

Hieher.

(X, Y) : Distributed as $X_i = f_0(X_i) + \epsilon_i$, iid $X_i \in [0, 1]^d$, $\epsilon_i \sim N(0, 1)$

f_0 takes the form $f_0 = g_0 \circ g_{-1} \circ \dots \circ g_1 \circ g_0$ and each $g_i : [a_i, b_i]^{d_i} \rightarrow [a_{i+1}, b_{i+1}]^{d_{i+1}}$

$d_0 = d$, $d_{q+1} = 1$, and each g_i is t_i -variate where $t_i \leq d_i$ and $g_i \in \mathcal{C}_{t_i}^k([a_i, b_i]^{d_i}, K)$ $i \in [q]$.

Network : $f : \mathbb{R}^{p_0} \rightarrow \mathbb{R}^{t_{q+1}}$, $\vec{x} \mapsto f(\vec{x}) = W_L b_L W_{L-1} b_{L-1} \dots W_1 b_1 W_0 x$.

Class : s.t. $\max_{j=0, \dots, L} \|W_j\|_\infty \vee \|b_j\|_\infty \leq 1$, $\sum_{j=0}^L \|W_j\|_0 + \|b_j\|_0 \leq S$, $\|f\|_\infty \leq F$.

Required Network Complexity

- (i) $F \geq \max(K, 1)$
- (ii) $\sum_{i=0}^q \log_2(4t_i \vee 4p_i) \log_2 n \leq L \lesssim n \phi_n$
- (iii) $n \phi_n \approx \min_{j \in [0, L]} p_j$
- (iv) $S \lesssim n \phi_n \log n$

where $\beta_i^* := \beta_i \prod_{k=i+1}^L (p_k \wedge 1)$, $\phi_n := \max_{i=0, \dots, q} n^{-\frac{2\beta_i^*}{2\beta_i^* + t_i}}$.

Results

\exists constants C, C' dependent on $g, \vec{q}, \vec{t}, \vec{\beta}, K$ such that any \hat{f}_n in network class above satisfying required complexity,

(i) if $\Delta_n(f_n, f_0) := E_{f_0} \left[\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{f}_n(X_i))^2 - \inf_{f \in \mathcal{F}(t_{q+1}, \vec{t})} \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2 \right] \leq C \phi_n L \log^2 n$, then,

(ii) if $\Delta_n(f_n, f_0) \geq C \phi_n L \log^2 n$, then

$$\frac{1}{C} \Delta_n(f_n, f_0) \leq R(f_n, f_0) \leq C' \Delta_n(f_n, f_0).$$

* Under the assumption that if X has Lebesgue density with positive lower & upper bounds + $t_i \leq \min\{d_0, \dots, d_{i-1}\}$ $\forall i \in [q]$, then $\inf_{f_n} \sup_{f_0} R(\hat{f}_n, f_0) \geq C \phi_n$ where inf is over all estimators.

Andra

$$(X, y) : \|X\|_2 = 1, |y| \leq 1, (X, y) \sim D \text{ over } \mathbb{R}^d \times \mathbb{R}, \left(H_{10}^{\infty} = E_{w \sim M_{10,2}} [X^T X, \mathbb{1}(w^T X \geq 0, w^T X_1 \geq 0)] \right) \Rightarrow \lambda_{\min}(H^{\infty}) > 0$$

Network Class : $f_{w,a}(x) = \frac{1}{\sqrt{m}} \sum_{i=1}^m a_i \delta(w_i^T x), \quad c(z) = \max(0, z)$

Training : $\Phi(w) := \frac{1}{2} \sum_{i=1}^n (y_i - f_{w,a}(x_i))^2, \quad w_r(k+1) - w_r(k) = -\eta \frac{\partial \Phi(w)}{\partial w_r} = -\eta \frac{\alpha}{\sqrt{m}} \sum_{i=1}^n (f_{w(a),a}(x_i) - y_i) \mathbb{1}(w_r(k)^T x_i \geq 0) x_i$
 where $\eta > 0$ is the learning rate. $w(0) \sim \mathcal{N}(0, K^2 I), \quad a_r \sim \text{unif}(-1, 1)$

Required Network Complexity : (i) $K = O\left(\frac{\varepsilon d}{\sqrt{n}}\right), m = \Omega\left(\frac{n}{\lambda_0^4 K^2 \delta^4 \varepsilon^2}\right), \eta = O\left(\frac{\lambda_0}{n^2}\right)$
 (ii) $K = O\left(\frac{\lambda_0 \delta}{n}\right), m \geq K^2 \text{poly}(n, \lambda_0^{-1}, \delta^{-1})$

Results (i) With probability $1 - \delta, \forall k = 0, 1, 2, \dots, \|y - u(k)\|_2 \leq \sqrt{\sum_{i=1}^k (1 - \eta \lambda_i)^{2k} (u_i^T y)^2} \pm \varepsilon$
 where $u_i = f_{w,a}(x_i), v_i$ are from $H^{\infty} = \sum_{i=1}^n \lambda_i v_i v_i^T$ (orthonormal eigenvectors).

