# Implicit Regularization in Deep Learning: Lessons Learned from Matrix and Tensor Factorization

**Nadav Cohen**

Tel Aviv University

*UCLA Institute for Pure & Applied Mathematics*

*Workshop on Tensor Methods and their Applications in the Physical and Data Sciences*

31 March 2021

## Sources

**Implicit Regularization in Deep Matrix Factorization**
Arora + **C** + Hu + Luo (alphabetical order)
*NeurIPS 2019*

**Implicit Regularization in Deep Learning May Not Be Explainable by Norms**
Razin + **C**
*NeurIPS 2020*

**Implicit Regularization in Tensor Factorization**
Razin* + Maman* + **C**
*Preprint*

---

*Equal contribution

**Sanjeev Arora**   **Wei Hu**   **Yuping Luo**

**Noam Razin**   **Asaf Maman**

# Outline

# Generalization in Deep Learning

# Generalization in Deep Learning

Deep neural networks (NNs) are typically overparameterized



*# of
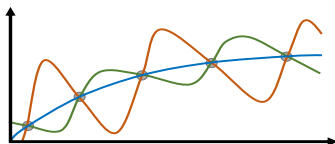learned weights*

$\gg$

*# of
training examples*

# Generalization in Deep Learning

Deep neural networks (NNs) are typically overparameterized



*# of learned weights*   $\gg$   *# of training examples*

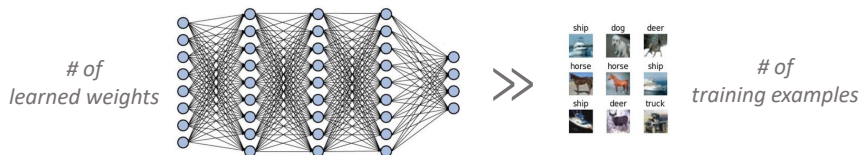$\implies$ many possible solutions (predictors) fit training data

# Generalization in Deep Learning

Deep neural networks (NNs) are typically overparameterized



*# of learned weights*   $\gg$   *# of training examples*

$\implies$   many possible solutions (predictors) fit training data
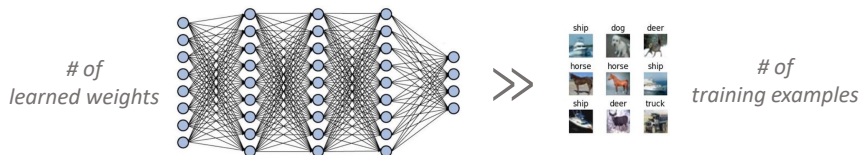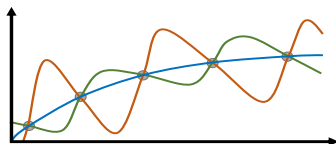


Variants of gradient descent (GD) usually find one of these solutions

# Generalization in Deep Learning

Deep neural networks (NNs) are typically overparameterized



*# of learned weights*  >>  *# of training examples*

$\implies$ many possible solutions (predictors) fit training data



Variants of gradient descent (GD) usually find one of these solutions

With "natural" data solution found often generalizes well

# Generalization in Deep Learning

Deep neural networks (NNs) are typically overparameterized

*# of learned weights*



$\gg$

*# of training examples*
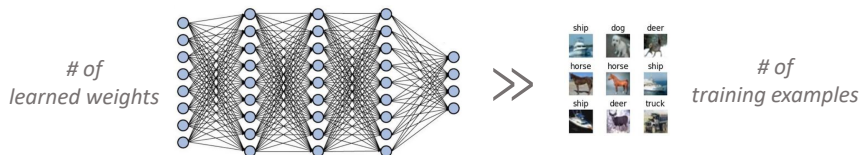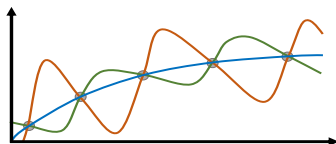
$\implies$ many possible solutions (predictors) fit training data



Variants of gradient descent (GD) usually find one of these solutions

With "natural" data solution found often generalizes well

$\uparrow$

Even without explicit regularization!

## Conventional Wisdom: Implicit Regularization
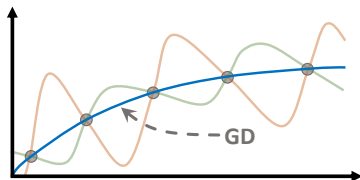
**Conventional Wisdom**

Implicit regularization minimizes "complexity":

# Conventional Wisdom: Implicit Regularization

**Conventional Wisdom**

Implicit regularization minimizes "complexity":

- GD fits training data with predictor of lowest possible complexity
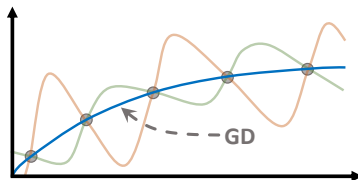
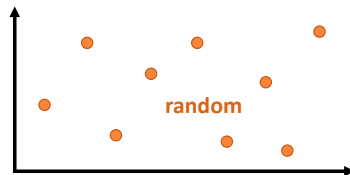# Conventional Wisdom: Implicit Regularization

**Conventional Wisdom**

Implicit regularization minimizes "complexity":

- GD fits training data with predictor of lowest possible complexity



- Natural data can be fit with low complexity, other data cannot

# Challenge: Formalizing Notion of Complexity

**Goal**

Mathematically formalize implicit regularization in deep learning (DL)

# Challenge: Formalizing Notion of Complexity

**Goal**

Mathematically formalize implicit regularization in deep learning (DL)

**Challenge**

We lack definitions for predictor complexity that are:

# Challenge: Formalizing Notion of Complexity

**Goal**

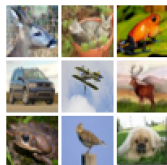Mathematically formalize implicit regularization in deep learning (DL)

**Challenge**

We lack definitions for predictor complexity that are:

- quantitative (admit generalization bounds)

  $$\text{test error} \leq \text{train error} + \mathcal{O}\Big(\text{complexity} / (\# \text{ of train examples})\Big)$$

# Challenge: Formalizing Notion of Complexity

**Goal**

Mathematically formalize implicit regularization in deep learning (DL)

**Challenge**

We lack definitions for predictor complexity that are:

- quantitative (admit generalization bounds)

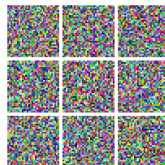    test error $\leq$ train error $+ \mathcal{O}\Big(\text{complexity} / (\#\text{ of train examples})\Big)$

- and capture essence of natural data (allow its fit with low complexity)



✔ **low complexity**        ✘ **high complexity**

# Outline

1. Implicit Regularization in Deep Learning

2. **Matrix Factorization**

3. CP Tensor Factorization

4. Tensor Rank as Measure of Complexity

5. Conclusion

# Matrix Completion $\longleftrightarrow$ Two-Dimensional Prediction

## Matrix Completion $\longleftrightarrow$ Two-Dimensional Prediction

Matrix completion: recover unknown matrix given subset of entries



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

# Matrix Completion $\longleftrightarrow$ Two-Dimensional Prediction

Matrix completion: recover unknown matrix given subset of entries



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

# Matrix Completion $\longleftrightarrow$ Two-Dimensional Prediction

Matrix completion: recover unknown matrix given subset of entries



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

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

## Matrix Completion $\longleftrightarrow$ Two-Dimensional Prediction

Matrix completion: recover unknown matrix given subset of entries



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

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

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

observed entries $\longleftrightarrow$ train data

# Matrix Completion $\longleftrightarrow$ Two-Dimensional Prediction

Matrix completion: recover unknown matrix given subset of entries



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

| value of entry $(i, j)$ | $\longleftrightarrow$ | label of input $(i, j)$ |
| --- | --- | --- |
| observed entries | $\longleftrightarrow$ | train data |
| unobserved entries | $\longleftrightarrow$ | test data |

