

MEASURING FREQUENCY NON-STATIONARITY FOR ROBUST TIME SERIES FORECASTING

Anonymous authors

Paper under double-blind review

ABSTRACT

Non-stationarity in time series has long posed a fundamental challenge for forecasting models, as it leads to distribution shifts between training and test data. A popular line of research, known as normalization methods, aims to measure and suppress non-stationarity by removing time-domain low-order statistics. Nevertheless, low-order statistics may inadequately address the underlying non-stationary structures manifested as a composition of frequencies. To tackle these issues, we propose to measure the degree of stationarity of each frequency component across distributions via spectral analysis. By identifying and downweighting frequencies that are more non-stationary, we re-represent the original time series to reduce distributional discrepancies between training and test sets. Concretely, we present FREMEN with threefold contributions. *Theoretically*, FREMEN is grounded in a principled formulation and we provide the first spectral analysis to support its validity. *Technically*, FREMEN is both novel and effective, incurring negligible additional computational cost. *Experimentally*, FREMEN is validated on four forecasting models across seven datasets, achieving 24 best results out of 28 settings and 28.4% average MSE improvements. Our code is publicly available¹.

1 INTRODUCTION

Time series forecasting is vital to decision-making in real-world applications like industrial system control and stock market tracking (Thompson & Wilson, 2016; Zhao et al., 2024). Recently, deep learning has shown some promise on benchmark datasets (Nie et al., 2023; Liu et al., 2024; Piao et al., 2024b; Wang et al., 2025). However, a challenge remains: *the non-stationary nature of time series such as seasonal fluctuations and irregular events often leads to poor generalization when forecasting models are applied to unseen test data*. The non-stationarity baffles training-patterns-driven forecasting models that assume consistency in the test dataset. Therefore, when the distribution shift occurs, these models show forecasting degeneration.

To tackle this issue, a recent popular line of research focuses on normalization methods that aim to *measure* and suppress non-stationarity in the input samples, thereby reducing distributional discrepancies between the training and test datasets (Kim et al., 2021; Liu et al., 2023; Fan et al., 2023; Han et al., 2024; Ye et al., 2024). These methods explicitly measure non-stationarity through statistical metrics (typically mean and variance) computed or learned from the training set. The metrics are then used to normalize the input, attempting to remove distribution shifts manifested in the location and scale. Since the normalization is applied consistently during both training and inference, it helps align the data distributions and thus improves generalization. Importantly, many of these methods learn to control the normalization strength through adaptive gates (Fan et al., 2023) or residual connections (Liu et al., 2023), balancing between preserving the original distributional information and measuring statistical stability for more robust forecasting.

Nevertheless, existing methods estimate statistics in the time domain, which may inadequately measure the non-stationarity due to the following reasons. First, these methods primarily rely on normalizing raw temporal values. While it may be effective for simple scale and temporal variations, it often fails to account for more complex *non-stationary structures* in the frequency domain, such as temporal drift in dominant frequencies, spectral reallocation, or shifts in periodicity (Ye et al., 2024;

¹<https://anonymous.4open.science/r/Fremen-code-82C8>

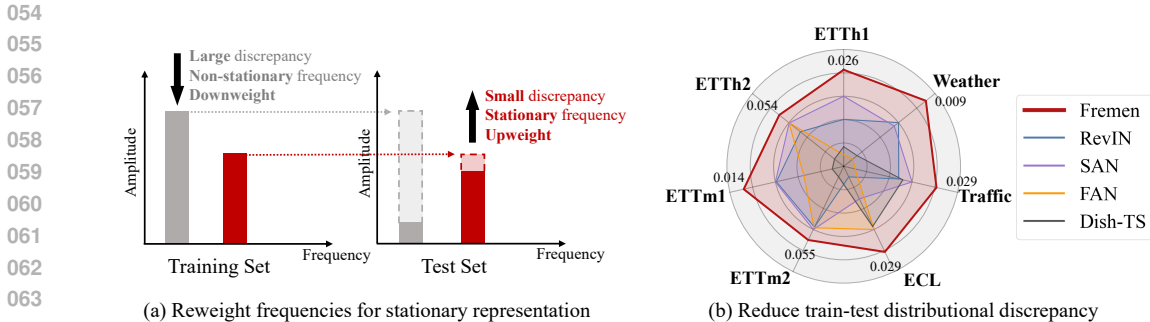


Figure 1: (a) The objective of this paper is to measure stationarity of frequencies and reweight them accordingly for stationary representation. (b) The learnt representation achieves the lowest Jensen-Shannon divergence between the training and test samples compared to existing methods. A larger shaded area indicates a smaller distributional discrepancy between the training and test set.

Piao et al., 2024b). As a result, these methods struggle to mitigate distribution shifts that consist of the underlying spectra of time series data. Second, the use of low-order statistics is insufficient to characterize complex distributional structures and intricate temporal dependencies (e.g., multimodality, higher-order statistics, or changes in its functional form) (Han et al., 2024). Consequently, the normalization becomes inaccurate, hindering the effectiveness in measuring distribution shifts.

To address these issues, this paper presents FREMEN, a frequency-space, non-stationarity-aware method to mitigate distribution shifts in time series forecasting. As shown in Figure 1, our idea is to perform spectral analysis to measure the degree of stationarity of each frequency component across samples. By identifying and downweighting frequency components that are more non-stationary, we re-represent the original time series to reduce distributional discrepancies between training and test sets. While a recent study also attempts to analyze non-stationarity in the frequency domain (Ye et al., 2024). It heuristically selects the top- k frequency magnitudes, running the risk of low-frequency dominance and inadequately characterizing the entire spectrum. By contrast, we introduce a kernel representation that is implicitly induced by the Fourier transform, integrating out all possible distribution shift patterns via the Yaglom’s theorem (Yaglom, 1987). By the one-one correspondence between kernel and spectral weights, learning the weights of frequencies is equivalent to learning a data-adaptive kernel representation itself, allowing the model to capture subtle distributional discrepancies and prioritize stationary frequency components for improved generalization.

Our contributions are threefold. Theoretically, we provide the first spectral analysis of frequency-domain non-stationarity to combat distribution shift. Methodologically, we present a simple yet effective algorithm for learning to weight the frequency components, thereby re-representing the distribution behind time series data. Experimentally, FREMEN is validated on four mainstream forecasting backbone models across seven benchmark datasets, achieving 24 best results out of 28 settings and 28.4% average MSE improvements in multivariate forecasting.

2 RELATED WORKS

Non-stationary Time Series Modeling. Existing methods mainly aim to find a way to measure the distribution shift in the time domain, thereby helping the model learn a robust data representation (Du et al., 2021; Kim et al., 2021; Piao et al., 2024a). RevIN (Kim et al., 2021) proposed using mean and variance to measure the distribution shift. They first set the mean and variance of each sample to a fixed value. After forecasting, the original values are returned to the forecasting model outputs. A series of time-domain methods followed this idea: (i) Earlier works consider the evolution of mean and variance between inputs and outputs and explicitly model them (Fan et al., 2023; Liu et al., 2023). (ii) Recent works tend to learn a more expressive measure of non-stationarity than low-order statistics (Han et al., 2024; Liu et al., 2023). Recently, FAN (Ye et al., 2024) took an initial step toward addressing distribution shifts in the frequency domain. FAN identifies high-amplitude frequencies as unstable and seeks to mitigate distribution shifts by removing these frequencies. However, this heuristic strategy may risk discarding critical patterns and lead to suboptimal performance.

