{200} \min\{\delta, 8-\delta, \frac{8-\delta}{|A_0|}\}$. Then, as long*
 319 *as $\eta \leq \gamma/|A_0|$, the limiting bias b_∞^{MM} satisfies*

$$0 \geq b_\infty^{\text{MM}} \geq -(\eta/\gamma) |A_0| = -O_{A_0, \delta}(1/d^2).$$

320 **Theorem 4** (mean model, EoS regime; see §D). *Consider the mean model initialized at $(A_0, 0)$,*
 321 *with step size $\eta = \frac{(8+\delta)\pi}{d^2}$ for some $\delta > 0$. Furthermore, assume that for all $t \geq 1$ such that*
 322 *$\frac{1}{2} d^2 g(b_t)^2 > 2/\eta$, we have $A_t \neq 0$. Then, the limiting bias b_∞^{MM} satisfies*

$$b_\infty^{\text{MM}} \leq g^{-1}(2/\sqrt{(8+\delta)\pi}) \leq -\Omega_\delta(1).$$

323 4 Conclusion

324 In this paper, we present the first explanation for the emergence of threshold neuron (i.e., ReLU
 325 neurons with negative bias) in models such as the sparse coding model (1) through a novel connection
 326 with the ‘‘edge of stability’’ (EoS) phenomenon. Along the way, we obtain a detailed and rigorous
 327 understanding of the dynamics of GD in the EoS regime for a simple class of loss functions, thereby
 328 shedding light on the impact of large learning rates in non-convex optimization.

329 Many interesting questions remain, and we conclude with some directions for future research.

- 330 • **Extending the analysis of EoS to richer models.** Although the analysis we present in this work
 331 is restricted to simple models, the underlying principles can potentially be applied to more general
 332 settings. In this direction, it would be interesting to study models which capture the impact of the
 333 depth of the neural network on the EoS phenomenon.
- 334 • **The interplay between the EoS and the choice of optimization algorithm.** As discussed in §2.3,
 335 the bouncing phase of the EoS substantially slows down the convergence of GD (see [Figure 9](#)).
 336 Investigating how different optimization algorithm (e.g., SGD, or GD with momentum) interact
 337 with the EoS phenomenon could potentially lead to practical speed-ups or improved generalization.
- 338 • **An end-to-end analysis of the sparse coding model.** Finally, we have left open the motivating
 339 question of analyzing how two-layer ReLU networks learn to solve the sparse coding model (1).
 340 Despite the apparent simplicity of the problem, its analysis has thus far remained out of reach, and
 341 we believe that a resolution to this question would constitute compelling and substantial progress
 342 towards understanding neural network learning. We are hopeful that the insights in this paper
 343 provide the first step towards this goal.

344 References

- 345 [AL22] Z. Allen-Zhu and Y. Li. “Feature purification: how adversarial training performs robust
346 deep learning”. In: *2021 IEEE 62nd Annual Symposium on Foundations of Computer
347 Science (FOCS)*. IEEE. 2022, pp. 977–988.
- 348 [ALP22] S. Arora, Z. Li, and A. Panigrahi. “Understanding gradient descent on the edge of
349 stability in deep learning”. In: *Proceedings of the 39th International Conference on
350 Machine Learning*. Ed. by K. Chaudhuri et al. Vol. 162. Proceedings of Machine
351 Learning Research. PMLR, July 2022, pp. 948–1024.
- 352 [ALS19] Z. Allen-Zhu, Y. Li, and Z. Song. “A convergence theory for deep learning via over-
353 parameterization”. In: *Proceedings of the 36th International Conference on Machine
354 Learning*. Ed. by K. Chaudhuri and R. Salakhutdinov. Vol. 97. Proceedings of Machine
355 Learning Research. PMLR, June 2019, pp. 242–252.
- 356 [And+22] M. Andriushchenko, A. Varre, L. Pillaud-Vivien, and N. Flammarion. “SGD with large
357 step sizes learns sparse features”. In: *arXiv preprint arXiv:2210.05337* (2022).
- 358 [Aro+15] S. Arora, R. Ge, T. Ma, and A. Moitra. “Simple, efficient, and neural algorithms for
359 sparse coding”. In: *Proceedings of the 28th Conference on Learning Theory*. Ed. by P.
360 Grünwald, E. Hazan, and S. Kale. Vol. 40. Proceedings of Machine Learning Research.
361 Paris, France: PMLR, July 2015, pp. 113–149.
- 362 [Aro+19] S. Arora, S. Du, W. Hu, Z. Li, and R. Wang. “Fine-grained analysis of optimization
363 and generalization for overparameterized two-layer neural networks”. In: *Proceedings
364 of the 36th International Conference on Machine Learning*. Ed. by K. Chaudhuri and
365 R. Salakhutdinov. Vol. 97. Proceedings of Machine Learning Research. PMLR, June
366 2019, pp. 322–332.
- 367 [AZS22] K. Ahn, J. Zhang, and S. Sra. “Understanding the unstable convergence of gradient
368 descent”. In: *Proceedings of the 39th International Conference on Machine Learning*.
369 Ed. by K. Chaudhuri et al. Vol. 162. Proceedings of Machine Learning Research. PMLR,
370 July 2022, pp. 247–257.
- 371 [Ba+22] J. Ba et al. “High-dimensional asymptotics of feature learning: how one gradient step
372 improves the representation”. In: *Advances in Neural Information Processing Systems*.
373 Ed. by A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho. 2022.
- 374 [BMR22] G. Beugnot, J. Mairal, and A. Rudi. “On the benefits of large learning rates for kernel
375 methods”. In: *Proceedings of Thirty Fifth Conference on Learning Theory*. Ed. by P.-L.
376 Loh and M. Raginsky. Vol. 178. Proceedings of Machine Learning Research. PMLR,
377 July 2022, pp. 254–282.
- 378 [CB18] L. Chizat and F. Bach. “On the global convergence of gradient descent for over-
379 parameterized models using optimal transport”. In: *Advances in Neural Information
380 Processing Systems*. Ed. by S. Bengio et al. Vol. 31. Curran Associates, Inc., 2018.
- 381 [CB22] L. Chen and J. Bruna. “On gradient descent convergence beyond the edge of stability”.
382 In: *arXiv preprint arXiv:2206.04172* (2022).
- 383 [Chi22] L. Chizat. “Mean-field Langevin dynamics: exponential convergence and annealing”.
384 In: *Transactions on Machine Learning Research* (2022).
- 385 [COB19] L. Chizat, E. Oyallon, and F. Bach. “On lazy training in differentiable programming”. In:
386 *Advances in Neural Information Processing Systems*. Ed. by H. Wallach et al. Vol. 32.
387 Curran Associates, Inc., 2019.
- 388 [Coh+21] J. Cohen, S. Kaur, Y. Li, J. Z. Kolter, and A. Talwalkar. “Gradient descent on neural
389 networks typically occurs at the edge of stability”. In: *International Conference on
390 Learning Representations*. 2021.
- 391 [DNL23] A. Damian, E. Nichani, and J. D. Lee. “Self-stabilization: the implicit bias of gradient
392 descent at the edge of stability”. In: *The Eleventh International Conference on Learning
393 Representations*. 2023.
- 394 [Du+19] S. S. Du, X. Zhai, B. Póczos, and A. Singh. “Gradient descent provably optimizes
395 over-parameterized neural networks”. In: *International Conference on Learning Repre-
396 sentations*. 2019.

- 397 [Jas+18] S. Jastrzebski et al. “Width of minima reached by stochastic gradient descent is influ-
398 enced by learning rate to batch size ratio”. In: *Artificial Neural Networks and Machine*
399 *Learning – ICANN 2018*. Ed. by V. Kůrková, Y. Manolopoulos, B. Hammer, L. Iliadis,
400 and I. Maglogiannis. Cham: Springer International Publishing, 2018, pp. 392–402.
- 401 [Jas+19] S. Jastrzebski et al. “On the relation between the sharpest directions of DNN loss and
402 the SGD step length”. In: *International Conference on Learning Representations*. 2019.
- 403 [Jas+20] S. Jastrzebski et al. “The break-even point on optimization trajectories of deep neural
404 networks”. In: *International Conference on Learning Representations*. 2020.
- 405 [Jas+21] S. Jastrzebski et al. “Catastrophic Fisher explosion: early phase Fisher matrix impacts
406 generalization”. In: *International Conference on Machine Learning*. PMLR. 2021,
407 pp. 4772–4784.
- 408 [JGH18] A. Jacot, F. Gabriel, and C. Hongler. “Neural tangent kernel: convergence and general-
409 ization in neural networks”. In: *Proceedings of the 32nd Advances in Neural Information*
410 *Processing Systems*. 2018, pp. 8580–8589.
- 411 [Kar+21] S. Karp, E. Winston, Y. Li, and A. Singh. “Local signal adaptivity: provable feature
412 learning in neural networks beyond kernels”. In: *Advances in Neural Information*
413 *Processing Systems*. Ed. by M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and
414 J. W. Vaughan. Vol. 34. Curran Associates, Inc., 2021, pp. 24883–24897.
- 415 [KR18] F. Koehler and A. Risteski. “The comparative power of ReLU networks and polynomial
416 kernels in the presence of sparse latent structure”. In: *International Conference on*
417 *Learning Representations*. 2018.
- 418 [Lew+20] A. Lewkowycz, Y. Bahri, E. Dyer, J. Sohl-Dickstein, and G. Gur-Ari. “The large
419 learning rate phase of deep learning: the catapult mechanism”. In: *arXiv preprint*
420 *arXiv:2003.02218* (2020).
- 421 [LLA22] K. Lyu, Z. Li, and S. Arora. “Understanding the generalization benefit of normalization
422 layers: sharpness reduction”. In: *Advances in Neural Information Processing Systems*.
423 Ed. by S. Koyejo et al. Vol. 35. Curran Associates, Inc., 2022, pp. 34689–34708.
- 424 [LWM19] Y. Li, C. Wei, and T. Ma. “Towards explaining the regularization effect of initial
425 large learning rate in training neural networks”. In: *Advances in Neural Information*
426 *Processing Systems* 32 (2019).
- 427 [Ma+22] C. Ma, D. Kunin, L. Wu, and L. Ying. “Beyond the quadratic approximation: the
428 multiscale structure of neural network loss landscapes”. In: *Journal of Machine Learning*
429 1.3 (2022), pp. 247–267.
- 430 [MMM19] S. Mei, T. Misiakiewicz, and A. Montanari. “Mean-field theory of two-layers neural
431 networks: dimension-free bounds and kernel limit”. In: *Proceedings of the Thirty-*
432 *Second Conference on Learning Theory*. Ed. by A. Beygelzimer and D. Hsu. Vol. 99.
433 Proceedings of Machine Learning Research. PMLR, June 2019, pp. 2388–2464.
- 434 [MMS21] R. Mulayoff, T. Michaeli, and D. Soudry. “The implicit bias of minima stability: a
435 view from function space”. In: *Advances in Neural Information Processing Systems* 34
436 (2021), pp. 17749–17761.
- 437 [Nac+22] M. S. Nacson, K. Ravichandran, N. Srebro, and D. Soudry. “Implicit bias of the step size
438 in linear diagonal neural networks”. In: *International Conference on Machine Learning*.
439 PMLR. 2022, pp. 16270–16295.
- 440 [Nak20] P. Nakkiran. “Learning rate annealing can provably help generalization, even for convex
441 problems”. In: *arXiv preprint arXiv:2005.07360* (2020).
- 442 [NWS22] A. Nitanda, D. Wu, and T. Suzuki. “Convex analysis of the mean field Langevin dynam-
443 ics”. In: *Proceedings of the 25th International Conference on Artificial Intelligence and*
444 *Statistics*. Ed. by G. Camps-Valls, F. J. R. Ruiz, and I. Valera. Vol. 151. Proceedings of
445 Machine Learning Research. PMLR, Mar. 2022, pp. 9741–9757.
- 446 [OF04] B. A. Olshausen and D. J. Field. “Sparse coding of sensory inputs”. In: *Current Opinion*
447 *in Neurobiology* 14.4 (2004), pp. 481–487.
- 448 [OF97] B. A. Olshausen and D. J. Field. “Sparse coding with an overcomplete basis set: a
449 strategy employed by V1?” In: *Vision Research* 37.23 (1997), pp. 3311–3325.
- 450 [OS20] S. Oymak and M. Soltanolkotabi. “Toward moderate overparameterization: global
451 convergence guarantees for training shallow neural networks”. In: *IEEE Journal on*
452 *Selected Areas in Information Theory* 1.1 (2020), pp. 84–105.