# Matrix Completion $\longleftrightarrow$ Two-Dimensional Prediction

Matrix completion: recover unknown matrix given subset of entries



| | | | | |
|---|---|---|---|---|
| Bob | 4 | ? | ? | 4 |
| Alice | ? | 5 | 4 | ? |
| Joe | ? | 5 | ? | ? |

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

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

| | | |
|---|:---:|---|
| value of entry $(i, j)$ | $\longleftrightarrow$ | label of input $(i, j)$ |
| observed entries | $\longleftrightarrow$ | train data |
| unobserved entries | $\longleftrightarrow$ | test data |
| matrix | $\longleftrightarrow$ | predictor |

# Matrix Factorization $\longleftrightarrow$ Linear Neural Network

# Matrix Factorization ⟷ Linear Neural Network

**Matrix factorization (MF):**

Parameterize solution as product of matrices and fit observations via GD

# Matrix Factorization $\longleftrightarrow$ Linear Neural Network

**Matrix factorization (MF):**

Parameterize solution as product of matrices and fit observations via GD



$$\min_{W_1,\dots,W_N} \ \sum_{(i,j)\in\Omega} \left( [W_N W_{N-1}\cdots W_1]_{ij} - y_{ij} \right)^2$$

## Matrix Factorization $\longleftrightarrow$ Linear Neural Network

**Matrix factorization** (**MF**):

Parameterize solution as product of matrices and fit observations via GD



$$\min_{W_1,\dots,W_N} \ \sum_{(i,j)\,\in\,\Omega} \left([W_N W_{N-1} \cdots W_1]_{ij} - y_{ij}\right)^2$$

MF $\longleftrightarrow$ matrix completion via linear NN (with no explicit regularization!)

# Matrix Factorization $\longleftrightarrow$ Linear Neural Network

**Matrix factorization** (**MF**):

Parameterize solution as product of matrices and fit observations via GD



$$\min_{W_1,\ldots,W_N} \sum_{(i,j)\in\Omega} \left([W_N W_{N-1} \cdots W_1]_{ij} - y_{ij}\right)^2$$

MF $\longleftrightarrow$ matrix completion via linear NN (with no explicit regularization!)

**Empirical Phenomenon** *(Gunasekar et al. 2017)*

MF (with small init and step size) accurately recovers low rank matrices

# Implicit Regularization = Norm Minimization?

## Implicit Regularization = Norm Minimization?

**Classic Result** *(Candes & Recht 2008)*

If (i) unknown matrix has low rank; (ii) observations are sufficiently many, then fitting them while minimizing nuclear norm yields accurate recovery

# Implicit Regularization = Norm Minimization?

**Classic Result** *(Candes & Recht 2008)*
If (i) unknown matrix has low rank; (ii) observations are sufficiently many, then fitting them while minimizing nuclear norm yields accurate recovery

**Conjecture** *(Gunasekar et al. 2017)*
MF of depth 2 (with small init and step size) fits observations while minimizing nuclear norm

# Implicit Regularization = Norm Minimization?

**Classic Result** *(Candes & Recht 2008)*
If (i) unknown matrix has low rank; (ii) observations are sufficiently many, then fitting them while minimizing nuclear norm yields accurate recovery

**Conjecture** *(Gunasekar et al. 2017)*
MF of depth 2 (with small init and step size) fits observations while minimizing nuclear norm

**Experiment**



matrix completion (size 100x100, rank 5)

# Implicit Regularization = Norm Minimization?

**Classic Result** *(Candes & Recht 2008)*
If (i) unknown matrix has low rank; (ii) observations are sufficiently many,
then fitting them while minimizing nuclear norm yields accurate recovery

**Conjecture** *(Gunasekar et al. 2017)*
MF of depth 2 (with small init and step size) fits observations while
minimizing nuclear norm

**Experiment**



matrix completion (size 100x100, rank 5)

**MF gives up min nuclear norm for low rank (more so with depth)!**

# Dynamical Analysis of Implicit Regularization

# Dynamical Analysis of Implicit Regularization

Denote:

$W_e := W_N \cdots W_1$ — end matrix of MF

# Dynamical Analysis of Implicit Regularization

Denote:

$W_e := W_N \cdots W_1$ — end matrix of MF     $\{\sigma_r\}_r$ — singular vals of $W_e$

# Dynamical Analysis of Implicit Regularization

Denote:

$W_e := W_N \cdots W_1$ — end matrix of MF     $\{\sigma_r\}_r$ — singular vals of $W_e$

### Theorem

*In training MF of depth $N$ (with small init and step size): $\frac{d}{dt}\sigma_r \propto \sigma_r^{2-2/N}$*

# Dynamical Analysis of Implicit Regularization

Denote:

$W_e := W_N \cdots W_1$ — end matrix of MF    $\{\sigma_r\}_r$ — singular vals of $W_e$

### Theorem

*In training MF of depth N (with small init and step size):* $\frac{d}{dt}\sigma_r \propto \sigma_r^{2-2/N}$

Depth speeds up (slows down) large (small) singular vals!

# Dynamical Analysis of Implicit Regularization

Denote:

$W_e := W_N \cdots W_1$ — end matrix of MF $\quad \{\sigma_r\}_r$ — singular vals of $W_e$

### Theorem

*In training MF of depth $N$ (with small init and step size): $\frac{d}{dt}\sigma_r \propto \sigma_r^{2-2/N}$*

Depth speeds up (slows down) large (small) singular vals!

## **Experiment**

Completion of low rank matrix via MF

# Dynamical Analysis of Implicit Regularization

Denote:

$W_e := W_N \cdots W_1$ — end matrix of MF      $\{\sigma_r\}_r$ — singular vals of $W_e$

### Theorem

*In training MF of depth $N$ (with small init and step size):* $\frac{d}{dt}\sigma_r \propto \sigma_r^{2-2/N}$

Depth speeds up (slows down) large (small) singular vals!

## **Experiment**

Completion of low rank matrix via MF



**depth 1** (reconst error: 8e-01)   **depth 2** (reconst error: 6e-02)   **depth 3** (reconst error: 3e-05)

**MF depth leads to larger gaps between singular vals (lower rank)!**

# Dynamical Analysis of Implicit Regularization

Denote:

$W_e := W_N \cdots W_1$ — end matrix of MF      $\{\sigma_r\}_r$ — singular vals of $W_e$

### Theorem

*In training MF of depth N (with small init and step size):* $\frac{d}{dt}\sigma_r \propto \sigma_r^{2-2/N}$

Depth speeds up (slows down) large (small) singular vals!

## **Experiment**

Completion of low rank matrix via MF



**MF depth leads to larger gaps between singular vals (lower rank)!**

Further theoretical support provided in Li et al. 2021

# Dynamical Analysis of Implicit Regularization (2)

**Practical Application**

---

## Implicit Rank-Minimizing Autoencoder

---

**Li Jing**
Facebook AI Research
New York

**Jure Zbontar**
Facebook AI Research
New York

**Yann LeCun**
Facebook AI Research
New York

34th Conference on Neural Information Processing Systems (NeurIPS 2020), Vancouver, Canada.

