# Dosage Disparities: Medicaid Reimbursement Trends and Predictive Insights

Raheela Charania, Prabhu Shankar and Rajivini Tiruveedhula

May 21, 2025

---

## Abstract

**:**

This study investigates patterns in Medicaid drug reimbursement across U.S. states from 2021-2024, with a focus on dosage disparities and underlying socioeconomic influences. Using data from the Medicaid Drug Utilization dataset, we conducted exploratory data analysis and developed predictive models to better understand state-level variations. Visualizations revealed that sodium chloride is the most frequently reimbursed drug in multiple states, while scatterplots demonstrated a proportional relationship between the number of prescriptions and units reimbursed. By integrating external variables such as poverty rates, per capita income, and real GDP, our regression and random forest models provided deeper insights into reimbursement dynamics. The findings indicate that the number of prescriptions is the strongest predictor of units reimbursed, while geographic and socioeconomic variables offered limited predictive power. These results underscore the complexity of Medicaid reimbursement mechanisms and suggest potential pathways for improving cost transparency and resource allocation.

**Index Terms:** Medicaid Drug Utilization, Medicaid Reimbursement, Dosage Disparities, Drug Pricing Trends, Health Economics, Regression Analysis, Geographic Variation, NDC (National Drug Code).

---

## 1 Introduction

Identify and investigate the differences in Medicaid drug reimbursements across various U.S. states, focusing on identifying consistent trends and assessing the potential impact of socioeconomic factors. Since Medicaid operates under a shared federal-state funding model but is implemented differently by each state, substantial variation in reimbursement approaches is expected. To analyze these differences, the study incorporates state-level data alongside external variables such as gross domestic product, poverty rates, income per capita, population size, and climatic conditions. By combining Medicaid drug utilization data from 2019 to 2024 with these socioeconomic indicators, the study evaluates how these broader contextual elements may affect reimbursement behaviors.

### 1.1 Background

Rising prescription drug costs have imposed a growing burden on the U.S. healthcare system, with Medicaid—an essential program for low-income individuals—bearing a significant share of this pressure. Because Medicaid is administered at the state level within federal parameters, each state determines its own approach to drug reimbursement. There is an emerging influence of broader socioeconomic and environmental conditions on healthcare delivery. Variables such as the gross domestic product (GDP), poverty levels, and even climate indicators like temperature for every state differ and may play a role in shaping drug consumption and reimbursement behaviors. Despite this, there is limited empirical evidence that clearly quantifies these relationships across states. This study seeks to fill that void by addressing the core question: What factors contribute to regional differences in Medicaid drug reimbursement? A complementary objective examines whether disparities in reimbursed drug dosages can be anticipated based on socioeconomic indicators such as income, poverty, and climate conditions. Using state-level data from 2019 to 2024, this project combines exploratory analysis with predictive modeling to reveal hidden patterns and evaluate the explanatory power of these contextual variables. We also aim to assess whether these external factors significantly influence reimbursements.

## 2 Literature Review

Recent analyses show that generic medications make up the majority of Medicaid prescriptions, but brand-name and specialty medications drive the majority of program costs. Specialty drugs, although accounting for only 2% of prescriptions, represent nearly 28% of Medicaid's drug spending (Bruen et al.,2014). The emergence of expensive specialty drugs like Direct Acting Antivirals (DAAs) for Hepatitis C reshaped Medicaid drug expenditures, with the average price per DAA prescription exceeding $20,000 (Gari et al.,2023) . However, access to these high-cost medications was often limited by prior authorization policies, contributing to significant inequities in dosage access across different populations. Insurance coverage plays a pivotal role in improving access to medications and ensuring adherence to prescribed dosages.

Studies have found that individuals with prescription drug coverage fill significantly more prescriptions and face lower out-of-pocket expenses compared to those without coverage (U.S. Department of Health and Human Services). The expansion of Medicaid under the Affordable Care Act (ACA) led to a substantial increase in medication use, particularly among low-income adults. States that expanded Medicaid saw a 19% rise in Medicaid-paid prescriptions, equivalent to approximately nine additional fills per new enrollee annually (Ghosh et al., 2019). Despite this progress, dosage disparities persisted due to residual cost-sharing mechanisms and state-level variation in drug benefit administration. States with higher Medicaid copayments experienced smaller gains in prescription drug utilization, suggesting that even small financial barriers can significantly impact medication adherence (Ghosh et al., 2019).

