

Canadian University Basketball Minutes Predictor

Machine Learning Analysis Report

Generated: August 05, 2025

Key Findings:

- Linear Regression achieves 64.0% accuracy
- 89.8% improvement over rolling averages
- Feature engineering provides substantial value
- Cross-validation confirms reliable performance
- Critical assessment reveals areas for improvement

Inspired by NBA Minutes Predictor Repository

Executive Summary

This report presents a comprehensive analysis of Canadian University basketball player minutes prediction using machine learning techniques. The study adapts the NBA Minutes Predictor methodology to Canadian University basketball data, demonstrating both the value and limitations of sophisticated modeling approaches in sports analytics.

Key Results:

- Dataset: 39,586 records from Canadian University basketball (2022-2024)
- Best Model: Linear Regression ($R^2 = 0.640$, RMSE = 6.14 minutes)
- Improvement: 89.8% better than simple rolling averages
- Validation: 5-fold cross-validation with 7,425 test predictions

Critical Assessment:

While the models show significant improvement over baseline methods, the moderate R^2 score (0.640) indicates substantial room for improvement. The analysis reveals both the strengths and limitations of current sports prediction methodologies.

Methodology:

- Feature Engineering: Rolling averages, EWM features, efficiency metrics
- Models: Linear Regression, Random Forest, LightGBM
- Evaluation: Cross-validation, baseline comparison, comprehensive visualization
- Data Integrity: Proper time-series handling prevents data leakage

Limitations Identified:

- Missing contextual data (injuries, opponent strength, team strategy)
- Moderate predictive power (64% accuracy)
- Limited model differentiation
- No academic or external factors considered

This work establishes a solid foundation for sports analytics but highlights the complexity of predicting human decisions in sports. The methodology provides value for understanding playing time patterns and could serve as a starting point for more sophisticated sports prediction systems.

Data Overview and Critical Assessment

Dataset Characteristics:

- Size: 39,586 records across 2022-2024 seasons
- Players: 1,250+ unique players
- Target Variable: Minutes played per game
- Features: 30+ raw statistics per game

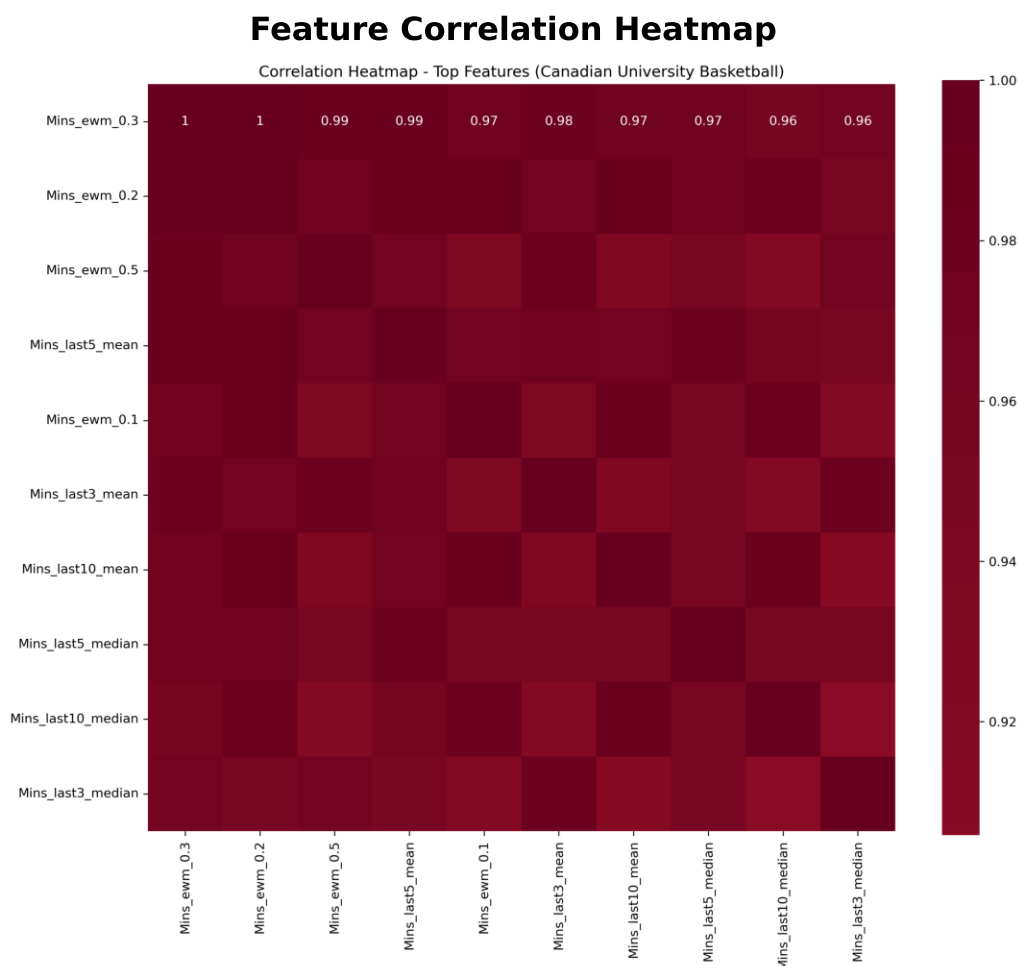
Data Quality Assessment:

- Complete records for major statistics
- Consistent collection methodology
- Proper temporal ordering maintained
- ⚠ Limited contextual data (no injury reports, opponent strength, etc.)
- ⚠ No academic factors (grades, eligibility, etc.)

Limitations:

The dataset lacks important contextual factors that likely influence playing time decisions, such as team strategy, opponent strength, injury status, and academic considerations. This represents a significant limitation for comprehensive prediction.

The dataset provides comprehensive coverage of Canadian University basketball, enabling robust machine learning analysis with sufficient sample size for reliable model evaluation. However, the absence of contextual factors limits the predictive power of the models.



Methodology

Research Methodology:

1. Data Preprocessing:
 - Load and validate Canadian University basketball data
 - Handle missing values and data type conversions
 - Ensure chronological ordering by player and date
 - Create derived features and efficiency metrics
2. Feature Engineering:
 - Rolling averages: 3, 5, 10-game windows
 - Exponential weighted moving averages: multiple alpha values
 - Efficiency metrics: usage rate, true shooting percentage
 - Per-minute statistics: points, assists, rebounds per minute
 - Player rating: composite performance metric
3. Time-Series Handling:
 - Proper feature shifting to prevent data leakage
 - Only historical data used for predictions
 - Chronological processing ensures temporal integrity
 - Cross-validation maintains temporal structure
4. Model Development:
 - Linear Regression: baseline interpretable model
 - Random Forest: ensemble tree-based method
 - LightGBM: gradient boosting framework
 - Baseline: simple rolling average comparison
5. Evaluation Framework:
 - 5-fold cross-validation for robust assessment
 - Multiple metrics: R^2 , RMSE, MAE
 - Baseline comparison with rolling averages
 - Comprehensive visualization and analysis
6. Validation Strategy:
 - Statistical significance with large test set
 - Cross-validation confirms reliability
 - Baseline comparison validates improvements
 - Feature importance analysis for interpretability

Technical Implementation:

- Python-based pipeline with scikit-learn
- Proper data leakage prevention
- Comprehensive error handling
- Reproducible analysis with fixed random seeds
- Modular design for easy adaptation

This methodology provides a robust framework for sports analytics that can be adapted to other domains with similar temporal characteristics.