- 453 [RV22] G. M. Rotskoff and E. Vanden-Eijnden. “Trainability and accuracy of artificial neural
454 networks: an interacting particle system approach”. In: *Comm. Pure Appl. Math.* 75.9
455 (2022), pp. 1889–1935.
- 456 [VG00] W. E. Vinje and J. L. Gallant. “Sparse coding and decorrelation in primary visual cortex
457 during natural vision”. In: *Science* 287.5456 (2000), pp. 1273–1276.
- 458 [Wan+22] Y. Wang, M. Chen, T. Zhao, and M. Tao. “Large learning rate tames homogeneity:
459 convergence and balancing effect”. In: *International Conference on Learning Representations*. 2022.
460
- 461 [WLL22] Z. Wang, Z. Li, and J. Li. “Analyzing sharpness along GD trajectory: progressive
462 sharpening and edge of stability”. In: *Advances in Neural Information Processing
463 Systems*. Ed. by S. Koyejo et al. Vol. 35. Curran Associates, Inc., 2022, pp. 9983–9994.
- 464 [WME18] L. Wu, C. Ma, and W. E. “How SGD selects the global minima in over-parameterized
465 learning: a dynamical stability perspective”. In: *Advances in Neural Information Pro-
466 cessing Systems* 31 (2018), pp. 8279–8288.
- 467 [Wu+21] J. Wu, D. Zou, V. Braverman, and Q. Gu. “Direction matters: on the implicit bias of
468 stochastic gradient descent with moderate learning rate”. In: *International Conference
469 on Learning Representations*. 2021.
- 470 [Xin+18] C. Xing, D. Arpit, C. Tsirigotis, and Y. Bengio. “A walk with SGD”. In: *arXiv preprint
471 arXiv:1802.08770* (2018).
- 472 [Yan+09] J. Yang, K. Yu, Y. Gong, and T. Huang. “Linear spatial pyramid matching using sparse
473 coding for image classification”. In: *2009 IEEE Conference on Computer Vision and
474 Pattern Recognition*. IEEE. 2009, pp. 1794–1801.
- 475 [Zhu+23] X. Zhu, Z. Wang, X. Wang, M. Zhou, and R. Ge. “Understanding edge-of-stability train-
476 ing dynamics with a minimalist example”. In: *The Eleventh International Conference
477 on Learning Representations*. 2023.

Appendix

478

479	A Further experimental results	13
480	A.1 Experiments for the full sparse coding model	13
481	A.2 Experiments on the CIFAR-10 dataset	14
482	B Proofs for the single-neuron linear network	15
483	B.1 Approximate conservation along GD	16
484	B.2 Gradient flow regime	17
485	B.3 EoS regime: proof outline	18
486	B.4 EoS regime: crossing the threshold and the convergence phase	18
487	B.5 EoS regime: quasi-static analysis	21
488	C Deferred derivations of mean model	25
489	D Proofs for the mean model	26
490	D.1 Deferred proofs	27
491	D.2 Gradient flow regime	27
492	D.3 EoS regime	28
493	E Additional figures	28

494 **A Further experimental results**

495 In this section, we report further experimental results which demonstrate that our theory, while limited
 496 to the specific models we study (namely, the single-neuron example and the mean model), is in fact
 497 indicative of behaviors commonly observed in more realistic instances of neural network training. In
 498 particular, we show that threshold neurons often emerge in the presence of oscillations in the other
 499 weight parameters of the network.

500 **A.1 Experiments for the full sparse coding model**

501 We provide the details for the top plot of [Figure 3](#). consider the sparse coding model in the form (1).
 502 Compared to (2), we assume that the basis vectors are unknown, and the neural network learn them
 503 through additional parameters $\mathbf{W} = (\mathbf{w}_i)_{i=1}^m$ together with m different weights $\mathbf{a} = (a_i)_{i=1}^m$ for the
 504 second layer as follows:

$$f(\mathbf{x}; \mathbf{a}, \mathbf{W}, b) = \sum_{i=1}^m a_i \text{ReLU}(\langle \mathbf{w}_i, \mathbf{x} \rangle + b). \quad (7)$$

505 We show results for $d = 100$, $m = 2000$. We generate $n = 20000$ data points according to the
 506 aforementioned sparse coding model with $\lambda = 5$. We use the He initialization, i.e., $\mathbf{a} \sim \mathcal{N}(0, I_m/m)$,
 507 $\mathbf{w} \sim \mathcal{N}(0, I_d/d)$, and $b = 0$. As shown in the top plot of [Figure 3](#), the bias decreases more with the
 508 large learning rate. Further, we report the behavior of the average of second layer weights in [Figure 6](#),
 509 and confirm that the sum oscillates.

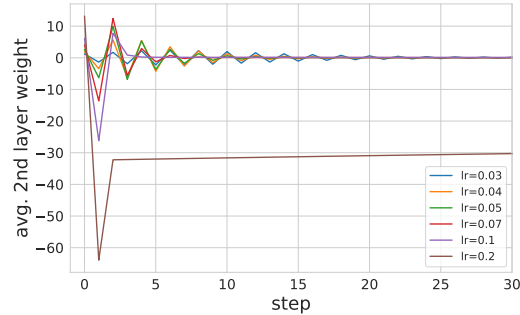


Figure 6: The average of the second layer weights of the ReLU network (7). Note that the average value oscillates similarly to our findings for the mean model.

510 A.2 Experiments on the CIFAR-10 dataset

511 Next, we provide the details for the bottom plot of Figure 3. We train ResNet-18 on a binarized
 512 version of the CIFAR-10 dataset formed by taking only the first two classes; this is done for the
 513 purpose of monitoring the average logit of the network. The average logit is measured over the entire
 514 training set. The median bias is measured at the last convolutional layer right before the pooling. For
 515 the optimizer, we use full-batch GD with no momentum or weight decay, plus a cosine learning rate
 516 scheduler where learning rates shown in the plots are the initial values.

517 **Oscillation of expected output (logit) of the network.** Bearing a striking resemblance to our
 518 two-layer models, here the expected mean of the output (logit) of the deep net also oscillates due to
 519 GD dynamics. As we have argued in the previous sections, this occurs as the bias parameters are
 520 driven towards negative values.

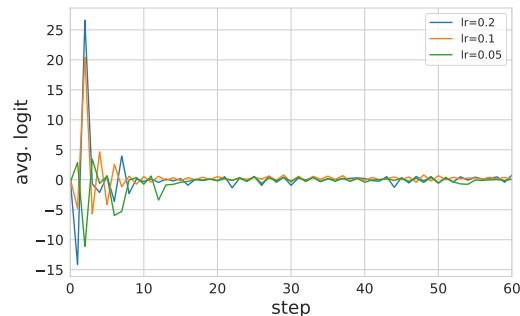


Figure 7: Oscillation of logit of ResNet18 model averaged over the (binary) CIFAR-10 training set. Since the dataset is binary, the logit is simply a scalar.

521 **Results for SGD.** In Figure 8, we report qualitatively similar phenomena when we instead train
 522 ResNet-18 with stochastic gradient descent (SGD), where we use all ten classes of CIFAR-10. Again,
 523 the median bias is measured at the last convolutional layer. We further report the average activation
 524 which is the output of the ReLU activation at the last convolutional layer, averaged over the neurons
 525 and the entire training set. The average activation statistics represent the hidden representations
 526 before the linear classifier part, and lower values represent sparser representations. Interestingly, the
 527 threshold neuron also emerges with larger step sizes similarly to the case of gradient descent.

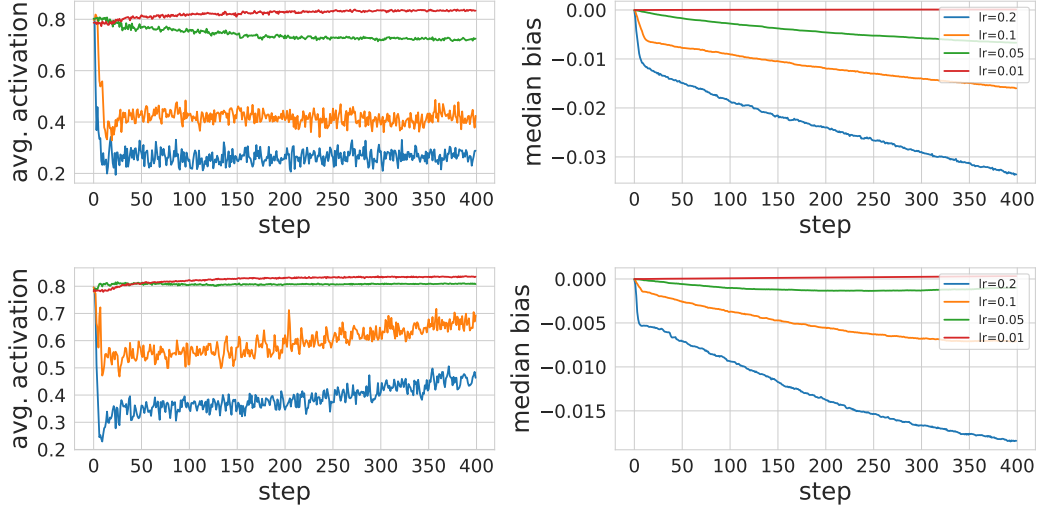


Figure 8: SGD dynamics of ResNet-18 on (multiclass) CIFAR-10 with various learning rates and batch sizes. (Top) batch size 100; (Bottom) batch size 1000. The results are consistent across different batch sizes.

528 B Proofs for the single-neuron linear network

529 We start by describing basic and relevant properties of the model and the assumptions on ℓ .

530 **Basic properties.** If ℓ is minimized at 0, then the *global minimizers* of f are the x - and y -axes. The
 531 *gradient and Hessian* of f are given by:

$$\begin{aligned} \nabla f(x, y) &= \ell'(xy) \begin{bmatrix} y \\ x \end{bmatrix}, \\ \nabla^2 f(x, y) &= \ell''(xy) \begin{bmatrix} y \\ x \end{bmatrix}^{\otimes 2} + \ell'(xy) \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}. \end{aligned} \quad (8)$$

532 This results in GD iterates x_t, y_t , for step size $\eta > 0$ and iteration $t \geq 0$:

$$\begin{aligned} x_{t+1} &= x_t - \eta \ell'(x_t y_t) y_t, \\ y_{t+1} &= y_t - \eta \ell'(x_t y_t) x_t. \end{aligned}$$

533 **Lemma 4** (invariant lines). Assume that ℓ is even, so that ℓ' is odd. Then, the lines $y = \pm x$ are
 534 invariant for gradient descent on f .

