

Sharp convergence rates for spectral methods via the Feature Space Decomposition method

Zong Shang

(joint with Guillaume Lécué, and Zhifan Li)

CREST-ENSAE, Institut Polytechnique de Paris

March 12, 2026

“The profound study of nature is the most fertile source of mathematical discoveries.”

— *Joseph Fourier, The Analytical Theory of Heat, 1822.*

Supervised learning problem and its solution

Benign Overfitting \Rightarrow Feature Space Decomposition method.

- ▶ A tool to analyze the excess risk, improving the classical uniform convergence argument.
- ▶ A theoretical framework.

Supervised Regression Problems: $\mathcal{R} = (\mu_X, f^*, \xi)$, s.t.

$$Y = f^*(X) + \xi, \quad f^* \in L^2(\mu_X), \quad \mathbb{E}[\xi] = 0, \quad \xi \perp\!\!\!\perp X, \quad \text{and } \mathbb{E}[\xi^2] = \sigma_\xi^2.$$

Solutions: $(\mathcal{F}, \{\hat{f}_N\}_{N \in \mathbb{N}_+})$.

- ▶ Feature space: $\mathcal{F} \subset \{f : \Omega_X \rightarrow \mathbb{R}, f \text{ measurable}\}$ (linear).
- ▶ Learning rule: $\{\hat{f}_N : (X_i, Y_i)_{i=1}^N \in (\Omega_X \times \mathbb{R})^N \mapsto \hat{f}_N((X_i, Y_i)_{i=1}^N, \cdot) \in \mathcal{F}\}_{N \in \mathbb{N}_+}$. Abbr. \hat{f} .

Objective: characterize $\|\hat{f} - f^*\|_{L^2(\mu_X)}^2$.

Feature Space Decomposition (FSD)

Definition

Any direct-sum decomposition $\mathcal{F} = V_J \oplus V_{J^c}$ is called a Feature Space Decomposition (FSD).

Equivalently, $I_{\mathcal{F}} = P_{V_J} + P_{V_{J^c}}$, $f = f_J + f_{J^c}$, where $f_J = P_{V_J} f$.

For any FSD, $\hat{f} - f^* = \hat{f}_J - f_J^* + \hat{f}_{J^c} - f_{J^c}^*$. Therefore,

$$\|\hat{f} - f^*\|_{L^2(\mu_X)}^2 \begin{cases} = \|\hat{f}_J - f_J^*\|_{L^2(\mu_X)}^2 + \|\hat{f}_{J^c} - f_{J^c}^*\|_{L^2(\mu_X)}^2, & \text{if } V_J \perp V_{J^c}, \\ \leq 2\|\hat{f}_J - f_J^*\|_{L^2(\mu_X)}^2 + 2\|\hat{f}_{J^c} - f_{J^c}^*\|_{L^2(\mu_X)}^2, & \text{otherwise.} \end{cases}$$

By triangular inequality for $\|\hat{f}_{J^c} - f_{J^c}^*\|_{L^2(\mu_X)}$,

$$\|\hat{f} - f^*\|_{L^2(\mu_X)}^2 \leq 2\|\hat{f}_J - f_J^*\|_{L^2(\mu_X)}^2 + 4\|\hat{f}_{J^c}\|_{L^2(\mu_X)}^2 + 4\|f_{J^c}^*\|_{L^2(\mu_X)}^2. \quad (1)$$

FSD method: seeking rate function r and deviation function δ , s.t., \forall FSD (V_J, V_{J^c}) ,

$$\mathbb{P}((1) \leq r^2(V_J, V_{J^c})) \geq 1 - \delta(V_J, V_{J^c}).$$

Two subspaces, two approaches

There are two fundamentally different ways to control $\|\hat{f} - f^*\|_{L^2(\mu_X)}^2$.

$$\|\hat{f}_J - f_J^*\|_{L^2(\mu_X)}^2 + \|\hat{f}_{J^c}\|_{L^2(\mu_X)}^2 + \|f_{J^c}^*\|_{L^2(\mu_X)}^2.$$

► V_J : estimation subspace.

\hat{f}_J estimates f_J^* , \Rightarrow contribution to excess risk: cancellation between \hat{f}_J and f_J^*

► V_{J^c} : free subspace.

\hat{f}_{J^c} fulfills certain tasks determined by the definition of \hat{f} , but **NOT** estimating $f_{J^c}^*$,
 \Rightarrow contribution to excess risk: smallness of $\|\hat{f}_{J^c}\|_{L^2(\mu_X)}^2$ and $\|f_{J^c}^*\|_{L^2(\mu_X)}^2$.

Optimal FSD

Define

$$(V_J^*, V_{J^c}^*) \in \operatorname{argmin} (r(V_J, V_{J^c}) : \mathcal{F} = V_J \oplus V_{J^c})$$

be the optimal FSD. Then,

$$\mathbb{P} \left(\|\hat{f} - f^*\|_{L^2(\mu_X)}^2 \leq r^2(V_J^*, V_{J^c}^*) \right) \geq 1 - \delta(V_J^*, V_{J^c}^*).$$

For some problems and solutions, e.g., ridge, gradient descent, gradient flow, $\exists 0 < c, \delta < 1$, s.t.,

$$\mathbb{P} \left(\|\hat{f} - f^*\|_{L^2(\mu_X)}^2 \geq cr^2(V_J^*, V_{J^c}^*) \right) \geq 1 - \delta.$$

Remarkable phenomenon: for those classes of (μ_X, f^*, ξ) and (\mathcal{F}, \hat{f}_N) , the estimation error $\|\hat{f} - f^*\|_{L^2(\mu_X)}^2$ is “characterized” by an interpolation between these two distinct approaches.

Improve the Classical Statistical Learning Theory via FSD

$\|\hat{f} - f^*\|_{L^2(\mu_X)}^2 \sim r^2(V_J^*, V_{J^c}^*)$, w.h.p. $\Rightarrow \hat{f}$ estimates **ONLY** $f_{J^*}^*$ (STRIKING!).

► Classical Statistical Learning Theory.

