

# Bunching estimation of elasticities using Stata

Marinho Bertanha  
University of Notre Dame  
Notre Dame, Indiana  
mbertanha@nd.edu

Andrew H. McCallum  
Board of Governors of the  
Federal Reserve System  
Washington, DC  
andrew.h.mccallum@frb.gov

Alexis Payne  
Stanford University  
Stanford, CA  
ampayne@stanford.edu

Nathan Seegert  
University of Utah  
Salt Lake City, Utah  
nathan.seegert@eccles.utah.edu

**Abstract.** Typical censoring models have mass-points at the upper, lower, or both tails of an otherwise continuous outcome distribution. In contrast, we consider a censoring model with a mass-point in the interior of the outcome distribution. We refer to this mass point as “bunching” and use it to estimate model parameters. For example, economic theory suggests that for increasing marginal income tax rates, many taxpayers will report income exactly at the threshold where the tax rate increases. This translates into a censoring model with bunching at the threshold. The size of this mass point of taxpayers can be used to estimate an elasticity parameter, which summarizes taxpayers responses to taxes. This article introduces the command `bunching`, which implements new non-parametric and semi-parametric identification methods for estimating elasticities developed by Bertanha, McCallum, and Seegert (2021). These methods rely on weaker assumptions than what are currently made in the literature and result in meaningfully different estimates of the elasticity.

**Keywords:** `st0001`, bunching, bunchbounds, bunchtobit, bunchfilter, partial identification, censored regression, income elasticity, tax

## 1 Introduction

Censoring models apply to distributions of an outcome variable that are continuous except for a mass-point at the upper, lower, or both tails of the distribution. This paper considers models where the mass-point occurs in the interior of the outcome distribution. We refer to this class of models as “mid-censoring models.” Although we use the adjective “mid-censoring”, the mass point may be at any point in the interior of the support of the distribution of outcomes.

Previously developed methods use such a mass point, often called “bunching”, to estimate model parameters. For example, economic theory suggests that for increasing marginal income tax rates, many taxpayers will report income exactly at the threshold where the tax rate increases. This translates to a mid-censoring model with a mass-point in the interior of the distribution of reported income. The size of this mass point

can be used to identify an important parameter of the censoring model, which is known to economists as an elasticity parameter. In this context, the elasticity parameter describes the percent change in reported income in response to a percentage point change in marginal income tax rate. More specifically, an elasticity of 0.5 means that taxpayers reduce their reported income (labor supply) by 0.5 percent for each 1 percentage point increase in marginal income tax rates. Section 3.1 provides simulated data and a numerical example interpreting this elasticity in more detail. In the rest of this paper, we use “bunching” to refer to a mass point in the interior of an outcome distribution, and “bunching methods” or “bunching estimator” to refer to the statistical methods that recover elasticity parameters from data that exhibit bunching.

Using bunching to estimate elasticities began with Saez (2010), Chetty et al. (2011), and Kleven and Waseem (2013). Following these influential papers, bunching methods became a popular way to estimate elasticities in a variety of settings ranging from electricity demand (Ito 2014), real estate taxes (Kopczuk and Munroe 2015), labor regulations (Garicano et al. 2016), and prescription drug insurance (Einav et al. 2017) to marathon finishing times (Allen et al. 2017), attribute-based regulations (Ito and Sallee 2018), education (Dee et al. 2019; Caetano et al. 2020a), minimum wage (Jales 2018; Cengiz et al. 2019), and air-pollution data manipulation (Ghanem et al. 2019), among others. Differences in mass point sizes across groups has been exploited as the first stage in a two-stage approach to control for endogeneity (Chetty et al. 2013; Caetano 2015; Grossman and Khalil 2020). Bunching has also been used for causal identification in Khalil and Yildiz (2020), Caetano and Maheshri (2018), Caetano et al. (2019), Caetano et al. (2020b), and Jales and Yu (2017) connects bunching to regression discontinuity (RD). Lastly, Kleven (2016) conducts a detailed review of the bunching literature.

This paper introduces a new Stata command, **bunching**, which utilizes assumptions that are weaker than current bunching methods. The command **bunching** is a wrapper function for three other commands. The first of those commands is **bunchbounds**, which estimates upper and lower bounds on the bunching elasticity using a partial-identification approach. The second is **bunchtobit**, which uses a semi-parametric method with covariates for point identification. The third is **bunchfilter**, which filters friction errors from the dependent variable before applying either **bunchbounds** or **bunchtobit**.

The statistical foundations for these commands are developed by Bertanha, McCallum, and Seegert (2021). That paper introduces multiple methods to recover elasticities from bunching. Each method relies on different assumptions to achieve identification of the elasticity. Since these are assumptions about an unobserved distribution, it is not possible to determine which assumption is correct. However, it is possible to check if estimates relying on different assumptions are robust across assumptions. In practice, we recommend that researchers use the **bunching** package to employ different estimation methods and check that elasticity estimates they recover are stable across those methods.

## 2 Bunching estimators

The application of bunching methods utilized by Bertanha, McCallum, and Seegert (2021) and this paper derives from bunching behavior caused by progressive marginal income taxes, as in Saez (2010). Formally, agents maximize an iso-elastic quasi-linear utility function of total consumption (or disposable income) and labor, which results in a data generating process (DGP) for optimal reported taxable income as follows

$$y_i = \begin{cases} \varepsilon s_0 + n_i^*, & \text{if } n_i^* < \underline{n}(k, \varepsilon, s_0) \\ k, & \text{if } \underline{n}(k, \varepsilon, s_0) \leq n_i^* \leq \bar{n}(k, \varepsilon, s_1) \\ \varepsilon s_1 + n_i^*, & \text{if } n_i^* > \bar{n}(k, \varepsilon, s_1). \end{cases} \quad (1)$$

in which  $y_i$  is the log of reported income,  $n_i^*$  is the log of unobserved heterogeneity of agent  $i$ ,  $\varepsilon$  is the elasticity parameter of interest, the log of the slope of the piecewise-linear constraint changes from  $s_0$  to  $s_1$  at the log of the kink point  $k$ , and  $s_1 < s_0$ . All logs in this paper are natural logs. The restriction  $s_1 < s_0$  guarantees concavity of the budget set, which is fundamental for the solution in Equation 1. In the original tax application,  $s_j = \log(1 - t_j)$ ,  $j \in \{0, 1\}$ , in which  $t_j$  is the marginal tax rate and  $t_0 < t_1$ . The expressions for the thresholds that determine the three cases in Equation 1 are  $\underline{n}(k, \varepsilon, s_0) = k - \varepsilon s_0$  and  $\bar{n}(k, \varepsilon, s_1) = k - \varepsilon s_1$ .

We use utility maximizing agents and income-taxes to motivate Equation 1 and for exposition of the command throughout the rest of this paper. Nevertheless, the methods developed by Bertanha, McCallum, and Seegert (2021), as well as the **bunching** package, apply to any data set generated by Equation 1. We emphasize that any data must be transformed into units that satisfy Equation 1. In the income-tax example, this is accomplished by taking logs of the outcome variable, kink, and slopes.

Our methods are applicable to non-tax data. For example, Bitler et al. (2021) study the Supplemental Nutrition Assistance Program (SNAP), in which low-income individuals receive benefits for food purchases as a function of labor income,  $y_i$ . The benefit is a constant amount for labor income less than a known value,  $k$ , but decreases linearly after that. This reduction in benefits creates a piece-wise linear budget set over total consumption and labor income with a kink. At  $y_i = k$ , the log of the slope changes from  $s_0$  to  $s_1$  with  $s_1 < s_0$  (see Bitler et al. 2021, Figure 1). In this case, bunching methods identify the elasticity of labor supply,  $\varepsilon$ , with respect to the benefit reduction rate.

Another non-income-tax application is Ito (2014), who studies consumption of electricity in Southern California. Electricity price per kilowatt-hour (kWh) changes as a function of quantity of consumption in kWh (see Figure 3 in his paper). This piece-wise linear pricing scheme creates a budget set over disposable income and electricity consumption with kinks, and bunching methods identify the demand elasticity with respect to electricity price.

Piece-wise linear constraints frequently exhibit several kinks at different locations. **bunching** can be applied to each kink separately as long as the constraint does not have a discontinuous jump —often called a “notch” —preceding the kink under study.

Appendix B of Bertanha et al. (2021) provides a general solution to a model with multiple kinks and notches and Section “3 Identification” of Bertanha et al. (2021) discusses inference for multiple kinks.

Our estimation methods rely on Equation 1, which maps the continuously distributed unobserved  $n_i^*$  into a mixed continuous-discrete observed distribution for  $y_i$  for given values of  $(s_0, s_1, k, \varepsilon)$ . For higher values of  $n_i^*$ , higher values of  $y_i$  will be observed except when  $n_i^*$  falls inside the bunching interval, that is,  $[\underline{n}(k, \varepsilon, s_0), \bar{n}(k, \varepsilon, s_1)]$ , in which case  $y_i$  remains constant and equal to  $k$ . Therefore, (1) leads to bunching in the distribution of  $y_i$  at the kink point  $k$ . In other words, the distribution of  $y_i$  has a mass point at  $k$ ,  $\mathbb{P}(y_i = k) > 0$ , but is continuous otherwise. The mass of the point at  $k$  depends on the size of the bunching interval according to

$$\begin{aligned} B &\equiv \mathbb{P}(y_i = k) = \mathbb{P}(\underline{n}(k, \varepsilon, s_0) \leq n_i^* \leq \bar{n}(k, \varepsilon, s_1)) \\ &= F_{n^*}(\bar{n}(k, \varepsilon, s_1)) - F_{n^*}(\underline{n}(k, \varepsilon, s_0)), \end{aligned} \quad (2)$$

in which  $F_{n^*}$  is the cumulative distribution function (CDF) of the unobserved  $n^*$ .