535 *Proof.* If $y_t = \pm x_t$, then

$$\begin{aligned} y_{t+1} &= y_t - \eta \ell'(x_t y_t) x_t = \pm x_t \mp \eta \ell'(x_t^2) x_t, \\ x_{t+1} &= x_t - \eta \ell'(x_t y_t) y_t = x_t - \eta \ell'(x_t^2) x_t, \end{aligned}$$

536 and hence $y_{t+1} = \pm x_{t+1}$. Note that the iterates $(x_t)_{t \geq 0}$ are the iterates of GD with step size η on the
 537 one-dimensional loss function $x \mapsto \frac{1}{2} \ell(x^2)$. \square

538 We focus instead on initializing away from these two lines. We now state our assumptions on ℓ .

539 We gather together some elementary properties of ℓ .

540 **Lemma 5** (properties of ℓ). Suppose that Assumption (A1) holds.

- 541 1. ℓ is minimized at the origin and $\ell'(0) = 0$.
- 542 2. Suppose that ℓ is four times continuously differentiable near the origin. If Assumption (A2) holds,
 543 then $\ell^{(4)}(0) \leq 0$. Conversely, if $\ell^{(4)}(0) < 0$, then Assumption (A2) holds for $\beta = 2$.

544 *Proof.* The first statement is straightforward. The second statement follows from Taylor expansion:
 545 for $s \neq 0$ near the origin,

$$\frac{\ell'(s)}{s} = \frac{\ell'(0) + \ell''(0)s + \int_0^s (s-r)\ell'''(r)dr}{s} = 1 + \int_0^s \left(1 - \frac{r}{s}\right)\ell'''(r)dr. \quad (9)$$

546 Since ℓ''' is odd, then Assumption (A2) and (9) imply that ℓ''' is non-positive on $(0, \varepsilon)$ for some
 547 $\varepsilon > 0$, which in turn implies $\ell^{(4)}(0) \leq 0$. Conversely, if $\ell^{(4)}(0) < 0$, then there exists $\varepsilon > 0$ such
 548 that $\ell'''(s) \leq -\varepsilon s$ for $s \in (0, \varepsilon)$. From (9), we see that $\ell'(s)/s \leq 1 - \varepsilon \int_0^s (s-r)dr \leq 1 - \varepsilon s^2/2$.
 549 By symmetry, we conclude that Assumption (A2) holds with $\beta = 2$ and some $c > 0$. \square

550 We give some simple examples of losses satisfying our assumptions.

551 **Example 1.** The examples below showcase several functions ℓ that satisfy Assumptions (A1) and (A2)
 552 with different values of β .

- 553 • *Rescaled and symmetrized logistic loss.* $\ell_{\text{rsym}}(s) := \frac{1}{2} \ell_{\text{logi}}(-2s) + \frac{1}{2} \ell_{\text{logi}}(+2s)$.
 554 Note $\ell'_{\text{rsym}}(s) = \tanh(s)$, thus $\ell'_{\text{rsym}}(s)/s \leq 1$ and $\ell'_{\text{rsym}}(s)/s \leq 1 - \frac{1}{4}|s|^2$, for $|s| < \frac{1}{4}$.
- 555 • *Square root loss.* $\ell_{\text{sqrt}}(s) := \sqrt{1+s^2}$.
 556 Note $\ell'_{\text{sqrt}}(s) = \frac{s}{\sqrt{1+s^2}}$, thus $\ell'_{\text{sqrt}}(s)/s \leq 1$ and $\ell'_{\text{sqrt}}(s)/s \leq 1 - \frac{2}{5}|s|^2$, for $|s| < \frac{2}{5}$.
- 557 • *Huber loss.* $\ell_{\text{Hub}}(s) := \frac{s^2}{2} \mathbb{1}\{s \in [-1, 1]\} + (|s| - \frac{1}{2}) \mathbb{1}\{s \notin [-1, 1]\}$.
 558 Note $\ell'_{\text{Hub}}(s) = s \mathbb{1}\{s \in [-1, 1]\} + \text{sgn}(s) \mathbb{1}\{s \notin [-1, 1]\}$, thus $\ell'_{\text{Hub}}(s)/s \leq 1$, i.e., we have
 559 Assumption (A2) with $\beta = +\infty$.
- *Higher-order.* For $\beta > 1$ let $c_\beta := \frac{1}{\beta+1} \left(\frac{\beta}{\beta+1}\right)^\beta$ and $r_\beta := \frac{\beta+1}{\beta}$. We define ℓ_β implicitly via its
 derivative

$$\ell'_\beta(s) := s(1 - c_\beta |s|^\beta) \mathbb{1}\{s^2 < r_\beta^2\} + \text{sgn}(s) \mathbb{1}\{s^2 \geq r_\beta^2\}.$$

560 By definition, $\ell'_\beta(s)/s \leq 1$ and $\ell'_\beta(s)/s \leq 1 - c_\ell |s|^\beta$, where $c_\ell = c_\beta \wedge r_\beta$.

561 We now prove our main results from §2.3 in order.

562 B.1 Approximate conservation along GD

563 We begin by stating and proving the approximate conservation of $y^2 - x^2$ for the GD dynamics.

Lemma 6 (approximately conserved quantity). *Let $(\tilde{x}, \tilde{y}) \in \mathbb{R}^2$ be such that $\tilde{y} > \tilde{x} > 0$ with
 $\tilde{y}^2 - \tilde{x}^2 = 1$. Suppose that we run GD on f with step size η with initial point $(x_0, y_0) := \sqrt{\frac{\tilde{y}}{\eta}}(\tilde{x}, \tilde{y})$,
 for some $\gamma > 0$. Then, there exists $t_0 = O(\frac{1}{\eta})$ such that $\sup_{t \geq t_0} |x_t| \leq O(\sqrt{(\gamma^{-1} \vee \gamma)\eta})$ and*

$$y_{t_0}^2 - x_{t_0}^2 = (1 - O(\eta))(y_0^2 - x_0^2),$$

564 where the implied constant depends on \tilde{x}, \tilde{y} , and ℓ .

565 *Proof.* Let $D_t := y_t^2 - x_t^2$ and note that

$$\begin{aligned} D_{t+1} &= (y_t - \eta \ell'(x_t y_t) x_t)^2 - (x_t - \eta \ell'(x_t y_t) y_t)^2 \\ &= (1 - \eta^2 \ell'(x_t y_t)^2) D_t. \end{aligned}$$

566 Since ℓ is 1-Lipschitz, then $D_{t+1} = (1 - O(\eta^2)) D_t$.

567 This shows that for $t \lesssim 1/\eta^2$, we have $y_t^2 - x_t^2 = D_t \gtrsim D_0 = y_0^2 - x_0^2 \asymp \gamma/\eta$. Since $\ell''(0) = 1$,
 568 there exist constants $c_0, c_1 > 0$ such that $\ell'(|xy|) \geq \ell'(c_0) \geq c_1$ whenever $|xy| \geq c_0$. Hence, for all
 569 $t \geq 1$ such that $t \lesssim 1/\eta^2$, $x_t > 0$, and $|x_t y_t| \geq c_0$, we have $y_t^2 \gtrsim \gamma/\eta$ and

$$x_{t+1} = x_t - \eta \ell'(x_t y_t) y_t = x_t - \Theta(\eta y_t) = x_t - \Theta(\sqrt{\gamma \eta}). \quad (10)$$

570 Since $x_0 \asymp \sqrt{\gamma/\eta}$, this shows that after at most $O(1/\eta)$ iterations, we must have either $x_t < 0$ or
 571 $|x_t y_t| \leq c_0$ for the first time. In the first case, (10) shows that $|x_t| \lesssim \sqrt{\gamma \eta}$. In the second case, since
 572 $y_t^2 \gtrsim \gamma/\eta$, we have $|x_t| \lesssim \sqrt{\eta/\gamma}$. Let t_0 denote the iteration at which this occurs.

573 Next, for iterations $t \geq t_0$, we use the dynamics (10) for x and the fact that $\ell'(x_t y_t)$ has the same sign
 574 as x_t to conclude that there are two possibilities: either x_{t+1} has the same sign as x_t , in which case
 575 $|x_{t+1}| \leq |x_t|$, or x_{t+1} has the opposite sign as x_t , in which case $|x_{t+1}| \leq \eta |\ell'(x_t y_t)| y_t \leq \eta y_t \leq$
 576 $O(\sqrt{\gamma \eta})$. This implies $\sup_{t \geq t_0} |x_t| \leq O(\sqrt{(\gamma^{-1} \vee \gamma)\eta})$ as asserted. \square

577 **B.2 Gradient flow regime**

578 In this section, we prove [Theorem 1](#).

579 *Proof of Theorem 1.* From [Lemma 6](#), there exists an iteration t_0 such that $|x_{t_0}| \lesssim \sqrt{\eta/(2-\delta)}$ and

$$\frac{2-\delta}{\eta} - O(2-\delta) \leq y_{t_0}^2 \leq \frac{2-\delta}{\eta} + x_{t_0}^2 \leq \frac{2-\delta}{\eta} + O\left(\frac{\eta}{2-\delta}\right).$$

580 In particular, $C := |x_{t_0} y_{t_0}| \lesssim 1$.

581 We prove by induction the following facts: for $t \geq t_0$,

- 582 1. $|x_t y_t| \leq C$.
 583 2. $|x_t| \leq |x_{t_0}| \exp(-\Omega(\alpha(t-t_0)))$, where $\alpha := \min\{\delta, 2-\delta\}$.

584 Suppose that these conditions hold up to iteration $t \geq t_0$. By [Assumption \(A2\)](#), we have $|\ell'(s)| \leq |s|$
 585 for all $s \neq 0$. Therefore,

$$\begin{aligned} y_{t+1} &= y_t - \eta \ell'(x_t y_t) x_t \geq (1 - \eta x_t^2) y_t \\ &\geq \exp\left(-O\left(\frac{\eta^2}{2-\delta}\right) \exp(-\Omega(\alpha(t-t_0)))\right) y_t \\ &\geq \exp\left(-O\left(\frac{\eta^2}{2-\delta}\right) \sum_{s=t_0}^t \exp(-\Omega(\alpha(s-t_0)))\right) y_{t_0} \geq \exp\left(-O\left(\frac{\eta^2}{\alpha(2-\delta)}\right)\right) y_{t_0}, \\ y_{t+1}^2 &\geq \frac{2-\delta}{\eta} - O(2-\delta) - O\left(\frac{\eta}{\alpha}\right). \end{aligned} \quad (11)$$

586 In particular, $\frac{1}{2} \frac{2-\delta}{\eta} \leq y_t^2 \leq \frac{2-\delta/2}{\eta}$ throughout. In order for these assertions to hold, we require
 587 $\eta^2 \lesssim \alpha(2-\delta)$, i.e., $\eta \lesssim \min\{\sqrt{\delta}, 2-\delta\}$.