Belief: estimation should happen over the entire \mathcal{F} . $\Rightarrow V_J = \mathcal{F}$.

Belief: If \hat{f} is consistent ($\|\hat{f} - f^*\|_{L^2(\mu_X)} \rightarrow 0$), \hat{f} should estimate f^* , **NOT** only a part.

Learning theory has to address the following four questions:

- (i) *What are (necessary and sufficient) conditions for consistency of a learning process based on the ERM principle?*
- (ii) *How fast is the rate of convergence of the learning process?*
- (iii) *How can one control the rate of convergence (the generalization ability) of the learning process?*
- (iv) *How can one construct algorithms that can control the generalization ability?*

The geometry of FSD

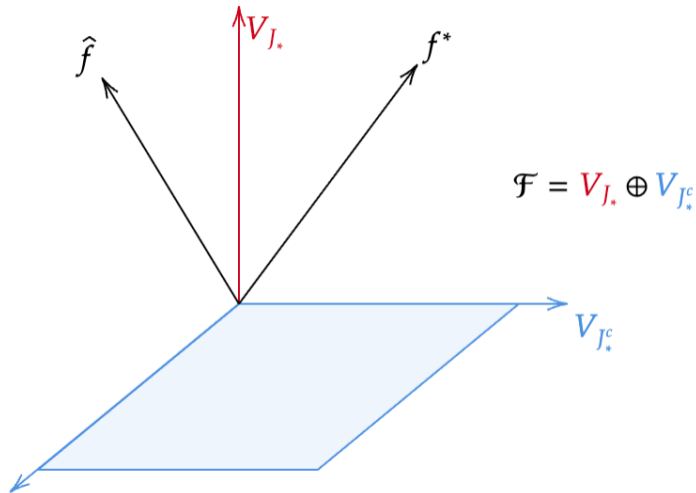


Figure: Only $f_{J_*}^*$ will be estimated by \hat{f}_J , $f_{J_*^c}^*$ will be regarded as the model noise of $f_{J_*}^*$.

\mathcal{R} : linear regression problem

- ▶ μ_X : centered, sub-Gaussian probability measure;
- ▶ $\beta^* \in \mathbb{R}^p$: the unknown signal, s.t., $f^* : \mathbf{x} \in \mathbb{R}^p \mapsto f^*(\mathbf{x}) = \langle \beta^*, \mathbf{x} \rangle \in \mathbb{R}$.
- ▶ $\Sigma = \mathbb{E}[X \otimes X] : \mathbf{v} \in \mathbb{R}^p \mapsto \Sigma \mathbf{v} = \mathbb{E}[X \langle X, \mathbf{v} \rangle] \in \mathbb{R}^p$.
- ▶ $\Sigma = \sum_{j=1}^p \sigma_j \mathbf{e}_j \otimes \mathbf{e}_j$: spectral decomposition, $\sigma_1 \geq \dots \geq \sigma_p > 0$;
- ▶ $\mathcal{F} = \{f(\cdot) : f \in \mathcal{F}\} \longleftrightarrow \{\langle \beta, \cdot \rangle : \beta \in \mathbb{R}^p\} \cong \mathbb{R}^p$;
- ▶ $\|\hat{f} - f^*\|_{L^2(\mu_X)} = \|\Sigma^{1/2}(\hat{\beta} - \beta^*)\|_2$.
- ▶ $\mathcal{R} = (\beta^*, \Sigma, \sigma_\xi)$.

- ▶ $(X_1, Y_1), \dots, (X_N, Y_N)$: i.i.d. copies of (X, Y) .
- ▶ $\mathbb{X} : \mathbf{v} \in \mathbb{R}^p \mapsto \mathbb{X} \mathbf{v} = (\langle X_i, \mathbf{v} \rangle)_{i=1}^N \in \mathbb{R}^N$;
- ▶ $\mathbb{X}^\top : \boldsymbol{\lambda} \in \mathbb{R}^N \mapsto \sum_{i=1}^N \lambda_i X_i \in \mathbb{R}^p$.
- ▶ $\hat{\Sigma} = \frac{1}{N} \mathbb{X}^\top \mathbb{X} = \frac{1}{N} \sum_{i=1}^N X_i \otimes X_i$.

For any FSD, $\Sigma_J = P_J \Sigma P_J$, and $\Sigma_{J^c} = P_{J^c} \Sigma P_{J^c}$, $\mathbb{X}_J = \mathbb{X} P_J$, and $\mathbb{X}_{J^c} = \mathbb{X} P_{J^c}$.

Solution to \mathcal{R} : spectral methods

Definition

Let $(\varphi_t)_{t \geq 1}$ be a family of real-valued functions defined on \mathbb{R}_+ , called the filter functions. For all $t \geq 1$, define $\hat{\beta}$ by

$$\hat{\beta} : \mathbf{y} \in \mathbb{R}^N \mapsto \hat{\beta}(\mathbf{y}) = \frac{1}{N} \varphi_t(\hat{\Sigma}) \mathbb{X}^\top \mathbf{y} = \frac{1}{N} \mathbb{X}^\top \varphi_t \left(\frac{1}{N} \mathbb{X} \mathbb{X}^\top \right) \mathbf{y}.$$

Here, $\varphi_t(\hat{\Sigma})$ is defined by spectral calculus. We abbreviate $\hat{\beta}(\mathbf{y})$ by $\hat{\beta}$. Define $\psi_t(x) = 1 - x\varphi_t(x)$, called residual function.

