# Fourier-domain Variational Formulation and Its Well-posedness for Supervised Learning

**Tao Luo**[*]                                          LUOTAO41@SJTU.EDU.CN
**Zheng Ma**                                          ZHENGMA@SJTU.EDU.CN
**Zhiwei Wang**                                      VICTORYWZW@SJTU.EDU.CN
**Zhi-Qin John Xu**                                  XUZHIQIN@SJTU.EDU.CN
**Yaoyu Zhang**                                      ZHYY.SJTU@SJTU.EDU.CN
*School of Mathematical Sciences, Institute of Natural Sciences, MOE-LSC and Qing Yuan Research Institute, Shanghai Jiao Tong University, Shanghai, 200240, P.R. China*

## Abstract

A supervised learning problem is to find a function in a hypothesis function space given values on isolated data points. Inspired by the frequency principle in neural networks, we propose a Fourier-domain variational formulation for supervised learning problem. This formulation circumvents the difficulty of imposing the constraints of given values on isolated data points in continuum modelling. Under a necessary and sufficient condition within our unified framework, we establish the well-posedness of the Fourier-domain variational problem, by showing a critical exponent depending on the data dimension. In practice, a neural network can be a convenient way to implement our formulation, which automatically satisfies the well-posedness condition.

**Keywords:** Fourier-domain variational problem, well-posedness, critical exponent, frequency principle, supervised learning

## 1. Introduction

Supervised learning is ubiquitous. In a supervised learning problem, the goal is to find a function in a hypothesis function space given values on isolated data points with labels. In practice, Deep neural network (DNN), although with limit understanding, has been a powerful method. A series of works provide a good explanation for the good generalization of DNNs by showing a *Frequency Principle (F-Principle), i.e., a DNN tends to learn a target function from low to high frequencies during the training* (Xu et al., 2019, 2020; Rahaman et al., 2019). The F-Principle shows a low-frequency bias of DNNs when fitting a given data set. In the neural tangent kernel regime (Jacot et al., 2018; Lee et al., 2019), later works show that the long-time training solution of a wide two-layer neural network is equivalent to the solution of a constrained Fourier-domain variational problem (Zhang et al., 2019; Luo et al., 2020).

Inspired by above works about the F-Principle, in this paper, we propose a general Fourier-domain variational formulation for supervised learning problem and study its well-posedness. In continuum modelling, it is often difficult to impose the constraint of given values on isolated data points in a function space without sufficient regularity, e.g., a $L^p$ space. We circumvent this difficulty by regarding the Fourier-domain variation as the primal problem and the constraint of isolated data points is imposed through a linear operator. Under a necessary and sufficient condition within

---

[*] Corresponding author

our unified framework, we establish the well-posedness of the Fourier-domain variational problem. We show that the well-posedness depends on a critical exponent, which equals to the data dimension. This is a stark difference compared with a traditional partial differential equation (PDE) problem. For example, in a boundary value problem of any PDE in a $d$-dimensional domain, the boundary data should be prescribed on the $(d-1)$-dimensional boundary of the domain, where the dimension $d$ plays an important role. However, in a well-posed supervised learning problem, the constraint is always on isolated points, which are 0-dimensional independent of $d$, while the model has to satisfy a well-posedness condition depending on the dimension. In practice, a neural network can be a convenient way to implement our formulation, which automatically satisfies the well-posedness condition. With a clear understanding of its posedness, the Fourier-domain variational formulation also provides insight for designing methods for supervised learning problems.

The rest of the paper is organized as follows. Section 2 shows some related work. In section 3, we propose a Fourier-domain variational formulation for supervised learning problems. The necessary and sufficient condition for the well-posedness of our model is presented in section 4. Section 5 is devoted to the numerical demonstration in which we solve the Fourier-domain variational problem using band-limited functions. Finally, we present a short conclusion and discussion in section 6.

## 2. Related Works

Our work, as a modelling for supervised learning, is related to the point cloud interpolation problem which belongs to semi-supervised learning. One of the most widely used methods for the point cloud problems is the 2-Laplacian method (Zhu et al., 2003), which is an approach based on a Gaussian random field and weighted graph model. But it has been observed (El Alaoui et al., 2016; Nadler et al., 2009) that when the number of unlabeled data point is large, the graph Laplacian method is usually ill-posed. A new weighted Laplace method was proposed to overcome this shortcoming of the original 2-Laplacian method (Shi et al., 2017). Calder and Slepev (2019) further considered a way to correctly set the weights in Laplacian regularization with a exponent $\alpha > d$ and proved the well-posedness of the corresponding continuum model in the large-sample limit. We remark that our continuum model is proposed for the case of finite number of data points, i.e., $n < +\infty$, not the large-sample limit case.

From extensive synthetic and realistic datasets, frequency principle is proposed to characterize the training process of deep neural networks (Xu et al., 2019; Rahaman et al., 2019; Xu et al., 2020). A series of theoretical works subsequently show that frequency principle holds in different settings, for example, a non-NTK (neural tangent kernel) regime with infinite samples (Luo et al., 2019) and the NTK regime with finite samples (Zhang et al., 2019; Bordelon et al., 2020; Luo et al., 2020) or infinite samples (Cao et al., 2019; Ronen et al., 2019). E et al. (2020) show that the integral equation would naturally leads to the frequency principle. The frequency principle inspires the design of deep neural networks to fast learn a function with high frequency (Liu et al., 2020; Wang et al., 2020b; Jagtap et al., 2020; Cai et al., 2019; Biland et al., 2020; Li et al., 2020; Wang et al., 2020a).

## 3. Fourier-domain Variational Problem for Supervised Learning

### 3.1. Motivation: Linear Frequency Principle

In the following, we consider the regression problem of fitting a target function $f \in C_c(\mathbb{R}^d)$. Clearly, $f \in L^2(\mathbb{R}^d)$. Specifically, we use a DNN, $h_{\mathrm{DNN}}(\boldsymbol{x}, \boldsymbol{\theta}(t))$, to fit the training dataset

$\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$ of $n$ sample points, where $\boldsymbol{x}_i \in \mathbb{R}^d$, $y_i = f(\boldsymbol{x}_i)$ for each $i$. For the convenience of notation, we denote $\boldsymbol{X} = (\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n)^\mathsf{T}$, $\boldsymbol{Y} = (y_1, \ldots, y_n)^\mathsf{T}$. It has been shown in (Jacot et al., 2018; Lee et al., 2019) that, if the number of neurons in each hidden layer is sufficiently large, then $\|\boldsymbol{\theta}(t) - \boldsymbol{\theta}(0)\| \ll 1$ for any $t \geq 0$. In such cases, the the following function

$$h(\boldsymbol{x}, \boldsymbol{\theta}) = h_{\mathrm{DNN}}(\boldsymbol{x}, \boldsymbol{\theta}_0) + \nabla_{\boldsymbol{\theta}} h_{\mathrm{DNN}}(\boldsymbol{x}, \boldsymbol{\theta}_0) \cdot (\boldsymbol{\theta} - \boldsymbol{\theta}_0),$$

is a very good approximation of DNN output $h_{\mathrm{DNN}}(\boldsymbol{x}, \boldsymbol{\theta}(t))$ with $\boldsymbol{\theta}(0) = \boldsymbol{\theta}_0$. Note that, we have the following requirement for $h_{\mathrm{DNN}}$ which is easily satisfied for common DNNs: for any $\boldsymbol{\theta} \in \mathbb{R}^m$, there exists a weak derivative of $h_{\mathrm{DNN}}(\cdot, \boldsymbol{\theta}_0)$ with respect to $\boldsymbol{\theta}$ satisfying $\nabla_{\boldsymbol{\theta}} h_{\mathrm{DNN}}(\cdot, \boldsymbol{\theta}_0) \in L^2(\mathbb{R}^d)$.

Inspired by the F-Principle and the linear dynamics in the kernel regime, (Zhang et al., 2019; Luo et al., 2020) derived a Linear F-Principle (LFP) dynamics to effectively study the training dynamics of a two-layer ReLU NN with the mean square loss in the large width limit. Up to a multiplicative constant in the time scale, the gradient descent dynamics of a sufficiently wide two-layer NN is approximated by

$$\partial_t \mathcal{F}[u](\boldsymbol{\xi}, t) = -(\gamma(\boldsymbol{\xi}))^2 \mathcal{F}[u_\rho](\boldsymbol{\xi}), \tag{1}$$

where $u(\boldsymbol{x}, t) = h(\boldsymbol{x}, t) - h_{\mathrm{target}}(\boldsymbol{x})$, $u_\rho(\boldsymbol{x}) = u(\boldsymbol{x}, t)\rho(\boldsymbol{x})$. We follow this work and further assume that $\rho(\boldsymbol{x}) = \frac{1}{n} \sum_{i=1}^n \delta(\boldsymbol{x} - \boldsymbol{x}_i)$, accounting for the real case of a finite training dataset $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$, and

$$(\gamma(\boldsymbol{\xi}))^2 = \mathbb{E}_{a(0), r(0)} \left[ \frac{r(0)^3}{16\pi^4 \|\boldsymbol{\xi}\|^{d+3}} + \frac{a(0)^2 r(0)}{4\pi^2 \|\boldsymbol{\xi}\|^{d+1}} \right],$$

where $r(0) = |\boldsymbol{w}(0)|$ and the two-layer ReLU NN parameters at initial $a(0)$ and $\boldsymbol{w}(0)$ are random variables with certain given distribution. In this work, for any function $g$ defined on $\mathbb{R}^d$, we use the following convention of the Fourier transform and its inverse:

$$\mathcal{F}[g](\boldsymbol{\xi}) = \int_{\mathbb{R}^d} g(\boldsymbol{x}) \mathrm{e}^{-2\pi \mathrm{i} \boldsymbol{\xi}^\mathsf{T} \boldsymbol{x}} \, \mathrm{d}\boldsymbol{x}, \quad g(\boldsymbol{x}) = \int_{\mathbb{R}^d} \mathcal{F}[g](\boldsymbol{\xi}) \mathrm{e}^{2\pi \mathrm{i} \boldsymbol{x}^\mathsf{T} \boldsymbol{\xi}} \, \mathrm{d}\boldsymbol{\xi}.$$

