

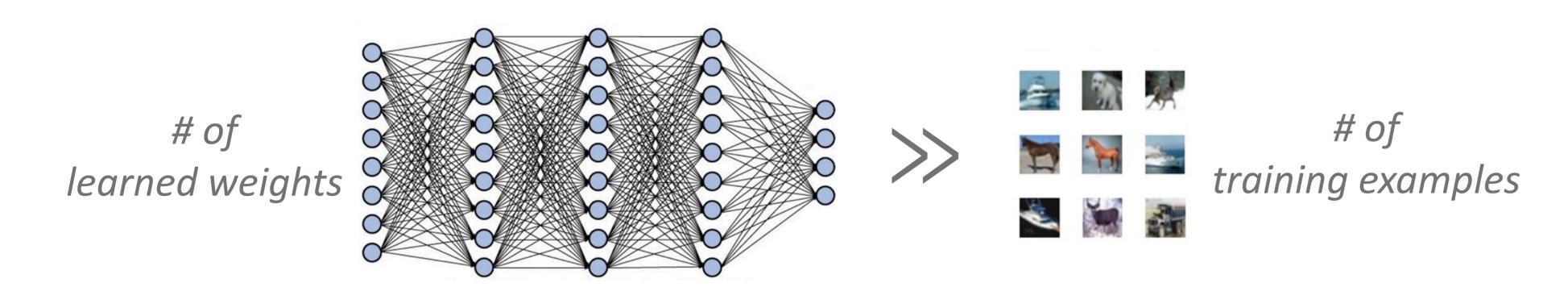
Implicit Regularization in Tensor Factorization

International Conference on Machine Learning (ICML) 2021 Asaf Maman* Nadav Cohen *Equal contribution



I) Implicit Regularization in Deep Learning (DL)

Deep neural networks (NNs) are typically overparameterized



With "natural data" predictors found by gradient descent (GD) generalize well

Conventional wisdom: implicit regularization towards low "complexity" predictors

Goal: mathematically understand this implicit regularization

Challenge: lack complexity measures that capture essence of natural data



Can we characterize the implicit regularization in concrete settings?

II) Common Testbed: Matrix Factorization (MF)

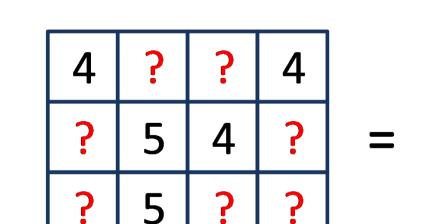
Matrix Completion: recover unknown matrix given subset of entries

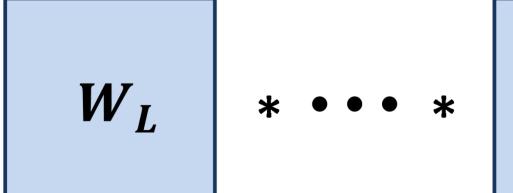
	Avenuens	THEPRESTIGE	NOW YOU SEE ME	THE WOLF OF WALL STREET
Bob	4	?	?	4
Alice	?	5	4	?
Joe	?	5	?	?

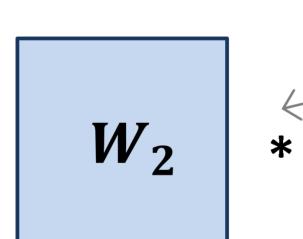
 d^{θ} matrix completion / prediction task from f1; ...; dg = f1; ...; $d^{\theta}g$ to R In many real-world scenarios matrices of interest have low rank

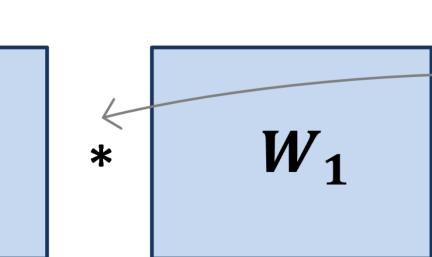
Matrix Factorization

Parameterize solution as product of matrices and fit observations with GD









hidden dims do not necessarily constrain rank

solving matrix completion via linear NN (w/o explicit regularization!)

