# Modeling for Insight or Accuracy? Contrasting Statistical Inference and Machine Learning in Predictive Analytics

### J. Prince Vijai
Department of Operations & Information Technology
ICFAI Business School
The ICFAI Foundation for Higher Education
(Deemed-to-be-University u/s 3 of the UGC Act, 1956)
Hyderabad - 501203, Telangana, India

### R.S. Chalapathi
Department of Operations & Information Technology
ICFAI Business School
The ICFAI Foundation for Higher Education
(Deemed-to-be-University u/s 3 of the UGC Act, 1956)
Hyderabad - 501203, Telangana, India

## ABSTRACT
Predictive analytics is pivotal in shaping strategic decision-making in today's business environment. Analytical projects, however, can be broadly categorized by two distinct objectives: explanation and prediction. This paper contrasts the two primary approaches of data modeling – statistical modeling for inference and machine learning for prediction – through a practical application in marketing analytics. First, a traditional multiple linear regression model is constructed using a publicly available advertising dataset to explain the relationship between different advertising channel expenditures and sales. The model's coefficients and statistical significance are interpreted to derive actionable insights for budget allocation. Second, a suite of machine learning models is developed and evaluated to identify the most accurate predictive engine for forecasting future sales. Lastly, by direct comparison, the study recommends employing explanatory statistical models to predict unseen data and advocates evaluating models' predictive accuracy using machine learning models designed explicitly for prediction tasks. This highlights the inherent trade-off between model interpretability and predictive performance, offering practical criteria for analysts to consider when selecting the most suitable modeling approach.

## General Terms
Statistical Learning, Machine Learning, Predictive Modeling

## Keywords
Statistical Modeling, Linear Regression, Causal Inference, R-Squared, Machine Learning, Predictive Analytics, Random Forest, Prediction, RMSE

## 1. INTRODUCTION
The proliferation of data has fundamentally altered the landscape of management, creating an imperative for leaders to ground strategic decisions in empirical evidence. Within the field of data analytics, two distinct modeling approaches have emerged to serve different, though often complementary, business objectives – statistical modeling for explanation and machine learning for prediction [1]. Statistical modeling has its roots in inference, seeking to explain the relationships between variables and test hypotheses about a data-generating process. Its primary value lies in answering "why" a particular outcome occurred. In contrast, machine learning is primarily concerned with predictive accuracy, building models that can generalize from historical data to make the most accurate forecasts possible on new, unseen data. Its value lies in answering "what" will happen next [2], [3].

While these approaches often use the same algorithms, such as linear regression, their philosophies, workflows, and evaluation criteria are fundamentally different [4], [5]. This distinction is critical for predictive analytics, as the choice of methodology directly impacts the nature and utility of the resulting insights. A model designed for explanation provides the causal understanding needed for strategic resource allocation, while a model intended for prediction provides the forecasting capability required for operational planning.

This paper aims to provide a clear, practical comparison of these two modeling approaches. It presents two parallel analyses using a classic marketing dataset on advertising expenditures and sales [6]. First, a traditional statistical regression model is constructed to infer the impact of different advertising channels on sales. Second, a few machine learning models are employed to predict future sales. By comparing the results and interpretations, the study provides clear guidelines for analysts to select the appropriate tool for the business question. Furthermore, it explores the synergistic potential of combining these approaches, demonstrating how inference and prediction can be integrated to form a more comprehensive and robust business analytics framework.

## 2. RELATED LITERATURE
The distinction between statistical modeling and machine learning has been a subject of extensive discussion, often framed as a choice between two different "approaches" of data analysis [4]. This paper builds upon that discourse by examining the practical implications of these two approaches in a business context.

The first approach, statistical modeling, primarily concerns inference and explanation [6]. Rooted in probability theory, its goal is to use a sample of data to understand the underlying data-generating process and test hypotheses about it [7]. This is a "model-driven" or "top-down" approach, where the analyst assumes a specific model from which the data have been generated. This translates to identifying the causal drivers of key outcomes in business applications. For instance, regression analysis is widely used to quantify the impact of marketing expenditures on sales, optimize operational processes, and support strategic decisions with empirical evidence. The interpretability of statistical models, where each coefficient has a precise meaning, is a significant advantage, particularly when results must be explained to stakeholders. However, these inferences' validity depends on strict assumptions, such as linearity, normality, and homoscedasticity, which must be rigorously tested.

