

Lecture 13: Test error of Random Features models

"Generalization error of random features and kernel methods: hypercontractivity and kernel matrix concentration"

[Mei, M., Montanari, 2021]

Here: $n = d^l$ "polynomial scaling"

Next lecture: $n = cd$ "proportional scaling"

Random Features model:

RKHS: $f(x; a) = \int_{\Omega} \sigma(x; \theta) a(\theta) \mathfrak{z}(d\theta)$

where $\|a\|_{L^2}^2 = \int_{\Omega} a(\theta)^2 \mathfrak{z}(d\theta) < \infty$

∞-dim space

Random features approx: $\theta_1, \dots, \theta_N \stackrel{iid}{\sim} \mathfrak{z}$

$$\hat{b}_{RF}(x, a) = \frac{1}{N} \sum_{j=1}^N \underbrace{a_j}_{\text{2nd layer weights}} \underbrace{\sigma(x; \theta_j)}_{\text{1st layer weights fixed}}$$

[Rahimi - Recht - 2008]

dim(RF) = N "finite dim. approx." of RKHS

$$M_N(x_1, x_2) = \frac{1}{N} \sum_{j=1}^N \sigma(x_1; \theta_j) \sigma(x_2; \theta_j)$$

$$\rightarrow \mathbb{E}_{\theta} [H_N(x_1, x_2)] = \int_{\mathcal{X}} \sigma(x_1, \theta) \sigma(x_2, \theta) z(d\theta)$$

"Modern" interpretation:

* 2-layer NN: train 2nd layer

* Linearized NNs:

$$b_{NN}(x; \beta) \approx b_{NN}(x; \beta_0) + \langle \beta - \beta_0, \nabla_{\beta} b_{NN}(x; \beta_0) \rangle$$

$$\beta = (a, \theta)$$

RF: linearization wrt 2nd layer weights

$$\hat{f}_{RF}(x; a) = \langle a, \nabla_a b_{NN}(x; \theta_0) \rangle$$

$$= \langle a, \varphi_N(x) \rangle$$

$$\varphi_N(x) = (\sigma(x; \theta_1), \dots, \sigma(x; \theta_N)) \in \mathbb{R}^N$$

Setting: - (X, ν) prob space $f_* \in L^2(X)$
 $X \subseteq \mathbb{R}^d$

$$\{(y_i, x_i)\}_{i \leq m} \quad y_i = f_*(x_i) + \varepsilon_i$$

$$x_i \underset{iid}{\sim} \nu \quad \varepsilon_i \underset{iid}{\sim} N(0, \sigma_\varepsilon^2)$$

- (Ω, τ) prob space weights

$$\{\theta_j\}_{j \leq N} \quad \theta_j \underset{iid}{\sim} \tau$$

RF model: $\hat{f}(x; a) = \frac{1}{N} \sum_{j=1}^N a_j \sigma(x; \theta_j)$

RF ridge regression (RFRR):

$$\hat{a}(\lambda) = \underset{a \in \mathbb{R}^N}{\operatorname{argmin}} \left\{ \frac{1}{m} \sum_{i=1}^m (y_i - \hat{f}(x_i, a))^2 + \frac{\lambda}{mN} \|a\|_2^2 \right\}$$

Test error:

$$R_{m,N}(f_*) = \mathbb{E}_{x_{\text{new}}} \left[(f_*(x_{\text{new}}) - \hat{f}(x_{\text{new}}, \hat{a}(\lambda)))^2 \right]$$

