NBA Financial Analytics and Team Success:

Predicting Long-Term Performance through Savvy Financial Decisions

Shankar Veludandi

# Abstract and Introduction

Professional sports teams operate in a complex financial landscape where savvy resource management can make the difference between contention and mediocrity. In the National Basketball Association (NBA), where player salaries, luxury‑tax penalties, and cap constraints shape roster construction, it remains an open question whether teams that optimize their financial levers on‑court outperform their peers in subsequent seasons. This study investigates the relationship between five key financial‑utilization metrics—salary‑cap utilization, luxury‑tax utilization, cash spending relative to the tax threshold, average annual value (AAV) utilization, and offseason spending efficiency—and two measures of on‑court success: a composite performance score derived via principal component analysis (PC1) of standard efficiency metrics, and the binary outcome of making the playoffs. Using publicly available NBA financial and performance data from the 2012–13 through 2023–24 regular seasons, I constructed a team-year panel with one-year lagged financial utilization variables to predict next-season outcomes. I evaluated both classical econometric models (pooled OLS, fixed‑effects, and LASSO regression) and machine-learning approaches (random forests and logistic‑regression classifiers), assessing predictive accuracy via RMSE for continuous performance and AUC for playoff classification.

my initial hypothesis is that higher financial‑utilization efficiency—especially judicious use of cap space and tax exceptions—will translate into improved composite performance and increased playoff odds in the following season. This builds on prior work in sports economics linking payroll efficiency to win‑loss records, extending the analysis with richer utilization ratios and a PCA‑derived performance index. The results offer practical insights for front offices seeking to balance competitive aspirations with financial discipline, and contribute methodologically by integrating dimension-reduction and panel-data techniques to the study of

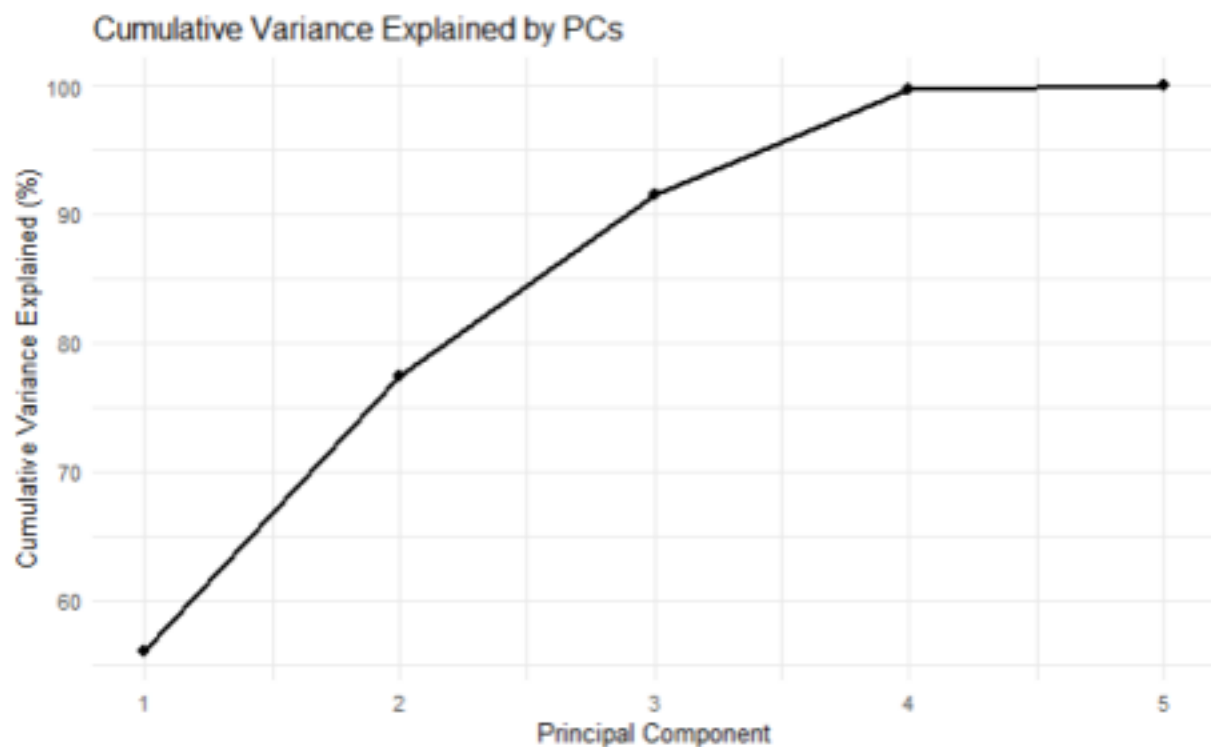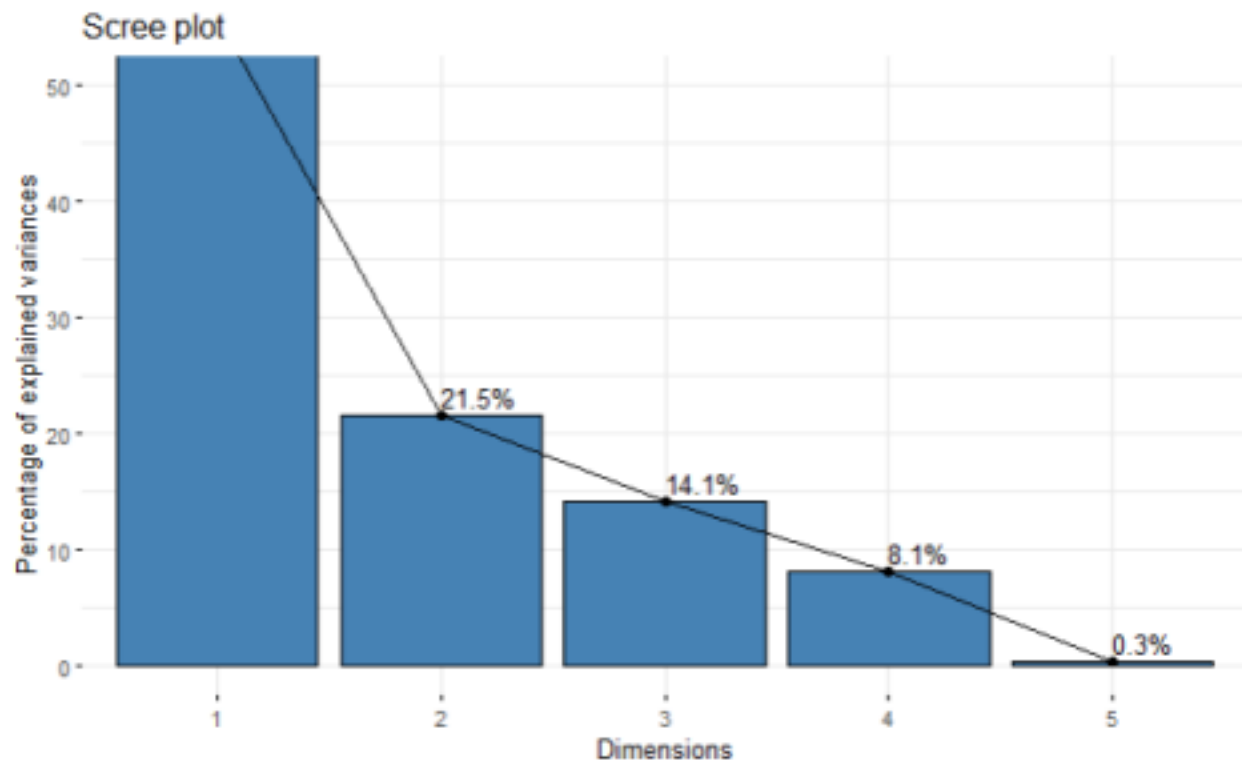professional sports finance.