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

Limitations of Matrix Factorization

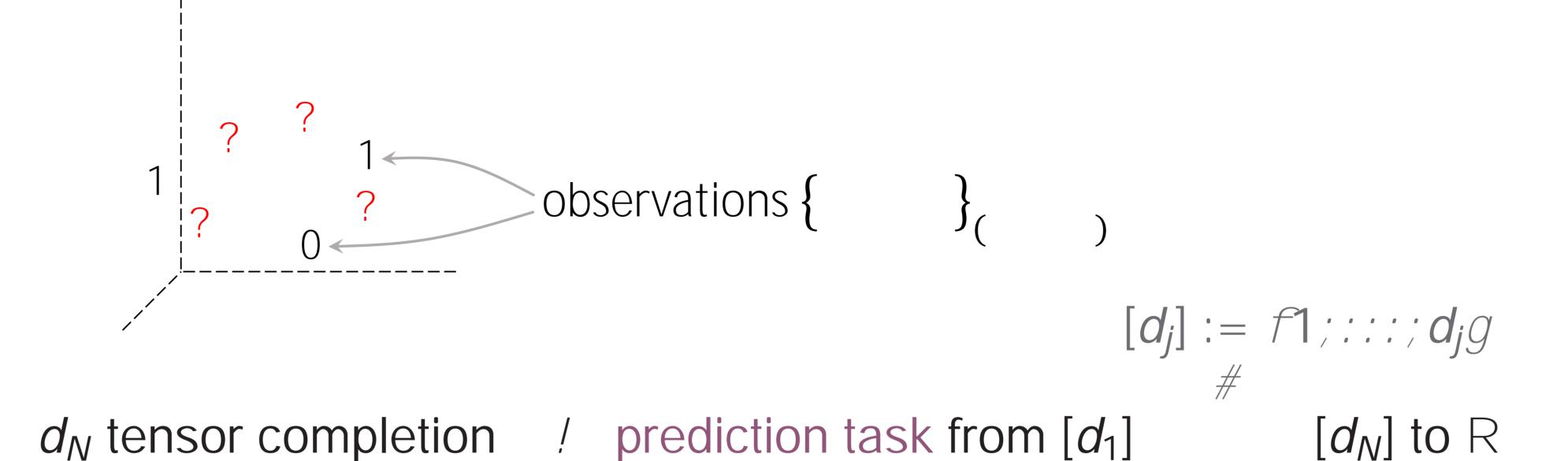
(1) Captures prediction over only 2 input variables (2) No non-linearity

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

III) Beyond Matrix Factorization: Tensor Factorization (TF)

Tensor Completion: recover unknown tensor given subset of entries

Tel Aviv University



Tensor Factorization

(1) d_1

Parameterize solution as sum of outer products and fit observations with GD:

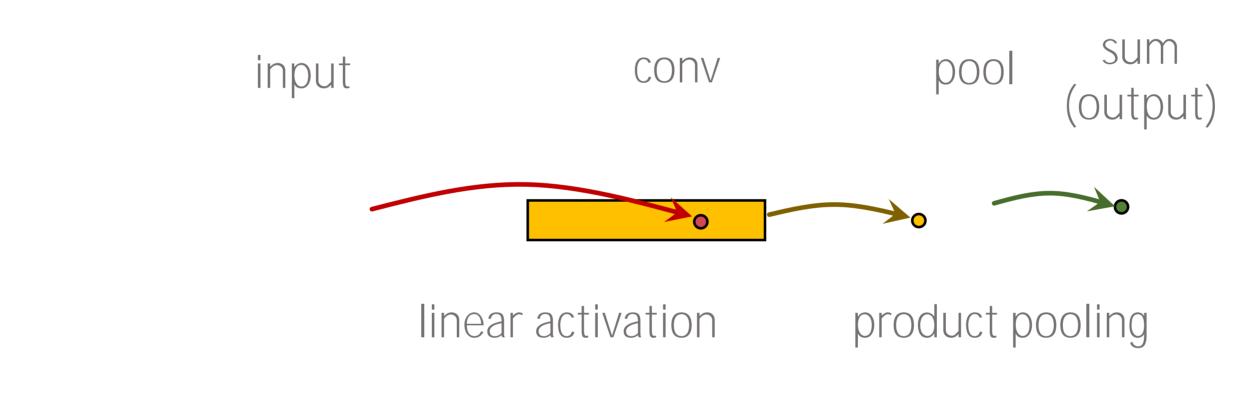
$$\min_{f \mathbf{w}_r^n g_{r;n}} \times \bigvee_{(i_1; \dots; i_N) 2} P_{R} \mathbf{w}_r^1 \qquad \mathbf{w}_r^N \bigvee_{i_1; \dots; i_N} y_{i_1; \dots; i_N}$$

solving tensor completion via NN with multiplicative non-linearity

Tensor Factorization

 $+ \bullet \bullet \bullet +$

Non-Linear Neural Network



Tensor rank: min # of components required to express a tensor

Razin & Cohen 2020: GD empirically minimizes tensor rank even when R is large

Question: can this empirical phenomenon be supported theoretically?