Frequency Domain Modeling. Non-stationary time series can be seen as a mix of frequencies that vary over time (Proakis & Manolakis, 1996). Earlier methods often aim to learn frequency features directly from the raw Fourier coefficients (Wu et al., 2021; Zhou et al., 2022b; Wang et al., 2022; Wu et al., 2023; Yi et al., 2023). However, the frequency features are often sparse and mixed with noise and time-varying features (Proakis & Manolakis, 1996; Piao et al., 2024b). Recent methods tend to learn more informative and robust representations via sparse selection (Zhou et al., 2022b; Woo et al., 2022; Zhou et al., 2022a; Ye et al., 2024) or normalization (Piao et al., 2024b). However, these methods often rely on heuristic strategies and do not consider cross-sample variations. To the best of our knowledge, we are the first to model the frequency variations across samples to mitigate the impact of non-stationary features on forecasting.

3 PROBLEM SETTING AND PRELIMINARY ANALYSIS

We formulate the problem of distribution shifts in non-stationary time series in Section 3.1, followed by our novel theoretical analysis in Section 3.2. Based on the analysis, we present a novel forecasting method FREMEN in Section 4. The key notations used in the paper are summarized in Table 1.

3.1 PROBLEM SETUP

Definition 1 (Time series data and forecasting.) We consider the multivariate time series forecasting problem on a given dataset $\{\mathcal{X}, \mathcal{Y}\}$, with $\mathcal{X} = \{x^{(i)}\}_{i=1}^N$, $\mathcal{Y} = \{y^{(i)}\}_{i=1}^N$ and N denotes the number of sequences. Let C, L_x, L_y respectively denote the number of variables, the input-sequence length and the model prediction length, then the goal can be formulated as that given an input time series data $x^{(i)} \in \mathbb{R}^{L_x \times C}$, predict the target values $y^{(i)} \in \mathbb{R}^{L_y \times C}$.

Definition 2 (Distribution shift issue in forecasting.) We consider the forecasting under distribution shift issue induced by non-stationarity in time series data. We assume that the data \mathcal{X} is generated from an evolving distribution over time $P_t(x)$. A time series is said to be *stationary* if its distribution remains invariant over time, i.e., $P_{t_1}(x) = P_{t_2}(x)$ for all t_1, t_2 . Conversely, *non-stationarity* refers to scenarios where the distribution changes with time: $\exists t_1 \neq t_2$ s.t. $P_{t_1}(x) \neq P_{t_2}(x)$. Such distribution shifts can manifest through changes in the mean, variance, feature correlation, or other latent structure of the input sequences. Formally, given a training set $\mathcal{X}_{\text{train}} = \{x^{(i)}\}_{i=1}^{N_{\text{train}}}$ drawn from $P_t(x)$ with $t \in \mathbb{T}_{\text{train}}$, the goal is to make accurate predictions on future inputs x drawn from a different distribution $P_{t'}(x)$ with $t' \in \mathbb{T}_{\text{test}}, \mathbb{T}_{\text{test}} \cap \mathbb{T}_{\text{train}} = \emptyset$, and $P_{t'}(x) \neq P_t(x)$.

Table 1: Key notations used in this paper.

| Notation | Description |
|----------------------------|--|
| \mathcal{X}, \mathcal{Y} | input and target time series |
| L_x, L_y | input sequence length, prediction length |
| N | number of variables |
| C | number of channels |
| P_t | time series generation distribution |
| t | sample index of time series |
| ω | frequency component |
| \mathcal{S} | power spectral density |
| k | valid kernel |
| λ | eigenvalue, frequency weight |
| F | frequency coefficient |

Problem (Statistical non-stationarity measure.) A common paradigm is the use of statistical normalization techniques applied directly to the input sequences. These methods normalize the observations across the temporal dimension before feeding into the model. Formally, it computes channel-wise μ_t, σ_t (e.g., mean and standard deviation at time t), and transforms x into a normalized time series \tilde{x} : $\tilde{x}_{t,c} = \frac{x_{t,c} - \mu_{t,c}}{\sigma_{t,c}}, \forall t \in [1, L_x], c \in [1, C]$. μ_t, σ_t can be empirically computed (Kim et al., 2021), learned (Fan et al., 2023), or vectorized using sliding windows (Liu et al., 2023).

However, this paradigm implicitly assumes that the data generating distribution $P_t(x)$ can be fully characterized by its low-order statistics. This assumption rules out the possibility of more complex distributions. Even in the location-scale family, members like the Student’s t distribution depend on additional parameters (Zhu et al., 2025). Existing normalization methods thus fail to reflect complex non-stationary patterns like frequency shift, temporal dynamics, or latent structural changes.

Therefore, *our goal is to develop more expressive, learnable measures of non-stationarity* that can adaptively characterize evolving dynamics in the input time series.

3.2 PRELIMINARY ANALYSIS: MEASURING NON-STATIONARITY IN FREQUENCIES

We investigate non-stationarity in time series from a frequency-domain perspective. The Fourier transform decomposes a time series into basis functions, disentangling temporal structure into interpretable frequency bands. Our theoretical analysis first shows that variations in the power spectrum reflect underlying distribution shifts (Lemma 1), and that such spectral differences provide a valid non-stationarity measure (Lemma 2). We further prove that modeling these differences induces a shift-invariant kernel in the frequency domain, offering a principled way to emphasize stationary components for robust forecasting (Lemma 3). We begin by formally stating our main theorem.

Theorem 1 *The Fourier transform on the timeseries dataset \mathcal{X} induces a similarity measure k that is invariant to the non-stationarity. This measure k can be learned in a data-driven manner by learning its frequency weights $\{\lambda_i\}$.*

To support this theorem, we present the following lemmas, which respectively establishes that: (i) spectral representations encode key differences to non-stationarity not apparent in the time domain; (ii) a measure capable of gauging the differences in the spectral domain holds the potential of distinguishing distribution shifts; (iii) kernel function is a valid measure that can be adapted to data by identifying its eigenvalues.

Lemma 1 (Spectral shift and energy redistribution) *Let x_t and x'_t be two sampled non-stationary time series, where the distribution shifts over time. Then their spectral representations \hat{x} and \hat{x}' exhibit distinct energy distributions across frequency bands:*

$$\exists \omega \text{ s.t. } |\hat{x}(\omega)|^2 \neq |\hat{x}'(\omega)|^2.$$

This redistribution of spectral energy reflects the underlying non-stationary behavior (e.g., seasonal transitions, structural drifts), which may not be apparent in the time domain.

Thus, power spectral density (PSD) analysis may provide a principled way to quantify time-varying distributions. We next formalize the discriminative capability of these spectral patterns, thereby validating the use of frequency-domain representations as a principled measure of non-stationarity.

