

On the Sparsity-Promoting Effect of Weight Decay in Deep Learning

Rahul Parhi

Biomedical Imaging Group

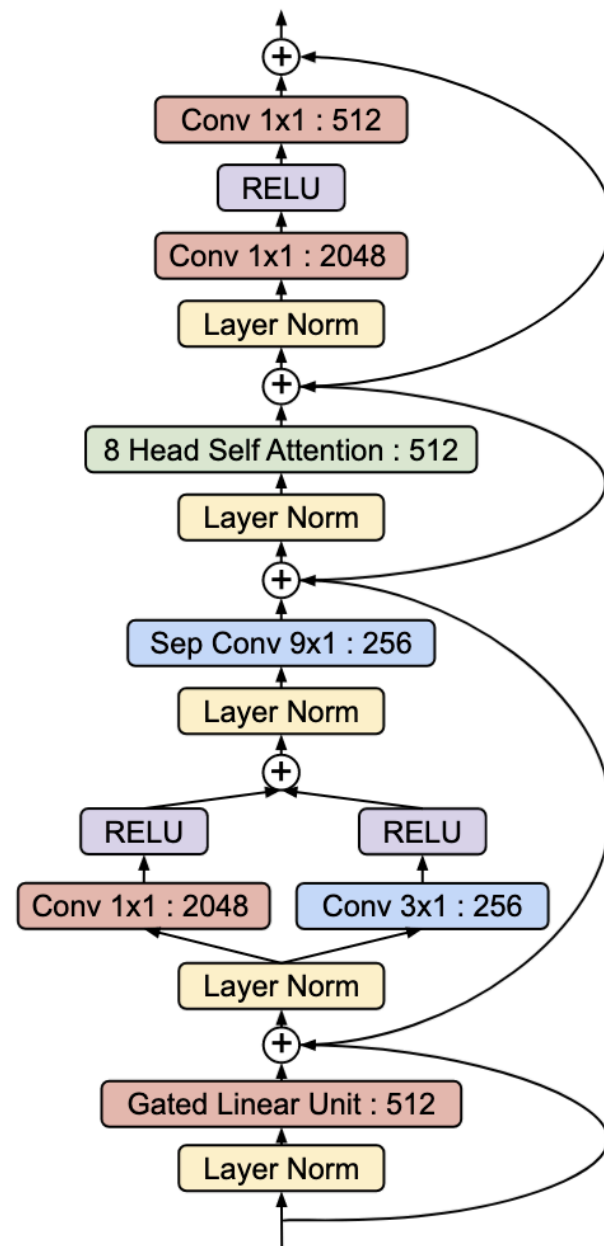
École polytechnique fédérale de Lausanne

Conference on Parsimony and Learning
4 January 2024

Deep Neural Network Architectures

The Evolved Transformer

David So, Quoc Le, Chen Liang *Proceedings of the 36th International Conference on Machine Learning*, PMLR 97:5877-5886, 2019.



