



# Optimal redistribution and monitoring of labor supply<sup>☆</sup>



Floris T. Zoutman<sup>a</sup>, Bas Jacobs<sup>b, c, d, \*</sup>

<sup>a</sup>NHH Norwegian School of Economics, Department of Business and Management Science, Norway

<sup>b</sup>Erasmus University Rotterdam, The Netherlands

<sup>c</sup>Tinbergen Institute, The Netherlands

<sup>d</sup>CESifo, Germany

## ARTICLE INFO

### Article history:

Received 17 September 2014

Received in revised form 2 October 2015

Accepted 18 December 2015

Available online 20 January 2016

### JEL classification:

H21

H26

H24

H31

### Keywords:

Optimal non-linear taxation

Monitoring

Costly verification of ability/labor supply

Optimal income redistribution

## ABSTRACT

This paper extends the Mirrlees (1971) model of optimal non-linear income taxation with a monitoring technology that allows the government to verify labor supply at a positive, but non-infinite cost. We analyze the joint determination of the non-linear monitoring and tax schedules, and the conditions under which these can be implemented. Monitoring of labor supply reduces the distortions created by income taxation and raises optimal marginal tax rates, possibly above 100%. The optimal intensity of monitoring increases with the marginal tax rate and the labor-supply elasticity. Our simulations demonstrate that monitoring strongly alleviates the trade-off between equity and efficiency. Welfare gains of monitoring are around 2.8% of total output. The optimal intensity of monitoring follows a U-shaped pattern, similar to that of optimal marginal tax rates. Our paper can explain why large welfare states optimally rely on work-dependent tax credits, active labor-market policies, benefit sanctions and work bonuses in welfare programs.

© 2016 Elsevier B.V. All rights reserved.

“Informational frictions are a specification of a particular type of technology. For example, when we say “effort is hidden”, we are really saying that it is infinitely costly for society to monitor effort. The desired approach would be to devise optimal tax systems for different specifications of the costs of monitoring different activities and/or individual attributes. To be able to implement

this approach, we need to... extend our modes of technical analysis to allow for costs of monitoring other than zero or infinity.”

Kocherlakota (2006 pp. 295–296)

## 1. Introduction

Redistribution of income is one of the most important tasks of modern welfare states. However, redistribution is expensive as it distorts the incentives to supply labor. As a result, there is a trade-off between equity and efficiency. On a fundamental level, Mirrlees (1971) demonstrates that the trade-off between equity and efficiency originates from an information problem. Earning ability and labor hours are private information. Hence, the government cannot condition redistributive taxes and transfers on earning ability, since the government cannot distinguish individuals that are unable to work from individuals that are unwilling to work. Therefore, redistribution from high-income to low-income earners inevitably distorts the incentives to supply labor hours.

In practice, labor supply is not completely non-verifiable, as assumed by Mirrlees (1971). Indeed, some welfare states do condition the tax burden on some measure of hours worked. For example,

<sup>☆</sup> The authors were awarded the IIPF Young Economist Award 2013 for economists under age 40 for this paper. The authors would like to thank Luca Micheletto, Jean-Marie Lozachmeur, Aart Gerritsen, Katherine Cuff, Dirk Schindler, Robin Boadway and two anonymous referees for very useful suggestions and comments on an earlier version of this paper. Furthermore, this paper benefited from comments and suggestions made by participants at the 66th and 69th IIPF Conferences in Uppsala and Taormina; the CPEG Conference, Quebec; and seminar participants at the Erasmus School of Economics. All remaining errors are our own. The Matlab programs used for the computations in this paper are available from the authors on request. Zoutman gratefully acknowledges financial support from the Netherlands Organisation for Scientific Research (NWO) under Open Competition grant 400-09-383.

\* Corresponding author at: Erasmus School of Economics, Erasmus University Rotterdam, PO Box 1738, 3000DR, Rotterdam, The Netherlands. Tel.: +31 10 4081481/1491.

E-mail addresses: [floris.zoutman@nhh.no](mailto:floris.zoutman@nhh.no) (F.T. Zoutman), [bjacobs@ese.eur.nl](mailto:bjacobs@ese.eur.nl) (B. Jacobs).

URL: <http://people.few.eur.nl/bjacobs>.

low-income individuals in the UK receive a tax credit if they work more than 30 hours. This policy can be implemented only if the government is able to verify hours worked. Similar restrictions apply to in-work tax credits in Ireland and New Zealand (see also OECD, 2011). Clearly, the assumption that hours worked and earning ability are not verifiable is a too strong assumption. In the real world, the government does verify hours worked of individuals to some extent, albeit at a cost. Consequently, the government can – to some extent – separate shirking high-ability individuals from hard-working low-ability individuals.

This paper extends Mirrlees (1971) by letting the government operate a monitoring technology. The monitoring technology allows the government to verify labor hours of an individual at a positive, but finite cost. The government optimally sets the monitoring schedule as a function of gross income. That is, the probability that an individual is monitored depends (possibly non-linearly) on his/her gross labor earnings. If an individual is monitored, the government perfectly verifies his/her labor supply and can deduce the individual's ability. By monitoring hours worked the government can thus provide incentives to individuals to change their labor supply in a direction that the government desires. In our model, individuals receive a work bonus when they work a minimum number of hours. This hours requirement is optimally determined by the government, and corresponds to the type of work bonuses observed in the UK, Ireland and New Zealand. Alternatively, we can formulate our model such that monitored individuals receive a penalty when their hours worked fall short of an optimally chosen minimum reference level, which can be thought of as work requirements or conditional welfare benefits that are observed in most advanced welfare states.

Each individual is aware of the monitoring and tax schedules before making labor-supply decisions. Hence, individuals can alter their monitoring probability by adjusting their hours worked. The total wedge on labor supply consists of the explicit income tax rate and an implicit subsidy on labor supply due to monitoring. Monitoring of hours worked acts as an implicit subsidy on labor supply for two reasons. First, the expected bonus increases (penalty decreases) in the difference between hours worked and the reference level of hours worked. Second, the monitoring intensity may decrease with gross earnings, depending on the shape of the monitoring schedule. For a given tax rate, monitoring can thus reduce the distortions of the income tax on labor supply, thereby increasing equity, efficiency or both.

In our model, the government maximizes social welfare by optimally setting the non-linear monitoring intensity and the minimum hours worked, alongside the optimal non-linear income tax.<sup>1</sup> We solve for the optimal non-linear tax and monitoring schedules by decentralizing the optimal, incentive-compatible direct mechanism that induces truthful revelation of ability types. We do not deviate from Mirrlees (1971) in that individuals always truthfully report earnings. We abstract from tax avoidance or evasion.<sup>2</sup> The schedule of optimal non-linear labor wedges is affected in two important ways in comparison to Mirrlees (1971). First, the monitoring intensity reduces the efficiency costs of the labor wedge, and allows for higher marginal tax rates. Second, a decrease in labor supply directly

increases the penalties (or decreases work bonuses). Monitoring generates inequality between monitored and non-monitored individuals at each ability level. Therefore, higher marginal taxes result in a distributional loss due to monitoring activities. The net effect of monitoring on the optimal wedge is thus theoretically ambiguous. In Mirrlees (1971), tax rates at, or above, 100% can never be optimal. In contrast, marginal tax rates could optimally be larger than 100% in our model due to optimal monitoring. Even if the marginal income tax rate is above 100%, individuals still supply labor as long as the total wedge on labor remains below 100%. Monitoring could explain why effective marginal tax rates of close to, or even higher than, 100% are observed in real-world tax-benefit systems in the phase-out range of means-tested benefits. See Immervoll (2004), Spadaro (2005), Brewer et al. (2010) and OECD (2011) for examples in OECD countries.

The non-linear monitoring schedule is set so as to equate the marginal cost of monitoring to the marginal efficiency gain associated with monitoring at each gross income level. The efficiency gain of monitoring is increasing in the distortion created by the wedge on labor. Therefore, the optimal monitoring intensity increases with both the total labor wedge and the labor-supply elasticity. The optimal minimum-hours requirement trades off the benefits of stronger work incentives and the costs of larger within-ability group inequality. A higher minimum-hours requirement raises the incentive for individuals to work more hours, which reduces income-tax distortions. However, a minimum-hours requirement also raises inequality between monitored and unmonitored individuals due to an increase in the penalties for untruthfully reporting hours worked.

Unfortunately, there is no closed-form solution for the optimal tax and monitoring schedules. Therefore, we resort to numerical simulations based on a realistic calibration of the model to US data. Our simulations demonstrate that the optimal tax schedule follows a U-shape, which closely resembles the simulations of Saez (2001). Moreover, the monitoring schedule also follows a U-shape. This confirms that the monitoring intensity should indeed be large when tax distortions on labor supply are large. The optimal minimum-hours requirement is typically around 95% of total labor hours, which corresponds to 38 hours based on a working week of 40 hours. The simulations demonstrate that the marginal tax rates with monitoring are generally higher than without monitoring. Hence, monitoring always results in more redistribution of income from high- to low-ability individuals, despite larger within-ability group inequality that results from monitoring and penalizing individuals.

Strikingly, our simulations demonstrate that the optimal tax rate at the bottom end of the income scale is substantially above 100%. This implies that the implicit subsidy on labor supply due to monitoring is very effective in reducing the total tax wedge on labor supply. Indeed, the optimal monitoring probability is close to one at low-income levels, but it drops substantially towards middle-income levels. There is a slight increase in the monitoring probability towards the top, since tax rates increase. We conclude from our simulations that monitoring is most important at the bottom of the income distribution. Strongly redistributive governments should therefore optimally employ a high monitoring intensity at the low end of the income scale, for example, via job-search requirements, benefit sanctions, work bonuses, and active labor-market programs. Moreover, our findings suggest that work-dependent tax credits for low-income earners, like those in the UK, Ireland and New Zealand, are indeed part of an optimal redistributive tax policy.

The welfare gains of monitoring are shown to be large. Compared to the optimal non-linear tax schedule without monitoring, monitoring increases total output by 1.3% in our baseline simulation. Moreover, the transfer increases by about 40%. The monetized welfare gain of monitoring is about 2.8% of total output. The optimal monitoring probability does not exceed 25% anywhere, except at the lower end of the income distribution. In our baseline simulations,

<sup>1</sup> In our model, first-best can generally not be obtained, because the penalty function is exogenous. If the government would be able to optimize the penalty function a trivial first-best outcome would result by either raising the penalty to infinity or adjusting the penalty function such that the implicit subsidy on work exactly off-sets the explicit tax on work.

<sup>2</sup> We realize that the assumption of truthful reporting of earnings is not always realistic due to, for example, tax evasion and avoidance. This issue has been discussed in, amongst others, Cremer and Gahvari (1996), Schroyen (1997) and Chander and Wilde (1998). In most developed countries, however, firms are required to report gross labor earnings directly to the tax authorities, which prevents underreporting of earnings for a very large fraction of the population (see e.g., Kleven et al., 2011).

the cost of monitoring equals 0.8% of average labor earnings. Extensive sensitivity analyses demonstrate that our results are robust to parameter changes in the monitoring technology, on which little empirical evidence exists.

The setup of the paper is the following. The next section gives a brief overview of the related literature. The third section introduces the model and derives the conditions for first- and second-order incentive compatibility. The fourth section derives the optimality conditions for monitoring and redistribution. The fifth section presents the simulations. Finally, the sixth section concludes.

## 2. Review of the literature

Our model builds upon two strands in the mechanism-design literature. *Mirrlees (1971), Diamond (1998) and Saez (2001)* develop the theory of the optimal non-linear income tax under the assumption that both hours worked and ability are completely private information, implicitly assuming that verification of either hours worked or ability is prohibitively costly. On the other hand, the literature on costly state verification develops principal-agent models where the outcome of a project is a function of both the state of the world and the action of the agent (see, e.g., *Mirrlees, 1999, 1976; Holmstrom, 1979; Townsend, 1979*). The outcome is observed, but the action and the state of the world can only be verified through costly monitoring. Monitoring can improve the ex-ante utility of both the principal and the agent. We apply the theory of costly state verification to the *Mirrlees (1971)* model and show that monitoring of labor supply can increase social welfare significantly.

In a related paper, *Armenter and Mertens (2013)* study the effect of optimal monitoring of ability types on the optimal tax schedule. They analyze a dynamic model of optimal taxation where the government can use a monitoring technology to establish the ability of an agent. In their model, the monitoring intensity is exogenous, while penalties are endogenous. In equilibrium, individuals do not misreport their ability, and are therefore never penalized. Indeed, the economy is shown to converge to first best in an infinite-horizon setting. We, instead, analyze the case where monitoring is endogenous and penalties are exogenously given. Because penalties are exogenously given, individuals may misreport their ability type in equilibrium. Consequently, our model does not converge to a first-best outcome. An advantage of allowing for an endogenous monitoring intensity is that we do not need to worry about a tax-riot equilibrium in which all individuals misreport their type when they expect other individuals to do the same (*Bassetto and Phelan, 2008*).

The effect of monitoring has also been studied in the literature on tax evasion and the literature on unemployment insurance. The literature on tax evasion (see, e.g., *Allingham and Sandmo, 1972; Sandmo, 1981; Mookherjee and Png, 1989; Slemrod, 1994; Cremer and Gahvari, 1994, 1996; Chander and Wilde, 1998; Slemrod and Kopczuk, 2002*) extends the *Mirrlees (1971)* framework by allowing individuals to underreport their earned income to the tax authorities.<sup>3</sup> Compared to the standard *Mirrlees (1971)* model, income taxation is more distortionary, because it not only reduces labor supply, but it also increases tax evasion. However, the government can monitor individuals by auditing their tax returns and fine them when they evade taxes. As a result, the equity-efficiency trade-off improves.

