

Implicit Regularization in Deep Learning May Not Be Explainable by Norms

Noam Razin

based on joint work with Nadav Cohen

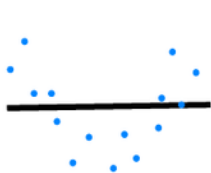
Tel Aviv University

Outline

- 1 Implicit Regularization in Deep Learning
- 2 Case Study: Matrix Factorization
- 3 Implicit Regularization Can Drive All Norms to Infinity
- 4 Implicit Regularization = Rank Minimization?
- 5 Conclusion

Generalization via Bias-Variance Tradeo

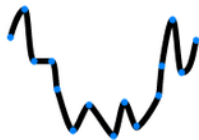
In classical learning theory generalization exhibits the bias-variance tradeoff



Underfitting



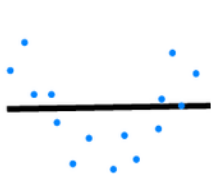
Desired



Overfitting

Generalization via Bias-Variance Tradeo

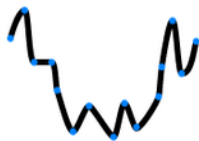
In classical learning theory generalization exhibits the bias-variance tradeoff



Underfitting



Desired

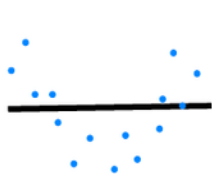


Overfitting

Tradeoff can be controlled through regularization:

Generalization via Bias-Variance Tradeoff

In classical learning theory generalization exhibits the bias-variance tradeoff



Underfitting



Desired



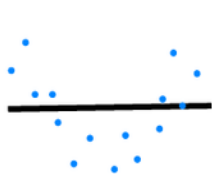
Overfitting

Tradeoff can be controlled through regularization:

- 1 Limiting model size

Generalization via Bias-Variance Tradeoff

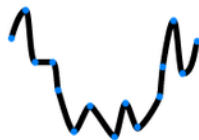
In classical learning theory generalization exhibits the bias-variance tradeoff



Underfitting



Desired



Overfitting

Tradeoff can be controlled through regularization:

- 1 Limiting model size
- 2 Adding term to loss (typically a norm)

Generalization in Deep Learning (DL)

DNNs In Practice

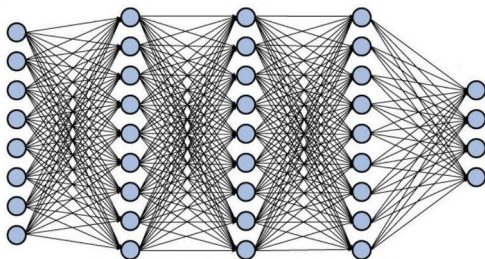
Generalize **without explicit regularization**:

Generalization in Deep Learning (DL)

DNNs In Practice

Generalize **without explicit regularization**:

- ① # of learned weights # of training examples

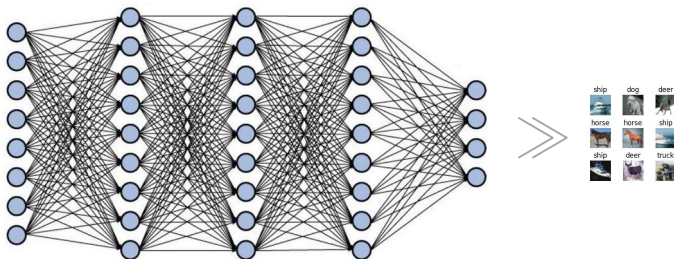


Generalization in Deep Learning (DL)

DNNs In Practice

Generalize **without explicit regularization**:

- ① # of learned weights # of training examples



- ② Loss unchanged (e.g. no weight decay/dropout)

Optimization Induces an Implicit Regularization

Multiple global minima: some generalize well, others don't

Optimization Induces an Implicit Regularization

Multiple global minima: some generalize well, others don't

Solution found by Gradient Descent (GD) often generalizes well

Multiple global minima: some generalize well, others don't

Solution found by Gradient Descent (GD) often generalizes well

Conventional Wisdom

Gradient-based optimization induces an implicit regularization

Multiple global minima: some generalize well, others don't

Solution found by Gradient Descent (GD) often generalizes well

Conventional Wisdom

Gradient-based optimization induces an implicit regularization

Question

Can we mathematically understand this effect in concrete settings?

Linear Regression

When # parameters > # training examples:

Linear Regression

When # parameters > # training examples:

GD initialized at 0 converges to $\min_w \|w\|_2$ norm solution

$$\operatorname{argmin}_w \|w\|_2 \text{ s.t. } Xw = y$$

Widespread Hope

GD in DL finds solutions with **minimal norm (or quasi-norm)**

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

Widespread Hope

GD in DL finds solutions with **minimal norm (or quasi-norm)**

$$\underset{w}{\operatorname{argmin}} w \text{ s.t. } w \text{ is global min}$$

Demonstrated in various settings, e.g.:

