

Mitigating False Data Injection Attacks Using Machine Learning Models

Jannavarapu Vani Akhila

*Department of Computer Science And Engineering
Vignan's Foundation for Science, Technology and Research
Vadlamudi, Guntur, AP
vaniakhilajannavarapu@gmail.com*

Mondem Manikanta

*Department of Computer Science And Engineering
Vignan's Foundation for Science, Technology and Research
Vadlamudi, Guntur, AP
manikantamondem@gmail.com*

Nagulapati Phanindra Raja Mithra

*Department of Computer Science And Engineering
Vignan's Foundation for Science, Technology and Research
Vadlamudi, Guntur, AP
nagulapatimithra@gmail.com*

Maridu Bhargavi

*Department of Computer Science And Engineering
Vignan's Foundation for Science, Technology and Research
Vadlamudi, Guntur, AP
bhargaviformal@gmail.com*

srigakolapu Sai Lakshmi

*Department of Computer Science And Engineering
Vignan's Foundation for Science, Technology and Research
Vadlamudi, Guntur, AP
sailakshmisrigakolapu@gmail.com*

Abstract—A vehicle that functions purely on electricity, as opposed to normal gasoline or diesel automobiles, is known as an electric vehicle (EV). Electric vehicles are becoming extremely popular. EVs come equipped with a built-in connection that allows for features like remote control via smartphone apps, over-the-air software upgrades, and navigation. Smart charging stations are available for electric vehicles (EVs) and are capable of utilizing energy consumption and cost data to improve charging times through communication with the car. However, there are cyber security problems nowadays everywhere. Monthly pricing is one option for charging at electric charging stations. False Data Injection Attacks are a frequent form of attack in which hackers alter data to their benefit. We are developing a machine learning algorithm to identify the site of the deceptive data injection in order to solve this issue. Random Forest Algorithms, RBF kernel, linear kernel, polynomial kernel, and Support Vector Machine were utilized for training and assessing the outcomes.

Keywords: False Data Injection Attack (FDIA), Electric Vehicle Charging Stations, SVM with RBF Kernel, Machine Learning

I. INTRODUCTION

Cars and other vehicles that run on electricity rather than gasoline or diesel are known as electric vehicles, or EVs. They are growing in popularity since they may lessen pollutants and are healthier for the environment. Electric motors in EVs are powered by batteries, which allow the vehicle to move. In addition, EVs are often smoother and quieter to drive than cars with internal combustion engines. Because they have fewer moving components, they can be more dependable over time and require less maintenance. In addition, as EVs may be charged at the comfort of home, a normal electrical socket, or at the community charging stations they are simple to use

on a regular schedule. But there are certain drawbacks to electric vehicles as well. One issue is that certain electric cars have a limited range, which means they can only go so far before needing to be recharged. Long excursions or locations with a lack of charging infrastructure may have this problem. Both bodily harm and financial damages may result from such attacks. The price of EVs, which might be more up front than that of conventional automobiles, is another factor to take into account. Nevertheless, over time, fuel and maintenance savings frequently outweigh this expense. In addition, costs are anticipated to decrease as EV adoption increases and technology progresses. Due to their many technological features, electric cars are susceptible to hackers. [1]

Consider an EV charging station located in a gated neighborhood. To utilize this station for charging their EVs, users must pay a monthly charge. The charging station keeps track of data, such as the quantity of power consumed by each user and the accompanying payments. Unauthorized access to the charging station's data system is obtained by an attacker. This may occur as a result of software flaws, weak passwords, or other security holes. The objective of the attacker is to use data manipulation to their advantage covertly. The charge data that is kept in the system is altered by the attacker.

For instance, they can make it seem as though they consumed more power than they actually did by increasing the documented energy consumption for their own EV. Alternatively, they might limit the consumption data for other customers, lowering their prices. By inserting fraudulent data, the attacker can profit monetarily or disrupt the charging station's

operations. The False Data Injection Attack caused numerous implications, including income loss for the charging station owner owing to inaccurate invoicing. The station's efficiency has been weakened, impacting overall operations. Alternatively, they might reduce usage data for other customers, lowering their fees. By inserting fraudulent data, the attacker can profit monetarily or disrupt the charging station's operations.

