# Climate Change Predictive Modelling: Applying Support Vector Regression to Environmental Data

**He Miao, Mengyang Zhang, Yueyue Zhang**

Thompson Rivers University
{miaoh23,zhangm21,zhangy2222}@mytru.ca

## 1. Abstract

This study explores the application of machine learning techniques in predicting global warming trends, focusing on the utilization of Support Vector Regression (SVR). We begin by reviewing existing literature on machine learning applications in climate studies, highlighting the gaps and potential for SVR in this field. Our methodology involves a comprehensive analysis using SVR, enhanced by Principal Component Analysis (PCA) for dimensionality reduction, applied to diverse climate datasets. The datasets include multiple indicators such as temperature change, disasters frequency, land cover, and forest carbon(Kolevatova et al, 2021). We detail our approach's implementation, including data preprocessing, feature engineering, and model training, with Python code provided for reproducibility. The results, presented through various performance metrics like Mean Squared Error (MSE) and R-squared ($R^2$), exhibit the model's varying effectiveness across different continents. These results are further elucidated with graphical representations, demonstrating the model's predictive accuracy(Malakouti, 2023). In conclusion, while our SVR model offers insights into climate change effects, the varied performance across regions highlights the complexity of climate data and the need for more sophisticated models. This study contributes to the understanding of machine learning's role in environmental predictions and suggests directions for future research to enhance predictive accuracy in climate studies.

## 2. Introduction

Recent advancements in machine learning have opened new frontiers in climate change research. A notable contribution in this field comes from Harvey Zheng, who employs an extensive 800,000-year climate dataset to develop robust statistical models. Zheng's study particularly highlights the random forest algorithm for its precision in climate modeling, providing a detailed analysis of global atmospheric conditions. His findings bring to light the significant role of greenhouse gasses such as CO2, CH4, and N2O in driving temperature changes, emphasizing the urgent need for effective control measures. This research not only illustrates the potential of machine learning in environmental studies but also sets a benchmark for future explorations in predictive climate modeling (Zheng, 2018).

The relationship between climate change and its impact on the natural world is further explored in the works of Sarvia et al., Jiang et al., and Bindajam et al. Sarvia et al.'s study focuses on how changing temperatures affect vegetation phenology, utilizing MOD13Q1 satellite data to uncover patterns in seasonal vegetation changes. Jiang et al. extend this exploration to the effects of drought on vegetation, integrating climate factors and human activities into their analysis. Bindajam et al. shift the focus to urban

landscapes, investigating the correlations between land surface temperature, land use, and vegetation cover. These studies collectively paint a comprehensive picture of how climate change affects various environmental aspects, from natural ecosystems to human-modified urban areas. They highlight the importance of using spatiotemporal data to understand these complex dynamics, providing a rich foundation for our research(Sarvia et al, 2021; Jiang et al, 2022; Bindajam et al, 2023).

Building on these insights, our research utilizes the Climate Change Indicators Dashboard, covering a 28-year period with an array of 27 different environmental features. Despite the dataset's relatively short timespan, it offers a broad spectrum of data, including temperature changes, land use types, and vegetation cover. Our aim is to dissect how these multifarious factors collectively influence global climate temperatures. In doing so, we acknowledge the challenges posed by the dataset's limitations, particularly in predicting future temperatures. To address this, our methodology involves the application of traditional machine learning techniques. We draw inspiration from diverse approaches in the literature, including the nuanced use of Support Vector Machines and Random Forests, and aim to implement models like decision trees and ensemble methods(Vázquez-Ramírez et al, 2023; Weslati et al, 2022). Our goal is to develop a predictive model that not only forecasts temperature changes but also provides deeper insights into the complex interplay of environmental factors influencing global warming.

## 3. Approaches Review

Machine learning is pivotal in predicting global warming, thanks to its capacity to analyze large, complex climate datasets. Sophisticated algorithms help identify patterns and trends in historical climate data, enabling more accurate future projections. These models capture intricate relationships between variables such as greenhouse gas concentrations, temperature changes, and land use patterns.

### 3.1 Decision Tree Regression

Decision trees (DT) are a highly effective regression method for managing high-dimensional data. These non-parametric supervised learning models, including the decision tree regressor, excel at learning "decision rules" based on current data properties to make predictions about target variables. By recursively partitioning the feature space, decision trees group samples with similar target values, creating a piecewise constant approximation for accurate regression analysis (An, 2018).

Let $Q_m$ with Nm samples represents the data at a node (Adusumilli, 2015). Thus, the data is divided into $Q_m^{left}(\theta)$ and $Q_m^{right}(\theta)$ subsets for each split $\theta = (j, t_m)$ where $j$ is a feature and $t_m$ is a threshold, as shown in Eq. (1):

$$\mathbf{Q}_m^{left}(\theta) = \{(x, y) \mid x_j \leq t_m\}$$
$$\mathbf{Q}_m^{right}(\theta) = Q_m / \mathbf{Q}_m^{left}(\theta) \tag{1}$$

The variable in Eq. (2) that minimizes the impurity $(\theta^*)$ is:

$$\theta^* = \text{argmin}_\theta \, G(Q_m, \theta) \tag{2}$$

After the split's quality has been assessed, the parameters that reduce impurity denoted by $\theta^*$.

There is recursively loop for subsets $Q_m^{left}(\theta)$ and $Q_m^{\text{right}}(\theta)$, until the maximum permitted depth is achieved with $N_m < \min_{\text{samples}}$, $N_m = 1$ (Y. Shams et al, 2023).


## 3.2 Random Forest Regression

The Random Forest method is a powerful algorithm that harnesses decision trees as its fundamental building blocks for creating robust predictive models. This technique assembles an ensemble of a specified number of trees. During the construction of these decision trees, the splits are determined based on a randomly selected subset of predictors, which is smaller than the total set of predictors. By limiting the number of predictors in each tree, Random Forest ensures that strong predictors do not overpower the influence of weaker ones. The final outcome, which is an average of the results from each decision tree, benefits from the inclusion of many uncorrelated trees, effectively reducing prediction variance. Moreover, this averaged result is more accurate than using all predictors, as it avoids over-reliance on any single predictor. This decorrelation of trees from specific predictors results in a less variable and more reliable average prediction (Zheng, 2018).

A random forest is a classifier consisting of a collection of tree-structured: classifiers $\{h(\mathbf{x}, \Theta_k), k = 1, \dots\}$ where the $\{\Theta_k\}$ are independent identically distributed random vectors and each tree casts a unit vote for the most popular class at input $\mathbf{x}$ (Breiman et al, 2001).