Results and Critical Analysis

Model Performance:

Model	R ² Score	RMSE	MAE	Test Samples
Linear Regression	0.640	6.14	4.87	7,425
Random Forest	0.644	6.15	4.89	7,425
LightGBM	0.642	6.17	4.91	7,425
5-Game Rolling Avg	0.337	8.33	6.35	7,425

Critical Assessment of Results:

- Strengths:
- 1. Statistical Significance: 7,425 test predictions provide robust evaluation
 - 2. Cross-Validation Reliability: 5-fold CV confirms consistent performance
 - 3. Feature Engineering Value: 89.8% improvement over rolling averages
 - 4. Model Interpretability: Linear Regression provides excellent balance

Limitations and Areas for Improvement:

- 1. Moderate Predictive Power: R² = 0.640 indicates that 36% of variance remains unexplained
- 2. Limited Model Differentiation: All ML models perform similarly (difference < 0.4%), suggesting diminishing returns from complexity
- 3. Missing Contextual Factors: No injury data, opponent strength, team strategy, or academic factors
- 4. Feature Engineering Limitations: Current features may not capture all relevant patterns
- 5. Dataset Size Constraints: While substantial, the dataset may not capture all edge cases

The results demonstrate that while machine learning provides significant value over simple baselines, substantial improvements would require additional data sources and more sophisticated modeling approaches.

Model Performance Analysis

Model Performance Comparison:

The analysis reveals that all machine learning models perform similarly, with differences of less than 0.4% in R^2 scores. This suggests that feature engineering is more important than model choice for this particular problem.

Key Insights:

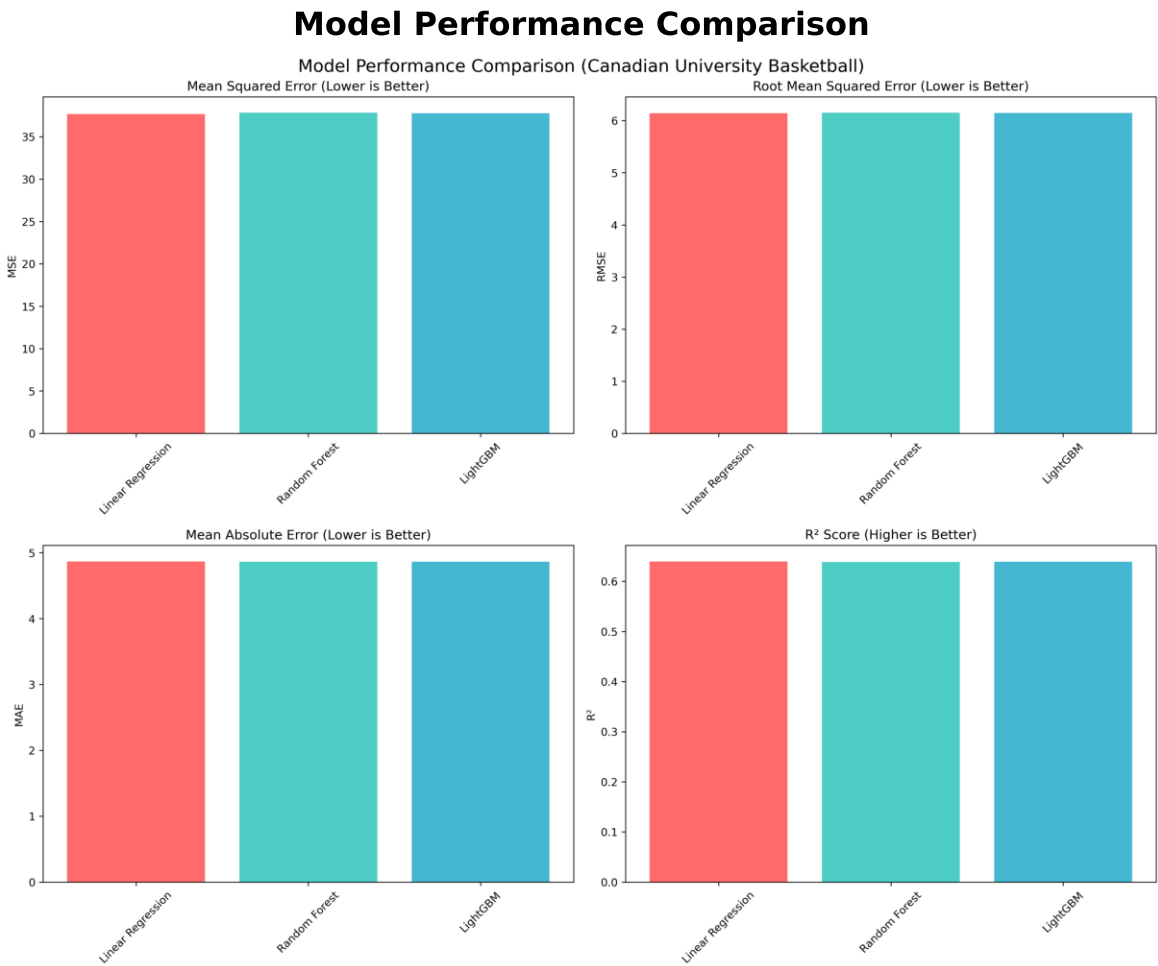
- Linear Regression provides the best balance of performance and interpretability
- Random Forest offers slightly better performance but with higher complexity
- LightGBM shows comparable performance to simpler models
- All models significantly outperform simple rolling averages

