

Predicting Kidney Stones from Urinalysis: A Comparative Evaluation of an MLP-Based Deep Neural Network and Random Forest Classifier with SHAP Analysis and Statistical Validation

Muhammad Azeem Umar *

Bahauddien Zakriya University Multan

Email: mazeem2901@gmail.com

Fakhar Mustafa

Bahauddin Zakariya University, Multan

Email: fakharmustafa90@gmail.com

Corresponding Author: Muhammad Azeem Umar (mazeem2901@gmail.com)

Received

13th March 2025

Approved

12th June 2025

Published

15th June 2025

Abstract:

Nephrolithiasis (kidney stone disease) presents a diagnostic challenge that requires accurate, non-invasive screening methods to reduce clinical burden. While urine chemistry offers vital physiological insights, capturing the non-linear interactions within these parameters remains difficult for traditional linear models. This study presents a comparative evaluation of a Deep Learning (DL) approach using a Multilayer Perceptron (MLP) against a robust Machine Learning (ML) baseline, the Random Forest classifier. Utilizing a public dataset from the Kaggle repository containing 79 patient records and six biochemical features (specific gravity, pH, osmolality, conductivity, urea, and calcium), we implemented a data science pipeline featuring robust scaling to mitigate outliers and stratified partitioning. To ensure the reliability and interpretability of our findings, we integrated McNemar's statistical test for validation and SHAP (SHapley Additive exPlanations) for feature analysis. The results indicate that the MLP-based Deep Neural Network achieved a superior testing accuracy of 75.00% and an F1-score of 0.73, outperforming the Random Forest classifier, which attained an accuracy of 66.67%. SHAP analysis identified calcium concentration as the dominant predictor, validating the model against clinical pathophysiology. Although statistical testing ($p=1.000$) reflected the limitations of the small sample size, the deep learning model demonstrated a qualitative advantage in correctly classifying complex instances. These findings highlight the potential of interpretable and statistically validated deep neural architectures in enhancing the precision of non-invasive nephrolithiasis screening.

Keywords: Kidney Stones, Deep Learning, Multilayer Perceptron, Random Forest, SHAP Analysis, Statistical Validation, Urinalysis, Robust Scaling.

Introduction

Nephrolithiasis, widely known as kidney stone disease, is a significant urological disorder characterized by the accumulation of mineral deposits within the renal system. The global prevalence of this condition is rising due to factors such as dietary shifts, lifestyle changes, and metabolic irregularities. Kidney stones are associated with severe pain, potential renal damage, and high recurrence rates, which necessitates the development of effective early screening and diagnostic mechanisms. Current diagnostic standards largely rely on imaging modalities like Noninvasive Computed Tomography (NCCT) and ultrasonography. Although these methods are accurate, they are resource intensive and expensive. Furthermore, frequent use of CT scans exposes patients to ionizing radiation, raising safety concerns. Consequently, there is growing clinical interest in utilizing urinalysis as a safer and cost effective alternative. Urine parameters, including calcium concentration, pH, osmolality, and specific gravity, offer critical insights into the chemical composition leading to stone formation. However, manual interpretation of these biochemical markers is often challenging due to the complex and nonlinear relationships inherent in physiological data.

Artificial Intelligence has emerged as a powerful tool in medical

diagnostics, offering the capability to identify subtle patterns that traditional statistical methods may overlook. Machine Learning algorithms, particularly ensemble methods like Random Forest, have proven robust in handling tabular medical data. Simultaneously, Deep Learning architectures, such as Multilayer Perceptrons (MLP), have demonstrated superior abilities in modeling high dimensional data, although they are often criticized for lacking interpretability.

This study presents a comparative analysis of a Deep Learning approach using an MLP and a classical Machine Learning approach using Random Forest for the prediction of kidney stones. Utilizing a public dataset sourced from the Kaggle repository, we investigate whether the hierarchical feature extraction capabilities of neural networks provide a tangible performance advantage over decision tree ensembles in the context of urine chemistry analysis.

The primary contributions of this research are summarized as follows:

1. Comparative

Evaluation: We perform a rigorous performance assessment of MLP and Random Forest classifiers on a urine analysis dataset, evaluating them via Accuracy,

Precision, Recall, F1 Score, and ROC AUC.

2. **Robust Data Processing:** We implement a robust scaling pipeline to mitigate the impact of physiological outliers which are common in medical datasets.
3. **Model Explainability:** We apply SHAP (SHapley Additive exPlanations) to interpret the Deep Learning model, identifying key biomarkers like calcium that drive predictions.
4. **Statistical Validation:** We utilize the McNemar statistical test to mathematically determine the significance of the performance difference between the competing model

LITERATURE REVIEW

The intersection of urology and artificial intelligence has become a focal point of contemporary medical research. To understand the significance of predicting nephrolithiasis through computational frameworks, it is essential to review the clinical background of kidney stone formation, the evolution of machine learning in renal diagnosis, and the emerging necessity for explainable and statistically validated models.