The data and model formally consist of five elements: (i) the CDF of the outcome  $F_y$ , (ii) the kink point  $k$ , (iii) the slopes of the budget constraint on the left,  $s_0$ , and right,  $s_1$ , of the kink point; (iv) the CDF of unobserved heterogeneity  $F_{n^*}$ , and (v) the elasticity  $\varepsilon$ . Equation 1 maps elements (ii)–(v) into the observed CDF,  $F_y$ . The researcher observes elements (i)–(iii), but not the last two elements,  $F_{n^*}$  and  $\varepsilon$ .

Original bunching estimators recover  $\varepsilon$  in two steps (Saez 2010; Chetty et al. 2011). First, they assume a specific function  $F_{n^*}$  over the bunching interval. Second, they invert Equation 2 to recover  $\varepsilon$  using their assumption about  $F_{n^*}$ . The methods developed by Bertanha, McCallum, and Seegert (2021) that are implemented by the **bunching** command are quite different than these original approaches.

**bunching** implements two novel identification strategies for the elasticity using a mass point at a kink.

The first strategy partially identifies the elasticity by assuming Lipschitz continuity and is implemented by **bunchbounds**. In other words, it assumes that the probability density function (PDF) of the unobserved heterogeneity has bounded slope magnitude. How this assumption recovers the elasticity is as follows. The observed bunching mass equals the area under the the heterogeneity PDF inside an interval. The size of this bunching interval is a function of the unknown elasticity parameter. The highest and lowest values for possible PDFs inside the bunching interval are set by the Lipschitz bound on the slope magnitude of the PDFs. With a fixed bunching mass, these PDF bounds determine the maximum and minimum widths of the bunching interval and imply lower and upper bounds for the elasticity. **bunchbounds** has two particularly valuable features. First, when bunching is observed the elasticity lower bound must be positive. Second, the bunching estimator based on the trapezoidal approximation (Saez 2010) is always within the bounds (partially identified set of elasticities).

The second strategy rewrites Equation 1 as a mid-censored regression model and is implemented by **bunchtobit**. The method assumes that the unobserved heterogeneity

conditional on covariates follows a normal distribution, but we prove that conditional normality is not required for consistency of the elasticity when the unconditional distribution of income is correctly specified. This approach effectively assumes that the unconditional distribution of heterogeneity belongs to a semi-parametric family of normal mixtures. Conditional normality implies a Tobit model that has a globally concave log likelihood that is easy to maximize. `bunchtobit` also truncates the sample using a sequence of smaller windows around the kink point. Consistency of the elasticity using these smaller windows requires weaker assumptions on the distribution of heterogeneity because the model tends to fit the unconditional distribution of income better as the window size decreases. To the best of our knowledge, this is the first bunching estimation strategy that utilizes covariates and semi-parametric assumptions to recover the elasticity. Covariates can control for a substantial amount of individual heterogeneity and `bunchtobit` only places assumptions on the remaining portion of heterogeneity that is unobserved. In general, researchers should prefer methods that control for observable heterogeneity using covariates over methods that omit covariates and instead restrict both observed and unobserved heterogeneity.

Many datasets have friction errors which imply that the bunching mass is dispersed in a small interval near, instead of exactly at, the kink. When friction errors are present, they must first be filtered out before a bunching estimation method can be applied. The procedure implemented by `bunchfilter` is a practical way of removing friction errors and works well when 1) the researcher has an accurate prior on the support of the friction error distribution, 2) the friction error affects non-bunching individuals more than it affects bunching individuals, or 3) the friction error has a small variance. A more general filtering method requires deconvolution theory, which is an active area of research.

## 2.1 The `bunchbounds` command

`bunchbounds` uses bunching to partially identify the elasticity of a response variable with respect to changes in the slope of the budget set. The syntax, options, and description of this command are as follows:

### Syntax for `bunchbounds`

```
bunchbounds depvar [if] [in] [weight] , kink(#) s0(#) s1(#) m(#) [ nopic
    savingbounds(filename[,replace]) ]
```

*depvar* must be one dependent variable (the response in logs in many applications).

`kink(#)` is the location of the kink point and must be a real number in the same units as the response variable.

`s0(#)` is a real number. In many applications, it is the log of the slope before the kink point.

`s1(#)` must be a real number that is strictly less than `s0(#)`. In many applications, it is the log of the slope after the kink point.

`m(#)` is the maximum magnitude of the heterogeneity PDF slope and must be a strictly positive real number.

Entries for `depvar`, `kink(#)`, `s0(#)`, `s1(#)`, and `m(#)` are required, whereas options inside the square brackets are not required.

### Options for `bunchbounds`

`if` and `in` restrict the working sample, like many other Stata commands.

`weight` follows Stata's `weight` syntax and only allows frequency weights, `fweight`.

`nopic` suppresses displaying graphs. The default is to display graphs.

`savingbounds(filename[,replace])` saves `filename.dta` with coordinates of the partially-identified set as a function of the slope magnitude of the heterogeneity distribution. Use `replace` if `filename.dta` already exists.

### Description for `bunchbounds`

The user enters the name of the response variable, the location of the kink point, the slopes before and after the kink point, and the maximum slope magnitude of the heterogeneity PDF. Before applying the command, all of these entries must be transformed into units that satisfy the DGP from Equation 1. For example, in the tax setting of Saez (2010), dollars of taxable income and the dollar value of the kink point are transformed by taking logs, and the slopes are the log of one minus the respective marginal tax.

`bunchbounds` computes the maximum and minimum values of the elasticity that are consistent with the slope restriction on the PDF specified in `m(#)`, the observed distribution of the response variable, and values of the PDF of the response variable evaluated at the left and right limits approaching the kink. These limits are computed non-parametrically using the method of Cattaneo et al. (2020) as implemented by their Stata package `lpdensity`, discussed by Cattaneo et al. (2021). Thus, the user needs to install `lpdensity` before using `bunchbounds`.

It is important to emphasize that the true value of the slope magnitude is unknowable but `bunchbounds` provides four sample values as suggestions for the user. The first two sample values are estimated using the continuous part of the distribution. Specifically, minimum and maximum slope magnitude sample values are constructed from a histogram of the dependent variable that excludes the kink point and uses a bin width that is half of the default bin width for the command `histogram`. The third sample value is the maximum slope magnitude that results in a finite upper bound on the elasticity. The fourth sample value is the minimum slope magnitude for which the elasticity bounds exist and are equal. This is the same elasticity estimate that one obtains with the trapezoidal approximation made by Saez (2010). `bunchbounds` outputs elasticity

bounds for three values of the slope: trapezoidal approximation, user-provided slope magnitude, `m(#)`, and the maximum slope magnitude that results in a finite upper bound.

## 2.2 The `bunchtobit` command

`bunchtobit` uses bunching, Tobit regressions, and covariates to point identify the elasticity of a response variable with respect to changes in the slope of the budget set. The syntax, options, and description of this command are as follows:

### Syntax for `bunchtobit`

```
bunchtobit depvar [indepvars] [if] [in] [weight] , kink(#) s0(#) s1(#) [
    binwidth(#) grid(numlist) nopic numiter(#)
    savingtobit(filename[,replace]) verbose ]
```

*depvar* must be one dependent variable (the response in logs in many applications).

kink(#) is the location of the kink point and must be a real number in the same units as the response variable.

s0(#) is a real number. In many applications, it is the log of the slope before the kink point.

s1(#) must be a real number that is strictly less s0(#). In many applications, it is the log of the slope after the kink point.

Entries for *depvar*, kink(#), s0(#), and s1(#) are required, whereas options inside the square brackets are not required.

### Options for `bunchtobit`

*indepvars* is a *varlist* of covariates. Heterogeneity is a linear function of these covariates and an unobserved error that is normally distributed conditional on these covariates.

*if* and *in* restrict the working sample, like many other Stata commands.

*weight* follows Stata's **weight** syntax and only allows frequency weights, **fweight**.

binwidth(#) is the width of the bins for the histograms. It must be a strictly positive real number. The default value is half of what is automatically produced by the command **histogram**.

grid(*numlist*) is a *numlist* of integers from 1 to 99. The values in the *numlist* correspond to percentages of the sample that define symmetric truncation windows around the kink point. The truncated Tobit model is estimated on each of these samples and also the full sample so that the number of estimates is always one more than the number

of entries in *numlist*. For example, if `grid(15 82)`, then `bunchtobit` estimates the Tobit model three times using 100, 82, and 15 percent of the data around the kink point. The default value for the *numlist* is 10(10)90, which provides 10 estimates.

