

# Machine learning for long-term monitoring and prediction of sub-surface temperatures for a campus-scale geothermal exchange field

## Aprendizaje automático para el monitoreo y la predicción a largo plazo de las temperaturas del subsuelo en un campo geotérmico a escala de campus

Shubham Dutt Attri

*Biological Systems Engineering, University of Wisconsin–Madison, USA, sattri@wisc.edu (Standards and Compliance Branch, Efficiency Division, California Energy Commission, California)*

Mingxue Jiang, Evan Heeg, and Dante Fratta, **James Tinjum**

*Department of Civil and Environmental Engineering, Geological Engineering Program, University of Wisconsin–Madison, USA, jmtinjum@wisc.edu*

Orhun Aydin

*Department of Earth and Atmospheric Sciences, Saint Louis University, USA*

David Hart

*Wisconsin Geological and Natural History Survey, University of Wisconsin–Madison, USA*

**ABSTRACT:** We apply machine learning (ML) models to predict subsurface temperature development with optimal hyperparameter tuning and less complexity. We use seven years of geothermal temperature data from a cooling-dominated campus-scale, instrumented geothermal heat exchange (GHX) field (2596, 152-m-deep boreholes over an area of 280 m by 360 m) in the Midwest region of the United States. The field temperatures are monitored using eight temperature monitoring wells, or TMWs, and this study discusses the analysis on one of the wells, TMW1. We use linear regression and two tree-based ML models (random forest regression, or RFR, and XGBoost) with five input features for training—two from weather data (air temperature and humidity) and three temperature parameters relating to the energy exchanged for heating and cooling the campus buildings and the geothermal field. The primary model fitting shows root mean square error (RMSE) values varying from 1.00° to 1.15° C. We then created lagged variables for each input variable (up to 6 months) and used them to make six-month predictions. The RMSE value was reduced to 0.71° C for an optimized RFR model. Findings also showcase a gradual, seasonal rise in subsurface temperature, offering valuable insights for designing more efficient GHX systems, conducting improved energy balance assessments, and creating long-term ground temperature change models.

**KEYWORDS:** machine learning, geothermal energy, data analytics, temperature predictions, decision trees

## 1. INTRODUCTION

Geothermal heat exchange systems (GHXs) are clean, renewable, energy-efficient technologies that transfer subsurface thermal energy to provide highly efficient heating and cooling for residential and commercial buildings by taking advantage of near-constant, year-round subsurface temperatures. The increasing development of GHXs has been portrayed to have significant economic and environmental benefits (Tinjum et al., 2023; Urcheguia et al. 2008, Michopoulos et al. 2013, Carvalho et al., 2015). GHXs are more energy efficient than conventional heating and cooling systems and have a long-term positive impact on the environment and economy (Bloom and Tinjum 2016). Fiber-optic Distributed Temperature Sensing (DTS) uses the interaction of laser pulses with imperfections in the silica in a fiber to sense the temperature in a medium. Distributed temperature data also provide insights into the spatial and temporal variations of temperature, allowing for qualitative and quantitative analyses of

the heterogeneous ground's thermal properties (Herrera et al. 2018, Attri et al. 2023). The determination of temperature is based on Stokes ( $P_s$ ) and anti-Stokes ( $P_{as}$ ) backscattered signals from position ( $z$ ) along the fiber at the time ( $t$ ). The instrumental equation is conveyed by (Van de Glesen 2012, McDaniel et al. 2018a, Tombe et al., 2020):

$$T(z, t) = \frac{\gamma}{\ln \frac{P_s(z, t)}{P_{as}(z, t)} + C - \int_0^z \Delta \alpha \cdot z' dz'} \quad (1)$$

where  $\gamma$  represents the energy difference between the incoming and backscattered Raman photons,  $C$  is a constant that depends on the input laser in the interrogator, and  $\Delta \alpha$  is the differential attenuation between the anti-stokes and the stokes signals in the fiber.

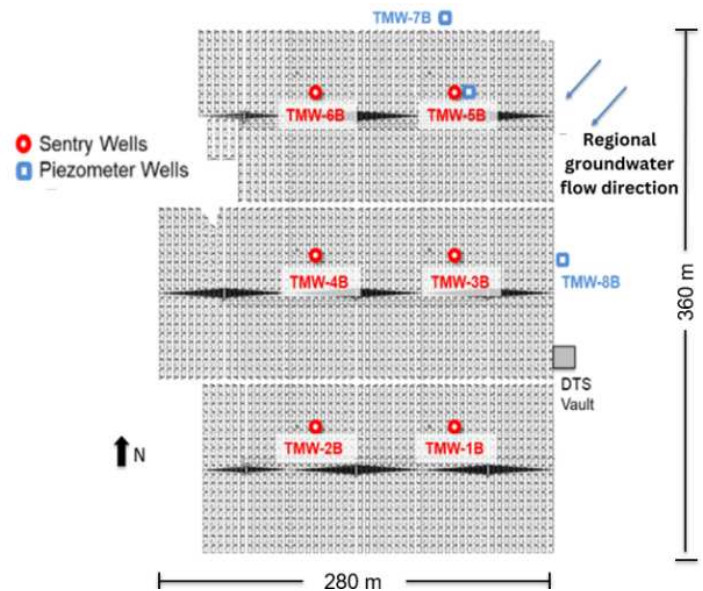
Machine Learning (ML) is a rapidly growing field that integrates computer science and statistical applications and is at the core of artificial intelligence and data science (Jordan and Mitchell 2015). The core concept of ML involves constructing models that can

