

Prev: dedicated algo.

Now: generic. [DLS22']

$\{(x_i, y_i)\}_{i=1 \dots n}$  samples,  $f^*(x) = g(\underline{w}^T x)$ . assume  $\deg(f^*) = p$ .

$$\mathbb{R}^d \rightarrow \mathbb{R}^p$$

$$(d \gg p)$$

$$y_i = f^*(x_i) + \epsilon_i.$$

$$\epsilon_i \sim \mathcal{N}(0, \sigma^2).$$

$f_\theta(x)$ : Shallow NN,  $f_\theta(x) = a^T \sigma(w^T x + b) = \sum_{j=1}^m a_j \sigma(w_j \cdot x + b)$ .

Squared loss:  $L_\theta(\theta) = \frac{1}{n} \sum_{i=1}^n (f_\theta(x_i) - f^*(x_i))^2$ ,  $a_j = \mathcal{Z}(-1, 1)$ ,  
 $w_j \sim \mathcal{N}(0, \frac{1}{d} \mathbb{I}_d)$ .

Kernel regime:  $n \times d^p$ . [GMMMG19']

This paper:  $n \leq d^2 r + dr^p$ ,

(Feature learning).

Idea: First grad. step (learn feature  $W$ ).

$$n \geq O(d^2)$$
,

$$r \text{ directions } d^2 r,$$

$w^{(1)}$  depends on  $(x, y)$ .

Remaining: Learn on.  $w^{(1)}$  and  $w^{(2)}$  not independent.

$dr^p$ .  $d$  can be removed by resampling. (transfer learning setting).  
subspace

(similar to exhaustiveness of PhD estimator).

Assumptions: ① Non-degeneracy of  $H = \mathbb{E}[\nabla^2 f^*(x)]$  ( $f^* \in \mathcal{Z}$ )

$$\text{rank}(H) = r, \text{span}(H) \supseteq f^*, \text{ denote } k = \frac{\|H^+\|}{\sqrt{r}}.$$

condition #.

Crit C.R. argument shows necessary,

② symmetric.

$$a_j = -a_{m-j}, b_j = b_{m-j}, \sigma_j = \sigma_{m-j}.$$

$$\text{s.t. } f_{\theta_0}(x) = 0,$$

Study gradient:

$$\nabla_{w_j} \ell_0(\theta) = \mathbb{E}_{x \sim D} [ \underbrace{2(f_0(x) - f^*(x))}_{=0 \text{ by symmetric init.}} \nabla_{w_j} f_0(x) ]$$

$\nabla_{w_j} f_0(x)$

$$= -2 \mathbb{E}_{x \sim D} [ f^*(x) \nabla_{w_j} f_0(x) ], \quad \text{Recall } \sum_{j=1}^m a_j b_l w_j \cdot x + b.$$

$$\begin{aligned} &= -2a_j \mathbb{E}_{x \sim D} [ f^*(x) x \sigma'(w_j \cdot x) ], \\ &= \mathbb{E}_{x \sim D} [ \nabla f^*(x) \sigma'(w \cdot x) + w f^*(x) \sigma''(w \cdot x) ], \end{aligned}$$

Stein's Lemma

Take Hermite expansion over  $w$ .

$$f^*(x) = \sum_{k=0}^p \frac{\langle c_k, H_k(x) \rangle}{k!} \quad \text{define } c_k = \mathbb{E}_x [\nabla_x^k f^*(x)]. \quad c_0 = H.$$

$$\sigma'(x) = \sum_{k \geq 0} \frac{c_k}{k!} H_k(x). \quad c_k : \text{Hermite coeff. of } \sigma'(x) = I_{x \geq 0}. \\ c_1 = \frac{1}{2}, \quad c_2 = \frac{1}{2\pi}.$$

$$= \sum_{k=0}^{p-1} \frac{c_{k+1} \mathbb{E}_x [\nabla_x^{k+1} f^*(x)] (w^{\otimes k})}{k!} + w \sum_{k=0}^p \frac{c_{k+2} \mathbb{E}_x [\nabla_x^k f^*(x)] (w^{\otimes k})}{k!}$$

$$= \sum_{k=0}^{p-1} \frac{1}{k!} [c_{k+1} c_{k+1} (w^{\otimes k}) + c_{k+2} w (w^{\otimes k})].$$

$$= \frac{Hw}{2\pi} + \frac{1}{2} [c_3 c_3 (w, w) + c_4 w c_2 (w, w)] + \frac{1}{6} [\dots] + \dots \\ O(d^{-1/2}) \quad O(d^{-1}) \quad O(d^{-3/2}) \quad \overset{T}{\text{higher-order}},$$