Our paper is most closely related to *Cremer and Gahvari (1996)*. They develop a two-type economy with non-linear taxation and monitoring. As in *Mirrlees (1971)*, the government cannot verify

ability and labor effort. In addition, the government can only verify labor earnings through monitoring. Hence, individuals make two choices: how much labor to supply and how much of their income to report to the tax authorities (tax evasion). *Cremer and Gahvari (1996)* show that the government optimally levies a positive marginal tax rate on the lowest type, provided monitoring is not prohibitively expensive, and monitoring and penalties are sufficient to deter individuals from misreporting their income. Marginal tax rates for the highest type are always zero. All individuals reporting an income below a given threshold should be monitored with positive probability. Our model, in contrast, assumes that labor earnings are perfectly verifiable and the individual makes only one choice: the amount of income to earn (labor to supply), or, equivalently, which type to report to the government. Since earnings are perfectly verifiable, individuals cannot evade taxes, as in *Cremer and Gahvari (1996)*. The monitoring instrument is aimed at measuring hours worked instead of tax evasion.<sup>4</sup> We contribute in a number of ways to the analysis in *Cremer and Gahvari (1996)*. First, in contrast to their two-type model, our model allows for a continuous distribution of ability types, which allows us to derive the complete non-linear tax and monitoring schedules for the entire population. By considering optimal non-linear tax and monitoring under a continuum of skill types we derive an elasticity-based formula for the optimal non-linear tax and monitoring schedule in the spirit of *Diamond (1998) and Saez (2001)*. Moreover, we provide the conditions under which these schedules are implementable. Second, we determine the shape of non-linear tax and monitoring schedules over the entire income distribution through simulations. Third, in *Cremer and Gahvari (1996)* prohibitively expensive monitoring would lead to a laissez-faire outcome, since individuals would report no income at all when there would not be any monitoring of earnings. In contrast, our model nests the standard *Mirrlees (1971)* model as a special case when monitoring is prohibitively expensive.

In the literature on unemployment insurance, *Ljungqvist and Sargent (1995a,b)* study the effect of monitoring on equilibrium employment in welfare states.<sup>5</sup> In their model, unemployed workers may receive a job offer each period. In the absence of monitoring, the benefits induce workers to decline an inefficiently large number of job offers. Monitoring can raise efficiency by punishing those workers that decline job offers. Simulations using Swedish data demonstrate that welfare states with large benefits and progressive taxation can have low equilibrium unemployment rates, provided that monitoring probabilities and sanctions are sufficiently large. In a model of optimal income redistribution with search, *Boadway and Cuff (1999)* determine the welfare-maximizing monitoring probability and demonstrate that it is increasing in the level of the benefits. *Boone and Van Ours (2006) and Boone et al. (2007)* develop a search model where the government can actively monitor and sanction job-search effort. They show that monitoring and sanctioning may be more effective in reducing unemployment than cutting the replacement rate. In addition, they show that monitoring may be effective, even when the duration of unemployment benefits is limited. This literature has focused on monitoring the search effort of unemployed workers. We contribute to this literature by studying the effect of monitoring on employed workers.

Finally, we contribute to the literature on optimal non-linear tax simulations (see, for example, *Mirrlees, 1971; Tuomala, 1984; Saez, 2001; Brewer et al., 2010; Zoutman et al., 2015*). We show

<sup>4</sup> An alternative interpretation would be that individuals exogenously supply labor, but can use a costly evasion technology.

<sup>5</sup> A large literature exists on optimal unemployment insurance. See *Fredriksson and Holmlund (2006)* for a survey of this literature. However, this literature typically does not consider monitoring of search effort.

<sup>3</sup> A comprehensive survey of the literature can be found in *Slemrod and Yitzhaki (2002)*.

that monitoring can lead to significant improvements in both equity and efficiency.<sup>6</sup>

### 3. Model

#### 3.1. Households

The setup of our model closely follows [Mirrlees \(1971\)](#). Individuals are heterogeneous in their earning ability  $n$ , which denotes the productivity per hour worked. Ability is distributed according to cumulative distribution function  $F(n)$  with support  $[\underline{n}, \bar{n}]$ , where  $\bar{n}$  could be infinite. The density function is denoted by  $f(n)$ . Workers are perfect substitutes in production and the wage rate per efficiency unit of labor is constant and normalized to one.  $n$  corresponds to the number of efficiency units of labor of each worker. Gross labor income of an individual is the product of his/her ability and his/her labor hours  $z_n = nl_n$ . Individuals derive utility from consumption  $c_n$  and disutility from hours worked  $l_n$ .

We introduce the model using a formulation where individuals may receive a work bonus when they work a certain number of hours. Then, we demonstrate that a tax implementation where individuals receive a penalty if they fail to meet this level of working hours is equivalent. The critical part of our analysis is therefore the monitoring of labor supply, not the particular tax implementation through bonuses or penalties. To fix ideas, we assume that the tax schedule consists of two parts. First, individuals pay income taxes  $\hat{T}(z_n)$  based on their earned income  $z_n$ . Second, individuals can apply for a working tax credit,  $\mathcal{T}$ , if their hours worked exceed a work requirement,  $l^*$ , which the government optimally chooses. The work requirement is the same for all individuals.<sup>7</sup> This tax schedule corresponds closely to what we observe in the UK, New Zealand and Ireland – see the remarks in the Introduction. Consequently, total tax payments for individuals applying for the tax credit are given by  $T(z_n) \equiv \hat{T}(z_n) - \mathcal{T}$ . Similarly, tax payments of the individuals who do not apply for the tax credit are simply  $\hat{T}(z_n)$ .

We make the technical assumption that all individuals apply for the tax credit. This assumption is nearly without loss of generality as we can always ensure that all individuals apply for the credit by simultaneously adjusting tax payments without the tax credit  $\hat{T}(z_n)$ , and the tax credit  $\mathcal{T}$  by similar amounts.<sup>8</sup> Below we demonstrate that such a policy is in the best interest of the government, since monitoring effectively alleviates the equity-efficiency trade-off and moves the optimal second-best allocation closer to the first-best allocation.

Individuals can misreport their hours worked to the tax authorities, and can claim the tax credit while not satisfying the minimum-hours requirement. The government, however, can operate a monitoring technology to verify actual hours worked of an individual applying for the tax credit.  $\pi(z_n)$  denotes the probability that an individual with earnings  $z_n$  is monitored by the government.  $\pi(z_n)$  is also

referred to as the monitoring intensity. We assume that the government receives a perfect signal of the individual's labor supply  $l_n$  if the individual is monitored.

Monitored individuals receive a penalty if they are found to misrepresent their hours worked. The size of the penalty depends on the difference between required working hours  $l^*$  and actual hours  $l_n$  worked:

$$P \equiv \begin{cases} P(l^* - l_n) & \text{if } l^* > l_n \\ 0 & \text{if } l^* \leq l_n \end{cases}, \quad P(\cdot), P'(\cdot), P''(\cdot) \geq 0. \quad (1)$$

We will refer to  $P(\cdot)$  as the penalty function. We restrict penalties to be non-negative. The penalty function  $P(\cdot)$  is exogenously given. We assume that the penalty function is differentiable once over its entire domain, and differentiable twice everywhere, except possibly at zero. Penalties and marginal penalties are both positive and marginal penalties are increasing when individuals are found to supply less labor than the hours requirement ( $P(\cdot), P'(\cdot), P''(\cdot) \geq 0$ ). Consequently, penalties decrease in hours worked. The government optimally determines the reference level of working hours  $l^*$ . By raising the the minimum-hours requirement  $l^*$ , the government effectively raises the penalty when individuals are monitored. For a given gross income level  $z_n$  penalties thus increase in ability, since higher ability individuals need to supply less hours in order to attain a given gross income level. Finally, we assume that the government does not penalize individuals that applied for the tax credit and supplied the required minimum amount of hours. [Fig. 1](#) displays an example of a penalty function. As can be seen, the penalty decreases quadratically in labor supply up to  $l_n = l^*$ , after which it remains constant at 0. Such a penalty function will be used in the simulations later.

We believe that constraining the penalty function  $P(\cdot)$  is realistic for two reasons.<sup>9</sup> First, the legal system practically imposes limitations on the government's ability to use infinite penalties. Second, we assume perfect monitoring as the labor supply of each individual is verified with perfect certainty. If we would more realistically assume that monitoring is imperfect, hard-working individuals could inadvertently be classified as shirking individuals. Then, we would be able to endogenize both the penalty function and the monitoring function, and infinite penalties would never be socially optimal (see e.g., [Stern, 1982](#); [Diamond and Sheshinski, 1995](#); [Jacquet, 2014](#)). We leave this extension for future research as it would severely complicate our analysis without affecting the main result: monitoring alleviates the equity-efficiency trade-off.

Rather than using work bonuses, the government could, equivalently, use a negative tax credit  $\mathcal{T}$ , i.e., a work penalty, for all individuals not supplying the minimum level of labor. Individuals would then be required to report to the government whether they met the work requirement to avoid having to pay  $\mathcal{T}$ , and the government then needs to verify whether these work reports are indeed truthful. The total tax schedule  $T(z_n)$  would remain the same. Hence, the particular tax implementation with either work bonuses or non-work penalties is immaterial to our main findings. In the remainder of this paper we focus on determining the optimal total tax schedule  $T(z_n)$  including the tax credit.

The consumption of an individual who is not monitored is given by  $c_n^U \equiv z_n - T(z_n)$ . The consumption of a monitored (and penalized if hours worked are less than required for the credit) individual is

<sup>6</sup> Our paper is also relevant for the literature on minimum wages and optimal taxation (see, for example, [Boadway and Cuff, 2001](#); [Lee and Saez, 2012](#); [Gerritsen and Jacobs, 2015](#)). In [Mirrlees \(1971\)](#), wage rates per hour worked are assumed to be non-verifiable, whereas wage rates need to be verifiable to implement and enforce the minimum wage. This informational inconsistency can be avoided by costly monitoring of wage rates so as to enforce the minimum wage. Combining optimal non-linear taxes with optimal minimum wages and optimal monitoring is an interesting route for further research.

<sup>7</sup> [Zoutman and Jacobs \(2014\)](#) show that it is straightforward to extend the analysis with a non-linear work requirement that is dependent on ability. However, no additional insights are gained and the analysis becomes more complex as incentive-compatibility constraints will be affected by the work-requirement schedule as well.

<sup>8</sup> To see this, suppose that the tax credit  $\mathcal{T}$  and the tax schedule  $\hat{T}(z_n)$  are given. Next, add an arbitrarily large number to both. The incentive to apply for the tax credit then increases, but it does not affect total tax payments  $T(z_n)$ . Consequently, there always exists a level of the tax credit  $\mathcal{T}$  beyond which everyone applies for it.

<sup>9</sup> A more thorough discussion on these issues can be found in [Schroyen \(1997\)](#), [Mirrlees \(1997\)](#), and [Mirrlees \(1999\)](#).



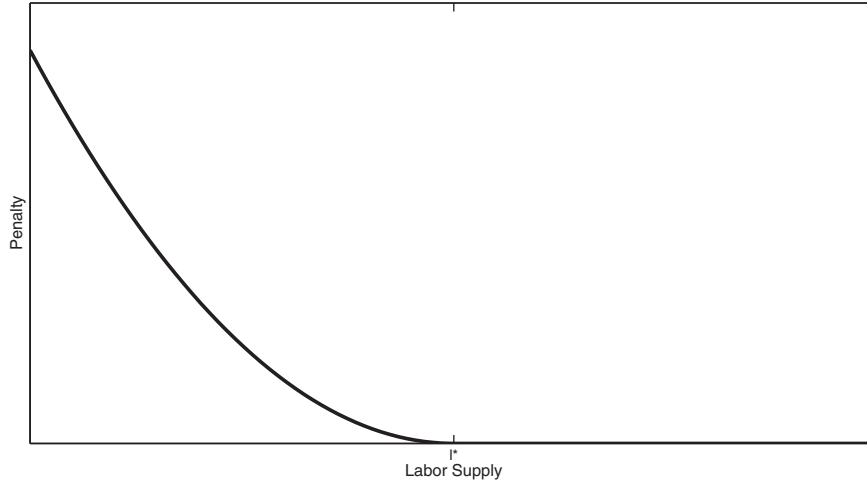


Fig. 1. Example of a penalty function.

given by  $c_n^p \equiv z_n - T(z_n) - P(l^* - l_n)$ .<sup>10</sup> Individuals are assumed to maximize expected utility subject to their budget constraints in monitored and unmonitored states. We follow Diamond (1998) by assuming that all individuals have an identical quasi-linear expected utility function:

$$u(z_n, n) \equiv \pi(z_n)c_n^p + (1 - \pi(z_n))c_n^u - v(l_n), \quad v'(\cdot) > 0, \quad v''(\cdot) < 0, \\ = z_n - T(z_n) - \pi(z_n)P(l^* - z_n/n) - v(z_n/n), \quad \forall n, \quad (2)$$

where we substituted the household budget constraint and  $l_n = z_n/n$  in the second line. An important analytical advantage of this quasi-linear-in-consumption utility function is that individuals are risk neutral.<sup>11</sup> The first term in the first line represents the non-monitoring probability times the consumption of an individual that is not monitored. The second term in the first line is the monitoring probability times the consumption of an individual that is monitored. The last term in the first line is the disutility of labor supply.

Individuals choose the optimal amount of gross income based on their productivity  $n$ , the tax function  $T(\cdot)$ , the monitoring function  $\pi(\cdot)$ , and the penalty function  $P(\cdot)$ . An income level  $z_n$  is incentive compatible if it maximizes  $u(z_n, n)$ . The first-order condition for optimal labor supply is given by:

$$v'(z_n/n) = (1 - T'(z_n) - \pi'(z_n)P(l^* - z_n/n))n + \pi(z_n)P'(l^* - z_n/n), \quad \forall n. \quad (3)$$

<sup>10</sup> Monitored and unmonitored individuals thus face the same tax schedule, but the tax function features a discontinuity, since the intercepts of the tax function for monitored and unmonitored individuals are different. This is similar to Jacquet et al. (2013) who analyze optimal taxes for employed and non-employed workers where intercept of the tax function for the non-employed is different from that for the employed.

