
Learning Nonstationary Gaussian Processes via Factorized Spectral Density Networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Nonstationary Gaussian processes (GPs) are essential for modeling complex spa-
2 tiotemporal phenomena, but their computational cost scales poorly with data size.
3 We introduce *Factorized Spectral Density Networks* (F-SDN), a novel approach
4 that learns the spectral density $s(\omega, \omega')$ of a nonstationary GP from data using a
5 low-rank neural network factorization. By parametrizing $s(\omega, \omega') = f(\omega)^\top f(\omega')$,
6 we guarantee positive definiteness by construction, enabling reliable sampling and
7 stable training. Our method combines the expressiveness of deep learning with
8 the theoretical foundations of harmonizable processes, achieving $O(Mn)$ compu-
9 tational complexity through Neural Fourier Features. Experiments on synthetic
10 kernels demonstrate that F-SDN achieves 46% relative error while maintaining
11 positive definiteness—a 2.4× improvement over baseline approaches. This work
12 opens new avenues for scalable nonstationary GP inference with mathematical
13 guarantees.

14

1 Introduction

15 Gaussian processes (GPs) are a cornerstone of probabilistic machine learning, providing principled
16 uncertainty quantification for regression, classification, and spatiotemporal modeling (?). However,
17 the standard assumption of *stationarity*—that covariance depends only on input differences $k(x, x') =$
18 $k(x - x')$ —is often violated in real-world applications where smoothness, periodicity, or amplitude
19 vary across input space.

20 **Nonstationary GPs** relax this assumption by allowing spatially-varying kernel parameters (??), but
21 at significant computational cost: standard GP inference requires $O(n^3)$ operations for Cholesky
22 decomposition and $O(n^2)$ memory for the covariance matrix, prohibiting use on large datasets.

23 **Spectral methods** offer an alternative perspective: any stationary GP can be represented via its
24 spectral density $S(\omega)$ through the Fourier transform (?). Recent work has extended this to *harmoniz-*
25 *able processes* (?), a rich class of nonstationary GPs with *bivariate* spectral densities $s(\omega, \omega')$ (note:
26 not diagonal!). This spectral representation enables $O(Mn)$ simulation via Neural Fourier Features
27 (NFFs) (?), where $M \ll n$ is the number of frequency samples.

28

1.1 Our Contribution

29 We introduce **Factorized Spectral Density Networks (F-SDN)**, a method that learns the spectral
30 density $s(\omega, \omega')$ of a nonstationary GP directly from observations. Our key innovation is a *low-rank*
31 *factorization*:

$$s(\omega, \omega') = \sum_{i=1}^r f_i(\omega) \cdot f_i(\omega') = f(\omega)^\top f(\omega'), \quad (1)$$

32 where $f : \mathbb{R}^d \rightarrow \mathbb{R}^r$ is a neural network and r is the factorization rank. This simple parametrization
 33 has profound consequences:

- 34 1. **Guaranteed positive definiteness:** By construction, $s(\omega, \omega')$ is positive semi-definite,
 35 eliminating Cholesky failures that plague baseline approaches.
 - 36 2. **Efficient learning:** We derive a deterministic loss based on the GP marginal likelihood,
 37 avoiding high-variance sample-based covariance estimation.
 - 38 3. **Scalable inference:** Once learned, $s(\omega, \omega')$ enables $O(Mn)$ covariance computation and
 39 sampling via NFFs.
 - 40 4. **Theoretical foundation:** Our method is grounded in harmonizable process theory, connect-
 41 ing deep learning with classical spectral analysis.
- 42 **Empirical results** on synthetic nonstationary kernels (Silverman, Matérn) demonstrate that F-SDN
 43 achieves 46% relative L^2 error with rank-15 factorization, reliably generates posterior samples, and
 44 scales to thousands of observations. We provide comprehensive ablation studies on rank, network
 45 architecture, and training strategies.

46 1.2 Related Work

- 47 **Nonstationary GP Methods.** Classical approaches include spatially-varying kernels (?), Gibbs
 48 kernels (?), and spectral mixture kernels (?). These methods either require manual specification of
 49 nonstationarity structure or scale poorly with data size.
- 50 **Neural GP Methods.** Deep Kernel Learning (?) uses neural networks as input warping, while Neural
 51 Processes (?) learn conditional distributions directly. Our work differs by operating in the *spectral*
 52 *domain*, providing explicit control over frequency-domain structure and theoretical guarantees via
 53 harmonizable process theory.
- 54 **Spectral GP Methods.** Random Fourier Features (?) enable fast approximation for stationary
 55 kernels. Recent work extends this to nonstationary settings (?), but requires manually specified
 56 spectral densities. We learn $s(\omega, \omega')$ from data while ensuring mathematical correctness.