Different from $\mathcal{F}[u](\boldsymbol{\xi}, t)$ on the left hand side, the formula on the right hand side reads as

$$\mathcal{F}[u_\rho](\boldsymbol{\xi}, t) = \mathcal{F}[u(\cdot, t)\rho(\cdot)](\boldsymbol{\xi}, t) = \frac{1}{n} \mathcal{F}\left[ \sum_{i=1}^n (h(\cdot, \boldsymbol{\theta}(t)) - y_i) \, \delta(\cdot - \boldsymbol{x}_i) \right](\boldsymbol{\xi}, t),$$

which incorporates the information of the training dataset. The solution of the LFP model (1) is equivalent to that of the following optimization problem in a proper hypothesis space $F_\gamma$,

$$\min_{h - h_{\mathrm{ini}} \in F_\gamma} \int_{\mathbb{R}^d} (\gamma(\boldsymbol{\xi}))^{-2} |\mathcal{F}[h](\boldsymbol{\xi}) - \mathcal{F}[h_{\mathrm{ini}}](\boldsymbol{\xi})|^2 \, \mathrm{d}\boldsymbol{\xi},$$

subject to constraints $h(\boldsymbol{x}_i) = y_i$ for $i = 1, \ldots, n$. The weight $(\gamma(\boldsymbol{\xi}))^{-2}$ grows as the frequency $\boldsymbol{\xi}$ increases, which means that a large penalty is imposed on the high frequency part of $h(\boldsymbol{x}) - h_{\mathrm{ini}}(\boldsymbol{x})$. As we can see, a random non-zero initial output of DNN leads to a specific type of generalization error. To eliminate this error, we use DNNs with an antisymmetrical initialization (ASI) trick (Zhang et al., 2020), which guarantees $h_{\mathrm{ini}}(\boldsymbol{x}) = 0$. Then the final output $h(\boldsymbol{x})$ is dominated by low frequency, and the DNN model possesses a good generalization.

3

### 3.2. Fourier-domain Variational Formulation

Inspired by the variational formulation of LFP model, we propose a new continuum model for the supervised learning. This is a variational problem with a parameter $\alpha > 0$ to be determined later:

$$\min_{h \in \mathcal{H}} Q_\alpha[h] = \int_{\mathbb{R}^d} \langle \boldsymbol{\xi} \rangle^\alpha |\mathcal{F}[h](\boldsymbol{\xi})|^2 \, \mathrm{d}\boldsymbol{\xi}, \tag{2}$$

$$\text{s.t.} \quad h(\boldsymbol{x}_i) = y_i, \quad i = 1, \cdots, n, \tag{3}$$

where $\langle \boldsymbol{\xi} \rangle = (1 + \|\boldsymbol{\xi}\|^2)^{\frac{1}{2}}$ is the "Japanese bracket" of $\boldsymbol{\xi}$ and $\mathcal{H} = \{h(x) | \int_{\mathbb{R}^d} \langle \boldsymbol{\xi} \rangle^\alpha |\mathcal{F}[h](\boldsymbol{\xi})|^2 \, \mathrm{d}\boldsymbol{\xi} < \infty\}$. Note that in the spatial domain, the evaluation on $n$ known data points is meaningless in the sense of $L^2$ functions. Therefore, we consider the problem in the frequency domain and define a linear operator $\mathcal{P}_{\boldsymbol{X}} : L^1(\mathbb{R}^d) \cap L^2(\mathbb{R}^d) \to \mathbb{R}^n$ for the given sample set $\boldsymbol{X}$ to transform the original constraints into the ones in the Fourier domain: $\mathcal{P}_{\boldsymbol{X}} \phi^* = \boldsymbol{Y}$. More precisely, we define for $\phi \in L^1(\mathbb{R}^d) \cap L^2(\mathbb{R}^d)$

$$\mathcal{P}_{\boldsymbol{X}} \phi := \left( \int_{\mathbb{R}^d} \phi(\boldsymbol{\xi}) \mathrm{e}^{2\pi \mathrm{i} \boldsymbol{\xi} \cdot \boldsymbol{x}_1} \, \mathrm{d}\boldsymbol{\xi}, \cdots, \int_{\mathbb{R}^d} \phi(\boldsymbol{\xi}) \mathrm{e}^{2\pi \mathrm{i} \boldsymbol{\xi} \cdot \boldsymbol{x}_n} \, \mathrm{d}\boldsymbol{\xi} \right)^{\mathsf{T}}. \tag{4}$$

The admissible function class reads as

$$\mathcal{A}_{\boldsymbol{X}, \boldsymbol{Y}} = \{\phi \in L^1(\mathbb{R}^d) \cap L^2(\mathbb{R}^d) \mid \mathcal{P}_{\boldsymbol{X}} \phi = \boldsymbol{Y}\}.$$

Notice that $\|\mathcal{F}^{-1}[\phi]\|_{H^{\frac{\alpha}{2}}} = \left( \int_{\mathbb{R}^d} \langle \boldsymbol{\xi} \rangle^\alpha |\phi(\boldsymbol{\xi})|^2 \, \mathrm{d}\boldsymbol{\xi} \right)^{\frac{1}{2}}$ is a Sobolev norm, which characterizes the regularity of the final output function $h(\boldsymbol{x}) = \mathcal{F}^{-1}[\phi](\boldsymbol{x})$. The larger the exponent $\alpha$ is, the better the regularity becomes.

For example, when $d = 1$ and $\alpha = 2$, by Parseval's theorem,

$$\|u\|_{H^1}^2 = \int_{\mathbb{R}} (1 + |\xi|^2) |\mathcal{F}[u](\xi)|^2 \, \mathrm{d}\xi = \int_{\mathbb{R}} u^2 + \frac{1}{4\pi^2} |\nabla u|^2 \, \mathrm{d}x.$$

Accordingly, the Fourier-domain variational problem reads as a standard variational problem in spatial domain. This is true for any quadratic Fourier-domain variational problem, but of course our Fourier-domain variational formulation is not necessarily being quadratic. The details for general cases (non-quadratic ones) are left to future work. For the quadratic setting with exponent $\alpha$, i.e., Problem (2), it is roughly equivalent to the following spatial-domain variational problem:

$$\min \int_{\mathbb{R}^d} (u^2 + |\nabla^{\frac{\alpha}{2}} u|^2) \, \mathrm{d}x.$$

This is clear for integer $\alpha/2$, while fractional derivatives are required for non-integer $\alpha/2$.

Back to our problem, after the above transformation, our goal is transformed into studying the following Fourier-domain variational problem,

**Problem 1** *Find a minimizer $\phi^*$ in $\mathcal{A}_{\boldsymbol{X}, \boldsymbol{Y}}$ such that*

$$\phi^* \in \arg \min_{\phi \in \mathcal{A}_{\boldsymbol{X}, \boldsymbol{Y}}} \|\mathcal{F}^{-1}[\phi]\|_{H^{\frac{\alpha}{2}}}^2. \tag{5}$$

This formulation is novel in the following two aspects:

1. We regard $\mathcal{F}[h]$ as the primal solution.

2. The evaluation on the sample points $\boldsymbol{x}_i$'s are imposed by the linear operator $\mathcal{P}_{\boldsymbol{X}}$.

Now we explain the importance of these new viewpoints for our task.

Traditionally, we all considered the problem in $\boldsymbol{x} - y$ space, the spatial domain. Recently, by the understanding of F-Principle (Xu et al., 2019, 2020; Luo et al., 2020; Zhang et al., 2019), we believe that if DNN is used to fit the data, it is more natural to consider the problem in the frequency domain. In particular, the functions defined on Fourier domain are assumed to be primal. And our variational problem is asked for such functions.

We remark that the operator $\mathcal{P}_{\boldsymbol{X}}$ is the inverse Fourier transform with evaluations on sample points $\boldsymbol{X}$. Actually, the linear operator $\mathcal{P}_{\boldsymbol{X}}$ projects a function defined on $\mathbb{R}^d$ to a function defined on 0-dimensional manifold $\boldsymbol{X}$. Just like the (linear) trace operator $T$ in a Sobolev space projects a function defined on $d$-dimensional manifold into a function defined on $(d-1)$-dimensional boundary manifold. Note that the only function space over the 0-dimensional manifold $\boldsymbol{X}$ is the $n$-dimensional vector space $\mathbb{R}^n$, where $n$ is the number of data points, while any Sobolev (or Besov) space over $d$-dimensional manifold ($d \geq 1$) is an infinite dimensional vector space.