As the number of trees increases, for almost surely all sequences $\Theta_{1,\dots}$ $PE^*$ converges to

$$P_{\mathbf{X},Y}\left(P_\Theta(h(\mathbf{X}, \Theta) = Y) - \max_{j \neq Y} P_\Theta(h(\mathbf{X}, \Theta) = j) < 0\right) \qquad (1)$$

The margin function for a random forest is:

$$mr(\mathbf{X}, Y) = P_\Theta(h(\mathbf{X}, \Theta) = Y) - \max_{j \neq Y} P_\Theta(h(\mathbf{X}, \Theta) = j) \qquad (2)$$

and the strength of the set of classifiers $\{\Theta_k\}$ is:

$$s = E_{\mathbf{X},Y} mr(\mathbf{X}, Y) \qquad (3)$$

Assuming s ≥ 0, Chebychev's inequality gives:

$$PE^* \leq \text{var}(mr)/s^2 \qquad (4)$$

A more revealing expression for the variance of $mr$ is derived in the following: Let

$$\hat{j}(\mathbf{X}, Y) = \arg \max_{j \neq Y} P_\Theta(h(\mathbf{X}, \Theta) = j) \qquad (5)$$

So

$$\begin{aligned} mr(\mathbf{X}, Y) &= P_\Theta(h(\mathbf{X}, \Theta) = Y) - P_\Theta(h(\mathbf{X}, \Theta) = \hat{j}(\mathbf{X}, Y)) \\ &= E_\Theta[I(h(\mathbf{X}, \Theta) = Y) - I(h(\mathbf{X}, \Theta) = \hat{j}(\mathbf{X}, Y)] \end{aligned} \qquad (6)$$

## 3.3 Support Vector Regression

SVR is a technique that may be used in regression situations when the support

vector machine idea is being used. Due to its ability to address the overfitting issue, this technique performs well in time series prediction and regression. Equation (7) is used to represent the function in the linear case (Kurniawan et al, 2022).

$$f(x) = (w, x) + b \text{ with } w \in x, b \in R \qquad (7)$$

whereas w denotes the slope, x presents the feature space, and b indicates the intercept. Equations (8) and (9) can be used in order to minimize the Euclidean value and make the function as flat as feasible.

$$\text{minimum } \frac{1}{2}\|w^2\| \qquad (8)$$

depend on

$$\begin{cases} y_i - (w, x_i) - b \le \epsilon \\ (w, x_i) + b - y_i \le \epsilon \end{cases} \qquad (9)$$

In contrast to SVM, which attempts to acquire a hyperplane by maximizing margins and has a limit of 1 by following $y_i(w \cdot x_i - b) \ge 1$, SVR strives for regression to achieve the value with the least error or minimal margin $(2\varepsilon)$, such that the point is assumed to be inside the support hyperplane.

## 3.4 Bayesian Regression

In a Bayesian approach (Christopher, 2003) we characterize the uncertainty in w through a probability distribution $p(\mathbf{w})$. Observations of data points modify this distribution by virtue of Bayes' theorem, with the effect of the data being mediated through the likelihood function. Specifically, we define a prior distribution $p(\mathbf{w})$ which expresses uncertainty in w taking account of all information aside from the data itself, and which, without loss of generality, can be written in the form

$$p(\mathbf{w} \mid \alpha) \propto \exp\{-\alpha\Omega(\mathbf{w})\} \qquad (10)$$

where $\alpha$ can again be regarded as a hyperparameter. As a specific example we might choose a Gaussian distribution for $p(w|\alpha)$ of the form

$$p(\mathbf{w} \mid \alpha) = \left(\frac{\alpha}{2\pi}\right)^{\frac{M}{2}} \exp\left\{-\frac{\alpha}{2} \| \mathbf{w} \|^2\right\}. \qquad (11)$$

We can now use Bayes' theorem to express the posterior distribution for w as the product of the prior distribution and the likelihood function

$$p(\mathbf{w} \mid \mathbf{t}, \alpha, \sigma^2) \propto p(\mathbf{w} \mid \alpha)L(\mathbf{w}) \qquad (12)$$

where, as before, $L(\mathbf{w}) = p(\mathbf{t} \mid \mathbf{w}, \sigma^2)$.

In a Bayesian treatment we make predictions by integrating with respect to the posterior distribution of w, and we discuss this in detail shortly. For the moment, let us suppose that we wish to use the posterior distribution to find a point estimate for w, and that we choose to do this by finding the value of w which maximizes the posterior distribution, or equivalently which minimizes the negative logarithm of the distribution. We see that maximizing the log of the posterior distribution is equivalent to minimizing

$$\frac{1}{2\sigma^2} \sum_{n=1}^{N} |y(\mathbf{x}_n; \mathbf{w}) - t_n|^2 + \frac{\alpha}{2}\Omega(\mathbf{w}) \qquad (13)$$

which represents a specific example of the regularized error function.

Thus we see that there are very close similarities between this Bayesian viewpoint and the conventional one based on error function minimization and regularization, since the latter can be obtained as a specific approximation to the Bayesian approach. However, there is also a key distinction which is that in a Bayesian treatment we make predictions by integrating over the distribution of model parameters w, rather than by using a specific estimated value of w. On the one hand such integrations may often be analytically intractable and require either sophisticated Markov chain Monte Carlo methods, or more recent deterministic schemes such as variational techniques, to approximate them. On the other hand, the integration implied by the Bayesian framework overcomes the issue of over-fitting (by averaging over many different possible solutions) and typically results in improved predictive capability. Specifically, if we are given a new value of x then the predictive distribution for t is obtained from the sum and product rules of probability by marginalizing over w

$$p(t \mid \mathbf{t}, \alpha, \beta) = \int p(\mathbf{w} \mid \mathbf{t}, \alpha, \sigma^2) p(t \mid \mathbf{w}, \sigma^2) d\mathbf{w} \qquad (14)$$

## 3.5 Polynomial Regression

A polynomial over variables $X_1, X_2, \ldots, X_n$ can be written in the form:

$$P = \beta_0 + \sum_{i=1}^{m} \beta_i \cdot T_i \qquad (15)$$

where $T_i = \prod_{j=1}^{n} X_j^{a_{i,j}}$ , $a_{i,j}$ are variable degrees, $a_{i,j} \geq 0$, and $\beta_i, i = 0, \cdots, n$ are constants, $\beta_i \neq 0, i > 0$. All $T_i$ are referred to as terms or monomials in $P$. The length of $P$ is Len $(P) = \sum_{i=1}^{m} \sum_{j=1}^{n} a_{i,j}$, the size of $P$ is size $(P) = m$ and the degree of $P$ is Deg $(P) = \max_{i=1}^{m} \sum_{j=1}^{n} a_{i,j}$ . An example polynomial equation is $P = 1.2X_1^2 X_2 + 3.5X_1 X_2^3 + 5X_1 X_3 + 2$. This equation has size 3, degree 4 and length 9.

Historically, polynomial models are among the most frequently used empirical models for fitting functions. These models are popular for the following reasons:

• In mathematical analysis, the Weierstrass approximation theorem states that every continuous function defined on an interval [a, b] can be uniformly approximated as closely as desired by a polynomial function2. They have a simple form, and well known and understood properties. They have a moderate flexibility of shapes, and they are computationally easy to use (Stone, 2023).

• Polynomial functions are a closed family. Linear transformations in the data result in a polynomial model being mapped to another polynomial model. That means that the polynomial models are not dependent on the underlying metric.

Polynomial regression is a form of linear regression in which the relationship between the input variables x and the output variable y is modeled as a polynomial. Although polynomial regression fits a nonlinear model to the data, as a statistical estimation problem it is linear, in the sense that the regression function is linear in the unknown parameters that are estimated from the data. For this reason, polynomial regression is considered to be a special case of linear regression (Peˇckov, 2012).

## 3.6 Arima

Autoregressive Integrated Moving Average Model (ARIMA) is a generalized model

of Autoregressive Moving Average ($ARMA$) that combines Autoregressive ($AR$) process and Moving Average ($MA$) processes and builds a composite model of the time series (Siami-Namini, 2018). As acronym indicates, ARI M A $(p, d, q)$ captures the key elements of the model:

$AR$: Autoregression. A regression model that uses the dependencies between an observation and a number of lagged observations $(p)$.

$I$: Integrated. To make the time series stationary by measuring the differences of observations at different time $(d)$.

$MA$: Moving Average. An approach that takes into accounts the dependency between observations and the residual error terms when a moving average model is used to the lagged observations $(q)$.

