

# Efficient Welfare Weights

Nathaniel Hendren\*

December, 2019

## Abstract

This paper provides a set of efficient welfare weights that modify the classic Kaldor-Hicks experiments that define economic efficiency to account for the distortionary cost of transfers. The weights are proportional to the marginal cost to the government of providing a \$1 transfer to a given income level, and are also known in previous literature as the implicit social preferences that rationalize the status quo tax schedule as optimal. I show that the Kaldor-Hicks Pareto-based experiments motivate their use even if one disagrees with the shape of the tax schedule. I estimate the weights using the universe of US income tax returns from 2012. The shape of the income distribution implies that it is efficient to weight surplus to the poor more than to the rich. Calibrations suggest the poor should receive 1.5-2x more weight than the rich: \$1 of surplus to the poor can be turned into \$1.5-\$2 of surplus to the rich through reductions in distortionary taxation. Using the weights to compare income distributions, I show that weighted income growth since 1979 is 15-20% lower than unweighted income growth. Extrapolations imply a social cost of increased income inequality since 1979 of roughly \$400B. Across countries, the U.S. is poorer than countries like Austria and the Netherlands, despite having higher national income per capita. I conclude with a simple welfare framework for assessing the desirability of economic policy changes that relies on the potential Pareto principle instead of a social welfare function.

## 1 Introduction

Comparisons across economic environments involve comparisons between winners and losers. Normative preferences over economic environments therefore require taking a stand on how to weigh the gains to the winners against the losses to the losers.

---

\*Harvard University, nhendren@fas.harvard.edu. This paper is a revised version of a paper that previously circulated under the title, "The Inequality Deflator: Interpersonal Comparisons without a Social Welfare Function", NBER WP No. 20351. I am deeply indebted to conversations with Louis Kaplow for the inspiration behind this paper, and to Sarah Abraham, Alex Bell, Alex Olssen, and Evan Storms for excellent research assistance. I also thank Daron Acemoglu, Raj Chetty, Amy Finkelstein, Ben Lockwood, Henrik Kleven, Patrick Kline, Kory Kroft, Matthew Notowodigdo, Jim Poterba, Emmanuel Saez, Matthew Weinzierl, Glen Weyl, Ivan Werning, and Floris Zoutman, along with seminar participants at Berkeley, Harvard, MIT, Michigan, and Stanford for very helpful comments. The opinions expressed in this paper are those of the author alone and do not necessarily reflect the views of the Internal Revenue Service or the U.S. Treasury Department. This work is a component of a larger project examining the effects of tax expenditures on the budget deficit and economic activity, and this paper in particular provides a general characterization of the welfare impact of changes in tax expenditures relative to changes in tax rates (illustrated in Section 6). The empirical results derived from tax data that are reported in this paper are drawn from the SOI Working Paper "The Economic Impacts of Tax Expenditures: Evidence from Spatial Variation across the U.S.", approved under IRS contract TIRNO-12-P-00374.

Kaldor and Hicks classic definition of economic efficiency was originally intended to resolve this tradeoff. They proposed constructing a sum of individuals' willingnesses to pay for the alternative environment. If this sum is negative, a set of individual-specific lump-sum transfers could reach an allocation that is *Pareto* superior to the alternative environment (Hicks (1940)). Everyone would prefer the transfers in the status quo relative to the alternative environment. Conversely, if the sum is positive, Kaldor (1939) argues that the winners could compensate the losers using individual-specific lump-sum transfers. Therefore, an alternative environment that is modified to include compensating transfers would generate a Pareto improvement relative to the status quo. As a result, the *Pareto* principle implies the alternative should be preferred if and only if the sum of willingness to pay for the alternative environment is positive. These two Pareto-based conceptual experiments underpin the standard definition of economic efficiency as the unweighted sum of individuals' willingnesses to pay.

Despite the aims of Kaldor and Hicks, measuring economic efficiency has become synonymous with a lack of consideration for interpersonal comparisons. Instead, researchers often employ a social welfare function to resolve the "equity-efficiency tradeoff". In doing so, one abandons the ability to provide normative guidance without relying on the subjective preferences of the researcher or policymaker.

This paper revisits these classic Kaldor-Hicks experiments. I argue that the criticism that Kaldor-Hicks does not account for equity concerns is related to the fact that the hypothetical policy changes in their conceptual experiments are infeasible. This is because the redistributive experiments envisioned by Kaldor and Hicks involve distortionary costs. Taxes are imposed on observable choices like incomes that respond to taxes and transfers (Mirrlees (1971)). Correctly measuring Kaldor and Hicks' original notion of economic efficiency requires accounting for these costs:

*"Since almost every conceivable kind of compensation (re-arrangement of taxation, for example) must itself be expected to have some influence on production, the task of the welfare economist is not completed until he has envisaged the total effects...If, as will often happen, the best methods of compensation feasible involve some loss in productive efficiency, this loss will have to be taken into account"* (Hicks (1939), p712).

This paper provides a set of "efficient" welfare weights that implement the Kaldor-Hicks tests for efficiency in a manner that accounts for the distortionary cost of redistribution through the tax schedule. The efficient social welfare weight at an income level  $y$  is equal to the marginal cost of providing \$1 of welfare to individuals earning near  $y$ . In a world without distortionary taxation, this cost is \$1. But, it differs from \$1 when individuals choose to change their earnings in response to a \$1 tax cut to those earning near  $y$ . For example, those earning below  $y$  might increase their incomes to  $y$ ; those above  $y$  may decrease their incomes to  $y$ . By the envelope theorem, those who change their earnings will not obtain a first-order utility improvement from the transfer; but, in the presence of taxes/transfers these responses have a first-order impact on government revenue. If taxes are positive (negative), increases in incomes create positive (negative) fiscal externalities. The efficient welfare weights differ from unity because of these fiscal externalities from redistributive taxes and transfers.

Weighting individuals' willingnesses to pay by efficient welfare weights searches for potential Pareto improvements in the spirit of Kaldor and Hicks' tests for economic efficiency. If weighted surplus is

negative, then modifications to the income tax schedule can lead to all points of the income distribution being better off than would be the case under the alternative environment. In this sense, the economist can suggest a better way to obtain the distribution of welfare gains offered by the alternative environment using a simple modification to the income tax schedule. Conversely, if the weighted surplus is positive, then a modified version of the alternative environment that redistributes the gains of the winners using feasible modifications to the tax schedule could make everyone better off. In this sense, the Pareto principle suggests preferring the alternative environment if and only if efficient-weighted surplus is positive.

Because the efficient welfare weights equal the marginal cost of providing \$1 of income to a given income level, they also reveal the implicit social preferences that rationalize the tax schedule as optimal. These “inverse-optimal” welfare weights have been derived and analyzed by a large and growing literature (see, e.g., Christiansen (1977); Christiansen and Jansen (1978); Blundell et al. (2009); Bargain et al. (2011); Bourguignon and Spadaro (2012); Lockwood and Weinzierl (2016); Zoutman et al. (2013); Bargain et al. (2014); Jacobs et al. (2017)). The core idea in this paper is that these implicit weights can be used to test for Kaldor-Hicks efficiency. This provides a normative justification for their use even if one personally does not agree with the social preferences that rationalize the status quo tax schedule as optimal. Weighting surplus by these weights searches for potential Pareto improvements.

What do the efficient welfare weights look like? I leverage a recent derivation provided in Jacobs et al. (2017) that shows the impact of the behavioral response to taxation on the government budget can be expressed as a function of (a) the joint distribution of taxable income and marginal tax rates and (b) a set of behavioral elasticities governing the response of income to changes in taxation.<sup>1</sup> I use the universe of US income tax returns from 2012 to estimate this joint distribution, and I begin by providing bounds on the efficient welfare weights (without assuming a magnitude of the behavioral response to taxation). I show that the shape of the income distribution - in particular the local Pareto parameter of the income distribution - plays a key role in determining the extent to which the weights are above or below 1.<sup>2</sup> The Pareto parameter rises from near -1 at the bottom of the income distribution to near 2 at the top of the income distribution, crossing zero around the 60th quantile of the income distribution (around \$43K in ordinary income). This means that weights are generally above one for those with incomes below \$43K, and below one for those with incomes above \$43K. Regardless of the size of the behavioral response to taxation, it is more costly to provide \$1 to the poor than to the rich. Thus, these bounds suggest it is efficient to weight surplus to the poor more than to the rich. Intuitively, it is more costly to move an additional \$1 from the top to the bottom of the income distribution through additional redistribution than it is to move \$1 from the bottom to the top of the income distribution through reduced redistribution.

Next, I construct point estimates using existing estimates of taxable income elasticities. The baseline specification suggests a \$1 tax cut to those with high incomes from a reduction in marginal

---

<sup>1</sup>I extend this result to allow for the presence of multiple tax schedule for those at the same income level, as occurs in the US.

<sup>2</sup>The Pareto parameter is given by  $-\left(1 + \frac{yf'(y)}{f(y)}\right)$  where  $f(y)$  is the density of the income distribution.

tax rates costs around \$0.65. At the other end of the income distribution, the estimates suggest that expansions of the earned income tax credit (EITC) by \$1 to low earners has a fiscal cost of around \$1.15 because additional transfers cause individuals to adjust their earnings to maximize their tax credits. Combining, efficient welfare weights decline from around 1.15 at the bottom of the income distribution to around 0.65 at the top. In other words, \$1 to the poor can be turned into \$1.5-2 to the rich through modifications to the tax schedule. As a result, it is efficient to value surplus to the poor 1.5-2x more than surplus to the rich.

Motivated by the original application to the comparison of income distributions in the work of Kaldor and Hicks, I apply the weights to two sets of comparisons of income distributions. First, I construct distributionally-adjusted measures of economic growth in the US. As is widely documented, growth in the US has been unequal across the income distribution. Because it is costly to redistribute from rich to poor, distributionally-adjusted measures of economic growth are 15-20% lower. Extrapolating across all economic growth between 1979 and 2012 suggests an increase in distributionally-adjusted growth of \$15K, in contrast to aggregate growth of \$18K. Multiplying by 119M households in the US suggests a social cost of increased income inequality in the US since 1979 of roughly \$400B.

Second, I compare the distribution of incomes across countries. Broadly, differences in aggregate income tend to accurately order comparisons across countries. But, there are several exceptions. Most notably, the income distributions of Austria and New Zealand would be preferred relative to the US income distribution, despite having a lower mean per capita income. Although the US has higher mean incomes, it is unable to replicate the distribution of income offered in those countries through modifications in the tax schedule.

Lastly, I nest the efficient welfare weights into a welfare framework that analyzes whether a policy change provides a potential Pareto improvement. Under conditions that restrict the degree of heterogeneity in the policy's impact, the framework motivates comparing a policy to a distributionally-equivalent tax cut. I discuss an empirical method to test for this by constructing a policy's marginal value of public funds (MVPF) as in Hendren (2016) and Hendren and Sprung-Keyser (2019). The resulting normative conclusions rely solely on the positive analyses of economic policies combined with the Pareto principle, as opposed to a social welfare function.

**Relation to Previous Literature** This is not the first paper to recognize that incorporating the distortionary cost of transfers leads to a modification of the Kaldor and Hicks compensation principle. The theoretical analysis in this paper is most closely related to Coate (2000), who proposes an approach that incorporates the costs of redistribution into the Hicks criterion by comparing the policy to feasible alternatives such as distortionary redistribution through the tax schedule. The analysis is also related to the methods discussed in Christiansen (1981) and Kaplow (2004), among others that suggest to compensate losers of policy changes through the tax schedule. Relative to this literature, the contribution of this paper is to provide a set of welfare weights that generate first-order implementation of this approach.<sup>3</sup>

---

<sup>3</sup>Coate (2000) writes: “An interesting problem for further research would be to investigate whether the efficiency approach might be approximately decentralised via a system of shadow prices which convey the cost of redistributing

As noted above, the efficient welfare weights are equivalent to the solution to the “inverse optimum” program in optimal tax: they reveal the implicit preferences of those who are indifferent to modifications in the tax schedule.<sup>4</sup> In this sense, the estimates suggest those indifferent to the current shape of the tax schedule in the US implicitly value an additional \$1 to the rich as equivalent to an additional \$1.5-2 to the poor.<sup>5</sup> The contribution of this paper is to justify the normative use of these weights even if one does not agree with the implicit social preferences that rationalize the tax schedule as optimal. The Kaldor-Hicks insight motivates their use by appealing to the Pareto principle. In this sense, the inverse-optimum weights studied in earlier literature have a use beyond assessing implicit social preferences embodied in the tax schedule: they can be used as efficient welfare weights to search for potential Pareto improvements for policies or comparisons across economic environments.

The rest of this paper proceeds as follows. Section 2 derives the efficient welfare weights in the context of a general model setup. Section 3 illustrates how the weights implement the modified Kaldor-Hicks efficiency experiments. Section 4 relates the weights to the solution to the inverse optimum program in optimal tax. Section 5 illustrates a representation of these weights using the distribution of income, tax rates, and behavioral elasticities. Section 6 provides estimates of the joint distribution of income and tax rates and discusses bounds on the shape of the efficient welfare weights. Section 7 provides point estimates by calibrating behavioral elasticities. Section 8 applies the weights to the comparison of income distributions over time in the US and across countries. Section 9 discusses a simple normative policy evaluation framework, and Section 10 concludes.

## 2 Model

I consider an economy with a population of agents, indexed by  $\theta$ , with population size normalized to 1. There is a status quo environment and an alternative environment. The alternative environment could be a world with greater spending on a public good, a more progressive tax schedule, or the distribution of income offered by another country. The model developed here will be used to derive the tests for whether the alternative environment provides a potential Pareto improvement when the transfers envisioned by Kaldor and Hicks are conducted through modifications to the tax schedule.

The model will be used to define two key variables that will be important for implementing the Kaldor-Hicks tests for efficiency. First, I use the model to define each person’s willingness to pay for an alternative environment,  $s(\theta)$ . And, I define the marginal cost to the government of providing a \$1 mechanical tax cut to those with incomes near  $y$  as  $g(y)$ . Foreshadowing the results discussed in Section 3, I label  $g(y)$  to be efficient welfare weights because they implement the test for economic efficiency. As noted in the introduction, they are also equal to the solution to the inverse optimum program in the optimal nonlinear income tax problem.

---

*between different types of citizens.*” The efficient welfare weights are the shadow prices envisioned by Coate (2000).

<sup>4</sup>See references noted above (Christiansen (1977); Christiansen and Jansen (1978); Blundell et al. (2009); Bargain et al. (2011); Bourguignon and Spadaro (2012); Bargain et al. (2014); Zoutman et al. (2013); Lockwood and Weinzierl (2016); Jacobs et al. (2017)).

<sup>5</sup>This is consistent with Lockwood and Weinzierl (2016), who estimate the solution to the inverse optimum program in the U.S. using aggregated data from the Congressional Budget Office.

The model builds heavily upon this earlier literature in optimal taxation and in particular the work of Jacobs et al. (2017). The main extension relative to earlier work is to allow for unrestricted heterogeneity across individuals. This is important when considering the Kaldor-Hicks experiments because it allows one to consider the potential case when there is not a one-to-one relationship between willingness to pay for an alternative environment and one’s income.<sup>6</sup>

## 2.1 Willingness to Pay

In the status quo environment, agents consumption,  $c$ , and earnings,  $y$ . I allow each agent of type  $\theta$  to have a potentially different utility function,  $u(c, y; \theta)$ , over consumption and earnings. I will not impose restrictions on the distribution of  $\theta$ , so without loss of generality one can think of  $\theta$  as simply indexing people in the population.

Agents choose  $c$  and  $y$  to maximize utility subject to a budget constraint,

$$c \leq y - T(y) + m$$

where  $T(y)$  is the taxes paid on earnings  $y$  and  $m$  is additional income beyond earnings.<sup>7</sup> With a slight abuse of notation, I let  $c(\theta)$  and  $y(\theta)$  denote the resulting choices of type  $\theta$ .

Let  $v^0(\theta)$  denote the utility level obtained by type  $\theta$  in the status quo environment. And, given a utility level  $v$ , define the expenditure function  $e(v; \theta)$  to be the smallest value of  $m$  that is required for a type  $\theta$  to obtain utility level  $v$  in the status quo environment.<sup>8</sup>

Let  $u^a(c, y; \theta)$  denote the utility function for type  $\theta$  in the alternative environment. I do not restrict any feature of the alternative environment – it could contain different wage distributions, better schools, less traffic, better restaurants, or simply different scenery – any of which can affect the level of  $u^a$  for any individual  $\theta$ . I also do not restrict that the tax schedule in the alternative environment be the same as the status quo. To that aim, let  $T^a(\circ)$  denote the tax schedule in the alternative environment so that the budget constraint is given by  $c \leq y - T^a(y) + m$ . Define  $v^a(\theta)$  to be the level of utility obtained and  $e^a(v; \theta)$  is the smallest value of  $m$  that is required for a type  $\theta$  to obtain utility level  $v$  in the alternative environment.<sup>9</sup> Individual  $\theta$ ’s willingness to pay (equivalent variation) for the alternative environment is then given by

$$s(\theta) = e(v^a(\theta); \theta) - e(v^0(\theta); \theta) \tag{1}$$

This is the amount of additional money a type  $\theta$  would need in the status quo to be as well off as in

---

<sup>6</sup>If two people with the same income have different willingnesses to pay for the alternative environment, then it will be difficult for the tax schedule modifications to implement the transfers envisioned by Kaldor and Hicks – the tax schedule will be a blunt instrument in attempting to make compensatory transfers. By allowing for arbitrary heterogeneity across individuals (beyond just their earnings capacities), the model is able characterize when this arises as problem. I discuss this further in Section 3.4 and in greater detail in Appendix D.

<sup>7</sup>For simplicity, I assume  $T(y)$  is the same for everyone. In the empirical implementation, I allow  $T$  to vary with individual characteristics, such as the number of dependents, and marital status. See Section E.1.

<sup>8</sup>Formally,  $e(v; \theta) = \inf \{m | \sup_{c, y} \{u(c, y; \theta) | c \leq y - T(y) + m\} \geq v\}$ . Duality implies that  $e(v^0(\theta); \theta) = m$ .

<sup>9</sup>Analogous to the definition in the status quo environment,  $e^a(v; \theta) = \inf \{m | \sup \{u^a(c, y; \theta) | c \leq y - T^a(y) + m\} \geq v\}$ .

the alternative environment.<sup>10</sup>

The simplest case of the analysis that follows is when the willingness to pay for the alternative environment,  $s(\theta)$ , is homogeneous in income. With an abuse of notation in this case, I let  $s(y)$  denote the willingness to pay for the alternative environment by a type  $\theta$  who chooses income  $y(\theta)$  in the status quo. Going forward, I consider this simpler case, and return to a discussion of its implications in Section 3.4 and Appendix D.

Given  $s(y)$ , the goal is to answer two questions: (1) can the surplus,  $s(y)$ , can be replicated through modifications in the tax schedule in the status quo environment (i.e. the experiment in Hicks (1940))? And, (2) does there exist a modification to the tax schedule in the alternative environment that makes everyone better off relative to the status quo (i.e. the experiment in Kaldor (1939))? The answer to these questions will depend on how changes to the tax schedule affect government revenue.

## 2.2 Marginal Cost of Taxation

The model of the marginal cost of taxation in this environment builds upon the recent work in optimal taxation and the inverse optimum program. Given the status quo tax schedule,  $T(y)$ , consider a tax cut of \$1 that is targeted to those with incomes near  $y = y^*$  in the status quo environment. Figure 1 depicts this tax cut to those with incomes in a bin of size  $\epsilon$  centered around  $y^*$ . The horizontal axis plots pre-tax income and the vertical axis plots after-tax income or consumption. The cost to the government equals not only the \$1 of mechanical spending to those with incomes near  $y^*$ , but also includes the impact of behavioral responses to the policy on the government budget, depicted in blue arrows. By the envelope theorem, those behavioral responses do not affect anyone's well-being: rather, the willingness to pay for this \$1 tax cut equals \$1 by those whose incomes are inside the bin of width  $\epsilon$  of  $y^*$ .

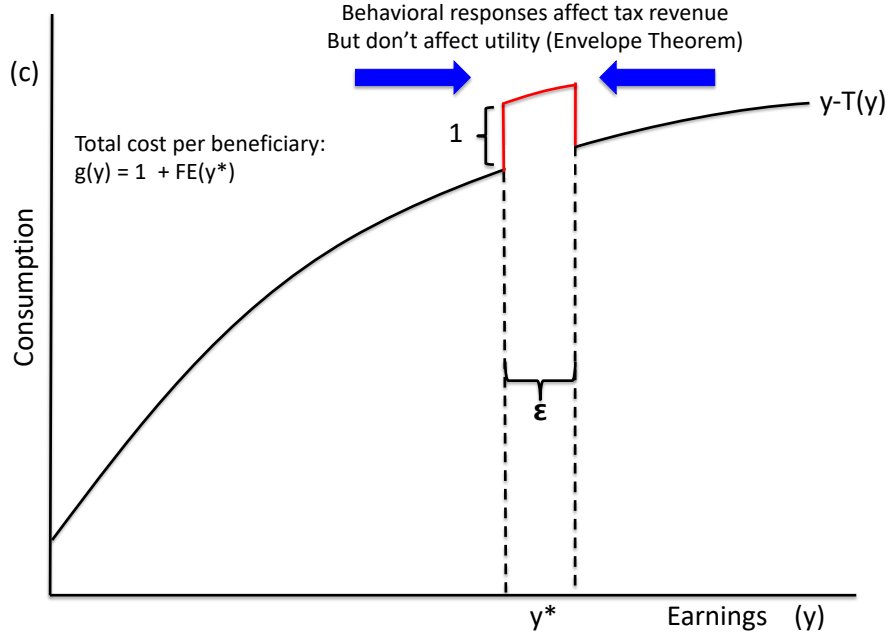
However, the behavioral do affect the cost to the government of the transfer. Let  $FE(y^*)$  denote the impact of these behavioral responses on the government budget. The marginal cost of taxation is then given by  $g(y^*) = 1 + FE(y^*)$ . Appendix A defines this calculus of variations argument and formally defines  $g(y)$  in the context of the model above. Yet the definition of  $g(y)$  is intuitive:  $g(y)$  equals the cost to the government of providing \$1 of a mechanical tax cut to those with incomes near  $y$ .

The next Section shows that one can use the weights  $g(y)$  as welfare weights to search for potential Pareto improvements in the spirit of Kaldor and Hicks.

---

<sup>10</sup>In addition to this equivalent variation definition of willingness to pay, one could also construct a compensating variation measure using the expenditure function in the alternative environment,  $cv(\theta) = e^a(v(\theta); \theta) - e^a(v^a(\theta); \theta)$ . Because the distinction between equivalent and compensating variation is second order (e.g. see Schlee (2013) for a recent discussion of the first-order equivalence of five common conceptualizations of willingness to pay, including compensating and equivalent variation.) and the approach below considers first-order adjustments, it will not be necessary to distinguish between equivalent or compensating variation in the analysis that follows.