`nopic` suppresses displaying graphs. The default is to display graphs.

`numiter(#)` is the maximum number of iterations allowed when maximizing the Tobit log likelihood. It must be a positive integer and the default is 500.

`savingtobit(filename[,replace])` saves *filename.dta* with Tobit estimates for each truncation window. The *filename.dta* file contains eight variables corresponding to the matrices that the code stores in `r()`. See Section 3.3 for more details. Use *replace* if *filename.dta* already exists.

`verbose` displays detailed output from the Tobit estimation including iterations of maximizing the log likelihood. Non-verbose mode is the default.

### Description for `bunchtobit`

The user enters the name of the response variable, the location of the kink point, and the slopes before and after the kink point. Before applying the command, all of these entries must be transformed into units that satisfy the DGP from Equation 1. For example, in the tax setting of Saez (2010), dollars of taxable income and the dollar value of the kink point are transformed by taking logs, and the slopes must be input as the log of one minus the marginal tax rates.

`bunchtobit` estimates multiple mid-censored Tobit regressions using specified subsamples of the data. It starts with the entire sample, then it truncates the sample to symmetric windows centered at the kink as specified by the user. The elasticity estimate is plotted as a function of the percentage of data used in each truncation window. The code also plots the histogram of the response variable along with the best-fit Tobit distribution for each truncation window.

The user has the option of entering covariates that help explain the unobserved heterogeneity. Lemma 2 by Bertanha, McCallum, and Seegert (2021) demonstrates that the distribution of the unobserved heterogeneity conditional on covariates does not need to be normal for the Tobit estimates to be consistent. Consistency requires (i) the unconditional distribution of heterogeneity is a semi-parametric mixture of normal distributions averaged over the included covariates; and (ii) the unconditional distribution of the response variable predicted by the Tobit model fits the observed distribution of the response variable well. If the user does not enter covariates, then the unconditional distribution of heterogeneity needs to be normal.

## 3 Examples for `bunchbounds` and `bunchtobit`

In this section, we use simulated data to illustrate `bunchbounds` and `bunchtobit`. These examples are motivated by the Earned Income Tax Credit that is investigated by Saez



(2010) and Bertanha, McCallum, and Seegert (2021). As such, sometimes we refer to the simulated outcome data as “earnings” and the slope of the incentive schedule as “marginal tax rates.” The units of the outcome also corresponds to log thousands of dollars.

### 3.1 Simulated data

We consider a data generating process from Equation 1 with one kink at  $k = \log(8) = 2.079$  given by

$$y_i = \begin{cases} 0.5 \log(1.3) + n_i^*, & \text{if } n_i^* < \log(8) - 0.5 \log(1.3) \\ \log(8), & \text{if } \log(8) - 0.5 \log(1.3) \leq n_i^* \leq \log(8) - 0.5 \log(0.9) \\ 0.5 \log(0.9) + n_i^*, & \text{if } n_i^* > \log(8) - 0.5 \log(0.9), \end{cases} \quad (3)$$

in which the elasticity is  $\varepsilon = 0.5$  and the slopes of the budget constraint to the left and right of the kink are  $s_0 = \log(1.3) = 0.2624$  and  $s_1 = \log(0.9) = -0.1054$  (representing tax rates of  $t_0 = -0.3$  and  $t_1 = 0.1$ ). To be concrete, the income tax rate changes from -30% to 10%, a 40 percentage point increase, and translates into a slope change in the budget set of  $-0.368 = \log(0.9) - \log(1.3)$ . The elasticity of 0.5 means that taxpayers respond to this marginal tax rate increase by decreasing their labor supply (and income) by about 18.4% ( $-0.184 = -0.368 \times 0.5$ ).

We assume that ability is a function of covariates and unobserved error given by  $n_i^* = 2 - 0.2x_{1i} + 2.5x_{2i} + 0.4x_{3i} + \nu_i$ ,  $\nu_i \sim N(0, 0.5)$ . The covariates  $x_1$ ,  $x_2$ , and  $x_3$ , are correlated binary variables with properties given in Table 1.

|       |      |           | Correlations |       |       |
|-------|------|-----------|--------------|-------|-------|
|       | Mean | Std. Dev. | $x_1$        | $x_2$ | $x_3$ |
| $x_1$ | 0.2  | 0.4       | 1            |       |       |
| $x_2$ | 0.5  | 0.5       | 0.2          | 1     |       |
| $x_3$ | 0.3  | 0.46      | 0.1          | 0.4   | 1     |

Table 1: Covariates’ proprieties

We simulate about one million weighted (100,000 unweighted) observations according to Equation 3. Frequency weights are drawn from a standard uniform distribution and we demonstrate how to employ weights throughout the **bunching** package.

In Figure 1, we graph the histogram of the one million observations in 100 bins. The simulated outcome variable is bimodal due to the covariates, which highlights that the unconditional distribution is not normally distributed. We graph the budget constraint (solid black) in  $(\log\text{-income}, \log\text{-consumption})$  space. That budget set has a kink, that is, a change in slope from 1.3 to 0.9 at the value of 2.079 (solid red) for log-income. The histogram in the same figure shows that individuals bunch exactly at the kink point (sand bar).

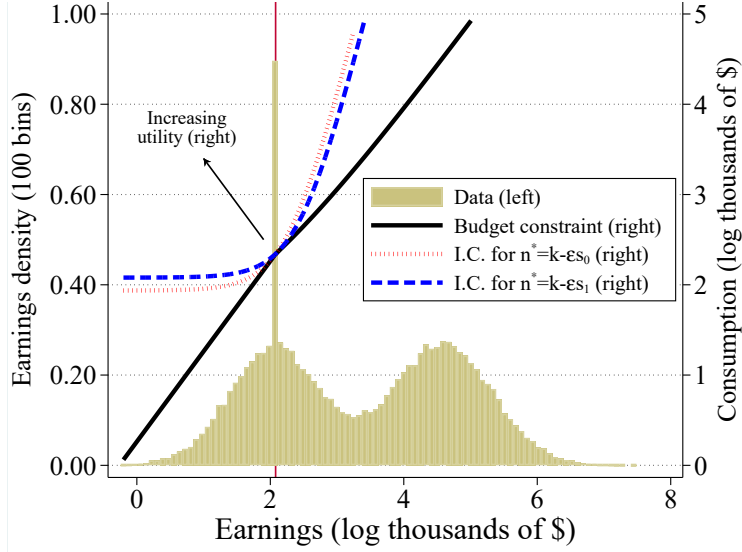


Figure 1: Histogram of simulated data

Bertanha, McCallum, and Seegert (2021) provide a complete description for how utility maximization with heterogeneous preferences and income tax brackets results in Figure 1, and we provide an overview here. The heterogeneity of agents' preferences is captured by  $n^*$ , and each value of  $n^*$  corresponds to a different indifference curve (I.C.). We graph two specific I.C.s, which correspond to the lower (dotted red) and upper (dashed blue) numerical thresholds in Equation 3, whose theoretical counterparts are  $\underline{n}(k, \varepsilon, s_0) = k - \varepsilon s_0$  and  $\bar{n}(k, \varepsilon, s_1) = k - \varepsilon s_1$  in Equation 1. Many I.C.s that are not graphed touch the budget set at the kink. In fact, the mass point at the kink corresponds to all agents whose preference heterogeneity,  $n^*$ , lies in the bunching interval, that is,  $n^* \in [\log(8) - 0.5 \log(1.3), \log(8) - 0.5 \log(0.9)]$ .

The simulated data also exhibit bunching exactly at the kink point. In many empirical applications, however, the bunching mass is dispersed in a small interval near, instead of exactly at, the kink. We provide a solution to this issue in Section 4.

### 3.2 Estimating elasticity bounds

We begin by estimating the elasticity bounds using the location of the kink,  $\log(8) = 2.0794$ ,  $k(2.0794)$ , tax rates on either side of the kink,  $s_0 = \log(1.3) = 0.2624$  and  $s_1 = \log(0.9) = -0.1054$ , and a choice of the maximum slope,  $m(2)$ . The `bunching` package and simulated data are available from the Boston College Statistical Software Components (SSC) archive provided by Research Papers in Economics (RePEc).

```
. ssc install bunching
checking bunching consistency and verifying not already installed...
```

```

installing into c:\ado\plus\...
installation complete.

. webuse set "http://fmwww.bc.edu/repec/bocode/b/"
(prefix now "http://fmwww.bc.edu/repec/bocode/b")
. webuse bunching.dta

. bunchbounds y [fw=w], k(2.0794) s0(0.2624) s1(-0.1054) m(2)

Your choice of M:
2.0000

Sample values of slope magnitude M
minimum value M in the data (continuous part of the PDF):
0.0000
maximum value M in the data (continuous part of the PDF):
0.3879
maximum choice of M for finite upper bound:
1.5930
minimum choice of M for existence of bounds:
0.0090

Elasticity Estimates
Point id., trapezoidal approx.:
0.4894
Partial id., M = 2.0000 :
[0.3913 , +Inf]
Partial id., M = 1.59 :
[0.4055 , 0.9374]

```

The `bunchbounds` command estimates the bounds for the elasticity using different slope values. First, the output shows that we entered a maximum slope of 2 and the bounds for this slope are  $[0.3912, \infty]$ . Second, the command also estimates the bounds using the maximum slope for a finite upper bound, when the maximum slope given is greater than that value. In this case, the maximum slope for a finite upper bound is 1.5933, resulting in the bounds  $[0.4055, 0.9353]$ . In both cases, the true elasticity estimate of 0.5 is within these bounds. The output also gives the estimated minimum and maximum slopes of the continuous portion of the probability density function of the data. These slopes are 0 and 0.3879. The point-identified elasticity using the trapezoidal approximation (which is the Saez (2010) estimator) of 0.4893 is also provided.

