

# **Using Machine Learning to Predict Health Insurance Fraud - A Literature Review**

Ben Zapka

12695179

Unit 7

Research Methods and Professional Practice

University of Essex Online

16 March 2025

## Table of Contents

1. Introduction.....	3
2. Critical Review of Existing Literature .....	4
3. Conclusion .....	11
4. List of References.....	12

## **1. Introduction:**

Health insurance fraud is the intentional act of deceiving or misrepresenting information to obtain unauthorised benefits or payments from health insurance providers (Liu et al., 2023; Villegas-Ortega et al., 2021). While it is no new phenomenon and companies spend large amounts yearly to prevent it, the numbers are still increasing (Insurance Europe, 2024; Association of Retired Persons, 2023; Galan, 2024). Companies increasingly use technology for prevention (Guo, 2024; Wang et al., 2025). Recent contributions in the academic literature have shown that machine learning algorithms present a superior way of predicting health insurance fraud compared to traditional, rule-based prediction algorithms (Cherkaoui, 2024; Sharma, 2024). This paper analyses how machine learning is used to predict health insurance fraud. A broad search on Google Scholar, including backward searches from key literature reviews, yielded a wide range of studies. It was ensured to only include studies from 2016 and after, i.e., from the last ten years and focusing on works from the last five years to include only the most recent research. The methods found were grouped into three categories: supervised learning, unsupervised learning, and ensemble methods. Each category is examined, highlighting strengths and weaknesses, to provide companies with a comprehensive overview of effective approaches to detect health insurance fraud.

## **2. Critical Review of Existing Literature:**

Supervised learning is a machine learning approach where models are trained on labelled datasets, learning to map inputs to known outputs. This technique enables models to make accurate predictions or classifications on new, unseen data (Jiang et al., 2020; Ono and Goto, 2022). Supervised learning is mainly used to classify claims as fraudulent or legitimate, detecting only already known fraudulent patterns. Algorithms used for this include gradient boosted trees, k-nearest neighbour, naïve Bayes and neural networks.

Gradient boosted trees (GBT) are highly effective in modelling non-linear relationships and can be used to classify health insurance claims with regards to fraud. They manage to uncover complex fraud patterns while staying computationally efficient, not requiring immense resources and creating an interpretable result (Hancock and Khoshgoftaar, 2021; Vandervorst et al., 2024). However, Gupta et al. (2021) found that GBT can be costly due to a relatively slow training process and they require careful hyperparameter tuning, resulting in relatively large resource requirements for labour and computation, contrasting the findings from the two other studies. Furthermore, using GBT risks overfitting if the hyperparameters are not tuned properly, e.g., by using cross-validation and early stopping (Parthasarathy et al., 2023).

K-Nearest neighbour (KNN) is non-parametric, can adapt to complex data distributions without requiring a predefined model and is straightforward to implement

and interpret (Bauder and Khoshgoftaar, 2018; Ekin et al., 2021). Yet, its reliance on distance metrics leads to poor performance in high-dimensional spaces or imbalanced datasets and it is relatively computationally intensive, especially for large datasets (Ekin et al., 2021; Lalithagayatri et al., 2018). Additionally, Agarwal (2023) emphasises the need for robust and adaptive models in fraud detection, suggesting that simpler models like KNN might not be sufficient for complex fraud patterns.

Naïve Bayes is computationally efficient, rather simple and easy to interpret (Bauder and Khoshgoftaar, 2018; Sadiq and Shyu, 2019). While these strengths prevail, Mambo and Moturi (2022) note that its performance can be limited by the complexity of real-world data and stress the necessity for robust feature selection. Bauder and Khoshgoftaar (2018) show that Naïve Bayes struggles with imbalanced datasets, as dependency patterns in minority classes are often insignificant, leading to misclassification and a classifier that achieved an area under the curve score of 0.63 at the maximum, showcasing a mediocre performance.

