

A Top-Down Approach Toward Understanding Deep Learning

Weijie Su

University of Pennsylvania

Otaru University of Commerce, November 4, 2021

A new paradigm of “science”: deep learning



- Collect data and buy GPU first

A new paradigm of “science”: deep learning



- Collect data and buy GPU first
- Scale model with data and computational resources

A new paradigm of “science”: deep learning



- Collect data and buy GPU first
- Scale model with data and computational resources
- End to end: Representation, computation, prediction

A new paradigm of “science”: deep learning



- Collect data and buy GPU first
- Scale model with data and computational resources
- End to end: Representation, computation, prediction

The Bitter Lesson

Rich Sutton

March 13, 2019

The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective, and by a large margin. The ultimate reason for this is Moore's law, or rather its generalization of continued exponentially falling cost per unit of computation. Most AI research has been conducted as if the computation available to the agent were constant (in which case leveraging human knowledge would be one of the only ways to improve performance) but, over a slightly longer time than a typical research project, massively more computation inevitably becomes available. Seeking an improvement that makes a difference in the shorter term, researchers seek to leverage their human knowledge of the domain, but the only thing that matters in the long run is the leveraging of computation. These two need not run counter to each other, but in practice they tend to. Time spent on one is time not spent on the other. There are psychological commitments to investment in one approach or the other. And the human-knowledge approach tends to complicate methods in ways that make them less suited to taking advantage of general methods leveraging computation. There were many examples of AI researchers' belated learning of this bitter lesson, and it is instructive to review some of the most prominent.

However, making deep learning a science requires...

- Why don't heavily parameterized neural networks overfit the data?

However, making deep learning a science requires...

- Why don't heavily parameterized neural networks overfit the data?
- What is the effective number of parameters?

However, making deep learning a science requires...

- Why don't heavily parameterized neural networks overfit the data?
- What is the effective number of parameters?
- Why doesn't backpropagation get stuck in poor local minima with low value of the loss function, yet bad test error?

However, making deep learning a science requires...

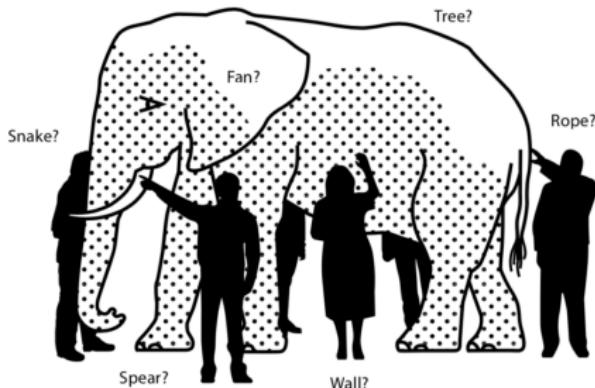
- Why don't heavily parameterized neural networks overfit the data?
- What is the effective number of parameters?
- Why doesn't backpropagation get stuck in poor local minima with low value of the loss function, yet bad test error?



Leo Breiman

Disclaimer: This talk doesn't attempt to answer these fundamental questions...

Have we really understood deep learning?



Limited scopes...

- Assume extremely large width and shallow depth
- Data assumed to be from Gaussian mixtures
- Linear activation
- Use gradient descent instead of stochastic gradient descent
-

Need “small” but useful surrogate models

A bitter lesson learned

Very difficult to build a comprehensive foundation for deep learning...

Need “small” but useful surrogate models

A bitter lesson learned

Very difficult to build a comprehensive foundation for deep learning...



Need “small” but useful surrogate models

A bitter lesson learned

Very difficult to build a comprehensive foundation for deep learning...



What is a good *surrogate* model?

- Mathematically tractable
- Yet maintains some characteristics of deep learning
- Insights into the practice of deep learning

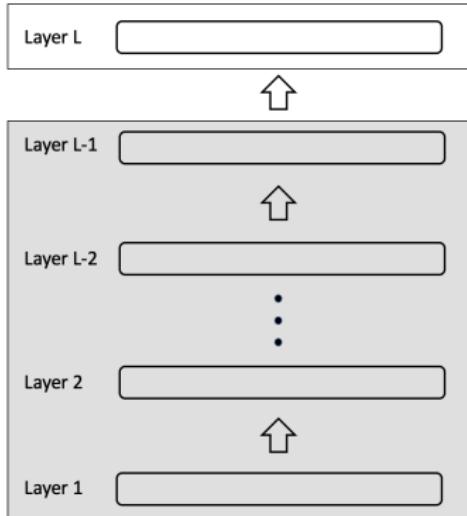
This talk: a top-down viewpoint



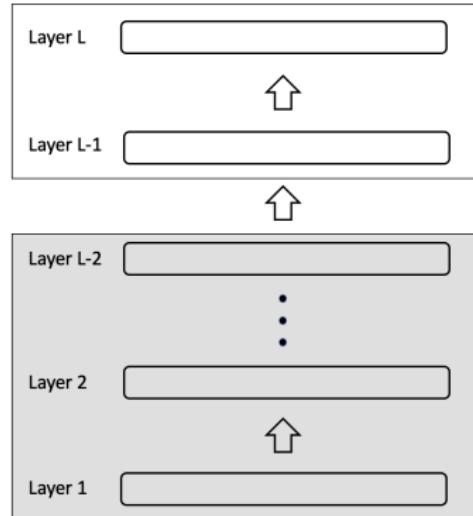
Collaborators

- Cong Fang (Penn CS)
- Hangfeng He (Penn CS)
- Qi Long (Penn Biostats)

Illustration of our top-down approach



(a) 1-Layer-Peeled Model



(b) 2-Layer-Peeled Model

Setup for deep learning

Neural network for K -class classification:

$$f(\mathbf{x}; \mathbf{W}_{\text{full}}) = \mathbf{W}_L \sigma (\mathbf{W}_{L-1} \sigma (\cdots \sigma (\mathbf{W}_1 \mathbf{x}) \cdots))$$

- $\sigma(\cdot)$ is a nonlinear activation function
- $\mathbf{W}_{\text{full}} := \{\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_L\}$ collects the weights
- Bias omitted

Setup for deep learning

Neural network for K -class classification:

$$f(\mathbf{x}; \mathbf{W}_{\text{full}}) = \mathbf{W}_L \sigma (\mathbf{W}_{L-1} \sigma (\cdots \sigma (\mathbf{W}_1 \mathbf{x}) \cdots))$$

- $\sigma(\cdot)$ is a nonlinear activation function
- $\mathbf{W}_{\text{full}} := \{\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_L\}$ collects the weights
- Bias omitted

Optimization problem:

$$\min_{\mathbf{W}_{\text{full}}} \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(f(\mathbf{x}_{k,i}; \mathbf{W}_{\text{full}}), \mathbf{y}_k) + \frac{\lambda}{2} \|\mathbf{W}_{\text{full}}\|^2$$

- \mathbf{y}_k is a one-hot vector denoting the k -th class
- λ weight decay parameter, \mathcal{L} cross-entropy loss

A peek at Layer-Peeled Model

$$f(\mathbf{x}; \mathbf{W}_{\text{full}}) = \mathbf{W}_L \sigma (\mathbf{W}_{L-1} \sigma (\cdots \sigma (\mathbf{W}_1 \mathbf{x}) \cdots))$$

$$\min_{\mathbf{W}_{\text{full}}} \quad \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(f(\mathbf{x}_{k,i}; \mathbf{W}_{\text{full}}), \mathbf{y}_k) + \frac{\lambda}{2} \|\mathbf{W}_{\text{full}}\|^2$$

- Difficult to pinpoint how any layer \mathbf{W}_l influences the output

A peek at Layer-Peeled Model

$$f(\mathbf{x}; \mathbf{W}_{\text{full}}) = \mathbf{W}_L \sigma (\mathbf{W}_{L-1} \sigma (\cdots \sigma (\mathbf{W}_1 \mathbf{x}) \cdots))$$

$$\min_{\mathbf{W}_L, \mathbf{H}} \quad \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}_L \mathbf{h}_{k,i}, \mathbf{y}_k) + \frac{\lambda}{2} \|\mathbf{W}_{\text{full}}\|^2$$

- Difficult to pinpoint how any layer \mathbf{W}_l influences the output
- $\mathbf{h}_{k,i}$ represents $\sigma (\mathbf{W}_{L-1} \sigma (\cdots \sigma (\mathbf{W}_1 \mathbf{x}_{k,i}) \cdots))$

A peek at Layer-Peeled Model

$$f(\mathbf{x}; \mathbf{W}_{\text{full}}) = \mathbf{W}_L \sigma(\mathbf{W}_{L-1} \sigma(\cdots \sigma(\mathbf{W}_1 \mathbf{x}) \cdots))$$

$$\begin{aligned} \min_{\mathbf{W}_L, \mathbf{H}} \quad & \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}_L \mathbf{h}_{k,i}, \mathbf{y}_k) \\ \text{s.t.} \quad & \frac{1}{K} \sum_{k=1}^K \|\mathbf{w}_k\|^2 \leq E_W \\ & \frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\mathbf{h}_{k,i}\|^2 \leq E_H \end{aligned}$$