## 4. Existence and Non-existence of Fourier-domain Variantional Problems

In this section, we consider the existence/non-existence dichotomy to Problem 1. In subsection 4.1, we prove that there is no solution to the Problem 1 in subcritical case $\alpha < d$. The supercritical case $\alpha > d$ will be investigated in subsection 4.2, where we prove the optimal function is a continuous and nontrivial solution. All proof of propositions and theorems in this section can be found in Appendix.

### 4.1. Subcritical Case: $\alpha < d$

In order to prove the nonexistence of the solution to the Problem 1 in $\alpha < d$ case, at first we need to find a class of functions that make the norm tend to zero. Let $\psi_\sigma(\boldsymbol{\xi}) = (2\pi)^{\frac{d}{2}} \sigma^d \mathrm{e}^{-2\pi^2 \sigma^2 \|\boldsymbol{\xi}\|^2}$ , then by direct calculation, we have $\mathcal{F}^{-1}[\psi_\sigma](\boldsymbol{x}) = \mathrm{e}^{-\frac{\|\boldsymbol{x}\|^2}{2\sigma^2}}$. For $\alpha < d$ the following proposition shows that the norm $\|\mathcal{F}^{-1}[\psi_\sigma]\|^2_{H^{\frac{\alpha}{2}}}$ can be sufficiently small as $\sigma \to 0$.

**Proposition 1 (critical exponent)** *For any input dimension d, we have*

$$\lim_{\sigma \to 0} \|\mathcal{F}^{-1}[\psi_\sigma]\|^2_{H^{\frac{\alpha}{2}}} = \begin{cases} 0, & \alpha < d, \\ C_d, & \alpha = d, \\ \infty, & \alpha > d. \end{cases} \tag{6}$$

*Here the constant $C_d = \frac{1}{2}(d-1)!(2\pi)^{-d}\frac{2\pi^{d/2}}{\Gamma(d/2)}$ only depends on the dimension d.*

**Remark 1** *The function $\mathcal{F}^{-1}[\psi]$ can be any function in the Schwartz space, not necessarily Gaussian. Proposition 1 still holds with (possibly) different $C_d$.*

For every small $\sigma$, we can use $n$ rapidly decreasing functions $\mathcal{F}^{-1}[\psi_\sigma](\boldsymbol{x} - \boldsymbol{x}_i)$ to construct the solution $\mathcal{F}^{-1}[\phi_\sigma](\boldsymbol{x})$ of the supervised learning problem. However, according to Proposition 1, when the parameter $\sigma$ tends to 0, the limit is the zero function in the sense of $L^2(\mathbb{R}^d)$. Therefore we have the following theorem:

**Theorem 1 (non-existence)** *Suppose that $\boldsymbol{Y} \neq \boldsymbol{0}$. For $\alpha < d$, there is no function $\phi^* \in \mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}$ satisfying*

$$\phi^* \in \arg \min_{\phi \in \mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}} \|\mathcal{F}^{-1}[\phi]\|^2_{H^{\frac{\alpha}{2}}}.$$

*In other words, there is no solution to the Problem 1.*

### 4.2. Supercritical Case: $\alpha > d$

In this section, we provide a theorem to establish the existence of the minimizer for Problem 1 in the case of $\alpha > d$.

**Theorem 2 (existence)** *For $\alpha > d$, there exists $\phi^* \in \mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}$ satisfying*

$$\phi^* \in \arg \min_{\phi \in \mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}} \|\mathcal{F}^{-1}[\phi]\|^2_{H^{\frac{\alpha}{2}}}.$$

*In other words, there exists a solution to the Problem 1.*

**Remark 2** *Note that, according to the Sobolev embedding theorem (Adams and Fournier, 2003; Evans, 1999), the minimizer in Theorem 2 has smoothness index no less than $[\frac{\alpha-d}{2}]$.*

## 5. Numerical Results

In this section, we illustrate our results by solving Fourier-domain variational problems numerically. We use uniform mesh in frequency domain with mesh size $\Delta\xi$ and band limit $M\Delta\xi$. In this discrete setting, the considered space becomes $\mathbb{R}^{(2M)^d}$. We emphasize that the numerical solution with this setup always exists even for the subcritical case which corresponds to the non-existence theorem. However, as we will show later, the numerical solution is trivial in nature when $\alpha < d$.

### 5.1. Special Case: One Data Point in One Dimension

To simplify the problem, we start with a single point $X = 0 \in \mathbb{Z}$ with the label $Y = 2$. Denote $\phi_j = \phi(\xi_j)$ for $j \in \mathbb{Z}$. We also assume that the function $\phi$ is an even function. Then according to the definition of $\mathcal{P}_{\boldsymbol{X}}$, we have the following problem:

**Example 1 (Problem 1 with a particular discretization)**

$$\min_{\phi \in \mathbb{R}^M} \sum_{j=1}^{M} (1 + j^2 \Delta\xi^2)^{\frac{\alpha}{2}} |\phi_j|^2, \tag{7}$$

$$\text{s.t.} \quad \sum_{j=1}^{M} \phi_j \Delta\xi = 1, \tag{8}$$

where we further assume $\phi_0 = \phi(0) = 0$. If we denote $\boldsymbol{\phi} = (\phi_1, \phi_2, \ldots, \phi_M)^\mathsf{T}$, $b = \frac{1}{\Delta\xi}$, $\boldsymbol{A} = (1, 1, \ldots, 1) \in \mathbb{R}^M$ and

$$\boldsymbol{\Gamma} = \sqrt{\lambda} \begin{pmatrix} (1 + 1^2 \Delta\xi^2)^{\frac{\alpha}{4}} & & & \\ & (1 + 2^2 \Delta\xi^2)^{\frac{\alpha}{4}} & & \\ & & \ddots & \\ & & & (1 + M^2 \Delta\xi^2)^{\frac{\alpha}{4}} \end{pmatrix}.$$

In fact this is a standard Tikhonov regularization (Tikhonov and Arsenin, 1977) also known as ridge regression problem with the Lagrange multiplier $\lambda$. The corresponding ridge regression problem is,

$$\min_{\phi} \|A\phi - b\|_2^2 + \|\Gamma\phi\|_2^2, \tag{9}$$

where we put $\lambda$ in the optimization term $\|\Gamma\phi\|_2^2$, instead of the constraint term $\|A\phi - b\|_2^2$. This problem admits an explicit and unique solution (Tikhonov and Arsenin, 1977),

$$\phi = (A^{\mathsf{T}}A + \Gamma^{\mathsf{T}}\Gamma)^{-1}A^{\mathsf{T}}b. \tag{10}$$

Here we need to point out that the above method is also applicable to the case that the matrix $\Gamma$ is not diagonal.

Back to our problem, in order to obtain the explicit expression for the optimal $\phi$ we need the following relation between the solution of the ridge regression and the singular-value decomposition (SVD).

By denoting $\tilde{\Gamma} = I$ and

$$\tilde{A} = A\Gamma^{-1} = \frac{1}{\sqrt{\lambda}}\left((1 + 1^2\Delta\xi^2)^{\frac{\alpha}{4}}, (1 + 2^2\Delta\xi^2)^{\frac{\alpha}{4}}, \ldots, (1 + M^2\Delta\xi^2)^{\frac{\alpha}{4}}\right),$$

where $I$ is the diagonal matrix, the optimal solution (10) can be written as

$$\phi = (\Gamma^{\mathsf{T}})^{-1}\left(\tilde{A}^{\mathsf{T}}\tilde{A} + I\right)^{-1}\Gamma^{-1}A^{\mathsf{T}}b = (\Gamma^{\mathsf{T}})^{-1}\left(\tilde{A}^{\mathsf{T}}\tilde{A} + I\right)^{-1}\tilde{A}^{\mathsf{T}}b = (\Gamma^{\mathsf{T}})^{-1}\tilde{\phi},$$

where $\tilde{\phi} = \left(\tilde{A}^{\mathsf{T}}\tilde{A} + I\right)^{-1}\tilde{A}^{\mathsf{T}}b$ is the solution of ridge regression with $\tilde{A}$ and $\tilde{\Gamma}$. In order to obtain the explicit expression for $\tilde{\phi}$ we need the following relation between the solution of the ridge regression and the singular-value decomposition (SVD).

