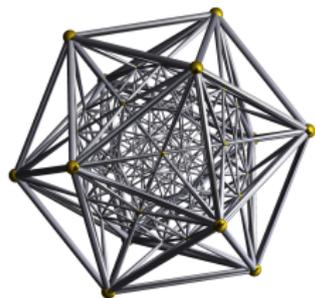


ICASSP22 Short Course One on Low-Dimensional Models for High-Dimensional Data

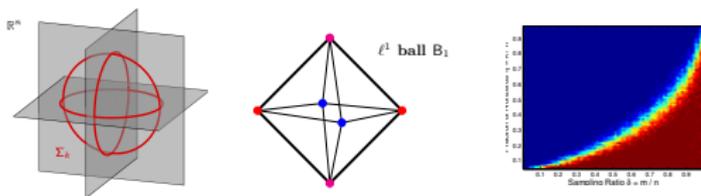
Lecture 4: Learning Low-Dimensional Structure via Deep Networks

**Sam Buchanan, Yi Ma, Qing Qu
John Wright, Yuqian Zhang, Zhihui Zhu**

May 26, 2022



Recap: Sparse Approximation (Linear, Convex)



Sparse approximation: **structured** signals, **linear** measurements

$$\mathbf{y} = \mathbf{A}\mathbf{x}_o, \quad \mathbf{x}_o \text{ sparse}, \quad \mathbf{A} \in \mathbb{R}^{m \times n} \text{ random}$$

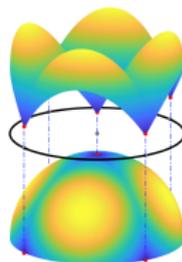
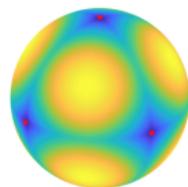
with **convex** optimization

$$\mathbf{x}_\star = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

and provable (high probability) guarantees

$$\mathbf{x}_\star = \mathbf{x}_o \text{ when } \text{measurements} \gtrsim \text{sparsity} \times \log \left(\frac{\text{measurements}}{\text{sparsity}} \right)$$

Recap: Dictionary Learning (Bilinear, Nonconvex)



Dictionary Learning: **structured** signals, **bilinear** measurements

$$Y = A_o X_o \in \mathbb{R}^{n \times p}, \quad X_o \text{ sparse and random}, \quad A_o^* A_o \approx I$$

with (efficient) **nonconvex** optimization

$$\mathbf{a}_\star = \arg \min_{\|\mathbf{a}\|_2=1} \|\mathbf{Y}^* \mathbf{a}\|_1$$

and provable (high probability) guarantees

$$\mathbf{a}_\star \approx (A_o)_j \text{ when } \text{observations} \geq \text{poly}(\text{expected sparsity})$$

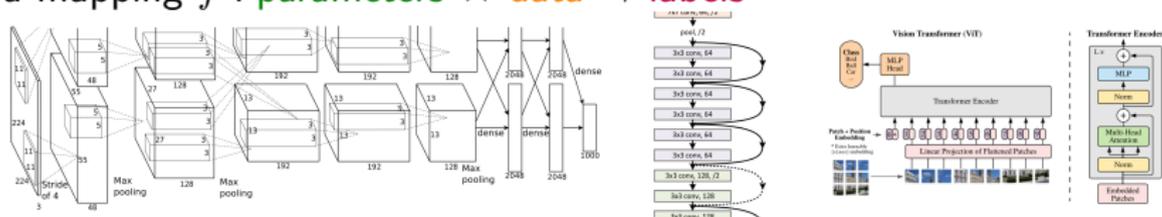
Today: Deep Learning (Very Nonlinear, Extra Nonconvex)

Supervised Deep Learning: Given labeled data

$$\left\{ \left(\underbrace{\mathbf{x}_i}_{\text{data (images, text, ...)}} , \underbrace{\mathbf{y}_i}_{\text{labels (classes, values, ...)}} \right) \right\}_{i=1}^N$$



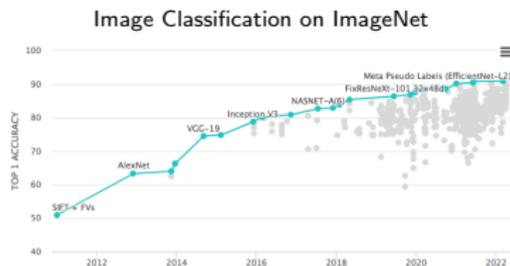
fit a mapping f : **parameters** \times **data** \rightarrow **labels**



using stochastic gradient descent on a task-appropriate loss

$$\underbrace{\theta_{\star} = \text{SGD}_{\theta} \left(\frac{1}{2} \sum_{i=1}^N \|f_{\theta}(\mathbf{x}_i) - \mathbf{y}_i\|_2^2 \right)}_{\text{regression}}$$

Today's Lectures



Big question: What role does **low-dimensional structure** play in the **practice** of deep learning?

TEXT PROMPT

an illustration of a baby daikon radish in a tutu walking a dog

AI-GENERATED IMAGES



TEXT DESCRIPTION

An astronaut **Teddy bears** A bowl of soup

mixing sparkling chemicals as mad scientists **shopping for groceries** working on new AI research

as kids' crayon art on the moon in the 1980s **underwater with 1990s technology**

DALL E 2



Outlook for Today's Lectures

Answer: A huge role!

Today:

- **Nonlinear** low-dimensional structures in practical data necessitate the use of **deep networks** over classical models;

Outlook for Today's Lectures

Answer: A huge role!

Today:

- **Nonlinear** low-dimensional structures in practical data necessitate the use of **deep networks** over classical models;
- **A mathematical model problem** helps understand **resource tradeoffs** between data geometry and network architecture (a **nonlinear generalization** of the sparse approximation analysis!);

Outlook for Today's Lectures

Answer: A huge role!

Today:

- **Nonlinear** low-dimensional structures in practical data necessitate the use of **deep networks** over classical models;
- **A mathematical model problem** helps understand **resource tradeoffs** between data geometry and network architecture (a **nonlinear generalization** of the sparse approximation analysis!);
- **For classification problems**, understand the features learned by deep neural networks, and improve training robustness using insights from low-dimensional structure;

Outlook for Today's Lectures

Answer: A huge role!

Today:

- **Nonlinear** low-dimensional structures in practical data necessitate the use of **deep networks** over classical models;
- **A mathematical model problem** helps understand **resource tradeoffs** between data geometry and network architecture (a **nonlinear generalization** of the sparse approximation analysis!);
- **For classification problems**, understand the features learned by deep neural networks, and improve training robustness using insights from low-dimensional structure;
- **Whitebox design of deep networks** for pursuing nonlinear low-dim structures. (Lecture 5)

Outline

Recap and Outlook

- 1 Motivating Examples for Low-Dim Structure in Deep Learning
- 2 Resource Tradeoffs in the Multiple Manifold Problem
 - Problem Formulation
 - Intrinsic Geometric Properties of Manifold Data
 - Network Architecture Resources and Training Procedure
 - Training Deep Networks with Gradient Descent
 - Resource Tradeoffs
- 3 Looking Inside: Neural Collapse in the Multiple Manifold Problem
 - Learned low-dimensional features—NC phenomena
 - Geometric analysis for understanding neural collapse
 - Exploit NC for improving training efficiency
 - Exploit NC for understanding the effect of loss functions
- 4 Exploit Sparse Model for Robust training

Low-Dimensional Structure in Deep Learning Problems



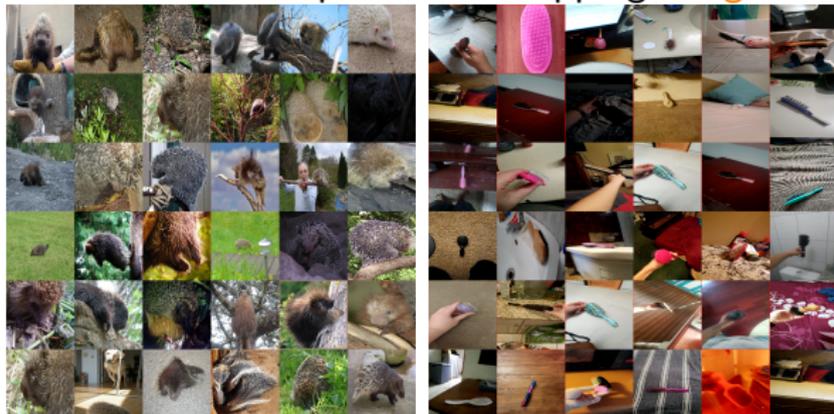
Questions:

What is an appropriate mathematical model for data with low-dimensional structure in deep learning applications?

What insights into practical deep learning can we get by studying low-dimensional structure?

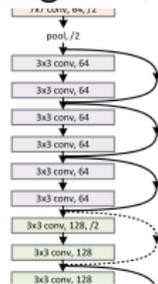
Vignette I: Large-Scale Image Classification

Task: Learn a deep network mapping **images** \rightarrow **object classes** from data.

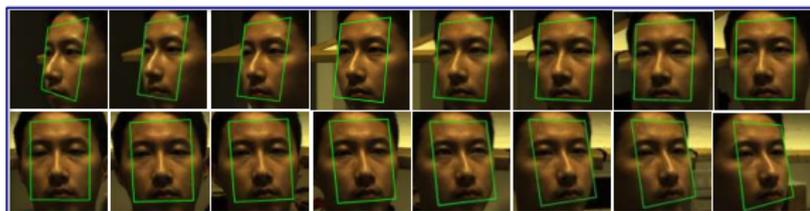


\rightarrow {hedgehog,
hairbrush}

Massive driver of innovation in the last 10 years (ImageNet, ResNet, ...)



Nonlinear Variabilities in Natural Images

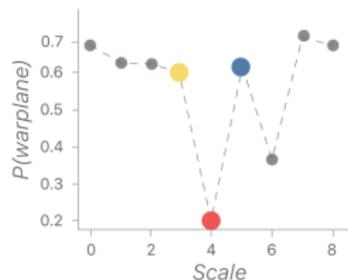
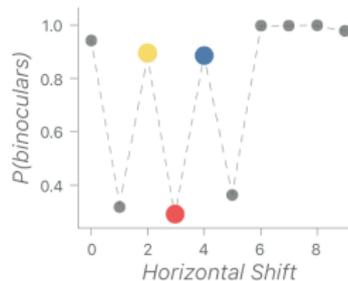
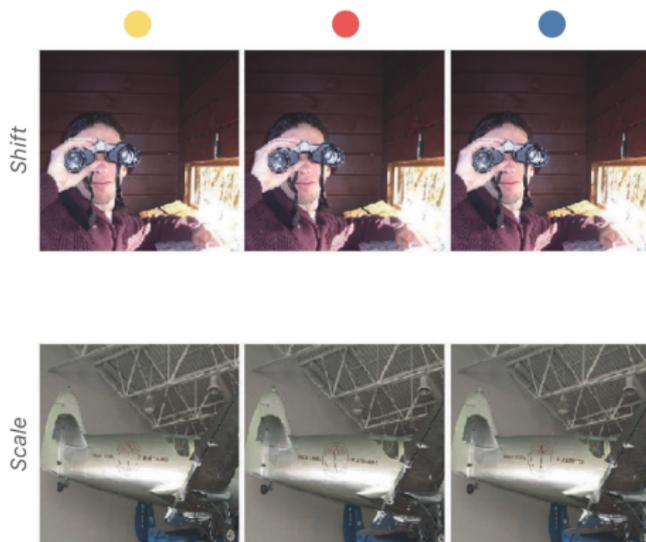


⇒ **nonlinear, geometric** structure

- 6D for 3D rigid pose; 8D for perspective; 9D for certain illumination...

Limitations of a Purely Data-Driven Approach?

Can fail to learn even simple invariances in the data:

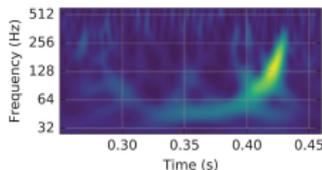
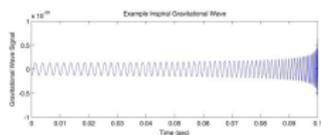
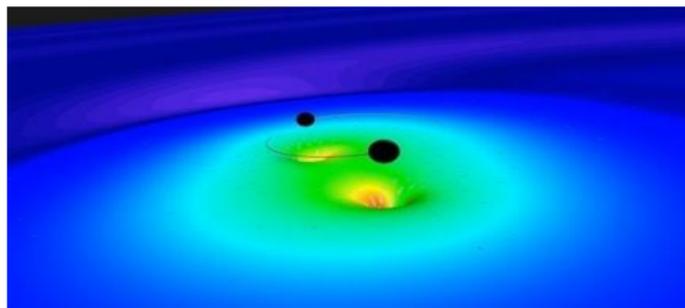


From [Azulay and Weiss, 2019]

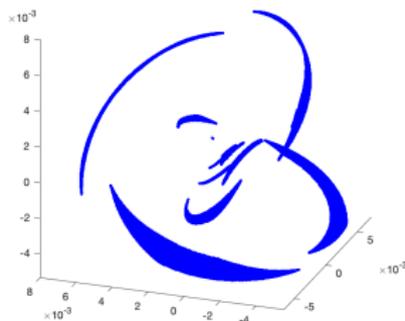
Vignette II: Deep Learning in Scientific Discovery

Gravitational Wave Astronomy

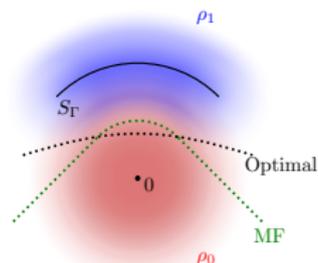
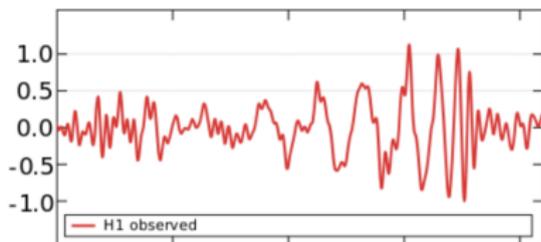
One binary black hole merger:



Many mergers
(varying mass M_1, M_2):
 \Rightarrow **low-dim manifold**



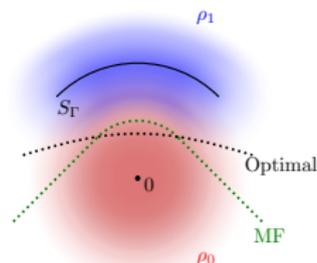
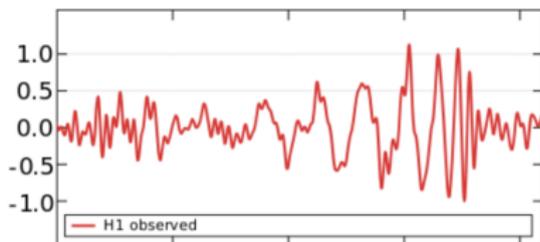
Gravitational Wave Astronomy as Parametric Detection



Is observation $x = s_\gamma + z$ or $x = z$?

\implies **two (noisy) manifolds!**

Gravitational Wave Astronomy as Parametric Detection

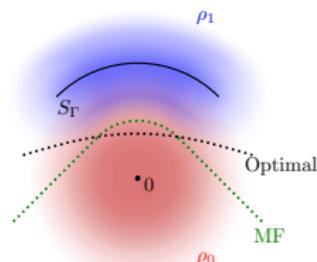
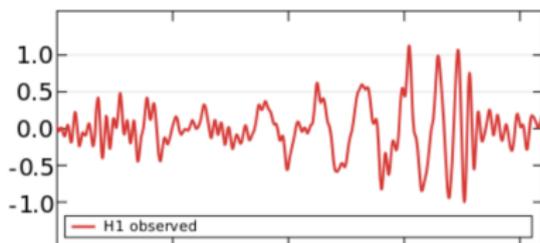


Is observation $x = s_\gamma + z$ or $x = z$?

\implies **two (noisy) manifolds!**

Classical approach: template matching $\max_\gamma \langle a_\gamma, x \rangle > \tau$

Gravitational Wave Astronomy as Parametric Detection



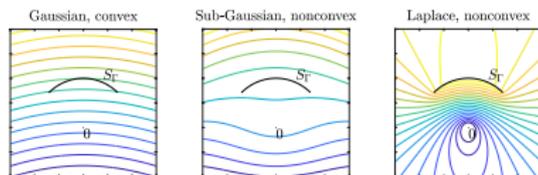
Is observation $x = s_\gamma + z$ or $x = z$?

\implies **two (noisy) manifolds!**

Classical approach: template matching $\max_\gamma \langle a_\gamma, x \rangle > \tau$

Issues: Optimality? Complexity?

Unknown unknowns? Unknown noise?

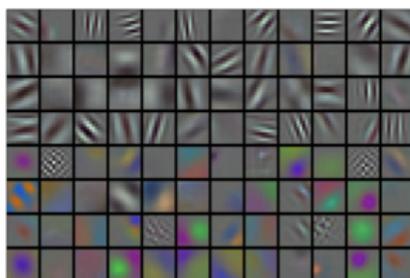


Ideally: Combine low-dim structure of Γ with data-driven for statistical structure...