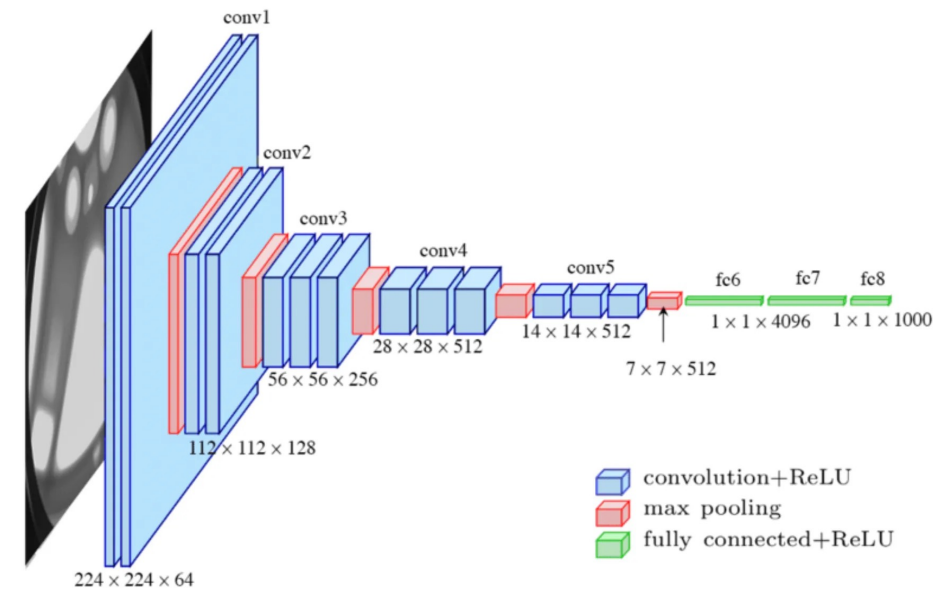
Google DeepMind

2023-10-26

ConvNets Match Vision Transformers at Scale

Samuel L Smith¹, Andrew Brock¹, Leonard Berrada¹ and Soham De¹

¹Google DeepMind

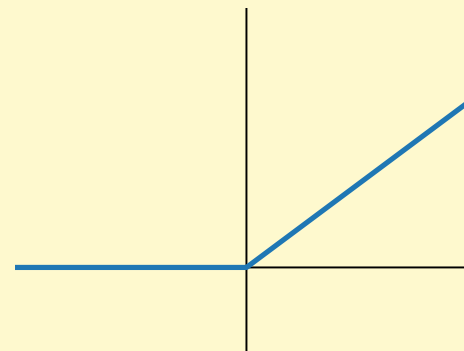


Very deep convolutional networks for large-scale image recognition

[K Simonyan, A Zisserman](#)

In this work we investigate the effect of the convolutional network depth on its accuracy in the large-scale image recognition setting. Our main contribution is a thorough evaluation of ...