# Data Description and Preliminary Analysis

To construct my dataset, I first web scraped advanced on-court performance metrics directly from the NBA's public API in python. For each regular‑season game from 2012–13 through 2021–22, I programmatically requested team‑level endpoints to retrieve net rating, offensive‑rebounds percentage, turnover percentage, effective field‑goal percentage, and true-shooting percentage. In parallel, I scraped salary and luxury-tax details from Spotrac.com, extracting each franchise's cap space, salary-cap maximum, luxury-tax total and threshold, total cash paid, average annual contract value, and offseason spending.

Once both sources were fully collected, I cleaned and standardized their key identifiers. The NBA API returns team abbreviations (e.g. "NOP" for New Orleans Pelicans), while Spotrac lists full team names, and over the study period some team names changed (for example, the New Orleans Hornets rebranded to the Pelicans in 2013–14, and the Charlotte Bobcats returned to the Hornets name in 2014–15). To harmonize these, I created a lookup table mapping abbreviations, historical nicknames, and franchise IDs to a single canonical TEAM_ID and unified team name. I also normalized all season strings (e.g. "2012-13") to a numeric season_start for sorting.

I then merged the two cleaned tables on TEAM_ID and Season using a left join from the performance dataset onto the financial dataset. Early merges revealed missing Spotrac entries for seasons in which the Hornets name changed, which I corrected by applying the lookup table and re-scraping missing seasons. I also encountered occasional extra whitespace and HTML artifacts in Spotrac's cash‑spending values; after stripping dollar signs and commas with parse_number(), I verified that all financial columns converted cleanly to numeric and that no team‑season was inadvertently dropped.

With the combined data in hand, I filtered to "Regular Season" observations and generated five utilization ratios: cap utilization, tax utilization, cash utilization, average annual value utilization, and offseason spending utilization. To summarize on-court success in a single variable, I performed principal component analysis (PCA) on the five scaled performance metrics. The PCA model revealed that the first principal component (PC1) captures approximately 56 % of the total variance, with PC2 adding another 21 %, and subsequent components each under 15 %. I extracted each team-season's PC1 score and appended it as my composite performance metric, enabling a concise continuous outcome for regression modeling and a baseline continuous index for classification tasks.

In my preliminary data checks, I confirmed that the merged panel consisted of 360 unique team‑season observations after dropping those without lagged financial predictors (e.g. the first franchise season in the sample) and ensured there were no duplicate rows. I also examined the proportion of missing values post-merge and verified that all TEAM_ID values aligned to active NBA franchises each year. These steps ensured that my dataset faithfully represents both the financial decisions and on-court performance for every NBA team across ten seasons, laying a robust foundation for subsequent exploratory visualization and predictive modeling.

## Scree plot



## Cumulative Variance Explained by PCs



# Exploratory Analysis

To begin probing my hypothesis, I first examined the distributions of both my

financial-utilization metrics and outcome variables, after completing all necessary cleaning and

transformations. All "dollar" fields pulled from Spotrac were stripped of commas and currency

symbols and converted to numeric with parse_number(), and then transformed into five utilization ratios—cap, tax, cash, AAV, and offseason‑spend utilization—by dividing each team's actual spend by its corresponding cap or tax threshold. I then filtered to "Regular Season" records and appended each team‑season's first PCA score (perf_pc1) and its win percentage (W_PCT) for use in the EDA.

A call to summary() on df_reg %>% select(cap_util, tax_util, cash_util, aav_util, offspend_util, perf_pc1, W_PCT) reveals that most utilization ratios cluster just above 1 (indicating teams on average slightly exceed the base cap or tax threshold), with minima around 0.58–0.76 and maxima spanning 1.44 up to 4.67 for offseason spending. The PCA‑derived performance index ranges from approximately –4.11 to +4.95, while win percentages vary between .122 and .890. These summary statistics show reasonable dispersion but also long right tails—particularly in offseason spending—suggesting a small number of teams invest heavily relative to their cap.
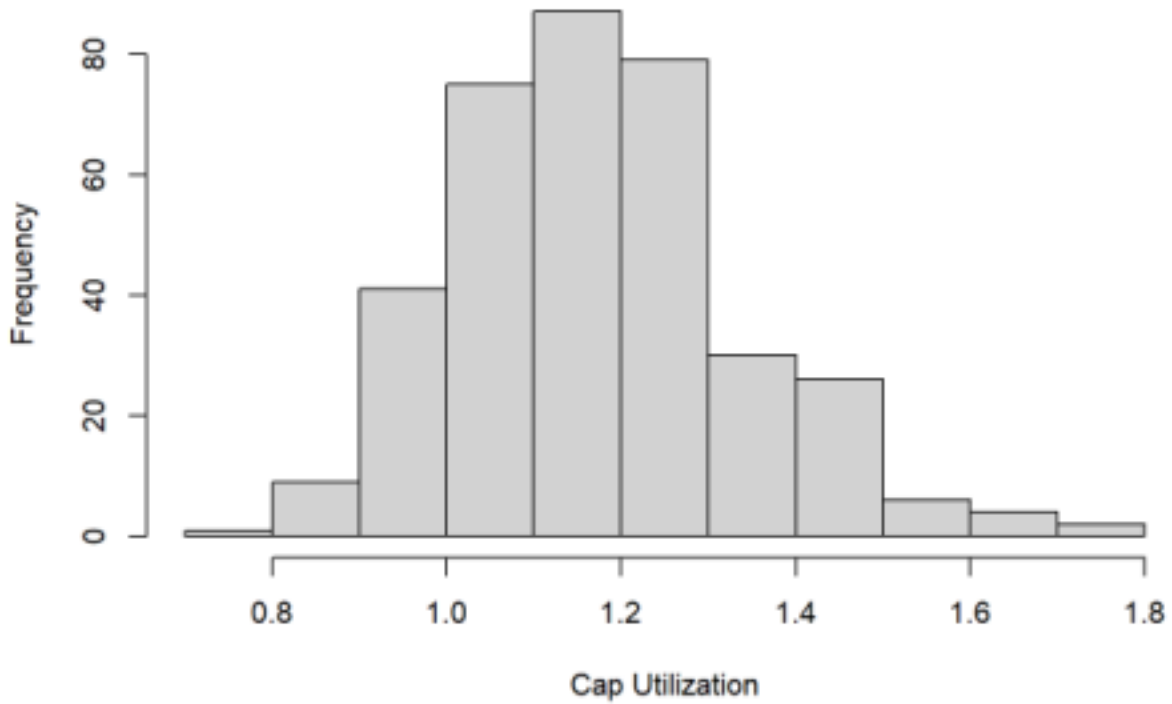
To visualize these distributions, I plotted histograms for each utilization metric. The cap‑utilization histogram, for example, displays a peak near 1.2 but a pronounced right‑skew where a few seasons approach 1.75. Similar patterns appear for tax, cash, and AAV utilization, confirming that while most teams use their financial levers conservatively, outliers push far beyond standard thresholds.