588 Next, we would like to show that $t \mapsto |x_t|$ is decaying exponentially fast. Since

$$|x_{t+1}| = |x_t - \eta \ell'(x_t y_t) y_t| = \left| |x_t| - \eta \ell'(|x_t| y_t) y_t \right|,$$

589 it suffices to consider the case when $x_t > 0$. [Assumption \(A2\)](#) implies that

$$x_{t+1} \geq (1 - \eta y_t^2) x_t \geq -\left(1 - \frac{\delta}{2}\right) x_t.$$

590 For the upper bound, we split into two cases. We begin by observing that since ℓ is twice continuously
 591 differentiable near the origin with $\ell''(0) = 1$, there is a constant ε_0 such that $|s| < \varepsilon_0$ implies
 592 $|\ell'(s)| \geq \frac{1}{2}|s|$. If $s_t := x_t y_t \leq \varepsilon_0$, then

$$x_{t+1} \leq \left(1 - \frac{\eta}{2} y_t^2\right) x_t \leq \left(1 - \frac{2-\delta}{4}\right) x_t.$$

593 Otherwise, if $s_t \geq \varepsilon_0$, then

$$x_{t+1} \leq x_t - \eta \ell'(\varepsilon_0) y_t \leq x_t - \eta \ell'(\varepsilon_0) \frac{y_t^2}{s_t} \leq x_t - \eta \ell'(\varepsilon_0) \frac{2-\delta}{2C\eta} x_t \leq (1 - \Omega(2-\delta)) x_t.$$

594 Combining these inequalities, we obtain

$$|x_{t+1}| \leq |x_t| \exp(-\Omega(\alpha)).$$

595 This verifies the second statement in the induction. The first statement follows because both $t \mapsto |x_t|$
 596 and $t \mapsto y_t$ are decreasing.

597 This shows in particular that $|x_t| \searrow 0$, i.e., we have global convergence. To conclude the proof,
 598 observe that [\(11\)](#) gives a bound on the final sharpness. \square

599 *Remark 1.* The proof also gives us estimates on the convergence rate. Namely, from [Lemma 6](#), the
 600 initial phase in which we approach the y -axis takes $O(\frac{1}{\eta})$ iterations. For the convergence phase, in
 601 order to achieve ε error, we need $|x_t| \lesssim \frac{\sqrt{\varepsilon\eta}}{\sqrt{2-\delta}}$; hence, the convergence phase needs only $O(\frac{1}{\alpha} \log \frac{1}{\varepsilon})$
 602 iterations. Note that the rate of convergence in the latter phase does not depend on the step size η .

603 **B.3 EoS regime: proof outline**

604 We give a brief outline of the proof of [Theorem 2](#): As before, [Lemma 6](#) shows that GD reaches the
 605 y -axis approximately at $(0, \sqrt{y_0^2 - x_0^2})$. At this point, x starts bouncing while y steadily decreases,
 606 and we argue that unless $x_t = 0$ or $y_t^2 \leq 2/\eta$, the GD dynamics cannot stabilize (see [Lemma 7](#)).

607 To bound the gap $2/\eta - y_\infty^2$, we look at the first iteration \mathbf{t} such that $y_{\mathbf{t}}^2$ crosses $2/\eta$. By making use
 608 of [Assumption \(A2\)](#), we simultaneously control both the convergence rate of $|x_t|$ to zero and the
 609 decrease in y_t^2 in order to prove that

$$y_\infty^2 \geq \frac{2}{\eta} - O(|x_{\mathbf{t}}y_{\mathbf{t}}|), \quad (12)$$

610 see [Proposition 1](#). Therefore, to establish [Theorem 2](#), we must bound $|x_{\mathbf{t}}y_{\mathbf{t}}|$ at iteration \mathbf{t} .

611 Controlling the size of $|x_{\mathbf{t}}y_{\mathbf{t}}|$, however, is surprisingly delicate as it requires a fine-grained under-
 612 standing of the bouncing phase. The insight that guides the proof is the observation that during the
 613 bouncing phase, the GD iterates lie close to a certain envelope ([Figure 12](#)). This envelope is predicted
 614 by the quasi-static heuristic as described in [[Ma+22](#)]. Namely, suppose that after one iteration of
 615 GD, we have perfect bouncing: $x_{t+1} = -x_t$. Substituting this into the GD dynamics, we obtain the
 616 equation

$$\eta \ell'(x_t y_t) y_t = 2x_t. \quad (13)$$

617 According to [Assumption \(A2\)](#), we have $\ell'(x_t y_t) = x_t y_t (1 - \Omega(|x_t y_t|^\beta))$, Together with (13), if
 618 $y_t^2 = (2 + \delta_t)/\eta \geq 2/\eta$, where δ_t is sufficiently small, it suggests that

$$|x_t y_t| \lesssim \delta_t^{1/\beta}. \quad (14)$$

619 The quasi-static prediction (14) fails when δ_t is too small. Nevertheless, we show that it remains
 620 accurate as long as $\delta_t \gtrsim \eta^{\beta/(\beta-1)}$, and consequently we obtain $|x_{\mathbf{t}}y_{\mathbf{t}}| \lesssim \eta^{1/(\beta-1)}$. Combined
 621 with (12), it yields [Theorem 2](#).

622 **B.4 EoS regime: crossing the threshold and the convergence phase**

623 In this section, we prove [Theorem 2](#).

624 We first show that y_t^2 must cross $\frac{2}{\eta}$ in order for GD to converge, and we bound the size of the jump
 625 across $\frac{2}{\eta}$ once this happens.

626 Throughout this section and the next, we use the following notation:

- 627 • $s_t := x_t y_t$;
- 628 • $r_t := \ell'(s_t)/s_t$.

629 In this notation, we can write the GD equations as

$$\begin{aligned} x_{t+1} &= (1 - \eta r_t y_t^2) x_t, \\ y_{t+1} &= (1 - \eta r_t x_t^2) y_t. \end{aligned}$$

630 We also make a remark regarding [Assumption \(A2\)](#). If $\beta < +\infty$, then [Assumption \(A2\)](#) is equivalent
 631 to the following seemingly strongly assumption: for all $r > 0$, there exists a constant $c(r) > 0$ such
 632 that

$$\frac{\ell'(s)}{s} \leq 1 - c(r) |s|^\beta, \quad \text{for all } 0 < |s| \leq r. \quad (\text{A2}^+)$$

633 Indeed, [Assumption \(A2\)](#) states that [\(A2⁺\)](#) holds for *some* $r > 0$. To verify that [\(A2⁺\)](#) holds for
 634 some larger $r' > r$, we can split into two cases. If $|s| \leq r$, then $\ell'(s)/s \leq 1 - c|s|^\beta$. Otherwise, if
 635 $|s| > r$, then $\ell'(r)/r < 1$ and the 1-Lipschitzness of ℓ' imply that $\ell'(s)/s < 1$ for $r \leq |s| \leq r'$, and
 636 hence $\ell'(s)/s \leq 1 - c'|s|^\beta$, for a sufficiently small constant $c' > 0$; thus we can take $c(r') = c \wedge c'$.
 637 Later, we will invoke [\(A2⁺\)](#) with r chosen to be a universal constant, so that $c(r)$ can also be thought
 638 of as universal.

639 We begin with the following result about the limiting value of y_t .

640 **Lemma 7** (threshold crossing). *Let $(\tilde{x}, \tilde{y}) \in \mathbb{R}^2$ satisfy $\tilde{y} > \tilde{x} > 0$ with $\tilde{y}^2 - \tilde{x}^2 = 1$. Suppose we*
641 *initialize GD with step size η with initial point $(x_0, y_0) := \sqrt{\frac{2+\delta}{\eta}}(\tilde{x}, \tilde{y})$, where $\delta > 0$ is a constant.*
642 *Then either $x_t = 0$ for some t or*

$$\lim_{t \rightarrow \infty} y_t^2 \leq \frac{2}{\eta}.$$

643 *Proof.* Assume throughout that $x_t \neq 0$ for all t . Recall the dynamics for y :

$$y_{t+1} = y_t - \eta \ell'(x_t y_t) x_t.$$

644 By assumption $\ell'(s)/s \rightarrow 1$ as $s \rightarrow 0$, and ℓ' is increasing, so this equation implies that if
645 $\liminf_{t \rightarrow \infty} |x_t| > 0$ then y_t^2 must eventually cross $2/\eta$.

646 Suppose for the sake of contradiction that there exists $\varepsilon > 0$ with $y_t^2 > (2 + \varepsilon)/\eta$, for all t . Let
647 $\varepsilon' > 0$ be such that $1 - (2 + \varepsilon)(1 - \varepsilon') < -1$, i.e., $\varepsilon' < \frac{\varepsilon}{2 + \varepsilon}$. Then, there exists $\delta > 0$ such $|x_t| \leq \delta$
648 implies $r_t > 1 - \varepsilon'$, hence

$$\frac{|x_{t+1}|}{|x_t|} = |1 - \eta r_t y_t^2| > |(2 + \varepsilon)(1 - \varepsilon') - 1| > 1.$$

649 The above means that $|x_t|$ increases until it exceeds δ , i.e., $\liminf_{t \rightarrow \infty} |x_t| \geq \delta$. This is our desired
650 contradiction and it implies that $\lim_{t \rightarrow \infty} y_t^2 \leq 2/\eta$. \square

651 **Lemma 8** (initial gap). *Suppose that at some iteration \mathbf{t} , we have*

$$y_{\mathbf{t}+1}^2 < \frac{2}{\eta} \leq y_{\mathbf{t}}^2.$$

652 *Then, it holds that*

$$y_{\mathbf{t}+1}^2 \geq \frac{2}{\eta} - 2\eta s_{\mathbf{t}}^2.$$

653 *Proof.* We can bound

$$y_{\mathbf{t}+1}^2 = y_{\mathbf{t}}^2 - 2\eta \ell'(x_{\mathbf{t}} y_{\mathbf{t}}) x_{\mathbf{t}} y_{\mathbf{t}} + \eta^2 \ell'(x_{\mathbf{t}} y_{\mathbf{t}})^2 x_{\mathbf{t}}^2 \geq y_{\mathbf{t}}^2 - 2\eta |x_{\mathbf{t}} y_{\mathbf{t}}|^2,$$

654 where we used the fact that $|\ell'(s)| \leq |s|$ for all $s \in \mathbb{R}$, \square

655 The above lemma shows that the size of the jump across $2/\eta$ is controlled by the size of $|s_{\mathbf{t}}|$ at the
656 time of the crossing. From Lemma 6, we know that $|s_{\mathbf{t}}| \lesssim 1$, where the implied constant depends on
657 δ . Hence, the size of the jump is always $O(\eta)$.

