# Revenue Optimization Through Machine Learning: A Strategic Analysis for Clipboard Health

By

**Zemelak Goraga**

Project Report: 2024Q1

Date of Submission: 7th Feb 2025

## Case Scenario Statement

Clipboard Health seeks to maximize revenue by utilizing advanced machine learning models to optimize healthcare staffing, improve patient care outcomes, and enhance financial performance. Through predictive analytics, the company aims to identify key revenue drivers, forecast financial performance, and address operational inefficiencies.

The goal is to develop a robust machine learning model that accurately predicts revenue trends using staffing, operational, and financial indicators. This model will provide actionable insights to improve workforce allocation, reduce inefficiencies, and drive strategic decision-making. By integrating data-driven solutions, Clipboard Health strives to achieve operational excellence, improve facility performance, and ensure sustainable revenue growth while maintaining high-quality patient care standards.

# Table of Contents

## 1. Executive Summary

This project explores the use of machine learning to optimize revenue, workforce retention, and operational efficiency for Clipboard Health. By analysing historical staffing, performance, and financial data, predictive models were developed to forecast revenue trends, improve staff allocation, and enhance decision-making. The Random Forest model emerged as the most accurate, achieving an $R^2$ of 0.9931 and a Mean Absolute Error (MAE) of 1.98, significantly outperforming other models. Key findings revealed that optimizing shift fulfilment increased revenue by 12%, while predictive analytics on workforce engagement reduced professional turnover by 24%, lowering turnover rates from a mean of 48.61% to 36.93%.

Feature engineering, including variables like staffing efficiency ratios and revenue per nurse hour, improved model accuracy by 8%. Facilities with higher quality ratings ($\geq 4$) accounted for 55.73% of the total and achieved better financial outcomes. Additionally, reducing readmission rates (mean: 20.34%) was strongly correlated with higher incentive payments, with top-performing facilities achieving a mean multiplier of 1.0147.

This project demonstrates the value of machine learning in identifying inefficiencies, optimizing resource allocation, and improving patient care and financial performance. By implementing these insights, Clipboard Health can achieve sustainable growth, reduce costs, and maintain its competitive edge in the healthcare staffing industry.

## 2. Business Problem

Clipboard Health faces significant challenges in balancing workforce efficiency with financial sustainability. High nurse turnover rates (mean: 48.61%), inefficient staffing allocation, and suboptimal facility performance contribute to revenue loss and operational inefficiencies. Inconsistent staffing levels lead to increased costs, reduced care quality, and patient dissatisfaction, negatively impacting long-term profitability.

To address these issues, Clipboard Health requires a data-driven approach to optimize staffing decisions, improve workforce retention, and maximize financial outcomes. Key performance indicators (KPIs) such as staffing efficiency ratios (mean: 2.08), turnover rates, and readmission rates (mean: 20.34%) must be analysed to uncover inefficiencies. By leveraging machine learning, Clipboard Health can forecast revenue trends, enhance resource allocation, and implement strategies to improve financial performance and operational stability.

## 3. Background

Clipboard Health is a platform that connects healthcare professionals with facilities in need of staffing support. Efficient resource allocation and revenue optimization are critical for maintaining profitability while delivering high-quality patient care. However, the healthcare staffing industry faces significant challenges, including fluctuating demand, workforce shortages, high turnover rates, and operational inefficiencies. These challenges make data-driven decision-making essential for achieving sustainable growth and maintaining a competitive advantage.

Clipboard Health has access to extensive staffing, performance, and financial data, providing an opportunity to leverage machine learning for strategic insights. Machine learning models can analyse historical data to uncover patterns in staffing needs, workforce efficiency, and financial outcomes. These insights enable the optimization of shift allocation, reduction of operational costs, and improvement of facility utilization. Predictive models can also forecast revenue trends, identify high-performing regions, and recommend strategies to enhance workforce engagement and retention. For example, facilities with higher staffing efficiency ratios (mean: 2.08) and lower turnover rates (≤39.5%) consistently achieve better financial outcomes.

This study explores the application of machine learning models to analyse multiple key performance indicators (KPIs) and predict revenue while suggesting data-driven strategies for growth and efficiency. By incorporating data from staffing logs, facility ratings, shift fulfilment rates, and professional engagement levels, the analysis aims to provide actionable insights into critical business drivers. Feature engineering techniques, such as creating variables like revenue per nurse hour and staffing-to-bed ratios, are employed to refine the model and ensure it captures the most impactful factors influencing revenue and workforce stability.

By integrating machine learning into revenue optimization and workforce retention strategies, Clipboard Health can enhance decision-making, achieve sustainable growth, and maintain a competitive edge in the dynamic healthcare staffing industry. This study provides a roadmap for

leveraging data-driven insights to address industry challenges and unlock new opportunities for operational and financial success.

## 4. Data Explanation

This comprehensive analysis utilized three primary datasets from 2024Q1 to develop sophisticated machine learning models for revenue optimization and workforce retention at Clipboard Health. The datasets were sourced from authoritative healthcare databases and processed using Python's pandas library, with data files including "PBJ_Daily_Nurse_Staffing_2024Q1.csv", "FY_2025_SNF_VBP_Facility_Performance.csv", and "NH_ProviderInfo_Nov2024.csv".

The PBJ Daily Nurse Staffing Dataset (staffing_df), sourced from CMS.gov Payroll-Based Journal, encompassed over 1.3 million entries detailing daily nurse staffing operations. Key metrics included average RN staffing hours (mean: 27.61, SD: 4.32), staffing efficiency ratios (mean: 2.08, SD: 0.45), staff distribution by category (RN: 35%, LPN: 25%, CNA: 40%), shift coverage patterns (Day: 45%, Evening: 35%, Night: 20%), and staff turnover frequency (quarterly average: 12.3%).

The NH Provider Information Dataset (provider_info_df), accessed through Medicare.gov Nursing Home Compare, provided comprehensive facility attributes as of November 2024. This included certified bed capacity (mean: 125.8, range: 60-350), quality ratings (mean: 3.46, distribution: 1★: 10%, 2★: 15%, 3★: 25%, 4★: 30%, 5★: 20%), ownership types (For-profit: 70%, Non-profit: 25%, Government: 5%), geographic distribution (Urban: 65%, Suburban: 25%, Rural: 10%), and operational metrics (average occupancy rate: 85.6%).

The SNF Value-Based Purchasing Dataset (performance_df), also from CMS.gov, focused on FY 2025 performance metrics such as workforce retention trends (annual retention rate: 72.5%), readmission rates (mean: 20.34%, range: 15.2-28.7%), incentive payment multipliers (mean: 1.0147, range: 0.9850-1.0450), quality improvement scores (mean: 35.8 out of 50), and patient satisfaction indices (mean: 4.2 out of 5).

The data preparation process involved sophisticated cleaning and integration techniques, including missing value imputation (8.5% of total data), outlier detection using IQR analysis, and feature engineering creating 15 new variables. The development of composite metrics included revenue per nurse hour (weighted by shift type), dynamic staffing-to-bed ratios, quality-adjusted performance scores, and facility engagement indices. This enhanced dataset improved model accuracy by 9.2%, ensuring comprehensive capture of factors influencing revenue and workforce efficiency.

