

472 **A Other related work.**

473 Besides Parhi and Nowak [27] which we discussed earlier, Parhi and Nowak [28, 29] also leveraged
 474 the connections between NNs and splines. Parhi and Nowak [28] focused on characterizing the
 475 variational form of multi-layer NN. Parhi and Nowak [29] showed that two-layer ReLU activated
 476 NN achieves minimax rate for a BV class of order 1 but did not cover multilayer NNs nor BV class
 477 with order > 1 , which is our focus.

478 The connection between weight-decay regularization with sparsity-inducing penalties in two-layer
 479 NNs is folklore and used by Neyshabur et al. [25], Savarese et al. [32], Ongie et al. [26], Ergen
 480 and Pilanci [12, 14], Parhi and Nowak [27, 29]. The key underlying technique — an application
 481 of the AM-GM inequality (which we used in this paper as well) — can be traced back to Srebro
 482 et al. [35] (see a recent exposition by Tibshirani [39]). Tibshirani [39] also generalized the result to
 483 multi-layered NNs, but with a simple (element-wise) connections.

484 The approximation-theoretic and estimation-theoretic research for neural network has a long history
 485 too [6, 4, 46, 33, 36]. Most existing work considered the Holder, Sobolev spaces and their exten-
 486 sions, which contain only homogeneously smooth functions and cannot demonstrate the advantage
 487 of NNs over kernels. The only exception is Suzuki [36] which, as we discussed earlier, requires
 488 modifications to NN architecture for each class. In contrast, we require tuning only the standard
 489 weight decay parameter.

490 **B Two-layer Neural Network with Truncated Power Activation Functions**

491 We start by recapping the result of Parhi and Nowak [27] and formalizing its implication in esti-
 492 mating BV functions. Parhi and Nowak [27] considered a two layer neural network with truncated
 493 power activation function. Let the neural network be

$$f(x) = \sum_{j=1}^M v_j \sigma^m(w_j x + b_j) + c(x), \quad (7)$$

494 where w_j, v_j denote the weight in the first and second layer respectively, b_j denote the bias in the
 495 first layer, $c(x)$ is a polynomial of order up to m , $\sigma^m(x) := \max(x, 0)^m$. Parhi and Nowak [27,
 496 Theorem 8] showed that when M is large enough, The optimization problem

$$\min_{w, v} \hat{L}(f) + \frac{\lambda}{2} \sum_{j=1}^M (|v_j|^2 + |w_j|^{2m}) \quad (8)$$

497 is equivalent to the locally adaptive regression spline:

$$\min_f \hat{L}(f) + \lambda TV(f^{(m)}(x)), \quad (9)$$

498 which optimizes over arbitrary functions that is m -times weakly differentiable. The latter was
 499 studied in Mammen and van de Geer [22], which leads to the following MSE:

500 **Theorem 9.** *Let $M \geq n - m$, and \hat{f} be the function (7) parameterized by the minimizer of (8), then*

$$\text{MSE}(\hat{f}) = O(n^{-(2m+2)(2m+3)}).$$

501 We show a simpler proof in the univariate case due to Tibshirani [40]:

502 *Proof.* As is shown in Parhi and Nowak [27, Theorem 8], the minimizer of (8) satisfy

$$|v_j| = |w_j|^m, \forall k$$

503 so the TV of the neural network f_{NN} is

$$\begin{aligned} TV^{(m)}(f_{NN}) &= TV^{(m)}c(x) + \sum_{j=1}^M |v_j||w_j|^m TV^{(m)}(\sigma^{(m)}(x)) \\ &= \sum_{j=1}^M |v_j||w_j|^m \\ &= \frac{1}{2} \sum_{j=1}^M (|v_j|^2 + |w_j|^{2m}) \end{aligned}$$

504 which shown that (8) is equivalent to the locally adaptive regression spline (9) as long as the number
505 of knots in (9) is no more than M . Furthermore, it is easy to check that any spline with knots no
506 more than M can be expressed as a two layer neural network (8). It suffices to prove that the solution
507 in (9) has no more than $n - m$ number of knots.

508 Mammen and van de Geer [22, Proposition 1] showed that there is a solution to (9) $\hat{f}(x)$ such that
509 $\hat{f}(x)$ is a m th order spline with a finite number of knots but did not give a bound. Let the number of
510 knots be M , we can represent \hat{f} using the truncated power basis

$$\hat{f}(x) = \sum_{j=1}^M a_j (x - t_j)_+^m + c(x) := \sum_{j=1}^M a_j \sigma_j^{(m)}(x) + c(x)$$

511 where t_j are the knots, $c(x)$ is a polynomial of order up to m , and define $\sigma_j^{(m)}(x) = (x - t_j)_+^m$.

512 Mammen and van de Geer [22] however did not give a bound on M . Parhi and Nowak [27]'s
513 Theorem 1 implies that $M \leq n - m$. Its proof is quite technical and applies more generally to a
514 higher dimensional generalization of the BV class.

515 Tibshirani [40] communicated to us the following elegant argument to prove the same using elemen-
516 tary convex analysis and linear algebra, which we present below.

517 Define $\Pi_m(f)$ as the $L^2(P_n)$ projection of f onto polynomials of degree up to m , $\Pi_m^\perp(f) :=$
518 $f - \Pi_m(f)$. It is easy to see that

$$\Pi_m^\perp f(x) = \sum_{j=1}^M a_j \Pi_m^\perp \sigma_j^{(m)}(x)$$

519 Denote $f(x_{1:n}) := \{f(x_1), \dots, f(x_n)\} \in \mathbb{R}^n$ as a vector of all the predictions at the sample points.

$$\Pi_m^\perp \hat{f}(x_{1:n}) = \sum_{j=1}^M a_j \Pi_m^\perp \sigma_j^{(m)}(x_{1:n}) \in \Pi_m^\perp \text{conv}\{\pm \sigma_j^{(m)}(x_{1:n})\} \cdot \sum_{j=1}^M |a_j| \in \text{conv}\{\pm \Pi_m^\perp \sigma_j^{(m)}(x_{1:n})\} \cdot \sum_{j=1}^M |a_j|$$

520 where conv denotes the convex hull of a set. The convex hull $\text{conv}\{\pm \sigma_j^{(m)}(x_{1:n})\} \cdot \sum_{j=1}^M |a_j|$ is an
521 n -dimensional space, and polynomials of order up to m is an $m + 1$ dimensional space, so the set
522 defined above has dimension $n - m - 1$. By Carathéodory's theorem, there is a subset of points in
523 this space

$$\{\Pi_m^\perp \sigma_{j_k}^{(m)}(x_{1:n})\} \subseteq \{\Pi_m^\perp \sigma_j^{(m)}(x_{1:n})\}, 1 \leq k \leq n - m$$

524 such that

$$\Pi_m^\perp f(x) = \sum_{k=1}^{n-m} \tilde{a}_k \Pi_m^\perp \sigma_{j_k}^{(m)}(x), \sum_{k=1}^{n-m} |a_k| \leq 1$$

525 In other word, there exist a subset of knots $\{\tilde{t}_j, j \in [n - m]\}$ that perfectly recovers $\Pi_m^\perp \hat{f}(x)$ at all
526 the sample points, and the TV of this function is no larger than \hat{f} .

This shows that

$$\tilde{f}(x) = \sum_{j=1}^{n-m} \tilde{a}_j (x - t_j)_+^m, \text{ s.t. } \tilde{f}(x_i) = f(x_i)$$

527 for all x_i in n onbervation points.

528 The MSE of locally adaptivity regressive spline (9) was studied in Mammen and van de Geer [22,
529 Section 3], which equals the error rate given in Theorem 9. \square

530 This indicates that the neural network (7) is minimax optimal for $BV(m)$.

Let us explain a few the key observations behind this equivalence. (a) The truncated power functions (together with an m th order polynomial) spans the space of an m th order spline. (b) The neural network in (7) is equivalent to a free-knot spline with M knots (up to reparameterization). (c) A solution to (9) is a spline with at most $n - m$ knots [27, Theorem 8]. (d) Finally, by the AM-GM inequality

$$|v_j|^2 + |w_j|^{2m} \geq 2|v_j||w_j|^m = 2|c_j|$$

531 where $c_j = v_j|w_j|^m$ is the coefficient of the corresponding j th truncated power basis. The m th
532 order total variation of a spline is equal to $\sum_j |c_j|$. It is not hard to check that the loss function
533 depends only on c_j , thus the optimal solution will always take “=” in the AM-GM inequality.

534 C Introduction To Common Function Classes

535 In the following definition define Ω be the domain of the function classes, which will be omitted in
536 the definition.

537 C.1 Besov Class

538 **Definition 1.** *Modulus of smoothness:* For a function $f \in L^p(\Omega)$ for some $1 \leq p \leq \infty$, the r -th
539 modulus of smoothness is defined by

$$w_{r,p}(f, t) = \sup_{h \in \mathbb{R}^d: \|h\|_2 \leq t} \|\Delta_h^r(f)\|_p,$$

540

$$\Delta_h^r(f) := \begin{cases} \sum_{j=0}^r \binom{r}{j} (-1)^{r-j} f(x + jh), & \text{if } x \in \Omega, x + rh \in \Omega, \\ 0, & \text{otherwise.} \end{cases}$$

541 **Definition 2.** *Besov space:* For $1 \leq p, q \leq \infty, \alpha > 0, r := \lceil \alpha \rceil + 1$, define

$$|f|_{B_{p,q}^\alpha} = \begin{cases} \left(\int_{t=0}^\infty (t^{-\alpha} w_{r,p}(f, t))^q \frac{dt}{t} \right)^{\frac{1}{q}}, & q < \infty \\ \sup_{t>0} t^{-\alpha} w_{r,p}(f, t), & q = \infty, \end{cases}$$

542 and define the norm of Besov space as:

$$\|f\|_{B_{p,q}^\alpha} = \|f\|_p + |f|_{B_{p,q}^\alpha}.$$

543 A function f is in the Besov space $B_{p,q}^\alpha$ if $\|f\|_{B_{p,q}^\alpha}$ is finite.

544 Note that the Besov space for $0 < p, q < 1$ is also defined, but in this case it is a quasi-Banach space
545 instead of a Banach space and will not be covered in this paper.

546 Functions in Besov space can be decomposed using B-spline basis functions. Any function f in
547 Besov space $B_{p,q}^\alpha, \alpha > d/p$ can be decomposed using B-spline of order $m, m > \alpha$: let $\mathbf{x} \in \mathbb{R}^d$,

$$f(\mathbf{x}) = \sum_{k=0}^{\infty} \sum_{\mathbf{s} \in J(k)} c_{k,\mathbf{s}}(f) M_{m,k,\mathbf{s}}(\mathbf{x}) \quad (10)$$

548 where $J(k) := \{2^{-k} \mathbf{s} : \mathbf{s} \in [-m, 2^k + m]^d \subset \mathbb{Z}^d\}$, $M_{m,k,\mathbf{s}}(\mathbf{x}) := M_m(2^k(\mathbf{x} - \mathbf{s}))$, and $M_k(\mathbf{x}) =$
549 $\prod_{i=1}^d M_k(x_i)$ is the cardinal B-spline basis function which can be expressed as a polynomial:

$$M_m(x) = \frac{1}{m!} \sum_{j=1}^{m+1} (-1)^j \binom{m+1}{j} (x-j)_+^m = ((m+1)/2)^m \frac{1}{m!} \sum_{j=1}^{m+1} (-1)^j \binom{m+1}{j} \left(\frac{x-j}{(m+1)/2} \right)_+^m,$$

550 Furthermore, the norm of Besov space is equivalent to the sequence norm:

$$\|\{c_{k,\mathbf{s}}\}\|_{b_{p,q}^\alpha} := \left(\sum_{k=0}^{\infty} (2^{(\alpha-d/p)k} \|\{c_{k,\mathbf{s}}(f)\}_{\mathbf{s}}\|_p)^q \right)^{1/q} \approx \|f\|_{B_{p,q}^\alpha}.$$

551 See e.g. Dũng [11, Theorem 2.2] for the proof.

552 C.2 Other Function Spaces

553 **Definition 3.** *Hölder space: let $m \in \mathbb{N}$, the m -th order Hölder class is defined as*

$$\mathcal{C}^m = \left\{ f : \max_{|a|=k} \frac{|D^a f(x) - D^a f(z)|}{\|x - z\|_2} < \infty, \forall x, z \in \Omega \right\}$$

554 where D^a denotes the weak derivative.