Neural Networks are highly beneficial in predicting health insurance fraud due to their ability to process complex, high-dimensional data, capture intricate patterns, and adapt to imbalanced datasets through techniques like random oversampling and algorithm-level adjustments, achieving high accuracy and robust performance in diverse fraud detection scenarios (Shamitha and Ilango, 2020; Johnson and Khoshgoftaar, 2019). For instance, Nabrawi and Alanazi (2023) utilised neural networks to classify fraudulent claims, emphasising the importance of robust feature selection and data preprocessing to improve model performance. Johnson and

Khoshgoftaar (2019) demonstrated that deep learning methods, including Neural Networks, can effectively handle class-imbalanced Medicare fraud detection tasks, achieving strong performance metrics using a combination of random oversampling and undersampling techniques. However, the use of Neural Networks also has weaknesses. Their computational complexity and the need for large, high-quality datasets present limitations (Nabrawi and Alanazi, 2023; Johnson and Khoshgoftaar, 2019). Additionally, Neural Networks' lack of interpretability limits their use in health insurance fraud, where companies must understand how predictions are made to understand fraud patterns and guide decisions (Wang et al., 2025).

Unsupervised learning is a machine learning technique that identifies patterns and structures in unlabelled data without predefined outputs. It aims to discover hidden relationships or groupings within datasets, typically using clustering methods (Glielmo et al., 2021; Weidel et al., 2021). Unsupervised learning techniques are mainly used to uncover anomalies in data that might be related to fraudulent activities. These techniques include anomaly detection via clustering algorithms like K-Means Clustering, Local Outlier Factor, One Class Support Vector Machines and Isolation Forest, which is an unsupervised version of the random forest algorithm. They offer several strengths in detecting health insurance fraud, particularly in identifying novel patterns and anomalies without requiring labelled data, contrasting the methodology of supervised learning where only already known fraudulent patterns can be analysed.

K-Means Clustering groups similar entities and identifies outliers indicative of fraudulent behaviour. Its simplicity and efficiency make it particularly suitable for large datasets, as it can quickly partition data into distinct clusters based on Euclidean distances (Massi et al., 2020; Zhang et al., 2020). Despite its strengths, K-Means clustering has notable limitations. Its reliance on Euclidean distance makes it less effective for high-dimensional or non-linear data, where relationships between variables may not be linear (Khan et al., 2024). Furthermore, the algorithm requires the number of clusters to be predefined, which can lead to suboptimal results if the chosen number does not align with the natural structure of the data (Fashoto et al., 2016). Also, K-Means Clustering struggles with imbalanced datasets (Massi et al., 2020). Thus, in comparative studies, K-Means Clustering has been shown to perform worse than other algorithms, especially Isolation Forest (Zhang et al., 2020; Kaur et al., 2021).

Local Outlier Factor calculates the local density deviation of a given data point with respect to its neighbours, making it effective for detecting outliers in datasets with varying densities and particularly useful in hierarchical datasets like health insurance claims, resulting in stronger performance than K Means Clustering or traditional rule-based fraud detection systems (Gao et al., 2018; Zhang et al., 2020). However, the computational complexity can be a limitation when applied to large-scale datasets (Cherkaoui et al., 2024). Furthermore, the results are hard to interpret and there is no clear definition which values for the local outlier factor should be considered anomalous (Bauder and Khoshgoftaar, 2017).

One-Class Support Vector Machines construct a boundary around normal data points in a high-dimensional space using a kernel function and flags points outside this boundary as anomalies which makes them effective in handling high-dimensional data and detecting anomalies without requiring labelled fraudulent examples (Sun et al., 2018; Zhao et al., 2021). Drawbacks are that the performance is sensitive to the choice of kernel and parameters and the model struggles with imbalanced data sets, thus still needing human validation of the outputs which counteracts the general idea of unsupervised learning (Li and Zhang, 2018; Bauder and Khoshgoftaar, 2018).

Isolation Forest isolates outliers by randomly partitioning the data and identifies fraudulent claims by leveraging its ability to isolate anomalies based on fewer splits in a tree structure (Zhang et al., 2020; Rosa and Khoshgoftaar, 2018). While Isolation Forest does a relatively good job in finding fraudulent claims, it suffers from a relatively high false positive rate, meaning that it declares relatively many legitimate claims as fraudulent which results in a high workload for human reviewers (Das et al., 2017; Rosa and Khoshgoftaar, 2018). Furthermore, although Isolation Forest is generally efficient with large datasets, it can be challenging to tune its parameters, such as the number of trees and the subsampling size, which might affect its performance in complex healthcare data (Naidoo and Marivate, 2020). Also, Tabassum et al. (2024) showed that Isolation Forest is effective in detecting anomalies but might not perform optimally with very high-dimensional data without proper feature engineering.