The non-parametric bounds are also graphed by `bunchbounds` for different maximum slope magnitudes of the unobserved heterogeneity PDF. These different slope magnitudes are plotted on the horizontal axis and the corresponding bounds are plotted on the vertical axis. For this example, these are given in Figure 2a. This figure shows how the upper bound, depicted as a dashed line, increases and the lower bound, depicted as a solid line, decreases as the maximum slope increases. The vertical lines in Figure 2a at 0.01 and 1.59 denote the minimum slope for the existence of the bounds and the maximum slope for a finite upper bound, respectively. The point identified elasticity using the trapezoidal approximation occurs where the bounds come together —the dash-dot horizontal red line in Figure 2a.

The `bunchbounds` command can also be combined with conditional statements that

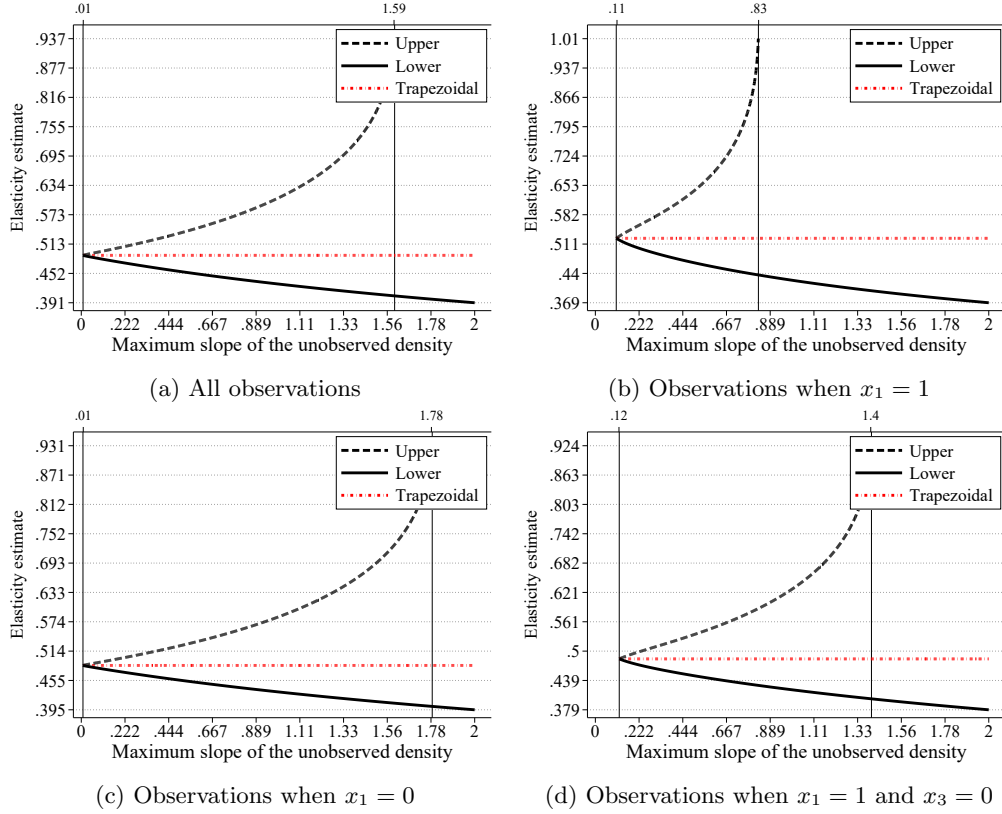


Figure 2: Estimating elasticity bounds

restricts to subsamples of the data based on values of different covariates but cannot otherwise be conditional on covariates. For example,

```
bunchbounds y if x1==1 & x3==0 [fw=w], k(2.0794) s0(0.2624) s1(-0.1054) m(2)
```

estimates the bounds when  $x_1 = 1$  and  $x_3 = 0$ . Restricting to subsamples when  $x_1 = 1$  or  $x_1 = 0$  have similar syntaxes. The output from these commands (not shown) is similar to the output without conditioning and the bound estimates for each subsample are graphed in Figures 2b, 2c, and 2d. The bounds shift only slightly for each subsample because the true elasticity is 0.5 for all subsamples and because the number of weighted observations is large.

### 3.3 Semi-parametric point estimates of the elasticity

We estimate the elasticity using a truncated Tobit model that allows for covariates. Truncation and covariates provide robust estimation that relies on semi-parametric as-

sumptions and does not require the unobserved heterogeneity PDF to be normally distributed (Bertanha, McCallum, and Seegert 2021). We demonstrate the robustness of this method by comparing estimates of the correctly specified model with estimates from a misspecified model that still recover the true elasticity.

### Correctly specified Tobit model

We begin by estimating the correctly specified model using `bunchtobit`.

```
. bunchtobit y x1 x2 x3 [fw=w], k(2.0794) s0(0.2624) s1(-0.1054) binwidth(0.084)
```

```
Obtaining initial values for ML optimization.
Truncation window number 1 out of 10, 100% of data.
Truncation window number 2 out of 10, 90% of data.
Truncation window number 3 out of 10, 80% of data.
Truncation window number 4 out of 10, 70% of data.
Truncation window number 5 out of 10, 60% of data.
Truncation window number 6 out of 10, 50% of data.
Truncation window number 7 out of 10, 40% of data.
Truncation window number 8 out of 10, 30% of data.
Truncation window number 9 out of 10, 20% of data.
Truncation window number 10 out of 10, 10% of data.
```

```
bunchtobit_out[10,5]
```

|    | data % | elasticity | std err   | # coll | cov | flag |
|----|--------|------------|-----------|--------|-----|------|
| 1  | 100    | .50938668  | .00218386 |        | 0   | 0    |
| 2  | 90     | .50756197  | .00224619 |        | 0   | 0    |
| 3  | 80     | .50898083  | .00227815 |        | 0   | 0    |
| 4  | 70     | .50808053  | .00229178 |        | 0   | 0    |
| 5  | 60     | .50848689  | .00231719 |        | 0   | 0    |
| 6  | 50     | .50660888  | .00236933 |        | 0   | 0    |
| 7  | 40     | .50975777  | .00251876 |        | 0   | 0    |
| 8  | 30     | .50959025  | .00273068 |        | 0   | 0    |
| 9  | 20     | .50463572  | .00317585 |        | 0   | 0    |
| 10 | 10     | .47913201  | .00419053 |        | 0   | 0    |

The command estimates the elasticity for ten different subsamples by default. The first uses all the data, the second uses 90% of the data around the kink, the third uses 80% around the kink, and so on. Estimation proceeds in 10 percentage point intervals declining down to the last subsample that uses only 10% of the data. Each subsample is truncated symmetrically, centered around the kink, and includes the observations at the kink. For the data simulated by Equation 3 and using the 90% truncated subsample as an example, about 42.5% of the data are from below the kink, about 42.5% of the data are from above the kink, and about 5% of the data are from the kink. The fraction of data at the kink does not change with this type of truncation. For example, the 10% subsample uses about 2.5% of the data above and below the kink and about 5% from the kink.

Because the model is correctly specified, the estimates reported in the `elasticity` column are always very close to the true value of 0.5 for any truncated subsample. Standard errors in column `std err` are small because the simulated data includes one million weighted observations. The standard errors increase as the percent of data used

decreases because we use fewer observations. The table also reports the number of covariates that were omitted because they were collinear in column `# coll cov` and when optimizing the log likelihood did not converge to a maximum in column `flag`.