> *"rank ... is implicitly minimized by relying on the fact that gradient descent ... in multi-layer linear networks leads to minimum-rank ..."*

# Implicit Regularization $\neq$ Norm Minimization

# Implicit Regularization $\neq$ Norm Minimization

## Theorem

*In training MF of depth N (with small init and step size):* $\frac{d}{dt}\sigma_r \propto \sigma_r^{2-2/N}$

# Implicit Regularization $\neq$ Norm Minimization

## Corollary

*In training MF of depth $N \geq 2$, $\det(W_e)$ does not change sign*

# Implicit Regularization $\neq$ Norm Minimization

### Corollary

*In training MF of depth $N \geq 2$, $\det(W_e)$ does not change sign*

Consider the matrix completion problem:

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

# Implicit Regularization $\neq$ Norm Minimization

### Corollary

*In training MF of depth $N \geq 2$, $\det(W_e)$ does not change sign*

Consider the matrix completion problem:

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

| quantity | minimizer |
|----------|-----------|
|          |           |

# Implicit Regularization $\neq$ Norm Minimization

**Corollary**

*In training MF of depth $N \geq 2$, $\det(W_e)$ does not change sign*

Consider the matrix completion problem:

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

| quantity | minimizer |
|---|---|
| *Schatten-p norm* | |

# Implicit Regularization $\neq$ Norm Minimization

### Corollary

*In training MF of depth $N \geq 2$, $\det(W_e)$ does not change sign*

Consider the matrix completion problem:

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

Special cases:
- nuclear norm
- Frobenius norm
- spectral norm

| quantity | minimizer |
|---|---|
| Schatten-p norm | |

# Implicit Regularization $\neq$ Norm Minimization

### Corollary

*In training MF of depth $N \geq 2$, $\det(W_e)$ does not change sign*

Consider the matrix completion problem:

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

Special cases:
- nuclear norm
- Frobenius norm
- spectral norm

| quantity | minimizer |
|---|---|
| Schatten-p norm | ? = 0 |

# Implicit Regularization $\neq$ Norm Minimization

**Corollary**

*In training MF of depth $N \geq 2$, $\det(W_e)$ does not change sign*

Consider the matrix completion problem:

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

Special cases:
- nuclear norm
- Frobenius norm
- spectral norm

| quantity | minimizer |
|---|---|
| Schatten-p norm | ? = 0 |
| arbitrary norm | |

# Implicit Regularization $\neq$ Norm Minimization

**Corollary**

*In training MF of depth $N \geq 2$, $\det(W_e)$ does not change sign*

Consider the matrix completion problem:

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

Special cases:
- nuclear norm
- Frobenius norm
- spectral norm

| quantity | minimizer |
|---|---|
| Schatten-p norm | ? = 0 |
| arbitrary norm | \|?\| < ∞ |

# Implicit Regularization $\neq$ Norm Minimization

**Corollary**

*In training MF of depth $N \geq 2$, $\det(W_e)$ does not change sign*

Consider the matrix completion problem:

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

Special cases:
- nuclear norm
- Frobenius norm
- spectral norm

| *quantity* | *minimizer* |
|---|---|
| *Schatten-p norm* | ? = 0 |
| *arbitrary norm* | \|?\| < ∞ |
| *rank* | |

# Implicit Regularization $\neq$ Norm Minimization

### Corollary

*In training MF of depth $N \geq 2$, $\det(W_e)$ does not change sign*

Consider the matrix completion problem:

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

Special cases:
- nuclear norm
- Frobenius norm
- spectral norm

| quantity | minimizer |
|---|---|
| *Schatten-p norm* | ? = 0 |
| *arbitrary norm* | $\|?\| < \infty$ |
| *rank* | $\|?\| \to \infty$ |

# Implicit Regularization $\neq$ Norm Minimization

### Corollary

*In training MF of depth $N \geq 2$, $\det(W_e)$ does not change sign*

Consider the matrix completion problem:

$$\begin{pmatrix} \textcolor{red}{?} & 1 \\ 1 & 0 \end{pmatrix}$$

Special cases:
- nuclear norm
- Frobenius norm
- spectral norm

| quantity | minimizer |
|---|---|
| *Schatten-p norm* | $\textcolor{red}{?} = 0$ |
| *arbitrary norm* | $|\textcolor{red}{?}| < \infty$ |
| *rank* | $|\textcolor{red}{?}| \to \infty$ |

Norm minimization contradicts rank minimization!

# Implicit Regularization $\neq$ Norm Minimization

### Corollary

*In training MF of depth $N \geq 2$, $\det(W_e)$ does not change sign*

Consider the matrix completion problem:

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

Special cases:
- nuclear norm
- Frobenius norm
- spectral norm

| quantity | minimizer |
|---|---|
| *Schatten-p norm* | ? = 0 |
| *arbitrary norm* | \|?\| < ∞ |
| *rank* | \|?\| → ∞ |

Norm minimization contradicts rank minimization!

By corollary, if $\det(W_e) > 0$ at init: fitting observations $\implies |?| \rightarrow \infty$

# Implicit Regularization $\neq$ Norm Minimization

**Corollary**

*In training MF of depth $N \geq 2$, $\det(W_e)$ does not change sign*

Consider the matrix completion problem:

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

Special cases:
- nuclear norm
- Frobenius norm
- spectral norm

| quantity | minimizer |
|---|---|
| *Schatten-p norm* | ? = 0 |
| *arbitrary norm* | $\lvert ? \rvert < \infty$ |
| *rank* | $\lvert ? \rvert \to \infty$ |

Norm minimization contradicts rank minimization!

By corollary, if $\det(W_e) > 0$ at init: fitting observations $\implies \lvert ? \rvert \to \infty$

## **Experiment**

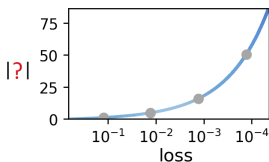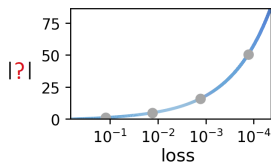# Implicit Regularization $\neq$ Norm Minimization

> **Corollary**
>
> *In training MF of depth $N \geq 2$, $\det(W_e)$ does not change sign*

Consider the matrix completion problem:

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

Special cases:
- nuclear norm
- Frobenius norm
- spectral norm

| quantity | minimizer |
|---|---|
| *Schatten-p norm* | ? = 0 |
| *arbitrary norm* | $\lvert ? \rvert < \infty$ |
| *rank* | $\lvert ? \rvert \to \infty$ |

Norm minimization contradicts rank minimization!

By corollary, if $\det(W_e) > 0$ at init: fitting observations $\implies \lvert ? \rvert \to \infty$

## **Experiment**



**There are settings where implicit regularization of MF drives all norms to $\infty$ while minimizing rank!**

# Outline

1. Implicit Regularization in Deep Learning

2. Matrix Factorization

3. **CP Tensor Factorization**

4. Tensor Rank as Measure of Complexity
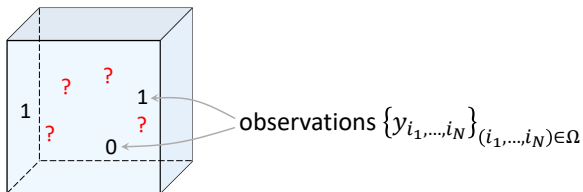
5. Conclusion

# Tensor Completion ⟷ Multi-Dimensional Prediction

# Tensor Completion ⟷ Multi-Dimensional Prediction

Tensor: multi-dim array

# Tensor Completion $\longleftrightarrow$ Multi-Dimensional Prediction
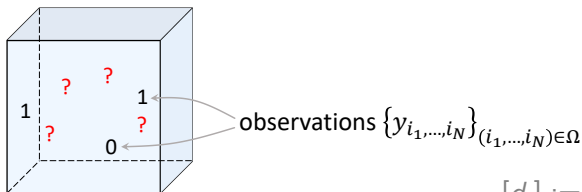
Tensor: multi-dim array

Tensor completion: recover unknown tensor given subset of entries



observations $\left\{ y_{i_1,\dots,i_N} \right\}_{(i_1,\dots,i_N) \in \Omega}$

# Tensor Completion $\longleftrightarrow$ Multi-Dimensional Prediction

Tensor: multi-dim array

Tensor completion: recover unknown tensor given subset of entries



observations $\left\{y_{i_1,\ldots,i_N}\right\}_{(i_1,\ldots,i_N)\in\Omega}$

$[d_j] := \{1,\ldots,d_j\}$
$\downarrow$

$d_1 \times \cdots \times d_N$ tensor completion $\longleftrightarrow$ prediction from $[d_1] \times \cdots \times [d_N]$ to $\mathbb{R}$