☆ Save 📄 Cite Cited by 112161 Related articles 🔗

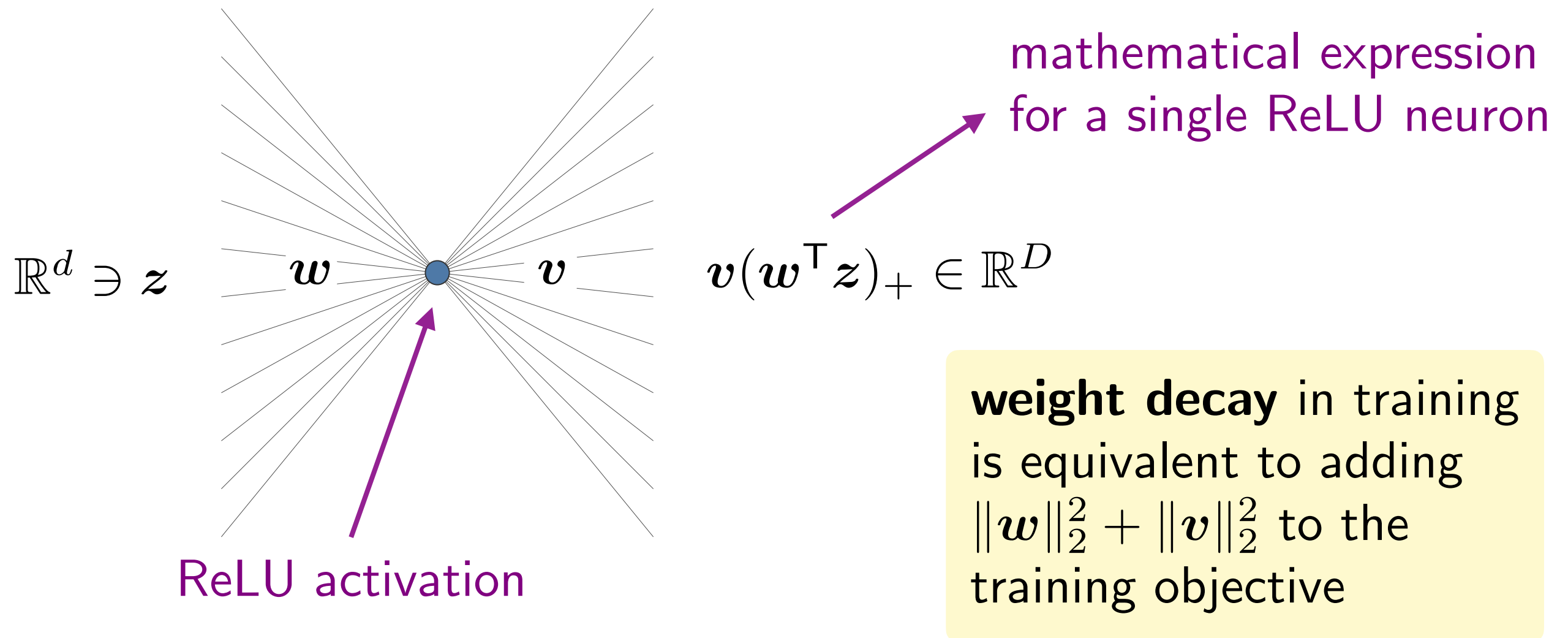


Rectified Linear Unit (ReLU)

$$\text{ReLU}(t) = \max\{0, t\} = t_+$$

+ weight decay in training

Neural Balance in Deep Neural Networks

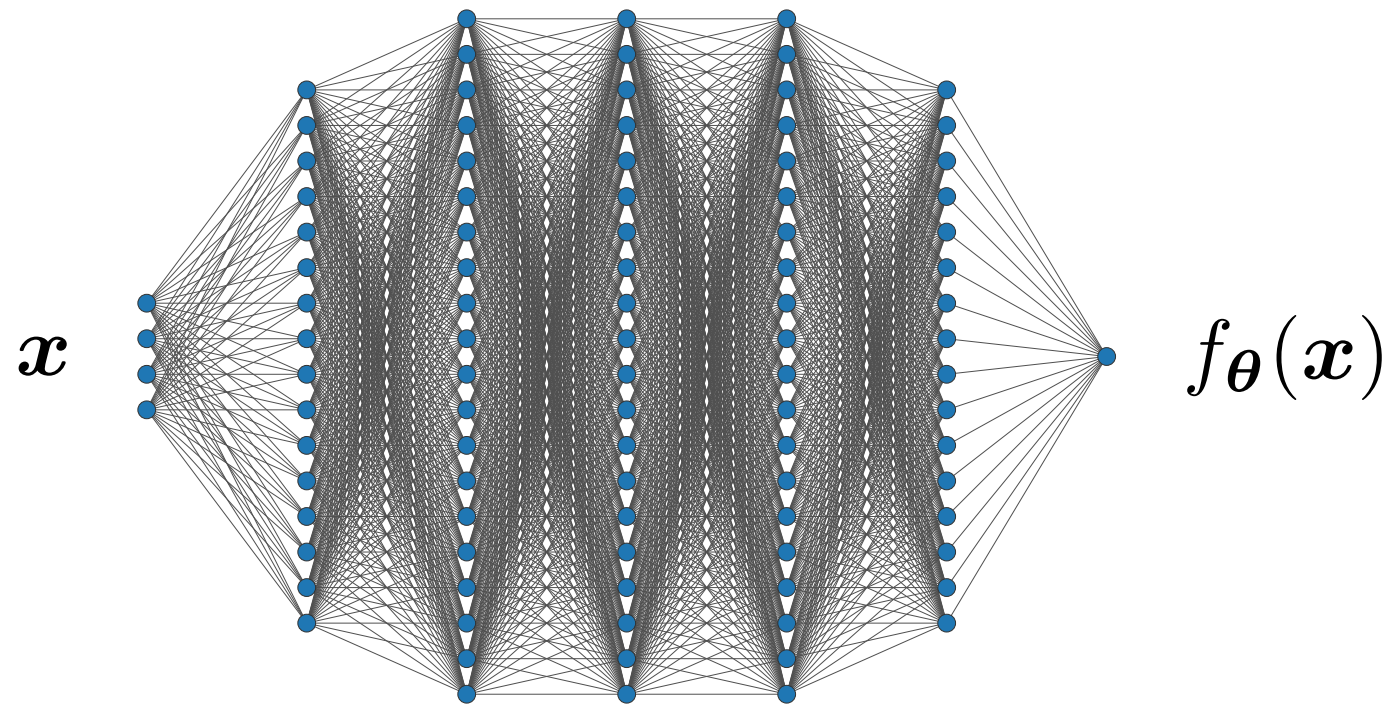


Neural Balance Theorem

If a DNN is trained with weight decay, then the 2-norms of the input and output weights to each ReLU neuron must be **balanced**.