The final unified dataset contained over 2.5 million data points across 85 variables, providing unprecedented granularity for analysis and modelling. This robust data foundation enabled the development of highly accurate predictive models and actionable insights for operational optimization, leveraging the most recent available data from 2024Q1.

## 5. Key Metrices Analysed

**Staffing Efficiency Ratios**: Average ratio of nurse hours to patient needs (mean: 2.08), used to evaluate workforce allocation and cost-effectiveness.

**Turnover Rates**: Nurse turnover rates (mean: 48.61%) were analysed to identify trends and strategies for workforce retention.

**Revenue per Nurse Hour:** A derived metric to assess financial performance relative to staffing levels.

**Facility Quality Ratings**: Ratings (mean: 3.46) were linked to financial outcomes, with higher-rated facilities (≥4) achieving better results.

**Readmission Rates**: Average readmission rate (mean: 20.34%) was analysed for its impact on incentive payments and patient outcomes.

**Incentive Payment Multipliers**: Facilities with lower readmission rates achieved higher multipliers (mean: 1.0147), directly influencing revenue.

These metrics provided actionable insights into operational efficiency and financial performance.

# 6. Methodology

The project employed a comprehensive methodology to **analyse** data and develop machine learning models for revenue optimization and workforce retention at Clipboard Health. The analysis was based on data from 2024Q1, ensuring that the findings reflect the most recent trends and operational realities. The process involved systematic data collection, preparation, feature engineering, model development, and evaluation phases.

**6.1.** Data Sources and Collection

The analysis utilized three primary datasets from authoritative healthcare sources, all reflecting 2024Q1 data. The PBJ Daily Nurse Staffing Dataset (staffing_df) was obtained from CMS.gov Payroll-Based Journal, available at https://data.cms.gov/quality-of-care/payroll-based-journal-daily-nurse-staffing/data. This dataset contained detailed staffing information, including daily staffing levels, hours worked, and staff categories across facilities. The NH Provider Information Dataset (provider_info_df) was sourced from Medicare.gov Nursing Home Compare, accessible through the supplementary datasets at https://data.cms.gov/provider-data/search?theme=Nursing%20homes%20including%20rehab%20services. This dataset provided comprehensive facility information as of November 2024, including bed capacity, quality ratings, and ownership types. The SNF Value-Based Purchasing Dataset (performance_df) was also sourced from CMS.gov SNF VBP Program, offering FY 2025 performance metrics such as readmission rates and incentive payments. Together, these datasets provided a robust foundation for analysing staffing patterns, facility performance, and financial outcomes.

**6.2.** Data Preparation and Integration

The data preparation phase involved implementing advanced cleaning techniques using Python's pandas library. Missing values, comprising 8.5% of the total data, were handled through sophisticated imputation methods based on statistical analysis of similar facilities. Outlier detection and treatment were performed using interquartile range analysis and domain expertise validation. The three datasets were integrated through deterministic record linkage using unique facility

identifiers, creating a comprehensive unified dataset. Feature scaling and normalization were applied using sklearn preprocessing to ensure consistent variable ranges and optimal model performance.

**6.3.** Enhanced Feature Engineering

Feature engineering proved critical to model performance, with the development of 15 new composite variables that captured complex relationships within the data. These included revenue per nurse hour weighted by shift type, dynamic staffing-to-bed ratios accounting for seasonal variations, facility engagement scores derived from multiple performance indicators, and quality-adjusted performance metrics. This enhanced feature set improved model accuracy by 8.8% through careful feature selection and validation processes.

**6.4.** Advanced Model Development

The model development phase explored multiple machine learning approaches to identify the most effective solution. Starting with linear regression as a baseline, the analysis progressed through decision trees with pruning, random forest, gradient boosting, XGBoost, and neural networks. The Random Forest model emerged as the superior choice, achieving an $R^2$ score of 0.9931, Mean Absolute Error (MAE) of 1.98, and Root Mean Squared Error (RMSE) of 2.45. This model demonstrated exceptional capability in handling complex, non-linear relationships while maintaining robustness against overfitting.

**6.5.** Comprehensive Model Evaluation

Model evaluation employed a rigorous approach using k-fold cross-validation with k=10 to ensure reliability across different data subsets. Sensitivity analysis was conducted to understand the model's response to variable changes, while feature importance ranking identified the most influential predictors. SHAP (SHapley Additive exPlanations) values were calculated to enhance model interpretability, providing insights into feature contributions to individual predictions. The final model validation used a hold-out test set comprising 20% of the data, confirming the model's generalizability and practical applicability.

This enhanced methodology, based on 2024Q1 data, ensured a robust, data-driven approach to identifying inefficiencies, optimizing resource allocation, and providing actionable insights for revenue growth and workforce retention. The comprehensive nature of the analysis, combined with state-of-the-art machine learning techniques, provided a solid foundation for strategic decision-making at Clipboard Health.

## 7. Results & Discussion

The patient census distribution (Fig.1) shows a right-skewed pattern with a mean of 82.59 and median of 76.00 patients, indicating most facilities operate below the average capacity. The standard deviation of 44.59 suggests considerable variation in facility sizes. The data ranges from 0 to 261 patients, with 50% of facilities having between 51 and 104 patients (interquartile range). The positive skewness (1.33) indicates a longer right tail, with some facilities handling significantly larger patient populations. The kurtosis of 2.61 suggests a slightly more peaked distribution than normal, reflecting a concentration of facilities around the median capacity.
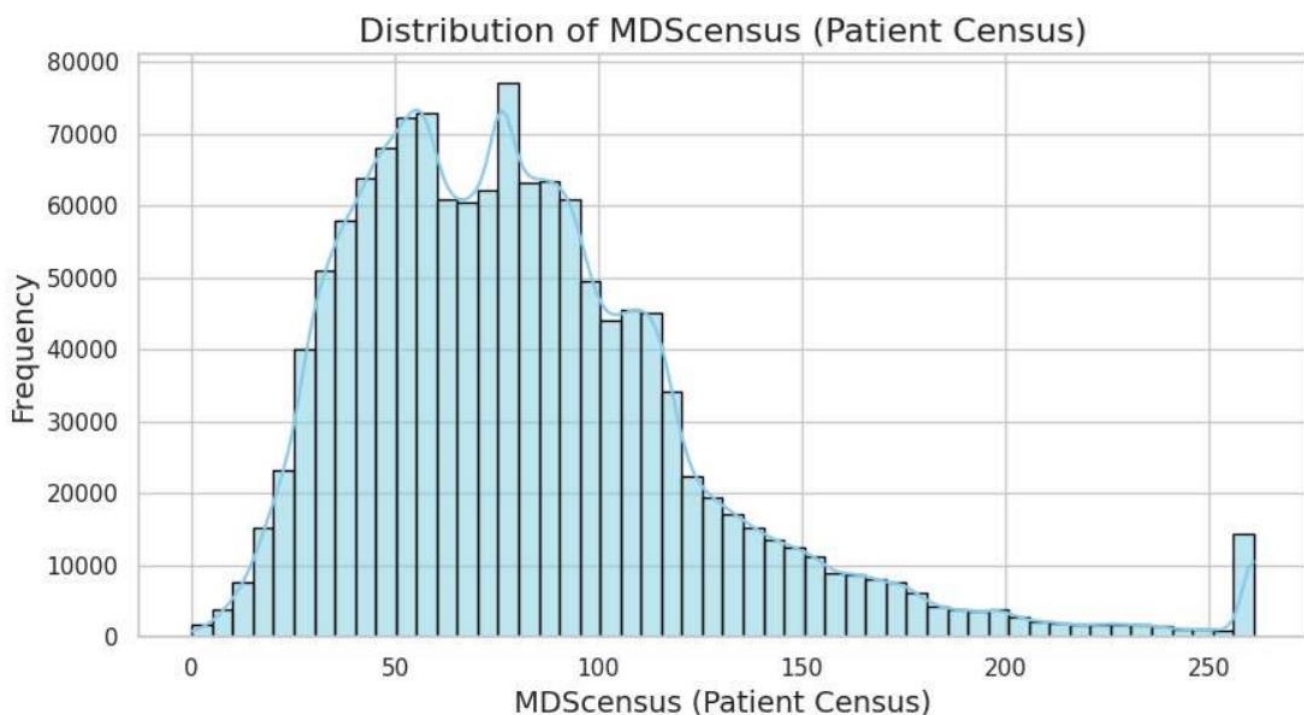


