

Towards Understanding Hierarchical Learning: Benefits of Neural Representations

Minshuo Chen*, Yu Bai^o, Jason D. Lee[†], Tuo Zhao^{*}, Huan Wang^o, Caiming Xiong^o, Richard Socher^o *Georgia Tech, *Salesforce Research, *Princeton University



Three-layer Network with Representation

▷ Prototypical trainable Model:

$$f_{\mathbf{W}}(\mathbf{x}) = \frac{1}{\sqrt{m}} \sum_{r=1}^{m} a_r \phi(\mathbf{w}_r^{\top} \mathbf{h}(\mathbf{x})).$$

- $\mathbf{x} \in \mathbb{R}^d$ is the raw input data;
- $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_m]^\top \in \mathbb{R}^{m \times D}$ is the weight matrix;
- ϕ is the activation function;
- $\mathbf{h} : \mathbb{R}^d \mapsto \mathbb{R}^D$ is a **fixed** representation function.

▷ One-layer Neural Representation Function:

$$\mathbf{h}(\mathbf{x}) = \sigma(\mathbf{V}\mathbf{x} + \mathbf{b}).$$

- $\mathbf{V} \in \mathbb{R}^{D \times d}$ and $\mathbf{b} \in \mathbb{R}^{D}$ are *fixed* weight parameters;
- σ is some nonlinear function.

Similar to a three-layer neural network:

$$\widetilde{f}_{\mathbf{W},\mathbf{V},\mathbf{b}}(\mathbf{x}) = \frac{1}{\sqrt{m}} \mathbf{a}^{\top} \phi(\mathbf{W}\sigma(\mathbf{V}\mathbf{x}+\mathbf{b})),$$

We let W be trainable and fix (V, b).

Linearized / Taylorized Models

▷ We consider the linearized / taylorized versions:

$$f_{\mathbf{W}}^{L}(\mathbf{x}) = \frac{1}{\sqrt{m}} \sum_{r=1}^{m} a_r \phi'(\mathbf{w}_{0,r}^{\top} \mathbf{h}(\mathbf{x}))(\mathbf{w}_r^{\top} \mathbf{h}(\mathbf{x})), \qquad (\text{NTK-h})$$

$$f_{\mathbf{W}}^{Q}(\mathbf{x}) = \frac{1}{2\sqrt{m}} \sum_{r=1}^{m} a_r \phi''(\mathbf{w}_{0,r}^{\top} \mathbf{h}(\mathbf{x})) (\mathbf{w}_r^{\top} \mathbf{h}(\mathbf{x}))^2. \quad (\text{Quad-h})$$

• \mathbf{W}_0 is the initialization;

• W is overloaded to denote the **increment** of weights.

Note: (NTK-h) and (Quad-h) are the linearized / quadratic expansions of the 3-layer neural network w.r.t. the \mathbf{W} layer. ▷ Supervised Learning — Minimizing Regularized Risk:

$$\min_{\mathbf{W}} \widehat{\mathcal{R}}_{\lambda} \left(f_{\mathbf{W}}^{Q} \right) \coloneqq \frac{1}{n} \sum_{i=1}^{n} \ell \left(f_{\mathbf{W}}^{Q}(\mathbf{x}_{i}), y_{i} \right) + \lambda \|\mathbf{W}\|_{2,4}^{4}.$$

- Data: $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ are data points with $\mathbf{x}_i \in \mathbb{S}^{d-1}$ and $y \in \mathcal{Y};$
- Loss: $\ell(\cdot, y)$ is a convex loss function, twice differentiable with bounded derivatives, and $|\ell(\mathbf{0}, y)| \leq 1$ for any $y \in \mathcal{Y}$, e.g., logistic loss and soft hinge loss;
- Regularization: $\|\mathbf{W}\|_{2,4}^4 \coloneqq \sum_{r=1}^m \|\mathbf{w}_r\|_2^4$.
- Initialization: $\mathbf{w}_{0,r} \stackrel{\text{i.i.d}}{\sim} \mathbb{N}(\mathbf{0}, \mathbf{I}_D)$ and $a_r \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\{\pm 1\})$.