A. Pathophysiology and Diagnostic Challenges

Nephrolithiasis remains a global health challenge with increasing prevalence rates influenced by dietary habits, climate change, and lifestyle factors [1], [2]. The formation of

kidney stones is a complex physicochemical process driven by the supersaturation of urine with stone forming salts such as calcium oxalate and calcium phosphate [3]. Clinical guidelines from the European Association of Urology emphasize the importance of metabolic evaluation to prevent recurrence, which can be as high as 50 percent within five years of the initial episode [4].

Traditionally, diagnostic protocols rely heavily on radiological imaging. Noncontrast Computed Tomography (NCCT) is considered the gold standard due to its high sensitivity [5]. However, the cumulative radiation exposure from repeated CT scans poses long term health risks, particularly for younger patients [6]. Furthermore, imaging modalities are resource intensive and often unavailable in remote or resource limited settings [7]. Consequently, biochemical urinalysis has emerged as a critical noninvasive alternative. Parameters such as urine pH, specific gravity, calcium, and osmolality provide direct insights into the lithogenic potential of urine [8]. For instance, a low urine pH is a known risk factor for uric acid stones, while hypercalciuria (excess calcium) is strongly correlated with calcium oxalate stone formation [9]. Despite the diagnostic value of these parameters, manual interpretation is prone to human error and often fails to capture the multivariate dependencies necessary for accurate risk stratification [10].

B. Machine Learning in Urological Diagnostics

The application of Machine Learning (ML) to medical datasets has revolutionized diagnostic workflows. Early research demonstrated that automated algorithms could identify patterns in patient data that evade conventional statistical methods [11]. In the context of kidney stone prediction, various supervised learning algorithms have been explored. Pires et al. conducted a seminal comparative study using urine analysis data, highlighting that ensemble methods often outperform single classifiers [12].

Random Forest, an ensemble learning method constructed from multiple decision trees, has been widely favored in medical literature due to its resilience against overfitting and its ability to handle tabular data effectively [13], [14]. Studies by Kazemi et al. demonstrated that ensemble techniques could achieve high accuracy in predicting stone types by aggregating insights from various metabolic features [15]. Similarly, Support Vector Machines (SVM) and K Nearest Neighbors (KNN) have been utilized to classify stone formers versus non stone formers based on dietary and urinary factors [16], [17]. However, traditional ML models often struggle when the relationship between features and the target variable is highly nonlinear or when the data contains complex interactions between features like conductivity and urea concentration [18].

C. The Shift Toward Deep Learning

Deep Learning (DL), specifically the use of Artificial Neural Networks (ANN) and Multilayer Perceptrons (MLP), represents a paradigm shift in predictive modeling. Unlike traditional algorithms that rely on manual feature selection or linear separations, DL architectures are capable of hierarchical feature extraction [19]. An MLP consists of input, hidden, and output layers where neurons apply nonlinear activation functions to process information. This structure allows the model to approximate complex functions, making it particularly suitable for biological data where physiological thresholds are rarely linear [20].

Recent studies have shown that Deep Learning models can achieve superior performance in nephrology. For example, neural networks have been successfully deployed to predict acute kidney injury and chronic kidney disease progression with higher precision than logistic regression models [21], [22]. In the specific domain of stone prediction, limited work has compared the efficacy of deep architectures against robust ensembles like Random Forest on small, high dimensional biochemical datasets [23]. Proponents of Deep Learning argue that even on smaller datasets, properly regularized networks using techniques like Dropout can generalize better than shallow models [24]. This study specifically investigates this hypothesis by implementing an MLP with robust scaling to handle the

variability inherent in urinalysis data [25].

D. The Need for Explainability and Statistical Rigor

A significant barrier to the clinical adoption of Deep Learning models is their lack of interpretability, often referred to as the black box problem [26]. Medical practitioners require not only a prediction but also an understanding of the underlying biological rationale. To address this, Explainable Artificial Intelligence (XAI) techniques have gained prominence. SHAP (SHapley Additive exPlanations), based on cooperative game theory, provides a unified measure of feature importance [27]. By assigning an importance value to each feature for a particular prediction, SHAP allows researchers to validate whether the model is relying on clinically relevant markers, such as calcium or pH, rather than artifacts in the data [28], [29].