As supervised and unsupervised learning techniques have different strengths and weaknesses, merging several models together could be promising. In an aim to profit from different strengths of both supervised and unsupervised learning, ensemble methods combine predictions from multiple base models to enhance accuracy, reduce bias and variance, and improve generalisability compared to individual models (Rane et al., 2024; Naderalvojud and Hernandez-Boussard, 2024). The ensemble techniques employed include Voting, Weighted, and Stacking methods. The Weighted method assigns different weights to model predictions based on their performance, ensuring a balanced contribution from each model (Patil et al., 2023; Ibraigheeth et al., 2024). The Stacking approach uses base models to generate predictions, which are then integrated by a meta-model to improve overall accuracy (Zhou and Jiao, 2022; Herianto et al., 2024). Voting aggregates predictions either through majority voting or averaging probabilities (Cui et al., 2023; Rojarath and Songpan, 2020).

Although a detailed comparison of individual ensemble methods is limited by scarce literature, a general assessment shows that ensemble methods enhance health insurance fraud detection as they by combining models improve accuracy and robustness, reducing bias and variance while effectively handling complex data. (Wang et al., 2025; Kotekani and Ilango, 2022; Kunickaitė et al., 2020). However, weaknesses include increased computational complexity and the need for careful model selection and weighting to avoid overfitting (Chaurasiya and Jain, 2025). Additionally, the results are often difficult to interpret as there is no exact specification of which features led to the classification of a given claim as fraudulent or legitimate (Kunickaite et al., 2020; Ekin et al., 2023). Furthermore, also ensemble methods

struggle with imbalanced datasets, which are common in healthcare fraud detection and still require a dedicated strategy to manage imbalanced data sets (Seshagiri and Prema, 2025).

After the evaluation of supervised, unsupervised, and ensemble methods for predicting health insurance fraud, several key gaps remain. Notably, there is a lack of rigorous comparisons between supervised and unsupervised approaches, which is essential for understanding their respective strengths in fraud detection. Although most studies report model performance metrics, they often use different datasets and varying definitions of health insurance fraud, making direct comparisons difficult. Comparisons of several methods on the same data set are limited and do not include all methods. A comprehensive study comparing a large range of supervised, unsupervised, and ensemble methods on the same dataset, with a consistent definition of fraud, would be highly valuable. In general, more studies considering ensemble methods would be beneficial as literature in this area is limited and the models seem to be promising. Additionally, the real-world applicability of these models is rarely tested. Beyond evaluating models on holdout data, future research should involve deploying models in partnership with insurance companies to assess performance in operational settings. This would enhance the generalisability and practical relevance of the findings. Moreover, there is limited research on integrating explainable AI into fraud detection systems, a crucial aspect for real-world adoption. Without transparency, the utility of such models is constrained, as companies must justify their predictions to meet auditing and regulatory requirements (Karangara et al., 2024; Farbmacher et al., 2022). Understanding the factors that drive fraud would also help companies prevent it.

### **3. Conclusion:**

In conclusion, the current body of literature provides a comprehensive yet fragmented understanding of the application of supervised, unsupervised, and ensemble learning techniques in the detection of health insurance fraud. Each approach exhibits distinct advantages and limitations. Supervised learning methods offer high accuracy but face challenges with overfitting and computational demands. Unsupervised methods excel in identifying novel anomalies but struggle with high-dimensional and imbalanced data. Ensemble methods demonstrate the potential to improve predictive performance by combining model strengths but introduce added complexity. A key challenge for all methods is class imbalance, as fraud occurs in only a small fraction of cases. This must be considered when building a predictive model. Still, significant gaps persist within the research landscape. The lack of standardised comparative studies across algorithms on uniform datasets limits the ability to make conclusive statements regarding their relative effectiveness. Also, ensemble methods should be further evaluated. Furthermore, the absence of real-world evaluations impedes understanding of how these models perform within the operational frameworks of insurance companies. Additionally, limited attention has been given to integrating explainable AI techniques, an essential factor for practical deployment where transparency and compliance with regulatory frameworks are critical. To advance the field, future research should prioritise systematic comparisons using consistent datasets and definitions of health insurance fraud, emphasise real-world implementation to validate generalisability, and explore explainable AI frameworks to enhance the transparency and trustworthiness of predictive models. Addressing these gaps will bridge the divide between academic insights and industry application, ultimately strengthening fraud detection efforts in the health insurance sector.

