# Optimization and Generalization for Deep Linear Neural Networks via Trajectories of Gradient Descent

**Nadav Cohen**

Tel Aviv University

*Princeton University, Computer Science Department*
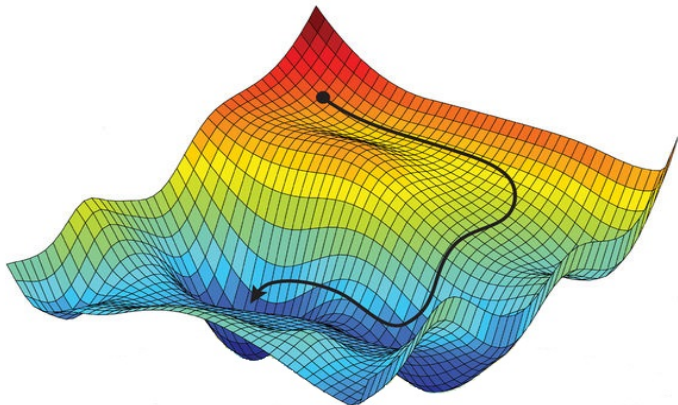*Theoretical Deep Learning Course (COS 597B)*

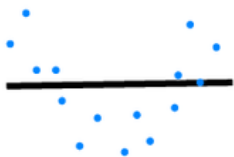6 December 2019

## Outline

## Optimization

Fitting training data by minimizing an objective (loss) function

## Generalization

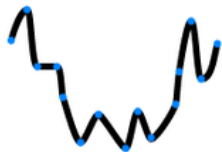Controlling gap between train and test errors, e.g. by adding regularization term/constraint to objective
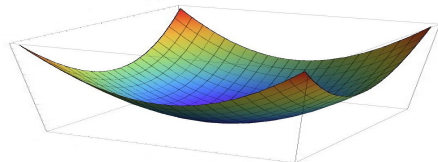


Underfitting     Desired     Overfitting

## Classical Machine Learning



**Theme:** make sure objective is convex!

## Classical Machine Learning



**Theme:** make sure objective is convex!

#### Optimization

- Single global minimum, efficiently attainable
- Choice of algorithm affects only speed of convergence

# Classical Machine Learning



**Theme:** make sure objective is convex!

## Optimization

- Single global minimum, efficiently attainable
- Choice of algorithm affects only speed of convergence

## Generalization

Bias-variance trade-off:

| regularization | train/test gap | train err |
|:---:|:---:|:---:|
| more | ↘ | ↗ |
| less | ↗ | ↘ |

# Classical Machine Learning



**Theme:** make sure objective is convex!

## Optimization

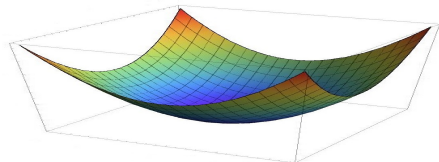- Single global minimum, efficiently attainable
- Choice of algorithm affects only speed of convergence

## Generalization

Bias-variance trade-off:

| regularization | train/test gap | train err |
|:---:|:---:|:---:|
| more | ↘ | ↗ |
| less | ↗ | ↘ |

# Deep Learning (DL)



**Theme:** allow objective to be non-convex

# Deep Learning (DL)



**Theme:** allow objective to be non-convex

## Optimization

- Multiple minima, a-priori not efficiently attainable
- Variants of gradient descent (GD) somehow reach global min

# Deep Learning (DL)



**Theme:** allow objective to be non-convex

## **Optimization**

- Multiple minima, a-priori not efficiently attainable
- Variants of gradient descent (GD) somehow reach global min

## **Generalization**

- Some global minima generalize well, others don't
- With typical data, solution found by GD often generalizes well
- No bias-variance trade-off — regularization implicitly induced by GD

# Deep Learning (DL)



**Theme:** allow objective to be non-convex

## Optimization

- Multiple minima, a-priori not efficiently attainable
- Variants of gradient descent (GD) somehow reach global min

## Generalization

- Some global minima generalize well, others don't
- With typical data, solution found by GD often generalizes well
- No bias-variance trade-off — regularization implicitly induced by GD

# Analysis via Trajectories of Gradient Descent

**Perspective**

- Language of classical learning theory may be insufficient for DL

# Analysis via Trajectories of Gradient Descent

**Perspective**

- Language of classical learning theory may be insufficient for DL

- Need to carefully analyze course of learning, i.e. trajectories of GD!

## Analysis via Trajectories of Gradient Descent

**Perspective**

- Language of classical learning theory may be insufficient for DL

- Need to carefully analyze course of learning, i.e. trajectories of GD!



We will demonstrate this for deep linear neural networks

## Outline

## Sources

**On the Optimization of Deep Networks:**
**Implicit Acceleration by Overparameterization**

   Arora + **C** + Hazan

   *International Conference on Machine Learning (ICML) 2018*

**A Convergence Analysis of Gradient Descent for Deep Linear Neural Networks**

   Arora + **C** + Golowich + Hu

   *International Conference on Learning Representations (ICLR) 2019*

**Implicit Regularization in Deep Matrix Factorization**

   Arora + **C** + Hu + Luo

   *Conference on Neural Information Processing Systems (NeurIPS) 2019*

## Collaborators



**Sanjeev Arora**  **Elad Hazan**

**Yuping Luo**  **Wei Hu**  **Noah Golowich**

# Linear Neural Networks

**Linear neural networks** (LNN) are fully-connected neural networks with linear (no) activation

$$\mathbf{x} \longrightarrow \boxed{W_1} \rightarrow \boxed{W_2} \rightarrow \cdots \longrightarrow \boxed{W_N} \longrightarrow \mathbf{y} = W_N \cdots W_2 W_1 \mathbf{x}$$

# Linear Neural Networks

**Linear neural networks** (LNN) are fully-connected neural networks with linear (no) activation

$$\boldsymbol{x} \longrightarrow \boxed{W_1} \longrightarrow \boxed{W_2} \longrightarrow \cdots \longrightarrow \boxed{W_N} \longrightarrow \boldsymbol{y} = W_N \cdots W_2 W_1 \boldsymbol{x}$$

LNN realize only linear mappings, but are highly non-trivial in terms of optimization and generalization

# Linear Neural Networks

**Linear neural networks** (LNN) are fully-connected neural networks with linear (no) activation

$$\boldsymbol{x} \longrightarrow \boxed{W_1} \longrightarrow \boxed{W_2} \longrightarrow \cdots \longrightarrow \boxed{W_N} \longrightarrow \boldsymbol{y} = W_N \cdots W_2 W_1 \boldsymbol{x}$$

LNN realize only linear mappings, but are highly non-trivial in terms of optimization and generalization

Studied extensively as surrogate for non-linear neural networks:

- Saxe et al. 2014
- Kawaguchi 2016
- Advani & Saxe 2017
- Hardt & Ma 2017

- Laurent & Brecht 2018
- Gunasekar et al. 2018
- Ji & Telgarsky 2019
- Lampinen & Ganguli 2019

## Outline

# Gradient Flow

**Gradient flow** (GF) is a continuous version of GD (step size $\rightarrow 0$):

$$\frac{d}{dt}\boldsymbol{\alpha}(t) = -\nabla f(\boldsymbol{\alpha}(t)) \quad , \ t \in \mathbb{R}_{>0}$$

# Gradient Flow

**Gradient flow** (GF) is a continuous version of GD (step size $\rightarrow$ 0):

$$\frac{d}{dt}\boldsymbol{\alpha}(t) = -\nabla f(\boldsymbol{\alpha}(t)) \ \ , \ t \in \mathbb{R}_{>0}$$



Gradient descent

**Gradient flow**

Admits use of theoretical tools from differential geometry/equations

# Balanced Trajectories



$$\boldsymbol{x} \rightarrow \boxed{W_1} \rightarrow \boxed{W_2} \rightarrow \cdots \rightarrow \boxed{W_N} \rightarrow \quad \boldsymbol{y} = W_N \cdots W_2 W_1 \boldsymbol{x}$$

## Balanced Trajectories

$$\boldsymbol{x} \longrightarrow \boxed{W_1} \rightarrow \boxed{W_2} \rightarrow \cdots \rightarrow \boxed{W_N} \rightarrow \boldsymbol{y} = W_N \cdots W_2 W_1 \boldsymbol{x}$$

Loss $\ell(\cdot)$ for linear model induces **overparameterized objective** for LNN:

$$\phi(W_1, \ldots, W_N) := \ell(W_N \cdots W_2 W_1)$$

## Balanced Trajectories

$$\boldsymbol{x} \longrightarrow \boxed{W_1} \longrightarrow \boxed{W_2} \longrightarrow \cdots \longrightarrow \boxed{W_N} \longrightarrow \boldsymbol{y} = W_N \dots W_2 W_1 \boldsymbol{x}$$

Loss $\ell(\cdot)$ for linear model induces **overparameterized objective** for LNN:

$$\phi(W_1, \dots, W_N) := \ell(W_N \cdots W_2 W_1)$$

### **Definition**

Weights $W_1 \dots W_N$ are **balanced** if $W_{j+1}^\top W_{j+1} = W_j W_j^\top \,, \forall j$.

## Balanced Trajectories

$$\boldsymbol{x} \rightarrow \boxed{W_1} \rightarrow \boxed{W_2} \rightarrow \cdots \rightarrow \boxed{W_N} \rightarrow \boldsymbol{y} = W_N \ldots W_2 W_1 \boldsymbol{x}$$

Loss $\ell(\cdot)$ for linear model induces **overparameterized objective** for LNN:

$$\phi(W_1, \ldots, W_N) := \ell(W_N \cdots W_2 W_1)$$

**Definition**

Weights $W_1 \ldots W_N$ are **balanced** if $W_{j+1}^\top W_{j+1} = W_j W_j^\top$ , $\forall j$.

                    ↑

Holds approximately under $\approx 0$ init, exactly under residual $(I_d)$ init

## Balanced Trajectories



$$\boldsymbol{x} \rightarrow \boxed{W_1} \rightarrow \boxed{W_2} \rightarrow \cdots \rightarrow \boxed{W_N} \rightarrow \boldsymbol{y} = W_N \cdots W_2 W_1 \boldsymbol{x}$$

Loss $\ell(\cdot)$ for linear model induces **overparameterized objective** for LNN:

$$\phi(W_1, \ldots, W_N) := \ell(W_N \cdots W_2 W_1)$$

### Definition

Weights $W_1 \ldots W_N$ are **balanced** if $W_{j+1}^\top W_{j+1} = W_j W_j^\top$, $\forall j$.