make informed predictions about future events or trends. This is achieved by training algorithms on historical data to discern patterns and relationships, which are then used to anticipate future outcomes within the scope and quality of the data provided (Tut Haklidir and Haklidir 2020). Fu Jiao Tang et al. (2022) successfully predicted the heat exchange capacity of a borehole heat exchanger (BHE) by using ML methods such as linear regression (LR), polynomial regression (PR), artificial neural work (ANN), and random forest (RF) to compute the annual heat exchange rates of 400 thermal performance tests (TPT) covering 12 factors. The results show that the PR approach performs best, with root-mean-square error (RMSE) less than  $1.74 \text{ W}\cdot\text{m}^{-1}$  and a Coefficient of Determination ( $R^2$ ) higher than 0.99. Tut Haklidir and Haklidir (2020) developed a deep-learning model to predict reservoir temperatures based on hydrogeochemical data while comparing their results to traditional regression approaches neural networks (DNN). The study showed that the DNN algorithm generated the lowest errors and provided accurate values close to geothermometer calculations. Another related research was conducted by Zhang et al. (2022) on the prediction of coefficient of performance (COP) models for heat pump units and ground-source heat pumps (GSHP) using ML methods. Zhang et al. (2022) used algorithms such as extreme learning machine (ELM), support vector machine (SVM), and back propagation neural network (BPNN) to predict the COP for heat pump units and a GSHP system. They used ten parameters for a feature-shrinkage and selection-operator (LASSO) approach. The results indicate that the ELM model has better prediction accuracy than the other models. In this paper, we expand the analysis of measured DTS data by applying machine learning for both causal and predictive analyses of the subsurface temperatures for a campus-scale, low-enthalpy GHX field. Predicting temperature variations is crucial for optimizing the performance and efficiency of geothermal systems. Accurate long-term predictions can reduce the costs in the maintenance and planning of geothermal energy extraction, making the application of ML techniques valuable.

Epic Systems (Epic) is an electronic health records company located in the Midwest of the United States with a corporate campus of over 13,000 employees. Epic uses geothermal reservoir fields to heat and cool its campus (Özdoğan Dölçek et al. 2017, McDaniel et al. 2018b). The system includes four borefields to deliver 48.5 MW of cooling power. Figure (1) shows the overview of the largest borehole field, borefield 4 (BF4). This field has 2,596 GHX wells in a volume of 360 m north to south, 280 m east to west, and 152 m deep for  $15.4 \cdot 10^6 \text{ m}^3$  of porous media available for thermal exchange (Attri et al. 2023). BF4 alone contributes to over fifty percent of the ground-based cooling capacity and is one of the world's largest, shallow, low-temperature GHX systems (Tinjum et al. 2023). It has temperature monitoring wells (TMWs) containing fiber-optic cables extending full depth. The fiber-optic cables used are OM2 ClearCurve Plenum Orange cables with a multi-mode 50/125- $\mu\text{m}$  core, 2-mm outer diameter, and E2000 APC connectors (McDaniel et al. 2018a). Temperature data has been consistently monitored since June 2016. The black solid circles indicate the locations of fiber optic loops in the field, which are used to detect ground temperature.

## 2. METHODOLOGY

This study is focused on BF4, which is equipped with temperature monitoring wells (TMWs) that include fiber-optic loops extending to the base. Figure 1 represents the 2596 GHX boreholes as the black rectangles as a frame and the red circles indicate the sentry wells or the TMWs, with fiber loops grouted directly in contact with the ground. The blue squares represent the piezometric wells in the field, which also have fiber-optic loops, as well as piezometer screens in both a shallow and a deep aquifer, and TMW-1B, or simply TMW1, has been selected for the preliminary stage of ML application since it provides the most optimal temperature data that reflect the overall geothermal behavior. The initial predictive analysis targets the 80-m-deep temperature data from TMW1, which will later be expanded to include various depths across all wells. Our analysis relies on five principal independent variables to forecast the primary variable, 'Well Temperature 1': relative humidity, air temperature, BF4 temperature differentials, chilled water temperature differentials, and hot water temperature differentials. The air temperature gauged near BF4 represents dry bulb temperature, and the relative humidity denotes the air's moisture content that affects the transfer of heat and moisture in the conditioned space.