Vignette III: Learning Features with Deep Learning for Downstream Tasks

Ubiquitous deep learning workflow (science/engineering/industry):

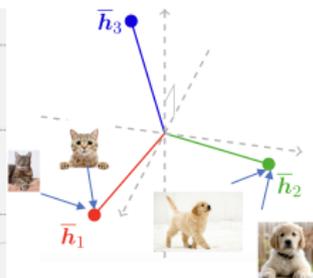
- 1 Data-driven pretraining to learn good features (ImageNet pretraining; masked prediction)
- 2 Fine-tuning on specific downstream tasks for performance (tracking; segmentation; object detection; ...)



First-layer filters from AlexNet



Inception v4 activations



Neural collapse visualization

Issues: What features are learned? Robustness to imperfect labeling? How to incorporate prior knowledge about data/task?

Takeaways from the Examples

Two key takeaways:

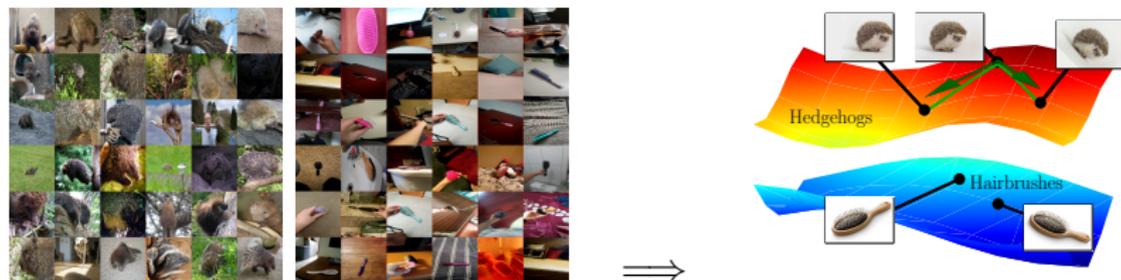
- Data with **nonlinear, geometric structure** pervade successful practical applications of deep learning
- Important practical issues (**robustness/invariance; resource efficiency; performance**) naturally linked to low-dim structure

Takeaways from the Examples

Two key takeaways:

- Data with **nonlinear, geometric structure** pervade successful practical applications of deep learning
- Important practical issues (**robustness/invariance; resource efficiency; performance**) naturally linked to low-dim structure

Next: Understanding mathematically when and why deep learning successfully classifies data with nonlinear geometric structure.



Outline

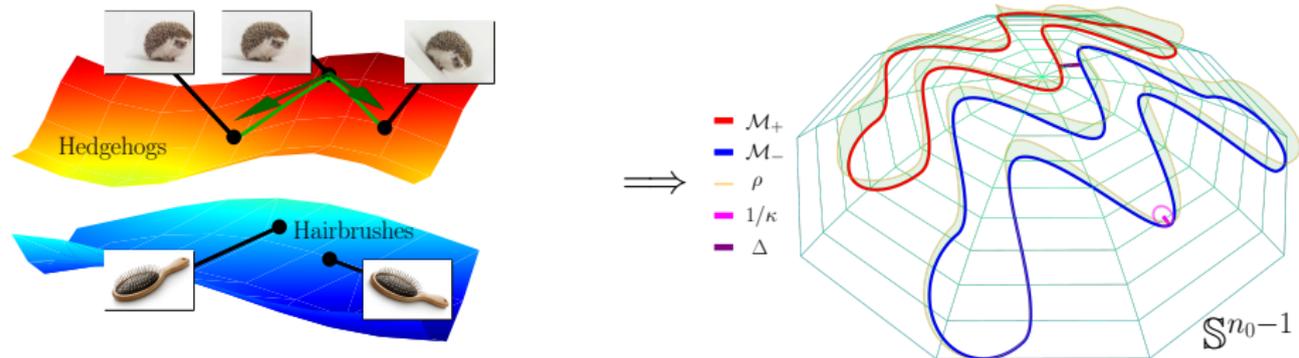
Recap and Outlook

- ① Motivating Examples for Low-Dim Structure in Deep Learning
- ② Resource Tradeoffs in the Multiple Manifold Problem
 - Problem Formulation
 - Intrinsic Geometric Properties of Manifold Data
 - Network Architecture Resources and Training Procedure
 - Training Deep Networks with Gradient Descent
 - Resource Tradeoffs
- ③ Looking Inside: Neural Collapse in the Multiple Manifold Problem
 - Learned low-dimensional features—NC phenomena
 - Geometric analysis for understanding neural collapse
 - Exploit NC for improving training efficiency
 - Exploit NC for understanding the effect of loss functions
- ④ Exploit Sparse Model for Robust training

A Mathematical Model Problem for Deep Learning + Low-Dimensional Structure

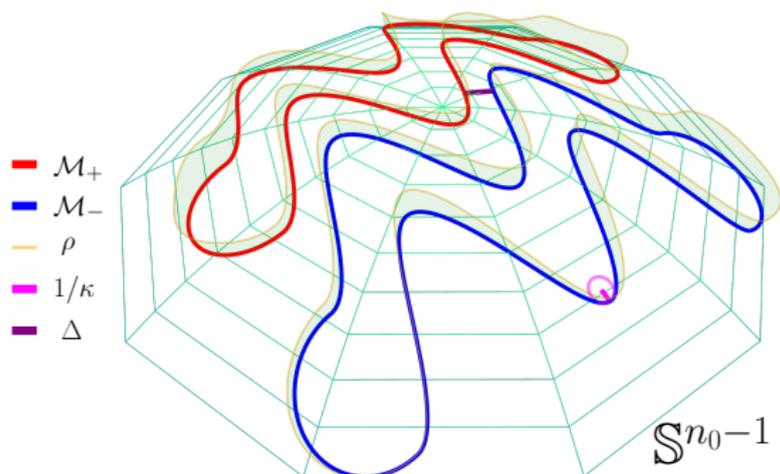
Formalizing data with nonlinear geometric structure:

Low-dimensional **Riemannian submanifolds** of high-dimensional space!



The multiple manifold problem: K -way classification of data on d -dimensional Riemannian manifolds in \mathbb{S}^{n_0-1} .

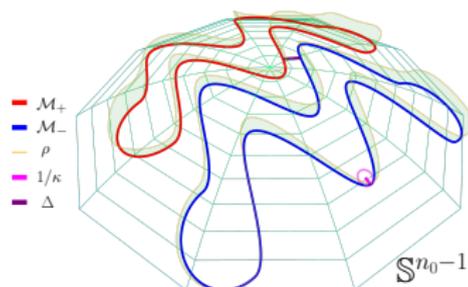
The Two Manifold Problem



Problem. Given N i.i.d. labeled samples $(\mathbf{x}_1, y(\mathbf{x}_1)), \dots, (\mathbf{x}_N, y(\mathbf{x}_N))$ from $\mathcal{M} = \mathcal{M}_+ \cup \mathcal{M}_-$, use gradient descent to train a deep network f_θ that *perfectly labels the manifolds*:

$$\text{sign}(f_\theta(\mathbf{x})) = y(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \mathcal{M}.$$

The Two Manifold Problem: Key Aspects



Problem. Given N i.i.d. labeled samples $(\mathbf{x}_1, y(\mathbf{x}_1)), \dots, (\mathbf{x}_N, y(\mathbf{x}_N))$ from $\mathcal{M} = \mathcal{M}_+ \cup \mathcal{M}_-$, use gradient descent to train a deep network f_θ that *perfectly labels the manifolds*:

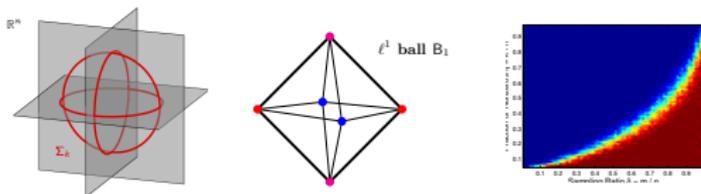
$$\text{sign}(f_\theta(\mathbf{x})) = y(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{M}.$$

- Binary classification with a deep neural network
- High-dimensional data with (unknown!) low-dimensional structure
- Statistical structure, and asking for “strong” generalization

We will focus on the case of one-dimensional manifolds (curves)

What Can We Hope to Understand Here?

Our “barometer”: compressed sensing.



$$\mathbf{y} = \mathbf{A}\mathbf{x}_o; \quad \mathbf{x}_\star = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

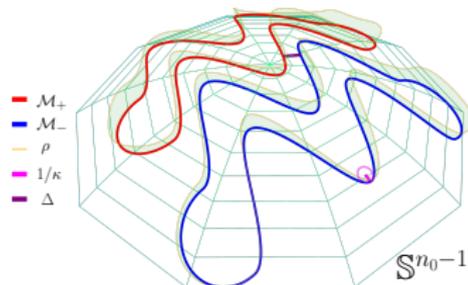
$$\mathbf{x}_\star = \mathbf{x}_o \text{ when } \text{measurements} \gtrsim \text{sparsity} \times \log \left(\frac{\text{measurements}}{\text{sparsity}} \right)$$

Questions:

What are our ‘measurement resources’ in the two manifold problem?

What are intrinsic structural properties of nonlinear manifold data?

The Two Manifold Problem: Geometric Parameters



Problem. Given N i.i.d. labeled samples $(\mathbf{x}_1, y(\mathbf{x}_1)), \dots, (\mathbf{x}_N, y(\mathbf{x}_N))$ from $\mathcal{M} = \mathcal{M}_+ \cup \mathcal{M}_-$, use gradient descent to train a deep network f_θ that *perfectly labels the manifolds*:

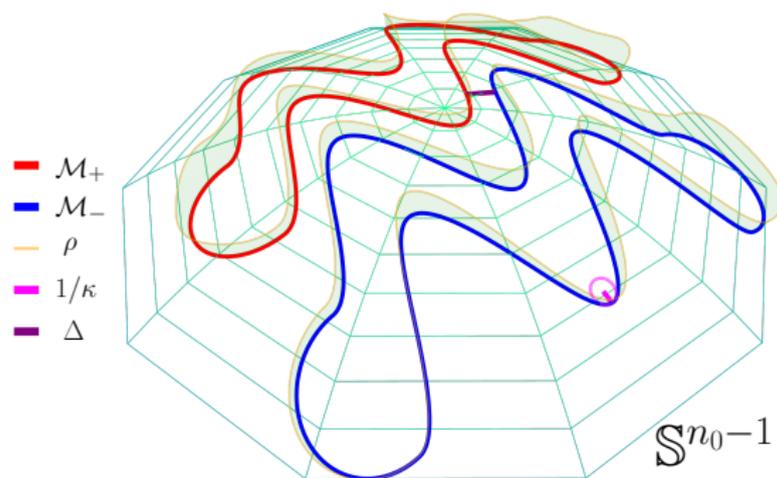
$$\text{sign}(f_\theta(\mathbf{x})) = y(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{M}.$$

A set of ‘sufficient’ intrinsic problem difficulty parameters:

- Curvature κ ;
- Separation Δ ;
- Separation ‘frequency’ ✿ .

Intrinsic Structural Properties I: Separation

Intuitively: How close are the class manifolds?

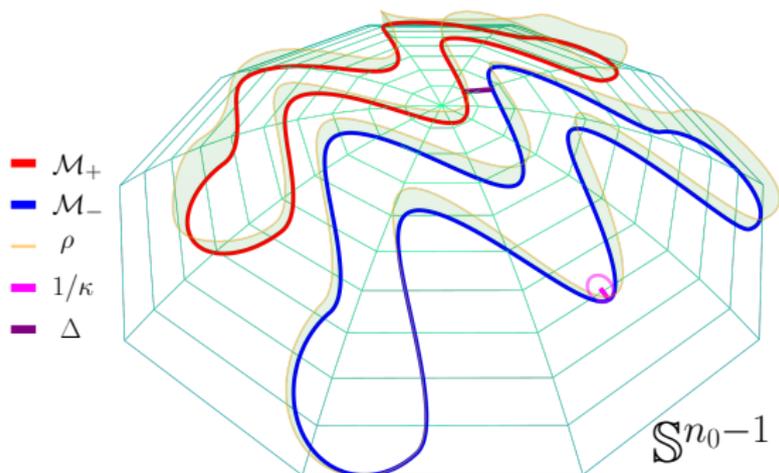


Mathematically:

$$\Delta = \inf_{\mathbf{x}, \mathbf{x}' \in \mathcal{M}} \{d_{\text{extrinsic}}(\mathbf{x}, \mathbf{x}')\}$$

Intrinsic Structural Properties II: Curvature

Intuitively: Local deviation from *flatness* of the manifold.

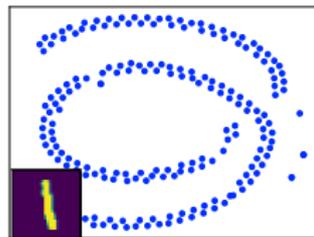
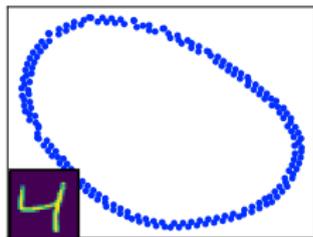
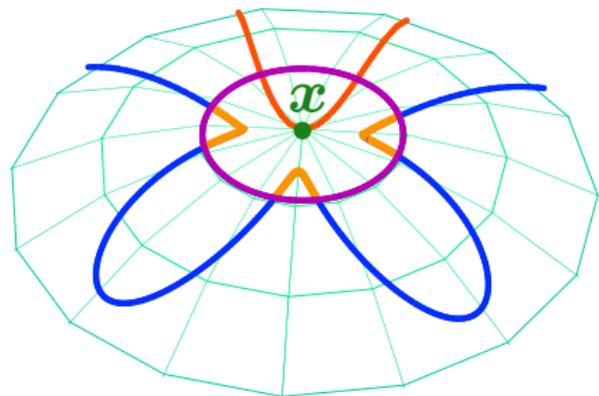


Mathematically:

$$\kappa = \sup_{\mathbf{x} \in \mathcal{M}} \left\| \left(\mathbf{I} - \frac{\mathbf{x}\mathbf{x}^*}{\|\mathbf{x}\|_2^2} \right) \ddot{\mathbf{x}} \right\|_2$$

Intrinsic Structural Properties III: ✿ -Number

Intuitively: How much do the class manifolds loop back on themselves?

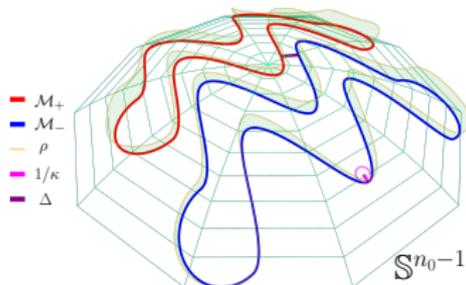


Mathematically:

$$\text{✿}(\mathcal{M}) = \sup_{\mathbf{x} \in \mathcal{M}} N_{\mathcal{M}} \left(\left\{ \mathbf{x}' \mid \begin{array}{l} d_{\text{intrinsic}}(\mathbf{x}, \mathbf{x}') > \tau_1 \\ d_{\text{extrinsic}}(\mathbf{x}, \mathbf{x}') < \tau_2 \end{array} \right\}, \frac{1}{\sqrt{1 + \kappa^2}} \right)$$

Here, $N_{\mathcal{M}}(T, \delta)$ is the covering number of $T \subseteq \mathcal{M}$ by δ balls in $d_{\text{intrinsic}}$.

The Two Manifold Problem: Geometric Parameters



Problem. Given N i.i.d. labeled samples $(\mathbf{x}_1, y(\mathbf{x}_1)), \dots, (\mathbf{x}_N, y(\mathbf{x}_N))$ from $\mathcal{M} = \mathcal{M}_+ \cup \mathcal{M}_-$, use gradient descent to train a deep network f_θ that *perfectly labels the manifolds*:

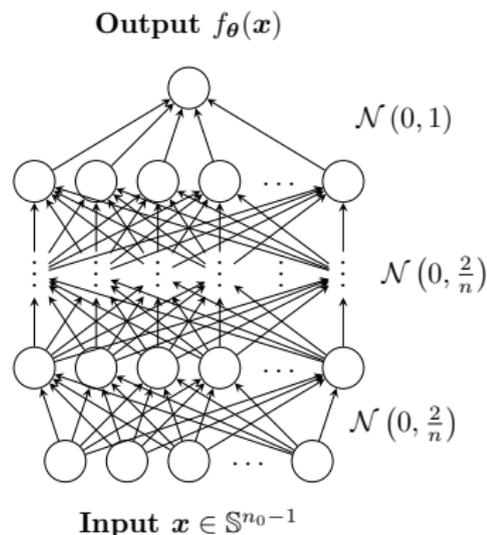
$$\text{sign}(f_\theta(\mathbf{x})) = y(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{M}.$$

A set of ‘sufficient’ intrinsic problem difficulty parameters:

- Curvature κ ;
- Separation Δ ;
- Separation ‘frequency’ $\text{\textcircled{+}}$.