- Difficult to pinpoint how any layer \mathbf{W}_l influences the output
- $\mathbf{h}_{k,i}$ represents $\sigma(\mathbf{W}_{L-1} \sigma(\cdots \sigma(\mathbf{W}_1 \mathbf{x}_{k,i}) \cdots))$
- Here $\mathbf{W}_L = [\mathbf{w}_1, \dots, \mathbf{w}_K]^\top$

Derivation

Rewrite the optimization problem as

$$\min_{\mathbf{W}_L, \mathbf{H}} \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}_L \mathbf{h}(\mathbf{x}_{k,i}; \mathbf{W}_{-L}), \mathbf{y}_k) + \frac{\lambda}{2} \|\mathbf{W}_L\|^2 + \frac{\lambda}{2} \|\mathbf{W}_{-L}\|^2$$

- Last-layer feature $\mathbf{h}(\mathbf{x}_{k,i}; \mathbf{W}_{-L}) := \sigma(\mathbf{W}_{L-1} \sigma(\cdots \sigma(\mathbf{W}_1 \mathbf{x}_{k,i}) \cdots))$

Derivation

Rewrite the optimization problem as

$$\min_{\mathbf{W}_L, \mathbf{H}} \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}_L \mathbf{h}(\mathbf{x}_{k,i}; \mathbf{W}_{-L}), \mathbf{y}_k) + \frac{\lambda}{2} \|\mathbf{W}_L\|^2 + \frac{\lambda}{2} \|\mathbf{W}_{-L}\|^2$$

- Last-layer feature $\mathbf{h}(\mathbf{x}_{k,i}; \mathbf{W}_{-L}) := \sigma(\mathbf{W}_{L-1} \sigma(\cdots \sigma(\mathbf{W}_1 \mathbf{x}_{k,i}) \cdots))$

From the *dual* viewpoint, a minimum is an optimal solution to

$$\begin{aligned} \min_{\mathbf{W}_L, \mathbf{H}} \quad & \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}_L \mathbf{h}_{k,i}, \mathbf{y}_k) \\ \text{s.t.} \quad & \|\mathbf{W}_L\|^2 \leq C_1 \\ & \|\mathbf{W}_{-L}\|^2 \leq C_2 \end{aligned}$$

Derivation

Rewrite the optimization problem as

$$\min_{\mathbf{W}_L, \mathbf{H}} \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}_L \mathbf{h}(\mathbf{x}_{k,i}; \mathbf{W}_{-L}), \mathbf{y}_k) + \frac{\lambda}{2} \|\mathbf{W}_L\|^2 + \frac{\lambda}{2} \|\mathbf{W}_{-L}\|^2$$

- Last-layer feature $\mathbf{h}(\mathbf{x}_{k,i}; \mathbf{W}_{-L}) := \sigma(\mathbf{W}_{L-1} \sigma(\cdots \sigma(\mathbf{W}_1 \mathbf{x}_{k,i}) \cdots))$

From the *dual* viewpoint, a minimum is an optimal solution to

$$\begin{aligned} \min_{\mathbf{W}_L, \mathbf{H}} \quad & \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}_L \mathbf{h}_{k,i}, \mathbf{y}_k) \\ \text{s.t.} \quad & \|\mathbf{W}_L\|^2 \leq C_1 \\ & \|\mathbf{W}_{-L}\|^2 \leq C_2 \end{aligned}$$

- Not a one-to-one mapping

Derivation

Rewrite the optimization problem as

$$\min_{\mathbf{W}_L, \mathbf{H}} \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}_L \mathbf{h}(\mathbf{x}_{k,i}; \mathbf{W}_{-L}), \mathbf{y}_k) + \frac{\lambda}{2} \|\mathbf{W}_L\|^2 + \frac{\lambda}{2} \|\mathbf{W}_{-L}\|^2$$

- Last-layer feature $\mathbf{h}(\mathbf{x}_{k,i}; \mathbf{W}_{-L}) := \sigma(\mathbf{W}_{L-1} \sigma(\cdots \sigma(\mathbf{W}_1 \mathbf{x}_{k,i}) \cdots))$

From the *dual* viewpoint, a minimum is an optimal solution to

$$\begin{aligned} \min_{\mathbf{W}_L, \mathbf{H}} \quad & \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}_L \mathbf{h}_{k,i}, \mathbf{y}_k) \\ \text{s.t.} \quad & \|\mathbf{W}_L\|^2 \leq C_1 \\ & \|\mathbf{W}_{-L}\|^2 \leq C_2 \Leftrightarrow \mathbf{H} \in \{\mathbf{H}(\mathbf{W}_{-L}) : \|\mathbf{W}_{-L}\|^2 \leq C_2\} \end{aligned}$$

- Not a one-to-one mapping
- $\mathbf{H}(\mathbf{W}_{-L}) := [\mathbf{h}(\mathbf{x}_{k,i}; \mathbf{W}_{-L}) : 1 \leq k \leq K, 1 \leq i \leq n_k]$

Derivation

Rewrite the optimization problem as

$$\min_{\mathbf{W}_L, \mathbf{H}} \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}_L \mathbf{h}(\mathbf{x}_{k,i}; \mathbf{W}_{-L}), \mathbf{y}_k) + \frac{\lambda}{2} \|\mathbf{W}_L\|^2 + \frac{\lambda}{2} \|\mathbf{W}_{-L}\|^2$$

- Last-layer feature $\mathbf{h}(\mathbf{x}_{k,i}; \mathbf{W}_{-L}) := \sigma(\mathbf{W}_{L-1} \sigma(\cdots \sigma(\mathbf{W}_1 \mathbf{x}_{k,i}) \cdots))$

From the *dual* viewpoint, a minimum is an optimal solution to

$$\begin{aligned} \min_{\mathbf{W}_L, \mathbf{H}} \quad & \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}_L \mathbf{h}_{k,i}, \mathbf{y}_k) \\ \text{s.t.} \quad & \|\mathbf{W}_L\|^2 \leq C_1 \\ & \mathbf{H} \in \{\mathbf{H}(\mathbf{W}_{-L}) : \|\mathbf{W}_{-L}\|^2 \leq C_2\} \end{aligned}$$

- Not a one-to-one mapping
- $\mathbf{H}(\mathbf{W}_{-L}) := [\mathbf{h}(\mathbf{x}_{k,i}; \mathbf{W}_{-L}) : 1 \leq k \leq K, 1 \leq i \leq n_k]$

Derivation: an *ansatz*

Assumption

$$\{\boldsymbol{H}(\boldsymbol{W}_{-L}) : \|\boldsymbol{W}_{-L}\|^2 \leq C_2\} \approx \left\{ \boldsymbol{H} : \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\boldsymbol{h}_{k,i}\|^2 \leq C'_2 \right\}$$

Derivation: an *ansatz*

Assumption

$$\{\boldsymbol{H}(\boldsymbol{W}_{-L}) : \|\boldsymbol{W}_{-L}\|^2 \leq C_2\} \approx \left\{ \boldsymbol{H} : \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\boldsymbol{h}_{k,i}\|^2 \leq C'_2 \right\}$$

$$\begin{aligned} & \min_{\boldsymbol{W}_L, \boldsymbol{H}} \quad \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\boldsymbol{W}_L \boldsymbol{h}_{k,i}, \boldsymbol{y}_k) \\ \text{s.t.} \quad & \|\boldsymbol{W}_L\|^2 \leq C_1 \\ & \boldsymbol{H} \in \{\boldsymbol{H}(\boldsymbol{W}_{-L}) : \|\boldsymbol{W}_{-L}\|^2 \leq C_2\} \end{aligned}$$

Derivation: an *ansatz*

Assumption

$$\{\mathbf{H}(\mathbf{W}_{-L}) : \|\mathbf{W}_{-L}\|^2 \leq C_2\} \approx \left\{ \mathbf{H} : \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\mathbf{h}_{k,i}\|^2 \leq C'_2 \right\}$$

$$\begin{aligned} & \min_{\mathbf{W}_L, \mathbf{H}} \quad \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}_L \mathbf{h}_{k,i}, \mathbf{y}_k) \\ \text{s.t.} \quad & \|\mathbf{W}_L\|^2 \leq C_1 \\ & \mathbf{H} \in \{\mathbf{H}(\mathbf{W}_{-L}) : \|\mathbf{W}_{-L}\|^2 \leq C_2\} \end{aligned}$$

$$\begin{aligned} & \min_{\mathbf{W}, \mathbf{H}} \quad \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W} \mathbf{h}_{k,i}, \mathbf{y}_k) \\ \text{s.t.} \quad & \frac{1}{K} \sum_{k=1}^K \|\mathbf{w}_k\|^2 \leq E_W \\ & \frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\mathbf{h}_{k,i}\|^2 \leq E_H \end{aligned}$$

- Self-duality of ℓ_2 spaces
- More justification for the ansatz later

More on Layer-Peeled Model

$$\min_{\mathbf{W}, \mathbf{H}} \quad \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}\mathbf{h}_{k,i}, \mathbf{y}_k)$$

$$\text{s.t.} \quad \frac{1}{K} \sum_{k=1}^K \|\mathbf{w}_k\|^2 \leq E_W$$

Prediction constraint

$$\frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\mathbf{h}_{k,i}\|^2 \leq E_H$$

