# Exploring Gray-Scott Model with RPCA (Robust Principal Component Analysis) and Machine Learning Integration

**Rusmia Sattar**[†]

[†] *Department of Mathematics and Statistics, Memorial University of Newfoundland,*
 *St. John's (NL) A1C 5S7, Canada*

E-mail: rsattar@mun.ca

### Abstract

Data-driven mathematical modeling plays a crucial role in understanding complex systems by leveraging the wealth of information contained within large-scale datasets. This study explores the integration of data-driven techniques, including Robust Principal Component Analysis (RPCA) and machine learning, to analyze and model complex phenomena. The Gray-Scott model, a well-established reaction-diffusion system, serves as a case study to demonstrate the application of these techniques. We first simulate the Gray-Scott model to generate synthetic data representing the spatiotemporal evolution of chemical concentrations. Subsequently, we apply RPCA to decompose the data into low-rank and sparse components, enabling the identification of underlying patterns and anomalous behavior. Furthermore, we employ machine learning algorithms, such as linear regression, support vector regression (SVR), and random forest regression, to predict system dynamics and evaluate model performance. Through comprehensive experimentation and model evaluation, we showcase the efficacy of data-driven approaches in capturing complex dynamics, extracting meaningful insights, and making accurate predictions. This study underscores the importance of data-driven mathematical modeling in elucidating intricate systems and highlights the potential of RPCA and machine learning techniques in advancing scientific understanding and decision-making in various domains.

Keywords: Data-driven, Mathematical Modeling, RPCA, Gray-Scott Model

## 1 Introduction

In today's new data-driven era, the integration of mathematical modeling techniques with advanced machine learning techniques has emerged as a powerful approach to solving complex scientific problems. Such integration facilitates the analysis of complex systems governed by feedback and diffusion equations, which involve deep, predictive insights, and open to potential. In this context, the Gray-Scott model, a classic feedback diffusion system, is a cornerstone for understanding pattern formation and dynamic behavior in diverse environments.[1]

The Gray-Scott model describes the spatial and temporal evolution of two chemical interactions, usually denoted by 'u' and 'v', under the reaction-diffusion equation. These forces are partial differential equations that capture the complex interactions of reaction-diffusion and spatial heterogeneity. Despite its shortcomings, the Gray-Scott model provides a basic framework for studying model-making phenomena found in biology, chemistry, and in natural systems.[2]

To take advantage of the insights provided by the Gray-Scott model and overcome the challenges posed by high-dimensional data and noise, data-driven statistical methods such as robust principal component analysis (RPCA) comes into play. RPCA decomposes spatiotemporal data into a low- degree of smoothness. It enables the underlying features, effectively separating the

model set from noise and outliers.[3] The application of RPCA to simulated or experimental data derived from the Gray-Scott model assists to identify underlying patterns, extract meaningful features, and optimize model interpretation.

A major part of the project is the usage of Robust Principal Component Analysis (RPCA). It is a powerful information analysis method that plays a vital role in uncovering underlying structures in excessive-dimensional datasets. Unlike conventional Principal Component Analysis (PCA), which assumes that the facts can be as it should be defined by using a low-dimensional subspace corrupted by Gaussian noise, RPCA extends this framework to address information containing outliers, sparse signals, or gross errors.[4] By decomposing the statistics matrix into low-rank and sparse components, RPCA permits researchers to perceive and separate structured styles from noise or anomalies, thereby enhancing the interpretability and robustness of data-driven models.[5] This functionality makes RPCA specifically properly perfect for analyzing complicated structures wherein the underlying dynamics can be obscured by means of noise or size mistakes.

Complex structures, characterised by nonlinear interactions, emergent houses, and intricate dynamics, pose significant demanding situations to standard modeling techniques. In many cases, the underlying mechanisms governing these structures aren't fully understood or amenable to specific mathematical description.[6] Data-driven mathematical models provide a complementary method, permitting researchers to discover hidden styles and relationships embedded inside complex datasets. By leveraging advanced statistical and computational strategies, inclusive of machine learning knowledge of algorithms and dimensionality reduction strategies, researchers can clarify meaningful records from noisy or high-dimensional facts, thereby improving our expertise of complicated structures and facilitating predictive modeling and choice-making.[7]

Linear Regression, Support Vector Regression (SVR), and Random Forest are machine learning models utilized in the project for data-driven mathematical modeling, specifically in the context of the Gray-Scott model.[7] It additionally enhances the project's functionality to understand the behavior of the Gray-Scott model by allowing the exploration of both linear and nonlinear relationships within the statistics. Linear Regression is a simple statistical technique used to model the linear courting between a dependent variable and one or greater impartial variables. SVR is a sort of Support Vector Machine (SVM) algorithm used for regression obligations. It is effective in shooting nonlinear relationships between entering features and the goal variable. Random Forest is a mastering approach that is useful for coping with high-dimensional facts and taking pictures of complex interactions among capabilities.[8] All these methods are used to conduct the experiment on the Gray-Scott model and evaluate the performances of RPCA and Machine Learning Algorithm on the model.

