

The Optimal Use of Fines and Imprisonment if Governments Don't Maximize Welfare[#]

Ingolf Dittmann^{*}

final version: 27 March 2005
forthcoming: Journal of Public Economic Theory

Abstract:

We consider a stylized model of crime and punishment in which the prosecution policy is defined by three variables: the size of punishment, the type of punishment and the detection probability. We derive the optimal type of punishment under the assumption that the detection probability is chosen by a government whose objective function places a higher weight on the government's budget than the social welfare function does. We show that for serious crimes exclusive imprisonment is welfare maximizing. If costs of imprisonment are taken into account, the optimal punishment is a prison term with an additional fine that is smaller or equal to the costs of the prison term. For less serious crimes, fines without imprisonment are welfare maximizing.

Therefore, this paper demonstrates that the standard result of the literature that fines should be used whenever feasible need not hold in the presence of a rent-seeking government. Moreover, it offers a new explanation for the wide-spread use of mandatory imprisonment for serious crimes.

JEL Classification Codes: K42, D72, H1

[#] I would like to thank Clive Fraser, Wolfgang Leininger, Ernst Maug, Alistair Munro, Dilip Mookherjee, Wolfram Richter, Mark Walker, and seminar participants at the University of California at Los Angeles and San Diego, the University of Dortmund, and the University of Leicester, and especially an anonymous referee for helpful discussions and comments. A previous version of this paper has been presented at the WEAI in Vancouver, the Econometric Society World Congress in Seattle and the meeting of the German Economic Association in Berlin. The first version of this paper was written while the author enjoyed the hospitality and stimulating atmosphere of the Public Sector Economic Research Centre at the University of Leicester. Financial support from the Rudolf Chaudoire Stiftung is gratefully acknowledged.

^{*} Humboldt-Universität Berlin, School of Business and Economics, Spandauer Str. 1, 10178 Berlin, Germany; e-mail: dittmann@wiwi.hu-berlin.de

1. Introduction

The literature on the economics of crime, pioneered by Becker (1968), Stigler (1970) and Polinsky and Shavell (1979, 1984), applies the economic approach of decisions under uncertainty to criminal behavior, assuming that the behavior of criminals does not differ in principle from the behavior of other economic agents. The main objective of this literature is to derive the socially optimal prosecution policy, which consists of the optimal size and type of punishment and of the optimal probability with which an offender is detected and punished. A prominent result of this literature is that imprisonment is not optimal if the offender is able to pay a fine instead. The economic intuition is straightforward: Imprisonment is costly for society, as prisons must be maintained and prisoners are hindered to take up legal employment. On the other hand, fines can be used to finance the police or to compensate victims. Therefore, Polinsky and Shavell (1984) conclude that it is desirable to use the fine to its maximum feasible extent before possibly supplementing it with an imprisonment term.

This result is at odds with reality, however, because serious crimes are typically punished with mandatory imprisonment, so that the offender must serve a prison term although he could pay a (higher) fine. Moreover, the fines that are used in addition to such a mandatory prison term are usually quite small. In the United States of America, to give an example, every prison sentence is accompanied by a fine whose range depends on the offense level, a measure of the seriousness of the crime.¹ For the lowest offense level that is associated with mandatory imprisonment, the maximum fine is \$30,000 which is smaller than the wealth of most Americans.² Polinsky and Shavell (2000) point out that most prisoners are poor and would not be able to pay a high fine. As long as there are some criminals who could pay the fine however, this fact cannot

¹ In contrast to the U.S., fines on top of a prison sentence are unusual in the UK and Germany. In these two countries, fines and prison terms are regarded as alternative forms of punishment with imprisonment being reserved for serious offenses.

² The median family net worth – the difference between families' gross assets and their liabilities – was \$71,600 in 1998 (see Kennickell, Starr-McCluer and Surette, 2000). The maximum fine depending on the offense level can be found in §5E1.2 of the United States Sentencing Commission Guidelines Manual.

explain *mandatory* imprisonment, because it is still optimal that wealthy criminals pay the fine instead of being imprisoned. Even crimes that are typically committed by wealthy individuals, like bribery, embezzlement or environmental crimes, are often punished with mandatory imprisonment.³ Waldfogel (1995) finds in an empirical study on U.S. federal fraud offenses that fines average less than 3 months' preconviction income for offenders punished with both a fine and a prison term.

This paper offers a new explanation for the wide-spread use of mandatory imprisonment for serious crimes that is based on agency conflicts.⁴ In particular, we assume that the government which decides on the police expenditures and thereby on the detection probability does not maximize social welfare but instead a weighted sum of social welfare and the residual budget. The residual budget is the government's budget that remains after the police has been paid and fines have been collected. Potential reasons for the government's interest in the residual budget are that it might have the right to decide what it is spent on or that the compensation of government officials increases with the size of the budget they oversee. We show that under this assumption mandatory imprisonment is welfare maximizing for serious crimes, i.e. for crimes associated with high harms and high punishments. If costs of imprisonment are taken into account, the optimal type of punishment is a combination of mandatory imprisonment and a fine, where the fine covers the costs of the prison term or a part of it. For less serious crimes, fines without imprisonment are optimal.

In order to develop an intuition for this result, consider the two opposite effects an increase in the detection probability has on the revenues from fines. First, the proportion of criminals that are detected and fined increases, which raises the revenues from fines. On the other hand, more individuals are deterred from committing the crime, so that the

³ For the following offenses, the United States Sentencing Commission Guidelines Manual provides for mandatory imprisonment even though such offenders can be expected to be wealthy: commercial bribery if the bribe or the improper benefit to be conferred exceeds \$5,000; criminal infringement of copyright or trademark if the retail value of the infringed items exceeds \$40,000; embezzlement if the loss exceeds \$120,000; evasion of export controls; repetitive discharge of a hazardous or toxic substance into the environment.