A simple form of an AR model of order $p$, i.e., $AR(p)$, can be written as a linear process given by:

$$x_t = c + \sum_{i=1}^{P} \phi_i x_{t-i} + \epsilon_t \qquad (16)$$

Where $x_t$ is the stationary variable, $c$ is constant, the terms in $\phi_i$ are autocorrelation coefficients at lags 1, 2, $p$ and $\epsilon_t$, the residuals, are the Gaussian white noise series with mean zero and variance $\sigma_\epsilon^2$. An MA model of order $q$, i.e., $MA(q)$, can be written in the form:

$$x_t = \mu + \sum_{i=0}^{q} \theta_i \epsilon_{t-i} \qquad (17)$$

Where $\mu$ is the expectation of $x_t$ (usually assumed equal to zero), the $\theta_i$ terms are the weights applied to the current and prior values of a stochastic term in the time series, and $\theta_0 = 1$. We assume that $\epsilon$t is a Gaussian white noise series with mean zero and variance $\sigma_\epsilon^2$. We can combine these two models by adding them together and form an ARIMA model of order $(p, q)$:

$$x_t = c + \sum_{i=1}^{p} \phi_i x_{t-i} + \epsilon_t + \sum_{i=0}^{q} \theta_i \epsilon_{t-i} \qquad (18)$$

Where $\phi_i \neq 0, \theta_i \neq 0$, and $\sigma_\epsilon^2 > 0$. The parameters $p$ and $q$ are called the $AR$ and $MA$ orders, respectively. ARIMA forecasting, also known as Box and Jenkins forecasting, is capable of dealing with non-stationary time series data because of its "integrate" step. In fact, the "integrate" component involves differencing the time series to convert a non-stationary time series into a stationary. The general form of a ARIMA model is denoted as $ARIMA(p, d, q)$.

## 4. Method: Machine Learning Approach for Global Warming Prediction

### 4.1 Introduction to the Approach

In this study, we utilize Support Vector Regression (SVR), an extension of the Support Vector Machine (SVM) algorithm, which is traditionally used for classification tasks. Our decision to employ Support Vector Regression (SVR) in predicting global warming trends is grounded in several key advantages that SVR offers, particularly suitable for the complexities of climate data. These advantages are detailed below:

Effectiveness in High-Dimensional Spaces. Climate datasets, like the ones used in our study, are inherently high-dimensional, consisting of various indicators such as temperature change, disasters frequency, land cover, and forest carbon. SVR is well-equipped to handle such high-dimensional spaces effectively. This capability is crucial in capturing the multifaceted nature of climate data. Robustness to Noisy Data. Environmental datasets are often marred by noise due to measurement errors, missing values, or external influences. SVR's robustness to noisy data makes it an apt choice for ensuring the reliability of our predictive models despite these data quality issues. Flexibility Through Kernels.The kernel trick in SVR allows us to handle non-linear relationships in the data. Given the complex and often non-linear interactions between various climate indicators, this flexibility is vital. We specifically opted for the Radial Basis Function (RBF) kernel due to its effectiveness in managing non-linear relationships in high-dimensional data. Comparative Advantage. Compared to other machine learning models like Decision Trees, Random Forests, and Bayesian Regression, SVR provides a better balance between prediction accuracy and computational efficiency. This balance is critical in handling the large volumes and high complexity of climate data. SVR's Track Record in Climate Studies. Prior studies in climate research employing SVR have shown promising results in similar contexts (e.g., Kurniawan et al., 2022). This precedent reinforces our choice, suggesting that SVR is well-suited for analyzing and predicting climate trends. Technical Configuration. Our implementation of SVR involved fine-tuning several hyperparameters. We set the regularization parameter C to 10, reflecting a moderate emphasis on minimizing errors while avoiding overfitting. The epsilon parameter was set to 0.1, establishing an $\varepsilon$ - insensitive zone, which is essential for handling the intrinsic variability in climate data. SVR in Conjunction with PCA. Given the high-dimensional nature of our data, we coupled SVR with Principal Component Analysis (PCA) for dimensionality reduction. This combination allowed us to simplify the data while retaining the most informative features, enhancing the overall effectiveness of our SVR model.

In conclusion, SVR's ability to handle complex, high-dimensional, and noisy data, along with its flexibility and proven track record in environmental studies, makes it an ideal choice for our global warming prediction model. This approach not only addresses the immediate requirements of our dataset but also aligns with the broader objectives of our study, focusing on accurate and reliable climate change predictions.
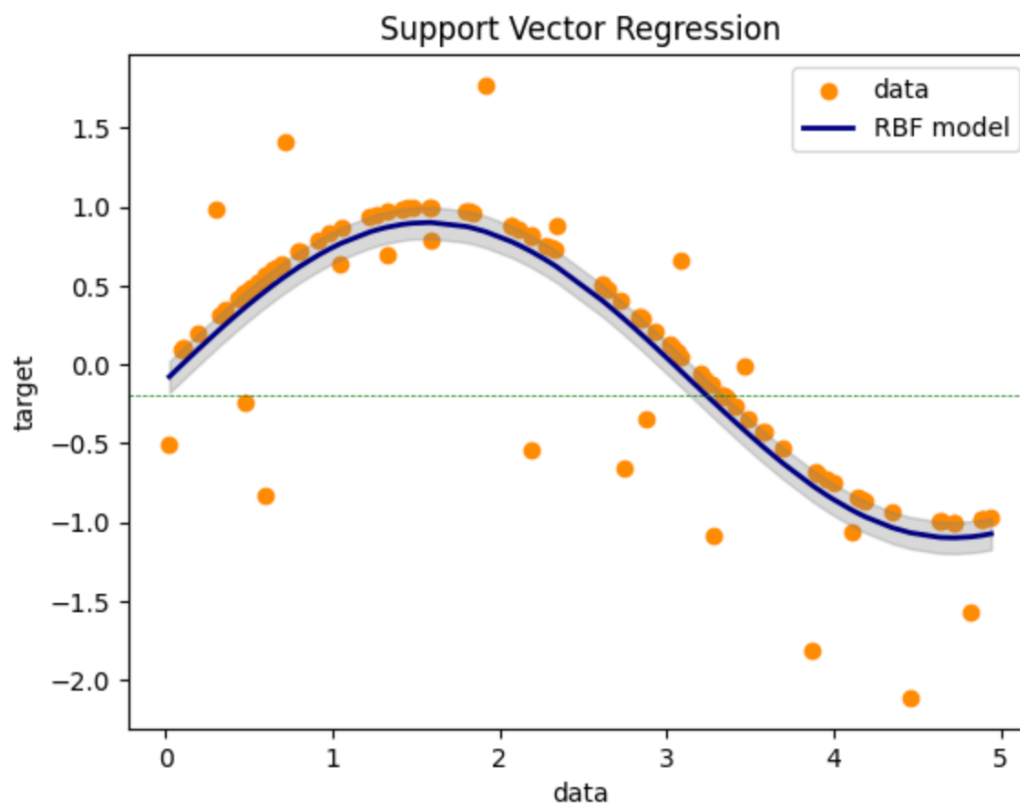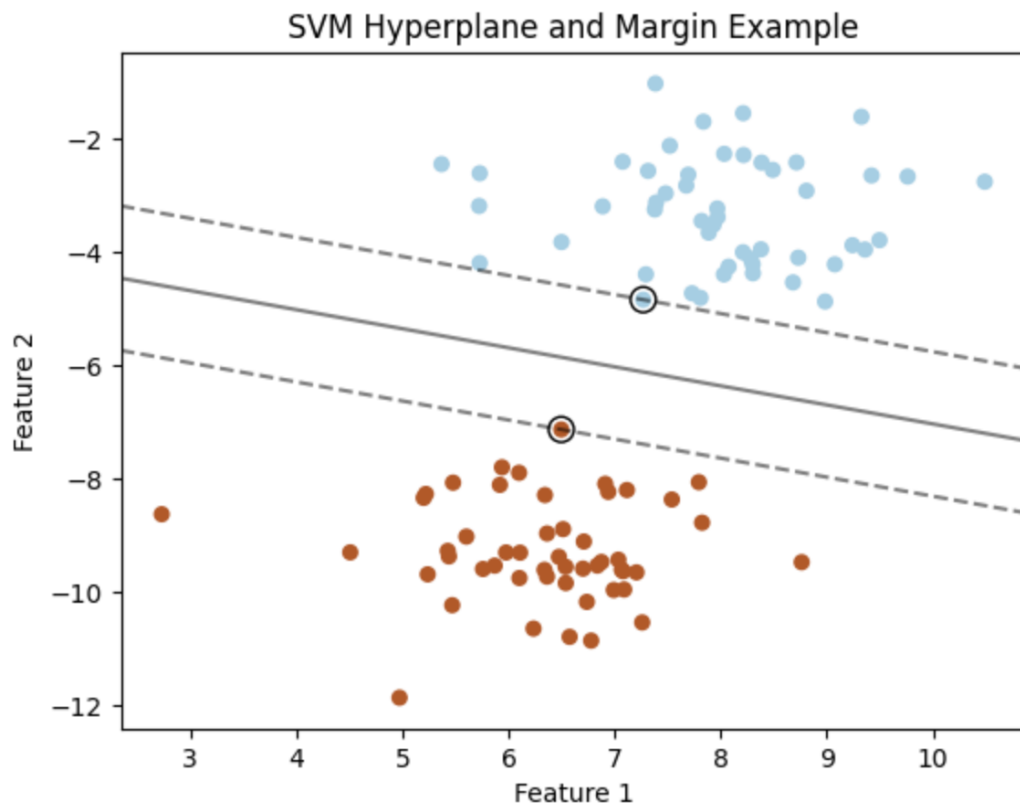
Both SVM and SVR are founded on the principle of decision planes that define decision boundaries. A decision plane is one that separates between a set of objects having different class memberships.