Lemma 2 (Discriminative power of spectral distribution) *Let $\mathcal{X}_1, \mathcal{X}_2$ be two subsets of sequences drawn from distributions P_{t_1} and P_{t_2} respectively, with $P_{t_1} \neq P_{t_2}$. Assume their average power spectral densities are $\mathcal{S}_1(\omega)$ and $\mathcal{S}_2(\omega)$. Then the total variation distance between them satisfies:*

$$\text{TV}(\mathcal{S}_1, \mathcal{S}_2) = \frac{1}{2} \int |\mathcal{S}_1(\omega) - \mathcal{S}_2(\omega)| d\omega > 0.$$

This implies that frequency-domain statistics can effectively distinguish different time-evolving distributions, and thus serve as a valid non-stationarity measure.

Having established that spectral differences can measure non-stationarity, we next explore how to learn a function to model the difference in a principled way. We turn to spectral analysis to show that a valid similarity measure as a kernel function is implicitly induced by Fourier transform.

Lemma 3 (Yaglom’s Theorem) *A continuous bounded function k on \mathbb{R}^{L_x} is a valid kernel if and only if it can be represented as*

$$k(x_1, x_2) = \int_{\mathbb{R}^{L_x} \times \mathbb{R}^{L_x}} e^{2\pi i(\omega_1^\top x_1 - \omega_2^\top x_2)} \mathcal{S}(\omega_1, \omega_2) d\omega_1 d\omega_2$$

where $\mathcal{S}(\omega_1, \omega_2)$ can be understood as a joint probability density function (Yaglom, 1987). Because a kernel can be fully characterized by its eigen-decomposition (Scholkopf & Smola, 2001), Yaglom’s theorem indicates that the measure k induced by Fourier transform on data can be adapted to data by learning its eigenvalues $\{\lambda_i\}$.

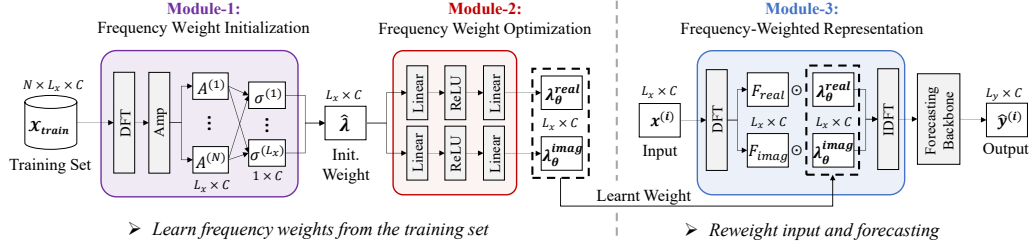


Figure 2: FREMEN begins by initializing frequency weights from the training set and refining them through neural network optimization (*Module-1,2*). Then, given an input, the weights are applied to produce a frequency-weighted input for the forecasting backbone to perform prediction (*Module-3*).

As time-sequential data can be derived from the integration of harmonic waves, the Yaglom’s theorem implies that a kernel as the result of integrating over the distribution of power spectra, is invariant to the time-varying statistical characteristics of time series (Xue et al., 2023).

In summary, our theorem and the lemmas manifest a key point: when we move to the frequency domain, the differences in power spectra directly capture non-stationarity. This means we can think of the kernel as a mathematical tool that measures similarity between time series in the frequency domain. Frequency components capture patterns at different scales, and kernels provide a principled way to weight them by their stability across training and test data. This naturally emphasizes stationary frequencies while suppressing non-stationary ones. Our goal is to identify relatively stationary frequency weights to mitigate distribution shift in non-stationary time series forecasting.

4 PROPOSED METHOD: FREMEN

Based on our analysis in Section 3.2, we present a novel forecasting framework FREMEN in this section. FREMEN can be employed as a representation layer that reweights frequency components of the input to produce non-stationarity-aware features.

4.1 OVERALL FORWARD PROCESS

The forward process is summarized in Figure 2. FREMEN first learns frequency weights from the training set $\mathcal{X}_{\text{train}}$. Specifically, the *frequency weight initialization* (*Module-1*) takes samples in $\mathcal{X}_{\text{train}}$ as input, outputs the empirical kernel eigenvalue $\hat{\lambda}$ as initial frequency weights. Then, $\hat{\lambda}$ is fed to the neural network in the *frequency weight optimization* (*Module-2*), producing weights for real and imaginary parts, i.e., $\lambda_{\theta}^{\text{real}}$ and $\lambda_{\theta}^{\text{imag}}$. Finally, when a new input x arrives, the *frequency-weighted representation* (*Module-3*) transforms it into the frequency domain, applies $\lambda_{\theta}^{\text{real}}$ and $\lambda_{\theta}^{\text{imag}}$ to the corresponding frequency components, and transforms it back to the time domain to obtain a weighted representation \tilde{x} , serving as the input for the forecasting model.

4.2 MODULE-1: FREQUENCY WEIGHT INITIALIZATION

As discussed in the preliminary analysis, frequency weights have one-one correspondence with kernels. To steer the learning process, we assume the commonly adopted RBF kernel (i.e., $\exp(-\gamma\|x_1 - x_2\|^2)$, where $\gamma := \frac{1}{2\sigma^2} > 0$ denotes the kernel width) to initialize the non-stationarity measure. Then, the corresponding eigenvalues of RBF kernel naturally serve as the starting point of frequency weights, which can be empirically estimated from the training set. Given $\mathcal{X}_{\text{train}} = \{x^{(i)}\}_{i=1}^N$, we first apply Discrete Fourier Transform (DFT) on each time series sample to obtain the amplitude spectrum $A^{(i)} = \text{Amp}(\text{DFT}(x^{(i)})) \in \mathbb{R}^{L_x \times C}$, where $\text{Amp}(\cdot)$ computes the amplitude. The frequency-wise RBF eigenvalue $\hat{\lambda} \in \mathbb{R}^{L_x \times C}$ is then measured by:

$$\hat{\lambda}(\omega) = \sqrt{\frac{\pi}{\hat{\gamma}}} \exp\left(-\frac{\omega^2}{4\hat{\gamma}}\right), \quad \text{where } \hat{\gamma} = \frac{1}{2\sigma^2(\{A^{(i)}\}_{i=1}^N)}. \quad (1)$$

Here, $\hat{\gamma}$ is the only value to be estimated. Frequency-wise standard deviation $\sigma(\cdot)$ is computed for frequency component ω over amplitudes of all training samples.

4.3 MODULE-2: FREQUENCY WEIGHT OPTIMIZATION

Following the initialization with RBF, one might consider learning the width γ for the frequency weights in a data-adaptive way. However, doing so is restricted to the RBF kernel as its functional form remains fixed. By contrast, learning adaptive weights λ_θ corresponds to learning diverse kernel representations itself, which has greater expressiveness than the fixed RBF. Therefore, we learn the frequency weights λ_θ from the RBF initialization by updating the neural networks with stochastic gradient descent. Specifically, we implement Multi-Layer Perceptrons (MLPs) to optimize frequency weights:

$$\begin{aligned}\lambda_\theta^{\text{real}}(\omega) &= \text{ReLU}(\hat{\lambda}(\omega)\mathbf{W}_1^{\text{real}} + \mathbf{b}_1^{\text{real}})\mathbf{W}_2^{\text{real}} + \mathbf{b}_2^{\text{real}}, \\ \lambda_\theta^{\text{imag}}(\omega) &= \text{ReLU}(\hat{\lambda}(\omega)\mathbf{W}_1^{\text{imag}} + \mathbf{b}_1^{\text{imag}})\mathbf{W}_2^{\text{imag}} + \mathbf{b}_2^{\text{imag}}.\end{aligned}\quad (2)$$

Here $\mathbf{W}_1^* \in \mathbb{R}^{C \times H}$, $\mathbf{W}_2^* \in \mathbb{R}^{H \times C}$, $\mathbf{b}_1^* \in \mathbb{R}^H$, and $\mathbf{b}_2^* \in \mathbb{R}^C$, where $*$ \in {real, imag}. H is the hidden dimension. Notably, we model the weights for real and imaginary frequency components separately as $\lambda_\theta^{\text{real}}$ and $\lambda_\theta^{\text{imag}} \in \mathbb{R}^{L_x \times C}$. The rationale behind this design is that assigning distinct weights to the real and imaginary components allows the model to capture non-stationarity arising from shifts in both amplitude and phase. In contrast, a single unified weight per frequency can only modulate amplitude, leaving phase-related non-stationarity structures unaddressed.