⁴ Chu and Jiang (1993) and Levitt (1997) present alternative models in which imprisonment is employed before fines are used to their feasible extent. Both models, however, cannot explain mandatory imprisonment with small or no fines for serious crimes.

number of criminals decreases and revenues from fines drop. For less serious crimes which are associated with small punishments, the crime rate is comparatively high, so that the first effect dominates the second one. Here, the government will choose a detection probability that is higher than optimal in order to increase the residual budget. For more serious crimes on the other hand, the crime rate is comparatively small so that the second effect dominates and the government will choose a smaller-than-optimal detection probability; it allows a higher-than-optimal crime rate in order to reap higher revenues from fines.⁵ For severe crimes, this policy will lead to a substantial loss to social welfare so that it becomes optimal to switch from a fine to mandatory imprisonment. Thereby, the perverse monetary incentives disappear and the government will be eager to choose a much higher detection probability in order to effectively deter criminals.

Garoupa and Klerman (2001, 2002) also present a model with a rent-seeking government. They show that the prosecution policy chosen by this government differs significantly from the socially optimal prosecution policy. In our paper, we investigate whether a law that restricts the type of punishment to mandatory imprisonment can mitigate this distortion and lead to higher social welfare than letting the government choose the type of punishment as well. Formally, our paper therefore resembles Polinsky (1980), Friedman (1984), Besanko and Spulber (1989) and Garoupa (1997) who consider the delegation of law enforcement from a welfare maximizing principal to a rent-seeking agent. Note, however, that there is no such principal in our model and we therefore do not consider elaborate incentive contracts. Instead, we ask whether a simple constitutional restriction of the rent-seeking government is socially optimal.

In spirit, our paper is similar to Friedman (1999), who points out that ‘inefficient’ punishments, in particular imprisonment, can be optimal if they can reduce the rent-seeking behavior of the enforcers. Friedman argues that rent-seeking leads to over-enforcement, i.e. to too many and possibly false convictions. Indeed, our model produces this result for less serious crimes. For serious crimes, however, the model

⁵ If fines are interpreted as taxes on the criminal activity, lowering the detection probability is tantamount to increasing the tax-base. In this sense, our argument is similar to Gordon and Wilson (1999) who show that officials with preferences for a high government’s budget will favor policies that increase the tax-base.

implies that enforcers have an incentive to *under-enforce*, thereby causing lower deterrence and higher crime rates. Note that under-enforcement is potentially a more serious problem than over-enforcement, because it cannot be corrected ex-post. Given the presence of an independent judiciary, over-enforcement can be easily offset by discarding some of the evidence. If enforcers know that evidence beyond the socially optimal threshold will not be considered, they have no incentive to over-enforce. A similar argument can be made for false accusations. In contrast, under-enforcement cannot be mitigated ex-post. The only way to tackle under-enforcement is to change the incentives ex-ante, which can be achieved by inefficient punishments like mandatory imprisonment.

The rest of the paper is organized as follows. The next section presents the basic framework of our analysis. Section 3 describes the government's choice of the detection probability given the type of punishment and Section 4 derives the optimal type of punishment. Costs of imprisonment are introduced to the model and analyzed in Section 5, and Section 6 contains conclusions and further discussions.

2. A model of crime and prosecution

2.1 The population of potential criminals

We consider a population of infinitely many risk-neutral individuals who are homogenous in all respects except the extra utility A_i individual i can gain by committing a crime. This extra utility, whose size is the individual's private knowledge, is the sum of all gains and costs associated with committing the crime, possibly including moral costs, but it does not include expected costs from punishment. Let $f(A)$ be the publicly known probability density function of A_i in the population and $F(A)$ the corresponding cumulative distribution function. If an individual commits the crime, social welfare is reduced by the harm H , either because the crime results in such a damage to public property or because another randomly chosen individual suffers the harm H .

An individual who committed the crime is convicted with probability r . If convicted he incurs a loss P , which is measured in dollars (just as A_i). Depending on the type of punishment τ , this loss can be due to a fine (in which case the dollar amount P is paid to the government), an equivalent prison term or a mixture of the two. The variable $\tau \in [0, 1]$ denotes the proportion of the punishment P that is paid as a fine. The

remaining part of the punishment, $(1 - \tau)P$, is a prison sentence. We implicitly assume that for each dollar amount P there is a prison term of a certain length such that all individuals are indifferent between paying the fine P and serving this equivalent prison term. We do not require that the relation between the fine P and the equivalent length of imprisonment is linear. If the punishment is a fine ($\tau = 1$), we suppose that all individuals can pay this fine, so there are no wealth restrictions. Therefore, there is mandatory imprisonment if $\tau < 1$.⁶ Arguably, these assumptions are restrictive, but they enable us to study the type of punishment τ independently from the size of punishment P . In Section 6, we will show that our qualitative results continue to hold if wealth restrictions are introduced.

It follows that individual i commits the crime if and only if $A_i > rP$. Consequently, the proportion of criminals in the population is $q = \int_{rP}^{\infty} f(A)dA = 1 - F(rP)$. Note that q only depends on the deterrence rP and does not directly depend on the type of punishment τ .

2.2 The prosecution policy and social welfare

The prosecution policy is described by three variables: The detection probability r , the type of punishment τ , and the size of punishment P . We focus our attention on the optimal choice of r and τ and assume that P is fixed exogenously. If P were endogenous, it would always be optimal to choose it as high as possible in our model. This standard ‘maximum-punishment’ result (see, e.g., Becker, 1968), which is clearly at odds with practical evidence, often does not hold in more detailed models. See Polinsky and Shavell (2000) for a detailed review of major justifications for not enforcing maximum punishments.

The welfare function consists of the loss from criminal activity $-qH$, the revenues from fines τqrP , and the costs of policing $-c(r)$, which have to be paid out of the government’s budget:

$$W(r, \tau) = -qH + \tau qrP - c(r) \quad (1)$$

⁶ In other words we define mandatory imprisonment as a prison term that must be served even if the offender has enough wealth to pay a higher fine.

This is a Benthamite social welfare function with zero weight given to criminal utility and the population being normalized to one. Since the emphasis of our model is on serious crimes, we follow Stigler's (1970) argument and do not include the criminal utility $A_i - rP$ in the social welfare function, in order to preclude socially beneficial crimes. Note that (1) does not include costs of imprisonment, since our results are easier to derive and to understand in the absence of these costs. Section 5 will introduce costs of imprisonment and analyze their effect on our findings.