To address this issue, we proposed applying a machine learning approach called False Data Injection to identify the cyberattack. Our major objective is to detect cyber attacks with Machine Learning models. Our approach was implemented using the Electric Vehicle Charging Dataset from Kaggle. Techniques for machine learning, such support vector machines (SVM), can detect anomalous charging patterns in data. These algorithms use previous data to identify potentially suspicious transactions. Regular audits, secure authentication, and encryption are crucial for preventing such crimes. False Data Injection Attacks modify charge data for personal advantage, potentially resulting in financial losses and operational problems. Detecting and blocking such attacks is critical for the safety of EV charging stations.

II. LITERATURE REVIEW

Dhaou et al. [1] This research study discusses cyber assaults on connected electric vehicles. The study focused on False Data Injection Attacks that may occur in parking lots where connected electric vehicles make energy transactions. The Effects of Injecting False Data cyberattack on pricing and power signal anomaly is the main subject of this study. The solution involved combining blockchain technology and machine learning algorithms. The approach being proposed aims to securely and efficiently transfer energy across connected electric vehicles.

Hongyang Li1 et al. [2] In this research paper, the author analyzes the advancement of technology in numerous disciplines and how industries are adapting to contemporary innovations. This study focuses on protecting digital systems from various sorts of cyber threats. Fault diagnosis is crucial for ensuring safe and reliable manufacturing. This research focuses on model-based fault diagnosis, which yields more accurate findings. The suggested approach involves analyzing data using a three-tank system. The author efficiently identifies crimes using a particle filter.

K Jonath Kwizera et al. [3] This study explains how to identify cyber threats using various data mining approaches. Cyber attacks can expose sensitive information, causing several challenges. Cybersecurity protects against cyber threats and keeps data safe. To combat more advanced and challenging threats, more efficient technologies are required. Datamining approaches were employed in this article to address the cyberattacks and problems that are prevalent worldwide.

Gaoqi Liang et al. [4] This research paper discusses the vital importance of power systems. Modern grid systems provide power for the necessities of life, which mainly depend on computer systems. To ensure a seamless and safe procedure, we must examine the security of the computing systems. The power security system consists of two components. Cyber

security protects many systems and identifies attacks that could occur. In addition, physical security describes how they operate while encountering various challenges. State estimation is an important feature in power control systems. It displays system status and enhances SCADA measures. Managing the electricity systems we use will be extremely challenging if the attacker targets the SCADA system.

Md. Ashfaqur Rahman et al. [5] The detection of fraudulent data injections that impact on nonlinear estimates in power systems is examined in this article. The author emphasizes the importance of addressing weaknesses, particularly in power systems. With this kind of assault, the culprit manipulates the data of the system as it leads that operator of the system may changes the settings in the power systems. The attackers inject the data into the system without being caught by the operators of the system. This paper focuses on the issue non linear state estimation. This will give assumption how the power system actually works. It resolves the problem step by step, monitoring on both active and reactive workflow. This process is common in every power system industry but they did not focus on the risk of False data injection attack.

III. METHODOLOGY

We propose a method to identify and detect false data injection in the datasets of electric vehicle charging stations.

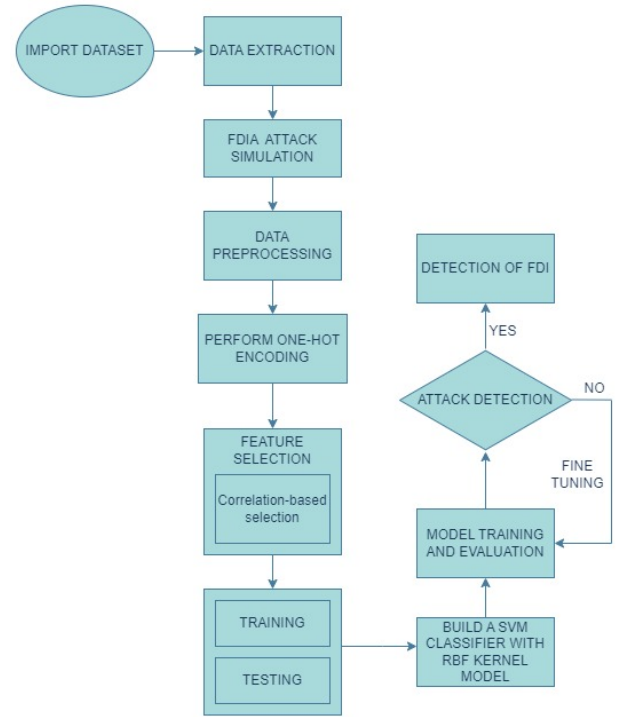


Fig. 1. Proposed Model diagram

The shown proposed algorithm (1) uses an SVM model and then incorporates an RBF kernel to detect false data injection in the data related to electric vehicle charging

station, calculating performance and saving for future use.