Neyshabur et al. 2015

Gunasekar et al. 2017

Soudry et al. 2018

Gunasekar et al. 2018a

Gunasekar et al. 2018b

Li et al. 2018

Jacot et al. 2018

Mei et al. 2019

Ji & Telgarsky 2019a

Ji & Telgarsky 2019b

Wu et al. 2019

Oymak & Soltanolkotabi 2019

Nacson et al. 2019a

Nacson et al. 2019b

Woodworth et al. 2020

Lyu & Li 2020

Ali et al. 2020

Chizat & Bach 2020

Belabbas 2020

- 1 Implicit Regularization in Deep Learning
- 2 **Case Study: Matrix Factorization**
- 3 Implicit Regularization Can Drive All Norms to Infinity
- 4 Implicit Regularization = Rank Minimization?
- 5 Conclusion

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

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

observed entries !	training data
unobserved entries !	test data

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

observed entries !	training data
unobserved entries !	test data

Denote observations b_{ij} $g_{(i;j)2}$

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

observed entries !	training data
unobserved entries !	test data

Denote observations b_{ij} $g_{(i;j)2}$

Convex Programming Approach

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

observed entries !	training data
unobserved entries !	test data

Denote observations b_{ij} $g_{(i;j)}^2$

Convex Programming Approach

Find minimal **nuclear norm** solution:

$$\min \|W\|_{\text{nuclear}} \quad \text{s.t.} \quad W_{ij} = b_{ij} \quad g_{(i;j)}^2$$

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

observed entries !	training data
unobserved entries !	test data

Denote observations b_{ij} $g_{(i;j)2}$

Convex Programming Approach

Find minimal **nuclear norm** solution:

$$\min \|W\|_{\text{nuclear}} \quad \text{s.t.} \quad W_{ij} = b_{ij} \quad \delta_{(i;j)2}$$

Perfectly recovers **if observations are sufficiently many** (Kandès & Recht 2008)

Deep Learning Approach

Deep Learning Approach

Parameterize solution as **LNN** and fit observations using GD (overfit)

Deep Learning Approach

Parameterize solution as **LNN** and fit observations using GD (over loss)

GD is run w.r.t. $W_1; \dots; W_L$ over:

$$\sum_{(i,j)} \frac{1}{2} (W_{L:1})_{ij} - b_{ij}^2$$

Deep Learning Approach

Parameterize solution as **LNN** and fit observations using GD (over loss)

GD is run w.r.t. $W_1; \dots; W_L$ over:

$$\ell(W_{L:1}) = \frac{1}{2} \sum_{(i;j)} (W_{L:1})_{ij} - b_{ij}^2$$

To which solutions does **product matrix** $W_{L:1}$ converge?

Deep Learning Approach

Parameterize solution as **LNN** and t observations using GD (over \mathcal{L})

GD is run w.r.t. $W_1; \dots; W_L$ over:

$$\mathcal{L}(W_{L:1}) = \frac{1}{2} \sum_{(i;j) \in \mathcal{I}} (W_{L:1})_{ij} - b_{ij}^2$$

To which solutions does **product matrix** $W_{L:1}$ converge?

Empirical phenomenon low-rank matrix often recovered accurately

Conjecture (Gunasekar et al. 2017)

With small learning rate and init close to the origin, GD over depth matrix factorization converges to **min nuclear norm** solution.

Conjecture (Gunasekar et al. 2017)

With small learning rate and init close to the origin, GD over depth matrix factorization converges to **min nuclear norm** solution.

GD implicitly solves convex programming approach?

Conjecture (Gunasekar et al. 2017)

With small learning rate and init close to the origin, GD over depth matrix factorization converges to **min nuclear norm** solution.

GD implicitly solves convex programming approach?

Gunasekar et al. supported conjecture with:

- Experiments

- Proof for **certain restricted case**

Conjecture (Gunasekar et al. 2017)

With small learning rate and init close to the origin, GD over depth matrix factorization converges to **min nuclear norm** solution.

GD implicitly solves convex programming approach?

Gunasekar et al. supported conjecture with:

- Experiments

- Proof for **certain restricted case**

Conjecture established under other restricted conditions:

- Li et al. 2018

- Belabbas 2020

Conjecture (Arora et al. 2019)

For any norm k , exist observations for which small learning rate and initialization **can not ensure** GD converges to **mirk solution**.

Conjecture (Arora et al. 2019)

For any norm k , exist observations for which small learning rate and initialization **can not ensure** GD converges to **mkirk solution**.

Arora et al. supported conjecture with:

Experiments: **nuclear norm** not always minimized, bias to **low rank**

Conjecture (Arora et al. 2019)

For any norm k , exist observations for which small learning rate and initialization **can not ensure** GD converges to **mirk solution**.

Arora et al. supported conjecture with:

Experiments: **nuclear norm** not always minimized, bias to **low rank**

Dynamical analysis of singular values: GD promotes **sparse spectrum**

Conjecture (Arora et al. 2019)

For any norm k , exist observations for which small learning rate and initialization **can not ensure** GD converges to **mkirk solution**.

Arora et al. supported conjecture with:

Experiments: **nuclear norm** not always minimized, bias to **low rank**

Dynamical analysis of singular values: GD promotes **sparse spectrum**

"
suggests tendency to **low rank**

Conjecture (Arora et al. 2019)

For any norm k , exist observations for which small learning rate and initialization **can not ensure** GD converges to **mirk solution**.

Arora et al. supported conjecture with:

Experiments: **nuclear norm** not always minimized, bias to **low rank**

Dynamical analysis of singular values: GD promotes **sparse spectrum**

" suggests tendency to **low rank**

Open Question

Does the implicit regularization in matrix factorization minimize a norm?

