

Online Appendix

Optimal Internality Taxation of Product Attributes

Andreas Gerster and Michael Kramm

A Proofs

A.A Proof of Proposition 1: Optimal Non-Linear Tax

We first use the mechanism design approach to characterize the optimal tax in terms of types and then employ the perturbation approach to characterize the optimal tax as a function of the attribute level. Finally, we show that the results we obtain under either approach are equivalent.

A.A.1 Equation (2) in Types (Mechanism Design Approach)

For our mechanism design approach, we employ the Revelation Principle for dominant-strategy implementation (Gibbard, 1973) to solve for a direct mechanism, where the consumer truthfully reveals information about his perceived valuation.

We start by discussing the dimensionality of mechanisms for internality taxation. In our model, a consumer decides based on his perceived valuation $\hat{v}(v, b)$ rather than v . Even though a consumer may or may not know his bias, i.e., be sophisticated or naive, we can, without loss of generality, neglect that distinction and restrict our analysis to one-dimensional mechanisms in the perceived valuation \hat{v} . Naive consumers are unaware of their bias and thus cannot report it, so that a social planner can only employ a one-dimensional mechanism in \hat{v} to correct them. Sophisticated consumers know their bias and can, in principle, report it. Yet, as biases do not influence decision utility, truth-telling in a two-dimensional mechanism is not incentive-compatible.

The formal argument of the above is as follows. We want to show that two-dimensional mechanisms, where sophisticated consumers report both their perceived valuation and their bias, lead to a violation of truth-telling. Without loss of generality, we assume that there exists at least one realization of perceived valuations $\hat{v}_1 \in [\underline{\hat{v}}, \bar{\hat{v}}]$ for which biases differ, so that some consumers are characterized by (\hat{v}_1, b_1) and others by (\hat{v}_1, b_0) , where $b_1 \neq b_0$.

The policy maker wants to implement a direct two-dimensional mechanism where the allocation $\zeta(\tilde{v}, \tilde{b})$ and the tax $\tau(\tilde{v}, \tilde{b})$ depend on reported valuations \tilde{v} and biases \tilde{b} . Consumers choose their reports to maximize decision utility:

$$(\tilde{v}^*(\hat{v}), \tilde{b}^*(\hat{v})) = \arg \max_{\tilde{v}, \tilde{b}} u^d(\zeta(\tilde{v}, \tilde{b}), \tau(\tilde{v}, \tilde{b}) | \hat{v}).$$

Importantly, this maximization problem depends only on reported biases \tilde{b} and not on actual biases b . As a consequence, every sophisticated consumer with perceived valuation \hat{v} will report the same bias $\tilde{b}^*(\hat{v})$. As biases differ for \hat{v}_1 , truth-telling is violated.

Hence, we can apply the Revelation Principle to our setting, where the space of reports for a consumer is given by the space of his perceived valuation \hat{v} . For simplicity, we refer to \hat{v} as a consumer's type in the following. The policy maker confines herself to designing a direct mechanism $(\zeta, \tau) : [\underline{\hat{v}}, \bar{\hat{v}}] \rightarrow Q \times \mathbb{R}$ under truth-telling in order to implement the welfare maximizing outcome. Based on the consumer's strategic report \tilde{v} , the allocation rule of the direct mechanism assigns the consumed attribute level, $\zeta(\tilde{v}) \in Q$, and transfer rule the amount of taxes to be paid, $\tau(\tilde{v}) \in \mathbb{R}$. Because the consumer has unit demand for the good, participation constraints are not relevant in our setting.

Under the direct mechanism, the decision utility for report \tilde{v} for a given perceived valuation \hat{v} is:

$$u^d(\zeta(\tilde{v}), \tau(\tilde{v}) | \hat{v}) = m + \hat{v} \cdot \zeta(\tilde{v}) - \tau(\tilde{v}) - c(\zeta(\tilde{v})).$$

Since the consumer may strategically misreport his perceived valuation, truth-telling must be ensured by implementing an incentive compatible mechanism. This implies that the tax schedule must satisfy:

$$u^d(\zeta(\hat{v}), \tau(\hat{v}) | \hat{v}) \geq u^d(\zeta(\tilde{v}), \tau(\tilde{v}) | \hat{v}) \quad \forall \hat{v}, \tilde{v} \in [\underline{\hat{v}}, \bar{\hat{v}}]. \quad (\text{IC})$$

Optimal strategic reporting of a consumer implies that the solution v^* to the problem $\max_{\tilde{v}} u^d(\zeta(\tilde{v}), \tau(\tilde{v}) | \hat{v})$ has to satisfy:

$$\hat{v} \zeta'(v^*) - \tau'(v^*) - \zeta'(v^*) c'(\zeta(v^*)) \stackrel{!}{=} 0. \quad (8)$$

As incentive compatibility requires that $v^* = \hat{v}$, equilibrium decision utility in an incentive-compatible direct mechanism is given by $\hat{u}^d(\hat{v}) := u^d(\zeta(\hat{v}), \tau(\hat{v}) | \hat{v})$, while equilibrium normative utility is given by $\hat{u}^n(\hat{v}, b) := u^n(\zeta(\hat{v}), \tau(\hat{v}) | v) = \hat{u}^d(\hat{v}) - b \zeta(\hat{v})$.

Put differently, incentive compatibility implies that, for all $\hat{v} \in [\underline{\hat{v}}, \bar{\hat{v}}]$, the following equation has to hold:

$$\frac{\partial \hat{u}^d(\hat{v})}{\partial \hat{v}} = \zeta(\hat{v}) + \hat{v} \zeta'(\hat{v}) - \tau'(\hat{v}) - \zeta'(\hat{v}) c'(\zeta(\hat{v})) \stackrel{(8)}{=} \zeta(\hat{v}). \quad (9)$$

To determine the optimal tax schedule, the policy maker solves an optimization problem which can be analyzed using an optimal control approach. Note that determining the equilibrium values of $\zeta(\hat{v})$ and $\hat{u}^d(\hat{v})$ for all \hat{v} pins down the equilibrium value of $\tau(\hat{v})$ for all \hat{v} . Hence, the mechanism design problem of the policy maker is given by:

$$\max_{\zeta \in \mathcal{Q}} \int_{\hat{v}} \alpha(\hat{v}) \cdot E[\hat{u}^n(\hat{v}, b) | \hat{v}] dF(\hat{v}) + \lambda \left(\int_{\hat{v}} \tau(\hat{v}) dF(\hat{v}) - B \right), \quad (10)$$

subject to the condition from Equation (9), where $\mathcal{Q} := \{f | f : [\underline{\hat{v}}, \bar{\hat{v}}] \rightarrow Q\}$ is the function space containing all functions with domain $[\underline{\hat{v}}, \bar{\hat{v}}]$ and codomain Q . The boundary conditions of the problem are given by $\hat{u}^d(\underline{\hat{v}}) = \underline{u}$ and $\hat{u}^d(\bar{\hat{v}}) \geq \underline{u}$. The control variable is ζ and the law of motion of the state variable \hat{u}^d is determined by incentive compatibility and optimal strategic reporting, as given by Equation (9).

Using the definition of decision utility to replace the tax and rewriting equilibrium normative utility in terms of equilibrium decision utility, the Hamiltonian for the problem stated in Equation (10) for all $\hat{v} \in [\underline{\hat{v}}, \bar{\hat{v}}]$ is given by:

$$\mathfrak{H}(\hat{v}, \zeta, \hat{u}^d) = \left[\alpha(\hat{v}) \cdot \underbrace{\left(\hat{u}^d(\hat{v}) - E[b | \hat{v}] \zeta(\hat{v}) \right)}_{=E[\hat{u}^n(\hat{v}, b) | \hat{v}]} + \lambda \underbrace{\left(m + \hat{v} \zeta(\hat{v}) - \hat{u}^d(\hat{v}) - c(\zeta(\hat{v})) \right)}_{=\tau(\hat{v})} \right] f(\hat{v}) + \mu(\hat{v}) \zeta(\hat{v}).$$

Following the standard solution procedure for such mechanism design problems, we employ Pontryagin's Maximum Principle, which yields the following necessary conditions for the optimal tax.²⁰

$$\text{FOC on control: } \frac{\partial \mathfrak{H}}{\partial \zeta} = [-E[b | \hat{v}] \cdot \alpha(\hat{v}) + \lambda (\hat{v} - c'(\cdot))] f(\hat{v}) + \mu(\hat{v}) \stackrel{!}{=} 0, \quad (\text{FOC}_\zeta)$$

$$\text{FOC on state: } \frac{\partial \mathfrak{H}}{\partial \hat{u}^d} = [\alpha(\hat{v}) - \lambda] f(\hat{v}) \stackrel{!}{=} -\mu'(\hat{v}), \quad (\text{FOC}_u)$$

$$\text{Transversality cond.: } \mu(\underline{\hat{v}}) \cdot \hat{u}^d(\underline{\hat{v}}) = \mu(\bar{\hat{v}}) \cdot \hat{u}^d(\bar{\hat{v}}) = 0. \quad (\text{TVC})$$

²⁰In addition, sufficiency is given if the control region is convex and the Hamiltonian is concave in (ζ, \hat{u}^d) for every \hat{v} . Both conditions are satisfied in our setup.

The consumer's first-order condition characterizing optimal consumption q^d is given by

$$\left. \frac{\partial u^d(q, t, \hat{v})}{\partial q} \right|_{q=q^d} = \hat{v} - c'(q^d) - t'(q^d) \stackrel{!}{=} 0 \Leftrightarrow c'(q^d) = \hat{v} - t'(q^d). \quad (11)$$

The second order condition is satisfied if $c''(q) + t''(q) \geq 0$ for all $q \in Q$. Since the costs are convex in q by assumption, this condition is satisfied if the optimal tax schedule is convex in q as well. We find that convex optimal tax schedules are optimal for many behavioral biases (see discussion after Proposition 4). More generally, this condition also holds if the optimal tax schedule is concave. In that case, the cost function may not be "too concave" in the sense that $c''(q) + t''(q) \geq 0$ holds.

In our setting, we abstract from participation constraints. Hence, we can w.l.o.g. assume that $\hat{u}(\hat{v}) = \underline{u} > 0$ and $\hat{u}(\bar{v}) \geq \underline{u}$, so that the transversality condition immediately implies $\mu(\hat{v}) = 0$ and $\mu(\bar{v}) = 0$. Integrating Equation (FOC_u) and using $\mu(\bar{v}) = 0$, we obtain

$$\int_{\hat{v}}^{\bar{v}} -\mu'(n)dn = -\mu(\bar{v}) - [-\mu(\hat{v})] = \mu(\hat{v}) \stackrel{!}{=} \int_{\hat{v}}^{\bar{v}} [\alpha(m) - \lambda] f(m)dm. \quad (12)$$

Using the Equations (11) and (12), we rearrange Equation (FOC_q), to obtain the result:

$$\begin{aligned} \lambda(\hat{v} - c'(\cdot)) &\stackrel{!}{=} -\frac{\mu(\hat{v})}{f(\hat{v})} + E[b|\hat{v}] \cdot \alpha(\hat{v}) \\ \Leftrightarrow (11) \quad t'(q) &= -\frac{\mu(\hat{v}_q)}{\lambda f(\hat{v}_q)} + E[b|\hat{v}_q] \cdot \frac{\alpha(\hat{v}_q)}{\lambda} \\ \Leftrightarrow (12) \quad t'(q) &= \frac{\int_{\hat{v}_q}^{\bar{v}} [1 - \hat{g}(m)] f(m)dm}{f(\hat{v}_q)} + E[b|\hat{v}_q] \cdot \hat{g}(\hat{v}_q). \end{aligned}$$

The average marginal social welfare weight for consumers in the upper part of the distribution (in comparison to \hat{v}) is defined as:

$$\hat{G}(\hat{v}) := \frac{\int_{\hat{v}}^{\bar{v}} \hat{g}(m)dF(m)}{1 - F(\hat{v})} = E[g(\theta)|\theta \geq \hat{v}]. \quad (13)$$