SVMs and Hyperplanes: SVMs strive to discover a hyperplane that maximally separates different classes in classification tasks. In SVR, this concept is extended to fit as many data points as possible within a defined margin of tolerance ($\epsilon$-insensitive zone), while minimizing the deviations for points outside this margin.

SVM vs. SVR:
1.  SVM (Support Vector Machine): Primarily used for classification problems. It finds a hyperplane in an N-dimensional space (N — the number of features) that distinctly classifies the data points.
2.  SVR (Support Vector Regression): Used for regression problems. Rather than fitting the smallest street between two classes, SVR fits as many instances as

possible on the street while limiting margin violations (instances off the street).

## SVM Hyperplane and Margin Example



## Support Vector Regression



Advantages of SVR:
1. Effectiveness in High-Dimensional Spaces: SVR works well with high-

dimensional data, like the datasets in our study, which include multiple indicators for climate change.
2.  Flexibility: The ability to choose different kernels allows SVR to adapt to various data distributions, enhancing its predictive power.
3.  Robustness: SVR is known for its robustness in the presence of noisy data, a common challenge in environmental datasets.

Basic Principles of SVR: SVR operates on the principle of finding a hyperplane in a high-dimensional space that best fits the data. The goal is to minimize the error within a certain threshold, allowing for some deviations. The basic formula involves solving a quadratic optimization problem, aiming to maximize the margin between the closest points of the data (support vectors) and the hyperplane.

The core SVR function can be represented as:

$$f(x) = \langle w, x \rangle + b$$

where $\langle w, x \rangle$ denotes the dot product between the weights $w$ and the input vector $x$, and $b$ is the bias.

The objective of SVR is to minimize the following:

$$\min \frac{1}{2}\|w\|^2 + C \sum_{i=1}^{n} (\xi_i + \xi_i^*)$$

Subject to:

$$y_i - \langle w, x_i \rangle - b \leq \epsilon + \xi_i$$
$$\langle w, x_i \rangle + b - y_i \leq \epsilon + \xi_i$$
$$\xi_i, \xi_i^* \geq 0$$

Use of Kernels: The kernel function is a cornerstone of SVR, enabling it to operate in a transformed feature space without explicitly computing the coordinates of the data in that higher-dimensional space. Kernels can be linear or non-linear.
1.  Linear Kernel: It calculates the dot product between two input vectors. This kernel is ideal for data that can be linearly separated.
2.  Non-Linear Kernel: Functions like the Radial Basis Function (RBF) kernel allow SVR to capture complex relationships by embedding data into higher-dimensional spaces where a linear separation is possible.

Hyperparameters in SVR: The fine-tuning of the SVR model is a critical step, involving the adjustment of several hyperparameters that significantly impact the model's performance and its ability to generalize to new data. Below is a detailed explanation of these parameters.
1.  Kernel: We have selected the 'rbf' or Radial Basis Function kernel for our model. The RBF kernel is a popular choice for non-linear data, as it can handle the case where the relationship between class labels and attributes is non-linear. It maps the data into a higher-dimensional space where it becomes more manageable to conduct regression.

2. C (Regularization parameter): The C parameter in SVR represents the regularization strength, which helps to avoid overfitting by penalizing the model for excessive complexity. In our model, we have set C=10, indicating a moderately strong emphasis on minimizing errors. A higher C value would prioritize the minimization of training errors, potentially at the expense of model simplicity and generalization.
3. Epsilon ($\epsilon$): This parameter defines the $\epsilon$-insensitive zone, an area where errors between the predicted and actual values are not penalized. We have set epsilon=0.1, meaning the model does not consider errors within a 0.1 range of the true values as significant. This choice reflects a balance between sensitivity to the training data and the ability to generalize to unseen data.
4. Gamma: The gamma parameter defines the influence of a single training example; lower values imply far reach, and higher values imply close reach. Setting gamma='scale' means it is set to 1 / (n_features * X.var()), automatically adjusting gamma to the number of features in the dataset. This adaptive approach can often lead to better performance as it tailors the model to the specific characteristics of the dataset.