Theorem (informal)

There exist matrix factorization settings where:

Theorem (informal)

There exist matrix factorization settings where:

- † All norms (and quasi-norms) are driven towards 1

Theorem (informal)

There exist matrix factorization settings where:

- 1 All norms (and quasi-norms) are driven towards 1
- 2 Rank is essentially minimized

Theorem (informal)

There exist matrix factorization settings where:

- 1 All norms (and quasi-norms) are driven towards 1
- 2 Rank is essentially minimized

A rms conjecture of Arora et al. 2019

Theorem (informal)

There exist matrix factorization settings where:

- 1 All norms (and quasi-norms) are driven towards 1
- 2 Rank is essentially minimized

A rms conjecture of Arora et al. 2019

Result stronger than conjecture:

Theorem (informal)

There exist matrix factorization settings where:

- 1 All norms (and quasi-norms) are driven towards 1
- 2 Rank is essentially minimized

A rms conjecture of Arora et al. 2019

Result stronger than conjecture:

Settings jointly disqualify all norms

Norms driven towards 1

- 1 Implicit Regularization in Deep Learning
- 2 Case Study: Matrix Factorization
- 3 **Implicit Regularization Can Drive All Norms to Infinity**
- 4 Implicit Regularization = Rank Minimization?
- 5 Conclusion

Common surrogate for GD with small learning rate and initialization

¹(e.g. used by Gunasekar et al. 2017, Arora et al. 2019)

Common surrogate for GD with small learning rate and initial
Gradient flow (GF) is a continuous version of GD (step size 0):

$$\frac{d}{dt} w(t) = -\eta \nabla f(w(t)) ; t \geq 0$$

¹(e.g. used by Gunasekar et al. 2017, Arora et al. 2019)

Common surrogate for GD with small learning rate and ϵ init

Gradient flow (GF) is a continuous version of GD (step size 0):

$$\frac{d}{dt} w(t) = -\eta \nabla f(w(t)) ; t \geq 0$$

Weights $W_1 ::: W_L$ are **balanced** at init: $W_{j+1}^2 = W_j W_j^2 ; \forall j$.

¹(e.g. used by Gunasekar et al. 2017, Arora et al. 2019)

Common surrogate for GD with small learning rate and ϵ init

Gradient flow (GF) is a continuous version of GD (step size 0):

$$\frac{d}{dt} w(t) = -\eta \nabla f(w(t)) ; t \geq 0$$

Weights $W_1 :::: W_L$ are **balanced** at init: $W_{j+1}^2 = W_j W_j^2 ; \forall j$.

Holds approximately under 0 init, exactly under residual 0 init

¹(e.g. used by Gunasekar et al. 2017, Arora et al. 2019)

Common surrogate for GD with small learning rate and ϵ init

Gradient flow (GF) is a continuous version of GD (step size 0):

$$\frac{d}{dt} w(t) = -\eta \nabla f(w(t)) ; t \geq 0$$

Weights $W_1 ::: W_L$ are **balanced** at init: $W_{j+1}^2 = W_j W_j^2 ; \forall j$.

Holds approximately under ϵ init, exactly under residual ϵ init

Closely matches **GD** in practice for **linear neural networks**

¹(e.g. used by Gunasekar et al. 2017, Arora et al. 2019)

Consider the following 2-by-2 matrix completion problem:

$$\begin{array}{cc} & ! \\ & 1 \\ 1 & 0 \end{array}$$

Consider the following 2-by-2 matrix completion problem:

$$\begin{array}{c} ! \\ 1 \\ 1 \quad 0 \end{array}$$

Min Frobenius norm $\| \cdot \|_F = 0$ (depth $L = 1$ solution)

Consider the following 2-by-2 matrix completion problem:

$$\begin{array}{c} ! \\ 1 \\ 1 \quad 0 \end{array}$$

Min Frobenius norm() = 0 (depth L = 1 solution)

Min nuclear norm() = 0

Consider the following 2-by-2 matrix completion problem:

$$\begin{array}{c} \\ \\ 1 \end{array} \quad \begin{array}{c} \\ \\ 0 \end{array}$$

Min Frobenius norm $\| \cdot \|_F = 0$ (depth $L = 1$ solution)

Min nuclear norm $\| \cdot \|_* = 0$

Generally: Schatten p (quasi-)norm $\| \cdot \|_p := \left(\sum_i \sigma_i^p \right)^{1/p}$ (σ_i - singular values)

Consider the following 2-by-2 matrix completion problem:

$$\begin{array}{c} ! \\ 1 \\ 1 \quad 0 \end{array}$$

Min Frobenius norm $\| \cdot \|_F = 0$ (depth $L = 1$ solution)

Min nuclear norm $\| \cdot \|_* = 0$

Generally: Schatten- p (quasi-)norm $\| \cdot \|_p := \left(\sum_i \sigma_i^p \right)^{1/p}$ (σ_i - singular values)

Proposition

Schatten- p (quasi-)norms are minimized $\| \cdot \|_p = 0$

Consider the following 2-by-2 matrix completion problem:

$$\begin{matrix} & 1 \\ 1 & 0 \end{matrix}$$

Min Frobenius norm $\|M\|_F = 0$ (depth $L = 1$ solution)