**Figure 1: Borefield 4 (BF4) map with fiber-optic temperature monitoring wells (after McDaniel et al. 2018b).**

The BF4 temperature differentials represent the temperature change in the water circulating within BF4 before and after it circulates through the field, showcasing the energy exchange within the borefield. The terms 'chilled water difference' and 'hot water difference' denote the temperature disparities between the inflow and outflow of water in one of the campus's central energy plants, illustrating the energy interactions facilitated by the building's heat exchangers. We refined the ML algorithms (Linear Regression, XGBoost, and Random Forest) to predict our target variable using a set of those six independent variables. The dataset was compiled with measurements taken at six-hour intervals from

June 2016 to December 2022. Recognizing that these variables might not immediately affect subsurface temperatures, we improved the models, specifically XGBoost and Random Forest, by incorporating lagged versions of the original variables. We created new predictors by introducing delays of up to six months with one-month increments for each variable, resulting in six additional lagged variables per original variable to capture their delayed influence on the target variable. Finally, we used all 36 variables to predict up to six months of well temperature.

Linear regression, random forest regression (RFR), and XGBoost were selected for this analysis because of their demonstrated effectiveness in modeling the complex relationships found in the data, subsurface temperature variations in our case. Linear regression serves as a baseline model with straightforward interpretation, aiding in understanding fundamental trends and relationships within the data. Random forest regression was chosen for its ability to mitigate overfitting and manage high-dimensional data, making it ideal for capturing intricate patterns in the geothermal temperature dataset. XGBoost, known for its scalability and efficiency, was selected for its exceptional predictive performance and ability to process large datasets with numerous features, including lagged variables. Together, these methods provide a robust framework for modeling and predicting subsurface temperature dynamics.

### 2.1 Linear Regression

Linear Regression is one of the most common ML approaches (Tang et al. 2022). Linear regression analysis allows us to predict the future by discerning linear relationships between the dependent and independent variables (Ansari and Nassif 2022). The most straightforward format of linear regression, univariate linear regression, involves one independent variable and one dependent variable:

$$y_1 = b_l + a_1 x_1 \quad (2)$$

This equation represents a simple linear regression where  $y_1$  is the dependent variable,  $b_l$  is the constant term intercepting the y-axis,  $x$  is the independent variable, and  $a_1$  is the coefficient for the predictor variable. Linear regressions involving more than one independent and dependent variable are called multiple linear regressions. The function is defined as:

$$y_i = b_l + [a_1 \ a_2 \ \dots \ a_{m-1} \ a_m] * \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_{m-1} \\ x_m \end{bmatrix} + e \quad (3)$$

where  $y_i$  is the regressed,  $m$  represents the number of the variables, and  $e$  is the regression model error.

The main objective of linear regression modeling is to determine the optimal line that minimizes the difference between predicted and observed values. This line's slope represents the rate at which the dependent variable changes in response to changes in the independent features. Model evaluation is necessary to determine whether the linear regression model has the best-fit line. The most common approach is  $R^2$ :

$$R^2 = \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y}_i)^2} \quad (4)$$

which is defined as a ratio of variation to the total variation. The closer  $R^2$  is to 1, the more accurate the regression is. Root Mean Squared Error (RMSE) is another evaluation criterion to determine our linear regression model's fit. It indicates the average difference between predicted and actual values and describes how well the data points match the expected values. We aim to minimize the RMSE value, and for a perfect model, where the predicted values are identical to the actual values for all instances in the dataset, the RMSE would be 0.

$$RMSE = \sqrt{\frac{RSS}{n}} = \sqrt{\frac{\sum_{i=2}^n (y_i^{actual} - y_i^{predicted})^2}{(n-2)}} \quad (5)$$

### 2.2 XGBOOST

Decision trees are binary trees where each internal node represents a test on an attribute, each branch represents the outcome of the test, and each leaf node represents a class label (Navada et al., 2011). Decision trees have transformed into versatile tools applicable across various disciplines, including artificial intelligence, machine learning, knowledge discovery, and data mining (de Ville, 2013). They are now considered highly cross-disciplinary, general-purpose methods that are computationally intensive and used for prediction and classification (de Ville 2013). Different decision trees are available depending on the specific situation and desired outcome; these trees include Classification, Regression, Decision Forests, and Classification and Regression; (Navada et al. 2011). The general idea of a decision tree is that it repeatedly divides the data into smaller groups based on the values of input features, creating branches that contain similar data within them but different data between them at each tree level (de Ville 2013).

