

Benign Overfitting Property of the Minimum ℓ_q -norm interpolant estimator in linear regression problem via Feature Space Decomposition

Zong Shang

(joint with Radosław Adamczak, Guillaume Lecué, and Marta Strzelecka)

CREST-ENSAE, Institut Polytechnique de Paris

March 24, 2026

“The unknown thing to be known appeared to me as some stretch of earth or hard marl, resisting penetration. . . the sea advances insensibly in silence, nothing seems to happen, nothing moves, the water is so far off you hardly hear it. . . yet it finally surrounds the resistant substance.”

— Alexandre Grothendieck, *Récoltes et Semailles*(1985), pp. 552–553.

Supervised learning problem and its solution

Feature Space Decomposition method.

- ▶ A tool to analyze the excess risk, improving the classical uniform convergence argument.
- ▶ A theoretical framework (not in this talk).

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}) ,

$$\mathbb{P}\left(\text{R.H.S. of (1)} \leq r^2(V_J, V_{J^c})\right) \geq 1 - \delta(V_J, V_{J^c}).$$

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, see [LLS25, arXiv: 2512.14473], $\exists 0 < c, \delta < 1$, s.t.,

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

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

► Classical Statistical Learning Theory.

Belief: An estimator should use all the features to estimate (no waste).

Belief: Estimation should happen over the entire \mathcal{F} ($V_J = \mathcal{F}$).

Belief: If \hat{f} is consistent ($\|\hat{f} - f^*\|_{L^2(\mu_X)} \rightarrow 0$), \hat{f} should estimate f^* .

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?*

Rethink consistency: Benign Overfitting

Benign Overfitting: motivation for rethinking consistency.

Fix canonical basis $\{e_1, \dots, e_p\}$ of \mathbb{R}^p . Let $\beta^* \in \mathbb{R}^p$, s.t.,

$$Y = \langle \beta^*, X \rangle + \xi.$$

The minimum ℓ_q -norm interpolant estimator is defined as

$$\hat{\beta} = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \left(\|\beta\|_q : \mathbb{X}\beta = \mathbf{y} \right),$$

where $\mathbb{X} = [X_1 | \dots | X_N]^\top \in \mathbb{R}^{N \times p}$, $\mathbf{y} = (Y_1, \dots, Y_N)$.

Belief: $\hat{\beta}$ should not be consistent, because of overfitting/interpolation.

Definition. We say $\hat{\beta}$ exhibits benign overfitting (B.O.) under some limit, if

$$\|\langle X, \hat{\beta} - \beta^* \rangle\|_{L^2(\mu_X)} \rightarrow 0.$$

Overfitting estimator and Benign Overfitting estimator

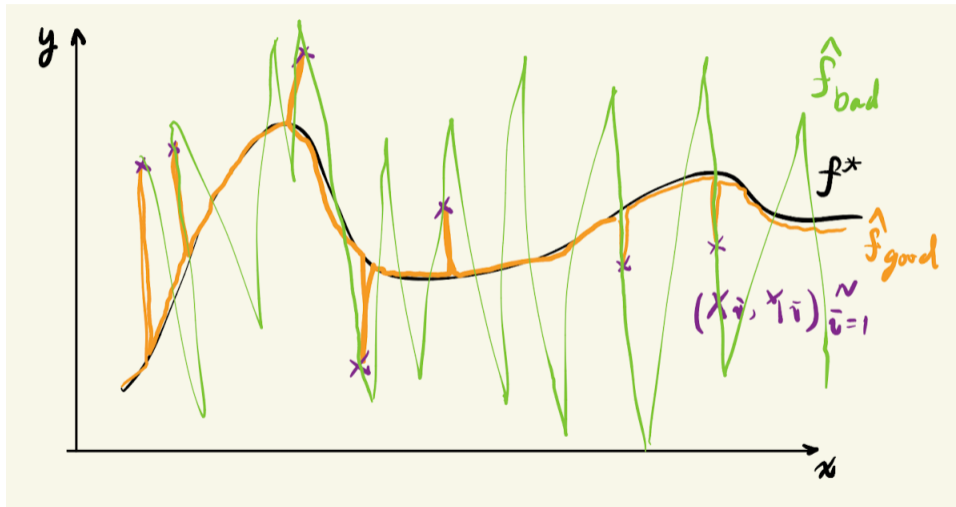


Figure: Overfitting estimator (green curve) and Benign Overfitting estimator (orange)

Classical Stats. Learning Theory fails to penetrate

Notation: $P_N \ell_f$ empirical risk; $P \ell_f = \mathbb{E}[P_N \ell_f]$ population risk; $P_N \mathcal{L}_f = P_N \ell_f - P_N \ell_{f^*}$ empirical excess risk; and $P \mathcal{L}_f = P \ell_f - P \ell_{f^*}$ population excess risk.

1. Localization: Take any deterministic $\mathcal{G} \subset \mathcal{F}$, s.t. $\hat{f} \in \mathcal{G}$, w.h.p..
2. Uniform Convergence argument.

$$P \mathcal{L}_{\hat{f}} = P_N \mathcal{L}_{\hat{f}} + (P - P_N) \mathcal{L}_{\hat{f}} \leq P_N \mathcal{L}_{\hat{f}} + \sup((P - P_N) \mathcal{L}_f : f \in \mathcal{G}),$$

Take $P_N \ell_{\beta} = \frac{1}{N} \sum_{i=1}^N (Y_i - \langle \beta, X_i \rangle)^2$. Then $P \mathcal{L}_{\beta} = \|\langle \beta - \beta^*, X \rangle\|_{L^2(\mu_X)}^2$.

- ▶ Def: $P_N \mathcal{L}_{\hat{\beta}} = P_N \ell_{\hat{\beta}} - P_N \ell_{\beta^*} = -(1 + o(1)) \sigma_{\xi}^2$, because $\mathbf{y} = \mathbb{X} \hat{\beta} = \mathbb{X} \beta^* + \xi$.
- ▶ We want: $P \mathcal{L}_{\hat{\beta}} = o(1) \sigma_{\xi}^2$.
- ▶ We need: $\sup((P - P_N) \mathcal{L}_{\beta} : \beta \in \mathcal{G}) = (1 + o(1)) \sigma_{\xi}^2$ to use U.C.