## 2 Literature Review

Data-driven machine learning and mathematical modeling represent two pivotal pillars in contemporary scientific and technological advancements. The roots of machine learning trace back to the early days of artificial intelligence and pattern recognition, evolving significantly with the exponential growth of digital data and computational capabilities.[7]Machine learning techniques, ranging from classical algorithms like linear regression to sophisticated deep learning architectures, enable computers to discern intricate patterns and make predictions or decisions autonomously across various domains.[9]

Concurrently, mathematical modeling has been integral to scientific inquiry for centuries, providing a formal framework for representing real-world phenomena through mathematical structures and techniques. Mathematical models, often based on differential equations, optimization principles, or statistical inference, facilitate abstraction and simplification, distilling complex systems into tractable representations. These models are foundational in fields such as physics, engineering, biology, and economics, aiding in understanding phenomena, making

predictions, and optimizing processes.[8]

The intersection of data-driven machine learning and mathematical modeling heralds a synergy that transcends traditional disciplinary boundaries. Machine learning techniques offer powerful tools for predictive modeling and pattern recognition, while mathematical models provide a conceptual framework for understanding system dynamics and relationships.[10]This convergence is particularly evident in fields like system identification, where empirical data inform the development of mathematical models, and in model validation and interpretation, where machine learning methods augment traditional analytical approaches.

**Linear Regression:** Linear regression is a simple but effective statistical approach used to version the relationship between a dependent variable and independent variables. It assumes a linear relationship between the enter capabilities and the goal variable.[11] In the context of statistic based mathematical modeling, linear regression can be useful for figuring out linear patterns or trends in the statistics. In this project, linear regression is applied to the concentrations of U primarily based on the generated records from the Gray-Scott version. By fitting a linear version to the data, it depicts the underlying linear relationships among the input parameters and the concentration of U.

**Support Vector Regression (SVR)**: SVR is a kind of Support vector device (SVM) that is used for regression responsibilities. It works by finding the hyperplane that suits the records points at the same time as maximizing the margin between the hyperplane and the facts points.[9] SVR is particularly beneficial while handling nonlinear relationships among the enter capabilities and the goal variable. It can capture complicated styles and relationships inside the information by using exceptional kernel features.[12] In the context of information-driven mathematical modeling, SVR can help in taking pictures of the nonlinear dynamics of the Gray-Scott model. It allows for modeling of the connection between the input parameters and the concentration of U, potentially taking pictures of complex styles that may not be captured with the aid of linear models.

**Random Forest:** Random Forest is known for its robustness and potential to deal with excessive-dimensional statistics with complicated interactions between capabilities. Random Forests are much less vulnerable to overfitting and can capture nonlinear relationships efficaciously.[13] In data-driven mathematical modeling, Random Forest can be useful for shooting complex interactions and nonlinear dynamics gift in the Gray-Scott model. It can provide insights into the underlying relationships among the input parameters and the concentration of U by thinking about more than one selection bushes and their collective predictions.

In this project, we embark on a journey to explore the connection between data-driven mathematical modeling, Gray-Scott Model dynamics, RPCA, and machine learning models. Through a comprehensive analysis of simulated data from the Gray-Scott model and incorporating machine learning techniques, we undertake to push the limits of scientific inquiry and make contributions to the advancement of interdisciplinary studies on the intersection of mathematics, and computational technology.

# 3  Methods

The objective of the project is to showcase the integration of data-driven mathematical modeling techniques with machine learning algorithms to analyze and predict complex dynamical systems. The primary goal is to demonstrate the effectiveness of this integration in enhancing the understanding and predictive capabilities of mathematical models, particularly in scenarios where traditional modeling approaches may fall short.