$\uparrow$

Holds approximately under $\approx 0$ init, exactly under residual ($I_d$) init

### Claim

*Trajectories of GF over LNN preserve balancedness: if $W_1 \ldots W_N$ are balanced at init, they remain that way throughout GF optimization*

# Balanced Trajectories — Proof

### Claim

*Trajectories of GF over LNN preserve balancedness: if $W_1 \ldots W_N$ are balanced at init, they remain that way throughout GF optimization*

**Proof**

# Balanced Trajectories — Proof

### Claim

*Trajectories of GF over LNN preserve balancedness: if $W_1 \ldots W_N$ are balanced at init, they remain that way throughout GF optimization*

**Proof**

GF over LNN:

$$\tfrac{d}{dt} W_j(t) \;\; = \;\; -\tfrac{\partial}{\partial W_j} \phi\Big( W_1(t), \ldots, W_N(t) \Big)$$

## Balanced Trajectories — Proof

### Claim

*Trajectories of GF over LNN preserve balancedness: if $W_1 \ldots W_N$ are balanced at init, they remain that way throughout GF optimization*

### **Proof**

GF over LNN:

$$
\begin{aligned}
\tfrac{d}{dt} W_j(t) &= -\tfrac{\partial}{\partial W_j} \phi\Big( W_1(t), \ldots, W_N(t) \Big) \\
&= -\prod_{i=j+1}^{N} W_i(t)^{\top} \cdot \nabla\ell\Big( W_N(t) \cdots W_1(t) \Big) \cdot \prod_{i=1}^{j-1} W_i(t)^{\top}
\end{aligned}
$$

## Balanced Trajectories — Proof

### Claim

*Trajectories of GF over LNN preserve balancedness: if $W_1 \ldots W_N$ are balanced at init, they remain that way throughout GF optimization*

### **Proof**

GF over LNN:

$$
\begin{aligned}
\tfrac{d}{dt} W_j(t) &= -\tfrac{\partial}{\partial W_j} \phi\Big( W_1(t), \ldots, W_N(t) \Big) \\
&= -\prod_{i=j+1}^{N} W_i(t)^\top \cdot \nabla\ell\Big( W_N(t) \cdots W_1(t) \Big) \cdot \prod_{i=1}^{j-1} W_i(t)^\top
\end{aligned}
$$

$$
\implies \Big( \tfrac{d}{dt} W_j(t) \Big) W_j(t)^\top \equiv W_{j+1}(t)^\top \Big( \tfrac{d}{dt} W_{j+1}(t) \Big)
$$

## Balanced Trajectories — Proof

### Claim

*Trajectories of GF over LNN preserve balancedness: if $W_1 \ldots W_N$ are balanced at init, they remain that way throughout GF optimization*

**Proof**

GF over LNN:

$$
\begin{aligned}
\tfrac{d}{dt} W_j(t) &= -\tfrac{\partial}{\partial W_j} \phi\Big( W_1(t), \ldots, W_N(t) \Big) \\
&= -\prod_{i=j+1}^{N} W_i(t)^\top \cdot \nabla \ell\Big( W_N(t) \cdots W_1(t) \Big) \cdot \prod_{i=1}^{j-1} W_i(t)^\top
\end{aligned}
$$

$$
\implies \Big( \tfrac{d}{dt} W_j(t) \Big) W_j(t)^\top \equiv W_{j+1}(t)^\top \Big( \tfrac{d}{dt} W_{j+1}(t) \Big)
$$

Take transpose of eq, add to itself, and integrate (w.r.t. $t$):

$$
W_j(t) W_j(t)^\top \equiv W_{j+1}(t)^\top W_{j+1}(t) + const
$$

## Balanced Trajectories — Proof

### Claim

*Trajectories of GF over LNN preserve balancedness: if $W_1 \ldots W_N$ are balanced at init, they remain that way throughout GF optimization*

### **Proof**

GF over LNN:

$$
\begin{aligned}
\tfrac{d}{dt} W_j(t) &= -\tfrac{\partial}{\partial W_j} \phi\Big( W_1(t), \ldots, W_N(t) \Big) \\
&= -\prod_{i=j+1}^{N} W_i(t)^\top \cdot \nabla\ell\Big( W_N(t) \cdots W_1(t) \Big) \cdot \prod_{i=1}^{j-1} W_i(t)^\top
\end{aligned}
$$

$$
\implies \quad \Big( \tfrac{d}{dt} W_j(t) \Big) W_j(t)^\top \equiv W_{j+1}(t)^\top \Big( \tfrac{d}{dt} W_{j+1}(t) \Big)
$$

Take transpose of eq, add to itself, and integrate (w.r.t. $t$):

$$
W_j(t) W_j(t)^\top \equiv W_{j+1}(t)^\top W_{j+1}(t) + const
$$

Balance at init $\implies const = 0$ □

# Implicit Preconditioning

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

How does **end-to-end matrix** $W_{1:N} := W_N \cdots W_1$ move on GF trajectories?

*Linear Neural Network*                    *Equivalent Linear Model*



Gradient flow over $\phi(W_1, ..., W_N)$                    **?**

# Implicit Preconditioning

## Question

How does **end-to-end matrix** $W_{1:N} := W_N \cdots W_1$ move on GF trajectories?

**Linear Neural Network**



Gradient flow over $\phi(W_1, \ldots, W_N)$

**Equivalent Linear Model**

**Preconditioned**
gradient flow over $\ell(W_{1:N})$

### Theorem

If $W_1 \ldots W_N$ are balanced at init, $W_{1:N}$ follows **end-to-end dynamics**:

$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$

where $P_{W_{1:N}(t)}$ is a preconditioner (PSD matrix) that "reinforces" $W_{1:N}(t)$

# Implicit Preconditioning

**Question**

How does **end-to-end matrix** $W_{1:N} := W_N \cdots W_1$ move on GF trajectories?



**Linear Neural Network**

Gradient flow over $\phi(W_1, ..., W_N)$

**Equivalent Linear Model**

**Preconditioned**
gradient flow over $\ell(W_{1:N})$

## Theorem

*If $W_1 \ldots W_N$ are balanced at init, $W_{1:N}$ follows **end-to-end dynamics**:*

$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla \ell(W_{1:N}(t))\right]$$

*where $P_{W_{1:N}(t)}$ is a preconditioner (PSD matrix) that "reinforces" $W_{1:N}(t)$*

$$P_{W_{1:N}(t)} \cdot vec\left[\nabla \ell(W_{1:N}(t))\right] =$$
$$vec\left[\sum_{j=1}^{N}\left[W_{1:N}(t)W_{1:N}(t)^{\top}\right]^{\frac{N-j}{N}} \cdot \nabla \ell(W_{1:N}(t)) \cdot \left[W_{1:N}(t)^{\top}W_{1:N}(t)\right]^{\frac{j-1}{N}}\right]$$

# Implicit Preconditioning

## Question

How does **end-to-end matrix** $W_{1:N} := W_N \cdots W_1$ move on GF trajectories?



**Linear Neural Network**

Gradient flow over $\phi(W_1, ..., W_N)$

**Equivalent Linear Model**

**Preconditioned** gradient flow over $\ell(W_{1:N})$

### Theorem

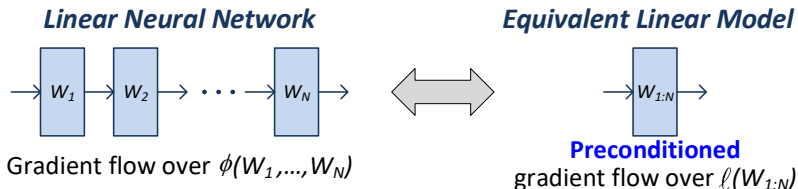If $W_1 \ldots W_N$ are balanced at init, $W_{1:N}$ follows **end-to-end dynamics**:

$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$

where $P_{W_{1:N}(t)}$ is a preconditioner (PSD matrix) that "reinforces" $W_{1:N}(t)$

**Adding (redundant) linear layers to classic linear model induces preconditioner promoting movement in directions already taken!**

# Implicit Preconditioning — Proof Sketch

## Theorem

*If $W_1 \ldots W_N$ are balanced at init, $W_{1:N}$ follows **end-to-end dynamics***:
$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla \ell(W_{1:N}(t))\right]$$
*where $P_{W_{1:N}(t)}$ is a preconditioner (PSD matrix) that "reinforces" $W_{1:N}(t)$*

**Proof Sketch**

# Implicit Preconditioning — Proof Sketch

### Theorem

*If $W_1 \ldots W_N$ are balanced at init, $W_{1:N}$ follows **end-to-end dynamics**:*
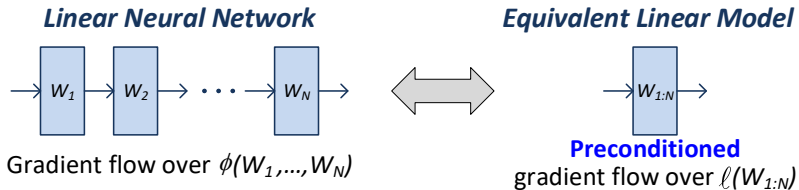$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla \ell(W_{1:N}(t))\right]$$
*where $P_{W_{1:N}(t)}$ is a preconditioner (PSD matrix) that "reinforces" $W_{1:N}(t)$*

**Proof Sketch**

SVD: $W_j(t) = U_j(t)S_j(t)V_j(t)^\top$

# Implicit Preconditioning — Proof Sketch

## Theorem

*If $W_1 \ldots W_N$ are balanced at init, $W_{1:N}$ follows **end-to-end dynamics**:*
$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$
*where $P_{W_{1:N}(t)}$ is a preconditioner (PSD matrix) that "reinforces" $W_{1:N}(t)$*

## **Proof Sketch**

SVD: $W_j(t) = U_j(t)S_j(t)V_j(t)^\top$

Balance $(W_j(t)W_j(t)^\top \equiv W_{j+1}(t)^\top W_{j+1}(t)) \implies S_j(t) \equiv S_{j+1}(t) \wedge U_j(t) \equiv V_{j+1}(t)$

# Implicit Preconditioning — Proof Sketch

### Theorem

If $W_1 \ldots W_N$ are balanced at init, $W_{1:N}$ follows **end-to-end dynamics**:
$$\tfrac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$
where $P_{W_{1:N}(t)}$ is a preconditioner (PSD matrix) that "reinforces" $W_{1:N}(t)$

### **Proof Sketch**

SVD: $W_j(t) = U_j(t)S_j(t)V_j(t)^\top$

Balance $(W_j(t)W_j(t)^\top \equiv W_{j+1}(t)^\top W_{j+1}(t)) \implies S_j(t) \equiv S_{j+1}(t) \wedge U_j(t) \equiv V_{j+1}(t)$

Products of weights thus simplify, yielding:

$$\tfrac{d}{dt}W_{1:N}(t) =$$
$$\sum_{j=1}^{N} \prod_{j+1}^{i=N} W_i(t) \cdot \tfrac{d}{dt}W_j(t) \cdot \prod_{1}^{i=j-1} W_i(t)$$

# Implicit Preconditioning — Proof Sketch

### Theorem

If $W_1 \ldots W_N$ are balanced at init, $W_{1:N}$ follows **end-to-end dynamics**:
$$\tfrac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla \ell(W_{1:N}(t))\right]$$
where $P_{W_{1:N}(t)}$ is a preconditioner (PSD matrix) that "reinforces" $W_{1:N}(t)$

### **Proof Sketch**

SVD: $W_j(t) = U_j(t)S_j(t)V_j(t)^\top$

Balance $(W_j(t)W_j(t)^\top \equiv W_{j+1}(t)^\top W_{j+1}(t)) \implies S_j(t) \equiv S_{j+1}(t) \wedge U_j(t) \equiv V_{j+1}(t)$

Products of weights thus simplify, yielding:

$$\tfrac{d}{dt} W_{1:N}(t) =$$
$$\sum_{j=1}^{N} \prod_{j+1}^{i=N} W_i(t) \cdot \left( -\tfrac{\partial}{\partial W_j} \phi(W_1(t), \ldots, W_N(t)) \right) \cdot \prod_{1}^{i=j-1} W_i(t)$$

# Implicit Preconditioning — Proof Sketch

### Theorem

If $W_1 \ldots W_N$ are balanced at init, $W_{1:N}$ follows **end-to-end dynamics**:
$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$
where $P_{W_{1:N}(t)}$ is a preconditioner (PSD matrix) that "reinforces" $W_{1:N}(t)$

### **Proof Sketch**

SVD: $W_j(t) = U_j(t)S_j(t)V_j(t)^\top$

Balance $(W_j(t)W_j(t)^\top \equiv W_{j+1}(t)^\top W_{j+1}(t)) \implies S_j(t) \equiv S_{j+1}(t) \wedge U_j(t) \equiv V_{j+1}(t)$

Products of weights thus simplify, yielding:

$$\frac{d}{dt} W_{1:N}(t) =$$
$$\sum_{j=1}^{N} \prod_{j+1}^{i=N} W_i(t) \cdot \left( - \prod_{i=j+1}^{N} W_i(t)^\top \nabla\ell(W_{1:N}(t)) \prod_{i=1}^{j-1} W_i(t)^\top \right) \cdot \prod_{1}^{i=j-1} W_i(t)$$

# Implicit Preconditioning — Proof Sketch

### Theorem

*If $W_1 \ldots W_N$ are balanced at init, $W_{1:N}$ follows **end-to-end dynamics**:*
$$\tfrac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$
*where $P_{W_{1:N}(t)}$ is a preconditioner (PSD matrix) that "reinforces" $W_{1:N}(t)$*

### **Proof Sketch**

SVD: $W_j(t) = U_j(t)S_j(t)V_j(t)^\top$

Balance $(W_j(t)W_j(t)^\top \equiv W_{j+1}(t)^\top W_{j+1}(t)) \implies S_j(t) \equiv S_{j+1}(t) \wedge U_j(t) \equiv V_{j+1}(t)$