Algorithm 1: Detecting False Data Injection Using SVM Model

Input: Dataset $D = \{(x_i, y_i)\}_{i=1}^n$, where x_i are features and $y_i \in \{0, 1\}$ is the label (0 = legitimate, 1 = fraud).

Output: Predicted false data injections, accuracy, precision, and saved model.

Step 1: Load and Prepare Dataset Load the dataset 'station_data_dataverse.csv'. Add a column 'FalseDataInjection' with random binary values (0 or 1). Drop rows with null values. Perform one-hot encoding on the 'Week-day' feature. Features selected: {kwhTotal, dollars, chargeTimeHrs, distance, Mon, Tues, Wed, Thurs, Fri, Sat, Sun, managerVehicle}. **Step 2: Split dataset into training (80%) and testing (20%) sets.** Divide dataset D into training and testing sets:

$$D_{\text{train}}, D_{\text{test}} = \text{split}(D, 0.8, 0.2)$$

Step 3: Train SVM Classifier Train SVM classifier with RBF kernel on D_{train} :

$$\text{clf} = \text{SVM}(\text{kernel}='rbf')$$

Fit the model on $X_{\text{train}}, y_{\text{train}}$:

$$\text{clf.fit}(X_{\text{train}}, y_{\text{train}})$$

Step 4: Make Predictions and Evaluate Model Predict on D_{test} :

$$y_{\text{pred}} = \text{clf.predict}(X_{\text{test}})$$

Generate confusion matrix:

$$\text{cm} = \text{confusion_matrix}(y_{\text{test}}, y_{\text{pred}})$$

Extract values from the confusion matrix:

$$T1 = \text{cm}_{00}, \quad T0 = \text{cm}_{11}, \quad F1 = \text{cm}_{01}, \quad F0 = \text{cm}_{10}$$

Compute accuracy:

$$\text{Accuracy} = \frac{T1 + T0}{T1 + T0 + F1 + F0}$$

Compute precision:

$$\text{Precision} = \frac{T1}{T1 + F1}$$

Step 5: Save Model Save the trained model to a file:

$$\text{pickle.dump}(\text{clf}, \text{open}('svm_model.pkl', 'wb'))$$

Return: Predicted results, accuracy, and precision.

A. Data Collection

Data from the Electric Vehicle Charging Dataset for FDIA identification, which was obtained through Kaggle, must be gathered in the first step. A new column, FalseDataInjection, is created with random binary values to replicate both valid (0) and fraudulent (1) data. This innovation helps to emulate real-world settings in which both forms of data exist, allowing the model to be trained more effectively. While investigating other datasets, a manual analysis found that many of them were either incomplete, of poor quality, or contained irrelevant data.

B. Data Preprocessing

We preprocessed the dataset by removing null or empty cells, added further feature cells, tested for any null or empty values, and converted the text to numerical representation using one-hot encoding, giving it numerical binary features.

C. Train/Test Split

The 3,395-entry The dataset is divided into two sections: One is utilized for training models, while the other is employed for testing them. The test set, which is not seen during training, is used to assess how well the model can predict current information, whereas the training set enhances the model's capacity to find patterns in the data.

D. Predicted Machine Learning Path

To create our model, we use the dataset and apply the Support Vector Machine (SVM) algorithm with a Radial Basis Function (RBF) kernel. SVM is a reliable and productive machine learning tool, especially for classification tasks, where we need to sort data into categories. The RBF kernel is particularly helpful because it can handle complex, non-linear patterns in the data. It accomplishes this by projecting the data onto a place with more dimensions, making it easier for the model to draw more accurate boundaries between different categories, which helps improve the overall classification accuracy.