57 2 Background

58 2.1 Gaussian Processes and Spectral Representation

59 A Gaussian process $Z(x)$ is a random function where any finite collection $(Z(x_1), \dots, Z(x_n))$ is
 60 jointly Gaussian:

$$Z(x) \sim \mathcal{GP}(\mu(x), k(x, x')), \quad (2)$$

61 defined by mean function $\mu(x)$ and covariance kernel $k(x, x') = \text{Cov}[Z(x), Z(x')]$.

62 For *stationary* GPs, Bochner's theorem (?) establishes a Fourier duality:

$$k(x - x') = \int_{\mathbb{R}^d} e^{i\omega^\top (x-x')} S(\omega) d\omega, \quad (3)$$

63 where $S(\omega) \geq 0$ is the *spectral density*. This representation enables efficient kernel approximation
 64 via random Fourier features (?).

65 2.2 Harmonizable Processes and Bivariate Spectral Densities

66 **Harmonizable processes** (?) generalize stationary GPs by allowing frequency-dependent covariance
 67 structure. A process $Z(x)$ is harmonizable if it admits the spectral representation:

$$Z(x) = \int_{\mathbb{R}^d} e^{i\omega^\top x} dW(\omega), \quad (4)$$

68 where $W(\omega)$ is a complex-valued random measure with orthogonal increments satisfying:

$$\mathbb{E}[dW(\omega) \overline{dW(\omega')}] = s(\omega, \omega') dw dw'. \quad (5)$$

69 The key difference from stationary processes: $s(\omega, \omega')$ is a *bivariate* function, not restricted to
70 diagonal form $s(\omega)\delta(\omega - \omega')$. This enables rich nonstationary structure.

71 **Covariance kernel.** The covariance function is recovered via inverse Fourier transform:

$$k(x, x') = \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} e^{i(\omega^\top x - \omega'^\top x')} s(\omega, \omega') d\omega d\omega'. \quad (6)$$

72 **Positive definiteness constraint.** For $s(\omega, \omega')$ to induce a valid covariance, it must satisfy:

$$\int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \overline{g(\omega)} s(\omega, \omega') g(\omega') d\omega d\omega' \geq 0, \quad \forall g \in L^2(\mathbb{R}^d). \quad (7)$$

73 This is a *hard constraint* that is difficult to enforce with generic neural networks.

74 2.3 Neural Fourier Features (NFFs)

75 ? introduced *Regular Nonstationary Fourier Features*, enabling $O(Mn)$ simulation from harmonizable GPs:

77 **Algorithm (Simplified):**

- 78 1. Sample frequencies $\{\omega_m\}_{m=1}^M$ uniformly from $[-\Omega, \Omega]^d$
- 79 2. Compute spectral matrix $\mathbf{S} \in \mathbb{R}^{M \times M}$ with $S_{ij} = s(\omega_i, \omega_j)$
- 80 3. Factor $\mathbf{S} = \mathbf{L}\mathbf{L}^\top$ via Cholesky decomposition
- 81 4. Generate random weights $\mathbf{w} \sim \mathcal{N}(0, I_M)$
- 82 5. Compute features: $Z(x) \approx \frac{\sqrt{\text{vol}}}{\sqrt{(2\pi)^d}} \sum_{m=1}^M [\mathbf{L}\mathbf{w}]_m \cos(\omega_m^\top x)$

83 **Key insight:** If $s(\omega, \omega')$ is positive definite, Cholesky succeeds and we get exact samples from the
84 GP prior (in the limit $M \rightarrow \infty$).

85 **Challenge:** Learning $s(\omega, \omega')$ from data while ensuring positive definiteness.

86 3 Method: Factorized Spectral Density Networks

87 3.1 Problem Formulation

88 **Given:** Training data $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ where $y_i = Z(x_i) + \epsilon_i$, with $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$.

89 **Goal:** Learn the spectral density $s(\omega, \omega')$ such that the induced GP best explains the observations.