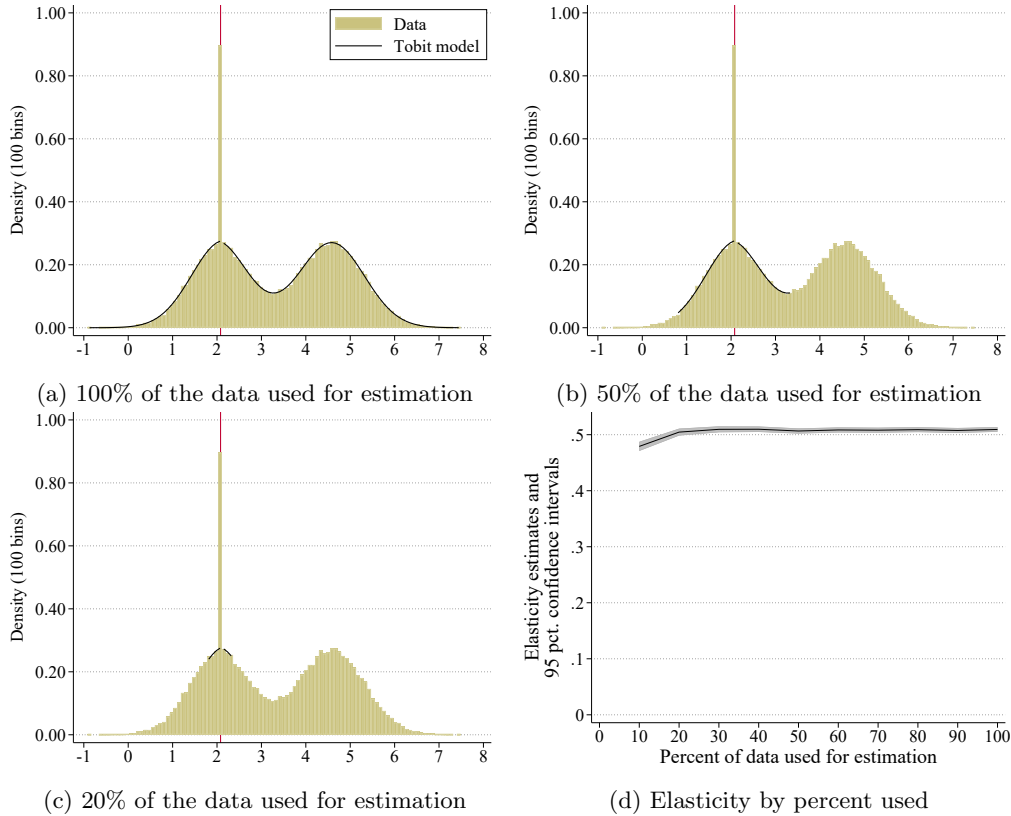


Figure 3: Correctly specified truncated Tobit estimates

Along with this numeric output, `bunchtobit` also produces a best-fit graph for each subsample and a graph of the elasticity estimate for all subsamples. Figures 3a, 3b, and 3c display these best-fit graphs for the 100%, 50%, and 20% truncation subsamples, respectively. Each of these panels presents a histogram of  $y_i$  (sand colored bars) and the estimate of the correctly specified and truncated Tobit model implied outcome variable (black line). The model is correctly specified and so it fits the data well for all truncated subsamples. Figure 3d plots the estimate (black line) and 95% confidence interval (gray shading) for each truncated subsample corresponding to the `elasticity` column. The elasticity is the main parameter of interest but the covariate coefficients for the smallest value in the `numlist` provided in `grid(numlist)` can be obtained by using the `estimates replay` command. For example, truncating to 77% of the data for the correctly specified model and then using `estimates replay` provides the following output:

```
. bunchtobit y x1 x2 x3 [fw=w], k(2.0794) s0(0.2624) s1(-0.1054) binwidth(0.084)
> grid(77)
```

```
Obtaining initial values for ML optimization.
Truncation window number 1 out of 2, 100% of data.
Truncation window number 2 out of 2, 77% of data.
```

```
bunchtobit_out[2,5]
      data % elasticity      std err # coll cov      flag
1      100   .50938668   .00218386      0      0
2       77   .50849786   .00228162      0      0
```

```
. estimates replay
```

```
-----
active results
-----
```

```
Log pseudolikelihood = -.96353496                Number of obs = 770,197
```

```
( 1) [eq_l]x1 - [eq_r]x1 = 0
( 2) [eq_l]x2 - [eq_r]x2 = 0
( 3) [eq_l]x3 - [eq_r]x3 = 0
```

|         |        | Coefficient | Robust<br>std. err. | z      | P> z  | [95% conf. interval] |           |
|---------|--------|-------------|---------------------|--------|-------|----------------------|-----------|
| eq_l    | x1     | -.2876614   | .0035942            | -80.03 | 0.000 | -.2947059            | -.2806168 |
|         | x2     | 3.541998    | .0038313            | 924.49 | 0.000 | 3.534488             | 3.549507  |
|         | x3     | .5509258    | .0036639            | 150.37 | 0.000 | .5437448             | .5581069  |
|         | _cons  | 3.022123    | .0033913            | 891.13 | 0.000 | 3.015476             | 3.02877   |
| eq_r    | x1     | -.2876614   | .0035942            | -80.03 | 0.000 | -.2947059            | -.2806168 |
|         | x2     | 3.541998    | .0038313            | 924.49 | 0.000 | 3.534488             | 3.549507  |
|         | x3     | .5509258    | .0036639            | 150.37 | 0.000 | .5437448             | .5581069  |
|         | _cons  | 2.757436    | .0035784            | 770.58 | 0.000 | 2.750422             | 2.764449  |
| lngamma | _cons  | .347303     | .001056             | 328.87 | 0.000 | .3452331             | .3493728  |
|         | sigma  | .7065912    | .0014946            |        |       | .7051302             | .7080553  |
|         | cons_l | 2.135406    | .0030205            |        |       | 2.129486             | 2.141326  |
|         | cons_r | 1.94838     | .0033687            |        |       | 1.941778             | 1.954983  |
|         | eps    | .5084979    | .0022816            |        |       | .504026              | .5129697  |

Olsen (1978) introduces a reparameterization that is discussed in (Hayashi 2000, Ch. 8.3) that ensures the log likelihood of a classical Tobit model is globally concave. That reparameterization divides each coefficient of the covariates by the standard deviation of the errors and we use the same reparameterization in our log likelihood. The results output by `estimates replay` report these reparameterized coefficients instead of the original coefficients. The reparameterization can be reversed by multiplying the reparameterized coefficients by the standard deviation. For example, the estimate of the coefficient on  $x_2$  from Equation 3 can be recovered as  $3.54 \times .71 = 2.51$ .

The elasticity reported in column `elasticity` for the 77% subsample is from the estimate `eps` in the `active results` table shown by `estimates replay`. The first equation, `eq_l`, coefficient estimates on  $x_1$ ,  $x_2$ , and  $x_3$  are from the left-hand side of the kink and are the same as the estimates from the second equation, `eq_r`, on the right of the kink. These coefficients are constrained to be the same on the left and right sides of the kink as reflected by the three constraints ( 1), ( 2), and ( 3), at the top of the table and consistent with Equation 3. Because the model is correctly specified, the covariate coefficient estimates are consistent and the estimates shown by `estimates replay` are close to the (reparameterized) truth for each coefficient.

### Incorrectly specified Tobit model

The correctly specified Tobit model from the previous section satisfies the assumption that  $\nu_i$  is normal and therefore always fits the observed distribution of  $y_i$ . A misspecified model that does not have normally distributed errors will not always fit the distribution of  $y_i$  well. However, Bertanha, McCallum, and Seegert (2021) prove that if the Tobit model's best-fit distribution matches the observed distribution of  $y_i$ , then the Tobit model estimates the elasticity consistently whether or not the distribution of  $\nu_i$  is normal. This section demonstrates this robustness property using a misspecified model that does not have normal errors. Specifically, we omit the covariate  $x_2$  and estimate the following model.

```
. bunchtobit y x1 x3 [fw=w], k(2.0794) s0(0.2624) s1(-0.1054) binwidth(0.084)
```

```
Obtaining initial values for ML optimization.
Truncation window number 1 out of 10, 100% of data.
Truncation window number 2 out of 10, 90% of data.
Truncation window number 3 out of 10, 80% of data.
Truncation window number 4 out of 10, 70% of data.
Truncation window number 5 out of 10, 60% of data.
Truncation window number 6 out of 10, 50% of data.
Truncation window number 7 out of 10, 40% of data.
Truncation window number 8 out of 10, 30% of data.
Truncation window number 9 out of 10, 20% of data.
Truncation window number 10 out of 10, 10% of data.
```

| bunchtobit_out[10,5] |        |            |           |            |      |
|----------------------|--------|------------|-----------|------------|------|
|                      | data % | elasticity | std err   | # coll cov | flag |
| 1                    | 100    | .6426979   | .00284279 | 0          | 0    |
| 2                    | 90     | .7643775   | .00347177 | 0          | 0    |
| 3                    | 80     | .74113379  | .00338469 | 0          | 0    |
| 4                    | 70     | .68969718  | .00316174 | 0          | 0    |
| 5                    | 60     | .61191988  | .00282291 | 0          | 0    |
| 6                    | 50     | .52858458  | .00248579 | 0          | 0    |
| 7                    | 40     | .51255963  | .00253649 | 0          | 0    |
| 8                    | 30     | .51034751  | .00273715 | 0          | 0    |
| 9                    | 20     | .50446083  | .0031749  | 0          | 0    |
| 10                   | 10     | .48045869  | .00528865 | 0          | 0    |

The misspecified model returns an elasticity estimate of 0.642 using 100% of the data. This is a substantially biased estimate of the true elasticity of 0.5 and Figure 4a shows that the misspecified model does not fit well.



We can truncate the sample to use data only local to the kink, however, to attenuate the effect of omitting  $x_2$ . In Bertanha et al. (2021, Lemma 2), we show that if the Tobit distribution of the fitted outcome (the black solid lines in Figures 4a to 4c) matches the true distribution of the outcome variable (the sand bars in those figures), and the unconditional distribution of  $n^*$  is a mixture of normals, then the elasticity estimated by the Tobit is consistent for the true elasticity, regardless of whether the conditional unobserved distribution,  $F_{n^*|X}$ , is normal.