658 We now provide an analysis of the convergence phase, i.e., after $y_{\mathbf{t}}^2$ crosses $2/\eta$.

659 **Proposition 1** (convergence phase). *Suppose that $y_{\mathbf{t}}^2 < 2/\eta \leq y_{\mathbf{t}-1}^2$. Then, GD converges to $(0, y_{\infty})$*
660 *satisfying*

$$\frac{2}{\eta} - O(|s_{\mathbf{t}}|) \leq y_{\infty}^2 \leq \frac{2}{\eta}.$$

661 *Proof.* Write $y_t^2 = (2 - \rho_t)/\eta$, so that $\rho_t = 2 - \eta y_t^2$. We write down the update equations for x and
662 for ρ . First, by the same argument as in the proof of Theorem 1, we have

$$|x_{t+1}| \leq |x_t| \exp(-\Omega(\rho_t)). \quad (15)$$

663 Next, using $r_t \leq 1$,

$$\begin{aligned} y_{t+1} &= (1 - \eta r_t x_t^2) y_t \geq (1 - \eta x_t^2) y_t, \\ y_{t+1}^2 &\geq (1 - 2\eta x_t^2) y_t^2, \end{aligned}$$

664 which translates into

$$\rho_{t+1} \leq \rho_t + 2\eta^2 x_t^2 y_t^2 \leq \rho_t + 4\eta x_t^2. \quad (16)$$

665 Using these two inequalities, we can conclude as follows. Let $q > 0$ be a parameter chosen later, and
 666 let t be the first iteration for which $\rho_t \geq q$ (if no such iteration exists, then $\rho_t \leq q$ for all t). Note that
 667 $\rho_t \leq q + O(\eta |x_t|)$ due to (15) and (16). By (15), we conclude that for all $t' \geq t$,

$$|x_{t'}| \leq |x_t| \exp(-\Omega(q(t' - t))) \leq |x_t| \exp(-\Omega(q(t' - t))).$$

668 Substituting this into (16),

$$\begin{aligned} \rho_{t'} &\leq \rho_t + 4\eta \sum_{s=t}^{t'-1} x_s^2 \leq q + O(\eta |x_t|) + O(\eta |x_t|^2) \sum_{s=1}^{t'-1} \exp(-\Omega(q(s - t))) \\ &\leq q + O(\eta |x_t|) + O\left(\frac{\eta |x_t|^2}{q}\right). \end{aligned}$$

669 By optimizing this bound over q , we find that for all t ,

$$\rho_t \lesssim \sqrt{\eta} |x_t| \lesssim \eta |s_t|.$$

670 Translating this result back into y_t^2 yields the result. \square

671 Let us take stock of what we have established thus far.

- 672 • According to Lemma 6, $|s_t|$ is bounded for all t by a constant.
- 673 • Then, from Lemma 7 and Lemma 8, we must have either $y_t^2 \rightarrow 2/\eta$, or $2/\eta - O(\eta) \leq y_t^2 \leq 2/\eta$
 674 for some iteration t .
- 675 • In the latter case, Proposition 1 shows that the limiting sharpness is $2/\eta - O(1)$.

676 Note also that the analyses thus far have not made use of Assumption (A2), i.e., we have established
 677 the $\beta = +\infty$ case of Theorem 2. Moreover, for all $\beta > 1$, the asymptotic $2/\eta - O(1)$ still shows
 678 that the limiting sharpness is close to $2/\eta$, albeit with suboptimal rate. The reader who is satisfied
 679 with this result can then skip ahead to subsequent sections. The remainder of this section and the next
 680 section are devoted to substantial refinements of the analysis.

681 To see where improvements are possible, note that both Lemma 8 and Proposition 1 rely on the
 682 size of $|s_t|$ at the crossing. Our crude bound of $|s_t| \lesssim 1$ does not capture the behavior observed
 683 in experiments, in which $|s_t| \lesssim \eta^{1/(\beta-1)}$. By substituting this improved bound into Lemma 7, we
 684 would deduce that the gap at the crossing is $O(\eta^{1+2/(\beta-1)})$, and then Proposition 1 would imply that
 685 the limiting sharpness is $2/\eta - O(\eta^{1/(\beta-1)})$. Another weakness of our proof is that it provides nearly
 686 no information about the dynamics during the bouncing phase, which constitutes an incomplete
 687 understanding of the EoS phenomenon. In particular, we experimentally observe that during the
 688 bouncing phase, the iterates lie very close to the quasi-static envelope (Figure 12). In the next section,
 689 we will rigorously prove all of these observations.

690 Before doing so, however, we show that Proposition 1 can be refined by using Assumption (A2),
 691 which could be of interest in its own right. It shows that even if the convergence phase begins
 692 with a large value of $|s_t|$, the limiting sharpness can be much closer to $2/\eta$ than what Proposition 1
 693 suggests. The following proposition combined with Lemma 6 implies Theorem 2 for all $\beta > 2$, but it
 694 is insufficient for the case $1 < \beta \leq 2$. From now on, we assume $\beta < +\infty$.

695 **Proposition 2** (convergence phase; refined). *Suppose that $y_t^2 < 2/\eta \leq y_{t-1}^2$. Then, GD converges to*
 696 $(0, y_\infty)$ *satisfying*

$$\frac{2}{\eta} \geq y_\infty^2 \geq \frac{2}{\eta} - O(\eta |s_t|^2) - \begin{cases} O(\eta^{1/(\beta-1)}), & \beta > 2, \\ O(\eta \log(|s_t|/\eta)), & \beta = 2, \\ O(\eta |s_t|^{2-\beta}), & \beta < 2. \end{cases}$$

Proof. Let $y_t^2 = (2 - \rho_t)/\eta$ as before. We quantify the decrease of $|x_t|$ in terms of ρ_t and conversely
 the increase of ρ_t in terms of $|x_t|$ by tracking the half-life of $|x_t|$, i.e., the number of iterations it takes
 $|x_t|$ to halve. We call these epochs: at the i -th epoch, we have

$$2^{-(i+1)} \sqrt{\eta} < |x_t| \leq 2^{-i} \sqrt{\eta}.$$

697 Let i_0 be the index of the first epoch, i.e., $i_0 = \lfloor \log_2(\sqrt{\eta}/|x_t|) \rfloor$. Due to Lemma 6, we know that
 698 $i_0 \geq -O(1)$. From (15), $|x_t|$ is monotonically decreasing and consequently $|s_t|$ is decreasing as

699 well. Also, our bound on the limiting sharpness implies that $y_t^2 > 1/\eta$ for all t , provided that η is
700 sufficiently small.

701 Let us now compute the dynamics of ρ_t and $|x_t|$. At epoch i , $|x_t| > 2^{-(i+1)}\sqrt{\eta}$ hence $|s_t| > 2^{-(i+1)}$.
702 Assumption (A2⁺) with $r = |s_t| \lesssim 1$ implies that

$$\frac{\ell'(s_t)}{s_t} \leq 1 - c2^{-\beta(i+1)}, \quad (17)$$

703 where $c = c(|s_t|)$. This allows to refine (15) on the decrease of $|x_t|$ to

$$\begin{aligned} \frac{|x_{t+1}|}{|x_t|} &= \eta r_t y_t^2 - 1 \leq (2 - \rho_t)(1 - c2^{-\beta(i+1)}) - 1 \\ &\leq 1 - \rho_t - c2^{-\beta(i+1)}, \end{aligned}$$

704 where the first inequality follows from (17) and the second from $\rho_t = 2 - \eta y_t^2 < 1$. In turn, this
705 inequality shows that the i -th phase only requires $O(2^{\beta i})$ iterations.

706 Hence, if $t(i)$ denotes the start of the i -th epoch, then (16) shows that

$$\rho_{t(i+1)} \leq \rho_{t(i)} + 4\eta^2 \cdot 2^{-2i} \cdot O(2^{\beta i}) \leq \rho_{t(i)} + O(\eta^2 2^{(\beta-2)i}).$$

707 Summing this up, we have

$$\rho_{t(i)} \leq \rho_{\mathbf{t}} + \eta^2 \times \begin{cases} O(2^{(\beta-2)i}), & \beta > 2, \\ O(i - i_0), & \beta = 2, \\ O(2^{(\beta-2)i_0}) = O(|s_{\mathbf{t}}|^{2-\beta}), & \beta < 2. \end{cases}$$

708 In the case of $\beta < 2$, the final sharpness satisfies $2/\eta - O(\rho_{\mathbf{t}}/\eta) - O(\eta |s_{\mathbf{t}}|^{2-\beta}) \leq y_{\infty}^2 \leq 2/\eta$.

709 In the other two cases, suppose that we use this argument until epoch i_* such that $2^{-i_*} \asymp \eta^\gamma$. Then,
710 we have $|x_{t(i_*)}| \asymp \eta^{\gamma+1/2}$, $|s_{t(i_*)}| \asymp \eta^\gamma$, and by using the argument from Proposition 1 from iteration
711 $t(i_*)$ onward we obtain

$$\rho_{\infty} = \rho_{t(i_*)} + \rho_{\infty} - \rho_{t(i_*)} \leq \rho_{\mathbf{t}} + O(\eta^{\gamma+1}) + \eta^2 \times \begin{cases} O(2^{(\beta-2)i_*}) = O(\eta^{-\gamma(\beta-2)}), & \beta > 2, \\ O(i_* - i_0), & \beta = 2. \end{cases}$$

712 We optimize over the choice of γ , obtaining $\gamma = 1/(\beta - 1)$ and thus

$$\rho_{\infty} \leq \rho_{\mathbf{t}} + \begin{cases} O(\eta^{1+1/(\beta-1)}), & \beta > 2, \\ O(\eta^2 \log(|s_{\mathbf{t}}|/\eta)), & \beta = 2. \end{cases}$$

713 By collecting together the three cases and using Lemma 8 to bound $\rho_{\mathbf{t}}$, we finish the proof. \square

714 Using the crude bound $|s_{t_0}| \lesssim 1$ from Lemma 6, it yields

$$\frac{2}{\eta} \geq y_{\infty}^2 \geq \frac{2}{\eta} - O(\eta) - \begin{cases} O(\eta^{1/(\beta-1)}), & \beta > 2, \\ O(\eta \log(1/\eta)), & \beta = 2, \\ O(\eta), & \beta < 2, \end{cases}$$

715 which is optimal for $\beta > 2$.

716 B.5 EoS regime: quasi-static analysis

717 We supplement Assumption (A2) with a corresponding lower bound on $\ell'(s)/s$:

$$\text{there exists } C > 0 \text{ such that } \frac{\ell'(s)}{s} \geq 1 - C |s|^\beta \quad \text{for all } s \neq 0. \quad (\text{A3})$$

718 Under these assumptions, we prove the following result which is also of interest as it provides detailed
719 information for the bouncing phase of the EoS.

720 **Theorem 5** (quasi-static principle). Suppose we run GD on f with step size $\eta > 0$, where $f(x, y) :=$
721 $\ell(xy)$ and ℓ satisfies Assumptions (A1), (A2), and (A3). Write $y_t^2 := (2 + \delta_t)/\eta$ and suppose that at
722 some iteration t_0 , we have $|x_{t_0}y_{t_0}| \asymp \delta_{t_0}^{1/\beta}$ and $\delta_{t_0} \lesssim 1$. Then, for all $t \geq t_0$ with $\delta_t \gtrsim \eta^{\beta/(\beta-1)}$, we
723 have

$$|x_t y_t| \asymp \delta_t^{1/\beta},$$

724 where all implied constants depend on ℓ but not on η .

725 In this section, we show that the GD iterates lie close to the quasi-static trajectory and give the full
726 proof of [Theorem 2](#). Recall from (13) that the quasi-static analysis predicts

$$\eta r_t y_t^2 \approx 2, \tag{18}$$

727 and that during the bouncing phase, this closely agrees with experimental observations ([Figure 12](#)).
728 We consider the phase where y_t^2 has not yet crossed the threshold $2/\eta$ and we write $y_t^2 := (2 + \delta_t)/\eta$,
729 thinking of δ_t as small. Then, (18) can be written $(2 + \delta_t) r_t \approx 2$. If we have the behavior
730 $\ell'(s)/s = 1 - \Theta(|s_t|^\beta)$ near the origin, then $r_t \approx 1 - \Theta(\delta_t)$ implies that

$$|s_t|^\beta \approx \delta_t. \tag{19}$$

731 Our goal is to rigorously establish (19). However, we first make two observations. First, in order to
732 establish [Theorem 2](#), we only need to prove an upper bound on $|s_t|$, which only requires Assump-
733 tion (A2) (to prove a lower bound on $|s_t|$, we need a corresponding lower bound on $\ell'(s)/s$). Second,
734 even if we relax (19) to read $|s_t|^\beta \lesssim \delta_t$, this fails to hold when δ_t is too small, because the error
735 terms (the deviation of the dynamics from the quasi-static trajectory) begin to dominate. With this in
736 mind, we shall instead prove $|s_t|^\beta \lesssim \delta_t + C' \eta^\gamma$, where the added η^γ handles the error terms and the
737 exponent $\gamma > 0$ emerges from the proof.