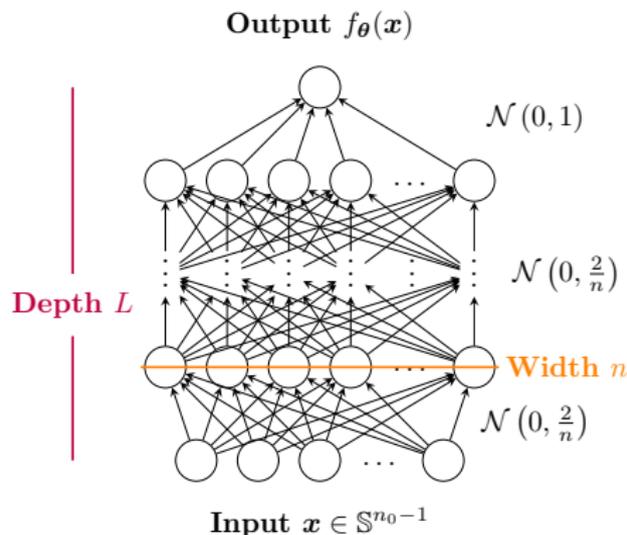
Network Architecture and Training Procedure

- Fully connected with ReLUs
- Gaussian initialization θ_0
- Trained with N i.i.d. samples from measure μ of density ρ



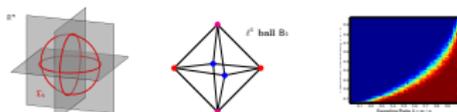
Network Architecture and Training Procedure

- Fully connected with ReLUs
- Gaussian initialization θ_0
- Trained with N i.i.d. samples from measure μ of density ρ



Resource Tradeoffs: From Linear to Nonlinear

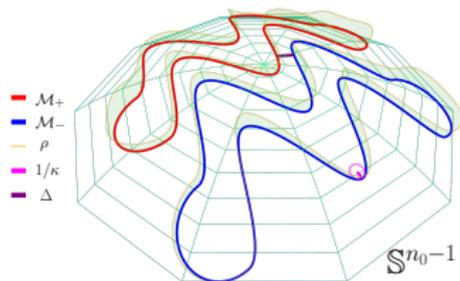
The “linear” case (compressed sensing):



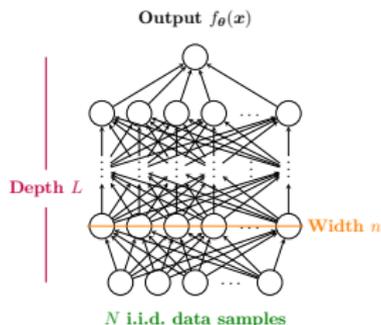
$$\mathbf{y} = \mathbf{A}\mathbf{x}_o; \quad \mathbf{x}_\star = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

$$\mathbf{x}_\star = \mathbf{x}_o \text{ when } \text{measurements} \gtrsim \text{sparsity} \times \log \left(\frac{\text{measurements}}{\text{sparsity}} \right)$$

Our current **nonlinear** setting:

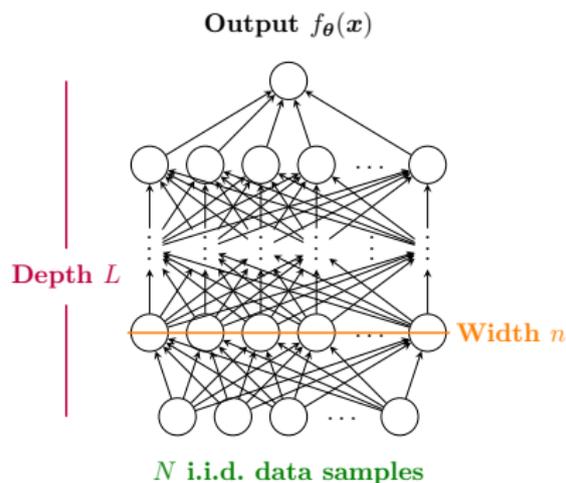


Data structure



Architectural resources ▶

The Two Manifold Problem: Resource Tradeoffs



Theory question: How should we set resources (depth L , width n , samples N) relative to data structure (separation Δ , $\color{purple}{\otimes}$; curvature κ ; density ρ) so that *gradient descent succeeds*?

Gradient Descent Training

Objective: Square Loss on Training Data

$$\min_{\boldsymbol{\theta}} \varphi(\boldsymbol{\theta}) \equiv \frac{1}{2} \int_{\mathcal{M}} (f_{\boldsymbol{\theta}}(\mathbf{x}) - y(\mathbf{x}))^2 d\mu_N(\mathbf{x}).$$

Does gradient descent correctly label the manifolds?

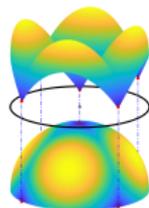
Gradient Descent Training

Objective: Square Loss on Training Data

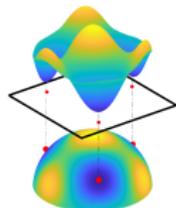
$$\min_{\theta} \varphi(\theta) \equiv \frac{1}{2} \int_{\mathcal{M}} (f_{\theta}(\mathbf{x}) - y(\mathbf{x}))^2 d\mu_N(\mathbf{x}).$$

Does gradient descent correctly label the manifolds?

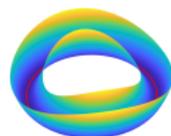
One Approach: Geometry (from symmetry!) in **parameter space**:



Dictionary Learning



Sparse Blind Deconvolution



Matrix Recovery

See [Gilboa, B., Wright '18], survey [Zhang, Qu, Wright 20] (Lecture 3!)

Gradient Descent Training

Objective: Square Loss on Training Data

$$\min_{\theta} \varphi(\theta) \equiv \frac{1}{2} \int_{\mathcal{M}} (f_{\theta}(\mathbf{x}) - y(\mathbf{x}))^2 d\mu_N(\mathbf{x}).$$

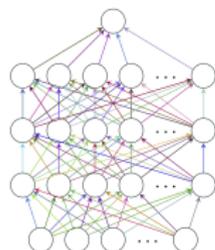
Does gradient descent correctly label the manifolds?

Today's talk: Dynamics in **input-output space**:

Neural Tangent Kernel

$$\Theta(\mathbf{x}, \mathbf{x}') = \left\langle \frac{\partial f_{\theta}(\mathbf{x})}{\partial \theta}, \frac{\partial f_{\theta}(\mathbf{x}')}{\partial \theta} \right\rangle$$

Measures ease of independently adjusting $f_{\theta}(\mathbf{x})$, $f_{\theta}(\mathbf{x}')$



Follows [Jacot et. al. 18], many recent works.

Dynamics of Gradient Descent

Objective: Square Loss on Training Data

$$\min_{\boldsymbol{\theta}} \varphi(\boldsymbol{\theta}) \equiv \frac{1}{2} \int_{\mathcal{M}} (f_{\boldsymbol{\theta}}(\mathbf{x}) - y(\mathbf{x}))^2 d\mu_N(\mathbf{x}).$$

Signed error: $\zeta(\mathbf{x}) = f_{\boldsymbol{\theta}}(\mathbf{x}) - y(\mathbf{x})$.

Gradient flow: $\dot{\boldsymbol{\theta}}_t = -\nabla_{\boldsymbol{\theta}} \varphi(\boldsymbol{\theta}_t) = -\int_{\mathcal{M}} \frac{\partial f_{\boldsymbol{\theta}}}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t}(\mathbf{x}) \zeta_t(\mathbf{x}) d\mu_N(\mathbf{x})$.

Dynamics of Gradient Descent

The error evolves according to the NTK:

$$\dot{\zeta}_t(\mathbf{x}) = \left. \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x})}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t}^* \dot{\boldsymbol{\theta}}_t$$

Dynamics of Gradient Descent

The error evolves according to the NTK:

$$\begin{aligned}\dot{\zeta}_t(\mathbf{x}) &= \left. \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x})}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t}^* \dot{\boldsymbol{\theta}}_t \\ &= - \left. \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x})}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t}^* \int_{\mathcal{M}} \left. \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x}')}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t} \zeta_t(\mathbf{x}') d\mu_N(\mathbf{x}')\end{aligned}$$

Dynamics of Gradient Descent

The error evolves according to the NTK:

$$\begin{aligned}
 \dot{\zeta}_t(\mathbf{x}) &= \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t}^* \dot{\boldsymbol{\theta}}_t \\
 &= - \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t}^* \int_{\mathcal{M}} \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x}')}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t} \zeta_t(\mathbf{x}') d\mu_N(\mathbf{x}') \\
 &= - \int_{\mathcal{M}} \left\langle \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t}, \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x}')}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t} \right\rangle \zeta_t(\mathbf{x}') d\mu_N(\mathbf{x}')
 \end{aligned}$$

Dynamics of Gradient Descent

The error evolves according to the NTK:

$$\begin{aligned}
 \dot{\zeta}_t(\mathbf{x}) &= \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t}^* \dot{\boldsymbol{\theta}}_t \\
 &= - \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t}^* \int_{\mathcal{M}} \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x}')}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t} \zeta_t(\mathbf{x}') d\mu_N(\mathbf{x}') \\
 &= - \int_{\mathcal{M}} \left\langle \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t}, \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x}')}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t} \right\rangle \zeta_t(\mathbf{x}') d\mu_N(\mathbf{x}')
 \end{aligned}$$

Dynamics of Gradient Descent

The error evolves according to the NTK:

$$\begin{aligned}
 \dot{\zeta}_t(\mathbf{x}) &= \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t}^* \dot{\boldsymbol{\theta}}_t \\
 &= - \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t}^* \int_{\mathcal{M}} \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x}')}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t} \zeta_t(\mathbf{x}') d\mu_N(\mathbf{x}') \\
 &= - \int_{\mathcal{M}} \left\langle \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t}, \frac{\partial f_{\boldsymbol{\theta}}(\mathbf{x}')}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t} \right\rangle \zeta_t(\mathbf{x}') d\mu_N(\mathbf{x}') \\
 &= - \int_{\mathcal{M}} \Theta_t(\mathbf{x}, \mathbf{x}') \zeta_t(\mathbf{x}') d\mu_N(\mathbf{x}') \\
 &= - \Theta_t[\zeta_t](\mathbf{x}).
 \end{aligned}$$

Dynamics of Gradient Descent (“NTK Regime”)

When **width** and **number of data samples** are large, we have (whp)

$$\sup_t \|\Theta_t - \Theta\|_{L^2 \rightarrow L^2} = o_{\text{width}}(1)$$

throughout training.

\implies *LTI dynamics*

$$\dot{\zeta}_t = -\Theta[\zeta_t]$$

\implies **Fast decay** if ζ_t is aligned with lead eigenvectors of Θ !

Implicit Error-NTK Alignment with Certificates

Challenge: For nonlinear \mathcal{M} , eigenvectors of Θ are intractable!

Definition. $g : \mathcal{M} \rightarrow \mathbb{R}$ is called a *certificate* if for all $\mathbf{x} \in \mathcal{M}$

$$f_{\theta_0}(\mathbf{x}) - y(\mathbf{x}) \underset{\text{square}}{\overset{\text{mean}}{\approx}} \int_{\mathcal{M}} \Theta(\mathbf{x}, \mathbf{x}') g(\mathbf{x}') d\mu(\mathbf{x}')$$

and $\int_{\mathcal{M}} (g(\mathbf{x}'))^2 d\mu(\mathbf{x}')$ is small.

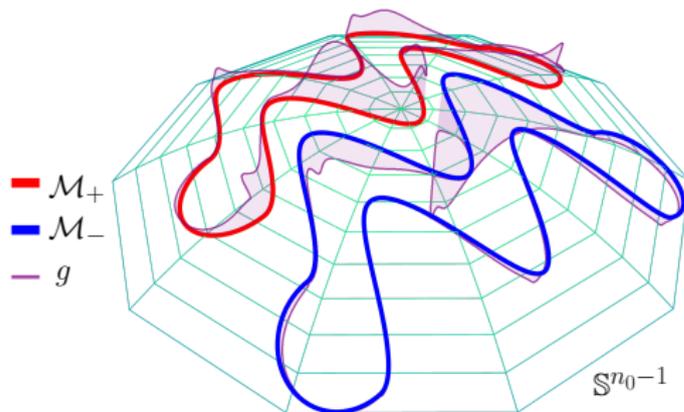
Implicit Error-NTK Alignment with Certificates

Challenge: For nonlinear \mathcal{M} , eigenvectors of Θ are intractable!

Definition. $g : \mathcal{M} \rightarrow \mathbb{R}$ is called a *certificate* if for all $\mathbf{x} \in \mathcal{M}$

$$f_{\theta_0}(\mathbf{x}) - y(\mathbf{x}) \underset{\text{square}}{\overset{\text{mean}}{\approx}} \int_{\mathcal{M}} \Theta(\mathbf{x}, \mathbf{x}') g(\mathbf{x}') d\mu(\mathbf{x}')$$

and $\int_{\mathcal{M}} (g(\mathbf{x}'))^2 d\mu(\mathbf{x}')$ is small.



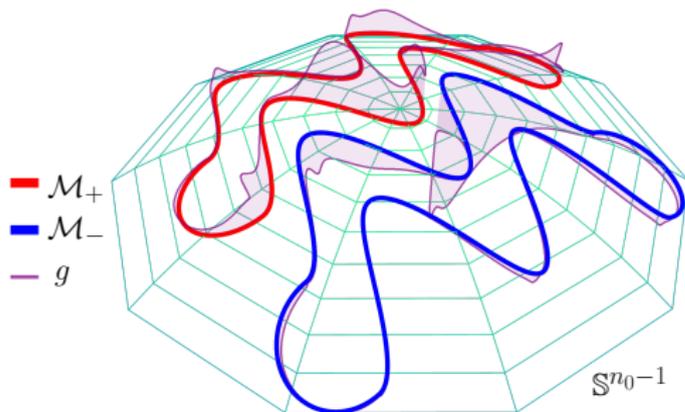
Implicit Error-NTK Alignment with Certificates

Challenge: For nonlinear \mathcal{M} , eigenvectors of Θ are intractable!

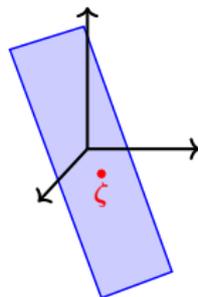
Definition. $g : \mathcal{M} \rightarrow \mathbb{R}$ is called a *certificate* if for all $\mathbf{x} \in \mathcal{M}$

$$f_{\theta_0}(\mathbf{x}) - y(\mathbf{x}) \underset{\text{square}}{\overset{\text{mean}}{\approx}} \int_{\mathcal{M}} \Theta(\mathbf{x}, \mathbf{x}') g(\mathbf{x}') d\mu(\mathbf{x}')$$

and $\int_{\mathcal{M}} (g(\mathbf{x}'))^2 d\mu(\mathbf{x}')$ is small.



Function space $L^2_{\mu_N}$



Error ζ near **stable range**
of random operator Θ

Implicit Error-NTK Alignment with Certificates

Challenge: For nonlinear \mathcal{M} , eigenvectors of Θ are intractable!

Definition. $g : \mathcal{M} \rightarrow \mathbb{R}$ is called a *certificate* if for all $\mathbf{x} \in \mathcal{M}$

$$f_{\theta_0}(\mathbf{x}) - y(\mathbf{x}) \underset{\text{square}}{\overset{\text{mean}}{\approx}} \int_{\mathcal{M}} \Theta(\mathbf{x}, \mathbf{x}') g(\mathbf{x}') d\mu(\mathbf{x}')$$

and $\int_{\mathcal{M}} (g(\mathbf{x}'))^2 d\mu(\mathbf{x}')$ is small.

Lemma. (informal) If a certificate g exists for \mathcal{M} , then

$$\|\zeta_t\|_{L^2_{\mu}} \lesssim \frac{L \log L}{t}.$$

Roles of Width, Depth, and Data

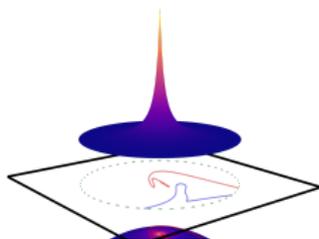
$$\dot{\zeta}_t = -\Theta[\zeta_t]$$

Questions:

How do **width**, **depth**, and **samples** affect Θ ?

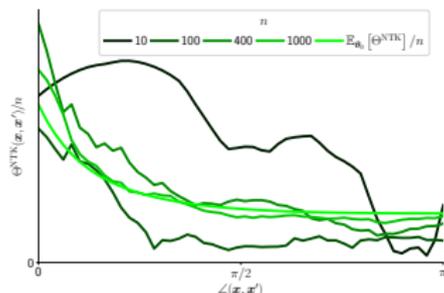
How does Θ depend on the geometry of the data?