Note:  $\|c_{k+1} w^{\otimes k}\|, \|c_k w^{\otimes k}\| = O(d^{-k/2})$ ,

$$= O(d^{-1/2}).$$

$$\|\nabla_{w_j} \hat{\ell}_0(\theta) - \nabla_{w_j} \ell_0(\theta)\| \leq \sqrt{\frac{d}{n}},$$

$$d^{-1/2} \geq \sqrt{\frac{d}{n}} \Rightarrow n \geq d^2,$$

---

**Algorithm 1:** Gradient-based training

---

**Input:** Learning rates  $\eta_t$ , weight decay  $\lambda_t$ , number of steps  $T$

**preprocess data**

$$\left| \begin{array}{l} \alpha \leftarrow \frac{1}{n} \sum_{i=1}^n y_i, \beta \leftarrow \frac{1}{n} \sum_{i=1}^n y_i x_i \\ y_i \leftarrow y_i - \alpha - \beta \cdot x_i \text{ for } i = 1, \dots, n \end{array} \right. \quad \text{Normalizing } y,$$

**end**

$$W^{(1)} \leftarrow W^{(0)} - \eta_1 [\nabla_W \mathcal{L}(\theta) + \lambda_1 W]$$

$$\text{re-initialize } b_j \sim N(0, 1) \quad \lambda_1 = \frac{1}{b_1},$$

**for**  $t = 2$  to  $T$  **do**

$$\left| \begin{array}{l} a^{(t)} \leftarrow a^{(t-1)} - \eta_t [\nabla_a \mathcal{L}(\theta^{(t-1)}) + \lambda_t a^{(t-1)}] \end{array} \right.$$

**end**

**return** Prediction function  $x \rightarrow \underline{\alpha + \beta \cdot x + a^T \sigma(Wx + b)}$

---

$$w^{(1)} = -y, \nabla_{w_0} L(\theta),$$

dominated by first order  $\frac{Hw}{\sqrt{2x}}$ .

dominated by  $Hw$ . ( $w \in \mathbb{R}^{d-1}$ ).

$Cf^*$ . ( $\because \text{span}(v) = f^*$ ).

In  $f^*$  where  $\text{rank}(f^*) = r$ , choose parameter with  $\alpha$  ( $\kappa$ ),

repeat (inequality regime)  $n \geq r^p$ , width  $m \geq r^p$ .

Remaining steps (learn  $a$ ),

build KKT to approx.  $x^{(k)}$   
 $\Rightarrow$  approx. high-dim.

1. Contract  $a^*$  s.t.  $L(a^*, w^{(1)}, \zeta) \ll 1$ .  $\|a^*\| = \tilde{\mathcal{O}}\left(\frac{r^p k^{2p}}{\sqrt{m}}\right)$ . make use of  $w^{(1)}$ .

2.  $\exists \lambda > 0$  s.t.  $L(a^{(T)}, w^{(1)}, \zeta) \ll 1$ .  $\|a^{(T)}\| \leq \|a^*\|$  where  $T = \tilde{\mathcal{O}}(y^{-1} \lambda^{-1})$

3. Rademacher Gen. Bound. (Standard),  $\tilde{\mathcal{O}}\left(\sqrt{\frac{dr^p k^{2p}}{n}} + \sqrt{\frac{r^p k^{2p}}{m}} + \frac{1}{n^{1/4}}\right)$ ,

$r$  directions.

$n \geq dr^p$   $m \geq r^p$ . uniform conv.

$$n \geq \underbrace{d^2 r}_{\text{after 1 rep.}} + \boxed{dr^p}.$$

first grad. step

$n \geq d^2 r$ . (full iteration).

$$m \geq r^p.$$

Lower bound: (Necessity of span(H)) =  $\ell^*$ .