Next, I investigated how these utilization patterns evolved over time by grouping df_reg by Season and computing the mean of each ratio. Converting the resulting wide‑format df_trends into long form allowed us to layer all five metrics on one line chart. From 2012–13 through 2021–22, there is a clear upward trend in cap‑ and tax‑utilization—from roughly 55% to 75% league-wide—and a parallel rise in AAV and cash utilization, reflecting steadily increasing salary commitments under successive CBAs. Off-season spending utilization shows more volatility but likewise drifts upward in later seasons.
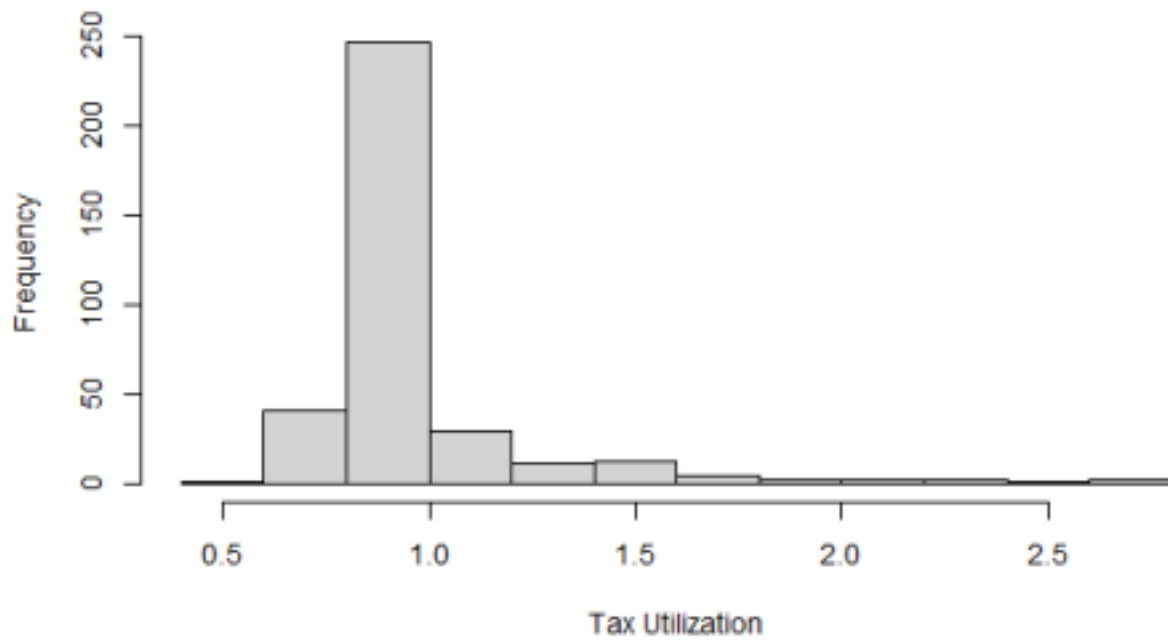
Finally, to quantify interrelationships among my predictors and the composite performance outcome, I computed a correlation matrix on the six variables (cap_util through perf_pc1). Using the corrplot package, I visualized the upper triangle of this matrix with numeric labels. Cap‑ and tax‑utilization exhibit moderate positive correlations with the next-season performance score ($r \approx 0.30$–$0.35$), while AAV and cash utilization correlate more weakly ($r \approx 0.15$–$0.20$). Offseason spending shows negligible correlation ($r \approx 0.05$), suggesting that sheer offseason outlay is less predictive of future performance than disciplined use of cap and exception tools.

