# Exploring Autoencoders and XGBoost for Predictive Maintenance in Geothermal Power Plants

Rudraksh Nanavaty

Department of Computer Science and Engineering, Pandit Deendayal Energy University, Gandhinagar, Gujarat, India

rudrakshnanavaty@gmail.com

**Keywords:** Geothermal Power Plants; Machine Learning; Predictive Maintenance; Autoencoders; XGBoost; Anomaly Detection; Remaining Useful Life (RUL) Prediction; Industry 4.0

## ABSTRACT

Geothermal power plants, harnessing the Earth's internal heat, are an essential catalyst for the transition to clean energy. However, the reliable operation of these plants requires effective maintenance strategies. This paper explores predictive maintenance in geothermal power plants, focusing on the application of Autoencoders and XGBoost. Autoencoders, mainly used for unsupervised learning, are apt for detecting anomalies in normal equipment behavior. XGBoost is an efficient ensemble model suitable for predicting the remaining useful life of critical components. Both algorithms can be used to make interpretable models that can help pinpoint the root cause of equipment failure. This paper reviews research applying these strategies to analogous systems, emphasizing the need for solutions tailored to geothermal contexts. It also highlights challenges, including limited data and literature, that underscore the need for collaboration between geothermal experts and computer scientists to unlock the full potential of predictive maintenance in geothermal power plants.

## INTRODUCTION

The scorching depths of our planet hold a vast, untapped reservoir of clean energy waiting to be harnessed. Capitalizing on the Earth's internal heat, geothermal power plants have demonstrated remarkable potential in providing clean, reliable, and continuous electricity. Geothermal energy is projected to account for 2-3% of global electricity generation by 2050 (van der Zwaan & Dalla Longa, 2019). Hence, geothermal energy has the potential to contribute significantly to the transition to clean energy.

However, geothermal power plants are capital-intensive projects and demand significant upfront investment for proper functioning. The machinery operates under extreme conditions, exposed to high temperatures, corrosive fluids, and abrasive steam. The consequences of failure are far-reaching. Unplanned outages can not only cripple energy production but also lead to expensive repairs, environmental damage, and even safety hazards. Hence, rigorous maintenance strategies are imperative to ensure their longevity and optimal performance.

Maintenance in geothermal power plants can be broadly classified into Corrective, Preventive, and Predictive. Corrective (run-to-failure) maintenance involves correcting faults and executing repairs after a failure event. This approach not only results in higher downtime but also leads to reduced equipment life. Preventive maintenance involves scheduled services to prevent potential failures. This mitigates the concerns of the previous approach. However, there remains a chance of premature or late maintenance, meaning an alternate approach could lead to more optimal performance and better equipment life.

In contrast, Predictive Maintenance (PdM) leverages advanced technologies like Machine Learning (ML) and Deep Learning (DL) to forecast equipment failures before they happen. The Remaining Useful Life (RUL) of equipment can be predicted using a regression model, while a classification model can predict whether a failure event is imminent. These approaches can enhance operational efficiency, reduce downtime, and optimize the return on investment in geothermal power infrastructure.

The subject of PdM for geothermal power plants is still relatively new, with little research to depend on. To address this disparity, this paper explores some scenarios that extend beyond the geothermal landscape. It investigates discoveries from adjacent fields, like PdM for water pumps in areas other than geothermal power plants or concrete pumps. While working settings vary, the fundamental principles of wear and tear, vibration measurement, and data-driven failure prediction are frequently applicable. By critically reviewing this research and applying their findings to the unique problems of geothermal systems, we can gain valuable knowledge and speed the development of viable PdM solutions for geothermal power plants.

This paper explores the application of Autoencoders and XGBoost, two powerful ML tools, for building robust PdM models in geothermal power plants. By deploying such advanced predictive methodologies, this paper paves the way for a new era of optimized operations, maximizing the potential of geothermal energy as a key pillar of our sustainable future.

## ANOMALY DETECTION USING AUTOENCODERS

The effects of equipment faults/failures on geothermal power plants can range from poor performance to catastrophic failures. Due to the time scale of the problems, they may often go undetected. Hence, sophisticated models are needed to aid in detecting possible faults/anomalies (Liu et al., 2022).