Examples:

1. Ridge regression with tuning parameter t^{-1} , i.e., $\hat{\beta} = \frac{1}{N} (\frac{1}{N} \mathbb{X}^\top \mathbb{X} + t^{-1} I_p)^{-1} \mathbb{X}^\top \mathbf{y}$ is a spectral method with

$$\varphi_t(x) = \frac{1}{t^{-1} + x}, \text{ and } \psi_t(x) = \frac{1}{xt + 1}.$$

Examples of spectral methods

2. Gradient flow $\dot{\beta}_t = -\nabla \frac{1}{2N} \|\mathbf{y} - \mathbb{X} \cdot \beta_t\|_2^2$ for any $t \geq 1$ with $\beta_1 = \mathbf{0}$. Then $\hat{\beta} = \beta_t$ is a spectral method with

$$\varphi_t(x) = \frac{1 - \exp(-tx)}{x}, \quad \varphi_t(0) = t, \quad \text{and} \quad \psi_t(x) = \exp(-tx).$$

3. Gradient Descent $\beta_t = \beta_{t-1} - \eta \nabla \frac{1}{2N} \|\mathbf{y} - \mathbb{X} \cdot \beta_{t-1}\|_2^2$, $\beta_1 = \mathbf{0}$. Then $\hat{\beta} = \beta_t$ is a spectral method when η is small enough with

$$\varphi_t(x) = \frac{1 - (1 - \eta x)^t}{x}, \quad \text{and} \quad \psi_t(x) = (1 - \eta x)^t.$$

4. Principle Components Regression (PCR) in dimension $d \leq p$,
 $\hat{\beta} \in \operatorname{argmin}(\|\mathbf{y} - \mathbb{X}\beta\|_2^2 : \beta \in \hat{V}_{\leq d})$, where $\hat{V}_{\leq d}$ is the subspace spanned by the first d eigenvectors of $\hat{\Sigma}$. Then

$$\varphi_t(x) = \frac{1}{x} \mathbb{1}(bt^{-1} \leq x), \quad \text{and} \quad \psi_t(x) = \mathbb{1}(bt^{-1} > x).$$

FSD for spectral methods

Definition (estimation dimension)

Let $b > 0$, $t \geq 1$. The estimation dimension of the spectral method $\hat{\beta}$ with filter function φ_t is defined as

$$k^* = k_{t-1,b}^* = \min \{k \in [p] : \sigma_{k+1} \leq bt^{-1}\}.$$

Define the optimal FSD be

$$V_J^* = \text{span}(e_j : 1 \leq j \leq k^*), \text{ and } V_{J^c}^* = \text{span}(e_j : k^* < j \leq p).$$

V_J^* is the subspace used for estimating β_J^* , giving name for “estimation dimension”.

$V_J^* \perp V_{J^c}^*$, hence

$$\|\Sigma^{1/2}(\hat{\beta} - \beta^*)\|_2^2 = \|\Sigma_J^{1/2}(\hat{\beta}_J - \beta_J^*)\|_2^2 + \|\Sigma_{J^c}^{1/2}(\hat{\beta}_{J^c} - \beta_{J^c}^*)\|_2^2.$$

The meaning of filter function

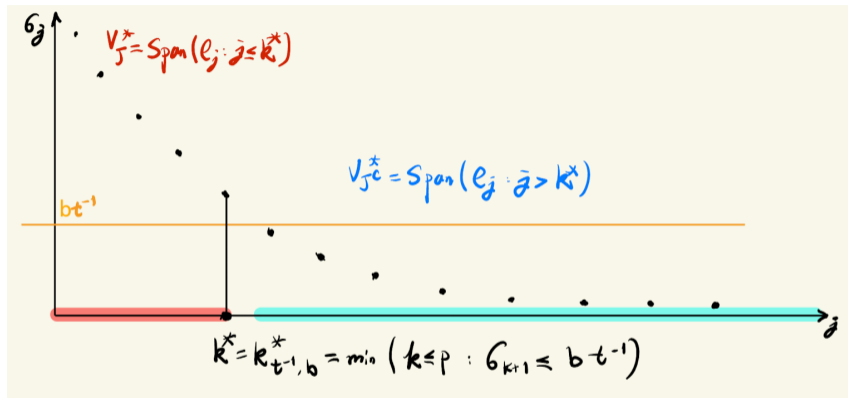


Figure: Given t^{-1} , $\hat{\beta}$ uses $V_J^* = \text{span}(e_j : j \leq k^*)$ to estimate β_J^* , and use $V_{J^c}^*$ to absorb noise. This also explains why φ_t is called a filter function: the feature subspace corresponding to eigenvalues smaller than bt^{-1} is effectively filtered out from being used for estimation.

Definition of rate function $r(\cdot, \cdot)$

Let $J_* = \{1, \dots, k^*\}$. Recall that $\hat{\beta} = \hat{\beta}_J + \hat{\beta}_{J^c}$.

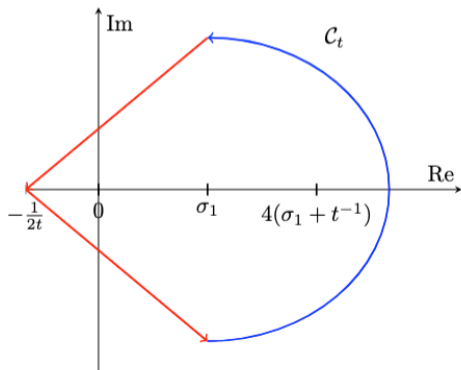
$$r(V_J^*, V_{J^c}^*) = \underbrace{\|\Sigma_{J_*}^{1/2} \psi_t(\Sigma) \beta_{J_*}^*\|_2}_{\hat{\beta}_{J_*} \text{ bias and variance}} + \underbrace{\|\Sigma_{J^c}^{1/2} \beta_{J^c}^*\|_2}_{\hat{\beta}_{J^c} \text{ bias and variance}} \quad \text{bias } \hat{\beta}$$
$$+ \underbrace{\sigma_\xi \sqrt{\frac{|J_*|}{N}}}_{\hat{\beta}_{J_*} \text{ bias and variance}} + \underbrace{\sigma_\xi t \sqrt{\frac{\text{Tr}(\Sigma_{J^c}^2)}{N}}}_{\hat{\beta}_{J^c} \text{ bias and variance}} \quad \text{variance } \hat{\beta}$$

1. As the terminology “residual” suggests, $\|\Sigma_{J_*}^{1/2} \psi_t(\Sigma) \beta_{J_*}^*\|_2$ is the bias.
2. $\|\Sigma_{J^c}^{1/2} \beta_{J^c}^*\|_2$ is the cost incurred by the un-estimated part $\beta_{J^c}^*$ and also the approximation error of V_J^* .