$$\|\mathbf{w}\|_2 = \|\mathbf{v}\|_2$$

Neural Network Training



parameterized by a vector $\theta \in \mathbb{R}^P$
of neural network **weights**

Neural network training problem for the data $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$.

$$\min_{\theta \in \mathbb{R}^P} \underbrace{\sum_{n=1}^N \mathcal{L}(y_n, f_{\theta}(\mathbf{x}_n))}_{\text{data fidelity}} + \underbrace{\frac{\lambda}{2} \|\theta\|_2^2}_{\text{regularization}}$$

Neural Balance

The ReLU activation is **homogeneous**

$$\boldsymbol{v}(\boldsymbol{w}^\top \boldsymbol{z})_+ = \gamma^{-1} \boldsymbol{v}(\gamma \boldsymbol{w}^\top \boldsymbol{z})_+, \quad \text{for any } \gamma > 0.$$

At a global minimizer of the weight decay objective, $\|\boldsymbol{v}\|_2 = \|\boldsymbol{w}\|_2$.

Proof. The solution to

$$\min_{\gamma > 0} \|\gamma^{-1} \boldsymbol{v}\|_2 + \|\gamma \boldsymbol{w}\|_2$$

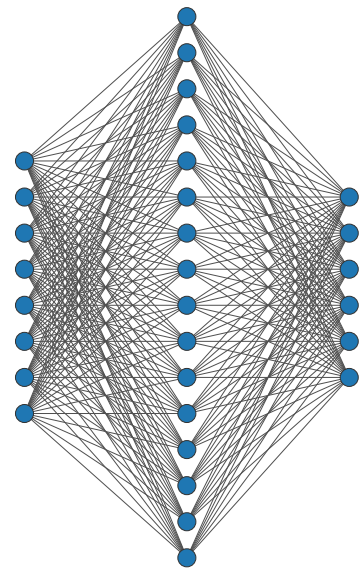
is $\gamma = \sqrt{\|\boldsymbol{v}\|_2 / \|\boldsymbol{w}\|_2}$. □

$$\text{At a global minimizer, } \frac{\|\boldsymbol{v}\|_2^2 + \|\boldsymbol{w}\|_2^2}{2} = \|\boldsymbol{v}\|_2 \|\boldsymbol{w}\|_2.$$

Grandvalet (1998, ICANN)

Neyshabur et al. (2015, ICLR Workshop)

Secret Sparsity/Implicit Parsimony of Weight Decay



$$f_{\theta}(x) = \sum_{k=1}^K v_k (w_k^T x)_+$$

$$\theta = \{(w_k, v_k)\}_{k=1}^K$$

$$\min_{\theta = \{(w_k, v_k)\}_{k=1}^K} \sum_{n=1}^N \mathcal{L}(y_n, f_{\theta}(x_n)) + \frac{\lambda}{2} \sum_{k=1}^K \|v_k\|_2^2 + \|w_k\|_2^2$$

weight decay

$$\min_{\theta = \{(w_k, v_k)\}_{k=1}^K} \sum_{n=1}^N \mathcal{L}(y_n, f_{\theta}(x_n)) + \lambda \sum_{k=1}^K \|v_k\|_2 \|w_k\|_2$$

path-norm

$$\min_{\substack{\theta = \{(w_k, v_k)\}_{k=1}^K \\ \|w_k\|_2 = 1}} \sum_{n=1}^N \mathcal{L}(y_n, f_{\theta}(x_n)) + \lambda \sum_{k=1}^K \|v_k\|_2$$

multitask lasso

Rebalancing

Secret Sparsity/Implicit Parsimony of Weight Decay

$$\text{weight decay} \iff \min_{\substack{\boldsymbol{\theta} = \{(\mathbf{w}_k, \mathbf{v}_k)\}_{k=1}^K \\ \|\mathbf{w}_k\|_2 = 1}} \sum_{n=1}^N \mathcal{L}(\mathbf{y}_n, f_{\boldsymbol{\theta}}(\mathbf{x}_n)) + \lambda \sum_{k=1}^K \|\mathbf{v}_k\|_2$$

