

Unlocking Cost Savings in Healthcare: How Difference-in-Differences (DID) Can Measure the Impact of Interventions

Vidya Rajasekhara Reddy Tetala

Abstract: *Difference-in-Differences is a powerful means of obtaining an estimate for the causal effect of interventions from observational data, especially in health care, economics and social sciences where random controlled trials (RCT) are often impracticable. This paper explores how DID can be applied to healthcare cost-saving interventions by balancing treatment and control groups with propensity score matching and weighting with the inverse propensity score. It also discusses the recent development of this method, integrating machine learning into DID, which strengthens the capability of DID. Case study on cost reduction of hospital readmission illustrates the usefulness of the methodology. It is further elaborated with a detailed explanation of the calculation process and the application of the propensity score to remove the confounding biases.*

Keywords: Difference-in-Differences, DID, healthcare treatment or interventions, cost reduction, calculation of propensity score, inverse of the propensity score, machine learning, causal inference, cloud computing, SRE.

1. Introduction

1.1 DID Methodology Overview

The Difference-in-Differences is a quasi-experimental design aimed at comparing changes in outcomes that are observed over time across a treatment-altered, or study, group and a control group. Because DID controls for the presence of confounding factors, the major strength of DID is its ability to assume that both groups would have followed a parallel trend in the absence of treatment. In general, DID has been widely used in healthcare to assess policy interventions at reducing costs and improving patient outcomes.

1.2 DID in Healthcare Application

The most significant applications of DID in health involve the analysis of the effects of interventions in situations where randomization is impossible or harmful. The DID method permits intervention examination through the comparison of cost and outcome trends in both treated and untreated populations before and after a policy change or treatment. Recent improvements in matching methods, specifically the propensity score match and the inverse propensity score weighted methods, have tended to reduce many of the selection biases found in applications using DID and therefore making causal inferences far more sure.

2. Literature Review

Diff-in-diff has been one of the most widely used methods in health studies to evaluate the impacts of various interventions such as changing hospital payment policies, insurance reforms, and the implementation of clinical guidelines (Bertrand et al., 2004). These two studies really nail the uses of DID in terms of cost evaluation and resource optimization. The researchers affirmed the issue of selection bias in DID,

as in non-randomized settings, there is a great use for the advanced matching technique propensity score matching and inverse propensity score estimations. Other recent works have explored the integration of machine learning into DID to further advance the performance of the methodology when dealing with large data volumes and complex, possible nonlinear relationships between variables. Works by among others Athey & Imbens, 2016 fall into this category. Successful completion of such tasks can have significant improvements in healthcare applications where precision is a key factor in making causal estimates.

3. DID Methodology

3.1 Study Group and Control Group

In DID, the treatment group-or study group-consists of those individuals or institutions exposed to some sort of intervention, while the control group consists of those that are not. A fundamental assumption is that, in the absence of the intervention, the trends of the study and control groups would have moved-that is, changed-in a parallel manner over time, or in other words, the parallel trends assumption.

3.2 Pre- and Post-intervention Measurement

While DID does require the outcome data to be measured at two points in time-before and after the intervention-for both study and control groups, these differences in outcomes between the groups after accounting for pre-existing trends enable the researcher to tease out the effect of the intervention.

3.3 DID Estimator

The DID estimator measures the difference between the changes in outcomes in the study group and the control group. It is expressed as:

$$DID = (Y_{post,treatment} - Y_{pre,treatment}) - (Y_{post,control} - Y_{pre,control})$$

where $Y_{post, treatment}$ and $Y_{pre, treatment}$ are the outcomes of the study group after and before the intervention, respectively, and $Y_{post, control}$ and $Y_{pre, control}$ represent the corresponding outcomes for the control group.

3.4 Propensity Score Matching (PSM) and Inverse Propensity Score

3.4.1 What is Propensity Score?

The propensity score is the probability of an individual or unit to which treatment is assigned given a set of observed covariates. Key idea related to propensity scores is to decrease the bias due to confounding variables by making the treated and control units comparable with respect to their observable characteristics.