Products of weights thus simplify, yielding:

$$\tfrac{d}{dt} W_{1:N}(t) =$$
$$-\sum_{j=1}^{N} \prod_{j+1}^{i=N} W_i(t) \prod_{i=j+1}^{N} W_i(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot \prod_{i=1}^{j-1} W_i(t)^\top \prod_{1}^{i=j-1} W_i(t)$$

# Implicit Preconditioning — Proof Sketch

### Theorem

If $W_1 \ldots W_N$ are balanced at init, $W_{1:N}$ follows **end-to-end dynamics**:
$$\tfrac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$
where $P_{W_{1:N}(t)}$ is a preconditioner (PSD matrix) that "reinforces" $W_{1:N}(t)$

### **Proof Sketch**

SVD: $W_j(t) = U_j(t)S_j(t)V_j(t)^\top$

Balance $(W_j(t)W_j(t)^\top \equiv W_{j+1}(t)^\top W_{j+1}(t)) \implies S_j(t) \equiv S_{j+1}(t) \land U_j(t) \equiv V_{j+1}(t)$

Products of weights thus simplify, yielding:
$$\tfrac{d}{dt} W_{1:N}(t) =$$
$$-\sum_{j=1}^{N} \prod_{j+1}^{i=N} W_i(t) \prod_{i=j+1}^{N} W_i(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot \prod_{i=1}^{j-1} W_i(t)^\top \prod_{1}^{i=j-1} W_i(t)$$

# Implicit Preconditioning — Proof Sketch

## Theorem

If $W_1 \ldots W_N$ are balanced at init, $W_{1:N}$ follows **end-to-end dynamics**:
$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$
where $P_{W_{1:N}(t)}$ is a preconditioner (PSD matrix) that "reinforces" $W_{1:N}(t)$

## **Proof Sketch**

SVD: $W_j(t) = U_j(t) S_j(t) V_j(t)^\top$

Balance $(W_j(t) W_j(t)^\top \equiv W_{j+1}(t)^\top W_{j+1}(t)) \implies S_j(t) \equiv S_{j+1}(t) \wedge U_j(t) \equiv V_{j+1}(t)$

Products of weights thus simplify, yielding:
$$\frac{d}{dt} W_{1:N}(t) =$$
$$-\sum_{j=1}^{N} \left[W_{1:N}(t) W_{1:N}(t)^\top\right]^{\frac{N-j}{N}} \cdot \nabla\ell(W_{1:N}(t)) \cdot \left[W_{1:N}(t)^\top W_{1:N}(t)\right]^{\frac{j-1}{N}}$$

# Implicit Preconditioning — Proof Sketch

## Theorem

*If $W_1 \ldots W_N$ are balanced at init, $W_{1:N}$ follows **end-to-end dynamics**:*
$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$
*where $P_{W_{1:N}(t)}$ is a preconditioner (PSD matrix) that "reinforces" $W_{1:N}(t)$*

## **Proof Sketch**

SVD: $W_j(t) = U_j(t)S_j(t)V_j(t)^\top$

Balance $(W_j(t)W_j(t)^\top \equiv W_{j+1}(t)^\top W_{j+1}(t)) \implies S_j(t) \equiv S_{j+1}(t) \wedge U_j(t) \equiv V_{j+1}(t)$

Products of weights thus simplify, yielding:

$$\frac{d}{dt} W_{1:N}(t) =$$
$$-\sum_{j=1}^{N}\left[W_{1:N}(t)W_{1:N}(t)^\top\right]^{\frac{N-j}{N}} \cdot \nabla\ell(W_{1:N}(t)) \cdot \left[W_{1:N}(t)^\top W_{1:N}(t)\right]^{\frac{j-1}{N}}$$

Vectorizing gives end-to-end dynamics (with closed-form expression for $P_{W_{1:N}(t)}$)

# Trajectories Cannot Be Emulated via Regularization

End-to-end dynamics (implicit preconditioning):

$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$

# Trajectories Cannot Be Emulated via Regularization

End-to-end dynamics (implicit preconditioning):

$\frac{d}{dt} vec \left[ W_{1:N}(t) \right] = -P_{W_{1:N}(t)} \cdot vec \left[ \nabla \ell(W_{1:N}(t)) \right]$

### Theorem

*If $\nabla \ell(0) \neq 0$ then $\nexists$ function $F(W)$ s.t. $vec \left[ \nabla F(W) \right] = P_W \cdot vec[\nabla \ell(W)]$*

# Trajectories Cannot Be Emulated via Regularization

End-to-end dynamics (implicit preconditioning):

$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right] \not\equiv -vec\left[\nabla F(W_{1:N}(t))\right]$$

### Theorem

*If $\nabla\ell(0) \neq 0$ then $\nexists$ function $F(W)$ s.t. $vec\left[\nabla F(W)\right] = P_W \cdot vec[\nabla\ell(W)]$*

# Trajectories Cannot Be Emulated via Regularization

End-to-end dynamics (implicit preconditioning):

$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right] \not\equiv -vec\left[\nabla F(W_{1:N}(t))\right]$

### Theorem

*If $\nabla\ell(0) \neq 0$ then $\nexists$ function $F(W)$ s.t. $vec\left[\nabla F(W)\right] = P_W \cdot vec[\nabla\ell(W)]$*

**Trajectories with LNN cannot be emulated by regularizing objective!**

# Trajectories Cannot Be Emulated via Regularization

End-to-end dynamics (implicit preconditioning):
$$\frac{d}{dt}vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla \ell(W_{1:N}(t))\right] \neq -vec\left[\nabla F(W_{1:N}(t))\right]$$

## Theorem

If $\nabla \ell(0) \neq 0$ then $\nexists$ function $F(W)$ s.t. $vec\left[\nabla F(W)\right] = P_W \cdot vec[\nabla \ell(W)]$

**Trajectories with LNN cannot be emulated by regularizing objective!**



contradicts gradient theorem!

$$\int_{\Gamma} P_W \cdot vec\left[\nabla \ell(W)\right] \neq 0$$

# Outline

# Classic Approach: Characterization of Critical Points

Prominent approach for analyzing optimization in DL (in spirit of classical learning theory) is via critical points in the objective



*Good local minimum*
*(≈ global minimum)*

*Poor local minimum*

*Strict saddle*

*Non-strict saddle*

# Classic Approach: Characterization of Critical Points

Prominent approach for analyzing optimization in DL (in spirit of classical learning theory) is via critical points in the objective



**Good local minimum** (≈ global minimum)

**(1)** **Poor local minimum**

**Strict saddle**

**(2)** **Non-strict saddle**

**Result** *(cf. Ge et al. 2015; Lee et al. 2016)*

If: **(1)** there are no poor local minima; and **(2)** all saddle points are strict, then GD converges to global min

# Classic Approach: Characterization of Critical Points

Prominent approach for analyzing optimization in DL (in spirit of classical learning theory) is via critical points in the objective



**(1)**

**(2)**

*Good local minimum (≈ global minimum)*  *Poor local minimum*  *Strict saddle*  *Non-strict saddle*

**Result** *(cf. Ge et al. 2015; Lee et al. 2016)*

If: **(1)** there are no poor local minima; and **(2)** all saddle points are strict, then GD converges to global min

Motivated by this, many [1] studied the validity of **(1)** and/or **(2)**

---

[1] e.g. Haeffele & Vidal 2015; Kawaguchi 2016; Soudry & Carmon 2016; Safran & Shamir 2018
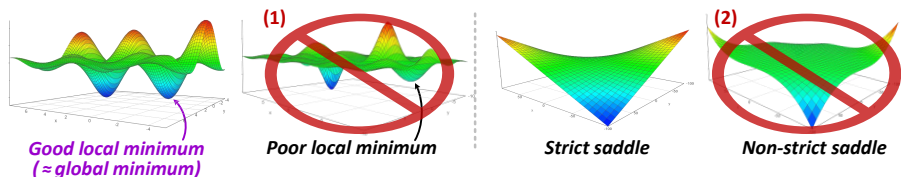
# Classic Approach: Characterization of Critical Points

Prominent approach for analyzing optimization in DL (in spirit of classical learning theory) is via critical points in the objective
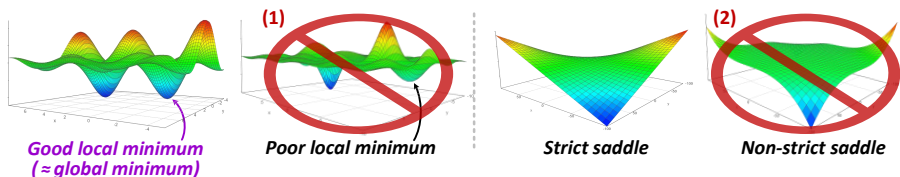


**Good local minimum** ($\approx$ global minimum)

**Poor local minimum**

**Strict saddle**

**Non-strict saddle**

**Result** *(cf. Ge et al. 2015; Lee et al. 2016)*

If: **(1)** there are no poor local minima; and **(2)** all saddle points are strict, then GD converges to global min

Motivated by this, many [1] studied the validity of **(1)** and/or **(2)**

**Limitation:** deep ($\geq 3$ layer) models violate **(2)** (consider all weights $= 0$)!

---

[1] e.g. Haeffele & Vidal 2015; Kawaguchi 2016; Soudry & Carmon 2016; Safran & Shamir 2018

# Applying Our Trajectory Analysis

# Applying Our Trajectory Analysis

Trajectory analysis revealed implicit preconditioning on end-to-end matrix:

$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$

# Applying Our Trajectory Analysis

Trajectory analysis revealed implicit preconditioning on end-to-end matrix:

$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$

$P_{W_{1:N}(t)} \succ 0$ when $W_{1:N}(t)$ has full rank

# Applying Our Trajectory Analysis

Trajectory analysis revealed implicit preconditioning on end-to-end matrix:

$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$

$P_{W_{1:N}(t)} \succ 0$ when $W_{1:N}(t)$ has full rank $\implies$ loss decreases until:

(1) $\nabla\ell(W_{1:N}(t)) = 0$   **or**   (2) $W_{1:N}(t)$ is singular

# Applying Our Trajectory Analysis

Trajectory analysis revealed implicit preconditioning on end-to-end matrix:

$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla \ell(W_{1:N}(t))\right]$$

$P_{W_{1:N}(t)} \succ 0$ when $W_{1:N}(t)$ has full rank $\implies$ loss decreases until:

**(1)** $\nabla \ell(W_{1:N}(t)) = 0$ **or** **(2)** $W_{1:N}(t)$ is singular

$\ell(\cdot)$ is typically convex $\implies$ **(1)** means global min was reached

# Applying Our Trajectory Analysis

Trajectory analysis revealed implicit preconditioning on end-to-end matrix:

$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$

$P_{W_{1:N}(t)} \succ 0$ when $W_{1:N}(t)$ has full rank $\implies$ loss decreases until:

(1) $\nabla\ell(W_{1:N}(t)) = 0$ **or** (2) $W_{1:N}(t)$ is singular

$\ell(\cdot)$ is typically convex $\implies$ (1) means global min was reached

### Corollary

*Assume $\ell(\cdot)$ is convex and LNN is init such that:*

1. $\ell(W_{1:N}) < \ell(W)$ *for any singular $W$*

2. $W_1 \ldots W_N$ *are balanced*

*Then, GF converges to global min*

# From Gradient Flow to Gradient Descent

# From Gradient Flow to Gradient Descent

**Corollary**

*Assume $\ell(\cdot)$ is convex and LNN is init such that:*

1. *$\ell(W_{1:N}) < \ell(W)$ for any singular $W$*

2. *$W_1 \ldots W_N$ are balanced*

*Then, GF converges to global min*

# From Gradient Flow to Gradient Descent

**Corollary**

*Assume $\ell(\cdot)$ is convex and LNN is init such that:*

1. $\ell(W_{1:N}) < \ell(W)$ , $\forall W$ *s.t.* $\sigma_{min}(W) = 0$

2. $W_1 \ldots W_N$ *are balanced*

*Then, GF converges to global min*

# From Gradient Flow to Gradient Descent

## Corollary

Assume $\ell(\cdot)$ is convex and LNN is init such that:

1. $\ell(W_{1:N}) < \ell(W)$ , $\forall W$ s.t. $\sigma_{min}(W) = 0$

2. $W_{j+1}^\top W_{j+1} = W_j W_j^\top$ , $\forall j$

Then, GF converges to global min

# From Gradient Flow to Gradient Descent

## Corollary

*Assume $\ell(\cdot)$ is convex and LNN is init such that:*

1. $\ell(W_{1:N}) < \ell(W)$ , $\forall W$ s.t. $\sigma_{min}(W) = 0$

2. $\|W_{j+1}^\top W_{j+1} - W_j W_j^\top\|_F = 0$ , $\forall j$

*Then, GF converges to global min*

# From Gradient Flow to Gradient Descent

## Theorem

*Assume $\ell(\cdot)$ is convex and LNN is init such that:*

1. $\ell(W_{1:N}) < \ell(W)$ , $\forall W$ s.t. $\sigma_{min}(W) = 0$

2. $\|W_{j+1}^\top W_{j+1} - W_j W_j^\top\|_F = 0$ , $\forall j$

*Then, GF converges to global min*

# From Gradient Flow to Gradient Descent

## Theorem

*Assume $\ell(\cdot) = \ell_2$ loss and LNN is init such that:*

1. $\ell(W_{1:N}) < \ell(W)$  ,$\forall W$ *s.t.* $\sigma_{min}(W) = 0$

2. $\|W_{j+1}^\top W_{j+1} - W_j W_j^\top\|_F = 0$  ,$\forall j$

*Then, GF converges to global min*

# From Gradient Flow to Gradient Descent

## Theorem

Assume $\ell(\cdot) = \ell_2$ loss and LNN is init such that:

1. $\ell(W_{1:N}) < \ell(W)$  $,\forall W$ s.t. $\sigma_{min}(W) \leq c$

2. $\|W_{j+1}^\top W_{j+1} - W_j W_j^\top\|_F = 0$  $,\forall j$

Then, GF converges to global min

# From Gradient Flow to Gradient Descent

## Theorem

*Assume $\ell(\cdot) = \ell_2$ loss and LNN is init such that:*

1. $\ell(W_{1:N}) < \ell(W)$ , $\forall W$ s.t. $\sigma_{min}(W) \leq c$

2. $\|W_{j+1}^\top W_{j+1} - W_j W_j^\top\|_F \leq \mathcal{O}(c^2)$ , $\forall j$

*Then, GF converges to global min*

# From Gradient Flow to Gradient Descent

## Theorem

Assume $\ell(\cdot) = \ell_2$ loss and LNN is init such that:

1. $\ell(W_{1:N}) < \ell(W)$ , $\forall W$ s.t. $\sigma_{min}(W) \leq c$

2. $\|W_{j+1}^\top W_{j+1} - W_j W_j^\top\|_F \leq \mathcal{O}(c^2)$ , $\forall j$

Then, GD with step size $\eta \leq \mathcal{O}(c^4)$ gives: loss(iteration t) $\leq e^{-\Omega(c^2 \eta t)}$

# From Gradient Flow to Gradient Descent

### Theorem

Assume $\ell(\cdot) = \ell_2$ loss and LNN is init such that:

1. $\ell(W_{1:N}) < \ell(W)$ , $\forall W$ s.t. $\sigma_{min}(W) \leq c$

2. $\|W_{j+1}^\top W_{j+1} - W_j W_j^\top\|_F \leq \mathcal{O}(c^2)$ , $\forall j$

Then, *GD* with step size $\eta \leq \mathcal{O}(c^4)$ gives: *loss(iteration t)* $\leq e^{-\Omega(c^2\eta t)}$

### Claim

Our assumptions on init:

# From Gradient Flow to Gradient Descent

## Theorem

*Assume $\ell(\cdot) = \ell_2$ loss and LNN is init such that:*

1. $\ell(W_{1:N}) < \ell(W)$ , $\forall W$ *s.t.* $\sigma_{min}(W) \leq c$

2. $\|W_{j+1}^\top W_{j+1} - W_j W_j^\top\|_F \leq \mathcal{O}(c^2)$ , $\forall j$

*Then, GD with step size $\eta \leq \mathcal{O}(c^4)$ gives: loss(iteration t) $\leq e^{-\Omega(c^2 \eta t)}$*

## Claim

*Our assumptions on init:*

- *Are necessary (violating any of them can lead to divergence)*

# From Gradient Flow to Gradient Descent

## Theorem

Assume $\ell(\cdot) = \ell_2$ loss and LNN is init such that:

1. $\ell(W_{1:N}) < \ell(W)$  , $\forall W$ s.t. $\sigma_{min}(W) \leq c$

2. $\|W_{j+1}^\top W_{j+1} - W_j W_j^\top\|_F \leq \mathcal{O}(c^2)$  , $\forall j$

Then, *GD* with step size $\eta \leq \mathcal{O}(c^4)$ gives: *loss(iteration t)* $\leq e^{-\Omega(c^2 \eta t)}$

## Claim

Our assumptions on init:

- Are necessary (violating any of them can lead to divergence)

- For out dim 1, hold with const prob under random "balanced" init

# From Gradient Flow to Gradient Descent

### Theorem

*Assume $\ell(\cdot) = \ell_2$ loss and LNN is init such that:*

1. $\ell(W_{1:N}) < \ell(W)$ , $\forall W$ s.t. $\sigma_{min}(W) \leq c$

2. $\|W_{j+1}^\top W_{j+1} - W_j W_j^\top\|_F \leq \mathcal{O}(c^2)$ , $\forall j$

*Then, GD with step size $\eta \leq \mathcal{O}(c^4)$ gives: loss(iteration t) $\leq e^{-\Omega(c^2 \eta t)}$*

### Claim

*Our assumptions on init:*

- *Are necessary (violating any of them can lead to divergence)*

- *For out dim 1, hold with const prob under random "balanced" init*

**Guarantee of efficient (linear rate) convergence to global min!
Most general guarantee to date for GD efficiently training deep net.**

# Effect of Depth on Optimization

# Effect of Depth on Optimization

**Viewpoint of classical learning theory:**

- Convex optimization is easier than non-convex

# Effect of Depth on Optimization

**Viewpoint of classical learning theory:**

- Convex optimization is easier than non-convex



- Hence depth complicates optimization

# Effect of Depth on Optimization

**Viewpoint of classical learning theory:**

- Convex optimization is easier than non-convex



- Hence depth complicates optimization



**Our trajectory analysis reveals:** not always true...

# Acceleration by Depth

# Acceleration by Depth

Discrete version of end-to-end dynamics for LNN:

$$vec\left[W_{1:N}(t+1)\right] \leftarrow vec\left[W_{1:N}(t)\right] - \eta \cdot P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$

# Acceleration by Depth

Discrete version of end-to-end dynamics for LNN:

$$vec[W_{1:N}(t+1)] \hookleftarrow vec[W_{1:N}(t)] - \eta \cdot P_{W_{1:N}(t)} \cdot vec[\nabla\ell(W_{1:N}(t))]$$

### Claim

$\forall p > 2$, $\exists$ settings where $\ell(\cdot) = \ell_p$ loss (i.e. $\ell(W) = \frac{1}{m}\sum_{i=1}^{m}\|W\mathbf{x}_i - \mathbf{y}_i\|_p^p$) and disc end-to-end dynamics reach global min arbitrarily faster than GD

# Acceleration by Depth

Discrete version of end-to-end dynamics for LNN:

$$vec[W_{1:N}(t+1)] \hookleftarrow vec[W_{1:N}(t)] - \eta \cdot P_{W_{1:N}(t)} \cdot vec[\nabla\ell(W_{1:N}(t))]$$

## Claim

$\forall p > 2, \exists$ settings where $\ell(\cdot) = \ell_p$ loss (i.e. $\ell(W) = \frac{1}{m}\sum_{i=1}^{m}\|W\mathbf{x}_i - \mathbf{y}_i\|_p^p$) and disc end-to-end dynamics reach global min arbitrarily faster than GD

**Experiment**

# Acceleration by Depth

Discrete version of end-to-end dynamics for LNN:

$$vec\big[W_{1:N}(t+1)\big] \hookleftarrow vec\big[W_{1:N}(t)\big] - \eta \cdot P_{W_{1:N}(t)} \cdot vec\big[\nabla\ell(W_{1:N}(t))\big]$$

### Claim

$\forall p > 2$, $\exists$ settings where $\ell(\cdot) = \ell_p$ loss (i.e. $\ell(W) = \frac{1}{m}\sum_{i=1}^{m}\|W\mathbf{x}_i - \mathbf{y}_i\|_p^p$) and disc end-to-end dynamics reach global min arbitrarily faster than GD

**Experiment**

Regression problem from UCI ML Repository ; $\ell_4$ loss

# Acceleration by Depth

Discrete version of end-to-end dynamics for LNN:

$$vec\big[W_{1:N}(t+1)\big] \hookleftarrow vec\big[W_{1:N}(t)\big] - \eta \cdot P_{W_{1:N}(t)} \cdot vec\big[\nabla\ell(W_{1:N}(t))\big]$$

### Claim

$\forall p > 2, \exists$ *settings where* $\ell(\cdot) = \ell_p$ *loss* (*i.e.* $\ell(W) = \frac{1}{m}\sum_{i=1}^{m}\|W\mathbf{x}_i - \mathbf{y}_i\|_p^p$) *and disc end-to-end dynamics reach global min arbitrarily faster than GD*

### **Experiment**

Regression problem from UCI ML Repository ; $\ell_4$ loss

# Acceleration by Depth

Discrete version of end-to-end dynamics for LNN:

$$vec\left[W_{1:N}(t+1)\right] \leftarrow vec\left[W_{1:N}(t)\right] - \eta \cdot P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$

## Claim

$\forall p > 2, \exists$ settings where $\ell(\cdot) = \ell_p$ loss (i.e. $\ell(W) = \frac{1}{m}\sum_{i=1}^{m}\|W\mathbf{x}_i - \mathbf{y}_i\|_p^p$) and disc end-to-end dynamics reach global min arbitrarily faster than GD

## Experiment

Regression problem from UCI ML Repository ; $\ell_4$ loss



**Depth can speed-up GD, even without any gain in expressiveness, and despite introducing non-convexity!**

## Outline

1 Optimization and Generalization in Deep Learning via Trajectories

2 Case Study: Linear Neural Networks
  - Trajectory Analysis
  - Optimization
  - **Generalization**

3 Conclusion

# Setting: Matrix Completion

**Matrix completion:** recover matrix given subset of entries



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

# Setting: Matrix Completion

**Matrix completion:** recover matrix given subset of entries



Can be viewed as classification (regression) problem:

$$observed\ entries \quad \longleftrightarrow \quad training\ data$$
$$unobserved\ entries \quad \longleftrightarrow \quad test\ data$$

# Setting: Matrix Completion

**Matrix completion:** recover matrix given subset of entries



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

Can be viewed as classification (regression) problem:

$$observed\ entries \quad \longleftrightarrow \quad training\ data$$
$$unobserved\ entries \quad \longleftrightarrow \quad test\ data$$

**<u>Standard Assumption</u>**

Matrix to recover (ground truth) has low rank

# Setting: Matrix Completion

**Matrix completion:** recover matrix given subset of entries



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

Can be viewed as classification (regression) problem:

$$observed\ entries \quad \longleftrightarrow \quad training\ data$$
$$unobserved\ entries \quad \longleftrightarrow \quad test\ data$$

## Standard Assumption
Matrix to recover (ground truth) has low rank

## Classical Result *(cf. Candes & Recht 2008)*
Nuclear norm minimization (convex program) perfectly recovers ("almost any") low rank matrix if observations are sufficiently many

# Two-Layer Network ⟷ Matrix Factorization

Matrix completion via two-layer LNN:

- Parameterize ground truth as $W_2 W_1$

# Two-Layer Network $\longleftrightarrow$ Matrix Factorization

Matrix completion via two-layer LNN:

- Parameterize ground truth as $W_2 W_1$



- Known as **matrix factorization** (MF)

## Two-Layer Network $\longleftrightarrow$ Matrix Factorization

Matrix completion via two-layer LNN:

- Parameterize ground truth as $W_2 W_1$

$$
\begin{array}{|c|c|c|c|}
\hline
4 & ? & ? & 4 \\
\hline
? & 5 & 4 & ? \\
\hline
? & 5 & ? & ? \\
\hline
\end{array}
=
\boxed{\quad W_2 \quad}
*
\boxed{\quad W_1 \quad}
$$

- Known as **matrix factorization** (MF)

**Empirical Phenomenon**

GD (with step size $\ll 1$ and init $\approx 0$) over MF recovers low rank matrices, even when shared dim of $W_1, W_2$ doesn't constrain rank!