Moreover, the smaller the truncation window, the easier it is to fit the unconditional distribution of the outcome variable with a Tobit, and the stronger is our conviction that the estimate of the elasticity is consistent.

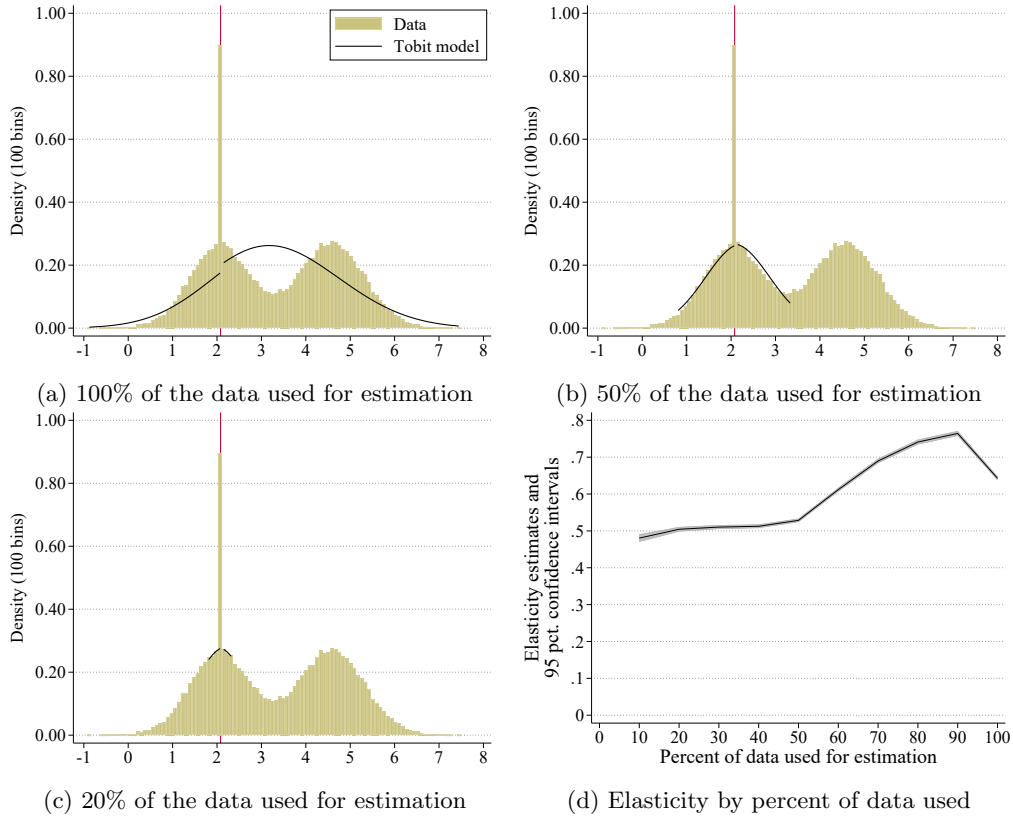


Figure 4: Incorrectly specified truncated Tobit estimates

Figure 4 demonstrates that using smaller truncation windows around the kink improves the estimated distribution fit. Figure 4b uses 50% of the data and fits much better than the estimate that uses all of the data in Figure 4a. Figure 4c uses 20% of the data local to the kink and fits even better than the 50% subsample. Figure 4d shows that for all subsamples that use 50% of the data or less, we recover an estimate

that is close to the true elasticity of 0.5. The largest truncation region for which the estimated distribution fits the observed distribution is context specific. For the example given in Figure 4, using 50% of the data around the kink is the largest subsample of data that provides a good fit to the outcome distribution. But for other datasets, the largest truncation region that fits the outcome distribution well could use any fraction of the data, and could be very small indeed.

## 4 Friction errors

Many datasets have friction errors which are defined as when the bunching mass is dispersed in a small interval near, instead of exactly at, the kink. Friction errors can be caused by measurement error, optimizing frictions (Chetty et al. 2011), or other distortions. When friction errors are present, they must first be filtered out before a bunching estimation method can be applied.

The procedure implemented by `bunchfilter` is a practical way of filtering out friction errors. It works by fitting a polynomial to the empirical CDF of the response variable with friction errors,  $y_{fric_i}$ . It excludes observations in a specified interval around the kink during estimation and allows the intercepts to differ to the left and right of that interval. The estimated CDF is then extrapolated into the excluded interval, which constitutes an estimate of the CDF of the response variable without friction errors,  $y_i$ . The inverse of the extrapolated CDF evaluated at each observation produces the filtered variable and the difference between the intercepts at the kink provides the estimate of the bunching mass.

This filtering method produces consistent estimates of the distribution of the response variable without frictions under three conditions. First, the friction error,  $e_i$ , must be independent and identically distributed (*iid*) with known and bounded support. We emphasize that it is not necessary for the friction error to be mean zero, or for the distribution of friction error,  $f(e_i)$ , to be symmetric or parametric. Second, friction errors must only affect bunching individuals. Third, the CDF of  $y_i$  without friction error must equal a polynomial in a known neighborhood of the kink that is bigger than the support of the friction error.

### 4.1 The `bunchfilter` command

`bunchfilter` removes friction errors from data generated by a mixed continuous-discrete distribution with one mass point plus a continuously distributed friction error. The distribution of the data with friction error is continuous and does not have a mass point. This type of data is common in economic bunching applications. For example, the distribution of taxable income usually has a hump around the kink where the marginal tax rate changes, instead of a mass point at the kink. The syntax, options, and description of this command are as follows:

**Syntax for `bunchfilter`**

```
bunchfilter depvar [if] [in] [weight] , kink(#) deltam(#) deltap(#)
      generate(varname) [ binwidth(#) nopic pctobs(#) polorder(#) ]
```

*depvar* must be one dependent variable (the response in logs in many applications).

**kink**(#) is the location of the kink point and must be a real number in the same units as the response variable.

**deltam**(#) is the distance between the kink point and the lower bound of the support of the friction error to be filtered. It must be a real number in the same units as the response variable.

**deltap**(#) is the distance between the kink point and the upper bound of the support of the friction error to be filtered. It must be a real number in the same units as the response variable.

**generate**(*varname*) generates the filtered variable with a user-specified name of *varname*.

Entries for *depvar*, **kink**(#), **deltam**(#), **deltap**(#), and **generate**(*varname*) are required, whereas options inside the square brackets are not required.

**Options for `bunchfilter`**

*if* and *in* restrict the working sample, like many other Stata commands.

*weight* follows Stata's **weight** syntax and only allows frequency weights, **fweight**.

**binwidth**(#) is the width of the bins for the histograms. It must be a strictly positive real number. The default value is half of what is automatically produced by the command **histogram**.

**nopic** suppresses displaying graphs. The default is to display graphs.

**pctobs**(#) for better fit, the polynomial regression uses observations in a symmetric window around the kink point that contains **pctobs**(#) percent of the sample. It must be a positive integer between 1 and 99 and the default is 40.

**polorder**(#) order of polynomial for the filtering regression. It must be a positive integer between 1 and 7 and the default is 7.

**Description for `bunchfilter`**

The user enters the variable to be filtered (for example, the log of income), the location of the kink, and size of a region around the mass point that contains the hump (in other words, **kink** - **deltam**, **kink** + **deltap**). **bunchfilter** fits a polynomial regression to the empirical CDF of the variable observed with error. This regression excludes

points in the hump window and has a dummy for observations on the left or right of the kink. The fitted regression is used to predict values of the empirical CDF in the hump window with a jump discontinuity at the mass point. The filtered variable is then recovered from the inverse of the predicted CDF evaluated at the empirical CDF value for each observation in the sample.

This procedure works well for cases where the friction error has bounded support and only affects observations that would be at the kink in the absence of error. A proper deconvolution theory still needs to be developed for a filtering procedure with general validity.

## 4.2 Example for `bunchfilter`

We show how to remove the friction errors as a precursor to estimating the relevant elasticity in this example. We simulate the outcome variable with friction errors as

$$y_{fric_i} = y_i + e_i \mathbb{I}(y_i = \log(8)), \quad (4)$$

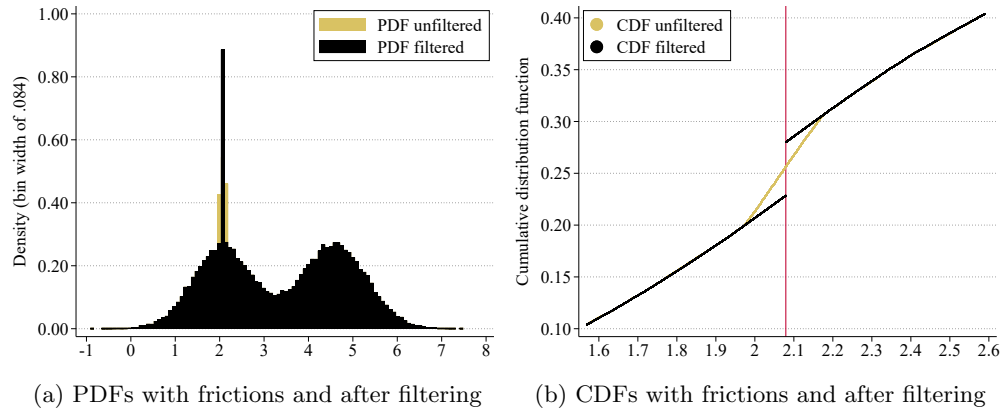
in which  $y_i$  is from Equation 3,  $e_i$  are *iid* truncated normal from  $f(e_i) = \phi(e_i) / [\Phi(\log(1.1)) - \Phi(\log(0.9))]$ , the standard normal PDF is  $\phi(\cdot)$ , and  $\Phi(\cdot)$  is the standard normal CDF. The errors have known and bounded support  $[\log(0.9), \log(1.1)]$ , which ensures frictions never add to or subtract from  $y_i$  by more than log 10 percent. The three conditions needed for `bunchfilter` to consistently estimate  $y_i$  discussed in Section 4 are satisfied by Equation 4.