<sup>11</sup> We could allow for risk-aversion in the utility function. In that case, we are only able to solve for the optimal non-linear tax and monitoring schedules if the social welfare function is utilitarian. Intuitively, the problem becomes analytically untractable if the government has a different degree of risk-aversion – which is implied by a non-utilitarian social welfare function – than households have. Without risk aversion, this problem is always absent and we can allow for any degree of inequality aversion in the social welfare function.

On the right-hand side, we see that policy drives a wedge between the private and social benefits of labor supply. The total labor wedge  $\mathcal{W}_n$  is given by:

$$\mathcal{W}_n \equiv \frac{n - v'(z_n/n)}{n} = \underbrace{T'(z_n)}_{\text{explicit tax}} + \underbrace{\pi'(z_n)P(l^* - z_n/n) - \frac{\pi(z_n)}{n}P'(l^* - z_n/n)}_{\text{implicit tax}}, \quad \forall n. \quad (4)$$

In a laissez-faire equilibrium, the right-hand side of Eq. (3) equals  $n$  and the total labor wedge  $\mathcal{W}_n$  is zero. The total labor wedge consists of the explicit marginal tax on labor ( $T'$ ) and the implicit marginal tax (subsidy) on labor due to monitoring ( $\pi'P - \pi P'/n$ ). If  $T' + \pi'P - \pi P'/n > 0$ , the redistributive tax and monitoring policy reduces optimal labor supply below the laissez-faire level, and vice versa if it is smaller than zero. The wedge is naturally increasing in the explicit marginal rate  $T'$ . Furthermore, it increases in the marginal monitoring probability,  $\pi'$ , if penalties are positive, i.e.,  $P > 0$ .  $\pi'$  gives the marginal increase in the monitoring probability as a function of gross earnings. If the monitoring probability increases (decreases) with income, this reduces (increases) the incentive to supply labor, because a higher labor income increases (decreases) the probability of receiving a penalty. Therefore, an increase in the marginal monitoring probability decreases the incentive to supply labor.

Proposition 1 shows that without loss of generality we can assume that expected consumption,  $C(z_n) \equiv z_n - T(z_n) - \pi(z_n)P(l^* - z_n/n)$ , is non-decreasing in earnings  $z_n$ . Consequently, the total labor wedge  $\mathcal{W}_n$  can never be larger than one, i.e., larger than 100%.

**Proposition 1.** All implementable continuous allocations can be implemented through a continuous non-decreasing expected consumption function  $C(z_n)$ ,  $\forall n$ . If  $C(z_n)$  is continuous and differentiable, the wedge  $\mathcal{W}_n$  can never exceed 1.

**Proof.** The proof directly follows Mirrlees (1971). Let  $\tilde{C}(z)$  be any continuous expected consumption function. The individual maximization problem is given by:

$$z_n = \arg \max_{z_n} \tilde{C}(z_n) - v(z_n/n), \quad \forall n. \quad (5)$$

Now, consider the function  $C(z_n) = \max_{\bar{z}_n \leq z_n} \tilde{C}(\bar{z}_n)$ . Clearly,  $C(\cdot)$  is non-decreasing and continuous, because  $\tilde{C}(\cdot)$  is continuous. Next, consider the maximization problem:

$$\max_{z_n} C(z_n) - v(z_n/n) = \max_{z_n} \left[ \max_{\bar{z}_n \leq z_n} \tilde{C}(\bar{z}_n) \right] - v(z_n/n), \quad \forall n. \quad (6)$$

Assume  $z_n$  is the solution to Problem (5). The solution to this second maximization problem must also be  $z_n$ . To see this, evaluate  $C(\cdot)$  at  $z_n$ :  $C(z_n) = \max_{\bar{z}_n \leq z_n} \tilde{C}(\bar{z}_n)$ . Either  $C(z_n) = \tilde{C}(z_n)$  or  $C(z_n) = \tilde{C}(\bar{z}_n)$  with  $\bar{z}_n < z_n$ . In the first case, maximization Problems (6) and (5) are equivalent, and hence, they must have the same solution. In the second case, because  $v(\cdot)$  is strictly increasing in  $z_n$ ,  $\bar{z}_n$  must give a higher value to the objective function in Eq. (5) than does  $z_n$ . Hence, we arrive at a contradiction, because  $z_n$  could not have been the solution to Problem (5) in the first place. Therefore, without loss of generality we can focus on non-decreasing functions  $C(\cdot)$ . Now, suppose  $C(\cdot)$  is differentiable and consider its derivative:

$$C'(z_n) = 1 - T'(z_n) - \pi'(z_n)P(l^* - z_n/n) + \frac{\pi(z_n)}{n} P'(l^* - z_n/n) = 1 - \mathcal{W}_n, \quad \forall n. \quad (7)$$

$C(z_n)$  is non-decreasing if its derivative is greater than or equal to zero:  $C'(z_n) \geq 0 \iff \mathcal{W}_n \leq 1$ . ■

**Proposition 1** has an intuitive interpretation. Suppose an individual has a budget constraint such that expected consumption is decreasing in gross income over some interval. Then, this individual will never choose gross income in this interval, because he/she can work less and consume more, both yielding higher utility. Consequently, the government can never increase social welfare by setting the wedge  $\mathcal{W}_n$  above 1. The explicit marginal tax rate  $T'(z_n)$ , however, could be above 1, provided that monitoring implies a sufficiently large implicit marginal subsidy on work, i.e.  $\pi'P - \pi P'/n < 0$ , such that the overall wedge remains below 1. This is the case if the expected penalty increases sufficiently fast in the difference between expected and required labor supply such that  $\pi P'/n > \pi'P$ . Therefore, monitoring can improve the incentives to supply labor.

### 3.2. Government

The government designs an optimal income tax system, monitoring schedule and work requirement so as to maximize social welfare, subject to resource and incentive constraints. The government's objective function is an expected concave sum of individual utilities:

$$\int_{\underline{n}}^{\bar{n}} (1 - \pi(z_n))G(u_n^U) + \pi(z_n)G(u_n^P) dF(n), \quad G'(\cdot) > 0, \quad G''(\cdot) < 0, \quad (8)$$

where  $G(\cdot)$  is the social welfare function.  $u_n^U \equiv c_n - v(z_n/n)$  and  $u_n^P \equiv u_n^U - P(l^* - z_n/n)$  denote the utility levels of the unpenalized and penalized individuals, respectively. Our specification for social welfare treats all sources of inequality the same. In particular, the government is averse to inequality generated by differences in earning ability, and the government is averse to income inequality

among individuals of the same earning ability caused by penalties.<sup>12</sup> Due to quasi-linearity of private utility there is no social desire to redistribute income if the social welfare function is utilitarian.

The government is constrained in its ability to redistribute income, because the ability of individuals is private information. However, the government can infer the ability of an individual from costly monitoring activities or it can induce self-selection by sacrificing on redistribution. The total cost of monitoring is given by:

$$\int_{\underline{n}}^{\bar{n}} k(\pi(z_n))dF(n), \quad k(0) = 0, \quad k'(\cdot), k''(\cdot) > 0. \quad (9)$$

The cost of monitoring is increasing and convex in the monitoring probability  $\pi$ . Since there is a perfect mapping between ability  $n$  and labor earnings  $z_n$ , we can also write  $\pi(\cdot)$  as a function of the ability level  $n$ , where we use the short-hand notation  $\pi(z_n) = \pi_n$ . However,  $\pi'(z_n) \equiv \frac{d\pi_n}{dz_n}$  always denotes the derivative of monitoring with respect to gross earnings. We assume that monitoring is sufficiently costly so that the government will not choose to monitor all individuals, i.e.,  $\pi_n < 1$  for all  $n$ . Hence, the government will rely on both self-selection and monitoring to obtain information on earning ability.

The economy's resource constraint implies that total labor earnings equal aggregate consumption plus monitoring costs:

$$\int_{\underline{n}}^{\bar{n}} z_n dF(n) = \int_{\underline{n}}^{\bar{n}} \left( (1 - \pi(z_n))c_n^U + \pi(z_n)c_n^P + k(\pi(z_n)) \right) dF(n). \quad (10)$$

By defining unpenalized consumption as  $c_n \equiv c_n^U = c_n^P + P(z_n, n)$ , we can write for aggregate consumption:

$$\int_{\underline{n}}^{\bar{n}} \left( (1 - \pi(z_n))c_n^U + \pi(z_n)c_n^P \right) dF(n) = \int_{\underline{n}}^{\bar{n}} (c_n - \pi(z_n)P(l^* - z_n/n)) dF(n). \quad (11)$$

Hence, using Eq. (11) the economy's resource constraint (Eq. (10)) can be rewritten as:

$$\int_{\underline{n}}^{\bar{n}} (z_n + \pi(z_n)P(l^* - z_n/n)) dF(n) = \int_{\underline{n}}^{\bar{n}} (c_n + k(\pi(z_n))) dF(n). \quad (12)$$

We do not need to consider the government budget constraint, since it is automatically satisfied by Walras' law if the individual budget constraints and the economy's resource constraint are satisfied.

<sup>12</sup> Our specification for social welfare is conceptually closest to the case where individuals would be risk averse and the government would be utilitarian. However, for analytical and numerical tractability we had to assume risk-neutrality of individuals. If, in contrast, we would let the government maximize a concave transformation of ex ante expected individual utilities, i.e.,  $G((1 - \pi(z_n))u_n^U + \pi(z_n)u_n^P)$ , with  $G'(\cdot) > 0$ , and  $G''(\cdot) < 0$ , the government would no longer value any inequality within ability groups caused by monitoring, only inequality generated by differences in ability levels. We deem this an undesirable property for our analysis.

In line with Mirrlees (1971), we assume that the government fully commits to the tax and monitoring schedules before individuals make their decisions.<sup>13</sup> The timing of the model is as follows:

1. The government announces the exogenously given penalty function  $P(l^* - l_n)$ , the optimal non-linear income tax  $T(z_n)$ , the optimal non-linear monitoring schedule  $\pi(z_n)$ , and the optimal work requirement  $l^*$ .
2. Each individual  $n$  optimally chooses hours worked  $l_n$ .
3. The government observes the labor incomes  $z_n$  chosen by each individual  $n$ , and taxes income and monitors individuals accordingly. The government penalizes all monitored individuals according to the penalty function.
4. Individuals receive utility from consumption and leisure.

By the revelation principle, any indirect mechanism can be replicated with an incentive-compatible direct mechanism (Myerson, 1979; Harris and Townsend, 1981). Therefore, we can find the optimal second-best allocation by maximizing welfare subject to feasibility and incentive-compatibility constraints. We can decentralize the optimal second-best allocation as a competitive market outcome through the non-linear tax and monitoring schedules.

### 3.3. First-order incentive compatibility

By using the envelope theorem we can derive a differential equation for the utility function  $u_n$  which is a necessary condition for incentive compatibility. The next subsection derives the conditions under which the first-order condition is also sufficient. The incentive-compatibility constraint is found by totally differentiating Eq. (2) with respect to  $n$ :

$$\frac{du_n}{dn} = \frac{\partial u(z_n, n)}{\partial n} + \frac{\partial u(z_n, n)}{\partial z_n} \frac{dz_n}{dn} = \frac{l_n}{n} (v'(l_n) - \pi(z_n)P'(l^* - l_n)), \quad \forall n, \tag{13}$$

where  $\frac{\partial u(z_n, n)}{\partial z_n} = 0$  due to the individual's first-order condition in Eq. (3). Thus, if the optimal allocation satisfies Eq. (13), individuals' first-order conditions for utility maximization are also satisfied.

### 3.4. Second-order incentive compatibility

Without further restrictions we cannot be certain that the optimal allocation derived under the first-order incentive-compatibility constraint (Eq. (13)) is also implementable. An implementable allocation should satisfy additional requirements to ensure that the first-order approach also respects the second-order conditions for utility maximization. The next Lemma summarizes the requirements for second-order incentive compatibility.

**Lemma 1.** *Second-order conditions for utility maximization are satisfied under the first-order approach if the following conditions hold at the optimal allocation for all  $n$ :*

- i) *Single-crossing conditions on the utility and penalty functions are satisfied:*

$$\frac{\partial (v'(l_n)/n)}{\partial n} - \frac{\pi(z_n)P'(l^* - l_n)}{n^2} (\varepsilon_n^P - 1) + \frac{l_n \pi'(z_n)}{n} P'(l^* - l_n) \leq 0, \tag{14}$$

where  $\varepsilon_n^P \equiv \frac{P'(l^* - l_n)l_n}{P(l^* - l_n)}$  is the elasticity of the penalty function,

- ii)  $z_n$  is non-decreasing in ability:

$$\frac{dz_n}{dn} \geq 0. \tag{15}$$

**Proof.** The second-order condition for the utility-maximization problem (Eq. (2)) is given by:

$$\frac{\partial^2 u(z_n, n)}{\partial z_n^2} \leq 0, \quad \forall n. \tag{16}$$

This second-order condition can be rewritten in a number of steps. Totally differentiating the first-order condition (Eq. (3)) gives:

$$\frac{\partial^2 u(z_n, n)}{\partial z_n^2} \frac{dz_n}{dn} + \frac{\partial^2 u(z_n, n)}{\partial z_n \partial n} = 0, \quad \forall n. \tag{17}$$

Substitution of this result in Eq. (16) implies that the second-order condition is equivalent to:

$$\frac{\partial^2 u(z_n, n)}{\partial z_n \partial n} \left( \frac{dz_n}{dn} \right)^{-1} \geq 0, \quad \forall n. \tag{18}$$

Differentiating the first-order condition (Eq. (3)) with respect to  $n$  and substituting the result yields:

$$\left( \frac{\partial (v'/n)}{\partial n} + \frac{\pi P'}{n^2} \left( 1 - \frac{P' l_n}{P'} \right) + \frac{l_n \pi' P'}{n} \right) \left( \frac{dz_n}{dn} \right)^{-1} \leq 0, \quad \forall n. \tag{19}$$

The inequality holds if all conditions of the Lemma are satisfied. ■

The single-crossing condition and the monotonicity of gross earnings are well-known from the Mirrlees model (Mirrlees, 1971; Ebert, 1992). The single-crossing condition ensures that – at the same income-consumption bundle – individuals with a higher ability have a larger marginal willingness to work. In our model, the single-crossing condition contains three elements. The first is the standard Spence-Mirrlees condition on the utility function, i.e.,  $\frac{\partial (v'(l_n)/n)}{\partial n} < 0$ . If this term is negative, the marginal disutility of work for individuals with a higher ability level is lower. Most utility functions considered in the literature exhibit this property, including the utility function we adopt in our simulations. The sign of the second term is determined by  $\pi P'(\varepsilon_n^P - 1)/n^2$ . Intuitively, it is more costly for high-ability individual to mimic a low-ability individual if  $\frac{\partial (P'/n)}{\partial n} > 0$ . That is, the marginal penalty of earning a lower income increases with ability.  $\frac{\partial (P'/n)}{\partial n} > 0$  is equivalent to  $\varepsilon_n^P > 1$ . Intuitively, if the elasticity of the marginal penalty is larger, penalties become increasingly more severe for high-ability individuals mimicking low-ability individuals. The third term,  $l_n \pi' P'/n$ , concerns the slope of the monitoring schedule, which determines its sign, since  $P' > 0$ . If the marginal monitoring probability decreases in gross earnings ( $\pi' < 0$ ) individuals will work harder in order to decrease the probability of being monitored and penalized. The sign of the last term is determined by the endogenous monitoring schedule. Hence, high-ability individuals

<sup>13</sup> Roberts (1984) shows that a time-consistency problem may emerge in Mirrlees (1971) when the government cannot credibly commit to its announced income-tax schedule. If all types truthfully reveal their ability, the government wants to renege on its announced income-tax schedule and levy individualized lump-sum taxes based on ability instead. However, rational individuals anticipating that the government will do this will not reveal any information in the first place. Hence, the optimal redistribution problem degenerates.

can be induced to self-select into higher income-consumption bundles, unless the monitoring probability increases too fast with ability.

