

Treatment Effects in Bunching Designs: The Impact of the Federal Overtime Rule on Hours

Leonard Goff*

This version: January 27, 2022
For current version [click here](#)

Abstract

The Fair Labor Standards Act mandates overtime premium pay for most U.S. workers, but a lack of changes to the rule has hindered assessment of its impacts on the labor market. I use bunching at 40 hours to estimate the effect of the federal overtime rule on work hours, pairing an extension of the “bunching design” identification strategy with data from individual weekly paychecks. I show that bunching at a choice-set kink partially identifies the average causal effect of the kink among bunchers, under weak assumptions on preferences and nonparametric restrictions on heterogeneity. The bounds are informative in the overtime context.

*Department of Economics, University of Georgia. I thank my PhD co-advisors Simon Lee and Suresh Naidu, as well as Michael Best, Daniel Hamermesh, and Bernard Salanié for their for their gracious advice and support throughout this project. I have also benefited from discussions with Christopher Ackerman, Doug Almond, Joshua Angrist, Jushan Bai, Iain Bamford, Marc Bellemare, Sandra Black, Carol Caetano, Brant Callaway, Ivan Canay, Gregory Cox, Junlong Feng, Bhargav Gopal, Jonas Hjort, Wojciech Kopczuk, Bentley MacLeod, Matthew Masten, Serena Ng, José Luis Montiel Olea, Dilip Ravindran, Miguel Urquiola, and seminar participants at Columbia, Duke, Lehigh, Microsoft, the University of Georgia, the University of North Carolina Greensboro, the US Census Bureau, and Washington State University, as well as audiences at the Canadian Economics Association and Southern Economics Association meetings, and the European Winter Meeting of the Econometric Society. I thank my main data provider as well as the Bureau of Labor Statistics. Online Appendices available [here](#).

1 Introduction

Many countries require premium pay for long work hours, in an effort to limit excessive work schedules and encourage hours to be spread over more workers. In the U.S., such regulation comes through the “time-and-a-half” rule of the Fair Labor Standards Act (FLSA): workers must be paid one and a half times their normal hourly wage for any hours they work in excess of 40 within a single week. Although some workers are exempt, the time-and-a-half rule applies to a majority of the U.S. workforce, including nearly all of its over 80 million hourly workers (U.S. Department of Labor, 2019). Workers in many industries average multiple overtime hours per week, making overtime the largest form of supplemental pay in the U.S. (Hart, 2004; Bishow, 2009).

Nevertheless, only a small literature has empirically examined the effects of the FLSA overtime rule on the U.S. labor market. This stands in marked contrast to the large body of work on the minimum wage, which was also introduced at the federal level by the FLSA in 1938. A key reason for this gap is that the overtime rule has varied little since then: the policy has remained as time-and-a-half after 40 hours in a week, for now more than 80 years. Reforms to overtime policy have been rare and have focused on eligibility, leaving the central parameters of the rule unaffected.¹ This lack of variation has afforded few opportunities to leverage research designs that exploit policy changes to identify causal effects.

This paper assesses the effect of the FLSA overtime rule on hours of work, taking a new approach that makes use of variation *within* the rule itself. The policy introduces a sharp discontinuity in the marginal cost of a worker-hour—a convex “kink” in firms’ costs—which provides firms with an incentive to set workers’ hours exactly at 40. Optimizing behavior on the part of firms predicts that the resulting mass of workers with hours at 40 in a given week will be larger or smaller depending on how responsive firms are to the wage variation imposed by time-and-a-half rule. Combining this observation with assumptions about the shape of the distribution of hours that would be chosen absent the FLSA, I use the bunching mass to identify the effect of the overtime rule on hours.

To do so, I develop a generalization of the “bunching design” identification strategy, which has used bunching at kinks in income tax liability to identify the elasticity of labor supply (Saez 2010; Chetty et al. 2011).² In particular, I give new identification results that hold under a set of weakened assumptions that may be suitable to a variety of empirical contexts, showing that the bunching design can be useful for program-evaluation questions such as assessing the FLSA.

¹Quach (2021) looks at very recent reforms to eligibility criteria for exemption from the FLSA, estimating effects of the expansion on employment and the incomes of salaried workers, but not on hours of work. My results complement the few studies that have used difference-in-differences approaches for effects on hours: Hamermesh and Trejo (2000) consider the expansion of a daily overtime rule in California to men in 1980, while Johnson (2003) use a supreme court decision on the eligibility of public-sector workers in 1985. Costa (2000) studies the initial phase-in of the FLSA in the years following 1938. See footnotes 32 and 33 for a comparison of my results to these papers.

²The same basic model has since been applied in a range of settings beyond income taxation. This paper considers only the bunching design for kinks, and not a related method for bunching at *notches* (e.g. Kleven and Waseem 2013).

In tax settings, the promise of the bunching design is to overcome endogeneity in the marginal tax rates that apply to different individuals while relying only the cross-sectional distribution of income near a threshold between tax brackets for identification. Analogously, the starting point in the overtime setting is to construct the distribution of hours for which workers are paid in a single week. This is itself a challenge, and previous studies have relied on self-reported integer hours of work from surveys such as the Current Population Survey. I obtain administrative data at the required level of detail via individual paycheck records from a large payroll processing company. These paychecks report the exact number of hours that a worker was paid for in a given week, allowing me to construct the distribution of hours-of-pay without measurement error or rounding.

With this data in hand, the goal is to translate features of the observed hours distribution into estimates of the overtime rule’s causal effect, under credible assumptions about how weekly working hours are determined. This requires moving beyond the standard bunching-design model popularized in public-finance applications, which assumes a stylized labor supply model with parametric “isoelastic” preferences and strong restrictions on heterogeneity.

The identifying assumptions of the bunching design can be separated into two parts: i) assumptions about how individual agents would make choices given counterfactual choice sets—a *choice model*, and ii) assumptions about the distribution of heterogeneity in choices across agents. I first show that the class of choice models under which the bunching-design can be applied is considerably more general than the benchmark isoelastic model and its variants. In particular, I find that the method does not require the researcher to suppose any explicit functional form for decision-makers’ utility; rather, the core behavioral prediction driving identification rests on *convexity* of individual preferences. In my formulation, agents in the bunching design can also have multiple underlying margins of choice, which might be unobserved to the researcher and vary by observational unit.³ These findings establish an important robustness property for the bunching design: it rests on a prediction about behavior that remains broadly valid even when the parametric utility model typically used to motivate the design is misspecified.

This generality is accomplished by recasting the bunching design in a potential outcomes framework, defining the parameter of interest in terms of counterfactual *choices* rather than as a preference parameter from a tightly specified choice model. I show that choice from a kinked choice set can be fully characterized by two such counterfactuals, and that bunching directly identifies a feature of their joint distribution. In the overtime setting I take firms to choose the hours of workers, and these potential outcomes correspond to: a) the number of hours the firm would choose for the worker this week if the worker’s normal wage rate applied to all such hours; and b) the number that the firm would choose if the worker’s overtime rate applied to all hours this week. Comparing

³This property of my choice model also generalizes Blomquist et al. (2021), who consider a bunching-design setup with nonparametric utility but with a scalar choice variable.

these counterfactuals speaks directly to the effect of the FLSA policy on hours of work.

In addition to generalizing the choice model underlying the bunching design, I show that the assumptions about heterogeneity typically used in the bunching design can also be relaxed. In my formulation, these take the form of assumptions about the marginal distributions of the two potential outcomes, which are observed in a censored manner. As emphasized by Blomquist and Newey (2017), the bunching design requires extrapolating such distributions beyond where they are observed to estimate causal effects. To perform this extrapolation I impose a weak nonparametric shape constraint—*bi-log-concavity*—on the distribution of each potential outcome. Bi-log-concavity nests many previously proposed assumptions for bunching analyses, and leads to a natural falsification test. The restriction affords partial identification of a conditional average treatment effect among individuals who are in fact located at the kink, a parameter I call the “buncher ATE”. In the overtime context, the buncher ATE represents an average reduced-form wage elasticity of hours demand, which I then use to assess the overall effect of the FLSA.

This result supplements other partial-identification approaches recently proposed for the bunching design. Importantly, the bounds I derive for the buncher ATE are substantially narrowed by making extrapolation assumptions separately for *each* of the two counterfactuals. By contrast, existing approaches operate by constraining the distribution of a single scalar heterogeneity parameter, a simplification afforded by the isoelastic choice model. In the context of that model, Bertanha et al. (2020) and Blomquist et al. (2021) obtain bounds on the elasticity when the researcher is willing to put an explicit limit on how quickly the density of heterogeneous choices can change. My approach based on bi-log-concavity avoids the need to choose any such tuning parameters, and unlike previous approaches is applicable in the general choice model.

The empirical setting of overtime pay does involve two challenges that are not typical of existing bunching-design analyses. Firstly, 40 hours is not an “arbitrary” point and bunching there could arise in part from factors other than it being the location of the kink. I use two strategies to estimate the amount of bunching that would exist at 40 absent the FLSA, and deliver clean estimates of the rule’s effect. My preferred strategy exploits the fact that when a worker makes use of paid-time-off hours, these do not count towards that week’s overtime threshold. This shifts the location of the kink in a plausibly idiosyncratic way, letting me identify the distinct contribution of the FLSA to bunching. A second feature of the overtime setting is that work hours may not always be set unilaterally by one party—in principle either the firm or the worker could have control over a worker’s schedule. I provide evidence that week-by-week variation in hours tends to be driven by the firm, and outline a conceptual framework to motivate this observation. I also extend the results to a model in which bargaining weight between workers and firms is arbitrary and heterogeneous.

I find that the FLSA overtime rule does in fact reduce hours of work among hourly workers, despite theoretical reasons to doubt the existence of any such effect (Trejo, 1991). My preferred

estimate suggests that about one quarter of the bunching observed at 40 among hourly workers is due to the FLSA, and those working at least 40 hours work, on average, about 30 minutes less than they would absent the time-and-a-half rule. Across specifications I estimate that the local wage elasticity of hours demand close to 40 falls in the range -0.04 to -0.19 , indicating that firms are fairly resistant to changing hours to avoid overtime payments. While the bunching design is only directly informative about the hours effects of the overtime rule—and not its employment effects—a back-of-the-envelope calculation using this estimate suggests that FLSA regulation creates about 700,000 jobs. I also use these estimates to evaluate proposed reforms to the FLSA: for example lowering the overtime threshold below 40 hours,⁴ or increasing the premium pay factor from 1.5 to 2. I find that even in my generalized bunching-design model the data can be informative about counterfactual policies that change the location or “sharpness” of a kink.

The structure of the paper is as follows. Section 2 lays out a motivating conceptual framework for work hours that relates my approach to existing literature on overtime. Section 3 introduces the payroll data I use in the empirical analysis. In Section 4 I develop the generalized bunching-design approach, with Appendix A developing some of the supporting formal results. Section 5 applies these results to estimate effect of the FLSA overtime rule on hours worked, as well as the effects of hypothetical reforms to the FLSA. Section 6 discusses the empirical findings from the standpoint of policy objectives, and 7 concludes.

2 Conceptual framework

This section outlines a framework for thinking about the determination of weekly hours, which motivates the identification strategy of Section 4. Readers primarily interested in the bunching design may wish to skip directly to that section.

The conceptual framework is centered around two observations from the data described in Section 3: weekly hours vary considerably between pay periods for an individual hourly worker, and a given worker’s hourly wage tends to change infrequently. I propose to view this as a two stage-process. First, workers are hired with an hourly wage set along with an “anticipated” number of hours they will work per week. Then, with that hourly wage fixed in the short-run, final scheduling of hours is controlled by the firm and varies by week given shocks to the firm’s demand for labor. Given the FLSA overtime rule and a worker’s fixed wage, their employer thus faces a kinked cost schedule when choosing hours in a given week, as pictured in Figure 1.

⁴For example, HR4721 introduced in 2021 would establish a 32-hour workweek for the purposes of overtime pay: [www.congress.gov/bill/117th-congress/house-bill/4728](https://www.congress.gov/bills/117/congress/house-bills/4728)

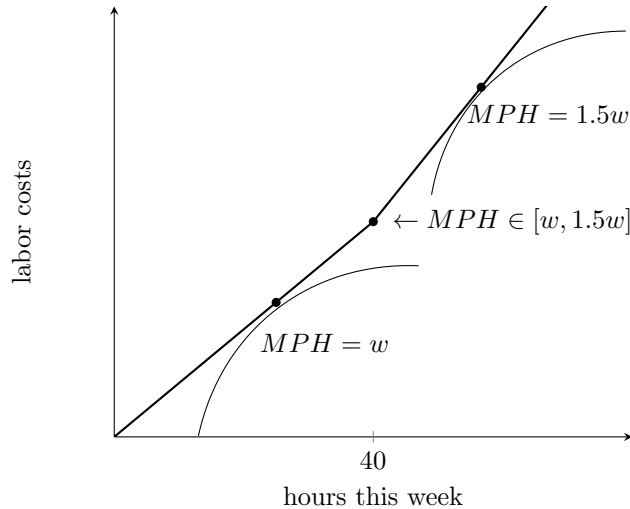


FIGURE 1: With a given worker's straight-time wage fixed at w , labor costs as a function of hours have a convex kink at 40 hours, given the overtime rule. Simple models of week-by-week hours choice (see Section 4.2) yield bunching when the marginal product of an hour at 40 is between w and $1.5w$ for a mass of workers.

Wages and anticipated hours set at hiring

We begin with the hiring stage, which pins down a worker's wage. Throughout the analysis, I focus on workers paid on an hourly basis. Following the literature I refer to the hourly rate of pay that applies to the first 40 of a worker's hours as their *straight-time wage* or simply *straight wage*. This section discusses a benchmark model to endogenize these straight wages, which is spelled out formally in Online Appendix 4.1. The basic bunching design strategy of Section 4 only requires that *some* straight-time wage is agreed upon and fixed in the short-run for each worker, as is indeed observed in the data. However the conceptual framework of this section will play a role in my final evaluation of the FLSA, in Section 4.4.

Suppose that firms hire by posting an earnings-hours pair (z, h) , where z is total weekly compensation offered to each worker, and h is the number of hours they are each expected to work per week. The firm faces a labor supply function $N(z, h)$ determined by workers' preferences over the labor-leisure tradeoff,⁵ and makes a choice of (z^*, h^*) given this labor supply function and their production technology. For simplicity, workers are here taken to be homogeneous in production, all paid hourly, and all covered by the FLSA.

While labor supply is viewed as a function over *total* compensation z and hours, there is always a unique straight wage associated with a particular (z, h) pair, such that h hours of work yields

⁵This labor supply function can be viewed as an equilibrium object that reflects both worker preferences and the competitive environment for labor. In the Online Appendix 4.2, I endogenize this function in a simple extension of the imperfectly competitive Burdett and Mortensen (1998) search model.

earnings of z , given the FLSA overtime rule:

$$w_s(z, h) = \frac{z}{h + 0.5 \cdot \mathbb{1}(h > 40)(h - 40)} \quad (1)$$

We can distinguish between the two main views on the likely effects of overtime policy by supposing that a workers' straight-time wage is set according to Eq. (1), given values z^* and h^* that the firm and worker agree upon at the time of hiring. Trejo (1991) calls these two views the *fixed-job* and the *fixed-wage* models of overtime.

The *fixed-job* view observes that for a generic smooth labor supply function $N(z, h)$ (and revenue production function with respect to hours), the optimal job package (z^*, h^*) for the firm to post will be the same as the optimal one absent the FLSA, as the hourly wage rate simply adjusts to fully neutralize the overtime premium.⁶ Suppose for the moment that all workers then in fact work exactly h^* hours in all weeks (abstracting away from any reasons for the firm to deviate from h^* in any given week). Then the FLSA would have no effect on earnings, hours or employment, provided that $w_s(z^*, h^*)$ is above any applicable minimum wage.

On the *fixed-wage* view, the firm faces an exogenous straight-time wage when determining hours. Versions of this idea are considered in Brechling (1965), Rosen (1968), Ehrenberg (1971), Hamermesh (1996), Hart (2004) and Cahuc and Zylberberg (2004). In a static model, this can be captured by a discontinuous labor supply function $N(z, h)$: one that exhibits perfect competition on the quantity $w_s(z, h)$. I show in Online Appendix 4.1 that in this case h^* and z^* are pinned down by the concavity of production with respect to hours and the scale of fixed costs (e.g. training) that do not depend on hours. The fixed-wage job makes the clear prediction that the FLSA will cause a reduction in hours, and bunching at 40.⁸

Trejo (1991) and Barkume (2010) investigate whether the fixed-job or fixed-wage model better accords with the observed joint distribution of hourly wages and hours. They find that wages do tend to be lower among jobs with overtime pay provisions and more overtime hours, but by a magnitude smaller than would be predicted by the fixed-jobs model. However, these estimates could be driven by selection, e.g. of lower-skilled workers into covered jobs with longer hours.

In Online Appendix 1, I construct an empirical test of Equation (1) that is instead based on assuming that the conditional distribution (across individual paychecks) of z is smooth across $h = 40$. I find that roughly one quarter of paychecks around 40 hours reflect the wage/hours relationship predicted by the fixed-job model. This finding, along with the observation that hours change much

⁶In Appendix 4.1 I give a closed-form expression for (z^*, h^*) when both labor supply and production are iso-elastic: hours and earnings are each⁷ increasing in the elasticity of labor supply with respect to earnings, and decreasing in the magnitude of the elasticity of labor supply with respect to pay.

⁸A fixed-wage model tends to predict an overall positive effect on employment given plausible assumptions on substitution between labor and capital (Cahuc and Zylberberg, 2004), though the total number of labor-hours will decrease (Hamermesh, 1996).

more frequently than wages, is consistent with a model in which once straight-wages are set according to Equation (1), they remain fairly static over time.⁹ In common with the fixed wage model, this two-stage framework allows for the possibility that the overtime rule affects hours, and predicts bunching at 40. However, this is driven by short-run rigidity in straight-wages, rather than by perfect competition.

Dynamic adjustment to hours by week

Confronted with the observation that workers' hours vary considerably week-to-week in my sample of hourly workers, I assume that this week-level variation reflects choices made by their employers. There are many reasons to expect variation in firms' demand for hours over time. Shocks to product demand or productivity change the number of weekly hours that would be optimal that week from the firm's perspective. If demand for the firm's products is seasonal or volatile, it may not be worthwhile to hire additional workers only to reduce employment later. Similarly, variation in productivity across workers may only become apparent to supervisors after their straight wages have been set, and it may be profitable to increase the hours of the most productive workers.

Throughout Section 4, I maintain a strong version of the assumption that a firm—rather than a worker—chooses the hours I observe on each paycheck. Workers' preferences *do* matter in the determination of each worker's straight wage at hiring, but I assume the firm has final scheduling rights week by week.¹⁰ This eases notation and emphasizes the intuition behind my identification strategy. Online Appendix 2 presents a generalization in which some fraction of workers choose their hours, along with intermediate cases in which the firm and worker bargain over hours each week. If only some workers have full control of their weekly hours, then the bunching-design strategy will only be informative about effects of the FLSA among workers whose final hours are chosen by the firm.

Available survey evidence suggests that this group is the dominant one: a relatively small share of workers report that they choose their own schedules (despite the ongoing increase in flexible work arrangements). For example, the 2017-2018 Job Flexibilities and Work Schedules Supplement of the American Time Use Survey asks workers whether they have some input into their schedule, or whether their firm decides it. Only 17% report that they have some input. In a survey of firms, about 10% report that most of their employees have control over their shifts (Society for Human Resource Management, 2018).

⁹This dovetails other recent evidence of uniformity and discretion in wage-setting, e.g. nominal wage rigidity (Grigsby et al. 2020), wage standardization (Hjort et al., 2020) and bunching at round numbers (Dube et al., 2020).

¹⁰This can be rationalized on the basis of workers generally having less bargaining power: if the worker and firm fail to agree on a worker's hours, the worker's outside option may be unemployment while the firm's outside option is having one less worker (Stole and Zwiebel, 1996).

3 Data and descriptive patterns

The main dataset I use comes from a large payroll processing company. They provided anonymized paychecks for the employees of 10,000 randomly sampled employers, for all pay periods in the years 2016 and 2017. At the paycheck level, I observe the check date, straight wage, and amount of pay and hours corresponding to itemized pay types, including normal (straight-time) pay, overtime pay, sick leave, holiday pay, and paid time off. The data also include state and industry for each employer and for employees: age, tenure, gender, state of residence, pay frequency and their salary if one is stored in the system.

3.1 Sample description

I construct a final sample for analysis based on two desiderata: a) the ability to observe hours within a single week; and b) a sample only of workers who are not exempt from the FLSA overtime rule. For the purposes of a), I keep paychecks from workers who are paid on a weekly basis (roughly half of the workers in the sample). To achieve b) I focus on hourly workers, since nearly all workers who are paid hourly are subject to FLSA regulation. The final sample includes 630,217 paychecks for 12,488 workers across 566 firms. Further details of the sample construction are provided in Online Appendix 1.1.

Table 1 shows how the sample compares to survey data that is constructed to be representative of the U.S. labor force. Column (1) reports means from the final sample used in estimation, while (2) reports means before sampling restrictions. Column (3) reports means from the Current Population Survey (CPS) for the same years 2016–2017, among individuals reporting hourly employment. The “gets overtime” variable for the CPS sample indicates that the worker usually receives overtime, tips, or commissions. Column (4) reports means for 2016–2017 from the National Compensation Survey (NCS), a representative establishment-level dataset accessed on a restricted basis from the Bureau of Labor Statistics. The NCS reports typical overtime worked at the quarterly level for each job in an establishment, drawn from administrative data when possible.¹¹

The sample I use is somewhat more male, earns lower straight-time wages, and works more overtime than a typical hourly worker in the U.S. Column (2) in Table 1 reveals that my sampling restrictions can explain why the estimation sample tilts male and has higher overtime hours than the workforce as a whole. In particular, conditioning on workers that are paid on a weekly basis oversamples industries that tend to have more men, and tend to pay somewhat lower wages. Appendix 1 compares the industry and regional distributions of the estimation sample to the CPS.