Assumption

Suppose the family of filter functions $(\varphi_t)_{t \geq 1}$ is such that for all $t \geq 1$, φ_t has an holomorphic extension to an open subset of \mathbb{C} containing some contour \mathcal{C}_t . Furthermore, there are two absolute constants $0 \leq c \leq C$ such that for all $t \geq 1$ and $0 \leq x \leq 8$,

$$\frac{c}{x + t^{-1}} \leq \varphi_t(x) \leq \frac{C}{x + t^{-1}}.$$



Fundamental theorem of spectral methods: informal

FSD method serves as a tool to help theorists analyze the excess risk.

Recall that $r(V_J^*, V_{J^c}^*) = \|\Sigma_{J_*}^{1/2} \psi_t(\Sigma) \beta_{J_*}^*\|_2 + \|\Sigma_{J^c}^{1/2} \beta_{J^c}^*\|_2 + \sigma_\xi \sqrt{\frac{|J_*|}{N}} + \sigma_\xi t \sqrt{\frac{\text{Tr}(\Sigma_{J^c}^2)}{N}}$.

Under some assumptions, for any \mathcal{R} characterized by $(\beta^*, \Sigma, \sigma_\xi)$. Let $(\varphi_t)_{t \geq 1}$ be a family of filter functions satisfying the aforementioned assumption.

$$\|\langle X, \hat{\beta} - \beta^* \rangle\|_{L^2(\mu_X)} \lesssim r(V_J^*, V_{J^c}^*).$$

Moreover, if $X = \Sigma^{1/2} Z$ for some Z having independent and centered coordinates, then w.h.p.,

$$\|\langle X, \hat{\beta} - \beta^* \rangle\|_{L^2(\mu_X)} \gtrsim r(V_J^*, V_{J^c}^*).$$

Why do I call it fundamental?

Fundamental theorem is sharp

Recall that $r(V_J^*, V_{J^c}^*) = \|\Sigma_{J_*}^{1/2} \psi_t(\Sigma) \beta_{J_*}^*\|_2 + \|\Sigma_{J_*^c}^{1/2} \beta_{J_*^c}^*\|_2 + \sigma_\xi \sqrt{\frac{|J_*|}{N}} + \sigma_\xi t \sqrt{\frac{\text{Tr}(\Sigma_{J_*^c}^2)}{N}}$.

Under some assumptions, for any \mathcal{R} characterized by $(\beta^*, \Sigma, \sigma_\xi)$. Let $(\varphi_t)_{t \geq 1}$ be a family of filter functions satisfying the aforementioned assumption.

$$\|\langle X, \hat{\beta} - \beta^* \rangle\|_{L^2(\mu_X)} \sim r(V_J^*, V_{J^c}^*).$$

1. For any $\mathcal{R} = (\Sigma, \beta^*, \sigma_\xi)$, $r(V_J^*, V_{J^c}^*)$ characterizes $\|\langle X, \hat{\beta} - \beta^* \rangle\|_{L^2(\mu_X)}$ precisely, with no assumptions imposed on Σ nor on β^*, σ_ξ .
2. Only $\|\Sigma_{J_*}^{1/2} \psi_t(\Sigma) \beta_{J_*}^*\|_2$ depends on the choice of filter function.
3. Intrinsic drawback of spectral methods: J_* is independent with β^* .

Partial order: comparing $\|\langle X, \hat{\beta} - \beta^* \rangle\|_{L^2(\mu_X)}^2 \leftrightarrow$ comparing $r^2(V_J^*, V_{J^c}^*)$

FSD serves also as a theoretical framework.

Definition (Partial Order w.r.t. \mathcal{R} on the set of spectral methods)

For any \mathcal{R} , and any filter functions $(\varphi_t^{(A)})_{t \geq 1}$, $(\varphi_t^{(B)})_{t \geq 1}$, any two tuning parameters t_A, t_B , we write $\hat{\beta}_{t_A}^{(A)} \preceq_{\mathcal{R}} \hat{\beta}_{t_B}^{(B)}$ if $r_{t_A}^{(A)}(V_J^*, V_{J^c}^*) = O(r_{t_B}^{(B)}(V_J^*, V_{J^c}^*))$ as $N, p \rightarrow \infty$.

Corollary

For any \mathcal{R} and any $t \geq 1$, $\hat{\beta}_t^{(A)} \preceq_{\mathcal{R}} \hat{\beta}_t^{(B)}$, if and only if,

$$\|\Sigma_{J_*}^{1/2} \psi_t^{(A)}(\Sigma) \beta_{J_*}^*\|_2 = O(\|\Sigma_{J_*}^{1/2} \psi_t^{(B)}(\Sigma) \beta_{J_*}^*\|_2).$$

Simple but important: $\psi_t^{(\text{GF})}(x) = \exp(-tx)$, and $\psi_t^{(\text{Ridge})}(x) = \frac{1}{xt+1}$. Then

$$\forall \mathcal{R}, \forall t \geq 1, \quad \hat{\beta}_t^{(\text{GF})} \preceq_{\mathcal{R}} \hat{\beta}_t^{(\text{Ridge})}.$$

Is it fair to compare the same t ?

Generalized saturation effect

Definition (Generalized Saturation Effect w.r.t. \mathcal{R})

For any interval $I \subset [1, \infty)$, we write $\{\varphi_t^{(A)}\}_{t \in I} \preceq_{\mathcal{R}} \{\varphi_t^{(B)}\}_{t \in I}$, if as $N, p \rightarrow \infty$,

$$\inf_{t_A \in I} \left(r_{t_A}^{(A)}(V_{J^*}, V_{J^c}^*) \right) = O \left(\inf_{t_B \in I} \left(r_{t_B}^{(B)}(V_{J^*}, V_{J^c}^*) \right) \right).$$

Similarly, $\prec_{\mathcal{R}}$ if big-O is replaced by small-o.