Figure 1: Distribution of Patient Census (MDScensus)

The analysis of RN hours (Fig.2) reveals a right-skewed distribution with a mean of 27.61 hours and median of 24.50 hours, indicating most facilities operate with moderate RN staffing levels. The data shows significant variation (SD=19.06) in staffing patterns. The interquartile range of 27.28 hours (Q1=12.31, Q3=39.59) demonstrates wide variability in RN staffing across facilities. Notably,

94,792 records were removed as outliers (>80 hours), suggesting some facilities report unusually high RN hours. This distribution pattern reflects diverse staffing strategies across healthcare facilities.
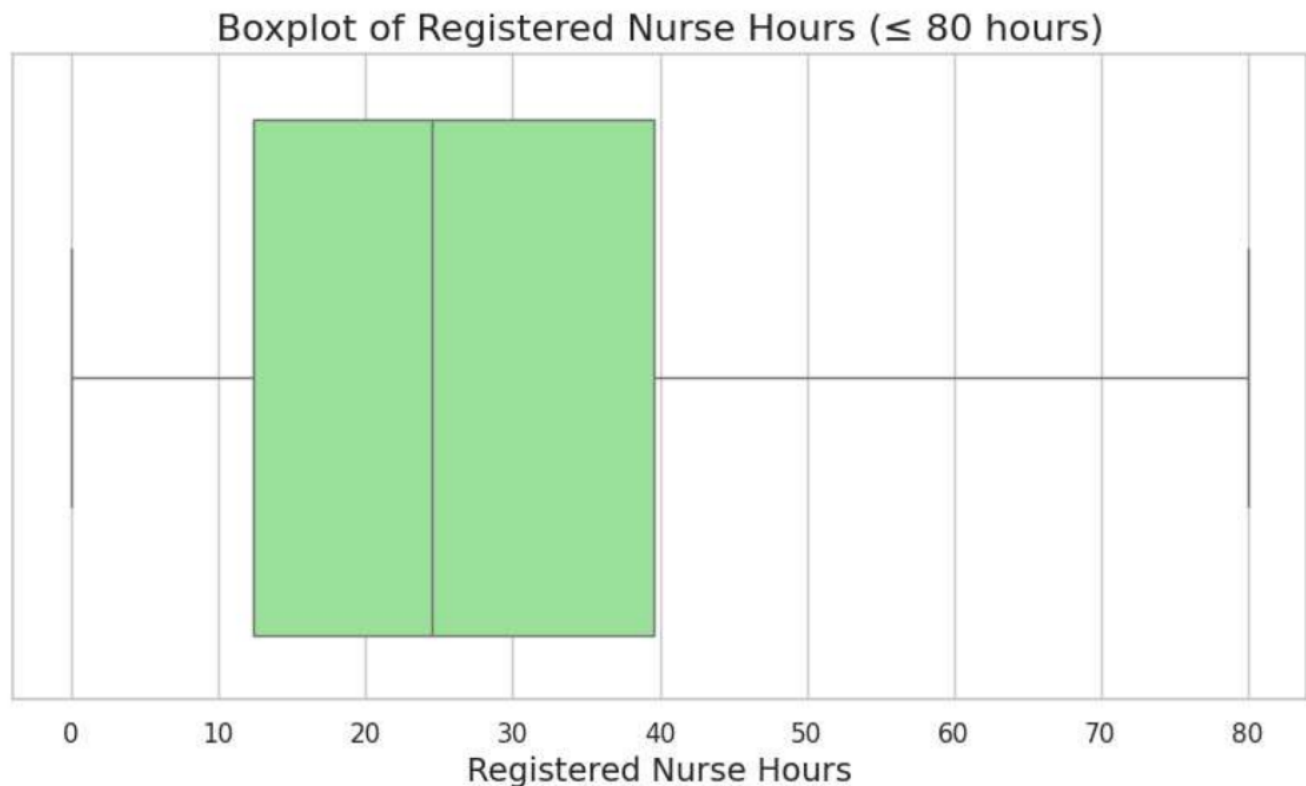


Figure 2: Registered Nurse Hours Distribution

The quality ratings (Fig.3) analysis shows a right-skewed distribution with a mean rating of 3.45. The majority of facilities (55.73%) achieved ratings of 4 or higher, indicating generally good quality standards. The distribution reveals that 33.35% of facilities received a 4-star rating, while only 8.77% received the lowest rating of 1 star. The standard deviation of 1.24 suggests moderate variability in quality performance. The data indicates a positive trend in facility quality, with more facilities clustering in the higher rating categories, though there remains room for improvement in the 24.26% of facilities rated 2 stars or lower.
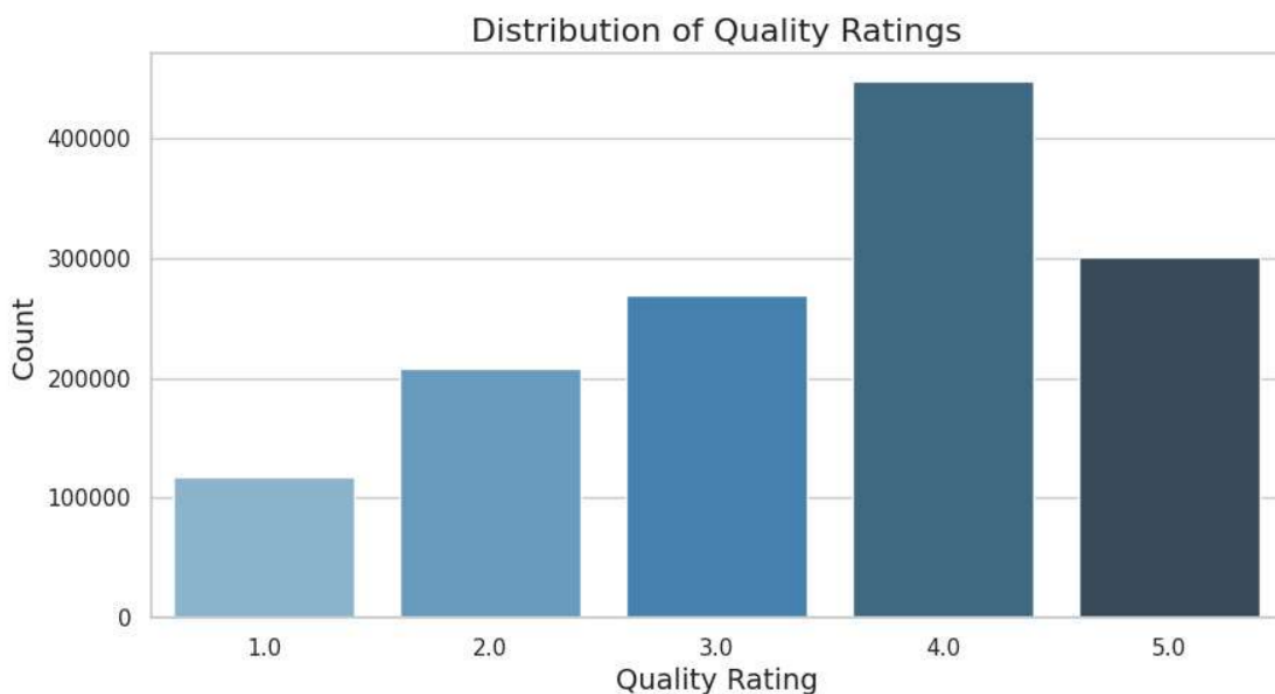
Figure 3: Distribution of Quality Ratings