¹¹The hourly wage variable for the CPS may mix straight-time and overtime rates, and is only present in the outgoing rotation group sample. The tenure variable comes from the 2018 Job Tenure Supplement. The NCS does not distinguish between hourly and salaried workers, reporting an average hourly rate that includes salaried workers, who tend to be paid more. This likely explains the higher value than the CPS and payroll samples.

	(1)	(2)	(3)	(4)
	Estimation sample	Initial sample	CPS	NCS
Tenure (years)	3.21	2.81	6.34	.
Age (years)	37.15	35.89	39.58	.
Female	0.23	0.46	0.50	.
Weekly hours	38.92	27.28	36.31	35.70
Gets overtime	1.00	0.37	0.17	0.52
Straight-time wage	16.16	22.17	18.09	23.31
Weekly overtime hours	3.56	0.94	.	1.04
Number of workers in sample	12488	149459	63404	228773

TABLE 1: Comparison of the sample with representative surveys. Columns 1 and 2 average across periods within worker from the administrative payroll sample, and then present means across workers. Column 2 presents means of worker-level data from the Current Population Survey and Column 3 averages representative job-level data from the National Compensation Survey.

3.2 Hours and wages in the sample

I turn now to the main variables to be used in the analysis. Figure 2 reports the distribution of hours of work in the pooled sample of paychecks. The graphs indicate a large mass of individuals who were paid for exactly 40 hours that week, amounting to about 11.6% of the sample.¹² Appendix Figure 8 shows that overtime pay is present in nearly all weekly paychecks that report more than 40 hours, in line with the presumption that workers in the final sample are not FLSA-exempt.

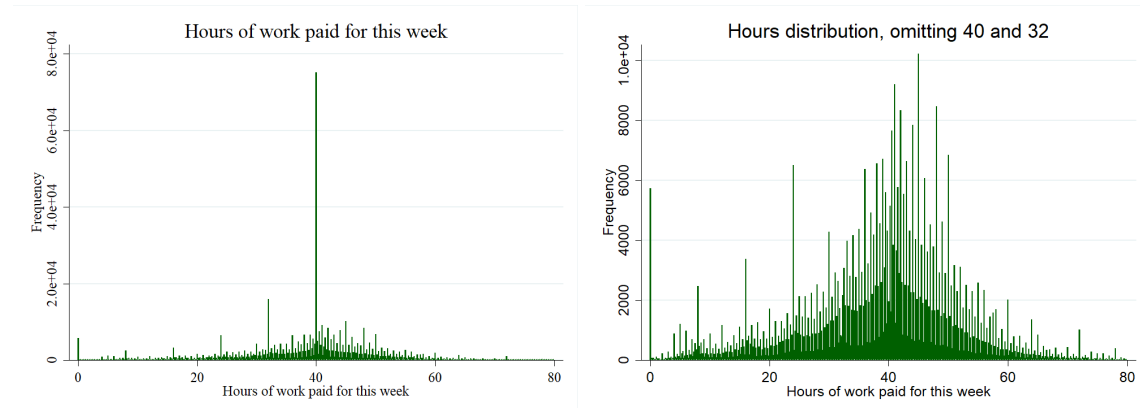


FIGURE 2: Empirical densities of hours worked pooling all paychecks in final estimation sample. Sample is restricted to hourly workers receiving overtime pay at some point (to ensure nearly all are non-exempt from FLSA, see text), and workers having hours variation. The right panel omits the points 40 and 32 to improve visibility elsewhere. Bins have a width of 1/8 of an hour, based on the observed granularity of hours (see Appendix Figure 10 for details).

Table 2 documents that while the hours paid in 70% of all pay checks in the final estimation sample differ from those of the last paycheck by at least one hour, just 4% of all paychecks record

¹²The second largest mass occurs at 32 hours, and is explained by paid-time-off, holiday, and sick pay hours as discussed in Section 5.

a different straight-time wage than the previous paycheck for the same worker. Among the roughly 22,500 wage change events, the average change is about a 45 cent raise per hour. When hours change the magnitude is about 7 hours on average (see Appendix Figure 10 for the distribution of hours changes), with no average secular increase in hours over time.

Online Appendix 1 reports some further details on the variation in hours and wages present in the data. Appendix Table 2 regresses hours, overtime hours, and an indicator for bunching on worker observables, and shows that after controlling for worker and date fixed effects bunching and overtime hours are both predicted by recent hiring at the firm, lending further evidence for the assumption that shocks to labor demand drive variation in hours. Appendix Table 3 shows that overall, about 63% of variation in total hours can be explained by worker and employer-by-date fixed effects. Appendix Figure 1 documents heterogeneity in the prevalence of overtime pay across industry classifications. Appendix Table 1 reports a direct estimate of the proportion of wage-hours pairs that are related according to Equation (1), among workers near 40.

	Mean	Std. dev.	N
Indicator for hours changed from last period	0.84	0.37	630,217
Indicator for hours changed by at least 1 hour	0.70	0.46	630,217
Indicator for wage changed from last period	0.04	0.19	630,217
Indicator for wage changed, if hours changed	0.04	0.19	529,791
Absolute value of hours difference, if hours changed	6.83	8.23	529,791
Difference in wage, if wage changed	0.45	26.46	22,501

TABLE 2: Changes in hours or straight wages between a worker’s consecutive paychecks.

4 Empirical strategy: a generalized kink bunching design

Let us now turn to the firm choosing the hours of a given worker in a particular week, with that worker’s wage fixed and costs a kinked function of hours as depicted in Figure 1. This section shows that under weak assumptions, firms facing such kinks will make choices that can be completely characterized by choices they *would* make under two counterfactual linear cost schedules that differ with respect to wage. I relate the observable bunching at 40 hours to a treatment effect defined from these two counterfactuals, which I then use to estimate the impact of the FLSA on hours.

The identification results in this section hold in a much more general setting in which a decision-maker faces a generic “kinked” choice set and has convex preferences. I present this general model in Appendix A. Throughout this section I refer to a worker i in week t as a *unit*: an observation of h_{it} for unit it is thus the hours recorded on a single paycheck. Probability statements are to be understood with respect to such paycheck-level units.

4.1 A general choice model

Let us start from the conceptual framework introduced in Section 2. In choosing the hours h_{it} of worker i in week t , worker i 's employer faces a kinked cost schedule, given the worker's straight-time wage this week w_{it} . If the firm chooses less than 40 hours, it will pay $w = w_{it}$ for each hour, and if the firm chooses $h > 40$ it will pay $40w$ for the first 40 hours and $1.5w(h - 40)$ for the remaining hours, giving the convex shape to Figure 1. We can write the kinked pay schedule for unit it , as a function of hours this week h , as

$$B_{kit}(h) = w_{it}h + .5w_{it}\mathbb{1}(h > 40)(h - 40) = \max\{B_{0it}(h), B_{1it}(h)\}$$

where $B_{0it}(h) = w_{it}h$ and $B_{1it}(h) = 1.5w_{it}h - 20w_{it}$. The kinked pay schedule $B_{kit}(h)$ is equal to $B_{0it}(h)$ for values $h \leq 40$ and $B_{kit}(h)$ is equal to $B_{1it}(h)$ for values $h \geq 40$. The functions B_0 and B_1 recover the two segments in Figure 1 when restricted to these domains, and are depicted separately in Appendix Figure A.1.

Definition (potential outcomes). Let h_{0it} denote the hours of work that of unit it would be paid for if instead of $B_{kit}(h)$, the pay schedule for this week's hours were $B_{0it}(h)$. Similarly, let h_{1it} denote the hours of pay that would occur for unit it if the pay schedule were $B_{1it}(h)$.

By contrast, let h_{it} denote the actual hours of work for which unit it is paid. Our first assumption is that actual hours and potential outcomes reflect choices made by the firm:

Assumption CHOICE. Each of h_{0it} , h_{1it} and h_{it} reflect choices the firm would make under the pay schedules $B_{0it}(h)$, $B_{1it}(h)$, and $B_{kit}(h)$ respectively.

CHOICE reflects the assumption that hours are perfectly manipulable by firms. Note that if firm preferences over a unit's hours are quasi-linear with respect to costs (e.g. if they maximize weekly profits), the term $-20w_{it}$ appearing in B_{1it} plays no role in firm choices. As such, I will often refer to h_{1it} as choice made under linear pay at the overtime rate $1.5w_{it}$, keeping in mind that the definition given above is necessary for the analysis if preferences are not quasi-linear.

Our second assumption is that each unit's firm optimizes some vector \mathbf{x} of choice variables that pin down that unit's hours. As a leading case, we may think of hours of work as a single component of firms' choice vector \mathbf{x} (Appendix A.3 gives some examples). Firm preferences are taken to be convex in \mathbf{x} and the unit's wage costs z :

Assumption CONVEX. Firm choices for unit it maximize some $\pi_{it}(z, \mathbf{x})$, where π_{it} is strictly quasiconcave in (z, \mathbf{x}) and decreasing in z . Hours are a continuous function of \mathbf{x} for each unit.

For ease of notation, I here state a version of CONVEX that is a bit stronger than necessary for the identification results below. In particular, Appendix A relaxes CONVEX to allow "double-peaked"

preferences with one peak located exactly at the kink. The appendix also shows that bunching still has some identifying power under no assumptions about convexity of preferences. That firms rather than workers choose hours enters in the assumption that π is decreasing (rather than increasing) in z , but Appendix 2 relaxes this to allow some workers to set their hours.

The starting point for our analysis of identification in the bunching design is the following mapping between actual hours h_{it} and the counterfactual hours choices h_{0it} and h_{1it} . Appendix Lemma 1 shows that Assumptions CHOICE AND CONVEX imply that:

$$h_{it} = \begin{cases} h_{0it} & \text{if } h_{0it} < 40 \\ 40 & \text{if } h_{1it} \leq 40 \leq h_{0it} \\ h_{1it} & \text{if } h_{1it} > 40 \end{cases} \quad (2)$$

That is, a worker will work h_{0it} hours when the counterfactual choice h_{0it} is less than 40, and h_{1it} hours when h_{1it} is greater than 40. They will be located at the corner solution of 40 if and only if the two counterfactual outcomes fall on either side, “straddling” the kink.¹³ Figure 3 depicts the implications of Eq. (2) for what is observable by the researcher in the bunching design: censored distributions of both h_0 and of h_1 , and a point-mass of size $\mathcal{B} = P(h_{1it} \leq 40 \leq h_{0it})$ at the kink.

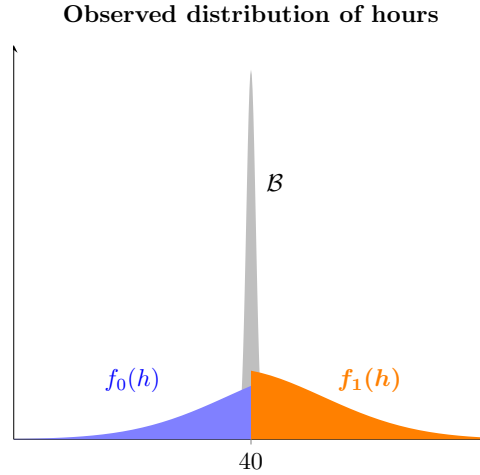


FIGURE 3: Observables in the bunching design, given Equation (2). To the left of the kink at 40, the researcher observes the density $f_0(h)$ of the counterfactual h_{0it} , up to values $h = 40$. To the right of the kink, the researcher observes the density $f_1(h)$ of h_{1it} for values $h > 40$. At the kink, one observes a point-mass of size $\mathcal{B} := P(h_{it} = 40) = P(h_{1it} \leq 40 \leq h_{0it})$.

Equation (2) represents a departure from previous approaches to the bunching design, which char-

¹³“Straddling” can only occur in one direction, with $h_{1it} \leq k \leq h_{0it}$. The other direction: $h_{0it} \leq k \leq h_{1it}$ with at least one inequality strict, is ruled out by the weak axiom of revealed preference (see Appendix A).

acterize bunching in terms of the counterfactual h_0 only. We will see that this is a simplification afforded by the benchmark isoelastic utility model, but in a generic choice model, both h_0 and h_1 are necessary to pin down actual choices h_{it} . Appendix A shows that Eq. (2) also holds in settings with possibly non piecewise-linear kinked choice sets of the form: $z \geq \max\{B_0(\mathbf{x}), B_1(\mathbf{x})\}$ where B_0 and B_1 are weakly convex in the full vector \mathbf{x} , and z any “cost” decision-makers dislike.

Intuition for Equation (2)

As an intuitive illustration of Equation (2), suppose that firms balance the cost $B_{kit}(h)$ against the value of h hours of the worker’s labor, in order to maximize that week’s profits. Then Eq. (2) can be written:

$$h_{it} = \begin{cases} MPH_{it}^{-1}(w_{it}) & \text{if } MPH_{it}(40) < w_{it} \\ 40 & \text{if } MPH_{it}(40) \in [w_{it}, 1.5w_{it}] \\ MPH_{it}^{-1}(1.5w_{it}) & \text{if } MPH_{it}(40) > 1.5w_{it} \end{cases} \quad (3)$$

where denotes $MPH_{it}(h)$ is the marginal product of an hour of labor for unit it , as a function of hours h . Assuming that production is strictly concave, the function $MPH_{it}(h)$ will be strictly decreasing in h , and we have that $h_{0it} = MPH_{it}^{-1}(w_{it})$ and $h_{1it} = MPH_{it}^{-1}(1.5w_{it})$.

Figure 1 depicts Eq. (3) visually. Consider for example a worker with a straight-wage of \$10 an hour. If there exists a value $h < 40$ such that the worker’s MPH is equal to \$10, then the firm will choose this point of tangency. This happens if and only if the marginal product of an hour at 40 hours is less than \$10. If instead, the marginal product of an hour is still greater than \$15 at $h = 40$, the firm will choose the value $h > 40$ such that MPH equals \$15. The third possibility is that the MPH at $h = 40$ is *between* the straight and overtime rates \$10 and \$15. In this case, the firm will choose the corner solution $h = 40$, contributing to bunching at the kink.

While Eq. (3) provides a natural nonparametric characterization of when the firm will ask a worker to work overtime (when the ratio of productivity to wages is high), it is still more restrictive than necessary for the purposes of the bunching design. Appendix A.3 provides some examples that use the full generality of CONVEX, in which firms simultaneously consider *multiple* margins of choice aside from a given unit’s hours. For example, the firm may attempt to mitigate the added cost of overtime by reducing bonuses when a worker works many overtime hours. Eq. (2) remains valid even when such additional margins of choice are unmodeled and unobserved by the econometrician, varying possibly by unit.

Note: if production depends jointly on the hours of all workers within a firm, we may expect the function $MPH_{it}(h)$ in Eq. (3) to depend on the hours of worker i ’s colleagues in week t . In this case the quantities h_{0it} and h_{1it} hold the hours of i ’s colleagues fixed at their *realized* values: they contemplate counterfactuals in which the pay schedule for a single unit it is varied, and nothing

else. This affects the interpretation of our treatment effects, as discussed in Section 4.4. In the baseline isoelastic model that we consider in the next section, such interdependencies between workers' hours are ruled out by virtue of assuming that production is linearly separable across units. Online Appendix 3 discusses the case of a general non-separable production function.

4.2 The benchmark isoelastic model

The canonical approach from the bunching-design literature (Saez, 2010; Chetty et al., 2011; Kleven, 2016), strengthens Assumption CONVEX to suppose that $\mathbf{x} = h$ and decision-makers' utility features an isoelastic functional form, with preferences common between units up to a scalar heterogeneity parameter. This corresponds to a situation in which firm profits from unit it take the form:

$$\pi_{it}(z, h) = a_{it} \cdot \frac{h^{1+\frac{1}{\epsilon}}}{1+\frac{1}{\epsilon}} - z \quad (4)$$

where $\epsilon < 0$ is common across units, and z represents wage costs for worker i in week t . Eq. (4) is analogous to the isoelastic, quasilinear labor *supply* model used in the context of tax kinks.

Under a linear pay schedule $z = wh$, the profit maximizing number of hours is $\left(\frac{w}{a_{it}}\right)^\epsilon$, so ϵ yields the elasticity of hours demand with respect to a linear wage. Let $\eta_{it} = a_{it}/w_{it}$ denote the ratio of a unit's current productivity factor to their straight wage. In the isoelastic model $h_{0it} = \eta_{it}^{-\epsilon}$ and $h_{1it} = 1.5^\epsilon \cdot \eta_{it}^{-\epsilon}$, and by Eq. (3) actual hours h_{it} are ranked across units in order of η_{it} .¹⁴ If η_{it} is continuously distributed over a region containing the interval $[40^{-1/\epsilon}, 1.5 \cdot 40^{-1/\epsilon}]$, then the observed distribution of h_{it} will feature a point mass at 40—"bunching"—and a density elsewhere.

The isoelastic model yields a case of Eq. (3) in which $MPH_{it}(h)$ depends on it only through the heterogeneity parameter a_{it} , which reflects a productivity shock for that unit. Whether a worker has overtime hours in a given week is determined by the ratio η_{it} : a worker with a wage w_{it} fixed throughout the year may for example work overtime in periods when a_{it} is relatively high due to seasonally elevated productivity.

Identification in the isoelastic model

In the context of the isoelastic model, a natural starting place for evaluating the FLSA is to estimate the parameter ϵ . Ignoring for the moment any effects of the policy on straight-wages, the effect of the time-and-a-half rule on unit it 's hours will simply be the difference $h_{it} - h_{0it}$, what we might call the *effect of the kink*. It follows from the above that the effect of the kink is $h_{it} \cdot (1 - 1.5^{-\epsilon})$ for any unit such that $h_{it} > 40$. Provided the value of ϵ , we could thus evaluate the effect of the kink

¹⁴In particular $h_{0it} < 40$ whenever $\eta_{it} < 40^{-1/\epsilon}$, $h_{1it} > 40$ whenever $\eta_{it} > 1.5 \cdot 40^{-1/\epsilon}$ and $h_{1it} \leq 40 \leq h_{0it}$ when η_{it} falls in the intermediate range $[40^{-1/\epsilon}, 1.5 \cdot 40^{-1/\epsilon}]$.

for any paycheck recording overtime using that unit's observed hours.

The classic bunching-design method pioneered by Saez (2010) identifies ϵ by relating it to the observable bunching probability:

$$\mathcal{B} := P(h_{it} = 40) = \int_{40}^{1.5^{|\epsilon|} \cdot 40} f_0(h) \cdot dh \quad (5)$$

where f_0 is the density of h_0 . If the function f_0 were known, the value of ϵ could be pinned down in a straightforward way from Eq. (5). However, f_0 is not globally identified from the data: by Figure 3 we can see that f_0 is only identified to the left of the kink, while the density of h_1 is identified to the right of the kink. In the isoelastic model, it is convenient to analyze observables after applying a log transformation to hours. Since $h_{1it} = 1.5^\epsilon \cdot h_{0it}$, treatment effects in logs: $\delta = \ln h_{0it} - \ln h_{1it} = |\epsilon| \cdot \ln 1.5$, are homogeneous across all units it , and the density of $\ln h_{1it}$ is thus a simple leftward shift of the density of $\ln h_{0it}$, by δ , as shown in Figure 4.

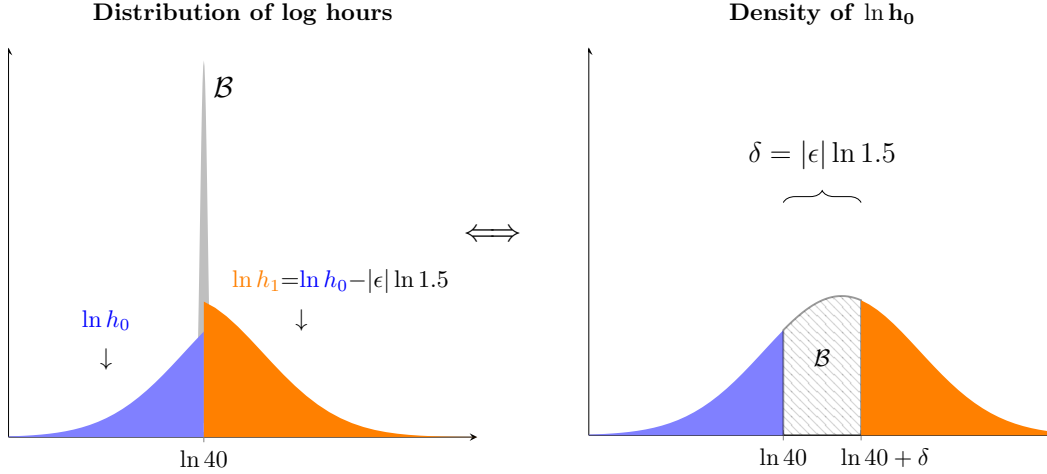


FIGURE 4: The left panel depicts the distribution of observed log hours $\ln h_{it}$ in the isoelastic model, while the right panel depicts the underlying full density of $\ln h_{0it}$. Specializing from the general setting of Figure 3, we have in the isoelastic model that $f_1(h) = f_0(h + |\epsilon| \cdot \ln 1.5)$. Thus, the full density of f_0 is related to the observed distribution by “sliding” the observed distribution for $h > 40$ out by the unknown distance $\delta = |\epsilon| \ln 1.5$, leaving a missing region in which f_0 is unobserved. The total area in the missing region from $\ln 40$ to $\ln 40 + \delta$ must equal the observed bunching mass \mathcal{B} .

