Implicit Regularization in Tensor Factorization

Noam Razin* Asaf Maman* Naday Cohen

*Equal contribution

Tel Aviv University



ICMLi 2021



Neural networks generalize with no explicit regularization even when:



of learned weights

Neural networks generalize with no explicit regularization even when:



Conventional Wisdom

Neural networks generalize with no explicit regularization even when:

Conventional Wisdom

GD induces implicit regularization towards low "complexity" predictors

<u>Goal</u>

Mathematically understand this implicit regularization

<u>Goal</u>

Mathematically understand this implicit regularization

<u>Why?</u>

• Explain why some models generalize while others don't

<u>Goal</u>

Mathematically understand this implicit regularization

<u>Why?</u>

- Explain why some models generalize while others don't
- Select a model with implicit bias that captures data properties

<u>Goal</u>

Mathematically understand this implicit regularization

<u>Why?</u>

- Explain why some models generalize while others don't
- Select a model with implicit bias that captures data properties
- Design models with new types of implicit regularization¹

¹e.g. Jing, Zbontar and LeCun 2020

<u>Goal</u>

Mathematically understand this implicit regularization

<u>Why?</u>

- Explain why some models generalize while others don't
- Select a model with implicit bias that captures data properties
- Design models with new types of implicit regularization¹

Challenge

Lack complexity measures that capture essence of natural data

X high complexity

¹e.g. Jing, Zbontar and LeCun 2020

Matrix completion: recover unknown matrix given subset of entries

	Avenuens	THEPRESTIGE	NOW YOU SEE ME	THE WOLF	
Bob	4	?	?	4 🦟	observations $\{v_{ij}\}$
Alice	?	5	4 🗸	?	$(j,i) \in \Omega$
Joe	?	5	?	?	

Matrix completion: recover unknown matrix given subset of entries

	Avenuens	THEPRESTIGE	NOW YOU SEE ME	THE WOLF	
Bob	4	?	?	4 🦟	observations $\{v_{ij}\}$
Alice	?	5	4 🗸	?	$(j,i) \in \Omega$
Joe	?	5	?	?	

 $d \times d'$ matrix completion \longleftrightarrow prediction from $\{1, ..., d\} \times \{1, ..., d'\}$ to \mathbb{R}

Matrix completion: recover unknown matrix given subset of entries

	Avenuens	THEPRESTIGE	NOW YOU SEE ME	THE WOLF	
Bob	4	?	?	4 ~	observations $\{v_{ij}\}$
Alice	?	5	4 👡	?	$(j_i,j) \in \Omega$
Joe	?	5	?	?	

 $d \times d'$ matrix completion \longleftrightarrow prediction from $\{1, ..., d\} \times \{1, ..., d'\}$ to \mathbb{R}

value of entry $(i,j) \iff$ label of input (i,j)

Matrix completion: recover unknown matrix given subset of entries

	Avenuens	THEPRESTIGE	NOW YOU SEE ME	THE WOLF	
Bob	4	?	?	4 ~	observations $\{v_{ij}\}$
Alice	?	5	4 👡	?	$(j_i,j) \in \Omega$
Joe	?	5	?	?	

 $d \times d'$ matrix completion \longleftrightarrow prediction from $\{1,...,d\} \times \{1,...,d'\}$ to \mathbb{R}

value of entry $(i,j) \iff$ label of input (i,j)

observed entries \iff train data

Matrix completion: recover unknown matrix given subset of entries

	Avenuens	THEPRESTIGE	NOW YOU SEE ME	THE WOLF	
Bob	4	?	?	4 🦟	observations $\{v_{ij}\}$
Alice	?	5	4 🗸	?	$(j,i) \in \Omega$
Joe	?	5	?	?	