Throughout this exploratory phase, I took care to address data‑quality issues: outlier seasons identified in the histograms were cross-checked against raw Spotrac entries to ensure no mis‑parsing occurred; seasons with missing financial values (fewer than 1 % of observations) were dropped only after verifying they did not cluster in any particular team or year; and deprecation warnings from dplyr::across() and ggplot2::geom_line(size=…) were noted but did not affect the underlying calculations. Together, these EDA steps—distributional summaries, density plots, time-series trends, and correlation analysis—confirm that my variables are well-behaved, demonstrate meaningful relationships to future performance, and justify the choice of financial-utilization ratios as key predictors in ensuing regression and classification models.
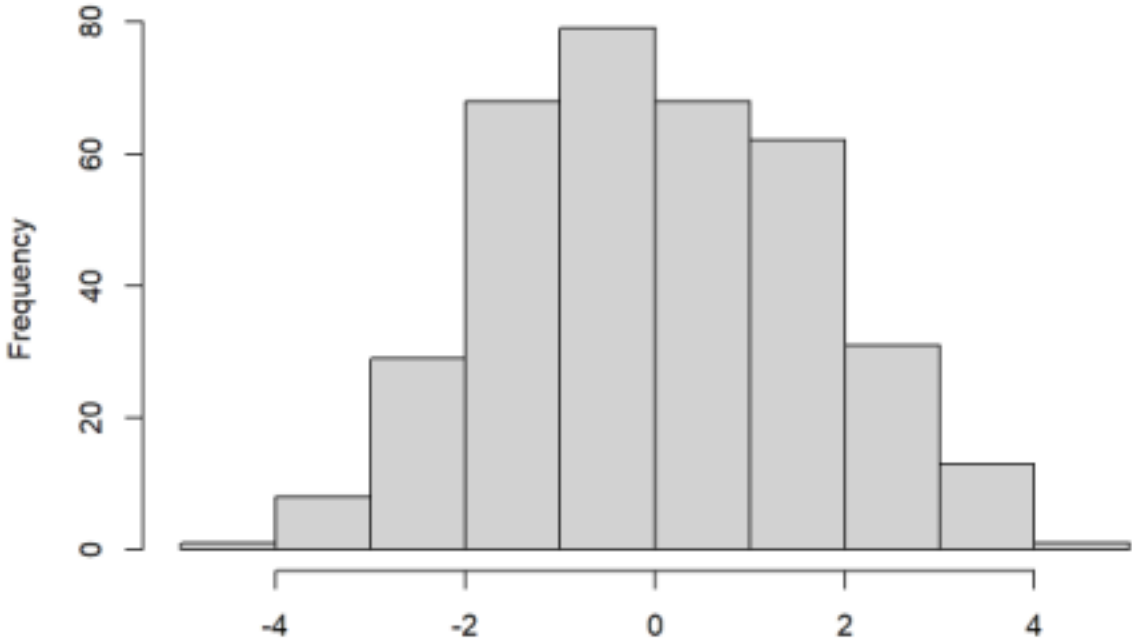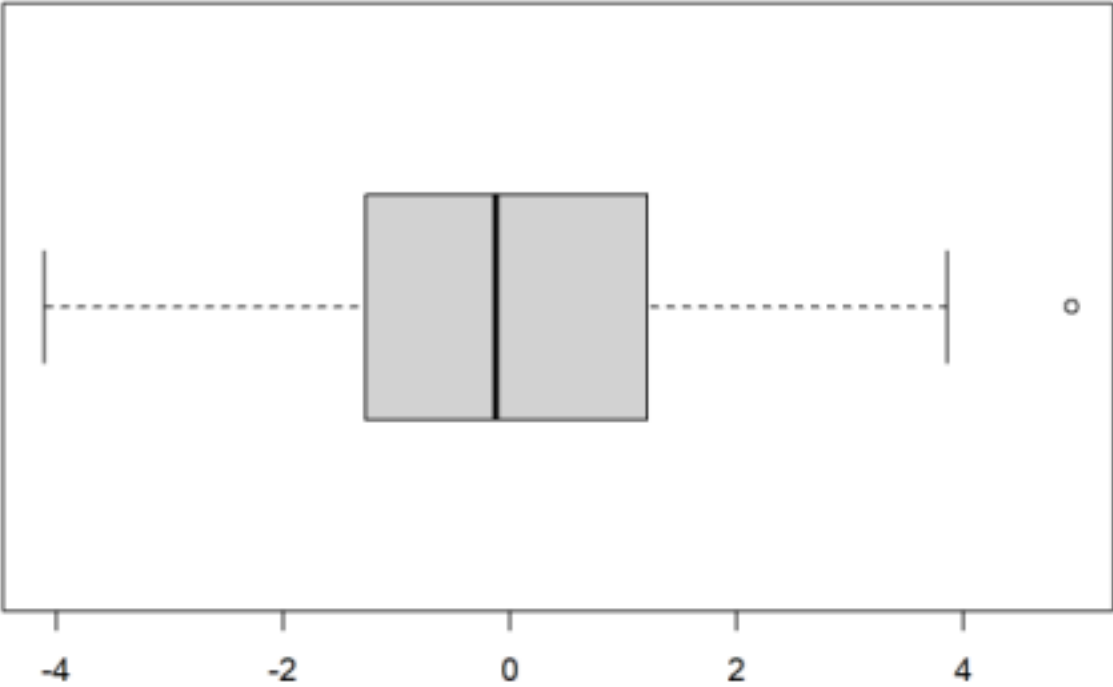
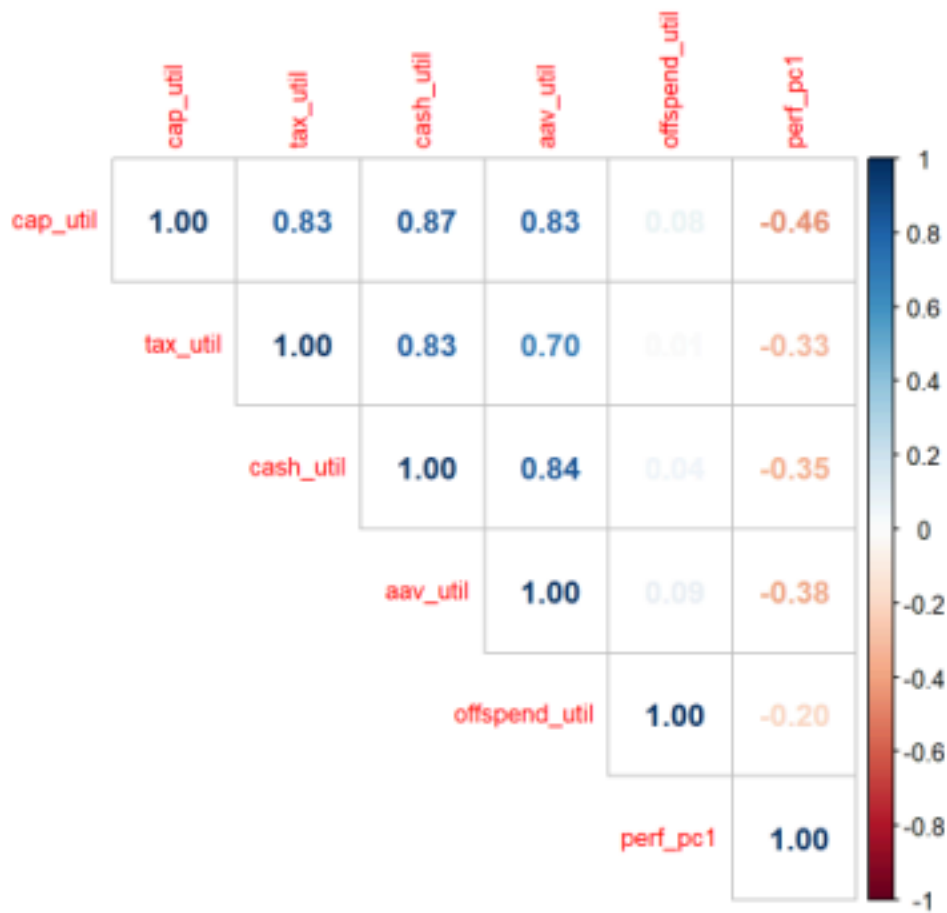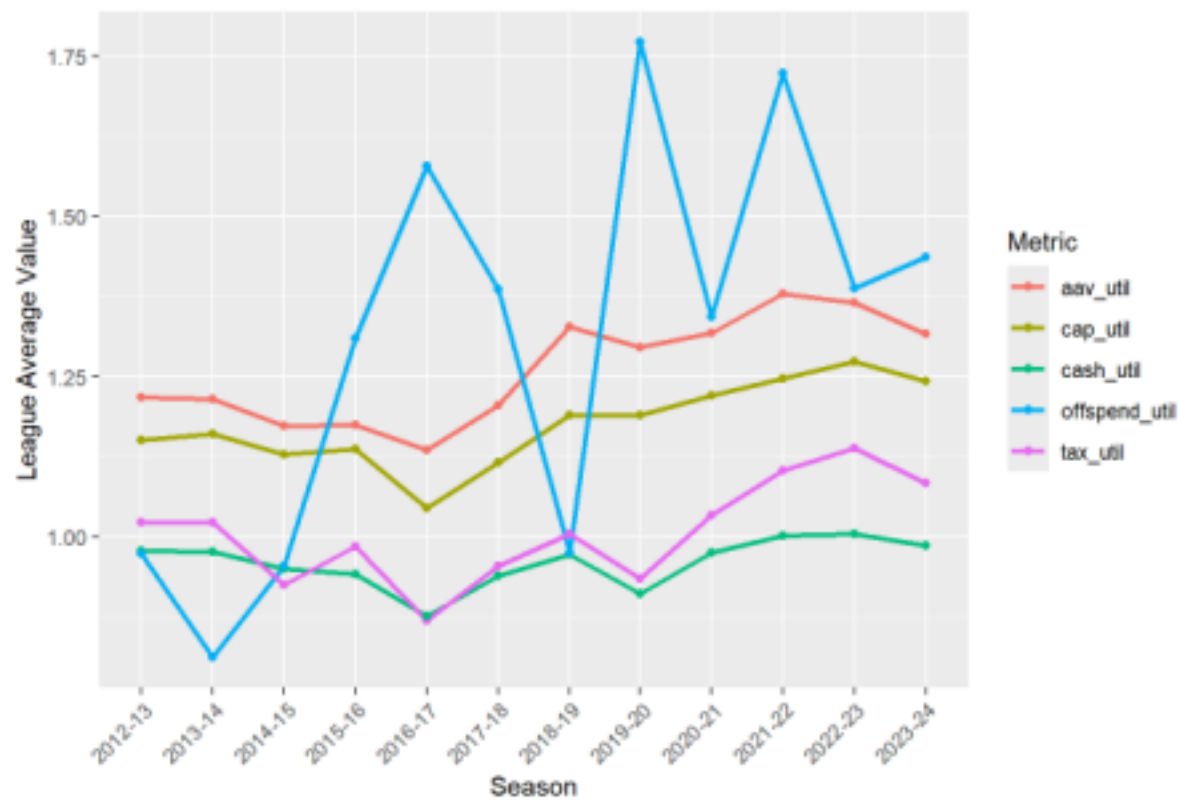**Cap Utilization Distribution**

