

# Greed is Good: Estimating Forward Difference-in-Differences in Stata

Jared Greathouse  
Georgia State University  
Atlanta, GA  
jgreathouse3@student.gsu.edu

Jason Coupet  
Georgia State University  
Atlanta, GA  
jcoupet@gsu.edu

Eric Sevigny  
Georgia State University  
Atlanta, GA  
esevigny@gsu.edu

**Abstract.** Difference-in-differences designs build counterfactuals by invoking a parallel trend assumption, but this may be violated in the presence of invalid control units. Thus, selecting a control group is vital to ensure proper identification. We introduce `fdid` based on (Li, Kathleen T. "A Simple Forward Difference-in-Differences Method." *Marketing Science* 43, no. 2 (2024): 267-279). We discuss estimation and inference, document `fdid`'s syntax, and apply it empirically.

**Keywords:** difference-in-differences, synthetic control methods, causal inference

## 1 Introduction

For identification, Difference-in-Differences designs (DiD) make some form of a parallel trends assumption (PTA), assuming a constant difference between the average of the control group and treated outcome trajectories absent treatment. Unfortunately, DiD's PTA is invalid in many realistic scenarios, such as retail Costa et al. (2023), where the control group average may differ substantially from the treatment group due to inappropriate controls. Control group selection has become of interest recently to researchers. Shi and Huang (2023) extend Hsiao et al. (2012) by developing a forward selected panel data approach. Synthetic control methods (SCMs, Abadie [2021]) typically rely on a (usually) convex average of some controls to impute the counterfactual.

To better justify DiD's PTA, Li (2024) develops the forward DiD method (FDID), showing that a forward-selection method can be used to select the control group. We introduce the `fdid` method for Stata. `fdid` fits in with Stata's pantheon of program evaluation tools. Like `rcm` and `synth2` by Yan and Chen (2022, 2023), `scul` by Greathouse (2022), and `sdid` by Clarke et al. (n.d.), `fdid` uses a subset of controls to estimate the causal impact. Also, `fdid` returns the list of selected controls, graphics, and fit statistics. However, `fdid` is more user-friendly. `rcm`, `synth2`, and `allsynth` by Wiltshire (2021) all require users to specify the panel id for the treatment unit and treatment date, whereas `fdid` simply requires a dummy variable. `fdid` has more flexible data requirements, only requiring outcome data. This is in contrast to SCMs, for example, which frequently depends on covariates for acceptable pre-treatment fit (Yan and Chen 2022; Amjad et al. 2018). `fdid` is also fast, relying on bivariate OLS for estimation. In contrast, methods such as `fekt` by Liu et al. (2024) or `scul` by Greathouse (2022) employ cross validation or LASSO penalization.

## 2 Forward Difference-in-Differences

---

**Algorithm 1:** Forward Difference-in-Differences (FDID)

---

```

for  $k = 1$  to  $N_0$  do
  for  $i \in \mathcal{N}_0 \setminus \widehat{U}_k$  do
    | Estimate  $y_{1t} = \hat{\alpha}_{\mathcal{N}_0} + \bar{y}_{\widehat{U}_k \cup \{i\}}$   $t \in \mathcal{T}_1$  and calculate  $R_k^2(\widehat{U}_k \cup \{i\})$ 
  end
  Update  $\widehat{U}_k \leftarrow \widehat{U}_k \cup \left\{ \operatorname{argmax}_{i \in \mathcal{N}_0 \setminus \widehat{U}_k} R_k^2(\widehat{U}_k \cup \{i\}) \right\};$ 
end
Set  $\widehat{U}^* \leftarrow \operatorname{argmax}_{k \in \{1, \dots, N_0\}} R_k^2(\widehat{U}_k);$ 
Compute  $\hat{y}_{1t}^0 = \hat{\alpha}_{\widehat{U}^*} + \bar{y}_{\widehat{U}^*};$ 
return  $\widehat{U}^*$  and  $\widehat{ATT}_{\widehat{U}^*} = \frac{1}{T_2} \sum_{t \in \mathcal{T}_2} (y_{1t} - \hat{y}_{1t}^0(\widehat{U}))$ 

