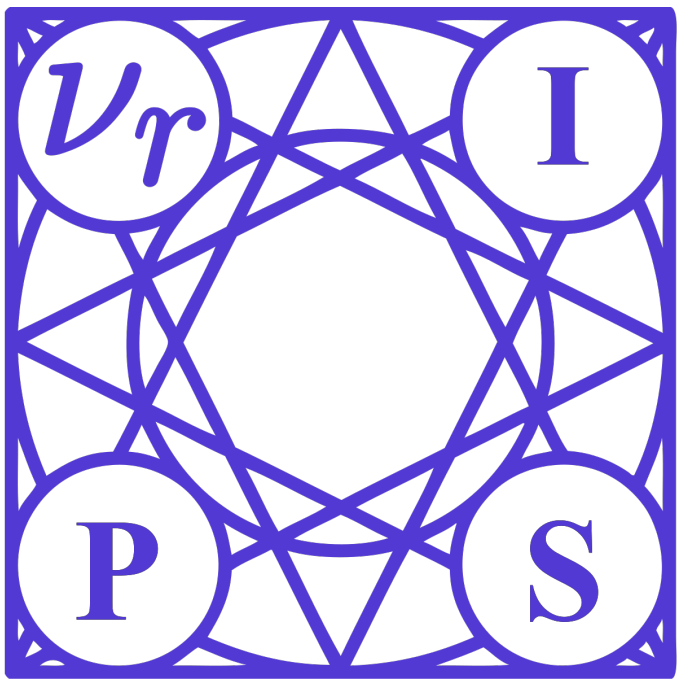




Algorithm-Dependent Generalization Bounds for Overparameterized Deep Residual Networks



Spencer Frei^o and Yuan Cao[†] and Quanquan Gu[†]

Department of Statistics^o and Department of Computer Science[†], UCLA

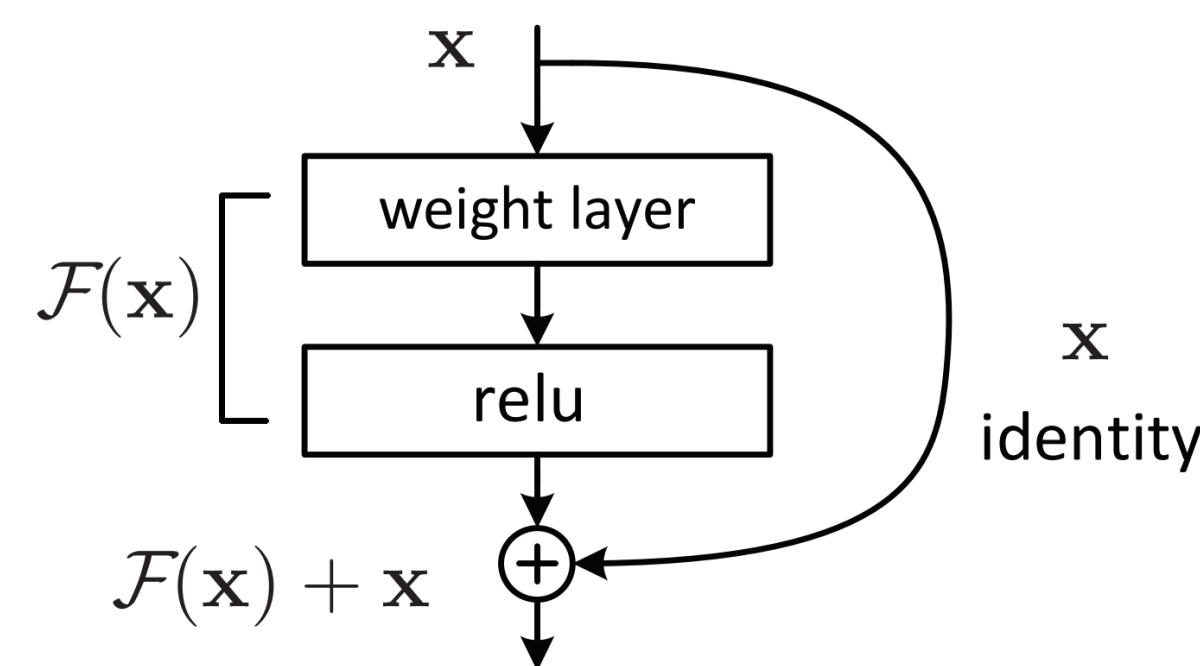
Background

- **Residual connections** common in modern NNs: but theoretical justifications lacking.
- **Fewer parameters, better generalization** observed empirically in many residual architectures.