While substantial evidence documents the general impact of insurance coverage on drug utilization, few studies focus specifically on dosage disparities — differences in patients' ability to access the full, recommended amount of a prescribed medication. Most prior work centres on whether a prescription is filled at all, rather than whether it is filled consistently and completely. Additionally, much of the research aggregates all Medicaid beneficiaries, without examining how dosage access may differ by socioeconomic factors such as income level, or state of residence

## 3 Data Preperation

### 3.1 Data Dictionary

| Name | Type | Description |
|---|---|---|
| utilization_type | string | Fee-For-Service Utilization & Managed Care Organization Utilization |
| state | string | State Name |
| ndc | string | LabelerCode + ProductCode + PackageSize |
| labeler_code | string | Name of Manufacturer/Supplier |
| product_code | string | Manufacturer, labeler, packager or distributor of the drug |
| package_size | string | Size of the package of the drug |
| year | integer | Year of dataset |
| quarter | integer | Quarter of dataset |
| product_name | string | Name of product |
| units_reimbursed | number | Units of drug reimbursed |
| number_of_prescriptions | integer | Number of prescriptions reimbursed |
| total_amount_reimbursed | number | Total amount reimbursed |
| medicaid_amount_reimbursed | number | Medicaid amount reimbursed |
| non_medicaid_amount_reimbursed | number | Non-Medicaid amount reimbursed |
| State_GDP | numeric | State-level real GDP |
| Poverty_Index | numeric | Percentage of population below poverty line |
| Temperature_Index | numeric | Average annual temperature per state |
| Per_Capita_Income | numeric | Average income per person in the state |
| Population_Estimates | numeric | Total population of the state |

## 3.2 Data Overview

This analysis uses publicly available data from the Medicaid Drug Utilization database on Data.Medicaid.gov, which provides detailed, quarterly records of prescription drug use and reimbursement across all U.S. states. The dataset offers a robust foundation for longitudinal analysis and supports the examination of trends over multiple years. The primary goal is to identify state-level differences in Medicaid reimbursement and explore how additional socioeconomic factors may influence these patterns. To strengthen the dataset's analytical capabilities, five additional socioeconomic indicators were incorporated from external sources. These variables—including state-level metrics like per capita personal income, % people below poverty line, Real GDP, and climate conditions—were chosen to add contextual depth and allow for a more comprehensive investigation into the factors influencing Medicaid reimbursement patterns. With these enhancements, the dataset becomes well-suited for longitudinal analysis aimed at uncovering regional disparities in Medicaid reimbursements and evaluating the role of socioeconomic conditions in shaping these outcomes

## 3.3 Data Cleaning and preprocessing

The dataset was fairly clean and did not contain any missing values. Although a few outliers were present, they were removed—rather than imputed— in order to avoid inconsistencies. Because the dataset also indicates whether a drug claim was processed through a Managed Care Organization (MCOs) or Fee-For-Service Utilization (FFSU), we first filtered the data to include only MCOs claims. These claims are known as Managed Care Organizaiton Utilization (MCOU). This decision reflects the fact that over 75% of Medicaid enrollees are covered through MCOs nationally (Centers for Medicare & Medicaid Services, 2023), and focusing on MCOU claims provides a more accurate representation of current Medicaid operations. States have increasingly expanded the use of managed care to control costs and improve care coordination, with MCOs implementing measures like formulary restrictions and prior authorization requirements that directly influence drug utilization patterns. After filtering, the dataset was grouped by state, NDC code, product name, and quarter, and unnecessary columns such as package size,

product code, labeler code, and suppression code were removed as they are all included in the NDC code to simplify the dataset and avoid redundancies. The final data used for modeling included observations from 2019 to 2024, with 2024 reserved as the test dataset and 2019–2023 used for training. Exploratory data analysis was performed, and Min-Max standardization was applied to the relevant variables to prepare the data for modeling.

