
An empirical study of the (L_0, L_1) -smoothness condition in deep learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The (L_0, L_1) -smoothness condition was introduced by Zhang-He-Sra-Jadbabai in
2 2020, who both proved convergence bounds under this assumption and provided
3 empirical evidence it is satisfied in deep learning. Since then, many groups have
4 proven convergence guarantees for functions which satisfy this condition, motivated
5 by the expectation that loss functions arising in deep learning satisfy it. In this paper
6 we provide further empirical study of this condition in the setting of feedforward
7 neural networks of depth at least 2, with L_2 or cross entropy loss. The results
8 suggest that the (L_0, L_1) -smoothness condition is not satisfied in this setting.

9 1 Introduction

10 The two properties most prized in the field of optimization are convexity and smoothness. Many
11 desirable convergence results have been proved for functions satisfying both of these conditions.

12 Meanwhile, in the past decade, deep learning has astonished experts and laypeople alike with its
13 power and versatility. Understanding how deep neural networks work is of great interest, and it would
14 be very desirable to apply the tools of optimization theory to understand the convergence properties
15 of deep neural networks.

16 At present this is usually not possible. This is because many powerful results in optimization theory
17 assume the function being optimized is smooth and convex, while it is well established that the loss
18 function of deep neural networks is neither smooth nor convex. Therefore, it would be useful to
19 develop weaker notions of smoothness and convexity, that both hold for deep neural networks and for
20 which good convergence guarantees can be proved.

21 This motivated Zhang-He-Sra-Jadbabai Zhang et al. [2020b] to introduce the (L_0, L_1) -smoothness
22 condition in 2020. In this promising approach, they succeeded both to provide empirical evidence that
23 this form of smoothness is satisfied by deep neural networks, as well as later prove good convergence
24 bounds for functions that satisfy this condition.

25 Motivated by the expectation that deep neural networks satisfy this condition, expressed for example
26 in Chen et al. [2023] and Crawshaw et al. [2022], a considerable amount of work has been done
27 proving convergence bounds under the (L_0, L_1) condition, for several different optimizers. Groups
28 including Wang-Zhang-Zhang-Meng-Ma-Chen in Wang et al. [2022] and Faw-Rout-Caramanis-
29 Shakkottai in Faw et al. [2023] have extended the work of Zhang-He-Sra-Jadbabai, and succeeded in
30 proving good convergence guarantees for functions that satisfy this condition.

31 However there has been limited study of whether deep neural networks satisfy the (L_0, L_1) -
32 smoothness condition. In 2022 Patel-Zhang-Tian gave two examples of neural network loss functions
33 L that are not (L_0, L_1) -smooth, in Patel et al. [2022]. Their work showed it is not the case that all
34 loss functions arising in deep learning satisfy the (L_0, L_1) condition, but left open the question of

35 whether the networks they found were exceptional, or reflected a general property of loss functions
36 arising in deep learning.

37 In this paper we undertake further empirical study of (L_0, L_1) -smoothness in the setting of deep
38 feedforward neural networks. Our experiments suggest that for deep feedforward neural networks,
39 the failure of the (L_0, L_1) -smoothness condition that Patel-Zhang-Tian observed in their examples is
40 not the exception but the rule.

41 Our **first contribution** is to compute the magnitudes of the gradient and hessian of L along a fixed
42 line \mathcal{R} , with L the loss function of a feedforward neural network. Under either L_2 loss or cross
43 entropy loss, we observe a quadratic relationship, suggesting a failure of the (L_0, L_1) -smoothness
44 condition.

45 Our **second contribution** is again to look at the loss function L of a feedforward neural network with
46 L_2 loss. We sample a number of randomly generated initializations, and compute the the magnitudes
47 of the gradient and hessian of L along a radial line segment through those lines. We observe a range
48 of behaviors, some consistent with the (L_0, L_1) -smoothness condition, others suggesting a failure of
49 the (L_0, L_1) -smoothness condition near initialization.

50 In the next section, we discuss related work. In Section 3 we describe several empirical results.
51 In Appendix A we provide background and notation, and in Appendix B we provide supporting
52 materials.

53 2 Related work

54 To date, none of the weaker smoothness conditions discussed in Section A.1 have been successfully
55 verified for deep neural networks. In response, in 2020 Zhang-He-Sra-Jadbabai proposed a novel
56 relaxation of the classical notion of smoothness, along with reasons one may hope that this condition
57 is satisfied by deep neural networks, in Zhang et al. [2020b]. Their innovation is to allow the local
58 smoothness constant to increase with the gradient norm.

59 **Definition 1** (see Zhang et al. [2020b]). Given two real numbers L_0 and L_1 , a twice differentiable
60 function $L : \mathbb{R}^d \rightarrow \mathbb{R}$ is (L_0, L_1) -smooth if the Hessian of L , denoted $\mathbf{H}L$, satisfies the inequality

$$\|\mathbf{H}L(\rho)\| \leq L_0 + L_1 \|\nabla L(\rho)\| \quad (1)$$

61 for all $\rho \in \mathbb{R}^d$, where the norm on the left is taken to be the operator norm of the matrix and the norm
62 on the right is the L^2 norm of the vector.

Remark 1. All matrix norms are equivalent up to constants, that is, for any two matrix norms $\|\cdot\|_A$
and $\|\cdot\|_B$, there exist constants q and r such that

$$q\|M\|_A \leq \|M\|_B \leq r\|M\|_A,$$

63 for all matrices M . Therefore, one can equivalently use any matrix norm on the left hand side of
64 Equation 1, up to rescaling the constants L_0 and L_1 . In this paper, we will take the Frobenius norm
65 on the left side and the L^2 norm on the right side of Equation 1.

66 After introducing this condition, Zhang-He-Sra-Jadbabai went on to prove upper and lower bounds
67 convergence for functions L that satisfy the (L_0, L_1) -smoothness condition for some choice of L_0
68 and L_1 , for several different optimizers. They consider gradient descent, clipped gradient descent, as
69 well as stochastic versions of both.

70 They go on to provide empirical evidence that this condition is satisfied by deep neural networks.
71 In experiments on a variety of architectures and tasks, including image recognition and language
72 generation, they consider the (L_0, L_1) -smoothness condition in the regions of the loss landscape
73 traversed during training and find evidence that it is satisfied.

74 Because a condition which both allows us to prove convergence results and is satisfied by deep neural
75 networks has been long desired, this work is very appealing and a number of works have expanded
76 on this seminal work.