A second requirement to induce self-selection is that gross earnings are indeed increasing with ability at the optimal schedule. Consequently, a tax schedule that provides higher income to higher ability individuals induces self-selection of higher ability types into higher income-consumption bundles. In the remainder we assume that all the conditions derived in Lemma 1 hold at the optimal allocation. In our simulations, we check the second-order sufficiency conditions ex-post and we always confirm that they are respected.<sup>14</sup>

#### 4. Optimal second-best allocation

The optimization problem with monitoring can be specified formally as follows:

$$\max \int_{\underline{n}}^{\bar{n}} [(1 - \pi_n) G(u_n + \pi_n P(l^* - z_n/n)) + \pi_n G(u_n - (1 - \pi_n) P(l^* - z_n/n))] f(n) dn, \quad (20)$$

$$\text{s.t.} \int_{\underline{n}}^{\bar{n}} [z_n + \pi_n P(l^* - z_n/n) - c_n - k(\pi_n)] f(n) dn = 0, \quad (21)$$

$$\frac{du_n}{dn} = \frac{l_n}{n} (v'(l_n) - \pi_n P'(l^* - z_n/n)), \quad \forall n, \quad (22)$$

$$u_n = c_n - \pi_n P(l^* - z_n/n) - v(z_n/n), \quad \forall n, \quad (23)$$

$$\pi_n \geq 0, \quad \forall n, \quad (24)$$

where utility of unpenalized and penalized individuals is, respectively, written as  $u_n^U = u_n + \pi_n P(l^* - z_n/n)$  and  $u_n^P = u_n - (1 - \pi_n) P(l^* - z_n/n)$ . The final constraint assumes that the probability of monitoring cannot be smaller than zero. We assume that the cost of monitoring is sufficiently large to ensure that the constraint  $\pi_n \leq 1$  is never binding.

We formulate a Lagrangian for this optimization problem, where  $\lambda$  is the multiplier of the economy's resource constraint (Eq. (21)),  $\theta_n$  denotes the multiplier on the incentive compatibility constraint (Eq. (22)),  $\mu_n$  is the multiplier for the definition of utility (Eq. (23)), and  $\eta_n$  is the Kuhn–Tucker multiplier of the non-negativity constraint on  $\pi_n$  (Eq. (24)). After integrating  $\theta_n \frac{du_n}{dn}$  by parts, we can write the Lagrangian function for this problem as:

$$\begin{aligned} \mathcal{L} \equiv & \int_{\underline{n}}^{\bar{n}} [(1 - \pi_n) G(u_n + \pi_n P(l^* - z_n/n)) \\ & + \pi_n G(u_n - (1 - \pi_n) P(l^* - z_n/n))] f(n) dn \\ & + \lambda \int_{\underline{n}}^{\bar{n}} [z_n + \pi_n P(l^* - z_n/n) - c_n - k(\pi_n)] f(n) dn \\ & - \int_{\underline{n}}^{\bar{n}} \frac{\theta_n z_n}{n^2} [v'(z_n/n) - \pi_n P'(l^* - z_n/n)] dn + \theta_{\bar{n}} u_{\bar{n}} - \theta_{\underline{n}} u_{\underline{n}} \\ & - \int_{\underline{n}}^{\bar{n}} \frac{d\theta_n}{dn} u_n dn + \int_{\underline{n}}^{\bar{n}} \mu_n [u_n - c_n + v(z_n/n) + \pi_n P(l^* - z_n/n)] dn \\ & + \eta_n \pi_n dn, \end{aligned} \quad (25)$$

where  $c_n$ ,  $z_n$ ,  $\pi_n$ ,  $u_n$  and  $l^*$  are the control variables.

The necessary first-order conditions are given by:

$$\frac{\partial \mathcal{L}}{\partial c_n} = 0 : -\lambda f(n) - \mu_n = 0, \quad \forall n, \quad (26)$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial z_n} = 0 : & \left[ (1 - \pi_n) \pi_n \frac{P'(\cdot)}{n} (G(u_n^P) - G(u_n^U)) + \lambda \left( 1 - \frac{\pi_n P'(\cdot)}{n} \right) \right] f(n) \\ & - \theta_n \left( \frac{v'(\cdot) + z_n v''(\cdot)/n - \pi_n (P'(\cdot) - z_n P''(\cdot)/n)}{n^2} \right) \\ & + \mu_n \left( \frac{v'(\cdot) - \pi_n P'(\cdot)}{n} \right) = 0, \quad \forall n, \end{aligned} \quad (27)$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \pi_n} = 0 : & \left[ -G(u_n^U) + (1 - \pi_n) P(\cdot) G'(u_n^U) \right. \\ & \left. + \pi_n P(\cdot) G'(u_n^P) + G(u_n^P) - \lambda (k'(\pi_n) - P(\cdot)) \right] f(n) \\ & + \frac{z_n \theta_n}{n^2} P'(\cdot) + \mu_n P(\cdot) + \eta_n = 0, \quad \forall n, \end{aligned} \quad (28)$$

$$\eta_n \pi_n = 0, \quad \eta_n \geq 0, \quad \pi_n \geq 0, \quad \forall n, \quad (29)$$

$$\frac{\partial \mathcal{L}}{\partial u_n} = 0 : \frac{d\theta_n}{dn} = \left[ (1 - \pi_n) G'(u_n^U) + \pi_n G'(u_n^P) \right] f(n) + \mu_n, \quad \forall n, \quad (30)$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial l^*} = 0 : & \int_{\underline{n}}^{\bar{n}} \left[ (1 - \pi_n) \pi_n P'(\cdot) (G'(u_n^U) - G'(u_n^P)) + \lambda \pi_n P'(\cdot) \right] f(n) dn \\ & + \int_{\underline{n}}^{\bar{n}} \left[ \frac{\pi_n \theta_n z_n}{n^2} P''(\cdot) + \mu_n \pi_n P'(\cdot) \right] dn = 0, \end{aligned} \quad (31)$$

$$\frac{\partial \mathcal{L}}{\partial u_n} = \frac{\partial \mathcal{L}}{\partial u_{\bar{n}}} = 0 : \theta_{\underline{n}} = \theta_{\bar{n}} = 0. \quad (32)$$

Compared to the analysis of Mirrlees there are three new first-order conditions. Eq. (28) states the optimal monitoring condition, Eq. (29) state the Kuhn–Tucker conditions for the non-negativity constraint on  $\pi_n$ , and Eq. (31) is the condition for the optimal labor requirement.

##### 4.1. Optimal wedge on labor

Proposition 2 gives the conditions for optimal income redistribution.

**Proposition 2.** The optimal net marginal wedge on labor  $\mathcal{W}_n$  at each ability level satisfies:

$$\frac{\mathcal{W}_n}{1 - \mathcal{W}_n} = A_n B_n C_n - D_n, \quad \forall n, \quad (33)$$

where

$$A_n \equiv 1 + \frac{1}{\varepsilon_n} + \pi_n \frac{P'(\cdot)}{v'(\cdot)} (\varepsilon_n^P - 1), \quad (34)$$

$$B_n \equiv \frac{\int_{\underline{n}}^{\bar{n}} (1 - g_m) f(m) dm}{1 - F(n)}, \quad (35)$$

$$C_n \equiv \frac{1 - F(n)}{n f(n)}, \quad (36)$$

$$D_n \equiv \frac{P'(\cdot)}{v'(\cdot)} \sigma_n, \quad (37)$$

<sup>14</sup> If the Spence–Mirrlees or monotonicity conditions are violated, then bunching generally occurs, possibly also at zero earnings (see also Ebert, 1992). Bunching at zero labor income due to income effects cannot occur, since we assumed quasi-linear utility (see also Seade, 1977). In all our simulations we allow for an atom of individuals with zero earning ability that are bunched at zero labor earnings so as to avoid counterfactual zero optimal marginal tax rates at the bottom.



$\sigma_n \equiv \frac{(1-\pi_n)\pi_n(G'(u_n^p)-G'(u_n^u))}{\lambda} > 0$  is a measure for the welfare cost of inequality between penalized and unpenalized individuals at ability level  $n$ ,  $\varepsilon_n \equiv \left(\frac{l_n v'(l_n)}{v'(l_n)}\right)^{-1} > 0$  is the compensated wage elasticity of labor supply, and  $g_n \equiv \frac{(1-\pi_n)G'(u_n^p)+\pi_n G'(u_n^u)}{\lambda} > 0$  is the average, social marginal value of income, expressed in money units, for individuals at ability level  $n$ .

**Proof.** Integrate Eq. (30) using a transversality condition from Eq. (32). It follows that  $\theta_n = \lambda \int_n^{\bar{n}} (1 - g_m) f(m) dm$ . Substitute this result and Eq. (26) in Eq. (27), use Eq. (4), and simplify to obtain the proposition. ■

The  $A_n$ -term is related to the inverse of the efficiency cost of the labor wedge at income level  $z_n$ . The second term in  $A_n$ ,  $1/\varepsilon_n$ , is the inverse of the labor-supply elasticity and it enters because the deadweight loss of the wedge increases in the labor-supply elasticity. The third term represents the efficiency gain of monitoring. As noted before, penalties are useful in separating high- and low-ability individuals if the elasticity of the penalty function  $\varepsilon^p$  is larger than 1. Penalties are more effective if the elasticity increases. The latter effect is stronger if the monitoring intensity  $\pi$  is larger. Finally, penalties are better at providing work incentives if the marginal penalty becomes more important relative to the marginal disutility of labor,  $\frac{P'}{v}$ . Hence, in comparison to the optimal wedge without monitoring (cf. Diamond, 1998; Saez, 2001) monitoring reduces the efficiency cost of taxation provided the elasticity of the penalty function is larger than 1.

The  $B_n$ -term measures the equity gain of an increase in the labor wedge at income level  $z_n$ . The first term, 1, captures the revenue gain of a larger marginal labor wedge at  $n$ , such that individuals with an income level above  $z_n$  pay one unit of extra income tax. The welfare loss of extracting one unit of income from the individuals above  $n$  is  $g_m$  for all individuals  $m \geq n$ . Therefore,  $\int_n^{\bar{n}} (1 - g_m) dF(m)$  measures the average redistributive gain of the labor wedge at  $n$ . The  $B_n$ -term is not directly affected by monitoring. Since welfare weights  $g_n$  are always declining with income,  $B_n$  rises with income (see also Diamond, 1998).

$C_n$  is the inverse relative hazard rate of the skill distribution. Its numerator is the fraction of the population whose net income is decreased by increasing the wedge and its denominator captures the size of the tax base that is distorted by the wedge. Hence, the numerator in  $C_n$  gives weights to average equity gains in  $B_n$  and the denominator to average efficiency losses in  $A_n$ — as in the Mirrlees (1971) model without monitoring. The numerator of  $C_n$  always declines with income; there are fewer individuals paying the marginal tax rate if the tax rate is increased at a higher income level. Hence, for a given  $B_n$  the total distributional benefits of raising the labor wedge fall as the income level rises. For a unimodal skill distribution the denominator of  $C_n$  always increases with income before the mode, since both  $n$  and  $f(n)$  are rising. Thus, labor wedges always decrease with income before modal income. After the mode,  $f(n)$  falls, although  $n$  continues to rise with income. Hence, it depends on the empirical distribution of  $n$  whether  $C_n$  rises or falls with income after modal income. For most empirical distributions,  $C_n$  appears to rise after the mode and converges to a constant at the top (see also Diamond, 1998; Saez, 2001; Zoutman et al., 2015).

Finally,  $D_n$  measures the welfare loss associated with within-ability group inequality. Earnings at  $n$  decrease if the labor wedge increases. Therefore, the penalty at  $n$  increases, which in turn increases inequality between monitored and unmonitored individuals.  $\sigma_n$  measures the marginal welfare cost of this within-ability group inequality. The effect of a wedge on within-ability group inequality is increasing in the relative importance of the penalty function with respect to the marginal disutility of labor (expressed

in monetary units),  $\frac{P'}{v}$ .  $D_n$  increases in the monitoring probability for  $\pi_n < .5$  because the within-ability group variance of utility is increasing in  $\pi_n$  for  $\pi_n < .5$ . Finally,  $D_n$  is increasing in the concavity of the welfare function, because the difference in welfare weights between penalized and unpenalized individuals,  $\frac{G'(u_n^p)-G'(u_n^u)}{\lambda}$ , is larger if the government is more inequality averse.

