

# Online Appendix for “Optimal Taxation with Behavioral Agents”

Emmanuel Farhi and Xavier Gabaix

May 2018

Section 9 contains additional results on the paper, in particular on setups with heterogeneous agents and endogenous attention. 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 (13), with

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

This yields:

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

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  $\frac{\tau_i^*}{p_i} = 0.0729$ , or 7.29%. Under homogeneity with the same average misperception  $m_i^h = .25$  for all agents,  $\frac{\tau_i^*}{p_i} = 20.3\%$ , for a ratio  $.203/.0729 = 2.78$ . When the taxes are fully salient, so  $m_i^h = 1$  for all agents, then the optimal tax is 1.27%, giving a ratio  $.0127/.0729 = .174$ .

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

**Nonzero default tax** It is sometimes important to introduce a distinction between the misperception of marginal tax changes and the misperception of the average level of taxes. To capture this possibility, we assume that perceived taxes are given by  $\tau_i^s = m_i \tau_i + (1 - m_i) \tau_i^d$ , where  $\tau_i^d$  is a default tax (which could be for instance the average tax in the economy, or an average of past tax rates). This change introduces a new additive term  $-\frac{\tau_i^d (1-m_i)(1-\Lambda)}{p_i m_i + (1-m_i)\Lambda}$  in the optimal tax formula (17) to correct for this new form of misperception, in the case of isoelastic utility.

To take a concrete example, suppose that we start from an equilibrium where taxes are optimal and default taxes are equal to true taxes. Imagine that there is a reduction in the need for public funds  $\Lambda$ , but that default taxes  $\tau_i^d$  remain high at the pre-change level. Then lowering taxes induces agents to over-perceive the average level taxes, and creates a force for the government to lower taxes even further to correct this new bias.

**Heterogeneous attention and default taxes** 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:

$$\tau_i = \frac{\Lambda' \mathbb{E}[a^h - \Psi^h p_i] - \mathbb{E}[\Psi^h (1 - m_i^h) (m_i^h + \Lambda') \tau_i^{d,h}]}{\mathbb{E}[\Psi^h m_i^h (m_i^h + 2\Lambda')]} \quad (47)$$

### 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 (11).

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

$$\chi = \frac{\mathbb{E}[\xi^h \eta^h]}{\mathbb{E}[\eta^{h2}]} = \frac{\mathbb{E}[\xi^h] \mathbb{E}[\eta^h] + cov(\xi^h, \eta^h)}{\mathbb{E}[\eta^h]^2 + var[\eta^h]}, \quad (48)$$

where again  $\mathbb{E}$  denotes the average over agents  $h$ .<sup>74</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  $cov(\xi^h, \eta^h)$ ). The optimal nudge is weaker when

<sup>74</sup>The intermediate steps are as follows. Using  $c_\chi^h = -\Psi \eta^h$ ,  $\tau = 0$ ,  $\tau^{b,h} = \tau^{X,h} - \chi \eta^h$ , we get  $\frac{\partial L}{\partial \chi}(\tau, \chi) = \sum_h [\lambda \tau - \lambda \tau^{\xi,h} - \beta^h \tau^{b,h}] \cdot c_\chi^h = \lambda \sum_h [0 - \tau^{X,h} + \chi \eta^h] \Psi \eta^h$ .

there is more heterogeneity in nudgeability (higher  $\text{var} [\eta^h]$ ).<sup>75</sup>

#### 9.1.4 Nudges vs. Taxes with Redistributive Concerns

**Jointly optimal nudges and taxes** We normalize  $p = 1$ . Agent  $h$  has utility  $u^h(\mathbf{c}) = c_0^h + \frac{a^h c^h - \frac{1}{2}(c^h)^2}{\Psi}$ , so that  $\beta^h = \gamma^h$  and  $c^h(\tau, \tau^{X,h}) = a^h - \Psi(m^h \tau + \tau^{X,h})$ . We investigate the optimal joint policy using both nudges and taxes on goods  $c$ . One can show that

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

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  $\eta^h$  are positively correlated. Loosely speaking, if poor agents (with a high  $\gamma^h$ ) are highly nudgeable, 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  $\gamma^h = \lambda$ .

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

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

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>76</sup> Similarly, the more heterogeneity there is in nudgeability (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  $\tau$ .

<sup>75</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.

<sup>76</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.

In the general case, we assume heterogenous welfare weights with  $\mathbb{E} [\gamma^h] = \lambda$ .

**Proposition 9.2** *The optimal tax and nudge satisfy*

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

**Proof.** 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. Individual  $h$  creates an externality plus internality. We have successively,

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

Proposition 2.2 gives optimal tax, using  $S^{H,h} = -\Psi m^h$ :

$$\begin{aligned}\frac{\partial L}{\partial \tau} &= \mathbb{E} \sum_h (\lambda - \gamma_h) c^h - \lambda (\tau - \tau^{\xi,h} - \tilde{\tau}^{\xi,h}) \Psi m^h \\ &= \mathbb{E} \sum_h (\lambda - \gamma_h) c^h - [(\lambda - \gamma^h (1 - m^h)) \tau - \lambda \tau^{X,h} + \gamma^h \tau^{X,h}] \Psi m^h.\end{aligned}\tag{49}$$

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

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

Proposition 2.3 gives optimal nudge:

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

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 formulas.

**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.

We augment the example of Section 3.4 with aversive nudges. We use quadratic utility functions. 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 9.3** (“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*

$$\frac{\lambda - \gamma^h}{m^h} > \frac{-\iota^h \gamma^h}{\eta^h}. \quad (50)$$

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).

**Proof.** We apply our tax formulas (11):

$$\begin{aligned} \frac{\partial L}{\partial \chi} &= - \sum_h \iota^h \gamma^h c^h - \sum_h [(\lambda - \gamma^h (1 - m^h)) \tau - \lambda \tau^{X,h} + \gamma^h \chi \eta^h] \Psi \tau_x^{X,h} \\ &= - \sum_h \iota^h \gamma^h c^h - \Psi \sum_h [(\lambda - \gamma^h (1 - m^h)) \tau - \lambda \tau^{X,h} + \gamma^h \chi \eta^h] \eta^h \end{aligned} \quad (51)$$

$$= \sum_h [-\iota^h \gamma^h c^h - \Psi x^h \eta^h], \quad (52)$$

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  $x^h < 0$ .

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

so

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

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 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 result is a bit knife-edge.

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

### 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  $\boldsymbol{\tau} \cdot \mathbf{S}_i^{r, h} = \sum_j \tau_j S_{ji}^{r, h}$  as  $\boldsymbol{\tau} \cdot \mathbf{S}_i^{r, h} = \sum_j \tau_j S_{ij}^{r, h}$ . We can then rewrite the optimal tax formula of Proposition 2.1 in ‘‘discouragement’’ form as

$$\frac{-\sum_{h, j} \tau_j S_{ij}^{r, h}}{c_i} = 1 - \frac{\bar{\gamma}}{\lambda} - \text{cov}\left(\frac{\gamma^h}{\lambda}, \frac{Hc_i^h}{c_i}\right), \quad (53)$$

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 decisions 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>77</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>78</sup> For instance we might consider a

<sup>77</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>78</sup>The first order condition characterizing the optimal allocation of attention can be written as  $\tau^b \cdot \mathbf{c}_{m_j}(\mathbf{q}, w, \mathbf{m}) = 0$  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  $k \in I^{\mathbf{m}}$  corresponding to dimension  $i \in \{1, \dots, A\}$  of attention. We then get  $\sum_{i \in I^{\mathbf{C}}} \tau_i^b \mathbf{C}_{m_j}(\mathbf{q}, w, \mathbf{m}) + \tau_{k(j)}^b = 0$  for all  $j \in \{1, \dots, A\}$ .

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 (8) has a term  $(\boldsymbol{\tau} - \tilde{\boldsymbol{\tau}}^{b,h}) \cdot \mathbf{S}_i^{C,h} = \sum_{k \in I^m \cup I^c} (\tau_k - \tilde{\tau}_k^{b,h}) S_{ki}^{C,h}$ , 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  $(\boldsymbol{\tau} - \tilde{\boldsymbol{\tau}}^{b,h}) \cdot \mathbf{S}_i^{C,h} = \sum_{k \in I^c} (\tau_k - \tilde{\tau}_k^{b,h}) S_{ki}^{C,h}$  is the sum restricted to commodities. Second, consider the “optimally allocated attention” case. Then (see Proposition 9.5)  $(\boldsymbol{\tau} - \tilde{\boldsymbol{\tau}}^{b,h}) \cdot \mathbf{S}_i^{C,h} = \sum_{k \in I^c} (\tau_k S_{ki}^{C,h} - \tilde{\tau}_k^{b,h} S_{ki|m}^{C,h})$ , where  $\mathbf{S}_{i|m}^{C,h}$  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  $\tau_k S_{ki}^{C,h}$ ). However, for welfare, when attention is assumed to be optimally allocated, it is the Slutsky matrix holding attention constant that matters (the term  $\tilde{\tau}_k^{b,h} S_{ki|m}^{C,h}$ ). 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>79</sup>

**Proposition 9.4** (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:*

$$\boldsymbol{\tau}^b \cdot \mathbf{c}_{\theta}(\mathbf{p}, w, \theta) = 0. \quad (54)$$

**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 9.5** (Value of  $D_j$  when attention is optimal). *When attention is of the form  $\mathbf{c}(\mathbf{p}, w, \theta) =$*

<sup>79</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.

$(\mathbf{C}(\mathbf{p}, w, \mathbf{m}(\theta)), \mathbf{m}(\theta))$ , and is optimally chosen, then

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

where  $\tau_{\mathbf{C}}^b = B_{\mathbf{C}}(\mathbf{C}, \mathbf{p}) - \frac{u_{\mathbf{C}}(\mathbf{C}, \mathbf{m})}{v_w(\mathbf{p}, w)}$  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

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

□

**“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,  $\mathbf{p}_{\mathbf{m}} = 0$ . 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 9.6** (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*

$$-D_j = (\tau_{\mathbf{C}}^b, 0) \cdot \mathbf{S}_j^H(\mathbf{p}, w) = \sum_{i=1}^n \tau_{\mathbf{C}_i}^b S_{ij}^H = (\tau_{\mathbf{C}}^b, 0) \cdot \mathbf{S}_j^C(\mathbf{p}, w) = (\tau_{\mathbf{C}}^b, 0) \cdot \mathbf{c}_j(\mathbf{p}, w).$$

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

**Proof**

We have  $\tau^b = (\tau_{\mathbf{c}}^b, \tau_{\mathbf{m}}^b) = (\tau_{\mathbf{c}}^b, 0)$  as  $u_{\mathbf{m}} = 0$ . So,  $-D_j = \tau^b \cdot \mathbf{S}_j^H(\mathbf{p}, w) = \tau_{\mathbf{C}}^b \cdot \mathbf{S}_{\mathbf{C}, j}^H$ . □

**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 9.5 with  $\mathbf{m}(\theta) = \theta$ . This gives:  $-D_j = \tau_C^b \cdot \mathbf{S}_{C_j}^{H,m}$  with

$$\mathbf{S}_{C_j|\mathbf{m}}^H = \mathbf{S}^r \mathbf{p}_{p_j}^s(\mathbf{p}, w, \mathbf{m}), \quad (55)$$

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

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

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

When attention is not necessarily optimal, we also have (from (44)), using again decomposition  $\tau^b = (\tau_c^b, \tau_m^b)$ :

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

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

$$\tau^{m,*} = \frac{\Lambda/\psi}{m(\tau)^2 + \tau m'(\tau)}. \quad (56)$$

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

$$\tau^{m,0} = \frac{\Lambda/\psi}{m(\tau)^2 + \tau m(\tau) m'(\tau)}. \quad (57)$$

When attention is fixed, the optimal tax is

$$\tau^{m,F} = \frac{\Lambda/\psi}{m(\tau)^2}. \quad (58)$$

**Proposition 9.7** *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>80</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 9.2.3 for completely worked out linear-quadratic and isoelastic examples.

### 9.2.3 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  $q_1^s(p_1, m) = p_1^d + m\tau_1$ , 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  $U'(c_1) = q^s$ , so  $c_1^r(p^s) = a - \Psi q^s$ . The perceived tax is:

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

and demand is

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

---

<sup>80</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 losses from inattention are  $\frac{1}{2}u_{c_1c_1}(c_1^r - c_1)^2 = -\frac{1}{2}\Psi\tau^2(1-m)^2$ . (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:

$$g'(m) = \Psi(1-m)\tau_1^2. \quad (59)$$

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  $U(c) = \frac{c^{1-1/\psi} - 1}{1-1/\psi}$  and attention technology  $q_1^s(p_1, m) = p_1^d + m\tau_1$ , 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:

$$S_{11}^H = \frac{dc_1(p_1^d + m\tau_1, m(\tau_1))}{d\tau_1} = -\Psi(m + \tau m'(\tau_1))$$

$$S_{21}^H = \frac{\partial m}{\partial \tau_1} = m'(\tau_1).$$

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

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

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 (17):

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

i.e.  $\frac{\tau_1}{q_1^s} = \frac{\Lambda}{\psi_1(m + \Lambda(1 - m))m}$ , which is equivalent to (17).