Autoencoders are a category of Artificial Neural Networks (ANNs) designed for unsupervised learning. Their ability to identify the essential features of the given data also enables them to detect the anomalies or outliers present in it. They follow an encoder-decoder

architecture. The encoder compresses the input data into a lower-dimensional representation (the "latent space"). Then, the decoder reconstructs the original input from this latent representation. The training involves minimizing reconstruction error and encouraging the model to learn a meaningful data representation. The network is prompted to learn the most essential features that represent the "normal" operating state of the system (Goodfellow et al., 2016).
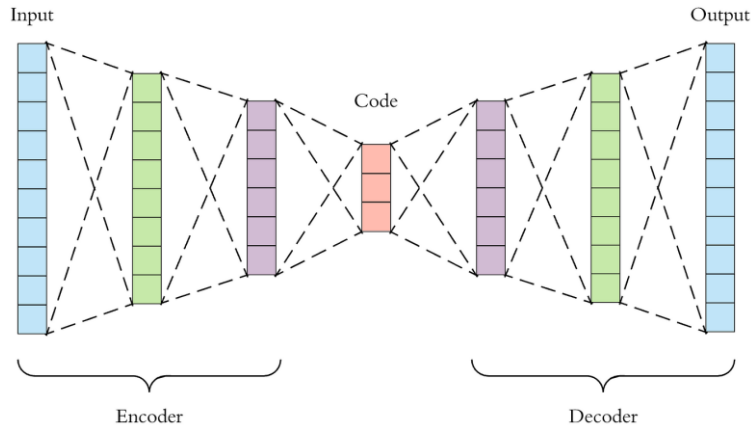


**Figure 1: Basic autoencoder structure, showing the input, latent space (or code), and the reconstructed output. The Encoder neural network compresses the input to code, which is then reconstructed by the decoder. The reconstruction error is used to identify anomalies. Image from "Applied Deep Learning - Part 3: Autoencoders", via Towards Data Science https://towardsdatascience.com/applied-deep-learning-part-3-autoencoders-1c083af4d798**

A simple threshold-based detection could predict anomalies in the data. Thresholding involves setting a threshold for the reconstruction error. Any data point that exhibits an error exceeding this threshold is flagged as an anomaly. This simple approach requires careful threshold value selection based on the dataset and application. However, manual interpretation for threshold anomaly detection tends to be subjective and varies among experts. This approach is prone to false positives. Givnan et al. (2021) presented an ML strategy for modeling machine behavior and automatically detecting anomalies to overcome this. The model learns generalizations and assigns thresholds based on fault severity. They classify using a traffic light system, with green, amber, and red indicating normal behavior, concerning behavior, and machine fault, respectively. Their findings showed that detecting anomalies in the amber range and raising alerts before the system fails is possible.

Liu et al. (2022) employed a Latent-Space Dynamic Neural Network (LSDNN) coupled with Principal Component Analysis (PCA) for fault detection in a binary cycle geothermal power plant. Time-series field data was collected from a plant operated by Cyrq Energy Inc. in Salt Lake City, UT, USA, for 5 years with an hourly resolution. The model detected abnormalities before the production pump and power generation unit field operators. The authors also suggest extending the model to incorporate fault diagnosis, potentially pinpointing the root causes of failures.

Hajgató et al. (2022) developed an explainable deep convolutional autoencoder PredMaX for anomaly detection in high-dimensional, unlabeled temporal data. They apply an explainable deep convolutional autoencoder followed by PCA. They perform automatic clustering in the latent space rather than the principal components space to ensure higher accuracy. Explainability allows their reasoning module to help identify the root cause of failure, aiding in deciding suitable maintenance measures. They demonstrated the framework using gearbox operation data collected by 98 sensors for over 1.5 years with a 1-minute resolution in a MW load range in a petrochemical plant. Their model successfully identified imminent maintenance, and appropriate measures were carried out timely.

Tian et al. (2022) proposed an anomaly detection method leveraging a convolutional autoencoder that can train only on normal operation data. Their approach incorporates a sliding window algorithm to generate sensor readings for inputs, which account for the dynamic characteristics of the data. They employed confusion matrices to determine the reconstruction error threshold for anomaly detection and assess the classification algorithm's effectiveness. The proposed CAE-AD model was assessed using a real-world benchmark dataset comprising sensor data from real-world water pumps. The model successfully identified all pump malfunctions and 98.8% of abnormal data.