```

---

**The Model** We follow Li (2024)’s exposition, observing  $\mathcal{N} = \{1, 2, \dots, N\}$  units where  $\mathcal{N}$  has cardinality  $N = |\mathcal{N}|$ .  $j = 1$  is treated and controls are  $\mathcal{N}_0 = \mathcal{N} \setminus \{1\}$ . Time is indexed by  $t$ . Denote pre-post-policy periods as  $\mathcal{T}_1 = \{1, 2, \dots, T_0\}$  and  $\mathcal{T}_2 = \{T_0 + 1, \dots, T\}$ , where  $\mathcal{T} = \mathcal{T}_1 \cup \mathcal{T}_2$ . We use Algorithm 1 to select controls for DiD, denoted as  $\widehat{U} \subset \mathcal{N}_0$ .<sup>1</sup> Example 1 offers a stylized explanation of Algorithm 1, but we also summarize it below by quoting verbatim from `fdid`’s help file.

To begin, we first estimate  $N_0$  one-unit DiD models using each control. We calculate the pre-treatment R-squared statistic for each of these DiD models. Whichever control/model has the highest pre-treatment R-squared is the first selected control. This model is also the first “candidate” DiD model.

Next, we estimate  $N_0 - 1$  DiD submodels, where we use the first selected control along with each one of the remaining  $N_0 - 1$  controls in  $N_0 - 1$  two-unit DiD submodels. Whichever of these  $N_0 - 1$  two-unit DiD submodels has the highest pre-intervention R-squared statistic is the second candidate DiD model. We then add the maximizing control unit to the list of selected controls. Now there are two selected controls.

For the third iteration, we estimate  $N_0 - 2$  submodels, one for each of the remaining  $N_0 - 2$  controls unit plus the two already selected controls. We calculate R-squared statistics for these  $N_0 - 2$  submodels, and see which third control plus the two selected ones maximizes the pre-treatment R-squared statistic. The best submodel of these is now the third candidate DiD model. This process continues until there are  $N_0$  candidate DiD models. The final model/control group FDID uses is whichever one of the  $N_0$  candidate DID models that has the highest pre-intervention R-squared statistic. Post-selection, Li (2024) estimates FDID like

$$y_{1t} = \hat{\alpha}_{\widehat{U}} + \bar{y}_{\widehat{U}_t} \quad t \in \mathcal{T}_1 \tag{1}$$

---

1.  $\widehat{U}$  from Algorithm 1 is used as the subscript to emphasize that we are including only the selected controls. We omit the asterisk for simplicity.

where  $\bar{y}_{\hat{U}t} := \frac{1}{|\hat{U}|} \sum_{j \in \hat{U}} y_{jt}$  is the average of the selected controls. The estimated least-squares intercept is computed like  $\hat{\alpha}_{\hat{U}} := T_1^{-1} \sum_{t \in \mathcal{T}_1} (y_{1t} - \bar{y}_{\hat{U}t})$ . Denote the FDID predictions as  $\hat{y}_{1t}^0 = \hat{\alpha}_{\hat{U}} + \bar{y}_{\hat{U}t}$ , where the pre-treatment periods corresponds to the in-sample fit and the opposite denotes the out-of-sample counterfactual. Our causal estimand is:  $\widehat{ATT}_{\hat{U}} = \frac{1}{T_2} \sum_{t \in \mathcal{T}_2} (y_{1t} - \hat{y}_{1t}^0)$ . From Assumption 2.1 of Li (2024) and Arkhangelsky et al. (2021, 4094), FDID assumes *parallel* trends,  $\hat{y}_{1t}^0 - \bar{y}_{\hat{U}t} = \hat{\alpha}_{\hat{U}} + \epsilon$ .<sup>2</sup>