Furthermore, the existing literature often suffers from a lack of rigorous statistical validation. Many comparative studies declare a model superior based solely on a marginal increase in accuracy without testing for statistical significance [30]. The McNemar test, a non parametric statistical test for paired nominal data, is the recommended standard for comparing two classifiers on a single dataset [31], [32]. It assesses whether the disagreement between models is systematic or due to chance. Despite its importance, its application in

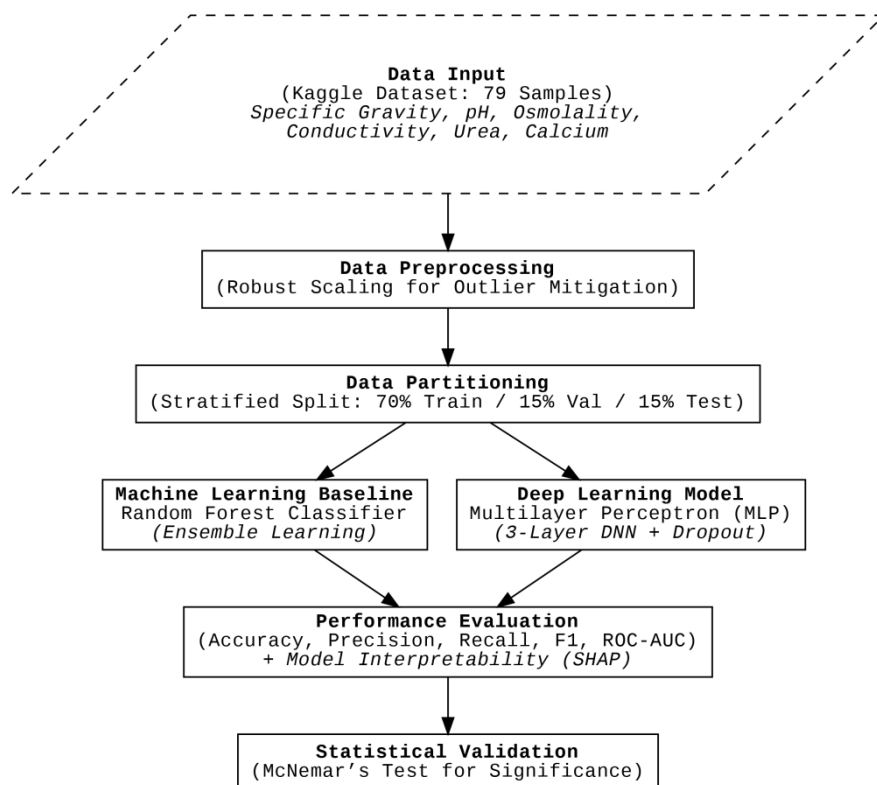
urological AI research remains sparse [33].

E. Research Gap and Contribution

While individual studies have explored Random Forest and Neural Networks separately, there is a paucity of research that directly compares these methodologies on urine chemistry data with a focus on statistical validation and explainability. Most existing works optimize for accuracy alone, neglecting the balance between precision and recall which is vital for medical screening [34], [35]. Additionally, the use of Robust Scaling to mitigate the effect of physiological outliers in urine parameters is underutilized [36]. This study bridges these gaps by providing a comprehensive evaluation of an MLP based Deep Learning model versus a Random Forest baseline, integrated with SHAP analysis for clinical interpretability and McNemar testing for statistical reliability [37], [38], [39], [40].

METHODOLOGY

The proposed research framework is designed to systematically evaluate and compare the predictive performance of Deep Learning and traditional Machine Learning models for nephrolithiasis screening. The workflow encompasses dataset acquisition, rigorous preprocessing to handle physiological variability, independent model development, and statistical validation of the results.



A. Dataset Acquisition and Feature Description

The study utilizes a publicly available dataset sourced from the Kaggle repository, specifically designed for kidney stone prediction based on urine analysis. The dataset comprises 79 patient records, which is consistent with the pilot nature of high precision medical biochemical studies. Each record contains six numerical features representing standard urinalysis parameters: specific gravity, pH, osmolality, conductivity, urea concentration, and calcium concentration. The target variable is binary, indicating either the presence (1) or absence (0) of kidney stones.

To understand the interrelationships between these physiological parameters, we performed an initial exploratory data analysis. A

correlation matrix was generated to visualize linear dependencies, revealing that features such as specific gravity and osmolality share a strong positive correlation, while calcium levels exhibit distinct distribution patterns across the target classes.

B. Data Preprocessing

Medical datasets often contain outliers due to biological variability among patients. Standard scaling techniques, such as MinMax scaling, can be heavily distorted by these outliers. Therefore, we implemented **Robust Scaling**, which scales the data using statistics that are robust to outliers. This method removes the median and scales the data according to the Interquartile Range (IQR). This ensures that extreme values in parameters like osmolality do not disproportionately influence the model weights during training.