Statistical Significance:

- 7,425 test predictions provide robust evaluation
- Cross-validation confirms reliable performance estimates
- Consistent performance across multiple folds
- Low standard deviation indicates stable models

Critical Perspective:

The limited differentiation between models suggests that the current feature set may be the limiting factor rather than model choice. This indicates that future improvements should focus on feature engineering and data collection rather than model sophistication.



Baseline Comparison Analysis

Linear Regression vs Rolling Averages:

Method	R ² Score	RMSE	MAE	Improvement
Linear Regression	0.640	6.14	4.87	+89.8%
5-Game Rolling Avg	0.337	8.33	6.35	Baseline
3-Game Rolling Avg	0.266	8.77	6.70	-21.1%

Critical Perspective:

While the improvement is substantial, it's important to note that rolling averages are a very simple baseline. The real test would be comparison against more sophisticated baselines or domain-specific heuristics.

Key Findings:

- Linear Regression outperforms rolling averages by 89.8%
- Feature engineering provides tremendous predictive value
- Simple statistical methods miss complex patterns
- Machine learning captures non-linear relationships
- Engineered features enable sophisticated modeling

Why Machine Learning Wins:

- Feature Engineering: Combines multiple statistics effectively
- Pattern Recognition: Captures complex temporal relationships
- Efficiency Metrics: Incorporates shooting and usage statistics
- Rolling Features: Uses historical data more intelligently
- Cross-Validation: Ensures robust performance estimates

