
Invariant Low-Dimensional Subspaces in Gradient Descent for Learning Deep Matrix Factorizations

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 An extensively studied phenomenon of the past few years in training deep networks
2 is the implicit bias of gradient descent towards parsimonious solutions. In this
3 work, we further investigate this phenomenon by narrowing our focus to deep
4 matrix factorization, where we reveal surprising low-dimensional structures in the
5 learning dynamics when the target matrix is low-rank. Specifically, we show that
6 the evolution of gradient descent starting from arbitrary orthogonal initialization
7 only affects a minimal portion of singular vector spaces across all weight matrices.
8 In other words, the learning process happens only within a small invariant subspace
9 of each weight matrix, despite the fact that all parameters are updated throughout
10 training. From this, we provide rigorous justification for low-rank training in a
11 specific, yet practical setting. In particular, we demonstrate that we can construct
12 compressed factorizations that are equivalent to full-width, deep factorizations
13 throughout training for solving low-rank matrix completion problems efficiently.

14 1 Introduction

15 In recent years, deep learning has demonstrated remarkable success across a wide range of applications
16 [1]. Many recent works attempt to explain the exceptional generalization capabilities of deep networks
17 by studying the implicit bias of gradient-based methods, showing that deep networks trained with
18 such algorithms tend to learn simple functions [2–6]. Similarly, it has been shown that gradient
19 descent induces max-margin [6, 7] or low-rank solutions [8–12] in deep networks, to name a few.

20 In another vein, recent work has explored the increasingly important problem of training deep
21 networks more efficiently via *low-rank training* [13–17], where the number of trainable parameters is
22 effectively reduced by replacing the original network weights with low-rank factorizations. While such
23 methods have shown promising empirical results for reducing training time, theoretical justifications
24 for these approaches remain deficient. The aforementioned works on implicit bias characterize
25 low-rank structure in the limit of gradient descent – they do not address whether the trajectory of
26 the original overparameterized network (along with its generalization/convergence properties) is
27 achievable via low-rank factorization *from initialization* throughout training, which is what low-rank
28 training necessitates.

29 **Contributions.** In this work, we draw theoretical connections between the implicit bias of gradient
30 descent in deep networks and the practice of low-rank training. Utilizing deep matrix factorizations
31 as a testbed (commonly assumed for analyzing the complex optimization dynamics of deep networks
32 [10, 18–20]) we demonstrate that for low-rank data, all weight matrices are only updated within
33 a low-dimensional subspace that is *invariant throughout training*, which can be determined from
34 their arbitrary orthogonal initialization. To our knowledge, this is the first work identifying such
35 invariant structures in gradient descent dynamics from random initialization. From this, we show
36 that we can construct a compressed, low-rank factorization that is nearly equivalent to the original

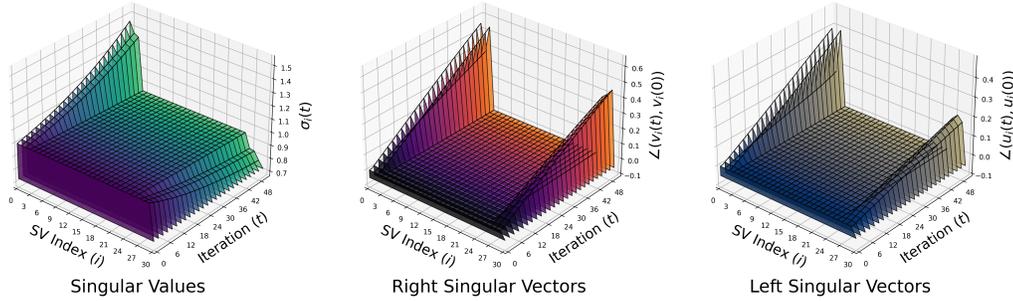


Figure 1: **Evolution of SVD of weight matrices.** We visualize the SVD dynamics of the first layer weight matrix of an $L = 3$ layer deep matrix factorization with $d = 30$, $r = 3$, $\sigma_l = 1$ throughout GD without weight decay. *Left:* Magnitude of the i -th singular value $\sigma_i(t)$ at iteration t . *Middle:* Angle $\angle(\mathbf{v}_i(t), \mathbf{v}_i(0))$ between the i -th right singular vector at iteration t and initialization. *Right:* Angle $\angle(\mathbf{u}_i(t), \mathbf{u}_i(0))$ between the i -th left singular vector at iteration t and initialization.

37 overparameterized network, thereby providing some rigorous foundations for low-rank training.
 38 Although deep matrix factorizations are mostly of theoretical interest, they are adopted for low-rank
 39 matrix sensing problems - therefore, we demonstrate that our theory can be applied (with slight
 40 modifications) towards accelerating practical low-rank deep matrix completion problems.

41 2 Analysis

42 **Setup.** We study the training dynamics of L -layer deep matrix factorizations $f(\Theta)$ given by

$$f(\Theta) := \mathbf{W}_L \mathbf{W}_{L-1} \cdots \mathbf{W}_2 \mathbf{W}_1$$

43 where $\Theta = (\mathbf{W}_l)_{l=1}^L$ are the parameters or weights with $\mathbf{W}_l \in \mathbf{R}^{d_l \times d_{l-1}}$ for $l \in [L]$. For a given
 44 target matrix $\Phi \in \mathbf{R}^{d \times d}$, we learn parameters Θ with $d_0 = d_1 = \cdots = d_L = d$ by minimizing the
 45 square loss

$$\ell(\Theta) = \frac{1}{2} \|f(\Theta) - \Phi\|_F^2 \quad (1)$$

46 via gradient descent (GD) from scaled *orthogonal* initialization, i.e., we initialize parameters $\Theta(0)$
 47 such that all singular values of $\mathbf{W}_l(0)$ are equal to some $\sigma_l > 0$ for each $l \in [L]$. Then, we update all
 48 weights for $t = 0, 1, 2, \dots$ as

$$\mathbf{W}_l(t+1) = (1 - \eta\lambda)\mathbf{W}_l(t) - \eta \nabla_{\mathbf{w}_l} \ell(\Theta(t)), \quad l \in [L] \quad (2)$$

49 where $\lambda \geq 0$ is an optional weight decay parameter and $\eta > 0$ is the learning rate.

50 **Main Result.** Under the setting described above, we show learning only occurs within an invariant
 51 low-dimensional subspace of the weight matrices, provided that the target matrix Φ is low-rank.

52 **Theorem 1.** *Suppose $\Phi \in \mathbf{R}^{d \times d}$ is at most rank r where $m := d - 2r > 0$. Then there exist*
 53 *orthogonal matrices $(\mathbf{U}_l)_{l=1}^L \subset \mathbf{R}^{d \times d}$ and $(\mathbf{V}_l)_{l=1}^L \subset \mathbf{R}^{d \times d}$ satisfying $\mathbf{V}_{l+1} = \mathbf{U}_l$ for $l \in [L - 1]$,*
 54 *such that $\mathbf{W}_l(t)$ admits the decomposition*

$$\mathbf{W}_l(t) = \mathbf{U}_l \begin{bmatrix} \widetilde{\mathbf{W}}_l(t) & \mathbf{0} \\ \mathbf{0} & \rho_l(t) \mathbf{I}_m \end{bmatrix} \mathbf{V}_l^\top \quad (3)$$

55 for all $l \in [L]$ and $t \geq 0$, where $\widetilde{\mathbf{W}}_l(t) \in \mathbf{R}^{2r \times 2r}$ with $\mathbf{W}_l(0) = \sigma_l \mathbf{I}_{2r}$, and

$$\rho_l(t) = \rho_l(t-1) \cdot (1 - \eta\lambda - \eta \cdot \prod_{k \neq l} \rho_k^2(t-1)) \quad (4)$$

56 for all $l \in [L]$ and $t \geq 1$ with $\rho_l(0) = \sigma_l$.