XGBOOST is a scalable ensemble method based on gradient boosting and built based on decision trees; it creates a forest of by sequentially optimizing decision trees with respect to the loss function. (Chen and Guestrin 2016). Gradient boosting sequentially adds new trees to an ensemble, with each sequential tree reducing the errors of the previous ensemble (Natekin and Knolls 2013). It aims to construct new base-learner models highly correlated with the negative gradient of the loss function associated with the ensemble, leading to improved accuracy in predicting the response variable (Natekin and Knolls, 2013). Since XGBoost exclusively utilizes decision trees as base classifiers, it employs a modified loss function to manage the trees' complexity (Bentéjac et al. 2021). Equation (6) is the loss function:

$$L_{xgb} = \sum_{i=1}^N L(y_i, F(X_i)) + \sum_{m=1}^M \Omega(h_m) \quad (6)$$

which is the sum of the individual losses and regularization terms. It measures the discrepancy between the predicted and actual values (Chen and Guestrin 2016). In the equation,  $N$  is the total number of data points,  $L(y_i, F(X_i))$  is the loss function that measures the difference between the predicted value  $F(X_i)$  and

the actual value  $y_i$ .  $M$  is the total leaves in all trees, and  $\Omega(h_m)$  is the regularization term for the  $m$ -th tree that can be defined using equation (7) (Chen and Guestrin 2016):

$$\Omega(h_m) = \gamma T + \frac{1}{2} \lambda \|w\|^2 \quad (7)$$

Here,  $\gamma$  is the regularization parameter that controls the complexity of the trees,  $T$  is the number of the tree leaves,  $\lambda$  controls the magnitude of the leaf weights, and  $w$  is the output score of the leaves in tree  $m$ , which are influenced by the modified loss function (Chen and Guestrin 2016). Higher values of gamma ( $\gamma$ ) lead to the creation of simpler trees by setting the threshold for the minimum loss reduction required to split an internal node, and tree complexity can be controlled by limiting the depth of the trees, which additionally speeds up model training and reduces the storage space needed (Bentéjac et al. 2021).

### 2.3 Random Forest

Random Forest (RF) analysis is a common ensemble-based method in ML (Breiman, 2001). Unlike XGBoost, RF is a bootstrap aggregating algorithm that builds many decision trees on various sub-samples of the dataset and averages the results to improve the predictive accuracy and control over-fitting (Breiman, 2001). Bootstrap aggregating enhances the stability and accuracy of ML algorithms used for regression and classification tasks (Breiman, 1996). RF approach decision forests can adapt to non-linear solutions, predicting better than LR models (Schonau and Zou, 2020). The algorithm builds individual decision trees on bootstrap samples and averages their predictions. The algorithms can be used for classification and regression models (Schonau and Zou, 2020). Decision trees are models constructed from training data by making a series of binary splits at each tree node. Each split divides the data into two child nodes based on an inequality query on one of the variables. The tree grows until each leaf node contains exactly one data point or until a stopping condition is met. A new data point is passed through the tree to make predictions, following the queries at each node. The predicted output is the data point's value in the leaf node where the new point ends up or the average value of the data points in that leaf. The key to the algorithm is the optimization process at each node, which selects a variable and a threshold to create the split (Breiman, 2001).

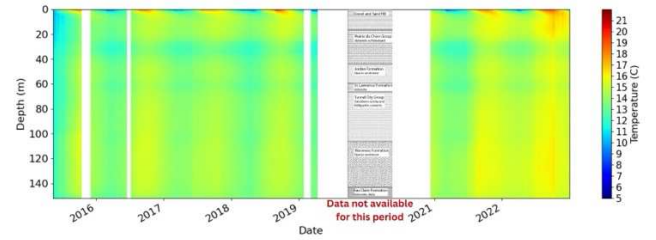
## 3. RESULTS AND DISCUSSION

### 3.1 Data Collection and Analysis

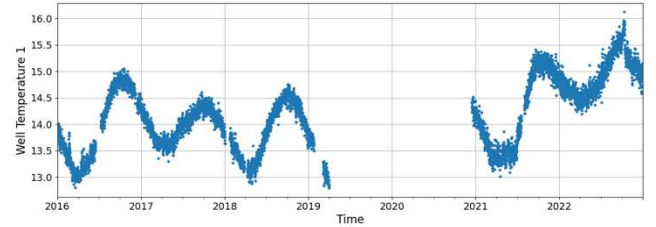
Our approach employs Fiber-Optic Distributed Temperature Sensing (FO-DTS) with dynamic, double-ended calibration for ongoing field monitoring. As shown in Figure (2), the temperature distribution for TMW-1 is depicted on a depth-time plot alongside the site's geological strata. Interruptions in data collection during the COVID-19 pandemic were due to malfunctions and operational disruptions. We have also shown our work on imputing this 2-year gap in the data using time series models like ARIMA and Holt Winters' Exponential Smoothing (Attri et al. 2024). The data demonstrate the relative stability of subsurface temperatures compared to surface temperatures, which fluctuate with

atmospheric conditions. Predominantly a cooling system, our Borefield 4 network stores more heat during the cooling season than it extracts during the annual heating season, leading to a general ground temperature increase over seven-plus years. A notable cooler area at approximately 30 m depth corresponds to groundwater flow within a karstic dolomite formation.