# Tensor Completion $\longleftrightarrow$ Multi-Dimensional Prediction

Tensor: multi-dim array

Tensor completion: recover unknown tensor given subset of entries



observations $\left\{ y_{i_1,\dots,i_N} \right\}_{(i_1,\dots,i_N)\in\Omega}$
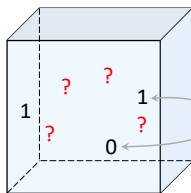
$$[d_j] := \{1,\dots,d_j\}$$
$\downarrow$

$d_1 \times \cdots \times d_N$ tensor completion $\longleftrightarrow$ prediction from $[d_1] \times \cdots \times [d_N]$ to $\mathbb{R}$

value of entry $(i_1,\dots,i_N)$ $\longleftrightarrow$ label of input $(i_1,\dots,i_N)$

## Tensor Completion $\longleftrightarrow$ Multi-Dimensional Prediction

Tensor: multi-dim array

Tensor completion: recover unknown tensor given subset of entries



observations $\left\{ y_{i_1,\ldots,i_N} \right\}_{(i_1,\ldots,i_N)\in\Omega}$
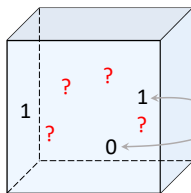
$[d_j] := \{1,\ldots,d_j\}$

$d_1 \times \cdots \times d_N$ tensor completion $\longleftrightarrow$ prediction from $[d_1] \times \cdots \times [d_N]$ to $\mathbb{R}$
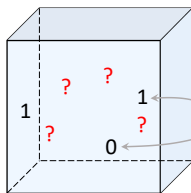
value of entry $(i_1,\ldots,i_N)$ $\longleftrightarrow$ label of input $(i_1,\ldots,i_N)$

observed entries $\longleftrightarrow$ train data

# Tensor Completion $\longleftrightarrow$ Multi-Dimensional Prediction

Tensor: multi-dim array

Tensor completion: recover unknown tensor given subset of entries



observations $\left\{ y_{i_1,\dots,i_N} \right\}_{(i_1,\dots,i_N)\in\Omega}$

$$[d_j] := \{1,\dots,d_j\}$$
$$\downarrow$$

$d_1 \times \cdots \times d_N$ tensor completion $\longleftrightarrow$ prediction from $[d_1] \times \cdots \times [d_N]$ to $\mathbb{R}$

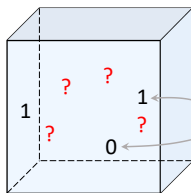| value of entry $(i_1,\dots,i_N)$ | $\longleftrightarrow$ | label of input $(i_1,\dots,i_N)$ |
|:---:|:---:|:---:|
| observed entries | $\longleftrightarrow$ | train data |
| unobserved entries | $\longleftrightarrow$ | test data |

# Tensor Completion $\longleftrightarrow$ Multi-Dimensional Prediction

Tensor: multi-dim array

Tensor completion: recover unknown tensor given subset of entries



observations $\left\{ y_{i_1,\ldots,i_N} \right\}_{(i_1,\ldots,i_N)\in\Omega}$

$[d_j] := \{1,\ldots,d_j\}$
$\downarrow$

$d_1 \times \cdots \times d_N$ tensor completion $\longleftrightarrow$ prediction from $[d_1] \times \cdots \times [d_N]$ to $\mathbb{R}$

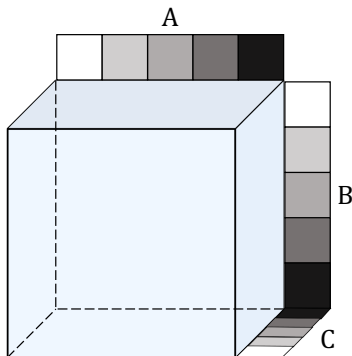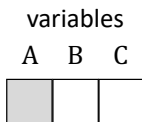| value of entry $(i_1,\ldots,i_N)$ | $\longleftrightarrow$ | label of input $(i_1,\ldots,i_N)$ |
|---|---|---|
| observed entries | $\longleftrightarrow$ | train data |
| unobserved entries | $\longleftrightarrow$ | test data |
| tensor | $\longleftrightarrow$ | predictor |

# Tensor Completion $\longleftrightarrow$ Multi-Dimensional Prediction (2)

Standard prediction tasks can be seen as tensor completion problems

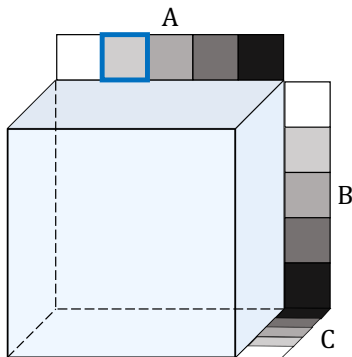# Tensor Completion ⟷ Multi-Dimensional Prediction (2)

Standard prediction tasks can be seen as tensor completion problems

## **Illustration — Image Recognition**

# Tensor Completion ⟷ Multi-Dimensional Prediction (2)

Standard prediction tasks can be seen as tensor completion problems

## **Illustration — Image Recognition**

# Tensor Completion ⟷ Multi-Dimensional Prediction (2)

Standard prediction tasks can be seen as tensor completion problems

## Illustration — Image Recognition

# Tensor Completion ⟷ Multi-Dimensional Prediction (2)

Standard prediction tasks can be seen as tensor completion problems

## **Illustration — Image Recognition**

# Tensor Completion ⟷ Multi-Dimensional Prediction (2)

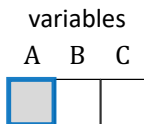Standard prediction tasks can be seen as tensor completion problems
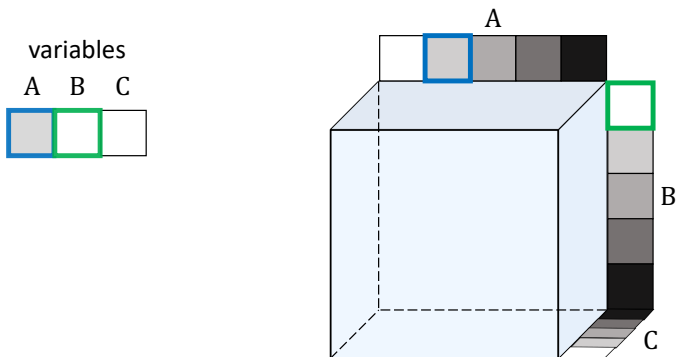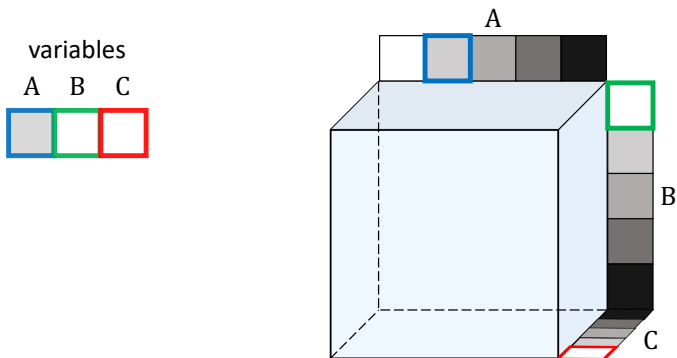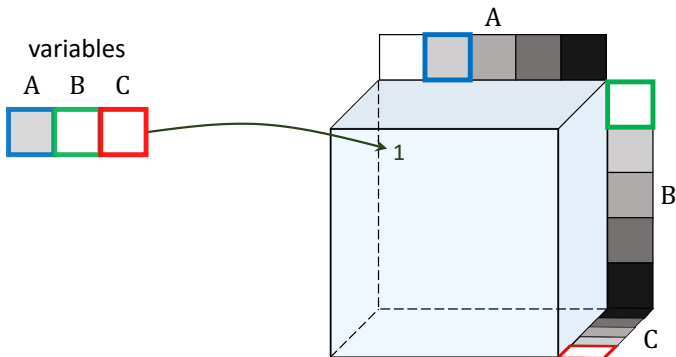
## Illustration — Image Recognition