57 We defer the proof of Theorem 1 to Appendix A.1. In the following, we discuss several implications
 58 of our result and its relationship to previous work.

59 • **SVD dynamics of weight matrices.** The decomposition (3) implies that $\mathbf{W}_l(t)$ has m identical
 60 singular values that follow the updates given in (4), whose corresponding singular vectors are
 61 stationary from initialization throughout GD – this is portrayed in Figure 1. By this, we can
 62 decompose the total space \mathbf{R}^d into two invariant singular subspaces: a $2r$ -dimensional space within
 63 which learning takes place, and an m -dimensional space corresponding to repeated singular values.

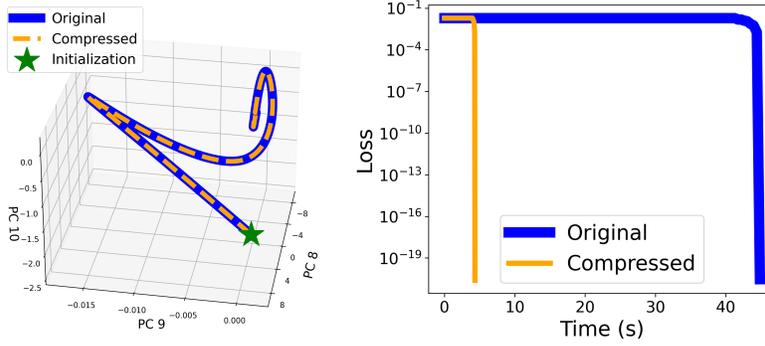


Figure 2: **Network compression for deep matrix factorization.** Comparison of trajectories for optimizing the original problem (1) vs. the compressed problem (6) with $L = 3$, $d = 1000$, $\hat{r} = r = 5$, and $\sigma_l = 10^{-3}$. *Left:* Principal components of end-to-end GD trajectories. *Right:* Training loss vs. wall-time comparison.

- 64 • **Low-rank bias.** From (4), we see that the GD trajectory either remains or tends towards a rank of
 65 at most $2r$ when we employ implicit or explicit regularization respectively. Indeed, if we use small
 66 initialization $\sigma_l \approx 0$, then the fact that ρ_l is a decreasing sequence implies that $\mathbf{W}_l(t)$ can be no
 67 more than rank $2r$ throughout the entire trajectory; whereas if $\lambda > 0$, then $\rho_l(t) \rightarrow 0$ as $t \rightarrow \infty$,
 68 which forces $\mathbf{W}_l(t)$ towards a solution of rank at most $2r$, regardless of initialization.
- 69 • **Comparison to prior arts.** In contrast to existing work that demonstrates the tendency of GD to
 70 find low nuclear-norm solutions [9, 11], our result directly shows that GD tends to find low-rank
 71 solutions. Moreover, unlike previous work on implicit bias [11, 21–23], we carefully examine the
 72 effect of weight decay, which is commonly employed during the training of deep networks. We
 73 note that our analysis is distinct from that of [18, 19], where continuous time dynamics are studied
 74 with the special (separable) setting $\mathbf{W}_{L:1}(0) = \mathbf{U}\mathbf{V}^\top$ with $\Phi = \mathbf{U}\Sigma\mathbf{V}^\top$. In contrast, our result
 75 applies to discrete time GD and holds for initialization that is agnostic to the target matrix. We also
 76 note that our result is unrelated to balanced initialization (as in [24]), since the σ_l can be arbitrarily
 77 different from one another.

78 **Compressed Deep Matrix Factorization.** We now show that, as a consequence of Theorem 1, we
 79 can run gradient descent on dramatically fewer parameters to achieve a near *identical* end-to-end
 80 trajectory to the original (full-width) factorization. More specifically, given an initialization $\Theta(0)$ of
 81 the original parameters and an upper bound on the rank $\hat{r} \geq r$ such that $d - 2\hat{r} > 0$, we define the
 82 *compressed* factorization

$$\hat{f}(\hat{\Theta}, \mathbf{U}_{L,1}, \mathbf{V}_{1,1}) := \mathbf{U}_{L,1} \hat{\mathbf{W}}_L \hat{\mathbf{W}}_{L-1} \cdots \hat{\mathbf{W}}_1 \mathbf{V}_{1,1}^\top \quad (5)$$

83 where $\hat{\Theta} = (\hat{\mathbf{W}}_l)_{l=1}^L$ are compressed weights with $\hat{\mathbf{W}}_l \in \mathbf{R}^{2\hat{r} \times 2\hat{r}}$ and $\mathbf{U}_{L,1}, \mathbf{V}_{1,1} \in \mathbf{R}^{d \times 2\hat{r}}$ are the
 84 first $2\hat{r}$ columns of $\mathbf{U}_L, \mathbf{V}_1 \in \mathbf{R}^{d \times d}$ respectively from Theorem 1 (depends on $\Theta(0)$ and Φ). Then,
 85 initializing $\hat{\Theta}(0)$ such that $\hat{\mathbf{W}}_l(0) = \mathbf{U}_{l,1}^\top \mathbf{W}_l(0) \mathbf{V}_{l,1}$ for all $l \in [L]$ and running gradient descent on
 86 the loss

$$\hat{\ell}(\hat{\Theta}) = \frac{1}{2} \|\hat{f}(\hat{\Theta}, \mathbf{U}_{L,1}, \mathbf{V}_{1,1}) - \Phi\|_F^2 \quad (6)$$

87 gives an almost equivalent network in the following sense.

88 **Proposition 1.** For $\hat{r} \geq r$ such that $\hat{m} := d - 2\hat{r} > 0$, running gradient descent on the compressed
 89 weights $\hat{\Theta}$ as described above for the loss (6) satisfies

$$\left\| f(\Theta(t)) - \hat{f}(\hat{\Theta}(t), \mathbf{U}_{L,1}, \mathbf{V}_{1,1}) \right\|_F^2 \leq \hat{m} \cdot \prod_{l=1}^L \sigma_l^2$$

90 for all iterates $t = 0, 1, 2, \dots$

91 We defer the proof of Proposition 1 to Appendix A.2. When we start from small initialization ($\sigma_l \approx 0$),
 92 Proposition 1 demonstrates that we only need to optimize $4L \cdot \hat{r}^2$ many parameters as opposed to the
 93 original $L \cdot d^2$ number of parameters to achieve an almost identical end-to-end trajectory, see Figure 2
 94 (left). Since it is often the case that $r \leq \hat{r} \ll d$, this results in an order of magnitude reduction in
 95 time to reach an optimal solution compared to the original network, see Figure 2 (right). In the next
 96 section, we demonstrate how this idea can be leveraged (with slight modification) to accelerate a
 97 more practical problem.