B.O.: a resistant substance that classical uniform convergence fails to penetrate.

This talk: non-trivial upper bound for $\|\langle X, \hat{\beta} - \beta^* \rangle\|_{L^2(\mu_X)}$ via FSD \Rightarrow suff. cond. for B.O.

- J defines new estimator: $\hat{\beta}_J$ is a RERM.

Let $\Sigma = \mathbb{E}[X \otimes X]$. Then $\|\langle \hat{\beta} - \beta^*, X \rangle\|_{L^2(\mu_X)} = \|\Sigma^{1/2}(\hat{\beta} - \beta^*)\|_2$. By FSD,

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

We first study $\|\Sigma^{1/2}(\hat{\beta}_J - \beta_J^*)\|_2$.

Key: $\hat{\beta}_J$ is RERM.

Proposition For any $J \subset [p]$, let $V_J = \text{span}(e_j : j \in J)$, and $\mathbb{X}_J = \mathbb{X}P_J$, $\mathbb{X}_{J^c} = \mathbb{X}P_{J^c}$. Define the non-linear metric embedding operator

$$\mathcal{A} : \mu \in \mathbb{R}^N \mapsto \mathcal{A}[\mu] \in \operatorname{argmin}(\|\nu\|_q : \mathbb{X}_{J^c}\nu = \mu), \text{ then} \quad (2)$$

$$\hat{\beta}_J \in \operatorname{argmin} (P_N \ell_{\beta_J} + \|\beta_J\|_q^q), \text{ where } P_N \ell_{\beta_J} = \|\mathcal{A}[\mathbf{y} - \mathbb{X}_J \beta_J]\|_q^q$$

is a random loss determined by \mathbb{X}_{J^c} . This is called **self-regularization**. Moreover,

$$\hat{\beta}_{J^c} = \mathcal{A}[\mathbf{y} - \mathbb{X}_J \hat{\beta}_J].$$

• J defines new f_J^*

β_J^* is consistent with the random loss function $P_N \ell_{\bullet}$.

Proposition (β_J^* is oracle) Suppose $\mathbb{X}_{J^c} \beta_{J^c}^* + \xi$ is independent with \mathbb{X}_J , and $\mathbb{E}[\mathbb{X}_J] = \mathbf{0}$, then

$$\beta_J^* \in \operatorname{argmin}(P \ell_{\beta_J} : \beta_J \in V_J), \text{ where } P \ell_{\beta_J} = \mathbb{E}_{\mathbb{X}_J, \xi}[P_N \ell_{\beta_J}], \quad (3)$$

and $P_N \ell_{\beta_J} = \|\mathcal{A}[\mathbf{y} - \mathbb{X}_J \beta_J]\|_q^q$.

Implication: $\|\Sigma^{1/2}(\hat{\beta}_J - \beta_J^*)\|_2^2$ is the estimation error of RERM

$$\hat{\beta}_J \in \operatorname{argmin}(P_N \ell_{\beta_J} + \|\beta_J\|_q^q), \text{ where } P_N \ell_{\beta_J} = \|\mathcal{A}[\mathbf{y} - \mathbb{X} \beta_J]\|_q^q.$$

We need to understand this **random loss function**. This seems to be desperate!

• J_c : new tools from Asymptotic Geometric Analysis

Theorem (Dvoretzky-Milman) There exist $\kappa_{DM}, c_1 \leq 1$ s.t. the following holds. Let $\|\cdot\|$ be some norm on \mathbb{R}^p and denote by B° the unit ball of the dual norm of $\|\cdot\|$. Denote by \mathbb{G} , the $N \times p$ standard Gaussian matrix. Given any $0 < \varepsilon_1 \leq 1$. Assume that $N \leq \kappa_{DM} \varepsilon_1^2 d_*(B)$. Then w.h.p.,

$$\forall \boldsymbol{\lambda} \in \mathbb{R}^N, \quad (1 - \varepsilon_1) \|\boldsymbol{\lambda}\|_2 \ell_*(B^\circ) \leq \|\mathbb{G}^\top \boldsymbol{\lambda}\| \leq (1 + \varepsilon_1) \|\boldsymbol{\lambda}\|_2 \ell_*(B^\circ), \quad \text{where} \quad (4)$$

$\ell_*(B^\circ) = \mathbb{E}\|G\|$, $\text{diam}(B^\circ) = \max(\|\mathbf{v}\|_2 : \mathbf{v} \in B^\circ)$, and $d_*(B) = \left(\frac{\ell_*(B^\circ)}{\text{diam}(B^\circ, \ell_2)}\right)^2$.
Let $\Sigma_{J_c} = P_{J_c} \Sigma P_{J_c}$. When $X \sim \mathcal{N}(\mathbf{0}, \Sigma)$ is Gaussian,

$$\mathbb{X}_{J_c}^\top = [\Sigma_{J_c}^{1/2} G_1 | \cdots | \Sigma_{J_c}^{1/2} G_N] = \Sigma_{J_c}^{1/2} \mathbb{G}^\top.$$

By DM applied to $\|\cdot\| = \|\Sigma_{J_c}^{1/2} \cdot\|_{q'}$, when $N \leq \kappa_{DM} \varepsilon_1^2 d_*(\Sigma_{J_c}^{-1/2} B_{q'}^p)$, w.h.p.,

$$\Omega_{DM, \text{reg}}(\varepsilon_1) = \left\{ \forall \boldsymbol{\lambda} \in \mathbb{R}^N : (1 - \varepsilon_1) \ell_*(\Sigma_{J_c}^{1/2} B_q^p) \|\boldsymbol{\lambda}\|_2 \leq \|\mathbb{X}_{J_c}^\top \boldsymbol{\lambda}\|_{q'} \leq (1 + \varepsilon_1) \ell_*(\Sigma_{J_c}^{1/2} B_q^p) \|\boldsymbol{\lambda}\|_2 \right\}.$$

How does it help understanding $P_N \ell_\bullet$?

