

SPARSE-SMOOTH DECOMPOSITION FOR NONLINEAR INDUSTRIAL TIME SERIES FORECASTING

Anonymous authors

Paper under double-blind review

ABSTRACT

Industrial time series forecasting faces unique challenges: hundreds of correlated sensors, complex nonlinear dynamics, and the critical need for interpretable models that engineers can trust. We introduce nonlinear causal sparse-smooth network, a framework that decomposes high-dimensional industrial forecasting into interpretable sparse-smooth feature extraction followed by nonlinear prediction. Unlike black-box deep learning approaches that use all sensors indiscriminately, our method automatically identifies critical sensor subsets while learning smooth temporal filters that reflect physical process dynamics. We cast this as a structured optimization problem with sparsity penalties for sensor selection and smoothness regularization for temporal patterns, unified within an identifiable Wiener model architecture. Theoretically, we prove convergence guarantees, establish sensor selection consistency, and derive generalization bounds that explicitly account for the interplay between sparsity, smoothness, and nonlinearity. On an industrial refinery benchmark, our structured approach achieves a 25.2% lower error rate than state-of-the-art Transformer models, while simultaneously identifying a sparse subset of critical sensors and their interpretable dynamic modes. Our work demonstrates that incorporating strong, domain-aware inductive biases into a structured architecture offers a powerful alternative to monolithic black-box models for real-world industrial forecasting.

1 INTRODUCTION

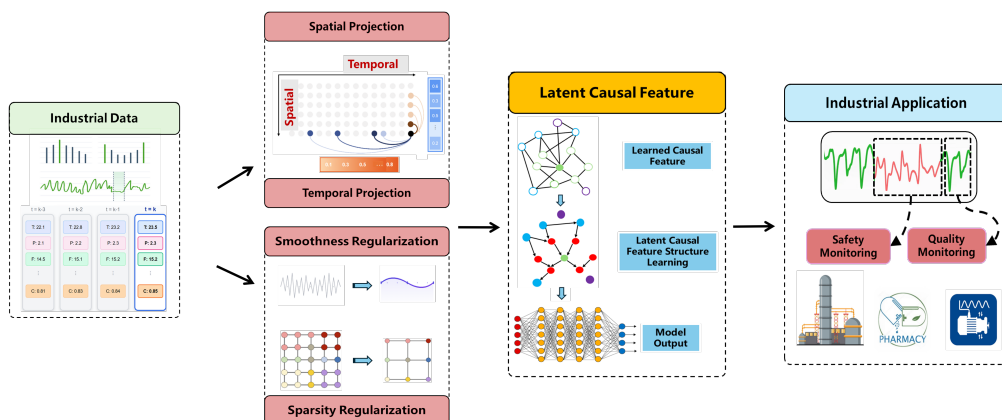
Industrial processes generate vast amounts of sensor data, yet paradoxically, the most economically important variables—product quality indicators—often remain unmeasured in real-time (Qin, 2012; Ge, 2017). Hardware analyzers for variables such as distillation column compositions, polymer melt indices, or catalyst activity levels typically require laboratory analysis with delays ranging from hours to days, creating a fundamental control challenge (Kadlec et al., 2009; Souza et al., 2016). Soft sensors address this gap by constructing mathematical models that estimate these hard-to-measure variables from readily available process measurements such as temperatures, pressures, and flow rates (Fortuna et al., 2007; Kano & Ogawa, 2008). While conceptually straightforward, developing effective soft sensors faces multiple challenges: the high dimensionality of modern sensor arrays, complex nonlinear process dynamics, time-varying operating conditions, and the industrial requirement for interpretable models that operators can trust and maintain (Jiang et al., 2021; Shang et al., 2014).

A critical yet underexplored aspect of industrial soft sensing is the inherent redundancy in sensor networks and the smooth nature of process dynamics (Sun & Ge, 2021; Yuan et al., 2020). Manufacturing facilities often install redundant sensors for safety and reliability, leading to highly correlated measurements that complicate model identification (Rasheed et al., 2020). Simultaneously, physical processes governed by conservation laws, reaction kinetics, and transport phenomena naturally exhibit smooth temporal behavior rather than abrupt changes (Seborg et al., 2016). Traditional soft sensing approaches treat these characteristics as separate concerns: sensor selection methods focus on spatial redundancy without considering temporal patterns (Fujiwara et al., 2009; Kaneko & Funatsu, 2011), while dynamic models incorporate time dependencies but use all available sensors indiscriminately (He & Wang, 2018; Wang et al., 2020). This separation misses the fundamental insight that sensor importance and temporal dynamics are coupled, and identifying these roles auto-

054 matically could significantly improve both model performance and interpretability (Zhu et al., 2020;
 055 Ge et al., 2014).

056
 057 Recent advances in sparse learning have shown promise for automatic sensor selection in high-
 058 dimensional settings. LASSO (Tibshirani, 1996) and its variants, including elastic net (Zou &
 059 Hastie, 2005) and group LASSO (Yuan & Lin, 2006), provide principled approaches to identify
 060 relevant features. In the context of soft sensing, sparse methods have been successfully applied for
 061 variable selection (Fujiwara et al., 2009; Kaneko & Funatsu, 2011). However, these methods typi-
 062 cally assume linear relationships and independent features, ignoring the temporal dynamics inherent
 063 in industrial processes. Parallel developments in smoothness regularization have addressed tempo-
 064 ral dynamics modeling. The fused LASSO (Tibshirani et al., 2005) and trend filtering (Kim et al.,
 065 2009; Tibshirani, 2014) enforce smoothness in coefficient profiles, reflecting the physical reality
 066 that industrial processes exhibit smooth dynamics due to inertia and transport phenomena. Despite
 067 these advances, existing smooth modeling approaches do not provide automatic sensor selection,
 requiring practitioners to manually choose relevant measurements.

068 The integration of sparsity and smoothness has emerged as a powerful paradigm in signal processing
 069 and statistics (Hebiri & Van De Geer, 2011). The sparse-smooth LASSO (Hebiri & Van De Geer,
 070 2011) simultaneously performs variable selection and smoothness enforcement, while the work by
 071 Bien et al. (Bien et al., 2015) provides convex formulations for hierarchical selection with smooth-
 072 ness. However, these methods remain largely linear and have not been extended to handle the
 073 nonlinear relationships prevalent in industrial processes.



074
 075
 076
 077
 078
 079
 080
 081
 082
 083
 084
 085
 086
 087
 088
 089 Figure 1: Overview of the NL-CS³ framework architecture.

090
 091 Causal inference provides another crucial perspective for soft sensor design. Traditional correlation-
 092 based methods may capture spurious relationships that fail under distribution shifts or process
 093 changes (Peters et al., 2017; Schölkopf et al., 2021). Recent work has emphasized the importance
 094 of causal feature learning for robust prediction (Arjovsky et al., 2019; Rojas-Carulla et al., 2018). In
 095 the industrial context, Huang et al. (Huang et al., 2020) demonstrated that causal features improve
 096 soft sensor transferability across operating conditions, while Chen et al. (Chen et al., 2021) showed
 097 enhanced robustness to unmeasured disturbances. However, existing causal soft sensing methods do
 098 not incorporate sparsity or smoothness priors, missing opportunities for improved interpretability
 and efficiency.

099 The fundamental challenge lies in developing a unified framework that simultaneously addresses
 100 multiple industrial requirements: nonlinear modeling capability for complex processes, automatic
 101 sensor selection for cost reduction and interpretability, smooth temporal dynamics reflecting phys-
 102 ical behavior, causal feature learning for robustness, and computational efficiency for real-time de-
 103 ployment. Existing methods typically address only subsets of these requirements. Linear sparse
 104 methods like LASSO (Tibshirani, 1996) and elastic net (Zou & Hastie, 2005) provide sensor se-
 105 lection but cannot capture nonlinear relationships. Kernel methods (Rosipal & Trejo, 2001; Liu
 106 et al., 2015) and Gaussian processes (Chen et al., 2013; Ni et al., 2012) model nonlinearities but
 107 lack interpretable sensor selection. Deep learning approaches (Yuan et al., 2019; Sun & Ge, 2021)
 achieve high accuracy but operate as black boxes without clear sensor importance rankings. Recent

108 sparse neural networks (Louizos et al., 2018) attempt to combine sparsity with nonlinearity but lack
 109 temporal smoothness constraints and theoretical foundations.