- Weight decay is equivalent to a **non-convex** multitask lasso.

What Kinds of Functions Do Neural Networks Learn?

Why Do Neural Networks Work Well in High-Dimensional Problems?

Principled Compression Algorithms for Pre-Trained DNNs.

A Banach Space Representer Theorem

Neural Network Representer Theorem (P. and Nowak 2021)

For any data set $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$ and lower semicontinuous $\mathcal{L}(\cdot, \cdot)$, there exists a solution to

$$\min_{f \in \mathcal{R} \text{BV}^2} \sum_{n=1}^N \mathcal{L}(y_n, f(\mathbf{x}_n)) + \lambda \|f\|_{\mathcal{R} \text{BV}^2}, \quad \lambda > 0,$$

that admits a representation of the form

$$f_{\text{ReLU}}(\mathbf{x}) = \sum_{k=1}^K v_k \underbrace{(\mathbf{w}_k^T \mathbf{x} - b_k)_+}_{\text{ReLU neurons}} + \underbrace{\mathbf{w}_0^T \mathbf{x} + b_0}_{\text{skip connection}}, \quad \underbrace{K < N}_{\text{sparse solution}}.$$

Training a **sufficiently parameterized** neural network ($K \geq N$) with weight decay (to a global minimizer) is a solution to the Banach space problem.

Neural networks learn $\mathcal{R} \text{BV}^2$ -functions.

Minimax Optimality of Neural Networks

Suppose that $\{\mathbf{x}_n\}_{n=1}^N$ are i.i.d. uniform on a bounded domain $\Omega \subset \mathbb{R}^d$. If $y_n = f^*(\mathbf{x}_n) + \varepsilon_n$ with $\|f^*\|_{\mathcal{R}BV^2} < \infty$, then any solution to

$$f_{\text{ReLU}} \in \arg \min_{\boldsymbol{\theta}} \sum_{n=1}^N \mathcal{L}(y_n, f_{\boldsymbol{\theta}}(\mathbf{x}_n)) + \frac{\lambda}{2} \sum_{k=1}^K |v_k|^2 + \|\mathbf{w}_k\|_2^2$$

weight decay
objective

satisfies

$$\mathbf{E} \|f^* - f_{\text{ReLU}}\|_{L^2(\Omega)}^2 = \tilde{O}(N^{-\frac{d+3}{2d+3}}) = \tilde{O}(N^{-\frac{1}{2}}).$$

minimax rate

no curse

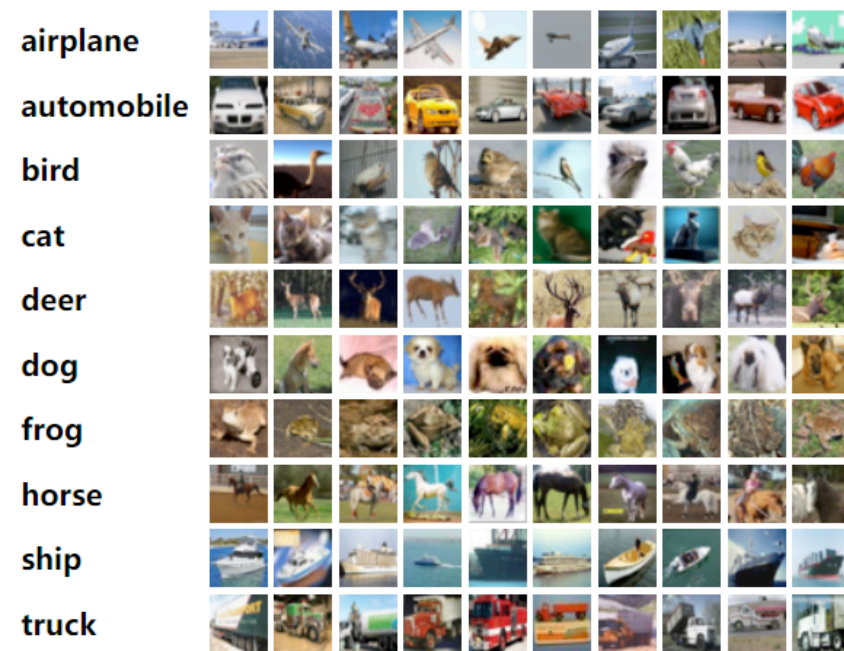
Linear methods (thin-plate splines, kernel methods, neural tangent kernels, etc.) **necessarily** suffer the curse of dimensionality.