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

$$\tau^0 = \frac{\Lambda/\psi}{m^2 + \tau m m'(\tau) + \Lambda\left((1 - m)(m + \tau m') - \frac{m}{\psi}\right)}. \quad (60)$$



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

$$\tau^1 = \frac{\Lambda/\psi}{m^2 + \tau m'(\tau) + \Lambda \left( (1 - m) m - \frac{m}{\psi} \right)}. \quad (61)$$

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

$$\tau^0 = \frac{\Lambda/\psi}{m^2 + \tau_1 m'(\tau_1) + \Lambda \left( (1 - m) m - \frac{m}{\psi} \right) - \tau (1 - m) m' + \Lambda (1 - m) \tau m'} \quad (62)$$

$$= \frac{\Lambda/\psi}{m^2 + \tau_1 m'(\tau_1) + \Lambda \left( (1 - m) m - \frac{m}{\psi} \right) - (1 - \Lambda) \tau (1 - m) m'}. \quad (63)$$

**Proposition 9.8** *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  $m(\tau)$ ; its 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 9.1** *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

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

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

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 9.9** *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

$$m(\tau) = 1_{\tau > \tau_*}, \quad \tau_* = \frac{\kappa}{\sqrt{\psi}}. \tag{64}$$

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

**Proposition 9.10** *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

$\Lambda > \sqrt{2} + 1. \square$

### 9.3 Other extensions

#### 9.3.1 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  $\tau^X = 0$ , and in the limit of small taxes. Defining  $\Lambda = \frac{\lambda}{\gamma} - 1$ , we can rewrite formula (13), 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 9.11** (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>81</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>82</sup>

#### 9.3.2 Tax instruments with differential saliences

We elaborate on a remark we made at the end of section 3.5. 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

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

<sup>81</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>82</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.

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  $-\tau \frac{m}{m'}$  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.

### 9.3.3 A different budget adjustment rule

The specific formulation of misperception that we have used in this section assumes that the budget adjustments required when agents misperceive taxes are all absorbed by the consumption a good (good 0) with a constant marginal utility. This renders these adjustments relatively painless.

We now explore a variant which increases their costs. We assume that the budget adjustments are concentrated on a “shock absorber” good with a sharply decreasing marginal utility. This increases the distortionary costs of non-salient taxes and reduces optimal taxes in a way that we characterize precisely below. Of course, it is difficult to know a priori which good is the “shock absorber” good (or set of goods) – this is one more place where more empirical evidence is needed to address a behavioral enrichment of the traditional model. It could be some luxury goods (e.g. some restaurant meals), or perhaps more pessimistically investments that can be postponed, e.g. health investments. Our purpose here is only to show how this possibility matters for the results.

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

$$c^{i,s}(\mathbf{q}, \mathbf{q}^s, w) = c^{i,r}(\mathbf{q}^s, w) \text{ for } i \neq n \quad (65)$$

$$c^{n,s}(\mathbf{q}, \mathbf{q}^s, w) = \frac{1}{q_n} \left( w - \sum_{i \neq n} q_i c^{i,s}(\mathbf{q}, \mathbf{q}^s, w) \right). \quad (66)$$

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>83</sup> We shall see in the next proposition that we can also write

$$c^{n,s}(\mathbf{q}, \mathbf{q}^s, w) = c^{n,r}(\mathbf{q}^s, w) - \frac{1}{q_n} (\mathbf{q} - \mathbf{q}^s) \cdot \mathbf{c}^r(\mathbf{q}^s, w), \quad (67)$$

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

---

<sup>83</sup>[Chetty et al. \(2009\)](#) consider such a rule in a 2-good context. [Gabaix \(2017b\)](#) considers such a rule when doing dynamic programming, and the last good is “next period wealth”.

good  $n$ .

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

**Proposition 9.12** (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:*

$$\partial_w c_i^s(\mathbf{q}, \mathbf{q}^s, w) = \partial_w c_i^r(\mathbf{q}, w). \quad (68)$$

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

$$S_{ij}^s = S_{ij}^r m_j + \left( \partial_w c_i^r - \frac{1}{q_n} 1_{i=n} \right) (1 - m_j) c^j, \quad (69)$$

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

**Proof.** The term  $\partial_w c_i^s$  is trivial, as it’s evaluate at  $\mathbf{q}^s = \mathbf{q}$ . 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  $S_{ij}^r + c_w^i c^j (1 - m_j)$  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.  $\square$

**A simple particular case** We next present a simple particular case. Utility is separable,  $u(\mathbf{c}) = \sum_{i=0}^n u_i(c_i)$  with  $u'_0(c_0) = 1$ ,  $u'_i(c_i) = c_i^{-1/\psi_i}$  for  $i = 1, \dots, n-1$  and the “shock absorber” good  $n$  has constant marginal utility of  $u'_n(c_n) = 1 - \nu < 1$  if  $c_n \geq 1$  and  $1 + \mu > 1$  if  $c_n < 1$ .<sup>84</sup> We call  $\mu > 0$  the marginal distortionary cost of budget adjustment. Goods 0 and  $n$  cost \$1, and they are untaxed.

The agent chooses his consumption of goods  $c_0, \dots, c_{n-1}$  based on the perceived prices  $q_i^s = 1 + m_i \tau_i$  and the rest of his money is spent on the last good. Specifically, the demands are as follows. For goods  $i = 1, \dots, n-1$ ,  $c_i = (q_i^s)^{-\psi_i}$  (as the consumer solves  $u'(c_i) = q_i^s$ ). The demand for good 0 is  $c_0 = w - \sum_{i=1}^{n-1} (q_i^s)^{1-\psi_i} - 1$ , as the consumer plans to consume  $c_i = (q_i^s)^{-\psi_i}$  for all good  $i = 1, \dots, n-1$ , and 1 of good  $n$ . Once goods 0 through  $n-1$  have been purchased, the remaining disposable income for good  $n$  is  $c_n = w - \sum_{i=0}^{n-1} q_i c_i$ .

Then (as derived shortly) the optimal tax on good  $i < n$  is as in (17), replacing  $\Lambda$  by

$$\Lambda_i = \frac{\Lambda - (1 - \Lambda)(1 - m_i)\mu}{1 - (1 - \Lambda)(1 - m_i)\mu}. \quad (70)$$

A direct consequence is that the optimal tax  $\tau_i$  is lower than in the baseline case and is decreasing in  $\mu$ , particularly for less salient taxes with a small  $m_i$ . Indeed, the measure of the social marginal cost of public funds  $\Lambda_i$  is decreasing in the marginal distortionary cost of budget adjustment  $\mu$  (recall  $\Lambda < 1$ ), coincides with its baseline value of  $\Lambda$  when  $\mu = 0$ , and is lower than  $\Lambda$  for all  $\mu > 0$ . Furthermore  $\mu$  enters the formula through the  $\mu(1 - m_i)$  so that these effects are particularly pronounced when attention  $m_i$  is low.

We next proceed in greater detail. 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  $n$ , for simplicity.

Utility is:

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

---

<sup>84</sup>The level of  $\nu$  is unimportant provided it is between 0 and 1.



Good 0 has marginal utility of 1, which absorbs income effects, so  $v_w = 1$ . Hence,  $\boldsymbol{\tau}^b = \mathbf{q} - \frac{u_c}{v_w}$  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 (67)).

**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  $p_i = 1$  and use normalized  $\lambda = \frac{\lambda}{\gamma}$  to denote the relative benefit of raising revenue compared to welfare weight. The distortion on good  $n$  is, from (67)

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

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 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 (8) 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:

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

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  $\mu (1 - m_i) < \lambda$ .

We can also investigate the case in which the utility associated with good  $n$ , the shock absorber, is isoelastic. In this case, a utility loss equal to  $L^D$  is 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  $p_i = 1$ , i.e.  $c_i^* = 1$ , for all  $1 \leq i < n$ )

$$L = -\frac{1}{2} \sum_{i=1}^{n-1} \tau_i^2 \psi_i m_i^2 c_i^* - \frac{1}{2} \frac{1}{\psi_n c_n^*} \left( \sum_{i=1}^{n-1} c_i^* (1 - m_i) \tau_i \right)^2 + \lambda \sum_i \tau_i c_i^*. \quad (71)$$

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 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  $n$ .

**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  $U(c_1) = Qc_1 - \frac{c_1^2}{2\Psi}$ , 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.

### 9.3.4 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:

$$g(\{\tau^{s,h}\}) = \max_{\{w^h\}} L(\{\tau^{s,h}\}, \{\tau^{s,h}\}, \{w^h\}), \quad (72)$$

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

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

with

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

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

**Proof.** 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

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

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,  $w^* = w^{**}(y^*)$ .

$$\begin{aligned} L^{\text{tot}} &= L(w, y) - L(w^*, y^*) \\ &= [L(w, y) - L(w^{**}(y), y)] + [L(w^{**}(y), y) - L(w^{**}(y^*), y^*)] \\ &= \frac{1}{2} (w - w^{**}(y)) \cdot L_{ww}(w^{**}(y), y) \cdot (w - w^{**}(y)) + \frac{1}{2} (y - y^*) g_{yy}(y - y^*) \text{ by Lemma 11.3} \\ &= L^{\text{distribution}} + L^{\text{distortion}} \end{aligned}$$

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

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

*Redistribution terms*

From Lemma 11.2, the expression of the loss involves  $L_{w^h}(w, \tau) = \gamma^{\xi,h} - \lambda$ , the social marginal

utility. Applying that Lemma 11.2 gives a loss:

$$L^{\text{distribution}} = \frac{1}{2} \sum_{h,h'} (\gamma^{\xi,h} - \bar{\gamma}) (L_{ww}(w, \tau)^{-1})_{h,h'} (\gamma^{\xi,h} - \bar{\gamma}). \quad (73)$$

*Tax distortion terms*

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

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

$$g_{\tau^s \tau^s} = L_{\tau^s \tau^s} - L_{\tau^s w} L_{ww}^{-1} L_{\tau^s w} \text{ by Lemma 11.3.}$$

□

**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 9.2** *We have*

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

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

$$g_{\tau^s, h \tau^s, h'} = \lambda S^{r,h} 1_{h=h'} + S^{r,h} \xi_{c^h c^{h'}} S^{r,h'}. \quad (74)$$

**Proof.** We observe that a tax  $\tau^{s,h}$  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

$$\frac{d\xi(\{w^h\}, \{\tau^{s,h}\})}{d\tau^{s,h}} = \frac{\xi_{c^h} c_{\tau^{s,h}}^h}{1 - \sum_{h'} \xi_{c^{h'}} c_{\xi}^{h'}}. \quad (75)$$

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

$$\frac{d\xi}{dw^h} = \frac{\xi_{c^h} c_{w^h}^h}{1 - \sum_{h'} \xi_{c^{h'}} c_{\xi}^{h'}}. \quad (76)$$

We note that the FOC of (72) in  $w^h$  is

$$\begin{aligned}
0 &= v_{w^h}^h + \lambda \left( \tau^{s,h} c_{w^h}^{r,h} - 1 \right) + \frac{d\xi}{dw^h} \sum_{h'} v_{\xi}^{h'} + \lambda \tau^{s,h'} \cdot \bar{c}_{\xi}^{r,h'} \\
&= v_{w^h}^h + \lambda \left( \tau^{s,h} c_{w^h}^{r,h} - 1 \right) + \frac{\xi_{ch} c_{w^h}^h}{1 - \sum_{h'} \xi_{ch'} c_{\xi}^{h'}} \sum_{h'} v_{\xi}^{h'} + \lambda \tau^{s,h'} \cdot \bar{c}_{\xi}^{r,h'} \text{ by (76)} \\
&= v_{w^h}^h + \lambda \left( \tau^{s,h} c_{w^h}^{r,h} - 1 \right) + \xi_{ch} c_{w^h}^h \Xi \\
&= v_{w^h}^h + \lambda \left( \tau^{s,h} c_{w^h}^{r,h} - 1 \right) - \lambda \tau^{\xi,h} c_{w^h}^h \\
&= \gamma^{\xi,h} - \lambda.
\end{aligned}$$

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

$$\begin{aligned}
g(\{\tau^{s,h}\}) &= \max_{\{w^h\}} \sum_h v^r(p + \tau^{s,h}, w^h, \xi) + \lambda \sum_h [\tau^{s,h} \cdot \bar{c}^r(p + \tau^{s,h}, w^h, \xi) - w^h] \\
&= \max_{\{w^h\}} \sum_h v^r(p + \tau^{s,h}, w^h, \xi) + \lambda \sum_h [p \cdot \bar{c}^r(p + \tau^{s,h}, w^h, \xi)] \\
&\quad (p + \tau^h) c = w^h.
\end{aligned}$$

We take the derivatives (72):

$$\begin{aligned}
g_{\tau^{s,h}}(\{\tau^{s,h'}\}) &= (\lambda - v_w^h) c^{r,h} + \lambda \tau^{s,h} c_p^{r,h} + \frac{d\xi}{d\tau^{s,h}} \sum_{h'} [v_{\xi}^{h'} + \lambda \tau^{s,h'} \cdot \bar{c}_{\xi}^{r,h'}] \\
&= (\lambda - v_w^h) c^{r,h} + \lambda \tau^{s,h} c_p^{r,h} + \frac{\xi_{ch} c_{\tau^{s,h}}^h}{1 - \sum_{h'} \xi_{ch'} c_{\xi}^{h'}} \sum_{h'} v_{\xi}^{h'} + \lambda \tau^{s,h'} \cdot \bar{c}_{\xi}^{r,h'} \text{ by (75)} \\
&= (\lambda - v_w^h) c^{r,h} + \lambda \tau^{s,h} c_p^{r,h} + \xi_{ch} c_{\tau^{s,h}}^h \Xi \\
&= (\lambda - v_w^h) c^{r,h} + \lambda \tau^{s,h} c_p^{r,h} - \lambda \tau^{\xi,h} c_p^h \\
&= (\lambda - v_w^h) c^{r,h} + \lambda (\tau^{s,h} - \tau^{\xi,h}) c_p^{r,h} \\
&= (\lambda - v_w^h) c^{r,h} + \lambda (\tau^{s,h} - \tau^{\xi,h}) [S^{r,h} - c_w^r c^h] \\
&= [\lambda - v_w^h - \lambda (\tau^{s,h} - \tau^{\xi,h}) c_w] c^h + \lambda (\tau^{s,h} - \tau^{\xi,h}) S^{r,h} \\
&= \lambda (\tau^{s,h} - \tau^{\xi,h}) S^{r,h}.
\end{aligned} \tag{77}$$

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

$$g_{\tau^{s,h} \tau^{s,h'}} = \lambda S^{r,h} \left( 1_{h=h'} - \frac{d\tau^{\xi,h}(\{\tau^{s,h''}\})}{d\tau^{s,h'}} \right). \tag{78}$$

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

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

so

$$g_{\tau^{s, h}, \tau^{s, h'}} = \lambda S^{r, h} 1_{h=h'} + H S^{r, h} \xi_{c^h c^{h'}} S^{r, h'}. \quad (79)$$

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  $\tau^{\xi, h} = -\frac{1}{\lambda} \xi_{c^h} (\{\tau^{s, h'}\})$ . 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  $\frac{dc^{h''}}{d\tau^{h'}} = 0$  for  $h'' \neq h'$ ,

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

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

## 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

---

<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 the NBER working paper version of this paper.