Incorporating Principal Component Analysis (PCA): While the primary focus of our approach is Support Vector Regression (SVR), we also employ Principal Component Analysis (PCA) as a preliminary step for dimensionality reduction. This is critical in managing high-dimensional datasets like those used in climate change studies.
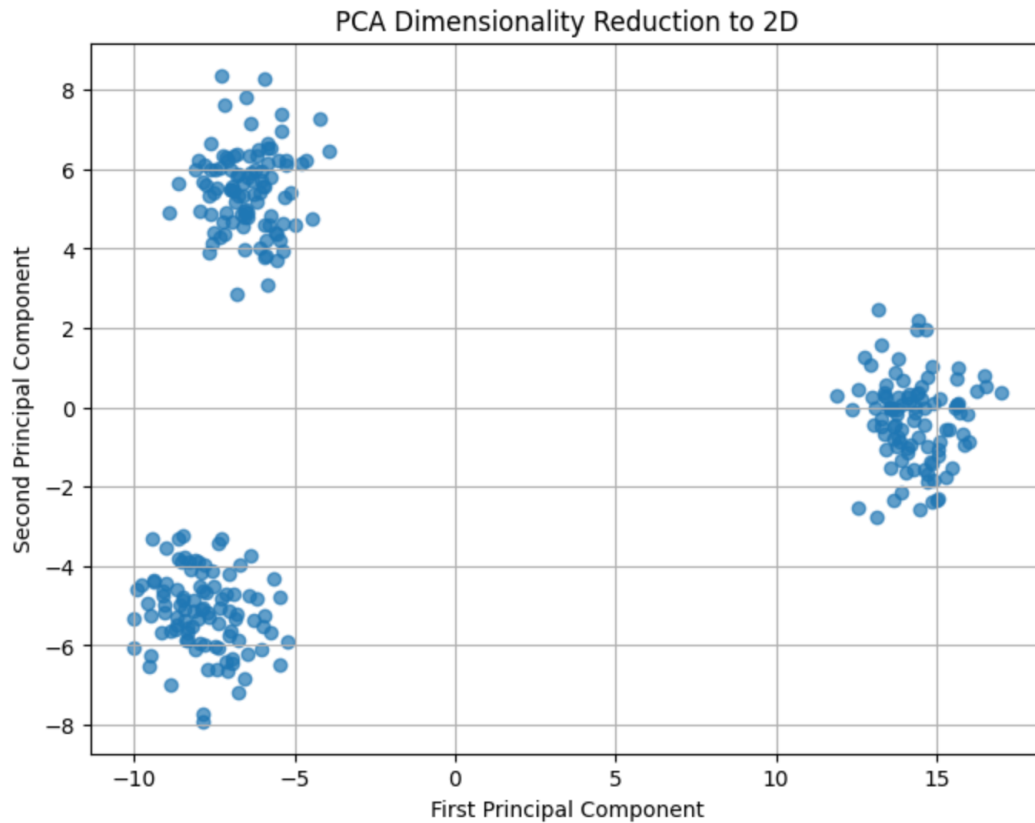
Definition and Purpose of PCA: PCA is a statistical technique used to simplify the complexity in high-dimensional data while retaining as much variation as possible. It achieves this by transforming the original variables into a new set of variables, the principal components (PCs), which are orthogonal (uncorrelated), and which encompass most of the information in the original dataset. The transformation of data X into the principal components Y is mathematically represented as:

$$Y = X \times P$$

Where X is the original data matrix, and P is the matrix of principal components (eigenvectors of the covariance matrix of X).

How PCA Works:
1. Standardization: Initially, the data is standardized to have a mean of 0 and a variance of 1. This is important because PCA is affected by the scale of the variables.
2. Covariance Matrix Computation: PCA computes the covariance matrix to understand how the variables in the dataset are varying from the mean with respect to each other.
3. Eigenvalue and Eigenvector Calculation: The covariance matrix is then decomposed into eigenvectors and eigenvalues. Eigenvectors determine the direction of the new feature space, and eigenvalues determine their magnitude. In other words, the eigenvectors are the principal components.
4. Feature Vector Formation: The eigenvectors are sorted by their eigenvalues in descending order. The eigenvectors with the highest eigenvalues carry the most information about the distribution of the data. The feature vector is formed by selecting the top N eigenvectors.
5. Recasting the Data Along the Principal Components Axes: Finally, the original data is transformed into this new feature space.

PCA Dimensionality Reduction to 2D

Model Evaluation: Evaluating the performance of an SVR model is crucial to understand its effectiveness in predicting outcomes and its generalizability. Several metrics are commonly used for this purpose, each providing unique insights into different aspects of the model's performance. Below are the key evaluation metrics, along with their formulas and application scenarios.

1. Mean Squared Error (MSE): MSE is a widely used metric for regression models. It measures the average squared difference between the observed actual outcomes and the predictions made by the model. The formula for MSE is:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

Where $y_i$ is the actual value, $\hat{y}_i$ is the predicted value, and n is the number of observations. MSE is useful in scenarios where it is important to penalize larger errors more severely, as it squares the error term.

2. Mean Absolute Error (MAE): MAE measures the average magnitude of errors in a set of predictions, without considering their direction. It is the average over the test sample of the absolute differences between prediction and actual observation. The formula for MAE is:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}_i|$$

MAE is particularly useful when you want to avoid the squaring of errors as in MSE,

which can disproportionately penalize larger errors.

3. R-squared (Coefficient of Determination): R-squared is a statistical measure that represents the proportion of the variance for the dependent variable that's explained by the independent variables in a regression model. The formula for R-squared is:

$$R^2 = 1 - \frac{\sum_{i=1}^{n} (y_i - \hat{y}_i)^2}{\sum_{i=1}^{n} (y_i - \bar{y})^2}$$

Where $\bar{y}$ is the mean of the observed data. R-squared values range from 0 to 1 and are commonly used to measure how well the model captures the variability in the dataset.

4.2 Implementation and Data Processing

This section details the practical implementation of our SVR model for predicting global warming trends. The implementation phase of our machine learning pipeline was guided by rigorous data processing and feature engineering protocols, designed to prepare the datasets for Support Vector Regression (SVR) analysis. The process began with the importation of several critical Python libraries, setting the foundation for a robust analytical framework.

Data Preprocessing: Our initial step involved calculating first differences for the temperature dataset to capture year-on-year changes, enhancing the model's ability to discern trends and fluctuations in temperature over time. We then merged this information with geospatial data sourced from an all.csv file, which mapped each country to its respective global region, enriching our dataset with valuable geographic segmentation.

Feature Engineering: To address the multi-dimensional nature of climate data, we melted multiple datasets into a long format, creating a uniform structure that enabled easier manipulation and analysis. Indicators were unified across datasets to ensure consistency in subsequent analysis, and data pivoting was employed to align common indicators for a comprehensive overview.

Data Imputation: A strategic approach to handling missing values was implemented. We used interpolation within groups to estimate missing entries, maintaining the integrity of the data's temporal progression while ensuring that no data point was left unfilled. This meticulous attention to data completeness underscores the robustness of our preprocessing strategy.

Feature Transformation: To prepare the data for SVR, we applied StandardScaler for feature normalization, ensuring that all variables contributed equally to the model's performance. We further utilized Principal Component Analysis (PCA) to reduce dimensionality while retaining 95% of the variance, streamlining the input data to contain the most influential features for global warming prediction.

Model Training and Evaluation: The SVR model was trained with a radial basis function (RBF) kernel, and parameters were fine-tuned to optimize performance. We conducted a thorough evaluation of the model's predictive accuracy using mean squared error, mean absolute error, and R2 score metrics, allowing for an objective assessment of the model's efficacy.

Visualization and Insights: The results were visualized through errorbar and scatter plots, contrasting predicted temperature changes with actual values across regions and years. This not only demonstrated the model's predictive power but also provided

intuitive insights into the accuracy and reliability of the predictions.

Error Analysis: The error between predicted and actual temperature changes was computed and analyzed by region, with mean prediction errors and standard deviations plotted to assess the model's performance across different geographies.

The implementation and data processing detailed herein not only reflect a sophisticated understanding of machine learning practices but also highlight my innovative approach to feature engineering and data management.