We generate the filtered variable, `yfiltered`, and Figure 5 by applying `bunchfilter` to the outcome variable with friction errors using the following command (output not shown)

```
. bunchfilter yfric [fw=w], kink(2.0794) deltam(0.12) deltap(0.12) generate(yfiltered)
> binwidth(0.084) pctobs(30)
```

We exclude log 12 percent below the kink, `deltam(0.12)`, and log 12 percent, `deltap(0.12)`, above the kink because we know this excluded region will capture the support of the friction errors because the example frictions in Equation 4 never add to or subtracts from  $y_i$  by more than log 10 percent.

Without the friction errors, 5.17% of the responses bunch at the kink in the simulated data from Equation 3. Including friction errors lowers this fraction to zero because no observation are exactly at the kink in Equation 4. After removing the frictions with `bunchfilter`, the filtered data has 5.15% of the responses at the kink. The histogram of  $y_{fric_i}$  is shown in Figure 5a. The unfiltered data (sand colored bars) exhibits diffuse bunching around the kink point. The filtered data is saved in the variable `yfiltered` by invoking the option `generate(yfiltered)`. The histogram for the filtered data is depicted in the (black bars) with evident reassignment of original dispersed observations around the kink to the kink point exactly. This reassignment can also be seen in the contrast between the filtered and unfiltered CDFs in Figure 5b. Both of these figures are produced by the `bunchfilter` command.

Figure 5: Effect of `bunchfilter` on data with friction errors

## 5 Automatic estimation

Despite friction errors and model misspecification, `bunching` provides multiple estimates of the true elasticity by implementing `bunchbounds`, `bunchtobit`, and `bunchfilter` automatically. The user can provide outcome data with friction errors and a misspecified model and `bunching` can still recover estimates that are close to the true elasticity.

### 5.1 The `bunching` command

The Stata command `bunching` is a wrapper function for three other commands: `bunchbounds`, `bunchtobit`, and `bunchfilter`.

#### Syntax

```
bunching depvar [indepvars] [if] [in] [weight] , kink(#) s0(#) s1(#) m(#)
[ nopic savingbounds(filename[,replace]) binwidth(#) grid(numlist)
numiter(#) savingtobit(filename[,replace]) verbose deltam(#)
deltap(#) generate(varname) pctobs(#) polorder(#) ]
```

The syntax and options for `bunching` are inherited from the three commands for which it is a wrapper function, and so we do not repeat them here. Entries for the first four options, `kink(#)`, `s0(#)`, `s1(#)`, and `m(#)` are required whereas options inside the square brackets are not required. `bunching` always implements `bunchbounds` and `bunchtobit`. In contrast, `bunchfilter` is only called by `bunching` if all the required entries for `bunchfilter`, namely, `deltam(#)`, `deltap(#)`, and `generate(varname)`, are specified.

## 5.2 Example using bunching

The following example uses `bunching` with the outcome data from Equation 4 but omits weights and the covariate  $x_2$  in order to demonstrate the robustness of this package.

```
. bunching yfric x1 x3, kink(2.0794) s0(0.2624) s1(-0.1054) m(2) binwidth(0.084)
> deltam(0.12) deltap(0.12) gen(ybunching) pctobs(30)
*****
Bunching - Filter
*****
[ 10% 20% 30% 40% 50% 60% 70% 80% 90% 100% ]
*****
Bunching - Bounds
*****
Your choice of M:
2.0000

Sample values of slope magnitude M
minimum value M in the data (continuous part of the PDF):
0.0000
maximum value M in the data (continuous part of the PDF):
0.3334
maximum choice of M for finite upper bound:
1.5530
minimum choice of M for existence of bounds:
0.0792

Elasticity Estimates
Point id., trapezoidal approx.:
0.4930
Partial id., M = 2.0000 :
[0.3926 , +Inf]
Partial id., M = 1.55 :
[0.4087 , 0.9480]
*****
Bunching - Tobit
*****
Obtaining initial values for ML optimization.
Truncation window number 1 out of 10, 100% of data.
Truncation window number 2 out of 10, 90% of data.
Truncation window number 3 out of 10, 80% of data.
Truncation window number 4 out of 10, 70% of data.
Truncation window number 5 out of 10, 60% of data.
Truncation window number 6 out of 10, 50% of data.
Truncation window number 7 out of 10, 40% of data.
Truncation window number 8 out of 10, 30% of data.
Truncation window number 9 out of 10, 20% of data.
Truncation window number 10 out of 10, 10% of data.

bunchtobit_out[10,5]
```

|   | data % | elasticity | std err   | # coll | cov | flag |
|---|--------|------------|-----------|--------|-----|------|
| 1 | 100    | .63579158  | .00894356 |        | 0   | 0    |
| 2 | 90     | .75808395  | .01094832 |        | 0   | 0    |
| 3 | 80     | .73437667  | .01066292 |        | 0   | 0    |
| 4 | 70     | .68368544  | .00996446 |        | 0   | 0    |
| 5 | 60     | .60786248  | .00891428 |        | 0   | 0    |
| 6 | 50     | .52680042  | .00787451 |        | 0   | 0    |
| 7 | 40     | .50716668  | .00796645 |        | 0   | 0    |
| 8 | 30     | .5045792   | .00858102 |        | 0   | 0    |

|    |    |           |           |   |   |
|----|----|-----------|-----------|---|---|
| 9  | 20 | .50167298 | .01001394 | 0 | 0 |
| 10 | 10 | .50828368 | .0303952  | 0 | 0 |

**bunching** first filters the data using **bunchfilter**. It then implements **bunchbounds** on the filtered outcome using the full sample and maximum slope magnitude as specified. Finally, it uses **bunchtobit** on the filtered outcome with the covariates specified,  $x_1$  and  $x_3$ , for each of the 10 default truncated subsamples.

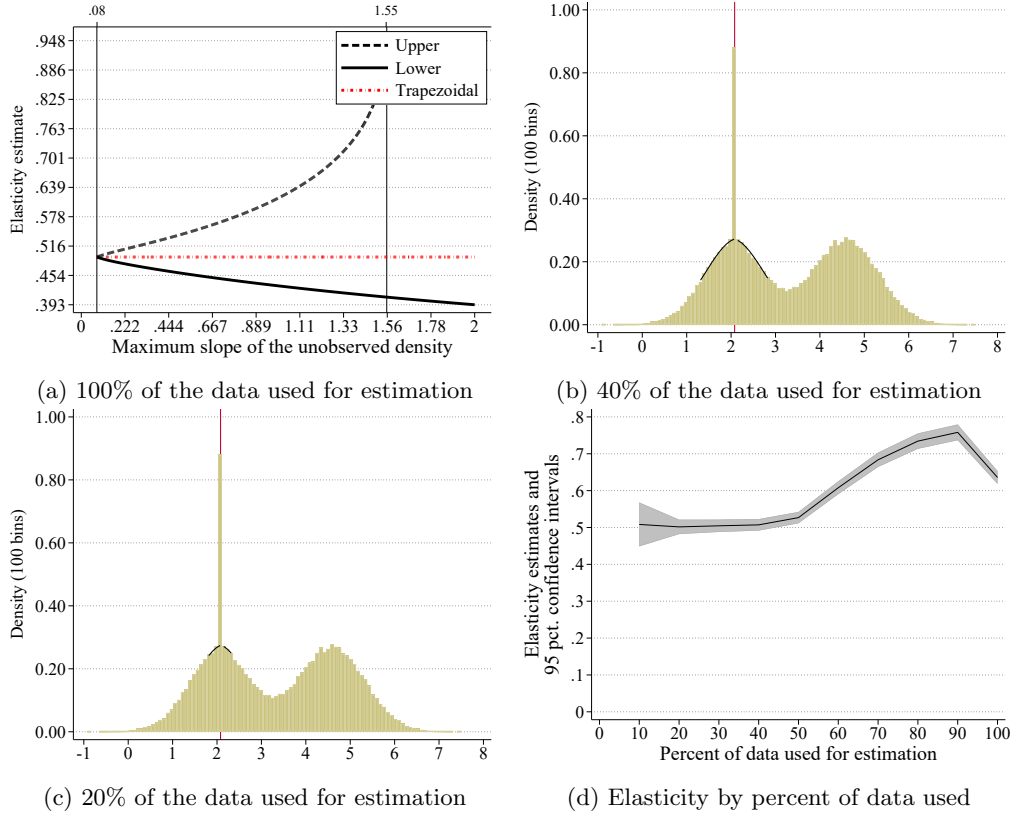


Figure 6: Elasticity estimates with friction errors and model misspecification

Along with numeric output, **bunching** produces the graphs produced by each of **bunchfilter**, **bunchbounds**, and **bunchtobit** commands. Selections from these graphs are shown in Figure 6. The output from **bunching** shows that after we filter the data, the bounds contain the true value of 0.5 (Figure 6a). Likewise, estimates from the Tobit model in the numeric output show that using a 40% subsample or less recovers the true elasticity of 0.5 despite friction errors and model misspecification. Truncating to 40% of the data provides a good fit as shown in Figure 6b, and Figure 6c shows that truncating to 20% also provides a good fit. Figure 6d shows that estimates with confidence intervals include the true elasticity of 0.5 for subsamples with 40% of the data and less.