Construct  $F_p$  of poly. with deg p. Each function depends on single relevant direction.  
not varying assumption 2.  $q_{\text{queries}} \leq Cr = 1$ .

$$\text{tolerance } \tau \leq \frac{\log^{p/4}(\text{card})}{d^{p/4}}$$

to output  $f \in F_p$  with  $\text{err} \leq \tau$ .

$$\boxed{\tau \propto \frac{1}{m} \Rightarrow m \geq d^{p/2}, \quad \text{optimal.}}$$

Pf: Recall s.t.  $g(x, y) \rightarrow \text{output } \hat{g}$  with  $|\hat{g} - \mathbb{E}_{x,y} g(x, y)| \leq \tau$ .

CfQ  $g(x, y) = \underbrace{y h(x)}_{\substack{\text{inner} \\ \text{convolution}}}.$

$$\mathbb{E}_{w_j} g_j(\theta) = \mathbb{E}_{x,y} [y h(x)] \text{ where } h(x) = -2w_j \times g'(w_j \cdot x),$$

Idea: Construct function class with small pairwise corrections.

Lemma:  $\left| \mathbb{E}_{x \sim D} [f(x)^2] - 1 \right| \leq \varepsilon$ .  $\left| \mathbb{E}_{x \sim D} [f(x)g(x)] \right| \leq \varepsilon$ . If  $f, g \in F$ .  $f \neq g$ .  
 $\Rightarrow$  Require  $q \geq \frac{|F|(\tau^2 - \varepsilon)}{2}$  queries to have  $\text{err} \leq 2 - 2\varepsilon$ .

How to find such  $f$ ?

Fact:  $\exists e^{C\varepsilon^2 d}$  vectors in  $\mathbb{R}^d$  s.t. their inner product  $\leq \varepsilon$ .

$$f_u(x) = \frac{He^P(u \cdot x)}{\sqrt{P!}}. \quad \text{Construct}$$

$$\mathbb{E}_{x \sim D} [f_u(x) \cdot f_v(x)] = (u \cdot v).$$

$$\therefore |u \cdot v| \leq \varepsilon \Rightarrow \left| \mathbb{E}[f_u(x) \cdot f_v(x)] \right| \leq \varepsilon^P.$$

By lemma.  $q \geq \frac{e^{C\varepsilon^2 d} (\tau^2 - \varepsilon^P)}{2}$ .

$$e^{c\varepsilon^2 d} \leq \frac{2a}{\varepsilon^{2-\zeta^\alpha}}.$$

$$\text{take } \varepsilon = \sqrt{\frac{\log(2a(cad)^{\alpha/2})}{cad}},$$

$$\tau^2 \leq \frac{(xy^{\alpha/2}(qd))}{d^{\alpha/2}},$$

Last time, Kernel  $n \asymp d^p$ . ( $f^*(x) = g(Lx)$ ,  $L: \mathbb{R}^d \rightarrow \mathbb{R}^r$  ( $d \gg r$ )).  $\deg(f^*) = p$ .  
 [DW22] One grad. step + learn second layer.  
 $n \asymp d^2 r + dr^p$ .

Today [BEJ+22]. how one grad. step improves rep. [Based on observation "non-kernel" behavior often occurs in early phase, especially in large lr]  
 High-dim Regime. (Follow not from [BEJ+22]).

Asymptotic: data size  $n$ , input dimension  $d$ . NN width  $N \rightarrow \infty$ .

FJGF+20, FDP+20

$$n/d \rightarrow \psi_1, N/d \rightarrow \psi_2. \quad \psi_1 \uparrow \rightarrow \text{sample size } T \quad n/N \text{ comparable.}$$

$$\psi_2 T \rightarrow \text{NN width } T. \quad \text{so } N \asymp n.$$

$$f(x) = \frac{1}{\sqrt{N}} \sum_{i=1}^N a_i \sigma(w^T x, w_i),$$

$$= \frac{1}{\sqrt{N}} a^T \sigma(w^T x) \quad \text{with required vars.}$$

$$x \in \mathbb{R}^d, \quad \underbrace{\sim}_{1 \times N} \quad \text{and } \underbrace{a \in}_{d \times 1}$$

Limitations of prev method v:

Build kernel on  $x \mapsto \sigma(w_0^T x)$ : Conjugate kernel.