Standard approaches in the bunching design make parametric assumptions that interpolate f_0 through the missing region of Figure 4 to point-identify ϵ .¹⁵ The approach of Saez (2010), for example, assumes for example that the density of h_0 is linear through the missing region of Figure 4. The popular method of Chetty et al. (2011) instead fits a global polynomial, using the distribution

¹⁵(Bertanha et al., 2020) note that given a full parametric model for f_0 , the entire model could be estimated by maximum likelihood. This approach would enforce (5) automatically while enjoying the efficiency properties of MLE.

of hours outside the missing region to impute the density of h_0 within it. Neither approach is particularly suitable in the overtime context. The linear method of Saez (2010) implies monotonicity of the density in the missing region, which is unlikely to hold given that 40 appears to be near the mode of the h_0 latent hours distribution. The method of Chetty et al. (2011) ignores the “shift” by δ in the right panel of Figure 4.¹⁶

If on the other hand, the researcher is unwilling to assume anything about the density of h_0 in the missing region of Figure 4, then the data are compatible with any finite $\epsilon < 0$, a point emphasized by Blomquist et al. (2021) and Bertanha et al. (2020). In particular, given the integration constraint (5), an arbitrarily small $|\epsilon|$ could be rationalized by a density that spikes sufficiently high just to the right of 40, while an arbitrarily large $|\epsilon|$ can be reconciled with the data by supposing that the density of h_0 drops quickly to some very small level throughout the missing region. I find a middle ground by imposing a nonparametric shape constraint on h_0 : *bi-log-concavity* (BLC), leading to a partial identification result. A detailed discussion of BLC is given in Section 4.3.

Limitations of the isoelastic model

Compared with the isoelastic model, the general choice model from Section 4.1 allows a wide range of underlying structural choice models that might drive a firm’s hours response to the FLSA. This robustness over structural models is important in the overtime context. As reported in Online Appendix 1.5, assuming either that h_0 is BLC or using the linear density method of Saez (2010) alongside the isoelastic model suggests that $\epsilon \approx -0.2$.¹⁷ This is implausible when interpreted through the lens of Equation (4): it implies an hours production function of $f(h) = -\frac{1}{4}h^{-4}$ (up to an affine transformation), which features an unreasonable degree of concavity. Allowing a more general non-exponential production function $f(h)$ (separable between units) is also not much help, as the standard bunching design approach then estimates an averaged local inverse elasticity of $f(h)$ (see Online Appendix 1.5). In short, the observed bunching is too small to be reconciled with a model in which ϵ parameterizes the concavity of weekly production with respect to hours. The estimand of the bunching design should instead be interpreted as a *reduced-form* elasticity of the demand for hours, which may reflect adjustment by firms along additional margins that can attenuate the hours response. Appendix A.3) discusses some examples of this type.

¹⁶This is perhaps less problematic in typical settings where the bunching is somewhat diffuse around the kink. However in my setting bunching is exact, and the slope of the density is far from zero near 40.

¹⁷This estimate is from the pooled sample across all industries, and attributes all of the bunching observed at 40 to the FLSA. Note that attributing just a portion of the observed bunching at 40 to the FLSA (as I do in Section 5.1) would only further reduce the magnitude of ϵ . This is also not driven by aggregating industries: estimation by industry yields bounds on ϵ ranging from -0.26 to -0.06 , which are similarly implausible as estimates of concavity of production.

4.3 Identifying treatment effects in the general choice model

In this section I turn to identification in the general choice model of Section 4.1. Without a single preference parameter like ϵ that characterizes responsiveness to incentives for all units, we face the following question: what parameter of interest might be identifiable from the data without the restrictive isoelastic model, but still help us to evaluate the effect of the FLSA on hours?

I refer to the difference $\Delta_{it} := h_{0it} - h_{1it}$ between h_0 and h_1 as unit it 's *treatment effect*. Recall that h_0 and h_1 are interpreted as potential outcomes, indicating what *would* have happened had the firm faced either of two counterfactual pay schedules instead of the kink. Δ_{it} thus represents the causal effect of a one-period 50% increase in worker i 's wage on their hours in week t : the difference between the hours that unit's firm would choose if the worker were paid at their straight-time rate versus at their overtime rate for all hours in that week. A unit's treatment effect can be contrasted with the "effect of the kink" quantity $h_{it} - h_{0it}$ introduced before, but importantly the two are related: by Eq. (2) the effect of the kink is $-\Delta_{it}$ for all units working overtime.

In the isoelastic model $\Delta_{it} = h_{0it} \cdot (1 - 1.5^\epsilon)$, representing a special case in which treatment effects are constant across all units after a log transformation of the outcome: $\ln h_{0it} - \ln h_{1it} = |\epsilon| \cdot \ln 1.5$. In general we can expect Δ_{it} to vary much more flexibly across units, and a reasonable parameter of interest becomes a summary statistic of Δ_{it} of some kind. In particular, Eq. (2) suggests that bunching is informative about the distribution of Δ_{it} among units "near" the kink. To see this, let $k = 40$ denote the location of the kink, and write the bunching probability as:

$$\mathcal{B} = P(h_{1it} \leq k \leq h_{0it}) = P(h_{0it} \in [k, k + \Delta_{it}]) = P(h_{1it} \in [k - \Delta_{it}, k]), \quad (6)$$

i.e. units bunch when their h_0 potential outcome lies to the right of the kink, but within that unit's individual treatment effect of it. Note that by Eq. (2) we can also write bunching in terms of the marginal distributions of h_0 and h_1 : $\mathcal{B} = F_1(k) - F_0(k)$, provided that they are continuously distributed and with F_0 and F_1 their cumulative distribution functions.

Parameter of interest: the buncher ATE

I focus my identification analysis on the average treatment effect among units who locate at exactly 40 hours, a parameter I call the "buncher ATE". In the overtime setting some care is needed in defining this parameter, allowing for the possibility that a mass of units would still work exactly 40 hours, even absent the FLSA. Let us indicate such "counterfactual bunchers" by an (unobserved) binary variable $K_{it}^* = 1$, and define the buncher ATE to be:

$$\Delta_k^* = \mathbb{E}[\Delta_{it} | h_{it} = k, K_{it}^* = 0],$$

That is, Δ_k^* is the average value of Δ_{it} among bunchers who bunch in response to the FLSA kink, and would not locate at 40 hours otherwise. In evaluating the FLSA, I will suppose that all counterfactual bunchers have a zero treatment effect, such that $h_{0it} = h_{1it} = k$. Since $\Delta_{it} = 0$ for these units by assumption, we can move between Δ_k^* and $\mathbb{E}[\Delta_{it}|h_{it} = k]$, provided the counterfactual bunching mass $p := P(K_{it}^* = 1)$ is known. In this section, I treat p as given, and present a strategy estimate it empirically in Section 5.1.

To simplify the discussion, suppose for now that there are no counterfactual bunchers, so that $\Delta_k^* = \mathbb{E}[\Delta_{it}|h_{it} = k]$. Our goal is to invert (6) in some way to learn about the buncher ATE from the observable bunching probability \mathcal{B} . In Figure 4, we've seen the intuition for this exercise in the context isoelastic model, in which there is only one dimension of heterogeneity and $h_{1it} = h_{0it} \cdot 1.5^\epsilon$. The key implication of the isoelastic model that aids in identification is *rank invariance* between h_0 and h_1 . Rank invariance (Chernozhukov and Hansen 2005) says that $F_0(h_{0it}) = F_1(h_{1it})$ for all units, i.e. increasing each unit's wage by 50% does not change their rank in the hours distribution (for example, a worker at the median of the h_0 distribution also has a median value of h_1). Rank invariance is satisfied by models in which there is perfect positive co-dependence between the potential outcomes (see left panel of Figure 5).

Rank invariance is useful because it allows us to translate statements about Δ_{it} into statements about the *marginal* distributions of h_{0it} and h_{1it} . In particular, under rank invariance the buncher ATE is equal to the quantile treatment effect $Q_0(u) - Q_1(u)$ averaged across all u between $F_0(k)$ and $F_1(k) = F_0(k) + \mathcal{B}$, with Q_d the quantile function of h_{dit} , i.e.:

$$\Delta_k^* = \frac{1}{\mathcal{B}} \int_{F_0(k)}^{F_1(k)} [Q_0(u) - Q_1(u)] du, \quad (7)$$

so long as $F_0(y)$ and $F_1(y)$ are continuous and strictly increasing. I focus on partial identification of the buncher ATE, for which it is sufficient to place point-wise bounds on the quantile functions $Q_0(u)$ and $Q_1(u)$ throughout the range $u \in [F_0(k), F_1(k)]$ as depicted in Figure 6.

While rank invariance already relaxes the isoelastic model used thus far in the literature, a still weaker assumption proves sufficient for Eq. (7) to hold:

Assumption RANK. *There exist fixed values Δ_0^* and Δ_1^* such that $h_{0it} \in [k, k + \Delta_{it}]$ iff $h_{0it} \in [k, k + \Delta_0^*]$, and $h_{1it} \in [k - \Delta_{it}, k]$ iff $h_{1it} \in [k - \Delta_1^*, k]$.*

Unlike (strict) rank invariance, Assumption RANK allows ranks to be reshuffled by treatment among bunchers and among the group of units that locates on each side of the kink.¹⁸ For example,

¹⁸When $p = 0$ Assumption RANK is equivalent to an instance of the *rank-similarity* assumption of Chernozhukov and Hansen (2005), in which the conditioning variable is which of the three cases of Equation (2) hold for the unit. Specifically, for both $d = 0$ and $d = 1$: $U_d| (h < k) \sim Unif[0, F_0(k)]$, $U_d| (h = k) \sim Unif[F_0(k), F_1(k)]$, and $U_d| (h > k) \sim Unif[F_1(k), 1]$.

suppose that a 50% increase in the wage of worker i would result in their hours being reduced from $h_{0it} = 50$ to $h_{1it} = 45$. If another worker j 's hours are instead reduced from $h_{0jt} = 48$ to $h_{1jt} = 46$ under a 50% wage increase, workers i and j will switch ranks, without violating RANK. Note that RANK is also compatible with the existence of counterfactual bunchers $p > 0$.

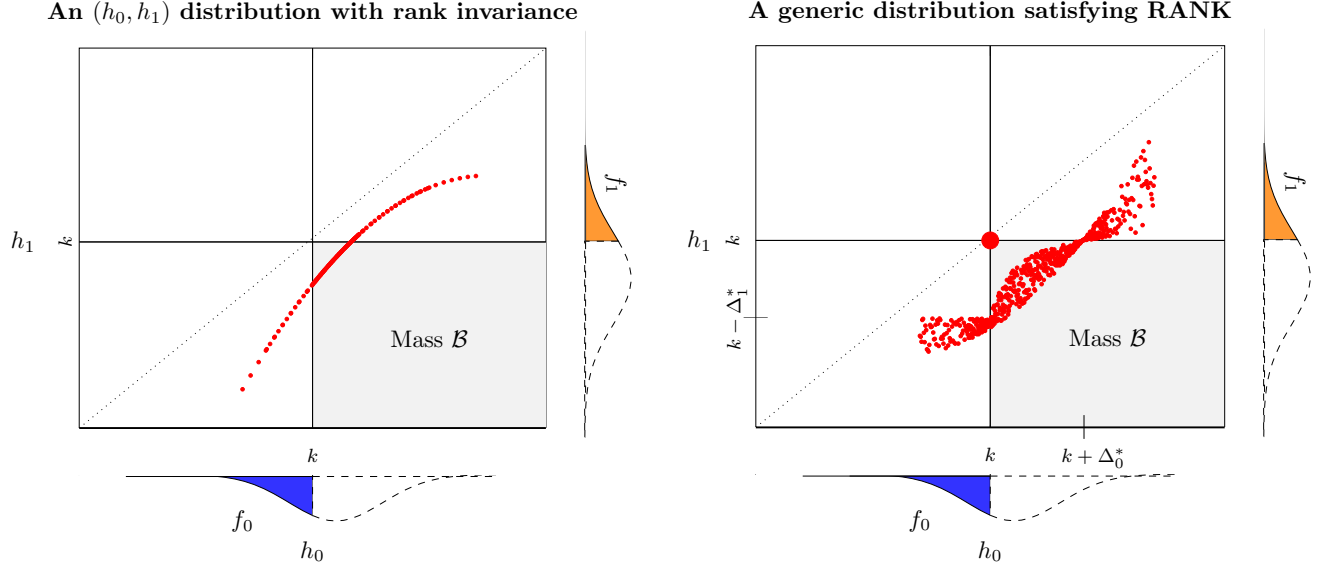


FIGURE 5: The joint distribution of (h_{0it}, h_{1it}) (in red), comparing an example satisfying rank invariance (left) to a case satisfying Assumption RANK (right). RANK allows the support of the joint distribution to “fan-out” from perfect co-dependence of h_0 and h_1 , except when either outcome is equal to k . The large dot in the right panel indicates a possible mass p of counterfactual bunchers. The observable data identifies the shaded portions of each outcome’s marginal distribution (depicted along the bottom and right edges), as well as the total mass \mathcal{B} in the (shaded) south-east quadrant.

The right panel of Figure 5 shows an example of a distribution satisfying RANK, which requires the support of (h_0, h_1) to narrow to a point as it crosses $h_0 = k$ or $h_1 = k$. When this is not perfectly satisfied, Appendix Figure A.3 demonstrates how the RHS of Equation (7) will then yield a lower bound on the true buncher ATE (and can still be interpreted as an averaged quantile treatment effect). Appendix Figure 11 generalizes RANK to case in which some workers choose their hours, resulting in mass appearing in the north-west quadrant of Figure 5.

4.3.1 Bounds on the buncher ATE via bi-log-concavity

Given Eq. (7), I obtain bounds on the buncher ATE by assuming that both h_0 and h_1 have *bi-log-concave* distributions. Bi-log-concavity is a nonparametric shape constraint that generalizes log-concavity, a property of many common parametric distributions:

Definition (BLC). A distribution function G is *bi-log-concave (BLC)* if both $\ln G$ and $\ln(1 - G)$ are concave functions.

If G is BLC then it admits a strictly positive density g that is itself differentiable with the locally bounded derivative: $\frac{-g(h)^2}{1-G(h)} \leq g'(h) \leq \frac{g(h)^2}{G(h)}$ (Dümbgen et al., 2017). Intuitively, this will rule out cases in which the density of h_0 or h_1 ever spikes or falls *too* quickly on the interior of its support, leading to non-identification of the type discussed in Section 4.2.¹⁹

The family of BLC distributions includes parametric distributions assumed by previous bunching design studies, such as those with uniform or linear densities Saez (2010), or those with polynomial densities as in Chetty et al. 2011 (provided they have real roots). All globally log-concave distributions are BLC.²⁰ Importantly, the BLC property is partially testable in the bunching design, since $F_0(y)$ is identified for all $h < k$ and $F_1(h)$ is identified for all $h > k$. Appendix Figure 7 shows that the observable portions of F_0 and F_1 indeed satisfy BLC. Identification requires us to believe that BLC also holds in the unobserved portions of F_0 and F_1 .

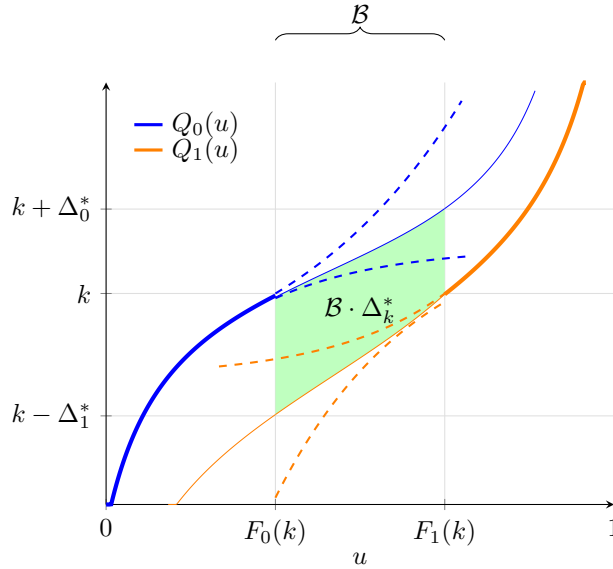


FIGURE 6: Extrapolating the quantile functions for h_0 and h_1 (blue and orange, respectively) to place bounds on the buncher ATE (case depicted has no counterfactual bunchers). The observed portions of each quantile function are depicted by thick curves, while the unobserved portions are indicated by thinner curves. The dashed curves represent upper and lower bounds for this unobserved portion coming from an assumption like bi-log-concavity (see text below). The buncher ATE is equal to the area shaded in green, divided by the bunching probability \mathcal{B} . The quantities Δ_0^* and Δ_1^* are defined in Assumption RANK below.

We are now ready to state the main identification result. Its logic is summarized by Figure 6: given the general choice model, RANK converts identification of the buncher ATE into a pair of extrapolation problems, each of which are approached by assuming bi-log-concavity of the corresponding

¹⁹Bertanha et al. (2020) propose partial identification in an isoelastic model by specifying a Lipschitz constant on the density of $\ln \eta_{it}$. This yields global rather than local bounds on g' .

²⁰However unlike log-concave densities, BLC distributions need not be unimodal (Dümbgen et al., 2017).

marginal potential outcome distribution.²¹ Let $F(h) := P(h_{it} \leq h)$ be the CDF of observed hours.

Theorem 1 (bi-log-concavity bounds on the buncher ATE). *Assume CHOICE, CONVEX, RANK and that h_{0it} and h_{1it} have bi-log-concave distributions conditional on $K_{it}^* = 0$. Then:*

1. $F(h)$, $F_0(h)$ and $F_1(h)$ are continuously differentiable for $h \neq k$. $F_0(k) = \lim_{h \uparrow k} F(h) + p$, $F_1(k) = F(k)$, $f_0(k) = \lim_{h \uparrow k} f(h)$ and $f_1(k) = \lim_{h \downarrow k} f(h)$, where if $p > 0$ we define the density of h_{dit} at $y = k$ to be $f_d(k) = \lim_{h \rightarrow k} f_d(h)$, for each $d \in \{0, 1\}$.
2. The buncher ATE Δ_k^* lies in the interval $[\Delta_k^L, \Delta_k^U]$, where:

$$\Delta_k^L := g(F_0(k) - p, f_0(k), \mathcal{B} - p) + g(1 - F_1(k), f_1(k), \mathcal{B} - p)$$

and

$$\Delta_k^U := -g(1 - F_0(k), f_0(k), p - \mathcal{B}) - g(F_1(k) - p, f_1(k), p - \mathcal{B})$$

with $g(a, b, x) = \frac{a}{bx} (a + x) \ln(1 + \frac{x}{a}) - \frac{a}{b}$, and the bounds are sharp.

Proof. See Appendix B. □

Combining Items 1 and 2 of Theorem 1, it follows that the sharp bounds Δ_k^L and Δ_k^U on the buncher ATE are identified, given the CDF of the data $F(h)$ and p .²² Inspection of the expressions appearing in Theorem 1 reveals that the bounds become wider the larger the net bunching probability $\mathcal{B} - p$. When $f_0(k) \approx f_1(k)$ and $p = 0$, the bounds will tend to be narrower when $F_0(k)$ is closer to $(1 - \mathcal{B})/2$, i.e. the kink is close to the median of the latent hours distribution. This helps explain why the estimated bounds in Section 5 turn out to be quite informative.

Comparison to existing results. The existing bunching design literature does contain a few identification results that circle the common intuition that bunching is informative about a local average responsiveness, when responsiveness to incentives varies by observational unit. For instance, Saez (2010) and Kleven (2016) consider a “small-kink” approximation that $\mathbb{E}[\Delta_{it} | h_{0it} = k] \approx \mathcal{B} / f_0(k)$, where the RHS is observable and referred to as the “excess mass” quantity in Chetty et al., 2011.²³

²¹It is worth noting that BLC of h_1 and h_0 implies bounds on the treatment effect $Q_1(u) - Q_0(u)$ at any quantile u . But these bounds widen quickly as one moves away from the kink. When $f_0(k) \approx f_1(k)$, the narrowest bounds for a single rank u are obtained for a “median” buncher roughly halfway between $F_0(k)$ and $F_1(k)$. However, averaging over a larger group is more useful for meaningful ex-post evaluation of the FLSA (Sec. 4.4), and reduces the sensitivity to departures from RANK (see Figure A.3). In the other extreme, one could drop RANK entirely and bound $\mathbb{E}[h_{0it} - h_{it}]$ directly via BLC of h_0 alone, but the bounds are very wide. The buncher ATE balances this tradeoff.

²²Since the bounds depend only on the density around k and the total amount mass to its left/right, point masses elsewhere in the distributions of h_0 and h_1 have no effect on the bounds provided that they are well-separated from k .

²³See Appendix A for a derivation in my generalized framework. The uniform density assumption is hard to justify except in the limit that the distribution of Δ_{it} concentrates around zero. Appendix Proposition 7 makes this claim precise, while connecting the approach from Saez (2010) and Kleven (2016) to results from Blomquist et al. (2015).

The result requires f_0 to be constant throughout the region $[k, k + \Delta_{it}]$ conditional on each value of Δ_{it} . This is most naturally justified by a kink that produces only tiny responses, an approximation that is likely to be quite poor in a context in which treatment corresponds to a 50% increase in the hourly cost of labor. Nevertheless, even in a small-kink setting, Theorem 1 offers a refinement to this result: a second-order approximation to $\ln(1 + \frac{x}{a})$ shows that when \mathcal{B} is small, the bounds converge and $\Delta_k^* \approx \frac{\mathcal{B}-p}{2f_0(k)} + \frac{\mathcal{B}-p}{2f_1(k)}$. Another existing result comes from Blomquist et al. (2015), who show in a generic labor supply model that bunching identifies a certain weighted average of compensated elasticities, if the density of choices at a kink is assumed to be linear across counterfactual tax rates. But such a parametric assumption is difficult to motivate, as these authors acknowledge.²⁴

4.4 Estimating policy relevant parameters

The buncher ATE yields the answer to a particular causal question, among a well-defined subgroup of the population. Namely: how would hours among workers bunched at 40 hours by the overtime rule be affected by a counterfactual change from linear pay at their straight-time wage to linear pay at their overtime rate? This section discusses how we may now use this quantity to both evaluate the overall ex-post effect of the FLSA on hours, as well as forecast the impacts of proposed changes to the FLSA. This requires some additional assumptions, which I continue to approach from a partial identification perspective.