Representation constraint



More on Layer-Peeled Model

$$\min_{\mathbf{W}, \mathbf{H}} \quad \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}\mathbf{h}_{k,i}, \mathbf{y}_k)$$

$$\text{s.t.} \quad \frac{1}{K} \sum_{k=1}^K \|\mathbf{w}_k\|^2 \leq E_W$$

Prediction constraint

Representation constraint

$$\frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\mathbf{h}_{k,i}\|^2 \leq E_H$$

- Terminal phase of deep learning training



More on Layer-Peeled Model

$$\min_{\mathbf{W}, \mathbf{H}} \quad \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}\mathbf{h}_{k,i}, \mathbf{y}_k)$$

$$\text{s.t.} \quad \frac{1}{K} \sum_{k=1}^K \|\mathbf{w}_k\|^2 \leq E_W$$

Prediction constraint

$$\frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\mathbf{h}_{k,i}\|^2 \leq E_H$$

Representation constraint

- Terminal phase of deep learning training
- E_W, E_H depend on weight decay λ



More on Layer-Peeled Model

$$\min_{\mathbf{W}, \mathbf{H}} \quad \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}\mathbf{h}_{k,i}, \mathbf{y}_k)$$

$$\text{s.t.} \quad \frac{1}{K} \sum_{k=1}^K \|\mathbf{w}_k\|^2 \leq E_W$$

Prediction constraint

Representation constraint

$$\frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\mathbf{h}_{k,i}\|^2 \leq E_H$$

- Terminal phase of deep learning training
- E_W, E_H depend on weight decay λ
- Nonconvex but analytically tractable



More on Layer-Peeled Model

$$\min_{\mathbf{W}, \mathbf{H}} \quad \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}\mathbf{h}_{k,i}, \mathbf{y}_k)$$

$$\text{s.t.} \quad \frac{1}{K} \sum_{k=1}^K \|\mathbf{w}_k\|^2 \leq E_W$$

Prediction constraint

Representation constraint

$$\frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\mathbf{h}_{k,i}\|^2 \leq E_H$$

- Terminal phase of deep learning training
- E_W, E_H depend on weight decay λ
- Nonconvex but analytically tractable
- Does not explicitly depend on data



More on Layer-Peeled Model

$$\min_{\mathbf{W}, \mathbf{H}} \quad \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}\mathbf{h}_{k,i}, \mathbf{y}_k)$$

$$\text{s.t.} \quad \frac{1}{K} \sum_{k=1}^K \|\mathbf{w}_k\|^2 \leq E_W$$

Prediction constraint

Representation constraint

$$\frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\mathbf{h}_{k,i}\|^2 \leq E_H$$

- Terminal phase of deep learning training
- E_W, E_H depend on weight decay λ
- Nonconvex but analytically tractable
- Does not explicitly depend on data
 - Cons: information lost



More on Layer-Peeled Model

$$\min_{\mathbf{W}, \mathbf{H}} \quad \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}\mathbf{h}_{k,i}, \mathbf{y}_k)$$

$$\text{s.t.} \quad \frac{1}{K} \sum_{k=1}^K \|\mathbf{w}_k\|^2 \leq E_W$$

Prediction constraint

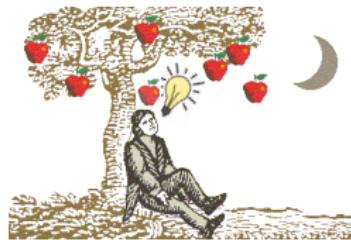
Representation constraint

$$\frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\mathbf{h}_{k,i}\|^2 \leq E_H$$

- Terminal phase of deep learning training
- E_W, E_H depend on weight decay λ
- Nonconvex but analytically tractable
- Does not explicitly depend on data
 - Cons: information lost
 - Pros: robust conclusion



Ask me anything about this “apple”

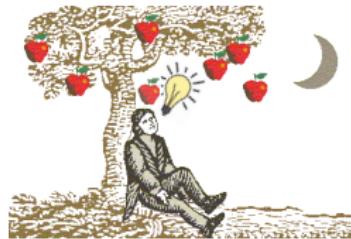


Ask me anything about this “apple”



Is it mathematically tractable?

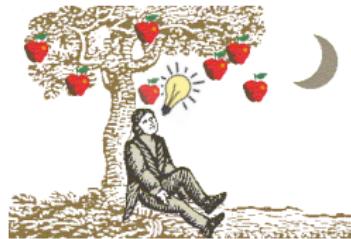
Ask me anything about this “apple”



Is it mathematically tractable?

Yes

Ask me anything about this “apple”



Is it mathematically tractable?

Yes

Does it maintain some characteristics of deep learning?

Ask me anything about this “apple”



Is it mathematically tractable?

Yes

Does it maintain some characteristics of deep learning?

Yes

Ask me anything about this “apple”



Is it mathematically tractable?

Yes

Does it maintain some characteristics of deep learning?

Yes

Can it provide insights into the practice of deep learning?

Ask me anything about this “apple”



Is it mathematically tractable?

Yes

Does it maintain some characteristics of deep learning?

Yes

Can it provide insights into the practice of deep learning?

I think so

Ask me anything about this “apple”



Is it mathematically tractable?

Yes

Does it maintain some characteristics of deep learning?

Yes

Can it provide insights into the practice of deep learning?

I think so

Does it answer Leo Breiman's questions?

Ask me anything about this “apple”



Is it mathematically tractable?

Yes

Does it maintain some characteristics of deep learning?

Yes

Can it provide insights into the practice of deep learning?

I think so

Does it answer Leo Breiman's questions?

Unfortunately, not

Outline

1. Explaining Neural Collapse
2. Predicting Minority Collapse
3. How to Mitigate Minority Collapse?

Balanced training

All class sizes are equal: $n_1 = n_2 = \dots = n_K$

Balanced training

All class sizes are equal: $n_1 = n_2 = \dots = n_K$

What can the Layer-Peeled Model say?

Balanced training

All class sizes are equal: $n_1 = n_2 = \dots = n_K$

What can the Layer-Peeled Model say?

Theorem

Any global minimizer $\mathbf{W}^* \equiv [\mathbf{w}_1^*, \dots, \mathbf{w}_K^*]^\top$, $\mathbf{H}^* \equiv [\mathbf{h}_{k,i}^* : 1 \leq k \leq K, 1 \leq i \leq n]$ with cross-entropy loss obeys

$$\mathbf{h}_{k,i}^* = C\mathbf{w}_k^* = C'\mathbf{m}_k^*,$$

where $[\mathbf{m}_1^*, \dots, \mathbf{m}_K^*]$ forms a K -simplex equiangular tight frame (ETF)

- $\mathbf{h}_{k,i}^*$ depends only on the class membership!
- $C = \sqrt{E_H/E_W}, C' = \sqrt{E_H}$

Balanced training

All class sizes are equal: $n_1 = n_2 = \dots = n_K$

What can the Layer-Peeled Model say?

Theorem

Any global minimizer $\mathbf{W}^* \equiv [\mathbf{w}_1^*, \dots, \mathbf{w}_K^*]^\top$, $\mathbf{H}^* \equiv [\mathbf{h}_{k,i}^* : 1 \leq k \leq K, 1 \leq i \leq n]$ with cross-entropy loss obeys

$$\mathbf{h}_{k,i}^* = C\mathbf{w}_k^* = C'\mathbf{m}_k^*,$$

where $[\mathbf{m}_1^*, \dots, \mathbf{m}_K^*]$ forms a K -simplex equiangular tight frame (ETF)

- $\mathbf{h}_{k,i}^*$ depends only on the class membership!
- $C = \sqrt{E_H/E_W}, C' = \sqrt{E_H}$
- What is a K -simplex ETF?

K-simplex ETF

K equal-length vectors form the *largest* possible equal-sized angles between any pair

Equivalently, random variables ξ_1, \dots, ξ_K of mean 0 and variance 1. If $\mathbb{E}\xi_i\xi_j = \rho$ for all $i \neq j$, what's the min of ρ ?

K -simplex ETF

K equal-length vectors form the *largest* possible equal-sized angles between any pair

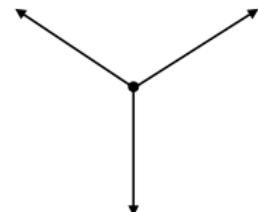
Equivalently, random variables ξ_1, \dots, ξ_K of mean 0 and variance 1. If $\mathbb{E}\xi_i\xi_j = \rho$ for all $i \neq j$, what's the min of ρ ?

$$\text{largest angle} = \arccos\left(-\frac{1}{K-1}\right)$$

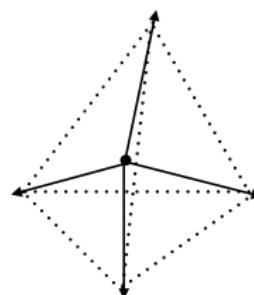
$K = 2$



$K = 3$



$K = 4$



Return to the theorem for balanced training

All class sizes are equal: $n_1 = n_2 = \dots = n_K$

Theorem

The solution to the Layer-Peeled Model in balanced training satisfies

$$\mathbf{h}_{k,i}^* = C\mathbf{w}_k^* = C'\mathbf{m}_k^*$$

- German shepherd, husky, chihuahua, rottweiler are all dogs!