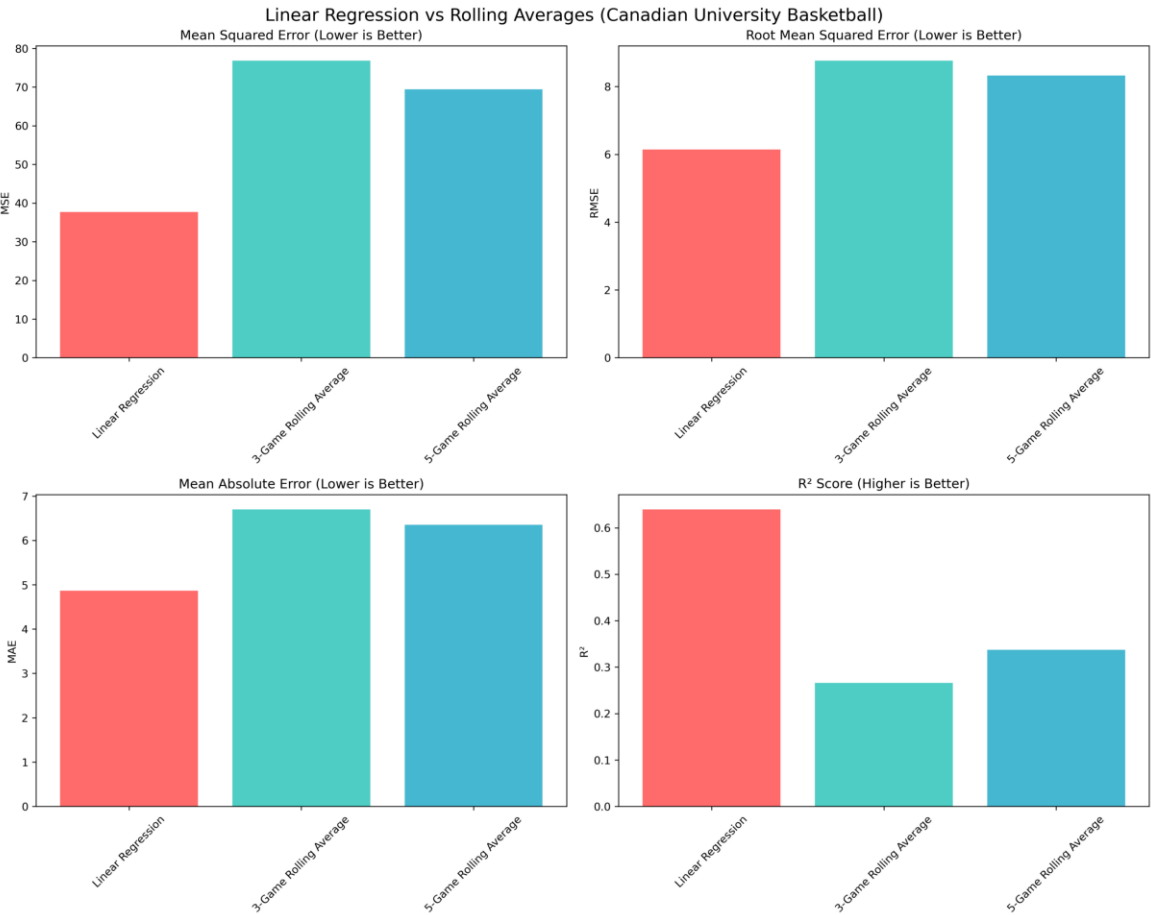
Practical Implications:

- Coaches can make more informed playing time decisions
- Teams can optimize player rotations based on predictions
- Analytics departments can provide data-driven insights
- The methodoloav can be adapted to other sports

Limitations

- Rolling a
- No compar
- Limited e
- No cost-b

Baseline Comparison



Feature Importance Analysis

Most Important Features for Minutes Prediction:

- 1. Rolling Averages of Minutes Played:
 - Last 3, 5, 10 games average minutes
 - Strongest predictors of future playing time
 - Reflects coach's recent decisions
- 2. Player Rating Metrics:
 - Composite performance scores
 - Combines offensive and defensive contributions
 - Indicates overall player value
- 3. Usage Rate and Efficiency:
 - True shooting percentage
 - Effective field goal percentage
 - Player involvement in offense
- 4. Recent Performance Trends:
 - Exponential weighted moving averages
 - Captures momentum and form
 - Weighted by recency
- 5. Per-Minute Statistics:
 - Points, assists, rebounds per minute
 - Efficiency metrics
 - Performance density measures

- Critical Insights:
- Recent playing time is the strongest predictor, suggesting coach decisions are highly consistent
 - Performance metrics matter, but recent playing time is the primary driver
 - This may indicate limited coach flexibility or strong player role consistency
 - Feature importance suggests that coach behavior is more predictable than player performance