77 Following Zhang et al. [2020b], several groups have gone on to provide analyses of the convergence
78 properties of additional optimizers, and under a wider range of assumptions, for functions that
79 satisfy the (L_0, L_1) -smoothness condition. In Li et al. [2023a], and Li et al. [2023b], Li-Qian-Tian-
80 Rakhlin-Jadbabai proved convergence results for additional optimizers, such as Adam, and under

81 a wider range of assumptions, including a generalization of the (L_0, L_1) -smoothness condition. In
 82 related work, Faw-Rout-Caramanis-Shakkottai developed new techniques allowing them to derive
 83 convergence bounds for SGD without assuming uniform bounds on the noise support in Faw et al.
 84 [2023].

85 There has been less attention on studying this new condition in the specific context of deep neural
 86 networks. Since the (L_0, L_1) -smoothness condition was introduced in Zhang et al. [2020b], fewer
 87 groups have analyzed the motivating hope that loss functions arising from deep neural networks
 88 satisfy the (L_0, L_1) -smoothness condition.

89 One group that did consider this question is Patel-Zhang-Tian, who gave a theoretical analysis of
 90 the geometry of several loss functions and in doing so produced two examples of loss functions L
 91 that do not satisfy the (L_0, L_1) -smoothness condition in Patel et al. [2022]. The first is a very simple
 92 feedforward network with three linear layers and a single nonlinear layer, learning the simple dataset
 93 input 0 output 0 with probability $1/2$ and input 1 output 1 with probability $1/2$.

94 The second is a 1-dimensional linear recurrent neural network, learning a similar dataset. In both
 95 cases, they give a complete mathematical analysis of the smoothness of the resulting loss function
 96 and conclude not only that the loss functions are not m -smooth for any m , but also do not satisfy the
 97 (L_0, L_1) -smoothness condition for any choice of L_0 and L_1 .

98 In this paper, we find that these examples are not isolated. Our experiments suggest that they are in
 99 fact representative of the general case.

100 3 Experiments

101 3.1 Fixed line

102 In the first experiment, we fix a line motivated by the work of Patel et al. [2022], compute the norms
 103 of the Hessian and gradient along the line, and observe the relative growth rates. We made the
 104 following computations in Mathematica, any other programming platform that can compute neural
 105 networks can be used as well.

106 We begin with the case of L2 loss. We begin by initializing a feedforward neural network of layer
 107 widths $(1, 4, 7, 1)$. We take the activation function σ to be \tanh . Next, we choose 20 data points,
 108 with x and y drawn uniformly at random from the interval $[-1, 1]$. Having made these choices, we
 109 can compute the corresponding loss function.

110 In the case of L2 loss, the direction in the loss landscape we sample from is the region near the line
 111 \mathcal{R} specified in Appendix B. We draw points at random from a tubular neighborhood of this line, by
 112 taking points on the line for values of t between 100 and 700, incrementing by 3 each time, and
 113 adding noise drawn uniformly at random to each point, with width $\epsilon(t) = 1/t^2$. Finally, we compute
 114 the norms of the Hessian and gradient and the sampled points.

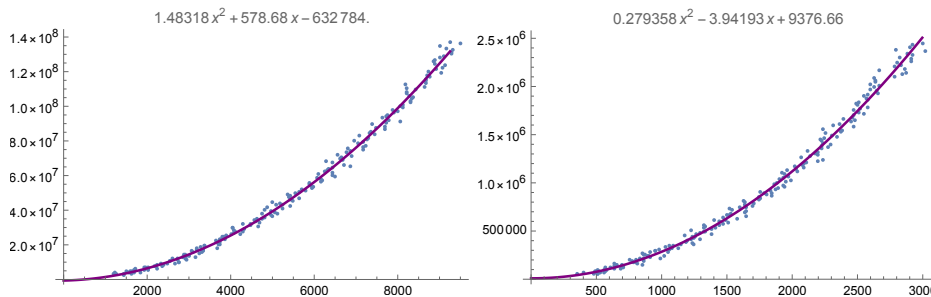


Figure 1: Left: our calculations for the example with L2 loss are shown in a scatter plot of the pairs (norm of gradient, norm of Hessian). Right: the same, for the example with cross-entropy loss.

115 For the case of cross-entropy loss, we proceed similarly. We begin by initializing a feedforward
 116 neural network of layer widths $(1, 3, 3, 2)$. We take the activation function σ to be \tanh . Next, we
 117 choose 10 data points, at random, with x drawn uniformly at random from the interval $[-1, 1]$ and
 118 y assigned to be $(1, 0)$ for the first 8 points, and $(0, 1)$ for the remaining. (The width of the tubular

119 neighborhood is chosen based on the proportions of each label, so fixing the proportions is easiest.)
 120 Having made these choices, we can compute the corresponding loss function.

121 In the case of cross-entropy loss, the direction in the loss landscape we sample from is the region
 122 near the line \mathcal{R} described in the proof. We draw points at random from a tubular neighborhood of
 123 this line, by taking points on the line for values of t between 100 and 700, incrementing by 3 each
 124 time, and adding noise drawn uniformly at random to each point, with width $\epsilon(t) = 20/t^2$. Finally,
 125 we compute the norms of the Hessian and gradient and the sampled points.

126 The resulting plots are shown in Figure 1, together with the degree 2 polynomial of best fit for each.

127 At first glance, these plots look different than the ones shown in Zhang et al. [2020b]. Here we note
 128 that the two sets of plots are consistent, as the scatter plots in Figure 1 are plotted directly, while the
 129 scatter plots shown in the body of Zhang et al. [2020b] are shown on a log-log plot.

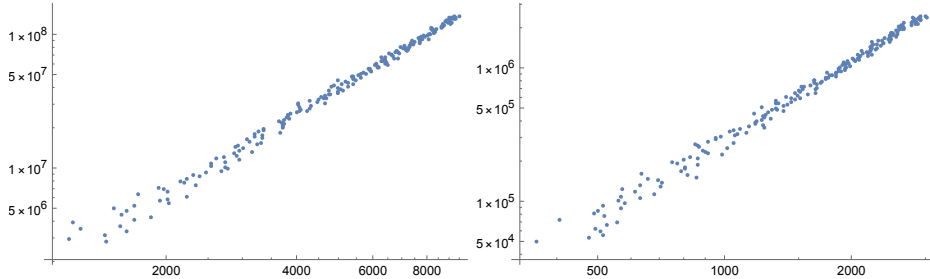


Figure 2: Left: our calculations for the example with L2 loss are shown in a log-log scatter plot of the pairs (norm of gradient, norm of Hessian). Right: the same, for the example with cross-entropy loss.

130 Note that a polynomial of any degree shown on a log-log plot will look linear. In Figure 2, we
 131 redisplay the information on Figure 1 on log-log plots. In this format, the figures look similar to the
 132 figures in Zhang et al. [2020b].

133 3.2 Random segments

134 In the previous experiment, we studied the relationship between the magnitudes of the gradient and
 135 the Hessian of L along a fixed line. One might ask what plots of the magnitude of the gradient against
 136 the magnitude of the Hessian look like in regions encountered when training a neural network. One
 137 place to look is at initialization.