An intuition from Dvoretzky-Milman theorem

By duality, $\|\mathcal{A}[\boldsymbol{\mu}]\|_q = \min(\|\boldsymbol{\nu}\|_q : \mathbb{X}_{J^c}\boldsymbol{\nu} = \boldsymbol{\mu}) = \max(\langle \boldsymbol{\mu}, \boldsymbol{\lambda} \rangle : \|\mathbb{X}_{J^c}^\top \boldsymbol{\lambda}\|_{q'} \leq 1)$. Therefore

$$\begin{aligned} \Omega_{\text{DM,reg}}(\varepsilon_1) &= \left\{ \forall \boldsymbol{\lambda} \in \mathbb{R}^N : (1 - \varepsilon_1)l_*(\Sigma_{J^c}^{1/2} B_q^p) \|\boldsymbol{\lambda}\|_2 \leq \|\mathbb{X}_{J^c}^\top \boldsymbol{\lambda}\|_{q'} \leq (1 + \varepsilon_1)l_*(\Sigma_{J^c}^{1/2} B_q^p) \|\boldsymbol{\lambda}\|_2 \right\} \\ &\subset \left\{ \forall \boldsymbol{\mu} \in \mathbb{R}^N : \frac{\|\boldsymbol{\mu}\|_2}{(1 + \varepsilon_1)l_*(\Sigma_{J^c}^{1/2} B_q^p)} \leq \|\mathcal{A}[\boldsymbol{\mu}]\|_q \leq \frac{\|\boldsymbol{\mu}\|_2}{(1 - \varepsilon_1)l_*(\Sigma_{J^c}^{1/2} B_q^p)} \right\}. \end{aligned}$$

Implication: Denote l_* by $l_*(\Sigma_{J^c}^{1/2} B_q)$. From $\|\mathcal{A}[\boldsymbol{\mu}]\|_q \sim l_*^{-1} \|\boldsymbol{\mu}\|_2$, and

$$P_N \ell_{\boldsymbol{\beta}_J} = \|\mathcal{A}[\mathbf{y} - \mathbb{X}_J \boldsymbol{\beta}_J]\|_q^q,$$

$$\hat{\boldsymbol{\beta}}_J \in \operatorname{argmin} (P_N \ell_{\boldsymbol{\beta}_J} + \|\boldsymbol{\beta}_J\|_q^q) \approx \operatorname{argmin} (\|\mathbf{y} - \mathbb{X}_J \boldsymbol{\beta}_J\|_2^q + \ell_*^q \|\boldsymbol{\beta}_J\|_q^q).$$

We have to extend Dvoretzky-Milman theorem.

New result in AGA: Generalized Dvoretzky-Milman theorem

- ▶ Y. Gordon, A. Litvak, A. Pajor and N. Tomczak-Jaegermann. "Random epsilon-nets and embeddings in ℓ_N^∞ ", Studia Mathematica, 2007.
- ▶ S. Mendelson, N. Tomczak-Jaegermann. "A subgaussian embedding theorem", Israel Journal of Mathematics, 2008.
- ▶ D. Bartl, S. Mendelson. "Random embeddings with an almost Gaussian distortion", Advances in Mathematics, 2022.
- ▶ S. Mendelson. "An isomorphic Dvoretzky-Milman Theorem using general random ensembles", Journal of Functional Analysis, 2024.
- ▶ D. Bartl, S. Mendelson. "Optimal Non-Gaussian Dvoretzky-Milman Embeddings", International Mathematics Research Notices, 2024.
- ▶ D. Bartl, S. Mendelson. "Structure preservation via the Wasserstein distance", Journal of Functional Analysis, to appear.

Theorem (this paper) Suppose $X = \Sigma^{1/2}\zeta$ for some ζ that has i.i.d. coordinates, satisfying $L^4(\mu_X) - L^2(\mu_X)$ norm equivalence condition, and for some diagonal Σ . The following hold w.h.p.,

$$\forall \boldsymbol{\lambda} \in \mathbb{R}^N, \begin{cases} cl_* \leq \|\mathbb{X}^\top \boldsymbol{\lambda}\|_{q'} \leq C \log(p) \ell_*, & \text{if } q \geq 2, \text{ and } N \leq \kappa_{\text{DM}} d_* \log^{-2}(p^{1/q'} / d_*) \\ cl_* \leq \|\mathbb{X}^\top \boldsymbol{\lambda}\|_{q'} \leq C \ell_*, & \text{if } 1 < q < 2, \text{ and } N \leq \kappa_{\text{DM}} d_*. \end{cases}$$

Come back to classics: fixed points

Recap (intuition): $\hat{\beta}_J \in \operatorname{argmin} (P_N \ell_{\beta_J} + \|\beta_J\|_q^q) \approx \operatorname{argmin} (\|\mathbf{y} - \mathbb{X}_J \beta_J\|_2^q + \ell_*^q \|\beta_J\|_q^q)$.
Reg. $\Psi : \mathcal{F} \rightarrow \mathbb{R}$, convex; empirical risk: $P_N \ell_{\bullet} : f \in \mathcal{F} \mapsto P_N \ell_f \in \mathbb{R}$.

$$\hat{f}_N \in \operatorname{argmin} (P_N \ell_f + \lambda \Psi(f) : f \in \mathcal{F}).$$

Definition (Multiplier & Quadratic fixed points) Let $\kappa, \square, \Delta > 0$, $\Delta \geq 4\square$, $0 < \delta_M, \delta_Q < \frac{1}{2}$. For any $\mathcal{G} \subset \mathcal{F}$, s.t., $f^* \in \mathcal{G}$, define

$$r_M(\mathcal{G}, \delta_M, \kappa, \square) = \min_{r > 0} \left(\mathbb{P} \left(\sup_f \inf_{\mathbf{g} \in \partial^- P_N \ell_{f^*}} |\langle \mathbf{g}, f - f^* \rangle| \leq \square r^{\frac{2}{\kappa}} \right) \geq 1 - \delta_M \right),$$

where sup is over $f \in \mathcal{G} \cap B_{L^2(\mu_X)}(f^*, r)$ and $r_Q(\mathcal{G}, \delta_Q, \kappa, \Delta)$ be

$$\min_{r > 0} \left(\mathbb{P} \left(\forall f \in \mathcal{G} \cap S_{L^2(\mu_X)}(f^*, r), P_N \mathcal{L}_f \geq \Delta r^{\frac{2}{\kappa}} + \sup_{\mathbf{g} \in \partial^- P_N \ell_{f^*}} (\langle \mathbf{g}, f - f^* \rangle) \right) \geq 1 - \delta_Q \right).$$