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 \(2001\)](#).

## 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.

---

<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  $\mathbf{S}^H$ .

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

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

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 Section 11.3.2):

$$\frac{v_q}{v_w} = z - \frac{\tau^b z}{q} \zeta^c, \quad \frac{v_{Q_{z^*}}}{v_w} = - \frac{\tau^b z}{q} \zeta_{Q_{z^*}}^c. \quad (80)$$

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

$$\text{smax}_{c, z | R^{s,n}(\cdot)} u^n(c, z) \text{ s.t. } c = R(z). \quad (81)$$

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

$$q^{s,n}(z) u_c^n(c, z) + u_z^n(c, z) = 0, \quad (82)$$

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

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

<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.

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  $z'$ .

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 11.3.2](#)):

$$\zeta^c = \left( 1 - \eta \frac{\tau - \tau^s}{q} \right) \zeta^{cr} m_{zz}, \quad \zeta_{Q_{z^*}}^c = \left( 1 - \eta \frac{\tau - \tau^s}{q} \right) \zeta^{cr} m_{zz^*}, \quad (83)$$

$$\tau^b = \frac{\tau - \tau^s}{1 - \eta \frac{\tau - \tau^s}{q}}. \quad (84)$$

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

$$\tau^b = \frac{\frac{u_c}{u_c^s} u_z^s - u_z}{v_r}. \quad (85)$$

**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

$$\gamma(z) = g(z) + \frac{\eta(z)}{1 - T'(z)} \left[ \tilde{\tau}^b(z) + (T'(z) - \tilde{\tau}^b(z)) \frac{h^*(z)}{h(z)} \right]. \quad (86)$$

This definition is the analogue of the corresponding definition in the Ramsey model. It is motivated by [Lemma 11.5](#), 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. [Section 11.3.1](#) presents the intermediary steps used in the derivation of this formula.

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

$$\begin{aligned} \frac{T'(z^*) - \tilde{\tau}^b(z^*)}{1 - T'(z^*)} &= \frac{1}{\zeta^c(z^*)} \frac{1 - H(z^*)}{z^* h^*(z^*)} \int_{z^*}^{\infty} (1 - \gamma(z)) \frac{h(z)}{1 - H(z^*)} dz \\ &\quad - \int_0^{\infty} \frac{\zeta_{Q_{z^*}}^c(z)}{\zeta^c(z^*)} \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} \frac{zh^*(z)}{z^* h^*(z^*)} dz. \end{aligned} \quad (87)$$

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

$$\begin{aligned} &\frac{T'(z^*) - \tilde{\tau}^b(z^*)}{1 - T'(z^*)} + \int_0^{\infty} \omega(z^*, z) \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} dz \\ &= \frac{1}{\zeta^c(z^*)} \frac{1 - H(z^*)}{z^* h^*(z^*)} \int_{z^*}^{\infty} e^{-\int_{z^*}^z \rho(s) ds} \left( 1 - g(z) - \eta \frac{\tilde{\tau}^b(z)}{1 - T'(z)} \right) \frac{h(z)}{1 - H(z^*)} dz, \end{aligned} \quad (88)$$

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 (87) 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 (86)). The second term  $-\frac{1}{z^*} \int_0^{\infty} \frac{\zeta_{Q_{z^*}}^c(z) T'(z) - \tilde{\tau}^b(z)}{\zeta^c(z^*)} z \frac{h^*(z)}{h^*(z^*)} dz$  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 (88) 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 (87) using equation (86). When there are no income effects so that  $\eta = \rho(z) = 0$ , the optimal tax formula (87) and the modified Saez formula (88) 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 (88) 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 (88) 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  $\tilde{\tau}^b(z) = g(z) \frac{u_c \frac{u_z^s}{u_c^s} - u_z}{v_r}$ , 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 (87) 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).

We now derive a formula for the marginal rate at very high incomes when the income distribution is unbounded at the top. It proves convenient to consider a (high)  $z_0$  above which we consider

that incomes are “top incomes”, and the marginal rate is constant. We consider tax systems with constant marginal tax rates for  $z \geq z_0$ . We assume that  $g(z) = \bar{g}$  for  $z \geq z_0$ . We call  $\zeta_{\bar{q}}^c(z) = \int_{z_0}^{\infty} \zeta_{Q_{z^*}}^c(z) dz^*$  the sensitivity to the asymptotic tax rate. This is the elasticity of earnings of an individual at earnings  $z < z_0$  to an increase to the top rate, arising perhaps because of a misperception of the tax environment. Concretely, think of the recent case of France where increasing the top rate to 75% might have created an adverse general climate with the perception that even earners the top income would pay higher taxes.

We call  $\bar{\eta}$ ,  $\bar{g}$ ,  $\bar{\zeta}^c$  the asymptotic values for large incomes and  $\pi$  the Pareto exponent of the earnings distribution (i.e. when  $z$  is large,  $1 - H(z) \propto z^{-\pi}$ ). We define the weighted means:  $\mathbb{E}^z[\phi(z)] = \frac{\int \phi(z)h(z)zdz}{\int \phi(z)zdz}$  and  $\mathbb{E}^*[\zeta_{\bar{q}}^c] = \frac{\int \zeta_{\bar{q}}^c(z) \frac{T'(z) - \bar{\tau}^b(z)}{1 - T'(z)} h^*(z) dz}{\int \frac{T'(z) - \bar{\tau}^b(z)}{1 - T'(z)} h^*(z) dz}$ .

**Proposition 10.2** (Optimal tax rate for top incomes) *The optimal marginal rate  $\bar{\tau}$  for top incomes is*

$$\bar{\tau} = \frac{1 - \bar{g} - \beta + \bar{\zeta}^c \pi \bar{g} \tau^b}{1 - \bar{g} - \beta + \bar{\zeta}^c \pi + \bar{\eta}}, \quad (89)$$

where

$$\beta = \mathbb{E}^z \left[ \frac{T'(z) - \bar{\tau}^b(z)}{1 - T'(z)} \right] \pi \frac{\mathbb{E}[z]}{\mathbb{E}[z1_{z \geq z_0}]} \mathbb{E}^*[\zeta_{\bar{q}}^c].$$

This generalizes the [Saez \(2001\)](#) formula which can be recovered in the particular case where  $\beta = \tau^b = 0$ . The intuition is as follows—the  $\beta$  term reflects not only the fact that the top marginal tax rate affects not only top earners, but also the tax perceived by agents at all points of the income distribution with associated effects on tax revenues. The more increasing the top tax rate lowers all incomes (the higher  $\zeta_{\bar{q}}^c(z)$ ), the higher  $\beta$ , and the lower the optimal top tax rate.

The  $\tau^b$  terms are positive (resp. negative) when top earners overperceive (resp. underperceive) the marginal benefits of effort or underperceive (resp. overperceive) taxes. These terms lead to higher (resp. lower) optimal top rates compared to the Saez formula.

Consider the typical Saez calibration with  $\zeta^c(\infty) = 0.2$ ,  $\eta = 0$  and  $\pi = 2$ . If the typical tax is  $T'(z) \simeq \frac{1}{3}$  so that  $\mathbb{E}^z \left[ \frac{T'(z)}{1 - T'(z)} \right] \simeq \frac{1}{2}$ , we take  $z_0$  to be at the top 1% quantile of the income distribution. [Piketty and Saez \(2003\)](#) (updated 2015) report that the income share of the top 1% is 20%, so that  $\frac{\mathbb{E}[z]}{\mathbb{E}[z1_{z \geq z_0}]} = \frac{1}{0.2}$ . This implies that  $\beta = \mathbb{E}^z \left[ \frac{T'(z) - \bar{\tau}^b(z)}{1 - T'(z)} \right] \pi \frac{\mathbb{E}[z]}{\mathbb{E}[z1_{z \geq z_0}]} \mathbb{E}^*[\zeta_{\bar{q}}^c] = \frac{1}{2} 2 \frac{1}{0.2} \mathbb{E}^*[\zeta_{\bar{q}}^c] = 5 \mathbb{E}^*[\zeta_{\bar{q}}^c]$ . Also, we take top earnings to be well calibrated, i.e.  $\tau^b = 0$ .

The average cross-influence  $\mathbb{E}^*[\zeta_{\bar{q}}^c]$  does not appear to have ever been measured. It is assumed to be 0 in the traditional model. We propose the following thought experiment to gauge its potential magnitude. Suppose that increasing the top rate by 10% will decrease earnings outside the top bracket by  $x = 1\%$ . Then,  $\mathbb{E}^*[\zeta_{\bar{q}}^c] = (1 - T'(z)) \frac{z_{\bar{q}}}{z} = \left(1 - \frac{1}{3}\right) \frac{x}{0.1} = 6.7x$ , which gives an interpretable benchmark that we now use.

Take first the case where  $\bar{g} = 0$ , i.e. where the top optimal tax rate maximizes revenues raised from top earners. With rational agents ( $x = 0$ ), the top marginal tax rate is  $\bar{\tau} = 71\%$ . If  $x = 1\%$ ,

then  $\bar{\tau} = 62\%$ , and if  $x = 2\%$ , then  $\bar{\tau} = 45\%$ . If  $x = -1\%$ , then  $\bar{\tau} = 77\%$ .<sup>91</sup> When the weight on top earnings is higher, say  $\bar{g} = 0.2$ , the corresponding numbers for the top rate are: 67%, 53%, 25%, and 74%. This illustrates the potentially large importance of the behavioral cross-impact of the top tax rate, a sufficient statistic that is assumed to be zero in traditional analyses.

The behavioral wedge  $\tau^b$  does not affect the optimal tax rate when  $\bar{g} = 0$ . When  $\bar{g} = 0.5$ , the optimal top rate increases from 56% to 67% when the externality goes from no misperception of taxes by top earners ( $\tau^b = 0$ ) to underperception of taxes by top earners ( $\tau^b = 0.5$ ).

#### 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  $u^s(c, z) = c - \phi\kappa(z)$ . We take experienced utility to be  $u(c, z) = \theta c - \phi(z)$ . Then the modified Saez formula (88) 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 (85). 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 also provide a behavioral treatment of the [Saez \(2002\)](#) extensive margin of labor supply model next.

<sup>91</sup>We thank Thomas Piketty for suggesting to us that if workers are happier, and strike less, because taxes on the wealthy are high, then  $x < 0$ .

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

### 10.3 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  $B_i^- = \{m \text{ s.t. } \partial i^*(h, \{T_j\}) / \partial T_i > 0\}$  (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,

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

## 11 Further Proofs and Derivations

### 11.1 General proofs and derivations

**Derivation of behavioral wedges in  $\beta$ - $\delta$  model in Section 2.1** Price of good 1 is normalized to be 1. For the consumer,  $\max_{c_0, c_1} \ln c_0 + \beta\delta \ln c_1$  subject to  $c_0 + \frac{c_1}{1+r} \leq w$ . FOC:  $\frac{1}{c_0} = \frac{\beta\delta/c_1}{1/(1+r)}$ . Thus demand functions are  $c_0 = \frac{1}{1+\beta\delta} w$ ,  $c_1 = \frac{\beta\delta}{1+\beta\delta} (1+r) w$ .

Subjectively perceived indirect utility and experienced indirect utility are

$$\begin{aligned}
v^s(w) &= \ln c_0 + \beta\delta \ln c_1 = \ln \frac{w}{1 + \beta\delta} + \beta\delta \ln \frac{\beta\delta (1+r) w}{1 + \beta\delta}, \\
v(w) &= \ln c_0 + \delta \ln c_1 = \ln \frac{w}{1 + \beta\delta} + \delta \ln \frac{\beta\delta (1+r) w}{1 + \beta\delta}.
\end{aligned}$$

To compute behavioral wedges, we have now  $u_{c_0}^s = \frac{1}{c_0}$ ,  $u_{c_1}^s = \frac{\beta\delta}{c_1}$ ,  $u_{c_0} = \frac{1}{c_0}$ ,  $u_{c_1} = \frac{\delta}{c_1}$ ,  $v_w^s = \frac{1+\beta\delta}{w}$ ,  $v_w = \frac{1+\delta}{w}$ .