I. Support Vector Machine

A support vector machine finds the hyperplane best differentiating classes in the high-dimensional space, using methods under supervised learning that can be applied both for regression and classification problems. Its objective is to optimize the margin between classes so as to increase the classification confidence of fresh data points. It is flexible, as different kernel functions can be applied to map the data into higher dimensions as needed.

$$\min_{x,y} \frac{1}{2} \|x\|^2 + D \sum_{i=1}^m \eta_i \quad (1)$$

subject to:

$$y_i(x \cdot x_i + y) \geq 1 - \eta_i, \quad \forall i \quad (2)$$

II. Radial Basis Function Support Vector Machine (RBF SVM)

Regression and classification issues may be effectively handled using the extremely potent machine learning method known as Radial Basis Function Support Vector Machine. This is a powerful non-parametric model that performs well on non-linear and high-dimensional data. It projects the input

data onto a higher-dimensional feature space in which the classes can be divided using a hyperplane. The technique computes similarity measurements between pairs of data points in the feature space using a kernel function, such as the Radial Basis Function.

An often utilized kernel function with RBF SVM is the Radial Basis Function. It is characterized as:

$$K(a, b) = \exp(-\delta \|a - b\|^2) \quad (3)$$

Here δ denotes the strength of the sampling effect of a training example, and a and b denote two input data points. Larger δ makes a pretty good fit to the data but results in broader decision boundaries compared to their smaller counterparts, which induce a very powerful RBF kernel and can sometimes become computationally exorbitant for RBF transitions over larger datasets. When comparing two data points, the kernel function calculates how similar they are in terms of the distance or relationship that exists in the high-dimensional feature space.

Linear Kernel The linear kernel is another often utilized kernel. For data that can be split by a straight line (or hyperplane in higher dimensions), this kernel is straightforward and efficient.

$$K(a_i, b_j) = a_i \cdot b_j \quad (4)$$

Polynomial Kernel The polynomial kernel allows the model to learn polynomial decision boundaries.

$$K(a_i, b_j) = (a_i \cdot b_j + c)^d \quad (5)$$

d is a degree of polynomial and c is a constant. Compared with a linear kernel, this kernel can handle more complex separation of data; However, the parameters must be carefully adjusted to avoid overfitting, especially when working with high polynomial degrees.

III. Random Forest

An ensemble learning method called Random Forest is applied to regression and classification problems. During the training process, It generates many decision trees using various random subsets of the given data. Random Forest combines the outputs of these distinct trees to improve prediction accuracy and resilience. In the case of classification problems, majority voting selects the mode of the trees as the final output. For regression problems, a single continuous output is generated by averaging the predictions made by each tree. This technique dramatically lowers the chance of overfitting, especially for large, multidimensional datasets considered to be noisy.

Random Forest can be mathematically represented as

$$\tilde{y} = \frac{1}{T} \sum_{t=1}^T g_t(a) \quad (6)$$

here $g_t(a)$ is the forecast provided by the t -th decision tree, T is the forest's total number of trees, and \tilde{y} is the expected output.. In classification tasks, the prediction is determined by majority voting:

$$\tilde{y} = \text{mode}(g_1(a), g_2(a), \dots, g_T(a)) \quad (7)$$

E. Metrics

Metrics are quantitative metrics utilized to assess the performance of machine learning models, providing information about their accuracy, precision, recall, and other important characteristics. These measures aid in determining how well the model predicts and its performance in various categorization tasks.

I. Accuracy

Accuracy is a metric for how often the model properly identifies occurrences in the dataset. It is determined as the ratio of accurately predicted cases (positive and negative) to total instances. The model's overall performance is given by accuracy.

$$\text{Accuracy} = \frac{T1 + T0}{T1 + T0 + F1 + F0} \quad (8)$$

Here $T1$ stand for the number of true positives, $T0$ for the number of true negatives, $F1$ for the number of false positives, and $F0$ for the number of false negatives.

II. Precision

Precision focuses specifically on how well the optimistic predictions made by the model. It illustrates the percentage of actual positives among all instances predicted as positive, providing insight into how reliable those positive predictions are. When precision is high, it means that when the model indicates a positive class, it is usually correct. This is especially crucial in scenarios where false positives can lead to disastrous results, such as fraud detection or medical diagnosis, where precision becomes a vital metric to be monitored.

$$\text{Precision} = \frac{T1}{T1 + F1} \quad (9)$$

where $F1$ stands for false positives and $T1$ for true positives.