138 In this experiment, we again compare the gradient and the Hessian of L , this time near points
 139 initialized according to Kaiming initialization, along random line segments through those points.

140 We used the same architectures as in Subsection 3.1. Namely, we consider a feedforward network
 141 with layer widths $[1, 4, 7, 1]$, L2 loss, and using tanh for the activation function. We then generated 20
 142 data points at random to define the loss function L . We record the randomly chosen data in Appendix
 143 B.

144 We then generated 15 random initialization parameters p_1, \dots, p_{15} using the Kaiming initialization
 145 procedure. We chose a radial line segment S_i through each parameter p_i , and computed the magni-
 146 tudes of the gradient and hessian at 50 equidistributed points along each S_i . In 6 cases the relationship
 147 appeared approximately linear. In 5 cases the relationship appeared superlinear. In 4 cases, other
 148 nonlinear graphs were observed.

149 In the following, we display some of these graphs, numbered in the order of appearance, not the order
 150 they were generated in. In Appendix B we record the endpoints of the line segments S_i , which we
 151 call $start_i$ and end_i .

152 At the randomly initialized point p_1 on the left in Figure 3, the resulting graph appears approximately
 153 linear to begin, then approximately like an upward-facing semicircle. At the randomly initialized
 154 point p_2 on the right, the resulting graph appears approximately linear for the first quarter, then
 155 approximately linear to begin, then approximately linear but with a steeper slope, then approximately
 156 linear with an even steeper slope in the last stretch.

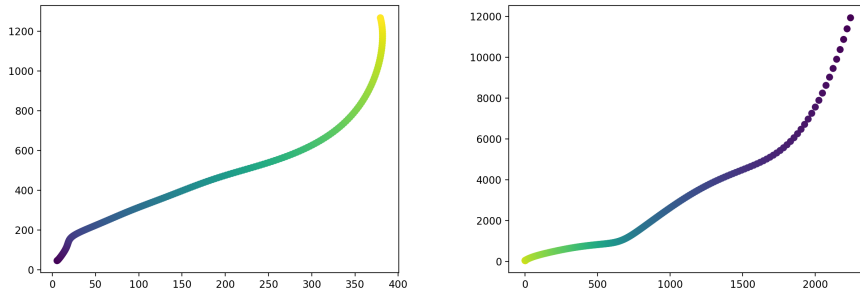


Figure 3: We show $|\nabla L|$ along the x-axis and the $|hess(L)|$ along the y-axis.

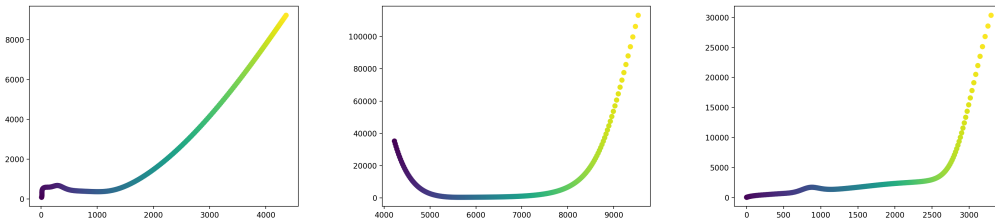


Figure 4: We show $|\nabla L|$ along the x-axis and the $|hess(L)|$ along the y-axis.

157 At the randomly initialized point p_3 on the left in Figure 4, the resulting graph appears approximately
 158 quadratic. At the randomly initialized point p_4 in the center, the resulting graph looks U-shaped.
 159 Finally at the randomly initialized point p_5 on the right, the resulting graph grows quickly at the end.
 160 The last two calculations show examples of initializations near which, in the radial direction, the
 161 magnitude of the hessian does appear to be bounded by a linear function of the magnitude of the
 162 gradient, as seen in Figure 5.

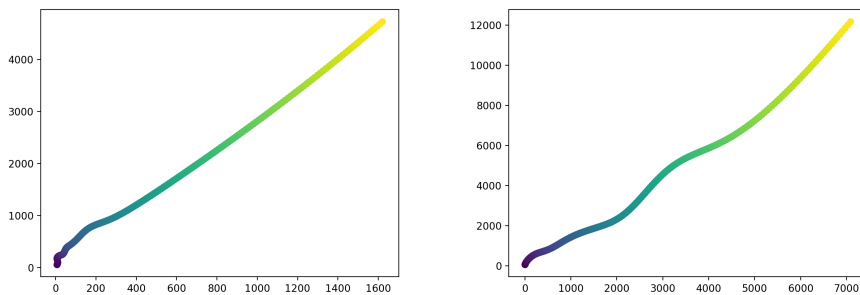


Figure 5: We show $|\nabla L|$ along the x-axis and the $|hess(L)|$ along the y-axis.

163 We note that empirical measurements cannot prove if the loss function L satisfies an (L_0, L_1) -
 164 smoothness condition or not. Indeed, on any compact set such as this spherical shell we are studying
 165 in this experiment the loss function L not only will satisfy a (L_0, L_1) -smoothness condition for some
 166 choice of L_0 and L_1 but will satisfy m -smoothness for some choice of m .
 167 That being said, on a compact region one may ask if the hessian of L appears bounded by a linear
 168 function of the gradient of L . Near some of the random initializations we generated the hessian of L
 169 does appear bounded by a linear function of the gradient of L , such as the plots in Figure 5.

170 However, near other random initializations we generated, the answer appears to be no, such as in the
171 plots in Figure 4. This provides empirical evidence that in the region near random initializations, the
172 loss function does not satisfy a (L_0, L_1) -smoothness condition.

173 4 Conclusion

174 In this paper, we made an empirical study of the (L_0, L_1) -smoothness condition in the setting of
175 feedforward networks, with either L_2 or cross-entropy loss. The results suggest that the (L_0, L_1) -
176 smoothness condition is not in general satisfied.

177 Thus the convergence guarantees that have been proved for (L_0, L_1) -smoothness might not be directly
178 applicable to the loss functions arising from deep feedforward networks. Though we take a different
179 conclusion from Zhang-He-Sra-Jadbabai, our results are not in contradiction with the empirical
180 studies in the original paper by Zhang et al. [2020a]. Note that in their work, they compute the
181 magnitude of the gradient of L and the magnitude of the hessian of L along gradient trajectories, but
182 do not compute those quantities in transverse directions.

183 In contrast, we compute the magnitude of the gradient of L and the magnitude of the hessian of
184 L in radial directions. So it is not contradictory that we observe different relationships. We note
185 that the geometry near a gradient trajectory, in directions transverse to the trajectory, are relevant in
186 theoretical bounds on convergence. So the additional empirical study here provides useful further
187 information.