 $d \times d'$ matrix completion \longleftrightarrow prediction from $\{1,...,d\} \times \{1,...,d'\}$ to \mathbb{R}

value of entry $(i,j) \iff$ label of input (i,j)

observed entries \iff train data

unobserved entries \longleftrightarrow test data

Matrix completion: recover unknown matrix given subset of entries

	Avenuens	THEPRESTIGE	NOW YOU SEE ME	THE WOLF	
Bob	4	?	?	4 🦟	observations $\{v_{ij}\}$
Alice	?	5	4 🗸	?	$(j,i) \in \Omega$
Joe	?	5	?	?	

 $d \times d'$ matrix completion \longleftrightarrow prediction from $\{1,...,d\} \times \{1,...,d'\}$ to $\mathbb R$

value of entr	y (i,j)	\longleftrightarrow	label of	input	(<i>i</i> , <i>j</i>)
---------------	---------	-----------------------	----------	-------	-----------------------	---

observed entries \iff train data

unobserved entries \iff test data

matrix \longleftrightarrow

predictor

Matrix completion: recover unknown matrix given subset of entries

	Avenuens	THEPRESTIGE	NOW YOU SEE ME	THE WOLF	
Bob	4	?	?	4 🦟	observations $\{v_{ij}\}$
Alice	?	5	4 🗸	?	$(j,i) \in \Omega$
Joe	?	5	?	?	

 $d \times d'$ matrix completion \longleftrightarrow prediction from $\{1,...,d\} \times \{1,...,d'\}$ to $\mathbb R$

value of e	entry ((i,j)	\longleftrightarrow	label	of	input	(1	i, J	j)	ĺ
------------	---------	-------	-----------------------	-------	----	-------	----	------	----	---

observed entries \iff train data

unobserved entries \iff test data

matrix \longleftrightarrow predictor

Natural complexity measure: matrix rank

Asaf Maman (TAU)

Matrix Factorization ↔ Linear Neural Network

Matrix factorization (MF):

Parameterize solution as product of matrices and fit observations via GD

Matrix factorization (MF):

Parameterize solution as product of matrices and fit observations via GD

 $MF \longleftrightarrow matrix completion via linear NN$

Matrix factorization (MF):

Parameterize solution as product of matrices and fit observations via GD

$\mathsf{MF}\longleftrightarrow\mathsf{matrix}\ \mathsf{completion}\ \mathsf{via}\ \mathsf{linear}\ \mathsf{NN}$

Past Work (e.g. Arora et al. 2019, Razin & Cohen 2020, Li et al. 2021) In MF (with small init and step size) implicit regularization minimizes rank

Matrix factorization (MF):

Parameterize solution as product of matrices and fit observations via GD

$\mathsf{MF}\longleftrightarrow\mathsf{matrix}\ \mathsf{completion}\ \mathsf{via}\ \mathsf{linear}\ \mathsf{NN}$

Past Work (e.g. Arora et al. 2019, Razin & Cohen 2020, Li et al. 2021) In MF (with small init and step size) implicit regularization minimizes rank

 $\begin{array}{rl} \mbox{Implicit regularization to low rank } + \mbox{ data is low rank} \\ \implies \mbox{ generalization} \end{array}$

As a surrogate for deep learning, MF is inherently limited:

As a surrogate for deep learning, MF is inherently limited:

(1) Misses crucial aspect of non-linearity

As a surrogate for deep learning, MF is inherently limited:

- (1) Misses crucial aspect of non-linearity
- (2) Does not capture prediction with more than 2 input variables

As a surrogate for deep learning, MF is inherently limited:

- (1) Misses crucial aspect of non-linearity
- (2) Does not capture prediction with more than 2 input variables

We study tensor factorization — accounts for both (1) and (2)

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

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

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

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

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

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

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

Tensor Factorization ↔ Non-Linear Neural Network

Tensor Factorization ↔ Non-Linear Neural Network

Tensor factorization (TF):

Parameterize solution as sum of outer products and fit observations via GD

$$\sum_{r=1}^{R} \mathbf{w}_{r}^{1} \otimes \cdots \otimes \mathbf{w}_{r}^{N}$$

Tensor factorization (TF):