Therefore,

$$\begin{aligned}\tau_0^b &= \frac{1/c_0}{\frac{1+\beta\delta}{w}} - \frac{1/c_0}{\frac{1+\delta}{w}} = \frac{\delta(1-\beta)}{1+\delta}, \\ \tau_1^b &= \frac{\beta\delta/c_1}{\frac{1+\beta\delta}{w}} - \frac{\delta/c_1}{\frac{1+\delta}{w}} = -\frac{1-\beta}{\beta(1+r)(1+\delta)}.\end{aligned}$$

□

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

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

so  $\frac{d\xi}{d\tau_i} = \frac{\sum_h \xi_{c^h} \cdot \mathbf{c}_{q_i}^h}{1 - \sum_h \xi_{c^h} \cdot \mathbf{c}_{\xi}^h}$ . The term  $\frac{1}{1 - \sum_h \xi_{c^h} \cdot \mathbf{c}_{\xi}^h}$  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)

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

Using Proposition 2.1,

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

□

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

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

hence

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

Hence, Proposition 2.2 implies:

$$\sum_h \left(1 - \frac{\gamma^{\xi,h}}{\lambda}\right) \mathbf{C}^h = - \sum_h (\mathbf{S}^{C,h})' (\boldsymbol{\tau} - \boldsymbol{\tau}^{\xi,h} - \tilde{\boldsymbol{\tau}}^{b,h}) = - \sum_h \mathbf{M}^{h'} \mathbf{S}^{h,r} \left[ \left[ \mathbf{I} - (\mathbf{I} - \mathbf{M}^h) \frac{\beta^h}{\lambda} \right] \boldsymbol{\tau} - \boldsymbol{\tau}^{X,h} \right]. \quad (90)$$

□

**Derivation of (46)** To make use of Proposition 2.1, we need to express demand function, Slutsky matrix, behavioral wedge and social marginal utility of income for each consumer in terms of taxes, in the limit of small taxes.

First, we approximate demand function up to the second order and decompose it,

$$\mathbf{c}^h(\mathbf{p} + \boldsymbol{\tau}, w) = \mathbf{c}^{r,h}(\mathbf{p} + \boldsymbol{\tau}, w) + \check{\mathbf{c}}^h(\mathbf{p}, \boldsymbol{\tau}, w), \quad (91)$$

with

$$\check{\mathbf{c}}^h(\mathbf{p}, \boldsymbol{\tau}, w) = \hat{\mathbf{c}}^{u,h}(\mathbf{p}, w) + \hat{\mathbf{c}}^{M,h}(\mathbf{p}, w) \boldsymbol{\tau} + O\left(\left(\tilde{\eta}_1^h\right)^2\right),$$

and  $\tilde{\eta}_1^h = \max(\|\boldsymbol{\tau}\|, \|\hat{\mathbf{c}}^{u,h}(\mathbf{p}, w)\|)$ . We can also write:

$$\mathbf{c}^h(\mathbf{p} + \boldsymbol{\tau}, w) = \mathbf{c}^h(\mathbf{p}, w) + \mathbf{c}_p^{r,h}(\mathbf{p}, w) \mathbf{M}^h \boldsymbol{\tau} + O\left(\left(\tilde{\eta}_1^h\right)^2\right),$$

where  $\mathbf{M}^h$  is the attention matrix of consumer  $h$ :

$$\mathbf{M}^h = \mathbf{I} + \left(\mathbf{c}_p^{r,h}(\mathbf{p}, w)\right)^{-1} \hat{\mathbf{c}}^{M,h}(\mathbf{p}, w).$$

We also observe that, thanks to envelope theorem (noting that  $(\mathbf{p} + \boldsymbol{\tau}) \cdot \check{\mathbf{c}}^h(\mathbf{p}, \boldsymbol{\tau}, w) = 0$  by the budget constraint)

$$\begin{aligned} v^h(\mathbf{q}, w) &= u^h\left(\mathbf{c}^{r,h}(\mathbf{q}, w) + \check{\mathbf{c}}^h(\mathbf{p}, \boldsymbol{\tau}, w) + O\left(\left(\tilde{\eta}_1^h\right)^2\right)\right) \\ &= u^h\left(\mathbf{c}^{r,h}(\mathbf{q}, w)\right) + O\left(\left(\tilde{\eta}_1^h\right)^2\right). \end{aligned}$$

With these, we derive Slutsky matrices

$$\mathbf{S}^{C,h} = \frac{\partial \mathbf{c}^h}{\partial \mathbf{q}} + \mathbf{c}_w^h \mathbf{c}^{h'} = \mathbf{c}_p^{r,h} + \hat{\mathbf{c}}^{M,h} + \mathbf{c}_w^h \mathbf{c}^{h'} + O\left(\tilde{\eta}_1^h\right) = \mathbf{S}^{r,h} + \hat{\mathbf{c}}^{M,h} + O\left(\tilde{\eta}_1^h\right), \quad (92)$$

and behavioral wedges

$$\begin{aligned} \boldsymbol{\tau}^{b,h} &= \mathbf{q} - \frac{u_c^h(\mathbf{c}^h(\mathbf{q}, w))}{v_w^h(\mathbf{q}, w)} = \mathbf{q} - \frac{u_c^h(\mathbf{c}^{r,h}(\mathbf{q}, w)) + u_{cc}^h(\mathbf{q}, w) \check{\mathbf{c}}^h}{v_w^h(\mathbf{q}, w)} + O\left(\left(\tilde{\eta}_1^h\right)^2\right) \\ &= \Omega^h \check{\mathbf{c}}^h + O\left(\left(\tilde{\eta}_1^h\right)^2\right) = \boldsymbol{\tau}^{b,h,nat} + \Omega^h \hat{\mathbf{c}}^{M,h} \boldsymbol{\tau} + O\left(\left(\tilde{\eta}_1^h\right)^2\right), \end{aligned} \quad (93)$$

in which we define

$$\boldsymbol{\tau}^{b,h,nat} = \Omega^h \hat{\mathbf{c}}^{u,h}(\mathbf{p}, w), \quad \Omega^h = -\frac{u_{\mathbf{c}\mathbf{c}}^h(\mathbf{p}, w)}{v_w^h(\mathbf{p}, w)},$$

which are constants evaluated at 0 taxes.

For small  $\lambda - b^h$ , up to the second order, social welfare function, social marginal welfare weight and social marginal utility of income are

$$\begin{aligned} W &= \sum_h \frac{b^h}{v_w^h(\mathbf{p}, w)} v^h(\mathbf{p} + \boldsymbol{\tau}, w), \\ \beta^h &= W_{v^h} v_w^h = b^h \frac{v_w^h(\mathbf{p} + \boldsymbol{\tau}, w)}{v_w^h(\mathbf{p}, w)} = b^h \left( 1 + \frac{v_{w,\mathbf{p}}^h(\mathbf{p}, w)}{v_w^h(\mathbf{p}, w)} \cdot \boldsymbol{\tau} \right), \\ \gamma^h &= \beta^h + \lambda \boldsymbol{\tau} \cdot \mathbf{c}_w^h = b^h \left( 1 + \left( \frac{v_{w,\mathbf{p}}^h(\mathbf{p}, w)}{v_w^h(\mathbf{p}, w)} + \frac{\lambda}{b^h} \mathbf{c}_w^h(\mathbf{p}, w) \right) \cdot \boldsymbol{\tau} \right). \end{aligned}$$

By Roy's identity  $v_{\mathbf{p}}^h(\mathbf{p}, w) = -\mathbf{c}^h(\mathbf{p}, w) v_w^h(\mathbf{p}, w)$ , we have  $\frac{v_{w,\mathbf{p}}^h(\mathbf{p}, w)}{v_w^h(\mathbf{p}, w)} = -\mathbf{c}_w^h(\mathbf{p}, w) - \mathbf{c}^h(\mathbf{p}, w) \frac{v_{ww}^h(\mathbf{p}, w)}{v_w^h(\mathbf{p}, w)}$ ,

$$\begin{aligned} \frac{\gamma^h}{b^h} &= 1 + \left[ \left( -1 + \frac{\lambda}{b^h} \right) \mathbf{c}_w^h(\mathbf{p}, w) - \mathbf{c}^h(\mathbf{p}, w) \frac{v_{ww}^h(\mathbf{p}, w)}{v_w^h(\mathbf{p}, w)} \right] \cdot \boldsymbol{\tau} \\ &= 1 - \frac{v_{ww}^h(\mathbf{p}, w)}{v_w^h(\mathbf{p}, w)} \mathbf{c}^h(\mathbf{p}, w) \cdot \boldsymbol{\tau} + O\left(\tilde{\eta}_2^h\right)^2, \end{aligned} \quad (94)$$

where we define  $\tilde{\eta}_2^h = \max(|b^h - \lambda|, \|\boldsymbol{\tau}\|, \|\hat{\mathbf{c}}^{u,h}(\mathbf{p}, w)\|)$ . In addition, we have  $\tilde{\boldsymbol{\tau}}^{b,h} = \frac{\beta^h}{\lambda} \boldsymbol{\tau}^b = \boldsymbol{\tau}^b + O\left(\tilde{\eta}_2^h\right)^2$ .

We remember that good 0 is not taxed. So, in all expressions below, the subscript  $> 0$  indicates the selection of the  $(N-1) \times (N-1)$  sub-matrix corresponding to all goods except good 0. Plugging these results in the behavioral multi-person Ramsey formula (8),

$$\begin{aligned} \mathbf{0} &= \sum_h \left[ \left( 1 - \frac{\gamma^h}{\lambda} \right) \mathbf{c}^h(\mathbf{q}, w) + \mathbf{S}^{C,h'}(\boldsymbol{\tau} - \tilde{\boldsymbol{\tau}}^{b,h}) \right]_{>0} \\ &= \sum_h \left( 1 - \frac{b^h}{\lambda} + \frac{b^h}{\lambda} \frac{v_{ww}^h(\mathbf{p}, w)}{v_w^h(\mathbf{p}, w)} \mathbf{c}^h(\mathbf{p}, w) \cdot \boldsymbol{\tau} \right) [\mathbf{c}^h(\mathbf{p}, w)]_{>0} \\ &\quad + \left[ \sum_h (\mathbf{S}^{r,h} + \hat{\mathbf{c}}^{M,h}(\mathbf{p}, w))' \left( (\mathbf{I} - \Omega^h \hat{\mathbf{c}}^{M,h}(\mathbf{p}, w)) \boldsymbol{\tau} - \Omega^h \hat{\mathbf{c}}^{u,h}(\mathbf{p}, w) \right) \right]_{>0} + O\left(\sum_h (\tilde{\eta}_2^h)^2\right) \\ &= \sum_h \left[ \frac{v_{ww}^h(\mathbf{p}, w)}{v_w^h(\mathbf{p}, w)} \mathbf{c}^h(\mathbf{p}, w) (\mathbf{c}^h(\mathbf{p}, w))' + (\mathbf{S}^{r,h} + \hat{\mathbf{c}}^{M,h}(\mathbf{p}, w))' (\mathbf{I} - \Omega^h \hat{\mathbf{c}}^{M,h}(\mathbf{p}, w)) \right]_{>0} \boldsymbol{\tau} \\ &\quad + \left[ \sum_h \left[ \left( 1 - \frac{b^h}{\lambda} \right) \mathbf{c}^h(\mathbf{p}, w) - (\mathbf{S}^{r,h} + \hat{\mathbf{c}}^{M,h}(\mathbf{p}, w))' \Omega^h \hat{\mathbf{c}}^{u,h}(\mathbf{p}, w) \right] \right]_{>0} + O\left(\sum_h (\tilde{\eta}_2^h)^2\right), \end{aligned}$$

which is solved by optimal tax

$$\boldsymbol{\tau} = - \left[ \sum_h \frac{v_{ww}}{v_w} \mathbf{c}^h \mathbf{c}^{h'} + (\mathbf{S}^{r,h} + \hat{\mathbf{c}}^{M,h})' (\mathbf{I} - \Omega^h \hat{\mathbf{c}}^{M,h}) \right]_{>0}^{-1} \left[ \sum_h \left( 1 - \frac{b^h}{\lambda} \right) \mathbf{c}^h - (\mathbf{S}^{r,h} + \hat{\mathbf{c}}^{M,h})' \Omega^h \hat{\mathbf{c}}^{u,h} \right] + O(\tilde{\eta}^2) \quad (95)$$

with  $\tilde{\eta} = \sum_h |b^h - \lambda| + \|\hat{\mathbf{c}}^{u,h}(\mathbf{p}, w)\|$ ,  $\mathbf{S}^{r,h} = \mathbf{c}_p^h + \mathbf{c}_w^h \mathbf{c}^{h'}$ , in which all the variables are evaluated at  $(\mathbf{p}, w)$ . <sup>93</sup>□

**Proof of Proposition 3.1** We start from the Ramsey planning problem in (16). 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 (18), the approximate loss from taxation of inattentive agents** We define the behavioral elasticity  $\alpha_i = m_i \psi_i$  and denote  $y_i$  as the expenditure on good  $i$  at zero tax

