

# Online Appendix for “Optimal Taxation with Behavioral Agents”

Emmanuel Farhi and Xavier Gabaix

August 22, 2019

Section 9 contains additional results on the paper, in particular on setups with heterogeneous agents, endogenous attention, and mental accounts. Section 10 gives much more detail on the Mirrlees model. Section 11 contains proofs not included in the main paper. Sections 12 and 13 gives complements to consumer theory, with linear and nonlinear budget constraints respectively.

## 9 Additional Results

### 9.1 Complements on optimal tax with heterogeneous agents

#### 9.1.1 Calibration: Optimal Ramsey tax with heterogeneous agents

Here we provide details to the calibration done in Section 3.1

With heterogeneous agents, the misperception is distributed as a 2-point distribution with the following properties:

$$m_i^h = \begin{cases} 1 & \text{with probability } p \\ a & \text{with probability } 1 - p \end{cases}$$

with  $a \in [0, 1]$ , and

$$\begin{aligned} \mathbb{E}[m_i^h] &= p \times 1 + (1 - p) \times a = 0.25 \\ \mathbb{E}[(m_i^h)^2] &= p \times 1 + (1 - p) \times a^2 = 0.25^2 + 0.13. \end{aligned}$$

These equations are satisfied at  $p = .1877$  and  $a = .0767$ . We then take equation (11), with

$$\begin{aligned} S_i^h &= - \frac{c_i^h \psi_i}{q_i^h} \\ q_i^h &= 1 + m_i^h \tau_i \\ c_i^h &= (q_i^h)^{-\psi_i}. \end{aligned}$$

This yields:

$$\frac{\tau_i^*}{p_i} = \frac{(\sum_h \pi_h c_i^h)(1 - \frac{\gamma}{\lambda})}{\psi_i \sum_h \pi_h [m_i^h \frac{c_i^h}{q_i^h} (1 - (1 - m_i^h) \frac{\gamma}{\lambda})]},$$

where  $\pi_h \in \{p, (1-p)\}$  is the fraction of agents of each type. Assume values  $1 - \frac{\gamma}{\lambda} = \Lambda = 1.25\%$  and  $\psi_i = 1$ . Then, under the case with heterogeneity ( $p = .1877$  and  $a = .0767$ ), we have  $\frac{\tau_i^*}{p_i} = 0.073$ , or 7.3%. Under homogeneity with the same average misperception  $m_i^h = .25$  for all agents,  $\frac{\tau_i^*}{p_i} = 20.28$ , for a ratio  $20.28/.073 = 2.78$ . When the taxes are fully salient, so  $m_i^h = 1$  for all agents, then the optimal tax is 1.27%, giving a ratio  $.0127/.073 = .174$ .

### 9.1.2 Optimal taxes with default tax perceptions and heterogeneous agents

Agent  $h$  has utility  $u^h(\mathbf{c}) = c_0^h + \sum_i U^h(c_i^h)$ , with quadratic utility  $U^h(c_i^h) = \frac{a^h c_i^h - \frac{1}{2}(c_i^h)^2}{\Psi^h}$ . Agents are heterogeneous in attention  $m_i^h$  and default taxes  $\tau_i^{d,h}$ . In particular, agent  $h$  perceives tax as  $\tau_i^{s,h} = m_i^h \tau_i + (1 - m_i^h) \tau_i^{d,h}$ . Each agent has the same social welfare weight  $\gamma$ .

The demand for good  $i$  is  $c_i^h(\tau_i) = a^h - \Psi^h(p_i + \tau_i^{s,h}(\tau_i))$ .

The Ramsey planning problem is

$$\max_{\tau} L(\tau)$$

where

$$L(\tau) = \sum_{h=1}^H \gamma \sum_{i=1}^n (U^h(c_i^h(\tau_i)) - (p_i + \tau_i) c_i^h(\tau_i) + \lambda \tau_i c_i^h(\tau_i))$$

First-order condition

$$\frac{\partial L}{\partial \tau_i} = \gamma \sum_{h=1}^H \left( U_{c_i^h}^h \frac{\partial c_i^h}{\partial \tau_i} - (p_i + \tau_i) \frac{\partial c_i^h}{\partial \tau_i} - c_i^h(\tau_i) + \frac{\lambda}{\gamma} c_i^h(\tau_i) + \frac{\lambda}{\gamma} \tau_i \frac{\partial c_i^h}{\partial \tau_i} \right) = 0$$

Let  $\Lambda' \equiv \lambda/\gamma - 1$ , and note that  $\partial c_i^h / \partial \tau_i = \Psi^h m_i^h$ , we can rewrite the FOC as:

$$\begin{aligned} \frac{\partial L}{\partial \tau_i} &= \gamma \sum_{h=1}^H \left( \left( \frac{a^h - c_i^h(\tau_i)}{\Psi^h} - p_i + \Lambda' \tau_i \right) (-\Psi^h m_i^h) + \Lambda' c_i^h(\tau_i) \right) \\ &= \gamma \sum_{h=1}^H \left( \left( \tau_i^{s,h} + \Lambda' \tau_i \right) (-\Psi^h m_i^h) + \Lambda' [a^h - \Psi^h (p_i + \tau_i^{s,h}(\tau_i))] \right) \\ &= \gamma \sum_{h=1}^H \left( -\Psi^h m_i^h (m_i^h + \Lambda') \tau_i - \Psi^h m_i^h (1 - m_i^h) \tau_i^{d,h} + \Lambda' [a^h - \Psi^h p_i] - \Lambda' \Psi^h m_i^h \tau_i - \Lambda' \Psi^h (1 - m_i^h) \tau_i^{d,h} \right) \\ &= \gamma \sum_{h=1}^H \left( -\Psi^h m_i^h (m_i^h + 2\Lambda') \tau_i - \Psi^h (1 - m_i^h) (m_i^h + \Lambda') \tau_i^{d,h} + \Lambda' [a^h - \Psi^h p_i] \right) \\ &= -\gamma H \left( \mathbb{E}[\Psi^h m_i^h (m_i^h + 2\Lambda')] \tau_i + \mathbb{E}[\Psi^h (1 - m_i^h) (m_i^h + \Lambda') \tau_i^{d,h}] - \Lambda' \mathbb{E}[a^h - \Psi^h p_i] \right) = 0 \end{aligned}$$

We can solve explicitly for optimal Ramsey tax in this case:

$$\tau_i = \frac{\Lambda' \mathbb{E}[a^h - \Psi^h p_i] - \mathbb{E}[\Psi^h (1 - m_i^h)(m_i^h + \Lambda') \tau_i^{d,h}]}{\mathbb{E}[\Psi^h m_i^h (m_i^h + 2\Lambda')]} \quad (45)$$

### 9.1.3 Pigouvian Nudges with heterogeneous agents

We start with the Pigouvian example of Section 3.2. There is only one taxed good  $n = 1$ . We use the specialization of the general model developed in Section 2.7. We assume no redistribution or revenue-raising motives ( $\gamma^h = \beta^h = \lambda$ ).

We model the nudge as a psychological tax, as in Section 2.4. Agent  $h$ 's demand is given by  $\arg \max_c U(c) - (p + \eta^h \chi) c$ , where  $\chi$  is the nudge and  $\eta^h$  is the agent's nudgeability. We use quadratic utilities, exactly as in the Pigouvian taxes of Section 3.2. The demand of a consumer can then be expressed as  $c^h(\tau, \chi) = a^h - \Psi(p + \eta^h \chi)$ , where  $\eta^h$  is the nudgeability of agent  $h$ . We apply the optimal nudge formula (9).

When the nudge is the only instrument, the optimal nudge is

$$\chi = \frac{\mathbb{E}[\xi^h \eta^h]}{\mathbb{E}[\eta^{h2}]} = \frac{\mathbb{E}[\xi^h] \mathbb{E}[\eta^h] + \text{cov}(\xi^h, \eta^h)}{\mathbb{E}[\eta^h]^2 + \text{var}[\eta^h]}, \quad (46)$$

where again  $\mathbb{E}$  denotes the average over agents  $h$ .<sup>70</sup>

Heterogeneities in nudgeability determine how well targeted the nudge is to the internality/externality. The optimal nudge is stronger when it is well-targeted, in the sense that nudgeable agents are also those with high internality/externality (higher  $\text{cov}(\xi^h, \eta^h)$ ). The optimal nudge is weaker when there is more heterogeneity in nudgeability (higher  $\text{var}[\eta^{h2}]$ ).<sup>71</sup>

### 9.1.4 Nudges vs. Taxes with Redistributive Concerns

**Jointly optimal nudges and taxes** We next consider the optimal joint policy using both nudges and taxes. We only highlight a few results; more results can be found in the online appendix (Section 9.1.4). We again normalize  $p_i = 1$ . One can show that

$$\frac{\partial^2 L}{\partial \tau \partial \chi} = -\frac{1}{\Psi} \mathbb{E}[(\lambda - \gamma^h (1 - m^h)) \eta^h].$$

As a result, if  $\gamma^h = \lambda$  so that there are no revenue raising or redistributive motives, then taxes and nudges are substitutes. Taxes and nudges are complements if and only if  $\mathbb{E}[(\lambda - \gamma^h (1 - m^h)) \eta^h] \leq 0$ . Nudges and taxes can be complement if social marginal utility of income  $\gamma^h$  and nudgeability

<sup>70</sup>The intermediate steps are as follows. Using  $\mathbf{c}_\chi^h = -\Psi \eta^h$ ,  $\tau = 0$ ,  $\tau^{b,h} = \tau^{X,h} - \chi \eta^h$ , we get  $\frac{\partial L}{\partial \chi}(\tau, \chi) = \sum_h [\lambda \tau - \lambda \tau^{\xi,h} - \beta^h \tau^{b,h}] \cdot \mathbf{c}_\chi^h = \lambda \sum_h [0 - \tau^{X,h} + \chi \eta^h] \Psi \eta^h$ .

<sup>71</sup>Some recent studies study the demographic covariates of nudgeability (Chetty et al. (2014), Beshears et al. (2016)), and it would be good to measure the covariance between nudgeability and internality.

$\eta^h$  are positively correlated. Loosely speaking, if poor agents (with a high  $\gamma^h$ ) are highly nudgable, then taxes and nudges can become complements, because in that case, nudges reduces the consumption of poor nudged agents, thereby improving the redistributive incidence of the tax. We next state the exact values of taxes and nudges, in the case  $\gamma^h = \lambda$ .<sup>72</sup>

**Proposition 9.1** *Assume  $\gamma^h = \lambda$ . Then jointly optimal nudges and taxes are given by the following formulas*

$$\tau = \frac{\mathbb{E}[(\eta^h)^2] \mathbb{E}[\tau^{X,h} m^h] - \mathbb{E}[\eta^h m^h] \mathbb{E}[\tau^{X,h} \eta^h]}{\mathbb{E}[(\eta^h)^2] \mathbb{E}[(m^h)^2] - (\mathbb{E}[\eta^h m^h])^2},$$

$$\chi = \frac{\mathbb{E}[\tau^{X,h} \eta^h] \mathbb{E}[(m^h)^2] - \mathbb{E}[\tau^{X,h} m^h] \mathbb{E}[\eta^h m^h]}{\mathbb{E}[(\eta^h)^2] \mathbb{E}[(m^h)^2] - (\mathbb{E}[\eta^h m^h])^2}.$$

The more powerful the nudge is for high-internality agents (the higher is  $\mathbb{E}[\tau^{X,h} \eta^h]$ , keeping all other moments constant), the more optimal policy relies on the nudge and the less it relies on the tax (the higher is  $\chi$ , the lower is  $\tau$ ). Symmetrically, if the better perceived is the tax by high-internality people (the higher is  $\mathbb{E}[\tau^{X,h} m^h]$ ), the more optimal policy relies on the tax and the less it relies on the nudge.

The more heterogeneity there is in the perception of taxes (the higher is  $\mathbb{E}[(m^h)^2]$ , holding all other moments constant), the less targeted the tax is to the internality/externality, and, as a result, the lower is the optimal tax  $\tau$ , and under certain conditions, the higher the optimal nudge  $\chi$ .<sup>73</sup> Similarly, the more heterogeneity there is in nudgability (the higher is  $\mathbb{E}[(\eta^h)^2]$ , holding all other moments constant), then lower is the optimal nudge  $\chi$ , and, under similar conditions, the higher is the optimal tax  $\tau$ .

**Nudges vs. taxes** We now ask how to choose, if one must, between nudges and taxes. We could analyze this question using the model outlined just above, comparing the relative merits of nudges and taxes in terms of internality targeting and redistributive incidence. Instead, we choose to investigate this question in the context of a model with no heterogeneity, but where the nudges are potentially aversive.

---

<sup>72</sup>In the general case, with the notation  $\sigma_{Y,Z} = cov(Y_h, Z_h)$  :

$$\tau = \frac{\mathbb{E}[\gamma^h \eta^{h^2}] \mathbb{E}[\lambda \tau^{X,h} m^h - \sigma_{\gamma,c/\Psi}] - \mathbb{E}[\gamma^h \eta_h m^h] \mathbb{E}[\lambda \tau^{X,h} \eta^h]}{\mathbb{E}[\gamma^h \eta^{h^2}] \mathbb{E}[\gamma^h m^{h^2} - \sigma_{\gamma m}] - \mathbb{E}[\gamma^h \eta^h m^h] \mathbb{E}[\gamma^h \eta^h m^h - \sigma_{\gamma,\eta}]},$$

$$\chi = \frac{\mathbb{E}[\lambda \tau^{X,h} \eta^h] \mathbb{E}[\gamma^h m^{h^2} - \sigma_{\gamma m}] - \mathbb{E}[\lambda \tau^{X,h} m^h - \sigma_{\gamma,c/\Psi}] \mathbb{E}[\gamma^h \eta^h m^h - \sigma_{\gamma,\eta}]}{\mathbb{E}[\gamma^h \eta^{h^2}] \mathbb{E}[\gamma^h m^{h^2} - \sigma_{\gamma m}] - \mathbb{E}[\gamma^h \eta^h m^h] \mathbb{E}[\gamma^h \eta^h m^h - \sigma_{\gamma,\eta}]}$$

<sup>73</sup>The condition is  $\mathbb{E}[\tau^{X,h} m^h] \mathbb{E}[\eta^{h^2}] \geq \mathbb{E}[\tau^{X,h} \eta^h] \mathbb{E}[\eta^h m^h]$ . It is verified if  $\eta^h, m^h, \tau^{X,h}$  are independent.

We augment the example of Section 3.4 with aversive nudges. We use the same quadratic utility functions as in Section 3.5. We use the nudge as a tax model developed in Section 2.4. We again normalize  $p_i = 1$ .

For concreteness, we interpret the harmful good (good 1) as cigarettes. We extend the model to account for the possibility that the nudge may directly create an aversive reaction (perhaps via a disgusting image of a cancerous lung), which we capture as a separable utility cost  $\iota^h \chi c_i$  so that experienced utility is now

$$u^h(\mathbf{c}, \chi) = u^h(\mathbf{c}) - \iota^h \chi c_i,$$

where  $\iota^h \chi c_i$  is the nudge aversion term. And we assume that there is no heterogeneity across agents.

The next proposition formalizes how nudge aversion changes the relative attractiveness of nudges vs. taxes. The planner must choose between two instruments to discourage cigarette consumption: a weakly positive tax ( $\tau \geq 0$ ) or an aversive nudge ( $\chi \geq 0$ ).

**Proposition 9.2** (“Nudge the poor, tax the rich”) *Consider a good with a “bad” externality (e.g. cigarettes). Suppose that at most one of two instruments (nudges and nonnegative taxes) can be used to correct this externality. And suppose that there is no heterogeneity across agents. Then an optimal tax is superior to an optimal nudge if and only if*

$$\frac{\lambda - \gamma^h}{m^h} > \frac{-\iota^h \gamma^h}{\eta^h}. \quad (47)$$

This proposition captures a new interesting trade-off between taxes and nudges. Both taxes and nudges correct externalities. But taxes also raise revenues on the agents consuming the good under consideration, which is desirable if  $\lambda > \gamma^h$  but undesirable if  $\lambda < \gamma^h$ . Nudges do not raise revenues, and instead directly reduce utility.

When  $\lambda > \gamma^h$ , taxes dominate nudges as taxes have desirable side effects by raising revenues while nudges have adverse side effects by reducing utility. But when  $\lambda < \gamma^h$  taxes and nudges both have undesirable side effects. Taxes dominate nudges when the desire to redistribute income towards agents consuming the good associated with the externality is weak ( $\gamma^h - \lambda$  is low), and when these agents are attentive to the tax ( $m^h$  is high). Nudges dominate taxes when nudge aversion is low ( $\iota^h$  is low) and when agents are easily nudged ( $\eta^h$  is high).

See section 9.5.4 for more details on optimal nudges and taxes.

### 9.1.5 Discouragement formula

In the traditional model without behavioral biases we can use the symmetry of the Slutsky matrix  $\mathbf{S}^{r,h}$  to write  $\boldsymbol{\tau} \cdot \mathbf{S}_i^{r,h} = \sum_j \tau_j \mathbf{S}_{ji}^{r,h}$  as  $\boldsymbol{\tau} \cdot \mathbf{S}_i^{r,h} = \sum_j \tau_j \mathbf{S}_{ij}^{r,h}$ . We can then rewrite the optimal tax

formula of Proposition 2.1 in “discouragement” form as

$$\frac{-\sum_{h,j} \tau_j \mathbf{S}_{ij}^{r,h}}{c_i} = 1 - \frac{\bar{\gamma}}{\lambda} - \text{cov}\left(\frac{\gamma^h}{\lambda}, \frac{Hc_i^h}{c_i}\right), \quad (48)$$

The left-hand side is the discouragement index of good  $i$ , which loosely captures how much the consumption of good  $i$  is discouraged by the taxes  $\tau_j$  on all the different commodities  $j$ . The right-hand side indicates that in the absence of distributive concerns (homogenous  $\gamma^h = \gamma$ ), all goods should be uniformly discouraged in proportion to the relative intensity  $1 - \frac{\bar{\gamma}}{\lambda}$  of the raising revenue objective. With redistributive concerns (heterogenous  $\gamma^h$ ), goods that are disproportionately consumed by agents that society tries to redistribute towards (agents with a high  $\gamma^h$ ) should be discouraged less.

## 9.2 Complements on Endogenous Attention: Attention as a good

### 9.2.1 Interpreting attention as a good

To capture attention and its costs, we propose the following reinterpretation of the general framework. We imagine that we have the decomposition  $\mathbf{c} = (\mathbf{C}, \mathbf{m})$ , where  $\mathbf{C}$  is the vector of traditional goods (champagne, leisure), and  $\mathbf{m}$  is the vector of attention (e.g.  $m_i$  is attention to good  $i$ ). We call  $I^{\mathbf{C}}$  (respectively  $I^{\mathbf{m}}$ ) the set of indices corresponding to traditional goods (respectively attention). Then, all the analyses and propositions apply without modification.

This flexible modeling strategy allows us to capture many potential interesting features of attention. The framework allows (but does not require) attention to be chosen and react endogenously to incentives in a general way (optimally or not). It also allows (but does not require) attention to be produced, purchased and taxed.

We find it most natural to consider the case where attention is not produced, cannot be purchased, and cannot be taxed. This case can be captured in the model by imposing that  $p_i = \tau_i = 0$  for  $i \in I^{\mathbf{m}}$ .

It is useful to consider two benchmarks. The first benchmark is “no attention cost in welfare,” where attention is endogenous (given by a function  $\mathbf{m}(\mathbf{q}, w)$ ) but its cost is assumed not to directly affect welfare so that  $u(\mathbf{C}, \mathbf{m}) = U(\mathbf{C})$ . For instance, as a decision vs. experienced utility generalization of the example of the previous paragraph, we could have  $\mathbf{m}(\mathbf{q}, w) = \arg \max_{\mathbf{m}} u^s(\mathbf{C}(\mathbf{q}, w, \mathbf{m}), \mathbf{m})$ , where  $u^s(\mathbf{C}, \mathbf{m}) = U(\mathbf{C}) - g(\mathbf{m})$ , but still  $u(\mathbf{C}, \mathbf{m}) = U(\mathbf{C})$ . In that view, people use decision heuristics that can respond to incentives, but the cost of those decision heuristics is not counted in the utility function. In this benchmark, we have  $\tau_i^b = 0$  for  $i \in I^{\mathbf{m}}$ .

The second benchmark is “attention cost in welfare”. For simplicity, we outline this case under the extra assumption, which is easy to relax, that attention is allocated optimally. We suppose that there is a primitive choice function  $\mathbf{C}(\mathbf{q}, w, \mathbf{m})$  for traditional goods that depends on attention

$\mathbf{m} = (m_1, \dots, m_A)$  so that  $\mathbf{c}(\mathbf{q}, w, \mathbf{m}) = (\mathbf{C}(\mathbf{q}, w, \mathbf{m}), \mathbf{m})$ .<sup>74</sup> Attention  $\mathbf{m} = \mathbf{m}(\mathbf{q}, w)$  is then chosen to maximize  $u(\mathbf{C}(\mathbf{q}, w, \mathbf{m}), \mathbf{m})$ . This generates a function  $\mathbf{c}(\mathbf{q}, w) = (\mathbf{C}(\mathbf{q}, w, \mathbf{m}(\mathbf{q}, w)), \mathbf{m}(\mathbf{q}, w))$ . In that benchmark, attention costs are incorporated in welfare.<sup>75</sup> For instance we might consider a separable utility function  $u(\mathbf{C}, \mathbf{m}) = U(\mathbf{C}) - g(\mathbf{m})$  for some cost function  $g(\mathbf{m})$ . A non-separable  $u$  might capture that attention is affected by consumption (e.g., of coffee) and attention affects consumption (by needing aspirin).

The tax formula (7) has a term  $(\boldsymbol{\tau} - \tilde{\boldsymbol{\tau}}^{b,h}) \cdot \mathbf{S}_i^{C,h} = \sum_{k \in I^m \cup I^c} (\tau_k - \tilde{\tau}_k^{b,h}) \mathbf{S}_{ki}^{C,h}$ , a sum that includes the ‘‘attention’’ goods  $k \in I^m$ . As attention is assumed to have zero tax, we have  $\tau_k = 0$  for  $k \in I^m$ . The term  $\tilde{\tau}_k^{b,h}$ , which accounts for potential misoptimization in the allocation of attention, requires no special treatment. However, two polar special cases are worth considering that simplify the calculations. First, consider the ‘‘no attention cost in welfare’’ case. In this case we saw that  $\tilde{\tau}_k^{b,h} = 0$  for  $k \in I^m$ . Together with  $\tau_k = 0$  for  $k \in I^m$ , this implies that  $(\boldsymbol{\tau} - \tilde{\boldsymbol{\tau}}^{b,h}) \cdot \mathbf{S}_i^{C,h} = \sum_{k \in I^c} (\tau_k - \tilde{\tau}_k^{b,h}) \mathbf{S}_{ki}^{C,h}$  is the sum restricted to commodities. Second, consider the ‘‘optimally allocated attention’’ case. Then (see Proposition 9.4)  $(\boldsymbol{\tau} - \tilde{\boldsymbol{\tau}}^{b,h}) \cdot \mathbf{S}_i^{C,h} = \sum_{k \in I^c} (\tau_k \mathbf{S}_{ki}^{C,h} - \tilde{\tau}_k^{b,h} \mathbf{S}_{ki|\mathbf{m}}^{C,h})$ , where  $\mathbf{S}_{i|\mathbf{m}}^{C,h}$  is a Slutsky matrix holding attention constant, which is in general different from  $\mathbf{S}_i^{C,h}$ . For tax revenues, the full Slutsky matrix, including changes in attention, matters (the term  $\tau_k \mathbf{S}_{ki}^{C,h}$ ). However, for welfare, when attention is assumed to be optimally allocated, it is the Slutsky matrix holding attention constant that matters (the term  $\tilde{\tau}_k^{b,h} \mathbf{S}_{ki|\mathbf{m}}^{C,h}$ ). This is a version of the envelope theorem.

**Characterizing optimal allocation of attention** Suppose that we have a constraint:  $\mathbf{c} = \mathbf{c}(\mathbf{p}, w, \theta)$  for some parameter  $\theta$ . For instance, suppose that  $\mathbf{c}(\mathbf{p}, w, \theta) = (\mathbf{C}(\mathbf{p}, w, \mathbf{m}(\theta)), \mathbf{m}(\theta))$ ; when  $\mathbf{m}(\theta) = \theta$ , we’re considering the potentially optimal allocation of attention, as attention affects directly the choice of goods. If  $\mathbf{m} = (m_1, m_2, m_3) = (\theta_1, \theta_2, \theta_2)$ , we capture that the attention to goods 2 and 3 have to be the same.<sup>76</sup>

**Proposition 9.3** (Characterizing optimal allocation of attention) *The first order condition for the optimal allocation of parameter  $\theta$  (i.e.,  $\theta(\mathbf{p}, w) = \arg \max_{\theta} u(\mathbf{c}(\mathbf{p}, w, \theta))$ ) is:*

$$\boldsymbol{\tau}^b \cdot \mathbf{c}_{\theta}(\mathbf{p}, w, \theta) = 0. \quad (49)$$

<sup>74</sup>For instance, in a misperception model, attention operates by changing the perceived price  $\mathbf{q}^s(\mathbf{q}, w, \mathbf{m})$  which in turn changes consumption as  $\mathbf{C}(\mathbf{q}, w, \mathbf{m}) = \mathbf{C}^s(\mathbf{q}, \mathbf{q}^s(\mathbf{q}, w, \mathbf{m}), w)$ .

<sup>75</sup>The first order condition characterizing the optimal allocation of attention can be written as  $\boldsymbol{\tau}^b \cdot \mathbf{c}_{m_j}(\mathbf{q}, w, \mathbf{m}) = 0$  for all  $j \in \{1, \dots, A\}$ . This condition can be re-expressed more conveniently by introducing the following notation: we call  $k(i)$  the index  $k \in I^m$  corresponding to dimension  $i \in \{1, \dots, A\}$  of attention. We then get  $\sum_{i \in I^c} \tau_i^b \mathbf{C}_{m_j}(\mathbf{q}, w, \mathbf{m}) + \tau_{k(j)}^b = 0$  for all  $j \in \{1, \dots, A\}$ .

<sup>76</sup>In a model of noisy decision-making à la Sims (2003), the same logic exactly applies, except that quantities are generally stochastic. The consumption is a random variable  $c(p, w, \tilde{\varepsilon})$ , where  $\tilde{\varepsilon}$  indexes noise, rather than a deterministic function. Then, utility is  $U(c(p, w)) = \mathbb{E}[u(c(p, w, \tilde{\varepsilon}))]$ ,  $\mathbf{S}^H(p, w)$  is likewise a random variable. We do not pursue this framework further here, at it is hard to solve beyond linear-quadratic settings, e.g. with Gaussian distribution of prices – which in turn generates potentially negative prices.

**Proof** The FOC is  $u_c \mathbf{c}_\theta = 0$ . We note that  $B_c \cdot \mathbf{c}_\theta = 0$  by budget constraint:  $B(\mathbf{c}(\mathbf{p}, w, \theta)) = w$ . So,

$$\boldsymbol{\tau}^b \cdot \mathbf{c}_\theta = \left( B_c - \frac{u_c(\mathbf{c}, \mathbf{p})}{v_w(\mathbf{p}, w)} \right) \cdot \mathbf{c}_\theta = -\frac{u_c(\mathbf{c}, \mathbf{p}) \cdot \mathbf{c}_\theta}{v_w(\mathbf{p}, w)},$$

so that  $\boldsymbol{\tau}^b \cdot \mathbf{c}_\theta = 0$  if and only if  $u_c \cdot \mathbf{c}_\theta = 0$ .  $\square$

**Proposition 9.4** (Value of  $D_j$  when attention is optimal). *When attention is of the form  $\mathbf{c}(\mathbf{p}, w, \theta) = (\mathbf{C}(\mathbf{p}, w, \mathbf{m}(\theta)), \mathbf{m}(\theta))$ , and is optimally chosen, then*

$$\begin{aligned} -D_j &= \boldsymbol{\tau}_C^b \cdot \mathbf{C}_{p_j}(\mathbf{p}, w, \mathbf{m})|_{\mathbf{m}=\mathbf{m}(\theta(\mathbf{p}, w))} \\ &= \boldsymbol{\tau}_C^b \cdot \mathbf{S}_{j|\mathbf{m}}^H(\mathbf{p}, w, \mathbf{m})|_{\mathbf{m}=\mathbf{m}(\theta(\mathbf{p}, w))} = \boldsymbol{\tau}_C^b \cdot \mathbf{S}_{j|\mathbf{m}}^C(\mathbf{p}, w, \mathbf{m})|_{\mathbf{m}=\mathbf{m}(\theta(\mathbf{p}, w))}. \end{aligned}$$

where  $\boldsymbol{\tau}_C^b = B_C(\mathbf{C}, \mathbf{p}) - \frac{u_C(\mathbf{C}, \mathbf{m})}{v_w(\mathbf{p}, w)}$  is the behavioral wedge restricted to goods consumption, and  $\mathbf{S}_{j|\mathbf{m}}^H$  and  $\mathbf{S}_{j|\mathbf{m}}^C$  are the Slutsky matrices  $\mathbf{S}_j^H$  and  $\mathbf{S}_j^C$  holding attention constant, i.e. associated to decision  $\mathbf{C}(\mathbf{p}, w, \mathbf{m})$  with constant  $\mathbf{m} = \mathbf{m}(\theta(\mathbf{p}, w))$ .

**Proof** We have

$$\begin{aligned} -D_j &= \boldsymbol{\tau}^b \cdot \mathbf{c}_{p_j}(\mathbf{p}, w, \theta) = \boldsymbol{\tau}^b \cdot [(\mathbf{C}_{p_j}(\mathbf{p}, w, \mathbf{m}), \mathbf{0}) + \mathbf{c}_\theta(\mathbf{p}, w, \theta) \theta_{p_j}(\mathbf{p}, w)] \\ &= \boldsymbol{\tau}^b \cdot (\mathbf{C}_{p_j}(\mathbf{p}, w, \mathbf{m}), \mathbf{0}) \text{ as } \boldsymbol{\tau}^b \cdot \mathbf{c}_\theta = 0 \\ &= (\boldsymbol{\tau}_c^b, \boldsymbol{\tau}_m^b) \cdot (\mathbf{C}_{p_j}(\mathbf{p}, w, \mathbf{m}), \mathbf{0}) = \boldsymbol{\tau}_c^b \cdot \mathbf{C}_{p_j}(\mathbf{p}, w, \mathbf{m}) = \boldsymbol{\tau}_c^b \cdot \mathbf{S}_{j|\mathbf{m}}^C = \boldsymbol{\tau}_c^b \cdot \mathbf{S}_{j|\mathbf{m}}^H. \end{aligned}$$

**“No attention cost in welfare” benchmark** Another benchmark is the “no attention cost in welfare”, i.e. the cost of attention is not taken into account in the welfare analysis. Suppose that attention  $\mathbf{m}$  just moves with prices, but as an automatic process whose “cost” is not counted: that is,  $u(\mathbf{C}, \mathbf{m}) = u(\mathbf{C})$  and attention has 0 price,  $\mathbf{p}_m = 0$ . This is the way it is often done in behavioral economic (see however [Bernheim and Rangel \(2009\)](#)): people choose using heuristics, but the “cognitive cost” associated with a decision procedure isn’t taken into account in the agent’s welfare (largely, because it is very hard to measure, and that revealed preference techniques do not apply).

**Proposition 9.5** (Value of  $D_j$  in the case of fixed attention, and the case of “No attention cost in welfare”). *In the “fixed attention” case and the “No attention cost in welfare” case*

$$-D_j = (\boldsymbol{\tau}_C^b, \mathbf{0}) \cdot \mathbf{S}_j^H(\mathbf{p}, w) = \sum_{i=1}^n \boldsymbol{\tau}_{C_i}^b \mathbf{S}_{ij}^H = (\boldsymbol{\tau}_C^b, \mathbf{0}) \cdot \mathbf{S}_j^C(\mathbf{p}, w) = (\boldsymbol{\tau}_C^b, \mathbf{0}) \cdot \mathbf{c}_j(\mathbf{p}, w).$$

*This is, only the components of  $\boldsymbol{\tau}^b$  and the Slutsky matrix linked to commodities matter.*

**Proof**



We have  $\boldsymbol{\tau}^b = (\boldsymbol{\tau}_c^b, \boldsymbol{\tau}_m^b) = (\boldsymbol{\tau}_c^b, 0)$  as  $u_m = 0$ . So,  $-D_j = \boldsymbol{\tau}^b \cdot \mathbf{S}_j^H(\mathbf{p}, w) = \boldsymbol{\tau}_C^b \cdot \mathbf{S}_{C_j}^H$ .  $\square$

**Misperception example** In the misperception model with attention policy  $\mathbf{m}(\mathbf{p}, w)$ , we have:

$$\mathbf{c}(\mathbf{p}, w) = (\mathbf{C}^s[\mathbf{p}, \mathbf{p}^s(\mathbf{p}, w, \mathbf{m}(\mathbf{p}, w)), v(\mathbf{p}, w)], \mathbf{m}(\mathbf{p}, w)).$$

When attention is optimally chosen, we can apply Proposition 9.4 with  $\mathbf{m}(\theta) = \theta$ . This gives:  $-D_j = \boldsymbol{\tau}_C^b \cdot \mathbf{S}_{C_j}^{H, \mathbf{m}}$  with

$$\mathbf{S}_{C_j | \mathbf{m}}^H = \mathbf{S}^r \mathbf{p}_{p_j}^s(\mathbf{p}, w, \mathbf{m}), \quad (50)$$

i.e. the Slutsky matrix has the sensitivity with *fixed* attention. Hence, we have both  $-D_j = \boldsymbol{\tau}_C^b \cdot \mathbf{S}_j^{H, \mathbf{m}}$  when attention is optimal.

In the “no attention cost in welfare” case,  $\boldsymbol{\tau}_m^b = 0$  and

$$-D_j = \boldsymbol{\tau}_C^b \cdot \mathbf{S}_{C_j}^H.$$

When attention is not necessarily optimal, we also have (from (43)), using again decomposition  $\boldsymbol{\tau}^b = (\boldsymbol{\tau}_c^b, \boldsymbol{\tau}_m^b)$ :

$$-D_j = \boldsymbol{\tau}^b \cdot \mathbf{S}_j = \boldsymbol{\tau}_C^b \cdot \mathbf{S}_{C_j}^H + \boldsymbol{\tau}_m^b \frac{\partial \mathbf{m}}{\partial p_j},$$

where  $\mathbf{S}_{C_j}^H = \mathbf{S}^r \cdot \mathbf{p}_{p_j}^s(\mathbf{p}, w)$ , where now the total derivative matters, including the variable attention.

### 9.2.2 Attention as a good: examples

In this subsection we normalize the pre-tax price to 1.