4.4.1 From the buncher ATE to the ex-post hours effect of the FLSA

To consider the overall ex-post hours effect of the FLSA among covered workers, I proceed in two steps. I first relate the buncher ATE to the overall average effect of introducing the overtime kink, holding fixed the distributions of counterfactual hours h_{0it} and h_{1it} . Then, I allow straight-time wages to be affected by the FLSA, using the buncher ATE again to bound the additional effect of these wage changes on hours.

To motivate this strategy, let us first define the parameter of interest to be the difference in average weekly hours with and without the FLSA: $\theta := \mathbb{E}[h_{it}] - \mathbb{E}^*[h_{it}^*]$, where h_{it}^* indicates the hours unit it would work absent the FLSA, and the second expectation \mathbb{E}^* is over units corresponding to workers that would exist in the no-FLSA counterfactual and be eligible were it introduced.²⁵ Defining θ in this way allows us to remain agnostic as to whether the FLSA changes employment, and hence the population of workers it applies to. However, I assume that the hours among any workers who enter or exit employment due to the FLSA are not systematically different from those who

²⁴In particular, the data identifies the density at the kink for two particular tax rates only (in the tax application), so cannot provide evidence of such linearity.

²⁵Note that h_{it}^* in this section differs from the “anticipated” hours quantity h^* in Section 2.

would exist anyways, so that we may rewrite θ as $\theta = \mathbb{E}[h_{it} - h_{it}^*]$, averaging over individual-level causal effects in the population that does exist given the FLSA.

Next, decompose θ as:

$$\begin{aligned} \theta = \mathbb{E}[h_{it}(w_{it}, \mathbf{h}_{-i,t}) - h_{0it}(w_{it}^*, \mathbf{h}_{-i,t}^*)] &= \underbrace{\mathbb{E}[h_{it}(w_{it}, \mathbf{h}_{-i,t}) - h_{0it}(w_{it}, \mathbf{h}_{-i,t})]}_{\text{“effect of the kink”}} \\ &+ \underbrace{\mathbb{E}[h_{0it}(w_{it}, \mathbf{h}_{-i,t}) - h_{0it}(w_{it}^*, \mathbf{h}_{-i,t})]}_{\text{“wage effects”}} + \underbrace{\mathbb{E}[h_{0it}(w_{it}^*, \mathbf{h}_{-i,t}) - h_{0it}(w_{it}^*, \mathbf{h}_{-i,t}^*)]}_{\text{“interdependencies”}}, \end{aligned} \quad (8)$$

where the notation makes explicit the dependence of h and h_0 on the worker’s straight-time wage w_{it} , and possibly the hours \mathbf{h}_{-i} of other workers in their firm this week. In the notation of the last section: $h_{it} = h_{it}(w_{it}, \mathbf{h}_{-i,t})$, $h_{0it} = h_{0it}(w_{it}, \mathbf{h}_{-i,t})$ and $h_{1it} = h_{1it}(w_{it}, \mathbf{h}_{-i,t})$. I have used that $h_{it}^* = h_{0it}(w_{it}^*, \mathbf{h}_{-i,t}^*)$, since pay is linear in hours in the no-FLSA counterfactual.

The first term in Equation (8) reflects the “effect of the kink” quantity $h_{it} - h_{0it}$ examined in Section 4.2, and I view it as the first-order object of interest. The second term reflects that straight-time wages w_{it} may differ from those that workers would face without the FLSA, denoted by w_{it}^* . The third term is zero when firms’ choice of hours for its workers decomposes into separate optimization problems for each unit, as in the benchmark model from Section 4.2. More generally, it will capture any interdependencies in hours across units, for instance due to different workers’ hours being not linearly separable in production. In Online Appendix 3 I provide evidence that such effects do not play a large role in θ , and I thus treat this term as zero when estimating θ .²⁶

Turning first to the “effect of the kink” term, note that with straight-wages and the hours of other units fixed, the kink only has such direct effects on those units working at least $k = 40$ hours:

$$h_{it} - h_{0it} = \begin{cases} 0 & \text{if } h_{it} < k \\ k - h_{0it} & \text{if } h_{it} = k \\ -\Delta_{it} & \text{if } h_{it} > k \end{cases} \quad (9)$$

and thus $\mathbb{E}[h_{it} - h_{0it}] = \mathcal{B} \cdot \mathbb{E}[k - h_{0it} | h_{it} = k] - P(h_{it} > k) \mathbb{E}[\Delta_{it} | h_{it} > k]$. To identify this quantity we must extrapolate from the buncher ATE to obtain an estimate of $\mathbb{E}[\Delta_{it} | h_{it} > k]$, the average effect for units who work overtime. To do this, I assume that Δ_{it} of units working more than 40 hours are at least as large on average as those who work 40, but that the reduced-form *elasticity* of their response is no greater than that of the bunchers. Assuming a constant percentage change

²⁶In particular, I fail to find evidence of contemporaneous hours substitution in response to colleague sick pay, in an event study design. Another piece of evidence comes from obtaining similar “effect of the kink” estimates across small, medium and large firms, which suggests that a firm’s ability to reallocate hours between existing workers does not tend to drive their hours response to the FLSA. See Online Appendix 3.

between h_0 and h_1 over units would imply responses that grow in proportion to h_1 , eventually becoming implausibly large. On the other hand, it would be an underestimate to assume high-hours workers, say at 60 hours, have the same effect in levels $h_0 - h_1$ as those closer to 40. Finally, to put bounds on the average effect of the kink among bunchers $\mathcal{B} \cdot \mathbb{E}[k - h_{0it} | h_{it} = k]$, I use bi-log-concavity of h_0 . Details are provided in Online Appendix 6.11.

The “wage effects” term in Equation (8) arises because the straight-time wages observed in the data may reflect some adjustment to the FLSA, as we would expect on the basis of the conceptual framework in Section 2. While the “effect of the kink” term is expected to be negative, this second term will be positive if the FLSA causes a reduction in the straight-time wages set at hiring on the basis of expected hours. However, both terms ultimately depend on the same thing: responsiveness of hours to the cost of an hour of work. I thus use the buncher ATE to compute an approximate upper bound on wage effects by assuming that all straight-time wages are adjusted according to Equation (1) and that the hours response is iso-elastic in wages, with anticipated hours approximated by h_{it} . Online Appendix 6.11 provides a visual depiction of these definitions. A lower bound on the “wage effects” term, on the other hand, is zero. In practice, the estimated size of the wage effect $\mathbb{E}[h_{0it} - h_{0it}^*]$ is appreciable but still small relative to $\mathbb{E}[h_{it} - h_{0it}]$ (cf. Appendix Table 11).

4.4.2 Forecasting the effects of policy changes

Apart from ex-post evaluation of the overtime rule, policymakers may also be interested in predicting what would happen if the parameters of overtime regulation were modified. Reforms that have been discussed in the U.S. include decreasing “standard hours” k at which overtime pay begins from 40 hours to 35 hours,²⁷ or increasing the overtime premium from time-and-a-half to “double-time” (Brown and Hamermesh, 2019). This section builds upon Sections 4.1 and 4.3 to show that the bunching-design model is also informative about the impact of such reforms on hours.

Let us begin by considering changes to standard hours k , for now holding the distributions of h_0 and h_1 fixed across the policy change. Inspection of Equation (2) reveals that as the kink is moved upwards, say from $k = 40$ hours to $k' = 44$ hours, some workers who were previously bunching at k now work h_{0it} hours: namely those for whom $h_{0it} \in [k, k']$. By the same token, some individuals with values of $h_{1it} \in [k, k']$ now bunch at k' . Some individuals who were bunching at k now bunch at k' —namely those for whom $h_{1it} \leq k$ and $h_{0it} \geq k'$. In the case of a reduction in overtime hours, say to $k' = 35$ this logic is reversed. Figure 8 depicts both cases. I assume that the mass of counterfactual bunchers p remains at $k = 40$ after the shift.²⁸

²⁷Several countries have implemented changes to standard hours; Brown and Hamermesh (2019) provides a review.

²⁸It is conceivable that some or all counterfactual bunchers locate at 40 because it is the FLSA threshold, while still being non-responsive to the incentives introduced there by the kink. In this case, we might imagine that they would all coordinate on k' after the change. The effects here should thus be seen as short-run effects before that occurs.

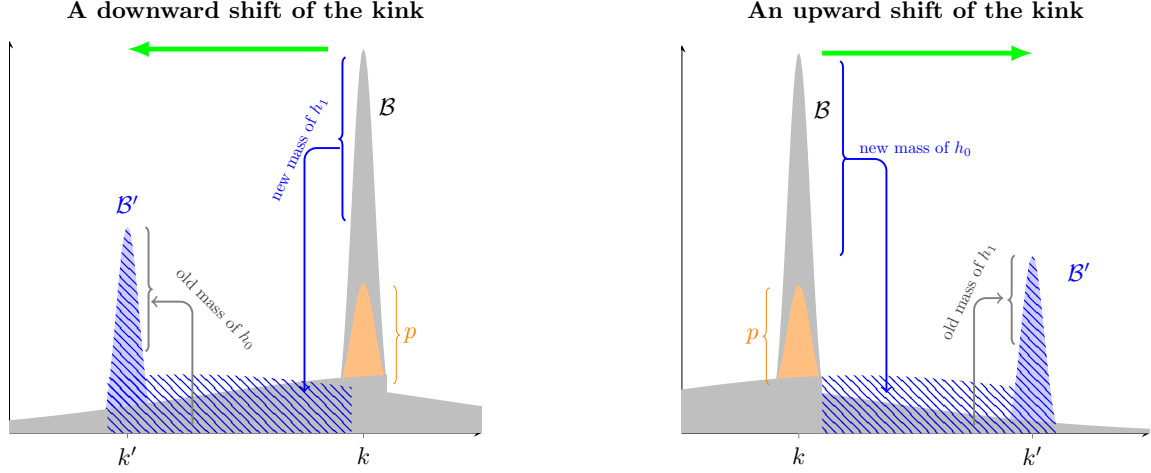


FIGURE 7: The left panel depicts a shift of the kink point downwards from k to k' , while right panel depicts a shift of the kink point upwards. See text for details.

Quantitatively assessing a change to double-time pay requires us to move beyond the two counterfactual choices h_{0it} and h_{1it} : hours that would be worked under straight-wages or under time-and-a-half pay. Let $h_{it}(\rho)$ be the hours that it would work if their employer faced a linear pay schedule at rate $\rho \cdot w_{it}$ (with w_{it} and hours of other units fixed at their realized levels). In this notation, $h_{0it} = h_{it}(1)$ and $h_{0it} = h_{it}(1.5)$. Now consider a new overtime policy in which a premium pay factor of ρ_1 is required for hours in excess of k , e.g. $\rho_1 = 2$ for a “double-time” policy. Let $h_{it}^{[k, \rho_1]}$ denote realized hours for unit it under this overtime policy, and let $\mathcal{B}^{[k, \rho_1]} := P(h_{it}^{[k, \rho_1]} = k)$ the observable bunching that would occur.

Theorem 2 allows one to discuss the effects of small changes to k or ρ_1 on hours. Some results make use of an explicit assumption that firm preferences are quasi-linear with respect to costs:

Assumption SEPARABLE. $\pi_{it}(z, \mathbf{x})$ is additively separable and linear in z for all units it .

I continue to assume that counterfactual bunchers $K_{it}^* = 1$ stay at $k^* := 40$, regardless of ρ and k . Let $p(k) = p \cdot \mathbb{1}(k = k^*)$ denote the possible mass of counterfactual bunchers as a function of k .

Theorem 2 (marginal comparative statics in the bunching design). *Under Assumptions CHOICE, CONVEX, SEPARABLE and SMOOTH:*

1. $\partial_k \left\{ \mathcal{B}^{[k, \rho_1]} - p(k) \right\} = f_1(k) - f_0(k)$
2. $\partial_k \mathbb{E}[h_{it}^{[k, \rho_1]}] = \mathcal{B}^{[k, \rho_1]} - p(k)$
3. $\partial_{\rho_1} \mathcal{B}^{[k, \rho_1]} = -k f_{\rho_1}(k) \mathbb{E} \left[\frac{dh_{it}(\rho_1)}{d\rho} \middle| h_{it}(\rho_1) = k \right]$
4. $\partial_{\rho_1} \mathbb{E}[h_{it}^{[k, \rho_1]}] = - \int_k^\infty f_{\rho_1}(h) \mathbb{E} \left[\frac{dh_{it}(\rho_1)}{d\rho} \middle| h_{it}(\rho_1) = h \right] dh$

Proof. See Appendix A. □

Assumption SMOOTH is a set of regularity conditions which imply that $h_{it}(\rho)$ admits a density $f_\rho(h)$ for all ρ ; see Appendix A for details. The above also uses a version of Assumption CHOICE given in Appendix A, which applies to all ρ rather than just ρ_0 and ρ_1 .

Beginning from the actual FLSA policy of $k = 40, \rho_1 = 1.5$, the RHS of Items 1 and 2 are point identified from the data, provided that p is known. Item 1 says that if the location of the kink is changed marginally, the kink-induced bunching probability will change according to the difference between the densities of h_{1i} and h_{0i} at k^* , which are in turn equal to the left and right limits of the observed density $f(h)$ at the kink. This result is intuitive: given continuity of each potential outcome's density, a small increase in k will result in a mass proportional to $f_1(k)$ being “swept in” to the mass point at the kink, while a mass proportional to $f_0(k)$ is left behind. Item 2 aggregates this change in bunching with the changes to non-bunchers' hours as k is increased: the combined effect turns out to be to simply transport the mass of inframarginal bunchers to the new value of k .²⁹ Making use of Theorem 2 for a discrete policy change like reducing standard hours to 35 requires integrating across the actual range of hypothesized policy variation. We lose point identification, but I use bi-log-concavity of the marginal distributions of h_0 and h_1 to retain bounds.

Now consider the effect of moving from time-and-a-half to double time on average hours worked, in light of Item 4. This scenario, similar to the effect of the kink term in Eq. (8), requires making assumptions about the response of individuals who may locate far above the kink, and for whom the buncher ATE is less directly informative. First, we integrate Item 4 over ρ to obtain an expression for the average effect on hours from a move to double-time, which can be written in terms of local average elasticities of response:

$$\mathbb{E}[h_{it}^{[k, \rho_1]} - h_{it}^{[k, \bar{\rho}_1]}] = \int_{\rho_1}^{\bar{\rho}_1} d \ln \rho \int_k^\infty f_\rho(h) \cdot h \cdot \mathbb{E} \left[\frac{d \ln h_{it}(\rho)}{d \ln \rho} \middle| h_{it}(\rho) = h \right] dh$$

Recall that in the isoelastic model the elasticity quantity $\frac{d \ln h_{it}(\rho)}{d \ln \rho} = \frac{dh_{it}(\rho)}{d\rho} \frac{\rho}{h_{it}(\rho)}$ is constant across ρ and across units, and it is partially identified under BLC. I argue that just as an isoelastic response is likely to overstate responsiveness at large values of hours, it is likely to *understate* responsiveness to larger values of ρ , thus yielding a lower bound on the effect of moving to double-time. For an upper bound on the magnitude of the effect, I assume rather than in levels $\mathbb{E}[h_{it}(\rho_1) - h_{it}(\bar{\rho}_1) | h_{1it} > k]$ is at least as large as $\mathbb{E}[h_{0it} - h_{1it} | h_{1it} > k]$, and that the increase in bunching from a change of ρ_1 to $\bar{\rho}_1$ is as large as the increase from ρ_0 to ρ_1 . Additional details

²⁹Intuitively, “marginal” bunchers who would choose exactly k under one of the two cost functions B_0 or B_1 cease to “bunch” as k increases, but in the limit of a small change they also do not change their realized h . Moore (2021) gives a closely-related result, derived independently of this work, that shows that bunching is a sufficient statistic for the effect of a marginal change in k on revenue in the context of a tax kink.

are provided in Online Appendix 6.11.

5 Implementation and Results

This section implements the empirical strategy described in the last section with the sample of administrative payroll data described in Section 3.

5.1 Identifying counterfactual bunching at 40 hours

To deliver final estimates of the effect of the FLSA overtime rule on hours, it is necessary to first return to an issue raised in the introduction and alluded to in Section 4: that there are other reasons to expect bunching at 40 hours, in addition to being the location of the FLSA kink. For one, 40 may reflect a status-quo choice, being chosen even when it is not exactly profit maximizing for the firm. This effect could be amplified by firms synchronizing the schedules of different workers, requiring some common number of hours per week to coordinate around. Finally, for any salaried workers who were not successfully removed from the sample, firms might record the number of hours in a pay period as 40 even as actual hours worked vary.

In terms of the empirical strategy from Section A.2, all of these alternative explanations manifest in the same way: a point mass p at 40 in the distribution of hours that would occur even if workers' pay did not feature a kink at 40. In the notation introduced in Section 4.3, these “counterfactual bunchers” are demarcated by $K_{it}^* = 1$. Let us refer to the $K_{it}^* = 0$ individuals who also locate at the kink as “active bunchers”. The mass of active bunchers is $\mathcal{B} - p$. Theorem 1 shows that we can still partially identify the buncher ATE in the presence of counterfactual bunchers, so long as we know what portion of the total bunchers are active and how many are counterfactual.

I leverage two strategies to provide plausible estimates for the mass of counterfactual bunchers p . My preferred estimate uses of the fact that when an employee is paid for hours that are not actually worked—including sick time, paid time off (PTO) and holidays—these hours do not contribute to the 40 hour overtime threshold of the FLSA. For example, if a worker applies PTO to miss a six hour shift, then they are not required to be paid overtime until they reach 46 total paid hours in that week. Thus while the kink remains at 40 hours *worked*, non-work hours like PTO shift the location of the kink in hours of *pay*.

The identifying assumption that I rely on is that individuals who still work 40 hours a week, even when they have non-work hours (and are hence paid for more than 40), are all active bunchers: they would not locate at forty hours in the counterfactuals h_{0it} and h_{1it} . This reflects the idea that additional explanations for bunching at 40 hours operate at the level of hours paid, rather than hours worked. Letting n_{it} indicate non-work hours of pay for worker i in week t , I make two assumptions:

1. $P(h_{it} = 40 | n_{it} > 0) = P(h_{it} = 40 \text{ and } K_{it}^* = 0 | n_{it} > 0)$
2. $P(h_{it} = 40 \text{ and } K_{it}^* = 0 | n_{it} > 0) = P(h_{it} = 40 \text{ and } K_{it}^* = 0 | n_{it} = 0)$

The first item allows me to identify the mass of active bunchers in the $n_{it} > 0$ conditional distribution of hours. The second item says that this conditional mass is representative of the unconditional mass of active bunchers. To increase the plausibility of this assumption, I focus on η as paid time off because it is generally planned in advance, yet has somewhat idiosyncratic timing. By contrast, sick pay is often unanticipated so the firm may not be able to re-optimize total hours within the week in which a worker calls in sick. Holiday pay is known in advance, but holidays are unlikely to be representative in terms of other factors important for hours determination (e.g. product demand).

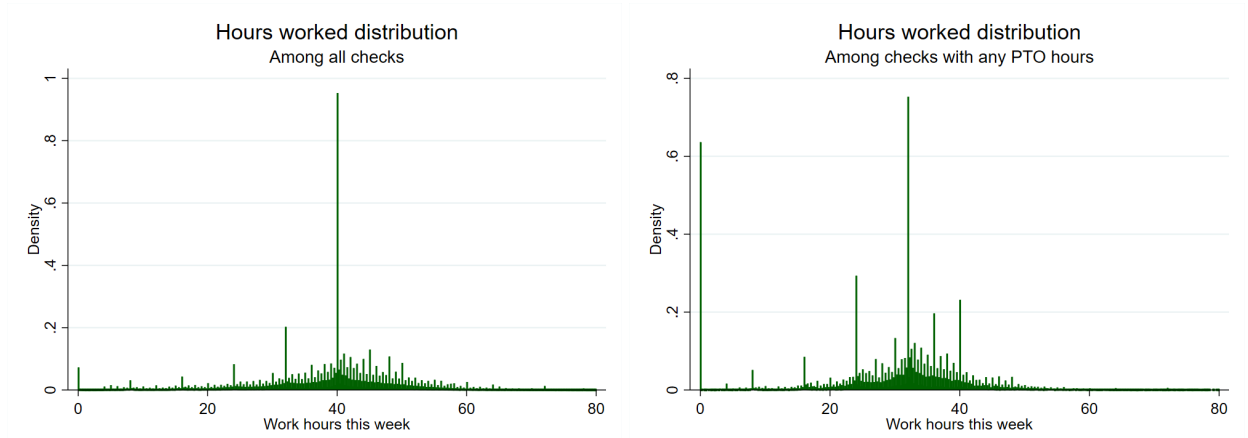


FIGURE 8: The right panel shows a histogram of hours worked when paid time off hours are positive ($\eta_{it} > 0$). The left panel shows the unconditional distribution. While $\mathcal{B} \approx 11.6\%$, $P(h_{it} = 40 | n_{it} > 0) \approx 2.7\%$.

Together, the two assumptions above imply that $p = P(K_{it}^* = 1 \text{ and } h_{it} = 40)$ is identified as $\mathcal{B} - P(h_{it} = 40 | \eta_{it} > 0)$. Figure 8 shows the conditional distribution of hours paid for work when the paycheck contains a positive number of PTO hours ($n_{it} > 0$). The figure reveals that when moving from the unconditional (left panel) to positive-PTO conditional (right panel) distribution, most of the point mass at 40 hours moves away, largely concentrating now at 32 hours (corresponding to the PTO covering a single eight hour shift). Of the total bunching of $\mathcal{B} \approx 11.6\%$ in the unconditional distribution, I estimate that only about $P(h_{it} = 40 | n_{it} > 0) \approx 2.7\%$ are active bunchers, leaving $p \approx 8.9\%$. Thus roughly three quarters of the individuals at 40 hours are counterfactual rather than active bunchers.