# Two-Layer Network ⟷ Matrix Factorization

Matrix completion via two-layer LNN:

- Parameterize ground truth as $W_2 W_1$

$$
\begin{array}{|c|c|c|c|}
\hline
4 & ? & ? & 4 \\
\hline
? & 5 & 4 & ? \\
\hline
? & 5 & ? & ? \\
\hline
\end{array}
= \boxed{W_2} * \boxed{W_1}
$$

- Known as **matrix factorization** (MF)

**Empirical Phenomenon**

GD (with step size $\ll 1$ and init $\approx 0$) over MF recovers low rank matrices, even when shared dim of $W_1, W_2$ doesn't constrain rank!

---

Conjecture (Gunasekar et al. 2017)

*GD (with step size $\ll 1$ and init $\approx 0$) over MF converges to solution with min nuclear norm (among those fitting observations)*

---

# Two-Layer Network $\longleftrightarrow$ Matrix Factorization

Matrix completion via two-layer LNN:

- Parameterize ground truth as $W_2 W_1$



- Known as **matrix factorization** (MF)

**Empirical Phenomenon**

GD (with step size $\ll 1$ and init $\approx 0$) over MF recovers low rank matrices, even when shared dim of $W_1, W_2$ doesn't constrain rank!

> Conjecture (Gunasekar et al. 2017)
>
> *GD (with step size $\ll 1$ and init $\approx 0$) over MF converges to solution with min nuclear norm (among those fitting observations)*

Gunasekar et al. proved conjecture for certain restricted setting

# $N$-Layer Network $\longleftrightarrow$ "Deep Matrix Factorization"

Matrix completion via $N$-layer LNN:

- Parameterize ground truth as $W_N \cdots W_2 W_1$

$$\begin{bmatrix} 4 & ? & ? & 4 \\ ? & 5 & 4 & ? \\ ? & 5 & ? & ? \end{bmatrix} = \boxed{W_N} \; * \; \bullet \; \bullet \; \bullet \; * \; \boxed{W_2} \; * \; \boxed{W_1}$$

# $N$-Layer Network $\longleftrightarrow$ "Deep Matrix Factorization"

Matrix completion via $N$-layer LNN:

- Parameterize ground truth as $W_N \cdots W_2 W_1$



- We refer to this as **deep matrix factorization** (DMF)

# $N$-Layer Network $\longleftrightarrow$ "Deep Matrix Factorization"

Matrix completion via $N$-layer LNN:

- Parameterize ground truth as $W_N \cdots W_2 W_1$

$$
\begin{array}{|c|c|c|c|}
\hline
4 & ? & ? & 4 \\
\hline
? & 5 & 4 & ? \\
\hline
? & 5 & ? & ? \\
\hline
\end{array}
= \boxed{W_N} * \cdots * \boxed{W_2} * \boxed{W_1}
$$

- We refer to this as **deep matrix factorization** (DMF)

### Experiment

Completion of low rank matrix via GD over DMF

# $N$-Layer Network $\longleftrightarrow$ "Deep Matrix Factorization"

Matrix completion via $N$-layer LNN:

- Parameterize ground truth as $W_N \cdots W_2 W_1$



- We refer to this as **deep matrix factorization** (DMF)

**Experiment**

Completion of low rank matrix via GD over DMF



**Depth enhanced implicit regularization towards low rank!**

# Can the Implicit Regularization Be Captured by Norms?

## Can the Implicit Regularization Be Captured by Norms?

Conjecture of Gunasekar et al. 2017 (in spirit of classical learning theory):

<table>
<tr>
<td><i>implicit regularization<br>with depth 2 LNN (MF)</i></td>
<td>$\longleftrightarrow$</td>
<td><i>minimizing nuclear norm<br>(surrogate for rank)</i></td>
</tr>
</table>

## Can the Implicit Regularization Be Captured by Norms?

Conjecture of Gunasekar et al. 2017 (in spirit of classical learning theory):

*implicit regularization*          *minimizing nuclear norm*
*with depth 2 LNN (MF)*    $\longleftrightarrow$    *(surrogate for rank)*

In light of our experiment, natural to hypothesize:

*implicit regularization*          *minimizing other norm*
*with deeper LNN (DMF)*    $\longleftrightarrow$    *closer to rank*

## Can the Implicit Regularization Be Captured by Norms?

Conjecture of Gunasekar et al. 2017 (in spirit of classical learning theory):

*implicit regularization*        *minimizing nuclear norm*
*with depth 2 LNN (MF)* $\longleftrightarrow$ *(surrogate for rank)*

In light of our experiment, natural to hypothesize:

*implicit regularization*        *minimizing other norm*
*with deeper LNN (DMF)* $\longleftrightarrow$ *closer to rank*

**Example**

## Can the Implicit Regularization Be Captured by Norms?

Conjecture of Gunasekar et al. 2017 (in spirit of classical learning theory):

$$\begin{array}{ccc} \textit{implicit regularization} & & \textit{minimizing nuclear norm} \\ \textit{with depth 2 LNN (MF)} & \longleftrightarrow & \textit{(surrogate for rank)} \end{array}$$

In light of our experiment, natural to hypothesize:

$$\begin{array}{ccc} \textit{implicit regularization} & & \textit{minimizing other norm} \\ \textit{with deeper LNN (DMF)} & \longleftrightarrow & \textit{closer to rank} \end{array}$$

### **Example**

Schatten-$p$ quasi-norm to the power of $p$:

- $\|W\|_{S_p}^p := \sum_r \sigma_r^p(W)$ where $\sigma_r(W)$ are singular vals of $W$

# Can the Implicit Regularization Be Captured by Norms?

Conjecture of Gunasekar et al. 2017 (in spirit of classical learning theory):

$$\begin{array}{ccc} \textit{implicit regularization} & & \textit{minimizing nuclear norm} \\ \textit{with depth 2 LNN (MF)} & \longleftrightarrow & \textit{(surrogate for rank)} \end{array}$$

In light of our experiment, natural to hypothesize:

$$\begin{array}{ccc} \textit{implicit regularization} & & \textit{minimizing other norm} \\ \textit{with deeper LNN (DMF)} & \longleftrightarrow & \textit{closer to rank} \end{array}$$

## **Example**

Schatten-$p$ quasi-norm to the power of $p$:

- $\|W\|^p_{S_p} := \sum_r \sigma^p_r(W)$ where $\sigma_r(W)$ are singular vals of $W$

- $p = 1$: nuclear norm, corresponds to depth 2 by Gunasekar et al. 2017

## Can the Implicit Regularization Be Captured by Norms?

Conjecture of Gunasekar et al. 2017 (in spirit of classical learning theory):

*implicit regularization*                    *minimizing nuclear norm*
*with depth 2 LNN (MF)*    $\longleftrightarrow$    *(surrogate for rank)*

In light of our experiment, natural to hypothesize:

*implicit regularization*                    *minimizing other norm*
*with deeper LNN (DMF)*    $\longleftrightarrow$    *closer to rank*

### **Example**

Schatten-$p$ quasi-norm to the power of $p$:

- $\|W\|_{S_p}^p := \sum_r \sigma_r^p(W)$ where $\sigma_r(W)$ are singular vals of $W$

- $p = 1$: nuclear norm, corresponds to depth 2 by Gunasekar et al. 2017

- $0 < p < 1$: closer to rank, may correspond to higher depths

# Current Theory is Oblivious to Depth

# Current Theory is Oblivious to Depth

### Theorem

*In restricted setting where Gunasekar et al. proved depth 2 minimizes nuclear norm, any depth > 2 does so as well*

# Current Theory is Oblivious to Depth

### Theorem

*In restricted setting where Gunasekar et al. proved depth 2 minimizes nuclear norm, any depth > 2 does so as well*

### Proposition

*∃ instances of this setting where nuclear norm minimization contradicts Schatten-p quasi-norm minimization (even locally) $\forall p \in (0, 1)$*

# Current Theory is Oblivious to Depth

### Theorem

*In restricted setting where Gunasekar et al. proved depth 2 minimizes nuclear norm, any depth > 2 does so as well*

### Proposition

*$\exists$ instances of this setting where nuclear norm minimization contradicts Schatten-p quasi-norm minimization (even locally) $\forall p \in (0, 1)$*

This implies:

implicit regularization for any depth $\not\equiv$ Schatten quasi-norm minimization

# Current Theory is Oblivious to Depth

### Theorem

*In restricted setting where Gunasekar et al. proved depth 2 minimizes nuclear norm, any depth > 2 does so as well*

### Proposition

*$\exists$ instances of this setting where nuclear norm minimization contradicts Schatten-p quasi-norm minimization (even locally) $\forall p \in (0, 1)$*

This implies:

implicit regularization for any depth $\not\equiv$ Schatten quasi-norm minimization

Instead, adopting lens of Gunasekar et al. leads to conjecturing:

implicit regularization for all depths $\equiv$ nuclear norm minimization

# Current Theory is Oblivious to Depth

### Theorem

*In restricted setting where Gunasekar et al. proved depth 2 minimizes nuclear norm, any depth > 2 does so as well*

### Proposition

*∃ instances of this setting where nuclear norm minimization contradicts Schatten-p quasi-norm minimization (even locally) $\forall p \in (0, 1)$*

This implies:

 implicit regularization for any depth $\not\equiv$ Schatten quasi-norm minimization

Instead, adopting lens of Gunasekar et al. leads to conjecturing:

 implicit regularization for all depths $\equiv$ nuclear norm minimization

**But our experiment shows depth changes implicit regularization!**

# Experiments Testing Nuclear Norm Conjecture

# Experiments Testing Nuclear Norm Conjecture

**Setup**:

- Completion of $100 \times 100$ rank 5 matrix
- Observed entries chosen uniformly at random

# Experiments Testing Nuclear Norm Conjecture

**Setup**:

- Completion of $100 \times 100$ rank 5 matrix
- Observed entries chosen uniformly at random

**Many (5K) Observations**:

|                     | reconst err | nuclear norm | effective rank |
|---------------------|-------------|--------------|----------------|
| nuclear norm min    |             |              |                |
| depth 2 LNN         |             |              |                |
| depth 3 LNN         |             |              |                |

# Experiments Testing Nuclear Norm Conjecture

**Setup**:

- Completion of $100 \times 100$ rank 5 matrix
- Observed entries chosen uniformly at random

**Many (5K) Observations**:

|  | *reconst err* | *nuclear norm* | *effective rank* |
|---|---|---|---|
| *nuclear norm min* | 8 e –07 | 221 | 5 |
| *depth* 2 *LNN* | | | |
| *depth* 3 *LNN* | | | |

- Nuclear norm min recovers ground truth

# Experiments Testing Nuclear Norm Conjecture

**Setup**:

- Completion of $100 \times 100$ rank 5 matrix
- Observed entries chosen uniformly at random

**Many (5K) Observations**:

|                  | reconst err | nuclear norm | effective rank |
|------------------|-------------|--------------|----------------|
| *nuclear norm min* | 8 e −07    | 221          | 5              |
| *depth* 2 *LNN*    | 5 e −06    | 221          | 5              |
| *depth* 3 *LNN*    | 4 e −06    | 221          | 5              |

- Nuclear norm min recovers ground truth
- LNN do so too

# Experiments Testing Nuclear Norm Conjecture

**Setup**:

- Completion of $100 \times 100$ rank 5 matrix
- Observed entries chosen uniformly at random