IV) Dynamics of Learning: Theoretical Analysis

Theorem (Dynamical Characterization of Component Norms) GD (with small init and step size) over TF leads component norms to evolve by:

$$\frac{d}{dt} \qquad \underset{n=1}{\overset{N}{\text{w}_r^n(t)}} \qquad \underset{n=1}{\overset{N}{\text{w}_r^n(t)}} \qquad 2 \quad 2=N$$

Interpretation

 $_{n=1}^{N}\mathbf{w}_{r}^{n}(t)$ evolves at a rate proportional to its size exponentiated by 2 2=N>1

Momentum-like effect: components move slower when small and faster when large

Small init = incremental learning of components = low tensor rank

Theorem above leads to:

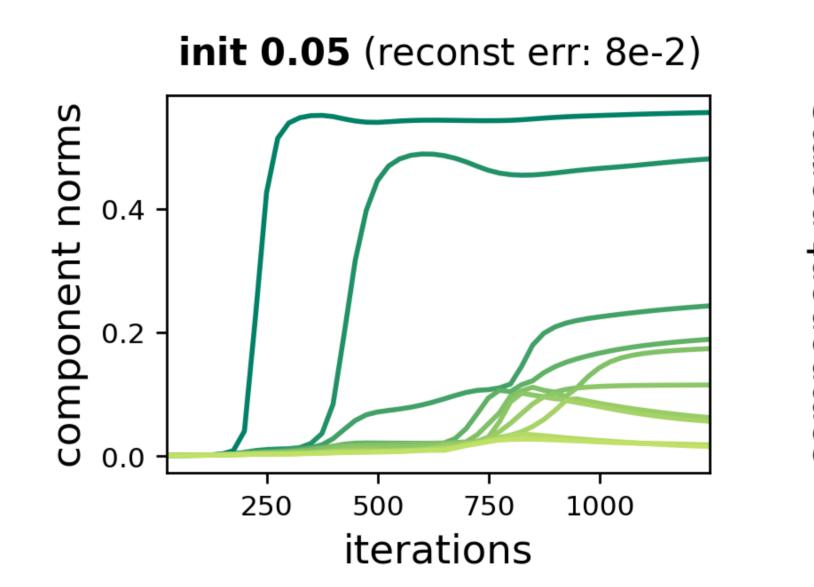
Theorem (Rank 1 Trajectory)

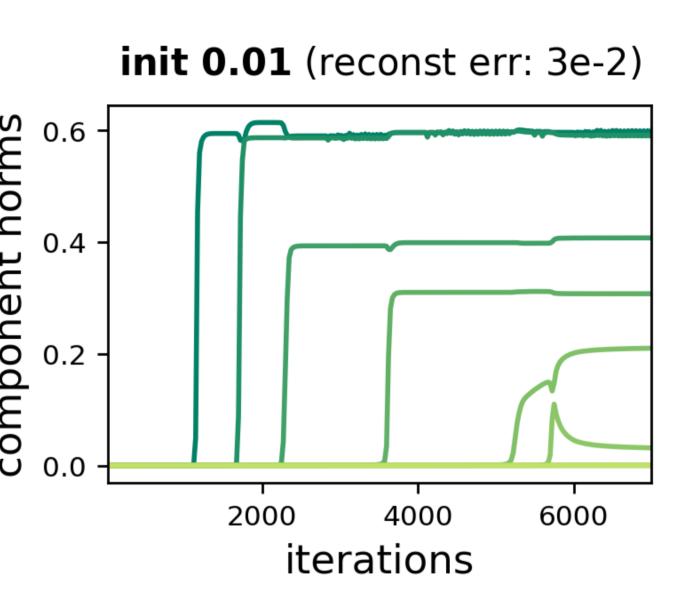
If tensor completion has a rank 1 solution, then under certain technical conditions and a sufficiently small init TF will reach it