Return to the theorem for balanced training

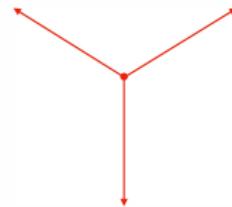
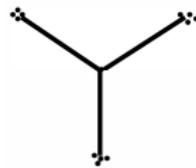
All class sizes are equal: $n_1 = n_2 = \dots = n_K$

Theorem

The solution to the Layer-Peeled Model in balanced training satisfies

$$\mathbf{h}_{k,i}^* = C\mathbf{w}_k^* = C'\mathbf{m}_k^*$$

- German shepherd, husky, chihuahua, rottweiler are all dogs!



Return to the theorem for balanced training

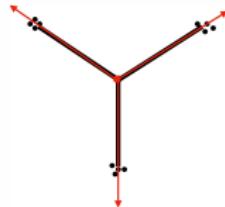
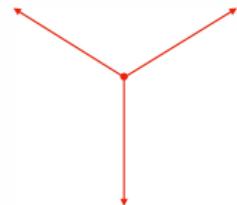
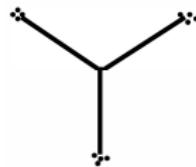
All class sizes are equal: $n_1 = n_2 = \dots = n_K$

Theorem

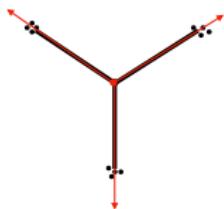
The solution to the Layer-Peeled Model in balanced training satisfies

$$\mathbf{h}_{k,i}^* = C\mathbf{w}_k^* = C'\mathbf{m}_k^*$$

- German shepherd, husky, chihuahua, rottweiler are all dogs!

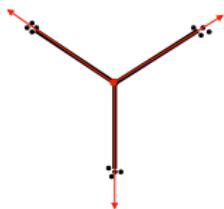


This is simply neural collapse



Papyan, Han, and Donoho discovered *neural collapse* in 2020:

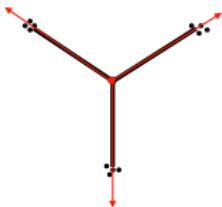
This is simply neural collapse



Papyan, Han, and Donoho discovered *neural collapse* in 2020:

- ① Variability collapse: features collapse to their class means

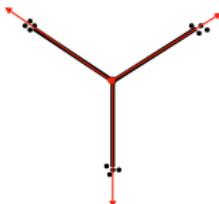
This is simply neural collapse



Papyan, Han, and Donoho discovered *neural collapse* in 2020:

- ① Variability collapse: features collapse to their class means
- ② Class means centered at their global mean collapse to ETF

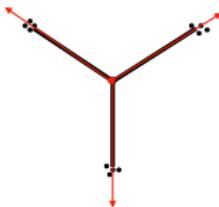
This is simply neural collapse



Papyan, Han, and Donoho discovered *neural collapse* in 2020:

- ① Variability collapse: features collapse to their class means
- ② Class means centered at their global mean collapse to ETF
- ③ Up to scaling, last-layer classifiers each collapse to class means

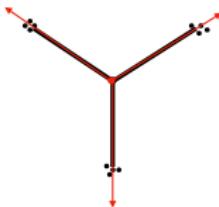
This is simply neural collapse



Papyan, Han, and Donoho discovered *neural collapse* in 2020:

- ① Variability collapse: features collapse to their class means
- ② Class means centered at their global mean collapse to ETF
- ③ Up to scaling, last-layer classifiers each collapse to class means
- ④ Classifier's decision collapses to choosing the closest class mean

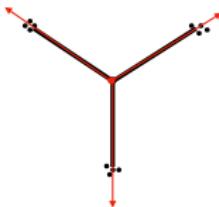
This is simply neural collapse



Papyan, Han, and Donoho discovered *neural collapse* in 2020:

- ① Variability collapse: features collapse to their class means
- ② Class means centered at their global mean collapse to ETF
- ③ Up to scaling, last-layer classifiers each collapse to class means
- ④ Classifier's decision collapses to choosing the closest class mean

This is simply neural collapse



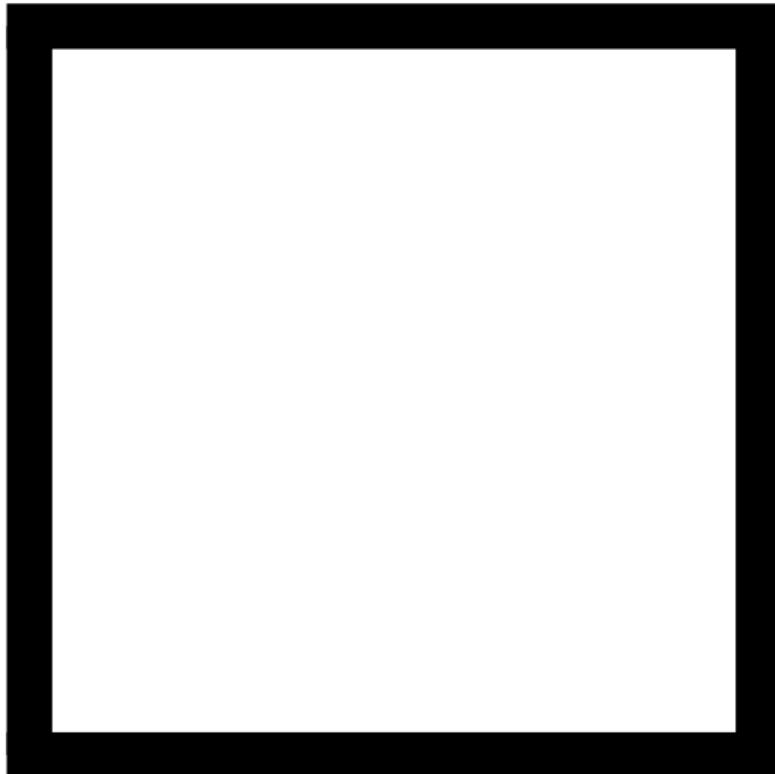
Papyan, Han, and Donoho discovered *neural collapse* in 2020:

- ① Variability collapse: features collapse to their class means
- ② Class means centered at their global mean collapse to ETF
- ③ Up to scaling, last-layer classifiers each collapse to class means
- ④ Classifier's decision collapses to choosing the closest class mean

Implications on better generalization, large margin, and robustness

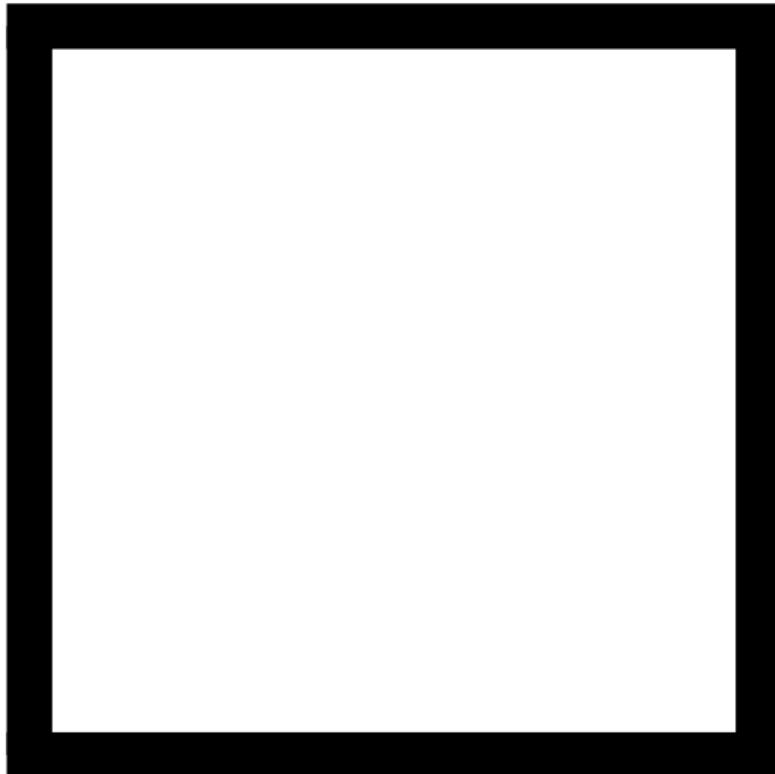
Concurrent works [MPP20, EW20, LS20] also justified neural collapse using different models