**Lemma 1** *If $\tilde{\Gamma} = I$, then this least-squares solution can be solved using SVD. Given the singular value decomposition*

$$\tilde{A} = U\Sigma V^{\mathsf{T}},$$

*with singular values $\sigma_i$, the Tikhonov regularized solution can be expressed aspects*

$$\tilde{\phi} = VDU^{\mathsf{T}}b,$$

*where $D$ has diagonal values*

$$D_{ii} = \frac{\sigma_i}{\sigma_i^2 + 1},$$

*and is zero elsewhere.*

**Proof** In fact, $\tilde{\phi} = (\tilde{A}^{\mathsf{T}}\tilde{A} + \tilde{\Gamma}^{\mathsf{T}}\tilde{\Gamma})^{-1}\tilde{A}^{\mathsf{T}}b = V(\Sigma^{\mathsf{T}}\Sigma + 1I)^{-1}V^{\mathsf{T}}V\Sigma^{\mathsf{T}}U^{\mathsf{T}}b$
$= VDU^{\mathsf{T}}b$, which completes the proof. ∎

Since $\tilde{A}\tilde{A}^{\mathsf{T}} = \frac{1}{\lambda}\sum_{j=1}^{M}(1 + j^2\Delta\xi^2)^{-\frac{\alpha}{2}}$, we have $\tilde{A} = U\Sigma V^{\mathsf{T}}$ with

$$U = 1, \quad \Sigma = \frac{1}{\sqrt{\lambda}}\left(\sum_{j=1}^{M}(1 + j^2\Delta\xi^2)^{-\frac{\alpha}{2}}\right)^{\frac{1}{2}} := Z/\sqrt{\lambda},$$

$$\boldsymbol{V} = \left( (1 + 1^2 \Delta \xi^2)^{-\frac{\alpha}{2}}/Z, (1 + 2^2 \Delta \xi^2)^{-\frac{\alpha}{2}}/Z, \ldots, (1 + M^2 \Delta \xi^2)^{-\frac{\alpha}{2}}/Z \right)^{\mathsf{T}}.$$

Then we get the diagonal value

$$D = \frac{Z/\sqrt{\lambda}}{Z^2/\lambda + 1}.$$

Therefore, by Lemma 1

$$\tilde{\phi} = \boldsymbol{V}DUb = \frac{1/\sqrt{\lambda}}{Z^2/\lambda + 1} \left( (1 + 1^2 \Delta \xi^2)^{-\frac{\alpha}{2}}, (1 + 2^2 \Delta \xi^2)^{-\frac{\alpha}{2}}, \ldots, (1 + M^2 \Delta \xi^2)^{-\frac{\alpha}{2}} \right)^{\mathsf{T}} b.$$

Finally, for the original optimal solution

$$\phi = (\boldsymbol{\Gamma}^{\mathsf{T}})^{-1}\tilde{\phi} = \frac{1}{(Z^2 + \lambda)\Delta \xi} \left( (1 + 1^2 \Delta \xi^2)^{-\frac{\alpha}{2}}, (1 + 2^2 \Delta \xi^2)^{-\frac{\alpha}{2}}, \ldots, (1 + M^2 \Delta \xi^2)^{-\frac{\alpha}{2}} \right)^{\mathsf{T}},$$

which means

$$\phi_j = \frac{(1 + j^2 \Delta \xi^2)^{-\frac{\alpha}{2}}}{(Z^2 + \lambda)\Delta \xi}.$$

To derive the function in $x$ space, say $h(x)$ then

$$
\begin{aligned}
h(x) &= \frac{1}{(Z^2 + \lambda)} \sum_{j=-M}^{M} (1 + j^2 \Delta \xi^2)^{-\frac{\alpha}{2}} e^{2\pi i j x} \\
&= \frac{2}{(Z^2 + \lambda)} \sum_{j=1}^{M} (1 + j^2 \Delta \xi^2)^{-\frac{\alpha}{2}} \cos(2\pi j x).
\end{aligned}
\tag{11}
$$

Fig. 1 shows that for this special case with a large $M$, $h(x)$ is not an trivial function in $\alpha > d$ case and degenerates to a trivial function in $\alpha < d$ case.

### 5.2. General Case: $n$ Points in $d$ Dimension

Assume that we have $n$ data points $\boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_n \in \mathbb{R}^d$ and each data point has $d$ components:

$$\boldsymbol{x}_i = (x_{i1}, x_{i2}, \ldots, x_{id})^{\mathsf{T}}$$

and denote the corresponding label as $(y_1, y_2, \ldots, y_n)^{\mathsf{T}}$. For the sake of simplicity, we denote the vector $(j_1, j_2, \cdots, j_d)^{\mathsf{T}}$ by $\boldsymbol{J}_{j_1 \ldots j_d}$. Then our problem becomes

**Example 2 (Problem 1 with general discretization)**

$$\min_{\phi \in \mathbb{R}^{(2M)^d}} \sum_{j_1, \ldots, j_d = -M}^{M} (1 + \|\boldsymbol{J}_{j_1 \ldots j_d}\|^2 \Delta \xi^2)^{\frac{\alpha}{2}} |\phi_{j_1 \ldots j_d}|^2, \tag{12}$$

$$\text{s.t.} \quad \sum_{j_1, \ldots, j_d = -M}^{M} \phi_{j_1 \ldots j_d} e^{2\pi i \Delta \xi \boldsymbol{J}_{j_1 \ldots j_d}^{\mathsf{T}} \boldsymbol{x}_k} = y_k, \quad k = 1, 2, \ldots, d \tag{13}$$
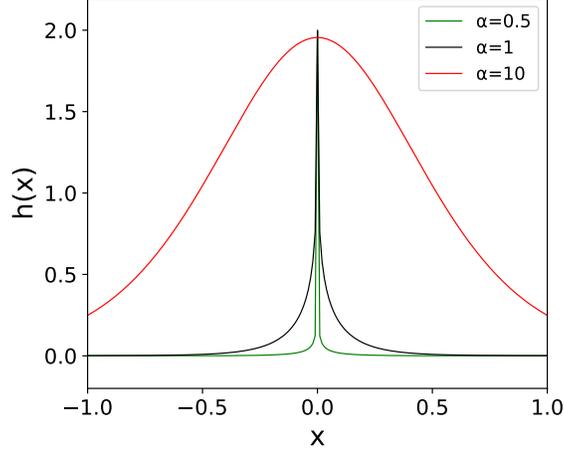
Figure 1: Fitting the function $h(x)$ shown in equation (11) with different exponent $\alpha$'s. Here we take $M = 10^6, \Delta\xi = 0.01, \lambda = 1$ and different $\alpha$ and observe that $h(x)$ is not an trivial function in $\alpha > d$ case and degenerates to a trivial function in $\alpha < d$ case.

The calculation of this example can be completed by the method analogous to the one used in subsection 5.1. Let

$$\boldsymbol{A}_j = \left(\mathrm{e}^{2\pi\mathrm{i}\Delta\xi \boldsymbol{J}^{\mathsf{T}}_{-M-M\ldots-M}\boldsymbol{x}_j}, \ldots, \mathrm{e}^{2\pi\mathrm{i}\Delta\xi \boldsymbol{J}^{\mathsf{T}}_{j_1 j_2\ldots j_d}\boldsymbol{x}_j}, \ldots, \mathrm{e}^{2\pi\mathrm{i}\Delta\xi \boldsymbol{J}^{\mathsf{T}}_{MM\ldots M}\boldsymbol{x}_j}\right)^{\mathsf{T}}, \ j = 1, 2, \ldots, n, \tag{14}$$

$$\boldsymbol{A} = (\boldsymbol{A}_1, \boldsymbol{A}_2, \ldots, \boldsymbol{A}_n)^{\mathsf{T}} \in \mathbb{R}^{n \times (2M)^d}, \quad \boldsymbol{b} = (y_1, y_2, \ldots, y_n)^{\mathsf{T}} \in \mathbb{R}^{n \times 1}, \tag{15}$$

$$\boldsymbol{\Gamma} = \lambda \begin{pmatrix} \ddots & & \\ & (1 + \|\boldsymbol{J}_{j_1 j_2\ldots j_d}\|^2 \Delta\xi^2)^{\frac{\alpha}{4}} & \\ & & \ddots \end{pmatrix} \in \mathbb{R}^{(2M)^d \times (2M)^d}. \tag{16}$$

We just need to solve the following equation:

$$\boldsymbol{\phi} = (\boldsymbol{A}^{\mathsf{T}}\boldsymbol{A} + \boldsymbol{\Gamma}^{\mathsf{T}}\boldsymbol{\Gamma})^{-1}\boldsymbol{A}^{\mathsf{T}}b. \tag{17}$$

Then we can get the output function $h(x)$ by using inverse Fourier transform:

$$h(\boldsymbol{x}) = \sum_{j_1,\ldots,j_d=-M}^{M} \phi_{j_1\ldots j_d}\mathrm{e}^{2\pi\mathrm{i}\Delta\xi \boldsymbol{J}_{j_1\ldots j_d}\cdot\boldsymbol{x}} \tag{18}$$

Since the size of the matrix is too large, it is difficult to solve $\phi$ by an explicit calculation. Thus we choose special $n, d$ and $M$ and show that $h(x)$ is not a trivial solution (non-zero function).

In our experiment, we set the hyper-parameter $M, \alpha, \lambda, \Delta\xi$ in advance. We set $\lambda = 0, 5, \Delta\xi = 0.1$ in 1-dimensional case and $\lambda = 0.2, \Delta\xi = 0.1$ in 2-dimensional case. We select two data points $\{(-0.5, 0.9), (0.5, 0.9)\}$ as the given points in 1-dimensional case and four points as given points in 2-dimensional case whose second coordinates are 0.5 so that it is convenient to observe the phenomenon. At first, we use formula (14), (15) and (16) to calculate matrix $\boldsymbol{A}, \boldsymbol{\Gamma}$ and vector $\boldsymbol{b}$.