#### **4. List of References:**

Agarwal, S. (2023) 'An Intelligent Machine Learning Approach for Fraud Detection in Medical Claim Insurance: A Comprehensive Study', *Scholars Journal of Engineering and Technology*. Available at:

[https://saspublishers.com/media/articles/SJET\\_119\\_191-200\\_FT.pdf](https://saspublishers.com/media/articles/SJET_119_191-200_FT.pdf)

American Association of Retired Persons (2023) Medicare Fraud. Available at:

<https://www.aarp.org/money/scams-fraud/medicare/> (Accessed 15 March 2025)

Bauder, R. A. and Khoshgoftaar, T. M. (2018) 'The Detection of Medicare Fraud Using Machine Learning Methods with Excluded Provider Labels', *Proceedings of the Thirty-First International Florida Artificial Intelligence Research Society Conference (FLAIRS 2018)*. Available at: <https://aaai.org/papers/404-flairs-2018-17617/>

Bauder, R. A. and Khoshgoftaar, T. M. (2017) 'Multivariate outlier detection in medicare claims payments applying probabilistic programming methods' *Health Services and Outcomes Research Methodology* 17(1), pp. 256-289. Available at: <https://link.springer.com/article/10.1007/s10742-017-0172-1>

Chaurasiya, R. & Jain, K. (2025) 'Healthcare Fraud Detection Using Machine Learning Ensemble Methods', *South Eastern European Journal of Public Health* 26(1), pp. 4789-4795. Available at: <https://seejph.com/index.php/seejph/article/view/4988>

Cherkaoui, O., Anoun, H. & Maizate, A. (2024) 'A benchmark of health insurance fraud detection using machine learning techniques', *International Journal of Artificial Intelligence* 13(2), pp. 1925-1934. Available at: <http://doi.org/10.11591/ijai.v13.i2.pp1925-1934>

Das, S., Wong, W.-K., Fern, A., Dietterich, T. G. & Siddiqui, M. A. (2017) 'Incorporating Feedback into Tree-based Anomaly Detection', *arXiv*. Available at: <https://arxiv.org/abs/1708.09441>

Das, A. (2024) 'Healthcare Fraud Detection Using Machine Learning', *International Journal of Creative Research Thoughts* 12(4), pp. 491-494. Available at: <https://www.ijcrt.org/papers/IJCRT2404517.pdf>

Ekin, T., Frigau, L. & Conversano, C. (2021) 'Health Care Fraud Classifiers in Practice', *Applied Stochastic Models in Business and Industry* 37(6), pp. 1182-1199. Available at: <https://doi.org/10.1002/asmb.2633>

Ekin, S., Han, Y., Duan, Y., Li, Y., Zhu, S. & Song, C. (2023) 'A Two-Stage Voting-Boosting Technique for Ensemble Learning in Social Network Sentiment Classification', *Entropy* 25(4). Available at: <https://doi.org/10.3390/e25040555>

Farbmacher, H., Löw, L. & Spindler, M. (2022) 'An explainable attention network for fraud detection in claims management', *Journal of Econometrics* 228(2), pp. 244-258. Available at: <https://doi.org/10.1016/j.jeconom.2020.05.021>

Fashoto, S. G., Adekoya, A., Gbadeyan, J. A., Sadiku, J. S. & Yahya, W. B. (2016) 'Development of Improved K-Means Clustering to Partition Health Insurance Claims', *GESJ: Computer Science and Telecommunications* 1(47). Available at: [https://www.academia.edu/79358761/Development\\_of\\_improved\\_K\\_means\\_clustering\\_for\\_health\\_insurance\\_claims](https://www.academia.edu/79358761/Development_of_improved_K_means_clustering_for_health_insurance_claims)