To achieve this objective, the project begins by generating gray-scott model datasets for different timesteps. Next, we implement Robust Principal Component Analysis (RPCA) to decompose high-dimensional data generated from a mathematical model, such as the Gray-Scott reaction-diffusion system, into low-rank and sparse components. This decomposition allows for the extraction of essential features and patterns from the data, which can then be utilized for further analysis and modeling. Subsequently, machine learning models, including linear regression, support vector regression (SVR), and random forest regression, are applied to the decomposed data for various tasks such as model calibration, parameter estimation, and prediction. By leveraging the power of machine learning algorithms, the project aims to capture the underlying dynamics of the mathematical model more accurately and make robust predictions about system behavior.

Through rigorous evaluation and comparison of different machine learning approaches, the project seeks to assess the performance of each algorithm in capturing the complex relationships inherent in the data. This analysis provides insights into the strengths and limitations of data-driven techniques in modeling dynamical systems and informs best practices for future applications.

The reaction-diffusion equation for the Gray-Scott model are as follows:

$$\frac{\partial U}{\partial t} = \Delta U - UV^2 + F(1 - U)$$
$$\frac{\partial V}{\partial t} = \Delta V + UV^2 - (F + k)V$$

where:

$U$ : Concentration of substance $U$

$V$ : Concentration of substance $V$

$\Delta$ : Laplacian operator

$Du$ : Diffusion rate for $U$

$Dv$ : Diffusion rate for $V$

$F$ : Feed rate of $U$

$k$ : Kill rate of $V$

To simulate the Gray-Scott model numerically, various numerical methods can be employed, such as finite difference methods or finite element methods. These methods discretize the spatial domain into a grid and approximate the partial derivatives using finite difference approximations. Time integration techniques like the explicit Euler method or implicit methods are then applied to evolve the concentrations over time. The workflow of the project starts with data generation from the Gray-Scott model, reshaping the synthetic dataset, application of RPCA for decomposition, preprocessing of data, and lastly, selection and evaluation of machine learning models. This methodology enables the exploration and understanding of the underlying dynamics of the Gray-Scott model data using data-driven approaches.

**Gray-Scott Model Simulation and Data Generation:** The Gray-Scott model simulation function is a pivotal component of the project as it generates synthetic data that mirrors the behavior of chemical concentrations over time according to the Gray-Scott reaction-diffusion equations. This data serves as the foundation for subsequent analysis and modeling. The function iteratively applies the equations to update the concentrations of the U and V variables

across a spatial grid over multiple time steps, capturing the dynamics of chemical reactions and diffusion processes. By simulating the Gray-Scott model, researchers gain insights into complex spatial-temporal patterns that emerge in reaction-diffusion systems, which are crucial for understanding phenomena like pattern formation and morphogenesis in biological and chemical systems. A function named gray_scott_model is defined to simulate the Gray-Scott model and generate synthetic data. This function initializes concentrations for variables U and V, applies the reaction-diffusion equations over a specified number of time steps, and returns the simulated data as a 4-dimensional array representing concentrations of U and V over time.

**Robust Principal Component Analysis (RPCA):** Robust Principal Component Analysis (RPCA) is employed to decompose the generated data from the Gray-Scott model into low-rank and sparse components. The RPCA algorithm utilizes Singular Value Thresholding (SVT) to iteratively update the low-rank and sparse components until convergence. SVT operates by performing a singular value decomposition (SVD) on the input data matrix and applying a thresholding operation to the singular values, effectively promoting sparsity in the resulting components.

This function takes an input data matrix X and performs Robust PCA on it. It initializes the low-rank component L, sparse component S, and Lagrange multiplier Y as zero matrices of appropriate dimensions. It calculates the penalty parameter mu based on the input data matrix X. Inside the main iteration loop, it updates L, S, and Y until convergence or reaching the maximum number of iterations. It uses Singular Value Thresholding (SVT) for updating L. Finally, it returns the low-rank component L and the sparse component S.