9

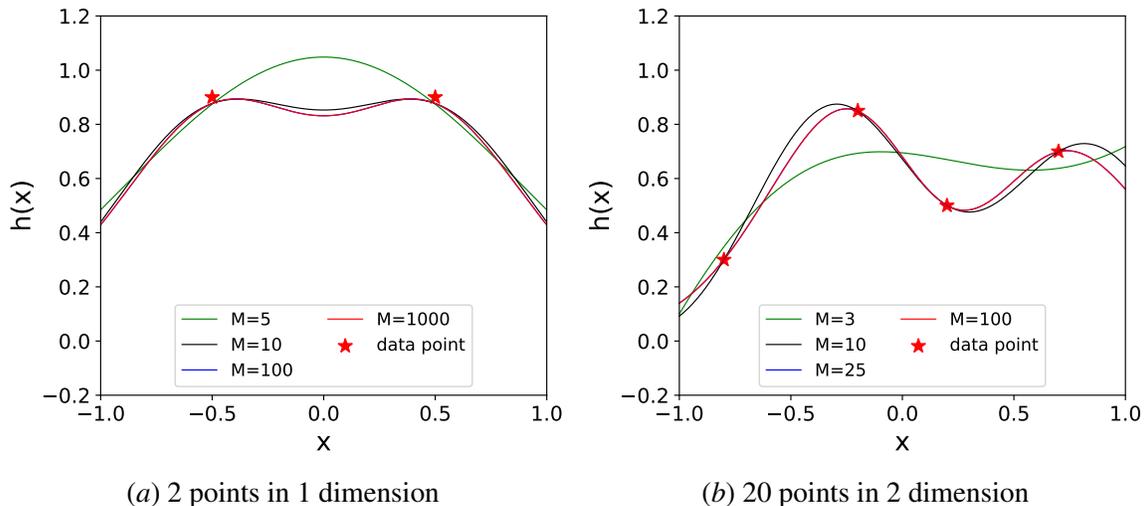(a) 2 points in 1 dimension      (b) 20 points in 2 dimension

Figure 2: Fitting data points in different dimensions with different band limit $M$. We use a proper $\alpha$ ($\alpha > d$) and observe that even for a large $M$, the function $h(x)$ does not degenerate to a trivial function. Note that the blue curve and the red one overlap with each. Here the trivial function represents a function whose value decays rapidly to zero except for the given training points.

Then from the equation (17) we can deduce vector $\phi$. The final output function $h(\boldsymbol{x})$ is obtained by inverse discrete Fourier transform (18).

In Fig.2, we set $\alpha = 10$ in both cases to ensure $\alpha > d$ and change the band limit $M$. We observe that as $M$ increases, the fitting curve converges to a non-trivial curve. In Fig.3, we set $M = 1000$ in 1-dimensional case and $M = 100$ in 2-dimensional case. By changing exponent $\alpha$, we can see in all cases, the fitting curves are non-trivial when $\alpha > d$, but degenerate when $\alpha < d$.

## 6. Conclusion

In this paper, we study the supervised learning problem by proposing a Fourier-domain variational formulation motivated by the frequency principle in deep learning. We establish the sufficient and necessary conditions for the well-posedness of the Fourier-domain variational problem, followed by numerical demonstration.

Our Fourier-domain variational formulation provides a novel viewpoint for modelling machine learning problem, that is, imposing more constraints, e.g., higher regularity, on the model rather than the data (always isolated points in practice) can give us the well-posedness as dimension of the problem increases. This is different from the modelling in physics and traditional point cloud problems, in which the model is independent of dimension in general. Our work suggests a potential approach of algorithm design by considering a dimension-dependent model for data modelling.

In contrast to the natural sciences, where models are usually derived from fundamental physical laws, in data science, the story may be totally different: we may propose mathematical models based on the algorithms of great practical success like DNNs. Afterwards, the continuum formulation and its mathematical validation can be analyzed. This seems to be a new scientific paradigm, along which our work plays a role as one step.
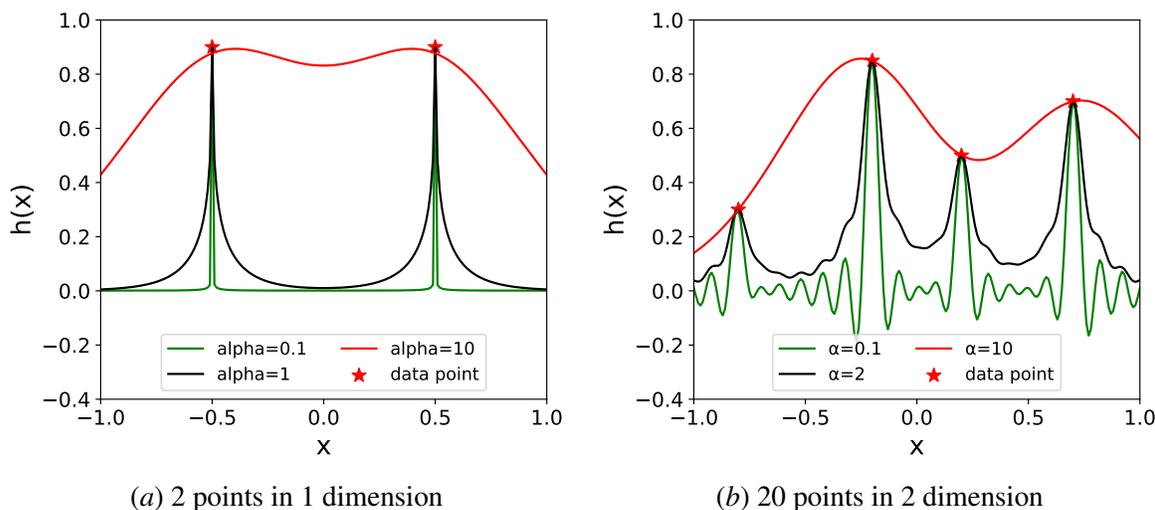
(*a*) 2 points in 1 dimension        (*b*) 20 points in 2 dimension

Figure 3: Fitting data points in different dimensions with different exponent $\alpha$'s. We observe that with a proper $M$, the function $h(x)$ is not a trivial function for $\alpha > d$ case and degenerates to a trivial function for $\alpha > d$ case.

## References

Robert A. Adams and John J.F. Fournier. *Sobolev spaces*. Elsevier Science, 2003.

Simon Biland, Vinicius C. Azevedo, Byungsoo Kim, and Barbara Solenthaler. Frequency-aware reconstruction of fluid simulations with generative networks. In Alexander Wilkie and Francesco Banterle, editors, *Eurographics 2020 - Short Papers*. The Eurographics Association, 2020. ISBN 978-3-03868-101-4. doi: 10.2312/egs.20201019.

Blake Bordelon, Abdulkadir Canatar, and Cengiz Pehlevan. Spectrum dependent learning curves in kernel regression and wide neural networks. *arXiv preprint arXiv:2002.02561*, 2020.

Wei Cai, Xiaoguang Li, and Lizuo Liu. A phase shift deep neural network for high frequency approximation and wave problems. *Accepted by SISC, arXiv:1909.11759*, 2019.

Jeff Calder and Dejan Slepev. Properly-weighted graph Laplacian for semi-supervised learning. *Applied Mathematics and Optimization*, (4), 2019.

Yuan Cao, Zhiying Fang, Yue Wu, Ding-Xuan Zhou, and Quanquan Gu. Towards understanding the spectral bias of deep learning. *arXiv preprint arXiv:1912.01198*, 2019.

Weinan E, Chao Ma, and Lei Wu. Machine learning from a continuous viewpoint, I. *Science China Mathematics*, pages 1–34, 2020.

Ahmed El Alaoui, Xiang Cheng, Aaditya Ramdas, Martin J Wainwright, and Michael I Jordan. Asymptotic behavior of $l_p$-based Laplacian regularization in semi-supervised learning. In *Conference on Learning Theory*, pages 879–906, 2016.

Lawrence C. Evans. Partial differential equations. *Mathematical Gazette*, 83(496):185, 1999.

Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018.

Ameya D Jagtap, Kenji Kawaguchi, and George Em Karniadakis. Adaptive activation functions accelerate convergence in deep and physics-informed neural networks. *Journal of Computational Physics*, 404:109136, 2020.

Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. In *Advances in neural information processing systems*, pages 8572–8583, 2019.

Xi-An Li, Zhi-Qin John Xu, and Lei Zhang. A multi-scale DNN algorithm for nonlinear elliptic equations with multiple scales. *Communications in Computational Physics*, 28(5):1886–1906, 2020.

Ziqi Liu, Wei Cai, and Zhi-Qin John Xu. Multi-scale deep neural network (MscaleDNN) for solving Poisson-Boltzmann equation in complex domains. *Communications in Computational Physics*, 28(5):1970–2001, 2020.

Tao Luo, Zheng Ma, Zhi-Qin John Xu, and Yaoyu Zhang. Theory of the frequency principle for general deep neural networks. *arXiv preprint arXiv:1906.09235*, 2019.

