

A Generalized Weighted Optimization Method for Computational Learning and Inversion

Björn Engquist¹, Kui Ren², **Yunan Yang**³

¹Department of Mathematics and Oden Institute, The University of Texas at Austin, USA

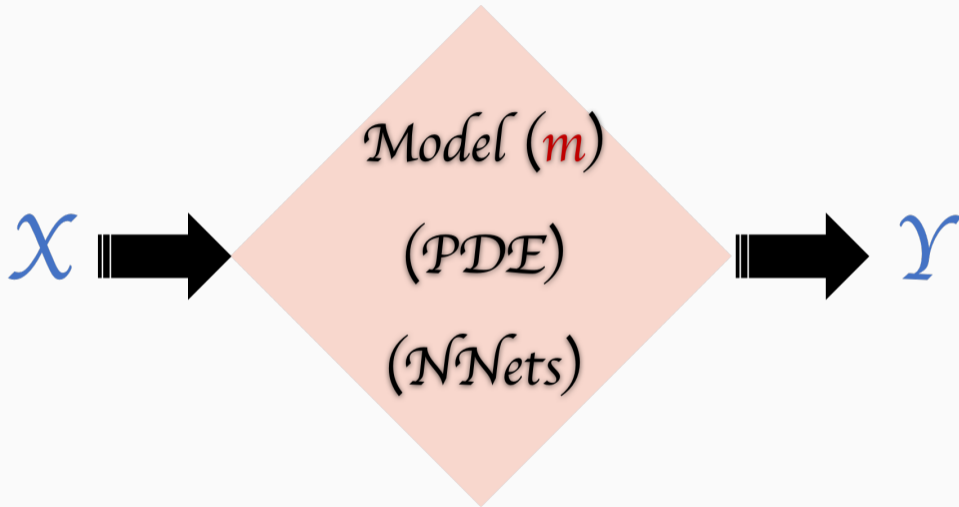
²Department of Applied Physics and Mathematics, Columbia University, USA

³Institute for Theoretical Studies, ETH Zürich, Switzerland



ICLR

General Inverse Problems



Aim at finding m such that the predicted outputs $(X, F(m))$ match given measured data (X, Y) .

Inverse Problem \rightarrow Inverse Data Matching Problem

Most of such inverse problems are solved as data-matching problems:

$$m^* = \operatorname{argmin}_{m \in (\mathcal{M}, d_m)} J(m) = \operatorname{argmin}_{m \in (\mathcal{M}, d_m)} d_g^2(f, g), \quad f = F(m), \quad f, g \in (\mathcal{D}, d_g).$$

(The map $f = F(m)$ is often *implicitly* given through a system of constraints.)

Two Metric Spaces

- (\mathcal{M}, d_m) is the metric space for the **parameters**.
- (\mathcal{D}, d_g) is the metric space for the **data**.

Inverse Problem \rightarrow Inverse Data Matching Problem

Most of such inverse problems are solved as data-matching problems:

$$m^* = \operatorname{argmin}_{m \in (\mathcal{M}, d_m)} J(m) = \operatorname{argmin}_{m \in (\mathcal{M}, d_m)} d_g^2(f, g), f = F(m), f, g \in (\mathcal{D}, d_g).$$

(The map $f = F(m)$ is often *implicitly* given through a system of constraints.)

Two Metric Spaces

- (\mathcal{M}, d_m) is the metric space for the **parameters**.
- (\mathcal{D}, d_g) is the metric space for the **data**.

How to choose these two spaces may affect the learning process greatly!

In this work, we propose to use **weighted** metrics for both spaces.

Zoom into Feature Regression

Given training data $\{x_j, y_j\}_{j=1}^N$, where $x_j \in \mathbb{R}$, $y_j \in \mathbb{C}$, we are interested in learning a random Fourier feature (RFF) model

$$f_{\theta}(x) = \sum_{k=0}^{P-1} \theta_k e^{ikx}, \quad x \in [0, 2\pi],$$

where we aim to find $\theta := (\theta_0, \dots, \theta_{P-1})^T$.