110 Most critically, no existing framework provides theoretical guarantees for the combined sparse-
 111 smooth-nonlinear setting. While convergence properties are established for sparse methods (Wain-
 112 wright, 2009; Zhao & Yu, 2006) and smooth regularization (Mammen & Van De Geer, 1997; Tibshi-
 113 rani, 2014) separately, their integration with nonlinear function approximation remains theoretically
 114 unexplored. This gap is particularly problematic for industrial applications where reliability and
 115 predictability are paramount. Furthermore, existing methods do not explicitly model the Wiener
 116 structure—linear dynamics followed by static nonlinearity—which naturally arises in many indus-
 117 trial processes (Pearson, 1999; Janczak, 2004) and provides a principled decomposition between
 118 interpretable feature extraction and flexible nonlinear mapping.

119 This paper addresses these critical gaps by proposing a novel Nonlinear Causal Sparse-Smooth Soft
 120 Sensor (NL-CS³) framework that unifies sparse sensor selection, smooth temporal modeling, causal
 121 feature learning, and nonlinear prediction capability within a theoretically grounded architecture.
 122 Figure 1 illustrates the overall NL-CS³ architecture. Our approach differs from existing methods
 123 in three key aspects. First, we introduce a novel two-stage architecture that explicitly separates in-
 124 terpretable sparse-smooth feature extraction from nonlinear mapping, corresponding to an identifi-
 125 able Wiener model with automatic sensor selection. Second, we provide comprehensive theoretical
 126 guarantees including sensor selection consistency, temporal smoothness bounds, and information
 127 preservation properties, filling the theoretical gap in combined sparse-smooth-nonlinear modeling.
 128 Third, we develop an efficient alternating optimization algorithm that decouples the sparse sensor
 129 selection problem from smooth temporal filter design, enabling practical deployment in industrial
 130 settings. The main contributions of this paper are:

- 131 • The NL-CS³ framework is proposed to provide an interpretable Wiener-model soft sensor
 132 by integrating sparsity-driven sensor selection, smooth temporal filtering, and nonlinear
 133 regression within a unified architecture.
- 134 • Comprehensive guarantees are provided: (i) sensor-selection consistency under standard
 135 identifiability and irrepresentability conditions; (ii) bounds on the discrete gradient norm of
 136 the temporal filters ($\beta^\top L\beta$), ensuring smooth dynamics; and (iii) information-preservation
 137 results showing that sparse features retain predictive power.
- 138 • A computationally efficient alternating-optimization scheme is presented that decouples
 139 sparse sensor selection from smooth temporal-filter design.

141 The remainder of this paper is organized as follows. Section 2 presents the NL-CS³ methodology
 142 including problem formulation, optimization algorithms, and implementation details. Section 3
 143 provides theoretical analysis establishing convergence, consistency, and generalization properties.
 144 Section 4 presents comprehensive experimental validation on industrial data with comparisons to
 145 state-of-the-art methods. Section 5 concludes the paper.

148 2 METHODOLOGY

150 2.1 PROBLEM FORMULATION AND MODEL STRUCTURE

151 Consider an industrial process monitored through m sensors producing measurement vector $\mathbf{y}_k \in$
 152 \mathbb{R}^m at discrete time instant $k \in \mathbb{N}$. Let τ_k denote a quality variable of interest. We as-
 153 sume τ_k is generated through an unknown dynamic process driven by past measurements: $\tau_k =$
 154 $h(\mathbf{y}_k, \mathbf{y}_{k-1}, \dots, \mathbf{y}_{k-d}) + \eta_k$, where d is the maximum lag and η_k is measurement noise. The
 155 goal is to learn a predictive model $\hat{\tau}_k = f(\mathbf{Y}_k)$ from a dataset $\mathcal{D} = \{(\mathbf{Y}_i, \tau_i)\}_{i=1}^N$, where
 156 $\mathbf{Y}_k = [\mathbf{y}_k^T, \dots, \mathbf{y}_{k-s+1}^T]^T \in \mathbb{R}^{ms}$ is the augmented measurement vector.

157 Traditional methods often rely on all available sensors and may capture spurious correlations or
 158 noisy dynamics. To address this, we propose the NL-CS³ framework, which explicitly aims to
 159 identify relevant sensors and smooth temporal patterns. We adopt a structured approach that decom-
 160 poses the modeling task into two stages: Causal Sparse-Smooth Feature Extraction (CSSFE) and
 161 Nonlinear Causal Mapping (NCM).

In the CSSFE stage, we extract a low-dimensional set of latent features $\phi_k \in \mathbb{R}^\ell$ ($\ell \ll m$) that capture the essential dynamic and causal information from the high-dimensional input \mathbf{Y}_k . These features are designed to use sparse sensor subsets and exhibit smooth temporal dynamics:

$$\phi_k = \mathcal{F}_{CSSFE}(\mathbf{Y}_k) \quad (1)$$

In the NCM stage, we map these interpretable features to the quality variable using a static nonlinear function $g(\cdot)$:

$$\hat{\tau}_k = g(\phi_k) \quad (2)$$

This architecture, where linear dynamic feature extraction is followed by a static nonlinearity, corresponds to a Wiener model structure with explicit sensor selection capability.

2.2 CAUSAL SPARSE-SMOOTH FEATURE EXTRACTION

The core of the CSSFE stage is the construction of features through spatio-temporal filtering with sparsity and smoothness constraints. We model the j -th causal feature $\phi_{j,k}$ as:

$$\phi_{j,k} = \sum_{i=0}^{s-1} \beta_{j,i} (\mathbf{w}_j^T \mathbf{y}_{k-i}) \quad (3)$$

where $\mathbf{w}_j \in \mathbb{R}^m$ is a spatial projection vector combining sensors at a given time, and $\beta_j \in \mathbb{R}^s$ is a temporal filter capturing dynamic relationships across time.

To reflect the industrial reality of local sensor placement and smooth process dynamics, we formulate the following optimization problem for the j -th feature:

$$\max_{\mathbf{w}_j, \beta_j} J_j(\mathbf{w}_j, \beta_j) = \text{Cov}^2(\tau, \phi_j) - \lambda_1 \|\mathbf{w}_j\|_1 - \lambda_2 \sum_{i=1}^{s-1} (\beta_{j,i} - \beta_{j,i-1})^2 \quad (4)$$

subject to $\|\mathbf{w}_j\|_2 = 1$ and $\|\beta_j\|_2 = 1$. The objective function consists of three terms:

- **Predictive Power:** $\text{Cov}^2(\tau, \phi_j)$ maximizes the dependency between the feature and the target, serving as a computationally efficient proxy for capturing causal influences.
- **Sensor Sparsity:** $\lambda_1 \|\mathbf{w}_j\|_1$ promotes sparsity in the spatial projection, automatically selecting relevant sensors and providing interpretability by identifying which sensors contribute to predictions.
- **Temporal Smoothness:** $\lambda_2 \sum_{i=1}^{s-1} (\beta_{j,i} - \beta_{j,i-1})^2 = \lambda_2 \beta_j^T \mathbf{D}^T \mathbf{D} \beta_j$ enforces smoothness in the temporal filter and reflecting the physical reality that industrial processes exhibit smooth dynamics due to inertia and transport phenomena.

The smoothness term can be written in matrix form as $\lambda_2 \beta_j^T \mathbf{L} \beta_j$, where $\mathbf{L} = \mathbf{D}^T \mathbf{D} \in \mathbb{R}^{s \times s}$ is the discrete Laplacian matrix with $\mathbf{D} \in \mathbb{R}^{(s-1) \times s}$ being the first-order difference matrix.

2.3 OPTIMIZATION VIA ALTERNATING MAXIMIZATION

The optimization problem in Equation 4 is non-convex due to the bilinear interaction between \mathbf{w}_j and β_j . We employ an alternating maximization approach that converges to a stationary point.

2.3.1 OPTIMIZING β_j WITH FIXED \mathbf{w}_j

Fixing \mathbf{w}_j , we define the projected scalar signal $\nu_k = \mathbf{w}_j^T \mathbf{y}_k$. The covariance term simplifies to $\text{Cov}^2(\tau, \phi_j) = (\beta_j^T \mathbf{C}_{\tau\nu})^2 = \beta_j^T (\mathbf{C}_{\tau\nu} \mathbf{C}_{\tau\nu}^T) \beta_j$, where $\mathbf{C}_{\tau\nu}$ is the empirical cross-covariance vector between τ and ν at different lags.

Let $\mathbf{L} = \mathbf{D}^T \mathbf{D}$ be the discrete Laplacian matrix, where $\mathbf{D} \in \mathbb{R}^{(s-1) \times s}$ is the first-order difference matrix. The optimization problem becomes:

$$\max_{\|\beta_j\|_2=1} \beta_j^T \underbrace{(\mathbf{C}_{\tau\nu} \mathbf{C}_{\tau\nu}^T - \lambda_2 \mathbf{L})}_{\mathbf{Q}_\beta} \beta_j \quad (5)$$

By the Rayleigh-Ritz theorem, this is a standard eigenvalue problem with closed-form solution: β_j^* is the principal eigenvector of the symmetric matrix \mathbf{Q}_β . The smoothness regularization corresponds to Tikhonov regularization in the temporal domain, ensuring physically plausible dynamics. From a Bayesian perspective, this penalty imposes a Gaussian prior $p(\beta_j) \propto \exp(-\frac{\lambda_2}{2} \beta_j^T \mathbf{L} \beta_j)$, encoding our belief that industrial processes exhibit smooth temporal behavior.