1) General assumptions

2) Example: $x \sim \text{Unif}(\mathbb{S}^{d-1})$

General assumptions:

$$(X_d, \nu_d) \quad (\Omega_d, z_d) \quad X \subseteq \mathbb{R}^d$$

Featurization map:

$$\sigma \in L^2(X \times \Omega)$$

$$\sigma : X \times \Omega \rightarrow \mathbb{R}$$

$$(x, \theta) \mapsto \sigma(x; \theta)$$

$$\bullet \Pi : L^2(\Omega) \rightarrow L^2(X)$$

$$a \mapsto f(x; a) = \int_{\Omega} \sigma(x; \theta) a(\theta) z(d\theta)$$

Π compact operator

$$\left(\Pi = \sum_{j=1}^{\infty} \lambda_j \Psi_j \Phi_j^* \right)$$

$$\mathbb{E}_x [\Psi_j(x) \Psi_k(x)] = \delta_{jk}$$

$\bullet \{ \Psi_j \}_{j \geq 1}$ orthonormal basis of $L^2(X)$

$\bullet \{ \Phi_j \}_{j \geq 1}$ ————— $L^2(\Omega)$

$$\sigma(x; \theta) = \sum_{j=1}^{\infty} \lambda_j \psi_j(x) \phi_j(\theta)$$

$$\sigma \in L^2 \quad \|\sigma\|_{L^2}^2 = \sum_{j=1}^{\infty} \lambda_j^2 < \infty$$

$$\bullet \quad H = \Pi \Pi^* : L^2(X) \rightarrow L^2(X)$$

$$H f(x) = \int_X H(x, x') f(x') \nu(dx')$$

$$H(x_1, x_2) = \int_{\Omega} \sigma(x_1; \theta) \sigma(x_2; \theta) z(d\theta)$$

$$H = \sum_{j=1}^{\infty} \lambda_j^2 \psi_j \psi_j^*$$

$$\bullet \quad U = \Pi^* \Pi : L^2(\Omega) \rightarrow L^2(\Omega)$$

$$U g(\theta) = \int_{\Omega} U(\theta, \theta') g(\theta') z(d\theta')$$

$$U(\theta_1, \theta_2) = \int_X \sigma(x; \theta_1) \sigma(x; \theta_2) \nu(dx)$$

$$U = \sum_{j=1}^{\infty} \lambda_j^2 \phi_j \phi_j^*$$

(simplified) assumptions at level (M, m, u)

$M(d)$ associated No $N(d)$

$m(d)$ ————— $m(d)$

$$u \geq \max(M, m)$$

① [Hypercontractivity] $\mathcal{X}_{\leq u} = \text{span}\{\psi_j : j \leq u\}$

$$\forall k \geq 1, \exists C > 0$$

$$\forall g \in \mathcal{X}_{\leq u}, \|g\|_{L^{2k}(X)} \leq C \cdot \|g\|_{L^2(X)}$$

$$\|g\|_{L^p} = \mathbb{E}[g^p]^{\frac{1}{p}} \quad p < q: \|g\|_{L^p} \leq \|g\|_{L^q}$$

(Jensen)

→ hypercontractivity: reverse inequ

space G "hypercontractive" $\forall g \in G, \|g\|_{L^p} \leq C \cdot \|g\|_{L^q}$

functions in G are "delocalized". $p > q$

Same thing $\Omega_{\leq u} = \text{span} \{ \Phi_j : j \leq u \}$

② [Concentration of diagonal elements]

$$H_{> m}(\alpha_1, \alpha_2) = \sum_{j > m} \lambda_j^2 \Phi_j(\alpha_1) \Phi_j(\alpha_2)$$

$$\sup_{i \leq m} \left| H_{> m}(\alpha_i, \alpha_i) - \mathbb{E}_{\alpha} [H_{> m}(\alpha, \alpha)] \right| = \underbrace{o_{d, \mathbb{P}}(1)}_{\text{small } o \text{ in probability}} \cdot \mathbb{E}_{\alpha} [H_{> m}(\alpha, \alpha)]$$

$o_{d, \mathbb{P}}$ "small o in probability"

$$X_d = o_{d, \mathbb{P}}(Y_d) \quad \frac{X_d}{Y_d} \rightarrow 0 \text{ in } \mathbb{P}.$$

Same for $U_{> m}(\theta_i, \theta_i)$

③ [Spectral gap]

• Underparametrized regime: $N \leq m$

$$\underbrace{\frac{1}{\lambda_M^2}}_{\text{gap: } \lambda_M} \sum_{k=M+1}^{\infty} \lambda_k^2 \ll N \ll \frac{1}{\lambda_{M+1}^2} \sum_{k=M+1}^{\infty} \lambda_k^2$$

→ M : subspace estimated accurately in under.