## 6 Concluding remarks

Our new `bunching` package provides a series of estimation methods that enable researchers to examine the sensitivity of their elasticity estimates to different identification assumptions. The new techniques include bounds based on non-parametric assumptions and a mid-censored regression based on semi-parametric assumptions and covariates. The non-parametric assumptions are the most flexible of the two approaches and nest the trapezoidal approximation assumption, which was the method utilized by the original bunching estimator. These methods can be applied in cases with multiple kinks at each kink separately if the constraint is continuous preceding the kink under study. `bunching` has broad applicability because budget constraints with kinks occur in a variety of fields within economics and other social sciences.

## 7 Acknowledgements

The views expressed in this paper represent the views of the authors and do not indicate concurrence either by the Board of Governors of the Federal Reserve System or other members of the Federal Reserve System. We gratefully acknowledge the contributions of Andrey Ampilogov. Michael A. Navarrete provided excellent research assistance. Bertanha thanks the Kenneth C. Griffin Department of Economics at the University of Chicago for the financial support and hospitality received during his sabbatical year, when part of this work was conducted.

## 8 References

- Allen, E. J., P. M. Dechow, D. G. Pope, and G. Wu. 2017. Reference-Dependent Preferences: Evidence from Marathon Runners. *Management Science* 63(6): 1657–1672.
- Bertanha, M., A. H. McCallum, and N. Seegert. 2021. Better Bunching, Nicer Notching. Finance and Economics Discussion Series 2021-002, Board of Governors of the Federal Reserve System (U.S.). <https://doi.org/10.17016/FEDS.2021.002>.
- Bitler, M., J. Cook, and J. Rothbaum. 2021. Working for Your Bread: The Labor Supply Effects of SNAP. *American Economic Association Papers and Proceedings* 111: 1–5.
- Caetano, C. 2015. A Test of Exogeneity Without Instrumental Variables in Models With Bunching. *Econometrica* 83(4): 1581–1600. <https://onlinelibrary.wiley.com/doi/abs/10.3982/ECTA11231>.
- Caetano, C., G. Caetano, and E. R. Nielsen. 2020a. Should Children Do More Enrichment Activities? Leveraging Bunching to Correct for Endogeneity. Technical Report 2020-036, Board of Governors of the Federal. <https://doi.org/10.17016/FEDS.2020.036>.



- . 2020b. Correcting for Endogeneity in Models with Bunching. Finance and Economics Discussion Series 2020-080, Board of Governors of the Federal Reserve System (U.S.). <https://ideas.repec.org/p/fip/fedgfe/2020-80.html>.
- Caetano, G., J. Kinsler, and H. Teng. 2019. Towards causal estimates of children's time allocation on skill development. *Journal of Applied Econometrics* 34(4): 588–605. <https://onlinelibrary.wiley.com/doi/abs/10.1002/jae.2700>.
- Caetano, G., and V. Maheshri. 2018. Identifying dynamic spillovers of crime with a causal approach to model selection. *Quantitative Economics* 9(1): 343–394.
- Cattaneo, M. D., M. Jansson, and X. Ma. 2020. Simple Local Polynomial Density Estimators. *Journal of the American Statistical Association* 115(531): 1449–1455.
- . 2021. lpdensity: Local Polynomial Density Estimation and Inference. *Journal of Statistical Software, Articles* Forthcoming.
- Cengiz, D., A. Dube, A. Lindner, and B. Zipperer. 2019. The Effect of Minimum Wages on Low-wage Jobs. *Quarterly Journal of Economics* 134(3): 1405–1454.
- Chetty, R., J. N. Friedman, T. Olsen, and L. Pistaferri. 2011. Adjustment Costs, Firm Responses, and Micro vs. Macro Labor Supply Elasticities: Evidence from Danish Tax Records. *Quarterly Journal of Economics* 126(2): 749–804.
- Chetty, R., J. N. Friedman, and E. Saez. 2013. Using Differences in Knowledge across Neighborhoods to Uncover the Impacts of the EITC on Earnings. *American Economic Review* 103(7): 2683–2721. <http://www.aeaweb.org/articles?id=10.1257/aer.103.7.2683>.
- Dee, T. S., W. Dobbie, B. A. Jacob, and J. Rockoff. 2019. The Causes and Consequences of Test Score Manipulation: Evidence from the New York Regents Examinations. *American Economic Journal: Applied Economics* 11(3): 382–423. <http://www.aeaweb.org/articles?id=10.1257/app.20170520>.
- Einav, L., A. Finkelstein, and P. Schrimpf. 2017. Bunching at the Kink: Implications for Spending Responses to Health Insurance Contracts. *Journal of Public Economics* 146: 27–40.
- Garicano, L., C. Lelarge, and J. Van Reenan. 2016. Firm Size Distortions and the Productivity Distribution: Evidence from France. *American Economic Review* 106(11): 3439–3479. <https://ideas.repec.org/a/aea/aecrev/v106y2016i11p3439-79.html>.
- Ghanem, D., S. Shen, and J. Zhang. 2019. A Censored Maximum Likelihood Approach to Quantifying Manipulation in China's Air Pollution Data. Working paper, University of California - Davis.
- Grossman, D., and U. Khalil. 2020. Neighborhood networks and program participation. *Journal of Health Economics* 70: 102257. <http://www.sciencedirect.com/science/article/pii/S0167629618306830>.

- Hayashi, F. 2000. *Econometrics*. Princeton University Press.
- Ito, K. 2014. Do Consumers Respond to Marginal or Average Price? Evidence from Nonlinear Electricity Pricing. *American Economic Review* 104(2): 537–563.
- Ito, K., and J. M. Sallee. 2018. The Economics of Attribute-Based Regulation: Theory and Evidence from Fuel Economy Standards. *Review of Economics and Statistics* 100(2): 319–336.
- Jales, H. 2018. Estimating the effects of the minimum wage in a developing country: A density discontinuity design approach. *Journal of Applied Econometrics* 33(1): 29–51. <https://onlinelibrary.wiley.com/doi/abs/10.1002/jae.2586>.
- Jales, H., and Z. Yu. 2017. Identification and Estimation Using a Density Discontinuity Approach. In *Regression Discontinuity Designs: Theory and Applications*, ed. M. D. Cattaneo and J. C. Escanciano, 29–72. Vol. 38. Emerald Publishing Limited. <https://www.emerald.com/insight/content/doi/10.1108/S0731-905320170000038003/full/html>.
- Khalil, U., and N. Yildiz. 2020. A Test of the Selection-on-Observables Assumption Using a Discontinuously Distributed Covariate. Technical report, Monash University. <https://www.dropbox.com/s/o9bgdua6kcut7a8/selnonobsvbls200715.pdf?dl=0>.
- Kleven, H. J. 2016. Bunching. *Annual Review of Economics* 8: 435–464.
- Kleven, H. J., and M. Waseem. 2013. Using Notches to Uncover Optimization Frictions and Structural Elasticities: Theory and Evidence from Pakistan. *Quarterly Journal of Economics* 128(2): 669–723.
- Kopczuk, W., and D. Munroe. 2015. Mansion Tax: The Effect of Transfer Taxes on the Residential Real Estate Market. *American Economic Journal: Economic Policy* 7(2): 214–57.
- Olsen, R. J. 1978. Note on the Uniqueness of the Maximum Likelihood Estimator for the Tobit Model. *Econometrica* 46(5): 1211–1215. <http://www.jstor.org/stable/1911445>.
- Saez, E. 2010. Do Taxpayers Bunch at Kink Points? *American Economic Journal: Economic Policy* 2(3): 180–212.

#### **About the authors**

Marinho Bertanha is the Gilbert F. Schaefer Assistant Professor of Economics at the University of Notre Dame.

Andrew H. McCallum is a Principal Economist in the International Finance Division of the Board of Governors of the Federal Reserve System. He also teaches econometrics as an adjunct professor at the McCourt School of Public Policy at Georgetown University.

Alexis M. Payne is an economics Ph.D. student at Stanford University and was a senior research assistant in the International Finance Division of the Board of Governors of the Federal Reserve System. She received her B.A. in Economics from William & Mary in 2019.

Nathan Seegert is an assistant professor of finance in the Eccles School of Business at the University of Utah.