Using this we get

$$\begin{aligned} t'(q_{\hat{v}}) &= \hat{g}(\hat{v})E[b|\hat{v}] + \frac{\int_{\hat{v}}^{\bar{v}} [1 - \hat{g}(n)] f(n)dn}{f(\hat{v})} \\ &= \hat{g}(\hat{v})E[b|\hat{v}] + \frac{\int_{\hat{v}}^{\bar{v}} f(n)dn - \int_{\hat{v}}^{\bar{v}} \hat{g}(n)f(n)dn}{f(\hat{v})} \end{aligned}$$

$$\begin{aligned}
&\stackrel{(13)}{=} \hat{g}(\hat{\vartheta})E[b|\hat{\vartheta}] + \frac{[1 - F(\hat{\vartheta})] - \hat{G}(\hat{\vartheta}) [1 - F(\hat{\vartheta})]}{f(\hat{\vartheta})} \\
&= \hat{g}(\hat{\vartheta})E[b|\hat{\vartheta}] + \frac{[1 - F(\hat{\vartheta})] [1 - \hat{G}(\hat{\vartheta})]}{f(\hat{\vartheta})}
\end{aligned}$$

A.A.2 Equation (3) in Attribute Level q (Perturbation Approach)

The perturbation approach uses insights from a local perturbation of the tax schedule to derive the optimal tax. In particular, it examines a change in the slope of the consumer's budget set in a small band $(q, q + dq)$ from $1 - c'(q) - t'(q)$ to $1 - c'(q) - t'(q) - d\tau$. Such a change has four effects on welfare. The mechanical effect represents an increase in collected tax money, since all consumers with consumption levels higher than $q + dq$ are affected by the increased tax rate. This effect is captured by

$$d\tau dq [1 - H(q)].$$

The consumer welfare effect captures the social value of the decrease in consumption of the numeraire due to the local tax perturbation:

$$-d\tau dq [1 - H(q)] G(q).$$

Since we abstract from income effects, there is no change of the consumption of q for these consumers. This is different for the consumers inside the small band $(q, q + dq)$. Here, relative prices are changed and thus consumers change their consumption of the attribute level by

$$\delta q = \frac{dq}{d(c' + t')} d\tau = -e \frac{q}{c' + t'} d\tau. \quad (14)$$

where the elasticity is defined as

$$e := -\frac{dq}{d(t' + c')} \frac{t' + c'}{q}. \quad (15)$$

We define the elasticity at the net price $t' + c'$, since the slope of the consumer's budget constraint is determined by $t' + c'$.

The fiscal externality captures the decrease in collected taxes

$$\delta q t'(q) h(q) dq,$$

while the bias correction effect

$$-g(q)\delta q E [b|q] h(q) dq$$

captures the change in welfare due to the induced change in consumer utility, which is caused by the potential correction of an externality.

In optimum, a local perturbation of the tax may not change social welfare. Hence, the four effects have to sum to zero, so that

$$\begin{aligned} 0 &= d\tau dq [1 - H(q)] [1 - G(q)] + \delta q [t'(q) - E [b|q] g(q)] h(q) dq \\ \Leftrightarrow t'(q) &= g(q) E [b|q] - \frac{d\tau [1 - H(q)] [1 - G(q)]}{\delta q h(q)} \\ &\stackrel{(14)}{=} g(q) E [b|q] - \frac{d\tau [1 - H(q)] [1 - G(q)]}{-e(q) \frac{q}{t'(q) + c'(q)} d\tau h(q)} \\ &= g(q) E [b|q] + \frac{1 - G(q)}{e(q) a(q)} (t'(q) + c'(q)). \end{aligned}$$

A.A.3 Equivalence of Equations (2) and (3)

We first linearize the budget constraint of the consumer in the tax so that, using the virtual income R , it can be rewritten as $z = R - c(q) - \tau q$. Then, the decision utility is given by $u = \hat{v}q + R - c(q) - \tau q$.

The goal is now to write the change of consumption q in type \hat{v} using elasticity e , which is the elasticity of q with respect to net-price $\tau + c'$. The first-order condition of the consumer with respect to consumption gives

$$0 \stackrel{!}{=} \hat{v} - \tau - c'(q). \quad (16)$$

Applying the implicit function theorem twice to (16) yields

$$\begin{aligned} \frac{dq}{d\hat{v}} &= -\frac{1}{-c''(q)} = \frac{1}{c''(q)} \\ \frac{dq}{d(\tau + c')} &= -\frac{-1}{-c''(q)} = -\frac{1}{c''(q)} \\ \Rightarrow \frac{dq}{d\hat{v}} &= -\frac{dq}{d(\tau + c')}. \end{aligned} \quad (17)$$

We use the above results and the definition of the elasticity from (15):

$$e = -\frac{dq}{d(\tau + c')} \frac{\tau + c'}{q}$$

$$\begin{aligned}
& \stackrel{(17)}{=} \frac{dq}{d\hat{\vartheta}} \frac{\tau + c'}{q} \\
& \Leftrightarrow \frac{dq}{d\hat{\vartheta}} = \frac{qe}{\tau + c'}.
\end{aligned} \tag{18}$$

For density $h(q)$ the following must hold

$$\begin{aligned}
& h(q_{\hat{\vartheta}})dq_{\hat{\vartheta}} = f(\hat{\vartheta})d\hat{\vartheta} \\
& \Leftrightarrow h(q_{\hat{\vartheta}})\frac{dq(\hat{\vartheta})}{d\hat{\vartheta}} = f(\hat{\vartheta}) \\
& \stackrel{(18)}{\Leftrightarrow} h(q_{\hat{\vartheta}})\frac{qe}{t' + c'} = f(\hat{\vartheta}).
\end{aligned} \tag{19}$$

The average marginal social welfare weight for consumers in the upper part of the distribution (in comparison to q) can be rewritten using a change of variables (Ch.o.V.):

$$\hat{G}(\hat{\vartheta}) := \frac{\int_{\hat{\vartheta}}^{\bar{\vartheta}} \hat{g}(m)dF(m)}{1 - F(\hat{\vartheta})} \stackrel{\text{Ch.o.V.}}{=} \frac{\int_{q_{\hat{\vartheta}}}^{\bar{q}} g(m)dH(q_m)}{1 - H(q_{\hat{\vartheta}})} =: G(q_{\hat{\vartheta}}). \tag{20}$$

The thinness of the top tail (individuals with consumption above q in relation to those at q) is given by

$$a(q) := \frac{h(q)q}{1 - H(q)}. \tag{21}$$

Using the above results we get

$$\begin{aligned}
t'(q_{\hat{\vartheta}}) &= \hat{g}(\hat{\vartheta})E[b|\hat{\vartheta}] + \frac{[1 - F(\hat{\vartheta})][1 - \hat{G}(\hat{\vartheta})]}{f(\hat{\vartheta})} \\
& \stackrel{(19), \text{C.o.V.}}{=} g(q)E[b|q] + \frac{[1 - H(q)][1 - G(q)]}{h(q)\frac{qe(q)}{t'(q) + c'(q)}} \\
&= g(q)E[b|q] + \frac{(t'(q) + c'(q))[1 - H(q)][1 - G(q)]}{h(q)qe(q)} \\
& \Leftrightarrow t'(q) \stackrel{(21)}{=} g(q)E[b|q] + \frac{1 - G(q)}{a(q)e(q)} (t'(q) + c'(q)).
\end{aligned}$$

A.B Derivation of the Optimal Linear Tax

The policy maker's problem of setting the optimal linear tax can be written as

$$\max_{t \in \mathbb{R}} \int_v \int_b \alpha(v, b) u^n(q^d, t, v) dF_b(b|v) dF_v(v) + \lambda \left(\int_v \int_b t \cdot q^d dF_b(b|v) dF_v(v) - B \right) =: V(t).$$

We evaluate the derivative with respect to the linear tax t :

$$\begin{aligned}\frac{\partial V(t)}{\partial t} &= \int_v \int_b \alpha(v,b) \left[-q^d + (v-t-c'(\cdot)) \frac{\partial q^d}{\partial t} \right] dF_b(b|v) dF_v(v) + \lambda \int_v \int_b \left[q^d + t \cdot \frac{\partial q^d}{\partial t} \right] dF_b(b|v) dF_v(v) \\ &= \int_v \int_b \left[\alpha(v,b) (v-c'(\cdot)) \frac{\partial q^d}{\partial t} - (\alpha(v,b) - \lambda) (q^d + t \frac{\partial q^d}{\partial t}) \right] dF_b(b|v) dF_v(v).\end{aligned}$$

The individually optimal consumption is again characterized by Equation (11), i.e., $c'(\cdot) = \hat{v} - t'(q) = (v+b) - t$, where the last equality holds since t is linear. Thus,

$$\begin{aligned}\frac{\partial V}{\partial t} &= \int_v \int_b \left[\alpha(v,b) (t-b) \frac{\partial q^d}{\partial t} - (\alpha(v,b) - \lambda) (q^d + t \frac{\partial q^d}{\partial t}) \right] dF_b(b|v) dF_v(v) \\ &= \int_v \int_b \left[(\lambda t - \alpha(v,b)b) \frac{\partial q^d}{\partial t} - (\alpha(v,b) - \lambda) q^d \right] dF_b(b|v) dF_v(v).\end{aligned}$$

Using that t is constant, we can rewrite the equation as follows:

$$\frac{\partial V}{\partial t} = \lambda t \frac{\partial \bar{q}^d}{\partial t} - \int_v \int_b \left[\alpha(v,b)b \frac{\partial q^d}{\partial t} + (\alpha(v,b) - \lambda) q^d \right] dF_b(b|v) dF_v(v),$$

where the change in average demand \bar{q}^d in response to a tax increase is given by $\frac{\partial \bar{q}^d}{\partial t} = \int_v \int_b \left[\frac{\partial q^d}{\partial t} \right] dF_b(b|v) dF_v(v)$. The optimal tax t^* is given by $\frac{\partial V}{\partial t} |_{t=t^*} \stackrel{!}{=} 0$, which gives:

$$\frac{\partial V}{\partial t} = t \frac{\partial \bar{q}^d}{\partial t} - \int_v \int_b \left[\hat{g}(v,b)b \frac{\partial q^d}{\partial t} - (1 - \hat{g}(v,b)) q^d \right] dF_b(b|v) dF_v(v) = 0$$

Abstracting from redistributive motives ($\hat{g}(v,b) = 1$ for all (v,b)) implies

$$t^* = \int_v \int_b b \left[\frac{\partial q^d}{\partial t} / \frac{\partial \bar{q}^d}{\partial t} \right] dF_b(b|v) dF_v(v), \quad (22)$$

where $\frac{\partial q^d}{\partial t} / \frac{\partial \bar{q}^d}{\partial t}$ denotes the relative responsiveness of a consumer type (v,b) , i.e., the change in demand for that consumer type in response to a tax increase, relative to change in total demand.

If $c'''(\cdot) = 0$, the optimal linear tax simplifies to $t^* = E[b]$. To see this, differentiate Equation (11) with respect to t , which yields $\frac{\partial q^d}{\partial t} = -\frac{1}{c''(q^d)}$. This term is constant if $c'''(\cdot) = 0$, so that $\partial q^d / \partial t = \partial \bar{q}^d / \partial t$. Hence, the optimal linear tax is $t^* = E[b]$.

A.C Proof of Proposition 2: Condition for Implementability

Definition 2 (Internality Tax Implementability). *The allocation function $\xi : [\hat{\vartheta}, \bar{\vartheta}] \rightarrow Q$ is internality-tax implementable if there exists a tax function $\tau : [\hat{\vartheta}, \bar{\vartheta}] \rightarrow \mathbb{R}$ such that $\{(\xi(\hat{\vartheta}), \tau(\hat{\vartheta})) | \hat{\vartheta} \in [\hat{\vartheta}, \bar{\vartheta}]\}$ satisfy incentive compatibility according to*

$$u^d(\xi(\hat{\vartheta}), \tau(\hat{\vartheta}) | \hat{\vartheta}) \geq u^d(\xi(\bar{\vartheta}), \tau(\bar{\vartheta}) | \hat{\vartheta}) \quad \forall \hat{\vartheta}, \bar{\vartheta} \in [\hat{\vartheta}, \bar{\vartheta}].$$

Implementability hinges on two necessary conditions, which concern the *consumer preferences* or directly stem from them. First, the consumer's utility function must satisfy a single-crossing condition. Second, ξ must be monotonic in $\hat{\vartheta}$. Since, in our setting, the single-crossing condition is satisfied via $\partial (u_q^d / u_t^d) / \partial \hat{\vartheta} < 0$, a necessary condition for implementability is that the allocation is non-decreasing in perceived types, i.e., $\partial \xi / \partial \hat{\vartheta} \geq 0$ (Proof: See Appendix A.C.1).