Problem Description

- Input $(x, y) \in \mathbb{R}^d \times \{\pm 1\}$, binary classification under cross-entropy loss $\ell(z) := \log(1 + \exp(-z))$.
- $f_W(x)$ = output of $L + 1$ hidden layer residual network,

$$\begin{aligned} x_1 &= W_1^\top x, \\ x_l &= x_{l-1} + \theta \sigma(W_l^\top x_{l-1}), \\ & \quad l = 2, \dots, L, \end{aligned}$$



$$\begin{aligned} x_{L+1} &= W_{L+1}^\top x_L, \\ f_W(x) &= v^\top x_{L+1}, \\ \sigma &= \text{ReLU}. \end{aligned}$$

- Layer weights $W_l \in \mathbb{R}^{m_{l-1} \times m_l}$ trained by G.D.:

$$\begin{aligned} W_l^{(t+1)} &= W_l^{(t)} - \eta \nabla_{W_l} L_S(W_1, \dots, W_{L+1}), \\ L_S(W) &:= L_S(W_1, \dots, W_{L+1}) = \frac{1}{n} \sum_{i=1}^n \ell(y_i \cdot f_W(x_i)), \\ \mathcal{E}_S(W) &:= \frac{1}{n} \sum_{i=1}^n -\ell'(y_i \cdot f_W(x_i)) = \text{surrogate error}. \end{aligned}$$

Assumptions

- **Gaussian initialization:** $[W_l^{(0)}]_{i,j} \stackrel{\text{i.i.d.}}{\sim} N(0, 2/m_l)$.
- **Separability by random feature model:** there exists $f(x) = \mathbb{E}_{u \sim N(0,1)} [c(u) \sigma(u^\top x)]$, $\|c(\cdot)\|_\infty \leq 1$, such that $y \cdot f(x) \geq \gamma > 0$ for all $(x, y) \in \text{supp } \mathcal{D}$.
- **Normalized input data:** $\|x\|_2 = 1 \forall x$.
- **Widths of same order:** $m_{L+1} = \Theta(m_L)$; denote smallest layer width as $m = m_L \wedge m_{L+1}$.
- **Residual scaling:** $\theta = 1/\Omega(L)$.

Main Theorems

- **Weights stay close to init.** and **bound for surrogate error:** denote τ -neighborhood of init. by $\mathcal{W}(\tau) := \{(W_1, \dots, W_{L+1}) : \|W_l - W_l^{(0)}\|_F \leq \tau \forall l\}$.

Let \tilde{O} hide logarithmic terms in L, n, δ^{-1} .
If $\tau = \tilde{O}(\gamma^{12})$, $\eta = \tilde{O}(\tau m^{-\frac{1}{2}} \wedge \gamma^4 m^{-1})$, and $K\eta = \tilde{O}(\tau^2 \gamma^4)$, then provided $m \geq \tilde{O}(\tau^{-\frac{4}{3}} d \gamma^{-2})$, w.h.p.

- (i) $W^{(k)} \in \mathcal{W}(\tau)$ for all $k \in [K]$.
- (ii) There exists $k \in \{0, \dots, K-1\}$ with $\mathcal{E}_S(W^{(k)}) \leq C \cdot m^{-\frac{1}{2}} \cdot (K\eta)^{-\frac{1}{2}} (\log \frac{n}{\delta})^{\frac{1}{4}} \cdot \gamma^{-2}$.