**Feature Engineering:** To encode categorical variables, target mean encoding was applied within a 5-fold cross-validation framework. For each fold, the mean of the target variable was computed using only the training subset, aggregated by the categories of the respective feature. These fold-specific averages were then used to encode the categorical values in the corresponding validation subset.

Integrating cross-validation into the encoding process is essential for mitigating data leakage. If the encoding were based on the full dataset, it would introduce information about the target variable into the predictors, thereby compromising model validity and leading to overfitting. By restricting the calculation to the training portion of each fold, the approach ensures that encodings remain independent of the target values in the validation data.

In instances where a category in the validation set does not appear in the training fold, the overall mean of the target variable across the training data is used as a substitute. This default assignment prevents the introduction of missing values and supports the robustness of the model. Collectively, this method preserves the underlying association between categorical inputs and the target variable while enhancing the model's capacity to generalize to unseen data.

## 3.4 Method selection

The initial phase of the project involved conducting exploratory data analysis (EDA) to better understand the structure and relationships within the dataset. Key steps included:

- Generating a correlation matrix to assess the relationships between variables.
- Creating a line graph to evaluate relation between units reimbursed and population over the years.

- Identifying total units reimbursed by different temperature groups.
- Identifying total units reimbursed by region.
- Developing U.S. map visualizations to identify regional variation in drug utilization and spending, including:
  - The top reimbursed drug by average total units per state.
  - The total amount reimbursed per 100k individuals
  - Per Capita Personal Income
  - The top manufacturer per state over the years (2019-2021)

These visual tools helped identify meaningful trends and outliers, guiding the modeling process.

Preliminary predictive models were built using only the original Medicaid dataset. To improve the model's explanatory power, several external variables were added to the dataset, including:

- Average temperature
- Poverty rate by state
- Population estimate by state
- Per capita personal income
- Real GDP

These additional predictors were chosen to provide a broader context and increase the robustness of the models. Intermediate runs were performed with the expanded dataset to assess improvements in model performance.

## 4 Results and Analysis

### 4.1 Exploratory Data Analysis

The correlation matrix (Figure 1) reveals several noteworthy relationships among variables. Notably, Medicaid Amount Reimbursed and Total Amount Reimbursed exhibit a perfect positive correlation, indicating that Medicaid reimbursements constitute a major component of total reimbursements. A similarly strong positive correlation exists between GDP per state and personal income (0.98), suggesting that economically stronger states tend to have higher average incomes. Conversely, personal income and poverty index display a strong negative correlation (-0.63), consistent with expectations. Moderate correlations are observed between Units Reimbursed and Number of Prescriptions (0.39), as well as between poverty index and Temp In-

dex (0.62), indicating potential geographic or demographic influences. The target variable Units Reimbursed prob shows only weak correlations with most features, such as Number of Prescriptions(0.44) and Units Reimbursed (0.27), and is largely uncorrelated with socioeconomic variables. Lastly, variables like NDC and population estimates exhibit minimal correlation with others, suggesting limited direct predictive value.
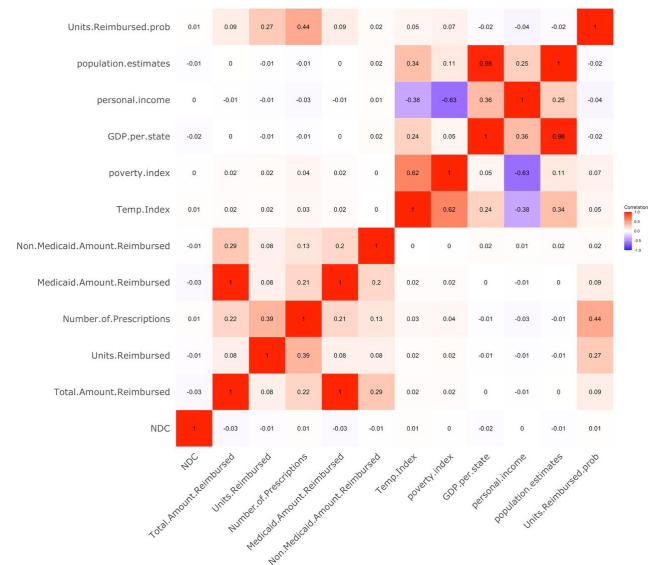