**Optimal taxes with endogenous attention: the case of small taxes** Given attention  $m(\tau)$ , the perceived tax is  $\tau^s(\tau) = \tau m(\tau)$ , and demand is  $c(\tau) = y(1 - \psi m(\tau)\tau)$ . We assume that attention comes from an optimal cost-benefit analysis:

$$m(\tau) = \arg \max_m -\frac{1}{2}\psi y \tau^2 (1 - m)^2 - g(m).$$

The first term represents the private costs of misunderstanding taxes,  $-\frac{1}{2}\psi y (\tau - \tau^s)^2$ , while the term  $-g(m)$  is the psychic cost of attention,  $g(m)$  (see Gabaix (2014)). The planner’s problem is  $\max_\tau L(\tau)$  with

$$L(\tau) = -\frac{1}{2}\psi y m^2(\tau) \tau^2 - A g(m(\tau)) + \Lambda \tau y,$$

where  $A = 1$  in the “optimally allocated attention” case and  $A = 0$  in the “no attention cost in welfare” case. In the “fixed attention” case,  $m(\tau)$  is fixed with  $m'(\tau) = 0$ , and  $g(m) = 0$ . The

optimal tax satisfies

$$L'(\tau) = -\psi y m(\tau) \tau (m(\tau) + \tau m'(\tau)) - A g'(m(\tau)) m'(\tau) + \Lambda y = 0.$$

In the “optimally allocated attention” case, we use the agent’s first order condition  $g'(m(\tau)) = \psi y \tau^2 (1 - m(\tau))$  and  $A = 1$ , and the optimal tax is

$$\tau^{m,*} = \frac{\Lambda/\psi}{m(\tau)^2 + \tau m'(\tau)}. \quad (51)$$

In the “no attention cost in welfare case,”  $A = 0$ , the optimal tax is

$$\tau^{m,0} = \frac{\Lambda/\psi}{m(\tau)^2 + \tau m(\tau) m'(\tau)}. \quad (52)$$

When attention is fixed, the optimal tax is

$$\tau^{m,F} = \frac{\Lambda/\psi}{m(\tau)^2}. \quad (53)$$

**Proposition 9.6** *In the interior region where attention has an increasing cost ( $\tau m(\tau) m'(\tau) > 0$ ), the optimal tax is lowest when attention is chosen optimally and its cost is taken into account in welfare; intermediate in the “no attention cost in welfare” case; and largest with fixed attention— $\tau^{m,*} < \tau^{m,0} < \tau^{m,F}$ .*

When attention’s cost is taken into account, the planner chooses lower taxes  $\tau^{m,*} < \tau^{m,0}$  to minimize both consumption distortions and attention costs.<sup>77</sup> Plainly, the tax is higher when attention is variable than when attention is fixed—this is basically because demand is more elastic then ( $-\frac{p}{c} \frac{\partial c}{\partial \tau} = -\psi (m(\tau) + \tau m'(\tau))$ ).

For more illustrations, see section 9.5.5 for completely worked out linear-quadratic and isoelastic examples.

### 9.3 Complements on Mental Accounts

The optimal tax formulas in Propositions 2.1 and 2.2 corresponding to the many-person Ramsey problem without and with externalities can be applied without modifications to this simple model of mental accounting. However, it is also enlightening to write these formulas in a slightly different way by leveraging the specific structure of the simple mental accounting model. We define the

<sup>77</sup>The example allows to appreciate the Slutsky matrix with or without constant attention. The Slutsky matrix with constant  $m$  has  $S_{11|m}^C = \frac{\partial c(1+\tau, m)}{\partial \tau} = -\psi c m$ , while the Slutsky matrix with variable  $m$  has  $S_{11}^C = \frac{dc(1+\tau, m(\tau))}{d\tau} = -\psi c (m + \tau m'(\tau))$ . The online appendix (section 9.2.2) provides other illustrations.

“income  $k$ -compensated” Slutsky matrix for the extended demand function as

$$S_j^{C,k}(\mathbf{q}, \boldsymbol{\omega}) = \mathbf{c}_{q_j}(\mathbf{q}, \boldsymbol{\omega}) + \mathbf{c}_{\omega^k}(\mathbf{q}, \boldsymbol{\omega})c_j(\mathbf{q}, \boldsymbol{\omega}). \quad (54)$$

This Slutsky matrix corresponds to a decomposition of price effects into income of substitution effects where the latter are compensated using with an adjustment of mental account  $k$ . In the traditional model without behavioral biases, this decomposition is independent of the mental account  $k$  which is used for this decomposition, since the marginal utility of income is equalized across all accounts:  $\mathbf{c}_{\omega^k}(\mathbf{q}, \boldsymbol{\omega}^r(\mathbf{q}, w)) = \mathbf{c}_{\omega^{k'}}(\mathbf{q}, \boldsymbol{\omega}^r(\mathbf{q}, w))$ . It follows that the “income  $k$ -compensated” Slutsky matrix is also independent of  $k$ .<sup>78</sup> By contrast, with behavioral biases in mental accounting, the marginal utility of income is not equalized across all accounts so that in general  $\mathbf{c}_{\omega^k}(\mathbf{q}, \boldsymbol{\omega}(\mathbf{q}, w)) \neq \mathbf{c}_{\omega^{k'}}(\mathbf{q}, \boldsymbol{\omega}(\mathbf{q}, w))$ . As a result, the decomposition of price effects on income effects and substitution effects depends on which mental account  $k$  is used for this income compensation, and the “income  $k$ -compensated” Slutsky matrix depends on  $k$ .<sup>79</sup>

Reintroducing  $h$  superscripts to index agent heterogeneity, we define the social marginal utility of  $k$ -income for agent  $h$  as

$$\gamma^{k,h} = \beta^{k,h} + \lambda \boldsymbol{\tau} \cdot \mathbf{c}_{\omega^{k,h}}^h \quad \text{where} \quad \beta^{k,h} = W_{v^h} v_{\omega^{k,h}}^h.$$

We also define the income- $k$  based behavioral wedges for the extended demand and utility function as

$$\boldsymbol{\tau}^{b,k} = \mathbf{q} - \frac{u_{\mathbf{c}}}{v_{\omega^k}}, \quad \tilde{\boldsymbol{\tau}}^{b,k,h} = \beta^{k,h} \boldsymbol{\tau}^{b,k}.$$

Finally, for every commodity  $i$ , we denote by  $k(i)$  the mental account to which this commodity is associated with. We can then rewrite the tax formula in the following way. Note that this is simply a re-expression of Proposition 2.1.

**Proposition 9.7** (Many-person Ramsey with mental accounting) *If commodity  $i$  can be taxed, then*

---

<sup>78</sup>However the “income  $k$ -compensated” Slutsky matrix  $S_j^{C,k,r}(\mathbf{q}, \boldsymbol{\omega}^r(\mathbf{q}, w))$  of the extended demand function is in general different from the “income compensated” Slutsky matrix  $S_j^{C,r}(\mathbf{q}, w)$  of the demand function, which is defined as in Section 2.1. Indeed, the latter also reflects the substitution effects associated with the adjustments  $\omega_{q_j}^{k,r}(\mathbf{q}, w)$  in the mental accounts in response to changes in the price  $q_j$  of commodity  $j$ .

<sup>79</sup>The price theory concepts introduced in Section 2.1 are still defined in the same way. They can be related to the corresponding concepts that we have introduced in this section. In particular, we have

$$c_w(\mathbf{q}, w) = \sum_k \omega_w^k \mathbf{c}_{\omega^k}(\mathbf{q}, \boldsymbol{\omega}),$$

and

$$S_j^C(\mathbf{q}, w) = \mathbf{c}_{q_j}(\mathbf{q}, \boldsymbol{\omega}) + \sum_k \omega_{q_j}^k \mathbf{c}_{\omega^k}(\mathbf{q}, \boldsymbol{\omega}) + \sum_k \omega_w^k \mathbf{c}_{\omega^k}(\mathbf{q}, \boldsymbol{\omega})c_j(\mathbf{q}, \boldsymbol{\omega}).$$

at the optimum

$$\frac{\partial L(\boldsymbol{\tau})}{\partial \tau_i} = 0 \quad \text{with} \quad \frac{\partial L(\boldsymbol{\tau})}{\partial \tau_i} = \sum_h [(\lambda - \gamma^{k(i),h}) c_i^h + \lambda(\boldsymbol{\tau} - \tilde{\boldsymbol{\tau}}^{b,k(i),h}) \cdot \mathbf{S}_i^{C,k(i),h} + \sum_k \gamma^{k,h} \omega_{q_i}^{k,h}]. \quad (55)$$

This alternative expression of the many-person Ramsey optimal tax formula for commodity  $i$  features the “income  $k(i)$ -compensated” Slutsky matrix corresponding the mental account to which commodity  $i$  is associated, the social marginal utilities of  $k$ -income  $\gamma^{k,h}$ , the  $k(i)$  based behavioral wedges, and the price derivatives of the mental accounting functions  $\omega_{q_i}^{k,h}$ . Writing the optimal tax formula in this way will prove useful to derive specific results below in the context of further specializations of the model.

We could also derive a similar alternative expression for the many-person Ramsey optimal tax formula in the presence of externalities along very similar lines. In the interest of space, we do not include it in the paper.

### 9.3.1 Roy’s identity with mental accounts

We consider the extended indirect utility function  $v(\mathbf{p}, \boldsymbol{\omega}) = u(\mathbf{c}(\mathbf{p}, \boldsymbol{\omega}))$ . The budget constraint is  $B(\mathbf{c}, \mathbf{p}, \boldsymbol{\omega}) \leq 0$ . A leading case is the linear budget constraint,  $B(\mathbf{c}, \mathbf{p}, \boldsymbol{\omega}) = \max_k \mathbf{C}^k \cdot \mathbf{p}^k - \omega^k$ . We define the behavioral wedge linked to account  $k$  as:

$$\boldsymbol{\tau}^{b,k} = -\frac{u_{\mathbf{c}}}{v_{\omega^k}} - \frac{B_{\mathbf{c}}(\mathbf{c}, \mathbf{p}, \boldsymbol{\omega})}{B_{\omega^k}(\mathbf{c}, \mathbf{p}, \boldsymbol{\omega})}.$$

With the linear budget constraint

$$\boldsymbol{\tau}^{b,k} = \mathbf{p} - \frac{u_{\mathbf{c}}}{v_{\omega^k}}.$$

**Proposition 9.8** (Roy’s identity with mental accounts) *With mental account, the modified Roy’s identity is:*

$$\frac{v_{p_i}(\mathbf{p}, \boldsymbol{\omega})}{v_{\omega^k}(\mathbf{p}, \boldsymbol{\omega})} = \frac{B_{p_i}}{B_{\omega^k}} - \boldsymbol{\tau}^{b,k} \cdot \mathbf{c}_{p_i} = \frac{B_{p_i}}{B_{\omega^k}} - \boldsymbol{\tau}^{b,k} \cdot \mathbf{S}_i^{C,k}. \quad (56)$$

With a linear budget constraint,

$$\frac{v_{p_i}(\mathbf{p}, \boldsymbol{\omega})}{v_{\omega^k}(\mathbf{p}, \boldsymbol{\omega})} = -c_i - \boldsymbol{\tau}^{b,k} \cdot \mathbf{S}_i^{C,k}. \quad (57)$$

### 9.3.2 Optimal taxes with rigid mental accounts: small taxes case

We consider the basic setup in Section 3.1 with no misperceptions ( $m_i = 1$  for all  $i$ ) but with rigid mental accounts instead. We make the further simplification that there is one commodity per mental account. Consumption is therefore given by  $c_i = \frac{\omega^i}{q_i} = \frac{\omega^i}{1+\tau_i}$ . We assume that before

taxes, the optimal amount  $\omega^i$  is allocated to good  $i$ , so that  $U^{i'}(\omega^i) = p_i$ , and that the rigid mental account  $\omega^i$  does not adjust after the introduction of taxes.

We first derive the optimal Ramsey and Pigou tax rules with this rigid mental account with one good per account. Recall that we denote by  $\psi_i = -\frac{U^{i''}(c_i)}{c_i U^{i'}(c_i)}$  the inverse of the curvature of the utility function  $U^i$  for good  $i$ , which coincides with the demand elasticity of a rational agent.

**Proposition 9.9** (Ramsey and Pigou formulas with rigid mental accounts) *Suppose that agents use a rigid mental account for good  $i$ . and the limit of small taxes. In the basic Ramsey problem, the optimal tax is*

$$\frac{\tau_i}{p_i} = \Lambda \psi_i, \quad (58)$$

while in the basic Pigou problem, it is

$$\tau_i = \xi_i \psi_i. \quad (59)$$

The formula for the Ramsey problem is in stark contrast with the traditional Ramsey case where  $\frac{\tau_i}{p_i} = \frac{\Lambda}{\psi_i}$ , and the misperception case where  $\frac{\tau_i}{p_i} = \frac{\Lambda}{m_i^2 \psi_i}$ . With rigid mental accounts, a low (rational) elasticity  $\psi_i$  leads to low taxes, not to high taxes, as in the basic Ramsey. The intuition is as follows: if a good is very “necessary”, rational demand is very inelastic:  $\psi_i$  is low. But with a rigid mental accounts, a tax  $\tau_i$  leads to a consumption  $c_i = \frac{\omega^i/p_i}{1+\tau_i}$ . So, a high tax leads to a high distortion. Hence, when (rational) demand is very inelastic, the tax should be low.

Likewise, the modified Pigou formula  $\tau_i = \xi_i \psi_i$  now features the rational elasticity of demand  $\psi_i$ . This is in contrast to the traditional case, where  $\tau_i = \xi_i$ , and to the case with misperception  $m_i$  where  $\tau_i = \frac{\xi_i}{m_i}$  (Proposition 3.2).

To derive this result and understand it fully, it is useful to generalize it. From now on, in this subsection, we normalize  $p_i = 1$ . We denote by  $\alpha_i$  the elasticity of the demand for good  $i$ . In the traditional model without behavioral biases, we have  $\alpha_i = \psi_i$ . But in the model with attention  $m_i$  to the tax, we had  $\alpha_i = m_i \psi_i$ . With a rigid mental account for commodity  $i$ , given demand is  $c_i = \frac{\omega^i}{1+\tau_i}$ , the elasticity of the demand for good  $i$  is  $\alpha_i = 1$ .<sup>80</sup>

**Proposition 9.10** (Ramsey and Pigou formulas with arbitrary behavioral elasticity) *Suppose that the rational demand elasticity for good  $i$  is  $\psi_i$ , and that the behavioral demand elasticity is  $\alpha_i$ . Consider the limit of small taxes. Then, in the basic Ramsey problem, the optimal tax is*

$$\tau_i = \Lambda \frac{\psi_i}{\alpha_i^2}, \quad (60)$$

while in the basic Pigou problem, it is

$$\tau_i = \frac{\xi_i}{\alpha_i} \psi_i. \quad (61)$$

---

<sup>80</sup>Proposition 9.9 is a consequence of Proposition 9.10 when  $\alpha_i = 1$  of the following result. Propositions 3.1 (in the limit of small taxes) and 3.2 are also an application, when  $\alpha_i = m_i \psi_i$ .

**Proof.** We could use the general formulas, but to gain intuition we proceed as follows, in the limit of small taxes. In the Ramsey problem, welfare can then be expressed as

$$L = -\frac{1}{2} \sum_i \frac{\alpha_i^2}{\psi_i} y_i \tau_i^2 + \Lambda \sum_i \tau_i y_i, \quad (62)$$

Indeed, a small tax  $\tau_i$  changes consumption by  $\delta c_i = -\alpha_i c_i \tau_i$ . The associated distortion is  $\frac{1}{2} (\delta c_i)^2 U^{i''}(c_i) = \frac{1}{2} (-\alpha_i c_i \tau_i)^2 \frac{-U^{i''}(c_i)}{c_i \psi_i} = \frac{-1}{2} \frac{\alpha_i^2}{\psi_i} y_i \tau_i^2$  (recall that  $\psi_i = -\frac{U^{i''}(c_i)}{c_i U^{i''}(c_i)}$ , and  $U^{i'} = p_i = 1$  at the optimum, with  $y_i = c_i$ ). Hence, the optimal tax is given by  $L_{\tau_i} = 0$ , i.e.  $\tau_i = \Lambda \frac{\psi_i}{\alpha_i^2}$ .

In the Pigou problem, at the first best, the planner would like  $U^{i'}(c_i) = 1 + \xi_i$ , as in the traditional tax. This means that consumption should change by  $\delta c_i = -\psi_i \xi_i$  after the tax. But as the actual elasticity of demand is  $\alpha_i$ , the tax should satisfy:  $\delta c_i = -\alpha_i \tau_i = -\psi_i \xi_i$ , and  $\tau_i = \frac{\xi_i}{\alpha_i} \psi_i$ .  $\square$

In the Ramsey problem, for a given demand elasticity  $\alpha_i$ , a higher value of  $\psi_i$  pushes for higher tax, while for a given  $\psi_i$ , a higher value of  $\alpha_i$  pushes for a lower tax. In the traditional model without behavioral biases,  $\alpha_i = \psi_i$  and the resulting effect of a higher  $\psi_i$  is a lower tax. By contrast, in the behavioral model with a rigid mental account,  $\alpha_i = 1$  so that a higher  $\psi_i$  results in a higher tax.

### 9.3.3 How mental accounts modify demand elasticities

We take the quasilinear case  $u(\mathbf{c}) = c_0 + U^1(c_1) + U(c_2, \dots, c_n)$  with good 1 in its own mental account,  $\omega^1$ , and default  $\omega_1^d$ . How much will be attributed to the mental account? We will have  $c_1 = \frac{\omega^1}{q_1}$

$$\begin{aligned} \omega^1 &= \arg \max_{\omega^1} U^1\left(\frac{\omega^1}{q_1}\right) - \omega^1 - g(\omega^1 - \omega_1^d) \\ c_1 &= \arg \max_{c_1} U^1(c_1) - q_1 c_1 - g(c_1 q_1 - \omega_1^d). \end{aligned} \quad (63)$$

Then, we can calculate the sensitivity to the tax.

**Lemma 9.1** *With a flexible mental account, the empirical elasticity is:*

$$\alpha_1 = -\frac{q_1}{c_1} \frac{\partial c_1}{\partial q_1} = \frac{1 + g'(\omega^1 - \omega_1^d) + \omega^1 g''(\omega^1 - \omega_1^d)}{\frac{1}{\psi_1} + \omega^1 g''(\omega^1 - \omega_1^d)},$$

with  $\omega^1 = q_1 c_1$ .

**Proof** The first order condition for consumption is:

$$f(c_1, q_1) = U^{1'}(c_1) - q_1 - q_1 g'(c_1 q_1 - \omega_1^d) = 0.$$

Hence

$$\begin{aligned}\alpha_1 &= -\frac{q_1}{c_1} \frac{\partial c_1}{\partial q_1} = -\frac{q_1}{c_1} \frac{f_{q_1}}{-f_{c_1}} = \frac{q_1}{c_1} \frac{1 + g' + q_1 c_1 g''}{-U'' + q_1^2 g''} \\ &= \frac{1 + g' + q_1 c_1 g''}{-\frac{c_1 U''}{q_1} + q_1 c_1 g''}.\end{aligned}$$

The rational elasticity is the one that would occur with  $g = 0$ ,

$$\psi_1 = \frac{1}{-\frac{c_1 U''(c_1)}{q_1}}, \quad (64)$$

so

$$\alpha_1 = \frac{1 + g' + q_1 c_1 g''}{\frac{1}{\psi_1} + q_1 c_1 g''}.$$

□

Next, we suppose that at  $q_1 = p_1$ , the account is optimal:  $\omega_1^d = \arg \max_{\omega_1} U^1\left(\frac{\omega_1}{q_1}\right) - \omega^1$ . We suppose that we are near  $\omega^1 = \omega_1^d$ , and  $g'(0) = 0$ . We have

$$\alpha^1 = \frac{1 + \omega^1 g''(\omega^1 - \omega_1^d)}{\frac{1}{\psi_1} + \omega^1 g''(\omega^1 - \omega_1^d)}. \quad (65)$$

In the traditional case,  $g'' = 0$ , so  $\alpha^1 = \psi_1$ . In the completely rigid case,  $g'' = +\infty$ , so  $\alpha^1 = 1$  (indeed, we have then  $c_1 = \frac{\omega_1^d}{q_1}$ , so the elasticity of demand is 1).

This allows to calculate the  $\omega_{q_1}^1$ , the derivative of the account value as a function of the price. Starting from  $\omega^1 = q_1 c_1$ , we have

$$\omega_{q_1}^1 = c_1 + q_1 \frac{\partial c_1}{\partial q_1} = c_1 - c_1 \alpha_1 = c_1 \left[ 1 - \frac{1 + \omega^1 g''(\omega_1 - \omega_1^d)}{\frac{1}{\psi_1} + \omega^1 g''(\omega_1 - \omega_1^d)} \right],$$

so finally

$$\omega_{q_1}^1 = c_1 \frac{\frac{1}{\psi_1} - 1}{\frac{1}{\psi_1} + \omega^1 g''(\omega^1 - \omega_1^d)}. \quad (66)$$

By the budget constraint ( $\omega^0 + \omega^1 = w$ ) we have:

$$\omega_{q_1}^0 = -\omega_{q_1}^1.$$

### 9.3.4 Summarizing the effects of misperceptions and mental accounts

We again normalize  $p_i = 1$ . We call  $\alpha_i = -\frac{q_i}{c_i} c_{\tau_i}^i$  the empirical elasticity, which is  $\alpha_i = \psi_i$  in the traditional model, and  $\alpha_i = m_i \psi_i$  in the misperception model with attention  $m_i$  to the tax. We call  $\psi_i = -\frac{U^{ii}(c_i)}{c_i U^{iii}(c_i)}$ , which is simply the inverse of the curvature of the utility function  $U^i$  for good  $i$ .

This is also the “rational” elasticity, the demand elasticity that the agent would have if he was fully attentive; it might be the elasticity elicited in a careful procedure that makes the agent attentive to the tax.

We have the following Lemma.

**Lemma 9.2** *As explained just above, call  $\psi_i$  the “rational” demand elasticity, and  $\alpha_i$  the “behavioral” elasticity. In the limit of small taxes, welfare is:*

$$L(\tau) - L(0) = -\frac{1}{2} \sum_i -\frac{\alpha_i^2}{\psi_i} y_i \tau_i^2 + \Lambda \sum_i \tau_i y_i + o(\|\tau\|^2) + o(\|\tau\| \Lambda). \quad (67)$$

This implies the following.

**Proposition 9.11** *In the basic Ramsey model, the optimal tax is*

$$\tau_i = \Lambda \frac{\psi_i}{\alpha_i^2}. \quad (68)$$

where  $\psi_i$  is the underlying elasticity of true preferences, and  $\alpha_i$  is the behavioral elasticity.

**Proof** We can also use the general formulas (Proposition 2.1) to verify the result. However, it is also instructive to use the following derivation. Maximizing over  $\tau_i$ , the result from Lemma 9.2

$$L = -\frac{1}{2} \sum_i \frac{\alpha_i^2}{\psi_i} y_i \tau_i^2 + \Lambda \sum_i \tau_i y_i,$$

we find:  $\tau_i = \frac{\Lambda \psi_i}{\alpha_i^2}$ .  $\square$

For instance, in the traditional case  $\alpha_i = \psi_i$ , and we recover the traditional formula  $\tau_i = \frac{\Lambda}{\psi_i}$ .

**Proposition 9.12** *In the basic Pigou model, the optimal tax is  $\tau_i = \xi_i \frac{\psi_i}{\alpha_i}$ .*

**Proof** We would like this to be the first best allocation, so that  $u'(c_i) = 1 + \xi_i$ , i.e.  $c_i = c_i^d(1 - \psi_i \xi_i)$ . The response to the tax is:  $c_i = c_i^d(1 - \alpha_i \tau_i)$ . So optimal tax satisfies:  $\alpha_i \tau_i = \psi_i \xi_i$ , i.e.  $\tau_i = \xi_i \frac{\psi_i}{\alpha_i}$ .  $\square$

The Table shows the link between different models. We use  $\omega_i = c_i q_i$ .



	Ramsey problem	Pigou problem	Elasticity $\alpha_i$ to tax rate
General	$\tau_i = \Lambda \frac{\psi_i}{\alpha_i^2}$	$\tau_i = \xi_i \frac{\psi_i}{\alpha_i}$	$\alpha_i = -\frac{q^i}{c^i} c^i \tau_i$
Traditional model	$\tau_i = \frac{\Lambda}{\psi_i}$	$\tau_i = \xi_i$	$\alpha_i = \psi_i$
Misperception model	$\tau_i = \frac{\Lambda}{m_i^2 \psi_i}$	$\tau_i = \frac{\xi_i}{m_i}$	$\alpha_i = m_i \psi_i$
Mental account: rigid	$\tau_i = \Lambda \psi_i$	$\tau_i = \xi_i \psi_i$	$\alpha_i = 1$
Mental account: flexible	$\tau_i = \frac{\Lambda}{\psi_i} \left( \frac{1 + \psi_i \omega_i g_i''}{1 + \omega_i g_i''} \right)^2$	$\tau_i = \xi_i \frac{1 + \psi_i \omega_i g_i''}{1 + \omega_i g_i''}$	$\alpha_i = \frac{1 + \omega_i g_i''}{\frac{1}{\psi_i} + \omega_i g_i''}$
Hybrid model: Flexible mental account with misperceptions	$\tau_i = \frac{\Lambda}{\psi_i} \frac{1}{m_i^2} \left( \frac{1 + \psi_i \omega_i g_i''}{1 + \omega_i g_i''} \right)^2$	$\tau_i = \frac{\xi_i}{m_i} \frac{1 + \psi_i \omega_i g_i''}{1 + \omega_i g_i''}$	$\alpha_i = m_i \frac{1 + \omega_i g_i''}{\frac{1}{\psi_i} + \omega_i g_i''}$

### 9.3.5 Derivation of the agent's consumption in the mental accounting model of Section 3.6

The agent maximizes his perceived utility  $u^s(c_1, c_2) = \frac{c_1^{\alpha_1^s} c_2^{\alpha_2^s}}{\alpha_1^s \alpha_2^s}$  subject to the perceived budget constraint  $B(c_1, c_2) = \kappa_1 |\omega_1^d - c_1| + \sum_{i=1}^2 c_i \leq w$ , with  $\omega_1^d = \alpha_1^s w + \beta b$ . The first order conditions are

$$\begin{aligned} u_{c_1}^s &= \mu (1 - \kappa_1 \eta_1) \\ u_{c_2}^s &= \mu, \end{aligned}$$

where  $\eta_1 = \text{sign}(\omega_1^d - c_1)$  is the sign of  $\omega_1^d - c_1$  if that quantity is non-zero, and otherwise is some number in  $[-1, 1]$ . There are two cases.

Case 1. If  $\omega_1^d - c_1 = 0$  – this is the rigid mental account region. Consumption is:

$$c_1 = \omega_1^d, \quad c_2 = w - \omega_1^d.$$

Case 2. If  $\omega_1^d \neq c_1$ , then the agent has de facto a perceived price  $p_1^s = 1 - \kappa_1 \eta_1$  and  $p_2^s = 1$ . Consumptions are  $c_i = w \frac{\alpha_i^s / p_i^s}{\sum_j \alpha_j^s p_j / p_j^s}$  (Gabaix 2014, Example 4). In particular,

$$c_1 = \frac{\alpha_1^s}{1 - \alpha_2^s \kappa_1 \eta_1} \omega. \quad (69)$$

We summarize the results, in the case  $b \geq 0$ .

**Proposition 9.13** (Consumption with mental accounts) *Consumption of good 1 is as follows. For  $0 \leq b < b^*$ ,  $c_1 = \omega_1^d = \alpha_1^s (w^* + b) + \beta b$ . For  $b \geq b^*$ ,  $c_1$  is given by (69) with  $\eta_1 = 1$ . The cutoff  $b^*$  is the value at which those two expressions are equal, i.e. it is the solution of:*

$$\alpha_1^s (w^* + b) + \beta b = \frac{\alpha_1^s}{1 - \alpha_2^s \kappa_1} (w^* + b). \quad (70)$$

For a given voucher  $b$ , we are in the rigid account region if and only if  $\kappa_1 \geq \kappa_1^*$  where  $\kappa_1^*$  satisfies

(70).

## 9.4 Complements on Diamond-Mirrlees and Atkinson-Stiglitz (1972)

### 9.4.1 Diamond-Mirrlees: Concrete examples

To illustrate Proposition 5.1, consider the separable case  $u(\mathbf{c}) = c_0 + u(c_1)$  in the misperception case with  $\tau_1^s = \tau_1^p + m_1\tau_1^c$ ,  $0 \leq m_1 \leq 1$  and  $\tau_1^p$  is exogenous (perhaps set to 0).

We represent the production function as follows—it takes  $C(y_1)$  units of good 0 to produce  $y_1$  units of good 1. We define supply and demand to be  $S(p_1) = C'^{-1}(p_1)$  and  $D(p_1 + \tau_1^p + m_1\tau_1^c) = u'^{-1}(p_1 + \tau_1^p + m_1\tau_1^c)$ . We denote the corresponding supply and demand elasticities (corresponding to a fully perceived change in  $p_1$ ) by  $\varepsilon_S > 0$  and  $\varepsilon_D > 0$ . Differentiating the equilibrium condition  $S(p_1) = D(p_1 + \tau_1^p + m_1\tau_1^c)$  yields

$$\varepsilon_{11}^c = -\frac{\varepsilon_D}{\varepsilon_S + \varepsilon_D} m_1,$$

with  $\varepsilon_{11}^c = \frac{dp_1}{d\tau_1^c}$ . Compared to the traditional incidence analysis, because consumers are not fully attentive to the tax on good 1 ( $m_1 < 1$ ), the burden of the tax is shifted to the consumer. This echoes a result in Chetty, Looney and Kroft (2009).

We now turn to optimal taxes. We work in the limit of small taxes when  $\Lambda = \lambda - 1$  is close to 0 as in Section 3.1. Then, the optimal tax  $\tau_1^c$  satisfies

$$0 = \left( \Lambda c_1 - \tau_1^s \frac{\psi_1}{p_1} c_1 m_1 \right) + \left( \Lambda c_1 - \tau_1^s \frac{\psi_1}{p_1} c_1 \right) \varepsilon_{11}^c,$$

which we can rewrite as

$$\frac{\Lambda}{\psi_1} = \frac{\tau_1^p + m_1\tau_1^c}{p_1} \frac{m_1 + \varepsilon_{11}^c}{1 + \varepsilon_{11}^c}.$$

As long as  $m_1 < 1$ , the higher is the supply elasticity  $\varepsilon_S$ , the more the burden of the tax is shifted to the consumer, the higher is  $\varepsilon_{11}^c < 0$ , and the lower is the optimal tax.<sup>81</sup>

We next provide an example to illustrate Proposition 5.2.

We now show that production efficiency can fail with a restricted set of commodity taxes  $\boldsymbol{\tau}^p$ , even if there is a full set of commodity taxes  $\boldsymbol{\tau}^c$ . Consider the following example. There are two consumption goods, 0 and 1, two types of labor,  $a$  and  $b$ , a representative agent with decision utility  $u^s(c_0, c_1, l_a, l_b) = c_0 + U^s(c_1) - l_a - l_b$ , and experienced utility  $u^e(c_0, c_1, l_a, l_b) = u(c_0, c_1, l_a, l_b) - \xi_* c_1$ , where  $\xi_* > 0$  indicates an externality. For instance,  $c_1$  could be cigarette consumption. Hence, the government would like to discourage consumption of good 1.

<sup>81</sup>Another way to see this is as follows. Consider the optimal tax with infinitely elastic supply  $\varepsilon_S = \infty$  (a constant price  $p_1$ ). It satisfies  $(\Lambda c_1 - \tau_1^s \frac{\psi_1}{p_1} c_1 m_1) = 0$ . Now imagine that  $\varepsilon_S < \infty$ . Then at this tax  $(\Lambda c_1 - \tau_1^s \frac{\psi_1}{p_1} c_1) < 0$  so that  $(\Lambda c_1 - \tau_1^s \frac{\psi_1}{p_1} c_1) \varepsilon_{11}^c > 0$  and by implication  $(\Lambda c_1 - \tau_1^s \frac{\psi_1}{p_1} c_1 m_1) + (\Lambda c_1 - \tau_1^s \frac{\psi_1}{p_1} c_1) \varepsilon_{11}^c > 0$ . This implies that increasing the tax improves welfare.

The production function for good  $i$  is  $y_i = \left(\frac{l_{ia}}{\alpha_i}\right)^{\alpha_i} \left(\frac{l_{ib}}{1-\alpha_i}\right)^{1-\alpha_i}$ , with  $\alpha_i \in (0, 1)$ . As before, 0 is the untaxed good,  $\tau_0 = 0$ . The government can set taxes  $\tau_1$ ,  $\tau_a$  and  $\tau_b$  on good 1, labor of type  $a$  and labor of type  $b$ , and tax the employment of type  $a$  labor in sector 1. We assume that the consumer perfectly perceives taxes  $\tau_a, \tau_b$ , and prices  $p_0, p_1, p_a, p_b$  (the latter being the price of labor of type  $a, b$ ). In addition, the government can set a tax  $\tau_{1a}$  for the use of input  $a$  in the production of good 1. Note that production efficiency is equivalent to  $\tau_{1a} = 0$ .

**Proposition 9.14** *If the consumer is fully inattentive to the tax  $\tau_1$ , then the optimal tax system features production inefficiency:  $\tau_{1a} > 0$ . If the consumer is fully attentive to the tax  $\tau_1$ , then the optimal tax system features production efficiency:  $\tau_{1a} = 0$ .*

The essence is the following—the government would like to lower consumption of good 1, which has a negative externality. However, agents do not pay attention to the tax  $\tau_1$  on good 1, therefore a tax on good 1 will not be effective. We assume that the government cannot use producer taxes. Hence, the government uses a tax  $\tau_{1a} > 0$  on the input use in the production of good 1 (lowering production efficiency) to discourage the production of good 1, increase its price and discourage its consumption.

#### 9.4.2 Atkinson-Stiglitz (1972)

**Atkinson and Stiglitz (1972)** show uniform commodity taxation is optimal if preferences have the form  $u^h(c_0, \phi(\mathbf{C}))$ , with  $\mathbf{C} = (c_1, \dots, c_n)$ ,  $\phi$  homogeneous of degree 1, and  $c_0$  (the untaxed good) might be leisure. We now investigate how to generalize this result with behavioral agents.

**Proposition 9.15** *Consider the decision vs. experienced utility model. Assume that decision utility is of the form  $u^{s,h}(c_0, \phi^s(\mathbf{C}))$  and that experienced utility is of the form  $u^h(c_0, \phi(\mathbf{C}))$  with  $\phi^s$  and  $\phi$  homogeneous of degree 1. Then, if  $\phi^s = \phi$ , then uniform ad valorem commodity taxes are optimal (even though decision and experienced utility represent different preference orderings), but, if  $\phi^s \neq \phi$ , then uniform ad valorem commodity taxes are not optimal in general.*

The bottom line is that with behavioral biases, it is no longer sufficient to establish empirically that expenditure elasticities for  $(c_1, \dots, c_n)$  are unitary.

Another relevant consideration has to do with time horizons. Consider a tax reform and assume away any link between periods for simplicity (say because agents do not have access to asset markets). Imagine a situation where, in the long-run, choices can be represented by a decision utility  $u^{s,h}(c_0, \phi^s(\mathbf{C}))$ , and welfare can be evaluated with an experienced utility  $u^h(c_0, \phi(\mathbf{C}))$  with  $\phi = \phi^s$ . But, in the short-run as the tax code changes, agents misperceive taxes and, hence, make different choices. Then optimal time-varying taxes might be uniform in the long run but not in the short run. Likewise, if agents pay differential attention to taxes (at least in the short run), the **Atkinson and Stiglitz (1972)** neutrality result will fail.

## 9.5 Other extensions

### 9.5.1 Cross-Effects of Attention

We again normalize  $p_i = 1$ . How does attention to one good affect the optimal tax on another? To answer this question, we use the specialization of the general model developed in Section 2.7, assuming a representative consumer (so that we drop the index  $h$ ), no internality/externality so that  $\boldsymbol{\tau}^X = 0$ , and in the limit of small taxes. Defining  $\Lambda = \frac{\lambda}{\gamma} - 1$ , we can rewrite formula (11), in the limit of small  $\Lambda$ , as

$$\boldsymbol{\tau} = -\Lambda (\mathbf{M}' \mathbf{S}^r \mathbf{M})^{-1} \mathbf{c}.$$

This is a generalization of Proposition 3.1, which assumed a diagonal matrix  $\mathbf{S}^r$ .

To gain intuition, we take  $n = 2$  goods,  $\mathbf{M} = \text{diag}(m_1, m_2)$ , we normalize prices to  $p_1 = p_2 = 1$ , and we write the rational Slutsky matrix as  $S_{ii}^r = -c_i \psi_i$  for  $i = 1, 2$ , and  $S_{12}^r = S_{21}^r = -\sqrt{c_1 c_2 \psi_1 \psi_2} \rho$ .