2.3 The government's role and objectives

We assume that the government can only choose the detection probability r , but not the type of punishment τ which is fixed by law. The objective function of the government is a weighted sum of social welfare (1) and the government's residual budget $RB(r, \tau) = \tau qrP - c(r)$:

$$U_G(r) = \alpha RB(r, \tau) + (1 - \alpha) W(r, \tau) = -(1 - \alpha)qH + \tau qrP - c(r) \quad (2)$$

where $\alpha \in [0, 1]$ is the 'budget preference parameter'. $RB(r, \tau)$ is the remaining budget after the police has been paid and fines have been collected. The assumption that the government has an interest in generating a large residual budget (i.e., $\alpha > 0$) can be justified by a number of arguments (see Gordon and Wilson, 1999). First, government officials have obviously more power the larger the budget is they oversee. In addition, their compensation typically depends on their responsibilities, i.e. it increases with the size of their budget. A further argument is that the residual budget can be used to reward swing voters and thereby to increase the chance of re-election. A part of it might even be used for perquisites, i.e. for the government's own rather than the public benefit. These arguments rely on another implicit assumption, namely that the government's budget is not tight. This means that government savings or additional income are not redistributed to the taxpayers immediately, which seems to be a reasonable assumption given that tax rates are typically quite stable. The idea of a voting restriction is also included in the formulation of (2), because social welfare still enters the government's utility function. We only claim that voting restrictions do not

prevent the government from partially maximizing its own welfare.⁷ Note that an alternative interpretation of the utility function (2) is that the government only bears a fraction $(1 - \alpha)$ of the social damage.

Given that the type of punishment is fixed by law, the question arises why the detection probability is not also fixed by law, but instead determined by the rent-seeking government. The main reason is that it is easy to enforce the type of punishment as it is observable, whereas it is practically impossible to enforce the government's effort that determines the detection probability. In addition, the optimal choice of detection probability is likely to change over time depending on the technologies of crime and prosecution, so that a fixed probability might not be optimal.

Formally, the model can be considered as a three-stage game. In the first stage, the socially optimal type of punishment τ is determined given the expected behavior of the government and the individuals in later stages. In the second stage, the government chooses the detection probability r , and in the third stage individuals decide whether or not to commit a crime. Finally, criminals are prosecuted and punished as announced in the first two stages.

2.4 Technical assumptions

For the derivation of our results, four additional technical assumptions on the functions $f(\cdot)$ and $c(\cdot)$ are needed:

- (A1) The cost function $c(r)$ is convex and continuously differentiable with $c'(0) = 0$ and $\lim_{r \rightarrow 1} c(r) = \infty$.

The last part of this assumption states that, no matter how much money is spent on police, it is never possible to clear up all crimes with certainty. Technically, the last two conditions prevent the corner solutions $r = 0$ and $r = 1$, respectively, and thereby a number of tedious case distinctions.

- (A2) The density function $f(A)$ is continuous with $f(A) > 0$ for all $A \geq 0$.

⁷ In most countries, general voting is not on individual policies or individual government officials but rather on complex bundles of different issues, only one of which is the considered prosecution policy. Together with the fact that elections are rather infrequent, this suggests that citizens are not able to fully control all government policies.

Note that we do not require that $f(A) = 0$ for negative A . Therefore, the model allows for individuals who never commit the crime even in the absence of policing. The assumption that $f(A)$ is positive for all positive A generates the feature that full deterrence ($q = 0$) cannot be achieved for finite punishments P . Apart from being realistic, this saves us some otherwise necessary case distinctions. An implication of this assumption is that there are always some individuals whose extra utility A_i exceeds the harm H . As the welfare function (1) does not include illicit gains from crime, a crime-rate of zero is nevertheless welfare maximizing (if costs of policing and revenues from fines are not taken into account).

$$(A3) \quad \text{Monotone hazard rate: } \frac{d}{dA} \frac{f(A)}{1-F(A)} \geq 0, \text{ for all } A \geq 0.$$

Intuitively this means that a marginal increase in the deterrence rP leads to a larger proportion of the remaining criminals being deterred the larger the deterrence rP is. It is a standard assumption in contract theory and is satisfied by a wide range of distributions, including the normal and the exponential distribution (See Fudenberg and Tirole, 1991, p. 267).

$$(A4) \quad \frac{d}{dr} \frac{c'(r)}{f(rP)} \geq 0 \text{ for all } r \in [0, 1) \text{ and all } P > 0.$$

Note that this assumption always holds if $f(A)$ monotonically decreases for positive A (e.g., for the normal distribution with non-positive expectation or the exponential distribution). Even if (A4) does not hold globally, it must hold for sufficiently large P and $r > 0$, because every density function eventually decreases. Hence, the results for large P will still hold if (A4) is violated.

3. The government's behavior

An important feature of the model is the complete separation of the two variables 'size of punishment' and 'type of punishment', which enables us to carry out comparative statics with respect to the type of punishment while keeping the size of punishment constant. In this section, we will use this tool to describe the detection probability r which the government chooses for a given proportion of fines τ .

We start by formulating three basic results which will help to develop an understanding of how the model works. All three results have already been derived in similar models of previous papers. The proofs of all results can be found in the appendix.

Lemma 1: If the government chooses the detection probability \hat{r} that maximizes its utility function $U_G(r)$ for a given proportion of fines τ , we obtain:

- a) \hat{r} strictly increases in the harm H .
- b) \hat{r} strictly decreases in the budget preference parameter α .
- c) The deterrence $\hat{r}P$ strictly increases in the size of punishment P , so that the crime rate $q = q(\hat{r}P)$ strictly decreases in P . For P large enough, \hat{r} eventually decreases in P and converges to zero.