**Figure 1.** Correlation Matrix showing the correlation between selected variables.
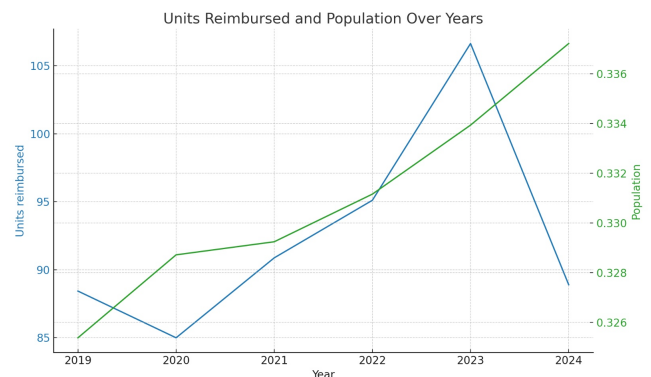


**Figure 2.** Line graph of Units Reimbursed vs. Population over the years.

The line chart (Figure 2) presents the progression of units reimbursed and population size from 2019 to
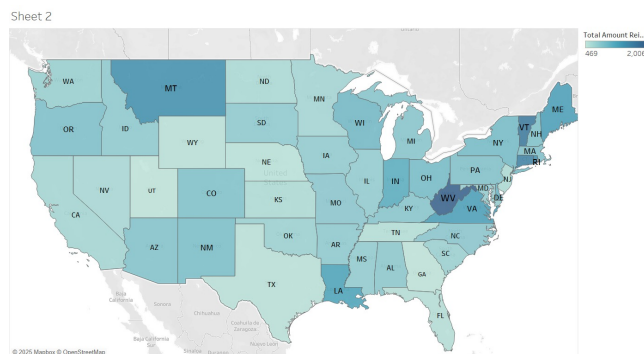
**Figure 3.** Total amount reimbursed per 100k Individuals (2019).



**Figure 4.** Personal Income per Capita by state (2019).

2024. While population figures show a consistent upward trend throughout the period, the number of units reimbursed exhibits greater variability—marked by a sharp increase between 2021 and 2023, followed by a significant drop in 2024. This downturn is largely the result of the conclusion of the COVID-19 continuous coverage requirement, which had previously restricted states from removing individuals from Medicaid rolls during the public health emergency. With the policy's expiration in March 2023, states reinstated eligibility reviews, leading to widespread exits in 2024.

This map (Figure 3) illustrates state-level variation in total Medicaid reimbursements, where darker shades represent higher amounts. States such as West Virginia, Vermont, and Montana show the highest totals, while lighter areas like Wyoming and Utah indicate lower reimbursements. These differences point to regional disparities.