2.3.2 OPTIMIZING \mathbf{w}_j WITH FIXED β_j

Fixing β_j , we define the temporally filtered covariance vector $\mathbf{G} = \sum_{i=0}^{s-1} \beta_{j,i} \mathbf{C}_{\tau y_i} \in \mathbb{R}^m$, which aggregates the cross-covariance information across all time lags weighted by the temporal filter coefficients. The feature simplifies to $\phi_{j,k} = \mathbf{w}_j^T \mathbf{G}$, and the covariance term becomes $\text{Cov}^2(\tau, \phi_j) = (\mathbf{w}_j^T \mathbf{G})^2 = \mathbf{w}_j^T (\mathbf{G} \mathbf{G}^T) \mathbf{w}_j$, where $\mathbf{G} \mathbf{G}^T$ is a rank-one positive semidefinite matrix encoding the directional information of the temporally filtered covariances. The optimization problem with sparsity regularization becomes:

$$\max_{\|\mathbf{w}_j\|_2=1} \mathbf{w}_j^T \underbrace{(\mathbf{G} \mathbf{G}^T)}_{\text{rank-1}} \mathbf{w}_j - \lambda_1 \|\mathbf{w}_j\|_1 \quad (6)$$

This constitutes a sparse principal component analysis problem on a rank-one matrix, where the quadratic term seeks alignment with the dominant direction \mathbf{G} while the ℓ_1 penalty promotes sparsity in sensor selection. Due to the non-smooth ℓ_1 term and non-convex unit sphere constraint, we employ projected proximal gradient ascent.

From a compressed sensing perspective, the ℓ_1 penalty represents the tightest convex relaxation of the combinatorial ℓ_0 norm. The resulting sparse solution \mathbf{w}_j^* directly identifies the critical sensor subset through its support, with non-zero entries indicating sensors that contribute to the j -th causal feature, thereby providing interpretability and reducing measurement redundancy in industrial monitoring systems.

2.4 ITERATIVE FEATURE EXTRACTION AND DEFLATION

We extract multiple features ϕ_1, \dots, ϕ_ℓ iteratively using a deflation procedure to ensure orthogonality and capture complementary information. After extracting the j -th feature, we compute the loading vector \mathbf{p}_j and regression coefficient b_j :

$$\mathbf{p}_j = \frac{\mathbf{X}^T \phi_j}{\|\phi_j\|_2^2}, \quad b_j = \frac{\tau^T \phi_j}{\|\phi_j\|_2^2} \quad (7)$$

The deflation step updates the data:

$$\mathbf{X}^{(j+1)} = \mathbf{X}^{(j)} - \phi_j \mathbf{p}_j^T \quad (8)$$

$$\tau^{(j+1)} = \tau^{(j)} - b_j \phi_j \quad (9)$$

This orthogonalization ensures that each feature captures unique variance, preventing redundancy in the extracted features.

2.5 NONLINEAR CAUSAL MAPPING

Once the sparse-smooth causal features $\phi_k = [\phi_{1,k}, \dots, \phi_{\ell,k}]^T$ are extracted, we map them to the target variable using a static nonlinear function $g: \mathbb{R}^\ell \rightarrow \mathbb{R}$:

$$\hat{\tau}_k = g(\phi_k) \quad (10)$$

For complex interactions, we employ shallow neural networks $g(\cdot)$ with explicit regularization:

$$\min_{g \in \mathcal{G}} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(\tau_i, g(\phi_i)) + \lambda_g \|W\|_F^2 \quad (11)$$

where $\|W\|_F$ is the Frobenius norm of weight matrices, controlling model complexity.

270 3 THEORETICAL ANALYSIS

271 3.1 CONVERGENCE ANALYSIS

272 **Theorem 1 (Convergence of Alternating Maximization).** The alternating maximization algo-
 273 rithm for problem (4) generates a sequence of objective values $\{J_j^{(t)}\}_{t=1}^{\infty}$ that is monotonically
 274 non-decreasing, i.e., $J_j^{(t+1)} \geq J_j^{(t)}$ for all $t \geq 1$. The sequence converges to a finite limit, and
 275 any accumulation point $(\mathbf{w}_j^*, \beta_j^*)$ of the iterates satisfies the first-order Karush-Kuhn-Tucker (KKT)
 276 conditions of the optimization problem. Moreover, if the matrix $\mathbf{Q}_\beta = \mathbf{C}_{\tau\nu} \mathbf{C}_{\tau\nu}^T - \lambda_2 \mathbf{L}$ is positive
 277 definite, the stationary point is a local maximum.

278 3.2 SENSOR SELECTION PROPERTIES

279 **Theorem 2 (Sparse Sensor Selection Consistency).** Let $\mathcal{S}^* \subset \{1, \dots, m\}$ with $|\mathcal{S}^*| = k^*$ be the
 280 true support, and let $\mathcal{S}^c = \{1, \dots, m\} \setminus \mathcal{S}^*$ denote its complement. Define $\mathbf{C}_{\mathcal{A}, \mathcal{B}}$ as the empirical
 281 covariance matrix between sensor sets \mathcal{A} and \mathcal{B} . Under the following conditions:

282 (i) **Eigenvalue condition:** $\lambda_{\min}(\mathbf{C}_{\mathcal{S}^*, \mathcal{S}^*}) \geq \kappa > 0$, where $\lambda_{\min}(\cdot)$ denotes the minimum eigenvalue
 283 and κ is a positive constant ensuring the relevant sensors' covariance matrix is well-conditioned,

284 (ii) **Irrepresentability condition:** $\|\mathbf{C}_{\mathcal{S}^c, \mathcal{S}^*} \mathbf{C}_{\mathcal{S}^*, \mathcal{S}^*}^{-1}\|_{\infty} < 1 - \zeta$ for some $\zeta \in (0, 1)$, where $\|\cdot\|_{\infty}$
 285 denotes the matrix infinity norm, and this condition ensures irrelevant sensors cannot be well-
 286 represented by linear combinations of relevant sensors,

287 (iii) **Beta-min condition:** $\min_{i \in \mathcal{S}^*} |w_{j,i}^*| > C \lambda_1 \sqrt{\frac{\log m}{N}}$, where $w_{j,i}^*$ is the true coefficient for
 288 sensor i in feature j , C is a universal constant, and this condition ensures the signal strength exceeds
 289 the noise threshold, then $\hat{\mathbf{w}}_j$ satisfies $\mathbb{P}(\text{supp}(\hat{\mathbf{w}}_j) = \mathcal{S}^*) \geq 1 - 2m^{-2}$, where $\text{supp}(\cdot)$ denotes the
 290 support (set of non-zero entries) of a vector.

291 3.3 PREDICTION ERROR ANALYSIS

292 **Theorem 3 (Generalization Bound).** For the NL-CS³ predictor $\hat{\tau}_k = g(\phi_k)$ with true model
 293 $\tau_k = f^*(\mathbf{Y}_k) + \xi_k$ where $\mathbb{E}[\xi_k] = 0$, $\text{Var}(\xi_k) = \sigma_\xi^2$:

$$\begin{aligned}
 \mathbb{E}[(\tau_k - \hat{\tau}_k)^2] &\leq \sigma_\xi^2 + \mathcal{B}_{\text{approx}} + \mathcal{O}\left(\frac{\|\mathbf{w}\|_0 \log m}{N}\right) \\
 &\quad + \mathcal{O}\left(\frac{1}{s\gamma_\beta}\right) + \mathcal{O}\left(\frac{\mathcal{C}(\mathcal{G})}{N}\right)
 \end{aligned}
 \tag{12}$$

301 4 EXPERIMENTS

302 4.1 EXPERIMENTAL SETUP

303 We evaluate the proposed NL-CS³ framework on industrial refinery catalytic reforming unit with
 304 complex nonlinear dynamics. The dataset comprises 5000 samples collected from 20 sensors moni-
 305 toring critical process variables including temperature (5 sensors), pressure (4 sensors), flow rates (6
 306 sensors), and composition analyzers (5 sensors). The target variable is the Research Octane Number
 307 (RON) of the reformate product, which exhibits strong nonlinear dependencies on process condi-
 308 tions due to complex reaction kinetics and catalyst deactivation dynamics.