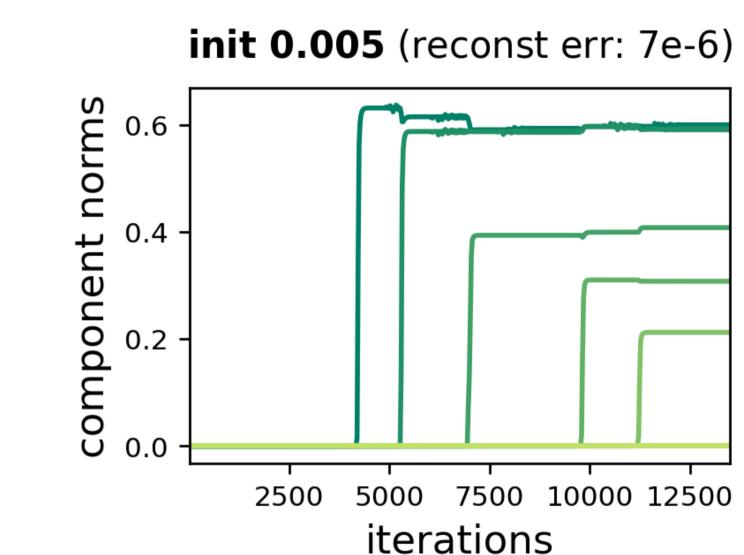
V) Dynamics of Learning: Experiments

Experiment (Rank 5 Order 4 Tensor Completion)

Component norms during GD over TF with different init scales:







0 fewer components depart from zero

Incremental learning of components leads to low tensor rank!

VI) Tensor Rank as Measure of Complexity

Our analysis: tensor rank captures the implicit regularization of a non-linear NN

Can tensor rank serve as a measure of complexity for predictors?

Experiment (Fitting Standard Datasets With Predictors of Low Tensor Rank)



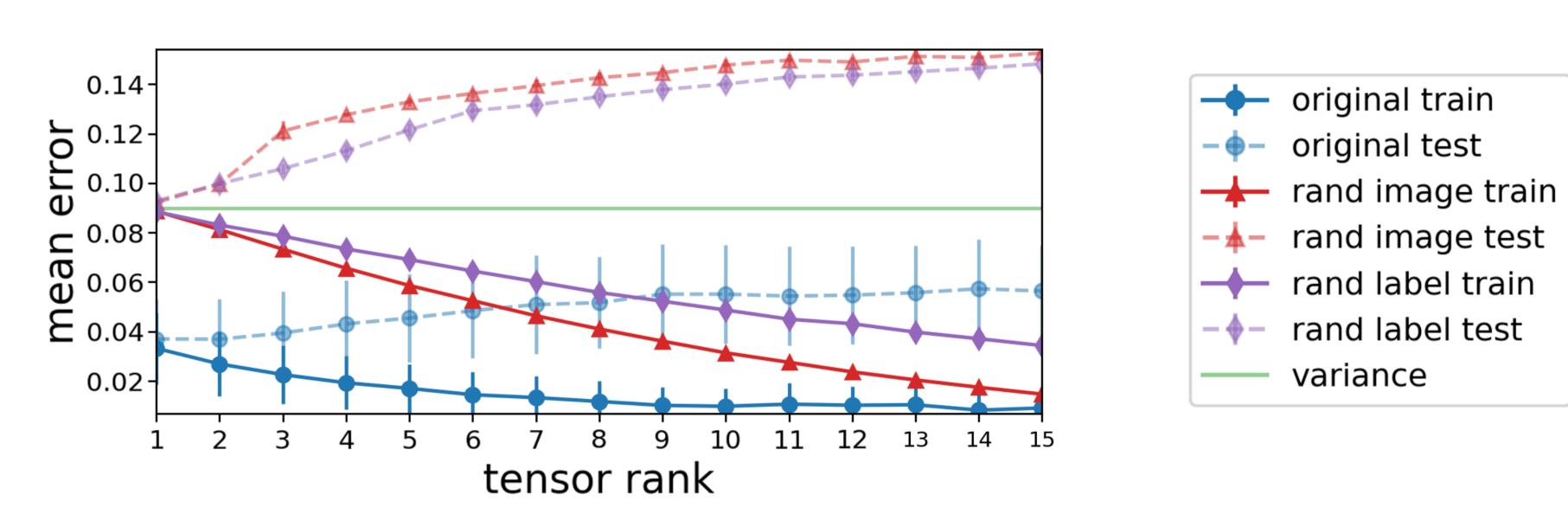




Dataset: Fashion-MNIST (similar results for MNIST)

Compared against two randomized variants:

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



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

Standard datasets can be fit with predictors of low tensor rank!

VII) Takeaway

Tensor rank may pave way to understanding:

- Implicit regularization of neural networks
- Properties of real-world data translating it to generalization