Depth L : **fitting resource**



$$\frac{1}{L} \Theta(e_1, \mathbf{x}'), L = 125$$

Width n : **statistical resource**

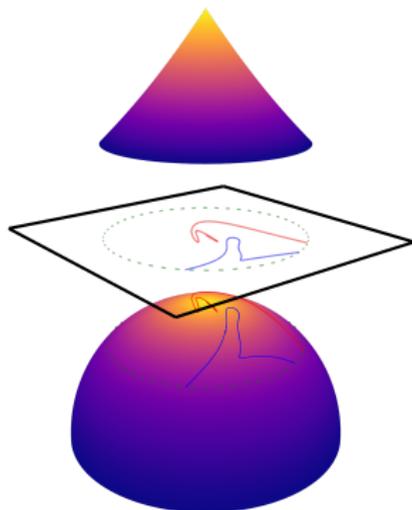


Resource Tradeoffs I: Depth as a Fitting Resource

Key insights:

- 1 Θ decays with angle.
- 2 Faster decay as depth increases.

\implies Set depth based on geometry!



$$\frac{1}{L} \Theta(\mathbf{e}_1, \mathbf{x}'), L = 5$$

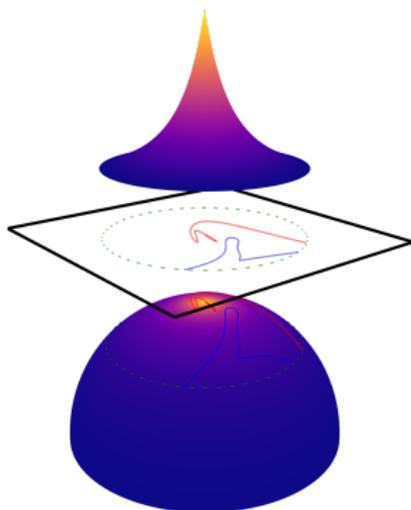
Deeper networks fit more complicated geometries.

Resource Tradeoffs I: Depth as a Fitting Resource

Key insights:

- 1 Θ decays with angle.
- 2 Faster decay as depth increases.

\implies Set depth based on geometry!



$$\frac{1}{L} \Theta(\mathbf{e}_1, \mathbf{x}'), L = 25$$

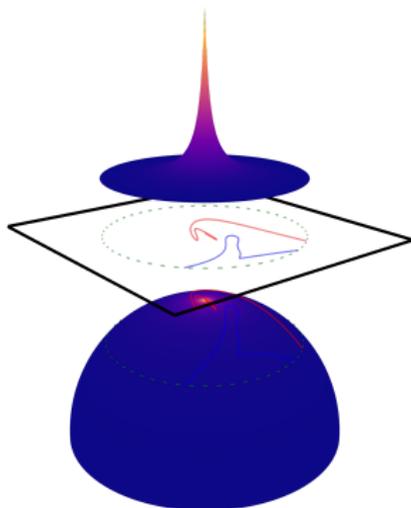
Deeper networks fit more complicated geometries.

Resource Tradeoffs I: Depth as a Fitting Resource

Key insights:

- 1 Θ decays with angle.
- 2 Faster decay as depth increases.

\implies Set depth based on geometry!



$$\frac{1}{L} \Theta(\mathbf{e}_1, \mathbf{x}'), L = 125$$

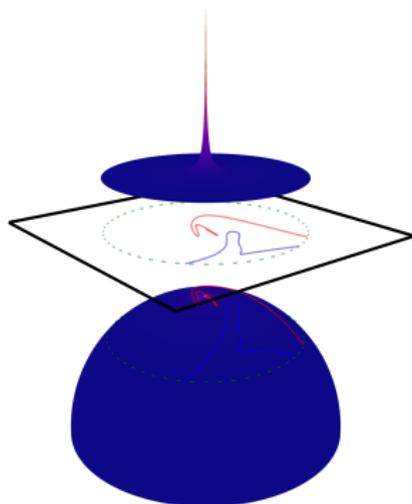
Deeper networks fit more complicated geometries.

Resource Tradeoffs I: Depth as a Fitting Resource

Key insights:

- 1 Θ decays with angle.
- 2 Faster decay as depth increases.

\implies Set depth based on geometry!

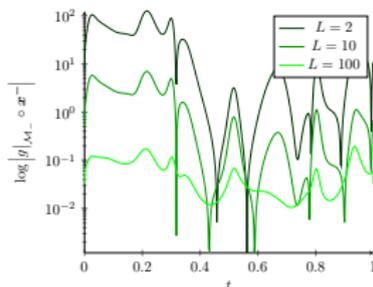
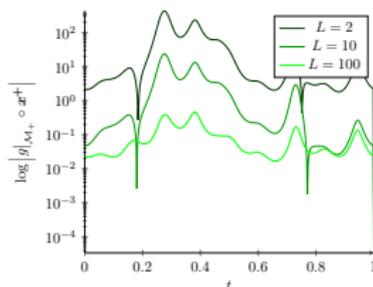
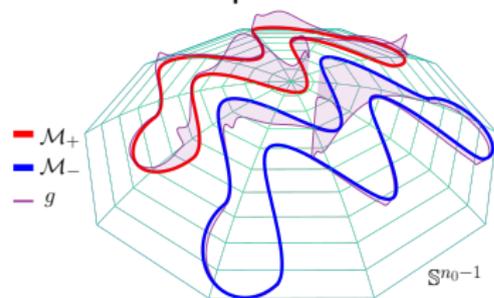


$$\frac{1}{L} \Theta(e_1, \mathbf{x}'), L = 625$$

Deeper networks fit more complicated geometries.

Resource Tradeoffs I: Certificates from Depth

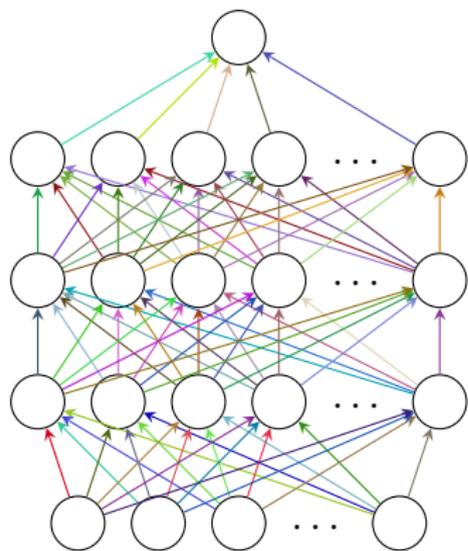
Numerical experiment:



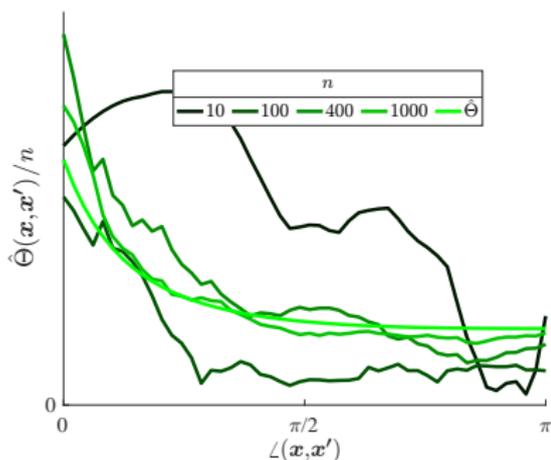
Depth as a fitting resource: Larger **depth** L leads to a sharper kernel Θ and a smaller certificate g
 \implies Easier fitting!

Resource Tradeoffs II: Width as a Statistical Resource

Output $f_{\theta}(x)$



Input $x \in \mathbb{S}^{n_0-1}$



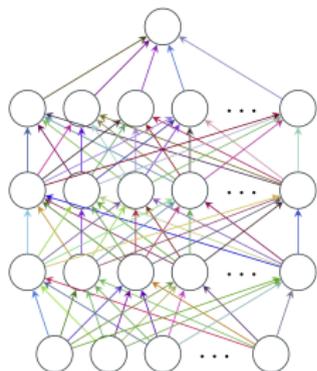
As **width** increases, $\Theta(x, x')$ concentrates about $\mathbb{E}_{\text{init weights}}[\Theta(x, x')]$

Resource Tradeoffs II: Width as a Statistical Resource

Proposition. Suppose that $n > L \text{polylog}(Ln_0)$. Then (whp)

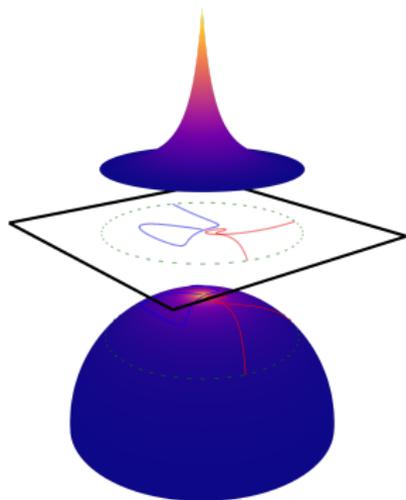
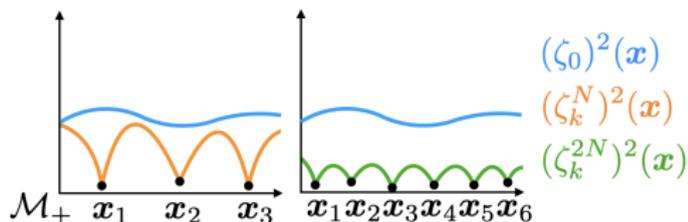
$$\left| \Theta(\mathbf{x}, \mathbf{x}') - \frac{n}{2} \sum_{\ell} \cos(\varphi^{\ell} \nu) \prod_{\ell'=\ell}^{L-1} \left(1 - \frac{\varphi^{\ell'} \nu}{\pi} \right) \right|$$

is small (simultaneously) for all $(\mathbf{x}, \mathbf{x}') \in \mathcal{M} \times \mathcal{M}$.



\Rightarrow **set width n based on depth L
and implicitly based on κ, Δ**

Resource Tradeoffs III: Data as a Statistical Resource



Depth $L = 50$

\Rightarrow **Sample complexity N is dictated by kernel “aperture”, which depends on geometry (κ, Δ) via L**

End-to-End Generalization Guarantee

Theorem [B., Wang, Gilboa, Wright 2021]: For sufficiently regular one-dimensional manifolds and ReLU networks, when

$$\text{depth} \geq \text{geometry}, \text{width} \geq \text{poly}(\text{depth}), \text{data} \geq \text{poly}(\text{depth}),$$

randomly-initialized small-stepping gradient descent perfectly classifies the two manifolds!

Upshot:

- We understand the role each resource plays in solving the classification problem.
- We understand how intrinsic geometric properties of the data drive these resource requirements.

Outline

Recap and Outlook

- ① Motivating Examples for Low-Dim Structure in Deep Learning
- ② Resource Tradeoffs in the Multiple Manifold Problem
 - Problem Formulation
 - Intrinsic Geometric Properties of Manifold Data
 - Network Architecture Resources and Training Procedure
 - Training Deep Networks with Gradient Descent
 - Resource Tradeoffs
- ③ Looking Inside: Neural Collapse in the Multiple Manifold Problem
 - Learned low-dimensional features—NC phenomena
 - Geometric analysis for understanding neural collapse
 - Exploit NC for improving training efficiency
 - Exploit NC for understanding the effect of loss functions
- ④ Exploit Sparse Model for Robust training

Image Classification Problem I

So far, Sam has talked about resources needed to ensure correctly classify two manifolds.

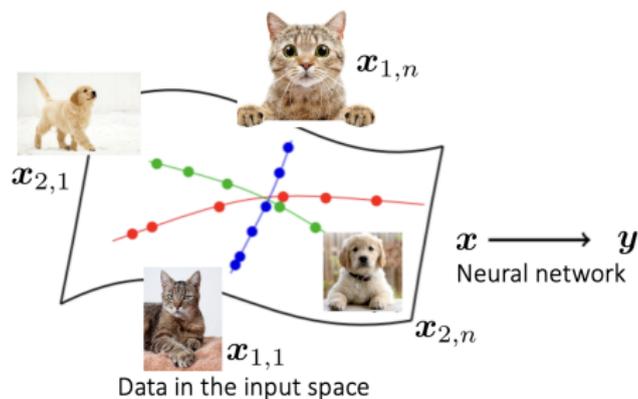
We will now focus on the general classification of K manifolds.

Instead of just on the output, we will focus more on the learned features and classifiers.

Image Classification Problem II

Labels: $k = 1, \dots, K$

- $K = 10$ classes (MNIST, CIFAR10, etc)
- $K = 1000$ classes (ImageNet)



$$\begin{array}{c} \text{Cat} \\ \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \end{array}
 \quad
 \begin{array}{c} \text{Dog} \\ \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \end{array}
 \quad
 \dots
 \quad
 \begin{array}{c} \text{Truck} \\ \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix} \end{array}$$

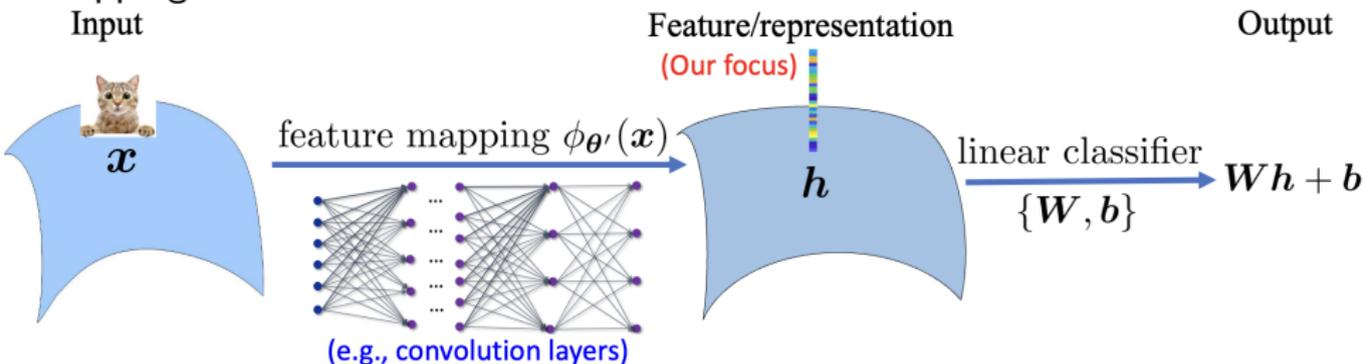
One-hot labeling vectors in \mathbb{R}^K

Assume balanced dataset where each class has n training samples

- If not, we can use data augmentation to make them balanced

Deep Neural Network Classifiers I

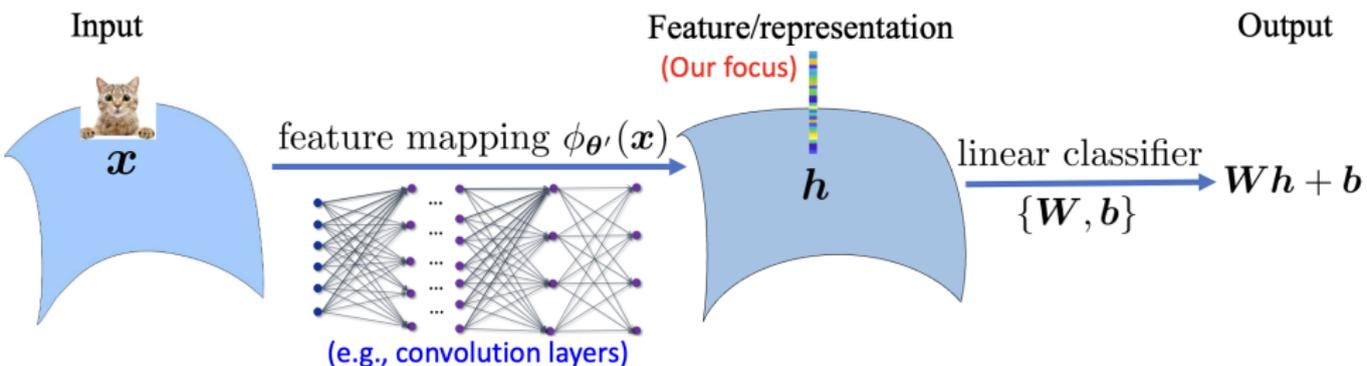
A deep neural network classifier often contains two parts: a feature mapping and a linear classifier



- Output: $f(\mathbf{x}; \theta) = W\phi_{\theta'}(\mathbf{x}) + \mathbf{b}$ with $\theta = (\theta', W, b)$.
- Training problem:

$$\min_{\theta', W, b} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \underbrace{\mathcal{L}_{\text{CE}}(W\phi_{\theta'}(\mathbf{x}_{k,i}) + \mathbf{b}, \mathbf{y}_k)}_{\text{cross-entropy (CE) loss}} + \lambda \underbrace{\|(\theta', W, b)\|_F^2}_{\text{weight decay}}$$

Deep Neural Network Classifiers II



Output: $f(\mathbf{x}; \theta) = \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} \xrightarrow{\text{Softmax function}} \begin{bmatrix} 0.6 \\ 0.3 \\ 0.1 \end{bmatrix}$

Cat $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$

Dog

Panda

Target

Prediction (probability)

$\text{CE}(\text{Cat}) := -q(\text{Cat}) \cdot \log p(\text{Cat})$

$= -1 \cdot \log 0.6$

$= 0.51\dots$

Neural Collapse in Classification I

Prevalence of neural collapse during the terminal phase of deep learning training

 Vardan Papyan,  X. Y. Han, and David L. Donoho

[+ See all authors and affiliations](#)