We can summarize the impact of monitoring on optimal labor wedges as follows. Monitoring decreases the efficiency cost of setting a higher labor wedge, but introduces within-ability group inequality. Therefore, the total effect of monitoring on the optimal labor wedge is theoretically ambiguous. Our simulations below demonstrate that the efficiency gains of monitoring outweigh the distributional loss due to inequality between monitored and non-monitored individuals.

We can derive the non-linear tax function, which implements the second-best allocation as the outcome of decentralized decision making in a competitive labor market. Substituting Eq. (3) into Eq. (33) yields:

$$\frac{T'(z_n) + \pi'(z_n)P(l^* - z_n/n) - \pi(z_n)P'(l^* - z_n/n)/n}{1 - T'(z_n) - \pi'(z_n)P(l^* - z_n/n) + \pi(z_n)P'(l^* - z_n/n)/n} = A_n B_n C_n - D_n, \quad \forall n. \tag{38}$$

Thus, when we know the optimal monitoring schedule  $\pi(z_n)$ , this equation implicitly defines the optimal non-linear income tax function  $T(z_n)$ .

#### 4.2. Optimal monitoring

The next proposition derives the optimal monitoring schedule.

**Proposition 3.** *The optimal level of monitoring at each ability level follows from*

$$k'(\pi_n) + \Delta_n - g_n P(\cdot) \geq \left( \frac{\frac{v_n}{1-v_n} + D_n}{A_n} \right) l_n P'(\cdot) \quad \forall n, \tag{39}$$

where  $\Delta_n \equiv \frac{G'(u_n^u)-G'(u_n^p)}{\lambda}$  is the welfare difference between a penalized and an unpenalized individual expressed in money units. If  $\pi_n > 0$ , the equation holds with equality.

**Proof.** Substitute Eq. (26) into Eq. (28), rearrange terms, employ the definitions for  $B_n$  in Eq. (35) and  $C_n$  in Eq. (36), and use the fact that  $\eta_n \geq 0$ . Finally, substitute Eq. (33) for  $B_n C_n$  to obtain the expression. By Eq. (29)  $\eta_n$  only equals zero if  $\pi_n > 0$  and therefore the equation holds with equality if  $\pi_n > 0$ . ■

The first term on the left-hand side in Eq. (39) is the marginal cost of raising the monitoring intensity. The second and third terms on the left-hand side jointly represent the welfare effect of a compensated increase in the monitoring probability. That is, the welfare effect of an increase in the monitoring probability, while keeping expected utility at skill level  $n$  unchanged. The second term represents the uncompensated, direct welfare loss of an increase in the monitoring probability. If the monitoring probability increases, there will be more penalized and less unpenalized individuals. Therefore, the loss is equal to the welfare difference between penalized and unpenalized individuals. The third term represents the welfare gain associated with the compensation to keep expected utility unchanged if the monitoring probability is increased. The compensation at ability level  $n$  requires a transfer of  $P$  and its associated welfare effect is thus given by  $g_n P$ . In Lemma 2 we derive how the compensated welfare effect of monitoring changes with the monitoring

probability for given levels of utility in monitored and unmonitored states.

**Lemma 2.** *The compensated welfare effect of the monitoring probability  $\Delta_n - g_n P(\cdot)$  is decreasing in  $\pi_n$ , positive if  $\pi_n = 0$ , and negative if  $\pi_n = 1$  for given levels of utility in penalized and unpenalized states.*

**Proof.** By a first-order Taylor expansion around  $u_n^U$  we can write  $\Delta_n$  as:

$$\Delta_n = \frac{G(u_n^U) - G(u_n^P)}{\lambda} = \frac{G'(u_n^U)(u_n^U - u_n^P)}{\lambda} + R(P) = \frac{G'(u_n^U)P}{\lambda} + R(P). \tag{40}$$

where  $R(P)$  is a second-order remainder term. Similarly, a first-order Taylor expansion around  $u_n^P$  yields:

$$\Delta_n = \frac{G'(u_n^P)P}{\lambda} - \hat{R}(P), \tag{41}$$

where  $\hat{R}(P)$  is again a second-order remainder term. By concavity of  $G$  both remainder terms are positive for  $P > 0$ :  $R(P), \hat{R}(P) > 0$ . Now multiply Eq. (40) with  $(1 - \pi_n)$  and Eq. (41) with  $\pi_n$  and add them to find

$$\Delta_n - g_n P = (1 - \pi_n) R(P) - \pi_n \hat{R}(P). \tag{42}$$

The right-hand side gives the compensated welfare effect of the monitoring probability, which is, *ceteris paribus*, decreasing in  $\pi_n$ , always positive if  $\pi_n = 0$ , and always negative if  $\pi_n = 1$ . ■

The right-hand side of Eq. (39) represents the marginal benefits of monitoring. The benefits of monitoring increase in the marginal penalty  $P(\cdot)$ , which can be interpreted as the power of the penalty function. In addition, the marginal benefits of monitoring increase if labor-supply distortions are larger, i.e., if the labor wedge  $\frac{\mathcal{W}_n}{1-\mathcal{W}_n}$  is larger or if the efficiency cost of taxation is larger, as captured by  $1/A_n$ . The benefits of monitoring also increase in within-ability group inequality  $D_n$ . Intuitively, as more monitoring leads to larger labor supply, the expected penalty decreases. Hence, monitoring helps to reduce within-ability group inequality.

Turning back to the optimal monitoring condition, from Proposition 3 it follows that the government does not engage in monitoring if and only if (evaluated at a no-monitoring equilibrium with  $\pi_n = 0$ ):

$$k'(0) + \Delta_n - g_n P(\cdot) > \left( \frac{\mathcal{W}_n}{1-\mathcal{W}_n} + D_n \right) I_n P'(\cdot), \quad \forall n. \tag{43}$$

That is, if the marginal costs of monitoring are higher than the marginal benefits for all types. By evaluating Proposition 1 at  $\pi_n = 0$  it easily follows that the optimal allocation is the allocation derived in Mirrlees (1971). Mirrlees (1971) is thus a special case of our model where monitoring is prohibitively expensive.

### 4.3. Optimal work requirement

The next proposition derives the optimal work requirement  $l^*$ .

**Proposition 4.** *The optimal work requirement  $l^*$  is implicitly determined by:*

$$\int_{\mathcal{P}} \sigma_n f(n) dn = \int_{\mathcal{P}} \left( \frac{\mathcal{W}_n}{1-\mathcal{W}_n} + D_n \right) \pi_n e_n^P f(n) dn, \tag{44}$$

where  $\mathcal{P}$  is the set of ability levels  $n$  where hours worked are smaller than the work requirement:  $\mathcal{P} \equiv \{n \in [\underline{n}, \bar{n}] : l^* < l_n\}$ .

**Proof.** Substitute Eq. (26) into Eq. (31), divide by  $\lambda$ , use  $\sigma_n \equiv \frac{(1-\pi_n)\pi_n(G'(u_n^U) - G'(u_n^P))}{\lambda} > 0$  to find

$$\int_{\underline{n}}^{\bar{n}} \sigma_n P'(\cdot) f(n) dn = \int_{\underline{n}}^{\bar{n}} \frac{\theta_n/\lambda}{nf(n)} \pi_n P'(\cdot) I_n f(n) dn. \tag{45}$$

Use  $\theta_n/\lambda = \int_n^{\bar{n}} (1 - g_m) f(m) dm$ , employ the definitions for  $B_n$  and  $C_n$  and substitute Eq. (33) for  $B_n C_n$  to obtain the expression. For ability levels where  $l_n > l^*$  the penalty function, as well as its first and second derivatives, equal zero:  $P(l^* - l_n) = P'(l^* - l_n) = P''(l^* - l_n) = 0$ . It follows that for those ability levels all terms equal zero and the integral  $\int_{\underline{n}}^{\bar{n}}$  can be replaced by  $\int_{\mathcal{P}}$ . In subdomain  $\mathcal{P}$ , we have  $l_n \geq l^*$ , and hence, the marginal penalty is strictly positive, so that  $P'(l^* - l_n) > 0$ . Therefore, we can divide both sides by  $P'(l^* - l_n)$ , and substitute  $e_n^P \equiv \frac{P'(l^* - l_n) l_n}{P(l^* - l_n)}$  on the right-hand side. ■

The left-hand side of Eq. (44) represents the marginal welfare cost of increasing the work requirement  $l^*$ . A marginal increase in the work requirement – *ceteris paribus* – increases the penalty provided to monitored individuals, resulting in an increase in within-ability group inequality. This is represented by the marginal welfare cost of within-ability group inequality  $\sigma_n$ . If the government would care little for within-ability group inequality, the minimum-work requirement  $l^*$  would be set at high levels.

The right-hand side of Eq. (44) provides the marginal welfare gain of increasing the work requirement  $l^*$ . A marginal increase in the work requirement decreases the burden of taxation, since it increases the incentive to work. This effect is stronger if the penalty function is more effective to boost labor supply (i.e.,  $e_n^P$  is larger), and if more individuals are monitored (i.e.,  $\pi_n$  is larger). Moreover, the welfare gain is larger if labor-supply distortions are larger. These distortions increase in the labor wedge  $\frac{\mathcal{W}_n}{1-\mathcal{W}_n}$ , and the efficiency cost of taxation  $1/A_n$ . Additionally, increasing the work requirement induces individuals to work harder, which reduces their expected penalty, and hence, within-group inequality, as captured by  $D_n$ . The work requirement thus serves to enhance the efficiency of the income tax.

If at a given ability level  $n$ , individuals work more than the work requirement, i.e.,  $l_n > l^*$ , the penalty, the marginal penalty, and the change in the marginal penalty are all zero, i.e.,  $P(\cdot) = P'(\cdot) = P''(\cdot) = 0$ . Hence, for these individuals there are neither costs nor benefits from setting a marginally higher work requirement  $l^*$ . Consequently, Eq. (44) sums marginal costs and the marginal benefits of raising the labor requirement  $l^*$  only for those individuals whose work requirement is binding ( $P(\cdot) > 0$ ).

### 4.4. Boundary results

In the next Proposition we derive the optimal wedge and monitoring probability at the bottom and the top of the ability distribution.<sup>15</sup>

**Proposition 5.** *If the income distribution is bounded at the top,  $\bar{n} < \infty$ , the optimal wedge and monitoring probabilities at the bounds of the ability distribution are:*

$$\mathcal{W}_{\underline{n}} = \mathcal{W}_{\bar{n}} = \pi_{\underline{n}} = \pi_{\bar{n}} = 0. \tag{46}$$

<sup>15</sup> Due to the absence of income effects in labor supply, bunching at zero labor earnings is not an issue in deriving the boundary results (see also Seade, 1977).

If non-distorted labor supply at the endpoints is above the work requirement  $l^*$ , then the penalties are zero at the end-points, and marginal tax rates are also zero at the endpoints:

$$T'(z_n) = T'(z_{\bar{n}}) = 0. \tag{47}$$

**Proof.** From Eq. (33) follows that  $(\frac{\mathcal{W}_n}{1-\gamma\mathcal{V}_n} + D_n)/A_n = B_n C_n$ . The transversality conditions (Eq. (32)) imply  $B_n C_n = B_{\bar{n}} C_{\bar{n}} = 0$ . At the extremes, the optimal monitoring condition (Eq. (39)), therefore simplifies to  $\Delta_n - g_n P + k'(\pi_n) \geq 0$ . Evaluate this expression at  $\pi_n = 0$ :

$$\Delta_n - g_n P + k'(0) = R(P) + k'(0) \geq 0, \tag{48}$$

where  $R(P) > 0$  is a second-order remainder term, and the second step follows from Lemma 2. The condition is always satisfied at  $\pi_n = 0$ . Hence,  $\pi_n = 0$  is optimal at the extremes. The optimal wedges in Eq. (33) at the extremes are zero, because the product  $B_n C_n$  is zero by the transversality conditions, and  $D_n$  is zero, since  $\pi_n = 0$ . When the wedge is zero, labor supply is at its non-distorted level. Therefore, if the work requirement  $l^*$  is lower than non-distorted labor supply, it follows that  $P(\cdot) = 0$  at the endpoints. In that case, using  $\pi_n = P(\cdot) = 0$  in Eq. (4) demonstrates that  $\mathcal{W}_n = \mathcal{W}_{\bar{n}} = T'(z_n) = T'(z_{\bar{n}}) = 0$ . ■

Proposition 4 establishes that the optimal zero wedge at the bottom and top of the model without monitoring carries over to the model with monitoring (Sadka,1976; Seade,1977). Intuitively, the wedge at  $n$  redistributes income from individuals above  $n$  to the government, and, hence indirectly to individuals below  $n$ . There are no individuals above  $n$ – and no individuals below  $\bar{n}$ . Therefore, there are no benefits associated to a positive wedge at these points of the ability distribution. However, the wedge does distort the labor-supply decision. Hence, the optimal wedge must be zero. Because the wedge is zero, there is no efficiency gain from monitoring. As a result, the optimal monitoring probability is also zero.

However, marginal tax rates at the endpoints do not necessarily need to be zero. This critically depends on the penalty function and the optimal work requirement  $l^*$ . In particular, if the work requirement  $l^*$  is larger than non-distorted labor supply, the marginal monitoring probability is non-zero at the end-points ( $\pi'(z_n) \neq 0$ ) and the expected penalty is positive. In that case, marginal tax rates at the endpoints have to be non-zero in order to compensate for the distortion caused by the change in monitoring intensity. In particular, marginal tax rates at the endpoints should be positive (negative) if  $\pi'(z_n)P(\cdot) < 0 (> 0)$ . Marginal penalties and tax rates at the endpoints are zero only if non-distorted labor supply at the end-points is higher than the work requirement  $l^*$ , so that (marginal) penalties are zero.