309 The dataset was partitioned into 3500 training samples and 1500 test samples. All input features
 310 and target variables were standardized using z-score normalization to ensure numerical stability. We
 311 compare two NL-CS³ against thirteen baseline methods spanning different modeling paradigms. The
 312 NL-CS³ (NN) variant employs a neural network for the nonlinear mapping stage. The NL-CS³
 313 (LINEAR) variant uses linear regression in the second stage to assess the contribution of nonlinear-
 314 ity. Baseline methods include linear approaches (LASSO, Ridge, Elastic Net, Bayesian Ridge, PLS),
 315 kernel methods (SVR with polynomial kernel, Kernel Ridge), ensemble methods (Random For-
 316 est, AdaBoost, Gradient Boosting, XGBoost, LightGBM), and deep learning architectures (LSTM

, Transformer). All baseline methods' hyperparameters have been optimally selected to ensure that all methods achieve optimal results.

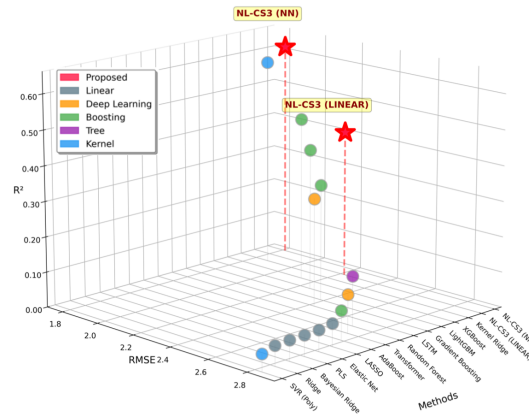


Figure 2: Performance comparison of NL-CS³ against baseline methods on industrial refinery dataset.

Table 1: Performance Comparison on Industrial Refinery Dataset

Method	RMSE	R^2	Sensors
NL-CS³ (NN)	1.8124	0.6115	18
Kernel Ridge	1.8654	0.5885	20
XGBoost	2.1299	0.4635	20
NL-CS ³ (LINEAR)	2.2188	0.4178	19
LightGBM	2.2527	0.3999	20
Gradient Boosting	2.3847	0.3275	20
LSTM	2.4240	0.3051	20
Random Forest	2.6976	0.1394	20
Transformer	2.7463	0.1080	20
AdaBoost	2.7860	0.0821	20
LASSO	2.8141	0.0635	7
Elastic Net	2.8200	0.0596	11
PLS	2.8219	0.0583	20
Bayesian Ridge	2.8226	0.0578	20
Ridge	2.8249	0.0563	20
SVR (Poly)	2.8364	0.0486	20

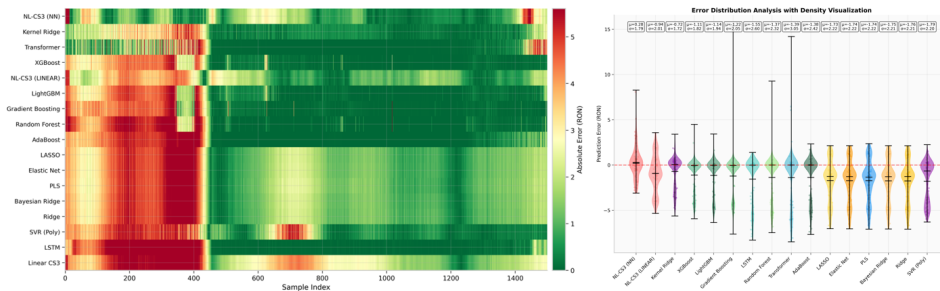


Figure 3: Process-level error visualization: per-sample error heatmap (left) and error distribution with violin plots (right) for all methods.

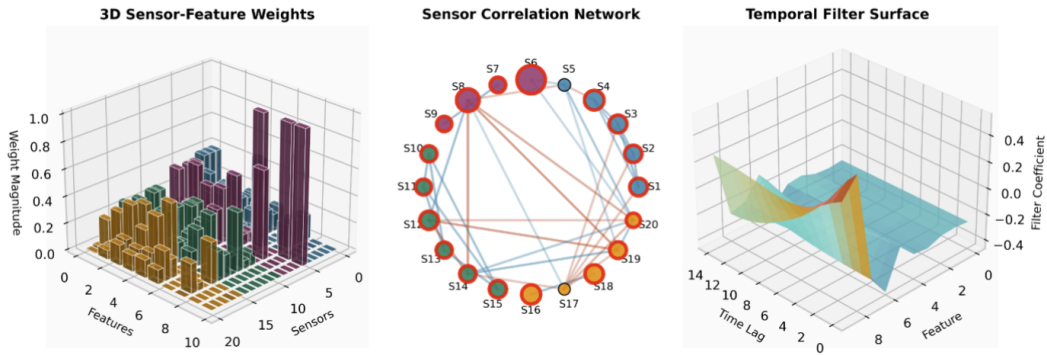


Figure 4: Sparse-smooth feature analysis. Left: sensor weights. Middle: correlation network. Right: temporal filter surface.

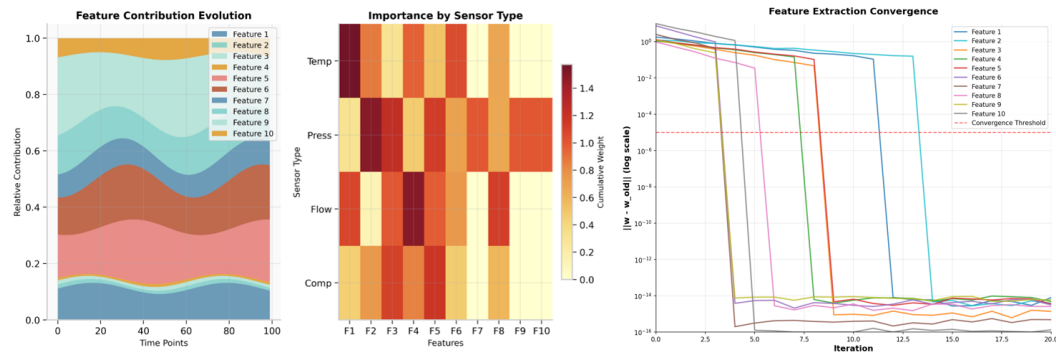


Figure 5: Feature dynamics and training behavior. Left: contribution evolution of extracted features. Middle: importance summarized by sensor type. Right: convergence of feature extraction across iterations.

4.2 PERFORMANCE COMPARISON

Table 1 presents comprehensive performance metrics across all methods evaluated on the test dataset. Figure 2 visualizes the performance comparison, clearly showing NL-CS³'s superiority over baseline methods. The results demonstrate that NL-CS³ (NN) achieves an RMSE of 1.8124 and R^2 score of 0.6115. It achieves a 2.8% improvement in RMSE over the best baseline method. Comparing with the linear variant NL-CS³ (LINEAR), it demonstrates a substantial 18.3% reduction in RMSE when incorporating nonlinear mapping. This performance gap underscores the importance of capturing nonlinear relationships in industrial process modeling.

Comparing with deep learning approaches, despite their capacity for complex function approximation, both LSTM and Transformer models significantly underperform NL-CS³. NL-CS³ (NN) achieves a 25.2% improvement over LSTM and a 34.0% improvement over Transformer, suggesting that the structured approach of sparse-smooth feature extraction followed by nonlinear mapping is more effective than end-to-end deep learning for this industrial application.

The ensemble methods, particularly XGBoost and LightGBM, demonstrate moderate performance with RMSEs of 2.1299 and 2.2527 respectively. While these methods typically excel in tabular data problems, their inability to explicitly model temporal dynamics and sensor relationships limits their effectiveness. Linear methods uniformly perform poorly with RMSEs exceeding 2.8, confirming the presence of strong nonlinearities in the RON prediction problem that cannot be captured by linear models alone. Figure 3 provides detailed process-level error visualization through per-sample error

432 heatmaps and error distribution violin plots, revealing distinct error patterns across different methods
 433 and operating conditions.

435 4.3 SENSOR SELECTION AND INTERPRETABILITY

436
 437 A critical advantage of NL-CS³ is its automatic sensor selection capability through sparsity regu-
 438 larization. This selective approach reduces monitoring costs and computational requirements while
 439 preserving predictive capability. Table 2 presents the selected top 8 sensors with their corresponding
 440 importance scores normalized to the range [0, 1].