188 In recent work, Li-Quian-Tian-Rakhlin-Jadbabaie Li et al. [2023a] introduce a class of conditions
189 generalizing the (L_0, L_1) -smoothness condition, which they call ℓ -smoothness conditions, for any
190 function ℓ . The (L_0, L_1) -smoothness condition is recovered in the special case that ℓ is an affine
191 linear function.

192 The rates of growth of the magnitude of the hessian as a function of the magnitude of the gradients
193 we observe suggest that not only does the loss function L of a deep neural network not satisfy the
194 (L_0, L_1) -smoothness condition, that is ℓ -smoothness for a linear function, but that L also does not
195 satisfy ℓ -smoothness for any subquadratic function ℓ . This is worth noting because in Li et al. [2023a],
196 convergence guarantees proven in cases when ℓ is subquadratic, and in the thorough analysis given,
197 examples are also provided illustrating that similar guarantees are not possible in cases when ℓ is
198 quadratic or superquadratic. Our work shows that the loss functions of deep feedforward networks lie
199 in this more challenging setting.

200 Our work suggests that in order to develop similar convergence arguments that can be applied directly
201 to the loss functions arising in deep learning, different generalizations of the (L_0, L_1) -smoothness
202 condition may be needed.

203 One could also study (L_0, L_1) -smoothness with the approach used to study weak convexity in the
204 setting of deep networks by Liu-Zhu-Belkin Liu et al. [2022]. It may be that while the (L_0, L_1) -
205 smoothness condition does not hold uniformly over the loss landscapes of deep feedforward networks,
206 that it is possible to identify regions of the loss landscape on which (L_0, L_1) -smoothness holds.

207 This work is an invitation to interesting directions for future work. The study of (L_0, L_1) -smoothness
208 is an exciting nexus where new techniques in optimization are being developed with inspiration from
209 the geometries that arise in deep learning.

210 5 Broader Impacts

211 This work focuses on the mathematical understanding of a technical aspect of deep learning. While
212 this may feel removed from the machine learning systems that are beginning to be integrated into our
213 daily lives, advances in this and similar papers are expected to improve the performance of machine
214 learning systems over time. Therefore this work may have greater societal impact than is initially
215 apparent.

216 As the authors of this work, we have a responsibility to make our technical advancements understand-
217 able to the broadest range of people, to promote the beneficial uses of these technologies, and to work
218 to mitigate the risks of these technologies.

References

- 219
- 220 A. Agarwal, S. Negahban, and M. J. Wainwright. Fast global convergence of gradient methods for
221 high-dimensional statistical recovery. *The Annals of Statistics*, 40(5):2452 – 2482, 2012. doi:
222 10.1214/12-AOS1032. URL <https://doi.org/10.1214/12-AOS1032>.
- 223 Z. Chen, Y. Zhou, Y. Liang, and Z. Lu. Generalized-smooth nonconvex optimization is as efficient as
224 smooth nonconvex optimization. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato,
225 and J. Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*,
226 volume 202 of *Proceedings of Machine Learning Research*, pages 5396–5427. PMLR, 23–29 Jul
227 2023. URL <https://proceedings.mlr.press/v202/chen23ar.html>.
- 228 M. Crawshaw, M. Liu, F. Orabona, W. Zhang, and Z. Zhuang. Robustness to unbounded smoothness
229 of generalized signsgd. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh,
230 editors, *Advances in Neural Information Processing Systems*, volume 35, pages 9955–9968. Curran
231 Associates, Inc., 2022. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2022/file/40924475a9bf768bdac3725e67745283-Paper-Conference.pdf)
232 [2022/file/40924475a9bf768bdac3725e67745283-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/40924475a9bf768bdac3725e67745283-Paper-Conference.pdf).
- 233 M. Faw, L. Rout, C. Caramanis, and S. Shakkottai. Beyond uniform smoothness: A stopped analysis
234 of adaptive sgd. In G. Neu and L. Rosasco, editors, *Proceedings of Thirty Sixth Conference*
235 *on Learning Theory*, volume 195 of *Proceedings of Machine Learning Research*, pages 89–160.
236 PMLR, 12–15 Jul 2023. URL <https://proceedings.mlr.press/v195/faw23a.html>.
- 237 E. Hazan, K. Levy, and S. Shalev-Shwartz. Beyond convexity: Stochastic quasi-convex
238 optimization. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors,
239 *Advances in Neural Information Processing Systems*, volume 28. Curran Associates,
240 Inc., 2015. URL [https://proceedings.neurips.cc/paper_files/paper/2015/file/](https://proceedings.neurips.cc/paper_files/paper/2015/file/934815ad542a4a7c5e8a2dfa04fea9f5-Paper.pdf)
241 [934815ad542a4a7c5e8a2dfa04fea9f5-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2015/file/934815ad542a4a7c5e8a2dfa04fea9f5-Paper.pdf).
- 242 H. Li, J. Qian, Y. Tian, A. Rakhlin, and A. Jadbabaie. Convex and non-convex optimization under
243 generalized smoothness. In *Thirty-seventh Conference on Neural Information Processing Systems*,
244 2023a. URL <https://openreview.net/forum?id=8aunGrXdK1>.
- 245 H. Li, A. Rakhlin, and A. Jadbabaie. Convergence of adam under relaxed assumptions. In
246 A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances*
247 *in Neural Information Processing Systems*, volume 36, pages 52166–52196. Curran Associates,
248 Inc., 2023b. URL [https://proceedings.neurips.cc/paper_files/paper/2023/file/](https://proceedings.neurips.cc/paper_files/paper/2023/file/a3cc50126338b175e56bb3cad134db0b-Paper-Conference.pdf)
249 [a3cc50126338b175e56bb3cad134db0b-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/a3cc50126338b175e56bb3cad134db0b-Paper-Conference.pdf).
- 250 C. Liu, L. Zhu, and M. Belkin. Loss landscapes and optimization in over-parameterized non-
251 linear systems and neural networks. *Applied and Computational Harmonic Analysis*, 59:85–
252 116, 2022. ISSN 1063-5203. doi: <https://doi.org/10.1016/j.acha.2021.12.009>. URL <https://www.sciencedirect.com/science/article/pii/S106352032100110X>. Special Issue on
253 Harmonic Analysis and Machine Learning.
254
- 255 H. Lu, R. M. Freund, and Y. Nesterov. Relatively smooth convex optimization by first-order methods,
256 and applications. *SIAM Journal on Optimization*, 28(1):333–354, 2018. doi: 10.1137/16M1099546.
257 URL <https://doi.org/10.1137/16M1099546>.
- 258 V. Patel, S. Zhang, and B. Tian. Global convergence and stability of stochastic gradient descent.
259 In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in*
260 *Neural Information Processing Systems*, volume 35, pages 36014–36025. Curran Associates,
261 Inc., 2022. URL [https://proceedings.neurips.cc/paper_files/paper/2022/file/](https://proceedings.neurips.cc/paper_files/paper/2022/file/ea05e4fc0299c27648c9985266abad47-Paper-Conference.pdf)
262 [ea05e4fc0299c27648c9985266abad47-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/ea05e4fc0299c27648c9985266abad47-Paper-Conference.pdf).
- 263 B. Wang, Y. Zhang, H. Zhang, Q. Meng, Z.-M. Ma, T.-Y. Liu, and W. Chen. Provable adaptivity in
264 Adam, 2022. arXiv:2208.09900.
- 265 B. Zhang, J. Jin, C. Fang, and L. Wang. Improved analysis of clipping algorithms for non-convex
266 optimization. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances*
267 *in Neural Information Processing Systems*, volume 33, pages 15511–15521. Curran Associates,
268 Inc., 2020a. URL [https://proceedings.neurips.cc/paper_files/paper/2020/file/](https://proceedings.neurips.cc/paper_files/paper/2020/file/b282d1735283e8eea45bce393cefe265-Paper.pdf)
269 [b282d1735283e8eea45bce393cefe265-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/b282d1735283e8eea45bce393cefe265-Paper.pdf).