Figure 1: Tax Cut to Those Earning Near  $y^*$



Notes: This figure illustrates the modification to the tax schedule that provides a tax cut of \$1 to those with earnings in a region of  $y^*$  of width  $\epsilon$ . To first order, those whose earnings would lie in  $[y^* - \frac{\epsilon}{2}, y^* + \frac{\epsilon}{2}]$  will value the tax cut at \$1. But, the costs will result from both this mechanical cost and the impact of behavioral responses to the tax cut (loosely illustrated by the blue arrows). So, the total cost per unit of mechanical beneficiary will be  $g(y) = 1 + FE(y)$ .

### 3 Using $g(y)$ as Efficient Welfare Weights

To set the stage, let  $s(y)$  denote the willingness to pay for some alternative environment. Is the alternative environment preferred to the status quo? Before outlining the approach suggested here, first consider the social welfare function approach. This would weight willingness to pay by a generalized social welfare weight (Saez and Stantcheva (2016)) that balance the gains to the winners (e.g. for whom  $s(y) > 0$ ) against the losses to the losers (e.g. for whom  $s(y) < 0$ ). In contrast to this approach, it will turn out that the Kaldor-Hicks experiments motivate weighting this willingness to pay by  $g(y)$  equation (8), regardless of one's own social preferences. To this aim, I refer to  $g(y)$  as "Efficient Welfare Weights" and refer to the weighted sum of surplus as "Efficient Surplus". Efficient surplus equals the amount of resources the government needs to replicate the surplus offered by the alternative environment using modifications to the tax schedule,

$$S = E [s(y) g(y)] \quad (2)$$

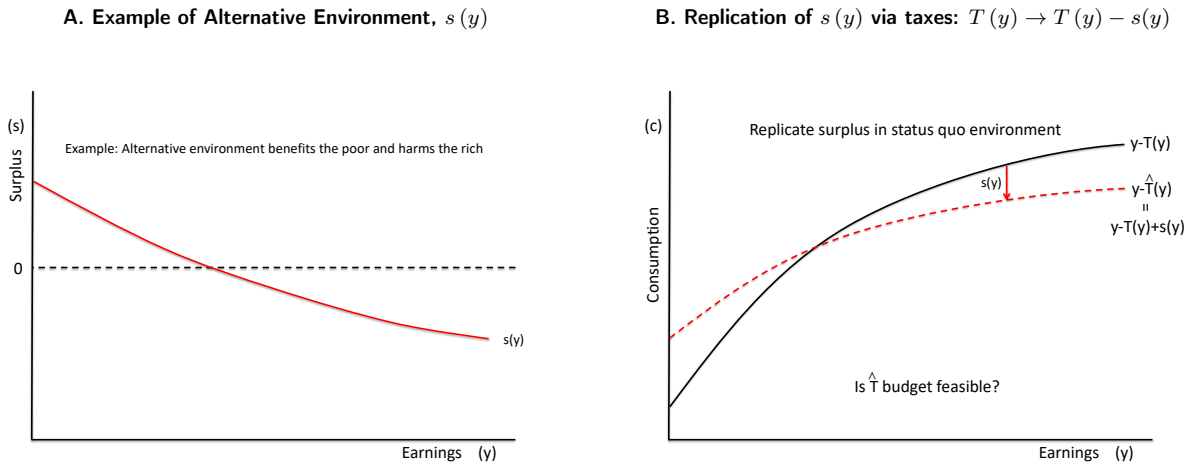
The next subsections illustrate how these weights and surplus measures implement the classic tests for efficiency in Kaldor and Hicks, modified to account for the distortionary cost of redistribution.

Section 3.1 implements a first-order test of efficiency in the spirit of Hicks (1940):  $S > 0$  if and only if one cannot replicate the surplus allocation offered by the alternative environment using modifications to the tax schedule. Section 3.2, implements a first order test of efficiency in the spirit of Kaldor (1939):  $S > 0$  if and only if one can modify the tax schedule in the alternative environment to generate a Pareto superior allocation that is preferred by everyone relative to the status quo. Combined, these tests motivate testing for whether  $S > 0$  to assess the desirability of the alternative environment.

### 3.1 Testing for Efficiency in the Spirit of Hicks (1940) and Coate (2000)

Can the benefits offered by the alternative environment,  $s(y)$ , be more efficiently provided through modifications in the tax schedule? To assess this, imagine replacing the current tax schedule,  $T(y)$ , with a new tax schedule,  $\hat{T}(y) = T(y) - s(y)$ , that offers a tax cut of size  $s(y)$  to those earning  $y$ . Figure 2 provides an illustration. Panel A presents a hypothetical alternative environment that is preferred by the poor but not by the rich. Panel B then modifies the tax schedule from  $T(y)$  to  $T(y) - s(y)$ . To first order, the envelope theorem implies that the tax cut of  $s(y)$  is valued at  $s(y)$  by those earning  $y$ . Therefore, everyone is approximately indifferent between the alternative environment and the status quo environment with the modified tax schedule, as depicted by the dashed red line in Figure 2, Panel B. The Hicks test for efficiency asks: Is this tax modification in the status quo world feasible?

**Figure 2: Hicks (1940) Efficiency Experiment**

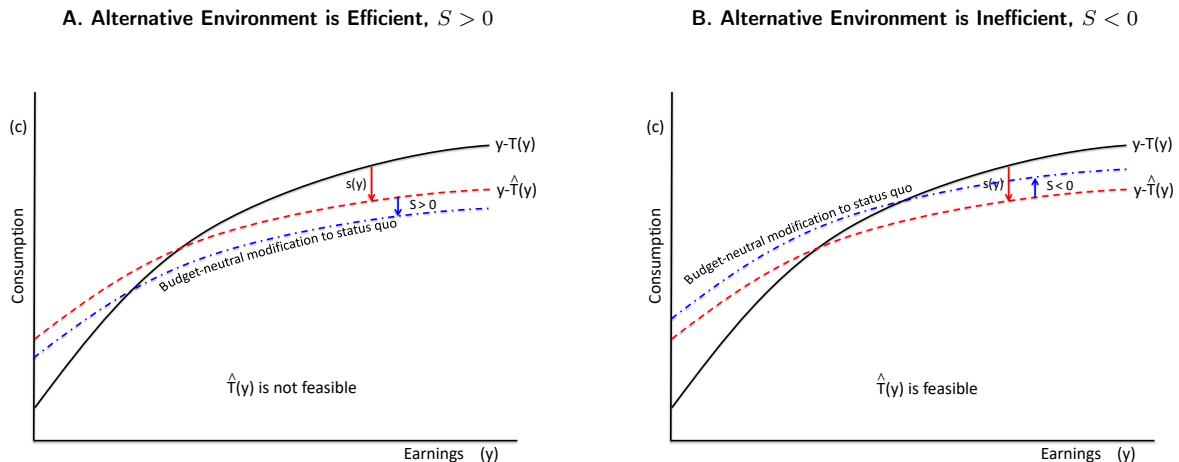


*Notes:* This figure illustrates the efficiency experiment of Hicks (1940) for a hypothetical alternative environment. Panel A presents the hypothetical willingness to pay for each person at different points of the income distribution. In this example, those with low incomes prefer the alternative environment, but those with higher incomes prefer the status quo. Panel B illustrates modifying the tax schedule in the status quo world to attempt to replicate the surplus offered by the alternative environment. To first order, everyone is indifferent between the alternative environment and the modified status quo with tax schedule  $T(y) - s(y)$ .

To first order, the marginal cost of providing \$1 of welfare to those earning  $y$  is given by  $g(y) =$

$1 + FE(y)$ . Therefore, the cost of this tax cut is given by  $E[g(y)s(y)]$ . If this quantity is positive, then providing surplus  $s(y)$  through the tax schedule would not be feasible. Closing the budget constraint by raising taxes on everyone would lead to the blue line in Figure 3, Panel A. In this sense, the alternative environment would be efficient relative to what is feasible through modifications to the tax schedule in the status quo. In contrast, if  $S < 0$ , then it is possible to replicate the alternative environment through modifications to the tax schedule. Redistributing the government surplus to everyone equally leads to the blue line in Figure 3, Panel B, which is preferred by all relative to the alternative environment. In this sense, the alternative environment is efficient if and only if  $S > 0$ .

**Figure 3: Testing for Hicks Efficiency**



*Notes:* This figure illustrates the efficiency test of Hicks (1939). The blue line illustrates the conceptual after-tax income that is feasible through modifications to the tax schedule but has the same distributional incidence as the alternative environment. Panel A illustrates the case in which the modified status quo tax schedule would deliver lower welfare to all points of the income distribution, so that the alternative environment is efficient relative to the status quo,  $S > 0$ . In contrast, Panel B illustrates the case in which replicating the surplus offered by the alternative environment through the tax schedule leads to higher welfare for all, so that the alternative environment is inefficient.

The formal version of these statements are valid up to first order, as they rely on the envelope theorem to ensure indifference between the modified status quo (the dashed red line in Figure 2, Panel B) and the alternative environment. Proposition 1 provides one method of formalizing this idea by considering a scaled surplus function.

**Proposition 1.** *For any  $\epsilon > 0$  define the scaled surplus by  $s_\epsilon(y) = \epsilon s(y)$  and  $S_\epsilon = E[s_\epsilon(y)g(y)] = \epsilon S$ . If  $S < 0$ , there exists an  $\tilde{\epsilon} > 0$  such that for any  $\epsilon < \tilde{\epsilon}$  there exists an augmentation to the tax schedule in the status quo environment that generates surplus,  $s_\epsilon^t(y)$ , that is uniformly greater than the surplus offered by the alternative environment,  $s_\epsilon^t(y) > s_\epsilon(y)$  for all  $y$ . Conversely, if  $S > 0$ , no such  $\tilde{\epsilon}$  exists.*

*Proof.* See Appendix B.2. □

Proposition 1 shows that the conclusions in Figures 2 and 3 hold to first order.<sup>11</sup> In this sense, the efficient welfare weights provide the right direction for adjusting for the marginal cost of redistribution. However, the costs and benefits of redistribution through the tax schedule could differ for larger movements in the tax schedule. Moving beyond this first-order approach is a difficult but important direction for future work. In the meantime, testing  $S > 0$  provides first-order guidance on how to correctly implement Hicks’ original experiment in a way that accounts for the distortionary cost of redistribution.

If  $S < 0$ , then the alternative environment is Pareto dominated by a modification to the tax schedule. In this sense, alternative environments for which  $S < 0$  are not desirable. But, what about policies for which  $S > 0$ ? Should these be pursued?

Armed with only the result in Proposition 1, it is unclear. While Hicks (1940) originally suggested yes, moving to the alternative environment does not generate a Pareto improvement relative to the status quo. Rather, it generates a Pareto improvement relative to a modified status quo that attempts to replicate the distributional incidence of the alternative environment. Actually moving to the alternative environment would generate winners and losers. Hence,  $S > 0$  suggests it is a useful policy to consider (it’s an “efficient” policy in the sense of Coate (2000)). But, it is not clear whether it is desirable relative to the status quo if  $s(y) < 0$  for some  $y$ .

In order to provide guidance in the case when efficient surplus is positive, it is useful to consider a different conceptual experiment: that of Kaldor (1939).

### 3.2 Finding Pareto Improvements in the Spirit of Kaldor (1939)

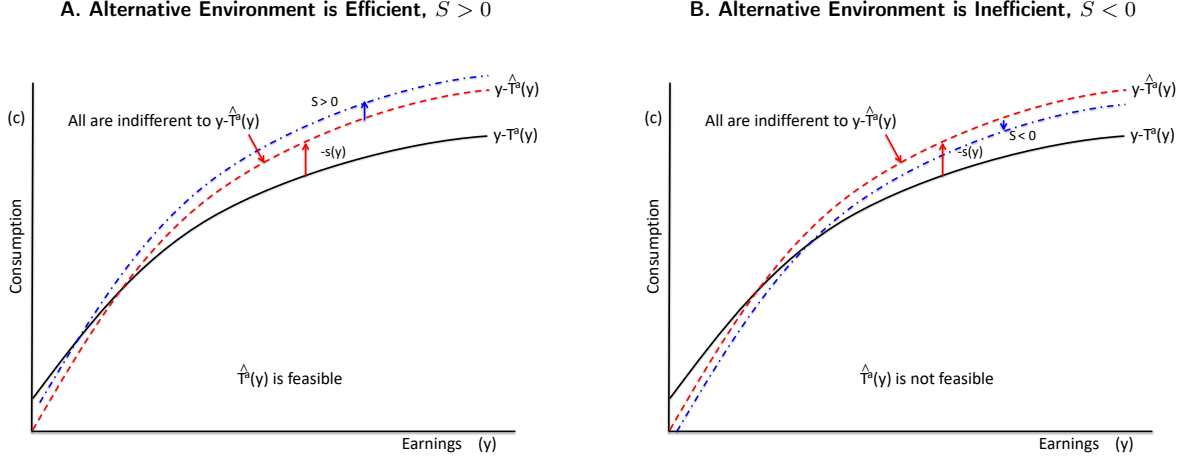
When can everyone be made better off relative to the status quo environment? Consider modifying the tax schedule in the alternative environment,  $T^a(y)$ , so that the winners compensate the losers,  $T^a(y) \rightarrow T^a(y) + s(y)$ .

Figure 4 presents this modified tax schedule in the alternative environment. Those with incomes  $y$  are better off by  $s(y)$  relative to the status quo. The dashed red line in Figure 4 taxes back these gains. The envelope theorem suggests that to first order individuals earning  $y$  in the alternative environment are worse off by  $s(y)$  when we tax back these benefits. Everyone is approximately indifferent between the status quo environment and the alternative environment with the modified income tax schedule. Therefore, the question becomes: Is this modification to the tax schedule in the alternative environment budget feasible?

---

<sup>11</sup>Proposition 1 formalizes the first order approach by scaling the surplus function. Alternatively, one could formalize the approach by directly modeling a continuum of alternative environments in the utility function. For example, suppose  $a$  is a continuous number indexing alternative environments (e.g. level of a public goods, trade policy, etc). Let  $a = 0$  corresponds to the status quo and assume one can write individuals’ utility functions,  $u(c, y, a; \theta)$ . In this case, one can define  $s(y)$  to be individuals marginal willingness to pay out of their own income for a marginal change in  $a$ :  $s(y) = \frac{\partial u}{\partial a} / \frac{\partial u}{\partial c}$  evaluated at  $a = 0$ . In this case, a modification to the tax schedule can make everyone better off relative to a world with a slightly higher value of  $a$  if and only if  $E[g(y) s(y)] < 0$ .

Figure 4: Testing for (Kaldor) Efficiency



*Notes:* This figure illustrates the test of efficiency in Kaldor (1939) that modifies the tax schedule in the alternative environment to attempt to find a Pareto improvement in the modified alternative environment relative to the status quo. The dashed red line presents the after-tax schedule that adds the surplus offered by the alternative environment to the tax schedule,  $T(y) + s(y)$ . To first order, everyone is indifferent between the status quo and the modified alternative environment illustrated by the dashed red line in Panels A and B. The dash-dot blue line then illustrates the after tax income curve that results from closing the government budget constraint. Panel A illustrates the case that the alternative environment is efficient, so that after modifying the tax schedule in the alternative environment there is a Pareto improvement relative to the status quo. Panel B illustrates the case where the alternative environment is inefficient, so that after taxing back the benefits of the alternative environment and closing the budget constraint everyone is worse off relative to the status quo.

To first order, the modification to the tax schedule generates revenue  $S^a = E[g^a(y)s(y)]$ , where  $g^a(y)$  is the cost to the government of providing \$1 to those earning near \$ $y$  in the alternative environment. If  $S^a > 0$ , a modified alternative environment in which the winners compensate the losers through modifications to the tax schedule can make everyone better off relative to the status quo.

In practice,  $S^a$  could differ from  $S$  because the marginal cost of a tax cut may differ in the alternative and status quo environment,  $g(y)$ . If this is the case, it could be that the alternative environment dominates all feasible modifications to the status quo tax schedule ( $S > 0$ ) but there does not exist a modified alternative environment that delivers a Pareto improvement relative to the status quo ( $S^a < 0$ ). But, many applications involve sufficiently small changes to the structure of the economy, in which case it seems reasonable to assume that the marginal cost of taxation is similar in the status quo and alternative environments,  $g^a(y) \approx g(y)$ . I state this formally in Assumption 2.

**Assumption 1.** For sufficiently small  $\tilde{\epsilon}$ , the marginal cost of taxation,  $g(y)$ , in the alternative environment is the same as in the status quo. Specifically, there exists  $\tilde{\epsilon}$  such that if  $\epsilon \in (0, \tilde{\epsilon})$ , then (1)  $y^\epsilon(\theta)$  is the same for all types  $\theta$  that had the same income in the status quo world,  $y^\epsilon(\theta) = y^\epsilon(\theta')$  iff  $y(\theta) = y(\theta')$ , and (2)  $g(y(\theta)) = g(y^\epsilon(\theta))$  for all  $\theta$ .

If Assumption 2 holds, then  $S > 0$  provides a first-order test of whether those with  $s(y) > 0$  can compensate those with  $s(y) < 0$  through modifications to the tax schedule in the alternative envi-

ronment. Proposition 2 states this formally using the same scaled surplus function as in Proposition 1.

**Proposition 2.** *Suppose Assumption 1 holds. For  $\epsilon > 0$ , let  $s_\epsilon = \epsilon s(y)$ . If  $S > 0$ , there exists  $\tilde{\epsilon} > 0$  such that for any  $\epsilon < \tilde{\epsilon}$ , there exists an augmentation to the tax schedule in the alternative environment that delivers surplus  $s_\epsilon^t(y)$  that is positive at all points along the income distribution,  $s_\epsilon^t(y) > 0$  for all  $y$ . Conversely, if  $S < 0$ , then no such  $\tilde{\epsilon}$  exists.*

*Proof.* See Appendix B.3. □

In this sense, testing whether  $S > 0$  provides a first-order approximation to searching for potential Pareto improvements as suggested by Kaldor.

Table 1		
	$S > 0$	$S < 0$
<b>Hicks Experiment:</b>		
Possible to replicate $s(y)$ using tax cut in status quo?	No	Yes
<b>Kaldor Experiment:</b>		
Possible to modify alternative environment tax schedule to make everyone better off relative to status quo?	Yes	No

**Summary** Table 1 summarizes the main results. When efficient surplus is negative,  $S < 0$ , the alternative environment is inefficient in the sense that a feasible modification to the tax schedule in the status quo environment can lead to a Pareto superior allocation to the alternative environment. In this sense, alternative environments for which  $S < 0$  can be rejected by the logic of Hicks (1940) and Coate (2000). When efficient surplus is positive,  $S > 0$ , a modified alternative environment in which the winners compensate the losers through modifications to the tax schedule offers a Pareto superior allocation relative to the status quo. In this sense, the alternative can be preferred using the compensation principle in Kaldor (1939).

Of course, whether such Pareto comparisons are realized depends on whether the modifications to the tax schedule are actually implemented. But, in contrast to the traditional Kaldor-Hicks individual-specific lump-sum transfers, these modifications are feasible.

### 3.3 Non-Marginal Comparisons

The formal results above show that weighting surplus by the efficient welfare weights search for potential Pareto improvements for small surplus comparisons. In practice however, many comparisons

of interest are likely not best thought of as “small”. In these instances, one can continue to construct efficient surplus, but whether this corresponds to a potential Pareto comparison is not guaranteed. Broadly, there are two potential pitfalls that can arise.

First, the efficient welfare weights,  $g(y)$ , are not “structural parameters”. As one implements the transfers, it could be that the marginal cost of the first dollar of the transfers does not equal the marginal cost of the last dollar of the transfers. In this case,  $E[s(y)g(y)]$  would not accurately measure the revenue that the government is able to one would prefer to use the weight that measures the average cost of providing  $s(y)$  to each level of income.

Second, if the alternative environment is sufficiently distinct from the status quo, then an individuals’ willingness to pay will depend on whether it is paid out of income in the status quo or alternative environment. The definition of  $s(\theta)$  above is an “equivalent variation” definition of willingness to pay because it imagines this amount being paid out of income in the status quo. Another method for measuring willingness to pay would be to consider a “compensating variation” definition, which would imagine a willingness to pay out of income in the alternative environment. To first order, these two definitions of willingness to pay are always equivalent. But, they generally differ to second order. This can induce the well-known “cycling” problems associated with compensating variation measures of willingness to pay.

However, there is one important case where compensating and equivalent variation are always equivalent. This is when the comparisons solely involves the willingness to pay for a difference in incomes. An individual is always willing to pay \$10 to receive \$10 of additional income – this is true whether one conceptualizes willingness to pay as an amount of income needed to give someone in the status quo world to make them indifferent to receiving \$10 (equivalent variation), or as the amount of income one can take away in the alternative environment to make them indifferent to not receiving the additional income (compensating variation). For example, the canonical application of the Kaldor-Hicks efficiency test is to use per-capita GDP to compare income distributions across countries or within a country over time. These conceptual comparisons imagine giving each individual in the economy a different level of income, and thus compensating and equivalent variation definitions of willingness to pay are identical. For these reasons, Section 8 applies the weights to revisit these classic experiments in the applications below. But, future work implementing these tests for efficiency could conduct robustness analyses to using both compensating and equivalent variation definitions of willingness to pay.

### 3.4 Additional Limitations

In addition to the issue of non-marginal comparisons, the approach above has several other potential limitations that are worth noting.

**General equilibrium effects** Second, the approach assumes that tax changes have no general equilibrium or spillover effects. Targeting a \$1 tax cut to those earning near  $y$  is assumed to have a willingness to pay of \$1 for the beneficiaries of the tax cut. But, if their wages change in response to

the tax cut, their willingness to pay may differ from \$1. Indeed, with spillovers and general equilibrium effects, the benefits of the tax cut may extend beyond those who are the direct target of the tax cut. But while taxation is not allowed to have GE effects, the approach does allow GE effects to drive the valuation of the alternative environment,  $s(y)$ . For example, the alternative environment could be a policy that makes more land available for agriculture, which in turn lowers food prices. One can still generate individuals' willingness to pay for this alternative environment,  $s(y)$ , and use the efficient welfare weights to ask whether this policy is efficient. In this sense, the efficient welfare weights,  $g(y)$ , are valid even if the policy change or alternative environment has GE effects; but it has ruled out the case where changes in the tax schedule,  $T(y)$ , has GE effects. I leave the incorporation of such effects for future work. Indeed, recent work by Tsyvinski and Werquin (2018) provide one path forward using a structural model of taxation with GE effects.

**Heterogeneity in  $s(\theta)$  conditional on  $y$ .** Third, alternative environments may generate willingness to pay that is heterogeneous conditional on income. In this case, Pareto comparisons are more difficult. To test for Hicks efficiency, one needs to construct the maximum willingness to pay at each income level,  $\bar{s}(y)$ , and test whether  $E[\bar{s}(y)g(y)] > 0$ . If it is negative, then it would be feasible for the government to replicate the surplus offered by the alternative environment and make everyone better off. Intuitively, the government can feasibly provide a tax cut that covers even the maximal willingness to pay at each income level,  $\bar{s}(y)$ . In this sense, the alternative environment would be inefficient. Conversely, to test for Kaldor efficiency, one needs to construct the minimum willingness to pay at each income level,  $\underline{s}(y)$ , and test whether  $E[\underline{s}(y)g(y)] > 0$ . If it is positive, then it would be feasible for the government to redistribute income in the alternative environment so that everyone prefers the modified alternative environment relative to the status quo. Appendix D provides formal statements and proofs of these claims. Often, one might find that  $E[\underline{s}(y)g(y)] < 0$  and  $E[\bar{s}(y)g(y)] > 0$ . In this instance, the alternative environment cannot not be Pareto-ranked relative to the status quo. Nonetheless, the efficient welfare weights,  $g(y)$ , continue to be the key component required to measure  $E[\underline{s}(y)g(y)]$  and  $E[\bar{s}(y)g(y)]$  that facilitates the search for these Pareto comparisons.