The second approach, machine learning, prioritizes prediction above all else. It treats the underlying data-generating mechanism as a "black box" and uses algorithms to find patterns that can accurately forecast future outcomes. This is a "data-driven" or "bottom-up" approach, where no particular model is assumed beforehand; instead, the algorithm develops a model with prediction as the primary goal. This approach is compelling for handling large, complex, and high-dimensional datasets where traditional statistical assumptions may not hold [6]. In demand forecasting, for example, machine learning models like decision trees, random forests, and gradient boosting often outperform traditional statistical methods by capturing complex, non-linear relationships in the data. However, this predictive power frequently comes at the cost of interpretability. Models like random forest, while highly accurate, are often opaque, making it difficult to understand the rationale behind their predictions [8].

The distinction between explanatory and predictive modeling is vital for real-world applications and decision-making, yet it is often mixed in various disciplines [4], [5], [9]. Statistical modeling is a powerful tool for developing and testing theories through causal explanation, prediction, and description. Many fields use statistical modeling for causal explanation, assuming that high explanatory power inherently implies high predictive power. While this distinction has been acknowledged in the philosophy of science, statistical literature lacks a comprehensive discussion of the differences in the modeling process for explanatory versus predictive goals. Shmueli's [1] article aims to clarify this distinction, discuss its origins, and reveal the practical implications for each step in the modeling process. Understanding these differences is essential for improving the application of statistical modeling across disciplines.

Shmueli [1] compares a conventional statistical linear regression model and a machine learning linear regression algorithm for prediction. The findings indicate that the machine learning algorithm, which learns from features within the data, produced predictions closer to the actual data than the conventional model, which relies on a fixed mathematical formula. This study builds upon the earlier comparison by evaluating a broader suite of machine learning models and explicitly testing the predictive accuracy of the statistical model on unseen data. This allows for a direct, "apples-to-apples" comparison that highlights the trade-off between the high interpretability of statistical models and the superior predictive accuracy of machine learning models. Furthermore, this study explores the growing synergy between the two fields, where hybrid approaches are being developed to leverage both strengths, combining the explanatory power of statistics with the predictive prowess of machine learning [9].

## 3. THE STATISTICAL LEARNING APPROACH

The primary goal of statistical modeling is inference – to understand and quantify the relationship between variables, answering the "why" questions that drive strategy. Which marketing channels are effective? What is the specific return on investment for each dollar spent? A multiple linear regression model is constructed using the entire Advertising.csv dataset [6] to answer these questions.

### 3.1 Model Output

A multiple linear regression model was fitted to the data to explain the relationship between advertising expenditures and sales. The resulting model is:

$$Sales = 2.939 + (0.046 * TV) + (0.189 * Radio) - (0.001 * Newspaper)$$

The detailed output of the statistical model is given below in Table 1.

**Table 1. Statistical regression model output**

| Variable | Coefficient | Std. Error | t-value | p-value |
|---|---|---|---|---|
| Intercept | 2.939 | 0.312 | 9.42 | < 0.001 |
| TV | 0.046 | 0.001 | 32.81 | < 0.001 |
| Radio | 0.189 | 0.009 | 21.89 | < 0.001 |
| Newspaper | -0.001 | 0.006 | -0.18 | 0.859 |

*Model Performance Metrics: R-squared = 0.897, F-statistic p-value < 0.001*

### 3.2 Model Insights

The R-squared value of 0.897 indicates that the statistical model explains 89.7% of the variability in sales, showing a very strong fit. The extremely low p-value for the F-statistic confirms that the model is statistically significant.

The coefficients represent the impact of each advertising channel while holding the others constant. Coefficients for TV and radio are positive and p-value at less than 0.001, making them statistically significant drivers of sales. For every additional $1,000 spent on radio, sales are predicted to increase by approximately 189 units. Meanwhile, the p-value of 0.859 is very high for the newspaper, indicating that there is no statistically significant relationship between newspaper advertising expenditure and sales in this model.

For strategic budget allocation, this model provides clear guidance. It suggests that marketing funds are most effectively spent on TV and radio, and that investment in newspaper advertising may not yield a significant return.

## 4. THE MACHINE LEARNING APPROACH

When the primary goal is to create the most accurate forecast possible, a machine learning approach is adopted. The objective here is not to interpret coefficients but to minimize prediction error on new, unseen data. This requires a different workflow and a broader set of algorithms [8].