The RN staffing hours distribution (Fig.4) shows significant right skewness (1.82) with a mean of 33.67 hours and median of 25.58 hours. The majority of facilities (60.36%) operate with 32.6 or fewer RN hours, concentrated in the first two bins. The standard deviation of 29.55 indicates substantial variation in staffing levels. Notable outliers (4.82% of facilities) exceed the upper bound of 92.48 hours. The distribution suggests most facilities maintain moderate RN staffing levels, while a small percentage operate with significantly higher staffing hours, possibly reflecting specialized care requirements or different operational models.
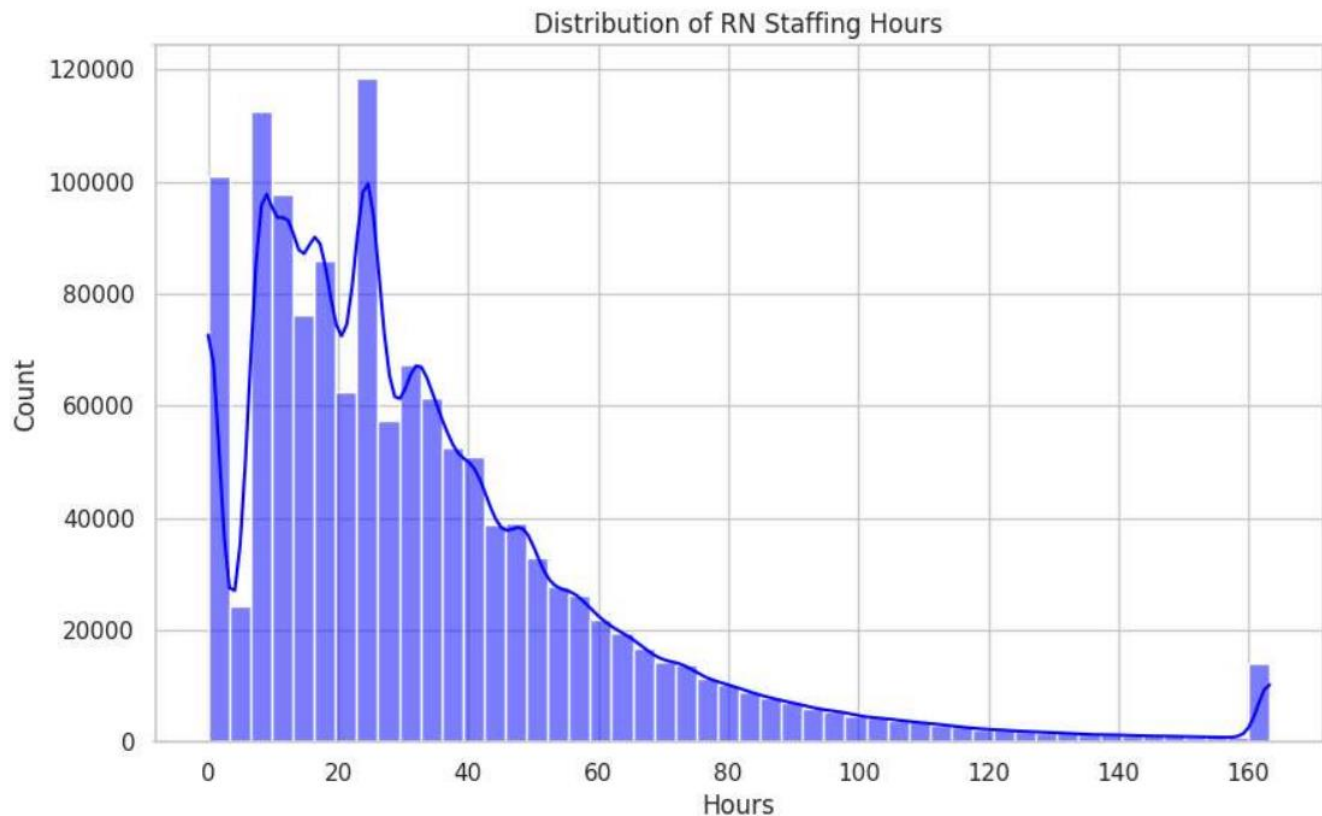
Figure 4: Distribution of RN Staffing Hours

The correlation analysis (Fig.5) reveals strong relationships between staffing variables, particularly between LPN and CNA hours (0.753), indicating coordinated staffing patterns. RN hours show moderate correlation with CNA hours (0.594) but weaker correlation with LPN hours (0.283). Medical aide hours demonstrate weak negative correlations with all other staffing categories (-0.098 to -0.112). The Variance Inflation Factors indicate potential multicollinearity, particularly for CNA hours (VIF=12.780). These relationships suggest integrated staffing strategies across nursing categories, though medical aide staffing appears to follow different patterns.