# Tensor Completion $\longleftrightarrow$ Multi-Dimensional Prediction (2)

Standard prediction tasks can be seen as tensor completion problems

## **Illustration — Image Recognition**

# Tensor Completion ⟷ Multi-Dimensional Prediction (2)

Standard prediction tasks can be seen as tensor completion problems

## Illustration — Image Recognition

# Tensor Completion ⟷ Multi-Dimensional Prediction (2)

Standard prediction tasks can be seen as tensor completion problems

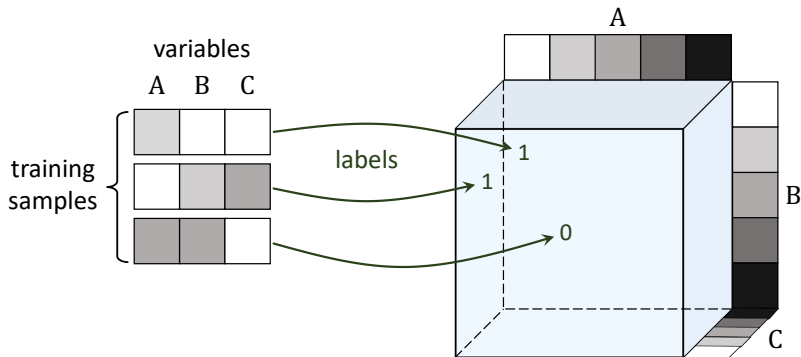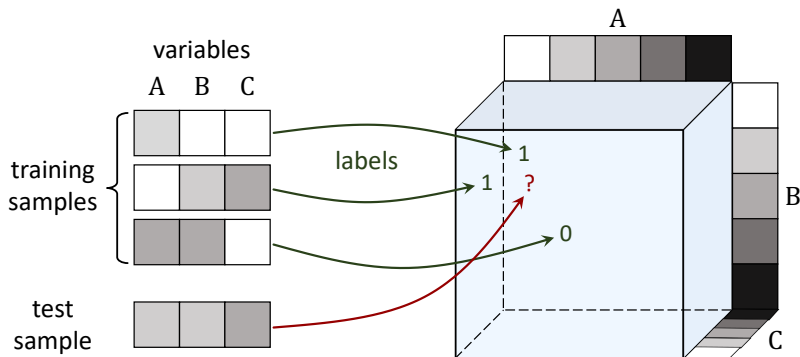## Illustration — Image Recognition

# Tensor Completion ⟷ Multi-Dimensional Prediction (2)

Standard prediction tasks can be seen as tensor completion problems

## <u>Illustration — Image Recognition</u>

# CP Tensor Factorization $\longleftrightarrow$ Non-Linear Neural Network

# CP Tensor Factorization $\longleftrightarrow$ Non-Linear Neural Network

**CP tensor factorization** (**TF**):

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

# CP Tensor Factorization $\longleftrightarrow$ Non-Linear Neural Network

**CP tensor factorization** (**TF**):

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

$$\min_{\{\mathbf{w}_r^n\}_{r,n}} \ell(\{\mathbf{w}_r^n\}_{r,n}) := \sum_{(i_1,\ldots,i_N) \in \Omega} \left( \left[ \sum_{r=1}^{R} \otimes_{n=1}^{N} \mathbf{w}_r^n \right]_{i_1,\ldots,i_N} - y_{i_1,\ldots,i_N} \right)^2$$

# CP Tensor Factorization $\longleftrightarrow$ Non-Linear Neural Network

**CP tensor factorization** (**TF**):

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

$$\min_{\{\mathbf{w}_r^n\}_{r,n}} \ \ell(\{\mathbf{w}_r^n\}_{r,n}) := \sum_{(i_1,\ldots,i_N) \in \Omega} \left( \left[ \sum_{r=1}^{R} \otimes_{n=1}^{N} \mathbf{w}_r^n \right]_{i_1,\ldots,i_N} - y_{i_1,\ldots,i_N} \right)^2$$

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

# CP Tensor Factorization ⟷ Non-Linear Neural Network

**CP tensor factorization** (**TF**):

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

$$\min_{\{\mathbf{w}_r^n\}_{r,n}} \ell(\{\mathbf{w}_r^n\}_{r,n}) := \sum_{(i_1,\dots,i_N)\in\Omega}\left(\left[\sum_{r=1}^{R}\otimes_{n=1}^{N}\mathbf{w}_r^n\right]_{i_1,\dots,i_N} - y_{i_1,\dots,i_N}\right)^2$$
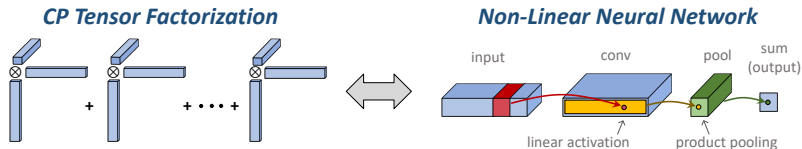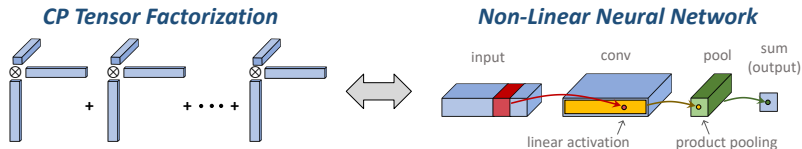
TF ⟷ tensor completion via NN with multiplicative non-linearity



*CP Tensor Factorization*    *Non-Linear Neural Network*

input    conv    pool    sum (output)

linear activation    product pooling

## Experiment

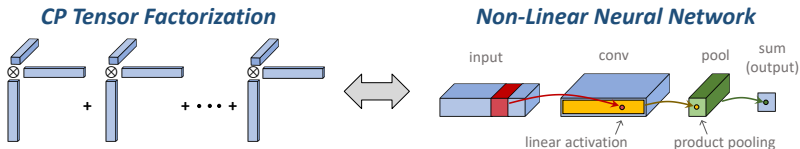TF (with small init and step size) accurately recovers low rank tensors

# CP Tensor Factorization $\longleftrightarrow$ Non-Linear Neural Network

**CP tensor factorization** (**TF**):

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

$$\min_{\{\mathbf{w}_r^n\}_{r,n}} \ell(\{\mathbf{w}_r^n\}_{r,n}) := \sum_{(i_1,\ldots,i_N)\in\Omega} \left( \left[ \sum_{r=1}^R \otimes_{n=1}^N \mathbf{w}_r^n \right]_{i_1,\ldots,i_N} - y_{i_1,\ldots,i_N} \right)^2$$

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



*CP Tensor Factorization*       *Non-Linear Neural Network*

### Experiment

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

                                        $\uparrow$

Tensor rank: min # of components in CP representation

# Dynamical Analysis of Implicit Regularization

# Dynamical Analysis of Implicit Regularization

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

# Dynamical Analysis of Implicit Regularization

### Theorem

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

Component norms accelerate (decelerate) when large (small)!

# Dynamical Analysis of Implicit Regularization

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

Component norms accelerate (decelerate) when large (small)!

**Proof Sketch**

# Dynamical Analysis of Implicit Regularization

### Theorem

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

Component norms accelerate (decelerate) when large (small)!

### **Proof Sketch**

GD with step size $\to 0$ (gradient flow): $\frac{d}{dt}\mathbf{w}_r^n(t) = -\frac{\partial}{\partial \mathbf{w}_r^n}\ell(\{\mathbf{w}_{r'}^{n'}(t)\}_{r',n'})$

# Dynamical Analysis of Implicit Regularization

### Theorem

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

Component norms accelerate (decelerate) when large (small)!