(ii) For $\delta > 0$, loss function $\mathcal{L} : \mathbb{R} \times \mathbb{R} \rightarrow [0, 1]$ that's 1-Lipschitz, in first argument with $\mathcal{L}(y, y) = 0$. with prob $1 - \delta$,
 for $k \geq \Omega\left(\frac{1}{\eta \lambda_0} \log \frac{n}{\delta}\right), E_{(x,y) \sim D} [\mathcal{L}(f_{w(a),a}(x), y)] \leq \sqrt{\frac{2 \mathbb{I}(H^{\infty})^{-1}}{n}} y + O\left(\sqrt{\frac{\log(\frac{n}{\lambda_0 \delta})}{n}}\right)$ ↓
over train sample
& initialization.

Under the assumption that $P_0(\lambda_{\min}(H^{\infty}) \geq \lambda_0) \geq 1 - \frac{\delta}{3}$ for fixed η_0 .
 (iii) If $y_i = g(x_i) = \alpha (g^T x_i)^p, \forall x_i \in [0, 1],$ for $p=1$ or $p=2$ for $\mathcal{L} \in \mathcal{M}_t, p \in \mathbb{R}^d, \alpha \in \mathbb{R}$, Then we have

$$\sqrt{y^T (H^{\infty})^{-1} y} \leq 3p |\alpha| \cdot \|p\|_2$$

Chen (2020)

(i) $H_0(p) = E_p[u^2 \langle \nabla_\theta h(\theta, x) \rangle, \nabla_\theta h(\theta, x) \rangle]^+ + E_p[h(\theta, x) h(\theta, x)]$ is positive definite.

$(x, \tau) : \|x\|_2 \leq 1, |\tau| \leq 1, x \in \mathbb{R}^d, (ii) \exists P_{true}$ with $D_p(P_{true} \| p_0) < \infty$ s.t. $y = \int u h(\theta, x) P_{true}(\theta, u) d\theta du$

Network : $S = \{(x_i, y_i), \dots, (x_n, y_n)\}$.

Class : $f(p, x) = \alpha \int_{\mathbb{R}^{d_H}} u h(\theta, x) p(\theta, u) d\theta du, p_\tau$ is generated from the PDE

Training : $\frac{dp_\tau}{d\tau} = -\nabla_u [p_\tau(\theta, u) g_1(t, \theta, u)] - \nabla_\theta [p_\tau(\theta, u) g_2(t, \theta, u)] + \lambda \Delta p_\tau(\theta, u)$

which is the infinite-width, continuous-time version of noisy gradient descent.

Required Network Complexity

$$\lambda_0 := \sqrt{\frac{\Lambda}{n}} \text{ where } \Lambda = \lambda_{\max}(H(p_0)).$$

Activation $h(\theta, x)$ is s.t. $h(\theta, x) = \tilde{h}(\theta^T x)$ where \tilde{h} satisfies $|\tilde{h}(z)| \leq G, |\tilde{h}'(z)| \leq G, |\tilde{h}''(z)| \leq G$.

$$(i) \alpha \geq \partial \sqrt{A_2^2 + \lambda_0^2}, \lambda_0^2 R^{-1}, R = \min\{\sqrt{d+1}, p_0 y(G, \log(\frac{1}{\lambda_0}))^{-1} \lambda_0^2\}, A_1 = 2G(d+1) + 4G\sqrt{d+1}, A_2 = 16G\sqrt{d+1} + 4G$$

$$(ii) \alpha \geq \sqrt{n} \lambda_0 > 0.$$

Results : (i) $L(p_0) = E_S[\phi(f(p, x), y)] \leq 2 \exp(-2\alpha^2 \lambda_0^2 t) + 2A_1^2 \lambda_0^2 \alpha^{-2} \lambda_0^{-4}$

$$D_K(p_\tau \| p_0) \leq 4A_2^2 \alpha^{-2} \lambda_0^{-4} + 4A_1^2 \lambda_0^2 \alpha^{-2} \lambda_0^{-4}.$$

$$(ii) \text{ For any } \delta > 0, \text{ with probability } 1-\delta, E_S[\phi(f(p^*, \tau), y)] \leq (\delta G + 1) \sqrt{\frac{D_{KL}(P_{true} \| p_0)}{n}} + 6 \sqrt{\frac{\log(\frac{2}{\delta})}{2n}}$$

where p^* is the minimizer of $L(p) + \lambda D_{KL}(p \| p_0)$, $x^{p^*}(y^*, y) = \mathbb{1}(y^* y < 0)$

Here, (x_i, y_i) are iid sampled from unknown distribution D .

Training : Noisy GD, $\{(\theta_j, u_j)\}_{j=1, \dots, n}$ initialized by $P_0(\theta, u) \sim \mathcal{N}(0, I_{d_H+d+1})$