Figure (3) shows the temperature changes at an 80-m depth within the TMW-1. The analysis indicates a slight temperature decrease leading up to 2019, followed by an increase from 2021 onwards. Before 2016, the data were more variable, likely due to the initial stages of field startup and calibration of the FO-DTS technology. Data before 2016 was excluded from the analyses to ensure a more accurate analysis.



**Figure 2: Borehole temperature variation for TMW-1 with depth from 2015-2022.**



**Figure 3: Borehole temperature variation at 80 m depth for TMW-1 from 2015-2022.**

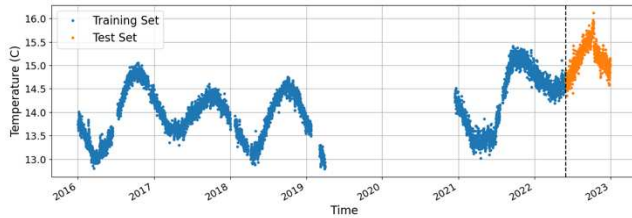
### 3.2 Linear Regression

Linear Regression (described in Section 2.1) is a simple and commonly used algorithm in ML. It fits a linear equation between dependent (target) and independent (features) variables. A correlation analysis was conducted before our ML data training, which helped extract the highly correlated features and select the six features we used for our study. The next step was to clean data for any outliers and null values. It is to be noted that we did not use the values for the period 2019 and 2020 to train our models since the data were missing due to system failures and restrictions due to COVID-19 pandemic. Our dataset was then divided into training and testing sets, with the training set comprising data from early 2016 to June 2022, as shown by the blue points in Figure (4), while the test set included data from June 2022 to the end of 2022, represented by the orange points. The black dashed line in Figure (4) separates the training and test sets. We trained the multivariate linear regression model using the testing set subsequently to make predictions for the period of the last six months of 2022.

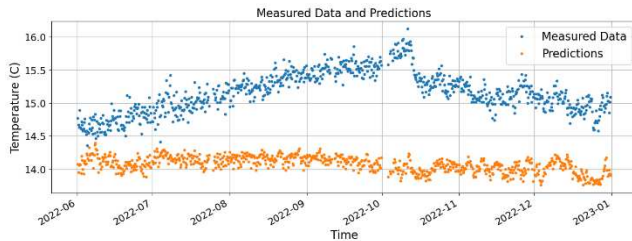
Figure (5) compares the model's forecasted values and the actual 'Well Temperature 1' measurements from the test dataset. The



model's predictions, presented in Figure (5), show a mismatch to the observed data, evidenced by an RMSE of 1.13. Such mismatch highlights the intricacies of the modeled system and underscores the constraints inherent in using linear regression. Linear regression assumes a linear relationship, constant unit change, between the dependent and independent variables, thus is not applicable to nonlinear temperature profile observed over time. Additionally, the model's vulnerability to outliers is notable; such data points can impact the regression coefficients and the model's overall predictive accuracy. This led us to use more complex models for predictions.



**Figure 4: The dataset is divided into training and testing sets.**



**Figure 5: Measured temperature data at 80-m depth for TMW1 and predicted temperatures for the test set using Linear Regression.**

### 3.3 XGBoost

The linear regression model, often considered a fundamental benchmark, yielded underwhelming results, leading to the exploration of more complex, tree-based models. These models are celebrated for their adeptness in handling multiple features and complex interactions within the data. Initially, our focus was on a model incorporating six primary features. This model underwent a rigorous full-factorial hyperparameter optimization to determine the optimal model parameters. The optimal XGBoost model performance on a six-month forecast resulted in an RMSE of 1.16, which fell short of the benchmarks set by the linear regression model. XGBoost generally performs better than the linear regression models, but this slightly decreased accuracy on the test set might be due to the increasing trend not being captured by the XGBoost model. The low R-squared value of 0.03 from the training set further corroborated that the select features did not explain the variance in the measured borehole temperatures.