### 5. Simulations

In this section we use numerical simulations to establish the shape of the optimal tax and monitoring schedules. The simulations require four main ingredients: the ability distribution, the individual preferences, the social preferences, and the monitoring technology. First, we use the skill distribution from Mankiw et al. (2009). The hourly wage is used as a proxy for earning ability. We follow Mankiw et al. (2009) by assuming that wage rates follow a log-normal distribution, which is extended with a Pareto-distribution for the top tail of the wage distribution. In addition, we assume that 5% of individuals are disabled and have zero earning ability ( $\bar{n} = 0$ ), which is also based on Mankiw et al. (2009). The earnings distribution is estimated from March 2007 CPS data. This resulted in a mean log-ability of  $m = 2.76$  and a standard deviation of log ability of  $s = 0.56$ .

The Pareto-tail starts at the top 1% of the earning distribution and features a Pareto-parameter of  $\alpha = 2$ . The latter is in accordance with estimates of Saez (2001).

Second, a description of individual preferences is needed. We assume the following utility function:

$$u(c_n, l_n) = c_n + \gamma \frac{(1 - l_n)^{1-1/\varepsilon}}{1 - 1/\varepsilon}, \quad \gamma, \varepsilon > 0, \tag{49}$$

where  $\varepsilon$  is the (un)compensated elasticity of leisure demand with respect to the net marginal wage rate. The elasticity of taxable income of individual  $n$  equals  $\varepsilon(1 - l_n)/l_n$ . Labor supply  $l_n$  generally increases in individual earning ability. Hence, high-ability types are more likely to satisfy the work requirement  $l^*$ , which is in accordance with empirical evidence that high-ability types supply more labor. Labor supply is bounded above at 1 by the Inada-conditions on utility. Therefore, it is theoretically possible to have a binding minimum-hours requirement  $l^*$  for all individuals by simply setting  $l^* = 1$  for all  $n$ .<sup>16</sup> We follow the empirical literature estimating the elasticity of taxable income and calibrated the elasticity of taxable income at 0.25 (see, e.g., Saez et al., 2012). To that end, we assumed  $\varepsilon = 0.25$  and we calibrated average labor supply at  $\bar{l} = 0.5$  in the baseline.<sup>17</sup>

The third ingredient is the social welfare function. We assume an Atkinson social welfare function featuring a constant elasticity of relative inequality aversion  $\beta$ :

$$G(u_n) = \frac{u_n^{1-\beta}}{1-\beta}, \quad \beta \geq 0, \quad \beta \neq 1, \\ G(u_n) = \ln(u_n), \quad \beta = 1. \tag{50}$$

The utilitarian objective is obtained by assuming  $\beta = 0$ . A Rawlsian social welfare function results if  $\beta \rightarrow \infty$ . The baseline assumes a moderately redistributive government with  $\beta = 0.99 \approx 1$ . In the robustness analysis we also consider less redistributive governments ( $\beta = 0.25$ ) and more redistributive governments ( $\beta = 2$ ).

Finally, we need to make specific assumptions on the monitoring technology and the penalty function. Unfortunately, no empirical evidence is available that guides us to calibrate these functions. However, our theoretical model provides some restrictions on the choice of the functions. Also, we perform robustness checks on the parameter choices we have made for these functions. The cost of monitoring should be increasing and convex in the monitoring intensity  $\pi$ . We assume that the cost of monitoring is quadratic:

$$k(\pi_n) = \frac{\kappa}{2} \pi_n^2, \quad \kappa > 0, \tag{51}$$

where  $\kappa$  is a cost parameter indicating the marginal cost of a higher monitoring probability. In the baseline we assume  $\kappa = 1$ . In the robustness analysis we vary  $\kappa$  between 0.5 and 2. We provide economic justification for these parameter values by showing that the change in the monitoring probability induced by the different values of  $\kappa$  is relatively large. In addition, we show that total monitoring costs in our calibration are a small, but significant and plausible fraction of total income earned in the economy.

<sup>16</sup> We prefer our specification over the utility function used by Diamond (1998) (i.e.,  $u(c_n, l_n) = c_n - \gamma l_n^{1+1/\varepsilon}/(1+1/\varepsilon)$ ,  $\gamma, \varepsilon > 0$ ), since labor supply is potentially unbounded in the latter specification. This is implausible empirically and gives rise to numerical complications.

<sup>17</sup> The working paper version of this paper simulated the model with a slightly different utility function and an exogenous work requirement. The results do not substantively differ from the ones reported here (see also Zoutman and Jacobs, 2014).

**Table 1**  
Calibration for simulations.

Parameter	Description	Base value	High value	Low value
$m$	Mean log ability	2.76	N/A	N/A
$s$	Standard deviation log ability	0.56	N/A	N/A
$\alpha$	Pareto-parameter	2.00	N/A	N/A
$d$	Fraction of disabled individuals	0.05	N/A	N/A
$\varepsilon$	Compensated elasticity	0.25	N/A	N/A
$r$	Government revenue as fraction of GDP	0.10	N/A	N/A
$\kappa$	Cost of monitoring	1.00	2.00	0.50
$p$	Penalty parameter	3.50	4.50	2.50
$\beta$	Relative inequality aversion	1.00	2.00	0.25

We assume that the penalty function is quadratic in labor hours  $l_n$  and is given by:<sup>18</sup>

$$P = \frac{p}{2} (\max\{0, l^* - l_n\})^2, \quad p > 0, \quad (52)$$

where  $p$  is a parameter determining the severity of the penalty. The penalty is a function of the reference level of labor hours  $l^*$ . If individuals work less than the reference level they are subject to a penalty when monitored, and increasingly so if their hours worked deviate more from the reference level of hours. When the work requirement is binding, monitoring is effective in boosting labor supply. In the baseline we set  $p = 3.5$ . In the robustness checks we employ values of  $p = 2.5$  and  $p = 4.5$ .

The government-revenue requirement is exogenous and set to 10% of labor earnings in the baseline specification without monitoring, following Tuomala (1984) and Zoutman et al. (2015). The choices for all the parameters can be found in Table 1.

In the table, the first column on the right-hand side gives the base value of the parameter. In addition, we perform robustness checks with high and low parameter values for the welfare function, all parameters of the penalty function, and all parameters of the monitoring technology.<sup>19</sup>

### 5.1. Results

Fig. 2 gives the optimal wedge, tax and monitoring schedules as a function of yearly income in US dollars. The fat solid line represents the optimal tax schedule with monitoring. The dashed line is the optimal tax schedule without monitoring. The circled line is the optimal total labor wedge with monitoring. And, the thin solid line is the optimal monitoring schedule. Recall that the optimal tax schedule coincides with the optimal labor wedge if there is no monitoring.

As can be seen, the optimal labor wedge follows a U-shape both with and without monitoring. Marginal wedges are extremely large at the bottom of the labor market, relatively small for middle-income levels, and somewhat higher at the top. The shape of these schedules is largely explained by the  $B_n$  and  $C_n$  terms in Eq. (33). The  $B_n$ -term increases with income as the welfare loss of taxing away one unit of income unit from individuals above  $z_n$  decreases in  $z_n$  (see our previous discussion). The  $C_n$ -term follows a U-shape. At the bottom of the earning distribution, the density of tax payers is small, and hence, efficiency costs of marginal taxes are low. In addition, the redistributive benefits of a higher marginal tax rate are large as it is paid

by almost the entire population. Towards middle-income levels, the efficiency cost increases as the population density increases, whereas the redistributive benefits decrease as fewer individuals are paying a higher tax rate. After modal income, marginal tax distortions decline more rapidly than distributional benefits of marginal taxes, hence marginal taxes increase. These results are entirely in line with previous simulations performed in e.g., Saez (2001), Brewer et al. (2010), and Zoutman et al. (2015).

The effect of monitoring on the labor wedge is theoretically ambiguous as we derived in the previous section. However, in our simulations we see that the efficiency gain of monitoring in reducing labor distortions outweighs the distributional cost of raising within-skill group inequality. The optimal monitoring schedule also follows a U-shape.

Table 2 reports the optimal labor requirements under the different scenarios we consider in our simulations. The optimal labor requirement equals 0.94 in the baseline scenario, which is close to the maximum amount of labor that individuals can supply under our utility function. Therefore, the work requirement is binding for all individuals in the baseline, and monitoring is effective in boosting labor supply at all income levels. Hence, the optimal monitoring intensity is always positive.

The labor wedge determines the shape of the monitoring schedule, see also Eq. (39) in Proposition 3. The monitoring intensity decreases very steeply at the bottom of the income distribution. This gives individuals a strong incentive to increase their labor supply. At middle-income levels the monitoring intensity is relatively low. The monitoring intensity increases towards top-incomes. However, the effect of monitoring on the labor wedge and the tax schedule is very small at these high income levels.

The optimal tax schedule exhibits extremely large tax rates at the bottom of the earning distribution. Indeed, the government can levy tax rates above 100% at the lowest income earners. The sharp decrease in the monitoring intensity works as an implicit subsidy on labor supply and partially offsets the high explicit tax on labor supply. The poverty trap found in many countries (see, e.g., Spadaro, 2005, Brewer et al., 2010 and OECD, 2011) can thus be optimal in the presence of monitoring. Indeed, there may not be a poverty trap if the monitoring schedule provides sufficient incentives, even if the tax-benefit system itself does not provide incentives to supply labor.

Note that the optimal wedge at the top does not equal zero, as was derived in Proposition 5 for a bounded income distribution. Mirrlees (1971), Diamond (1998), and Saez (2001) show theoretically that the optimal wedge converges to a constant if the right tail of the ability distribution is Pareto distributed. In the Pareto tail of the earning distribution, the ratio of marginal distributional benefits and marginal efficiency costs of taxes becomes constant, and the tax wedge converges to a constant. Our simulations confirm that this result holds as well in the model with monitoring. However, the monitoring probability does converge to zero. With our utility function, labor supply converges to 1 when ability approaches infinity. Since the optimal work requirement  $l^*$  is smaller than 1, monitoring

<sup>18</sup> Note that with this specification of the penalty function the elasticity  $\varepsilon^p$  is not unambiguously larger than 1 so that violations of second-order conditions might occur (see Lemma 1). However, in none of our reported simulations is this the case.

<sup>19</sup> The numerical procedure we use to solve for the optimal allocation is a so-called shooting method. We solve the differential Eqs. (13) and (30) numerically for given initial values  $\theta_n$ ,  $u_n$ , and  $\lambda$ . Subsequently, we 'shoot' for initial values until we meet boundary Conditions (12) and (32). The wedge, tax, and monitoring schedule can be found using Eq. (38). A more detailed explanation of the numerical procedure can be found in the Simulation algorithm in the Appendix.



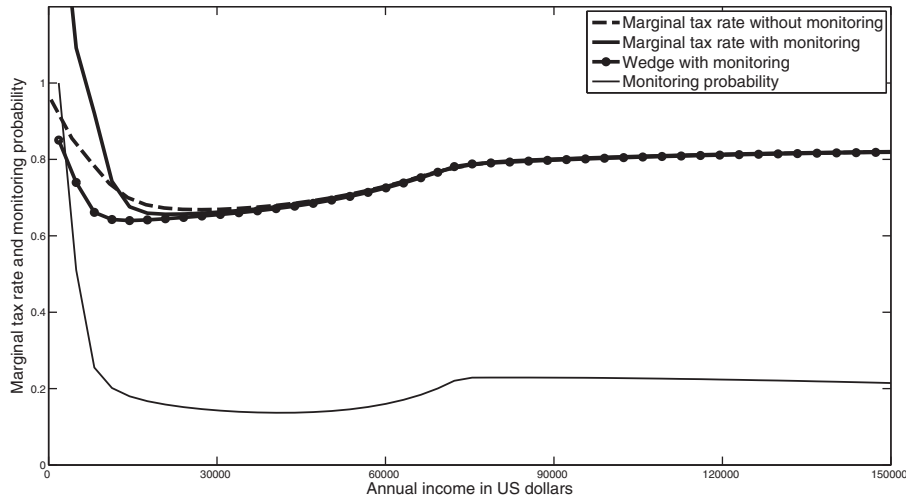


Fig. 2. The optimal wedge, tax and monitoring schedules in the baseline scenario. Baseline parameter values of the model can be found in Table 1.

at very high income levels no longer affects the incentive to work. Hence, there will be no more monitoring. It is difficult to observe the decrease in the monitoring intensity at top income levels in Fig. 2, since labor supply converges only slowly to 1.

5.2. Sensitivity analysis

In this subsection we present the sensitivity analysis of the results obtained in the previous subsection. We especially explore the sensitivity of our simulation outcomes with respect to the monitoring technology and the penalty function.

Fig. 3 summarizes the simulations when the cost of monitoring is decreased ( $\kappa = 0.5$ ) or increased ( $\kappa = 2$ ). As expected, the monitoring schedule moves up if the monitoring cost decreases and down if the cost increases. In addition, in the high-cost scenario the monitoring schedule no longer has a U-shape, since the government no longer monitors individuals with an income above 30,000 dollars. The reason is that the monitoring intensity decreases when monitoring costs increase. As a result, the government optimally reduces the work requirement when monitoring costs increase to  $l^* = 0.66$ , as can be seen from Table 2. Intuitively, the marginal benefits of the work requirement increase in the monitoring intensity (see the right-hand side of Eq. (44)). When the work requirement is no longer binding at middle- and high-income levels, it is not optimal to monitor labor effort at these income levels.

The optimal tax schedule largely remains unaffected. From the optimal tax expression in Eq. (38) we can infer that monitoring increases the optimal tax rate if the allocation remains unchanged. However, the allocation changes, since an increase in the monitoring probability increases revenue from taxation for any given tax rate. Therefore, the redistributive benefit of a marginal tax decreases at the same time. In our simulations, these two effects roughly cancel out and the optimal tax rates remain largely unaffected.

Table 2  
Optimal work requirements.