PNAS October 6, 2020 117 (40) 24652-24663; first published September 21, 2020;
<https://doi.org/10.1073/pnas.2015509117>

Contributed by David L. Donoho, August 18, 2020 (sent for review July 22, 2020; reviewed by Helmut Boelsckei and Stéphane Mallat)

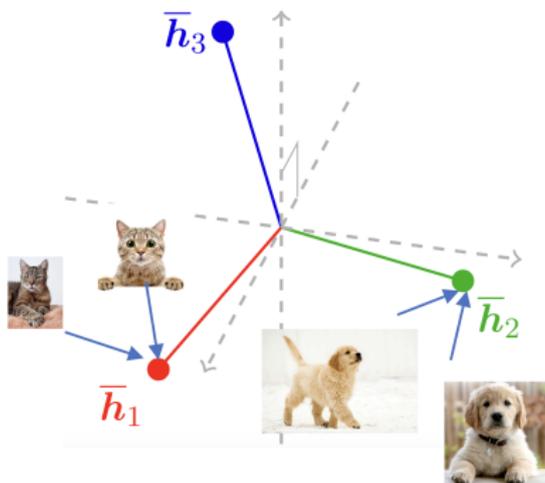
- Reveals common outcome of learned features and classifiers across a variety of architectures and dataset
- Precise mathematical structure within the features and classifier

Neural Collapse in Classification II

Neural Collapse (NC) refers to

- **NC1: Within-Class Variability Collapse:** features of each class collapse to class-mean with zero variability (*low-dimensional features*):

k -th class, i -th sample : $\mathbf{h}_{k,i} \rightarrow \bar{\mathbf{h}}_k$,

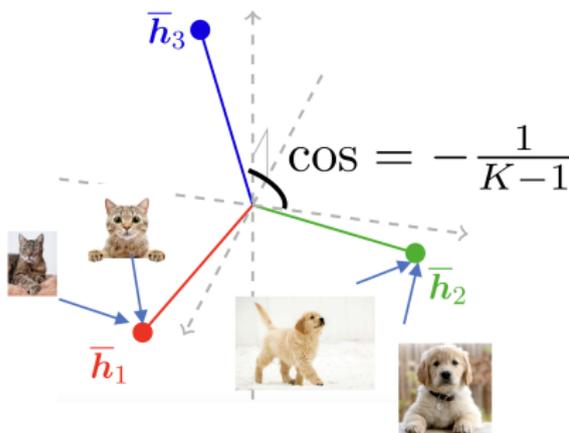


Neural Collapse in Classification III

Neural Collapse (NC) refers to

- **NC2: Convergence to Simplex Equiangular Tight Frame (ETF):** the class means are linearly separable, have same length, and maximal angle between each other

$$\frac{\langle \bar{\mathbf{h}}_k, \bar{\mathbf{h}}_{k'} \rangle}{\|\bar{\mathbf{h}}_k\| \|\bar{\mathbf{h}}_{k'}\|} \rightarrow \begin{cases} 1, & k = k' \\ -\frac{1}{K-1}, & k \neq k' \end{cases}$$



- If K vectors have equal angle between each other, then the largest possible cosine angle between each pair is $-\frac{1}{K-1}$.

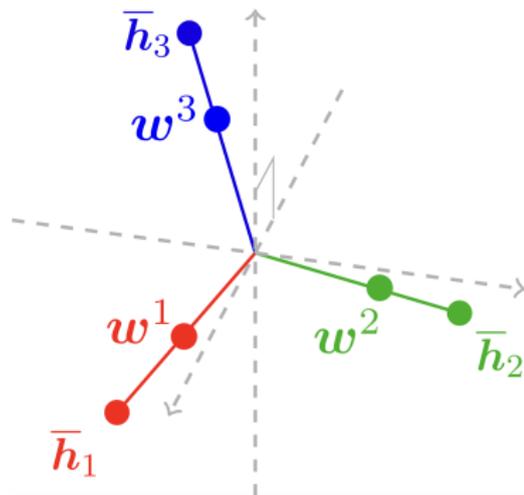
Neural Collapse in Classification IV

Neural Collapse (NC) refers to

- **NC3: Convergence to Self-Duality**: the last-layer classifiers are perfectly matched with the class-means of features

$$\frac{w^k}{\|w^k\|} \rightarrow \frac{\bar{h}_k}{\|\bar{h}_k\|},$$

where w^k represents the k -th row of W .



Neural Collapse in Classification V

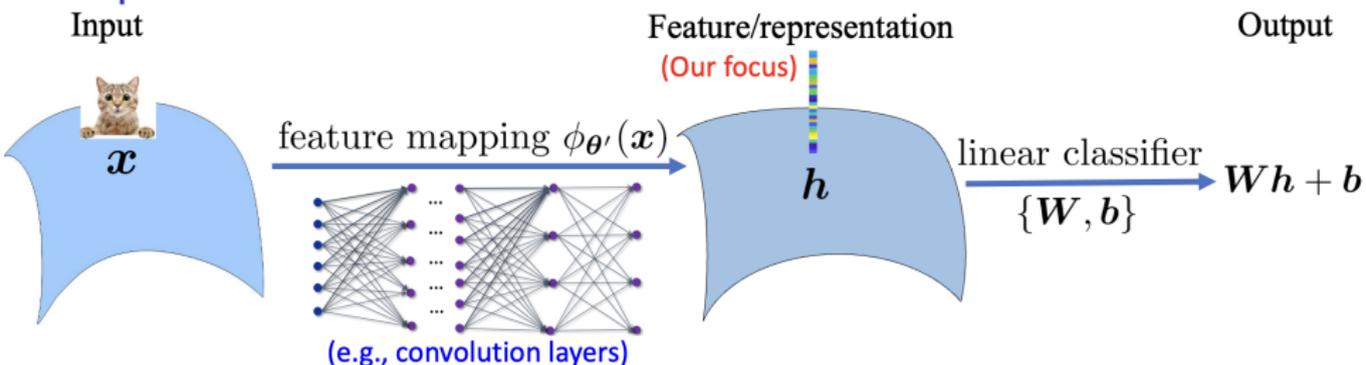
NC is preferred among every successful exercise in feature engineering
[Papayan et al.'20]

- Information Theory: Simplex ETF is the optimal Shannon code
- Classification: Simple ETF features \Rightarrow Simplex ETF max-margin classifier

Q: Why iterative training algorithm learns low-dimensional NC features and classifiers?

A: We will use tools developed in nonconvex optimization in Lecture 3 to understand NC phenomenon

Simplification: Unconstrained Features I

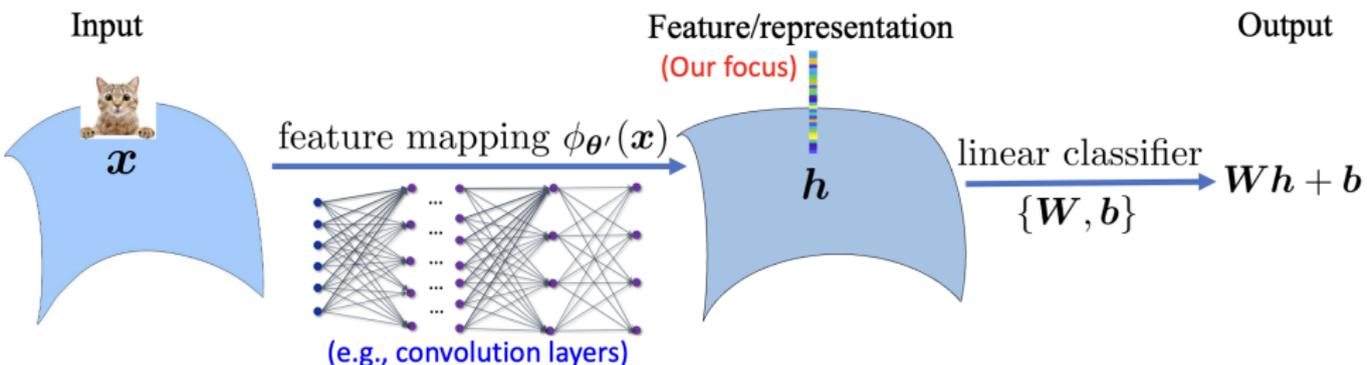


Training problem is highly nonconvex [Li et al.'18]:

$$\min_{\theta', W, b} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(W \phi_{\theta'}(\mathbf{x}_{k,i}) + b, \mathbf{y}_k) + \lambda \|(\theta', W, b)\|_F^2$$

- Neural Tangent Kernel focuses on output, and thus hardly provides much insights about features
- Neural Collapse is about the classifier W and the features $\phi_{\theta'}(\mathbf{x}_{k,i})$

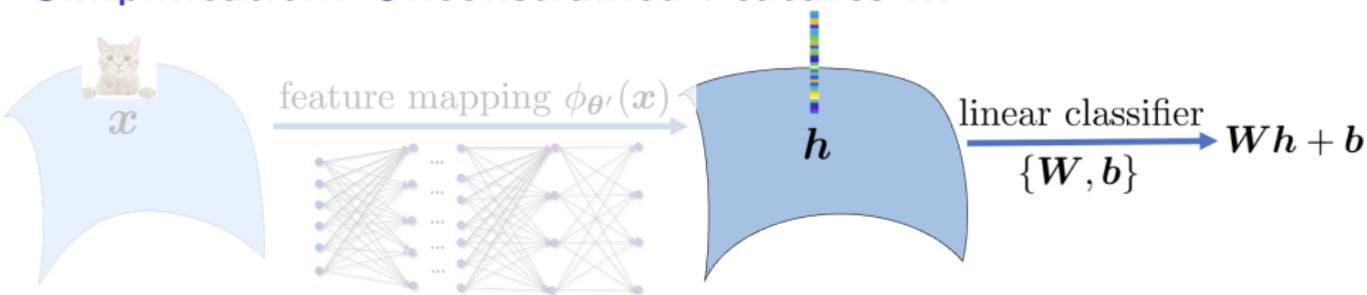
Simplification: Unconstrained Features II



- Neural Collapse is about the classifier W and the features $\phi_{\theta'}(\mathbf{x}_{k,i})$
- To understand NC, we treat the features $\mathbf{h}_{k,i} = \phi_{\theta'}(\mathbf{x}_{k,i})$ as free optimization variables (unconstrained features model [Mixon et al.'21])

$$\min_{\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W}\mathbf{h}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \lambda \|(\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b})\|_F^2$$

Simplification: Unconstrained Features III



$$\min_{\{h_{k,i}\}, \mathbf{W}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W}h_{k,i} + \mathbf{b}, \mathbf{y}_k) + \lambda \|(\{h_{k,i}\}, \mathbf{W}, \mathbf{b})\|_F^2$$

- **Validity:** Modern networks are highly **over-parameterized**, that can approximate any point in the feature space
- Also called **layer-peeled model** and has been studied recently to understand NC
- We will show such simplification preserves the core properties of last-layer classifiers and features—the NC phenomenon

Simplification: Unconstrained Features IV

[Lu et al.'20] study the following one-example-per class model

$$\min_{\{\mathbf{h}_k\}} \frac{1}{K} \sum_{k=1}^K \mathcal{L}_{\text{CE}}(\mathbf{h}_k, \mathbf{y}_k), \text{ s.t. } \|\mathbf{h}_{k,i}\|_2 = 1$$

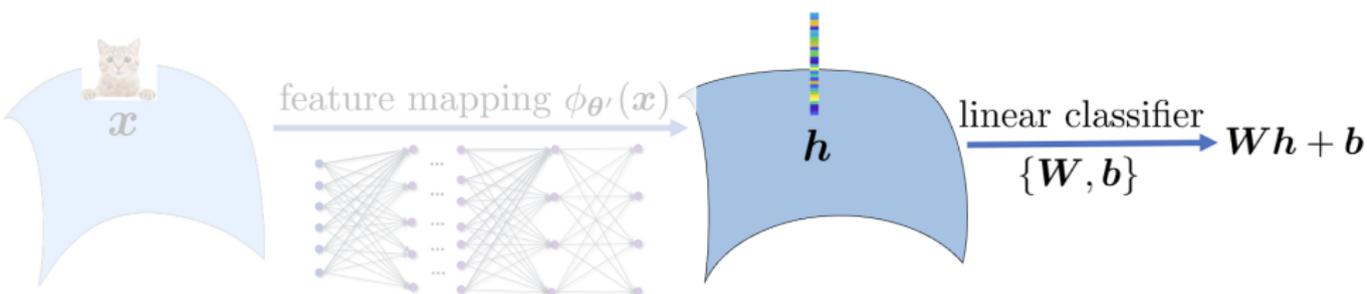
[E et al.'20, Fang et al.'21, Gal et al.'21, etc.] study constrained formulation

$$\min_{\{\mathbf{h}_k\}, \mathbf{W}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W}\mathbf{h}_{k,i}, \mathbf{y}_k), \text{ s.t. } \|\mathbf{W}\|_F \leq 1, \|\mathbf{h}_{k,i}\|_2 \leq 1$$

These work show that any global solution has NC, but

- What about [local minima/saddle points](#)?
- The constrained formulations are not aligned with practice

Geometric Analysis for Unconstrained Features Model I



$$\min_{\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W}\mathbf{h}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \lambda \|(\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b})\|_F^2$$

- Closely related to the matrix factorization problem in Lecture 3: bilinear form $\mathbf{W}\mathbf{h}_{k,i}$
- We will study its global/local minima and saddle points

Geometric Analysis for Unconstrained Features Model II

$$\min_{\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W}\mathbf{h}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \lambda \|(\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b})\|_F^2$$

Theorem (global optimality) [Zhu et al. 2021] Let feature dim. $d \geq \# \text{class } K - 1$. Then any global solution $(\{\mathbf{h}_{k,i}^*, \mathbf{W}^*, \mathbf{b}^*\})$ must satisfy NC: $\mathbf{b}^* = 0$ and

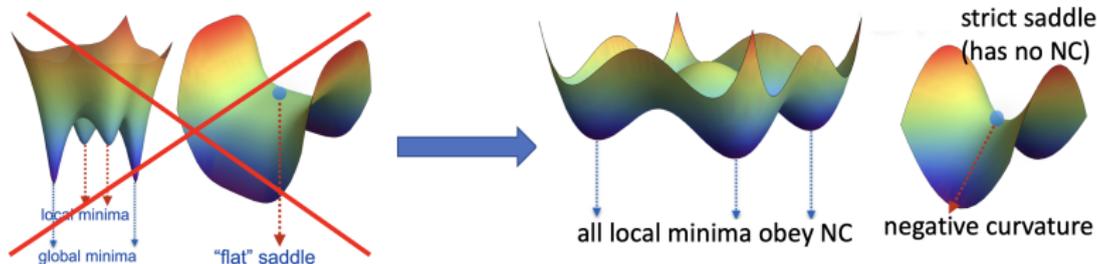
$$\underbrace{\mathbf{h}_{k,i}^* = \bar{\mathbf{h}}_k^*}_{\text{NC1}}, \quad \underbrace{\frac{\langle \bar{\mathbf{h}}_k^*, \bar{\mathbf{h}}_{k'}^* \rangle}{\|\bar{\mathbf{h}}_k^*\| \|\bar{\mathbf{h}}_{k'}^*\|} = \begin{cases} 1, & k = k' \\ -\frac{1}{K-1}, & k \neq k' \end{cases}}_{\text{NC2}}, \quad \underbrace{\frac{\mathbf{w}^{k*}}{\|\mathbf{w}^{k*}\|} = \frac{\bar{\mathbf{h}}_k^*}{\|\bar{\mathbf{h}}_k^*\|}}_{\text{NC3}}$$

- $d \geq K - 1$ is required to make K class-mean features equal angle and with cosine angle $-\frac{1}{K-1}$ (the largest possible) between each pair.

Geometric Analysis for Unconstrained Features Model III

$$\min_{\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W}\mathbf{h}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \lambda \|(\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b})\|_F^2$$

Theorem (benign global landscape) [Zhu et al. 2021] Let feature dim. $d > \# \text{class } K$. Then the above objective function (i) has no spurious local minima, and (ii) any non-global critical point is a strict saddle with negative curvature. Conjecture: $d \geq K - 1$ is sufficient.



General nonconvex problems

Our training problem

Geometric Analysis for Unconstrained Features Model IV

$$\min_{\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W}\mathbf{h}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \lambda \|(\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b})\|_F^2 \quad (\text{NVX})$$

Theorem (benign global landscape) [Zhu et al. 2021] Let feature dim. $d > \# \text{class } K$. Then the above objective function (i) has no spurious local minima, and (ii) any non-global critical point is a strict saddle with negative curvature.

- Proof idea: let $\mathbf{z}_{k,i} = \mathbf{W}\mathbf{h}_{k,i}$. Then (NVX) is equivalent to the following convex problem [Haeffele & Vidal'15, Li et al.'17, Ciliberto et al.'17]

$$\min_{\mathbf{Z}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{z}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \lambda \|\mathbf{Z}\|_* + \lambda \|\mathbf{b}\|_2^2 \quad (\text{CVX})$$

where $\|\cdot\|_*$ is the nuclear norm (sum of singular values).