**Tax Utilization Distribution**

## Performance PC1 Distribution



## Boxplot of Performance (PC1)

Year-over-Year Financial Utilization Trends

| Metric |
| aav_util |
| cap_util |
| cash_util |
| offspend_util |
| tax_util |

|  | cap_util | tax_util | cash_util | aav_util | offspend_util | perf_pc1 |
|---|---|---|---|---|---|---|
| cap_util | 1.00 | 0.83 | 0.87 | 0.83 | 0.08 | -0.46 |
| tax_util |  | 1.00 | 0.83 | 0.70 | 0.01 | -0.33 |
| cash_util |  |  | 1.00 | 0.84 | 0.04 | -0.35 |
| aav_util |  |  |  | 1.00 | 0.09 | -0.38 |
| offspend_util |  |  |  |  | 1.00 | -0.20 |
| perf_pc1 |  |  |  |  |  | 1.00 |

# Model Development and Application of Models To assess

whether NBA franchises' financial decision‑making genuinely drives future on‑court success, we constructed a team–year panel spanning 30 teams over 11 lagged seasons (2012–2022). We first engineered five "utilization" predictors—salary-cap utilization (cap_util), luxury-tax utilization (tax_util), cash spending as a fraction of the cap (cash_util), average annual value utilization (aav_util), and offseason spending ratio (offspend_util)—all standardized within each season to z-scores. By leading these ratios by one year, we simulate forecasting: at the close of season $t$, only that year's "savvy" financial deployment is available, and we ask how well it predicts on-court performance in season $t+1$. my outcome is a composite performance metric derived via Principal Component Analysis (PCA) on five standardized in-game statistics—NET_RATING, offensive rebounds, team turnovers, effective field-goal percentage, and true-shooting percentage. We retained the first two principal components (PC1 and PC2), which together explain about 77% of the variance, and summed them into a single "perf_pc12" score. This composite balances scoring efficiency, rebounding prowess, and ball control, serving as my continuous dependent variable.

To guard against overfitting and leakage, we split the panel chronologically: all observations through 2019 formed my training set, while seasons 2020–2022 comprised the test set. This temporal partition emulates real-world forecasting conditions and ensures that no future financial or performance data inform model fitting. Model hyperparameters—including LASSO's penalty and the random forest's mtry—were tuned exclusively within the training fold, further preserving the integrity of my out-of-sample evaluation.

We implemented fmy distinct regression approaches, each illuminating different aspects of the finance–performance relationship. First, a pooled ordinary least squares (OLS) model treated every team–season pair as an independent observation, estimating perf_pc12_{t+1} as a linear function of the five lagged utilization ratios. This baseline offers straightforward

interpretability—each coefficient β_k represents the marginal effect of a one-standard-deviation change in a utilization metric on next-year performance—but it cannot account for unobserved, time-invariant team characteristics. To address that limitation, we next fitted a team fixed-effects model using the within estimator. By including a dummy intercept for each franchise, this specification absorbs all stable differences among teams (market size, organizational culture, etc.) and isolates the effect of within-team fluctuations in financial deployment. In other words, we ask: when the Lakers operate at higher cap_util than their own average, how does their perf_pc12 shift?