270 J. Zhang, T. He, S. Sra, and A. Jadbabaie. Why gradient clipping accelerates training: A theoretical
 271 justification for adaptivity. In *International Conference on Learning Representations, 2020b*. URL
 272 <https://openreview.net/forum?id=BJgnXpVYwS>.

273 A Background and notation

274 A.1 Weak smoothness

275 **Definition 2.** Let $m > 0$. A function $L: \mathbb{R}^d \rightarrow \mathbb{R}$ is m -smooth if for every $\alpha, \beta \in \mathbb{R}^d$, we have

$$\|\nabla L(\beta) - \nabla L(\alpha)\| \leq m\|\alpha - \beta\|.$$

276 When L is twice differentiable, this can alternatively be stated as the condition that the magnitude of
 277 the second derivative of L is uniformly bounded by m .

278 Because many functions that one would like to minimize are not m -smooth for any m , researchers
 279 have long proposed weaker notions of smoothness and tried to prove convergence results under
 280 such alternate definitions of smoothness, in an effort to expand the range of functions that we can
 281 confidently optimize. We begin by noting a few of the popular definitions.

282 **Definition 3** (see Hazan et al. [2015]). Let $\tau, \epsilon > 0, \gamma \in \mathbb{R}^d$. A function $L: \mathbb{R}^d \rightarrow \mathbb{R}$ is (τ, ϵ, γ) -
 283 locally-smooth if for every $\alpha, \beta \in \mathbb{R}^d$ such that $\|\alpha - \gamma\| \leq \epsilon$ and $\|\beta - \gamma\| \leq \epsilon$, we have

$$|L(\beta) - L(\alpha) - \langle \nabla L(\beta), \alpha - \beta \rangle| \leq \frac{\tau}{2}\|\alpha - \beta\|^2.$$

284 **Definition 4** (see Agarwal et al. [2012]). Let $\tau, \epsilon > 0$ and let $R: \mathbb{R}^d \rightarrow \mathbb{R}^+$ be a regularizer. A
 285 function $L: \mathbb{R}^d \rightarrow \mathbb{R}$ satisfies restricted smoothness with respect to R with parameters (τ, ϵ) if for
 286 every $\alpha, \beta \in \mathbb{R}^d$, we have

$$|L(\beta) - L(\alpha) - \langle \nabla L(\beta), \alpha - \beta \rangle| \leq \frac{\tau}{2}\|\alpha - \beta\|^2 + \epsilon R^2(\alpha - \beta).$$

287 **Definition 5** (see Lu et al. [2018]). Let $h: \mathbb{R}^d \rightarrow \mathbb{R}$ be a differentiable convex “reference function”,
 288 and m a positive real number. A function $L: \mathbb{R}^d \rightarrow \mathbb{R}$ is m -smooth relative to h if for any $\alpha, \beta \in \mathbb{R}^d$
 289 we have

$$L(\beta) \leq L(\alpha) + \langle \nabla L(\alpha), \beta - \alpha \rangle + m(h(\beta) - h(\alpha) - \langle \nabla h(\alpha), \beta - \alpha \rangle).$$

290 A.2 Weak convexity

291 In the Introduction, we noted that there is interest both in weaker notions of smoothness and weaker
 292 notions of convexity, and in determining whether the loss functions arising from deep neural networks
 293 satisfy any of them.

294 While in this paper we will focus on smoothness conditions, here we point to work considering these
 295 questions for convexity conditions.

296 In Liu et al. [2022], Liu-Zhu-Belkin considered the PL* condition, a condition related to the classical
 297 Polyak-Lojasiewicz condition. They showed that for the loss function of neural networks, if the
 298 network satisfies some conditions including that they are sufficiently wide, one can construct many
 299 balls within the parameter space \mathbb{R}^d on which the PL* condition holds.

300 While it is known that the loss functions arising in deep learning are not convex, this result shows
 301 that there is a weaker type of convexity that is satisfied in some regions for some neural networks.

302 At this time, we do not know of analogous results for relaxed smoothness conditions. In this work, we
 303 will give a negative result, for most neural networks, it is not the case that the (L_0, L_1) -smoothness
 304 condition holds over the entire parameter space \mathbb{R}^d . Perhaps an analog of Liu-Zhu-Belkin’s result
 305 for the PL* convexity condition is possible - perhaps there are regions in the parameter space \mathbb{R}^d on
 306 which the (L_0, L_1) -smoothness condition holds. We leave that question to future work.

307 **A.3 Notation**

308 **A.3.1 Fully connected feedforward neural networks**

309 To define a fully connected feedforward neural network, we begin by specifying the number of layers
 310 ℓ of the network, and the widths $d_{\text{in}}, c_1, \dots, c_\ell, d_{\text{out}}$ of the layers, ordered from “earliest” to “latest”.
 311 For each adjacent pair of layers $i, i + 1$, we will have the space of affine linear maps from \mathbb{R}^{c_i} to
 312 $\mathbb{R}^{c_{i+1}}$. Such a map is given by the choice of a $c_i \times c_{i+1}$ matrix we will call M^i , and a vector in $\mathbb{R}^{c_{i+1}}$
 313 we will call b^i . The entries of M^i we call weights, the entries of b^i we call biases, and the choice of
 314 weights and biases for all the layers we call the choice of a parameter vector $\rho \in \mathbb{R}^p$.

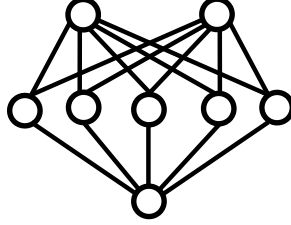


Figure 6: This is a diagram representing a feedforward neural network with one hidden layer, with the input width $d_{\text{in}} = 2$, the width of the hidden layer $k_1 = 5$, and output width $d_{\text{out}} = 1$.