738 **Proposition 3** (quasi-static analysis; upper bound). For all t such that $0 \leq \delta_{t-1} \lesssim 1/(\beta \vee 1)$ (for a
739 sufficiently small implied constant), it holds that

$$|s_t|^\beta \leq C (\delta_t + C' \eta^{\beta/(\beta-1)}),$$

740 where $C, C' > 0$ are constants which may depend on the problem parameters but not on η .

741 We first show that [Theorem 2](#) now follows.

742 *Proof of [Theorem 2](#).* As previously noted, the $\beta = +\infty$ case is handled by the arguments of the
743 previous section, so we focus on $\beta < +\infty$. From [Lemma 7](#), we either have $y_t^2 \rightarrow 2/\eta$ and $|x_t| \rightarrow 0$,
744 in which case we are done, or there is an iteration \mathbf{t} such that $y_{\mathbf{t}}^2 < 2/\eta \leq y_{\mathbf{t}-1}^2$. From [Proposition 3](#),
745 since $\delta_{t-1} \geq 0$ and $\delta_t \leq 0$, it follows that $|s_{\mathbf{t}}|^\beta \lesssim \eta^{1/(\beta-1)}$. The theorem now follows, either
746 from [Proposition 1](#) or from the refined [Proposition 2](#). \square

747 We now prove [Proposition 3](#). In the proof, we use asymptotic notation $O(\cdot)$, \lesssim , etc. in order to hide
748 constants that depend on ℓ (including β), but not on δ_t and η . However, the proof also involves
749 choosing parameters $C, C' > 0$, and we keep the dependence on these parameters explicit for clarity.

750 *Proof of [Proposition 3](#).* The proof goes by induction; namely, if $|s_t|^\beta \leq C (\delta_t + C' \eta^\gamma)$ and $\delta_t \geq 0$
751 at some iteration t , we prove that the same holds one iteration later, where the constants $C, C' > 0$ as
752 well as the exponent $\gamma > 0$ are chosen later in the proof.

753 For the base case, observe that the approximate conservation lemma ([Lemma 6](#)) gives $|s_t| \lesssim 1$, and
754 $\delta_t \gtrsim 1/(\beta \vee 1)$ at the beginning of the induction, so the bound is satisfied initially if we choose C
755 sufficiently large enough.

756 Throughout, we also write $\hat{\delta}_t := \delta_t + C' \eta^\gamma$ as a convenient shorthand. The strategy is to prove the
757 following two statements:

- 758 1. If $|s_t|^\beta = C_t \hat{\delta}_t$ for some $C_t > \frac{C}{2}$, then $|s_{t+1}|^\beta \leq C_{t+1} \hat{\delta}_{t+1}$ for some $C_{t+1} \leq C_t$.
- 759 2. If $|s_t|^\beta = C_t \hat{\delta}_t$ for some $C_t \leq \frac{C}{2}$, then $|s_{t+1}|^\beta \leq C \hat{\delta}_{t+1}$.

760 **Proof of 1.** The dynamics for x give

$$|x_{t+1}| = |1 - \eta y_t^2 r_t| |x_t|.$$

761 By Assumption (A2⁺) and $|s_t| \lesssim 1$,

$$r_t \leq 1 - \Omega(|s_t|^\beta) = 1 - \Omega(C\hat{\delta}_t)$$

762 and hence

$$\eta y_t^2 r_t = (2 + \delta_t) (1 - \Omega(C\hat{\delta}_t)) = 2 - \Omega(C\hat{\delta}_t)$$

763 for large C . Also, $\ell''(0) = 1$ and a similar argument as in the proof of [Theorem 1](#) yields the reverse
764 inequality $\eta y_t^2 r_t \gtrsim 1$. We conclude that

$$|x_{t+1}| = (1 - \Omega(C\hat{\delta}_t)) |x_t|$$

765 and hence

$$|s_{t+1}|^\beta \leq (1 - \Omega(C\hat{\delta}_t)) |s_t|^\beta = C_t (1 - \Omega(C\hat{\delta}_t)) \hat{\delta}_t.$$

766 Since we need a bound in terms of $\hat{\delta}_{t+1}$, we use the dynamics of y ,

$$\begin{aligned} y_{t+1} &= (1 - \eta x_t^2 r_t) y_t \geq (1 - \eta x_t^2) y_t, \\ y_{t+1}^2 &\geq (1 - 2\eta x_t^2) y_t^2, \\ \delta_{t+1} &= \eta y_{t+1}^2 - 2 \geq \delta_t - 2\eta^2 s_t^2 \geq \delta_t - 2\eta^2 (C\hat{\delta}_t)^{2/\beta}. \end{aligned} \quad (20)$$

767 Substituting this in,

$$\begin{aligned} |s_{t+1}|^\beta &\leq C_t (1 - \Omega(C\hat{\delta}_t)) (\hat{\delta}_{t+1} + 2\eta^2 (C\hat{\delta}_t)^{2/\beta}) \\ &= C_t \hat{\delta}_{t+1} - \Omega(C^2 \hat{\delta}_t \hat{\delta}_{t+1}) + 2C\eta^2 (C\hat{\delta}_t)^{2/\beta}. \end{aligned} \quad (21)$$

768 Let us show that

$$\hat{\delta}_{t+1} \geq \frac{3}{4} \hat{\delta}_t. \quad (22)$$

769 From (20), we have $\hat{\delta}_{t+1} \geq \hat{\delta}_t - 2\eta^2 (C\hat{\delta}_t)^{2/\beta}$, so we want to prove that $\eta^2 (C\hat{\delta}_t)^{2/\beta} \leq \hat{\delta}_t/8$. If
770 $\beta \leq 2$ this is obvious by taking η small, and if $\beta > 2$ then this is equivalent to $C^{2/\beta} \eta^2 \lesssim \hat{\delta}_t^{1-2/\beta}$.
771 It suffices to have $C^{2/\beta} \eta^2 \lesssim (C')^{1-2/\beta} \eta^\gamma (1-2/\beta)$, which is achieved by taking C' large relative to
772 C and by taking $\gamma \leq 2/(1-2/\beta)$; this constraint on γ will be satisfied by our eventual choice of
773 $\gamma = \beta/(\beta-1)$.

774 Returning to (21), in order to finish the proof and in light of (22), we want to show that $C^2 \hat{\delta}_t^2 \gtrsim$
775 $C^{1+2/\beta} \eta^2 \hat{\delta}_t^{2/\beta}$. Rearranging, it suffices to have $\hat{\delta}_t^{2-2/\beta} \gtrsim C^{2/\beta-1} \eta^2$, or $\hat{\delta}_t^{1-1/\beta} \gtrsim C^{1/\beta-1/2} \eta$.
776 Since by definition $\hat{\delta}_t \geq C' \eta^\gamma$, by choosing C' large it suffices to have $\gamma \leq 1/(1-1/\beta) = \beta/(\beta-1)$,
777 which leads to our choice of γ .

778 **Proof of 2.** Using the simple bound $\eta y_t^2 r_t \leq 2 + \delta_t$, we have

$$\begin{aligned} |s_{t+1}| &\leq (1 + \delta_t) |s_t|, \\ |s_{t+1}|^\beta &\leq \exp(\beta\delta_t) |s_t|^\beta = C_t \exp(\beta\delta_t) \hat{\delta}_t \leq \frac{4}{3} C_t \exp(\beta\delta_t) \hat{\delta}_{t+1} \end{aligned}$$

779 where we used (22). If $\exp(\beta\delta_t) \leq 4/3$, which holds if $\delta_t \lesssim 1/\beta$, then from $C_t \leq C/2$ we obtain
780 $|s_{t+1}|^\beta \leq C\hat{\delta}_{t+1}$ as desired. \square

781 By following the same proof outline but reversing the inequalities, we can also show a corresponding
782 lower bound on $|s_t|^\beta$, as long as $\delta_t \gtrsim \eta^{\beta/(\beta-1)}$. Although this is not needed to establish [Theorem 2](#),
783 it is of interest in its own right, as it shows (together with [Proposition 3](#)) that the iterates of GD do in
784 fact track the quasi-static trajectory.

785 **Proposition 4** (quasi-static analysis; lower bound). *Suppose additionally that (A3) holds and that*
 786 $\beta < +\infty$. *Also, suppose that at some iteration t_0 , we have $\delta_{t_0} \lesssim 1$ and that*

$$|s_t| \geq c \delta_t^{1/\beta} \quad (23)$$

787 *holds at iteration $t = t_0$, where c is a sufficiently small constant (depending on the problem parameters*
 788 *but not on η). Then, (23) also holds for all iterations $t \geq t_0$ such that $\delta_t \gtrsim \eta^{\beta/(\beta-1)}$.*

789 *Proof.* The proof mirrors that of [Proposition 3](#). Let $\delta_t \gtrsim \eta^{\beta/(\beta-1)}$ for a sufficiently large implied
 790 constant. We prove the following two statements:

- 791 1. If $|s_t| = c_t \delta_t^{1/\beta}$ for some $c_t < 2c$, then $|s_{t+1}| \geq c_{t+1} \delta_{t+1}^{1/\beta}$ for some $c_{t+1} \geq c_t$.
 792 2. If $|s_t| = c_t \delta_t^{1/\beta}$ for some $c_t \geq 2c$, then $|s_{t+1}| \geq c \delta_{t+1}^{1/\beta}$.

793 Throughout the proof, due to [Proposition 3](#), we also have $|s_t| \lesssim \delta_t^{1/\beta}$.

794 **Proof of 1.** The dynamics for x give

$$|x_{t+1}| = |1 - \eta y_t^2 r_t| |x_t|.$$

795 By Assumption (A3),

$$r_t \geq 1 - O(|s_t|^\beta) \geq 1 - O(c \delta_t).$$

796 If c is sufficiently small, then

$$\eta y_t^2 r_t \geq (2 + \delta_t) (1 - O(c \delta_t)) \geq 2 + \Omega(\delta_t).$$

797 Therefore, we obtain

$$|x_{t+1}| \geq (1 + \Omega(\delta_t)) |x_t|.$$

798 On the other hand,

$$y_{t+1} \geq (1 - \eta x_t^2) y_t \geq (1 - O(\eta^2 s_t^2)) y_t \geq (1 - O(\eta^2 \delta_t^{2/\beta})) y_t \quad (24)$$

799 and hence

$$\begin{aligned} |s_{t+1}| &\geq (1 + \Omega(\delta_t)) (1 - O(\eta^2 \delta_t^{2/\beta})) |s_t| \geq c_t (1 + \Omega(\delta_t) - O(\eta^2 \delta_t^{2/\beta})) \delta_t^{1/\beta} \\ &\geq c_t (1 + \Omega(\delta_t) - O(\eta^2 \delta_t^{2/\beta})) \delta_{t+1}^{1/\beta}. \end{aligned}$$

800 To conclude, we must prove that $\eta^2 \delta_t^{2/\beta} \lesssim \delta_t$, but since $\delta_t \gtrsim \eta^{\beta/(\beta-1)}$ (with sufficiently large
 801 implied constant), then this holds, as was checked in the proof of [Proposition 3](#).