KF method (C, k, NTK). Rotation invariant kernel.

[EK10, HC20, NTZ20, Buz21].

$$\inf_{\lambda \gg 0} R_{\text{NTK}}(\lambda) \gtrsim \|P_{\geq 1} f^*\|_{L^2}^2 + o_d(1),$$

$\swarrow$   $\smile$

*projection to* *cannot perform better than linear method.*  
*deg > 1 poly,*

This paper: small lr  $y = \Theta(1)$ . Better than KF, still linear regime.

$$\text{large lr } y = \Theta(\sqrt{N}) \quad \text{rank } < \|P_{\geq 1} f^*\|_{L^2}^2,$$

Assumption: Recall  $f(x) = \frac{1}{\sqrt{N}} a^T \sigma(w^T x)$ .

① Init:  $\sqrt{d} \cdot [w_0]_{ij} \stackrel{iid}{\sim} N(0, 1)$ ,

$\sqrt{N} \cdot [a]_j \stackrel{iid}{\sim} N(0, 1)$ .

$\Rightarrow \frac{1}{n}$  Mean-field scaling.

②  $y_i = f^*(x) + \varepsilon_i$ .  $\varepsilon_i \sim N(0, \sigma_\varepsilon^2)$ .

$f^*$ : Lipschitz.  $\|f^*\|_{C^2} = \Theta_d(1)$ .

In particular, we study single-index.

$f^*(x) = \sigma^*(\langle x, p^* \rangle)$ ,  $p^* \in \mathbb{R}^d$ .

(where  $E[\sigma^*(z)] \neq 0 \Rightarrow \ell^* = 1$ )

③  $\delta: E[\sigma(z)] = 0$ .  $E[\sigma(z_t)] \neq 0$ .

first three deriv. bounded a.s.

This paper:

First step: Learn  $w$ . (Goal: Gradient matrix is of rank-1, which contains info. of labels  $y$ ).

$$w_{t+1} = w_t - \gamma \frac{\partial L}{\partial w_t} = w_t + \gamma \sqrt{N} G_t,$$

where  $G_t = \frac{1}{n} X^T \left[ \underbrace{\left( \frac{1}{\sqrt{N}} (y - \frac{1}{\sqrt{N}} \sigma(Xw_t a)) a^T \right)}_{n \times 1} \otimes \underbrace{\sigma'(Xw_t)}_{n \times n} \right]$

(neglect deviation).

Orthogonal decomposition of  $\sigma$ :

$\sigma(z) = \mu_1 z + \sigma_\perp(z)$ , where  $\mu_1 = E[z\sigma(z)] \neq 0$ ,

$E[\sigma_\perp(z)] = E[z\sigma_\perp(z)] = 0$ .

$E[\sigma_\perp(z)^2] = \mu_2^2$ . where  $\mu_2 = \sqrt{E[\sigma_\perp(z)^2] - \mu_1^2}$ ,

Denote  $G_0 := \frac{1}{\sqrt{nN}} (w_1 - w_0)$ , rank-1 matrix  $A := \frac{1}{n\sqrt{N}} X^T y a^\top$ .

w.h.p.,  $\|G_0 - A\| \leq \frac{C \alpha^2 n}{\sqrt{n}} \|G_0\|$ .

(Depends on the decomposition above. App. B.1.1).

Expect: (Want  $y, w_1$  should have alignment with  $f^*$ ,

$$f^*(x) = \mu_0^* + \mu_1^* \langle x, \beta^* \rangle + P_{\perp} f^*(x)$$

denote  $\|P_{\perp} f^*\|_{L^2} \rightarrow \mu_2^*$  as  $d \rightarrow \infty$ .

Learning rate regime:

$$\sqrt{d} \cdot \mathbb{E}[w_0]_{ij} \stackrel{iid}{\sim} N(0, 1),$$

Small  $\text{lr}$ :  $\eta = \Theta(1) \Rightarrow \|w_1 - w_0\| \asymp \|w_0\|$

Large  $\text{lr}$ :  $\eta = \Theta(\sqrt{N}) \underset{\text{Adhere to } \mu P\text{-scaling}}{\Rightarrow} \|w_1 - w_0\|_F \asymp \|w_0\|_F$ .