3.4.2 How is Propensity Score Calculated?

Propensity scores are most commonly calculated by the application of logistic regression or other binary classification techniques. It is an effort to model the probability of getting treated as a function of observed covariates. In a healthcare-based setting, these may be patient demographics such as age or gender, or health status, including chronic conditions and previous medical history, or may also refer to hospital characteristics: size or location. The propensity score for an individual i , in this respect, is defined in:

$$P(X_i) = P(T_i = 1|X_i)$$

where $T_i=1$ indicates that individual i belongs to the treatment group, and X_i represents the vector of covariates. The logistic regression model for calculating the propensity score can be expressed as:

$$P(T_i = 1|X_i) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})}}$$

Here, $\beta_0, \beta_1, \dots, \beta_n$ are the estimated coefficients for the covariates. Once the propensity scores are calculated, individuals from the study and control groups can be matched based on their propensity scores, creating balanced groups that are comparable in terms of observed characteristics.

3.4.3 Propensity Score Matching (PSM)

After estimating the propensity scores, the Propensity Score Matching method will be used to match those in the study and control groups with similar propensity scores. This is based on the fact that matching individuals from both groups will reduce their bias due to observed covariates. Matching can be done using different methods such as nearest neighbor matching, caliper matching, or kernel matching.

3.4.4 Inverse Propensity Score Weighting

Another balancing technique for study and control groups is inverse propensity score weighting. Rather than weighing subjects directly against a match, subjects are weighted inversely by their propensity scores. The idea here is to give more weight to subjects in the control group that resemble the study group subjects, which accommodates any residual imbalance. The weight for the individual i is given by:

$$w_i = \frac{1}{P(X_i)} \quad \text{if } T_i = 1$$

$$w_i = \frac{1}{1 - P(X_i)} \quad \text{if } T_i = 0$$

By weighting individuals this way, the treatment and control groups can be made comparable, even in the presence of selection bias.

3.5 Balancing Study and Control Populations

Diagnostic tools—such as standardized mean differences or Love plots—can be used to check whether, after matching or weighting, the study and control populations are balanced. That is, it checks if the covariates are well balanced across the matching/weighting groups so that the subsequent DID analysis will yield valid causal estimates.

4. DID Methodology Improvements

4.1 Synthetic Control Methods

The synthetic control method represents an extension of DID applied when there is only a small number of treated units. It builds the weighted combination of untreated units to form the synthetic counterpart of the treated unit. It has been very productive in policy analysis when one or a few units make up the treatment group.

4.2 Dynamic DID Models

Traditional models of DID consider that the effect of treatment remains constant over time, but dynamic DID models relax this assumption, with treatment effects varying across time. Such models would be useful in health, where the interventions may have short-term and long-term impacts. Dynamic DID captures this time-varying heterogeneity in treatment effects, thus enabling a nuanced understanding of how the dynamics of an intervention play out over time.

4.3 Integration of Machine Learning: Enhancing DID with ML Techniques

One of the recent exciting developments in enhancing the robustness of estimated causal effects is the integration of techniques from machine learning into DID. In healthcare settings, where data are often complex and their inter-relationships nonlinear, ML models can give more accurate propensity scores and allow better balancing between the study and control groups.

4.3.1 Machine Learning for Propensity Score Estimation

While traditional methods, such as logistic regression, have remained standard ways to estimate the propensity score, the method has continued to be improved by the application of machine learning models that include random forests, gradient boosting machines, and neural networks. These newer models now capture more complex relationships that covariates may take on in the data, undetected by the linear model, accounting for interactions and nonlinearities. The random forest model can be trained on covariates, including patient demographics, prior medical history, and hospital characteristics, to predict the probability of being treated. This

would give a better propensity score, as it leverages the power of machine learning for accurate predictions and gives out better-matched treatment and control groups.

4.3.2 Double Machine Learning

Another recent development takes the combination of DID with the estimation of machine learning models for estimating causal effects. If so, this is what Double Machine Learning can attain, using machine learning to do the control for confounding and correcting biases that may emerge when using complex high-dimensional data. In healthcare, where one may have large datasets including many covariates and interactions, DML refines estimates by allowing the DID framework to handle nonlinear and high-dimensional confounding effects. Integrating machine learning techniques into the DID framework thus allows for the accurate and reliable estimation of such effects, especially when analyzing large-scale and complexly interacted healthcare data.