To enhance the model's predictive capabilities, we expanded the feature set by incorporating lagged variables, taking the total count to 36 while maintaining the original target variable. Figure (6)

depicts the feature importance of the XGBoost model with lagged variables, elaborating the impact of each independent variable on the model's predictions. Meanwhile, Figure (7) showcases the predictions of the XGBoost model with the lagged predictor variables. The RMSE for these predictions dropped to 0.80, a marked improvement over all previously attempted models. This notable advancement can be credited to the synergistic effect of a more intricate modeling technique coupled with a substantial and strategically selected set of input variables that more accurately reflected the target variable's variations.  $R^2 = 0.84$  for the training set supports the view that the model could account for a significant portion of the temperature variance with the extended set of features.

Upon a detailed analysis of Figure (6), the feature importance plot sheds light on the predictors that significantly drive borehole temperatures. The leading indicator was the temperature differential of the inlet and outlet waters from the borefield measured three months prior, implying a delayed influence of energy transfer on the temperature. This was closely followed by the energy exchange metrics within the borefield and the cooling loop of a campus building, both from two months earlier. The strong influence of these features likely stems from the cooling-dominant nature of the borefield system, which exerts a pronounced effect on the thermal dynamics within the field. Remarkably, the top seven determinants for temperature variability were linked to the energy transactions within the borefield and the cooling loads from the buildings, highlighting the intricate interplay between these factors in shaping the borefield's thermal profile.

### 3.4 Random Forest

Since XGBoost performed better for our dataset with more features, we also decided to use the RF method since it is a simpler model and can perform better on datasets where the relationship between features is more complex and nonlinear. It is also less sensitive to the scale of features and can handle unnormalized or standardized features. We performed hyperparameter optimization to define the optimal random forest model, and our primary model with only five features gave an RMSE value of 1.08, which is better than both LR and XGBoost with these five features. Next, we took the same steps for the lagged features as we did with XGBoost, and the RMSE obtained using RF was 0.71, the best RMSE we have achieved with our data. Figure (8) shows the relative feature importance of the features using RF, and again, the top six factors affecting temperature fluctuations were associated with the energy exchanges occurring in the borefield and the thermal demands from the buildings, with five of them with lags. Figure (9) shows a scatter plot of the measured temperature values and the RF model's predictions.

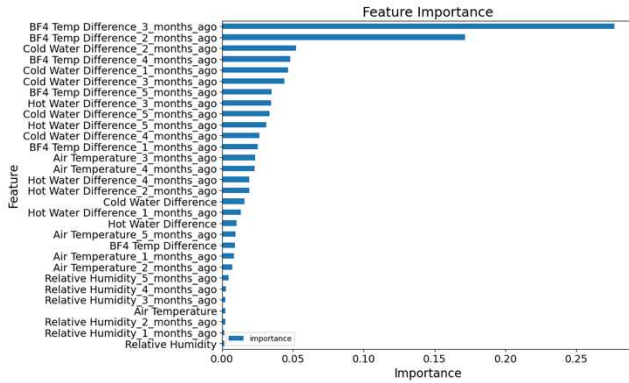


Figure 6: Feature importance of independent variables using XGBoost.

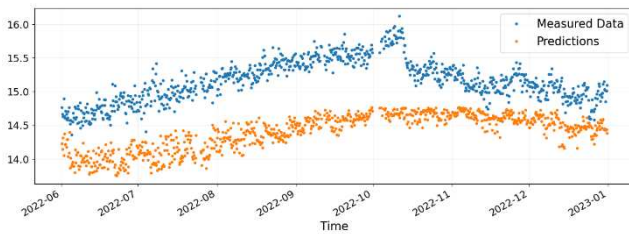


Figure 7: Measured temperature data at 80-m depth for TMW1 and predicted temperatures for the test set using XGBoost with lagged variables.

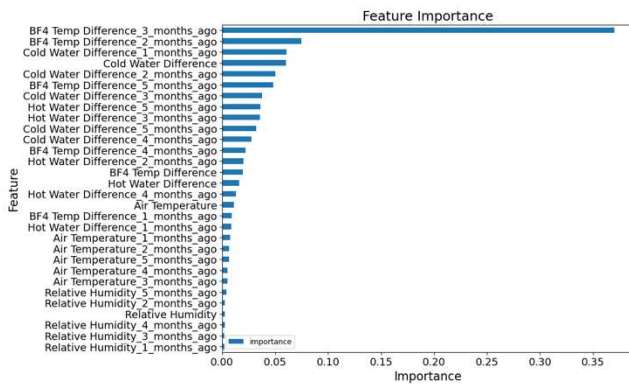


Figure 8: Feature importance of independent variables using Random Forest.

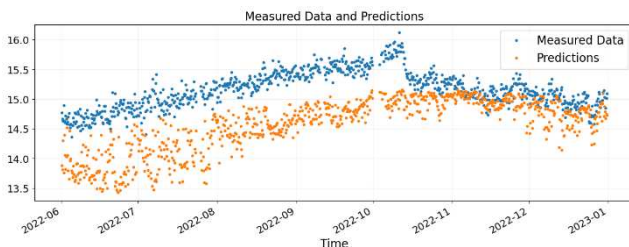


Figure 9: Measured temperature data at 80-m depth for TMW1 and predicted temperatures for the test set using Random Forest with lagged variables.