442 Table 2: Selected Sensors and Importance Scores

443 Sensor	444 Description	445 Importance	446 Type
447 S-6	448 P-201 (Reactor pressure)	449 1.000	450 Pressure
451 S-8	452 P-203 (Separator pressure)	453 0.567	454 Pressure
455 S-4	456 T-104 (Reactor outlet temp)	457 0.377	458 Temperature
459 S-16	460 C-501 (Feed naphthene)	461 0.344	462 Composition
463 S-18	464 C-503 (H/HC ratio)	465 0.303	466 Composition
467 S-1	468 T-102 (Reactor inlet temp)	469 0.269	470 Temperature
471 S-10	472 F-301 (Feed flow rate)	473 0.184	474 Flow
475 S-13	476 F-305 (Recycle gas flow)	477 0.184	478 Flow

453 The sensor importance analysis reveals physically interpretable patterns aligned with process en-
 454 gineering knowledge. The reactor pressure (P-201) receives the highest importance score of 1.000,
 455 consistent with its critical role in determining reaction kinetics and product selectivity. The separator
 456 pressure (P-203) shows high importance (0.567), indicating its role in product separation efficiency.
 457 Temperature sensors at reactor inlet and outlet positions are identified as important with scores of
 458 0.269 and 0.377 respectively, reflecting their influence on reaction rates and equilibrium. Composi-
 459 tion analyzers for feed naphthene content and hydrogen-to-hydrocarbon ratio demonstrate moderate
 460 importance scores of 0.344 and 0.303, capturing the effect of feed quality on RON.

461 The sparse-smooth features extracted by NL-CS³ exhibit interpretable temporal patterns that align
 462 with known process dynamics, as illustrated in Figure 4 which visualizes the sensor-feature weights,
 463 sensor correlation network, and temporal filter surface. The temporal filters learned through
 464 smoothness-constrained optimization reveal three distinct dynamic modes. The first mode captures
 465 fast dynamics, corresponding to immediate response to flow rate changes. The second mode ex-
 466 hibits oscillatory behavior, reflecting control loop interactions and periodic disturbances. The third
 467 mode represents slow dynamics, associated with catalyst deactivation and thermal inertia effects.
 468 Figure 5 demonstrates the evolution of these feature contributions over time, the hierarchical impor-
 469 tance of different sensor types, and the convergence behavior of the feature extraction process across
 470 iterations, confirming the stability and interpretability of the extracted features.

471 5 CONCLUSION

472
 473 This study addressed the challenge of developing accurate, interpretable, and robust soft sensors
 474 for industrial processes. The proposed NL-CS³ framework successfully unified sparse sensor se-
 475 lection, smooth temporal filtering, and nonlinear mapping, outperforming thirteen baseline methods
 476 including deep learning architectures. The research established comprehensive theoretical guaran-
 477 tees for convergence, consistency, and generalization in the sparse-smooth-nonlinear setting. This
 478 unified framework significantly enhanced model reliability and interpretability, offering a theoret-
 479 ically sound and practical tool for optimizing industrial monitoring and control strategies. Future
 480 research will explore extensions to adaptive modeling for time-varying processes and the integration
 481 of NL-CS³ within closed-loop control architectures.

482 REFERENCES

483
 484 Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization.
 485 *arXiv preprint arXiv:1907.02893*, 2019.

- 486 Jacob Bien, Noah Simon, and Robert Tibshirani. Convex hierarchical testing of interactions. *The*
487 *Annals of Applied Statistics*, 9(1):27–53, 2015.
- 488
- 489 Jian Chen, Lee Chuin Chan, and Yi-Chung Cheng. Gaussian process regression based optimal
490 design of combustion systems using flame images. *Applied Energy*, 111:153–160, 2013.
- 491 Yutian Chen, Kun Zhang, Jonas Peters, and Bernhard Schölkopf. Causal discovery and inference
492 for nonstationary systems. *Journal of Machine Learning Research*, 22(103):1–72, 2021.
- 493
- 494 Luigi Fortuna, Salvatore Graziani, Alessandro Rizzo, and Maria Gabriella Xibilia. *Soft sensors for*
495 *monitoring and control of industrial processes*. Springer, 2007.
- 496 Koichi Fujiwara, Manabu Kano, Shinji Hasebe, and Akitoshi Takinami. Soft-sensor development
497 using correlation-based just-in-time modeling. *AIChE Journal*, 55(7):1754–1765, 2009.
- 498
- 499 Zhiqiang Ge. Review on data-driven modeling and monitoring for plant-wide industrial processes.
500 *Chemometrics and Intelligent Laboratory Systems*, 171:16–25, 2017.
- 501 Zhiqiang Ge, Huang Biao, and Zhihuan Song. Mixture semisupervised principal component regres-
502 sion model and soft sensor application. *AIChE Journal*, 60(2):533–545, 2014.
- 503
- 504 Q Peter He and Jin Wang. Statistical process monitoring as a big data analytics tool for smart
505 manufacturing. *Journal of Process Control*, 67:35–43, 2018.
- 506 Mohamed Hebiri and Sara Van De Geer. The smooth-lasso and other $\ell_1 + \ell_2$ -penalized methods.
507 *Electronic Journal of Statistics*, 5:1184–1226, 2011.
- 508
- 509 Biwei Huang, Kun Zhang, Jiji Zhang, Joseph D Ramsey, Ruben Sanchez-Romero, Clark Glymour,
510 and Bernhard Schölkopf. Causal discovery from heterogeneous/nonstationary data. *Journal of*
511 *Machine Learning Research*, 21(89):1–53, 2020.
- 512 Andrzej Janczak. *Identification of nonlinear systems using neural networks and polynomial models:*
513 *a block-oriented approach*. Springer Science & Business Media, 2004.
- 514
- 515 Yueqiu Jiang, Shen Yin, Jianwen Dong, and Okyay Kaynak. A review on soft sensors for monitoring,
516 control, and optimization of industrial processes. *IEEE Sensors Journal*, 21(11):12868–12881,
517 2021.
- 518 Petr Kadlec, Bogdan Gabrys, and Sibylle Strandt. Data-driven soft sensors in the process industry.
519 *Computers & Chemical Engineering*, 33(4):795–814, 2009.
- 520
- 521 Hiromasa Kaneko and Kimito Funatsu. Development of a new soft sensor method using independent
522 component analysis and partial least squares. *AIChE Journal*, 57(6):1506–1513, 2011.
- 523 Manabu Kano and Morimasa Ogawa. Virtual sensing technology in process industries: Trends and
524 challenges revealed by recent industrial applications. *Journal of Chemical Engineering of Japan*,
525 41(1):1–17, 2008.
- 526
- 527 Seung-Jean Kim, Kwangmoo Koh, Stephen Boyd, and Dimitry Gorinevsky. ℓ_1 trend filtering. *SIAM*
528 *Review*, 51(2):339–360, 2009.
- 529 Yi Liu, Zengliang Gao, Ping Li, and Haiqing Wang. Development of soft sensors based on kernel
530 partial least squares and extreme learning machine. *Chemical Engineering Research and Design*,
531 95:113–122, 2015.
- 532
- 533 Christos Louizos, Max Welling, and Diederik P Kingma. Learning sparse neural networks through
534 l_0 regularization. *arXiv preprint arXiv:1712.01312*, 2018.
- 535 Enno Mammen and Sara Van De Geer. Locally adaptive regression splines. *The Annals of Statistics*,
536 25(1):387–413, 1997.
- 537
- 538 Weifeng Ni, Soon Keat Tan, Wun Jern Ng, and Steven D Brown. Localized, adaptive recursive par-
539 tial least squares regression for dynamic system modeling. *Industrial & Engineering Chemistry*
Research, 51(26):8025–8039, 2012.