315 Next, we choose an **activation function**

$$\sigma: \mathbb{R} \rightarrow \mathbb{R}. \quad (2)$$

316 In this paper, we assume that σ is twice differentiable.

317 Given the choices above of an architecture and σ , this neural network provides a way to input a set ρ
 318 of weights and biases for the network and output a function $f_\rho: \mathbb{R}^{d_{\text{in}}} \rightarrow \mathbb{R}^{d_{\text{out}}}$.

319 Given a vector

$$\rho = (\mathbf{w}, \mathbf{b}) \in \mathbb{R}^p \quad (3)$$

320 in the parameter space, we define the function

$$f_\rho = f_{\mathbf{w}, \mathbf{b}}: \mathbb{R}^{d_{\text{in}}} \rightarrow \mathbb{R}^{d_{\text{out}}} \quad (4)$$

321 by composing the following sequence of maps specified by the neural network and the choice of the
 322 weights and biases in all the layers \mathbf{w}, \mathbf{b} :

$$\mathbb{R}^{d_{\text{in}}} \xrightarrow{M^1 \mathbf{x} + \mathbf{b}^1} \mathbb{R}^{c_1} \xrightarrow{\sigma} \mathbb{R}^{c_1} \xrightarrow{M^2 \mathbf{x} + \mathbf{b}^2} \dots \xrightarrow{M^\ell \mathbf{x} + \mathbf{b}^\ell} \mathbb{R}^{c_\ell} \xrightarrow{\sigma} \mathbb{R}^{c_\ell} \xrightarrow{M^{\ell+1} \mathbf{x} + \mathbf{b}^{\ell+1}} \mathbb{R}^{d_{\text{out}}}. \quad (5)$$

323 In this construction, the arrow $\mathbb{R}^{k_i} \xrightarrow{\sigma} \mathbb{R}^{k_i}$ indicates that we apply σ componentwise.

324 *Example 1.* Consider a fully connected feedforward graph with layers of widths 1, 3, 1 and $\sigma =$
 325 $(u \mapsto u^2 + 1)$. The corresponding function space consists of those functions of the form

$$f_\alpha: x \mapsto w_{11}^2 ((w_{11}^1 x + b_1^1)^2 + 1) + w_{12}^2 ((w_{21}^1 x + b_2^1)^2 + 1) + w_{13}^2 ((w_{31}^1 x + b_3^1)^2 + 1) + b_1^2. \quad (6)$$

326 In our calculations we will find it useful to have the following notation for the stages of the neural
 327 network. We define recursively

$$f_\rho^1(\mathbf{x}) = M^1 \mathbf{x} + \mathbf{b}^1 \quad (7)$$

328

$$f_\rho^i(\mathbf{x}) = M^i \sigma(f_\rho^{i-1}(\mathbf{x})) + \mathbf{b}^i, \quad (8)$$

329 so that the previously defined function $f_\rho(\mathbf{x})$ equals $f_\rho^{\ell+1}(\mathbf{x})$.

330 **A.3.2 The loss function L**

331 In deep learning, one starts with a data set, chooses an architecture for a neural network, and then
 332 wishes to find a parameter vector, in other words a set of weights and biases for the network, such
 333 that with that choice, the function expressed by the network predicts well on similar data.

334 To find such a parameter vector, a key step is to define a loss function

$$L: \mathbb{R}^d \rightarrow \mathbb{R} \tag{9}$$

335 from the set of all parameters to the real numbers. This function L is constructed in such a way that
 336 parameter vectors ρ on which the loss function achieves a low value are good choices for the network.
 337 In today's implementations, a gradient descent based algorithm is used to find such ρ that minimize
 338 L .

339 In this paper, we consider two ways of constructing L , and we define each in this section. In both
 340 cases, we fix

- 341 • a neural network,
- 342 • a choice of activation function σ ,
- 343 • and a data set

$$D := \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i \in \{1, \dots, n\}} \subset \mathbb{R}^{d_{\text{in}}} \times \mathbb{R}^{d_{\text{out}}}. \tag{10}$$

344 **Definition 6.** The L_2 loss is defined by:

$$L(\rho) := \sum_{s=1}^n (f_{\alpha}(\mathbf{x}_s) - \mathbf{y}_s)^2. \tag{11}$$

345 **Definition 7.** The cross-entropy loss is defined by:

$$L(\rho) = - \sum_{m=1}^n \sum_{j=1}^{d_{\text{out}}} [y_m]_j \log \frac{e^{[f_{\rho}(x_m)]_j}}{\sum_{k=1}^b e^{[f_{\rho}(x_m)]_k}}. \tag{12}$$

346 where $[v]_j$ denotes the j^{th} entry of a vector v .

347 **B Supporting material**

348 First, the description of the line \mathcal{R} appearing in Section 3.1.

349 We define \mathcal{R} to be the image of the following linear map $\rho: \mathbb{R} \rightarrow \mathbb{R}^d$, where \mathbb{R}^d is the parameter
 350 space of the neural network. Given a real number $t \in \mathbb{R}$, $\rho(t)$ is the following choice of weights and
 351 biases.

352 For all but the last layer, we take the weights to be zero and the biases all equal to zero.

$$M^i = 0, \quad b^i = (c \quad \dots \quad c)^T \quad \text{if } 1 \leq i \leq \ell$$

353 In the last layer, we take $M^{\ell+1}$ to be t times a constant matrix M . We choose M carefully, depending
 354 on the loss. Finally, we take all the biases equal to 0.

$$M = \begin{pmatrix} m_{11} & \dots & m_{1k} \\ \vdots & \ddots & \vdots \\ m_{b1} & \dots & m_{bk} \end{pmatrix} \tag{13}$$

$$M^{\ell+1} = tM, \quad b^{\ell+1} = \vec{0} \tag{14}$$

355 Now, the randomly chosen data x and y for the experiment in Section 3.2.

$$\vec{x} = \begin{pmatrix} -0.51 & -0.32 & 0.45 & -0.42 & -0.51 & -0.30 & 0.50 & -0.98 & -0.22 & -0.51 \\ & 0.40 & 0.79 & 0.20 & 0.96 & 0.90 & -0.04 & 0.60 & -0.27 & 0.01 & 0.41 \end{pmatrix}$$

356

$$\vec{y} = \begin{pmatrix} -0.70 & -0.86 & -0.20 & 0.88 & 0.88 & -0.33 & -0.92 & 0.06 & 0.89 & 0.21 \\ -0.35 & 0.32 & 0.27 & 0.97 & 0.86 & -0.20 & 0.49 & -0.05 & -0.75 & 0.90 \end{pmatrix}$$