• J reduces fixed points

Proposition. Under some assumptions, $\exists 0 < \delta_M < \frac{1}{100}$, $c, c' < 1$, $\ell_* > 0$ and $c'' = c''(c, c', \delta_M) > 1$ such that for any $\mathcal{G} \subset V_J$ s.t. $f^* \in \mathcal{G}$, w.h.p.,

$$\begin{cases} r_M \left(\mathcal{G}, \delta_M, \frac{2}{q}, 4c \frac{N^{\frac{q}{2}}}{\ell_*^q} \right) \leq c'' \sigma_\xi \left(\frac{|J|}{N} \right)^{\frac{1}{2(q-1)}}, & \text{when } q \geq 2, \text{ and} \\ r_M \left(\mathcal{G}, \delta_M, 1, 4c' \sigma_\xi^{q-2} \frac{N^{\frac{q}{2}}}{\ell_*^q} \right) \leq c'' \sigma_\xi^{q-1} \left(\frac{|J|}{N} \right)^{\frac{1}{2}}, & \text{when } 1 < q < 2. \end{cases}$$

Proposition. Suppose $|J| \leq cN$. Under some assumptions, $\exists 0 < \delta_Q < \frac{1}{100}$, $c = c(q)$, and $c' = c'(q)$, s.t., when $q \geq 2$. Then $\forall r > 0$, $\forall \mathcal{G} \subset V_J$ s.t. $f^* \in \mathcal{G}$, w.p.a.l. $1 - \delta_Q$,

$$r_Q \left(\mathcal{G}, \delta_Q, \frac{2}{q} \right) = 0.$$

FSD improves classical SLT: applying classical theory on V_J .

What does FSD provide?

V_J is morphism on the category of supervised regression problems and solutions.

• $J : (\mu_X, f^*, \xi, \mathcal{F}, \hat{f}) \mapsto (\mu_X, f_J^*, \zeta, V_J, \hat{f}_J)$, preserves input-output relation,

since $Y = f^*(X) + \xi = f_J^*(X) + \underline{f_{J^c}^*}(X) + \xi = f_J^*(X) + \zeta$, where $\zeta = \xi + f_{J^c}^*$.

The new problem-solution is either easier to study.

- ▶ $\hat{\beta}_J$ is RERM, via self-regularization;
- ▶ β_J^* is new oracle;
- ▶ r_M and r_Q are reduced.

Other role: V_J^* is the true subspace where estimation happens (not this paper).

Classical SLT: identity map, $V_J = \mathcal{F}$.

V_{J^c} provides stochastic properties of $\hat{\beta}_J$.

$$\text{Dvoretzky - Milman : } \forall \mu \in \mathbb{R}^N, \frac{\|\mu\|_2}{(1 + \varepsilon_1)\ell_*(\Sigma_{J^c}^{1/2} B_q^p)} \leq \|\mathcal{A}[\mu]\|_q \leq \frac{\|\mu\|_2}{(1 - \varepsilon_1)\ell_*(\Sigma_{J^c}^{1/2} B_q^p)}.$$

\Rightarrow upper bounds for $\|\Sigma^{1/2}(\hat{\beta}_J - \beta_J^*)\|_2^2$. We are left with $\|\Sigma^{1/2}\beta_{J^c}^*\|_2 + \|\Sigma^{1/2}\hat{\beta}_{J^c}\|_2$.

Free subspace and the energy of $\hat{\beta}_{J^c}$

Remember that $\hat{\beta}_{J^c} = \mathcal{A}[\mathbf{y} - \mathbb{X}_J \hat{\beta}_J]$ in regression. By $N \lesssim d_*(\Sigma_{J^c}^{1/2} B_{q'}^p)$, and the fact

$\|\Sigma_{J^c}^{1/2}\|_{\ell_q \rightarrow \ell_2} = \text{diam}(\Sigma_{J^c}^{1/2} B_q^p)$, there holds

$$\|\Sigma^{1/2} \hat{\beta}_{J^c}\|_2 \leq \|\Sigma_{J^c}^{1/2}\|_{\ell_q \rightarrow \ell_2} \|\mathcal{A}[\mathbf{y} - \mathbb{X}_J \hat{\beta}_J]\|_q \lesssim \frac{\text{diam}(\Sigma_{J^c}^{1/2} B_q^p)}{\ell_*(\Sigma_{J^c}^{1/2} B_q^p)} \|\mathbf{y} - \mathbb{X}_J \hat{\beta}_J\|_2.$$

Since $\|\mathbf{y} - \mathbb{X}_J \hat{\beta}_J\|_2 = \|\mathbb{X}_J(\beta_J^* - \hat{\beta}_J) + \mathbb{X}_{J^c} \beta_{J^c}^* + \boldsymbol{\xi}\|_2 \lesssim \sqrt{N} \sigma_\xi$, and

$$d_*(\Sigma_{J^c}^{1/2} B_{q'}^p) = \left(\frac{\ell_*(\Sigma_{J^c}^{1/2} B_q^p)}{\text{diam}(\Sigma_{J^c}^{1/2} B_q^p)} \right)^2,$$

$$\|\Sigma^{1/2} \hat{\beta}_{J^c}\|_2 \lesssim \sqrt{\frac{N}{d_*(\Sigma_{J^c}^{1/2} B_{q'}^p)}} \sigma_\xi. \quad (5)$$

However, this general bound is not always sharp.

Open direction: non-linear metric embedding

When $q = 2$, $\mathcal{A}[\boldsymbol{\mu}] = \operatorname{argmin}(\|\boldsymbol{\nu}\|_2 : \mathbb{X}_{J^c} \boldsymbol{\nu} = \boldsymbol{\mu}) = \mathbb{X}_{J^c}^\top (\mathbb{X}_{J^c} \mathbb{X}_{J^c}^\top)^{-1} \boldsymbol{\mu}$. Linear!

