

# Revealing the Structure of Deep Neural Networks via Convex Duality

ICML 2021

---

Tolga Ergen & Mert Pilanci

July 19, 2021

Stanford University



# Deep Learning Revolution

## Deep learning models:

- | often provide the best performance due to their large capacity
  - challenging to train

## Deep learning models:

- | often provide the best performance due to their large capacity
  - challenging to train
- | are complex black-box systems based on non-convex optimization
  - hard to interpret what the model is actually learning

# Prior Work on Regularized Deep Learning Training Problems

	Width ( $m$ )	Assumption	Depth ( $L$ )	# of outputs ( $K$ )
(Savarese et al., 2019)	1	1D data ( $d = 1$ )	2	7 ( $K = 1$ )
(Parhi and Nowak, 2019)	1	1D data ( $d = 1$ )	2	7 ( $K = 1$ )
(Ergen and Pilanci, 2020a,b)	nite	rank-one/whitened	2	3 ( $K = 1$ )
<b>Our results</b>	nite	rank-one/whitened or BatchNorm	$L = 2$	3 ( $K = 1$ )

# Prior Work on Regularized Deep Learning Training Problems

	Width ( $m$ )	Assumption	Depth ( $L$ )	# of outputs ( $K$ )
(Savarese et al., 2019)	1	1D data ( $d = 1$ )	2	7 ( $K = 1$ )
(Parhi and Nowak, 2019)	1	1D data ( $d = 1$ )	2	7 ( $K = 1$ )
(Ergen and Pilanci, 2020a,b)	nite	rank-one/whitened	2	3 ( $K = 1$ )
<b>Our results</b>	nite	rank-one/whitened or BatchNorm	$L = 2$	3 ( $K = 1$ )

**Optimal solution for  
 $L$ -layer ReLU networks is  
given by piecewise linear  
splines for any  $L \geq 2$ .**

**Figure 1:** One dimensional interpolation using  
 $L$ -layer ReLU networks

## Warmup: Two-layer Linear Networks

$\mathbf{X} \in \mathbb{R}^{n \times d}$  : Data matrix;  $\mathbf{y} \in \mathbb{R}^n$  : Label vector

$\mathbf{W}_l \in \mathbb{R}^{m_{l-1} \times m_l}$  :  $l^{\text{th}}$  layer weight matrix

$L(\cdot)$  : Arbitrary convex loss function

$\lambda > 0$  : Regularization coefficient

$f_{\lambda, L}(\mathbf{X})$  : Output of an L-layer network

| **Model:**  $f_{\lambda, 2}(\mathbf{X}) = \mathbf{XW}_1\mathbf{w}_2$

## Warmup: Two-layer Linear Networks

$\mathbf{X} \in \mathbb{R}^{n \times d}$ : Data matrix;  $\mathbf{y} \in \mathbb{R}^n$ : Label vector

$\mathbf{W}_l \in \mathbb{R}^{m_{l-1} \times m_l}$ :  $l^{\text{th}}$  layer weight matrix

$L(\cdot)$ : Arbitrary convex loss function

$\lambda > 0$ : Regularization coefficient

$f_{\cdot;L}(\mathbf{X})$ : Output of an L-layer network

| **Model:**  $f_{\cdot;2}(\mathbf{X}) = \mathbf{X}\mathbf{W}_1\mathbf{w}_2$

| **Optimization problem:**

$$\min_{\mathbf{W}_1, \mathbf{w}_2} L(f_{\cdot;2}(\mathbf{X}); \mathbf{y}) + (\lambda \|\mathbf{W}_1\|_F^2 + \lambda \|\mathbf{w}_2\|_2^2)$$



## Warmup: Two-layer Linear Networks

$\mathbf{X} \in \mathbb{R}^{n \times d}$ : Data matrix;  $\mathbf{y} \in \mathbb{R}^n$ : Label vector

$\mathbf{W}_l \in \mathbb{R}^{m_l \times m_{l-1}}$ :  $l^{\text{th}}$  layer weight matrix