**Adding Random Noise to the Data**:Noise Generation: Random Gaussian noise, $f \sim N(0, \sigma^2)$ with mean zero and configurable standard deviation simulates common data disturbances. This step aligns the study with real-world data scenarios, ensuring our findings are practically applicable. Utilizing the numpy library, we developed the add_noise(data, noise_level) function. This function blends the original dataset with generated noise, based on the specified noise intensity enhancing the dataset's realism.

```
data_noisy_U = add_noise(U, noise_level=0.2)
data_noisy_V = add_noise(V, noise_level=0.2)
```

The RPCA function takes the input data matrix(X), along with optional parameters such as regularization (lam), convergence tolerance (tol), and maximum iterations (max_iter), and returns the low-rank component (L) and sparse component (S). In the context of the provided code, Robust Principal Component Analysis (RPCA) serves as a powerful technique for decomposing high-dimensional data generated from the Gray-Scott model into its constituent components: low-rank and sparse. The primary motivation behind using RPCA lies in its ability to separate structured patterns (low-rank) from sparse noise or outliers in the data. Robust Principal Component Analysis (RPCA) is applied to decompose high-dimensional data, like that generated from the Gray-Scott model, into low-rank and sparse components. RPCA is valued for its ability to distinguish between structured patterns and sporadic noise or outliers within the data. The low-rank component signifies the main structures or underlying patterns, while the sparse component isolates noise, outliers, or deviations. Specifically for the Gray-Scott model, RPCA aims to clearly differentiate the model's essential dynamics from any present irregularities, thereby aiding in a deeper understanding of the system's behavior and improving the precision of further analysis or predictions.

**Reshape Data for RPCA:**
The generated synthetic data is reshaped into a 2D matrix data_2d suitable for input to the RPCA function.

```
data_combined = np.vstack((U.reshape(1, -1), V.reshape(1, -1)))
```

**Apply RPCA to Data:** The RPCA function (robust_pca) is applied to the reshaped data, resulting in the low-rank component (low_rank) and the sparse component (sparse).

```
data_noisy = add_noise(data_combined, noise_level=0.2)
low_rank, sparse = robust_pca(data_noisy, lam=1.0, tol=1e-7, max_iter=100, noise_st=0.2)
```

**Visualization of Results** A function visualize_components is defined to visualize the low-rank and sparse components for a specific time step. This function takes the low-rank component, sparse component, grid dimensions (nx, ny), and a specific time step as input. It plots the low-rank and sparse components side by side using matplotlib.

```
def visualize_components(original, low_rank, sparse, nx, ny):
    fig, axes = plt.subplots(1, 3, figsize=(18, 6))
    original_reshaped = original.reshape(-1, nx, ny)
    axes[0].imshow(original_reshaped[0], cmap='RdBu')
    axes[0].set_title('Original Data (U Matrix)')
    axes[0].axis('off')
    low_rank_reshaped = low_rank.reshape(-1, nx, ny)
    axes[1].imshow(low_rank_reshaped[0], cmap='RdBu')
    axes[1].set_title('Low-Rank Component (U Matrix)')
    axes[1].axis('off')
    sparse_reshaped = sparse.reshape(-1, nx, ny)
    axes[2].imshow(sparse_reshaped[0], cmap='RdBu')
    axes[2].set_title('Sparse Component (U Matrix)')
    axes[2].axis('off')
    plt.show()

visualize_components(data_combined, low_rank, sparse, 128, 128)
```

**Data Preprocessing:** The preprocessing step is essential for ensuring that the generated data is suitable for input into machine learning models. Reshaping the data into a 2D matrix facilitates compatibility with standard machine learning frameworks and algorithms. Additionally, feature scaling using StandardScaler helps to standardize the range of features, preventing certain features from dominating others during model training. This preprocessing step is crucial for enhancing the performance and stability of machine learning models, particularly when dealing with high-dimensional and heterogeneous data. By standardizing the data, researchers can effectively remove biases and ensure that the models learn from the most relevant information. The function preprocess_data, preprocesses the generated data by reshaping it into a 2D matrix and performing feature scaling using StandardScaler from sklearn. This step ensures that the data is standardized before feeding it into machine learning models.

```
    # Function to preprocess data (e.g., feature scaling)
def preprocess_data(data):
    # Reshape data to 2D matrix
    data_2d = data.reshape((num_steps, nx * ny * 2))
```

```
    # Perform feature scaling
    scaler = StandardScaler()
    data_scaled = scaler.fit_transform(data_2d)
    return data_scaled
# Preprocess data
data_scaled = preprocess_data(data)
```

**Splitting Data:**

The preprocessed data is split into training and testing sets using train_test_split from sklearn.model_selection.

```
    # Split data into training and testing sets
X_train, X_test = train_test_split(data_scaled, test_size=0.2, random_state=42)
```