Implementability of an incentive compatible mechanism additionally hinges on conditions, which stem from the *policy maker's preferences*. In our setting, these involve (paternalistic) corrective and redistributive motives. In a setting without corrective motives, a sufficient condition for this requirement involves the hazard rate of F , that is, $f(\hat{\vartheta}) / (1 - F(\hat{\vartheta}))$, which is a measure of the thinness of the tail of the distribution. For our setting, Proposition 2 shows that the condition is more restrictive involving terms that stem from the corrective motive (Proof: See Appendix A.C.2).

A.C.1 Proof: Single-Crossing and Monotonicity

Incentive compatibility requires

$$\hat{\vartheta} = \arg \max_{\bar{\vartheta}} u^d(\xi(\bar{\vartheta}), \tau(\bar{\vartheta}), \hat{\vartheta}) \quad \forall \bar{\vartheta} \in \hat{V}$$

The first-order condition implies

$$\begin{aligned} \frac{\partial u^d(\xi(\bar{\vartheta}), \tau(\bar{\vartheta}), \hat{\vartheta})}{\partial \bar{\vartheta}} \Big|_{\bar{\vartheta}=\hat{\vartheta}} &= \frac{\partial u^d(\bar{\vartheta}, \hat{\vartheta})}{\partial \bar{\vartheta}} \Big|_{\bar{\vartheta}=\hat{\vartheta}} = 0 \\ \Leftrightarrow \frac{\partial u^d(\bar{\vartheta}, \hat{\vartheta})}{\partial \xi} \cdot \frac{\partial \xi(\bar{\vartheta})}{\partial \bar{\vartheta}} + \frac{\partial u^d(\bar{\vartheta}, \hat{\vartheta})}{\partial \tau} \cdot \frac{\partial \tau(\bar{\vartheta})}{\partial \bar{\vartheta}} &= 0 \\ \Leftrightarrow \frac{\partial \tau(\bar{\vartheta})}{\partial \bar{\vartheta}} &= - \frac{\partial u^d(\bar{\vartheta}, \hat{\vartheta}) / \partial \xi}{\partial u^d(\bar{\vartheta}, \hat{\vartheta}) / \partial \tau} \cdot \frac{\partial \xi(\bar{\vartheta})}{\partial \bar{\vartheta}}. \end{aligned} \quad (23)$$

Differentiating the first-order condition with respect to $\hat{\vartheta}$ yields

$$\frac{\partial^2 u^d(\bar{\vartheta}, \hat{\vartheta})}{\partial \bar{\vartheta}^2} \Big|_{\bar{\vartheta}=\hat{\vartheta}} + \frac{\partial^2 u^d(\bar{\vartheta}, \hat{\vartheta})}{\partial \bar{\vartheta} \partial \hat{\vartheta}} \Big|_{\bar{\vartheta}=\hat{\vartheta}} = 0. \quad (24)$$

The second-order condition implies that

$$\left. \frac{\partial^2 u^d(\tilde{v}, \hat{v})}{\partial \tilde{v}^2} \right|_{\tilde{v}=\hat{v}} \leq 0. \quad (25)$$

Equations (25) and (24) imply

$$\begin{aligned} & \left. \frac{\partial u^d(\tilde{v}, \hat{v})}{\partial \tilde{v} \partial \hat{v}} \right|_{\tilde{v}=\hat{v}} \geq 0 \\ \Leftrightarrow & \frac{\partial \left(\frac{\partial u^d(\tilde{v}, \hat{v})}{\partial \xi} \right)}{\partial \hat{v}} \cdot \frac{\partial \xi(\tilde{v})}{\partial \tilde{v}} + \frac{\partial \left(\frac{\partial u^d(\tilde{v}, \hat{v})}{\partial \tau} \right)}{\partial \hat{v}} \cdot \frac{\partial \tau(\tilde{v})}{\partial \tilde{v}} \geq 0 \\ \stackrel{(23)}{\Leftrightarrow} & \frac{\partial \left(\frac{\partial u^d}{\partial \xi} \right)}{\partial \hat{v}} \cdot \frac{\partial \xi(\tilde{v})}{\partial \tilde{v}} - \frac{\partial \left(\frac{\partial u^d}{\partial \tau} \right)}{\partial \hat{v}} \cdot \frac{\partial u^d / \partial \xi}{\partial u^d / \partial \tau} \cdot \frac{\partial \xi(\tilde{v})}{\partial \tilde{v}} \geq 0 \\ \Leftrightarrow & \frac{\frac{\partial \left(\frac{\partial u^d}{\partial \xi} \right)}{\partial \hat{v}} \partial u^d / \partial \tau - \frac{\partial \left(\frac{\partial u^d}{\partial \tau} \right)}{\partial \hat{v}} \cdot \partial u^d / \partial \xi}{\partial u^d / \partial \tau} \cdot \frac{\partial \xi(\tilde{v})}{\partial \tilde{v}} \geq 0 \\ \Leftrightarrow & \frac{\frac{\partial \left(\frac{\partial u^d}{\partial \xi} \right)}{\partial \hat{v}} \partial u^d / \partial \tau - \frac{\partial \left(\frac{\partial u^d}{\partial \tau} \right)}{\partial \hat{v}} \cdot \partial u^d / \partial \xi}{(\partial u^d / \partial \tau)^2} \cdot \frac{\partial u^d}{\partial \tau} \cdot \frac{\partial \xi(\tilde{v})}{\partial \tilde{v}} \geq 0 \\ \Leftrightarrow & \frac{\partial \left(\frac{\partial u^d / \partial \xi}{\partial u^d / \partial \tau} \right)}{\partial \hat{v}} \cdot \frac{\partial u^d}{\partial \tau} \cdot \frac{\partial \xi(\tilde{v})}{\partial \tilde{v}} \geq 0, \end{aligned}$$

where the second-but-last step follows from the quotient rule. Since we know that

$\frac{\partial u^d}{\partial \tau} < 0$ and via single crossing $\frac{\partial \left(\frac{\partial u^d / \partial \xi}{\partial u^d / \partial \tau} \right)}{\partial \hat{v}} = \frac{\partial \left(\frac{\hat{v} - c' - t'}{-1} \right)}{\partial \hat{v}} = -1 < 0$, it follows that we need monotonicity of the form $\frac{\partial \xi(\tilde{v})}{\partial \tilde{v}} \geq 0$.

A.C.2 Proof: Internality Tax Implementability

Inserting the optimal smooth tax from Proposition 1, the consumer's first-order condition can be rewritten as

$$c'[\xi(\hat{v})] \stackrel{!}{=} \hat{v} - \underbrace{\hat{g}(\hat{v}) E[b|\hat{v}] - (1 - \hat{G}(\hat{v})) \frac{1-F(\hat{v})}{f(\hat{v})}}_{=: \phi(\hat{v})}. \quad (26)$$

An incentive compatible mechanism must guarantee that ξ is non-decreasing in \hat{v} (see Appendix A.C.1). Since c is convex, ξ is non-decreasing in \hat{v} if and only if the left-hand side of Equation (26) is non-decreasing in \hat{v} . Accordingly, the right-hand side of Equation (26) must be non-decreasing in \hat{v} as well, which implies

$$1 - \frac{\partial \hat{g}(\hat{v})}{\partial \hat{v}} E[b|\hat{v}] - \frac{\partial E[b|\hat{v}]}{\partial \hat{v}} \hat{g}(\hat{v}) - \frac{\partial [(1 - \hat{G}(\hat{v})) \cdot \{1 - F(\hat{v})\} / f(\hat{v})]}{\partial \hat{v}} \geq 0. \quad (27)$$

A.D Proof of Proposition 3: Bunching

In this section, we show that a failure of Internality Tax implementability leads to bunching and discuss how bunching at the top or at the bottom is equivalent to a ban for high or low attribute levels. To determine the optimal policy when bunching occurs, we apply the approach by Guesnerie and Laffont (1984) to conduct “ironing” using optimal control theory.

Before we present the formal method of ironing, we illustrate the intuition of this procedure and explain why it implies different forms of standards (uniform, maximum, and minimum). For the ease of exposition we abstract from redistributive motives ($g = 1 = G$). Suppose the smooth tax derived in Proposition 1 implies that the allocation rule ζ is decreasing in the interval $[\hat{\vartheta}^h, \hat{\vartheta}^l]$, where $\hat{\vartheta}^h$ denotes the level of perceived valuation that leads to a higher ζ , while $\hat{\vartheta}^l$ leads to a lower ζ ($\hat{\vartheta}^h < \hat{\vartheta}^l$, but $\zeta(\hat{\vartheta}^h) > \zeta(\hat{\vartheta}^l)$). The goal of ironing is to determine over which interval $[a, b]$ of the type space bunching will take place.

Recall that an allocation rule that decreases in $\hat{\vartheta}$ implies a fundamental misalignment of preferences violating the condition in Proposition 2: the policy maker would like to allocate a lower allocation level to a consumer with a higher perceived valuation. However, incentive compatibility requires a non-decreasing allocation rule (see Appendix A.C.1): a consumer will not reveal a higher perceived valuation, if he obtains a lower attribute level than a consumer with a lower perceived valuation.

To ensure incentive compatibility, the policy maker resorts to bunching and assigns the same product attribute $\zeta^* = \hat{\zeta}$ to a subset of consumers. Suppose to the contrary that the implemented allocation rule ζ were partly increasing in the interval $[\hat{\vartheta}^h, \hat{\vartheta}^l]$. Then, the welfare of types that underconsume compared to ζ can be increased by increasing the allocation, while the welfare of types that overconsume can be increased by decreasing the allocation. This is true up to the point where the consumption of all types $\hat{\vartheta} \in [a, b]$ is equal to $\hat{\zeta}$.

The standard $\hat{\zeta}$ is set such that it minimizes the loss in welfare compared to the solution candidate ζ . Intuitively, the standard $\hat{\zeta}$ lies in between the extremes of ζ^l and ζ^h so that some consumers overconsume compared to candidate ζ and some underconsume. If all consumers underconsume, i.e., $\hat{\zeta} = \zeta^l$, then increasing the standard marginally would raise the welfare of all consumers that still underconsume after the increase, while it would only decrease the welfare of the consumers who overconsume after the increase, that is, for types with an ideal consumption $\zeta = \zeta^l$. An analogous statement holds for the case when all consumers overconsume, i.e., $\hat{\zeta} = \zeta^h$.

The exact procedure of determining the types that over- or underconsume compared to ζ is described below for the general case of interior bunching, which sub-

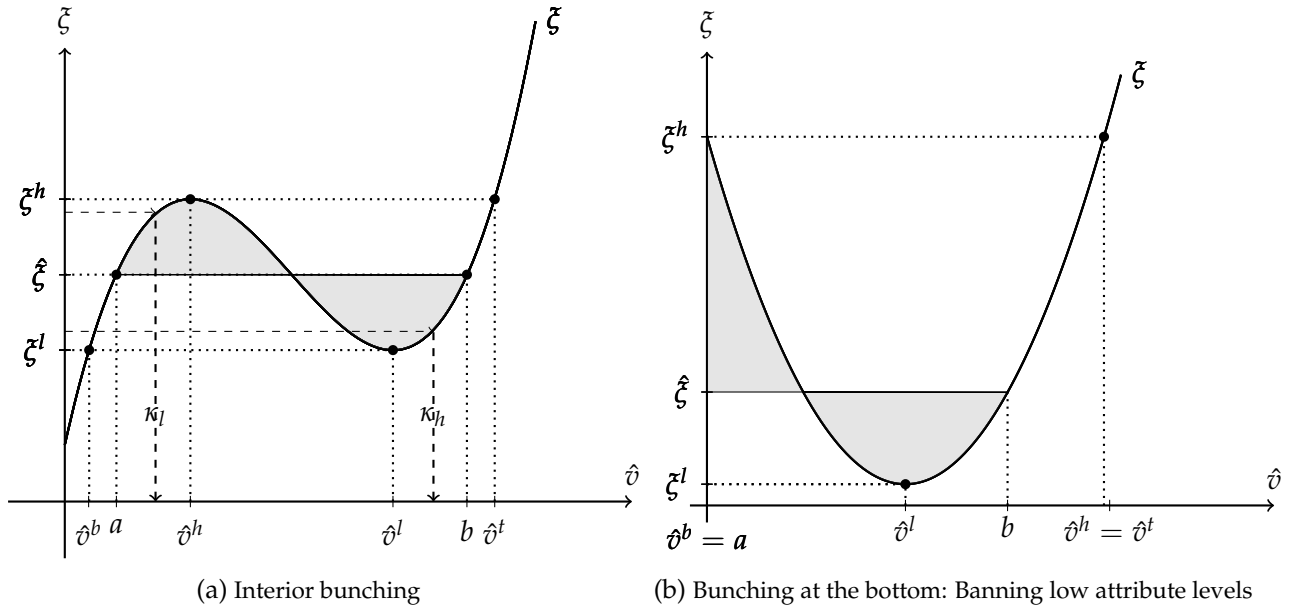


Figure 1: Ironing yields the optimal bunching regions

sumes all cases relevant to our scenario (uniform, minimum and maximum standards). Figure 1 b) visualizes the specific case of minimum standards, i.e., bunching at the bottom.