	$l^*$
Base scenario	0.94
Low monitoring cost	0.98
High monitoring cost	0.66
Low penalty	0.94
High penalty	0.88
Low inequality aversion	0.50
High inequality aversion	0.91

Fig. 4 shows the optimal tax and monitoring schedules when the penalty parameter is decreased ( $p = 2.5$ ) or increased ( $p = 4.5$ ). As can be seen, the differences in both the optimal monitoring and tax schedules with the baseline are minor. From the optimal tax formula in Eq. (33) it follows that an increase in the penalty raises the marginal tax rate for a given wedge. Second, an increase in the penalty itself may increase or decrease the optimal marginal tax rate for a given labor wedge depending on the sign of  $\pi'(z_n)$ . Third, an increase in the convexity of the penalty function decreases the efficiency cost of a wedge. Fourth, the penalty affects the monitoring probability, although the effect is ambiguous. Fifth, an increase in the penalty increases within-ability group inequality, which decreases the optimal wedge. Sixth, the allocation itself is affected, but it is a priori unclear whether higher penalties lead to more or less redistribution. Seventh, an increase the penalty parameter may either increase or decrease the labor requirement, as both the benefits and the cost of the higher labor requirement increase. The simulation outcomes confirm these theoretical ambiguities, and show that the seven effects roughly cancel out along the entire income distribution.

Finally, in Fig. 5 we simulate the optimal tax and monitoring schedules for a higher degree of inequality aversion ( $\beta = 2$ ) and a lower degree ( $\beta = 0.25$ ) of inequality aversion. As can be seen, the optimal tax rate increases in inequality aversion as should be expected, although at the bottom of the income distribution the difference is small. Intuitively, monitoring decreases the distortion of a higher tax rate, but it also creates within skill-group inequality. The poorest individuals in society are the low-income individuals who are penalized. Hence, within-ability group inequality is particularly costly if the government is strongly inequality-averse. For very low levels of income, the optimal monitoring intensity decreases, both when inequality aversion increases, and when inequality aversion decreases. At higher levels of income, within-skill group inequality aversion is less important, and the monitoring intensity unambiguously increases with inequality aversion as labor wedges are set higher when redistributive desires are stronger.

The optimal work requirement decreases (increases) when inequality aversion increases (decreases) (see Table 2). The effect of inequality aversion on the optimal labor requirement is ambiguous. On the one hand, an increase in inequality aversion increases the social cost of within-ability group inequality, which reduces the optimal work requirement, cf. the left-hand side of Eq. (44). On the other hand, an increase in inequality aversion increases the

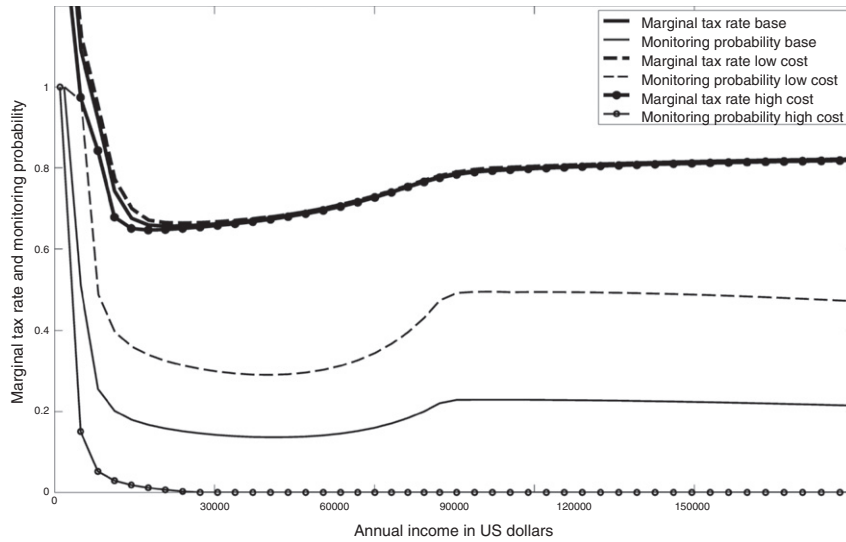


Fig. 3. Optimal tax and monitoring schedules for high ( $\kappa = 2$ ) and low ( $\kappa = 0.5$ ) marginal costs of monitoring. All other parameters take baseline values (see Table 1).

optimal labor wedge, which raises the optimal work requirement, cf. the right-hand side of Eq. (44). In our simulations, the first effect dominates the second effect when inequality aversion increases, and vice versa when inequality aversion decreases. Interestingly, the decrease in the work requirement is very large when inequality aversion decreases. As a result, the government only monitors individuals with very low incomes and the work requirement is no longer binding at higher income levels.

5.3. Allocations and welfare

Clearly, monitoring is part of the optimal redistributive tax-benefit system. But, how important is monitoring for the optimal second-best allocation and welfare? Table 3 reports the average monitoring cost  $\bar{k}/\bar{z}$ , the average penalty  $\bar{P}/\bar{z}$ , the penalty for the lowest working individual,  $P(n)/\bar{z}$ , the transfer paid out to individuals

having zero earnings,  $-T(0)/\bar{z}$ , and the change in average earnings,  $\Delta\bar{z}/\bar{z}$ . All table entries are denoted in percentages of average earnings.

The first column shows that the average monitoring cost  $\bar{k}/\bar{z}$  is relatively small: about 0.77% of average earnings in the baseline. An increase in the marginal cost of monitoring actually decreases total monitoring costs. The increase in marginal monitoring costs is accompanied by a strong decrease in the optimal monitoring intensity. Monitoring costs decrease slightly when the penalty parameter decreases, since the government relies less on monitoring when penalties are lower. However, monitoring costs are very sensitive with respect to inequality aversion, since a more inequality-averse government relies more heavily on monitoring to alleviate the equity-efficiency trade-off.

The second column represents the average penalty given to monitored individuals as a percentage of average labor earnings  $\bar{P}/\bar{z}$ . As can be seen, penalties are relatively small throughout all simulations. In the baseline, the average penalty equals 1.4% of average

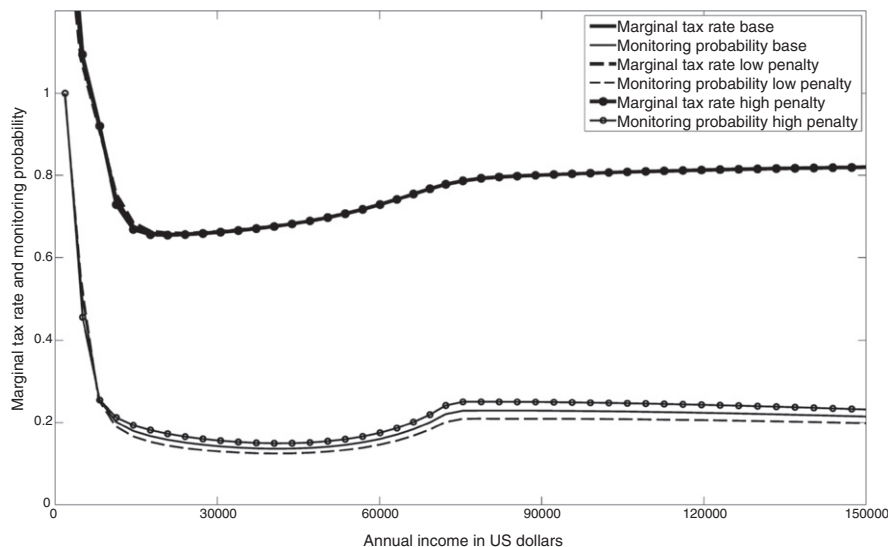


Fig. 4. Optimal tax and monitoring schedules for strong ( $p = 4.5$ ) and weak ( $p = 2.5$ ) penalties. All other parameters take baseline values (see Table 1).

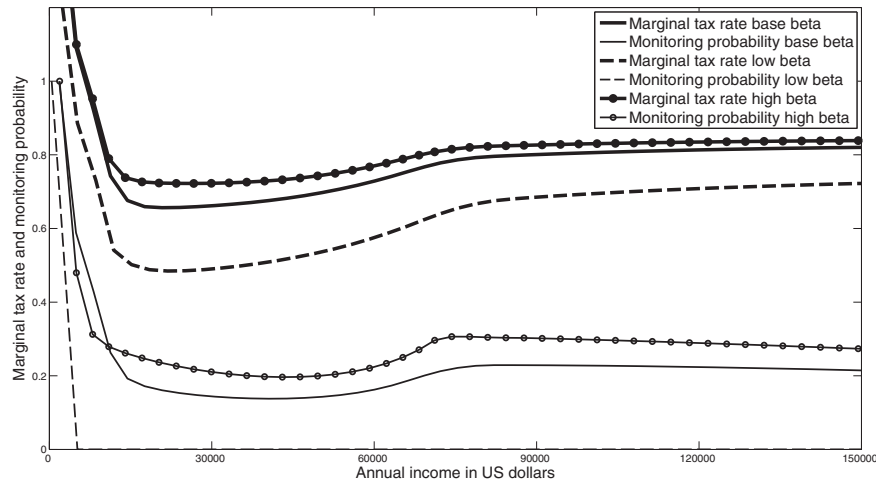


Fig. 5. The tax and monitoring schedule for a higher ( $\beta = 2$ ) and a lower ( $\beta = 0.25$ ) degree of inequality aversion. All other parameters take baseline values (see Table 1).

earnings. Penalties decrease with the monitoring cost, as the labor requirement declines when monitoring costs increase. As expected, penalties decrease when the penalty parameter decreases. However, average penalties also decrease when the penalty parameter increases, since the optimal labor requirement declines. The penalty increases (decreases) with stronger (weaker) inequality aversion, since a more (less) inequality-averse government optimally sets higher (lower) labor wedges.

The third column gives the average penalty at the bottom of the earning distribution  $P(\underline{n})/\bar{z}$ . Penalties at the bottom are relatively large, because the wedge at the bottom is large. The comparative-static effects of the penalty at the bottom are similar to the comparative statics of the average penalty.

The fourth column shows the transfer as a fraction of average earnings,  $-T(0)/\bar{z}$ , and the fifth column shows the change in average labor earnings  $\Delta\bar{z}/\bar{z}$  in comparison to the optimal tax system without monitoring. In almost all simulations, both the transfer and average labor earning increase, indicating an improvement in both equity and efficiency of the tax-transfer system. As expected, this effect decreases in the cost of monitoring. However, the increase in both the transfer and the monitoring cost remains substantial even when monitoring cost are large. These outcomes are largely explained by the fact that monitoring is most effective at the bottom of the skill distribution. At low end of the earning distribution, monitoring costs are relatively unimportant, since the density of monitored individuals is low.

The allocation is quite sensitive to the size of the penalty, despite the fact that the change in the penalty parameter does not have

a large impact on optimal tax and monitoring schedules. Intuitively, when the penalty parameter increases the government optimally reduces the labor requirement, which leaves average penalties approximately unchanged. The average marginal penalty thus increases, since the penalty function is quadratic. As a result, a larger penalty is more effective in reducing labor-supply distortions. Even if the government does not alter the tax and monitoring schedules, both equity and efficiency increase. Finally, a change in inequality aversion changes the weight given to either equity (higher transfers  $T(0)$ ) or efficiency (higher average labor earnings  $\bar{z}$ ). Remarkably, the optimal transfer increases compared to the case without monitoring even in our scenario with low inequality aversion, where average tax rates decrease significantly, and monitoring only occurs at the very bottom of the income distribution. However, in the scenario with large inequality aversion average labor earnings decrease slightly.

Finally, Table 4 reports the welfare effects of monitoring. The first column represents the income-weighted average of the marginal deadweight loss of increasing the marginal tax rate by 1%. As can be seen, monitoring decreases the marginal deadweight loss by about 1% in our baseline simulation from 0.28 to 0.26. This result is robust across all our sensitivity analyses. The last column reports the monetized welfare gain of monitoring. We compute the compensating variation by calculating the amount of resources that have to be injected into an economy without monitoring in order to attain the same social welfare as the economy with optimal monitoring. In our baseline scenario, the welfare gain is about 2.8% of average labor earnings, i.e., 2.8% of total output. The welfare gain increases if the cost of monitoring is lower and if penalties are higher. Also, an

Table 3  
Change in allocation due to monitoring. (All numbers are in percentages of average earnings).

	$\bar{k}$ $\bar{z}$	$\bar{P}$ $\bar{z}$	$\frac{P(\underline{n})}{\bar{z}}$	$\frac{-T(0)}{\bar{z}}$	$\frac{\Delta\bar{z}}{\bar{z}}$
No monitoring	0.00	0.00	0.00	35.74	0.00
Base scenario	0.77	1.40	10.77	49.33	1.29
Low monitoring cost	0.76	1.65	11.12	50.32	1.48
High monitoring cost	0.24	0.08	7.19	44.79	1.17
Low penalty	0.64	1.04	8.99	46.11	0.90
High penalty	0.76	1.12	11.14	50.64	1.60
Low inequality aversion	0.03	0.02	6.49	37.25	5.84
High inequality aversion	0.95	1.25	9.45	50.66	-0.81

Note:  $\bar{z}$  is per capita labor income in the specified calibration,  $\bar{k}$  is the per capita monitoring cost,  $\bar{P}$  is the average penalty over the monitored population,  $P(\underline{n})$  is the penalty at the lowest skill level,  $-T(0)$  is the transfer and  $\Delta\bar{z}$  is the change in average labor earnings as compared to the model without monitoring.

Table 4  
Welfare effects of monitoring.

	Marginal dead weight loss	Welfare gain
No monitoring	0.28	0.00
Base scenario	0.26	2.78
Low monitoring cost	0.26	2.86
High monitoring cost	0.27	1.72
Low penalty	0.27	0.89
High penalty	0.26	2.89
Low inequality aversion	0.16	0.70
High inequality aversion	0.30	3.01

Note: The marginal deadweight loss refers to the income-weighted average of the marginal deadweight loss of all households as a consequence of increasing the labor wedge on labor one percent by 1%. Welfare gains are obtained by calculating the compensating variation as a percentage of average earnings in the specified simulation.

increase in inequality aversion increases the welfare gain of monitoring, because the efficiency gain of monitoring is increasing in the optimal labor wedges, which are larger when the government is more inequality averse. We find quantitatively substantial social welfare gains of monitoring in all scenarios.