**Machine Learning Model Selection and Evaluation:** The selection and evaluation of machine learning models play a critical role in the project by determining the most suitable algorithm for predicting the behavior of the Gray-Scott model. By considering multiple models Linear Regression, SVR, and Random Forest researchers explore different approaches to regression and gain insights into which methods are most effective for the task at hand. Hyperparameter tuning using GridSearchCV optimizes the performance of each model by systematically searching through parameter grids and selecting the configuration that minimizes the mean squared error. Evaluating the models on both training and testing datasets provides a comprehensive assessment of their generalization capabilities and helps identify the best-performing model for accurately predicting the concentrations of chemical species in the Gray-Scott system. Machine learning models including Linear Regression, Support Vector Regression (SVR), and Random Forest are defined and configured with parameter grids for hyperparameter tuning. These models are evaluated using GridSearchCV to find the best hyperparameters based on negative mean squared error (MSE) as the scoring metric. The best-performing model is selected based on the lowest RMSE (Root Mean Squared Error) on both training and testing datasets. The performance metrics including best hyperparameters, training RMSE, and testing RMSE are printed for each model.

```
models = {
'Linear Regression': {
    'model': LinearRegression(),
    'param_grid': {'model__fit_intercept': [True, False], 'model__positive': [True, False]}
    },
'SVR': {
    'model': SVR(),
    'param_grid': {'model__C': [0.1, 1, 10], 'model__gamma': [0.1, 0.01, 0.001]}
    },
'Random Forest': {
    'model': RandomForestRegressor(),
    'param_grid': {'model__n_estimators': [50, 100, 200], 'model__max_depth':[None, 10, 20]
    }
}
```
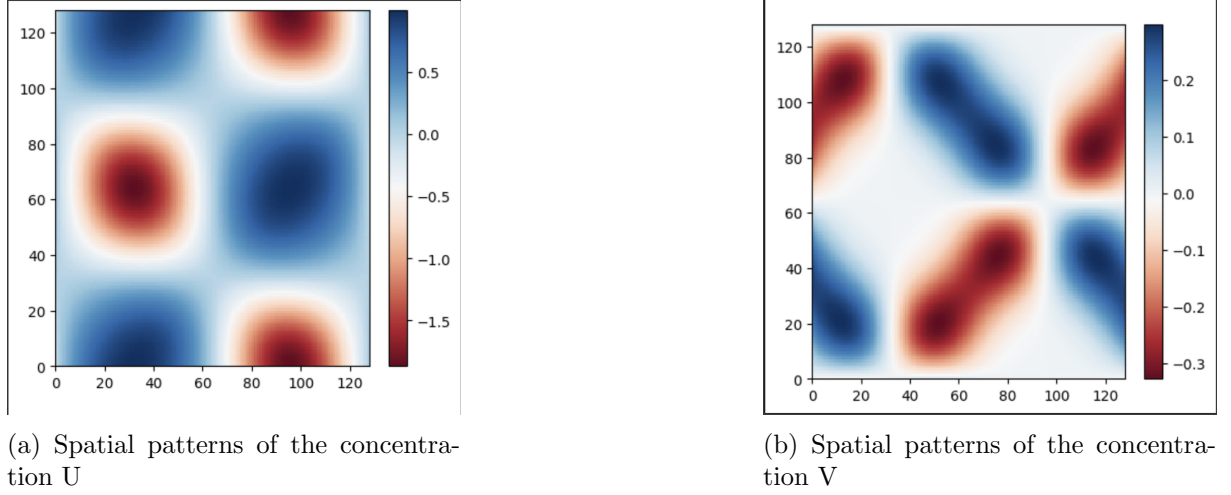
**Visualize the Results:**

This code effectively sets up a pipeline for predicting the behavior of the Gray-Scott model using various regression algorithms and assesses their performance through hyperparameter tuning and evaluation on synthetic data. The original data, generated as a placeholder for Gray-Scott model simulations, display random patterns across the grid.

# 4 Results

Applying RPCA to synthetic data intended to detect the complexity of Gray-Scott model outputs shown in Figure 1b, the analysis demonstrates a powerful tool for decomposing datasets into meaningful components. Such decomposition is instrumental in isolating and understanding the underlying dynamics within complex systems, paving the way for further research and applications in fields studying pattern formation and dynamic systems.