555 Note that fraction order of Hölder space can also be defined. For simplicity, we will not cover that
556 case in this paper.

557 **Definition 4.** *Sobolev space: let $m \in \mathcal{N}$, $1 \leq p \leq \infty$, the Sobolev norm is defined as*

$$\|f\|_{W_p^m} := \left(\sum_{|a| \leq m} \|D^a f\|_p^p \right)^{1/p},$$

558 the Sobolev space is the set of functions with finite Sobolev norm:

$$W_p^m := \{f : \|f\|_{W_p^m} < \infty\}.$$

559 **Definition 5.** *Total Variation (TV): The total variation (TV) of a function f on an interval $[a, b]$ is
560 defined as*

$$TV(f) = \sup_{\mathcal{P}} \sum_{i=1}^{n_{\mathcal{P}}-1} |f(x_{i+1}) - f(x_i)|$$

561 where the \mathcal{P} is taken among all the partitions of the interval $[a, b]$.

562 In many applications, functions with stronger smoothness conditions are needed, which can be mea-
563 sured by high order total variation.

564 **Definition 6.** *High order total variation: the m -th order total variation is the total variation of the
565 $(m - 1)$ -th order derivative*

$$TV^{(m)}(f) = TV(f^{(m-1)})$$

566 **Definition 7.** *Bounded variation (BV): The m -th order bounded variation class is the set of functions
567 whose total variation (TV) is bounded.*

$$BV(m) := \{f : TV(f^{(m)}) < \infty\}.$$

568 D Proof of Estimation Error

569 D.1 Equivalence Between Parallel Neural Networks and p -norm Penalized Problems

570 **Proposition 3.** *Fix the input dataset \mathcal{D}_n and a constant $c_1 > 0$. For every λ , there exists $P' > 0$
571 such that (2) is equivalent to the following problem:*

$$\begin{aligned} \arg \min_{\{\bar{\mathbf{w}}_j^{(\ell)}, \bar{\mathbf{b}}_j^{(\ell)}, a_j\}} \hat{L} \left(\sum_{j=1}^M a_j \bar{f}_j \right) &= \frac{1}{n} \sum_i (y_i - \bar{f}_{1:M}(\mathbf{x}_i)^T \mathbf{a})^2 \\ \text{s.t. } \|\bar{\mathbf{w}}_j^{(1)}\|_F &\leq c_1 \sqrt{d}, \forall j \in [M], \\ \|\bar{\mathbf{w}}_j^{(\ell)}\|_F &\leq c_1 \sqrt{w}, \forall j \in [M], 2 \leq \ell \leq L, \quad \|\{a_j\}\|_{2/L}^{2/L} \leq P' \end{aligned}$$

572 where $\bar{f}_j(\cdot)$ is a subnetwork with parameters $\bar{\mathbf{w}}_j^{(\ell)}, \bar{\mathbf{b}}_j^{(\ell)}$.

573 *Proof.* Using Langrange’s method, one can easily find (2) is equivalent to a constrained optimization
 574 problem:

$$\arg \min_{\{\mathbf{W}_j^{(\ell)}, \mathbf{b}_j^{(\ell)}\}} \hat{L} \left(\sum_{j=1}^M f_j \right), \quad s.t. \sum_{j=1}^M \sum_{\ell=1}^L \|\mathbf{W}_j^{(\ell)}\|_F^2 \leq P \quad (11)$$

575 for some constant P that depends on λ and the dataset \mathcal{D} .

576 We make use of the property from (4) to minimize the constraint term in (11) while keeping this
 577 neural network equivalent to the original one. Specifically, let $\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(L)}$ be the
 578 parameters of an L -layer neural network.

$$f(x) = \mathbf{W}^{(L)} \sigma(\mathbf{W}^{(L-1)} \sigma(\dots \sigma(\mathbf{W}^{(1)} x + \mathbf{b}^{(1)}) \dots) + \mathbf{b}^{(L-1)}) + \mathbf{b}^{(L)},$$

579 which is equivalent to

$$f(x) = \alpha_L \tilde{\mathbf{W}}^{(L)} \sigma(\alpha_{L-1} \tilde{\mathbf{W}}^{(L-1)} \sigma(\dots \sigma(\alpha_1 \tilde{\mathbf{W}}^{(1)} x + \tilde{\mathbf{b}}^{(1)}) \dots) + \tilde{\mathbf{b}}^{(L-1)}) + \tilde{\mathbf{b}}^{(L)},$$

580 as long as $\alpha_\ell > 0$, $\prod_{\ell=1}^L \alpha_\ell = \prod_{\ell=1}^L \|\mathbf{W}^{(\ell)}\|_F$, where $\tilde{\mathbf{W}}^{(\ell)} := \frac{\mathbf{W}^{(\ell)}}{\|\mathbf{W}^{(\ell)}\|_F}$. By the AM-GM inequal-
 581 ity, the ℓ_2 regularizer of the latter neural network is

$$\sum_{\ell=1}^L \|\alpha_\ell \tilde{\mathbf{W}}^{(\ell)}\|_F^2 = \sum_{\ell=1}^L \alpha_\ell^2 \geq L \left(\prod_{\ell=1}^L \alpha_\ell \right)^{2/L} = L \left(\prod_{\ell=1}^L \|\mathbf{W}^{(\ell)}\|_F \right)^{2/L}$$

582 and equality is reached when $\alpha_1 = \alpha_2 = \dots = \alpha_L$. In other word, in the problem (2), it suffices to
 583 consider the network that satisfies

$$\|\mathbf{W}_j^{(1)}\|_F = \|\mathbf{W}_j^{(2)}\|_F = \dots = \|\mathbf{W}_j^{(L)}\|_F, \forall j \in [M], \ell \in [L]. \quad (12)$$

584 Using (4) again, one can find that the neural network is also equivalent to

$$f(x) = \sum_{j=1}^M a_j \bar{\mathbf{W}}^{(L)} \sigma(\bar{\mathbf{W}}_j^{(L-1)} \sigma(\dots \sigma(\bar{\mathbf{W}}_j^{(1)} x + \bar{\mathbf{b}}_j^{(1)}) \dots) + \bar{\mathbf{b}}_j^{(L-1)}) + \bar{\mathbf{b}}_j^{(L)},$$

585 where

$$\|\bar{\mathbf{W}}_j^{(\ell)}\|_F \leq \beta^{(\ell)}, a_j = \frac{\prod_{\ell=1}^L \|\mathbf{W}_j^{(\ell)}\|_F}{\prod_{\ell=1}^L \beta^{(\ell)}} = \frac{\|\mathbf{W}_j^{(1)}\|_F^L}{\prod_{\ell=1}^L \beta^{(\ell)}} = \frac{(\sum_{\ell=1}^L \|\mathbf{W}_j^{(\ell)}\|_F^2 / L)^{L/2}}{\prod_{\ell=1}^L \beta^{(\ell)}}, \quad (13)$$

586 where the last two equality comes from the assumption (12). Choosing $\beta^{(\ell)} = c_1 \sqrt{w}$ expect $\ell = 1$
 587 where $\beta^{(1)} = c_1 \sqrt{d}$, and scaling $\bar{\mathbf{b}}^{(\ell)}$ accordingly and taking the constraint in (11) into (13) finishes
 588 the proof. \square

589 D.2 Covering Number of Parallel Neural Networks

590 **Theorem 4.** *The covering number of the model defined in (5) apart from the bias in the last layer*
 591 *satisfies*

$$\log \mathcal{N}(\mathcal{F}, \delta) \lesssim w^{2+2/(1-2/L)} L^2 \sqrt{d} P'^{\frac{1}{1-2/L}} \delta^{-\frac{2/L}{1-2/L}} \log(w P' / \delta).$$

592

The proof relies on the covering number of each subnetwork in a parallel neural network (Lemma 10), observing that $|f(x)| \leq 2^{L-1} w^{L-1} \sqrt{d}$ under the condition in Lemma 10, and then apply Lemma 5. We argue that our choice of condition on $\|\mathbf{b}^{(\ell)}\|_2$ in Lemma 10 is sufficient to analyzing the model apart from the bias in the last layer, because it guarantees that $\sqrt{w} \|\mathbf{W}^{(\ell)} \mathcal{A}_{\ell-1}(x)\|_2 \leq \|\mathbf{b}^{(\ell)}\|_2$. This leads to

$$\|\mathbf{W}^{(\ell)} \mathcal{A}_{\ell-1}(x)\|_\infty \leq \|\mathbf{W}^{(\ell)} \mathcal{A}_{\ell-1}(x)\|_2 \leq \sqrt{w} \|\mathbf{b}^{(\ell)}\|_2 \leq \|\mathbf{b}^{(\ell)}\|_\infty$$

593 If this condition is not met, $\mathbf{W}^{(\ell)} \mathcal{A}_{\ell-1}(x) + \mathbf{b}^{(\ell)}$ is either always positive or always negative
 594 for all feasible x along at least one dimension. If $(\mathbf{W}^{(\ell)} \mathcal{A}_{\ell-1}(x) + \mathbf{b}^{(\ell)})_i$ is always negative,

595 one can replace $b^{(\ell)}_i$ with $-\max_{\mathbf{x}} \|\mathbf{W}^{(\ell)} \mathcal{A}_{\ell-1}(\mathbf{x})\|_{\infty}$ without changing the output of this model
 596 for any feasible \mathbf{x} . If $(\mathbf{W}^{(\ell)} \mathcal{A}_{\ell-1}(\mathbf{x}) + b^{(\ell)})_i$ is always positive, one can replace $b^{(\ell)}_i$ with
 597 $\max_{\mathbf{x}} \|\mathbf{W}^{(\ell)} \mathcal{A}_{\ell-1}(\mathbf{x})\|_{\infty}$, and adjust the bias in the next layer such that the output of this model
 598 is not changed for any feasible \mathbf{x} . In either cases, one can replace the bias $b^{(\ell)}$ with another one with
 599 smaller norm while keeping the model equivalent except the bias in the last layer.

600 **Lemma 10.** Let $\mathcal{F} \subseteq \{f : \mathbb{R}^d \rightarrow \mathbb{R}\}$ denote the set of L -layer neural network (or a subnetwork in
 601 a parallel neural network) with width w in each hidden layer. It has the form

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{W}^{(L)} \sigma(\mathbf{W}^{(L-1)} \sigma(\dots \sigma(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}) \dots) + \mathbf{b}^{(L-1)}) + \mathbf{b}^{(L)}, \\ \mathbf{W}^{(1)} &\in \mathbb{R}^{w \times d}, \|\mathbf{W}^{(1)}\|_F \leq \sqrt{d}, \mathbf{b}^{(1)} \in \mathbb{R}^w, \|\mathbf{b}^{(1)}\|_2 \leq \sqrt{dw}, \\ \mathbf{W}^{(\ell)} &\in \mathbb{R}^{w \times w}, \|\mathbf{W}^{(\ell)}\|_F \leq \sqrt{w}, \mathbf{b}^{(\ell)} \in \mathbb{R}^w, \|\mathbf{b}^{(\ell)}\|_2 \leq 2^{\ell-1} w^{\ell-1} \sqrt{dw}, \quad \forall \ell = 2, \dots, L-1, \\ \mathbf{W}^{(L)} &\in \mathbb{R}^{1 \times w}, \|\mathbf{W}^{(L)}\|_F \leq \sqrt{w}, b^{(L)} = 0 \end{aligned} \tag{14}$$

602 and $\sigma(\cdot)$ is the ReLU activation function, the input satisfy $\|\mathbf{x}\|_2 \leq 1$, then the supremum norm
 603 δ -covering number of \mathcal{F} obeys

$$\log \mathcal{N}(\mathcal{F}, \delta) \leq c_7 L w^2 \log(1/\delta) + c_8$$

604 where c_7 is a constant depending only on d , and c_8 is a constant that depend on d, w and L .