As a secondary strategy, I estimate an upper bound for p by using the assumption that the potential outcomes of counterfactual bunchers are relatively “sticky” over time. If the hours of counterfactual bunchers are at forty for behavioral or administrative reasons, it may be reasonable to assume that these external considerations are fairly static, preventing latent hours h_{0it} from changing much between adjacent weeks. In particular, assume that in a given week t nearly all of

the counterfactual bunchers are also non-movers from week $t - 1$. Then:

$$p = P(h_{0it} = 40) \approx P(h_{0it} = h_{0it-1} = 40) \leq P(h_{it} = h_{i,t-1} = 40),$$

where the inequality follows from $(h_{0it} = 40) \implies (h_{it} = 40)$ by Lemma 1. The probability $P(h_{it} = h_{i,t-1} = 40)$ can be directly estimated from the data, yielding $p \leq 6\%$.

5.2 Estimation and inference

Given Theorem 1 and a value of p , computing bounds on the buncher ATE requires estimates of the right and left limits of the CDF and density of hours at the kink. I use the local polynomial density estimator of Cattaneo, Jansson and Ma (2020) (CJM), which is well-suited to estimating a CDF and its derivatives at boundary points. For instance, a local-linear CJM estimator provides a smoothed estimate of the left limit of the CDF and density at k as:

$$(\hat{F}_-(k), \hat{f}_-(k)) = \underset{(b_1, b_2)}{\operatorname{argmin}} \sum_{it: h_{it} < k} (F_n(h_{it}) - b_1 - b_2 h_{it})^2 \cdot K\left(\frac{h_{it} - k}{\alpha}\right) \quad (10)$$

where $F_n(y) = \frac{1}{n} \sum_{it} \mathbb{1}(h_{it} \leq y)$ is the empirical CDF of a sample of size n , $K(\cdot)$ is a kernel function, and α is a bandwidth. The right limits $F_+(k)$ and $f_+(k)$ are estimated analogously using observations for which $h_{it} > k$. I use a triangular kernel, and choose h as follows: first, I use CJM's mean-squared error minimizing bandwidth selector to produce a bandwidth choice using the data on either side of $k = 40$ (for the left and right limits, respectively). I then average the two bandwidths, and use this as the bandwidth in the final calculation of both the right and left limits. In the full sample, the bandwidth chosen by this procedure is about 1.7 hours, and is somewhat larger for subsamples that condition on a single industry.

To construct confidence intervals for parameters that are partially identified (e.g. the buncher ATE), I use adaptive critical values proposed by Imbens and Manski (2004) and Stoye (2009) that are valid for the underlying parameter. In each case, estimators of bounds or point identified quantities are functions of inputs that are \sqrt{n} -asymptotically normal. To easily incorporate sampling uncertainty in $\hat{F}_-(k)$, $\hat{f}_-(k)$, $\hat{F}_+(k)$, $\hat{f}_+(k)$ and \hat{p} , I estimate the variances by a cluster nonparametric bootstrap that resamples at the firm level. This allows arbitrary autocorrelation in hours across pay periods for a single worker, and between workers within a firm. All standard errors use 500 bootstrap replications.

5.3 Results of the bunching estimator: the buncher ATE

Table 3 reports treatment effect estimates based on Theorem 1, when p is either assumed zero or estimated by one of the two methods described in Section 5.1. These estimates use a sample that pools across workers in all industries. The first row reports the corresponding estimate of the net bunching probability $\mathcal{B} - p$, while the second row reports the bounds on the buncher ATE $\mathbb{E}[h_{0it} - h_{1it}|h_{it} = k, K_{it}^* = 0]$. Within a fixed estimate of p , the bounds on the buncher ATE based on bi-log-concavity are quite informative: the upper and lower bounds are always close to each other and precisely estimated. Online Appendix 1 reports estimates based on alternative shape constraints and assumptions about effect heterogeneity, which deliver similar results.³⁰

	$p=0$	p from non-changers	p from PTO
Net bunching:	0.116 [0.112, 0.120]	0.057 [0.055, 0.058]	0.027 [0.024, 0.030]
Buncher ATE	[2.614, 3.054] [2.493, 3.205]	[1.324, 1.435] [1.264, 1.501]	[0.640, 0.666] [0.574, 0.736]
Num observations	630217	630217	630217
Num clusters	566	566	566

TABLE 3: Estimates of net bunching $\mathcal{B} - p$ and the buncher ATE: $\Delta_k^* = \mathbb{E}[h_{0it} - h_{1it}|h_{it} = k, K_{it}^* = 0]$, across various strategies to estimate counterfactual bunching $p = P(K_{it}^* = 1)$. Unit of analysis is a paycheck, and 95% bootstrap confidence intervals (in gray) are clustered by firm.

The PTO-based estimate of p provides the most conservative treatment effect estimates, attributing roughly one quarter of the observed bunching to active rather than counterfactual bunchers. Nevertheless, this estimate still yields a highly statistically significant buncher ATE of about 2/3 of an hour, or 40 minutes. This estimate can be interpreted as follows: consider the group of workers that are in fact working 40 hours in a given pay period and are not counterfactual bunchers. This group would work on average about 40 minutes more that week if they were paid their straight-time wage for all hours, compared with a counterfactual in which they are paid their overtime rate for all hours. If we instead attribute all of the observed bunching mass to active bunchers ($p = 0$),

³⁰In particular, I present a point estimate based on Appendix Proposition 3, which assumes that treatment effects are constant and that the density is linear in the missing region, as well as results under a weaker assumption that the density is monotonic in the missing region. Monotonicity is not likely to hold in the overtime context, but the bounds based on monotonicity do not deliver vastly different results. See Appendix Tables 9 and 10, which applies the same assumptions to the distribution of log hours rather than hours.

then this buncher ATE parameter is estimated to be at least 2.6 hours. In Online Appendix 1 I report estimates of the buncher ATE for each of the largest industries in the sample, and also plot estimates directly as a function of the assumed mass p of counterfactual bunchers at 40 hours.

5.4 Estimates of policy effects

I now use estimates of the buncher ATE and the results of Section 4.4 to estimate the overall causal effect of the FLSA overtime rule, and simulate changes based on modifying standard hours or the premium pay factor. Table 4 first reports an estimate of the buncher ATE expressed as a reduced-form hours demand elasticity,³¹ which I use as an input in these calculations. The next two rows report bounds on $\mathbb{E}[h_{it} - h_{it}^*]$ and $\mathbb{E}[h_{it} - h_{it}^* | h_{1it} \geq 40, K_{it}^* = 0]$, respectively. The second row is the overall ex-post effect of the FLSA on hours, averaged over workers and pay periods, and the third row conditions on paychecks reporting at least 40 hours (omitting counterfactual bunchers). The final row reports an estimate of the effect of moving to double-time pay. I provide details of the calculations in Online Appendix 6.11.

Taking the PTO-based estimate of p as yielding a lower bound on treatment effects, the estimates suggest that FLSA eligible workers work at least 1/5 of an hour less in any given week than they would absent overtime regulation: about one third the magnitude of the buncher ATE in levels. When I focus on those eligible workers that are directly affected in a given week, the figure is about twice as high: roughly 30 minutes. I estimate that a move to double-time pay would introduce a further reduction that may be comparable to the existing overall ex-post effect, but with substantially wider bounds. These estimates include the effects of possible adjustments to straight-time wages, which tend to attenuate the effects of the policy change. Appendix Table 11 replicates Table 4 neglecting these wage adjustments, which might be viewed as a short-run response to the FLSA before wages have time to adjust.

Figure 9 breaks down estimates of the ex-post effect of the overtime rule by major industries, revealing considerable heterogeneity between them. The estimates suggest that Real Estate & Rental and Leasing as well as Wholesale Trade see the highest average reduction in hours. The least-affected industries are Health Care and Social Assistance and Professional Scientific and Technical, with the average worker working just about 6 minutes less per week. Appendix Figure 6 compares the hours distribution for Real Estate & Rental and Leasing with the distribution for Professional Scientific and Technical, showing that the difference in their effects is explained both by a larger value of $\mathcal{B} - p$ and a lower density of hours close to the kink for Real Estate & Rental and Leasing. Appendix Table 5 reports the numerical estimates and confidence intervals by industry. Online Ap-

³¹ This is $\hat{\Delta}_k^* / (40 \ln(1.5))$ where $\hat{\Delta}_k$ is the estimate of the buncher ATE presented in Table 3, which is numerically equivalent to the elasticity implied by the buncher ATE in logs $\mathbb{E}[\ln h_{0it} - \ln h_{1it} | h_{it} = k, K_{it}^* = 0] / (\ln 1.5)$ estimated under assumption that $\ln h_0$ and $\ln h_1$ are BLC.

	$p=0$	p from non-changers	p from PTO
Buncher ATE as elasticity	[-0.188,-0.161] [-0.198,-0.154]	[-0.088,-0.082] [-0.093,-0.078]	[-0.041,-0.039] [-0.045,-0.035]
Average effect of FLSA on hours	[-1.466, -1.026] [-1.535, -0.977]	[-0.727, -0.486] [-0.762, -0.463]	[-0.347, -0.227] [-0.384, -0.203]
Avg. effect among directly affected	[-2.620, -1.833] [-2.733, -1.750]	[-1.453, -0.972] [-1.518, -0.929]	[-0.738, -0.483] [-0.812, -0.434]
Double-time, average effect on hours	[-2.604, -0.569] [-2.707, -0.547]	[-1.239, -0.314] [-1.285, -0.300]	[-0.580, -0.159] [-0.638, -0.143]

TABLE 4: Estimates of the buncher ATE expressed as an elasticity, the average ex-post effect of the FLSA $\mathbb{E}[h_{it} - h_{it}^*]$,³¹ the effect among directly affected units $\mathbb{E}[h_{it} - h_{it}^* | h_{it} \geq k]$ and predicted effects of a change to double-time. 95% bootstrap confidence intervals in gray, clustered by firm.

pendix 1 reports estimates broken down by gender, finding that the FLSA has considerably higher effects on the hours of men compared with women.

Appendix Figure 4 looks at the effect of changing the threshold for overtime hours k from 40 to alternative values k' . The left panel reports estimates of the identified bounds on $\mathcal{B}^{[k', \rho_1]}$ as well as point-wise 95% confidence intervals (gray) across values of k' between 35 and 45, for each of the three approaches to estimating p . In all cases, the upper bound on bunching approaches zero as k' is moved farther from 40. This is sensible if the h_0 and h_1 distributions are roughly unimodal with modes around 40: straddling of potential outcomes becomes less and less likely as one moves away from where most of the mass is. Appendix 5 shows these bounds as k' ranges all the way from 0 to 80, for the $p = 0$ case. These estimates should be viewed as short-run responses, as they do not account for adjustment to straight-time wages.

When p is estimated using PTO or non-changers between periods, we see that the upper bound of the identified set for $\mathcal{B}^{[k', \rho_1]}$ in fact reaches zero quite quickly in k' . Moving standard hours to 35 is predicted to completely eliminate bunching due to the overtime kink in the short run, before any adjustment to latent hours (e.g. through changes to straight-time wages). The right panel of Appendix Figure 4 shows estimates for the average effect on hours of changing standard hours, inclusive of wage effects (see Online Appendix 6.11 for details). Increases to standard hours cause an increase in hours per worker, as overtime policy becomes less stringent, and reductions to standard hours reduce hours.³² The size of these effects is not precisely estimated for changes

³² The magnitudes are consistent with estimates by Costa (2000), that hours fell by 0.2-0.4 on average during the phased introduction of the FLSA in which standard hours declined by 2 hours in 1939 and 1940.

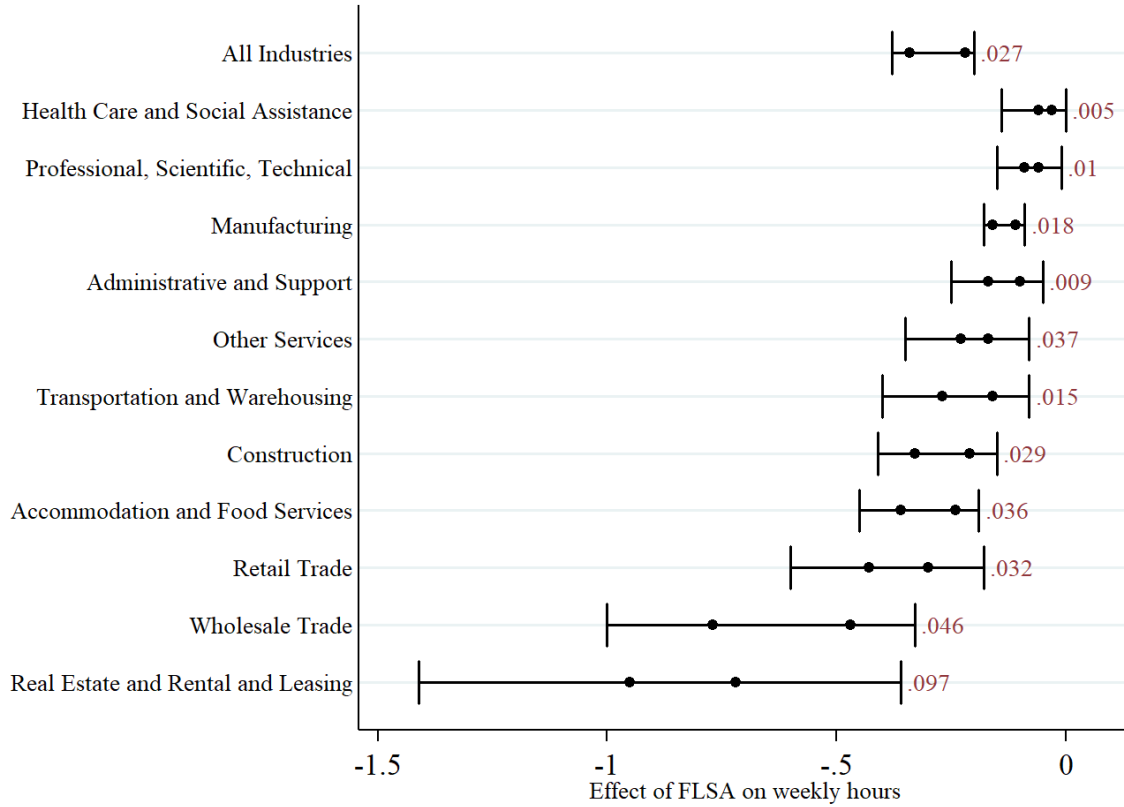


FIGURE 9: 95% confidence intervals for the effect of the FLSA on hours by industry, using PTO-based estimates of p for each. Dots are point estimates of the upper and lower bounds. The number to the right of each range is the point estimate of the net bunching $B - p$ for that industry.

larger than a couple of hours, however the range of statistically significant effects depends on p . Even for the preferred estimate of p from PTO, increasing the overtime threshold as high as 43 hours is estimated to increase average working hours by an amount distinguishable from zero.

6 Implications of the estimates for overtime policy

The estimates from the preceding section suggest that FLSA regulation indeed has real effects on hours worked, in line with labor demand theory when wages do not fully adjust to absorb the added cost of overtime hours. When averaged over affected workers and across pay periods, I find that hourly workers in my sample work at least 30 minutes less per week than they would without the overtime rule. This lower bound is broadly comparable to the few causal estimates that exist in the literature, including Hamermesh and Trejo (2000) who assess the effects of expanding California's daily overtime rule to cover men in 1980, and Brown and Hamermesh (2019) who use the erosion

of the real value of FLSA exemption thresholds over the last several decades.³³ By contrast, my estimates carry the strengths of an approach to identification that does not require focusing on the sub-population affected by a natural experiment, and use much more recent data.

These estimates speak to the substitutability of hours of labor between workers. The primary justifications for overtime regulation have been to reduce excessive workweeks, while encouraging hours to be distributed over more workers (Ehrenberg and Schumann, 1982). How well this—and related policies such as work-sharing programs—play out in practice hinges on how easily an hour of work can be moved from one worker to another or across time, from the perspective of the firm. The results of this paper find hours demand to be relatively inelastic: hours cannot be easily so reallocated between workers or weeks. This suggests that ongoing efforts to expand coverage of the FLSA overtime rule (by increasing the earnings threshold at which some salaried workers are exempt) may have limited scope to drastically affect the hours of U.S. workers.

Nevertheless, the overall impact of the FLSA overtime rule on workers could be substantial. The data suggest that at least about 3% and as many as about 12% of workers’ hours are adjusted to the threshold introduced by the policy, indicating that it may have distortionary impacts for a significant portion of the labor force. The policy may also have important effects on unemployment. While a full assessment of the employment effects of the FLSA overtime rule is beyond the scope of this paper, my estimates of the hours effect can be used to build a back-of-the-envelope calculation. Following Hamermesh (1996), I assume a value for the rate at which firms substitute labor for capital to obtain a “best-guess” estimate that the FLSA overtime rule creates about 700,000 jobs (see Online Appendix 1.6 for details). To get an overall upper bound on the size of employment effects, I attribute all of the bunching at 40 to the FLSA and assume that the total number of worker-hours is not reduced by the FLSA. By this estimate the FLSA increases employment by at most 3 million jobs, or 3% among covered workers. A reasonable range of parameter values rules out negative overall employment effects from the FLSA.

7 Conclusion

This paper has provided a reinterpretation of the bunching-design method in the language of treatment effects, showing that the basic identifying power of the method is robust to a variety of underlying structural choice models. Across such modeling choices, the parameter of interest remains a reduced-form local average treatment effect between two appropriately-defined counterfactual choices, which can be partially identified via nonparametric assumptions on the distributions of

³³ Hamermesh and Trejo (2000) and Brown and Hamermesh (2019) report estimates of -0.5 and -0.18 for the elasticity of overtime hours with respect to the overtime rate. My preferred estimate of -0.04 for the buncher ATE as an elasticity is the elasticity of *total* hours, including the first 40. An elasticity of overtime hours can be computed from this using the ratio of mean hours to mean overtime hours in the sample, resulting in an estimate of roughly -0.45 .

counterfactuals. This provides conditions under which the bunching design can be useful to answer program evaluation questions in a broad variety of contexts, particularly beyond those in which the researcher is prepared to posit a parametric model of decision-makers' preferences.

By leveraging these insights with a new payroll dataset recording exact weekly hours paid at the individual level, I estimate that U.S. workers subject to the Fair Labor Standard Act work shorter hours due to its overtime provision, which may lead to positive employment effects. Given the large amount of within-worker variation in hours observed in the data, the modest size of the FLSA effects estimated in this paper suggest that firms do face significant incentives to maintain longer working hours, countervailing against the ones introduced by policies intended to reduce them.

References

- BARKUME, A. (2010). "The Structure of Labor Costs with Overtime Work in U.S. Jobs". *Industrial and Labor Relations Review* 64 (1).
- BERTANHA, M., MCCALLUM, A. H. and SEEGER, N. (2020). "Better Bunching , Nicer Notching". *SSRN Working Paper*.
- BEST, M. C., BROCKMEYER, A., KLEVEN, H. J., SPINNEWIJN, J. and WASEEM, M. (2015). "Production vs Revenue Efficiency With Limited Tax Capacity: Theory and Evidence From Pakistan". *Journal of Political Economy* 123 (6), p. 48.
- BISHOW, J. L. (2009). "A Look at Supplemental Pay: Overtime Pay, Bonuses, and Shift Differentials". *Monthly Labor Review*. Publisher: Bureau of Labor Statistics, U.S. Department of Labor.
- BLOMQUIST, S., KUMAR, A., LIANG, C.-Y. and NEWHEY, W. (2021). "On Bunching and Identification of the Taxable Income Elasticity". *Journal of Political Economy* 129 (8).
- BLOMQUIST, S., KUMAR, A., LIANG, C.-Y. and NEWHEY, W. K. (2015). "Individual heterogeneity, nonlinear budget sets and taxable income". *The Institute for Fiscal Studies Working Paper* CWP21/15.
- BLOMQUIST, S. and NEWHEY, W. (2017). "The Bunching Estimator Cannot Identify the Taxable Income Elasticity". *The Institute for Fiscal Studies Working Paper* CWP40/17.
- BRECHLING, F. P. R. (1965). "The Relationship Between Output and Employment in British Manufacturing Industries". *The Review of Economic Studies* 32 (3), p. 187.
- BROWN, C. and HAMERMESH, D. S. (2019). "Wages and Hours Laws: what do we know? what can be done?" *The Russell Sage Foundation Journal of the Social Sciences* 5 (5), pp. 68–87.