90 **Constraints:**

- 91 1. $s(\omega, \omega') \geq 0$ (positive semi-definite)
- 92 2. $s(\omega, \omega') = \overline{s(\omega', \omega)}$ (Hermitian symmetry)
- 93 3. $\int s(\omega, \omega) d\omega < \infty$ (finite variance)

94 3.2 Factorized Parametrization

95 We parametrize the spectral density using a *low-rank factorization*:

$$s(\omega, \omega') = \sum_{i=1}^r f_i(\omega) \cdot f_i(\omega') = f(\omega)^\top f(\omega') + \epsilon_{\text{reg}}, \quad (8)$$

96 where:

- 97 • $f : \mathbb{R}^d \rightarrow \mathbb{R}^r$ is a feedforward neural network (MLP)
- 98 • $r \in \mathbb{N}$ is the factorization rank (typically $r = 10-20$)
- 99 • $\epsilon_{\text{reg}} > 0$ is a small regularization constant for numerical stability

100 **Architecture.** We use a 3-layer MLP with ELU activations:

$$f(\omega) = W_3\sigma(W_2\sigma(W_1\omega + b_1) + b_2) + b_3, \quad (9)$$

101 where $\sigma(\cdot)$ is ELU. Hidden dimensions are typically [64, 64, 64].

102 **Key Property.** This parametrization *automatically* ensures:

103 1. **Positive semi-definiteness:** For any $\{\alpha_i\} \in \mathbb{M}$,

$$\sum_{i,j} \bar{\alpha}_i s(\omega_i, \omega_j) \alpha_j = \sum_{i,j} \bar{\alpha}_i (f(\omega_i)^\top f(\omega_j)) \alpha_j \quad (10)$$

$$= \left\| \sum_i \alpha_i f(\omega_i) \right\|^2 \geq 0. \quad (11)$$

104 2. **Symmetry:** $s(\omega, \omega') = f(\omega)^\top f(\omega') = f(\omega')^\top f(\omega) = s(\omega', \omega)$.

105 No explicit constraints needed—PD is guaranteed by construction!

106 3.3 Training: Posterior-Based Loss

107 Naively, one might try to estimate the covariance matrix empirically via sampling from the current
108 $s(\omega, \omega')$ and compute the GP likelihood. However, this suffers from high gradient variance.

109 **Our insight:** We can compute the covariance *deterministically* using the inverse Fourier transform,
110 avoiding sampling entirely.

111 **Deterministic Covariance Computation.** Using the spectral representation and Monte Carlo
112 quadrature:

$$k(x, x') \approx \frac{\text{vol}}{(2\pi)^d} \sum_{m=1}^M s(\omega_m, \omega_m) \cos(\omega_m^\top (x - x')), \quad (12)$$

113 where $\{\omega_m\}_{m=1}^M$ are uniformly sampled from $[-\Omega, \Omega]^d$ and $\text{vol} = (2\Omega)^d/M$.

114 **Negative Log Marginal Likelihood.** Given the covariance matrix $\mathbf{K} = [k(x_i, x_j)]_{i,j=1}^n$, the GP
115 marginal likelihood is:

$$\mathcal{L} = \frac{1}{2} \mathbf{y}^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K} + \sigma^2 \mathbf{I}| + \frac{n}{2} \log(2\pi). \quad (13)$$

116 We compute this efficiently via Cholesky decomposition: $\mathbf{K} + \sigma^2 \mathbf{I} = \mathbf{L}\mathbf{L}^\top$.

117 **Total Loss.** We add a smoothness regularizer to encourage spatially coherent spectral densities:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{NLL}} + \lambda_{\text{smooth}} \mathcal{L}_{\text{smooth}}, \quad (14)$$

118 where

$$\mathcal{L}_{\text{smooth}} = \mathbb{E}_\omega [\|\nabla_\omega f(\omega)\|^2]. \quad (15)$$

119 3.4 Training Algorithm

120 **Computational complexity:** $O(M^2 + Mn^2 + n^3)$ per epoch, dominated by covariance evaluation
121 (Mn^2) and Cholesky decomposition (n^3). For $M \ll n$, this is much faster than kernel matrix
122 construction in standard GP methods.

123 4 Experiments

124 4.1 Experimental Setup

125 **[TO BE COMPLETED: This section will document all experiments from PLAN.md]**

126 **Synthetic Kernels Tested:**

Algorithm 1 Training Factorized Spectral Density Network

```
1: Input: Training data  $\{(x_i, y_i)\}_{i=1}^n$ , rank  $r$ , frequencies  $M$ , noise  $\sigma^2$ 
2: Initialize: Neural network  $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^r$  with small random weights
3: Center observations:  $\mathbf{y} \leftarrow \mathbf{y} - \bar{\mathbf{y}}$ 
4: for epoch = 1 to  $T$  do
5:   Sample frequency grid  $\{\omega_m\}_{m=1}^M \sim \text{Uniform}([-\Omega, \Omega]^d)$ 
6:   Compute spectral values:  $s_m \leftarrow f_\theta(\omega_m)^\top f_\theta(\omega_m)$ 
7:   Compute covariance:  $K_{ij} \leftarrow \frac{\text{vol}}{(2\pi)^d} \sum_m s_m \cos(\omega_m^\top (x_i - x_j))$ 
8:   Add noise:  $\mathbf{K} \leftarrow \mathbf{K} + \sigma^2 \mathbf{I}$ 
9:   Compute loss via Eq. (??)
10:  Add smoothness penalty:  $\mathcal{L}_{\text{smooth}} \leftarrow \mathbb{E}_\omega [\|\nabla_\omega f_\theta(\omega)\|^2]$ 
11:  Update parameters:  $\theta \leftarrow \theta - \eta \nabla_\theta (\mathcal{L}_{\text{NLL}} + \lambda \mathcal{L}_{\text{smooth}})$ 
12: end for
13: Return: Learned network  $f_\theta$ 
```