Min nuclear norm $\|M\|_* = 0$

Generally: Schatten- p (quasi-)norm $\|M\|_p := \left(\sum_i \sigma_i^p \right)^{1/p}$ (σ_i - singular values)

Proposition

Schatten- p (quasi-)norms are minimized $\|M\|_p = 0$

Arbitrary norms (or quasi-norms):

Proposition

For minimal k solutions, $\|M\|_p$ is bounded

$$\begin{array}{c} ! \\ 1 \\ 1 \ 0 \end{array}$$

$$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

(i) All solutions are of rank 2 (ii) $\|j\|_1$ essentially rank 1

$$\begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}$$

(i) All solutions are of **rank 2** (ii) $(j, j+1)$ essentially **rank 1**

Continuous measures for rank:

$$\begin{pmatrix} 1 & \\ & 1 \end{pmatrix}$$

(i) All solutions are of rank 2 (ii) $\sum_{j=1}^n \sigma_j^2 > 0$ essentially rank 1

Continuous measures for rank:

Distance from rank 1 := value of smallest singular value

$$\begin{pmatrix} 1 & \\ & 1 \end{pmatrix}$$

(i) All solutions are of rank 2 (ii) $\sum_{j=1}^n \sigma_j$ essentially rank 1

Continuous measures for rank:

Distance from rank 1 := value of smallest singular value

Effective rank (erank) = entropy of singular values

$$\begin{matrix} & & ! \\ & & 1 \\ 1 & & 0 \end{matrix}$$

(i) All solutions are of **rank 2** (ii) $\sum_{j=1}^n \sigma_j$ essentially **rank 1**

Continuous measures for rank:

Distance from rank 1 := value of smallest singular value

Effective rank (erank) = entropy of singular values

Definition (Rank suboptimality)

rank-subopt := $\max_f \text{Distance from rank 1}; \text{erank} \inf \text{erank}_g$

$$\begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}$$

(i) All solutions are of rank 2 (ii) $\sum_{j=1}^n \sigma_j$ essentially rank 1

Continuous measures for rank:

Distance from rank 1 := value of smallest singular value

Effective rank (erank) = entropy of singular values

Definition (Rank suboptimality)

rank-subopt := $\max_f \text{Distance from rank 1}; \text{erank} = \inf_{\text{rank } k} \text{erank}_k$

Proposition

rank-subopt is maximized when $\sigma = 0$, and is minimized as $\sum_{j=1}^n \sigma_j \rightarrow 1$

$$\begin{matrix} & & ! \\ & & 1 \\ 1 & & 0 \end{matrix}$$

(i) All solutions are of **rank 2** (ii) $\sum_{j \neq 1} \sigma_j$ essentially **rank 1**

Continuous measures for rank:

Distance from rank 1 := value of smallest singular value

Effective rank (erank) = entropy of singular values

Definition (Rank suboptimality)

rank-subopt := $\max_f \text{Distance from rank 1}; \text{erank} = \inf \text{erank}_g$

Proposition

rank-subopt is maximized when $\sigma = 0$, and is minimized as $\sum_{j \neq 1} \sigma_j \rightarrow 1$

Contradiction between norm and rank minimization

Theorem

Theorem

If $\det(W_{L:1}(0)) > 0$ at init, then for any (quasi-)norm $\|\cdot\|_k$:

Theorem

If $\det(W_{L:1}(0)) > 0$ at init, then for any (quasi-)norm $\|\cdot\|_k$:

$$\|W_{L:1}(t)\|_k = \left(\frac{1}{\det(W_{L:1}(t))} \right)^{\frac{1}{k}}$$

Theorem

If $\det(W_{L:1}(0)) > 0$ at init, then for any (quasi-)norm $\|\cdot\|_k$:

- 1 $\|W_{L:1}(t)\|_k = O(t^{-\frac{p}{k}})$
- 2 $\text{rank-subopt}(W_{L:1}(t)) = O(t^{-\frac{p}{k}})$

Theorem

If $\det(W_{L:1}(0)) > 0$ at init, then for any (quasi-)norm $\|\cdot\|_k$:

- 1 $\|W_{L:1}(t)\|_k = O(\frac{1}{t^{\frac{1}{k}}})$
- 2 $\text{rank-subopt}(W_{L:1}(t)) = O(\frac{1}{t^{\frac{1}{k}}})$

$\frac{1}{t^{\frac{1}{k}}} \rightarrow 0$ implies:

All norms (and quasi-norms) driven towards 1

Rank is essentially minimized

Theorem

If $\det(W_{L:1}(0)) > 0$ at init, then for any (quasi-)norm $\|\cdot\|_k$:

- $\|W_{L:1}(t)\|_k = (1 - \frac{t}{T})^{\frac{1}{k}}$
- $\text{rank-subopt}(W_{L:1}(t)) = O\left(\frac{1}{(1 - \frac{t}{T})^{\frac{1}{k}}}\right)$

$\frac{1}{k} \rightarrow 0$ implies:

All **norms (and quasi-norms)** driven towards 1

Rank is essentially minimized

Claim

Assumption on $\det(W_{L:1}(0)) > 0$ holds with prob 0.5 under standard inits

Theorem

If $\det(W_{L:1}(0)) > 0$ at init, then for any (quasi-)norm $\|\cdot\|_k$:

- 1 $\|W_{L:1}(t)\|_k = O(\frac{1}{\sqrt{t}})$
- 2 $\text{rank-subopt}(W_{L:1}(t)) = O(\frac{1}{\sqrt{t}})$

$\frac{1}{\sqrt{t}} \rightarrow 0$ implies:

All **norms (and quasi-norms)** driven towards 1

Rank is essentially minimized

Claim

Assumption on $\det(W_{L:1}(0))$ holds with prob 0.5 under standard inits

Implicit regularization $\hat{=}$ norm minimization!