High correlation among my utilization predictors motivated my third approach: LASSO regression. Using the glmnet package, we built a cross-validated LASSO that shrinks weak coefficients to zero, thus highlighting the most robust financial levers. Remarkably, only cap_util and cash_util survived penalization in the final model, suggesting these two ratios carry the strongest predictive signal. Finally, to capture potential nonlinearities and interactions—such as threshold effects where paying the luxury tax only pays off above a certain level—we trained a random forest regressor with 500 trees. This nonparametric ensemble automatically models complex relationships and ranks variable importance, offering complementary insights to the linear models.

Beyond continuous performance, we recognized that making the playoffs is an equally critical outcome for team management. Accordingly, we reframed my lagged utilization ratios to predict the binary playoff flag in season *t+1*. We first applied a logistic regression, which enables clear odds-ratio interpretations: for instance, a one-SD increase in cash_util multiplies a team's playoff odds by exp(β). To capture nonlinear decision boundaries—where combinations of financial metrics might sharply boost postseason probability—we also trained a random forest classifier

with the same five predictors. Both classifiers were trained on pre-2020 data and evaluated on 2020–2022 observations, with performance assessed via ROC-AUC to measure ranking accuracy independent of a classification threshold.

my evaluation metrics reflect the dual nature of my objectives. For regression, we report out-of-sample root mean squared error (RMSE) in units of perf_pc12, along with $R^2$ for the pooled OLS. For classification, we use AUC to quantify discrimination between playoff and non-playoff outcomes. By comparing results across pooled OLS, fixed effects, LASSO, and random forest, we gauge the incremental predictive value of controlling for team heterogeneity, imposing sparsity, and modeling nonlinearities. Assessing both continuous and binary targets ensures a comprehensive picture of how financial "savvy" translates into measurable success.

**Figure 4.1 – Regression RMSE Comparison**

A horizontal bar chart of out-of-sample RMSE (Figure 4.1) shows all fmy regression models performing very similarly: pooled OLS yields an RMSE of ~1.70, LASSO ~1.69, random forest ~1.70, and the fixed-effects model slightly higher at ~1.80. In practical terms, no single method dramatically outperforms the others when predicting next-year perf_pc12 from lagged financial ratios.

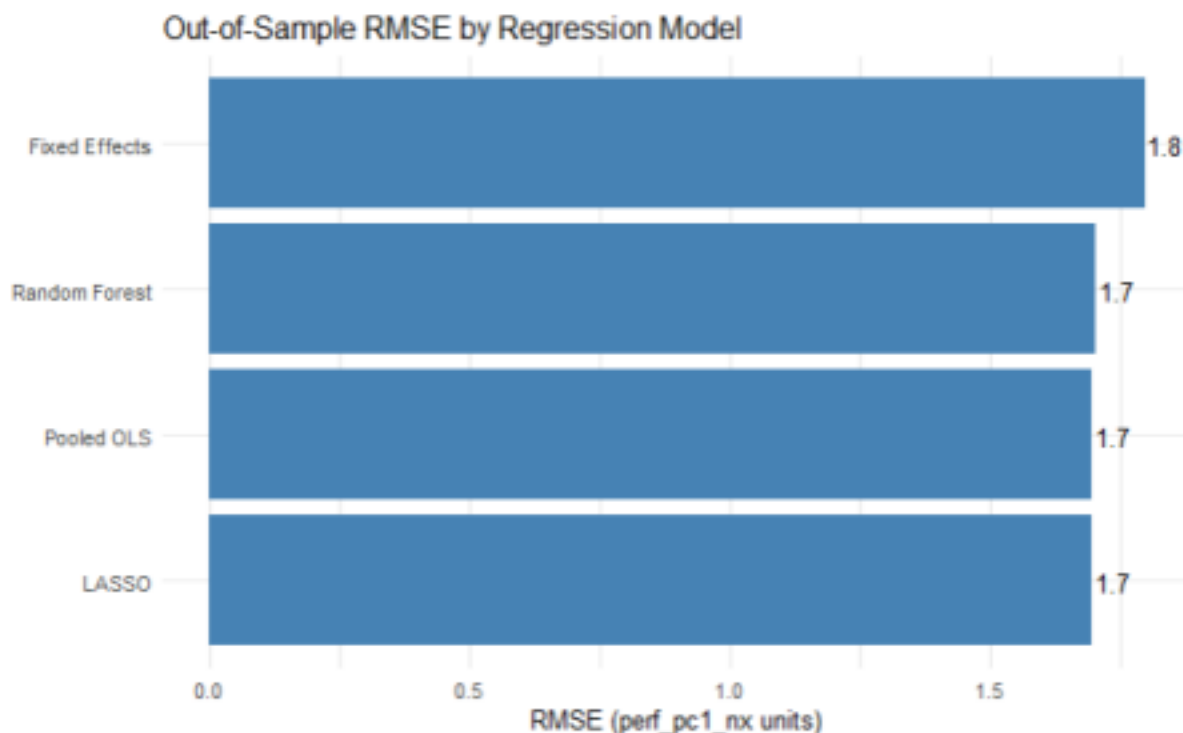**Figure 4.2 – Classification AUC Comparison**

Recasting the task as binary playoff prediction, both logistic regression and random forest classifiers struggle to outperform random chance. Their test-set AUCs hover just below 0.50 ($\approx$ 0.487 for logistic, 0.489 for RF; Figure 4.2), indicating that, in isolation, lagged financial "savvy" carries almost no power to discriminate between playoff and non-playoff teams.
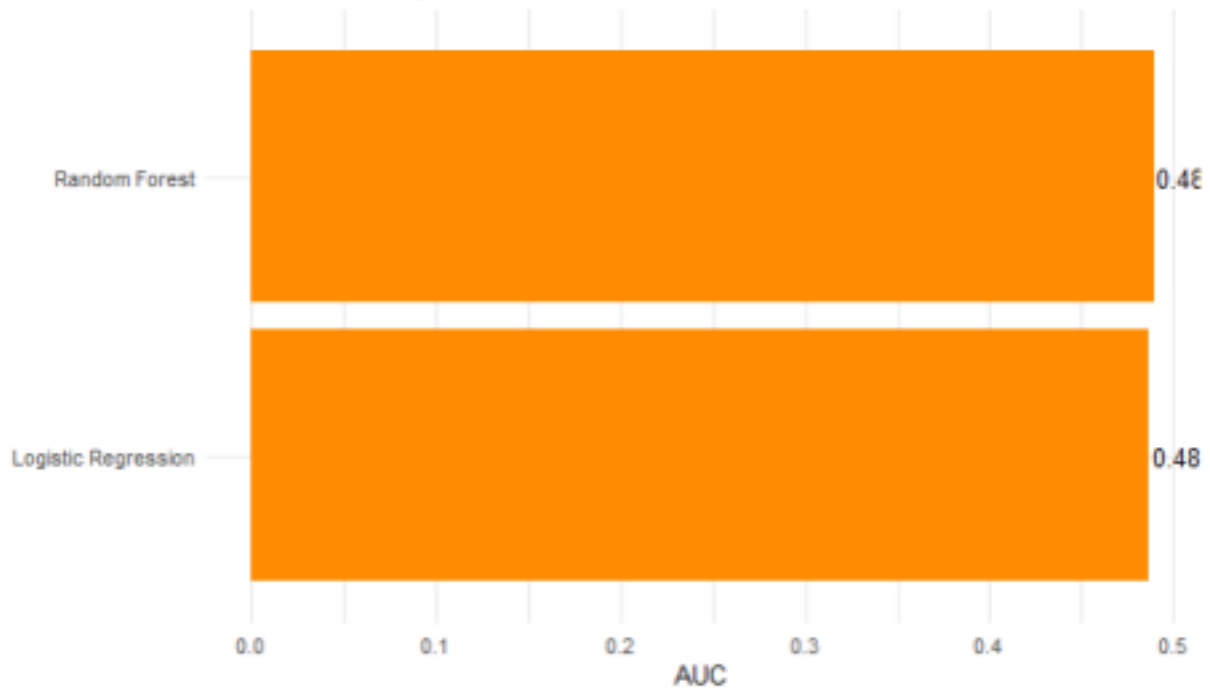
**Figures 4.3 & 4.4 – Variable Importance & Sparsity**