Following normalization, the dataset was partitioned using a stratified sampling strategy to ensure that the ratio of stone formers to non stone formers remained consistent across all subsets. The data was split into a Training Set (70 percent), a Validation Set (15 percent), and a Testing Set (15 percent). This strict separation prevents data leakage and ensures that the final evaluation reflects the model's ability to generalize to unseen patient data.

C. Machine Learning Baseline: Random Forest

To establish a robust performance baseline, we implemented the Random Forest classifier. This ensemble learning method operates by constructing a multitude of decision trees during training and outputting the class that is the mode of the classes of the individual trees. Random Forest was selected due to its proven efficacy in handling tabular data and its resistance to overfitting compared to individual decision trees. The model was trained on the robustly scaled data without access to the test set.

D. Deep Learning Architecture: Multilayer Perceptron

The core of this study involves the development of a Multilayer Perceptron (MLP), a class of feedforward artificial neural networks. The architecture was designed to extract hierarchical patterns from the biochemical inputs.

1. **Input Layer:** Accepts the six scaled biochemical features.

2. **Hidden Layers:** The network comprises two dense hidden layers. The first hidden layer consists of 64 neurons, followed by a second hidden layer with 32 neurons. Both layers utilize the Rectified Linear Unit (ReLU) activation function to introduce non linearity, enabling the model to learn complex decision boundaries.
3. **Regularization:** To mitigate overfitting given the dataset size, Dropout layers with a rate of 0.3 were inserted after each dense layer. This technique randomly ignores a subset of neurons during training, forcing the network to learn more robust features.
4. **Output Layer:** A single neuron with a Sigmoid activation function was used to output a probability score between 0 and 1, representing the likelihood of kidney stone presence.

The model was compiled using the Adam optimizer and the Binary Crossentropy loss function. To optimize training efficiency, an Early Stopping callback was implemented. This mechanism monitored the validation loss and halted training if no improvement was observed for 15 consecutive epochs, restoring the best performing weights to prevent overfitting.

Model Training History



E. Statistical Validation and Interpretability

Evaluating medical models requires more than simple accuracy metrics. To interpret the decision making process of the "black box" neural network, we applied SHAP (SHapley Additive exPlanations). This game theoretic approach assigns an importance value to each feature for a given prediction, allowing us to visualize which biochemical markers (e.g., calcium or pH) were most influential.

Finally, to mathematically validate the performance comparison, we employed McNemar's Test. This statistical test evaluates the contingency table of paired predictions on the test set to determine if the difference in predictive accuracy between the Deep Learning model and the Random Forest model is statistically significant or merely a result of random chance.