## **Proof Sketch**

GD with step size $\to 0$ (gradient flow): $\frac{d}{dt}\mathbf{w}_r^n(t) = -\frac{\partial}{\partial\mathbf{w}_r^n}\ell(\{\mathbf{w}_{r'}^{n'}(t)\}_{r',n'})$

For any $n, \bar{n}$: $\|\mathbf{w}_r^n(t)\|^2 - \|\mathbf{w}_r^{\bar{n}}(t)\|^2$ is constant through time

# Dynamical Analysis of Implicit Regularization

### Theorem

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

Component norms accelerate (decelerate) when large (small)!

**Proof Sketch**

GD with step size $\to 0$ (gradient flow): $\frac{d}{dt}\mathbf{w}_r^n(t) = -\frac{\partial}{\partial \mathbf{w}_r^n}\ell(\{\mathbf{w}_{r'}^{n'}(t)\}_{r',n'})$

For any $n, \bar{n}$: $\|\mathbf{w}_r^n(t)\|^2 - \|\mathbf{w}_r^{\bar{n}}(t)\|^2$ is constant through time

$\implies$ under small init $\|\mathbf{w}_r^n(t)\|^2 \approx \|\mathbf{w}_r^{\bar{n}}(t)\|^2 \approx \|\otimes_{n'=1}^{N}\mathbf{w}_r^{n'}(t)\|^{\frac{2}{N}}$

# Dynamical Analysis of Implicit Regularization

---

### Theorem

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

---

Component norms accelerate (decelerate) when large (small)!

**Proof Sketch**

GD with step size $\to 0$ (gradient flow): $\frac{d}{dt}\mathbf{w}_r^n(t) = -\frac{\partial}{\partial \mathbf{w}_r^n}\ell(\{\mathbf{w}_{r'}^{n'}(t)\}_{r',n'})$

For any $n, \bar{n}$: $\|\mathbf{w}_r^n(t)\|^2 - \|\mathbf{w}_r^{\bar{n}}(t)\|^2$ is constant through time

$\implies$ under small init $\|\mathbf{w}_r^n(t)\|^2 \approx \|\mathbf{w}_r^{\bar{n}}(t)\|^2 \approx \|\otimes_{n'=1}^N \mathbf{w}_r^{n'}(t)\|^{\frac{2}{N}}$

Denote:

$\mathcal{W}_e := \sum_{r=1}^R \otimes_{n=1}^N \mathbf{w}_r^n$ — end tensor

# Dynamical Analysis of Implicit Regularization

### Theorem

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

Component norms accelerate (decelerate) when large (small)!

**Proof Sketch**

GD with step size $\to 0$ (gradient flow): $\frac{d}{dt}\mathbf{w}_r^n(t) = -\frac{\partial}{\partial \mathbf{w}_r^n}\ell(\{\mathbf{w}_{r'}^{n'}(t)\}_{r',n'})$

For any $n, \bar{n}$: $\|\mathbf{w}_r^n(t)\|^2 - \|\mathbf{w}_r^{\bar{n}}(t)\|^2$ is constant through time

$\implies$ under small init $\|\mathbf{w}_r^n(t)\|^2 \approx \|\mathbf{w}_r^{\bar{n}}(t)\|^2 \approx \|\otimes_{n'=1}^{N}\mathbf{w}_r^{n'}(t)\|^{\frac{2}{N}}$

Denote:

$\mathcal{W}_e := \sum_{r=1}^{R}\otimes_{n=1}^{N}\mathbf{w}_r^n$ — end tensor , $\mathcal{L}(\cdot) :=$ loss w.r.t. $\mathcal{W}_e$

# Dynamical Analysis of Implicit Regularization

### Theorem

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

Component norms accelerate (decelerate) when large (small)!

### **Proof Sketch**

GD with step size $\to 0$ (gradient flow): $\frac{d}{dt}\mathbf{w}_r^n(t) = -\frac{\partial}{\partial\mathbf{w}_r^n}\ell(\{\mathbf{w}_{r'}^{n'}(t)\}_{r',n'})$

For any $n, \bar{n}$: $\|\mathbf{w}_r^n(t)\|^2 - \|\mathbf{w}_r^{\bar{n}}(t)\|^2$ is constant through time

$\implies$ under small init $\|\mathbf{w}_r^n(t)\|^2 \approx \|\mathbf{w}_r^{\bar{n}}(t)\|^2 \approx \|\otimes_{n'=1}^{N}\mathbf{w}_r^{n'}(t)\|^{\frac{2}{N}}$

Denote:

$\mathcal{W}_e := \sum_{r=1}^{R}\otimes_{n=1}^{N}\mathbf{w}_r^n$ — end tensor , $\mathcal{L}(\cdot) :=$ loss w.r.t. $\mathcal{W}_e$ , $\widehat{\mathbf{w}}_r^n := \frac{\mathbf{w}_r^n}{\|\mathbf{w}_r^n\|}$

# Dynamical Analysis of Implicit Regularization

### Theorem

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

Component norms accelerate (decelerate) when large (small)!

## **Proof Sketch**

GD with step size $\to 0$ (gradient flow): $\frac{d}{dt}\mathbf{w}_r^n(t) = -\frac{\partial}{\partial \mathbf{w}_r^n}\ell(\{\mathbf{w}_{r'}^{n'}(t)\}_{r',n'})$

For any $n, \bar{n}$: $\|\mathbf{w}_r^n(t)\|^2 - \|\mathbf{w}_r^{\bar{n}}(t)\|^2$ is constant through time

$\implies$ under small init $\|\mathbf{w}_r^n(t)\|^2 \approx \|\mathbf{w}_r^{\bar{n}}(t)\|^2 \approx \|\otimes_{n'=1}^{N}\mathbf{w}_r^{n'}(t)\|^{\frac{2}{N}}$

Denote:

$\mathcal{W}_e := \sum_{r=1}^{R}\otimes_{n=1}^{N}\mathbf{w}_r^n$ — end tensor , $\mathcal{L}(\cdot) :=$ loss w.r.t. $\mathcal{W}_e$ , $\widehat{\mathbf{w}}_r^n := \frac{\mathbf{w}_r^n}{\|\mathbf{w}_r^n\|}$

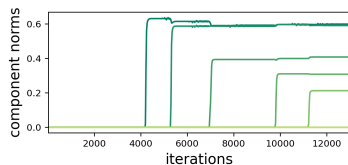Differentiate w.r.t. time:

$\frac{d}{dt}\|\otimes_{n=1}^{N}\mathbf{w}_r^n(t)\| = \sum_{n=1}^{N}\prod_{n'\neq n}\|\mathbf{w}_r^{n'}(t)\|^2 \cdot \left\langle -\nabla\mathcal{L}(\mathcal{W}_e(t)), \otimes_{n=1}^{N}\widehat{\mathbf{w}}_r^n(t)\right\rangle$
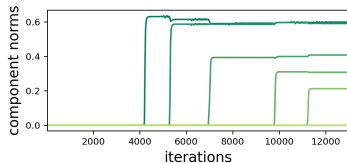
# Dynamical Analysis of Implicit Regularization

### Theorem

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

Component norms accelerate (decelerate) when large (small)!

**Proof Sketch**

GD with step size $\to 0$ (gradient flow): $\frac{d}{dt}\mathbf{w}_r^n(t) = -\frac{\partial}{\partial\mathbf{w}_r^n}\ell(\{\mathbf{w}_{r'}^{n'}(t)\}_{r',n'})$

For any $n, \bar{n}$: $\|\mathbf{w}_r^n(t)\|^2 - \|\mathbf{w}_r^{\bar{n}}(t)\|^2$ is constant through time

$\implies$ under small init $\|\mathbf{w}_r^n(t)\|^2 \approx \|\mathbf{w}_r^{\bar{n}}(t)\|^2 \approx \|\otimes_{n'=1}^{N}\mathbf{w}_r^{n'}(t)\|^{\frac{2}{N}}$

Denote:

$\mathcal{W}_e := \sum_{r=1}^{R}\otimes_{n=1}^{N}\mathbf{w}_r^n$ — end tensor , $\mathcal{L}(\cdot) := $ loss w.r.t. $\mathcal{W}_e$ , $\widehat{\mathbf{w}}_r^n := \frac{\mathbf{w}_r^n}{\|\mathbf{w}_r^n\|}$

Differentiate w.r.t. time:

$\frac{d}{dt}\|\otimes_{n=1}^{N}\mathbf{w}_r^n(t)\| = \sum_{n=1}^{N}\prod_{n'\neq n}\|\mathbf{w}_r^{n'}(t)\|^2 \cdot \left\langle -\nabla\mathcal{L}(\mathcal{W}_e(t)), \otimes_{n=1}^{N}\widehat{\mathbf{w}}_r^n(t)\right\rangle$

# Dynamical Analysis of Implicit Regularization

## Theorem

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

Component norms accelerate (decelerate) when large (small)!