5. Limitations of DID Methodology

5.1 Parallel Trends Assumption

One of the important identifying assumptions of DID is the so-called parallel trends assumption. It assumes that, in the absence of the treatment, the study and control groups would have followed the same trends in outcomes. This means the DID estimates could potentially be biased if this is a violated assumption. In healthcare, the parallel trend assumption might be very vulnerable to violation due to exogenous factors such as technological changes or new clinical guidelines that may affect both groups at different times and rates.

5.2 Time-Varying Confounders

DID assumes that confounders affecting the outcome remain constant over time. However, in healthcare settings, time-varying confounders like new treatments, changes in patient demography, or shift in policy by hospitals can distort the results. Dynamic DID models and machine learning techniques can help to overcome this weakness by specifically allowing for time-varying confounding.

5.3 Heterogeneity of Treatment Effects

The most standard DID estimate assumes homogeneity in treatment effects at the individual level. Clearly, in healthcare, treatment effects may well vary across patients due to variation in genetics, socio-economic status, or other pre-existing health conditions. Dynamic DID models and machine learning methods that incorporate treatment effect heterogeneity can give more precise estimates by capturing such variation.

6. Application of DID in Estimating Cost Savings for Healthcare

6.1 Case Study: Reduction in Hospital Readmission Costs

As a practical application of DID within healthcare, we discuss here a case study on the reduction in hospital readmission costs. A policy intervention imposed financial

penalties on every hospital that showed high readmission rates with the aim of decreasing healthcare costs.

Study Group: The group of hospitals which were under financial penalty.

Control Group: It includes those hospitals for which the policy did not take effect.

6.2 Data Collection and Design

We have collected data with respect to the readmission rates of the hospitals along with their costs, both for study and control groups, for a period of 5 years. The first 2 years represent the pre-intervention period and the remaining 3 years represent the post-intervention period. Propensity scores were calculated based on different attributes of the hospitals such as size, patient demographics, and geographic location. Inverse propensity score weighting was applied to adjust for the differences between the groups.

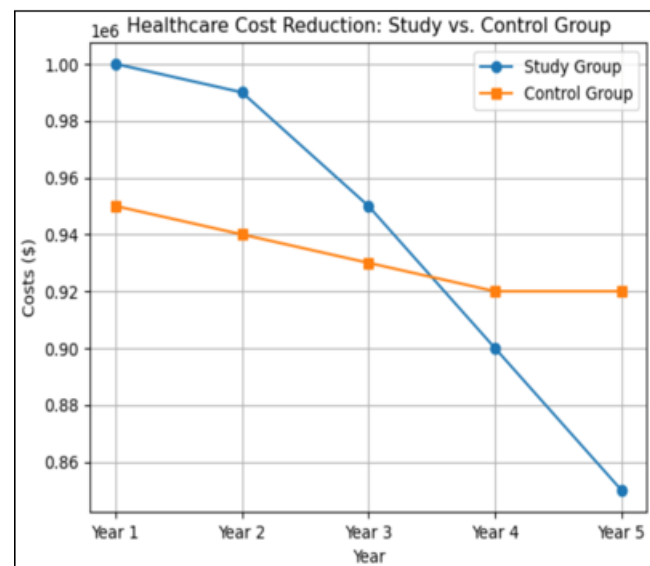
6.3 Results and Analysis

The DID analysis showed a significant reduction in readmission rates and costs related to readmission in the study group in comparison with the control group. The results of this are summarized below:

Group	Pre-Intervention Costs	Post-Intervention Costs	Cost Reduction (%)
Study Group	\$1,000,000	\$850,000	15%
Control Group	\$950,000	\$920,000	3%

6.4 Graphical Representation

The following graph shows the reduction in costs over time for the study and control groups.



6.5 Interpretation

The results on the costs of intervention proved that, because of the intervention, the costs were considerably reduced. For the treatment group hospitals, there was a reduction in 15% of the costs compared to the control group, which had a cost

reduction of about 3%, thus proving that the policy was very effective in reducing the readmission costs.

7. Technologies for Implementing DID

Effective implementation of the Difference-in-Differences (DID) methodology requires a combination of technological tools to handle large healthcare datasets and ensure reliable, scalable operations.