### 4.1 Predictive Modeling

The dataset was split into a training set (80%) and a test set (20%) to evaluate predictive performance properly. The models are built on the training data, and their final accuracy is measured on the unseen test data. Several supervised learning algorithms were trained on the data, such as a statistical model, linear regression, ridge regression, lasso regression, decision tree, random forest, and gradient boosting. To ensure robustness and prevent overfitting, 10-fold cross-validation is performed on the training set for each model.

In the statistical model, the equation from the inferential model is used to make predictions on the test set to provide a direct comparison. A linear regression model is trained only on the training data. Regularized regression models, such as ridge and lasso regressions, are used to penalize large coefficients that prevent overfitting. On the other hand, a decision tree model is a non-linear model that uses if-else rules for learning. Powerful ensemble methods, such as random forest and gradient boosting, are used to combine multiple decision trees to

improve prediction accuracy and robustness.

## 4.2  Predictive Model Performance

The performance of each model is evaluated on the test dataset. The key metrics are root mean squared error (RMSE) and mean absolute error (MAE). Both measure the average prediction error, with lower values indicating a better model. MAE represents the average absolute error, while RMSE penalizes large errors more. Table 2 compares the predictive performance, based on RMSE and MAE, of all models on the test data.

**Table 2. Comparison of the predictive performance of all models**

| Model | RMSE (in Sales Units) | MAE (in Sales Units) |
|---|---|---|
| Statistical Model (for Prediction) | 1.41 | 1.22 |
| Linear Regression (ML) | 1.41 | 1.22 |
| Ridge Regression | 1.41 | 1.22 |
| Lasso Regression | 1.41 | 1.22 |
| Decision Tree | 1.89 | 1.45 |
| **Random Forest** | **0.98** | **0.81** |
| Gradient Boosting | 1.02 | 0.84 |

Figure 1 presents the spider chart illustrating the predictive performance of all models, enabling clearer visualization and comparison.
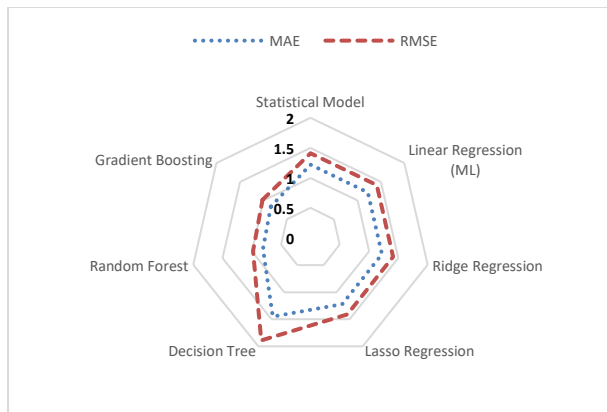


**Fig 1: Comparing the predictive performance statistical model with machine learning models**

## 4.3  Model Insights

To evaluate the predictive performance of each machine learning (ML) model, the dataset is split into a training set (80%) and a test set (20%). When used for prediction on the test set, the statistical model performs identically to ML's linear regression model. The outcome is expected, given that both are linear models, and the coefficients learned from the full dataset (statistical) versus 80% of the data (ML) are nearly identical. The random forest model emerged as the most effective for prediction, outperforming all other models by achieving the lowest error rates (RMSE = 0.98; MAE = 0.81) on this dataset. This means its forecasts are, on average, off by only 810 sales units, which is a significant improvement over the other linear models considered here. The superior performance of ensemble

models like random forest and gradient boosting shows that they can capture complex, non-linear relationships in the data that the simpler linear models cannot.

## 5.  RESULTS AND DISCUSSION

The direct comparison highlights the fundamental trade-off between the two modeling approaches. The core difference lies in their primary goals, which dictate their strengths and weaknesses, particularly concerning interpretability. Table 3 contrasts the two modeling approaches. It effectively illustrates the trade-off between interpretability and predictive accuracy, helping analysts and stakeholders choose the right tool based on business objectives – whether to understand the "why" or to optimize the "what will happen."