**Proof Sketch**

GD with step size $\to 0$ (gradient flow): $\frac{d}{dt}\mathbf{w}_r^n(t) = -\frac{\partial}{\partial\mathbf{w}_r^n}\ell(\{\mathbf{w}_{r'}^{n'}(t)\}_{r',n'})$

For any $n, \bar{n}$: $\|\mathbf{w}_r^n(t)\|^2 - \|\mathbf{w}_r^{\bar{n}}(t)\|^2$ is constant through time

$\implies$ under small init $\|\mathbf{w}_r^n(t)\|^2 \approx \|\mathbf{w}_r^{\bar{n}}(t)\|^2 \approx \|\otimes_{n'=1}^{N}\mathbf{w}_r^{n'}(t)\|^{\frac{2}{N}}$

Denote:

$\mathcal{W}_e := \sum_{r=1}^{R}\otimes_{n=1}^{N}\mathbf{w}_r^n$ — end tensor , $\mathcal{L}(\cdot) :=$ loss w.r.t. $\mathcal{W}_e$ , $\widehat{\mathbf{w}}_r^n := \frac{\mathbf{w}_r^n}{\|\mathbf{w}_r^n\|}$

Differentiate w.r.t. time:

$\frac{d}{dt}\|\otimes_{n=1}^{N}\mathbf{w}_r^n(t)\| \approx \|\otimes_{n=1}^{N}\mathbf{w}_r^n(t)\|^{2-\frac{2}{N}} \cdot N\langle -\nabla\mathcal{L}(\mathcal{W}_e(t)), \otimes_{n=1}^{N}\widehat{\mathbf{w}}_r^n(t)\rangle$

# Dynamical Analysis of Implicit Regularization (2)

**Experiment**

Completion of low rank tensor via TF

# Dynamical Analysis of Implicit Regularization (2)

## Experiment

Completion of low rank tensor via TF



> **Training TF leads to gaps between component norms (low tensor rank)!**

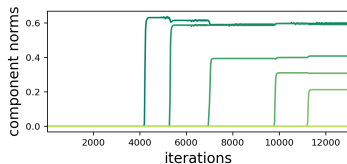# Dynamical Analysis of Implicit Regularization (2)

## **Experiment**

Completion of low rank tensor via TF



> **Training TF leads to gaps between component norms (low tensor rank)!**

### Proposition

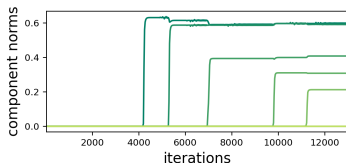*If tensor completion has rank 1 solution, then under technical conditions TF will reach it*

# Dynamical Analysis of Implicit Regularization (2)

**Experiment**

Completion of low rank tensor via TF



> **Training TF leads to gaps between component norms (low tensor rank)!**

## Proposition

*If tensor completion has rank 1 solution, then under technical conditions TF will reach it*

**Proof Sketch**

# Dynamical Analysis of Implicit Regularization (2)

## Experiment

Completion of low rank tensor via TF



> **Training TF leads to gaps between component norms (low tensor rank)!**

### Proposition

*If tensor completion has rank 1 solution, then under technical conditions TF will reach it*
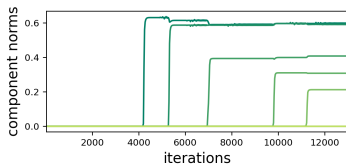
## Proof Sketch

Denote: $\alpha > 0$ — init scale

# Dynamical Analysis of Implicit Regularization (2)

## Experiment

Completion of low rank tensor via TF



**Training TF leads to gaps between component norms (low tensor rank)!**

## Proposition

*If tensor completion has rank 1 solution, then under technical conditions TF will reach it*

## Proof Sketch

Denote: $\alpha > 0$ — init scale

$$\frac{d}{dt}\|\otimes_{n=1}^{N}\mathbf{w}_r^n\| \propto \|\otimes_{n=1}^{N}\mathbf{w}_r^n\|^{2-\frac{2}{N}}$$

# Dynamical Analysis of Implicit Regularization (2)

**Experiment**

Completion of low rank tensor via TF



**Training TF leads to gaps between component norms (low tensor rank)!**

### Proposition

*If tensor completion has rank 1 solution, then under technical conditions TF will reach it*

**Proof Sketch**

Denote: $\alpha > 0$ — init scale

$$\frac{d}{dt}\|\otimes_{n=1}^{N}\mathbf{w}_r^n\| \propto \|\otimes_{n=1}^{N}\mathbf{w}_r^n\|^{2-\frac{2}{N}} \implies \text{one component } \mathcal{O}(1) \text{ while others } \mathcal{O}(\alpha^N)$$

# Dynamical Analysis of Implicit Regularization (2)

## Experiment

Completion of low rank tensor via TF



**Training TF leads to gaps between component norms (low tensor rank)!**

## Proposition

*If tensor completion has rank 1 solution, then under technical conditions TF will reach it*

## Proof Sketch

Denote: $\alpha > 0$ — init scale

$$\frac{d}{dt}\|\otimes_{n=1}^{N}\mathbf{w}_r^n\| \propto \|\otimes_{n=1}^{N}\mathbf{w}_r^n\|^{2-\frac{2}{N}} \implies \text{one component } \mathcal{O}(1) \text{ while others } \mathcal{O}(\alpha^N)$$

$\alpha \to 0$

# Dynamical Analysis of Implicit Regularization (2)

## **Experiment**

Completion of low rank tensor via TF



**Training TF leads to gaps between component norms (low tensor rank)!**

## Proposition

*If tensor completion has rank 1 solution, then under technical conditions TF will reach it*

## **Proof Sketch**

Denote: $\alpha > 0$ — init scale

$\frac{d}{dt}\|\otimes_{n=1}^{N}\mathbf{w}_r^n\| \propto \|\otimes_{n=1}^{N}\mathbf{w}_r^n\|^{2-\frac{2}{N}} \implies$ one component $\mathcal{O}(1)$ while others $\mathcal{O}(\alpha^N)$

$\alpha \to 0 \implies$ end tensor $\mathcal{W}_e$ follows rank 1 trajectory until convergence

# Outline

# Challenge: Formalizing Notion of Complexity

**<u>Goal</u>**

Mathematically formalize implicit regularization in deep learning (DL)

**<u>Challenge</u>**

We lack definitions for predictor complexity that are:

- quantitative (admit generalization bounds)

    $$\text{test error} \leq \text{train error} + \mathcal{O}\Big(\text{complexity} / (\# \text{ of train examples})\Big)$$

- and capture essence of natural data (allow its fit with low complexity)



✔️ **low complexity**       ❌ **high complexity**

# Tensor Rank Captures Non-Linear Neural Network

We saw:

# Tensor Rank Captures Non-Linear Neural Network

We saw:

- Tensor completion $\longleftrightarrow$ multi-dim prediction

# Tensor Rank Captures Non-Linear Neural Network

We saw:

- Tensor completion $\longleftrightarrow$ multi-dim prediction



- CP tensor factorization $\longleftrightarrow$ non-linear NN

# Tensor Rank Captures Non-Linear Neural Network

We saw:

- Tensor completion $\longleftrightarrow$ multi-dim prediction



- CP tensor factorization $\longleftrightarrow$ non-linear NN



- Implicit regularization favors tensors (predictors) of low rank

# Tensor Rank Captures Non-Linear Neural Network

We saw:

- Tensor completion $\longleftrightarrow$ multi-dim prediction



- CP tensor factorization $\longleftrightarrow$ non-linear NN



- Implicit regularization favors tensors (predictors) of low rank

**<u>Question</u>**

Can tensor rank serve as measure of complexity for predictors?