Linear minimax lower bound: $N^{-\frac{3}{d+3}}$

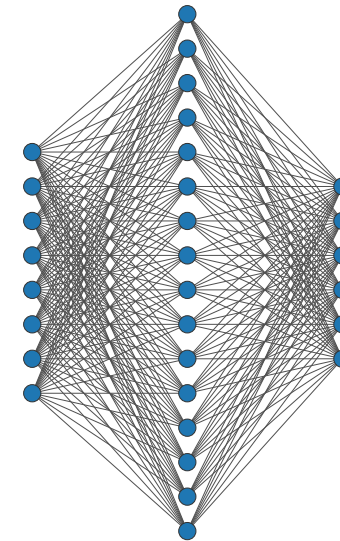
the curse

Principled DNN Compression

VGG-19 trained with weight decay on CIFAR-10.



final ReLU layer
 $K = 512$ neurons



output dimension
 $D = 10$

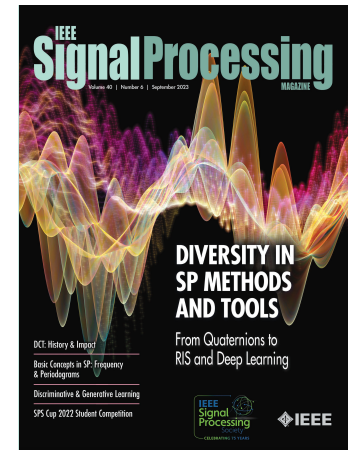
Theory: There exists a representation with ≤ 100 neurons that can be found by solving a **convex** multitask lasso problem.

	original network	compressed network
active neurons	512	47
test accuracy	93.92%	93.88%
train loss	0.0104	0.0112

10× compression!
no change in
performance!

Conclusion

- General audience overview article:
 - *Deep Learning Meets Sparse Regularization: A signal processing perspective*
Rahul Parhi and Robert D. Nowak
IEEE Signal Processing Magazine, vol. 40, no. 6, pp. 63–74, Sept. 2023.
- Technical articles:
 - *Banach Space Representer Theorems for Neural Networks and Ridge Splines*
Rahul Parhi and Robert D. Nowak
Journal of Machine Learning Research, vol. 22, no. 43, pp. 1–40, 2021.
 - *What Kinds of Functions Do Deep Neural Networks Learn? Insights from Variational Spline Theory*
Rahul Parhi and Robert D. Nowak
SIAM Journal on Mathematics of Data Science, vol. 4, no. 2, pp. 464–489, 2022.
 - *Near-Minimax Optimal Estimation With Shallow ReLU Neural Networks*
Rahul Parhi and Robert D. Nowak
IEEE Transactions on Information Theory, vol. 69, no. 2, pp. 1125–1140, Feb. 2023.
 - *Vector-Valued Variation Spaces and Width Bounds for DNNs: Insights on Weight Decay Regularization*
Joseph Shenouda, **Rahul Parhi**, Kangwook Lee, Robert D. Nowak
arXiv preprint arXiv:2305.16534, 2023+