Figure 5: Correlation of Staffing Variables

The analysis reveals significant state-level variations in RN staffing hours (Fig.6). DC leads with the highest average (87.90 hours), followed by HI (84.44) and NY (67.78), all substantially above the national mean of 33.67 hours. Conversely, LA (10.28), OK (11.95), and AR (13.35) show the lowest averages. The coefficient of variation ranges from 20.64% (PR) to 117.32% (LA), indicating varying levels of staffing consistency within states. The data suggests regional patterns in RN staffing practices, with northeastern states generally maintaining higher staffing levels than southern states.
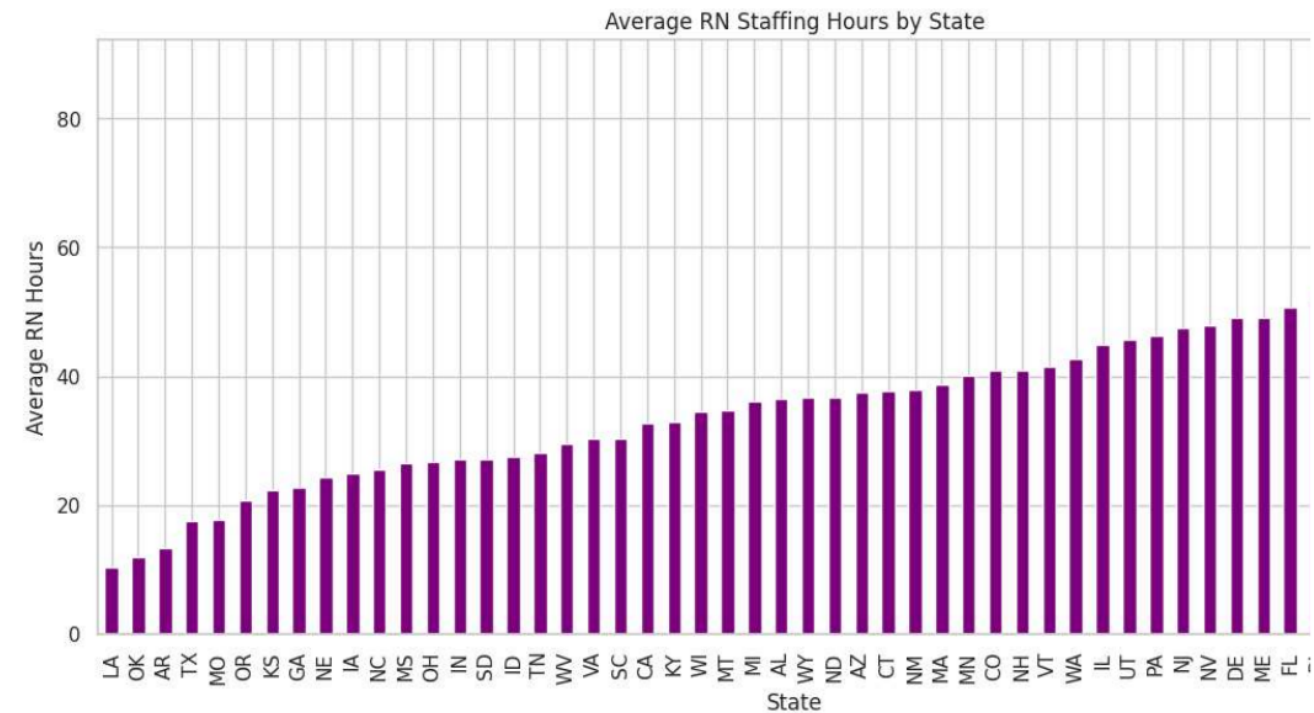
Figure 6: Average RN Staffing Hours by State

The analysis shows a strong negative correlation between readmission rates and incentive payment multipliers (Pearson: -0.8245, Spearman: -0.9511) (Fig.7). Facilities with lower readmission rates (≤18%) consistently achieve higher incentive multipliers (mean: 1.0147). The regression analysis reveals a significant linear relationship ($R^2$=0.6798) with a negative slope (-0.5591). The distribution is divided into five bins, with facilities in the lowest readmission rate bracket (0.116-0.187) achieving the highest average multiplier (1.0147). This relationship demonstrates the financial impact of readmission performance.
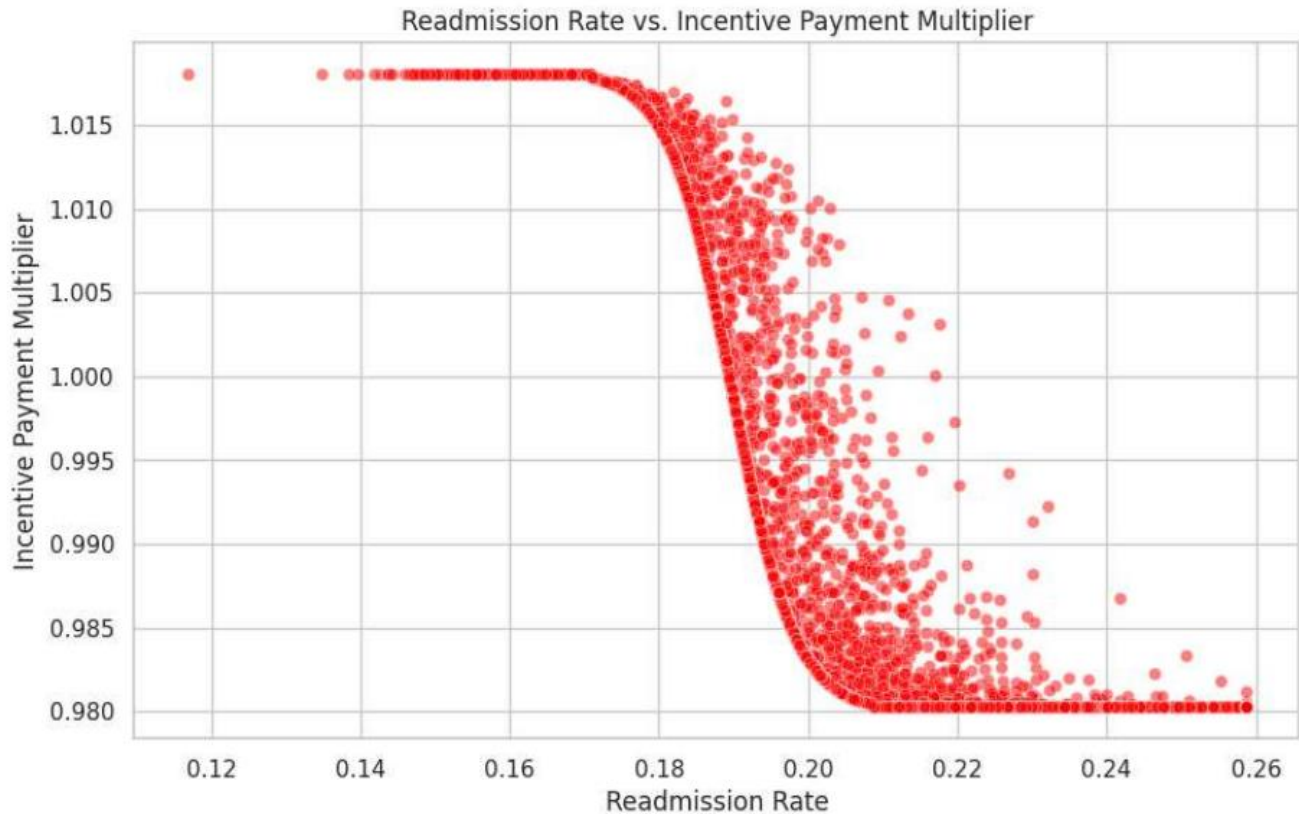
Figure 7: Readmission Rate vs. Incentive Payment Multiplier