This map (Figure 4) displays personal income levels across U.S. states, with darker shades representing higher income. Northeastern states like Massachusetts, New Jersey, and Connecticut show the highest income levels, while states in the South and parts of the Midwest have lower values. The variation highlights regional economic differences, likely influenced by factors such as cost of living, industry presence, and employment opportunities.
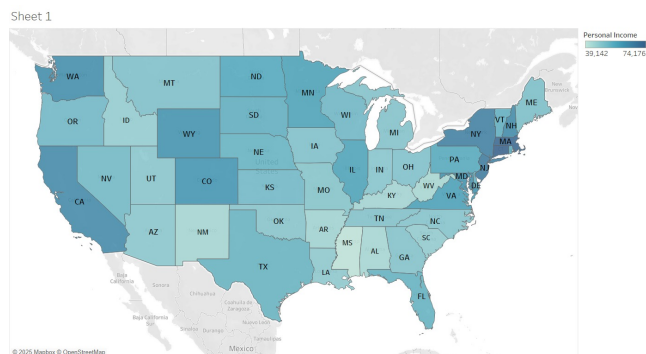
This map (Figure 5) highlights the top reimbursed drug by state, with different colors representing distinct medications. Sodium Chloride is the most commonly reimbursed drug across a majority of states, especially in the South and Midwest. GaviLyte-G and Polyethylene Glycol are dominant in several Western and Northeastern states, while Amoxicillin appears frequently in the Southwest. A few states, such as Missouri and Colorado, show other leading drugs like Frisium and category-labeled "Others." The variation reflects differences in regional prescribing trends, population health needs, and treatment practices.
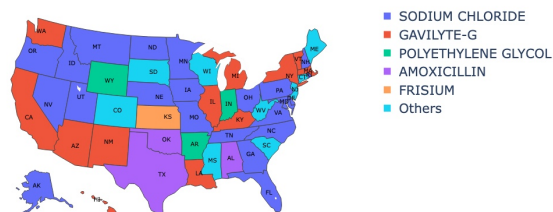


**Figure 5.** Top Drugs by State. Sorted by Units Reimbursed

Next, we conducted a visual analysis of the leading pharmaceutical manufacturers by state in the U.S. from 2021 to 2024, using price per unit as the key metric. In 2021, Genentech, Inc. (red) held a dominant presence across much of the country (Figure 6). By 2022, a notable shift occurred, with Genentech, Inc. (green) emerging as the top manufacturer in the majority of states, particularly in the Midwest, South, and Western regions (Figure 6). This trend remained consistent through 2023 and 2024, with Genentech, Inc. (blue) continuing to lead across most states (Figure 6 & 7). Janssen Biotech, Inc. and Regeneron Pharmaceuticals,

**5**

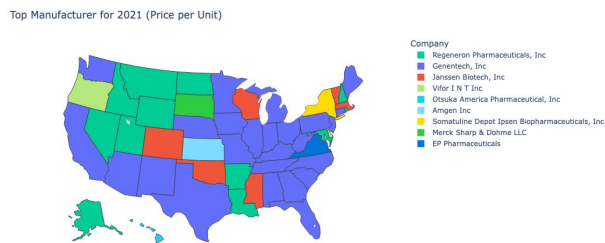Inc. also maintained a significant presence, ranking as the next most prominent manufacturers in several regions.
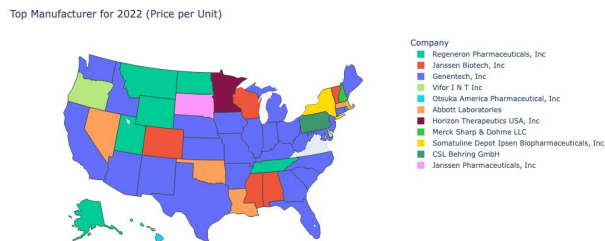


**Figure 6.** Top Manufacturer by Price per Unit (2021)



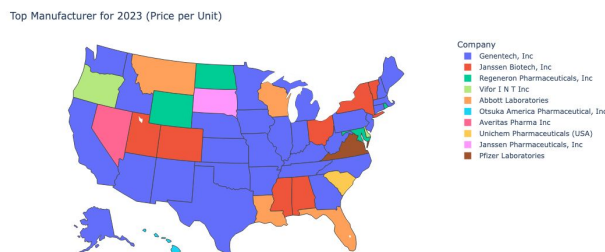**Figure 7.** Top Manufacturer by Price per Unit (2022)



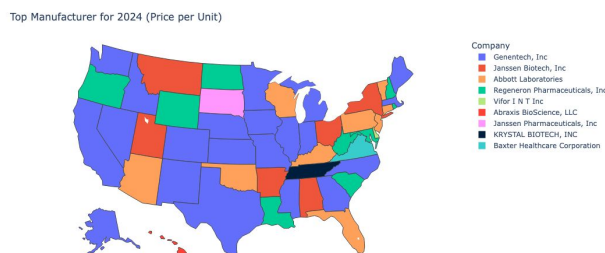**Figure 8.** Top Manufacturer by Price per Unit (2023)