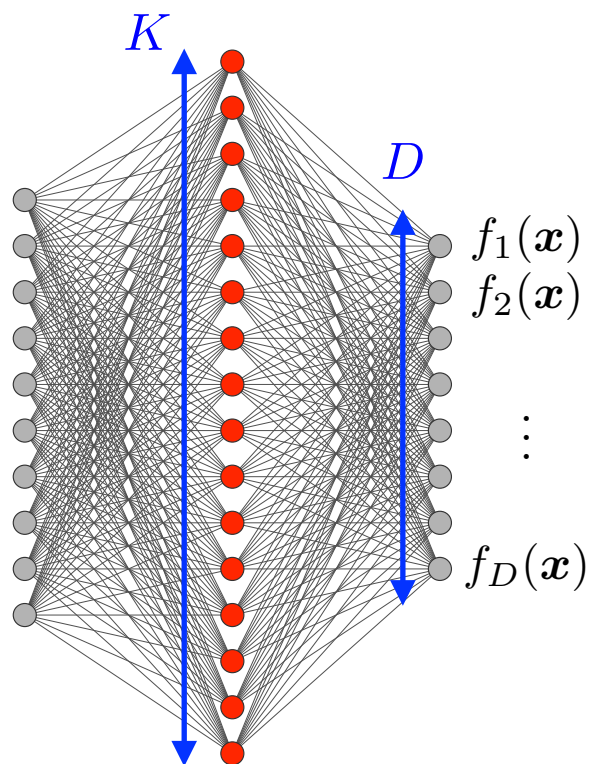
Questions?

The Structured Sparsity of Weight Decay

$$\min_{\theta = \{(\mathbf{w}_k, \mathbf{v}_k)\}_{k=1}^K, \|\mathbf{w}_k\|_2=1} \sum_{n=1}^N \mathcal{L}(\mathbf{y}_n, f_{\theta}(\mathbf{x}_n)) + \lambda \sum_{k=1}^K \|\mathbf{v}_k\|_2$$

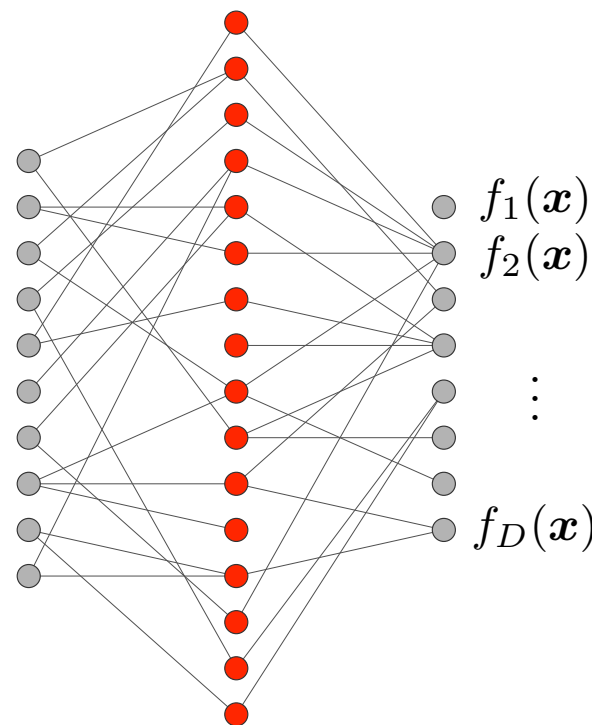
weight decay
 \longleftrightarrow
 non-convex multitask lasso