Animation of neural collapse



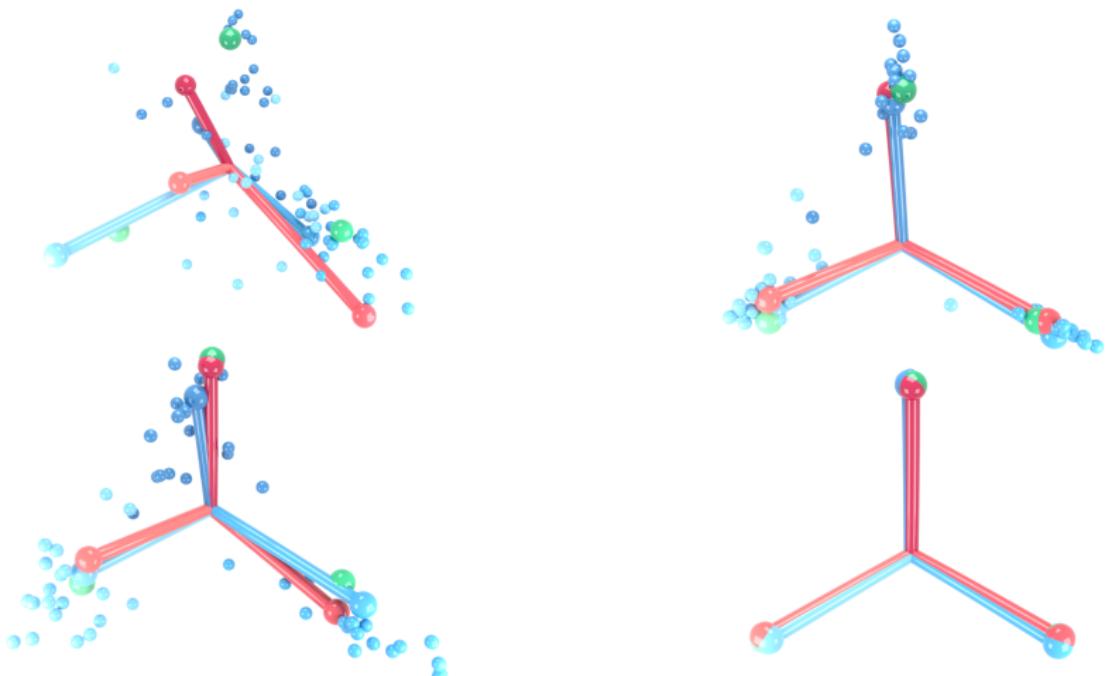
Credit: Popyan, Han, and Donoho

Animation of neural collapse



Credit: Popyan, Han, and Donoho

Snapshot of neural collapse



Credit: Popyan, Han, and Donoho

Neural collapse can justify the Layer-Peeled Model

About the ansatz

Recall

$$\{\mathbf{H}(\mathbf{W}_{-L}) : \|\mathbf{W}_{-L}\|^2 \leq C_2\} \approx \left\{ \mathbf{H} : \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\mathbf{h}_{k,i}\|^2 \leq C'_2 \right\}$$

This gives

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{H}} \quad & \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W} \mathbf{h}_{k,i}, \mathbf{y}_k) \\ \text{s.t.} \quad & \frac{1}{K} \sum_{k=1}^K \|\mathbf{w}_k\|^2 \leq E_W \\ & \frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\mathbf{h}_{k,i}\|^2 \leq E_H \end{aligned}$$

What happens without the ansatz?

Without the ansatz:

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{H}} \quad & \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}(\mathbf{W}\mathbf{h}_{k,i}, \mathbf{y}_k) \\ \text{s.t.} \quad & \frac{1}{K} \sum_{k=1}^K \|\mathbf{w}_k\|^2 \leq E_W \\ & \frac{1}{K} \sum_{k=1}^K \frac{1}{n} \sum_{i=1}^n \|\mathbf{h}_{k,i}\|_q^q \leq E_H \end{aligned}$$

What happens without the ansatz?

Without the ansatz:

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{H}} \quad & \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}(\mathbf{W}\mathbf{h}_{k,i}, \mathbf{y}_k) \\ \text{s.t.} \quad & \frac{1}{K} \sum_{k=1}^K \|\mathbf{w}_k\|^2 \leq E_W \\ & \frac{1}{K} \sum_{k=1}^K \frac{1}{n} \sum_{i=1}^n \|\mathbf{h}_{k,i}\|_q^q \leq E_H \end{aligned}$$

Theorem

Assume $K \geq 3$ and $p \geq K$. For any $q \in (0, 2) \cup (2, \infty)$, neural collapse does **not** emerge in the model above

Is the Layer-Peeled Model satisfactory?

*Is the Layer-Peeled Model satisfactory?
A higher standard: can it predict new stuff?*

Outline

1. Explaining Neural Collapse
2. Predicting Minority Collapse
3. How to Mitigate Minority Collapse?

Imbalanced training

Datasets often have a disproportionate ratio of observations in each class

Imbalanced training

Datasets often have a disproportionate ratio of observations in each class

As a simple starting point, assume

- The first K_A majority classes each contain n_A training examples
 $(n_1 = n_2 = \dots = n_{K_A} = n_A)$

Imbalanced training

Datasets often have a disproportionate ratio of observations in each class

As a simple starting point, assume

- The first K_A majority classes each contain n_A training examples
($n_1 = n_2 = \dots = n_{K_A} = n_A$)
- The remaining $K_B := K - K_A$ minority classes each contain n_B examples
($n_{K_A+1} = n_{K_A+2} = \dots = n_K = n_B$)

Imbalanced training

Datasets often have a disproportionate ratio of observations in each class

As a simple starting point, assume

- The first K_A majority classes each contain n_A training examples
($n_1 = n_2 = \dots = n_{K_A} = n_A$)
- The remaining $K_B := K - K_A$ minority classes each contain n_B examples
($n_{K_A+1} = n_{K_A+2} = \dots = n_K = n_B$)
- Call $R := n_A/n_B > 1$ the imbalance ratio

*No closed-form expression for
the solutions to LPM...*

Technique: Convex relaxation

- Define \mathbf{h}_k as the feature mean of the k -th class

$$\mathbf{h}_k := \frac{1}{n_k} \sum_{i=1}^{n_k} \mathbf{h}_{k,i}$$

- Introduce a new decision variable

$$\mathbf{X} := [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K, \mathbf{W}^\top]^\top [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K, \mathbf{W}^\top] \in \mathbb{R}^{2K \times 2K}$$

Technique: Convex relaxation

- Define \mathbf{h}_k as the feature mean of the k -th class

$$\mathbf{h}_k := \frac{1}{n_k} \sum_{i=1}^{n_k} \mathbf{h}_{k,i}$$

- Introduce a new decision variable

$$\mathbf{X} := [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K, \mathbf{W}^\top]^\top [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K, \mathbf{W}^\top] \in \mathbb{R}^{2K \times 2K}$$

Then

- \mathbf{X} is positive semidefinite

Technique: Convex relaxation

- Define \mathbf{h}_k as the feature mean of the k -th class

$$\mathbf{h}_k := \frac{1}{n_k} \sum_{i=1}^{n_k} \mathbf{h}_{k,i}$$

- Introduce a new decision variable

$$\mathbf{X} := [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K, \mathbf{W}^\top]^\top [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K, \mathbf{W}^\top] \in \mathbb{R}^{2K \times 2K}$$

Then

- \mathbf{X} is positive semidefinite
-

$$\frac{1}{K} \sum_{k=1}^K \mathbf{X}(k, k) = \frac{1}{K} \sum_{k=1}^K \|\mathbf{h}_k\|^2 \leq \frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\mathbf{h}_{k,i}\|^2 \leq E_H$$

Technique: Convex relaxation

- Define \mathbf{h}_k as the feature mean of the k -th class

$$\mathbf{h}_k := \frac{1}{n_k} \sum_{i=1}^{n_k} \mathbf{h}_{k,i}$$

- Introduce a new decision variable

$$\mathbf{X} := [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K, \mathbf{W}^\top]^\top [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K, \mathbf{W}^\top] \in \mathbb{R}^{2K \times 2K}$$

Then

- \mathbf{X} is positive semidefinite
-

$$\frac{1}{K} \sum_{k=1}^K \mathbf{X}(k, k) = \frac{1}{K} \sum_{k=1}^K \|\mathbf{h}_k\|^2 \leq \frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\mathbf{h}_{k,i}\|^2 \leq E_H$$

-

$$\frac{1}{K} \sum_{k=K+1}^{2K} \mathbf{X}(k, k) = \frac{1}{K} \sum_{k=1}^K \|\mathbf{w}_k\|^2 \leq E_W$$

Technique: Convex relaxation

$$\begin{aligned} \min_{\boldsymbol{X} \in \mathbb{R}^{2K \times 2K}} \quad & \sum_{k=1}^K \frac{n_k}{N} \mathcal{L}(\boldsymbol{z}_k, \boldsymbol{y}_k) \\ \text{s.t.} \quad & \boldsymbol{z}_k = [\boldsymbol{X}(k, K+1), \boldsymbol{X}(k, K+2), \dots, \boldsymbol{X}(k, 2K)]^\top \\ & \frac{1}{K} \sum_{k=1}^K \boldsymbol{X}(k, k) \leq E_H, \quad \frac{1}{K} \sum_{k=K+1}^{2K} \boldsymbol{X}(k, k) \leq E_W \\ & \boldsymbol{X} \succeq 0 \end{aligned}$$

Technique: Convex relaxation

$$\begin{aligned} \min_{\mathbf{X} \in \mathbb{R}^{2K \times 2K}} \quad & \sum_{k=1}^K \frac{n_k}{N} \mathcal{L}(\mathbf{z}_k, \mathbf{y}_k) \\ \text{s.t.} \quad & \mathbf{z}_k = [\mathbf{X}(k, K+1), \mathbf{X}(k, K+2), \dots, \mathbf{X}(k, 2K)]^\top \\ & \frac{1}{K} \sum_{k=1}^K \mathbf{X}(k, k) \leq E_H, \quad \frac{1}{K} \sum_{k=K+1}^{2K} \mathbf{X}(k, k) \leq E_W \\ & \mathbf{X} \succeq 0 \end{aligned}$$