<sup>93</sup>  $\mathbf{S}^{r,h} + \hat{\mathbf{c}}^{M,h}$  can also be expressed as

$$\mathbf{S}^{r,h} + \hat{\mathbf{c}}^{M,h} = \mathbf{S}^{r,h} + \mathbf{c}_p^{r,h} (\mathbf{M}^h - \mathbf{I}) = \mathbf{S}^{r,h} + (\mathbf{S}^{r,h} - \mathbf{c}_w^h \mathbf{c}^{h'}) (\mathbf{M}^h - \mathbf{I}) = \mathbf{S}^{r,h} \mathbf{M}^h - \mathbf{c}_w^h \mathbf{c}^{h'} (\mathbf{M}^h - \mathbf{I}).$$

rate. Demand is:

$$c_i(\tau_i) = \frac{y_i}{p_i} \left( 1 - \alpha_i \frac{\tau_i}{p_i} + O \left( \left( \frac{\tau_i}{p_i} \right)^2 \right) \right), \quad (96)$$

Hence have:

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

with

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

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

$$\begin{aligned} f_i'(0) &= 0 \\ f_i''(0) &= F^{i''}(c_i(0)) c_i'(0)^2 \\ &= U_i'' \left( \frac{y_i}{p_i} \right) \frac{y_i}{p_i^2} \alpha_i^2 \text{ using (96)} \\ &= -\frac{\alpha_i^2 y_i}{\psi_i p_i^2}, \end{aligned}$$

so

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

So the objective function is:

$$L = -\frac{1}{2} \sum_i \frac{\alpha_i^2 y_i}{\psi_i} \left( \frac{\tau_i}{p_i} \right)^2 + \Lambda \sum_i \frac{\tau_i y_i}{p_i} + o(\|\tau\|^2) + o(\|\tau\| \Lambda). \quad (97)$$

□

**Derivation of (22), the endogenous social cost of public funds** Now the objective function of the government is

$$\begin{aligned} W(\boldsymbol{\tau}) &= \gamma u(\mathbf{C}(\boldsymbol{\tau})) + V(G(\boldsymbol{\tau})) \\ &= \gamma \sum_i \left[ \frac{\left[ (p_i + m_i \tau_i)^{-\psi_i} \right]^{1-1/\psi_i} - 1}{1 - 1/\psi_i} - (p_i + \tau_i) (p_i + m_i \tau_i)^{-\psi_i} \right] + V \left( \sum_i (p_i + m_i \tau_i)^{-\psi_i} \tau_i \right), \end{aligned}$$

where  $V(G)$  is the utility from public goods  $G$ ,

$$G(\boldsymbol{\tau}) = \sum_i (p_i + m_i \tau_i)^{-\psi_i} \tau_i.$$

A Taylor expansion gives a second order approximation:

$$\begin{aligned} W(\boldsymbol{\tau}) - W(0) &= \sum_i c_i(0) p_i \left[ -\frac{1}{2} \psi_i m_i (m_i - 2) \left( \frac{\tau_i}{p_i} \right)^2 - \frac{\tau_i}{p_i} \right] \\ &\quad + V'(0) g + \frac{1}{2} V''(0) g^2 - V'(0) \sum_i c_i(0) p_i \psi_i m_i \left( \frac{\tau_i}{p_i} \right)^2 + o(\|\boldsymbol{\tau}\|^2), \\ g &= \sum_i \tau_i c_i(0) = \sum_i c_i(0) p_i \frac{\tau_i}{p_i}. \end{aligned}$$

We take  $V'(0) = 1 + \Lambda_0$  with  $\Lambda_0$  small, as we are taking the limit of small taxes. So,

$$\begin{aligned} W(\boldsymbol{\tau}) - W(0) &= -\frac{1}{2} \sum_i c_i(0) p_i \psi_i m_i^2 \left( \frac{\tau_i}{p_i} \right)^2 + \Lambda_0 \sum_i c_i(0) p_i \frac{\tau_i}{p_i} + \frac{1}{2} V''(0) \left( \sum_i c_i(0) p_i \frac{\tau_i}{p_i} \right)^2 \\ &\quad + o(\|\boldsymbol{\tau}\|^2) + O(\Lambda_0 \|\boldsymbol{\tau}\|^2). \end{aligned}$$

Optimizing over  $\tau_i$  gives  $\frac{\tau_i}{p_i} = \frac{\Lambda}{\psi_i m_i^2}$  with  $\Lambda = \Lambda_0 + V''(0) g$ . Note that  $\Lambda_0$  and  $g$  are of the same order of magnitude.

Then we can solve for  $\Lambda$ :

$$\begin{aligned} g &= \sum_i c_i(0) p_i \frac{\tau_i}{p_i} = \alpha \Lambda, \\ \alpha &= \sum_i \frac{c_i(0) p_i}{\psi_i m_i^2}, \end{aligned}$$

which gives that

$$\Lambda = \frac{\Lambda_0}{1 - \alpha V''(0)} = \frac{V'(0) - 1}{1 - \alpha V''(0)}.$$

When  $V''(0) < 0$ , we find  $\frac{\partial \Lambda}{\partial m_i} > 0$ . So, when attention is lower, the marginal utility of public fund (at the optimum) is lower, mitigating the increase in taxes.  $\square$

**Proof of Proposition 3.3** The government's planning problem is

$$\sum_h U^h(c^h) - (p + \xi^h) c^h. \quad (98)$$

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>94</sup>

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

$$L(\tau) = -\frac{\Psi}{2} \sum_h (m^h \tau - \xi^h)^2. \quad (99)$$

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,  $U_c^h(c^{h*}) = p + \xi^h$ . 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  $\Psi = \frac{1}{U_{cc}^h(c^{h*})}$  then. The social welfare is  $L = \sum_h L^h = -\frac{\Psi}{2} \sum_h (m^h \tau - \xi^h)^2$  by (99).

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

$$\tau^* = \frac{\sum_h \xi^h m^h}{\sum_h m^{h2}} = \frac{\mathbb{E}[\xi^h m^h]}{\mathbb{E}[m^{h2}]}.$$

---

<sup>94</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  $\tau = \tau^*$ ,

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

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

If there is no tax, the loss is (from equation (99))

$$L^{\text{no tax}} = -\frac{\Psi}{2} \sum_h (m^h \cdot 0 - \xi^h)^2 = -\frac{\Psi}{2} \sum_h \xi^{h2} = -\frac{1}{2} H \Psi \mathbb{E} \left[ \xi^{h2} \right].$$

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

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

$$\frac{1}{H} \sum_h U_c^h(c^*) = p + \frac{1}{H} \sum_h \xi^h. \quad (100)$$

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

$$L^h = \frac{1}{2} U_{cc}^h(c) \left( c^{h*} - c^* \right)^2 = -\frac{1}{2} \frac{1}{\Psi} \left( c^{h*} - c^* \right)^2.$$

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

The loss is:

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

The inequality  $\frac{1}{2\Psi} \text{var} \left( c^{h*} \right) < \Psi \frac{\mathbb{E} \left[ \xi^{h2} \right] \mathbb{E} \left[ m^{h2} \right] - \left( \mathbb{E} \left[ \xi^h m^h \right] \right)^2}{2\mathbb{E} \left[ m^{h2} \right]}$  can also be written as  $\frac{1}{2\Psi} \text{var} \left( c^{h*} \right) < \frac{\Psi}{2} \text{var}_{m^{h2}} \left( \xi^h / m^h \right) \mathbb{E} \left[ m^{h2} \right]$ , since

$$\begin{aligned}\mathbb{E}_{m^{h^2}}(\xi^h/m^h) &= \frac{\sum_h m^{h^2} \xi^h/m^h}{\sum_h m^{h^2}} = \frac{\mathbb{E}[\xi^h m^h]}{\mathbb{E}[m^{h^2}]}, \\ \text{var}_{m^{h^2}}(\xi^h/m^h) &= \frac{\sum_h m^{h^2} (\xi^h/m^h - \mathbb{E}_{m^{h^2}}(\xi^h/m^h))^2}{\sum_h m^{h^2}} = \frac{\sum_h m^{h^2} \left( \frac{\xi^h}{m^h} - \frac{\mathbb{E}[\xi^h m^h]}{\mathbb{E}[m^{h^2}]} \right)^2}{\sum_h m^{h^2}} \\ &= \frac{\sum_h \xi^{h^2} - 2 \sum_h \xi^h m^h \frac{\mathbb{E}[\xi^h m^h]}{\mathbb{E}[m^{h^2}]} + \sum_h m^{h^2} \left( \frac{\mathbb{E}[\xi^h m^h]}{\mathbb{E}[m^{h^2}]} \right)^2}{\sum_h m^{h^2}} = \frac{\mathbb{E}[\xi^{h^2}] \mathbb{E}[m^{h^2}] - (\mathbb{E}[\xi^h m^h])^2}{(\mathbb{E}[m^{h^2}])^2}.\end{aligned}$$

□

**Proof of Proposition 3.4** Equation (13) then yields the optimal tax:

$$\boldsymbol{\tau} = (\mathbb{E}[\mathbf{M}^{h'} \mathbf{S}^r \mathbf{M}^h])^{-1} \mathbb{E}[\mathbf{M}^{h'}] \mathbf{S}^r \boldsymbol{\tau}^X. \quad (101)$$

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

When agents have uniform misperceptions ( $\mathbf{M}^h = \mathbf{M}$ ), the optimal tax is  $\boldsymbol{\tau} = \mathbf{M}^{-1} \boldsymbol{\tau}^X$ . 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  $(\mathbb{E}[\mathbf{M}^{h'} \mathbf{S}^r \mathbf{M}^h])_{ij} = S_{ij}^r \mathbb{E}[m_i^h m_j^h]$  and  $(\mathbb{E}[\mathbf{M}^{h'} \mathbf{S}^r])_{ij} = \mathbb{E}[m_i^h] S_{ij}^r$ . Matrix inversion gives:

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

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  $\mathbf{S}^r$ . □

**Proof of Proposition 3.5** We apply Proposition 2.4. Here,

$$\begin{aligned}\mathbf{M}^h &= \mathbf{I}, \mathbf{S}^{r,h^*} = \begin{pmatrix} -\psi_1 (p_1 + \tau_1)^{-\psi_1 - 1} & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{C_2}^{r,h^*} \end{pmatrix}, \mathbf{S}^{r,h \neq h^*} = \begin{pmatrix} 0 & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_{C_2}^{r,h \neq h^*} \end{pmatrix}, \\ c_1^{h^*} &= (p_1 + \tau_1)^{-\psi_1}, \boldsymbol{\tau}^{X,h^*} = \begin{pmatrix} \frac{\gamma^{h^*}}{\lambda} \xi^{h^*} \\ \vdots \end{pmatrix}, \boldsymbol{\tau}^{X,h \neq h^*} = \begin{pmatrix} 0 \\ \vdots \end{pmatrix}.\end{aligned}$$

We plug the results above into (13). Suppose  $\rho_p, \rho_r$  are the portions of agents  $h^*$  and agents  $h \neq h^*$  in the population with  $\rho_p + \rho_r = 1$ . Then

$$-\left[ \sum_h \mathbf{M}^{h'} \mathbf{S}^{r,h} \left( \mathbf{I} - (\mathbf{I} - \mathbf{M}^h) \frac{\gamma^h}{\lambda} \right) \right]^{-1} = \frac{1}{H} \begin{pmatrix} \rho_p^{-1} \psi_1^{-1} (p_1 + \tau_1)^{\psi_1+1} & \mathbf{0} \\ \mathbf{0} & -(\rho_p \mathbf{S}_{\mathbf{C}_2}^{r,h^*} + \rho_r \mathbf{S}_{\mathbf{C}_2}^{r,h \neq h^*})^{-1} \end{pmatrix},$$

$$\sum_h \left[ \left( 1 - \frac{\gamma^h}{\lambda} \right) \mathbf{C}^h - \mathbf{M}^{h'} \mathbf{S}^{r,h} \boldsymbol{\tau}^{X,h} \right] = H \begin{pmatrix} \rho_p \left[ \left( 1 - \frac{\gamma^{h^*}}{\lambda} \right) (p_1 + \tau_1)^{-\psi_1} + \frac{\gamma^{h^*}}{\lambda} \xi^{h^*} \psi_1 (p_1 + \tau_1)^{-\psi_1-1} \right] \\ \vdots \end{pmatrix}$$

and get

$$\tau_1 = \left( 1 - \frac{\gamma^{h^*}}{\lambda} \right) \frac{p_1 + \tau_1}{\psi_1} + \frac{\gamma^{h^*}}{\lambda} \xi^{h^*},$$

$$\tau_1 = \frac{\frac{\gamma^{h^*}}{\lambda} \xi^{h^*} + \left( 1 - \frac{\gamma^{h^*}}{\lambda} \right) \frac{p_1}{\psi_1}}{1 + \left( \frac{\gamma^{h^*}}{\lambda} - 1 \right) \frac{1}{\psi_1}}.$$

□

## 11.2 Helpful tools for the calculations

**Lemma 11.1** (*Losses from lack of equalization of marginal utilities*). Suppose we have  $F(x) = \sum_i f^i(x_i)$ , and  $F^{FB} = \max_y F(y)$  s.t.  $\sum y_i = \sum x_i$ . Then,

$$F(x) - F^{FB} = \frac{1}{2} \frac{\text{var}(f^{i'}(x_i))}{f''(x)} \quad (102)$$

If the budget constraint is  $F^{FB} = \max_{y_i} F(y)$  s.t.  $\sum B(y) = \sum x_i$ .