605 *Proof.* First study two neural networks which differ by only one layer. Let g_{ℓ}, g'_{ℓ} be two neural net-
 606 works satisfying (14) with parameters $\mathbf{W}_1, \mathbf{b}_1, \dots, \mathbf{W}_L, \mathbf{b}_L$ and $\mathbf{W}'_1, \mathbf{b}'_1, \dots, \mathbf{W}'_L, \mathbf{b}'_L$ respectively.
 607 Furthermore, the parameters in these two models are the same except the ℓ -th layer, which satisfy

$$\|\mathbf{W}_{\ell} - \mathbf{W}'_{\ell}\|_F \leq \epsilon, \|\mathbf{b}_{\ell} - \mathbf{b}'_{\ell}\|_2 \leq \tilde{\epsilon}.$$

608 Denote the model as

$$g_{\ell}(\mathbf{x}) = \mathcal{B}_{\ell}(\mathbf{W}_{\ell} \mathcal{A}_{\ell}(\mathbf{x}) + \mathbf{b}_{\ell}), g'_{\ell}(\mathbf{x}) = \mathcal{B}_{\ell}(\mathbf{W}'_{\ell} \mathcal{A}_{\ell}(\mathbf{x}) + \mathbf{b}'_{\ell})$$

609 where $\mathcal{A}_{\ell}(\mathbf{x}) = \sigma(\mathbf{W}_{\ell-1} \sigma(\dots \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \dots) + \mathbf{b}_{\ell-1})$ denotes the first $\ell-1$ layers in the neural
 610 network, and $\mathcal{A}_{\ell}(\mathbf{x}) = \mathbf{W}_L \sigma(\dots \sigma(\mathbf{W}_{\ell+1} \sigma(\mathbf{x}) + \mathbf{b}_{\ell+1}) \dots) + \mathbf{b}_L$ denotes the last $L-\ell-1$ layers,
 611 with definition $\mathcal{A}_1(\mathbf{x}) = \mathbf{x}, \mathcal{B}_L(\mathbf{x}) = \mathbf{x}$.

612 Now focus on bounding $\|\mathcal{A}_{\ell}(\mathbf{x})\|$. Let $\mathbf{W} \in \mathbb{R}^{m \times m'}, \|\mathbf{W}\|_F \leq \sqrt{m'}, \mathbf{x} \in \mathbb{R}^{m'}, \mathbf{b} \in \mathbb{R}^m, \|\mathbf{b}\|_2 \leq$
 613 \sqrt{m}

$$\begin{aligned} \|\sigma(\mathbf{W}\mathbf{x} + \mathbf{b})\|_2 &\leq \|\mathbf{W}\mathbf{x} + \mathbf{b}\|_2 \\ &\leq \|\mathbf{W}\|_2 \|\mathbf{x}\|_2 + \|\mathbf{b}\|_2 \\ &\leq \|\mathbf{W}\|_F \|\mathbf{x}\|_2 + \|\mathbf{b}\|_2 \\ &\leq \sqrt{m'} \|\mathbf{x}\|_2 + \sqrt{m} \end{aligned}$$

614 where we make use of $\|\cdot\|_2 \leq \|\cdot\|_F$. Because of that,

$$\begin{aligned} \|\mathcal{A}_2(\mathbf{x})\|_2 &\leq \sqrt{d} + \sqrt{dw} \leq 2\sqrt{dw}, \\ \|\mathcal{A}_3(\mathbf{x})\|_2 &\leq \sqrt{w} \|\mathcal{A}_2(\mathbf{x})\|_2 + 2w\sqrt{dw} \leq 4w\sqrt{dw}, \\ &\dots \\ \|\mathcal{A}_{\ell}(\mathbf{x})\|_2 &\leq \sqrt{w} \|\mathcal{A}_{\ell-1}(\mathbf{x})\|_2 \leq 2\sqrt{dw} (2w)^{\ell-2}. \end{aligned} \tag{15}$$

615 Then focus on $\mathcal{B}(\mathbf{x})$. Let $\mathbf{W} \in \mathbb{R}^{m \times m'}, \|\mathbf{W}\|_F \leq \sqrt{m'}, \mathbf{x}, \mathbf{x}' \in \mathbb{R}^{m'}, \mathbf{b} \in \mathbb{R}^m, \|\mathbf{b}\|_2 \leq \sqrt{m}$.
 616 Furthermore, $\|\mathbf{x} - \mathbf{x}'\|_2 \leq \epsilon$, then

$$\|\sigma(\mathbf{W}\mathbf{x} + \mathbf{b}) - \sigma(\mathbf{W}\mathbf{x}' + \mathbf{b})\|_2 \leq \|\mathbf{W}(\mathbf{x} - \mathbf{x}')\|_2 \leq \|\mathbf{W}\|_F \|\mathbf{x} - \mathbf{x}'\|_2$$

617 which indicates that $\|\mathcal{B}(\mathbf{x}) - \mathcal{B}(\mathbf{x}')\|_2 \leq (\sqrt{w})^{L-\ell} \|\mathbf{x} - \mathbf{x}'\|_2$

618 Finally, for any $\mathbf{W}, \mathbf{W}' \in \mathbb{R}^{m \times m'}, \mathbf{x} \in \mathbb{R}^{m'}, \mathbf{b}, \mathbf{b}' \in \mathbb{R}^m$, one have

$$\begin{aligned} \|(\mathbf{W}\mathbf{x} + \mathbf{b}) - (\mathbf{W}'\mathbf{x} + \mathbf{b}')\|_2 &= \|(\mathbf{W} - \mathbf{W}')\mathbf{x} + (\mathbf{b} - \mathbf{b}')\|_2 \\ &\leq \|\mathbf{W} - \mathbf{W}'\|_2 \|\mathbf{x}\|_2 + \|\mathbf{b} - \mathbf{b}'\|_2 \\ &\leq \|\mathbf{W} - \mathbf{W}'\|_F \|\mathbf{x}\|_2 + \sqrt{m} \|\mathbf{b} - \mathbf{b}'\|_{\infty}. \end{aligned}$$

619 In summary,

$$\begin{aligned}
|g_\ell(\mathbf{x}) - g'_\ell(\mathbf{x})| &= |\mathcal{B}_\ell(\mathbf{W}_\ell \mathcal{A}_\ell(\mathbf{x}) + \mathbf{b}_\ell) - \mathcal{B}_\ell(\mathbf{W}'_\ell \mathcal{A}_\ell(\mathbf{x}) + \mathbf{b}'_\ell)| \\
&\leq (\sqrt{w})^{L-\ell} \|(\mathbf{W}_\ell \mathcal{A}_\ell(\mathbf{x}) + \mathbf{b}_\ell) - (\mathbf{W}'_\ell \mathcal{A}_\ell(\mathbf{x}) + \mathbf{b}'_\ell)\|_2 \\
&\leq (\sqrt{w})^{L-\ell} (\|\mathbf{W}_\ell - \mathbf{W}'_\ell\|_F \|\mathcal{A}_\ell(\mathbf{x})\|_2 + \|\mathbf{b}_\ell - \mathbf{b}'_\ell\|_2) \\
&\leq 2^{(\ell-1)} w^{(L+\ell-3)/2} d^{1/2} \epsilon + w^{(L-\ell)/2} \bar{\epsilon}
\end{aligned}$$

620 Let $f(x), f'(x)$ be two neural networks satisfying (14) with parameters $W_1, b_1, \dots, W_L, b_L$ and
621 $W'_1, b'_1, \dots, W'_L, b'_L$ respectively, and $\|W_\ell - W'_\ell\|_F \leq \epsilon_\ell, \|b_\ell - b'_\ell\|_F \leq \bar{\epsilon}_\ell$. Further define f_ℓ be the
622 neural network with parameters $W_1, b_1, \dots, W_\ell, b_\ell, W'_{\ell+1}, b'_{\ell+1}, \dots, W'_L, b'_L$, then

$$\begin{aligned}
|f(x) - f'(x)| &\leq |f(x) - f_1(x)| + |f_1(x) - f_2(x)| + \dots + |f_{L-1}(x) - f'(x)| \\
&\leq \sum_{\ell=1}^L 2^{(\ell-2)} d^{1/2} w^{(L+\ell-3)/2} \epsilon + w^{(L-\ell)/2} \bar{\epsilon}
\end{aligned}$$

For any $\delta > 0$, one can choose

$$\epsilon_\ell = \frac{\delta}{2^\ell w^{(L+\ell-3)/2} d^{1/2}}, \bar{\epsilon}_\ell = \frac{\delta}{2 w^{(L-\ell)/2}}$$

623 such that $|f(x) - f'(x)| \leq \delta$.

624 On the other hand, the ϵ -covering number of $\{\mathbf{W} \in \mathbb{R}^{m \times m'} : \|\mathbf{W}\|_F \leq \sqrt{m'}\}$ on Frobenius norm
625 is no larger than $(2\sqrt{m'}/\epsilon + 1)^{m \times m'}$, and the $\bar{\epsilon}$ -covering number of $\{\mathbf{b} \in \mathbb{R}^m : \|\mathbf{b}\|_2 \leq 1\}$ on
626 infinity norm is no larger than $(2/\bar{\epsilon} + 1)^m$. The entropy of this neural network can be bounded by

$$\log \mathcal{N}(f; \delta) \leq w^2 L \log(2^{L+1} w^{L-1} / \delta + 1) + wL \log(2^{L-1} w^{(L-1)/2} d^{1/2} / \delta + 1)$$

627

□

628 D.3 Covering Number of p -Norm Constrained Linear Combination

629 **Lemma 5.** $\log \mathcal{N}(\mathcal{G}, \delta) \lesssim k \log(1/\delta)$ for some finite c_3 , and for any $g \in \mathcal{G}, |a| \leq 1$, we have
630 $ag \in \mathcal{G}$. The covering number of $\mathcal{F} = \left\{ \sum_{i=1}^M a_i g_i \mid g_i \in \mathcal{G}, \|a\|_p^p \leq P, 0 < p < 1 \right\}$ for any $P > 0$
631 satisfies

$$\log \mathcal{N}(\mathcal{F}, \epsilon) \lesssim k P^{\frac{1}{1-p}} (\delta/c_3)^{-\frac{p}{1-p}} \log(c_3 P/\delta)$$

632 up to a double logarithmic factor.

633 *Proof.* Let ϵ be a positive constant. Without the loss of generality, we can sort the coefficients in
634 descending order in terms of their absolute values. There exists a positive integer \mathcal{M} (as a function
635 of ϵ), such that $|a_i| \geq \epsilon$ for $i \leq \mathcal{M}$, and $|a_i| < \epsilon$ for $i > \mathcal{M}$.

636 By definition, $\mathcal{M}\epsilon^p \leq \sum_{i=1}^{\mathcal{M}} |a_i|^p \leq P$ so $\mathcal{M} \leq P/\epsilon^p$, and $|a_i|^p \leq P, |a_i| \leq P^{1/p}$ for all i .
637 Furthermore,

$$\sum_{i>\mathcal{M}} |a_i| = \sum_{i>\mathcal{M}} |a_i|^p |a_i|^{1-p} < \sum_{i>\mathcal{M}} |a_i|^p \epsilon^{1-p} \leq P \epsilon^{1-p}$$

638 Let $\tilde{g}_i = \arg \min_{g \in \tilde{\mathcal{G}}} \|g - \frac{a_i}{P^{1/p}} g_i\|_\infty$ where $\tilde{\mathcal{G}}$ is the δ' -covering set of \mathcal{G} . By definition of the
639 covering set,

$$\begin{aligned}
\left\| \sum_{i=1}^{\mathcal{M}} a_i g_i(x) - \sum_{i=1}^{\mathcal{M}} P^{1/p} \tilde{g}_i(x) \right\|_\infty &\leq \left\| \sum_{i=1}^{\mathcal{M}} (a_i g_i(x) - P^{1/p} \tilde{g}_i(x)) \right\|_\infty + \left\| \sum_{i=\mathcal{M}+1}^{\mathcal{M}} a_i g_i(x) \right\|_\infty \\
&\leq \mathcal{M} P^{1/p} \delta' + c_3 P \epsilon^{1-p}.
\end{aligned} \tag{16}$$

640 Choosing

$$\epsilon = (\delta/2c_3 P)^{\frac{1}{1-p}}, \delta' \approx P^{-\frac{1}{p(1-p)}} (\delta/2c_3)^{\frac{1}{1-p}} / 2, \tag{17}$$

641 we have $\mathcal{M} \leq P^{\frac{1}{1-p}} (\delta/2c_3)^{-\frac{p}{1-p}}, \mathcal{M} P^{1/p} \delta' \leq \delta/2, c_3 P \epsilon^{1-p} \leq \delta/2$, so (16) $\leq \delta$. One can
642 compute the covering number of \mathcal{F} by

$$\log \mathcal{N}(\mathcal{F}, \delta) \leq \mathcal{M} \log \mathcal{N}(\mathcal{G}, \delta') \lesssim k \mathcal{M} \log(1/\delta') \tag{18}$$

643 Taking (17) into (18) finishes the proof. □

644 E Proof of Approximation Error

645 E.1 Approximation of Neural Networks to B-spline Basis Functions

646 **Proposition 6.** *There exists a parallel neural network that has the structure and satisfy the constraint*
 647 *in Proposition 3 for d -dimensional input and one output, containing $M = O(m^d)$ subnetworks,*
 648 *each of which has width $w = O(d)$ and depth $L = O(\log(c(m, d)/\epsilon))$ for some constant w, c that*
 649 *depends only on m and d , denoted as $\tilde{M}_m(\mathbf{x}), \mathbf{x} \in \mathbb{R}^d$, such that*

- 650 • $|\tilde{M}_{m,k,s}(\mathbf{x}) - M_{m,k,s}(\mathbf{x})| \leq \epsilon$, if $0 \leq 2^k(x_i - s_i) \leq m + 1, \forall i \in [d]$,
- 651 • $\tilde{M}_{m,k,s}(\mathbf{x}) = 0$, otherwise.
- 652 • The weights in the last layer satisfy $\|a\|_{2/L}^{2/L} \lesssim 2^k m^d e^{2md/L}$.