**The weights,  $g(y)$  are not structural** Lastly, as noted above, the weights  $g(y)$  are not structural parameters. They are endogenous to the economic environment. In addition to weights changing as one implements transfers, there is also no reason to expect efficient welfare weights identified in one setting or country to readily translate to another setting. This also has potentially interesting implications. For example, some have suggested that top tax rates in France are close to the top of the Laffer curve (Bourguignon and Spadaro (2012)), which implies that reductions in tax rates nearly pay for themselves,  $FE(y) \approx -1$ . In contrast, the point estimates below for the US will suggest smaller fiscal externalities. This would suggest testing for efficiency involves placing less weight on surplus accruing to the rich in France than in the US.

## 4 Relation to Inverse Optimum Program

There is a large and recently-growing literature estimating the solution to the inverse optimum program in optimal taxation. This literature solves for the implicit social preferences that rationalize indifference to the status quo tax schedule.<sup>12</sup> It is straightforward to see that  $g(y)$  is equivalent to the implicit social welfare weights that rationalize indifference to modifications to the tax schedule.

Appendix C provides a formal derivation. To see the logic, let  $\chi(y)$  denote the social marginal utilities of income for those earning near  $y$ , so that an additional \$1 to an individual earning  $y$  has an impact of  $\chi(y)$  on social welfare. Suppose one provides a small tax cut of \$1 to those earning near  $y$ . Those with incomes near  $y$  will be willing to pay \$1 for this tax cut, and it will generate a social welfare impact of  $1 * \chi(y)$ . But, it will have a cost of  $1 + FE(y) = g(y)$ . Hence, the marginal value of additional government spending on a tax cut to those earning near  $y$  will be given by  $\frac{\chi(y)}{g(y)}$ . If the tax schedule is set to maximize social welfare, then the government must be indifferent between raising \$1 from those earning  $y'$  to finance a tax cut to those earning  $y$ . In other words,  $\frac{\chi(y)}{g(y)}$  must be constant for all  $y$ ;

$$\frac{\chi(y)}{g(y)} = \kappa \quad \forall y$$

So,  $\chi(y) = \kappa g(y)$ . Since social welfare weights are only defined up to a constant,  $g(y)$  is the unique set of social welfare weights that rationalize the tax schedule as optimal. Relative to the literature on the inverse optimum program, the core contribution of the present paper is not to re-derive these weights, but rather to show that the Kaldor-Hicks experiments provide a normative foundation for using these weights – even by those whose own social preferences do not rationalize the status quo tax schedule as optimal.

**Indifference to Transfers in Kaldor-Hicks Experiment** The efficient welfare weights,  $g(y)$  measures the cost of implementing the transfers envisioned in the Kaldor-Hicks experiments. In general, one will not be indifferent to whether or not these transfers are implemented. Invoking the Pareto principle requires implementing the transfers. If they are not implemented, then one is back in the world where one must specify a social welfare function to resolve interpersonal comparisons.

However, there is one case in which one is indifferent to whether or not the transfers are implemented: this occurs if and only if one's own social preferences equal the efficient welfare weights (and thus those that rationalize the tax schedule as optimal). To see this, consider an individual with social preferences  $\eta(y)$ . The social welfare impact of taxing back the benefits  $s(y')$  from income level  $y'$  is given by  $-s(y')\eta(y')$ . The government receives revenue of  $s(y')g(y')$ , which can provide  $\frac{s(y')g(y')}{g(y'')}$  dollars of welfare to those earning  $y''$ . The social welfare impact of providing those benefits to those earning  $y''$  is  $\eta(y'') * \frac{s(y')g(y')}{g(y'')}$ . Hence, the net social welfare impact of the transfer,  $\Delta$ , to someone

---

<sup>12</sup>See, e.g., Christiansen (1977); Christiansen and Jansen (1978); Blundell et al. (2009); Bargain et al. (2011); Bourguignon and Spadaro (2012); Lockwood and Weinzierl (2016); Zoutman et al. (2013); Bargain et al. (2014); Jacobs et al. (2017)

with social preferences  $\eta(y)$  is given by

$$\begin{aligned}\Delta &= \chi(y'') s(y) \frac{g(y')}{g(y'')} - s(y') \chi(y') \\ &= s(y') \left( \chi(y'') \frac{g(y')}{g(y'')} - \chi(y') \right)\end{aligned}$$

which equals zero when  $g(y) = \chi(y)$  for all  $y$ . However, when one's social preferences differ from those that rationalize the tax schedule as optimal (i.e.  $g \neq \chi$ ), one will not generally be indifferent to whether the transfers are undertaken. Nonetheless, even though one's own social preferences differ, the feasible transfers envisioned by Kaldor and Hicks continues to motivate their use: implementing these transfers translates a comparison about which many people may disagree depending on their social preferences into a comparison over which universal agreement can be possible.

## 5 Representing Fiscal Externalities using Estimable Parameters

As illustrated in Figure 1, the marginal cost of providing a \$1 tax cut to those with earnings near  $y$  is given by  $g(y) = 1 + FE(y)$ , where  $FE(y)$  is the impact of the behavioral response to the tax cut on government tax revenue. To estimate  $FE(y)$ , I build upon recent work by Jacobs et al. (2017) who provide an expression for  $FE(y)$  as the sum of three components: a participation response, income effect, and substitution effect. Here, I extend the results in Jacobs et al. (2017) to allow for multi-dimensional heterogeneity.

The core assumption required for the representation of  $FE(y)$  is that intensive margin responses to taxation are continuous in the tax rate (note this does not restrict responses at the extensive margin). Assumption 3 states this more precisely in the context of the general model developed in Section 2.

**Assumption 2.** *Fix a type  $\theta$ . For any  $\kappa > 0$ , let  $B(\kappa) = [u(y(\theta) - T(y(\theta)), y(\theta); \theta) - \kappa, u(y(\theta) - T(y(\theta)), y(\theta); \theta) + \kappa]$  denote an interval of width  $\kappa$  near the status quo utility level. For any level of earnings  $y$  and utility level  $w$ , let  $c(y; w, \theta)$  trace out a type  $\theta$ 's indifference curve that is defined implicitly by:*

$$u(c(y; w, \theta), y; \theta) = w$$

*I assume each indifference curve,  $c(y; w, \theta)$ , satisfies the following conditions:*

1. *(Continuously differentiable in utility) For each  $y \geq 0$ , there exists  $\kappa > 0$  such that  $c(y; w, \theta)$  is continuously differentiable in  $w$  for all  $w \in B(\kappa)$*
2. *(Convex in  $y$  for positive earnings, but arbitrary participation decision) For each  $y > 0$ , there exists  $\kappa > 0$  such that  $c(y; w, \theta)$  is twice continuously differentiable in  $y$  for all  $w \in B(\kappa)$  and  $c_y > 0$  and  $c_{yy} > 0$ .*

Part (1) imposes the standard assumption that indifference curves vary smoothly with utility changes. Part (2) requires that indifference curves are convex on the region  $y > 0$  (but not at  $y = 0$ ). Importantly, it allows extensive margin responses: small changes in the tax schedule to cause jumps between

$y = 0$  and some positive income level. It is important to emphasize that Assumption 2 imposes very weak assumptions on utility functions and also allows for arbitrary distributions of unobserved heterogeneity,  $\theta$ .<sup>13</sup>

When Assumption 2 holds, then three behavioral elasticities determine the response to taxation: a compensated elasticity, income elasticity, and participation elasticity. To define these, let  $\tau(y) = T'(y)$  denote the marginal tax rate faced by an individual earning  $y$ . The average intensive margin compensated elasticity of earnings with respect to the marginal keep rate,  $1 - \tau(y)$ , for those earning  $y(\theta) = y$  is given by the percent change in earnings from a percent change in the price of consumption,

$$\epsilon^c(y) = E \left[ \frac{1 - \tau(y(\theta))}{y(\theta)} \frac{dy}{d(1 - \tau)} \Big|_{u=u(c,y;\theta)} \Big| y(\theta) = y \right].$$

The average income elasticity of earnings,  $\zeta(y)$ , is given by the percentage response in earnings to a percent increase consumption,

$$\zeta(y) = E \left[ \frac{dy(\theta) y(\theta) - T(y(\theta))}{dm} \Big|_{y(\theta) = y} \right]$$

The extensive margin (participation) elasticity with respect to net of tax earnings,  $\epsilon^P(y)$ , is given by

$$\epsilon^P(y) = \frac{d[f(y)]}{d[y - T(y)]} \frac{y - T(y)}{f(y)}$$

where  $f(y)$  is the density of income at  $y$ .

Proposition 3 shows how these three elasticities along with the joint distribution of tax rates and income characterize the fiscal externality,  $FE(y)$ .

**Proposition 3.** *For any point  $y$  such that  $\tau(y)$  and  $\epsilon^c(y)$  are constant in  $y$  and the distribution of  $y$  is continuous with density  $f(y)$ , the fiscal externality of providing additional resources to individuals near  $y$  is given by*

$$FE(y) = \underbrace{-\epsilon_c^P(y) \frac{T(y) - T(0)}{y - T(y)}}_{\text{Participation Effect}} - \underbrace{\zeta(y) \frac{\tau(y)}{1 - \frac{T(y)}{y}}}_{\text{Income Effect}} - \underbrace{\epsilon^c(y) \frac{\tau(y)}{1 - \tau(y)} \alpha(y)}_{\text{Substitution Effect}} \quad (3)$$

where  $\alpha(y) = -\left(1 + \frac{yf'(y)}{f(y)}\right)$  is the local Pareto parameter of the income distribution.

*Proof.* Proof provided in Appendix B The appendix also provides a generalized formula for points  $y$  such that  $\epsilon^c(y)$  is not constant in  $y$ .<sup>14</sup>  $\square$

<sup>13</sup>See Kleven and Kreiner (2006) for a particular utility specification that satisfies Assumption 2 and captures these features of intensive and extensive margin labor supply responses.

<sup>14</sup>As noted above, Proposition 3 is a generalization of the formula in Jacobs et al. (2017) to the case of multi-dimensional heterogeneity. Consistent with the intuition provided by Saez (2001), Proposition 3 shows that the relevant empirical elasticities in the case of potentially multi-dimensional heterogeneity are the population average elasticities conditional on income.

The fiscal externality associated with providing an additional dollar resources to an individual earning  $y$  is the sum of three effects. First, people may enter the labor force.  $\epsilon^P(y)$  measures the size of this effect. It's impact on tax revenue depends on the difference between the average taxes received at  $y$ ,  $T(y)$ , and the taxes/transfers received from those out of the labor force,  $T(0)$ .

Second, the increased transfer may change the labor supply of those earning  $y$  due to an income effect. The size of this effect is measured by  $\zeta(y)$ . The impact of this change in earnings on the government budget depends on the marginal tax rate,  $\tau(y)$ .

Finally, people earning close to  $y$  may change their earnings towards  $y$  in order to get the transfer. The elasticity,  $\epsilon^c(y)$ , measures how much people move their earnings towards  $y$  in response to the tax cut. The tax ratio,  $\frac{\tau(y)}{1-\tau(y)}$ , captures the fiscal impact of these responses. However, the net impact on government revenue is the sum of two effects. Some people will decrease their earnings towards  $y$ ; others will increase their earnings towards  $y$ , as depicted by the blue arrows in Figure 1. When  $\tau(y) > 0$ , the former effect increases tax revenue and the latter effect decreases tax revenue. The extent to which the losses outweigh the gains depends on the elasticity of the income distribution,  $\frac{yf'(y)}{f(y)}$ . When  $\frac{yf'(y)}{f(y)} < -1$  (as is the case with the Pareto upper tails in the US income distribution), more people increase rather than decrease their taxable earnings. This means  $\alpha(y) > 0$ . Conversely, if  $\frac{yf'(y)}{f(y)} > -1$  (e.g. if  $f$  is a uniform distribution so that  $f'(y) = 0$ ), then more people decrease than increase their earnings so that  $\alpha(y) < 0$ . This increases the marginal cost of the tax cut. Importantly, this shows that even if elasticities and tax rates are constant, the shape of the income distribution plays a key role in determining the marginal cost of taxation and the shape of efficient welfare weights.

## 6 Bounds on Efficient Welfare Weights in the U.S.

At first glance, equation (3) suggests one requires precise estimates of the size of behavioral responses to taxation in order to quantify the efficient welfare weights. This is potentially problematic because of the general lack of consensus on the size of behavioral responses to taxation (Saez et al. (2012)). Fortunately, under fairly plausible assumptions outlined below, the shape of the income distribution provides insights into the shape of the efficient welfare weights.

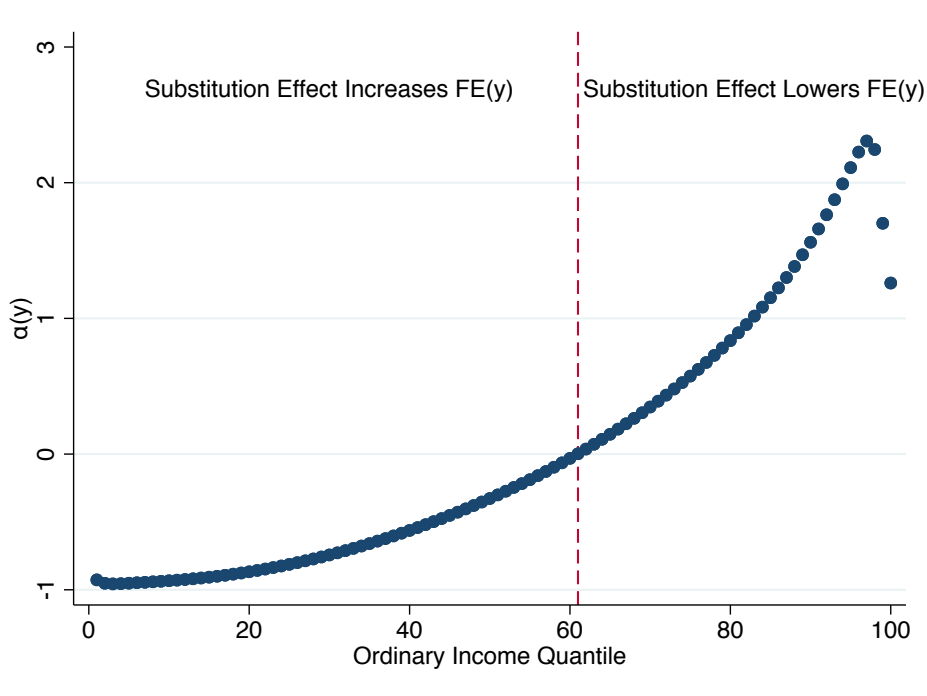
I use the universe of income tax returns from 2012 to estimate the value of  $\alpha(y)$  at each quantile of the income distribution.<sup>15</sup> Figure 5 presents the mean value of  $\alpha(y)$  at each quantile of the ordinary income distribution. The average  $\alpha(y)$  reaches around 1.5 at the top of the income distribution, consistent with findings in previous literature focusing on top incomes (Diamond and Saez (2011) and Piketty and Saez (2013)). However, the key point on Figure 5 is that  $\alpha(y)$  exhibits considerable heterogeneity across the income distribution. It is negative below the 60th percentile of the income distribution,  $\frac{yf'(y)}{f(y)} > -1$ . This implies that the substitution effect increases the marginal cost of a tax cut (assuming a positive elasticity).<sup>16</sup> Conversely, it crosses zero around the 60th percentile, and is

<sup>15</sup>Formally, I construct this by separately estimating  $\alpha(y)$  for each tax schedule using the information in the tax returns on filing status and other determinants of the tax schedule. As noted in the Appendix, throughout I estimate  $g(y)$  using a method that correctly accounts for the heterogeneity in tax schedules faced by those at the same level of income. The details of this procedure are provided in Appendix E

<sup>16</sup>This is consistent with the findings of Werning (2007) who estimates the marginal cost of taxation using the SOI

then positive. This means that  $\frac{yf'(y)}{f(y)} < -1$  for values of  $y$  above the 60th quantile. For those earning more than about \$43K in ordinary income, the substitution effect reduces the cost of providing a tax cut. As long as  $\tau(y) > 0$  and  $\epsilon^c(y) > 0$ , the substitution effect,  $-\epsilon^c(y) \frac{\tau(y)}{1-\tau(y)} \alpha(y)$ , in equation (3) is positive for incomes below \$43K (60th quantile of 2012 ordinary income) and negative for incomes above \$43K.

**Figure 5: Shape of the Income Distribution,  $\alpha(y)$**



*Notes:* This figure presents the estimates of the average  $\alpha(y)$  for each quantile of the income distribution. This function is given by  $\alpha(y) = -\left(1 + \frac{yf'(y)}{f(y)}\right)$ , where  $f(y)$  is the density of the income distribution. For values of  $y$  below the 60th quantile,  $\alpha(y) < 0$  so that the substitution effect in equation (3) raises the marginal cost of taxation. In contrast, for values of  $y$  above the 61st quantile,  $\alpha(y) > 0$  so that the substitution effect lowers the marginal cost of taxation.

In addition to the substitution effect, it is also possible to put bounds on the natural shape of the impact of the participation effect on the government budget. For those with low incomes, the EITC offers transfers for those who enter the labor force; this renders  $T(y) < 0$  so that those who enter the labor force in response to an increased tax cut actually increase the budgetary cost because they obtain the EITC benefits. In contrast, for higher values of  $y$  individuals contribute positive tax revenue so that  $T(y) > 0$ ; thus any increase in labor force participation for those at higher income levels will result in a positive fiscal externality. This suggests the participation effect in equation (3) is also declining in  $y$ .

Lastly, most empirical works suggests income effects effects are either small (Gruber and Saez (2002); Saez et al. (2012)) or declining in income (Cesarini et al. (2015)). As a result, one has a

---

public use file.

natural bound on the shape of efficient welfare weights: Efficient welfare weights put greater weight on those with lower incomes (i.e. below \$43K) than those with higher incomes (i.e. above \$43K). This means that it is costly to redistribute an additional dollar from rich to poor, but cheap to redistribute from poor to rich.

## 7 Using Elasticities to Quantify $g(y)$ in the U.S.

Point estimates of  $g(y)$  require estimates of the behavioral responses to taxation. For those subject to the EITC, I draw upon Chetty et al. (2013) who calculate elasticities of 0.31 in the phase-in region (income below \$9,560) and 0.14 in the phase-out region (income between \$22,870 and \$43,210). Using the income tax return data, I assign these elasticities to EITC filers in these regions of the income distribution. Second, for filers subject to the top marginal income tax rate, I assign a compensated elasticity of 0.3. This is consistent with the midpoint of estimates estimated from previous literature studying the behavioral response to changes in the top marginal income tax rate (Saez et al. (2012)). For those not on EITC and not subject to the top marginal income tax rate, I assign a compensated elasticity of 0.3, consistent with Chetty (2012) who shows such an estimate can rationalize the large literature on the response to taxation. I assess the robustness to alternative elasticities such as 0.1 and 0.5.

In addition to these intensive margin responses, there is also significant evidence of extensive margin behavioral responses, especially for those subject to the EITC. This literature suggests EITC expansions are roughly 9% more costly to the government due to extensive margin behavioral responses.<sup>17</sup> Therefore, I assume the participation effect in equation (3) is equal to 0.09 for income groups subject to the EITC. Above the EITC range, there is mixed evidence of participation responses to taxation. Liebman and Saez (2006) find no statistically significant impact of tax changes on women’s labor supply of women married to higher-income men. Indeed, higher tax rates can reduce participation from a price effect but increase participation due to an income effect. As a result, I assume a zero participation elasticity for those not subject to the EITC.

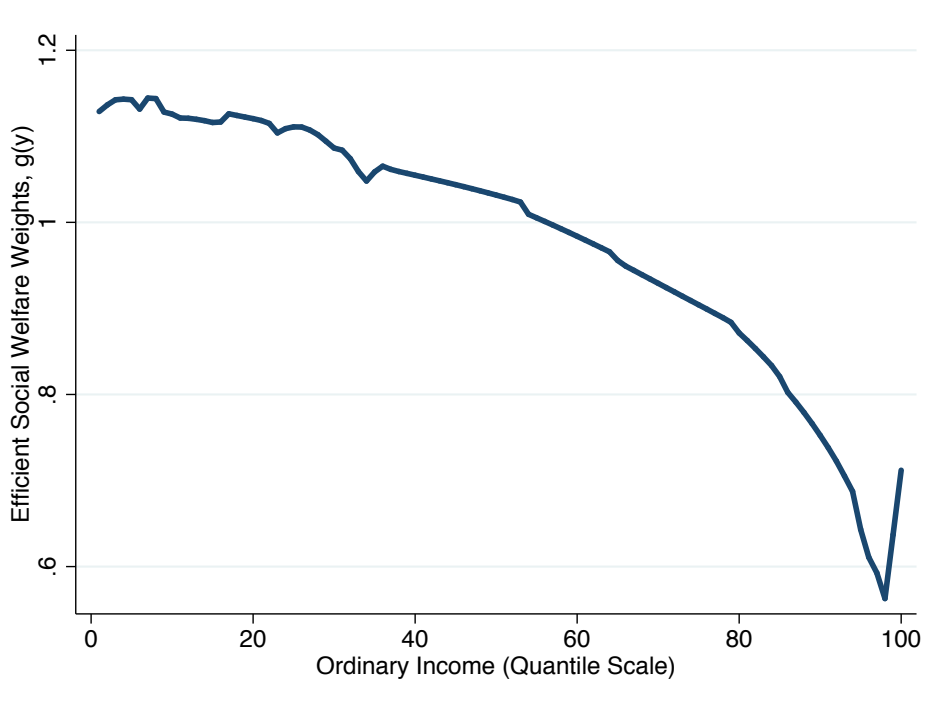
Lastly, I assume away intensive margin income effects, consistent with a large literature suggesting such effects are small (Gruber and Saez (2002); Saez et al. (2012)). Cesarini et al. (2015) find evidence of income effects using Swedish lotteries; however a large portion of these effects are driven by extensive margin responses and arguably already captured by the EITC responses measured above.<sup>18</sup>

**Results** I use equation (3) to combine the estimates of the shape of the income distribution, marginal tax rates, and elasticity calibrations, which generates an estimate of  $FE(y)$  for each filer. I then bin the income distribution into 100 quantile bins and construct the mean fiscal externality,

<sup>17</sup>See Hotz and Scholz (2003) for a summary of elasticities and Hendren (2016) for the 9% calculation.

<sup>18</sup>Nonetheless, Appendix F reports the robustness of the results to an alternative specification that incorporates income effects assuming that the estimates from Cesarini et al. (2015) are entirely along the intensive margin and correspond to an elasticity of  $\zeta = 0.15$ . As shown in Appendix Figure 3, income effects tend to increase the marginal cost of taxation at all income levels; but in contrast to the compensated elasticity they do not affect the relative difference in the weights to low versus high income individuals.