Alamu et al. (2020) employed autoencoders to predict anomalous behavior in Electrical Submersible Pumps (ESPs). They used 2 years of historical data from 97 sensors on an Electrical Submersible Pump (ESP). The autoencoder consisted of 7 layers and used an exponential linear unit as the activation function. During the data collection period, the pump saw 5 significant trips, 3 due to gas locks and 2 from electrical issues. The model accurately predicted the gas locks on average 5 hrs and electrical issues several days before the actual events. The relative differences in reconstruction errors from individual sensors helped pinpoint the top ten sensors responsible for each event. They also consulted a Subject Matter Expert to validate the outputs.

While anomaly detection through autoencoders offers a promising path for predicting impending failures in geothermal power plants, the "best choice" model remains elusive due to the nascent research stage and limited literature. Nonetheless, the potential for autoencoders

for PdM for equipment in geothermal power plants is significant. Their ability to learn latent representations of normal operating conditions allows for identifying deviations and triggering timely maintenance interventions. Reconstruction losses further offer insights into the root cause of anomalies, guiding targeted troubleshooting. However, choosing the optimal autoencoder architecture requires careful consideration of project-specific factors. Data availability, complexity, and computational resources all play crucial roles. Vanilla autoencoders, for instance, may suffice for simple data patterns, while more intricate architectures like variational autoencoders or recurrent neural networks might be necessary for capturing complex temporal dependencies. Choosing the suitable model demands a fine-tuned approach that balances performance with resource constraints and data characteristics. Only by understanding these nuances can we unlock the full potential of autoencoders for robust and effective PdM in geothermal power plants.

### PREDICTIVE MAINTENANCE WITH XGBOOST

eXtreme Gradient Boosting (XGBoost), is a powerful ML algorithm noted for its efficiency and accuracy, particularly in predictive modeling and classification tasks. It is a supervised ensemble model that sequentially trains a series of weak learners, typically decision trees, and combines their outputs to form a robust predictive model. The algorithm employs a gradient boosting framework, where each subsequent tree corrects the errors of its predecessors, gradually improving the overall predictive performance. XGBoost incorporates regularization techniques to control overfitting and assigns weights to individual data points, emphasizing those with higher errors during training. Its ability to handle missing data and flexibility in handling various data types contribute to its versatility (Chen & Guestrin, 2016). XGBoost has gained widespread recognition due to its success in diverse domains, from finance (Dalal et al., 2022) to healthcare (Zhang et al., 2020) to geothermal energy (Alqahtani et al., 2023).
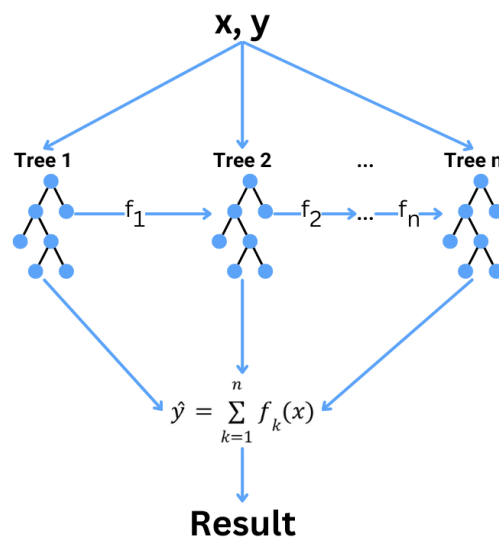


**Figure 2: General architecture of XGBoost.**

PdM often involves numerous sensor measurements, potentially containing redundant or irrelevant information. XGBoost offers feature importance scores, which can help alleviate this challenge. The same feature importance scores and decision tree visualization tools improve the model's interpretability, allowing engineers to understand the factors affecting the predictions, facilitating trust, and enabling the identification of critical degradation indicators. Another common challenge for industry applications is the size of the datasets. XGBoost's distributed learning capabilities make it well-suited for large datasets. Following are examples of studies conducted on pumps and equipment, the basic operation of which is similar to that in geothermal power plants, indicating the potential benefits of using XGBoost for PdM in Geothermal power plants.