**Proposition 9.16** (Impact of cross-elasticities on optimal taxes with inattentive agents) *With two*

*taxed goods, the optimal tax on good 1 is  $\tau_1 = \frac{\Lambda}{m_1^2 \psi_1} \frac{1 - \rho \sqrt{\frac{m_1^2 \psi_1 c_2}{m_2^2 \psi_2 c_1}}}{1 - \rho^2}$ . When attention to the tax of good 2  $m_2$  falls, the optimal tax on good 1 increases (respectively decreases) if goods 1 and 2 are substitutes (respectively complements).*

Suppose for example that the goods are substitutes with  $\rho < 0$ .<sup>82</sup> When  $m_2$  falls, the optimal tax on good 2 increases by the effects in Proposition 3.1, and optimal taxes on substitute goods also increase.<sup>83</sup>

### 9.5.2 Tax instruments with differential saliences

We elaborate on a remark we made at the end of section 3.7. As an extreme example, consider again the basic Ramsey example outlined above, and assume that the two tax systems with salience  $m$  and  $m'$  can be used jointly. Consider the case where there is only one agent and only one (taxed) good. With  $m' > m$ , we get

$$0 = (\lambda - \gamma) c + [\lambda \tau + \gamma(\bar{\tau}^s - \bar{\tau})] m \mathbf{S}^r, \quad 0 = (\lambda - \gamma) c + [\lambda \tau + \gamma(\bar{\tau}^s - \bar{\tau})] m' \mathbf{S}^r,$$

where  $\bar{\tau}^s$  is the total perceived tax arising from the joint perception of the two tax instruments. This requires  $\lambda = \gamma$  and with  $\bar{\tau}^s = 0$ . In other words, the solution is the first best. This is because a planner can replicate a lump sum tax by combining a tax  $\tau$  with low salience  $m$  and a tax  $-\tau \frac{m}{m'}$

<sup>82</sup>We have  $\rho^2 < 1$  since  $\mathbf{S}^r$  is a  $2 \times 2$  negative definite matrix so that  $0 < \det \mathbf{S}^r = c_1 c_2 \psi_1 \psi_2 (1 - \rho^2)$ .

<sup>83</sup>Perhaps curiously, we can have  $\frac{\partial \tau_1}{\partial m_1} > 0$  with complement goods  $\rho > 0$ . This happens if and only if  $2 < \frac{m_1}{m_2} \rho \sqrt{\frac{\psi_1 c_2}{\psi_2 c_1}}$ . That latter condition is quite extreme, and would imply that  $\tau_1 < 0$  even though the planner wants to raise revenues. This is because the planner wants to increase consumption of the low elasticity (low  $m_2, \psi_2$ ), good 2, he wants to subsidize good 1 if it is a strong complement of good 2.

with high salience  $m' > m$ , generating tax revenues  $\tau \frac{m'-m}{m'}$  per unit of consumption of the taxed good with no associated distortion. This is an extreme result, already derived by [Goldin \(2015\)](#). In general, with more than one agent and heterogeneities in the misperceptions of the two taxes, the first best might not be achievable.

### 9.5.3 A different budget adjustment rule

When perceived prices  $q_j^s$  are different from the true prices  $q_j$ , some adjustment is needed for the budget constraint. Let us study a different rule, where a certain good  $n$  (“the last good”, imagining a temporal order) bears the brunt of the budget adjustment (it’s a “shock absorber”). This leads to

$$c^{i,s}(\mathbf{q}, \mathbf{q}^s, w) = c^{i,r}(\mathbf{q}^s, w) \text{ for } i \neq n \quad (71)$$

$$c^{n,s}(\mathbf{q}, \mathbf{q}^s, w) = \frac{1}{q_n} \left( w - \sum_{i \neq n} q_i c^{i,s}(\mathbf{q}, \mathbf{q}^s, w) \right). \quad (72)$$

This is: for all goods but the last one, the consumer only pays attention to perceived prices. Only for the last one does she see the budget constraint.<sup>84</sup> We shall see in the next proposition that we can also write

$$c^{n,s}(\mathbf{q}, \mathbf{q}^s, w) = c^{n,r}(\mathbf{q}^s, w) - \frac{1}{q_n} (\mathbf{q} - \mathbf{q}^s) \cdot \mathbf{c}^r(\mathbf{q}^s, w), \quad (73)$$

i.e. actual consumption of good  $n$  is planned consumption  $c^{n,r}(\mathbf{q}^s, w)$  minus the adjustment for the surprise  $(\mathbf{q} - \mathbf{q}^s) \cdot \mathbf{c}^r(\mathbf{q}^s, w)$  in the actual cost of the goods  $i < n$  that have been purchased before good  $n$ .

For completeness, we record the Slutsky matrix properties of that rules. (Here we consider the income-compensated matrix  $S^C$ ).

**Proposition 9.17** (With the “last good adjusting for the budget” rule) *Consider the model above, with attention  $m_j$  to price  $j$ . Evaluating at  $\mathbf{q}^s = \mathbf{q}$ , the marginal propensity to consume out of wealth isn’t changed:*

$$\partial_w c_i^s(\mathbf{q}, \mathbf{q}^s, w) = \partial_w c_i^r(\mathbf{q}, w). \quad (74)$$

However, the Slutsky matrix  $S_{ij}^s$  is changed as follows:

$$S_{ij}^s = S_{ij}^r m_j + \left( \partial_w c_i^r - \frac{1}{q_n} 1_{i=n} \right) (1 - m_j) c^j, \quad (75)$$

where  $S_{ij}^r$  is the rational Slutsky matrix.

<sup>84</sup>[Chetty, Looney and Kroft \(2009\)](#) consider such a rule in a 2-good context. [Gabaix \(2016\)](#) consider such a rule when doing dynamic programming, and the last good is “next period wealth”.

**Proof** The term  $\partial_w c_i^s$  is trivial, as it's evaluate a  $\mathbf{q}^s = \mathbf{q}$ . We move on to the  $S_{ij}$ . First, take  $i \neq n$ . Then,  $c^{i,s}(\mathbf{q}, \mathbf{q}^s, w) = c^{i,r}(\mathbf{q}^s, w)$ , hence:

$$\begin{aligned} S_{ij}^s &= \partial_{q_j} c^{i,r}(\mathbf{q}^s, w) + c_w^i c^j \\ &= c_{q_j}^{i,r}(\mathbf{q}, w) m_j + c_w^i c^j = (S_{ij}^r - c_w^i c^j) m_j + c_w^i c^j \\ &= S_{ij}^r + c_w^i c^j (1 - m_j), \end{aligned}$$

which gives the announced result.

For good  $n$ , we rewrite:

$$\begin{aligned} q_n c^{n,s}(\mathbf{q}, \mathbf{q}^s, w) &= w - \sum_{i \neq n} q_i c^{i,s}(\mathbf{q}, \mathbf{q}^s, w) \\ &= \mathbf{q}^s \cdot \mathbf{c}^r(\mathbf{q}^s, w) - (\mathbf{q} \cdot \mathbf{c}^r(\mathbf{q}^s, w) - q_n c^{n,r}(\mathbf{q}^s, w)) \\ &= (\mathbf{q}^s - \mathbf{q}) \cdot \mathbf{c}^r(\mathbf{q}^s, w) + q_n c^{n,r}(\mathbf{q}^s, w), \end{aligned}$$

i.e. another useful expression:

$$c^{n,s}(\mathbf{q}, \mathbf{q}^s, w) = c^{n,r}(\mathbf{q}^s, w) + \frac{1}{q_n} (\mathbf{q}^s - \mathbf{q}) \cdot \mathbf{c}^r(\mathbf{q}^s, w).$$

Its interpretation is that the consumption of the last good is the planned consumption (the first term,  $c^{n,r}(\mathbf{q}^s, w)$ ), plus an adjustment for the “surprise” difference between planned and actual expenditure (the last term).

Now, differentiate w.r.t.  $q_j$ :

$$S_{nj} = (\partial_{q_j} c^{n,r}(\mathbf{q}^s, w) + c_w^n c^j) + \partial_{q_j} \left( \frac{1}{q_n} (\mathbf{q}^s - \mathbf{q}) \cdot \mathbf{c}^r(\mathbf{q}^s, w) \right).$$

By the earlier calculation of  $S_{ij}$ , the first term is  $S_{nj}^r + c_w^n c^j (1 - m_j)$  with  $i = n$ , by the earlier result, and the last term is (as we evaluate at  $\mathbf{q}^s = \mathbf{q}$ )

$$\partial_{q_j} \left( \frac{1}{q_n} \sum_i (q_i^s - q_i) c^i(\mathbf{q}^s, w) \right) = \frac{1}{q_n} (m_j - 1) c^j(\mathbf{q}^s, w).$$

This gives the announced result.

□

**Behavioral wedges** We take a particular case, which is particularly tractable. There are  $n - 2$  goods, and good  $n$  is the “shock absorber” good. The price of goods 0 and  $n$  is normalized to 1. There’s no tax on goods 0 and 2, for simplicity.

Utility is:

$$u(c_0, \dots, c_n) = c_0 + \sum_{i=1}^n u^i(c_i).$$

Good 0 has marginal utility of 1, which absorbs income effects, so  $v_w = 1$ . Hence,  $\tau^b = \mathbf{q} - \frac{u_c}{v_w}$  is:

$$\begin{aligned}\tau_0^b &= 0 \\ \tau_i^b &= q_i - q_i^s \text{ for } i = 1, \dots, n-1 \\ \tau_n^b &= 1 - u'_n(c_n)\end{aligned}$$

for  $c_i = c_i^r(q_i^s)$  for  $1 \leq i < n$  and  $c_n = c_n^r + \sum_{i < n} (q_i^s - q_i) c_i$  (from 73).

**Derivation of the optimal tax and**  $\Lambda_i = \frac{\Lambda - (1-\Lambda)(1-m_i)\mu}{1 - (1-\Lambda)(1-m_i)\mu}$  We first provide an intuitive proof. We again normalize  $p_i = 1$ . The distortion on good  $n$  is, from 73

$$c_n - c_n^* = - \sum_{i < n} (1 - m_i) \tau_i c_i,$$

and the distortion on good 0 is  $-(c_n - c_n^*)$ . Hence in the objective function we have

$$\begin{aligned}L &= W + \sum_{i < n} \lambda \tau_i c_i - \mu \sum_{i < n} (1 - m_i) \tau_i c_i \\ &= W + \sum_{i < n} (\lambda - \mu(1 - m_i)) \tau_i c_i,\end{aligned}$$

where  $W$  =utility distortion all goods except 0 and  $n$ . Hence, we just replace  $\lambda$  by  $\lambda' = \lambda - \mu(1 - m_i)$ .

Remember that we write  $\lambda = \frac{1}{1-\Lambda}$ . Hence, this which corresponds to

$$\Lambda' = 1 - \frac{1}{\lambda'} = 1 - \frac{1}{\frac{1}{1-\Lambda} - (1 - m_i)\mu} = 1 - \frac{1 - \Lambda}{1 - (1 - \Lambda)(1 - m_i)\mu} = \frac{\Lambda - (1 - \Lambda)(1 - m_i)\mu}{1 - (1 - \Lambda)(1 - m_i)\mu}.$$

We also provide a more computational proof, which we found also instructive. Take an  $i = 1, \dots, n-1$ . We have  $\tau_i^b = (1 - m_i) \tau_i$  and  $\tau_n^b = 1 - u'(c_n) = -\mu$ . We have  $S_{ii} = -\frac{\psi_i c_i}{q_i^s} m_i$ , while

$$S_{ni} = -(1 - m_i) c_i \left( 1 - \psi_i \frac{\tau_i}{q_i^s} m_i \right).$$

Plugging this into the general Ramsey optimal tax formula (7) gives:

$$\begin{aligned}
0 &= (\lambda - \gamma) c_i + \lambda(\boldsymbol{\tau} - \tilde{\boldsymbol{\tau}}^b) \cdot \mathbf{S}_i^C \\
&= (\lambda - 1) c_i - \lambda \left( \tau_i - \frac{1}{\lambda} (1 - m_i) \tau_i \right) \frac{\psi_i c_i}{q_i^s} m_i - \lambda \left( 0 + \frac{\mu}{\lambda} \right) (1 - m_i) c_i \left( 1 - \psi_i \frac{\tau_i}{q_i^s} m_i \right) \\
&= (\lambda - 1 - \mu(1 - m_i)) c_i - \frac{\psi_i c_i}{q_i^s} m_i \tau_i (\lambda - (1 - m_i) - \mu(1 - m_i)),
\end{aligned}$$

which is the expression with  $\mu = 0$ , if we replace  $\lambda$  by  $\lambda' = \lambda - (1 - m_i) \mu$ .

**Analysis of small taxes** We analyze the case of small taxes. Compared to the first best, distortions are:

$$\begin{aligned}
c_i - c_i^* &= -\psi_i c_i^* m_i \tau_i \\
c_n - c_n^* &= -\sum_{i=1}^{n-1} c_i (1 - m_i) \tau_i,
\end{aligned}$$

and a utility loss equal to  $L^D$  such that:

$$-2L^D = \sum_i \frac{1}{\psi_i c_i^{*2}} (c_i - c_i^*)^2.$$

Hence we have the following generalization of the objective function in the simple Ramsey case with small taxes (we normalized prices to  $p_i = 1$ )

$$L = -\frac{1}{2} \sum_{i=1}^{n-1} \tau_i^2 \psi m_i^2 c_i^* - \frac{1}{2} \frac{1}{\psi_n c_n^*} \left( \sum_{i=1}^{n-1} c_i^* (1 - m_i) \tau_i \right)^2 + \lambda \sum_i \tau_i c_i^*. \quad (76)$$

In particular, now the distortion is not just  $\psi_i m_i^2$  as before, but there is another term, multiplied by  $\frac{1}{\psi_n}$ . Hence, attention is beneficial only if risk aversion ( $\frac{1}{\psi_n}$ ) for the shock absorber good is small enough (in the baseline model it is 0). The optimal tax is

$$\tau_i = \frac{\Lambda_i}{m_i^2 \psi_i},$$

with

$$\Lambda_i = \Lambda - \mu(1 - m_i),$$

and

$$\mu = \frac{1}{\psi_n c_n^*} \sum_{i=1}^{n-1} c_i (1 - m_i) \tau_i,$$

which is the marginal distortion on good  $n$ .



We can also study the variant in the done in the main body of the paper. A variant:  $u'_2(c_2) \leq 1$  for  $c_2 \geq c_2^d$  and  $1 + \mu > 1$  for  $c_2 < c_2^d$ . Then losses are:

$$\begin{aligned} L &= - \sum_{i=1}^{n-1} \left( \frac{1}{2} \psi_i m_i^2 c_i^* \tau_i^2 + \mu c_i^* (1 - m_i) \tau_i \right) + \lambda \sum_i \tau_i y_i \\ &= - \sum_{i=1}^{n-1} \left( \frac{1}{2} \psi_i m_i^2 c_i^* \tau_i^2 \right) + \sum_i (\lambda - \mu (1 - m_i)) \tau_i c_i, \end{aligned}$$

so, optimal tax on good  $i$  is 0 iff  $\mu \psi_i (1 - m_i) > \lambda$ .

**Impact on Pigouvian taxes** We revisit our simple model of Section 3.2, with an externality on good 1. We have  $\lambda = 1$ , so that the government's objective function is:

$$L = U(c_1) - (p + \xi) c_1 + u_2(c_2^* - (1 - m) c_1 \tau) + (1 - m) c_1 \tau.$$

i.e. utility from good 1, utility from good 2 (which absorbs the shock  $(1 - m) c_1 \tau$ ), and consumption of good 0 is increased by the lump-sum rebate, which accounts for the last term. The consumer chooses  $c_1$  according to  $U'(c_1) = p + m\tau$ .

We take utility  $U(c) = Qc - \frac{c^2}{2\Psi}$ , so that demand is  $c_1 = \Psi(Q - p - m\tau)$ . We keep  $u'_2(c_2) = 1 + \mu$ . We have:

$$\begin{aligned} L'(\tau) &= [p + m\tau - (p + \xi)](-\Psi m) + [-(1 - m)(1 + \mu) + (1 - m)] \frac{d}{d\tau}(c_1 \tau) \\ &= -(m\tau - \xi) \Psi m - (1 - m) \mu (c_1 - \Psi m \tau) \\ &= -(m\tau - \xi) \Psi m - (1 - m) \mu (\Psi(Q - p - 2m\tau)), \end{aligned}$$

which leads to:

$$\begin{aligned} \tau &= \frac{\frac{\xi}{m} - \mu \left( \frac{1-m}{m} \right) (Q - p)}{1 - 2\mu \left( \frac{1-m}{m} \right)} \\ &= \frac{\frac{\xi}{m} - \mu \left( \frac{1-m}{m} \right) \frac{c_1^0}{\Psi}}{1 - 2\mu \left( \frac{1-m}{m} \right)}, \end{aligned}$$

where  $c_1^0 = \Psi(Q - p)$  is the consumption of good 1 if there is no tax.

Hence, the government doesn't tax the good if:  $\xi \Psi < \mu (1 - m) c_1^0$ , i.e. if the externality is too small.

### 9.5.4 Precisions on jointly optimal taxes and nudges

We assume that  $(\tau - \tau^{\xi h}) \cdot \mathbf{c}_w^h = 0$ , so that  $\beta^h = \gamma^h$ . As in Section 2.7, we call  $\xi^h = \tau^{X,h} = \frac{\gamma^{\xi,h}}{\lambda} \tau^{I,h} + \tau^{\xi,h}$  as the sum of the internality plus externality.

$$c^h(\tau, \tau^{X,h}) = c_0^h - \Psi(m^h \tau + \tau^{X,h}).$$

Individual  $h$  creates an externality plus internality. We have successively:

$$\begin{aligned} \Xi &= -\mathbb{E}[\beta^{h'}] = -\mathbb{E}[\gamma^{h'}] = -\lambda \\ \tau^{\xi h} &= \frac{-\Xi}{\lambda} \tau^{\xi,h} = \tau^{\xi,h} \\ \tau^{bh} &= (1 - m^h) \tau + \tau^{I,h} - \tau^{X,h} \\ \tilde{\tau}^{b,\xi,h} &= \frac{\gamma^{\xi,h}}{\lambda} \tau^{b,h} \\ \lambda(\tau - \tau^{\xi,h} - \tilde{\tau}^{b,\xi,h}) &= \lambda(\tau - \tau^{\xi,h}) - \gamma^h((1 - m^h) \tau + \tau^{I,h} - \tau^{X,h}) \\ &= (\lambda - \gamma^h(1 - m^h)) \tau - \lambda \tau^{X,h} + \gamma^h \tau^{X,h}. \end{aligned}$$

Proposition 2.2 gives, using  $\mathbf{S}_i^{H,h} = -\Psi m^h$ ,

$$\begin{aligned} \frac{\partial L}{\partial \tau} &= \mathbb{E} \sum_h (\lambda - \gamma_h) c^h - \lambda(\tau - \tau^{\xi,h} - \tilde{\tau}^{b,\xi,h}) \Psi m^h \\ \frac{\partial L}{\partial \tau} &= \mathbb{E} \sum_h (\lambda - \gamma_h) c^h - [(\lambda - \gamma^h(1 - m^h)) \tau - \lambda \tau^{X,h} + \gamma^h \tau^{X,h}] \Psi m^h. \end{aligned} \quad (77)$$

We use the notation

$$\sigma_{Y,Z} = \text{cov}(Y_h, Z_h).$$

Using  $\mathbb{E}[\gamma^h] = \lambda$ , we have:

$$-\mathbb{E}[(\lambda - \gamma^h) m^h] = \mathbb{E}[\gamma^h m^h] - \mathbb{E}[\gamma^h] \mathbb{E}[m^h] = \sigma_{\gamma,m}.$$

Hence, using  $\tau^{X,h} = \chi \eta^h$

$$\begin{aligned} \frac{1}{\Psi} \frac{\partial L}{\partial \tau} &= -\mathbb{E}[(\lambda - \gamma^h(1 - m^h)) m^h] \tau - \mathbb{E}[\gamma^h \eta^h m^h] \chi + \mathbb{E} \left[ (\lambda - \gamma_h) \frac{c^h}{\Psi} + \lambda \tau^{X,h} m^h \right] \\ &= -\mathbb{E}[(\lambda - \gamma^h(1 - m^h)) m^h] \tau - \mathbb{E}[\gamma^h \eta^h m^h] \chi + \mathbb{E} \left[ (\lambda - \gamma^h) \frac{c^h}{\Psi} + \lambda \tau^{X,h} m^h \right] \\ &= -\mathbb{E}[\gamma^h m^{h^2} - \sigma_{\gamma m}] \tau - \mathbb{E}[\gamma^h \eta^h m^h] \chi + \mathbb{E}[\lambda \tau^{X,h} m^h - \sigma_{\gamma, \frac{c}{\Psi}}]. \end{aligned}$$

Proposition 2.3 gives:

$$\begin{aligned}\frac{\partial L}{\partial \chi} &= \sum_h [\lambda (\tau - \tau^{\xi h}) - \beta^h \tau^{b,\xi,h}] \mathbf{c}_\chi^h \\ &= -\mathbb{E} \sum_h [\lambda (\tau - \tau^{X,h}) - \gamma^h ((1 - m^h) \tau - \tau^{X,h})] \Psi \tau_\chi^{X,h} \\ \frac{\partial L}{\partial \chi} &= -\mathbb{E} \sum_h [(\lambda - \gamma^h (1 - m^h)) \tau - \lambda \tau^{X,h} + \gamma^h \tau^{X,h}] \Psi \tau_\chi^{X,h}.\end{aligned}$$

Using  $\tau^{X,h} = \chi \eta^h$  gives  $\tau_\chi^{X,h} = \eta^h$  hence:

$$\begin{aligned}\frac{1}{\Psi} \frac{\partial L}{\partial \chi} &= -\mathbb{E} \sum_h [(\lambda - \gamma^h (1 - m^h)) \tau - \lambda \tau^{X,h} + \gamma^h \tau^{X,h}] \eta^h \\ &= -\mathbb{E} [\gamma^h \eta^h m^h - \sigma_{\gamma,\eta}] \tau - \mathbb{E} [\gamma^h \eta^{h^2}] \chi + \mathbb{E} [\lambda \tau^{X,h} \eta^h].\end{aligned}$$

This implies

$$\frac{1}{\Psi} \frac{\partial^2 L}{\partial \tau \partial \chi} = -\mathbb{E} \sum_h [(\lambda - \gamma^h (1 - m^h)) \eta^h].$$

Hence, at the optimum:

$$\begin{aligned}\mathbb{E} [\gamma^h m^{h^2} - \sigma_{\gamma m}] \tau + \mathbb{E} [\gamma^h \eta^h m^h] \chi &= \mathbb{E} [\lambda \tau^{X,h} m^h - \sigma_{\gamma,c/\Psi}] \\ \mathbb{E} [\gamma^h \eta^h m^h - \sigma_{\gamma,\eta}] \tau + \mathbb{E} [\gamma^h \eta^{h^2}] \chi &= \mathbb{E} [\lambda \tau^{X,h} \eta^h].\end{aligned}$$

Solving for the two unknowns  $\tau$  and  $\chi$  gives the following.

**Proposition 9.18** *The optimal tax and nudge satisfy*

$$\begin{aligned}\tau &= \frac{\mathbb{E} [\gamma^h \eta^{h^2}] \mathbb{E} [\lambda \tau^{X,h} m^h - \sigma_{\gamma,c/\Psi}] - \mathbb{E} [\gamma^h \eta_h m^h] \mathbb{E} [\lambda \tau^{X,h} \eta^h]}{\mathbb{E} [\gamma^h \eta^{h^2}] \mathbb{E} [\gamma^h m^{h^2} - \sigma_{\gamma m}] - \mathbb{E} [\gamma^h \eta^h m^h] \mathbb{E} [\gamma^h \eta^h m^h - \sigma_{\gamma,\eta}]} \\ \chi &= \frac{\mathbb{E} [\lambda \tau^{X,h} \eta^h] \mathbb{E} [\gamma^h m^{h^2} - \sigma_{\gamma m}] - \mathbb{E} [\lambda \tau^{X,h} m^h - \sigma_{\gamma,c/\Psi}] \mathbb{E} [\gamma^h \eta^h m^h - \sigma_{\gamma,\eta}]}{\mathbb{E} [\gamma^h \eta^{h^2}] \mathbb{E} [\gamma^h m^{h^2} - \sigma_{\gamma m}] - \mathbb{E} [\gamma^h \eta^h m^h] \mathbb{E} [\gamma^h \eta^h m^h - \sigma_{\gamma,\eta}]}.\end{aligned}$$

### 9.5.5 Worked out examples of endogenous attention

**A linear-quadratic example** To illustrate the situation, we work out completely a linear-quadratic example. Take decision utility have  $u^s(c_0, c_1, m) = c_0 + U(c_1) - g(m)$ , with  $U(c) = \frac{ac - \frac{1}{2}c^2}{\Psi}$  and attention technology  $p_1^s(p_1, m) = p_1^d + m\tau_1$ , where  $\tau_1$  is a tax. Full utility is  $u(c_0, c_1, m) = c_0 + U(c_1) - Ag(m)$ , where  $A = 0$  in the “no attention cost in welfare” case, and  $A = 1$  in the “optimally allocated attention” case.

We assume  $p_0 = 1$ ,  $\Psi > 0$ . Given attention  $m$ , demand satisfies  $U'(c_1) = p^s$ , so  $c_1^r(p^s) = a - \Psi p^s$ . The perceived tax is:

$$\tau_1^s = m(\tau_1) \tau_1,$$

and demand is

$$c_1 = a - \Psi (p_1^d + m(\tau_1) \tau_1).$$

The losses from inattention are  $\frac{1}{2} u_{c_1 c_1} (c_1^r - c_1)^2 = -\frac{1}{2} \Psi \tau^2 (1 - m)^2$ . (This is always true to the leading order, and here this is exact as the function is quadratic). Hence, the optimal attention problem is:

$$m(\tau_1) = \arg \max_m -\frac{1}{2} \Psi \tau_1^2 (1 - m)^2 - g(m),$$

whose first order condition is:

$$g'(m) = \Psi (1 - m) \tau_1^2. \quad (78)$$

The Slutsky matrix with constant  $m$  has:

$$S_{11|m}^H = \frac{\partial c_1(p_1^d + m\tau_1, m)}{\partial \tau_1} = -\Psi m,$$

while with variable attention  $m(p)$ , we have:

$$\begin{aligned} S_{11}^H &= \frac{dc_1(p_1^d + m\tau_1, m(\tau_1))}{d\tau_1} = -\Psi (m + \tau m'(\tau_1)) \\ S_{21}^H &= \frac{\partial m}{\partial \tau_1} = m'(\tau_1). \end{aligned}$$

Then, we have:  $\tau^b = (0, q - q^s, Ag'(m)) = (0, \tau_1(1 - m), Ag'(m))$ , and given  $\tau = (0, \tau_1, 0)$ , so

$$\begin{aligned} \tau - \tau^b &= (0, \tau_1 m, -Ag'(m)) \\ S_1^H &= (0, -\Psi (m + \tau_1 m'(\tau_1)), m'(\tau_1)). \end{aligned}$$

Applying Proposition 2.1 gives:

$$\begin{aligned} \frac{\partial L(\boldsymbol{\tau}, w)}{\partial \tau_1} &= (\lambda - \gamma) c_1 + \lambda (\tau - \tau^b) \cdot S_1^H \\ &= (\lambda - \gamma) c_1 - \Psi \tau_1 m (m + \tau_1 m'(\tau_1)) - Ag'(m) m'(\tau_1). \end{aligned}$$

Normalize  $\lambda = 1$ ,  $\gamma = 1 - \Lambda$ , and define  $\psi_1(c_1) = \Psi/c_1$ . First, when  $m_1$  is exogenous, we verify our formula from Section 2

$$\frac{\partial L(\boldsymbol{\tau}, w)}{\partial \tau_1} = \Lambda c_1 - \Psi m \tau_1^s.$$

i.e.  $\tau_1^s = \frac{\Lambda}{m\psi_1}$ ,  $\tau_1 = \frac{\Lambda}{m^2\psi_1}$ .

Next, in the “no attention cost in welfare” case,  $A = 0$

$$\frac{\partial L(\boldsymbol{\tau}, w)}{\partial \tau_1} = \Lambda c_1 - \Psi \tau_1^s \frac{d\tau_1^s(\tau_1, m_1(\tau_1))}{d\tau_1} = \Lambda c_1 - \Psi \tau_1^s (m_1 + \tau_1 m'(\tau_1)),$$

so

$$\tau_1^s = \frac{-c_1 \Lambda}{S_{11}^H} = \frac{\Lambda}{(m + \tau_1 m'(\tau_1)) \psi_1}, \quad \tau_1 = \frac{\Lambda}{(m^2 + \tau_1 m m'(\tau_1)) \psi_1}.$$

Finally, in the “optimally allocated attention” case,  $A = 1$ . First, we verify:

$$\begin{aligned} -D_1 &= \tau^b \cdot \mathbf{S}^H = (0, \tau_1(1-m), g'(m)) \cdot (0, -\Psi(m + \tau_1 m'(p)), m'(p_1)) \\ &= -\Psi(m + \tau_1 m'(p)) \tau_1(1-m) + g'(m) m'(p_1) = -\Psi m \tau_1(1-m) = -(\tau_1 - \tau_1^s) \Psi m \\ &= \tau_C^s \cdot \mathbf{S}_{j|m}^H(\mathbf{p}, w, m), \end{aligned}$$

with  $\tau_C^s = \tau_1 - \tau_1^s = (1-m)\tau_1$  and  $\mathbf{S}_{j|m}^H(\mathbf{p}, w, m) = -\Psi m$ .

$$\begin{aligned} \frac{\partial L(\boldsymbol{\tau}, w)}{\partial \tau_1} - \Lambda c_1 &= (\tau - \tau^b) \cdot \mathbf{S}_1^H = \boldsymbol{\tau} \cdot \mathbf{S}_1^H - \tau^b \cdot \mathbf{S}_1^H = \boldsymbol{\tau} \cdot \mathbf{S}_1^H - \tau_C^s \cdot \mathbf{S}_{1|m}^H \\ &= -\Psi \tau (m + \tau_1 m'(\tau_1)) + \Psi \tau (1-m)m = -\Psi \tau (m^2 + \tau m'(\tau)) \end{aligned}$$

$$\frac{\partial L(\boldsymbol{\tau}, w)}{\partial \tau_1} = \Lambda c_1 - \Psi \tau (m^2 + \tau m'(\tau)),$$

so

$$\tau = \frac{\Lambda / \psi_1}{m(\tau)^2 + \tau m'(\tau)}.$$

**An isoelastic example** We now work out completely an isoelastic example. Take decision utility have  $u^s(c_0, c_1, m) = c_0 + U(c_1) - g(m)$ , with  $U(c) = \frac{c^{1-1/\psi}}{1-1/\psi}$  and attention technology  $p_1^s(p_1, m) = p_1^d + m\tau_1$ , where  $\tau_1$  is a tax. Full utility is  $u(c_0, c_1, m) = c_0 + U(c_1) - Ag(m)$ , where  $A = 0$  in the “no attention cost in welfare” case, and  $A = 1$  in the “optimally allocated attention” case.

We assume  $p_0 = 1$ . The perceived tax is:

$$\tau_1^s = m(\tau_1) \tau_1,$$

and demand is

$$c_1 = (p_1 + m(\tau_1) \tau_1)^{-\psi}.$$

The Slutsky matrix with constant  $m$  has:

$$S_{11|m}^H = \frac{\partial c_1(p_1 + m\tau_1, m)}{\partial \tau_1} = -\Psi m,$$

where (to leverage the calculations already done for the quadratic utility case) we define:

$$\Psi = \psi \frac{c_1}{q_1^s},$$

with  $q_1^s = p_1 + m_1(\tau_1)$ , while with variable attention  $m(\tau)$ , we have:

$$\begin{aligned} S_{11}^H &= \frac{dc_1(p_1^d + m\tau_1, m(\tau_1))}{d\tau_1} = -\Psi(m + \tau m'(\tau_1)) \\ S_{21}^H &= \frac{\partial m}{\partial \tau_1} = m'(\tau_1). \end{aligned}$$

Then, we have:  $\tau^b = (0, q - q^s, Ag'(m)) = (0, \tau_1(1 - m), Ag'(m))$ , and given  $\tau = (0, \tau_1, 0)$ , and  $\tilde{\tau}^b = (1 - \frac{\beta}{\lambda})\tau = (1 - \Lambda)\tau$  so

$$\begin{aligned} \tau - \tilde{\tau}^b &= \tau - (1 - \Lambda)\tau^b = (0, \tau_1(1 - (1 - m)(1 - \Lambda)), -(1 - \Lambda)Ag'(m)) \\ &= (0, \tau_1(m + \Lambda(1 - m)), -(1 - \Lambda)Ag'(m)) \\ S_1^H &= (0, -\Psi(m + \tau_1 m'(\tau_1)), m'(\tau_1)). \end{aligned}$$

Applying Proposition 2.1 gives (with  $\lambda = 1, \gamma = 1 - \Lambda$ )

$$\begin{aligned} \frac{\partial L(\tau, w)}{\partial \tau_1} &= (\lambda - \gamma)c_1 + \lambda(\tau - \tilde{\tau}^b) \cdot S_1^H \\ &= \Lambda c_1 - \Psi \tau_1(m + \Lambda(1 - m))(m + \tau_1 m'(\tau_1)) - Ag'(m)m'(\tau_1). \end{aligned}$$

Define

$$\psi_1(c_1) = \frac{\Psi}{c_1} = \frac{\psi}{q_1^s}.$$

First, when  $m_1$  is exogenous, we verify our formula (13):

$$0 = \Lambda - \frac{\psi}{q_1^s} \tau_1(m + \Lambda(1 - m))m,$$

i.e.  $\frac{\tau_1}{q_1^s} = \frac{\Lambda}{\psi_1(m + \Lambda(1 - m))m}$ , which is equivalent to (13).