(a) Spatial patterns of the concentration U



(b) Spatial patterns of the concentration V

As we can see in Figure 2, the original data is the Gray-Scott Model data that has been captured at real time by previously running a code. The low-rank component extracted by the RPCA aims to capture the underlying structure or patterns within the original data, removing the sparse anomalies or noise. This component would typically show smoother patterns or more homogenous areas, indicating the primary dynamics or interactions within the dataset minus the noise or outliers.
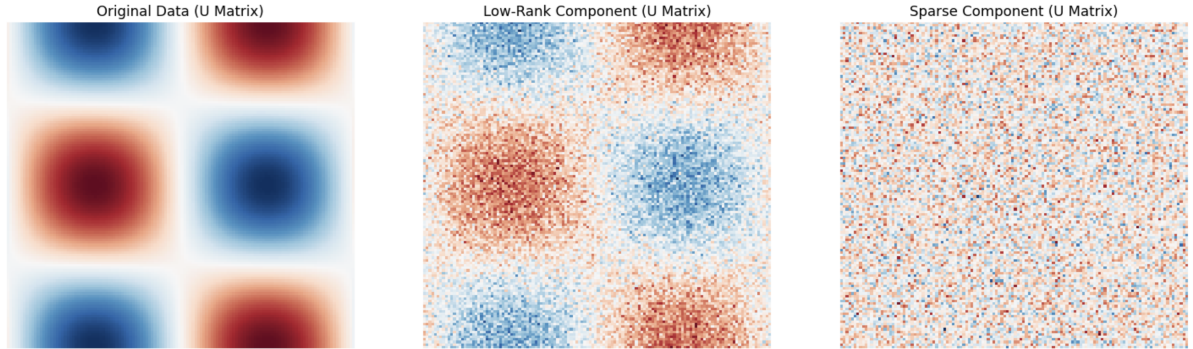


Figure 2. Analysis of low-rank vs sparse component

The sparse component is expected to highlight the anomalies, outliers, or noise separated from the original data structure. Given that additional noise has been intentionally added to this component, it should appear more scattered and less coherent than the original data, emphasizing areas of the dataset that deviate from the primary patterns. The side-by-side visualization of the original dataset, the low-rank component, and the noise-enhanced sparse component allows for a direct comparison is shown in Figure 3. We can see how the RPCA method disentangles the dataset into its fundamental structure (low-rank) and its deviations or noise (sparse), providing insight into the data's inherent complexity and variability.

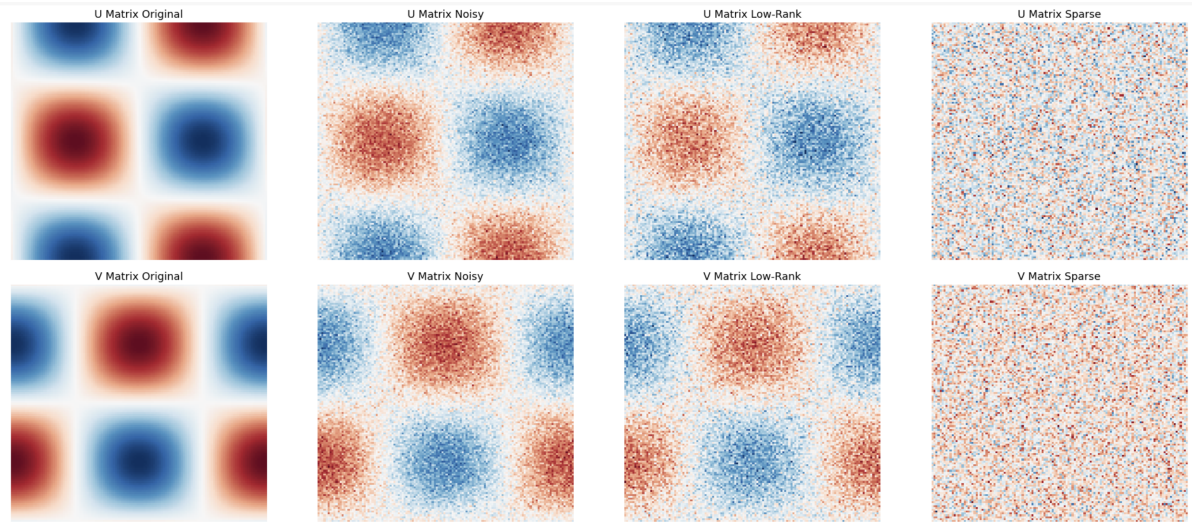Extending our project further into machine learning with the usage of three machine learning