Geometric Analysis for Unconstrained Features Model V

$$\min_{\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W}\mathbf{h}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \lambda \|(\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b})\|_F^2 \quad (\text{NVX})$$

$$\min_{\mathbf{Z}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{z}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \lambda \|\mathbf{Z}\|_* + \lambda \|\mathbf{b}\|_2^2 \quad (\text{CVX})$$

- Step 1: (NVX) and (CVX) have the "same" global solutions: if $(\mathbf{H}^*, \mathbf{W}^*, \mathbf{b}^*)$ is a global solution of (NVX), then $(\mathbf{W}^* \mathbf{H}^*, \mathbf{b}^*)$ is a global solution of (CVX); vice versa.

$$\text{variational form } \|\mathbf{Z}\|_* = \min_{\mathbf{Z}=\mathbf{W}\mathbf{H}} \frac{1}{2} (\|\mathbf{W}\|_F^2 + \|\mathbf{H}\|_F^2)$$

Geometric Analysis for Unconstrained Features Model VI

$$\min_{\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W}\mathbf{h}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \lambda \|(\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b})\|_F^2 \quad (\text{NVX})$$

$$\min_{\mathbf{Z}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{z}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \lambda \|\mathbf{Z}\|_* + \lambda \|\mathbf{b}\|_2^2 \quad (\text{CVX})$$

- Step 2: if $(\mathbf{H}, \mathbf{W}, \mathbf{b})$ is a critical point but not a global min of (NVX)
 - (\mathbf{Z}, \mathbf{b}) with $\mathbf{Z} = \mathbf{W}\mathbf{H}$ is not a critical point to (CVX)
 - (\mathbf{Z}, \mathbf{b}) does not satisfy the first-order optimality condition of (CVX)
 - Exploiting this, we show the Hessian at $(\mathbf{H}, \mathbf{W}, \mathbf{b})$ has a negative eigenvalue, i.e., it is a strict saddle of (NVX)

Geometric Analysis for Unconstrained Features Model VII

$$\min_{\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W}\mathbf{h}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \lambda \|(\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b})\|_F^2 \quad (\text{NVX})$$

$$\min_{\mathbf{Z}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{z}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \lambda \|\mathbf{Z}\|_* + \lambda \|\mathbf{b}\|_2^2 \quad (\text{CVX})$$

- Step 1: (NVX) and (CVX) have the "same" global solutions.
- Step 2: if $(\mathbf{H}, \mathbf{W}, \mathbf{b})$ is a critical point but not a global min of (NVX)
 - the Hessian at $(\mathbf{H}, \mathbf{W}, \mathbf{b})$ has a negative eigenvalue, i.e., it is a strict saddle
- Step 2 holds for any non-global critical point \Rightarrow (NVX) has benign global landscape (no spurious local minima & strict saddle function)

Geometric Analysis for Unconstrained Features Model VIII

$$\min_{\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W}\mathbf{h}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \lambda \|(\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b})\|_F^2$$

Theorem (global optimality & benign global landscape) Let feature dim. $d > \# \text{class } K$.

- Any global solution $(\{\mathbf{h}_{k,i}^*, \mathbf{W}^*, \mathbf{b}^*\})$ obeys Neural Collapse.
- The objective function (i) has no spurious local minima, and (ii) any non-global critical point is a strict saddle with negative curvature.

Message. Iterative algorithms such as (stochastic) gradient descent always learns Neural Collapse features and classifiers.

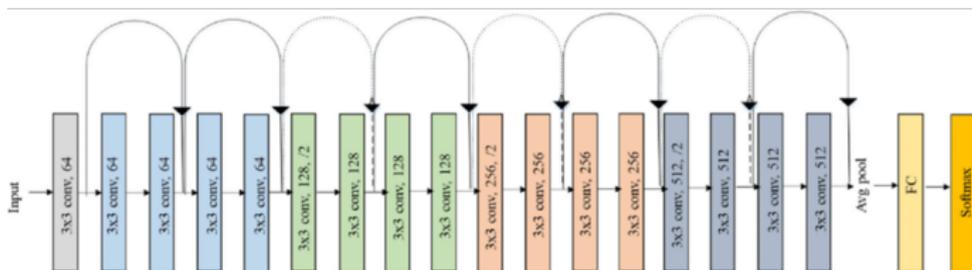
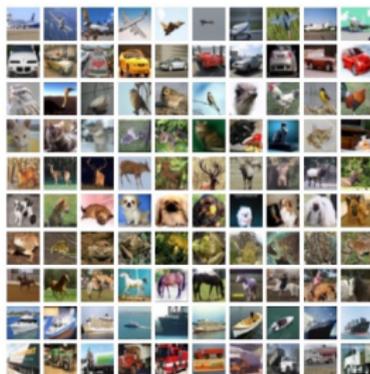
Experiments on Practical Neural Networks

Conduct experiments with **practical networks** to verify our findings on Unconstrained Features Model

Use a Residual Neural Network (ResNet) on CIFAR-10 Dataset:

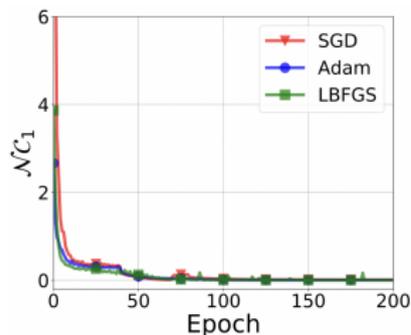
- $K = 10$ classes
- 50K training images
- 10K testing images

airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

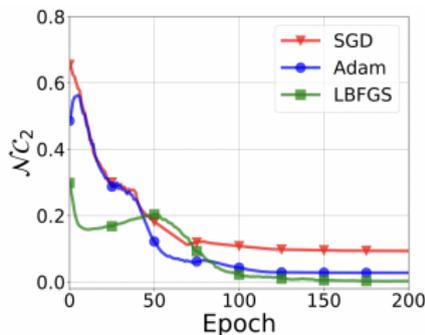


Experiments: NC is algorithm independent

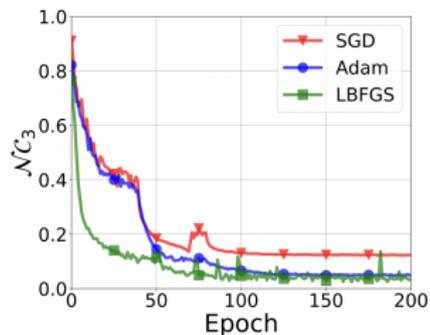
ResNet18 on CIFAR-10 with **different training algorithms**



Within-Class Variability (NC1)



Between-Class Separation (NC2)

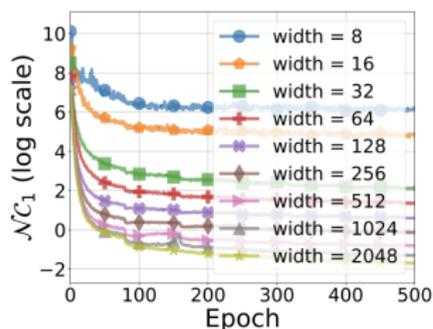


Self-Duality Collapse (NC3)

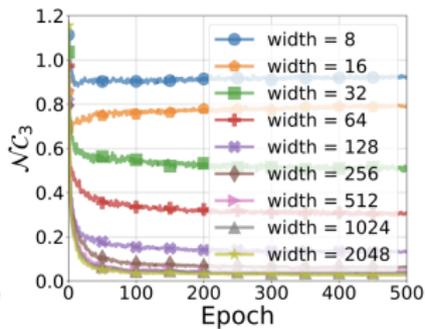
- The smaller the quantities, the severer NC
- NC across **different training algorithms**

Experiments: NC Occurs on Random Labels/Inputs

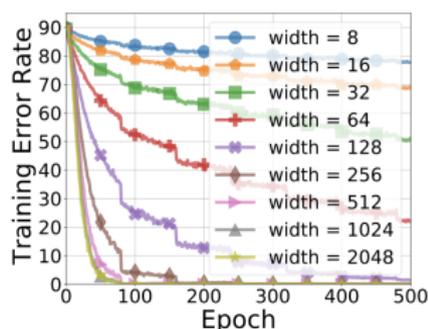
CIFAR-10 with **random** labels, multi-layer perceptron (MLP) with **varying** network widths



Within-Class Variability (NC1)



Self-Duality Collapse (NC2)



Training Error

- **Validity of unconstrained features model:** Learn NC last-layer features and classifiers for any inputs
- The network memorizes training data in a very special way: NC
- We observe similar results on **random inputs (random pixels)**

Exploit NC

Experiments in [Papayan, Han Donoho] shows NC leads to better

- Generalization performance
- Robustness

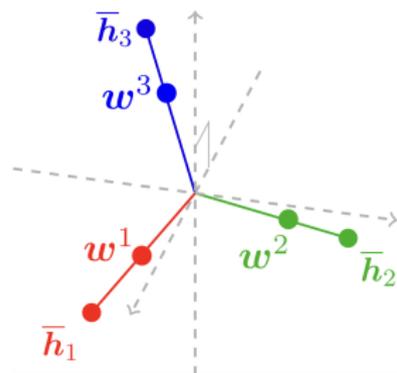
We can also exploit NC for

- Improving training efficiency & memory cost (covered later)
- Understanding the effect of loss functions (covered later)
- Understanding transferability
- etc.

Exploit NC for Improving Training & Memory I

NC is prevalent, and classifier always converges to a Simplex ETF

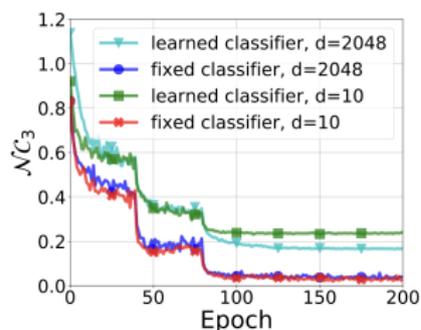
- **Implication 1: No need to learn the classifier** [Hoffer et al. 2018]
 - Just fix it as a Simplex ETF
 - Save **8%, 12%, and 53%** parameters for ResNet50, DenseNet169, and ShuffleNet!
- **Implication 2: No need of large feature dimension d**
 - Just use feature dim. $d = \# \text{class } K$ (e.g., $d = 10$ for CIFAR-10)
 - Further saves **21% and 4.5%** parameters for ResNet18 and ResNet50!



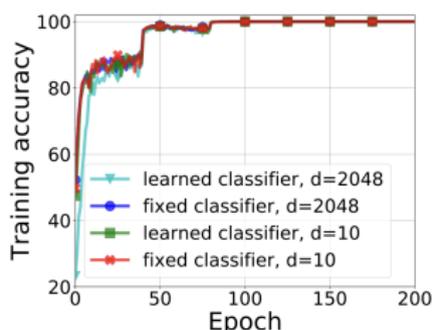
Exploit NC for Improving Training & Memory II

ResNet50 on CIFAR-10 with different settings

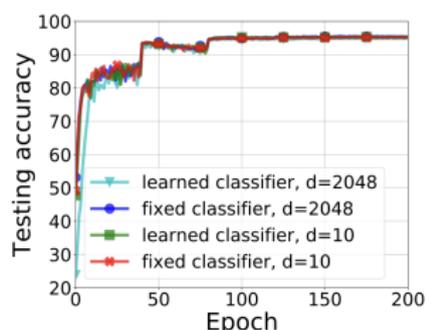
- Learned classifier (default) VS fixed classifier as a simplex ETF
- Feature dim $d = 2048$ (default) VS $d = 10$



Self-Duality Collapse (NC3)



Training Accuracy

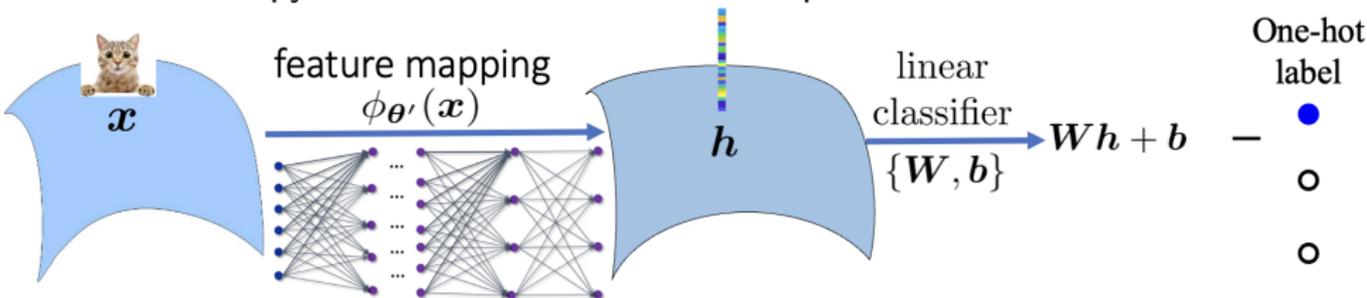


Testing Accuracy

- Training with **small** dimensional features and **fixed** classifiers achieves on-par performance with **large** dimensional features and **learned** classifiers.

Is Cross-entropy Loss Essential?

Is cross-entropy loss essential to neural collapse?



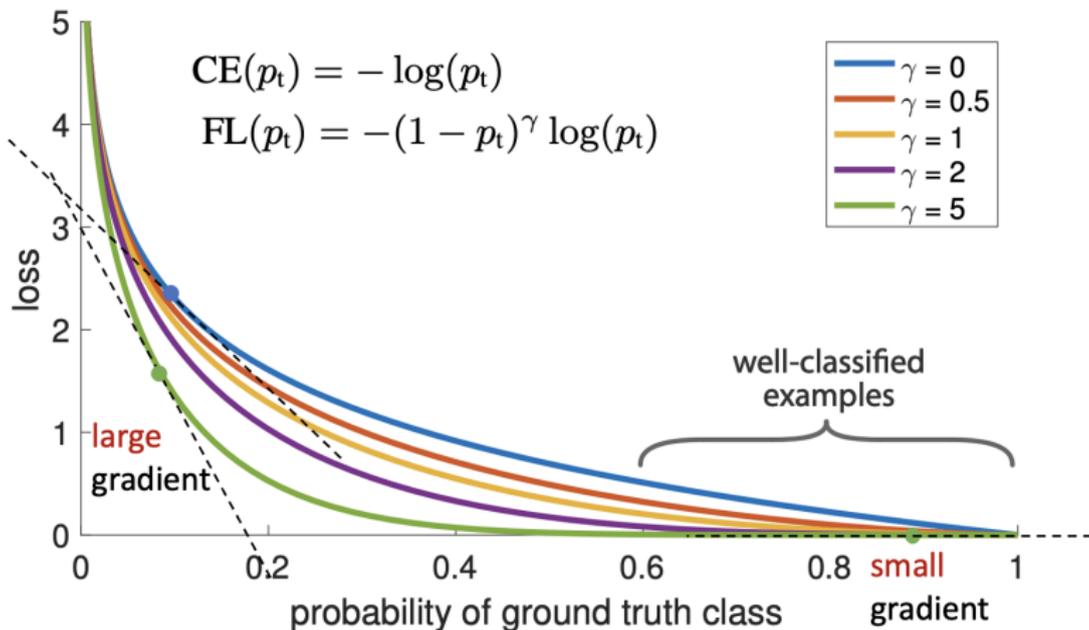
We can measure the mismatch between the network output and the one-hot label in many ways.

Various losses and tricks (e.g., label smoothing, focal loss) have been proposed to improve network training and performance¹

¹He et al., Bag of tricks for image classification with convolutional neural networks, CVPR'19.

Focal Loss (FL)

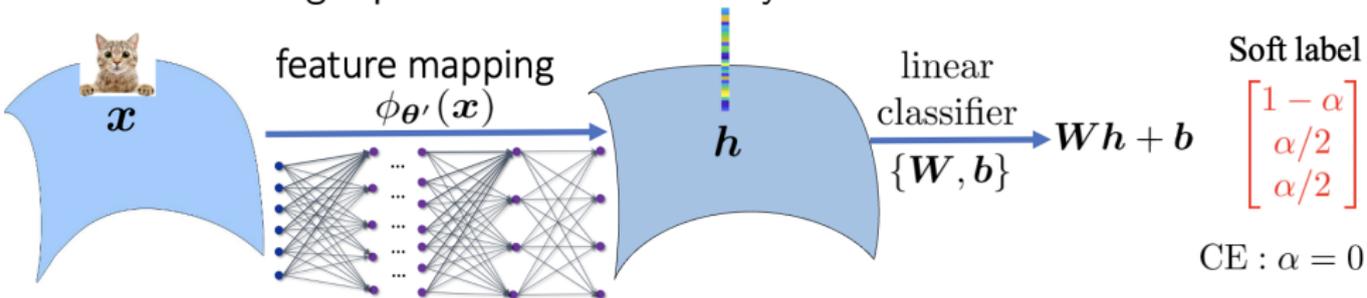
Focal loss puts more focus on hard, misclassified examples²



²Lin et al., Focal Loss for Dense Object Detection, CVPR'18.