Let $\boldsymbol{\zeta} = \mathbb{X}_J (\boldsymbol{\beta}_J^* - \hat{\boldsymbol{\beta}}_J) + \mathbb{X}_{J^c} \boldsymbol{\beta}_{J^c}^*$. Then $\|\boldsymbol{\zeta}\|_2 = o(\sqrt{N})$. As a result,

$$\|\Sigma^{1/2} \hat{\boldsymbol{\beta}}_{J^c}\|_2 = \|\Sigma^{1/2} \mathcal{A}[\mathbb{X}_J (\boldsymbol{\beta}_J^* - \hat{\boldsymbol{\beta}}_J) + \mathbb{X}_{J^c} \boldsymbol{\beta}_{J^c}^* + \boldsymbol{\xi}]\|_2 \leq \|\Sigma_{J^c}^{1/2} \mathcal{A}[\boldsymbol{\zeta}]\|_2 + \|\Sigma_{J^c}^{1/2} \mathcal{A}[\boldsymbol{\xi}]\|_2,$$

where $\|\Sigma_{J^c}^{1/2} \mathcal{A}[\boldsymbol{\zeta}]\|_2$ is much smaller, and

$$\mathbb{E}_{\boldsymbol{\xi}} \|\Sigma_{J^c}^{1/2} \mathcal{A}[\boldsymbol{\xi}]\|_2 \leq \sigma_{\boldsymbol{\xi}} (\operatorname{Tr}((\mathbb{X}_{J^c} \mathbb{X}_{J^c}^\top)^{-2} \mathbb{X}_{J^c} \Sigma_{J^c} \mathbb{X}_{J^c}^\top))^{1/2} \leq \sigma_{\boldsymbol{\xi}} \frac{\sqrt{N \operatorname{Tr}(\Sigma_{J^c}^2)}}{\operatorname{Tr}(\Sigma_{J^c})}. \quad (6)$$

However, (5) provides only

$$\sigma_{\boldsymbol{\xi}} \sqrt{\frac{N \|\Sigma_{J^c}\|_{\text{op}}}{\operatorname{Tr}(\Sigma_{J^c})}} \geq \sigma_{\boldsymbol{\xi}} \frac{\sqrt{N \operatorname{Tr}(\Sigma_{J^c}^2)}}{\operatorname{Tr}(\Sigma_{J^c})}. \quad \Rightarrow \quad \text{The tool we used is not correct!}$$

What is the correct mathematical tool for studying $\mathcal{A} : (\mathbb{R}^N, \|\cdot\|_2) \rightarrow (V_{J^c}, \|\cdot\|_q)$'s $\|\Sigma_{J^c}^{1/2} \cdot\|_2$ -norm, that is, $\|\Sigma_{J^c}^{1/2} \mathcal{A}[\boldsymbol{\zeta} + \boldsymbol{\xi}]\|_2^2$, so that it recovers (6) when $q = 2$?

Gaussian case

We say that an FSD is an admissible FSD in the Gaussian case (abbreviated as admissible FSD) if the following conditions are satisfied:

1. $V_J = \text{span}(e_j : j \in J)$, where $J \subset [p]$, and X_J is independent with X_{J^c} .
2. $|J| \lesssim N \lesssim d_*(\Sigma_{J^c}^{-1/2} B_{q'}^p)$.
3. When $1 < q < 2$, $\beta_{J^c}^* = \mathbf{0}$.

For any admissible FSD, and a positive parameter σ_ξ we define

$$R(V_J, V_{J^c}) = \begin{cases} \sigma_\xi \left(\frac{|J|}{N}\right)^{\frac{1}{2(q-1)}} + \|\Sigma_{J^c}^{1/2} \beta_{J^c}^*\|_2 + \|\beta_J^* \odot |\beta_J^*|^{\odot(q-2)}\| \ell_*^{\frac{1}{q-1}} \frac{\ell_*^{\frac{q}{q-1}} (\Sigma_{J^c}^{1/2} B_q^p)}{N^{\frac{q}{2(q-1)}}} & \text{if } q \geq 2, \\ \sigma_\xi^{\frac{1}{3}} \left(\frac{|J|}{N}\right)^{\frac{1}{6}} + \sigma_\xi^{q-2} \|\beta_J^* \odot |\beta_J^*|^{\odot(q-2)}\| \frac{\ell_*^q (\Sigma_{J^c}^{1/2} B_q^p)}{N^{\frac{q}{2}}} & \text{if } 1 < q < 2, \\ \sqrt{\frac{|J|}{N}} + \|\Sigma_{J^c}^{1/2} \beta_{J^c}^*\|_2 + \varepsilon_1 \sigma_\xi & \text{if } q = 1, \end{cases}$$

where ε_1 is the distortion of Dvoretzky-Milman.

Upper bound of the estimation error

Theorem (Gaussian case, informal)

Suppose ξ is independent of X , $\xi \sim \mathcal{N}(0, \sigma_\xi^2)$, and $X \sim \mathcal{N}(\mathbf{0}, \Sigma)$. Let $\hat{\beta}$ be the minimum ℓ_q -norm interpolant estimator in the linear regression problem. For any admissible FSD in the Gaussian case, when N is large enough, the following hold w.h.p.,

1. When $q > 1$,

$$\|\Sigma^{1/2}(\hat{\beta} - \beta^*)\|_2 \lesssim r(V_J, V_{J^c}) := R(V_J, V_{J^c}) + \sqrt{\frac{N}{d_*(\Sigma_{J^c}^{-1/2} B_{q'}^p)}} \sigma_\xi.$$

2. When $q = 1$, suppose there exists an absolute constant $c_2 \in (0, 1)$ such that

$$(\ell^*(\Sigma_{J^c}^{1/2} B_1^p))^2 \leq c_2 N \left(\sum_{j \in J \cap \text{supp}(\beta^*)} \sigma_j^{-1} \right)^{-1}.$$

Then the same upper bound holds.

Rmk: Gaussian case is universal. **Rmk:** to obtain B.O., we need $r(V_J, V_{J^c}) \rightarrow 0$.

Sufficient conditions for B.O.