- **Bound for Rademacher complexity in $\mathcal{W}(\tau)$:** for $m \geq \tilde{O}(\tau^{-\frac{4}{3}} d)$ and $f_{\mathcal{W}(\tau)} := \{f_W : W \in \mathcal{W}(\tau)\}$,

$$\mathfrak{R}_n(f_{\mathcal{W}(\tau)}) \leq C_2 \left(\tau^{\frac{4}{3}} \sqrt{m \log m} + \frac{\tau \sqrt{m}}{\sqrt{n}} \right).$$

- **Bound for test error:** for $m \geq m_{\text{res}}^*$,

$$\begin{aligned} m_{\text{res}}^* &= \tilde{O}(\text{poly}(\gamma^{-1})) \cdot \max(d, \varepsilon^{-14}), \\ n_{\text{res}} &= \tilde{O}(\text{poly}(\gamma^{-1})) \cdot \varepsilon^{-4}, \\ \eta_{\text{res}} &= O(\gamma^4 \cdot m^{-1}), \quad K_{\text{res}} = \tilde{O}(\text{poly}(\gamma^{-1})) \cdot \varepsilon^{-2}, \end{aligned}$$

G.D. with step size η_{res} finds $W^{(k^*)}$, $k^* \leq K_{\text{res}}$, s.t.
 $\mathbb{E}[\mathbb{1}(y \neq \text{sign}(f_{W^{(k^*)}}(x)))] \leq 2\mathbb{E}[-\ell'(y f_{W^{(k^*)}}(x))] \leq \varepsilon$.

Comparison with Non-Residual Results

- **Above results at most log dependence on L .**
- **Width and sample requirement is reduced:** $m_{\text{nonres}}^* > \text{poly}(L) m_{\text{res}}^*$, $n_{\text{nonres}} > \text{poly}(L) n_{\text{res}}$.
- **Step size and iterations required are better:** $\eta_{\text{nonres}} < \text{poly}(L^{-1}) \eta_{\text{res}}$, $K_{\text{nonres}} > \text{poly}(L) K_{\text{res}}$.
- **Distance from initialization** is key to above bounds. $\tau_{\text{nonres}} > \text{poly}(L) \tau_{\text{res}}$, $\tau_{\text{res}} = \tilde{O}(\gamma^{-4} \varepsilon^{-1} m^{-1/2})$.

Key Ingredients for Proof

- **Backpropagation and forward propagations are bounded independent of depth:** if x_l represents layers from input x to l -th layer, and $b_l(x)$ represents layers l to final layer, then

$$\|x_l\|_2 \leq C, \quad \|b_l(x)\|_2 \leq C.$$

- **Network output is almost linear in $\mathcal{W}(\tau)$:** for m large and $\widehat{W}, \widetilde{W} \in \mathcal{W}(\tau)$,

$$f_{\widehat{W}}(x) \approx f_{\widetilde{W}}(x) + \langle \widehat{W} - \widetilde{W}, \nabla_{\widetilde{W}} f_{\widetilde{W}}(x) \rangle.$$

- **Loss is Lipschitz and almost convex in $\mathcal{W}(\tau)$:** for m large and $\widehat{W}, \widetilde{W} \in \mathcal{W}(\tau)$,

$$\begin{aligned} \|\nabla_{W_l} L_S(\widehat{W})\|_F &\leq C \theta^{1(2 \leq l \leq L)} \sqrt{m}, \\ L_S(\widehat{W}) - L_S(\widetilde{W}) &\gtrsim \langle \widehat{W} - \widetilde{W}, \nabla_{\widetilde{W}} L_S(\widetilde{W}) \rangle. \end{aligned}$$

- **Width at most logarithmic in L** is required for above approximations, rather than usual poly/exponential.

How Does Residual Architecture Help?

- **Skip connections and scaling factor θ prevents blowup** by forcing Lipschitz constant of network output to be bounded independent of depth:

$$\begin{aligned} \|x_l\|_2 &= \|(I + \theta \Sigma_l W_l^\top) x_{l-1}\|_2 \leq (1 + C\theta) \|x_{l-1}\|_2 \\ &\leq \dots \leq (1 + C\theta)^L \|x\|_2. \end{aligned}$$

- **Representations from earlier layers are not lost in forward propagation**, allowing separability of R.F. model in first layer to persist through all layers. If α can separate at margin γ in first layer, then

$$y \cdot \langle \alpha, x_l \rangle = \langle \alpha, x_1 \rangle + \theta \sum_{l'=2}^l \langle \alpha, \sigma(W_{l'}^\top x_{l'-1}) \rangle \gtrsim \gamma.$$