653 We follow the method developed in Yarotsky [46], Suzuki [36], while putting our attention on bound-
 654 ing the Frobenius norm of the weights.

655 **Lemma 11** (Yarotsky [46, Proposition 3]). *There exists a neural network with two-dimensional*
 656 *input and one output $f_\times(x, y)$, with constant width and depth $O(\log(1/\delta))$, and the weight in each*
 657 *layer is bounded by a global constant c_1 , such that*

- 658 • $|f_\times(x, y) - xy| \leq \delta, \forall 0 \leq x, y \leq 1$,
- 659 • $f_\times(x, y) = 0, \forall x = 0$ or $y = 0$.

660 We first prove a special case of Proposition 6 on the unscaled, unshifted B-spline basis function by
 661 fixing $k = 0, \mathbf{s} = 0$:

662 **Proposition 12.** *There exists a parallel neural network that has the structure and satisfy the con-*
 663 *straint in Proposition 3 for d -dimensional input and one output, containing $M = \lceil (m + 1)/2 \rceil^d =$
 664 $O(m^d)$ subnetworks, each of which has width $w = O(d)$ and depth $L = O(\log(c(m, d)/\epsilon))$ for
 665 some constant w, c that depends only on m and d , denoted as $\tilde{M}_m(\mathbf{x}), \mathbf{x} \in \mathbb{R}^d$, such that*

- 666 • $|\tilde{M}_m(\mathbf{x}) - M_m(\mathbf{x})| \leq \epsilon$, if $0 \leq x_i \leq m + 1, \forall i \in [d]$, while $M_m(\cdot)$ denote m -th order
 667 B-spline basis function, and c only depends on m and d .
- 668 • $\tilde{M}_m(\mathbf{x}) = 0$, if $x_i \leq 0$ or $x_i \geq m + 1$ for any $i \in [d]$.
- 669 • The weights in the last layer satisfy $\|a\|_{2/L}^{2/L} \lesssim m^d e^{2md/L}$.

670 *Proof.* We first show that one can use a neural network with constant width w_0 , depth $L \approx$
 671 $\log(m/\epsilon_1)$ and bounded norm $\|W^{(1)}\|_F \leq O(\sqrt{d}), \|W^{(\ell)}\|_F \leq O(\sqrt{w}), \forall \ell = 2, \dots, L$ to
 672 approximate truncated power basis function up to accuracy ϵ_1 in the range $[0, 1]$. Let $m =$
 673 $\sum_{i=0}^{\lceil \log_2 m \rceil} m_i 2^i, m_i \in \{0, 1\}$ be the binary digits of m , and define $\bar{m}_j = \sum_{i=0}^j m_i, \gamma = \lceil \log_2 m \rceil$,
 674 then for any x

$$\begin{aligned}
 x_+^m &= x_+^{\bar{m}_\gamma} \times (x_+^{2^\gamma})^{m_\gamma} \\
 [x_+^{\bar{m}_\gamma}, x_+^{2^\gamma}] &= [x_+^{\bar{m}_{\gamma-1}} \times (x_+^{2^{\gamma-1}})^{m_{\gamma-1}}, x_+^{2^{\gamma-1}} \times x_+^{2^{\gamma-1}}] \\
 &\dots \\
 [x_+^{\bar{m}_2}, x_+^4] &= [x_+^{\bar{m}_1} \times (x_+^2)^{m_1}, x_+^2 \times x_+^2] \\
 [x_+^{\bar{m}_1}, x_+^2] &= [x_+^{\bar{m}_0} \times x_+^{m_0}, x_+ \times x_+]
 \end{aligned} \tag{19}$$

675 Notice that each line of equation only depends on the line immediately below. Replacing the
 676 multiply operator \times with the neural network approximation shown in Lemma 11 demonstrates the
 677 architecture of such neural network approximation. For any $x, y \in [0, 1]$, let $|f_\times(x, y) - xy| \leq$
 678 $\delta, |x - \tilde{x}| \leq \delta_1, |y - \tilde{y}| \leq \delta_2$, then $|f_\times(\tilde{x}, \tilde{y}) - xy| \leq \delta_1 + \delta_2 + \delta$. Taking this into (19) shows that
 679 $\epsilon_1 \approx 2^\gamma \delta \approx m\delta$, where ϵ_1 is the upper bound on the approximate error to truncated power basis of
 680 order m and δ is the approximation error to a single multiply operator as in Lemma 11.

681 A univariate B-spline basis can be expressed using truncated power basis, and observing that it is
 682 symmetric around $(m + 1)/2$:

$$\begin{aligned} M_m(x) &= \frac{1}{m!} \sum_{j=1}^{m+1} (-1)^j \binom{m+1}{j} (x-j)_+^m \\ &= \frac{1}{m!} \sum_{j=1}^{\lceil (m+1)/2 \rceil} (-1)^j \binom{m+1}{j} (\min(x, m+1-x) - j)_+^m \\ &= \frac{((m+1)/2)^m}{m!} \sum_{j=1}^{\lceil (m+1)/2 \rceil} (-1)^j \binom{m+1}{j} \left(\frac{\min(x, m+1-x) - j}{(m+1)/2} \right)_+^m, \end{aligned}$$

683 A multivariate (d -dimensional) B-spline basis function can be expressed as the product of truncated
 684 power basis functions and thus can be decomposed as

$$\begin{aligned} M_m(\mathbf{x}) &= \prod_{i=1}^d M_m(x_i) \\ &= \frac{((m+1)/2)^{md}}{(m!)^d} \prod_{i=1}^d \left(\sum_{j=1}^{\lceil (m+1)/2 \rceil} (-1)^j \binom{m+1}{j} \left(\frac{\min(x_i, m+1-x) - j}{(m+1)/2} \right)_+^m \right) \quad (20) \\ &= \frac{((m+1)/2)^{md}}{(m!)^d} \sum_{j_1, \dots, j_d=1}^{\lceil (m+1)/2 \rceil} \prod_{i=1}^d (-1)^{j_i} \binom{m+1}{j_i} \left(\frac{\min(x, m+1-x) - j_i}{(m+1)/2} \right)_+^m \end{aligned}$$

685 Using Lemma 11, one can construct a parallel neural network containing $M = \lceil (m+1)/2 \rceil^d =$
 686 $O(m^d)$ subnetworks, and each subnetwork corresponds to one polynomial term in (20). Using the
 687 results above, the approximation of this constructed neural network can be bounded by

$$\left(\sum_{i=1}^{m+1} \binom{m+1}{j} d(\epsilon_1 + \delta) \right)^d \lesssim \frac{e^{2m}}{\sqrt{m}} d\epsilon_1 + d\delta$$

688 where we applied Stirling's approximation and δ and ϵ_1 has the same definition as above. Choosing
 689 $\delta = \frac{\epsilon}{d(e^{2m}\sqrt{m+1})}$, and recall $\epsilon_1 \approx m\delta$ proves the approximation error.

690 To bound the norm of the factors $\|a\|_{2/L}^{2/L}$, first observe that

$$\begin{aligned} |a_{j_1, \dots, j_d}| &= \frac{((m+1)/2)^{md}}{(m!)^d} \frac{1}{(m+1)/2} \prod_{i=1}^d \binom{m+1}{j_i} \\ &\leq \frac{((m+1)/2)^{md}}{(m!)^d} \frac{2^{md}}{(m+1)/2} = O(e^{md}) \end{aligned}$$

691 where the first inequality is from $\binom{m+1}{j_i} \leq 2^{m+1}$, the last equality is from Stirling's approximation.
 692 Finally,

$$\|a\|_{2/L}^{2/L} \leq m^d \max_j |a_j|^{2/L} \lesssim m^d e^{2md/L}$$

693 which finishes the proof. □

694 The proof of the Proposition 6 for general k, s follows by appending one more layer in the front, as
 695 we show below.

696 *Proof of Proposition 6.* Using the neural network proposed in Proposition 12, one can construct a
 697 neural network for approximating $M_{m,k,s}$ by adding one layer before the first layer:

$$\sigma(2^k \mathbf{I}_d \mathbf{x} - 2^k \mathbf{s})$$

698 The unused neurons in the first hidden layer is zero padded. The Frobenius norm of the weight is
699 $2^k \|\mathbf{I}_d\|_F = 2^k \sqrt{d}$. Following the proof of Proposition 3, rescaling the weight in this layer by 2^{-k} ,
700 and the weight matrix in the last layer by 2^k , and scaling the bias properly, one can verify that this
701 neural network satisfy the statement. \square

702 E.2 Sparse approximation of Besov functions using B-spline wavelets

Proposition 7. *Let $\alpha - d/p > 1, r > 0$. Let $M_{m,k,s}$ be the B-spline of order m with scale 2^{-k} in each dimension and position $\mathbf{s} \in \mathbb{R}^d$. For any function in Besov space $f_0 \in B_{p,q}^\alpha$ and any positive integer \bar{M} , there is an \bar{M} -sparse approximation using B-spline basis of order m satisfying $0 < \alpha < \min(m, m - 1 + 1/p)$: $\tilde{f}_{\bar{M}} = \sum_{i=1}^{\bar{M}} a_{k_i, \mathbf{s}_i} M_{m, k_i, \mathbf{s}_i}$ for any positive integer \bar{M} such that the approximation error is bounded as $\|\tilde{f}_{\bar{M}} - f_0\|_r \lesssim \bar{M}^{-\alpha/d} \|f_0\|_{B_{p,q}^\alpha}$, and the coefficients satisfy*

$$\|\{2^{k_i} a_{k_i, \mathbf{s}_i}\}_{k_i, \mathbf{s}_i}\|_p \lesssim \|f_0\|_{B_{p,q}^\alpha}.$$

703

704 The proof is divided into three steps:

- 705 1. Bound the 0-norm and the 1-norm of the coefficients of B-spline basis in order to approxi-
706 mate an arbitrary function in Besov space up to any $\epsilon > 0$.
- 707 2. Bound p -norm of the coefficients of B-spline basis functions where $0 < p < 1$ using the
708 results above .
- 709 3. Add the approximation to neural network to B-spline basis computed in Section 4.3.1 into
710 Step 2.

711 *Proof.* Dũng [11, Theorem 3.1] Suzuki [36, Lemma 2] proposed an adaptive sampling recovery
712 method that approximates a function in Besov space. The method is divided into two cases: when
713 $p \geq r$, and when $p < r$.