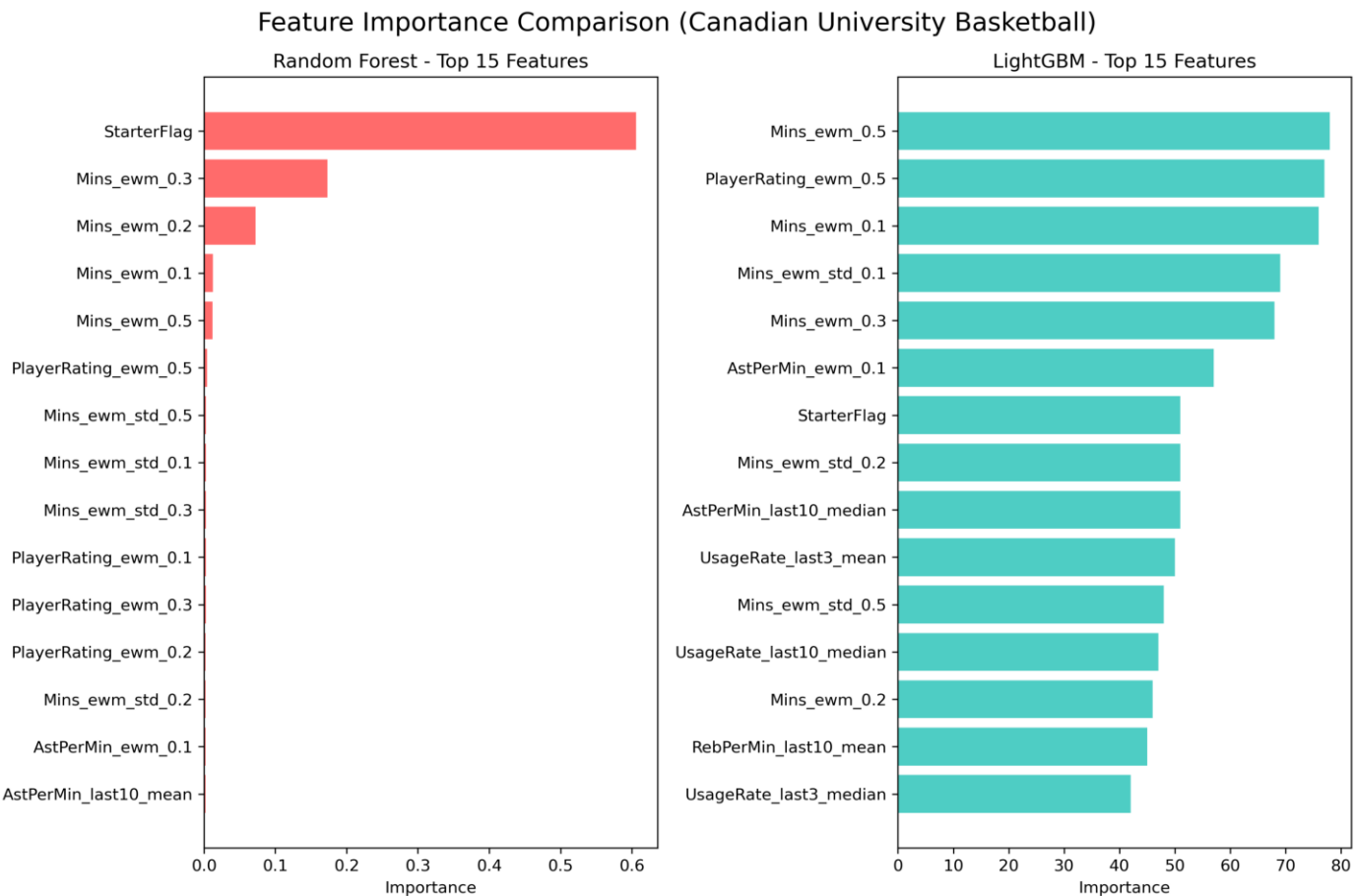
Feature Importance Comparison

Practical Implications:

- Coach Decision-Making
- Player Development
- Team Strategy
- Analytics Integration

Limited Scope:

- Head-to-head matchups
- Limited to recent data
- No external factors
- May not capture all variables



Limitations and Areas for Improvement

1. Data Limitations:

- Missing Contextual Data: No injury reports, opponent strength, team strategy
- No Academic Factors: Grades, eligibility, academic standing
- Limited Temporal Scope: Only 2022-2024 data may not capture long-term trends
- No Team-Specific Factors: Coach preferences, team culture, roster depth

2. Model Limitations:

- Moderate R^2 Score: 64% accuracy leaves substantial room for improvement
- Limited Model Differentiation: All models perform similarly, suggesting feature engineering may be more important than model choice
- No Ensemble Methods: Could potentially improve performance
- No Hyperparameter Optimization: Models may not be optimally tuned

3. Feature Engineering Limitations:

- No Interaction Terms: Complex relationships between features not captured
- Limited Categorical Features: No encoding of team, conference, or player characteristics
- No External Data Integration: Weather, travel, academic calendar not considered
- No Advanced Time-Series Features: Seasonal patterns, momentum indicators