We now present the formal method of ironing. Bunching occurs if ζ (or equivalently, the term $\phi(\hat{\vartheta})$ as defined in Equation (26)) is decreasing over an interval $[\hat{\vartheta}^h, \hat{\vartheta}^l] \subseteq [\hat{\vartheta}, \bar{\vartheta}]$. Let ζ denote the non-ironed “solution candidate” and let ζ^* denote the correct solution involving ironing. Since we cannot guarantee that the allocation is strictly increasing, write

$$v(\hat{\vartheta}) := \frac{\partial \zeta}{\partial \hat{\vartheta}} \geq 0. \quad (28)$$

We now have a control problem with an inequality constraint on the control v . The Hamiltonian $\mathfrak{H}(\hat{\vartheta}, \zeta, v, \delta)$ of the general optimization problem without implicitly assuming that the incentive constraints hold, can be written as

$$\mathfrak{H}(\cdot) = \left[\alpha(\hat{\vartheta}) \cdot \underbrace{\left(\hat{u}^d(\hat{\vartheta}) - E[b|\hat{\vartheta}]\zeta(\hat{\vartheta}) \right)}_{=E[\hat{u}^n(\hat{\vartheta}, b)|\hat{\vartheta}]} + \lambda \underbrace{\left(m + \hat{\vartheta}\zeta(\hat{\vartheta}) - \hat{u}^d(\hat{\vartheta}) - c(\zeta(\hat{\vartheta})) \right)}_{=\tau(\hat{\vartheta})} \right] f(\hat{\vartheta}) + \delta(\hat{\vartheta})v(\hat{\vartheta}),$$

where δ denotes the multiplier for the constraint given by Inequality (28). The first-order condition on the state ζ yields

$$\frac{-\delta'(\hat{\vartheta})}{\lambda} \stackrel{!}{=} [\hat{\vartheta} - \hat{g}(\hat{\vartheta})E[b|\hat{\vartheta}] - c'(\zeta^*(\hat{\vartheta}))] f(\hat{\vartheta}). \quad (29)$$

Without redistributive motives ($\hat{G} \equiv 1$), the term $\phi(\hat{v})$ as defined in Equation (26) is given by $\phi(\hat{v}) = \hat{v} - \hat{g}(\hat{v})E[b|\hat{v}]$. For the ease of exposition, we discuss the case without redistributive motives. Thus, Equation (29) can be rewritten as

$$\frac{-\delta'(\hat{v})}{\lambda} \stackrel{!}{=} [\phi(\hat{v}) - c'(\zeta^*(\hat{v}))] f(\hat{v}). \quad (30)$$

For simplicity, let us assume there is one unique interval $[a, b]$ of bunching.²¹ Then, integrating (30) between a and b , and using the transversality conditions $\delta(a) = \delta(b) = 0$, which hold since the monotonicity constraint is non-binding at the boundaries, yields

$$\int_a^b [\phi(\theta) - c'(\zeta^*(\theta))] f(\theta) d\theta \stackrel{!}{=} \frac{1}{\lambda} \int_a^b -\delta'(\theta) f(\theta) d\theta = \frac{-\delta(b) + \delta(a)}{\lambda} = 0, \quad (31)$$

which implies that the average difference between the ϕ and the marginal costs (i.e., the average distortion of the "virtual surplus") is zero over the bunching interval.²² Equation (31) together with the fact that $\zeta^*(a) = \zeta^*(b) = \zeta(a) = \zeta(b)$ characterizes the allocation $\zeta^*(\hat{v}) = \hat{\zeta}$ given to types $\hat{v} \in [a, b]$ in the optimal mechanism.

We now use the above characterization of the bunching region to determine its boundaries a and b and the allocation $\hat{\zeta}$. The procedure is illustrated in Figure 1. We first consider a case where the bunching region is in the middle of the range of \hat{v} . If $a > \underline{\hat{v}}$ and $b < \bar{\hat{v}}$, then there exists $\hat{v}^l := \arg \min_{\hat{v}} \zeta(\hat{v})$ s.t. $\hat{v} \in [a, b]$ and $\hat{v}^h := \arg \max_{\hat{v}} \zeta(\hat{v})$ s.t. $\hat{v} \in [a, b]$. Denote $\zeta^h := \zeta(\hat{v}^h)$ and $\zeta^l := \zeta(\hat{v}^l)$. Let κ_l denote the inverse function of the increasing part of ζ defined on $[\hat{v}^b, \hat{v}^h]$ with $\hat{v}^b := \min_{\hat{v}} \hat{v}$ s.t. $\zeta(\hat{v}) = \zeta^l$. Analogously, κ_h denotes the inverse function of the increasing part of ζ defined on $[\hat{v}^l, \hat{v}^t]$ with $\hat{v}^t := \max_{\hat{v}} \hat{v}$ s.t. $\zeta(\hat{v}) = \zeta^h$. For $\tilde{\zeta} \in [\zeta^h, \zeta^l]$, define

$$\Delta(\tilde{\zeta}) := \int_{\kappa_l(\tilde{\zeta})}^{\kappa_h(\tilde{\zeta})} [\phi(\theta) - c'(\tilde{\zeta})] f(\theta) d\theta. \quad (32)$$

Since $\zeta^h > \zeta(\hat{v})$, $\forall \hat{v} \in (\hat{v}^h, \hat{v}^t)$, it holds that $\Delta(\zeta^h) < 0$. Analogously, since $\zeta^l < \zeta(\hat{v})$, $\forall \hat{v} \in (\hat{v}^b, \hat{v}^l)$, it holds that $\Delta(\zeta^l) > 0$. Therefore, by the intermediate value theorem, there must exist some $\hat{\zeta}$, such that $\Delta(\hat{\zeta}) = 0$ as required by Equation (31). Thus, $a = \kappa_l(\hat{\zeta})$ and $b = \kappa_h(\hat{\zeta})$. As a consequence, all individuals with $\hat{v} \in [a, b]$ will be assigned the same level of the product attribute $\hat{\zeta}$.

²¹The extension to a setting with several bunching regions is only subject to minor caveats. See Guesnerie and Laffont (1984) for details.

²²Note that in a mechanism which optimally does not involve bunching, the difference between ϕ and the marginal costs is zero at each point as can be seen in Equation (26).

If $a = \hat{v}$ (analogous reasoning applies for $b = \bar{\hat{v}}$), bunching occurs at the bottom of the type distribution (see Panel b of Figure 1). In that case, we need to define and evaluate

$$\Delta(\tilde{\xi}) := \int_{\hat{v}}^{\kappa_h(\tilde{\xi})} [\phi(\theta) - c'(\tilde{\xi})] f(\theta) d\theta. \quad (33)$$

As a consequence, all individuals with $\hat{v} \in [\underline{\hat{v}}, b]$ will be assigned the same level of the product attribute $\hat{\xi}$. As bunching occurs at the bottom, this outcome corresponds to a minimum standard, i.e., a ban of all realizations of the product attribute below $\hat{\xi}$. The rationale is equivalent for bunching at the top.

A.E Proof of Equation (5): First-Order Approximation

As a consequence of the Regression Conditional Expectation Function Theorem (e.g., Angrist and Pischke, 2009), the Minimum Mean Squared Error (MMSE) linear approximation of the conditional expectation $E[b|\hat{v}]$ is given by:

$$\hat{E}[b|\hat{v}] = E(b|\hat{v} = \mu_{\hat{v}}) + \frac{\text{cov}(b, \hat{v})}{\sigma_{\hat{v}}^2} \cdot [\hat{v} - \mu_{\hat{v}}]. \quad (34)$$

Using that $\text{cov}(b, \hat{v}) = \text{cov}(b, v) + \sigma_b^2$, $\sigma_{\hat{v}}^2 = \sigma_v^2 + \sigma_b^2 + 2\text{cov}(b, v)$, and $\text{cov}(b, v) = \rho\sigma_v\sigma_b$, we obtain after rearranging:

$$\hat{E}[b|\hat{v}] = E[b|\mu_{\hat{v}}] + \frac{\rho + (\sigma_b/\sigma_v)}{(\sigma_b/\sigma_v) + (\sigma_v/\sigma_b) + 2\rho} \cdot [\hat{v} - \mu_{\hat{v}}]. \quad (35)$$

The existence of the conditional expectation as a function of \hat{v} is guaranteed by the Factorization Lemma and the Radon-Nikodym theorem.

In the bivariate normal case, (v, b) is jointly normal distributed with:

$$(v, b) \sim N \left(\begin{bmatrix} \mu_v \\ \mu_b \end{bmatrix}, \begin{bmatrix} \sigma_v^2 & \rho\sigma_v\sigma_b \\ \rho\sigma_v\sigma_b & \sigma_b^2 \end{bmatrix} \right). \quad (36)$$

Hence, the perceived valuation $\hat{v} = v + b$ has the following normal distribution

$$\hat{v} \sim N(\mu_v + \mu_b, \sigma_v^2 + \sigma_b^2 + 2\rho\sigma_v\sigma_b) =: (\mu_{\hat{v}}, \sigma_{\hat{v}}^2). \quad (37)$$

The conditional expectation $E[b|\hat{v}]$ of the bias can be calculated as

$$E[b|\hat{v}] = \frac{\text{cov}(b, \hat{v})}{\sigma_{\hat{v}}^2} \cdot \hat{v} + \left[\mu_b - \frac{\text{cov}(b, \hat{v})}{\sigma_{\hat{v}}^2} \cdot \mu_{\hat{v}} \right] = \mu_b + \frac{\text{cov}(b, \hat{v})}{\sigma_{\hat{v}}^2} \cdot [\hat{v} - \mu_{\hat{v}}], \quad (38)$$

with $cov(b, \hat{v}) = cov(b, v) + \sigma_b^2$. Rearranging gives that

$$E[b|\hat{v}] = \mu_b + \frac{\rho + (\sigma_b/\sigma_v)}{(\sigma_b/\sigma_v) + (\sigma_v/\sigma_b) + 2\rho} \cdot [\hat{v} - \mu_{\hat{v}}]. \quad (39)$$

Hence, the approximation error $R(\hat{v})$ equals zero for all \hat{v} if v and b follow a bivariate normal distribution.

A.F Proof of Proposition 6: Prices Vs. Quantities

The proof proceeds in two steps. In a first step, we explore the welfare gains from non-linear taxation relative to linear taxes and to standards. The corresponding results are summarized in Proposition 8. In a second step, we then use these results to compare linear taxes with standards. We analyze a setting in which the local bias heterogeneity A is constant.

Proposition 8 (Welfare Gain from Non-Linear Taxation). *The expected welfare gain of the optimal non-linear tax relative to the optimal linear tax and the optimal standard is weakly positive. It depends on the local bias heterogeneity A as follows:*

- a) *The expected welfare gain over the optimal linear tax is zero when $A = 0$ and increases in the absolute value of A ;*
- b) *The expected welfare gain over the optimal standard is zero when $A \geq 1$ and decreases in A .*

We now provide the proof of Proposition 8. We analyze the impact of a change in the local bias heterogeneity A on expected welfare under the optimal non-linear tax, the optimal linear tax, and the optimal standard. To separate the impact of changes in A from changes in the population of consumers, we hold the distribution of \hat{v} as well as the first moments of v and b constant. In addition, we consider a scenario in which redistribution does not matter, i.e., $\hat{g}(\hat{v}) = 1$.

To evaluate welfare implications of a change in A , we need to evaluate the derivative with respect to A of the expected equilibrium normative utility (net of taxes) for consumers with \hat{v} (see also the Hamiltonian of the policy maker's problem):

$$\int_{\hat{v}} \hat{u}^n(A) dF(\hat{v}) = \int_{\hat{v}} E[v|\hat{v}](A) q^M(A) - c(q^M(A)) dF(\hat{v}), \quad (40)$$

where q^M denotes the allocation under the non-linear, price (linear tax) or standard (consumption ban) mechanism, $q^M(A) \in \{\xi(A), q^P(A), q^S\}$, and $\hat{u}^n(A, q^M)$ denotes the equilibrium normative utility in the respective mechanism.