**Figure 6: Efficient Welfare Weights,  $g(y)$**



*Notes:* This figure presents the baseline estimates of the efficient welfare weights,  $g(y)$ , estimated using equation (3) for each quantile of the income distribution.

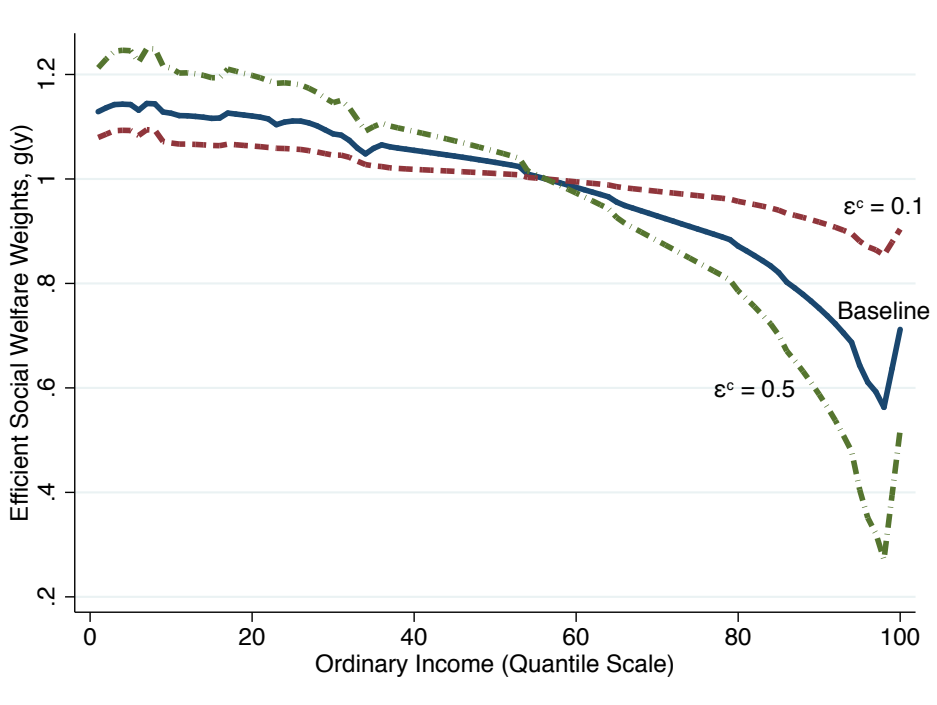
$FE(y)$ , for each quantile of income. The efficient social welfare weight at each income quantile is then given by  $g(y) = 1 + FE(y)$ . Figure 6 presents the resulting estimates for  $g(y)$ . Figure 7 presents the results for the alternative calibrations of the compensated elasticity of  $\epsilon^c = 0.1$  and  $\epsilon^c = 0.5$ .

The weights have several key features. First, consistent with the bounds shown in the previous section, the results suggest it is efficient to place higher weight on surplus to the poor than to the rich. Under the baseline specification, these weights fall from around 1.15 for those at the bottom of the income distribution to 0.65 for those at the top. Transferring \$1 from the top of the distribution can generate around  $0.65/1.15 = \$0.57$  of welfare to someone at the bottom of the distribution. Conversely, transferring \$1 from the bottom of the income distribution can generate around  $1.15/0.65 = \$1.77$  of welfare to the those at the top of the income distribution.

Second, although the weights place more weight on low versus high income individuals, the weights never differ by more than a factor of 2. In other words,  $\left| \frac{g(y)}{g(y')} \right| < 2$  for all  $y$  and  $y'$ . This means that it is not efficient to discount surplus more than 50%, regardless of where it falls in the income distribution. For example, the consumer surplus standard in merger analysis (which gives no weight to producer surplus) would still not be efficient even after accounting for the distortionary cost of taxation.

Third, while the weights generally decline in income, there is an increase in the top 1%. For the baseline specification, it is cheaper to provide additional transfers to the upper middle class than to the top 1%. However, Figure 7 illustrates that this non-monotonicity is not robust to plausible assumptions about how elasticities might change across the income distribution. In particular, if the

Figure 7: Robustness to Alternative Elasticities



Notes: This figure presents the baseline specification for the efficient welfare weights alongside with estimates under alternative constant compensated elasticity scenarios of  $\epsilon^c(y) = 0.1$  and  $\epsilon^c(y) = 0.5$ .

elasticity moves from 0.3 to 0.5 as one goes from the top 2% to the top 1%, the weights would again be monotonically declining in income.

Fourth, all the weights are positive,  $g(y) > 0$  for all  $y$  for the baseline and alternative specifications. This means that it is always costly to provide a tax cut. This implements a Pareto efficiency test suggested by Werning (2007), and suggests there are no Pareto improvements solely from modifying the tax schedule.

Lastly, as foreshadowed by the bounding exercise in the previous Section, there is a similarity between the estimates of  $\alpha(y)$  in Figure 5 and the shape of the efficient welfare weights,  $g(y)$ . Higher elasticities,  $\epsilon^c(y)$ , increase the difference between the weights on the low- versus high-income individuals. But, they do not affect the general conclusion that  $g(y) > 1$  for those with low incomes and  $g(y) < 1$  for those with high incomes.

## 8 Applications: Comparison of Income Distributions

*[Using transfers], “it is always possible for the Government to ensure that the previous income-distribution should be maintained intact” (Kaldor (1939)).*

Kaldor and Hicks’ original motivation was the comparison of different distributions of endowments. Motivated by this classic comparison, I use the efficient welfare weights to compare distributions of

income. I begin with an analysis of changes in the U.S. income distribution over time; I then explore cross-country differences in income distributions.

To compare income distributions, one needs to define a conceptual experiment that clarifies where an individual in one distribution would fall in the alternative distribution. This experiment then defines the surplus function,  $s(y)$ , that can be used to compare the distributions.

In general, one could consider any number of potential experiments. Perhaps people who are at the top of the distribution stay at the top of the distribution in the alternative environment; conversely one could imagine an experiment where people at the top switch with those at the bottom. More generally, to each individual at quantile  $\alpha$  in the status quo world, one can be assigned to a quantile  $r(\alpha)$  in the alternative environment, where  $r(\alpha)$  is a permutation function on  $[0, 1]$ . This generates the surplus function:

$$s^r(\alpha) = Q_a(r(\alpha)) - Q_0(\alpha)$$

For simplicity, I will define the surplus experiment as one that maintains quantile stability,  $r(\alpha) = \alpha$ , so that each person's relative position in the income distribution is maintained intact. This minimizes the size of each individual surplus, which helps make the first-order approximation for the marginal cost of taxation more appropriate. Moreover, choosing  $r(\alpha) = \alpha$  minimizes the estimated surplus of the status quo relative to the alternative environment. Intuitively, having  $r(\alpha) \neq \alpha$  adds an additional redistributive component to the distributional comparison that has value because of the desire for redistribution; but is arguably not relevant for making distributional comparisons.

## 8.1 Income Growth in the U.S.

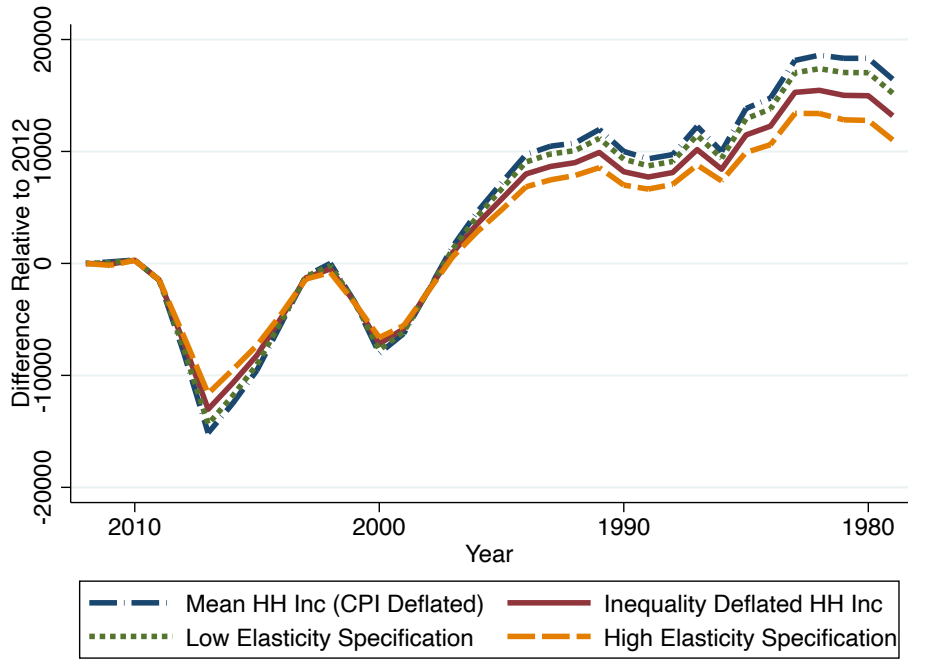
It is well-known that income inequality in the U.S. has increased in recent decades, especially at the top of the distribution (Piketty and Saez (2003)). Appendix Figure 1 plots several quantiles of the household after-tax income distribution over time using data from the Congressional Budget Office (CBO) from 1979-2009.<sup>19</sup> As is well-known, incomes have increased significantly in the top portions of the income distribution, especially the top 20% and top 1%; in contrast, income for the bottom 80% has experienced smaller growth.

Here, I use the efficient welfare weights to calculate how much richer all points of the income distribution would be relative to a given previous year if the tax schedule were augmented in order to hold changes in income inequality constant over time. Let  $Q_0(\alpha)$  denote the  $\alpha$ -quantile of the 2012 income distribution; let  $Q_t(\alpha)$  denote the  $\alpha$ -quantile of an alternative income distribution in year  $t$ . I

---

<sup>19</sup>The data is constructed using Table 7 from CBO publication 43373. I take market income minus federal taxes to construct after-tax income shares across the population. To account for the fact that government spending may have value, I assign net tax collection back to each household in proportion to their after-tax income. This assumes each individuals' willingness to pay for government expenditure is proportional to after-tax income. The CBO also reports an "after-tax" measure of income that includes government transfers. Unfortunately, the bottom portion of the income distribution for these transfers disproportionately falls on the non-working elderly, through social security and Medicare payments. Since these would be affected by modifications to the nonlinear income tax schedule, I do not use this measure of income.

Figure 8: Raw and Deflated Household Income Change Relative to 2012



Notes: This figure presents the un-weighted growth in incomes (relative to 2012) and the distributionally-adjusted growth in incomes using the low elasticity ( $\epsilon = 0.1$ ), baseline elasticity ( $\epsilon = 0.3$ ), and high elasticity ( $\epsilon = 0.5$ ) specifications.

define the efficient surplus in household income by

$$S_t = \int_0^1 [Q_0(\alpha) - Q_t(\alpha)] g^H(Q_0(\alpha)) d\alpha \quad (4)$$

where  $g^H(y)$  are the efficient welfare weights. Intuitively,  $S_t$  is the first-order approximation to the amount by which the U.S. would be richer in 2012 relative to year  $t$  if the 2012 income tax schedule were augmented in to hold constant the changes to the income distribution relative to year  $t$ . All incomes are in units of 2012 income using the CPI-U deflator.

Figure 8 reports the change in mean household income (dashed blue line), along with the efficient surplus under the baseline specification and two alternative elasticity specifications. Mean household income has increased by roughly \$18,300 relative to 1979, but if these benefits were redistributed equally across the population, growth would have increased \$15,000 under the baseline specification (\$13K and \$17K under the high and low elasticity specifications, respectively). From a normative perspective, this lowers the overall growth rate of the U.S. economy by roughly 15-20%: if the U.S. were to make a tax adjustment so that everyone shared equally in the after tax earnings increases, roughly 15-20% of the growth since 1979 would be evaporated.

Figure 9 provides an estimate of the social cost of increased income inequality. To do so, I multiply

the per-household social cost by the total number of households in the U.S.<sup>20</sup> This suggests the social cost of increased income inequality since 1979 is roughly \$400B. From an equivalent variation perspective, undoing the increased inequality would cost roughly \$400B; from a compensating variation perspective, if the U.S. had not experienced the increased inequality, it could have replicated the social surplus provided by the 2012 after tax income distribution even if aggregate economic growth were \$400B less than actually occurred. These numbers depend on the behavioral responses to taxation – if one believes behavioral responses to taxes are larger (e.g. a compensated elasticity of 0.5), then the social cost of increased income inequality is in excess of \$600B.

To be sure, the comparison of the income distribution in 2012 to the income distribution in 1979 is perhaps not best thought of as a “marginal” policy comparison. To that aim, the most robust conclusion that can be drawn from the analysis above is the following: if the distribution of economic growth continued from today to follow the average trend in the US since 1979, then unweighted measures of economic growth will over-state the growth in societal well-being by roughly 15-20%. This 15-20% statistic holds exactly when considering small amounts of economic growth (i.e. short time windows), but as noted in Section 3.3, it could differ when considering larger differences in the income distribution if the marginal cost of taxation changes as one modifies the tax schedule. An important direction for future work is understanding how changes in the tax schedule lead to changes in the efficient welfare weights, which could then be used to adjust for these second-order effects.

## 8.2 Comparisons of Income Distributions: Cross-Country Analysis

It is often noted that the U.S. has a higher degree of income inequality than many other countries of similar income per capita levels. In this subsection, I use the efficient welfare weights to ask how much richer or poorer the U.S. would be relative to these countries if it attempted to replicate their income distributions using modifications to the tax schedule.

The efficient surplus associated with moving from the status quo income distribution to the income distribution in country  $a$  is given by

$$S_a^{ID} = \int_0^1 [Q_a(\alpha) - Q_0(\alpha)] g^H(Q_0(\alpha)) d\alpha \quad (5)$$

I form estimates of  $Q_a(\alpha)$  using data from the World Bank Development Indicators and UN World Income Inequality Database. These sources aggregate household survey data from various countries and to provide measures of the shape of the income distribution.

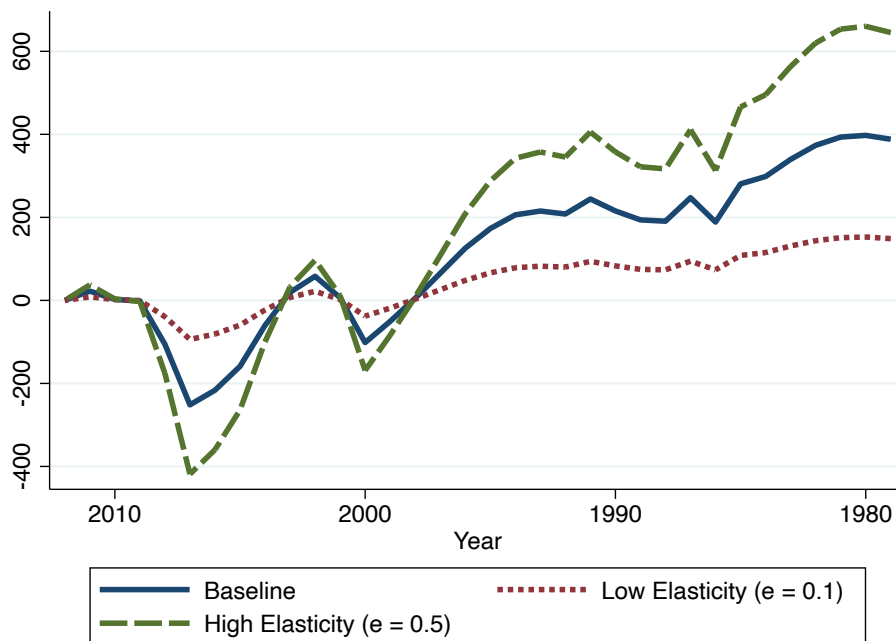
Figure 10 plots deflated surplus against the GNI per capita of each country within \$10,000 of the U.S. GNI per capita. The dots represent the estimates for the baseline specification and the brackets plot the estimates for the low and high elasticity specifications.

The results suggest that a couple of cross-country comparisons based on mean incomes are reversed when using the efficient welfare weights to control for differences in inequality. The U.S. is richer in

---

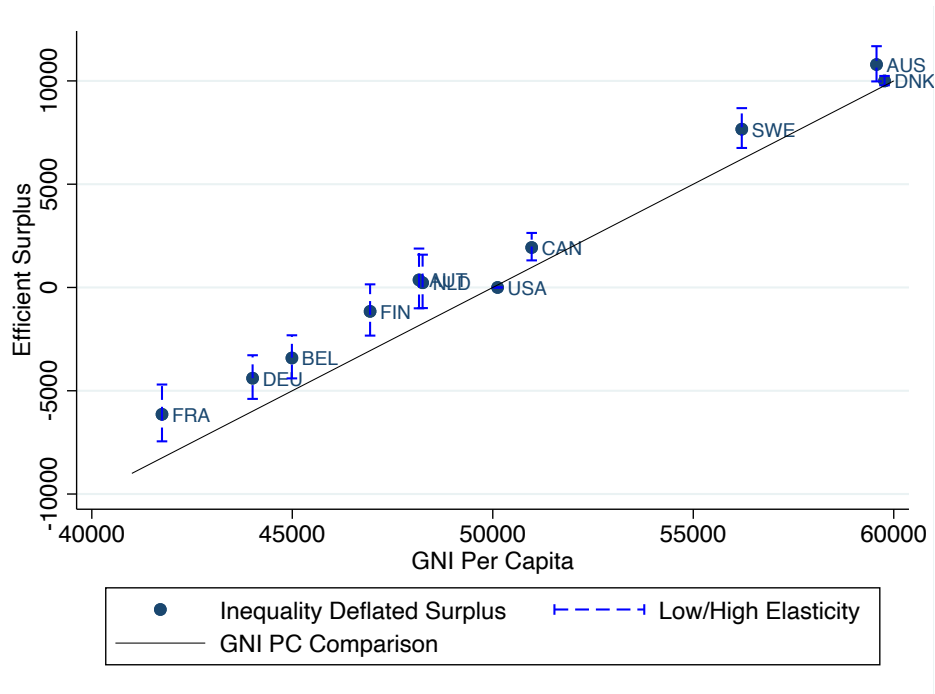
<sup>20</sup>The census reports 117.6M households in 2009, with an annual increase over the years 2006-2009 of roughly 500 households per year, implying roughly 119M households in 2012.

**Figure 9: Social Cost of Increased Income Inequality**



*Notes:* This figure presents a measure of the social cost of income inequality, which is defined as the difference between un-weighted growth in incomes since 2012 and distributionally-adjusted growth in incomes, multiplied by the number of households in the US of 117.6M. The results are presented using the low elasticity ( $\epsilon = 0.1$ ), baseline elasticity ( $\epsilon = 0.3$ ), and high elasticity ( $\epsilon = 0.5$ ) specifications.

Figure 10: Comparisons of Income Distributions Across Countries



Notes: This figure plots efficient surplus and GNI per capita for a selection of countries with gross national incomes (GNI) near that of the US. For each country, the efficient surplus (defined in equation 5) is presented for the baseline elasticity specification against GNI per capita on the horizontal axis; vertical bars representing the high and low elasticity specifications. If all countries had the same degree of inequality, then all countries would align on the 45 degree line. The fact that other countries lie above this 45 degree line reflects the greater degree of income inequality in the U.S. relative to these countries.

mean per capita terms than Austria (AUT) and New Zealand (NLD) by roughly \$2,000. But despite it's higher income level, if the U.S. were to try to provide the distribution of purchasing power offered by these countries, each point of the income distribution would be made worse off relative to these countries under the baseline elasticity specification. Under the high elasticity specification, it would be efficient to take Finland's income distribution over the US's income distribution, even though it has \$3,180 less in per capita national income.

## 9 Welfare Evaluation of Policy Changes

*“All that economics can, and should, do in this field, is to show, given the pattern of income-distribution desired, which is the most convenient way of bringing it about?”*

*(Kaldor (1939, p552))*

While many may disagree with Kaldor about whether this is all that economics should do, the efficient welfare weights provide a path to answering the classic question posed by Kaldor about how one can

most efficiently provide a given distribution of income. Should we increase food stamp spending? Reduce Medicaid spending? Provide free public transportation?

Relative to the comparison of income distributions discussed in the previous section, the key additional complexity in these examples is that the policy changes envisioned are generally not budget-neutral. Willingness to pay may be positive for all the beneficiaries of the policy, but one also needs a method to account for the cost to the government of the policy.

To account for this, this section shows how one can search for potential Pareto improvements in the spirit of Kaldor’s question by constructing a policy’s marginal value of public funds (MVPF). The MVPF of a policy is the willingness to pay of the policy divided by the net cost to the government of the policy.<sup>21</sup> The MVPF measures the “bang for the buck” of the policy.

Given the MVPF of a policy, the Kaldor-Hicks search for efficiency suggests comparing the MVPF of a policy to the MVPF of a tax cut with the same distributional incidence. If the MVPF is higher than the MVPF of a distributionally-equivalent tax cut, then spending money on the policy, financed by increased taxes on those individuals, leads to a Pareto improvement (as long as one retains the assumption that there is no heterogeneity in willingness to pay conditional on income).

To see this, consider a policy that affects those with incomes near  $y^*$ . Let  $s^*$  denote individuals’ willingness to pay out of their own income for the policy change and let  $c$  denote the net cost to the government of the policy. Importantly,  $c$  should incorporate any fiscal externalities from the policy change. For example, if the policy builds roads that increase labor earnings, it should incorporate the resulting increase in tax revenue.

Now, consider the Hicks’ experiment in which the government tries to replicate  $s$  through modifications to the tax schedule. This would cost  $s^*g(y^*)$ . It would be cheaper to replicate this surplus through the tax schedule if and only if

$$s^*g(y^*) \geq c \tag{6}$$

If equation (6) holds, it is more efficient to provide a tax cut to those earning near  $y$  than it is to increase spending on the policy.<sup>22</sup> In this sense, one can prefer the policy using the Pareto principle: one could raise revenue from those individuals themselves to pay for the policy, and still make them better off. Re-writing equation (6) as:

$$MVPF = \frac{s^*}{c} \geq \frac{1}{g(y^*)} \tag{7}$$

---

<sup>21</sup>See Mayshar (1990) for an original definition and more recently Slemrod and Yitzhaki (1996, 2001); Kleven and Kreiner (2006); Eissa et al. (2008); Immervoll et al. (2007, 2011); Hendren (2016); Hendren and Sprung-Keyser (2019).