1. Classical saturation effect in Sobolev regression setup.

$\{\varphi_t^{(\text{GF})}\}_{t \in \mathbb{R}_+} \prec_{\mathcal{R}} \{\varphi_t^{(\text{Ridge})}\}_{t \in \mathbb{R}_+}$, for any $\mathcal{R} \in \mathfrak{R}_{\text{Sob}}(s, \alpha)$, $s > 2$. Here, for any $\alpha > 1, s \geq 0$, define

$\mathfrak{R}_{\text{Sob}}(s, \alpha) = \{(\Sigma, \beta^*, \sigma_{\xi}) : \sigma_j \sim j^{-\alpha}, \|\Sigma^{\frac{1-s}{2}} \beta^*\|_2 < \infty, \sigma_{\xi} \text{ is constant}\}$ be the class of Sobolev regression problems.

Example: saturation effect in spiked covariance model

2. $\sigma_1 = \dots = \sigma_k = \sigma$, $\sigma_{k+1} = \dots = \sigma_p = \varepsilon$ for some $\sigma > \varepsilon$ and $k \lesssim N \lesssim p - k$. Let $|\langle \beta^*, e_j \rangle| = \alpha_* \mathbb{1}(j \leq k)$ for some $\alpha_* > 0$. Let $\text{SNR} = \frac{\|\Sigma^{1/2} \beta^*\|_2}{\sigma_\varepsilon} \cdot \frac{\sigma \sqrt{N}}{\sqrt{\text{Tr}(\Sigma_{J^c}^2)}}$. Suppose $4 < \text{SNR} \leq b \frac{\sigma}{\varepsilon}$. Let $I = \{t > 1 : b^{-1} \varepsilon \leq t^{-1} < \sigma\}$. Then $\{\varphi_t^{(\text{GF})}\}_{t \in I} \prec_{\mathcal{R}} \{\varphi_t^{(\text{Ridge})}\}_{t \in I}$ as $\text{SNR}, \sigma/\varepsilon \rightarrow \infty$.

- ▶ The saturation effect is widely present.
- ▶ Developing a theory for an arbitrary \mathcal{R} is important.

Let us go back to the method that produces the fundamental theorem.

How does FSD produce the main result?

$$\|\Sigma^{1/2}(\hat{\beta} - \beta^*)\|_2 \lesssim \|\Sigma_{J_*}^{1/2}(\hat{\beta}_{J_*} - \beta_{J_*}^*)\|_2 + \|\Sigma_{J_*^c}^{1/2} \hat{\beta}_{J_*^c}\|_2 + \|\Sigma_{J_*^c}^{1/2} \beta_{J_*^c}^*\|_2.$$

Upper bound for $\|\Sigma_{J_*}^{1/2}(\hat{\beta}_{J_*} - \beta_{J_*}^*)\|_2$

1. Identify new supervised regression problem.

$\mathbf{y} = \mathbb{X}\beta^* + \xi = \mathbb{X}_{J_*}\beta_{J_*}^* + \mathbb{X}_{J_*^c}\beta_{J_*^c}^* + \xi$, and treat $\mathbb{X}_{J_*^c}\beta_{J_*^c}^* + \xi$ as the noise in a new regression problem in which $\hat{\beta}_{J_*}$ estimates $\beta_{J_*}^*$.

$$\hat{\beta}_{J_*}(\mathbf{y}) - \beta_{J_*}^* = \hat{\beta}_{J_*}(\mathbb{X}_{J_*}\beta_{J_*}^*) - \beta_{J_*}^* + \hat{\beta}_{J_*}(\mathbb{X}_{J_*^c}\beta_{J_*^c}^* + \xi)$$

2. Apply classical statistical analysis to $\hat{\beta}_{J_*}$.

Recall that $\hat{\beta}_{J_*}(\mathbb{X}_{J_*}\beta_{J_*}^*) = \varphi_t(\hat{\Sigma})\hat{\Sigma}\beta_{J_*}^*$. Let $\tilde{\beta}_{J_*} = \varphi_t(\Sigma)\Sigma\beta_{J_*}^*$. Then

$$\begin{aligned} \left\| \Sigma_{J_*}^{1/2}(\hat{\beta}_{J_*}(\mathbf{y}) - \beta_{J_*}^*) \right\|_2 &\leq \left\| \Sigma_{J_*}^{1/2}(\tilde{\beta}_{J_*} - \beta_{J_*}^*) \right\|_2 \\ &+ \left\| \Sigma_{J_*}^{1/2}(\hat{\beta}_{J_*}(\mathbb{X}_{J_*}\beta_{J_*}^*) - \tilde{\beta}_{J_*}) \right\|_2 + \left\| \Sigma_{J_*}^{1/2}\hat{\beta}_{J_*}(\mathbb{X}_{J_*^c}\beta_{J_*^c}^* + \xi) \right\|_2. \end{aligned}$$

Classical analysis for spectral methods applied to $\hat{\beta}_J$

1. By definition, $\tilde{\beta}_{J_*} = \Sigma \varphi_t(\Sigma) \beta_{J_*}^*$, so

$$\|\Sigma_{J_*}^{1/2}(\beta_{J_*}^* - \tilde{\beta}_{J_*})\|_2 = \|\Sigma_{J_*}^{1/2}(I - \Sigma \varphi_t(\Sigma))\beta_{J_*}^*\|_2 = \|\Sigma_{J_*}^{1/2}\psi_t(\Sigma)\beta_{J_*}^*\|_2.$$

2. Using the method in [Li, Gan, Shi, and Lin \[2024\]](#), w.h.p.,

$$\|\Sigma_{J_*}^{1/2}(\hat{\beta}_{J_*}(\mathbb{X}_{J_*}\beta_{J_*}^*) - \tilde{\beta}_{J_*})\|_2 \lesssim \|\Sigma_{J_*}^{1/2}\psi_t(\Sigma)\beta_{J_*}^*\|_2 + \text{higher order small.}$$