Sufficient conditions (informal): there exists admissible FSD, s.t.

$$\begin{array}{ccc} \left| \left| \beta_J^* \odot |\beta_J^*|^{\odot(q-2)} \right| \right| = o\left(\frac{N^{q/2}}{\ell_*^q (\Sigma_{J^c}^{1/2} B_q^p)}\right) & , & \|\Sigma_{J^c}^{1/2} \beta_{J^c}^*\|_2 = o(1) , & \text{bias } \hat{\beta} \\ |J| = o(N) & , \text{ and } & N = o(d_*(\Sigma_{J^c}^{-1/2} B_{q'}^p)) & \text{variance } \hat{\beta} \\ \downarrow & & \downarrow & \\ \hat{\beta}_J \text{ bias and variance} & & \hat{\beta}_{J^c} \text{ "bias" and "variance"} & \end{array}$$

Interpretation:

- ▶ $\hat{\beta}_J \approx \operatorname{argmin}\left(\frac{1}{N^{q/2}} \|\mathbf{y} - \mathbb{X}_J \beta_J\|_2^q + \frac{\ell_*^q}{N^{q/2}} \|\beta_J\|_q^q\right).$
 - ▶ Bias $\left| \left| \beta_J^* \odot |\beta_J^*|^{\odot(q-2)} \right| \right| = o\left(\frac{N^{q/2}}{\ell_*^q}\right),$
 - ▶ Variance $|J| = o(N).$
- ▶ $\hat{\beta}_{J^c} = \mathcal{A}[\mathbf{y} - \mathbb{X}_J \hat{\beta}_J].$
 - ▶ $\beta_{J^c}^*$ is not estimated \Rightarrow model noise $\|\Sigma_{J^c}^{1/2} \beta_{J^c}^*\|_2.$
 - ▶ Variance $N = o(d_*(\Sigma_{J^c}^{-1/2} B_{q'}^p)).$

Take-home messages

- ▶ **Improvement of classical SLT via FSD.**

- ▶ estimation ONLY happens on V_J instead of the entire \mathcal{F} .
- ▶ the estimation error $\|\Sigma^{1/2}(\hat{\beta} - \beta^*)\|_2$ can be controlled by interpolating cancellation $\|\Sigma^{1/2}(\hat{\beta}_J - \beta_J^*)\|_2$, and smallness $\|\Sigma^{1/2}\hat{\beta}_J\|_2 + \|\Sigma^{1/2}\beta_J^*\|_2$.

- ▶ **Reason of B.O.: Self-regularization**, $\hat{\beta}_J$ is a RERM with random loss function.

$$\hat{\beta}_J \in \operatorname{argmin}_{\beta_J \in V_J} (P_N \ell_{\beta_J} + \|\beta_J\|_q^q), \text{ where } P_N \ell_{\beta_J} = \|\mathcal{A}[\mathbf{y} - \mathbb{X}\beta_J]\|_q^q, \text{ and } \beta_J^* \in \operatorname{argmin}_{\beta_J \in V_J} P \ell_{\beta_J}.$$

- ▶ **AGA results for stats**: D-M theorem provides stochastic properties of $P_N \ell_{\bullet}$.

- ▶ **AGA results from stats**: Generalized Dvoretzky-Milman.

- ▶ **Implications in B.O.**

- ▶ B.O. depends on the alignment between β^* and canonical basis.
- ▶ Gaussian case is universal.

- ▶ **Mathematical challenge**: non-linear metric embedding.

End

The 'hard marl' of Benign Overfitting has finally been surrounded and dissolved by the geometric structure we built. What was a paradox in the classical uniform convergence world becomes a lucid geometric reality here.

Thank you for your attention!

FSD in classification

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

$$P\mathcal{L}_{\hat{f}}^{(0,1)} = \mathbb{P}\left(Y\hat{f}(X) < 0 \mid (X_i, Y_i)_{i=1}^N\right) - \mathbb{P}\left(Y\hat{f}_J(X) < 0 \mid (X_i, Y_i)_{i=1}^N\right) \quad (7)$$

$$+ \mathbb{P}\left(Y\hat{f}_J(X) < 0 \mid (X_i, Y_i)_{i=1}^N\right) - \mathbb{P}\left(Yf_J^*(X) < 0\right) \quad (8)$$

$$+ \mathbb{P}\left(Yf_J^*(X) < 0\right) - \mathbb{P}\left(Y\left(\eta(X) - \frac{1}{2}\right) < 0\right), \quad (9)$$

- ▶ (7): error caused by free part \hat{f}_{J^c} ;
- ▶ (8): prediction error caused by \hat{f}_J compared to that of f_J^* ; and
- ▶ (9): 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.

$\hat{\beta}_J$ is “almost” a support vectors machine with squared hinge loss

Proposition. Recall $P_N \ell_{\beta_J} = \|\mathcal{B}[\mathbb{1} - \mathbb{X}_y \beta_J]\|_2^2$ in classification. Let $P_N \ell_{\beta_J}^{(\text{sh})} = \frac{1}{N} \sum_{i=1}^N (1 - Y_i \langle \beta_J, X_i \rangle)_+^2$ be the squared hinge loss. Then w.h.p.,

$$\hat{\beta}_J \in \operatorname{argmin} (P_N \ell_{\beta_J} + \|\beta_J\|_2^2) \approx \operatorname{argmin} \left(\frac{N}{\operatorname{Tr}(\Sigma_{J^c})} P_N \ell_{\beta_J}^{(\text{sh})} + \|\beta_J\|_2^2 \right).$$

Proof. The dual problem of $\|\mathcal{B}[\mu]\|_2$ is $\max(\langle \lambda, \mu \rangle : \lambda \succeq \mathbf{0}, \|\mathbb{X}_{J^c}^\top \lambda\|_2 \leq 1)$. Let $H(\mu) := \{i \in [N] : \mu_i < 0\}$ and let

$$\lambda^- \in \operatorname{argmax}(\langle \lambda, \mu \rangle : \lambda \succeq \mathbf{0}, (1 + \delta) \sqrt{\operatorname{Tr}(\Sigma_{J^c})} \|\lambda\|_2 \leq 1). \quad (10)$$