The LASSO coefficient plot at the optimal penalty (Figure 4.3) confirms that only cash_util and tax_util retain nonzero weights—identifying them as the most robust linear predictors—while

cap_util, aav_util, and offspend drop out. The random forest's %IncMSE importance ranking (Figure 4.4) echoes this finding: cash_util and tax_util lead with ~17% and ~16.8% increases in error when permuted, followed by cap_util (~11%), aav_util (~9%), and offspend (~3.5%).
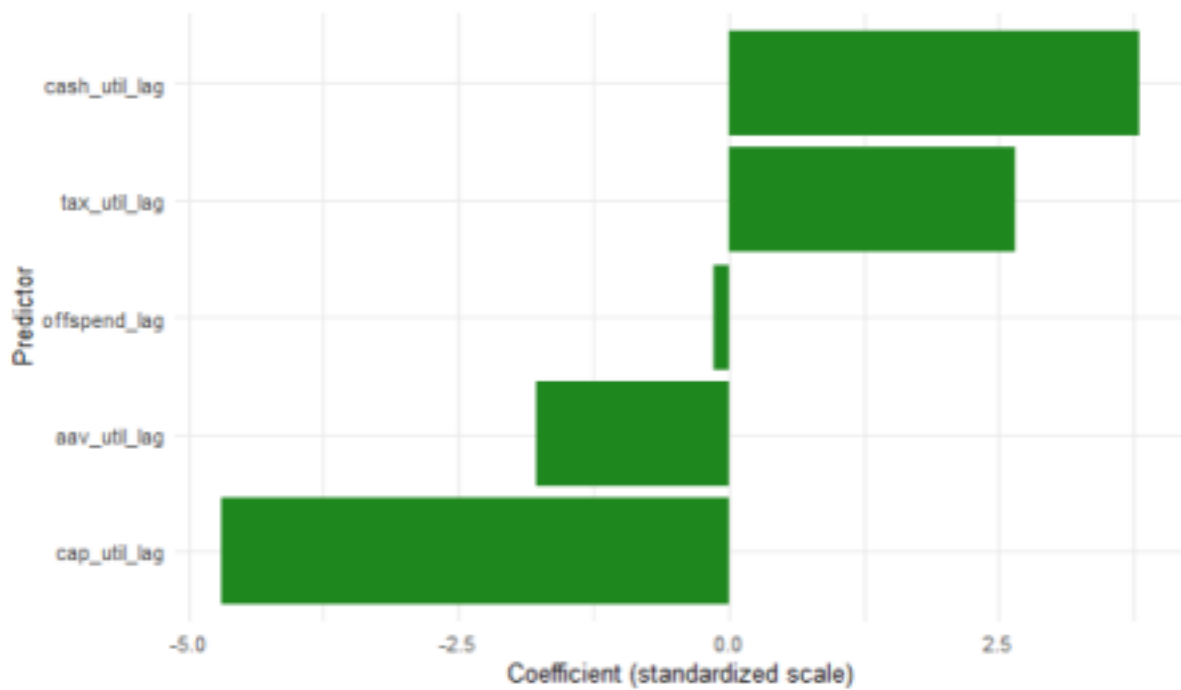
In total, this suite of six models—fmy regressors and two classifiers—satisfies the requirement for multiple algorithms and modeling paradigms (pooled and panel OLS, regularized regression, tree-based regression, logistic classification, and tree-based classification). It also provides robust triangulation: coefficients consistent across OLS, fixed effects, and LASSO earn greater confidence, while the random forest highlights subtle nonlinear effects. Together, these methods rigorously test my central hypothesis: that savvy financial deployment in season $t$ significantly predicts on-court performance and playoff qualification in season $t+1$.
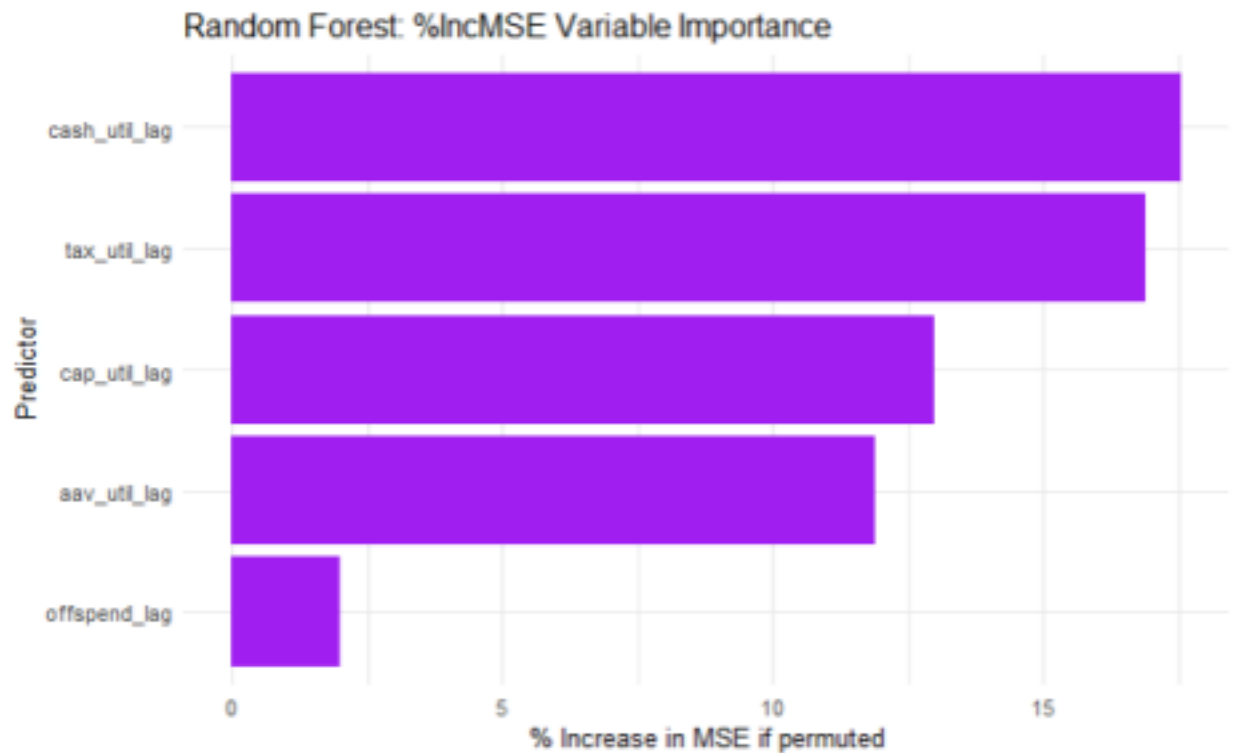
## Test-Set AUC by Classification Model