Zoom into Feature Regression

Given training data $\{x_j, y_j\}_{j=1}^N$, where $x_j \in \mathbb{R}$, $y_j \in \mathbb{C}$, we are interested in learning a random Fourier feature (RFF) model

$$f_{\theta}(x) = \sum_{k=0}^{P-1} \theta_k e^{ikx}, \quad x \in [0, 2\pi],$$

where we aim to find $\theta := (\theta_0, \dots, \theta_{P-1})^T$.

Let $\Psi \in \mathbb{C}^{N \times P}$ be the feature matrix

$$(\Psi)_{jk} = e^{ikx_j}, \quad 0 \leq j \leq N-1, \quad 0 \leq k \leq P-1.$$

The learning problem is recast as an optimization problem

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \|\Psi\theta - \mathbf{y}\|_2^2, \quad \mathbf{y} = (y_0, \dots, y_{N-1})^T.$$

Previous Work on Weighted Optimization [Xie et al., 2020]

Prior assumption about the coefficient vector θ :

$$\mathbb{E}_{\theta}[\theta] = \mathbf{0}, \quad \mathbb{E}_{\theta}[\theta\theta^*] = c_{\gamma}\mathbf{\Lambda}_{[P]}^{-2\gamma}, \quad \gamma > 0$$

where $\mathbf{\Lambda}_{[P]} = \text{diag}\{t_0, \dots, t_k, \dots, t_{P-1}\}$, $t_k := 1 + k$.

The normalization constant c_{γ} enforces $\mathbb{E}_{\theta}[\|\theta\|^2] = 1$.

This is assuming that the forward model is smoothing!

Previous Work on Weighted Optimization [Xie et al., 2020]

Prior assumption about the coefficient vector θ :

$$\mathbb{E}_{\theta}[\theta] = \mathbf{0}, \quad \mathbb{E}_{\theta}[\theta\theta^*] = c_{\gamma}\mathbf{\Lambda}_{[p]}^{-2\gamma}, \quad \gamma > 0$$

where $\mathbf{\Lambda}_{[p]} = \text{diag}\{t_0, \dots, t_k, \dots, t_{p-1}\}$, $t_k := 1 + k$.

The normalization constant c_{γ} enforces $\mathbb{E}_{\theta}[\|\theta\|^2] = 1$.

This is assuming that the forward model is smoothing!

[Xie et al., 2020] proposed to learn a model with $p \leq P$ features through the *weighted* least-squares formulation

$$\hat{\theta}_p = \mathbf{\Lambda}_{[p]}^{-\beta}\hat{\mathbf{w}}, \quad \text{with } \hat{\mathbf{w}} = \text{argmin}_{\theta} \|\Psi_{[N \times p]}\mathbf{\Lambda}_{[p]}^{-\beta}\mathbf{w} - \mathbf{y}\|_2^2,$$

when the learning problem is overparameterized, i.e., $p > N$.

Main Proposal of Our Work

A **generalized weighted** least-squares formulation for feature regression:

$$\hat{\theta}_p^\delta = \Lambda_{[p]}^{-\beta} \hat{\mathbf{w}}, \quad \text{with } \hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \|\Lambda_{[N]}^{-\alpha} \left(\Psi_{[N \times p]} \Lambda_{[p]}^{-\beta} \mathbf{w} - \mathbf{y}^\delta \right)\|_2^2,$$

where δ denotes that the training data is polluted with noise.

Main Proposal of Our Work

A **generalized weighted** least-squares formulation for feature regression:

$$\hat{\theta}_\rho^\delta = \Lambda_{[p]}^{-\beta} \hat{\mathbf{w}}, \quad \text{with } \hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \|\Lambda_{[N]}^{-\alpha} \left(\Psi_{[N \times p]} \Lambda_{[p]}^{-\beta} \mathbf{w} - \mathbf{y}^\delta \right)\|_2^2,$$

where δ denotes that the training data is polluted with noise.