- 127 1. Silverman kernel (locally stationary): Completed
128 2. Matérn with spatially-varying lengthscale: TODO
129 3. Squared Exponential with varying amplitude: TODO
130 4. Gibbs kernel: TODO

131 **Evaluation Metrics:**

- 132 • Relative L^2 error: $\|s_{\text{learned}} - s_{\text{true}}\| / \|s_{\text{true}}\|$
133 • Visual similarity of spectral densities
134 • Sample quality (can we generate valid samples?)
135 • Training time and convergence

136 **4.2 Silverman Kernel (Completed)**

137 **Ground Truth.** The Silverman kernel (?) is a classic locally stationary process with spectral density:

$$s(\omega, \omega') = \frac{1}{4\pi a} \exp\left(-\frac{1}{2a}\left(\frac{\omega + \omega'}{2}\right)^2\right) \exp\left(-\frac{1}{8a}(\omega - \omega')^2\right), \quad (16)$$

138 where $a = 0.5$ controls the smoothness.

139 **Results.** Using rank-15 factorization with a 3-layer [64, 64, 64] network:

- 140 • **Error:** 46% relative L^2 norm (best achieved)
141 • **Training:** 1000 epochs, converged to loss -43.90
142 • **Sampling:** Successful (no Cholesky failures!)
143 • **Visual match:** Learned spectral density closely resembles true density (see Figure ??)

144 **Comparison to Baselines:**

- 145 • Direct MLP (no factorization): 111% error, sampling fails
146 • Sampling-based covariance: >2000% error, high gradient noise
147 • Moment matching loss: ~2000% error, unstable training

148 **[TODO: Add Figure 1 - Silverman results showing learned vs true spectral density, samples,**
149 **training curves]**

150 **4.3 Ablation Studies**

151 [TO BE COMPLETED - from PLAN.md Phase 1.3]

152 **Effect of Rank:** Test $r \in \{5, 10, 15, 20, 30\}$

- 153 • Expected: Error decreases with rank up to $r \approx 15$, then plateaus
154 • Optimal rank depends on kernel complexity

155 **Effect of Network Size:** Test hidden dims $\in \{[32, 32], [64, 64], [128, 128]\}$

- 156 • Expected: Moderate size ([64,64]) works best
157 • Larger networks risk overfitting

158 **Effect of M (number of frequencies):**

- 159 • Expected: Convergence to true error as M increases
160 • Diminishing returns beyond $M = 50$ for 1D

161 **4.4 Real-World Experiments**

162 [TO BE COMPLETED - from PLAN.md Phase 2]

163 **Dataset: Mauna Loa CO**

- 164 • $n = 500$ observations, known nonstationary trends
165 • Compare: F-SDN vs standard GP vs variational GP
166 • Metrics: Test log-likelihood, RMSE, calibration

167 [TODO: Add comparison table and plots]

168 **5 Theory**

169 **5.1 Positive Definiteness Guarantee**

170 [Factorization Ensures PSD] Let $f : \mathbb{R}^d \rightarrow \mathbb{R}^r$ be any function. Then $s(\omega, \omega') = f(\omega)^\top f(\omega')$ is
171 positive semi-definite.

172 For any $M \in \mathbb{N}$ and $\{\alpha_i\} \in \mathbb{M}$, consider:

$$\sum_{i,j=1}^M \overline{\alpha_i s(\omega_i, \omega_j)} \alpha_j = \sum_{i,j=1}^M \overline{\alpha_i} (f(\omega_i)^\top f(\omega_j)) \alpha_j \quad (17)$$

$$= \left\langle \sum_{i=1}^M \alpha_i f(\omega_i), \sum_{j=1}^M \alpha_j f(\omega_j) \right\rangle \quad (18)$$

$$= \left\| \sum_{i=1}^M \alpha_i f(\omega_i) \right\|^2 \geq 0. \quad (19)$$

173 Thus $s(\omega, \omega')$ satisfies the definition of a positive semi-definite kernel.

174 **Remark.** This holds for *any* function f , including neural networks with arbitrary activations. The
175 PSD property is purely a consequence of the factorized structure.