- Overp. regime: $m \leq N$

$$\frac{1}{\lambda_m^2} \sum_{k=m+1}^{\infty} \lambda_k^2 \ll m \ll \frac{1}{\lambda_{m+1}^2} \sum_{k=m+1}^{\infty} \lambda_k^2$$

$$b_* \in L^2(X) \quad b_* = \sum_{k=1}^{\infty} \langle b_*, \Psi_k \rangle_{L^2(X)} \Psi_k$$

$$P_{>m} b_* = \sum_{k=m+1}^{\infty} \langle b_*, \Psi_k \rangle_{L^2(X)} \Psi_k$$

Thm: [M, M, M, 21] Assumpt^o at level (M, m, u)

- Overp. regime: $N \geq d^{\delta} m$ for $\delta > 0$

$\lambda \in [0, \lambda_0]$, any fixed $\eta > 0$:

$$\left[R_{m,N}(b_*) = \| P_{>m} b_* \|_{L^2}^2 + o_{d,P}(1) \cdot \| b_* \|_{L^{2+\eta}}^2 \right]$$

- Under regime: $m \geq d^{\delta} N$ $\delta > 0$

$\lambda \in [0, \lambda_u]$, any fixed $\eta > 0$

$$\left[R_{m,N}(b_*) = \| P_{>M} b_* \|_{L^2}^2 + o_{d,P}(1) \dots$$

$$\hat{b}_{\text{RF}} \approx \begin{cases} P_{\leq m} b_* & N \geq d^\delta m \\ P_{\leq M} b_* & m \geq d^\delta N \end{cases}$$