**Example 1.** Let  $\mathcal{N}_0 = \{i_1 \text{ (Chicago)}, i_2 \text{ (Miami)}, i_3 \text{ (Phoenix)}\}$  be the controls for a generic treated unit. For ( $k = 1$ ), we estimate DiD for each control unit in  $\mathcal{N}_0$  individually, yielding pre-treatment  $R^2$  values:  $R_{1,1}^2 = 0.60$ ,  $R_{2,1}^2 = 0.50$ , and  $R_{3,1}^2 = 0.23$ . Since  $R_{1,1}^2 = 0.60$  is the highest, we update the control set to  $\hat{U}_1 = \{i_1\}$  and  $R_k^2 = 0.60$ . For ( $k = 2$ ), we estimate two DiD models using  $i_1$  with the remaining controls from  $\{i_2, i_3\}$ , yielding  $R_{2,2}^2 = 0.88$  and  $R_{3,2}^2 = 0.68$ . We select  $i_2$  (Miami) and update the control set to  $\hat{U}_2 = \{i_1, i_2\}$  since  $R_{2,2}^2 = 0.88$  is the highest. For ( $k = 3$ ), using all controls, we get  $R_{3,3}^2 = 0.55$ . The final control set is  $\hat{U}_2 = \{i_1, i_2\}$ , as  $\max_k R_k^2 = 0.88$ .

**Inference** Per Li (2024), our default standard error for the ATT is:

$$\hat{\Omega} = \left[ \left( \frac{T_2}{T_1} \right) \cdot T_1^{-1} \sum_{t \in \mathcal{T}_1} \hat{v}_{1t}^2 \right]^{0.5}, \quad \hat{v} = y_{1t} - \bar{y}_{\hat{U}} - \hat{\alpha}_{\hat{U}} \quad (2)$$

Li (2024) establishes the normal inference theory of the FDID method (see appendices B and D for theoretical derivations). Users of `fdid` may use the placebo algorithm from Arkhangelsky et al. (2021) to generate the standard errors.

### 3 The fdid command

Users need strongly balanced panel data (see [XT] `xtset`). `sdid_event` must be installed. Users also need Stata 16 or later.

#### 3.1 Syntax

```
fdid depvar [if] [in] treated(varname) [unitnames(string) gr1opts(string)
gr2opts(string) placebo ]
```

where *depvar* is our dependent variable and *treated* is our dummy for treatment.

---

2. SCMs generally attempt to *match* the counterfactual to the pre-treatment trajectory.

**gr1opts:** Edits the display options of the observed versus predicted plot.

**gr2opts:** See the above, except for the plotted pointwise-treatment effect.

**unitnames:** The string variable that serves as the value labels (required if the panel id is not already labeled). Note each string value pair must be uniquely identified.

**placebo:** Uses the placebo standard error of the ATT from Arkhangelsky et al. (2021).

Matrices			
<b>e(results)</b>	DID/FDID results	<b>e(b)</b>	Coefficients
<b>e(V)</b>	variance-covariance matrix	<b>e(dyneff)</b>	dynamic effects
<b>e(series)</b>	means/counterfactuals	<b>e(setting)</b>	pre-treatment periods, treatment date, post-treatment periods, number of time periods
Macros			
<b>e(U)</b>	selected controls	<b>e(depvar)</b>	dependent variable
<b>e(properties)</b>	list of properties		

We replicate Abadie et al. (2010) for two reasons: firstly, the basic results of DiD are not in dispute, being quite popular in the econometrics literature for introducing the SCM or shortcomings of DiD. More importantly, Abadie et al. (2010) explicitly say DiD’s PTA is invalid. Since the point of FDID is to choose controls such that standard PTA is more credible, Abadie et al. (2010) presents a good avenue to demonstrate how `fdid` is useful for Stata users. We begin with loading in the dataset, obtained from the syntax from section 6.