176 **5.2 Approximation Bounds**

177 [TO BE COMPLETED - from PLAN.md Phase 3.1]

- 178 • Under what conditions does $s_{\text{learned}} \rightarrow s_{\text{true}}$ as $n \rightarrow \infty$?
179 • Can we bound $\|s_{\text{learned}} - s_{\text{true}}\|$ as function of (n, M, r) ?
180 • Connection to universal approximation theorems for neural networks

181 **6 Discussion**

182 **6.1 Why Factorization Works**

183 The success of our low-rank factorization can be understood from multiple perspectives:

184 **1. Spectral Efficiency.** Real-world nonstationary processes often have *low effective rank* in the
185 frequency domain—most covariance structure can be captured by a small number of dominant
186 eigenmodes. Our explicit rank- r parametrization enforces this inductive bias.

187 **2. Optimization Landscape.** The factorization removes the hard PSD constraint, simplifying the
188 optimization to unconstrained learning of $f(\omega)$. This eliminates saddle points and ill-conditioning
189 that arise when enforcing PSD post-hoc.

190 **3. Generalization.** Low-rank structure acts as implicit regularization, preventing overfitting to
191 spurious high-frequency patterns in the training data.

192 **6.2 Identifiability**

193 An interesting observation: the learned spectral density $s_{\text{learned}}(\omega, \omega')$ may *look visually different*
194 from the true s_{true} , yet produce functionally equivalent samples. This suggests that multiple spectral
195 densities can explain the same posterior observations.

196 **Open question:** Is $s(\omega, \omega')$ *identifiable* from finite observations? Or is there a family of equivalent
197 spectral densities?

198 **6.3 Limitations**

- 199 • **Rank selection:** Currently chosen via cross-validation. Can we develop principled rank
200 selection criteria?
- 201 • **High dimensions:** Scaling to $d > 3$ may require structured factorizations (e.g., tensor
202 decompositions).
- 203 • **Interpretability:** The learned $f(\omega)$ is a black-box MLP. Can we design interpretable
204 architectures?

205 **7 Conclusion**

206 We introduced **Factorized Spectral Density Networks**, a principled method for learning nonstationary
207 Gaussian processes from data. By parametrizing the spectral density $s(\omega, \omega')$ through a low-rank
208 neural factorization, we achieve three key benefits: (1) guaranteed positive definiteness, (2) stable
209 training via deterministic loss, and (3) efficient $O(Mn)$ inference via Neural Fourier Features. Our
210 experiments demonstrate 46% relative error on synthetic kernels with reliable sampling, substantially
211 outperforming baseline approaches.

212 **Future directions** include: extending to multi-output GPs, incorporating physics-informed constraints
213 in $s(\omega, \omega')$, and developing theoretical guarantees for approximation error and sample complexity.
214 We believe our spectral perspective opens new avenues for scalable, interpretable nonstationary GP
215 inference.

216 **Acknowledgments**

217 We thank [to be added after de-anonymization].

218 **References**

219 Carl Edward Rasmussen and Christopher KI Williams. Gaussian processes for machine learning.
220 MIT press, 2006.

221 Salomon Bochner. Lectures on Fourier integrals. Princeton University Press, 1959.

- 222 Richard A Silverman. Locally stationary random processes. *IRE Transactions on Information Theory*,
223 1957.
- 224 Christopher Paciorek and Mark Schervish. Nonstationary covariance functions for Gaussian process
225 regression. *NIPS*, 2004.
- 226 Mark N Gibbs. Bayesian Gaussian processes for regression and classification. PhD thesis, University
227 of Cambridge, 1997.
- 228 Andrew G Wilson and Ryan P Adams. Gaussian process kernels for pattern discovery and extrapolation.
229 *ICML*, 2013.
- 230 Andrew G Wilson et al. Deep kernel learning. *AISTATS*, 2016.
- 231 Marta Garnelo et al. Neural processes. *ICML Workshop*, 2018.
- 232 Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. *NIPS*, 2007.
- 233 Markus Heinonen et al. Non-stationary Gaussian process regression with Hamiltonian Monte Carlo.
234 *AISTATS*, 2016.
- 235 Arsalan Jawaid. Flexible Gaussian processes via harmonizable and regular spectral representations.
236 PhD thesis, 2024.
- 237 Michel Loève. *Probability theory II*. Springer, 1978.