Rule: ① m, N symmetric role

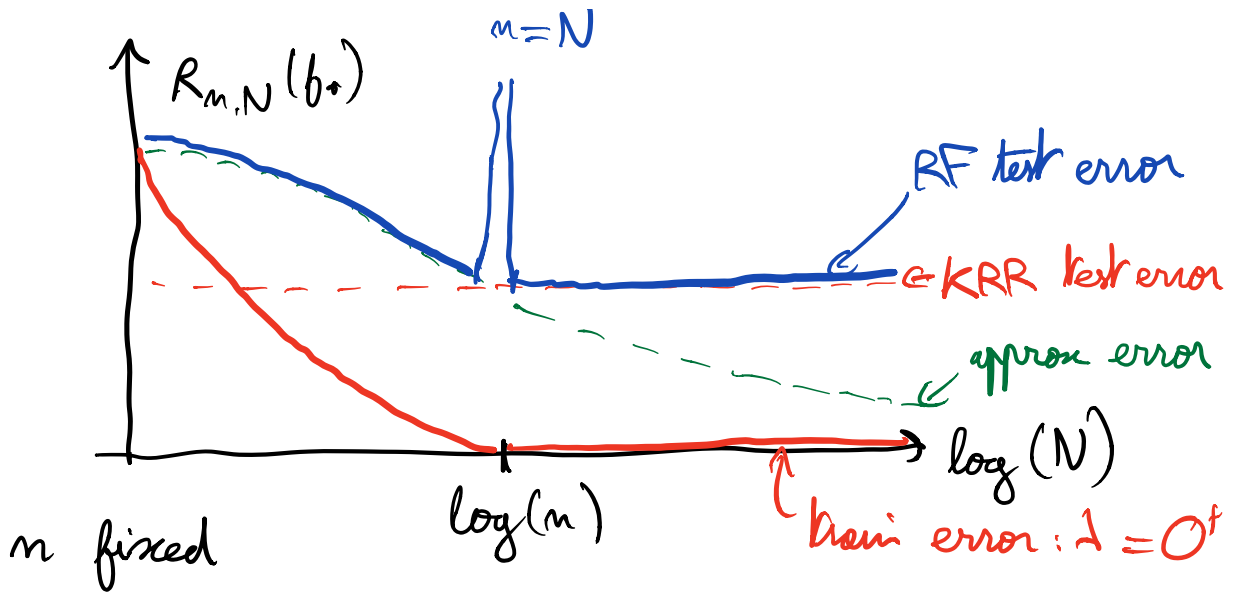
Approx. error: N finite, $m = \infty$

$$\begin{aligned} R_N^{\text{app}}(b_*) &= \inf_a \left\| b_* - \sum_{j=1}^N a_j \sigma(\cdot; \theta_j) \right\|_{L^2}^2 \\ &= \| P_{>M} b_* \|_{L^2}^2 + o_{d,P}(1). \end{aligned}$$

Statistical error: m finite, $N = \infty$ (KRR)

$$R_m^{\text{KRR}}(b_*) = \| P_{>m} b_* \|_{L^2}^2 + o_{d,P}(1)$$

$$R_{m,N}(b_*) \approx \begin{cases} R_N^{\text{app}}(b_*) & N \ll m \\ R_m^{\text{KRR}}(b_*) & N \gg m \end{cases}$$



② Optimal overparametrization: $N \geq d^\delta m$

③ Optimality interpolators: $\lambda \in (0, \lambda_*]$

$$\lambda = 0_+$$

$$y_i = b_* \phi(x_i) + \varepsilon_i$$

Example: $(X, \nu) = (\Omega, \mathbb{z}) = (\mathbb{S}^{d-1}(\sqrt{d}), \text{Unif})$
 $= \{ \alpha \in \mathbb{R}^d, \|\alpha\|_2 = \sqrt{d} \}$

$$\sigma(\alpha; \theta) = \sigma(\langle \alpha, \theta \rangle / \sqrt{d})$$

$$\sigma: \mathbb{R} \rightarrow \mathbb{R} \quad \rightarrow \quad \mathcal{O}(1)$$

$$\mathbb{E}_{\alpha, \theta}[\sigma(\langle \alpha, \theta \rangle / \sqrt{d})] = \mathbb{E}_{\alpha_1}[\sigma(\alpha_1)]$$

$$\downarrow d$$

$$N(0, 1)$$

$$\leq C$$

Thm [MMM 21] σ satisfying "genericity" conditions
 $d^{\delta+\delta} \leq m \leq d^{\delta+1-\delta}$ $d^{\delta+\delta} \leq N \leq d^{\delta+1-\delta}$

• Over regime: $N \geq d^\delta m$

$$R_{m, N}(b_*) = \|\bar{P}_{> \Delta} b_*\|_{L^2}^2 + o_{d, P}(1)$$

$\bar{P}_{> \Delta}$: project^o orthogonal to polynomials $d^0 \leq \Delta$

$\hat{b}_{\text{RF}} \rightarrow$ fit exactly polynomials $d^0 \leq \Delta$

- Under require: $n \geq d^2 N$

$$R_{n,N}(b_*) = \|\overline{P}_{>S} b_*\|_{L^2}^2 + o_{d,P}(1)$$

□

Proof: checking the conditions

Funct^o space: $L^2(S^{d-1}(\sqrt{d}), \text{Unif})$

$$L^2(S^{d-1}) = \bigoplus_{l=0}^{\infty} V_{d,l}$$

↓
linear subspace of d^l
polynomials

- $\dim(V_{d,l}) = B_l = \frac{2l+d-2}{d-2} \binom{l+d-3}{l} = \Theta(d^l)$

- orthonormal basis: $\{Y_{lj}\}_{j \in [B_l]}$
spherical harmonics

$$\mathbb{E}[Y_{lj} Y_{k j'}] = \delta_{lk} \delta_{j j'}$$

Integral operator: T commute with $SO(d)$

$$* \sigma(\langle \alpha, \theta \rangle / d) = \sum_{k=0}^{\infty} \zeta_k \sum_{j \in [B_k]} Y_{k_j}(r) Y_{k_j}(\theta)$$