Abdalla et al. (2022) attempted PdM for Electrical Submersible Pumps using XGBoost and PCA. They used time-series data from sensors on electrical submersible pumped wells. Their model achieved a mean AUC of 0.95. The model demonstrated a high precision (0.8) but a relatively lower recall (0.6), indicating a small no. of false alarms, but not all events were detected 7 days in advance. They stated the potential cause of this to be the fact that signs of an upcoming workover will not be exhibited on all days before a workover event. Their study could serve as an initial step for further research on improving the model, aiming for an accurate and timely detection of more events.

Salim et al. (2022) coupled XGBoost with the Genetic Algorithm for PdM for compressor 103J and water pump systems. They used the compressor 103J and the water pump dataset from Kaggle (*Pump_Sensor_Data*, n.d.). Their model achieved an accuracy of 99.08%. However, they noted that processing different genetic operations required significant execution time. They suggested carrying out further research to include optimization methods like Cuckoo search and developing parallel versions of the proposed GA-XGBoost algorithm to improve performance.

Xie et al. (2023) compared 4 approaches for predicting the remaining service life of a concrete pump. They explored XGBoost, Support Vector Regression (SVR), BP Neural Networks, and an Ensemble of the three. They utilized IoT monitoring data covering the entire life cycle of a piston in a concrete pump. Of the individual models, XGBoost performed second best with an $R^2$ of 0.9875, only bested by the Ensemble approach with an $R^2$ of 0.9924. Hence, XGBoost could prove to be a useful tool for predicting the RUL of equipment, and more comprehensive research should be carried out to explain the ensemble of XGBoost with other algorithms.

Hence, the properties and features of XGBoost, along with examples of its applications for PdM, point towards it being a robust, scalable algorithm capable of providing an explainable model, making it a viable candidate for PdM. The ensemble of XGBoost with other algorithms should also be explored. However, research on its applications for PdM in geothermal power plants is scant, and more studies are required to derive an objective conclusion.

## CHALLENGES AND FUTURE SCOPE

In pursuing advancing PdM strategies for geothermal power plants, this study encountered a significant hurdle in the form of limited access to sufficient data regarding the operation, maintenance, and failure of equipment in geothermal power plants. The lack of diverse and extensive datasets hampered the possibility of training and proposing a robust model for PdM. The scarcity of relevant data not only restricted the scope of our study but also underscored a broader challenge facing the field.

Another noteworthy challenge, possibly due to the first one, was the paucity of literature on PdM in geothermal power plants. The existing body of work in this domain is noticeably limited, with only a few studies delving into the unique intricacies and challenges associated with geothermal facilities. This scarcity constrained the breadth of our literature review and underscored the pressing need for further research endeavors in PdM for geothermal power plants.

The limited studies in this domain hinder the development and validation of effective models and methodologies. Further research is imperative to advance the field and address the critical needs of geothermal power plants. This scarcity of literature necessitates a collective effort from the energy and computer science communities to bridge the gap between domain-specific knowledge and advanced ML techniques.

The future scope of research in PdM for geothermal power plants involves overcoming data challenges and fostering interdisciplinary collaboration. Establishing partnerships between geothermal energy experts and computer scientists will be instrumental in harnessing the full potential of this natural resource. Combining domain expertise with advanced data-driven methodologies can pave the way for more resilient and efficient geothermal power plants, ultimately contributing to the broader goal of sustainable and reliable energy sources.

## CONCLUSION

In conclusion, autoencoders and XGBoost both hold the potential to be essential tools in the PdM for geothermal power plants. Autoencoders can identify the key features of the data. They are ideal for anomaly detection, a step invaluable for early detection of potential equipment failures. On the other hand, XGBoost's robust and scalable nature positions it as a formidable tool for predicting the RUL of critical components, providing a foundation for more proactive and efficient maintenance practices. XGBoost is a robust and scalable algorithm that can predict the RUL of equipment to better guide maintenance efforts. Both approaches for predictive maintenance can be used to make an explainable model, which would be helpful to point to the root cause of failure and further optimize the maintenance efforts. However, sufficient data regarding equipment operation, maintenance, and failure in geothermal power plants is required to develop, train, and test models. Computer Science and geothermal fields have to collaborate closely to bridge the former's lack of expertise in the latter. However, it is crucial to highlight the prerequisite of substantial data encompassing the operation, maintenance, and failure scenarios in geothermal power plants. The close collaboration of Computer Science and geothermal energy experts is essential to overcoming challenges and establishing a comprehensive framework for developing AI tailored to the unique requirements of the geothermal energy sector, helping us make the most out of the natural resources Mother Earth provides.