4.4 MODULE-3: FREQUENCY-WEIGHTED REPRESENTATION

Given an input time series $x \in \mathbb{R}^{L_x \times C}$, we first transform it to the frequency domain via DFT, producing real and imaginary coefficients as $F_{\text{real}}, F_{\text{imag}} = \text{DFT}(x)$. Then, the learned frequency weights are applied to the corresponding coefficient using Hadamard product:

$$\tilde{F}_{\text{real}}(\omega) = F_{\text{real}}(\omega) \odot \lambda_\theta^{\text{real}}(\omega), \quad \tilde{F}_{\text{imag}}(\omega) = F_{\text{imag}}(\omega) \odot \lambda_\theta^{\text{imag}}(\omega).\quad (3)$$

The weighted coefficients, i.e., \tilde{F}_{real} and \tilde{F}_{imag} , are supposed to establish a more stationary representation with enhanced stationary frequencies and suppressed non-stationary ones, accommodating robust forecasting under distribution shifts. By aggregating all weighted frequencies, the final time-domain representation is obtained via Inverse Discrete Fourier Transform (IDFT) as $\tilde{x} = \text{IDFT}(\tilde{F}_{\text{real}} + i\tilde{F}_{\text{imag}}) \in \mathbb{R}^{L_x \times C}$, serving as the input for the downstream forecasting backbone model to perform prediction. The whole framework is trained jointly with the forecasting backbone using mean squared error (MSE) loss in an end-to-end manner.

5 EXPERIMENTS

We conduct various experiments on widely adopted benchmark datasets to answer the following questions: **RQ1**: How does FREMEN enhance the performance of existing time series forecasting backbone models? **RQ2**: Does FREMEN mitigate distribution shift issue? **RQ3**: How effective are the learned frequency weights λ_θ ? **RQ4**: How does each design choice of FREMEN contribute to its performance? **RQ5**: How does the inclusion of FREMEN affect the efficiency of backbone models?

5.1 EXPERIMENT SETUP

Datasets. We use seven widely adopted datasets in multivariate time series forecasting, including: (1) **ETT** (Electricity Transformer Temperature) with four subsets of oil temperature and electrical load recorded at hourly (ETTh1, ETTh2) and 15-minute (ETTM1, ETTM2) resolutions from July 2016 to July 2018; (2) **ECL** contains 15-minute-level electricity consumption of 321 clients from 2012 to 2014. (3) **Weather** includes 21 meteorological features collected every 10 minutes in 2020. (4) **Traffic** is comprised of hourly-recorded traffic load by 862 sensors in San Francisco freeways from 2015 to 2016. All datasets have been published in (Wu et al., 2021). We adopt the split ratio setting in (Wu et al., 2021), which is 6:2:2 for four ETT datasets and 7:1:2 for the other datasets. A global normalization is applied to transform the whole dataset to a fixed scale. Note that this normalization keeps the statistics unchanged; thus, it is unable to handle non-stationarity.

Baselines. We compare FREMEN with state-of-the-art normalization methods for non-stationary time series forecasting including: RevIN (Kim et al., 2021), SAN (Liu et al., 2023), Dish-TS (Fan

Table 2: Forecasting results of backbone models with and without FREMEN. Results are averaged with prediction length $L_y \in \{96, 192, 336, 720\}$. The best results are highlighted in **bold**.

| Methods Metrics | iTransformer | | + FREMEN | | PatchTST | | + FREMEN | | DLinear | | + FREMEN | | RLinear | | + FREMEN | |
|-----------------|-----------------|-----------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|-----------------|-----------------|------------------------|------------------------|-----------------|-----------------|------------------------|------------------------|
| | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE |
| ETTh1 | 0.511 ±0.001 | 0.500 ±0.001 | 0.451 ±0.017 | 0.448 ±0.011 | 0.495 ±0.024 | 0.490 ±0.023 | 0.438 ±0.003 | 0.439 ±0.002 | 0.425 ±0.002 | 0.440 ±0.004 | 0.407 ±0.001 | 0.423 ±0.002 | 0.535 ±0.014 | 0.504 ±0.005 | 0.466 ±0.028 | 0.455 ±0.013 |
| ETTh2 | 0.786 ±0.068 | 0.642 ±0.035 | 0.379 ±0.002 | 0.405 ±0.002 | 0.649 ±0.102 | 0.526 ±0.058 | 0.373 ±0.014 | 0.401 ±0.008 | 0.489 ±0.012 | 0.476 ±0.009 | 0.335 ±0.005 | 0.383 ±0.003 | 0.618 ±0.018 | 0.553 ±0.008 | 0.393 ±0.004 | 0.413 ±0.002 |
| ETTh1 | 0.449 ±0.004 | 0.454 ±0.004 | 0.400 ±0.001 | 0.407 ±0.001 | 0.419 ±0.009 | 0.430 ±0.009 | 0.382 ±0.010 | 0.396 ±0.007 | 0.357 ±0.001 | 0.378 ±0.001 | 0.355 ±0.001 | 0.375 ±0.001 | 0.419 ±0.003 | 0.419 ±0.001 | 0.416 ±0.003 | 0.415 ±0.002 |
| ETTh2 | 0.562 ±0.024 | 0.523 ±0.026 | 0.289 ±0.001 | 0.333 ±0.002 | 0.392 ±0.145 | 0.412 ±0.097 | 0.281 ±0.002 | 0.326 ±0.003 | 0.291 ±0.011 | 0.352 ±0.014 | 0.256 ±0.005 | 0.313 ±0.005 | 0.362 ±0.006 | 0.407 ±0.006 | 0.289 ±0.002 | 0.332 ±0.002 |
| ECL | 0.182 ±0.002 | 0.282 ±0.003 | 0.172 ±0.002 | 0.264 ±0.002 | 0.211 ±0.004 | 0.309 ±0.008 | 0.202 ±0.003 | 0.294 ±0.005 | 0.173 ±0.001 | 0.274 ±0.001 | 0.167 ±0.001 | 0.260 ±0.001 | 0.214 ±0.004 | 0.304 ±0.004 | 0.211 ±0.004 | 0.290 ±0.004 |
| Traffic | 0.571 ±0.007 | 0.314 ±0.007 | 0.429 ±0.013 | 0.285 ±0.009 | 0.594 ±0.006 | 0.315 ±0.006 | 0.512 ±0.009 | 0.329 ±0.009 | 0.453 ±0.000 | 0.318 ±0.001 | 0.436 ±0.001 | 0.298 ±0.002 | 0.629 ±0.008 | 0.390 ±0.003 | 0.622 ±0.010 | 0.376 ±0.007 |
| Weather | 0.252 ±0.003 | 0.300 ±0.007 | 0.251 ±0.003 | 0.276 ±0.002 | 0.248 ±0.004 | 0.301 ±0.007 | 0.256 ±0.005 | 0.279 ±0.003 | 0.245 ±0.001 | 0.298 ±0.002 | 0.226 ±0.001 | 0.265 ±0.001 | 0.269 ±0.001 | 0.319 ±0.001 | 0.258 ±0.003 | 0.283 ±0.004 |

et al., 2023), and FAN (Ye et al., 2024). RevIN, SAN, and Dish-TS conduct time-domain normalization, while FAN focuses on modeling dominant frequencies to overcome non-stationarity.