The following output from `xtdescribe` displays the panel setup for `smoking.dta`.

```

id: 1, 2, ..., 39
year: 1970, 1971, ..., 2000
Delta(year) = 1 year
Span(year) = 31 periods
(id*year uniquely identifies each observation)

Distribution of T_i:  min      5%      25%      50%      75%      95%      max
                    31      31      31      31      31      31      31

  Freq.  Percent  Cum. | Pattern
-----+-----
    39   100.00  100.00 | 11111111111111111111111111111111
-----+-----
    39   100.00      | XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX

```

California is treated in 1989, compared  $N_0 = 38$  states that remain untreated. Time extends from 1970 to 2000, so  $T_1 = 19$  and  $T_2 = 12$ . Our outcome is the rate of tobacco consumption. We estimate `fdid` like

```
fdid cigsale, tr(treated) unitnames(state)
```

```
Forward Difference-in-Differences          TO R2:      0.988      TO RMSE:      1.282
```

cigsale	ATT	Std. Err.	t	P> t	[95% Conf. Interval]
treated	-13.64671	0.46016	29.66	0.000	-14.54861 -12.74481

Treated Unit: California

FDID selects Montana, Colorado, Nevada, Connecticut, as the optimal donors.

See Li (2024) for technical details.

We plot the predictions from both DiD and FDID as well as their control group means; for both FDID and DiD, the pre-1989 predictions are the in sample fit and the post-1989 values are the out-of-sample counterfactual predictions.<sup>3</sup> The results appear in Figure 1. DiD's in-sample prediction misses the observed in-sample values

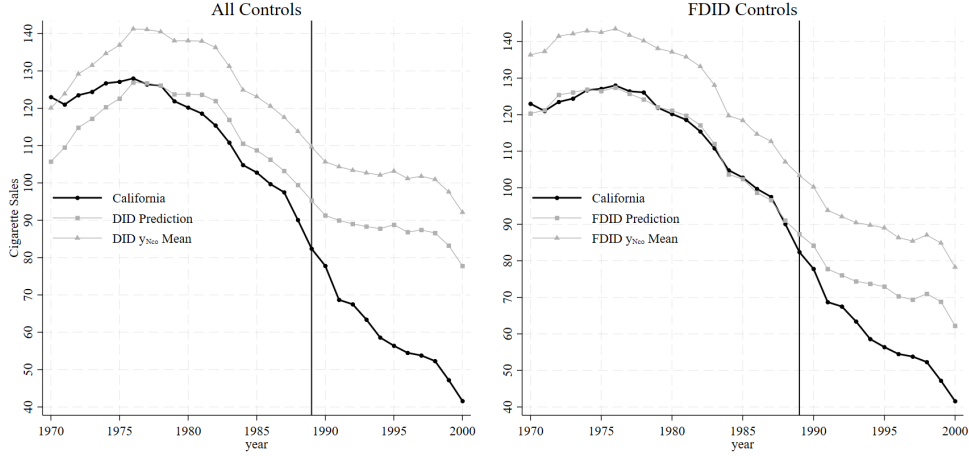


Figure 1: Observed, Predicted, and Average Curves

of California for the first 5 years of the time series and overestimates them from the mid-1970s until 1989. Abadie et al. (2010) also remark that in 1988, the rest of the United States has a 27% higher consumption rate than California. Given DiD's  $R^2$  is equal to 60%, this comports with Abadie et al. (2010)'s conclusion that that PTA is not viable for all controls in this case. The DiD ATT is  $\widehat{ATT}_{N_0} = -27.349$ , a value

3. We omit the code in order to save space, but see `FDID_SJ_Rep.do` at the first author's GitHub.

which is likely overestimated. When we view the results of `xtddidregress` by doing `mat 1 e(results)` after running `fdid`, we get a 95% CI for DiD of  $[-33.02, -21.68]$ . Note that `xtddidregress` uses the robust standard error.