Galan, S. (2024) 'Total value of health insurance frauds in France 2005-2018', *Statista*. Available at: <https://www.statista.com/statistics/1149217/global-amount-health-insurance-fraud-france/> (Accessed 15 March 2025)

Gao, Y., Sun, C., Li, R., Li, Q., Cui, L. & Gong, B. (2018) 'An Efficient Fraud Identification Method Combining Manifold Learning and Outliers Detection in Mobile Healthcare Services', *IEEE Access* 6. Available at: <https://ieeexplore.ieee.org/document/8489846>

Glielmo, A., Husic, B. E., Rodriguez, A., Clementi, C., Noé, F. & Laio, A. (2021) 'Unsupervised learning methods for molecular simulation data', *Chemical Reviews* 121(16), pp. 9722-9758. Available at: <https://pubs.acs.org/doi/10.1021/acs.chemrev.0c01195>

Guo, Y. (2024) 'Application of Machine Learning in Insurance Fraud Detection: Achievements and Future Prospects', *Proceedings of the 2024 International Conference on Artificial Intelligence and Communication*. Available at: <https://www.atlantis-press.com/proceedings/icaic-24/126003412>

Gupta, R. Y., Mudigonda, S. S. & Baruah, P. K. (2021) 'A Comparative Study of Using Various Machine Learning and Deep Learning-Based Fraud Detection Models For Universal Health Coverage', *International Journal of Engineering Trends and Technology* 69(3), pp. 96-102. Available at: <https://ijettjournal.org/archive/ijett-v69i3p216>

Hancock, J. T. & Khoshgoftaar, T. M. (2021) 'Gradient Boosted Decision Tree Algorithms for Medicare Fraud Detection', *SN Computer Science* 2. Available at: <https://link.springer.com/article/10.1007/s42979-021-00655-z>

Herianto, Kurniawan, B., Hartomi, Z. H., Irawan, Y. & Anam, M. K. (2024) 'Machine Learning Algorithm Optimization using Stacking Technique for Graduation Prediction', *Journal of Applied Data Sciences* 5(3). Available at: <https://bright-journal.org/Journal/index.php/JADS/article/view/316/226>

Ibraigheeth, M. A., Abu Eid, A. I., Alsariera, Y. A., Awwad, W. F. & Nawaz, M. (2024) 'A New Weighted Ensemble Model to Improve the Performance of Software Project

Failure Prediction', *International Journal of Advanced Computer Science and Applications* 15(2), pp. 359-365. Available at:  
<http://dx.doi.org/10.14569/IJACSA.2024.0150238>

Insurance Europe (2024) 'Insurance Europe Annual Report 2023 - 2024'. Available from: <https://insurancefraud.org/fraud-stats/> (Accessed 15 March 2025)

Jiang, T., Gradus, J. L. & Roselli, A. J. (2020) 'Supervised machine learning: A brief primer', *Behavior Therapy* 51(5), pp. 675-687. Available at:  
<https://doi.org/10.1016/j.beth.2020.05.002>

Johnson, J. M. and Khoshgoftaar, T. M. (2019) 'Medicare fraud detection using neural Networks', *Journal of Big Data* 6. Available at: <https://doi.org/10.1186/s40537-019-0225-0>

Karangara, R., Devineni, S. K. & Challa, N. (2024) 'Enhancing Explainability in AI Fraud Detection', *International Journal of Computer Techniques* 11(1). Available at: [https://www.researchgate.net/publication/377329151\\_Enhancing\\_Explainability\\_in\\_AI\\_Fraud\\_Detection](https://www.researchgate.net/publication/377329151_Enhancing_Explainability_in_AI_Fraud_Detection)

Khan, I. K., Daud, H. B., Zainuddin, N. B., Sokkalingam, R., Museeb, A. & Inayat, A. (2024) 'Addressing limitations of the K-means clustering algorithm: outliers, non-



spherical data, and optimal cluster selection', *AIMS Mathematics* 9(9), pp. 25070-25097. Available at:

<https://www.aimspress.com/article/doi/10.3934/math.20241222?viewType=HTML>

Kotekani, S. S. & Ilango, V. (2022) 'HEMClust: An Improved Fraud Detection Model for Health Insurance using Heterogeneous Ensemble and K-prototype Clustering', *International Journal of Advanced Computer Science and Applications* 13(3), pp. 127-139. Available at:

<https://thesai.org/Publications/ViewPaper?Volume=13&Issue=3&Code=IJACSA&SerialNo=18>

Kunickaite, R., Zdanaviciute, M. & Krilavičius, T. (2020) 'Fraud Detection in Health Insurance Using Ensemble Learning Methods', *International Conference on Information Technology*. Available at: <https://www.semanticscholar.org/paper/Fraud-Detection-in-Health-Insurance-Using-Ensemble-Kunickaite-Zdanaviciute/ac028d54c9f5bedf8d22dc1e2cb41d4d4095ec06>

Kaur, K., Bhaskar, A., Pande, S. D., Malik, R. & Khamparia, A. (2021) 'An intelligent unsupervised technique for fraud detection in health care systems', *Intelligent Decision Technologies* 15(1), pp. 127-139. Available at: <https://journals.sagepub.com/doi/full/10.3233/IDT-200052>

Lalithagayatri, T., Priyanka, T. & Pavate, A. (2018) 'Fraud Detection in Health Insurance using Hybrid System', *International Journal of Engineering Research & Technology (IJERT)* 5(1). Available at: <https://www.ijert.org/fraud-detection-in-health-insurance-using-hybrid-system>

Liu, J., Wang, Y. & Yu, J. (2023) 'A study on the path of governance in health insurance fraud considering moral hazard', *Frontiers in Public Health*. Available at: <https://pmc.ncbi.nlm.nih.gov/articles/PMC10543491/>

Mambo, S. R. and Moturi, C. A. (2022) 'Towards Better Detection of Fraud in Health Insurance Claims in Kenya: Use of Naïve Bayes Classification Algorithm', *East African Journal of Information Technology* 5(1). Available at: <https://doi.org/10.37284/eajit.5.1.1023>

Massi, M. C., Ieva, F. & Lettieri, E. (2020) 'Data mining application to healthcare fraud detection: a two-step unsupervised clustering method for outlier detection with administrative databases', *BMC Medical Informatics and Decision Making* 20. Available at: <https://doi.org/10.1186/s12911-020-01143-9>

Nabrawi, E. and Alanazi, A. (2023) 'Fraud Detection in Healthcare Insurance Claims Using Machine Learning', *Risks* 11(9). Available at: <https://doi.org/10.3390/risks11090160>

Naidoo, K. & Marivate, V. (2020) 'Unsupervised Anomaly Detection of Healthcare Providers Using Generative Adversarial Networks', *Responsible Design, Implementation and Use of Information and Communication Technology*, pp. 419-30. Available at: <https://pmc.ncbi.nlm.nih.gov/articles/PMC7134221/>

Ono, S. & Goto, T. (2022) 'Introduction to supervised machine learning in clinical epidemiology', *Annals of Clinical Epidemiology* 4(3), pp. 63-71. Available at: <https://pubmed.ncbi.nlm.nih.gov/38504945/>

Parthasarathy, S., Lakshminarayanan, A. R., Khan, A. A. A., Sathick, K. J., & Jayaraman, V. (2023) 'Detection of Health Insurance Fraud using Bayesian Optimized XGBoost', *International Journal of Safety and Security Engineering* 13(5), pp. 853-861. Available at: <https://doi.org/10.18280/ijssse.130509>

Patil, P., Kale, G., Bivalkar, N. & Kolhatkar, A. (2023) 'Comparative Analysis of Weighted Ensemble and Majority Voting Algorithms for Intrusion Detection in OpenStack Cloud Environments', *International Journal of Advanced Computer Science and Applications* 14(12), pp. 741-747. Available at: <http://dx.doi.org/10.14569/IJACSA.2023.0141276>

Rojarath, A. & Songpan, W. (2020) Probability-Weighted Voting Ensemble Learning for Classification Model. *Journal of Advances in Information Technology* 11(4), pp.