Backbones. FREMEN is a model-agnostic framework that can be applied to any time series forecasting model. To validate its effectiveness, we select four mainstream backbones, including: MLP-based DLinear (Zeng et al., 2023) and RLinear (Li et al., 2023), and Transformer-based iTransformer (Liu et al., 2024) and PatchTST (Nie et al., 2022). The normalization baselines and our FREMEN method are deployed on these backbones for the following experiments.

Experiments Details. The prediction length is set as $L_y \in \{96, 192, 336, 720\}$ for all backbones. The input-sequence length L_x is set to 336 for DLinear and 96 for the other backbones. We use the Adam optimizer and report the mean absolute error (MAE) and mean squared error (MSE) as the evaluation metrics. All experiments are implemented with PyTorch 2.3.0 and conducted on a single NVIDIA A100 40GB GPU. Details of setup and full experiment results are in the Appendix.

5.2 MAIN RESULTS

Effectiveness on Time Series Forecasting Backbones. To answer **RQ1**, we present the multivariate forecasting results in Table 2. Here, the MSE and MAE are presented in the form of mean \pm std for five runs across four prediction lengths. It is evident that FREMEN consistently enhances the performance of backbone models by a substantial margin under nearly all experimental settings. For example, the average MSE improvements for iTransformer are notable on ETTh2 (51.78%), ETTh1 (48.57%), and Traffic (24.86%), with an average MSE reduction of 28.43% among all datasets. Comparable improvements are observed for PatchTST, DLinear, and RLinear, with average MSE reductions of 18.75%, 10.31%, and 12.83%, respectively, over all benchmark datasets. The superior results can be primarily attributed to the adaptive frequency weights applied, which yields representations with reduced non-stationarity, thereby benefiting learning of forecasting backbones.

Comparison with Baseline Methods. Continuing the investigation of **RQ1**, we present the evaluation result of different normalization methods on iTransformer and DLinear in Table 3. We observe that FREMEN generally outperforms baselines for different forecasting backbones, achieving the best forecasting results in 24 out of 28 experiment settings on average. Specifically, on the ECL dataset, FREMEN achieves MSE values of 0.172 and 0.167 for iTransformer and DLinear, outperforming the best baseline results (i.e., 0.176 by Dish-TS and 0.171 by SAN). Similarly, on the Traffic dataset, the MSE of FREMEN averaged across backbones is 0.432, compared to 0.450, 0.467, 0.491, and 0.456 for RevIN, SAN, FAN, and Dish-TS, respectively. These results may be attributed to the inadequacy of existing methods in modeling non-stationary structures involving spectral shifts.

5.3 DETAILED ANALYSIS

Distribution Shift Analysis. To answer **RQ2**, we compare the frequency-domain distributional distance between the training and test sets for each normalization method. Specifically, the distance is quantified using the Jensen-Shannon Divergence (JSD) between the empirical distributions of the training and test samples. As shown in Figure 3 (a), we begin by analyzing the JSD across all channels. Overall, FREMEN exhibits the best performance in reducing the train-test distributional gaps.

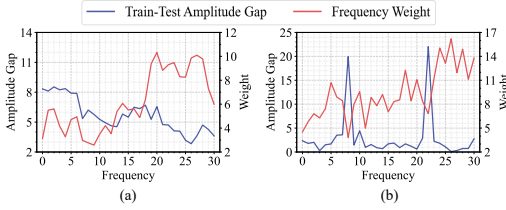


Figure 4: Comparison between the learned frequency weight and the Train-Test amplitude gap.

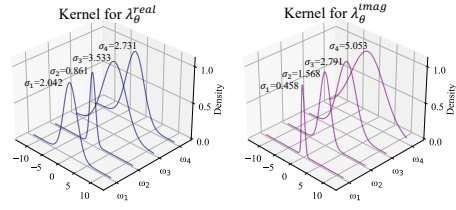


Figure 5: Kernel representations corresponded to the learned frequency weights.

Figure 6: Ablation study with $L_y \in \{96, 720\}$. Best MSE results are highlighted in bold.

| Method | Variant | ETTh1 | | ECL | | Weather | |
|--------------|-----------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | | 96 | 720 | 96 | 720 | 96 | 720 |
| iTransformer | + FREMEN | 0.389 | 0.497 | 0.143 | 0.207 | 0.165 | 0.351 |
| | + random init | 0.406 | 0.532 | 0.164 | 0.235 | 0.174 | 0.357 |
| | + fixed kernel | 0.416 | 0.548 | 0.166 | 0.237 | 0.177 | 0.359 |
| | + single weight | 0.396 | 0.534 | 0.148 | 0.217 | 0.168 | 0.357 |
| DLLinear | + FREMEN | 0.369 | 0.428 | 0.138 | 0.207 | 0.149 | 0.319 |
| | + random init | 0.379 | 0.440 | 0.142 | 0.211 | 0.154 | 0.323 |
| | + fixed kernel | 0.380 | 0.437 | 0.154 | 0.221 | 0.160 | 0.325 |
| | + single weight | 0.376 | 0.433 | 0.141 | 0.209 | 0.152 | 0.321 |

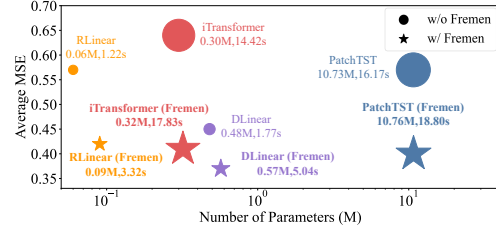


Figure 7: Efficiency analysis of forecasting backbones with (circle) and without FREMEN (star).

kernel representations differ significantly from one another, demonstrating the ability of FREMEN to capture diverse non-stationary patterns in time series data by learning distinct kernels.