**Many (5K) Observations**:

|                    | *reconst err* | *nuclear norm* | *effective rank* |
| ------------------ | ------------- | -------------- | ---------------- |
| *nuclear norm min* | 8 e −07       | 221            | 5                |
| *depth 2 LNN*      | 5 e −06       | 221            | 5                |
| *depth 3 LNN*      | 4 e −06       | 221            | 5                |

- Nuclear norm min recovers ground truth
- LNN do so too
- Correspondence, but can't distinguish between nuclear norm min and any bias leading to low rank

# Experiments Testing Nuclear Norm Conjecture (cont')

**Few (2K) Observations**:

| | reconst err | nuclear norm | effective rank |
|---|---|---|---|
| nuclear norm min | | | |
| depth 2 LNN | | | |
| depth 3 LNN | | | |

# Experiments Testing Nuclear Norm Conjecture (cont')

**Few (2K) Observations**:

| | reconst err | nuclear norm | effective rank |
|---|---|---|---|
| **nuclear norm min** | 2 e −01 | 217 | 8 |
| **depth 2 LNN** | | | |
| **depth 3 LNN** | | | |

- Nuclear norm min doesn't recover ground truth

# Experiments Testing Nuclear Norm Conjecture (cont')

**Few (2K) Observations**:

|                    | *reconst err* | *nuclear norm* | *effective rank* |
|--------------------|:-------------:|:--------------:|:----------------:|
| *nuclear norm min* |    2 e −01    |      217       |        8         |
|   *depth* 2 *LNN*  |    6 e −02    |      220       |        6         |
|   *depth* 3 *LNN*  |    3 e −05    |      221       |        5         |

- Nuclear norm min doesn't recover ground truth
- LNN focus on lowering effective rank at expense of nuclear norm

# Experiments Testing Nuclear Norm Conjecture (cont')

**Few (2K) Observations**:

|                    | reconst err | nuclear norm | effective rank |
|--------------------|-------------|--------------|----------------|
| *nuclear norm min* | 2 e −01     | 217          | 8              |
| *depth* 2 *LNN*    | 6 e −02     | 220          | 6              |
| *depth* 3 *LNN*    | 3 e −05     | 221          | 5              |

- Nuclear norm min doesn't recover ground truth

- LNN focus on lowering effective rank at expense of nuclear norm

- Discrepancy!

# Experiments Testing Nuclear Norm Conjecture (cont')

**Few (2K) Observations**:

|  | *reconst err* | *nuclear norm* | *effective rank* |
|---|---|---|---|
| *nuclear norm min* | 2 e −01 | 217 | 8 |
| *depth 2 LNN* | 6 e −02 | 220 | 6 |
| *depth 3 LNN* | 3 e −05 | 221 | 5 |

- Nuclear norm min doesn't recover ground truth

- LNN focus on lowering effective rank at expense of nuclear norm

- Discrepancy!

**LNN implicitly minimize nuclear norm sometimes but not always!**

# Experiments Testing Nuclear Norm Conjecture (cont')

**Few (2K) Observations**:

|                  | reconst err | nuclear norm | effective rank |
|------------------|-------------|--------------|----------------|
| **nuclear norm min** | 2 e −01    | 217          | 8              |
| **depth 2 LNN**      | 6 e −02    | 220          | 6              |
| **depth 3 LNN**      | 3 e −05    | 221          | 5              |

- Nuclear norm min doesn't recover ground truth
- LNN focus on lowering effective rank at expense of nuclear norm
- Discrepancy!

**LNN implicitly minimize nuclear norm sometimes but not always!**

## Hypothesis

Single norm (or quasi-norm) not enough to capture implicit regularization, detailed account for trajectories is needed

# Trajectory Analysis $\longrightarrow$ Dynamics of Singular Values

## Trajectory Analysis $\longrightarrow$ Dynamics of Singular Values

Trajectory analysis gave dynamics for end-to-end matrix of $N$-layer LNN:

$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$

## Trajectory Analysis $\longrightarrow$ Dynamics of Singular Values

Trajectory analysis gave dynamics for end-to-end matrix of $N$-layer LNN:

$$\frac{d}{dt} vec \left[ W_{1:N}(t) \right] = -P_{W_{1:N}(t)} \cdot vec \left[ \nabla \ell (W_{1:N}(t)) \right]$$

Denote:

- $\{\sigma_r(t)\}_r$ — singular vals of $W_{1:N}(t)$
- $\{\mathbf{u}_r(t)\}_r / \{\mathbf{v}_r(t)\}_r$ — corresponding left/right singular vecs

# Trajectory Analysis $\longrightarrow$ Dynamics of Singular Values

Trajectory analysis gave dynamics for end-to-end matrix of $N$-layer LNN:

$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$

Denote:

- $\{\sigma_r(t)\}_r$ — singular vals of $W_{1:N}(t)$
- $\{\mathbf{u}_r(t)\}_r/\{\mathbf{v}_r(t)\}_r$ — corresponding left/right singular vecs

> ### Theorem
>
> $$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \left\langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t) \right\rangle$$

# Trajectory Analysis $\longrightarrow$ Dynamics of Singular Values

Trajectory analysis gave dynamics for end-to-end matrix of $N$-layer LNN:

$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$

Denote:

- $\{\sigma_r(t)\}_r$ — singular vals of $W_{1:N}(t)$
- $\{\mathbf{u}_r(t)\}_r / \{\mathbf{v}_r(t)\}_r$ — corresponding left/right singular vecs

> **Theorem**
>
> $$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \left\langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^{\top}(t) \right\rangle$$

**Interpretation**

- Given $W_{1:N}(t)$, depth affects evolution only via factors $N \cdot \sigma_r^{2-\frac{2}{N}}(t)$

# Trajectory Analysis $\longrightarrow$ Dynamics of Singular Values

Trajectory analysis gave dynamics for end-to-end matrix of $N$-layer LNN:

$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$

Denote:

- $\{\sigma_r(t)\}_r$ — singular vals of $W_{1:N}(t)$
- $\{\mathbf{u}_r(t)\}_r/\{\mathbf{v}_r(t)\}_r$ — corresponding left/right singular vecs

> ### Theorem
>
> $$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \left\langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t) \right\rangle$$

**Interpretation**

- Given $W_{1:N}(t)$, depth affects evolution only via factors $N \cdot \sigma_r^{2-\frac{2}{N}}(t)$
- $N = 1$ (classic linear model): factors reduce to 1

# Trajectory Analysis $\longrightarrow$ Dynamics of Singular Values

Trajectory analysis gave dynamics for end-to-end matrix of $N$-layer LNN:

$$\frac{d}{dt} vec\left[W_{1:N}(t)\right] = -P_{W_{1:N}(t)} \cdot vec\left[\nabla\ell(W_{1:N}(t))\right]$$

Denote:

- $\{\sigma_r(t)\}_r$ — singular vals of $W_{1:N}(t)$
- $\{\mathbf{u}_r(t)\}_r / \{\mathbf{v}_r(t)\}_r$ — corresponding left/right singular vecs

> **Theorem**
>
> $$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \left\langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t) \right\rangle$$

**Interpretation**

- Given $W_{1:N}(t)$, depth affects evolution only via factors $N \cdot \sigma_r^{2-\frac{2}{N}}(t)$
- $N = 1$ (classic linear model): factors reduce to 1
- $N \geq 2$: factors speed up (slow down) large (small) singular vals,
  more so for larger $N$ (higher depth)

# Dynamics of Singular Values — Proof Sketch

### Theorem

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$$

**Proof Sketch**

# Dynamics of Singular Values — Proof Sketch

### Theorem

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla \ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t) \rangle$$

**Proof Sketch**

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$    $\left( S = diag(\sigma_1, \sigma_2, ...) \quad U = [\mathbf{u}_1, \mathbf{u}_2, ...] \quad V = [\mathbf{v}_1, \mathbf{v}_2, ...] \right)$

# Dynamics of Singular Values — Proof Sketch

### Theorem

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla \ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t) \rangle$$

**Proof Sketch**

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$    $\left( S = diag(\sigma_1, \sigma_2, ...) \quad U = [\mathbf{u}_1, \mathbf{u}_2, ...] \quad V = [\mathbf{v}_1, \mathbf{v}_2, ...] \right)$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t) \cdot S(t) \cdot V(t)^\top + U(t) \cdot \frac{d}{dt}S(t) \cdot V(t)^\top + U(t) \cdot S(t) \cdot \frac{d}{dt}V(t)^\top$

# Dynamics of Singular Values — Proof Sketch

### Theorem

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$$

**Proof Sketch**

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$ $\quad \left(S = diag(\sigma_1, \sigma_2, ...) \quad U = [\mathbf{u}_1, \mathbf{u}_2, ...] \quad V = [\mathbf{v}_1, \mathbf{v}_2, ...]\right)$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t) \cdot S(t) \cdot V(t)^\top + U(t) \cdot \frac{d}{dt}S(t) \cdot V(t)^\top + U(t) \cdot S(t) \cdot \frac{d}{dt}V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t) \cdot V(t) = U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

# Dynamics of Singular Values — Proof Sketch

### Theorem

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla \ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t) \rangle$$

## Proof Sketch

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$  $\quad (S = diag(\sigma_1, \sigma_2, ...) \quad U = [\mathbf{u}_1, \mathbf{u}_2, ...] \quad V = [\mathbf{v}_1, \mathbf{v}_2, ...])$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t) \cdot S(t) \cdot V(t)^\top + U(t) \cdot \frac{d}{dt}S(t) \cdot V(t)^\top + U(t) \cdot S(t) \cdot \frac{d}{dt}V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t) \cdot V(t) = U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

End-to-end dynamics:

$\frac{d}{dt}W_{1:N}(t) = -\sum_{j=1}^{N} \left[ W_{1:N}(t)W_{1:N}(t)^\top \right]^{\frac{N-j}{N}} \cdot \nabla \ell(W_{1:N}(t)) \cdot \left[ W_{1:N}(t)^\top W_{1:N}(t) \right]^{\frac{j-1}{N}}$

# Dynamics of Singular Values — Proof Sketch

**Theorem**

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$$

## Proof Sketch

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$ $\quad$ ($S = diag(\sigma_1, \sigma_2, ...)$ $\quad$ $U = [\mathbf{u}_1, \mathbf{u}_2, ...]$ $\quad$ $V = [\mathbf{v}_1, \mathbf{v}_2, ...]$)

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t) \cdot S(t) \cdot V(t)^\top + U(t) \cdot \frac{d}{dt}S(t) \cdot V(t)^\top + U(t) \cdot S(t) \cdot \frac{d}{dt}V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t) \cdot V(t) = U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

End-to-end dynamics:

$\frac{d}{dt}W_{1:N}(t) = -\sum_{j=1}^{N}\left[W_{1:N}(t)W_{1:N}(t)^\top\right]^{\frac{N-j}{N}} \cdot \nabla\ell(W_{1:N}(t)) \cdot \left[W_{1:N}(t)^\top W_{1:N}(t)\right]^{\frac{j-1}{N}}$

# Dynamics of Singular Values — Proof Sketch

### Theorem

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t) \rangle$$

## Proof Sketch

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$ $\quad \left(S = diag(\sigma_1, \sigma_2, ...) \quad U = [\mathbf{u}_1, \mathbf{u}_2, ...] \quad V = [\mathbf{v}_1, \mathbf{v}_2, ...]\right)$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t) \cdot S(t) \cdot V(t)^\top + U(t) \cdot \frac{d}{dt}S(t) \cdot V(t)^\top + U(t) \cdot S(t) \cdot \frac{d}{dt}V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t) \cdot V(t) = U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

End-to-end dynamics:

$\frac{d}{dt}W_{1:N}(t) = -\sum_{j=1}^{N} U(t)\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}} V(t)^\top$

# Dynamics of Singular Values — Proof Sketch

### Theorem

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla \ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t) \rangle$$

## Proof Sketch

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$ $\quad \left(S = diag(\sigma_1, \sigma_2, ...) \quad U = [\mathbf{u}_1, \mathbf{u}_2, ...] \quad V = [\mathbf{v}_1, \mathbf{v}_2, ...]\right)$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t) \cdot S(t) \cdot V(t)^\top + U(t) \cdot \frac{d}{dt}S(t) \cdot V(t)^\top + U(t) \cdot S(t) \cdot \frac{d}{dt}V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t) \cdot V(t) = U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