<sup>22</sup>Equation (6) can be readily extended to the case where there are multiple beneficiaries with willingness to pay  $s(y)$ . In this case, the LHS of equation (6) would be  $E[s(y)g(y)]$ , instead of  $s^*g(y^*)$ . The bias from using the average WTP,  $s^*$ , and the average value of  $g$ ,  $g(y^*)$  instead of  $E[g(y)s(y) | y \in Y]$  comes from two sources: nonlinearities in  $g(y)$  and covariance between  $g(y)$  and  $s(y)$  amongst the beneficiaries,

$$E[g(y)s(y) | y \in Y] = s^*g(y^*) + \underbrace{(E[g(y) | y \in Y] - g(y^*)) s^*}_{\text{Nonlinearity in } g(y)} + \underbrace{\text{cov}(g(y), s(y) | y \in Y)}_{\text{Cov of WTP with } g(y)}$$

These biases are small when the income of the target population is concentrated around a particular  $y^*$  or if the efficient welfare weights are relatively constant within the beneficiary population.

yields an expression in which the LHS of equation (7) is the marginal value of public funds (MVPF) of the policy change, defined as the benefits each policy provides to its beneficiaries,  $s^*$ , normalized by the net cost to the government of the policy  $c$ . The RHS of equation (7) is the MVPF of a tax cut targeted to those with the same incomes as the beneficiaries of the policy. If there is no heterogeneity in willingness to pay conditional on income, then one can search for potential Pareto improvements by comparing the MVPF of the policy in question to the MVPF of a tax cut with the same distributional incidence, which is given by  $1/g(y^*)$ . As a result, the efficient welfare weights allow one to provide precise guidance on the desirability of a policy given (a) its MVPF and (b) the incomes of its beneficiaries.<sup>23</sup>

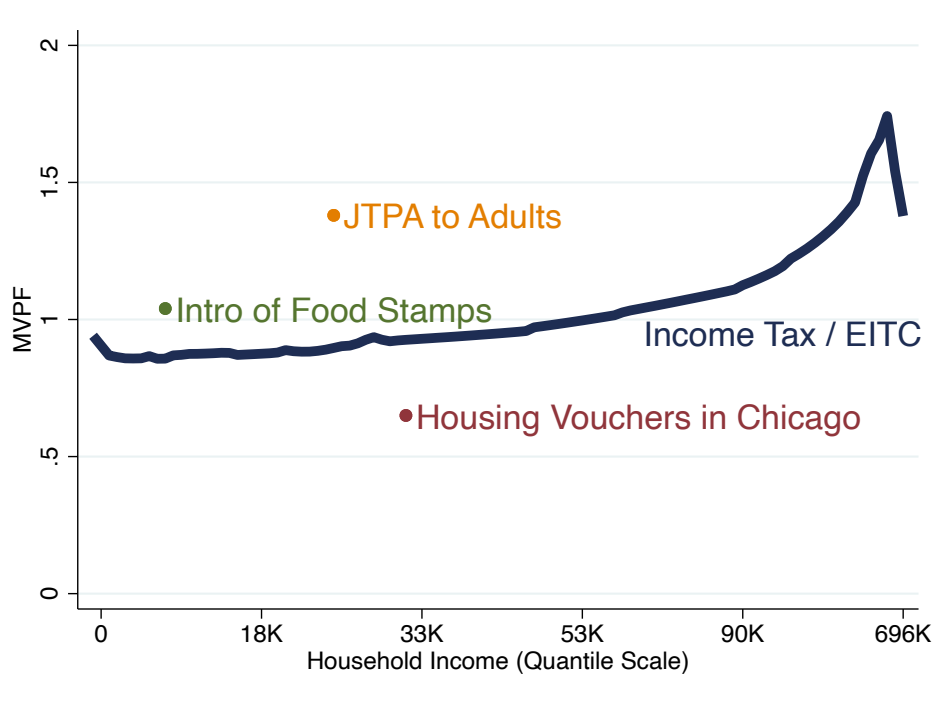
**Application** I illustrate the welfare framework by studying the efficiency of three policies whose MVPFs are computed in Hendren and Sprung-Keyser (2019). These policies include Section 8 housing vouchers, food stamps, and job training. Hendren and Sprung-Keyser (2019) provides details on the construction of the MVPF for each of these policies. The MVPF for Section 8 vouchers draws upon estimates of the fiscal externality inferred from the causal effects in Jacob and Ludwig (2012) on adult labor supply using estimates from Chicago. The MVPF for food stamps draws upon estimates of the labor supply fiscal externalities from Hoynes and Schanzenbach (2012), along with subsequent work documenting spillovers onto children and health outcomes. And, the MVPF for job training draws upon estimates of the impact of the Job Training Partnership Act (JTPA) on labor earnings and benefit substitution from Bloom et al. (1997).

Figure 11 presents the MVPFs of these three policy changes and compares them to the MVPF of a distributionally-equivalent tax,  $1/g(y)$ . The horizontal axis corresponds to the quantile corresponding to the mean income,  $\bar{y}$ , of the policy beneficiaries. The point estimates suggest that housing vouchers are slightly less efficient forms of redistribution than modifications to the income tax schedule. Put differently, the beneficiaries of these policies would prefer the government instead spend the same amount of money on a tax cut (e.g. EITC expansion) instead of housing vouchers. In contrast, the estimates suggest the JTPA may be a more efficient policy than a tax cut. This is because of positive fiscal externalities generated through this program through increased taxable income and reductions in other social programs, which lowers the net cost of the policy to the government,  $c$ . However, as is noted in Hendren and Sprung-Keyser (2019), each of these estimates contains considerable sampling variation, and thus these conclusions should be thought of as illustrative of the methods, not definitive policy conclusions. The key advantage of the framework is that it provides normative conclusions about policies without relying on a social welfare function. In the spirit of Kaldor and Hicks, the efficient welfare weights help replace normative preferences over policies with positive assessments about causal effects and individuals' willingnesses to pay, combined with the Pareto principle.

---

<sup>23</sup>Appendix G discusses how comparing the MVPF of a policy change to the MVPF of a tax cut with the same distributional incidence relates to a test of the weak separability assumption in the Atkinson-Stiglitz and Hylland-Zeckhauser theorems (Atkinson and Stiglitz (1976); Hylland and Zeckhauser (1979)).

Figure 11: Illustrating the Test for Efficiency of Policy Changes



*Notes:* This figure illustrates the use of the efficient welfare weights for assessing the efficiency of government policy changes. The line presents the value of  $\frac{1}{g(y)}$ , which represents the amount of welfare that can be delivered to each portion of the income distribution per dollar of government spending. The dots present estimates of the marginal value of public funds (MVPF) for three policy examples: the job training partnership act (JTPA) from Bloom et al. (1997), food stamps from Hoynes and Schanzenbach (2012), and Section 8 housing vouchers from Jacob and Ludwig (2012). The vertical axis presents the estimated MVPF from Table 1 of Hendren (2016); the horizontal axis presents the estimated income quantiles of the beneficiaries of each policy (normalized to 2012 income using the CPI-U). An MVPF that falls above (below) the the Income/EITC line correspond to policies that can(not) generate Pareto improvements.

## 10 Conclusion

In their original work, Kaldor and Hicks hoped to provide a method to avoid the inherent subjectivity involved in resolving interpersonal comparisons. Weighting surplus using efficient welfare weights measures the economic efficiency of an alternative environment (or policy change) using the Kaldor-Hicks redistributive experiments but accounting for the distortionary cost of taxation. Estimates for the US suggest that redistribution from rich to poor is more costly than from poor to rich. Thus, it is efficient to place greater weight on the poor than on the rich. Regardless of one's own social preferences, surplus to the poor can be turned into greater welfare for everyone than surplus to the rich. The shape of the efficient welfare weights is largely driven by the shape of the income distribution, as opposed to assumptions about behavioral responses to taxation. As a result, the broad conclusion of declining efficient welfare weights is robust to a wide range of assumptions about behavioral elasticities.

There are many important directions for future work, including incorporating the general equilibrium effects of taxation (as in ongoing work by Tsyvinski and Werquin (2018)). Additionally, one could extend the analysis here to construct weights that involve redistribution not just through the tax schedule but also via other means, such as health insurance subsidies or other policies. By expanding the dimensionality of the weights, it could help deal with settings where surplus varies conditional on income. Lastly, implementing the approach requires implementing the transfers that were envisioned by Kaldor and Hicks. Future work could discuss the implications of political economy or other constraints that might prevent such transfers in practice.

In the end, reasonable people and economists will always disagree about the optimal degree of redistribution in society. But, such debates need not lead to paralysis for debates about how best to bring about this degree of redistribution. To that aim, I hope the efficient welfare weights can be a tool to help generate greater consensus about the desirability of economic policies.

## References

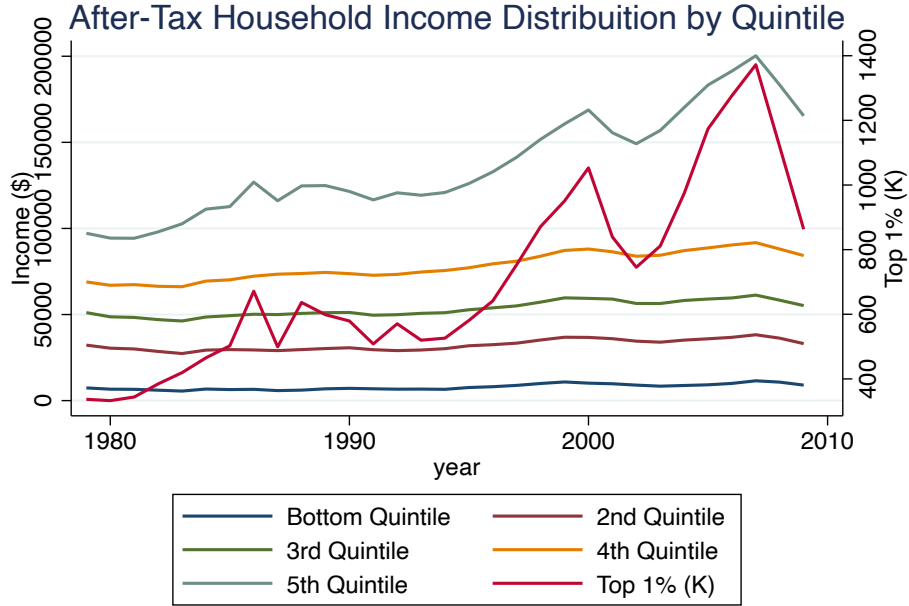
- Atkinson, A. B. and J. E. Stiglitz (1976). The design of the tax structure: Direct versus indirect taxation. *Journal of Public Economics* 6(1-2). 23
- Bargain, O., M. Dolls, D. Neumann, A. Peichl, and S. Siegloch (2011). Tax-benefit systems in europe and the us: between equity and efficiency. *CESifo Working Paper No. 3534*. 1, 4, 12
- Bargain, O., M. Dolls, D. Neumann, A. Peichl, and S. Siegloch (2014). Comparing inequality aversion across countries when labor supply responses differ. *International Tax and Public Finance* 21, 845â873. 1, 4, 12
- Bloom, H., L. L. Orr, S. H. Bell, G. Cave, F. Doolittle, W. Lin, and J. M. Boss (1997). The benefits and costs of jtpa title ii-a programs: Key findings from the national job training partnership act study. *The Journal of Human Resources* 32(3), 549â576. 9

- Blundell, R., M. Brewer, P. Haan, and A. Shephard (2009). Optimal income taxation of lone mothers: An empirical comparison of the uk and germany. *The Economic Journal* 119, 101–121. 1, 4, 12
- Bourguignon, F. and A. Spadaro (2012). Tax-benefit revealed social preferences. *Journal of Economic Inequality* 10, 75–108. 1, 4, 3.4, 12
- Cesarini, D., E. Lindqvist, M. J. Notowidigdo, and R. Östling (2015). The effect of wealth on individual and household labor supply: evidence from swedish lotteries. Technical report, National Bureau of Economic Research. 6, 7, 18, F
- Chetty, R. (2012). Bounds on elasticities with optimization frictions: A synthesis of micro and macro evidence on labor supply. *Econometrica* 80(3), 969–1018. 7
- Chetty, R., J. Friedman, and E. Saez (2013). Using differences in knowledge across neighborhoods to uncover the impacts of eitc on earnings. *American Economic Review (Forthcoming)*. 7
- Chetty, R., N. Hendren, P. Kline, and E. Saez (2014). Where is the land of opportunity: The geography of intergenerational mobility in the united states. *The Quarterly Journal of Economics* 129(4), 1553–1623. E.1
- Christiansen, V. (1977). The theoretical basis for deriving distributive weights to be used in cost-benefit analysis. Technical report, Memorandum from the Institute of Economics, University of Oslo, 25 April. 1, 4, 12
- Christiansen, V. (1981, 07). Evaluation of Public Projects under Optimal Taxation. *The Review of Economic Studies* 48(3), 447–457. 1
- Christiansen, V. and E. S. Jansen (1978). Implicit social preferences in the norwegian system of indirect taxation. *Journal of Public Economics* 10(2), 217 – 245. 1, 4, 12
- Coate, S. (2000). An efficiency approach to the evaluation of policy changes. *Economic Journal* 110(463), 437–455. 1, 3, 3.1, 3.2
- Diamond, P. A. and E. Saez (2011). The case for a progressive tax: From basic research to policy recommendations. *Journal of Economic Perspectives* 25(4), 165–90. 6, E.2
- Eissa, N., H. J. Kleven, and C. T. Kreiner (2008). Evaluation of four tax reforms in the united states: Labor supply and welfare effects for single mothers. *Journal of Public Economics* 92(3-4), 795–816. 21
- Gruber, J. and E. Saez (2002, April). The elasticity of taxable income: evidence and implications. *Journal of Public Economics* 84(1), 1–32. 6, 7
- Hendren, N. (2016). The policy elasticity. *Tax Policy and the Economy* 30. 1, 17, 21, 9
- Hendren, N. and B. D. Sprung-Keyser (2019). A unified welfare analysis of government policies. Technical report, National Bureau of Economic Research. 1, 21, 9

- Hicks, J. R. (1939). The foundations of welfare economics. *Economic Journal* 49, 696–712. 1
- Hicks, J. R. (1940). The valuation of social income. *Economica* 7(26), 105–124. 1, 2.1, 3.1, 3.2
- Hotz, J. and K. Scholz (2003). The earned income tax credit. In R. A. Moffit (Ed.), *Tax Policy and the Economy*. University of Chicago Press. 17
- Hoynes, H. W. and D. W. Schanzenbach (2012). Work incentives and the food stamp program. *Journal of Public Economics* 96(1-2), 151–162. 9
- Hylland, A. and R. Zeckhauser (1979). Distributional objectives should affect taxes but not program choice or design. *The Scandinavian Journal of Economics* 81(2), 264–284. 23, G
- Immervoll, H., H. J. Kleven, C. T. Kreiner, and E. Saez (2007). Welfare reform in european countries: A microsimulation analysis. *The Economic Journal* 117, 1–44. 21
- Immervoll, H., H. J. Kleven, C. T. Kreiner, and N. Verdelin (2011). Optimal tax and transfer programs for couples with extensive labor supply responses. *Journal of Public Economics* 95, 1485–1500. 21
- Jacob, B. A. and J. Ludwig (2012). The effects of housing assistance on labor supply: Evidence from a voucher lottery. *American Economic Review* 102(1), 272–304. 9
- Jacobs, B. A., E. L. W. Jongen, and F. T. Zoutman (2017). Revealed social preferences of dutch political parties. *Journal of Public Economics* 165. 1, 4, 2, 12, 5, 14
- Kaldor, N. (1939). Welfare propositions of economics and interpersonal comparisons of utility. *The Economic Journal* 49(195), 549–552. 1, 2.1, 3.1, 3.2, 8, 9
- Kaplow, L. (1996). The optimal supply of public goods and the distortionary cost of taxation. *National Tax Journal* 49, 513–534. G
- Kaplow, L. (2004). On the (ir)relevance of the distribution and labor supply distortion to government policy. *Journal of Economic Perspectives* 18, 159–175. 1, G
- Kaplow, L. (2008). *The Theory of Taxation and Public Economics*. Princeton University Press. G
- Kleven, H. J. and C. T. Kreiner (2006). The marginal cost of public funds: Hours of work versus labor force participation. *Journal of Public Economics*, 1955–1973. 13, 21
- Liebman, J. and E. Saez (2006). Earnings responses to increases in payroll taxes. *Working Paper*. 7
- Lockwood, B. B. and M. Weinzierl (2016). The evolution of revealed social preferences in the united states and the costs of unequal growth and recessions. *Journal of Monetary Economics* 77, 30–47. 1, 4, 5, 12
- Mayshar, J. (1990). On measures of excess burden and their applications. *Journal of Public Economics* 43(3), 263–89. 21

- Mirrlees, J. A. (1971). An exploration into the theory of optimal income taxation. *The Review of Economic Studies* 38(2), 175–208. 1
- Piketty, T. and E. Saez (2003). Income inequality in the united states, 1913-1998. *Quarterly Journal of Economics* 118(1), 1–41. 8.1
- Piketty, T. and E. Saez (2013). Optimal labor income taxation. *Handbook of Public Economics* 5, 391. 6, E.2
- Saez, E. (2001). Using elasticities to derive optimal income tax rates. *Review of Economic Studies* 68, 205–229. 14, 31, E.2
- Saez, E., J. Slemrod, and S. H. Giertz (2012). The elasticity of taxable income with respect to marginal tax rates: A critical review. *Journal of Economic Literature* 50(1), 3–50. 6, 7, E.1
- Saez, E. and S. Stantcheva (2016). Generalized social marginal welfare weights for optimal tax theory. *American Economic Review* 106(1), 24–45. 3, C
- Schlee, E. E. (2013). Radner’s cost-benefit analysis in the small: An equivalence result. *Working Paper*. 10
- Slemrod, J. and S. Yitzhaki (1996). The social cost of taxation and the marginal cost of funds. *International Monetary Fund Staff Papers* 43(1), 172–98. 21
- Slemrod, J. and S. Yitzhaki (2001). Integrating expenditure and tax decisions: The marginal cost of funds and the marginal benefit of projects. *National Tax Journal*. 21
- Tsyvinski, A. and N. Werquin (2018). Generalized compensation principle. Technical report, National Bureau of Economic Research, Working Paper 23509. 3.4, 10
- Werning, I. (2007). Pareto efficient income taxation. *MIT Working Paper*. 16, 7
- Zoutman, F. T., B. Jacobs, and E. L. W. Jongen (2013). Optimal redistributive taxes and redistributive preferences in the netherlands. *Mimeo*. 1, 4, 12

Appendix Figure 1: Distribution of After-Tax Income in the US (1979-2009)



Source: CBO; Supplemental Tables 43373, Table 7

Notes: This figure the distribution of income in the US by quintile and for the top 1% using data from the Congressional Budget Office. .

## Appendix (Not for Publication)

### A Marginal Cost of Taxation

This section formally defines the marginal cost of taxation,  $g(y)$ . To begin, for any tax schedule,  $\tilde{T}(\circ)$ , let  $R(\tilde{T})$  denote government revenue in the status quo environment with tax schedule  $\tilde{T}$ ,

$$R(\tilde{T}) = E \left[ \tilde{T} \left( \tilde{y}(\theta; \tilde{T}) \right) \right]$$

where  $\tilde{y}(\theta; \tilde{T})$  denotes the earnings choice of a type  $\theta$  when facing tax schedule,  $\tilde{T}$ , in the status quo environment.

I impose the regularity condition that tax revenue is continuously differentiable with respect to changes in the tax schedule. This would be immediately satisfied if individual behavioral responses were continuously differentiable (e.g. imposing standard quasi-convexity of the utility function). But, the regularity assumption allows for people to enter/exit the labor force, change jobs, or conduct other behavior that has discrete impacts on their earnings in response to tax changes. It only assumes these discrete responses average out across the population so that the magnitude of the aggregate behavioral response on tax revenue is smooth. Assumption 3 states this more formally.

**Assumption 3.** For any function  $h(y)$  of taxable income  $y$ , let  $\tilde{T}_\epsilon(y) = T(y) + \epsilon h(y)$ . Then  $R$  is continuously differentiable in  $\epsilon$  for any function  $h(y)$ .

In the alternative environment, let  $R^a(\tilde{T}) = E\left[\tilde{T}\left(\tilde{y}^a(\theta; \tilde{T})\right)\right]$  denote the revenue raised from a tax schedule  $\tilde{T}(\circ)$ , where  $\tilde{y}^a(\theta; \tilde{T})$  is the earnings choice of type  $\theta$  in the alternative environment facing tax schedule  $\tilde{T}(\circ)$ . I assume Assumption 1 holds for  $R^a$ .

Suppose individuals with income  $y^*$  are willing to pay  $s(y) = \$1$  for the alternative environment. How much does it cost to replicate this \$1 benefit through a tax cut in the status quo environment? Figure 1 depicts a small tax deduction to those with earnings near  $y^*$ . To be precise, let  $\eta, \epsilon > 0$  and fix a given income level  $y^*$ . Consider providing an additional  $\eta$  to individuals in an  $\epsilon$ -region near  $y^*$ . Define  $\hat{T}(y; y^*, \epsilon, \eta)$  by

$$\hat{T}(y; y^*, \epsilon, \eta) = \begin{cases} T(y) & \text{if } y \notin (y^* - \frac{\epsilon}{2}, y^* + \frac{\epsilon}{2}) \\ T(y) - \eta & \text{if } y \in (y^* - \frac{\epsilon}{2}, y^* + \frac{\epsilon}{2}) \end{cases}$$

so that  $\hat{T}$  provides  $\eta$  additional resources to an  $\epsilon$ -region of individuals earning between  $y^* - \epsilon/2$  and  $y^* + \epsilon/2$ .<sup>24</sup> By the envelope theorem, to first order individuals with earnings between  $y^* - \epsilon/2$  and  $y^* + \epsilon/2$  will be willing to pay  $\eta$  to have a tax schedule given by  $\hat{T}(y; y^*, \epsilon, \eta)$  instead of  $T(y)$ .<sup>25</sup>

If there were no behavioral responses to changes in taxes, then the cost to the government of providing  $\eta$  to those earning near  $y^*$  would be  $\eta$  per person, so that the mechanical marginal cost in the absence of behavioral responses is 1. But, the presence of behavioral responses induce a fiscal externality on the government. To capture this, consider the derivative of revenue,  $R(\hat{T}(\circ; y^*, \epsilon, \eta))$ , with respect to the size of the tax cut,  $\eta$ , and evaluate at  $\eta = 0$ . This yields the function  $\frac{d[R(\hat{T}(\circ; y^*, \epsilon, \eta))]}{d\eta}\Big|_{\eta=0}$ , which is the marginal cost of providing an additional dollar through the tax code to individuals with earnings in an  $\epsilon$ -region of  $y^*$ . Then, taking the limit as  $\epsilon \rightarrow 0$ , one arrives at the marginal cost to the government of providing an additional dollar of resources to an individual earning  $y^*$ :

$$g(y^*) \equiv \lim_{\epsilon \rightarrow 0} \frac{d\left[R\left(\hat{T}(\circ; y^*, \epsilon, \eta)\right)\right]}{d\eta}\Big|_{\eta=0} \quad (8)$$

Since the choice of  $y^*$  was arbitrary, I use  $y$  to denote the argument of  $g(y)$  instead of  $y^*$  going forward. The cost of the tax cut to those earning near  $y$  is comprised of two components: a mechanical cost of \$1 and a fiscal externality,  $FE(y) = g(y) - 1$

$$g(y) = 1 + FE(y)$$

<sup>24</sup>Note this is a discontinuous modification to the tax schedule. However, nothing in the analysis requires either  $T$  (or the modified schedule,  $\hat{T}$ ) to be continuous or differentiable.