Theorem

If $\det(W_{L:1}(0)) > 0$ at init, then for any (quasi-)norm $\|\cdot\|_k$:

- 1 $\|W_{L:1}(t)\|_k = O(t^{-\frac{p}{k}})$
- 2 $\text{rank-subopt}(W_{L:1}(t)) = O(t^{-\frac{p}{k}})$

Proof Sketch

¹based on singular values differential equations from Arora et al. 2019

Theorem

If $\det(W_{L:1}(0)) > 0$ at init, then for any (quasi-)norm $\|\cdot\|_k$:

- 1 $\|W_{L:1}(t)\|_k = O(\sqrt[p]{t})$
- 2 $\text{rank-subop}(W_{L:1}(t)) = O(\sqrt[p]{t})$

Proof Sketch

GF trajectory analysis $\det(W_{L:1}(t))$ does not change sign¹

¹based on singular values differential equations from Arora et al. 2019

Theorem

If $\det(W_{L:1}(0)) > 0$ at init, then for any (quasi-)norm $\|\cdot\|_k$:

- 1 $\|W_{L:1}(t)\|_k = O(\sqrt[p]{t})$
- 2 $\text{rank-subopt}(W_{L:1}(t)) = O(\sqrt[p]{t})$

Proof Sketch

GF trajectory analysis $\det(W_{L:1}(t))$ does not change sign¹

$$\det(W_{L:1}(t)) = w_{1;1}(t)w_{2;2}(t) \quad w_{1;2}(t)w_{2;1}(t) > 0$$

¹based on singular values differential equations from Arora et al. 2019

Theorem

If $\det(W_{L:1}(0)) > 0$ at init, then for any (quasi-)norm k :

- 1 $\|W_{L:1}(t)\|_k = O(t^{-\frac{p}{k}})$
- 2 $\text{rank-subopt}(W_{L:1}(t)) = O(t^{-\frac{p}{k}})$

Proof Sketch

GF trajectory analysis $\det(W_{L:1}(t))$ does not change sign¹

$$\det(W_{L:1}(t)) = w_{1;1}(t) \underbrace{w_{2;2}(t)}_{! \ 0} \underbrace{w_{1;2}(t)w_{2;1}(t)}_{! \ 1} > 0$$

¹based on singular values differential equations from Arora et al. 2019

Theorem

If $\det(W_{L:1}(0)) > 0$ at init, then for any (quasi-)norm k :

- 1 $\|W_{L:1}(t)\|_k = O(\sqrt[p]{t})$
- 2 $\text{rank-subopt}(W_{L:1}(t)) = O(\sqrt[p]{t})$

Proof Sketch

GF trajectory analysis $\det(W_{L:1}(t))$ does not change sign

$$\det(W_{L:1}(t)) = w_{1;1}(t) \underbrace{w_{2;2}(t)}_{! \ 0} \underbrace{w_{1;2}(t)w_{2;1}(t)}_{! \ 1} > 0$$

$$\Rightarrow \|w_{1;1}(t)\|_1 \quad (\text{behaves as } O(\sqrt[p]{t}))$$

¹based on singular values differential equations from Arora et al. 2019

Theorem

If $\det(W_{L:1}(0)) > 0$ at init, then for any (quasi-)norm k :

- $kW_{L:1}(t)k = O(\sqrt[p]{t})$
- $\text{rank-subopt}(W_{L:1}(t)) = O(\sqrt[p]{t})$

Proof Sketch

GF trajectory analysis $\det(W_{L:1}(t))$ does not change sign¹

$$\det(W_{L:1}(t)) = w_{1;1}(t) \underbrace{w_{2;2}(t)}_{! \ 0} \underbrace{w_{1;2}(t)w_{2;1}(t)}_{! \ 1} > 0$$

$\Rightarrow |w_{1;1}(t)| \neq 0$ (behaves as $O(\sqrt[p]{t})$)

Bound on $|w_{1;1}(t)|$ implies bounds for norms and rank suboptimality

¹based on singular values differential equations from Arora et al. 2019

Customary

Separating aspects of convergence to global min and implicit regularization

Customary

Separating aspects of convergence to global min and implicit regularization

"

commonly observed in practice

Customary

Separating aspects of convergence to global min and implicit regularization

"

commonly observed in practice

Convergence to Zero Loss (in our setting)

Customary

Separating aspects of convergence to global min and implicit regularization

"

commonly observed in practice

Convergence to Zero Loss (in our setting)

Experiments: GD consistently finds global min

Customary

Separating aspects of convergence to global min and implicit regularization

"

commonly observed in practice

Convergence to Zero Loss (in our setting)

Experiments: GD consistently finds global min

Proof for depth 2 with scaled identity init

Customary

Separating aspects of convergence to global min and implicit regularization

"

commonly observed in practice

Convergence to Zero Loss (in our setting)