Suppose $i \in H(\mu)$ but $\lambda_i^- > 0$. Let $\tilde{\lambda}^- = (\lambda_1^-, \dots, \lambda_{i-1}^-, 0, \lambda_{i+1}^-, \dots, \lambda_N^-)$, then $\|\tilde{\lambda}^-\|_2 < \|\lambda^-\|_2 \leq \frac{1}{(1+\delta)\sqrt{\operatorname{Tr}(\Sigma_{J^c})}}$. Since $\mu_i \lambda_i^- < 0$,

$$\langle \mu, \tilde{\lambda}^- \rangle = \sum_{i' \neq i} \mu_{i'} \lambda_{i'}^- > \sum_{i'=1}^N \mu_{i'} \lambda_{i'}^- = \langle \mu, \lambda^- \rangle, \Rightarrow \lambda^- \text{ is not the maximizer of (10).}$$

Necessarily, $\forall i \in H(\mu), \lambda_i^- = 0$. Therefore, $\lambda^- = \left(\frac{\mu}{(\|\mu\|_2 (1+\delta) \sqrt{\operatorname{Tr}(\Sigma_{J^c})})} \right)_+$. ■

Free subspace and the energy of $\hat{\beta}_{J^c}$ in classification

Let $\mathcal{F} = V_J \oplus V_{J^c}$ be any FSD and \hat{f}_N be any estimator. There holds $\mu^{\otimes N}$ -a.s.,

$$\begin{aligned} & \mathbb{P} \left(Y \hat{f}(X) < 0 \mid (X_i, Y_i)_{i=1}^N \right) - \mathbb{P} \left(Y \hat{f}_J(X) < 0 \mid (X_i, Y_i)_{i=1}^N \right) \\ & \leq \mathbb{P} \left(|\hat{f}_{J^c}(X)| > |\hat{f}_J(X)| \mid (X_i, Y_i)_{i=1}^N \right). \end{aligned}$$

Classical theory on RERM

We say Ψ has non-trivial Bregman div. around f^* , if $D_\Psi(\cdot, f^*)$ is bounded from below by convex, non-negative func., where $D_\Psi(f, g) = \Psi(f) - \Psi(g) - \langle \nabla \Psi(g), f - g \rangle$. For any $\rho > 0$, let $B_\Psi(f^*, \rho) = \{f \in \mathcal{F} : D_\Psi(f, f^*) \leq \rho\}$.

E.g., $\Psi(\cdot) = \|\cdot\|_q^q$, then $D_\Psi(\mathbf{v}_1, \mathbf{v}_2) \geq c\alpha_q(|\mathbf{v}_2|, \mathbf{v}_1 - \mathbf{v}_2)$, where

$$\alpha_q(x, y) = \begin{cases} \frac{q}{2}x^{q-2}y^2, & |y| \leq x, \\ |y|^q + (\frac{q}{2} - 1)x^q, & \text{otherwise.} \end{cases}$$

Theorem

Suppose Ψ has non-trivial Breg. div. around f^* . For any $\rho > 0$, let $\mathcal{G} = B_\Psi(f^*, \rho)$ for any $\lambda > 0$, let $r_{\text{iso}}(\rho) = r_{\text{iso}}(\mathcal{G}, \kappa)$. Let r_* and ρ_* be the smallest r and its corresponding ρ , s.t.,

$$r \geq r_{\text{iso}}(\rho), \quad 3\Box r^{\frac{2}{\kappa}} > \lambda \|\|\nabla \Psi(f^*)\|\|_{(r, \rho)}, \quad \text{and} \quad \rho \geq \frac{1}{\lambda} \Delta r^{\frac{2}{\kappa}}$$

where $\|\|\nabla \Psi(f^*)\|\|_{(r, \rho)} = \sup (\langle \nabla \Psi(f^*), f - f^* \rangle : f \in B_\Psi(f^*; \rho) \cap B_{L^2(\mu_X)}(f^*; r))$.

Then $\|\hat{f} - f^*\|_{L^2(\mu_X)}^2 \leq r_*^2$, w.p.a.l. $1 - \delta_M - \delta_Q$.

Notations

For any $x \geq 0$ and $y \in \mathbb{R}$, let

$$\alpha_q(x, y) = \begin{cases} \frac{q}{2}x^{q-2}y^2, & \text{if } |y| \leq x \\ |y|^q + \left(\frac{q}{2} - 1\right)x^q, & \text{otherwise.} \end{cases}$$

For any $\rho > 0$, we define

$$\rho K_{\text{model}} = \begin{cases} \{\mathbf{v} \in V_J : \sum_{j \in J} \alpha_q(|\beta_j^*|, v_j) \leq \rho^q\} & \text{when } q < 2, \\ \rho B_q^J & \text{when } q \geq 2. \end{cases}$$

Define

$$\|\|\beta\|\| = \sup \left(\langle \beta, \mathbf{u} \rangle : \mathbf{u} \in \frac{C_1 \sqrt{N}}{\ell_*(\Sigma_{J^c}^{1/2} B_q^p)} K_{\text{model}} \cap \Sigma_J^{-1/2} B_2^J \right),$$

where $C_1 = C_1(q)$ is some absolute constant.

Weak moments case

Suppose $\Sigma^{-1/2}X$ has i.i.d. coordinates that satisfy $L^{8+\varepsilon} - L^2$ norm equivalence assumption.

We say that an FSD $\mathbb{R}^p = V_J \oplus V_{J^c}$ is admissible in the heavy-tailed case if the following conditions are satisfied:

1. $V_J = \text{span}(e_j : j \in J)$, where $J \subset [p]$.
2. There exist absolute constants $0 < \varepsilon_1, \kappa_{RIP}, \kappa_{DM} < 1$ and $C_2 > 1$ such that
 - ▶ when $1 < q < 2$, there holds $N^{\frac{4}{4+\varepsilon}} \lesssim |J| \lesssim N \lesssim d_*(\Sigma_{J^c}^{-1/2} B_{q'}^p)$, and
 - ▶ when $q \geq 2$, there holds

$$N^{\frac{4}{4+\varepsilon}} \lesssim |J| \leq N \lesssim \frac{d_*(\Sigma_{J^c}^{-1/2} B_{q'}^p)}{\log^2(|J^c|^{1/q'} / d_*(\Sigma_{J^c}^{-1/2} B_{q'}^p))}.$$

In particular, when X is sub-Gaussian, the \log factor can be removed.

Universality of Gaussian case

Theorem (Universality, informal)

Suppose Σ is diagonal. For any admissible FSD in the heavy-tailed case, when N is large enough, and the following hold w.h.p..