$L(\cdot)$ : Arbitrary convex loss function

$\lambda > 0$ : Regularization coefficient

$f_{\cdot;L}(\mathbf{X})$ : Output of an L-layer network

| **Model:**  $f_{\cdot;2}(\mathbf{X}) = \mathbf{X}\mathbf{W}_1\mathbf{w}_2$

| **Optimization problem:**

$$\min_{\mathbf{W}_1, \mathbf{w}_2} L(f_{\cdot;2}(\mathbf{X}); \mathbf{y}) + \lambda (\|\mathbf{W}_1\|_F^2 + \|\mathbf{w}_2\|_2^2)$$

| **Optimal hidden layer weight:**  $\mathbf{w}_1 = \frac{\mathbf{X}^T P_{\mathbf{X}; (\mathbf{y})}}{\|\mathbf{X}^T P_{\mathbf{X}; (\mathbf{y})}\|_2}$   
where  $P_{\mathbf{X}; (\mathbf{y})}$  projects to  $\mathbf{u} \in \mathbb{R}^n$   $\|\mathbf{X}^T \mathbf{u}\|_2 = \|\mathbf{y}\|_2$ .

# Deep Linear Networks

| Model:  $f_{;L}(\mathbf{X}) = \prod_{j=1}^m \mathbf{X} \mathbf{W}_{1;j} \mathbf{W}_{2;j} \dots \mathbf{w}_{L;j}$

# Deep Linear Networks

- | **Model:**  $f_{;L}(\mathbf{X}) = \prod_{j=1}^m \mathbf{XW}_{1;j} \mathbf{W}_{2;j} \dots \mathbf{w}_{L;j}$
- | **Optimization problem:**

$$\min_{\mathbf{W}_1, \mathbf{W}_2} L(f_{;L}(\mathbf{X}); \mathbf{y}) + \sum_{j=1}^m \sum_{l=1}^L \lambda \|\mathbf{W}_{l;j}\|_F^2$$

# Deep Linear Networks

| **Model:**  $f_{;L}(\mathbf{X}) = \prod_{j=1}^m \mathbf{X} \mathbf{W}_{1;j} \mathbf{W}_{2;j} \dots \mathbf{w}_{L;j}$

| **Optimization problem:**

$$\min_{\mathbf{W}_1, \mathbf{W}_2} L(f_{;L}(\mathbf{X}); \mathbf{y}) + \sum_{j=1}^m \sum_{l=1}^L k \mathbf{W}_{l;j} k_F^2$$

| **Optimal hidden layer weights:**

$$\mathbf{W}_{l;j} = \begin{cases} t_j \frac{\mathbf{x}^T \mathbf{P}_{\mathbf{x};}(\mathbf{y})}{k \mathbf{X}^T \mathbf{P}_{\mathbf{x};}(\mathbf{y}) k_2} \mathbf{1}_j^T & \text{if } l = 1 \\ t_j \mathbf{1}_{1;j} \mathbf{1}_j^T & \text{if } 1 < l < L - 2 \\ \mathbf{1}_{L-2;j} & \text{if } l = L - 1 \end{cases}$$

where  $k_{1;j} k_2 = 1$ ,  $\mathbf{P}_{\mathbf{x};}(\cdot)$  projects to  $\mathbf{u} \in \mathbb{R}^n$   $j k \mathbf{X}^T \mathbf{u} k_2$   $t_j^2 \leq \frac{1}{L}$   
 and  $t_j = k \mathbf{W}_{l;j} k_F$ .

# Deep ReLU Networks

- | **Model:**  $f_{;L}(\mathbf{X}) = \mathbf{A}_{L-1} \mathbf{w}_L$ , where  $\mathbf{A}_{l;j} = (\mathbf{A}_{l-1;j} \mathbf{W}_{l;j})_+$ ,  $\mathbf{A}_{0;j} = \mathbf{X}$ ,  $g_{l;j}$ , and  $(x)_+ = \max\{0; x\}$

# Deep ReLU Networks

- Model:  $f_{;L}(\mathbf{X}) = \mathbf{A}_{L-1} \mathbf{w}_L$ , where  $\mathbf{A}_{l;j} = (\mathbf{A}_{l-1;j} \mathbf{W}_{l;j})_+$ ;  $\mathbf{A}_{0;j} = \mathbf{X}$ ,  $8l;j$ , and  $(x)_+ = \max\{0; x\}$