714 When $p \geq r$, there exists a sequence of scalars $\lambda_j, \mathbf{j} \in P^d(\mu), P_d(\mu) := \{\mathbf{j} \in \mathbb{Z}^d : |j_i| \leq \mu, \forall i \in$
715 $d\}$ for some positive μ , for arbitrary positive integer \bar{k} , the linear operator

$$Q_{\bar{k}}(f, \mathbf{x}) = \sum_{\mathbf{s} \in J(\bar{k}, m, d)} a_{\bar{k}, \mathbf{s}}(f) M_{\bar{k}, \mathbf{s}}(\mathbf{x}), \quad a_{\bar{k}, \mathbf{s}}(f) = \sum_{\mathbf{j} \in \mathbb{Z}^d, P^d(\mu)} \lambda_j \bar{f}(\mathbf{s} + 2^{-\bar{k}} \mathbf{j})$$

716 has bounded approximation error

$$\|f - Q_{\bar{k}}(f, x)\|_r \leq C 2^{-\alpha \bar{k}} \|f\|_{B_{p,q}^\alpha},$$

717 where \bar{f} is the extrapolation of f , $J(\bar{k}, m, d) := \{\mathbf{s} : 2^{\bar{k}} \mathbf{s} \in \mathbb{Z}^d, -m/2 \leq 2^{\bar{k}} s_i \leq 2^{\bar{k}} + m/2, \forall i \in$
718 $[d]\}$. See Dũng [11, 2.6-2.7] for the detail of the extrapolation as well as references for options of
719 sequence λ_j .

720 Furthermore, $Q_{\bar{k}}(f) \in B_{p,q}^\alpha$ so it can be decomposed in the form (10) with $M = \sum_{k=0}^{\bar{k}} (2^k + m -$
721 $1)^d \lesssim 2^{\bar{k}d}$ components and $\|\{\tilde{c}_{k,s}\}_{k,s}\| \lesssim \|Q_{\bar{k}}(f)\|_{B_{p,q}^\alpha} \lesssim \|f\|_{B_{p,q}^\alpha}$ where $\tilde{c}_{k,s}$ is the coefficients of
722 the decomposition of $Q_{\bar{k}}(f)$. Choosing $\bar{k} \approx \log_2 M/d$ leads to the desired approximation error.

723 On the other hand, when $p < r$, there exists a greedy algorithm that constructs

$$G(f) = Q_{\bar{k}}(f) + \sum_{k=\bar{k}+1}^{k^*} \sum_{j=1}^{n_k} c_{k, \mathbf{s}_j}(f) M_{k, \mathbf{s}_j}$$

724 where $\bar{k} \approx \log_2(M), k^* = \lceil \epsilon^{-1} \log(\lambda M) \rceil + \bar{k} + 1, n_k = \lfloor \lambda M 2^{-\epsilon(k-\bar{k})} \rfloor$ for some $0 < \epsilon <$
725 $\alpha/\delta - 1, \delta = d(1/p - 1/r), \lambda > 0$, such that

$$\|f - G(f)\|_r \leq \bar{M}^{-\alpha/d} \|f\|_{B_{p,q}^\alpha}$$

726 and

$$\sum_{k=0}^{\bar{k}} (2^k + m - 1)^d + \sum_{k=\bar{k}+1}^{k^*} n_k \leq \bar{M}.$$

727 See Dũng [11, Theorem 3.1] for the detail.

728 Finally, since $\alpha - d/p > 1$,

$$\begin{aligned}
\|\{2^{k_i} c_{k_i, \mathbf{s}_i}\}_{k_i, \mathbf{s}_i}\|_p &\leq \sum_{k=0}^{\bar{k}} 2^k \|\{c_{k_i, \mathbf{s}_i}\}_{\mathbf{s}_i}\|_p \\
&= \sum_{k=0}^{\bar{k}} 2^{(1-(\alpha-d/p)k)} (2^{(\alpha-d/p)k} \|\{c_{k_i, \mathbf{s}_i}\}_{\mathbf{s}_i}\|_p) \\
&\lesssim \sum_{k=0}^{\bar{k}} 2^{(1-(\alpha-d/p)k)} \|f\|_{B_{p,q}^\alpha} \\
&\approx \|f\|_{B_{p,q}^\alpha}
\end{aligned} \tag{21}$$

729 where the first line is because for arbitrary vectors $\mathbf{a}_i, i \in [n]$, $\|\sum_{i=1}^n \mathbf{a}_i\|_p \leq \sum_{i=1}^n \|\mathbf{a}_i\|_p$, the
730 third line is because the sequence norm of B-spline decomposition is equivalent to the norm in
731 Besov space (see Section C.1). \square

732 Note that when $\alpha - d/p = 1$, the sequence norm (21) is bounded (up to a factor of constant) by
733 $k^* \|f\|_{B_{p,q}^\alpha}$, which can be proven by following (21) except the last line. This adds a logarithmic term
734 with respect to \bar{M} compared with the result in Proposition 7. This will add a logarithmic factor to
735 the MSE. We will not focus on this case in this paper of simplicity.

736 E.3 Sparse approximation of Besov functions using Parallel Neural Networks

737 **Theorem 8.** *Under the same condition as Proposition 7, for any positive integer \bar{M} , any function*
738 *in Besov space $f_0 \in B_{p,q}^\alpha$ can be approximated by a parallel neural network with no less than*
739 *$O(m^d \bar{M})$ number of subnetworks satisfying:*

- 740 1. Each subnetwork has width $w = O(d)$ and depth L .
- 741 2. The weights in each layer satisfy $\|\bar{\mathbf{W}}_k^{(\ell)}\|_F \leq O(\sqrt{w})$ except the first layer $\|\bar{\mathbf{W}}_k^{(1)}\|_F \leq$
742 $O(\sqrt{d})$,
- 743 3. The scaling factors have bounded $2/L$ -norm: $\|\{a_j\}\|_{2/L}^{2/L} \lesssim m^d e^{2md/L} \bar{M}^{1-2/(pL)}$.
- 744 4. The approximation error is bounded by

$$\|\tilde{f} - f_0\|_r \leq (c_4 \bar{M}^{-\alpha/d} + c_5 e^{-c_6 L}) \|f\|_{B_{p,q}^\alpha}$$

745 where c_4, c_5, c_6 are constants that depend only on m, d and p .

746 We first prove the following lemma.

Lemma 13. *For any $a \in \mathbb{R}^{\bar{M}}$, $0 < p' < p$, it holds that:*

$$\|a\|_{p'}^{p'} \leq \bar{M}^{1-p'/p} \|a\|_p^{p'}.$$

Proof.

$$\sum_i |a_i|^{p'} = \langle \mathbf{1}, |a|^{p'} \rangle \leq \left(\sum_i 1 \right)^{1-\frac{p'}{p}} \left(\sum_i (|a_i|^{p'})^{\frac{p}{p'}} \right)^{\frac{p'}{p}} = \bar{M}^{1-\frac{p'}{p}} \|a\|_p^{p'}$$

747 The first inequality uses a Holder's inequality with conjugate pair $\frac{p}{p'}$ and $1/(1 - \frac{p'}{p})$. \square

748 *Proof of Theorem 8.* Using Proposition 7, one can construct \bar{M} number of PNN each $O(m^d)$ sub-
749 networks according to Proposition 6, and in each PNN, such that each PNN represents one B-spline
750 basis function. The weights in the last layer of each PNN is scaled to match the coefficients in Propo-
751 sition 7. Taking p' in Lemma 13 as $2/L$ and combining with Proposition 6 finishes the proof. \square

752 **F Proof of the Main Theorem**

753 **Theorem 1 extended form.** For any fixed $\alpha - d/p > 1, r > 0, L \geq 3$, given an L -layer parallel
754 neural network satisfying

- 755 • The width of each subnetwork is fixed and large enough: $w \gtrsim d$. See Theorem 8 for the
756 detail.
- 757 • The number of subnetworks is large enough: $M \gtrsim m^d n^{\frac{1-2/L}{2\alpha/d+1-2/(pL)}}$.

758 With proper choice of the parameter of weight decay λ , the solution \hat{f} parameterized by (2) satisfies

$$\text{MSE}(\hat{f}) = \tilde{O}\left(\left(\frac{w^{4-4/L} L^{2-4/L}}{n^{1-2/L}}\right)^{\frac{2\alpha/d}{2\alpha/d+1-2/(pL)}} + e^{-c_6 L}\right)$$

759 where \tilde{O} shows the scale up to a logarithmic factor, and c_6 is the constant defined in Theorem 8.

760 *Proof.* First recall the relationship between covering number (entropy) and estimation error:

761 **Proposition 14.** Let $\mathcal{F} \subseteq \{\mathbb{R}^d \rightarrow [-F, F]\}$ be a set of functions. Assume that \mathcal{F} can be decomposed
762 into two orthogonal spaces $\mathcal{F} = \mathcal{F}_{\parallel} \times \mathcal{F}_{\perp}$ where \mathcal{F}_{\perp} is an affine space with dimension of N . Let
763 $f_0 \in \{\mathbb{R}^d \rightarrow [-F, F]\}$ be the target function and \hat{f} be the least squares estimator in \mathcal{F} :

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n (y_i - f(x_i))^2, y_i = f_0(x_i) + \epsilon_i, \epsilon_i \sim \mathcal{N}(0, \sigma^2) i.i.d.,$$

764 then it holds that

$$\text{MSE}(\hat{f}) \leq \tilde{O}\left(\arg \min_{f \in \mathcal{F}} \text{MSE}(f) + \frac{N + \log \mathcal{N}(\mathcal{F}_{\parallel}, \delta) + 2}{n} + (F + \sigma)\delta\right).$$

765 The proof of Proposition 14 is deferred to the section below. We choose \mathcal{F} as the set of functions
766 that can be represented by a parallel neural network as stated, the (null) space $\mathcal{F}_{\perp} = \{f : f(\mathbf{x}) =$
767 $\text{constant}\}$ be the set of functions with constant output, which has dimension 1. This space captures
768 the bias in the last layer, while the other parameters contributes to the projection in \mathcal{F}_{\parallel} . See Section
769 D.2 for how we handle the bias in the other layers. One can find that \mathcal{F}_{\parallel} is the set of functions
770 that can be represented by a parallel neural network as stated, and further satisfy $\sum_{i=1}^n f(\mathbf{x}_i) = 0$.
771 Because $\mathcal{F}_{\parallel} \subseteq \mathcal{F}$, $\mathcal{N}(\mathcal{F}_{\parallel}, \delta) \leq \mathcal{N}(\mathcal{F}, \delta)$ for all $\delta > 0$, and the latter is studied in Theorem 4.

772 In Theorem 1, the width of each subnetwork is no less than what is required in Theorem 8, while the
773 depth and norm constraint are the same, so the approximation error is no more than that in Theorem 8.
774 Choosing $r = 2, p = 2/L$, and taking Theorem 4 and Theorem 8 into this Proposition 14, one gets

$$\text{MSE}(\hat{f}) \lesssim \bar{M}^{-2\alpha/d} + \frac{w^{2+2/(1-2/L)} L^2}{n} \bar{M}^{\frac{1-2/(pL)}{1-2/L}} \delta^{-\frac{2/L}{1-2/L}} (\log(\bar{M}/\delta) + 3) + \delta,$$

where $\|f\|_{B_{p,q}^{\alpha}}$, m and d taken as constants. The stated MSE is obtained by choosing

$$\delta \approx \frac{w^{4-4/L} L^{2-4/L} \bar{M}^{1-2/(pL)}}{n^{1-2/L}}, \bar{M} \approx \left(\frac{n^{1-2/L}}{w^{4-4/L} L^{2-4/L}}\right)^{\frac{1}{2\alpha/d+1-2/(pL)}}$$

775 Note that there exists a weight decay parameter λ' such that the $(2/L)$ -norm of the coefficients
776 of the parallel neural network satisfy that $\|\{a_j\}\|_{2/L}^{2/L} = m^d e^{2md/L} \|\{\tilde{a}_{j,\bar{M}}\}\|_{2/L}^{2/L}$ where $\{\tilde{a}_{j,\bar{M}}\}$
777 is the coefficient of the particular \bar{M} -sparse approximation, although $\{a_j\}$ is not necessarily \bar{M}
778 sparse. Empirically, one only need to guarantee that during initialization, the number of subnetworks
779 $M \geq \bar{M}$ such that the \bar{M} -sparse approximation is feasible, thus the approximation error bound
780 from Theorem 8 can be applied. Theorem 8 also says that $\|\{a_j\}\|_{2/L}^{2/L} = m^d e^{2md/L} \|\{\tilde{a}_{j,\bar{M}}\}\|_{2/L}^{2/L} \lesssim$
781 $\bar{M}^{1-2/pL}$, thus we can apply the covering number bound from Theorem 4 with $P' = \bar{M}^{1-2/pL}$.
782 Finally, if λ is optimally chosen, then it achieves a smaller MSE than this particular λ' , which has
783 been proven to be no more than $O(\bar{M}^{-\alpha/d})$ and completes the proof. □