<sup>25</sup>This is true as long as the incidence of the tax cut falls entirely on the beneficiaries and does not result in changes in wages. For example, if firms respond to the tax cut of \$1 by lowering wages by \$0.50, then the individual would only be willing to pay \$0.50 for a \$1 tax cut. Here, I assume no general equilibrium responses, but this could be incorporated into future work. I discuss this further in Section 3.4.

If taxable income did not respond to changes in taxes, the marginal cost would be \$1 per beneficiary,  $g(y) = 1$ . The difference,  $FE(y) = g(y) - 1$ , equals the size of the fiscal externality from the behavioral response to the tax cut. The tax cut could cause people to work less and reduce tax revenue ( $FE(y) > 0$ ); conversely, it could cause others to increase their tax payments,  $FE(y) < 0$ . The size of the fiscal externality is an empirical question: it depends on the causal impact of tax changes on the government budget. In Section 5, I provide additional assumptions that enable one to identify  $FE(y)$  using behavioral elasticities, the shape of the income distribution, and the shape of the tax schedule.

The definition of  $g(y)$  above is the marginal cost of providing an additional \$1 to those earning  $y$  in the status quo environment. For the alternative environment, I define  $g^a(y)$  analogously.<sup>26</sup> Note that it may not be the case that  $g^a(y) = g(y)$  since the alternative environment may have different tax schedules, distributions of income, and responses to taxation.

## B Proofs

### B.1 Preliminaries

For all proofs below, let  $\mu(\theta)$  denote the measure over the type distribution and let  $F(y)$  denote the cumulative distribution of income in the status quo,  $F(x) = \int 1\{y(\theta) \leq x\} d\mu(\theta)$ . The function  $\hat{T}(y; y^*, \epsilon, \eta)$  is given by

$$\hat{T}(y; y^*, \epsilon, \eta) = \begin{cases} T(y) & \text{if } y \notin (y^* - \frac{\epsilon}{2}, y^* + \frac{\epsilon}{2}) \\ T(y) - \eta & \text{if } y \in (y^* - \frac{\epsilon}{2}, y^* + \frac{\epsilon}{2}) \end{cases}$$

so that  $\hat{T}$  provides  $\eta$  additional resources to an  $\epsilon$ -region of individuals earning between  $y^* - \epsilon/2$  and  $y^* + \epsilon/2$ . Fix  $\epsilon$  and let  $\hat{q}(y^*, \epsilon, \eta)$  denote the net government resources expended under  $\hat{T}(y; y^*, \epsilon, \eta)$ . Given the tax schedule,  $\hat{T}(y; y^*, \epsilon, \eta)$ , let  $\hat{y}(\theta; y^*, \epsilon, \eta)$  denote the individual  $\theta$ 's choice of earnings,  $y$ . The net resources expended is given by:

$$\hat{q}(y^*, \epsilon, \eta) = \frac{-1}{F(y^* + \frac{\epsilon}{2}) - F(y^* - \frac{\epsilon}{2})} \int_{\theta} \hat{T}(\hat{y}(\theta; y^*, \epsilon, \eta); y^*, \epsilon, \eta) d\mu(\theta) \quad (9)$$

One needs to evaluate this derivative with respect to  $\eta$  at  $\eta = 0$ . WLOG, I assume  $\hat{q}(y^*, \epsilon, 0) = 0$  so that the status quo tax schedule is budget neutral.

<sup>26</sup>Formally, let  $\tilde{T}(\circ; y^*, \epsilon, \eta)$  denote a modified tax schedule in the alternative environment. Then,

$$g^a(y^*) \equiv \lim_{\epsilon \rightarrow 0} \frac{d \left[ R \left( \tilde{T}(\circ; y^*, \epsilon, \eta) \right) \right]}{d\eta} \Big|_{\eta=0}$$

is the marginal cost of providing additional resources to those with earnings  $y^*$  in the alternative environment.

## B.2 Proof of Proposition 1

**Statement of Proposition** For any  $\epsilon > 0$  define the scaled surplus by  $s_\epsilon(y) = \epsilon s(y)$  and  $S_\epsilon = E[s_\epsilon(y)g(y)] = \epsilon S$ . If  $S < 0$ , there exists an  $\tilde{\epsilon} > 0$  such that for any  $\epsilon < \tilde{\epsilon}$  there exists an augmentation to the tax schedule in the status quo environment that generates surplus,  $s_\epsilon^t(y)$ , that is uniformly greater than the surplus offered by the alternative environment,  $s_\epsilon^t(y) > s_\epsilon(y)$  for all  $y$ . Conversely, if  $S > 0$ , no such  $\tilde{\epsilon}$  exists.

**Proof** The strategy of the proof is to consider a modification to the tax schedule that gives a discrete tax cut to each interval of the income distribution that makes everyone in the interval better off relative to the alternative environment. I show that when  $S < 0$  and for sufficiently small  $\epsilon$ , one can find sufficiently fine partitions that lead to feasible modifications of the tax schedule that make everyone better off relative to the alternative environment.

More formally, suppose  $S < 0$ . Then,

$$\int s(y)g(y(\theta))d\mu(\theta) < 0$$

so that

$$\int s_\epsilon(y)g(y(\theta))d\mu(\theta) = \epsilon \int s(y)g(y(\theta))d\mu(\theta) < 0$$

For any tax schedule  $\hat{T}$ , let  $y(\theta; \hat{T})$  denote the choice of earnings by type  $\theta$  facing tax schedule  $\hat{T}$ . Given these choices, total tax revenue is given by

$$R(\hat{T}) = \int \hat{T}(y(\theta; \hat{T}))d\mu(\theta)$$

Now, consider an augmented tax schedule. Let  $P = \{P_j\}_{j=1}^{N_P}$  denote a partition of the income distribution into intervals and let  $\eta_j^P$  denote transfers provided to each such region of the income distribution:

$$\hat{T}_\epsilon^P(y) = T(y) - \epsilon \sum_{j=1}^{N_P} \eta_j^P 1\{y \in P_j\}$$

and let

$$\hat{T}_\epsilon^j(y) = T(y) - \epsilon \eta_j^P 1\{y \in P_j\}$$

Note that aggregate government revenue is the sum of revenue from the partition. Assumption 1, combined with the linearity properties of the derivative implies that

$$\frac{d}{d\epsilon}\Big|_{\epsilon=0} R(\hat{T}_\epsilon^P) = \sum_{j=1}^{N_P} \frac{d}{d\epsilon}\Big|_{\epsilon=0} R(\hat{T}_\epsilon^j) \quad (10)$$

Note that each partition can be represented as

$$P_j = [y_j^* - \epsilon_j, y_j^* + \epsilon_j)$$

so that

$$\frac{d}{d\epsilon}|_{\epsilon=0} R(\hat{T}_\epsilon^j) = -\eta_j^P \frac{d}{d\eta}|_{\eta=0} \hat{q}(y_j^*, \epsilon_j, \eta) \Pr\{y(\theta) \in P_j\}$$

where  $\Pr\{y(\theta) \in P_j\} = \mu\{y^{-1}(P_j)\}$  and  $\hat{q}$  is defined in equation (9) above.

Now, define  $\eta_j^P$  as

$$\eta_j^P = \sup\{s(y) | y \in P_j\} - \frac{SID}{2} \bar{g}$$

where  $\bar{g} = E[g(y)]$  is the average value of the marginal cost of taxation. Let  $s_\epsilon^t(y)$  denote the surplus the individual earning  $y$  obtains when facing tax schedule  $\hat{T}_\epsilon^j$ . By the envelope theorem (and the assumption of no externalities / GE effects), there exists  $\tilde{\epsilon}$  such that for all  $\epsilon < \tilde{\epsilon}$ , the individual obtains surplus at least as large as  $\epsilon(\sup\{y | y \in P_j\})$

$$s_\epsilon^t(y) > \epsilon(\sup\{y | y \in P_j\}) \quad (11)$$

for all  $\epsilon \in (0, \tilde{\epsilon})$  (to see this, note that the tax augmentation not only gives people surplus  $s^y$  but also provides  $-\frac{S}{2} > 0$ ; so this inequality is made strict).

Now, consider the marginal cost of the policy. By construction

$$\frac{d}{d\epsilon}|_{\epsilon=0} R(\hat{T}_\epsilon^P) = -\sum_{j=1}^{N_P} \eta_j^P \frac{d}{d\eta}|_{\eta=0} \hat{q}(y_j^*, \epsilon_j, \eta) \Pr\{y(\theta) \in P_j\}$$

and taking the limit as partition widths go to zero,

$$\lim_{width(P) \rightarrow 0} \frac{d}{d\epsilon}|_{\epsilon=0} R(\hat{T}_\epsilon^P) = -\lim_{width(P) \rightarrow 0} \sum_{j=1}^{N_P} \left( (\sup\{y | y \in P_j\}) - \frac{S}{2} \right) \frac{d}{d\eta}|_{\eta=0} \hat{q}(y_j^*, \epsilon_j, \eta) \Pr\{y(\theta) \in P_j\}$$

Note that the terms inside the sum have limits that exist and are unique (because  $g(y)$  is assumed to be continuous and the mean surplus function is assumed to be continuous). Note in principle this limit existing does not require continuity of either the surplus function or the marginal cost function  $g(y)$  – some suitable integrability condition would work – but this is sufficient. So,

$$\begin{aligned} \lim_{width(P) \rightarrow 0} \frac{d}{d\epsilon}|_{\epsilon=0} R(\hat{T}_\epsilon^P) &= -\int s(y(\theta)) g(y(\theta)) d\mu(\theta) + \frac{S}{2} \bar{g} \\ &= -S\bar{g} + \frac{S}{2} \bar{g} \\ &= -\frac{S}{2} \bar{g} \end{aligned}$$

which is positive. Therefore, there exists  $\epsilon^* < \tilde{\epsilon}$  such that for all  $\epsilon \in (0, \epsilon^*)$  we have  $R(\hat{T}_\epsilon^P) > 0$  and  $s_\epsilon^t(y) > s_\epsilon(y)$  for all  $y$ .

**Converse** Now suppose  $S > 0$ . Then,

$$\int s(y(\theta))g(y(\theta))d\mu(\theta) > 0$$

And, suppose for contradiction that some  $\tilde{\epsilon}$  exists so that there are a set of tax schedules,  $\hat{T}_\epsilon$ , that deliver greater surplus along the income distribution,  $\frac{d}{d\epsilon}|_{\epsilon=0}s_\epsilon^t(y) \geq s(y)$ . I will show that this implies the tax schedule modification is not budget neutral for sufficiently small  $\epsilon$ .

Note that the envelope theorem implies  $\frac{d}{d\epsilon}|_{\epsilon=0}\hat{T}_\epsilon(y) = \frac{d}{d\epsilon}|_{\epsilon=0}s_\epsilon^t(y)$  for all  $y$ . For any  $\epsilon > 0$  and  $\gamma > 0$ , one can approximate the revenue function using a partition  $P^\gamma = \{P_j^\gamma\}_{j=1}^{N_{P^\gamma}}$  and a step tax function  $T_\epsilon^{P^\gamma}$  that provides exactly  $E[s_\epsilon^t(y(\theta)) | y(\theta) \in P_j^\gamma]$  units of tax reduction. Therefore, the marginal cost of the policy is approximated by  $\frac{d}{d\epsilon}|_{\epsilon=0}R(\hat{T}_\epsilon^{P^\gamma})$ ,

$$\left| \frac{d}{d\epsilon}|_{\epsilon=0}R(\hat{T}_\epsilon) - \frac{d}{d\epsilon}|_{\epsilon=0}R(\hat{T}_\epsilon^{P^\gamma}) \right| < \gamma$$

where

$$\frac{d}{d\epsilon}|_{\epsilon=0}R(\hat{T}_\epsilon^{P^\gamma}) = - \sum_{j=1}^{N_{P^\gamma}} \frac{d}{d\epsilon}|_{\epsilon=0}E[s_\epsilon^t(y(\theta)) | y(\theta) \in P_j^\gamma] \frac{d}{d\eta}|_{\eta=0}\hat{q}(y_j^*, \epsilon_j, \eta) \Pr\{y(\theta) \in P_j^\gamma\}$$

where  $P_j^\gamma = [y_{\gamma,j}^* - \epsilon_{\gamma,j}, y_{\gamma,j}^* + \epsilon_{\gamma,j})$ . For sufficiently small  $\epsilon$  we know that  $E[s_\epsilon^t(y(\theta)) | y(\theta) \in P_j^\gamma] > E[s_\epsilon(y(\theta)) | y(\theta) \in P_j^\gamma]$ . Therefore,

$$\frac{d}{d\epsilon}|_{\epsilon=0}E[s_\epsilon^t(y(\theta)) | y(\theta) \in P_j^\gamma] > \frac{d}{d\epsilon}|_{\epsilon=0}E[s_\epsilon(y(\theta)) | y(\theta) \in P_j^\gamma] - \gamma$$

since both have values of zero when  $\epsilon = 0$ .

Therefore,

$$\frac{d}{d\epsilon}|_{\epsilon=0}R(\hat{T}_\epsilon^P) < - \sum_{j=1}^N \frac{d}{d\epsilon}|_{\epsilon=0}E[s_\epsilon(y(\theta)) | y(\theta) \in P_j^\gamma] \frac{d}{d\eta}|_{\eta=0}\hat{q}(y_j^*, \epsilon_j, \eta) \Pr\{y(\theta) \in P_j^\gamma\} + \gamma$$

and taking the limit as the partition widths converge towards zero (so that  $\gamma \rightarrow 0$ ), we arrive as

$$\frac{d}{d\epsilon}|_{\epsilon=0}R(\hat{T}_\epsilon) \leq -SE[g(\theta)] < 0$$

so that the policy is not budget neutral.

**Discussion** The proof relied on two key assumptions. First, I assume that providing a small amount of money through modifications in the tax schedule generates surplus of at least the mechanical amount of money provided in the absence of any behavioral response. This follows from the envelope theorem,

combined with the assumption that infinitesimal tax changes in one portion of the income distribution do not affect the welfare of anyone at other points of the distribution. This was implicitly assumed by writing the utility function as a function of one's own consumption and earnings, and not a function of anyone else's choices of labor supply or earnings. For example, if taxing the rich caused them to reduce their earnings which in turn increased the wages of the poor, then equation (11) would no longer hold, since individuals outside of the intended target of the tax transfers would have surplus impacts. Accounting for such general equilibrium effects is an interesting and important direction for both theoretical and empirical work.

Second, I assume that the revenue function is continuously differentiable and additive in modifications to the tax schedule. This is primarily a technical assumption that rules out types that are indifferent to many points along the income distribution (which would cause them to be double-counted as costs in equation (10)).

### B.3 Proof of Proposition 2

**Statement of Proposition** *Suppose Assumption 1 holds. For  $\epsilon > 0$ , let  $s_\epsilon = \epsilon s(y)$ . If  $S > 0$ , there exists  $\tilde{\epsilon} > 0$  such that for any  $\epsilon < \tilde{\epsilon}$ , there exists an augmentation to the tax schedule in the alternative environment that delivers surplus  $s_\epsilon^t(y)$  that is positive at all points along the income distribution,  $s_\epsilon^t(y) > 0$  for all  $y$ . Conversely, if  $S < 0$ , then no such  $\tilde{\epsilon}$  exists.*

**Proof** I provide the brief sketch here that does not go through the formality of defining the partitions as in the proof above, but one can do so analogously to the proof of Proposition 1. Let  $y(\theta)$  continue to denote the choice of income of a type  $\theta$  in the status quo environment, which may differ from their choice of  $y$  in the alternative environment. To capture this, let  $y_\epsilon^\alpha(y)$  denote the choice of income in the alternative environment made by those who chose  $y$  in the status quo environment. Per Assumption 1, this function is a bijection. Given the surplus function,  $s(y)$ , consider a modification to the income distribution that taxes away all but  $\epsilon \frac{S}{2}$  of this surplus to those earning  $y$  in the status quo (i.e. those earning  $y_\epsilon^\alpha(y)$  in the  $\epsilon$ -alternative environment). If  $\tilde{T}_\epsilon$  is the tax schedule in the  $\epsilon$ -alternative environment, then the modified tax schedule is

$$\hat{T}_\epsilon(y) = \tilde{T}_\epsilon(y) + \epsilon \left( s((y_\epsilon^\alpha)^{-1}(y)) - \frac{S}{2} \right)$$

Let  $s_\epsilon^t(y)$  denote the surplus of the tax-modified  $\epsilon$ -alternative environment with tax schedule  $\hat{T}_\epsilon(y)$ . For sufficiently small  $\epsilon$ , the off-setting transfer ensures everyone is better off relative to the status quo (note this relies on the fact that  $S > 0$ , so that there is aggregate surplus to spread around). Hence,  $s_\epsilon^t(y) > s_\epsilon(y) - E[s_\epsilon(y)]$  for sufficiently small  $\epsilon$ ; and taking the expectation conditional on  $y(\theta) = y$  yields

$$E[s_\epsilon^t(y)] > 0 \quad \forall y$$

Now, one needs to show that, for sufficiently small  $\epsilon$ , the cost of the modification to the tax

schedule is not budget-negative. Note that for each  $y$ , the tax modification provides a transfer of  $s_\epsilon \left( (y_\epsilon^\alpha)^{-1}(y) \right)$ . Note that Assumption 1, the marginal cost of implementing these surplus transfers is the same as in the status quo environment.

$$\frac{dR}{d\epsilon}|_{\epsilon=0} = \int s(y) g(y) dF(y) - \frac{S}{2} = \frac{S}{2} > 0$$

so that the transfer scheme is feasible for sufficiently small  $\epsilon$ .

## B.4 Proof of Proposition 3

This subsection provides the derivation of the marginal cost of taxation using elasticities when the utility function satisfies Assumption 3.

### B.4.1 Continuous Responses Only

To begin, I assume there is no participation margin response. Specifically, I assume that preferences are convex in consumption-earnings space so that  $\hat{y}(\theta; y^*, \epsilon, \eta)$  is continuously differentiable in  $\eta$ . Below, I add back in extensive margin responses that allow types  $\theta$  to move to/from 0 and a point of interior earnings,  $y > 0$ , in response to a change in the size of the tax cut,  $\eta$ .

A key source of complexity is that individuals may have different curvatures of their utility function. To capture this, define  $c(y; \theta)$  to be the individual  $\theta$ 's indifference curve in consumption-earnings space at the baseline utility level. Given an agent  $\theta$ 's choice  $y(\theta)$  facing the baseline tax schedule  $T(y)$ , the indifference curve solves

$$u(c(y; \theta), y; \theta) = u(T(y(\theta)) - y(\theta), y(\theta); \theta)$$

Note that the individual's first order condition requires:

$$c'(y(\theta); \theta) = -\frac{u_y}{u_c} = 1 - T'(y(\theta)) \quad (12)$$

so that the slope of this indifference curve equals the marginal keep rate,  $1 - T'$ .

In addition, the curvature of this indifference curve governs the size of the fraction of people who change their behavior in order to obtain the transfer,  $\eta$ . Let  $k(\theta) = c''(y(\theta); \theta)$  denote the curvature of the indifference curve of type  $\theta$  in the status quo world. First, consider those whose baseline income is just above  $y^* + \frac{\epsilon}{2}$  but the opportunity to obtain the  $\eta$  transfer induces them to drop their income down to  $y^* + \frac{\epsilon}{2}$ . For individuals with curvature  $k$ , a second-order expansion of  $c$  (i.e. first order expansion of  $c'$ ) shows that anyone between  $y^* + \frac{\epsilon}{2}$  and  $y^* + \frac{\epsilon}{2} + \gamma(\eta; k)$  will choose incomes at  $y^* + \frac{\epsilon}{2}$ , where  $\gamma(\eta; k)$  solves

$$\frac{(\gamma(\eta; k))^2}{2} k = \eta$$

or

$$\gamma(\eta; k) = \sqrt{\frac{2\eta}{k}}$$

Similarly, for individuals with curvature  $k$ , those with incomes between  $y^* - \frac{\epsilon}{2} - \gamma(\eta; k)$  and  $y^* - \frac{\epsilon}{2}$  will choose to increase their incomes to  $y^* - \frac{\epsilon}{2}$ .

Given these definitions, one can write the budget cost as the sum of four terms:

$$\int_{\theta} \hat{T}(\hat{y}(\theta; y^*, \epsilon, \eta); y^*, \epsilon, \eta) d\mu(\theta) = A + B + C + D + o(\eta)$$

where  $\lim_{\eta \rightarrow 0} \frac{o(\eta)}{\eta} = 0$  (so that  $\frac{do}{d\eta}|_{\eta=0} = 0$ , so that one can ignore this term in the calculation of  $\frac{d\hat{q}}{d\eta}|_{\eta=0}$ ).

The first term,  $A$  is the mechanical cost that must be paid to all those who receive the  $\eta$  transfer.

$$A = \eta \int 1 \left\{ y(\theta) \in \left( y^* - \frac{\epsilon}{2} - \sqrt{\frac{2\eta}{k(\theta)}}, y^* + \frac{\epsilon}{2} + \sqrt{\frac{2\eta}{k(\theta)}} \right) \right\} d\mu(\theta)$$

The second term is the cost from those with baseline earnings above  $y^* + \frac{\epsilon}{2}$  who drop their income down to  $y^* + \frac{\epsilon}{2}$ ,

$$B = \int \left( T\left(y^* + \frac{\epsilon}{2}\right) - T(y(\theta)) \right) 1 \left\{ y(\theta) \in \left( y^* + \frac{\epsilon}{2}, y^* + \frac{\epsilon}{2} + \sqrt{\frac{2\eta}{k(\theta)}} \right) \right\} d\mu(\theta)$$

And, conversely, the third term is from those with baseline earnings below  $y^* - \frac{\epsilon}{2}$  who increase their incomes to  $y^* - \frac{\epsilon}{2}$ ,

$$C = \tau \left( y - \frac{\epsilon}{2} \right) \int \left[ y(\theta) - \left( y^* - \frac{\epsilon}{2} \right) \right] 1 \left\{ y(\theta) \in \left( y^* - \frac{\epsilon}{2} - \sqrt{\frac{2\eta}{k(\theta)}}, y^* - \frac{\epsilon}{2} \right) \right\} d\mu(\theta)$$

and finally the fourth term is the income effect on earnings for those with baseline earnings in the  $\epsilon$ -region near  $y^*$ ,

$$D = \int [T(\hat{y}(\theta; y^*, \epsilon, \eta)) - T(y)] 1 \left\{ y(\theta) \in \left( y^* - \frac{\epsilon}{2}, y^* + \frac{\epsilon}{2} \right) \right\} d\mu(\theta)$$

The remaining term,  $o(\eta)$ , captures the bias from approximating the  $B$  and  $C$  terms using the second-order expansion for  $c(y; \theta)$ .

Clearly,

$$\frac{d \left[ \int_{\theta} \hat{T}(\hat{y}(\theta; y^*, \epsilon, \eta); y^*, \epsilon, \eta) d\mu(\theta) \right]}{d\eta} \Big|_{\eta=0} = \frac{dA}{d\eta} \Big|_{\eta=0} + \frac{dB}{d\eta} \Big|_{\eta=0} + \frac{dC}{d\eta} \Big|_{\eta=0} + \frac{dD}{d\eta} \Big|_{\eta=0}$$

I characterize each of these terms. After doing so, one can divide by  $F(y^* + \frac{\epsilon}{2}) - F(y^* - \frac{\epsilon}{2})$  and take the limit as  $\epsilon \rightarrow 0$  to arrive at the expression for  $\lim_{\epsilon \rightarrow 0} \frac{d\hat{q}}{d\eta} \Big|_{\eta=0}$ .