Ablation Study. To answer **RQ4**, this section systematically evaluates the key components of FREMEN. We consider three variants to assess the contribution of each part: “*random init*” initializes $\hat{\lambda}$ with a random vector; “*fixed kernel*” fixes λ_θ as the RBF eigenvalue; “*single weight*” uses shared weights for real and imaginary frequency coefficients. As shown in Table 6, FREMEN consistently outperforms all variants. Specifically, the increased forecasting error observed in “*random init*” highlights the critical role of initializing with the empirical RBF kernel eigenvalue $\hat{\lambda}$, which provides a meaningful prior and facilitates more effective frequency weight learning. Similar performance degradation are observed for “*single weight*” due to inability in handling phase shifts via separately modeling non-stationarity for real and imaginary frequency components. The most pronounced decline in performance is observed for “*fixed kernel*”, which is expected since fixing the kernel form severely restricts the model’s expressiveness and adaptability to diverse non-stationarity patterns.

Model Efficiency Analysis. To address **RQ5**, we evaluate the efficiency of forecasting backbones integrated with FREMEN, as illustrated in Figure 7 on the ETTh1 and ETTh2 datasets. In this figure, each pattern represents the outcome of a specific experimental setting, with the size of the pattern reflecting the corresponding running time. The results clearly demonstrate that FREMEN is a lightweight yet highly effective enhancement, introducing only a slight increase in the number of parameters (approximately 0.04M) and computational overhead (averaging 2.85 seconds per epoch), while delivering substantial improvements in forecasting performance.

6 CONCLUSION

This paper studies the problem of learning robust representations for non-stationary time series forecasting. Existing methods mainly focus on measuring non-stationarity in the time domain using low-order statistics. This paper proposed a novel, non-stationarity-aware representation learning method to capture complex temporal structures and variations. We provided theoretical analysis to show that learning a valid non-stationarity measure in frequencies induced a kernel representation, which can be further represented as an orthonormal set of frequency components weights. We introduced FREMEN, that applied frequency weighting on the input time series to learn a more robust representation for forecasting. FREMEN demonstrated effectiveness via extensive experiments. The results confirmed that FREMEN improved mainstream forecasting models by a large margin and outperformed other state-of-the-art normalization methods.

7 REPRODUCIBILITY STATEMENT

We make our code publicly available² which contains detailed implementation of our method. The code and hyperparameters for forecasting backbones adopted in this paper are based on the Time-Series-Library³. For normalization baseline methods, we utilize the code from their official github repository together with the configurations.

REFERENCES

- Yuntao Du, Jindong Wang, Wenjie Feng, Sinno Pan, Tao Qin, Renjun Xu, and Chongjun Wang. Adarnn: Adaptive learning and forecasting of time series. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pp. 402–411, 2021.
- Wei Fan, Pengyang Wang, Dongkun Wang, Dongjie Wang, Yuanchun Zhou, and Yanjie Fu. Dish-ts: A general paradigm for alleviating distribution shift in time series forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 7522–7529, 2023.
- Lu Han, Han-Jia Ye, and De-Chuan Zhan. SIN: Selective and interpretable normalization for long-term time series forecasting. In *Forty-first International Conference on Machine Learning*, 2024.
- Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Reversible instance normalization for accurate time-series forecasting against distribution shift. In *International Conference on Learning Representations*, 2021.
- Zhe Li, Shiyi Qi, Yiduo Li, and Zenglin Xu. Revisiting long-term time series forecasting: An investigation on linear mapping. *arXiv preprint arXiv:2305.10721*, 2023.
- Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long. itransformer: Inverted transformers are effective for time series forecasting. In *The Twelfth International Conference on Learning Representations*, 2024.
- Zhiding Liu, Mingyue Cheng, Zhi Li, Zhenya Huang, Qi Liu, Yanhu Xie, and Enhong Chen. Adaptive normalization for non-stationary time series forecasting: A temporal slice perspective. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Diego Martín Mateos, Leonardo Esteban Riveaud, and Pedro Walter Lamberti. Detecting dynamical changes in time series by using the jensen shannon divergence. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 27(8), 2017.
- Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. *arXiv preprint arXiv:2211.14730*, 2022.
- Yuqi Nie, Nam H. Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In *International Conference on Learning Representations*, 2023.
- Xihao Piao, Zheng Chen, Yushun Dong, Yasuko Matsubara, and Yasushi Sakurai. Frednormer: Frequency domain normalization for non-stationary time series forecasting, 2024a. URL <https://arxiv.org/abs/2410.01860>.
- Xihao Piao, Zheng Chen, Taichi Murayama, Yasuko Matsubara, and Yasushi Sakurai. Fredformer: Frequency debiased transformer for time series forecasting. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD ’24, 2024b.
- John G. Proakis and Dimitris G. Manolakis. Digital signal processing (3rd ed.): Principles, algorithms, and applications. 1996.
- Bernhard Scholkopf and Alexander J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, Cambridge, MA, USA, 2001.

²<https://anonymous.4open.science/r/Fremen-code-82C8>

³<https://github.com/thuml/Time-Series-Library>

- 540 James R. Thompson and James R. Wilson. Multifractal detrended fluctuation analysis: Practical
541 applications to financial time series. In *Mathematics and Computers in Simulation*, pp. 63–88,
542 2016.
- 543 Hao Wang, Licheng Pan, Zhichao Chen, Degui Yang, Sen Zhang, Yifei Yang, Xinggao Liu, Haoxuan
544 Li, and Dacheng Tao. Fredf: Learning to forecast in the frequency domain. In *ICLR*, 2025.
- 545 Zhiyuan Wang, Xovee Xu, Weifeng Zhang, Goce Trajcevski, Ting Zhong, and Fan Zhou. Learning
546 latent seasonal-trend representations for time series forecasting. In *Advances in Neural Informa-*
547 *tion Processing Systems*, 2022.
- 548 Gerald Woo, Chenghao Liu, Doyen Sahoo, Akshat Kumar, and Steven C. H. Hoi. Etsformer: Expo-
549 nential smoothing transformers for time-series forecasting. 2022.
- 550 Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition trans-
551 formers with auto-correlation for long-term series forecasting. In *Advances in Neural Information*
552 *Processing Systems*, 2021.
- 553 Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet:
554 Temporal 2d-variation modeling for general time series analysis. In *International Conference on*
555 *Learning Representations*, 2023.
- 556 Yanfang Xue, Pengfei Fang, Jinyue Tian, Shipeng Zhu, and hui xue. Cosnet: A generalized spectral
557 kernel network. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- 558 A. M. Yaglom. *Correlation Theory of Stationary and Related Random Functions, Volume 1*.
559 Springer, Springer New York, NY, 1987.
- 560 Weiwei Ye, Songgaojun Deng, Qiaosha Zou, and Ning Gui. Frequency adaptive normalization
561 for non-stationary time series forecasting. In *The Thirty-eighth Annual Conference on Neural*
562 *Information Processing Systems*, 2024.
- 563 Kun Yi, Qi Zhang, Wei Fan, Shoujin Wang, Pengyang Wang, Hui He, Ning An, Defu Lian, Longbing
564 Cao, and Zhendong Niu. Frequency-domain MLPs are more effective learners in time series
565 forecasting. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- 566 Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series
567 forecasting? In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp.
568 11121–11128, 2023.
- 569 Dafang Zhao, Xihao Piao, Zheng Chen, Zhengmao Li, and Ittetsu Taniguchi. A unified en-
570 ergy management framework for multi-timescale forecasting in smart grids. *arXiv preprint*
571 *arXiv:2411.15254*, 2024.
- 572 Tian Zhou, Ziqing Ma, Qingsong Wen, Liang Sun, Tao Yao, Wotao Yin, Rong Jin, et al. Film:
573 Frequency improved legendre memory model for long-term time series forecasting. In *Advances*
574 *in neural information processing systems*, pp. 12677–12690, 2022a.
- 575 Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer: Frequency
576 enhanced decomposed transformer for long-term series forecasting. In *Proc. 39th International*
577 *Conference on Machine Learning*, pp. 1–12, 2022b.
- 578 Lingwei Zhu, Haseeb Shah, Han Wang, Yukie Nagai, and Martha White. $\$q$ -exponential family
579 for policy optimization. In *The Thirteenth International Conference on Learning Representations*,
580 2025.