A short informal summary:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \|\Psi \theta - \mathbf{y}\|_2^2, \quad (\text{unweighted}) \tag{1}$$

$$\hat{\mathbf{w}} = \underset{\theta}{\operatorname{argmin}} \|\Psi \Lambda_{[p]}^{-\beta} \mathbf{w} - \mathbf{y}\|_2^2, \quad ([\text{Xie et al., 2020}]) \tag{2}$$

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \|\Lambda_{[N]}^{-\alpha} \left(\Psi \Lambda_{[p]}^{-\beta} \mathbf{w} - \mathbf{y} \right)\|_2^2, \quad (\text{our proposal}) \tag{3}$$

$\hat{\theta}_\rho = \Lambda_{[p]}^{-\beta} \hat{\mathbf{w}}$ for the last two.

The Benefits of Weighting

$$\hat{\theta}_p^\delta = \Lambda_{[p]}^{-\beta} \hat{\mathbf{w}}, \quad \text{with } \hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \left\| \Lambda_{[N]}^{-\alpha} \left(\Psi_{[N \times p]} \Lambda_{[p]}^{-\beta} \mathbf{w} - \mathbf{y}^\delta \right) \right\|_2^2$$

Define the generalization error (thanks to RFF):

$$\mathcal{E}_{\alpha, \beta}^\delta(P, p, N) := \mathbb{E}_\theta \left[\left\| f_\theta(\mathbf{x}) - f_{\hat{\theta}_p^\delta}(\mathbf{x}) \right\|_{L^2([0, 2\pi])}^2 \right] = \mathbb{E}_\theta \left[\left\| \hat{\theta}_p^\delta - \theta \right\|_2^2 \right].$$

The Benefits of Weighting

$$\widehat{\boldsymbol{\theta}}_p^\delta = \boldsymbol{\Lambda}_{[p]}^{-\beta} \widehat{\mathbf{w}}, \quad \text{with } \widehat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \left\| \boldsymbol{\Lambda}_{[N]}^{-\alpha} \left(\boldsymbol{\Psi}_{[N \times p]} \boldsymbol{\Lambda}_{[p]}^{-\beta} \mathbf{w} - \mathbf{y}^\delta \right) \right\|_2^2$$

Define the generalization error (thanks to RFF):

$$\mathcal{E}_{\alpha, \beta}^\delta(P, p, N) := \mathbb{E}_\theta \left[\left\| f_\theta(\mathbf{x}) - f_{\widehat{\boldsymbol{\theta}}_p^\delta}(\mathbf{x}) \right\|_{L^2([0, 2\pi])}^2 \right] = \mathbb{E}_\theta \left[\left\| \widehat{\boldsymbol{\theta}}_p^\delta - \boldsymbol{\theta} \right\|_2^2 \right].$$

In this paper, we analyze the **improved** generalization error in three cases:

1. Training with noise-free data in RFF model (Section 3)
2. Training with noisy data in RFF model (Section 4)
3. Beyond Random Fourier Feature (RFF) model (Section 5)

Summary

We solve an inverse problem through the optimization format.

$$m^* = \underset{m \in (\mathcal{M}, d_m)}{\operatorname{argmin}} J(m) = \underset{m \in (\mathcal{M}, d_m)}{\operatorname{argmin}} d_g^2(f, g), f = F(m), f, g \in (\mathcal{D}, d_g).$$

Two Metric Spaces

- (\mathcal{M}, d_m) is the metric space for the **parameters**.
- (\mathcal{D}, d_g) is the metric space for the **data**.

Two Weight Matrices

$$\hat{\theta}_p^\delta = \Lambda_{[p]}^{-\beta} \hat{\mathbf{w}}, \text{ with } \hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \|\Lambda_{[N]}^{-\alpha} \left(\Psi_{[N \times p]} \Lambda_{[p]}^{-\beta} \mathbf{w} - \mathbf{y}^\delta \right)\|_2^2.$$

The Impact of the Two Metric Spaces/Weighting Matrices

1. Improve Generalization Error/Resolution
2. Robustness to Noise (Mitigate Overfitting)
3. Improve Optimization Landscape
4. Change Convergence Speed/Trajectory