Next, in the “no attention cost in welfare” case,  $A = 0$

$$\begin{aligned} \frac{\partial L(\tau, w)}{\partial \tau_1} &= \Lambda c_1 - \Psi \tau_1(m + \Lambda(1 - m))(m + \tau_1 m'(\tau_1)) \\ &= \Lambda c_1 - \frac{\psi}{1 + m\tau} c\tau(m + \Lambda(1 - m))(m + \tau m'(\tau)), \end{aligned}$$

so

$$\tau^0 = \frac{\Lambda/\psi}{m^2 + \tau m m'(\tau) + \Lambda \left( (1 - m)(m + \tau m') - \frac{m}{\psi} \right)}. \quad (79)$$

Finally, in the “optimally allocated attention” case,  $A = 1$ . First, we verify:

$$\begin{aligned}
-D_1 &= \tau^b \cdot \mathbf{S}^H = (0, \tau_1 (1 - m), g'(m)) \cdot (0, -\Psi (m + \tau_1 m'(p)), m'(p_1)) \\
&= -\Psi (m + \tau_1 m'(p)) \tau_1 (1 - m) + g'(m) m'(p_1) = -\Psi m \tau_1 (1 - m) = -(\tau_1 - \tau_1^s) \Psi m \\
&= \tau_C^s \cdot \mathbf{S}_{j|m}^H(\mathbf{p}, w, m).
\end{aligned}$$

with  $\tau_C^s = \tau_1 - \tau_1^s = (1 - m) \tau_1$  and  $\mathbf{S}_{j|m}^H(\mathbf{p}, w, m) = -\Psi m$ .

$$\begin{aligned}
\frac{\partial L(\boldsymbol{\tau}, w)}{\partial \tau_1} - \Lambda c_1 &= (\tau - \tilde{\tau}^b) \cdot \mathbf{S}_1^H = \boldsymbol{\tau} \cdot \mathbf{S}_1^H - \tilde{\tau}^b \cdot \mathbf{S}_1^H = \boldsymbol{\tau} \cdot \mathbf{S}_1^H - (1 - \Lambda) \tau_C^s \cdot \mathbf{S}_{1|m}^H \\
&= -\Psi \tau (m + \tau_1 m'(\tau_1)) + (1 - \Lambda) \Psi \tau (1 - m) m = -\Psi \tau (m (m + \Lambda (1 - m)) + \tau m'(\tau))
\end{aligned}$$

$$\frac{\partial L(\boldsymbol{\tau}, w)}{\partial \tau_1} = \Lambda c_1 - \frac{\psi c_1}{1 + m\tau} \tau (m (m + \Lambda (1 - m)) + \tau m'(\tau)),$$

so  $\frac{\partial L(\boldsymbol{\tau}, w)}{\partial \tau_1} = 0$  gives

$$\tau^1 = \frac{\Lambda/\psi}{m^2 + \tau m'(\tau) + \Lambda \left( (1 - m) m - \frac{m}{\psi} \right)}. \quad (80)$$

We can compare this to the following re-write of the optimal tax in the no-attention in welfare case:

$$\tau^0 = \frac{\Lambda/\psi}{m^2 + \tau_1 m'(\tau_1) + \Lambda \left( (1 - m) m - \frac{m}{\psi} \right) - \tau (1 - m) m' + \Lambda (1 - m) \tau m'} \quad (81)$$

$$= \frac{\Lambda/\psi}{m^2 + \tau_1 m'(\tau_1) + \Lambda \left( (1 - m) m - \frac{m}{\psi} \right) - (1 - \Lambda) \tau (1 - m) m'}. \quad (82)$$

**Proposition 9.19** *The optimal tax is lower in the “attention in welfare” case than in the “no attention in welfare” case.*

**Proof.** Suppose the opposite, i.e.  $\tau^1(m_1) \geq \tau^0(m_0)$ .

We observe that, for all  $m$ , (i)  $\tau^0(m) \geq \tau^1(m)$  (ii)  $\tau^1(m)$  is decreasing in  $m$ , and (iii)  $m(\tau)$  is weakly increasing in  $\tau$ ,

Then

$$\tau^1(m_1) \geq \tau^0(m_0) \geq \tau^1(m_0),$$

hence, as  $\tau^1$  is decreasing in  $m$ , we have  $m_1 < m_0$ . As  $m(\tau)$  is increasing, this implies  $\tau_1 < \tau_0$ . We’ve reached a contradiction.

There's a function  $m(\tau)$ ; it's inverse is  $\tau(m)$  we define

$$\tau^1(m) = \frac{\Lambda/\psi}{m^2 + \tau m'(\tau(m)) + \Lambda \left( (1-m)m - \frac{m}{\psi} \right)}.$$

For (ii) a sufficient condition is  $\psi \geq 1$  and  $\tau(m)m'(\tau)$  weakly increasing in  $\tau$ : then we have

$$m^2 + \tau_1 m'(\tau_1) + \Lambda \left( (1-m)m - \frac{m}{\psi} \right) = m^2(1-\Lambda) + \tau m' + m\Lambda \left( 1 - \frac{1}{\psi} \right)$$

increasing in  $m$ .

For (iii), the problem is

$$\max_m u(c(p+m\tau)) - (p+\tau)c(p+m\tau) - g(m)$$

$$g'(m) = (u'(c) - q)c'(q^s)\tau = (q^s - q)c'(q^s)\tau = -c'(q^s)(1-m)\tau^2.$$

In the isoelastic case,

$$f(m, \tau) = \psi(1+m\tau)^{-\psi-1}(1-m)\tau^2 - g'(m).$$

We have

$$f(m(\tau), \tau) = 0.$$

**Optimal tax with endogenous, optimally chosen attention** There is just one taxed good. The case with many, independent taxed goods follows.

Recall that the consumer chooses:  $m(\tau) = \arg \min \frac{1}{2}\psi y \tau^2 (1-m)^2 - \kappa g^1(m)$  and the planner chooses:  $\tau(\Lambda) = \arg \max_{\tau} L(\tau, \Lambda)$  with

$$L(\tau, \Lambda) = -\frac{1}{2}\psi y m(\tau)^2 \tau^2 - \kappa g(m(\tau)) + \Lambda y \tau.$$

**A lemma on scaling** We show that it is enough to compute the solution in the case  $\psi = y = \kappa = 1$ .

**Lemma 9.3** *Suppose that when  $\psi = y = \kappa = 1$  the optimal tax is  $\tau' = f(\Lambda)$  and optimal attention is  $m^1(\tau')$ . Then, in the general case it is:*

$$\tau(\Lambda) = \sqrt{\frac{\kappa}{\psi y}} f\left(\Lambda \sqrt{\frac{y}{\kappa \psi}}\right),$$

and the attention is  $m(\tau) = m^1\left(\tau \sqrt{\frac{\psi y}{\kappa}}\right)$ .



For instance, in the basic rational case,  $f(\Lambda) = \Lambda$  and  $m^1(\tau') = 1$ .

**Proof** This is a simple scaling argument. We define

$$\begin{aligned} m^1(\tau') &= \arg \min \frac{-1}{2} \tau'^2 (1 - m)^2 - g^1(m) \\ L^1(\tau', \Lambda') &= -\frac{1}{2} m^1(\tau')^2 \tau'^2 - g(m^1(\tau')) + \Lambda' \tau' \end{aligned}$$

$$\begin{aligned} \tau' &= \tau \sqrt{\frac{\psi y}{\kappa}} \\ \Lambda' &= \Lambda \sqrt{\frac{y}{\kappa \psi}} = \frac{\Lambda y \tau}{\kappa \tau'} \end{aligned}$$

Then, we have:

$$\begin{aligned} m(\tau) &= \arg \min \frac{-1}{2} \frac{\psi y}{\kappa} \tau^2 (1 - m)^2 - g^1(m) \\ &= m^1(\tau') \\ L(\tau, \Lambda) &= -\frac{1}{2} \psi y m(\tau)^2 \tau^2 - \kappa g(m) + \Lambda y \tau \\ &= \kappa \left[ -\frac{1}{2} \frac{\psi y}{\kappa} \tau^2 m(\tau)^2 - g(m) + \frac{\Lambda y}{\kappa} \tau \right] \\ &= \kappa L^1(\tau', \Lambda'). \end{aligned}$$

Hence, as  $\tau' = f(\Lambda')$  at the optimum.  $\square$

**Example with continuously adjusting attention** We have  $g(m) = -\kappa \ln(1 - m)$ , so that attention is  $m(\tau) = \left(1 - \frac{1}{\sqrt{\psi y \tau}}\right)_+$ . Indeed,  $\arg \min \frac{\sigma^2}{2} (1 - m)^2 + g(m)$  is  $m = \left(1 - \frac{1}{\sigma}\right)_+$ .

**Proposition 9.20** *In the above setup with optimal attention, the optimal tax is  $\tau_i = \sqrt{\frac{\kappa}{\psi_i y_i}} f\left(\Lambda \sqrt{\frac{y_i}{\kappa \psi_i}}\right)$ , for the continuous function*

$$\begin{aligned} f(\Lambda) &= \frac{\Lambda + 1 + \sqrt{(\Lambda + 1)^2 - 4}}{2} \text{ for } \Lambda \geq 1 \\ &= 1 \text{ for } \Lambda < 1. \end{aligned}$$

Also,  $m^1(\tau') = \left(1 - \frac{1}{\tau'}\right)_+$ .

**Proof.** We first reason in the case  $\psi = y = \kappa = 1$ . Then,  $m(\tau) = \left(1 - \frac{1}{\tau}\right)_+$  and

$$L(\tau) = -\frac{1}{2}m(\tau)^2\tau^2 - g(m) + \Lambda\tau.$$

Then, for  $\tau > 1$ ,

$$L'(\tau) = 1 - \frac{1}{\tau} - \tau + \Lambda,$$

so  $\tau$  is the greater root of:

$$\tau + \frac{1}{\tau} = \Lambda + 1,$$

which exists provided  $\Lambda \geq 1$ , i.e.:

$$\begin{aligned} f(\Lambda) &= \frac{\Lambda + 1 + \sqrt{(\Lambda + 1)^2 - 4}}{2} \text{ for } \Lambda \geq 1 \\ &= 1 \text{ for } \Lambda < 1. \end{aligned}$$

□

**An example with 0-1 attention** A concrete example of attention choice is:

$$m(\tau) = \arg \max_m -\frac{1}{2}\psi\tau^2(1-m)^2 - g(m),$$

with

$$g(m) = \frac{1}{2}\kappa^2 [1 - (1-m)^2].$$

Then, the solution is

$$m(\tau) = 1_{\tau > \tau_*}, \quad \tau_* = \frac{\kappa}{\sqrt{\psi}}. \tag{83}$$

As an aside, a fixed cost  $g(m) = \frac{\kappa^2}{2}1_{m>0}$  gives the same result.

**Proposition 9.21** *The optimal tax is  $\tau_i = \sqrt{\frac{\kappa}{\psi_i y_i}} f\left(\Lambda \sqrt{\frac{y_i}{\kappa \psi_i}}\right)$ , for  $f(\Lambda) = 1$  if  $\Lambda \leq \sqrt{2} + 1$  and  $f(\Lambda) = \Lambda$  if  $\Lambda > \sqrt{2} + 1$ . Also,  $m^1(\tau') = 1_{\tau' > 1}$ .*

In that case, the optimal tax has a discontinuity. When  $\Lambda$  is low enough, the planner keeps the taxes at  $\tau_i = \sqrt{\frac{\kappa}{\psi_i y_i}}$ , just below the “detectability threshold” and agents do not pay attention to the tax.

**Proof** We start with the case  $\psi = y = \kappa = 1$ . Then,  $m(\tau) = 1_{\tau > 1}$ . For  $\tau \leq 1$ ,  $L(\tau) = \Lambda\tau$ , so the optimum for  $\tau \in [0, 1]$  is  $\tau = 1$ .

$$L(1) = \Lambda.$$

For  $\tau > \tau_*$ ,  $m(\tau) = 1$ , so  $L(\tau) = -\frac{1}{2}\tau^2 - g(1) + \Lambda\tau$ , and the optimum is  $\tau = \Lambda$ . We have

$$L(\tau) = \frac{\Lambda^2 - 1}{2}.$$

So  $L(\tau) > L(\tau_*)$  if and only if  $\frac{\Lambda^2 - 1}{2} > \Lambda$ , i.e. if and only if  $\frac{\Lambda^2 - 1}{2} > \Lambda$ , i.e. if and only if  $\Lambda > \sqrt{2} + 1$ .  $\square$

### 9.5.6 Quadratic losses from imperfect tax instruments

We introduce the Lagrangian that allows for agent-specific lump-sum transfers  $w^h$  and taxes  $\tau^h, \tau^{s,h}$  (we normalize  $p_i = 1$ )

$$L(\{\tau^h\}, \{\tau^{s,h}\}, \{w^h\}) = W(v^h(p + \tau, p + \tau^{s,h}, w^h, \xi)) + \lambda \sum_h [\tau \cdot c^h(p + \tau, p + \tau^{s,h}, w^h, \xi) - w^h],$$

with  $\xi = \xi(\{c^h\})$  as a fixed point. We also define:

$$g(\{\tau^{s,h}\}) = \max_{\{\tau^h\}} L(\{\tau^{s,h}\}, \{\tau^h\}, \{w^h\}), \quad (84)$$

which is the Lagrangian with rational agents perceiving  $\tau^{s,h}$  and with optimum agent-specific lump-sum transfer.

The social utility achieved with agent-specific taxes  $\tau^{s,h}$ , and optimum agent-specific lump-sum transfers, with a rational agent.

**Proposition 9.22** *In general, in the Ramsey problem with externalities and redistribution, the social loss (realized social minus first best) is:*

$$L = L^{distribution} + L^{distortions},$$

with

$$L^{distribution} = \frac{1}{2} \sum_{h,h'} (\gamma^{\xi,h} - \bar{\gamma}) (L_{ww}(w, \tau)^{-1})_{h,h'} (\gamma^{\xi,h} - \bar{\gamma})$$

$$L^{distortions} = \frac{1}{2} \sum_{h,h'} (\tau^{s,h} - \tau^{*s,h}) g_{\tau^{s,h}\tau^{s,h'}} (\tau^{s,h'} - \tau^{*s,h'}).$$

This reflects that at the optimum, the  $\gamma^{\xi,h}$  should be the same (and equal to  $\lambda$ ), and we should have  $\tau^{s,h} - \tau^{*s,h}$ .

**Understanding the redistribution term** For instance, take the case:  $W = \sum v^h (q, q^{s,h}, w^h, \xi)$  and  $\xi$  is independent of  $w^h$ , then  $L_{w^h w^{h'}} = v_{ww}^h$ , so that

$$L^{\text{distribution}} = \frac{1}{2} \sum_{h,h'} \frac{(\gamma^h - \bar{\gamma})^2}{v_{ww}^h}.$$

The losses come from the lack of equalization of  $\gamma$ 's.

### Understanding the $g_{\tau^s \tau^s}$ better

**Lemma 9.4** *We have*

$$g_{\tau^s, h \tau_s, h'} = \lambda S^{hr} \left( 1_{h=h'} - \frac{d\tau^{\xi, h}}{d\tau^s, h'} \right).$$

When utility is quasi linear and the externality enters additively,  $u(c, \xi) = u(c_1, \dots, c_n) + \lambda c_0 + \frac{1}{H} \xi$ , we have:

$$g_{\tau_s^h \tau_s, h'} = \lambda S^{rh} 1_{h=h'} + S^{rh} \xi_{c^h c^{h'}} S^{rh'}. \quad (85)$$

## 10 The Nonlinear income tax problem

Here are the notations we shall use.

$g(z)$  :social welfare weight

$h(z)$  (resp.  $h^*(z)$ ): density (resp. virtual density) of earnings  $z$

$H(z)$ : cumulative distribution function of earnings

$n$  :agent's wage, also the index of his type

$q(z) = R'(z)$ : marginal retention rate, locally perceived

$\mathbf{Q} = (q(z))_{z \geq 0}$ : vector of marginal retention rates

$r_0$ : tax rebate at 0 income

$r(z)$  :virtual income

$R(z) = z - T(z)$ : retained earnings

$T(z)$ : tax given earnings  $z$

$z$ : pre-tax earnings

$\gamma(z)$ : marginal social utility of income

$\eta$ : income elasticity of earnings

$\pi$  :Pareto exponent of the earnings distribution

$\zeta^c$ : compensated elasticity of earnings

$\zeta_{Q_{z^*}}^c(z)$ : compensated elasticity of earnings when the tax rate at  $z^*$  changes.

$\zeta^u$  : uncompensated elasticity of earnings

## 10.1 Setup

**Agent’s behavior** There is a continuum of agents indexed by skill  $n$  with density  $f(n)$  (we use  $n$  rather than  $h$ , the conventional index in that literature). Agent  $n$  has a utility function  $u^n(c, z)$ , where  $c$  is his one-dimensional consumption,  $z$  is his pre-tax income, and  $u_z \leq 0$ .<sup>85</sup>

The total income tax for income  $z$  is  $T(z)$ , so that disposable income is  $R(z) = z - T(z)$ . We call  $q(z) = R'(z) = 1 - T'(z)$  the local marginal “retention rate”,  $\mathbf{Q} = (q(z))_{z \geq 0}$  the ambient vector of all marginal retention rates, and  $r_0 = R(0)$  the transfer given by the government to an agent earning zero income. We define the “virtual income” to be  $r(z) = R(z) - zq(z)$ . Equivalently  $R(z) = q(z)z + r(z)$ , so that  $q(z)$  is the local slope of the budget constraint, and  $r(z)$  its intercept.

We use a general behavioral model in a similar spirit to Section 2. The primitive is the income function  $z^n(q, \mathbf{Q}, r_0, r)$ , which depends on the local marginal retention rate  $q$ , the ambient vector of all marginal retention rates  $\mathbf{Q}$ ,  $r_0 = R(0)$  the transfer given by the government to an agent earning zero income, and the virtual income  $r$ . In the traditional model without behavioral biases we have  $z^n(q, \mathbf{Q}, r_0, r) = \arg \max_z u^n(qz + r, z)$ , so that  $z^n$  does not depend on  $\mathbf{Q}$  and  $r_0$ . With behavioral biases, this is no longer true in general. The income function is associated with the indirect utility function  $v^n(q, \mathbf{Q}, r_0, r) = u^n(qz + r, z)|_{z=z^n(q, \mathbf{Q}, r_0, r)}$ . The earnings  $z(n)$  of agent  $n$  facing retention schedule  $R(z)$  is then the solution of the fixed point problem  $z = z^n(q(z), \mathbf{Q}, r_0, r(z))$ . His consumption is  $c(n) = R(z(n))$  and his utility is  $v(n) = u^n(c(n), z(n))$ .

**Planning problem** The objective of the planner is to design the tax schedule  $T(z)$  in order to maximize the following objective function

$$\int_0^\infty W(v(n)) f(n) dn + \lambda \int_0^\infty (z(n) - c(n)) f(n) dn.$$

Like Saez (2001), we normalize  $\lambda = 1$ . We call  $g(n) = W'(v(n)) v_r^n(q(z(n)), \mathbf{Q}, r_0, r(z(n)))$  the marginal utility of income. This is the analogue of  $\beta^h$  in the Ramsey problem of Section 2, and we identify agents with their income level  $z(n)$  instead of their skill  $n$ . Most of the time, we leave implicit the dependence of  $n(z)$  on  $z$  to avoid cluttering the notations. We now derive a behavioral version of the optimal tax formula in Saez (2001).

---

<sup>85</sup>If the agent’s pre-tax wage is  $n$ ,  $L$  is his labor supply, and utility is  $U^n(c, L)$ , then  $u^n(c, z) = U(c, \frac{z}{n})$ . Note that this assumes that the wage is constant (normalized to one). We discuss the impact of relaxing this assumption in Sections 5.1 and 10.3.2.

## 10.2 Saez Income Tax Formula with Behavioral Agents

### 10.2.1 Elasticity Concepts

Recall that the marginal retention rate is  $q(z) = 1 - T'(z)$ . Given an income function  $z(q, \mathbf{Q}, r_0, r)$ , we introduce the following definitions. We define the income elasticity of earnings

$$\eta = qz_r(q, \mathbf{Q}, r_0, r).$$

We also define the uncompensated elasticity of labor (or earnings) supply with respect to the actual marginal retention rate

$$\zeta^u = \frac{q}{z} z_q(q, \mathbf{Q}, r_0, r).$$

Finally, we define the compensated elasticity of labor supply with respect to the actual marginal retention rate

$$\zeta^c = \zeta^u - \eta.$$

We also introduce two other elasticities, which are zero in the traditional model without behavioral biases. We define the compensated elasticity of labor supply at  $z$  with respect to the marginal retention rate  $q(z^*)$  at a point  $z^*$  different from  $z$ :

$$\zeta_{Q_{z^*}}^c = \frac{q}{z} z_{Q_{z^*}}(q, \mathbf{Q}, r_0, r).$$

We also define the earnings sensitivity to the lump-sum rebate at zero income<sup>86</sup>

$$\zeta_{r_0}^c = \frac{q}{z} z_{r_0}(q, \mathbf{Q}, r_0, r).$$

We shall call  $\zeta_{Q_{z^*}}^c$  a “behavioral cross-influence” of the marginal tax rate at  $z^*$  on the decision of an agent earning  $z$ . In the traditional model with no behavioral biases,  $\zeta_{Q_{z^*}}^c = \zeta_{r_0}^c = 0$ , not so with behavioral agents.<sup>87’88</sup>

All these elasticities a priori depend on the agent earnings  $z$ . As mentioned above, we leave this dependence implicit most of the time.

Just like in the Ramsey model, we define the “behavioral wedge”

$$\tau^b(q, \mathbf{Q}, r_0, r) = - \frac{qu_c(c, z) + u_z(c, z)}{v_r(q, \mathbf{Q}, r_0, r)} \Big|_{z=z(q, \mathbf{Q}, r_0, r), c=qz+r}.$$

<sup>86</sup>Formulas would be cleaner without the multiplication by  $q$  in those elasticities, but here we follow the public economics tradition.

<sup>87</sup>For instance, in the misperception model, in general, the marginal tax rate at  $z^*$  affects the default tax rate and therefore the perceived tax rate at earnings  $z$ .

<sup>88</sup>In the language of Section 2.1, we use income-compensation based notion of elasticity,  $\mathbf{S}^C$ , rather than the utility-compensation based notion  $\mathbf{S}^H$ .

We also define the renormalized behavioral wedge

$$\tilde{\tau}^b(z) = g(z) \tau^b(z).$$

In the traditional model with no behavioral biases, we have  $\tau^b(q, \mathbf{Q}, r_0, r) = \tilde{\tau}^b(z) = 0$ . But this is no longer true with behavioral agents.

We have the following behavioral version of Roy's identity (proven in the online appendix, Section 11.2.2):

$$\frac{v_q}{v_w} = z - \frac{\tau^b z}{q} \zeta^c, \quad \frac{v_{Q_{z^*}}}{v_w} = -\frac{\tau^b z}{q} \zeta_{Q_{z^*}}^c. \quad (86)$$

As in Section 2, the general model can be particularized to a decision vs. experienced utility model, or to a misperception model.

**Misperception model** The agent may misperceive the tax schedule, including her marginal tax rate. We call  $T^{s,n}(q, \mathbf{Q}, r_0)(z)$  the perceived tax schedule,  $R^{s,n}(z) = z - T^{s,n}(q, \mathbf{Q}, r_0)(z)$  the perceived retention schedule, and  $q^{s,n}(q, \mathbf{Q}, r_0)(z) = \frac{dR^{s,n}(q, \mathbf{Q}, r_0)(z)}{dz}$  the perceived marginal retention rate. Faced with this tax schedule, the behavior of the agent can be represented by the following problem

$$\text{smax}_{c, z | R^{s,n}(\cdot)} u^n(c, z) \text{ s.t. } c = R(z).$$

This formulation implies that the agent's choice  $(c, z)$  satisfies  $c = R(z)$  and

$$q^{s,n}(z) u_c^n(c, z) + u_z^n(c, z) = 0,$$

instead of the traditional condition  $q(z) u_c^n(c, z) + u_z^n(c, z) = 0$ . This means that the agent correctly perceives consumption and income  $(c, z)$  but misperceives his marginal retention rate  $q^{s,n}(z)$ . Together with  $c = R(z)$ , this characterizes the behavior of the agent.<sup>89</sup>

Accordingly, we define  $z^n(q, q^s, r)$  to be the solution of  $q^{s,n} u_c^n(c, z) + u_z^n(c, z) = 0$  with  $c = qz + r$ .<sup>90</sup> The income  $z(n)$  of agent  $n$  is then the solution of the fixed point equation

$$z = z^n(q(z), q^{n,s}(q, \mathbf{Q}, r_0)(z), r(z)),$$

his consumption is  $c(n) = R(z(n))$  and his utility is  $v(n) = u^n(c(n), z(n))$ .

Summing up, in the misperception model, the primitives are a utility function  $u$  and a perception function  $q^s(q, \mathbf{Q}, r_0)(z)$ . This yields an income function  $z(q, q^s, r)$ . The general function  $z(q, \mathbf{Q}, r_0, r)$  is then  $z(q(z'), \mathbf{Q}, r_0, r) = z(q(z'), q^s(q, \mathbf{Q}, r_0)(z'), r)$  for any earnings  $z'$ .

<sup>89</sup>This is a sparse max problem with a non-linear budget constraint, which generalizes the sparse max with a linear budget constraint we analyzed in section 3.1. The true constraint is  $c = R(z)$ , but the perceived constraint is  $c = R^{s,n}(q, \mathbf{Q}, r_0)(z)$ .

<sup>90</sup>If there are several solutions, we choose the one that yields the greatest utility.

One concrete example of misperception is  $q^{s,n}(q, \mathbf{Q}, r_0) = q^s(q, \mathbf{Q}, r_0)$  with

$$q^s(q, \mathbf{Q}, r_0)(z) = mq(z) + (1 - m) \left[ \alpha q^d(\mathbf{Q}) + (1 - \alpha) \frac{r_0 + \int_0^z q(z') dz'}{z} \right],$$

where  $m \in [0, 1]$  is the attention to the true tax (hence retention) rate,  $\frac{r_0 + \int_0^z q(z') dz'}{z}$  is the average retention rate (as in [Liebman and Zeckhauser \(2004\)](#)), and  $\alpha \in [0, 1]$ . The default perceived retention rate might be a weighted average of marginal rates, e.g.  $q^d(\mathbf{Q}) = \int q(z) \omega(z) dz$  for some weights  $\omega(z)$ .

As in the Ramsey case, it is useful to express behavioral elasticities as a function of an agent without behavioral biases. Call  $z^r(q^s, r') = \arg \max_z u(q^s z + r', z)$  the earnings of a rational agent facing marginal tax rate  $q^s$  and extra non-labor income  $r'$ . Then,  $z(q, q^s, r) = z^r(q^s, r')$  where  $r'$  solves  $r' + q^s z^r(q^s, r') = r + q z^r(q^s, r')$ . We call  $S^r(q^s, r') = \frac{\partial z^r}{\partial q^s}(q^s, r') - \frac{\partial z^r}{\partial r'}(q^s, r') z^r(q^s, r')$  the rational compensated sensitivity of labor supply (it is just a scalar). We also define  $\zeta^{cr} = \frac{q S^r}{z}$  as the compensated elasticity of labor supply of the agent if he were rational.

We define  $m_{zz} = q_q^s(q, \mathbf{Q}, r_0)(z)$  as the attention to the own marginal retention rate and  $m_{zz^*} = q_{Q_{z^*}}^s(q, \mathbf{Q}, r_0)(z)$  as the marginal impact on the perceived marginal retention rate at  $z$  of an increase in the marginal retention rate at  $z^*$ . Then, we have the following concrete values for the elasticities of the general model (the derivation is in Section 11.2.2 of the appendix):

$$\zeta^c = \left( 1 - \eta \frac{\tau - \tau^s}{q} \right) \zeta^{cr} m_{zz}, \quad \zeta_{Q_{z^*}}^c = \left( 1 - \eta \frac{\tau - \tau^s}{q} \right) \zeta^{cr} m_{zz^*}, \quad (87)$$

$$\tau^b = \frac{\tau - \tau^s}{1 - \eta \frac{\tau - \tau^s}{q}}. \quad (88)$$

If the behavioral agent overestimates the tax rate ( $\tau - \tau^s < 0$ ), the term  $\tau^b$  is negative. Loosely, we can think of  $\tau^b$  as indexing an “underperception” of the marginal tax rate. In the traditional model without behavioral biases,  $m_{zz^*} = 1_{z=z^*}$ ,  $\tau^s = \tau$  and  $\tau^b = 0$ .

**Decision vs. experienced utility model** In the decision vs. experienced utility model, behavior is represented by the maximization of a subjective decision utility  $u^s(c, z)$  subject to the budget constraint  $c = R(z)$ . We then have  $\zeta_{Q_{z^*}}^c = 0$ , and  $\zeta^c$  and  $\eta$  are the elasticities associated with decision utility  $u^s$ . The behavioral wedge is

$$\tau^b = \frac{\frac{u_c}{u_c^s} u_z^s - u_z}{v_r}. \quad (89)$$

**Other useful concepts and notations** We next study the impact of the above changes on welfare. Following [Saez \(2001\)](#), we call  $h(z)$  the density of agents with earnings  $z$  at the optimum and  $H(z) = \int_0^z h(z') dz'$ . We also introduce the virtual density  $h^*(z) = \frac{q(z)}{q(z) - \zeta^c z R''(z)} h(z)$ .



We define the social marginal utility of income

$$\gamma(z) = g(z) + \frac{\eta(z)}{1 - T'(z)} \left[ \tilde{\tau}^b(z) + (T'(z) - \tilde{\tau}^b(z)) \frac{h^*(z)}{h(z)} \right]. \quad (90)$$

This definition is the analogue of the corresponding definition in the Ramsey model. It is motivated by Lemma 11.2 in the online appendix, which shows that, if the government transfers a lump-sum  $\delta K$  to an agent previously earning  $z$ , the objective function of the government increases by  $\delta L(z) = (\gamma(z) - 1) \delta K$ . The social marginal utility of income  $\gamma(z)$  reflects a direct effect  $g(z)$  of that transfer to the agent's welfare, and an indirect effect on labor supply captured—to the leading order as the agent receives  $\delta K$ , his labor supply changes by  $\frac{\eta(z)}{1 - T'(z)} \delta K$ , which impacts tax revenues by  $\frac{\eta(z)}{1 - T'(z)} T'(z) \delta K$  and welfare by  $\frac{\eta(z)}{1 - T'(z)} \tilde{\tau}^b(z) \delta K$ ; the terms featuring  $\frac{h^*(z)}{h(z)}$  (in practice often close to 1) capture the fact that the agent's marginal tax rate changes as the agent adjusts his labor supply, which impacts tax revenues and welfare because misoptimization.

## 10.2.2 Optimal Income Tax Formula

We next present the optimal income tax formula. The online appendix (section 11.2.1) presents the intermediary steps used in the derivation of this formula.

**Proposition 10.1** *Optimal taxes satisfy the following formulas (for all  $z^*$ )*

$$\begin{aligned} \frac{T'(z^*) - \tilde{\tau}^b(z^*)}{1 - T'(z^*)} &= \frac{1}{\zeta^c(z^*)} \frac{1 - H(z^*)}{z^* h^*(z^*)} \int_{z^*}^{\infty} (1 - \gamma(z)) \frac{h(z)}{1 - H(z^*)} dz \\ &\quad - \int_0^{\infty} \frac{\zeta_{Q_{z^*}}^c(z)}{\zeta^c(z^*)} \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} \frac{zh^*(z)}{z^* h^*(z^*)} dz. \end{aligned} \quad (91)$$

*This formula can also be expressed as a modification of the Saez (2001) formula*

$$\begin{aligned} \frac{T'(z^*) - \tilde{\tau}^b(z^*)}{1 - T'(z^*)} &+ \int_0^{\infty} \omega(z^*, z) \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} dz \\ &= \frac{1}{\zeta^c(z^*)} \frac{1 - H(z^*)}{z^* h^*(z^*)} \int_{z^*}^{\infty} e^{-\int_{z^*}^z \rho(s) ds} \left( 1 - g(z) - \eta \frac{\tilde{\tau}^b(z)}{1 - T'(z)} \right) \frac{h(z)}{1 - H(z^*)} dz, \end{aligned} \quad (92)$$

where  $\rho(z) = \frac{\eta(z)}{\zeta^c(z)} \frac{1}{z}$  and

$$\omega(z^*, z) = \left( \frac{\zeta_{Q_{z^*}}^c(z)}{\zeta^c(z^*)} - \int_{z'=z^*}^{\infty} e^{-\int_{z^*}^{z'} \rho(s) ds} \rho(z') \frac{\zeta_{Q_{z'}}^c(z)}{\zeta^c(z^*)} dz' \right) \frac{zh^*(z)}{z^* h^*(z^*)}.$$

The first term  $\frac{1}{\zeta^c(z^*)} \frac{1 - H(z^*)}{z^* h^*(z^*)} \int_{z^*}^{\infty} (1 - \gamma(z)) \frac{h(z)}{1 - H(z^*)} dz$  on the right-hand side of the optimal tax formula (91) is a simple reformulation of Saez's formula, using the concept of social marginal utility of income  $\gamma(z)$  rather than the marginal social welfare weight  $g(z)$ . The link between the two is

in equation (90)). The second term  $-\frac{1}{z^*} \int_0^\infty \frac{\zeta_{Q_{z^*}}^c(z)}{\zeta^c(z^*)} \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} z \frac{h^*(z)}{h^*(z^*)} dz$  on the right-hand side is new and captures a misoptimization effect together with the term  $\frac{-\tilde{\tau}^b(z^*)}{1 - T'(z^*)}$  on the left-hand side.

The intuition is as follows. First, suppose for concreteness that  $\zeta_{Q_{z^*}}^c(z) > 0$ , then increasing the marginal tax rate at  $z^*$  leads the agents at another income  $z$  to perceive higher taxes on average, which leads them to decrease their labor supply and reduces tax revenues. Ceteris paribus, this consideration pushes towards a lower tax rate, compared to the Saez optimal tax formula. Second, suppose for concreteness that  $\tilde{\tau}^b(z) < 0$ , then increasing the marginal tax rate at  $z^*$  further reduces welfare. This, again, pushes towards a lower tax rate.

The modified Saez formula (92) uses the concept of the social marginal welfare weight  $g(z)$  rather than the social marginal utility of income  $\gamma(z)$ . It is easily obtained from formula (91) using equation (90). When there are no income effects so that  $\eta = \rho(z) = 0$ , the optimal tax formula (91) and the modified Saez formula (92) are identical. They coincide with the traditional Saez formula when there are no behavioral biases so that  $\zeta_{Q_{z^*}}^c(z) = \omega(z^*, z) = \tilde{\tau}^b(z) = 0$ . In this case, the left-hand side of (92) is simply  $\frac{T'(z^*)}{1 - T'(z^*)}$  so that the formula solves for the optimal marginal tax rate  $T'(z^*)$  at  $z^*$ .

The formula is expressed in terms of endogenous objects or “sufficient statistics”: social marginal welfare weights  $g(z)$ , elasticities of substitution  $\zeta^c(z)$ , income elasticities  $\eta(z)$ , and income distribution  $h(z)$  and  $h^*(z)$ . With behavioral agents, there are two differences. First, there are two additional sufficient statistic, namely the behavioral wedge  $\tilde{\tau}^b(z)$  and the behavioral cross-elasticities  $\zeta_{Q_{z^*}}^c(z)$ . Second, it is not possible to solve out the optimal marginal tax rate in closed form. Instead, the modified Saez formula (92) at different values of  $z^*$  form a system of linear equations in the optimal marginal tax rates  $T'(z)$  for all  $z$ . The formula simplifies greatly in the case where behavioral biases can be represented by a decision vs. experienced utility model. Indeed, we then have  $\omega(z^*, z) = 0$  and  $\tilde{\tau}^b(z) = g(z) \frac{u_c \frac{u_z}{u_c} - u_z}{v_r}$ , so that there is no linear system of equations to solve out to recover  $T'(z)$ .

### 10.2.3 Marginal Tax Rate for Top Incomes

We start by revisiting the classic result that if the income distribution is bounded at  $z_{\max}$ , then the top marginal income tax rate should be zero. In our model, this needs not be the case. One simple way to see that is to consider the case of decision vs. experienced utility. The tax formula (91) prescribes  $T'(z_{\max}) = \tilde{\tau}^b(z_{\max})$  which is positive or negative depending on whether top earners overperceive or underperceive the benefits of work (underperceive or overperceive the costs of work).