357 $start_1 =$

$$\begin{pmatrix} 0.048 & 0.026 & -0.0063 & 0.025 & -0.014 & -0.0014 & -0.059 & 0.024 \\ -0.026 & 0.0021 & -0.022 & 0.013 & -0.012 & -0.0093 & -0.021 & 0.0052 \\ -0.026 & -0.03 & -0.022 & -0.014 & -0.011 & 0.018 & 0.00088 & -0.024 \\ -0.023 & -0.021 & 0.026 & 0.019 & -0.019 & 0.025 & 0.0017 & -0.028 \\ 0.012 & 0.014 & -0.0012 & 0.012 & -0.017 & -0.019 & 0.022 & 0.022 \\ 0.009 & 0.015 & -0.027 & 0.02 & -0.01 & -0.0078 & -0.0032 & -0.016 \\ -0.0082 & -0.0049 & -0.0077 \end{pmatrix}$$

358 $end_1 =$

$$\begin{pmatrix} 9.6 & 5.2 & -1.3 & 4.9 & -2.8 & -0.29 & -12 & 4.8 & -5.3 & 0.43 \\ -4.3 & 2.7 & -2.5 & -1.9 & -4.2 & 1 & -5.2 & -6 & -4.4 & -2.8 \\ -2.2 & 3.6 & 0.18 & -4.9 & -4.6 & -4.3 & 5.3 & 3.8 & -3.8 & 5.1 \\ 0.34 & -5.7 & 2.3 & 2.9 & -0.25 & 2.3 & -3.3 & -3.8 & 4.5 & 4.5 \\ 1.8 & 2.9 & -5.4 & 4 & -2.1 & -1.6 & -0.65 & -3.3 & -1.7 & -0.98 & -1.6 \end{pmatrix}$$

359 $start_2 =$

$$\begin{pmatrix} 0.068 & -0.059 & -0.035 & -0.034 & 0.096 & -0.055 & 0.055 & -0.02 \\ -0.03 & -0.007 & 0.036 & 0.048 & -0.027 & -0.04 & -0.029 & 0.0052 \\ -0.041 & -0.014 & -0.019 & 0.025 & 0.017 & 0.042 & 0.017 & -0.018 \\ -0.044 & -0.022 & 0.046 & 0.041 & 0.021 & 0.046 & 0.049 & -0.041 \\ 0.016 & 0.022 & 0.024 & 0.017 & 0.036 & -0.0045 & 0.0011 & 0.048 \\ -0.033 & -0.05 & 0.041 & -0.013 & -0.026 & 0.018 & 0.031 & 0.013 \\ -0.024 & 0.019 & 0.015 \end{pmatrix}$$

360 $end_2 =$

$$\begin{pmatrix} 14 & -12 & -7.1 & -6.8 & 19 & -11 & 11 & -4 & -6 & -1.4 \\ 7.2 & 9.6 & -5.4 & -8 & -5.8 & 1 & -8.2 & -2.8 & -3.9 & 5 \\ 3.5 & 8.4 & 3.4 & -3.6 & -8.9 & -4.5 & 9.3 & 8.2 & 4.3 & 9.2 \\ 9.9 & -8.1 & 3.3 & 4.4 & 4.9 & 3.5 & 7.2 & -0.91 & 0.23 & 9.7 \\ -6.6 & -10 & 8.3 & -2.7 & -5.3 & 3.6 & 6.2 & 2.6 & -4.8 & 3.9 & 3 \end{pmatrix}$$

361 $start_3 =$

$$\begin{pmatrix} 0.085 & 0.052 & 0.4 & 0.021 & 0.24 & 0.23 & 0.29 & -0.064 \\ -0.098 & 0.024 & 0.079 & 0.12 & -0.14 & 0.12 & 0.13 & 0.042 \\ -0.083 & -0.076 & -0.047 & -0.055 & 0.13 & 0.02 & 0.097 & -0.0038 \\ -0.19 & -0.015 & 0.088 & -0.048 & 0.089 & -0.16 & -0.13 & 0.013 \\ -0.12 & -0.011 & 0.2 & 0.13 & 0.053 & -0.2 & 0.062 & -0.014 \\ 0.15 & -0.013 & -0.035 & 0.042 & -0.0064 & -0.14 & -0.087 & -0.14 \\ -0.099 & 0.082 & -0.024 \end{pmatrix}$$

362 $end_3 =$

$$\begin{pmatrix} 17 & 10 & 80 & 4.3 & 49 & 46 & 59 & -13 & -20 & 4.7 \\ 16 & 24 & -28 & 24 & 26 & 8.5 & -17 & -15 & -9.5 & -11 \\ 25 & 4 & 20 & -0.76 & -37 & -3 & 18 & -9.7 & 18 & -32 \\ -27 & 2.7 & -24 & -2.3 & 40 & 26 & 11 & -39 & 12 & -2.9 \\ 30 & -2.6 & -6.9 & 8.5 & -1.3 & -28 & -17 & -28 & -20 & 17 & -4.8 \end{pmatrix}$$

363 $start_4 =$

$$\begin{pmatrix} -16 & 27 & 5.7 & 0.032 & 15 & 42 & -35 & -22 \\ -5.9 & 20 & 19 & 0.22 & -13 & -13 & -5.4 & -15 \\ 21 & 3.8 & 10 & 17 & -6 & -21 & 17 & -11 \\ 11 & -15 & 4.4 & 13 & 6.5 & -15 & -11 & 0.92 \\ -22 & -16 & 3.1 & -0.93 & 16 & 0.43 & 19 & -12 \\ -16 & -11 & 15 & 17 & -13 & 8.3 & -0.28 & -12 \\ 11 & 12 & 13 \end{pmatrix}$$

364 $end_4 =$

$$\begin{pmatrix} -35 & 60 & 13 & 0.071 & 32 & 92 & -78 & -49 & -13 & 45 \\ 42 & 0.5 & -28 & -28 & -12 & -34 & 46 & 8.4 & 23 & 38 \\ -13 & -48 & 37 & -24 & 25 & -32 & 9.8 & 28 & 14 & -34 \\ -24 & 2 & -48 & -35 & 6.8 & -2.1 & 36 & 0.96 & 42 & -27 \\ -36 & -25 & 32 & 38 & -29 & 18 & -0.62 & -26 & 24 & 27 & 28 \end{pmatrix}$$

365 $start_5 =$

$$\begin{pmatrix} 0.37 & -0.29 & 0.045 & 0.35 & -0.16 & 0.2 & 0.16 & 0.14 \\ -0.047 & -0.11 & 0.16 & -0.15 & 0.066 & 0.11 & -0.22 & -0.095 \\ 0.053 & 0.028 & 0.068 & 0.22 & -0.13 & -0.22 & -0.092 & 0.023 \\ -0.097 & -0.11 & -0.22 & -0.19 & 0.17 & 0.014 & 0.18 & 0.18 \\ 0.011 & 0.015 & -0.13 & -0.17 & 0.22 & 0.051 & 0.17 & 0.14 \\ 0.13 & -0.071 & 0.22 & 0.054 & -0.084 & 0.036 & 0.037 & 0.035 \\ -0.1 & -0.16 & 0.14 \end{pmatrix}$$