Tao Luo, Zheng Ma, Zhi-Qin John Xu, and Yaoyu Zhang. On the exact computation of linear frequency principle dynamics and its generalization. *arXiv preprint arXiv:2010.08153*, 2020.

B. Nadler, Nathan Srebro, and Xueyuan Zhou. Semi-supervised learning with the graph Laplacian: the limit of infinite unlabelled data. In *NIPS 2009*, 2009.

Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *International Conference on Machine Learning*, pages 5301–5310, 2019.

Basri Ronen, David Jacobs, Yoni Kasten, and Shira Kritchman. The convergence rate of neural networks for learned functions of different frequencies. In *Advances in Neural Information Processing Systems*, pages 4763–4772, 2019.

Zuoqiang Shi, Stanley Osher, and Wei Zhu. Weighted nonlocal Laplacian on interpolation from sparse data. *Journal of Scientific Computing*, 73(2-3):1–14, 2017.

Andrej Nikolaevich Tikhonov and Vasiliy Yakovlevich Arsenin. Solutions of ill-posed problems. *Mathematics of Computation*, 32(144):491–491, 1977.

Bo Wang, Wenzhong Zhang, and Wei Cai. Multi-scale deep neural network (mscalednn) methods for oscillatory stokes flows in complex domains. *Communications in Computational Physics*, 28 (5):2139–2157, 2020a.

Feng Wang, Alberto Eljarrat, Johannes Müller, Trond R Henninen, Rolf Erni, and Christoph T Koch. Multi-resolution convolutional neural networks for inverse problems. *Scientific reports*, 10(1):1–11, 2020b.

Zhi-Qin John Xu, Yaoyu Zhang, and Yanyang Xiao. Training behavior of deep neural network in frequency domain. *International Conference on Neural Information Processing*, pages 264–274, 2019.

Zhi-Qin John Xu, Yaoyu Zhang, Tao Luo, Yanyang Xiao, and Zheng Ma. Frequency principle: Fourier analysis sheds light on deep neural networks. *Communications in Computational Physics*, 28(5):1746–1767, 2020.

Yaoyu Zhang, Zhi-Qin John Xu, Tao Luo, and Zheng Ma. Explicitizing an implicit bias of the frequency principle in two-layer neural networks. *arXiv:1905.10264*, May 2019.

Yaoyu Zhang, Zhi-Qin John Xu, Tao Luo, and Zheng Ma. A type of generalization error induced by initialization in deep neural networks. In Jianfeng Lu and Rachel Ward, editors, *Proceedings of The First Mathematical and Scientific Machine Learning Conference*, volume 107, pages 144–164, 2020.

Xiaojin Zhu, Zoubin Ghahramani, and John D Lafferty. Semi-supervised learning using Gaussian fields and harmonic functions. In *Proceedings of the 20th International conference on Machine learning (ICML-03)*, pages 912–919, 2003.

## Appendix A. Lemma 2

**Lemma 2** *Let the function $\psi_\sigma(\boldsymbol{\xi}) = (2\pi)^{\frac{d}{2}}\sigma^d \mathrm{e}^{-2\pi^2\sigma^2\|\boldsymbol{\xi}\|^2}$, $\boldsymbol{\xi} \in \mathbb{R}^d$. We have*

$$\lim_{\sigma\to 0}\int_{\mathbb{R}^d}\|\boldsymbol{\xi}\|^\alpha |\psi_\sigma(\boldsymbol{\xi})|^2\,\mathrm{d}\boldsymbol{\xi} = \begin{cases} 0, & \alpha < d, \\ C_d, & \alpha = d, \\ \infty, & \alpha > d. \end{cases} \tag{19}$$

*Here the constant $C_d = \frac{1}{2}(d-1)!(2\pi)^{-d}\frac{2\pi^{d/2}}{\Gamma(d/2)}$ only depends on the dimension $d$.*

**Proof** In fact,

$$\lim_{\sigma\to 0}\int_{\mathbb{R}^d}\|\boldsymbol{\xi}\|^\alpha |\psi_\sigma(\boldsymbol{\xi})|^2\,\mathrm{d}\boldsymbol{\xi} = \lim_{\sigma\to 0}\int_{\mathbb{R}^d}\|\boldsymbol{\xi}\|^\alpha (2\pi)^d\sigma^{2d}\mathrm{e}^{-4\pi^2\sigma^2\|\boldsymbol{\xi}\|^2}\,\mathrm{d}\boldsymbol{\xi}$$

$$= \lim_{\sigma\to 0}(2\pi)^d\sigma^{d-\alpha}\int_{\mathbb{R}^d}\|\sigma\boldsymbol{\xi}\|^\alpha \mathrm{e}^{-4\pi^2\|\sigma\boldsymbol{\xi}\|^2}\,\mathrm{d}(\sigma\boldsymbol{\xi})$$

$$= \lim_{\sigma\to 0}(2\pi)^d\sigma^{d-\alpha}\int_0^\infty r^{\alpha+d-1}\mathrm{e}^{-4\pi^2 r^2}\,\mathrm{d}r \cdot \omega_d,$$

where $\omega_d = \frac{2\pi^{\frac{d}{2}}}{\Gamma\left(\frac{d}{2}\right)}$ is the surface area of a unit $(d-1)$-sphere.

Notice that

$$\int_0^\infty r^{\alpha+d-1}\mathrm{e}^{-4\pi^2 r^2}\,\mathrm{d}r = \int_0^1 r^{\alpha+d-1}\mathrm{e}^{-4\pi^2 r^2}\,\mathrm{d}r + \int_1^\infty r^{\alpha+d-1}\mathrm{e}^{-4\pi^2 r^2}\,\mathrm{d}r$$

$$\leq \int_0^\infty \mathrm{e}^{-4\pi^2 r^2}\,\mathrm{d}r + \int_0^\infty r^{[\alpha]+d}\mathrm{e}^{-4\pi^2 r^2}\,\mathrm{d}r$$

$$= \frac{1}{8\pi^{\frac{3}{2}}} + \int_0^\infty r^{[\alpha]+d}\mathrm{e}^{-4\pi^2 r^2}\,\mathrm{d}r$$

and

$$\int_0^\infty r^{[\alpha]+d} e^{-4\pi^2 r^2} \, dr = \begin{cases} \frac{1}{2}\left(\frac{[\alpha]+d-1}{2}\right)!(2\pi)^{-([\alpha]+d+1)}, & [\alpha]+d \text{ is odd}, \\ \frac{\sqrt{\pi}}{2}(2\pi)^{-([\alpha]+d+1)}\left(\frac{1}{2}\right)^{\frac{[\alpha]+d}{2}}([\alpha]+d-1)!!, & [\alpha]+d \text{ is even}. \end{cases}$$

Therefore, in both cases, the integral $\int_0^\infty r^{\alpha+d-1} e^{-4\pi^2 r^2} \, dr$ is finite. Then we have

$$\lim_{\sigma \to 0} \int_{\mathbb{R}^d} \|\boldsymbol{\xi}\|^\alpha |\psi_\sigma(\boldsymbol{\xi})|^2 \, d\boldsymbol{\xi} = \lim_{\sigma \to 0} (2\pi)^d \sigma^{d-\alpha} \int_0^\infty r^{\alpha+d-1} e^{-4\pi^2 r^2} \, dr \cdot \omega_d$$

$$= \begin{cases} 0, & \alpha < d, \\ \infty, & \alpha > d. \end{cases}$$

When $\alpha = d$, it follows that

$$\int_0^\infty r^{\alpha+d-1} e^{-4\pi^2 r^2} \, dr = \frac{1}{2}(2\pi)^{-2d}(d-1)!.$$

Therefore

$$\lim_{\sigma \to 0} \int_{\mathbb{R}^d} \|\boldsymbol{\xi}\|^\alpha |\psi_\sigma(\xi)|^2 \, d\xi = \frac{1}{2}(d-1)!(2\pi)^{-d} \frac{2\pi^{\frac{d}{2}}}{\Gamma\left(\frac{d}{2}\right)},$$

which completes the proof. ∎

## Appendix B. Proof of Proposition 1

**Proof** Similar to the proof of Lemma 2, we have

$$\lim_{\sigma \to 0} \|\mathcal{F}^{-1}[\psi_\sigma]\|_{H^{\frac{\alpha}{2}}}^2 = \lim_{\sigma \to 0} (2\pi)^d \sigma^{d-\alpha} \int_{\mathbb{R}^d} (\sigma^2 + \|\sigma\boldsymbol{\xi}\|^2)^{\frac{\alpha}{2}} e^{-4\pi^2 \|\sigma\boldsymbol{\xi}\|^2} \, d(\sigma\boldsymbol{\xi})$$

$$= \lim_{\sigma \to 0} (2\pi)^d \sigma^{d-\alpha} \int_0^\infty r^{d-1} (\sigma^2 + r^2)^{\frac{\alpha}{2}} e^{-4\pi^2 r^2} \, dr \cdot \omega_d.$$