F. Detecting FDIA

We test the model to determine its effectiveness in identifying fake data injection attacks (FDIA) once it has been trained. The model produces a result for each data point: a 1 indicates that an attack has happened due to the injection of false data, whereas a 0 indicates that everything is normal and there has been no attack. This allows us to detect any security concerns in the information gathered from the charging stations for electric vehicles. The model's accuracy in making these detections is determined by how well it learns during the training process.

IV. EXPERIMENTAL RESULTS AND DISCUSSION:

Data was preprocessed, key features were extracted, and subsets for training and testing were created in order to evaluate the model. In order to assess how well each kernel could classify the data, three kernel functions—the polynomial, linear, and radial basis function (RBF)—were used to test out SVM models with their performance criteria of accuracy and precision.

```
print(FDI.head())
```

	sessionId	kwhTotal	dollars	created	ended
0	1366563	7.78	0.00	0014-11-18 15:40:26	0014-11-18 17:11:04
1	3075723	9.74	0.00	0014-11-19 17:40:26	0014-11-19 19:51:04
2	4228788	6.76	0.58	0014-11-21 12:05:46	0014-11-21 16:46:04
3	3173284	6.17	0.00	0014-12-03 19:16:12	0014-12-03 21:02:18
4	3266500	0.93	0.00	0014-12-11 20:56:11	0014-12-11 21:14:06

	startTime	endTime	chargeTimeHrs	weekday	platform	...	managerVehicle
0	15	17	1.510556	Tue	android	...	0
1	17	19	2.177222	Wed	android	...	0
2	12	16	4.671667	Fri	android	...	0
3	19	21	1.768333	Wed	android	...	0
4	20	21	0.298611	Thu	android	...	0

	facilityType	Mon	Tues	Wed	Thurs	Fri	Sat	Sun	reportedZip
0	3	0	1	0	0	0	0	0	0
1	3	0	0	1	0	0	0	0	0
2	3	0	0	0	0	1	0	0	0
3	3	0	0	1	0	0	0	0	0
4	3	0	0	0	1	0	0	0	0

Fig. 2. DATASET

Figure 2 describes the dataset used to evaluate the model obtained from Kaggle.

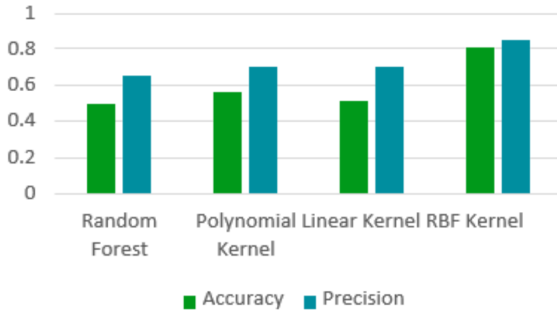


Fig. 3. Accuracy and Precision of various kernels in SVM

The Figure-3 demonstrates the random forest classifier's accuracy.

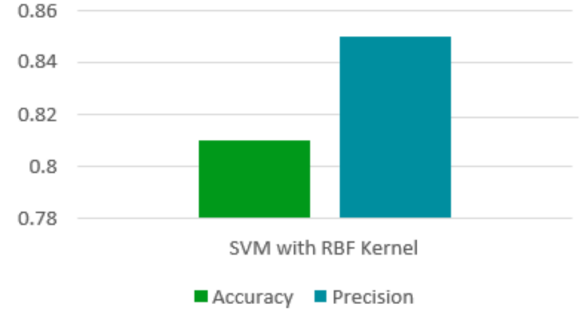


Fig. 4. Accuracy and Precision of the SVM with RBF Kernel

The Figure-4 describes the precision and accuracy of SVM with RBF Kernel.

```
In [9]: print(y_pred)
result= x_session.assign(predict=y_pred)
false_detect=result.loc[result['predict']==1]
print(len(false_detect))
print(false_detect)
```

```
[0 0 0 ... 0 0 0]
5
```

	sessionId	predict
303	6962786	1
1567	8094323	1
1909	9185227	1
2723	6470613	1
2774	4087293	1

Fig. 5. Final Results

Figure-5 represents the final results of the model that 1 represents that FDIA is detected.