**Proof of Proposition 4.2** For high incomes,  $\tilde{\tau}^b$  approaches zero (as high-earnings agents asymptotically accurately perceive their marginal rate, see (27) and (87)), and  $\frac{h^*(z)}{h(z)}$  approaches 1.<sup>91</sup> Hence,

---

<sup>91</sup>Recall that  $\frac{h^*(z)}{h(z)} = \frac{q(z)}{q(z) + \zeta^c z T''(z)}$ . Calling  $a = \lim_{z \rightarrow \infty} z T''(z)$ , if we had  $a \neq 0$ , we'd have  $T'(z) \sim a \ln z$ , which would diverge for large  $z$ . So  $a = 0$ . Hence, for large  $z$ ,  $\frac{h^*(z)}{h(z)} \rightarrow 1$ .

from (90), we have

$$\bar{\gamma} = \bar{g} + \bar{\eta}^r \frac{\bar{\tau}}{1 - \bar{\tau}}.$$

We observe that  $\frac{1-H(z^*)}{z^*h^*(z^*)} \rightarrow \frac{1}{\pi}$ , that  $\zeta^c(z^*) = m\zeta^{c,r}(z^*)$ , and  $\zeta_{Q_{z^*}}^c(z) = (1-m)\frac{\psi(z^*/z)}{z}\zeta^{c,r}(z)$  (see (87)). Taking the large  $z^*$  limit in (26) gives:

$$\frac{\bar{\tau}}{1 - \bar{\tau}} = \frac{1}{m\bar{\zeta}^{c,r}} \frac{1}{\pi} (1 - \bar{\gamma}) - \lim_{z^* \rightarrow \infty} C(z^*),$$

where

$$\begin{aligned} C(z^*) &:= \int_0^\infty \frac{\zeta_{z^*}^c(z)}{\zeta^c(z^*)} \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} \frac{zh^*(z)}{z^*h^*(z^*)} dz = \frac{1-m}{m} \int_0^\infty \frac{\frac{\psi(z^*/z)}{z}\zeta^{c,r}(z)}{\zeta^{c,r}(z^*)} \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} \frac{zh^*(z)}{z^*h^*(z^*)} dz \\ &= \frac{1-m}{m} \int_0^\infty \psi(a) \frac{\zeta^{c,r}(\frac{z^*}{a})}{\zeta^{c,r}(z^*)} \frac{T'(\frac{z^*}{a}) - \tilde{\tau}^b(\frac{z^*}{a})}{1 - T'(\frac{z^*}{a})} \frac{\frac{z^*}{a}h^*(\frac{z^*}{a})}{z^*h^*(z^*)} \frac{da}{a}, \end{aligned}$$

where we used the change in variables  $z = \frac{z^*}{a}$ , so  $\frac{da}{a} = -\frac{dz}{z}$ . Observing that  $\frac{T'(\frac{z^*}{a}) - \tilde{\tau}^b(\frac{z^*}{a})}{1 - T'(\frac{z^*}{a})} \rightarrow \frac{\bar{\tau}}{1 - \bar{\tau}}$  and  $\frac{\frac{z^*}{a}h^*(\frac{z^*}{a})}{z^*h^*(z^*)} \rightarrow a^\pi$ , we get

$$\lim_{z^* \rightarrow \infty} C(z^*) = \frac{1-m}{m} \int_0^\infty \psi(a) \frac{\bar{\tau}}{1 - \bar{\tau}} a^\pi \frac{da}{a} = \frac{1-m}{m} \frac{\bar{\tau}}{1 - \bar{\tau}} A$$

where  $A = \int_0^\infty a^{\pi-1} \psi(a) da$ . So

$$\frac{\bar{\tau}}{1 - \bar{\tau}} = \frac{1}{m\bar{\zeta}^c\pi} \left( 1 - \bar{g} - \bar{\eta}^r \frac{\bar{\tau}}{1 - \bar{\tau}} \right) - \frac{1-m}{m} \frac{\bar{\tau}}{1 - \bar{\tau}} A.$$

Rearranging we get the top marginal rate announced,

$$\bar{\tau} = \frac{1 - \bar{g}}{1 - \bar{g} + \bar{\eta}^r + \bar{\zeta}^{c,r}\pi(m + (1-m)A)}.$$

#### 10.2.4 Possibility of Negative Marginal Income Tax Rates

In the traditional model with no behavioral biases, negative marginal income tax rates can never arise at the optimum. With behavioral biases negative marginal income tax rates are possible at the optimum. To see this, consider for example the decision vs. experienced utility model with decision utility  $u^s$  and assume that  $u^s$  is quasilinear so that there are no income effects  $u^s(c, z) = c - \phi z$ . We take experienced utility to be  $u(c, z) = \theta c - \phi(z)$ . Then the modified Saez formula (92) becomes

$$\frac{T'(z^*) - \tilde{\tau}^b(z^*)}{1 - T'(z^*)} = \frac{1}{\zeta^c(z^*)} \frac{1 - H(z^*)}{z^*h^*(z^*)} \int_{z^*}^\infty (1 - g(z)) \frac{h(z)}{1 - H(z^*)} dz,$$

where  $\tilde{\tau}^b(z) = -g(z)\phi'(z)\frac{\theta-1}{\theta}$  by (89). When  $\theta > 1$ , we have  $\tilde{\tau}^b(z^*) < 0$ , and it is possible for this formula to yield  $T'(z^*) < 0$ . This occurs if agents undervalue the benefits or overvalue the costs from higher labor supply. For example, it could be the case that working more leads to higher human capital accumulation and higher future wages, but that these benefits are underperceived by agents, which could be captured in reduced form by  $\theta > 1$ . Such biases could be particularly relevant at the bottom of the income distribution (see [Chetty and Saez \(2013\)](#) for a review of the evidence). If these biases are strong enough, the modified Saez formula could predict negative marginal income tax rates at the bottom of the income distribution. This could provide a behavioral rationale for the EITC program.<sup>92</sup> In parallel and independent work, [Gerritsen \(2016\)](#) and [Lockwood \(2017\)](#) derive a modified Saez formula in the context of decision vs. experienced utility model. [Lockwood \(2017\)](#) zooms in on the EITC program and provides an empirical analysis documenting significant present-bias among EITC recipients and shows that a calibrated version of the model goes a long way towards rationalizing the negative marginal tax rates associated with the EITC program.

This differs from alternative rationales for negative marginal income tax rates that have been put forth in the traditional literature. For example, [Saez \(2002\)](#) shows that if the Mirrlees model is extended to allow for an extensive margin of labor supply, then negative marginal income tax rates can arise at the optimum. We refer the reader to the online appendix (section 10.3.1) for a behavioral treatment of the [Saez \(2002\)](#) extensive margin of labor supply model.

## 10.3 Complements on the Mirrlees problem

### 10.3.1 Mirrlees problem with extensive margin

We provide a behavioral enrichment to [Saez \(2002\)](#). We take his simplest framework (Proposition 1). Activity 0 is unemployment, and there are  $I$  other activities. One type  $i$  of agent chooses between working and not working: working gives utility  $u^h(c_i, i)$ , not working utility  $u^h(c_0, 0)$ , where  $c_i = z_i - T(z_i)$ . If the agent is rational, he solves

$$i^* = \arg \max_{i^* \in \{0, i\}} u^h(c_i, i),$$

but our behavioral agent may make a mistake. E.g., in the misperception model, he might perceive  $c_i^s$ , so that he decides according to

$$i^* = \arg \max_{i^* \in \{0, i\}} u^{h,b}(c_i^s, i).$$

In general, we will simply model the choice as some  $i^*(h, \{T_j\})$ . We say that an agent is “at the margin for tax  $i$ ” if the agent changes activity as tax  $i$  changes  $B_i^+ = \{m \text{ s.t. } \partial i^*(h, \{T_j\}) / \partial T_i < 0\}$  (which is the set of agent moving into active employment if the tax rate on activity  $i$  falls) and

---

<sup>92</sup>The EITC program itself could be misperceived, see [Chetty, Friedman and Saez \(2013\)](#).

$B_i^- = \{m \text{ s.t. } \partial i^* (h, \{T_j\}) / \partial T_i > 0\}$  (which is the set of agent moving out of active employment if the tax rate on activity  $i$  falls). The normal case is that  $B_i^-$  is an empty set. The derivative is in the sense of distributions, and simply indicates a change in agent's behavior.

Suppose that the government increases tax  $T_i$  on activity  $i$  by  $dT_i$ . That induces a quantity  $dH_j$  of people to switch to employment, where

$$dH_j = -H_j \eta_{ji} \frac{dT_i}{c_i - c_0}.$$

We have  $h_j(\{T_k\}) =$  number of agents of type  $j$  who work.

Each  $h$  has a potential earnings level  $j(h)$ . We call

$$\tau_{ji}^b = - \sum_{\varepsilon \in \{-, +\}} \mathbb{E} [\varepsilon \mu^m (u^h(c_j, j) - u^h(c_0, 0)) \mid j(h) = j \text{ and } h \in B_i^\varepsilon].$$

We have

$$\frac{\partial H^j}{\partial T_i} = - \sum_{\varepsilon \in \{-, +\}} \int \varepsilon 1_{\{j(h)=j \text{ and } h \in B_i^\varepsilon\}} d\nu(m).$$

The change in welfare from  $dT_i$  is then

$$\begin{aligned} dL &= (1 - g_i) H_i dT_i - \sum_j \sum_{\varepsilon \in \{-, +\}} \int \varepsilon [T_j - T_0 + \mu^m (u^h(c_j, j) - u^h(c_0, 0))] 1_{\{j(h)=j \text{ and } h \in B_i^\varepsilon\}} d\nu(m) \\ &= (1 - g_i) H_i dT_i + \sum_j (T_j - T_0) \frac{\partial H^j}{\partial T_i} - \sum_j \sum_{\varepsilon \in \{-, +\}} \int \varepsilon [\mu^m (u^h(c_j, j) - u^h(c_0, 0))] 1_{\{j(h)=j \text{ and } h \in B_i^\varepsilon\}} d\nu(m) \\ &= (1 - g_i) H_i dT_i + \sum_j (T_j - T_0 - \tau_{ji}^b) \frac{\partial H^j}{\partial T_i} \\ &= (1 - g_i) H_i dT_i - \sum_j (T_j - T_0 - \tau_{ji}^b) H_j \eta_{ji} \frac{dT_i}{c_i - c_0}. \end{aligned}$$

Hence, at the optimum:

$$\sum_j \frac{T_j - T_0 - \tau_{ji}^b}{c_i - c_0} \frac{H_j}{H_i} \eta_{ji} = (1 - g_i).$$

For instance, suppose that people overestimate taxes, i.e. underperceive the benefits from working:  $c_i^s < c_i$ , and no cross-effects. Then,

$$\tau_{ji}^b = -1_{j=i} \sum_{\varepsilon \in \{-, +\}} \mathbb{E} [\varepsilon \mu^m (u^h(c_i, i) - u^h(c_0, 0)) \mid j(h) = j \text{ and } h \in B_i^\varepsilon].$$

### 10.3.2 Supply Elasticities: Mirrlees case

We now verify that the logic of section 5.1 applies to the Mirrlees case: with behavioral agent, the supply elasticities generally are featured in the optimal income tax formula. To make the point, take the case where the aggregate constraint is:  $\Phi \left( \int nL(n) f(n) dn \right) \leq g + \int c(n) f(n) dn$ , where  $\Phi(L)$  is an aggregate production function. Indeed, recall that in the Mirrlees framework an agent of productivity  $n$  is considered to supply  $n$  units of effective labor. Call  $w = \Phi'(L)$  the wage rate. We can extend the analysis of section 4, with an index  $w$  (which can be thought as being normalized to  $w = 1$  in that section). Then, we can write the agents' utility function problem as  $u^n(c, z, w) = U\left(c, \frac{z}{nw}\right)$  and the earnings supply as  $z^n(q, Q, r_0, r, w) = wnL^n(q, Q, r_0, r, w)$ . Other functions acquire a  $w$  term, e.g. indirect utility becomes  $v^n(q, Q, r_0, r, w)$ . Given a tax system, the equilibrium wage  $w$  satisfies:

$$w = \Phi' \left( \int \frac{z^n(q, Q, r_0, r, w)}{w} f(n) dn \right), \quad (93)$$

which defines an equilibrium wage  $w(Q, r_0)$

The objective function is:

$$L(Q, r_0, w) = \int W(v(n)) f(n) dn + \lambda \left[ \Phi \left( \int \frac{z(n)}{w} f(n) dn \right) - \int c(n) f(n) dn \right],$$

and we can define  $L(Q, r_0) = L(Q, r_0, w(Q, r_0))$  when taking into account the equilibrium wage  $w$ .

**Proposition 10.2** *In the Mirrlees model, suppose that the production function is imperfectly elastic. Then, the optimum tax features  $L_{qz^*}(Q, r_0) = 0$ , with*

$$L_{Qz^*}(Q, r_0) = L_{Qz^*}(Q, r_0, w) + L_w(Q, r_0, w) w_{Qz^*}(Q, r_0). \quad (94)$$

*The term  $L_{Qz^*}(Q, r_0, w)$ , with fixed wage, was calculated in Proposition 11.1 (with the normalization  $w = 1$ ). Hence, the optimal tax formula  $w_{Qz^*}(Q, r_0)$  generally depends on production elasticity, and does not coincide with the one in Section 5.1.*

When agents are rational, one can verify that  $L_w(Q, r_0, w) = 0$  at the optimum (see the proof of Proposition 10.2). Hence, in the traditional analysis, the supply elasticity (captured by  $w_{Qz^*}(Q, r_0)$ ) doesn't appear in the optimal tax formula. This is not true any more with a behavioral model.

**Proof of Proposition 10.2** The tax formula in the Proposition follows from the Chain rule. Next, we verify that when agents are rational,  $L_w = 0$  at  $w = w_0$ . We normalize  $w_0 = 1$  for simplicity. Suppose a given value of  $w$  and  $R(z)$ . Define  $\tilde{R}(z', w) = R(z'w)$ . Then, as  $z^n =$

$\arg \max_z U^n \left( R(z), \frac{z}{nw} \right)$ , i.e.

$$\frac{z}{w} = \arg \max U^n \left( \tilde{R} \left( \frac{z}{w}, w \right), \frac{z}{nw} \right).$$

So  $L(R(\cdot), w) = L(\tilde{R}(\cdot, w), 1)$ . That is, the welfare is the same as if we had a different tax system  $\tilde{R}$ , and a wage  $w = 1$ . Thus, given we started at an optimum tax system ( $R^0(\cdot) = \arg \max_{R(\cdot)} L(R(\cdot), 1)$ ), we have  $L_{\tilde{R}}(\tilde{R}, 1) = 0$ , hence  $L_w = 0$ .  $\square$

## 11 Proofs not included in the paper

### 11.1 General proofs

**Proof of Proposition 2.2** We observe that a tax  $\tau_i$  modifies the externality as:

$$\frac{d\xi}{d\tau_i} = \sum_h \xi_{c^h} \left[ c_{q_i}^h(\mathbf{q}, w, \xi) + c_\xi^h \frac{d\xi}{d\tau_i} \right],$$

so  $\frac{d\xi}{d\tau_i} = \frac{\sum_h \xi_{c^h} c_{q_i}^h}{1 - \sum_h \xi_{c^h} c_\xi^h}$ . The term  $\frac{1}{1 - \sum_h \xi_{c^h} c_\xi^h}$  represents the “multiplier” effect of one unit of pollution on consumption, then on more pollution. So, calling  $\frac{\partial L}{\partial \tau_i}^{\text{no } \xi}$  the value of  $\frac{\partial L}{\partial \tau_i}$  without the externality (that was derived in Proposition 2.1)

$$\begin{aligned} \frac{\partial L}{\partial \tau_i} - \frac{\partial L^{\text{no } \xi}}{\partial \tau_i} &= \frac{d\xi}{d\tau_i} \left\{ \sum_h W_{v^h} v_w^h \frac{v_\xi^h}{v_w^h} + \lambda \sum_h \boldsymbol{\tau} \cdot c_\xi^h(\mathbf{q}, w, \xi) \right\} = \frac{d\xi}{d\tau_i} \sum_h \left[ \beta^h \frac{v_\xi^h}{v_w^h} + \lambda \boldsymbol{\tau} \cdot c_\xi^h \right] \\ &= \frac{\sum_h \xi_{c^h} c_{q_i}^h}{1 - \sum_h \xi_{c^h} c_\xi^h} \sum_h \left[ \beta^h \frac{v_\xi^h}{v_w^h} + \lambda \boldsymbol{\tau} \cdot c_\xi^h \right] = \Xi \sum_h \xi_{c^h} c_{q_i}^h. \end{aligned}$$

Using Proposition 2.1,

$$\begin{aligned} \frac{\partial L}{\partial \tau_i} &= \sum_h [(\lambda - \gamma^h) c_i^h + \lambda \boldsymbol{\tau} \cdot \mathbf{S}_i^{C,h} - \beta^h \boldsymbol{\tau}^{b,h} \cdot \mathbf{S}_i^{C,h} + \Xi \xi_{c^h} \cdot (-c_w^h c_i^h + \mathbf{S}_i^{C,h})] \\ &= \sum_h [(\lambda - \gamma^h - \Xi \xi_{c^h} \cdot c_w^h) c_i^h + \lambda (\boldsymbol{\tau} + \frac{\Xi}{\lambda} \xi_{c^h}) \cdot \mathbf{S}_i^{C,h} - \beta^h \boldsymbol{\tau}^{b,h} \cdot \mathbf{S}_i^{C,h}]. \end{aligned}$$

**Proof of Proposition 2.4** We have, from Proposition 13.5

$$\boldsymbol{\tau}^{b,h} = u_C^{s,h}(\mathbf{C}^h) - u_C^h(\mathbf{C}^h) + \mathbf{p} - \mathbf{p}^{s,h} + \boldsymbol{\tau} - \boldsymbol{\tau}^{s,h} = \boldsymbol{\tau}^{I,h} + \boldsymbol{\tau} - \boldsymbol{\tau}^{s,h} = \boldsymbol{\tau}^{I,h} + (I - \mathbf{M}^h) \boldsymbol{\tau},$$

hence

$$\boldsymbol{\tau} - \boldsymbol{\tau}^{\xi,h} - \frac{\gamma^{\xi,h}}{\lambda} \boldsymbol{\tau}^{b,h} = \boldsymbol{\tau} - \boldsymbol{\tau}^{\xi,h} - \frac{\gamma^{\xi,h}}{\lambda} (\boldsymbol{\tau}^{I,h} + (I - \mathbf{M}^h) \boldsymbol{\tau}) = [I - (I - \mathbf{M}^h) \frac{\gamma^{\xi,h}}{\lambda}] \boldsymbol{\tau} - \boldsymbol{\tau}^{X,h}.$$

Hence, Proposition 2.2 implies:

$$\sum_h (1 - \frac{\gamma^{\xi,h}}{\lambda}) \mathbf{c}^h = - \sum_h (\mathbf{S}^{H,h})' (\boldsymbol{\tau} - \boldsymbol{\tau}^{\xi,h} - \tilde{\boldsymbol{\tau}}^{b,\xi,h}) = - \sum_h \mathbf{M}^{h,r} \mathbf{S}^{h,r} [[I - (I - \mathbf{M}^h) \frac{\gamma^{\xi,h}}{\lambda}] \boldsymbol{\tau} - \boldsymbol{\tau}^{X,h}]. \quad (95)$$

□

### Proof of Proposition 3.1

We start from the Ramsey planning problem in (12). Define

$$L = \gamma \sum_{i=1}^n \left[ \frac{(c_i(\tau_i))^{1-1/\psi_i} - 1}{1 - 1/\psi_i} - (p_i + \tau_i) c_i(\tau_i) \right] + \lambda \sum_{i=1}^n \tau_i c_i(\tau_i)$$

where  $c_i = (p_i + m_i \tau_i)^{-\psi_i}$ . The first-order condition with respect to  $\tau_i$  is:

$$L_{\tau_i} = \gamma \left[ [(c_i(\tau_i))^{-1/\psi_i} - (p_i + \tau_i)] \frac{\partial c_i}{\partial \tau_i} - c_i(\tau_i) \right] + \lambda \left[ c_i(\tau_i) + \tau_i \frac{\partial c_i}{\partial \tau_i} \right] = 0$$

Note that  $c_i(\tau_i)^{-1/\psi_i} = p_i + m_i \tau_i$  and  $\partial c_i / \partial \tau_i = -\psi_i \frac{c_i}{p_i + m_i \tau_i} m_i$ , we can rewrite the FOC as:

$$\begin{aligned} L_{\tau_i} &= \gamma \left[ \left( \frac{\lambda}{\gamma} - 1 + m_i \right) \tau_i \frac{-\psi_i c_i(\tau_i) m_i}{p_i + m_i \tau_i} \right] + (\lambda - \gamma) c_i(\tau_i) \\ &= -\lambda \left( \Lambda + \frac{\gamma}{\lambda} m_i \right) \frac{\psi_i \tau_i c_i(\tau_i) m_i}{p_i + m_i \tau_i} + \lambda \Lambda c_i(\tau_i) = 0 \end{aligned}$$

Simplifying gives us:

$$\left( \Lambda + \frac{\gamma}{\lambda} m_i \right) \psi_i \tau_i m_i = \Lambda (p_i + m_i \tau_i)$$

which gives an explicit expression for  $\tau_i$ :

$$\frac{\tau_i}{p_i} = \frac{\Lambda}{\psi_i m_i} \frac{1}{\Lambda + (1 - \Lambda) m_i - \Lambda / \psi_i} = \frac{\Lambda}{\psi_i m_i^2} \frac{1}{1 + \Lambda \left( \frac{1 - m_i - 1/\psi_i}{m_i} \right)}$$

**Derivation of (14), the approximate loss from taxation of inattentive agents** We give two proofs of this result. The first, and most elementary one, is that this is a Corollary to Lemma 9.2, using the behavioral elasticity  $\alpha_i = m_i \psi_i$ .

The other proof is as follow. Because utility is quasilinear,  $v(\mathbf{p}, \mathbf{p}^s, w) = w + v(\mathbf{p}, \mathbf{p}^s, 0)$ , so



$e(\mathbf{p}, \mathbf{p}^s, 0) = -v(\mathbf{p}, \mathbf{p}^s, 0)$ . We have

$$\begin{aligned}\mathcal{L}(\boldsymbol{\tau}) &= u(\mathbf{c}) + \lambda \boldsymbol{\tau} \cdot \mathbf{c} \\ &= v(\mathbf{p} + \boldsymbol{\tau}, \mathbf{p} + \boldsymbol{\tau}^s, w) + \lambda \boldsymbol{\tau} \cdot \mathbf{c}(\mathbf{p} + \boldsymbol{\tau}, \mathbf{p} + \boldsymbol{\tau}^s, w) \\ &= w - e(\mathbf{p} + \boldsymbol{\tau}, \mathbf{p} + \boldsymbol{\tau}^s, 0) + (1 + \Lambda) \boldsymbol{\tau} \cdot \mathbf{c}(\mathbf{p} + \boldsymbol{\tau}, \mathbf{p} + \boldsymbol{\tau}^s, w).\end{aligned}$$

By Taylor expansion, around  $(\boldsymbol{\tau}, \boldsymbol{\tau}^s) = (0, 0)$ , using Propositions 12.1 and 12.5, we have:

$$\begin{aligned}e(\mathbf{p} + \boldsymbol{\tau}, \mathbf{p} + \boldsymbol{\tau}^s, 0) - e(\mathbf{p}, \mathbf{p}) &= [e_{\mathbf{p}} \boldsymbol{\tau} + e_{\mathbf{p}^s} \boldsymbol{\tau}^s] + \left[ \frac{1}{2} \boldsymbol{\tau} e_{\mathbf{p}\mathbf{p}} \boldsymbol{\tau} + \boldsymbol{\tau} e_{\mathbf{p}, \mathbf{p}^s} \boldsymbol{\tau}^s + \frac{1}{2} \boldsymbol{\tau}^{s'} e_{\mathbf{p}^s \mathbf{p}^s} \boldsymbol{\tau}^s \right] + o(\|\boldsymbol{\tau}\|^2) \\ &= [\mathbf{c}^d \boldsymbol{\tau} + 0 \cdot \boldsymbol{\tau}^s] + \left[ 0 + \boldsymbol{\tau} \mathbf{S}^r \boldsymbol{\tau}^s + \frac{1}{2} \boldsymbol{\tau}^{s'} \mathbf{S}^r \boldsymbol{\tau}^s \right] + o(\|\boldsymbol{\tau}\|^2) \\ &= \mathbf{c}^d \boldsymbol{\tau} + \boldsymbol{\tau} \mathbf{S}^r \boldsymbol{\tau}^s + \frac{1}{2} \boldsymbol{\tau}^{s'} \mathbf{S}^r \boldsymbol{\tau}^s + o(\|\boldsymbol{\tau}\|^2).\end{aligned}$$

Using Proposition 12.3,

$$\begin{aligned}\boldsymbol{\tau} \cdot \mathbf{c}(\mathbf{p} + \boldsymbol{\tau}, \mathbf{p} + \boldsymbol{\tau}^s, w) &= \boldsymbol{\tau} \cdot (\mathbf{c}^d + \mathbf{c}_{\mathbf{p}} \boldsymbol{\tau} + \mathbf{c}_{\mathbf{p}^s} \boldsymbol{\tau}^s) \\ &= \boldsymbol{\tau} \cdot (\mathbf{c}^d + 0 + \mathbf{S}^r \boldsymbol{\tau}^s + o(\|\boldsymbol{\tau}\|)) = \boldsymbol{\tau} \mathbf{c}^d + \boldsymbol{\tau}' \mathbf{S}^r \boldsymbol{\tau}^s + o(\|\boldsymbol{\tau}\|^2)\end{aligned}$$

$$\begin{aligned}\mathcal{L}(\boldsymbol{\tau}) - \mathcal{L}(0) &= -[e(\mathbf{p} + \boldsymbol{\tau}, \mathbf{p} + \boldsymbol{\tau}^s, 0) - e(\mathbf{p}, \mathbf{p})] + (1 + \Lambda) \boldsymbol{\tau} \cdot \mathbf{c}(\mathbf{p} + \boldsymbol{\tau}, \mathbf{p} + \boldsymbol{\tau}^s, w) \\ &= -\mathbf{c}^d \boldsymbol{\tau} - \boldsymbol{\tau} \mathbf{S}^r \boldsymbol{\tau}^s - \frac{1}{2} \boldsymbol{\tau}^{s'} \mathbf{S}^r \boldsymbol{\tau}^s + (1 + \Lambda) (\boldsymbol{\tau} \mathbf{c}^d + \boldsymbol{\tau}' \mathbf{S}^r \cdot \boldsymbol{\tau}^s) + o(\|\boldsymbol{\tau}\|^2) \\ &= \Lambda (\boldsymbol{\tau} \mathbf{c}^d + \boldsymbol{\tau}' \mathbf{S}^r \cdot \boldsymbol{\tau}^s) - \frac{1}{2} \boldsymbol{\tau}^{s'} \mathbf{S}^r \boldsymbol{\tau}^s + o(\|\boldsymbol{\tau}\|^2) \\ &= \Lambda \boldsymbol{\tau} \mathbf{c}^d - \frac{1}{2} \boldsymbol{\tau}^{s'} \mathbf{S}^r \boldsymbol{\tau}^s + o(\|\boldsymbol{\tau}\|^2) + O(\|\boldsymbol{\tau}\|^2 \Lambda).\end{aligned}$$

□

**Proof of Proposition 3.3** We now assume that there are several consumers, indexed by  $h = 1 \dots H$ . Agent  $h$  maximizes  $u^h(c_0^h, c^h) = c_0^h + U^h(c^h)$ . The associated externality/internality is  $\xi^h c^h$ . He pays an attention  $m^h$  to the tax so that perceived taxes are  $\tau_h^s = m^h \tau$ . The government is utilitarian, so that the government planning problem is

$$\sum_h U^h(c^h) - (p + \xi^h) c^h. \quad (96)$$

We call  $c^{*h} = \arg \max_{c^h} U^h(c^h) - (p + \xi^h) c^h$  the quantity consumed by the agent at the first best.

To make things transparent, we specify

$$U^h(c) = \frac{a^h c - \frac{1}{2}c^2}{\Psi},$$

which using  $U_c^h = \frac{a^h - c}{\Psi} = q^s$ , implies a demand function  $c^h(q^s) = a^h - \Psi q^s$ .<sup>93</sup>

After some algebraic manipulations, social welfare compared to the first best can be written as

$$L(\tau) = -\frac{\Psi}{2} \sum_h (m^h \tau - \xi^h)^2. \quad (97)$$

The first best cannot be implemented unless all agents have the same ideal Pigouvian tax,  $\xi^h/m^h$ . Heterogeneity in attention creates welfare losses.

*Optimal Pigouvian tax.* At the optimum,  $U^{h'}(c^{h*}) = \mathbf{p} + \xi$ . If the agent perceives only  $m^h \tau$ , his demand is off the ideal  $c^{h*}$  (up to second order terms) as:

$$c^h = c^{h*} - \Psi (m^h \tau - \xi^h).$$

This expression is exact in the quadratic functional form about, and otherwise the leading term of a Taylor expansion of a general function, with now the interpretation  $\Psi = \psi^h c^{h*}$  then. So the welfare loss is:

$$L^h = W^h - W^{h*} = \frac{1}{2} u^{h''} \cdot (-\Psi \cdot (m^h \tau - \xi^h))^2 = -\frac{1}{2} \Psi (m^h \tau - \xi^h)^2,$$

and social welfare is  $L = \sum_h L^h = -\frac{\Psi}{2} \sum_h (m^h \tau - \xi^h)^2$

Because  $L_\tau = -\Psi \sum_h m^h (m^h \tau - \xi^h)$ , the optimal tax is

$$\tau^* = \frac{\sum_h \xi_h m^h}{\sum_h m_h^2} = \frac{\mathbb{E}[\xi_h m^h]}{\mathbb{E}[m^h]^2}.$$

---

<sup>93</sup>The expressions in the rest of this section are exact with this quadratic utility specification. For general utility functions, they hold provided that they are understood as the leading order terms in a Taylor expansion around an economy with no heterogeneity.

Let us calculate  $V = \mathbb{E} \left[ (m^h \tau - \xi^h)^2 \right]$  at this optimum  $\tau = \tau^*$ ,

$$\begin{aligned} V &= \mathbb{E} \left[ m^{h2} \right] \tau^{*2} - 2\mathbb{E} \left[ m^h \xi^h \right] \tau^* + \mathbb{E} \left[ \xi^{h2} \right] \\ &= \mathbb{E} \left[ m^{h2} \right] \frac{\mathbb{E} \left[ \xi_h m^h \right]^2}{\mathbb{E} \left[ m^{h2} \right]^2} - 2\mathbb{E} \left[ m^h \xi_h \right] \frac{\mathbb{E} \left[ \xi_h m^h \right]}{\mathbb{E} \left[ m^{h2} \right]} + \mathbb{E} \left[ \xi_h^2 \right] = -\frac{\mathbb{E} \left[ \xi_h m^h \right]^2}{\mathbb{E} \left[ m^{h2} \right]} + \mathbb{E} \left[ \xi_h^2 \right] \\ &= \frac{\mathbb{E} \left[ \xi_h^2 \right] \mathbb{E} \left[ m^{h2} \right] - \mathbb{E} \left[ \xi_h m^h \right]^2}{\mathbb{E} \left[ m^{h2} \right]}. \end{aligned}$$

hence the welfare loss is:  $L = -\frac{1}{2} H \Psi \frac{\mathbb{E} \left[ \xi_h^2 \right] \mathbb{E} \left[ m^{h2} \right] - (\mathbb{E} \left[ \xi_h m^h \right])^2}{\mathbb{E} \left[ m^{h2} \right]}$ .

If there is no tax, the loss is (from equation 97):

$$L^{\text{no tax}} = -\frac{\Psi}{2} \sum_h (m^h \cdot 0 - \xi^h)^2 = -\frac{\Psi}{2} \sum_h \xi_h^2 = -\frac{1}{2} H \Psi \mathbb{E} \left[ \xi_h^2 \right].$$

So,  $L = L^{\text{no tax}} \frac{\mathbb{E} \left[ \xi_h^2 \right] \mathbb{E} \left[ m^{h2} \right] - (\mathbb{E} \left[ \xi_h m^h \right])^2}{\mathbb{E} \left[ m^{h2} \right] \mathbb{E} \left[ \xi_h^2 \right]}$ .

*Optimal quantity mandate.* Welfare is  $\sum_h [U^h(c^*) - (p + \xi^h) c^*]$ . The optimal quantity restriction  $c^*$  is characterized by:

$$\frac{1}{H} \sum_h U^{h'}(c^*) = p + \frac{1}{H} \sum_h \xi^h. \quad (98)$$

The welfare loss compared to the first best, which entails  $U^{h'}(c^{h*}) = p + \xi^h$  is

$$L^h = \frac{1}{2} U^{h''}(c) (c^{h*} - c^*)^2 = -\frac{1}{2} \frac{1}{\Psi} (c^{h*} - c^*)^2.$$

The best consumption satisfies:  $L_{c^*}^Q = \frac{1}{2} \sum_h \frac{1}{\Psi} (c^{h*} - c^*) = 0$ , i.e.  $c^* = \mathbb{E} [c^{h*}]$

The loss is:

$$L^Q = -\frac{1}{2} \frac{H}{\Psi} \mathbb{E} \left[ (c^{h*} - c^*)^2 \right] = -\frac{1}{2} \frac{H}{\Psi} \text{var} (c^{h*}).$$

□

**Proof of Proposition 3.4** Equation (11) then yields the optimal tax:

$$\tau = (\mathbb{E} [M^{h'} S^r M^h])^{-1} \mathbb{E} [M^{h'}] S^r \tau^X. \quad (99)$$

with  $\tau^X = (\xi_*, 0)'$ .

When agents have uniform misperceptions ( $M^h = M$ ), the optimal tax is  $\tau = M^{-1} \tau^X$ . This

implies  $\tau_1 = \frac{\xi_*}{m_1} > 0$  and  $\tau_2 = 0$ . The principle of targeting applies. This is no longer true when misperceptions are not uniform.

We have  $(\mathbb{E}[\mathbf{M}^{h'} \mathbf{S}^r \mathbf{M}^h])_{ij} = \mathbf{S}_{ij}^r \mathbb{E}[m_i^h m_j^h]$  and  $(\mathbb{E}[\mathbf{M}^{h'} \mathbf{S}^r])_{ij} = \mathbb{E}[m_i^h] S_{ij}^r$ . Matrix inversion gives:

$$\tau_2 = \frac{S_{11}^r S_{12}^r \mathbb{E}[m_1] (\mathbb{E}[m_1^2] \mathbb{E}[m_2] - \mathbb{E}[m_1 m_2] \mathbb{E}[m_1])}{\det \mathbb{E}[\mathbf{M}^{h'} \mathbf{S}^r \mathbf{M}^h]} \xi_*.$$

Because  $\mathbb{E}[\mathbf{M}^{h'} \mathbf{S}^r \mathbf{M}^h]$  is a dimension  $2 \times 2$  and has negative roots (there is a good 0, so that  $\mathbf{S}^r$  is the block matrix excluding good 0, and has only negative root),  $\det \mathbb{E}[\mathbf{M}^{h'} \mathbf{S}^r \mathbf{M}^h] > 0$ . The condition in the Proposition is that  $\mathbb{E}[m_1^2] \mathbb{E}[m_2] - \mathbb{E}[m_1 m_2] \mathbb{E}[m_1] > 0$ . Hence,  $\text{sign}(\tau_2) = -\text{sign}(S_{12})$ .

The quadratic case simply gives a constant matrix  $\mathbf{S}^r$ .  $\square$

### Proof of Proposition 3.8

$$\begin{aligned} 0 &= q_i \frac{\partial L}{\partial \tau_i} = \sum_h [(\lambda - \gamma^{k,h}) q_i c_i^h + \lambda \boldsymbol{\tau} \cdot \mathbf{S}_i^{C,k,h} q_i] \\ &= \sum_h [(\lambda - \gamma^{k,h}) q_i c_i^h + \lambda \sum_j \tau_j S_{ji}^{C,k,h} q_i]. \end{aligned}$$

Summing over  $i$  in the account gives:

$$0 = \sum_h [(\lambda - \gamma^{k,h}) \sum_i q_i c_i^h + \lambda \sum_j \tau_j \sum_i S_{ji}^{C,k,h} q_i].$$

By the traditional Slutsky relation with account  $\sum_i S_{ji}^{C,k,h} q_i = 0$  for all  $j, h$ , so

$$0 = \sum_h [(\lambda - \gamma^{k,h}) \sum_i q_i c_i^h] = 0.$$

With just one type of agent  $h$ , this gives  $\lambda - \gamma^k = 0$ . This implies, for all  $i$ :

$$\sum_j \tau_j S_{ji}^{C,k} q_i = 0.$$

This first order condition is verified by  $\tau_j = \tau$  for all  $j$ . For a generic Slutsky matrix, the only solution of  $x \cdot S^C = 0$  is  $x = tq$  for some real  $t$  (we do not have a proof of this, but this is highly likely). This implies that  $\tau_j = \tau$  for some  $\tau$ .  $\square$

**Proof of Proposition 3.9** *Ramsey case.* We have  $c_i(\tau_i) = \frac{\omega_i}{p_i + \tau_i}$ . The planner's problem entails:

$$\max_{\tau_i} u_i(c_i(\tau_i)) - (p_i + \tau_i) c_i(\tau_i) + \lambda \tau_i c_i(\tau_i),$$

which gives, using  $c'_i(\tau_i) = -\frac{c_i(\tau_i)}{p_i + \tau_i}$

$$[u'_i(c_i) - (p_i + \tau_i(1 - \lambda))] \frac{c_i}{p_i + \tau_i} + (\lambda - 1)c_i = 0,$$

i.e.

$$u'_i(c_i) = \lambda. \quad (100)$$

When  $u'_i(c_i) = c_i^{1/\psi_i}$  this gives the announced result,  $\frac{\tau_i}{p_i} = \lambda^{\psi_i} - 1$ .