802 **Proof of 2.** Using Assumption (A3),

$$1 - O(\delta_t) \leq 1 - O(|s_t|^\beta) \leq r_t \leq 1.$$

803 Therefore,

$$2 - O(\delta_t) \leq (2 + \delta_t) (1 - O(\delta_t)) \leq \eta y_t^2 r_t \leq 2 + \delta_t$$

804 and

$$-1 + O(\delta_t) \geq 1 - \eta y_t^2 r_t \geq -1 - \delta_t.$$

805 Together with the dynamics for x and (24),

$$|s_{t+1}| \geq (1 - O(\delta_t)) (1 - O(\eta^2 \delta_t^{2/\beta})) |s_t| \geq c_t (1 - O(\delta_t)) (1 - O(\eta^2 \delta_t^{2/\beta})) \delta_{t+1}^{1/\beta}.$$

806 Since $c_t \geq 2c$, if δ_t and η are sufficiently small it implies $|s_{t+1}| \geq c \delta_{t+1}^{1/\beta}$. □

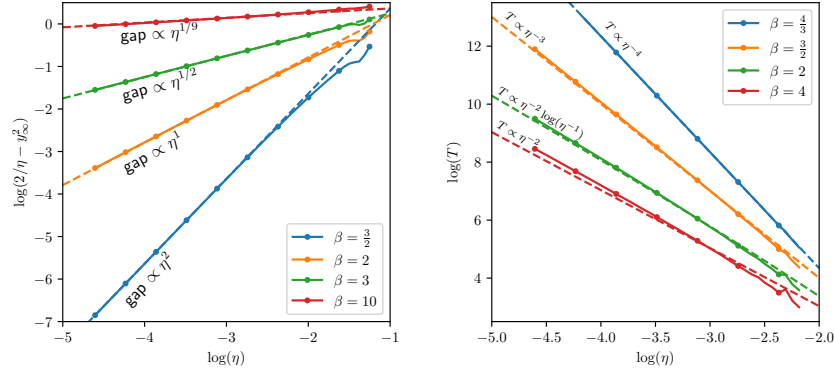


Figure 9: (Left) Log-log plot of the **sharpness gap** as a function of η , for ℓ_β in Example 1 and $\beta = \frac{3}{2}, 2, 3, 10$. (Right) Log-log plot of the **iteration count** for the bouncing region with $y_t^2 \in [\frac{2}{\eta}, \frac{3}{\eta}]$ as a function of η , for ℓ_β in Example 1 and $\beta = \frac{4}{3}, \frac{3}{2}, 2, 4$. **The dashed lines show the predicted sharpness gap and iteration count with an offset computed via linear regression of the data for $\eta < e^{-2}$.**

807 **Convergence rate estimates.** Our analysis also provides estimates for the convergence rate of GD in
 808 both regimes. Namely, in the gradient flow regime, we show that GD converges in $O(1/\eta)$ iterations,
 809 whereas in the EoS regime, GD typically spends $\Omega(1/\eta^{(\beta/(\beta-1))\vee 2})$ iterations ($\Omega(\log(1/\eta)/\eta^2)$
 810 iterations when $\beta = 2$) in the bouncing phase (Figure 9). Hence, the existence of the bouncing phase
 811 dramatically slows down the convergence of GD.

812 *Remark 2.* Suppose that at iteration t_0 , we have $\delta_{t_0} \asymp 1$. Then, the assumption of Proposition 4 is
 813 that $|s_{t_0}| \gtrsim 1$. If this is not satisfied, i.e., $|s_{t_0}| \ll 1$, then the first claim in the proof of Proposition 4
 814 shows that $|s_{t_0+1}| \geq (1 + \Omega(\delta_t)) |s_{t_0}| = (1 + \Omega(1)) |s_{t_0}|$. Therefore, after $t' = O(\log(1/|s_{t_0}|))$
 815 iterations, we obtain $|s_{t_0+t'}| \gtrsim 1$ and then Proposition 4 applies thereafter.

816 *Remark 3.* From the quasi-static analysis, we can also derive bounds on the length of the bouncing
 817 phase. Namely, suppose that t_0 is such that $\delta_{t_0} \asymp 1$ and for all $t \geq t_0$, we have $|s_t| = \delta_t^{1/\beta}$. If δ_{t_0} is
 818 sufficiently small so that $r_t \gtrsim 1$ for all $t \geq t_0$, then the equation for y yields

$$\delta_{t+1} \leq \delta_t - \Theta(\eta^2 s_t^2) = \delta_t - \Theta(\eta^2 \delta_t^{2/\beta}).$$

819 We declare the k -th phase to consist of iterations t such that $2^{-k} \leq \delta_t \leq 2^{-(k-1)}$. During this phase,
 820 $\delta_{t+1} \leq \delta_t - \Theta(\eta^2 2^{-2k/\beta})$, so the number of iterations in phase k is $\asymp 2^{k(2/\beta-1)}/\eta^2$. We sum over
 821 the phases until $\delta_t \asymp \eta^{\beta/(\beta-1)}$, since after this point the quasi-static analysis fails and y_t^2 crosses
 822 over $2/\eta$ shortly afterwards. This yields

$$\frac{1}{\eta^2} \sum_{\substack{k \in \mathbb{Z} \\ \eta^{\beta/(\beta-1)} \lesssim 2^{-k} \lesssim 1}} 2^{k(2/\beta-1)} \asymp \begin{cases} 1/\eta^2, & \beta > 2, \\ \log(1/\eta)/\eta^2, & \beta = 2, \\ 1/\eta^{\beta/(\beta-1)}, & \beta < 2. \end{cases}$$

823 The time spent in the bouncing phase increases dramatically as $\beta \searrow 1$.

824 C Deferred derivations of mean model

In this section, we provide the details for the derivations of the mean model in §3.1. Recall

$$f(\mathbf{x}; a^-, a^+, b) = a^- \sum_{i=1}^d \text{ReLU}(-\mathbf{x}[i] + b) + a^+ \sum_{i=1}^d \text{ReLU}(\mathbf{x}[i] + b),$$

825 where $\mathbf{x} = \lambda y \mathbf{e}_j + \boldsymbol{\xi}$.

We first approximate

$$\sum_{i=1}^d \text{ReLU}(\pm \mathbf{x}[i] + b) \approx \sum_{i=1}^d \text{ReLU}(\pm \boldsymbol{\xi}[i] + b).$$

826 In other words, we can ignore the contribution of the signal $\lambda y e_j$. This approximation holds because
 827 (i) initially, the bias b is not yet negative enough to threshold out the noise, and hence the summation
 828 $\sum_{i=1}^d \text{ReLU}(\pm \boldsymbol{\xi}[i] + b)$ is of size $O(d)$, and (ii) the difference between the left- and right-hand
 829 sides above is simply $\text{ReLU}(\pm \lambda y \pm \boldsymbol{\xi}[j] + b) - \text{ReLU}(\pm \boldsymbol{\xi}[j] + b)$, which is of size $O(1)$ and hence
 830 negligible compared to the full summation.

831 Next, letting $g(b) := \mathbb{E}_{z \sim \mathcal{N}(0,1)} \text{ReLU}(z + b)$ be the ‘smoothed’ ReLU (see [Figure 5](#)), concentration
 832 of measure implies

$$\begin{aligned} 833 & \bullet \sum_{i=1}^d \text{ReLU}(\pm \boldsymbol{\xi}[i] + b) \approx d \mathbb{E}_{\xi \sim \mathcal{N}(0,1)} \text{ReLU}(\xi + b) =: dg(b) \text{ and} \\ 834 & \bullet \sum_{i=1}^d \mathbb{1}\{\pm \mathbf{x}[i] + b \geq 0\} \approx d \mathbb{E}_{\xi \sim \mathcal{N}(0,1)} \mathbb{1}\{\xi + b \geq 0\} = dg'(b). \end{aligned}$$

835 Indeed, the summations above are sums of d i.i.d. non-negative random variables, and hence its
 836 mean is $\Omega(d)$ (as long as $b \geq -O(1)$) and its standard deviation is $O(\sqrt{d})$. Now, using these
 837 approximations, one can rewrite the GD dynamics on the population loss $\mathbb{E}[\ell_{\log i}(yf(\mathbf{x}; a^-, a^+, b))]$.

Using these approximations, the output of the ReLU network (2) can be written as

$$f(\mathbf{x}; a^-, a^+, b) \approx d(a^- + a^+)g(b),$$

838 which in turn leads to an approximation of the GD dynamics on the population loss $(a^-, a^+, b) \mapsto$
 839 $\mathbb{E}[\ell_{\log i}(yf(\mathbf{x}; a^-, a^+, b))]$:

$$\begin{aligned} a_{t+1}^\pm &= a_t^\pm - \eta \mathbb{E} \left[\underbrace{\ell'_{\log i}(y f(\mathbf{x}; a_t^-, a_t^+, b_t))}_{\approx d(a_t^- + a_t^+)g(b_t)} \underbrace{\sum_{i=1}^d \text{ReLU}(\pm \mathbf{x}[i] + b_t)}_{\approx dg(b_t)} \right] \\ &\approx a_t^\pm - \eta \ell'_{\text{sym}}(d(a_t^- + a_t^+)g(b_t)) dg(b_t), \\ b_{t+1} &= b_t - \eta \mathbb{E} \left[\underbrace{\ell'_{\log i}(y f(\mathbf{x}; a_t^-, a_t^+, b_t))}_{\approx d(a_t^- + a_t^+)g(b_t)} \right. \\ &\quad \times \left(\underbrace{a_t^- \sum_{i=1}^d \mathbb{1}\{-\mathbf{x}[i] + b_t \geq 0\}}_{\approx dg'(b_t)} + \underbrace{a_t^+ \sum_{i=1}^d \mathbb{1}\{+\mathbf{x}[i] + b_t \geq 0\}}_{\approx dg'(b_t)} \right) \Big] \\ &\approx b_t - \eta \ell'_{\text{sym}}(d(a_t^- + a_t^+)g(b_t)) d(a_t^- + a_t^+)g'(b_t), \end{aligned}$$

840 where $\ell_{\text{sym}}(s) := \frac{1}{2}(\log(1 + \exp(-s)) + \log(1 + \exp(+s)))$ is the symmetrized logistic loss. Hence
 841 we arrive at the following dynamics on a^\pm and b that we call the mean model:

$$\begin{aligned} a_{t+1}^\pm &= a_t^\pm - \eta \ell'_{\text{sym}}(d(a_t^- + a_t^+)g(b_t)) dg(b_t), \\ b_{t+1} &= b_t - \eta \ell'_{\text{sym}}(d(a_t^- + a_t^+)g(b_t)) d(a_t^- + a_t^+)g'(b_t). \end{aligned}$$

842 Now, we can write the above dynamics more compactly in terms of the parameter $A_t := d(a_t^- + a_t^+)$.

$$\begin{aligned} A_{t+1} &= A_t - 2d^2 \eta \ell'_{\text{sym}}(A_t g(b_t)) g(b_t), \\ b_{t+1} &= b_t - \eta \ell'_{\text{sym}}(A_t g(b_t)) A_t g'(b_t). \end{aligned}$$

843 D Proofs for the mean model

844 In this section, we prove the main theorems for the mean model. We first recall the mean model for
 845 the reader’s convenience.

$$\begin{aligned} A_{t+1} &= A_t - 2d^2 \eta \ell'(A_t g(b_t)) g(b_t), \\ b_{t+1} &= b_t - \eta \ell'(A_t g(b_t)) A_t g'(b_t). \end{aligned}$$

846 **D.1 Deferred proofs**

847 In this section, we collect together deferred proofs from §3.2.

848 *Proof of Lemma 2.* By definition, $g(b) = \int_{-b}^{\infty} (\xi + b) \varphi(\xi) d\xi = \int_{-b}^{\infty} \xi \varphi(\xi) d\xi + b\Phi(b)$. Recalling
849 $\varphi'(\xi) = -\xi \varphi(\xi)$, the first term equals $\varphi(b)$. Moreover, $g'(b) = -b \varphi(b) + \Phi(b) + b \varphi(b) = \Phi(b)$. \square

850 *Proof of Lemma 3.* Note that $\partial_t(\frac{1}{2} A^2) = A\dot{A} = -2d^2 \ell'(Ag(b)) Ag(b)$ and also that $\partial_t \kappa(b) =$
851 $-\ell'(Ag(b)) \kappa'(b) Ag'(b) = -\ell'(Ag(b)) Ag(b)$ since $\kappa' = g/g'$. Hence, $\partial_t(\frac{1}{2} A^2 - 2d^2 \kappa(b)) = 0$
852 and the proof is completed. \square

853 **D.2 Gradient flow regime**

854 *Proof of Theorem 3.* The following proof is analogous to the proof of Theorem 1. We first list several
855 facts we use in the proof:

856 (i) $|g'(b)| = |\Phi(b)| \leq 1$ for all $b \in \mathbb{R}$.