3. One may also prove that the variance part w.h.p. satisfies

$$\|\Sigma_{J_*}^{1/2}\hat{\beta}_{J_*}(\mathbb{X}_{J_*^c}\beta_{J_*^c}^* + \xi)\|_2 \lesssim \|\Sigma_{J_*^c}^{1/2}\beta_{J_*^c}^*\|_2 + \sigma_\xi \sqrt{\frac{|J_*|}{N}}.$$

Combining, without additional effort, w.h.p.,

$$\left\| \Sigma_{J_*}^{1/2}(\hat{\beta}_{J_*} - \beta_{J_*}^*) \right\|_2 \lesssim \|\Sigma_{J_*}^{1/2}\psi_t(\Sigma)\beta_{J_*}^*\|_2 + \|\Sigma_{J_*^c}^{1/2}\beta_{J_*^c}^*\|_2 + \sigma_\xi \sqrt{\frac{|J_*|}{N}}.$$

In spectral methods, FSD gives “correct” oracle $\beta_{J_*}^*$ and bias $\|\Sigma_{J_*^c}^{1/2}\beta_{J_*^c}^*\|_2$.

In this work, FSD serves as an outer analytical layer built on top of the classical analysis of spectral methods.

Upper bound for $\|\Sigma_{J^c}^{1/2} \hat{\beta}_{J^c}\|_2$

This part remains a mystery and requires creative insight!

$$\begin{aligned}\|\Sigma_{J^c}^{1/2} \hat{\beta}_{J^c}(\mathbf{y})\|_2 &\leq \|\Sigma_{J^c}^{1/2} \varphi_t(\hat{\Sigma}) \hat{\Sigma} \beta_J^*\|_2 \\ &\quad + \|\Sigma_{J^c}^{1/2} \varphi_t(\hat{\Sigma}) \hat{\Sigma} \beta_{J^c}^*\|_2 + \|\Sigma_{J^c}^{1/2} \varphi_t(\hat{\Sigma}) \frac{1}{N} \mathbb{X}^\top \boldsymbol{\xi}\|_2.\end{aligned}$$

By $\varphi_t(\hat{\Sigma}) \hat{\Sigma} = I_p - \psi_t(\hat{\Sigma})$, and $\Sigma_{J^c}^{1/2} \beta_J^* = \mathbf{0}$,

$$\|\Sigma_{J^c}^{1/2} \varphi_t(\hat{\Sigma}) \hat{\Sigma} \beta_J^*\|_2 = \|\Sigma_{J^c}^{1/2} (\psi_t(\hat{\Sigma}) - \psi_t(\Sigma)) \beta_J^*\|_2.$$

Inserting $(\Sigma + t^{-1})(\Sigma + t^{-1})^{-1}$ to

$$\psi_t(\hat{\Sigma}) - \psi_t(\Sigma) = \frac{1}{2\pi i} \oint_{\mathcal{C}_t} (\hat{\Sigma} - zI_p)^{-1} (\hat{\Sigma} - \Sigma) (\Sigma - zI_p)^{-1} \psi_t(z) dz$$

yields $\frac{1}{t} \|\Sigma_J^{-1/2} \beta_J^*\|_2$. The other two terms are not difficult.

Going beyond Spectral Methods: Feature Learning Property

This part is based on joint work with G. Lecué, T. Suzuki, and T. Wakayama (on-going).

How the good features are automatically learned by a neural network?

- ▶ Given (μ_X, f^*, σ_ξ) , $(X_i, Y_i)_{i=1}^N$, and an estimator \hat{f}_N .
 - ▶ There exist a data-dependent RKHS $\mathcal{H}_{\text{fea}} \subset L^2(\mu_X)$, called learned feature subspace, a latent estimator $\hat{g}_N \in \mathcal{H}_{\text{fea}}$ and the oracle $g_{\mathcal{H}_{\text{fea}}}^* \in \operatorname{argmin}(\|f^* - g\|_{L^2(\mu_X)} : g \in \mathcal{H}_{\text{fea}})$, such that the following hold.
 - ▶ $\|f^* - g_{\mathcal{H}_{\text{fea}}}^*\|_{L^2(\mu_X)}$ is small (approximation error);
 - ▶ $\|g_{\mathcal{H}_{\text{fea}}}^* - \hat{g}_N\|_{L^2(\mu_X)}$ is small (estimation error);
 - ▶ $\|\hat{g}_N - \hat{f}_N\|_{L^2(\mu_X)}$ is small, (\hat{g}_N can explain \hat{f}_N)
- and $\|\hat{g}_N - g_{\mathcal{H}_{\text{fea}}}^*\|_{L^2(\mu_X)}^2$ decreases when $g_{\mathcal{H}_{\text{fea}}}^*$ gets aligned with the top k eigenvectors of $\Sigma = \mathbb{E}[\phi_{\text{fea}}(X) \otimes \phi_{\text{fea}}(X) | (X_i, Y_i)_{i=1}^N]$ for some $k = o(N)$.

For almost any supervised regression problem, mean-field Langevin dynamics trained shallow neural network has feature learning property.

Take-home messages

- ▶ FSD as analytical tool: how features are utilized.
 - ▶ Estimation occurs only on V_J .
 - ▶ For a large class of spectral methods and almost any $\mathcal{R} = (\beta^*, \Sigma, \sigma_\xi)$, w.h.p., $\|\langle X, \hat{\beta} - \beta^* \rangle\|_{L^2(\mu_X)} \sim r(V_J^*, V_{J^c}^*)$, where

$$r(V_J^*, V_{J^c}^*) = \|\Sigma_{J_*}^{1/2} \psi_t(\Sigma) \beta_{J_*}^*\|_2 + \|\Sigma_{J_*^c}^{1/2} \beta_{J_*^c}^*\|_2 + \sigma_\xi \sqrt{\frac{|J_*|}{N}} + \sigma_\xi t \sqrt{\frac{\text{Tr}(\Sigma_{J_*^c}^2)}{N}}.$$

- ▶ For the same \mathcal{R} , all spectral methods use the same features for estimation.
 - ▶ For the same \mathcal{R} and the same t , the estimation error of different spectral methods depends on the convergence rate of ψ_t .
 - ▶ The saturation effect is widespread and not limited to Sobolev regression.
- ▶ FSD as theoretical framework: how features are constructed (feature learning).
 - ▶ The learned feature subspace \mathcal{H}_{fea} approaches f^* ;
 - ▶ The oracle in \mathcal{H}_{fea} is “aligned” with top k eigenfunctions of the integral operator of \mathcal{H}_{fea} .
 - ▶ Mean-field Langevin dynamics has feature learning property.