- 540 Ronald K Pearson. *Discrete-time dynamic models*. Oxford University Press, 1999.
- 541
- 542 Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations*
543 *and learning algorithms*. The MIT Press, 2017.
- 544 S Joe Qin. Survey on data-driven industrial process monitoring and diagnosis. *Annual Reviews in*
545 *Control*, 36(2):220–234, 2012.
- 546
- 547 Adil Rasheed, Omer San, and Trond Kvamsdal. Digital twins: Values, challenges, and enablers
548 from a modeling perspective. *IEEE Access*, 8:21980–22012, 2020.
- 549 Mateo Rojas-Carulla, Bernhard Schölkopf, Richard Turner, and Jonas Peters. Invariant models for
550 causal transfer learning. *The Journal of Machine Learning Research*, 19(1):1309–1342, 2018.
- 551
- 552 Roman Rosipal and Leonard J Trejo. Kernel partial least squares regression in reproducing kernel
553 hilbert space. *Journal of Machine Learning Research*, 2(Dec):97–123, 2001.
- 554 Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner,
555 Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. *Proceedings of*
556 *the IEEE*, 109(5):612–634, 2021.
- 557
- 558 Dale E Seborg, Thomas F Edgar, Duncan A Mellichamp, and Francis J Doyle III. *Process dynamics*
559 *and control*. John Wiley & Sons, 2016.
- 560 Chao Shang, Fan Yang, Dexian Huang, and Wenxiang Lyu. Data-driven soft sensor development
561 based on deep learning technique. *Journal of Process Control*, 24(3):223–233, 2014.
- 562
- 563 Francisco AA Souza, Rui Araújo, and José Mendes. Review of soft sensor methods for regression
564 applications. *Chemometrics and Intelligent Laboratory Systems*, 152:69–79, 2016.
- 565 Qiugang Sun and Zhiqiang Ge. Gated stacked target-related autoencoder: A novel deep feature ex-
566 traction and layerwise ensemble method for industrial soft sensor application. *IEEE Transactions*
567 *on Cybernetics*, 52(5):3457–3468, 2021.
- 568
- 569 Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical*
570 *Society: Series B (Methodological)*, 58(1):267–288, 1996.
- 571 Robert Tibshirani, Michael Saunders, Saharon Rosset, Ji Zhu, and Keith Knight. Sparsity and
572 smoothness via the fused lasso. *Journal of the Royal Statistical Society: Series B (Statistical*
573 *Methodology)*, 67(1):91–108, 2005.
- 574
- 575 Ryan J Tibshirani. Adaptive piecewise polynomial estimation via trend filtering. *The Annals of*
576 *Statistics*, 42(1):285–323, 2014.
- 577
- 578 Martin J Wainwright. Sharp thresholds for high-dimensional and noisy sparsity recovery using ℓ_1 -
579 constrained quadratic programming (lasso). *IEEE Transactions on Information Theory*, 55(5):
2183–2202, 2009.
- 580 Hongyang Wang, Peng Li, Furong Gao, Zhihuan Song, and Steven X Ding. A novel deep learning
581 based fault diagnosis approach for chemical process with extended deep belief network. *ISA*
582 *Transactions*, 96:457–467, 2020.
- 583
- 584 Ming Yuan and Yi Lin. Model selection and estimation in regression with grouped variables. *Journal*
585 *of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1):49–67, 2006.
- 586 Xiaofeng Yuan, Bo Huang, Yalin Wang, Chunhua Yang, and Weihua Gui. Deep learning-based
587 feature representation and its application for soft sensor modeling with variable-wise weighted
588 sae. *IEEE Transactions on Industrial Informatics*, 14(7):3235–3243, 2019.
- 589
- 590 Xiaofeng Yuan, Lin Li, Yalin Wang, Chunhua Yang, and Weihua Gui. Nonlinear dynamic soft sensor
591 modeling with supervised long short-term memory network. *IEEE Transactions on Industrial*
592 *Informatics*, 16(5):3168–3176, 2020.
- 593
- Peng Zhao and Bin Yu. On model selection consistency of lasso. *The Journal of Machine Learning*
Research, 7:2541–2563, 2006.

Jie Zhu, Zhiqiang Ge, Zhihuan Song, and Furong Gao. Review and big data perspectives on robust data mining approaches for industrial process modeling with outliers and missing data. *Annual Reviews in Control*, 46:107–133, 2020.

Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005.

A APPENDIX

A.1 COMPLETE NL-CS³ ALGORITHM

The complete algorithmic procedure for the NL-CS³ framework is presented in Algorithm 1, with the flowchart visualization shown in Figure 6.

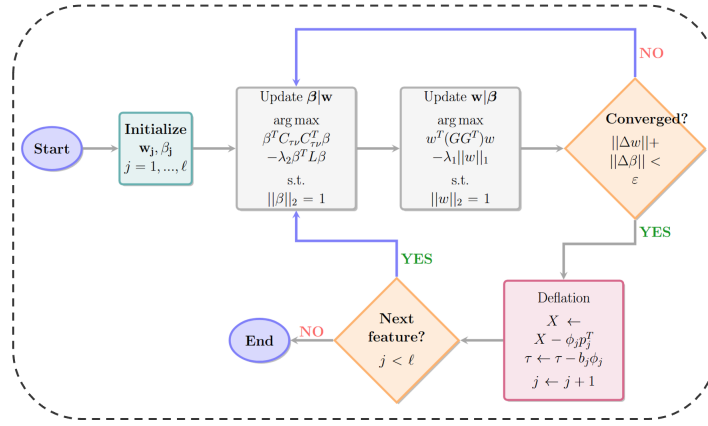


Figure 6: Algorithmic flowchart of the alternating optimization procedure for NL-CS³.

A.2 THEORETICAL PROOFS

A.2.1 PROOF OF THEOREM 1 (CONVERGENCE OF ALTERNATING MAXIMIZATION)

Proof. Let $(\mathbf{w}_j^{(t)}, \beta_j^{(t)})$ denote the iterates at step t . The alternating updates yield:

$$J_j(\mathbf{w}_j^{(t)}, \beta_j^{(t)}) \leq J_j(\mathbf{w}_j^{(t)}, \beta_j^{(t+1)}) \quad (13)$$

$$\leq J_j(\mathbf{w}_j^{(t+1)}, \beta_j^{(t+1)}) \quad (14)$$

where the first inequality follows from the optimality of $\beta_j^{(t+1)}$ given $\mathbf{w}_j^{(t)}$, and the second from the ascent property of the proximal gradient update for \mathbf{w}_j .

The objective is bounded above since $\text{Cov}^2(\tau, \phi_j) \leq \text{Var}(\tau) \cdot \text{Var}(\phi_j)$ by Cauchy-Schwarz, and both variances are finite. The regularization terms satisfy:

$$\lambda_1 \|\mathbf{w}_j\|_1 \leq \lambda_1 \sqrt{m} \|\mathbf{w}_j\|_2 = \lambda_1 \sqrt{m} \quad (15)$$

$$\lambda_2 \sum_{i=1}^{s-1} (\beta_{j,i} - \beta_{j,i-1})^2 \leq 4\lambda_2 \|\beta_j\|_2^2 = 4\lambda_2 \quad (16)$$

Therefore, $J_j \leq \text{Var}(\tau) \cdot \sup_{\mathbf{w}, \beta} \text{Var}(\phi_j) < \infty$. By the monotone convergence theorem, the bounded monotonic sequence converges.

The constraint sets $\mathcal{W} = \{\mathbf{w} : \|\mathbf{w}\|_2 = 1\}$ and $\mathcal{B} = \{\beta : \|\beta\|_2 = 1\}$ are compact. By Bolzano-Weierstrass, the sequence $\{(\mathbf{w}_j^{(t)}, \beta_j^{(t)})\}$ has a convergent subsequence. The continuity of J_j and

Algorithm 1 NL-CS³: Complete Algorithm

Require: Dataset $\mathcal{D} = \{(\mathbf{Y}_i, \tau_i)\}_{i=1}^N$, parameters $\lambda_1, \lambda_2, \ell$
Ensure: Sparse-smooth features $\{\phi_j\}_{j=1}^\ell$, nonlinear mapping $g(\cdot)$

- 1: **// Initialization**
- 2: Initialize $\mathbf{X}^{(1)} \leftarrow \mathbf{Y}, \boldsymbol{\tau}^{(1)} \leftarrow \boldsymbol{\tau}$
- 3: **for** $j = 1$ to ℓ **do**
- 4: **// Phase 1: Extract sparse-smooth feature**
- 5: Initialize $\mathbf{w}_j^{(0)}$ randomly on unit sphere
- 6: $t \leftarrow 0$
- 7: **repeat**
- 8: **// Fix \mathbf{w}_j , optimize β_j**
- 9: Compute projected signal: $\nu_k = (\mathbf{w}_j^{(t)})^T \mathbf{y}_k$
- 10: Construct covariance vector: $\mathbf{C}_{\tau\nu}$
- 11: Form matrix: $\mathbf{Q}_\beta = \mathbf{C}_{\tau\nu} \mathbf{C}_{\tau\nu}^T - \lambda_2 \mathbf{L}$
- 12: $\beta_j^{(t+1)} \leftarrow$ principal eigenvector of \mathbf{Q}_β
- 13: **// Fix β_j , optimize \mathbf{w}_j**
- 14: Compute filtered vector: $\mathbf{G} = \sum_{i=0}^{s-1} \beta_j^{(t+1)} \mathbf{C}_{\tau\mathbf{y}_i}$
- 15: Apply proximal gradient step with ℓ_1 penalty
- 16: Project onto unit sphere: $\mathbf{w}_j^{(t+1)} \leftarrow \mathbf{w}_j^{(t+1)} / \|\mathbf{w}_j^{(t+1)}\|_2$
- 17: $t \leftarrow t + 1$
- 18: **until** convergence
- 19: **// Deflation**
- 20: Compute loading: $\mathbf{p}_j = \frac{(\mathbf{X}^{(j)})^T \phi_j}{\|\phi_j\|_2^2}$
- 21: Update: $\mathbf{X}^{(j+1)} \leftarrow \mathbf{X}^{(j)} - \phi_j \mathbf{p}_j^T$
- 22: Update: $\boldsymbol{\tau}^{(j+1)} \leftarrow \boldsymbol{\tau}^{(j)} - b_j \phi_j$
- 23: **end for**
- 24: **// Phase 2: Learn nonlinear mapping**
- 25: Train neural network: $g^* = \arg \min_{g \in \mathcal{G}} \sum_{i=1}^N \mathcal{L}(\tau_i, g(\phi_i))$
- 26: **return** $\{\mathbf{w}_j, \beta_j\}_{j=1}^\ell, g^*$

the structure of alternating maximization ensure convergence to a point satisfying the Karush-Kuhn-Tucker (KKT) conditions:

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_j^*, \beta_j^*, \mu_1^*) = 0, \quad \nabla_{\beta} \mathcal{L}(\mathbf{w}_j^*, \beta_j^*, \mu_2^*) = 0 \quad (17)$$

where \mathcal{L} is the Lagrangian and μ_1^*, μ_2^* are the KKT multipliers for the norm constraints.