581 APPENDIX/SUPPLEMENTARY MATERIALS

582 A LARGE LANGUAGE MODELS (LLMs) USAGE STATEMENT

583 Large Language Models (LLMs) were used solely to assist with the English writing and language
584 polishing of this manuscript. All research ideas, experimental design, data analysis, and scientific
585 content were conceived and executed by the authors without the involvement of LLMs.
586
587
588
589

B ADDITIONAL EXPERIMENTAL DETAILS

B.1 DATASET DETAILS

The statistical properties of the seven benchmark datasets are summarized in Table 4. To quantify the distribution shift between training and test samples, we employ the Jensen-Shannon Divergence (JSD), which serves as a measure of dataset non-stationarity (Mateos et al., 2017). The JSD is computed between the empirical distributions derived from the training and test sets. The computational procedure for JSD between two sample sets is formally described in Algorithm 1. Our implementation utilizes the “jensenshannon” function from the SciPy library, which computes the square root of the JSD. Then, the squared output is considered as the distributional divergence measure.

Table 4: Statistics of benchmark datasets.

| Datasets | ETTh1 | ETTh2 | ETTm1 | ETTm2 | ECL | Traffic | Weather |
|------------------|-----------|-----------|-----------|-----------|-----------|-----------|---------|
| # Channels | 7 | 7 | 7 | 7 | 321 | 862 | 21 |
| # Timesteps | 17,420 | 17,420 | 69,680 | 69,680 | 26,304 | 17,544 | 52,969 |
| Sample Frequency | 1h | 1h | 15mins | 15mins | 1h | 1h | 10mins |
| Time Range | 2016-2017 | 2017-2018 | 2016-2017 | 2017-2018 | 2012-2014 | 2015-2016 | 2020 |
| JSD* | 0.2091 | 0.2839 | 0.1225 | 0.3138 | 0.0716 | 0.0627 | 0.1524 |

* A smaller JSD indicates a more stationary time series dataset

Algorithm 1: Computation of Jensen-Shannon Divergence between two sets of samples

Input: arrays a, b ; number of bins B

Output: Jensen-Shannon Divergence D_{JSD}^2

- 1: $v_{\min} \leftarrow \min(\min(a), \min(b))$
 - 2: $v_{\max} \leftarrow \max(\max(a), \max(b))$
 - 3: $h_a \leftarrow \text{histogram}(a; B, [v_{\min}, v_{\max}])$
 - 4: $h_b \leftarrow \text{histogram}(b; B, [v_{\min}, v_{\max}])$
 - 5: $p \leftarrow h_a / \sum h_a$
 - 6: $q \leftarrow h_b / \sum h_b$
 - 7: $D_{\text{JSD}} \leftarrow \text{jensenshannon}(p, q, \text{base} = 2)$
 - 8: **return** D_{JSD}^2
-

B.2 BASELINE METHOD DETAILS

In this study, we study four state-of-the-art normalization methods as baselines: RevIN, SAN, Dish-TS, and FAN. The technical details of each approach are presented below:

Reversible Instance Normalization (RevIN) (Kim et al., 2021). RevIN proposes a symmetric normalization-denormalization framework to address distribution shifts in time series data. The method first applies instance-wise z-score normalization to input samples, effectively eliminating non-stationary components. The normalized data is then fed into the forecasting model for forecasting. After which the original statistical properties (i.e., mean and variance) are restored to the forecasting result through a denormalization process. This reversible transformation maintains crucial distributional characteristics while enabling models to operate on stationary representations.

Slice-level Adaptive Normalization (SAN) (Liu et al., 2023). SAN introduces a fine-grained normalization approach that operates at the sub-series level instead of the whole input and output series. Unlike direct statistical transfer, SAN employs a dedicated statistics prediction module to explicitly model the evolution of mean and variance of the data distribution. During training, this module is first pre-trained to predict future statistics. Then, it is frozen and used to produce normalized input data for the downstream forecasting models training.

Dish-TS (Fan et al., 2023). Dish-TS provides a systematic framework that classifies distribution shifts into intra-space (within input/output spaces) and inter-space (between input/output spaces) variations. The method introduces a specialized network architecture for input and output distribution estimation, augmented with learnable adaptive distribution statistics. Notably, Dish-TS incorporates empirical mean values as prior knowledge to enhance statistical learning.

Frequency Adaptive Normalization (FAN) (Ye et al., 2024). FAN represents the first normalization approach that addresses non-stationarity through frequency-domain analysis. The method classifies the whole frequency spectrum into two sets where the top- k dominant frequency components are considered as non-stationary components which are fed into an MLP network to model the future statistical variations. The remaining stationary frequency components are directly fed into the forecasting models for prediction. Combining the two parts of outputs, FAN effectively captures potential variants in frequencies to mitigate non-stationarity.

B.3 FORECASTING BACKBONE DETAILS

In this work, we evaluate our method on four prominent time series forecasting architectures: the MLP-based DLinear and RLinear, and the Transformer-based iTransformer and PatchTST. We provide an overview of their key design principles:

DLinear (Zeng et al., 2023). DLinear establishes a decomposition-based linear architecture. The method first decomposes the input series into trend and seasonal components using moving average smoothing. These components are then processed independently through dedicated linear layers, with their outputs aggregated to produce the final prediction result.

RLinear (Li et al., 2023). RLinear adopts an extremely light-weight architecture comprising a single linear layer enhanced with reversible normalization. The approach capitalizes on the inherent capability of linear mappings to capture periodic patterns, while the normalization scheme transforms trend components into seasonality-like representations.

iTransformer (Liu et al., 2024). iTransformer reconfigures the standard Transformer architecture for time series analysis. Rather than tokenizing multivariate points at each timestep, it represents entire univariate series as individual tokens. This inverted paradigm enables self-attention mechanisms to focus on cross-variate dependencies while feed-forward networks handle temporal patterns, better accommodating the unique characteristics of time series data.

PatchTST (Nie et al., 2022). PatchTST employs a Transformer encoder architecture with two core modifications: First, time series are divided into overlapping or non-overlapping patches that serve as input tokens, reducing sequence length while preserving local patterns. Second, it processes each channel independently with shared weights, enabling efficient multivariate forecasting.

B.4 OTHER EXPERIMENTAL DETAILS

Implementation details. We mainly tune the value of the hidden dimension H within the range of $\{16, 64, 128, 256, 512\}$ and select the one with the best forecasting accuracy on the validation set as the final hyperparameter for each experimental setting. We repeat each experiment for five times with fixed seed and report the average evaluation results.