**Proof.** At the FB,  $f^{i'}(x_i) = \lambda$

$$F(x) = \sum_i f^i(x_i) = \sum_i f^i(x_i^{FB}) + f^{i'}(x_i^{FB})(x_i - x_i^{FB}) + \frac{1}{2} f^{i''}(x_i^{FB})(x_i - x_i^{FB})^2$$

$$F(x) - F^{FB} = \sum_i \lambda (x_i - x_i^{FB}) + \frac{1}{2} f^{i''}(x_i^{FB})(x_i - x_i^{FB})^2$$

$$= \sum_i \frac{1}{2} f^{i''}(x_i^{FB})(x_i - x_i^{FB})^2$$

and

$$\begin{aligned} F(x) - F^{FB} &= \frac{1}{2} \sum_i \frac{1}{f^{i''}(x_i^{FB})} (f^{i'}(x_i) - f^{i'}(x_i^{FB}))^2 \text{ as } f^{i'}(x_i) - f^{i'}(x_i^{FB}) = f'' \cdot (x_i - x_i^{FB}) \\ &= \frac{1}{2} \sum_i \frac{\text{var}(f^{i'}(x_i))}{f^{i''}(x_i^{FB})} \text{ as } \sum_i f^{i'}(x_i) - f^{i'}(x_i^{FB}) = f'' \sum_i x_i - x_i^{FB} = 0 \end{aligned}$$

□

**Lemma 11.2** *Suppose that  $F(x)$  is more general, then the loss is still linked to the variance of marginal utilities:*

$$\begin{aligned} F(x) - F^{FB} &= \frac{1}{2} [F_x(x) - \lambda' I] F_{xx}^{-1} [F_x(x) - \lambda' I] \\ &= \frac{1}{2} \sum_{i,j} [F_{x_i}(x) - \lambda' I] (F_{xx}^{-1})_{i,j} [F_{x_j}(x) - \lambda' I] \end{aligned}$$

with  $\lambda' = \langle F_{x_i} \rangle$ .

**Proof.** At the FB,  $F_{x_i} = \lambda$

$$\begin{aligned} F(x) - F(x^{FB}) &= \sum_i F_{x_i}(x_i - x_i^{FB}) + \frac{1}{2} (x_i - x_i^{FB}) F_{x_i x_j}(x_j - x_j^{FB}) \\ &= 0 + \frac{1}{2} (x_i - x_i^{FB}) F_{x_i x_j}(x_j - x_j^{FB}) \\ &= \frac{1}{2} y F_{xx} y \end{aligned}$$

and  $F_x(x) - F_x(x^{FB}) = F_{xx}(x - x^{FB}) = F_{xx} y$ , so  $y = F_{xx}^{-1} \cdot (F_x(x) - F_x(x^{FB}))$ . Hence,

$$\begin{aligned} F(x) - F^{FB} &= \frac{1}{2} y F_{xx} y = \frac{1}{2} [F_x(x) - F_x(x^{FB})] F_{xx}^{-1} [F_x(x) - F_x(x^{FB})] \\ &= \frac{1}{2} [F_x(x) - \lambda' I] F_{xx}^{-1} [F_x(x) - \lambda' I] + o(\varepsilon^2) \text{ with } \lambda' = \langle F_{x_i} \rangle \\ &= \frac{1}{2} \sum_{i,j} [F_{x_i}(x) - \lambda' I] (F_{xx}^{-1})_{i,j} [F_{x_j}(x) - \lambda' I] \end{aligned}$$

which is a variance term. □

**Lemma 11.3** *Suppose we have  $f(x, y)$  and a function  $X(y) = \arg \max_x f(x, y)$ . Call  $(x^*, y^*) = \arg \max_{x,y} f(x, y)$  and  $g(y) = f(X(y), y)$ . We have:*

$$\begin{aligned} f(x, y) - f(x^*, y^*) &= [f(x, y) - f(X(y), y)] + [f(X(y), y) - f(X(y^*), y^*)] \\ &= \frac{1}{2} (x - X(y)) \cdot f_{xx}(X(y), y) \cdot (x - X(y)) + \frac{1}{2} (y - y^*) \cdot g_{yy} \cdot (y - y^*) \end{aligned}$$

Also,

$$\begin{aligned} g_y &= f_y(X(y), y) \\ g_{yy} &= f_{yy} + f_{xy}X'(y) \\ &= f_{yy} - f_{xy}f_{xx}^{-1}f_{xy} \end{aligned}$$

with  $X'(y) = -f_{xx}^{-1}f_{xy}$ .

## 11.3 Additional derivations for the Mirrlees Problem

### 11.3.1 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:

$$\bar{\zeta}_{Q_{z^*}}^c(z) = \zeta_{Q_{z^*}}^c(z) + \zeta^c(z) \delta_z(z^*), \quad (103)$$

where  $\delta_z$  is a Dirac distribution at point  $z$ . Hence, as  $\zeta_{Q_{z^*}}^c(z)$  was a potentially smooth function of  $z^*$ ,  $\bar{\zeta}_{Q_{z^*}}^c(z)$  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 11.4** (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:*

$$\delta z = \frac{\eta \delta K + z \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c \delta q_{z^*} dz^*}{q - \zeta^c z R''}, \quad (104)$$

$$\frac{\delta v}{v_r} = \delta K - z \frac{\tau^b}{q} \left( \zeta^c R'' \delta z + \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c \delta q_{z^*} dz^* \right). \quad (105)$$

In these equations, the integrals involving  $z \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c \delta q_{z^*} dz^*$  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

$$\begin{aligned} \delta z &= \frac{\eta(z) \delta K + z \zeta^c(z) \delta q_z + \int_0^\infty z \zeta_{Q_{z^*}}^c(z) \delta q_{z^*} dz^*}{q(z) - \zeta^c(z) z R''(z)}, \\ \frac{\delta v}{v_r} &= \delta K - z \frac{\tau^b}{q} \zeta^c(z) R''(z) \delta z - z \frac{\tau^b(z)}{q(z)} \zeta^c(z) \delta q_z - \int_{z^*=0}^\infty z \frac{\tau^b(z)}{q(z)} \zeta_{Q_{z^*}}^c(z) \delta q_{z^*} dz^*. \end{aligned}$$

To interpret the economics of (104), 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 (104) 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  $\frac{\zeta^c}{q}zR''\delta z$ . 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 (105), 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 11.5** *Under the conditions of the Lemma 11.4, the change in the government objective function associated with the agent is*

$$\delta L(z) = (\gamma(z) - 1) \delta K + \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} \frac{h^*(z)}{h(z)} z \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c(z) \delta q_{z^*} dz^*, \quad (106)$$

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

$$\gamma(z) = g(z) + \eta(z) \frac{\tilde{\tau}^b(z)}{1 - T'(z)} + \eta(z) \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} \frac{h^*(z)}{h(z)}. \quad (107)$$

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 (106) 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 11.6** *At the optimum,*

$$\int_0^\infty \left( 1 - \gamma(z) - \frac{T'(z) - \tilde{\tau}^b(z) h^*(z)}{1-T'(z) h(z)} z \zeta_{r_0}^c(z) \right) h(z) dz = 0. \quad (108)$$

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*

$$\frac{\partial L}{\partial \tau_{z^*}} = \int_{z^*}^\infty (1 - \gamma(z)) h(z) dz - \zeta^c(z^*) \frac{T'(z^*) - \tilde{\tau}^b(z^*)}{1-T'(z^*)} z^* h^*(z^*) - \int_0^\infty \zeta_{Q_{z^*}}^c(z) \frac{T'(z) - \tilde{\tau}^b(z)}{1-T'(z)} z h^*(z) dz. \quad (109)$$

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 (109) 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 (109) can be written in a more compact way as:

$$\frac{\partial L}{\partial \tau_{z^*}} = \int_{z^*}^\infty (1 - \gamma(z)) h(z) dz - \int_0^\infty \bar{\zeta}_{Q_{z^*}}^c(z) \frac{T'(z) - \tilde{\tau}^b(z)}{1-T'(z)} z h^*(z) dz. \quad (110)$$

### 11.3.2 Proofs for the Mirrlees results

**Notations and Derivation of relation (80)** 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 (37) gives

$$\boldsymbol{\tau}^b = (1, -q) - \frac{(u_c, u_z)}{v_r}. \quad (111)$$

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

$$\tau^b = \boldsymbol{\tau}^b(\mathbf{q}, r) \cdot (q, 1) = -\frac{qu_c + u_z}{v_r}. \quad (112)$$

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

$$\begin{aligned} \mathbf{S}_q^C(\mathbf{q}, r) &= \mathbf{c}_q - \mathbf{c}_r z = (c_q - c_r z, z_q - z_r z) \\ &= (qz_q + z - (qz_r + 1)z, z_q - z_r z) = (q, 1) (z_q - z_r z) \\ &= (q, 1) z \frac{\zeta^c}{q}. \end{aligned} \quad (113)$$

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 (83) 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

$$B_c(\mathbf{c}, \mathbf{q}) = (1, -q). \quad (114)$$



We use (43)

$$\begin{aligned}
\boldsymbol{\tau}^b &= B_{\mathbf{c}}(\mathbf{c}, \mathbf{q}) - \frac{B_{\mathbf{c}}(\mathbf{c}, \mathbf{q}^s)}{B_{\mathbf{c}}(\mathbf{c}, \mathbf{q}^s) \cdot \mathbf{c}_r(\mathbf{q}, r)} \\
&= (1, -q) - \frac{(1, -q^s)}{(1, -q^s) \cdot (qz_r + 1, z_r)} \text{ as } c(q, r) = qz(q, r) + r \text{ gives } c_r = qz_r + 1. \\
&= (1, -q) - \frac{(1, -q^s)}{1 + (q - q^s)z_r} \\
\boldsymbol{\tau}^b &= (1, -q) - \frac{(1, -q^s)}{1 + (q - q^s)\frac{\eta}{q}}. \tag{115}
\end{aligned}$$

Next, recall (112),

$$\begin{aligned}
\tau^b &= \boldsymbol{\tau}^b(\mathbf{q}, r) \cdot (q, 1) \\
&= \left[ (1, -q) - \frac{(1, -q^s)}{1 + (q - q^s)\frac{\eta}{q}} \right] \cdot (q, 1) = \frac{q^s - q}{1 + (q - q^s)\frac{\eta}{q}} \\
&= \frac{\tau - \tau^s}{1 - (\tau - \tau^s)\frac{\eta}{q}}. \tag{116}
\end{aligned}$$

using  $q = 1 - \tau$ ,  $q^s = 1 - \tau^s$ . Thus we have proven (84).

Next, we calculate  $\zeta^c$ . We call  $\mathbf{e} = (0, 1)$  the vector singling earnings on the vector  $\mathbf{c} = (c, z)$ . We apply (42) with  $p_j = q$ , the price of earnings. We have:

$$\begin{aligned}
\mathbf{e} \cdot \mathbf{S}_j^H &= \mathbf{e} \cdot \mathbf{S}^r(\mathbf{p}, r) \cdot \mathbf{p}_{p_j}^s(\mathbf{p}, r) = z_q^r \frac{\partial q^s}{\partial q} = \frac{z}{q} \zeta^{c,r} m_{zz} \\
\mathbf{e} \cdot \mathbf{S}_j^H &= \frac{z}{q} \zeta^{c,r} m_{zz}. \tag{117}
\end{aligned}$$

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

$$\begin{aligned}
D_j &= -\boldsymbol{\tau}^b \cdot \mathbf{S}_j^H \text{ by (36)} \\
&= [B_{\mathbf{c}}(\mathbf{p}, \mathbf{c}) - B_{\mathbf{c}}(\mathbf{p}^s, \mathbf{c})] \cdot \mathbf{S}_j^H \text{ by (44)} \\
&= [(1, q) - (1, q^s)] \cdot \mathbf{S}_j^H \text{ by (114)} \\
&= (q - q^s) \mathbf{e} \cdot \mathbf{S}_j^H \text{ as } \mathbf{e} = (0, 1) \\
&= (q - q^s) \frac{z}{q} \zeta^{c,r} m_{zz} \text{ by (117)}.
\end{aligned}$$

We record:

$$D_j = (q - q^s) \frac{z}{q} \zeta^{c,r} m_{zz}. \tag{118}$$

Next, we apply (40):  $\mathbf{S}_j^C = \mathbf{S}_j^H + \mathbf{c}_r D_j$ , 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. (113)):

$$\begin{aligned} \zeta^c &= \frac{q}{z} \mathbf{e} \cdot \mathbf{S}_j^C \\ \zeta^c &= \left( 1 - \eta \frac{\tau - \tau^s}{q} \right) \zeta^{c,r} m_{zz}. \end{aligned} \tag{119}$$

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

$$\zeta_{Q_{z^*}}^c = \left( 1 - \eta \frac{\tau - \tau^s}{q} \right) \zeta^{c,r} m_{zz^*}. \tag{120}$$

Hence, we have proven (83).