**Figure 9.** Top Manufacturer by Price per Unit (2024)

The bar chart (Figure 10) shows total units reimbursed by U.S. region, with the South leading significantly, followed by the Northeast and Midwest, while the West records the lowest total. This suggests higher Medicaid drug utilization in the southern states, which may relate to population size, health conditions, or coverage policies. The accompanying pie chart (Figure 11) breaks down total units reimbursed by temperature groups, with nearly half coming from cold regions (49.14%), followed by warm (32.52%) and mild areas (16.89%). Very cold and hot regions contribute minimally. Together, these visuals highlight that both geographic and climatic factors may play a role in influencing
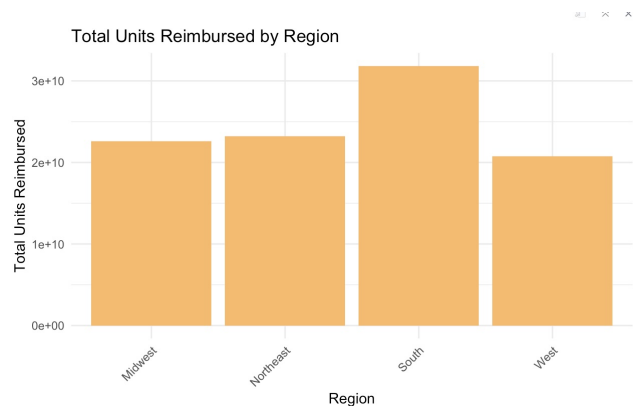


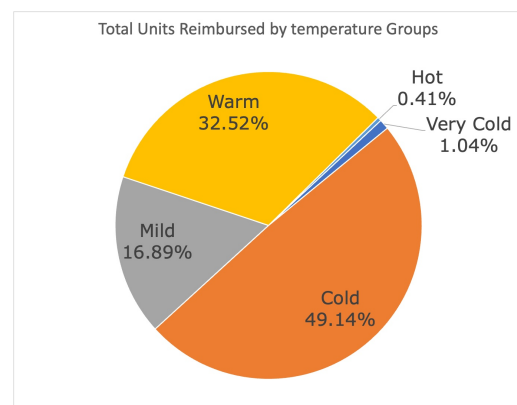**Figure 10.** Total Units Reimbursed by Region.



**Figure 11.** Total Units Reimbursed by Temperature groups.

### 4.2 Models

To identify the primary factors influencing Medicaid drug reimbursements and dosage variations, three predictive models were constructed and assessed: Multiple Linear Regression (MLR), Light Gradient Boosting

Machine (LightGBM), and Extreme Gradient Boosting (XGBoost). These modeling techniques were selected to strike a balance between transparency and predictive strength, while also testing the added value of incorporating socioeconomic and environmental variables. All models were built using standardized inputs and applied target mean encoding within a cross-validation framework to avoid data leakage. The objective was to forecast the volume of units reimbursed based on prescription data and state-level indicators, including poverty rates, population estimates, per capita income, and climate metrics.

### 4.2.1 Model 1: Multiple Linear Regression

The Multiple Linear Regression (MLR) model identified several significant factors influencing the number of Medicaid drug units reimbursed. The most influential positive predictor was the number of prescriptions, underscoring a strong correlation between prescription volume and reimbursed units. Medicaid Amount Reimbursed also had a large positive coefficient, indicating that increased Medicaid spending is closely tied to greater reimbursement volumes. Additionally, higher values in the Poverty Index and Temperature Index were modestly associated with increased reimbursements, suggesting that socioeconomic and environmental factors may play a role. On the other hand, variables such as personal income, GDP per state, and Non-Medicaid Amount Reimbursed exhibited negative relationships with the target variable. These inverse associations may suggest that states with higher income levels and economic output tend to have lower Medicaid utilization, potentially due to reduced program dependency or differing prescribing behaviors in more affluent areas.