Experiments: GD consistently finds global min

Proof for depth 2 with scaled identity init

Proposition

If at init $W_{L:1}(0) = I$, for depth $L = 2$ and $0 < \epsilon < 1$, then $\|W(t)\| \leq \epsilon$

Customary

Separating aspects of convergence to global min and implicit regularization

"

commonly observed in practice

Convergence to Zero Loss (in our setting)

Experiments: GD consistently finds global min

Proof for depth 2 with scaled identity init

Proposition

If at init $W_{L:1}(0) = I$, for depth $L = 2$ and $0 < \eta < 1$, then $\| \nabla \ell(t) \| \rightarrow 0$

Proof Approach

Careful analysis of GF differential equations

What happens when observations are perturbed?

$$\begin{array}{c} ! \\ 1 \\ 1 \quad 0 \end{array}$$

What happens when observations are **perturbed**?

$$\begin{matrix} & ! \\ & 1 \\ 1 & 0 \end{matrix} =) \begin{matrix} & ! \\ & z \\ z^0 & \end{matrix} \begin{matrix} \text{non-zero } z; z^0 \\ \text{arbitrary} \end{matrix}$$

What happens when observations are **perturbed**?

$$\begin{pmatrix} 1 & z^0 \\ 0 & z^0 \end{pmatrix} = \begin{pmatrix} 1 & z^0 \\ 0 & z^0 \end{pmatrix} \quad \begin{matrix} \text{non-zero } z^0; z^0 \\ \text{arbitrary} \end{matrix}$$

Theorem (original setting)

If $\det(W_{L:1}(0)) > 0$ at init, then for any (quasi-)norm $\|\cdot\|_k$:

- $\|W_{L:1}(t)\|_k = O(\|1 - \overline{\lambda}(t)\|_k^p)$
- $\text{rank-subopt}(W_{L:1}(t)) = O(\|1 - \overline{\lambda}(t)\|_k^p)$

What happens when observations are **perturbed**?

$$\begin{pmatrix} 1 & z^0 \\ 0 & z^0 \end{pmatrix} = \begin{pmatrix} 1 & z^0 \\ 0 & z^0 \end{pmatrix} \quad \begin{matrix} \text{non-zero } z; z^0 \\ \text{arbitrary} \end{matrix}$$

Theorem (original setting)

If $\text{sign}(\det(W_{L:1}(0))) = \text{sign}(1 - 1)$ at init, then for any (quasi-)norm $\|\cdot\|_k$:

- $\|W_{L:1}(t)\|_k = O(\sqrt[p]{t})$
- $\text{rank-subopt}(W_{L:1}(t)) = O(\sqrt[p]{t})$

What happens when observations are **perturbed**?

$$\begin{pmatrix} 1 & z \\ 0 & z^0 \end{pmatrix} \Rightarrow \begin{pmatrix} 1 & z \\ z^0 & z^0 \end{pmatrix} \quad \begin{array}{l} \text{non-zero } z; z^0 \\ \text{arbitrary} \end{array}$$

Theorem (intermediate)

If $\text{sign}(\det(W_{L:1}(0))) = \text{sign}(z - z^0)$ at init, then for any (quasi-)norm k :

- $\|W_{L:1}(t)\|_k = O(\|z - z^0\|_k^{\frac{1}{p}})$
- $\text{rank-subopt}(W_{L:1}(t)) = O(\|z - z^0\|_k^{\frac{1}{p}})$

What happens when observations are **perturbed**?

$$\begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix} \Rightarrow \begin{pmatrix} z^0 & z^0 \\ z^0 & z^0 \end{pmatrix} \quad \begin{array}{l} \text{non-zero } z; z^0 \\ \text{arbitrary} \end{array}$$

Theorem (intermediate)

If $\text{sign}(\det(W_{L:1}(0))) = \text{sign}(z - z^0)$ at init, then for any (quasi-)norm k :

- $kW_{L:1}(t)k = \min_j |z_j| = (|z - z^0| + \frac{P}{2\sqrt{t}})$
- $\text{rank-subop}(W_{L:1}(t)) = O(\frac{P}{\sqrt{t}})$

What happens when observations are **perturbed**?

$$\begin{pmatrix} 1 & z^0 \\ 0 & z^0 \end{pmatrix} = \begin{pmatrix} 1 & z^0 \\ z^0 & z^0 \end{pmatrix} \quad \begin{array}{l} \text{non-zero } z^0; z^0 \\ \text{arbitrary} \end{array}$$

Theorem

If $\text{sign}(\det(W_{L:1}(0))) = \text{sign}(z^0)$ at init, then for any (quasi-)norm k :

- $kW_{L:1}(t)k = \min_j \{z_j; |z_j^0|\} = (j + j + \frac{p}{2}(t))$
- $\text{rank-subop}(W_{L:1}(t)) = O(j + j + \frac{p}{2}(t)) = \min_j \{z_j; |z_j^0|\}$

What happens when observations are **perturbed**?

$$\begin{matrix} & & ! & & & & ! \\ & & 1 & & & & z \\ 1 & & 0 & =) & & & z^0 \end{matrix} \quad \begin{matrix} & & & & & & \text{non-zero } z; z^0 \\ & & & & & & \text{arbitrary} \end{matrix}$$