1. When $q \geq 2$, assume there exists an absolute constant C_3 such that $(\mathbb{E}[|\xi|^4])^{1/4} \leq C_3 \sigma_\xi$.
Then

$$\|\Sigma^{1/2}(\hat{\beta} - \beta^*)\|_2 \lesssim R(V_J, V_{J^c}) \log^{\frac{q}{q-1}}(|J^c|) + \sqrt{\frac{N}{d_*(\Sigma_{J^c}^{-1/2} B_{q'}^p)}} \sigma_\xi.$$

In particular, when X is sub-Gaussian, the \log factor can be removed.

2. When $1 < q < 2$, assume additionally that $X_J \sim \mathcal{N}(\mathbf{0}, \Sigma_J)$, $\xi \sim \mathcal{N}(0, \sigma_\xi^2)$, and $\beta_{J^c}^* = \mathbf{0}$. If X_{J^c} satisfies the $L^{8+\epsilon}$ - L^2 equivalence condition, then the same result holds.

• J reduces r_M

Recall that in regression, $P_N \ell_{\beta_J} = \|\mathcal{A}[\mathbf{y} - \mathbb{X}\beta_J]\|_q^q$ and

$$\partial^- P_N \ell_{\beta_J} = \left\{ -q \|\mathcal{A}[\mathbf{y} - \mathbb{X}_J \beta_J]\|_q^{q-1} \mathbb{X}_J^\top \boldsymbol{\lambda}^*[\mathbf{y} - \mathbb{X}_J \beta_J] : \boldsymbol{\lambda}^*[\cdot] \text{ is dual solution} \right\}. \quad (11)$$

Proposition. Under some assumptions, $\exists 0 < \delta_M < \frac{1}{100}$, $c, c' < 1$, $\ell_* > 0$ and $c'' = c''(c, c', \delta_M) > 1$ such that for any local. subset $\mathcal{G} \subset V_J$,

$$\begin{cases} r_M \left(\mathcal{G}, \delta_M, \frac{2}{q}, 4c \frac{N^{\frac{q}{2}}}{\ell_*^q} \right) \leq c'' \sigma_\xi \left(\frac{|J|}{N} \right)^{\frac{1}{2(q-1)}}, & \text{when } q \geq 2, \text{ and} \\ r_M \left(\mathcal{G}, \delta_M, 1, 4c' \sigma_\xi^{q-2} \frac{N^{\frac{q}{2}}}{\ell_*^q} \right) \leq c'' \sigma_\xi^{q-1} \left(\frac{|J|}{N} \right)^{\frac{1}{2}}, & \text{when } 1 < q < 2. \end{cases}$$

Proof. For any $\mathbf{g} \in \partial^- P_N \ell_{\beta_J^*}$ from (11), by $B_{L^2(\mu_X)}(\beta_J^*, r) = \beta_J^* + r \Sigma_J^{-1/2} B_2$,

$$\begin{aligned} \sup_{\beta_J \in \mathcal{G} \cap B_{L^2(\mu_X)}(\beta_J^*, r)} |\langle \mathbf{g}, \beta_J - \beta_J^* \rangle| &= q \|\mathcal{A}[\mathbb{X}_{J^c} \beta_J^* + \boldsymbol{\xi}]\|_q^{q-1} \sup_{\mathbf{v} \in \mathcal{G} \cap r \Sigma_J^{-1/2} S_2} \sum_{i=1}^N \lambda_i^*[\mathbb{X}_{J^c} \beta_{J^c}^* + \boldsymbol{\xi}] \langle X_i, \mathbf{v} \rangle \\ &\leq C_q \|\mathcal{A}[\mathbb{X}_{J^c} \beta_J^* + \boldsymbol{\xi}]\|_q^{q-1} \|\boldsymbol{\lambda}^*[\mathbb{X}_{J^c} \beta_{J^c}^* + \boldsymbol{\xi}]\|_2 \gamma_2(\Sigma_J^{1/2} \mathcal{G} \cap r B_2^J, \|\cdot\|_2). \end{aligned}$$

- J reduces r_Q

Proposition

Suppose $|J| \leq cN$. Under the assumptions, $\exists 0 < \delta_Q < \frac{1}{100}$, $c = c(q)$, and $c' = c'(q)$, s.t., when $q \geq 2$. Then $\forall r > 0, \forall \mathcal{G}$, w.p.a.l. $1 - \delta_Q$, for any local. sub. $\beta_J \in \mathcal{G} \cap (\beta_J^* + r\Sigma_J^{-1/2}S_2)$,

$$P_N \mathcal{L}_{\beta_J} = \|\mathcal{A}[\mathbf{y} - \mathbb{X}\beta_J]\|_q^q - \|\mathcal{A}[\mathbf{y} - \mathbb{X}\beta_J^*]\|_q^q \geq \langle \mathbf{g}, \beta_J - \beta_J^* \rangle + \Delta r^q, \text{ where } \Delta = c \frac{N^{\frac{q}{2}}}{\ell_*^q}.$$

That is, $r_Q(\mathcal{G}, \delta_Q, \frac{2}{q}) = 0, \forall \mathcal{G}$.

Proof.

$$\begin{aligned} P_N \mathcal{L}_{\beta_J} &\geq \langle \mathbf{g}, \beta_J - \beta_J^* \rangle + c_q \|\mathcal{A}[\mathbf{y} - \mathbb{X}\beta_J] - \mathcal{A}[\mathbf{y} - \mathbb{X}\beta_J^*]\|_q^q \\ &\geq \langle \mathbf{g}, \beta_J - \beta_J^* \rangle + c_q \|\mathcal{A}[\mathbb{X}(\beta_J - \beta_J^*)]\|_q^q \\ &\geq \langle \mathbf{g}, \beta_J - \beta_J^* \rangle + c_q \frac{\|\mathbb{X}(\beta_J - \beta_J^*)\|_2^q}{\ell_*^q} \geq \langle \mathbf{g}, \beta_J - \beta_J^* \rangle + c'_q \frac{N^{\frac{q}{2}}}{\ell_*^q} \|\Sigma_J^{1/2}(\beta_J - \beta_J^*)\|_2^q. \end{aligned}$$

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.