# Experiment: Fitting Data with Low Tensor Rank

# Experiment: Fitting Data with Low Tensor Rank

**Experiment**

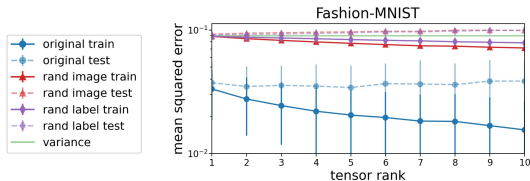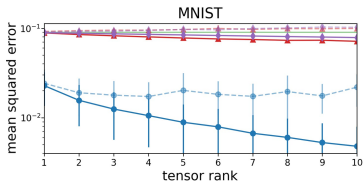Fitting data with predictors of low tensor rank

# Experiment: Fitting Data with Low Tensor Rank

**Experiment**

Fitting data with predictors of low tensor rank

Datasets:

- MNIST  and Fashion-MNIST  (one-vs-all)

# Experiment: Fitting Data with Low Tensor Rank

**Experiment**

Fitting data 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 Data with Low Tensor Rank

**Experiment**

Fitting data 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 Data with Low Tensor Rank

**Experiment**

Fitting data 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**

# Outline

# Recap

# Recap

Understanding implicit regularization in DL:

# Recap

Understanding implicit regularization in DL:

- <u>Challenge</u>: lack measures of complexity that capture natural data

# Recap

Understanding implicit regularization in DL:

- <u>Challenge</u>: lack measures of complexity that capture natural data

**Matrix factorization**:

## Recap

Understanding implicit regularization in DL:

- <u>Challenge</u>: lack measures of complexity that capture natural data

**Matrix factorization**:

- Equivalent to two-dim prediction via linear NN

## Recap

Understanding implicit regularization in DL:

- <u>Challenge</u>: lack measures of complexity that capture natural data

**Matrix factorization**:

- Equivalent to two-dim prediction via linear NN

- <u>Conjecture</u>: implicit regularization minimizes norm

## Recap

Understanding implicit regularization in DL:

- <u>Challenge</u>: lack measures of complexity that capture natural data

**Matrix factorization**:

- Equivalent to two-dim prediction via linear NN

- <u>Conjecture</u>: implicit regularization minimizes norm

- <u>Dynamical analysis</u>: implicit regularization minimizes rank (not norm)

# Recap

Understanding implicit regularization in DL:

- <u>Challenge</u>: lack measures of complexity that capture natural data

**Matrix factorization**:

- Equivalent to two-dim prediction via linear NN

- <u>Conjecture</u>: implicit regularization minimizes norm

- <u>Dynamical analysis</u>: implicit regularization minimizes rank (not norm)
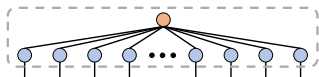
**CP tensor factorization**:

## Recap

Understanding implicit regularization in DL:

- <u>Challenge</u>: lack measures of complexity that capture natural data

**Matrix factorization**:

- Equivalent to two-dim prediction via linear NN

- <u>Conjecture</u>: implicit regularization minimizes norm

- <u>Dynamical analysis</u>: implicit regularization minimizes rank (not norm)

**CP tensor factorization**:

- Equivalent to multi-dim prediction via non-linear NN

## Recap

Understanding implicit regularization in DL:

- <u>Challenge</u>: lack measures of complexity that capture natural data

**Matrix factorization**:

- Equivalent to two-dim prediction via linear NN

- <u>Conjecture</u>: implicit regularization minimizes norm

- <u>Dynamical analysis</u>: implicit regularization minimizes rank (not norm)

**CP tensor factorization**:

- Equivalent to multi-dim prediction via non-linear NN

- <u>Dynamical analysis</u>: implicit regularization minimizes tensor rank

## Recap

Understanding implicit regularization in DL:

- <u>Challenge</u>: lack measures of complexity that capture natural data

**Matrix factorization**:

- Equivalent to two-dim prediction via linear NN

- <u>Conjecture</u>: implicit regularization minimizes norm

- <u>Dynamical analysis</u>: implicit regularization minimizes rank (not norm)

**CP tensor factorization**:

- Equivalent to multi-dim prediction via non-linear NN

- <u>Dynamical analysis</u>: implicit regularization minimizes tensor rank

Tensor rank as measure of complexity may capture natural data!

# Ongoing Work: Adding Depth via Hierarchy

# Ongoing Work: Adding Depth via Hierarchy

**CP Tensor Factorization**



**Shallow Non-Linear Neural Network**



input    conv    pool   output

linear activation    product pooling

# Ongoing Work: Adding Depth via Hierarchy

**CP Tensor Factorization**

**Shallow Non-Linear Neural Network**



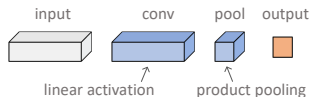**Implicit regularization = minimization of tensor rank**

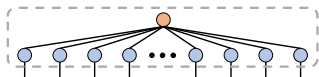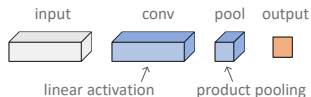# Ongoing Work: Adding Depth via Hierarchy

**CP Tensor Factorization**

**Shallow Non-Linear Neural Network**



**Implicit regularization = minimization of tensor rank**

❌ **Oblivious to input ordering**

# Ongoing Work: Adding Depth via Hierarchy
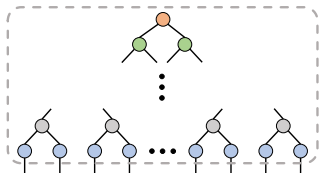
**CP Tensor Factorization**

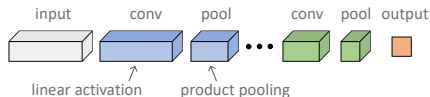**Shallow Non-Linear Neural Network**



*Implicit regularization = minimization of tensor rank*

❌ **Oblivious to input ordering**
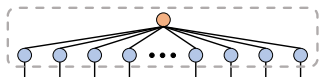
**Hierarchical Tensor Factorization**

**Deep Non-Linear Neural Network**

# Ongoing Work: Adding Depth via Hierarchy

**CP Tensor Factorization**

**Shallow Non-Linear Neural Network**



*Implicit regularization = minimization of tensor rank*

❌ **Oblivious to input ordering**

**Hierarchical Tensor Factorization**

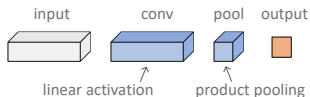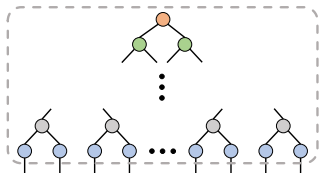**Deep Non-Linear Neural Network**



**?**

*Implicit regularization = minimization of hierarchical tensor rank*

# Ongoing Work: Adding Depth via Hierarchy

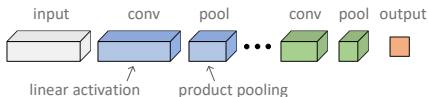**CP Tensor Factorization**                    **Shallow Non-Linear Neural Network**



*Implicit regularization = minimization of <u>tensor rank</u>*

❌ **Oblivious to input ordering**

**Hierarchical Tensor Factorization**          **Deep Non-Linear Neural Network**



*Implicit regularization = minimization of <u>hierarchical tensor rank</u>*

✔ **Accounts for input ordering**

# Outline

# Thank You