**Characterizing**  $\frac{dA}{d\eta}|_{\eta=0}$  First, I show that  $\frac{dA}{d\eta}|_{\eta=0} = F\left(y^* + \frac{\epsilon}{2}\right) - F\left(y^* - \frac{\epsilon}{2}\right)$ .

To see this, first write  $A$  by conditioning on  $k(\theta)$ . Formally, recall that  $\mu(\theta)$  is the measure on the type space. Let  $\mu_{\theta|k}(\theta|k)$  denote the measure of  $\theta$  conditional on having curvature  $k$  (i.e.  $c''(y(\theta)) = k$ ) and let  $\mu_k(k)$  denote the measure of those having curvature  $k$ .<sup>27</sup> Then,

$$A = -\eta \int_k \int_{\theta|k} 1 \left\{ y(\theta) \in \left( y^* - \frac{\epsilon}{2} - \sqrt{\frac{2\eta}{k}}, y^* + \frac{\epsilon}{2} + \sqrt{\frac{2\eta}{k}} \right) \right\} d\mu_{\theta|k}(\theta|k(\theta) = k) d\mu_k(k)$$

Taking a derivative yields

$$\frac{dA}{d\eta} = -F\left(y^* + \frac{\epsilon}{2} - \sqrt{\frac{2\eta}{k}}\right) + F\left(y^* - \frac{\epsilon}{2} + \sqrt{\frac{2\eta}{k}}\right) - \int_k \eta \left[ f_{y|k}\left(y^* + \frac{\epsilon}{2} + \sqrt{\frac{2\eta}{k}}|k\right) \sqrt{\frac{1}{2\eta k}} - f_{y|k}\left(y^* - \frac{\epsilon}{2} - \sqrt{\frac{2\eta}{k}}|k\right) \sqrt{\frac{1}{2\eta k}} \right] d\mu_k(k)$$

where  $f_{y|k}(y|k)$  is the density of  $y(\theta)$  given  $k(\theta)$ . Note that one can re-write the second term in a manner that makes it clear that it is proportional to  $\sqrt{\eta}$ :

$$\frac{dA}{d\eta} = -F\left(y^* + \frac{\epsilon}{2} - \sqrt{\frac{2\eta}{k}}\right) + F\left(y^* - \frac{\epsilon}{2} + \sqrt{\frac{2\eta}{k}}\right) - \sqrt{\eta} \left[ \int_k \left[ f_{y|k}\left(y^* + \frac{\epsilon}{2} + \sqrt{\frac{2\eta}{k}}|k\right) \sqrt{\frac{1}{2k}} - f_{y|k}\left(y^* - \frac{\epsilon}{2} - \sqrt{\frac{2\eta}{k}}|k\right) \sqrt{\frac{1}{2k}} \right] d\mu_k(k) \right]$$

Therefore, evaluating at  $\eta = 0$  yields

$$\frac{dA}{d\eta}|_{\eta=0} = - \left[ F\left(y^* + \frac{\epsilon}{2}\right) - F\left(y^* - \frac{\epsilon}{2}\right) \right]$$

**Characterizing**  $\frac{dB}{d\eta}|_{\eta=0}$  To see this, note that

$$\begin{aligned} \frac{dB}{d\eta} &= \frac{d}{d\eta} \int_k \int_{\theta|k} \left( T\left(y^* + \frac{\epsilon}{2}\right) - T(y(\theta)) \right) 1 \left\{ y(\theta) \in \left( y^* + \frac{\epsilon}{2}, y^* + \frac{\epsilon}{2} + \sqrt{\frac{2\eta}{k}} \right) \right\} d\mu_{\theta|k}(\theta|k(\theta) = k) d\mu_k(k) \\ &= \int_k \left( T\left(y^* + \frac{\epsilon}{2}\right) - T\left(y^* + \frac{\epsilon}{2} + \sqrt{\frac{2\eta}{k}}\right) \right) \sqrt{\frac{1}{2\eta k}} f_{y|k}\left(y^* + \frac{\epsilon}{2} + \sqrt{\frac{2\eta}{k}}|k\right) d\mu_k(k) \end{aligned}$$

which follows from differentiating at the upper endpoint  $y^* + \frac{\epsilon}{2} + \sqrt{\frac{2\eta}{k}}$  after conditioning on curvature  $k$ . Re-writing yields

$$\frac{dB}{d\eta} = \int_k \frac{T\left(y^* + \frac{\epsilon}{2}\right) - T\left(y^* + \frac{\epsilon}{2} + \sqrt{\frac{2\eta}{k}}\right)}{\sqrt{\frac{2\eta}{k}}} \frac{1}{k} f_{y|k}\left(y^* + \frac{\epsilon}{2} + \sqrt{\frac{2\eta}{k}}|k\right) d\mu_k(k)$$

<sup>27</sup>In other words, for any function of the type space and level of curvature,  $r(\theta, k(\theta))$ , one has

$$\int \int r(\theta, k(\theta)) d\mu_{\theta|k}(\theta|k(\theta)) d\mu_k(k(\theta)) = \int r(\theta, k(\theta)) d\mu(\theta)$$

so that one can either integrate over  $\theta$  (RHS) or one can first condition on curvature (and integrate over  $\theta$  given curvature  $k(\theta)$ ) and then integrate over curvature,  $k(\theta)$ .

Now, evaluating as  $\eta \rightarrow 0$ , yields

$$\begin{aligned}\frac{dB}{d\eta}\Big|_{\eta=0} &= -T' \left( y^* + \frac{\epsilon}{2} \right) \int_k \frac{f_{y|k} \left( y^* + \frac{\epsilon}{2} | k \right)}{k} d\mu_k(k) \\ &= -T' \left( y^* + \frac{\epsilon}{2} \right) E \left[ \frac{1}{k(\theta)} | y(\theta) = y^* + \frac{\epsilon}{2} \right] f \left( y^* + \frac{\epsilon}{2} \right)\end{aligned}$$

so that tax revenue is decreased by individuals decreasing their income down to  $y^* + \frac{\epsilon}{2}$  in order to get the  $\eta$  transfer.

**Characterizing**  $\frac{dC}{d\eta}\Big|_{\eta=0}$  Analogous to the calculation for  $\frac{dB}{d\eta}\Big|_{\eta=0}$ , it is possible to show that

$$\frac{dC}{d\eta}\Big|_{\eta=0} = T' \left( y^* - \frac{\epsilon}{2} \right) E \left[ \frac{1}{k(\theta)} | y(\theta) = y^* - \frac{\epsilon}{2} \right] f \left( y^* - \frac{\epsilon}{2} \right)$$

so that tax revenue is increased because individuals move from below  $y^* - \frac{\epsilon}{2}$  up to  $y^* - \frac{\epsilon}{2}$  in order to get the  $\eta$  transfer.

**Characterizing**  $\frac{dD}{d\eta}\Big|_{\eta=0}$  Finally, I show that

$$\frac{dD}{d\eta}\Big|_{\eta=0} = \left[ F \left( y^* + \frac{\epsilon}{2} \right) - F \left( y^* - \frac{\epsilon}{2} \right) \right] E \left[ \frac{dy}{d\eta} T'(y(\theta)) | y(\theta) \in \left[ y^* - \frac{\epsilon}{2}, y^* + \frac{\epsilon}{2} \right] \right]$$

so that  $\frac{dD}{d\eta}\Big|_{\eta=0}$  is proportional to the average income effects near  $y^*$ .

To see this, note that

$$\frac{dD}{d\eta} = \frac{d}{d\eta} \int [T(\hat{y}) - T(y)] 1 \left\{ y(\theta) \in \left( y^* - \frac{\epsilon}{2}, y^* + \frac{\epsilon}{2} \right) \right\} dF(\theta)$$

Note that for these individuals in the  $\epsilon$  region near  $y^*$  they only receive an income effect from the policy change. Therefore, we have

$$\frac{dD}{d\eta}\Big|_{\eta=0} = \int T'(y(\theta)) \frac{d\hat{y}}{d\eta}\Big|_{\eta=0} 1 \left\{ y(\theta) \in \left( y^* - \frac{\epsilon}{2}, y^* + \frac{\epsilon}{2} \right) \right\} dF(\theta)$$

where  $\frac{d\hat{y}}{d\eta}\Big|_{\eta=0}$  is the effect of an additional dollar of after-tax income on labor supply. One can define the income elasticity by multiplying by the after-tax price,

$$\zeta(\theta) = (1 - T'(y)) \frac{d\hat{y}}{d\eta}\Big|_{\eta=0}$$

so that

$$\frac{dD}{d\eta}\Big|_{\eta=0} = \int \frac{T'(y(\theta))}{1 - T'(y(\theta))} \zeta(\theta) 1 \left\{ y(\theta) \in \left( y^* - \frac{\epsilon}{2}, y^* + \frac{\epsilon}{2} \right) \right\} dF(\theta)$$

**Taking**  $\epsilon \rightarrow 0$  Now, to take the limit as  $\epsilon \rightarrow 0$ , note that

$$\frac{dB}{d\eta}|_{\eta=0} + \frac{dC}{d\eta}|_{\eta=0} = -T'(y^* + \frac{\epsilon}{2}) E \left[ \frac{1}{k(\theta)} |y(\theta) = y^* + \frac{\epsilon}{2} \right] f(y^* + \frac{\epsilon}{2}) + T'(y^* - \frac{\epsilon}{2}) E \left[ \frac{1}{k(\theta)} |y(\theta) = y^* - \frac{\epsilon}{2} \right] f(y^* - \frac{\epsilon}{2})$$

so that

$$= \lim_{\epsilon \rightarrow 0} \left( \frac{\frac{dB}{d\eta}|_{\eta=0} + \frac{dC}{d\eta}|_{\eta=0}}{F(y^* + \frac{\epsilon}{2}) - F(y^* - \frac{\epsilon}{2})} \right) = \lim_{\epsilon \rightarrow 0} \frac{\frac{dB}{d\eta}|_{\eta=0} + \frac{dC}{d\eta}|_{\eta=0}}{F(y^* + \frac{\epsilon}{2}) - F(y^* - \frac{\epsilon}{2})} = \lim_{\epsilon \rightarrow 0} \left( \frac{\epsilon}{F(y^* + \frac{\epsilon}{2}) - F(y^* - \frac{\epsilon}{2})} \right) \left( \frac{-T'(y^* + \frac{\epsilon}{2}) E \left[ \frac{1}{k(\theta)} |y(\theta) = y^* + \frac{\epsilon}{2} \right] f(y^* + \frac{\epsilon}{2}) + T'(y^* - \frac{\epsilon}{2}) E \left[ \frac{1}{k(\theta)} |y(\theta) = y^* - \frac{\epsilon}{2} \right] f(y^* - \frac{\epsilon}{2})}{\epsilon} \right)$$

or

$$\lim_{\epsilon \rightarrow 0} \frac{\frac{dB}{d\eta}|_{\eta=0} + \frac{dC}{d\eta}|_{\eta=0}}{F(y^* + \frac{\epsilon}{2}) - F(y^* - \frac{\epsilon}{2})} = \frac{1}{f(y^*)} \left( -\frac{d}{dy} \Big|_{y=y^*} \left[ T'(y) E \left[ \frac{1}{k(\theta)} |y(\theta) = y \right] f(y) \right] \right)$$

Now, note also that

$$\lim_{\epsilon \rightarrow 0} \frac{-\frac{dA}{d\eta}|_{\eta=0}}{F(y^* + \frac{\epsilon}{2}) - F(y^* - \frac{\epsilon}{2})} = 1$$

and

$$\lim_{\epsilon \rightarrow 0} \frac{-\frac{dD}{d\eta}|_{\eta=0}}{F(y^* + \frac{\epsilon}{2}) - F(y^* - \frac{\epsilon}{2})} = -T'(y^*) E \left[ \frac{d\hat{y}}{d\eta} \Big|_{\eta=0} |y(\theta) = y^* \right]$$

which is given by the average income effect at  $y^*$  multiplied by the marginal tax rate.

Combining,

$$\lim_{\epsilon \rightarrow 0} \frac{d\hat{q}(y^*, \epsilon, \eta)}{d\eta} \Big|_{\eta=0} = 1 + \frac{1}{f(y^*)} \frac{d}{dy} \Big|_{y=y^*} \left[ T'(y) E \left[ \frac{1}{k(\theta)} |y(\theta) = y \right] f(y) \right] - \frac{T'(y)}{1 - T'(y)} E[\zeta(\theta) |y(\theta) = y]$$

**Replacing curvature with compensated elasticity** Now, note that the curvature,  $k$ , is related to the compensated elasticity of earnings. To see this, note that

$$c'(y(\theta); \theta) = 1 - \tau$$

where  $\tau$  is the marginal tax rate faced by the individual,  $\tau = T'(y(\theta))$ . Totally differentiating with respect to one minus the marginal tax rate yields

$$c''(y(\theta)) \frac{dy^c}{d(1-\tau)} = 1$$

where  $\frac{dy^c}{d(1-\tau)}$  is the compensated response to an increase in the marginal keep rate,  $1 - \tau$ . Re-writing,

$$\frac{dy^c}{d(1-\tau)} = \frac{1}{c''(y(\theta))}$$

Intuitively, the size of a compensated response to a price change is equal to the inverse of the curvature of the indifference curve.

Now, by definition, the compensated elasticity of earnings is given by

$$\epsilon^c(\theta) = \frac{dy^c}{d(1-\tau)} \frac{(1-\tau)}{y(\theta)} = \frac{1}{c''(y)} \frac{1-\tau}{y}$$

or

$$\frac{1}{k(\theta)} = \epsilon^c(\theta) \frac{y(\theta)}{1-T'(y(\theta))}$$

where  $\epsilon^c(\theta)$  is the compensated elasticity of type  $\theta$  defined locally around the status quo tax schedule.

Replacing  $\frac{1}{k(\theta)}$  in the main equation yields

$$\lim_{\epsilon \rightarrow 0} \frac{d\hat{q}(y^*, \epsilon, \eta)}{d\eta} \Big|_{\eta=0} = 1 + \frac{1}{f(y^*)} \frac{d}{dy} \Big|_{y=y^*} \left[ T'(y) E \left[ \epsilon^c(\theta) \frac{y(\theta)}{1-T'(y(\theta))} \Big|_{y(\theta)=y} \right] f(y) \right] - \frac{T'(y)}{1-T'(y)} E[\zeta(\theta) \Big|_{y(\theta)=y}]$$

#### B.4.2 Adding a Participation Margin

Heretofore, I have ignored the potential for extensive margin responses. Put differently, I assumed everyone's intensive margin first order condition (equation (12)) held. Now, I show how one can overlay participation margin responses for people who move in and out of the labor force in response to changes in the tax schedule.

For simplicity, consider an alternative world where  $y = 0$  was removed from individuals' feasibility set. Let  $y^P(\theta)$  denote the earnings choice of type  $\theta$  in this restricted world. Clearly,  $y^P(\theta)$  solves

$$y^P(\theta) = \operatorname{argmax}_{y>0} u(y - T(y), y; \theta)$$

For all types in the labor force in the status quo world,  $y^P(\theta) = y(\theta)$ . For those out of the labor force,  $y(\theta) = 0$ . I retain the assumption that preferences are convex over the region  $y > 0$ . Therefore,  $y^P(\theta)$  is continuously differentiable in response to changes in the tax schedule,  $T$ . So, I allow for discrete moves between 0 and  $y > 0$ , but do not allow discrete moves across two different labor supply points in response to small changes in the tax schedule.

Given  $y^P(\theta)$ , let  $c^P(\theta)$  denote the consumption level required by type  $\theta$  to enter into the labor force to earn  $y^P(\theta)$ :

$$u(c^P(\theta), y^P(\theta); \theta) = u(y - T(0), 0; \theta)$$

Given  $y^P(\theta)$  and  $c^P(\theta)$ , one can define the labor force participation rate at each point along the income distribution. Note that an individual of type  $\theta$  chooses to work whenever

$$c^P(\theta) \leq y^P(\theta) - T(y^P(\theta))$$

For any consumption and income level,  $(c, y)$ , let  $LFP(c, y)$  denote the fraction of individuals with  $y^P(\theta) = y$  who choose to work,  $y(\theta) = y$ :

$$LFP(c, y) = \int \mathbf{1}\{c \geq c^P(\theta)\} d\mu(\theta | y^P(\theta) = y)$$

With this definition, one can write

$$\hat{q}(y^*, \epsilon, \eta) = A + B + C + D + P + o(\eta)$$

where  $P$  is the cost resulting from non-marginal changes in labor supply and  $\frac{dP}{d\eta}|_{\eta=0, \epsilon=0}$  is given by

$$\begin{aligned} \frac{dP}{d\eta}|_{\eta=0} &= \frac{d}{d\eta}|_{\eta=0} \int_{y^P(\theta) \in [y^* - \frac{\epsilon}{2}, y^* + \frac{\epsilon}{2}]} [(T(y^P(\theta)) - \eta - T(0)) LFP\{y^P - T(y^P) + \eta, y^P(\theta)\} dF(\theta)] \\ &= \int_{y(\theta) \in [y^* - \frac{\epsilon}{2}, y^* + \frac{\epsilon}{2}]} \left[ (T(y(\theta)) - T(0)) \frac{dLFP\{y(\theta) - T(y(\theta)), y(\theta)\}}{dc} dF(\theta) \right] \end{aligned}$$

so that

$$\begin{aligned} \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \frac{dP}{d\eta}|_{\eta=0} &= E \left[ (T(y) - T(0)) \frac{dLFP(y)}{dc} |_{y^P(\theta) = y} \right] LFP(y) \\ &= \frac{T(y) - T(0)}{y - T(y)} \hat{\epsilon}(y^P) \\ &= \frac{T(y) - T(0)}{y - T(y)} \epsilon_c^{LFP}(y) \end{aligned}$$

where  $\hat{\epsilon}_c^{LFP}(y)$  is the semi-elasticity of labor force participation at  $y$  off of the base of all potential people who have  $y^P(\theta)$  as their most preferred earnings point. To align with A-D, we need to replace the distribution of  $y^P$  with the distribution of  $y$ , so that we must divide by  $LFP$ . Dividing by  $LFP(y)$ , this is equal to the elasticity of labor force participation at  $y^P(\theta)$

$$\epsilon_c^{LFP}(y) = \frac{1}{LFP(y - T(y), y)} \frac{\partial LFP(y - T(y), y)}{\partial c}$$

Therefore, we have

$$\lim_{\epsilon \rightarrow 0} \frac{d\hat{q}(y^*, \epsilon, \eta)}{d\eta}|_{\eta=0} = 1 + \frac{1}{f(y^*)} \frac{d}{dy}|_{y=y^*} \left[ \frac{T'(y)}{1 - T'(y)} \epsilon^c(y) y f(y) \right] - \frac{T'(y)}{1 - T'(y)} \zeta(y) + \frac{T(y) - T(0)}{y - T(y)} \epsilon_c^{LFP}(y) \quad (13)$$

where

$$\epsilon^c(y) = E[\epsilon^c(\theta) | y(\theta) = y]$$

and

$$\zeta(y) = E[\zeta(\theta) | y(\theta) = y]$$

## C Inverse Optimum Derivation

Efficient social welfare weights correspond to the implicit welfare weights that rationalize the status quo tax schedule as optimal. To see this, let  $\chi(\theta)$  denote the social marginal utility of income of individual  $\theta$ , so that the marginal impact on social welfare of providing an additional \$1 of resources to type  $\theta$  is  $\chi(\theta)$ , which is normalized so that  $E[\chi(\theta)] = 1$ . Ratios of social marginal utilities of income,  $\frac{\chi(\theta_1)}{\chi(\theta_2)}$ , characterize the social willingness to pay to transfer resources from  $\theta_2$  to  $\theta_1$  and provide a generic local representation of social preferences (Saez and Stantcheva (2016)).

**Proposition 4.** *Suppose the income tax schedule in the status quo,  $T(y)$ , maximizes social welfare and let  $\chi(\theta)$  denote the local social marginal utilities of income. Then, the efficient welfare weights  $g(y)$ , equals the average social marginal utilities of income for those earning  $y(\theta) = y$ ,*

$$g(y) = E[\chi(\theta) | y(\theta) = y]$$

*Proof.* Given a tax function  $\hat{T}(y; y^*, \epsilon, \eta)$ , let  $\hat{v}(\theta, \epsilon, \eta)$  denote the utility to type  $\theta$ . By the envelope theorem, we have

$$\frac{d\hat{v}}{d\eta}\Big|_{\eta=0} = \begin{cases} 0 & \text{if } y \notin (y^* - \frac{\epsilon}{2}, y^* + \frac{\epsilon}{2}) \\ \frac{\partial v(\theta)}{\partial m} & \text{if } y \in (y^* - \frac{\epsilon}{2}, y^* + \frac{\epsilon}{2}) \end{cases}$$

so that the impact on the social welfare function is  $\int \chi(\theta) 1\{y(\theta) \in (y^* - \frac{\epsilon}{2}, y^* + \frac{\epsilon}{2})\} d\mu(\theta)$ , where  $\chi(\theta)$  equals  $\frac{\partial v(\theta)}{\partial m}$  multiplied by the local social welfare weight. Taking the limit as  $\epsilon \rightarrow 0$ , we have that the benefit of a small increase in  $\eta$  is  $E[\chi(\theta) | y(\theta) = y]$ ; moreover, by definition the cost of a small increase in  $\eta$  is  $g(y)$ . Optimality of the tax code implies that the welfare benefit per unit cost is equated for all  $y$ :

$$\frac{E[\chi(\theta) | y(\theta) = y_1]}{E[\chi(\theta) | y(\theta) = y_2]} = \frac{g(y_1)}{g(y_2)}$$

Finally, note that  $g(y) = \frac{E[\chi(\theta) | y(\theta) = y]}{E[\chi(\theta) | y(\theta) = y_2]} g(y_2)$ , so that  $E[g(y)] = \frac{E[\chi(\theta)]}{E[\chi(\theta) | y(\theta) = y_2]} g(y_2)$ . Now, by construction  $E[g(y)] = 1$  and  $E[\chi(\theta)] = 1$ , so replacing notation of  $y_2$  with  $y$  yields  $g(y) = E[\chi(\theta) | y(\theta) = y]$ .  $\square$

## D Heterogeneity

If two people earning the same income,  $y(\theta)$ , have different surplus,  $s(\theta)$ , then undoing the distributional incidence through the tax schedule will necessarily make one of the two people strictly better off. Fortunately, with a slight modification of the surplus function, one can use the efficient welfare weights to characterize the existence of local Pareto improvements.