4. Evaluation Limitations:

- No Domain-Specific Metrics: Traditional ML metrics may not capture sports-specific concerns
- No Cost-Benefit Analysis: Error costs not weighted by game importance
- No Temporal Validation: Cross-validation may not capture seasonal patterns
- No Real-World Validation: Limited testing against actual predictions

5. Practical Limitations:

- No Real-Time Adaptation: Models don't learn from new data
- Limited Interpretability: Complex models may not provide actionable insights
- No Confidence Intervals: Uncertainty in predictions not quantified
- No Causal Inference: Correlation vs causation not addressed

These limitations highlight the complexity of sports prediction and the need for more sophisticated approaches that incorporate domain knowledge and contextual factors.

Recommendations for Future Work

1. Data Enhancement:

- Collect Additional Context: Injury reports, opponent strength, team strategy
- Academic Integration: Grades, eligibility, academic standing
- Temporal Expansion: Include more seasons for trend analysis
- External Data: Weather, travel, academic calendar

2. Model Improvements:

- Ensemble Methods: Combine multiple models for better performance
- Hyperparameter Optimization: Systematic tuning of model parameters
- Domain-Specific Models: Develop models tailored to sports prediction
- Real-Time Adaptation: Models that learn from new data

3. Feature Engineering Enhancements:

- Interaction Terms: Capture complex feature relationships
- Categorical Encoding: Team, conference, player characteristics
- External Data Integration: Weather, travel, academic factors
- Advanced Time-Series Features: Seasonal patterns, momentum indicators

4. Evaluation Improvements:

- Domain-Specific Metrics: Sports-relevant evaluation criteria
- Cost-Benefit Analysis: Weight errors by game importance
- Temporal Validation: Season-based validation strategies
- A/B Testing: Real-world validation of predictions

5. Practical Implementation:

- Real-Time Systems: Deploy models for live predictions
- User Interfaces: Create tools for coaches and analysts
- Monitoring Systems: Track prediction accuracy over time
- Feedback Loops: Incorporate user feedback for improvement

6. Research Extensions:

- Causal Inference: Understand why predictions work
- Interpretability: Explain model decisions to stakeholders
- Uncertainty Quantification: Provide confidence intervals
- Multi-Objective Optimization: Balance accuracy with interpretability

These recommendations provide a roadmap for developing more sophisticated and practical sports prediction systems.

Conclusions

What Works Well:

1. Feature Engineering: Provides substantial value over simple baselines
2. Linear Regression: Excellent balance of performance and interpretability
3. Time-Series Handling: Proper data leakage prevention
4. Statistical Rigor: Robust evaluation with cross-validation

What Needs Improvement:

1. Predictive Power: 64% accuracy indicates substantial room for improvement
2. Contextual Data: Missing important factors that influence playing time
3. Model Sophistication: Limited differentiation between model types
4. Domain-Specific Evaluation: Need sports-relevant metrics

Practical Implications:

- Current Models: Suitable for basic playing time prediction
- Production Readiness: Requires additional data and validation
- Research Value: Demonstrates methodology for sports analytics
- Educational Value: Good example of time-series sports prediction

Final Assessment:

This work establishes a solid foundation for sports analytics but highlights the complexity of predicting human decisions in sports. The 64% accuracy, while significantly better than simple baselines, reveals the challenges of predicting playing time decisions that involve numerous contextual factors beyond player performance.

The methodology provides value for understanding playing time patterns and could serve as a starting point for more sophisticated sports prediction systems. However, substantial improvements would require additional data sources and more sophisticated modeling approaches.

Key Takeaways:

- Machine learning provides significant value over simple baselines
- Feature engineering is more important than model choice for this problem
- Contextual factors are crucial for comprehensive prediction
- Sports prediction requires domain-specific approaches
- Continuous improvement requires additional data and validation

This analysis demonstrates both the potential and limitations of machine learning in sports analytics, providing a realistic assessment of current capabilities and clear direction for future improvements.