## Theorem

Let  $\mathbf{X}$  be a rank-one matrix such that  $\mathbf{X} = \mathbf{c} \mathbf{a}_0^T$ , where  $\mathbf{c} \in \mathbb{R}_+^n$  and  $\mathbf{a}_0 \in \mathbb{R}^d$ , then strong duality holds and the optimal weights are

$$\mathbf{W}_{l;j} = \frac{\mathbf{1}_{1;j}}{k_{l-1;j} k_2} \mathbf{1}_{l;j}^T; \quad 8l \in [L-2]; \quad \mathbf{w}_{L-1;j} = \frac{\mathbf{1}_{L-2;j}}{k_{L-2;j} k_2};$$

where  $\mathbf{1}_{0;j} = \mathbf{a}_0$  and  $\mathbf{1}_{l;j} \in \mathbb{R}_+^{l-1}$  is a set of vectors such that  $\mathbf{1}_{l;j} \in \mathbb{R}_+^{m_l}$  and  $k_{l-1;j} k_2 = \mathbf{1}_{l;j}^T \mathbf{1}_{l;j}$ ;  $8l \in [L-2]; 8j \in [m]$ .

# Deep ReLU Networks

- Model:**  $f_{;L}(\mathbf{X}) = \mathbf{A}_{L-1} \mathbf{w}_L$ , where  $\mathbf{A}_{l;j} = (\mathbf{A}_{l-1;j} \mathbf{W}_{l;j})_+$ ;  $\mathbf{A}_{0;j} = \mathbf{X}$ ,  $g_{l;j}$ , and  $(x)_+ = \max\{0; x\}$

## Theorem

Let  $\mathbf{X}$  be a rank-one matrix such that  $\mathbf{X} = \mathbf{c} \mathbf{a}_0^T$ , where  $\mathbf{c} \in \mathbb{R}_+^n$  and  $\mathbf{a}_0 \in \mathbb{R}^d$ , then strong duality holds and the optimal weights are

$$\mathbf{W}_{l;j} = \frac{l-1;j}{k_{l-1;j} k_2} \mathbf{t}_{l;j}; \quad g_{l;j} \in [L-2]; \quad \mathbf{w}_{L-1;j} = \frac{L-2;j}{k_{L-2;j} k_2};$$

where  $\mathbf{a}_{0;j} = \mathbf{a}_0$  and  $f_{l-1;j} \mathbf{g}_{l=1}^L$  is a set of vectors such that  $\mathbf{t}_{l;j} \in \mathbb{R}_+^m$  and  $k_{l-1;j} k_2 = \mathbf{t}_{l;j}; \quad g_{l;j} \in [L-2]; \quad g_{j} \in [m]$ .

## Corollary

For 1D data, i.e.,  $\mathbf{x} \in \mathbb{R}^n$ , the optimal network output has kinks only at the input data points, i.e., the output function is in the following form:  $f_{;L}(\hat{\mathbf{x}}) = \sum_i (\hat{\mathbf{x}} - x_i)_+$ . Therefore, the optimal network output is a linear spline interpolation.

# Vector-output ReLU Networks

## Theorem

Let  $(\mathbf{X}; \mathbf{Y})$  be a dataset such that  $\mathbf{X}\mathbf{X}^T = \mathbf{I}_n$  and  $\mathbf{Y} \in \mathbb{R}^{n \times K}$  is one-hot encoded, then a set of optimal solutions for the following regularized training problem

$$\min_{\Theta} \frac{1}{2} \|\mathbf{f}_{\Theta}(\mathbf{X}) - \mathbf{Y}\|_F^2 + \frac{\lambda}{2} \sum_{j=1}^m \sum_{l=1}^L \|\mathbf{w}_{l,j}\|_F^2$$

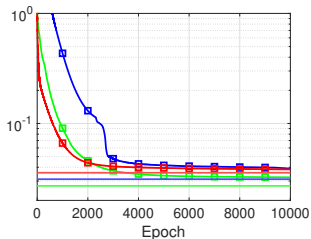
can be formulated as follows

$$\mathbf{w}_{l,j} = \begin{cases} \frac{\phi_{l-1,j}}{\|\phi_{l-1,j}\|_2} \mathbf{t}_{l,j} & \text{if } l \in [L-1]; \\ (\|\phi_{l-1,j}\|_2)_+ \mathbf{e}_j^T & \text{if } l = L \end{cases};$$

