

Implicit Regularization of Discrete Gradient Dynamics in Deep Linear Neural Networks

Gauthier Gidel¹, Francis Bach² and Simon Lacoste-Julien¹

¹Mila, Université de Montréal; ²SIERRA Project-Team, INRIA and ENS Paris.

Overview

Takeaways

- The choice of the optimization algorithm introduces **biases** that will lead to convergence to specific minimizers of the objective.
The path matters more than the destination.
- Different optimization algorithms or parametrizations of the model changes the **optimization path**.
- Study the **discrete gradient dynamics** of the training of a **two-layer linear network**. Can be related to matrix factorization [1].
- **Sequentially learns** the solutions of a reduced-rank regression with a gradually increasing rank.

Setting

A Simple Deep Linear Model

Deep linear model: $\hat{y}^d(x) := \mathbf{W}_L^\top \cdots \mathbf{W}_1^\top \mathbf{x}$,

Trained with MSE [2]:

$$\min_{\substack{\mathbf{W}_l \in \mathbb{R}^{r_{l-1} \times r_l} \\ 1 \leq l \leq L}} \frac{1}{2n} \|\mathbf{Y} - \mathbf{X} \mathbf{W}_1 \cdots \mathbf{W}_L\|_2^2.$$

- Thin matrices: **Low rank** constraint, reduced-rank regression,

$$\mathbf{W}^{k,*} \in \arg \min_{\substack{\mathbf{W} \in \mathbb{R}^{p \times d} \\ \text{rank}(\mathbf{W}) \leq r}} \frac{1}{2n} \sum_{i=1}^n \|\mathbf{Y} - \mathbf{X} \mathbf{W}\|_2^2.$$

- Large matrices: **overparametrized model**. Same expressivity as a linear model but **different** dynamics.

Gradient Dynamics

Discrete gradient dynamics,

$$\mathbf{W}_l^{(t+1)} = \mathbf{W}_l^{(t)} - \eta \nabla_{\mathbf{W}_l} f(\mathbf{W}_{[L]}^{(t)})$$

Continuous version,

$$\dot{\mathbf{W}}_l(t) = -\nabla_{\mathbf{W}_l} f(\mathbf{W}_{[L]}(t))$$

where $\mathbf{W}_{[L]}^{(t)} := (\mathbf{W}_1^{(t)}, \dots, \mathbf{W}_L^{(t)})$.

Assumption

The matrices $\mathbf{X}^\top \mathbf{X}$ and $\mathbf{X}^\top \mathbf{Y}$ are close to have **common decomposition**:

$$\begin{aligned} \mathbf{X}^\top \mathbf{X} &= \mathbf{U}(\mathbf{D}_x + \mathbf{B})\mathbf{U}^\top \\ \mathbf{X}^\top \mathbf{Y} &= \mathbf{U} \mathbf{D}_{xy} \mathbf{V}^\top, \end{aligned}$$

where \mathbf{B} is such that $\|\mathbf{B}\|_2 \leq \epsilon$ and $\mathbf{D}_x, \mathbf{D}_{xy}$ are matrices only with diagonal coefficients.

Sequential Learning of Components

Continuous Case

- Can find a **closed form** solution.
- Let $\mathbf{W}_1(t)$ and $\mathbf{W}_2(t)$ be these solutions.