857 (ii) $\ell'(s) = \frac{1}{2} \frac{\exp(s)-1}{\exp(s)+1}$. Hence, $|\ell'(s)| \leq \frac{1}{2}$ for all $s \in \mathbb{R}$, and we have

$$\frac{\ell'(s)}{s} \geq \frac{1}{8} \times \begin{cases} 1, & \text{if } |s| \leq 2, \\ 2/|s|, & \text{if } |s| > 2. \end{cases}$$

858 (iii) $\ell''(0) = 1/4$.

859 (iv) $\ell'''(s) = -\frac{\exp(s)(\exp(s)-1)}{(\exp(s)+1)^3}$. Hence, $\ell'''(s) < 0$ for $s > 0$ and $\ell'''(s) > 0$ for $s < 0$. In particular,

860 $|\ell'(s)| \leq \frac{1}{4} |s|$ for all $s \in \mathbb{R}$.

861 Throughout the proof, we assume that $A_0 > 0$ without loss of generality. We prove by induction the
862 following claim: for $t \geq 0$ and

$$\gamma := \frac{1}{200} \min\left\{\delta, 8 - \delta, \frac{8 - \delta}{A_0}\right\},$$

863 it holds that

$$|A_t| \leq A_0 \exp(-\gamma t).$$

864 This clearly holds at initialization.

865 Suppose that the claim holds up to iteration t . Using the bounds on $|g'|$ and $|\ell'|$, it follows that

$$\begin{aligned} b_{t+1} &\geq b_t - |\ell'(A_t g(b_t))| |A_t| g'(b_t) \geq b_t - \frac{1}{2} \eta |A_t| \\ &\geq b_t - \frac{1}{2} \eta A_0 \exp(-\gamma t) \geq \dots \geq b_0 - \frac{1}{2} \eta A_0 \sum_{s=0}^t \exp(-\gamma s) \geq -\frac{\eta A_0}{\gamma}. \end{aligned}$$

866 In particular, $b_t \geq -1$ and $g(b_t) > 0.08$, since $\eta \leq \frac{\gamma}{A_0}$. Also, the bound shows that if the claim holds
867 for all t , then we obtain the desired conclusion.

868 It remains to establish the inductive claim; assume that it holds up to iteration t . For the dynamics of
869 A , by symmetry we may suppose that $A_t > 0$. From $\ell'(A_t g(b_t)) \leq A_t g(b_t)/4$ and $g(b_t) \leq g(0) =$
870 $\frac{1}{\sqrt{2\pi}}$,

$$\begin{aligned} A_{t+1} &= A_t - 2\eta d^2 \ell'(A_t g(b_t)) g(b_t) \geq \left(1 - \frac{\eta d^2}{2} g(b_t)^2\right) A_t \\ &\geq \left(1 - \frac{\eta d^2}{2} g(0)^2\right) A_t = -\left(1 - \frac{\delta}{4}\right) A_t. \end{aligned}$$

871 This shows that $A_{t+1} \geq -(1 - \gamma) A_t$. Next, we show that $A_{t+1} \leq (1 - \gamma) A_t$. First, if $A_t g(b_t) \leq 2$,

$$\begin{aligned} A_{t+1} &= A_t - 2\eta d^2 \ell'(A_t g(b_t)) g(b_t) \leq A_t - \frac{1}{4} \eta d^2 A_t g(b_t)^2 \\ &= \left(1 - \frac{(8 - \delta)\pi}{4} g(b_t)^2\right) A_t \leq \left(1 - \frac{(8 - \delta)}{4} \pi \cdot 0.08^2\right) A_t \leq (1 - \gamma) A_t, \end{aligned}$$

872 since we have $g(b_t) > 0.08$. Next, if $A_t g(b_t) \geq 2$, then

$$\begin{aligned} A_{t+1} &= A_t - 2\eta d^2 \ell'(A_t g(b_t)) g(b_t) \leq A_t - \frac{1}{2} \eta d^2 g(b_t) = \left(1 - \frac{(8-\delta)\pi}{2} \frac{g(b_t)}{A_t}\right) A_t \\ &\leq \left(1 - \frac{(8-\delta)\pi}{2} \cdot \frac{0.08}{A_0}\right) A_t \leq (1-\gamma) A_t. \end{aligned}$$

873 This shows that $|A_{t+1}| \leq (1-\gamma)|A_t|$ for the case $A_t > 0$. A similar conclusion is obtained for the
874 case $A_t < 0$. The induction is complete. \square

875 D.3 EoS regime

876 *Proof of Theorem 4.* The following proof is analogous to the proof of Lemma 7. Assume throughout
877 that $A_t \neq 0$ for all t . Recall the dynamics for b :

$$b_{t+1} = b_t - \eta \ell'(A_t g(b_t)) A_t g'(b_t).$$

878 Since $\ell'(s)/s \rightarrow 1/4$ as $s \rightarrow 0$, and ℓ' is increasing, this equation implies that if $\liminf_{t \rightarrow \infty} |A_t| > 0$
879 then b_t must keep decreasing until $\frac{1}{2} d^2 g(b_t)^2 < 2/\eta$.

880 Suppose for the sake of contradiction that there exists $\varepsilon > 0$ with $\frac{1}{2} d^2 g(b_t)^2 > (2+\varepsilon)/\eta$, for all
881 t . Let $\varepsilon' > 0$ be such that $1 - (2+\varepsilon)(1-\varepsilon') < -1$, i.e., $\varepsilon' < \frac{\varepsilon}{2+\varepsilon}$. Then, there exists $\delta > 0$ such
882 $|A_t| \leq \delta$ implies $\ell'(A_t g(b_t))/(A_t g(b_t)) > \frac{1}{4}(1-\varepsilon')$, hence

$$\frac{|A_{t+1}|}{|A_t|} = \left|1 - 4 \cdot \frac{1}{4}(1-\varepsilon') \cdot \frac{1}{2} \eta d^2 g(b_t)^2\right| > |(2+\varepsilon)(1-\varepsilon') - 1| > 1.$$

883 The above means that $|A_t|$ increases until it exceeds δ , i.e., $\liminf_{t \rightarrow \infty} |A_t| \geq \delta$. This is our desired
884 contradiction and it implies that $\lim_{t \rightarrow \infty} \frac{1}{2} d^2 g(b_t)^2 \leq 2/\eta$. \square

885 *Remark 4.* A straightforward calculation yields that when $(a_\star^-, a_\star^+, b_\star)$ is a global minimizer (i.e.,
886 $a_\star^- + a_\star^+ = 0$), then $\lambda_{\max}(\nabla^2 f(a_\star^-, a_\star^+, b_\star)) = \frac{1}{2} d^2 g(b_\star)^2$. The mean model initialized at $(A_0, 0)$
887 approximately reaches $(0, 0)$ whose sharpness is $d^2 g(0)^2/2 = d^2/4\pi$. Hence, the bias learning
888 regime $2/\eta < d^2/(4\pi)$ precisely corresponds to the EoS regime, $2/\eta < \lambda_{\max}(\nabla^2 f(a_\star^-, a_\star^+, b_\star))$.

889 E Additional figures

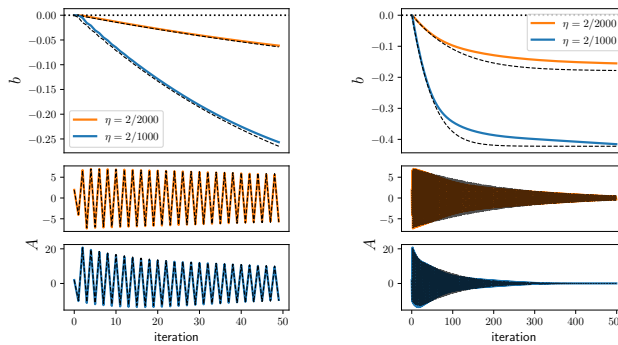


Figure 10: Under the same setting as Figure 1, we compare the mean model with the GD dynamics of the ReLU network. The mean model is plotted with black dashed line. Note that **the mean model tracks the GD dynamics quite well during the initial phase of training.**

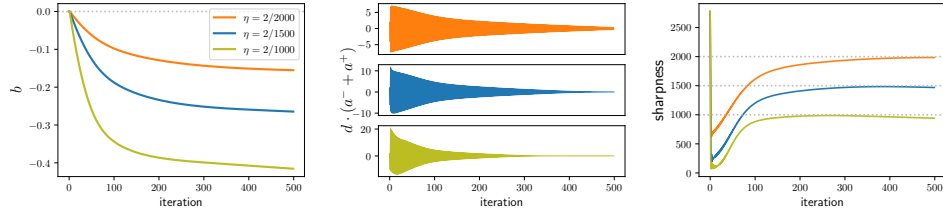


Figure 11: **Understanding our main question is surprisingly related to the EoS.** Under the same setting as Figure 1, we report the largest eigenvalue of the Hessian (“sharpness”), and observe that GD iterates lie in the EoS during the initial phase of training when there is a fast drop in the bias.

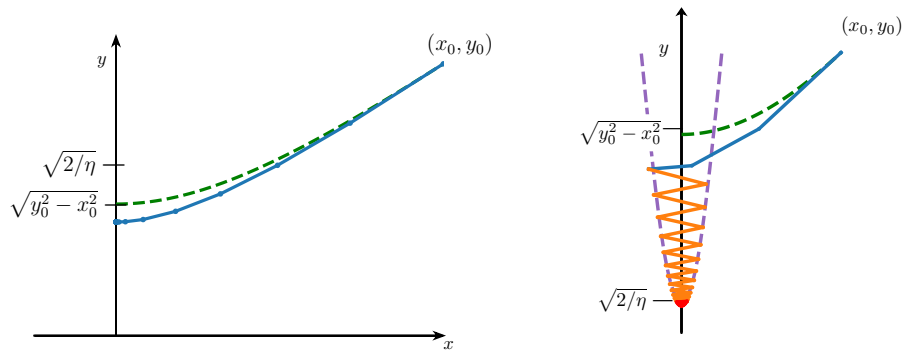


Figure 12: **Two regimes for GD.** We run GD on the square root loss with step size $\frac{1}{4}$. The gradient flow regime is illustrated on the left for $(x_0, y_0) = (3, 4)$. GD (blue) tracks the gradient flow (green) when $\eta < 2/(y_0^2 - x_0^2)$. Otherwise, as illustrated on the right for $(x_0, y_0) = (3, 6)$, GD is in the EoS regime and goes through a gradient flow phase (blue), an intermediate bouncing phase (orange) that tracks the quasi-static envelope (purple), and a converging phase (red).