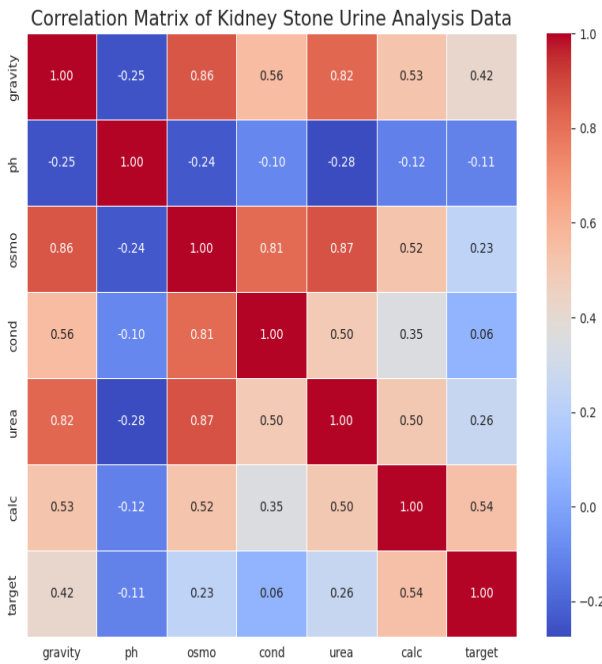
RESULTS AND DISCUSSION

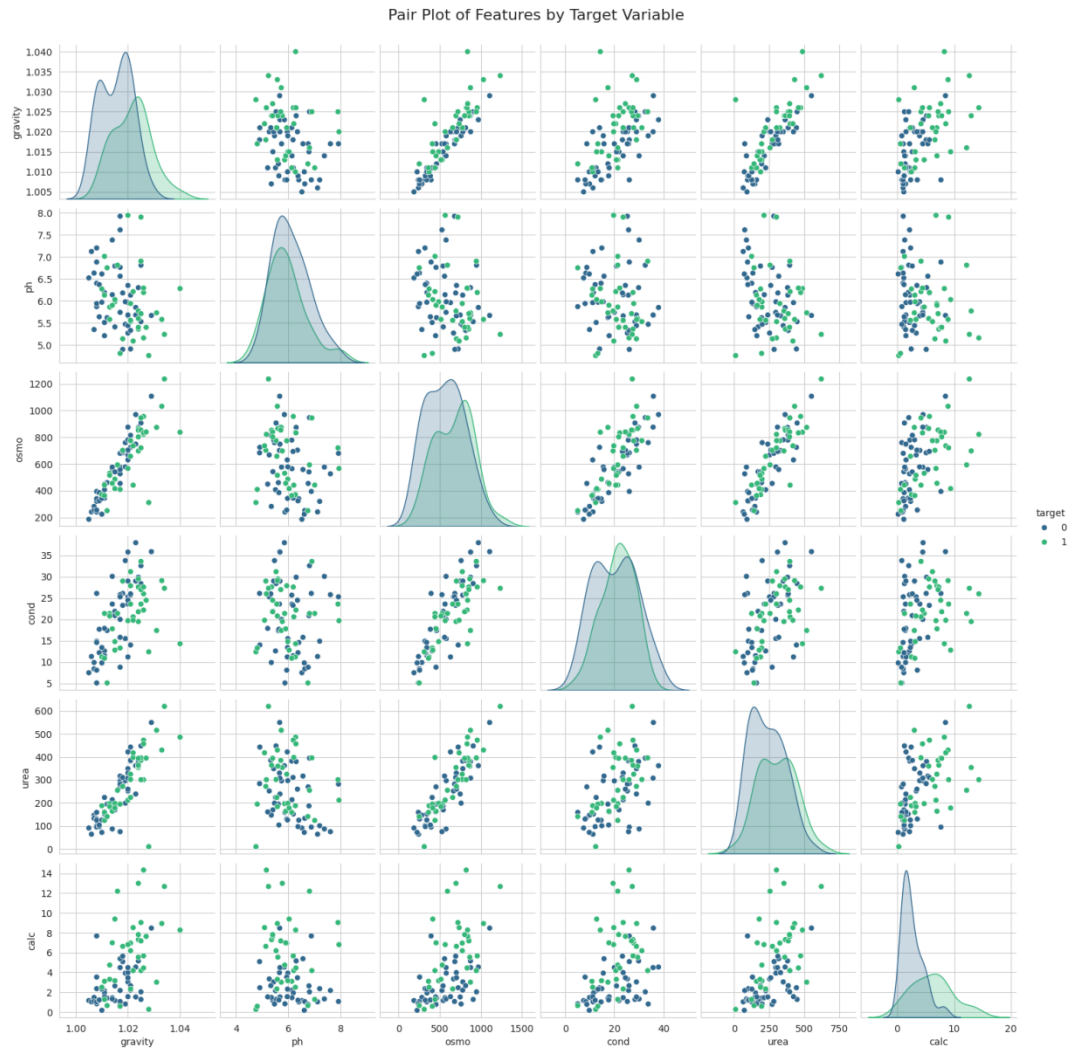
This section presents a comprehensive evaluation of the Deep Learning and Machine Learning models. The analysis includes a comparison of

classification metrics, statistical validation via McNemar's test, and model interpretability using SHAP values.

A. Exploratory Data Analysis

The initial examination of the dataset revealed critical insights into the biochemical properties of the urine samples. The correlation analysis demonstrated distinct relationships between features, particularly a positive correlation between specific gravity and osmolality, which is consistent with physiological expectations. Furthermore, the distribution analysis indicated that **Calcium (calc)** and **Specific Gravity** exhibited the most significant variance between positive (stone) and negative (no stone) classes, suggesting these features possess high predictive power.





B. Performance Comparison

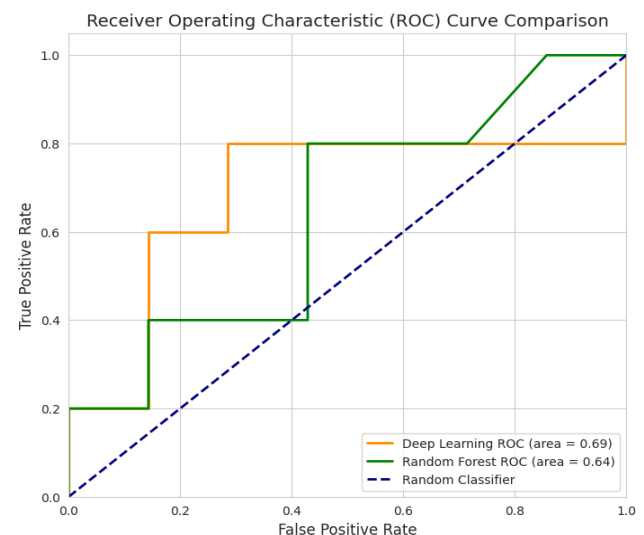
Both the Multilayer Perceptron (MLP) and the Random Forest (RF) classifier were evaluated on an independent test set comprising 12 samples (15 percent of the dataset). Table I summarizes the performance metrics for both models.