*Pigou case.* The first best features  $u'_i(c_i) = p_i + \xi_i$ , and the rigid mental account gives  $c_i(\tau_i) = \frac{\omega_i}{p_i + \tau_i}$ , where  $u'_i(\omega_i) = p_i$ . Hence, we have:

$$c_i = \frac{u_i'^{-1}(p_i)}{p_i + \tau_i} = u_i'^{-1}(p_i + \xi_i),$$

i.e.

$$\frac{\tau_i}{p_i} = \frac{u_i'^{-1}(p_i)}{u_i'^{-1}(p_i + \xi_i)} - 1.$$

When  $u'_i(c_i) = c_i^{1/\psi_i}$  this gives the announced result,  $\frac{\tau_i}{p_i} = \left(1 + \frac{\xi_i}{p_i}\right)^{\psi_i} - 1$ .  $\square$

**Proof of Proposition 5.1** We compute the derivatives of the Lagrangian:

$$\frac{\partial L}{\partial \tau_i^\kappa} = \sum_h \left[ W_{v^h} \left( v_{\tau_i^\kappa}^h + v_{q^p}^h \cdot p_{\tau_i^\kappa} \right) - \mu p \cdot \left( c_{\tau_i^\kappa}^h + c_{q^p}^h \cdot p_{\tau_i^\kappa} \right) \right].$$

To calculate this, let us make an analogy with our basic Ramsey model with fixed prices. We expressed it  $L = W + \lambda \sum_h (\boldsymbol{\tau} \cdot c^h - w)$ , and it can be re-expressed:

$$\begin{aligned} L &= W + \lambda \sum_h (\boldsymbol{\tau} \cdot c^h - w) = W + \lambda \sum_h [(p + \tau) \cdot c^h - w - p \cdot c^h] \\ &= W - \lambda \sum_h p \cdot c^h, \end{aligned}$$

as  $q \cdot c^h - w = 0$ . So

$$\begin{aligned} L_{\tau_i^\kappa}^{\text{fixed price}} &= \frac{\partial (W - \lambda \sum_h p \cdot c^h)}{\partial \tau_i^\kappa} = \sum_h \left[ W_{v^h} v_{\tau_i^\kappa}^h - \mu p \cdot c_{\tau_i^\kappa}^h \right] \\ &= \sum_h [(\lambda - \gamma^h) c_i^h + \lambda(\bar{\boldsymbol{\tau}} - \tilde{\boldsymbol{\tau}}^{b,h}) \cdot S_i^{C,\kappa,h}]. \end{aligned} \quad (101)$$

We then have

$$L_{\tau_i^\kappa} = L_{\tau_i^\kappa}^{\text{fixed price}} + \sum_j L_{\tau_j^p} \varepsilon_{ji}^\kappa = L_{\tau_i^\kappa}^{\text{fixed price}} + L_{\tau^p} \cdot \varepsilon_i^\kappa. \quad (102)$$

□

The intuition is as follows. With a full set of commodity taxes  $\boldsymbol{\tau}^p$ , we can rewrite the objective function and the resource constraint in the planning problem as a function of  $\mathbf{q} = \mathbf{p} + \boldsymbol{\tau}^p$ . We can then relax the planning problem by dropping the competitive pricing equation, which is slack—this equation can then simply be used to find  $\boldsymbol{\tau}^p$  given the desired value of  $\mathbf{q}$ . As a result, only the first derivatives of the production function  $\mathbf{p} = F'$  enter the optimal tax formulas and not the second derivatives  $F''$  (and hence do not depend on supply elasticities). With a restricted set of commodity taxes  $\boldsymbol{\tau}^p$ , this relaxation of the planning problem fails, the competitive pricing equation cannot be dropped, and the optimal tax formulas depend on the second derivatives  $F''$  (and hence depend on supply elasticities).□

**Proof of Lemma 9.4** We observe that a tax  $\tau^{sh}$  modifies the externality as:

$$\frac{d\xi(\{w^h\}, \{\tau^{s,h}\})}{d\tau^{s,h}} = \xi_{c^h} c_{\tau^{s,h}}^h(q, w, \xi) + \sum_{h'} \xi_{c^{h'}} c_{\xi}^{h'} \frac{d\xi}{d\tau^{s,h}},$$

so

$$\frac{d\xi(\{w^h\}, \{\tau^{s,h}\})}{d\tau^{s,h}} = \frac{\xi_{c^h} c_{\tau^{s,h}}^h}{1 - \sum_{h'} \xi_{c^{h'}} c_{\xi}^{h'}}. \quad (103)$$

Also  $\frac{d\xi}{dw^h} = \xi_{c^h} c_{w^h}^h(q, w, \xi) + \sum_{h'} \xi_{c^{h'}} c_{\xi}^{h'} \frac{d\xi}{dw^h}$ , so

$$\frac{d\xi}{dw^h} = \frac{\xi_{c^h} c_{w^h}^h}{1 - \sum_{h'} \xi_{c^{h'}} c_{\xi}^{h'}}. \quad (104)$$

We note that the FOC of (84) in  $w^h$  is

$$\begin{aligned} 0 &= v_{w^h}^h + \lambda(\tau^{s,h} \cdot c_{w^h}^{rh} - 1) + \frac{d\xi}{dw^h} \sum_{h'} v_{\xi}^{h'} + \lambda \tau^{s,h'} \cdot \bar{c}_{\xi}^{rh'} \\ &= v_{w^h}^h + \lambda(\tau^{sh} \cdot c_{w^h}^{rh} - 1) + \frac{\xi_{c^h} c_{w^h}^h}{1 - \sum_{h'} \xi_{c^{h'}} c_{\xi}^{h'}} \sum_{h'} v_{\xi}^{h'} + \lambda \tau^{sh'} \cdot \bar{c}_{\xi}^{rh'} \text{ by (104)} \\ &= v_{w^h}^h + \lambda(\tau^{sh} \cdot c_{w^h}^{rh} - 1) + \xi_{c^h} c_{w^h}^h \Xi = v_{w^h}^h + \lambda(\tau^{s,h} \cdot c_{w^h}^{rh} - 1) - \lambda \tau^{\xi h} c_{w^h}^h \\ &= \gamma^{h,\xi} - \lambda. \end{aligned}$$

which confirms that at the optimum  $\gamma^{h,\xi} = \lambda$  for all agents – even if the tax hasn't been optimized upon.

$$\begin{aligned}
g(\{\tau_{s,h}\}) &= \max_{\{w^h\}} \sum_h v^r(p + \tau^{s,h}, w^h, \xi) + \lambda \sum_h [\tau^{s,h} \cdot \bar{c}^r(p + \tau^{s,h}, w^h, \xi) - w^h] \\
&= \max_{\{w^h\}} \sum_h v^r(p + \tau^{s,h}, w^h, \xi) + \lambda \sum_h [p \cdot \bar{c}^r(p + \tau^{s,h}, w^h, \xi)] \quad [\text{probably not useful}] \\
&\quad (p + \tau^h) c = w^h.
\end{aligned}$$

We take the derivatives (84):

$$\begin{aligned}
g_{\tau_{s,h}}(\{\tau_{s,h'}\}) &= (\lambda - v_w^h) c^{hr} + \lambda \tau^{s,h} c_p^{hr} + \frac{d\xi}{d\tau_{s,h}} \sum_{h'} [v_\xi^{h'} + \lambda \tau^{s,h'} \cdot \bar{c}_\xi^r] \\
&= (\lambda - v_w^h) c^{hr} + \lambda \tau^{s,h} c_p^{hr} + \frac{\xi_{c^h} c_{\tau_{s,h}}^h}{1 - \sum_{h'} \xi_{c^{h'}} c_\xi^{h'}} \sum_{h'} v_\xi^{h'} + \lambda \tau^{s,h'} \cdot \bar{c}_\xi^r \quad \text{by (103)} \\
&= (\lambda - v_w^h) c^{hr} + \lambda \tau^{s,h} c_p^{hr} + \xi_{c^h} c_{\tau_{s,h}}^h \Xi \\
&= (\lambda - v_w^h) c^{hr} + \lambda \tau^{s,h} c_p^{hr} - \lambda \tau^{\xi h} c_p^h \\
&= (\lambda - v_w^h) c^{hr} + \lambda (\tau^{sh} - \tau^{\xi h}) c_p^{hr} \\
&= (\lambda - v_w^h) c^{hr} + \lambda (\tau^{sh} - \tau^{\xi h}) [S^{hr} - c_w^r c^h] \\
&= [\lambda - v_w^h - \lambda (\tau^{s,h} - \tau^{\xi h}) c_w] c^h + \lambda (\tau^{s,h} - \tau^{\xi h}) S^{hr} \\
&= \lambda (\tau^{s,h} - \tau^{\xi h}) S^{hr}. \tag{105}
\end{aligned}$$

Hence, observing that  $\tau^{s,h} - \tau^{\xi h} = 0$  at the optimum,

$$g_{\tau_{s,h} \tau_{s,h'}} = \lambda S^{hr} \left( 1_{h=h'} - \frac{d\tau^{\xi h}(\{\tau^{s,h''}\})}{d\tau_{s,h'}} \right). \tag{106}$$

*Example with quasi-linear utility, additive externality*

When  $u(c, \xi) = u(c_1, \dots, c_n) + \lambda c_0 + \frac{1}{H} \xi$ , we have  $c(p, \xi)$  independent of  $\xi$ , and  $\Xi = \frac{\sum_h \left[ \beta^h \frac{v_\xi^h}{v_w^h} + \lambda \tau \cdot \bar{c}_\xi^h \right]}{1 - \sum_h \xi_{c^h} c_\xi^h} = 1$ , and  $\tau^{\xi h} = -\frac{1}{\lambda} \xi_{c^h}$ .

So,  $\frac{d\tau^{\xi,h}}{d\tau_{s,h'}} = -\frac{1}{\lambda} \xi_{c^h c^{h'}} \frac{dc^{h'}}{d\tau_{s,h'}} = -\frac{1}{\lambda} \xi_{c^h c^{h'}} S^{rh'}$ :

$$\frac{d\tau^{\xi,h}}{d\tau_{s,h'}} = -\frac{1}{\lambda} \xi_{c^h c^{h'}} S^{rh'},$$

so

$$g_{\tau_{s,h} \tau_{s,h'}} = \lambda S^{rh} 1_{h=h'} + H S^{rh} \xi_{c^h c^{h'}} S^{rh'}. \tag{107}$$

That should generalize to additive externality:  $u(c, \xi) = u(c) + \frac{1}{H} \xi$ . Then  $\Xi = \frac{\sum_h \left[ \beta^h \frac{v_\xi^h}{v_w^h} + \lambda \tau \cdot \bar{c}_\xi^h \right]}{1 - \sum_h \xi_{c^h} c_\xi^h} = 1$ . And  $\tau^{\xi h} = -\frac{1}{\lambda} \xi_{c^h}(\{\tau^{sh}\})$ . When  $\{\tau^{-h'}\}$  are held constant, varying  $\tau^{h'}$  changing  $c^{h'}$  and  $\xi$ , but doesn't change the marginal utility of the agent  $-h'$ , so doesn't change their consumption, so

$$\frac{dc^{h''}}{d\tau^{h'}} = 0 \text{ for } h'' \neq h',$$

$$\begin{aligned} \frac{d\tau^{\xi,h}}{d\tau^{s,h'}} &= -\frac{1}{\lambda} \xi_{c^h c^{h'}} \frac{dc^{h'}}{d\tau^{h'}} = -\frac{1}{\lambda} \xi_{c^h c^{h'}} S^{rh'} + \sum_{h''} -\frac{1}{\lambda} \xi_{c^h c^{h''}} \frac{\partial c^{h''}}{dw^{h''}} \frac{dw^{h''}}{d\tau^{h'}} \\ &= -\frac{1}{\lambda} \xi_{c^h c^{h'}} \left( S^{rh'} - c_{w^{h'}}^{h'} c + c_{w^{h'}}^{h'} \frac{dw^{h'}}{d\tau^{h'}} \right) \\ &= -\frac{1}{\lambda} \xi_{c^h c^{h'}} \left[ S^{rh'} - c_{w^{h'}}^{h'} \left( c - \frac{dw^{h'}}{d\tau^{h'}} \right) \right]. \end{aligned}$$

as  $\gamma^{h,\xi} = v_w^h + \lambda (\tau^{sh} - \tau^{\xi h}) \cdot c_w^h = \lambda$  implies  $dv_w^h + \lambda (d\tau^{sh}) \cdot c_w = 0$ .  $\square$

**Proof of Proposition 9.2** We have apply our tax formulas (9):

$$\begin{aligned} \frac{\partial L}{\partial \chi} &= -\sum_h t^h \gamma^h c^h - \sum_h [(\lambda - \gamma^h (1 - m^h)) \tau - \lambda \tau^{X,h} + \gamma^h \chi \eta^h] \Psi \tau_x^{\chi,h} \\ &= -\sum_h t^h \gamma^h c^h - \Psi \sum_h [(\lambda - \gamma^h (1 - m^h)) \tau - \lambda \tau^{X,h} + \gamma^h \chi \eta^h] \eta^h \end{aligned} \quad (108)$$

$$= \sum_h [-t^h \gamma^h c^h - \Psi x^h \eta^h], \quad (109)$$

where

$$x^h = (\lambda - \gamma^h (1 - m^h)) \tau - \lambda \tau^{X,h} + \gamma^h \chi \eta^h.$$

Likewise,

$$\frac{\partial L}{\partial \tau} = \sum_h [(\lambda - \gamma_h) c^h - \Psi x^h m^h].$$

The problem is  $\max_{\chi, \tau} L(\tau, \chi)$  s.t.  $\chi \geq 0, \tau \geq 0$ . The Lagrangian is:

$$L^*(\tau, \chi) = L(\tau, \chi) + \pi \chi + \pi' \tau,$$

where  $\pi, \pi'$  are Lagrangian multipliers.

We observe that when there is no intervention ( $\tau = \chi = 0$ ), then  $x^h < 0$ .

$$\begin{aligned} \frac{L_\chi}{\eta^h} &= -\frac{t^h \gamma^h}{\eta^h} c^h - \Psi x^h \\ \frac{L_\tau}{m^h} &= \frac{\lambda - \gamma^h}{m^h} c^h - \Psi x^h, \end{aligned}$$

so

$$\frac{L_\tau}{m^h} - \frac{L_\chi}{\eta^h} = \left[ \frac{\lambda - \gamma^h}{m^h} + \frac{t^h \gamma^h}{\eta^h} \right] c^h = \Delta.$$



If the optimum features  $\chi > 0, \tau = 0$ , then  $L_\tau = -\pi' < 0 = L_\chi$ , which implies  $\Delta < 0$ .

If the optimum features  $\chi = 0, \tau > 0$ , then  $L_\chi = -\pi < 0 = L_\tau$ , which implies  $\Delta > 0$ .

If the optimum features  $\chi = \tau = 0$ , then  $L_\chi = -\pi < 0, L_\tau = -\pi' < 0$ , so  $\Psi x^h > \max(-\iota^h \gamma^h c^h, (\lambda - \gamma^h) c^h)$ . This implies in particular that  $\lambda < \gamma^h$ .

We note that if the problem had with no inequality constraints, and just one type of agent, then an interior solution features:  $\lambda - \gamma^h = -\iota^h \gamma^h$ . there is a large subsidy in place, (to help the agent), and the excess consumption is corrected via the nudge. That is, the policy is to ‘‘Subsidize the poor, and nudge them away from the good at the same time’’. This results is a bit knife-edge.

**Proof of Proposition 9.22** We note that for any tax system,

$$L(\{\tau^h\}, \{\tau^{s,h}\}, \{w^h\}) = L(\{\tau^{s,h}\}, \{\tau^{s,h}\}, \{w^h + (\tau^{s,h} - \tau^h) \cdot c^h(p + \tau^h, p + \tau^{s,h}, w^h, \xi)\}).$$

and

$$L(\{\tau^{s,h}\}, \{\tau^{s,h}\}, \{w^h\}) = W(v^{h,r}(p + \tau^{s,h}, w^h, \xi)) + \lambda \sum_h \left[ \tau \cdot \bar{c}^{\tau,h}(p + \tau^{s,h}, w^h, \xi) - w^h \right].$$

Here  $w, w^* \in \mathbb{R}^H$ . Call  $y = (\{\tau^h\}, \{\tau^{s,h}\}) \in \mathbb{R}^{2nH}$  (with  $n$  the number of goods). The first best (in a world with externalities) has  $(w^*, y^*)$ . We call  $w^{**}(y)$  the optimal redistribution given a tax system  $y$ . So,  $w^* = w^{**}(y^*)$ .

$$\begin{aligned} L^{\text{tot}} &= L(w, y) - L(w^*, y^*) \\ &= [L(w, y) - L(w^{**}(y), y)] + [L(w^{**}(y), y) - L(w^{**}(y^*), y^*)] \\ &= \frac{1}{2}(w - w^{**}(y)) \cdot L_{ww}(w^{**}(y), y) \cdot (w - w^{**}(y)) + \frac{1}{2}(y - y^*) g_{yy}(y - y^*) \text{ by Lemma 15.3} \\ &= L^{\text{distribution}} + L^{\text{distortion}} \end{aligned}$$

$$L^{\text{distribution}} = \frac{1}{2}(w - w^{**}(y)) \cdot L_{ww}(w^{**}(y), y) \cdot (w - w^{**}(y))$$

$$L^{\text{distortion}} = \frac{1}{2}(y - y^*) g_{yy}(y - y^*).$$

*Redistribution terms*

From Lemma 15.2, the expression of the loss involves  $L_{w_h}(w, \tau) = \gamma^{\xi h} - \lambda$ , the social marginal utility. Applying that Lemma 15.2 gives a loss:

$$L^{\text{distribution}} = \frac{1}{2} \sum_{h,h'} (\gamma^{\xi h} - \bar{\gamma}) (L_{ww}(w, \tau)^{-1})_{h,h'} (\gamma^{\xi h} - \bar{\gamma}). \quad (110)$$

*Tax distortion terms*

We have  $L^{\text{distortion}} = \frac{1}{2} (y - y^*) g_{yy} (y - y^*)$ . Note that  $g(y) = g(\{\tau^s\})$ .

$$g(\tau^s) = \max_{w_1, \dots, w_n} L(w, \tau^s) = L(w^*(\tau^s), \tau^s)$$

$$g_{\tau^s \tau^s} = L_{\tau^s \tau^s} - L_{\tau^s w} L_{ww}^{-1} L_{\tau^s w} \text{ by Lemma 15.3.}$$

□

**Proof of Lemma 9.2** Demand is:

$$c_i(\tau_i) = y_i (1 - \alpha_i \tau_i + O(\tau_i^2)), \quad (111)$$

Hence have:

$$\begin{aligned} L &= \sum_i U^i(c_i(\tau_i)) - (1 + \tau_i) c_i(\tau_i) + (1 + \Lambda) \tau_i c_i(\tau_i) \\ &= \sum_i U^i(c_i(\tau_i)) - c_i(\tau_i) + \Lambda \tau_i c_i(\tau_i) \\ &= \sum_i f_i(\tau_i) + \Lambda \tau_i c_i(\tau_i), \end{aligned}$$

with

$$f_i(\tau_i) = F^i(c_i(\tau_i)), \quad F^i(c) = U^i(c) - c.$$

We have

$$f_i(\tau_i) - f_i(0) = f'_i(0) \tau_i + \frac{1}{2} f''_i(0) \tau_i^2 + o(\tau^2)$$

$$\begin{aligned} f'_i(\tau_i) &= F^{i'}(c_i(\tau_i)) c'_i(\tau_i) \\ f''_i(\tau_i) &= F^{i''}(c_i(\tau_i)) c'_i(\tau_i)^2 + F^{i'}(c_i(\tau_i)) c''_i(\tau_i). \end{aligned}$$

As  $F'_i(0) = 0$ , we have

$$\begin{aligned} f'_i(0) &= 0 \\ f''_i(0) &= F''(c_i(0)) c'_i(0)^2 = U''_i(y_i) y_i^2 \alpha_i^2 \text{ using (111)} \\ &= \frac{-1}{\psi_i} q_i y_i \alpha_i^2 \text{ using (64)} \\ &= -\frac{\alpha_i^2}{\psi_i} y_i, \end{aligned}$$

so

$$\begin{aligned} L(\tau) - L(0) &= \sum_i \left[ \frac{1}{2} f''(0) \tau_i^2 + \Lambda \tau_i c_i \right] + o(\tau^2) + o(\tau \Lambda) \\ &= \sum_i \left[ -\frac{1}{2} \frac{\alpha_i^2}{\psi_i} y_i \tau_i^2 + \Lambda \tau_i c_i \right] + o(\tau^2) + o(\tau \Lambda). \end{aligned}$$

So the objective function is:

$$L = -\frac{1}{2} \sum_i -\frac{\alpha_i^2}{\psi_i} y_i \tau_i^2 + \Lambda \sum_i \tau_i c_i + o(\tau^2) + o(\tau \Lambda). \quad (112)$$

□

**Proof of Proposition 9.7** We use the extended utility function  $v^h(\mathbf{q}, \boldsymbol{\omega})$  and demand function  $\mathbf{c}^h(\mathbf{q}, \boldsymbol{\omega})$ . We use the Roy's identify with mental accounts, Proposition 9.8:

$$\begin{aligned} \frac{\partial L}{\partial \tau_i} &= \sum_h [W_{v^h} v_{\omega^{k(i)}}^h \frac{v_{q_i}^h}{v_{\omega^{k(i)}}^h} + \sum_k W_{v^h} v_{\omega^k}^h \omega_{q_i}^k + \lambda c_i^h + \lambda \tau \cdot [c_{q_i}^h + \sum_k \mathbf{c}_{\omega^k}^h \omega_{q_i}^k]], \\ &= \sum_h [-\beta^{k(i),h} (c_i^h + \tau^{b,k(i)} \cdot \mathbf{S}_i^{C,k(i),h}) + \sum_k \beta^{k,h} \omega_{q_i}^k + \lambda c_i^h + \lambda \tau \cdot \left( (\mathbf{S}_i^{C,k(i),h} - \mathbf{c}_{\omega^{k(i)}}^{k(i),h} c_i) + \sum_k \mathbf{c}_{\omega^k}^h \omega_{q_i}^k \right)], \\ &= \sum_h [(\lambda - \beta^{k(i),h} - \lambda \tau \cdot \mathbf{c}_{\omega^{k(i)}}^{k(i),h}) c_i^h + \sum_k [\beta^{k,h} + \lambda \tau \cdot \mathbf{c}_{\omega^k}^h] \omega_{q_i}^k + \lambda (\tau - \tilde{\tau}^{b,k(i),h}) \cdot \mathbf{S}_i^{C,k(i),h}], \\ &= \sum_h [(\lambda - \gamma^{k(i),h}) c_i^h + \sum_k \gamma^{k,h} \omega_{q_i}^k + \lambda (\tau - \tilde{\tau}^{b,k(i),h}) \cdot \mathbf{S}_i^{C,k(i),h}], \\ \frac{\partial L}{\partial \tau_i} &= \sum_h [(\lambda - \gamma^{k(i),h}) c_i^h + \lambda (\tau - \tilde{\tau}^{b,k(i),h}) \cdot \mathbf{S}_i^{C,k(i),h} + \sum_k \gamma^{k,h} \omega_{q_i}^k]. \end{aligned}$$

□

**Proof of Proposition 9.8** We first note a few simple identities. As  $B(\mathbf{c}(\mathbf{p}, \boldsymbol{\omega}), \mathbf{p}, \boldsymbol{\omega}) = 0$  and  $v(\mathbf{p}, \boldsymbol{\omega}) = u(\mathbf{c}(\mathbf{p}, \boldsymbol{\omega}))$ , we have:

$$B_{\mathbf{c}} \mathbf{c}_{p_i} + B_{p_i} = 0, \quad B_{\mathbf{c}} \mathbf{c}_{\omega^k} + B_{\omega^k} = 0, \quad v_{\omega^k} = u_{\mathbf{c}} \mathbf{c}_{\omega^k}. \quad (113)$$

We calculate:

$$\begin{aligned}\boldsymbol{\tau}^{b,k} \cdot \mathbf{c}_{\omega_l} &= \left( -\frac{u_{\mathbf{c}}}{v_{\omega_k}} - \frac{B_{\mathbf{c}}}{B_{\omega_k}} \right) \cdot \mathbf{c}_{\omega_l} \\ &= -\frac{v_{\omega_l}}{v_{\omega_k}} + \frac{B_{\omega_l}}{B_{\omega_k}}.\end{aligned}$$

using (113). So typically  $\boldsymbol{\tau}^{b,k} \cdot \mathbf{c}_{\omega_l} \neq 0$ , except when  $l = k$ :

$$\boldsymbol{\tau}^{b,k} \cdot \mathbf{c}_{\omega^k} = 0. \quad (114)$$

We are now ready to study Roy's identity. We have:

$$\begin{aligned}\frac{v_{p_i}(\mathbf{p}, \boldsymbol{\omega})}{v_{\omega_k}(\mathbf{p}, \boldsymbol{\omega})} &= \frac{u_{\mathbf{c}} \mathbf{c}_{p_i}(\mathbf{p}, \boldsymbol{\omega})}{v_{\omega_k}} = \left( \frac{u_{\mathbf{c}}}{v_{\omega_k}} + \frac{B_{\mathbf{c}}(\mathbf{c}, \mathbf{p}, \boldsymbol{\omega})}{B_{\omega_k}(\mathbf{c}, \mathbf{p}, \boldsymbol{\omega})} - \frac{B_{\mathbf{c}}(\mathbf{c}, \mathbf{p}, \boldsymbol{\omega})}{B_{\omega_k}(\mathbf{c}, \mathbf{p}, \boldsymbol{\omega})} \right) \mathbf{c}_{p_i} \\ &= \left( \frac{u_{\mathbf{c}}}{v_{\omega_k}} + \frac{B_{\mathbf{c}}}{B_{\omega_k}} \right) \mathbf{c}_{p_i} - \frac{B_{\mathbf{c}}}{B_{\omega_k}} \mathbf{c}_{p_i} \\ \frac{v_{p_i}(\mathbf{p}, \boldsymbol{\omega})}{v_{\omega_k}(\mathbf{p}, \boldsymbol{\omega})} &= -\boldsymbol{\tau}^{b,k} \cdot \mathbf{c}_{p_i} + \frac{B_{p_i}}{B_{\omega_k}}.\end{aligned} \quad (115)$$

using (113).

Using (114) gives:

$$\boldsymbol{\tau}^{b,k} \cdot \mathbf{c}_{p_j} = \boldsymbol{\tau}^{b,k} \cdot \mathbf{S}_j^{C,k},$$

so

$$\frac{v_{p_i}(\mathbf{p}, \boldsymbol{\omega})}{v_{\omega_k}(\mathbf{p}, \boldsymbol{\omega})} = -\boldsymbol{\tau}^{b,k} \cdot \mathbf{S}_j^{C,k} + \frac{B_{p_i}}{B_{\omega_k}}.$$

□

**Proof of Proposition 9.14** *Case of an inattentive consumer.* Call  $q_a = p_a + \tau_a$ . Equilibrium requires  $q_a = q_b = 1$ . Competitive pricing in good 1 requires that firms choose inputs according to:  $\max_{l_{ia}, l_{ib}} p_i \left( \frac{l_{ia}}{\alpha_i} \right)^{\alpha_i} \left( \frac{l_{ib}}{1-\alpha_i} \right)^{1-\alpha_i} - (1 + \tau_{ia}) l_{ia} - l_{ib}$  with  $\tau_{0a} = 0$ . Hence, the equilibrium price is  $p_i = (1 + \tau_{ia})^{\alpha_i}$ , and input use features  $(1 + \tau_{ia}) l_{ia} = \alpha_i p_i y_i$ , so  $l_{ia} = \alpha_i (1 + \tau_{ia})^{\alpha_i - 1} y_i$  and  $l_{ib} = (1 - \alpha_i) (1 + \tau_{ia})^{\alpha_i} y_i$

The planning problem is  $\max_{\tau_{1a}} L(p_1)$  with  $p_1 = (1 + \tau_{1a})^{\alpha_1}$ , so that:

$$\begin{aligned}L(p_1) &= c_0 + U^s(c_1(p_1)) - \xi_* c_1(p_1) - \sum_{i=0}^1 (l_{ia}(p_1) + l_{ib}(p_1)) \\ &= U^s(c_1(p_1)) - \left( \xi_* + \alpha_1 (1 + \tau_{1a})^{\alpha_1 - 1} + (1 - \alpha_1) (1 + \tau_{1a})^{\alpha_1} \right) c_1(p_1) \text{ as } c_0 = (l_{0a}(p_1) + l_{0b}(p_1)) \\ &= U^s(c_1(p_1)) - \left( \xi_* + \alpha_1 p_1^{1 - \frac{1}{\alpha_1}} + (1 - \alpha_1) p_1 \right) c_1(p_1).\end{aligned}$$

with  $c_1(p_1) = U'^{-1}(p_1)$ .

Hence, as  $U^{s'}(c_1(p_1)) = p_1$ ,

$$L_{p_1} = c'_1(p_1) \left[ p_1 - \left( \xi_* + \alpha_1 p_1^{1-\frac{1}{\alpha_1}} + (1-\alpha_1)p_1 \right) \right] - \left( (\alpha_1 - 1) p_1^{-\frac{1}{\alpha_1}} + (1-\alpha_1) \right) c_1(p_1).$$

When there is production efficiency,  $p_1 = 1$  and,

$$\begin{aligned} L_{p_1|\tau_{1a}=0} &= c'_1(p_1) [U^{s'}(c_1(p_1)) - (\xi_* + 1)] \\ &= -\xi_* c'_1(p_1) > 0. \end{aligned}$$

Hence, production efficiency is not an optimum. Starting from it, it is optimal to increase the tax  $\tau_{1a}$  to discourage the production of good 1, increase its price, and discourage its consumption.

*Case of an attentive consumer.* It is enough to do a Pigouvian tax  $\tau_1^c = \xi_*$ , and restore production efficiency ( $\tau_{1a} = 1$ ). Then, we achieve the first best.  $\square$

**Proof of Proposition 9.15** Suppose that  $\phi = \phi^s$ . Let  $e(\phi, q)$  be the expenditure function associated with  $\phi(c_1, \dots, c_n)$ . Since  $\phi$  is homogeneous of degree 1, we have  $e(\phi, q) = \phi e(1, q)$ . Consider a non homogeneous tax system with associated prices  $q$ . Tax revenues are

$$\sum_h \phi^h \sum_{i=1}^n (q_i - p_i) e_{q_i}(1, q).$$

Now consider a reformed uniform tax system with associated prices  $\hat{q}_i = x p_i$  for some scalar  $x$ , which delivers the same  $c_0$  and the same  $\phi^h$  for all  $h$ . We just need to solve in  $x$  the following equation

$$e(1, q) = e(1, xp).$$

The reformed tax system leaves the experienced utility of all agents identical (this step crucially uses  $\phi = \phi^s$ ). We claim that the reformed tax system also raises more revenues. This concludes the proof that the optimal tax system must be uniform. Laroque (“Indirect taxation is superfluous under separability and taste homogeneity: A simple proof,” *Economics Letters* 2005) presents related arguments). This amounts to showing that

$$\sum_{i=1}^n (q_i - p_i) e_{q_i}(1, q) < \sum_{i=1}^n (\hat{q}_i - p_i) e_{q_i}(1, \hat{q}),$$

or using  $\sum_{i=1}^n q_i e_{q_i}(1, q) = e(1, q) = e(1, \hat{q}) = \sum_{i=1}^n \hat{q}_i e_{q_i}(1, \hat{q})$ , this amounts to showing that

$$0 < \sum_{i=1}^n p_i [e_{q_i}(1, q) - e_{q_i}(1, \hat{q})],$$

or equivalently since  $\hat{q}_i = xp_i$ , to showing that

$$0 < \sum_{i=1}^n \hat{q}_i [e_{q_i}(1, q) - e_{q_i}(1, \hat{q})],$$

which holds by a straight revealed preference argument.

□

## 11.2 Additional Derivations for the Mirrlees Problem

### 11.2.1 Intermediary results for the Mirrlees problem

**Impact of a change in taxes on earnings and individual utility** We first study the impact of a small change  $\delta q_{z^*}$  of the marginal retention rate at  $z^*$  and how it affects labor supply at  $z$  (e.g. via misperceptions). We simultaneously study the impact of a lump-sum (independent of  $z$ ) virtual income change  $\delta K$ . It will prove conceptually and notationally useful to define:

$$\overline{\zeta}_{Q_{z^*}}^c(z) = \zeta_{Q_{z^*}}^c(z) + \zeta^c(z) \delta_z(z^*), \quad (116)$$

where  $\delta_z$  is a Dirac distribution at point  $z$ . Hence, as  $\zeta_{Q_{z^*}}^c(z)$  was a potentially smooth function of  $z^*$ ,  $\overline{\zeta}_{Q_{z^*}}^c(z)$  is a generalized function of  $z^*$ , in the sense of the theory of distributions. From now on, we mostly use our notation convention of dropping the dependency on  $z$ .

**Lemma 11.1** (Impact of changes in taxes on behavior and welfare) *Suppose that there is a change  $(\delta q_{z^*})_{z^* \geq 0}$  to marginal retention rate schedule and a lump sum increase in revenue  $\delta K$ . The impact on earnings and agent's welfare is:*

$$\delta z = \frac{\eta \delta K + z \int_0^\infty \overline{\zeta}_{Q_{z^*}}^c \delta q_{z^*} dz^*}{q - \zeta^c z R''}, \quad (117)$$

$$\frac{\delta v}{v_r} = \delta K - z \frac{\tau^b}{q} \left( \zeta^c R'' \delta z + \int_0^\infty \overline{\zeta}_{Q_{z^*}}^c \delta q_{z^*} dz^* \right). \quad (118)$$

In these equations, the integrals involving  $z \int_0^\infty \overline{\zeta}_{Q_{z^*}}^c \delta q_{z^*} dz^*$  should be understood in the sense of the theory of distributions as  $z \zeta^c(z) \delta q_z + \int_{z^*=0}^\infty z \zeta_{Q_{z^*}}^c(z) \delta q_{z^*} dz^*$  (reintroducing in these equations the dependency on  $z$ ), leading to

$$\delta z = \frac{\eta(z) \delta K + z \zeta^c(z) \delta q_z + \int_0^\infty z \zeta_{Q_{z^*}}^c(z) \delta q_{z^*} dz^*}{q(z) - \zeta^c(z) z R''(z)},$$

$$\frac{\delta v}{v_r} = \delta K - z \frac{\tau^b}{q} \zeta^c(z) R''(z) \delta z - z \frac{\tau^b(z)}{q(z)} \zeta^c(z) \delta q_z - \int_{z^*=0}^\infty z \frac{\tau^b(z)}{q(z)} \zeta_{Q_{z^*}}^c(z) \delta q_{z^*} dz^*.$$

To interpret the economics of (117), start with an increase in income  $\delta K$ . It has, first, an impact

on labor supply: it creates a direct change in earnings supply equal to  $\frac{\eta}{q}\delta K$ . The additional term  $\zeta^c z R''$  in the denominator of (117) is more subtle and arises from the fact that as the agent adjusts his labor supply, he experiences a different marginal tax rate (which changes as  $R''\delta z$ ), leading to an additional change in income  $\frac{\zeta^c}{q}z R''\delta z$ . The final expression solves for  $\delta z$  as a fixed point. The term  $z \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c \delta q_{z^*} dz^*$  reflects the impact of a change in the marginal tax rate on earnings. The difference with Saez (2001) is that it is non-zero even when the change in the tax schedule occurs at  $z^* \neq z$ . This is because when agents have behavioral biases, a change of the marginal rate at  $z^*$  potentially affects the perceived tax at  $z$ .