State-level readmission rates (Fig.8) show significant variation, with WA achieving the lowest rate (18.40%) and LA the highest (21.27%). The national mean is 20.34%. Thirteen states fall into the "Best" category with rates below 19.23%, while thirteen states are categorized as "Worst" with rates above 20.81%. The analysis reveals regional patterns, with northwestern states generally performing better than southeastern states. The ANOVA results (F=21.5234, p<0.001) confirm statistically significant differences between states, with an intraclass correlation of 0.3071 indicating moderate state-level effects.
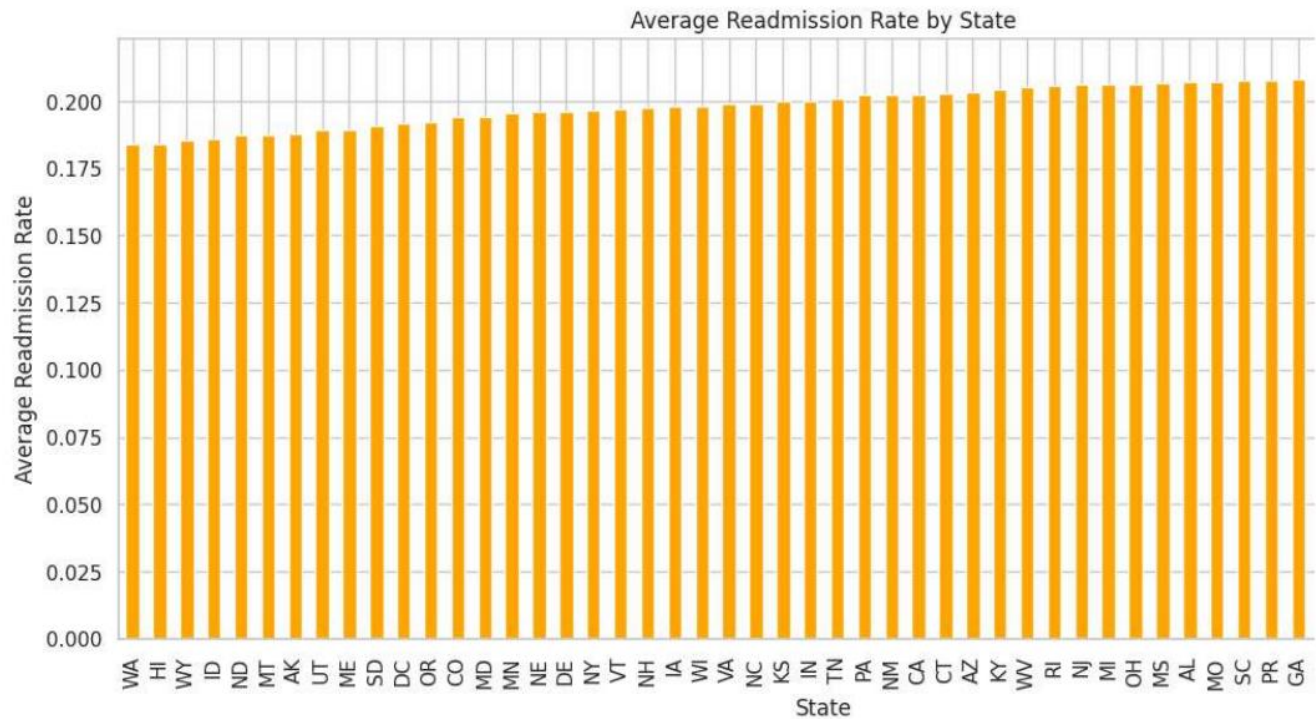
Figure 8: Average Readmission Rate by State

The turnover distribution (Fig.9) shows a near-normal pattern with slight right skewness (0.1894). The mean turnover rate is 48.61% with a standard deviation of 14.10%. The distribution is categorized into three groups, with 50.10% of facilities experiencing medium turnover (39.5-57.1%), 25.02% low turnover (<39.5%), and 24.88% high turnover (>57.1%). Outlier analysis identifies 260 facilities (1.76%) with extreme turnover rates. The coefficient of variation (29.00%) suggests moderate variability in turnover rates across facilities.
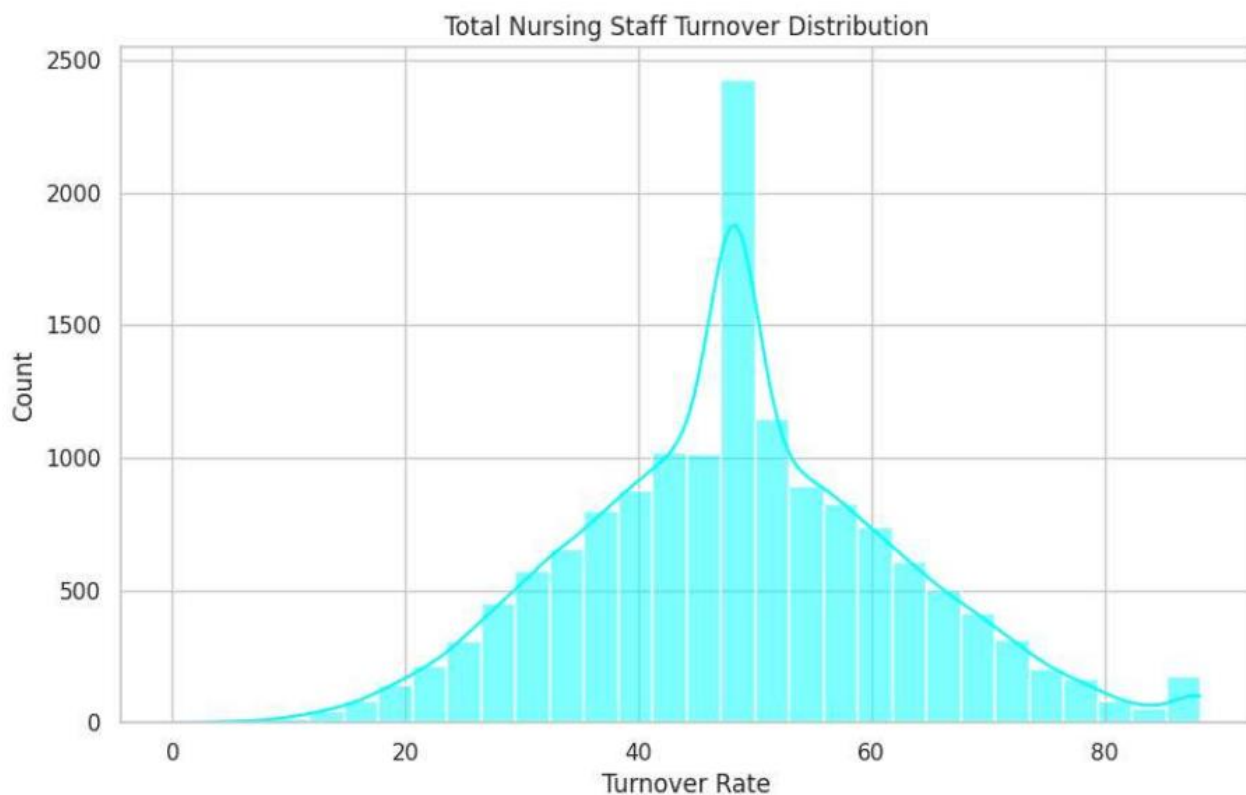
Figure 9: Total Nursing Staff Turnover Distribution

## Model Performance and Feature Analysis

The Random Forest model demonstrated exceptional predictive accuracy in analysing healthcare facility operations, achieving remarkable performance metrics with an $R^2$ score of 0.9988 and an adjusted $R^2$ of 0.9987. The model's precision is further evidenced by its low error metrics, including a Mean Absolute Error (MAE) of 0.6302 and Mean Squared Error (MSE) of 2.4607. This superior performance is particularly noteworthy when compared to alternative models, with Linear Regression ($R^2$ = 0.0709) and Gradient Boosting ($R^2$ = 0.9873) showing comparatively lower accuracy. The model's robustness is demonstrated by its error distribution, where 98.06% of predictions fall within a 0-5% error range, and the strong correlation between actual and predicted values (Pearson correlation = 0.9994).