Loss Functions. We employ mean squared error (MSE) as the loss function for all forecasting backbones, which quantifies the averaged squared difference between the predicted and actual target time series. Mathematically, the MSE loss is expressive as: $\mathcal{L}_{\text{MSE}} = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2$, where N is the number of samples, \hat{y}_i represents the predicted value, and y_i represents the actual value.

C ADDITIONAL EVALUATION RESULTS

C.1 EFFECTIVENESS ON TIME SERIES FORECASTING BACKBONES

Table 5 presents comprehensive evaluation results across all prediction lengths and benchmark datasets for the four forecasting backbones. The results demonstrate that FREMEN consistently enhances the performance of all four baseline architectures. Quantitative analysis reveals that FREMEN outperforms the original models in 210 out of 224 evaluation scenarios (4 backbones \times 7 datasets \times 4 prediction lengths \times 2 metrics). More specifically, FREMEN achieves average MSE reductions of 12.62%, 12.47%, 17.41%, and 25.97% for prediction lengths of 96, 192, 336, and 720 steps, respectively. Notably, the performance gains become increasingly pronounced as the prediction horizon lengthens. These findings provide strong evidence for the effectiveness of the frequency weighting mechanism in FREMEN, particularly in the context of long-term time series forecasting tasks.

C.3 FREQUENCY-WISE DISTRIBUTION SHIFT ANALYSIS

Figure 8 provides additional examples illustrating the frequency-wise train-test distribution discrepancies. Although variables across different datasets exhibit diverse spectral distributions, FREMEN demonstrates effectiveness in addressing distributional differences between the training and test set.

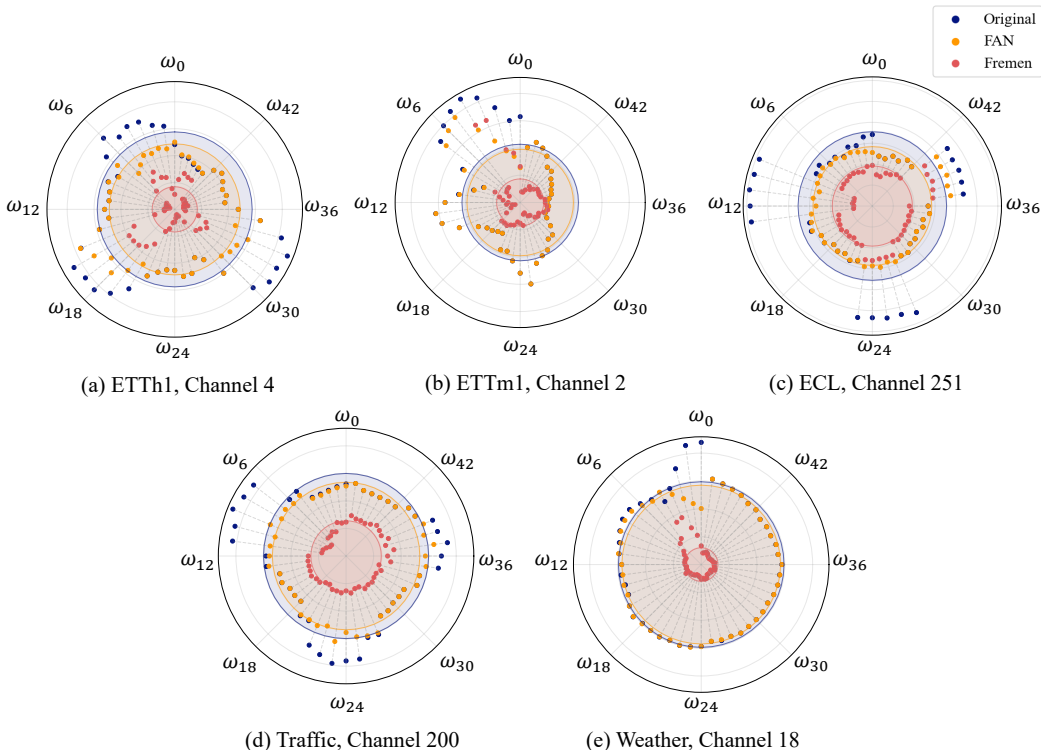


Figure 8: Additional showcases on the Train-Test distribution distance over all frequencies. Each scatter corresponds to the JSD value of a frequency component. A smaller distance to the center indicates a smaller JSD value.

C.4 FREQUENCY WEIGHT ANALYSIS

To investigate the dynamics of the learned frequency weights, we visualize both the training loss and the evolution of $\lambda_{\theta}^{\text{real}}$ in Figure 9, using DLinear on the Traffic dataset. The upper subfigure demonstrates a clear positive correlation between increasing $\lambda_{\theta}^{\text{real}}$ and improved model accuracy. As discussed in our preliminary analysis from the main paper, $\lambda_{\theta}^{\text{real}}$ governs kernel representation in the frequency domain. While its exact closed-form kernel mapping remains analytically intractable, we hypothesize its membership within common kernel families and empirically analyze the induced representational transformations. The lower subfigure illustrates the kernel evolution under two constrained settings: RBF and Cauchy kernels, sampled across five training epochs. Notably, we observe divergent trends in the scale parameters σ : the RBF kernel’s σ exhibits monotonic growth, while the Cauchy kernel’s σ demonstrates consistent decay. This antithetical behavior underscores the expressive flexibility of $\lambda_{\theta}^{\text{real}}$, validating its capacity to adaptively model heterogeneous non-stationary patterns through implicit kernel learning.

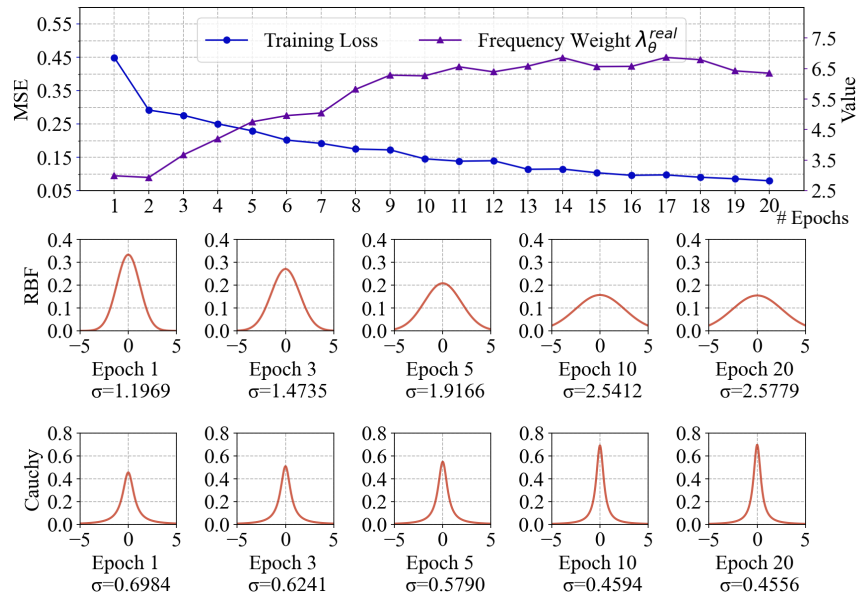


Figure 9: Training process of DLinear with FREMEN on the Traffic dataset. Upper subgraph: the change of the training loss and weights for frequency components. Lower subgraph: the evolution of kernels corresponding to the frequency weights.