## REFERENCES

Abdalla, R., Samara, H., Perozo, N., Carvajal, C. P., & Jaeger, P. (2022, May 19). Machine Learning Approach for Predictive Maintenance of the Electrical Submersible Pumps (ESPs). ACS Omega, 7(21), 17641–17651. https://doi.org/10.1021/acsomega.1c05881

Alamu, O. A., Pandya, D. A., Warner, O., & Debacker, I. (2020, May 4). ESP Data Analytics: Use of Deep Autoencoders for Intelligent Surveillance of Electric Submersible Pumps. Day 1 Mon, May 04, 2020. https://doi.org/10.4043/30468-ms

Alqahtani, F., Ehsan, M., Abdulfarraj, M., Aboud, E., Naseer, Z., El-Masry, N. N., & Abdelwahed, M. F. (2023, August 22). Machine Learning Techniques in Predicting Bottom Hole Temperature and Remote Sensing for Assessment of Geothermal Potential in the Kingdom of Saudi Arabia. Sustainability, 15(17), 12718. https://doi.org/10.3390/su151712718

Chen, T., & Guestrin, C. (2016, August 13). XGBoost. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. https://doi.org/10.1145/2939672.2939785

Dalal, S., Seth, B., Radulescu, M., Secara, C., & Tolea, C. (2022, December 9). Predicting Fraud in Financial Payment Services through Optimized Hyper-Parameter-Tuned XGBoost Model. Mathematics, 10(24), 4679. https://doi.org/10.3390/math10244679

Givnan, S., Chalmers, C., Fergus, P., Ortega, S., & Whalley, T. (2021, October 1). Real-Time Predictive Maintenance using Autoencoder Reconstruction and Anomaly Detection. arXiv preprint. https://doi.org/10.48550/arXiv.2110.01447

Goodfellow, I., Bengio, Y., & Courville, A. (2016, November 18). Deep Learning.

Liu, Y., Ling, W., Young, R., Zia, J., Cladouhos, T. T., & Jafarpour, B. (2022, March 31). Latent-Space Dynamics for Prediction and Fault Detection in Geothermal Power Plant Operations. Energies, 15(7), 2555. https://doi.org/10.3390/en15072555

Hajgató, G., Wéber, R., Szilágyi, B., Tóthpál, B., Gyires-Tóth, B., & Hős, C. (2022, October). PredMaX: Predictive maintenance with explainable deep convolutional autoencoders. Advanced Engineering Informatics, 54, 101778. https://doi.org/10.1016/j.aei.2022.101778

pump_sensor_data. (n.d.). Pump_Sensor_Data | Kaggle. https://www.kaggle.com/datasets/nphantawee/pump-sensor-data

Salim, K., Hebri, R. S. A., & Besma, S. (2022, December 31). Classification Predictive Maintenance Using XGboost with Genetic Algorithm. Revue D'Intelligence Artificielle, 36(6), 833–845. https://doi.org/10.18280/ria.360603

Tian, R., Liboni, L., & Capretz, M. (2022, November 26). Anomaly Detection with Convolutional Autoencoder for Predictive Maintenance. 2022 9th International Conference on Soft Computing & Machine Intelligence (ISCMI). https://doi.org/10.1109/iscmi56532.2022.10068441

van der Zwaan, B., & Dalla Longa, F. (2019, November). Integrated assessment projections for global geothermal energy use. Geothermics, 82, 203–211. https://doi.org/10.1016/j.geothermics.2019.06.008

Xie, X., Meng, X., Du, J., & Peng, Z. (2023, October 17). Predicting Remaining Service Life of Concrete Pump Pistons Using Ensemble Learning. Proceedings of the 7th International Conference on Computer Science and Application Engineering. https://doi.org/10.1145/3627915.3628078

Zhang, X., Yan, C., Gao, C., Malin, B. A., & Chen, Y. (2020, August 3). Predicting Missing Values in Medical Data Via XGBoost Regression. Journal of Healthcare Informatics Research, 4(4), 383–394. https://doi.org/10.1007/s41666-020-00077-1