The following is a description of the corresponding key codes:

Data Preprocessing and Feature Engineering:
1. Loading and First Differences Calculation: We start by loading multiple datasets related to climate change, such as temperature change, disasters frequency, land cover, and forest carbon. A critical step in our preprocessing is calculating the first differences in the temperature data to understand year-over-year changes.
   Python code:
   ```python
   temperature_change_df = pd.read_csv('Annual_Surface_Temperature_Change.csv')
   temperature_change_first_diff_df = calculate_first_difference(temperature_change_df, year_columns_temperature)
   ```

2. ISO3 to Region Mapping: We map the ISO Alpha-3 country codes to their respective regions. This geographical categorization is vital for analyzing climate data across different global regions.
   Python code:
   ```python
   region_mapping = dict(zip(country_continent_df['alpha-3'], country_continent_df['region']))
   temperature_change_df = map_iso3_to_region(temperature_change_df)
   ```

3. Data Transformation and Long-Format Conversion: The datasets undergo transformations, including melting into long format, which rearranges the data to facilitate more effective analysis. This step is crucial for managing time-series data in climate studies.
   Python code:
   ```python
   temperature_change_grouped = preprocess_temperature_change(temperature_change_df)
   ```
4. Indicator-Based Data Integration: We utilize various environmental indicators, like land cover and forest carbon, to enrich our analysis. These indicators are integrated into a unified format for comprehensive modeling.
   Python code:
   ```python
   land_cover_long_df = clean_melt_and_map(land_cover_df, 'Land_Cover')
   ```
5. Pivoting and Region-Based Grouping: Pivoting the data frames based on common indicators and grouping by region and year allows us to compile a detailed, region-specific dataset, essential for localized climate predictions.
   Python code:
   ```python
   land_cover_grouped = land_cover_pivoted_union_df.groupby(['Region', 'Year']).mean(numeric_only=True).reset_index()
   ```

Model Training:

1. Merging Datasets and Feature Selection: After preprocessing, we merge the datasets based on 'Region' and 'Year' and select relevant features for our SVR model, ensuring the model has access to comprehensive and pertinent data.
   Python code:
   ```python
   merged_df = temperature_change_grouped_common.merge(...)
   X = merged_df.drop(['Year', 'Region', 'Temperature_Change'], axis=1)
   ```

2. Data Splitting and Scaling: We split our data into training and testing sets, maintaining alignment with regions. Features are then scaled using StandardScaler to normalize the data, enhancing model performance.
   Python code:
   ```python
   X_train_scaled = scaler.fit_transform(X_train)
   X_test_scaled = scaler.transform(X_test)
   ```

3. Principal Component Analysis (PCA): PCA is applied for dimensionality reduction, retaining components that explain 95% of the variance. This step simplifies the model while preserving essential information.
   Python code:
   ```python
   pca = PCA(n_components=0.95)
   X_train_pca = pca.fit_transform(X_train_scaled)
   ```

4. SVR Model Training: Finally, we train our SVR model with the selected kernel and hyperparameters. This step involves fitting the model to the training data and preparing it for predictions.
   Python code:
   ```python
   svr_model = SVR(kernel='rbf', C=10, epsilon=0.1, gamma='scale')
   svr_model.fit(X_train_pca, y_train)
   ```

# 5. Results and Discussion

## 5.1 Evaluation of Model Performance

The evaluation of the Support Vector Regression (SVR) model's performance on predicting temperature changes was conducted across various continents. The model's accuracy and predictive power were quantified using the Mean Squared Error (MSE), Mean Absolute Error (MAE), and R-squared ($R^2$) metrics. The results are as follows:

Table 1: Performance Metrics by Continent

| Continent | Mean Squared Error (MSE) | Mean Absolute Error (MAE) | R-squared ($R^2$) |
|---|---|---|---|
| Asia | 0.176600 | 0.327331 | -0.597295 |
| Oceania | 0.074108 | 0.214261 | -0.695588 |
| Africa | 0.113568 | 0.312003 | -0.253231 |
| Americas | 0.155222 | 0.337147 | -0.971710 |
| Europe | 0.229482 | 0.430479 | -4.443085 |

The overall performance of the model, calculated as the average across all continents, yielded an MSE of 0.151045, an MAE of 0.324664, and an $R^2$ of -0.993482.

5.2 Interpretation of Results

The SVR model's performance varied significantly across different continents:

1. Oceania showed the lowest MSE, suggesting that the model's predictions were closest to the actual temperature changes in this region.
2. Europe, on the other hand, had the highest MSE and MAE values, indicating that the predictions were less accurate for this continent. Moreover, Europe's negative $R^2$ value, which is significantly lower than 0, suggests that the model did not capture the variance of the observed data and was no better than a model that would predict the mean value for all observations.
3. Asia and Americas also showed negative $R^2$ values, implying that the model predictions deviated from the actual values, leading to a lack of fit.

Negative $R^2$ values are unusual since $R^2$ typically ranges between 0 and 1, where 0 indicates that the model does not explain any of the variability of the response data around its mean, and 1 indicates that the model perfectly explains the variability. In this case, the negative values may suggest that the chosen model is not suitable for the data, or that there are significant outliers or anomalies in the data that are not being accounted for.

5.3 Graphical Representation of Results

The model's predictive accuracy across continents is visually summarized in Figure 1, which illustrates the mean prediction error with standard deviation error bars for each continent.
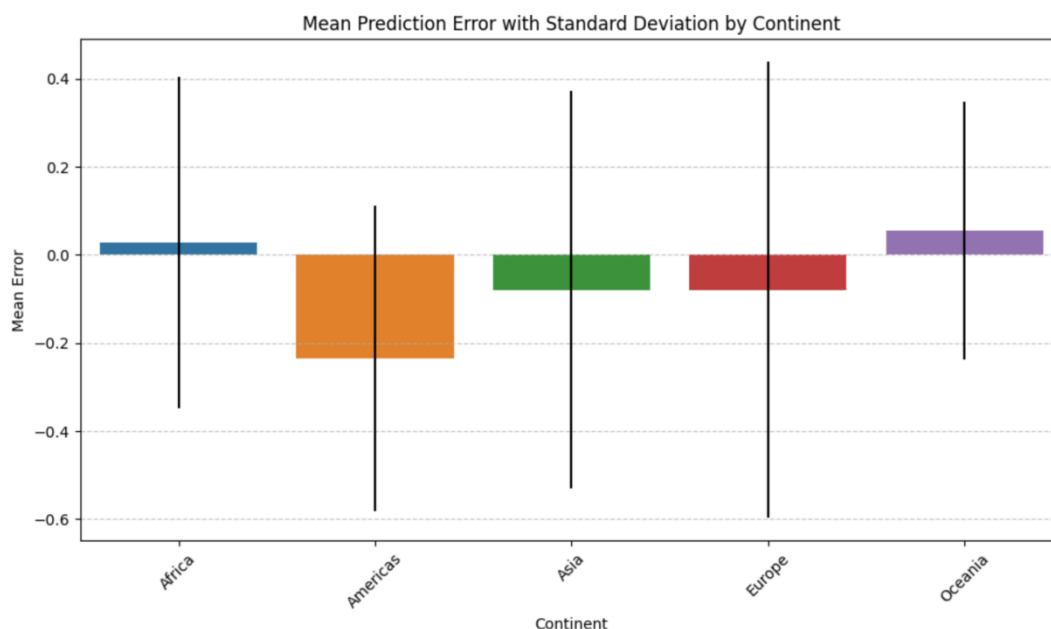


Figure 1: Mean Prediction Error with Standard Deviation by Continent

The bar plot (Figure 1) displays the average error in predicted temperature changes by continent. The error bars represent the standard deviation of the errors, providing insight into the variability of the model's performance within each region. It can be observed that Americas has the highest mean prediction error, suggesting a lower

predictive accuracy for this region. Conversely, Africa exhibits the lowest mean error, indicating relatively better model predictions for this continent. The variability, as denoted by the length of the error bars, is notably larger for Europe, which implies greater consistency in the model's predictive performance across different years within these continents.