In (118), the term  $\delta K$  is a mechanical income effect and is the only term present in the traditional model of Saez (2001). The term  $-z \frac{\tau^b}{q} \left( \zeta^c R'' \delta z + \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c \delta q_{z^*} dz^* \right)$  represent the welfare impact arising from changes in behavior (as the envelope theorem no longer applies) because of misoptimization, respectively, because movements in labor supply change the marginal tax rate  $(-z \frac{\tau^b}{q} \zeta^c R'' \delta z)$  along the initial schedule and because of changes in the tax schedule itself  $(-z \frac{\tau^b}{q} \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c \delta q_{z^*} dz^*)$ .

**Impact of a change in taxes on social welfare** We next study the impact of the above changes on welfare. Following Saez (2001), we call  $h(z)$  the density of agents with earnings  $z$  at the optimum, and  $H(z) = \int_0^z h(z') dz'$ . We also define the virtual density  $h^*(z) = \frac{q(z)}{q(z) - \zeta^c z R''(z)} h(z)$ , which can also be written as  $\frac{1 - T'(z)}{1 - T'(z) + \zeta^c z T''(z)} h(z)$ .

**Lemma 11.2** *Under the conditions of the Lemma 11.1, the change in the government objective function associated with the agent is*

$$\delta L(z) = (\gamma(z) - 1) \delta K + \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} \frac{h^*(z)}{h(z)} z \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c(z) \delta q_{z^*} dz^*, \quad (119)$$

where  $\gamma(z)$  is the marginal social utility of income:

$$\gamma(z) = g(z) + \eta(z) \frac{\tilde{\tau}^b(z)}{1 - T'(z)} + \eta(z) \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} \frac{h^*(z)}{h(z)}. \quad (120)$$

This definition of the social marginal utility of income  $\gamma(z)$  is similar to the one we encountered in the Ramsey problem. It encompasses the direct impact of one extra dollar on the agent's welfare (the  $g(z)$  term) and the impact coming from a change in labor supply on tax revenues  $(\frac{T'(z)}{1 - T'(z)} \eta(z) \frac{h^*(z)}{h(z)})$ . Compared to Saez (2001), it features a new term arising from the failure of the envelope theorem,  $\eta \frac{\tilde{\tau}^b(z)}{1 - T'(z)} \left( 1 - \frac{h^*(z)}{h(z)} \right)$ .

The effect on the government objective function (119) is much like in the many-person Ramsey of Proposition 2.1. The term  $(\gamma(z) - 1) \delta K$  is a mechanical effect, abstracting from changes in behavior. As the government gives (back)  $\delta K$  to agent, the impact on revenues is  $-\delta K$ , while the impact on the agent is valued as  $\gamma(z) \delta K$ . Next, there is a substitution effect  $\frac{T'(z)}{1 - T'(z)} \frac{h^*}{h} z \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c(z) \delta q_{z^*} dz^*$ :

as the agent changes his labor supply, there is a change in tax revenues proportional to

$$\frac{T'(z)}{1-T'(z)} \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c(z) \delta q_{z^*} dz^*.$$

Third, there is a misoptimization term,  $\frac{-\tilde{\tau}^b(z) h^*(z)}{1-T'(z) h(z)} z \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c(z) \delta q_{z^*} dz^*$ .

We also note the following first order condition for the intercept of the tax schedule,  $r_0$ .

**Lemma 11.3** *At the optimum,*

$$\int_0^\infty \left( 1 - \gamma(z) - \frac{T'(z) - \tilde{\tau}^b(z) h^*(z)}{1-T'(z) h(z)} z \zeta_{r_0}^c(z) \right) h(z) dz = 0. \quad (121)$$

We next state the impact of a marginal change in the tax rate,  $\frac{\partial L}{\partial \tau_{z^*}} \equiv -\frac{\partial L}{\partial q_{z^*}}$ .

**Proposition 11.1** (Impact of a local change on the marginal tax rate on the government objective function) *We have*

$$\frac{\partial L}{\partial \tau_{z^*}} = \int_{z^*}^\infty (1 - \gamma(z)) h(z) dz - \zeta^c(z^*) \frac{T'(z^*) - \tilde{\tau}^b(z^*)}{1-T'(z^*)} z^* h^*(z^*) - \int_0^\infty \zeta_{Q_{z^*}}^c(z) \frac{T'(z) - \tilde{\tau}^b(z)}{1-T'(z)} z h^*(z) dz. \quad (122)$$

This equation involves an equality between two generalized functions of  $z^*$ . This is the income tax equivalent of the formula in Proposition 2.1 for the many-person Ramsey. The three terms in (122) correspond to the, by now familiar, mechanical ( $\int_{z^*}^\infty (1 - \gamma(z)) h(z) dz$ ), substitution ( $-\zeta^c(z^*) \frac{T'(z^*)}{1-T'(z^*)} z^* h^*(z^*)$ ), and misoptimization ( $\zeta^c(z^*) \frac{\tilde{\tau}^b(z^*)}{1-T'(z^*)} z^* h^*(z^*) - \int_0^\infty \zeta_{Q_{z^*}}^c(z) \frac{T'(z) - \tilde{\tau}^b(z)}{1-T'(z)} z h^*(z) dz$ ) effects. The first two terms are exactly as in Saez (2001), and the third one is new as it is present only with behavioral agents. We will describe its meaning shortly. We also note that formula (122) can be written in a more compact way as:

$$\frac{\partial L}{\partial \tau_{z^*}} = \int_{z^*}^\infty (1 - \gamma(z)) h(z) dz - \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c(z) \frac{T'(z) - \tilde{\tau}^b(z)}{1-T'(z)} z h^*(z) dz. \quad (123)$$

### 11.2.2 Proofs for the Mirrlees results

**Notations and Derivation of relation (86)** We take the material from section 7.1. The extended good is the two-dimensional  $\mathbf{c} = (c, z)$ , the (generalized) price vector is  $\mathbf{q} = (1, q, \mathbf{Q}, r_0)$ . The budget function is  $B(\mathbf{c}, \mathbf{q}) = q_1 c_1 - q_2 c_2 = c - qz$ , so that the budget constraint is  $B(\mathbf{c}, \mathbf{q}) \leq r$ . Note that the Saez  $r$  is also the  $w$  in the rest of the paper (as the budget constraint is generally expressed as  $B(\mathbf{c}, \mathbf{q}) \leq r$ ); we still found useful to stick here to the Saez notations; so in the derivations of the Mirrlees case, we will use  $r$  and  $w$  interchangeably, depending on what the context calls for.



Applying definition (36) gives

$$\boldsymbol{\tau}^b = (1, -q) - \frac{(u_c, u_z)}{v_r}. \quad (124)$$

We know that  $c_{Q_{z^*}} = qz_{Q_{z^*}}$  (which comes from differentiating  $c = qz + r$  w.r.t.  $Q_{z^*}$ ), so

$$\mathbf{S}_{Q_{z^*}}^C = (c_{Q_{z^*}}, z_{Q_{z^*}})' = (q, 1)' z_{Q_{z^*}}.$$

Proposition 7.1 implies:

$$\begin{aligned} \frac{v_{Q_{z^*}}(\mathbf{q}, r)}{v_r(\mathbf{q}, r)} &= -\boldsymbol{\tau}^b(\mathbf{q}, r) \cdot \mathbf{S}_{Q_{z^*}}^C(\mathbf{q}, r) = -\boldsymbol{\tau}^b(\mathbf{q}, r) (q, 1)' z_{Q_{z^*}} = -\tau^b z_{Q_{z^*}} \\ &= -\tau^b \frac{z}{q} \zeta_{Q_{z^*}}^c, \end{aligned}$$

as we defined

$$\tau^b = \boldsymbol{\tau}^b(\mathbf{q}, r) \cdot (q, 1) = -\frac{qu_c + u_z}{v_r}. \quad (125)$$

Likewise,  $c - qz = r$  implies (taking the derivative w.r.t.  $q$ ):  $c_q - qz_q - z = 0$  and (taking the derivative w.r.t.  $r$ )  $c_r - qz_r = 1$ , so

$$\begin{aligned} \mathbf{S}_q^C(\mathbf{q}, r) &= \mathbf{c}_q - \mathbf{c}_r z = (c_q - c_r z, z_q - z_r z) \\ &= (qz_q + z - (qz_r + 1)z, z_q - z_r z) = (q, 1) (z_q - z_r z) \\ &= (q, 1) z \frac{\zeta^c}{q}. \end{aligned} \quad (126)$$

Proposition 7.1 implies:

$$\begin{aligned} \frac{v_q(\mathbf{q}, r)}{v_r(\mathbf{q}, r)} &= z - \boldsymbol{\tau}^b(\mathbf{q}, r) \cdot \mathbf{S}_q^C(\mathbf{q}, r) = z - \boldsymbol{\tau}^b(\mathbf{q}, r) \cdot (q, 1) z \frac{\zeta^c}{q} \\ &= z - \frac{z}{q} \tau^b \zeta^c. \end{aligned}$$

□

**Proof of Elasticity relations (87) in the Mirrlees framework: Concrete values of the general model in the misperception case** Now consider the model with misperception. As above, the extended good is  $\mathbf{c} = (c, z)$ , and the (generalized) price  $\mathbf{q} = (1, q, \mathbf{Q}, r_0)$ , and the budget function is  $B(\mathbf{c}, \mathbf{q}) = c_1 q_1 - c_2 q_2 = c - qz$ , so that

$$B_c(\mathbf{c}, \mathbf{q}) = (1, -q). \quad (127)$$

We use (42)

$$\begin{aligned}
\tau^b &= B_c(\mathbf{c}, \mathbf{q}) - \frac{B_c(\mathbf{c}, \mathbf{q}^s)}{B_c(\mathbf{c}, \mathbf{q}^s) \cdot \mathbf{c}_r(\mathbf{q}, r)} \\
&= (1, -q) - \frac{(1, -q^s)}{(1, -q^s) \cdot (qz_r + 1, z_r)} \text{ as } c(q, r) = qz(q, r) + r \text{ gives } c_r = qz_r + 1. \\
&= (1, -q) - \frac{(1, -q^s)}{1 + (q - q^s)z_r} \\
\tau^b &= (1, -q) - \frac{(1, -q^s)}{1 + (q - q^s)\frac{\eta}{q}}. \tag{128}
\end{aligned}$$

Next, recall (125),

$$\begin{aligned}
\tau^b &= \tau^b(\mathbf{q}, r) \cdot (q, 1) \\
&= \left[ (1, -q) - \frac{(1, -q^s)}{1 + (q - q^s)\frac{\eta}{q}} \right] \cdot (q, 1) = \frac{q^s - q}{1 + (q - q^s)\frac{\eta}{q}} \\
&= \frac{\tau - \tau^s}{1 - (\tau - \tau^s)\frac{\eta}{q}}. \tag{129}
\end{aligned}$$

using  $q = 1 - \tau$ ,  $q^s = 1 - \tau^s$ . Thus we have proven (88).

Next, we calculate  $\zeta^c$ . We call  $\mathbf{e} = (0, 1)$  the vector singling earnings on the vector  $\mathbf{c} = (c, z)$ . We apply (41) with  $p_j = q$ , the price of earnings. We have:

$$\begin{aligned}
\mathbf{e} \cdot \mathbf{S}_j^H &= \mathbf{e} \cdot \mathbf{S}^r(\mathbf{p}, r) \cdot \mathbf{p}_{p_j}^s(\mathbf{p}, r) = z_q^r \frac{\partial q^s}{\partial q} = \frac{z}{q} \zeta^{c,r} m_{zz} \\
\mathbf{e} \cdot \mathbf{S}_j^H &= \frac{z}{q} \zeta^{c,r} m_{zz}. \tag{130}
\end{aligned}$$

Next, using the notation  $D_j$  of Proposition 7.1,

$$\begin{aligned}
D_j &= -\tau^b \cdot \mathbf{S}_j^H \text{ by (35)} \\
&= [B_c(\mathbf{p}, \mathbf{c}) - B_c(\mathbf{p}^s, \mathbf{c})] \cdot \mathbf{S}_j^H \text{ by (43)} \\
&= [(1, q) - (1, q^s)] \cdot \mathbf{S}_j^H \text{ by (127)} \\
&= (q - q^s) \mathbf{e} \cdot \mathbf{S}_j^H \text{ as } \mathbf{e} = (0, 1) \\
&= (q - q^s) \frac{z}{q} \zeta^{c,r} m_{zz} \text{ by (130)}.
\end{aligned}$$

We record:

$$D_j = (q - q^s) \frac{z}{q} \zeta^{c,r} m_{zz}. \tag{131}$$

Next, we apply (39):  $\mathbf{S}_j^C = \mathbf{S}_j^H + \mathbf{c}_r D_j$ , which implies:

$$\begin{aligned}
\mathbf{e} \cdot \mathbf{S}_j^C &= \mathbf{e} \cdot \mathbf{S}_j^H + \mathbf{e} \cdot \mathbf{c}_r D_j \\
&= \frac{z}{q} \zeta^{c,r} m_{zz} + \frac{\eta}{q} D_j \text{ as } \mathbf{e} \cdot \mathbf{c}_r = z_r = \frac{\eta}{q} \\
&= \frac{z}{q} \zeta^{c,r} m_{zz} \left( 1 + \frac{\eta}{q} (q - q^s) \right) \\
&= \frac{z}{q} \left( 1 - \eta \frac{\tau - \tau^s}{q} \right) \zeta^{c,r} m_{zz} \text{ as } q = 1 - \tau, q^s = 1 - \tau^s.
\end{aligned}$$

Now, as  $\zeta^c = \frac{q}{z} \mathbf{e} \cdot \mathbf{S}_j^C$  is the compensated earnings elasticity (see e.g. (126)):

$$\begin{aligned}
\zeta^c &= \frac{q}{z} \mathbf{e} \cdot \mathbf{S}_j^C \\
\zeta^c &= \left( 1 - \eta \frac{\tau - \tau^s}{q} \right) \zeta^{c,r} m_{zz}.
\end{aligned} \tag{132}$$

Exactly the same reasoning (using then  $p_j = q_{z^*}$ , and  $\mathbf{e} \cdot \mathbf{S}_j^H = \frac{z}{q} \zeta^{c,r} m_{zz^*}$ ) shows

$$\zeta_{Q_{z^*}}^c = \left( 1 - \eta \frac{\tau - \tau^s}{q} \right) \zeta^{c,r} m_{zz^*}. \tag{133}$$

Hence, we have proven (87).

**Proof of (89): Decision vs. Experienced utility model** The agent's optimization gives  $q u_c^s + u_z^s = 0$ . Equation (125) gives:

$$\tau^b = -\frac{q u_c + u_z}{v_r} = \frac{u_c \frac{u_z^s}{u_c^s} - u_z}{v_r}.$$

**Dirac / Double bar Notation for the proofs in the Mirrlees framework** We define:

$$\bar{\bar{\zeta}}_{Q_{z^*}}^c(z) = \zeta_{Q_{z^*}}^c(z) + \zeta^c(z) \delta_z(z^*).$$

Informally, this definition means that  $\bar{\bar{\zeta}}_{Q_{z^*}}^c$  is like  $\zeta_{Q_{z^*}}^c(z)$ , but with an extra Dirac term when  $z = z^*$ .

**Proof of Lemma 11.1** We have

$$\begin{aligned}
z &= z(q(z), \mathbf{Q}, r(z)) \\
r(z) &= R(z) - zq(z),
\end{aligned}$$

so

$$\begin{aligned}\delta r &= r'(z) \delta z + \delta r|_{\text{constant } z} = -zq'(z) \delta z + (\delta K - z\delta q_z) \\ &= -zR''\delta z + \delta K - z\delta q_z.\end{aligned}$$

$$\begin{aligned}\delta z &= z_q(q'(z) \delta z + \delta q_z) + \int_0^\infty z_{Q_{z^*}} \delta q_{z^*} dz^* + z_r \delta r \\ &= \frac{z}{q} \zeta^u q'(z) \delta z + \frac{z}{q} \zeta^u \delta q_z + \frac{z}{q} \int_0^\infty \zeta_{Q_{z^*}}^c \delta q_{z^*} dz^* + \frac{\eta}{q} (-zR''\delta z + \delta K - z\delta q_z) \\ &= \frac{z}{q} \zeta^u R'' \delta z + \frac{z}{q} \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c \delta q_{z^*} dz^* + \frac{\eta}{q} (-zR''\delta z + \delta K),\end{aligned}$$

so

$$\delta z = \frac{z \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c \delta q_{z^*} dz^* + \eta \delta K}{q + (\eta - \zeta^u) z R''} = \frac{z \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c \delta q_{z^*} dz^* + \eta \delta K}{q - \zeta^c z R''}.$$

For welfare  $v(q, \mathbf{Q}, r_0, r)$ , we have:

$$\begin{aligned}\frac{\delta v}{v_r} &= \frac{v_q}{v_r} (q'(z) \delta z + \delta q_z) + \int_0^\infty \frac{v_{Q_{z^*}}}{v_r} \delta q_{z^*} dz^* + \delta r, \\ &= z \left( 1 - \frac{\tau^b \zeta^c}{q} \right) (R'' \delta z + \delta q_z) + z \int_0^\infty \left( -\frac{\tau^b \zeta_{Q_{z^*}}^c}{q} \right) \delta q_{z^*} dz^* - zR''\delta z + \delta K - z\delta q_z, \\ \frac{\delta v}{v_r} &= z \left( -\frac{\tau^b \zeta^c}{q} \right) R'' \delta z + z \left( \int_0^\infty \left( -\frac{\tau^b \zeta_{Q_{z^*}}^c}{q} \right) \delta q_{z^*} dz^* + \left( -\frac{\tau^b \zeta^c}{q} \right) \delta q_z \right) + \delta K, \\ &= \delta K - z \frac{\tau^b}{q} \zeta^c R''(z) \delta z - \frac{\tau^b}{q} z \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c \delta q_{z^*} dz^*, \\ &= \delta K - z \frac{\tau^b}{q} \left( \zeta^c R''(z) \delta z + \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c \delta q_{z^*} dz^* \right).\end{aligned}$$

□

**Proof of Lemma 11.2** Observe that

$$\frac{h^*(z)}{q} = \frac{h(z)}{q - \zeta^c z R''(z)},$$

so that  $q - \zeta^c z R''(z) = \frac{qh}{h^*}$  and

$$zR'' = q \frac{h^* - h}{\zeta^c h^*}. \quad (134)$$

We have

$$\begin{aligned}
\delta T &= \delta(z(1-q(z)) - r), \\
&= (1-q_z)\delta z - zq'(z)\delta z - z\delta q_z - \delta r, \\
&= (1-q_z)\delta z - zR''\delta z - z\delta q_z - (-zR''\delta z + \delta K - z\delta q_z), \\
&= T'(z)\delta z - \delta K.
\end{aligned}$$

We also have

$$\begin{aligned}
\delta L &= \delta T + g(z)\frac{\delta v}{v_r}, \\
&= T'(z)\delta z + g(z)\frac{\delta v}{v_r} - \delta K.
\end{aligned}$$

Using Lemma 11.1, we can rewrite this as

$$\delta L = T'(z)\frac{z\int_0^\infty \bar{\bar{\zeta}}_{Q_{z^*}}^c \delta q_{z^*} dz^* + \eta \delta K}{q - \zeta^c z R''} + g(z)\left(\delta K - z\frac{\tau^b}{q}\left(\zeta^c R''(z)\delta z + \int_0^\infty \bar{\bar{\zeta}}_{Q_{z^*}}^c \delta q_{z^*} dz^*\right)\right) - \delta K.$$

Using equation (134) and Lemma 11.1, we can rewrite this as

$$\begin{aligned}
\delta L &= \left[-1 + g(z) + \eta\frac{T'(z)h^*}{q} + g(z)\left(-\frac{\tau^b\zeta^c}{q}\right)\eta\frac{h^* - h}{\zeta^c h}\right]\delta K, \\
&+ \left(-\frac{g(z)\tau^b}{q} + \frac{T'(z)h^*}{q} + g(z)\left(-\frac{\tau^b\zeta^c}{q}\right)\frac{h^* - h}{\zeta^c h}\right)z\int_0^\infty \bar{\bar{\zeta}}_{Q_{z^*}}^c \delta q_{z^*} dz^*, \\
&= (\gamma(z) - 1)\delta K + \frac{h^*}{h}z\frac{T'(z) - \tilde{\tau}^b(z)}{q}\int_0^\infty \bar{\bar{\zeta}}_{Q_{z^*}}^c \delta q_{z^*} dz^*,
\end{aligned}$$

where

$$\gamma(z) = g(z) + \eta\frac{\tilde{\tau}^b(z)}{q} + \frac{T'(z) - \tilde{\tau}^b(z)}{q}\eta\frac{h^*(z)}{h(z)}.$$

□

**Proof of Proposition 11.1** We use the following notations:

$$\begin{aligned}
\bar{\bar{F}}(z_*) &= \int_0^\infty \left[-\bar{\bar{\zeta}}_{Q_{z^*}}^c \frac{\tilde{\tau}^b(z)}{q(z)} + \zeta_{Q_{z^*}}^c \frac{T'(z)}{q}\right] zh^*(z) dz, \\
J(z^*) &= \zeta^c z^* \frac{T'(z^*)}{q} h^*(z^*). \\
\bar{\bar{\bar{F}}}(z_*) &= \int_0^\infty \bar{\bar{\zeta}}_{Q_{z^*}}^c \frac{T'(z) - \tilde{\tau}^b(z)}{q(z)} zh^*(z) dz = \bar{\bar{F}}(z_*) + J(z_*),
\end{aligned}$$

We consider a change  $\delta q_{z^*}$  at  $z^*$ . This leads to a lump-sum change  $\delta K = 1_{z > z^*} \delta q_{z^*}$ . Hence,

Lemma 11.2 gives the change in the government objective function

$$\delta L(z) = (\gamma(z) - 1) 1_{z > z^*} \delta q_{z^*} + \frac{h^*}{h} z \frac{T'(z) - \tilde{\tau}^b(z)}{q} \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c \delta q_{z^*} dz^*.$$

The total change is

$$\begin{aligned} \delta L &= \int_0^\infty \delta L(z) h(z) dz, \\ \frac{\delta L}{\delta q_{z^*}} &= \int_{z^*}^\infty (\gamma(z) - 1) h(z) dz + \int_0^\infty \left[ \frac{T'(z) - \tilde{\tau}^b(z)}{q} \bar{\zeta}_{Q_{z^*}}^c \right] \frac{h^*}{h} z h(z) dz, \\ \frac{\partial L}{\partial \tau_{z^*}} &= -\bar{F}(z_*) + \int_{z^*}^\infty (1 - \gamma(z)) h(z) dz. \end{aligned} \quad (135)$$

We also have

$$\frac{\partial L}{\partial \tau_{z^*}} \equiv -\frac{\partial L}{\partial Q_{z^*}} = -\bar{F}(z^*) + \int_{z^*}^\infty (1 - \gamma(z)) h(z) dz. \quad (136)$$

□

**Proof of Lemma 11.3** Using Lemma 11.2, applied to a change  $\delta r_0$  to all agents, and slightly generalizing, we find

$$\delta L(z) = (\gamma(z) - 1) \delta r_0 + \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} \frac{h^*(z)}{h(z)} z \zeta_{r_0}^c(z) \delta r_0,$$

and  $\delta L = \int \delta L(z) h(z) dz$  should be 0. □

**Proof of Proposition 10.1** Let us now solve for the optimal  $J$ , which ensures  $\frac{\partial L}{\partial Q_{z^*}} = 0$ . We can write

$$-\frac{\partial L}{\partial \tau_{z^*}} = J(z^*) - \int_{z^*}^\infty a(z) dz - b(z^*) + \int_{z^*}^\infty J(z) \rho(z) dz. \quad (137)$$

We use the notations

$$\rho(z) = \frac{\eta}{\zeta^c} \frac{1}{z}, \quad (138)$$

$$a(z) = (1 - g(z)) h(z) - \rho g(z) z \left( -\tau^b(z) \frac{\zeta^c}{q(z)} \right) (h^*(z) - h(z)), \quad (139)$$

$$b(z) = -\bar{F}(z_*). \quad (140)$$

$a(z)$  is the effect of giving \$1 to agent  $z$  (that's the  $(1 - g(z)) h(z)$  term), corrected from distortions from the non-linearity of the income tax.

$\bar{F}(z^*)$  is the part impact on the government's objective function of increase  $\delta q_{z^*}$ , coming from

the distortions from perceptions

$$\begin{aligned}\bar{\bar{F}}(z^*) &= \int_0^\infty \left[ -\bar{\zeta}_{Q_{z^*}}^c \frac{\tilde{\tau}^b(z)}{q(z)} + \zeta_{Q_{z^*}}^c \frac{T'(z)}{q} \right] z h^*(z) dz = F(z^*) + z g(z) \left( -\tau^b \frac{\zeta^c}{q(z)} \right) h^*(z), \\ F(z^*) &= \int_0^\infty \left[ \zeta_{Q_{z^*}}^c \frac{T'(z) - \tilde{\tau}^b(z)}{q} \right] z h^*(z) dz.\end{aligned}$$

We note that

$$\begin{aligned}a &= (1 - \gamma) h + \eta \frac{T'(z) h^*}{q} \frac{h^*}{h}, \\ a &= (1 - \gamma) h + \rho J.\end{aligned}\tag{141}$$

We also have

$$\begin{aligned}J(z^*) &= \int_{z^*}^\infty a(z) dz + b(z^*) - \int_{z^*}^\infty J(z) \rho(z) dz, \\ \dot{J} &= -a + \dot{b} + J\rho, \\ \frac{d}{dz} \left[ J(z) e^{-\int_0^z \rho(s) ds} \right] &= e^{-\int_0^z \rho(s) ds} \left( -a(z) + \dot{b}(z) \right), \\ J(z) e^{-\int_0^z \rho(s) ds} &= C + \int_z^\infty e^{-\int_0^{z'} \rho(s) ds} \left( a(z') - \dot{b}(z') \right) dz', \\ J(z) &= \int_z^\infty e^{-\int_z^{z'} \rho(s) ds} \left( a(z') - \dot{b}(z') \right) dz'.\end{aligned}\tag{142}$$

Integrating by parts, we get

$$\begin{aligned}\int_z^\infty e^{-\int_z^{z'} \rho(s) ds} \dot{b}(z') dz' &= \left[ e^{-\int_z^{z'} \rho(s) ds} b(z') \right]_z^\infty + \int_z^\infty e^{-\int_z^{z'} \rho(s) ds} \rho(z') b(z') dz' \\ &= -b(z) + \int_z^\infty e^{-\int_z^{z'} \rho(s) ds} \rho(z') b(z') dz',\end{aligned}\tag{143}$$

$$J(z) = b(z) + \int_z^\infty e^{-\int_z^{z'} \rho(s) ds} \left( a(z') - \rho(z') b(z') \right) dz'.\tag{144}$$

We can rewrite this as

$$\begin{aligned}J(z^*) &= -\bar{\bar{F}}(z^*) + \int_{z^*}^\infty e^{-\int_{z^*}^z \rho(s) ds} \left( a(z) + \rho \bar{\bar{F}}(z) \right) dz \\ &= -z^* g(z^*) \left( -\tau^b(z^*) \frac{\zeta^c(z^*)}{q(z^*)} \right) h^*(z^*) - F(z^*) \\ &\quad + \int_{z^*}^\infty e^{-\int_{z^*}^z \rho(s) ds} \left[ (1 - g(z)) h(z) + g(z) \rho(z) z \left( -\tau^b(z) \frac{\zeta^c(z)}{q(z)} \right) h(z) + \rho(z) F(z) \right] dz.\end{aligned}\tag{145}$$

Using

$$J(z^*) = \zeta^c(z^*) z^* h^*(z^*) \frac{T'(z^*)}{1 - T'(z^*)}$$

and rearranging gives

$$\begin{aligned} & \zeta^c(z^*) z^* h^*(z^*) \frac{T'(z^*) - \tilde{\tau}^b(z^*)}{1 - T'(z^*)} + F(z_*) - \int_{z^*}^{\infty} e^{-\int_{z^*}^z \rho(s) ds} \rho(z) F(z) dz \\ &= \int_{z^*}^{\infty} e^{-\int_{z^*}^z \rho(s) ds} \left[ (1 - g(z)) h(z) + g(z) \rho(z) z \left( -\tau^b(z) \frac{\zeta^c(z)}{q(z)} \right) h(z) \right] dz, \end{aligned}$$

which can be rewritten to get the announced formula

$$\begin{aligned} & \frac{T'(z^*) - \tilde{\tau}^b(z^*)}{1 - T'(z^*)} + \frac{1}{\zeta^c z^* h^*(z^*)} \int_{z=0}^{\infty} \left( \zeta_{Q_{z^*}}^c(z) - \int_{z'=z^*}^{\infty} e^{-\int_{z^*}^{z'} \rho(s) ds} \rho(z') \zeta_{Q_{z'}}^c(z) dz' \right) \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} z h^*(z) dz, \\ &= \frac{1}{\zeta^c} \frac{1 - H(z^*)}{z^* h^*(z^*)} \int_{z^*}^{\infty} e^{-\int_{z^*}^z \rho(s) ds} \left( 1 - g(z) - \eta \frac{\tilde{\tau}^b(z)}{q(z)} \right) \frac{h(z)}{1 - H(z^*)} dz. \end{aligned}$$

□

## 12 Basic behavioral consumer theory with linear budget constraints

### 12.1 Traditional theory: Recap

The objects in the traditional theory are  $e(\mathbf{p}, u)$ ,  $v(\mathbf{p}, w)$ ,  $\mathbf{h}(\mathbf{p}, u) = \arg \min_{\mathbf{c}} \mathbf{p} \cdot \mathbf{c}$  s.t.  $u(\mathbf{c}) = u$ . Let us prove the traditional relations – a warm up for the proof in the behavioral case.

Roy's identity is proven as follows:  $v(\mathbf{p}, w) = \max_{\mathbf{c}} u(\mathbf{c}) + \lambda(w - \mathbf{p} \cdot \mathbf{c})$ , so  $v_{p_j} = -\lambda c_j$ ,  $v_w = \lambda$ , so:

$$v_{p_j} + v_w c^j = 0. \quad (146)$$

Shepard's lemma is proven as follows: The envelope theorem gives  $e_p(\mathbf{p}, u) = \mathbf{h}(\mathbf{p}, u)$ , i.e.

$$e_{p_i}(\mathbf{p}, u) = h^i, \quad (147)$$

and differentiating once more gives:

$$e_{p_i p_j} = h_{p_j}^i(\mathbf{p}, u) = S_{ij}. \quad (148)$$



We have  $\mathbf{c}(\mathbf{p}, w) = \mathbf{h}(\mathbf{p}, v(\mathbf{p}, w))$ , which implies

$$\begin{aligned}\mathbf{c}_{p_j} &= \mathbf{h}_{p_j} + \mathbf{h}_u v_{p_j} \\ \mathbf{c}_w &= \mathbf{h}_u v_w,\end{aligned}$$

and because of Roy, we have Slutsky's relation:

$$\mathbf{c}_{p_j} + \mathbf{c}_w c_j = \mathbf{h}_{p_j} = \mathbf{S}_j,$$

i.e.

$$c_{p_j}^i + c_w^i c_j = h_{p_j}^i = S_{ij}. \quad (149)$$

## 12.2 Behavioral version with perceived prices

The sparse max demand

$$\operatorname{smax}_{\mathbf{c}|\mathbf{p}^s} u(\mathbf{c}) \quad \text{s.t.} \quad \mathbf{p} \cdot \mathbf{c} \leq w$$

of a behavioral agent perceiving prices  $\mathbf{p}^s$  (while true prices are  $\mathbf{p}$  and the true budget is  $w$ ) is:

$$\mathbf{c}(\mathbf{p}, \mathbf{p}^s, w) = \mathbf{c}^r(\mathbf{p}^s, w'(\mathbf{p}, \mathbf{p}^s, w)), \quad (150)$$

where perceived budget  $w'$  satisfies:

$$\mathbf{p} \cdot \mathbf{c}^r(\mathbf{p}^s, w'(\mathbf{p}, \mathbf{p}^s, w)) = w. \quad (151)$$

We call  $\mathbf{c}^r(\mathbf{p}^s, w')$  the rational Marshallian demand under prices  $\mathbf{p}^s$  and budget  $w'$ .

We define  $v(\mathbf{p}, \mathbf{p}^s, w) = u(\mathbf{c}(\mathbf{p}, \mathbf{p}^s, w))$ . The expenditure function is  $e(\mathbf{p}, \mathbf{p}^s, u) = \min_w w$  s.t.  $v(\mathbf{p}, \mathbf{p}^s, w) \geq u$ . We define the Hicksian demand  $\mathbf{h}(\mathbf{p}, \mathbf{p}^s, \bar{u}) = \operatorname{argmax}_{\mathbf{c}|\mathbf{p}^s} -\mathbf{p} \cdot \mathbf{c}$  s.t.  $u(\mathbf{c}) = \bar{u}$  with perception  $\mathbf{p}^s$  by the agent, which gives  $\mathbf{h}(\mathbf{p}, \mathbf{p}^s, u) = \mathbf{h}^r(\mathbf{p}^s, u)$ . So  $\mathbf{c}(\mathbf{p}, \mathbf{p}^s, w) = \mathbf{h}^r(\mathbf{p}^s, v(\mathbf{p}, \mathbf{p}^s, w))$ .

Here we derive Shepard, Roy etc. for this behavioral model. This generalizes [Gabaix \(2014\)](#), which derives similar relations under the assumption that  $\mathbf{p}^s = \mathbf{M}\mathbf{p} + (1 - \mathbf{M})\mathbf{p}^d$ .

We call

$$\mathbf{S}_j^r = \mathbf{h}_{p_j^s}^r(\mathbf{p}^s, u)$$

the rational Slutsky matrix, and  $\mathbf{S}_j^r = (S_{ij}^r)_{i=1\dots n}$  the vector of Slutsky sensitivities with respect to price  $p_j$ .

**Proposition 12.1** (Generalized Shepard's lemma) *Given the function  $e(\mathbf{p}, \mathbf{p}^s, u) = \mathbf{p} \cdot \mathbf{h}^r(\mathbf{p}^s, u)$ ,*

we have:

$$\begin{aligned} e_{p_j} &= c_j \\ e_{p_j^s} &= (\mathbf{p} - \mathbf{p}^s) \cdot \mathbf{S}_j^r. \end{aligned}$$

**Proof.** We have:  $e(\mathbf{p}, \mathbf{p}^s, u) = \mathbf{p} \cdot \mathbf{h}^r(\mathbf{p}^s, u)$ , so

$$\begin{aligned} e_{p_j} &= \mathbf{h}_j^r = c_j \\ e_{p_j^s} &= \mathbf{p} \cdot \mathbf{h}_{p_j^s}^r(\mathbf{p}^s, u) = (\mathbf{p} - \mathbf{p}^s) \cdot \mathbf{h}_{p_j^s}^r(\mathbf{p}^s, u). \end{aligned}$$

Indeed, we have  $\mathbf{p}^s \cdot \mathbf{h}^r(\mathbf{p}^s, u) = 0$ . To prove this, observe that

$$\begin{aligned} \mathbf{q} \cdot \mathbf{h}_{q_j}^r(\mathbf{q}, u) &= \sum_i q_i h_{q_j}^i = \sum_i q_i h_{q_i}^j \text{ by symmetry} \\ &= 0 \text{ as } h^j(q, u) \text{ is homogeneous of degree 0.} \end{aligned}$$

□

**Proposition 12.2** (Generalized Roy's identity). *Given the function  $v(\mathbf{p}, \mathbf{p}^s, w)$ , we have:*

$$\begin{aligned} \frac{v_{p_j}}{v_w} &= -c_j \\ \frac{v_{p_j^s}}{v_w} &= (\mathbf{p}^s - \mathbf{p}) \cdot \mathbf{S}_j^r = D_j^s, \end{aligned}$$

*i.e.*  $\frac{v_{p_j^s}}{v_w} = \sum_i (p_i^s - p_i) \mathbf{S}_{ij}^r$ .