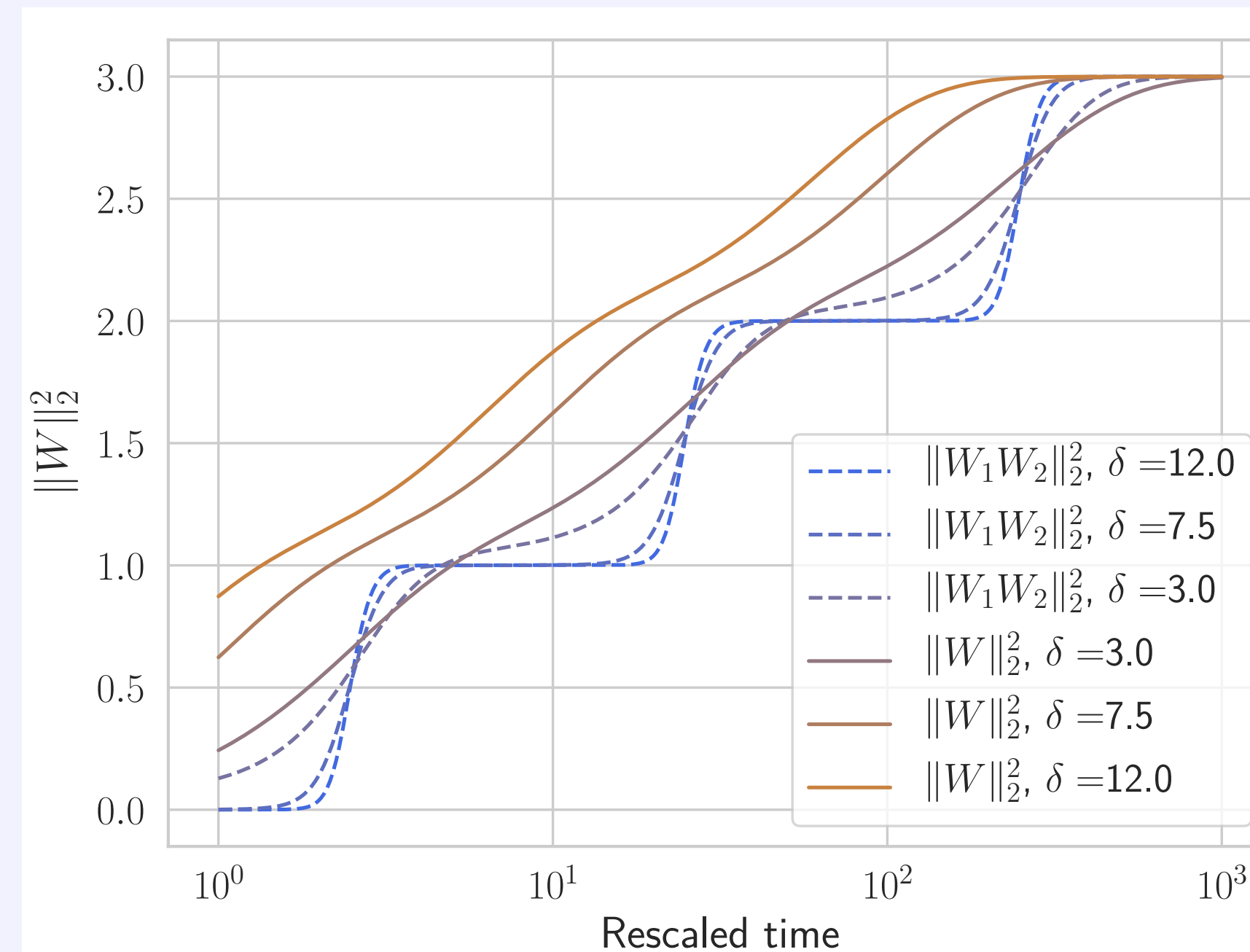
$$\frac{1}{\sigma_k} < t < \frac{1}{\sigma_{k+1}} \Rightarrow \mathbf{W}_1(\delta t) \mathbf{W}_2(\delta t) \xrightarrow{\delta \rightarrow \infty} \mathbf{W}^{k,*}$$

where $\mathbf{W}^{k,*}$ is the minimum ℓ_2 norm solution of the reduced-rank- k regression problem.

- Two-layer linear model sequentially find **min-norm low-rank** solutions.
- Notion of *phase transition time*:

$$T_i := \frac{1}{\sigma_i}$$

- Not the case for a one-layer linear model.
- Linear auto-encoder $\mathbf{X} = \mathbf{Y}$, trace norm is witness of the increasing rank of the solution.