We first consider the welfare implications of a change in A under the optimal non-linear tax. We have that:

$$\frac{\partial \hat{u}^n(A, \xi)}{\partial A} = \frac{\partial E[v|\hat{v}]}{\partial A} \xi + \frac{\partial \xi}{\partial A} \underbrace{\{E[v|\hat{v}] - c'(\xi)\}}_{=0}.$$

By an envelope theorem argument, the second summand is zero since $E[v|\hat{v}] = \hat{v} - E[b|\hat{v}] = c'(\xi)$ due to optimal consumer behavior (see Equation (11)).

The optimal standard q^S can be calculated by solving the following utility maximization problem: $\max_{q \in \mathbb{R}} \int_v u^n(q, v) dF_v(v)$. Inserting $u^n(q, v) = m + vx - c(q)$ and solving for the maximum yields the following implicit solution: $c'(q^S) = E[v]$. Hence, q^S does not depend on A and we obtain:

$$\frac{\partial \hat{u}^n(A, q^S)}{\partial A} = \frac{\partial E[v|\hat{v}]}{\partial A} q^S.$$

As shown in Section A.B, the optimal linear tax is given by the following expression: $t^*(A) = \int_{\hat{v}} E(b|\hat{v})(A) \left[\frac{\partial q^d}{\partial t} / \frac{\partial \bar{q}^d}{\partial t} \right] dF(\hat{v})$. From the first order condition of consumer maximization, we know that $c'(q^d) = \hat{v} - t^*(A) = E[v|\hat{v}] + E[b|\hat{v}] - t^*(A)$. Using this equation, we obtain:

$$\frac{\partial \hat{u}^n(A, q^P)}{\partial A} = \frac{\partial E[v|\hat{v}]}{\partial A} q^P + \frac{\partial q^P}{\partial A} \{E[v|\hat{v}] - c'(q^P)\} = \frac{\partial E[v|\hat{v}]}{\partial A} q^P + \frac{\partial q^P}{\partial A} \{E[b|\hat{v}] - t^*(A)\}.$$

The second summand vanishes when taking the integral over all types, which can be seen as follows:

$$\begin{aligned} \int_{\hat{v}} \frac{\partial q^P}{\partial A} \{E[b|\hat{v}] - t^*(A)\} dF(\hat{v}) &= \frac{\partial t^*(A)}{\partial A} \int_{\hat{v}} \frac{\partial q^P}{\partial t^*} \{E[b|\hat{v}] - t^*(A)\} dF(\hat{v}) \\ &= \frac{\partial t^*(A)}{\partial A} \left(\int_{\hat{v}} \frac{\partial q^P}{\partial t^*} E[b|\hat{v}] dF(\hat{v}) - t^*(A) \int_{\hat{v}} \frac{\partial q^P}{\partial t^*} dF(\hat{v}) \right) = 0, \end{aligned}$$

where the last equality follows from the definition of $t^*(A)$ and the fact that $\frac{\partial \bar{q}^d}{\partial t^*} = \int_{\hat{v}} \frac{\partial q^P}{\partial t^*} dF(\hat{v})$. Hence, we obtain:

$$\frac{\partial \hat{u}^n(A, q^P)}{\partial A} = \frac{\partial E[v|\hat{v}]}{\partial A} q^P.$$

We now evaluate how a change in A changes the relative advantage of non-linear taxation relative to a standard in terms of expected welfare:

$$\frac{\partial \Delta u_{NL,S}^n}{\partial A} = \int \frac{\partial E[v|\hat{v}]}{\partial A} (\xi - q^S) dF(\hat{v}),$$

where $\Delta u_{NL,S}^n := \int \hat{u}^n(A, \zeta) - \hat{u}^n(A, q^S) dF(\hat{\nu})$. Using $\frac{\partial E[v|\hat{\nu}]}{\partial A} = (\mu_{\hat{\nu}} - \hat{\nu})$, we obtain:

$$\begin{aligned} \frac{\partial \Delta u_{NL,S}^n}{\partial A} &= \int (\mu_{\hat{\nu}} - \hat{\nu}) [\zeta(\hat{\nu}) - q^S] dF(\hat{\nu}) \\ &= \mu_{\hat{\nu}} \int \zeta(\hat{\nu}) dF(\hat{\nu}) - \int \hat{\nu} \zeta(\hat{\nu}) dF(\hat{\nu}) - q^S \left[\mu_{\hat{\nu}} - \int \hat{\nu} dF(\hat{\nu}) \right] \\ &= E[\hat{\nu}] \cdot E[\zeta(\hat{\nu})] - E[\hat{\nu} \cdot \zeta(\hat{\nu})] - 0 \\ &= E[\hat{\nu}] \cdot E[\zeta(\hat{\nu})] - E[\hat{\nu}] \cdot E[\zeta(\hat{\nu})] - cov(\hat{\nu}, \zeta(\hat{\nu})) \leq 0, \end{aligned}$$

where the last inequality holds since ζ is increasing in $\hat{\nu}$ for an incentive compatible mechanism, so that $cov(\hat{\nu}, \zeta(\hat{\nu}))$ is positive.

The relative advantage in terms of expected welfare of non-linear taxation relative to the optimal linear tax is given by:

$$\frac{\partial \Delta u_{NL,P}^n}{\partial A} = \int \frac{\partial E[v|\hat{\nu}]}{\partial A} (\zeta(\hat{\nu}) - q^P(\hat{\nu})) dF\hat{\nu}.$$

We approximate $\zeta - q^P$ by a first-order Taylor approximation in the marginal tax rate t' around $t' = t^*$, which yields:

$$\zeta - q^P \approx \left. \frac{dq^P}{dt} \right|_{t'=t^*} [t'(\hat{\nu}) - t^*],$$

where $t'(\hat{\nu}) = E[b|\hat{\nu}]$ is the optimal non-linear tax rate from Proposition 1 and t^* is the optimal linear tax rate from Equation (22).

Using that $E[b|\hat{\nu}] = E[b|\mu_{\hat{\nu}}] + A(\hat{\nu} - \mu_{\hat{\nu}})$, rearranging, and using $W(\hat{\nu}) := \frac{dq(\hat{\nu})/dt}{d\bar{q}/dt}$ with $\frac{d\bar{q}}{dt} := \int \hat{\nu} \frac{dq(\hat{\nu})}{dt} dF(\hat{\nu})$ to denote the relative responsiveness of a consumer type $\hat{\nu}$ we obtain:

$$\zeta - q^P \approx \left. \frac{dq^P}{dt} \right|_{t'=t^*} A \left(\hat{\nu} - \int \hat{\nu} W(\hat{\nu}) dF(\hat{\nu}) \right).$$

Denoting $v^P = \int \hat{\nu} W(\hat{\nu}) dF(\hat{\nu})$, we have:

$$\begin{aligned} \frac{\partial \Delta u_{NL,P}^c}{\partial A} &= A \int (\mu_{\hat{\nu}} - \hat{\nu})(\hat{\nu} - v^P) \left. \frac{dq^P}{dt} \right|_{t'=t^*} dF(\hat{\nu}) \\ &= A \int (\mu_{\hat{\nu}} \hat{\nu} - \mu_{\hat{\nu}} v^P - \hat{\nu}^2 + \hat{\nu} v^P) \left. \frac{dq^P}{dt} \right|_{t'=t^*} dF(\hat{\nu}) \\ &= A \left((\mu_{\hat{\nu}} + v^P) \int \hat{\nu} \left. \frac{dq^P}{dt} \right|_{t'=t^*} dF(\hat{\nu}) - (\mu_{\hat{\nu}} v^P) \int \left. \frac{dq^P}{dt} \right|_{t'=t^*} dF(\hat{\nu}) - \int \hat{\nu}^2 \left. \frac{dq^P}{dt} \right|_{t'=t^*} dF(\hat{\nu}) \right) \\ &= A \frac{d\bar{q}^P}{dt} \left((\mu_{\hat{\nu}} + v^P) \int \hat{\nu} W(\hat{\nu}) dF(\hat{\nu}) - \mu_{\hat{\nu}} v^P - \int \hat{\nu}^2 W(\hat{\nu}) dF(\hat{\nu}) \right) \end{aligned}$$

$$\begin{aligned}
&= A \frac{d\bar{q}^P}{dt} \left((v^P)^2 - \int \hat{v}^2 W(\hat{v}) dF(\hat{v}) \right) \\
&= \underbrace{A \frac{d\bar{q}^P}{dt}}_{\leq 0 \text{ if } A \geq 0} \underbrace{\left(\left(\int \hat{v} W(\hat{v}) dF(\hat{v}) \right)^2 - \int \hat{v}^2 W(\hat{v}) dF(\hat{v}) \right)}_{\leq 0} \geq 0,
\end{aligned}$$

where $\frac{d\bar{q}^P}{dt} = \int \frac{dq^P}{dt} \Big|_{t'=t^*} dF(\hat{v})$. The fact that the third factor is non-positive follows from Jensen's inequality because $d(\hat{v}) = W(\hat{v})f(\hat{v})$ satisfies the properties of a density function, with $d(\hat{v}) \geq 0 \forall \hat{v}$ and $\int d(\hat{v})d\hat{v} = \int W(\hat{v})dF(\hat{v}) = 1$.

Comparison of Linear Price Instruments with Standards

In Proposition 8, we have shown that a linear price instrument is welfare-optimal for $A = 0$, a standard is optimal for $A = 1$, and the advantage of a standard compared to a linear price instrument increases in A . Hence, the existence of \hat{A} from Proposition 6 is guaranteed by applying the intermediate value theorem.

To complete the proof of Proposition 6, we need to show that the welfare advantage of a price instrument over a quantity instrument becomes more pronounced as A decreases below zero. For that purpose, it is sufficient to show that the difference in the expected equilibrium normative utility between the linear tax and the quantity regulation, $\partial \Delta u_{P,S}^n / \partial A = \int \hat{u}^n(A, q^P) - \hat{u}^n(A, q^S) dF(\hat{v})$, decreases in A . Let \bar{v} be the \hat{v} such that $q^P(\bar{v}) = q^S$. We approximate the demand function by a first-order Taylor approximation at \bar{v} , which yields $q^P(\hat{v}) \approx q^P(\bar{v}) + \frac{\partial q^P}{\partial \hat{v}} \Big|_{\bar{v}} (\hat{v} - \bar{v})$. As the first-order conditions of utility maximization imply $\hat{v} - t^* = c'(q^P(\hat{v}))$ for linear taxation and $c'(q^S) = E[v]$ for quantity regulation, we have that $\bar{v} = E[v] + t^*$, where $t^* = \int_{\hat{v}} E[b|\hat{v}] W(\hat{v}) dF(\hat{v})$. We can now show that the welfare difference between prices and quantities decreases in A , since

$$\begin{aligned}
\frac{\partial \Delta u_{P,S}^n}{\partial A} &= \int_{\hat{v}} \frac{\partial E[v|\hat{v}]}{\partial A} [q^P(\hat{v}) - q^S] dF(\hat{v}) \\
&= \int_{\hat{v}} (\mu_{\hat{v}} - \hat{v}) [q^P(\hat{v}) - q^S] dF(\hat{v}) \\
&= - \frac{\partial q^P}{\partial \hat{v}} \Big|_{\hat{v}=\bar{v}} \int_{\hat{v}} (\mu_{\hat{v}} - \hat{v}) (\hat{v} - \bar{v}) dF(\hat{v}) \\
&= - \frac{\partial q^P}{\partial \hat{v}} \Big|_{\hat{v}=\bar{v}} \int_{\hat{v}} (\mu_{\hat{v}} - \hat{v}) [(v - E[v]) (b - t^*)] dF(\hat{v}) \\
&= \underbrace{- \frac{\partial q^P}{\partial \hat{v}} \Big|_{\hat{v}=\bar{v}}}_{<0} \underbrace{[E[\hat{v}^2] - \mu_{\hat{v}}^2]}_{>0} < 0,
\end{aligned}$$

where the evaluation of the second factor in the last line follows from Jensen's inequality.