| Variable | Estimate |
|---|---|
| (Intercept) | 3217.3 |
| Personal Income (Scaled) | -1734.4 |
| GDP per State (Scaled) | -12113.5 |
| Poverty Index (Scaled) | 7617.1 |
| Temperature Index (Scaled) | 6653 |
| Non-Medicaid Amt Reimbursed (Scaled) | -64398.6 |
| Medicaid Amt Reimbursed (Scaled) | 255289.5 |
| Number of Prescriptions (Scaled) | 2963643.9 |

### 4.2.2 Model 2: Light Gradient Boosting Machine

The Light Gradient Boosting Machine (LightGBM) model provided insights into the relative importance of each feature in predicting the number of units reimbursed. The most influential variable was Non-Medicaid Amount Reimbursed, suggesting a notable distinction in utilization patterns between Medicaid and other payment sources. This was closely followed by Medicaid Amount Reimbursed and the Number of Prescriptions, both of which were strong contributors to the model's predictive performance. Additional variables—including GDP per state, Temperature Index, Poverty Index, and Personal Income—also demonstrated meaningful influence, albeit to a lesser extent. The inclusion of both economic and demographic indicators among the top predictors highlights the multifactorial nature of drug reimbursement patterns and reinforces the utility of gradient boosting techniques in capturing complex relationships across diverse features.

| Variable | Importance |
|---|---|
| Non-Medicaid Amount Reimbursed (Scaled) | 2638 |
| Medicaid Amount Reimbursed (Scaled) | 2346 |
| Number of Prescriptions (Scaled) | 2315 |
| GDP per State (Scaled) | 2036 |
| Temperature Index (Scaled) | 1944 |
| Poverty Index (Scaled) | 1867 |
| Personal Income (Scaled) | 1854 |

### 4.2.3 Model 3: Extreme Gradient Boosting

The Extreme Gradient Boosting (XGBoost) model highlighted the Number of Prescriptions as the most dominant predictor of units reimbursed, with a significantly higher importance score than all other variables. This finding underscores the strong link between prescription volume and reimbursement levels. Personal Income and Medicaid Amount Reimbursed followed as the next most influential features, indicating that both individual financial capacity and Medicaid spending play meaningful roles in shaping reimbursement outcomes. Other variables—such as Non-Medicaid Reimbursement, GDP per state, Poverty Index, and Temperature Index—had comparatively lower importance scores but still contributed to the model's performance. Overall, the XGBoost model effectively captured the complex interactions among socioeconomic, economic, and healthcare-specific variables in predicting Medicaid reimbursement trends.

| Variable | Importance |
|---|---|
| Number of Prescriptions (Scaled) | 8232 |
| Personal Income (Scaled) | 2695 |
| Medicaid Amount Reimbursed (Scaled) | 2277 |
| Non-Medicaid Amount Reimbursed (Scaled) | 1617 |
| GDP per State (Scaled) | 1393 |
| Poverty Index (Scaled) | 1244 |
| Temperature Index (Scaled) | 966 |

### 4.3 Model performance

Among the four models evaluated, the XGBM model demonstrated the strongest performance, achieving the highest R² value of 0.93 and the lowest RMSE of 41,991.16, indicating superior predictive accuracy and minimal error. Following closely, the LGBM model also performed well, with an R² of 0.89 and an RMSE of 61,231.16, reflecting a substantial improvement in model fit and accuracy over traditional techniques. In contrast, the Linear Model and Stepwise regression produced identical results, each with an R² of 0.67 and an RMSE of 74,401.75. These results suggest that while the simpler models offer moderate predictive power, the gradient boosting approaches—particularly XGBM—are far more effective for forecasting in this context.

| Model | $R^2$ | RMSE |
|---|---|---|
| Linear Model | 0.67 | 74,401.75 |
| Stepwise | 0.67 | 74,401.75 |
| LGBM | 0.89 | 61,231.16 |
| XGBM | 0.93 | 41,991.16 |