Closed form solutions for a vanishing initialization.

Discrete Case

Why it is interesting:

- Want to understand implicit regularization in gradient based ML.
- In practice *discrete updates*.

Why it is more challenging:

- No *closed form* solution.
- *Infinite horizon* \rightarrow cannot consider discrete as an approximation of the continuous.
- New analysis required.

Results:

- Notion of *phase transition time*:

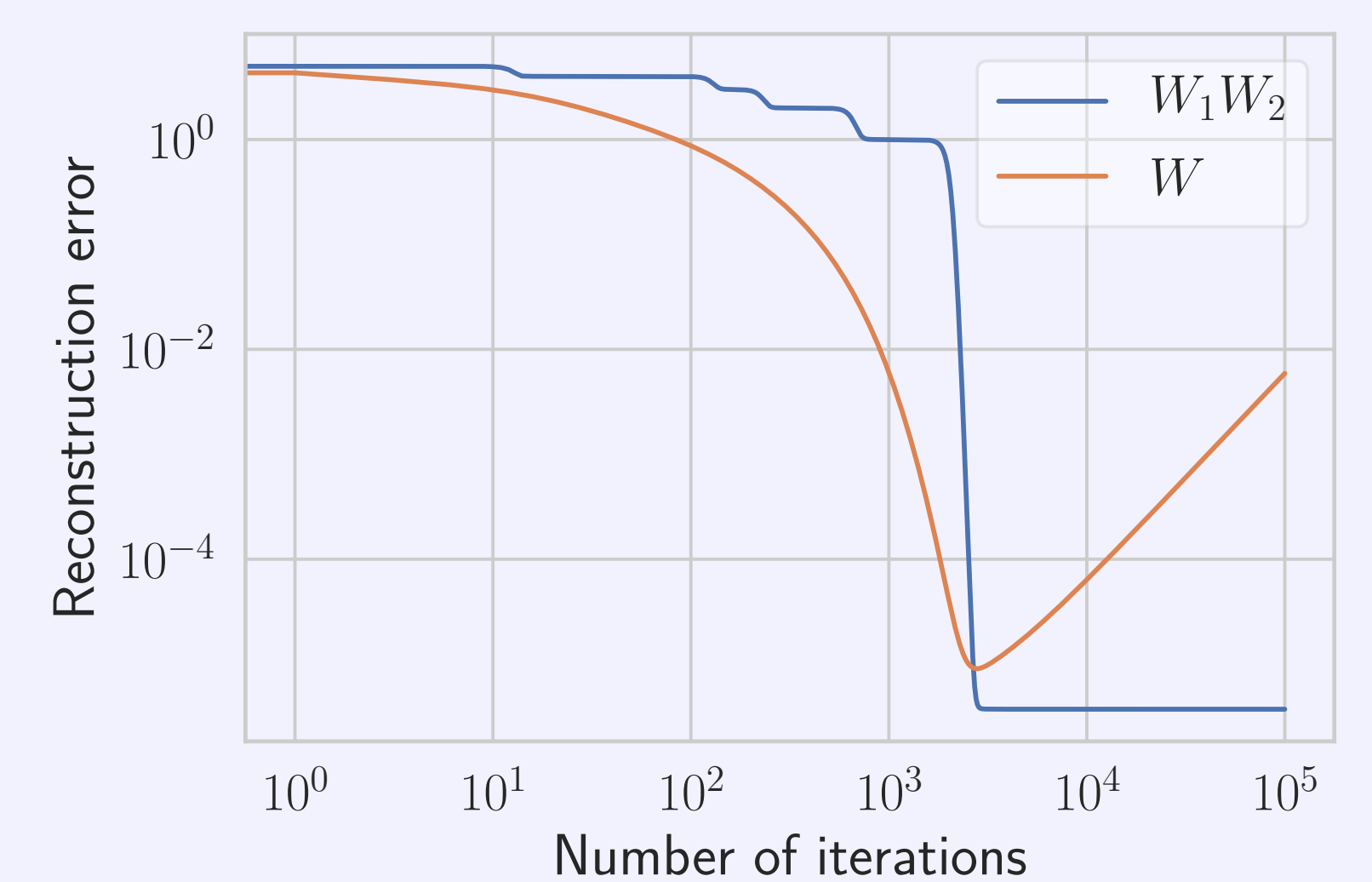
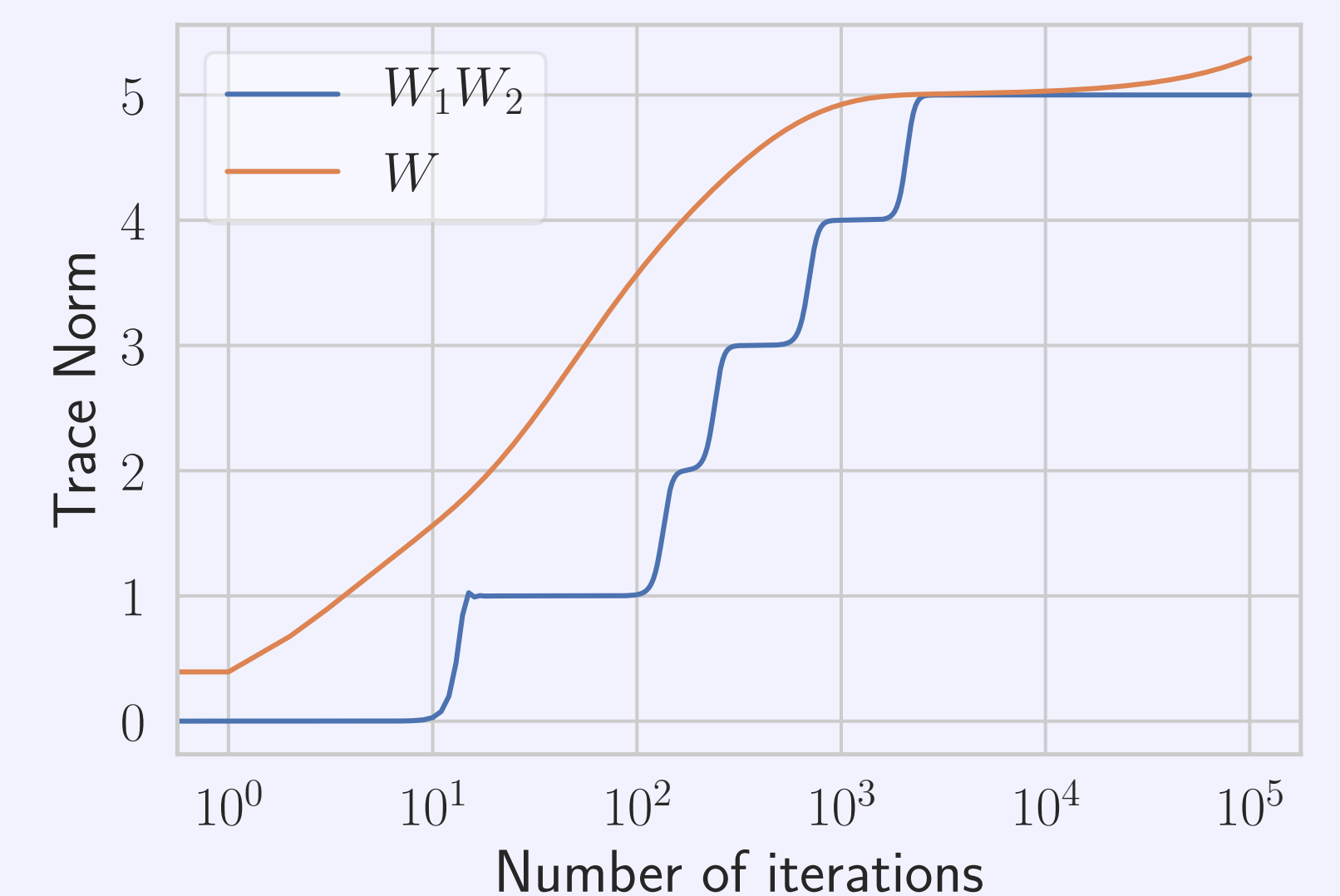
$$T_i := \frac{1}{\eta \sigma_i}$$