## 6. Conclusions

In this paper we demonstrate that redistributive governments should optimally monitor labor hours in order to redistribute income at the lowest efficiency cost. Monitoring of labor supply alleviates the equity–efficiency trade-off and raises equity, efficiency, or both. The reason is that distortions from redistribution derive from the informational problem that earning ability is private information. By using a monitoring technology this informational asymmetry is reduced. A first-best outcome cannot be reached, however, because monitoring is costly. Mirrlees (1971) is a special case of our model when monitoring is prohibitively costly.

We demonstrated that monitoring labor supply works as an implicit subsidy on labor supply, which partially offsets the explicit tax on labor supply. We derived conditions on the desirability of monitoring and demonstrated that the optimal non-linear monitoring schedule generally follows the optimal labor wedge. Monitoring is more desirable when redistributive taxation creates larger distortions in labor supply. Moreover, optimal labor taxes can optimally be above 100% when monitoring is allowed for. At the endpoints of the earnings distribution labor wedges – including taxes and the implicit subsidy on work due to monitoring – are zero in the absence of bunching and with a finite skill level.

Simulations confirmed that the optimal monitoring intensity features a U-shaped pattern with income; very high at the lower end of the earnings distribution, declining towards the middle-income groups, increasing again towards the high-income groups, and becoming constant at the top-income groups. Our simulations demonstrated that marginal tax rates will be higher if the government monitors labor supply, while the labor wedges – including the explicit tax and implicit subsidy of monitoring – decreases. Indeed, monitoring is very effective to alleviate the equity efficiency trade-off.

In practice, monitoring is not infinitely costly as in Mirrlees (1971). By allowing for a monitoring technology we can explain why work-dependent tax credits for low-income earners, that are employed in the UK, Ireland and New Zealand, are part of an optimal redistributive tax policy. Our findings also show that sanctions for welfare recipients, bonuses for low-income workers, and extensive monitoring of labor effort or working ability of low-earning individuals are especially desirable in more generous welfare states. Moreover, we can also explain why (large) penalties on hours worked (or high bonuses on hours worked) are more desirable when the government desires to redistribute more income. Finally, we find that marginal tax rates higher than 100% at the lower end of the earnings distribution, as commonly observed in many countries, can be optimal in the presence of monitoring of labor supply.

In future research, monitoring technologies to verify hidden behaviors of tax payers may be fruitfully applied in other areas of optimal taxation. For example, one can study optimal income taxation and minimum wages as in e.g., *Boadway and Cuff (2001)*. A monitoring technology would allow the government to verify wage rates per hour worked, which is needed to enforce a minimum-wage policy. Similarly, our analysis may be applied to models with an extensive labor-supply margin. Then, the monitoring technology might allow the government to monitor participation costs, rather than wage rates. Doing so alleviates the trade-off between equity and participation distortions. Our analysis could also be applied to generalize the study of *Cremer and Gahvari (1996)* to allow for a

continuum of skill types and study the consequences of tax evasion for the setting of optimal non-linear taxes and non-linear monitoring probabilities.

## Appendix A. Simulation algorithm

The algorithm we use to solve for the optimal allocation consists of two steps. First, we find the optimal allocation using a shooting method. Second, we calculate the implied wedge, tax, and monitoring schedules.<sup>20</sup>

### A.1. Finding the optimal allocation

We find the optimal allocation through five nested loops:

1. The first loop chooses the labor requirement  $l^*$  that maximizes social welfare.
2. The second loop solves the resource constraint (Eq. (12)) for  $\lambda$ . A higher value of  $\lambda$  implies a higher shadow value of resources, and thus, a lower resource deficit, and vice versa. Therefore, we can satisfy the resource constraint arbitrarily by altering the value of  $\lambda$ .
3. The third loop solves the transversality condition at the top (Eq. (32)) for a given utility level at the bottom  $u_n$ , and  $\lambda$ . The most important determinant in  $u_n$  is the transfer implied by  $T(0)$ . Therefore, one can think of this procedure as finding the intercept of the tax function  $T(0)$ . If the intercept is too low, the distortion at the top has to be positive to finance the transfer, and vice versa if the intercept is set too high. As a consequence, by varying the transfer  $T(0)$  we can satisfy the transversality condition arbitrarily closely.
4. The fourth loop solves the differential Eqs. (13) and (30) for given  $u_n$ ,  $\lambda$ , and  $\theta_n$  using a Runge–Kutta method to integrate over  $n$ .
5. The fifth loop maximizes the Lagrangian (Eq. (25)) with respect to  $\pi_n$  and  $z_n$  for a given state  $u_n$  and costate variable  $\theta_n$  at each  $n$ .

The above algorithm is known as a shooting method because it shoots for the initial values of the differential equations that satisfy the boundary condition.

### A.2. Finding the optimal wedge, tax, and monitoring schedules

The above algorithm gives us a numerical approximation of the allocation  $\{u_n, \theta_n, z_n, \pi_n\}$  at each  $n$ .  $\pi'(z_n)$  can be approximated by taking the first difference:

$$\pi'(z_n) \approx \frac{\Delta \pi_n}{\Delta z_n}. \quad (53)$$

With  $\pi'(z_n)$  we have all the information we need to find the optimal tax schedule using Eq. (38).

## References

- Allingham, Micheal G., Sandmo, Agnar, 1972. *Income tax evasion: a theoretical analysis*. *J. Public Econ.* 1 (3–4), 323–338.
- Armenter, Roc, Mertens, Thomas M., 2013. *Fraud deterrence in dynamic Mirrleesian economies*. *J. Monet. Econ.* 60 (2), 139–151.
- Bassetto, Marco, Phelan, Christopher, 2008. *Tax riots*. *Rev. Econ. Stud.* 75 (3), 649–669.

<sup>20</sup> All Matlab programs used in the computations are available from the authors upon request.



- Boadway, Robin, Cuff, Katherine, 1999. Monitoring job search as an instrument for targeting transfers. *Int. Tax Public Financ.* 6 (3), 317–337.
- Boadway, Robin, Cuff, Katherine, 2001. A minimum wage can be welfare-improving and employment-enhancing. *Eur. Econ. Rev.* 45 (3), 553–576.
- Boone, Jan, Fredriksson, Peter, Holmlund, Bertil, Van Ours, Jan C., 2007. Optimal unemployment insurance with monitoring and sanctions. *Econ. J.* 117 (518), 399–421.
- Boone, Jan, Van Ours, Jan C., 2006. Modeling financial incentives to get the unemployed back to work. *J. Inst. Theor. Econ.* 162 (2), 227–252.
- Brewer, Mike, Saez, Emmanuel, Shephard, Andrew, 2010. Means-testing and tax rates on earnings. In: Mirrlees, James A., Adam, Stuart, Besley, Timothy J., Blundell, Richard, Bond, Steven, Chote, Robert, Gammie, Malcolm, Johnson, Paul, Myles, Gareth D., Poterba, James M. (Eds.), *The Mirrlees Review – Dimensions of Tax Design*. Oxford University Press, Oxford, pp. 202–274.
- Chander, Parkash, Wilde, Louis L., 1998. A general characterization of optimal income tax enforcement. *Rev. Econ. Stud.* 65 (1), 165–183.
- Cremer, Helmuth, Gahvari, Firouz, 1994. Tax evasion, concealment and the optimal linear income tax. *Scand. J. Econ.* 2 (96), 219–239.
- Cremer, Helmuth, Gahvari, Firouz, 1996. Tax evasion and the optimum general income tax. *J. Public Econ.* 60 (2), 235–249.
- Diamond, Peter, Sheshinski, Eytan, 1995. Economic aspects of optimal disability benefits. *J. Public Econ.* 57 (1), 1–23.
- Diamond, Peter A., 1998. Optimal income taxation: an example with a U-shaped pattern of optimal marginal tax rates. *Am. Econ. Rev.* 88 (1), 83–95.
- Ebert, Udo, 1992. A reexamination of the optimal non-linear income tax. *J. Public Econ.* 49 (1), 47–73.
- Fredriksson, Peter, Holmlund, Bertil, 2006. Improving incentives in unemployment insurance: a review of recent research. *J. Econ. Surv.* 20 (3), 357–386.
- Gerritsen, Aart, Jacobs, Bas, 2015. Is a minimum wage an appropriate instrument for redistribution? Mimeo: Max Planck Institute Munich/Erasmus University Rotterdam Munich/Rotterdam.
- Harris, Milton, Townsend, Robert M., 1981. Resource allocation under asymmetric information. *Econometrica* 49 (1), 33–64.
- Holmstrom, Bengt, 1979. Moral hazard and observability. *Bell J. Econ.* 10 (1), 74–91.
- Immervoll, Herwig, 2004. Average and marginal effective tax rates facing workers in the EU: A micro-level analysis of levels, distributions and driving factors. OECD Social, Employment and Migration Working Papers, No. 19. OECD Publishing, Paris.
- Jacquet, Laurence, 2014. Tagging and redistributive taxation with imperfect disability monitoring. *Soc. Choice Welf.* 42 (2), 403–435.
- Jacquet, Laurence, Lehmann, Etienne, Van der Linden, Bruno, 2013. Optimal redistributive taxation with both extensive and intensive responses. *J. Econ. Theory* 148 (5), 1770–1805.
- Kleven, Henrik J., Knudsen, Martin B., Kreiner, Claus T., Pedersen, Søren, Saez, Emmanuel, 2011. Unwilling or unable to cheat? Evidence from a tax audit experiment in Denmark. *Econometrica* 79 (3), 651–692.
- Kocherlakota, Narayana R., 2006. Advances in dynamic optimal taxation. In: Blundell, Richard, Newey, Whitney K., Persson, Torsten (Eds.), *Econometric Society Monographs*. Cambridge University Press, Cambridge, pp. 269–297.
- Lee, David, Saez, Emmanuel, 2012. Optimal minimum wage policy in competitive labor markets. *J. Public Econ.* 96 (9–10), 739–749.
- Ljungqvist, Lars, Sargent, Thomas J., 1995. The Swedish unemployment experience. *Eur. Econ. Rev.* 39 (5), 1043–1070.
- Ljungqvist, Lars, Sargent, Thomas J., 1995. Welfare states and unemployment. *Econ. Theory* 6 (1), 143–160.
- Mankiw, N. Gregory, Weinzierl, Matthew, Yagan, Danny, 2009. Optimal taxation in theory and practice. *J. Econ. Perspect.* 23 (4), 147–174.
- Mirrlees, James A., 1971. An exploration in the theory of optimum income taxation. *Rev. Econ. Stud.* 38 (2), 175–208.
- Mirrlees, James A., 1976. Optimal tax theory: a synthesis. *J. Public Econ.* 6 (4), 327–358.
- Mirrlees, James A., 1997. Information and incentives: The economics of carrots and sticks. *Econ. J.* 107 (444), 1311–1329.
- Mirrlees, James A., 1999. The theory of moral hazard and unobservable behaviour: part I. *Rev. Econ. Stud.* 66 (1), 3–21.
- Mookherjee, Dilip, Png, Ivan, 1989. Optimal auditing, insurance, and redistribution. *Q. J. Econ.* 104 (2), 399–415.
- Myerson, Roger B., 1979. Incentive compatibility and the bargaining problem. *Econometrica* 47 (1), 61–73.
- OECD, 2011. Taxation and employment. OECD Tax Policy Studies, No. 21, OECD, Paris.
- Roberts, Kevin, 1984. The theoretical limits of redistribution. *Rev. Econ. Stud.* 51 (2), 177–195.
- Sadka, Efraim, 1976. On income distribution, incentive effects and optimal income taxation. *Rev. Econ. Stud.* 43 (2), 261–267.
- Saez, Emmanuel, 2001. Using elasticities to derive optimal income tax rates. *Rev. Econ. Stud.* 68 (1), 205–229.
- Saez, Emmanuel, Slemrod, Joel B., Giertz, Seth H., 2012. The elasticity of taxable income with respect to marginal tax rates: a critical review. *J. Econ. Lit.* 50 (1), 3–50.
- Sandmo, Agnar, 1981. Income tax evasion, labour supply, and the equity-efficiency tradeoff. *J. Public Econ.* 16 (3), 265–288.
- Schroyen, Fred, 1997. Pareto efficient income taxation under costly monitoring. *J. Public Econ.* 65 (3), 343–366.
- Seade, Jesus K., 1977. On the shape of optimal tax schedules. *J. Public Econ.* 7 (2), 203–235.
- Slemrod, Joel, 1994. Fixing the leak in Okun's bucket: Optimal tax progressivity when avoidance can be controlled. *J. Public Econ.* 55 (1), 41–51.
- Slemrod, Joel, Kopczuk, Wojciech, 2002. The optimal elasticity of taxable income. *J. Public Econ.* 84 (1), 91–112.
- Slemrod, Joel, Yitzhaki, Shlomo, 2002. Tax avoidance, evasion, and administration. In: Auerbach, Alan J., Feldstein, Martin (Eds.), *Handbook of Public Economics*. 3. Elsevier, Amsterdam, pp. 1423–1470.
- Spadaro, Amedeo, 2005. Micro-simulation and normative policy evaluation: an application to some EU tax benefits systems. *J. Public Econ. Theory* 7 (4), 593–622.
- Stern, Nicholas, 1982. Optimum taxation with errors in administration. *J. Public Econ.* 17 (2), 181–211.
- Townsend, Robert M., 1979. Optimal contracts and competitive markets with costly state verification. *J. Econ. Theory* 21 (2), 265–293.
- Tuomala, Matti, 1984. On the optimal income taxation: some further numerical results. *J. Public Econ.* 23, 351–366.
- Zoutman, Floris T., Jacobs, Bas, 2014. Optimal redistribution and monitoring of ability. CES-ifo Working Paper No. 4646. CESifo, Munich.
- Zoutman, Floris T., Jacobs, Bas, Jongen, Egbert L.W., 2015. Optimal redistributive taxes and redistributive preferences in the Netherlands. Rotterdam. Mimeo: Erasmus University Rotterdam/CPB Netherlands Bureau for Economic Policy Analysis.