- BURDETT, K. and MORTENSEN, D. T. (1998). “Wage Differentials, Employer Size, and Unemployment”. *International Economic Review* 39 (2), p. 257.
- CAHUC, P. and ZYLBERBERG, A. (2004). *Labor economics*. OCLC: 265445233. Cambridge, Mass.: MIT Press.
- CATTANEO, M. D., JANSSON, M. and MA, X. (2020). “Simple Local Polynomial Density Estimators”. *Journal of the American Statistical Association* 115 (531), pp. 1449–1455.
- CHERNOZHUKOV, V. and HANSEN, C. (2005). “An IV Model of Quantile Treatment Effects”. *Econometrica* 73 (1), pp. 245–261.
- CHETTY, R., FRIEDMAN, J. N., OLSEN, T. and PISTAFERRI, L. (2011). “Adjustment Costs, Firm Responses, and Micro vs. Macro Labor Supply Elasticities: Evidence from Danish Tax Records.” *Quarterly Journal of Economics* 126 (2), pp. 749–804.
- COSTA, D. L. (2000). “Hours of Work and the Fair Labor Standards Act: A Study of Retail and Wholesale Trade, 1938–1950”. *Industrial and Labor Relations Review*, p. 17.
- DÜMBGEN, L., KOLESNYK, P. and WILKE, R. A. (2017). “Bi-log-concave distribution functions”. *Journal of Statistical Planning and Inference* 184, pp. 1–17.
- DUBE, A., MANNING, A. and NAIDU, S. (2020). “Monopsony, Misoptimization, and Round Number Bunching in the Wage Distribution”. *NBER Working Paper* w24991.
- EHRENBERG, R. and SCHUMANN, P. (1982). *Longer hours or more jobs? : an investigation of amending hours legislation to create employment*. New York State School of Industrial and Labor Relations, Cornell University.
- EHRENBERG, R. G. (1971). “The Impact of the Overtime Premium on Employment and Hours in U . S . Industry”. *Economic Inquiry* 9 (2).
- EINAV, L., FINKELSTEIN, A. and SCHRIMPF, P. (2017). “Bunching at the kink: Implications for spending responses to health insurance contracts”. *Journal of Public Economics* 146, pp. 27–40.
- GRIGSBY, J., HURST, E. and YILDIRMAZ, A. (2020). *Aggregate Nominal Wage Adjustments: New Evidence from Administrative Payroll Data*. American Economic Review, forthcoming.
- HAMERMESH, D. S. (1996). *Labor demand*. Princeton, NJ: Princeton Univ. Press.

- HAMERMESH, D. S. and TREJO, S. J. (2000). “The Demand for Hours of Labor : Direct Evidence from California”. *The Review of Economics and Statistics* 82 (1), pp. 38–47.
- HART, R. A. (2004). *The economics of overtime working*. OCLC: 704550114. Cambridge, UK: Cambridge University Press.
- HJORT, J., LI, X. and SARSONS, H. (2020). “Across-Country Wage Compression in Multinationals”. *NBER Working Paper* w26788.
- IMBENS, G. W. and MANSKI, C. F. (2004). “Confidence Intervals for Partially Identified Parameters”. *Econometrica* 72, p. 14.
- JOHNSON, J. (2003). “The Impact of Federal Overtime Legislation on Public Sector Labor Markets”. *Journal of Labor Economics* 21 (1), pp. 43–69.
- KASY, M. (2017). “Who wins, who loses? Identification of the welfare impact of changing wages”. *Working Paper*, pp. 1–26.
- KLEVEN, H. J. (2016). “Bunching”. *Annual Review of Economics* 8 (June), pp. 435–464.
- KLEVEN, H. J. and WASEEM, M (2013). “Using notches to uncover optimization frictions and structural elasticities: Theory and evidence from Pakistan”. *The Quarterly Journal of Economics* 128 (2), pp. 669–723.
- MILGROM, P. and ROBERTS, J. (1996). “The LeChatelier Principle”. *American Economic Review* 1 (86), pp. 173–179.
- MOORE, D. T. (2021). “Evaluating Tax Reforms without Elasticities: What Bunching Can Identify”, p. 61.
- QUACH, S. (2021). “The Labor Market Effects of Expanding Overtime Coverage”. *MPRA Paper* 100613.
- ROSEN, S. (1968). “Short-Run Employment Variation on Class-I Railroads in the U.S., 1947-1963”. *Econometrica* 36 (3), p. 511.
- SAEZ, E. (2010). “Do Taxpayers Bunch at Kink Points?” *American Economic Journal: Economic Policy* 2 (3). ISBN: 1945-7731 _eprint: arXiv:1011.1669v3, pp. 180–212.
- SOCIETY FOR HUMAN RESOURCE MANAGEMENT (2018). “National Study of Employers”, p. 79.

- STOLE, L. A. and ZWIEBEL, J. (1996). “Intra-Firm Bargaining under Non-Binding Contracts”. *The Review of Economic Studies* 63 (3). Publisher: [Oxford University Press, Review of Economic Studies, Ltd.], pp. 375–410.
- STOYE, J. (2009). “More on Confidence Intervals for Partially Identified Parameters”. *Econometrica* 77 (4), pp. 1299–1315.
- TREJO, B. S. J. (1991). “The Effects of Overtime Pay Regulation on Worker Compensation”. *American Economic Review* 81 (4), pp. 719–740.
- U.S. DEPARTMENT OF LABOR (2019). “Defining and Delimiting the Exemptions for Executive, Administrative, Professional, Outside Sales and Computer Employees”. *Federal Register* 84 (188).

A Identification in a generalized bunching design

This section presents some generalization of the bunching-design model used in the main text. While the FLSA will provide a running example throughout, I largely abstract from the overtime context to emphasize the general applicability of the results.

To facilitate comparison with the existing literature on bunching at kinks – which has mostly considered cross-sectional data – I throughout this section suppress time indices and use the single index i to refer to each unit of observation (a paycheck in the overtime case). Further, the “running variable” of the bunching design is typically denoted by Y rather than h . This is done to emphasize the link to the treatment effects literature, while allowing a distinction that is in some cases important (e.g. models in which hours of pay for work differ from actual hours of work).

A.1 A generalized bunching-design model

Consider a population of observational units indexed by i . For each i , a decision-maker $d(i)$ chooses a point (z, \mathbf{x}) in some space $\mathcal{X} \subseteq \mathbb{R}^{m+1}$ where z is a scalar and \mathbf{x} a vector of m components, subject to a constraint of the form:

$$z \geq \max\{B_{0i}(\mathbf{x}), B_{1i}(\mathbf{x})\} \quad (\text{A.1})$$

The functions $B_{0i}(\mathbf{x})$ and $B_{1i}(\mathbf{x})$ are continuous and weakly convex functions of the vector \mathbf{x} , and that there exist continuous scalar functions $y_i(\mathbf{x})$ and a scalar k such that:

$$B_{0i}(\mathbf{x}) > B_{1i}(\mathbf{x}) \text{ whenever } y_i(\mathbf{x}) < k \quad \text{and} \quad B_{0i}(\mathbf{x}) < B_{1i}(\mathbf{x}) \text{ whenever } y_i(\mathbf{x}) > k$$

The value k is taken to be common to all units i , and is assumed to be known by the researcher.³⁴ In the overtime setting, $y_i(\mathbf{x})$ represents the hours of work for which a worker is paid in a given week, and $k = 40$. In most applications of the bunching design, the decision-maker $d(i)$ is simply i themselves, for example a worker choosing their labor supply subject to a tax kink. In the overtime application however i is a worker-week pair, and $d(i)$ is the worker’s firm.

Let X_i be i ’s realized outcome of \mathbf{x} , and $Y_i = y_i(X_i)$. I assume that Y_i is observed by the econometrician, but not that X_i is.

In a typical example, the functions B_{0i} , B_{1i} will represent a schedule of some kind of “cost” as a function of the choice vector \mathbf{x} , with two regimes of costs that are separated by the condition $y_i(\mathbf{x}) = k$, characterizing the locus of points at which the two cost functions cross. Let $B_{ki}(\mathbf{x}) := \max\{B_{0i}(\mathbf{x}), B_{1i}(\mathbf{x})\}$. Online Appendix 5 discusses a case of this from the literature in which the functions B_0 and B_1 depend on a vector \mathbf{x} of two components.³⁵ Budget constraints like Eq. $z \geq B_{ki}(\mathbf{x})$ are typically “kinked” because while the function $B_{ki}(\mathbf{x})$ is continuous, it will generally be non-differentiable at the \mathbf{x} for which $y_i(\mathbf{x}) = k$.³⁶ While the functions B_0 , B_1 and y can all depend on i , I will often suppress this dependency for clarity of notation.

In the most common cases from the literature, \mathbf{x} is assumed to be the scalar $y_i(x) = x$, i.e. there is no distinction between the “kink variable” y and underlying choice variables \mathbf{x} . For example, the seminal bunching design papers Saez (2010) and Chetty et al. (2011) considered progressive taxation with z being tax liability (or credits), both $y = x$ corresponding to taxable income, and B_0 and B_1 linear tax functions on either side of a threshold y between two adjacent tax/benefit brackets. Similarly, in the overtime context, the functions B_0 and B_1 are linear and only depend on hours $y_i(\mathbf{x})$, as depicted in Figure A.1. However, even when the functions B_0 and B_1 only depend on \mathbf{x} through $y_i(\mathbf{x})$, the bunching design is compatible with models in which multiple margins of choice respond to the incentives provided by the kink. In fact, the econometrician may be agnostic as to even what the full set of components of \mathbf{x} are, with $y(\cdot)$, $B_0(\cdot)$ or $B_1(\cdot)$ depending only on various subsets of them. The next section will discuss how the bunching design allows us to conduct causal inference on the variable Y_i , but not directly on the underlying choice variables X_i .

³⁴This comes at little cost of generality since with heterogeneous k_i this could be subsumed as a constant into the function $y_i(\mathbf{x})$, so long as the k_i are observed by the researcher.

³⁵An example from the literature in which a distinction between y and \mathbf{x} cannot be avoided is Best et al. (2015). These authors study firms in Pakistan, who pay either a tax on output or a tax on profit, whichever is higher. The two tax schedules cross when the ratio of profits to output crosses a certain threshold that is pinned down by the two respective tax rates. In this case, the variable y depends both on production and on reported costs, leading to two margins of response to the kink: one from choosing the scale of production and the other from choosing whether and how much to misreport costs.

³⁶In particular, the subgradient of $\max\{B_{0i}(\mathbf{x}), B_{1i}(\mathbf{x})\}$ will depend on whether one approaches from the $y_i(\mathbf{x}) > k$ or the $y_i(\mathbf{x}) < k$ side. For example with a scalar x and linear B_0 and B_1 , the derivative of $B_{ki}(x)$ discontinuously rises at \mathbf{x} for which $y_i(\mathbf{x}) = k$.

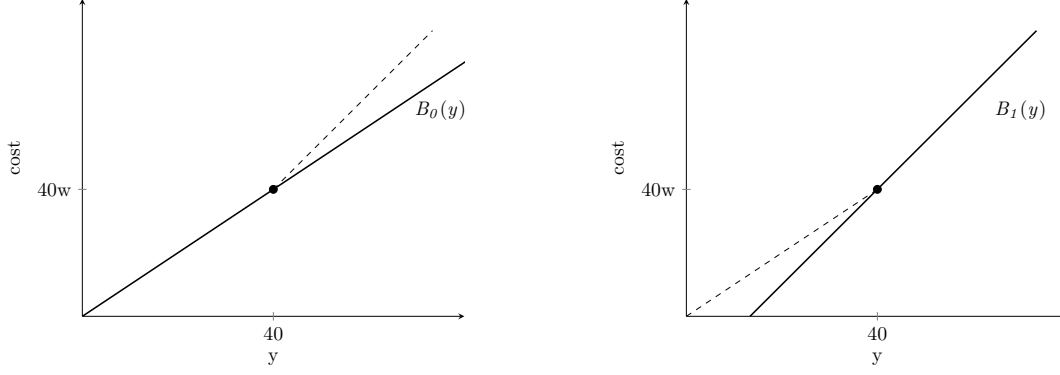


FIGURE A.1: Definition of counterfactual cost functions B_0 and B_1 that firms could have faced, absent the overtime kink. Dashed lines show the rest of actual cost function in comparison to the counterfactual as a solid line. Note that we use the notation y here to indicate hours, rather than the h used in the main text.

A.2 Potential outcomes as counterfactual choices

To introduce a notion of treatment effects in the bunching design, I define a pair of potential outcomes as what would occur if the decision-maker faced either of the functions B_0 or B_1 globally, without the kink.

Definition (potential outcomes). Let Y_{0i} be the value of $y_i(\mathbf{x})$ that would occur for unit i if $d(i)$ faced the constraint $z \geq B_0(\mathbf{x})$, and let Y_{1i} be the value that would occur under the constraint $z \geq B_1(\mathbf{x})$.

To relate these potential outcomes to choices of the decision-maker, we make explicit the assumption that they control the value of $y_i(\mathbf{x})$. For any function B let Y_{Bi} be the outcome that would occur under the choice constraint $z \geq B(\mathbf{x})$, with Y_{0i} and Y_{1i} shorthands for Y_{B_0i} and Y_{B_1i} , respectively.³⁷

Assumption CHOICE (perfect manipulation of y). For any function $B(\mathbf{x})$, $Y_{Bi} = y_i(\mathbf{x}_{Bi})$, where $(z_{Bi}, \mathbf{x}_{Bi})$ is the choice that $d(i)$ would make under the constraint $z \geq B(\mathbf{x})$.

Assumption CHOICE rules out for example optimization error, which could limit the decision-maker's ability to exactly manipulate values of \mathbf{x} and hence y . It also takes for granted that counterfactual choices are unique, and rules out some kinds of extensive margin effects in which a decision-maker would not choose any value of Y at all under B_1 or B_0 . Note that CHOICE here differs from the version given in the main text in that it applies to all functions B , not just B_0 , B_1 and B_k (this is useful for Theorem 2).

³⁷In this notation Assumption CHOICE implies that the actual outcome Y_i observed by the econometrician is equal to $Y_{B_{ki}i}$.

The central behavioral assumption that allows us to reason about the counterfactuals Y_0 and Y_1 is that decision-makers have convex preferences over (c, \mathbf{x}) and dislike costs z :

Assumption CONVEX (strictly convex preferences except at kink, monotonic in z). For each i and any function $B(\mathbf{x})$, choice is $(z_{Bi}, \mathbf{x}_{Bi}) = \operatorname{argmax}_{z, \mathbf{x}} \{u_i(z, \mathbf{x}) : z \geq B(\mathbf{x})\}$ where $u_i(z, \mathbf{x})$ is weakly decreasing in z and satisfies

$$u_i(\theta z + (1 - \theta)z^*, \theta \mathbf{x} + (1 - \theta)\mathbf{x}^*) > \min\{u_i(z, \mathbf{x}), u_i(z^*, \mathbf{x}^*)\}$$

for any $\theta \in (0, 1)$ and points $(z, \mathbf{x}), (z^*, \mathbf{x}^*)$ such that $y_i(\mathbf{x}) \neq k$ and $y_i(\mathbf{x}^*) \neq k$.

Note: The function $u_i(\cdot)$ represents preferences over choice variables for unit i , but the preferences are those of the decision maker $d(i)$. I avoid more explicit notation like $u_{d(i), i}(\cdot)$ for brevity. In the overtime setting with firms choosing hours, $u_i(z, \mathbf{x})$ corresponds to the firm's profit function π as a function of the hours of a particular worker this week, and costs this week z for that worker.

The second part of Assumption CONVEX is implied by strict quasi-concavity of the function (z, \mathbf{x}) , corresponding to strictly convex preferences. However it also allows for decision-makers preferences to have “two peaks”, provided that one of the peaks is located exactly at the kink. This is useful in cases in which the kink is located at a point that has particular value to decision-makers, such as firms setting weekly hours. For example, suppose that firms choose hours only $\mathbf{x} = h$, and have preferences of the form:

$$u_i(z, h) = af(h) + \phi \cdot \mathbb{1}(h = 40) - z \tag{A.2}$$

where $f(h)$ is strictly concave. This allows firms to have a behavioral “bias” towards 40 hours, or to extract extra profits when $h = 40$ exactly.³⁸ Figure A.2 depicts an example of such preferences, given an arbitrary linear budget function $B(h)$.

Note: The notation of Assumption CONVEX does not make explicit any dependence of the functions $u_i(\cdot)$ on the choices made for other observational units $i' \neq i$. When the functions $u_i(\cdot)$ are indeed invariant over such counterfactual choices, we have a version of the no-interference condition of the stable unit treatment values assumption (SUTVA). Maintaining SUTVA is not necessary to define treatment effects in the bunching design, provided that the variables y and z can be coherently defined at the individual unit i level (see Example 3 in Section A.3, and Online Appendix 3). Nevertheless, the interpretation of the treatment effects identified by the bunching design is most

³⁸If a mass of firms were to have preferences of this form, then it would be natural to expect bunching in the distributions of h_{0it} and h_{1it} , which I allow in Section 5.

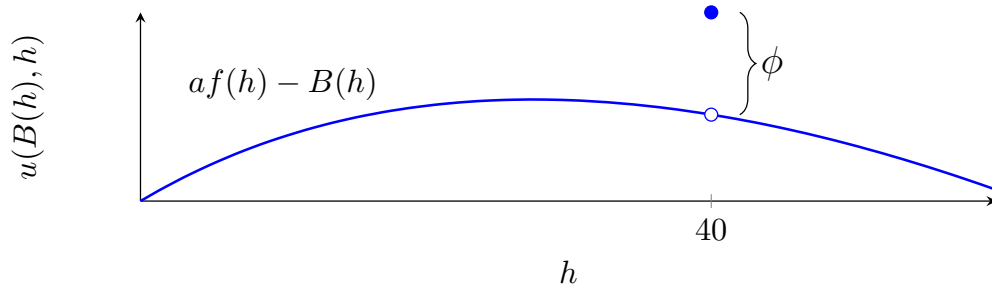


FIGURE A.2: An example of preferences that satisfy CONVEX but are not strictly convex, cf. Eq. (A.2).

straightforward when SUTVA does hold. This assumption is standard in the bunching design.³⁹

A weaker assumption than CONVEX that will still have identifying power is simply that decision-makers' choices do not violate the weak axiom of revealed preference:

Assumption WARP (rationalizable choices). *Consider two budget functions B and B' and any unit i . If $d(i)$'s choice under B' is feasible under B , i.e. $z_{B'i} \geq B(\mathbf{x}_{B'i})$, then $(z_{Bi}, \mathbf{x}_{Bi}) = (z_{B'i}, \mathbf{x}_{B'i})$.*

I make the stronger assumption CONVEX for most of the identification results, but Assumption WARP still allows a version of many of them in which equalities become weak inequalities, indicating a degree of robustness with respect to departures from convexity. Note that the monotonicity assumption in CONVEX implies that choices will always satisfy $z = B(\mathbf{x})$, i.e. agents' choices will lay on their cost functions (despite Eq. A.1 being an inequality, indicating “free-disposal”).

Further notes on the general model

I conclude this section with some further remarks on the generality of Eq. (A.1) given the above assumptions. The first is that the budget functions B_0 and B_1 can depend on a subset of the variables that enter into the function for y , and vice versa. In the former case, this is because the only restriction on $B_{0i}(\mathbf{x})$ and $B_{1i}(\mathbf{x})$ is that they are continuous and weakly convex in all components of \mathbf{x} ; thus, having zero dependence on a component of \mathbf{x} is permissible. This is of particular interest because while the variables entering into the budget functions are generally known from the empirical context generating the kink, the model can allow additional choice variables to enter into the threshold-crossing variable y , that may not even be known to the econometrician. The next section gives some illustrative examples for the overtime setting.

³⁹I note that SUTVA issues like those addressed in Online Appendix 3 could also occur in canonical bunching designs: for example if spouses choose their labor supply jointly, the introduction of a tax kink may cause one spouse to increase labor supply while the other decreases theirs.

A.3 Examples from the general choice model in the overtime setting

To demonstrate the flexibility of the general choice model, I below present some examples in the overtime context. These examples are illustrative, and each could apply to a different subset of units in the population. In these examples we continue to take the decision-maker for a given unit to be the firm employing that worker.⁴⁰

Example 1: Substitution from bonus pay

Let the firm's choice vector be $\mathbf{x} = (h, b)'$, where $b \geq 0$ indicates a bonus (or other fringe benefit) paid to the worker. Firms may find it optimal to offer bonuses to improve worker satisfaction and reduce turnover. Suppose firm preferences are: $\pi(z, h, b) = f(h) + g(z + b - \nu(h)) - z - b$, where z continues to denote wage compensation this week, $z + b - \nu(h)$ is the worker's utility with $\nu(h)$ a convex disutility from labor h , and $g(\cdot)$ increasing and concave. In this model firms will choose the surplus maximizing choice of hours $h_m := \arg\max_h f(h) - \nu(h)$, provided that the corresponding optimal bonus is non-negative. Bonuses fully adjust to counteract overtime costs, and $h_0 = h_1 = h_m$.

Example 2: Off-the-clock hours and paid breaks

Suppose firms choose a pair $\mathbf{x} = (h, o)'$ with h hours worked and o hours worked "off-the-clock", such that $y(\mathbf{x}) = h - o$ are the hours for which the worker is ultimately paid. Evasion is harder the larger o is, which could be represented by firms facing a convex evasion cost $\phi(o)$, so that firm utility is $\pi(z, h, o) = f(h) - \phi(o) - z$.⁴¹ This model can also include some firms voluntarily offering paid breaks by allowing o to be negative.

Example 3: Complementaries between workers or weeks

Suppose the firm simultaneously chooses the hours $\mathbf{x} = (h, g)$ of two workers according to production that is isoelastic in a CES aggregate (g could also denote planned hours next week): $\pi(z, h, g) = a \cdot ((\gamma h^\rho + g^\rho)^{1/\rho})^{1+\frac{1}{\epsilon}} - z$ with γ a relative productivity shock. Let g^* denote the firm's optimal choice of hours for the second worker. Optimal h then maximizes $\pi(z, h, g^*)$ subject to $z = B_k(h)$, as if the firm faced a single-worker production function of $f(h) = a \cdot ((\gamma h^\rho + g^{*\rho})^{1/\rho})^{1+\frac{1}{\epsilon}}$. This function is more elastic than $a \cdot h^{1+\frac{1}{\epsilon}}$ provided that $\rho < 1 + 1/\epsilon$,

⁴⁰Appendix 2 discusses a further example in which the firm and worker bargain over this week's hours. This model can attenuate the wage elasticity of chosen hours since overtime pay gives the parties opposing incentives.

⁴¹Note that the data observed in our sample are of hours of work $y(\mathbf{x})$ for which the worker is paid, when this differs from h . Appendix A describes how Equation 2 still holds, but for counterfactual values of hours paid $y = h - o$ rather than hours worked h . The bunching design lets us investigate treatment effects on paid hours, without observing off-the-clock hours or break time o .

attenuating the response to an increase in w implied by a given ϵ .⁴² Section 4.4 discusses how complementarities affect the final evaluation of the FLSA.