- Similar as the continuous case. (new proof)
- Additional constraints on the step-size.
 - Smaller than a notion of eigen-gap.
 - Smaller than the Lipschitz constant of the gradient.

Experiments and details

Linear auto-encoder

For linear auto encoder $\mathbf{X} = \mathbf{Y}$ and the trace norm is a witness of the sequentially increasing rank of the solution:



Trace norm and reconstruction errors of $\mathbf{W}^{(t)}$ for $L = 1$ and 2 as a function of t .

Detailed Theorems

Continuous Dynamics: if we initialize with,

- $\mathbf{W}_1(0) = \mathbf{U} \text{diag}(e^{-\delta_1}, \dots, e^{-\delta_p}) \mathbf{Q}$
 - $\mathbf{W}_2(0) = \mathbf{Q}^{-1} \text{diag}(e^{-\delta_1}, \dots, e^{-\delta_d}) \mathbf{V}^\top$
- \mathbf{Q} is an arbitrary invertible matrix. Then,

$$\mathbf{W}_1(t) = \mathbf{W}_1^0(t) + \mathbf{W}_1^\epsilon(t)$$

$$\mathbf{W}_2(t) = \mathbf{W}_2^0(t) + \mathbf{W}_2^\epsilon(t)$$

$$\mathbf{W}_1^0(t) := \mathbf{U} \text{diag}(\sqrt{w_1(t)}, \dots, \sqrt{w_p(t)}) \mathbf{Q}$$

$$\mathbf{W}_2^0(t) := \mathbf{Q}^{-1} \text{diag}(\sqrt{w_1(t)}, \dots, \sqrt{w_d(t)}) \mathbf{V}^\top$$

where $\|\mathbf{W}_i^\epsilon(t)\| \leq \epsilon \cdot e^{ct^2}$ and,

$$w_i(t) = \frac{\sigma_i e^{2\sigma_i t - 2\delta_i}}{\lambda_i (e^{2\sigma_i t - 2\delta_i} - e^{-2\delta_i}) + \sigma_i}$$

(σ_i) and (λ_i) are the diagonals of \mathbf{D}_x and \mathbf{D}_{xy} .

Discrete dynamics: under $\epsilon = 0$, we have

$$w_i^{(t)} \geq \frac{w_i^{(0)}}{(\sigma_i - \lambda_i w_i^{(0)}) e^{(-2\eta \sigma_i + 4\eta^2 \sigma_i^2)t} + w_i^{(0)} \lambda_i}$$

$$w_i^{(t)} \leq \frac{w_i^{(0)}}{(\sigma_i - \lambda_i w_i^{(0)}) e^{(-2\eta \sigma_i - \eta^2 \sigma_i^2)t} + w_i^{(0)} \lambda_i},$$

Differences with the continuous case are in red.

References

- [1] S. Gunasekar, B. E. Woodworth, S. Bhojanapalli, B. Neyshabur, and N. Srebro. Implicit regularization in matrix factorization. In *NIPS*, 2017.
- [2] A. M. Saxe, J. L. McClelland, and S. Ganguli. A mathematical theory of semantic development in deep neural networks. *arXiv preprint arXiv:1810.10531*, 2018.