End

Thank you for your attention!

FSD in classification

For any FSD and any $f_J^* \in V_J$,

$$P\mathcal{L}_{\hat{f}}^{(0,1)} = \mathbb{P}\left(Y\hat{f}(X) < 0 \mid (X_i, Y_i)_{i=1}^N\right) - \mathbb{P}\left(Y\hat{f}_J(X) < 0 \mid (X_i, Y_i)_{i=1}^N\right) \quad (2)$$

$$+ \mathbb{P}\left(Y\hat{f}_J(X) < 0 \mid (X_i, Y_i)_{i=1}^N\right) - \mathbb{P}\left(Yf_J^*(X) < 0\right) \quad (3)$$

$$+ \mathbb{P}\left(Yf_J^*(X) < 0\right) - \mathbb{P}\left(Y\left(\eta(X) - \frac{1}{2}\right) < 0\right), \quad (4)$$

- ▶ (2): error caused by free part \hat{f}_{J^c} ;
- ▶ (3): prediction error caused by \hat{f}_J compared to that of f_J^* ; and
- ▶ (4): approximation error compared with that of the Bayes rule.

Counterparts of $\|\hat{f}_{J^c}\|_{L^2(\mu_X)}^2$, $\|\hat{f}_J - f_J^*\|_{L^2(\mu_X)}^2$, and $\|f_{J^c}^*\|_{L^2(\mu_X)}^2$ in regression.

NTK does not have feature learning property

Example. Consider the NTK parameterization of the shallow neural network

$f_{\theta}(\cdot) = \frac{1}{\sqrt{M}} \sum_{j=1}^M a_j \sigma(\langle \mathbf{w}_j, \cdot \rangle)$ where $\theta = (\mathbf{a}, \mathbb{W})$. Define

$\hat{f}_N \in \operatorname{argmin}(P_N \ell_{f_{\theta}}^{(2)} : \theta \in \mathbb{R}^M \times \mathbb{R}^{M \times d})$, where $P_N \ell_{f_{\theta}}^{(2)} = \frac{1}{N} \sum_{i=1}^N (Y_i - f_{\theta}(X_i))^2$, and

suppose \hat{f}_N is computed via gradient flow $(f_t)_{t=0}^{\infty}$ initialized at

$f_0(\cdot) = \frac{1}{\sqrt{M}} \sum_{j=1}^M a_j^{(0)} \sigma(\langle \mathbf{w}_j^{(0)}, \cdot \rangle)$ where $a_j^{(0)} \sim \mathcal{N}(0, 1)$ and $\mathbf{w}_j^{(0)} \sim \mathcal{N}(\mathbf{0}, I_d)$ for each

$1 \leq j \leq M$. Then when $M \rightarrow \infty$, f_t converges in probability to the solution of the kernel gradient flow

$$\partial_t f_t(\cdot) = -\frac{2}{N} \sum_{i=1}^N K_{\text{NTK}}(\cdot, X_i) (f_t(X_i) - Y_i),$$

where $K_{\text{NTK}} : (\mathbf{x}_1, \mathbf{x}_2) \in \Omega_X \times \Omega_X \mapsto \mathbb{E}[\langle \nabla_{\theta} f_0(\mathbf{x}_1), \nabla_{\theta} f_0(\mathbf{x}_2) \rangle]$. As a result, $\hat{f}_N \in \mathcal{H}_{\text{NTK}}$, where \mathcal{H}_{fea} is the RKHS generated by K_{NTK} and is independent with μ and with $(X_i, Y_i)_{i=1}^N$, hence is **not** problem-specific. Therefore, NTK parameterization shallow neural network does not have feature learning property.

Two stages trained NN has feature learning property

Example. Train NN by gradient flow $(f_t)_{t=0}^T$ in rich regime stopped at certain time $T > 0$. Denote $\mathbf{w}_1^{(T)}, \dots, \mathbf{w}_M^{(T)}$ be the learned features of f_T .

- ▶ Define $\phi_{\text{fea}} : \mathbf{x} \in \Omega_X \mapsto (\sigma(\langle \mathbf{w}_j^{(T)}, \mathbf{x} \rangle))_{j=1}^M$; let \mathcal{H}_{fea} be the RKHS generated by K_{fea} , called learned feature kernel.
- ▶ Let \hat{f}_N be a ridge regression with some tuning parameter λ on top of the RKHS generated by $K_{\text{fea}} : (\mathbf{x}_1, \mathbf{x}_2) \in \Omega_X \times \Omega_X \mapsto K_{\text{fea}}(\mathbf{x}_1, \mathbf{x}_2) = \langle \phi_{\text{fea}}(\mathbf{x}_1), \phi_{\text{fea}}(\mathbf{x}_2) \rangle$.
- ▶ $\hat{f}_N \in \mathcal{H}_{\text{fea}}$, and \hat{g}_N is ridge on \mathcal{H}_{fea} , hence alignment property.
- ▶ \mathcal{H}_{fea} is problem-specific. For single- and multi-index regression problems, it can be shown that \mathcal{H}_{fea} has a small approximation error.
- ▶ The \hat{f}_N obtained in this way is called a two-stage trained shallow neural network, \Rightarrow feature learning property for solving those supervised regression problems.