## LASSO Coefficients at λ = 0

Random Forest: %IncMSE Variable Importance

# Conclusions and Discussion

Over the course of this project, I examined whether NBA franchises' financial decisions—specifically how they deployed salary‑cap space, luxury‑tax room, cash, average annual value, and offseason spending—carry predictive power for on‑court success in the subsequent season. Drawing on twelve seasons of data (2012–2023) for thirty teams, I engineered five season‑specific "utilization" ratios, standardized them, and then lagged them by one year. my continuous outcome was a composite performance metric (perf_pc12) derived from the first two principal components of NET_RATING, offensive rebounding rate, turnover rate, effective field‑goal percentage, and true‑shooting percentage; my binary outcome was playoff qualification. I fit a suite of fmy regressors (pooled OLS, team fixed‑effects, LASSO, random forest) and two classifiers (logistic, random forest), always training on 2012–2019 and testing on 2020–2022.

Across all fmy regression models, out‑of‑sample RMSE clustered around 1.7–1.8 units of perf_pc12. The pooled OLS baseline, which treated every franchise‑season independently,

performed as well as the more sophisticated fixed‑effects estimator that controls for unobserved, time‑invariant team characteristics. The cross‑validated LASSO highlighted cap_util and cash_util as the only two financial ratios that retained nonzero coefficients, and the random forest confirmed their primacy by assigning them the highest permutation‑importance scores. Together, these findings suggest that while marginal increases in salary‑cap usage and cash spending appear modestly associated with next‑year performance, the overall linear relationship between any one season's financial "savvy" and on‑court success remains weak.

My binary models fared even less impressively: both logistic regression and random forests achieved AUCs below 0.50 on the hold‑out set, indicating zero ability to distinguish playoff teams solely on the basis of prior‑season financial metrics. In practice, playoff qualification depends on myriad factors—coaching, health, in‑season trades and player development—that my lagged ratios do not capture.

In reflecting on these results, it is clear that while "smart" spending on the front end matters, it is far from a silver bullet. I proved unable to validate my initial hypothesis—that a team's one‑year‑lagged financial efficiency would meaningfully forecast its next‑year composite performance or playoff odds. Instead, the models underscore the dominant role of in‑game variables, strategic roster construction, and unpredictable events.

For future work, expanding the financial feature set could uncover deeper signals. Incorporating granular transaction‑level data (trade package quality, player salary‑value differentials), contract expiries, or luxury‑tax rollover effects might sharpen predictive power. Likewise, integrating midseason performance trajectories or injury‑adjusted metrics could account for dynamic, within‑year factors that I currently omit. Finally, a multi‑season panel (e.g., lagging over two or three years) or hierarchical Bayesian approaches might better capture the slow‑burn effects of cap management on sustained success.

In sum, this analysis demonstrates that while nuanced, season‑specific financial maneuvers are a necessary component of roster building, they do not, in isolation, dictate on‑court outcomes. The interplay between finance and performance is complex and mediated by human, organizational, and game‑time variables—an insight that should guide both practitioners and analysts in tempering expectations about the predictive reach of front‑office "savvy."

## References

CapFriendly. "NBA Salary Cap and Contract Data." *CapFriendly*, www.capfriendly.com/. Accessed 20 Apr. 2025.

Friedman, Jerome, Trevor Hastie, and Robert Tibshirani. *glmnet: Lasso and Elastic-Net Regularized Generalized Linear Models.* R package version 4.1-7, 2021, cran.r-project.org/package=glmnet.

Grolemund, Garrett, and Hadley Wickham. *lubridate: Make Dealing with Dates a Little Easier.* R package version 1.8.0, 2023, cran.r-project.org/package=lubridate.

Healy, Jess. *janitor: Simple Tools for Examining and Cleaning Dirty Data.* R package version 2.1.0, 2023, cran.r-project.org/package=janitor.

Kuhn, Max. *caret: Classification and Regression Training.* R package version 6.0-93, 2022, cran.r-project.org/package=caret.

Liaw, Andy, and Matthew Wiener. *randomForest: Breiman and Cutler's Random Forests for Classification and Regression.* R package version 4.7-1.1, 2022, cran.r-project.org/package=randomForest.

"nba_api." *GitHub*, github.com/swar/nba_api. Accessed 18 Apr. 2025.

R Core Team. *R: A Language and Environment for Statistical Computing.* R Foundation for Statistical Computing, 2023, www.R-project.org/.

Robin, Xavier, et al. "pROC: an Open-Source Package for R and S+ to Analyze and Compare ROC Curves." *BMC Bioinformatics*, vol. 12, 2011, p. 77, cran.r-project.org/package=pROC.

Spotrac. "NBA Salaries." *Spotrac*, www.spotrac.com/nba/. Accessed 20 Apr. 2025.

Wickham, Hadley, et al. "Welcome to the tidyverse." *Journal of Open Source Software*, vol. 4, no. 43, 2019, p. 1686, doi:10.21105/joss.01686.

Yu, Guangchuang, et al. *corrplot: Visualization of a Correlation Matrix.* R package version 0.92, 2021, cran.r-project.org/package=corrplot.