217-227. Available at:

<https://www.jait.us/index.php?m=content&c=index&a=show&catid=203&id=1121>

Rosa, R. & Khoshgoftaar, T. (2018) 'Identifying Medicare Provider Fraud with Unsupervised Machine Learning', *2018 IEEE International Conference on Information Reuse and Integration (IRI)*, pp. 285-292. Available at:

<https://ieeexplore.ieee.org/document/8424722>

Sadiq, S. & Shyu, M.-L. (2019) 'Cascaded Propensity Matched Fraud Miner: Detecting Anomalies in Medicare Big Data', *Journal of Innovative Technology* 1(1), pp. 51-61. Available at: [http://doi.org/10.29424/JIT.201903\\_1\(1\).0007](http://doi.org/10.29424/JIT.201903_1(1).0007)

Seshagiri, S. & Prema, K. V. (2025) 'Efficient Handling of Data Imbalance in Health Insurance Fraud Detection Using Meta-Reinforcement Learning', *IEEE Access* 13, pp. 23482-23497. Available at: <https://ieeexplore.ieee.org/document/10858064>

Sharma, A. (2024) 'Predictive Accuracy of Machine Learning Models in Fraud Detection for Health Insurance in India', *American Journal of Statistics and Actuarial Science* 5(2), pp. 1-12. Available at:

<https://ajpojournals.org/journals/index.php/AJSAS/article/view/2253>

Sun, C., Li, Q., Li, H., Shi, Y., Zhang, S., Guo, W. (2018) 'Patient Cluster Divergence Based Healthcare Insurance Fraudster Detection', *IEEE Access*. Available at: <https://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=8576507>

Tabassum, M., Mahmood, S., Bukhari, A., Alshemaimri, B., Daud, A. & Khalique, F. (2024) 'Anomaly-based threat detection in smart health using machine learning', *BMC Medical Informatics and Decision Making* 24(1). Available at: <https://pmc.ncbi.nlm.nih.gov/articles/PMC11577804/>

Villegas-Ortega, J., Bellido-Boza, L. & Mauricio, D. (2021) 'Fourteen years of manifestations and factors of health insurance fraud, 2006–2020: a scoping review', *Health and Justice* 9. Available at: <https://healthandjusticejournal.biomedcentral.com/articles/10.1186/s40352-021-00149-3>

Wang, Z., Chen, X., Wu, Y., Jiang, L., Lin, S. & Qiu, G. (2025) 'A robust and interpretable ensemble machine learning model for predicting healthcare insurance fraud', *Scientific Reports* 15(1). Available at: <https://pmc.ncbi.nlm.nih.gov/articles/PMC11697271/>

Weidel, P., Duarte, R. & Morrison, A. (2021) 'Unsupervised Learning and Clustered Connectivity Enhance Reinforcement Learning in Spiking Neural Networks', *Frontiers*

*in Computational Neuroscience* 15. Available at:

<https://doi.org/10.3389/fncom.2021.543872>

Zhang, C., Xiao, X. & Wu, C. (2020) 'Medical Fraud and Abuse Detection System Based on Machine Learning', *International Journal of Environmental Research and Public Health* 17(19). Available at: <https://doi.org/10.3390/ijerph17197265>

Zhao et al. (2021) 'SUOD: Accelerating Large-Scale Unsupervised Heterogeneous Outlier Detection', *Proceedings of the 4th Conference on Machine Learning and Systems (MLSys)*. Available at:  
[https://proceedings.mlsys.org/paper\\_files/paper/2021/file/37385144cac01dff38247ab11c119e3c-Paper.pdf](https://proceedings.mlsys.org/paper_files/paper/2021/file/37385144cac01dff38247ab11c119e3c-Paper.pdf)

Zhou, T. & Jiao, H. (2022) 'Exploration of the Stacking Ensemble Machine Learning Algorithm for Cheating Detection in Large-Scale Assessment' *Educational and Psychological Measurement* 83(4), pp. 831-854. Available at:  
<https://pmc.ncbi.nlm.nih.gov/articles/PMC10311957/>