A.G Proof of Proposition 7: Minimum and Maximum Standards

The proof proceeds as follows. From Proposition 3, we know that $\partial E[v|\hat{v}]/\partial \hat{v} < 0$ for some \hat{v} implies that a policy maker will use bunching as part of the optimal policy. Furthermore, we know from Appendix A.D that bans are part of the optimal policy mix if bunching occurs at the top or the bottom of the perceived valuation distribution. We use a second-order Taylor approximation to approximate $E[v|\hat{v}]$ as a function of \hat{v} and \hat{v}^2 . This allows to derive the conditions under which mixed policies optimally involve bans for high or low levels of the product attribute in terms of higher-order moments of the (joint) distributions of the random variables v , b , and \hat{v} .

The second-order Taylor approximation of $E[v|\hat{v}]$ is:

$$\hat{E}[v|\hat{v}] = c_0 + c_1\hat{v} + c_2\hat{v}^2. \quad (41)$$

Bunching occurs whenever $\partial E[v|\hat{v}]/\partial \hat{v} = c_1 + 2c_2\hat{v} \leq 0$ for some \hat{v} . For the purpose of this proof, we are interested in bunching at the top and at the bottom. With $\hat{v} \in [0, \infty)$, we have:

1. Bunching at the Top if $c_1 > 0$ and $c_2 < 0$.
2. Bunching at the Bottom if $c_1 < 0$ and $c_2 > 0$.

Otherwise, bunching takes place over the entire support of \hat{v} ($c_1 < 0, c_2 \leq 0$) or not at all ($c_1 > 0, c_2 \geq 0$). We proceed by first calculating c_1 and c_2 . In a subsequent step, we derive sufficient conditions for bunching at the top and at the bottom.

The Frisch-Waugh-Lovell Theorem demonstrates that c_r for $r \in \{1, 2\}$, as defined by Equation (41), can be calculated as

$$c_r = \frac{\text{cov}(\epsilon_r, \hat{\epsilon}_r)}{\text{var}(\hat{\epsilon}_r)}, \quad (42)$$

where $\hat{\epsilon}_r$ is the "residual" of \hat{v}^r after regressing it on all covariates except \hat{v}^r , and ϵ_r is the "residual" of v after regressing it on all covariates except \hat{v}^r .

We first calculate c_1 . Regressing \hat{v} on \hat{v}^2 and a constant yields:

$$\begin{aligned} \hat{v} &= d_1 + d_2\hat{v}^2 + \hat{\epsilon}_1 \\ &= E[\hat{v}|\hat{v}^2 = 0] + \frac{\text{cov}(\hat{v}, \hat{v}^2)}{\text{var}(\hat{v}^2)}\hat{v}^2 + \hat{\epsilon}_1. \end{aligned}$$

The residual $\hat{\epsilon}_1$ is then given by:

$$\hat{\epsilon}_1 = \hat{v} - \left\{ E[\hat{v}|\hat{v}^2 = 0] + \frac{\text{cov}(\hat{v}, \hat{v}^2)}{\text{var}(\hat{v}^2)} \hat{v}^2 \right\}.$$

Regressing v on \hat{v}^2 and a constant yields:

$$\begin{aligned} v &= e_0 + e_1 \hat{v}^2 + \epsilon_1 \\ &= E[v|\hat{v}^2 = 0] + \frac{\text{cov}(v, \hat{v}^2)}{\text{var}(\hat{v}^2)} \hat{v}^2 + \epsilon_1 \end{aligned}$$

The residual ϵ_1 is then given by:

$$\epsilon_1 = v - \left\{ E[v|\hat{v}^2 = 0] + \frac{\text{cov}(v, \hat{v}^2)}{\text{var}(\hat{v}^2)} \hat{v}^2 \right\}.$$

Therefore, by Equation (42), the coefficient c_1 is given by:

$$\begin{aligned} c_1 &= \frac{\text{cov}(\epsilon_1, \hat{\epsilon}_1)}{\text{var}(\hat{\epsilon}_1)} \\ &= \frac{\text{cov}\left(v - \left\{ E[v|\hat{v}^2 = 0] + \frac{\text{cov}(v, \hat{v}^2)}{\text{var}(\hat{v}^2)} \hat{v}^2 \right\}, \hat{v} - \left\{ E[\hat{v}|\hat{v}^2 = 0] + \frac{\text{cov}(\hat{v}, \hat{v}^2)}{\text{var}(\hat{v}^2)} \hat{v}^2 \right\}\right)}{\text{var}(\hat{\epsilon}_1)} \\ &= \frac{\text{cov}\left(v, \hat{v} - \frac{\text{cov}(\hat{v}, \hat{v}^2)}{\text{var}(\hat{v}^2)} \hat{v}^2\right)}{\text{var}(\hat{\epsilon}_1)} + \frac{\text{cov}\left(-\frac{\text{cov}(v, \hat{v}^2)}{\text{var}(\hat{v}^2)} \hat{v}^2, \hat{v} - \frac{\text{cov}(\hat{v}, \hat{v}^2)}{\text{var}(\hat{v}^2)} \hat{v}^2\right)}{\text{var}(\hat{\epsilon}_1)} \\ &= \frac{\text{cov}\left(v, \hat{v} - \frac{\text{cov}(\hat{v}, \hat{v}^2)}{\text{var}(\hat{v}^2)} \hat{v}^2\right)}{\text{var}(\hat{\epsilon}_1)} - \frac{\frac{\text{cov}(\hat{v}, \hat{v}^2)}{\text{var}(\hat{v}^2)} \text{cov}(v, \hat{v}^2)}{\text{var}(\hat{\epsilon}_1)} + \frac{\frac{\text{cov}(\hat{v}, \hat{v}^2) \text{cov}(v, \hat{v}^2)}{[\text{var}(\hat{v}^2)]^2} \text{cov}(\hat{v}^2, \hat{v}^2)}{\text{var}(\hat{\epsilon}_1)} \\ &= \frac{\text{cov}\left(v, \hat{v} - \frac{\text{cov}(\hat{v}, \hat{v}^2)}{\text{var}(\hat{v}^2)} \hat{v}^2\right)}{\text{var}(\hat{\epsilon}_1)}. \end{aligned} \tag{43}$$

Next, we calculate c_2 . Regressing \hat{v}^2 on \hat{v} and a constant yields:

$$\begin{aligned} \hat{v}^2 &= d_0 + d_1 \hat{v} + \hat{\epsilon}_2 \\ &= E[\hat{v}^2|\hat{v} = 0] + \frac{\text{cov}(\hat{v}^2, \hat{v})}{\text{var}(\hat{v})} \hat{v} + \hat{\epsilon}_2. \end{aligned}$$

The residual $\hat{\epsilon}_2$ is then given by:

$$\hat{\epsilon}_2 = \hat{\vartheta}^2 - \left\{ E [\hat{\vartheta}^2 | \hat{\vartheta} = 0] + \frac{\text{cov}(\hat{\vartheta}^2, \hat{\vartheta})}{\text{var}(\hat{\vartheta})} \hat{\vartheta} \right\}.$$

Regressing v on $\hat{\vartheta}$ and a constant yields:

$$\begin{aligned} v &= d_0 + d_1 \hat{\vartheta} + \epsilon_2 \\ &= E[v | \hat{\vartheta} = 0] + \frac{\text{cov}(v, \hat{\vartheta})}{\text{var}(\hat{\vartheta})} \hat{\vartheta} + \epsilon_2. \end{aligned}$$

The residual ϵ_2 is then given by:

$$\epsilon_2 = v - \left\{ E[v | \hat{\vartheta} = 0] + \frac{\text{cov}(v, \hat{\vartheta})}{\text{var}(\hat{\vartheta})} \hat{\vartheta} \right\}.$$

Therefore, by Equation (42), the coefficient c_2 is given by:

$$\begin{aligned} c_2 &= \frac{\text{cov}(\epsilon_2, \hat{\epsilon}_2)}{\text{var}(\hat{\epsilon}_2)} \\ &= \frac{\text{cov}\left(v - \left\{ E[v | \hat{\vartheta} = 0] + \frac{\text{cov}(v, \hat{\vartheta})}{\text{var}(\hat{\vartheta})} \hat{\vartheta} \right\}, \hat{\vartheta}^2 - \left\{ E[\hat{\vartheta}^2 | \hat{\vartheta} = 0] + \frac{\text{cov}(\hat{\vartheta}^2, \hat{\vartheta})}{\text{var}(\hat{\vartheta})} \hat{\vartheta} \right\}\right)}{\text{var}(\hat{\epsilon}_2)} \\ &= \frac{\text{cov}(v, \hat{\vartheta}^2) - \frac{\text{cov}(v, \hat{\vartheta})}{\text{var}(\hat{\vartheta})} \text{cov}(\hat{\vartheta}^2, \hat{\vartheta})}{\text{var}(\hat{\epsilon}_2)} + \frac{\text{cov}\left(-\frac{\text{cov}(v, \hat{\vartheta})}{\text{var}(\hat{\vartheta})} \hat{\vartheta}, \hat{\vartheta}^2 - \frac{\text{cov}(\hat{\vartheta}^2, \hat{\vartheta})}{\text{var}(\hat{\vartheta})} \hat{\vartheta}\right)}{\text{var}(\hat{\epsilon}_2)} \\ &= \frac{\text{cov}(v, \hat{\vartheta}^2) - \frac{\text{cov}(v, \hat{\vartheta})}{\text{var}(\hat{\vartheta})} [\text{cov}(\hat{\vartheta}^2, v) + \text{cov}(\hat{\vartheta}^2, b)]}{\text{var}(\hat{\epsilon}_2)} + \frac{-\frac{\text{cov}(v, \hat{\vartheta})}{\text{var}(\hat{\vartheta})} \text{cov}(\hat{\vartheta}, \hat{\vartheta}^2) + \text{cov}(v, \hat{\vartheta}) \frac{\text{cov}(\hat{\vartheta}^2, \hat{\vartheta})}{\text{var}(\hat{\vartheta})}}{\text{var}(\hat{\epsilon}_2)} \\ &= \frac{\text{cov}(v, \hat{\vartheta}^2) - \frac{\text{cov}(v, \hat{\vartheta})}{\text{var}(\hat{\vartheta})} [\text{cov}(\hat{\vartheta}^2, v) + \text{cov}(\hat{\vartheta}^2, b)]}{\text{var}(\hat{\epsilon}_2)} \tag{44} \\ &= \frac{\text{cov}(v, \hat{\vartheta}^2) \text{cov}(b, \hat{\vartheta}) - \text{cov}(v, \hat{\vartheta}) \text{cov}(\hat{\vartheta}^2, b)}{\text{var}(\hat{\vartheta}) \text{var}(\hat{\epsilon}_2)} \\ &= \frac{\text{cov}(\hat{\vartheta}, \hat{\vartheta}^2) \text{cov}(b, \hat{\vartheta}) - \text{cov}(\hat{\vartheta}, \hat{\vartheta}) \text{cov}(\hat{\vartheta}^2, b)}{\text{var}(\hat{\vartheta}) \text{var}(\hat{\epsilon}_2)} \\ &\stackrel{(*)}{=} \frac{[2\mu_{\hat{\vartheta}} \sigma_{\hat{\vartheta}}^2 + \text{sk}(\hat{\vartheta}) \sigma_{\hat{\vartheta}}^3] \text{cov}(b, \hat{\vartheta}) - \sigma_{\hat{\vartheta}}^2 \text{cov}(\hat{\vartheta}^2, b)}{\text{var}(\hat{\vartheta}) \text{var}(\hat{\epsilon}_2)} \\ &= \frac{\sigma_{\hat{\vartheta}}^2 \overbrace{\{ [2\mu_{\hat{\vartheta}} + \text{sk}(\hat{\vartheta}) \sigma_{\hat{\vartheta}}] \text{cov}(b, \hat{\vartheta}) - \text{cov}(\hat{\vartheta}^2, b) \}}{=:K}}{\text{var}(\hat{\vartheta}) \text{var}(\hat{\epsilon}_2)}, \end{aligned}$$

where the denominator is always positive and (*) follows from the fact that, for every random variable Y , we have:²³

$$\text{cov}(Y, Y^2) = [\text{sk}(Y)\sigma_Y + 2\mu_Y] \sigma_Y^2. \quad (45)$$

We now analyze the term K :