To gain intuition for the term in  $v_{p_j^s}$ , observe that:

$$v_{p_j^s} \cdot \delta p^s \geq 0 \text{ with } \delta p^s = 0.01(p - p^s).$$

This is, the agent is better off if his perceived price goes towards the true price.

**Proof of Proposition 12.2**

For a number  $\bar{u}$ , we have the identity  $\bar{u} = v(\mathbf{p}, \mathbf{p}^s, e(\mathbf{p}, \mathbf{p}^s, \bar{u}))$  for all  $\mathbf{p}, \mathbf{p}^s, \bar{u}$ . Deriving w.r.t.  $p_j$  gives:

$$0 = v_{p_j} + v_w e_{p_j} = v_{p_j} + v_w c_j,$$

by the behavioral Shepard's lemma (Proposition 12.1).

Deriving w.r.t.  $p_j^s$  gives:

$$0 = v_{p_j^s} + v_w e_{p_j^s} = v_{p_j^s} + v_w (\mathbf{p} - \mathbf{p}^s) \cdot \mathbf{S}_j^r$$

again by the behavioral Shepard's lemma (Proposition 12.1). □

**Proposition 12.3** (Marshallian demand) Given the consumption function  $\mathbf{c}(\mathbf{p}, \mathbf{p}^s, w)$ , we have:

$$\mathbf{c}_{p_j} = -\mathbf{c}_w c_j \quad (152)$$

$$\mathbf{c}_{p_j^s} = \mathbf{S}_j^r + \mathbf{c}_w D_j^s =: S \quad (153)$$

$$= \mathbf{S}_j^r + \mathbf{c}_w [(\mathbf{p}^s - \mathbf{p}) \cdot \mathbf{S}_j^r] \quad (154)$$

$$= (1 + \mathbf{c}_w (\mathbf{p}^s - \mathbf{p})') \mathbf{S}_j^r. \quad (155)$$

i.e.  $c_{p_j}^i = -c_w^i c_j$  and  $c_{p_j^s}^i = S_{ij}^r + c_w^i D_j^s$ . In addition,  $\mathbf{c}_w = \frac{1}{\mathbf{p} \cdot \mathbf{c}_{w'}^r(\mathbf{p}^s, w')} \mathbf{c}_{w'}^r = \frac{v_w}{\lambda} \mathbf{c}_{w'}^r$ .

The new term is  $c_w^i D_j^s$ . To interpret it, consider again what happens if the agent's perceived price goes towards the true price.  $dp^s = \chi(p - p^s)$ ,  $\chi > 0$ . Then,

$$\begin{aligned} d\mathbf{c} &= \mathbf{c}_{p^s} dp^s = S dp^s + \mathbf{c}_w dE \\ dE &= [(\mathbf{p}^s - \mathbf{p}) \cdot \mathbf{S}^r \cdot dp^s] \geq 0. \end{aligned}$$

The extra term  $dE$  is positive: it's as if the agent became richer. That creates an income effect, and she shifts her consumption  $\mathbf{c}_w dE$ . We can summarize: "If the agent's perceived price goes towards the true price, the agent is better off, and the consumer consumes as if she was richer".

**Proof.** We have  $\mathbf{c}(\mathbf{p}, \mathbf{p}^s, w) = \mathbf{h}^r(\mathbf{p}^s, v(\mathbf{p}, \mathbf{p}^s, w))$ , which implies

$$\mathbf{c}_w = \mathbf{h}_u^r v_w, \quad \mathbf{c}_{p_j} = \mathbf{h}_u^r v_{p_j}. \quad (156)$$

Because of Roy ( $v_{p_j} + c_j v_w = 0$ ), we have:  $\mathbf{c}_{p_j} + \mathbf{c}_w c_j = 0$ .

Also,  $\mathbf{c}(\mathbf{p}, \mathbf{p}^s, w) = \mathbf{h}(\mathbf{p}^s, v(\mathbf{p}, \mathbf{p}^s, w))$  gives:

$$\begin{aligned} \mathbf{c}_{p_j^s}(\mathbf{p}, \mathbf{p}^s, w) &= \mathbf{h}_{p_j^s} + \mathbf{h}_u v_{p_j^s} \\ &= \mathbf{S}_j^r + \mathbf{c}_w \frac{v_{p_j^s}}{v_w} \text{ using } \mathbf{c}_w = \mathbf{h}_u^r v_w \\ &= \mathbf{S}_j^r + \mathbf{c}_w D_j^s \text{ using Proposition 12.2.} \end{aligned}$$

We have (151): so  $\mathbf{p} \cdot \mathbf{c}_{w'}^r \frac{\partial w'}{\partial w} = 1$ , and  $\frac{\partial w'}{\partial w}(\mathbf{p}, \mathbf{p}^s, w) = \frac{1}{\mathbf{p} \cdot \mathbf{c}_{w'}^r(\mathbf{p}^s, w')}$ . So  $\mathbf{c}(\mathbf{p}, \mathbf{p}^s, w) = \mathbf{c}^r(\mathbf{p}^s, w'(\mathbf{p}, \mathbf{p}^s, w))$ , we have  $\mathbf{c}_w^s = \mathbf{c}_{w'}^r \frac{\partial w'}{\partial w}$ .

□

In the traditional model  $v(\mathbf{p}, w) = \max_{\mathbf{c}} u(\mathbf{c}) + \lambda(w - \mathbf{p} \cdot \mathbf{c})$  implies  $v_w = \lambda$ . There is a deviation here, as indicated below.

**Proposition 12.4** (Envelope theorem, modified) Call  $\lambda$  the Lagrange multiplier such that  $u'(\mathbf{c}) = \lambda \mathbf{p}^s$ . We have:

$$\frac{v_w}{\lambda} = \mathbf{p}^s \cdot \mathbf{c}_w(\mathbf{p}, \mathbf{p}^s, w) = 1 + (\mathbf{p}^s - \mathbf{p}) \cdot \mathbf{c}_w(\mathbf{p}, \mathbf{p}^s, w).$$

**Proof.** We have:

$$\mathbf{p} \cdot \mathbf{c}^r(\mathbf{p}^s, w'(\mathbf{p}, \mathbf{p}^s, w)) = w,$$

so  $\mathbf{p} \cdot \mathbf{c}_{w'}^r(\mathbf{p}^s, w') = 1$ , and

$$\frac{\partial w'}{\partial w} = \frac{1}{\mathbf{p} \cdot \mathbf{c}_{w'}^r(\mathbf{p}^s, w')}.$$

Also, given  $\mathbf{c}(\mathbf{p}, \mathbf{p}^s, w) = \mathbf{c}^r(\mathbf{p}^s, w'(\mathbf{p}, \mathbf{p}^s, w))$

$$\mathbf{c}_w = \mathbf{c}_{w'}^r(\mathbf{p}^s, w') \frac{\partial w'}{\partial w} = \frac{\mathbf{c}_{w'}^r(\mathbf{p}^s, w')}{\mathbf{p} \cdot \mathbf{c}_{w'}^r(\mathbf{p}^s, w')}.$$

Next, given  $v(\mathbf{p}, \mathbf{p}^s, w) = u(\mathbf{c}^s(\mathbf{p}, \mathbf{p}^s, w))$  we have:

$$\begin{aligned} v_w &= u'(\mathbf{c}^s) \cdot \mathbf{c}_w = \lambda \mathbf{p}^s \cdot \mathbf{c}_w = \lambda \mathbf{p}^s \cdot \left( \frac{\mathbf{c}_{w'}^r(\mathbf{p}^s, w')}{\mathbf{p} \cdot \mathbf{c}_{w'}^r(\mathbf{p}^s, w')} \right) = \lambda \frac{\mathbf{p}^s \cdot \mathbf{c}_{w'}^r(\mathbf{p}^s, w')}{\mathbf{p} \cdot \mathbf{c}_{w'}^r(\mathbf{p}^s, w')} \\ v_w &= \frac{\lambda}{\mathbf{p} \cdot \mathbf{c}_{w'}^r(\mathbf{p}^s, w')}. \end{aligned}$$

□

We can check that things are consistent: with  $u_c = \lambda \mathbf{p}^s$ ,

$$v_{p_j^s} = u_c \mathbf{c}_{p_j^s} = \lambda \mathbf{p}^s (1 + \mathbf{c}_w p') \mathbf{S}_j^r = \lambda \mathbf{p}^s \mathbf{c}_w p' \mathbf{S}_j^r = v_w p' \mathbf{S}_j^r = v_w D_j^s.$$

**Proposition 12.5** (*Expenditure function – second derivatives*) Given  $e^s(\mathbf{p}, \mathbf{p}^s, u) = \mathbf{p} \cdot \mathbf{h}(\mathbf{p}^s, u)$ , we have

$$\begin{aligned} e_{p_i p_j}^s &= 0 \\ e_{p_i p_j^s}^s &= S_{ij}^r \\ e_{p_i^s p_j^s}^s &= -S_{ij}^r + (\mathbf{p} - \mathbf{p}^s) \cdot \mathbf{h}_{p_i^s p_j^s} = -S_{ij}^r + \sum_k (p_k - p_k^s) h_{p_i^s p_j^s}^k. \end{aligned}$$

The first derivatives of the expenditure functions were calculated in Proposition 12.1.

**Proof.** Given  $e^s(\mathbf{p}, \mathbf{p}^s, \bar{u}) = \mathbf{p} \cdot \mathbf{h}(\mathbf{p}^s, \bar{u})$ , we saw earlier in Proposition 12.1.

$$\begin{aligned} e_{p_j}^s &= h^j(\mathbf{p}^s, \bar{u}) \\ e_{p_j^s}^s &= (\mathbf{p} - \mathbf{p}^s) \cdot \mathbf{h}_{p_j^s}(\mathbf{p}^s, \bar{u}) = \mathbf{p} \cdot \mathbf{h}_{p_j^s}(\mathbf{p}^s, \bar{u}). \end{aligned}$$

Differentiating more,

$$e_{p_i p_j}^s = 0 \quad (157)$$

$$e_{p_i p_j}^s = h_{p_j^s}^i(\mathbf{p}^s, u) = S_{ij}^r, \quad (158)$$

and as  $e_{p_j^s}^s = (\mathbf{p} - \mathbf{p}^s) \cdot \mathbf{h}_{p_j^s}^r(\mathbf{p}^s, u)$ ,

$$e_{p_i^s p_j^s}^s = -h_{p_j^s}^i + \sum_k (p_k - p_k^s) h_{p_i^s p_j^s}^k.$$

□

### 12.3 Representation lemma for behavioral models

The following Lemma means that the demand function of a general abstract consumer can be represented as that of a misperceiving consumer with perceived prices  $\mathbf{p}^s(\mathbf{p}, w)$ .

**Lemma 12.1** (*Representing an abstract demand by a misperception*). *Given an abstract demand  $\mathbf{c}(\mathbf{p}, w)$ , and a utility function  $u(\mathbf{c})$ , we can define the function:*

$$\mathbf{p}^s(\mathbf{p}, w) = \frac{u_{\mathbf{c}}(\mathbf{c}(\mathbf{p}, w))}{v_w(\mathbf{p}, w)}. \quad (159)$$

*Then, the demand function can be represented as that of a sparse agent with perceived prices  $\mathbf{p}^s(\mathbf{p}, w)$ .*

$$\mathbf{c}(\mathbf{p}, w) = \mathbf{c}^s(\mathbf{p}, \mathbf{p}^s(\mathbf{p}, w), w). \quad (160)$$

**Proof** The demand of a sparse agent  $\mathbf{c}^s(\mathbf{p}, \mathbf{p}^s, w)$  is characterized by  $u_{\mathbf{c}}(\mathbf{c}^s(\mathbf{p}, \mathbf{p}^s, w)) = \lambda \mathbf{p}^s$  for some  $\lambda$ , and  $\mathbf{p} \cdot \mathbf{c} = w$ . By construction, we have  $u_{\mathbf{c}}(\mathbf{c}^s(\mathbf{p}, \mathbf{p}^s, w)) = \lambda \mathbf{p}^s$  for  $\mathbf{p}^s = \mathbf{p}^s(\mathbf{p}, w)$ . Hence, the representation is valid. We make a mild assumption, namely that given a  $\mathbf{c} = \mathbf{c}(\mathbf{p}, w)$ , there's no other  $\mathbf{c}'$  with  $\mathbf{p} \cdot \mathbf{c}' = w$ ,  $u_{\mathbf{c}}(\mathbf{c}') = u_{\mathbf{c}}(\mathbf{c})$ , and  $u(\mathbf{c}') > u(\mathbf{c})$ . Otherwise, we would need to consider another “branch” of the sparse max, namely a solution  $u_{\mathbf{c}}(\mathbf{c}^s) = \lambda \mathbf{p}^s$  with  $\mathbf{c}^s \cdot \mathbf{p} = w$  with  $\lambda$  not necessarily the lowest value possible. □

We note that for any  $\mathbf{p}^s(\mathbf{p}, w) = k u_{\mathbf{c}}(\mathbf{c}(\mathbf{p}, w))$  for some  $k > 0$ , we have  $\frac{u_{\mathbf{c}}(\mathbf{c}(\mathbf{p}, w))}{v_w(\mathbf{p}, w)} = \frac{\mathbf{p}^s}{\mathbf{p}^s \cdot \mathbf{c}_w}$  (indeed, both are equal to  $\frac{u_{\mathbf{c}}}{u_{\mathbf{c}} \cdot \mathbf{c}_w}$ ).

By contrast, the general model cannot in general be represented by a decision vs. experienced utility model. Indeed, a decision utility model always generates a symmetric Slutsky matrix  $\mathbf{S}^H(\mathbf{q}, w)$ , and this property does not hold in general for the general model. For example, the misperception model with exogenous perception  $M_{ij}(\mathbf{q}, w) = m_j 1_{\{i=j\}}$  features  $S_{ij}^H(\mathbf{q}, w) = S_{ij}^r(\mathbf{q}, w) m_j$ . Since  $S^r(\mathbf{q}, w)$  is symmetric,  $\mathbf{S}^H(\mathbf{q}, w)$  is not symmetric as long as there exists  $i$  and  $j$  with  $m_i \neq m_j$ .

## 13 Complements on basic consumer theory with nonlinear budget constraints

We give here complements to Section 7.

### 13.1 Rational agent

Primal is:  $v(\mathbf{p}, w) = \max_{\mathbf{c}} u(\mathbf{c})$  s.t.  $B(\mathbf{p}, \mathbf{c}) \leq w$  and demand  $\mathbf{c}(\mathbf{p}, w)$ . We can also define  $e(\mathbf{p}, \bar{u}) = \arg \min_{\mathbf{c}} B(\mathbf{p}, \mathbf{c})$  s.t.  $u(\mathbf{c}) \geq \bar{u}$  and Hicksian demand  $\mathbf{h}(\mathbf{p}, \bar{u})$ . We next derive the traditional consumer relation with that non-linear budget constraint.

*Shepard's lemma:* The envelope theorem gives

$$\begin{aligned} e_{p_i}(\mathbf{p}, \bar{u}) &= B_{p_i}(\mathbf{h}(\mathbf{p}, \bar{u}), \mathbf{p}) \\ e_{p_i p_j} &= B_{p_i p_j}(\mathbf{h}(\mathbf{p}, \bar{u}), \mathbf{p}) + B_{p_i c} \cdot \mathbf{h}_{p_j}(\mathbf{p}, \bar{u}). \end{aligned}$$

(the last term is to be read:  $B_{p_i c} \cdot \mathbf{h}_{p_j} = \sum_k B_{p_i c_k} \cdot \mathbf{h}_{p_j}^{c_k}$ ). We note that  $B_{p_i c} \cdot \mathbf{h}_{p_j}$  is symmetric.

*Roy's identity:*  $\bar{u} = v(\mathbf{p}, e(\mathbf{p}, \bar{u}))$ , so  $0 = v_{p_i} + v_w e_{p_i}$ , i.e.

$$v_{p_i} = -v_w B_{p_i}. \quad (161)$$

Given  $\mathbf{c}(\mathbf{p}, w) = \mathbf{h}(\mathbf{p}, v(\mathbf{p}, w))$ , we have  $\mathbf{c}_w = \mathbf{h}_u v_w$ ,  $\mathbf{c}_{p_i} = \mathbf{h}_{p_i} + \mathbf{h}_u v_{p_i} = \mathbf{h}_{p_i} + \mathbf{c}_w \frac{v_{p_i}}{v_w}$  and because of Roy,  $\mathbf{c}_{p_i} = \mathbf{h}_{p_i} - \mathbf{c}_w B_{p_i}$ , i.e. the Slutsky relation for nonlinear budget constraints:

$$\mathbf{h}_{p_i} = \mathbf{c}_{p_i} + \mathbf{c}_w B_{p_i}. \quad (162)$$

Finally, given  $B(\mathbf{c}(\mathbf{p}, w), \mathbf{p}) = w$ ,

$$B_c \mathbf{c}_{p_i} = -B_{p_i}, B_c \mathbf{c}_w = 1. \quad (163)$$

Premultiplying (162) by  $B_c$  gives:  $B_c \mathbf{h}_{p_i} = B_c \mathbf{c}_{p_i} + B_c \mathbf{c}_w B_{p_i} = -B_{p_i} + B_{p_i} = 0$ ,

$$B_c \mathbf{h}_{p_i} = 0. \quad (164)$$

All in all, traditional consumer theory holds, replacing  $p$  by  $B_c$ .

### 13.2 Misperceiving Agent

We now study

$$v(\mathbf{p}, \mathbf{p}^s, w) = \operatorname{smax}_{\mathbf{c} | \mathbf{p}^s} u(\mathbf{c}) \text{ s.t. } B(\mathbf{c}, \mathbf{p}) \leq w. \quad (165)$$

We call  $\mathbf{c}^s(\mathbf{p}, \mathbf{p}^s, w)$  the demand function. Recall that's it's characterized by  $u_{\mathbf{c}}(\mathbf{c}) = \lambda B_{\mathbf{c}}(\mathbf{c}, \mathbf{p}^s)$  for some  $\lambda$ , and  $\mathbf{c}(\mathbf{p}, \mathbf{p}^s, w) = \mathbf{c}^r(\mathbf{p}^s, w')$  for a  $w'$  that ensures

$$B(\mathbf{c}^s(\mathbf{p}, \mathbf{p}^s, w), \mathbf{p}) = w.$$

Given  $B(\mathbf{c}(\mathbf{p}, \mathbf{p}^s, w), \mathbf{p}) = w$ , we have (taking the derivatives w.r.t.  $\mathbf{p}, \mathbf{p}^s, w$ ):

$$\begin{aligned} B_{\mathbf{c}}\mathbf{c}_{\mathbf{p}} &= -B_{\mathbf{p}} \\ B_{\mathbf{c}}\mathbf{c}_{\mathbf{p}^s} &= 0 \\ B_{\mathbf{c}}\mathbf{c}_w &= 1. \end{aligned}$$

where  $B_{\mathbf{c}} = B_{\mathbf{c}}(\mathbf{c}, \mathbf{p}^s)$ ,  $B_{\mathbf{p}} = B_{\mathbf{p}}(\mathbf{c}, \mathbf{p}^s)$ .

Likewise,  $B(\mathbf{c}^r(\mathbf{p}^s, w'), \mathbf{p}^s) = w'$  gives

$$\begin{aligned} B_{\mathbf{c}}^s\mathbf{c}_{\mathbf{p}^s}^r &= -B_{\mathbf{p}^s}^s \\ B_{\mathbf{c}}\mathbf{c}_{w'}^r &= 1. \end{aligned}$$

Define the rational Hicksian action

$$\mathbf{h}^r(\mathbf{p}^s, \bar{u}) = \arg \min_{\mathbf{c}} B(\mathbf{c}, \mathbf{p}^s) \text{ s.t. } u(\mathbf{c}) \geq \bar{u}, \quad (166)$$

and the corresponding Slutsky matrix, and the perceived  $\mathbf{p}^s$

$$\mathbf{S}_j^r = \mathbf{h}_{\mathbf{p}_j^s}^r(\mathbf{p}^s, \bar{u})|_{\bar{u}=v(\mathbf{p}, \mathbf{p}^s, w)}. \quad (167)$$

Define the dual expenditure function

$$e(\mathbf{p}, \mathbf{p}^s, \bar{u}) = \text{smin}_{\mathbf{c}|\mathbf{p}^s} B(\mathbf{c}, \mathbf{p}) \text{ s.t. } u(\mathbf{c}) \geq \bar{u}. \quad (168)$$

We have the simpler representation:

$$e(\mathbf{p}, \mathbf{p}^s, \bar{u}) = B(\mathbf{h}^r(\mathbf{p}^s, \bar{u}), \mathbf{p}). \quad (169)$$

**Proposition 13.1** (Shepard's lemma, nonlinear and behavioral). *Given the expenditure function  $e(\mathbf{p}, \mathbf{p}^s, \bar{u}) = B(\mathbf{h}^r(\mathbf{p}^s, \bar{u}), \mathbf{p})$ , we have*

$$\begin{aligned} e_{\mathbf{p}_j} &= B_{\mathbf{p}_j} \\ e_{\mathbf{p}_j^s} &= (B_{\mathbf{c}} - B_{\mathbf{c}}^s) \cdot \mathbf{S}_j^r. \end{aligned}$$

**Proof.**  $e(\mathbf{p}, \mathbf{p}^s, \bar{u}) = B(\mathbf{h}^r(\mathbf{p}^s, \bar{u}), \mathbf{p})$  gives:  $e_{\mathbf{p}_j} = B_{\mathbf{p}_j}$ , and

$$e_{\mathbf{p}_j^s} = B_c \mathbf{h}_{\mathbf{p}_j^s}^r = (B_c - B_c^s) \mathbf{h}_{\mathbf{p}_j^s}^r.$$

as (164) gives  $B_c^s \mathbf{h}_{\mathbf{p}_j^s}^r = 0$ .  $\square$

**Proposition 13.2** (Generalized Roy's identity). *Given the function  $v(\mathbf{p}, \mathbf{p}^s, w)$ , we have:*

$$\begin{aligned} \frac{v_{\mathbf{p}_j}}{v_w} &= -B_{\mathbf{p}_j} \\ \frac{v_{\mathbf{p}_j^s}}{v_w} &= (B_c^s - B_c) \cdot \mathbf{S}_j^r = D_j^s. \end{aligned}$$

$$i.e. \frac{v_{\mathbf{p}_j^s}}{v_w^s} = \sum_i (B_i^s - B_i) \mathbf{S}_{ij}^r.$$

**Proof of Proposition 13.2**

For a number  $\bar{u}$ , we have the identity  $\bar{u} = v(\mathbf{p}, \mathbf{p}^s, e(\mathbf{p}, \mathbf{p}^s, \bar{u}))$  for all  $\mathbf{p}, \mathbf{p}^s, \bar{u}$ . Deriving w.r.t.  $\mathbf{p}_j$  gives:

$$0 = v_{\mathbf{p}_j} + v_w e_{\mathbf{p}_j} = v_{\mathbf{p}_j} + v_w B_{\mathbf{p}_j}$$

by the behavioral Shepard's lemma (Proposition 13.1).

Deriving w.r.t.  $\mathbf{p}_j^s$  gives:

$$0 = v_{\mathbf{p}_j^s} + v_w e_{\mathbf{p}_j^s} = v_{\mathbf{p}_j^s} + v_w (B_c - B_c^s) \cdot \mathbf{S}_j^r$$

again by the behavioral Shepard's lemma (Proposition 13.1).  $\square$

**Proposition 13.3** (Marshallian demand) *Given the Marshallian action  $\mathbf{c}(\mathbf{p}, \mathbf{p}^s, w)$ , we have:*

$$\mathbf{c}_{\mathbf{p}_j} = -\mathbf{c}_w B_{\mathbf{p}_j} \tag{170}$$

$$\mathbf{c}_{\mathbf{p}_j^s} = \mathbf{S}_j^r + \mathbf{c}_w D_j^s, \tag{171}$$

$$i.e. \mathbf{c}_{\mathbf{p}_j}^i = -B_{\mathbf{p}_j}^i \cdot \mathbf{c}_w \text{ and } \mathbf{c}_{\mathbf{p}_j^s}^i = \mathbf{S}_j^{i,r} + \mathbf{c}_w^i D_j^s.$$

**Proof.** We have  $\mathbf{c}(\mathbf{p}, \mathbf{p}^s, w) = \mathbf{h}(\mathbf{p}^s, v(\mathbf{p}, \mathbf{p}^s, w))$ , which implies

$$\mathbf{c}_w = \mathbf{h}_u v_w, \quad \mathbf{c}_{\mathbf{p}_j} = \mathbf{h}_u v_{\mathbf{p}_j}. \tag{172}$$

Because of Roy ( $v_{\mathbf{p}_j} = -B_{\mathbf{p}_j} v_w$ ), we have:  $\mathbf{c}_w B_{\mathbf{p}_j} = \mathbf{h}_u v_w B_{\mathbf{p}_j} = -\mathbf{h}_u v_{\mathbf{p}_j} = -\mathbf{c}_{\mathbf{p}_j}$ , hence  $\mathbf{c}_{\mathbf{p}_j} = -\mathbf{c}_w B_{\mathbf{p}_j}$ .



Also,  $\mathbf{c}_{p_j^s}(\mathbf{p}, \mathbf{p}^s, w) = \mathbf{h}(\mathbf{p}^s, v(\mathbf{p}, \mathbf{p}^s, w))$  gives:

$$\begin{aligned}\mathbf{c}(\mathbf{p}, \mathbf{p}^s, w) &= \mathbf{h}_{p_j^s} + \mathbf{h}_u v_{p_j^s} \\ &= \mathbf{S}_j^r + \mathbf{c}_w \frac{v_{p_j^s}}{v_w} \text{ using } \mathbf{c}_w = \mathbf{h}_u v_w \\ &= \mathbf{S}_j^r + \mathbf{c}_w [(B_c^s - B_c) \cdot \mathbf{S}_j^r] \text{ using Proposition 13.2.}\end{aligned}$$

□

In the traditional model  $v(\mathbf{p}, w) = \max_{\mathbf{c}} u(\mathbf{c}) + \lambda(w - B(\mathbf{c}, \mathbf{p}))$  implies  $v_w = \lambda$ . There is a deviation here, as indicated below.

**Proposition 13.4** (*Envelope theorem, modified*) Call  $\lambda$  the Lagrange multiplier such that  $u_{\mathbf{c}}(\mathbf{c}) = \lambda B_{\mathbf{c}}(\mathbf{c}, \mathbf{p}^s)$ . We have:

$$\frac{v_w}{\lambda} = B_{\mathbf{c}}^s \mathbf{c}_w = 1 + (B_{\mathbf{c}}^s - B_{\mathbf{c}}) \mathbf{c}_w, \quad (173)$$

where  $B_{\mathbf{c}}^s = B_{\mathbf{c}}(\mathbf{c}, \mathbf{p}^s)$  and  $B_{\mathbf{c}} = B_{\mathbf{c}}(\mathbf{c}, \mathbf{p})$ .

**Proof.** We have:  $v(\mathbf{p}, \mathbf{p}^s, w) = u(\mathbf{c}(\mathbf{p}, \mathbf{p}^s, w))$ , so

$$\begin{aligned}\frac{v_w}{\lambda} &= \frac{1}{\lambda} u_{\mathbf{c}} \mathbf{c}_w = B_{\mathbf{c}}^s \mathbf{c}_w \\ &= B_{\mathbf{c}}^s \mathbf{c}_w + 1 - B_{\mathbf{c}} \mathbf{c}_w \text{ using } B_{\mathbf{c}} \mathbf{c}_w = 1 \text{ from } B(\mathbf{c}, \mathbf{p}) = w \\ &= 1 + (B_{\mathbf{c}}^s - B_{\mathbf{c}}) \mathbf{c}_w.\end{aligned}$$

□

### 13.3 Hybrid Model: Agent maximizing the wrong utility function with the wrong prices

Suppose now an agent with true problem  $\max_{\mathbf{c}} u(\mathbf{c})$  s.t.  $B(\mathbf{p}, \mathbf{c}) \leq w$  but maximizes instead  $\max_{\mathbf{c}|\mathbf{p}^s} u^s(\mathbf{c})$  s.t.  $B(\mathbf{p}, \mathbf{c}) \leq w$  with both the wrong utility and the wrong prices. This is hybrid of the two previous models.

In terms of decision (if not welfare), the agent is a misperceiving agent with utility  $u^s$  and perceived prices  $\mathbf{p}^s$ . Call  $v^s(\mathbf{p}, w) = u^s(\mathbf{c}(\mathbf{p}, w))$  and  $\mathbf{h}^{r,s}(\mathbf{p}^s, \hat{w}) = \arg \min_{\mathbf{c}} B(\mathbf{p}^s, \mathbf{c})$  s.t.  $u^s(\mathbf{c}) \geq \hat{w}$  the indirect utility function (of that misperceiving agent) and the rational compensated demand of that agent with utility  $u^s$ . Then, our agent has demand:

$$\mathbf{c}(\mathbf{p}, w) = \mathbf{h}^{r,s}(\mathbf{p}^s(\mathbf{p}, w), v^s(\mathbf{p}, w)). \quad (174)$$

**Proposition 13.5** (Agent misperceiving both utility and prices) Take the model of an agent maximizing the wrong utility function  $u^s(\mathbf{c})$ , with the wrong perceived prices  $\mathbf{p}^s$ . Call  $\mathbf{S}^{r,s}(\mathbf{p}, w) =$

$\mathbf{h}_{\mathbf{p}^s}^{r,s}(\mathbf{p}^s(\mathbf{p}, w), v^s(\mathbf{p}, w))$  the Slutsky matrix of the underlying rational agent who has utility  $u^s$ , and define

$$\mathbf{S}_j^s(\mathbf{p}, w) = \mathbf{S}^{r,s}(\mathbf{p}, w) \cdot \mathbf{p}_{p_j}^s(\mathbf{p}, w). \quad (175)$$

i.e.  $S_{ij}^s = \sum_k S_{ik}^{r,s} \frac{\partial p_k^s(\mathbf{p}, w)}{\partial p_j}$ , where  $\frac{\partial p_k^s(\mathbf{p}, w)}{\partial p_j}$  is the matrix of marginal perception. Then,

$$\begin{aligned} \mathbf{S}_j^C(\mathbf{p}, w) &= \mathbf{S}_j^s(\mathbf{p}, w) + \mathbf{c}_w \left( \frac{v_{p_j}^s}{v_w^s} + B_{p_j} \right) \\ \mathbf{S}_j^H(\mathbf{p}, w) &= \mathbf{S}_j^s + \mathbf{c}_w \left( \frac{v_{p_j}^s}{v_w^s} - \frac{v_{p_j}}{v_w} \right) \\ \mathbf{S}_j^s(\mathbf{p}, w) &= \mathbf{c}_{p_j} - \mathbf{c}_w \frac{v_{p_j}^s}{v_w^s}. \end{aligned}$$

We can write

$$-D_j = \bar{\tau}^b \cdot \mathbf{S}_j^s,$$

with:

$$\bar{\tau}^b = (B_c(\mathbf{p}, \mathbf{c}) - B_c(\mathbf{p}^s, \mathbf{c})) + \left( \frac{u_c^s}{v_w^s} - \frac{u_c}{v_w} \right). \quad (176)$$

Finally,  $B_c(\mathbf{p}^s, \mathbf{c}) \cdot \mathbf{S}_j^s = 0$ .

This tax  $\bar{\tau}^b$  is the sum of two gaps: between the prices and perceived prices ( $B_c(\mathbf{c}, \mathbf{p}) - B_c(\mathbf{p}^s, \mathbf{c})$ ), and between true utility and perceived utility ( $\frac{u_c^s}{v_w^s} - \frac{u_c}{v_w}$ ).

**Proof.** So, with  $\mathbf{M}_j = \frac{\partial \mathbf{p}^s(\mathbf{p}, w)}{\partial p_j}$ , and use  $\mathbf{c}(\mathbf{p}, w) = \mathbf{h}^{r,s}(\mathbf{p}^s, v^s(\mathbf{p}, w))$ :

$$\begin{aligned} \mathbf{c}_w(\mathbf{p}, w) &= \mathbf{h}_u^{r,s} v_w^s \\ \mathbf{c}_{p_j}(\mathbf{p}, w) &= \mathbf{h}_{p_j^s}^{r,s} \cdot \mathbf{M}_j + \mathbf{h}_u^{r,s} v_j^s = \mathbf{S}_j^s + \mathbf{c}_w \frac{v_j^s}{v_w^s}. \end{aligned}$$

Using (31) and (32) gives:

$$\begin{aligned} \mathbf{S}_j^C &= \mathbf{c}_{p_j}(\mathbf{p}, w) + \mathbf{c}_w(\mathbf{p}, w) B_{p_j}(\mathbf{c}, \mathbf{p}) = \mathbf{S}_j^s + \mathbf{c}_w \left( \frac{v_j^s}{v_w^s} + B_{p_j} \right) \\ \mathbf{S}_j^H &= \mathbf{c}_{p_j}(\mathbf{p}, w) - \mathbf{c}_w(\mathbf{p}, w) \frac{v_{p_j}(\mathbf{p}, w)}{v_w(\mathbf{p}, w)} = \mathbf{S}_j^s + \mathbf{c}_w \left( \frac{v_j^s}{v_w^s} - \frac{v_{p_j}(\mathbf{p}, w)}{v_w(\mathbf{p}, w)} \right). \end{aligned}$$

We have

$$B_c(\mathbf{p}^s, \mathbf{c}) \cdot \mathbf{S}_j^s = B_c(\mathbf{p}^s, \mathbf{c}) \cdot \mathbf{h}^r(\mathbf{p}^s, v^s) \cdot \mathbf{M}_j = 0 \text{ as } B_c(\mathbf{p}^s, \mathbf{c}) \cdot \mathbf{h}^r(\mathbf{p}^s, v^s) = 0.$$

Recall also that  $\frac{u_c^s}{v_w^s} = \Lambda B_c(\mathbf{p}^s, \mathbf{c})$  as the agent maximizes with perceived prices  $\mathbf{p}^s$ . Hence,

$$\begin{aligned}
 -D_j &= \left( B_c(\mathbf{c}, \mathbf{p}) - \frac{u_c}{v_w} \right) \mathbf{S}_j^s \\
 &= \left( (B_c(\mathbf{c}, \mathbf{p}) - B_c(\mathbf{p}^s, \mathbf{c})) - \left( \frac{u_c}{v_w} - \Lambda B_c(\mathbf{p}^s, \mathbf{c}) \right) \right) \mathbf{S}_j^s \text{ as } B_c(\mathbf{p}^s, \mathbf{c}) \cdot \mathbf{S}_j^s = 0 \\
 &= \left( (B_c(\mathbf{c}, \mathbf{p}) - B_c(\mathbf{p}^s, \mathbf{c})) - \left( \frac{u_c}{v_w} - \frac{u_c^s}{v_w^s} \right) \right) \cdot \mathbf{S}_j^s = \bar{\tau}_b \cdot \mathbf{S}_j^H.
 \end{aligned}$$

□