- Not a semidefinite program in the strict sense because a semidefinite program uses a linear objective function

Nonconvex optimization via convex optimization

Lemma

Assume $p \geq 2K$ and \mathcal{L} is convex in its first argument. Let \mathbf{X}^* be a minimizer of the convex relaxation. Define $(\mathbf{H}^*, \mathbf{W}^*)$ as

$$\begin{aligned} [\mathbf{h}_1^*, \mathbf{h}_2^*, \dots, \mathbf{h}_K^*, (\mathbf{W}^*)^\top] &= \mathbf{P}(\mathbf{X}^*)^{1/2} \\ \mathbf{h}_{k,i}^* &= \mathbf{h}_k^*, \text{ for all } 1 \leq i \leq n, 1 \leq k \leq K \end{aligned}$$

Then $(\mathbf{H}^*, \mathbf{W}^*)$ is a minimizer of the Layer-Peeled Model

Nonconvex optimization via convex optimization

Lemma

Assume $p \geq 2K$ and \mathcal{L} is convex in its first argument. Let \mathbf{X}^* be a minimizer of the convex relaxation. Define $(\mathbf{H}^*, \mathbf{W}^*)$ as

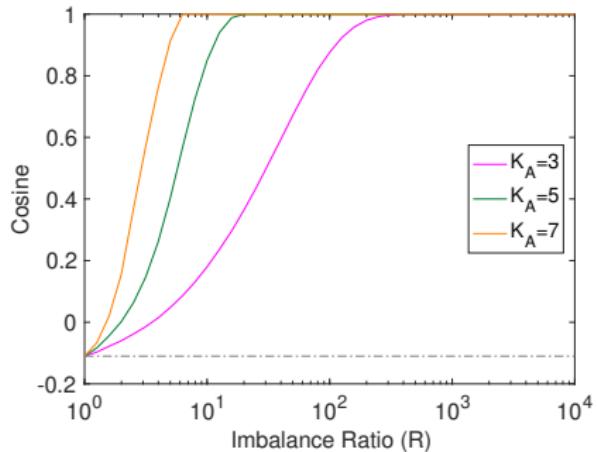
$$\begin{aligned} [\mathbf{h}_1^*, \mathbf{h}_2^*, \dots, \mathbf{h}_K^*, (\mathbf{W}^*)^\top] &= \mathbf{P}(\mathbf{X}^*)^{1/2} \\ \mathbf{h}_{k,i}^* &= \mathbf{h}_k^*, \text{ for all } 1 \leq i \leq n, 1 \leq k \leq K \end{aligned}$$

Then $(\mathbf{H}^*, \mathbf{W}^*)$ is a minimizer of the Layer-Peeled Model

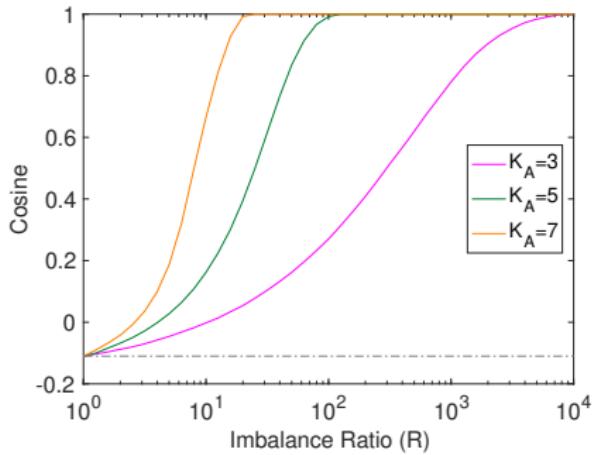
- No loss of information when we study the Layer-Peeled Model through a convex program
- But class means *no longer* collapse to classifiers
- Alternatives of convex relaxation exist [BMPO8, HV19]

A numerical surprise

Average cosine of between-minority-class angles



(c) $E_W = 1, E_H = 5$



(d) $E_W = 1, E_H = 10$

- ① When $R < R_0$ for some $R_0 > 0$, average between-minority-class angle becomes smaller as R increases
- ② Once $R \geq R_0$, average between-minority-class angle becomes 0: implying that all minority classifiers collapse!

Minority Collapse

- ① When $R < R_0$ for some $R_0 > 0$, average between-minority-class angle becomes smaller as R increases
- ② Once $R \geq R_0$, average between-minority-class angle becomes $\mathbf{0}$: implying that all minority classifiers collapse!

Proposition

Let $(\mathbf{H}^*, \mathbf{W}^*)$ be any global minimizer of the Layer-Peeled Model. As $R \equiv n_A/n_B \rightarrow \infty$, we have

$$\lim \mathbf{w}_k^* - \mathbf{w}_{k'}^* = \mathbf{0}_p, \text{ for all } K_A < k < k' \leq K$$

- The prediction on the minority classes becomes *completely at random*

Illustration of Minority Collapse

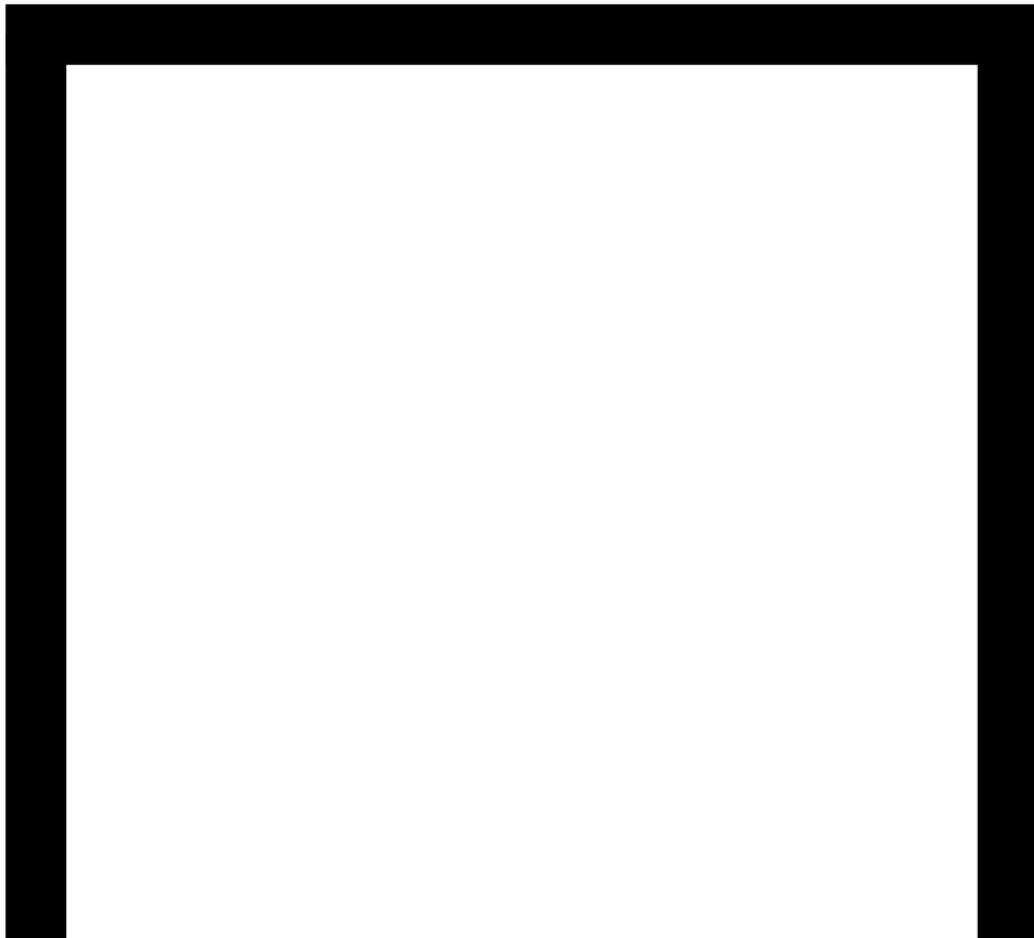
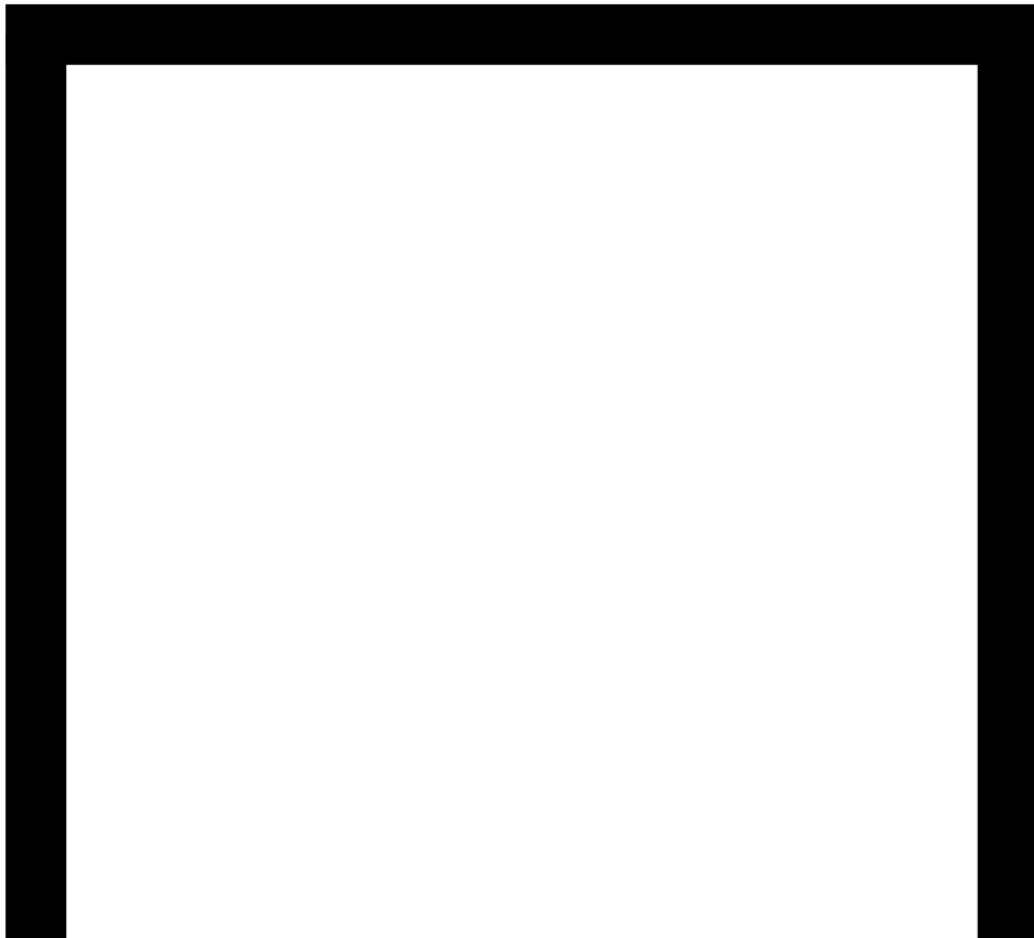