The forward selection method chooses 4 control units: Montana, Colorado, Nevada, and Connecticut (*all* of which were given weight by the original SCM). The pre-intervention average of these units is obviously parallel to the pre-trends of California, per the plot. This fact is supported by  $R^2_{\hat{\theta}}$ , which says 98.8% of the pre-intervention variance is explained by the intercept shifted average of the new control group. For FDID,  $\widehat{ATT}_{\hat{\theta}} = -13.647$ , a reduction of the original DiD ATT by half. FDID’s 95% CI is  $[-14.55, -12.74]$ . Another point to note is how FDID’s in-sample PTA seems to hold without any covariates or predictors, suggesting that FDID’s data requirements are, in some cases, more relaxed compared to SCM whose methods typically rely on predictors for convergence (Amjad et al. 2018; Vives-i-Bastida 2022), or DiD where analysts sometimes make a conditional PTA.

## 5 Conclusion

We wish to make clear the central limitation of `fdid`: its PTA must still be valid. As per usual, researchers should check if the standard DiD PTA is plausible first. Researchers who have found DiD’s PTA to be invalid should then check if PTA holds for FDID in the pre-intervention period. Li (2024) notes that if researchers have a treated unit whose trend is much steeper than control units, for example, then use of `fdid` is invalid. Researchers may then consider more flexible methods, such as factor models or modified SCMs (Li and Shankar 2024). However, even if `fdid`’s PTA is plausible, other methods such as `synth2` may also serve as a robustness check.

While `fdid` is useful, we now highlight FDID’s limitations and opportunities for development. For staggered adoption, Li (2024) is silent on whether using the not yet treated controls vs. never treated controls would be preferable, or on how to weight ATTs across multiple units (Wing et al. 2024). While `fdid` use the never treated by default and reports Cohort ATTs, we more formal investigation of the FDID method and staggered adoption is warranted. Also, some newer methods invoke *conditional* parallel trends assumptions, where covariates are included (Callaway and Sant’Anna 2021). FDID does not adjust for covariates at present. FDID also does not account for settings where units may be treated and then untreated, or where units receive non-binary treatments as discussed in de Chaisemartin and D’Haultfœuille (2024) and D’Haultfœuille et al. (2023). In terms of selection algorithms, Li (2024) notes other control group selection methods may be used such as the recently user-written `classifylasso` by Huang et al. (2024) (naturally, a comparison is outside the scope of our paper).

We introduced the `fdid` command whose algorithm selects a control group for DiD. We overviewed the syntax of `fdid` and applied it empirically to a setting where the classical PTA was too restrictive to deliver satisfactory results. Given `fdid`’s practical benefits, we believe `fdid` is of use to Stata users who are interested in treatment effect estimation.

## 6 Program Installation

```
net from "https://raw.githubusercontent.com/jgreathouse9/FDIDTutorial/main"
net install fdid
net get fdid, replace
```

## 7 References

- Abadie, A. 2021. Using Synthetic Controls: Feasibility, Data Requirements, and Methodological Aspects. *Journal of Economic Literature* 59(2): 391–425. [10.1257/jel.20191450](https://doi.org/10.1257/jel.20191450).
- Abadie, A., A. Diamond, and J. Hainmueller. 2010. Synthetic Control Methods for Comparative Case Studies: Estimating the Effect of California’s Tobacco Control Program. *Journal of the American Statistical Association* 105(490): 493–505. <https://doi.org/10.1198/jasa.2009.ap08746>.
- Amjad, M., D. Shah, and D. Shen. 2018. Robust synthetic control. *The Journal of Machine Learning Research* 19(1): 802–852.
- Arkhangelsky, D., S. Athey, D. A. Hirshberg, G. W. Imbens, and S. Wager. 2021. Synthetic Difference-in-Differences. *American Economic Review* 111(12): 4088–4118. <https://doi.org/10.1257/aer.20190159>.
- Callaway, B., and P. H. Sant’Anna. 2021. Difference-in-Differences with multiple time periods. *Journal of Econometrics* 225(2): 200–230. <https://doi.org/10.1016/j.jeconom.2020.12.001>.
- de Chaisemartin, C., and X. D’Haultfœuille. 2024. Difference-in-Differences Estimators of Intertemporal Treatment Effects. *The Review of Economics and Statistics* 1–45. [https://doi.org/10.1162/rest\\_a\\_01414](https://doi.org/10.1162/rest_a_01414).
- Clarke, D., D. Pailańir, S. Athey, and G. Imbens. n.d. On Synthetic Difference-in-Differences and Related Estimation Methods in Stata. [working paper]. <https://doi.org/10.48550/arXiv.2301.11859>.
- Costa, L., V. F. Farias, P. Foncea, J. D. Gan, A. Garg, I. R. Montenegro, K. Pathak, T. Peng, and D. Popovic. 2023. Generalized Synthetic Control for TestOps at ABI: Models, Algorithms, and Infrastructure. *INFORMS Journal on Applied Analytics* 53(5): 336–349. <https://doi.org/10.1287/inte.2023.0028>.
- D’Haultfœuille, X., S. Hoderlein, and Y. Sasaki. 2023. Nonparametric difference-in-differences in repeated cross-sections with continuous treatments. *Journal of Econometrics* 234(2): 664–690. <https://doi.org/10.1016/j.jeconom.2022.07.003>.
- Greathouse, J. 2022. Scul: Regularized Synthetic Controls in Stata. [working paper]. <https://doi.org/10.2139/ssrn.4196189>.