Lemma 1a implies that, *ceteris paribus*, more harmful crimes are policed more thoroughly. Lemma 1b states that the government spends less on policing the stronger its interest in the residual budget is. In particular, the detection probability chosen by a rent-seeking government is always lower than the socially optimal detection probability. Polinsky (1980) and Garoupa and Klerman (2002) already established result (a) and, for high harms, result (b). In contrast to our model, these two papers include criminal utility in the government's utility function, which is the reason why they also find that \hat{r} might increase in α for small harms. Lemma 1c implies that, if the punishment P is already large, the government uses an additional increase in P to achieve two objectives: First it lowers the crime rate by increasing the deterrence, and second it economizes on the police by lowering \hat{r} .⁸ Basically, this is just the rationale of Becker's (1968) high-fine-low-probability result: As the fine increases, the costly detection probability can be lowered without reducing deterrence.

Proposition 1: If the government chooses the detection probability \hat{r} that maximizes its utility $U_G(r)$ for a given proportion of fines τ , we obtain:

- a) If the size of punishment P is large enough, \hat{r} strictly decreases in τ . Consequently, fines lead to more crime than imprisonment if punishment is harsh.
- b) If P is small, the optimal detection probability \hat{r} strictly increases in τ , so that fines lead to less crime than imprisonment.

⁸ A similar result to Lemma 1c has been derived by Garoupa (2001).

c) As P increases without bound, the crime rate q converges to zero only if $\tau = 0$. If the government collects some fines from convicted criminals ($\tau > 0$), the crime rate q converges to a strictly positive number as P increases without bound.

The main statement of Proposition 1 is that more fines increase the government's effort to fight crime only if the size of punishment is low. If punishment is harsh, the government's effort is lower the larger the proportion of fines is.⁹ To develop an intuition for this result, note that an increase in the detection probability r has two effects on revenues from fines: On the one hand, the proportion of criminals that are fined increases, which leads to a rise in revenues from fines. On the other hand, more criminals are deterred, so that the proportion of criminals decreases and the revenues from fines drop. If the punishment is small, the crime rate is relatively high and the first effect dominates the second one. Here, an increase in the proportion of fines τ leads to more policing and less crime. If the punishment is large on the other hand, the crime rate is low and the second effect dominates. Here, an increase in τ leads to less policing and more crime.

This argument can be substantiated by investigating the first-order condition of the government's maximization problem:

$$(1 - \alpha)HPf(rP) + \tau Pq(rP) = c'(r) + \tau rP \cdot Pf(rP) \quad (3)$$

The two terms on the left hand side are the marginal utility from less harm and the marginal revenue from more frequent fines due to an increase in r . On the right hand side, we have the marginal costs of increasing r and the marginal loss of fines due to less criminals. Since the inverse hazard rate $q(rP)/f(rP)$ is well-behaved according to assumption (A3), the loss due to less criminals dominates the utility from more frequent fines if P is large enough. If τ increases, the net loss of these two effects increases and the optimal r decreases. Proposition (1c) demonstrates that this result persists if the government can punish criminals harsher and harsher. It implies that full deterrence can be achieved asymptotically only with an exclusive prison term.

⁹ Polinsky (1980) and Garoupa and Klerman (2002) derive a similar result in a model in which prosecution is delegated to a private agency.

An initially surprising result is that all statements of Proposition 1 also hold for $\alpha = 0$, i.e., if the government maximizes social welfare. So even the socially optimal detection probability decreases in the proportion of fines if the punishment is high.

A consequence of Proposition 1 is that, although the fines and prison terms we compare are equivalent in the sense that individuals are indifferent between the two, the two types of punishment have different ‘net’ deterrent effects: For small punishments, fines lead to higher deterrence, whereas for large punishments imprisonment deters more potential criminals.

4. The optimality of mandatory imprisonment

After having derived the government’s choice of the detection probability given the type of punishment, we can now shift our attention to the socially optimal type of punishment given the government’s behavior.

Proposition 2: *The optimal type of punishment*

- a) If the government is benevolent ($\alpha = 0$), an exclusive fine, i.e. $\tau^* = 1$ is optimal for all P and H .
- b) If the government favors a high residual budget ($\alpha > 0$), exclusive imprisonment, i.e. $\tau^* = 0$, is welfare maximizing if the punishment P and the harm H are large.

Proposition 2a contains the standard result that fines are optimal whenever feasible, but Proposition 2b demonstrates that this need not be the case if the government does not maximize social welfare. Even a small bias towards over-valuing the residual budget can overturn this result for sufficiently large P and H .

The statement of Proposition 2b is striking: Mandatory imprisonment without any additional fines is optimal if the crime is serious. An intuition can be easily developed with our insights from the previous section: If the government places more weight on the residual budget, it will choose a lower-than-optimal detection probability \hat{r} (Lemma 1b). As long as H is small, the loss from this sub-optimal policy is small. However, if the crime is serious enough, i.e., if H and P are high, the harm from this sub-optimal policy can be substantial while revenues from fines are comparatively small, because most of the criminals are already deterred due to the high P (Lemma 1c). As a result, it becomes optimal to restrict the government to punish criminals with an exclusive prison term.

Technically, the driving force for this result is Proposition 1c, that the crime rate converges to zero only if the proportion of fines is zero. Analogous to Lemma 1b, the asymptotic crime rate increases with the budget preference parameter α if fines are used, so that it becomes optimal to force the government to lower the crime rate by prescribing mandatory imprisonment if punishment and harm are high. Thereby, the perverse monetary incentives to generate revenues from fines by allowing a high crime rate are eliminated and the government's objective function is better aligned with social welfare.

Proposition 2a provides the optimal choice of detection probability if there are no conflicts of interest between government and social welfare. This is the ideal first-best solution which is unattainable if $\alpha > 0$. For this case, Proposition 2b gives the second-best solution, which is clearly better than letting the government also choose the type of punishment τ . If the government could choose both, r and τ , it would always choose an exclusive fine and, consequently, a low detection probability, resulting in sub-optimally low social welfare for serious crimes.

5. Costs of Imprisonment

Up to now, we did not include costs of imprisonment in the analysis and the reader might wonder if our results continue to hold if these costs are taken into account. The answer is, that mandatory imprisonment is still optimal but that it is combined with a fine that covers the costs of imprisonment or a part of it. We did not include costs of imprisonment from the very beginning in order to isolate and better understand the effect of fines on the decision of the government. If these costs are taken into account, the problem becomes more involved and only partial results can be derived.