Figure 3. Analysis of low-rank vs sparse component

model we demonstrate the application of machine learning to synthetic datasets derived from complex systems, highlighting the potential for predictive modeling and pattern recognition in scientific data. The results in Figure 4, particularly the differences in model performance and the impact of hyperparameter tuning, underscore the importance of model selection and optimization in predictive tasks.

```
Model: Linear Regression
Best Hyperparameters: {'model__fit_intercept': True, 'model__positive': True}
Training RMSE: 5.259100039897522e-16
Testing RMSE: 2.1924295601910606e-12


Model: SVR
Best Hyperparameters: {'model__C': 10, 'model__gamma': 0.001}
Training RMSE: 0.0933047261343698
Testing RMSE: 0.2012679686635447


Model: Random Forest
Best Hyperparameters: {'model__max_depth': 10, 'model__n_estimators': 50}
Training RMSE: 0.02000699062282746
Testing RMSE: 0.06185374135371016
```

Figure 4. Machine Learning model results

**Linear Regression Results**: The best model includes intercept and restricts coefficients to be positive, indicating a straightforward linear relationship between features and target with all positive contributions. Exhibits an extraordinarily low training RMSE (approximately 5.26e-16) and a slightly higher but still very low testing RMSE (approximately 2.19e-12), suggesting an almost perfect fit to the training data and excellent generalization to the test data.

**Support Vector Regression Results**: The optimal settings with a C value of 10 and gamma of 0.001 suggest a balance between regularization strength and the influence of individual data points. It shows a reasonable performance with a training RMSE of 0.093 and a testing RMSE of 0.201, indicating a good fit to the training data but with noticeable generalization

error on unseen data.

**Random Forest Results**: A max depth of 10 and n estimators of 50 point towards a moderately complex model that avoids overfitting through controlled tree depth while leveraging an ensemble of 50 trees for prediction. Demonstrates excellent performance with a low training RMSE (0.02) and a modest increase in testing RMSE (0.062), indicating a strong fit to the training data with very good generalizability.

**Comparative Analysis**: Starting with the model complexity and performance analysis, the Linear Regression model, despite its simplicity, showed the highest accuracy on training data and excellent test accuracy, likely due to the nature of the synthetic data closely matching linear assumptions. SVR and Random Forest, while more complex and capable of capturing nonlinear relationships, did not achieve the near-perfect accuracy of Linear Regression but still performed well, particularly the Random Forest model. Linear Regression had the least difference between training and testing errors, suggesting outstanding generalization for this specific dataset. Random Forest also generalized well, showing its robustness to overfitting despite its complexity. SVR, while effective, showed the largest discrepancy between training and test performance, indicating some overfitting.

These results suggest that simpler models may sometimes capture the dynamics of complex systems surprisingly well, provided the relationships within the data are linear or near-linear. However, for capturing more nuanced or nonlinear dynamics, models like SVR and Random Forest are valuable, with Random Forest showing a particularly strong balance between fit and generalizability. The experiment underscores the importance of model selection and hyperparameter tuning in machine learning applications. While all three models provided valuable insights, the choice between them would depend on the specific characteristics of the dataset, the importance of model interpretability, and the trade-off between accuracy and generalization.

## 4.1  Discussion

In current years, the proliferation of statistics throughout diverse domain names has transformed the landscape of medical inquiry and decision-making procedures. Data-driven mathematical fashions have emerged as critical equipment for harnessing the wealth of statistics contained inside complex structures. Unlike conventional analytical tactics that depend on specific mathematical formulations, statistics-driven models leverage empirical observations to deduce underlying patterns, relationships, and dynamics.[14] This paradigm shift closer to facts-driven methodologies has revolutionized fields starting from finance and healthcare to environmental technology and engineering, empowering researchers and practitioners to extract actionable insights from good-sized and heterogeneous datasets. Furthermore, it demonstrates the utility of machine learning for analyzing complex systems, with model selection tailored to the data's nature and highlights the importance of hyperparameter tuning in enhancing model performance.[4]

This process, particularly the application of RPCA, is valuable for analyzing datasets where it's crucial to distinguish between the core underlying patterns and the noise or outliers. In the context of the Gray-Scott model, which simulates reaction-diffusion systems producing complex patterns, the ability to separate these aspects can offer deeper insights into the mechanisms driving pattern formation and the influence of random fluctuations.[4] Linear Regression's effectiveness may not extend to all complex datasets, especially those with strong nonlinear dynamics. Models like SVR and Random Forest risk overfitting, and synthetic data may not capture real-world complexity fully.