$$\begin{aligned} K &\stackrel{\text{BG1969}}{=} [2\mu_{\hat{\vartheta}} + \text{sk}(\hat{\vartheta})\sigma_{\hat{\vartheta}}] \text{cov}(b, \hat{\vartheta}) - [2\mu_{\hat{\vartheta}}\text{cov}(\hat{\vartheta}, b) + E[(\hat{\vartheta} - \mu_{\hat{\vartheta}})^2(b - \mu_b)]] \\ &= \text{sk}(\hat{\vartheta})\sigma_{\hat{\vartheta}}\text{cov}(b, \hat{\vartheta}) - E[(\hat{\vartheta} - \mu_{\hat{\vartheta}})^2(b - \mu_b)] \\ &= E\left[\left(\frac{\hat{\vartheta} - \mu_{\hat{\vartheta}}}{\sigma_{\hat{\vartheta}}}\right)^3\right] \sigma_{\hat{\vartheta}} E[(\hat{\vartheta} - \mu_{\hat{\vartheta}})(b - \mu_b)] - E[(\hat{\vartheta} - \mu_{\hat{\vartheta}})^2(b - \mu_b)] \\ &= \frac{E[(\hat{\vartheta} - \mu_{\hat{\vartheta}})^3] E[(\hat{\vartheta} - \mu_{\hat{\vartheta}})(b - \mu_b)]}{\sigma_{\hat{\vartheta}}^2} - E[(\hat{\vartheta} - \mu_{\hat{\vartheta}})^2(b - \mu_b)] \\ &= \frac{E[(\hat{\vartheta} - \mu_{\hat{\vartheta}})^3] E[(\hat{\vartheta} - \mu_{\hat{\vartheta}})(b - \mu_b)]}{\sigma_{\hat{\vartheta}}^2} - E[(\hat{\vartheta} - \mu_{\hat{\vartheta}})^2(b - \mu_b)] \\ &= \sigma_{\hat{\vartheta}}^2 \sigma_b \left\{ \frac{E[(\hat{\vartheta} - \mu_{\hat{\vartheta}})^3] E[(\hat{\vartheta} - \mu_{\hat{\vartheta}})(b - \mu_b)]}{\sigma_{\hat{\vartheta}}^4 \sigma_b} - \frac{E[(\hat{\vartheta} - \mu_{\hat{\vartheta}})^2(b - \mu_b)]}{\sigma_{\hat{\vartheta}}^2 \sigma_b} \right\} \\ &= \sigma_{\hat{\vartheta}}^2 \sigma_b \{ \text{sk}(\hat{\vartheta}) \rho(\hat{\vartheta}, b) - \text{cosk}(\hat{\vartheta}, b) \}, \end{aligned}$$

where BG1969 refers to Bohrnstedt and Goldberger (1969). Hence, a sufficient (and necessary) condition for $c_2 < 0$ is:

$$c_2 < 0 \iff \rho(\hat{\vartheta}, b) \text{sk}(\hat{\vartheta}) < \text{cosk}(\hat{\vartheta}, b) \quad (46)$$

and a sufficient (and necessary) condition for $c_2 > 0$ is:

$$c_2 > 0 \iff \rho(\hat{\vartheta}, b) \text{sk}(\hat{\vartheta}) > \text{cosk}(\hat{\vartheta}, b). \quad (47)$$

²³This can be seen as follows:

$$\begin{aligned} \text{cov}(Y, Y^2) &= E(Y^3) - \mu_Y^2 \mu_Y = E(Y^3) - \mu_Y [E(Y^2) - \mu_Y^2 + \mu_Y^2] = E(Y^3) - \mu_Y \sigma_Y^2 - \mu_Y^3 \\ &= \frac{[E(Y^3) - \mu_Y^3 - 3\mu_Y \sigma_Y^2 + 3\mu_Y \sigma_Y^2] \sigma_Y^3}{\sigma_Y^3} - \mu_Y \sigma_Y^2 = \text{sk}(Y) \sigma_Y^3 + 3\mu_Y \sigma_Y^2 - \mu_Y \sigma_Y^2 = \text{sk}(Y) \sigma_Y^3 + 2\mu_Y \sigma_Y^2. \end{aligned}$$

Sufficient Conditions for Bunching at the Top

To obtain sufficient conditions for Bunching at the Top, we need to find sufficient conditions for $c_1 > 0$, given that $c_2 < 0$. Using the results for c_1 from Equation (43), we obtain:

$$\begin{aligned}
& c_1 > 0 \\
& \Leftrightarrow \frac{\text{cov} \left(v, \hat{v} - \frac{\text{cov}(\hat{v}, \hat{v}^2)}{\text{var}(\hat{v}^2)} \hat{v}^2 \right)}{\text{var}(\hat{\epsilon}_1)} > 0 \\
& \Leftrightarrow \underbrace{\text{cov}(v, \hat{v})}_{(1)} - \underbrace{\text{cov}(v, \hat{v}^2)}_{(2)} \underbrace{\frac{\text{cov}(\hat{v}, \hat{v}^2)}{\text{var}(\hat{v}^2)}}_{(3)} > 0. \tag{48}
\end{aligned}$$

Term (1) of Inequality (48) is positive by assumption because we restrict ourselves to settings where partial bans become relevant (note that $A < 1$ from Proposition 5 is equivalent to $\text{cov}(v, \hat{v}) > 0$). A sufficient condition for Term (3) to be positive is that the skewness is positive, which follows from:

$$\begin{aligned}
\text{cov}(\hat{v}, \hat{v}^2) > 0 & \Leftrightarrow 2\mu_{\hat{v}}\sigma_{\hat{v}}^2 + \text{sk}(\hat{v})\sigma_{\hat{v}}^3 > 0 \\
& \Leftrightarrow \text{sk}(\hat{v}) > -\frac{2\mu_{\hat{v}}}{\sigma_{\hat{v}}} \\
& \Leftrightarrow \boxed{\text{sk}(\hat{v}) > 0}.
\end{aligned}$$

Next, we show that Inequality (48) holds under the assumptions made so far. When $\text{cov}(v, \hat{v}^2) \leq 0$, the inequality holds because Term (1) and Term (3) are positive. When $\text{cov}(v, \hat{v}^2) > 0$, we can derive an upper bound for Term (2) by rearranging Equation (44):

$$c_2 < 0 \iff \text{cov}(v, \hat{v}^2) < \text{cov}(v, \hat{v}) \frac{\text{cov}(\hat{v}^2, \hat{v})}{\text{var}(\hat{v})}. \tag{49}$$

Then, Inequality (48) holds if it holds after substituting the upper bound from Inequality (49) for Term (2):

$$\begin{aligned}
\text{cov}(v, \hat{v}) - \text{cov}(v, \hat{v}) \frac{\text{cov}(\hat{v}^2, \hat{v})}{\text{var}(\hat{v})} \frac{\text{cov}(\hat{v}, \hat{v}^2)}{\text{var}(\hat{v}^2)} & > 0 \\
& \Leftrightarrow \text{cov}(v, \hat{v}) \left(1 - \rho_{\hat{v}^2, \hat{v}}^2 \right) > 0,
\end{aligned}$$

where $\rho_{\hat{v}^2, \hat{v}}^2 \in (-1, 1)$ is the square of the correlation between \hat{v}^2 and \hat{v} . As we consider settings with $\text{cov}(v, \hat{v}) > 0$, this inequality holds.

To close the proof, note that $sk(\hat{\vartheta}) > 0$ in combination with Equation (46) implies that:

$$\boxed{\frac{cosk(\hat{\vartheta}, b)}{sk(\hat{\vartheta})} > \rho(\hat{\vartheta}, b)}.$$

Hence, this condition and the condition that $\boxed{sk(\hat{\vartheta}) > 0}$ are sufficient for Bunching at the Top if $\hat{\vartheta} \in [0, \infty)$ and $cov(v, \hat{\vartheta}) > 0$.

Sufficient Conditions for Bunching at the Bottom

To obtain sufficient conditions for Bunching at the Bottom, we need to find sufficient conditions for $c_1 < 0$, given that $c_2 > 0$. Using the results for c_1 from Equation (48), we obtain:

$$c_1 < 0 \iff cov(v, \hat{\vartheta}) - cov(v, \hat{\vartheta}^2) \frac{cov(\hat{\vartheta}, \hat{\vartheta}^2)}{var(\hat{\vartheta}^2)} < 0. \quad (50)$$

Using Equation (45), we obtain that

$$cov(\hat{\vartheta}, \hat{\vartheta}^2) = [sk(\hat{\vartheta})\sigma_{\hat{\vartheta}} + 2\mu_{\hat{\vartheta}}] \sigma_{\hat{\vartheta}}^2.$$

Furthermore,

$$\begin{aligned} cov(v, \hat{\vartheta}^2) &= E[(v - \mu_v)(\hat{\vartheta}^2 - \mu_{\hat{\vartheta}^2})] \\ &= E(v\hat{\vartheta}^2) - \mu_{\hat{\vartheta}^2}\mu_v - \mu_v E(\hat{\vartheta}^2) + \mu_v\mu_{\hat{\vartheta}^2} \\ &= \frac{[E(v\hat{\vartheta}^2) - \mu_v E(\hat{\vartheta}^2) - 2E(v\hat{\vartheta})\mu_{\hat{\vartheta}} + 2\mu_{\hat{\vartheta}}^2\mu_v + 2E(v\hat{\vartheta})\mu_{\hat{\vartheta}} - 2\mu_{\hat{\vartheta}}^2\mu_v] \sigma_{\hat{\vartheta}}^2 \sigma_v}{\sigma_{\hat{\vartheta}}^2 \sigma_v} \\ &\stackrel{(*)}{=} cosk(\hat{\vartheta}, v)\sigma_{\hat{\vartheta}}^2\sigma_v + 2E(v\hat{\vartheta})\mu_{\hat{\vartheta}} - 2\mu_{\hat{\vartheta}}^2\mu_v \\ &= cosk(\hat{\vartheta}, v)\sigma_{\hat{\vartheta}}^2\sigma_v + 2\mu_{\hat{\vartheta}}[E(v\hat{\vartheta}) - \mu_{\hat{\vartheta}}\mu_v + \mu_{\hat{\vartheta}}\mu_v - \mu_{\hat{\vartheta}}\mu_v] \\ &= cosk(\hat{\vartheta}, v)\sigma_{\hat{\vartheta}}^2\sigma_v + 2\mu_{\hat{\vartheta}}cov(v, \hat{\vartheta}) \\ &= [cosk(\hat{\vartheta}, v)\sigma_{\hat{\vartheta}} + 2\mu_{\hat{\vartheta}}\rho(v, \hat{\vartheta})] \sigma_{\hat{\vartheta}}\sigma_v, \end{aligned}$$

where in (*) we use that $cosk(\hat{\vartheta}, v) = E(v\hat{\vartheta}^2) - 2E(v\hat{\vartheta})\mu_{\hat{\vartheta}} - 2E(\hat{\vartheta}^2)\mu_v + 2\mu_{\hat{\vartheta}}^2\mu_v$. Hence, we can rewrite Inequality (50) as follows

$$\begin{aligned} cov(v, \hat{\vartheta}) - cov(v, \hat{\vartheta}^2) \frac{cov(\hat{\vartheta}, \hat{\vartheta}^2)}{var(\hat{\vartheta}^2)} &< 0 \\ \Leftrightarrow cov(v, \hat{\vartheta}) - [cosk(\hat{\vartheta}, v)\sigma_{\hat{\vartheta}} + 2\mu_{\hat{\vartheta}}\rho(v, \hat{\vartheta})] \sigma_{\hat{\vartheta}} \frac{[sk(\hat{\vartheta})\sigma_{\hat{\vartheta}} + 2\mu_{\hat{\vartheta}}] \sigma_{\hat{\vartheta}}^2}{var(\hat{\vartheta}^2)} &< 0 \end{aligned}$$