Let us assume that imprisoning an individual for a length of time that results in a loss of P to the individual costs the government $J \cdot P$, where J (for 'jail') is a positive constant. Then social welfare is given by

$$\tilde{W}(r, \tau) = -qH + \tau qrP - c(r) - (1 - \tau)JqrP = -qH + \tilde{\tau}qrP - c(r) \quad (4)$$

$$\text{with} \quad \tilde{\tau} = \tau - (1 - \tau)J, \quad (5)$$

and the government's utility function becomes

$$\tilde{U}_G(r) = -(1 - \alpha)qH + \tilde{\tau}qrP - c(r). \quad (6)$$

A simple comparison of the utility and welfare functions (4) and (6) with the according functions without prison costs (1) and (2) immediately implies that all

previous results continue to hold for $\tilde{\tau}$ instead of τ if $\tilde{\tau} \geq 0$. For negative $\tilde{\tau}$, the first order condition can have more than one solution so that only the corresponding results for large punishments P can be derived.

Note that $\tilde{\tau} = 0$ is equivalent to $\tau = J/(1+J)$. In this case, the proportion of the total punishment achieved with imprisonment is $1 - \tau = 1/(1+J)$. This prison term costs $J \cdot 1/(1+J) P$ which is equal to the fine τP , so the fine exactly covers the costs of the respective prison term.

Lemma 2: Let $J > 0$. If the government chooses the detection probability \hat{r} that maximizes its utility $\tilde{U}_G(r)$ for a given proportion of fines τ , we obtain:

- a) If $\tau > J/(1+J)$, all the results of Lemma 1 and Proposition 1 continue to hold.
- b) If $\tau \leq J/(1+J)$ and P large enough, \hat{r} eventually decreases in P and converges to zero.
- c) If $\tau \leq J/(1+J)$, the crime rate q converges to zero as P increases without bound.

Proposition 3: *The optimal type of punishment in the presence of costs of imprisonment*

- a) If the government is benevolent ($\alpha = 0$), an exclusive fine, i.e. $\tau^* = 1$ is optimal for all P and H .
- b) For $\alpha > 0$, a combination of mandatory imprisonment and a fine is optimal for large P and H , with the fine being smaller or equal to the costs of the prison term, i.e., $\tau^* \leq J/(1+J)$.

The formal result which we prove in the Appendix is that every $\tau \leq J/(1+J)$ results in higher social welfare than any $\tau > J/(1+J)$. The reason is that according to Lemma 2c, the crime rate converges to zero for all $\tau \leq J/(1+J)$. This implies that social welfare also converges to zero, so that all $\tau \leq J/(1+J)$ are asymptotically optimal. For finite numbers P and H , there is a unique $\tau^* \in [0, J/(1+J)]$ which depends on the parameters P , H , and α and on the functions $c(\cdot)$ and $f(\cdot)$, so that it cannot be derived explicitly under the weak assumptions of our model.

As noted above, the fine exactly covers the costs of imprisonment if $\tau = J/(1+J)$. In this case, the revenues from fines exactly offset the costs of the respective prison term and the government has no gain from punishing an individual apart from deterring criminals. Here, the perverse monetary incentives of fines are completely eliminated.

Still, the government will choose an inefficiently low detection probability in order to economize on the costs of policing. This inefficiency vanishes asymptotically, but for finite P it might be optimal to choose τ even smaller than $J/(1+J)$ in order to induce the government to increase the detection probability.

The welfare function (1) used in this paper gives zero weight to the utility of criminals. Following Kaplow and Shavell (2001), one could argue that this welfare function is inconsistent with the Pareto principle. This is true if, *ceteris paribus*, the utility of an honest individual is independent of the utility of criminals. However, if honest individuals prefer criminals to have low utility, our welfare function (1) might very well be consistent with the Pareto principle. The more standard welfare function gives the same weight to all individuals and therefore includes gains from the illegal activity and disutility from imprisonment:

$$\tilde{W}(r, \tau) = \int_{rP}^{\infty} (A - H)f(A)dA - (1 - \tau)qrP - c(r) \quad (7)$$

As $f(A) > 0$ for all $A > 0$, (7) implies that there are socially beneficial crimes that should not be deterred. Under the first-best solution, the crime rate q will therefore converge to a positive, efficient rate \bar{q} as the size of the punishment P increases. Proposition 1 (a) and (b) continue to hold when (1) is replaced by (7). Moreover along the line of argument of Proposition 3, one can show that now $\tau^* \leq 1 - \alpha$, i.e. mandatory imprisonment (possibly complemented by a fine) is optimal when $\alpha > 0$. Hence, our main qualitative results continue to hold if the welfare function (7) is used instead of (1).

6. Conclusions and discussion

This paper offers a new economic explanation of why serious crimes are punished with mandatory imprisonment even if offenders could pay an equivalent fine. We argue that fines change the incentives of the government. In the absence of fines, the government minimizes the social harm from criminal activity by choosing a policy that deters most potential criminals. In the presence of fines on the other hand, the government does not only want to reduce social harm but also to increase revenues from fines, two objectives that are in conflict with each other if the considered crime is serious. In that case, the government will – compared to the no-fine case – spend less on police in order to lower the deterrence and to increase the crime rate, a policy that increases revenues from fines but also the social harm from the criminal activity. If the government puts too much weight on revenue maximization as opposed to harm minimization, it will be optimal to force the government not to use fines by requiring a mandatory prison term. Thereby, the objective of revenue maximization is reduced and the government will concentrate more on the reduction of harm.

This paper also provides a justification for the separation of powers, a common practice in modern states. Typically, the executive (the government in our model) is in charge of the police and the detection probability, whereas the size and the type of punishment are determined by the legislature (in countries with civil law) or the judiciary (in countries with common law). From the point of view of the standard literature on the economics of crime this arrangement looks cumbersome and redundant. If we assume, however, that the government does not maximize social welfare but puts a higher weight on its own budget, this arrangement looks quite intelligent because it ensures that the second-best outcome is actually attained. If the government could also choose the type of punishment, it would always use a fine and, as a consequence, allow higher-than-optimal crime rates for serious crimes.