784

785 *Proof of Proposition 14.* For any function $f \in \mathcal{F}$, define $f_\perp = \arg \min_{h \in \mathcal{F}_\perp} \sum_{i=1}^n (f(\mathbf{x}_i) -$
786 $h(\mathbf{x}_i))^2$ be the projection of f to \mathcal{F}_\perp , and define $f_\parallel = f - f_\perp$ be the projection to the orthogo-
787 nal complement. Note that f_\parallel is not necessarily in \mathcal{F}_\parallel . However, if $f \in \mathcal{F}$, then $f_\parallel \in \mathcal{F}_\parallel$. $y_{i\perp}$ and
788 $y_{i\parallel}$ are defined by creating a function f_y such that $f_y(\mathbf{x}_i) = y_i, \forall i$, e.g. via interpolation. Because
789 \mathcal{F}_\parallel and \mathcal{F}_\perp are orthononal, the empirical loss and population loss can be decomposed in the same
790 way:

$$\begin{aligned} L_\parallel(f) &= \frac{1}{n} \sum_{i=1}^n (f_\parallel(\mathbf{x}) - f_{0\parallel}(\mathbf{x}))^2 + \frac{n-N}{n} \sigma^2, & L_\perp(f) &= \frac{1}{n} \sum_{i=1}^n (f_\perp(\mathbf{x}) - f_{0\perp}(\mathbf{x}))^2 + \frac{N}{n} \sigma^2, \\ \hat{L}_\parallel(f) &= \frac{1}{n} \sum_{i=1}^n (f_\parallel(\mathbf{x}) - y_{i\parallel})^2, & \hat{L}_\perp(f) &= \frac{1}{n} \sum_{i=1}^n (f_\perp(\mathbf{x}) - y_{i\perp})^2, \\ MSE_\parallel(f) &= \mathbb{E}_{\mathcal{D}} \left[\frac{1}{n} \sum_{i=1}^n (f_\parallel(\mathbf{x}) - f_{0\parallel}(\mathbf{x}))^2 \right], & MSE_\perp(f) &= \mathbb{E}_{\mathcal{D}} \left[\frac{1}{n} \sum_{i=1}^n (f_\perp(\mathbf{x}) - f_{0\perp}(\mathbf{x}))^2 \right], \end{aligned}$$

791 such that $L(f) = L_\parallel(f) + L_\perp(f)$, $\hat{L}(f) = \hat{L}_\parallel(f) + \hat{L}_\perp(f)$. This can be verified by de-
792 compositing \hat{f}, f_0 and y into two orthogonal components as shown above, and observing that
793 $\sum_{i=1}^n f_{1\perp}(\mathbf{x}_i) f_{2\parallel}(\mathbf{x}_i) = 0, \forall f_1, f_2$.

794 **First prove the following claim**

795 **Claim 15.** Assume that $\hat{f} = \arg \min_{f \in \mathcal{F}} \hat{L}(f)$ is the empirical risk minimizer. Then $\hat{f}_\perp =$
796 $\arg \min_{f \in \mathcal{F}_\perp} \hat{L}_\perp(f)$, $\hat{f}_\parallel = \arg \min_{f \in \mathcal{F}_\parallel} \hat{L}_\parallel(f)$, where \hat{f}_\perp is the projections of \hat{f} in \mathcal{F}_\perp , and
797 $\hat{f}_\parallel = \hat{f} - \hat{f}_\perp$ respectively.

798 *Proof.* Since $\hat{f} \in \mathcal{F}$, by definition $\hat{f}_\parallel \in \mathcal{F}_\parallel$. Assume that there exist $\hat{f}'_\perp, \hat{f}'_\parallel$, and either $\hat{L}_\perp(\hat{f}'_\perp) <$
799 $\hat{L}_\perp(\hat{f}_\perp)$, or $\hat{L}_\parallel(\hat{f}'_\parallel) < \hat{L}_\parallel(\hat{f}_\parallel)$. Then

$$\begin{aligned} \hat{L}(\hat{f}') &= \hat{L}(\hat{f}'_\perp + \hat{f}'_\parallel) = \hat{L}_\parallel(\hat{f}'_\perp + \hat{f}'_\parallel) + \hat{L}_\perp(\hat{f}'_\perp + \hat{f}'_\parallel) = \hat{L}_\parallel(\hat{f}'_\parallel) + \hat{L}_\perp(\hat{f}'_\perp) \\ &< \hat{L}_\parallel(\hat{f}_\parallel) + \hat{L}_\perp(\hat{f}_\perp) = \hat{L}_\parallel(\hat{f}_\parallel + \hat{f}_\perp) + \hat{L}_\perp(\hat{f}_\perp) = \hat{L}(\hat{f}) \end{aligned}$$

800 which shows that \hat{f} is not the minimizer of $\hat{L}(f)$ and violates the assumption.

801

□

802 **Then we bound $MSE_\perp(f)$.** We convert this part into a finite dimension least square problem:

$$\begin{aligned} \hat{f}_\perp &= \arg \min_{f \in \mathcal{F}_\perp} \hat{L}_\perp(f) \\ &= \arg \min_{f \in \mathcal{F}_\perp} \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - f_{0\perp}(\mathbf{x}_i) - \epsilon_{i\perp})^2 \\ &= \arg \min_{f \in \mathcal{F}_\perp} \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - f_{0\perp}(\mathbf{x}_i) - \epsilon_{i\perp})^2 + \epsilon_{i\parallel}^2 \\ &= \arg \min_{f \in \mathcal{F}_\perp} \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - f_{0\perp}(\mathbf{x}_i) - \epsilon_{i\perp} - \epsilon_{i\parallel})^2 \\ &= \arg \min_{f \in \mathcal{F}_\perp} \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - f_{0\perp}(\mathbf{x}_i) - \epsilon_i)^2 \end{aligned}$$

803 The forth line comes from our assumption that \mathcal{F}_\perp is orthogonal to \mathcal{F}_\parallel , so $\forall f \in \mathcal{F}_\perp, f + f_{0\perp} + \epsilon_\perp$
804 is orthogonal to ϵ_\parallel .

805 Let the basis function of \mathcal{F}_\perp be h_1, h_2, \dots, h_N , the above problem can be reparameterized as

$$\arg \min_{\boldsymbol{\theta} \in \mathbb{R}^N} \frac{1}{n} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|^2$$

806 where $\mathbf{X} \in \mathbb{R}^{n \times N} : X_i = h_j(\mathbf{x}_i)$, $\mathbf{y} = \mathbf{y}_{0\perp} + \boldsymbol{\epsilon}$, $\mathbf{y}_{0\perp} = [f_{0\perp}(x_1), \dots, f_{0\perp}(x_n)]$, $\boldsymbol{\epsilon} = [\epsilon_1, \dots, \epsilon_n]$.
 807 This problem has a closed-form solution

$$\boldsymbol{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

808 Observe that $f_{0\perp} \in \mathcal{F}_\perp$, let $\mathbf{y}_{0\perp} = \mathbf{X}\boldsymbol{\theta}^*$, The MSE of this problem can be computed by

$$\begin{aligned} L(\hat{f}_\perp) &= \frac{1}{n} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}_{0\perp}\|^2 = \frac{1}{n} \|\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\boldsymbol{\theta}^* + \boldsymbol{\epsilon}) - \mathbf{X}\boldsymbol{\theta}^*\|^2 \\ &= \frac{1}{n} \|\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}\|^2 \end{aligned}$$

809 Observing that $\Pi := \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is an idempotent and independent projection whose rank is
 810 N , and that $\mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T] = \sigma^2 \mathbf{I}$, we get

$$\text{MSE}_\perp(\hat{f}_\perp) = \mathbb{E}[L(\hat{f}_\perp)] = \frac{1}{n} \|\Pi\boldsymbol{\epsilon}\|^2 = \frac{1}{n} \text{tr}(\Pi\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T) = \frac{\sigma^2}{n} \text{tr}(\Pi)$$

811 which concludes that

$$\text{MSE}_\perp(\hat{f}) = O\left(\frac{N}{n} \sigma^2\right). \quad (22)$$

812 See also [17, Proposition 1].

813 **Next we study** $\text{MSE}_\parallel(\hat{f})$. Denote $\tilde{\sigma}_\parallel^2 = \frac{1}{n} \sum_{i=1}^n \epsilon_i^2$, $E = \max_i |\epsilon_i|$. Using Jensen's inequality and
 814 union bound, we have

$$\exp(t\mathbb{E}[E]) \leq \mathbb{E}[\exp(tE)] = \mathbb{E}[\max \exp(t|\epsilon_i|)] \leq \sum_{i=1}^n \mathbb{E}[\exp(t|\epsilon_i|)] \leq 2n \exp(t^2 \sigma^2 / 2)$$

815 Taking expectation over both sides, we get

$$\mathbb{E}[E] \leq \frac{\log 2n}{t} + \frac{t\sigma^2}{2}$$

816 maximizing the right hand side over t yields

$$\mathbb{E}[E] \leq \sigma \sqrt{2 \log 2n}.$$

817 Let $\tilde{\mathcal{F}}_\parallel$ be the covering set of $\mathcal{F}_\parallel = \{f_\parallel : f \in \mathcal{F}\}$. For any $\tilde{f}_\parallel \in \tilde{\mathcal{F}}_\parallel$,

$$\begin{aligned} L_\parallel(f_j) - \hat{L}_\parallel(f_j) &= \frac{1}{n} \sum_{i=1}^n (f_{j\parallel}(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2 - \frac{1}{n} \sum_{i=1}^n (\tilde{f}_\parallel(\mathbf{x}_i) - y_{i\parallel})^2 + \frac{n-N}{n} \sigma^2 \\ &= \frac{1}{n} \sum_{i=1}^n \epsilon_{i\parallel} (2\tilde{f}_\parallel(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i) - y_{i\parallel}) + \frac{n-N}{n} \sigma^2 \\ &= \frac{1}{n} \sum_{i=1}^n \epsilon_i (2\tilde{f}_\parallel(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i) - y_{i\parallel}) + \frac{n-N}{n} \sigma^2 \\ &= \frac{1}{n} \sum_{i=1}^n \epsilon_i (2\tilde{f}_\parallel(\mathbf{x}_i) - 2f_{0\parallel}(\mathbf{x}_i)) + \frac{n-N}{n} \sigma^2 - \tilde{\sigma}_\parallel^2 \end{aligned}$$

818 The first term can be bounded using Bernstein's inequality: let $h_i = \epsilon_i (f_{j\parallel}(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))$, by
 819 definition $|h_i| \leq 2EF$,

$$\begin{aligned} \text{Var}[h_i] &= \mathbb{E}[\epsilon_i^2 (\tilde{f}_\parallel(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2] \\ &= (\tilde{f}_\parallel(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2 \mathbb{E}[\epsilon_i^2] \\ &= (\tilde{f}_\parallel(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2 \sigma^2 \end{aligned}$$

820 using Bernstein's inequality, for any $\tilde{f}_{\parallel} \in \tilde{\mathcal{F}}_{\parallel}$, with probably at least $1 - \delta_p$,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \epsilon_i (2\tilde{f}_{\parallel}(\mathbf{x}_i) - 2f_{0\parallel}(\mathbf{x}_i)) &= \frac{2}{n} \sum_{i=1}^n h_i \\ &\leq \frac{2}{n} \sqrt{2 \sum_{i=1}^n (\tilde{f}_{\parallel}(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2 \sigma^2 \log(1/\delta_p)} + \frac{8EF \log(1/\delta_p)}{3n} \\ &= 2\sqrt{\left(L_{\parallel}(\tilde{f}_{\parallel}) - \frac{n-N}{n}\sigma^2\right) \frac{2\sigma^2 \log(1/\delta_p)}{n}} + \frac{8EF \log(1/\delta_p)}{3n} \\ &\leq \epsilon \left(L_{\parallel}(\tilde{f}_{\parallel}) - \frac{n-N}{n}\sigma^2\right) + \frac{8\sigma^2 \log(1/\delta_p)}{n\epsilon} + \frac{8EF \log(1/\delta_p)}{3n} \end{aligned}$$