Given the surplus function  $s(\theta)$  of interest, I define the min and max surplus at each point of the income distribution. First, for any  $\hat{y}$  let  $\underline{s}(\hat{y}) = \inf \{s(\theta) | y(\theta) = \hat{y}\}$  be the smallest surplus obtained by a type  $\theta$  that earns  $\hat{y}$  (note this number may be negative). Second, let  $\bar{s}(\hat{y}) = \sup \{s(\theta) | y(\theta) = \hat{y}\}$  be the largest surplus obtained by a type  $\theta$  that earns  $\hat{y}$ . The search for local Pareto improvements involves weighting not actual surplus,  $s(\theta)$ , but rather these min and max surplus functions conditional on income. In particular, let

$$\underline{S} = \int \underline{s}(y) g(y(\theta)) d\mu(\theta)$$

and

$$\bar{S} = \int \bar{s}(y) g(y(\theta)) d\mu(\theta)$$

If  $\bar{S} < 0$ , then there exists a modification to the existing tax schedule such that everyone locally prefers the modified status quo to the alternative environment.

**Proposition 5.** *Suppose  $\bar{S} < 0$ . Then, there exists an  $\tilde{\epsilon} > 0$  such that, for each  $\epsilon < \tilde{\epsilon}$  there exists a modification to the income tax schedule that delivers a Pareto improvement relative to  $s_\epsilon(\theta)$ . Conversely, if  $\bar{S} > 0$ , there exists an  $\tilde{\epsilon} > 0$  such that for each  $\epsilon < \tilde{\epsilon}$  any budget-neutral modification to the tax schedule results in lower surplus for some  $\theta$  relative to  $s_\epsilon(\theta)$ .*

*Proof.* The proof follows immediately by providing surplus  $\bar{s}_\epsilon(y) = \sup \{s_\epsilon(\theta) | y(\theta) = y\}$  instead of  $E[s_\epsilon(\theta) | y(\theta) = y]$  in the proof of Proposition 1.  $\square$

When  $\bar{S} < 0$ , a change in the tax schedule within the status quo locally Pareto dominates the alternative environment. Clearly,  $\bar{S} \geq S$  so that this is a more restrictive test of whether the status quo should be preferred to the alternative environment.

Conversely, using Assumption 1, one can test whether the alternative environment, modified with a change to the tax schedule, provides a local Pareto improvement relative to the status quo.

**Proposition 6.** *Suppose Assumption 1 holds. Suppose  $\underline{S} > 0$ . Then, there exists an  $\tilde{\epsilon} > 0$  such that, for each  $\epsilon < \tilde{\epsilon}$  there exists a modification to the income tax schedule in the alternative environment such that the modified alternative environment delivers positive surplus to all types relative to the status quo,  $s_\epsilon^t(\theta) > 0$  for all  $\theta$ .*

*Proof.* The proof follows immediately by providing surplus  $\underline{s}_\epsilon(y) = \inf \{s_\epsilon(\theta) | y(\theta) = y\}$  instead of  $s_\epsilon(y)$  in the proof of Proposition 2.  $\square$

In general, it can be the case that  $\bar{S} > 0 > \underline{S}$ , so that the potential Pareto criterion cannot lead to a sharp comparison between the status quo and the alternative environment.

**Corollary 1.** *Suppose  $s(\theta)$  does not vary with  $\theta$  conditional on income,  $y(\theta)$  (i.e.  $s(\theta) = \tilde{s}(y(\theta))$ ). Then,  $\bar{S} = \underline{S} = S$ .*

**Dealing with Heterogeneity in Practice** When surplus is heterogeneous conditional on income, it may be the case that  $\bar{S} > 0 > \underline{S}$ . In this case, there does not exist a modification to the tax schedule in the alternative or status quo environment that can render a Pareto comparisons between the status quo and alternative environment. Here, there are several options. First, one could bias the status quo, choosing the alternative environment iff  $\underline{S} > 0$ . Of course, this might be overly conservative. Second, one can use average surplus,  $S = E[s(\theta)g(y(\theta))]$ , and decide if the alternative environment brings sufficient benefits to each point of the income distribution to warrant the lack of Pareto improvement. This approach does not rely on the Pareto principle, but may be a useful application in cases with important sources of heterogeneity conditional on income.

Third, one could consider additional compensation instruments, such as capital taxation, commodity taxation, Medicaid eligibility, etc. Intuitively, when  $\bar{S} > 0 > \underline{S}$ , the income tax alone is too blunt an instrument to conduct compensating transfers. For example, if surplus is a function of both health and income, one could imagine making compensating transfers through modifications to both income and Medicaid / Medicare generosity and eligibility. Here, one requires estimates of  $FE(\mathbf{X})$  (e.g. if  $\mathbf{X} = (y, m)$  where  $m$  is Medicaid expenditures  $m$ , one requires the causal effect of the behavioral response to a transfer directed towards those not only with income  $y$  but also with Medicaid expenditures  $m$ . The key requirement is empirical estimation of the fiscal externalities.

Finally, one can consider policies that have smaller variations in surplus conditional on income. Intuitively, it is likely easier to find Pareto improvements for policies of the form “approve mergers of type X” as opposed to policies of the form “approve merger X”, since the willingness to pay can be thought of as ex-ante to the set of mergers that will be approved. Efficient surplus is well-suited to addressing comparisons where the key source of heterogeneity is income.

Appendix Table I  
Summary Statistics

Marginal Federal Tax Rate (1)	Number of Filers (2)	Percent of Filers (3)	Mean Ordinary Income (4)	Mean Family Income (5)
-35.0%	589,750	0.6%	79	11,622
-30.0%	1,413,594	1.4%	229	11,598
-25.0%	6,604	0.0%	16,004	34,459
-24.0%	1,366,424	1.4%	346	9,435
-19.0%	10,905	0.0%	15,088	30,532
2.4%	2,738,247	2.7%	170	6,320
7.4%	18,321	0.0%	19,049	35,938
10.0%	19,355,436	19.3%	(7,140)	18,754
15.0%	35,974,064	35.9%	32,510	52,689
17.4%	620	0.0%	58,739	65,642
17.7%	2,736,669	2.7%	1,120	11,585
22.7%	927	0.0%	12,936	16,493
25.0%	19,217,503	19.2%	74,974	97,153
26.0%	3,035,148	3.0%	6,416	25,402
28.0%	6,266,531	6.3%	193,243	234,677
31.0%	1,078,152	1.1%	15,822	32,212
31.1%	4,868,642	4.9%	4,622	28,992
33.0%	225,353	0.2%	243,716	234,349
35.0%	400,306	0.4%	1,977,424	1,154,254
36.1%	795,541	0.8%	16,523	38,355
Other	549	0.0%	15,934,759	54,534
<b>Total</b>	<b>100,099,286</b>	<b>100.0%</b>	<b>46,110</b>	<b>64,745</b>

*Notes.* This table presents summary statistics for the universe of income tax returns in 2012 for U.S. citizens aged 25-60. This sample is used to construct the elasticity of the density of the income distribution (i.e. alpha) for each marginal tax rate. The table presents the number of filers, mean ordinary income, and mean family income by each federal marginal tax rate. The mean federal tax rate is the effective marginal tax rate each filer would face on an additional dollar of income. This equals the tax on ordinary income for most filers, but includes additional tax rates generated by the earned income tax credit (EITC) for EITC filers and the alternative minimum tax (AMT) for filers subject to the AMT.

## E Sample Estimation Details

### E.1 Sample and Variables

To estimate the joint distribution of income and tax rates, I use the universe of de-identified 2012 tax returns taken from the 2012 IRS-SOI Databank maintained under the Statistics of Income Division at the IRS. I focus on primary filers aged 25-60 and their married spouses, if applicable.<sup>28</sup> Following Chetty et al. (2014), I restrict the sample to households with positive family income. Details of the sample and data construction are provided in Appendix E; Appendix Table I presents the summary statistics. The resulting sample has roughly 100 million filing units.

I define  $y$  to be the tax filer’s ordinary income in 2012.<sup>29</sup> This equals taxable income (f1040, line 43) minus income not subject to the ordinary income tax (long-term capital income (line 13) and qualified dividends (line 9b)). Ordinary income is primarily comprised of labor income, but subtracts deductions for things like the number of children and charitable donations.

To each tax return, I assign the marginal tax rate faced by the 2012 federal income tax schedule. The federal rate schedule on ordinary income provides the marginal tax rate for the filer as long as s/he did not have any additional tax credits, such as the earned income tax credit, and was not subject to the alternative minimum tax. If the individual was subject to the alternative minimum tax (AMT), I record their marginal tax rate at the 28% AMT level. If the filer received EITC, I add the marginal tax rate on the EITC schedule using information on the number of EITC-eligible children (reported in the tax return), filing status, and the size of the EITC benefit claimed. This provides a precise measure of the federal marginal tax rate faced by each filer on an additional dollar of ordinary income.

In addition to federal taxes, I account for state and local taxes. For state taxes, I assume a constant tax rate of 5% and account for the fact that state taxes are deductible from federal tax liability when calculating the total marginal tax rate. For Medicare and sales taxes, I follow Saez et al. (2012) and assume a 2.9% tax rate for Medicare and a 2.3% sales tax rate. Finally, some states provide additional EITC benefits. To account for this, I assume a 10% “top-up” EITC rate for EITC filers.<sup>30</sup> In the end, this generates a marginal tax rate,  $\tau(y)$ , faced by each filer on an additional dollar of income.

### E.2 Summary statistics and Estimation Approach

Appendix Table I presents the summary statistics of the sample used to construct the estimates of the shape of the income distribution conditional on the marginal income tax rate. Overall, there are roughly 100M filers aged 25-60 used in the analysis, with mean family incomes of roughly \$65M, and mean ordinary incomes of \$46M.

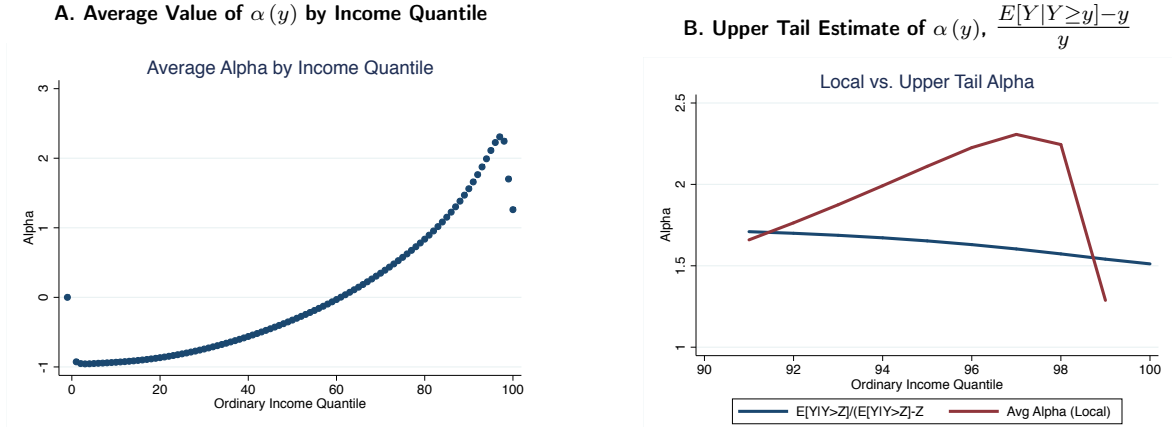
---

<sup>28</sup>I exclude individuals below age 25 because of the likelihood they still live at home and are part of another household. I exclude people above 60, the age at which many begin exiting the labor force and begin collecting unearned income such as social security income or savings withdrawals.

<sup>29</sup>Because ordinary income determines the federal tax, it is the notion of income that most closely aligns with the theory.

<sup>30</sup>Choosing alternative values for the state tax rates or the EITC rates do not significantly alter the results; as discussed below, the primary driver of the shape of the weights is the Pareto parameter combined with the assumption of a constant elasticity, not the shape of tax rates,  $\tau(y)$ .

## Appendix Figure 1: Estimation of Shape of Income Distribution



Notes: Panel A of this figure presents estimates of the average value of  $\alpha(y)$  by income quantile. Panel B presents estimates using an alternative method of estimating  $\frac{E[Y|Y \geq y] - y}{y}$  in each quantile.

To estimate the Pareto parameter of the income distribution, I proceed as follows. First, for computational simplicity, I define 1000 equally sized bins of ordinary income. I then collapse the data to generate counts of returns in each of these 1000 bins separately for returns facing different tax schedules,  $j$ . I generate these groups as the intersection of filing status, EITC status (marital status + number of qualified EITC dependents), and those subject to the alternative minimum tax rate.

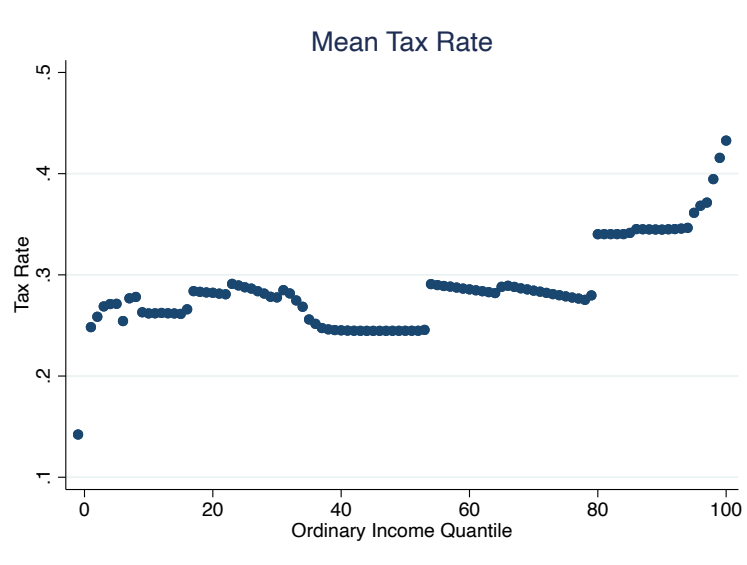
Given these groupings, I estimate the shape of the income distribution,  $\alpha$ , in a manner that allows it to vary with the marginal tax rate for a majority of the population. Let  $j$  index the set of tax schedules. For tax schedules with at least at least 500,000 observations with earnings between the 10th and 99th percentile of the income distribution, I estimate the elasticity of the income distribution separately for each filing characteristic, which I denote  $\alpha_j(y)$ .<sup>31</sup> To do so, I construct the log density of the income distribution measuring the number of households in each bin divided by the width of the bin. I then regress this on a fifth order polynomial of log income in the bin (where income is the mean income within the bin). The estimated slope at each bin generates an estimate of  $\alpha_j$  for each income bin in tax group  $j$ . I verify that the results are virtually identical when increasing or decreasing the number of bins or changing the number of polynomials in the regression.

For the remaining smaller tax groups (~25% of the sample) with fewer than 500,000 returns, I impose the assumption that the elasticity of the income distribution is the same across these less-populated tax schedules at a given level of income.<sup>32</sup> I then take advantage of the fact that the aggregate elasticity can be written as a weighted average of the elasticities of the income distribution

<sup>31</sup>I do not include the returns below the 10th quantile of the income distribution because of the large fraction of returns posting exactly \$3k in ordinary income, which introduces significant nonlinearities in some of these groups. Above the 99th percentile, I follow a strategy from Saez (2001) described below.

<sup>32</sup>This 500,000 threshold is chosen for computational simplicity on the remaining groups, but the results are similar to lowering it to 250,000.

Appendix Figure 2: Average Tax Rates by Income Quantile



Notes: This figure presents the average tax rate by income quantile. Each tax rate is the sum of the federal income tax rate, state taxes, Medicare, sales taxes, and EITC top-up, as discussed in Section E.1.

for each marginal tax rate,  $\alpha_j$ . So, I estimate the elasticity of the aggregate income distribution and then construct the implied elasticity for these smaller groups as the population weighted difference between the total elasticity and the elasticities of the larger tax groups. To estimate the elasticity of the aggregate income distribution, I regress the log density on a tenth order polynomial in log income for each bin (again, results are nearly identical if one includes additional polynomials) and compute the slope at each bin.

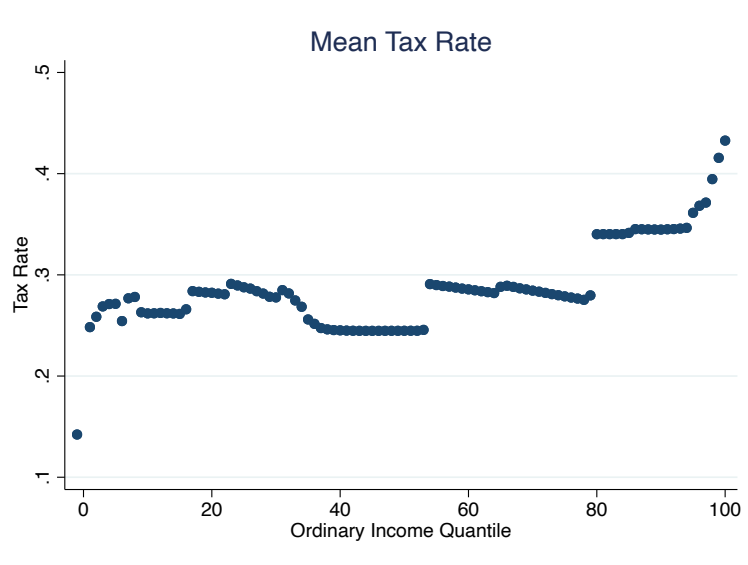
The advantage of this estimation approach is that it allows the elasticity of the income distribution to vary non-parametrically with the tax rates for  $\sim 75\%$  of the sample. This allows for correlation between the shape of the income distribution and the marginal tax rate, as is potentially required for accurate estimation of the substitution effect in the presence of multiple tax schedules.<sup>33</sup>

For individuals near the top of the income distribution, the local calculation of the elasticity of the income distribution becomes difficult and potentially biased because of endpoint effects. Intuitively, the binning of incomes into 1,000 bins ignores the fact that the U.S. income distribution has a fairly thick upper tail. Fortunately, it is well documented that the upper tail of the income distribution is Pareto, and hence has a constant elasticity so that  $\alpha(y) = \frac{E[Y|Y \geq y] - y}{y}$  (Saez (2001)). Hence, I also compute an “upper tail” value of  $\alpha$  given by  $\frac{E[Y|Y \geq y] - y}{y}$  for each income bin. Appendix Figure 1 (left panel) plots the average local estimate of  $\alpha$  (using the fifth order polynomial) across the income distribution and Appendix Figure 1 (right panel) plots both this estimate and the upper tail value of  $\alpha$ ,  $\frac{E[Y|Y \geq y] - y}{y}$ , for the upper decile of the income distribution.

<sup>33</sup>In practice, this degree of generality turns out not to matter in the estimation: one could arrive at a similar set of weights using the average Pareto parameter at each income level instead of estimating its heterogeneity across tax schedules.

For the upper regions of the income distribution, the value of  $\frac{E[Y|Y \geq y] - y}{y}$  converges to around 1.5, consistent with the findings of Diamond and Saez (2011) and Piketty and Saez (2013). Conversely, the local estimate of the elasticity of the income distribution arguably becomes downwardly biased in the upper region because the fifth order polynomial does not capture the size of the thick tail in the top-most income bucket. Hence, for incomes in this upper region with earnings above \$250,000, I assign the maximum value of these two estimates.

Appendix Figure 3: Incorporating Income Effects



*Notes:* This figure presents the efficient social welfare weights using both the baseline specification (solid blue line) and a modified specification that incorporates an income effect (dashed red line). To calculate the modified specification with the income effect, I assume a constant elasticity of labor supply with respect to income of -0.15, similar to the estimate in Cesarini et al. (2015).

## F Income Effects

The baseline specification assumes no income effects on labor supply. This section illustrates how income effects increase the marginal cost of taxation,  $g(y)$ , but do so similarly at all points of the income distribution (assuming a constant elasticity). To illustrate, Appendix Figure 3 presents the baseline specification for  $g(y)$  combined with an alternative specification that incorporates income effects. For simplicity, I approximate the income effect as  $\zeta(y) \frac{\tau(y)}{1-\tau(y)}$  where  $\tau(y)$  is the average marginal tax rate for those in each quantile of ordinary income. For  $\zeta(y)$ , I take an estimate of 0.15 from Cesarini et al. (2015) who study the impact of winning the lottery in Sweden on labor supply.

As shown in Appendix Figure 3, incorporating income effects raises the marginal cost of taxation at all income levels. But, in contrast to the substitution effect and the compensated elasticity, it does not differentially affect the marginal cost of taxation at different income levels. In this sense, the broad set of conclusions that one should apply greater weight to surplus to the poor than to the rich remains true if one incorporates income effects into the analysis.

## G Testing Weak Separability: Relation to Kaplow (2000, 2006, 2008, etc.) and Hylland and Zeckhauser (1979)

There is an long debate about whether or not one should weight the willingness to pay for publicly provided goods for the poor differently than the rich. Most influentially, Hylland and Zeckhauser (1979) followed by Kaplow (1996, 2004, 2008) provide a weak separability assumption on the utility function that, if satisfied, implies that additional spending on the publicly provided good increases utility if and only if the sum of individuals' willingness to pay exceeds the mechanical cost of the publicly provided good. This Appendix shows how this theoretical result is nested in the model of Section 9, and thus the welfare framework provides a test of the weak separability assumptions employed in this literature.

Consider a policy of spending \$1 per capita on a publicly provided good,  $G$ . This will have a net cost to the government of \$ $c$  that may differ from \$1 because of any fiscal externalities from behavioral responses to the provision of the public good,  $FE^G = c - 1$ . Assume that individuals of income level  $y$  are willing to pay  $s(y)$  for this additional expenditure so that the average willingness to pay is  $E[s(y)]$ . Individuals are thus willing to pay the mechanical cost of the expenditure if and only if  $E[s(y)] \geq 1$ . In contrast, equation (7) (generalized to the case of willingness to pay that varies with  $y$ ) suggests that additional spending on  $G$  is efficient if and only if  $E[g(y)s(y)] \geq c$ . How are these different?

It turns out that the weak separability assumption in Hylland and Zeckhauser (1979) and Kaplow (1996, 2004, 2008) implies that the behavioral response to \$1 of a tax cut to those earning near \$ $y$  should be the same as the behavioral response to a policy that provides \$1's worth of additional  $G$  to those earning near \$ $y$ . Weak separability imposes that the behavioral response to a tax cut scaled by the willingness to pay for the policy equals the behavioral response to the policy. Hence,

$$FE^G = E[s(y)FE(y)] \tag{14}$$

where  $s(y)$  is the willingness to pay of individuals with income  $y$  for the additional spending on  $G$  and  $FE(y)$  is the fiscal externality associated with the tax cut. Hence, the total cost of the additional spending on  $G$  is equal to

$$\begin{aligned} c &= 1 + FE^G \\ &= 1 + E[s(y)FE(y)] \end{aligned}$$

Hence, testing whether  $E[s(y)g(y)] \geq c$  is equivalent to testing whether

$$E[s(y)g(y)] \geq c \iff E[s(y)(1 + FE(y))] \geq 1 + E[s(y)FE(y)] \iff E[s(y)] \geq 1$$

So, the test for efficiency reduces to

$$E[s(y)] \geq 1$$

which asks whether the aggregate willingness to pay exceeds the mechanical cost of the policy (that

does not include the fiscal externalities). In this sense, if equation (14) holds for the policy change in question, one need not know either the efficient welfare weights,  $g(y)$ , or the fiscal externalities induced by the policy change,  $FE^G$ . One can simply compare unweighted aggregate willingness to pay to the mechanical cost of the policy. But more generally, testing whether  $E[g(y)s(y)] \geq c$  provides a general method for asking whether the policy is efficient, even if weak separability does not hold.