Illustration of Minority Collapse



Intuition for Minority Collapse

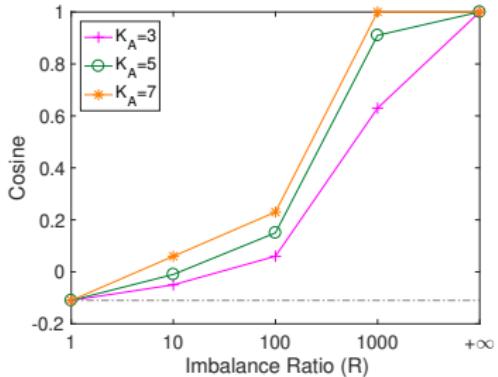
$$\begin{aligned} \min_{\mathbf{W}, \mathbf{H}} \quad & \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}\mathbf{h}_{k,i}, \mathbf{y}_k) \\ \text{s.t.} \quad & \frac{1}{K} \sum_{k=1}^K \|\mathbf{w}_k\|^2 \leq E_W \\ & \frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\mathbf{h}_{k,i}\|^2 \leq E_H \end{aligned}$$



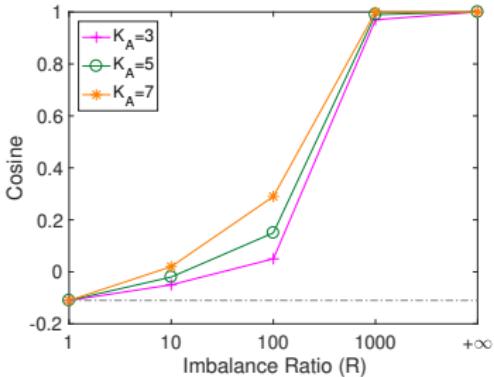
Competition for space!

Is Minority Collapse a real thing?

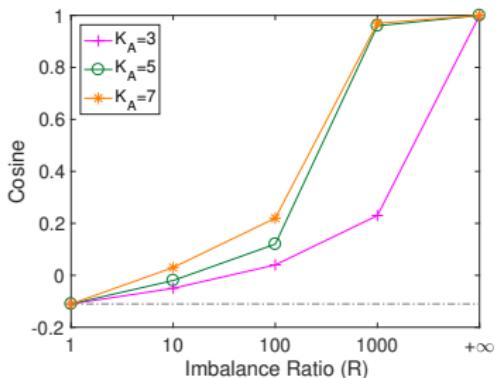
Minority Collapse in experiments



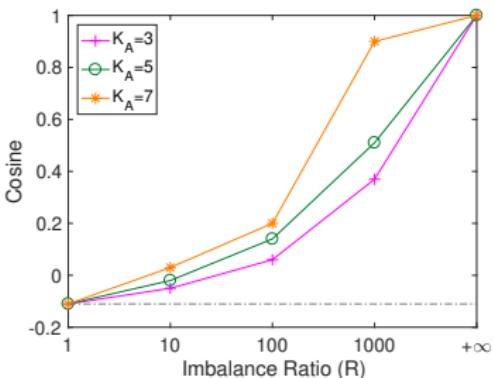
(e) VGG11 on FashionMNIST



(f) VGG13 on CIFAR10

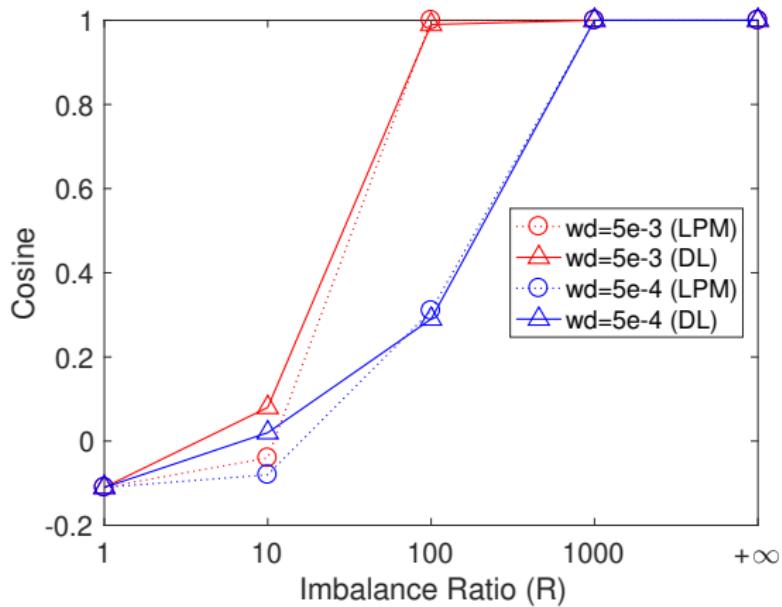


(g) ResNet18 on FashionMNIST



(h) ResNet18 on CIFAR10

LPM predictions match experiments



Layer-Peeled Model (LPM, in dotted lines) and real DNNs (DL, in solid lines) with VGG on CIFAR10

Outline

1. Explaining Neural Collapse
2. Predicting Minority Collapse
3. How to Mitigate Minority Collapse?

Idea: make the minority stronger!

Oversample minority classes

Oversampling duplicates training example from minority classes [JK09]

Oversample minority classes

Oversampling duplicates training example from minority classes [JK09]

The adjusted optimization problem:

$$\frac{1}{n_A K_A + w_r n_B K_B} \left[\sum_{k=1}^{K_A} \sum_{i=1}^{n_A} \mathcal{L}(\mathbf{f}(\mathbf{x}_{k,i}; \mathbf{W}_{\text{full}}), \mathbf{y}_k) \right. \\ \left. + w_r \sum_{k=K_A+1}^K \sum_{i=1}^{n_B} \mathcal{L}(\mathbf{f}(\mathbf{x}_{k,i}; \mathbf{W}_{\text{full}}), \mathbf{y}_k) \right]$$

while keeping the penalty term $\frac{\lambda}{2} \|\mathbf{W}_{\text{full}}\|^2$

Layer-Peeled Model with oversampling

$$\begin{aligned} \min_{\mathbf{H}, \mathbf{W}} \quad & \frac{1}{n_A K_A + w_r n_B K_B} \left[\sum_{k=1}^{K_A} \sum_{i=1}^{n_A} \mathcal{L}(\mathbf{W} \mathbf{h}_{k,i}, \mathbf{y}_k) + w_r \sum_{k=K_A+1}^K \sum_{i=1}^{n_B} \mathcal{L}(\mathbf{W} \mathbf{h}_{k,i}, \mathbf{y}_k) \right] \\ \text{s.t.} \quad & \frac{1}{K} \sum_{k=1}^K \|\mathbf{w}_k\|^2 \leq E_W \\ & \frac{1}{K} \sum_{k=1}^{K_A} \frac{1}{n_A} \sum_{i=1}^{n_A} \|\mathbf{h}_{k,i}\|^2 + \frac{1}{K} \sum_{k=K_A+1}^K \frac{1}{n_B} \sum_{i=1}^{n_B} \|\mathbf{h}_{k,i}\|^2 \leq E_H \end{aligned}$$