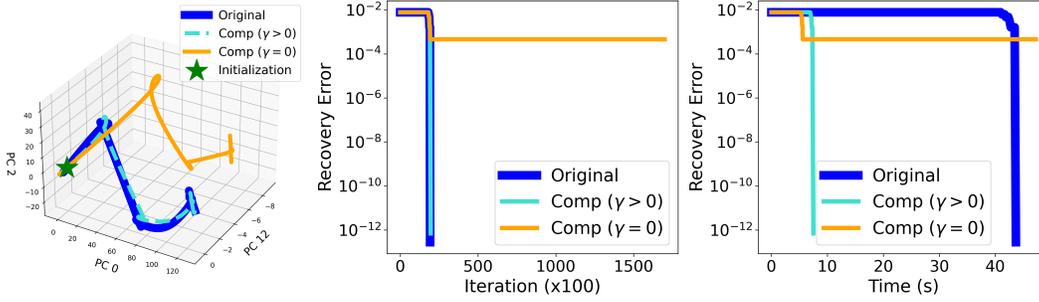


Figure 3: **Network compression for deep matrix completion.** Comparison of trajectories for optimizing the original problem (7) vs. the compressed problem (8) with γ discrepant updates ($\gamma = 0.01$) and ablating γ ($\gamma = 0$) with $L = 3$, $d = 1000$, $r = 5$, $\sigma_l = 10^{-3}$ and 20% of entries observed. *Left*: Principal components of end-to-end trajectories of each factorization. *Middle*: Recovery error vs. iteration comparison. *Right*: Recovery error vs wall-time comparison.

98 3 Application: Accelerating Deep Low-Rank Matrix Completion

99 **Problem Setup.** We consider the low-rank matrix completion problem [25–27] with ground-truth
 100 $\Phi \in \mathbf{R}^{d \times d}$ with rank $r \ll d$, where the goal is to recover Φ from only a few number of observations
 101 encoded by a mask $\Omega \in \{0, 1\}^{d \times d}$. Adopting a deep matrix factorization approach [11], we minimize
 102 the objective

$$\ell_{\text{mc}}(\Theta) = \frac{1}{2} \|\Omega \odot (f(\Theta) - \Phi)\|_F^2 \quad (7)$$

103 which simplifies to (1) when $\Omega = \mathbf{1}_d \mathbf{1}_d^\top$ in the full observation case. In practice, the true rank r is
 104 not known – instead, we assume to have an upper bound \hat{r} of the same order as r , i.e., $r \leq \hat{r} \ll d$.

105 **Compressed Deep Matrix Completion.** In the setting described above, it is advantageous to *over-*
 106 *parameterize* along both the depth and width of the factorization, particularly for accelerating GD
 107 convergence to well-generalizing solutions – see Appendix B for a more detailed discussion alongside
 108 evidence. Nonetheless, the advantages of over-parameterization are hindered by the fact that depth
 109 and width incur much higher *per-iteration* costs – for an L -layer factorization of (full)-width d , we
 110 require $O(L \cdot d^3)$ multiplications to evaluate gradients, where d is often very large. However, using
 111 ideas from the previous section, we can effectively reduce the computation to $O(\hat{r}^2 \cdot (L\hat{r} + d))$
 112 multiplications via a compressed factorization that emulates the trajectory of a (full)-width d network,
 113 thereby enjoying accelerated GD convergence with heavily reduced per-iteration computational cost.

114 In the full observation case ($\Omega = \mathbf{1}_d \mathbf{1}_d^\top$), we have already seen via Proposition 1 that the compressed
 115 factorization (5) with small initialization stays close to the trajectory of the full-width factorization.
 116 However, applying this directly to the projection $\Omega \odot \Phi$ will result in the compressed factorization’s
 117 trajectory diverging from that of the original – see the orange trace in Figure 3. Intuitively, this
 118 is because the factors $U_{L,1}, V_{1,1}$ are initialized from incomplete measurement of Φ – instead, we
 119 optimize the modified objective

$$\hat{\ell}_{\text{mc}}(\hat{\Theta}, U_{L,1}, V_{1,1}) = \frac{1}{2} \|\Omega \odot (\hat{f}(\hat{\Theta}, U_{L,1}, V_{1,1}) - \Phi)\|_F^2 \quad (8)$$

120 where $\hat{\Theta}$ are updated with learning rate η while the $U_{L,1}, V_{1,1}$ factors are updated with a *discrepant*
 121 learning rate $\gamma\eta$ where $\gamma > 0$ is small. While this results in an additional $2d\hat{r}$ parameters to be
 122 tracked, the trajectory of this compressed factorization will ultimately align with that of the original
 123 while converging roughly $5\times$ faster w.r.t. wall-time, as demonstrated in Figure 3. Moreover, the
 124 accelerated convergence induced by the full-width trajectory results in the compressed factorization
 125 being $3\times$ faster than randomly initialized factorizations of similar width – see Appendix C for more
 126 details.

127 4 Conclusion

128 This paper offers novel insights into simple structures in gradient descent for learning deep matrix
 129 factorizations, by which we derive some rigorous justification for the practice of low-rank training.
 130 Through this work, we hope to ultimately inspire more principled approaches to designing efficient
 131 and effective deep models by exploiting low-dimensional aspects of their training dynamics.

References

- 132
- 133 [1] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- 134 [2] Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The pitfalls of
135 simplicity bias in neural networks. *Advances in Neural Information Processing Systems*, 33:9573–9585,
136 2020.
- 137 [3] Guillermo Valle-Perez, Chico Q. Camargo, and Ard A. Louis. Deep learning generalizes because the
138 parameter-function map is biased towards simple functions. In *International Conference on Learning
139 Representations*, 2019.
- 140 [4] Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent on linear
141 convolutional networks. *Advances in neural information processing systems*, 31, 2018.
- 142 [5] Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. *arXiv preprint
143 arXiv:1810.02032*, 2018.
- 144 [6] Daniel Kunin, Atsushi Yamamura, Chao Ma, and Surya Ganguli. The asymmetric maximum margin bias
145 of quasi-homogeneous neural networks. *arXiv preprint arXiv:2210.03820*, 2022.
- 146 [7] Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit
147 bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878,
148 2018.
- 149 [8] Minyoung Huh, Hossein Mobahi, Richard Zhang, Brian Cheung, Pulkit Agrawal, and Phillip Isola. The
150 low-rank simplicity bias in deep networks. *Transactions on Machine Learning Research*, 2023.
- 151 [9] Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro.
152 Implicit regularization in matrix factorization. *Advances in Neural Information Processing Systems*, 30,
153 2017.
- 154 [10] Gauthier Gidel, Francis Bach, and Simon Lacoste-Julien. Implicit regularization of discrete gradient
155 dynamics in linear neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- 156 [11] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization.
157 *Advances in Neural Information Processing Systems*, 32, 2019.
- 158 [12] Zhiyuan Li, Yuping Luo, and Kaifeng Lyu. Towards resolving the implicit bias of gradient descent for
159 matrix factorization: Greedy low-rank learning, 2020.
- 160 [13] Samuel Horvath, Stefanos Laskaridis, Shashank Rajput, and Hongyi Wang. Maestro: Uncovering low-rank
161 structures via trainable decomposition, 2023.
- 162 [14] Jiawei Zhao, Yifei Zhang, Beidi Chen, Florian Schäfer, and Anima Anandkumar. Inrank: Incremental
163 low-rank learning, 2023.
- 164 [15] Hongyi Wang, Saurabh Agarwal, Pongsakorn U-chupala, Yoshiki Tanaka, Eric P. Xing, and Dimitris
165 Papaliopoulos. Cuttlefish: Low-rank model training without all the tuning, 2023.
- 166 [16] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and
167 Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*,
168 2021.
- 169 [17] Albert Gural, Phillip Nadeau, Mehul Tikekar, and Boris Murmann. Low-rank training of deep neural
170 networks for emerging memory technology, 2020.
- 171 [18] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of
172 learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- 173 [19] Andrew M Saxe, James L McClelland, and Surya Ganguli. A mathematical theory of semantic development
174 in deep neural networks. *Proceedings of the National Academy of Sciences*, 116(23):11537–11546, 2019.
- 175 [20] Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit acceleration
176 by overparameterization. In *International Conference on Machine Learning*, pages 244–253. PMLR, 2018.
- 177 [21] Hancheng Min, Salma Tarmoun, René Vidal, and Enrique Mallada. Convergence and implicit bias of
178 gradient flow on overparametrized linear networks. *arXiv preprint arXiv:2105.06351*, 2022.
- 179 [22] Daniel Gissin, Shai Shalev-Shwartz, and Amit Daniely. The implicit bias of depth: How incremental
180 learning drives generalization. *arXiv preprint arXiv:1909.12051*, 2019.
- 181 [23] Gal Vardi and Ohad Shamir. Implicit regularization in relu networks with the square loss. In *Conference
182 on Learning Theory*, pages 4224–4258. PMLR, 2021.
- 183 [24] Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient descent
184 for deep linear neural networks. *arXiv preprint arXiv:1810.02281*, 2018.
- 185 [25] Emmanuel Candes and Benjamin Recht. Exact matrix completion via convex optimization. *Communica-
186 tions of the ACM*, 55(6):111–119, 2012.