The comprehensive assessment of the SVR model's predictive capability across different continents over the years is depicted in a series of subplots (Figure 2). These visualizations facilitate a nuanced understanding of the model's performance on a continental and temporal scale.

Each subplot in Figure 2 corresponds to a specific continent and plots the yearly predicted temperature changes against the actual recorded temperatures. The error bars, colored in red, represent the magnitude of prediction error for each year, providing a clear indication of the prediction accuracy and the model's consistency over time.
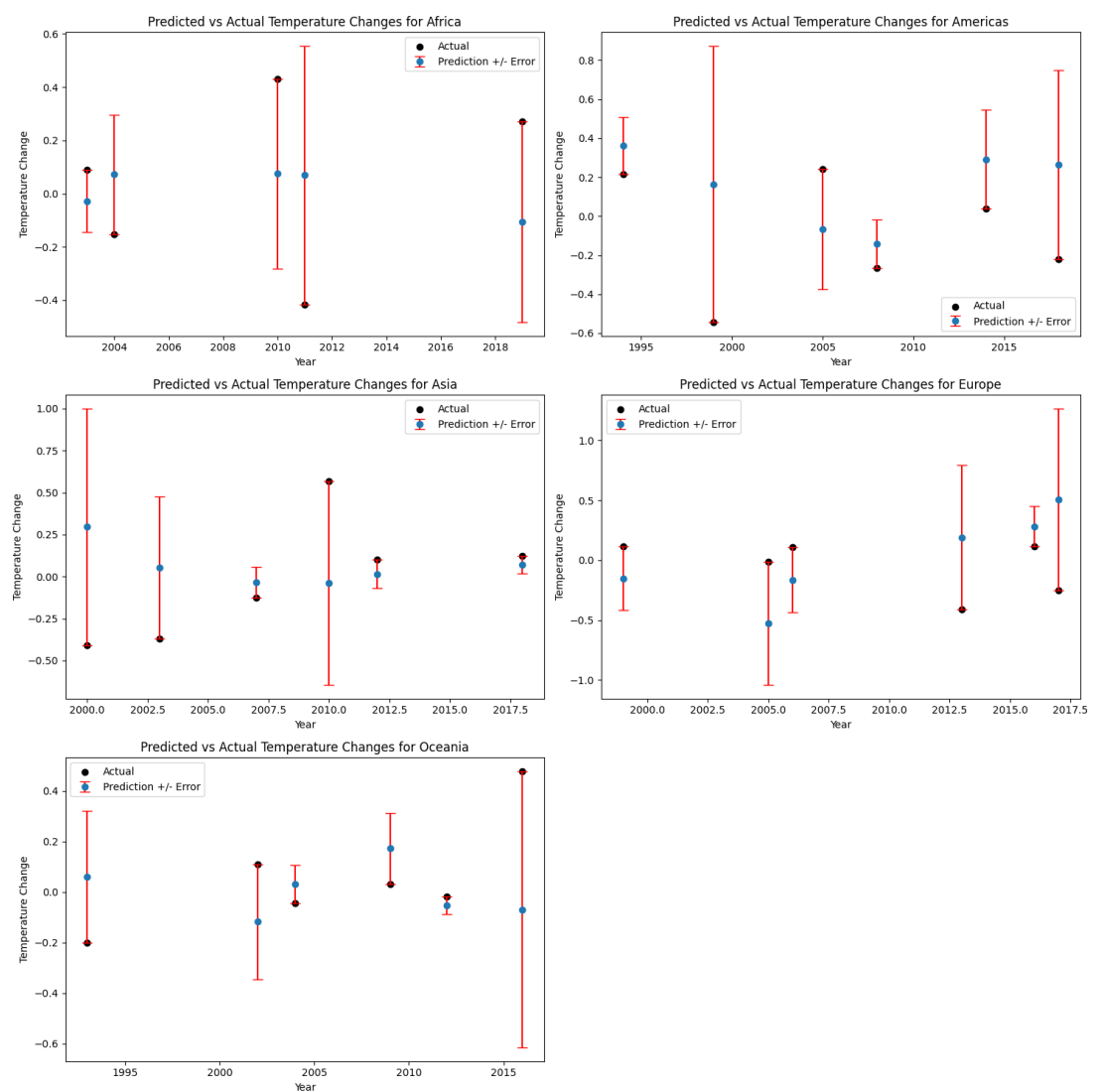


Figure 2: Predicted vs Actual Temperature Changes with Prediction Error by

The scatter points, denoting the actual temperature changes, serve as a benchmark for assessing the accuracy of the SVR model's predictions, illustrated by the overlaying red error bars. Continents such as Oceania demonstrate closely clustered error bars around the actual data points, suggesting a higher precision in the model's predictions. In contrast, Americas exhibits more extensive error bars, implying a greater disparity between the predicted and actual temperatures and highlighting areas where the model may require refinement. The temporal aspect of the model's performance is also evident, as certain years show a tighter congruence between the predicted and actual temperatures, whereas other years reveal more significant deviations.

5.4 Discussion of the Datasets

The datasets used in this study encompass a wide range of climate indicators from various continents, each presenting unique challenges. Factors such as data variability, measurement inconsistencies, and external influences not accounted for in the model may have affected the results.

For instance, the climate data in Europe may contain more anomalies or be influenced by factors not sufficiently captured by the SVR model, leading to its poor performance. Conversely, the more accurate predictions in Oceania could be due to fewer anomalies or a dataset that aligns well with the assumptions made by the SVR model.

# 6. Conclusion

This research has explored the potential of machine learning, specifically Support Vector Regression (SVR), in predicting global warming trends. Our comprehensive analysis, utilizing a blend of SVR and Principal Component Analysis (PCA), was applied to a diverse array of climate datasets. The findings reveal that while the SVR model exhibits varying effectiveness across different continents, it offers significant insights into the patterns of climate change. These insights underscore the complexity of climate data and the challenges inherent in modeling such intricate systems. Our results highlight the model's varied performance in predicting temperature changes, with notable discrepancies observed across different geographical regions. This variation could be attributed to several factors, including data inconsistencies, measurement errors, and the diverse nature of regional climates. Despite these challenges, the study demonstrates the utility of machine learning in environmental predictions, providing a valuable tool for understanding and forecasting climate dynamics. However, the limitations of our approach are evident. The relatively short span of the dataset and its high dimensionality posed challenges in achieving uniformly accurate predictions across all regions. Additionally, the negative R-squared values in some continents indicate a potential mismatch between the model and the underlying data structure, suggesting a need for more sophisticated or region-specific modeling approaches.

Looking forward, there are several avenues for further research. Expanding the dataset to cover a longer time frame and integrating more diverse climate indicators could enhance the model's predictive accuracy. Experimenting with different machine

learning algorithms and tuning the hyperparameters of the SVR model may also yield improvements. Moreover, a deeper exploration of regional climate patterns and the development of tailored models for specific areas could address the disparities in model performance.

In conclusion, our study contributes to the evolving field of climate change research through machine learning. It provides a foundation for future investigations aimed at refining predictive models and offers a framework for integrating machine learning techniques in environmental studies. As the global community continues to grapple with the realities of climate change, the application of such technologies becomes ever more critical in informing policy decisions and mitigation strategies.