In some countries, including the United States, the local law-enforcement agency is elected directly and independently from the government. This arrangement can be interpreted as another constitutional way to overcome the rent-seeking problem depicted in our model, although it clearly introduces other incentive problems. Our model does not account for such a direct election. Instead it should rather be interpreted as a description of those countries in which the law-enforcement agency is not directly elected but run by the executive. This is the case for most European countries.

A unique feature of our model is the absence of a welfare maximizing principal. In contrast, the literature on delegated law enforcement (e.g., Polinsky, 1980, Friedman, 1984, Garoupa, 1997, Garoupa and Klerman, 2002) assumes the presence of a welfare maximizing principal who delegates law enforcement to a rent-seeking agent. In this setting, the principal can take advantage of many mechanisms (e.g., competition between agents, or decoupling the reward from the fine) that are not available in our model. While delegated law enforcement is an appropriate model for minor crimes like traffic violations or shop lifting, it remains unconvincing for serious crimes which our model focuses on.

In our model, the only difference between fines and imprisonment is that fines enter the government's utility function with a positive sign whereas prison sentences do not enter or enter with a negative sign if costs of imprisonment are taken into account. One could argue that the whole problem could be resolved if the fine is not paid to the government but to a third party, e.g., a charitable organization (which is indeed a common practice in many European countries). However, since the government is involved in virtually all aspects of public life, it can easily extract at least a part of such a fine, e.g., by lowering transfers to the recipient of the fine or by reducing services that are also provided by this organization.

A distinct feature of the model is the separation of the two variables 'type of punishment' and 'size of punishment', which allows us to study comparative statics with respect to the type of punishment without changing the size of punishment. In order to achieve this separation we assumed that each individual is able to pay the punishment P if it is a fine. Although this does not necessarily mean that all individuals are equally wealthy, this assumption seems unrealistic if the fine P is large. Fortunately however, all the results are robust to the inclusion of wealth restrictions as long as a non-negligible proportion of individuals is able to pay the fine. If we assume that an individual with wealth w_i smaller than P pays a fine of w_i and is imprisoned for a period that corresponds to a monetary value of $P - w_i$ (if there is no mandatory prison sentencing), then the government's incentive to punish more individuals is still present although somewhat reduced compared to the situation without wealth restrictions. Assume, for instance, that the wealth of a proportion of $\theta < 1$ of individuals is equal to w and that the wealth of the remaining $(1 - \theta)$ is unrestricted (or higher than the relevant range of P). For $P > w$, the expected revenues from fines are $\tau qr[(1 - \theta)P + \theta w]$, i.e.,

the government's incentive to increase the crime rate q still increases with P . If the crime is sufficiently harmful and the punishment sufficiently harsh, mandatory imprisonment will still be optimal.

Finally, note that the present paper only investigates the deterrent effect of imprisonment. It does not take into account that imprisonment also prevents crime, because notorious criminals cannot commit any further crimes while they are imprisoned. One could argue that this 'incapacitation' effect of imprisonment alone can explain the use of mandatory imprisonment for serious crimes. Ehrlich (1981) argues however that incapacitation effects are rather small because other individuals will replace most of the convicted and successfully removed offenders, in order to exploit the prevailing opportunities for illegitimate rewards. This view is also supported by recent empirical evidence (Levitt, 1998a, 1998b, Kessler and Levitt, 1999) which shows that deterrence is more important than incapacitation in reducing crime, particularly in the case of property crime.

Appendix

Proof of Lemma 1

The first-order condition is given by equation (3) in the main text. Using (A2), it can be rewritten as

$$\frac{c'(r)}{Pf(rP)} + \tau rP = (1 - \alpha)H + \tau \frac{q(rP)}{f(rP)}. \quad (8)$$

Due to (A3) and (A4), the LHS strictly increases in r and the RHS monotonically decreases in r . In addition, (A1) ensures that there is always a unique solution to (8).

A marginal increase in H increases the RHS of (8), which can be compensated only by an increase in r . Hence the optimal \hat{r} strictly increases with H , and analogously strictly decreases in α . This proves Lemma 1a and 1b.

We relabel $D = rP$, and assume that the government chooses D rather than r . Then (8) becomes

$$\frac{c'(D/P)}{Pf(D)} + \tau D = (1 - \alpha)H + \tau \frac{q(D)}{f(D)}. \quad (9)$$

As before, the LHS strictly increases in D , whereas the RHS decreases in D . A marginal increase in P now leads to a drop in the LHS which can be compensated only by an increase in D . Hence, $D = \hat{r}P$ strictly increases in P , which proves the first part of Lemma 1c.

In order to show that \hat{r} converges to zero as P diverges, we assume the opposite, namely $\hat{r} \rightarrow r_\infty > 0$ as $P \rightarrow \infty$, which implies that $D \rightarrow \infty$. If $\tau > 0$, this immediately leads to a contradiction, as the RHS of (9) converges (the inverse hazard rate is monotonically decreasing and positive) whereas the LHS of (9) diverges. If $\tau = 0$, (8) reduces to $c'(r) = (1 - \alpha)HPf(rP)$. Here we obtain a contradiction, because the LHS diverges or converges to a strictly positive number, whereas the RHS converges to zero. The reason for the latter is that $Df(D)$ converges to zero as $D \rightarrow \infty$, because otherwise $f(D)$ would converge to zero at a rate lower or equal to $1/D$ for $D \rightarrow \infty$; in that case it would not be integrable and could not be a density function. This proves Lemma 1c.