- 187 [26] Emmanuel J Candès and Terence Tao. The power of convex relaxation: Near-optimal matrix completion.
188 *IEEE Transactions on Information Theory*, 56(5):2053–2080, 2010.
- 189 [27] Mark A Davenport and Justin Romberg. An overview of low-rank matrix recovery from incomplete
190 observations. *IEEE Journal of Selected Topics in Signal Processing*, 10(4):608–622, 2016.
- 191 [28] Samuel Burer and Renato DC Monteiro. A nonlinear programming algorithm for solving semidefinite
192 programs via low-rank factorization. *Mathematical Programming*, 95(2):329–357, 2003.
- 193 [29] Prateek Jain, Praneeth Netrapalli, and Sujay Sanghavi. Low-rank matrix completion using alternating
194 minimization. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages
195 665–674, 2013.
- 196 [30] Qinqing Zheng and John Lafferty. Convergence analysis for rectangular matrix completion using burer-
197 monteiro factorization and gradient descent. *arXiv preprint arXiv:1605.07051*, 2016.
- 198 [31] Ruoyu Sun and Zhi-Quan Luo. Guaranteed matrix completion via non-convex factorization. *IEEE*
199 *Transactions on Information Theory*, 62(11):6535–6579, 2016.
- 200 [32] Rong Ge, Jason D Lee, and Tengyu Ma. Matrix completion has no spurious local minimum. *Advances in*
201 *neural information processing systems*, 29, 2016.
- 202 [33] Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Global optimality of local search for low rank
203 matrix recovery. *Advances in Neural Information Processing Systems*, 29, 2016.
- 204 [34] Rong Ge, Chi Jin, and Yi Zheng. No spurious local minima in nonconvex low rank problems: A unified
205 geometric analysis. In *International Conference on Machine Learning*, pages 1233–1242. PMLR, 2017.
- 206 [35] Qiuwei Li, Zhihui Zhu, and Gongguo Tang. The non-convex geometry of low-rank matrix optimization.
207 *Information and Inference: A Journal of the IMA*, 8(1):51–96, 2019.
- 208 [36] Yuejie Chi, Yue M Lu, and Yuxin Chen. Nonconvex optimization meets low-rank matrix factorization: An
209 overview. *IEEE Transactions on Signal Processing*, 67(20):5239–5269, 2019.
- 210 [37] Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in over-parameterized matrix
211 sensing and neural networks with quadratic activations. In *Conference On Learning Theory*, pages 2–47.
212 PMLR, 2018.
- 213 [38] Mahdi Soltanolkotabi, Dominik Stöger, and Changzhi Xie. Implicit balancing and regularization: Gener-
214 alization and convergence guarantees for overparameterized asymmetric matrix sensing. *arXiv preprint*
215 *arXiv:2303.14244*, 2023.

Appendix

Notation. Given any $L \in \mathbb{N}$, we use $[L]$ to denote the index set $\{1, \dots, L\}$. We use $I_n \in \mathbb{R}^n$ to denote the identity matrix of size n , and $\mathbf{1}_n$ to denote a vector of length n with all entries equal to 1. We denote by $\|\mathbf{A}\|_F^2$ the squared *Frobenius* norm of matrix \mathbf{A} , i.e., the sum of squares of all entries of \mathbf{A} . For convenience, whenever $j > i$ we adopt the abbreviations $\mathbf{W}_{j:i} = \mathbf{W}_j \cdots \mathbf{W}_i$ and $\mathbf{W}_{j:i}^\top = \mathbf{W}_i^\top \cdots \mathbf{W}_j^\top$, whereas both are identity if $j < i$.

224 A Proofs in Section 2

225 Substituting the analytic form of the gradient into (2), we have the update rule

$$\mathbf{W}_l(t+1) = (1 - \eta\lambda)\mathbf{W}_l(t) - \eta\mathbf{W}_{L:l+1}^\top(t)\mathbf{E}(t)\mathbf{W}_{l-1:1}^\top(t), \quad l \in [L] \quad (9)$$

226 for $t = 0, 1, 2, \dots$, where $\mathbf{E}(t) = f(\Theta(t)) - \Phi$.

227 We first establish the following Lemma 1 – the claim in Theorem 1 then follows in a relatively
228 straightforward manner. We note that all statements quantified by i in this section implicitly hold for
229 all $i \in [m]$ (as defined in Theorem 1) for the sake of notational brevity.

230 A.1 Proof of Theorem 1

231 **Lemma 1.** *Under the setting of Theorem 1, there exist orthonormal sets $\{\mathbf{u}_i^{(l)}\}_{i=1}^m \subset \mathbb{R}^d$ and*
232 *$\{\mathbf{v}_i^{(l)}\}_{i=1}^m \subset \mathbb{R}^d$ for $l \in [L]$ satisfying $\mathbf{v}_i^{(l+1)} = \mathbf{u}_i^{(l)}$ for all $l \in [L-1]$ such that the following hold*
233 *for all $t \geq 0$:*

$$\begin{aligned} \mathcal{A}(t) : \mathbf{W}_l(t)\mathbf{v}_i^{(l)} &= \rho_l(t)\mathbf{u}_i^{(l)} & \forall l \in [L], \\ \mathcal{B}(t) : \mathbf{W}_l^\top(t)\mathbf{u}_i^{(l)} &= \rho_l(t)\mathbf{v}_i^{(l)} & \forall l \in [L], \\ \mathcal{C}(t) : \Phi^\top \mathbf{W}_{L:l+1}(t)\mathbf{u}_i^{(l)} &= \mathbf{0} & \forall l \in [L], \\ \mathcal{D}(t) : \Phi \mathbf{W}_{l-1:1}^\top(t)\mathbf{v}_i^{(l)} &= \mathbf{0} & \forall l \in [L], \end{aligned}$$

234 where $\rho_l(t) = \rho_l(t-1) \cdot (1 - \eta\lambda - \eta \cdot \prod_{k \neq l} \rho_k(t-1)^2)$ for all $t \geq 1$ with $\rho_l(0) = \sigma_l > 0$.