eigenvalues (λ_k) ζ_k degeneracy B_k

$$* \sum_k \zeta_k^2 B_k \leq \mathbb{E}[\sigma^2] < \infty$$

$$\zeta_k^2 = O(d^{-k})$$

$$\rightarrow \sum_k \zeta_k^2 = O(d^{-k})$$

$$\mu_k(\sigma) = \mathbb{E}_G[\sigma(G) M_k(G)] \\ \neq 0$$

Check the assumptions:

$$d^0 \leq m \leq d^{0+1}$$

$$m = \sum_{l \leq 0} B_l$$

$$d^S \leq N \leq d^{S+1}$$

$$N = \sum_{l \leq S} B_l$$

$$u \geq \text{mon}(m, M)$$

$$\sum_k \alpha_k d^{-\frac{k}{2}}$$

① Spaces of low- d^0 pol. on the sphere are hypercontractive

[Beckner, 92] f degree l polynomial

$$\|f\|_{L^q(S^{d-1})}^2 \leq (q-1)^l \|f\|_{L^2(S^{d-1})}^2$$

$$\begin{aligned} \textcircled{2} \quad H_{>m}(\alpha_i, \alpha_i) &= \sum_{k=S+1}^{\infty} \underbrace{\sum_{j \in [B_k]} Y_{kj}(\alpha_i)^2}_{= B_k} \\ &> m \\ &= \sum_{k=S+1}^{\infty} \sum_k B_k \end{aligned}$$

$$U_{>M}(\theta_i, \theta_i) = \sum_{k=S+1}^{\infty} \sum_k B_k$$

$$\textcircled{3} \quad \sum_{k=l+1}^{\infty} \lambda_k^2 = O(1)$$

$$= \Theta(1)$$

6 not polynomial

Underparam. regime:

$$\textcircled{1} \quad \frac{1}{d_M^2} \sum_{k>M} d_k^2 \approx \frac{1}{\sum_S^2} \approx d^S$$

$$\textcircled{2} \quad \frac{1}{d_{M+1}^2} \sum_{k>M} d_k^2 \approx \frac{1}{\sum_{S+1}^2} \approx d^{S+1}$$

$$\textcircled{1} \ll N \ll \textcircled{2}$$

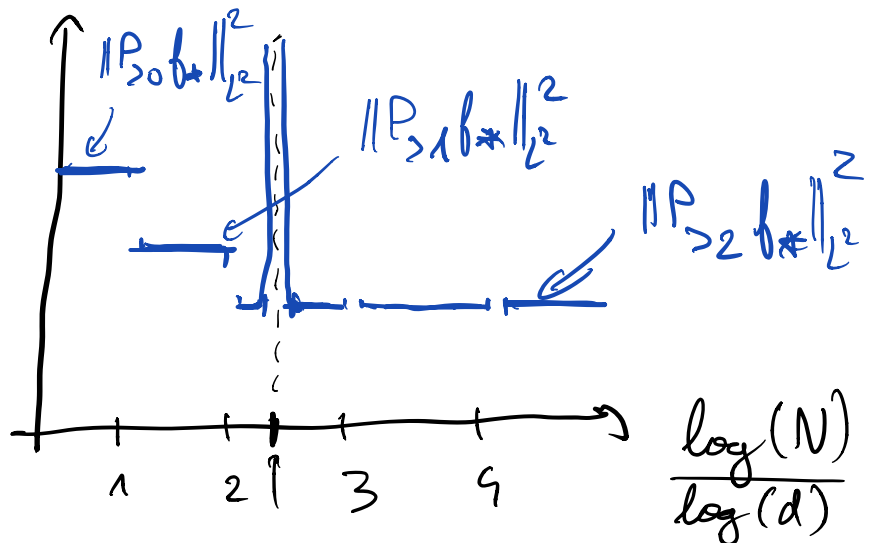
$$d^S \ll N \ll d^{S+1}$$

Overp. regime $d^S \ll n \ll d^{S+1}$



Figure:

$$n = d^{2.4}$$



$$\frac{\log(m)}{\log(d)} = \frac{\log(N)}{\log(d)}$$

$$N = c_1 d$$

$$m = c_2 d$$

$$N = m$$

Double descent
phenomenon