Optimization and Generalization

radius $B_{w,\star} > 0$, suppose the width $m \ge \widetilde{O}(B_h^4 B_{w,\star}^4 \epsilon^{-1})$ and coefficient $\lambda > 0$ chosen properly. Then any second-order stationary point (SOSP) $\widehat{\mathbf{W}}$ of the regularized risk $\widehat{\mathcal{R}}_{\lambda}(f^Q_{\mathbf{W}})$ satisfies $\|\widehat{\mathbf{W}}\|_{2,4} \leq O(B_{w,\star})$, and achieves

 $\mathbb{E}_{(\mathbf{x}_i, j)}$







Theorem 1. Assume

Bounded representation: $\|\mathbf{h}(\mathbf{x})\|_2 \leq B_h$ almost surely;

Differentiable activation: $\sup_{t \in \mathbb{R}} |\phi''(t)| \leq const.$

(1) (Optimization) Given any $\epsilon > 0$, $\tau = \Theta(1)$, and some

$$\widehat{\mathcal{R}}_{\lambda}(f_{\widehat{\mathbf{W}}}^{Q}) \le (1+\tau) \min_{\|\mathbf{W}\|_{2,4} \le B_{w,\star}} \widehat{\mathcal{R}}(f_{\mathbf{W}}^{Q}) + \epsilon.$$

(2) (Generalization) For any radius $B_w > 0$, we have with high probability (over $(\mathbf{a}, \mathbf{W}_0)$) that

$$y_{i}\left|\sup_{\|\mathbf{W}\|_{2,4}\leq B_{w}}\left|\mathcal{R}(f_{\mathbf{W}}^{Q})-\widehat{\mathcal{R}}(f_{\mathbf{W}}^{Q})\right|\right| \leq \widetilde{O}\left(\frac{B_{h}^{2}B_{w}^{2}M_{h,op}}{\sqrt{n}}+\frac{1}{\sqrt{n}}\right),$$

where $M_{h,op}^2 = B_h^{-2} \mathbb{E}_{\mathbf{x}} \left[\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{h}(\mathbf{x}_i) \mathbf{h}(\mathbf{x}_i)^\top \right\|_2 \right]$ and \mathcal{R} , $\widehat{\mathcal{R}}$ are unregularized population and empirical risks, respectively.

Implications of Theorem 1:

- Efficient optimization: escaping-saddle type algorithms, e.g., noisy SGD, can efficiently optimize $\mathcal{R}(f^Q_{\mathbf{W}})$;
- Feature isotropicity: isotropic $h(x) \Rightarrow$ small $M_{h,op} \Rightarrow$ good generalization.

Whitened Representations

▷ Whitened Representation:

$$\mathbf{x} = \widehat{\mathbf{\Sigma}}^{-1/2} \mathbf{g}(\mathbf{x}), \quad \text{with} \quad \mathbf{g}(\mathbf{x}) = \sigma(\mathbf{V}\mathbf{x} + \mathbf{b}).$$

• Sample convariance matrix $\widehat{\Sigma} = \frac{1}{n_0} \sum_{i=1}^{n_0} \mathbf{g}(\widetilde{\mathbf{x}}_i) \mathbf{g}(\widetilde{\mathbf{x}}_i)^{\top}$

• $\{\widetilde{\mathbf{x}}_i\}_{i=1}^{n_0}$ are additional unlabeled samples;

• $\mathbf{v}_r \stackrel{\text{i.i.d.}}{\sim} \mathbb{N}(\mathbf{0}, \mathbf{I}_d)$ and $\mathbf{b} \sim \mathbb{N}(\mathbf{0}, \mathbf{I}_D)$; fixed during training.

Assumption 1 (Lower bounded covariance). For any k and $D \leq O(d^k)$, with high probability over \mathbf{V}, \mathbf{b} (as $d \to \infty$), we have the minimum eigenvalue

 $\lambda_{\min}(\mathbf{\Sigma}) \geq \lambda_k$

for some constant $\lambda_k > 0$ that only depends on k but not d_k where $\Sigma = \mathbb{E}_{\mathbf{x}}[\sigma (\mathbf{V}\mathbf{x} + \mathbf{b}) \sigma (\mathbf{V}\mathbf{x} + \mathbf{b})^{\top}]$

•
$$D = 6$$

> Benefit of Representations — Hierarchical Learning: Quad-h model can express polynomials hierarchically, using weight matrices with much smaller norms than that of a shallow learner.



Discussions

• Unwhitened representation: using a data dependent regularizer can also achieve the improved sample complexity.

• Implications of Assumption 1: $h(x_i)$'s are not too correlated, which roughly requires that the distribution of \mathbf{x} spans all directions in \mathbb{R}^d .

Quad-h significantly outperforms Quad-Raw (when $p \geq 4$), while NTK-h does not improve on NTK-Raw.

 $\min_{\mathbf{W}} \ \widehat{\mathcal{R}}_{\lambda}^{\mathrm{dreg}}(f_{\mathbf{W}}^{Q}) = \frac{1}{n} \sum_{i=1}^{n} \ell(f_{\mathbf{W}}^{Q}(\mathbf{x}_{i}), y_{i}) + \lambda \left\| \mathbf{W} \widehat{\boldsymbol{\Sigma}}^{1/2} \right\|_{2, 4}^{4}.$