235 *Proof.* Define $\Psi := \mathbf{W}_{L:2}^\top(0)\Phi$. Since the rank of Φ is at most r , we have that the rank of $\Psi \in \mathbb{R}^{d \times d}$
236 is at most r , which implies that $\dim \mathcal{N}(\Psi) = \dim \mathcal{N}(\Psi^\top) \geq d - r$. We define the subspace

$$\mathcal{S} := \mathcal{N}(\Psi) \cap \mathcal{N}(\Psi^\top \mathbf{W}_1(0)) \subset \mathbb{R}^d.$$

237 Since $\mathbf{W}_1(0) \in \mathbb{R}^{d \times d}$ is nonsingular, we have

$$\dim(\mathcal{S}) \geq 2(d - r) - d = m.$$

238 Let $\{\mathbf{v}_i^{(1)}\}_{i=1}^m$ denote an orthonormal set contained in \mathcal{S} and set $\mathbf{u}_i^{(1)} := \mathbf{W}_1(0)\mathbf{v}_i^{(1)}/\sigma_1$, where
239 $\sigma_1 > 0$ is the scale of $\mathbf{W}_1(0)$ – since $\mathbf{W}_1(0)/\sigma_1$ is orthogonal, $\{\mathbf{u}_i^{(1)}\}_{i=1}^m$ is also an orthonormal
240 set. Then we trivially have $\mathbf{W}_1(0)\mathbf{v}_i^{(1)} = \sigma_1\mathbf{u}_i^{(1)}$, which implies $\mathbf{W}_1^\top(0)\mathbf{u}_i^{(1)} = \sigma_1\mathbf{v}_i^{(1)}$. It follows
241 from $\mathbf{v}_i^{(1)} \in \mathcal{S}$ that $\Psi\mathbf{v}_i^{(1)} = \mathbf{0}$ and $\Psi^\top \mathbf{W}_1(0)\mathbf{v}_i^{(1)} = \mathbf{0}$, which is equivalent to $\mathbf{W}_{L:2}^\top(0)\Phi\mathbf{v}_i^{(1)} = \mathbf{0}$
242 and $\Phi^\top \mathbf{W}_{L:2}(0)\mathbf{W}_1(0)\mathbf{v}_i^{(1)} = \sigma_1\Phi^\top \mathbf{W}_{L:2}(0)\mathbf{u}_i^{(1)} = \mathbf{0}$ respectively. Since $\mathbf{W}_{L:2}^\top(0)$ is full column
243 rank, we further have that $\Phi\mathbf{v}_i^{(1)} = \mathbf{0}$.

244 Now let $\mathcal{E}(l)$ be the event that we have orthonormal sets $\{\mathbf{u}_i^{(l)}\}_{i=1}^m$ and $\{\mathbf{v}_i^{(l)}\}_{i=1}^m$ satisfying
245 $\mathbf{W}_l(0)\mathbf{v}_i^{(l)} = \sigma_l\mathbf{u}_i^{(l)}$, $\mathbf{W}_l^\top(0)\mathbf{u}_i^{(l)} = \sigma_l\mathbf{v}_i^{(l)}$, $\Phi^\top \mathbf{W}_{L:l+1}(0)\mathbf{u}_i^{(l)} = \mathbf{0}$, and $\Phi \mathbf{W}_{l-1:1}^\top(0)\mathbf{v}_i^{(l)} = \mathbf{0}$.
246 From the above arguments, we have that $\mathcal{E}(1)$ holds – now suppose $\mathcal{E}(k)$ holds for some $1 \leq k < L$.

247 Set $\mathbf{v}_i^{(k+1)} := \mathbf{u}_i^{(k)}$ and $\mathbf{u}_i^{(k+1)} := \mathbf{W}_{k+1}(0)\mathbf{v}_i^{(k+1)}/\sigma_{k+1}$. This implies that $\mathbf{W}_{k+1}(0)\mathbf{v}_i^{(k+1)} =$
 248 $\sigma_{k+1}\mathbf{u}_i^{(k+1)}$ and $\mathbf{W}_{k+1}^\top(0)\mathbf{u}_i^{(k+1)} = \sigma_{k+1}\mathbf{v}_i^{(k+1)}$. Moreover, we have

$$\begin{aligned}\Phi^\top \mathbf{W}_{L:(k+1)+1}(0)\mathbf{u}_i^{(k+1)} &= \Phi^\top \mathbf{W}_{L:k+1}(0)\mathbf{W}_{k+1}^\top(0)\mathbf{u}_i^{(k+1)}/\sigma_{k+1}^2 \\ &= \Phi^\top \mathbf{W}_{L:k+1}(0)\mathbf{v}_i^{(k+1)}/\sigma_{k+1} \\ &= \Phi^\top \mathbf{W}_{L:k+1}(0)\mathbf{u}_i^{(k)}/\sigma_{k+1} = \mathbf{0},\end{aligned}$$

249 where the first two equalities follow from orthogonality and $\mathbf{u}_i^{(k+1)} = \mathbf{W}_{k+1}(0)\mathbf{v}_i^{(k+1)}/\sigma_{k+1}$, and
 250 the last equality is due to $\mathbf{v}_i^{(k+1)} = \mathbf{u}_i^{(k)}$. Similarly, we have

$$\begin{aligned}\Phi \mathbf{W}_{(k+1)-1:1}^\top(0)\mathbf{v}_i^{(k+1)} &= \Phi \mathbf{W}_{k-1:1}^\top(0)\mathbf{W}_k^\top(0)\mathbf{v}_i^{(k+1)} \\ &= \Phi \mathbf{W}_{k-1:1}^\top(0)\mathbf{W}_k^\top(0)\mathbf{u}_i^{(k)} \\ &= \sigma_k \Phi \mathbf{W}_{k-1:1}^\top(0)\mathbf{v}_i^{(k)} = \mathbf{0},\end{aligned}$$

251 where the second equality follows from $\mathbf{v}_i^{(k+1)} = \mathbf{u}_i^{(k)}$ and the third equality is due to $\mathbf{W}_k^\top(0)\mathbf{u}_i^{(k)} =$
 252 $\sigma_k\mathbf{v}_i^{(k)}$. Therefore $\mathcal{E}(k+1)$ holds, so we have $\mathcal{E}(l)$ for all $l \in [L]$. As a result, we have shown the
 253 base cases $\mathcal{A}(0)$, $\mathcal{B}(0)$, $\mathcal{C}(0)$, and $\mathcal{D}(0)$.

254 Now we proceed by induction on $t \geq 0$. Suppose that $\mathcal{A}(t)$, $\mathcal{B}(t)$, $\mathcal{C}(t)$, and $\mathcal{D}(t)$ hold for some
 255 $t \geq 0$. First, we show $\mathcal{A}(t+1)$ and $\mathcal{B}(t+1)$. We have

$$\begin{aligned}\mathbf{W}_i(t+1)\mathbf{v}_i^{(l)} &= [(1-\eta\lambda)\mathbf{W}_i(t) - \eta\mathbf{W}_{L:l+1}^\top(t)\mathbf{E}(t)\mathbf{W}_{l-1:1}^\top(t)]\mathbf{v}_i^{(l)} \\ &= [(1-\eta\lambda)\mathbf{W}_i(t) - \eta\mathbf{W}_{L:l+1}^\top(t)(\mathbf{W}_{L:1}(t) - \Phi)\mathbf{W}_{l-1:1}^\top(t)]\mathbf{v}_i^{(l)} \\ &= (1-\eta\lambda)\mathbf{W}_i(t)\mathbf{v}_i^{(l)} - \eta\mathbf{W}_{L:l+1}^\top(t)\mathbf{W}_{L:1}(t)\mathbf{W}_{l-1:1}^\top(t)\mathbf{v}_i^{(l)} \\ &= (1-\eta\lambda)\mathbf{W}_i(t)\mathbf{v}_i^{(l)} - \eta \cdot \left(\prod_{k \neq l} \rho_k^2(t)\right)\mathbf{W}_i(t)\mathbf{v}_i^{(l)} \\ &= \rho_l(t) \cdot (1-\eta\lambda - \eta \cdot \prod_{k \neq l} \rho_k^2(t))\mathbf{u}_i^{(l)} = \rho_l(t+1)\mathbf{u}_i^{(l)}\end{aligned}$$