Label Smoothing (LS)

Label smoothing replaces the hard label by a soft label ³



Output: $Wh + b = \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix}$ $\xrightarrow{\text{Softmax function}}$ $\begin{bmatrix} 0.6 \\ 0.3 \\ 0.1 \end{bmatrix}$ **Prediction**

Target $\begin{bmatrix} 1 - \alpha \\ \alpha/2 \\ \alpha/2 \end{bmatrix}$

Cat Dog Panda

$$\text{LS} = -q(\text{Cat}) \cdot \log p(\text{Cat})$$

$$-q(\text{Dog}) \cdot \log p(\text{Dog})$$

$$-q(\text{Panda}) \cdot \log p(\text{Panda})$$

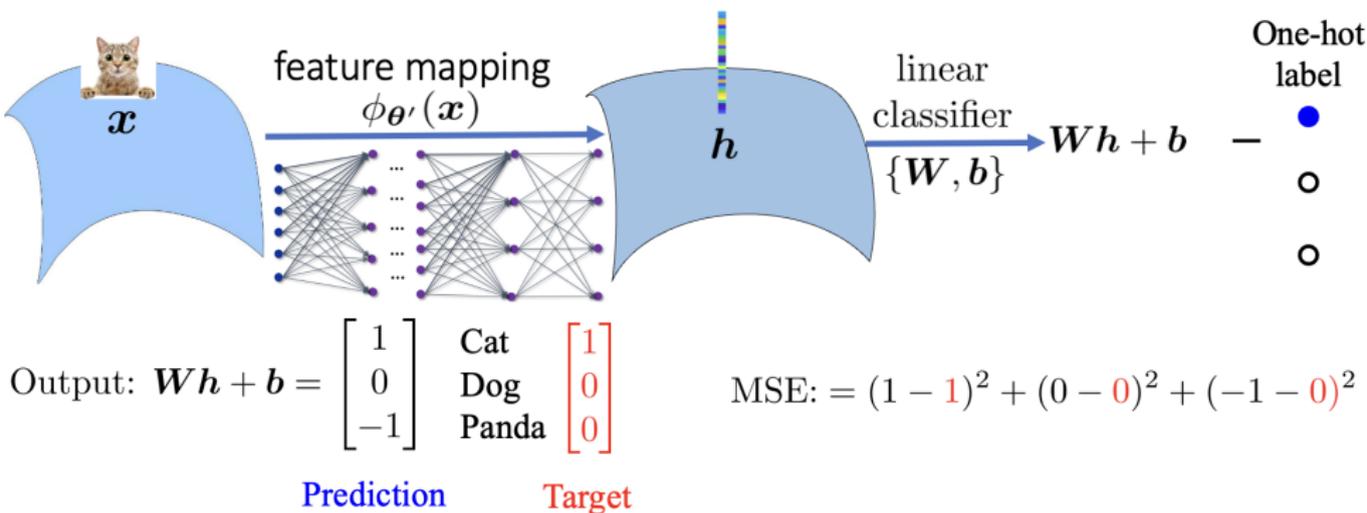
$$= -(1 - \alpha) \log(0.6)$$

$$- \frac{\alpha}{2} \log(0.3)$$

$$- \frac{\alpha}{2} \log(0.1)$$

³Szegedy et al., Rethinking the inception architecture for computer vision, CVPR'16.
Muller, Kornblith, Hinton, When does label smoothing help?, NeurIPS'19.

Mean-squared Error (MSE) Loss?



Compared with CE, (rescaled) MSE loss produces on par/slightly worse results for computer vision tasks and on par/slightly better results for NLP tasks.⁴

⁴Hui & Belkin, Evaluation of neural architectures trained with square loss vs cross-entropy in classification tasks, ICLR 2021.

Are All Losses Created Equal?—A NC Perspective I

Do all these losses make difference?

We study them under the unconstrained feature model:

$$\min_{\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}(\mathbf{W}\mathbf{h}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \lambda \|(\{\mathbf{h}_{k,i}\}, \mathbf{W}, \mathbf{b})\|_F^2$$

Theorem (informal) [Zhou et al.'22] With feature dim. $d > \#class K$, all the one-hot labeling based losses (e.g., CE, FL, LS, MSE) lead to (almost) the same NC features and classifiers [Han et al'21, Tirer & Bruner'22, Zhou'22].

Are All Losses Created Equal?—A NC Perspective II

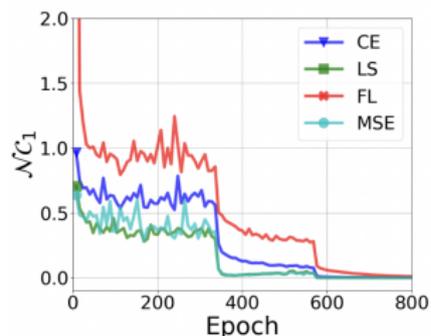
Theorem (informal) [Zhou et al.'22] With feature dim. $d > \# \text{class } K$, all the one-hot labeling based losses (e.g., CE, FL, LS, MSE) lead to (almost) the same NC features and classifiers [Han et al'21, Tixer & Bruner'22, Zhou'22].

Implication If network is large enough and trained longer enough

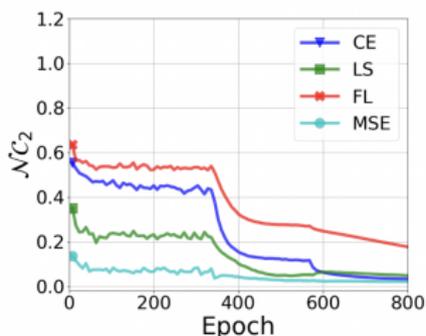
- All losses lead to largely identical features on training data—NC phenomena
- All losses lead to largely identical performance on test data (experiments in the following slides)

Are All Losses Created Equal?—A NC Perspective III

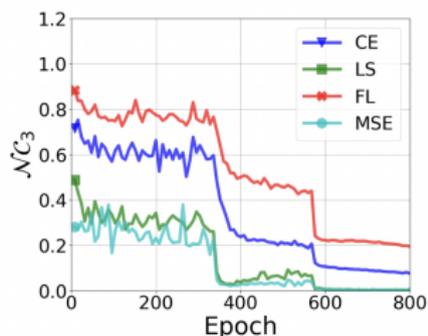
ResNet50 on CIFAR-10 with **different training losses**



Within-Class Variability (NC1)



Between-Class Separation (NC2)

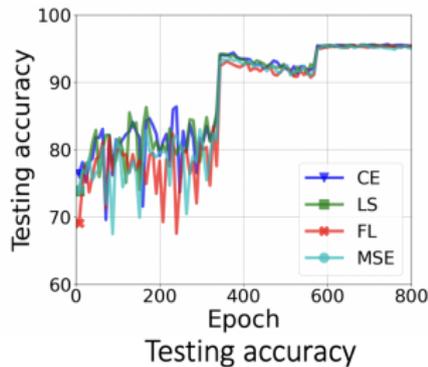
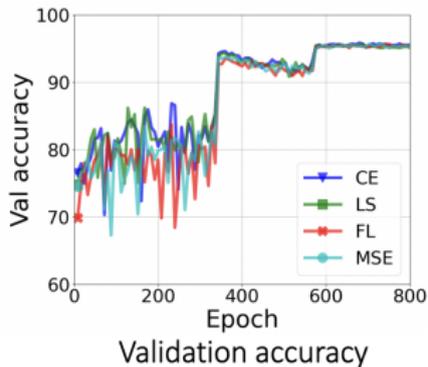
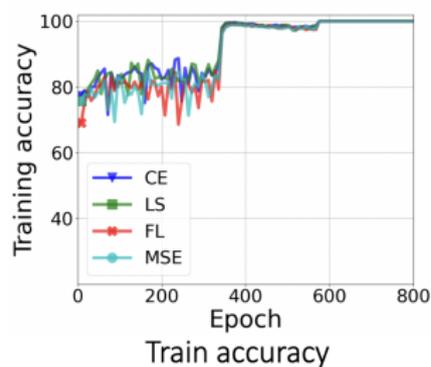


Self-Duality Collapse (NC3)

- The smaller the quantities, the severer NC
- NC across **different training losses**

Are All Losses Created Equal?—A NC Perspective IV

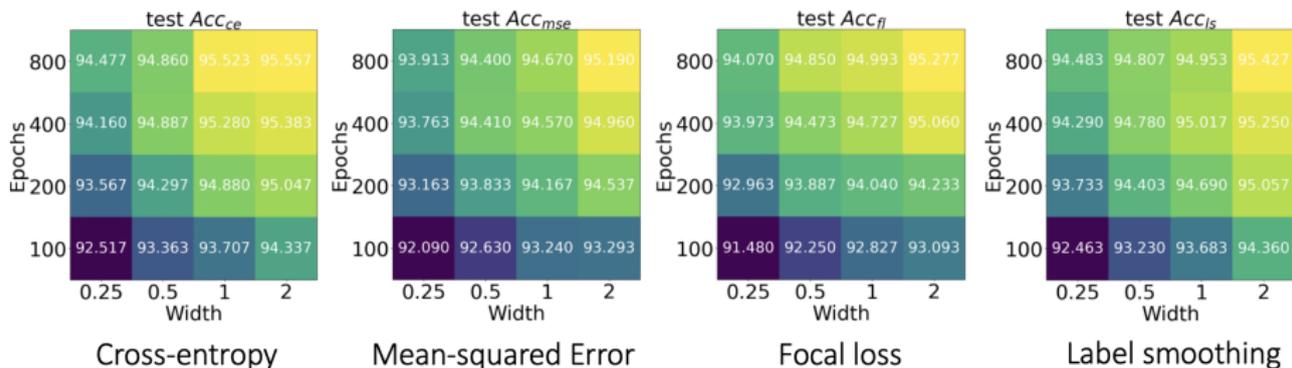
ResNet50 on CIFAR-10 with **different training losses**



- All losses lead to largely identical performance on training, validation, and test data

Are All Losses Created Equal?—A NC Perspective V

ResNet50 (with different network widths and training epochs) on CIFAR-10 with **different training losses**



- If network is large enough and trained longer enough, all losses lead to largely identical performance on test data

Outline

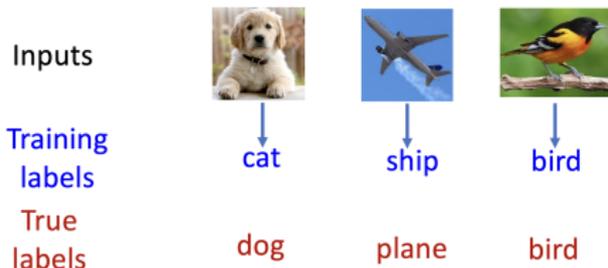
Recap and Outlook

- ① Motivating Examples for Low-Dim Structure in Deep Learning
- ② Resource Tradeoffs in the Multiple Manifold Problem
 - Problem Formulation
 - Intrinsic Geometric Properties of Manifold Data
 - Network Architecture Resources and Training Procedure
 - Training Deep Networks with Gradient Descent
 - Resource Tradeoffs
- ③ Looking Inside: Neural Collapse in the Multiple Manifold Problem
 - Learned low-dimensional features—NC phenomena
 - Geometric analysis for understanding neural collapse
 - Exploit NC for improving training efficiency
 - Exploit NC for understanding the effect of loss functions
- ④ Exploit Sparse Model for Robust training

NC \rightarrow Overfitting to Corruptions!

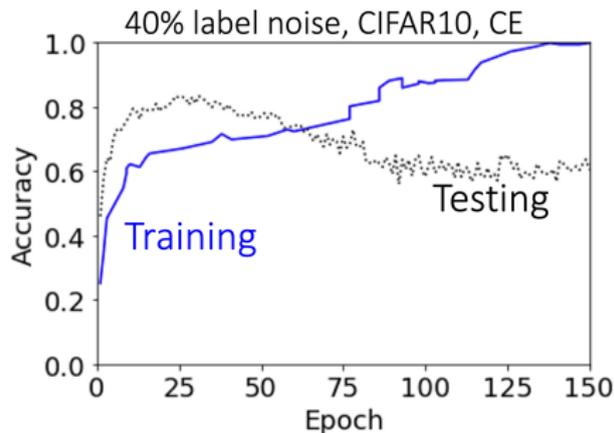
Label noise is common and often unavoidable

- Some proportion of the labels are incorrect (5-80%?)
- We don't know which labels are correct/incorrect



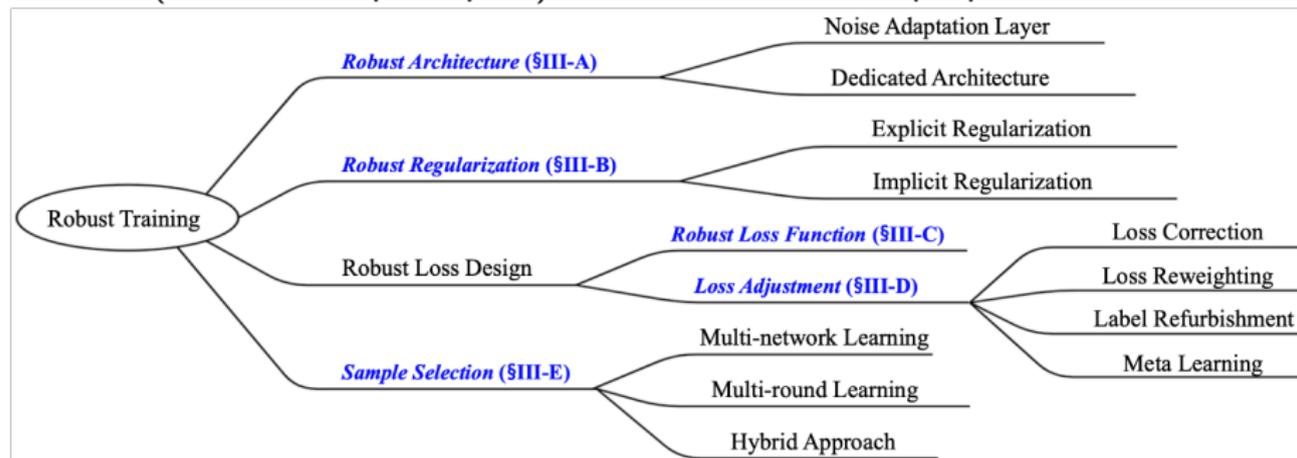
NC always happens

- Perfectly fits noisy labels (overfitting)
- Can't predict well on new images



Prior Work on Robust Deep Learning for Noisy Labels

Various (heuristic or principled) methods have been proposed⁵



⁵Song et al., Learning from noisy labels with deep neural networks: A survey, IEEE TNNLS, 2022.

A Sparse Over-Parameterization (SOP) Method

We model the label noise and (hopefully) correct it. Only a fraction of the labels are corrupted (sparse), and the corruption in each label is also sparse



$$\begin{array}{c}
 \text{"cat"} \\
 \text{"dog"} \\
 \vdots \\
 0
 \end{array}
 \begin{bmatrix}
 0 \\
 1 \\
 \vdots \\
 0
 \end{bmatrix}
 =
 f(x; \theta)
 +
 \begin{bmatrix}
 1 \\
 0 \\
 \vdots \\
 0
 \end{bmatrix}
 +
 \begin{bmatrix}
 -1 \\
 1 \\
 \vdots \\
 0
 \end{bmatrix}$$

Corrupted label
True label
Sparse noise

Lecture 1 introduced principled methods for dealing with sparse corruption in compressive sensing, robust PCA⁶

⁶Candes & Tao, Decoding by linear programming, TIT 2005.

Wright et al., Robust face recognition via sparse representation, TPAMI, 2008.

Candes et al., Robust principal component analysis? JACM, 2011.

