

Implicit Regularization in Deep Learning May Not Be Explainable by Norms Neural Information Processing Systems (NeurIPS) 2020 Noam Razin Nadav Cohen

I) Implicit Regularization in Deep Learning (DL)

DNNs generalize w/o explicit regularization when # of weights \gg training set size





Conventional wisdom: gradient descent (GD) induces an implicit regularization

Motivating example: in linear regression, when # weights > # training examples GD initialized at 0 converges to min ℓ_2 norm solution

Widespread hope: GD in DL finds solutions minimizing some norm

 $\operatorname{argmin}_{\mathbf{w}} \|\mathbf{w}\|$ s.t. w is global min

II) Case Study: Matrix Factorization

Matrix completion: recover low-rank matrix given subset of entries

	Avenuens	THEPRESTIGE	NOW YOU SEE ME	LEONARDO DICAPRIO THE WOLF OF WALL STREET	
Bob	4	?	?	4 ←	Observations
Alice	?	5	4 🔶	?	
Joe	?	5	?	?	

Parameterize solution as linear neural network and minimize ℓ_2 loss with GD:



Product matrix $W_{L:1} \coloneqq W_L \cdots W_2 W_1$

Objective: $\min_{W_1...W_L} \sum_{(i,j)\in\Omega} \left((W_L W_{L-1} \cdots W_1)_{ij} - b_{ij} \right)^2$

III) Open Question

Does the implicit regularization in matrix factorization minimize a norm?

Opposing Conjectures

- ► Gunasekar et al. 2017: nuclear norm¹ is minimized
- Arora et al. 2019: no norm is minimized

¹ Sum of singular values

 $\left\{b_{ij}\right\}_{(i,j)\in\Omega}$

hidden dims do not necessarily constrain rank

IV) Our Main Contribution: Resolving Open Question

Theorem

There exist matrix factorization settings where:

- \blacktriangleright All norms (and quasi-norms) are driven towards ∞
- Rank is essentially minimized

Implicit regularization in matrix factorization \neq norm minimization

V) Analysis: Implicit Regularization Can Drive All Norms to Infinity

A Simple Matrix Completion Problem

unobserved entry

*	1
1	0

What are the min norm solutions?

- Min nuclear/Frobenius/spectral norm solution $\iff * = 0$

When is rank minimized?

- When $|*| \rightarrow \infty$ distance from rank 1 is minimized
- When * = 0 distance from rank 1 is maximal

Theorem

over standard inits), then $|*| \ge \Omega(1/\sqrt{loss(t)})$. This implies:

- $\|W_{L:1}(t)\| \ge \Omega(1/\sqrt{loss(t)})$, for any $\|\cdot\|$
- ▶ Distance from rank 1 of $W_{L:1} \leq \mathcal{O}(\sqrt{loss(t)})$

Experiment

GD over the construction above, generalized to different matrix dimensions:



Tel Aviv University

-

construction is easily generalized to $\frac{1}{2}$ arbitrary dimensions and different configurations of observed entries

Minimizing an arbitrary norm (or quasi-norm) requires value of * to be bounded

We construct settings where norm and rank minimization are contradictory

Under gradient flow (GD with learning rate \rightarrow 0), if det($W_{L:1}(0)$) > 0 (holds w.p. 0.5)

GD drives all norms to infinity in favor of minimizing rank



Theory transfers to practice: $|*|
ightarrow \infty$



Language of standard norm regularizers might not suffice



Linear Neural Networks In our construction: implicit regularization provably minimizes matrix rank

Arora et al. 2019: theory and experiments suggest this holds in general

Does rank minimization extend beyond matrix factorization (linear NN)?





linear model \Rightarrow high rank tensor factorizations (tf) \Rightarrow low rank

Implicit rank minimization occurs in tensor factorization (non-linear NN)

Implicit regularization in DL minimizes rank of input-output mapping

Rank minimization may be key to explaining generalization in DL

To understand implicit regularization in deep learning:

Developing notions of rank for input-output mappings of NNs may be key