A.4 Observables in the kink bunching design

Lemma 1 outlines the core consequence of Assumption CONVEX for the relationship between observed Y_i and the potential outcomes introduced in the last section:

Lemma 1 (realized choices as truncated potential outcomes). *Under Assumptions CONVEX and CHOICE, the outcome observed given the constraint $z \geq \max\{B_{0i}(\mathbf{x}), B_{1i}(\mathbf{x})\}$ is:*

$$Y_i = \begin{cases} Y_{0i} & \text{if } Y_{0i} < k \\ k & \text{if } Y_{1i} \leq k \leq Y_{0i} \\ Y_{1i} & \text{if } Y_{1i} > k \end{cases}$$

Proof. See Appendix B. □

Lemma 1 says that the pair of counterfactual outcomes (Y_{0i}, Y_{1i}) is sufficient to pin down actual choice Y_i , which can be seen as an observation of one or the other potential outcome, or k , depending on how the potential outcomes relate to the kink point k .

Note that the “straddling” event $Y_{0i} \leq k \leq Y_{1i}$ from Lemma 1 can be written as $Y_{0i} \in [k, k + \Delta_i]$, where $\Delta_i = Y_{0i} - Y_{1i}$. Equivalently, we could also write $Y_{1i} \leq k \leq Y_{0i}$ as $Y_i \in [k - \Delta_i, k]$. This forms the basic link between bunching and treatment effects, where Δ_i can be thought of as the causal effect of a counterfactual change from the choice set under B_1 to the choice set under B_0 . Let $\mathcal{B} := P(Y_i = k)$ be the observable probability that the decision-maker chooses to locate exactly at $Y = k$. This relationship holds in a weakened form under WARP rather than CONVEX:

Proposition 1 (relation between bunching and straddling). *a) Under CONVEX and CHOICE: $\mathcal{B} = P(Y_{0i} \in [k, k + \Delta_i])$; b) under WARP and CHOICE: $\mathcal{B} \leq P(Y_{0i} \in [k, k + \Delta_i])$.*

Proof. See Appendix B. □

Consider a random sample of observations of Y_i . Under i.i.d. sampling of Y_i , the distribution $F(y)$ of Y_i is identified.⁴³ Let $F_1(y) = P(Y_{0i} \leq y)$ be the distribution function of the random

⁴²This expression overstates the degree of attenuation somewhat, since h_1 and h_0 maximize $f(h)$ above for different values g^* , which leads to a larger gap between h_0 and h_1 compared with a fixed g^* by the Le Chatelier principle (Milgrom and Roberts, 1996). However h_1/h_0 still increases on net given $\rho < 1 + 1/\epsilon$.

⁴³Note that in the overtime application sampling is actually at the firm level, which coincides with the level of decision-making units $d(i)$.

variable Y_0 , and $F_1(y)$ the distribution function of Y_1 . From Lemma 1 it follows immediately that $F_0(y) = F(y)$ for all $y < k$, and $F_1(y) = F(y)$ for $Y > k$. Thus observations of Y_i are also informative about the marginal distributions of Y_{0i} and Y_{1i} . Again, a weaker version of this also holds under WARP rather than CONVEX:

Proposition 2 (identification of truncated densities). *Suppose that F_0 and F_1 are continuously differentiable with derivatives f_0 and f_1 , and that F admits a derivative function $f(y)$ for $y \neq k$. Under WARP and CHOICE: $f_0(y) \leq f(y)$ for $y < k$ and $f_0(k) \leq \lim_{y \uparrow k} f(y)$, while $f_1(y) \leq f(y)$ for $y > k$ and $f_1(k) \leq \lim_{y \downarrow k} f(y)$, with equalities under CONVEX.*

Proof. See Appendix B. □

As an example of how WARP alone (without CONVEX) can still be useful for identification, suppose that $\Delta_i = \Delta$ were known to be homogenous across units,⁴⁴ and $f_0(y)$ were constant across the interval $[k, k + \Delta]$, then by Propositions 1 and 2 we have that $\Delta \geq \mathcal{B}/f_0(k)$ under WARP and CHOICE.

A.5 Treatment effects in the bunching design

Proposition 1 establishes that bunching can be informative about features of the distribution of treatment effects Δ_i . This section discusses the interpretation of these treatment effects and some additional identification results omitted in the main text.

Discussion of treatment effects vs. structural parameters:

The treatment effects Δ_i are “reduced form” in the sense that when the decision-maker has multiple margins of response \mathbf{x} to the incentives introduced by the kink, these may be bundled together in the treatment effect Δ_i . This clarifies a limitation sometimes levied against the bunching design, while also revealing a perhaps under-appreciated strength. On the one hand, it is not always clear “which elasticity” is elicited by bunching at a kink, complicating efforts to identify a elasticity parameter having a firm structural interpretation.

On the other hand, the bunching design can be useful for ex-post policy evaluation and even forecasting effects of small policy changes (as described in Section 4.4), without committing to a tightly parameterized underlying model of choice. This provides a response to Einav et al. (2017), who caution that alternative structural models calibrated from the bunching-design can yield very different predictions about counterfactuals. By focusing on the counterfactuals Y_{0i} and Y_{1i} , we can

⁴⁴One way to get homogenous treatment effects in levels in the overtime setting is to assume exponential production: $f(h) = \gamma(1 - e^{-h/\gamma})$ where $\gamma > 0$ and $h_{0it} - h_{1it} = \gamma \ln(1.5)$ for all units.

specify a *particular* type of counterfactual question that can be answered robustly across a broad class of models.

The “trick” of Lemma 1 is to express the observable data in terms of counterfactual choices, rather than of primitives of the utility function. The underlying utility function $u_i(z, \mathbf{x})$ is used only as an intermediate step, only requiring the nonparametric restrictions of convexity and monotonicity rather than knowing its functional form. The econometrician need not even know the full vector \mathbf{x} of choice variables underlying agents’ observed value of y , they simply need to believe that preferences are convex in them, and verify that B_0 and B_1 are convex in a subset of them (no dependence on some components of \mathbf{x} satisfies weak convexity). This greatly increases the robustness of the method to potential misspecification of the underlying choice model. Online Appendix 5 illustrates this point looking at the setting of Best et al. (2015).

Additional identification results for the bunching design:

While Theorem 1 of Section 4 develops the treatment effect identification result used to evaluate the FLSA, Supplemental Appendix 5 presents some further identification results for the bunching design that are not used in this paper, which can be considered alternatives to Theorem 1. This includes re-expressing various results in the general framework of this section, including the linear interpolation approach of Saez (2010), the polynomial approach of Chetty et al. (2011) and a “small-kink” approximation appearing in Saez (2010) and Kleven (2016). The Supplemental Appendix also outlines alternative shape constraints to bi-log-concavity, including monotonicity of densities. I also give there a result in which a lower bound to a certain local average treatment effect is identified under WARP, without requiring convexity of preferences.

The buncher ATE when Assumption RANK fails:

This section picks up from the discussion in Section 4.3, but continues with the notation of this Appendix. When RANK fails (and $p = 0$ for simplicity), the bounds from Theorem 1 are still valid under BLC of Y_0 and Y_1 for the averaged quantile treatment effect:

$$\frac{1}{\mathcal{B}} \int_{F_0(k)}^{F_1(k)} Q_0(u) - Q_1(u) = \mathbb{E}[Y_{0i} | Y_{0i} \in [k, k + \Delta_0^*]] - \mathbb{E}[Y_{1i} | Y_{1i} \in [k - \Delta_1^*, k]], \quad (\text{A.3})$$

where $\Delta_0^* := Q_0(F_1(k)) - Q_1(F_1(k)) = Q_0(F_1(k)) - k$ and $\Delta_1^* := Q_0(F_0(k)) - Q_1(F_0(k)) = k - Q_1(F_0(k))$. Thus, Δ_0^* is the value such that $F_0(k + \Delta_0^*) = F_0(k) + \mathcal{B}$, and Δ_1^* is the value such that $F_1(k - \Delta_1^*) = F_1(k) - \mathcal{B}$. The averaged quantile treatment effect of Eq. (A.3) yields a lower bound on the buncher ATE, as described in Fig. A.3.

Signing the bias when RANK fails

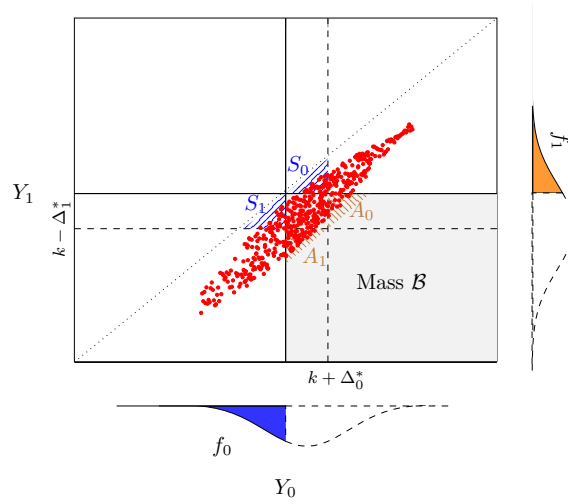


FIGURE A.3: When Assumption RANK fails, the average $\mathbb{E}[Y_{0i}|Y_{0i} \in [k, k + \Delta_0^*]]$ will include the mass in the region S_0 , who are not bunchers (NE lines) but will be missing the mass in the region A_0 (NW lines) who are. This causes an under-estimate of the desired quantity $\mathbb{E}[Y_{0i}|Y_{1i} \leq k \leq Y_{0i}]$. Similarly, $\mathbb{E}[Y_{1i}|Y_{1i} \in [k - \Delta_1^*, k]]$ will include the mass in the region S_1 , who are not bunchers but will be missing the mass in A_1 , who are. This causes an over-estimate of the desired quantity $\mathbb{E}[Y_{1i}|Y_{1i} \leq k \leq Y_{0i}]$.

A.6 Policy changes in the bunching-design

This section discusses the logic establishing Theorem 1 in the main text. Consider a bunching design setting in which the cost functions B_0 and B_1 can be viewed as members of family $B_i(\mathbf{x}; \rho, k)$ parameterized by a continuum of scalars ρ and k , where $B_{0i}(\mathbf{x}) = B_i(\mathbf{x}; \rho_0, k^*)$ and $B_{1i}(\mathbf{x}) = B_i(\mathbf{x}; \rho_1, k^*)$ for some $\rho_1 > \rho_0$ and value k^* of k . In the overtime setting ρ represents a wage-scaling factor, with $\rho = 1$ for straight-time and $\rho = 1.5$ for overtime:

$$B_i(y; \rho, k) = \rho w_i y - k w_i (\rho - 1) \quad (\text{A.4})$$

where work hours y may continue to be a function $y(\mathbf{x})$ of a vector of choice variables to the firm. Here ρ represents an arbitrary wage-scaling factor, while k controls the size of a lump-sum subsidy that keeps $B_i(k; \rho, k)$ invariant across ρ .

Assume that ρ takes values in a convex subset of \mathbb{R} containing ρ_0 and ρ_1 , and that for any k and $\rho' > \rho$ the cost functions $B_i(\mathbf{x}; \rho, k)$ and $B_i(\mathbf{x}; \rho', k)$ satisfy the conditions of the bunching design framework from Section 4, with the function $y_i(\mathbf{x})$ fixed across all such values. That is, $B_i(\mathbf{x}; \rho', k) > B_i(\mathbf{x}; \rho, k)$ iff $y_i(\mathbf{x}) > k$ with equality when $y_i(\mathbf{x}) = k$, the functions $B_i(\cdot; \rho, k)$ are weakly convex and continuous, and $y_i(\cdot)$ is continuous. It is readily verified that Equation (A.4)

satisfies these requirements with $y_i(h) = h$.⁴⁵

For any value of ρ , let $Y_i(\rho, k)$ be agent i 's realized value of $y_i(\mathbf{x})$ when a choice of (z, \mathbf{x}) is made under the constraint $c \geq B_i(\mathbf{x}; \rho, k)$. A natural restriction in the overtime setting that is that the function $Y_i(\rho, k)$ does not depend on k , and some of the results below will require this. A sufficient condition for $Y_i(\rho, k) = Y_i(\rho)$ is a family of cost functions that are linearly separable in k , as we have in the overtime setting with Equation (A.4), along with quasi-linearity of preferences:

Assumption SEPARABLE (invariance of potential outcomes with respect to k). *For all i, ρ and k , $B_i(\mathbf{x}; \rho, k)$ is additively separable between k and \mathbf{x} (e.g. $b_i(\mathbf{x}, \rho) + \phi_i(\rho, k)$ for some functions b_i and ϕ_i), and for all i $u_i(z, \mathbf{x})$ can be chosen to be additively separable and linear in z .*

Quasilinearity of preferences is a property of profit-maximizing firms when c represents a cost, and is thus a natural assumption in the overtime setting. However, additive separability of $B_i(\mathbf{x}; \rho, k)$ in k may be context specific: in the example from Best et al. (2015) described in Appendix A, quasi-linearity of preferences is not sufficient since the cost functions are not additively separable in k . To maintain clarity of exposition, I will keep k implicit in $Y_i(\rho)$ throughout the foregoing discussion, but the proofs make it clear when SEPARABLE is being used.

Below I state two intermediate results that allow us to derive expressions for the effects of marginal changes to ρ_1 or k on hours. Lemma 2 generalizes an existing result from Blomquist et al. (2021), and makes use of a regularity condition I introduce in the proof as Assumption SMOOTH.⁴⁶ Counterfactual bunchers $K_i^* = 1$ are assumed to stay at k^* , regardless of ρ and k . Let $p(k) = p \cdot \mathbb{1}(k = k^*)$ denote the possible counterfactual mass at the kink as a function of k . Let $f_\rho(y)$ be the density of $Y_i(\rho)$, which exists by SMOOTH and is defined for $y = k^*$ as a limit (see proof).

Lemma 2 (bunching from marginal responsiveness). *Assume CHOICE, SMOOTH and WARP. Then:*

$$\mathcal{B} - p(k) \leq \int_{\rho_0}^{\rho_1} f_\rho(k) \mathbb{E} \left[-\frac{dY_i(\rho)}{d\rho} \middle| Y_i(\rho) = k \right] d\rho$$

with equality under CONVEX.

Proof. See Online Appendices. □

The main tool in establishing Lemma 2 is relate the integrand in the above to the rate at which kink-induced bunching goes away as the “size” of the kink goes to zero.

⁴⁵As an alternative example, I construct in Appendix A functions $B_i(\mathbf{x}; \rho, k)$ for the bunching design setting from Best et al. (2015). In that case, ρ parameterizes a smooth transition between an output and a profit tax, where k enters into the rate applied to the tax base for that value of ρ .

⁴⁶Blomquist et al. (2021) derive the special case of Lemma 2 with CONVEX and $p = 0$, in the context of a more restricted model of labor supply under taxation. I establish it here for the general bunching design model where in particular, the $Y_i(\rho)$ may depend on an underlying vector \mathbf{x} which are not observed by the econometrician. I also use different regularity conditions.

Lemma SMALL (small kink limit). Assume *CHOICE**, *WARP*, and *SMOOTH*. Then:

$$\lim_{\rho' \downarrow \rho} \frac{P(Y_i(\rho') \leq k \leq Y_i(\rho)) - p(k)}{\rho' - \rho} = -f_\rho(k) \mathbb{E} \left[\frac{dY_i(\rho)}{d\rho} \middle| Y_i(\rho) = k \right]$$

Proof. See Online Appendices. □

Note that the quantity $P(Y_i(\rho') \leq k \leq Y_i(\rho)) - p(k)$ is an upper bound on the bunching that would occur due to a kink between budget functions $B_i(\mathbf{x}; \rho, k)$ and $B_i(\mathbf{x}; \rho', k)$ (with equality under CONVEX). Lemma SMALL thus shows that the small-kink approximation that has appeared in Saez (2010) and Kleven (2016) (stated in Theorem 7 of the Online Appendix) becomes exact in the limit that the two cost functions approach one another, since treatment effects are:

$$Y_i(\rho) - Y_i(\rho') = \frac{dY_i(\rho)}{d\rho}(\rho' - \rho) + O((\rho' - \rho)^2)$$

By Lemma 2, we can also see that the RHS in Lemma SMALL evaluated at $\rho = \rho_1$ is equal to the derivative of bunching as ρ_1 is increased, under CONVEX.

Lemma 2 is particularly useful when combined with a result from Kasy (2017), which considers how the distribution of a generic outcome variable changes as heterogeneous units flow to different values of that variable in response to marginal policy changes.

Lemma 3 (flows under a small change to ρ). Under *SMOOTH*:

$$\partial_\rho f_\rho(y) = \partial_y \left\{ f_\rho(y) \mathbb{E} \left[-\frac{dY_i(\rho)}{d\rho} \middle| Y_i(\rho) = y, K_i^* = 0 \right] \right\}$$

Proof. See Online Appendices or Kasy (2017) for original result. □

The intuition behind Lemma 3 comes from the physical dynamics of fluids. When ρ changes, a mass of units will “flow” out of a small neighborhood around any y , and this mass is proportional to the density at y and to the average rate at which units move in response to the change. When the magnitude of this net flow varies with y , the change to ρ will lead to a change in the density there.

With ρ_0 fixed at some value, let us index observed Y_i and bunching \mathcal{B} with the superscript $[k, \rho_1]$ when they occur in a kinked policy environment with cost functions $B_i(\cdot; \rho_0, k)$ and $B_i(\cdot; \rho_1, k)$. Lemmas 2 and 3 together imply Theorem 2, which in the notation of this section reads as:

1. $\partial_k \left\{ \mathcal{B}^{[k, \rho_1]} - p(k) \right\} = f_1(k) - f_0(k)$
2. $\partial_k \mathbb{E}[Y_i^{[k, \rho_1]}] = \mathcal{B}^{[k, \rho_1]} - p(k)$
3. $\partial_{\rho_1} \mathcal{B}^{[k, \rho_1]} = -k f_{\rho_1}(k) \mathbb{E} \left[\frac{dY_i(\rho_1)}{d\rho} \middle| Y_i(\rho_1) = k \right]$

$$4. \partial_{\rho_1} \mathbb{E}[Y_i^{[k, \rho_1]}] = - \int_k^\infty f_{\rho_1}(y) \mathbb{E} \left[\frac{dY_i(\rho_1)}{d\rho} \middle| Y_i(\rho_1) = y \right] dy$$

Note: Assumption SEPARABLE is only necessary for Items 1-2 in Theorem 2, Item 3 holds without it and with $\frac{\partial Y_i(\rho, k)}{\partial \rho}$ replacing $\frac{dY_i(\rho)}{d\rho}$.

B Main Proofs

B.1 Proof of Lemma 1

The proof proceeds in the following two steps:

- i) First, I show that $Y_{0i} \leq k$ implies that $Y_i = Y_{0i}$, and similarly $Y_{1i} \geq k$ implies that $Y_i = Y_{1i}$. This holds under CONVEX but also under the weaker assumption of WARP.
- ii) Second, I show that under CONVEX $Y_i < k \implies Y_i = Y_{0i}$ and $Y_i > k \implies Y_i = Y_{1i}$.

Item i) above establishes the first and third cases of Lemma 1. The only remaining possible case is that $Y_{1i} \leq k \leq Y_{0i}$. However, to finish establishing Lemma 1, we also need the reverse implication: that $Y_{1i} \leq k \leq Y_{0i}$ implies $Y_i = k$. This comes from taking the contrapositive of each of the two claims in item ii).

Proof of i): Let $\mathcal{X}_{0i} = \{\mathbf{x} : y_i(\mathbf{x}) \leq k\}$ and $\mathcal{X}_{1i} = \{\mathbf{x} : y_i(\mathbf{x}) \geq k\}$. If $Y_{0i} \leq k$, then by CHOICE \mathbf{x}_{B_0} is in \mathcal{X}_{0i} . Since $B_k(\mathbf{x}) = B_0(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{X}_{0i}$, it follows that $z_{B_0i} \geq B_k(\mathbf{x}_{B_0i})$, i.e. Y_{0i} is feasible under B_k . Note that $B_{ki}(\mathbf{x}) \geq B_{0i}(\mathbf{x})$ for all \mathbf{x} . By WARP then $(z_{B_{ki}}, \mathbf{x}_{B_{ki}}) = (z_{B_{0i}}, \mathbf{x}_{B_{0i}})$. Thus $Y_i = y_i(\mathbf{x}_{B_k}) = y_i(\mathbf{x}_{B_0}) = Y_{0i}$. So $Y_{0i} \leq k \implies Y_i = Y_{0i}$. By the same logic we can show that $Y_{1i} \geq k \implies Y_i = Y_{1i}$.