Proof of Proposition 1

We rewrite (9) as
$$\frac{c'(D/P)}{P} + \tau Df(D) = (1 - \alpha)Hf(D) + \tau q(D) \quad (10)$$

Case 1, $\tau = 0$: For every $\tilde{D} > 0$, we can find a \tilde{P} that satisfies the first-order condition $c'(\tilde{D}/\tilde{P}) = (1 - \alpha)H\tilde{P}f(\tilde{D})$, because – for \tilde{D} fixed – the LHS strictly decreases in P (from ‘ ∞ ’ to 0) and the RHS strictly increases in P (from 0 to ‘ ∞ ’). According to Lemma 1c, $D > \tilde{D}$ for all $P > \tilde{P}$. Consequently, D is not bounded and diverges as $P \rightarrow \infty$, which proves the first part of Proposition 1c.

Case 2, $\tau > 0$: Rewriting (10) yields
$$D = -\frac{c'(r)}{\tau Pf(D)} + \frac{(1 - \alpha)H}{\tau} + \frac{q(D)}{f(D)}, \quad (11)$$

Due to (A3) and Lemma 1c, the last term decreases monotonically in P and, being strictly positive, converges to some non-negative number. Since the first term on the RHS of (11) is strictly negative for all P , this implies that D cannot diverge as $P \rightarrow \infty$. As a result, D converges as $P \rightarrow \infty$ and the first term on the RHS of (11) converges to zero, due to Lemma 1c. Consequently,

$$D \xrightarrow{P \rightarrow \infty} D_\tau^\infty = \frac{(1-\alpha)H}{\tau} + \frac{q(D_\tau^\infty)}{f(D_\tau^\infty)} - D, \quad (12)$$

We rewrite (9) as

$$\frac{c'(D/P)}{Pf(D)} = (1-\alpha)H + \tau \left(\frac{q(D)}{f(D)} - D \right) \quad (13)$$

Assume that (13) holds and consider a marginal increase in τ . If $D < q(D)/f(D)$, the RHS of (13) increases which can be compensated by an increase in D only. On the other hand, if $D > q(D)/f(D)$, a marginal increase in τ results in a decrease in D . As $\hat{r} = D/P$, \hat{r} increases in τ if $D < q(D)/f(D)$, which is the case if and only if P is ‘small’, because $D(P)$ is zero for $P = 0$ and increases monotonically in P , whereas the inverse hazard rate $q(D)/f(D)$ decreases monotonically. This proves Proposition 1b.

As P diverges, the first term on the RHS in (11) converges to zero, so that the condition $D > q(D)/f(D)$ must hold if P is large enough. This argument works for $\tau > 0$ only. For $\tau = 0$, D diverges as $P \rightarrow \infty$ (Proposition 1c), so $D > q(D)/f(D)$ must hold for large P . Hence, \hat{r} decreases in τ for ‘large’ P , which proves Proposition 1a.

Proof of Proposition 2

If $\alpha = 0$, $U_G(r) = W(r)$. Since $\max_\tau[\max_r W(r, \tau)] = \max_r[\max_\tau W(r, \tau)]$, we only need to show that $\tau = 1$ maximizes $W(r, \tau)$ for all r . But this follows trivially from the welfare function (1), because q does not depend on τ . This proves Proposition 2a.

For part (b), we show that $\exists P_0, H_0 : \forall P > P_0 \quad \forall H > H_0 \quad \forall \tau \in (0, 1] : W(0) - W(\tau) > 0$:

Let $\tau > 0$. Then $W(0) - W(\tau) = -q_0 H - c(r_0) + q_\tau H - \tau q_\tau r_\tau P + c(r_\tau)$. Due to Proposition 1c and Lemma 1c, q_0 , r_0 and r_τ converge to zero as P diverges, whereas q_τ converges to $q_\tau^\infty > 0$ and $r_\tau P$ converges to $D_\tau^\infty = (1-\alpha)H/\tau + q(D_\tau^\infty)/f(D_\tau^\infty)$ as shown in (12). Therefore we obtain:

$$W(0) - W(\tau) \xrightarrow{P \rightarrow \infty} q_\tau^\infty H - \tau q_\tau^\infty D_\tau^\infty = q_\tau^\infty \left(\alpha H - \tau \frac{q(D_\tau^\infty)}{f(D_\tau^\infty)} \right). \quad (14)$$

As H increases, D_τ^∞ increases and $q(D_\tau^\infty)/f(D_\tau^\infty)$ decreases, so that for H large enough $W(0) - W(\tau) \xrightarrow{P \rightarrow \infty} c > 0$.

Note that if τ decreases, then D_τ^∞ increases and $q(D_\tau^\infty)/f(D_\tau^\infty)$ decreases. Consequently, if the RHS of (14) is positive for τ_0 , it must be also positive for all $\tau \in (0, \tau_0)$. Setting $\tau_0 = 1$, proves our claim.

Proof of Lemma 2

Part (a) follows directly from the government’s utility function in equation (6) in the main text.

The argument in the proof of Lemma 1c for $\tau > 0$ also holds for $\tau < 0$, which proves Lemma 2b.

Note that assumption (A3) implies that $[Dq(D)]/[Df(D)] = q(D)/f(D) \xrightarrow{D \rightarrow \infty} c \geq 0$. As $Df(D) \xrightarrow{D \rightarrow \infty} 0$ (see proof of Lemma 2d for an explanation), this implies $Dq(D) \xrightarrow{D \rightarrow \infty} 0$. For $\tau < 0$, all terms in $U_G(r) = -(1-\alpha)qH + \tau q r P - c(r)$ are negative, so that zero is an upper bound. As

$qrP = Dq(D) \xrightarrow{D \rightarrow \infty} 0$ and due to Lemma 2b, this upper bound can actually be obtained asymptotically if and only if $D \xrightarrow{P \rightarrow \infty} \infty$. This proves Lemma 2c.

Proof of Proposition 3

With the reformulation of $U_G(r)$ in equation (6), Proposition 2a immediately implies that $\tilde{\tau}^* = 1$ which is equivalent to $\tau^* = 1$. (Allowing negative $\tilde{\tau}$ does not change the maximum.) This proves Proposition 3a.