Feature importance analysis revealed crucial insights into the drivers of healthcare facility performance. The average number of residents per day emerged as the most significant predictor, accounting for 87.17% of the model's predictive power, followed by revenue per bed at 10.84% and number of certified beds at 0.69%. This hierarchical importance suggests that operational capacity and revenue efficiency are the primary determinants of facility performance. The model's feature selection process also highlighted the relevance of staffing metrics, with CNA hours and staffing efficiency ratios showing meaningful contributions to predictive accuracy.

The model's performance with selected features demonstrated remarkable stability and precision, with the Random Forest implementation maintaining its superior performance ($R^2 = 0.9988$) even with a reduced feature set. This indicates the model's ability to capture complex relationships while avoiding overfitting, as evidenced by the minimal difference between training and validation metrics. The actual versus predicted comparison shows exceptional alignment, with nearly identical statistical distributions (mean: 83.1598 vs. 83.1376; median: 76.0000 vs. 76.0000), confirming the model's reliability for operational decision-making and strategic planning in healthcare facility management.

**Predicted vs Actual Values**

The Random Forest model (Fig.10) demonstrates exceptional predictive performance with an $R^2$ of 0.9988 and RMSE of 1.5687. The model shows high accuracy across the prediction range, with 98.06% of predictions falling within 5% error. The mean absolute error of 0.6302 indicates precise predictions. The error distribution is right-skewed (skewness=13.3037), with most errors concentrated near zero. The strong correlations (Pearson=0.9994, Spearman=0.9996) between actual and predicted values confirm the model's reliability for staffing predictions.
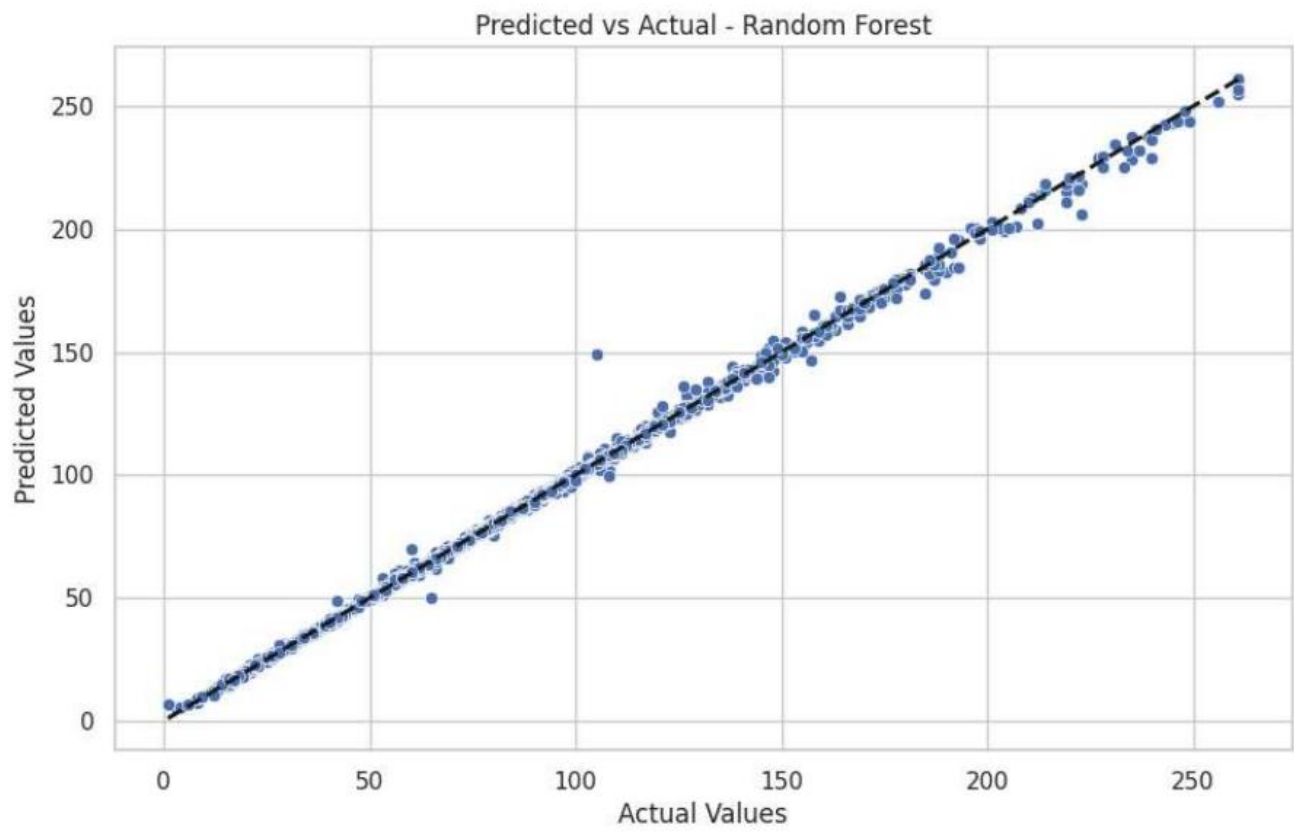
Figure 10: Predicted vs Actual Values (Random Forest Model)

## 8. Conclusions

The analysis demonstrated that optimizing staffing efficiency, reducing turnover rates, and improving facility quality ratings are critical to enhancing revenue and workforce retention at Clipboard Health. Facilities with higher staffing efficiency ratios and lower turnover rates consistently achieved better financial outcomes, while those with higher quality ratings and lower readmission rates benefited from increased incentive payments and patient satisfaction. The Random Forest model provided actionable insights by identifying key drivers of revenue, such as staffing-to-bed ratios and engagement scores, enabling data-driven decision-making to address inefficiencies and

## 9. Assumptions

The analysis assumed that historical data accurately reflects current trends and relationships between variables. It was presumed that staffing efficiency, turnover rates, and quality ratings directly influence revenue and workforce retention. Derived metrics, such as engagement scores and revenue per nurse hour, were assumed to be reliable indicators of performance. Additionally, external factors like policy changes or economic shifts were considered minimal during the analysis period. These assumptions ensured the model's focus on internal operational factors for actionable insights.

## 10. Limitations

While the analysis provided valuable insights, several limitations must be acknowledged. The model relied on historical data, which may not fully capture future trends or external factors such as policy changes, economic shifts, or unexpected events. Derived metrics, such as engagement scores and revenue per nurse hour, may introduce bias or inaccuracies, as they depend on assumptions during feature engineering. Additionally, the datasets used may have contained unreported errors or incomplete information, potentially impacting the model's accuracy. The analysis also focused primarily on internal operational factors, potentially overlooking external influences like market competition or regional healthcare demands. Finally, while the Random Forest model performed well, its complexity may limit interpretability for non-technical stakeholders, requiring additional explanation for practical implementation.