End-to-end dynamics:

$\frac{d}{dt}W_{1:N}(t) = -\sum_{j=1}^N U(t)\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla \ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}} V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t) \cdot V(t)$

$\qquad = -\sum_{j=1}^N \left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla \ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}}$

# Dynamics of Singular Values — Proof Sketch

**Theorem**

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$$

## Proof Sketch

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$    $\left(S = diag(\sigma_1, \sigma_2, ...)  \quad U = [\mathbf{u}_1, \mathbf{u}_2, ...] \quad V = [\mathbf{v}_1, \mathbf{v}_2, ...]\right)$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t) \cdot S(t) \cdot V(t)^\top + U(t) \cdot \frac{d}{dt}S(t) \cdot V(t)^\top + U(t) \cdot S(t) \cdot \frac{d}{dt}V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t) \cdot V(t) = U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

End-to-end dynamics:

$\frac{d}{dt}W_{1:N}(t) = -\sum_{j=1}^N U(t)\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}} V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t) \cdot V(t)$

$\qquad = -\sum_{j=1}^N \left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}}$

# Dynamics of Singular Values — Proof Sketch

**Theorem**

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$$

**Proof Sketch**

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$    $\left(S = diag(\sigma_1, \sigma_2, ...) \quad U = [\mathbf{u}_1, \mathbf{u}_2, ...] \quad V = [\mathbf{v}_1, \mathbf{v}_2, ...]\right)$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t) \cdot S(t) \cdot V(t)^\top + U(t) \cdot \frac{d}{dt}S(t) \cdot V(t)^\top + U(t) \cdot S(t) \cdot \frac{d}{dt}V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t) \cdot V(t) = U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

End-to-end dynamics:

$\frac{d}{dt}W_{1:N}(t) = -\sum_{j=1}^{N} U(t)\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}} V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

$\qquad = -\sum_{j=1}^{N} \left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}}$

# Dynamics of Singular Values — Proof Sketch

**Theorem**

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$$

## Proof Sketch

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$   $\left(S = diag(\sigma_1, \sigma_2, ...) \quad U = [\mathbf{u}_1, \mathbf{u}_2, ...] \quad V = [\mathbf{v}_1, \mathbf{v}_2, ...]\right)$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t)\cdot S(t)\cdot V(t)^\top + U(t)\cdot\frac{d}{dt}S(t)\cdot V(t)^\top + U(t)\cdot S(t)\cdot\frac{d}{dt}V(t)^\top$

$\implies U(t)^\top\cdot\frac{d}{dt}W_{1:N}(t)\cdot V(t) = U(t)^\top\cdot\frac{d}{dt}U(t)\cdot S(t) + \frac{d}{dt}S(t) + S(t)\cdot\frac{d}{dt}V(t)^\top\cdot V(t)$

End-to-end dynamics:

$\frac{d}{dt}W_{1:N}(t) = -\sum_{j=1}^N U(t)\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}}U(t)^\top\cdot\nabla\ell(W_{1:N}(t))\cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}}V(t)^\top$

$\implies U(t)^\top\cdot\frac{d}{dt}U(t)\cdot S(t) + \frac{d}{dt}S(t) + S(t)\cdot\frac{d}{dt}V(t)^\top\cdot V(t)$

$\qquad = -\sum_{j=1}^N\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}}U(t)^\top\cdot\nabla\ell(W_{1:N}(t))\cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}}$

Restrict attention to $r$'th diagonal element:

# Dynamics of Singular Values — Proof Sketch

## Theorem

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$$

## **Proof Sketch**

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$    $\left(S = diag(\sigma_1, \sigma_2, ...)\quad U = [\mathbf{u}_1, \mathbf{u}_2, ...]\quad V = [\mathbf{v}_1, \mathbf{v}_2, ...]\right)$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t) \cdot S(t) \cdot V(t)^\top + U(t) \cdot \frac{d}{dt}S(t) \cdot V(t)^\top + U(t) \cdot S(t) \cdot \frac{d}{dt}V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t) \cdot V(t) = U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

End-to-end dynamics:

$\frac{d}{dt}W_{1:N}(t) = -\sum_{j=1}^{N} U(t)\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}} V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

$\qquad = -\sum_{j=1}^{N}\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}}$

Restrict attention to $r$'th diagonal element:

$\mathbf{u}_r(t)^\top \cdot \frac{d}{dt}\mathbf{u}_r(t) \cdot \sigma_r(t) + \frac{d}{dt}\sigma_r(t) + \sigma_r(t) \cdot \frac{d}{dt}\mathbf{v}_r(t)^\top \cdot \mathbf{v}_r(t) =$

$\qquad -\sum_{j=1}^{N} \sigma_r^{2\frac{N-j}{N}}(t) \cdot \mathbf{u}_r(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot \mathbf{v}_r(t) \cdot \sigma_r^{2\frac{j-1}{N}}(t)$

# Dynamics of Singular Values — Proof Sketch

### Theorem

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$$

## **Proof Sketch**

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$ $\quad \left(S = diag(\sigma_1, \sigma_2, ...) \quad U = [\mathbf{u}_1, \mathbf{u}_2, ...] \quad V = [\mathbf{v}_1, \mathbf{v}_2, ...]\right)$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t) \cdot S(t) \cdot V(t)^\top + U(t) \cdot \frac{d}{dt}S(t) \cdot V(t)^\top + U(t) \cdot S(t) \cdot \frac{d}{dt}V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t) \cdot V(t) = U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

End-to-end dynamics:

$\frac{d}{dt}W_{1:N}(t) = -\sum_{j=1}^{N} U(t)\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}} V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

$\qquad = -\sum_{j=1}^{N}\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}}$

Restrict attention to $r$'th diagonal element:

$\mathbf{u}_r(t)^\top \cdot \frac{d}{dt}\mathbf{u}_r(t) \cdot \sigma_r(t) + \frac{d}{dt}\sigma_r(t) + \sigma_r(t) \cdot \frac{d}{dt}\mathbf{v}_r(t)^\top \cdot \mathbf{v}_r(t) =$

$\qquad\qquad -\sum_{j=1}^{N}\sigma_r^{2\frac{N-1}{N}}(t) \cdot \mathbf{u}_r(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot \mathbf{v}_r(t)$

# Dynamics of Singular Values — Proof Sketch

### Theorem

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$$

**Proof Sketch**

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$     $\left(S = diag(\sigma_1, \sigma_2, ...) \quad U = [\mathbf{u}_1, \mathbf{u}_2, ...] \quad V = [\mathbf{v}_1, \mathbf{v}_2, ...]\right)$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t) \cdot S(t) \cdot V(t)^\top + U(t) \cdot \frac{d}{dt}S(t) \cdot V(t)^\top + U(t) \cdot S(t) \cdot \frac{d}{dt}V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t) \cdot V(t) = U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

End-to-end dynamics:

$\frac{d}{dt}W_{1:N}(t) = -\sum_{j=1}^N U(t)\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}}U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}}V(t)^\top$

$\implies \quad U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

$\qquad = -\sum_{j=1}^N \left[S(t)S(t)^\top\right]^{\frac{N-j}{N}}U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}}$

Restrict attention to $r$'th diagonal element:

$\mathbf{u}_r(t)^\top \cdot \frac{d}{dt}\mathbf{u}_r(t) \cdot \sigma_r(t) + \frac{d}{dt}\sigma_r(t) + \sigma_r(t) \cdot \frac{d}{dt}\mathbf{v}_r(t)^\top \cdot \mathbf{v}_r(t) =$

$\qquad\qquad -N \cdot \sigma_r^{2\frac{N-1}{N}}(t) \cdot \mathbf{u}_r(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot \mathbf{v}_r(t)$

# Dynamics of Singular Values — Proof Sketch

### Theorem

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$$

### **Proof Sketch**

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$ $\quad (S = diag(\sigma_1, \sigma_2, ...) \quad U = [\mathbf{u}_1, \mathbf{u}_2, ...] \quad V = [\mathbf{v}_1, \mathbf{v}_2, ...])$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t)\cdot S(t)\cdot V(t)^\top + U(t)\cdot \frac{d}{dt}S(t)\cdot V(t)^\top + U(t)\cdot S(t)\cdot \frac{d}{dt}V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t)\cdot V(t) = U(t)^\top \cdot \frac{d}{dt}U(t)\cdot S(t) + \frac{d}{dt}S(t) + S(t)\cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

End-to-end dynamics:

$\frac{d}{dt}W_{1:N}(t) = -\sum_{j=1}^{N} U(t)\Big[S(t)S(t)^\top\Big]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t))\cdot V(t)\Big[S(t)^\top S(t)\Big]^{\frac{j-1}{N}} V(t)^\top$

$\implies \quad U(t)^\top \cdot \frac{d}{dt}U(t)\cdot S(t) + \frac{d}{dt}S(t) + S(t)\cdot \frac{d}{dt}V(t)^\top \cdot V(t) =$

$\qquad\qquad = -\sum_{j=1}^{N}\Big[S(t)S(t)^\top\Big]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t))\cdot V(t)\Big[S(t)^\top S(t)\Big]^{\frac{j-1}{N}}$

Restrict attention to $r$'th diagonal element:

$\mathbf{u}_r(t)^\top \cdot \frac{d}{dt}\mathbf{u}_r(t)\cdot \sigma_r(t) + \frac{d}{dt}\sigma_r(t) + \sigma_r(t)\cdot \frac{d}{dt}\mathbf{v}_r(t)^\top \cdot \mathbf{v}_r(t) =$

$\qquad\qquad -N\cdot \sigma_r^{2\frac{N-1}{N}}(t)\cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$

# Dynamics of Singular Values — Proof Sketch

### Theorem

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$$

## Proof Sketch

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$ $\quad \left(S = diag(\sigma_1, \sigma_2, ...) \quad U = [\mathbf{u}_1, \mathbf{u}_2, ...] \quad V = [\mathbf{v}_1, \mathbf{v}_2, ...]\right)$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t) \cdot S(t) \cdot V(t)^\top + U(t) \cdot \frac{d}{dt}S(t) \cdot V(t)^\top + U(t) \cdot S(t) \cdot \frac{d}{dt}V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t) \cdot V(t) = U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

End-to-end dynamics:

$\frac{d}{dt}W_{1:N}(t) = -\sum_{j=1}^{N} U(t)\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}} V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

$\qquad = -\sum_{j=1}^{N}\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}}$

Restrict attention to $r$'th diagonal element:

$\mathbf{u}_r(t)^\top \cdot \frac{d}{dt}\mathbf{u}_r(t) \cdot \sigma_r(t) + \frac{d}{dt}\sigma_r(t) + \sigma_r(t) \cdot \frac{d}{dt}\mathbf{v}_r(t)^\top \cdot \mathbf{v}_r(t) =$

$\qquad -N \cdot \sigma_r^{2\frac{N-1}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$

# Dynamics of Singular Values — Proof Sketch

### Theorem

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$$

## Proof Sketch

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$  $\left(S = diag(\sigma_1, \sigma_2, ...) \quad U = [\mathbf{u}_1, \mathbf{u}_2, ...] \quad V = [\mathbf{v}_1, \mathbf{v}_2, ...]\right)$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t) \cdot S(t) \cdot V(t)^\top + U(t) \cdot \frac{d}{dt}S(t) \cdot V(t)^\top + U(t) \cdot S(t) \cdot \frac{d}{dt}V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t) \cdot V(t) = U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

End-to-end dynamics:

$\frac{d}{dt}W_{1:N}(t) = -\sum_{j=1}^{N} U(t)\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}} V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

$\qquad = -\sum_{j=1}^{N}\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}}$

Restrict attention to $r$'th diagonal element:

$\frac{1}{2}\frac{d}{dt}\|\mathbf{u}_r(t)\|_2^2 \cdot \sigma_r(t) + \frac{d}{dt}\sigma_r(t) + \sigma_r(t) \cdot \frac{1}{2}\frac{d}{dt}\|\mathbf{v}_r(t)\|_2^2 = -N \cdot \sigma_r^{2\frac{N-1}{N}}(t) \cdot \langle\nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$