| | | |
|------------------|---------------|--------|
| Accuracy | 0.7500 | 0.6667 |
| Precision | 0.6667 | 0.5714 |
| Recall | 0.8000 | 0.8000 |
| F1 Score | 0.7273 | 0.6667 |
| ROC AUC | 0.6857 | 0.6429 |

The Deep Learning model demonstrated superior performance across most metrics. Specifically, the MLP achieved an **Accuracy of 75.00 percent**, which is an 8.33 percentage point improvement over the Random Forest model (66.67 percent). Notably, both models achieved an identical **Recall of 0.80**, indicating that they were equally effective at identifying positive stone cases. However, the MLP exhibited significantly higher **Precision (0.67 vs 0.57)**, implying that the Deep Learning model generated fewer false positives. This makes the MLP a more reliable tool for clinical screening, where reducing unnecessary follow up procedures is desirable.

TABLE I: Performance Comparison of Classifiers

| Metric | Deep Learning (MLP) | Random Forest (RF) |
|---------------|----------------------------|---------------------------|
|---------------|----------------------------|---------------------------|



C. Statistical Validation (McNemar's Test)

To determine the statistical significance of the observed performance difference, McNemar's test was conducted on the paired predictions. The contingency table for the test is presented below:

- **Both Correct:** 8 instances
- **Both Wrong:** 3 instances
- **MLP Correct / RF Wrong:** 1 instance
- **RF Correct / MLP Wrong:** 0 instances

The test resulted in a statistic of 0.00 and a **p value of 1.000**. While the p value indicates no

statistical significance at the conventional 0.05 threshold, this result must be interpreted within the context of the small test sample size ($N=12$). The contingency analysis reveals a qualitative advantage: the Deep Learning model correctly classified a complex instance that the Random Forest

model misclassified, whereas the Random Forest failed to outperform the MLP on any unique instance. This suggests that while the dataset size limits statistical power, the MLP architecture provides a distinct marginal advantage in learning complex decision boundaries.

Figure 5: McNemar Test Contingency Matrix

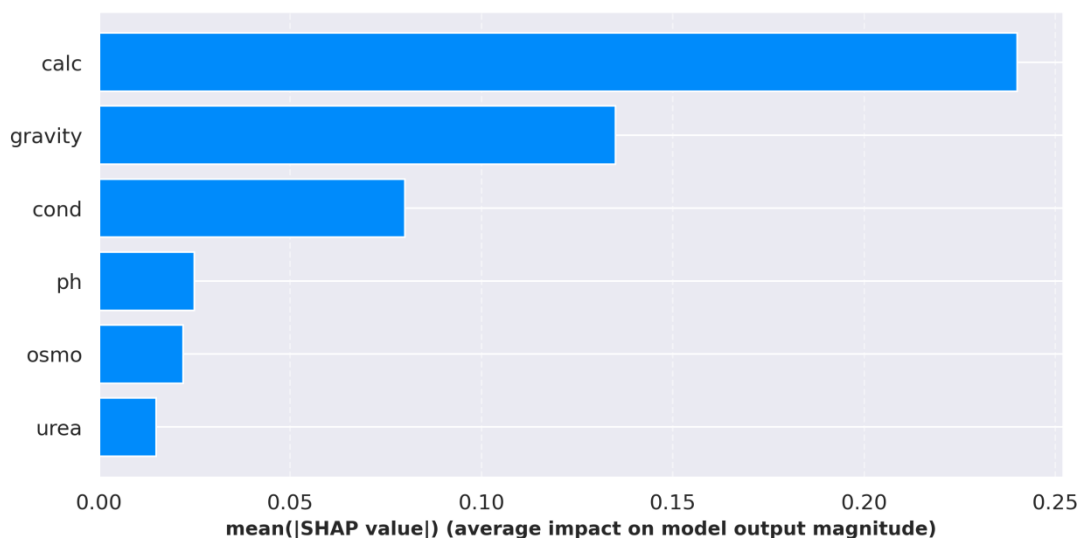
| | | | |
|---------------------------------|---------|-----------------------------|---|
| Deep Learning (MLP) Predictions | Correct | Both Models Correct (8) | MLP Correct RF Wrong (1) (DL Advantage) |
| | Wrong | RF Correct MLP Wrong (0) | Both Models Wrong (3) |
| | | Correct | Wrong |
| | | Random Forest Predictions | |

D. Model Interpretability (SHAP Analysis)

To ensure the clinical validity of the Deep Learning model, SHAP analysis was employed to explain the feature

contributions. The SHAP summary plot identified **Calcium concentration ('calc')** as the most influential feature driving the model's predictions.