## 4. CONCLUSIONS AND FUTURE SCOPE

The research explores applying ML strategies for predicting subsurface temperatures within a low-enthalpy geothermal exchange field. It evaluates three distinct models, noting that their effectiveness varied according to the data structure and the relationships between predictors and the target variable. Despite its simplicity, linear regression prompted the consideration of more complex models to improve prediction accuracy. Random Forest model performed best with both non-lagged and lagged features. It explains the delayed impact of the energy interaction with the field on the temperatures below the ground, which agrees with the basic scientific understanding of how heat transfer and storage works sub-surface. The results are satisfactory but also far from what can be achieved. The outcomes, while promising, also highlight a consistent trend of underestimating temperature values across all models. A possible reason for this could be the limited scope of input data, which currently includes only the energy consumption from one borefield and a portion of the campus. This limitation likely contributes to the models' tendency to predict less temperature variation, particularly as temperatures have risen.

Other improvements could incorporate additional ML techniques more suited to the dataset's characteristics. The study aims to extend its analysis to encompass the full depth of TMW 1 and across all such wells to deepen our comprehension of how various factors influence ground temperatures. Efforts to amass more comprehensive data, such as flow rates and energy consumption from other campus buildings, are expected to enhance the robustness and dependability of the predictions.

The ultimate objective is to fine-tune the performance of ML models across all monitoring wells and depths, capturing the full spectrum of temperature variation indicators. Success in this area could significantly contribute to optimizing system utilization and fostering sustainable operational design.

## 5. ACKNOWLEDGEMENTS

The authors would like to acknowledge Epic Systems Corporation for their generous contributions of time, talent, and resources in the ongoing development of this research and for giving us access to a unique facility. The Morse Company, Hooper Corporation., Salas O'Brien, JP Cullen, and Teel Plastics all donated time, materials, and expertise, without which this research would not have been possible. We thank Adam McDaniel (Westwood Professional Services) for preparing the initial code to process and calibrate FO-DTS data. We acknowledge the National Science Foundation for sponsoring undergraduate research opportunities during the early phases of this project under Award # 1156674. This material is based upon work supported by the National Science Foundation Graduate Research Fellowship under Grant No. 2137424. Any opinions, findings, conclusions, or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the funding organizations.

## 6. REFERENCES

Attri, S. D., Heeg, E., Yilmaz, M., Tinjum, J. M., Fratta, D., & Hart, D.