Parameterize solution as sum of outer products and fit observations via GD

 $\sum_{r=1}^{R} \mathbf{w}_{r}^{1} \otimes \cdots \otimes \mathbf{w}_{r}^{N}$

 $\mathsf{TF}\longleftrightarrow\mathsf{tensor}$ completion via NN with multiplicative non-linearity

Tensor factorization (TF):

Parameterize solution as sum of outer products and fit observations via GD

 $\sum_{r=1}^{R} \mathbf{w}_{r}^{1} \otimes \cdots \otimes \mathbf{w}_{r}^{N}$

 $\mathsf{TF}\longleftrightarrow\mathsf{tensor}$ completion via NN with multiplicative non-linearity

Tensor rank: min # of components (*R*) required to express a tensor

Tensor factorization (TF):

Parameterize solution as sum of outer products and fit observations via GD

 $\sum_{r=1}^{R} \mathbf{w}_{r}^{1} \otimes \cdots \otimes \mathbf{w}_{r}^{N}$

 $\mathsf{TF}\longleftrightarrow\mathsf{tensor}$ completion via NN with multiplicative non-linearity

Tensor rank: min # of components (*R*) required to express a tensor **Empirical Phenomenon** (*Razin & Cohen 2020*)

TF (with small init and step size) accurately recovers low rank tensors

Tensor factorization (TF):

Parameterize solution as sum of outer products and fit observations via GD

 $\sum_{r=1}^{R} \mathbf{w}_{r}^{1} \otimes \cdots \otimes \mathbf{w}_{r}^{N}$

 $\mathsf{TF}\longleftrightarrow\mathsf{tensor}$ completion via NN with multiplicative non-linearity

Tensor rank: min # of components (*R*) required to express a tensor **Empirical Phenomenon** (*Razin & Cohen 2020*)

TF (with small init and step size) accurately recovers low rank tensors

We provide the first theoretical analysis for that phenomenon

Theorem

In training TF (with small init and step size): $\frac{d}{dt} \|\otimes_{n=1}^{N} \mathbf{w}_{r}^{n}\| \propto \|\otimes_{n=1}^{N} \mathbf{w}_{r}^{n}\|^{2-\frac{2}{N}}$

Theorem

In training TF (with small init and step size): $\frac{d}{dt} \|\otimes_{n=1}^{N} \mathbf{w}_{r}^{n}\| \propto \|\otimes_{n=1}^{N} \mathbf{w}_{r}^{n}\|^{2-\frac{2}{N}}$

Components move slower when small and faster when large!

Theorem

In training TF (with small init and step size): $\frac{d}{dt} \|\otimes_{n=1}^{N} \mathbf{w}_{r}^{n}\| \propto \|\otimes_{n=1}^{N} \mathbf{w}_{r}^{n}\|^{2-\frac{2}{N}}$

Components move slower when small and faster when large! Small init

Theorem

In training TF (with small init and step size): $\frac{d}{dt} \|\otimes_{n=1}^{N} \mathbf{w}_{r}^{n}\| \propto \|\otimes_{n=1}^{N} \mathbf{w}_{r}^{n}\|^{2-\frac{2}{N}}$

Components move slower when small and faster when large! Small init \implies incremental growth of components

Theorem

In training TF (with small init and step size): $\frac{d}{dt} \|\otimes_{n=1}^{N} \mathbf{w}_{r}^{n}\| \propto \|\otimes_{n=1}^{N} \mathbf{w}_{r}^{n}\|^{2-\frac{2}{N}}$

Components move slower when small and faster when large! Small init \implies incremental growth of components \implies low tensor rank

Theorem

In training TF (with small init and step size): $\frac{d}{dt} \|\otimes_{n=1}^{N} \mathbf{w}_{r}^{n}\| \propto \|\otimes_{n=1}^{N} \mathbf{w}_{r}^{n}\|^{2-\frac{2}{N}}$

Components move slower when small and faster when large!