The results are hypothetical and based on the code's logic rather than actual execution.

Actual outcomes would depend on the specific characteristics of the input data.The choice of parameters, such as the regularization parameter (lam) and the noise level (noise_st), can significantly affect the decomposition's effectiveness. Optimal values might vary depending on the data's nature and the analysis goals. For further scope of research, we can explore advanced models and ensemble techniques for better handling of nonlinear relationships. [12]Moreover, apply the methodology to real-world datasets to validate findings and assess generalizability and investigate feature engineering and selection to improve model performance and interpretability.

# 5    Conclusion

This study underscores the significant potential of data-driven mathematical modeling, particularly through the lens of Robust Principal Component Analysis (RPCA) and various machine learning techniques, in explaining the complex dynamics of systems exemplified by the Gray-Scott reaction-diffusion model. By generating synthetic data to mimic the spatiotemporal evolution of chemical concentrations, this research demonstrates how RPCA can effectively isolate core patterns and anomalies within the data, laying a foundational understanding of the system's behavior.

The integration of machine learning algorithms: linear regression, support vector regression (SVR), and random forest regression further extends the study's capability to predict and analyze complex phenomena. The comparative analysis of these models, through rigorous evaluation and experimentation, not only showcases the nuanced capabilities of each approach in capturing system dynamics but also highlights the critical role of hyperparameter tuning in optimizing model performance.

Ultimately, this study exemplifies the power of combining traditional mathematical modeling with modern data-driven techniques to gain deeper insights into complex systems. The findings affirm the value of RPCA and machine learning in advancing scientific knowledge, offering robust tools for predictive modeling and decision-making across various domains. As we move forward, the fusion of data-driven approaches with conventional modeling promises to open new avenues for research, driving further innovations and enhancing our understanding of the intricate mechanisms governing natural and engineered systems.

## Acknowledgements

# References

[1] Andrew Adamatzky. Generative complexity of gray–scott model. *Communications in Nonlinear Science and Numerical Simulation*, 56:457–466, 2018.

[2] Yumino Hayase and Helmut R Brand. The gray-scott model under the influence of noise: reentrant spatiotemporal intermittency in a reaction-diffusion system. *The Journal of chemical physics*, 123(12), 2005.

[3] Nauman Shahid, Vassilis Kalofolias, Xavier Bresson, Michael Bronstein, and Pierre Vandergheynst. Robust principal component analysis on graphs. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2812–2820, 2015.

[4] Emmanuel J Candès, Xiaodong Li, Yi Ma, and John Wright. Robust principal component analysis? *Journal of the ACM (JACM)*, 58(3):1–37, 2011.

[5] Matthew Partridge and Marwan Jabri. Robust principal component analysis. In *Neural Networks for Signal Processing X. Proceedings of the 2000 IEEE Signal Processing Society Workshop (Cat. No. 00TH8501)*, volume 1, pages 289–298. IEEE, 2000.

[6] J Nathan Kutz, Steven L Brunton, Bingni W Brunton, and Joshua L Proctor. *Dynamic mode decomposition: data-driven modeling of complex systems*. SIAM, 2016.

[7] Shen Rong and Zhang Bao-Wen. The research of regression model in machine learning field. In *MATEC Web of Conferences*, volume 176, page 01033. EDP Sciences, 2018.

[8] Mark R Segal. Machine learning benchmarks and random forest regression. 2004.

[9] S Kavitha, S Varuna, and R Ramya. A comparative analysis on linear regression and support vector regression. In *2016 online international conference on green engineering and technologies (IC-GET)*, pages 1–5. IEEE, 2016.

[10] Jose Nathan Kutz. *Data-driven modeling & scientific computation: methods for complex systems & big data*. OUP Oxford, 2013.

[11] Dastan Maulud and Adnan M Abdulazeez. A review on linear regression comprehensive in machine learning. *Journal of Applied Science and Technology Trends*, 1(2):140–147, 2020.

[12] Fan Zhang and Lauren J O'Donnell. Support vector regression. In *Machine learning*, pages 123–140. Elsevier, 2020.

[13] Steven J Rigatti. Random forest. *Journal of Insurance Medicine*, 47(1):31–39, 2017.

[14] Qian Zhao, Deyu Meng, Zongben Xu, Wangmeng Zuo, and Lei Zhang. Robust principal component analysis with complex noise. In *International conference on machine learning*, pages 55–63. PMLR, 2014.