256 for all $l \in [L]$, where the first equality follows from (9), the second equality follows from definition
 257 of $\mathbf{E}(t)$, the third equality follows from $\mathcal{D}(t)$, and the fourth equality follows from $\mathcal{A}(t)$ and $\mathcal{B}(t)$
 258 applied repeatedly along with $\mathbf{v}_i^{(l+1)} = \mathbf{u}_i^{(l)}$ for all $l \in [L-1]$, proving $\mathcal{A}(t+1)$. Similarly, we have

$$\begin{aligned}\mathbf{W}_i^\top(t+1)\mathbf{u}_i^{(l)} &= [(1-\eta\lambda)\mathbf{W}_i^\top(t) - \eta\mathbf{W}_{l-1:1}(t)\mathbf{E}^\top(t)\mathbf{W}_{L:l+1}(t)]\mathbf{u}_i^{(l)} \\ &= [(1-\eta\lambda)\mathbf{W}_i^\top(t) - \eta\mathbf{W}_{l-1:1}(t)(\mathbf{W}_{L:1}^\top(t) - \Phi^\top)\mathbf{W}_{L:l+1}(t)]\mathbf{u}_i^{(l)} \\ &= (1-\eta\lambda)\mathbf{W}_i^\top(t)\mathbf{u}_i^{(l)} - \eta\mathbf{W}_{l-1:1}(t)\mathbf{W}_{L:1}^\top(t)\mathbf{W}_{L:l+1}(t)\mathbf{u}_i^{(l)} \\ &= (1-\eta\lambda)\mathbf{W}_i^\top(t)\mathbf{u}_i^{(l)} - \eta \cdot \left(\prod_{k \neq l} \rho_k^2(t)\right)\mathbf{W}_i^\top(t)\mathbf{u}_i^{(l)} \\ &= \rho_l(t) \cdot (1-\eta\lambda - \eta \cdot \prod_{k \neq l} \rho_k^2(t))\mathbf{v}_i^{(l)} = \rho_l(t+1)\mathbf{v}_i^{(l)}\end{aligned}$$

259 for all $l \in [L]$, where the third equality follows from $\mathcal{C}(t)$, and the fourth equality follows from $\mathcal{A}(t)$
 260 and $\mathcal{B}(t)$ applied repeatedly along with $\mathbf{v}_i^{(l+1)} = \mathbf{u}_i^{(l)}$ for all $l \in [L-1]$, proving $\mathcal{B}(t+1)$. Now, we
 261 show $\mathcal{C}(t+1)$. For any $k \in [L-1]$, it follows from $\mathbf{v}_i^{(k+1)} = \mathbf{u}_i^{(k)}$ and $\mathcal{A}(t+1)$ that

$$\mathbf{W}_{k+1}(t+1)\mathbf{u}_i^{(k)} = \mathbf{W}_{k+1}(t+1)\mathbf{v}_i^{(k+1)} = \rho_{k+1}(t+1)\mathbf{u}_i^{(k+1)}.$$

262 Repeatedly applying the above equality for $k = l, l+1, \dots, L-1$, we obtain

$$\Phi^\top \mathbf{W}_{L:l+1}(t)\mathbf{u}_i^{(l)} = \left[\prod_{k=l}^{L-1} \rho_{k+1}(t) \right] \cdot \Phi^\top \mathbf{u}_i^{(L)} = \mathbf{0}$$

263 which follows from $\mathcal{C}(t)$, proving $\mathcal{C}(t+1)$. Finally, we show $\mathcal{D}(t+1)$. For any $k \in \{2, \dots, L\}$, it
 264 follows from $\mathbf{v}_i^{(k)} = \mathbf{u}_i^{(k-1)}$ and $\mathcal{B}(t+1)$ that

$$\mathbf{W}_{k-1}^\top(t+1)\mathbf{v}_i^{(k)} = \mathbf{W}_{k-1}^\top(t+1)\mathbf{u}_i^{(k-1)} = \rho_{k-1}(t+1)\mathbf{v}_i^{(k-1)}.$$

265 Repeatedly applying the above equality for $k = l, l-1, \dots, 2$, we obtain

$$\Phi \mathbf{W}_{l-1:1}^\top(t)\mathbf{v}_i^{(l)} = \left[\prod_{k=2}^l \rho_{k-1}(t) \right] \cdot \Phi \mathbf{v}_i^{(1)} = \mathbf{0}$$

266 which follows from $\mathcal{D}(t)$. Thus we have proven $\mathcal{D}(t+1)$, concluding the proof. \square

267 *Proof of Theorem 1.* By $\mathcal{A}(t)$ and $\mathcal{B}(t)$ of Lemma 1, there exists orthonormal matrices $\{\mathbf{U}_{l,2}\}_{l=1}^L \subset$
 268 $\mathbf{R}^{d \times m}$ and $\{\mathbf{V}_{l,2}\}_{l=1}^L \subset \mathbf{R}^{d \times m}$ for $l \in [L]$ satisfying $\mathbf{U}_{l+1,2} = \mathbf{V}_{l,2}$ for all $l \in [L-1]$ as well as

$$\mathbf{W}_l(t)\mathbf{V}_{l,2} = \rho_l(t)\mathbf{U}_{l,2} \quad \text{and} \quad \mathbf{W}_l(t)^\top \mathbf{U}_{l,2} = \rho_l(t)\mathbf{V}_{l,2} \quad (10)$$

269 for all $l \in [L]$ and $t \geq 0$, where $\rho_l(t)$ satisfies (4) for $t \geq 1$ with $\rho_l(0) = \sigma_l$. First, complete $\mathbf{V}_{l,2}$ to
 270 an orthonormal basis for \mathbf{R}^d as $\mathbf{V}_l = [\mathbf{V}_{l,1} \ \mathbf{V}_{l,2}]$. Then for each $l \in [L-1]$, set $\mathbf{U}_l = [\mathbf{U}_{l,1} \ \mathbf{U}_{l,2}]$
 271 where $\mathbf{U}_{l,1} = \mathbf{W}_l(0)\mathbf{V}_{l,1}/\sigma_l$ and $\mathbf{V}_{l+1} = [\mathbf{V}_{l+1,1} \ \mathbf{V}_{l+1,2}]$ where $\mathbf{V}_{l+1,1} = \mathbf{U}_{l,1}$, and finally set
 272 $\mathbf{U}_L = [\mathbf{U}_{L,1} \ \mathbf{U}_{L,2}]$ where $\mathbf{U}_{L,1} = \mathbf{W}_L(0)\mathbf{V}_{L,1}/\sigma_L$. We note that $\mathbf{V}_{l+1} = \mathbf{U}_l$ for each $l \in [L-1]$.
 273 Then we have

$$\mathbf{U}_{l,1}^\top \mathbf{W}_l(t)\mathbf{V}_{l,2} = \rho_l(t)\mathbf{U}_{l,1}^\top \mathbf{U}_{l,2} = \mathbf{0} \quad (11)$$