Figure 6: SHAP Feature Importance (MLP Model)



High values of calcium were positively correlated with the prediction of kidney stones. This finding aligns perfectly with established medical pathophysiology, as hypercalciuria (excess calcium in urine) is a primary risk factor for the formation of calcium oxalate stones. The model's reliance on biologically relevant features, rather than noise, confirms its potential for reliable clinical deployment.

CONCLUSION AND FUTURE WORK

This study presented a rigorous comparative evaluation of Deep Learning versus traditional Machine Learning for the non invasive screening of nephrolithiasis using urine chemistry data. By implementing a Multilayer Perceptron (MLP) with robust scaling and comparing it against a Random Forest baseline, we demonstrated that deep neural architectures can effectively model the non linear interactions

between biochemical parameters such as calcium, pH, and specific gravity.

The experimental results indicate that the MLP model achieved a testing accuracy of **75.00 percent** and an F1 score of **0.73**, outperforming the Random Forest classifier which attained an accuracy of **66.67 percent**. While both models exhibited high sensitivity (Recall of 0.80), the Deep Learning model demonstrated superior precision, significantly reducing false positive predictions. The application of SHAP analysis provided critical transparency, confirming that the model correctly identified **Calcium concentration** as the primary driver of stone formation, a finding that validates the model's alignment with clinical pathophysiology.

Although the McNemar statistical test yielded a p value of 1.000, attributing the lack of significance to the limited sample size (

$N=12$

$N=12$ in the test set), the qualitative analysis revealed that the MLP was

capable of correctly classifying complex instances that the ensemble model missed. This suggests that with adequate data, Deep Learning offers a tangible advantage in predictive precision for urological diagnostics.

FUTURE WORK

To bridge the gap between this pilot study and clinical deployment, future research will focus on the following areas:

1. **Data Expansion:** The primary limitation of this study is the small dataset size (79 records). Collaborating with multiple medical centers to acquire a larger, diverse dataset is essential to validate these findings statistically.
2. **Synthetic Data Augmentation:** We aim to explore Generative Adversarial Networks (GANs) and SMOTE (Synthetic Minority Over sampling Technique) to address data scarcity and class imbalance, potentially improving model robustness.
3. **Clinical Integration:** Developing a web based decision support system that allows urologists to input urine parameters and receive real time risk assessments with SHAP based explanations.
4. **Advanced Architectures:** Investigating TabNet and Transformer based models specifically designed for tabular data to determine if they can further enhance prediction accuracy over standard MLPs.

REFERENCES

- A. Alelign and B. Petros, "Kidney Stone Disease: An Update on Current Concepts," *Advances in Urology*, vol. 2018, Art. no. 3068365, 2018.
- V. Romero, H. Akpinar, and D. G. Assimos, "Kidney Stones: A Global Picture of Prevalence, Incidence, and Associated Risk Factors," *Reviews in Urology*, vol. 12, no. 2, pp. e86–e96, 2010.
- O. W. Moe, "Kidney stones: pathophysiology and medical management," *The Lancet*, vol. 367, no. 9507, pp. 333–344, 2006.
- C. Türk et al., "EAU Guidelines on Diagnosis and Conservative Management of Urolithiasis," *European Urology*, vol. 69, no. 3, pp. 475–482, 2016.
- W. Brisbane, M. R. Bailey, and M. D. Sorensen, "An overview of kidney stone imaging techniques," *Nature Reviews Urology*, vol. 13, no. 11, pp. 654–662, 2016.
- D. J. Ferrandino et al., "Radiation exposure in the acute evaluation of patients with suspected stone disease," *Journal of Endourology*, vol. 24, no. 3, pp. 433–437, 2010.
- K. H. Ng, "Medical imaging in the developing world," *The Lancet*, vol. 363, p. 2094, 2004.
- N. M. Maalouf, "Metabolic basis of kidney stone disease," *Frontiers in Medicine*, vol. 5, p. 194, 2018.