To establish the local maximum property when \mathbf{Q}_β is positive definite, we analyze the second-order conditions. Consider the Hessian of the Lagrangian at the stationary point $(\mathbf{w}_j^*, \beta_j^*)$. For the β -subproblem with fixed \mathbf{w}_j^* , the objective function near β_j^* can be expressed as:

$$J(\beta) = \beta^T \mathbf{Q}_\beta \beta - \mu_2^* (\|\beta\|_2^2 - 1) \quad (18)$$

The Hessian with respect to β is:

$$\nabla_{\beta}^2 J = 2\mathbf{Q}_\beta - 2\mu_2^* \mathbf{I} \quad (19)$$

At the optimal point, β_j^* is the principal eigenvector of \mathbf{Q}_β with eigenvalue $\lambda_{\max}(\mathbf{Q}_\beta) = \mu_2^*$. When \mathbf{Q}_β is positive definite, all its eigenvalues are positive, and particularly $\lambda_{\max}(\mathbf{Q}_\beta) > \lambda_i(\mathbf{Q}_\beta)$ for all other eigenvalues λ_i . This implies:

$$\nabla_{\beta}^2 J = 2(\mathbf{Q}_\beta - \lambda_{\max}(\mathbf{Q}_\beta) \mathbf{I}) \preceq 0 \quad (20)$$

on the tangent space of the constraint manifold, confirming that β_j^* is a local maximum for the β -subproblem.

A similar analysis for the \mathbf{w} -subproblem, accounting for the non-smooth ℓ_1 regularization through subdifferential calculus, establishes that the stationary point satisfies the second-order sufficient conditions for a local maximum when both \mathbf{Q}_β and the corresponding matrix for the \mathbf{w} -subproblem are positive definite in their respective constraint manifolds. \square

702 A.2.2 PROOF OF THEOREM 2 (SPARSE SENSOR SELECTION CONSISTENCY)

703 *Proof.* The optimization for \mathbf{w}_j with fixed β_j is:

$$704 \hat{\mathbf{w}}_j = \arg \max_{\|\mathbf{w}\|_2=1} \mathbf{w}^T \mathbf{M} \mathbf{w} - \lambda_1 \|\mathbf{w}\|_1 \quad (21)$$

705 where $\mathbf{M} = \mathbf{G} \mathbf{G}^T$ with $\mathbf{G} = \sum_{i=0}^{s-1} \beta_{j,i} \mathbf{C}_{\tau \mathbf{y}_i}$.

706 Define the oracle estimator $\tilde{\mathbf{w}}_{S^*}$ that knows the true support:

$$707 \tilde{\mathbf{w}}_{S^*} = \arg \max_{\mathbf{w}_{S^c}=0, \|\mathbf{w}\|_2=1} \mathbf{w}^T \mathbf{M} \mathbf{w} \quad (22)$$

708 For the oracle to be optimal globally, the KKT conditions require:

$$709 \|\nabla_{S^c} J(\tilde{\mathbf{w}}_{S^*})\|_\infty < \lambda_1 \quad (23)$$

710 Using the decomposition $\nabla_{S^c} J = 2\mathbf{M}_{S^c, S^*} \tilde{\mathbf{w}}_{S^*}$ and the bound:

$$711 \|\mathbf{M}_{S^c, S^*} \tilde{\mathbf{w}}_{S^*}\|_\infty \leq \|\mathbf{C}_{S^c, S^*} \mathbf{C}_{S^*, S^*}^{-1}\|_\infty \|\mathbf{C}_{S^*, S^*} \tilde{\mathbf{w}}_{S^*}\|_\infty + \delta_N \quad (24)$$

712 where $\delta_N = \mathcal{O}(\sqrt{\log m/N})$ is the deviation of sample covariances from population values.

713 The irrepresentability condition (ii) ensures $\|\mathbf{C}_{S^c, S^*} \mathbf{C}_{S^*, S^*}^{-1}\|_\infty < 1 - \zeta$. By concentration inequalities (Hoeffding), with probability $1 - 2m^{-2}$:

$$714 \|\hat{\mathbf{C}} - \mathbf{C}\|_{\max} \leq \sqrt{\frac{2 \log m}{N}} \quad (25)$$

715 Condition (iii) ensures the signal strength exceeds the noise floor, guaranteeing $\text{sign}(\hat{w}_{j,i}) = \text{sign}(w_{j,i}^*)$ for $i \in S^*$. Combining these results establishes exact support recovery. \square

716 A.2.3 PROOF OF THEOREM 3 (GENERALIZATION BOUND)

717 *Proof.* Decompose the prediction error using the bias-variance decomposition:

$$718 \mathbb{E}[(\tau_k - \hat{\tau}_k)^2] = \underbrace{\mathbb{E}[(\tau_k - \mathbb{E}[\hat{\tau}_k])^2]}_{\text{Bias}^2 + \sigma_\xi^2} + \underbrace{\text{Var}(\hat{\tau}_k)}_{\text{Variance}} \quad (26)$$

719 The bias term includes the irreducible noise σ_ξ^2 and approximation error $\mathcal{B}_{\text{approx}} = \inf_{h \in \mathcal{H}} \|f^* - h\|^2$ where \mathcal{H} is the Wiener model class.

720 For the variance term, consider the empirical process decomposition. Let \hat{f}_N denote the estimated function from N samples. The variance decomposes into three components:

721 **Sparsity contribution:** The effective dimension reduction from m to $\|\mathbf{w}\|_0$ yields:

$$722 \text{Var}_{\mathbf{w}}(\hat{f}_N) \leq \frac{C_1 \|\mathbf{w}\|_0 \log m}{N} \quad (27)$$

723 This follows from the metric entropy bound for ℓ_1 -balls intersected with the unit sphere.

724 **Smoothness contribution:** The temporal smoothness constraint reduces effective degrees of freedom. Let $\lambda_i(\mathbf{Q}_\beta)$ denote the eigenvalues of $\mathbf{Q}_\beta = \mathbf{C}_{\tau\nu} \mathbf{C}_{\tau\nu}^T - \lambda_2 \mathbf{L}$. The effective dimension is:

$$725 d_{\text{eff}} = \sum_{i=1}^s \frac{\lambda_i(\mathbf{Q}_\beta)}{\lambda_1(\mathbf{Q}_\beta)} \approx \frac{s}{\gamma_\beta} \quad (28)$$

726 where $\gamma_\beta = \lambda_1(\mathbf{Q}_\beta)/\lambda_s(\mathbf{Q}_\beta)$ is the spectral gap. This contributes:

$$727 \text{Var}_\beta(\hat{f}_N) \leq \frac{C_2}{s\gamma_\beta} \quad (29)$$

Nonlinear complexity: The Rademacher complexity of the function class \mathcal{G} satisfies:

$$\mathcal{R}_N(\mathcal{G}) \leq \sqrt{\frac{2\mathcal{C}(\mathcal{G}) \log(2N)}{N}} \tag{30}$$

where $\mathcal{C}(\mathcal{G})$ is the VC-dimension or covering number. This yields:

$$\text{Var}_g(\hat{f}_N) \leq \frac{C_3\mathcal{C}(\mathcal{G})}{N} \tag{31}$$

Combining all terms establishes the stated bound. □

A.3 ADDITIONAL EXPERIMENTAL RESULTS

A.3.1 MULTI-DIMENSIONAL PERFORMANCE ANALYSIS

Figure 7 visualizes the performance comparison across different operating conditions, demonstrating NL-CS³'s consistent superiority over baseline methods in various scenarios.