(Adhere to  $\mu P$ -scaling),

$N$  is the width.

Thm 3: Can derive asymptotic limit of leading singular value  $s_i(w_1) \propto \frac{\eta}{\text{lr}}$ .  
 When  $\eta = \Theta(1)$ )

corr. left sing. vector  $u_1$

$$|\langle u_1, \beta^* \rangle|^2,$$

$$\forall i > 1, |s_i(w_1) - s_i(w_0)| = o_d(1).$$

spiked model, (Figure 3. of paper). buck doesn't change.

Book: KMT4ML, Cao et al. Ch 2.5, 2.6.

- After first rep., obtain feature map  $\mathbf{x} \mapsto \sigma(\mathbf{w}_1^T \mathbf{x})$ .

train conjugate kernel on top.

similar to LSSZ.

Remaining step on  $\alpha$ :

Do ridge regression on  $\alpha$  with train samples  $(\tilde{\mathbf{X}}, \tilde{\mathbf{y}})$ ,

$\therefore \mathbf{w}_1$  correlated with  $(\mathbf{X}, \mathbf{y})$ ,

$$\hat{\alpha} = \underset{\alpha}{\operatorname{argmin}} \left( \frac{1}{n} \|\tilde{\mathbf{y}} - \bar{\Phi} \alpha\|^2 + \frac{\lambda}{N} \|\alpha\|^2 \right).$$

$$\text{where } \bar{\Phi} = \frac{1}{\sqrt{N}} \sigma(\tilde{\mathbf{X}} \mathbf{w}_1),$$

$n \times N,$

[HLZ0, LSSZ]

- Gaussian Equivalence property, ( $y = \theta(1)$ ) precise asymptotics,

$$R_{CK}(\lambda) \asymp R_{GE}(\lambda)$$

$$\text{for } \phi_{GE}(x) = \frac{1}{\sqrt{N}} (\mu_1 \mathbf{w}^T \mathbf{x} + \mu_2 \underline{z}),$$

$$\underline{z} \sim N(0, 1),$$

nonlinear kernel behaves like noisy linear method.

Can be computed explicitly.

Conclusion: Can improve over initial  $c_k$ .

$$\text{but } R_{GE}(\lambda) \geq \|P_{S_1} f^*\|_{C^2}^2,$$

$\gamma = O(\sqrt{N})$ , (Cannot derive asymptotic behavior).

Given  $\mathbf{w}_1$ , construct second layer  $\tilde{\alpha}$  s.t. (Exist good sol.  $\tilde{\alpha}$ ),

$$\tilde{f}(\mathbf{x}) = \frac{1}{\sqrt{N}} \tilde{\alpha}^T \sigma(\mathbf{w}_1^T \mathbf{x}) \text{ has risk } \leq k^*,$$

$$\text{where } k^* \stackrel{\Delta}{=} \inf_{k \in \mathbb{R}} \left[ \underbrace{\sigma^*(z_1)}_{\text{feature}} - \underbrace{\mathbb{E}_{z_2} [\sigma(k z_1 + z_2)]}_{\text{student}} \right]$$

$$z_1, z_2 \sim N(0, 1).$$

(?)

For suff. large  $\psi_1 \triangleq \frac{n}{d}$ .  $f_1(\lambda) \leq 10k^* + C\left(\sqrt{k^*} \cdot \sqrt{\frac{d}{n}} + \frac{d}{n}\right)$ .

$\overset{t}{\underset{\text{ridge estimator}}{\wedge}}$  w.p. 1 when n.d.  $N \rightarrow \infty$ .  
ridge penalty  $n^{\varepsilon-1} < n^{-1}\lambda < n^{-\varepsilon}$   
for some  $\varepsilon > 0$ .

If  $\|P_{S_1} f^*\|_{L^2} > 10k^*$ , then it outperforms linear method.

e.g.:  $\sigma = \sigma^* = \tanh$ .