Small init \implies incremental growth of components \implies low tensor rank

Experiment

Completion of low rank tensor via TF

Theorem

In training TF (with small init and step size): $\frac{d}{dt} \|\otimes_{n=1}^{N} \mathbf{w}_{r}^{n}\| \propto \|\otimes_{n=1}^{N} \mathbf{w}_{r}^{n}\|^{2-\frac{2}{N}}$

Components move slower when small and faster when large!

Small init \implies incremental growth of components \implies low tensor rank

Experiment

Completion of low rank tensor via TF

GD over TF leads to low tensor rank!

Theorem

In training TF (with small init and step size): $\frac{d}{dt} \|\otimes_{n=1}^{N} \mathbf{w}_{r}^{n}\| \propto \|\otimes_{n=1}^{N} \mathbf{w}_{r}^{n}\|^{2-\frac{2}{N}}$

Components move slower when small and faster when large!

Small init \implies incremental growth of components \implies low tensor rank

Experiment

Completion of low rank tensor via TF

GD over TF leads to low tensor rank!

Proposition

If tensor completion has a rank 1 solution, then under certain technical conditions, tensor factorization will reach it.

Asaf Maman (TAU)

Implicit Regularization in TF

ICMLi July 11, 2021 9 / 12

We saw:

We saw:

• Tensor completion \longleftrightarrow multi-dimensional prediction

We saw:

• Tensor completion \longleftrightarrow multi-dimensional prediction

• Tensor factorization \longleftrightarrow non-linear NN

We saw:

• Tensor completion \longleftrightarrow multi-dimensional prediction

• Tensor factorization \longleftrightarrow non-linear NN

• Implicit regularization favors tensors (predictors) of low tensor rank

We saw:

• Tensor completion \longleftrightarrow multi-dimensional prediction

• Tensor factorization \longleftrightarrow non-linear NN

• Implicit regularization favors tensors (predictors) of low tensor rank

Can tensor rank serve as a measure of complexity?

Experiment

Fitting standard datasets with predictors of low tensor rank

Experiment

Fitting standard datasets with predictors of low tensor rank

Datasets:

Experiment

Fitting standard datasets with predictors of low tensor rank

Datasets:

- MNIST 🏂 🙋 🖌 and Fashion-MNIST 🌆 👕 📗 (one-vs-all)
- Each compared against:

(i) random images (same labels) (ii) random labels (same images)

Experiment

Fitting standard datasets with predictors of low tensor rank

Datasets:

- MNIST 🏂 🙋 🍟 and Fashion-MNIST 🍠 👕 📋 (one-vs-all)
 - Each compared against:

(i) random images (same labels) (ii) random labels (same images)

Experiment

Fitting standard datasets with predictors of low tensor rank

Datasets:

- MNIST 🏂 🙋 🚽 and Fashion-MNIST 🛃 🕋 📗 (one-vs-all)
 - Each compared against:

(i) random images (same labels) (ii) random labels (same images)

Original data fit far more accurately than random (leading to low test err)!

Experiment

Fitting standard datasets with predictors of low tensor rank

Datasets:

- MNIST 🏂 🙋 🚽 and Fashion-MNIST 🛃 👕 📗 (one-vs-all)
 - Each compared against:

(i) random images (same labels) (ii) random labels (same images)

Original data fit far more accurately than random (leading to low test err)!

Tensor rank may shed light on both implicit regularization of NNs and properties of real-world data translating it to generalization

Asaf Maman (TAU)

Implicit Regularization in TF

ICMLi July 11, 2021

11 / 12

Thank You

Work supported by: Amnon and Anat Shashua, Len Blavatnik and the Blavatnik Family Foundation, Yandex Initiative in Machine Learning, Google Research Gift

Asaf Maman (TAU)

Implicit Regularization in TF

ICMLi July 11, 2021

12/12