Theorem

If $\text{sign}(\det(W_{L:1}(0))) = \text{sign}(z \cdot z^0)$ at init, then for any (quasi-)norm k :

$$1 \quad \|k W_{L:1}(t)\| = \min_j \{z_j; |z^0_j|\} = (j \cdot j + \frac{P}{2^j(t)})$$

$$2 \quad \text{rank-subopt}(W_{L:1}(t)) = O(j \cdot j + \frac{P}{2^j(t)}) = \min_j \{z_j; |z^0_j|\}$$

= 0) all **norms** driven to 1 and **rank** is minimized

What happens when observations are **perturbed**?

$$\begin{pmatrix} 1 & z^0 \\ 0 & z^0 \end{pmatrix} \Rightarrow \begin{pmatrix} 1 & z^0 \\ z^0 & z^0 \end{pmatrix} \quad \begin{array}{l} \text{non-zero } z^0; z^0 \\ \text{arbitrary} \end{array}$$

Theorem

If $\text{sign}(\det(W_{L:1}(0))) = \text{sign}(z^0)$ at init, then for any (quasi-)norm k :

$$1 \quad \|W_{L:1}(t)\|_k = \min_j \{z_j; |z_j^0|\} = (j + j + \frac{p}{2} \overline{z}(t))$$

$$2 \quad \text{rank-subopt}(W_{L:1}(t)) = O(j + j + \frac{p}{2} \overline{z}(t)) = \min_j \{z_j; |z_j^0|\}$$

$= 0$) all **norms** driven to 1 and **rank** is minimized

Phenomenon **gracefully recedes** **perturbed** from 0

What happens when observations are **perturbed**?

$$\begin{matrix} 1 \\ 1 & 0 \end{matrix} \Rightarrow \begin{matrix} 1 \\ z^0 \end{matrix} \begin{matrix} z \\ z^0 \end{matrix} \quad \begin{matrix} \text{non-zero } z; z^0 \\ \text{arbitrary} \end{matrix}$$

Theorem

If $\text{sign}(\det(W_{L:1}(0))) = \text{sign}(z \ z^0)$ at init, then for any (quasi-)norm k :

$$1 \quad \|W_{L:1}(t)\|_k = \min_j \{z_j; |z^0_j|\} = (j \ j + \frac{p}{2} \overline{z}(t))$$

$$2 \quad \text{rank-subopt}(W_{L:1}(t)) = O \left((j \ j + \frac{p}{2} \overline{z}(t)) = \min_j \{z_j; |z^0_j|\} \right)$$

$= 0$) all **norms** driven to 1 and **rank** is minimized

Phenomenon **gracefully recedes** **perturbed** from 0

Same results hold when changing **unobserved entry location**

!

1

1 0

!

1

1 0

$$\begin{matrix} & & ! \\ & & 1 \\ 1 & & 0 \end{matrix}$$

$$\begin{matrix} & & ! \\ & & 2 \\ 0:5 & & \end{matrix}$$

$$\begin{matrix} & & ! \\ & & 1 \\ 1 & & 0 \end{matrix}$$

$$\begin{matrix} & & ! \\ & & 2 \\ 0:5 & & \end{matrix}$$

$$\begin{array}{c} ! \\ 1 \\ 1 \quad 0 \end{array}$$

$$\begin{array}{c} ! \\ 2 \\ 0:5 \end{array}$$

Theory transfers to practice: unobserved entry ! 1

- 1 Implicit Regularization in Deep Learning
- 2 Case Study: Matrix Factorization
- 3 Implicit Regularization Can Drive All Norms to Infinity
- 4 **Implicit Regularization = Rank Minimization?**
- 5 Conclusion

Analyzed Setting (our work)

Analyzed Setting (our work)

Contrast between **norm** and **rank** minimization

Analyzed Setting (our work)

Contrast between **norm** and **rank** minimization

Implicit regularization drives **norms** to ∞ to minimize **rank**

Analyzed Setting (our work)

Contrast between **norm** and **rank** minimization

Implicit regularization drives **norms** to ∞ to minimize **rank**

Past Work

Analyzed Setting (our work)

Contrast between **norm** and **rank** minimization

Implicit regularization drives **norms** to ∞ to minimize **rank**

Past Work

Empirical evidence: **low-rank** tendency in matrix factorization

Analyzed Setting (our work)

Contrast between **norm** and **rank** minimization

Implicit regularization drives **norms** to ℓ_1 to minimize **rank**

Past Work

Empirical evidence: **low-rank** tendency in matrix factorization

Theoretical analysis: **sparsity in singular values** (Arora et al. 2019)

Analyzed Setting (our work)

Contrast between **norm** and **rank** minimization

Implicit regularization drives **norms** to ∞ to minimize **rank**

Past Work

Empirical evidence: **low-rank** tendency in matrix factorization

Theoretical analysis: **sparsity in singular values** (Arora et al. 2019)

Better interpretation rank minimization?

Analyzed Setting (our work)

Contrast between **norm** and **rank** minimization

Implicit regularization drives **norms** to ∞ to minimize **rank**

Past Work

Empirical evidence: **low-rank** tendency in matrix factorization

Theoretical analysis: **sparsity in singular values** (Arora et al. 2019)

Better interpretation rank minimization?

Does this interpretation extend **beyond matrix factorization**?