**Proof of (85): Decision vs. Experienced utility model** The agent's optimization gives  $q u_c^s + u_z^s = 0$ . Equation (112) 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 Proposition 10.1** Let us now solve for the optimal  $J$ , which ensures  $\frac{\partial L}{\partial Q_{z^*}} = 0$ . We can write

$$- \frac{\partial L}{\partial \tau_{z^*}} = J(z^*) - \int_{z^*}^{\infty} a(z) dz - b(z^*) + \int_{z^*}^{\infty} J(z) \rho(z) dz. \tag{121}$$

We use the notations

$$\rho(z) = \frac{\eta}{\zeta^c} \frac{1}{z}, \quad (122)$$

$$a(z) = (1 - g(z)) h(z) - \rho g(z) z \left( -\tau^b(z) \frac{\zeta^c}{q(z)} \right) (h^*(z) - h(z)), \quad (123)$$

$$b(z) = -\overline{\overline{F}}(z_*). \quad (124)$$

$a(z)$  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.

$\overline{\overline{F}}(z^*)$  is the part impact on the government's objective function of increase  $\delta q_{z^*}$ , coming from the distortions from perceptions

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

We note that

$$\begin{aligned} a &= (1 - \gamma) h + \eta \frac{T'(z)}{q} \frac{h^*}{h} h, \\ a &= (1 - \gamma) h + \rho J. \end{aligned} \quad (125)$$

We also have

$$\begin{aligned} J(z^*) &= \int_{z^*}^\infty a(z) dz + b(z^*) - \int_{z^*}^\infty J(z) \rho(z) dz, \\ \dot{J} &= -a + \dot{b} + J\rho, \\ \frac{d}{dz} \left[ J(z) e^{-\int_0^z \rho(s) ds} \right] &= e^{-\int_0^z \rho(s) ds} \left( -a(z) + \dot{b}(z) \right), \\ J(z) e^{-\int_0^z \rho(s) ds} &= C + \int_z^\infty e^{-\int_0^{z'} \rho(s) ds} \left( a(z') - \dot{b}(z') \right) dz', \\ J(z) &= \int_z^\infty e^{-\int_z^{z'} \rho(s) ds} \left( a(z') - \dot{b}(z') \right) dz'. \end{aligned} \quad (126)$$

Integrating by parts, we get

$$\begin{aligned} \int_z^\infty e^{-\int_z^{z'} \rho(s) ds} b(z') dz' &= \left[ e^{-\int_z^{z'} \rho(s) ds} b(z') \right]_z^\infty + \int_z^\infty e^{-\int_z^{z'} \rho(s) ds} \rho(z') b(z') dz' \\ &= -b(z) + \int_z^\infty e^{-\int_z^{z'} \rho(s) ds} \rho(z') b(z') dz', \end{aligned} \quad (127)$$

$$J(z) = b(z) + \int_z^\infty e^{-\int_z^{z'} \rho(s) ds} (a(z') - \rho(z') b(z')) dz'. \quad (128)$$

We can rewrite this as

$$\begin{aligned} J(z^*) &= -\overline{F}(z^*) + \int_{z^*}^\infty e^{-\int_{z^*}^z \rho(s) ds} \left( a(z) + \rho \overline{F}(z) \right) dz \\ &= -z^* g(z^*) \left( -\tau^b(z^*) \frac{\zeta^c(z^*)}{q(z^*)} \right) h^*(z^*) - F(z^*) \\ &\quad + \int_{z^*}^\infty e^{-\int_{z^*}^z \rho(s) ds} \left[ (1 - g(z)) h(z) + g(z) \rho(z) z \left( -\tau^b(z) \frac{\zeta^c(z)}{q(z)} \right) h(z) + \rho(z) F(z) \right] dz. \end{aligned} \quad (129)$$

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

$$\begin{aligned} \frac{T'(z^*) - \tilde{\tau}^b(z^*)}{1 - T'(z^*)} + \frac{1}{\zeta^c z^* h^*(z^*)} \int_{z=0}^\infty \left( \zeta_{Q_{z^*}}^c(z) - \int_{z'=z^*}^\infty e^{-\int_{z^*}^{z'} \rho(s) ds} \rho(z') \zeta_{Q_{z'}}^c(z) dz' \right) \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} z h^*(z) dz, \\ = \frac{1}{\zeta^c} \frac{1 - H(z^*)}{z^* h^*(z^*)} \int_{z^*}^\infty e^{-\int_{z^*}^z \rho(s) ds} \left( 1 - g(z) - \eta \frac{\tilde{\tau}^b(z)}{q(z)} \right) \frac{h(z)}{1 - H(z^*)} dz. \end{aligned}$$

□

**Proof of Proposition 10.2** We use (11.1):

$$\begin{aligned}
\frac{\partial L}{\partial \tau_{z^*}} &= \int_{z^*}^{\infty} (1 - \gamma(z)) h(z) dz - \zeta^c(z^*) \frac{T'(z^*) - \tilde{\tau}^b(z^*)}{1 - T'(z^*)} z^* h^*(z^*) - \int_0^{\infty} \zeta_{Q_{z^*}}^c(z) \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} z h^*(z) dz \\
&= \int_{z^*}^{\infty} \left( 1 - g(z) - \eta(z) \frac{T'(z)}{1 - T'(z)} \right) h(z) dz - \zeta^c(z^*) \frac{T'(z^*) - \tilde{\tau}^b(z^*)}{1 - T'(z^*)} z^* h^*(z^*) \\
&\quad - \int_0^{\infty} \zeta_{Q_{z^*}}^c(z) \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} z h^*(z) dz \text{ using (86)} \\
&= (1 - H(z_*)) \left( 1 - \bar{g} - \bar{\eta} \frac{\bar{\tau}}{1 - \bar{\tau}} \right) - \bar{\zeta}^c \frac{\bar{\tau} - \bar{g}\bar{\tau}^b}{1 - \bar{\tau}} z^* h^*(z^*) - \int_0^{\infty} \zeta_{Q_{z^*}}^c(z) \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} z h^*(z) dz.
\end{aligned}$$

Recall that

$$\begin{aligned}
\int_{z_0}^{\infty} (1 - H(z_*)) dz_* &= [(1 - H(z_*))(z - z_*)]_{z_0}^{\infty} - \int_{z_0}^{\infty} h(z_*) (z - z_*) dz \\
&= \mathbb{E}[(z - z_0) \mathbf{1}_{z \geq z_0}] \\
&= \frac{1}{\pi} \mathbb{E}[z \mathbf{1}_{z \geq z_0}].
\end{aligned}$$

Given the constraint that  $T'(z) = \bar{\tau}$  for  $z > z_0$ , the FOC on  $\bar{\tau}$  is:  $0 = \int_{z_0}^{\infty} \frac{\partial L}{\partial \tau_{z^*}} dz_*$ , i.e.

$$\begin{aligned}
0 &= \int_{z_0}^{\infty} \frac{\partial L}{\partial \tau_{z^*}} dz_* = \left( 1 - \bar{g} - \bar{\eta} \frac{\bar{\tau}}{1 - \bar{\tau}} \right) \int_{z_0}^{\infty} (1 - H(z_*)) dz_* - \bar{\zeta}^c \frac{\bar{\tau} - \bar{g}\bar{\tau}^b}{1 - \bar{\tau}} \mathbb{E}[z \mathbf{1}_{z \geq z_0}] \\
&\quad - \int_0^{\infty} \left( \int_{z_0}^{\infty} \zeta_{Q_{z^*}}^c(z) dz_* \right) \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} z h^*(z) dz \\
&= \left( 1 - \bar{g} - \bar{\eta} \frac{\bar{\tau}}{1 - \bar{\tau}} \right) \frac{1}{\pi} \mathbb{E}[z \mathbf{1}_{z \geq z_0}] - \bar{\zeta}^c \frac{\bar{\tau} - \bar{g}\bar{\tau}^b}{1 - \bar{\tau}} \mathbb{E}[z \mathbf{1}_{z \geq z_0}] - \int_0^{\infty} \zeta_{\bar{q}}^c(z) \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} z h^*(z) dz \\
0 &= \left( 1 - \bar{g} - \bar{\eta} \frac{\bar{\tau}}{1 - \bar{\tau}} \right) - \bar{\zeta}^c \frac{\bar{\tau} - \bar{g}\bar{\tau}^b}{1 - \bar{\tau}} - \frac{\pi}{\mathbb{E}[z \mathbf{1}_{z \geq z_0}]} \int_0^{\infty} \zeta_{\bar{q}}^c(z) \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} z h^*(z) dz.
\end{aligned}$$

Hence, we have:

$$\frac{\bar{\tau}}{1 - \bar{\tau}} = \frac{1 - \bar{g} - \beta + \zeta^c \pi \bar{g} \frac{\bar{\tau}^b}{1 - \bar{\tau}}}{\zeta^c \pi + \bar{\eta}}, \tag{130}$$

with

$$\begin{aligned}
\beta &= \frac{\pi}{\mathbb{E}[z \mathbf{1}_{z \geq z_0}]} \int_0^{\infty} \zeta_{\bar{q}}^c(z) \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} z h^*(z) dz \\
&= \mathbb{E}^z \left[ \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} \right] \mathbb{E}^* [\zeta_{\bar{q}}^c(z)] \frac{\pi \mathbb{E}[z]}{\mathbb{E}[z \mathbf{1}_{z \geq z_0}]}.
\end{aligned}$$

where  $\mathbb{E}^z[\phi(z)] = \frac{\int \phi(z) h(z) z dz}{\int \phi(z) z dz}$  and  $\mathbb{E}^*[\zeta_{\bar{q}}^c(z)] = \frac{\mathbb{E}^*[\zeta_{\bar{q}}^c(z) \frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} h^*(z) dz]}{\mathbb{E}^*[\frac{T'(z) - \tilde{\tau}^b(z)}{1 - T'(z)} h^*(z) dz]}$  is a weighted average too.

We can rewrite the equation as :  $\frac{\bar{\tau}}{1-\bar{\tau}} = \frac{a+\frac{b}{c}}{c}$ , which gives  $\bar{\tau} = \frac{a+b}{a+c}$ , so that:

$$\bar{\tau} = \frac{1 - \bar{g} - \beta + \zeta^c \pi \bar{g} \tau^b}{1 - \bar{g} - \beta + \zeta^c \pi + \bar{\eta}}.$$

□

**Proof of Lemma 11.4** 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 11.5** 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

$$z R'' = q \frac{h^* - h}{\zeta^c h^*}. \quad (131)$$

We have

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

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 11.4, we can rewrite this as

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

Using equation (131) and Lemma 11.4, we can rewrite this as

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

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:

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

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 11.5 gives the change in the government objective function

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

The total change is

$$\begin{aligned}\delta L &= \int_0^\infty \delta L(z) h(z) dz, \\ \frac{\delta L}{\delta q_{z^*}} &= \int_{z^*}^\infty (\gamma(z) - 1) h(z) dz + \int_0^\infty \left[ \frac{T'(z) - \tilde{\tau}^b(z)}{q} \bar{\bar{\zeta}}_{Q_{z^*}}^c \right] \frac{h^*}{h} z h(z) dz, \\ \frac{\partial L}{\partial \tau_{z^*}} &= -\bar{\bar{\bar{F}}}(z_*) + \int_{z^*}^\infty (1 - \gamma(z)) h(z) dz.\end{aligned}\tag{132}$$

We also have

$$\frac{\partial L}{\partial \tau_{z^*}} \equiv -\frac{\partial L}{\partial Q_{z^*}} = -\bar{\bar{\bar{F}}}(z_*) + \int_{z^*}^\infty (1 - \gamma(z)) h(z) dz.\tag{133}$$

□

**Proof of Lemma 11.6** Using Lemma 11.5, 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. □



## 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:

$$v_{p_j} + v_w c_j = 0. \quad (134)$$

Shepard's lemma is proven as follows: The envelope theorem gives  $e_p(\mathbf{p}, u) = \mathbf{h}(\mathbf{p}, u)$ , i.e.

$$e_{p_i}(\mathbf{p}, u) = h^i, \quad (135)$$

and differentiating once more gives:

$$e_{p_i p_j} = h_{p_j}^i(\mathbf{p}, u) = S_{ij}. \quad (136)$$

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.

$$c_{p_j}^i + c_w^i c_j = h_{p_j}^i = S_{ij}. \quad (137)$$

### 12.2 Behavioral version with perceived prices

The sparse max demand

$$\operatorname{smax}_{\mathbf{c}|\mathbf{p}^s} u(\mathbf{c}) \text{ s.t. } \mathbf{p} \cdot \mathbf{c} \leq w$$

of a behavioral agent perceiving prices  $\mathbf{p}^s$  (while true prices are  $\mathbf{p}$  and the true budget is  $w$ ) is:

$$\mathbf{c}(\mathbf{p}, \mathbf{p}^s, w) = \mathbf{c}^r(\mathbf{p}^s, w'(\mathbf{p}, \mathbf{p}^s, w)), \quad (138)$$

where perceived budget  $w'$  satisfies:

$$\mathbf{p} \cdot \mathbf{c}^r(\mathbf{p}^s, w'(\mathbf{p}, \mathbf{p}^s, w)) = w. \quad (139)$$

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.** 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:*

$$\mathbf{c}_{p_j} = -\mathbf{c}_w c_j \tag{140}$$

$$\mathbf{c}_{p_j^s} = \mathbf{S}_j^r + \mathbf{c}_w D_j^s =: S \tag{141}$$

$$= \mathbf{S}_j^r + \mathbf{c}_w [(\mathbf{p}^s - \mathbf{p}) \cdot \mathbf{S}_j^r] \tag{142}$$

$$= (1 + \mathbf{c}_w (\mathbf{p}^s - \mathbf{p})') \mathbf{S}_j^r. \tag{143}$$

*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,

$$d\mathbf{c} = \mathbf{c}_{p^s} dp^s = S dp^s + \mathbf{c}_w dE$$

$$dE = [(\mathbf{p}^s - \mathbf{p}) \cdot \mathbf{S}^r \cdot dp^s] \geq 0.$$

The extra term  $dE$  is positive: it's as if the agent became richer. That creates an income effect, and shifts his 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

$$\mathbf{c}_w = \mathbf{h}_u^r v_w, \quad \mathbf{c}_{p_j} = \mathbf{h}_u^r v_{p_j}. \tag{144}$$