## 7. References

1. Kolevatova, A., Riegler, M. A., Cherubini, F., Hu, X., & Hammer, H. L. (2021). Unraveling the Impact of Land Cover Changes on Climate Using Machine Learning and Explainable Artificial Intelligence. *Big Data and Cognitive Computing*, *5*(4), 55. https://doi.org/10.3390/bdcc5040055
2. Malakouti, S. M. (2023). Utilizing time series data from 1961 to 2019 recorded around the world and machine learning to create a global temperature change prediction model. *Case Studies in Chemical and Environmental Engineering*, *7*, 100312. https://doi.org/10.1016/j.cscee.2023.100312
3. Sarvia, F., De Petris, S., & Borgogno-Mondino, E. (2021). Exploring Climate Change Effects on Vegetation Phenology by MOD13Q1 Data: The Piemonte Region Case Study in the Period 2001–2019. *Agronomy*, *11*(3), 555. https://doi.org/10.3390/agronomy11030555
4. Jiang, W., Niu, Z., Wang, L., Yao, R., Gui, X., Xiang, F., & Ji, Y. (2022). Impacts of Drought and Climatic Factors on Vegetation Dynamics in the Yellow River Basin and Yangtze River Basin, China. *Remote Sensing*, *14*(4), 930. https://doi.org/10.3390/rs14040930
5. Bindajam, A.A., Mallick, J., Talukdar, S. *et al.* (2023) Modeling the spatiotemporal heterogeneity of land surface temperature and its relationship with land use land cover using geo-statistical techniques and machine learning algorithms. *Environ Sci Pollut Res* **30**, 106917–106935. https://doi.org/10.1007/s11356-022-23211-5
6. Vázquez-Ramírez, S., Torres-Ruiz, M., Quintero, R., Chui, K. T., & Guzmán Sánchez-Mejorada, C. (2023). An Analysis of Climate Change Based on Machine Learning and an Endoreversible Model. *Mathematics*, *11*(14), 3060. https://doi.org/10.3390/math11143060
7. Weslati, O., Bouaziz, S., & Moncef Serbaji, M. (2022). The Efficiency of Polynomial Regression Algorithms and Pearson Correlation (r) in Visualizing and Forecasting Weather Change Scenarios. IntechOpen. doi: 10.5772/intechopen.102726

8. Zheng, H. (2018). Analysis of global warming using machine learning. *Computational Water, Energy, and Environmental Engineering*, *7*(03), 127.https://doi.org/10.4236/cweee.2018.73009

9. Roy, T. J., & Ashiq Mahmood, M. (2023). Global Warming and Bangladesh: A Machine Learning Approach to Analyze the Warming Rate Utilizing Neural Network. In *The Fourth Industrial Revolution and Beyond: Select Proceedings of IC4IR+* (pp. 19-30). Singapore: Springer Nature Singapore. https://doi.org/10.1007/978-981-19-8032-9_2

10. Hema, D. D., Pal, A., Loyer, V., & Gaurav, R. (2019). Global Warming Prediction in India using Machine Learning. *International Journal of Engineering and Advanced Technology (IJEAT)*, *9*(1).

11. Kaur, S., & Randhawa, S. (2018, July). Global land temperature prediction by machine learning combo approach. In *2018 9th International Conference on Computing, Communication and Networking Technologies (ICCCNT)* (pp. 1-8). IEEE. https://doi.org/10.1109/ICCCNT.2018.8494173

12. Bindajam, A. A., Mallick, J., Talukdar, S., Shahfahad, Shohan, A. A. A., & Rahman, A. (2022). Modeling the spatiotemporal heterogeneity of land surface temperature and its relationship with land use land cover using geo-statistical techniques and machine learning algorithms. *Environmental Science and Pollution Research*, 1-19. https://doi.org/10.1007/s11356-022-23211-5

13. An, Y., Wang, X., Qu, Z., Liao, T., & Nan, Z. (2018). Fiber Bragg grating temperature calibration based on BP neural network. *Optik*, *172*, 753-759. https://doi.org/10.1016/j.ijleo.2018.07.064

14. Adusumilli, S., Bhatt, D., Wang, H., Devabhaktuni, V., & Bhattacharya, P. (2015). A novel hybrid approach utilizing principal component regression and random forest regression to bridge the period of GPS outages. *Neurocomputing*, *166*, 185-192. https://doi.org/10.1016/j.neucom.2015.03.080

15. Shams, M. Y., Tarek, Z., Elshewey, A. M., Hany, M., Darwish, A., & Hassanien, A. E. (2023). A Machine Learning-Based Model for Predicting Temperature Under the Effects of Climate Change. In *The Power of Data: Driving Climate Change with Data Science and Artificial Intelligence Innovations* (pp. 61-81). Cham: Springer Nature Switzerland. https://doi.org/10.1007/978-3-031-22456-0_4

16. Breiman, L. (2001). Random forests. *Machine learning*, *45*, 5-32. https://doi.org/10.1023/A:1010933404324

17. Kurniawan, R., Setiawan, I. N., Caraka, R. E., & Nasution, B. I. (2022). Using Harris hawk optimization towards support vector regression to ozone prediction. *Stochastic Environmental Research and Risk Assessment*, *36*(2), 429-449. https://doi.org/10.1007/s00477-022-02178-2

18. Suykens, J. A. (Ed.). (2003). *Advances in learning theory: methods, models, and applications* (Vol. 190). IOS Press. Stone, M. H. (1948). The Generalized Weierstrass Approximation Theorem. Mathematics Magazine, 21(5), 237–254. https://doi.org/10.2307/3029337

19. Siami-Namini, S., Tavakoli, N., & Namin, A. S. (2018). A comparison of ARIMA and LSTM in forecasting time series. In *2018 17th IEEE international conference on machine learning and applications (ICMLA)* (pp. 1394-1401). IEEE.

20. Peckov, A. (2012). A machine learning approach to polynomial regression. *Ljubljana, Slovenia, URL: http://kt. ijs. si/theses/phd_aleksandar_peckov. pdf.*

21. Stone, M. H. (1948). The generalized Weierstrass approximation theorem. *Mathematics Magazine*, *21*(5), 237-254.

# 8. Supplementary Materials

The datasets and code supporting the findings of this study are openly available in the "global-climate-change-prediction" repository on GitHub at the following URL: https://github.com/flistz/global-climate-change-prediction.

Below is a description of the contents of the repository:

1. Annual_Surface_Temperature_Change.csv: Contains historical data on surface temperature changes. Used to correlate global warming trends with time.
2. Climate_related_Disasters_Frequency.csv: Tracks the frequency of various climate-related disasters. This data helps to analyze the impact of global warming on natural disaster occurrences.
3. Forest_and_Carbon.csv: Details the carbon sequestration data in global forest areas, pivotal for understanding the carbon cycle's role in climate change.
4. Land_Cover_Accounts.csv: Includes land cover classification data to study the effects of land cover changes on global warming.
5. all.csv: This file includes ISO Alpha-3 country codes and the corresponding global regions and sub-regions. It serves as a reference for mapping country-specific data to their geographic locations in analyses. This is particularly useful for merging datasets on the basis of country or region and for ensuring accurate geographical categorization in our global warming prediction models.
6. run.ipynb: A Jupyter notebook containing the SVR model implementation, data preprocessing steps, and prediction outputs. This notebook serves as the backbone for the machine learning analysis in the study.
7. README.md: Provides an overview of the repository, detailing how each file contributes to the study and instructions on how to replicate the analysis.

Each file is integral to the analysis presented in this paper, and their inclusion in the repository ensures transparency and reproducibility of the research findings.