366 $end_5 =$

$$\begin{pmatrix} -33 & 5.1 & 39 & -18 & 22 & 18 & 16 & -5.3 & -12 & 18 \\ -17 & 7.5 & 13 & -25 & -11 & 6 & 3.1 & 7.7 & 25 & -15 \\ -25 & -10 & 2.6 & -11 & -12 & -25 & -21 & 19 & 1.6 & 20 \\ 21 & 1.3 & 1.7 & -15 & -20 & 25 & 5.8 & 19 & 16 & 14 \\ -8 & 25 & 6.2 & -9.6 & 4.1 & 4.2 & 4 & -11 & -18 & 16 \end{pmatrix}$$

367 $start_6 =$

$$\begin{pmatrix} 0.27 & -0.38 & 0.37 & 0.35 & -0.31 & -0.055 & 0.19 & 0.12 \\ 0.031 & 0.16 & 0.1 & -0.00026 & 0.11 & 0.18 & 0.16 & 0.0043 \\ 0.14 & 0.073 & -0.13 & -0.0036 & -0.17 & 0.055 & -0.14 & 0.11 \\ -0.13 & -0.21 & -0.027 & 0.0049 & -0.14 & 0.068 & 0.12 & -0.12 \\ -0.093 & -0.083 & -0.086 & 0.11 & -0.2 & 0.029 & 0.2 & 0.22 \\ -0.12 & -0.2 & -0.093 & -0.074 & 0.0041 & -0.094 & -0.017 & 0.016 \\ -0.032 & -0.12 & -0.053 \end{pmatrix}$$

368 $end_6 =$

$$\begin{pmatrix} 30 & -42 & 41 & 39 & -35 & -6.1 & 21 & 14 & 3.5 & 18 \\ 11 & -0.029 & 12 & 20 & 18 & 0.48 & 16 & 8.2 & -14 & -0.4 \\ -19 & 6.1 & -16 & 12 & -15 & -24 & -3 & 0.55 & -16 & 7.7 \\ 14 & -13 & -10 & -9.3 & -9.7 & 12 & -22 & 3.3 & 22 & 25 \\ -14 & -22 & -10 & -8.3 & 0.46 & -11 & -1.9 & 1.8 & -3.6 & -13 & -5.9 \end{pmatrix}$$

369 $start_7 =$

$$\begin{pmatrix} 0.34 & 0.09 & -0.37 & -0.33 & 0.21 & -0.21 & 0.17 & -0.38 \\ 0.13 & -0.1 & -0.021 & -0.099 & -0.12 & -0.16 & 0.14 & 0.15 \\ 0.11 & -0.0072 & 0.073 & -0.12 & -0.039 & 0.18 & 0.18 & -0.15 \\ -0.16 & -0.18 & 0.12 & -0.18 & -0.016 & -0.049 & 0.024 & -0.01 \\ 0.044 & 0.13 & 0.12 & -0.073 & -0.12 & -0.077 & -0.19 & -0.11 \\ 0.15 & -0.023 & -0.12 & -0.12 & 0.089 & 0.15 & 0.038 & -0.061 \\ -0.096 & -0.0099 & 0.082 \end{pmatrix}$$

370 $end_7 =$

$$\begin{pmatrix} 33 & 8.9 & -37 & -32 & 21 & -21 & 17 & -38 & 13 & -10 \\ -2 & -9.7 & -12 & -15 & 14 & 15 & 11 & -0.71 & 7.2 & -11 \\ -3.9 & 18 & 18 & -15 & -16 & -18 & 12 & -18 & -1.6 & -4.8 \\ 2.3 & -1 & 4.4 & 12 & 11 & -7.2 & -12 & -7.6 & -19 & -11 \\ 15 & -2.2 & -11 & -12 & 8.7 & 14 & 3.8 & -6 & -9.5 & -0.98 & 8.1 \end{pmatrix}$$

371 The random vector

$$w_2 = \begin{pmatrix} -0.87 & -0.67 & 0.5 & -0.29 & -0.93 & 0.83 & -0.004 & 0.64 \\ 0.5 & 0.46 & -0.46 & 0.47 & 0.046 & -0.29 & -0.023 & 0.38 \\ -0.22 & -0.27 & 0.41 & -0.1 & 0.032 & -0.14 & -0.28 & -0.31 \\ 0.23 & 0.18 & -0.12 & -0.066 & 0.2 & 0.058 & -0.36 & 0.37 \\ 0.31 & 0.22 & -0.17 & -0.27 & -0.11 & -0.34 & 0.095 & 0.34 \\ -0.31 & -0.26 & 0.18 & 0.24 & -0.29 & 0.23 & 0.022 & -0.012 \\ -0.28 & 0.28 & 0.36 \end{pmatrix}$$

372 The random vector

$$w_4 = \begin{pmatrix} -0.72 & 0.74 & -0.86 & -0.19 & -0.15 & -0.094 & -0.23 & 0.41 \\ -0.32 & 0.35 & -0.095 & 0.46 & 0.39 & 0.1 & 0.037 & -0.38 \\ -0.3 & -0.39 & -0.44 & 0.1 & -0.21 & 0.084 & -0.5 & 0.12 \\ 0.31 & 0.22 & -0.41 & -0.29 & 0.057 & -0.39 & -0.048 & 0.091 \\ -0.094 & -0.079 & -0.44 & 0.31 & 0.27 & 0.16 & 0.35 & 0.48 \\ -0.22 & 0.23 & -0.12 & 0.3 & -0.078 & -0.31 & -0.13 & 0.31 \\ 0.011 & -0.2 & 0.34 \end{pmatrix}$$

373 The random vector

$$w_6 = \begin{pmatrix} 0.64 & -0.26 & -0.76 & 0.37 & 0.45 & 0.82 & 0.6 & 0.2 \\ -0.34 & -0.4 & -0.23 & -0.2 & -0.08 & -0.099 & 0.017 & 0.35 \\ -0.36 & -0.3 & 0.18 & -0.13 & -0.062 & 0.33 & 0.37 & -0.37 \\ 0.23 & -0.13 & -0.16 & -0.23 & 0.37 & 0.31 & 0.0082 & -0.047 \\ 0.37 & -0.18 & -0.32 & 0.037 & -0.25 & 0.44 & 0.088 & 0.2 \\ -0.048 & 0.17 & 0.44 & 0.12 & -0.11 & 0.23 & 0.11 & 0.22 \\ -0.063 & 0.14 & 0.12 \end{pmatrix}$$