# Dynamics of Singular Values — Proof Sketch

### Theorem

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$$

## **Proof Sketch**

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top \quad \left(S = diag(\sigma_1, \sigma_2, ...) \quad U = [\mathbf{u}_1, \mathbf{u}_2, ...] \quad V = [\mathbf{v}_1, \mathbf{v}_2, ...]\right)$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t)\cdot S(t)\cdot V(t)^\top + U(t)\cdot \frac{d}{dt}S(t)\cdot V(t)^\top + U(t)\cdot S(t)\cdot \frac{d}{dt}V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t)\cdot V(t) = U(t)^\top \cdot \frac{d}{dt}U(t)\cdot S(t) + \frac{d}{dt}S(t) + S(t)\cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

End-to-end dynamics:

$\frac{d}{dt}W_{1:N}(t) = -\sum_{j=1}^N U(t)\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t))\cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}} V(t)^\top$

$\implies \quad U(t)^\top \cdot \frac{d}{dt}U(t)\cdot S(t) + \frac{d}{dt}S(t) + S(t)\cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

$\qquad = -\sum_{j=1}^N \left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t))\cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}}$

Restrict attention to $r$'th diagonal element:

$\frac{1}{2}\frac{d}{dt}\underbrace{\|\mathbf{u}_r(t)\|_2^2}_{\equiv 1}\cdot\sigma_r(t) + \frac{d}{dt}\sigma_r(t) + \sigma_r(t)\cdot\frac{1}{2}\frac{d}{dt}\underbrace{\|\mathbf{v}_r(t)\|_2^2}_{\equiv 1} = -N\cdot\sigma_r^{2\frac{N-1}{N}}(t)\cdot\langle\nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$

# Dynamics of Singular Values — Proof Sketch

### Theorem

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t) \rangle$$

## Proof Sketch

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$    $\left( S = diag(\sigma_1, \sigma_2, ...) \quad U = [\mathbf{u}_1, \mathbf{u}_2, ...] \quad V = [\mathbf{v}_1, \mathbf{v}_2, ...] \right)$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t) \cdot S(t) \cdot V(t)^\top + U(t) \cdot \frac{d}{dt}S(t) \cdot V(t)^\top + U(t) \cdot S(t) \cdot \frac{d}{dt}V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t) \cdot V(t) = U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

End-to-end dynamics:

$\frac{d}{dt}W_{1:N}(t) = -\sum_{j=1}^{N} U(t)\left[ S(t)S(t)^\top \right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[ S(t)^\top S(t) \right]^{\frac{j-1}{N}} V(t)^\top$

$\implies \quad U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

$\qquad = -\sum_{j=1}^{N} \left[ S(t)S(t)^\top \right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[ S(t)^\top S(t) \right]^{\frac{j-1}{N}}$

Restrict attention to $r$'th diagonal element:

$0 \cdot \sigma_r(t) + \frac{d}{dt}\sigma_r(t) + \sigma_r(t) \cdot 0 = -N \cdot \sigma_r^{2\frac{N-1}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t) \rangle$

# Dynamics of Singular Values — Proof Sketch

### Theorem

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$$

## **Proof Sketch**

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$    $\left(S = diag(\sigma_1, \sigma_2, ...)\quad U = [\mathbf{u}_1, \mathbf{u}_2, ...]\quad V = [\mathbf{v}_1, \mathbf{v}_2, ...]\right)$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t) \cdot S(t) \cdot V(t)^\top + U(t) \cdot \frac{d}{dt}S(t) \cdot V(t)^\top + U(t) \cdot S(t) \cdot \frac{d}{dt}V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t) \cdot V(t) = U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

End-to-end dynamics:

$\frac{d}{dt}W_{1:N}(t) = -\sum_{j=1}^{N} U(t)\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}} V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

$\qquad = -\sum_{j=1}^{N}\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}}$

Restrict attention to $r$'th diagonal element:

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2\frac{N-1}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$$
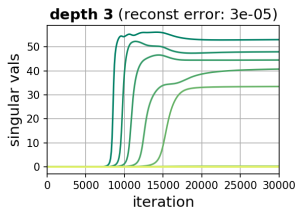
# Dynamics of Singular Values — Proof Sketch

> ### Theorem
> $$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2-\frac{2}{N}}(t) \cdot \langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\rangle$$
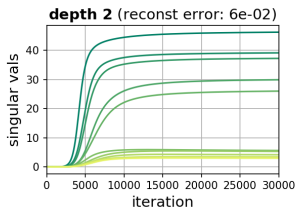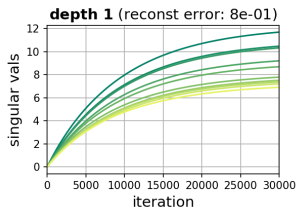
**Proof Sketch**

SVD: $W_{1:N}(t) = U(t)S(t)V(t)^\top$    $\left(S = diag(\sigma_1, \sigma_2, ...) \quad U = [\mathbf{u}_1, \mathbf{u}_2, ...] \quad V = [\mathbf{v}_1, \mathbf{v}_2, ...]\right)$

$\implies \frac{d}{dt}W_{1:N}(t) = \frac{d}{dt}U(t) \cdot S(t) \cdot V(t)^\top + U(t) \cdot \frac{d}{dt}S(t) \cdot V(t)^\top + U(t) \cdot S(t) \cdot \frac{d}{dt}V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}W_{1:N}(t) \cdot V(t) = U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

End-to-end dynamics:

$\frac{d}{dt}W_{1:N}(t) = -\sum_{j=1}^{N} U(t)\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}} V(t)^\top$

$\implies U(t)^\top \cdot \frac{d}{dt}U(t) \cdot S(t) + \frac{d}{dt}S(t) + S(t) \cdot \frac{d}{dt}V(t)^\top \cdot V(t)$

$\qquad = -\sum_{j=1}^{N}\left[S(t)S(t)^\top\right]^{\frac{N-j}{N}} U(t)^\top \cdot \nabla\ell(W_{1:N}(t)) \cdot V(t)\left[S(t)^\top S(t)\right]^{\frac{j-1}{N}}$

Restrict attention to $r$'th diagonal element:

$$\frac{d}{dt}\sigma_r(t) = -N \cdot \sigma_r^{2\frac{N-1}{N}}(t) \cdot \left\langle \nabla\ell(W_{1:N}(t)), \mathbf{u}_r(t)\mathbf{v}_r^\top(t)\right\rangle$$

$\square$

# Implicit Bias Towards Low Rank

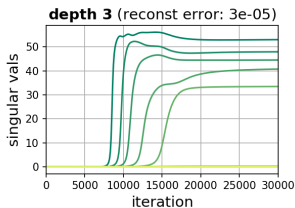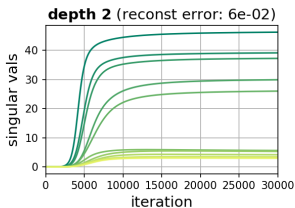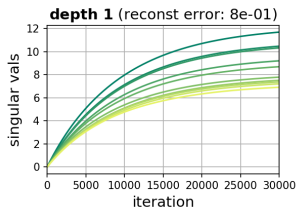# Implicit Bias Towards Low Rank

## **Experiment**

Completion of low rank matrix via GD over LNN

# Implicit Bias Towards Low Rank

## Experiment

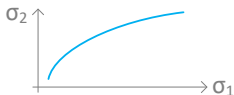Completion of low rank matrix via GD over LNN



## Theoretical Example

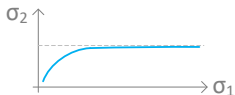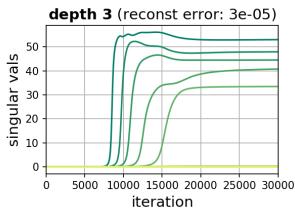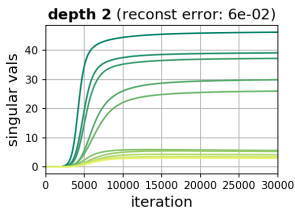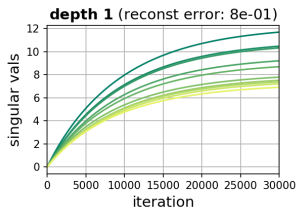For one observed entry and $\ell_2$ loss, relationship between singular vals is:

*depth 1: linear*          *depth 2: polynomial*          *depth $\geq 3$: asymptotic*

# Implicit Bias Towards Low Rank

## Experiment

Completion of low rank matrix via GD over LNN



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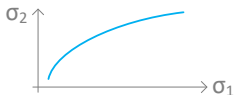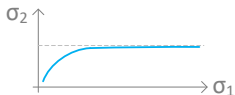
## Theoretical Example

For one observed entry and $\ell_2$ loss, relationship between singular vals is:

*depth 1: linear*          *depth 2: polynomial*          *depth $\geq 3$: asymptotic*



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

# Outline

1. Optimization and Generalization in Deep Learning via Trajectories

2. Case Study: Linear Neural Networks
   - Trajectory Analysis
   - Optimization
   - Generalization

3. **Conclusion**

# Recap

## Recap

**Perspective**

To understand optimization and generalization in deep learning:

## Recap

**Perspective**

To understand optimization and generalization in deep learning:

- Language of classical learning theory may be insufficient

## Recap

**Perspective**

To understand optimization and generalization in deep learning:

- Language of classical learning theory may be insufficient
- Might need to analyze trajectories of gradient descent

## Recap

**Perspective**

To understand optimization and generalization in deep learning:

- Language of classical learning theory may be insufficient
- Might need to analyze trajectories of gradient descent

**Case Study — Deep Linear Neural Networks**

## Recap

**Perspective**

To understand optimization and generalization in deep learning:

- Language of classical learning theory may be insufficient
- Might need to analyze trajectories of gradient descent

**Case Study — Deep Linear Neural Networks**

Trajectory analysis:

## Recap

**Perspective**

To understand optimization and generalization in deep learning:

- Language of classical learning theory may be insufficient
- Might need to analyze trajectories of gradient descent

**Case Study — Deep Linear Neural Networks**

Trajectory analysis:

- **Depth induces preconditioner** promoting movement in directions taken

## Recap

**Perspective**

To understand optimization and generalization in deep learning:

- Language of classical learning theory may be insufficient
- Might need to analyze trajectories of gradient descent

**Case Study — Deep Linear Neural Networks**

Trajectory analysis:

- **Depth induces preconditioner** promoting movement in directions taken

Optimization:

## Recap

**Perspective**

To understand optimization and generalization in deep learning:

- Language of classical learning theory may be insufficient
- Might need to analyze trajectories of gradient descent

**Case Study — Deep Linear Neural Networks**

Trajectory analysis:

- **Depth induces preconditioner** promoting movement in directions taken

Optimization:

- **Guarantee of efficient convergence to global min** (most general yet)

# Recap

**Perspective**

To understand optimization and generalization in deep learning:

- Language of classical learning theory may be insufficient
- Might need to analyze trajectories of gradient descent

**Case Study — Deep Linear Neural Networks**

Trajectory analysis:

- **Depth induces preconditioner** promoting movement in directions taken

Optimization:

- **Guarantee of efficient convergence to global min** (most general yet)
- **Depth can accelerate convergence** (w/o any gain in expressiveness)!

# Recap

**Perspective**

To understand optimization and generalization in deep learning:

- Language of classical learning theory may be insufficient
- Might need to analyze trajectories of gradient descent

**Case Study — Deep Linear Neural Networks**

Trajectory analysis:

- **Depth induces preconditioner** promoting movement in directions taken

Optimization:

- **Guarantee of efficient convergence to global min** (most general yet)
- **Depth can accelerate convergence** (w/o any gain in expressiveness)!

Generalization:

## Recap

**Perspective**

To understand optimization and generalization in deep learning:

- Language of classical learning theory may be insufficient
- Might need to analyze trajectories of gradient descent

**Case Study — Deep Linear Neural Networks**

Trajectory analysis:

- **Depth induces preconditioner** promoting movement in directions taken

Optimization:

- **Guarantee of efficient convergence to global min** (most general yet)
- **Depth can accelerate convergence** (w/o any gain in expressiveness)!

Generalization:

- **Depth enhances implicit regularization towards low rank**, yielding generalization for problems such as matrix completion

# Outline

1.

2.

3.

# Thank You