274 for all $l \in [L]$, where the first equality follows from (10). Similarly, we also have

$$\mathbf{U}_{l,2}^\top \mathbf{W}_l(t)\mathbf{V}_{l,1} = \rho_l(t)\mathbf{V}_{l,2}^\top \mathbf{V}_{l,1} = \mathbf{0} \quad (12)$$

275 for all $l \in [L]$, where the first equality also follows from (10). Therefore, combining (10), (11), and
 276 (12) yields

$$\mathbf{U}_l^\top \mathbf{W}_l(t)\mathbf{V}_l = [\mathbf{U}_{l,1} \ \mathbf{U}_{l,2}]^\top \mathbf{W}_l(t) [\mathbf{V}_{l+1,1} \ \mathbf{V}_{l+1,2}] = \begin{bmatrix} \widetilde{\mathbf{W}}_l(t) & \mathbf{0} \\ \mathbf{0} & \rho_l(t)\mathbf{I}_m \end{bmatrix}$$

277 for all $l \in [L]$, where $\widetilde{\mathbf{W}}_l(0) = \sigma_l \mathbf{I}_{2r}$ by construction of $\mathbf{U}_{l,1}$. This directly implies (3), completing
 278 the proof. \square

279 A.2 Proof of Proposition 1

280 *Proof.* First, it follows from Theorem 1 that for any $1 \leq i \leq j \leq L$ we have

$$\mathbf{W}_{j:i}(t) = \mathbf{U}_{j,1} \widetilde{\mathbf{W}}_{j:i}(t) \mathbf{V}_{i,1}^\top + \left(\prod_{k=i}^j \rho_k(t) \right) \cdot \mathbf{U}_{j,2} \mathbf{V}_{i,2}^\top \quad (13)$$

281 for all $t \geq 0$, where $\mathbf{U}_{l,1}, \mathbf{V}_{l,1} \in \mathbf{R}^{d \times 2\widehat{r}}$ and $\mathbf{U}_{l,2}, \mathbf{V}_{l,2} \in \mathbf{R}^{d \times \widehat{m}}$ are the first $2\widehat{r}$ and last \widehat{m} columns
 282 of $\mathbf{U}_l, \mathbf{V}_l \in \mathbf{R}^{d \times d}$ respectively.

283 The key claim to be shown here is that $\widehat{\mathbf{W}}_l(t) = \widetilde{\mathbf{W}}_l(t)$ for all $l \in [L]$ and $t \geq 0$. Afterwards, it
 284 follows straightforwardly from (13) that

$$\begin{aligned} & \left\| f(\Theta(t)) - \widehat{f}(\widehat{\Theta}(t), \mathbf{U}_{L,1}, \mathbf{V}_{1,1}) \right\|_F^2 \\ &= \left\| \mathbf{U}_{L,1} \widetilde{\mathbf{W}}_{L:1}(t) \mathbf{V}_{1,1}^\top + \left(\prod_{l=1}^L \rho_l(t) \right) \cdot \mathbf{U}_{L,2} \mathbf{V}_{1,2}^\top - \mathbf{U}_{L,1} \widehat{\mathbf{W}}_{L:1}(t) \mathbf{V}_{L,1}^\top \right\|_F^2 \\ &= \left\| \mathbf{U}_{L,1} (\widetilde{\mathbf{W}}_{L:1}(t) - \widehat{\mathbf{W}}_{L:1}(t)) \mathbf{V}_{1,1}^\top + \left(\prod_{l=1}^L \rho_l(t) \right) \cdot \mathbf{U}_{L,2} \mathbf{V}_{1,2}^\top \right\|_F^2 \\ &= \left\| \left(\prod_{l=1}^L \rho_l(t) \right) \cdot \mathbf{U}_{L,2} \mathbf{V}_{1,2}^\top \right\|_F^2 \leq \widehat{m} \cdot \prod_{l=1}^L \sigma_l^2. \end{aligned}$$

285 We proceed by induction. For $t = 0$, we have that

$$\widehat{\mathbf{W}}_l(0) = \mathbf{U}_{l,1}^\top \mathbf{W}_l(0) \mathbf{V}_{l,1} = \widetilde{\mathbf{W}}_l(0)$$

286 for all $l \in [L]$ by (13) and choice of initialization.

287 Now suppose $\widehat{\mathbf{W}}_l(t) = \widetilde{\mathbf{W}}_l(t)$ for all $l \in [L]$. Comparing

$$\widehat{\mathbf{W}}_l(t+1) = (1 - \eta\lambda)\widehat{\mathbf{W}}_l(t) - \eta\nabla_{\widehat{\mathbf{w}}_l} \widehat{\ell}(\widehat{\Theta}(t))$$

288 with

$$\begin{aligned} \widetilde{\mathbf{W}}_l(t+1) &= \mathbf{U}_{l,1}^\top \mathbf{W}_l(t+1) \mathbf{V}_{l,1} \\ &= \mathbf{U}_{l,1}^\top [(1 - \eta\lambda)\mathbf{W}_l(t) - \eta\nabla_{\mathbf{w}_l} \ell(\Theta(t))] \mathbf{V}_{l,1} \\ &= (1 - \eta\lambda)\widetilde{\mathbf{W}}_l(t) - \eta\mathbf{U}_{l,1}^\top \nabla_{\mathbf{w}_l} \ell(\Theta(t)) \mathbf{V}_{l,1} \end{aligned}$$

289 it suffices to show that

$$\nabla_{\widehat{\mathbf{w}}_l} \widehat{\ell}(\widehat{\Theta}(t)) = \mathbf{U}_{l,1}^\top \nabla_{\mathbf{w}_l} \ell(\Theta(t)) \mathbf{V}_{l,1}, \quad \forall l \in [L] \quad (14)$$

290 to yield $\widehat{\mathbf{W}}_l(t+1) = \widetilde{\mathbf{W}}_l(t+1)$ for all $l \in [L]$. Computing the right hand side of (14), we have

$$\begin{aligned} \mathbf{U}_{l,1}^\top \nabla_{\mathbf{w}_l} \ell(\Theta(t)) \mathbf{V}_{l,1} &= \mathbf{U}_{l,1}^\top \mathbf{W}_{L:l+1}^\top(t) (\mathbf{W}_{L:1}(t) - \Phi) \mathbf{W}_{l-1:1}^\top(t) \mathbf{V}_{l,1} \\ &= (\mathbf{W}_{L:l+1}(t) \mathbf{U}_{l,1})^\top (\mathbf{W}_{L:1}(t) - \Phi) (\mathbf{V}_{l,1}^\top \mathbf{W}_{l-1:1}(t))^\top \end{aligned}$$

291 where

$$\mathbf{W}_{L:l+1}(t) \mathbf{U}_{l,1} = \left(\mathbf{U}_{L,1} \widetilde{\mathbf{W}}_{L:l+1}(t) \mathbf{V}_{l+1,1}^\top + \left(\prod_{k=l+1}^L \rho_k(t) \right) \cdot \mathbf{U}_{L,2} \mathbf{V}_{l+1,2}^\top \right) \mathbf{U}_{l,1} = \mathbf{U}_{L,1} \widetilde{\mathbf{W}}_{L:l+1}(t)$$