$$\begin{aligned}
&\Leftrightarrow [\text{cosk}(\hat{v}, v)\sigma_{\hat{v}} + 2\mu_{\hat{v}}\rho(v, \hat{v})] \sigma_{\hat{v}}\sigma_v \frac{[\text{sk}(\hat{v})\sigma_{\hat{v}} + 2\mu_{\hat{v}}] \sigma_{\hat{v}}^2}{\text{var}(\hat{v}^2)} > \text{cov}(v, \hat{v}) \\
&\Leftrightarrow [\text{cosk}(\hat{v}, v)\sigma_{\hat{v}} + 2\mu_{\hat{v}}\rho(v, \hat{v})] \frac{[\text{sk}(\hat{v})\sigma_{\hat{v}} + 2\mu_{\hat{v}}] \sigma_{\hat{v}}^2}{\text{var}(\hat{v}^2)} \stackrel{(*)}{>} \rho(v, \hat{v}) \\
&\Leftrightarrow \left[\frac{\text{cosk}(\hat{v}, v)}{\rho(v, \hat{v})} \sigma_{\hat{v}} + 2\mu_{\hat{v}} \right] \frac{[\text{sk}(\hat{v})\sigma_{\hat{v}} + 2\mu_{\hat{v}}] \sigma_{\hat{v}}^2}{\text{var}(\hat{v}^2)} > 1 \\
&\Leftrightarrow \left[\frac{\text{cosk}(\hat{v}, v)}{\rho(v, \hat{v})} \sigma_{\hat{v}} + 2\mu_{\hat{v}} \right] [\text{sk}(\hat{v})\sigma_{\hat{v}} + 2\mu_{\hat{v}}] > \frac{\text{var}(\hat{v}^2)}{\sigma_{\hat{v}}^2} \\
&\Leftrightarrow \left[\frac{\text{cosk}(\hat{v}, v)}{\rho(v, \hat{v})} \sigma_{\hat{v}} + 2\mu_{\hat{v}} \right] [\text{sk}(\hat{v})\sigma_{\hat{v}} + 2\mu_{\hat{v}}] - \frac{\sigma_{\hat{v}}^2}{\sigma_{\hat{v}}^2} > 0 \\
&\Leftrightarrow \underbrace{\left[\frac{\text{cosk}(\hat{v}, v)}{\rho(v, \hat{v})} \sigma_{\hat{v}} + 2\mu_{\hat{v}} + \frac{\sigma_{\hat{v}}^2}{\sigma_{\hat{v}}^2} \right]}_{(1)} \underbrace{\left[\text{sk}(\hat{v})\sigma_{\hat{v}} + 2\mu_{\hat{v}} - \frac{\sigma_{\hat{v}}^2}{\sigma_{\hat{v}}^2} \right]}_{(2)} + \underbrace{\sigma_{\hat{v}}^2 \left[\frac{\text{cosk}(\hat{v}, v)}{\rho(v, \hat{v})} - \text{sk}(\hat{v}) \right]}_{(3)} > 0,
\end{aligned}$$

where in (*) we use that $\text{cov}(v, \hat{v}) > 0$. The last inequality holds if (1) < 0 , (2) < 0 , and (3) > 0 . Note that (3) > 0 follows from $c_2 > 0$. Furthermore, (1) < 0 if:

$$\frac{\text{cosk}(\hat{v}, v)}{\rho(v, \hat{v})} < -2\frac{\mu_{\hat{v}}}{\sigma_{\hat{v}}} - \frac{\sigma_{\hat{v}}^2}{\sigma_{\hat{v}}^2} =: k, \quad (51)$$

where $k < 0$. From (3) > 0 and (1) < 0 it also follows that (2) < 0 , which can be seen as follows:

$$\text{sk}(\hat{v}) \stackrel{(3)>0}{<} \frac{\text{cosk}(\hat{v}, v)}{\rho(v, \hat{v})} \stackrel{(1)<0}{<} -2\frac{\mu_{\hat{v}}}{\sigma_{\hat{v}}} - \frac{\sigma_{\hat{v}}^2}{\sigma_{\hat{v}}^2} < -2\frac{\mu_{\hat{v}}}{\sigma_{\hat{v}}} + \frac{\sigma_{\hat{v}}^2}{\sigma_{\hat{v}}^2}. \quad (52)$$

To close the proof, note that $\text{sk}(\hat{v}) < 0$ in combination with Equation (47) implies that:

$$\boxed{\frac{\text{cosk}(\hat{v}, b)}{\text{sk}(\hat{v})} > \rho(\hat{v}, b)}.$$

Hence, this condition and the condition that $\text{sk}(\hat{v}) < \frac{\text{cosk}(\hat{v}, v)}{\rho(v, \hat{v})} < k < 0$ are sufficient for Bunching at the Bottom if $\hat{v} \in [0, \infty)$ and $\text{cov}(v, \hat{v}) > 0$.

B Optimal Non-Linear Taxation in the Light Bulb Market

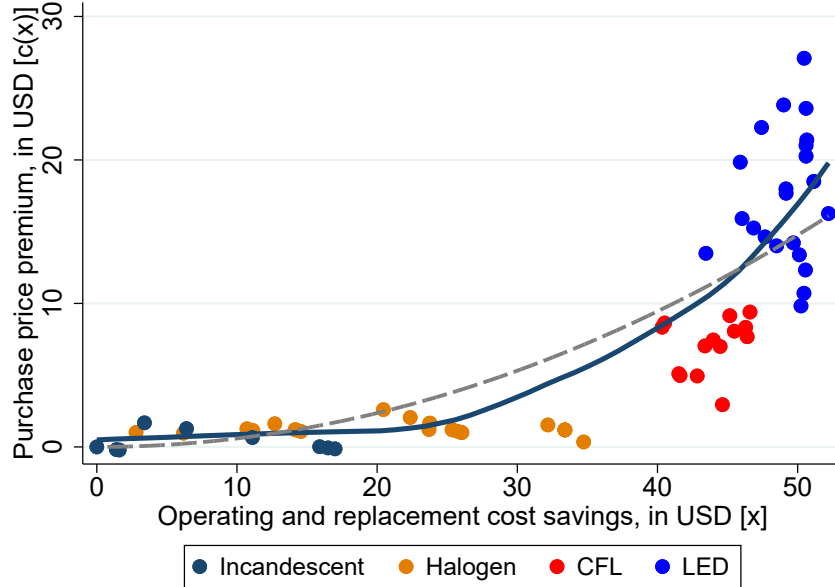
We first explain how we approximate the cost function for energy efficiency and consumers' normative valuations v . In a subsequent step, we present the aggregate welfare effects of the optimal linear and non-linear subsidy on energy efficiency, relative to no taxation.

Our supply data stems from a price comparison service, *geizhals.de*, which reports the cheapest price of a product offered on an online website. We focus on light bulbs that are typically purchased by households. In particular, we consider bulbs with an energy intensity of 25 to 75 Watt-equivalents and a warm light color of around 2700 Kelvin. To reduce the impact of branding effects, we focus on bulbs produced by one of the two large manufacturers, *Osram* and *Philips*, that offer bulbs both in the European Union and the United States. As in Allcott and Taubinsky (2015a), we express all prices in 2012 US dollars (USD) and collect product prices during that year. Some LED and CFL bulbs enter the market after 2012: in these cases, we extrapolate their 2012 price based on their aggregate annual price trends, which imply a 20% and 10% price decrease per annum for LED and CFL bulbs, respectively. For every bulb, we determine the operating and replacement cost (ORC) to consume 8.000 hours of light over eight years, which corresponds to three hours per day, assuming electricity prices of 0.1 USD per kWh (Allcott and Taubinsky, 2015a).

Based on this data, we determine the purchase price premiums and ORC savings relative to the most electricity-intensive bulb. In the following, we use ORC savings as the measure of attribute level q , i.e., of energy efficiency. Figure 2 plots the price premiums against ORC savings, which corresponds to the cost function $c(q)$ in our model. The least energy inefficient, yet cheapest, bulbs are incandescent bulbs, followed by halogen, CFL and LED bulbs. The cost curve is convex, which reflects that a one unit increase in ORC savings becomes increasingly more expensive as the level of energy efficiency increases. In 2012, the most energy efficient LED bulbs sold at a price premium of around 30 USD and yielded cost savings of about 50 USD over the course of eight years, compared to the most energy inefficient incandescent bulbs.

We use the elicitation of time preferences by Allcott and Taubinsky (2015a) to determine individual-specific discount factors. We assume that all other factors that influence normative valuations do not vary by participant and thus merely constitute a scaling factor. This assumption allows us to calibrate valuations to match the supply function from Figure 2. In particular, we set valuations to $v = s \cdot D(\delta)$, where s is a scaling factor that ensures that consumers demand every product variety offered on the market. To illustrate, consider a consumer with a discount rate of $\delta = 20\%$ and annual operating cost of 1/8 USD for eight years, which results in ORC savings of 1

Figure 2: Energy Efficiency Cost Function in the Light Bulb Market



Note: Price premiums, as well as operating and replacement cost savings are determined relative to the most electricity intensive bulb. Operating and replacement cost assume eight years of total usage (8,000 hours) and an electricity price of 0.1 USD per kWh, as in Allcott and Taubinsky (2015a). The solid line plots predictions from a local linear regression (bandwidth: 9), while the dashed line plots predictions from a regression on quadratic terms.

USD. For that consumer, the normative valuation of 1 USD in ORC savings is then $D(\delta) = (1 + 1/(1 + \delta) + \dots + 1/(1 + \delta)^7) \cdot (1/8) = 0.58$ USD. We consider this approach as a useful approximation of individuals' normative preferences v that isolates one source of heterogeneity in v and is consistent with observable market behavior.²⁴

We assume a quadratic cost function and estimate it based on the data from Figure 2. We derive consumers' choices in five scenarios: the absence of a corrective tax, the presence of the optimal linear tax, the optimal non-linear tax a) under full information, b) under a first-order approximation of the conditional bias, and c) under a second-order approximation of the conditional bias. The aggregate welfare effects are presented in Table 1.

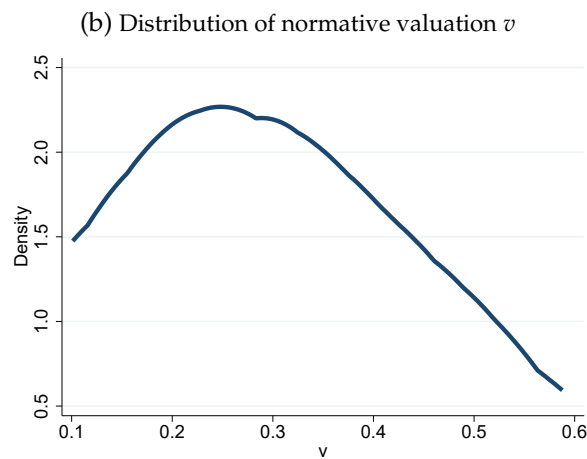
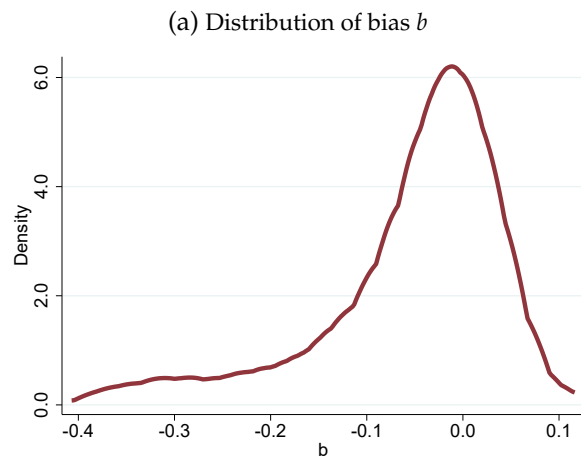
²⁴We set the scaling factor to 0.6, which ensures that the largest perceived valuations are equal to the highest gradient of the cost function from Figure 2. In addition, we impose some consistency restrictions on the data. Starting with a sample of 633 individuals with non-censored valuations in the treatment groups, we drop all observations with missing values on biases and discount rates (23 observations). We also drop observations where the elicitation of time preferences does not yield a discount factor between 0 and 1 (50 observations) and where biases or perceived valuations are above below the 1 or above the 99 percentile (14 observations). In addition, we drop all observations where perceived valuations would be negative (57 observations), which leaves us with 489 observations for our numerical example.

Table 1: Welfare Implications of Taxation

	Mean welfare, in USD/bulb	Mean welfare gain over status quo, in EUR/bulb	Mean welfare gain relative to linear tax, in %
Status quo (no tax)	3.79	0.00	—
Linear tax	3.91	0.12	0
Non-linear tax (first-order approx.)	3.96	0.16	36
Non-linear tax (second-order approx.)	3.97	0.17	44
Non-linear tax (full information)	3.98	0.19	57
<i>Welfare effects under a first-order approximation, (wrongly) setting $\rho = 0$</i>			
Non-linear tax (first-order approx, $\rho = 0$)	3.95	0.16	31

Note: Mean welfare is calculated under the optimal linear and non-linear tax schedules using the joint distribution of perceived valuations and biases, as well as the cost function estimated in Section 5 and Appendix Section B. “Non-linear tax (first-order approx, $\rho = 0$)” implements the optimal non-linear tax based on the first-order approximation to the expected bias, (wrongly) setting $\rho = 0$.

Figure 3: Distribution of bias and normative valuations



Notes for Figure a) and b): Densities estimated via kernel density estimation (Epanechnikov kernel, bandwidth: 0.03 and 0.1, respectively).