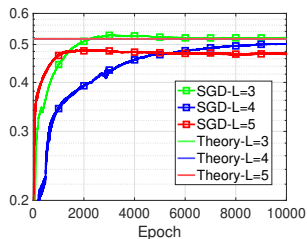
where  $\phi_{l-1,j} = \mathbf{X}^T \mathbf{y}_j$ ,  $\mathbf{t}_{l,j} \in \mathbb{R}_+^{m_l}$  are vectors such that  $\|\mathbf{t}_{l,j}\|_2 = 1$ ,  $\|\phi_{l-1,j}\|_2 = t_j$ , and  $\mathbf{e}_j^T \mathbf{t}_{l,j} = 0$ ;  $\delta_i \in j$ . Moreover,  $\mathbf{e}_j$  is the  $j^{\text{th}}$  ordinary basis vector.



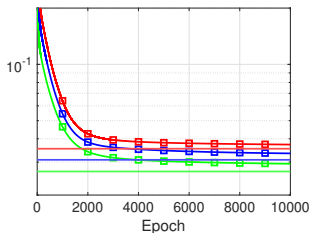
# Numerical Results



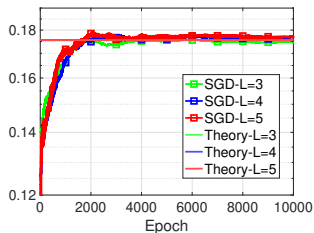
(a) MNIST-Training obj.



(b) MNIST-Test accuracy



(c) CIFAR10-Training obj.



(d) CIFAR10-Test accuracy

Figure 2: Training and test performance on whitened and sampled datasets.

# Takeaways and Open Problems

- | Optimal solutions to regularized deep neural network training problems can be explicitly characterized via convex analytic frameworks

## Takeaways and Open Problems

- | Optimal solutions to regularized deep neural network training problems can be explicitly characterized via convex analytic frameworks
- | When the input data is whitened or rank-one, optimal layer weights of an  $L$ -layer deep ReLU network can be found the closed-form

# Takeaways and Open Problems

- | Optimal solutions to regularized deep neural network training problems can be explicitly characterized via convex analytic frameworks
- | When the input data is whitened or rank-one, optimal layer weights of an  $L$ -layer deep ReLU network can be found the closed-form
- | For 1D datasets, kinks of ReLU occur exactly at the input data so that the optimal network outputs linear spline interpolations

# Takeaways and Open Problems

- | Optimal solutions to regularized deep neural network training problems can be explicitly characterized via convex analytic frameworks
- | When the input data is whitened or rank-one, optimal layer weights of an  $L$ -layer deep ReLU network can be found the closed-form
- | For 1D datasets, kinks of ReLU occur exactly at the input data so that the optimal network outputs linear spline interpolations
- | **Open problems:**
  - extension of the analysis to standard deep networks

# Takeaways and Open Problems

- | Optimal solutions to regularized deep neural network training problems can be explicitly characterized via convex analytic frameworks
- | When the input data is whitened or rank-one, optimal layer weights of an  $L$ -layer deep ReLU network can be found the closed-form
- | For 1D datasets, kinks of ReLU occur exactly at the input data so that the optimal network outputs linear spline interpolations
- | **Open problems:**
  - extension of the analysis to standard deep networks
  - generalization properties of the optimal solutions

## References

---

- Ergen, T. and Pilanci, M. (2020a). Convex geometry and duality of over-parameterized neural networks. *arXiv preprint arXiv:2002.11219*.
- Ergen, T. and Pilanci, M. (2020b). Convex geometry of two-layer relu networks: Implicit autoencoding and interpretable models. In Chiappa, S. and Calandra, R., editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 4024{4033, Online. PMLR.
- Parhi, R. and Nowak, R. D. (2019). Minimum  $\ell_1$  norm neural networks are splines. *arXiv preprint arXiv:1910.02333*.

Savarese, P., Evron, I., Soudry, D., and Srebro, N. (2019). How do finite width bounded norm networks look in function space? *CoRR*, abs/1902.05040.