292 by (13) and the fact that $\mathbf{U}_l = \mathbf{V}_{l+1}$, and similarly

$$\mathbf{V}_{l,1}^\top \mathbf{W}_{l-1:1}(t) = \mathbf{V}_{l,1}^\top \left(\mathbf{U}_{l-1,1} \widetilde{\mathbf{W}}_{l-1:1}(t) \mathbf{V}_{1,1}^\top + \left(\prod_{k=1}^{l-1} \rho_k(t) \right) \cdot \mathbf{U}_{l-1,2} \mathbf{V}_{1,2}^\top \right) = \widetilde{\mathbf{W}}_{l-1:1}(t) \mathbf{V}_{1,1}^\top.$$

293 We also have that

$$\begin{aligned} \mathbf{U}_{L,1}^\top (\mathbf{W}_{L:1}(t) - \Phi) \mathbf{V}_{1,1} &= \mathbf{U}_{L,1}^\top \left(\mathbf{U}_{L,1} \widetilde{\mathbf{W}}_{L:1}(t) \mathbf{V}_{1,1}^\top + \left(\prod_{k=1}^L \rho_k(t) \right) \cdot \mathbf{U}_{L,2} \mathbf{V}_{1,2}^\top - \Phi \right) \mathbf{V}_{1,1} \\ &= \widetilde{\mathbf{W}}_{L:1}(t) - \mathbf{U}_{L,1}^\top \Phi \mathbf{V}_{1,1} \end{aligned}$$

294 so putting together the previous four equalities yields

$$\begin{aligned} \mathbf{U}_{l,1}^\top \nabla_{\mathbf{w}_l} \ell(\Theta(t)) \mathbf{V}_{l,1} &= (\mathbf{W}_{L:l+1}(t) \mathbf{U}_{l,1})^\top (\mathbf{W}_{L:1}(t) - \Phi) (\mathbf{V}_{l,1}^\top \mathbf{W}_{l-1:1}(t))^\top \\ &= \widetilde{\mathbf{W}}_{L:l+1}^\top(t) \mathbf{U}_{L,1}^\top (\mathbf{W}_{L:1}(t) - \Phi) \mathbf{V}_{1,1} \widetilde{\mathbf{W}}_{l-1:1}^\top(t) \\ &= \widetilde{\mathbf{W}}_{L:l+1}^\top(t) (\widetilde{\mathbf{W}}_{L:1}(t) - \mathbf{U}_{L,1}^\top \Phi \mathbf{V}_{1,1}) \widetilde{\mathbf{W}}_{l-1:1}^\top(t). \end{aligned}$$

295 On the other hand, the left hand side of (14) gives

$$\begin{aligned} \nabla_{\widehat{\mathbf{w}}_l} \widehat{\ell}(\widehat{\Theta}(t)) &= \widehat{\mathbf{W}}_{L:l+1}(t)^\top \mathbf{U}_{L,1}^\top (\mathbf{U}_{L,1} \widehat{\mathbf{W}}_{L:1}(t) \mathbf{V}_{1,1}^\top - \Phi) \mathbf{V}_{1,1} \widehat{\mathbf{W}}_{l-1:1}(t)^\top \\ &= \widehat{\mathbf{W}}_{L:l+1}(t)^\top (\widehat{\mathbf{W}}_{L:1}(t) - \mathbf{U}_{L,1}^\top \Phi \mathbf{V}_{1,1}) \widehat{\mathbf{W}}_{l-1:1}(t)^\top \end{aligned}$$

296 so (14) holds by the fact that $\widehat{\mathbf{W}}_l(t) = \widetilde{\mathbf{W}}_l(t)$ for all $l \in [L]$, completing the proof. \square

297 **B Benefits of Over-Parameterization in Deep Matrix Completion**

298 In the setup described in Section 3, we claim that depth and width are beneficial for accelerating
 299 GD convergence to well-generalizing solutions, and therefore constructing more computationally
 300 efficient factorizations that share the same trajectory is a fruitful endeavour. Below, we give a more
 301 detailed explanation of these ideas:

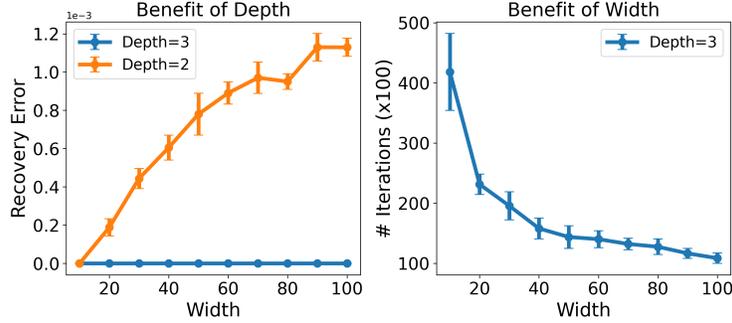


Figure 4: **Benefits of depth & width in overparameterized matrix completion** with $d = 100$, $r = 5$, $\sigma_l = 10^{-3}$ and 30% of entries observed. *Left*: Recovery error. *Right*: Number of GD iterations to converge to 10^{-10} error.

- 302 • **Benefits of depth.** When $L = 2$, (7) reduces to Burer-Monteiro factorization [28], whose global
 303 optimality and convergence under GD have been widely studied under various settings [9, 29–38].
 304 However, it has been demonstrated [11] that in the over-parameterized regime $\hat{r} > r$, deeper
 305 factorizations (starting from small random initialization) continue to generalize well beyond the
 306 exact parameterization $\hat{r} = r$ unlike their shallow counterparts, see Figure 4 (left).
- 307 • **Benefits of width.** On the other hand, increasing the width \hat{r} of the deep factorization beyond r
 308 results in accelerated convergence of GD in terms of iterations, see Figure 4 (right).

309 C Compressed vs. Narrow Factorizations

310 We compare the training efficiency of a $2\hat{r}$ -compressed factorization (with trajectory equivalent to
 311 a wide factorization of width $d \gg \hat{r}$) versus a narrow factorization with width $2\hat{r}$ under different
 312 over-parameterized estimates \hat{r} . As depicted in Figure 5 (left), the compressed factorization requires
 313 fewer iterations to reach convergence, and the number of iterations necessary is almost unaffected by
 314 \hat{r} . Consequently, **training compressed factorizations is considerably more time-efficient than**
 315 **training narrow ones of the same size**, provided that \hat{r} is not significantly larger than r .

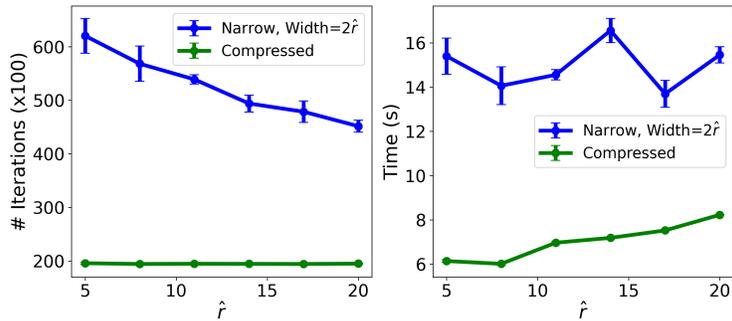


Figure 5: **Efficiency of compressed vs. narrow factorizations** for different overestimated \hat{r} with $L = 3$, $d = 1000$, $r = 5$, $\sigma_l = 10^{-3}$ and 20% of entries observed. *Left*: Number of iterations to converge to 10^{-10} error. *Right*: Time to converge.

316 The distinction between the compressed and narrow factorizations underscores the benefits of width,
 317 as previously demonstrated and discussed in Figure 4 (right), where increasing the width results in
 318 faster convergence. However, increasing the width alone also increases the number of parameters. By
 319 employing our compression methodology, we can achieve the best of both worlds.