V. CONCLUSION

Our research's objective is to improve the detection of False Data Injection Attacks in the dataset of electric vehicle charging stations. We proposed a model with an accuracy of 81%, which stands out from all prior instances that have only achieved an accuracy of approximately 60%, in contrast to existing systems that achieved minimal gain in accuracy.

In our support vector machine model, we tried different kernels like RBF, linear, polynomial, and random forest algorithms. However, we realized that SVM with an RBF kernel is working best to identify False Data which was inserted into the dataset at a high accuracy level.

The SVM model is developed to classify instances into two classes depending on whether there is false data injection (1) or not (0) based on various features: energy consumption ('kwhTotal'), charge time ('chargeTimeHrs'), day of the week, and others.

REFERENCES

- [1] D. Said, M. Elloumi and L. Khokhi, "Cyber-Attack on P2P Energy Transaction Between Connected Electric Vehicles: A False Data Injection Detection Based Machine Learning Model," in IEEE Access, vol. 10, pp. 63640-63647, 2022, doi: 10.1109/ACCESS.2022.3182689.

- [2] E. Drayer and T. Routtenberg, "Detection of False Data Injection Attacks in Power Systems with Graph Fourier Transform," 2018 IEEE Global Conference on Signal and Information Processing (Global-SIP), Anaheim, CA, USA, 2018, pp. 890-894, doi: 10.1109/Global-SIP.2018.8646454.
- [3] Abdelfatah, Ahmed and Ali, Abdelfatah and Shaaban, Mostafa and Osman, Ahmed. (2023). Optimal False Data Injection Attack on EV Chargers and DGs in Active Distribution Networks. 1-6. 10.1109/ICEC-CME57830.2023.10253138.
- [4] Y. Liu, O. Ardakanian, I. Nikolaidis and H. Liang, "False Data Injection Attacks on Smart Grid Voltage Regulation With Stochastic Communication Model," in IEEE Transactions on Industrial Informatics, vol. 19, no. 5, pp. 7122-7132, May 2023, doi: 10.1109/TII.2022.3209287.
- [5] A. Kumar, N. Saxena and B. J. Choi, "Machine Learning Algorithm for Detection of False Data Injection Attack in Power System," 2021 International Conference on Information Networking (ICOIN), Jeju Island, Korea (South), 2021, pp. 385-390, doi: 10.1109/ICOIN50884.2021.9333913.
- [6] R. Nawaz, M. A. Shahid, I. M. Qureshi and M. H. Mehmood, "Machine learning based false data injection in smart grid," 2018 1st International Conference on Power, Energy and Smart Grid (ICPESG), Mirpur Azad Kashmir, Pakistan, 2018, pp. 1-6, doi: 10.1109/ICPESG.2018.8384510.
- [7] P. L. Bhattar, N. M. Pindoriya, A. Sharma and R. T. Naayagi, "False Data Injection Attack Detection with Feedforward Neural Network in Electric Vehicle Aggregator Bidding Price," 2022 IEEE PES Innovative Smart Grid Technologies - Asia (ISGT Asia), Singapore, Singapore, 2022, pp. 665-669, doi: 10.1109/ISGTAsia54193.2022.10003474.
- [8] D. Said and M. Elloumi, "A New False Data Injection Detection Protocol based Machine Learning for P2P Energy Transaction between CEVs," 2022 IEEE International Conference on Electrical Sciences and Technologies in Maghreb (CISTEM), Tunis, Tunisia, 2022, pp. 1-5, doi: 10.1109/CISTEM55808.2022.10044067.
- [9] A. Kumar, N. Saxena and B. J. Choi, "Machine Learning Algorithm for Detection of False Data Injection Attack in Power System," 2021 International Conference on Information Networking (ICOIN), Jeju Island, Korea (South), 2021, pp. 385-390, doi: 10.1109/ICOIN50884.2021.9333913.
- [10] X. Tong and W. Qi, "False Data Injection Attack on Power System Data-Driven Methods Based on Generative Adversarial Networks," 2021 IEEE Sustainable Power and Energy Conference (iSPEC), Nanjing, China, 2021, pp. 4250-4254, doi: 10.1109/iSPEC53008.2021.9735442.