dense weights



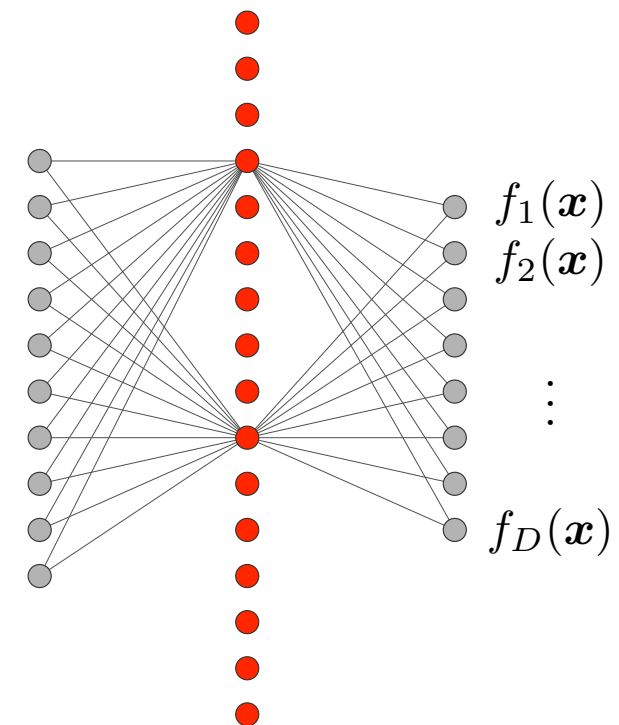
$$\|f\|_{\mathcal{V}} = O(K\sqrt{D})$$

sparse weights



$$\|f\|_{\mathcal{V}} = O(D)$$

sparse neurons

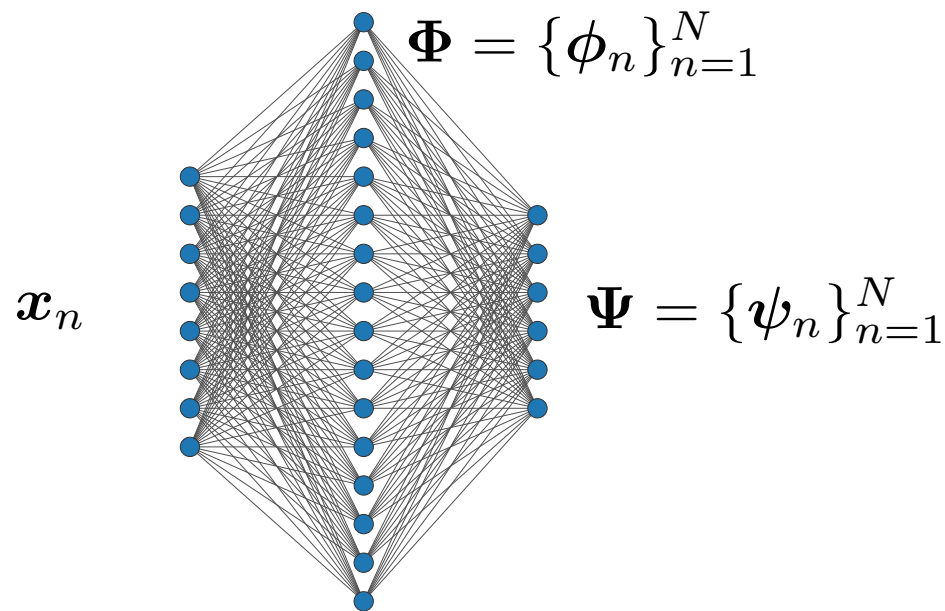


$$\|f\|_{\mathcal{V}} = O(\sqrt{D})$$

Weight decay favors variation in only a few directions (sparse weights)

Weight decay favors outputs that “share” neurons (sparse neurons)

Tight Bounds on Widths



$$\min_{\{\mathbf{v}_k\}_{k=1}^K} \sum_{k=1}^K \|\mathbf{v}_k\|_2 \quad \text{s.t.} \quad \Psi = \mathbf{V}\Phi.$$

Low-rank data embeddings have been observed empirically by [Huh et al. \(2022\)](#).

Layer Width Theorem (Shenouda, P., Lee and Nowak 2023+)

Let Φ denote the post-activation features and Ψ denote the neuron outputs of any ReLU layer in a **trained** DNN (minimizes the weight decay objective). Then, there exists a representation with

$$K \leq \text{rank}(\Phi) \text{rank}(\Psi) \leq N^2$$

Bound of [Jacot \(2023\)](#): $N(N + 1)$.

neurons. The representation can be found by solving a **convex multitask lasso** problem.