7.1 Site Reliability Engineering (SRE)

Site Reliability Engineering (SRE) applies software engineering principles to maintain system reliability, particularly in large-scale DID studies. Key practices include:

- **Automation:** Automating data cleaning, processing, and analysis tasks reduces human error and ensures repeatability.
- **Monitoring and Observability:** Tools like **Prometheus** and **Grafana** monitor system performance during data processing, allowing for real-time detection of issues.
- **Scalability:** SRE principles ensure systems can dynamically scale to handle the size and complexity of healthcare datasets without failures.

SRE helps ensure resilience and efficiency, making it a crucial component for handling large-scale DID implementations in healthcare.

7.2 Cloud Computing

Cloud computing platforms like **Amazon Web Services (AWS)**, **Google Cloud**, and **Microsoft Azure** provide scalable computational and storage resources for processing large datasets. These platforms support:

- **Elastic Computing:** Scalable resources can be provisioned as needed, allowing researchers to handle vast healthcare datasets efficiently.
- **Data Storage:** Secure, compliant storage solutions (e.g., **AWS S3**, **Google Cloud Storage**) ensure sensitive healthcare data is safely stored and accessible for analysis.

7.3 Data Modeling and ETL Pipelines

Data modeling and **ETL (Extract, Transform, Load)** pipelines streamline the preparation of large healthcare datasets for DID analysis. Tools such as **Apache NiFi** and **AWS Glue** automate the collection, cleaning, and structuring of data from diverse sources like hospital records and insurance claims, ensuring consistency and accuracy.

7.4 Machine Learning for Propensity Score Calculation

Machine learning techniques, including **random forests** and **gradient boosting machines (GBMs)**, can improve the accuracy of **propensity score calculation** by capturing non-linear relationships between variables. These models enhance the balance between study and control groups, increasing the reliability of DID results.

7.5 Data Visualization Tools

Data visualization tools like **Tableau**, **Power BI**, and programming languages such as **Python** and **R** help

researchers communicate DID results clearly. These tools generate visualizations that illustrate trends, causal effects, and cost-saving impacts derived from DID analysis, aiding in decision-making processes.

8. Conclusion

DID is a remarkable way of getting the causal effect of healthcare policies meant for cost-cutting analysis. Matching and weighting by propensity score balance the treatment and comparison groups to ameliorate some of the problems in DID. In addition, the introduction of machine learning methods has enhanced DID's ability to analyze large, complex health datasets. However, despite its several limitations, DID remains a major approach in healthcare policy analysis, particularly when applied with advanced treatments and controls for balancing.

References

- [1] Card, D., & Krueger, A. B. (1994). Minimum wages and employment: A case study of the fast-food industry in New Jersey and Pennsylvania. *American Economic Review*, 84(4), 772-793.
- [2] Finkelstein, A., Taubman, S., Wright, B., & Baicker, K. (2012). The Oregon Health Insurance Experiment: Evidence from the first year. *The Quarterly Journal of Economics*, 127(3), 1057-1106.
- [3] Bertrand, M., Duflo, E., & Mullainathan, S. (2004). How much should we trust differences-in-differences estimates? *Quarterly Journal of Economics*, 119(1), 249-275.
- [4] Stuart, E. A. (2010). Matching methods for causal inference: A review and a look forward. *Statistical Science*, 25(1), 1-21.
- [5] Wing, C., Simon, K., & Bello-Gomez, R. A. (2018). Designing difference-in-difference studies: Best practices for public health policy research. *Annual Review of Public Health*, 39, 453-469.
- [6] Angrist, J. D., & Pischke, J. S. (2008). *Mostly harmless econometrics: An empiricist's companion*. Princeton University Press.
- [7] Ryan, A. M., Burgess, J. F., & Dimick, J. B. (2015). Why we should not be indifferent to specification choices for difference-in-differences. *Health Services Research*, 50(4), 1211-1235.
- [8] Athey, S., & Imbens, G. W. (2016). Recursive Partitioning for Heterogeneous Causal Effects. *Proceedings of the National Academy of Sciences*, 113(27), 7353-7360.
- [9] Jayanna Hallur, "The Future of SRE: Trends, Tools, and Techniques for The Next Decade", *International Journal of Science and Research (IJSR)*, Volume 13 Issue 9, September 2024, pp. 1688-1698, <https://www.ijsr.net/getabstract.php?paperid=SR24927125336>