For $\sigma < 1$, the following integrals are bounded from below and above, respectively:

$$\int_0^\infty r^{d-1} (\sigma^2 + r^2)^{\frac{\alpha}{2}} e^{-4\pi^2 r^2} \, dr \geq \int_0^\infty r^{\alpha+d-1} e^{-4\pi^2 r^2} \, dr = C_1 > 0,$$

and

$$\int_0^\infty r^{d-1}(\sigma^2 + r^2)^{\frac{\alpha}{2}} e^{-4\pi^2 r^2} \, dr \leq \int_0^1 r^{d-1}(1+r^2)^{\frac{\alpha}{2}} e^{-4\pi^2 r^2} \, dr + \int_1^\infty r^{d-1}((2r)^2)^{\frac{\alpha}{2}} e^{-4\pi^2 r^2} \, dr$$

$$\leq \int_0^1 r^{d-1}(1+r^2)^{\frac{\alpha}{2}} e^{-4\pi^2 r^2} \, dr + 2^\alpha \int_0^\infty r^{\alpha+d-1} e^{-4\pi^2 r^2} \, dr$$

$$= C_2 < \infty,$$

where $C_1 = \int_0^\infty r^{\alpha+d-1} e^{-4\pi^2 r^2} \, dr$ and $C_2 = \int_0^1 r^{d-1} (1+r^2)^{\frac{\alpha}{2}} e^{-4\pi^2 r^2} \, dr + 2^\alpha \int_0^\infty r^{\alpha+d-1} e^{-4\pi^2 r^2} \, dr$.
Therefore, we obtain the results for the subcritical ($\alpha < d$) and supercritical ($\alpha > d$) cases

$$
\lim_{\sigma \to 0} \|\mathcal{F}^{-1}[\psi_\sigma]\|_{H^{\frac{\alpha}{2}}}^2 = \lim_{\sigma \to 0} (2\pi)^d \sigma^{d-\alpha} \int_0^\infty r^{d-1} (\sigma^2 + r^2)^{\frac{\alpha}{2}} e^{-4\pi^2 r^2} \, dr \cdot \omega_d
$$

$$
= \begin{cases} 0, & \alpha < d, \\ \infty, & \alpha > d. \end{cases}
$$

For the critical case $\alpha = d$, we have

$$
\lim_{\sigma \to 0} \|\mathcal{F}^{-1}[\psi_\sigma]\|_{H^{\frac{\alpha}{2}}}^2
$$

$$
= \lim_{\sigma \to 0} (2\pi)^d \int_0^\infty r^{d-1} (\sigma^2 + r^2)^{\frac{\alpha}{2}} e^{-4\pi^2 r^2} \, dr \cdot \omega_d
$$

$$
= \lim_{\sigma \to 0} (2\pi)^d \int_0^\infty r^{2d-1} e^{-4\pi^2 r^2} \, dr \cdot \omega_d + \lim_{\sigma \to 0} \left[ \frac{\alpha}{2} (2\pi)^d \sigma^2 \int_0^\infty r^{2d-3} e^{-4\pi^2 r^2} \, dr \cdot \omega_d + o(\sigma^2) \right]
$$

$$
= \lim_{\sigma \to 0} (2\pi)^d \int_0^\infty r^{2d-1} e^{-4\pi^2 r^2} \, dr \cdot \omega_d
$$

$$
= \frac{1}{2} (d-1)! (2\pi)^{-d} \frac{2\pi^{\frac{d}{2}}}{\Gamma\left(\frac{d}{2}\right)}.
$$

Therefore the proposition holds. ∎

## Appendix C. Proof of Theorem 1

**Proof** Given $X = (x_1, \ldots, x_n)^\mathsf{T}$ and $Y = (y_1, \ldots, y_n)^\mathsf{T}$, let $A = \left( \exp\left(-\frac{\|x_j - x_i\|^2}{2\sigma^2}\right) \right)_{n \times n}$ be an $n \times n$ matrix. For sufficiently small $\sigma$, the matrix $A$ is diagonally dominant, and hence invertible. So the linear system $Ag^{(\sigma)} = Y$ has a solution $g^{(\sigma)} = \left( g_1^{(\sigma)}, g_2^{(\sigma)}, \cdots, g_n^{(\sigma)} \right)^\mathsf{T}$. Let

$$
\phi_\sigma(\boldsymbol{\xi}) = \sum_i g_i^{(\sigma)} e^{-2\pi i \boldsymbol{\xi}^\mathsf{T} x_i} \psi_\sigma(\boldsymbol{\xi}),
$$

where $\psi_\sigma(\boldsymbol{\xi}) = (2\pi)^{\frac{d}{2}} \sigma^d e^{-2\pi^2 \sigma^2 \|\boldsymbol{\xi}\|^2}$ satisfying $\mathcal{F}^{-1}[\psi_\sigma](x) = e^{-\frac{\|x\|^2}{2\sigma^2}}$. Thus

$$
\mathcal{F}^{-1}[\phi_\sigma](x) = \sum_i g_i^{(\sigma)} \mathcal{F}^{-1}[\psi_\sigma](x - x_i) = \sum_i g_i^{(\sigma)} e^{-\frac{\|x - x_i\|^2}{2\sigma^2}}.
$$

In particular, for all $i = 1, 2, \cdots, n$

$$
\mathcal{F}^{-1}[\phi_\sigma](x_i) = \sum_j g_j^{(\sigma)} e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}} = (Ag^{(\sigma)})_i = y_i.
$$

Therefore, $\phi_\sigma \in \mathcal{A}_{X,Y}$ for sufficiently small $\sigma > 0$.

According to the above discussion, we can construct a sequence $\{\phi_{\frac{1}{m}}\}_{m=M}^{\infty} \subset \mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}$, where $M$ is a sufficiently large positive integer to make the matrix $\boldsymbol{A}$ invertible. As Proposition 1 shows,

$$\lim_{m \to +\infty} \|\mathcal{F}^{-1}[\phi_{\frac{1}{m}}]\|_{H^{\frac{\alpha}{2}}}^2 = 0.$$

Now, suppose that there exists a solution to the Problem 1, denoted as $\phi^* \in \mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}$. By definition,

$$\|\mathcal{F}^{-1}[\phi^*]\|_{H^{\frac{\alpha}{2}}}^2 \leq \min_{\phi \in \mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}} \|\mathcal{F}^{-1}[\phi]\|_{H^{\frac{\alpha}{2}}}^2 \leq \lim_{m \to +\infty} \|\mathcal{F}^{-1}[\phi_{\frac{1}{m}}]\|_{H^{\frac{\alpha}{2}}}^2 = 0.$$

Therefore, $\phi^*(\boldsymbol{\xi}) \equiv 0$ and $\mathcal{P}_{\boldsymbol{X}}\phi^* = \boldsymbol{0}$, which contradicts to the restrictive condition $\mathcal{P}_{\boldsymbol{X}}\phi^* = \boldsymbol{Y}$ for the situation that $\boldsymbol{Y} \neq \boldsymbol{0}$. The proof is completed. ∎

## Appendix D. Proof of Theorem 2

**Proof** 1. We introduce a distance for functions $\phi, \psi \in L^2(\mathbb{R}^d)$:

$$\text{dist}(\phi, \psi) = \|\mathcal{F}^{-1}[\phi] - \mathcal{F}^{-1}[\psi]\|_{H^{\frac{\alpha}{2}}}.$$

Under the topology induced by this distance, the closure of the admissible function class $\mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}$ reads as

$$\overline{\mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}} := \overline{\{\phi \in L^1(\mathbb{R}^d) \cap L^2(\mathbb{R}^d) \mid \mathcal{P}_{\boldsymbol{X}}\phi = \boldsymbol{Y}\}}^{\text{dist}(\cdot,\cdot)}.$$

2. We will consider an auxiliary minimization problem: to find $\phi^*$ such that

$$\phi^* \in \arg \min_{\phi \in \overline{\mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}}} \|\mathcal{F}^{-1}[\phi]\|_{H^{\frac{\alpha}{2}}}. \tag{20}$$

Let $m := \inf_{\phi \in \overline{\mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}}} \|\mathcal{F}^{-1}[\phi]\|_{H^{\frac{\alpha}{2}}}$. According to the proof of Proposition 1 and Theorem 1, for a small enough $\sigma > 0$, the inverse Fourier transform of function

$$\phi_\sigma(\boldsymbol{\xi}) = \sum_i g_i^{(\sigma)} e^{-2\pi i \boldsymbol{\xi}^\mathsf{T} \boldsymbol{x}_i} \psi_\sigma(\boldsymbol{\xi})$$

has finite Sobolev norm $\|\mathcal{F}^{-1}[\phi_\sigma]\|_{H^{\frac{\alpha}{2}}} < \infty$, where $\psi_\sigma(\boldsymbol{\xi})$ satisfies $\mathcal{F}^{-1}[\psi_\sigma](\boldsymbol{x}) = e^{-\frac{\|\boldsymbol{x}\|^2}{2\sigma^2}}$, $\boldsymbol{A} = \left( \exp(-\frac{\|\boldsymbol{x}_j - \boldsymbol{x}_i\|^2}{2\sigma^2}) \right)_{n \times n}$ and $\boldsymbol{g}^{(\sigma)} = \left( g_1^{(\sigma)}, g_2^{(\sigma)}, \cdots, g_n^{(\sigma)} \right)^\mathsf{T} = \boldsymbol{A}^{-1}\boldsymbol{Y}$. Thus $m < +\infty$.