Layer-Peeled Model with oversampling

Theorem

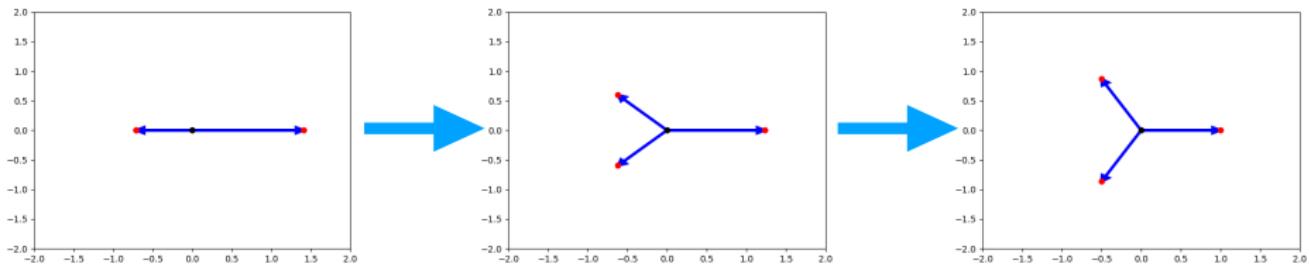
Assume $p \geq 2K$ and \mathcal{L} is convex in the first argument. Let \mathbf{X}^* be any minimizer of the convex relaxation with $n_1 = n_2 = \dots = n_{K_A} = n_A$ and $n_{K_A+1} = n_{K_A+2} = \dots = n_K = w_r n_B$. Define $(\mathbf{H}^*, \mathbf{W}^*)$ as

$$\begin{aligned} [\mathbf{h}_1^*, \mathbf{h}_2^*, \dots, \mathbf{h}_K^*, (\mathbf{W}^*)^\top] &= \mathbf{P}(\mathbf{X}^*)^{1/2} \\ \mathbf{h}_{k,i}^* &= \mathbf{h}_k^*, \text{ for all } 1 \leq i \leq n_A, 1 \leq k \leq K_A \\ \mathbf{h}_{k,i}^* &= \mathbf{h}_k^*, \text{ for all } 1 \leq i \leq n_B, K_A < k \leq K \end{aligned}$$

Then $(\mathbf{H}^*, \mathbf{W}^*)$ is a global minimizer of the oversampling-adjusted Layer-Peeled Model.

- The size of minority class is now in effect $w_r n_B$ instead of n_B
- If the oversampling rate $w_r = n_A/n_B \equiv R$, neural collapse is back!

Effect of oversampling, in theory

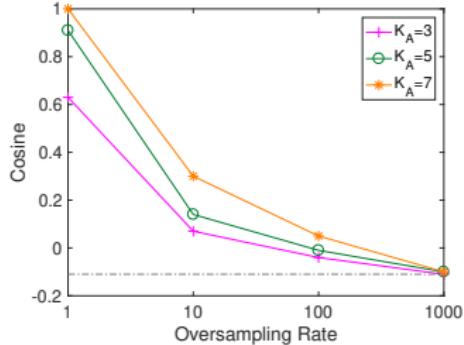




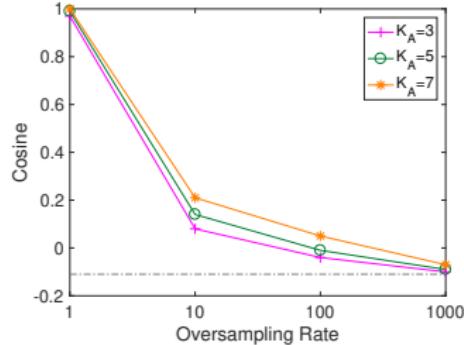
Can oversampling really resolve Minority Collapse?



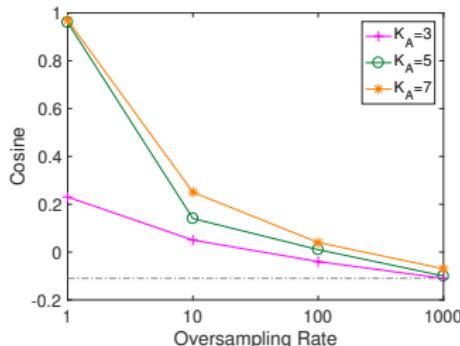
Oversampling mitigates Minority Collapse



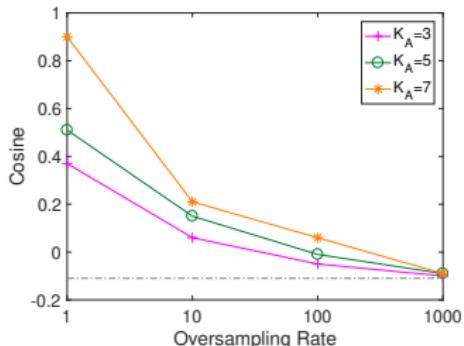
(i) VGG11 on FashionMNIST



(j) VGG13 on CIFAR10



(k) ResNet18 on FashionMNIST



(l) ResNet18 on CIFAR10

Test performance

Network architecture	VGG11			ResNet18		
No. of majority classes	$K_A = 3$	$K_A = 5$	$K_A = 7$	$K_A = 3$	$K_A = 5$	$K_A = 7$
Original (minority)	15.29	20.30	17.00	30.66	34.26	5.53
Oversampling (minority)	41.13	57.22	30.50	37.86	53.46	8.13
Improvement (minority)	25.84	36.92	13.50	7.20	19.20	2.60
Original (overall)	40.10	57.61	69.09	50.88	64.89	66.13
Oversampling (overall)	58.25	76.17	73.37	55.91	74.56	67.10
Improvement (overall)	18.15	18.56	4.28	5.03	9.67	0.97

Table: Test accuracy (%) on FashionMNIST when $R = 1000$. "Original (minority)" means that the test accuracy is evaluated only on the minority classes and oversampling is not used. When oversampling is used, we report the **best** test accuracy among four oversampling rates: 1, 10, 100, and 1000.

Test performance

Network architecture	VGG11			ResNet18		
	$K_A = 3$	$K_A = 5$	$K_A = 7$	$K_A = 3$	$K_A = 5$	$K_A = 7$
Original (minority)	15.29	20.30	17.00	30.66	34.26	5.53
Oversampling (minority)	41.13	57.22	30.50	37.86	53.46	8.13
Improvement (minority)	25.84	36.92	13.50	7.20	19.20	2.60
Original (overall)	40.10	57.61	69.09	50.88	64.89	66.13
Oversampling (overall)	58.25	76.17	73.37	55.91	74.56	67.10
Improvement (overall)	18.15	18.56	4.28	5.03	9.67	0.97

Table: Test accuracy (%) on FashionMNIST when $R = 1000$. "Original (minority)" means that the test accuracy is evaluated only on the minority classes and oversampling is not used. When oversampling is used, we report the **best** test accuracy among four oversampling rates: 1, 10, 100, and 1000.

The best test accuracy is **not** achieved at $w_r = 1000$, indicating that oversampling with a large w_r would impair the test performance

Remarks on oversampling

- Large value of w_r can **mitigate** Minority Collapse on the training set
- But might degrade test accuracy

Remarks on oversampling

- Large value of w_r can **mitigate** Minority Collapse on the training set
 - But might degrade test accuracy
-
- Remains open: how to select an oversampling rate?
 - Other approaches such as *fixing* the classifiers?

Concluding remarks

One-line summary



It's a small but useful surrogate model



Future directions

Immediate connections:

- Go diverse: general imbalanced datasets
- Try various loss functions
- Relate Minority Collapse to fairness

Future directions

Immediate connections:

- Go diverse: general imbalanced datasets
- Try various loss functions
- Relate Minority Collapse to fairness

More broadly:

- Multiple Layer-Peeled Model:



$$\begin{aligned} \min_{\mathbf{W}, \mathbf{H}} \quad & \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{f}(\mathbf{h}_{k,i}, \mathbf{W}_{(L-m+1):L}), \mathbf{y}_k) \\ \text{s.t.} \quad & \frac{1}{K} \|\mathbf{W}_{(L-m+1):L}\|^2 \leq E_W \\ & \frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\mathbf{h}_{k,i}\|^2 \leq E_H \end{aligned}$$

- Model the training dynamics and test performance
- Why does the ansatz yield reasonable prediction?

Take-home messages

Layer-Peeled Model = minimal integration of

$$\text{prediction}(\mathbf{W}) + \text{representation}(\mathbf{H})$$

- Nonconvex but analytical
- Explain neural collapse
- Predict Minority Collapse
- Practical insights into deep learning

Reference

Exploring Deep Neural Networks via Layer-Peeled Model: Minority Collapse in Imbalanced Training
with Cong Fang, Hangfeng He, Qi Long. *Proceedings of the National Academy of Sciences (PNAS)*, 2021

- Code: <https://github.com/HornHehhf/LPM>
- NSF CAREER and TRIPODS, and Sloan

Take-home messages

Layer-Peeled Model = minimal integration of

$$prediction (\mathbf{W}) + representation (\mathbf{H})$$

- Nonconvex but analytical
- Explain neural collapse
- Predict Minority Collapse
- Practical insights into deep learning



Reference

Exploring Deep Neural Networks via Layer-Peeled Model: Minority Collapse in Imbalanced Training

with Cong Fang, Hangfeng He, Qi Long. *Proceedings of the National Academy of Sciences (PNAS)*, 2021

- Code: <https://github.com/HornHehhf/LPM>
- NSF CAREER and TRIPODS, and Sloan