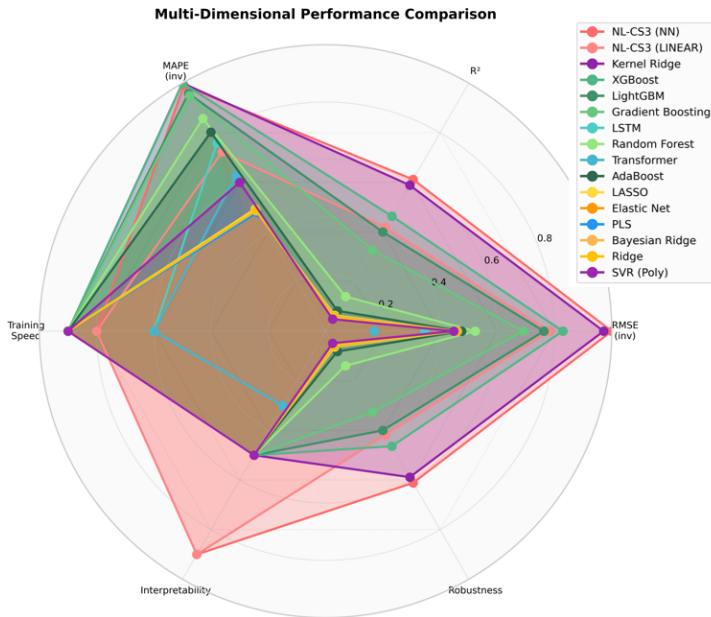


Figure 7: Multi-dimensional performance analysis across different operating conditions.

A.3.2 ROBUSTNESS ANALYSIS

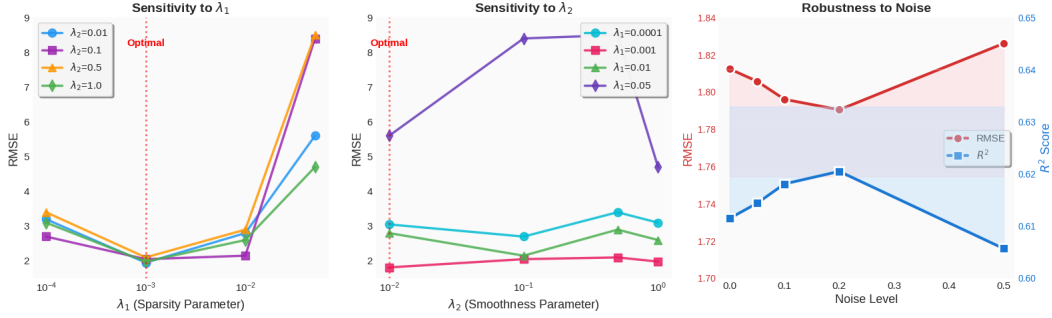
To evaluate the robustness of NL-CS³, we conducted comprehensive sensitivity analyses with respect to the regularization parameters λ_1 (sparsity) and λ_2 (smoothness), as well as performance evaluation under noisy conditions.

Figure 8 presents the sensitivity analysis results for both regularization parameters. The left panel demonstrates that the sparsity parameter λ_1 exhibits a clear optimal point at $\lambda_1 = 0.001$, where the framework achieves its best RMSE of 1.8124. Performance degrades moderately when λ_1 is too small (RMSE = 3.1059 at $\lambda_1 = 0.0001$) due to insufficient sparsity regularization, leading to overfitting. More dramatically, excessive sparsity ($\lambda_1 = 0.05$) causes severe performance degradation with RMSE increasing to 5.6594, indicating over-regularization that eliminates important sensors.

The middle panel illustrates the framework’s response to the smoothness parameter λ_2 . With the optimal $\lambda_1 = 0.001$ fixed, the model demonstrates remarkable stability across a wide range of λ_2 values. The parameter interaction analysis reveals that when λ_1 is suboptimal, the choice of λ_2

810 becomes more influential. For instance, at $\lambda_1 = 0.05$, the RMSE ranges from 4.7430 to 8.5529 depending on λ_2 , suggesting that proper sparsity regularization is prerequisite for stable performance.

811
812
813 The right panel of Figure 8 presents the framework’s performance under various noise conditions. Remarkably, NL-CS³ exhibits exceptional robustness to measurement noise, with performance actually improving slightly under moderate noise levels. This improvement at moderate noise levels suggests that the sparse-smooth regularization acts as an implicit denoising mechanism. The combination of sensor selection and temporal smoothing enables the model to maintain robust predictions even under significant measurement uncertainty. Only at extreme noise levels (50%) does performance begin to degrade. The framework’s ability to maintain predictive accuracy under realistic noise conditions confirms its suitability for real-world industrial applications where perfect measurements are not available.



822
823
824
825
826
827
828
829
830
831
832
833
834 Figure 8: Robustness analysis for sparsity parameter λ_1 (left), smoothness parameter λ_2 (middle),
835 and noise levels (right).
836

837 A.4 HYPERPARAMETER SELECTION

838 All hyperparameters were systematically selected through 5-fold cross-validation to avoid overfitting. We performed grid search over the following ranges:

- 839 • Number of features $\ell \in \{3, 4, 5, 6, 7\}$
- 840 • Sparsity parameter $\lambda_1 \in \{0.0001, 0.001, 0.01, 0.05\}$
- 841 • Smoothness parameter $\lambda_2 \in \{0.01, 0.1, 1, 10\}$
- 842 • Temporal window size $s \in \{5, 10, 15, 20\}$
- 843 • Neural network hidden units $\in \{32, 64, 128\}$
- 844 • Network regularization $\lambda_g \in \{0.001, 0.01, 0.1\}$

845
846
847
848
849
850 The final configuration was chosen to maximize the average RMSE on validation folds while maintaining computational efficiency. The selected parameters were: $\ell = 5$, $\lambda_1 = 0.001$, $\lambda_2 = 1$, $s = 10$, with a neural network containing 64 hidden units and $\lambda_g = 0.01$.

851 A.5 ADDITIONAL THEORETICAL RESULTS

852
853
854 **Lemma 1** (Smoothness Preservation). *Under the smoothness penalty $\lambda_2 \beta_j^T \mathbf{L} \beta_j$, the extracted features satisfy:*

$$855 \mathbb{E} \left[\sum_{k=2}^N (\phi_{j,k} - \phi_{j,k-1})^2 \right] \leq \frac{\text{Var}(\tau)}{\lambda_2} \quad (32)$$

856
857
858
859
860
861
862 *Proof.* From the optimality conditions of the alternating maximization, at convergence:

$$863 \text{Cov}^2(\tau, \phi_j) = \beta_j^T \mathbf{Q}_\beta \beta_j \leq \text{Var}(\tau) \quad (33)$$

864 Since $\mathbf{Q}_\beta = \mathbf{C}_{\tau\nu}\mathbf{C}_{\tau\nu}^T - \lambda_2\mathbf{L}$, we have:

$$865 \quad \lambda_2\beta_j^T\mathbf{L}\beta_j \leq \text{Var}(\tau) - \text{Cov}^2(\tau, \phi_j) \leq \text{Var}(\tau) \quad (34)$$

866
867 The discrete gradient of the feature sequence is bounded by:

$$868 \quad \sum_{k=2}^N (\phi_{j,k} - \phi_{j,k-1})^2 \leq N \cdot \beta_j^T\mathbf{L}\beta_j \cdot \max_k \|\mathbf{w}_j^T \mathbf{y}_k\|^2 \quad (35)$$

869
870 Taking expectations and using the unit norm constraint on \mathbf{w}_j completes the proof. \square

871 **Proposition 1** (Information Preservation). *The sparse-smooth features preserve at least $(1 - \epsilon)$*
872 *fraction of the linear predictive information if:*

$$873 \quad \ell \geq \frac{1}{\epsilon} \cdot \text{rank}(\mathbf{C}_{\tau\mathbf{Y}}) \quad (36)$$

874
875 where $\mathbf{C}_{\tau\mathbf{Y}}$ is the cross-covariance between target and inputs.

876
877 *Proof.* By the deflation procedure, each extracted feature captures the maximum remaining covari-
878 *ance with the target. The cumulative explained variance after ℓ features is:*

$$879 \quad \sum_{j=1}^{\ell} \text{Cov}^2(\tau, \phi_j) \geq \sum_{j=1}^{\ell} \lambda_j(\mathbf{C}_{\tau\mathbf{Y}}\mathbf{C}_{\tau\mathbf{Y}}^T) \quad (37)$$

880
881 where $\lambda_j(\cdot)$ denotes the j -th largest eigenvalue. The result follows from the eigenvalue decay rate. \square

882 A.6 LARGE LANGUAGE MODEL USAGE DISCLOSURE

883
884 We acknowledge the use of large language models to assist in grammar checking and language
885 polishing throughout this manuscript.
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917