ConvACs are competitive in practice, and admit algebraic structure

Extensively studied (e.g. Cohen et al. 2016, Cohen & Shashua 2016, Cohen & Shashua 2017)

ConvACs are competitive in practice, and admit algebraic structure

Extensively studied (e.g. Cohen et al. 2016, Cohen & Shashua 2016, Cohen & Shashua 2017)

Tensor factorizations correspond to non-linear NN

Tensor completion : recover low-rank tensor given subset of entries

Tensor completion : recover low-rank tensor given subset of entries

Natural extension of matrix completion

Tensor completion : recover low-rank tensor given subset of entries

Natural extension of matrix completion

Tensor Basics

Tensor completion : recover **low-rank** tensor given subset of entries

Natural extension of matrix completion

Tensor Basics

Tensor N-dimensional array (N = **order** of tensor)

Tensor completion : recover **low-rank** tensor given subset of entries

Natural extension of matrix completion

Tensor Basics

Tensor N-dimensional array ($N =$ **order** of tensor)

Tensor rank minimal R s.t. $W = \prod_{r=1}^R w_r^{(1)} \dots w_r^{(N)}$
 $:=$ outer product , $w_r^{(i)} \in \mathbb{R}^{d_i}$

Tensor completion : recover **low-rank** tensor given subset of entries

Natural extension of matrix completion

Tensor Basics

Tensor N-dimensional array ($N =$ **order** of tensor)

Tensor rank minimal R s.t. $W = \prod_{r=1}^R w_r^{(1)} \otimes \dots \otimes w_r^{(N)}$
 \otimes := outer product , $w_r^{(i)} \in \mathbb{R}^{d_i}$

For $N = 2$ this is exactly matrix rank

Parameterize solution as **tensor factorization** :

$$W = \sum_{r=1}^R w_r^{(1)} \otimes \dots \otimes w_r^{(N)}$$

Parameterize solution as **tensor factorization** :

$$W = \sum_{r=1}^R w_r^{(1)} \otimes \dots \otimes w_r^{(N)}$$

R taken large enough to **not constrain rank**

Parameterize solution as **tensor factorization** :

$$W = \prod_{r=1}^R w_r^{(1)} \quad w_r^{(N)}$$

R taken large enough to **not constrain rank**

Does W converge to **low-rank** tensor when running **GD** w.r.t. $w_{r;n}^{(n)}$?

Order 4 Rank 1 Tensor Completion

Order 4 Rank 1 Tensor Completion

"linear" baseline exactly t s observations, 0 elsewhere

Order 4 Rank 1 Tensor Completion

"linear" baseline exactly t s observations, 0 elsewhere

GD drives rank of a non-linear NN towards minimum!

Theory & Experiments implicit regularization minimizes **matrix rank**

Theory & Experiments implicit regularization minimizes **matrix rank**

Theory & Experiments: implicit regularization minimizes **matrix rank**

Experiments: implicit regularization minimizes **tensor rank**

Theory & Experiments: implicit regularization minimizes **matrix rank**

Experiments: implicit regularization minimizes **tensor rank**

Hypothesis

Implicit regularization in DL minimizes **rank of input-output mapping**

Theory & Experiments: implicit regularization minimizes **matrix rank**

Experiments: implicit regularization minimizes **tensor rank**

Hypothesis

Implicit regularization in DL minimizes **rank of input-output mapping**

If true, may be key to explaining generalization

- 1 Implicit Regularization in Deep Learning
- 2 Case Study: Matrix Factorization
- 3 Implicit Regularization Can Drive All Norms to Infinity
- 4 Implicit Regularization = Rank Minimization?
- 5 Conclusion

Implicit Regularization & Norm Minimization

Matrix factorization: exist cases where **all norms go to 1**

Implicit Regularization & Norm Minimization

Matrix factorization: exist cases where ℓ_1 norms go to 1

Unlikely implicit regularization in DL_{∞} norm minimization

Implicit Regularization & Norm Minimization

Matrix factorization: exist cases where ℓ_1 norms go to 1

Unlikely implicit regularization in DL_{∞} norm minimization

Better Interpretation: Bias to Low Rank?

Conclusion

Implicit Regularization = Norm Minimization

- Matrix factorization: exist cases where **all norms go to**
- Unlikely implicit regularization in DL = norm minimization

Better Interpretation: Bias to Low Rank?

- Matrix factorization: growing empirical and theoretical evidence

Conclusion

Implicit Regularization = Norm Minimization

- Matrix factorization: exist cases where **all norms go to**
- Unlikely implicit regularization in DL = norm minimization

Better Interpretation: Bias to Low Rank?

- Matrix factorization: growing empirical and theoretical evidence
- Extends to certain type of **non-linear NN**

Conclusion

Implicit Regularization = Norm Minimization

- Matrix factorization: exist cases where **all norms go to**
- Unlikely implicit regularization in DL = norm minimization

Better Interpretation: Bias to Low Rank?

- Matrix factorization: growing empirical and theoretical evidence
- Extends to certain type of **non-linear NN**

Looking Forward

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

Conclusion

Implicit Regularization = Norm Minimization

- Matrix factorization: exist cases where **all norms go to**
- Unlikely implicit regularization in DL = norm minimization

Better Interpretation: Bias to Low Rank?

- Matrix factorization: growing empirical and theoretical evidence
- Extends to certain type of **non-linear NN**

Looking Forward

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

Thank You