- Hsiao, C., H. Steve Ching, and S. Ki Wan. 2012. A PANEL DATA APPROACH FOR PROGRAM EVALUATION: MEASURING THE BENEFITS OF POLITICAL AND ECONOMIC INTEGRATION OF HONG KONG WITH MAINLAND CHINA. *Journal of Applied Econometrics* 27(5): 705–740.
- Huang, W., Y. Wang, and L. Zhou. 2024. Identify latent group structures in panel data: The classifylasso command. *The Stata Journal* 24(1): 46–71. <https://doi.org/10.1177/1536867X241233642>.
- Li, K. T. 2024. Frontiers: A Simple Forward Difference-in-Differences Method. *Marketing Science* 43(2): 239–468. <https://doi.org/10.1287/mksc.2022.0212>.
- Li, K. T., and V. Shankar. 2024. A Two-Step Synthetic Control Approach for Estimating Causal Effects of Marketing Events. *Management Science* 70(6): 3734–3747.
- Liu, L., Y. Wang, and Y. Xu. 2024. A Practical Guide to Counterfactual Estimators for Causal Inference with Time-Series Cross-Sectional Data. *American Journal of Political Science* 68(1): 160–176. <https://doi.org/10.1111/ajps.12723>.
- Shi, Z., and J. Huang. 2023. Forward-selected panel data approach for program evaluation. *Journal of Econometrics* 234(2): 512–535.
- Vives-i-Bastida, J. 2022. Predictor Selection for Synthetic Controls. [working paper]. <https://arxiv.org/abs/2203.11576>.
- Wiltshire, J. C. 2021. allsynth: Synthetic Control Bias-Correction Utilities for Stata. [working paper].
- Wing, C., S. M. Freedman, and A. Hollingsworth. 2024. Stacked Difference-in-Differences. Working Paper 32054, National Bureau of Economic Research. <http://www.nber.org/papers/w32054>.
- Yan, G., and Q. Chen. 2022. rcm: A command for the regression control method. *The Stata Journal* 22(4): 842–883. <https://doi.org/10.1177/1536867X221140960>.
- . 2023. synth2: Synthetic control method with placebo tests, robustness test, and visualization. *The Stata Journal* 23(3): 597–624. <https://doi.org/10.1177/1536867X231195278>.

#### About the authors

Jared Greathouse is a PHD Candidate of Public Policy at Georgia State University. He studies causal inference, econometrics, and machine-learning.

Jason Coupet is an Associate Professor of Public Management and Policy in the Andrew Young School of Policy Studies at Georgia State University. His research interests include strategic management, efficiency, applied microeconomics, organizational economics, public sector management science, and the political economy of organizations.

Eric Sevigny...