- (2023). Long-term temperature monitoring of a campus-scale geothermal exchange field using a fiber-optic sensing array. In *Proceedings of the 48th Workshop on Geothermal Reservoir Engineering*, SGP-TR-224.
- Attri, S. D., Heeg, E., Tinjum, J. M., Fratta, D., Hart, D. J., & Aydin, O. (2024). Time series analysis for long-term monitoring and forecasting subsurface temperatures for a campus-scale geothermal exchange field. *Proceedings of the 49th Workshop on Geothermal Reservoir Engineering*, Stanford University, SGP-TR-227.
- Bentéjac, C., Csörgő, A., & Martínez-Muñoz, G. (2021). A comparative analysis of gradient boosting algorithms. *Artificial Intelligence Review*, 54.
- Breiman, L. (2001). Random forests. *Machine Learning*, 45, 5–32.
- Breiman, L. (1996). Bagging predictors. *Machine Learning*, 24(2), 123–140.
- Breiman, L., Friedman, J., Olshen, R. A., & Stone, C. J. (1984). *Classification and regression trees* (1st ed.). Chapman and Hall/CRC.
- Bloom, E. F., & Tinjum, J. M. (2016). Fully instrumented life-cycle analyses for a residential geo-exchange system. In *Geo-Chicago 2016*.
- Carvalho, A. D., Moura, P., Vaz, G. C., & de Almeida, A. T. (2015). Ground source heat pumps as high efficient solutions for building space conditioning and for integration in smart grids. *Energy Conversion and Management*, 103, 991–1007.
- Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 785–794). San Francisco, California, USA.
- de Ville, B. (2013). Decision trees. *Wiley Interdisciplinary Reviews*, 5(6), 448–455.
- Heeg, E., Tinjum, J. M., Fratta, D., Attri, S. D., Hart, D. J., & Luebke, A. K. (2024). Quantifying the long-term performance of a district-scale geothermal exchange field. *Proceedings of the 49th Workshop on Geothermal Reservoir Engineering*, Stanford University, SGP-TR-227.
- Herrera, C., Nellis, G., Reindl, D., Klein, S., Tinjum, J.M., and McDaniel, A.: Use of a Fiber Optic Distributed Temperature Sensing System for Thermal Response Testing of Ground-coupled Heat Exchangers. *Geothermics*, 71, (2018), 331–338.
- Jordan, M. I., & Mitchell, T. M. (2015). Machine learning: Trends, perspectives, and prospects. *Science*, 349(6245), 255–260.
- McDaniel, A., Fratta, D., Tinjum, J.M., and Hart, D.J.: Long-term District-scale Geothermal Exchange Borefield Monitoring with Fiber Optic Distributed Temperature Sensing. *Geothermics*, 72, (2018a), 193–204.
- McDaniel, A., Tinjum, J.M., Hart, D.J., and Fratta, D.: Dynamic Calibration for Permanent Distributed Temperature Sensing Networks. *IEEE Sensors Journal*, 18(6), (2018b), 2342–2352.
- Michopoulos, A., Zachariadis, T., & Kyriakis, N. (2013). Operation characteristics and experience of a ground source heat pump system with a vertical ground heat exchanger. *Energy*, 51, 349–357.
- Natekin, A., & Knoll, A. (2013). Gradient boosting machines, a tutorial. *Frontiers in Neuroinformatics*, 7, Article 21.
- Navada, A., Ansari, A., Patil, S., & Sonkamble, B. (2011). Overview of the use of decision tree algorithms in machine learning. In *Proceedings - 2011 IEEE Control and System Graduate Research Colloquium, ICSGRC 2011* (pp. 37–42).
- Özdoğan Dölçek, A., Atkins, I., Harper, M. K., Tinjum, J. M., & Choi, C. Y. (2017). Performance and sustainability of district-scale ground coupled heat pump systems. *Geotechnical and Geological Engineering*, 35(2), 1–14.
- Schonlau, M., & Zou, R. Y. (2020). The random forest algorithm for statistical learning. *The Stata Journal*, 20(1), 3–29.
- Tang, F., Nowamooz, H., Wang, D., Luo, J., Wang, W., & Sun, X. (2022). Heat exchange capacity prediction of borehole heat exchanger (BHE) from infrastructure based on machine learning (ML) methods. *IEEE Transactions on Intelligent Transportation Systems*, 23(11), 22409.
- Tinjum, J. M., Yilmaz, M., Heeg, E., Fratta, D., Hart, D. J., & Attri, S. D. (2023). Energy efficiency and life cycle assessment of a district-scale geothermal exchange field. 48th Workshop on Geothermal Reservoir Engineering. Stanford University, Stanford, California, February 6–8.
- Tombe, B. des, Schilperoort, B., and Bakker, M.: Estimation of Temperature and Associated Uncertainty from Fiber-optic Raman spectrum Distributed Temperature Sensing. *Sensors* (Switzerland), 20(8), (2020)
- Tut Haklidir, F. S., & Haklidir, M. (2020). Prediction of reservoir temperatures using hydrogeochemical data, Western Anatolia Geothermal Systems (Turkey): A machine learning approach. *Natural Resources Research*, 29(4), 1–10.
- Urchueguía, J. F., Zacarés, M., Corberán, J. M., Montero, Á., Martos, J., & Witte, H. (2008). Comparison between the energy performance of a ground coupled water to water heat pump system and an air to water heat pump system for heating and cooling in typical conditions of the European Mediterranean coast. *Energy Conversion and Management*, 49(10), 2917–2923.
- van de Giesen, N., Steele-Dunne, S. C., Jansen, J., Hoes, O., Hausner, M. B., Tyler, S., & Selker, J. (2012). Double-ended calibration of fiber-optic Raman spectra distributed temperature sensing data. *Sensors*, 12, 5471–5485.
- Zhang, X., Wang, E., Liu, L., & Qi, C. (2022). Machine learning-based performance prediction for ground source heat pump systems. *Geothermics*, 105, 102509.

# INTERNATIONAL SOCIETY FOR SOIL MECHANICS AND GEOTECHNICAL ENGINEERING



*This paper was downloaded from the Online Library of the International Society for Soil Mechanics and Geotechnical Engineering (ISSMGE). The library is available here:*

<https://www.issmge.org/publications/online-library>

*This is an open-access database that archives thousands of papers published under the Auspices of the ISSMGE and maintained by the Innovation and Development Committee of ISSMGE.*

*The paper was published in the proceedings of the 17th Pan-American Conference on Soil Mechanics and Geotechnical Engineering (XVII PCSMGE) and was edited by Gonzalo Montalva, Daniel Pollak, Claudio Roman and Luis Valenzuela. The conference was held from November 12<sup>th</sup> to November 16<sup>th</sup> 2024 in Chile.*