Feature Learning reduces “approximation error”

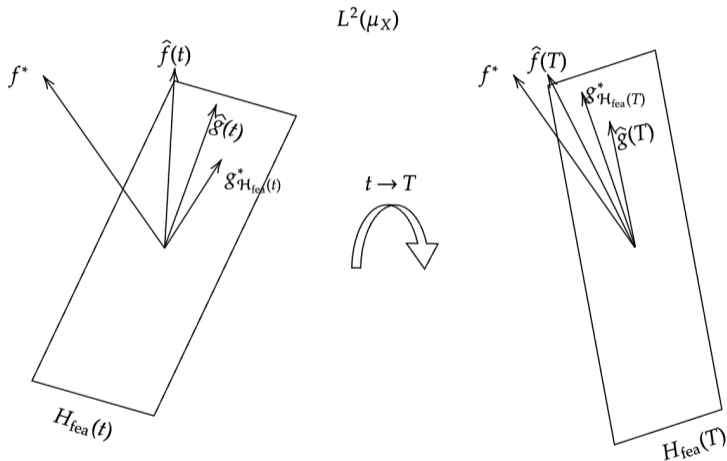


Figure: \mathcal{H}_{fea} is getting closer to f^* during training: $\|f^* - g_{\mathcal{H}_{\text{fea}}}^*\|_{L^2(\mu_X)}$ is getting small.

Feature Learning reduces “estimation error”

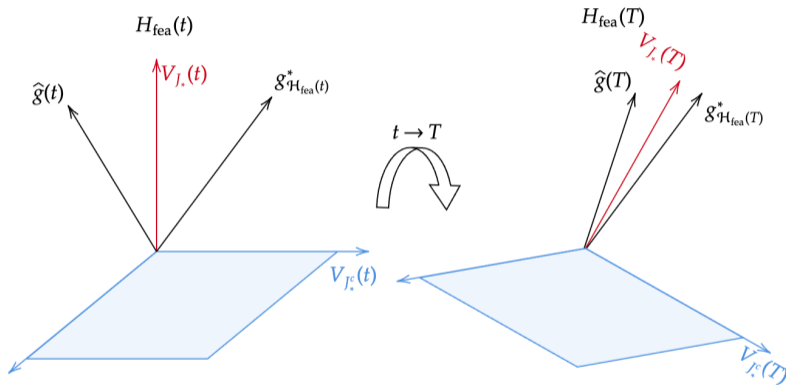


Figure: $g_{\mathcal{H}_{\text{fea}}^*}$ is getting aligned with the top k eigenvectors.

Mean-field Langevin dynamics has feature learning property

- ▶ Let $\Theta \subset \mathbb{R}^{d+1}$ compact, and $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ some continuous function. For any $\theta = (a, \mathbf{w}) \in \Theta$, $\phi(\cdot, \theta) = a\sigma(\langle \mathbf{w}, \cdot \rangle)$, $\mathcal{F} = \{f_\nu(\cdot) = \int_\Theta \phi(\cdot, \theta) d\nu(\theta) : \nu \in \mathcal{P}_{\text{ac}}(\Theta)\}$.
- ▶ Let $\text{Ent}^-(\nu) = \int_\Theta \log(\frac{d\nu}{d\text{Leb}}) d\nu$. For any $\lambda \geq 0$,

$$\hat{\nu}_\lambda \in \underset{\nu \in \mathcal{P}_{\text{ac}}(\Theta)}{\text{argmin}} \left(\frac{1}{N} \sum_{i=1}^N (Y_i - f_\nu(X_i))^2 + \lambda \text{Ent}^-(\nu) \right), \text{ and } \hat{f}_N = f_{\hat{\nu}_\lambda}$$

is called mean-field Langevin dynamics.

- ▶ Let $\hat{\varphi} = (P_W)_\# \hat{\nu}_\lambda$, where $P_W : \theta = (a, \mathbf{w}) \mapsto \mathbf{w}$. Let $\hat{a} : \mathbf{w} \mapsto \mathbb{E}_{\hat{\nu}_\lambda}[a|\mathbf{w}] = \int a d\hat{\nu}_\lambda(a|\mathbf{w})$.
- ▶ Let $\phi_{\text{fea}} : \mathbf{x} \mapsto \sigma(\langle \mathbf{x}, \cdot \rangle)$, $K_{\text{fea}} : (\mathbf{x}_1, \mathbf{x}_2) \mapsto \langle \phi_{\text{fea}}(\mathbf{x}_1), \phi_{\text{fea}}(\mathbf{x}_2) \rangle_{L^2(\hat{\varphi})}$, and

$$\mathcal{H}_{\text{fea}} = \{f_g(\cdot) = \langle \phi_{\text{fea}}(\cdot), g \rangle_{L^2(\hat{\varphi})} : g \in L^2(\hat{\varphi})\}.$$

- ▶ One can prove $\hat{g}_N : h \in \mathcal{H}_{\text{fea}} \mapsto \langle \hat{a}, h \rangle_{L^2(\hat{\varphi})}$ satisfies $\hat{f}_N(\cdot) = \hat{g}_N(\phi_{\text{fea}}(\cdot))$, and \hat{g}_N is a convex regularized M-estimator.
- ▶ For “almost any” supervised regression problem, \hat{f}_N has feature learning property.

Open problem: is PCR the minimal element?

Given any \mathcal{R} , what is the minimal element within $\preceq_{\mathcal{R}}$ (and $\preceq_{\mathcal{R}}$ for $I = [1, \infty)$ respectively)?

PCR: $\psi_t(x) = \mathbb{1}(bt^{-1} > x)$ and $k^* = \min\{k \in [p] : \sigma_{k+1} \leq bt^{-1}\}$, hence

$\|\Sigma_{J_*}^{1/2} \psi_t^{(\text{PCR})}(\Sigma) \beta_{J_*}^*\|_2 = 0$. However, PCR's filter function $\varphi_t(x) = \frac{1}{x} \mathbb{1}(bt^{-1} \leq x)$ cannot be holomorphic extended over \mathcal{C}_t .

Conjecture: PCR is the minimal element of the partial order $\preceq_{\mathcal{R}}$ for any \mathcal{R} .

References I

Yicheng Li, Weiye Gan, Zuoqiang Shi, and Qian Lin. Generalization Error Curves for Analytic Spectral Algorithms under Power-law Decay, July 2024. URL <http://arxiv.org/abs/2401.01599>. arXiv:2401.01599.