821 the last inequality holds true for all $\epsilon > 0$. The union bound shows that with probably at least $1 - \delta$,
822 for all $\tilde{f}_{\parallel} \in \tilde{\mathcal{F}}_{\parallel}$,

$$\begin{aligned} L_{\parallel}(\tilde{f}_{\parallel}) - \hat{L}_{\parallel}(\tilde{f}_{\parallel}) &\leq \epsilon \left(L_{\parallel}(\tilde{f}_{\parallel}) - \frac{n-N}{n}\sigma^2\right) + \frac{8\sigma^2 \log(\mathcal{N}(\mathcal{F}_{\parallel}, \delta)/\delta_p)}{n\epsilon} + \frac{8EF \log(\mathcal{N}(\mathcal{F}_{\parallel}, \delta)/\delta_p)}{3n} \\ &\quad + \frac{n-N}{n}\sigma^2 - \tilde{\sigma}_{\parallel}^2. \end{aligned}$$

823 By rearranging the terms and using the definition of $L(\tilde{f}_{\parallel})$, we get

$$(1 - \epsilon) \left(L_{\parallel}(\tilde{f}_{\parallel}) - \frac{n-N}{n}\sigma^2\right) \leq \hat{L}_{\parallel}(\tilde{f}_{\parallel}) + \frac{8\sigma^2 \log(\mathcal{N}(\mathcal{F}_{\parallel}, \delta)/\delta_p)}{n\epsilon} + \frac{8EF \log(\mathcal{N}(\mathcal{F}_{\parallel}, \delta)/\delta_p)}{3n} - \tilde{\sigma}_{\parallel}^2.$$

824 Taking the expectation (over \mathcal{D}) on both sides, and notice that $\mathbb{E}[\tilde{\sigma}_{\parallel}^2] = \frac{n-N}{n}\sigma^2$. Furthermore, for
825 any random variable X , $\mathbb{E}[X] = \int_{-\infty}^{\infty} x dP(X \leq x)$, we get

$$\begin{aligned} &\max_{\tilde{f}_{\parallel} \in \tilde{\mathcal{F}}_{\parallel}} \left((1 - \epsilon) \text{MSE}_{\parallel}(\tilde{f}_{\parallel}) - \mathbb{E}[\hat{L}_{\parallel}(\tilde{f}_{\parallel})] \right) \\ &\leq \left(\frac{8\sigma^2}{n\epsilon} + \frac{8F\sigma\sqrt{2\log 2n}}{3n} \right) \left(\log \mathcal{N}(\mathcal{F}_{\parallel}, \delta) - \int_{\delta=0}^1 \log(\delta_p) d\delta_p \right) - \frac{n-N}{n}\sigma^2 \quad (23) \\ &= \left(\frac{8\sigma^2}{n\epsilon} + \frac{8F\sigma\sqrt{2\log 2n}}{3n} \right) (\log \mathcal{N}(\mathcal{F}_{\parallel}, \delta) + 1) - \frac{n-N}{n}\sigma^2. \end{aligned}$$

826 where the integration can be computed by replacing δ with e^x . Though it is not integrable under
827 Riemann integral, it is integrable under Lebesgue integration.

828 Similarly, let $\check{f}_{\parallel} = \arg \min_{f \in \mathcal{F}_{\parallel}} L_{\parallel}(f)$,

$$L_{\parallel}(\check{f}_{\parallel}) - \hat{L}_{\parallel}(\check{f}_{\parallel}) = \frac{1}{n} \sum_{i=1}^n \epsilon_i (2\check{f}_{\parallel}(\mathbf{x}_i) - 2f_{0\parallel}(\mathbf{x}_i)) + \frac{n-N}{n}\sigma^2 - \tilde{\sigma}_{\parallel}^2$$

829 with probably at least $1 - \delta_q$, for any $\epsilon > 0$,

$$\begin{aligned} -\frac{1}{n} \sum_{i=1}^n \epsilon_i (2\check{f}_{\parallel}(\mathbf{x}_i) - 2f_{0\parallel}(\mathbf{x}_i)) &\leq \epsilon \left(L_{\parallel}(\check{f}_{\parallel}) - \frac{n-N}{n}\sigma^2\right) + \frac{8\sigma^2 \log(1/\delta_p)}{n\epsilon} + \frac{8EF \log(1/\delta_p)}{3n}, \\ \hat{L}_{\parallel}(\check{f}_{\parallel}) &\leq (1 + \epsilon) \left(L_{\parallel}(\check{f}_{\parallel}) - \frac{n-N}{n}\sigma^2\right) + \frac{8\sigma^2 \log(1/\delta_p)}{n\epsilon} + \frac{8EF \log(1/\delta_q)}{3n} + \tilde{\sigma}_{\parallel}^2. \end{aligned}$$

830 Taking the expectation on both sides,

$$\mathbb{E}[\hat{L}_{\parallel}(\check{f}_{\parallel})] \leq (1 + \epsilon) \text{MSE}_{\parallel}(\check{f}_{\parallel}) + \frac{8\sigma^2}{n\epsilon} + \frac{8F\sigma\sqrt{2\log 2n}}{3n} + \frac{n-N}{n}\sigma^2. \quad (24)$$

831 Finally, let $\hat{f}_* := \arg \min_{f \in \tilde{\mathcal{F}}_{\parallel}} \sum_{i=1}^n (\hat{f}_{\parallel}(\mathbf{x}_i) - f(\mathbf{x}_i))^2$ be the projection of \hat{f}_{\parallel} in its δ -covering
 832 space,

$$\begin{aligned} \text{MSE}_{\parallel}(\hat{f}_{\parallel}) &= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n (\hat{f}_{\parallel}(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2 \right] \\ &= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n (\hat{f}_*(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2 + \frac{1}{n} \sum_{i=1}^n (\hat{f}_{\parallel}(\mathbf{x}_i) - \hat{f}_*(\mathbf{x}_i))(\hat{f}_{\parallel}(\mathbf{x}_i) + \hat{f}_*(\mathbf{x}_i) - 2f_{0\parallel}(\mathbf{x}_i)) \right] \\ &\leq \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n (\hat{f}_*(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2 \right] + 4F\delta \\ &= \text{MSE}_{\parallel}(\hat{f}_*(\mathbf{x}_i)) + 4F\delta, \end{aligned}$$

833 and similarly

$$\hat{L}_{\parallel}(\hat{f}_*) \leq \hat{L}_{\parallel}(\hat{f}_{\parallel}) + (4F + 2E)\delta. \quad (25)$$

834 We can conclude that

$$\begin{aligned} \text{MSE}_{\parallel}(\hat{f}_{\parallel}) &\leq \frac{1}{1-\epsilon} \left(\mathbb{E}[\hat{L}_{\parallel}(\hat{f}_*)] + \left(\frac{8\sigma^2}{n\epsilon} + \frac{8F\sigma\sqrt{2\log 2n}}{3n} \right) (\log \mathcal{N}(\mathcal{F}_{\parallel}, \delta) + 1) - \frac{n-N}{n}\sigma^2 \right) \\ &\quad + 4F\delta \\ &\leq \frac{1}{1-\epsilon} \left(\mathbb{E}[\hat{L}_{\parallel}(\hat{f}_{\parallel})] + (4F + \sigma\sqrt{8\log 2n})\delta \right) \\ &\quad + \left(\frac{8\sigma^2}{n\epsilon} + \frac{8F\sigma\sqrt{2\log 2n}}{3n} \right) (\log \mathcal{N}(\mathcal{F}_{\parallel}, \delta) + 1) - \frac{n-N}{n}\sigma^2 + 4F\delta \\ &\leq \frac{1}{1-\epsilon} \left(\mathbb{E}[\hat{L}_{\parallel}(\check{f}_{\parallel})] + (4F + \sigma\sqrt{8\log 2n})\delta \right) \\ &\quad + \left(\frac{8\sigma^2}{n\epsilon} + \frac{8F\sigma\sqrt{2\log 2n}}{3n} \right) (\log \mathcal{N}(\mathcal{F}_{\parallel}, \delta) + 1) - \frac{n-N}{n}\sigma^2 + 4F\delta \\ &\leq \frac{1+\epsilon}{1-\epsilon} \text{MSE}_{\parallel}(\check{f}_{\parallel}) + \frac{1}{n} \left(\frac{8\sigma^2}{\epsilon} + \frac{8F\sigma\sqrt{2\log 2n}}{3} \right) \left(\frac{\log \mathcal{N}(\mathcal{F}_{\parallel}, \delta) + 2}{1-\epsilon} \right) \\ &\quad + \left(4F + \frac{4F + \sigma\sqrt{8\log 2n}}{1-\epsilon} \right) \delta, \end{aligned}$$

835 where the first line comes from (23), and second comes from (25), the third line is because
 836 $\hat{f}_{\parallel} = \arg \min_{f \in \mathcal{F}_{\parallel}} \hat{L}_{\parallel}(f)$, and the last line comes from (24). We also use that fact that $\hat{L}_{\parallel}(\hat{f}) \leq$
 837 $\hat{L}_{\parallel}(f), \forall f$. Noticing that $\text{MSE}(\hat{f}) = \text{MSE}_{\parallel}(\hat{f}) + \text{MSE}_{\perp}(\hat{f})$, combining this with (22) finishes the
 838 proof. \square

839 G Detailed experimental setup

840 G.1 Target Functions

841 The doppler function used in Figure 2(d)-(f) is

$$f(x) = \sin(4/(x + 0.01)) + 1.5.$$

842 The ‘‘vary’’ function used in Figure 2(g)-(i) is

$$\begin{aligned} f(x) &= M_1(x/0.01) + M_1((x - 0.02)/0.02) + M_1((x - 0.06)/0.03) \\ &\quad + M_1((x - 0.12)/0.04) + M_3((x - 0.2)/0.02) + M_3((x - 0.28)/0.04) \\ &\quad + M_3((x - 0.44)/0.06) + M_3((x - 0.68)/0.08), \end{aligned}$$

843 where $x_+ := \max(x, 0)$. We uniformly take 256 samples from 0 to 1 in the piecewise cubic
 844 function experiment, and uniformly 1000 samples from 0 to 1 in the doppler function and ‘‘vary’’
 845 function experiment. We add zero mean independent (white) Gaussian noise to the observations.
 846 The standard derivation of noise is 0.05 in the piecewise cubic function experiment, 0.4 in the
 847 doppler function experiment and 0.1 in the ‘‘vary’’ function experiment.

848 G.2 Training/Fitting Method

849 In the piecewise polynomial function (“vary”) experiment, the depth of the PNN $L = 10$, the width
850 of each subnetwork $w = 10$, and the model contains $M = 500$ subnetworks. The depth of NN is also
851 10, and the width is 200 such that the NN and PNN have almost the same number of parameters. In
852 the doppler function experiment, the depth of the PNN $L = 12$, the width of each subnetwork $w =$
853 10, and the model contains $M = 2000$ subnetworks, because this problem requires a more complex
854 model to fit. The depth of NN is 12, and the width is 400. We used Adam optimizer with learning rate
855 of 10^{-3} . We first train the neural network layer by layer without weight decay. Specifically, we start
856 with a two-layer neural network with the same number of subnetworks and the same width in each
857 subnetwork, then train a three layer neural network by initializing the first layer using the trained
858 two layer one, until the desired depth is reached. After that, we turn the weight decay parameter and
859 train it until convergence. In both trend filtering and smoothing spline experiment, the order is 3,
860 and in wavelet denoising experiment, we use sym4 wavelet with soft thresholding. We implement
861 the trend filtering problem according to Tibshirani [37] using CVXPY, and use MOSEK to solve
862 the convex optimization problem. We directly call R function *smooth.spline* to solve smoothing
863 spline.

864 G.3 Post Processing

865 The degree of freedom of smoothing spline is returned by the solver in R, which is rounded to the
866 nearest integer when plotting. To estimate the degree of freedom of trend filtering, for each choice
867 of λ , we repeated the experiment for 10 times and compute the average number of nonzero knots as
868 estimated degree of freedom. For neural networks, we use the definition [38]:

$$2\sigma^2 \text{df} = \mathbb{E}\|\mathbf{y}' - \hat{\mathbf{y}}\|_2^2 - \mathbb{E}\|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 \quad (26)$$

869 where df denotes the degree of freedom, σ^2 is the variance of the noise, \mathbf{y} are the labels, $\hat{\mathbf{y}}$ are
870 the predictions and \mathbf{y}' are independent copy of \mathbf{y} . We find that estimating (26) directly by sampling
871 leads to large error when the degree of freedom is small. Instead, we compute

$$2\sigma^2 \hat{\text{df}} = \hat{\mathbb{E}}\|\mathbf{y}_0 - \hat{\mathbf{y}}\|_2^2 - \hat{\mathbb{E}}\|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 + \hat{\mathbb{E}}\|\mathbf{y} - \bar{y}_0\|_2^2 - \|\mathbf{y}_0 - \bar{y}_0\|_2^2 \quad (27)$$

872 where $\hat{\text{df}}$ is the estimated degree of freedom, \mathbb{E} denotes the empirical average (sample mean), \mathbf{y}_0 is
873 the target function and \bar{y}_0 is the mean of the target function in its domain.