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 (139): 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  $\mathbf{c}_w^s = \mathbf{c}_{w'}^r \frac{\partial w'}{\partial w}$ .  $\square$

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  $\mathbf{p} \cdot \mathbf{c}_{w'}^r \frac{\partial w'}{\partial w} = 1$ , 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:

$$\begin{aligned} v_w &= u'(\mathbf{c}^s) \cdot \mathbf{c}_w = \lambda \mathbf{p}^s \cdot \mathbf{c}_w = \lambda \mathbf{p}^s \cdot \left( \frac{\mathbf{c}_{w'}^r(\mathbf{p}^s, w')}{\mathbf{p} \cdot \mathbf{c}_{w'}^r(\mathbf{p}^s, w')} \right) = \lambda \frac{\mathbf{p}^s \cdot \mathbf{c}_{w'}^r(\mathbf{p}^s, w')}{\mathbf{p} \cdot \mathbf{c}_{w'}^r(\mathbf{p}^s, w')} \\ v_w &= \frac{\lambda}{\mathbf{p} \cdot \mathbf{c}_{w'}^r(\mathbf{p}^s, w')}. \end{aligned}$$

$\square$

We can check that things are consistent: with  $u_{\mathbf{c}} = \lambda \mathbf{p}^s$ ,

$$v_{p_j^s} = u_{\mathbf{c}} \mathbf{c}_{p_j^s} = \lambda \mathbf{p}^s (1 + \mathbf{c}_w p_j') \mathbf{S}_j^r = \lambda \mathbf{p}^s \mathbf{c}_w p_j' \mathbf{S}_j^r = v_w p_j' \mathbf{S}_j^r = v_w D_j^s.$$

**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,

$$e_{p_i p_j}^s = 0 \tag{145}$$

$$e_{p_i p_j^s}^s = h_{p_j^s}^i(\mathbf{p}^s, u) = S_{ij}^r, \tag{146}$$

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:

$$\mathbf{p}^s(\mathbf{p}, w) = \frac{u_{\mathbf{c}}(\mathbf{c}(\mathbf{p}, w))}{v_w(\mathbf{p}, w)}. \tag{147}$$

Then, the demand function can be represented as that of a sparse agent with perceived prices  $\mathbf{p}^s(\mathbf{p}, w)$ .

$$\mathbf{c}(\mathbf{p}, w) = \mathbf{c}^s(\mathbf{p}, \mathbf{p}^s(\mathbf{p}, w), w). \tag{148}$$

**Proof.** The demand of an agent misperceiving prices,  $\mathbf{c}^s(\mathbf{p}, \mathbf{p}^s, w)$ , 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.  $\square$

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,  $S^H(\mathbf{q}, w)$  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

Here we give 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.

$$v_{p_i} = -v_w B_{p_i}. \quad (149)$$

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:

$$\mathbf{h}_{p_i} = \mathbf{c}_{p_i} + \mathbf{c}_w B_{p_i}. \quad (150)$$

Finally, given  $B(\mathbf{c}(\mathbf{p}, w), \mathbf{p}) = w$ ,

$$B_{\mathbf{c}}\mathbf{c}_{p_i} = -B_{p_i}, B_{\mathbf{c}}\mathbf{c}_w = 1. \quad (151)$$

Premultiplying (150) by  $B_{\mathbf{c}}$  gives:  $B_{\mathbf{c}}\mathbf{h}_{p_i} = B_{\mathbf{c}}\mathbf{c}_{p_i} + B_{\mathbf{c}}\mathbf{c}_w B_{p_i} = -B_{p_i} + B_{p_i} = 0$ ,

$$B_{\mathbf{c}}\mathbf{h}_{p_i} = 0. \quad (152)$$

All in all, traditional consumer theory holds, replacing  $p$  by  $B_{\mathbf{c}}$ .

## 13.2 Misperceiving Agent

We now study

$$v(\mathbf{p}, \mathbf{p}^s, w) = \operatorname{smax}_{\mathbf{c}|\mathbf{p}^s} u(\mathbf{c}) \text{ s.t. } B(\mathbf{c}, \mathbf{p}) \leq w. \quad (153)$$

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

$$\mathbf{h}^r(\mathbf{p}^s, \bar{u}) = \arg \min_{\mathbf{c}} B(\mathbf{c}, \mathbf{p}^s) \text{ s.t. } u(\mathbf{c}) \geq \bar{u}, \quad (154)$$

and the corresponding Slutsky matrix, and the perceived  $\mathbf{p}^s$

$$\mathbf{S}_j^r = \mathbf{h}_{\mathbf{p}_j^s}^r(\mathbf{p}^s, \bar{u})|_{\bar{u}=v(\mathbf{p}, \mathbf{p}^s, w)}. \quad (155)$$

Define the dual expenditure function

$$e(\mathbf{p}, \mathbf{p}^s, \bar{u}) = \underset{\mathbf{c} | \mathbf{p}^s}{\text{sm}} B(\mathbf{c}, \mathbf{p}) \text{ s.t. } u(\mathbf{c}) \geq \bar{u}. \quad (156)$$

We have the simpler representation:

$$e(\mathbf{p}, \mathbf{p}^s, \bar{u}) = B(\mathbf{h}^r(\mathbf{p}^s, \bar{u}), \mathbf{p}). \quad (157)$$

**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

$$e_{\mathbf{p}_j^s} = B_{\mathbf{c}} \mathbf{h}_{\mathbf{p}_j^s}^r = (B_{\mathbf{c}} - B_{\mathbf{c}}^s) \mathbf{h}_{\mathbf{p}_j^s}^r.$$

as (152) gives  $B_{\mathbf{c}}^s \mathbf{h}_{\mathbf{p}_j^s}^r = 0$ .  $\square$

**Proposition 13.2** (Generalized Roy's identity). *Given the function  $v(\mathbf{p}, \mathbf{p}^s, w)$ , we have:*

$$\begin{aligned} \frac{v_{\mathbf{p}_j}}{v_w} &= -B_{\mathbf{p}_j} \\ \frac{v_{\mathbf{p}_j^s}}{v_w} &= (B_{\mathbf{c}}^s - B_{\mathbf{c}}) \cdot \mathbf{S}_j^r = D_j^s. \end{aligned}$$

*i.e.*  $\frac{v_{\mathbf{p}_j^s}}{v_w^s} = \sum_i (B_i^s - B_i) \mathbf{S}_{ij}^r$ .

**Proof.** 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:

$$0 = v_{\mathbf{p}_j^s} + v_w e_{\mathbf{p}_j^s} = v_{\mathbf{p}_j^s} + v_w (B_{\mathbf{c}} - B_{\mathbf{c}}^s) \cdot \mathbf{S}_j^r$$

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:

$$\mathbf{c}_{\mathbf{p}_j} = -\mathbf{c}_w B_{\mathbf{p}_j} \quad (158)$$

$$\mathbf{c}_{\mathbf{p}_j^s} = \mathbf{S}_j^r + \mathbf{c}_w D_j^s, \quad (159)$$

i.e.  $\mathbf{c}_{\mathbf{p}_j}^i = -B_{\mathbf{p}_j}^i \cdot \mathbf{c}_w$  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

$$\mathbf{c}_w = \mathbf{h}_u v_w, \quad \mathbf{c}_{\mathbf{p}_j} = \mathbf{h}_u v_{\mathbf{p}_j}. \quad (160)$$

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,  $\mathbf{c}_{\mathbf{p}_j^s}(\mathbf{p}, \mathbf{p}^s, w) = \mathbf{h}(\mathbf{p}^s, v(\mathbf{p}, \mathbf{p}^s, w))$  gives:

$$\begin{aligned} \mathbf{c}(\mathbf{p}, \mathbf{p}^s, w) &= \mathbf{h}_{\mathbf{p}_j^s} + \mathbf{h}_u v_{\mathbf{p}_j^s} \\ &= \mathbf{S}_j^r + \mathbf{c}_w \frac{v_{\mathbf{p}_j^s}}{v_w} \text{ using } \mathbf{c}_w = \mathbf{h}_u v_w \\ &= \mathbf{S}_j^r + \mathbf{c}_w [(B_{\mathbf{c}}^s - B_{\mathbf{c}}) \cdot \mathbf{S}_j^r] \text{ using Proposition 13.2.} \end{aligned}$$

□

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:

$$\frac{v_w}{\lambda} = B_{\mathbf{c}}^s \mathbf{c}_w = 1 + (B_{\mathbf{c}}^s - B_{\mathbf{c}}) \mathbf{c}_w, \quad (161)$$

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  $\mathbf{h}^{r,s}(\mathbf{p}^s, \hat{u}) = \arg \min_{\mathbf{c}} B(\mathbf{p}^s, \mathbf{c})$  s.t.  $u^s(\mathbf{c}) \geq \hat{u}$  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:

$$\mathbf{c}(\mathbf{p}, w) = \mathbf{h}^{r,s}(\mathbf{p}^s(\mathbf{p}, w), v^s(\mathbf{p}, w)). \quad (162)$$

**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*

$$\mathbf{S}_j^s(\mathbf{p}, w) = \mathbf{S}^{r,s}(\mathbf{p}, w) \cdot \mathbf{p}_{p_j}^s(\mathbf{p}, w). \quad (163)$$

*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:

$$\bar{\tau}^b = (B_{\mathbf{c}}(\mathbf{p}, \mathbf{c}) - B_{\mathbf{c}}(\mathbf{p}^s, \mathbf{c})) + \left( \frac{u_{\mathbf{c}}^s}{v_w^s} - \frac{u_{\mathbf{c}}}{v_w} \right). \quad (164)$$

Finally,  $B_{\mathbf{c}}(\mathbf{p}^s, \mathbf{c}) \cdot \mathbf{S}_j^s = 0$ .

This tax  $\bar{\tau}^b$  is the sum of two gaps: between the prices and perceived prices ( $B_{\mathbf{c}}(\mathbf{c}, \mathbf{p}) - B_{\mathbf{c}}(\mathbf{p}^s, \mathbf{c})$ ), and between true utility and perceived utility ( $\frac{u_{\mathbf{c}}^s}{v_w^s} - \frac{u_{\mathbf{c}}}{v_w}$ ).

**Proof.**

Use  $\mathbf{c}(\mathbf{p}, w) = \mathbf{h}^{r,s}(\mathbf{p}^s, v^s(\mathbf{p}, w))$ :

$$\begin{aligned}\mathbf{c}_w(\mathbf{p}, w) &= \mathbf{h}_u^{r,s} v_w^s + \mathbf{h}_{\mathbf{p}^s}^{r,s} \mathbf{p}_w^s \\ \mathbf{c}_{p_j}(\mathbf{p}, w) &= \mathbf{h}_{\mathbf{p}^s}^{r,s} \mathbf{p}_{p_j}^s + \mathbf{h}_u^{r,s} v_j^s = \mathbf{S}_j^s + \mathbf{c}_w \frac{v_j^s}{v_w^s}.\end{aligned}$$

Using (32) and (33) 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

$$B_{\mathbf{c}}(\mathbf{p}^s, \mathbf{c}) \cdot \mathbf{S}_j^s = B_{\mathbf{c}}(\mathbf{p}^s, \mathbf{c}) \cdot \mathbf{h}_{\mathbf{p}^s}^r(\mathbf{p}^s, v^s) \mathbf{p}_{p_j}^s = 0 \text{ as } B_{\mathbf{c}}(\mathbf{p}^s, \mathbf{c}) \mathbf{h}_{\mathbf{p}^s}^r(\mathbf{p}^s, v^s) = 0.$$

Recall also that  $\frac{u_{\mathbf{c}}^s}{v_w^s} = \Lambda B_{\mathbf{c}}(\mathbf{p}^s, \mathbf{c})$  as the agent maximizes with perceived prices  $\mathbf{p}^s$ . Hence,

$$\begin{aligned}-D_j &= \left( B_{\mathbf{c}}(\mathbf{c}, \mathbf{p}) - \frac{u_{\mathbf{c}}}{v_w} \right) \mathbf{S}_j^s \\ &= \left( (B_{\mathbf{c}}(\mathbf{c}, \mathbf{p}) - B_{\mathbf{c}}(\mathbf{p}^s, \mathbf{c})) - \left( \frac{u_{\mathbf{c}}}{v_w} - \Lambda B_{\mathbf{c}}(\mathbf{p}^s, \mathbf{c}) \right) \right) \mathbf{S}_j^s \text{ as } B_{\mathbf{c}}(\mathbf{p}^s, \mathbf{c}) \cdot \mathbf{S}_j^s = 0 \\ &= \left( (B_{\mathbf{c}}(\mathbf{c}, \mathbf{p}) - B_{\mathbf{c}}(\mathbf{p}^s, \mathbf{c})) - \left( \frac{u_{\mathbf{c}}}{v_w} - \frac{u_{\mathbf{c}}^s}{v_w^s} \right) \right) \cdot \mathbf{S}_j^s = \bar{\tau}_b \cdot \mathbf{S}_j^H.\end{aligned}$$

□