3. Choose a minimizing sequence $\{\bar{\phi}_k\}_{k=1}^{\infty} \subset \overline{\mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}}$ such that

$$\lim_{k \to \infty} \|\mathcal{F}^{-1}[\bar{\phi}_k]\|_{H^{\frac{\alpha}{2}}} = m.$$

By definition of the closure, there exists a function $\phi_k \in \mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}$ for each $k$ such that

$$\|\mathcal{F}^{-1}[\bar{\phi}_k] - \mathcal{F}^{-1}[\phi_k]\|_{H^{\frac{\alpha}{2}}} \leq \frac{1}{k}.$$

Therefore $\{\phi_k\}_{k=1}^{\infty} \subset \mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}$ is also a minimizing sequence, i.e.,

$$\lim_{k \to \infty} \|\mathcal{F}^{-1}[\phi_k]\|_{H^{\frac{\alpha}{2}}} = m.$$

Then $\{\mathcal{F}^{-1}[\phi_k]\}_{k=1}^{\infty}$ is bounded in the Sobolev space $H^{\frac{\alpha}{2}}(\mathbb{R}^d)$. Hence there exist a weakly convergent subsequence $\{\mathcal{F}^{-1}[\phi_{n_k}]\}_{k=1}^{\infty}$ and a function $\mathcal{F}^{-1}[\phi^*] \in H^{\frac{\alpha}{2}}(\mathbb{R}^d)$ such that

$$\mathcal{F}^{-1}[\phi_{n_k}] \rightharpoonup \mathcal{F}^{-1}[\phi^*] \quad \text{in } H^{\frac{\alpha}{2}}(\mathbb{R}^d) \text{ as } k \to \infty.$$

Note that

$$m = \inf_{\phi \in \overline{\mathcal{A}_{X,Y}}} \|\mathcal{F}^{-1}[\phi]\|_{H^{\frac{\alpha}{2}}} \leq \|\mathcal{F}^{-1}[\phi^*]\|_{H^{\frac{\alpha}{2}}} \leq \liminf_{\phi_{n_k}} \|\mathcal{F}^{-1}[\phi_{n_k}]\|_{H^{\frac{\alpha}{2}}} = m,$$

where we have used the lower semi-continuity of the Sobolev norm of $H^{\frac{\alpha}{2}}(\mathbb{R}^d)$ in the third inequality. Hence $\|\mathcal{F}^{-1}[\phi^*]\|_{H^{\frac{\alpha}{2}}} = m$.

4. We further establish the strong convergence that $\mathcal{F}^{-1}[\phi_{n_k}] - \mathcal{F}^{-1}[\phi^*] \to 0$ in $H^{\frac{\alpha}{2}}(\mathbb{R}^d)$ as $k \to \infty$. In fact, since $\mathcal{F}^{-1}[\phi_{n_k}] \rightharpoonup \mathcal{F}^{-1}[\phi^*]$ in $H^{\frac{\alpha}{2}}(\mathbb{R}^d)$ as $k \to \infty$ and $\lim_{k\to\infty}\|\mathcal{F}^{-1}[\phi_{n_k}]\|_{H^{\frac{\alpha}{2}}} = m = \|\mathcal{F}^{-1}[\phi^*]\|_{H^{\frac{\alpha}{2}}}$, we have

$$\lim_{k\to\infty} \|\mathcal{F}^{-1}[\phi_{n_k}] - \mathcal{F}^{-1}[\phi^*]\|_{H^{\frac{\alpha}{2}}}^2 = \lim_{k\to\infty} \langle \mathcal{F}^{-1}[\phi_{n_k}] - \mathcal{F}^{-1}[\phi^*], \mathcal{F}^{-1}[\phi_{n_k}] - \mathcal{F}^{-1}[\phi^*] \rangle$$

$$= \lim_{k\to\infty} \langle \mathcal{F}^{-1}[\phi_{n_k}], \mathcal{F}^{-1}[\phi_{n_k}] \rangle + \langle \mathcal{F}^{-1}[\phi^*], \mathcal{F}^{-1}[\phi^*] \rangle - \langle \mathcal{F}^{-1}[\phi_{n_k}], \mathcal{F}^{-1}[\phi^*] \rangle - \langle \mathcal{F}^{-1}[\phi^*], \mathcal{F}^{-1}[\phi_{n_k}] \rangle$$

$$= m^2 + m^2 - \lim_{k\to\infty} \left( \langle \mathcal{F}^{-1}[\phi_{n_k}], \mathcal{F}^{-1}[\phi^*] \rangle + \langle \mathcal{F}^{-1}[\phi^*], \mathcal{F}^{-1}[\phi_{n_k}] \rangle \right)$$

$$= m^2 + m^2 - \langle \mathcal{F}^{-1}[\phi^*], \mathcal{F}^{-1}[\phi^*] \rangle - \langle \mathcal{F}^{-1}[\phi^*], \mathcal{F}^{-1}[\phi^*] \rangle = 0.$$

Here $\langle \cdot, \cdot \rangle$ is the inner product of the Hilbert space $H^{\frac{\alpha}{2}}$.

5. We have $\phi^* \in L^1(\mathbb{R}^d)$ because

$$\int_{\mathbb{R}^d} |\phi^*(\boldsymbol{\xi})| \, d\boldsymbol{\xi} = \int_{\mathbb{R}^d} \frac{\langle \boldsymbol{\xi} \rangle^{\frac{\alpha}{2}} |\phi^*(\boldsymbol{\xi})|}{\langle \boldsymbol{\xi} \rangle^{\frac{\alpha}{2}}} \, d\boldsymbol{\xi} \leq \|\mathcal{F}^{-1}[\phi^*]\|_{H^{\frac{\alpha}{2}}} \left( \int_{\mathbb{R}^d} \frac{1}{\langle \boldsymbol{\xi} \rangle^{\alpha}} \, d\boldsymbol{\xi} \right)^{\frac{1}{2}} = Cm < +\infty,$$

where $C := \left( \int_{\mathbb{R}^d} \frac{1}{\langle \boldsymbol{\xi} \rangle^{\alpha}} \, d\boldsymbol{\xi} \right)^{\frac{1}{2}} < +\infty$. Hence $\phi^* \in L^1(\mathbb{R}^d) \cap L^2(\mathbb{R}^d)$ and $\mathcal{P}_X \phi^*$ is well-defined.

6. Recall that $\mathcal{P}_X \phi_{n_k} = Y$. We have

$$|Y - \mathcal{P}_X \phi^*| = \lim_{k\to+\infty} |\mathcal{P}_X \phi_{n_k} - \mathcal{P}_X \phi^*|$$

$$= \lim_{k\to+\infty} \left| \int_{\mathbb{R}^d} (\phi_{n_k} - \phi^*) e^{2\pi i \boldsymbol{x}\boldsymbol{\xi}} \, d\boldsymbol{\xi} \right|$$

$$= \lim_{k\to+\infty} \left| \int_{\mathbb{R}^d} \frac{\langle \boldsymbol{\xi} \rangle^{\frac{\alpha}{2}} (\phi_{n_k} - \phi^*)}{\langle \boldsymbol{\xi} \rangle^{\frac{\alpha}{2}}} e^{2\pi i \boldsymbol{x}\boldsymbol{\xi}} \, d\boldsymbol{\xi} \right|$$

$$\leq \lim_{k\to+\infty} \|\mathcal{F}^{-1}[\phi_{n_k}] - \mathcal{F}^{-1}[\phi^*]\|_{H^{\frac{\alpha}{2}}} \left( \int_{\mathbb{R}^d} \frac{|e^{2\pi i \boldsymbol{x}\boldsymbol{\xi}}|^2}{\langle \boldsymbol{\xi} \rangle^{\alpha}} \, d\boldsymbol{\xi} \right)^{\frac{1}{2}}$$

$$= C \lim_{k\to+\infty} \|\mathcal{F}^{-1}[\phi_{n_k}] - \mathcal{F}^{-1}[\phi^*]\|_{H^{\frac{\alpha}{2}}} = 0.$$

Hence $\mathcal{P}_X \phi^* = Y$ and $\phi^* \in \mathcal{A}_{X,Y}$.

7. Note that

$$m = \inf_{\phi \in \overline{\mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}}} \|\mathcal{F}^{-1}[\phi]\|_{H^{\frac{\alpha}{2}}} \leq \inf_{\phi \in \mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}} \|\mathcal{F}^{-1}[\phi]\|_{H^{\frac{\alpha}{2}}} \leq \|\mathcal{F}^{-1}[\phi^*]\|_{H^{\frac{\alpha}{2}}} = m.$$

This implies that $\inf_{\phi \in \mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}} \|\mathcal{F}^{-1}[\phi]\|_{H^{\frac{\alpha}{2}}} = m$ and $\phi^* \in \arg\min_{\phi \in \mathcal{A}_{\boldsymbol{X},\boldsymbol{Y}}} \|\mathcal{F}^{-1}[\phi]\|_{H^{\frac{\alpha}{2}}}$, which completes the proof. ∎