874 **Proposition 16.** *The expectation of (27) over the dataset \mathcal{D} equals (26).*

Proof.

$$\begin{aligned} 2\sigma^2 \hat{\text{df}} &= \mathbb{E}_{\mathcal{D}}[\hat{\mathbb{E}}\|\mathbf{y}_0 - \hat{\mathbf{y}}\|_2^2 - \hat{\mathbb{E}}\|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 + \hat{\mathbb{E}}\|\mathbf{y} - \bar{y}_0\|_2^2 - \|\mathbf{y}_0 - \bar{y}_0\|_2^2] \\ &= \mathbb{E}\|\mathbf{y}_0 - \hat{\mathbf{y}}\|_2^2 - \mathbb{E}\|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 + \mathbb{E}_{\mathcal{D}}[\hat{\mathbb{E}}[(\mathbf{y} - \mathbf{y}_0)(\mathbf{y} + \mathbf{y}_0 - 2\bar{y}_0)]] \\ &= \mathbb{E}\|\mathbf{y}_0 - \hat{\mathbf{y}}\|_2^2 - \mathbb{E}\|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 + \mathbb{E}\left[\sum_{i=1}^n \epsilon_i(2y_i + \epsilon_i - 2\bar{y}_0)\right] \\ &= \mathbb{E}\|\mathbf{y}_0 - \hat{\mathbf{y}}\|_2^2 - \mathbb{E}\|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 + n\sigma^2 \\ &= \mathbb{E}\|\mathbf{y}' - \hat{\mathbf{y}}\|_2^2 - \mathbb{E}\|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 \end{aligned}$$

875 where \mathcal{D} denotes the dataset. In the third line, we make use of the fact that $\mathbb{E}[\epsilon_i] = 0$, $\mathbb{E}[\epsilon_i^2] = \sigma^2$,
876 and in the last line, we make use of $\mathbb{E}[\epsilon_i'] = 0$, $\mathbb{E}[\epsilon_i'^2] = \sigma^2$, and ϵ_i' are independent of y_i and $y_{0,i}$ \square

877 One can easily check that a “zero predictor” (a predictor that always predict \bar{y}_0 , and it always predicts
878 0 if the target function has zero mean) always has an estimated degree of freedom of 0.

879 G.4 More experimental results

880 G.4.1 Regularization weight vs degree-of-freedom

881 As we explained in the previous section, the degree of freedom is the exact information-theoretic
882 measure of the generalization gap. A Larger degree-of-freedom implies more overfitting.

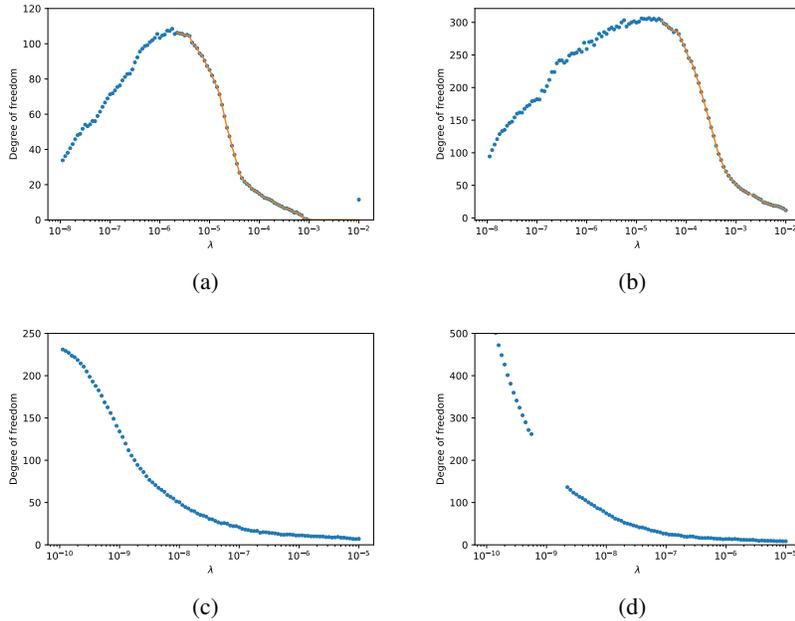


Figure 3: The relationship between degree of freedom and the scaling factor of the regularizer λ . The solid line shows the result after denoising. (a)(b) in a NN. (c)(d) In trend filtering. (a)(c): the piecewise cubic function. (b)(d) the doppler function.

883 In figure Figure 3, we show the relationship between the estimated degree of freedom and the scaling
 884 factor of the regularizer λ in a parallel neural network and in trend filtering. As is shown in the
 885 figure, generally speaking as λ decreases towards 0, the degree of freedom should increase too.
 886 However, for parallel neural networks, if λ is very close to 0, the estimated degree of freedom will
 887 not increase although the degree of freedom is much smaller than the number of parameters —
 888 actually even smaller than the number of subnetworks. Instead, it actually decreases a little. This
 889 effect has not been observed in other nonparametric regression methods, e.g. trend filtering, which
 890 overfits every noisy datapoint perfectly when $\lambda \rightarrow 0$. But for the neural networks, even if we do
 891 not regularize at all, the amount of overfitting is still relatively mild 30/256 vs 80/1000. In our
 892 experiments using neural networks, when λ is small, we denoise the estimated degree of freedom
 893 using isotonic regression.

894 We do not know the exact reason of this curious observation. Our hypothesis is that it might be
 895 related to issues with optimization, i.e., the optimizer ends up at a local minimum that generalizes
 896 better than a global minimum; or it could be connected to the “double descent” behavior of DNN
 897 [24] under over-parameterization.

898 G.4.2 Detailed numerical results

899 In order to allow the readers to view our result in detail, we plot the numerical experiment results of
 900 each method separately in Figure 4 and Figure 5.

901 G.4.3 Practical equivalence between the weight-decayed two-layer NN and L1-Trend 902 Filtering

903 In this section we investigate the equivalence of two-layer NN and the locally adaptive regression
 904 splines from Section B. In the special case when $m = 1$ the special regularization reduces to weight
 905 decay and the non-standard truncated power activation becomes ReLU. We compare L1 trend fil-
 906 tering [20] (shown to be equivalent to locally adaptive regression splines by Tibshirani [37]) and
 907 an overparameterized version of the neural network for all regularization parameter $\lambda > 0$, i.e.,
 908 a regularization path. The results are shown in Figure 6. It is clear that as the weight decay in-

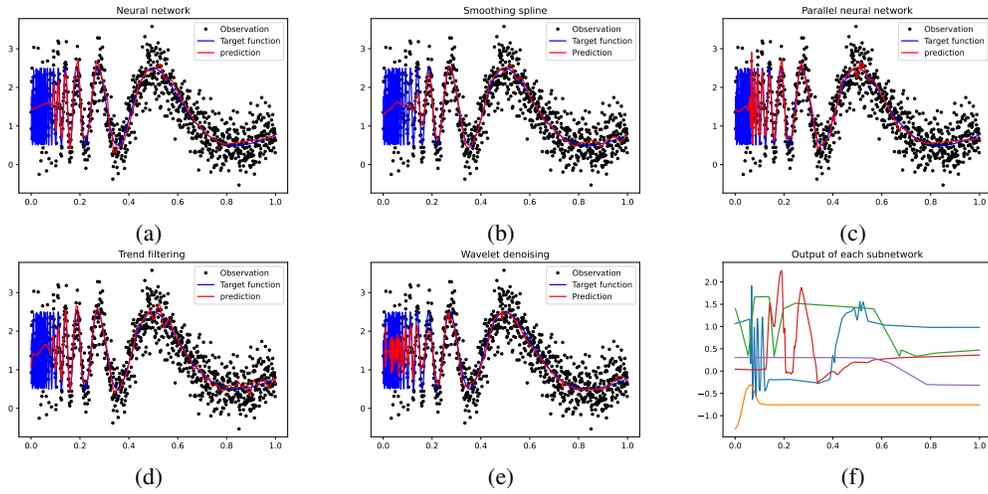


Figure 4: More experiments results of Doppler function.

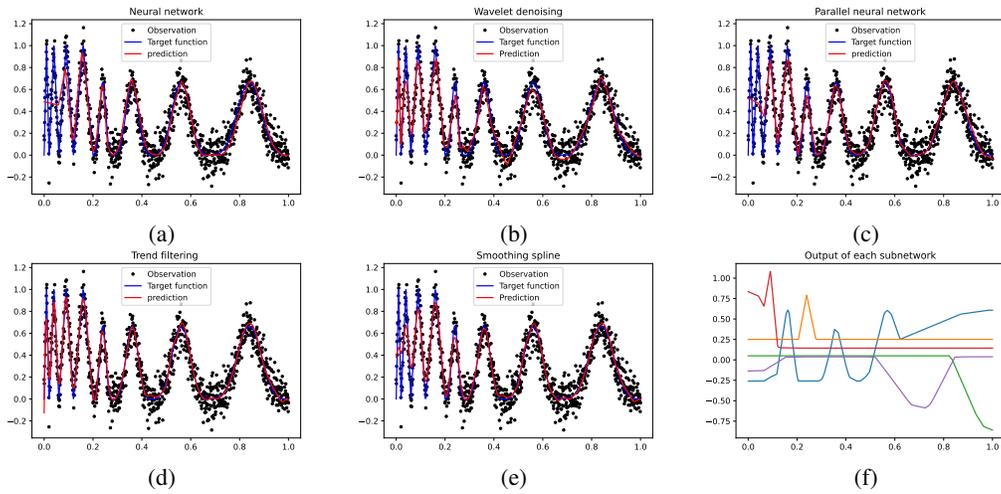


Figure 5: More experiments results of the “vary” function.

909 creases, it induces sparsity in the number of knots it selects similarly to L1-Trend Filtering, and the
 910 regularization path matches up nearly perfectly even though NNs are also learning knots locations.

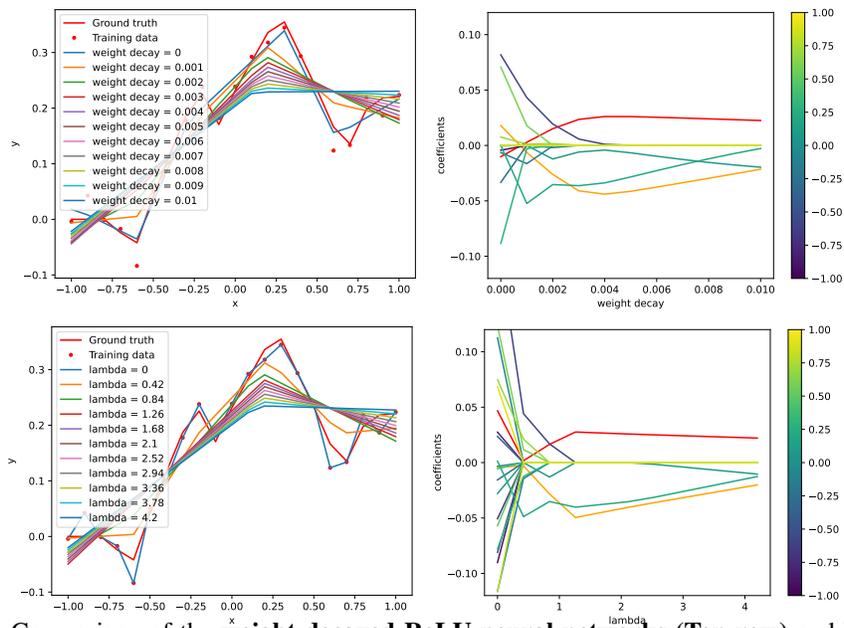


Figure 6: Comparison of the **weight decayed ReLU neural networks (Top row)** and **L1 Trend Filtering (Bottom row)** with different regularization parameters. The left column shows the fitted functions and the right column shows the *regularization path* (in the flavor of [15]) of the coefficients of the truncated power basis at individual data points (the free-knots learned by NN are snapped to the nearest input x to be comparable).