**Table 3. Statistical vs. machine learning models**

| Feature | Statistical Model (Linear Regression) | Best Predictive Model (Random Forest) |
|---|---|---|
| Primary Goal | Explanation | Prediction |
| Key Output | Coefficients and p-values | Prediction for a given input |
| Performance Metric | R-squared | RMSE |
| Interpretability | High | Low |
| Model Assumptions | To be validated | Distribution-free |
| Actionable Insight | Strategic | Operational |

Linear regression (a statistical model) aims to explain relationships between variables. For instance, it identifies how changes in advertising expenditure across different channels (like TV, radio, and newspaper) affect sales. The statistical model provides interpretable coefficients. For example, a $1k increase in radio advertising expenditure increases sales by 189 units. It evaluates the model fit using R-squared, which expresses how much of the variance in sales is explained by the model. It is highly interpretable and clearly provides estimates of each variable's impact; therefore, it is easily communicated to stakeholders. Statistical models like linear regression have strict assumptions (e.g., linearity, homoscedasticity, normal residuals) that must be validated for results to be trustworthy. It offers strategic insights; for example, it reveals which advertising channels drive sales and guides budget allocations.

On the other hand, random forest (the best machine learning model) is designed for high-accuracy prediction of future sales without explaining the underlying relationships. The predictive model outputs a numerical forecast. For example, given the advertising budget for each channel, one can expect 15,250 units to be sold. It uses MAE on new data to assess real-world prediction accuracy. It is less interpretable and aggregates predictions from many decision trees, making it a "black box" to understand and explain, though often more accurate. Predictive models like random forests require fewer assumptions, focusing on empirical accuracy. It supports operational decisions such as accurate sales forecasts, aiding in inventory planning, revenue forecasting, and setting performance targets.

## 6.  CONCLUSION

The study underscores the fundamental distinction between modeling for explanation and modeling for prediction. Neither approach is inherently superior; they are powerful tools

designed for different business objectives [1].

The statistical regression model excelled at providing clear, interpretable insights. Its strength lies in its transparency; the model's equation is easily understood, and each coefficient has a direct business interpretation. This makes it an invaluable tool for explaining "why" certain factors drive the business, which is essential for gaining stakeholder buy-in for strategic decisions like marketing budget allocation.

The machine learning approach, by testing a variety of algorithms, identified the random forest as the most accurate predictive engine. Its superior performance comes at the cost of interpretability. As a "black box" model, dissecting the exact reasons behind a specific forecast is difficult. However, its higher precision is invaluable for operational tasks where the prediction accuracy is paramount, such as managing inventory or setting financial targets.

The direct comparison on the test set makes the trade-off explicit. The statistical model provides invaluable strategic guidance due to its explainability; however, it is not the best tool for pure forecasting. The random forest model, despite its complexity, is the superior choice when predictive accuracy is the primary goal.

Ultimately, the most effective business analytics strategy leverages the synergy between both approaches or cultures [1], [4], [9]. The analyst can use statistical models to understand "why" certain factors drive the business and then use the best machine learning models to predict "what" will happen in the future, creating a comprehensive, data-driven foundation for strategy and operations.

# 7. REFERENCES

[1] Shmueli, G. 2010. To explain or to predict? Statistical Science, 25(3), 289-310.

https://doi.org/10.1214/10-STS330

[2] Efron, B. 2020. Prediction, estimation, and attribution. Journal of the American Statistical Association, 115(530), 636-655. https://doi.org/10.1080/01621459.2020.1762613

[3] Friedman, J., Hastie, T., and Tibshirani, R. 2020. Discussion of "Prediction, Estimation, and Attribution" by Bradley Efron. Journal of the American Statistical Association, 115(530), 665-666. https://doi.org/10.1080/01621459.2020.1762617

[4] Breiman, L. 2001. Statistical modeling: The two cultures (with comments and a rejoinder by the author). Statistical Science, 16(3), 199-231. https://doi.org/10.1214/ss/1009213726

[5] Shmueli, G. 2021. Comment on Breiman's "Two Cultures" (2002): From two cultures to multicultural. Observational Studies, 7(1), 197-201. https://doi.org/10.1353/obs.2021.0010

[6] James, G., Witten, D., Hastie, T., and Tibshirani, R. 2023. An Introduction to Statistical Learning with Applications in R, 2nd edition, Springer.

[7] Montgomery, D.C., Peck, E.A., and Vining, G.G. 2021. Introduction to Linear Regression Analysis, 6th edition, Wiley.

[8] Jordan, M.I., and Mitchell, T.M. 2015. Machine learning: Trends, perspectives, and prospects. Science, 349(6245), 255-260. https://www.science.org/doi/10.1126/science.aaa8415

[9] Daoud, A., and Dubhashi, D. 2023. Statistical modeling: The three cultures. Harvard Data Science Review, 5(1). https://doi.org/10.1162/99608f92.89f6fe66