## 5 Cost-Benefit Analysis

The cost-benefit assessment highlights the substantial financial upside of improving predictive accuracy in Medicaid reimbursement forecasting. By reducing the RMSE from 8% to 1%, potential misallocation costs dropped from approximately $274.4 million to $34.3 million—a variance reduction of $240.1 million. Even realizing just 10% of this as efficiency savings could yield $24 million in tangible benefits. These findings align with broader industry insights, such as those from McKinsey & Company and Health Affairs, which underscore how advanced analytics can deliver 8–20% cost reductions and enable more targeted, time-sensitive healthcare interventions. This reinforces the strategic value of predictive modeling in enhancing fiscal efficiency and data-driven decision-making in public healthcare systems. According to a McKinsey analysis, healthcare systems that applied advanced analytics for cost control achieved savings between 8% and 15%, all while maintaining or enhancing care quality. Based on this industry insight, we selected a moderate 10% estimate to represent potential efficiency gains in our cost-benefit analysis. This figure served as the basis for estimating realized savings, allowing our projections to remain grounded in established outcomes while reflecting a cautious and realistic approach to quantifying the fiscal benefits of enhanced predictive accuracy in Medicaid reimbursement forecasting.

| Item | Value |
|---|---|
| Baseline error cost (8% RMSE) | $274,424,620 |
| Model error cost (1% RMSE) | $34,303,080 |
| Variance Reduced | $240,121,540 |
| Estimated realized savings (10%) | $24,012,154 |

## 6 Conclusion and Recommendation

This analysis explored Medicaid drug reimbursement trends across U.S. states between 2019 and 2024, with a focus on dosage inconsistencies and the influence of socioeconomic conditions. The results indicate that the primary drivers of Medicaid reimbursements are the number of prescriptions and the total units reimbursed, whereas factors such as GDP, poverty rates, population size, and climate variables improved the model, though only some were significant predictors of Medicaid reimbursement variations. Furthermore, the most reimbursed drugs showed consistency across regions, suggesting a uniform pattern in Medicaid utilization. These findings highlight the importance of concentrating on prescription volume and drug-specific metrics when building predictive models, rather than relying heavily on broader economic or demographic indicators. Incorporating these core utilization variables into forecasting efforts can improve the precision of Medicaid planning and resource distribution, while contextual socioeconomic factors should be used selectively given their limited impact.

# 7 References

1. Centers for Medicare & Medicaid Services. (2023). Medicaid managed care. Medicaid.gov. https://www.medicaid.gov/medicaid/managed-care/index.html

2. National Centers for Environmental Information. (n.d.). Climate at a glance: National time series. National Oceanic and Atmospheric Administration. https://www.ncei.noaa.gov/access/monitoring/climate-at-a-glance/national/mapping

3. U.S. Census Bureau. (n.d.). Population and housing unit estimates. https://www.census.gov

4. U.S. Bureau of Economic Analysis. (n.d.). Regional data: GDP and personal income. https://www.bea.gov

5. Department of health and human services. (1999). Effects of prescription drug coverage on spending and utilization. https://aspe.hhs.gov/sites/default/files/private/pdf/172171/c2.pdf

6. Gari, M. H., Alsuhibani, A., Alashgar, A., & Guo, J. J. (2023). Utilization, reimbursement, and price trends for hepatitis C virus medications in the US Medicaid programs: 2001–2021. Exploratory Research in Clinical and Social Pharmacy, 12, 100383. https://doi.org/10.1016/j.rcsop.2023.100383

7. Ghosh, A., Simon, K., & Sommers, B. D. (2019). The effect of health insurance on prescription drug use among low-income Adults:Evidence from recent Medicaid expansions. Journal of Health Economics, 63, 64-80. https://doi.org/10.1016/j.jhealeco.2018.11.002

8. What drives spending and utilization on Medicaid drug benefits in states? (2014, December 10). KFF. https:// www.kff.org/medicaid/issue-brief/what-drives-spending-and-utilization-on-medicaid-drug-benefits-in-states/

9. Lee, S. (2025, March 27). 7 data-driven insights for healthcare cost optimization success. Number Analytics. https://www.numberanalytics.com/blog/7-data-driven-insights-healthcare-cost-optimization-success