A Sparse Over-Parameterization (SOP) Method

Our approach:⁷ minimize the distance between \mathbf{y} and $f(\boldsymbol{\theta}; \mathbf{x}) + \mathbf{s}$

$$\min_{\boldsymbol{\theta}, \mathbf{u}_i, \mathbf{v}_i} \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{\text{CE}}(f(\mathbf{x}_i; \boldsymbol{\theta}) + \underbrace{\mathbf{u}_i \odot \mathbf{u}_i - \mathbf{v}_i \odot \mathbf{v}_i}_{\text{over-parameterize } \mathbf{s}_i \text{ to promote sparsity}}, \mathbf{y}_i)$$

Here the over-parameterization $\mathbf{u}_i \odot \mathbf{u}_i - \mathbf{v}_i \odot \mathbf{v}_i$ introduces implicit algorithmic regularization [Vaskevicius et al.'19, Zhao et al.'19]

$$\text{variational form } \|\mathbf{s}\|_1 = \min_{\mathbf{s}=\mathbf{u}\odot\mathbf{u}-\mathbf{v}\odot\mathbf{v}} \frac{1}{2} (\|\mathbf{u}\|^2 + \|\mathbf{v}\|^2)$$

⁷Liu, Zhu, Qu, You, Robust Training under Label Noise by Over-parameterization, ICML'22

A Sparse Over-Parameterization (SOP) Method

Our approach:⁷ minimize the distance between \mathbf{y} and $f(\boldsymbol{\theta}; \mathbf{x}) + \mathbf{s}$

$$\min_{\boldsymbol{\theta}, \mathbf{u}_i, \mathbf{v}_i} \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{\text{CE}}(f(\mathbf{x}_i; \boldsymbol{\theta}) + \underbrace{\mathbf{u}_i \odot \mathbf{u}_i - \mathbf{v}_i \odot \mathbf{v}_i}_{\text{over-parameterize } \mathbf{s}_i \text{ to promote sparsity}}, \mathbf{y}_i)$$

Here the over-parameterization $\mathbf{u}_i \odot \mathbf{u}_i - \mathbf{v}_i \odot \mathbf{v}_i$ introduces implicit algorithmic regularization [Vaskevicius et al.'19, Zhao et al.'19]

$$\text{variational form } \|\mathbf{s}\|_1 = \min_{\mathbf{s}=\mathbf{u}\odot\mathbf{u}-\mathbf{v}\odot\mathbf{v}} \frac{1}{2} (\|\mathbf{u}\|^2 + \|\mathbf{v}\|^2)$$

Why not use explicit regularization?

$$\min_{\boldsymbol{\theta}, \{\mathbf{s}_i\}} \frac{1}{N} \sum_{i=1}^N \underbrace{\mathcal{L}_{\text{CE}}(f(\mathbf{x}_i; \boldsymbol{\theta}) + \mathbf{s}_i, \mathbf{y}_i)}_{\rightarrow 0} + \lambda \underbrace{\|\mathbf{s}_i\|_1}_{\rightarrow 0}$$

⁷Liu, Zhu, Qu, You, Robust Training under Label Noise by Over-parameterization, ICML'22

A Sparse Over-Parameterization (SOP) Method

A simple model: assume $f(\mathbf{x}; \boldsymbol{\theta})$ is a scalar function and can be approximated by first-order Taylor expansion

$$f(\mathbf{x}; \boldsymbol{\theta}) \approx f(\mathbf{x}; \boldsymbol{\theta}_0) + \langle \nabla f(\mathbf{x}; \boldsymbol{\theta}_0), \boldsymbol{\theta} - \boldsymbol{\theta}_0 \rangle$$

A Sparse Over-Parameterization (SOP) Method

A simple model: assume $f(\mathbf{x}; \boldsymbol{\theta})$ is a scalar function and can be approximated by first-order Taylor expansion

$$f(\mathbf{x}; \boldsymbol{\theta}) \approx f(\mathbf{x}; \boldsymbol{\theta}_0) + \langle \nabla f(\mathbf{x}; \boldsymbol{\theta}_0), \boldsymbol{\theta} - \boldsymbol{\theta}_0 \rangle$$

WLOG, assume $f(\mathbf{x}; \boldsymbol{\theta}_0) + \langle \nabla f(\mathbf{x}; \boldsymbol{\theta}_0), \boldsymbol{\theta}_0 \rangle = 0$. For N training samples,

$$\begin{bmatrix} f(\mathbf{x}_1; \boldsymbol{\theta}) \\ \vdots \\ f(\mathbf{x}_N; \boldsymbol{\theta}) \end{bmatrix} \approx \begin{bmatrix} \nabla f(\mathbf{x}_1; \boldsymbol{\theta}_0)^\top \\ \vdots \\ \nabla f(\mathbf{x}_N; \boldsymbol{\theta}_0)^\top \end{bmatrix} \boldsymbol{\theta} = \mathbf{J} \cdot \boldsymbol{\theta}$$

A Sparse Over-Parameterization (SOP) Method

A simple model: assume $f(\mathbf{x}; \boldsymbol{\theta})$ is a scalar function and can be approximated by first-order Taylor expansion

$$f(\mathbf{x}; \boldsymbol{\theta}) \approx f(\mathbf{x}; \boldsymbol{\theta}_0) + \langle \nabla f(\mathbf{x}; \boldsymbol{\theta}_0), \boldsymbol{\theta} - \boldsymbol{\theta}_0 \rangle$$

WLOG, assume $f(\mathbf{x}; \boldsymbol{\theta}_0) + \langle \nabla f(\mathbf{x}; \boldsymbol{\theta}_0), \boldsymbol{\theta}_0 \rangle = 0$. For N training samples,

$$\begin{bmatrix} f(\mathbf{x}_1; \boldsymbol{\theta}) \\ \vdots \\ f(\mathbf{x}_N; \boldsymbol{\theta}) \end{bmatrix} \approx \begin{bmatrix} \nabla f(\mathbf{x}_1; \boldsymbol{\theta}_0)^\top \\ \vdots \\ \nabla f(\mathbf{x}_N; \boldsymbol{\theta}_0)^\top \end{bmatrix} \boldsymbol{\theta} = \mathbf{J} \cdot \boldsymbol{\theta}$$

This leads to the following corrupted observation problem

$$\mathbf{y} = \mathbf{J} \cdot \boldsymbol{\theta}_\star + \mathbf{s}_\star$$

where $\boldsymbol{\theta}_\star$ is the underlying groundtruth parameter, and \mathbf{s}_\star is sparse.

A Sparse Over-Parameterization (SOP) Method

We over-parameterize the sparse noise by $\mathbf{u} \odot \mathbf{u} - \mathbf{v} \odot \mathbf{v}$ and solve

$$\min_{\boldsymbol{\theta}, \mathbf{u}, \mathbf{v}} g(\boldsymbol{\theta}, \mathbf{u}, \mathbf{v}) = \frac{1}{2} \|\mathbf{J} \cdot \boldsymbol{\theta} + \mathbf{u} \odot \mathbf{u} - \mathbf{v} \odot \mathbf{v} - \mathbf{y}\|_2^2$$

using gradient descent with *discrepant learning rates*

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \mu \nabla_{\boldsymbol{\theta}} g(\boldsymbol{\theta}_t, \mathbf{u}_t, \mathbf{v}_t) \\ \begin{bmatrix} \mathbf{u}_{t+1} \\ \mathbf{v}_{t+1} \end{bmatrix} &= \begin{bmatrix} \mathbf{u}_t \\ \mathbf{v}_t \end{bmatrix} - \alpha \mu \begin{bmatrix} \nabla_{\mathbf{u}} g(\boldsymbol{\theta}_t, \mathbf{u}_t, \mathbf{v}_t) \\ \nabla_{\mathbf{v}} g(\boldsymbol{\theta}_t, \mathbf{u}_t, \mathbf{v}_t) \end{bmatrix} \end{aligned}$$

A Sparse Over-Parameterization (SOP) Method

We over-parameterize the sparse noise by $\mathbf{u} \odot \mathbf{u} - \mathbf{v} \odot \mathbf{v}$ and solve

$$\min_{\boldsymbol{\theta}, \mathbf{u}, \mathbf{v}} g(\boldsymbol{\theta}, \mathbf{u}, \mathbf{v}) = \frac{1}{2} \|\mathbf{J} \cdot \boldsymbol{\theta} + \mathbf{u} \odot \mathbf{u} - \mathbf{v} \odot \mathbf{v} - \mathbf{y}\|_2^2$$

using gradient descent with *discrepant learning rates*

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \mu \nabla_{\boldsymbol{\theta}} g(\boldsymbol{\theta}_t, \mathbf{u}_t, \mathbf{v}_t) \\ \begin{bmatrix} \mathbf{u}_{t+1} \\ \mathbf{v}_{t+1} \end{bmatrix} &= \begin{bmatrix} \mathbf{u}_t \\ \mathbf{v}_t \end{bmatrix} - \alpha \mu \begin{bmatrix} \nabla_{\mathbf{u}} g(\boldsymbol{\theta}_t, \mathbf{u}_t, \mathbf{v}_t) \\ \nabla_{\mathbf{v}} g(\boldsymbol{\theta}_t, \mathbf{u}_t, \mathbf{v}_t) \end{bmatrix} \end{aligned}$$

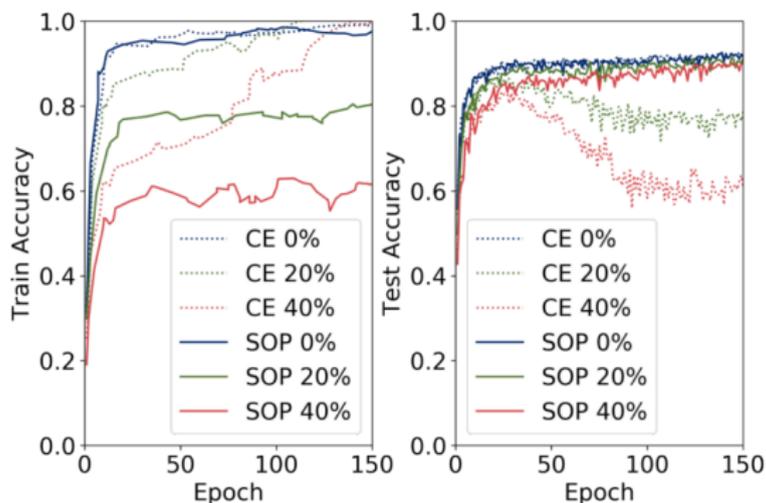
Theorem (informal) If gradient descent with infinitesimally small initialization and step size μ converges to $(\hat{\boldsymbol{\theta}}, \hat{\mathbf{u}}, \hat{\mathbf{v}})$, then $(\hat{\boldsymbol{\theta}}, \hat{\mathbf{u}} \odot \hat{\mathbf{u}} - \hat{\mathbf{v}} \odot \hat{\mathbf{v}})$ is an optimal solution to the following convex problem

$$\min_{\boldsymbol{\theta}, \mathbf{s}} \frac{1}{2} \|\boldsymbol{\theta}\|_2^2 + \frac{1}{\alpha} \|\mathbf{s}\|_1, \text{ s.t. } \mathbf{y} = \mathbf{J} \cdot \boldsymbol{\theta} + \mathbf{s}$$

Exactly recover $(\boldsymbol{\theta}_*, \mathbf{s}_*)$ when \mathbf{J} is incoherent [Candes & Tao'05].

A Sparse Over-Parameterization (SOP) Method

$\{0\%, 20\%, 40\%\}$ percent of labels for CIFAR-10 training data are randomly flipped uniformly to another class. Use ResNet34.

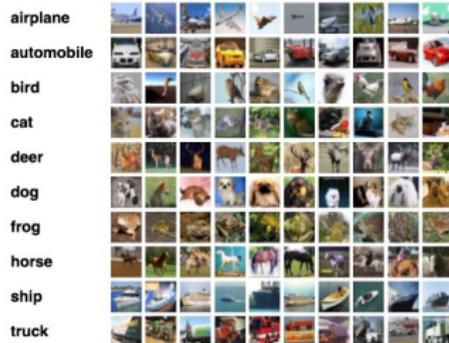


SOP trains a deep image classification networks without overfitting to wrong labels and obtain better generalization performance

SOP on CIFAR-10 with human annotated noisy labels

CIFAR-10N: provide CIFAR-10 with human annotated noisy labels⁸

- Annotated by 747 independent workers
- Provide 5 noisy label sets for CIFAR-10 train images:
- **Random** $i = 1, 2, 3$: the i -th submitted label for each image;
- **Aggregate**: aggregation of three noisy labels by majority voting
- **Worst**: label set with the highest noise rate



Label Set	CIFAR-10N Aggregate	CIFAR-10N Random 1	CIFAR-10N Random 2	CIFAR-10N Random 3	CIFAR-10N Worst
Noise Rate	9.03%	17.23%	18.12%	17.64%	40.21%

⁸Wei et al., Learning with noisy labels revisited: A study using real-world human annotations, ICLR 2022. [↗](#) [↻](#) [🔍](#)

SOP on CIFAR-10 with human annotated noisy labels

Method	CIFAR-10N					
	Clean	Aggregate	Random 1	Random 2	Random 3	Worst
CE (Standard)	92.92 ± 0.11	87.77 ± 0.38	85.02 ± 0.65	86.46 ± 1.79	85.16 ± 0.61	77.69 ± 1.55
Forward T (Patrini et al., 2017)	93.02 ± 0.12	88.24 ± 0.22	86.88 ± 0.50	86.14 ± 0.24	87.04 ± 0.35	79.79 ± 0.46
Backward T (Patrini et al., 2017)	93.10 ± 0.05	88.13 ± 0.29	87.14 ± 0.34	86.28 ± 0.80	86.86 ± 0.41	77.61 ± 1.05
GCE (Zhang & Sabuncu, 2018)	92.83 ± 0.16	87.85 ± 0.70	87.61 ± 0.28	87.70 ± 0.56	87.58 ± 0.29	80.66 ± 0.35
Co-teaching (Han et al., 2018)	93.35 ± 0.14	91.20 ± 0.13	90.33 ± 0.13	90.30 ± 0.17	90.15 ± 0.18	83.83 ± 0.13
Co-teaching+ (Yu et al., 2019)	92.41 ± 0.20	90.61 ± 0.22	89.70 ± 0.27	89.47 ± 0.18	89.54 ± 0.22	83.26 ± 0.17
T-Revision (Xia et al., 2019)	93.35 ± 0.23	88.52 ± 0.17	88.33 ± 0.32	87.71 ± 1.02	87.79 ± 0.67	80.48 ± 1.20
Peer Loss (Liu & Guo, 2020)	93.99 ± 0.13	90.75 ± 0.25	89.06 ± 0.11	88.76 ± 0.19	88.57 ± 0.09	82.00 ± 0.60
ELR (Liu et al., 2020)	93.45 ± 0.65	92.38 ± 0.64	91.46 ± 0.38	91.61 ± 0.16	91.41 ± 0.44	83.58 ± 1.13
ELR+ (Liu et al., 2020)	95.39 ± 0.05	94.83 ± 0.10	94.43 ± 0.41	94.20 ± 0.24	94.34 ± 0.22	91.09 ± 1.60
Positive-LS (Lukasik et al., 2020)	94.77 ± 0.17	91.57 ± 0.07	89.80 ± 0.28	89.35 ± 0.33	89.82 ± 0.14	82.76 ± 0.53
F-Div (Wei & Liu, 2020)	94.88 ± 0.12	91.64 ± 0.34	89.70 ± 0.40	89.79 ± 0.12	89.55 ± 0.49	82.53 ± 0.52
Divide-Mix (Li et al., 2020)	95.37 ± 0.14	95.01 ± 0.71	95.16 ± 0.19	95.23 ± 0.07	95.21 ± 0.14	92.56 ± 0.42
Negative-LS (Wei et al., 2021)	94.92 ± 0.25	91.97 ± 0.46	90.29 ± 0.32	90.37 ± 0.12	90.13 ± 0.19	82.99 ± 0.36
JoCoR (Wei et al., 2020)	93.40 ± 0.24	91.44 ± 0.05	90.30 ± 0.20	90.21 ± 0.19	90.11 ± 0.21	83.37 ± 0.30
CORES ² (Cheng et al., 2021)	93.43 ± 0.24	91.23 ± 0.11	89.66 ± 0.32	89.91 ± 0.45	89.79 ± 0.50	83.60 ± 0.53
CORES* (Cheng et al., 2021)	94.16 ± 0.11	95.25 ± 0.09	94.45 ± 0.14	94.88 ± 0.31	94.74 ± 0.03	91.66 ± 0.09
VolMinNet (Li et al., 2021)	92.14 ± 0.30	89.70 ± 0.21	88.30 ± 0.12	88.27 ± 0.09	88.19 ± 0.41	80.53 ± 0.20
CAL (Zhu et al., 2021a)	94.50 ± 0.31	91.97 ± 0.32	90.93 ± 0.31	90.75 ± 0.30	90.74 ± 0.24	85.36 ± 0.16
PES (Semi) (Bai et al., 2021)	94.76 ± 0.2	94.66 ± 0.18	95.06 ± 0.15	95.19 ± 0.23	95.22 ± 0.13	92.68 ± 0.22
SOP (Liu et al., 2022)	N/A	95.61 ± 0.13	95.28 ± 0.13	95.31 ± 0.10	95.39 ± 0.11	93.24 ± 0.21

Sparse modeling gives super performance again label noise⁹

⁹Wei et al., Learning with noisy labels revisited: A study using real-world human annotations, ICLR 2022.

Conclusion and Coming Attractions

Learning common deep networks for low-dim structure

- **Low-dimensional data:** understand resource tradeoffs between data structure and network architecture
- **Low-dimensional features:** understand low-dim. features (NC) learned in deep classifiers trained with one-hot labeling based losses
- **Robust training:** Exploit low-dim structure in the label noise to improve training robustness

Next lecture: New approach for learning diverse and discriminative features (beyond NC).

Designing deep network architectures for low-dimensional structures

Thank You! Questions?

Figure Credits I

- Slide 3: Dictionary learning figures from [Mairal, Elad, and Sapiro 2008]
- Slide 4: ImageNet classes from paperswithcode.com; AlexNet architecture: [Krizhevsky et al. 2012]; ResNet architecture: [He et al. 2015];
- Slide 5: ImageNet top1 from paperswithcode.com; DALL-E 1 and 2 from <https://openai.com/blog/dall-e/> and <https://openai.com/dall-e-2/>
- Slide 7: Right image from <https://www.cityscapes-dataset.com/dataset-overview/>
- Slide 8: Hairbrushes from <https://objectnet.dev/download.html>
- Slide 9: Illumination figure from [Basri and Jacobs 2003]
- Slide 13: Left figure from [Krizhevsky et al. 2012]; right from <https://openai.com/blog/microscope/>;