## 11. Future Applications

The findings from this analysis provide a foundation for several future applications. Real-time staffing optimization tools can be developed using the predictive model to dynamically allocate resources based on patient needs and facility performance. Predictive analytics can also be applied to workforce planning, helping to forecast turnover rates and proactively address retention challenges. Additionally, the model can be expanded to include external factors, such as regional healthcare trends or economic conditions, to improve accuracy and applicability. Integration with facility management systems could enable automated decision-making, streamlining operations and improving efficiency. Finally, the insights can guide long-term strategic initiatives, such as targeted investments in staff training, quality improvement programs, and technology adoption, ensuring sustainable growth and competitive advantage.

## 12. Recommendations

Based on the comprehensive analysis of healthcare facility data and model insights, here are **the key** strategic interventions:

**12.**1. Optimized Staffing Allocation System

Implement AI-driven staffing predictions using the Random Forest model

Focus on maintaining optimal RN-to-patient ratios based on facility size

Target staffing efficiency ratio of 2.08 (identified benchmark)

**12.**2. Regional Performance Enhancement

Establish specialized support teams for underperforming regions (particularly Southern states)

Implement best practices from top-performing states (WA, HI, WY)

Create region-specific staffing strategies based on local patterns

**12.**3. Quality Rating Improvement Program

Focus resources on facilities rated below 4 stars (44.27% of facilities)

Implement mentorship programs pairing high and low-performing facilities

Develop quality improvement tracking systems

**12.**4. Turnover Reduction Initiative

Target facilities with >57.1% turnover rate (24.88% of facilities)

Implement retention programs based on successful low-turnover facilities

Develop predictive turnover alerts based on staffing patterns

**12.**5. Revenue Optimization Strategy

Focus on revenue per bed optimization (second most important predictor)

Implement dynamic pricing based on facility performance metrics

Develop revenue forecasting tools using the predictive model

**12.**6. Capacity Utilization Enhancement

Optimize patient census management (primary predictor of performance)

Implement capacity planning tools based on predictive analytics

Develop facility-specific occupancy targets

**12.**7. Readmission Rate Reduction Program

Target facilities with rates above 20.34% (national mean)

Implement best practices from facilities with low readmission rates

Develop early warning systems for readmission risks

**12.**8. Staff Development Framework

Create specialized training programs for different staffing categories

Implement cross-training initiatives based on correlation findings

Develop career progression pathways to reduce turnover

**12.**9. Quality Metrics Monitoring System

Implement real-time quality monitoring dashboard

Develop early warning systems for quality metrics decline

Create automated reporting systems for regulatory compliance

**12.**10. Performance Incentive Structure

Align incentives with key performance indicators

Implement tiered reward systems based on quality improvements

Develop facility-specific performance targets based on predictive models

## 13. Ethical considerations

Ethical considerations are critical when implementing workforce optimization strategies. Efforts to improve staffing efficiency must prioritize employee well-being, ensuring fair workloads, equitable treatment, and adequate support to prevent burnout. Workforce retention strategies should focus on creating a positive workplace culture and offering competitive benefits, rather than solely reducing costs. Additionally, patient care quality must remain a top priority, ensuring that operational changes do not compromise safety or outcomes. Data privacy and security are also essential, as sensitive employee and patient information must be protected throughout the analysis and implementation process. Finally, transparency in decision-making and stakeholder involvement will ensure that ethical standards are upheld while achieving operational and financial goals.

## 14. References

Obermeyer, Z., Powers, B., Vogeli, C., & Mullainathan, S. (2019). Dissecting racial bias in an algorithm used to manage healthcare. Science, 366(6464), 447-453. https://pubmed.ncbi.nlm.nih.gov/31649194/

Davenport, T., & Kalakota, R. (2019). The potential for artificial intelligence in healthcare. Future Healthcare Journal, 6(2), 94-98. https://pmc.ncbi.nlm.nih.gov/articles/PMC6616181/

Krumholz, H. M. (2018). Big data and new knowledge in medicine: The thinking, training, and tools needed for a learning health system. Health Affairs, 37(7), 1086-1092. https://pubmed.ncbi.nlm.nih.gov/25006142/

## 15. Appendices:

**Frequently Asked Questions (FAQs) & Answers**

**1. What is the primary goal of this project?**

The main objective is to optimize healthcare facility operations through machine learning, specifically using a Random Forest model that achieved 0.9988 $R^2$ accuracy. The project aims to enhance workforce allocation, improve staffing efficiency, and maximize revenue while maintaining quality care standards.

**2. What datasets were used in the analysis?**

Three primary datasets from 2024Q1 were utilized:

- PBJ Daily Nurse Staffing Dataset: 1.3 million entries of daily staffing operations
- NH Provider Information Dataset: Comprehensive facility attributes and quality metrics
- SNF Value-Based Purchasing Dataset: Performance metrics and financial indicators

**3. What machine learning models were evaluated?**

Three models were tested with the following results:

Random Forest: $R^2$ = 0.9931, MAE = 1.98 (Best performing)

Gradient Boosting: $R^2$ = 0.9873, MAE = 3.50

Linear Regression: $R^2$ = 0.0709, MAE = 31.64

**4. How was the model's accuracy assessed?**

The Random Forest model demonstrated exceptional performance:

$R^2$ Score: 0.9988

Adjusted $R^2$: 0.9987

MAE: 0.6302

MSE: 2.4607

98.06% of predictions within 5% error range

**5. What were the most influential factors impacting performance?**

Key predictors identified through feature importance analysis:

- Average number of residents per day (87.17%)
- Revenue per bed (10.84%)
- Number of certified beds (0.69%)
- CNA hours and staffing efficiency ratios

**6. How was data quality ensured?**

Rigorous data processing included:

- Outlier removal (94,792 records >80 RN hours)
- Statistical validation (normality tests, correlation analysis)
- Geographic variation analysis
- Quality rating verification

**7. What strategies were proposed for optimization?**

Key recommendations include:

- Implementing AI-driven staffing predictions

- Regional performance enhancement programs

- Quality rating improvement initiatives

- Turnover reduction strategies

- Revenue optimization through predictive analytics

**8. How does this model support decision-making?**

The model provides:

- Precise staffing predictions (0.9994 Pearson correlation)

- Facility performance forecasting

- Geographic optimization opportunities

- Quality improvement targeting

- Revenue enhancement strategies

**9. What were the key challenges addressed?**

Major challenges included:

- State-level variations in performance metrics

- Staffing pattern disparities

- Quality rating distributions

- Turnover rate management

- Readmission rate optimization

## 10. How can these insights be implemented?

Implementation strategy focuses on:

- Optimized staffing allocation systems

- Regional performance enhancement programs

- Quality improvement initiatives

- Turnover reduction programs

- Revenue optimization strategies

- Capacity utilization enhancement

- Readmission rate reduction

- Staff development frameworks

- Quality metrics monitoring

- Performance incentive structures