For part (b), we show that $\forall \tilde{\tau}_0 \in [-J, 0] \exists P_0, H_0 : \forall P > P_0 \forall H > H_0 \forall \tilde{\tau} \in (0, 1] : W(\tilde{\tau}_0) - W(\tilde{\tau}) > 0$:

Let $\tilde{\tau}_0 \in [-J, 0]$ and $\tilde{\tau} \in (0, 1]$. Then $W(\tilde{\tau}_0) - W(\tilde{\tau}) = -q_0 H + \tilde{\tau}_0 q_0 r_0 P - c(r_0) + q_{\tilde{\tau}} H - \tilde{\tau} q_{\tilde{\tau}} r_{\tilde{\tau}} P + c(r_{\tilde{\tau}})$. Due to Lemma 2 and its proof, q_0 , r_0 , $r_{\tilde{\tau}}$ and $q_0 r_0 P$ converge to zero as P diverges, whereas $q_{\tilde{\tau}}$ converges to $q_{\tilde{\tau}}^\infty > 0$ and $r_{\tilde{\tau}} P$ converges to $D_{\tilde{\tau}}^\infty = (1 - \alpha)H / \tilde{\tau} + q(D_{\tilde{\tau}}^\infty) / f(D_{\tilde{\tau}}^\infty)$ as shown in (12). Therefore we obtain:

$$W(\tilde{\tau}_0) - W(\tilde{\tau}) \xrightarrow{P \rightarrow \infty} q_{\tilde{\tau}}^\infty H - \tilde{\tau} q_{\tilde{\tau}}^\infty D_{\tilde{\tau}}^\infty = q_{\tilde{\tau}}^\infty \left(\alpha H - \tilde{\tau} \frac{q(D_{\tilde{\tau}}^\infty)}{f(D_{\tilde{\tau}}^\infty)} \right). \quad (15)$$

As H increases, $D_{\tilde{\tau}}^\infty$ increases and $q(D_{\tilde{\tau}}^\infty) / f(D_{\tilde{\tau}}^\infty)$ decreases, so that for H large enough $W(\tilde{\tau}_0) - W(\tilde{\tau}) \xrightarrow{P \rightarrow \infty} c > 0$. If $\tilde{\tau}$ decreases, then $D_{\tilde{\tau}}^\infty$ increases and $q(D_{\tilde{\tau}}^\infty) / f(D_{\tilde{\tau}}^\infty)$ decreases. Consequently, if the RHS of (14) is positive for $\tilde{\tau}'$, it must be also positive for all $\tilde{\tau} \in (0, \tilde{\tau}')$. Setting $\tilde{\tau}' = 1$, proves Proposition 3b.

References

- Becker, G. S. (1968) "Crime and punishment: An economic approach" *Journal of Political Economy* **76**, 169-217.
- Besanko, D. and D. F. Spulber (1989) "Delegated law enforcement and noncooperative behavior" *Journal of Law, Economics, and Organization* **5**, 25-52.
- Chu, C. Y. C. and N. Jiang (1993) "Are fines more efficient than imprisonment?" *Journal of Public Economics* **51**, 391-413.
- Ehrlich, I. (1981) "On the usefulness of controlling individuals: An economic analysis of rehabilitation, incapacitation, and deterrence" *American Economic Review* **71**, 307-322.
- Friedman, D. (1984) "Efficient institutions for the private enforcement of law" *Journal of Legal Studies* **13**, 379-397.
- Friedman, D. (1999) "Why not hang them all: The virtues of inefficient punishments" *Journal of Political Economy* **107**, S259-S269.
- Fudenberg, D. and J. Tirole (1991) *Game Theory*, MIT Press: Cambridge.
- Garoupa, N. (1997) "A note on private enforcement and type-I error" *International Review of Law and Economics* **17**, 423-429.
- Garoupa, N. (2001) "Optimal magnitude and probability of fines" *European Economic Review* **45**, 1765-1771.
- Garoupa, N. and D. Klerman (2004) "Corruption and the optimal use of nonmonetary sanctions" *International Review of Law and Economics* **24**, 219-225.
- Garoupa, N. and D. Klerman (2002) "Optimal law enforcement with a rent-seeking government" *American Law and Economics Review* **4**, 116-140.
- Gordon, R. H. and J. D. Wilson (1999) "Tax structure and government behavior: Implications for tax policy" NBER working paper number W7244
- Kaplow, L. and S. Shavell (2001) "Any non-welfarist method of policy assessment violates the Pareto principle" *Journal of Political Economy* **109**, 281-286.
- Kessler, D. and S. D. Levitt (1999) "Using sentence enhancements to distinguish between deterrence and incapacitation" *Journal of Law and Economics* **42**, 343-363.

- Kennickell, A. B., Martha Starr-McCluer and Brian J. Surette (2000) "Recent changes in U.S. family finances: Results from the 1998 Survey of Consumer Finances" *Federal Reserve Bulletin* **86**, 1-29.
- Levitt, S. D. (1997) "Incentive compatibility constraints as an explanation for the use of prison sentences instead of fines" *International Review of Law and Economics* **17**, 179-192.
- Levitt, S. D. (1998a) "Why do increased arrest rates appear to reduce crime: deterrence, incapacitation, or measurement error?" *Economic Inquiry* **36**, 353-372.
- Levitt, S. D. (1998b) "Juvenile crime and punishment" *Journal of Political Economy* **106**, 1156-1185.
- Polinsky, A. M. (1980) "Private versus public enforcement of fines" *Journal of Legal Studies* **9**, 105-127.
- Polinsky, A. M. and S. Shavell (1979) "The optimal tradeoff between the probability and magnitude of fines" *American Economic Review* **69**, 880-891.
- Polinsky, A. M. and S. Shavell (1984) "The optimal use of fines and imprisonment" *Journal of Public Economics* **24**, 89-99.
- Polinsky, A. M. and S. Shavell (2000) "The economic theory of public enforcement of law" *Journal of Economic Literature* **38**, 45-76.
- Stigler, G. J. (1970) "The optimum enforcement of laws" *Journal of Political Economy* **78**, 526-536.
- Waldfogel, J. (1995) "Are fines and prison terms used efficiently? Evidence on federal fraud offenders" *Journal of Law and Economics* **38**, 107-139.
- United States Sentencing Commission (May 2000) *Guidelines Manual*, www.ussc.gov.