Proof of ii): For any convex budget function $B(\mathbf{x})$, $(z_{Bi}, \mathbf{x}_{Bi}) = \operatorname{argmax}_{z, \mathbf{x}} \{u_i(z, \mathbf{x}) \text{ s.t. } z \geq B(\mathbf{x})\}$. If $u_i(z, \mathbf{x})$ is strictly quasi-concave, then the RHS exists and is unique since it maximizes u_i over the convex domain $\{(z, \mathbf{x}) : z \geq B(\mathbf{x})\}$. Furthermore, by monotonicity of $u(z, \mathbf{x})$ in z we may substitute in the constraint $z = B(\mathbf{x})$ and write

$$\mathbf{x}_{Bi} = \operatorname{argmax}_{\mathbf{x}} u_i(B(\mathbf{x}), \mathbf{x})$$

Suppose that $y_i(\mathbf{x}_{Bi}) \neq k$, and consider any $\mathbf{x} \neq \mathbf{x}_{Bi}$ such that $y_i(\mathbf{x}) \neq k$. Let $\tilde{\mathbf{x}} = \theta \mathbf{x} + (1 - \theta) \mathbf{x}^*$ where $\mathbf{x}^* = \mathbf{x}_{Bi}$ and $\theta \in (0, 1)$. Since $B(\mathbf{x})$ is convex in \mathbf{x} and $u_i(z, \mathbf{x})$ is weakly decreasing in z :

$$u_i(B(\tilde{\mathbf{x}}), \tilde{\mathbf{x}}) \geq u_i(\theta B(\mathbf{x}) + (1 - \theta) B(\mathbf{x}^*), \tilde{\mathbf{x}}) > \min\{u_i(B(\mathbf{x}), \mathbf{x}), u_i(B(\mathbf{x}^*), \mathbf{x}^*)\} = u_i(B(\mathbf{x}), \mathbf{x}) \quad (\text{B.5})$$

where I have used CONVEX in the second step, and that \mathbf{x}^* is a maximizer in the third. This result implies that for any such $\mathbf{x} \neq \mathbf{x}^*$, if one draws a line between \mathbf{x} and \mathbf{x}^* , the function $u_i(B(\mathbf{x}), \mathbf{x})$ is strictly increasing as one moves towards \mathbf{x}^* . When \mathbf{x} is a scalar, this argument is used by Blomquist et al. (2015) (see Lemma A2 therein) to show that $u_i(B(\mathbf{x}), \mathbf{x})$ is strictly increasing to the left of \mathbf{x}^* , and strictly decreasing to the right of \mathbf{x}^* . Note that for any (binding) linear budget constraint $B(\mathbf{x})$, the result holds without monotonicity of $u_i(z, \mathbf{x})$ in z . This is useful for Theorem 2* in which some workers choose their hours.

For any function B , let $u_{Bi}(\mathbf{x}) = u_i(B(\mathbf{x}), \mathbf{x})$, and note that

$$u_{B_k i}(\mathbf{x}) = \begin{cases} u_{B_{0i}}(\mathbf{x}) & \text{if } \mathbf{x} \in \mathcal{X}_{0i} \\ u_{B_{1i}}(\mathbf{x}) & \text{if } \mathbf{x} \in \mathcal{X}_{1i} \end{cases}$$

Let \mathbf{x}_{ki} be the unique maximizer of $u_{B_k i}(\mathbf{x})$, where $Y_i = y_i(\mathbf{x}_{ki})$. Suppose that $Y_i < k$. Suppose furthermore that $Y_{0i} \neq Y_i$, with $Y_{0i} = y_i(\mathbf{x}_{0i})$ and \mathbf{x}_{0i} the maximizer of $u_{B_{0i}}(\mathbf{x})$. Note that we must have that $\mathbf{x}_{0i} \notin \mathcal{X}_{0i}$, because $B_0 = B_k$ in \mathcal{X}_{0i} so we can't have $u_{B_{0i}}(\mathbf{x}_{0i}) > u_{B_{0i}}(\mathbf{x}_{ki})$ (since \mathbf{x}_{ki} maximizes $u_{B_k i}(\mathbf{x})$). Thus $Y_{0i} > k$.

By continuity of $y_i(\mathbf{x})$, \mathcal{X}_{0i} is a closed set and \mathbf{x}_{ki} belongs to the interior of \mathcal{X}_{0i} . Thus, while \mathbf{x}_{0i} is not in \mathcal{X}_{0i} , there exists a point $\tilde{\mathbf{x}} \in \mathcal{X}_{0i}$ along the line between \mathbf{x}_{0i} to \mathbf{x}_{ki} . Since $Y_i \neq k$ and $Y_{0i} \neq k$, Eq. (B.5) with $B = B_k$ then implies that $u_{B_k i}(\tilde{\mathbf{x}}) > u_{B_k i}(\mathbf{x}_{0i})$. Since $u_{B_{0i}}(\mathbf{x}) = u_{B_k i}(\mathbf{x})$ for all \mathbf{x} in \mathcal{X}_{0i} , it follows that $u_{B_{0i}}(\tilde{\mathbf{x}}) > u_{B_{0i}}(\mathbf{x}_{0i})$. However, this contradicts the premise that \mathbf{x}_{0i} maximizes $u_{B_{0i}}(\mathbf{x})$. Thus, $Y_i < k$ implies $Y_i = Y_{0i}$. Figure B.4 depicts the logic visually. The proof that $Y_i > k$ implies $Y_i = Y_{1i}$ is analogous.

B.2 Proof of Propositions 1 and 2

Consider Proposition 1. Item i) in the proof of Lemma 1 establishes that under CHOICE and WARP $Y_i = k$ implies $Y_{1i} \leq k \leq Y_{0i}$, since taking contrapositives we have that $(Y_i \geq k \text{ and } Y_i \leq k)$ implies $Y_{1i} \leq k \leq Y_{0i}$. We have also seen from item ii) that under CHOICE and CONVEX $Y_{1i} \leq k \leq Y_{0i}$ also implies $Y_i = k$, thus $Y_{1i} \leq k \leq Y_{0i}$ and $Y_i = k$ are equivalent. Note that by adding $\Delta_{0i} = Y_{0i} - Y_{1i}$ to both sides of the inequality $Y_{1i} \leq k$, we have that $Y_{0i} \leq k + \Delta_{0i}$. Combining with the other inequality that $Y_{0i} \geq k$, we can thus rewrite the event $Y_{1i} \leq k \leq Y_{0i}$ as $Y_{0i} \in [k, k + \Delta_{0i}]$ (or equivalently to $Y_{1i} \in [k - \Delta_{0i}, k]$). We thus have that $\mathcal{B} \leq P(Y_{0i} \in [k, k + \Delta_{0i}])$ under CHOICE and WARP, and that $\mathcal{B} = P(Y_i = k) = P(Y_{1i} \leq k \leq Y_{0i})$ under CHOICE and CONVEX.

Now consider Proposition 2. By item i) in the proof of Proposition 1, we have that under WARP and CHOICE $Y_{0i} \leq k \implies Y_i = Y_{0i}$. Thus, for any $\delta > 0$ and $y < k$: $Y_{0i} \in [y - \delta, y]$ implies

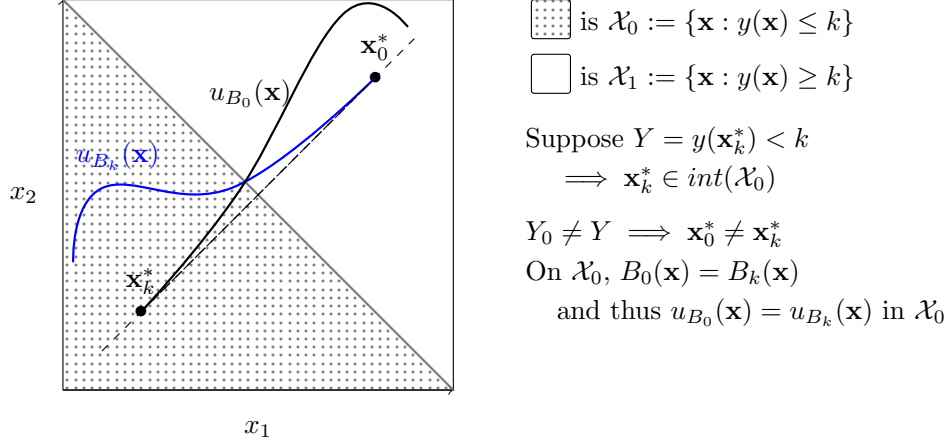


FIGURE B.4: Depiction of the step establishing $(Y < k) \implies (Y = Y_0)$ in the proof of Lemma 1. In this example $z = (x_1, x_2)$ and $y(\mathbf{x}) = x_1 + x_2$. We suppress indices i for clarity. Proof is by contradiction. If $Y_0 \neq Y$, then $\mathbf{x}_k^* \neq \mathbf{x}_0^*$, where \mathbf{x}_k^* and \mathbf{x}_0^* are the unique maximizers of $u_{B_k}(\mathbf{x})$ and $u_{B_0}(\mathbf{x})$, respectively. By Equation B.5, we have that the function $u_{B_0}(\mathbf{x})$, depicted heuristically as a solid black curve, is strictly increasing as one moves along the dotted line from \mathbf{x}_k^* towards \mathbf{x}_0^* . Similarly, the function $u_{B_k}(\mathbf{x})$, depicted as a solid blue curve, is strictly increasing as one moves in the opposite direction along the same line, from \mathbf{x}_0^* towards \mathbf{x}_k^* . By the assumption that $Y < k$, then using continuity of $y(\mathbf{x})$ it must be the case that \mathbf{x}_k^* lies in the interior of \mathcal{X}_0 , the set of \mathbf{x} 's that make $y(\mathbf{x}) \leq k$. This means that there is some interval of the dotted line that is within \mathcal{X}_0 . On this interval, the functions B_0 and B_k are equal, and thus so must be the functions u_{B_k} and u_{B_0} . Since the same function cannot be both strictly increasing and strictly decreasing, we have obtained a contradiction.

that $Y_i \in [y - \delta, y]$ and hence $P(Y_{0i} \in [y - \delta, y]) \leq P(Y_i \in [y - \delta, y])$. This implies that $f_0(y) - f(y) = \lim_{\delta \downarrow 0} \frac{P(Y_{0i} \in [y - \delta, y]) - P(Y_i \in [y - \delta, y])}{\delta} \leq 0$, i.e. that $f(y) \geq f_0(y)$. An analogous argument holds for Y_1 , where we consider the event $Y_{1i} \in [y, y + \delta]$ any $y > k$. Under CONVEX instead of WARP, the inequalities are all equalities, by Lemma 1.

B.3 Proof of Theorem 1

Theorem 1 of Dömbgen et al. (2017) gives a characterization of bi-log concavity in terms of a random variable's CDF and its density. In our case this reads as follows: for $d \in \{0, 1\}$ and any t ,

$$1 - (1 - F_{d|K^*=0}(k))e^{-\frac{f_{d|K^*=0}(k)}{1 - F_{d|K^*=0}(k)}t} \leq F_{d|K^*=0}(k + t) \leq F_{d|K^*=0}(k)e^{\frac{f_{d|K^*=0}(k)}{F_{d|K^*=0}(k)}t}$$

Defining $u = F_{0|K^*=0}(k+t)$, we can use the substitution $t = Q_{0|K^*=0}(u) - k$ to translate the above into bounds on the conditional quantile function of Y_{0i} , evaluated at u :

$$\frac{F_{0|K^*=0}(k)}{f_{0|K^*=0}(k)} \cdot \ln \left(\frac{u}{F_{0|K^*=0}(k)} \right) \leq Q_{0|K^*=0}(u) - k \leq -\frac{1 - F_{0|K^*=0}(k)}{f_{0|K^*=0}(k)} \cdot \ln \left(\frac{1 - u}{1 - F_{0|K^*=0}(k)} \right)$$

And similarly for Y_1 , letting $v = F_{1|K^*=0}(k-t)$:

$$\frac{1 - F_{1|K^*=0}(k)}{f_{1|K^*=0}(k)} \cdot \ln \left(\frac{1 - v}{1 - F_{1|K^*=0}(k)} \right) \leq k - Q_{1|K^*=0}(v) \leq -\frac{F_{1|K^*=0}(k)}{f_{1|K^*=0}(k)} \cdot \ln \left(\frac{v}{F_{1|K^*=0}(k)} \right)$$

Note that:

$$\begin{aligned} E[Y_{0i} - Y_{1i} | Y_i = k, K_i^* = 0] &= \frac{1}{\mathcal{B}^*} \int_{F_{0|K^*=0}(k)}^{F_{0|K^*=0}(k) + \mathcal{B}^*} \{Q_{0|K^*=0}(u) - Q_{0|K^*=0}(u)\} du \\ &= \frac{1}{\mathcal{B}^*} \int_{F_{0|K^*=0}(k)}^{F_{0|K^*=0}(k) + \mathcal{B}^*} \{Q_{0|K^*=0}(u) - k\} du + \frac{1}{\mathcal{B}^*} \int_{F_{1|K^*=0}(k) - \mathcal{B}^*}^{F_{1|K^*=0}(k)} \{k - Q_{1|K^*=0}(v)\} dv \end{aligned}$$

where $\mathcal{B}^* := P(h_{it} = k | K^* = 0)$. A lower bound for $E[Y_{0i} - Y_{1i} | Y_i = k, K_i^* = 0]$ is thus:

$$\begin{aligned} &\frac{F_{0|K^*=0}(k)}{f_{0|K^*=0}(k) \cdot \mathcal{B}^*} \int_{F_{0|K^*=0}(k)}^{F_{0|K^*=0}(k) + \mathcal{B}^*} \ln \left(\frac{u}{F_{0|K^*=0}(k)} \right) du + \frac{1 - F_{1|K^*=0}(k)}{f_{1|K^*=0}(k) \cdot \mathcal{B}^*} \int_{F_{1|K^*=0}(k) - \mathcal{B}^*}^{F_{1|K^*=0}(k)} \ln \left(\frac{1 - v}{1 - F_{1|K^*=0}(k)} \right) dv \\ &= g(F_{0|K^*=0}(k), f_{0|K^*=0}(k), \mathcal{B}^*) + h(F_{1|K^*=0}(k), f_{1|K^*=0}(k), \mathcal{B}^*) \end{aligned}$$

where

$$\begin{aligned} g(a, b, x) &:= \frac{a}{bx} \int_a^{a+x} \ln \left(\frac{u}{a} \right) du = \frac{a^2}{bx} \int_1^{1+\frac{x}{a}} \ln(u) du \\ &= \frac{a^2}{bx} \{u \ln(u) - u\} \Big|_1^{1+\frac{x}{a}} = \frac{a^2}{bx} \left\{ \left(1 + \frac{x}{a}\right) \ln \left(1 + \frac{x}{a}\right) - \frac{x}{a} \right\} \\ &= \frac{a}{bx} (a+x) \ln \left(1 + \frac{x}{a}\right) - \frac{a}{b} \end{aligned}$$

and

$$h(a, b, x) := \frac{1-a}{bx} \int_{a-x}^a \ln \left(\frac{1-v}{1-a} \right) dv = \frac{(1-a)^2}{bx} \int_1^{1+\frac{x}{1-a}} \ln(u) du = g(1-a, b, x)$$

Similarly, an upper bound is:

$$\begin{aligned}
& -\frac{1 - F_{0|K^*=0}(k)}{f_{0|K^*=0}(k)(\mathcal{B}^*)} \int_{F_{0|K^*=0}(k)}^{F_{0|K^*=0}(k) + \mathcal{B}^*} \ln \left(\frac{1 - u}{1 - F_{0|K^*=0}(k)} \right) du \\
& \quad - \frac{F_{1|K^*=0}(k)}{f_{1|K^*=0}(k)(\mathcal{B}^*)} \int_{F_{1|K^*=0}(k) - \mathcal{B}^*}^{F_{1|K^*=0}(k)} \ln \left(\frac{v}{F_{1|K^*=0}(k)} \right) dv \\
& = \tilde{g}(F_{0|K^*=0}(k), f_{0|K^*=0}(k), \mathcal{B}^*) + \tilde{h}(F_{1|K^*=0}(k), f_{1|K^*=0}(k), \mathcal{B}^*)
\end{aligned}$$

where

$$\begin{aligned}
\tilde{g}(a, b, x) &:= -\frac{1 - a}{bx} \int_a^{a+x} \ln \left(\frac{1 - u}{1 - a} \right) du = -\frac{(1 - a)^2}{bx} \int_{1 - \frac{x}{1 - a}}^1 \ln(u) du \\
&= \frac{(1 - a)^2}{bx} \{u - u \ln(u)\} \Big|_{1 - \frac{x}{1 - a}}^1 = \frac{1 - a}{b} + \frac{1 - a}{bx} (1 - a - x) \ln \left(1 - \frac{x}{1 - a} \right) \\
&= -g(1 - a, b, -x)
\end{aligned}$$

and

$$\tilde{h}(a, b, x) := -\frac{a}{bx} \int_{a-x}^a \ln \left(\frac{v}{a} \right) dv = -\frac{a^2}{bx} \int_{1 - \frac{x}{a}}^1 \ln(u) du = \tilde{g}(1 - a, b, x) = -g(a, b, -x)$$

Given p , we relate the $K^* = 0$ conditional quantities to their unconditional analogues:

$$\begin{aligned}
F_{0|K^*=0}(k) &= \frac{F_0(k) - p}{1 - p} \quad \text{and} \quad F_{1|K^*=0}(k) = \frac{F_1(k) - p}{1 - p} \quad \text{and} \quad \mathcal{B}^* = \frac{\mathcal{B} - p}{1 - p} \\
f_{0|K^*=0}(k) &= \frac{f_0(k)}{1 - p} \quad \text{and} \quad f_{1|K^*=0}(k) = \frac{f_1(k)}{1 - p}
\end{aligned}$$

Let $F(h) = P(h_{it} \leq h)$ be the CDF of the data, and define $f(h) = \frac{d}{dh} P(h_{it} \leq h)$ for $h \neq k$. By Lemma 1 and the BLC assumption, the above quantities are related to observables as:

$$F_0(k) = \lim_{h \uparrow k} F(h) + p, \quad F_1(k) = F(k), \quad f_0(k) = \lim_{h \uparrow k} f(h), \quad \text{and} \quad f_1(k) = \lim_{h \downarrow k} f(h)$$

As shown by Dümbgen et al. (2017), BLC implies the existence of a continuous density function, which assures that the required density limits exist, and delivers Item 1. of Theorem 1.

To obtain the final result, note that the function $g(a, b, x)$ is homogeneous of degree zero. Thus $\Delta_k^* \in [\Delta_k^L, \Delta_k^U :]$, with

$$\Delta_k^L := g(F_-(k), f_-(k), \mathcal{B} - p) + g(1 - F(k), f_+(k), \mathcal{B} - p)$$

$$\Delta_k^U := -g(1 - p - F_-(k), f_-(k), p - \mathcal{B}) - g(F(k) - p, f_+(k), p - \mathcal{B})$$

where $-$ and $+$ subscripts denote left and right limits. The bounds are sharp as CHOICE, CONVEX and RANK imply no further restrictions on the potential outcome distributions.

B.4 Proof of Theorem 2

This proof follows the notation of Appendix A. Throughout this proof we let $Y_i(\rho, k) = Y_i(\rho)$, given Assumption SEPARABLE. By Lemmas 2 and 3 the effect of changing k on bunching is:

$$\begin{aligned} \partial_k \{\mathcal{B} - p(k)\} &= -\frac{\partial}{\partial k} \int_{\rho_0}^{\rho_1} f_\rho(k) \mathbb{E} \left[\frac{Y_i(\rho)}{d\rho} \middle| Y_i(\rho) = k \right] d\rho \\ &= -\int_{\rho_0}^{\rho_1} \frac{\partial}{\partial k} \left\{ f_\rho(k) \mathbb{E} \left[\frac{Y_i(\rho)}{d\rho} \middle| Y_i(\rho) = k \right] \right\} d\rho = \int_{\rho_0}^{\rho_1} \partial_\rho f_\rho(k) d\rho = f_1(k) - f_0(k) \end{aligned}$$

Turning now to the total effect on average hours.

$$\begin{aligned} \partial_k E[Y_i^{[k, \rho_1]}] &= \partial_k \{P(Y_i(\rho_0) < k) \mathbb{E}[Y_i(\rho_0) | Y_i(\rho_0) < k]\} + k \partial_k (\mathcal{B}^{[k, \rho_1]} - p(k)) + \mathcal{B}^{[k, \rho_1]} - p(k) \\ &\quad + \partial_k \{P(Y_i(\rho_1) > k) \mathbb{E}[Y_i(\rho_1) | Y_i(\rho_1) > k]\} \\ &= \partial_k \int_{-\infty}^k y \cdot f_{\rho_0}(y) \cdot dy + k(f_0(k) - f_1(k)) + \mathcal{B}^{[k, \rho_1]} - p(k) + \partial_k \int_k^\infty y \cdot f_{\rho_1}(y) \cdot dy \\ &= \cancel{k f_0(k)} + \cancel{k(f_1(k) - f_0(k))} + \mathcal{B}^{[k, \rho_1]} - p(k) - \cancel{k f_1(k)} \end{aligned}$$

Meanwhile:

$$\begin{aligned} \partial_{\rho_1} E[Y_i^{[k, \rho_1]}] &= k \partial_{\rho_1} \mathcal{B}^{[k, \rho_1]} + \partial_{\rho_1} \{P(Y_i(\rho_1) > k) \mathbb{E}[Y_i(\rho_1) | Y_i(\rho_1) > k]\} = k \partial_{\rho_1} \mathcal{B}^{[k, \rho_1]} + \int_k^\infty y \cdot \partial_{\rho_1} f_{\rho_1}(y) \cdot dy \\ &= -k f_{\rho_1}(k) \mathbb{E} \left[\frac{Y_i(\rho_1)}{d\rho} \middle| Y_i(\rho_1) = k \right] - \int_k^\infty y \cdot \partial_y \left\{ f_{\rho_1}(y) \mathbb{E} \left[\frac{dY_i(\rho_1)}{d\rho} \middle| Y_i(\rho_1) = y \right] \right\} dy \\ &= \cancel{-k f_{\rho_1}(k) \mathbb{E} \left[\frac{Y_i(\rho_1)}{d\rho} \middle| Y_i(\rho_1) = k \right]} + \cancel{y f_{\rho_1}(y) \mathbb{E} \left[\frac{dY_i(\rho_1)}{d\rho} \middle| Y_i(\rho_1) = y \right] \Big|_k^\infty} \\ &\quad - \int_k^\infty f_{\rho_1}(y) \mathbb{E} \left[\frac{dY_i(\rho_1)}{d\rho} \middle| Y_i(\rho_1) = y \right] dy \end{aligned}$$

where I have used Lemma 2 with the Leibniz rule (establishing Item 3 in Theorem 2) as well as Lemma 3 in the third step, and then integration by parts along with the boundary condition that $\lim_{y \rightarrow \infty} y \cdot f_{\rho_1}(y) = 0$, implied by Assumption SMOOTH.