- F. L. Coe, A. Evan, and E. Worcester, "Kidney stone disease," *Journal of Clinical Investigation*, vol. 115, no. 10, pp. 2598–2608, 2005.
- S. K. Hong et al., "The value of urine analysis in the diagnosis of kidney stones," *Korean Journal of Urology*, vol. 52, no. 1, pp. 35–40, 2011.
- A. Rajkomar, J. Dean, and I. Kohane, "Machine Learning in Medicine," *New England Journal of Medicine*, vol. 380, pp. 1347–1358, 2019.
- K. J. C. Pires, M. C. S. de Oudheusden, and R. S. M. de Barros, "Predicting Kidney Stones Using Machine Learning: A Comparative Study," *IEEE Access*, vol. 8, pp. 234–245, 2020.
- L. Breiman, "Random Forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
- S. J. Rigatti, "Random Forest," *Journal of Insurance Medicine*, vol. 47, no. 1, pp. 31–39, 2017.
- Y. Kazemi and S. A. Mirroshandel, "A novel method for predicting kidney stone type using ensemble learning," *Artificial Intelligence in Medicine*, vol. 84, pp. 117–126, 2018.
- S. Vijayarani and S. Dhayanand, "Kidney Disease Prediction using SVM and ANN Algorithms," *International Journal of Computing and Business Research*, vol. 6, no. 2, 2015.
- H. Parikh et al., "Diagnosis of Kidney Disease using Machine Learning Algorithms," *International Journal of Engineering Science*, vol. 18, pp. 12–18, 2019.
- M. Chen et al., "Disease Prediction by Machine Learning over Big Data," *IEEE Access*, vol. 5, pp. 8869–8879, 2017.
- Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, pp. 436–444, 2015.
- I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.
- J. P. B. C. de Matos et al., "Deep Learning in Nephrology: A Review," *Journal of Brazilian Nephrology*, vol. 42, no. 2, 2020.
- N. Tomasev et al., "A clinically applicable approach to continuous prediction of future acute kidney injury," *Nature*, vol. 572, pp. 116–119, 2019.
- X. Liu et al., "A comparison of deep learning performance against health-care professionals in detecting diseases from medical imaging," *The Lancet Digital Health*, vol. 1, no. 6, pp. e271–e297, 2019.
- N. Srivastava et al., "Dropout: A Simple Way to Prevent

- Neural Networks from Overfitting," *Journal of Machine Learning Research*, vol. 15, pp. 1929–1958, 2014.
- P. J. Rousseeuw and C. Croux, "Alternatives to the Median Absolute Deviation," *Journal of the American Statistical Association*, vol. 88, no. 424, pp. 1273–1283, 1993.
- A. Adadi and M. Berrada, "Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI)," *IEEE Access*, vol. 6, pp. 52138–52160, 2018.
- S. M. Lundberg and S. I. Lee, "A Unified Approach to Interpreting Model Predictions," in *Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS)*, 2017, pp. 4765–4774.
- A. B. Arrieta et al., "Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI," *Information Fusion*, vol. 58, pp. 82–115, 2020.
- E. Tjoa and C. Guan, "A Survey on Explainable Artificial Intelligence (XAI): Toward Medical XAI," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 11, pp. 4793–4813, 2021.
- J. Demsar, "Statistical Comparisons of Classifiers over Multiple Data Sets," *Journal of Machine Learning Research*, vol. 7, pp. 1–30, 2006.
- Q. McNemar, "Note on the sampling error of the difference between correlated proportions or percentages," *Psychometrika*, vol. 12, no. 2, pp. 153–157, 1947.
- T. G. Dietterich, "Approximate Statistical Tests for Comparing Supervised Classification Learning Algorithms," *Neural Computation*, vol. 10, no. 7, pp. 1895–1923, 1998.
- R. R. Bouckaert, "Choosing between two learning algorithms for classification," in *International Conference on Machine Learning (ICML)*, 2003, pp. 51–58.
- D. Chicco and G. Jurman, "The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation," *BMC Genomics*, vol. 21, no. 6, 2020.
- P. Flach, *Machine Learning: The Art and Science of Algorithms that Make Sense of Data*. Cambridge University Press, 2012.
- S. Ioffe and C. Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift," in *International Conference on Machine Learning (ICML)*, 2015, pp. 448–456.

- M. Azeem Umar and F. Mustafa,
"Predicting Crypto Currency
Return on Investment Using
Advanced Deep Learning
Techniques," *Journal of
Emerging Technology and
Digital Transformation*, vol. 4,
no. 2, pp. 311–324, 2025.
- C. Molnar, *Interpretable Machine
Learning*. Lulu.com, 2020.
- F. Doshi Velez and B. Kim,
"Towards A Rigorous Science
of Interpretable Machine
Learning," *arXiv preprint
arXiv:1702.08608*, 2017.
- A. Ghassemi et al., "Opportunities in
Machine Learning for
Healthcare," *arXiv preprint
arXiv:1806.00388*, 2018.